

From TOWER to SPIRE: Adding the Speech Modality to a Text-Only LLM

Anonymous ACL submission

Abstract

We introduce SPIRE, a speech-augmented language model (LM) capable of both translating and transcribing speech input from English into 10 other languages as well as translating text input in both language directions. SPIRE integrates the speech modality into an existing multilingual LM (MLM) via speech discretization and continued pre-training using only 42.5K hours of speech. In particular, we adopt the pretraining framework of MLMs and treat discretized speech input as an additional *translation language*. This approach not only equips the MLM with speech capabilities, but also preserves its strong text-only performance. We achieve this using significantly less data than existing speech LMs, demonstrating that discretized speech input integration as an additional language is feasible during LM adaptation. We will make our code and models available to the community.

1 Introduction

Large language models (LLMs) have demonstrated remarkable success on various text-based natural language processing tasks (Achiam et al., 2023; Touvron et al., 2023; Yang et al., 2024; Alves et al., 2024; Martins et al., 2024), motivating research into extending them to other modalities. This has led to the development of multimodal LMs (MLMs) capable of processing speech, audio, images, and video (Team et al., 2023; Driess et al., 2023; Rubenstein et al., 2023; Liu et al., 2023; Tang et al., 2024; Défossez et al., 2024; Hu et al., 2024; Huang et al., 2024; Nguyen et al., 2025). However, the integration of new modalities often come at the cost of existing capabilities (Zhai et al., 2024).

For speech-LLM integration, a simple approach is to link the output of an automatic speech recognition (ASR) system to a text-only LLM (Huang et al., 2024). This solution, however, is prone to error propagation and depends largely on individual model quality. More popular are solutions that

investigate equipping LLMs natively with speech processing capabilities through modality projection (Shu et al., 2023; Radhakrishnan et al., 2023; Wu et al., 2023a; Tang et al., 2024; Xue et al., 2024; Hu et al., 2024). Typically, a speech foundation model generates speech representations that are mapped to the embedding space of the LLM, following which the model is then fine-tuned along with a projector on speech-to-text tasks to equip the LLM with speech processing capabilities. In this setting, key challenges include prompt overfitting and high training costs, as tuning these MLMs requires the adaptation of the speech projector module on vast amounts of raw speech data (Tang et al., 2024; Hu et al., 2024).

An alternative approach for MLMs is to use *speech discretization*, where continuous speech features are transformed prior to training into sequences of “discrete speech units” (DSUs), which can be processed similarly to text (Chou et al., 2023a; Zhang et al., 2023; Rubenstein et al., 2023; Chang et al., 2024; Défossez et al., 2024; Trinh et al., 2024; Maiti et al., 2024; Nguyen et al., 2025). This approach simplifies training by eliminating the need for additional parameters beyond extended embedding matrices. Finally, while both projector-based and discretization-based MLMs have shown promising results on text-to-speech and speech-to-text tasks, their development has prioritized speech-centric tasks at the expense of textual performance. Furthermore, limited research has focused on integrating speech while preserving the LLM’s original capabilities in textual tasks (Chou et al., 2023b; Huang et al., 2024).

In this work we present SPIRE, a speech-augmented LLM built from the open-weight multilingual model TOWER (Alves et al., 2024). SPIRE can perform English ASR and from-English speech translation (ST) while maintaining TOWER’s strong performance on machine translation (MT) across

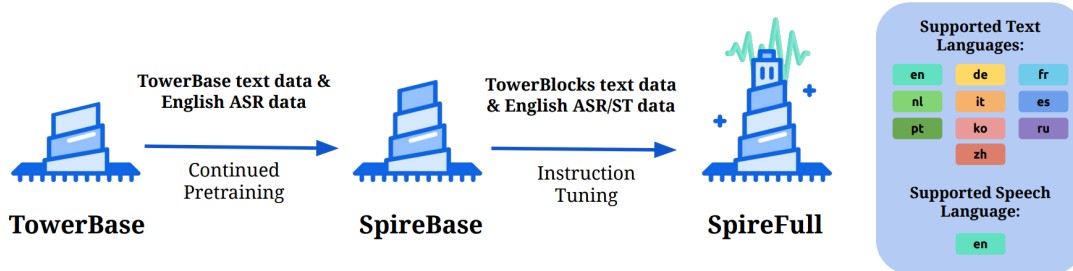


Figure 1: Illustration of the model training method for SPIREBASE and SPIREFULL.

all 10 languages¹ supported by TOWER. SPIRE encodes speech via HuBERT-based (Hsu et al., 2021) k-means clustering, as in previous work (Zhang et al., 2023; Rubenstein et al., 2023; Chang et al., 2024). We perform training in two stages: Continued Pre-Training (CPT) and Instruction Tuning (IT). For the CPT stage, we use a mixture of ASR data and a small fraction of TOWER’s text CPT data. For IT, we leverage TOWER’s task-specific MT data, as well as additional English ASR and ST data. SPIRE is trained using only **42.5K** hours of speech, differing from the large scale of data used by existing models (Radford et al., 2023; Nguyen et al., 2025; Chu et al., 2024). Figure 1 illustrates our training process. We make the following contributions:

- We present a pipeline for integrating speech as an additional modality into an existing LLM, enabling it to transcribe and translate English speech while preserving its original text-only capabilities across 10 languages;
- We analyze speech integration at two stages, namely CPT and IT, demonstrating the necessity of both stages to achieve optimal performance across both modalities;
- We make our models, datasets, and scripts available to the community.²

2 Related Work

Speech-to-Text Models An increasing number of studies have explored integrating speech into LLMs (Zhang et al., 2023; Rubenstein et al., 2023; Hassid et al., 2024). For discrete speech input, Hassid et al. (2024) demonstrate the benefits of initializing a speech LLM from a text-based LLM. SpeechGPT (Zhang et al., 2023) applies

IT on speech-to-text cross-modal ASR, text-to-speech (TTS), and text-based question answering. AudioPALM (Rubenstein et al., 2023) is trained in a multi-task fashion, similarly to SpeechGPT, but on multilingual input. Recently, VoxLM (Maiti et al., 2024) was trained jointly on DSUs and text data for ASR, TTS, and open-ended speech/text generation. Our work is most similar to Spirit-LM (Nguyen et al., 2025), which adapts an LLM with an interleaved mixture of DSU and text data, which requires an expensive DSU-to-transcript step to create. In contrast, we adopt a more cost-effective input representation that can be extended to any language, regardless of the availability of a speech aligner. Our focus is on successfully incorporating speech input while preserving the original competence of the model, so that the resulting model can successfully perform both speech-to-text and text-only tasks. None of the aforementioned models are trained to preserve the original model’s performance in text tasks.

Adapting LLMs Previous approaches involve training from scratch with task- and domain-specific data (Singhal et al., 2023; Lewkowycz et al., 2022), performing CPT with a diverse training data mix designed to broadly extend the model’s knowledge (Wu et al., 2023b), or instruction tuning on use-case-specific data (Chen et al., 2023). Recent work has explored combining the latter two approaches (Xu et al., 2024a; Alves et al., 2024; Wei et al., 2021; Roziere et al., 2023). In our approach to integrating DSUs into TOWER, we take inspiration from Alves et al. (2024) in adopting a two-step CPT+IT process. Our work differs in that we focus on adding the speech modality, whereas Alves et al. (2024) focused on increasing the multilingual capabilities of an LLM.

Continuous and Discrete Speech Representations Self-supervised speech representation models produce contextualized high-dimensional

¹en, de, fr, nl, it, es, pt, ko, ru, zh

²[REDACTED]

speech vectors directly from raw audio (Hsu et al., 2021; Baevski et al., 2020; Chen et al., 2022), largely outperforming statistical speech features on downstream tasks (Yang et al., 2021). These continuous representations can be used to derive DSUs that capture both linguistic content and prosody through clustering (Borsos et al., 2023; Kharitonov et al., 2022). DSUs provide better alignment with textual data, facilitating the transfer of successful training settings from the text domain (Cui et al., 2024). Building on Lakhotia et al. (2021), which demonstrated that HuBERT (Hsu et al., 2021) is a powerful feature extractor, several studies have adopted this approach, incorporating a k-means clustering step for discretization (Zhang et al., 2023; Rubenstein et al., 2023; Lam et al., 2024; Chang et al., 2024; Nguyen et al., 2025). Xu et al. (2024b) study the optimal settings to obtain DSUs in terms of cluster size and feature extraction layer. We use their findings to inform our initial choices.

3 SPIRE: A Speech-to-Text LLM

We introduce SPIRE, whose goal is to equip an LLM with speech capabilities while preserving its preexisting text capabilities. As our base LLM we choose TOWER (Alves et al., 2024), which was developed from Llama-2 (Touvron et al., 2023) with a two-step approach: CPT on a mixture of monolingual and parallel data (TOWERBASE), followed by IT on translation-related tasks (TOWERINSTRUCT). We use an approach similar to TOWER to extend the model to the speech modality. First, we perform CPT with a combination of text-only and aligned speech-to-text datasets, followed by IT using both text-only general-purpose and task-specific data curated in TOWERBLOCKS,³ alongside task-specific speech-to-text datasets.

We choose TOWER in particular due to its competitive performance compared to other open alternatives. TOWER-based models were among the best participating systems in the WMT24 general translation task (Kocmi et al., 2024). TOWER’s usage of open source data during the CPT phase along with the release of the TOWERBLOCKS dataset, used in the IT phase, further motivates our choice.

3.1 Speech Discretization

To easily transfer the training set-up of TOWER, we use DSUs as opposed to an auxiliary speech encoder. For all speech datasets that were used, we

follow recent discretization methodology (Zhang et al., 2023; Rubenstein et al., 2023; Chang et al., 2024) to produce DSUs by first extracting continuous speech representations for our speech data from the 22nd layer of an HuBERT-large model, trained on 60K hours of English speech (Hsu et al., 2021), and then using k-means clustering ($K = 5000$) to produce centroids that are used to convert our continuous speech representation into a discrete sequence of cluster IDs.⁴ We train our k-means model on a collection of 235K audio files (approximately 720 hours), drawn from three speech corpora: CoVoST-2 (Wang et al., 2021b), VoxPopuli (Wang et al., 2021a), and Multilingual LibriSpeech (MLS; Pratap et al., 2020). The CoVoST subset consists of 62K audio files from 10,049 speakers, with a maximum of 8 audio files per speaker. The VoxPopuli subset includes 65K audio files from 639 speakers, capped at 250 audio files per speaker. Finally, the MLS subset contains 107K audio files from 5,490 speakers.

3.2 SPIREBASE

The first CPT stage, yielding SPIREBASE, is trained from TOWERBASE-7B⁵ using both text-only and aligned speech-to-text datasets. Following previous work, we include a fraction of TOWER’s original training data to preserve its existing performance (Scialom et al., 2022; de Masson D’Autume et al., 2019).

3.2.1 Data

We use a mixture of monolingual and parallel text in Chinese (zh), Dutch (nl), English (en), French (fr), German (de), Italian (it), Korean (ko), Portuguese (pt), Russian (ru), and Spanish (es), that was sourced from the TOWER training data, as well as English ASR data sourced from popular open-source ASR datasets, as reported in Table 1. Both speech and text data are downsampled to create a 6B token data mixture (5B speech; 1B text), measured by the model tokenizer.⁶ Note that the 5B speech tokens include both DSUs (4.4B tokens) and their text transcriptions (0.6B tokens).

⁴Optimizing the layer selection for feature extraction is a complex research problem (Pasad et al., 2023; Mousavi et al., 2024). In this work we follow the insights from Gow-Smith et al. (2023) and Xu et al. (2024b).

⁵We used TOWER-7B models instead of the 13B or 70B versions due to its lower compute requirements

⁶Preliminary experiments on the data mixture led to this particular choice.

³<https://huggingface.co/datasets/Unbabel/TowerBlocks-v0.2>

Text Data The monolingual text data split corresponds to data from mC4 (Raffel et al., 2019), a multilingual web-crawled corpus which we uniformly sample from across all languages. The parallel data split includes uniformly sampled instances to and from English (en↔xx) for the 10 languages, sourced from various public sources. Further details can be found in Alves et al. (2024).

Speech Data We collect 35K hours of speech data from SPGI Speech (O’Neill et al., 2021), GigaSpeech (Chen et al., 2021), MLS, and VoxPopuli. We normalize as described in Appendix A.1.

3.2.2 CPT Setup

We train SPIREBASE using MegatronLLM (Cano et al., 2023) on 8 A100-80GB GPUs for 6 days. We use the same hyperparameters as TOWER, except for the effective batch size, which in our case is 2,304. To incorporate the DSUs in the CPT stage, we extend the model’s original vocabulary by 5000 types, *e.g.*, <extra_id_x>. This allows us to have a vocabulary that can encode both text in subword units and speech in DSUs. For the extended vocabulary, we initialize new embeddings from a multivariate Gaussian distribution. The mean of this distribution is set to the average of the original embeddings, while the covariance is derived from the empirical covariance of the original embeddings, scaled by a factor of 1×10^{-5} (Hewitt, 2021).

3.3 SPIREFULL

SPIREFULL is obtained by instruction tuning SPIREBASE on task-specific text and speech data.

3.3.1 Data

We use a mixture of text and speech instructions for ASR, MT, and ST. The prompt formats used during training are shown in Appendix A.2.

Text Data We use TOWERBLOCKS (Alves et al., 2024), which includes high quality translation bi-texts between English and the other languages supported by TOWER. It also includes instructions for the translation-related tasks of named entity recognition and automatic post-editing.

ASR Data We use 0.8K hours of ASR data from CommonVoice 18 (CV; Ardila et al., 2020), down-sampling strategy as described in Appendix A.1.

ST Data In our IT set, we use 842 hours of speech across three ST training sets: FLEURS (all nine language pairs; we filter out examples

Dataset	Task	Phase	# DSUs	# Hours
SPGI Speech	ASR	CPT	645M	5.1K
Gigaspeech	ASR	CPT	1.2B	9.9K
MLS	ASR	CPT	2.4B	19.2K
VoxPopuli	ASR	CPT	69M	0.5K
CV	ASR	IT	105M	0.8K
Europarl-ST	ST	IT	122M	1.0K
FLEURS	ST	IT	11M	0.09K
CoVoST-2	ST	IT	12M	0.09K
SPGI Speech	Pseudo-ST	IT	350M	2.8K
GigaSpeech	Pseudo-ST	IT	161M	1.3K
CV	Pseudo-ST	IT	212M	1.7K

Table 1: Statistics for speech training data. Hours are approximated from the number of deduplicated DSUs.

whose transcriptions overlap with the FLORES devtest set), Europarl-ST (Iranzo-Sánchez et al., 2020) (en → {de, es, fr, it, nl, pt}), and CoVoST-2 (en→zh). Since this amounts to far less data for ST than ASR, and since en→{ko, ru} have only examples from the tiny FLEURS set, we augment our speech collection with **pseudo-labeled** data, which has been effective for other ST systems (Barrault et al., 2023). We select 300k ASR examples each from CV, SPGI, and GigaSpeech and translate them to all nine target languages using TowerInstruct-13B.⁷ We then filter examples whose transcript-translation combination has a COMET-QE⁸ (Rei et al., 2022b) score under 85. Finally, for each language pair, we sample 60K examples to be used in direct ST prompts and another 60K to be used in multi-turn prompts. This results in 180K direct ST prompts and 180K multi-turn prompts for each language pair.⁹ The prompt formats are shown in Appendix A.2.

3.3.2 IT Training Setup

We use the chatml template (OpenAI, 2023) to format our instructions in dialogue form. We train models using Axolotl¹⁰ on 4 H100-80GB GPUs for 2.7 days. We use a learning rate of 7×10^{-6} and a cosine scheduler with 100 warm-up steps. We train for 4 epochs with an effective batch size of 576 and a weight decay of 0.01. We impose a maximum sequence length of 4096 and use the AdamW optimizer (Loshchilov and Hutter, 2019). Other hyperparameters are derived from TOWERINSTRUCT (Alves et al., 2024).

⁷<https://huggingface.co/Unbabel/TowerInstruct-13B-v0.1>

⁸<https://huggingface.co/Unbabel/wmt22-cometkiwi-da>

⁹Due to our aggressive filtering, we were left with slightly fewer examples for en → zh.

¹⁰<https://github.com/axolotl-ai-cloud/axolotl>

4 Experiments

We evaluate our models across three tasks: ASR, MT, and ST. First, we present our results for ASR (§4.1), confirming the new capabilities SPIRE has in the speech domain. We then present MT results (§4.2), demonstrating that the speech performance does not come at the expense of the original model’s MT performance. Finally, we present results for ST (§4.3) to investigate model performance on a task that requires both ASR and MT capabilities.

Evaluation Setup Across models and tasks, we perform inference with greedy decoding with a maximum of 256 generated tokens. For the TOWER and SPIRE models, we decode with vllm. However, since vllm does not support all of our baselines, we use alternative libraries (transformers) where necessary. Unless specified otherwise, we use zero-shot prompts for all models and tasks.

4.1 ASR

Datasets and Metrics We evaluate ASR performance across multiple test sets, in order to cover a variety of recording styles: Librispeech (LS) test-clean and test-other (Panayotov et al., 2015), FLEURS (Conneau et al., 2023), and VoxPopuli.¹¹ We report the Word Error Rate (WER) between the hypotheses and gold transcripts, after Whisper normalization (Radford et al., 2023).

Baselines We include the following models:

- **Whisper** (Radford et al., 2023) is an encoder-decoder transformer trained on over 5 million hours of labeled data that performs multilingual ASR and to-English ST. We report results for Whisper-base (74M parameters) and Whisper-large-v3 (1.5B parameters).
- **SeamlessM4T** (Barrault et al., 2023) is an encoder-decoder transformer trained on 406K hours of speech that performs ASR, ST and MT across 100 languages. We report results for SeamlessM4T-large-v2 (2.3B parameters).
- **SALMONN** (Tang et al., 2024) integrates a pre-trained text LLM with separate speech and audio encoders into a single multimodal

¹¹For CPT models, LS is an in-domain evaluation because its training set is part of MLS.

	LibriSpeech		FLEURS	VoxPopuli
	Clean	Other		
Whisper-base	5.0	11.9	12.1	9.8
Whisper-large-v3	1.8	3.7	5.8	9.2
SeamlessM4T	2.6	4.9	8.1	7.5
SALMONN	2.4	5.3	9.3	8.9
Qwen2-Audio	1.6	3.9	6.6	6.5
Spirit-LM	6.0*	11.0*	-	-
HuBERT-large+CTC	4.3	7.6	11.4	14.7
<i>Our models</i>				
SPIREBASE	28.9	56.3	11.0	13.7
SPIREFULL	4.2	7.1	10.7	15.8

*We were unable to reproduce Spirit-LM’s ASR performance; therefore, we report their self-reported LS results using ten-shot prompts.

Table 2: WER on various ASR test sets.

model.¹² SALMONN uses a LoRA adapter (Hu et al., 2022) to align the spaces.

- **Qwen2-Audio** (Chu et al., 2024) integrates audio into Qwen-7B (Bai et al., 2023) using a specialized encoder that is initialized from Whisper large-v3. The resulting model is pretrained on ~520K hours of data spanning speech, sound, and music.
- **Spirit-LM** (Nguyen et al., 2025) is a decoder-only model, trained from Llama-2 on 307B tokens of text, 458K hours of unlabeled speech, and 111K hours of labeled speech. As in SPIRE, it uses HuBERT DSUs.
- **HuBERT-large+CTC** is a CTC-based ASR model trained using the same speech representation model we use for DSU generation, and using the same ASR data from the IT stage (Section 3.3.1).¹³ Unlike SPIRE, this model has access to a very powerful speech representation backbone, however, lacks strong language modeling capabilities.

Results Our results are presented in Table 2. SPIREFULL’s performance demonstrates that performing both the CPT and IT stages is an effective strategy to give speech capabilities to a text LLM. On the other hand, SPIREBASE does not consistently show reasonable speech performance, however, on FLEURS and VoxPopuli we obtain somewhat strong results in the zero-shot settings, which is surprising given that non-instruction-tuned models often struggle to work out-of-domain without

¹²SALMONN uses 4400 hours of speech/audio data in the IT phase but does not specify the large amount of pre-training ASR and audio captioning data used.

¹³The hyperparameters are described in Appendix B.

	en→xx		xx→en	
	C22	spB	C22	spB
SeamlessM4T	87.22	39.0	87.42	39.9
TOWERBASE-7B	87.38	37.8	88.02	41.7
TOWERINSTRUCT-7B	88.45	38.8	88.27	42.0
<i>Our models</i>				
SPIREBASE	87.41	37.4	87.97	41.4
SPIREFULL	88.54	39.3	88.21	41.8

Table 3: COMET-22 (C22) and spBLEU (spB) on the FLORES devtest set between English and the other languages supported by TOWER And SPIRE.

in-context learning examples.¹⁴

Although SPIREFULL does not match the performance of SeamlessM4T, Whisper-large-v3, SALMONN, or Qwen2-Audio, these were trained on far more speech data than our models (around 10x for Qwen2-Audio and SeamlessM4T). Given this training data gap, it is notable that SPIREFULL *does* outperform Whisper-base on LS and FLEURS, and Spirit-LM on all benchmarks Spirit-LM reports at a fraction of the speech data.

SPIREFULL also outperforms the HuBERT-large+CTC baseline on three out of four datasets—an impressive result given that the CTC model has access to continuous features, which SPIREFULL lacks, showing that our compressed discrete representations *can* recover more powerful features.

4.2 MT

Having demonstrated that our training approach works well to initially equip TOWER with speech processing capabilities, we now turn to MT to investigate whether SPIRE can maintain TOWER’s strong performance on MT despite its speech-centric CPT.

Datasets and Metrics We evaluate on two datasets for MT: FLORES-200 (Team et al., 2024), which covers SPIRE’s languages, and the WMT23 test set (Kocmi et al., 2023), which covers en↔{de, ru, zh}. We report COMET-22 (COMET; Rei et al., 2022a) and spBLEU¹⁵ (Papineni et al., 2002) scores via the SacreBLEU toolkit (Post, 2018).

Baselines We compare the SPIRE models to the text-to-text translation performance of Seam-

¹⁴We also tried prompting SPIREBASE with few-shot examples, but the results were much worse, possibly because the length of the DSU sequences led to in-context examples that were too long for the model to handle effectively.

¹⁵nrefs:1|case:mixed|eff:no|tok:flores200|smooth:exp|version:2.5.1

	APE		NER
	en→xx	xx→en	Multilingual
TOWERINSTRUCT-7B	83.08	80.29	71.56
SPIREFULL	83.13	80.08	67.10

Table 4: Results on APE (COMET) and NER (seq. F1).

lessM4T. Additionally, we report the performance of TOWERBASE-7B and TOWERINSTRUCT-7B.

Results Our results show that even after the speech-centric CPT and mixed speech and text IT stage, the SPIRE models retain the original text-only performance of TOWER on both FLORES (Table 3) and WMT23 (Table 5). This indicates that neither CPT nor IT on speech data negatively impacts the model’s ability to perform MT. This is true for both SPIREBASE, which achieves performance comparable to TOWERBASE; and for IT models, where SPIREFULL slightly surpasses the performance of TOWERINSTRUCT on en→xx. SPIREFULL also outperforms SeamlessM4T by both metrics on all WMT23 language pairs, and for both en→xx and xx→en on FLORES.

Translation-related Tasks We follow the evaluation set-up from TOWER (Alves et al., 2024) to additionally evaluate SPIRE on translation-related tasks. In Table 4 we report our results on automatic post-edition (APE) for en↔{de, ru, zh} and named entity recognition (NER) for {de, en, es, fr, it, pt, zh}. SPIRE performs similarly to TOWERINSTRUCT across both tasks and all language directions, maintaining the original text-only capabilities even after training on speech data.

4.3 ST

As SPIRE has shown success at both ASR and MT, we now investigate its performance on ST.

Datasets For ST, we evaluate our models on FLEURS (Conneau et al., 2023), covering ST between all en→xx pairs, and CoVoST-2 (Wang et al., 2021b) for en→{de, zh}. For brevity, we report spBLEU and COMET-22 in Appendix C.

ST approaches As well as direct ST, we report self-cascades, in which each model transcribes the audio before translating its own output to the target language (*i.e.*, ASR followed by MT).

Baselines We compare SPIRE to SeamlessM4T in both direct and cascaded settings. We also report the results of SALMONN and Qwen2-Audio,

	en→de		en→ru		en→zh		de→en		ru→en		zh→en	
	C22	spB	C22	spB	C22	spB	C22	spB	C22	spB	C22	spB
SeamlessM4T	77.76	27.8	83.22	34.2	80.14	29.7	78.69	26.6	80.58	32.5	76.96	23.8
TOWERBASE-7B	79.96	36.1	83.08	34.2	83.49	33.3	83.56	41.1	80.06	32.7	78.48	23.5
TOWERINSTRUCT-7B	82.34	38.8	84.66	34.9	85.09	35.3	84.95	45.1	82.94	36.7	80.14	26.1
<i>Our models</i>												
SPIREBASE	79.88	34.7	83.04	33.7	83.85	32.4	83.19	40.5	80.20	32.4	78.65	23.1
SPIREFULL	82.50	39.5	84.60	34.9	85.37	37.3	85.24	45.2	82.58	36.4	79.92	26.3

Table 5: COMET-22 (C22) and spBLEU (spB) on the WMT23 test set.

which are both 7B parameter models, like SPIRE. However, SALMONN and Qwen2-Audio do not support text-to-text translation, so we use them only for direct ST.¹⁶ There are also coverage differences between the models: while SeamlessM4T can handle all of SPIRE’s language pairs, neither SALMONN nor Qwen2-Audio supports en→ko; SALMONN also does not support en→ru.

Results Our FLEURS ST results are reported in Table 7. SeamlessM4T performs best at direct ST for all language pairs except en→zh. Among the 7B parameter models, SPIREFULL is the best direct model on average, notably beating SALMONN on all language pairs except en→zh. It also outperforms Qwen2-Audio on 6 out of 8 language pairs that Qwen2-Audio supports, and ties or beats it for all except en→zh and en→de.

Performance on CoVoST-2 (Table 6) tells a different story. Although SPIREFULL maintains its advantage over SeamlessM4T in self-cascaded translation, it attains the worst performance on en→zh, while performing similarly to SALMONN for en→de. This shows that the direct ST performance of SPIREFULL is dataset-dependent, which could be a consequence of its relatively small training data.

SPIREFULL achieves the best self-cascaded performance by a significant margin for both datasets, outperforming SeamlessM4T by a large margin in this setting. This demonstrates that SPIREFULL maintains greater robustness to its own outputs than SeamlessM4T, supporting the insight that LLM-based translation models can be very robust to perturbations (Peters and Martins, 2025).

¹⁶Although Whisper is frequently used for ST, we exclude it because it only supports to-English translation, whereas SPIRE is a from-English ST model. Therefore ST comparisons between the two models are impossible.

	en→de		en→zh	
	C22	spB	C22	spB
<i>Self-cascade</i>				
SeamlessM4T	72.40	21.7	72.32	17.0
SPIREFULL	78.05	31.8	79.50	28.1
<i>Direct</i>				
SALMONN	74.98	22.7	80.92	27.8
Qwen2-Audio	82.29	34.5	85.27	38.7
SeamlessM4T	85.95	42.3	83.62	31.3
SPIREFULL	73.96	25.4	74.53	21.0

Table 6: ST results on CoVoST-2.

5 Analysis

The key innovation of our approach is the application of the CPT followed by IT paradigm to discretized speech allowing us to build upon existing text-only capabilities of our base model. Here, we analyze how the composition of these two training phases contributes overall to model performance across all tasks previously evaluated. To that end, we consider several variants of SPIREBASE and SPIREFULL which are described in Table 8 and whose results are reported in Table 9.

- *i*) no CPT was performed and IT was performed with the entire IT data mix (TOWERFULL);
- *ii*) CPT was performed and no data from TOWERBLOCKS was seen during IT (SPIRENObLOCKS), and
- *iii*) CPT was performed and pseudo-labeled ST data and FLEURS were omitted from the IT data mix (SPIRENOPSEUDO).

We report additional datasets in Appendix D.

Effectiveness of CPT and IT Our previous results demonstrated that using both CPT and IT was the most effective strategy. The performance

	de	es	fr	it	ko	nl	pt	ru	zh	avg ₇	avg _{all}
<i>Self-Cascade</i>											
SeamlessM4T	24.2	21.5	37.7	18.9	12.5	16.9	28.2	27.1	14.6	23.1	22.4
SPIREFULL	38.1	29.4	45.3	31.2	23.1	31.2	42.9	33.5	29.0	35.3	33.7
<i>Direct</i>											
SeamlessM4T	39.2	28.0	48.1	30.6	21.5	30.8	47.5	34.3	23.2	35.3	33.7
SALMONN	25.5	20.8	34.3	16.7	0.1	20.5	32.6	3.1	21.9	24.6	19.5
Qwen2-Audio	31.8	23.5	31.3	23.5	5.4	22.3	36.1	23.7	24.7	27.6	24.7
SPIREFULL	31.1	23.5	37.9	25.5	15.4	25.7	37.3	26.9	21.0	28.9	27.1

Table 7: FLEURS ST ex→xx results with self-cascade and direct models in terms of spBLEU. avg₇ covers the 7 language pairs that all models in the table support (excluding en→{ko, ru}).

Model	Base Model	CPT		IT		
		Speech	Text	Speech	Pseudo	Text
TOWERFULL	TowerBase-7B	✗	✗	✓	✓	✓
SPIREBASE	SpireBase	✓	✓	✗	✗	✗
SPIREFULL	SpireBase	✓	✓	✓	✓	✓
<i>SPIRE Variants</i>						
SPIRENOBLOCKS	SpireBase	✓	✓	✓	✓	✗
SPIRENOPSEUDO	SpireBase	✓	✓	✓	✗	✓

Table 8: Ablations of our models. The CPT and IT columns indicate which data was seen during training.

	ASR			MT		ST	
	en→xx			xx→en		en→xx	
	WER	C22	spB	C22	spB	C22	spB
SPIREFULL	4.2	88.54	39.3	88.21	41.8	81.33	27.1
TOWERFULL	9.5	88.57	39.4	88.17	41.7	79.10	26.1
SPIRENOBLOCKS	4.1	82.98	34.2	85.93	36.1	81.11	27.1
SPIRENOPSEUDO	3.9	88.40	38.9	88.22	42.0	62.80	27.1

Table 9: Ablation models and SPIREFULL on LS Clean for ASR, FLORES devtest for MT, and Fleurs for ST reporting WER, COMET-22 (C22), and spBLEU (spB).

gap between SPIREFULL and the TOWERFULL on ASR (5.3 points in LS test-clean) further shows that IT alone is also not as effective. However, for ST we observe that only performing IT leads to a strong model that is capable of performing speech translation unlike SPIREBASE where we also attempted direct ST but the model failed to produce output in the target language, even when given few-shot prompts. Despite the impressive results from TOWERFULL, we still observe the best performance by SPIREFULL showing that while the effect of CPT is not as drastic as in the case of ASR, we still observe gains with a speech-centric CPT phase.

Modality Interplay Our results show that text and speech modalities are orthogonal to each other. Specifically, the performances of TOWERFULL and SPIREFULL show that speech-centric CPT *does not* degrade the text performance of the base model. However, MT quality suffers when TOWERBLOCKS is removed from the IT data, as is shown by SPIRENOBLOCKS’s much weaker performance than SPIREFULL. Simultaneously, SPIREFULL performs on par with SPIRENOBLOCKS on both ASR and ST, indicating that adding text instructions also *does not* degrade performance on speech tasks. It is worth highlighting that a model strong at both MT and ASR (SPIRENOPSEUDO) does not lead to a strong ST

model, showing surprisingly that competence at MT is not very helpful for direct ST.

6 Conclusion

In this work we presented SPIRE, a simple and effective recipe for adapting a text-based, translation-specialist LLM to the speech modality while preserving the original performance on text-based tasks. We investigated the impact of speech integration on two stages of LLM adaptation, CPT and IT, finding that both contribute to the final model’s performance on speech tasks. Our results demonstrate that we are able to successfully integrate a new modality without compromising the original model’s capabilities. SPIRE achieves competitive performance on ASR, while its MT abilities remain on par with the original TOWER model. Finally, for the ST task, we find that the leveraging ASR and MT data does not directly transfer to ST performance. Nonetheless, the model achieves promising performance with both direct and self-cascaded ST.

As future work, we intend to extend this recipe to multilingual settings by replacing our English HuBERT speech component by the multilingual mHuBERT-147 (Boito et al., 2024). To benefit the community, we only use publicly available and licensed data to train our models, making our results reproducible.

Limitations

The downstream tasks we evaluate on are restricted to MT and ASR/ST, which provides an idea of the model performance but do not give us the full picture. We plan to address this by utilizing the LM-harness evaluation (Gao et al., 2024) to evaluate on a suite of text-based benchmarks such as MMLU (Multitask Language Understanding) (Hendrycks et al., 2021b,a), Arc (Commonsense Reasoning) (Clark et al., 2018), Belebele (Reading Comprehension) (Bandarkar et al., 2024), and HellaSwag (Sentence Completion) (Zellers et al., 2019). Lastly, our model handles speech and text on the input side but is currently limited to generating only text.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Duarte Miguel Alves, José Pombal, Nuno M Guerreiro, Pedro Henrique Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and Andre Martins. 2024. *Tower: An open multilingual large language model for translation-related tasks*. In *First Conference on Language Modeling*.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. *Common voice: A massively-Â-multilingual speech corpus*. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 4218–4222, Marseille, France. European Language Resources Association.

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. *Qwen technical report*. *Preprint*, arXiv:2309.16609.

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. 2024. *The belebele benchmark: a parallel reading comprehension dataset in 122 language variants*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 749–775, Bangkok, Thailand. Association for Computational Linguistics.

Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, et al. 2023. Seamlessm4t-massively multilingual & multimodal machine translation. *arXiv preprint arXiv:2308.11596*.

Marcely Zanon Boito, Vivek Iyer, Nikolaos Lagos, Laurent Besacier, and Ioan Calapodescu. 2024. mHuBERT-147: A Compact Multilingual HuBERT Model. In *Interspeech 2024*.

Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, et al. 2023. Audioldm: a language modeling approach to audio generation. *IEEE/ACM transactions on audio, speech, and language processing*, 31:2523–2533.

Alejandro Hernández Cano, Matteo Pagliardini, Andreas Köpf, Kyle Matoba, Amirkeivan Mohtashami, Xingyao Wang, Olivia Simin Fan, Axel Marmet, Deniz Bayazit, Igor Krawczuk, Zeming Chen, Francesco Salvi, Antoine Bosselut, and Martin Jaggi. 2023. *epfilm megatron-llm*.

Xuankai Chang, Brian Yan, Kwanghee Choi, Jee-Weon Jung, Yichen Lu, Soumi Maiti, Roshan Sharma, Jia-tong Shi, Jinchuan Tian, Shinji Watanabe, et al. 2024. Exploring speech recognition, translation, and understanding with discrete speech units: A comparative study. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 11481–11485. IEEE.

Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei Zou, Xiangang Li, Xuchen Yao, Yongqing Wang, Zhao You, and Zhiyong Yan. 2021. *GigaSpeech: An Evolving, Multi-Domain ASR Corpus with 10,000 Hours of Transcribed Audio*. In *Proc. Interspeech 2021*, pages 3670–3674.

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, and Xiong Xiao. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518.

Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco

699	Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. 2023. Meditron-70b: Scaling medical pretraining for large language models. <i>arXiv preprint arXiv:2311.16079</i> .	754
700		755
701		756
702		757
703	Ju-Chieh Chou, Chung-Ming Chien, Wei-Ning Hsu, Karen Livescu, Arun Babu, Alexis Conneau, Alexei Baevski, and Michael Auli. 2023a. Toward joint language modeling for speech units and text. <i>arXiv preprint arXiv:2310.08715</i> .	758
704		759
705		760
706		761
707		762
708	Ju-Chieh Chou, Chung-Ming Chien, Wei-Ning Hsu, Karen Livescu, Arun Babu, Alexis Conneau, Alexei Baevski, and Michael Auli. 2023b. Toward joint language modeling for speech units and text. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 6582–6593, Singapore. Association for Computational Linguistics.	763
709		764
710		765
711		766
712		767
713		768
714		769
715	Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. 2024. Qwen2-audio technical report. <i>arXiv preprint arXiv:2407.10759</i> .	770
716		771
717		772
718		773
719		774
720	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>arXiv:1803.05457v1</i> .	775
721		776
722		777
723		778
724		779
725	Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2023. Fleurs: Few-shot learning evaluation of universal representations of speech. In <i>2022 IEEE Spoken Language Technology Workshop (SLT)</i> , pages 798–805. IEEE.	780
726		781
727		782
728		783
729		784
730		785
731	Wenqian Cui, Dianzhi Yu, Xiaoqi Jiao, Ziqiao Meng, Guangyan Zhang, Qichao Wang, Yiwen Guo, and Irwin King. 2024. Recent advances in speech language models: A survey. <i>arXiv preprint arXiv:2410.03751</i> .	786
732		787
733		788
734		789
735	Cyprien de Masson D’Autume, Sebastian Ruder, Lingpeng Kong, and Dani Yogatama. 2019. Episodic memory in lifelong language learning. <i>Advances in Neural Information Processing Systems</i> , 32.	790
736		791
737		792
738		793
739	Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. 2024. Moshi: a speech-text foundation model for real-time dialogue. <i>arXiv preprint arXiv:2410.00037</i> .	794
740		795
741		796
742		797
743		798
744	Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. 2023. Palm-e: An embodied multimodal language model. <i>arXiv preprint arXiv:2303.03378</i> .	799
745		800
746		801
747		802
748		803
749	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. A framework for few-shot language model evaluation.	804
750		805
751		806
752		807
753		808
	Edward Gow-Smith, Alexandre Berard, Marcey Zanon Boito, and Ioan Calapodescu. 2023. NAVER LABS Europe’s multilingual speech translation systems for the IWSLT 2023 low-resource track. In <i>Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)</i> , pages 144–158, Toronto, Canada (in-person and online). Association for Computational Linguistics.	809
	Michael Hassid, Tal Remez, Tu Anh Nguyen, Itai Gat, Alexis Conneau, Felix Kreuk, Jade Copet, Alexandre Defossez, Gabriel Synnaeve, Emmanuel Dupoux, et al. 2024. Textually pretrained speech language models. <i>Advances in Neural Information Processing Systems</i> , 36.	810
	Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2021a. Aligning ai with shared human values. <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> .	811
	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021b. Measuring massive multitask language understanding. <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> .	812
	John Hewitt. 2021. Initializing new word embeddings for pretrained language models. https://nlp.stanford.edu/johnhew/vocab-expansion.html .	813
	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. <i>IEEE/ACM transactions on audio, speech, and language processing</i> , 29:3451–3460.	814
	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .	815
	Shujie Hu, Long Zhou, Shujie Liu, Sanyuan Chen, Lingwei Meng, Hongkun Hao, Jing Pan, Xunying Liu, Jinyu Li, Sunit Sivasankaran, et al. 2024. Wavllm: Towards robust and adaptive speech large language model. <i>arXiv preprint arXiv:2404.00656</i> .	816
	Rongjie Huang, Mingze Li, Dongchao Yang, Jiaotong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. 2024. Audiogpt: Understanding and generating speech, music, sound, and talking head. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 23802–23804.	817

- Javier Iranzo-Sánchez, Joan Albert Silvestre-Cerda, Javier Jorge, Nahuel Roselló, Adria Giménez, Albert Sanchis, Jorge Civera, and Alfons Juan. 2020. Europarl-st: A multilingual corpus for speech translation of parliamentary debates. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8229–8233. IEEE.
- Eugene Kharitonov, Ann Lee, Adam Polyak, Yossi Adi, Jade Copet, Kushal Lakhotia, Tu Anh Nguyen, Morgane Riviere, Abdelrahman Mohamed, Emmanuel Dupoux, and Wei-Ning Hsu. 2022. [Text-free prosody-aware generative spoken language modeling](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8666–8681, Dublin, Ireland. Association for Computational Linguistics.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinhórfur Steingrímsson, and Vilém Zouhar. 2024. [Findings of the WMT24 general machine translation shared task: The LLM era is here but MT is not solved yet](#). In *Proceedings of the Ninth Conference on Machine Translation*, pages 1–46, Miami, Florida, USA. Association for Computational Linguistics.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. [Findings of the 2023 conference on machine translation \(WMT23\): LLMs are here but not quite there yet](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.
- Kushal Lakhotia, Eugene Kharitonov, Wei-Ning Hsu, Yossi Adi, Adam Polyak, Benjamin Bolte, Tu-Anh Nguyen, Jade Copet, Alexei Baevski, Abdelrahman Mohamed, and Emmanuel Dupoux. 2021. [On generative spoken language modeling from raw audio](#). *Transactions of the Association for Computational Linguistics*, 9:1336–1354.
- Tsz Kin Lam, Alexandra Birch, and Barry Haddow. 2024. Compact speech translation models via discrete speech units pretraining. *arXiv preprint arXiv:2402.19333*.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. 2022. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. [Visual Instruction Tuning \(LLaVA\)](#). *arXiv preprint*. ArXiv:2304.08485 [cs].
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.
- Soumi Maiti, Yifan Peng, Shukjae Choi, Jee-weon Jung, Xuankai Chang, and Shinji Watanabe. 2024. Voxlm: Unified decoder-only models for consolidating speech recognition, synthesis and speech, text continuation tasks. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 13326–13330. IEEE.
- Pedro Henrique Martins, Patrick Fernandes, João Alves, Nuno M Guerreiro, Ricardo Rei, Duarte M Alves, José Pombal, Amin Farajian, Manuel Faysse, Mateusz Klimaszewski, et al. 2024. Eurollm: Multilingual language models for europe. *arXiv preprint arXiv:2409.16235*.
- Pooneh Mousavi, Jarod Duret, Salah Zaiem, Luca Della Libera, Artem Ploujnikov, Cem Subakan, and Mirco Ravanelli. 2024. How should we extract discrete audio tokens from self-supervised models? *arXiv preprint arXiv:2406.10735*.
- Tu Anh Nguyen, Benjamin Muller, Bokai Yu, Marta R Costa-Jussa, Maha Elbayad, Sravya Popuri, Christophe Ropers, Paul-Ambroise Duquenne, Robin Algayres, Ruslan Mavlyutov, et al. 2025. Spiritlm: Interleaved spoken and written language model. *Transactions of the Association for Computational Linguistics*, 13:30–52.
- Patrick K O’Neill, Vitaly Lavrukhin, Somshubra Majumdar, Vahid Noroozi, Yuekai Zhang, Oleksii Kuchaiev, Jagadeesh Balam, Yuliya Dovzhenko, Keenan Freyberg, Michael D Shulman, et al. 2021. Sgspispeech: 5,000 hours of transcribed financial audio for fully formatted end-to-end speech recognition. *arXiv preprint arXiv:2104.02014*.
- OpenAI. 2023. URL <https://github.com/openai/openai-python/blob/release-v0.28.1/chatml.md>.
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5206–5210. IEEE.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

924	Ankita Pasad, Bowen Shi, and Karen Livescu. 2023.	Paul K Rubenstein, Chulayuth Asawaroengchai,	982
925	Comparative layer-wise analysis of self-supervised	Duc Dung Nguyen, Ankur Bapna, Zalán Borsos,	983
926	speech models. In <i>ICASSP 2023-2023 IEEE Interna-</i>	Félix de Chaumont Quitry, Peter Chen, Dalia El	984
927	<i>tional Conference on Acoustics, Speech and Signal</i>	Badawy, Wei Han, Eugene Kharitonov, et al. 2023.	985
928	<i>Processing (ICASSP)</i> , pages 1–5. IEEE.	Audiopalm: A large language model that can speak	986
		and listen. <i>arXiv preprint arXiv:2306.12925</i> .	987
929	Ben Peters and André F. T. Martins. 2025. Did trans-	Thomas Scialom, Tuhin Chakrabarty, and Smaranda	988
930	lation models get more robust without anyone even	Muresan. 2022. Fine-tuned language models are	989
931	noticing? <i>Preprint</i> , arXiv:2403.03923.	continual learners . In <i>Proceedings of the 2022 Con-</i>	990
932	Matt Post. 2018. A call for clarity in reporting BLEU	<i>ference on Empirical Methods in Natural Language</i>	991
933	scores . In <i>Proceedings of the Third Conference on</i>	<i>Processing</i> , pages 6107–6122, Abu Dhabi, United	992
934	<i>Machine Translation: Research Papers</i> , pages 186–	Arab Emirates. Association for Computational Lin-	993
935	191, Brussels, Belgium. Association for Computa-	guistics.	994
936	tional Linguistics.		
937	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel	Yu Shu, Siwei Dong, Guangyao Chen, Wenhao Huang,	995
938	Synnaeve, and Ronan Collobert. 2020. MLS: A	Ruihua Zhang, Daochen Shi, Qiqi Xiang, and Yemin	996
939	Large-Scale Multilingual Dataset for Speech Re-	Shi. 2023. Llam: Large language and speech model.	997
940	search . In <i>Proc. Interspeech 2020</i> , pages 2757–2761.	<i>arXiv preprint arXiv:2308.15930</i> .	998
941	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mah-	999
942	man, Christine McLeavey, and Ilya Sutskever. 2023.	davi, Jason Wei, Hyung Won Chung, Nathan Scales,	1000
943	Robust speech recognition via large-scale weak su-	Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl,	1001
944	pervision. In <i>International conference on machine</i>	et al. 2023. Large language models encode clinical	1002
945	<i>learning</i> , pages 28492–28518. PMLR.	knowledge. <i>Nature</i> , 620(7972):172–180.	1003
946	Srijith Radhakrishnan, Chao-Han Yang, Sumeer Khan,	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao	1004
947	Rohit Kumar, Narsis Kiani, David Gomez-Cabrero,	Chen, Tian Tan, Wei Li, Lu Lu, MA Zejun, and Chao	1005
948	and Jesper Tegnér. 2023. Whispering LLaMA: A	Zhang. 2024. Salmonn: Towards generic hearing	1006
949	cross-modal generative error correction framework	abilities for large language models. In <i>The Twelfth</i>	1007
950	for speech recognition . In <i>Proceedings of the 2023</i>	<i>International Conference on Learning Representa-</i>	1008
951	<i>Conference on Empirical Methods in Natural Lan-</i>	<i>tions</i> .	1009
952	<i>guage Processing</i> , pages 10007–10016, Singapore.		
953	Association for Computational Linguistics.	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	1010
		Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	1011
954	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Schalkwyk, Andrew M Dai, Anja Hauth, Katie	1012
955	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	Millican, et al. 2023. Gemini: a family of	1013
956	Wei Li, and Peter J. Liu. 2019. Exploring the limits	highly capable multimodal models. <i>arXiv preprint</i>	1014
957	of transfer learning with a unified text-to-text trans-	<i>arXiv:2312.11805</i> .	1015
958	former . <i>arXiv e-prints</i> .		
		NLLB Team et al. 2024. Scaling neural machine trans-	1016
959	Ricardo Rei, José G. C. de Souza, Duarte Alves,	lation to 200 languages. <i>Nature</i> , 630(8018):841.	1017
960	Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova,		
961	Alon Lavie, Luisa Coheur, and André F. T. Martins.	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1018
962	2022a. COMET-22: Unbabel-IST 2022 submission	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1019
963	for the metrics shared task . In <i>Proceedings of the</i>	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1020
964	<i>Seventh Conference on Machine Translation (WMT)</i> ,	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	1021
965	pages 578–585, Abu Dhabi, United Arab Emirates	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	1022
966	(Hybrid). Association for Computational Linguistics.	Jude Fernandes, Jeremy Fu, Wenying Fu, Brian Fuller,	1023
		Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	1024
967	Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro,	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	1025
968	Chrysoula Zerva, Ana C Farinha, Christine Maroti,	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	1026
969	José G. C. de Souza, Taisiya Glushkova, Duarte	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	1027
970	Alves, Luisa Coheur, Alon Lavie, and André F. T.	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	1028
971	Martins. 2022b. CometKiwi: IST-unbabel 2022 sub-	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	1029
972	mission for the quality estimation shared task . In	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	1030
973	<i>Proceedings of the Seventh Conference on Machine</i>	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	1031
974	<i>Translation (WMT)</i> , pages 634–645, Abu Dhabi,	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	1032
975	United Arab Emirates (Hybrid). Association for Com-	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	1033
976	putational Linguistics.	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	1034
		lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	1035
977	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	1036
978	Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi,	Melanie Kambadur, Sharan Narang, Aurelien Ro-	1037
979	Jingyu Liu, Romain Sauvestre, Tal Remez, et al. 2023.	driguez, Robert Stojnic, Sergey Edunov, and Thomas	1038
980	Code llama: Open foundation models for code. <i>arXiv</i>	Scialom. 2023. Llama 2: Open foundation and fine-	1039
981	<i>preprint arXiv:2308.12950</i> .	tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	1040

1041	Viet Anh Trinh, Rosy Southwell, Yiwen Guan, Xinlu	1098
1042	He, Zhiyong Wang, and Jacob Whitehill. 2024. Dis-	1099
1043	crete multimodal transformers with a pretrained large	1100
1044	language model for mixed-supervision speech pro-	1101
1045	cessing. <i>arXiv preprint arXiv:2406.06582</i> .	1102
1046	Changhan Wang, Morgane Riviere, Ann Lee, Anne Wu,	1103
1047	Chaitanya Talnikar, Daniel Haziza, Mary Williamson,	1104
1048	Juan Pino, and Emmanuel Dupoux. 2021a. VoxPop-	1105
1049	uli: A large-scale multilingual speech corpus for rep-	1106
1050	resentation learning, semi-supervised learning and	
1051	interpretation . In <i>Proceedings of the 59th Annual</i>	
1052	<i>Meeting of the Association for Computational Lin-</i>	
1053	<i>guistics and the 11th International Joint Conference</i>	
1054	<i>on Natural Language Processing (Volume 1: Long</i>	
1055	<i>Papers)</i> , pages 993–1003, Online. Association for	
1056	Computational Linguistics.	
1057	Changhan Wang, Anne Wu, Jiatao Gu, and Juan Pino.	
1058	2021b. Covost 2 and massively multilingual speech	
1059	translation. <i>Interspeech 2021</i> .	
1060	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	
1061	Adams Wei Yu, Brian Lester, Nan Du, Andrew M	
1062	Dai, and Quoc V Le. 2021. Finetuned language mod-	
1063	els are zero-shot learners. In <i>International Confer-</i>	
1064	<i>ence on Learning Representations</i> .	
1065	T Wolf. 2019. Huggingface’s transformers: State-of-	
1066	the-art natural language processing. <i>arXiv preprint</i>	
1067	<i>arXiv:1910.03771</i> .	
1068	Jian Wu, Yashesh Gaur, Zhuo Chen, Long Zhou, Yi-	
1069	meng Zhu, Tianrui Wang, Jinyu Li, Shujie Liu,	
1070	Bo Ren, Linquan Liu, et al. 2023a. On decoder-only	
1071	architecture for speech-to-text and large language	
1072	model integration. In <i>2023 IEEE Automatic Speech</i>	
1073	<i>Recognition and Understanding Workshop (ASRU)</i> ,	
1074	pages 1–8. IEEE.	
1075	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski,	
1076	Mark Dredze, Sebastian Gehrmann, Prabhajan Kam-	
1077	badur, David Rosenberg, and Gideon Mann. 2023b.	
1078	Bloomberggpt: A large language model for finance.	
1079	<i>arXiv preprint arXiv:2303.17564</i> .	
1080	Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Has-	
1081	san Awadalla. 2024a. A paradigm shift in machine	
1082	translation: Boosting translation performance of	
1083	large language models. In <i>The Twelfth International</i>	
1084	<i>Conference on Learning Representations</i> .	
1085	Yaoxun Xu, Shi-Xiong Zhang, Jianwei Yu, Zhiyong	
1086	Wu, and Dong Yu. 2024b. Comparing discrete and	
1087	continuous space llms for speech recognition. In	
1088	<i>Proc. Interspeech 2024</i> .	
1089	Hongfei Xue, Wei Ren, Xuelong Geng, Kun Wei, Long-	
1090	hao Li, Qijie Shao, Linju Yang, Kai Diao, and Lei	
1091	Xie. 2024. Ideal-llm: Integrating dual encoders and	
1092	language-adapted llm for multilingual speech-to-text.	
1093	<i>arXiv preprint arXiv:2409.11214</i> .	
1094	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,	
1095	Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,	
1096	Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 tech-	
1097	nical report. <i>arXiv preprint arXiv:2412.15115</i> .	
	Shu-Wen Yang, Po-Han Chi, Yung-Sung Chuang,	
	Cheng-I Jeff Lai, Kushal Lakhotia, Yist Y. Lin,	
	Andy T. Liu, Jiatong Shi, Xuankai Chang, Guan-Ting	
	Lin, Tzu-Hsien Huang, Wei-Cheng Tseng, Ko tik	
	Lee, Da-Rong Liu, Zili Huang, Shuyan Dong, Shang-	
	Wen Li, Shinji Watanabe, Abdelrahman Mohamed,	
	and Hung yi Lee. 2021. SUPERB: Speech Process-	
	ing Universal PERformance Benchmark . In <i>Proc.</i>	
	<i>Interspeech 2021</i> , pages 1194–1198.	
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	
	Farhadi, and Yejin Choi. 2019. HellaSwag: Can a ma-	
	chine really finish your sentence? In <i>Proceedings of</i>	
	<i>the 57th Annual Meeting of the Association for Com-</i>	
	<i>putational Linguistics</i> , pages 4791–4800, Florence,	
	Italy. Association for Computational Linguistics.	
	Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing	
	Qu, Yong Jae Lee, and Yi Ma. 2024. Investigating the	
	catastrophic forgetting in multimodal large language	
	model fine-tuning. In <i>Conference on Parsimony and</i>	
	<i>Learning</i> , pages 202–227. PMLR.	
	Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan,	
	Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023.	
	SpeechGPT: Empowering large language models	
	with intrinsic cross-modal conversational abilities .	
	In <i>Findings of the Association for Computational</i>	
	<i>Linguistics: EMNLP 2023</i> , pages 15757–15773, Sin-	
	gapore. Association for Computational Linguistics.	

A Data

A.1 Speech Data Preprocessing

Normalization In order to make transcripts consistent across the different datasets, the following normalization is applied:

- **GigaSpeech (CPT):** we lower-case the text and replace punctuation tags: <COMMA>, <PERIOD>, QUESTIONMARK>, <EXCLAMATIONPOINT> with their appropriate punctuation.
- **MLS (CPT):** we apply a tail-end normalization step here which uniformly samples each speaker to have at maximum 13 transcriptions. This allows us to have a better distribution of speakers.
- **CV (IT):** we subsampled from CommonVoice to ensure a minimum duration of 3 seconds per sample. To enhance transcript diversity, we limit each transcript to 4 unique speakers.

Deduplication As in previous work (Zhang et al., 2023; Rubenstein et al., 2023; Chang et al., 2024), we merge consecutive repeated DSU tokens into a single token to reduce sequence length.

A.2 Prompt Format

Table 10 show the prompts used during both training stages.

ASR (CPT)
Speech:<extra_id_i>...<extra_id_j> English: {TRANSCRIPT}
MT (CPT)
Source_lang: Source-sentence Target_lang: {TRANSLATION}
ASR (IT)
Speech: <extra_id_i>...<extra_id_j> English: {TRANSCRIPT}
Direct ST (IT)
Speech: <extra_id_i>...<extra_id_j> TARGET_LANG: {TRANSLATION}
Multi-turn ST (IT)
Speech: <extra_id_i>...<extra_id_j> English:{TRANSCRIPT} TARGET_LANG: {TRANSLATION}

Table 10: Prompt formats for CPT and IT.

B CTC-based ASR model

We train a CTC-based ASR model using the HuggingFace Transformers library (Wolf, 2019), leveraging the ASR data from the IT stage (CommonVoice, Table 1) as training data. Our ASR model is made of the HuBERT-Large¹⁷ speech representation model, followed by three hidden layers and a vocabulary projection layer. We train for 50 epochs with a dropout of 0.3 and a learning rate of 1e-4 with a warm-up ratio of 0.15. The best checkpoint is selected using CER scores. This was obtained at step 220K (at epoch 12.8).

C ST results

Table 11 report results of ST on FLEURS across baseline models and SPIREFULL. We report COMET-22. We observe the same trend in scores as reported by spBLEU where in SPIREFULL obtains the best self-cascaded performance while beating Qwen2-Audio and SALMONN on direct ST across most language pairs. SeamlessM4T obtains the overall best performance in direct ST.

D Ablation results

Table 12 reports results from all remaining evaluation datasets across ASR, MT, and ST. We report the same metrics as in Section 4. Here as well, we note that in MT, the inclusion of speech data did not degrade text-only performance (SPIREFULL vs. TOWERFULL). Similarly, the inclusion of task-specific text data also did not harm performance on ASR (SPIRENOBLOCKS vs. SPIREFULL). Lastly, SPIREFULL has the best performing direct ST system, further showing that individual task competencies (in MT and ASR) do not contribute directly to a compositional task (ST) but rather the inclusion of task-specific data leads to the highest gains (SPIRENOPSEUDO vs SPIREFULL).

¹⁷<https://huggingface.co/facebook/hubert-large-1160k>

	de	es	fr	it	ko	nl	pt	ru	zh	avg ₇	avg _{all}
Self-Cascade											
SeamlessM4T	72.69	76.97	78.06	76.03	75.33	72.58	78.25	79.38	69.76	74.91	75.45
SPiREFULL	84.26	83.32	84.70	85.16	86.89	84.91	86.01	86.45	85.21	84.80	85.21
Direct											
SeamlessM4T	84.79	83.20	85.32	85.03	85.17	85.17	86.75	86.31	79.90	84.31	84.63
SALMONN	77.41	77.99	79.95	74.47	61.07	77.18	80.94	53.05	81.63	78.51	73.74
Qwen2-Audio	79.82	80.43	79.44	81.28	69.33	78.75	83.41	77.90	80.71	80.55	79.01
SPiREFULL	80.16	79.82	80.68	81.63	82.62	81.93	83.18	82.19	79.76	81.02	81.33

Table 11: FLEURS ST ex→xx results with self-cascade and direct models in terms of COMET-22. avg₇ covers the 7 language pairs that all models in the table support (excluding en→{ko, ru}).

	ASR			MT				ST	
	WER			C22	spB	C22	spB	C22	spB
	LS Other	Fleurs	VoxPopuli	en→xx		xx→en		en→xx	
SPiREFULL	7.1	10.7	15.8	84.16	37.2	82.58	41.8	81.33	27.1
TOWERFULL	13.8	14.3	40.7	84.19	36.9	82.25	35.6	71.52	20.1
SPiRENObLOCKS	7.4	10.4	15.8	73.12	26.9	74.78	25.1	74.02	23.2
SPiRENOPSEUDO	7.3	11.1	14.3	83.93	36.9	82.50	35.9	59.88	6.8

Table 12: Ablation models and SPiREFULL on LS Other, Fleur, VoxPopuli for ASR, WMT23 for MT, and CoVoST-2 for ST reporting WER, COMET-22 (C22), and spBLEU (spB).