



VOXReality

VOICE DRIVEN INTERACTION IN XR SPACES


Advanced AI multi-model for XR analysis

WP3

31-03-2025



Funded by
the European Union

Version	2.1
WP	WP3
Dissemination level	Public
Deliverable lead	UM
Authors	Yusuf Can Semerci, Pawel Maka, Abderrahmane Issam, Gerasimos Spanakis (UM), Georgios Papadopoulos, Athanasios Ntovas, Stefanos Biliouris, Sotiris Karavarsamis, Petros Drakoulis, Alexandros Doumanoglou, Konstantinos Konstantoudakis, Dimitris Zarpalas (CERTH), Apostolos Maniatis, Stavroula Bourou (SYN), Jiahuan Pei, Irene Viola, Pablo Cesar (NWO-I)
Reviewers	Olga Chatzifoti (MAG), Carina Pamminger, Leesa Joyce (HOLO)
Abstract	This document describes the work done in the first 30 months of the project regarding the natural language processing models.
Keywords	Natural Language Processing, Automatic Speech Recognition, Neural Machine Translation, Visual Language Models, Conversation Agents
License	 <p>This work is licensed under a Creative Commons Attribution-No Derivatives 4.0 International License (CC BY-ND 4.0). See: https://creativecommons.org/licenses/by-nd/4.0/</p>

Dissemination Level: PU

PU Public

PP	Restricted to other programme participants (Including the Commission Services)
RE	Restricted to a group specified by the consortium (Including the Commission Services)
CO	Confidential, only for members of the consortium (Including the Commission Services)

Nature : RE

PR	Prototype
RE	Report
SP	Specification
TO	Tool
OT	Other

Version History

Version	Date	Owner	Editor(s)	Changes to previous version
0.1	2023-11-10	UM	Yusuf Can Semerci	Outline
0.2	2023-12-07	UM	Yusuf Can Semerci	Related works
0.21	2023-12-09	SYN	Apostolos Maniatis	Upodated Related works
0.3	2023-12-12	UM	Pawel Maka	Added Sec. 2.2
0.4	2023-12-12	CERTH	Petros Drakoulis	Added Sec. 2.3
0.41	2023-12-13	SYN	Stavroula Bourou	Updated Related works
0.5	2023-12-14	SYN	Apostolos Maniatis	Added Sec. 2.4.1
0.6	2023-12-15	UM	Abderrahmane Issam	Added Sec. 2.1, Updated Sec. 2.2
0.61	2023-12-15	UM	Yusuf Can Semerci	Updated Sec. 2.1 and Sec 2.2
0.62	2023-12-19	CERTH	Petros Drakoulis	Updated Sec. 2.3
0.63	2023-12-19	SYN	Stavroula Bourou	Updated Sec. 2.4.1
0.7	2023-12-20	NWO-I	Jiahuan Pei	Added Sec. 2.4.2
1.0	2023-12-20	UM	Yusuf Can Semerci	Final formatting and controls
1.01	2023-12-22	MAG	Olga Chatzifoti	Internal review complete
1.02	2023-12-22	HOLO	Carina Pamminger	Internal review complete
1.1	2023-12-22	UM	Yusuf Can Semerci	Revisions complete
1.11	2025-02-13	UM	Yusuf Can Semerci	Updated outline
1.2	2025-02-17	UM	Apostolos Maniatis	Added new sections under 3.4
1.3	2025-03-13	SYN	Pawel Maka	Added new sections under 3.2
1.4	2025-03-20	UM	Abderrahmane Issam	Added new sections under 3.1
1.5	2025-03-20	CERTH	Georgios Papadopoulos, Petros Drakoulis	Added new sections under 3.3
1.51	2025-03-24	UM	Yusuf Can Semerci	Updated Chapter 3
2.0	2025-03-25	UM	Yusuf Can Semerci	Final formatting and controls
2.01	2025-03-27	HOLO	Leesa Joyce	Internal review complete
2.02	2025-03-31	MAG	Olga Chatzifoti	Internal review complete
2.1	2025-03-31	UM	Yusuf Can Semerci	Revisions complete

Table of Contents

List of Abbreviations & Acronyms.....	6
List of Figures	8
List of Tables	10
Executive Summary	11
1. Introduction	12
2. Related Works	13
2.1. Automatic Speech Recognition (ASR).....	13
2.2. Neural Machine Translation (NMT)	15
2.3. Visual Language Models (VL)	17
2.4. Conversation Agents (CA)	18
2.5. Evaluation Metrics.....	21
2.5.1. Bilingual Evaluation Understudy (BLEU)	21
2.5.2. Recall-Oriented Understudy for Gisting Evaluation (ROUGE)	22
2.5.3. Metric for Evaluation of Translation with Explicit Ordering (METEOR)	23
2.5.4. Character-level F-score (chrF)	23
3. VOXReality Natural Language Processing Models	25
3.1. Automatic Speech Recognition (ASR).....	25
3.1.1. Improving the Speech Recognition Models	25
3.1.2. Speech Translation for Non-native Speech.....	28
3.1.3. Streaming Speech Recognition and Translation	30
3.2. Neural Machine Translation (NMT)	33
3.2.1. Context-aware Machine Translation.....	33
3.2.2. Contextual Attention Heads	39
3.2.3. Robust Machine Translation	43
3.2.4. Robustness of Machine Translation Against Non-native Speaker Errors.....	46
3.2.5. Simultaneous Machine Translation	52
3.2.6. Subtitle Alignment.....	54
3.2.7. Deployment	55
3.3. Vision-Language Models (VL)	57
3.3.1. Early models delivery	58
3.3.2. Metrics and limitations	59
3.3.3. Image spatial captioning	62
3.3.4. Image spatial visual question answering	69
3.3.5. Video spatial captioning	71
3.3.6. Visual navigation.....	73
3.3.7. Alternative image spatial captioning.....	75
3.4. Conversation Agents (CA)	78
3.4.1. VR Conference Conversation Agent	78



3.4.2. AR Training Assistant	104
4. Conclusions	124
References	125
Appendix I: Endpoints for the ASR and NMT components	134
Appendix II: Endpoints for the VL components	139
Appendix III: Endpoints for the navigation assistant	142
Appendix IV: Endpoints for the training assistant	143



List of Abbreviations & Acronyms

AI	: Artificial Intelligence
API	: Application Programming Interface
AR	: Augmented Reality
ARTA	: Augmented Reality Training Assistant
ASR	: Automatic Speech Recognition
BART	: Bidirectional and Auto-Regressive Transformers
BEA	: Building Educational Applications
BERT	: Bidirectional Encoder Representations from Transformers
BLEU	: Bilingual Evaluation Understudy
C4	: Colossal Clean Crawled Corpus
CA	: Conversation Agent
CNN	: Convolutional Neural Network
COCO	: Common Objects in Context
CTC	: Connectionist Temporal Classification
DDPG	: Deep Deterministic Policy Gradients
DM	: Dialogue Management
DP	: Dialogue Policy
DST	: Dialogue State Tracking
FCE	: First Certificate in English corpus
FFN	: Feed-forward Network
GA	: Grant Agreement
GAtt	: Ghost Attention
GMM	: Gaussian Mixture Model
GPT	: Generative Pre-trained Transformer
GPU	: Graphics Processing Unit
GUG	: Grammatical/Ungrammatical
HMM	: Hidden Markov Model
IC	: Image Captioning
JFLEG	: JHU FLuency-Extended GUG corpus
JHU	: John Hopkins University
LASER3	: Language-Agnostic SEntence Representations
LR	: Learning Rate
LCPT	: Large Contrastive Pronoun Test
LLM	: Large Language Models
LSTM	: Long Short-Term Memory
LXMERT	: Language-visual Model for Efficient Representations of Transformers
LoRA	: Low-Rank Adaptation
M2M-100	: Many-to-Many Multilingual
METEOR	: Metric for Evaluation of Translation with Explicit Ordering
MLM	: Masked Language Modeling
MMS	: Massively Multilingual Speech
MRM	: Masked Region Modeling
MT	: Machine Translation
MoE	: Mixture of Experts
NAS	: Neural Architecture Search
NLG	: Natural Language Generation
NLLB	: No Language Left Behind
NLP	: Natural Language Processing



NLU	: Natural Language Understanding
NMT	: Neural Machine Translation
OC	: Open Calls
PEFT	: Parameter-efficient fine-tuning
PPO	: Proximal Policy Optimization
QLoRA	: Quantized Low Rank Adaptor
ReAct	: Reason and Act
RGB	: Red-Green-Blue
RL	: Reinforcement Learning
RLHF	: Reinforcement Learning with Human Feedback
RNN	: Recurrent Neural Network
ROUGE	: Recall-Oriented Understudy for Gisting Evaluation
ReLU	: Rectified Linear Unit
SD	: Scene Description
SKS	: Semantic Knowledge Similarity
SMT	: Statistical Machine Translation
SOTA	: State-of-the-art
T5	: Text-To-Text Transfer Transformer
TEACH	: Task-driven Embodied Agents that Chat
TEACH-EDH	: TEACH - Execution from Dialogue History
UNITER	: UNiversal Image-Text Representation
VL	: Vision-Language
VQA	: Visual Question Answering
VR	: Virtual Reality
WER	: Word Error Rate
WP	: Work Package
XR	: eXtended Reality
TTS	: Text-to-Speech
CoVoST2	: Common Voice Speech-To-Text 2
ST	: Speech Translation
VAD	: Voice Activity Detection
chrF	: Character-level F-score
EMMA	: Efficient Monotonic Multihead Attention
VC	: Video Captioning
VN	: Visual Navigation
GED	: Grammatical Error Detection
GEC	: Grammatical Error Correction
POS	: Part of Speech
CKA	: Centered Kernel Alignment
ViT	: Vision Transformer
BLIP	: Bootstrapped Language-Image Pretraining
FAISS	: Facebook AI Similarity Search
MRR	: Mean Reciprocal Rank
MRA	: Memory Retrieval Accuracy
RC	: Response Correctness
CC	: Contextual Coherence



List of Figures

Figure 2.1. The architectures of a) Deep Speech and b) Deep Speech 2.....	14
Figure 2.2. The architectures of a) wav2vec, b) XLSR, and c) HuBERT.....	14
Figure 2.3. OpenAI's Whisper architecture.....	15
Figure 2.4. Seq2Seq architecture.....	15
Figure 2.5. The architectures of a) BART, b) T5, and c) M2M-100.....	16
Figure 2.6. Show and Tell architecture.....	17
Figure 2.7. The architectures of a) LXMERT, b) VisualBERT, and c) UNITER.....	18
Figure 2.8. The architecture of GPT model [Radford et al., 2018].....	19
Figure 2.9. Training process of Llama 2-Chat.....	19
Figure 3.1. Whisper Model with Adapter Modules.....	26
Figure 3.2. The speech transcription endpoint.....	27
Figure 3.3. The speech translation endpoint.....	28
Figure 3.4. The context-aware speech translation endpoint.....	28
Figure 3.5: Robust Speech Translation Using Text Noise.....	29
Figure 3.6: Chunking Based Streaming Speech Recognition.....	31
Figure 3.7: Dynamic Speech Recognition and Translation.....	32
Figure 3.8: Streaming API configuration interface.....	32
Figure 3.9. The architecture of latent grouping and latent selecting.....	34
Figure 3.10. The architecture of VOXReality context-aware MT model.....	34
Figure 3.11. The terminology-constrained architecture.....	35
Figure 3.12. BLEU on the OpenSubtitles2018 test set in all consortium languages.....	38
Figure 3.13. An example marking the different relations between words in source and context sentences.....	39
Figure 3.14 Results of the analysis with the OpenMT En-De models.....	41
Figure 3.15 Changes in accuracy for OpusMT en-de models.....	42
Figure 3.16 Accuracy of the models on the CTXPRO dataset.....	43
Figure 3.17 The BLEU scores of the models on the OpenSubtitles 2018 dataset.....	43
Figure 3.18. Training a Robust Model Using Adapters.....	44
Figure 3.19: GED probing performance of the models on En-Es direction.....	48
Figure 3.20: CKA distances of word representations on En-Es direction.....	49
Figure 3.21: Attention to the 10 most common POS tags on En-Es direction.....	50
Figure 3.22: Accuracy of Influential Heads on En-Es direction.....	51
Figure 3.23. The architecture of wait-k adapter model.....	52
Figure 3.24. BLEU-Average lagging results for a) transformer-small, b) transformer-base, and c) transformer-big.....	54
Figure 3.25. The chrF between the subtitles and the transcriptions.....	55
Figure 3.26. The text translation endpoint.....	56
Figure 3.27. The context-aware translation endpoint.....	56
Figure 3.28 The context-aware terminological translation endpoint.....	56
Figure 3.29 The terminology configuration endpoint.....	57
Figure 3.30 The subtitle upload endpoint.....	57
Figure 3.31. Early endpoints for Visual Language Models.....	59
Figure 3.32. An example image from the spatial dataset.....	61
Figure 3.33. An example from the spatial dataset with its accompanying depth and object-annotated images.....	64
Figure 3.34: A histogram showcasing the frequency of all possible spatial relations occurring between objects in the dataset.....	65
Figure 3.35. Our architecture with ViT encoder and GPT2 decoder.....	66
Figure 3.36. Spatial IC model progression across training epochs in a) standard evaluation metrics, b) variants of OURS metric and c) the respective training loss.....	68
Figure 3.37. An example image for the generic and spatial VQA comparison.....	70
Figure 3.38. The BLIP model architecture for VQA.....	70

Figure 3.39. Example frame captioned as: “A laptop is to the right of a person”	72
Figure 3.40. An example capture from the virtual space.....	74
Figure 3.41. An example capture from the virtual space with noise (tilted view)	74
Figure 3.42. BLEU score of the pre-trained vs randomly initialized model weights	75
Figure 3.43. Our alternative IC model architecture	77
Figure 3.44. Data distribution across categories.....	80
Figure 3.45. Conference agent workflow	82
Figure 3.46. Illustration of LoRA during and after training.....	87
Figure 3.47. Training and validation loss per epoch	88
Figure 3.48. The intend and destination recognition endpoint	92
Figure 3.49. The response generation endpoint	92
Figure 3.50: Conference agent workflow	93
Figure 3.51: Sample PDF document for the trade show	96
Figure 3.52: Training and validation loss per epoch	101
Figure 3.53: Two example conversations with the agent.....	104
Figure 3.54. Autonomous agent powered by large language models	105
Figure 3.55. The workflow of ARTA.....	106
Figure 3.56. The illustration of conversational ReAct walkthrough.....	107
Figure 3.57. Interactive information flow between XR test application and AI agent	107
Figure 3.58. Comparison of full fine-tuning, LoRA, and QLoRA.....	108
Figure 3.59. Examples of TEACH-EDH dataset	109
Figure 3.60. Examples of VOX-ARTA-LEGO dataset.....	110
Figure 3.61. The learning curve of Llama-2-7b-chat-hf model finetuning on a) TEACH-EDH dataset and b) VOX-ARTA-LEGO	112
Figure 3.62. Interactive WebSocket API.....	114
Figure 3.63: Workflow of the free mode of Training Assistant.....	115

List of Tables

Table 3.1. Whisper Size Variants	26
Table 3.2. Hyperparameters of Whisper Finetuning	26
Table 3.3. WER Results of Finetuning vs. Adapter finetuning	27
Table 3.4. Robust Speech Translation Results.....	29
Table 3.5. The hyper-parameters of all models trained on IWSLT 2017 dataset	36
Table 3.6. Results of the models trained on the IWSLT 2017 with EN-DE dataset	37
Table 3.7. Results of the models trained on the IWSLT 2017 with EN-FR dataset	37
Table 3.8. The hyper-parameters of the terminology-constrained model	38
Table 3.9. Accuracy of the models on the ContraPro contrastive dataset.....	38
Table 3.10. The hyper-parameters of the robust MT model.....	45
Table 3.11. A Comparison of Finetuning and Adapter based finetuning on JFLEG-es	45
Table 3.12. COMET of OPUS-MT, M2M100, MBART, and NLLB on En-Es direction	47
Table 3.13. Hyperparameters of the Simultaneous Machine Translation Models	52
Table 3.14. Performance of the approaches in Hippolytus	55
Table 3.15. Our COCO spatial extension dataset specifications	64
Table 3.16. Spatial IC model's evaluation results across training epochs.....	67
Table 3.17. Pearson correlation coefficients between all IC metrics after 50 epochs.....	68
Table 3.18. spatial IC model's evaluation results after epoch extension	69
Table 3.19. Spatial VQA model's evaluation results across training epochs.....	71
Table 3.20. The hyper-parameters of the spatial VC model.....	73
Table 3.21. Spatial VC model's evaluation results across training epochs	73
Table 3.22. The hyper-parameters of the VN model.....	75
Table 3.23. VN model's evaluation with pre-trained and randomly initialized weights.....	75
Table 3.24. The hyper-parameters of the alternative IC model.....	77
Table 3.25. The alternative IC model's object-segmentation performance	78
Table 3.26. The alternative IC model's spatial-relations classification performance	78
Table 3.27. Experimental results of the NLU model	81
Table 3.28. Evaluation of sentence-transformers fine-tuning with sentence similarity.....	89
Table 3.29. Evaluation of navigation agent before and after fine-tuning	90
Table 3.30. Evaluation of intent recognition agent before and after fine-tuning	91
Table 3.31. Evaluation of destination recognition agent before and after fine-tuning	91
Table 3.32. Evaluation of intent recognition agent.....	94
Table 3.33. Evaluation of the trade show agent.....	97
Table 3.34. Evaluation of the presentation summarization agent	98
Table 3.35. Evaluation of the one-to-one summarization agent.....	99
Table 3.36. Evaluation of sentence-transformers fine-tuning process	102
Table 3.37. Evaluation of conference agent after fine-tuning.....	102
Table 3.38. Evaluation results of conversation history	104
Table 3.39. Descriptions of the functions for the XR tools in the test application	108
Table 3.40. Statistics of TEACH-EDH and VOX-ARTA-LEGO datasets	110
Table 3.41. Evaluation results on TEACH-EDH and VOX-ARTA-LEGO dataset	111
Table 3.42. Evaluation results of the first intent recognition phase	116
Table 3.43. Evaluation results of the second intent recognition phase	120
Table 3.44. Evaluation results of conversation history management	123



Executive Summary

This document corresponds to the deliverable D3.1 and D3.2 - Advanced AI multi-model for XR analysis, of work package 3 (WP3) and describes the work done in the first 30 months of the project regarding the natural language processing (NLP) models.

In VOXReality, the natural language processing models are developed for the following tasks:

- 1) Automatic speech recognition (ASR),
- 2) Machine translation (MT),
- 3) Speech translation (ST),
- 4) Image captioning (IC),
- 5) Video captioning (VC),
- 6) Visual Navigation (VN),
- 7) Visual question answering (VQA),
- 8) Conversation agents (CA) for navigation and training assistance.

The document is divided into four chapters:

Chapter 1 provides a brief introduction to the WP3 tasks and to this deliverable.

Chapter 2 describes the background knowledge required to understand the VOXReality NLP models and the related works from the literature for tasks listed above.

Chapter 3 presents the models developed for each of the five tasks listed above. Section 3.1 describes the work performed regarding automatic speech recognition, speech translation, and streaming speech recognition. Section 3.2 describes the context-aware machine translation, simultaneous machine translation, robust machine translation and script alignment models and their performances with benchmark datasets, Section 3.3 describes the models implemented for image and video captioning, visual question answering, and visual navigation, and Section 3.4 presents the work performed on developing the navigation assistant conversation agent and the training assistant conversation agent.

Chapter 4 summarizes the conclusions obtained with this deliverable and with the work done during the first 30 months of the project.

The changes from the first version to the final version are as follows: Executive Summary, Chapter 1, Section 3.2.7, Section 3.3, Chapter 4, and All Appendices are updated and Section 2.5.4, Section 3.1.2, Section 3.1.3, Section 3.2.2, Section 3.2.4, pages between 100 and 113 in Section 3.4.1, and pages between 124 and 134 in Section 3.4.2 are added.

1.Introduction

The eXtended Reality (XR) technologies, specifically Virtual Reality (VR) and Augmented Reality (AR), offer individuals immersive environments where they can engage with and manipulate both physical and virtual environments in multidimensional ways, presenting considerable opportunities for various sectors like education, training, collaboration, entertainment, and health. The technologies that facilitate this immersive experience offer ways to engage and interact with these virtual or augmented environments through hand-held point and click interfaces (e.g., controllers) or data-driven AI-based hand tracking. However, these technologies are limited for natural interactions, especially for human-to-human communication since human-to-human interaction is mainly driven by language (verbal communication) and enhanced with gesture-based physical cues. Furthermore, the rapid spread of the COVID-19 pandemic accelerated the progress of digital meeting platforms and digital assistants. As a result, human communication and interaction underwent a notable shift toward the digital medium. The ability to organize events, engage in meetings with colleagues, and conduct training activities became feasible through digital platforms in the midst of the pandemic allowing individuals all around the world to be connected. Moreover, recent advancements in Artificial Intelligence (AI) technologies have further elevated the interaction with machines, enabling users to seamlessly employ natural language in their day-to-day activities. These observations underscore the need for a coordinated strategy in implementing context-aware, multilingual, visually grounded tools to enhance natural communication among humans and between humans and machines on digital platforms such as XR environments.

One of the main objectives of VOXReality is to implement natural language processing models which are pre-trained, publicly available, optimized, multilingual, visually grounded and knowledgeable of the domain-specific needs of applications in XR environments. To this end, the VOXReality consortium invests in research and development activities in WP3 to provide NLP models for automatic speech recognition, machine translation, visual question answering, visual captioning, spatial scene description, and conversation agents for navigation and training assistance in 4 tasks.

“Task 3.1 - VOXReality Audio processing” focuses on the algorithms, models and tools to process audio streams or files for the automatic speech recognition, “Task 3.2 - VOXReality Multilingual translation” is responsible for the implementation, training and fine-tuning of the machine translation models, “Task 3.3 - Vision and language based XR models” implements visual language (VL) models that can answer questions about visual cues, describe an image or a video by summarizing their content (caption), and provide navigation instructions by leveraging spatial understanding in an environment, and “Task 3.4 - Context aware generative dialogue system” is responsible for the implementation of two generative conversation agents that are adapted to two different tasks: VR conference conversation agent and AR training assistant.

The technical work described in this document is the work performed in all four tasks (T3.1, T3.2, T3.3, and T3.4) of WP3 until the end of the 30th month of the VOXReality project.

2. Related Works

Natural Language Processing (NLP) is the study and development of algorithms and computational models that enable machines to understand, interpret, and generate human language. Early NLP research focused on the development of rule-based systems and symbolic approaches that relied on predefined grammatical rules and linguistic structures to analyse and generate text. However, these approaches faced limitations in handling the complexity and variability of natural language. The shift towards statistical methods, probabilistic models and machine learning techniques enabled the NLP systems to learn from data and adapt to linguistic patterns facilitating more robust language processing capabilities. Recent advancements in machine learning, deep learning, the availability of large datasets and more efficient and powerful hardware setups revolutionized the NLP.

Recent years have seen the rise of transformer-based models, such as BERT (Bidirectional Encoder Representations from Transformers) [Devlin et al., 2018] and GPT (Generative Pre-trained Transformer) [Radford et al., 2018], which have demonstrated remarkable performance in a wide range of NLP tasks. The following sections present the evolution of approaches for the NLP tasks that VOXReality is investing in, namely, automatic speech recognition, machine translation, visual language models, and conversation agents.

2.1. Automatic Speech Recognition (ASR)

Automatic Speech Recognition (ASR) focuses on the development of systems capable of transcribing spoken language into text. The approaches to ASR evolved from early rule-based approaches to the current state-of-the-art methods driven by deep learning and transformers. Early systems relied heavily on handcrafted rules and linguistic models that match acoustic features with phonetic units [Jelinek, 1976]. Gaussian Mixture Models (GMMs) and Hidden Markov Models (HMMs) were utilized to develop early systems that were capable of handling limited vocabulary tasks for specific domains. HMMs were used to handle the temporal dynamics of speech, while GMMs were used to model the probability distribution of acoustic features [Rabiner, 1989]. The first revolutionary shift happened with the advancements in machine learning methodologies. Deep Speech [Hannun et al, 2014] and Deep Speech 2 [Amodie et al, 2016] are two of the most prominent examples of the ASR models introduced in the era of rise of machine learning algorithms in NLP.

Deep Speech employs an end-to-end neural network architecture (Figure 2.1a) with connectionist temporal classification (CTC) loss to train an NLP model, where acoustic features, such as spectrograms, are mapped to a character sequence without the need for phonetic or linguistic modelling and without the need for aligned input-output pairs. The CTC loss, increased number of parameters (model size) and increased amount of data utilized to train the model contributed to the success of the model in achieving the state-of-the-art at the time of publication. Deep Speech 2 builds upon the previous iteration by introducing the Recurrent Neural Network (RNN) approach (Figure 2.1b), specifically Long Short-Term Memory (LSTM), to handle long-term dependencies in the sequential data (speech), and the batch normalization, which normalizes the inputs to a layer within a mini batch, addressing issues like vanishing or exploding gradients. Furthermore, Deep Speech 2 introduces the multilingual capabilities of deep neural network architectures and surpasses its predecessor in achieving the state-of-the-art.

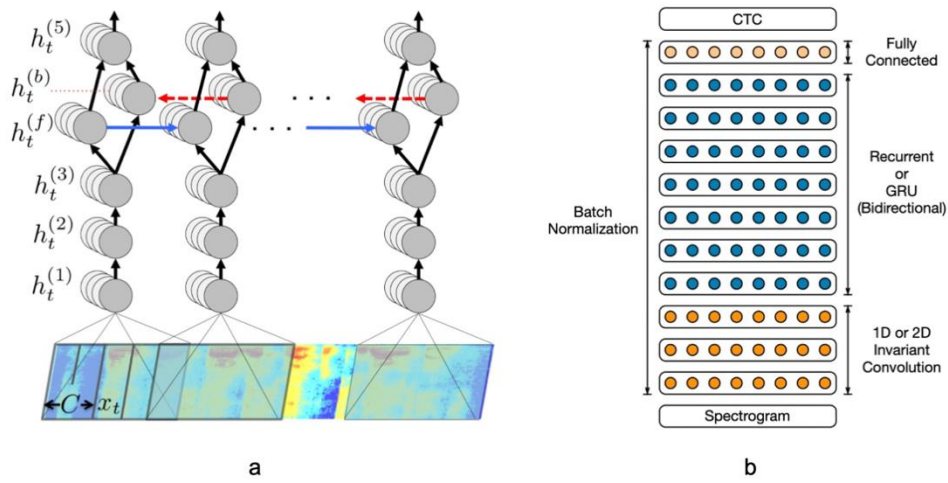


Figure 2.1. The architectures of a) Deep Speech and b) Deep Speech 2

The introduction of attention mechanisms and transformer models [Vaswani et al., 2017] accelerated the NLP research tremendously, especially with the success of models such as BERT. Recent ASR approaches employ speech representation models such as wav2vec [Baevski et al., 2020], XLSR [Conneau et al., 2020] and HuBERT [Hsu et al., 2021], which are three of the most prominent speech models benefited from the methodologies introduced with attention mechanisms, transformers and BERT. All three models are trained with a self-supervised learning approach where a model is trained on a task generated from its own input data that allow the model to learn meaningful representations without relying on externally labelled datasets. wav2vec introduces a framework (Figure 2.2a) that leverages contrastive learning, a self-supervised learning approach where a model is trained to maximize the similarity between augmented views of the same instance and minimize the similarity between views of different instances, while XLSR introduces the functionality of multilingual quantization (Figure 2.2b) to the contrastive learning approach, where the aim is to capture the shared phonetic representations across languages. HuBERT, on the other hand, employs the masked language modeling (MLM) approach (Figure 2.2c), where random tokens in a given language sequence are masked and the model is trained to predict the masked tokens based on the surrounding unmasked tokens, to capture contextual dependencies in the speech signal.

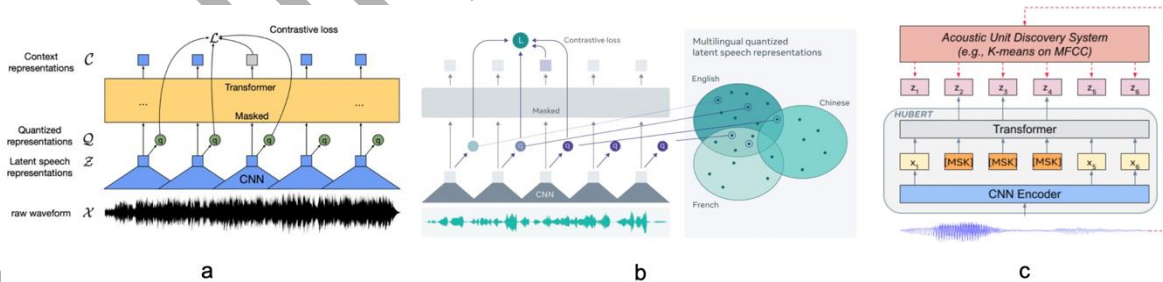


Figure 2.2. The architectures of a) wav2vec, b) XLSR, and c) HuBERT

One of the most recent advancements in ASR is Whisper [Radford et al., 2023], which is an end-to-end model based on transformers and utilizes log-Mel spectrogram of audio inputs. Whisper is trained to predict textual outputs with a multi-task approach where these tasks include language identification, multilingual speech transcription, and to-English speech translation (Figure 2.3). In the pre-training phase, the model is trained on a large dataset using the MLM approach. In the fine-tuning phase, the model is further trained on a smaller labelled dataset using supervised learning to improve its performance on specific tasks.

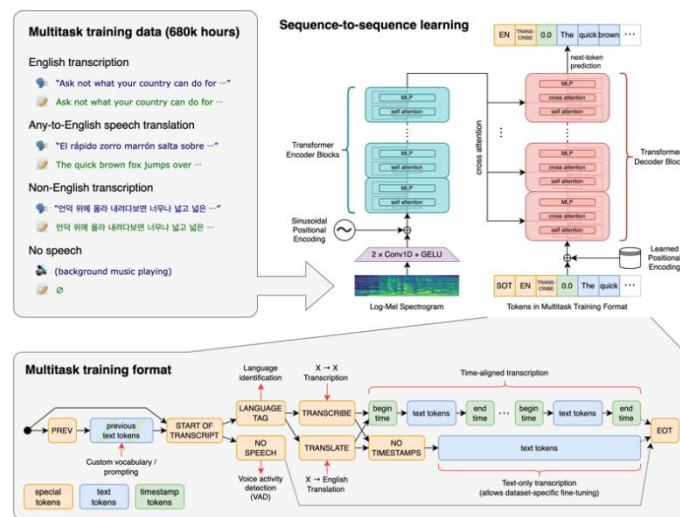


Figure 2.3. OpenAI's Whisper architecture

2.2. Neural Machine Translation (NMT)

Machine Translation (MT) refers to the automatic text or speech translation from one language to another using computational methods. Similar to the ASR evolution, MT transitioned from rule-based methodologies to transformer-based, deep neural network methodologies over the years. Early methods utilized linguistic rules crafted by linguists and dictionaries to translate text [Bar-Hillel, 1951]. The first automated machine translation tools were developed with the introduction of Statistical Machine Translation (SMT) approaches such as Moses [Koehn et al., 2007]. SMT models, especially Moses, utilize statistical alignment models to align words or phrases in source and target corpora and decoding strategies such as beam search or cube pruning.

The first breakthrough from the early statistical methods occurred with the introduction of deep learning techniques and the emergence of the Neural Machine Translation (NMT) field, which is the approach to machine translation that utilizes neural network architectures, particularly sequence-to-sequence models. In the first study that mentions NMT [Sutskever et al., 2014], the model proposed, called Seq2Seq, consists of an RNN encoder that processes the input sequence into a fixed-size vector, and an RNN decoder that generates the output sequence based on this vector. This architecture (Figure 2.4) paved the way for models that assign a sequence of data to a fixed size vector, which revolutionized the MT field.

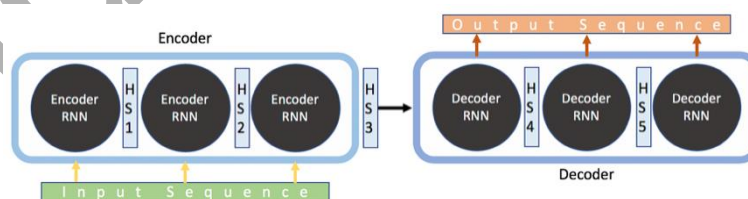


Figure 2.4. Seq2Seq architecture¹

The latest breakthrough occurred with the introduction of attention mechanisms, transformers and BERT and GPT architectures. Recent models such as BART [Lewis et al., 2019], T5 [Raffel et al., 2020], M2M-100 [Fan et al., 2021], and NLLB [Costa-jussà et al., 2022] all benefited from this breakthrough and they were considered to be the state-of-the-art at the time of their publication. BART (Bidirectional and Auto-Regressive Transformers) is an encoder-decoder model (Figure 2.5a), where the encoder is trained in a bidirectional manner using MLM objective with a denoising approach, meaning the input sequence is randomly masked, the encoder is trained to reconstruct the original sequence from the corrupted one

¹ <https://towardsdatascience.com/day-1-2-attention-seq2seq-models-65df3f49e263>

and the decoder is trained using an autoregressive approach, meaning the decoder generates one token at a time and the next token is generated based on the previously generated token. T5 (Text-To-Text Transfer Transformer) is an encoder-decoder structure (Figure 2.5b), where the model is pretrained on a large corpus using unsupervised learning with a denoising autoencoder objective and introduces task-specific prefixes during training and fine-tuning that enable the model to handle a plethora of down-stream tasks such as translation. M2M-100 (Many-to-Many Multilingual) is a many-to-many multilingual transformer-based model (Figure 2.5c) trained on a large multilingual corpus, which includes parallel data for translation tasks in a task-agnostic manner leveraging self-attention mechanisms to capture contextual relationships in the input sequence. NLLB (No Language Left Behind) is a multilingual transformer-based model, which introduces the Sparsely Gated Mixture of Experts approach to the network structure that is constructed based on another model called LASER3 (Language-Agnostic SEntence Representations) [Heffernan et al., 2022]. Sparsely Gated Mixture of Experts models activate a subset of model parameters as opposed to all parameters in traditional dense models.

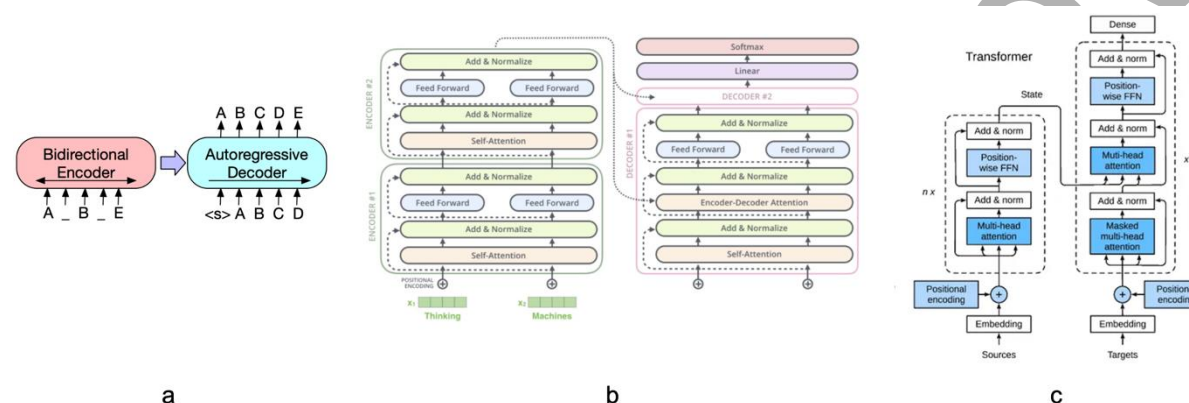


Figure 2.5. The architectures of a) BART, b) T5, and c) M2M-100

Various NMT studies focus on specific tasks for machine translation to tackle the challenges of the multilingual general NMT models. VOXReality focuses on three of these major challenges in NMT research: a) Context-aware Machine Translation, which refers to the integration of contextual information, such as surrounding context or terminology, into the process of machine translation to enhance the accuracy of translated text, especially to minimise the errors coming from ambiguities in the source text, b) Robust Machine Translation, which refers to the machine translation systems that can produce accurate and reliable translations across a variety of challenging conditions, including linguistic variations and noisy inputs, and c) Simultaneous Machine Translation, which refers to the machine translation systems that can translate in real-time as the source text is being generated, without waiting for the entire input to be available.

In Context-aware Machine Translation, the state-of-the-art architectures can be broadly categorized into single-encoder and multi-encoder approaches. In single-encoder architectures, the context sentences are concatenated with the current sentence and processed as a long sequence by a single encoder [Sun et al., 2020]. In multi-encoder architectures, the context sentences are processed by a separate encoder than the current sentence and may involve sharing parameters of encoders to reduce the number of parameters needed to be trained [Huo et al., 2020]. In multi-encoder approaches the features obtained from the encoders can be fused in multiple ways, such as by concatenation before passing the information to the decoder or by using cross-attention mechanisms.

In Robust Machine Translation, the methodologies can be categorized into two main approaches: 1) data augmentation and domain adaptation and 2) adversarial training. In data augmentation approaches, the model is exposed to a diverse set of linguistic patterns through

techniques such as back-translation, multi-domain training, and fine-tuning on in-domain and/or out-of-domain data [Xie et al., 2020]. In adversarial approaches, the model is trained to produce accurate results even in the presence of adversarial and noisy inputs using techniques such as regularization that penalize sensitivity to input variations or employing adversarial training to reduce the effect of domain shift [Miyato et al., 2016].

In Simultaneous Machine Translation, the methodologies employed can be categorized into fixed-wait, multiple-wait and adaptive-wait approaches. In fixed wait architectures, the models wait for k amount of tokens to be received to start the translation process [Ma et al., 2018]. In multiple-wait architectures, the model is trained to support a set of k values and can be utilized with one of these k values by using methods such as mixture-of-experts to wait for tokens to start translating [Zhang & Feng, 2021]. In the adaptive-wait architectures, the model is trained to decide when to translate and when to wait for more tokens by using methods such as monotonic multihead attention. The models trained with this approach also have the capability to adapt to different wait (k) values [Ma et al., 2019].

2.3. Visual Language Models (VL)

Visual language models refer to the neural network models that aim to jointly represent and understand cross-modal information in images and text. Early examples of visual language understanding approaches followed traditional computer vision methods and hand-crafted features and utilized separate traditional models for visual content and the natural language.

The introduction of Convolutional Neural Networks (CNNs), RNNs and large-scale datasets such as Common Objects in Context (COCO) [Lin et al., 2014] paved the way for data-driven approaches such as Show and Tell [Vinyals et al., 2015], which utilizes CNNs for image feature extraction and LSTMs for natural language generation (Figure 2.6).

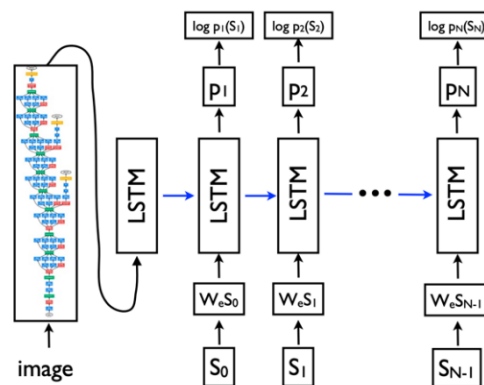


Figure 2.6. Show and Tell architecture

The latest breakthrough occurred with the introduction of attention mechanisms and transformer architectures. The most prominent models that utilize these methods are called LXMERT [Tan & Bansal, 2019], VisualBERT [Li et al., 2019] and UNITER [Chane et al., 2020]. These models are capable of a plethora of visual language tasks, such as image captioning and visual question answering, which are the tasks VOXReality is also investing in.

LXMERT (Language-visual Model for Efficient Representations of Transformers) is a visual language model that utilizes a transformer-based architecture (Figure 2.7a), a vision-answer bi-attention mechanism that focuses on relevant image regions and words in the textual context, and a positional encoding for the vision modality to capture the spatial relationships, and learns joint representations of both modalities during training for better understanding of the relationship between visual and textual context.

VisualBERT is a visual language model that extends the BERT architecture (Figure 2.7b) and it is pre-trained with an MLM objective and a sentence-image prediction objective using the COCO dataset. It also introduces a set of visual embeddings: a visual feature representation of the bounding region, a segment embedding to indicate that the representation is an image embedding, and a position embedding to align words and bounding regions.

UNITER (UNiversal Image-Text Representation) is a visual language model that employs a transformer-based architecture (Figure 2.7c) and is pre-trained using four objectives: MLM, Masked Region Modeling (MRM), where regions in an image are masked and the model's objective is to predict the masked region, Image-Text Matching, where the objective of the model is to predict if the given text and the image region is a match or not, and Word-Region Alignment, where the aim is to align words and image regions using a predetermined loss function instead of letting the model determine the alignment implicitly. In UNITER, the MLM objective is conditioned on the observation of the full text and the MRM objective is conditioned on the full image instead of the commonly used approach of applying joint random masking to both modalities.

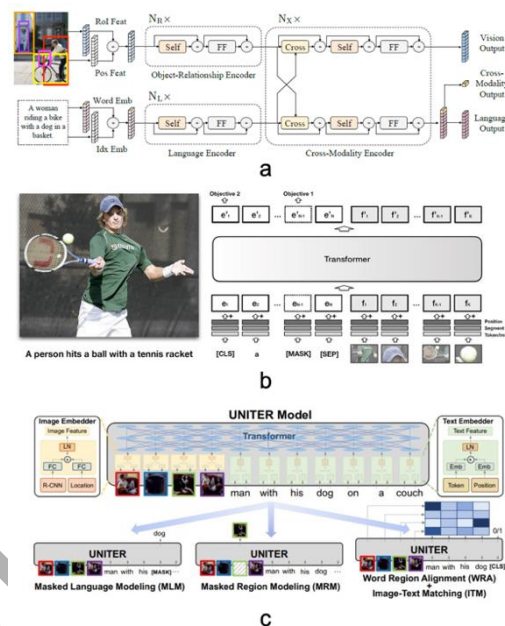


Figure 2.7. The architectures of a) LXMERT, b) VisualBERT, and c) UNITER

2.4. Conversation Agents (CA)

Conversation agents (CAs), commonly known as chatbots, are systems designed to simulate human-like conversations with users, often through text or voice-based interfaces. The history of CA began with the creation of ELIZA [Weizenbaum, 1966], a program that mimicked a Rogerian psychotherapist by using pattern matching and substitution methodology. This was followed by PARRY [Kochen, 1980], another milestone that simulated a person with paranoid schizophrenia.

The advancement of CAs received a considerable boost with the advent of deep learning technologies. Notably, the introduction of Seq2Seq [Sutskever et al., 2014] marked a major breakthrough in this evolution, significantly enhancing the capabilities and complexity of conversational agents. This model was primarily designed for machine translation but was quickly recognized for its potential in creating conversational agents. The significant advancements in CA technology took a big step forward with the introduction of the attention mechanism and transformer-based models [Vaswani et al., 2017]. Transformer models are

exceptionally efficient in handling sequential data and are particularly adept at capturing the intricacies of human language, making them ideal for creating highly sophisticated conversational agents.

Meanwhile, OpenAI's GPT (Generative Pre-trained Transformer) [Radford et al., 2018] series, including GPT-3 [Brown et al., 2020], utilized an architecture (Figure 2.8) that is characterized by its use of self-attention mechanisms, enabling the models to weigh the importance of each part of the input data in relation to the rest, a crucial factor in understanding and generating coherent, contextually relevant text. The combination of a transformer backbone with extensive pre-training and fine-tuning underpins the remarkable capabilities of GPT models in generating human-like, context-aware text.

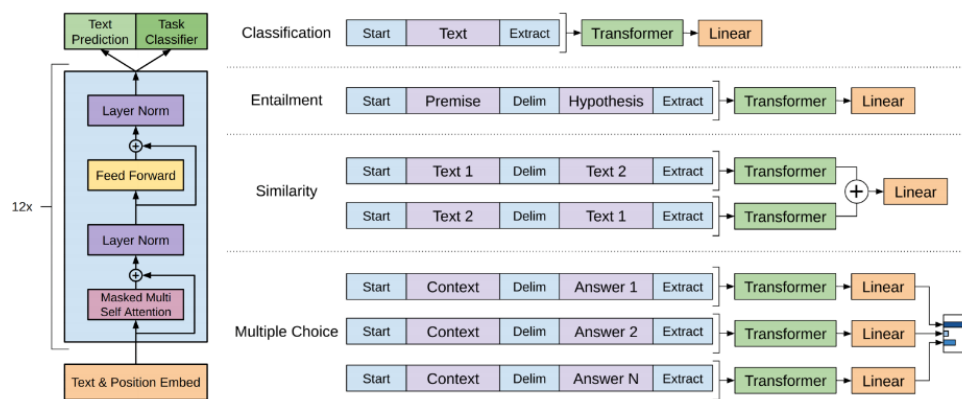


Figure 2.8. The architecture of GPT model [Radford et al., 2018]

Moreover, the Llama-2 [Touvron et al., 2023] model represents a significant advancement in the field with the utilization of the Ghost Attention (GAtT) method, which addresses the loss of context in multi-turn conversations. This method artificially concatenates the instruction to all user messages, leading to context-rich dialogues and improved attention. Furthermore, Llama-2 forms the backbone of instruction-tuned CA Llama2-Chat (Figure 2.9), which applies pretraining on the Llama 2 model using publicly available online sources and a further supervised fine-tuning. The model was then further refined through iterative processes using Reinforcement Learning with Human Feedback (RLHF) methodologies, including rejection sampling and Proximal Policy Optimization (PPO) [Schulman et al., 2017]. Crucially, throughout the RLHF stage, there was a parallel accumulation of iterative, reward modeling data, ensuring the reward models remained within distribution and aligned with the model enhancements. The Llama2-Chat models are available in three versions: 7 billion, 13 billion, and 70 billion parameters, each tailored to different levels of complexity and conversational requirements.

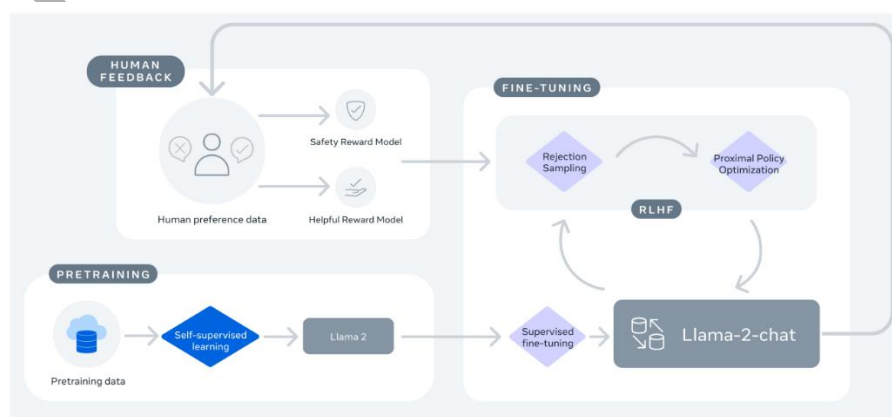


Figure 2.9. Training process of Llama 2-Chat

The conversation agents can be classified into two categories: task-oriented conversational agents, which focus on specific functions, and general conversational agents, designed for more dynamic and wide-ranging interactions. VOXReality focuses on the development of task-oriented conversational agents, specifically tailored to the requirements of the virtual conference and machine assembly training use cases. The agent in virtual conference use case is being developed with a focus on navigation assistance, providing detailed information about the conference program and facilitating interactions within the trade show, and the agent in machine assembly training use case focuses on the conversational training assistant. The task oriented CAs include three main components: Natural Language Understanding (NLU), Dialogue Management (DM), and Natural Language Generation (NLG).

1. Natural Language Understanding (NLU): NLU is the component responsible for extracting the meaning from user inputs in natural language. It serves as the bridge between the user's query and the agent's ability to comprehend and respond effectively.

2. Dialogue Management (DM): Following NLU's extraction of user intent and context, the extracted information is passed on to the DM, which is the component responsible for orchestrating the conversation and deciding how the agent should respond based on the user's input and the bot's task or purpose. It acts as the decision-maker, determining the appropriate course of action within the conversation flow. Dialogue Management is often split into two sub-components: Dialogue State Tracking and Dialogue Policy.

a) Dialogue State Tracking (DST): DST focuses on keeping track of the current state of the conversation by maintaining a record of all relevant information and context, including user queries, preferences, and any previous interactions. DST employs techniques such as slot filling to extract key pieces of information from user messages.

b) Dialogue Policy (DP): DP is responsible for deciding the agent's actions and responses based on the dialogue state using strategies and algorithms for selecting the most suitable response or action at each turn of the conversation. DP uses the information provided by DST to determine the next steps in the dialogue. DP ensures that the agent's responses are coherent, relevant, and aligned with the user's goals.

3. Natural Language Generation (NLG): NLG is responsible for converting structured data and decisions generated by the DM into human-readable and natural-sounding responses by generating text responses that are coherent, contextually relevant, and grammatically correct.

In the context of task-oriented CAs, it's essential to highlight that the utilization of state-of-the-art transformer-based models like T5, BERT and similar ones can be advantageous for all the components. T5 has reimaged NLP tasks by framing them uniformly as text-to-text problems, facilitating a more integrated and coherent approach to a variety of tasks, including translation, summarization, and question answering. T5 has undergone comprehensive pre-training on the "Colossal Clean Crawled Corpus" (C4)², which is instrumental in its ability to generate and understand human-like text. The model is available in several scalable variants, each tailored for different computational needs and performance requirements:

- **T5-Small:** Optimized for limited computational resources, ideal for mobile or edge computing.
- **T5-Base:** A balanced choice for general-purpose NLP tasks, offering a good compromise between efficiency and performance.

² <https://github.com/google-research/text-to-text-transfer-transformer#c4>

- **T5-Large:** Geared towards more complex NLP tasks, it delivers higher accuracy due to its larger size.
- **T5-3B and T5-11B:** These are the largest configurations, containing billions of parameters, and are typically used for high-end NLP research and applications with substantial computational resources.

However, DP often takes a different approach; instead of utilizing pre-trained models, DP commonly relies on Reinforcement Learning (RL) methods. PPO is a popular RL algorithm used for DP, which optimizes the policy of the agent by iteratively collecting data and improving the policy to maximize rewards. Another algorithm, Deep Deterministic Policy Gradients (DDPG) [Lillicrap et al., 2017], is used when the action space is continuous since the aim of the algorithm is to learn deterministic policies that map states to specific actions, making it suitable for precise control in task-oriented dialogues.

In recent advancements, the end-to-end approach for task-oriented dialogue agents, such as GPT-3 and Llama-2, has gained significant traction, offering a more streamlined and integrated methodology. This approach treats the entire dialogue process as a single holistic task, which includes understanding user intent, managing the dialogue, and generating responses. Unlike traditional models that handle these components separately, the end-to-end approach process and respond to user inputs directly. GODEL [Peng et al., 2022] is another end-to-end model that addresses the challenges of task-oriented dialogues.

2.5. Evaluation Metrics

There are three evaluation metrics widely adopted across the NLP field: BLEU [Papineni et al., 2002], ROUGE [Lin et al., 2004], METEOR [Banerjee et al., 2005], and chrF [Popović, 2015].

2.5.1. Bilingual Evaluation Understudy (BLEU)

BLEU (Bilingual Evaluation Understudy) is the most widely used metric for evaluating the output of machine translation systems with human-generated reference translations. It operates by comparing n-grams (contiguous sequences of n items, typically words) between the candidate (machine-generated) and reference (human-generated) translations. The precision of n-grams in the candidate translation is measured by how many of them overlap with n-grams in the reference translations. While this concept is sound for long translations, it has a tendency to over-rate short translations. In an effort to overcome this, Modified Precision is introduced, which considers the maximum number of times a particular n-gram appears in any single reference translation. This is done to prevent the metric from being overly optimistic about short and repetitive translations. Another measure to that direction is that it includes a brevity penalty if the length of the candidate translation is less than the length of the reference translations.

BLEU computes precision scores for different n-gram lengths (usually up to 4-grams), and then combines them using a weighted geometric mean. The weights are usually equal, but some variations may assign different weights. The final score is calculated by multiplying the combined precision by the brevity penalty. The brevity penalty is typically the exponential of 1 minus the ratio of the length of the reference translation to the length of the candidate translation. Its values range from 0 to 1, with 1 indicating a perfect match between the candidate and reference translations.

A notable limitation of BLEU is that it does not explicitly consider word order or capture semantic nuances and a high BLEU score doesn't necessarily mean the translation is of high

quality in all aspects. It doesn't capture semantic meaning or fluency and might favour translations that are close to the references but not necessarily correct or natural sounding.

2.5.2. Recall-Oriented Understudy for Gisting Evaluation (ROUGE)

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) is a metric commonly used for evaluating the quality of automatic summarization and machine translation outputs. It assesses the similarity between a generated summary and one or more reference summaries by measuring the overlap of n-grams (contiguous sequences of n items, usually words) and other units of text. There are several variants of ROUGE metrics, such as ROUGE-N, ROUGE-L, ROUGE-L_{sum}, ROUGE-W, ROUGE-S, ROUGE-SU and ROUGE-M. Each variant focuses on different aspects of similarity:

- ROUGE-N (Unigram, Bigram, etc.): Measures overlap of n-grams between the system summary and the reference summary. ROUGE-1 considers unigrams (single words), ROUGE-2 considers bigrams (two consecutive words), and so on.
- ROUGE-L (Longest Common Subsequence): Measures the longest common subsequence of words between the system and reference summaries. It is less sensitive to word order than n-gram-based metrics.
- ROUGE-L_{sum}: Similar to ROUGE-L but considers multiple occurrences of the same n-gram in the reference.
- ROUGE-W (Weighted Overlap): Considers the overlap of weighted words. It assigns higher weights to important words, giving more importance to content words.
- ROUGE-S (Skip-Bigram): Measures the overlap of pairs of words that have one or more words in between them. It helps capture some word order information.
- ROUGE-SU (Skip-Bigram with Unigram): An extension of ROUGE-S that considers unigrams as well. It combines unigrams and skip-bigrams to capture more information.
- ROUGE-M (ROUGE for Multigrams): Computes precision, recall, and F1 score for various lengths of n-grams.

ROUGE scores range from 0 to 1, where a higher score indicates better similarity between the generated and reference summaries. The basic calculations in all variants involve counting the number of overlapping units (n-grams, words, or sub-sequences) between the generated summary and the reference summary, based on the chosen metric (n-grams, longest common subsequence, etc.). The precision, recall, and F1 score are then computed based on these counts:

- Precision: Proportion of overlapping units in the generated summary with respect to the total number of units in the generated summary.
- Recall: Proportion of overlapping units in the generated summary with respect to the total number of units in the reference summary.
- F1 Score: Harmonic mean of precision and recall.

ROUGE, like all metrics described here, has innate limitations that need to be taken into consideration. Firstly, it is case-sensitive, so "cat" and "Cat" would be treated as different units. In order to alleviate this, case folding may need to be applied. Also, punctuation and stop-words (common words like "the," "and," etc.) are often excluded or given less importance. Additionally, ROUGE metrics may be sensitive to summary length. Extremely short or long summaries can affect the evaluation. It's important to ensure that the evaluation considers the length of the summaries appropriately. Finally, it should be noted that it is upon us to choose the correct variant of ROUGE, better suited to our case. For example, ROUGE-N is sensitive to word order, while ROUGE-L is more lenient in this regard. Also, ROUGE can be applied at



different levels, including word level, sentence level, or document level. The level chosen depends on the granularity of evaluation desired.

2.5.3. Metric for Evaluation of Translation with Explicit Ordering (METEOR)

METEOR (Metric for Evaluation of Translation with Explicit Ordering) is another metric used for evaluating the quality of machine translation outputs and was designed to address some of the limitations of BLEU and ROUGE. It incorporates explicit word matching, stemming, synonymy, and word order into its evaluation process. The evaluation process starts by aligning words in the candidate (system-generated) translation with words in the reference translation. It considers exact matches, stemming (reducing words to their root form), and synonymy (using WordNet) to find matches. Then, it computes precision, recall, and F1 score based on the number of matched unigrams, bigrams, and synonym pairs. These measures are calculated separately and then combined using a harmonic mean to get an overall score. METEOR incorporates a penalization mechanism for overgeneration (when the system generates more words than necessary) and undergeneration (when the system generates fewer words than necessary). These penalties are incorporated into the score to encourage more accurate translations. Another feature of it is that it employs a more flexible alignment approach compared to some other metrics. It allows for multiple reference words to be matched to a single candidate word and vice versa. This helps handle variations in expression.

METEOR provides scores for various components, such as unigram precision, unigram recall, precision for stemmed unigrams, recall for stemmed unigrams, synonym precision, synonym recall and penalty terms for overgeneration and undergeneration. It allows for different components to be weighted differently, based on their importance. This enables users to customize the metric according to their priorities. The final score is a combination of the individual component scores, with each component weighted according to user preferences. The flexibility of METEOR in handling stemming, synonymy, and various language-specific nuances makes it suitable for evaluating translations in a wide range of languages. It has shown good correlation with human judgment in various machine translation evaluations, making it a quite reliable metric for automatic assessment. Of course, like any automatic evaluation metric, it might not fully capture the intricacies of human judgment. There are cases where a translation may be technically correct but might not align perfectly with reference translations.

While METEOR is a comprehensive metric for the evaluation of machine translation outputs, it also has some limitations and considerations to be aware of. It can be sensitive to preprocessing choices, such as tokenization and stemming. Small variations in preprocessing can lead to differences in the evaluation results. Also, it relies on WordNet for synonym matching. While WordNet is a valuable resource, its coverage and accuracy may vary across different languages and domains. In some cases, it may not capture all relevant synonyms. Most notably though, METEOR places a strong emphasis on exact matching. While it considers stemming and synonyms, deviations from the reference translations still result in lower scores even if the meaning is preserved, mainly rooted to the fact that it takes word order into account. Lastly, in practice tuning METEOR for optimal performance can be challenging due to its multiple parameters and components. The need for careful parameter tuning makes it less straightforward to use.

2.5.4. Character-level F-score (chrF)

Character-level F-score (chrF) is an n-gram based F-score that measures translation quality by assessing the overlap between the hypothesis and reference at the character level. This fine-grained evaluation is particularly useful for languages with rich morphology, where token- or word-level metrics may fall short. It is useful because of its sensitivity to small errors and

morphological variations, but it is less sensitive to syntactic and contextual coherence. This can make it less suitable for evaluating overall sentence structure.

3.VOXReality Natural Language Processing Models

3.1. Automatic Speech Recognition (ASR)

In VOXReality, we provide a set of ASR models and corresponding functionalities: 1) improved Whisper, 2) Speech translation for non-native speech, and 3) streaming ASR for utilizing real-time speech in simultaneous translation.

3.1.1. Improving the Speech Recognition Models

We found that ASR models underperform on Greek language, and we worked on improving their performance on Greek through finetuning and using Adapter modules. Whisper [Radford et al., 2023] is a Multilingual Multi-task model that can perform ASR, translation into English and Language Identification. We integrate Whisper into our pipeline because of its robustness to speech variations that comes with training on large amounts of weakly supervised data. Although Whisper performs really well on high resource languages, it underperforms on low or medium resource languages like Greek, but the representations learned during the Multi-task pretraining can still be useful for finetuning a better model for Greek.

The Massively Multilingual Speech (MMS) Project is a recent effort from Meta to scale Speech models to 1000+ languages [Pratap et al., 2023]. The effort has resulted in pretrained speech representation models, ASR, Language Identification and Speech Synthesis models. The pretrained model is a Wav2Vec 2.0 based model [Baevski et al., 2020], which they finetune on downstream tasks. The ASR finetuning was done on a multilingual dataset of 1107 languages. This has led to MMS achieving better results than Whisper on languages that both models are trained on. We evaluated this model on the consortium languages, and we improved its results on Greek similarly to Whisper.

We use adapters [Rebuffi et al., 2017] [Houlsby et al., 2019] for efficient finetuning of the Whisper model. Whisper architecture is a normal transformer-based encoder-decoder architecture, therefore, we plugin adapter modules after the feed forward network of both the encoder and the decoder. We insert one adapter module per layer. The adapter architecture [Bapna et al., 2019] we use is composed of layer normalization, a linear layer that lowers the dimension of the input (down projection), and a linear layer that brings the input back to its original dimension (up projection). The two linear layers are separated by a Rectified Linear Unit (RELU). The architecture of the model that utilizes adapter modules is illustrated in Figure 3.1. We compare Whisper-Finetuned and Whisper-Adapters to MMS, which by default relies on language adapters to limit interference between languages.

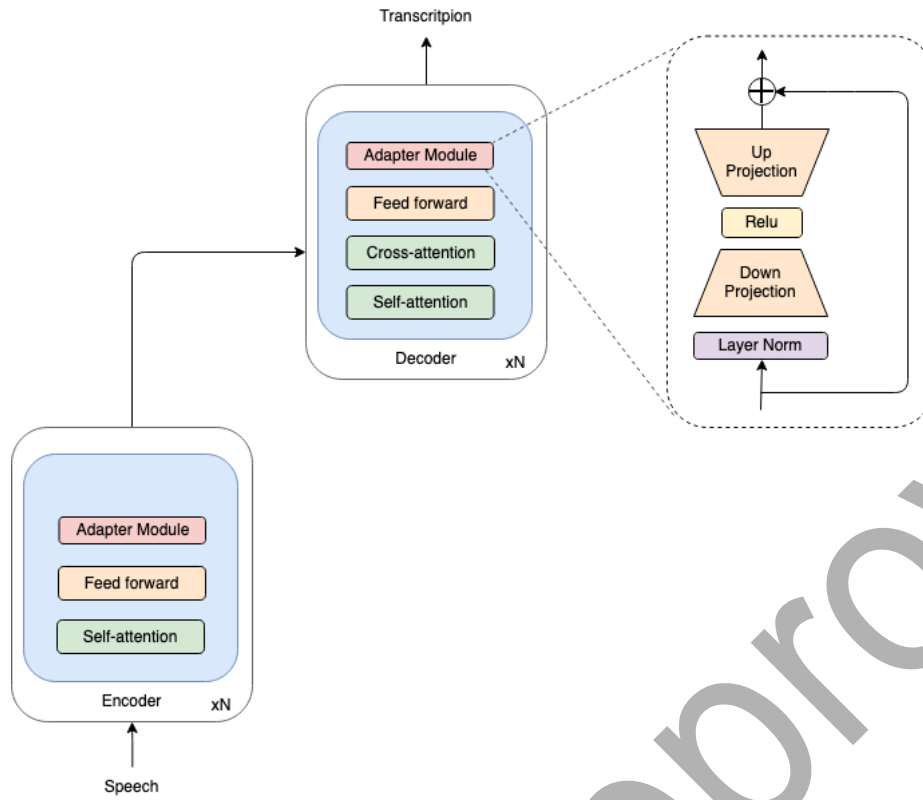


Figure 3.1. Whisper Model with Adapter Modules

Results

Whisper comes in multiple variants that can suit different quality and size requirements. As shown in Table 3.1, these models differ in terms of their number of parameters, and thus their performance is also different. The number of parameters can have implications on the latency and the amount of memory required to run the model for inference. We choose Whisper-small as a good trade-off between quality and size of the model.

Table 3.1. Whisper Size Variants

Model	Layers	Width	Heads	Parameters
Tiny	4	384	6	39M
Base	6	512	8	74M
Small	12	768	12	244M
Medium	24	1024	16	769M
Large	32	1280	20	1550M

For training and evaluation, we use the Common Voice dataset Version 16.1 [Ardila et al., 2019], which is a multilingual corpus for Speech Recognition and Language Identification. The dataset is collected through crowdsourcing, where contributors can record their own voice reading a sentence in the Common Voice website or app, hence, the dataset contains speech from different speakers and with Different levels of noise, and it can be used to evaluate or train robust speech recognition models. We use the Greek Language partition of the dataset, which contains 1.91k rows for training and 1.67k for validation. We train the model on the training set and use the validation set to report the results in Table 3.3. For Both Whisper-small and MMS, we report in Table 3.2 the hyperparameters of the finetuning.

Table 3.2. Hyperparameters of Whisper Finetuning

Hyper-parameter	Value
Number of train steps	5000
Warmup steps	500
Dropout	0.

Optimizer	Adam
Learning Rate	1e-5
LR Scheduler	Linear
Batch Size	32

We report the Word Error Rate (WER) metric results in Table 3.3. The original model (Whisper-small) achieves decent performance on all languages except Greek, where the WER is between 13.17 and 22.9 points higher than the other reported languages, a fact which motivated our work on finetuning Whisper for Greek. After finetuning, the performance improves by ~10 points, but it worsened on the other languages due to information forgetting. We found that finetuning is still better than using adapter, but the adapters guarantee a similar performance on the other languages because the model is frozen and it allows us to omit the use of the adapters during inference. MMS given that it is a larger model SOTA model achieves significantly better results on all languages. Since MMS already uses adapters, we only report the results of fine-tuning the adapters on Greek, which we found to not improve the results compared to MMS. This can be attributed to the fact that MMS already achieves good results on this dataset, therefore fine-tuning the adapters was not effective.

Table 3.3. WER Results of Finetuning vs. Adapter finetuning

Model	Greek	Dutch	German	Italian	Spanish	English
Whisper-small	43.4	24.73	24.5	30.23	20.51	25.27
Whisper-small-finetuned	33.83	26.74	25.5	35.67	22.02	27.32
Whisper-small-adapters	41.15	24.73	24.5	30.23	20.51	25.27
MMS	20.46	6.57	10.97	11.08	8.55	19.82
MMS-Adapters	20.49	6.57	10.97	11.08	8.55	19.82

Deployment

In the VOXReality pipeline the audio processing is handled through a set of endpoints in a REST API developed using FastAPI³ framework. The speech related endpoints are designed to generate transcriptions and translations from audio files in “wav” format and utilizes the Whisper-small-adapters model for ASR. There are three endpoints: speech recognition (transcription), speech translation and context-aware speech translation. The details of the API calls are presented in Appendix I.

In speech recognition, the function generates a textual response to a given audio file in the “wav” format. The endpoint does not require a source language to be specified since it is automatically recognized by the model itself. Figure 3.2 presents the FastAPI-powered endpoint that can be requested using the handle “transcribe_audio_files” from the deployed server.

The screenshot displays the Swagger UI for the endpoint `POST /transcribe_audio_files`. The interface includes a 'Parameters' section with the following fields:

- Name:** `source_language` (required), with a description '(query)'. The value is set to 'en' in a dropdown menu.
- Request body:** (required), set to 'multipart/form-data' in a dropdown menu.
- audio_files:** (required), an array of strings. There is an 'Add string item' button next to it.

Buttons for 'Cancel' and 'Reset' are located at the top right of the parameters section.

Figure 3.2. The speech transcription endpoint

³ <https://fastapi.tiangolo.com/>

In speech translation, the function generates a textual response to a given audio file in the “wav” format. The endpoint does not require a source language to be specified since it is automatically recognized by the ASR model itself. This endpoint utilizes the speech recognition model to transcribe the audio file and automatically utilizes the latest text translation model to translate the transcribed text to the given target language. Figure 3.3 presents the FastAPI powered endpoint that can be requested using the handle “translate_audio_files” from the deployed server.

The screenshot displays the configuration for the POST endpoint `/translate_audio_files`. The parameters section includes:

- `source_language` (query, required): dropdown menu set to `en`.
- `target_language` (query, required): dropdown menu set to `en`.
- `return_transcription` (query, boolean): dropdown menu set to `false`.

The request body section is set to `multipart/form-data` and includes:

- `audio_files` (array, required): a button labeled "Add string item".

Figure 3.3. The speech translation endpoint

In context-aware speech translation, the function generates a textual response to a given audio file in the “wav” format. The endpoint does not require a source language to be specified since it is automatically recognized by the ASR model itself. This endpoint utilizes the speech recognition model to transcribe the audio file and automatically utilizes the latest context-aware text translation model to translate the transcribed text to the given target language using the given contextual text. Figure 3.4 presents the FastAPI-powered endpoint that can be requested using the handle “contextual_translate_audio_files” from the deployed server.

The screenshot displays the configuration for the POST endpoint `/contextual_translate_audio_files`. The parameters section includes:

- `source_language` (query, required): dropdown menu set to `en`.
- `target_language` (query, required): dropdown menu set to `en`.
- `return_transcription` (query, boolean): dropdown menu set to `false`.

The request body section is set to `multipart/form-data` and includes:

- `audio_files` (array, required): a button labeled "Add string item".
- `context` (required): a text input field with the placeholder text `context`.

Figure 3.4. The context-aware speech translation endpoint

3.1.2. Speech Translation for Non-native Speech

In this section, we address the objective of building a direct speech translation model that is robust to non-native speech. For this end, we focus on errors that are common in non-native speaker language, such as preposition, noun number and article errors [Napoles et al., 2016]. Where we synthetically introduce them into the dataset in order to train as well as evaluate the models. We run our experiments on the CoVoST2 dataset (Common Voice Speech-To-



Text 2), which is commonly used for speech translation research [Wang et al., 2020]. It is a multilingual dataset that covers translations from 21 languages into English and translations from English to 15 languages. Furthermore, CoVoST2 is based on Common Voice project, which is a speech recognition crowd-sourcing project with speakers from diverse backgrounds [Ardila et al., 2019]. This makes CoVoST2 ideal for training and evaluating robustness to non-native speech translation. However, acquiring speech data from non-native speakers is not straightforward. Therefore, our research aims to leverage synthetic text data to build models that are robust to non-native speech phenomena such as preposition errors. Figure 3.5 illustrates the details of our architecture, which is composed mainly of a speech encoder, translation encoder, and a translation decoder. The speech encoder is initialized from HuBERT, a pre-trained self-supervised speech representation model. The translation encoder and decoder form a conventional transformer model that we pre-train on transcription-translation pairs. During training, we train the full model for speech translation while leveraging transcriptions for bridging the modality gap between speech and text, which we achieve with the Mixup strategy [Fang et al., 2022]. Furthermore, to train models to be robust to text level noise, we introduce the synthetic transcription noise into the Mixup. During inference, the model is used end-to-end and only require speech input while the Mixup and transcription components are dropped.

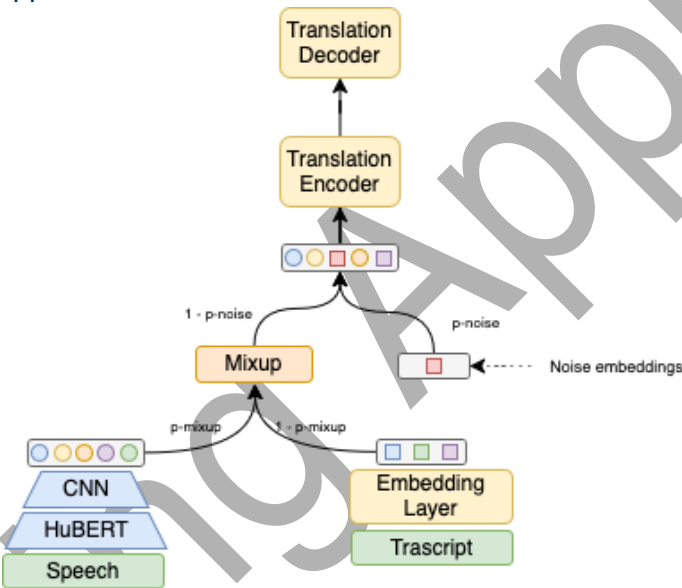


Figure 3.5: Robust Speech Translation Using Text Noise

Results

We train and evaluate all the models on En-De subset of the dataset, which contains 364 hours of speech for training, 26 hours for evaluation, and 25 hours for testing, which we use to report the results in Table 3.4 (CoVoST2-test column). To evaluate model robustness to the synthetic errors that we introduce in transcriptions, we use a Text-to-Speech (TTS) model to transform the noisy transcriptions of the test set into speech and produce the Noisy-TTS subset. For comparison reasons, we also apply TTS to the clean original data, which we refer to as Clean-TTS. Table 3.4 presents the BLEU score of our experiments. In Clean-TTS column, we can see that models perform better on TTS data than original CoVoST2-test data, which can be explained by the fact that CoVoST2 contains non-native accents and pronunciations that make the task more challenging than clean TTS inputs. When evaluating on the Noisy-TTS data, we notice a degradation of performance especially when comparing to Clean-TTS test set. Although all the models perform similarly on Clean-TTS, when evaluating on Noisy-TTS, our robust approach, which uses the Mixup strategy with synthetic text noise, achieves higher performance.

Table 3.4. Robust Speech Translation Results

Model	CoVoST2-test	Clean-TTS	Noisy-TTS
-------	--------------	-----------	-----------



HuBERT-Transformer	20.79	24.12	20.55
HuBERT-Mixup	20.28	24.22	21.01
HuBERT-Robust (Ours)	20.32	24.27	21.26

3.1.3. Streaming Speech Recognition and Translation

We employ two approaches for Streaming Speech Recognition task, one is based on chunking audio streams and feeding them into Whisper based on a fixed chunk size, and the other is an adaptive Streaming Speech Recognition and Translation model that decides the chunk adaptively depending on the context.

Chunk-based Streaming Approach

We deploy a chunking-based API that uses the Whisper model we fine-tuned previously (Sec 3.1.1). The conceptual flow of this approach is illustrated in Figure 3.6. The API continuously receives audio and chunks it based on a predefined chunk size parameter and *chunk length* that can be specified during the request. This has been set to a 2 second default as illustrated in Figure 3.6. To avoid chunking words in the input speech, we apply Voice Activity Detection (VAD) to make sure that the chunk ends with silence.

Voice Activity Detection is a technique for identifying segments in audio, where there is speech. This is important to avoid executing a transcription model on silence or on background noise. We use the VAD model from Pyannote [Bredin & Laurent, 2021]. The model architecture is composed of convolutional layers, LSTM recurrent layers and fully connected layers. Finally, a sigmoid activation function is used to predict voice activity.

If the chunk does not end with a silence, the chunk size is increased with the increment parameter, *accumulation chunk length*, and the process is repeated. Furthermore, if the accumulated chunk size surpasses a specified length in seconds, *maximum accumulation length*, the chunk is automatically sent to the ASR model, which returns the transcription of that chunk. The default value for this threshold is 5 seconds and it can be adjusted during the request.

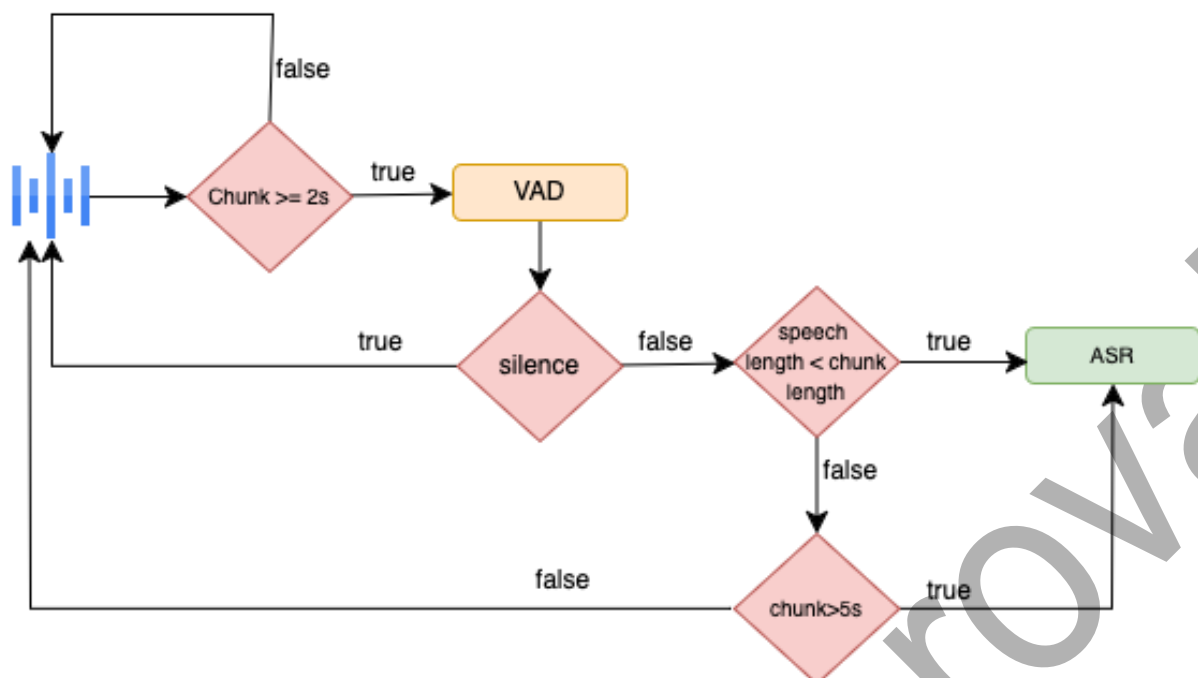


Figure 3.6: Chunking Based Streaming Speech Recognition

Adaptive Streaming Approach

The chunking-based strategy forcefully chunks the audio, which can negatively affect the quality of the transcriptions. Whisper is not specifically trained to handle speech segments of short length, which can lead to degradation in performance due to the mismatch between training and inference input. Furthermore, there are instances where the speech might be chunked in an instance where more context is required, and this is not taken into consideration in the chunking-based approach. In order to tackle this issue, we utilize Seamless Streaming [Communication et al., 2023], a state-of-the-art streaming Speech Recognition and Translation model that dynamically decides when to generate outcomes based on the input speech. The model performs feature extraction using Mel-Filterbanks extractor, which transforms audio features from frequency domain to Mel scale that better mimics how humans perceive sound. The features are passed to the speech encoder, which is a conformer model [Gulati et al., 2020] that combines convolutional layers to capture local features and transformer layers to capture global interactions. The streaming audio is transformed into hidden representation with the previous steps and passed to the simultaneous speech decoder which uses Efficient Monotonic Multihead Attention (EMMA) [Ma et al., 2023] to decide whether to wait for more speech or predict the next token.

The model is trained to generate transcription or translation with a low latency and good quality. Similar to Whisper, Seamless Streaming is a multilingual model and supports 101 source languages for speech input and 96 target languages in text output, which include all the languages represented in the consortium. The end-to-end nature of Seamless Streaming makes it a better solution for streaming speech translation since it does not require an extra step of translation, which can lead to error propagation and impacting the latency as well.

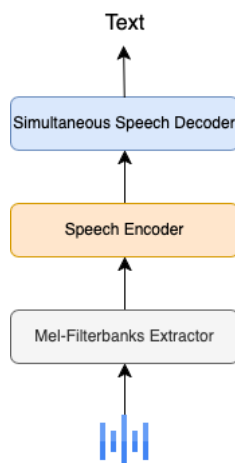


Figure 3.7: Dynamic Speech Recognition and Translation

Deployment

We deploy the streaming model using FastAPI and WebSockets. The deployment offers an interface for initializing the API and an endpoint to execute continuous requests. The interface is presented in Figure 3.8.

The interface is titled "Live Audio Transcription". It features a configuration section with four input fields for time parameters: "Chunk Length (s)" (set to 3), "Accumulation Chunk Length (s)" (set to 1), "Maximum Accumulation length (s)" (set to 5), and "Silence at the End of Chunk (s)" (set to 0.1). Below these are two dropdown menus for "Model" (set to "Whisper Small") and "Language" (set to "Multilingual"), followed by a blue "Connect" button. Below the configuration section are two blue buttons: "Start Streaming" and "Stop Streaming". A large white rectangular area is provided for the transcription output. At the bottom, the status is shown as "WebSocket: Not Connected" and "Last Processing Time: Undefined".

Figure 3.8: Streaming API configuration interface

The API accepts 6 parameters:

- **Chunk Length (s):** The size of the chunk that needs to be reached before sending the audio for processing.
- **Silence at the End of Chunk (s):** The silence length required for the VAD model.

- **Accumulation Chunk Length (s):** The length of the chunk to be added if the speech doesn't end before chunk length.
- **Maximum Accumulation Length (s):** The threshold for the maximum chunk length.
- **Model:** The translation model used
- **Language:** The language of the model (multimodal, if there are multiple languages are expected from the audio streams)

The API is initialized with the parameters using the interface and the API starts to accept streams with the configuration set.

3.2. Neural Machine Translation (NMT)

In VOXReality, we introduce three major categories of machine translation models: a) Context-aware Machine Translation, which includes surrounding sentences and/or terminology as context, b) Robust Machine Translation, which utilizes mechanisms to handle translation of noisy inputs, and c) Simultaneous Machine Translation, which translates source text in real-time as the input is being generated.

3.2.1. Context-aware Machine Translation

In VOXReality, we utilize multi-encoder architectures, where all the encoder parameters are shared. This allows caching the hidden representation of the current sentence and reusing it as the hidden representation of the context when translating subsequent sentences, and we refer to this architecture as caching. Furthermore, several techniques have been proposed to shorten the sequence of tokens where the tokens are combined in shortening modules that are added between a specified number of encoder layers. This can lead to the reduction of the computational and memory requirements in the subsequent layers; therefore, a compressed representation of the previously seen sentences should be enough to use as a context.

Consequently, we focused on the application of Sequence Shortening to Context-aware Machine Translation and we introduce two new shortening techniques (Latent Grouping and Latent Selecting), where the network can learn how to group or select tokens to form a shortened sequence. Shortening can be done by average/max pooling of the hidden representation of the tokens [Subramanian et al., 2020] or linear pooling of the concatenated representation of the tokens of the original sequence [Nawrot et al., 2021]. In the pooling-based shortening, the sequence is divided into non-overlapping groups of K neighbouring tokens.

In Latent Grouping, each token is categorized into a group by the feed-forward network (FFN) with the number of outputs equals to the number of groups (K). We obtain the categorization for the i^{th} token to K^{th} group by applying the SoftMax function to the outputs in the dimension of the groups. Subsequently, the groups are constructed as the sum of the hidden representations. The network learns how to softassign each token to the groups. A group representation is computed using the weighted average of tokens, which makes backpropagations into the categorizing network possible. Finally, the attention module is applied on the group representations. Latent Selecting enables the groups to select tokens to aggregate rather than to assign each token to a group and allows the model to ignore tokens entirely rather than to assign them to at least one group. Although Latent Selecting can be achieved by maintaining a categorizing feed-forward network for each group, we utilize the same network as described for Latent Grouping but apply SoftMax in the sequence dimension instead of the token dimension. Figure 3.9 illustrates the model utilized for both latent grouping and latent selecting where the number of groups (K) is set to three.

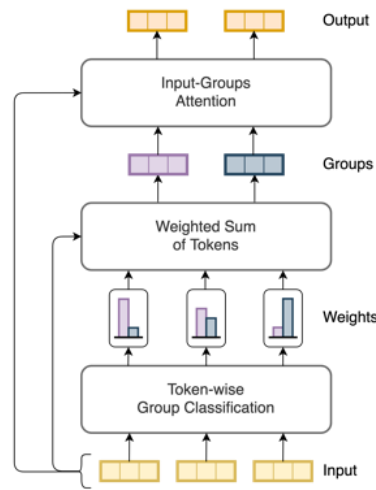


Figure 3.9. The architecture of latent grouping and latent selecting

The architecture we use, illustrated in Figure 3.10, is based on caching the hidden representations produced by the encoder, where the representations of the tokens of the current sentence are stored and can be reused as context when the subsequent sentences are translated. Although this architecture uses only a single encoder, it is different from the single-encoder models because the current sentence and the context sentences are processed separately. While in the standard caching architecture the hidden representation of all the tokens is stored, we introduce a sequence shortening module directly after the encoder, which returns the compressed hidden representation usually containing fewer tokens than the original sequence. The context is integrated in the decoder in the multi-encoder manner by using a separate cross-attention module for context tokens. To allow the decoder to distinguish between context sentences we employ learned segment embeddings [Devlin et al., 2018] and learned positional encoding for the shortened tokens inside context sentences.

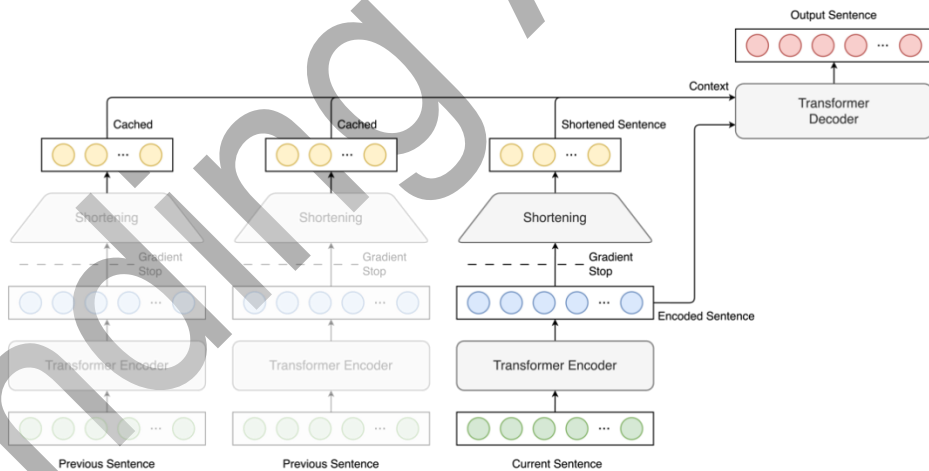


Figure 3.10. The architecture of VOXReality context-aware MT model

During training, caching is not used, meaning that the model receives tokenized context sentences and processes them using the same encoder, where the weights of the encoder receive the backpropagation gradient from multiple sources - the current sentence and each of the context sentences. In order to eliminate difficulties in training that would arise from this, we block the gradient after the encoder and before shortening where applicable. This is done by allowing the gradient information to flow for a specified number of context sentences, after which the gradient is blocked.

Moreover, to allow a machine translation model to use terminology, we adapted the work of [Jon et al., 2021]. The main differences of our work are threefold:

- 1) we apply this technique to the multi-lingual setting,
- 2) we train the model capable of context-aware translation by utilizing the previous sentence of the source language,
- 3) we finetune the pre-trained model (NLLB-600M [Costa-jussà et al., 2022]) instead of training from a random initialization.

The model requires the list of terminology constraints with each item in the form of phrases in all consortium languages (English, German, Italian, Dutch, Spanish and Greek), although the phrase is not required to be specified in all languages. Each pair formed from the two phrases in different languages is treated as a constraint in both directions. We pre-process the terminology list by lemmatizing the phrases using spaCy [Honnibal et al., 2020].

During inference, the source sentence is lemmatized, and the search is conducted for any of the phrases in the terminology list from the same language. If a phrase is found, the corresponding phrase in the target language is appended to the source sentence delimited with the separator token “</s>”. In case when no phrase is found, the token is still added at the end of the source sentence (apart from the usual end-of-sequence token). Additionally, we allow the model to accept the previous sentence from the source language as context. To this end, we prepend the context to the input of the network delimited with the separator token. NLLB-600M model uses the same tokens as a separator and end-of-sequence token and we do not change that in our finetuned model. The architecture is illustrated in Figure 3.11.

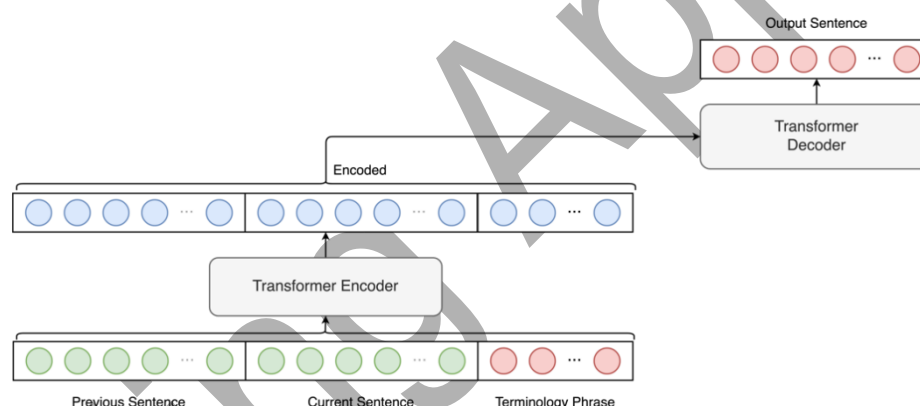


Figure 3.11. The terminology-constrained architecture

Results

It has been argued that the sentence-level metrics do not capture well the ability of the model to use the context. For this reason, to measure how well the model utilizes the contextual information, it is common to use a contrastive dataset such as Contrastive evaluation of Pronoun translation (ContraPro) [Müller et al., 2018] and Large Contrastive Pronoun Test (LCPT) [Lopes et al., 2020]. They are very similar to each other and both target pronoun disambiguation task. They differ in the language pair they use – ContraPro is English-German and LCPT is English-French dataset. The example from the datasets is in the form of context sentences + source sentence in the input and several (usually three) target sentences that differ only in the pronoun. Only one of the target sentences is correct and the others are erroneous. Here is an example from EN-DE dataset:

Context:	<i>“What’s your plan?”</i>
Source:	<i>“I forgot to confide it to you.”</i>
Correct Target:	<i>“Ich vergaß, es euch zu vertraun.”</i>
Incorrect Targets:	<i>“Ich vergaß, sie euch zu vertraun.”</i>

The model outputs the probabilities of each target sentence (as a factorized probability calculated by multiplying the probability of each token). When the model scores the correct target as the most probable, model counts it as a correct prediction. Otherwise, it is incorrect.

We trained the following models:

- **Sentence-level Transformer** - where context sentences are ignored
- **Single-encoder Transformer** - where context sentences are prepended to the current sentence and processed by the encoder
- **Multi-encoder Transformer** - with the separate encoder (without weights-sharing) used to encode the context sentences, where the context and the current sentence are concatenated in the decoder (Our experiments revealed that this integration yields better results than with the separate context-attention module)
- **Caching Tokens** - where the encoder representations of the context sentences are stored directly
- **Caching Sentences** - where the representations of the context sentences are averaged and stored
- **Shortening - Mean Pooling** - Sequence shortening with Mean Pooling applied to the outputs of the encoder
- **Shortening - Max Pooling** - shortening with Max Pooling
- **Shortening - Linear Pooling** - shortening with Linear Pooling
- **Shortening - Grouping** - shortening with Latent Grouping (our proposed method)
- **Shortening - Selecting** - shortening with Latent Selecting (our proposed method)

To test the caching and sequence shortening models we trained each model with the same hyper-parameters that are presented in Table 3.5. We used the English to German and English to French directions of the IWSLT 2017 [Cettolo et al., 2017] document-level dataset that is based on the subtitles of TED Talks. We trained each model on both language pairs with varied context sizes (in terms of the number of previous sentences) separately. We measured BLEU [Papineni et al., 2002] using the sacreBleu library, and the accuracy on the contrastive datasets (ContraPro for EN-DE and LCPT for EN-FR).

Table 3.5. The hyper-parameters of all models trained on IWSLT 2017 dataset

Hyper-parameter	Value
Encoder Layers	6
Decoder Layers	6
Attention Heads	8
Embed Dim	512
FFN Embed Dim	2048
Dropout	0.3
Optimizer	Adam
Learning Rate (LR)	5e-4
LR Scheduler	Inverse Sqrt
Batch Max Tokens	4096

The results for EN-DE are presented in Table 3.6. The BLEU score of the context-aware models is generally similar to or slightly higher than the sentence-level Transformer. BLEU does not correlate well with the contrastive accuracy, which is strictly higher for all context-

aware models. This confirms that sentence-level metrics do not reflect the context usage of the models. The highest contrastive dataset accuracy was achieved by the Grouping Shortening model for the context size of one, the Max Pooling Shortening model for the context size of two, and the Selecting Shortening model for the context size of three. The highest accuracy averaged over the tested context sizes was reached by the model employing Latent Grouping, followed by the Latent Selecting model. Caching Tokens architecture exhibits comparable BLEU scores to the Single- and Multi-encoder architectures while achieving higher accuracy on the contrastive dataset.

Table 3.6. Results of the models trained on the IWSLT 2017 with EN-DE dataset

Model	BLEU	Accuracy				
Sentence-level	28.11	43.67%				
Model	Context Size: 1		Context Size: 2		Context Size: 3	
	BLEU	Accuracy	BLEU	Accuracy	BLEU	Accuracy
Single-encoder	28.31	47.42%	27.95	48.18%	27.88	48.88%
Multi-encoder	28.67	44.93%	28.50	46.65%	28.26	45.00%
Caching Tokens	28.35	54.06%	28.50	54.13%	29.08	51.23%
Caching Sentence	28.38	45.72%	26.73	45.26%	26.70	44.91%
Shortening – Max Pooling	27.62	51.67%	27.88	55.08%	28.26	50.89%
Shortening – Avg Pooling	28.09	53.37%	27.85	54.81%	28.38	50.54%
Shortening – Linear Pooling	27.62	52.71%	28.03	52.13%	28.18	51.27%
Shortening – Grouping	28.21	56.98%	28.70	54.51%	28.49	51.16%
Shortening – Selecting	28.15	54.48%	28.55	54.21%	28.01	51.95%

The results of the models on the EN-FR datasets can be seen in Table 3.7. The BLEU scores of all models are comparable (apart from the Caching Sentence architecture). Latent Grouping achieved the highest accuracy on the contrastive dataset for the context size of one, and Latent Selecting and Single-encoder architectures for the context sizes of one and three, respectively. In general, Caching Tokens and Shortening models achieved higher accuracies than the Single- and Multi-encoder architectures (with the exception of the Single-encoder on English to French translation with a context size of three).

Table 3.7. Results of the models trained on the IWSLT 2017 with EN-FR dataset

Model	BLEU	Accuracy				
Sentence-level	37.64	75.92%				
Model	Context Size: 1		Context Size: 2		Context Size: 3	
	BLEU	Accuracy	BLEU	Accuracy	BLEU	Accuracy
Single-encoder	37.25	77.27%	37.18	78.98%	37.12	80.87%
Multi-encoder	37.44	75.72%	37.12	77.23%	37.34	75.76%
Caching Tokens	36.88	79.67%	37.29	80.14%	37.73	79.90%
Caching Sentence	36.50	77.33%	34.21	76.25%	34.78	75.71%
Shortening – Max Pooling	37.48	79.51%	36.72	80.59%	37.85	79.71%
Shortening – Avg Pooling	37.13	77.75%	37.12	80.16%	38.18	80.41%
Shortening – Linear Pooling	37.02	80.47%	37.12	79.37%	37.42	79.64%
Shortening – Grouping	37.05	79.91%	37.98	81.13%	37.18	79.54%
Shortening – Selecting	37.38	80.89%	37.83	80.32%	37.81	80.09%

We train the terminology-constrained model on the OpenSubtitles2018 dataset [Lison et al., 2018] for all 30 language pairs formed from the consortium languages in both directions. For each language, we use 50000 examples with the previous sentence as context and 50000 without context. Following [Jon et al., 2021] we augment the dataset with artificial terminology by selecting a random span in the lemmatized target sentence and using it as a target terminology constraint. This augmentation is done to the examples in both sentence-level and

context-aware datasets with the probability of 0.3. The word-length of the terminology item is drawn from the geometric distribution with $p=0.85$ with the maximum of three. The final training dataset is formed by interleaving the examples from the sentence-level and context-aware datasets as well as the language pair-specific datasets. Table 3.8 presents the hyper-parameters.

Table 3.8. The hyper-parameters of the terminology-constrained model

Hyper-parameter	Value
Encoder Layers	12
Decoder Layers	12
Attention Heads	16
Embed Dim	1024
FFN Embed Dim	4096
Dropout	0.1
Optimizer	Adam
Learning Rate	5e-5
LR Scheduler	Linear
Batch Size	12

We evaluated our terminology-constrained model and the baseline (NLLB-600M) on the test subset of the OpenSubtitles2018 dataset for all consortium language pairs. We used our model in four modes of operation:

- No context, no terminology (Figure 3.12 - Ours)
- Context, no terminology (Figure 3.12 – Ours with context)
- No context, terminology (Figure 3.12 – Ours with terminology)
- Context, terminology (Figure 3.12 – Ours with context and terminology)

To obtain the terminology for our model, we used the same procedure as used for training with the same probabilities. The results in terms of BLEU are illustrated in Figure 3.12. Our model outperforms the baseline in all models of operation. Using the terminology leads to the largest improvements in terms of BLEU.

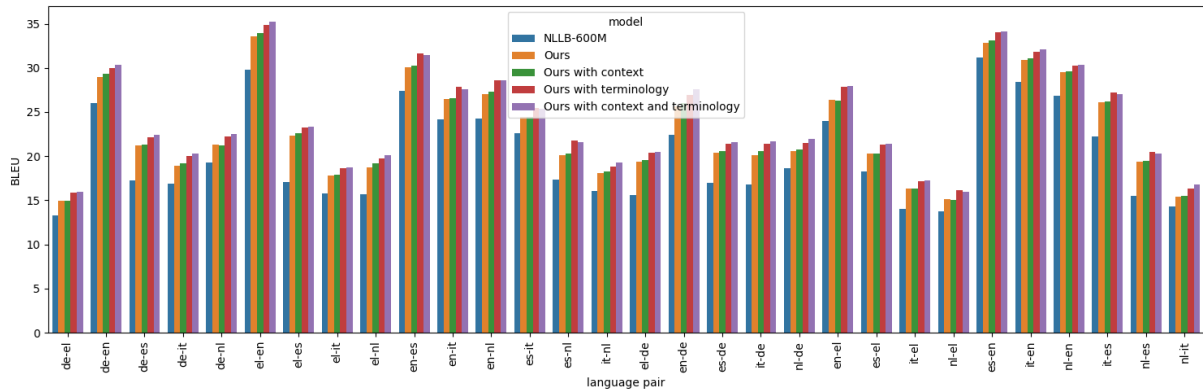


Figure 3.12. BLEU on the OpenSubtitles2018 test set in all consortium languages

Additionally, we tested the usage of the context of our model on the ContraPro (EN-DE) dataset. The results including the base model are presented in Table 3.9. The performance in terms of accuracy on ContraPro dataset of our model is increased by 4% compared to the baseline.

Table 3.9. Accuracy of the models on the ContraPro contrastive dataset

Model	Accuracy
NLLB-600M	46.44%

3.2.2. Contextual Attention Heads

In the previous section, the results show that while the models are relatively competent in context-utilization during machine translation, there are still room for improvement. For this reason, we wanted to find components in the model responsible for context integration. Motivated by the findings that attention heads can learn to perform seemingly specific functions [Voita et al., 2019; Olsson et al, 2022], we hypothesize that certain heads in a model might be crucial for the context utilization. We study the case of pronoun disambiguation, where the selection of the correct pronoun is dependent on the antecedent (contextual cue), which can be present in the previous (context) sentences. Therefore, we analyze the translation models through the lens of the attention given by the attention heads to the plausible relations that could influence the prediction of a pronoun: pronoun-antecedent on the source ($S_P \rightarrow S_C$) and target ($T_P \rightarrow T_C$) sides, and pronoun-pronoun ($T_P \rightarrow S_P$) and pronoun-antecedent ($T_P \rightarrow S_C$) between target and source sides. Because the input tokens have been found to retain their identity throughout the layers of the Transformer [Brunner et al., 2020], we add the relation between target pronoun and target antecedent on the input to the model ($T_P \rightarrow T_{C+1}$). Figure 13 illustrates the relations we consider. It should be noted that the relation pronoun-antecedent on the target side requires target-side context. Therefore, for this analysis we consider a single-encoder model that is trained to utilize context on both the source and target sides.

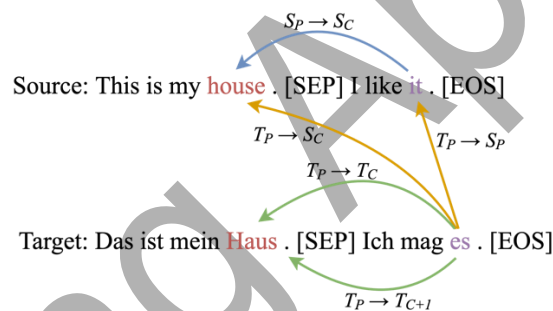


Figure 3.13. An example marking the different relations between words in source and context sentences

For each relation of interest and each attention head in the corresponding module, we evaluate the model on the ContraPro contrastive dataset and calculate three metrics:

- **Average attention scores** – we average the attention score corresponding to the relation of interest.
- **Score-Accuracy Correlation** - we calculate the point-biserial correlation coefficient between attention scores given by each head and the variable, which is defined to be 1 if the model correctly scored the example in the dataset and 0 otherwise.
- **Accuracy when Modifying Heads** – we modify the attention score corresponding to the relation of interest to a specified value and measure the model's accuracy on the dataset. This allows us to measure how the model behaves if a particular head is incapable of attending the relation (modifying to 0.01) or is perfect at attending the relation (modifying to 0.99).

Analysis

We fine-tune two publicly available sentence-level machine translation models (OpusMT en-de and NLLB-600M) for context-aware translation task. For OpusMT En-De, we train on a single direction (English-to-German) and context sizes (number of previous sentences) of one and three. For NLLB-600M, we train two separate models with context size of one for language

directions of English-to-German and English-to-French. In this report, we present the results for models based on OpusMT en-de. Other results and more detailed analysis can be found in our paper “Analyzing the Attention Heads for Pronoun Disambiguation in Context-aware Machine Translation Models” [Mařka et al, 2025].

Figure 3.14 presents the measured metrics for the models based on OpusMT En-De (sentence-level, context size 1, and context size 3). Rows represent OpenMT en-de models. X-axis presents the average attention scores. The correlations between attention scores and accuracy are at the first column and differences in accuracy when modifying to 0.01 and 0.99 are presented in second and third columns respectively. The correlation is only partially corresponding to the changes in accuracy when modifying the attention scores, meaning that we find heads that demonstrate low correlation but have a high impact on the model's performance when modified. We ignore the losses in accuracy when modified to 0.99 as possibly resulting from the decreased attention scores for other token-to-token relations. Similarly, we discard the increases in accuracy when modified to 0.01. We broadly categorized heads' behavior for a particular relation into the following groups (note that the same head can be in different groups for different relations of interest):

- **irrelevant** - heads that do not attend the relation of interest and do not respond to the modifications.
- **attending and fully responsive** - heads that attend the relation, the model's accuracy decreases when modified to 0.01 and increases when modified to 0.99.
- **attending and negatively responsive** - heads that attend the relation and respond negatively to modifying to 0.01 but do not show improvement when modified to 0.99; we interpret those heads as already at the peak of their ability to help the model with pronoun disambiguation.
- **attending and positively responsive** - similar to the previous category but with only the positive response to modifying to 0.99; we hypothesize that the heads in this category do not attend the relation to a sufficient degree for some examples and the correctly disambiguated examples are also supported by other heads in the model.
- **attending and non-responsive** - heads that attend the relation but are not responsive to any modification; our interpretation is that those heads are important for other tasks than pronoun disambiguation.
- **non-attending and positively responsive** - heads in this category are not attending the relation but positively influence the model's accuracy when modified to 0.99.

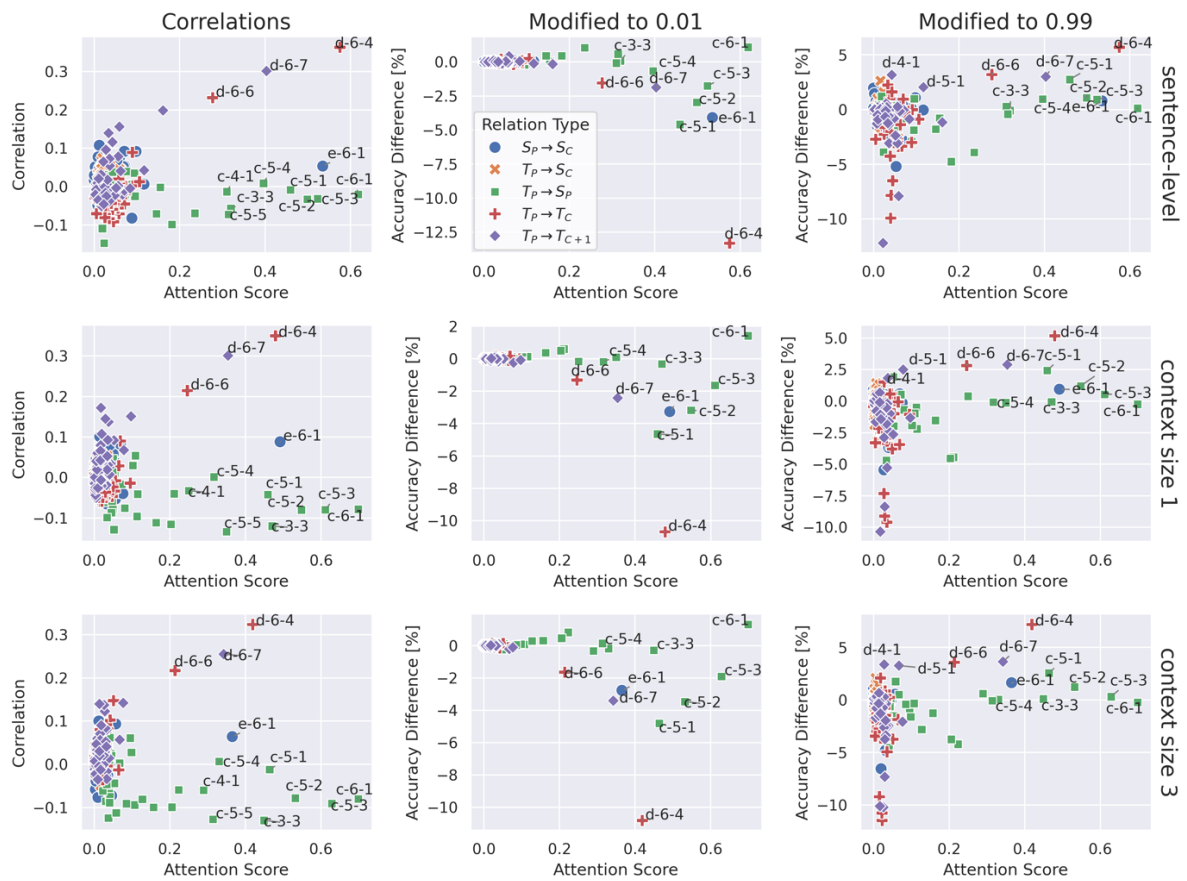


Figure 3.14 Results of the analysis with the OpenMT En-De models

Considering the results for all models (including those based on NLLB-600M) we make the following observations:

- 1) **Some heads appear to have a function identifiable through analysis.** This is confirmed by previous research [Voita et al., 2019; Olsson et al, 2022], however we additionally demonstrate that adjustments to the functioning of heads (whether improved or diminished) leads to noticeable changes in models' performance (in terms of the accuracy of pronoun disambiguation).
- 2) **The decoder-attention (corresponding to the target-side context) has the highest impact on the pronoun disambiguation accuracy.** Note that we used the gold context (provided with the examples). In a real-world system, the context would come from the model's predictions. It is interesting to note that the most relevant heads are located in higher layers. Our intuition is that in the decoder the output token is presumably decided in the layers closer to the output and only with this information can heads attempt to find the corresponding antecedent. In the decoder-attention, the important context tokens are the tokens corresponding to both the antecedent being predicted (the $T_P \rightarrow T_C$ relation) and being passed to the model as input (the $T_P \rightarrow T_{C+1}$ relation) with most heads specializing in attending one of the two and only some heads being able to utilize both.
- 3) **Attending to the relation by a head does not necessarily imply it has an impact on the pronoun disambiguation accuracy.** The most responsive heads already attend the relations of interest, but there exist heads that are not utilized by the model that could improve the performance if were attending the relations.

After identifying the most responsive heads, we fine-tuned them to assess to what extent the augmented behavior of the heads can be solidified into the models' parameters without compromising overall translation quality. To obtain the dataset containing the pronoun-antecedent pairs we applied the CTXPRO toolset [Wicks et al., 2023] to the IWSLT 2017 en-de dataset [Cettolo et al., 2017] (unrelated to the ContraPro dataset). We tuned selected heads of the OpusMT en-de models. The changes in accuracy for OpusMT en-de models when modifying and fine-tuning heads identified as influential in pronoun disambiguation are presented in Figure 3.15. The models exhibit higher accuracy for each tuned head without reducing the translation quality. The improvement is the highest for the sentence-level model and reduces with increased context size.

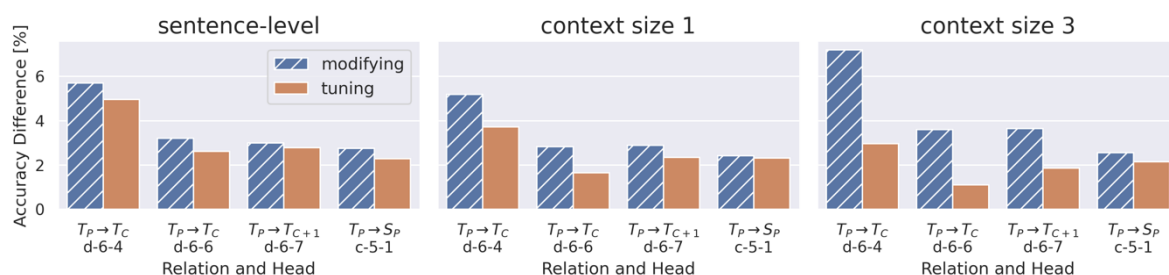


Figure 3.15 Changes in accuracy for OpusMT en-de models

Results

We first fine-tune a context-aware NLLB-600M model with the context size of one on the source and target sides on language direction between every consortium language. We randomly select 100,000 examples from the OpenSubtitles 2018 dataset for each language direction and train for 10 epochs. To take advantage of the analysis we performed, we fine-tuned the NLLB-600M model using the “dev” and “devtest” subsets of the CTXPRO dataset. We used not only pronoun disambiguation (referred to as “gender” in CTXPRO dataset) but also “animacy”, “formality”, “auxiliary” phenomena to obtain a more robust model. To avoid interference between phenomena during fine-tuning we freeze the model with only the selected heads being unfrozen. We fine-tune iteratively on one phenomenon to another and only heads corresponding to the phenomenon are trained. As a baseline, we fine-tuned the model without freezing. The training is done over 20 epochs with the learning rate of 0.0001 for head fine-tuning and 0.000001 for full fine-tuning. Notably tuning only heads leads to reduced training time roughly by half.

In Figure 3.16, we show the results of the two models as well as the sentence-level and fine-tuned context-aware models for the CTXPRO accuracy. Fully fine-tuned model improves over the context-aware model for most phenomena but decreases the performance for gender in English-to-German and English-to-Italian, and formality in English-to-Spanish. Head fine-tuning model improves over the context-aware model for all phenomena. It achieves better performance than fully fine-tuned model on 5 out of 12 tested phenomena but comes short on the other 7. It should be noted that fully fine-tuned model reduces the overall translation quality (measured by BLEU) as shown in Figure 3.17. Head tuning outperforms fully fine-tuning the model on the same metric proving the superiority of our method in terms of preventing the interference between tasks during fine-tuning.

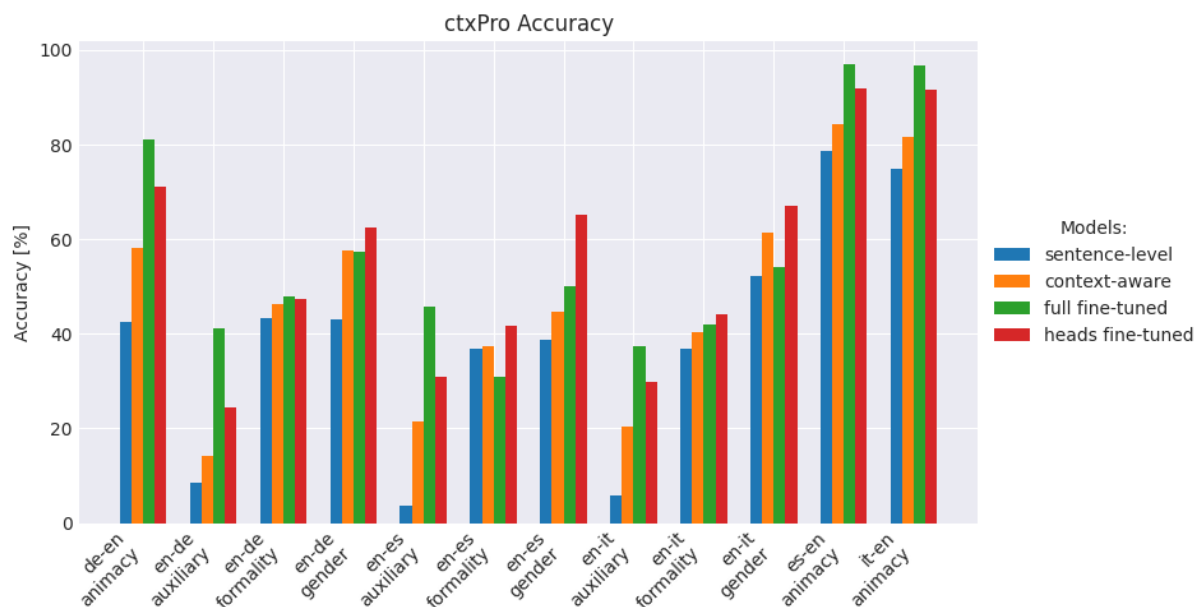


Figure 3.16 Accuracy of the models on the CTXPRO dataset

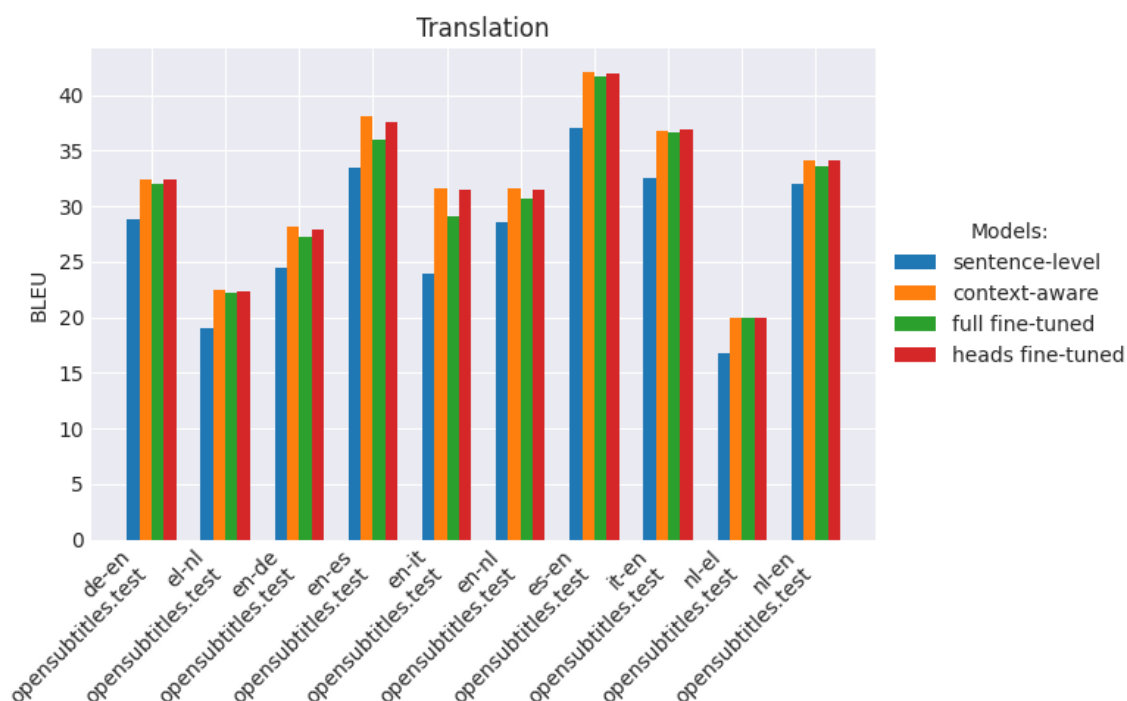


Figure 3.17 The BLEU scores of the models on the OpenSubtitles 2018 dataset

3.2.3. Robust Machine Translation

Part of our goals at VOXReality is to build robust and reliable models that are ready for production. We research techniques for building a Robust Machine Translation model that is able to produce a correct translation even when the input is noisy. We focus on non-native speaker noise and adapt the model to handle non-native input issues.

We worked on improving a model that underperforms on non-native speakers' noise. We experimented with finetuning and realised that it leads to information forgetting which negatively impacts the performance on clean data. This has led us to the solution of using adapter networks [Rebuffi et al., 2017] [Houlsby et al., 2019]. Adapter networks are small

modules that are inserted into the model and can be trained to adapt the model's representations to different tasks or domains. The model can be frozen or trained jointly with the adapters. In our case, we freeze the model to minimize information forgetting which can impact the performance, especially since we are finetuning on noisy data. We use a simple adapter architecture [Bapna et al., 2019] composed of layer normalization, followed by a down projection and up projection with a Rectified Linear Unit (ReLU) in between. The down projection is a linear layer that projects the input into a lower dimension, and the up projection brings the representations back to their original dimension. We insert the adapter network after the feedforward network of the transformer layer. The adapter output is added to the original feedforward network output, which forms a residual connection.

In an encoder-decoder architecture, the encoder takes care of building a representation of the input text that then the decoder can rely on to generate the output, which means that if we can nudge the encoder to output a clean representation even when the input is noisy, the decoder will be able to correctly generate the output. We insert the adapters into the encoder and train them to update the encoder representations while the encoder and decoder parameters are frozen. We train the adapters with two objectives: a translation objective, and a similarity objective where the decoder is motivated to produce representations that have high similarity to clean representations from the same frozen encoder. Figure 3.18 illustrates this architecture. The following equation shows the training objective of our model:

$$L = \lambda. \text{Similarity}(\text{encoder}(x_{\text{noisy}}), \text{encoder}'(x_{\text{clean}})) + (1 - \lambda). \text{CrossEntropy}(y, \hat{y})$$

where λ is a hyperparameter for controlling the impact of the similarity on the model. $\text{encoder}'$ refers to the encoder with the adapters. x_{noisy} and x_{clean} are the noisy and clean input respectively. y and \hat{y} are the target and predicted translation respectively.

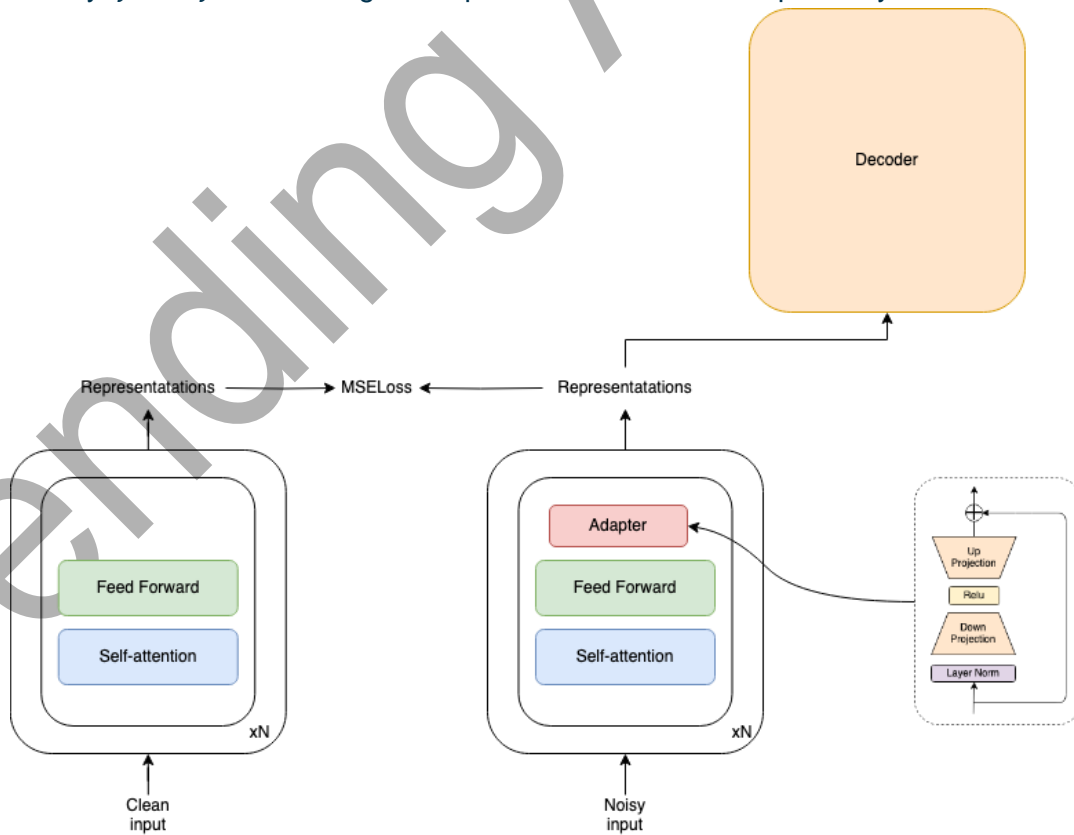


Figure 3.18. Training a Robust Model Using Adapters

Results

In order to properly evaluate our models, we make sure that we evaluate using both clean and noisy input to ensure we don't sacrifice performance on clean data for performance on the noisy one. We use The John Hopkins University (JHU) FLuency-Extended Grammatical/Ungrammatical (GUG) corpus (JFLEG) [Napoles et al., 2017], a dataset of non-native speaker text with corrections. More specifically, we use the JFLEG-es [Anastasopoulos et al., 2019] which is the same dataset paired with Spanish translations. For finetuning the model, we rely on Grammatical Error Correction datasets to acquire clean and noisy input pairs. Then we use the original model to translate the clean inputs to the target language and use the results as reference translations for finetuning. We evaluate the model using three datasets, namely: Lang-8 [Tajiri et al., 2012], First Certificate in English corpus (FCE) [Yannakoudakis et al., 2011] and Building Educational Applications (BEA) [Bryant et al., 2019]. The training dataset size of these datasets are 30543, 1.9 million, and 34309 for FCE, LANG-8, and BEA respectively. We finetune OPUS-mt-en-ROMANCE [Tiedemann & Thottingal, 2020] on each of the three datasets. We report the finetuning hyperparameters in Table 3.10.

Table 3.10. The hyper-parameters of the robust MT model

Hyper-parameter	Value
Encoder Layers	12
Decoder Layers	12
Attention Heads	16
Embed Dim	1024
FFN Embed Dim	4096
Dropout	0.1
Optimizer	Adam
Learning Rate	5e-5
LR Scheduler	Linear
Batch Size	16

We report BLEU on both noisy and clean data in Table 3.11. No finetuning refers to using the original model directly for translation. We can see that there is a gap of more than 2 BLEU points between the performance on clean and noisy data. After finetuning the whole parameters of the model, the performance on both clean and noisy data drops, which can be explained by information forgetting. Finetuning all parameters on noisy data leads the model to forget what it has learned during the initial training. Table 3.11 also reports the results of using adapters with and without the similarity loss and we can see that both techniques improve the results over the original model in both clean and noisy data, while adapters with similarity loss are better on noisy translation when finetuning on Lang-8 and BEA.

Table 3.11. A Comparison of Finetuning and Adapter based finetuning on JFLEG-es

Dataset	No finetuning		Finetuning		Adapters-Similarity Loss		Adapters	
	Clean	Noisy	Clean	Noisy	Clean	Noisy	Clean	Noisy
LANG-8	31.92	29.40	26.94	23.67	32.46	31.10	32.32	30.73
FCE	31.92	29.40	29.82	28.54	32.39	31.28	32.55	31.57
BEA	31.92	29.40	30.75	30.06	32.33	31.46	32.41	30.96

3.2.4. Robustness of Machine Translation Against Non-native Speaker Errors

In section 3.2.3, we presented that we could fine-tune translation models to be more robust by training them to push noisy and clean sentence representations to be similar. We also further analysed how robustness is represented inherently in translation models, which can offer insights building more robust models. Our hypothesis is that the encoder detects and corrects the representation of the ungrammatical word by moving its representation toward the correct form. For detection, we use Grammatical Error Detection (GED) probing to investigate how the accuracy of detecting the ungrammatical word changes through the encoder layers. For correction, we measure the representation distance between the ungrammatical word and its correct grammatical form for Grammatical Error Correction (GEC). We find that generally the probing performance increases in the first half layers of the model, then plateaus or decreases, while the representation distance decreases along model layers. To understand what contributes to correcting the representations, we turn to the attention mechanism due to its crucial contribution to transformer model performance. We identify *Robustness Heads*, which are attention heads that contribute to moving the ungrammatical word's representation toward its correct form. We find that after fine-tuning models on ungrammatical sentences, and thus, making them more robust, they learn to rely more on Robustness Heads for updating the ungrammatical word's representation especially in deeper layers where we hypothesize that the correction is happening.

To provide a representation level analysis of robustness, it is crucial to have granular control over grammatical errors. We achieve this by inserting synthetic grammatical errors into clean sentences. We focus on three types of grammatical errors that are common in non-native speaker language [Naploes et al., 2016] and we create an adversarial copy of the dataset for each type, where we insert one error per sentence when possible. We focus on grammatical errors with clear linguistic functions to be able to link our analysis to the linguistic features of the source language. These three types of errors are:

- **Nounnum:** Noun number error, where we change the noun from singular to plural or the opposite.
- **Article:** Article errors, where we replace an article with a different one.
- **Prep:** Preposition error, where we replace a preposition with a different one.

We use the Europarl-ST [Iranzo-Sánchez et al., 2019] dataset which contains official speech, transcriptions and translations of European Parliament debates of multiple European languages. In our investigation, we only use the transcriptions and translations of 5 directions: En-Es, En-De, En-It, En-Nl and Fr-Es. For conciseness, we only include the results for En-Es direction.

Our analysis focuses on four well-established NMT models, namely: OPUS-MT, M2M100, MBART and NLLB. We fine-tune these models on the adversarial dataset of one of the error types, and since this leads to improving their robustness to the error type, we also analyze the representations of fine-tuned models. To separate the effect of robustness from domain adaptation, we compare against a version of the model that is fine-tuned on the clean version of the data. We refer to these models as Base, Noise-Finetuned, and Clean-Finetuned respectively. However, we only fine-tune the encoder given that it is the source-side representation engine of the model.

Results

In Table 3.12, we present the COMET scores of Base, Clean-Finetuned and Noise-Finetuned models on clean and noisy test sets and their difference (Δ) for En-Es. Even in this simple setup, where we insert one error per sentence, we see a significant drop in performance in



Base models (at least 0.66). Furthermore, fine-tuning on clean or noisy data leads to better results but only fine-tuning on noisy data leads to improving robustness, which is seen in the reduced difference in COMET (e.g. from 0.74 to 0.01 for NLLB on Article errors). Surprisingly, fine-tuning on grammatical errors doesn't affect performance on clean data, and can even lead to better results compared to fine-tuning on clean data (77.51 compared to 77.14 for M2M100 on Prep errors).

Table 3.12: COMET of OPUS-MT, M2M100, MBART, and NLLB on En-Es direction

Model	Article			Nounnum			Prep		
	Clean	Noisy	Δ	Clean	Noisy	Δ	Clean	Noisy	Δ
opus-mt-base	78.72	77.71	1	78.72	77.74	0.97	78.72	77.67	1.05
opus-mt-clean	78.88	78.05	0.84	78.97	78.12	0.85	78.93	78.05	0.88
opus-mt-noise	78.94	78.89	0.06	78.99	78.82	0.17	79.06	78.83	0.23
m2m100-base	75.99	74.84	1.15	75.99	75.12	0.88	75.99	74.63	1.36
m2m100-clean	77.39	76.4	0.99	77.28	76.22	1.07	77.14	76.14	1
m2m100-noise	77.57	77.54	0.03	77.58	77.5	0.08	77.51	77.23	0.28
mbart-base	78.04	77.23	0.81	78.04	77.25	0.79	78.04	77.24	0.79
mbart-clean	78.51	77.77	0.74	78.64	77.74	0.9	78.49	77.73	0.76
mbart-noise	78.61	78.58	0.04	78.57	78.51	0.06	78.58	78.38	0.2
nllb-base	78.34	77.6	0.74	78.34	77.69	0.66	78.34	77.52	0.83
nllb-clean	78.82	78.14	0.68	78.85	78.16	0.69	78.81	78.21	0.6
nllb-noise	78.81	78.8	0.01	78.78	78.73	0.06	78.94	78.66	0.29

After we established that grammatical errors lead to a significant drop in performance and fine-tuning on them leads to increased robustness, we proceed to analyze Base, Noise-Finetuned and Clean-Finetuned models' representations when responding to the grammatical errors. In Figure 3.19, we show the result of probing each layer of the models in terms of F1 scores. According to these results, the fine-tuning negatively affects the GED probing performance in deeper layers, where our GEC setup hypothesis suggests that correction is happening. Noise-Finetuned models maintain their ability to detect errors in lower layers, then correct them in deeper layers, which makes the GED probing accuracy more challenging because the representation is corrected.

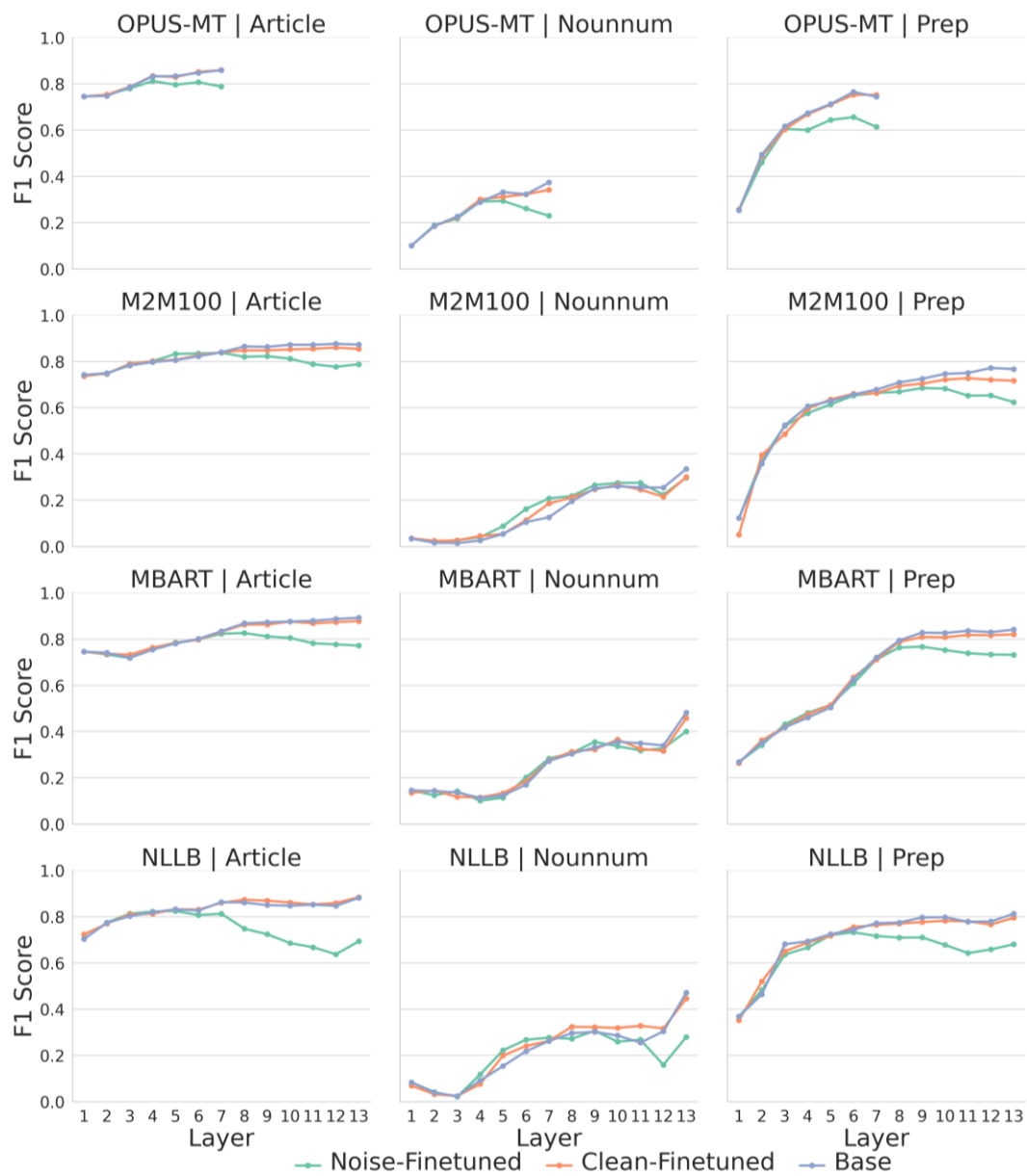


Figure 3.19: GED probing performance of the models on En-Es direction

Figure 3.20 further supports the correction hypothesis, where models, in general, attempt to push the representations of the noisy words towards the clean word representations since the distance between the two, calculated using the Centered Kernel Alignment (CKA) [Kornblith et al, 2019], becomes smaller as in each layer. Compared to Base and Clean-Finetuned models, the Noise-Finetuned model performs better at this, which explains its robustness as shown in Table 3.12.

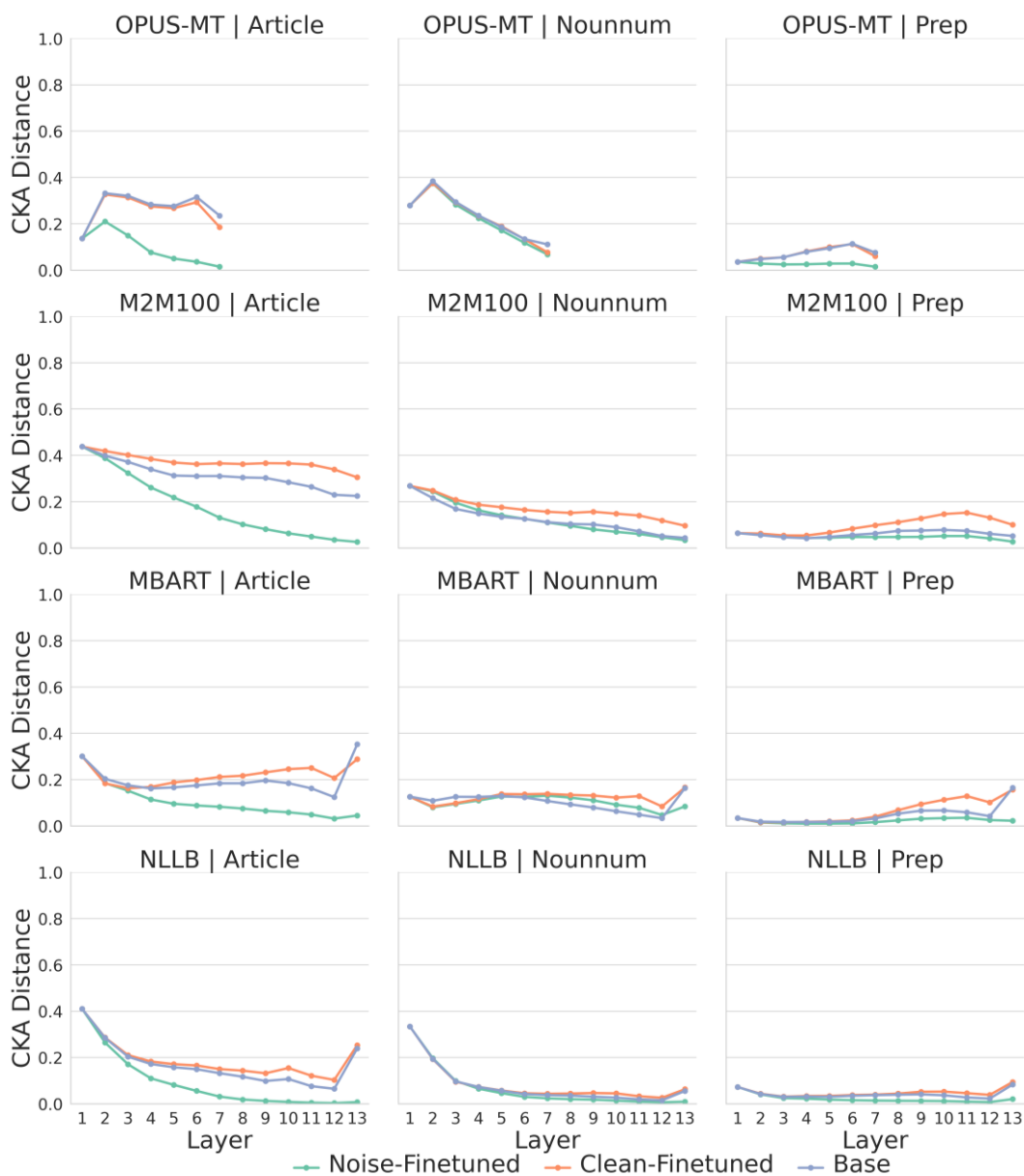


Figure 3.20: CKA distances of word representations on En-Es direction

The GED probing and representation distance results support the hypothesis that translation encoders detect then correct the representation of errors to achieve improved robustness. In Figure 3.21, we analyze Robustness Heads, which are attention heads that contribute the most to pushing the representation of the noisy word towards its clean form. We look into how these attention heads attend to important parts of the sentence by applying Part of Speech (POS) tagging, and collecting the attention scores assigned to each POS tag when resolving each of the three types of errors that we study. For clarity reasons, we only show the attention scores over the 10 most common POS tags in the dataset, and we keep their universal names⁴. The figure shows that Robustness Heads attend to important parts of the sentence depending on the error. For Article errors, the models focus on adjectives (ADJ) and nouns (NOUN and PROPN), which are important for handling article errors such as deciding between using “a” or “an” depending on whether the next word starts with a consonant or a vowel. When dealing with Nounnum errors we find that Robustness Heads attend mostly to adjectives and determinants (DET). Finally, for Prep errors, Robustness Heads attend to nouns, determinants and verbs (VERB).

⁴ <https://universaldependencies.org/u/pos/>

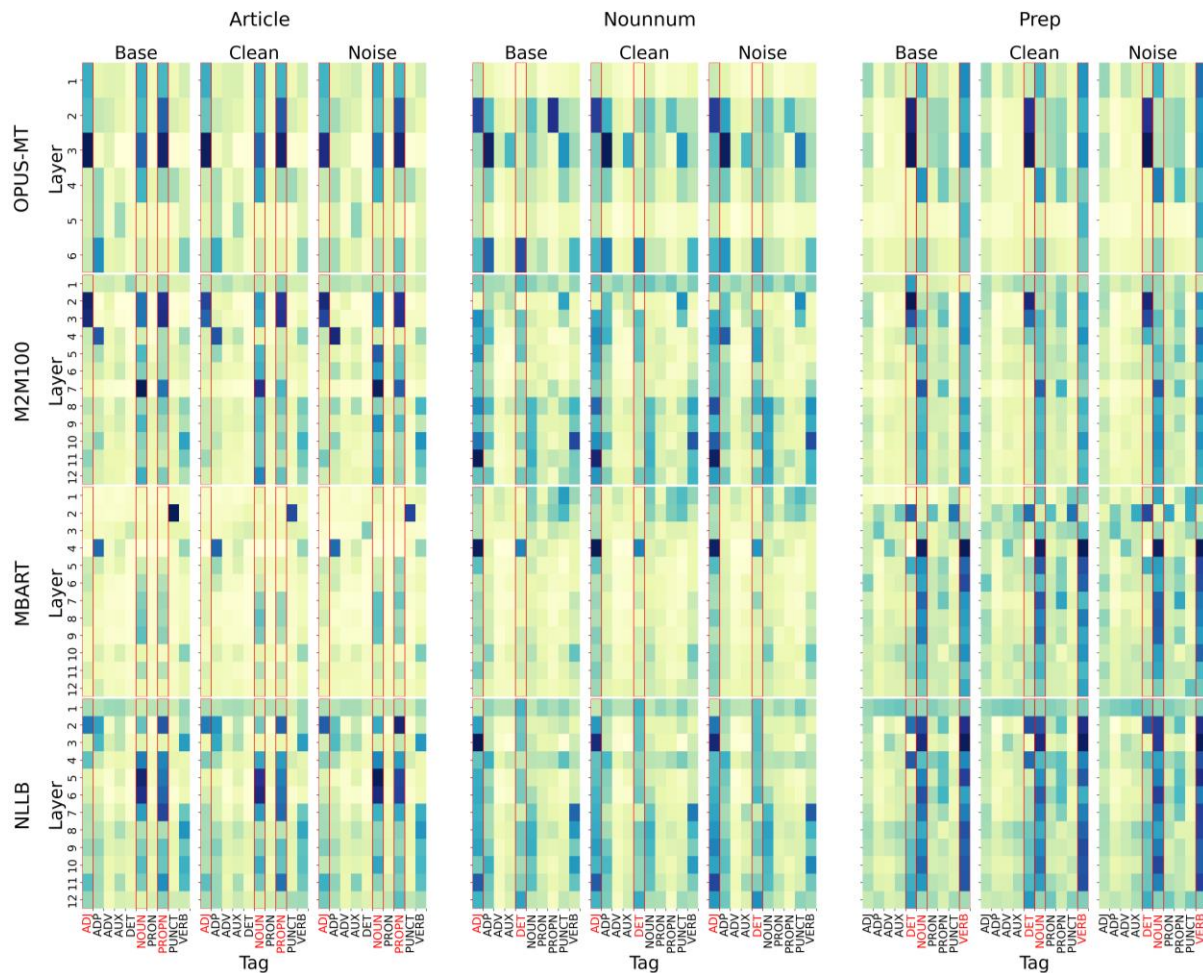


Figure 3.21: Attention to the 10 most common POS tags on En-Es direction

When comparing the Robustness Heads of Base, Clean-Finetuned and Noise-Finetuned in terms of their attention to POS tags, we find that although in some cases the Noise-Finetuned model has learned to attend more to the important POS tags, in general, they attend similarly to interpretable POS tags when dealing with grammatical errors. To investigate what makes the Noise-Finetuned models succeed in terms of robustness, we identify Influential Heads of the noisy word, as the heads that most influence it toward its current state. Then, we compare them with Robustness Heads. Figure 3.22 presents the accuracy between Influential Heads and Robustness Heads at each layer of the encoder. According to the results, this accuracy is higher in Noise-Finetuned models, especially in deeper layers, meaning that models after fine-tuning on noise tend to employ more Robustness Heads for updating the noisy word representation.

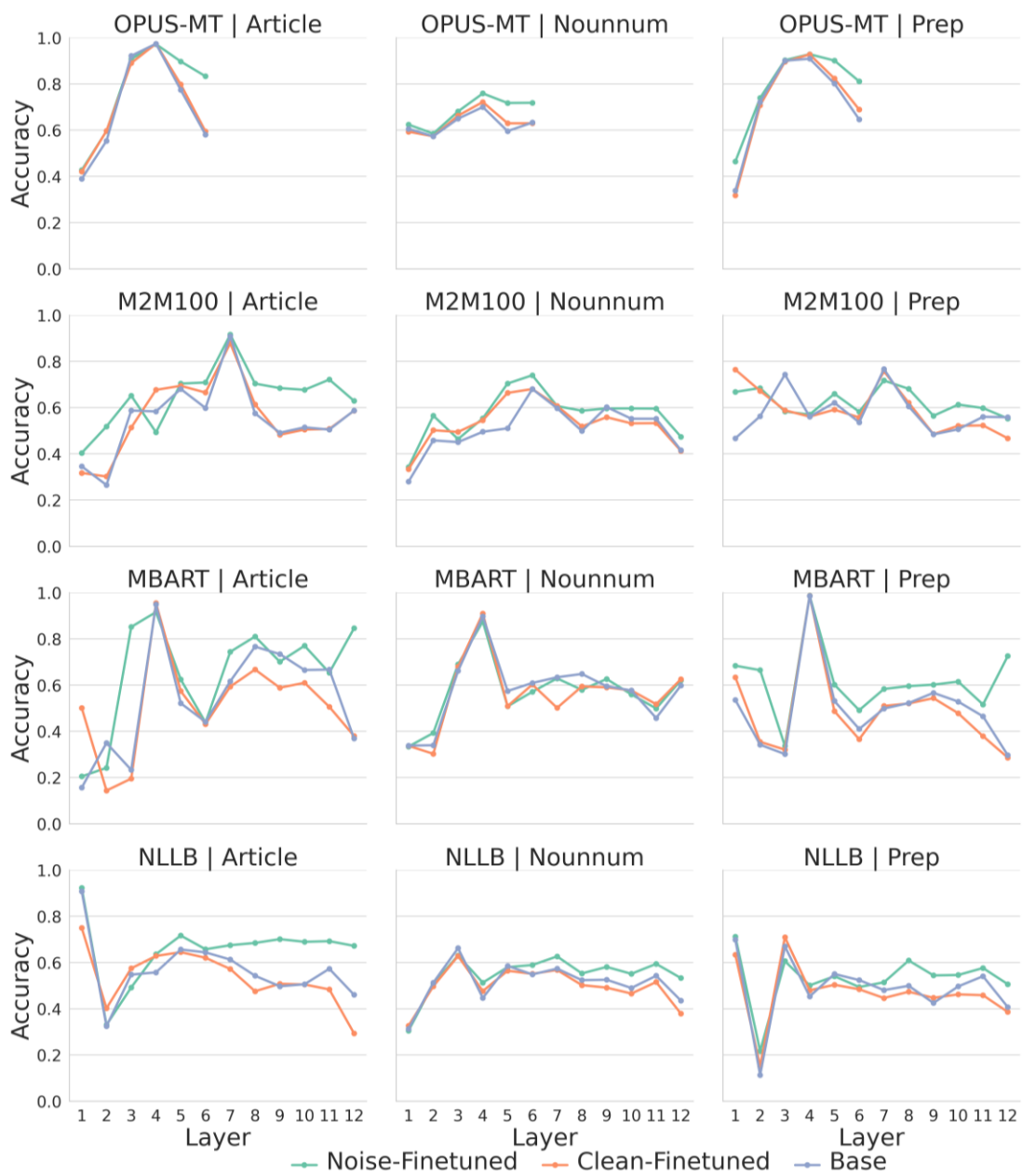


Figure 3.22: Accuracy of Influential Heads on En-Es direction

3.2.5. Simultaneous Machine Translation

In VOXReality we introduce two approaches for simultaneous machine translation: a) Multi-path training with adapters: a single model that can support multiple fixed wait-k values by using adapters and b) Adaptive strategy: utilizing the probabilities that the model assigns to the most likely token and a predefined probability threshold, we decide on either a “read” (the model continues to read new tokens) or a “write” (the model can continue with translation) action.

In multi-path training with adapters, we insert adapters into the decoder layer, following the implementation of [Bapna et al, 2019] and insert the residual adapter modules after the feed-forward layer. Each adapter is activated for a set of train or test wait-k values. During training, the wait-k values are sampled uniformly following multi-path training [Elbayad et al., 2020]. In the adaptive strategy, we follow the approach of [Zheng et al., 2020] to build an adaptive strategy by using adapters instead of different models for each wait-k value. At each decoding step, we activate one single adapter based on the lagging behind the current generation step. Then, we utilize the probability of the most likely token to decide whether to write or read a new token. If the probability is less than the set threshold, we read a new token, otherwise we write. Furthermore, if wait value (k) is lower than the minimum wait value, we force the model to read and if it is higher than the maximum wait-k, we force the model to write, which means that the choice of minimum and maximum values of wait also impacts latency. The architecture of the model is presented in Figure 3.23.

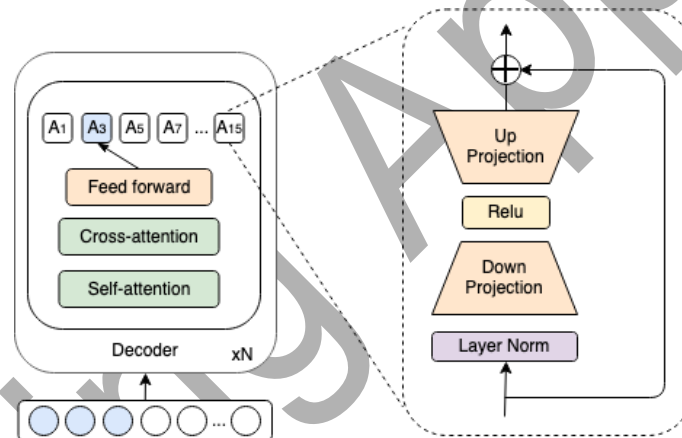


Figure 3.23. The architecture of wait-k adapter model

Results

We utilized our methods in three transformer versions: transformer-small, transformer-base and transformer-big. The hyperparameters for these models are presented in Table 3.13. The evaluation was performed with two public datasets: EN-VI (English - Vietnamese) dataset from IWSLT 2015 Evaluation Campaign [Cettolo et al., 2015] for evaluating transformer-small models, and DE-EN (German - English) dataset for evaluating transformer-base and transformer-big models. The EN-VI dataset consists of 133K pairs and DE-EN dataset consists of 4.5M pairs.

Table 3.13: Hyperparameters of the Simultaneous Machine Translation Models

Hyper-parameter	Transformer-small	Transformer-base	Transformer-big
Encoder Layers	6	6	6
Decoder Layers	6	6	6
Attention Heads	4	8	16

Embed Dim	512	512	1024
FFN Embed Dim	1024	2048	4096
Dropout	0.3	0.3	0.3
Optimizer	Adam	Adam	Adam
Learning Rate	5e-4	5e-4	5e-4
LR Scheduler	Inverse Sqrt	Inverse Sqrt	Inverse Sqrt
Batch Max Tokens	16000	8192	8192
Number of Adapters	8	8	8
Adapter Bottleneck	64	64	64

We trained the following models:

- **Offline Transformer** - where the model takes the full sentence
- **Wait-k** - where the model waits for k source tokens before starting to alternate between writing a target token and reading a source token [Ma et al., 2018]
- **Multipath Wait-k** - where the model supports multiple wait-k policies by randomly sampling k during training and the k value is fixed during inference [Elbayad et al., 2020]
- **Adaptive Wait-k** - where there are multiple models for different wait-k values and during inference the model is selected based on the lagging behind generation step and the decision to write or read is based on the output probabilities. [Zheng et al., 2020]
- **MoE Wait-k** - where the model uses multipath wait-k approach but it utilizes Mixture-of-Experts (MoE) method instead of random [Zhang & Feng, 2021]
- **MMA** - where the model follows the Monotonic Multi-head Attention (MMA) that jointly learns a Bernoulli variable that is used to decide read/write action [Ma et al., 2019]
- **Adapters Wait-k** - multi-path training with adapters (our proposed method)
- **Adaptive Adapters** - the adaptive strategy (our proposed method)

Figure 3.24 presents the results obtained from the comparison with average lagging (k-latency) and BLEU score of our proposed methods with the other models we trained, where Figure 3.24a is EN-VI transformer-small, Figure 3.24b is DE-EN transformer-base, and Figure 3.24c is DE-EN transformer-big model results. According to these results, our method improves or competes with other strategies on different latencies while using a single model. MMA, Wait-k, and Adaptive-waitk require the training of multiple models in order to support multiple latencies, while our method is more flexible in this regard. The number and capacity of the adapters can be adjusted depending on the complexity of the language direction and the latencies we are planning to support. Furthermore, using adapters alone (Adapters Wait-k) is competitive with other methods especially on EN-VI, but using an adaptive strategy further improves the results, especially in low latencies. In comparison to Adaptive Wait-k, where wait-k policy models are trained and composed during inference, we find that our method is better in all latencies while being more efficient. Compared to Moe-Waitk, which also aims at minimizing interference introduced by multi-path training, we find that our method is better in all latencies on EN-VI for transformer-small and on DE-EN with transformer-big, while achieving competitive results when using transformer-base.

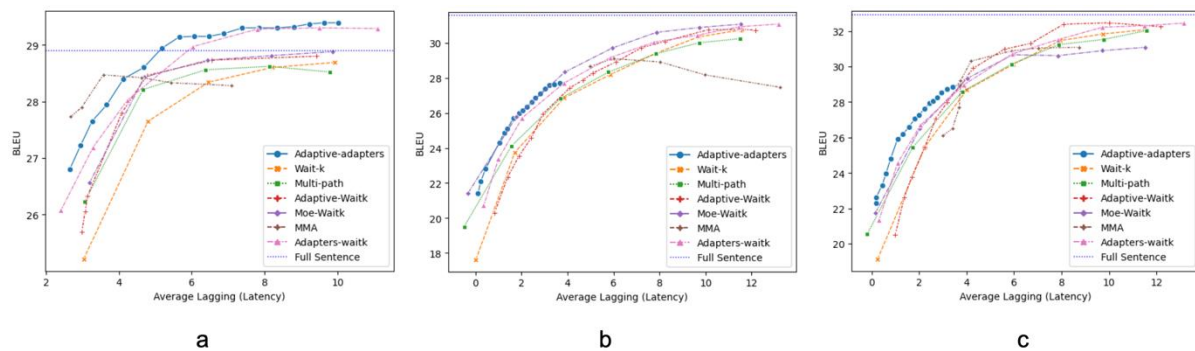


Figure 3.24. BLEU-Average lagging results for a) transformer-small, b) transformer-base, and c) transformer-big

3.2.6. Subtitle Alignment

Some use cases of the translation system might require a reliable and faithful translation that the current state-of-the-art AI systems do not achieve. One example is the Theatrical Texts, where the translations must not only be accurate but also faithful to the nature of the theatre play. For those cases, we developed the Subtitle Alignment system. In essence, the text (for example a theatrical play or presentation transcript) can be divided into utterances or subtitle lines and translated beforehand using our translation system. Furthermore, the translations can be inspected and corrected by a professional translator if needed. The file containing the subtitles in all (or selected) supported languages is uploaded to the system. Subsequently, the system will use the ASR functionality to obtain the raw transcription. Lastly, we match the transcribed speech with the provided subtitles using the chrF metric and return top-k lines in the original and the desired language as well as the corresponding metric scores. The score can be used to disregard the provided translation if the score is below a predefined threshold. This ensures that we maintain a high quality of subtitles while delivering low latency transcription and translation.

Results

To test the subtitle alignment system, we used “Hippolytus” by Euripides play containing 932 lines of subtitles in Greek. This text is difficult for the ASR system because of the literary language and uncommon names of the characters. We selected 90 lines from the original text and transformed it to speech using Google’s Text-To-Speech service and the voice “el-GR-Wavenet-A”. Next, we transcribed the resulting audio files back to text using OpenAI’s whisper-small model [Radford et al., 2023]. Figure 3.25 shows the chrF scores between transcriptions and original text of selected lines. According to the outcomes of this analysis, the metric successfully scores the highest with the corresponding transcription and original text (the diagonal in the figure).

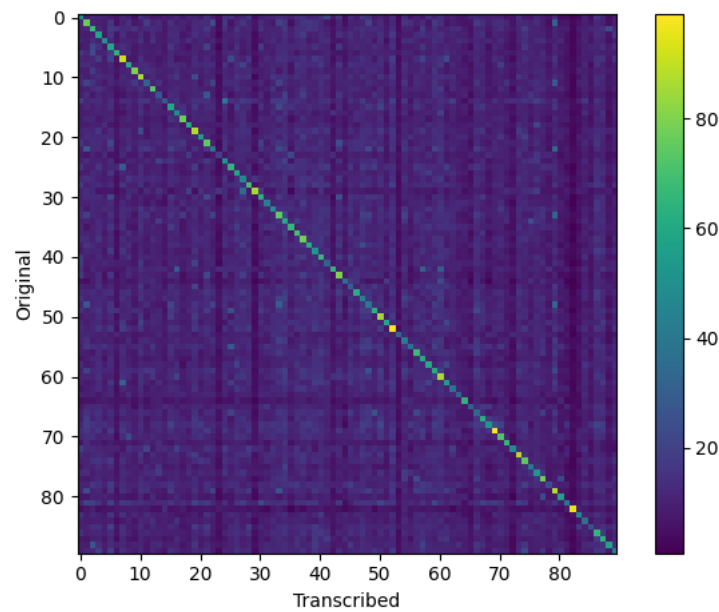


Figure 3.25. The chrF between the subtitles and the transcriptions

Next, we test the system on the full play. We first initialize the system with the file containing all subtitle lines. Next, we use the previously mentioned transcribed files and retrieve the matching lines. Table 3.14 shows the results in terms of BLEU and chrF with the true line as the reference for the raw transcriptions and our system. While not all lines have been matched, the improvement is high. Notably, combining subtitle alignment and the improvements in the ASR system can lead to even better performance.

Table 3.14: Performance of the approaches in Hippolytus

Method	BLEU	chrF
Transcription only	18.04	59.74
Transcription+ Subtitle Alignment	44.98	73.40

3.2.7. Deployment

In the VOXReality pipeline, the translation process is handled through a set of endpoints in a REST API developed using FastAPI framework. The machine translation endpoints are designed to generate translations from textual inputs. There are four endpoints: text translation, context-aware translation, context-aware terminological translation, and terminology configuration. The details of the API calls are presented in Appendix I.

In text translation, the function generates a textual response in the target language to a given text. The endpoint does not require a source language to be specified since it is automatically recognized by the model itself. Figure 3.26 presents the FastAPI powered endpoint that can be requested using the handle “translate_text” from the deployed server.

POST /translate_text Translate Text

Parameters Cancel

Name	Description
source_language * required (query)	en
target_language * required (query)	en

Request body required application/json

```
{
  "text": "string"
}
```

Figure 3.26. The text translation endpoint

In context-aware translation, the function generates a textual response in the target language to a given text and a given context, which is provided as textual information. The endpoint does not require a source language to be specified since it is automatically recognized by the model itself. Figure 3.27 presents the FastAPI powered endpoint that can be requested using the handle “contextual_translate_text” from the deployed server.

POST /contextual_translate_text Contextual Translate Text

Parameters Cancel

Name	Description
source_language * required (query)	en
target_language * required (query)	en

Request body required application/json

```
{
  "text": "string",
  "context": "string"
}
```

Figure 3.27. The context-aware translation endpoint

In context-aware terminological translation, the function generates a textual response in the target language to a given text and a given context, which is provided as textual information, as in the context-aware translation, but the model in this endpoint also utilizes a file that is uploaded to the server, which contains terminological terms and their translations in other languages. The endpoint does not require a source language to be specified since it is automatically recognized by the model itself. Figure 3.28 presents the FastAPI powered endpoint that can be requested using the handle “contextual_terminology_translate_text” from the deployed server.

POST /contextual_terminology_translate_text Contextual Terminology Translate Text

Parameters Cancel

Name	Description
source_language * required (query)	en
target_language * required (query)	en

Request body required application/json

```
{
  "text": "string",
  "context": "string"
}
```

Figure 3.28 The context-aware terminological translation endpoint

In terminology configuration, the function uploads a file that contains terminological terms and their translations in other languages to the server. Figure 3.29 presents the FastAPI powered

endpoint that can be requested using the handle “upload_terminology” from the deployed server.



Figure 3.29 The terminology configuration endpoint

In subtitle alignment system, the subtitles .csv file have to be uploaded through the “upload_subtitles” endpoint. The file should contain the subtitle lines (as rows) in the corresponding languages (columns, first row indicates the language: en, de, es, el, nl, it). Figure 3.30 shows the FastAPI endpoint.

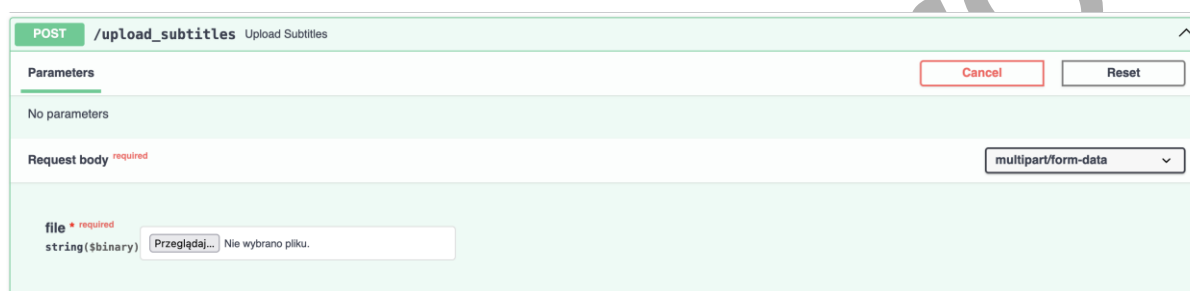


Figure 3.30 The subtitle upload endpoint

3.3. Vision-Language Models (VL)

In compliance with the task description in the project’s Grant Agreement (GA), we are developing a series of spatially aware Vision-Language (VL) models that can potentially act as “visual context” providers to the rest of the VOXReality’s pipeline. In this way, various other components like the Neural Machine Translation (NMT), the Automatic Speech Recognition (ASR) and the Conversation Agents (CA) of the use-cases may benefit from additional, current knowledge of the surrounding environment, adapting their behaviour to that.

Core-concept of any VL model is the apparent alignment of Vision and Text modalities, which the network implicitly has to learn in order to perform all the complex VL tasks. All these tasks require foremost the connection between the various entities of the visual media (i.e. image, video) and the words in the accompanying sentence. This is called “vision and language representations alignment”. Note that these connections can be so subtle that may refer to sub-word (minimum meaningful parts of words/stems) and sub-object levels (patterns in objects).

In literature, we commonly come across with two prominent VL tasks: the image/video captioning (IC/VC) and the visual question answering (VQA). In the first task, the model is asked to provide captions of an image, while in the second the model is asked to answer a question about the image. What differentiates our tasks to the typical VL ones is that our models need to understand, describe and answer questions related with the spatial aspects

of the word, meaning that generic descriptions about the nature of the scene will not be considered (e.g. the caption “a cat is playing with the ball” is not acceptable. We prefer something similar to “a cat is to the left of the ball”). This requirement puts extra strain on the devised model architecture, since the training and evaluation protocol and the data gathering and selection should all be able to facilitate such an out-of-envelope approach to VL. Apart from the two typical VL tasks, in the context of the VR Conferences use case we also explored a novel approach of repurposing a VQA architecture for providing turn-by-turn navigation instructions.

3.3.1. Early models delivery

In order to facilitate the timely and unobstructed development across the VOXReality pipeline, we deemed necessary to define early in the project a standardized model delivery mechanism and an inference API that, more or less, will remain constant throughout the project’s pilots and Open Calls (OC).

Instead of simply defining the templates, we preferred to adapt, package and deliver to the rest of the platform actual state-of-the-art (SOTA) VL models (not spatially aware though), to enable use-case developers to infer and use real working models, in a way similar to what they are about to experience with ours. So, we released four inferable SOTA VL models addressing all our intended tasks: Image Captioning⁵, Video Captioning [Lin et al., 2021], Visual Question Answering (single-word answer) [Tan & Bansal, 2019], Visual Question Answering (whole-sentence answer) [Li et al., 2022] and an extra composite service that provides a full spatial description of an image exploiting Object Detection (YOLO-NAS⁶), transformers-based monocular Depth Estimation [Ranftl et al., 2021] and our in-house heuristics for sentence creation.

We decided that all our services provided to the project will be in the form of Docker images, GPU-accelerated when needed, exposing RESTful API⁷ endpoints for inference. Details on the exact message formatting can be found in Appendix II. Additionally, the services will also be providing a web interface built around FastAPI⁸ enabling GUI-based interactive inference, ideal for demonstration and debugging purposes. Figure 3.31 presents the endpoints for the (a) early image captioning model, (b) early video captioning model, (c and d) the two early visual question answering models, and (e) the composite spatial scene description (scene graph) service used for the creation of our COCO extension dataset.

⁵ <https://qi-xin.github.io/image%20caption%20generation.pdf>

⁶ <https://github.com/Deci-AI/super-gradients/blob/master/YOLONAS.md>

⁷ <https://docs.github.com/en/rest?apiVersion=2022-11-28>

⁸ <https://fastapi.tiangolo.com/>

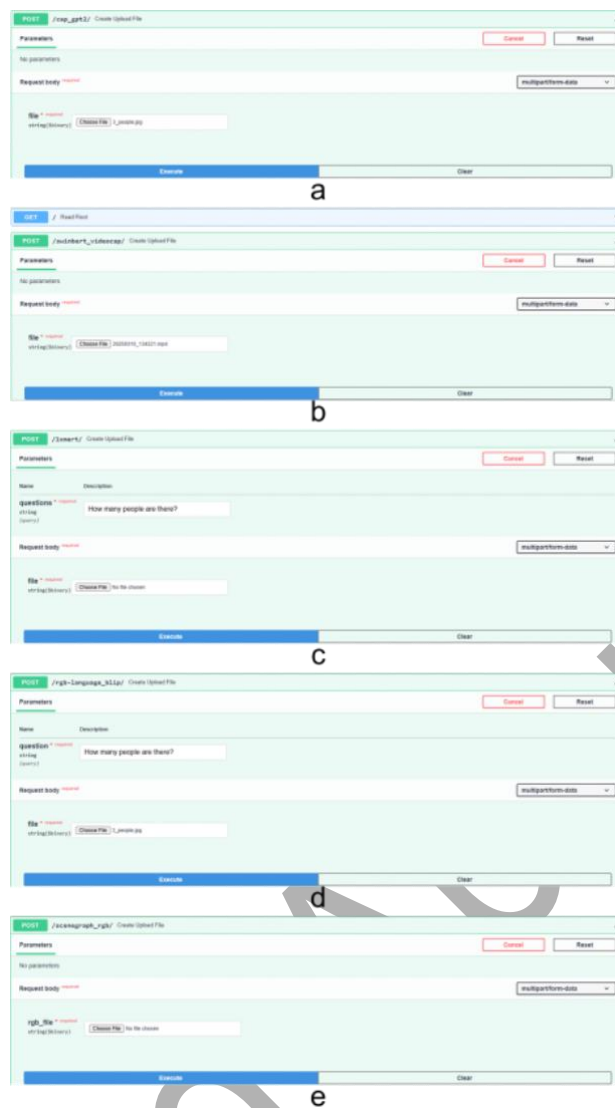


Figure 3.31. Early endpoints for Visual Language Models

3.3.2. Metrics and limitations

A general problem in evaluating NLP, and consequently VL, model predictions is that in human language it is common to find sentences that have almost the same meaning by utilizing synonyms or have the exact same meaning despite having different word ordering. On the other hand, it is also common having sentences that although contain the exact same words, word order changes can alter their meaning significantly. Moreover, sentences that have all words equal and in the same order, with the absence of even one word may totally lose their meaning or natural sound (e.g. the removal of one “not”). Others, with the addition of one redundant word can still have the exact same meaning. All these phenomena, in most cases, render the use of naive similarity or entropy metrics at least questionable for the task of fairly evaluating the output of an NLP model and even though they frequently participate in training loss functions, rarely should be used for actual evaluation.

There are multiple evaluation metrics proposed in literature, each one trying to alleviate different pitfalls of the aforementioned linguistic behaviours. So, it is safe to assume that there is no metric adequate to cover meticulously every aspect of the performance envelope and using multiple ones may usually be the preferred approach. Importantly, since our intended use is for the evaluation of spatiality-containing language, even the most widely used NLP-targeted metrics may not be able to capture our intentions. For this specific reason, we devised our own metric which in its simplicity we believe captures better the spatial accuracy of a

model and is most relevant to our case. In the following paragraphs, we explain more about our devised evaluation metric.

When planning for our developments, we soon realized that all these of the well-established metrics might not be suitable for judging a produced sentence for its spatial accuracy. The reason is that the correctness of a sentence we care about does not lie in the top-level semantics of the words contained (their individual meaning) but in the existence of the objects described (correct object identification) and the correct description of the relationships between them (their relative positions). BLEU, ROUGE and METEOR may capture these relationships to some extent, but only incidentally. There has been works like [Kritharoula et al., 2023] and [Liu et al., 2022] that face similar challenges and develop their own evaluation procedures, trying to attack the problem each one from another angle. Our devised approach is described below:

Concurring with the general consensus, we believe that a single number cannot capture all aspects of performance in this challenging task. So, we devised a parametrizable metric that when calculated for the whole dataset intends to answer 2 questions; A: “What is the probability of a caption to be correct” (to describe a valid relationship) and B: “What is the probability of an image to be captioned correctly?” (to have a caption that describes a valid relationship between object in it). We need to remind here that our models generate a small paragraph describing an image and not a single sentence. Consequently, many captions may refer to the same image. Additionally, a different number of captions may be generated for each image. Let us call our metric “ $OURS_{h-\{A,B\}}$ ”. For each variation, we define h to be the number of sentences we wish to consider for each image (the first $h[1,MAX]$ captions generated). We decided to consider only 1,2,3 and the MAX imum number of captions the models can generate for each image. A and B indicate the question answered we mentioned earlier. Note that, we omit the $\{1,2,3\}$ - B cases because, obviously, they are equal to their respective A -cases, as difference between A and B scores arises only when a different number of captions is considered between images (the case of MAX - B). In short notation, each of the variations can be defined as:

- $OURS_{\{1,2,3,MAX\}-A} = \frac{nTrue}{nTrue+nFalse}$

where $nTrue$ is the total number of captions that can be exact matched in their respective images’ ground truth paragraphs and $nTrue+nFalse$ is the total number of generated captions (in the case of $h=1$, $nTrue+nFalse=nImages$).

- $OURS_{MAX-B} = \frac{\sum_{i=1}^{nImages} \frac{nTrue_i}{MAX_i}}{nImages}$

where $nTrue_i$ is the number of exact matched captions of image i in its ground truth paragraph and MAX_i is the maximum number of captions the model can generate for image i and can vary between images.

A question that naturally arises is “Why the exact matching between a generated caption and the respective image’s ground truth is a robust and fair way to extract truthfulness?”. Anecdotaly, we realised that in our use case the models are trained on a very basic vocabulary that includes solely the object labels (“the cat”, “the car”, “the human”, etc.) and a limited number of phrases that indicate relative positioning (“is in front of”, “is behind”, “is above”, “is below”, “is to the left of”, “is to the right of”). In that context, the models negligibly fail to produce a grammatically valid sentence and any errors are almost always due to object misidentifications and/or wrong identification of their relative positioning, which is exactly what

we look for. So, in fair terms, we can assume that exact searching of the generated captions inside the ground truth paragraph is reliable.

Our metric, of course, comes with its own limitations and considerations. The most notable one is being dependent on the correct identification of the described objects. Consider the example in Figure 3.32 with a flower in a vase and a kitten next to it.



Figure 3.32. An example image from the spatial dataset

Prediction: "The cat is to the right of the vase?" Ground Truth: "The kitten is to the right of the vase."

The problem is that even if the relative position between the identified objects is correct and a human would totally make sense of the slight misidentification of the cat to be the kitten (and the kitten is also a cat), our metric would evaluate it being wrong. A more challenging example would be to misidentify the kitten as a leopard. In the context of only two objects being present and the feline nature of the leopard, a human would guess that it probably meant the kitten rather than the vase and potentially evaluate the relationship as correct. As a countermeasure to all these, we could try to add a pre-processing step similar to METEOR's that use an external source to identify stems and synonyms, in an effort to minimize these effects; but still, it is obvious that there is no clear answer to these challenges.

In the second phase of VOXReality, we delved into the visual spatial question answering task (spatial VQA). Accordingly, the need to develop a suitable spatiality-focused metric for this task as well, arose. Drawing inspiration from our previous work on the evaluation of spatially aware captions, we extended the family of OURS metrics to also consider the VQA task. In that manner, we created two more variants of the metric, treating the spatial VQA evaluation problem from two different angles, analogous to what we did for the spatial captioning. Note that, given the limited vocabulary of our datasets and the power of Transformer-based architectures, we chose the use of the exact-match criterion to be a good fit for our purposes, once more. The two new variants of OURS metric are:

- $OURS_{VQA-A} = \frac{nTrue}{nTrue+nFalse}$

where $nTrue$ is the total number of correct answers (exact matched in their respective images' ground truth question-answer pairs), and $nTrue+nFalse$ is the total number of answers generated.

- $OURS_{VQA-B} = \frac{\sum_{i=1}^{nImages} \frac{nTrue_i}{nQ_i}}{nImages}$

where $nTrue_i$ is the number of correct answers for image i with nQ_i being the number

of question-answer pairs available in the ground truth of image I and can vary between images.

Finally, note that for all variants of OURS metric, as a pre-process step, before applying any operation or calculation on the subject sentences, we ensure that they contain only lower-case characters and are trimmed of redundant whitespaces or special characters.

3.3.3. Image spatial captioning

Dataset

As we described in previous chapters, we aim to develop models for 3 cases: Vision (RGB)-Language, 3D (depth)-Language and Video-Language. To achieve spatial understanding in any case, having good amounts of data spatially annotated is of primal importance. We soon realized that there were not big repositories with such data openly available, at least in the form we imagined having them. So, we deemed necessary to extend a widely used VL dataset in the spatial domain to have the large amounts of data we needed to initiate the development of the models. COCO 2014-2015 dataset is considered one of the golden standards of the field, comprising about 164.000 RGB images paired with generic captions. In order to extend this massive dataset to the spatial domain, we devised a composite mechanism able to produce spatial captions, questions & answers and depth images for the whole dataset automatically. Note that, the depth images produced are needed internally by this mechanism to create the spatial captions and are not used directly anywhere by the models developed for the Vision-Language case. In the future, the 3D-Language models we are about to develop will probably process them, to profoundly understand 3D space. The core part of this mechanism, the “spatial scene description (scene-graph)” generator that produces for a given image a paragraph describing all the 1-1 relationships between the detected objects of the scene, is dockerized and released for use together with our early SOTA models. The algorithm we devised to automatically produce the spatial captions for all images is the following:

We take for granted that all positioning will be relative to the viewer’s point of view and not according to the perceived orientation of the respective objects. So, when we say “to the left of the gorilla” we mean relative to our view, not the gorilla’s view. We chose to evaluate all 1-1 relationships between the identified objects in a scene in 3-axes; lateral (left-right), longitudinal (in front of-behind) and height (above-below). We evaluate lateral and height relationships on the RGB images and longitudinal ones in their respective depth images. We produce the depth images using a SOTA transformers-based depth-estimation model [Ranftl et al., 2021] and since we are interested only in the objects’ depth-ordering and not their exact depth, we can safely assume that this model is powerful enough to produce almost perfect results. We randomly sampled and verified manually the truthiness of the inferred depth-ordering in about 200 out of the 164.000 created images. In any case, as we will describe in a moment, objects with close depth estimations are not even evaluated in this axis (remains inconclusive what is in front of what), so they cannot teach the model erroneously.

For each RGB image, the procedure initiates by detecting all the objects in the scene. This is done using the latest iteration of the SOTA Object Detector YOLO-NAS and results in a list containing the label, 2D-bounding box (y_{min} , y_{max} , x_{min} , x_{max}) and centre (c_x , c_y) coordinates of all detected entities. Then, for each item in this list we evaluate its relative position against all other items across the 3 aforementioned axes. To form the caption sentences we arrange and concatenate the objects labels with the definite article “the”, together with a fixed set of phrases indicative of the spatial relationship between the involved objects, using a heuristic set of rules.

Using as example the lateral axis (x) and 2 objects A and B, we evaluate whether B’s left-most bounding box boundary (Bx_{min}) is to the right of the right-most boundary of A’s (Ax_{max}). If it is,

then we form the sentences “The is to the right of the <A>” and equally “The <A> is to the left of the ”. If it is not, then we check whether their existing overlap is insignificant enough compared to the scale of the items, that it can still be considered that one item is to the right/left of another. Algorithmically, we evaluate if $(Bc_x > Ac_x)$ AND $((Bc_x - Ac_x) / (Ax_{max} - Bx_{min}) > 3)$. If yes, then we still form the aforementioned sentences. If not, we discard this relationship and do not create any sentence. To continue, we repeat the exact same process for the height axis (y) evaluating By_{min} , Ay_{max} , Bc_y , Ac_y respectively, forming sentences using the words “above” and “below” this time. As we mentioned previously, the longitudinal (depth) axis is treated differently. Given that both the RGB and the depth images have the same resolution, and we know the placement of the various objects in the RGB image, we use the same coordinates to sample their depth from the depth image. For simplicity, we consider the bounding box’s central pixel’s depth value representative of the whole object’s depth. To be on the safe side, we have set a threshold of at least 40 depth units [0-255] that two objects must differ in order to consider one in front of the other. So, we check if $depth(Ac_x, Ac_y) - depth(Bc_x, Bc_y) > 40$ to create the sentences “The <A> is in front of the ” and equally “The is behind the <A>”.

Another aspect we added to improve the quality of our dataset was the provision of changing the vocabulary depending on whether only one or multiple similar objects participate in a relation. In that manner, the typical response is only produced for two different-class objects. For images with only one apparent object, the template caption is “there is a <A>.”, while for relations involving two or more similar objects it is “{two, three, four} <A(plural)> {in a row, side by side, one on top of the other}.” depending on which axis does this relation regards. Note that, for all images we chose to consider only the four largest in area detected objects, and that sentences produced begin with a capital letter and end with a full stop.

Apart from the image spatial captioning task, in order to also enable the visual spatial question answering task we needed to create pairs of space-related questions & answers. An automated way to do that is to decompose each one of the captions into two self-evident questions & answers about the “What”, the “Where”, and the “How many” of the described relationship. The sentence-template may vary, depending on the following scenarios with example forms.

- For scenes with a single object, we can derive:

Question: What is there?	Answer: A cat.
Question: How many cats are there?	Answer: One.
Question: How many objects are there?	Answer: One.

- Respectively, for scenes with more than one objects and no same items we can derive (and more, extrapolating the same patterns):

Question: What is there?	Answer: A cat, a dog and a car.
Question: How many cats are there?	Answer: One.
Question: How many cars are there?	Answer: One.
Question: How many objects are there?	Answer: Three.
Question: What is above the car?	Answer: The cat.
Question: Where is the car?	Answer: Below the cat.

- Finally, for scenes with multiple instances of the same object, the scenario above deviates to (and more, extrapolating the same patterns):

Question: What is there?	Answer: A cat, a dog and two cars.
Question: How many cats are there?	Answer: One.
Question: How many cars are there?	Answer: Two.

Question: How many objects are there?	Answer: Four.
Question: What is above the car?	Answer: The cat.
Question: Where is the car?	Answer: Below the cat.

In this manner, we can create a huge number of spatial captions, questions and answers. Using the objects' labels, their bounding boxes, centres and the corresponding depth image, the automated mechanism can populate sentence-templates to create the set of captions and questions and answers we require. For instance, using the RGB image and the depth image presented in Figure 3.33, descriptions such as “The person is to the right of the fire hydrant” and “The fire hydrant is in front of the person” can be generated.

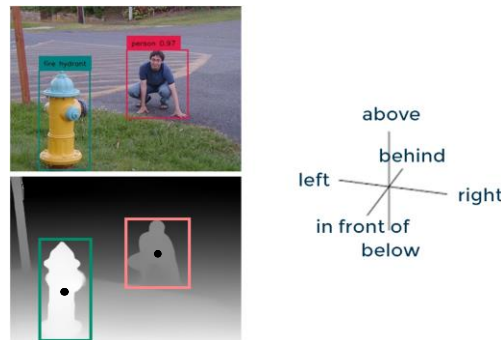


Figure 3.33. An example from the spatial dataset with its accompanying depth and object-annotated images

To give an estimate of the complete dataset's extend, it consists of 179,170 RGB images, each one accompanied by its respective depth-map, object-masks and object-annotated duplicate, together with tens of millions of spatial captions, questions and answers; all packed in about 283 GB of storage. The detailed dataset specifications are presented in Table 3.15 and Figure 3.34.

Table 3.15: Our COCO spatial extension dataset specifications

Dataset split	Disk Size (GB)	Number of images	Number of spatial captions	Number of spatial question-answer pairs
train	114.9	72,784	314,159	989,350
validation	55.9	35,359	149,487	472,004
test	112	71,027	294,422	933,211

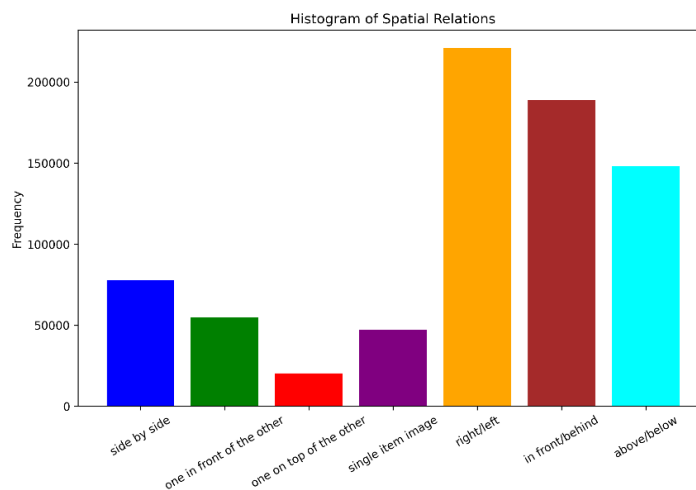


Figure 3.34: A histogram showcasing the frequency of all possible spatial relations occurring between objects in the dataset

Model architecture

Stepping on the theoretical foundations described in the project's GA, we first looked for potential solutions into a couple-years-old prominent family of VL pre-training encoders. We studied and experimented with works like LXMERT [Tan and Bansal, 2019], VisualBERT [Li et al., 2019] and UNITER [Chane et al., 2020]. Although these works have shown, for both IC and VQA tasks, remarkable performance in various benchmarks, the language output they generate is limited and shallow. It was soon realized that in order to generate more natural and richer textual representations we needed to migrate to a Vision Encoder-Decoder scheme. This versatile family of solutions can generate sophisticated and natural-sounding language that is currently considered the state-of-the-art of VL modelling.

Building on the Hugging Face's vision encoder-decoder model⁹, we experimented with various combinations of pre-trained transformer-based vision models as the encoder (i.e. ViT [Dosovitskiy et al., 2020], BEiT [Bao et al., 2021], DeiT [Touvron et al., 2020], and Swin [Liu et al., 2021]) and pre-trained transformer-based language models as the decoder (i.e. RoBERTa [Liu et al., 2019], GPT2 [Radford et al., 2019]). The effectiveness of initializing image-to-text-sequence models with pretrained checkpoints has been shown in various works like [Li et al., 2021]. Our SOTA experiments were conducted in order to enhance our understanding of the Hugging Face's framework and decide which combination of vision encoder and language decoder we should train to cover our needs. We need to mention here that the SOTA pre-trained models we tested were mostly trained on the standard COCO dataset and thus showcase no notion of spatiality in their generated captions and question answers. So, we chose to train our selected encoder-decoder combination on our spatially aware extension of the COCO dataset, from scratch (random initial model weights). Of course, when our pre-training is done, we will finetune our model to our specific use cases in Transfer-Learning manner [Zhuang et al., 2019]. Judging from the pre-trained versions' performance and weighting in our hardware limitations, mainly in GPU memory, we opted to develop a combination of ViT for the vision encoder and GPT2 for the language decoder.

Beginning from the encoder, ViT processes the input image by dividing it into fixed-size (16x16) non-overlapping patches (input image dimensions: 224x224). In ViT's context, the patch is treated as the unit for further processing. Each patch is linearly embedded to obtain a patch-embedding. The idea is to capture local information within each one of these, while enabling interactions between them in the subsequent transformer layers. To achieve that, the

⁹ https://huggingface.co/docs/transformers/v4.36.1/en/model_doc/vision-encoder-decoder#transformers.VisionEncoderDecoderModel

model needs to introduce positional embeddings to capture the spatial relationships between the different patches of the image. Unlike sequence data where the order of tokens is inherent, images lack a natural ordering of pixels, and therefore, positional information needs to be injected explicitly. In ViT, the patch's position is represented by a unique positional embedding vector created using a combination of trigonometric functions (sine and cosine), ensuring that the positional embeddings have a smooth and continuous pattern. These positional embeddings, added to the patch embeddings, form the input sequence for the transformer encoder. Then, the transformer-based encoder attends to the relationships between the different patches, capturing both local and global contextual information. The final output of the encoder is a set of contextualized embeddings for each patch, representing a deep and rich understanding of the image content.

As regards the decoder, the GPT-2 is a language model designed for generating coherent and contextually relevant sequences of language tokens. In this case, it takes as input the output embeddings from the ViT encoder and fuses them with its own contextual understanding of language during the decoding process in an autoregressive manner, generating one token at a time conditioned on the previously generated tokens. It uses a transformer-based architecture, similar to the original GPT-2 language model, consisting of multiple layers of transformer decoder blocks. Each block contains a multi-head self-attention mechanism and a feedforward neural network. The self-attention mechanism allows the model to attend to different parts of the input sequence, capturing dependencies and context. Similar to the ViT encoder, the GPT-2 decoder utilizes positional embeddings to provide information about the sequential order of tokens in the generated output. While the ViT uses sinusoidal positional embeddings, the GPT-2 decoder uses learned positional embeddings. During the decoding process, the GPT-2 decoder attends to the encoded visual information provided by the ViT encoder. The attention mechanism allows the model to focus on relevant parts of the image content when generating each token in the output sequence.

Note that the architectural analysis given above is indicative to the image captioning task. Slight adaptations in architecture and, of course, training it with different data (pairs of questions and answers) is needed to adapt our model to the visual question answering task. For VQA, the model takes both an image and a textual question as input. The image is processed by the ViT encoder to obtain visual embeddings and the textual question is typically encoded separately, either using a transformer-based architecture or other encoding mechanisms. The ViT-encoded visual embeddings and the encoded question must be combined to form a joint representation. This integration may involve concatenation, attention mechanisms or other fusion techniques to effectively merge the two modalities. This joint representation is then used as input to the GPT-2 decoder. We will adapt our model for use in the VQA task in the near future. To give a perspective of the developed model's size, it incorporates about 239.000.000 parameters. Figure 3.35 illustrates our architecture.

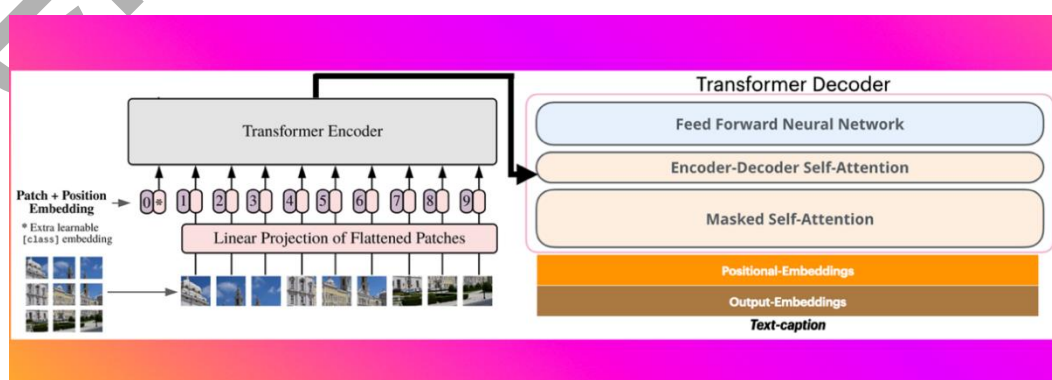


Figure 3.35. Our architecture with ViT encoder and GPT2 decoder

Training and results

Objective of the training process is to minimize the difference between the predicted caption generated by the GPT-2 decoder and the ground truth caption associated with the input image. This involves using a suitable loss function, in this case the cross-entropy loss, to measure the dissimilarity between the predicted and actual sequences of tokens. We need to mention here that, like in many other cases, this loss function (albeit differentiable) may not be very good at capturing the task's nuances and correlate highly with the better-suited metrics (like OURS). The optimization of the loss function is being handled by AdamW [Loshchilov & Hutter, 2017] with a linearly diminishing Learning Rate (LR) as the epochs progress. For training, we used 64.172 samples and for testing 31.156. Periodically, every 10 epochs, we evaluate our model on 10 metrics. The standard BLEU, ROUGE-1, ROUGE-2, ROUGE-L and ROUGE-L_{sum}, and our own devised metric in all its variations OURS_{1-A}, OURS_{2-A}, OURS_{3-A}, OURS_{MAX-A} and OURS_{MAX-B}. In Table 3.16, we present in detail the results of our model's evaluation study across a training session of 50 epochs, that took around 100 hours of processing on 2x Nvidia RTX 3090 24GB GPUs.

Table 3.16: Spatial IC model's evaluation results across training epochs

Evaluation Metric	Training Epoch				
	10	20	30	40	50
BLEU	14.99	15.16	15.15	15.21	15.33
ROUGE-1	51.83	51.96	51.95	51.96	52.02
ROUGE-2	45.70	45.94	45.97	46.01	46.07
ROUGE-L	49.96	50.14	50.18	50.17	50.22
ROUGE-Lsum	51.83	51.96	51.94	51.96	52.02
OURS1-A	0.430	0.590	0.600	0.590	0.620
OURS2-A	0.310	0.470	0.480	0.470	0.500
OURS3-A	0.180	0.340	0.350	0.340	0.370
OURS _{MAX-A}	0.177	0.337	0.347	0.337	0.367
OURS _{MAX-B}	0.156	0.316	0.326	0.316	0.346
Training Loss	0.008	0.0017	0.005	0.002	0.001

As it is evident from Figure 3.36, the model keeps improving until the designated completion of the experiment. Note that, this training session's length is chosen for the purposes of this deliverable document. The actual model to be released will be left to train for a couple-hundreds of epochs. Delving into the details, we observe the typical phenomenon where the model showcases a rapid initial improvement until around 20 epochs (both regarding the train-set loss and the test-set evaluation metrics), and then the progress slows-down, partially due the gradual decrease of the Learning Rate by the linear scheduler. Interesting is that, although the loss value fluctuates quite a lot between 20 and 40 epochs into the training, the evaluation metrics do not follow this trend, remaining almost constant for these 20 epochs and then improving further, together with the further decrease of the loss (possibly escaping out of a local minimum), hinting at the healthy behaviour of the session.

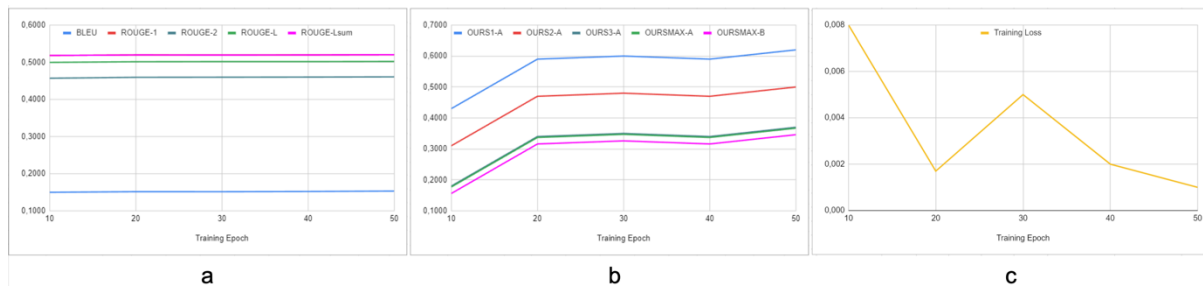


Figure 3.36. Spatial IC model progression across training epochs in a) standard evaluation metrics, b) variants of OURS metric and c) the respective training loss

Decomposing the evaluation metrics' behaviour, we observe that all the 5 standard and the 5 variants of OURS improve across training, though to different extend. The BLEU, ROUGE-1, ROUGE-2, ROUGE-L and ROUGE-L_{sum} metrics improve marginally, while OURS_{1-A}, OURS_{2-A}, OURS_{3-A}, OURS_{MAX-A} and OURS_{MAX-B} improve at a healthier rate. This is particularly aligned with our manual observations of the generated captions that, indeed, seem to improve. As regards the behaviour of the different variants of OURS metric, we conclude that they follow their theorized trend. As it was expected, the OURS_{1-A} variant achieves the highest accuracy of all (62% at epoch 50), since it evaluates only the first-generated caption of the model for each image. It is reasonable to expect that the more words we pump out of the model the more it will fail and, consequently, all other variants of OURS may showcase lower accuracy. This can be attributed to the autoregressive nature of GPT-2 where each new word is conditioned on the previous ones, and thus any errors accumulate leading to ever-increasing chance of generating nonsense. Indeed, as the results suggest, each other variant of OURS that considers more sentences (2, 3 and the maximum the model can produce for an image - the *MAX_i*) shows lower accuracy. Interestingly, we observe that past the first 3 generated captions per image considered, the accuracy of the model seems to not drop further, rather stays constant.

Furthermore, to study the difference between our proposed metric and the conventional metrics of the field and the potential usefulness of our metric, we conducted a small study examining all of its variants' correlation with the rest of the metrics, as well as the correlation between all of them for the sake of completeness. What can immediately be noticed is the highly positive correlation of all the standard metrics with each other. We can safely assume that, at least in our specialized dataset, good score in either of the BLEU, ROUGE-1, ROUGE-2, ROUGE-L and ROUGE-L_{sum} metrics would possibly translate to a good evaluation score in any of them. This behaviour is supported in literature, and is reasonable, given that they were all developed to evaluate similar linguistic scenarios. At the same time, all the variants of OURS metric showcase strong positive correlation with each other as well, behaviour equally reasonable given their similar workings. Reaching to the very core of our study, the most valuable finding and primal motive behind its conduction is that the OURS metric, in all of its variants, behaves independently to the rest of the existing established metrics. This fact strongly supports that, at least in the context of our specialized needs, the OURS metric succeeds in capturing a facet of the performance envelop which evades the rest of the well-known metrics; ultimately justifying its creation. Table 3.17 presents these results, where green coloured cells indicate positive correlation, red coloured cells indicate negative, and the white tones in-between suggest independence between the examined metrics.

Table 3.17. Pearson correlation coefficients between all IC metrics after 50 epochs

	BLEU	ROUGE -1	ROUGE -2	ROUGE -L	ROUGE -Lsum	OURS 1-A	OURS 2-A	OURS 3-A	OURS MAX-A	OURS MAX-B
BLEU		0,970	0,953	0,916	0,976	0,019	0,375	0,384	0,327	0,342
ROUGE-1	0,970		0,979	0,969	0,998	-0,081	0,210	0,185	0,227	0,305
ROUGE-2	0,953	0,979		0,989	0,971	-0,242	0,201	0,165	0,106	0,248
ROUGE-L	0,916	0,969	0,989		0,954	-0,252	0,182	0,113	0,112	0,305
ROUGE-Lsum	0,976	0,998	0,971	0,954		-0,063	0,204	0,194	0,228	0,276
OURS1-A	0,019	-0,081	-0,242	-0,252	-0,063		0,593	0,628	0,897	0,618
OURS2-A	0,375	0,210	0,201	0,182	0,204	0,593		0,970	0,829	0,823
OURS3-A	0,384	0,185	0,165	0,113	0,194	0,628	0,970		0,802	0,689
OURSMAX-A	0,327	0,227	0,106	0,112	0,228	0,897	0,829	0,802		0,875
OURSMAX-B	0,342	0,305	0,248	0,305	0,276	0,618	0,823	0,875	0,875	

In the second phase of VOXReality, as we were approaching the second round of pilots, amidst the Open Call sub-project developments, we decided to invest some more time in an effort to improve the quality of our spatial IC model, rendering it more useful for all purposes. After heavy experimentation with the training parameters and extending our training session from 50 to 130 epochs, our efforts yielded mixed results. Although for the majority of metrics performance improved drastically, two important for real-world applications variants of OURS metric, the 1-A and 2-A (which consider only the first 1 or 2 spatial captions of the input image), indicated a drop; behaviour that could be indicative of overfitting. The detailed results of our extensive training session are presented in Table 3.18.

Table 3.18: spatial IC model's evaluation results after epoch extension

Epochs	BLEU	ROUGE				MET EOR	OURS				
		1	2	L	Lsum		1-A	2-A	3-A	MAX-A	MAX-B
130	70.7	79.8	71.1	78.1	78.1	78.2	0.530	0.481	0.451	0.404	0.480

3.3.4. Image spatial visual question answering

Visual Question Answering (VQA) is an active area of research, evolving from early rule-based systems to modern deep learning approaches. Traditional VQA models rely on convolutional neural networks (CNNs) for visual feature extraction and recurrent networks or transformers for language understanding. More recently, fully transformer-based architectures have dominated the field. Models such as the LXMERT that pioneered two-stream processing of visual and textual inputs enabling improved contextual understanding arose, acting as fertile ground for even newer models like BLIP, that building upon these advances leverages bootstrapped pretraining on large-scale image-text datasets, significantly enhancing zero-shot and few-shot VQA capabilities. While most VQA models excel in general question-answering tasks, spatial reasoning remains challenging. Our work, using BLIP as the architecture of preference and our COCO spatial extension as foundational dataset, attempts to delve into the spatial visual question-answering problem, offering credible solutions to our VOXReality partners and beyond. Figure 3.37 presents an example image and the boxes below demonstrates the difference between generic and spatial answers in VQA using this image.



Figure 3.37. An example image for the generic and spatial VQA comparison

Question:	<i>"What is there?"</i>
Generic answer:	<i>"A man playing baseball in a field."</i>
Spatial answer:	<i>"A baseball bat to the left of a person."</i>

Question:	<i>"Where is the person?"</i>
Generic answer:	<i>"In a field."</i>
Spatial answer:	<i>"To the right of a baseball bat."</i>

Model architecture

The BLIP (Bootstrapped Language-Image Pretraining) [Li et al., 2022] architecture is a deep learning framework designed to tackle all prominent VL tasks, among them the Visual Question Answering (VQA) task. VQA requires a model to process an image and a corresponding natural language question and generate a coherent and contextually relevant answer. We trained a BLIP-based model with a spatial-reasoning focus to enhance its understanding of object relationships, positions, and spatial semantics within images.

The BLIP model, as illustrated in Figure 3.38, when configured for the visual question answering task, consists of three main components: A vision encoder, a text encoder and a language decoder. The vision encoder, which is in fact a Vision Transformer (ViT), processes the image by breaking it into small cells (patches) and extracting for each one high-dimensional meaningful features. The text encoder is a transformer-based model (similar to BERT) which processes the question and learns how words relate to each other. Finally, the language decoder generates an answer in a free-form text style, predicting words one-by-one (autoregressively).

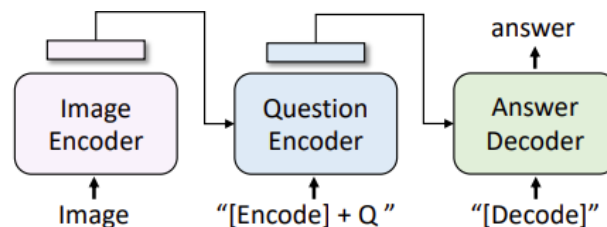


Figure 3.38. The BLIP model architecture for VQA

Dataset, training and results

As described in Section 3.3.3, our COCO dataset spatial extension comes with millions of question-answer pairs. We utilized a composite training strategy, different for the various modules of our model. In early stages of training, we kept the generically pre-trained vision

encoder frozen, so it retains knowledge of generic visual patterns before focusing on the nuances of our dataset, later on, once linguistic patterns and cross-modal connections have roughly developed. As training progresses, we gradually expose vision encoder to training, converging its individual learning rate with that of the rest of the network.

To evaluate our trained model, we used a subset of our very large dataset comprising of 1000 images and 9020 spatial questions. As we see in Table 3.19, our model manages to answer correctly around 60% of the designated spatial questions.

Table 3.19. Spatial VQA model's evaluation results across training epochs

Training epoch	OURS VQA-A	OURS VQA-B
1	0.554	0.645
2	0.574	0.662
3	0.574	0.651
4	0.604	0.684
5	0.576	0.661

3.3.5. Video spatial captioning

Apart from still images, we also put effort in transferring the same spatial-captioning functionality to the video domain on the condition that we would treat video fundamentally as a spatiotemporal medium and not merely as a sequence of images. For this purpose, we chose to base our model on the well-known SwinBERT architecture [Lin et al., 2021].

Model architecture

The system's design relies on a multi-modal Transformer that can model both visual information and textual information using a sparse attention map connecting together visual patch tokens and word tokens. The visual modelling part is based on the Swin Transformer [Liu et al., 2021]. A video sequence is modelled as a tensor of size $S \times W \times H \times 3$, where S is the video sequence length, W and H are the width and height of the video frames, and 3 corresponds to the three components of the RGB colour space. For video captioning systems, researchers have been long asking how many frames are necessary to be sampled from the video frame sequence in order to reliably model the visual information with textual descriptions. Authors of SwinBERT conducted a study to measure the performance of the model in terms of the frame sampling frequency, in a typical 30 fps video. Among T in $\{2, 4, 8, 16, 32, 64\}$, they found out that video captioning performance is optimal at $T = 32$, when the attention mask is sparse and learnable. As for the internals, SwinBERT has a visual and a textual learning component to map video sequences to word sequences. At first, to extract visual features, the authors used a VidSwin Transformer [Liu et al., 2021] in order to extract a tensor $(T/2) \times (H/32) \times (W/32) \times 8C$, where C is the channel dimension. In total, N textual tokens ($N = 50$) and M ($M = (T/2) \times (W/32) \times (H/32)$) visual tokens are extracted, forming a $(M+N) \times (M+N)$ sparse attention mask, where the top-left $M \times M$ block represents the relationships of visual tokens, and the bottom-right $N \times N$ block represents textual token relationships. The rest of two blocks represent the cross-associations of visual tokens and textual tokens. The authors used this learnable cross-attention mask and used a sparsity regularizer: $\lambda \sum_{i,j} \{M\} \{N\} |V_{i,j}|$. Using this regularization term, SwinBERT is forced to form only necessary associations between tokens, instead of very dense relationships that would be otherwise uncovered if a dense attention map was used. As for the linguistic component of SwinBERT, instead, it relies on Masked Language Modeling (MLM), which masks out several word tokens and asks the multimodal Transformer model to predict the masked tokens.

Dataset

In order to be useful for the VOXReality end-applications, being either for the project Pilots or the Open Call (OC) products, hard request of our solution was to be responsive and provide answers to requests after only a handful of seconds (given reasonably adept hardware). Thus, we deemed important to design for short length, approximately between 2.5 and 4 seconds, videos. Such a video is long enough for the model to extract all its useful temporal information, while at the same time being short enough to be reliably captioned with only one sentence, given that it is expected to contain only a single scene. Having that in mind, we began drafting the extraction of our train dataset. Drawing inspiration from our image spatial captioning dataset, our idea was to acquire the videos from a well-known generically captioned dataset and using some heuristics to automatically create for these spatial captions to be used for our training. We chose to use the videos from the Microsoft MSR-VTT dataset [Xu et al., 2016], containing 10.000 videos. These videos often last quite longer than our set requirement and may also contain more than one shot, breaking our assumptions. Thus, we devised the following heuristic mechanism to cut the videos into as many useful smaller ones as we could, effectively multiplying the number of samples we could train on.

First, we used the TransNet v2 [Souček et al., 2008] video-shot segmentation algorithm, to cut the original videos into their composing scenes, so that no video contains more than one. After this step, we cut each resulted video into equal-duration smaller (2.5 and 4 seconds long) samples. With this mechanism, we successfully created 39.963 short videos usable for training and evaluation. As for their textual counterparts, after extensive experimentation with prompt-engineering, we decided to just re-use our own spatial captioning mechanism that was used to create our COCO-dataset spatial extension and be applied to the middle frame of the video. In this way, we managed to gather quite sensible spatial captions for a very good portion of the videos (we randomly supervised 150 samples and they were humanly acceptable to correctly describe the video to a rate of about 85%). After our training was completed and the model was tested with current, real-life videos, we realized that our model is somewhat biased towards the “human” class, since it is evident that the vast majority of the videos used for training contained some human actions (TV shows etc.). Figure 3.39 presents an example frame, extracted from a 4-second clip.



Figure 3.39. Example frame captioned as: “A laptop is to the right of a person”

Training and results

Table 3.20 presents the training parameters. Various learnable components of the model have their own learning rate. The “Backbone coefficient learning rate” is a hyperparameter controlling the relation between the overall learning rate and the VidSwin’s one. The “Loss sparsity w ” is a sparsity-loss regularization term which controls the magnitude of the sum $\sum_{i=1}^M \sum_{j=1}^M |V_{i,j}|$. It controls how much the model should pay attention to forming sparse interconnections among textual-visual tokens.

Table 3.20. The hyper-parameters of the spatial VC model

Hyper-parameter	Value
Backbone coefficient learning rate	0.05
Loss sparsity w	0.5
Optimizer	Adam
Learning Rate	4e-4
LR Scheduler	Linear
Batch Size	6

Evaluating the trained model in unseen data suggests that the model is capable of predicting the correct spatial relations between apparent items of the scene at about 62-76% of the cases depending on the accuracy envelope captured by OURS metrics (presented in Table 3.21).

Table 3.21. Spatial VC model's evaluation results across training epochs

Training epochs	OURS1-A	OURSMAX-B	BLEU	METEOR	ROUGE-L
5	0.744	0.610	0.357	0.303	0.775
10	0.751	0.617	0.358	0.305	0.784
15	0.758	0.622	0.364	0.311	0.790

3.3.6. Visual navigation

Apart from the typical Vision-Language tasks, we took the chance to explore further on non-standard VL-models' applications. In that frame, we proposed the repurposing of a VQA architecture for generating turn-by-turn navigation instructions towards a destination. In respect with the standard VQA scheme, we consider the input image to be a capture from the user's point of view and the input question to be the desired destination room. The whole concept was developed and tested for the VR Conferences use case, where the user can ask our conversational agent for navigation instructions.

Dataset

The navigation model was developed for use inside the VR Conferences – Pilot 2 virtual environment. The VR space consists of a conference hall premises spanning across seven different areas/rooms, connected through a lobby. When wandering around in the lobby area, the user can ask the conversational agent (CA) for directions on how to go from the current position to the entrance of one of the rooms.

In order to train and evaluate in unseen data our intended model, we created a dataset of image captures from the viewpoint of the user (from their "virtual eyes"), each one accompanied with seven pairs of question-answer that correspond to the name of the destination-room (question) and the turn-by-turn navigation instructions on how to go there (answer). Figure 3.40 presents an example capture from this datasets including the questions and answers in the form of the Dijkstra algorithm's output denoted with A for each destination-room denoted with Q.

We sampled the virtual premises on a grid with nodes being 0.5m apart in both Z (length) and X (width) axes. The height of the camera (Y axis) was left constant at 1.7m. On each node, we captured one image for every 20° of rotation. To promote generalization of the model, we added uniform-random noise in [-15°, 15°] to the camera pitch (look up/look down) and in [-10°, 10°] to the camera roll (tilt left/tilt right). Figure 3.41 presents an example capture and its corresponding questions and answers from the dataset with a tilted view introduced as noise. As regards the textual part of the dataset, based on the Dijkstra¹⁰ algorithm, we devised a

¹⁰ https://www.w3schools.com/dsa/dsa_algo_graphs_dijkstra.php

heuristic mechanism to automatically create turn-by-turn navigation instructions from every node to each one of the destination-rooms. More on the instructions-generation heuristics can be found at VOXReality deliverable D5.2. This process resulted in a 405,559 train and 2.044 validation-sampled dataset which was utilized to train and evaluate our visual navigation model.

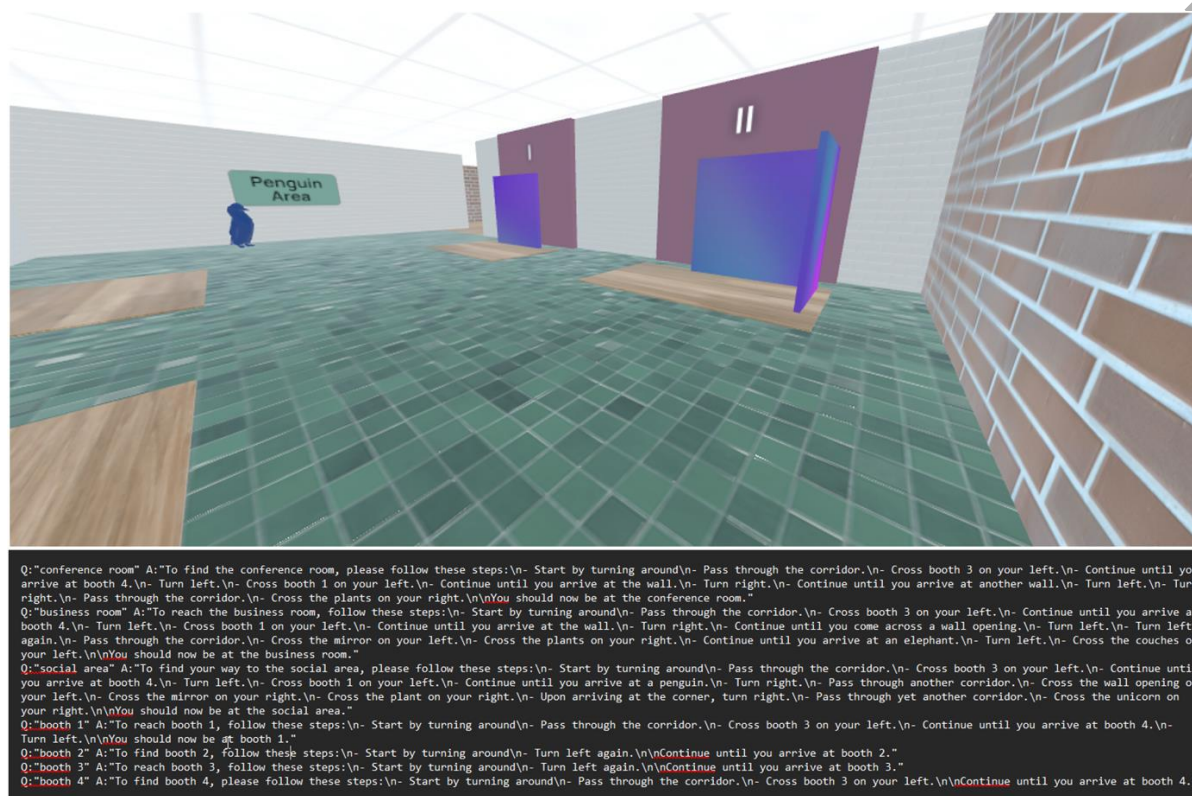


Figure 3.40. An example capture from the virtual space



Figure 3.41. An example capture from the virtual space with noise (tilted view)

Training and results

For the fitness/quality of the generated navigation instructions, in the absence of a truly suitable metric, we chose to employ the standard BLEU metric with setting the bar high enough, so that the acceptable instructions could only be near-perfect matches to the ground truth. In that context, the 385M parameter BLIP VQA base-model was set to be trained for as long as it takes, with both pre-trained and randomly initialized weights, until the model achieves an average validation BLEU score 89 (rounded to nearest integer, evaluated right after every training epoch). Table 3.22 presents the parameters used for training.

Table 3.22. The hyper-parameters of the VN model

Hyper-parameter	Value
GPU	Nvidia RTX 4090 24GB
Weight decay	0
Optimizer	AdamW
Learning Rate	8e-6
LR Scheduler	Cosine Annealing (T=30)
Batch Size	21

In the context of this endeavour, we also test the time to fine-tune a pretrained model to a specific scene and the resource requirements using models that need small finetuning instead of training from scratch. The results of our two training sessions, presented in Table 3.23, indicate that we successfully addressed both of them. In practical terms, a mean BLEU of ~89 suggests that in most of the time the model response is almost identical to the reference one. According to Figure 3.42 the trained-from-scratch version of the model struggles to achieve similar performance with the pre-trained one, even after a significantly longer training session, highlighting the importance and benefits of finetuning.

Table 3.23. VN model’s evaluation with pre-trained and randomly initialized weights.

Model weights’ initial status	BLEU	Training time (hours)	Training epochs
Pre-trained on various VL VQA datasets	89.4	3	2
Randomly initialized	88.9	27	18

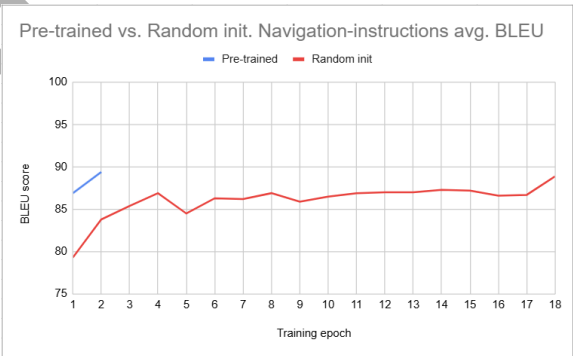


Figure 3.42. BLEU score of the pre-trained vs randomly initialized model weights

3.3.7. Alternative image spatial captioning

Apart from relying on the ViT encoder of our standard solution for learning the visual characteristics of the input image, with no sense of spatiality itself, and relying on the textual decoder to learn and produce spatially aware descriptions, we decided to additionally explore another, more intriguing idea. Core of our concept is to embed in the encoder’s



representations themselves, bits and pieces -genes- of its respective object spatial identity, in such a way that we could then train a classifier to classify / recognize any pairs of objects spatial relationship by simply inspecting them.

Model architecture

The complete model is a composition of two sequential parts trained in an end-to-end fashion. We base the first in sequence component of our architecture on the SegFormer segmentation architecture [Xie et al., 2021], which is actually a vision transformer curated for dense object segmentation (pixel-by-pixel object-class classification). In our setup, it accepts as input an RGB image which is scaled, preserving its aspect ratio, to 512 pixels for its long edge. In principle, the model processes the image on a 128x128 grid where for each cell it produces two vectors 81D and 768D respectively. The former one is the “logits” which contains the probability distribution of the cell to belong to an object-class (81 object classes) and the latter one is the semantic descriptor - “token” of the cell, which encodes the contents of the cell. By choosing for each cell its most probable object-class, per its logits, we can form clusters of cells that appear to belong to the same object. By applying mean pooling to each such cluster, we can have a single 768D deep-representation of each object of the image.

After mean-tokens for each object have been calculated, we concatenate them in pairs representing all possible 1-1 relations. This new set of feature vectors (Number of 1-1 relations x 1536D) is fed to the second module of our model, the spatial classifier. The classifier is responsible for accepting the concatenated tokens of two objects in the scene and predicting their spatial relation in all three axes (Left/Right/Uncertain, Back/Front/Uncertain and Above/Below/ Uncertain). The classifier itself consists of two fully connected hidden layers, each sized smaller than its feeding one in encoder-fashion. We experimented heavily on the ratio between the sizes of each subsequent block, with 1/3 yielding the best results. Note that this second hidden layer is actually split into three - no weights sharing - parallel components, one for each axis. The model architecture described can be seen in Figure 3.43.

In its normal operation, SegFormer’s tokens are shaped through training to encode information useful for segmenting the objects in the image. In our case, via conducting end-to-end training that involves minimizing a composite loss containing terms for both the segmentation and classification tasks, we drive SegFormer to embed useful to the classifiers spatial information deep into its semantic descriptors, fulfilling our initial goal to endow the visual encoder with the primitives of spatial reasoning.

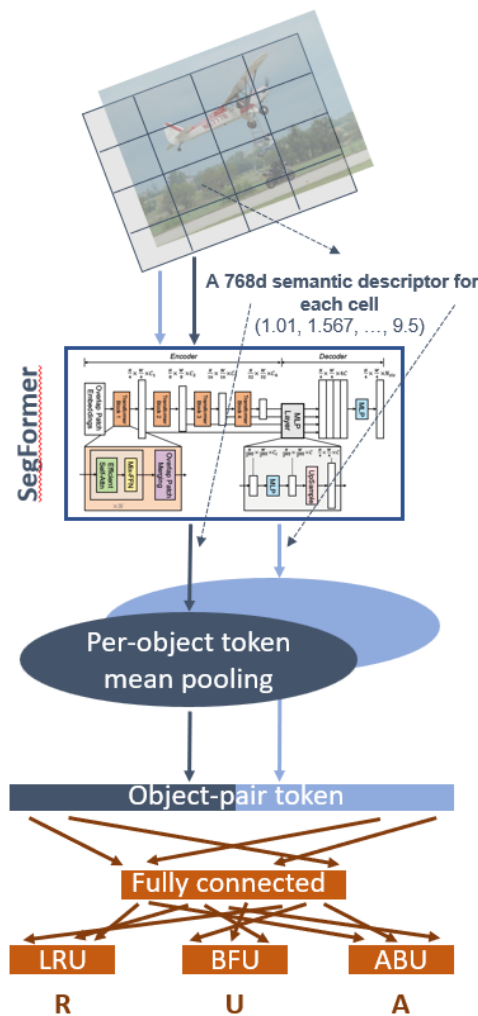


Figure 3.43. Our alternative IC model architecture

Dataset, training and results

The dataset used to develop on, and explore this idea was our COCO spatial extension dataset (a large portion of it, counting 90.000 train and 40.000 validation images) and the models were finetuned using the parameters presented in Table 3.24.

Table 3.24. The hyper-parameters of the alternative IC model

Hyper-parameter	Value
Weight decay	4e-4
Optimizer	Adam
Learning Rate	1e-5
LR Scheduler	Step (epochs=240, n_steps=3, steps_ratio=0.5)
Batch Size	17

For our model to properly function, it actually needs to perform well in two different sub-tasks. At first level, it needs to segment accurately the objects found in the scene. The more items it finds, the more of the apparent relations in the image it has the basis to correctly classify. Moreover, the better and crisper the object segmentation is, the better quality the extracted feature vector is expected to be, increasing also the possibility to successfully classify the spatial relations of the objects. In that sense, we divided the evaluation results into two tables.

Table 3.24 presents the model's performance in the first sub-task, that is the object segmentation. The metrics include the standard pixel-label accuracy and mean intersection over union (IoU) calculations, in both micro and macro integration modes, since for our heavily

class-imbalanced data, both are required for proper interpretation. “Seg miou per img” is the average mean IoU of all images in the dataset (we calculate mean IoU separately in each image and then average for all images). “Cla rels found” is the percentage of 1-1 object relations present in the ground truth of the image, the model was able to classify (because it recognized the involved objects and proceeded to the second task, no matter its result). In Table 3.25, we observe that the model manages to correctly assign each pixel to its belonging object-class (across 80 object classes, excluding the background class) in around 57% of the cases. This result is enough for the model to at least consider 70% of the existing 1-1 object relations in the dataset (it detected the existence of the involved objects). We consider the model to have under-performed in mean IoU terms, meaning that the average detected area of an object is frequently much bigger than the object itself (x4-5), hindering the correct classification of its relations due to the “noisy” selection of its feature vectors to be considered for classification.

Table 3.25. The alternative IC model's object-segmentation performance

Training epoch	Training steps	Seg acc micro	Seg acc macro	Seg miou micro	Seg miou macro	Seg miou per img	Cla rels found
40	39760	0.332	0.370	0.134	0.124	0.123	0.577
80	79520	0.323	0.365	0.136	0.123	0.123	0.552
120	119280	0.481	0.413	0.168	0.198	0.199	0.674
160	159040	0.501	0.399	0.172	0.214	0.215	0.674
200	198800	0.505	0.409	0.179	0.218	0.219	0.696
240	238560	0.566	0.415	0.193	0.259	0.262	0.704

Starting with this disadvantage, the model manages to correctly classify around 60% of the detected objects' 1-1 spatial relations. The Behind-Front axis, contrary to intuition and without explicitly using depth, seems somewhat easier to judge, as we see in Table 3.26, which presents the second sub-task, the spatial-relations classification, conducted only for the objects the model managed to identify. The metrics include the standard classification accuracy in both micro and macro integration modes, for all three axes (Left-Right, Above-Below, Behind-Front). Note that the model can also predict the “uncertain” class for a relation, meaning that for all results, random classification would approximate the 0.333 (1/3) threshold.

Table 3.26. The alternative IC model's spatial-relations classification performance

Training epoch	Training steps	Cla LR acc micro	Cla AB acc micro	Cla BF acc micro	Cla LR acc micro	Cla AB acc micro	Cla BF acc micro
40	39760	0.449	0.537	0.609	0.459	0.481	0.526
80	79520	0.454	0.537	0.610	0.464	0.479	0.542
120	119280	0.501	0.547	0.633	0.504	0.480	0.581
160	159040	0.525	0.559	0.630	0.526	0.497	0.580
200	198800	0.518	0.558	0.646	0.519	0.498	0.603
240	238560	0.569	0.568	0.650	0.570	0.503	0.603

3.4. Conversation Agents (CA)

In VOXReality, we aim to develop two conversation agents with two different tasks: 1) VR Conference Conversation Agent and 2) AR Training Assistant.

3.4.1. VR Conference Conversation Agent

The VOXReality VR Conference Conversation Agent is specifically designed to enhance the user experience at the conference, serving a multifaceted role. It is tasked with guiding users through the conference venue, providing essential information about the conference program, and offering insights into the trade show.

In this section we present the methodology followed for the implementation of this conversation agent. We explain in this section the creation of a robust Natural Language Understanding (NLU) model, creation of a dataset and extending it to the fine-tuning process of the model to ensure its effectiveness and accuracy in understanding and processing user queries, the mechanisms behind the conversation agent's workflow and the fine-tuning process applied specifically to the conference agent to provide accurate navigational assistance within the VR conference environment.

Natural Language Understanding (NLU) Model

The primary objective of the NLU model is to accurately determine whether a user's query is related to seeking navigation guidance within the conference venue, inquiring about the conference schedule, or requesting insights about the trade show. It's important to note that while this model plays a crucial role in the research phase, it is not directly integrated into the final workflow of the conference's conversation agent. However, it is planned to be made available for future open calls, offering broader applicability and utility.

Model Selection

In the dynamic landscape of NLU, the selection of an appropriate model is a critical decision that significantly impacts the success of language processing tasks. The T5 model, known for its exceptional adaptability and flexibility, has emerged as a leading choice for handling a wide array of NLU challenges. Its ability to seamlessly adjust to various linguistic tasks, coupled with a unified methodology, greatly eases the process of integrating it into different systems. This makes the T5 model not only highly efficient but also remarkably adaptable in addressing the diverse needs of NLU applications.

For the navigation assistant task, the model that was selected is the T5-Small. This choice was guided by the nature of our NLU task, which does not necessitate the advanced capabilities of more substantial models. The T5-Small perfectly meets our requirements, striking an optimal balance between performance and computational efficiency. It adeptly handles the language processing needs with the right level of complexity, ensuring both accuracy and resource economy. Thus, the T5-Small emerges as the most fitting solution for our NLU objectives, providing just the right mix of efficiency, affordability, and computational effectiveness.

Dataset Description

To develop a stable NLU model using T5-Small, the creation of a dataset tailored to specific use cases is essential. Following the dataset structure utilized by Convlab-3 [Zhu et al., 2022] for training, the Multiwoz [Budzianowski et al., 2018] dataset was beneficial for ensuring well-structured and appropriately formatted data. Each dataset sample starts with a user query, which is a question directed at the model. For example, a query might be: *"Can you find an Indian restaurant for me that is also in the town center?"* The corresponding label for this query is *"[inform][restaurant][area][center, food][Indian]"*. This labelling system denotes the general user intent, such as seeking information about a restaurant, and specific details like desired location and food preferences. A sample data from the NLU dataset is presented below:

```
{
  "context": "How can I go to the conference room?",
  "dialogue_acts_seq": "[request][direction]([destination][conference room])",
  "context": "Is there a company that is offering solutions for automation?",
  "dialogue_acts_seq": "[request][trade_show]([product_info][Automation])",
  "context": "What time does the 'Socializing' event ends?",
  "dialogue_acts_seq": "[request][program_info]([end_time][Socializing])",
  "context": "I need to get there. Can you guide me?",
  "dialogue_acts_seq": "[request][direction]([destination][unspecified])",
  "context": "What event will be held at 11:30?",
  "dialogue_acts_seq": "[request][program_info]([event_at_time][11:30])",
  "context": "Which company is present at booth B02?",
  "dialogue_acts_seq": "[request][trade_show]([exhibitor][B02])"
}
```

Inspired by the structure of Convlab-3's dataset, a new dataset is created with a similar format. This dataset is designed to train a model specifically for a conference assistance application, where it is expected to deliver navigational guidance within the conference venue, details about the conference program, and offer insights about the trade show, including information on the participants and their respective products. Each sample begins with the "context" field, containing the user's query. Following this is the "dialogue_acts_seq" field, which holds the corresponding label for the query. This label is crucial for the model to discern the user's initial intent, such as seeking directions, program details, or trade show information.

The structure of the rest of the label varies depending on this intent. For navigation-related queries, the label is designed to identify the user's desired destination. If the destination is not specified in the user's query, the label marks the destination as "*unspecified*". At this point, it should be mentioned that the destinations are unrelated to the rooms of the conference venue of our use case. This approach is adopted to provide the model with a generalized capability, ensuring it does not only learn to identify specific rooms but can adapt to a variety of scenarios. Similarly, for program information queries, the label is structured to extract all related details from the user's request. For instance, as shown in the example above under the program information category, if a user asks about when an event ends, the label includes "end_time" to indicate the user's interest in the event's end time and the event's name. The trade show examples are set up in a similar way to the program information, organizing labels to gather detailed information from questions about the trade show.

The dataset, comprising 498 human-annotated samples, is evenly distributed across three categories: navigation, program information, and trade show, with each category containing 166 distinct samples. This dataset is divided into training, validation, and test sets. For the training set, 100 samples from each category are selected, totaling 300 samples, while the test and validation set includes 99 samples. This approach not only ensures a balanced representation of each category within the dataset but also maintains the quality and relevance of the annotations. Detailed information of the data distribution is presented in Figure 3.44.

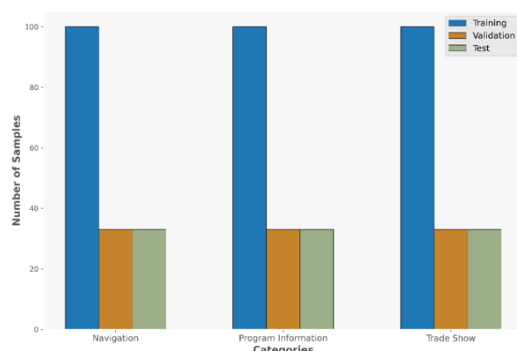


Figure 3.44. Data distribution across categories

Model Fine-tuning Process

The fine-tuning of the T5-small NLU model is carried out using the Convlab-3 platform, specifically designed to adapt different sizes of the T5 model to custom datasets. Convlab-3's robust framework provides an integrated solution for the entire fine-tuning process, ensuring that the model is optimally trained for the task.

The initial step in the model training process involves a comprehensive pre-processing phase, crucial for adapting the dataset for effective use with the T5 model. This phase includes tokenizing both the input queries and the target responses, vital for converting natural

language into a format suitable for the T5 model's processing. In addition to tokenization, the Convlab-3 platform efficiently manages the challenges of padding and truncation. Moreover, the platform also includes an evaluation process as part of the training pipeline, involving testing the model on unseen data to measure metrics such as the F1 score, which provides insights into the precision and recall of the model's predictions.

Experimental Results

To guarantee a structured and efficient training procedure, a step-by-step method is implemented. The model is first fine-tuned and evaluated using only the navigation category samples. This focused approach allows for detailed monitoring of the model's performance in a specific context. For the navigation category, 100 samples are used for training, with 33 samples each allocated to testing and validation. The results from this initial phase were highly encouraging, with the model achieving an F1 score of 90.4%. Following the success in the navigation category, the same fine-tuning process was applied to the program information and trade show categories. The model continued to exhibit strong performance, achieving F1 scores of 91.8% and 95% in these categories, respectively. This demonstrates the model's adaptability and effectiveness across different types of user queries.

In the final phase of training, samples from all categories are combined to provide a comprehensive dataset. This allows the model to be trained across a broader spectrum of queries, simulating more closely the variety it would encounter in real-world applications. The overall F1 score for the model, after being trained on this combined dataset is 95.2%. The detailed results of these training phases, illustrating the model's performance across different categories and in the comprehensive dataset, are presented in the Table 3.27.

Table 3.27. Experimental results of the NLU model

Category	Accuracy (%)	Precision (%)	Recall (%)	F1 score (%)
Navigation	89.9	90.4	90.4	90.4
Program information	88.4	90.3	93.3	91.8
Trade show	95	95	95	95
All Categories	94	94.3	96.1	95.2

Conversation Agent

For the Conference Dialogue Agent, the Llama2-Chat 7B model was selected. This model, is specifically designed for chat and dialogue scenarios, making it a natural fit for an agent focused on conversational interactions in a conference setting. The 7B version of the Llama2-Chat model stands out for its efficiency and effectiveness. It is particularly well-suited for environments where computational resources might be limited or for applications that need to be cost-effective. This is a crucial consideration for deployment in varied settings, where the availability of resources can significantly impact the functionality of the agent.

Another key aspect of the 7B model is its ability to provide fast response times. In a conference environment, where interactions are real-time and often require quick exchanges, the speed of the model is essential. The model is designed to maintain a fluid and natural conversation flow, a critical factor in ensuring user engagement and satisfaction. Despite its relatively smaller size compared to the other models in the Llama2-Chat range, the 7B version does not compromise on performance. It can handle the spectrum of queries and interactions, from basic informational requests to more complex dialogue. This model can efficiently manage basic to moderately complex dialogues, ensuring effective communication without the need for the computational resources demanded by larger models.



In summary, the Llama2-Chat 7B model is a well-rounded choice for the Conference Dialogue Agent. It strikes a balance between efficiency, speed, and performance, aligning well with the practical requirements of a conference environment. This model's capabilities ensure that the agent can handle a variety of interactions effectively, making it a suitable tool for enhancing the conference experience through intelligent and responsive communication.

Workflow Description

Designing the conference agent's workflow presented a significant challenge due to its complex, multi-stage response generation process. To address this, the LangChain¹¹ framework is utilized. LangChain is an innovative, open-source Python library that simplifies building applications with language models like GPT-3 from OpenAI. It employs "chains," which are sequences of modular components, enabling flexible and efficient construction of language-based applications. This framework stands out for its versatility, allowing for the integration of different processing stages, such as context management, response generation, and post-processing. Such modular design makes LangChain particularly advantageous for creating intricate language interactions and tailoring them to specific needs.

The workflow for the conference agent (Figure 3.45), is a multi-stage process. It consistently employs a single model, the Llama2, across its chains. Initially, a brief overview of each component is presented, followed by a more detailed explanation in the following subsections.

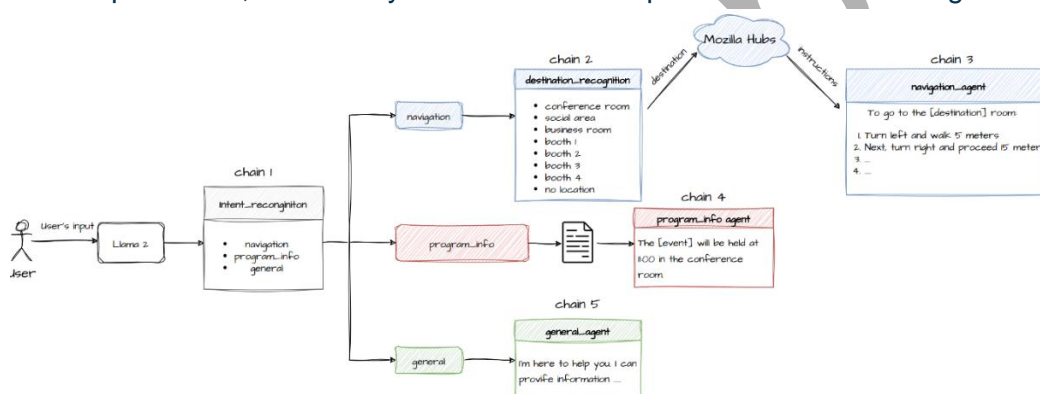


Figure 3.45. Conference agent workflow

The first chain employs the model to analyse user input, categorizing it into three distinct types: navigation, program information, and general queries. The 'navigation' type refers to user requests for directions to a specific room. In case that the model intent is navigation, the second chain follows that determines the user's intended destination. This step is crucial for integrating with Mozilla Hubs to identify the shortest path to the desired location. Once the shortest route is established, the third chain takes the navigation instructions and delivers the response to the user. In the second scenario, the 'program information' category involves queries about the conference schedule. In this case, Llama2 relates to a pdf file where the pdf file includes the program of the conference. In this case, chain 4 is activated by combining user's input and the information of the pdf file response to the user. Lastly, the 'general' category handles all other user inquiries not related to directions or program details.

Intent Recognition

When a user poses a question, the initial and crucial step for the agent is to accurately identify the user's input. This initial step is the key to ensuring that the subsequent process and response align correctly with the user's query. If the model misinterprets the user's intent, the resulting answer will likely be unrelated to the original question. To successfully build an intent recognition model using Llama 2, the capabilities offered by LangChain are utilized.

¹¹ <https://github.com/langchain-ai/langchain>

LangChain provides a specialized feature for developing prompt templates based on few-shot examples. This process involves creating several task-specific examples. These examples, coupled with a system prompt, enable the model to adeptly handle the desired task. In essence, this feature is provided by the *FewShotPromptTemplate* function, and it is a powerful tool for creating structured, example-based prompts that guide the model's responses in a specific and controlled manner, making it highly useful for complex querying and data processing tasks.

For the conference agent's intent recognition task, the model is designed to face a classification problem. Based on the user's input, it must categorize the query into one of three labels: 'navigation', 'program information', or 'general'. To assist the model in this process, specific examples were initially created to guide its comprehension and enhance its ability to accurately classify queries. A few examples for intent recognition task are presented below:

```
{  
  "input": "Where is the main auditorium located?",  
  "output": "navigation",  
},  
{  
  "input": "What time does the keynote speech start?",  
  "output": "program_information",  
},  
{  
  "input": "Hello how are you?",  
  "output": "general",  
},  
{  
  "input": "How do I get to the registration desk?",  
  "output": "navigation",  
},  
{  
  "input": "Who is the speaker for the afternoon session?",  
  "output": "program_information",  
},  
{  
  "input": "Can you help me find information about the conference?",  
  "output": "general",  
}
```

In addition to these examples, a carefully selected system prompt is employed to further guide the model's focus. This prompt was finalized after experimenting with various combinations to determine the most effective one and outlines the task for the model, ensuring that it focuses on categorizing the query into the correct label without providing extraneous details.

Destination Recognition

The approach used for intent recognition was similarly applied to destination recognition. In this scenario, the model's objective is to pinpoint the user's desired destination based on their input but also to have the ability to recognize if the requested destination is not part of the conference venue. Utilizing few-shot prompting examples, user inputs are paired with their corresponding correct outputs that the model is expected to generate. A crucial part in this process is the system prompt. Through a meticulous prompt engineering process aimed at identifying the most effective system prompt for this specific use case. A few examples for the destination recognition task are presented below:

```
{  
  "input": "Can you guide me to the booth 1?",  
  "output": "booth 1",  
},  
{  
  "input": "I want to go to the social area.",  
  "output": "social area",  
},  
{  
  "input": "I want to leave. Where is the exit?",  
  "output": "exit",  
},  
{  
  "input": "Where is the business room located?",  
  "output": "business room",  
},  
{  
  "input": "How do I get to the conference room?",  
  "output": "conference room",  
},  
{  
  "input": "How do I get to the gym?",  
  "output": "no location",  
}
```

In the prompt, the labels represent all the available locations within VOXReality's conference venue. It is designed to be adaptable, meaning that the prompt can be modified to suit different conference venues with their respective unique rooms. Additionally, the prompt specifies that if the user requests a destination not listed in the venue's specific configuration, the response should be classified as "no location." This feature ensures that the model remains adaptable and applicable across various settings, providing accurate and relevant responses based on the particular layout and offerings of each conference venue.

Navigation Agent

The navigation agent is a key component within the workflow of the conference dialogue agent. The agent handles two types of input: user queries and Mozilla Hubs' navigation instructions. To formulate an appropriate response, the model integrates these two inputs. Mozilla Hubs provides navigation instructions that are necessary for guiding users effectively. These instructions form the core of the navigation process, ensuring users reach their desired destinations accurately. The types of instructions are:

- **Turn around:** Typically, the initial instruction, this directive assists users in aligning themselves correctly at the onset of their journey. It ensures that the user starts from the correct orientation, which is crucial for the accuracy of subsequent directions.
- **Move [number] meters:** This command specifies a precise distance for the user to travel in a straight line. It provides clarity on how far to proceed before the next action, aiding in maintaining a clear and measured progression along the route.
- **Turn right/ left:** These instructions guide the user in making accurate turns. Whether a right or left turn, they play a critical role in keeping the user on the intended path and preventing deviations that could lead to confusion or longer routes.
- **Stairs up/ down:** These commands are particularly important for vertical navigation. They direct users when ascending or descending stairs, ensuring safe and efficient travel between different floors or levels within the environment.

The model's ability to seamlessly blend user queries with these structured navigation cues is what makes the navigation agent essential to the overall efficiency and effectiveness of the conference dialogue agent in a Mozilla Hubs environment.

To optimize a user-friendly agent's behaviour in handling navigation inputs, a specialized system prompt is implemented. This process, known as prompt engineering, involves rigorously testing various system prompts to ensure that the model generates the most effective responses. The chosen prompt aids the model in addressing situations where user instructions are either uncertain or involve requests for non-existent locations within the conference venue. It is especially helpful in guiding the model during scenarios where the user either fails to specify a destination or selects an unavailable one. In such instances, the prompt instructs the model to inform the user about the available rooms at the venue.

The navigation agent is a crucial component of the conference dialogue system, adeptly handling user queries and Mozilla Hubs' navigation instructions. Its advanced integration and prompt engineering process enable it to adapt to various conference venues, providing precise and relevant navigational assistance. This adaptability marks a significant advancement in virtual assistance, enhancing user experiences in complex virtual environments like Mozilla Hubs.

Program Information Agent

The program information agent functions as a part of a larger system designed to interpret and respond to user queries about conference program details. Initially, if the user's input is identified as program information from the intent recognition component, the program information agent activates.

Specifically, the agent is using the *Retrieval QA* function provided from Langchain. The agent is connected to a comprehensive database, typically a PDF file, which contains detailed information about the conference program. Using advanced search and retrieval techniques, the agent scans this document to locate the exact information relevant to the user's question. The strength of this system lies in its ability to parse large volumes of text data and extract

specific pieces of information. The agent, by understanding the context in which they appear, ensures that the responses are accurate and relevant to the user's query.

In line with the system's guiding prompt, the agent is directed to be a helpful navigator, focusing on delivering accurate, concise, and reliable information. This prompt ensures that responses are direct and contain only the necessary details asked by the users. The agent's primary mission is to assist users by providing detailed information about the schedule, speakers, events, locations, and other specific details included in the conference program.

This approach ensures that users receive quick, accurate, and contextually relevant answers to their questions about the conference program, enhancing their experience and making information access more efficient.

General Agent

The final element in the system's workflow is the general agent, a component designed to address a wide range of user inquiries that fall outside the specialized functions of other agents, such as navigation or specific conference program details. This agent comes into play when the intent recognition component classifies an inquiry as 'general'.

In its operation, guided by a carefully structured system prompt, the general agent's primary objective is to enhance the overall user experience. It achieves this by delivering clear, precise guidance, informing users about the types of assistance it can offer. This approach is strategically focused on ensuring that users receive assistance that is not only relevant but also precisely aligned with their immediate queries and requirements.

The essence of the general agent's role is to facilitate an interaction that is both efficient and user-friendly. By providing succinct, relevant responses, it helps users navigate their queries with ease, making the process of obtaining information seamless and straightforward. The general agent, with its broad scope and user-centric design, is therefore an indispensable part of the system, ensuring that users' general inquiries are addressed promptly and effectively.

Fine tuning of the navigation agent

While the model initially performed well in delivering navigation instructions, some problems appeared. Firstly, one of the primary concerns was the variability in answer format. This inconsistency not only creates a lack of uniformity across different responses but also potentially leads to confusion for users. For example, some answers may be presented in a detailed, step-by-step manner with clear numerical ordering, while others might be more narrative in style, lacking clear separation of steps.

Secondly, it included unnecessary navigation details that potentially misinformed the users. In the context of the example presented below **Error! Reference source not found.**, phrases such as “Face the direction you came from,” “until you reach the next landmark,” and “The height difference is approximately 5 meters” were extraneous and not pertinent to the conference venue, adding confusion. Lastly, the incorporation of too much non-essential information resulted in overly lengthy responses, which could be incomprehensible for users.

```
{
  "input": "How do I get to the tech booth?"
  "navigation_instructions": "turn around, move 8 meters, turn left, move 12 meters, stairs up, turn right, move 10 meters, turn left, move 3 meters",
  "output": "Great! Here are the step-by-step instructions on how to get to the tech booth based on the information provided:\n\n### Instruction 1: Turn Around (8 meters away)\n* Face the direction you came from and turn around to face the venue entrance again.\n\n### Instruction 2: Move 8 Meters Forward\n* Walk straight ahead for approximately 8 meters until you reach the next landmark.\n\n### Instruction 3: Turn Left (12 meters away)\n* Once you have reached the end of the hallway, turn left and walk down the new corridor.\n\n### Instruction 4: Climb the Stairs (up, 5 meters high)\n* Locate the staircase located nearby and climb the stairs to go up. The height difference is approximately 5 meters.\n\n### Instruction 5: Turn Right (10 meters away)\n* After reaching the top of the stairs, turn right and continue walking along the corridor.\n\n### Instruction 6: Move 10 Meters Forward\n* Walk straight ahead for approximately 10 meters until you reach the next landmark.\n\n### Instruction 7: Turn Left (3 meters away)\n* Once you have walked 10 meters forward, turn left and walk towards the tech booth.\n\nCongratulations! You have now arrived at the tech booth."
}
```

Considering these challenges, and after various system prompts were tested without achieving the desired format in navigation responses, a comprehensive fine-tuning process was initiated. The fine-tuning process was implemented not only to address the specific formatting issues identified in the responses but also to enhance the overall clarity and usability of the guidance provided by the agent. This step was essential for the system's evolution, ensuring that the navigation instructions provided were both user-friendly and in alignment with the intended navigational objectives.

The following subsections will delve into the specifics of the fine-tuning process, describing from the dataset creation process to the strategies and methodologies employed to refine the system's capabilities and achieve the desired level of performance and user experience.

Dataset Description

The creation of an extensive dataset was a fundamental step in the development of the navigation agent. This carefully structured dataset is comprised of three main components, each playing a vital role in ensuring the agent's effective response to user queries with accurate and useful navigation instructions.

Firstly, the dataset starts with the "input" field, which includes a wide variety of user questions. These queries, ranging from specific requests like "Guide me to the gym" to more general ones such as "I am completely lost, how do I reach the parking lot?", are designed to replicate a wide range of real-world scenarios. The inclusion of such a varied set of queries is crucial for training the agent to manage numerous types of interactions within a virtual environment. A key feature of the dataset is that it contains queries about a variety of rooms and places, not limited to those in the VOXReality conference venue. This strategy is intended to provide the model with a level of generalizability, equipping it to handle a diverse range of queries that might arise in different virtual environments.

The second component, the "knowledge" field, consists of real navigation instructions directly sourced from Mozilla Hubs. This essential part of the dataset includes actual navigation data, formatted as they appear in Mozilla Hubs, as mentioned in the previous subsection. Incorporating these real navigation instructions into the dataset is crucial to ensure that the agent is fine-tuned with accurate and authentic navigation patterns for effective and precise guidance.

The third field is the model's output, which is the agent's navigational guidance in response to the user's query. This output is designed to be clear and concise, prioritizing straightforward

navigation instructions and avoiding any unnecessary or confusing details. The primary aim here is to ensure that the agent's responses are user-friendly, helping users efficiently reach their intended destinations.

The dataset prepared for the navigation agent consists of 100 carefully selected samples for training purposes and an additional 40 samples specifically designated for testing model's performance. Each of these samples has undergone meticulous annotation by human experts, a process that is essential for the dataset's quality. The careful annotation process also serves to validate the relevance and applicability of each navigation instruction within the given "knowledge" field. The following examples showcases two examples from the dataset, providing a clear representation of how each sample is structured and annotated.

```
{
  "input": "How do I reach the parking lot?",
  "navigation_instructions": "move forward, move 10 meters, turn right, move 20 meters, turn left, move 30 meters, stairs down, turn right, move 10 meters",
  "output": "To find the parking lot:\n- Begin by moving forward 10 meters.\n- Turn right and continue 20 meters.\n- Then turn left and walk 30 meters.\n- Head downstairs.\n- Lastly, turn right and go an additional 10 meters.\n\nThis route will take you to the parking lot."
},
{
  "input": "Can you show the way to the exit?",
  "navigation_instructions": "turn around, move 11 meters, turn left, move 5 meters, turn right, move 9 meters, turn left, move 4 meters",
  "output": "To find the way to the exit:\n- Initially, turn around.\n- Go straight for 11 meters.\n- Make a left and travel 5 meters.\n- Turn right and proceed 9 meters.\n- Finally, turn left and move ahead 4 meters.\n\nYou will be at the exit."
}
```

In summary, the well-organized structure and intensive annotation of the dataset are instrumental in the fine-tuning process of the navigation agent. The objective is to enhance the agent's ability to process diverse user inputs, combine them with structured navigation instructions, and generate clear and effective responses. Such fine-tuning process is necessary for improving user experience in virtual space navigation, ensuring the agent offers precise and user-friendly assistance.

Fine-tuning process

In the fine-tuning process of the navigation agent, which utilizes the Llama-2 model, the Low-Rank Adaptation (LoRA) [Hu et al., 2021] method was used. The choice of LoRA for fine-tuning was based on its effectiveness and efficiency in adapting large-scale neural networks like Llama-2 to specific tasks. LoRA is particularly suited for this purpose due to its core principle of matrix factorization and low-rank approximations. This approach allows for the simplification of complex, highly parameterized layers in the neural network into more manageable structures. Figure 3.46 demonstrates LoRA during and after training, showcasing how this approach simplifies complex, highly parameterized layers in the neural network into more manageable structures. Such simplification is crucial for isolating and modifying the most influential features of the model without the need for extensive retraining.

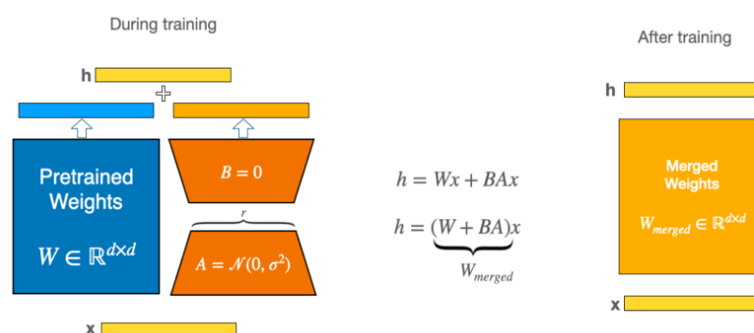


Figure 3.46. Illustration of LoRA during and after training

The use of the LoRA technique in refining the navigation agent is based on its effectiveness in tailoring neural networks for new challenges. LoRA accomplishes this by adding task-specific parameters but keeping the overall size manageable. This method is especially beneficial for the navigation agent because it allows for flexibility and efficiency in processing, which is vital for training substantial models like Llama-2 without excessive computational requirements.

After conducting experiments with various combinations of model parameters, the final selection for fine-tuning the model involved using 24% of all the model's parameters. This specific proportion of parameters was found to be the most effective in striking the right balance between model performance and computational efficiency. The fine-tuning process of the navigation agent yielded promising results (Figure 3.47). During the training epochs, both the training and validation loss showed a steady downward trend, signaling enhancements in the model's accuracy and its ability to generalize to new data. Specifically, the training loss exhibited a steep decrease, suggesting that the model was effectively learning the navigation tasks. Meanwhile, the validation loss also decreased, at a slower rate, which is a positive sign of the model's ability to adapt to new, unseen data. This reduction in loss across epochs demonstrates the effectiveness of the fine-tuning process, enabling the model to retain previously learned information while successfully acquiring new navigational knowledge.

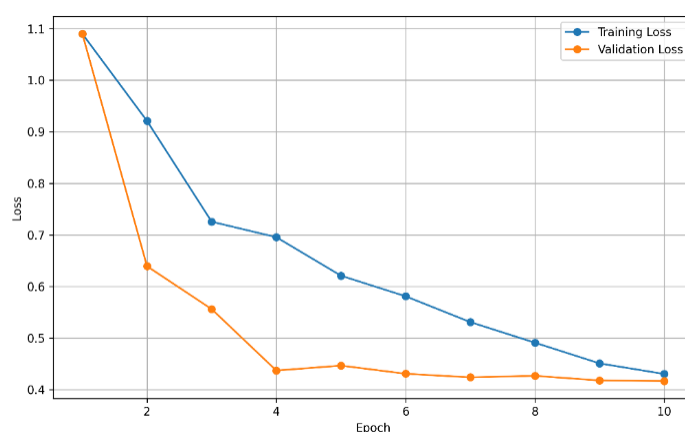


Figure 3.47. Training and validation loss per epoch

Evaluation Process

The evaluation of the navigation agent posed a complicate challenge, given the complex nature of the task. A comprehensive evaluation required careful consideration of two key aspects. Firstly, it was essential to confirm that the model accurately replicated each navigation instruction from Mozilla Hubs without omissions. This step is crucial to confirm the model's reliability in conveying complete navigational guidance. Secondly, the model's output needed to be compared with the ground truth in the test set to ascertain its accuracy. The proximity of the model's responses to this ground truth serves as a measure of its effectiveness in providing precise navigational directions.

So, the first evaluation process was to ensure that the model generated all the navigation instructions. To evaluate the performance of the model the semantic knowledge similarity (SKS) between the provided knowledge and the output of the model is compared. Specifically, the "all-MiniLM-L6-v2" sentence transformer was utilized. This model developed by the Sentence Transformers project [Reimers & Gurevych, 2019], which specializes in providing dense vector representations for sentences and paragraphs. These representations can be used in various NLP tasks such as clustering, semantic search, and information retrieval.

This model is based on the MiniLM architecture, which is known for its small size and efficient performance. The "L6" in the name indicates that the model has 6 layers, making it smaller and faster than larger models with similar capabilities. The model has been fine-tuned on a diverse and extensive dataset, consisting of over 1 billion sentence pairs, which contributes to its ability to understand and encode a wide range of sentences effectively. The model allows to encode batches of sentences, effectively converting textual information into numerical data that capture the semantic essence of the sentences.

To enhance the model's accuracy for the designated task, a fine-tuning process was implemented of the "all-MiniLM-L6-v2" model. This involved developing a specialized dataset comprised of three distinct sections: "anchor," "positive," and "negative." The "anchor" field contains the navigation instructions as issued by Mozilla Hubs. The "positive" field includes the correct rendition of the navigation instructions, ensuring completeness without omitting any details. Conversely, the "negative" field contains outputs that are similar with the "positive" instructions but with certain elements either missing or modified. For instance, in the example below, one "negative" field might involve misdirection such as "Turn left" with "Turn right" and changing the moving distance from "10 meters" to "8 meters." Another example of the figure shows a complete omission of a step. The construction of this dataset improves the capability of the model to recognize the errors, such as missing or incorrect navigation instructions.

```
{
  "anchor": "move 8 meters, turn right, move 10 meters, turn around, move 15 meters",
  "positive": "To get to the booth 2, follow these steps:\n- Move forward 8 meters.\n- Turn right and move 10 meters.\n- Turn around and proceed for 15 meters.\n\n You should now be in the booth 2.",
  "negative": "To get to the booth 2, follow these steps:\n- Move forward 8 meters.\n- Turn left and move 10 meters.\n- Turn around and proceed for 15 meters.\n\n You should now be in the booth 2."
},
{
  "anchor": "turn around, move 15 meters, turn left, move 20 meters, turn right, move 10 meters",
  "positive": "To reach the exit, follow these steps:\n- Start by turning around.\n- Move forward 15 meters.\n- Turn left and proceed for 20 meters.\n- Turn right and move 10 meters forward.\n\n You should now be at the exit.",
  "negative": "To reach the exit, follow these steps:\n- Start by turning around.\n- Move forward 15 meters.\n- Turn left and proceed for 20 meters.\n- Turn right.\n\n You should now be at the exit."
}
```

For the training and testing of this fine-tuning process, the dataset comprised 80 training samples and 30 test samples. This number of samples was chosen to provide a balance between comprehensive coverage of different scenarios and maintaining a manageable dataset size for efficient training and evaluation.

The effectiveness of the fine-tuning process was evaluated using three distinct metrics to measure accuracy, each providing a unique perspective on the model's ability to generate semantically similar sentence embeddings. The first metric employed was the cosine distance, which calculates the cosine of the angle between two vectors in a multi-dimensional space. In addition to cosine distance, the Manhattan and Euclidean distances were utilized. The Manhattan distance measures similarity by summing the absolute differences between coordinates of the vectors. On the other hand, the Euclidean distance measures the direct straight-line distance between two points in Euclidean space, giving a sense of the "ordinary" distance one would intuitively consider between points. The purpose of these evaluations is to understand how well the model is performing in terms of sentence similarity using different measures of distance. The usage of all three metrics helps to get a comprehensive view of the model's performance, as each distance measure captures different aspects of similarity. The evaluation results on the test dataset are presented in Table 3.28.

Table 3.28. Evaluation of sentence-transformers fine-tuning with sentence similarity

Cosine Similarity (%)	Manhattan Distance (%)	Euclidean Distance (%)
90	89.9	90



The evaluation of the navigation agent incorporated the fine-tuned sentence transformer model, specifically aimed at ensuring the agent's capability to generate all required navigation instructions comprehensively, without any omissions. For this purpose, the evaluation process involved a comparison of the provided navigation instructions with the outputs generated by the agent, followed by the calculation of their cosine similarity. Additionally, semantic textual similarity (STS) was evaluated using the original "all-MiniLM-L6-v2" model, by comparing the ground truth responses with those produced by the fine-tuned model. Moreover, the Rouge score was also employed as a metric in this evaluation process. The inclusion of the Rouge score provided an additional layer of assessment regarding the preciseness and completeness of the generated instructions. It helped in estimating how well the model's outputs aligned with the key information presented in the original navigation instructions.

To benchmark the effectiveness of the fine-tuning process, a comparative evaluation was conducted between the original Llama 2 model and the fine-tuned. This assessment aimed to highlight the enhancements achieved through the fine-tuning. The results of this comprehensive evaluation are presented in Table 3.29. The table offers a detailed comparison, showcasing the performance of both models across various metrics.

Table 3.29. Evaluation of navigation agent before and after fine-tuning

Phase	Rouge-1 (%)	Rouge-2 (%)	Rouge-L (%)	SKS (%)	STS (%)
Before	33.1	17.2	29.4	85.1	92.2
After	68.9	42.3	60.5	95.7	97.3

Evaluation of the Workflow before and after fine-tuning

Moving forward in the development process of our workflow with the fine-tuned Llama 2 model, we directed our focus towards a critical evaluation phase. This was particularly pertinent with the fine-tuning efforts that had been applied to the Llama2 model, as we aimed to comprehensively understand the influence of these modifications on the model's performance. It was imperative to understand the impact of the fine-tuning process on the model's performance in crucial workflow tasks. This phase of evaluation was especially focused on tasks such as intent and destination recognition.

To ensure the accurate classification of user intent by the model, an extensive evaluation process was implemented. This process began with the creation of an evaluation dataset, comprising 99 carefully selected samples, evenly distributed with 33 samples for each of the three categories. Once the evaluation dataset was established, each user query from the dataset was inputted into the model to generate responses. The effectiveness of the model was then determined by comparing these responses against the ground truth labels provided in the dataset. Key performance metrics such as accuracy, precision, recall, and F1 score were calculated to provide a comprehensive assessment of the model's performance.

To ensure a comprehensive assessment, the evaluation process was conducted on both the original Llama 2 model and its fine-tuned version. This was crucial to verify that the fine-tuning, while enhancing the model's performance for navigation instruction tasks, did not compromise its effectiveness in other areas. The goal was to maintain the model's versatility and ensure its proficiency across a broad spectrum of tasks, not just those it was specifically fine-tuned for.

In the evaluation of the Intent Recognition task, the results post-fine-tuning showed a decrease in key performance metrics compared to the pre-fine-tuning phase. Specifically, there was a reduction in accuracy, precision, recall, and F1 score. This decline suggests that while the fine-tuning process targeted certain aspects of the model's functionality, it may have had



unintended effects on its ability to perform in the intent recognition task. It's important to recognize that such outcomes are part of the iterative process of model development and refinement. Future iterations and enhancements are expected to address these areas, aiming to strike a balance between specialized improvements and overall task proficiency. The results are presented in Table 3.30.

Table 3.30. Evaluation of intent recognition agent before and after fine-tuning

Phase	Accuracy (%)	Precision (%)	Recall (%)	F1 score (%)
Before	85.8	86.1	85.8	85.8
After	79.7	80.4	79.7	79.7

For the destination recognition task, a similar evaluation method to the one used for intent recognition was implemented. A dataset with 100 samples, evenly distributed across various destination categories, was used to assess the model's performance. Additional samples were included in the 'no location' category to test scenarios involving non-existent locations within the building, as previously mentioned, and instances where users request navigation guidance without specifying a destination. This evaluation process involved analyzing the model's responses to the dataset queries and comparing them with the correct labels. The key metrics accuracy, precision, recall, and F1 score were calculated again for the original Llama 2 model and the fine-tuned version, and the results are succinctly presented in Table 3.31.

Table 3.31. Evaluation of destination recognition agent before and after fine-tuning

Phase	Accuracy (%)	Precision (%)	Recall (%)	F1 score (%)
Before	98	98.1	98	97.8
After	97	97.1	97	96.8

In the case of the destination recognition task, the performance metrics post-fine-tuning remained robust, showing no significant drop in effectiveness. The accuracy, precision, recall, and F1 score maintained high levels, with only marginal variations observed. This consistency indicates that the fine-tuning process did not adversely affect the model's ability to perform in the destination recognition task. Such stability in performance is encouraging, as it demonstrates the model's resilience and effectiveness in maintaining high standards of task proficiency even after undergoing specific modifications.

Deployment

The conference assistant provides a REST API, developed using FastAPI, interface for interacting with a conference-related virtual assistant. This assistant is designed to handle user queries related to conference navigation, program information, and general inquiries. The API is structured around two main endpoints, each catering to different aspects of the assistant's functionalities. The detailed explanation of the API is presented in Appendix III.

The first endpoint is the intent destination endpoint (Figure 3.48), which processes user queries to determine the user's intent and the corresponding destination within the conference context. This function takes a string input representing the user query and returns a response identifying the intent and destination.



POST /intent_dest/ Intent Destination

Parameters

No parameters

Request body *required* application/json

```
{
  "user_query": "How can I go to the social area?"
}
```

Cancel Reset

Figure 3.48. The intend and destination recognition endpoint

The second and the last endpoint is the response generation, which caters to generating responses based on user queries, and additional Mozilla inputs. It offers personalized assistance depending on whether the query is related to navigation, program information, or general conference queries. The Figure 3.49 demonstrates the graphical representation of the response generation endpoint in the FastAPI framework.

GET / Root

POST /intent_dest/ Intent Destination

POST /response/ Response

Parameters

No parameters

Request body *required* application/json

```
{
  "user_query": "How can I go to the social area?",
  "mozilla_input": "turn right, move 10, turn left, move 5, stairs up, turn right"
}
```

Cancel Reset

Figure 3.49. The response generation endpoint

Updating Virtual Conference Agent

In this version of the conference dialogue agent, we introduce a significant change in the underlying language model, transitioning from **Llama2-Chat 7B**¹² to **Mistral-7B-Instruct-v0.2**¹³. Over the past two years, the landscape of open-source large language models (LLMs) has evolved significantly, with numerous models becoming publicly available. After careful evaluation, Mistral-7B-Instruct-v0.2 was selected as it outperforms Llama2 in terms of efficiency, response quality, and adaptability to complex dialogue interactions.

Despite this model change, the overall workflow structure remains consistent with the previous version. The intent recognition process, responsible for categorizing user queries, remains largely unchanged except for the addition of 3 new intents (trade show, summary presentation agent and summary conversation agent). For the navigation assistance module, the logic and implementation remain the same, ensuring seamless virtual space guidance within the conference environment. Similarly, the program information retrieval system continues to function as in the previous version, providing users with real-time access to conference schedules, speaker details, and session locations. The general query handling remains unchanged, ensuring the system continues to handle questions that fall outside the primary scope of the agent.

In addition to maintaining the core functionalities, this version introduces several key enhancements. The agent is now capable of providing detailed information about the trade show, including participants, products, and exhibitor details. Users can inquire about specific

¹² <https://huggingface.co/voxreality/llama2-navigation>

¹³ <https://huggingface.co/voxreality/mistral-7B-navigation-new-instructions>

companies, showcased innovations, or booth locations, making it easier to navigate the trade show aspect of the event. Another enhancement is the ability to summarize presentations. When a participant exits a session, they can request a summary from the agent, allowing them to quickly review key points, highlights, and takeaways from the talk. This feature helps users retain and revisit important information more efficiently. Additionally, the agent can now generate summaries of one-on-one discussions between participants. After a private conversation, users can request a concise recap, capturing important topics, agreements, or action points, making post-discussion follow-ups more seamless and organized. The design of the updated workflow is illustrated in Figure 3.50 providing a visual representation of the agent's structure and processing logic.

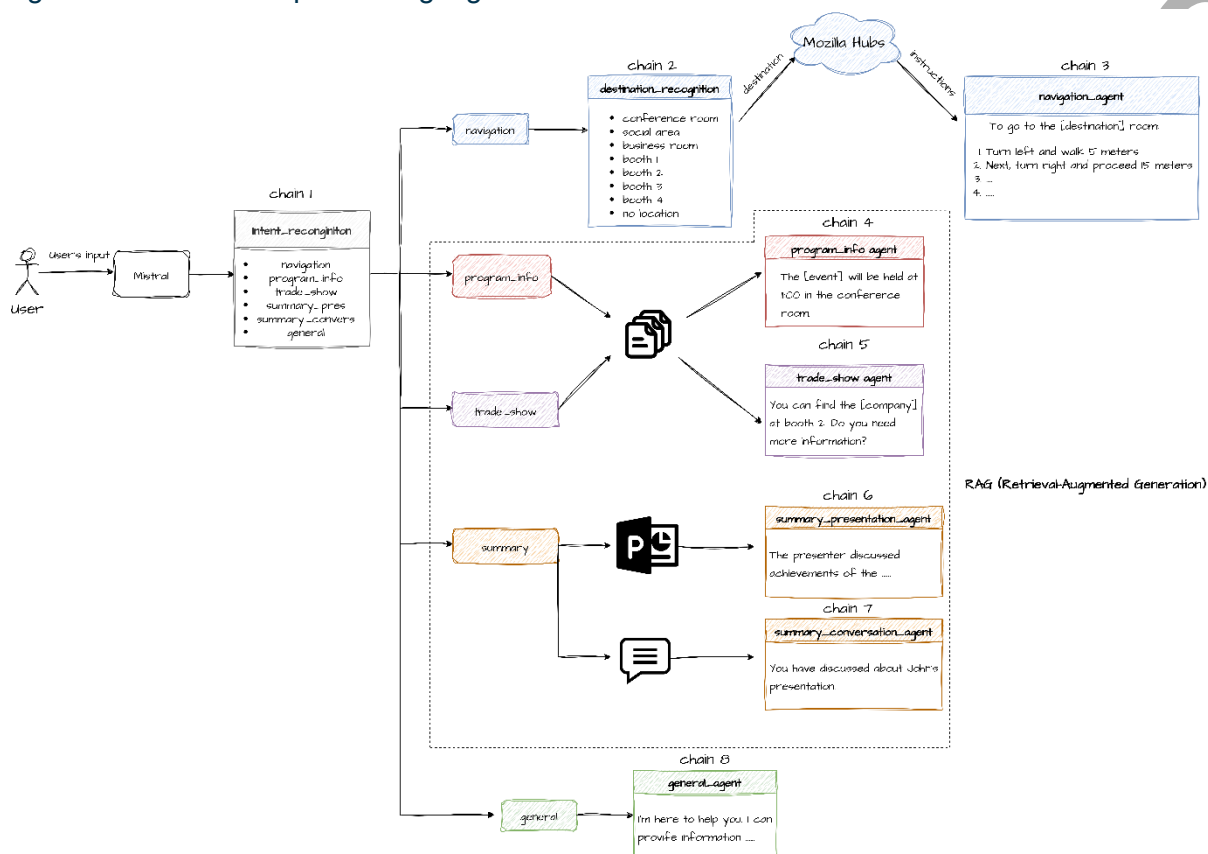


Figure 3.50: Conference agent workflow

These updates expand the agent's capabilities, making it a more comprehensive and intelligent assistant for conference attendees by providing richer information access and enhanced summarization functionalities. In the next sections, we will delve deeper into the modifications made to the existing features and the implementation of the newly introduced functionalities, providing a detailed analysis of their design, integration, and evaluation. Only the components that have been updated or newly introduced will be described in detail. Specifically, we will cover intent recognition, the navigation, the trade show, the presentation summarization, and the one-on-one communication summarization agents. The destination recognition, program information agent, and general agent remain unchanged from the previous version and will not be described further.

Intent Recognition

The intent recognition process follows the same implementation as in the previous version, utilizing few-shot examples to guide the model in classifying user queries. This approach involves providing the model with a set of predefined examples for each intent category, helping it learn how to correctly interpret and categorize new inputs. Based on the user's query, the model assigns it to one of the predefined categories, ensuring accurate routing to the appropriate response chain. The only modification in this version is the addition of new

intents: "trade show", "summary presentation agent" and "summary conversation agent" which enables the agent to identify when a user is requesting information about trade show participants, products, summary or exhibitors. Below there are examples of the few shot process (1 of each class):

```
{user_prompt: "Where is the main auditorium located?" "response": "navigation"},
{user_prompt: "What time does the keynote speech start?" "response": "program_info"},
{user_prompt: "Who is exhibiting at Booth Number A01?" "response": "trade_show"},
{user_prompt: "Can you give me a summary of the presentation?" "response": "presentation_summary"},
{user_prompt: "What we discussed with John?" "response": "conversation_summary"},
{user_prompt: "Can you assist me?" "response": "general"}
```

To assess the performance of the intent recognition component, an evaluation dataset of 192 samples was created, with 32 samples for each of the six intent categories. These categories include navigation, program information, trade show, presentation summarization, one-to-one conversation summarization and general. This dataset was used to test the model's ability to correctly classify user queries into their respective intents.

Each query from the dataset was input into the model, and its predicted intent was compared against the ground truth labels. The evaluation was conducted using standard classification metrics, including accuracy, precision, recall, and F1 score, providing a comprehensive assessment of the model's effectiveness in identifying user intents. The evaluation results are summarized in Table 3.32, demonstrating the model's performance across the different intent categories.

Table 3.32: Evaluation of intent recognition agent

Accuracy (%)	Precision (%)	Recall (%)	F1 score (%)
84.2	84.2	84.2	84.2

Navigation Agent

The main logic of the navigation agent remains unchanged in this version. As before, the agent accepts two inputs: user queries and Mozilla Hubs' navigation instructions. The agent processes user requests for directions and integrates the structured navigation instructions provided by Mozilla Hubs to generate an appropriate response. These instructions guide users efficiently within the virtual conference environment, ensuring they can reach their destinations accurately.

The only significant change in this version is the format of the navigation instructions. In the previous implementation, users often struggled to interpret distances measured in meters within the virtual space. To address this issue, we introduced a new navigation format that relies on visual markers and room-based references instead of numerical distance descriptions.

To make navigation more intuitive, visual objects such as yellow spheres and blue spheres have been strategically placed within the environment. These objects serve as reference points, allowing users to follow step-by-step directions based on recognizable landmarks rather than estimated distances. Additionally, navigation instructions now emphasize room names and spatial relationships, making it easier for users to understand and follow directions.

The updated navigation instructions maintain the core structure of the previous version while introducing visual references to improve user comprehension. The primary navigation commands now include:

- **Turn around:** This directive helps users align themselves correctly at the start of their journey, ensuring they face the correct direction before proceeding.
- **Turn right/left:** These commands guide users in making accurate turns at key points along their route, helping them stay on the intended path.
- **Continue until you arrive at [landmark]:** Instead of specifying distances in meters, users are instructed to move forward until they reach a recognizable reference point, such as a wall, a corridor, or a visual marker like a yellow sphere or blue sphere.
- **Cross [landmark] on your left/right:** This instruction helps users navigate by referencing objects or areas they will pass along the way, such as booths, social areas, or corridors.
- **Pass through the corridor/opening:** When the path involves moving through a hallway, doorway, or other passage, this command ensures users proceed correctly to the next section of the environment.

By replacing distance-based navigation with landmark-based instructions, the updated navigation system enhances clarity and usability, making it easier for users to orient themselves and reach their destinations more intuitively. However, it is important to note that the navigation agent described here, which is powered by **Mistral-7B-Instruct-v0.2**, is not the final model that will be used in the second phase of the pilots. Instead, it serves as a tool for generating the dataset needed to fine-tune the Vision-Language (VL) model, as outlined in Section 3.3. The VL model will ultimately replace the navigation system, enabling more advanced and multimodal navigation capabilities that integrate both visual and linguistic cues for improved user guidance in virtual environments.

Trade Show Agent

The Trade Show Agent provides users with detailed information about the booths, exhibitors, and products showcased at the trade show. To achieve this, the agent utilizes a Retrieval-Augmented Generation (RAG) system, allowing it to extract and retrieve relevant information from a structured PDF document containing exhibitor details.

The PDF document includes comprehensive information about the trade show, such as the event name, dates, and location, along with detailed descriptions of each exhibitor. For each booth, the document provides the exhibitor's name, description, relevant keywords, and a list of showcased products. Additionally, it contains contact details for exhibitor representatives, ensuring users can easily obtain further information or connect with company representatives. An example of the PDF file is illustrated in Figure 3.51.

The RAG system processes this information by first extracting text from the PDF and then dividing it into smaller chunks to facilitate efficient searching. A character splitter is used to segment the text while maintaining slight overlaps between sections, ensuring that key details remain contextually connected.

Event Details

Start Date: 2025-10-05

End Date: 2025-10-07

Location: Trade Show

Booth Number: 1

Exhibitor: Synelxis Solutions S.A.

Description: Software development company, specializing in the design and implementation of cutting-edge solutions for various industries.

Keywords: Software, AI, IoT, Smart Agriculture

Product: SynField (<https://www.synfield.gr/>)

Contact: John Papadopoulos (some.one@synelxis.com), Phone: +30 6905482952

Booth Number: 2

Exhibitor: Hololight

Description: Hololight is transforming how we work and collaborate using augmented and virtual reality.

Keywords: Virtual Reality, Augmented Reality, Immersive Technology

Product: Hololight Space - XR Application, Hololight Stream - SDK, Hololight Hub - XR Platform (<https://hololight.com/>)

Contact: Alisa Smith (alisasmith@holo-light.com), Phone: +1 (123) 456-789

Figure 3.51: Sample PDF document for the trade show

Once the text is structured into manageable sections, each chunk is converted into an embedding using a pre-trained sentence transformer model. These embeddings enable the system to understand the semantic relationships between different pieces of information, improving the accuracy of search results. The processed text chunks and their embeddings are then stored in a FAISS (Facebook AI Similarity Search) library¹⁴, an efficient database that allows for quick and scalable retrieval of relevant information.

When a user submits a query about the trade show, the system converts the query into an embedding and searches the FAISS index for the most relevant text segments. The retrieved content is then passed to the language model, which generates a structured response in a clear and user-friendly format. By leveraging RAG-based retrieval, the trade show agent ensures that users receive accurate, contextually relevant, and up-to-date information about exhibitors, products, and booth locations, enhancing their experience at the event.

In our evaluation of the RAG agent, we conducted two distinct assessments: one focused on the retrieval performance and the other on the generation quality. The evaluation process was carried out over a dataset of 20 carefully selected queries derived from the provided PDF. This comprehensive dataset allowed us to rigorously assess both components of the system.

The retrieval evaluation employed metrics such as Precision@3, Recall@3, and the mean Reciprocal Rank (MRR) to quantify how accurately the system returned relevant document passages. For instance, for a representative query—"What are the details for booth 2 at TechExpo 2025?"—the system correctly retrieved content that identified Hololight as the exhibitor, including product details and contact information. In this context, Precision@3 indicates the proportion of relevant documents among the top three retrieved results, while

¹⁴ <https://github.com/facebookresearch/faiss>

Recall@3 represents the fraction of all relevant documents that are captured within those top three results. Additionally, MRR calculates the average reciprocal of the rank position where the first relevant document appears across all queries, thereby providing insight into how early in the ranking the relevant content is found.

The generation evaluation employed metrics such as ROUGE-L, BLEU, and BERTScore, along with an expert review to verify factual accuracy. ROUGE-L measures the overlap of the longest common subsequence between the generated response and the reference text, thereby assessing both content preservation and fluency. BERTScore leverages contextual embeddings from pretrained language models to capture semantic similarity between the generated and reference responses, even when there is limited n-gram overlap. Factual accuracy was determined through expert evaluation, ensuring that the generated responses faithfully reflected the information contained solely within the provided PDF. The results of our evaluation are summarized in Table 3.33.

Table 3.33: Evaluation of the trade show agent

Component	Metric	Result
Retrieval	Precision@3	80%
Retrieval	Recall@3	90%
Retrieval	Mean Reciprocal Rank	95%
Generation	ROUGE-L	85%
Generation	BLEU	75%
Generation	BERTScore	90%
Generation	Factual Accuracy	94%

These results demonstrate that our RAG system exhibits robust retrieval capabilities, with high precision, recall, and early retrieval of relevant information as indicated by the MRR. Concurrently, the generation component consistently produces responses that closely match the gold standard, as shown by high ROUGE-L, BLEU, and BERTScore values, alongside a factual accuracy rate of 94%.

Presentation Summarization Agent

The Presentation Summarization Agent is designed to assist conference attendees by generating concise summaries of presentations they have attended. The system processes PowerPoint files, extracts textual content from the slides, and then generates a structured summary using a large language model (LLM). This approach ensures that users receive a clear and contextually relevant summary of the key points discussed during the session.

The summarization process begins with the extraction of textual content from the presentation slides. The system utilizes a PowerPoint document loader that systematically processes each slide, identifying and categorizing the different elements within the file. The text extraction process differentiates between multiple content types, including:

- **Titles:** The system detects the slide titles and assigns them a slide number, ensuring a structured representation of the presentation flow.
- **Narrative Text and Bullet Points:** The content of each slide is extracted and categorized as either main text or list items, preserving the logical structure of the slides.
- **Page Breaks:** When a new slide begins, a page break is identified, allowing for proper segmentation of the extracted text.

This segmentation process ensures that the extracted text maintains a coherent format, making it easier for the language model to process and summarize.



Once the PowerPoint content is converted into structured text, the extracted information is fed into an LLM to generate a concise and meaningful summary. The model is guided by a simple system prompt that instructs it to act as a helpful assistant summarizing a presentation for a user who briefly stepped out. The prompt tells the model to highlight the key points, main arguments, and critical details shared by the presenter while keeping the summary concise, clear, and free of unnecessary details or speculation. By following this structured approach, the Presentation Summarization Agent ensures that users receive accurate and well-organized summaries, enhancing their ability to stay informed even if they momentarily leave a session. Here is an example of how the agent generates a summary of a presentation based on a user's request:

User: "What did the presenter covered?"

Agent: "The presenter covered recent advancements in AI-driven cybersecurity."

Key Discussion Points:

- AI models are being used to detect cyber threats in real-time.
- Anomaly detection plays a crucial role in identifying security breaches.
- Case studies were presented on how AI could have prevented recent cybersecurity attacks.
- Ethical concerns were raised about data privacy and algorithmic bias in AI security systems."

In our evaluation of the presentation summarization agent, we focused exclusively on assessing the quality of the generated summaries. The system extracts the complete textual content from the input PDF and uses it as context to produce a concise summary. We evaluated the system on 20 presentation files, each summarized by experts. We compared the generated summaries with these references using ROUGE-L, BLEU, and BERTScore, and further assessed their factual accuracy through expert review. Table 3.34 presents the average results obtained across the dataset.

Table 3.34: Evaluation of the presentation summarization agent

Metric	Result
ROUGE-L	87%
BLEU	78%
BERTScore	89%
Factual Accuracy	96%

One to One Summarization Agent

The One-to-One Conversation Summarization Agent is designed to assist conference attendees by providing a structured summary of their discussions in the business meeting room. During the conference, two participants may engage in a conversation within the business meeting space, discussing various topics such as collaborations, project ideas, or future opportunities. Once the meeting concludes, the agent enables users to quickly retrieve a summary of what was discussed, ensuring that key points, agreements, and action items are preserved. This feature helps attendees efficiently recall important details from their discussions without the need to manually review or take notes during the meeting.

The summarization workflow begins with the ASR model, which captures and transcribes the spoken conversation between participants. This transcription process ensures that all dialogue exchanges, contextual nuances, and critical discussion points are preserved in text form. Once the ASR model processes the audio, the resulting text is saved in a structured format (TXT file), making it accessible for further processing.

After the transcription is completed, the text file is fed into the One-to-One Summarization Agent, where it undergoes further refinement and structuring using an LLM-based inference system. The summarization process is driven by a custom system prompt, which provides explicit instructions to the model on how to extract relevant information from the conversation. The model follows a structured approach to ensure clarity and completeness in the generated summary.

Once the model processes the transcript, it generates a well-structured, easy-to-read summary. The final output presents a concise yet informative overview of the conversation, ensuring that users can quickly retrieve essential information without manually reviewing the entire transcript. The structured format makes it particularly useful for business meetings, networking discussions, and collaborative sessions, where users need to keep track of important agreements and action items. Here is an example of a user query requesting a summary of the conversation from the agent:

User: ""What did we discuss with John?""

Agent: "You discussed about the upcoming financial and operational decisions for Q2.
Key Discussion Points:
-The Q2 budget needs to be finalized by the end of the week.
-The vendor selection for the new software tool is still pending.
-The IT department's assessment is required before making the final vendor decision.
-A cost-benefit analysis report should be prepared before the next meeting.
Action Items:
1. John will follow up with IT for their final vendor assessment.
2. The team will prepare a cost-benefit analysis report before the next meeting."

The evaluation process of the One-to-One Conversation Summarization Agent is similar with the process that we followed in the Presentation Summarization Agent. A dataset of 20 business meeting transcripts was used, each containing conversations between two participants on topics such as collaborations, project planning, and operational decisions. Human experts manually summarized these transcripts to create gold-standard reference summaries, serving as benchmarks for evaluating the system's performance. The evaluation results are presented in Table 3.35.

Table 3.35: Evaluation of the one-to-one summarization agent

Metric	Result
ROUGE-L	86%
BLEU	79%
BERTScore	80%
Factual Accuracy	95%

Fine Tuning the Model for the Updated Virtual Conference Agent

In this version of the conference dialogue agent, the fine-tuning process was applied again to enhance the system's overall performance and adaptability. There are two primary reasons for this fine-tuning iteration.

The first reason is the transition from Llama2-Chat 7B to Mistral-7B-Instruct-v0.2. As discussed earlier, this change was made to leverage advancements in open-source large language models (LLMs), as Mistral-7B has demonstrated superior efficiency, response quality, and adaptability to complex dialogue scenarios. Fine-tuning was necessary to ensure that the new model properly aligns with the requirements of the conference agent.

The second reason is the reconstruction of the dataset, particularly in the navigation agent. As detailed in Subsection 3.4.1, the response format was modified to replace distance-based navigation instructions with a landmark-based approach, improving user comprehension in the virtual environment. This restructuring required fine-tuning to ensure that the model could



accurately generate responses in the new format, integrating spatial references such as rooms and visual markers (e.g., yellow and blue spheres) instead of numerical distances.

In this subsection, we will provide a detailed overview of the dataset used for fine-tuning and the methodology applied to optimize the model's performance.

Dataset Description

The dataset¹⁵ used for fine-tuning the navigation agent retains the same structure as in the previous version, consisting of three main components: input queries, knowledge data, and model outputs. The input field contains a diverse set of user queries, including both specific requests (e.g., "Guide me to the gym") and more general navigation-related inquiries (e.g., "I am completely lost, how do I reach the parking lot?"). These variations ensure that the model is trained to handle a wide range of real-world scenarios within a virtual environment. The knowledge field consists of real navigation instructions extracted from Mozilla Hubs, providing authentic navigation patterns. This structured information allows the model to generate accurate and contextually relevant directions when assisting users. The output field, however, has been significantly modified in this version. Instead of generating distance-based instructions (e.g., "Move forward 10 meters, then turn left"), the model now produces landmark-based navigation guidance to improve user comprehension in the virtual space. The updated responses incorporate room names and visual markers such as coloured spheres to make navigation more intuitive. The following example illustrates how each dataset sample is structured and annotated, offering a clear representation of the data format.

```
{
  "input": "Where is the cafeteria located?",
  "navigation_instructions": "start, turn left, crossing yellow sphere left, arrive wall opening, turn left, turn right,
pass corridor, crossing magenta sphere left, arrive conference room, finish",
  "output": "To reach the cafe bar, follow these steps:
    - Start by turning left.
    - Cross the yellow sphere on your left.
    - Continue until you arrive at a wall opening.
    - Turn left.
    - Turn right.
    - Pass through the corridor.
    - Cross the magenta sphere on your left.
    You should now be at the conference room."
}
```

This approach eliminates confusion regarding numerical distances and enhances navigation clarity by referencing visible, easily recognizable landmarks. The dataset consists of 100 training samples and 40 test samples, the same as in the previous version. Each sample has been carefully annotated by human experts to ensure high-quality data and accurate navigation guidance.

Fine-tuning Process

To adapt the Mistral-7B-Instruct-v0.2 model to the updated requirements of the conference agent, a fine-tuning process was conducted. As in the previous version, Low-Rank Adaptation (LoRA) [Hu et al., 2021] was applied to efficiently modify the model while keeping computational costs manageable. This technique allowed for targeted improvements in navigation response generation without requiring full model retraining.

The fine-tuning process focused on optimizing the model's ability to generate landmark-based navigation instructions, ensuring that responses aligned with the newly structured dataset. Various configurations were tested, and the final tuning utilized 24% of the model's

¹⁵ https://huggingface.co/datasets/voxreality/navigation_instructions_v2

parameters, achieving a balance between performance enhancement and efficiency. Figure 3.52 presents the training and validation loss per epoch, demonstrating a steady decline, indicating improved model generalization and accuracy in generating structured navigation responses.

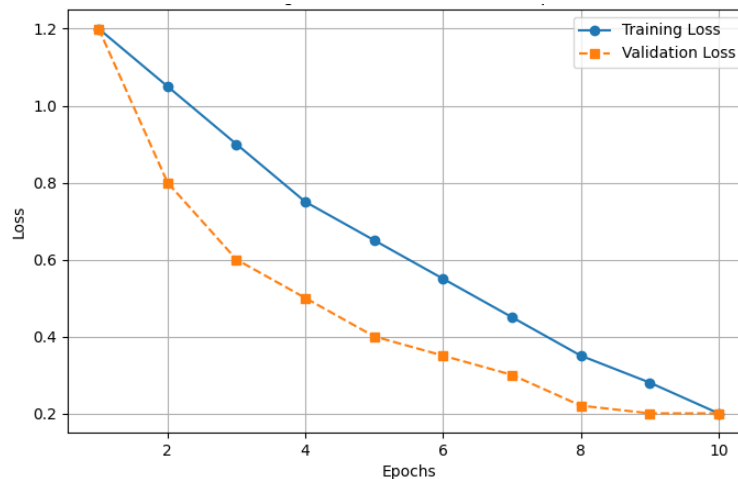


Figure 3.52: Training and validation loss per epoch

Evaluation of the Navigation Agent

The evaluation process for the navigation agent followed the same methodology as in the previous version, focusing on assessing the model's ability to generate accurate, complete, and structured navigation instructions. Since the dataset format was modified to incorporate landmark-based navigation rather than distance-based directions, it was necessary to re-fine-tune the sentence transformer model to align with these changes.

The first step in the evaluation involved verifying that the model generated all necessary navigation instructions without omissions. This was done by comparing the model's responses to the original navigation instructions provided by Mozilla Hubs. The objective was to ensure that the agent preserved all critical navigation steps and accurately conveyed them in the updated format. The updated dataset used for fine-tuning the all-MiniLM-L6-v2 sentence transformer [Reimers & Gurevych, 2019] consists of three key components: anchor, positive, and negative fields.

- The **"anchor"** field contains a simplified sequence of navigation steps, representing the core navigation path using key landmarks and actions. This structured format ensures a standardized way of describing movements within the virtual environment.
- The **"positive"** field provides the correct, well-structured navigation instructions based on the anchor sequence. These instructions are formatted in a clear and user-friendly manner, ensuring that the model learns to generate step-by-step guidance using the appropriate spatial references (e.g., "cross booth 1 on your left").
- The **"negative"** field contains incorrect or misleading instructions, which include errors such as wrong directional references (left instead of right), omitted steps, or unnecessary deviations. This field helps the model learn to differentiate between accurate and flawed navigation guidance, improving its ability to generate reliable responses.

These fields ensure that the model learns to generate navigation instructions that are clear, structured, and contextually accurate while distinguishing between correct and incorrect guidance. The dataset maintains the same number of training (80 samples) and test (30 samples) instances as the previous version. Below is an example of the updated dataset:

```

{
  "anchor": "start, turn around, pass corridor, crossing blue sphere left, crossing wall opening, crossing booth 1 left, arrive wall, turn left, crossing booth 2 left, arrive booth 3, finish",

  "positive": "To find the exit, follow these steps: - Start by turning around. - Pass through the corridor. - Cross the blue sphere on your left. - Cross the wall opening. - Cross booth 1 on your left. - Continue until you arrive at the wall. - Turn left. - Cross booth 2 on your left. You should now be at the exit.",

  "negative": "To find the exit, follow these steps: - Start by turning around. - Pass through the corridor. - Pass by the blue sphere on your left. - Cross the wall opening. - Cross booth 1 on your right. - Continue until you reach the wall. - Turn right. - Cross booth 2 on your left. You should now be at the exit."
}

```

By leveraging this dataset structure, the fine-tuning process refines the sentence transformer's ability to assess the correctness of generated instructions, ensuring that the navigation agent delivers precise and context-aware directions. The fine-tuned sentence transformer was evaluated using three distance-based metrics to measure sentence similarity. Cosine Similarity assesses how closely the generated response matches the expected output in semantic meaning. Manhattan Distance measures similarity based on absolute differences between sentence embeddings. Euclidean Distance evaluates similarity based on the direct distance between two embedding vectors. The results, presented in Table 3.36, confirm that the model achieves high similarity scores, demonstrating its ability to generate responses that closely match the reference navigation instructions.

Table 3.36. Evaluation of sentence-transformers fine-tuning process

Cosine Similarity (%)	Manhattan Distance (%)	Euclidean Distance (%)
92	90.2	92

Additionally, the evaluation process incorporated ROUGE scores (ROUGE-1, ROUGE-2, and ROUGE-L) as well as Semantic Knowledge Similarity (SKS) and Semantic Textual Similarity (STS). These metrics were used to assess how well the model preserved the key details and structure of the navigation instructions after fine-tuning. The final results, summarized in Table 3.37, show a substantial improvement in the model's ability to generate clear and precise navigation instructions. The higher ROUGE, SKS, and STS scores indicate that the fine-tuned model is more effective in delivering structured and user-friendly guidance based on the new dataset format.

Table 3.37 Evaluation of conference agent after fine-tuning

Rouge-1 (%)	Rouge-2 (%)	Rouge-L (%)	SKS (%)	STS (%)
71.8	49.3	66.4	96.5	98.2

Conversation History Management

Another important difference compared to the previous version is the addition of a conversation memory management system. This enhancement allows the agent to maintain context and coherence in multi-turn interactions, ensuring that responses remain relevant throughout an ongoing conversation. By implementing this memory mechanism, the agent can provide a more natural, context-aware dialogue experience, improving user engagement and overall interaction quality.

The memory management system uses a specific technique in which only the last 'k' number of conversations are passed to the agent, with 'k' set to 7 in our system. At the start of each interaction, the system initializes a memory buffer that stores the last seven exchanges. Each interaction, including both the user's input and the AI's response, is recorded in this buffer. This continuous recording allows the AI to reference recent exchanges, maintaining a coherent and contextually relevant dialogue. The memory buffer is updated as the

conversation goes on, removing earlier exchanges to make space for new ones after seven interactions are reached.

We chose to retain only the last seven exchanges to balance the need for context with practical limitations of memory and processing capacity. Without this approach, several significant challenges would arise. Firstly, providing the agent with the entire conversation history on each follow-up call would result in considerable computational overhead and hallucination to the model. Moreover, LLMs have token limits that restrict the maximum number of tokens they can process in a single call. An expanding conversation history could easily exceed these limits, leading to incomplete or failed responses. Additionally, handling more tokens increases the response time of the agent, introducing latency and degrading the user experience. Another critical issue is the relevance of older messages. As conversations evolve, initial exchanges may become irrelevant, introducing noise that can confuse the agent and potentially lead to incoherent or inappropriate responses.

By retaining these exchanges, we effectively allay these challenges. This approach ensures that the agent maintains a clear and relevant context for ongoing interactions while managing costs, adhering to token limits, and maintaining response efficiency. The memory buffer is continuously updated, discarding older exchanges to make room for new ones, thus preserving the most pertinent information for the agent's reference.

Overall, this conversation memory management system enhances the virtual navigation experience by ensuring that the agent's guidance and responses remain relevant and contextually appropriate. This method significantly improves user engagement and satisfaction within the virtual environment, enabling the AI to provide more personalized and effective assistance.

Conversation History Management Evaluation

Effective conversation history management is crucial for a conversational agent. To assess our agent's ability to recall and utilize recent conversational context, we conducted evaluations focusing on its short-term memory capabilities. Recognizing that interactions in such settings are typically concise and context-specific, we limited the conversation history to the most recent seven exchanges.

Several evaluation metrics have been proposed to assess conversation history management in conversational agents. In our case, we adopt the evaluation framework introduced by [Zhong et al., (2024)], which includes three key metrics: Memory Retrieval Accuracy (MRA), Response Correctness (RC), and Contextual Coherence (CC). MRA measures the model's ability to successfully retrieve relevant information from previous exchanges, RC evaluates whether the response correctly answers the user's inquiry, ensuring factual accuracy and relevance, and CC assesses how well the response maintains logical flow and consistency with the dialogue context. These metrics allowed us to comprehensively assess the agent's performance in recalling relevant information, providing accurate responses, and maintaining logical coherence within the conversation.

Our evaluation involved 10 different sessions, each containing chats designed to test various aspects of conversation history management. The scenarios ranged from simple queries, such as checking if the agent remembers the user's name, to more complex inquiries like asking the agent to recall all the rooms the user had visited in the conference venue. Examples of two such conversations are presented in Figure 3.53.

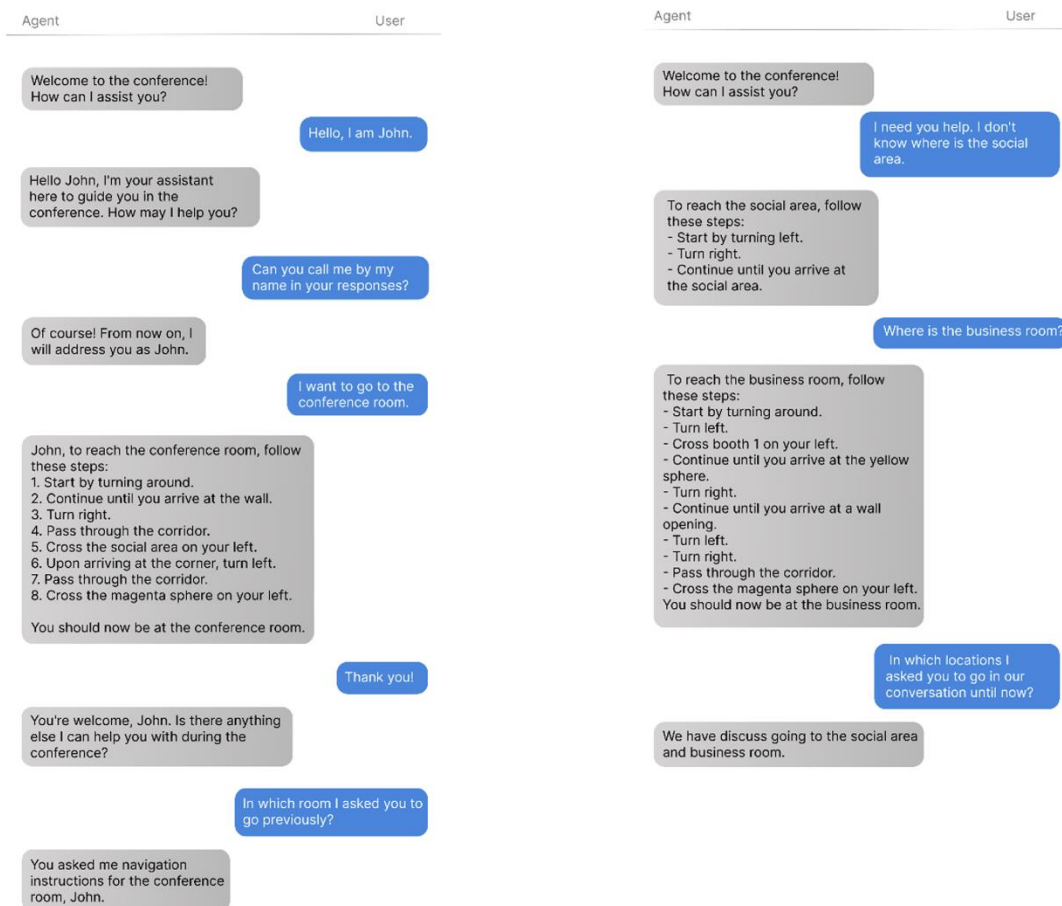


Figure 3.53: Two example conversations with the agent

The results of our model's evaluation in conversation history management are encouraging, demonstrating notable strengths and some areas for refinement. The model exhibited a high level of proficiency in MRA, achieving a score of 0.91, indicating its ability to effectively recall relevant information from previous exchanges. The RC score of 0.94 reflects the model's strong capability to provide correct answers to probing questions. However, there were instances where the model included unnecessary information in its responses. For example, in some cases when asked, "What was my first question?" the agent correctly identified the first question but also included additional information about the conversation flow, which was not required. These cases of excessive inclusivity had an impact on the RC score. The CC score of 0.9 suggests that the model can generate responses that are logically structured and contextually appropriate. Overall, the model achieved a robust performance score of 0.91, as detailed in Table 3.38, highlighting its effectiveness in maintaining context and coherence in short-term interactions, even though the evaluation was limited to the most recent five exchanges.

Table 3.38: Evaluation results of conversation history

Metric	Result
MRA	91%
RC	94%
CC	90%
Overall Performance	91%

3.4.2. AR Training Assistant

We first introduce the autonomous training agent architecture powered by large language models (LLMs) [Wang et al., 2023], aiming at a solution for training assistants in an Augmented Reality (AR) environment. The core idea is to adopt a new workflow (Figure 3.54) powered by

LLMs and driven by four functional submodules: memory, tools, planning, and actions. Specifically, a LLM, the brain of the dialogue agent, can plan the next action based on memory such as historical conversations (short-term) and knowledge base (long-term memory) and reaction of the executable tools that links to use case. To implement this architecture, we adopt LangChain¹⁶, the flexible abstractions and extensive toolkit for LLM applications.

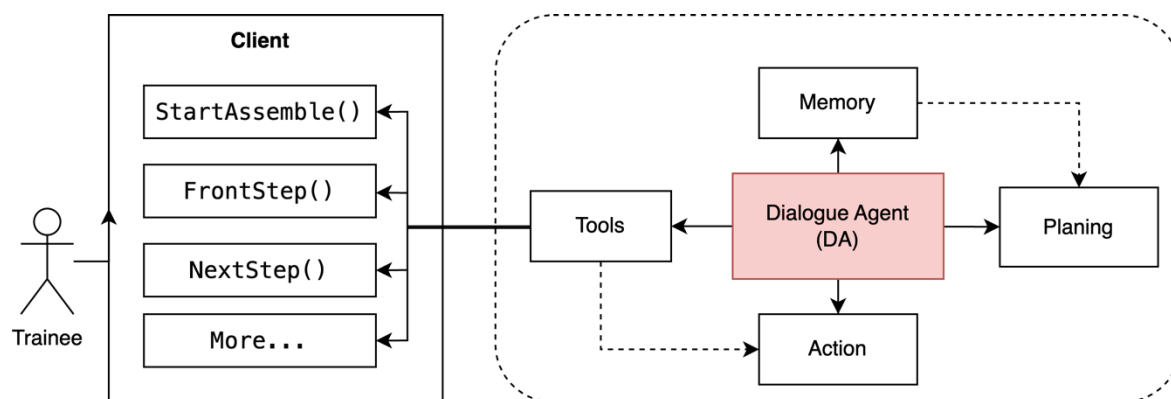


Figure 3.54. Autonomous agent powered by large language models

Workflow and Motivation

In this work, we aim to introduce a workflow (Figure 3.55): we use large language models (LLMs) as a human-like "brain" to augment XR systems with Artificial Intelligence; and vice versa, we enable XR systems to provide services by external API tools to augment general LLMs specifically for the AR Training Assistant (ARTA) use case, so that the AR system can better plan which tool to use and when to use it based on fully understanding the users' need.

The core idea of this workflow is the autonomous Dialogue Agent (DA), which can 1) utilize a language model such as Llama2 to conduct planning of a sequence of actions to take, 2) invoke a set of useful tools of real-world XR applications and observe the execution results and 3) make the decision of the next action based on all observations in the memory (e.g., historical utterances, response from XR environment, knowledge base, database, etc).

Particularly, the revolution of the autonomous agent powered by LLMs brings two benefits:

- Unlike general chains, where a sequence of actions is hardcoded (e.g., traditional dialogue pipeline with NLU followed NLG), the autonomous agents are powered by a language model and a prompt, which can serve as a reasoning engine to determine which actions to take and in which order automatically.
- The real-world applications can be learned to use automatically through reasoning, reducing the human's efforts on guidelines or defining training needs.

¹⁶ https://python.langchain.com/docs/get_started/introduction

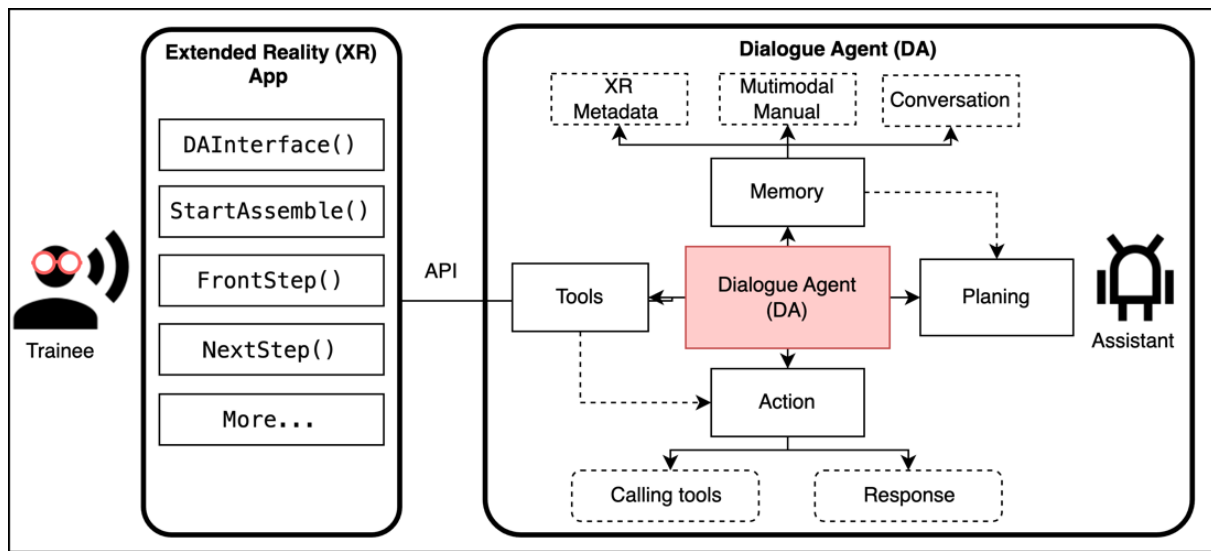


Figure 3.55. The workflow of ARTA

In this workflow, there are three key phrases to develop the model:

Phase 1: Develop the workflow and test the existing LLMs such as "gpt-3.5-turbo-16k-0613"¹⁷. It aims to investigate if this workflow based on the dialogue agent can integrate and immersively interact with a XR application seamlessly.

Phase 2: Create domain-specific datasets for the use case. We created datasets by two ways: (1) we reconstruct a TEACH-EDH dataset from human Commander-Follower conversations for household tasks in the VR simulation environment; (2) We generate and collect a VOX-ARTA dataset for LEGO brick assembly tasks in the XR simulation environment, using the workflow with a GPT3.5 and multimodal LEGO manuals.

Phase 3: Finetuning open-sourced LLMs such as "Llama2-chat-7b"¹⁸ with context-response pairs obtained from various datasets for teaching conversations (i.e., the reconstructed TEACH-EDH¹⁹ dataset), tool-calling and multimodal conversations (i.e., VOX-ARTA-LEGO dataset²⁰).

Model description

Following the workflow, we develop the model in terms of three aspects: 1) conversational agent walkthrough for planning and prompt engineering by fine-grained instructions 2) XR Application tools' API integration and 3) core LLM upgrade by finetuning the prevailing open-sourced LLM on use-case specific datasets.

Conversation agent walkthrough

We adopt a conversation agent walkthrough that optimizes the standard ReAct [Yao et al., 2023] with prompt engineering for conversation. This walkthrough demonstrates how to use an agent plan for conversation. LM can obtain a set of observations from the environment and decide the next action by reasoning the trace. The observations are a chain of information from diverse sources, e.g., historical conversations and prompts, XR tool responses, etc. The

¹⁷ <https://platform.openai.com/docs/models/gpt-3-5>

¹⁸ <https://huggingface.co/meta-llama/Llama-2-7b-chat>

¹⁹ https://huggingface.co/datasets/Jiahuan/teach_edh

²⁰ https://huggingface.co/datasets/Jiahuan/vox_arta_lego

final response is decided by observing the action of XR tools. Following is an example of one-loop chain-of-thought in the walkthrough:

Human: Can you show the detail of the roof?
AI: Action: ShowPieces
Action Input: roof
Response: Pieces available for roof: 2 x 2 red, 2 x 2 blue, 1 x 2 black, 1 x 2 white.
Please select which piece you want to use.

In contrast to many other agent walkthroughs that primarily focus on optimizing tool usage for determining the best response, our approach diverges, recognizing the importance of fostering conversational capabilities. This distinction becomes evident when comparing it to the conventional ReAct (Reason and Act) walkthrough, with the key disparity lying in the nature of the prompt. Our emphasis is on cultivating a much more conversational prompt to enhance the agent's ability to engage in meaningful dialogue with users. Specifically, we integrate task-specific prompt and XR tool responses in the observations so as to enable the agent to consider these in the reasoning traces for the decision of a next action. The walkthrough procedure is illustrated in Figure 3.56.

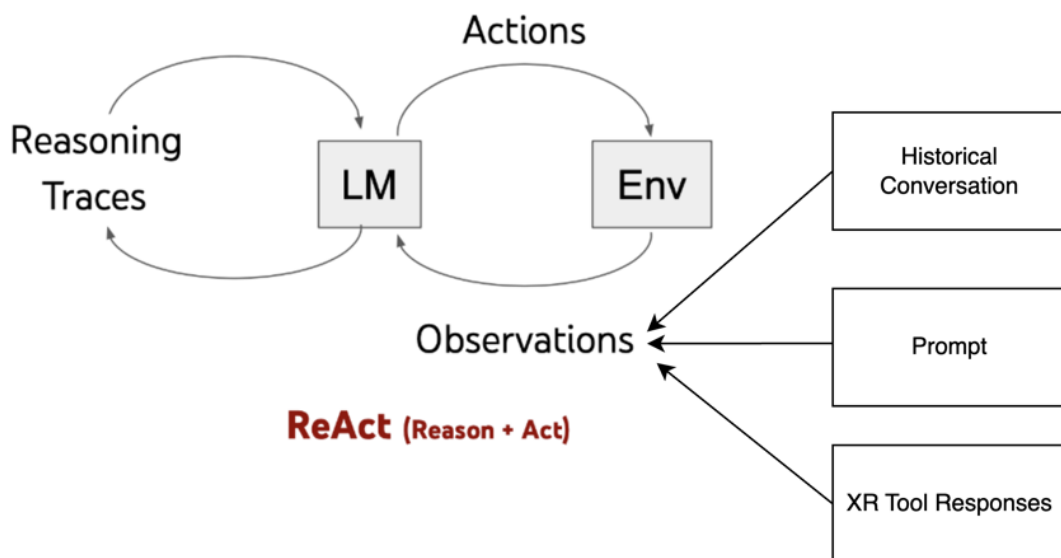


Figure 3.56. The illustration of conversational ReAct walkthrough

XR Application tools' API integration

We develop an XR test application and integrate the information interaction flow (Figure 3.57) between the application and the AI agent. This will be considered as a part of the whole workflow when the agent decides to call a specific tool for the user.

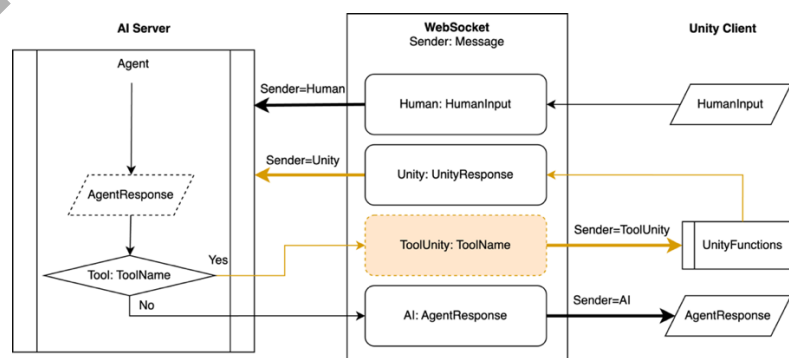


Figure 3.57. Interactive information flow between XR test application and AI agent

We develop a collection of executable interactive functions, presented in Table 3.39 within the XR test application, serving as external tools available for the agent. These tools play a crucial role in enabling agents to engage with the XR world and systematically navigate the XR environment. To facilitate the utilization of these functionalities, we have developed a toolkit—a set of API calls. This toolkit allows LLMs to access and invoke the services or functions outlined in the accompanying table, empowering them to strategize and determine the optimal for utilization within the XR system with minimum latency.

Table 3.39. Descriptions of the functions for the XR tools in the test application

No	Tool Name	Description
1	StartAssemble	Initiates the assembly process.
2	NextStep	Moves to the next assembly step.
3	FrontStep	Goes back to the previous assembly step.
4	Explode	Triggers an explosion for detailed viewing.
5	Recover	Restores the initial state of AR objects after an explosion.
6	FinishedVideo	Ends the assembly process and shows a video of the assembled LEGO bricks.
7	ReShow	Repeats the current assembly step.
8	Enlarge	Enlarges or zooms out the current object.
9	Shrink	Shrinks or zooms in the current object.
10	GoToStep	Goes to the given assembly step number.
11	Rotate	Rotates the current object to a direction.
12	ShowPieces	Shows all candidate LEGO pieces to be assembled.
13	HighlightCorrectComponents	Highlights correct attachment points and components.
14	GetCurrentStep	Gets the number of the current step.
15	GetRemainingStep	Gets the number of the remaining steps.
16	CheckStepStatusVR	Checks if the current step in Unity is accomplished correctly or not.
17	APICallObjectRecognitionAR	Calls the VLM agent to identify LEGO pieces based on the provided video streaming data from AR glasses and highlights the recognized pieces in the AR environment.
18	APICallCheckStepStatusAR	Calls the VLM agent to determine if the current assembly step is completed correctly or not, using the provided video streaming data from AR glasses as input.

Parameter-Efficient Fine-Tuning of Quantized LLMs

Fully finetuning an LLM requires powerful computation resources as a result of the recent growth of LLM sizes. For example, the smallest version of Llama2 contains 7 billion learnable parameters that require a GPU with ~28GB memory. To be able to utilize LLMs in the GPU we use to perform our developments (GeForce RTX 3090 with 24GB), we adopt 4-bit Quantized Low Raw Adaptor (QLoRA) [Dettmers et al, 2022] [Dettmers et al, 2023]. We use this approach for finetuning the quantization of Llama2-7b-chat and conduct parameter-efficient fine-tuning (PEFT) of the quantized model with Low Raw Adaptor (LoRA). The comparison of the graphical representation of full fine-tuning, LoRA and QLoRA is presented in Figure 3.58.

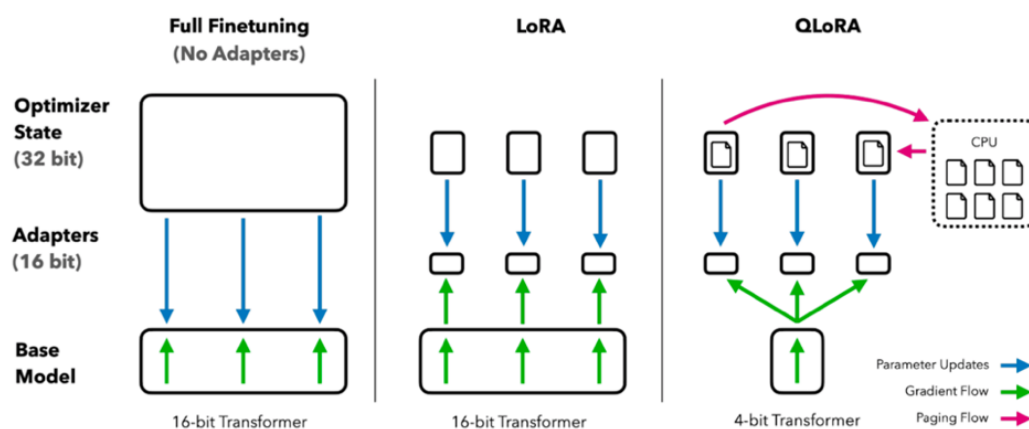


Figure 3.58. Comparison of full fine-tuning, LoRA, and QLoRA

Dataset Construction and Statistics

In VOXReality, we produced two datasets that can be used for teaching assistant conversation agents, TEACH-EDH and VOX-ARTA-LEGO.

TEACH-EDH

We reconstruct the TEACH-EDH (Task-driven Embodied Agents that Chat - Execution from Dialog History) dataset by curating conversations between human Commanders and human Followers (Drivers) who collaboratively use natural language to accomplish household tasks within the AI2-THOR simulation environment [Kolve et al., 2017]. Each data sample consists of all historical utterances as an “input” and the current utterance as an “output”. We construct 14,740 context-response pairs for TEACH-EDH dataset, including 9740 training, 1140 validation and 3860 test data. We represent examples of the dataset in Figure 3.59.

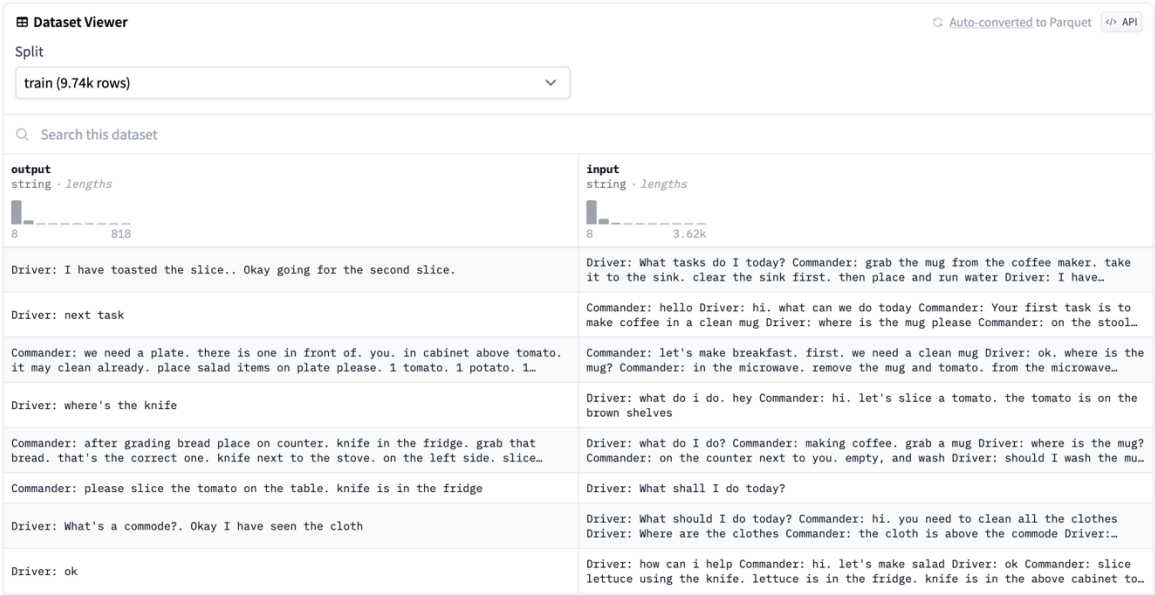


Figure 3.59. Examples of TEACH-EDH dataset

VOX-ARTA-LEGO

We generate and collect conversations between a user (trainee) and an AI teaching assistant (AI trainer) following the proposed workflow, which is powered by a commercially available LLM using GPT-3.5-turbo-16k-0613 and multimodal manuals for LEGO brick assembly. Specifically, we first performed web crawling in 65 multimodal manuals that are openly available on LEGO’s official website. Every manual consists of a summary and multiple assembly steps with both text instructions and images. This provides the fine-grained alignment between vision and language that are important to teach users how to assemble a LEGO brick set with a step-by-step guide. Then, we design prompts using the manuals and XR tools to generate textual simulated conversations. Similar to TEACH-EDH, each data sample consists of all historical utterances as an “input” and the current utterance as an “output”. We construct 866 context-response pairs for VOX-ARTA-LEGO dataset, including 622 training, 70 validation and 174 test data samples. We present examples of the dataset in Figure 3.60.



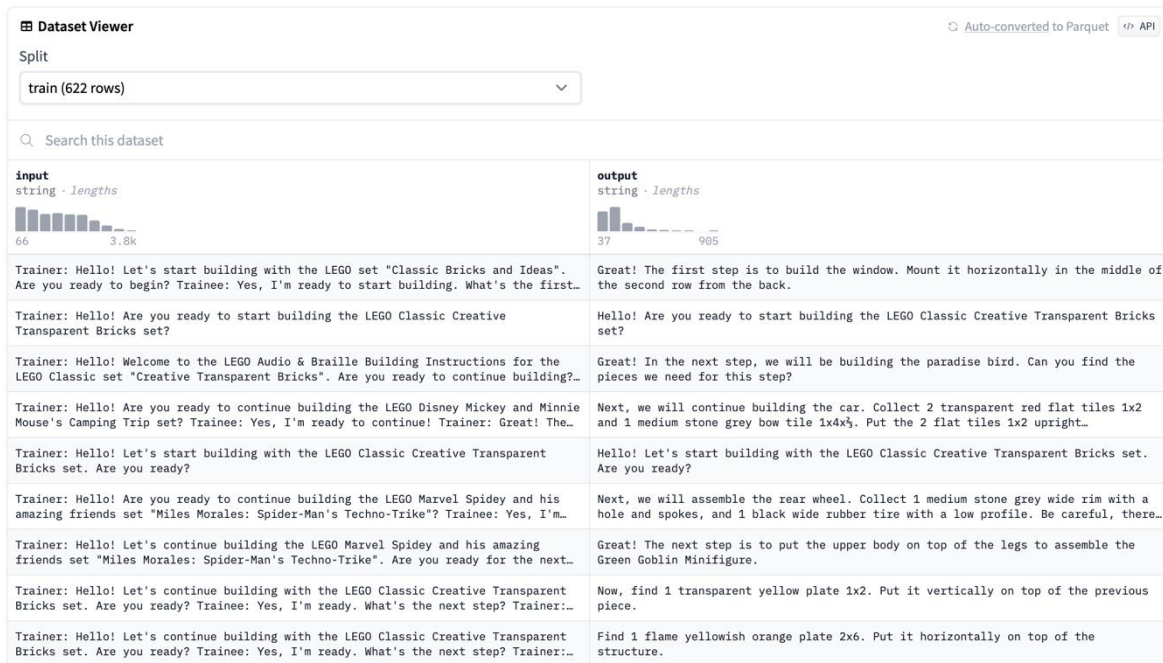


Figure 3.60. Examples of VOX-ARTA-LEGO dataset

Analysis of the datasets

We conduct a comparative analysis of the statistics, as presented in Table 3.40, and delve into the specifics of the two proposed datasets. Regarding the statistical aspects, the TEACH-EDH dataset boasts a larger number of data samples, featuring a higher count of dialogues and utterances. Conversely, the VOX-ARTA-LEGO dataset exhibits longer utterances from both users and systems. From a linguistic point of view, the utterances within the VOX-ARTA-LEGO dataset resemble a more human-like language, characterised by rich and diverse descriptions. On the other hand, the TEACH-EDH dataset leans more towards a machine-like language expression.

Table 3.40. Statistics of TEACH-EDH and VOX-ARTA-LEGO datasets

Phase	TEACH-EDH	VOX-ARTA-LEGO
Before	98	98.1
Context-Response	14,740	866
Dialogue	3320	171
Utterance	45,000	1,670
Avg Context Length	49.81	1,302
Avg Response Length	10.90	184
Avg Utterance Per Dialogue	13.67	9.77

Experimental Setup

Task definition

In the context of generative dialogue agents, the task is to generate responses to the user (trainee) given a sequence of diverse dialogue contexts, including historical utterances from trainer and trainee, multimodal manual, and XR metadata such as hand gesture input and head pose. In general, this task can be seen as a context-response generation task.

Evaluation metrics

We adopt the following metrics to evaluate models on both TEACH-EDH and VOX-ARTA-LEGO datasets.

- Perplexity is a commonly used measure of how well a language model predicts a sample. Lower perplexity indicates a better model, as it measures the average uncertainty of the model in predicting the next word in a sequence.
- ROUGE (described in Section 2.5) is a set of metrics used for the automatic evaluation of generated response by comparing it to a set of reference responses.
- METEOR (described in Section 2.5) incorporates various linguistic and semantic features to assess the quality of generated responses compared with reference responses.

Implementation details

The base LLM we use is the open-resourced Llama-2-7b-chat-hf²¹, which is a collection of pretrained and finetuned generative text models ranging in scale from 7 billion and up. We use 4-bit quantization of the model and conduct PEFT with LoRA to obtain finetuning of 4.43% of original trainable parameters. The batch size is 8 and the maximum step is 2,000 with early stopping when the validation loss does not decrease with continuous 10 steps. The maximum number of tokens is 512.

Results

We evaluate the Llama2 and its finetuned models on two proposed datasets and represent both quantitative results and qualitative analysis, respectively.

Quantitative Results

We present the evaluation results for the two proposed datasets both before and after finetuning with QLoRA, as outlined in Table 3.41. Perplexity, ROUGE, and METOR metrics were employed for assessment on the test set of both TEACH-EDH and VOX-ARTA-LEGO dataset. A significant reduction in perplexity is observed for both datasets, plummeting from 42.18 to 3.08 for the TEACH-EDH dataset and from 8.81 to 1.06 for the VOX-ARTA-LEGO dataset. This indicates an enhanced model performance in predicting the next token with reduced uncertainty. Lower perplexity corresponds to a heightened proficiency in next token prediction. However, the ROUGE and METOR metrics show a decline on both datasets following QLoRA finetuning. Specifically, ROUGE experiences a 3.42% drop in TEACH-EDH dataset and 5.63% drop in VOX-ARTA-LEGO dataset, indicating a reduction in the overlap of 1-gram tokens between generated and reference responses. METEOR exhibits a slight drop in both the TEACH-EDH dataset (2.01%) and the VOX-ARTA-LEGO dataset (1.09%), suggesting a decline in linguistic and semantic similarity between generated and reference responses. This decrease can be attributed to two potential factors: firstly, QLoRA finetunes only 4.43% of the original trainable parameters, necessitating a substantial volume of data to effectively influence the limited parameter adjustments. Secondly, the quality of fine-grained, human-like data is crucial; otherwise, the introduction of new data may introduce noise to the original model.

Table 3.41. Evaluation results on TEACH-EDH and VOX-ARTA-LEGO dataset

Model	TEACH-EDH			VOX-ARTA-LEGO		
	Perplexity	ROUGE	METEOR	Perplexity	ROUGE	METEOR
Llama2	42.18	8.17	12.08	8.81	23.36	12.42
Llama2+QLoRA	3.08	4.74	10.07	1.06	17.73	11.37

²¹ meta-llama/Llama-2-7b-chat-hf



Qualitative Analysis

We scrutinize the learning curves depicting the loss evolution for both training and validation datasets throughout the finetuning process, illustrated in Figure 3.61, where training loss is coloured blue and validation loss is coloured red. In the case of the TEACH-EDH dataset, the training loss decreases from 3.732 to 1.202, while the validation loss decreases from 3.742 to 1.188. The model employs early stopping, interrupting the training at the 67th step in the 2nd epoch, as the validation loss shows no continuous improvement over 10 consecutive steps. In the case of the VOX-ARTA-LEGO dataset, the training loss exhibits a notable reduction, decreasing 2.116 to 0.177 while the validation loss decreases from 2.176 to 0.060. The model adopts early stopping at 60th step in the 25th epoch.

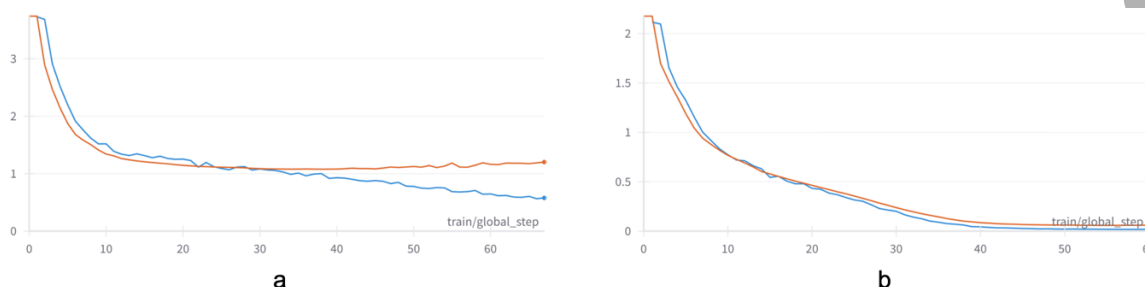


Figure 3.61. The learning curve of Llama-2-7b-chat-hf model finetuning on a) TEACH-EDH dataset and b) VOX-ARTA-LEGO

Case Study

We discuss the following showcases to analyse the responses generated by the model and provide an example of a walkthrough to illustrate how the reasoning of traces works with observations of actions. The following are two examples of the inference results from the test set, where the agent can observe actions, action inputs, and historical utterances from both human and assistant. In comparison with the reference responses, the model tends to generate longer responses with richer context and instructions. This motivates us to conduct a human evaluation to reconsider whether the decrease in ROUGE and METEOR scores is a fair metric for evaluating generation performance.

Observations #1

Action: Guide

Action Input: Yes

Response: Sure thing! To find another transparent bright orange roof tile 1x2, 45°, I suggest checking the nearby shelves or drawers. Once you locate it, put it horizontally on the left knob, overhanging to the left and sloping to the left. Could you please let me know once this item is found?

Human (Trainee): No, I followed the instructions and everything looks good so far. What's the next step?

Assistant: Great! In that case, the next step is to attach the control levers to the sides of the vehicle. Can you please find two transparent bright orange roof tiles 1x2, 45°? Place them horizontally on the left and right knobs respectively, overhanging to the left and sloping to the left. Once you have located these tiles, please let me know so I can assist you further.

Prediction #1

Great! In that case, the next step is to attach the control levers to the sides of the vehicle. Can you please find two transparent bright orange roof tiles 1x2, 45°? Place them horizontally on the left and right knobs respectively, overhanging to the left and sloping to the left. Once you have located these tiles, please let me know so I can assist you further.

Reference #1

That completes this step! You've successfully built the control levers and added the plates. Is there anything you're unsure about or any questions you have?

Observations #2

Action: Guide

Action Input: Yes

Human: Can you find a bright red plate 1x2 with a single stud on top?

Response: Of course! To find a bright red plate 1x2 with a single stud on top, I suggest checking the nearby shelves or drawers. Once located, please place it horizontally on the table. Could you please let me know once this item is found?

Prediction #2

Of course! To find a bright red plate 1x2 with a single stud on top, I suggest checking the nearby shelves or drawers. Once located, please place it horizontally on the table. Could you please let me know once this item is found?

Reference #2

Great! Let's proceed to the next step.

The following is an example instance of an interactive conversation session with actual human input. The current agent demonstrates an ability to comprehend the logic behind tool invocation and provides sensible responses in alignment with the overall context of the conversation.

Human: Hi can you start assemble the LEGO car?

Assistant: Sure thing! I'd love to help you with that. Which type of LEGO car would you like to build? We have different options available. Once you've made your choice, I can guide you through each step of the assembly process. Let me know what you prefer!

Human: OK, got it. Can you go to next step?

Assistant: Great! Let's proceed to the next step. Initiating the assembly process...

Action: StartAssemble()

Action Input: None

Response: Assembly process started. Please follow the instructions carefully. Please provide the actual instruction manual text for the LEGO 11001 Classic set so I can generate the appropriate responses.

Human: Can you go to step 3?

Action: NextStep()

Action Input: None

Response: Moving to step 3...

Human: Then what should I do next step?

Action: NextStep()

Action Input: None

Response: Moving to step 4...

Human: Can you tell me something not related to LEGO assembly task?

Assistant: Sorry, I cannot fulfill that request. My primary function is to assist users with their LEGO assembly tasks. If you have any questions or requests related to the task, please feel free to ask. Otherwise, I'll be here to help whenever you need it. How may I assist you today?

Human: Can you show the detail of the roof?

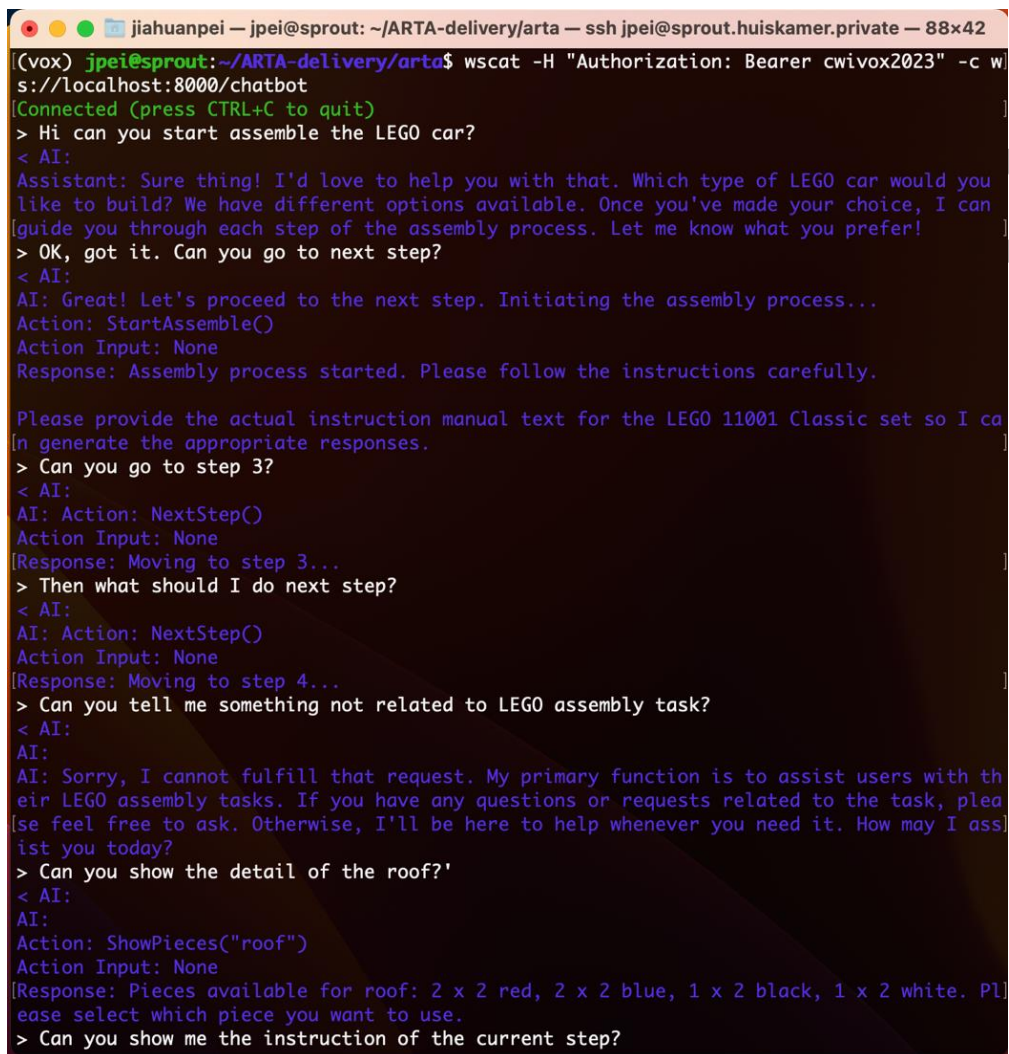
Action: ShowPieces("roof")

Action Input: None

Response: Pieces available for roof: 2 x 2 red, 2 x 2 blue, 1 x 2 black, 1 x 2 white. Please select which piece you want to use.

Deployment

We deploy our model through the WebSocket API (an example run can be seen in Figure 3.62), encapsulated within a Docker Image, to facilitate interactive communication between the user (trainee) and the assistant (AI trainer). We also prepared a deployment using FastAPI, a robust framework that enables seamless and efficient handling of the communication interface. The details can be found in Appendix IV.



```
jiahuanpei — jpei@sprout: ~/ARTA-delivery/arta — ssh jpei@sprout.huiskamer.private — 88x42
(jvox) jpei@sprout:~/ARTA-delivery/arta$ wscat -H "Authorization: Bearer cwivox2023" -c w
s://localhost:8000/chatbot
[Connected (press CTRL+C to quit)]
> Hi can you start assemble the LEGO car?
< AI:
Assistant: Sure thing! I'd love to help you with that. Which type of LEGO car would you
like to build? We have different options available. Once you've made your choice, I can
guide you through each step of the assembly process. Let me know what you prefer!
> OK, got it. Can you go to next step?
< AI:
AI: Great! Let's proceed to the next step. Initiating the assembly process...
Action: StartAssemble()
Action Input: None
Response: Assembly process started. Please follow the instructions carefully.

Please provide the actual instruction manual text for the LEGO 11001 Classic set so I ca
n generate the appropriate responses.
> Can you go to step 3?
< AI:
AI: Action: NextStep()
Action Input: None
Response: Moving to step 3...
> Then what should I do next step?
< AI:
AI: Action: NextStep()
Action Input: None
Response: Moving to step 4...
> Can you tell me something not related to LEGO assembly task?
< AI:
AI:
AI: Sorry, I cannot fulfill that request. My primary function is to assist users with th
eir LEGO assembly tasks. If you have any questions or requests related to the task, plea
se feel free to ask. Otherwise, I'll be here to help whenever you need it. How may I ass
ist you today?
> Can you show the detail of the roof?'
< AI:
AI:
Action: ShowPieces("roof")
Action Input: None
Response: Pieces available for roof: 2 x 2 red, 2 x 2 blue, 1 x 2 black, 1 x 2 white. Pl
ease select which piece you want to use.
> Can you show me the instruction of the current step?
```

Figure 3.62. Interactive WebSocket API

Updating AR Training Assistant

The updated AR Training Assistant introduces key improvements in its workflow architecture to enhance its accuracy and functionality. A major update in this version is the replacement of the fine-tuned **Llama-2-7B-Chat** model with **Mistral-7B-Instruct-v0.3**, which was selected due to its superior performance in handling dialogue interactions. In our implementation, Mistral-7B-Instruct-v0.3 demonstrated improved response quality, better task understanding, and more efficient handling of training-related queries.

The AR Training Assistant has evolved through two distinct versions: an initial prototype (Linear Mode) and the final version (Free Mode). The Linear Mode was developed as a structured, sequential training approach where users followed a predefined sequence of actions without deviation. However, after further refinement and testing, this mode was removed, and the final version of the system only includes Free Mode, which allows for a more flexible and dynamic training experience. The workflow structure is illustrated in Figure 3.63.

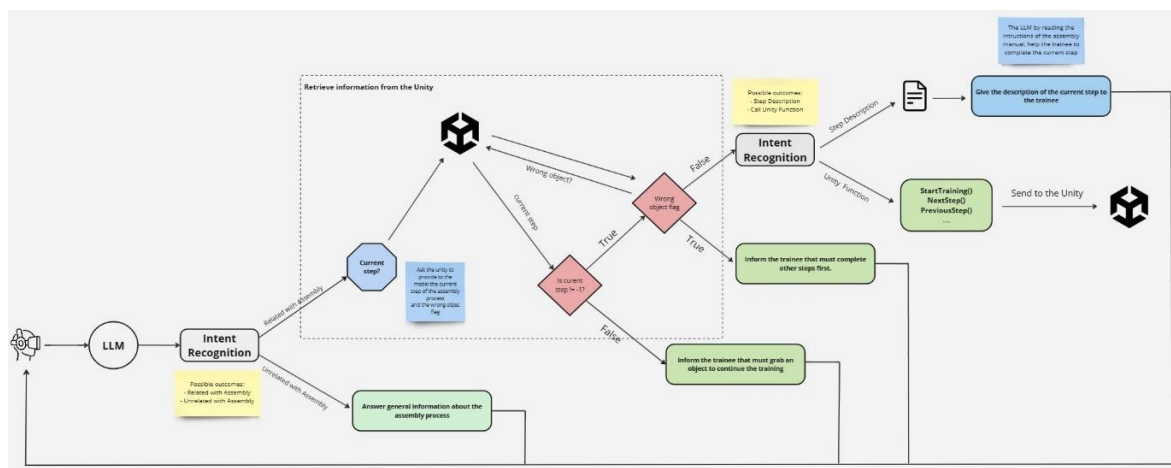


Figure 3.63: Workflow of the free mode of Training Assistant

In the Linear Mode, users interacted with the XR application in a strictly linear sequence. When a user issued a command, the model performed basic intent recognition to determine if the question was related to the assembly process or out of scope. If the query was unrelated, the model informed the user that it could only handle assembly-related questions. If the question was relevant, the XR application returned a flag indicating the user's current assembly step. The model then identified whether the user was requesting step instructions or wanted to trigger a Unity function (e.g., start assembly, go to the next step, etc.). If step instructions were requested, the model provided guidance; otherwise, it sent the appropriate command to Unity.

While Linear Mode provided a structured approach, it lacked flexibility, limiting users to a strictly predefined order of interactions. To address this, the system was redesigned to introduce a more adaptive workflow, Free Mode, which is the only mode included in the final version of the system. In Free Mode, the first and last phases remain similar to Linear Mode, where the system performs basic intent recognition at the start and either provides a step description or triggers a Unity function at the end. However, Free Mode introduces additional decision-making steps to create a more dynamic and interactive training experience.

After determining that a query is related to assembly, the XR application returns a flag indicating the user's current step. If this flag is -1, it means the user has not yet grabbed an object, and the model instructs them to do so before proceeding. If the user has grabbed an object, the system then verifies whether it is wrong or not. If the user selected the wrong object, the model notifies them that they must complete previous steps first. If the correct one is detected, the process continues as in Liner Mode, performing a second intent recognition step before either describing the step or triggering a Unity function.

By implementing Free Mode as the final version, the system now offers a more flexible and intelligent training experience, allowing users to navigate the assembly process dynamically rather than following a rigid sequence. In the next subsections, we will provide a detailed explanation of each component of the Free Mode of AR Training Assistant, covering their implementation, functionality, and role within the overall workflow.

First intent recognition phase

The intent recognition module plays a crucial role in the Free Mode workflow of the AR Training Assistant by distinguishing between assembly-related queries and out-of-scope questions. To achieve this, a **few-shot learning** approach is employed, enabling the model to generalize effectively with a limited set of annotated examples. When a user submits a question, the system determines whether it pertains to the assembly process or is unrelated. If the query is

related to assembly, the model assigns it a classification of "1", indicating that the system should proceed with further processing. If the query is unrelated, the model classifies it as "0", signaling that the assistant should inform the user that it only handles assembly-related tasks. This classification process of the few shot examples can be demonstrated through the following examples:

```
{user_prompt: "What I must do now?", "response": 1},
{user_prompt: "Show a hint for this", "response": 1},
{user_prompt: "What's the weather today?", "response": 0},
{user_prompt: "What's your name?", "response": 0}
```

To ensure consistency and accuracy, a system prompt is incorporated into the classification process. The system prompt explicitly instructs the model on how to handle different types of queries, ensuring that responses are generated in a structured and predictable manner. The prompt provides a clear directive to the model:

- If the user’s question is related to the assembly process, the response should be "1", allowing the system to continue processing the request.
- If the question is unrelated, the model must respond with "0", without adding any additional information.

This constraint-based prompting strategy ensures that the model's classification remains consistent and unambiguous, preventing unnecessary variations in response formatting. By keeping the output strictly binary, the system reduces processing complexity and minimizes errors, leading to a more robust and efficient interaction process. Additionally, the few-shot examples used in the prompt serve as a guiding reference for the model, showcasing a range of sample queries and their corresponding classifications. This structured example-based learning reinforces the intended behavior of the system, ensuring that the model reliably distinguishes between relevant and irrelevant queries across different conversational scenarios. By combining few-shot learning with a well-defined system prompt, the intent recognition module significantly enhances the overall performance of the AR Training Assistant, ensuring that users receive accurate and contextually appropriate assistance throughout the assembly process.

Evaluation Process

To assess the effectiveness of the intent recognition module, a structured evaluation process was conducted. A test dataset comprising 50 labeled examples (with 25 samples per class) was created to rigorously measure the model's ability to differentiate between assembly-related and unrelated queries. Each query from this dataset was presented to the model, and its classification output was compared against the predefined ground truth labels.

The evaluation utilized standard classification metrics, including accuracy, precision, recall, and F1-score, to quantify the model's performance. The results demonstrated that the model achieves high reliability in distinguishing between relevant and irrelevant queries, confirming the effectiveness of the few-shot learning approach and system prompt constraints. The model's F1-score reached 95%, indicating strong overall classification performance as demonstrated in Table 3.42.

Table 3.42: Evaluation results of the first intent recognition phase

Metric	Result
Accuracy	96%
Precision	94.8%
Recall	95.2%
F1-score	95%



These results confirm that the intent recognition module operates with high precision and consistency, effectively distinguishing between assembly-related and unrelated queries. Given the structured few-shot learning approach and the well-defined system prompt, the model demonstrates a strong ability to generalize and maintain robust classification accuracy across various conversational inputs.

Handling unrelated queries

This module ensures that the AR Training Assistant remains focused on its primary function guiding users through the Raptor Engine assembly process. This component is responsible for filtering out unrelated queries while maintaining a structured and contextually relevant interaction.

To ensure the assistant maintains a clear and structured interaction, the model is provided with a system prompt that defines its identity and purpose. This prompt explicitly instructs the model to introduce itself as "VOXY" and describe its role in assisting users with the Raptor Engine assembly process. Additionally, it is guided to encourage users to use imperative language when requesting actions, such as "Start the process" or "Go to the next step.". The reason that the model is designed to inform the user about the imperative format will be described in the second intent recognition phase module. The system prompt includes the following instruction:

"You are a Dialogue Agent powered by a large language model often addressed as 'VOXY'. Your purpose is to assist the user with the assembly process for the Raptor Engine, step by step. If the user wants to start the process or proceed to the next step, inform them to use imperative language, such as 'Start the process' or 'Go to the next step.' Always be polite and clear in your responses. If the user asks an unrelated question, kindly let them know that you are here solely to help with the assembly process."

This predefined instruction ensures that the model remains on task while maintaining a natural and engaging conversation. If a user inquires about VOXY's capabilities, the model responds by clearly outlining its function, ensuring that users understand how to interact with the system effectively. By implementing this structured approach, the assistant enhances user experience by providing relevant, informative, and task-focused responses while gracefully handling out-of-scope queries.

Retrieving the current step flag

A crucial component of the Free Mode workflow is determining the user's current step in the assembly process. Unlike intent recognition, which classifies user queries based on predefined categories, this information cannot be inferred directly by the AR Training Assistant. Instead, it must be retrieved from the Unity environment, which tracks the user's actions in real time.

Once the intent recognition module classifies a user query as assembly-related, the next step in the workflow is for the system to request the current step flag from Unity. This flag plays a critical role in guiding the model's decision-making and ensuring that responses align with the user's actual progress in the assembly process. The current step flag is essential for two key reasons:

- **Step-Specific Instructions:** When the assistant provides guidance, it must generate step-specific descriptions. Without knowing the user's exact position in the assembly sequence, the model would be unable to deliver precise instructions. By retrieving the current step flag, the assistant ensures that users receive relevant, context-aware guidance tailored to their progress. A more detailed explanation of why is so important for the agent to retrieve the current step information and how the agent generated

assembly instructions will be provided later in the final component, where the assistant delivers direct guidance to the user.

- **Object Handling Verification:** The assistant can only provide step descriptions if the user has already grabbed an object. If the user has not yet picked up an object, Unity returns a special "-1" flag, signaling to the model that it must first instruct the user to pick up an object before proceeding with assembly guidance.

By incorporating real-time step tracking via Unity, the AR Training Assistant enhances its ability to provide accurate, structured, and context-sensitive assembly guidance. This integration ensures a seamless and adaptive interaction, allowing users to progress through the training process efficiently while receiving relevant assistance at every step.

Handling the no object grabbed case

A critical challenge in the Free Mode workflow of the AR Training Assistant arises when the user has not yet picked up an object. Since object interaction is a prerequisite for meaningful assembly guidance, the system must detect this scenario and adjust its response accordingly. This is achieved through a dedicated model behavior, triggered when the current step flag returned by Unity is "-1". In this case, the model does not proceed with standard step descriptions. Instead, it follows a carefully designed system prompt, ensuring that users are first instructed to grab an object before continuing. This approach guarantees that the assistant maintains a logical and structured guidance process, preventing confusion or premature instructions.

To ensure consistency, the system prompt enforces several key constraints. The model must always begin by reminding the user to grab an object before proceeding, without providing detailed assembly instructions. Even if the user requests a hint or guidance on the next step, the response must prioritize object interaction before addressing any assembly-related queries. Additionally, responses are generated in a direct and concise manner, without unnecessary explanations or contextual details. Below is the final system prompt that used to the model:

"You are "VOXY", an XR assistant that helps users assemble a Raptor engine.

Important rules:

- 1. The user hasn't grabbed an object yet. You must always start your response by reminding them to grab an object. Do not provide detailed assembly instructions or hints until they have done so.*
- 2. Even if the user's question requests a hint, instruction, or guidance on the next step, your response should first emphasize the need to grab an object.*
- 3. Provide only the final answer message without any additional text, context, or explanations."*

By implementing this specialized model behavior, the system ensures that users remain aligned with the natural flow of the assembly process while preserving the flexibility of Free Mode. This mechanism enhances user experience by reinforcing structured guidance without imposing rigid constraints, allowing for intuitive and adaptive interaction within the XR environment.

Retrieving the wrong object flag

In the Free Mode workflow of the AR Training Assistant, users may sometimes grab an object that is not required for the current step of the assembly process. Since certain steps in the XR environment are object-dependent, some actions cannot be performed unless the necessary prerequisites have been completed. However, the agent itself does not have direct awareness of which object the user has picked up. To address this, it queries the Unity engine, which returns a "wrong object flag" with a Boolean value.

If the wrong object flag is “*True*”, it indicates that the user has selected an incorrect object that cannot be used at the current step. In this case, the assistant does not proceed with step instructions. Instead, it informs the user that they must complete previous steps first before continuing. This ensures that users follow the correct assembly sequence and do not attempt to perform actions that require prior steps to be completed.

If the wrong object flag is “*False*”, the user has grabbed the correct object, and the system proceeds as usual. The assistant then performs the second intent recognition step, determining whether the user is requesting step instructions or wants to trigger a Unity function related to the assembly process.

By dynamically validating object selection through Unity, the system maintains the integrity of the assembly process while preserving the flexibility of Free Mode. This approach ensures that users receive accurate, context-aware guidance, preventing errors and reinforcing a structured yet adaptive training experience within the XR environment.

Handling the wrong object grabbed case

The approach for handling cases where the user has grabbed the wrong object, the wrong object flag is *False*, follows the same logic as the “Handling the no object grabbed case” that we described before. Instead of proceeding with standard step instructions, a different system prompt is provided to the model to ensure users are first informed that they must complete a necessary prior step or select the correct component before continuing. This ensures that the assistant provides structured and relevant feedback, maintaining a logical assembly flow while preventing incorrect actions. The system prompt used in this case is:

*“You are an assistant named ‘VOXY’ guiding the assembly of the Raptor Engine.
It appears that a necessary previous step has not yet been completed before you can execute the current action.
Please inform the user in a friendly and clear manner that they need to finish the prior step or pick up the correct component before proceeding safely.”*

Second Intent recognition

The second intent recognition phase in the Free Mode workflow is responsible for classifying user queries into two categories: (1) requests for step instructions and (2) commands to trigger actions in Unity. This distinction is essential, as it determines whether the assistant should provide assembly guidance or communicate with the XR environment to execute a specific function. Classifying user intent in this phase proved to be particularly challenging. Through extensive testing of different few-shot examples, we observed that the model struggled to clearly differentiate between the two classes. This difficulty occurred because instructional queries and commands often had overlapping language patterns. For instance, a user might ask:

- “What must I do now?” or “Can you explain what to do?” seeking **step instructions**
- “Can you start the training?” or “I want you to go to the next step please.”—issuing a **command to Unity**

To overcome this, we refined our implementation by introducing a clear linguistic distinction by using imperative statements (direct commands) that indicate the user wants to trigger a Unity function and non-imperative queries that indicate that the user seeks assembly instructions. To ensure accurate classification, users are explicitly informed upon entering the XR environment that Unity actions must be triggered using imperative commands (e.g., “Start the process”, “Go to the next step”). This standardization significantly improves classification

accuracy. By enforcing these rules, the second intent recognition phase successfully differentiates between user inquiries and execution commands, allowing the system to provide accurate assembly instructions or seamlessly communicate with Unity to trigger specific functions. The final system prompt used to guide the model is as follows:

*"You are a classification system that determines whether a user's query is an IMPERATIVE command or not.
If the query starts with 'start', 'go to', 'take me back', 'show me', or 'give me', classify the query as an imperative command.
Non-imperative queries are descriptive or informational questions (for example, 'How do I start the training?' or 'What is the next step?').
Ignore any differences in letter casing and punctuation.

IMPORTANT: Do not include any explanations, chain-of-thought, or additional text.
If the query is an imperative command, classify it as 1; otherwise, classify it as 0.
Output ONLY the digit '0' or '1' with no extra text, punctuation, or explanation."*

To evaluate the effectiveness of the second intent recognition phase, a structured testing process was conducted using a 50-sample dataset (with 25 examples per class). Each test query was classified by the model and compared against its respective ground truth label to measure accuracy. This phase presented greater challenges than the first intent recognition phase, as imperative commands and step-related inquiries often share similar linguistic structures. However, by refining the few shot learning approach and explicitly defining classification rules in the system prompt, the model was able to significantly improve its accuracy in distinguishing between step instructions and commands to trigger Unity functions. To quantify performance, four key evaluation metrics were computed: accuracy, precision, recall, and F1-score. These results, presented in Table 3.43, demonstrate that despite the linguistic overlap between classes, the system successfully identifies imperative commands versus step-related queries with high accuracy and consistency. The F1-score of 92.3% confirms that the model effectively generalizes across different phrasing structures, ensuring that user commands are appropriately handled.

Table 3.43: Evaluation results of the second intent recognition phase

Metric	Result
Accuracy	93.5%
Precision	91.8%
Recall	92.9%
F1-score	92.3%

Triggering Unity Functions

Once the second intent recognition phase determines that the user's query is an imperative command, the next step is to identify the appropriate Unity function that must be triggered. In the XR application, a set of predefined functions handle different aspects of the training process, ensuring smooth user interaction within the virtual environment.

The available Unity functions include:

- **StartTraining():** Initiates the entire assembly training process, setting the system into training mode and guiding the user through the steps.
- **NextStep():** Advances the training sequence to the next step, updating the visual and instructional components accordingly.
- **PreviousStep():** Moves the process back to the previous step, allowing users to review and correct any mistakes.



- **ShowHint():** Highlights the relevant components and areas in the XR environment, providing visual assistance when the user is lost or unsure about what to do next.
- **ShowVideo():** Plays a step-specific instructional video, offering a detailed demonstration of the required action for the current stage of assembly.

To ensure accurate classification, we employed a few-shot learning approach, where the model was provided with carefully selected examples for each Unity function. By supplying the model with two distinct examples per function, we improved its ability to distinguish between similar commands and correctly match user intent with the corresponding Unity function. By leveraging this structured classification approach, the system ensures that user commands are processed efficiently, allowing for a seamless and interactive training experience. When the model recognizes an imperative command, it does not provide a direct response to the user but instead sends the corresponding Unity function, enabling real-time execution of training actions in the XR environment. The following examples illustrate how user queries are mapped to the appropriate Unity function:

```
"Hi Voxy, can you tell me how to start the assembly process?" → StartTraining()
"How to begin the assembly process?" → StartTraining()
"Great, how do I proceed?" → NextStep()
"Alright, I did that. What is the next step, Voxy?" → NextStep()
"Perfect Voxy, can you repeat the previous step?" → PreviousStep()
"Great Voxy, what was the previous step again?" → PreviousStep()
"Show a hint for this." → ShowHint()
"Show me a cue." → ShowHint()
"Show me what I should do with this." → ShowVideo()
"Video of this step." → ShowVideo()
```

Providing Step-Specific Assembly Instructions

The final stage of the Free Mode workflow involves delivering step-specific assembly instructions when the second intent recognition phase determines that the user is requesting guidance rather than triggering an action. For the model to generate accurate and relevant instructions, it must first determine which step the user is currently on. This critical information is retrieved during the "Retrieving the Current Step Flag" phase, ensuring that responses are directly aligned with the user's progress in the assembly process.

Once the system has identified the user's current step, it must retrieve the corresponding instructions. To facilitate this, a pre-processed assembly manual is provided to the model in a structured format. This manual, stored as a text file, contains step-by-step descriptions of the Raptor Engine assembly process. The content follows a standardized structure, ensuring clarity and consistency in instructional responses. An example of the manual's format is as follows:

- **Step 1:** Positioning the Pressure Transmitter Stacks. These pressure transmitters are crucial for monitoring engine pressure variations and ensuring stable performance. Locate the three pressure transmitter stacks on the workspace table. Identify the dedicated slots in the holder on the table. Grab a pressure transmitter stack and carefully place it into the left slot in the holder.
- **Step 2:** Grab the second pressure transmitter stack and place it into the middle slot in the holder on the table.
- **Step 3:** Grab the last pressure transmitter stack and place it into the rightmost slot of the holder on the table.

To efficiently match the user's current step with the correct instructions, a pre-processing step converts the text manual into a dictionary format. In this dictionary, each step number serves as a key, while the corresponding step description functions as its value. When a user requests guidance, the system accesses this dictionary, retrieves the correct step-specific instructions, and then forwards them to the model for response generation.

A key advantage of this approach is its adaptability to different instructional workflows. Since the model does not rely on hardcoded step descriptions, it can be applied to various assembly processes, training sequences, or even navigation guidance, as long as the instructions follow the same structured format. This means that by simply replacing the manual with another document, formatted similarly, the model can generate accurate instructions for different tasks or environments, making it a highly versatile solution for training and assistance applications.

To ensure structured and user-friendly responses, a system prompt is employed. This prompt defines the assistant's role and instructs it to provide clear, engaging, and concise step descriptions. The model is guided to avoid unnecessary elaboration while maintaining a helpful and encouraging tone. The system prompt follows this format:

```
"You are an assistant named 'Voxy' that provides step-by-step assembly instructions for the Raptor Engine. Respond in a kind, engaging, and clear manner. Begin your answers with a friendly and varied opening to make the user feel valued.  
Assembly Step Instruction (to follow exactly): {step_instruction}  
User's Question (for context only): {user_question}  
Respond with the exact instruction above for this specific assembly step, avoiding repetitive or unnecessary openings.  
Response:"
```

By implementing this structured retrieval and response approach, the AR Training Assistant ensures that users receive accurate, context-aware assembly guidance. This method significantly enhances the clarity, efficiency, and interactivity of the XR training experience, allowing users to seamlessly progress through the assembly process with precise and actionable step-by-step instructions. Additionally, the flexibility of this approach makes it possible to extend the model's capabilities to other training domains, as long as the provided manual adheres to the same structured format.

To assess the accuracy of the step-specific assembly instructions, we conducted an evaluation using six different assembly manuals. The goal of this evaluation was to determine whether the model-generated instructions precisely matched the reference instructions from the provided manuals. For this reason, we employed the BLEU score, a widely used metric for text generation tasks that measures n-gram overlap between the generated and reference text. Since the assembly instructions are predefined and must be followed word-for-word, BLEU serves as an ideal evaluation metric, ensuring that the model does not introduce unnecessary variations or modifications in its responses. Each of the six manuals contained structured step-by-step instructions, formatted in the same manner as the training manual provided to the model. The evaluation process involved:

- Providing each manual to the system as a structured text input.
- Requesting step-specific instructions at various stages of the assembly process.
- Comparing the model's output with the corresponding ground-truth step description from the manual.

The results demonstrated that the system was highly accurate in retrieving and delivering the correct instructions, achieving an average BLEU score of 94.8% across all tested manuals. This indicates that the model-generated responses were nearly identical to the reference instructions, ensuring high reliability and consistency in the assembly guidance process.

Conversation History Management

Effective conversation history management is essential for the AR Training Assistant to maintain coherent multi-turn interactions and ensure that responses remain contextually relevant throughout the training process. In this evaluation, we assess the model's ability to recall and appropriately utilize previous exchanges, ensuring that it effectively supports users through the Raptor Engine assembly workflow.

To evaluate the system's memory management capabilities, we adopt the same evaluation approach discussed in Section 3.4.1, where conversation history is restricted to the most recent seven interactions. This design ensures that the assistant can retain short-term memory of the user's instructions, clarifications, and status updates while avoiding excessive computational overhead and response degradation caused by outdated context.

The evaluation involved 10 testing sessions, each simulating a realistic training scenario where users engaged in multi-turn dialogues with the assistant. The test cases ranged from simple memory recall tasks, such as checking whether the assistant remembers the current assembly step, to more complex interactions, such as verifying the sequence of previously completed steps or asking the assistant to remind them of the last provided instruction.

The results, presented in Table 3.44, indicate that the AR Training Assistant demonstrates strong conversation history management capabilities. The model achieved an MRA score of 0.89, confirming its ability to accurately retrieve past information within the conversation window. The RC score of 0.92 suggests that the assistant correctly answered user inquiries based on its retained memory, though minor inconsistencies were observed in cases where users attempted to revisit earlier steps out of sequence. Additionally, the CC score of 0.88 highlights the model's ability to maintain logical coherence throughout extended interactions, ensuring that responses remained aligned with the overall training flow. Overall, the model achieved an average performance score of 0.90, demonstrating that the conversation history mechanism effectively supports structured and context-aware training interactions while maintaining logical consistency and user engagement.

Table 3.44: Evaluation results of conversation history management

Metric	Result
MRA	89%
RC	92%
CC	88%
Overall Performance	90%

These findings confirm that the AR Training Assistant effectively retains and utilizes relevant conversation history, enhancing user experience by maintaining continuity and accuracy in multi-turn interactions. Future improvements could focus on refining the system's ability to prioritize the most relevant past exchanges in cases where users attempt to deviate from the expected assembly sequence, ensuring even greater flexibility and adaptability in training environments.



4. Conclusions

The main objective of VOXReality WP3 is to implement natural language processing models that are context-aware, multilingual, visually grounded and knowledgeable of the domain specific needs of applications in XR environments. To this end, activities in WP3 invest in providing NLP models for automatic speech recognition, machine translation, vision language, and conversation agents.

This deliverable presents the methodology and the results obtained from the experiments of VOXReality natural language processing models that are developed in the first 30 months of the project timeline. These models include automatic speech recognition, simultaneous speech translation and machine translation, context-aware machine translation, robust machine translation, image and video spatial captioning, image spatial visual question answering, visual navigation, conversational navigation assistant and conversational training assistant tasks. This deliverable also describes the APIs developed for the inference capabilities of these models to be utilized in a plug and play approach.

The results present that the VOXReality models achieve better or competitive results in their respective tasks compared to pre-trained models that do not involve the added values of the proposed models. The experiments proved that the natural language processing models developed in the project timeline until month 15 showed promise and the further research that was conducted in the second phase of the WP3 tasks, until month 30, evolved the models to be better suited to their respective tasks. The deployment analysis that will be presented in a future deliverable (D4.2) also shows their applicability in XR applications.

The activities in WP3 resulted in a plethora of scientific and exploitable outcomes. The research results are published through 6 peer-reviewed scientific articles [Deng et al., 2023], [Pei et al., 2024a], [Maka et al., 2024], [Pei et al., 2024b], [Issam et al., 2024], [Maka et al., 2025] and additional articles are either under review or are being prepared at the time of writing this deliverable. A set of deployable APIs are published at Docker Hub²². 14 NLP models and 4 datasets are released at Hugging Face²³. Moreover, the open-source code repositories for each component are released at GitLab²⁴.

The consortium will focus on the deployment of the models in the VOXReality usecase applications in the coming months. The deployment results and optimization results for different deployment options will be presented in “D4.2 - Model Deployment analysis” on M32. Furthermore, the models and their performances will be closely monitored during the rest of the project timeline and adjustments, if needed, will be provided to the published models and their respective APIs. If such adjustments occur, they will be reported in the final project report and in the respective repositories of each component.

²² <https://hub.docker.com/u/voxreality>

²³ <https://huggingface.co/voxreality>

²⁴ <https://gitlab.com/horizon-europe-voxreality>

References

- [Amodie et al, 2016] Amodei, D., Ananthanarayanan, S., Anubhai, R., Bai, J., Battenberg, E., Case, C., ... & Zhu, Z. (2016, June). Deep speech 2: End-to-end speech recognition in english and mandarin. In International conference on machine learning (pp. 173-182). PMLR.
- [Anastasopoulos et al. 2019] Anastasopoulos, Antonios & Lui, Alison & Nguyen, Toan & Chiang, David. (2019). Neural Machine Translation of Text from Non-Native Speakers. 3070-3080. 10.18653/v1/N19-1311.
- [Ardila et al., 2019] Ardila, Rosana & Branson, Megan & Davis, Kelly & Henretty, Michael & Kohler, Michael & Meyer, Josh & Morais, Reuben & Saunders, Lindsay & Tyers, Francis & Weber, Gregor. (2019). Common Voice: A Massively-Multilingual Speech Corpus.
- [Baevski et al., 2020] Baevski, A., Zhou, Y., Mohamed, A., & Auli, M. (2020). wav2vec 2.0: A framework for self-supervised learning of speech representations. Advances in Neural Information Processing Systems, 33, 12449-12460.
- [Banerjee et al., 2005] Banerjee, S., & Lavie, A. (2005, June). METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization (pp. 65-72).
- [Bao et al., 2021] Bao, H., Dong, L., Piao, S., & Wei, F. (2021). Beit: Bert pre-training of image transformers. arXiv preprint arXiv:2106.08254.
- [Bapna et al., 2019] Bapna, Ankur & Arivazhagan, Naveen & Firat, Orhan. (2019). Simple, Scalable Adaptation for Neural Machine Translation.
- [Bar-Hillel, 1951] Bar-Hillel, Y. (1951). Fully Automatic High-Quality Translation of Indo-European Languages.
- [Barrault et al, 2023] Barrault, L., Chung, Y. A., Meglioli, M. C., Dale, D., Dong, N., Duppenhaler, M., ... & Williamson, M. (2023). Seamless: Multilingual Expressive and Streaming Speech Translation. arXiv preprint arXiv:2312.05187.
- [Bredin & Laurent, 2021] Bredin, Hervé & Laurent, Antoine. (2021). End-To-End Speaker Segmentation for Overlap-Aware Resegmentation. 3111-3115. 10.21437/Interspeech.2021-560.
- [Brown et al., 2020] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. Advances in neural information processing systems, 33, 1877-1901.
- [Brunner et al., 2020] Gino Brunner, Yang Liu, Damian Pascual, Oliver Richter, Massimiliano Ciaramita, and Roger Wattenhofer. 2020. On identifiability in transformers. In 8th International Conference on Learning Representations (ICLR 2020)(virtual). International Conference on Learning Representations.



[Bryant et al., 2019] Bryant, Christopher & Felice, Mariano & Andersen, Øistein & Briscoe, Ted. (2019). The BEA-2019 Shared Task on Grammatical Error Correction. 52-75. 10.18653/v1/W19-4406.

[Budzianowski et al., 2018] Budzianowski, P., Wen, T. H., Tseng, B. H., Casanueva, I., Ultes, S., Ramadan, O., & Gašić, M. (2018). Multiwoz--a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. arXiv preprint arXiv:1810.00278.

[Cettolo et al., 2015] Cettolo, M., Niehues, J., Stüker, S., Bentivogli, L., Cattoni, R., & Federico, M. (2015). The IWSLT 2015 evaluation campaign. In Proceedings of the 12th International Conference on Spoken Language Translation.

[Cettolo et al., 2017] Cettolo, M., Federico, M., Bentivogli, L., Niehues, J., Stüker, S., Sudoh, K., ... & Federmann, C. (2017). Overview of the IWSLT 2017 Evaluation Campaign. In Proceedings of the 14th International Conference on Spoken Language Translation (pp. 2-14).

[Chane et al., 2020] Chen, Y. C., Li, L., Yu, L., El Kholy, A., Ahmed, F., Gan, Z., ... & Liu, J. (2020, August). Uniter: Universal image-text representation learning. In European conference on computer vision (pp. 104-120). Cham: Springer International Publishing.

[Conneau et al., 2020] Conneau, A., Baevski, A., Collobert, R., Mohamed, A., & Auli, M. (2020). Unsupervised cross-lingual representation learning for speech recognition. arXiv preprint arXiv:2006.13979.

[Costa-jussà et al., 2022] Costa-jussà, M. R., Cross, J., Çelebi, O., Elbayad, M., Heafield, K., Heffernan, K., ... & NLLB Team. (2022). No language left behind: Scaling human-centered machine translation. arXiv preprint arXiv:2207.04672.

[Deng et al., 2023] Deng, W., Pei, J., Ren, Z., Chen, Z., & Ren, P. (2023). Intent-calibrated self-training for answer selection in open-domain dialogues. Transactions of the Association for Computational Linguistics, 11, 1232-1249.

[Dettmers et al., 2022] Dettmers, T., Lewis, M., Belkada, Y., & Zettlemoyer, L. (2022). Llm.int8(): 8-bit matrix multiplication for transformers at scale. arXiv preprint arXiv:2208.07339.

[Dettmers et al., 2023] Dettmers, T., Pagnoni, A., Holtzman, A., & Zettlemoyer, L. (2023). Qlora: Efficient finetuning of quantized llms. arXiv preprint arXiv:2305.14314.

[Devlin et al., 2018] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

[Dosovitskiy et al., 2020] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

[Elbayad et al., 2020] Elbayad, M., Besacier, L., & Verbeek, J. (2020). Efficient wait-k models for simultaneous machine translation. arXiv preprint arXiv:2005.08595.

[Fan et al., 2021] Fan, A., Bhosale, S., Schwenk, H., Ma, Z., El-Kishky, A., Goyal, S., ... & Joulin, A. (2021). Beyond english-centric multilingual machine translation. The Journal of Machine Learning Research, 22(1), 4839-4886.

- [Fang et al., 2022] Fang, Qingkai & Ye, Rong & Li, Lei & Feng, Yang & Wang, Mingxuan. (2022). STEMM: Self-learning with Speech-text Manifold Mixup for Speech Translation. 7050-7062. 10.18653/v1/2022.acl-long.486.
- [Gulati et al., 2020] Gulati, Anmol & Qin, James & Chiu, Chung-Cheng & Parmar, Niki & Zhang, Yu & Yu, Jiahui & Han, Wei & Wang, Shibo & Zhang, Zhengdong & Wu, Yonghui & Pang, Ruoming. (2020). Conformer: Convolution-augmented Transformer for Speech Recognition. 5036-5040. 10.21437/Interspeech.2020-3015.
- [Hannun et al, 2014] Hannun, A., Case, C., Casper, J., Catanzaro, B., Diamos, G., Elsen, E., ... & Ng, A. Y. (2014). Deep speech: Scaling up end-to-end speech recognition. arXiv preprint arXiv:1412.5567.
- [Heffernan et al., 2022] Heffernan, K., Çelebi, O., & Schwenk, H. (2022). Bitext mining using distilled sentence representations for low-resource languages. arXiv preprint arXiv:2205.12654.
- [Honnibal et al., 2020] Honnibal, M., Montani, I., Van Landeghem, S., & Boyd, A. (2020). spaCy: Industrial-strength Natural Language Processing in Python. <https://doi.org/10.5281/zenodo.1212303>
- [Houlsby et al., 2019] Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M. & Gelly, S.. (2019). Parameter-Efficient Transfer Learning for NLP. Proceedings of the 36th International Conference on Machine Learning, in Proceedings of Machine Learning Research 97:2790-2799
- [Hsu et al., 2021] Hsu, W. N., Bolte, B., Tsai, Y. H. H., Lakhotia, K., Salakhutdinov, R., & Mohamed, A. (2021). Hubert: Self-supervised speech representation learning by masked prediction of hidden units. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 29, 3451-3460.
- [Hu et al., 2021] Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., ... & Chen, W. (2021). Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685.
- [Huo et al., 2020] Huo, J., Herold, C., Gao, Y., Dahlmann, L., Khadivi, S., & Ney, H. (2020). Diving deep into context-aware neural machine translation. arXiv preprint arXiv:2010.09482.
- [Iranzo-Sánchez et al., 2019] Iranzo-Sánchez, J., Silvestre-Cerdà, J.A., Jorge, J., Roselló, N., Giménez, A., Sanchís, A., Civera Saiz, J., & Juan-Císcar, A. (2019). Europarl-ST: A Multilingual Corpus for Speech Translation of Parliamentary Debates. ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 8229-8233.
- [Issam et al., 2024] Issam, A., Semerci, Y. C., Scholtes, J., & Spanakis, G. (2024, August). Fixed and Adaptive Simultaneous Machine Translation Strategies Using Adapters. In Proceedings of the 21st International Conference on Spoken Language Translation (IWSLT 2024) (pp. 298-310).
- [Jelinek, 1976] Jelinek, F. (1976). Continuous speech recognition by statistical methods. Proceedings of the IEEE, 64(4), 532-556.



[Jon et al., 2021] Jon, J., Aires, J. P., Varis, D., & Bojar, O. (2021, August). End-to-End Lexically Constrained Machine Translation for Morphologically Rich Languages. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers) (pp. 4019-4033)

[Kochen, 1980] Kochen, M. (1981). On the generality of PARRY, Colby's paranoia model. Behavioral and Brain Sciences, 4(4), 540-541. <https://doi.org/10.1017/s0140525x00000091>

[Koehn et al., 2007] Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., ... & Herbst, E. (2007, June). Moses: Open source toolkit for statistical machine translation. In Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions (pp. 177-180). Association for Computational Linguistics.

[Kolve et al., 2017] Kolve, E., Mottaghi, R., Han, W., VanderBilt, E., Weihs, L., Herrasti, A., ... & Farhadi, A. (2017). Ai2-thor: An interactive 3d environment for visual ai. arXiv preprint arXiv:1712.05474.

[Kornblith et al, 2019] Kornblith, S., Norouzi, M., Lee, H., & Hinton, G. (2019, May). Similarity of neural network representations revisited. In International conference on machine learning (pp. 3519-3529). PMLR.

[Kritharoula et al., 2023] Kritharoula, A., Lymperaiou, M., & Stamou, G. (2023). Language Models as Knowledge Bases for Visual Word Sense Disambiguation. arXiv preprint arXiv:2310.01960.

[Lewis et al., 2019] Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., ... & Zettlemoyer, L. (2019). Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. arXiv preprint arXiv:1910.13461.

[Li et al., 2019] Li, L. H., Yatskar, M., Yin, D., Hsieh, C. J., & Chang, K. W. (2019). Visualbert: A simple and performant baseline for vision and language. arXiv preprint arXiv:1908.03557.

[Li et al., 2021] Li, M., Lv, T., Cui, L., Lu, Y., Florencio, D., Zhang, C., ... & Wei, F. (2021). Trocr: transformer-based optical character recognition with pre-trained models. arXiv 2021. arXiv preprint arXiv:2109.10282.

[Li et al., 2022] Li, J., Li, D., Xiong, C., & Hoi, S. (2022, June). Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In International conference on machine learning (pp. 12888-12900). PMLR.

[Lillicrap et al., 2017] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., ... & Wierstra, D. (2015). Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971.

[Lin et al., 2004] Lin, C. Y. (2004, July). Rouge: A package for automatic evaluation of summaries. In Text summarization branches out (pp. 74-81).

[Lin et al., 2021] Lin, K., Li, L., Lin, C. C., Ahmed, F., Gan, Z., Liu, Z., ... & Wang, L. (2022). Swinbert: End-to-end transformers with sparse attention for video captioning. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 17949-17958).

- [Lison et al. 2018] Lison, P., Tiedemann, J., & Kouylekov, M. (2018, May). OpenSubtitles2018: Statistical Rescoring of Sentence Alignments in Large, Noisy Parallel Corpora. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018).
- [Liu et al., 2019] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- [Liu et al., 2021] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the IEEE/CVF international conference on computer vision (pp. 10012-10022).
- [Liu et al., 2022] Liu, F., Emerson, G., & Collier, N. (2023). Visual spatial reasoning. Transactions of the Association for Computational Linguistics, 11, 635-651.
- [Lopes et al., 2020] Lopes, A., Farajian, M. A., Bawden, R., Zhang, M., & Martins, A. F. (2020, November). Document-level Neural MT: A Systematic Comparison. In Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (pp. 225-234).
- [Loshchilov & Hutter, 2017] Loshchilov, I., & Hutter, F. (2017). Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101.
- [Ma et al., 2018] Ma, M., Huang, L., Xiong, H., Zheng, R., Liu, K., Zheng, B., ... & Wang, H. (2018). STACL: Simultaneous translation with implicit anticipation and controllable latency using prefix-to-prefix framework. arXiv preprint arXiv:1810.08398.
- [Ma et al., 2019] Ma, X., Pino, J., Cross, J., Puzon, L., & Gu, J. (2019). Monotonic multihead attention. arXiv preprint arXiv:1909.12406.
- [Ma et al., 2023] Ma, X., Sun, A., Ouyang, S., Inaguma, H., & Tomasello, P. (2023). Efficient Monotonic Multihead Attention. ArXiv, abs/2312.04515.
- [Maka et al., 2024] Maka, P., Semerci, Y., Scholtes, J., & Spanakis, G. (2024, March). Sequence Shortening for Context-Aware Machine Translation. In Findings of the Association for Computational Linguistics: EACL 2024 (pp. 1874-1894).
- [Maka et al., 2025] Paweł Mała, Yusuf Can Semerci, Jan Scholtes, and Gerasimos Spanakis. 2025. Analyzing the Attention Heads for Pronoun Disambiguation in Context-aware Machine Translation Models. In Proceedings of the 31st International Conference on Computational Linguistics, pages 6348–6377, Abu Dhabi, UAE. Association for Computational Linguistics.
- [Miyato et al., 2016] Miyato, T., Dai, A. M., & Goodfellow, I. (2016). Adversarial training methods for semi-supervised text classification. arXiv preprint arXiv:1605.07725.
- [Müller et al., 2018] Müller, M., Gonzales, A. R., Voita, E., & Sennrich, R. (2018, October). A Large-Scale Test Set for the Evaluation of Context-Aware Pronoun Translation in Neural Machine Translation. In Proceedings of the Third Conference on Machine Translation: Research Papers (pp. 61-72).



[Napoles et al., 2016] Napoles, Courtney & Cahill, Aoife & Madnani, Nitin. (2016). The Effect of Multiple Grammatical Errors on Processing Non-Native Writing. 1-11. 10.18653/v1/W16-0501.

[Napoles et al., 2017] Napoles, Courtney & Sakaguchi, Keisuke & Tetreault, Joel. (2017). JFLEG: A Fluency Corpus and Benchmark for Grammatical Error Correction.

[Nawrot et al., 2021] Nawrot, P., Tworkowski, S., Tyrolski, M., Kaiser, Ł., Wu, Y., Szegedy, C., & Michalewski, H. (2021). Hierarchical transformers are more efficient language models. arXiv preprint arXiv:2110.13711.

[Olsson et al., 2022] Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A. and Conerly, T., 2022. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*.

[Padmakumar et al., 2022] Padmakumar, A., Thomason, J., Shrivastava, A., Lange, P., Narayan-Chen, A., Gella, S., ... & Hakkani-Tur, D. (2022, June). Teach: Task-driven embodied agents that chat. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 36, No. 2, pp. 2017-2025).

[Papineni et al., 2002] Papineni, K., Roukos, S., Ward, T., & Zhu, W. J. (2002, July). Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics (pp. 311-318).

[Pei et al., 2024a] Pei, J., Yan, G., De Rijke, M., & Ren, P. (2024). Mixture-of-Languages Routing for Multilingual Dialogues. *ACM Transactions on Information Systems*, 42(6), 1-33.

[Pei et al., 2024b] Pei, J., Viola, I., Huang, H., Wang, J., Ahsan, M., Ye, F., ... & Cesar, P. (2024, August). Autonomous Workflow for Multimodal Fine-Grained Training Assistants Towards Mixed Reality. In Findings of the Association for Computational Linguistics ACL 2024 (pp. 4051-4066).

[Peng et al., 2022] Peng, B., Galley, M., He, P., Brockett, C., Liden, L., Nouri, E., ... & Gao, J. (2022). Godel: Large-scale pre-training for goal-directed dialog. arXiv preprint arXiv:2206.11309.

[Popović, 2015] Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In Proceedings of the Tenth Workshop on Statistical Machine Translation, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

[Pratap et al., 2023] Pratap, Vineel & Tjandra, Andros & Shi, Bowen & Tomasello, Paden & Arun Babu, Thirunavukkarasu & Kundu, Sayani & Elkahky, Ali & Ni, Zhaoheng & Vyas, Apoorv & Fazel-Zarandi, Maryam & Baevski, Alexei & Adi, Yossi & Zhang, Xiaohui & Hsu, Wei-Ning & Conneau, Alexis & Auli, Michael. (2023). Scaling Speech Technology to 1,000+ Languages.

[Rabiner, 1989] Rabiner, L. R. (1989). A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE*, 77(2), 257-286.

[Radford et al., 2018] Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.

[Radford et al., 2019] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, 1(8), 9.

- [Radford et al., 2023] Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2023, July). Robust speech recognition via large-scale weak supervision. In International Conference on Machine Learning (pp. 28492-28518). PMLR.
- [Raffel et al., 2020] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140), 1-67.
- [Ranftl et al., 2021] Ranftl, R., Bochkovskiy, A., & Koltun, V. (2021). Vision transformers for dense prediction. In Proceedings of the IEEE/CVF international conference on computer vision (pp. 12179-12188).
- [Rebuffi et al., 2017] Rebuffi, Sylvestre-Alvise & Bilen, Hakan & Vedaldi, Andrea. (2017). Learning multiple visual domains with residual adapters.
- [Reimers & Gurevych, 2019] Reimers, N., & Gurevych, I. (2019). Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- [Schulman et al., 2017] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- [Souček et al., 2008] Souček, T., & Lokoč, J. (2008). TransNet V2: An effective deep network architecture for fast shot transition detection." *arXiv*, Aug. 11, 2020. doi: 10.48550. *arXiv*.
- [Subramanian et al., 2020] Subramanian, S., Collobert, R., Ranzato, M. A., & Boureau, Y. L. (2020). Multi-scale Transformer Language Models. *arXiv preprint arXiv:2005.00581*.
- [Sun et al., 2020] Sun, Z., Wang, M., Zhou, H., Zhao, C., Huang, S., Chen, J., & Li, L. (2020). Rethinking document-level neural machine translation. *arXiv preprint arXiv:2010.08961*.
- [Sutskever et al., 2014] Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.
- [Tajiri et al., 2012] Tajiri, Toshikazu & Komachi, Mamoru & Matsumoto, Yuji. (2012). Tense and Aspect Error Correction for ESL Learners Using Global Context. 50th Annual Meeting of the Association for Computational Linguistics, ACL 2012 - Proceedings of the Conference. 2.
- [Tan & Bansal, 2019] Tan, H., & Bansal, M. (2019). Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*.
- [Tiedemann & Thottingal, 2020] Tiedemann, J., & Thottingal, S. (2020). OPUS-MT – Building open translation services for the World. European Association for Machine Translation Conferences/Workshops.
- [Touvron et al., 2020] Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., & Jégou, H. (2020). Training data-efficient image transformers & distillation through attention. *arXiv 2020. arXiv preprint arXiv:2012.12877*.



[Touvron et al., 2023] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., ... & Scialom, T. (2023). Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.

[Vaswani et al., 2017] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

[Vinyals et al., 2015] Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 3156-3164).

[Voita et al., 2019] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019. Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 5797–5808, Florence, Italy. Association for Computational Linguistics.

[Wang et al., 2020] Wang, Changhan & Wu, Anne & Pino, Juan. (2020). CoVoST 2: A Massively Multilingual Speech-to-Text Translation Corpus. 10.48550/arXiv.2007.10310.

[Wang et al., 2023] Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., ... & Wen, J. R. (2023). A survey on large language model based autonomous agents. arXiv preprint arXiv:2308.11432.

[Wicks et al., 2023] Rachel Wicks and Matt Post. 2023. Identifying Context-Dependent Translations for Evaluation Set Production. In Proceedings of the Eighth Conference on Machine Translation, pages 452–467, Singapore. Association for Computational Linguistics.

[Xie et al., 2020] Xie, Q., Dai, Z., Hovy, E., Luong, T., & Le, Q. (2020). Unsupervised data augmentation for consistency training. Advances in neural information processing systems, 33, 6256-6268.

[Xie et al., 2021] Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., & Luo, P. (2021). SegFormer: Simple and efficient design for semantic segmentation with transformers. Advances in neural information processing systems, 34, 12077-12090.

[Xu et al., 2016] Xu, J., Mei, T., Yao, T., & Rui, Y. (2016). Msr-vtt: A large video description dataset for bridging video and language. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 5288-5296).

[Yannakoudakis et al., 2011] Yannakoudakis, Helen & Briscoe, Ted & Medlock, Ben. (2011). A New Dataset and Method for Automatically Grading ESOL Texts. 180-189.

[Yao et al., 2023] Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2022). React: Synergizing reasoning and acting in language models. ICLR 2023.

[Zhang & Feng, 2021] Zhang, S., & Feng, Y. (2021). Universal simultaneous machine translation with mixture-of-experts wait-k policy. arXiv preprint arXiv:2109.05238.

[Zheng et al., 2020] Zheng, B., Liu, K., Zheng, R., Ma, M., Liu, H., & Huang, L. (2020). Simultaneous translation policies: From fixed to adaptive. arXiv preprint arXiv:2004.13169.



[Zhu et al., 2022] Zhu, Q., Geishauser, C., Lin, H. C., van Niekerk, C., Peng, B., Zhang, Z., ... & Huang, M. (2022). ConvLab-3: A Flexible Dialogue System Toolkit Based on a Unified Data Format. arXiv preprint arXiv:2211.17148.

[Zhuang et al., 2019] Zhuang, F., Qi, Z., Duan, K., Xi, D., Zhu, Y., Zhu, H., ... & He, Q. (2019). A comprehensive survey on transfer learning. arxiv e-prints. arXiv preprint arXiv:1911.02685.

[Zhong et al., 2024] Zhong, W., Guo, L., Gao, Q., Ye, H., & Wang, Y. (2024, March). Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 38, No. 17, pp. 19724-19731).



Appendix I: Endpoints for the ASR and NMT components

1) Transcribe Audio: The method to transcribe given audio files. The source language is not required.

URL: <api_url>/transcribe_audio_files

Method: POST

Request parameters:

- audio_files [array]: the audio file to transcribe

Response:

- **Success:** Status Code: 200 Content: {"transcriptions": [String]}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

2) Translate Audio: The method to translate given audio files to the target language. The source language is not required.

URL: <api_url>/translate_audio_files

Method: POST

Request parameters:

- audio_files [array]: the audio file to translate
- target_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish
- return_transcription [bool]: true/false

Response:

- **Success:** Status Code: 200 Content: {"transcriptions": [String], "translations": [String]}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

3) Translate Audio with Context: The method to translate given audio files to the target language using additional textual context for better translation of the transcribed text. The source language is not required.

URL: <api_url>/contextual_translate_audio_files

Method: POST

Request parameters:

- audio_files [array]: the audio file to translate
- context [String]: the textual context to use in the translation
- target_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish
- return_transcription [bool]: true/false

Response:

- **Success:** Status Code: 200 Content: {"transcriptions": [String], "translations": [String]}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

4) Translate Text: The method to translate given text to the target language. The source language is not required.

URL: <api_url>/translate_text

Method: POST

Request parameters:



- text [String]: the text to translate
- target_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish

Response:

- **Success:** Status Code: 200 Content: {"translations": [String]}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

5) Translate Text with Context: The method to translate given text to the target language using additional textual context for better translation. The source language is not required.

URL: <api_url>/contextual_translate_text

Method: POST

Request parameters:

- text [String]: the text to translate
- context [String]: the textual context to use in the translation
- target_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish

Response:

- **Success:** Status Code: 200 Content: {"translations": [String]}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

6) Translate Text with Context and Terminology: The method to translate given text to the target language using additional textual context and terminological term pairs for better translation. The source language is not required but the terminology file is required.

URL: <api_url>/contextual_terminology_translate_text

Method: POST

Request parameters:

- text [String]: the text to translate
- context [String]: the textual context to use in the translation
- target_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish

Response:

- **Success:** Status Code: 200 Content: {"translations": [String]}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

7) Upload Terminology: The method to upload a file that contains the terminological pairs.

URL: <api_url>/upload_terminology

Method: POST

Request parameters:

- file [String]: the file that contains terminological term pairs for different languages

Response:

- **Success:** Status Code: 200 Content: {"Successfully uploaded terminology file: <file_name>"}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

8) Upload Subtitles: The method to upload the subtitles file. The first line in the file needs to be the language short codes.

URL: <api_url>/upload_terminology

Method: POST

Request parameters:

- file [String]: the file that contains subtitles for the selected languages; the lines in different languages should correspond to (be a translation of) each other

Response:

- **Success:** Status Code: 200 Content: {"Successfully uploaded subtitles file: <file_name>"}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

9) Transcribe Audio with Subtitles: The method to transcribe given audio files and returns the most closely matching subtitle line from the uploaded subtitles file. If target language is specified, the corresponding line in the target language is also returned.

URL: <api_url>/subtitle_transcribe_audio_files

Method: POST

Request parameters:

- audio_files [array]: the audio file to transcribe
- source_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish
- target_language [String]: Optional "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish

Response:

- **Success:** Status Code: 200 Content: {"transcriptions": [String]}
- **Error:** Status Code: 466 Content: {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

10) Translate Audio with Subtitles: The method to translate given audio files to the target language. Before translation, the transcription is matched with the subtitles and the most closely matching subtitle line from the uploaded subtitles file is used for translation. If the target language is present in the subtitles file the corresponding subtitle line in the target language is returned instead of the translation.

URL: <api_url>/subtitle_translate_audio_files

Method: POST

Request parameters:

- audio_files [array]: the audio file to translate
- source_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish
- target_language [String]: Optional "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish
- return_transcription [bool]: true/false

Response:

- **Success:** Status Code: 200 Content: {"transcriptions": [String], "translations": [String]}

- **Error:** *Status Code:* 466 *Content:* {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

11) Translate Audio with Context and Subtitles: The method to translate given audio files to the target language using additional textual context for better translation of the transcribed text. Before translation, the transcription is matched with the subtitles and the most closely matching subtitle line from the uploaded subtitles file is used for translation. If the target language is present in the subtitles file the corresponding subtitle line in the target language is returned instead of the translation.

URL: <api_url>/subtitle_contextual_translate_audio_files

Method: POST

Request parameters:

- audio_files [array]: the audio file to translate
- context [String]: the textual context to use in the translation
- source_language [String]: "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish
- target_language [String]: Optional "en"-English, "el"-Greek, "it"-Italian, "nl"-Dutch, "de"-German, "es"-Spanish
- return_transcription [bool]: true/false

Response:

- **Success:** *Status Code:* 200 *Content:* {"transcriptions": [String], "translations": [String]}
- **Error:** *Status Code:* 466 *Content:* {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

12) Transcribe Audio Stream: The method to transcribe audio in streaming mode.

URL: <api_url>/listen

Method: POST

Request parameters:

- audio_files [array]: the audio stream.
- model_name [String]: Name of the model to use for transcription.
- language [String]: The source language. if it is not specified then it is detected automatically.
- chunk_length_seconds [float]: The number of seconds of audio in each chunk.
- accumulation_chunk_length_seconds [float]: The chunk length to add to accumulate audio until speech ends.
- max_accumulation_seconds [float]: The maximum audio length for accumulation.

Response:

- **Success:** *Status Code:* 200 *Content:* {"transcriptions": [String]}
- **Error:** *Status Code:* 466 *Content:* {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

13) Transcribe/Translate E2E Audio Stream: The method to transcribe or translate audio in streaming mode using SeamlessStreaming.

URL: <api_url>/e2e_listen

Method: POST

Request parameters:

- `audio_files` [array]: the audio stream.
- `model_name` [String]: Name of the model to use for transcription.
- `language` [String]: The target language for translation or transcription. If it is the source language then the output is the transcription. The source language is detected automatically.

Response:

- **Success:** *Status Code:* 200 *Content:* {"translation": [String]}
- **Error:** *Status Code:* 466 *Content:* {"message": "An error occurred in <error_location>. Details: <error_type>-<error_details>"}

Appendix II: Endpoints for the VL components

1) Image Generic Captioning: The endpoint accepts as input an image file (png, jpg, jpeg) and returns as text a generic description of the image provided.

URL: <api_url>/cap_gpt2

Method: POST

Request parameters:

- file [String]: the image to be captioned.

Response:

- **Success:** Status Code: 200 Content: {"Image caption": [String]}

2) Image Visual Generic One-Word Question Answering: The endpoint accepts as input an image file (png, jpg, jpeg) and a question as a text and returns as text a single-word answer to the provided question.

URL: <api_url>/lxmert

Method: POST

Request parameters:

- file [String]: the image to be questioned about.
- questions [String]: the question on the image provided.

Response:

- **Success:** Status Code: 200 Content: {"Prediction by LXMERT": ["<provided answer>"] }

3) Video Generic Captioning: The endpoint accepts as input a video file (mp4, avi) and returns as text a generic description of the video provided.

URL: <api_url>/swinbert_videocap

Method: POST

Request parameters:

- file [String]: the video to be captioned.

Response:

- **Success:** Status Code: 200 Content: {"Video caption": "<video description>"}

4) Image Visual Generic Question Answering: The endpoint accepts as input an image file (png, jpg, jpeg) and a question as a text and returns as text an answer (whole sentence) to the provided question.

URL: <api_url>/rgb-language_blip

Method: POST

Request parameters:

- file [String]: the image to be questioned about.

- question [String]: the question on the image provided.

Response:

- **Success:** Status Code: 200 Content: {"Prediction": ["<provided answer>"]}

5) Image Scenegrph Generation: The endpoint accepts as input an image file (png, jpg, jpeg) and returns as text the complete spatial description of the scene, as well as all the metadata created in the process, which drove the generation of the returned description.

URL: <api_url>/scenegrph_rgb

Method: POST

Request parameters:

- rgb_file [String]: the image to be spatially described.

Response:

- **Success:** Status Code: 200 Content: {"Inference_time_sec": [float],
"Caps": [String],
"Meta": [{"class_names": [String],
"object_uids": [String],
"spatial_relationships": [String]
},
"Objects": [{"Class id": [String],
"bbox_xmin": [String],
"bbox_ymin": [String],
"bbox_xmax": [String]
}]
}

6) Image Spatial Captioning: The endpoint accepts as input an image file (png, jpg, jpeg) and returns as text a spatial description of the image provided.

URL: <api_url>/rgb-language_cap

Method: POST

Request parameters:

- file [String]): the image to be spatially captioned.

Response:

- **Success:** Status Code: 200 Content: {"Image caption": [String], "Inference_time_sec": [float]}

7) Image Visual Spatial Question Answering: The endpoint accepts as input an image file (png, jpg, jpeg) and a question related with spatiality as a text and returns as text an answer (whole sentence) to the provided question.

URL: <api_url>/rgb-language_vqa

Method: POST

Request parameters:

- file [String]: the image to be questioned about.
- question [String]: the spatiality-related question on the image provided.

Response:

- **Success:** Status Code: 200 Content: {"Prediction": [String]}

8) Video Spatial Captioning: The endpoint accepts as input a video file (mp4, avi) and returns as text a spatial description of the video provided.

URL: <api_url>/video-language_vqa

Method: POST

Request parameters:

- file [String]: the video to be captioned.

Response (success):

- **Success:** Status Code: 200 Content: {"Video caption": [String]}

9) VR Conferences Visual Navigation: The endpoint accepts as input an image captured from the viewpoint of the VR Conferences application user and a destination-room name as a text and returns the complete set of navigation instructions (turn-by-turn steps) to get there.

URL: <api_url>/rgb-language_vqa

Method: POST

Request parameters:

- file [String]: the image/view from the user's current position.
- question [String]: the destination-room name in ["conference room", "business room", "social area", "booth 1", "booth 2", "booth 3", "booth 4"].

Response:

- **Success:** Status Code: 200 Content: {"Prediction": [String]}

Appendix III: Endpoints for the navigation assistant

1) Intent and Destination Detection: This endpoint processes a user query to identify the intent and the destination (topic or area of interest) within the context of the conference.

URL: <api_url> /intent_dest/

Method: POST

Request parameters:

- user_query [String]: The user's query or question.

Response:

- **Success:** Status Code: 200 Content: {"intent": [String], "destination": [String]}
- **Error:** Status Code: 400 Content: {"message": "Invalid request: 'user_query' is required."}

2) Generate Response: This endpoint generates a response for the user's query using the provided user query and additional knowledge input (Mozilla input).

URL: <api_url> /response/

Method: POST

Request parameters:

- user_query [String]: The user's query or question.
- mozilla_input [String]: Additional knowledge or context to assist in generating the response.

Response:

- **Success:** Status Code: 200 Content: {"response": [String]}
- **Error:** Status Code: 400 Content: {"message": "Invalid request format. Please check user_query, intent and mozilla_input fields."}

Appendix IV: Endpoints for the training assistant

1) Textual Chat Agent: This endpoint generates a response for the user's query using the provided user query and simulated XR functional tools using textual requests.

URL: <api_url> /ask/

Method: POST

Request parameters:

- query [String]: The user's query or question.

Response:

- **Success:** Status Code: 200 Content: {"response": [String]}
- **Error:** Status Code: 400 Content: {"message": "Invalid request format."}

2) Unity Chat Agent: This endpoint generates a response for the user's query using the provided user query and simulated XR functional tools using requests coming from Unity app

URL: <api_url> /chatbot/

Method: WebSocket

Request parameters:

- websocket [WebSocket]: The websocket from Unity with user's query or question as messages.

Response:

- **Success:** Status Code: 200 Content: {"response": [String]}
- **Error:** Status Code: 400 Content: {"message": "Invalid request format."}

3) API Key Validation: This endpoint runs before every request and checks if the request includes the correct "Authorization" header. Returns an unauthorized error if the API key is invalid.

Response:

- **Error:** Status Code:403

4) Ask Question Endpoint: This endpoint generates a response using the provided user query and XR test application functional tools using requests coming from Unity.

URL: <api_url> /ask

Method: (POST)

Request parameters:

- question [String]: The user's query.
- current_step [int]: Current step flag
- wrong_object_flag [bool]: Wrong object flag.

Response:

- **Success:** Status Code: 200 Content: {"response": [String]}
- **Error:** Status Code: 400 Content: {"message": "Invalid request format. Please check user_query, current step or wrong object flags."}



VOXReality

Voice driven
interaction in XR spaces



**Funded by
the European Union**

Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or Directorate-General for Communications Networks, Content and Technology (DG CNECT). Neither the European Union nor the granting authority can be held responsible for them.