# Languages in Multilingual Speech Foundation Models Align Both Phonetically and Semantically

**Ryan Soh-Eun Shim**[1,2] **Domenico De Cristofaro**[3] **Chengzhi Martin Hu**[1]
**Alessandro Vietti**[3] **Barbara Plank**[1,2]
[1]MaiNLP, Center for Information and Language Processing, LMU Munich, Germany
[2]Munich Center for Machine Learning (MCML), Munich, Germany
[3]ALPS, Free University of Bozen-Bolzano, Bozen-Bolzano, Italy

## Abstract

Cross-lingual alignment in pretrained language models (LMs) has enabled efficient transfer in text-based LMs. Such an alignment has also been observed in speech foundation models. However, it remains an open question whether findings and methods from text-based cross-lingual alignment apply to speech. Building on prior work on spoken translation retrieval, we perform pronunciation-controlled experiments to observe if cross-lingual alignment can indeed occur in such models on a *semantic* basis, instead of relying on *phonetic* similarities. Our findings indicate that even in the absence of phonetic cues, spoken translation retrieval accuracy remains relatively stable. We follow up with a controlled experiment on a word-level dataset of cross-lingual synonyms and near-homophones, confirming the existence of both phonetic and semantic knowledge in the encoder. Finally, we qualitatively examine the transcriptions produced by early exiting the encoder, where we observe that speech translation produces semantic errors that are characterized by phonetic similarities to corresponding words in the source language. We apply this insight from early exiting to speech recognition in seven low-resource languages unsupported by the Whisper model, and achieve improved accuracy in all languages examined, particularly for languages with transparent orthographies.

## 1 Introduction

Text-based multilingual models such as mBERT (Devlin et al., 2018) and XLM (Lample and Conneau, 2019; Conneau et al., 2020) enable cross-lingual transfer by mapping language-specific input into a shared semantic space. In speech, a growing body of work has shown speech foundation models to exhibit emergent multilingual capabilities (Peng et al., 2023; Yang et al., 2024), implying the existence of cross-lingual alignment in such models. Prior work probes such alignment through spoken translation retrieval (Abdullah et al., 2024; Ma et al., 2025). However, *to what degree does spoken translation retrieval rely on semantic features?* A possible confounding factor in using speech retrieval as a proxy is the existence of *pronunciation*-level cues such as cognates, loanwords, and proper nouns, that could serve as shortcuts even in the absence of a shared semantic space for retrieving semantically equivalent utterances in a different language (Table 1). In this work, we explicitly control for such pronunciation-level cues by constructing a challenge set without pronunciation-level cues between typologically-distant languages, with the goal of understanding whether and to what extent such shortcuts impact the cross-lingual alignment in speech. In addition, we perform early exiting experiments in the encoder layers to observe how the input is iteratively refined as it maps to the shared space, where we observe the influence of pronunciation to be stronger in the earlier layers. Our work contributes to the following research areas:

i. **Speech Retrieval Explainability**: Extending the speech retrieval metric of SeqSim (Ma et al., 2025), we propose **SeqSimInterp** (Figure 2), a method to derive human-interpretable insights into the similarity scores computed by SeqSim. Concretely, SeqSimInterp highlights word-level contributions to these scores, enabling a direct analysis of which cross-lingual word pairs drive retrieval decisions. We find that semantically equivalent words in cross-lingual pairs indeed contribute significantly to the similarity score.

ii. **Pronunciation-controlled Insights**: We propose a challenge set devoid of *pronunciation*-level cues, where we show cross-lingual retrieval accuracy to remain stable, despite noticeable performance drops.

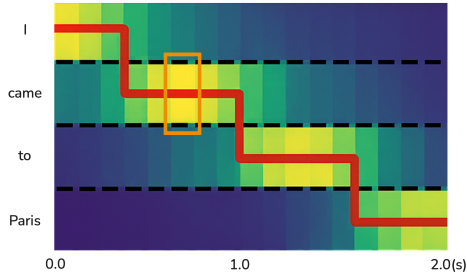iii. **Early Exiting**: To understand how represen-

Figure 1: Pasad et al. (2024) find the center frame of an audio embedding to retain word identity information. As such, we infer the center frame representation for each word in an utterance by way of word-level timestamps. The timestamps are obtained through applying dynamic time warping to cross-attention weights (Zusag et al., 2024).

tations are iteratively refined in the encoder layers, we propose early exiting experiments. We observe speech translation to produce semantic errors characterized by phonetic similarities to corresponding words in the source language (Table 1). We apply this finding to speech recognition in seven low-resource languages, achieving consistent zero-shot improvements.

| | |
|---|---|
| English: | You can **use** a **boda-boda** to get around **Goma**. The normal price is ∼**500 Congolese Francs** for a short ride. |
| Italian: | È possibile **usare** un **boda-boda** per muoversi a **Goma**. Il prezzo normale è di circa **500 franchi congolesi** per una corsa breve. |
| Chinese: | 你可以乘坐**boda-boda**游览戈马 (**gema**)。短途车程的正常价格是500 刚果法郎 (**gangguo falang**)。 |

Table 1: Example parallel utterances from the FLEURS dataset, where highlighted items share pronunciation similarities that spoken translation retrieval could rely on instead of semantic cues. Colors indicate phonetically similar items.

## 2 Related Work

**Cross-Lingual Alignment in Speech** Although in text-based models such as mBERT (Devlin et al., 2019) and XLM (Lample and Conneau, 2019) cross-lingual alignment is widely observed, for a model trained 90% on speech recognition, it is not fully known to what degree the task of speech recognition drives the model to acquire semantic knowledge, particularly given the monotonous alignment of speech and text. As such, very recent works has started to look at what audio representa-

tions capture. For instance, Mohebbi et al. (2024) disentangle acoustic features and textual content in speech models, but they look at supervised fine-tuning to explicitly disentangle, while we examine what their representations capture out of the box, without further training. In addition, a growing body of work has shown speech foundation models to exhibit emergent multilingual capabilities (Peng et al., 2023), suggesting the existence of cross-lingual alignment in such models. Recent work has found models trained on speech translation to exhibit strong capacity in spoken translation retrieval (Abdullah et al., 2024), implying speech translation models to map to a shared semantic subspace. Ma et al. (2025) also make use of spoken translation retrieval to analyze cross-lingual alignment in Whisper's encoder, which allows them to freeze the encoder and finetune only the decoder for speech translation to new target languages. In this work, we explicitly control for such pronunciation-level cues by constructing a challenge set devoid of proper nouns and loanwords between typologically-distant languages.

**Lexical Semantics in Speech** Moreover, the continuous nature of speech tokens and the lack of words boundaries complicates the adaptation of text-based methods. Work that probe the lexical semantics of speech models thus often rely on textual transcriptions to obtain corresponding timestamps in the audio, allowing embeddings for the words to be sliced out based on temporal information. Pasad et al. (2024) perform a large-scale analysis of what word-level knowledge self-supervised speech models contain across layers. Choi et al. (2024) extend this line of inquiry by studying whether self-supervised speech models encode more phonetic or semantic knowledge within and across languages. However, these approaches depend heavily on isolated or short-span word inputs, which limits the semantic context available to the model. In contrast, our approach mitigates this limitation by extracting word-level embeddings from utterance-level inputs, where contextual information is preserved.

**Early Exiting** Prior work has shown that representations at earlier layers of a neural network are often already adequate for making a correct prediction (Kaya et al., 2019), which allows for early exiting strategies that minimize compute (Schuster et al., 2022) and mitigate the influence of false demonstrations (Halawi et al., 2024). Recent work in mechanistic interpretability (Rai et al., 2025)

extend this line of research by using early exiting as a way to understand model internals. For instance, the logit lens method (nostalgebraist, 2020) interprets the incremental layer updates of large language models by projecting intermediate hidden states into the vocabulary space. Belrose et al. (2023) and Din et al. (2024) extend the logit lens approach by learning affine transformations that match earlier layer logits to final layer logits, which obtains more robust predictions. However, the work above has mainly been applied to decoder-based text models. Langedijk et al. (2024) adapt the method to encoder-decoder architectures and also apply it to Whisper. However, their work aims to identify at what encoder layer stable predictions emerge, and focuses on the tasks Whisper was explicitly trained on. We extend this body of work in applying early exiting to interpret emergent speech translation, and provide concrete applications as to how it can aid low-resource speech recognition.

## 3 Experimental Setup

### 3.1 Model

In our experiments, we follow Ma et al. (2025) in examining Whisper (Radford et al., 2022) as a case study. Whisper is an encoder-decoder Transformer model trained on the tasks of language identification, speech recognition, and X -> English speech translation. It is trained on 680,000 hours of data, where 10% is for the task of X -> English translation. Audio is first converted to a mel spectrogram, then passed through convolutional layers to extract features that are then passed to the Transformer's blocks. In this paper, unless otherwise stated, we use the whisper-large-v2 model, which has been shown to exhibit stronger cross-lingual alignment than later variants (Ma et al., 2025). We also perform additional experiments on Open Whisper-style Speech Models (OWSM), which we describe in more detail in section 6.

### 3.2 Dataset

| Language Pair | Full Test Set | Challenge Set |
|---|---|---|
| eng–zho | 427 | 101 |
| fra–zho | 427 | 110 |
| deu–zho | 427 | 94 |
| eng–jpn | 427 | 77 |
| fra–jpn | 427 | 72 |
| deu–jpn | 427 | 67 |

Table 2: Statistics for the FLEURS dataset employed in our study (original test set size vs. our challenge set).

For our speech retrieval experiments, we follow Ma et al. (2025) in employing FLEURS (Conneau et al., 2023), a multilingual parallel speech dataset. They build their test set by combining the dev and test sets of English, French, and German, Mandarin, Japanese, where only utterances common to all five languages are kept. This is followed by a deduplication step to remove duplicate utterances spoken by different speakers. We build our challenge set on top of this procedure, where we employ the spaCy (Honnibal et al., 2020) pipeline to filter out utterances containing proper nouns. We pair together typologically and geographically distant languages to avoid the influence of cognates between related varieties, resulting in six language pairs (Table 2). Proficient in-house speakers of the five languages then manually inspect the resulting dataset, removing utterances which still contain pronunciation similarities between utterance pairs.

As second task, we conduct word-level analyses, which provide a complementary perspective towards our utterance-level insights. For that, we employ the dataset constructed in Choi et al. (2024) who build on the basis of the Multilingual Spoken Words dataset (MSW) (Mazumder et al., 2021) a dataset of cross-lingual synonyms and near-homophones. The MSW dataset consists of 1-second sliced words in the CommonVoice (Ardila et al., 2020) dataset. Choi et al. (2024) obtain cross-lingual synonyms and near-homophones by way of the Open Multilingual Wordnet (OMW) (Bond and Foster, 2013) and Epitran (Mortensen et al., 2018), and create their dataset based on the intersection of languages supported by OMW and Epitran. This results in English, Chinese, Italian, Spanish, Indonesian, Polish, and Swedish. 2,000 word utterances per language are then randomly sampled, resulting in 14,000 total utterances. We refer the reader to Choi et al. (2024) for more details on the dataset.

## 4 Methodology

### 4.1 Cross-Lingual Speech Retrieval

Ma et al. (2025) propose SeqSim to quantify the similarity between two sequences of audio embeddings $X = \{\mathbf{x}_1, \ldots, \mathbf{x}_m\}$ and $Y = \{\mathbf{y}_1, \ldots, \mathbf{y}_n\}$ by measuring how well each frame in one sequence aligns with the most similar frame in the other, where the process is repeated in both directions. The benefit of this measure is that it is simple and allows frames in two audio embeddings to match one another regardless of their position. Formally,
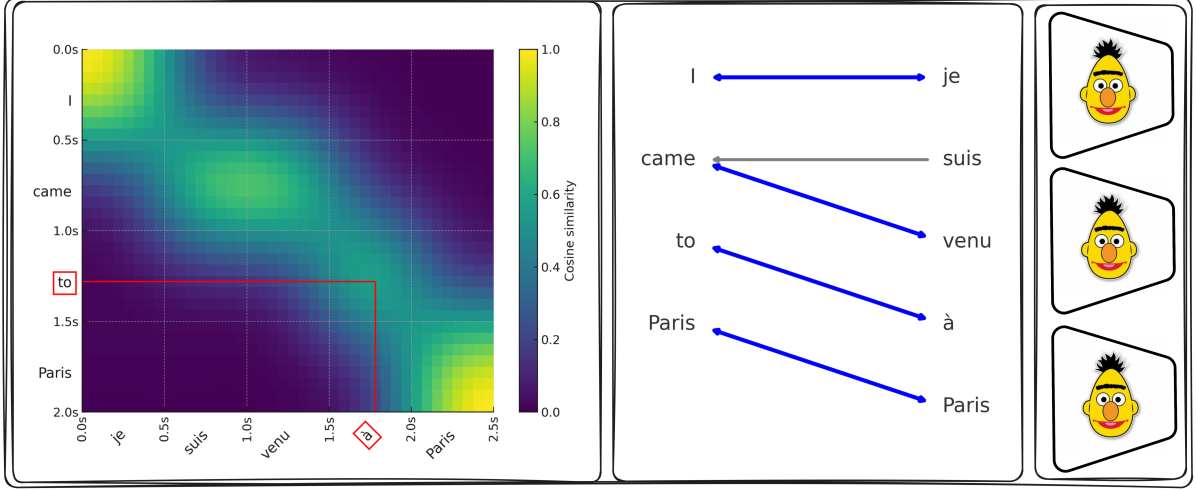
Figure 2: Illustration of our proposed method to determine whether cross-lingual speech retrieval relies on semantic features. Starting with the center frame embeddings obtained in Figure 1, we match each center frame to the most similar frame in the target utterance based on cosine similarity. We repeat this in the reverse direction. We then obtain words the frames belong to by inferred timestamps, and quantify to what degree words that mutually select each other as the most similar are semantically equivalent with a multilingual text encoder.

SeqSim is defined as:

$$\text{Re}_{\text{seq}}(X, Y) = \frac{1}{|X|} \sum_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y} \mathbf{x}^\top \mathbf{y},$$

$$\text{Pr}_{\text{seq}}(X, Y) = \frac{1}{|Y|} \sum_{\mathbf{y} \in Y} \max_{\mathbf{x} \in X} \mathbf{x}^\top \mathbf{y}, \qquad (1)$$

$$\text{SeqSim}(X, Y) = 2 \cdot \frac{\text{Pr}_{\text{seq}} \cdot \text{Re}_{\text{seq}}}{\text{Pr}_{\text{seq}} + \text{Re}_{\text{seq}}}.$$

As shown in Ma et al. (2025), SeqSim outperforms mean pooling and dynamic time warping for spoken translation retrieval. However, it was not designed with interpretability in mind.

To gain a qualitative understanding of what cues SeqSim is relying on, we propose **SeqSimInterp**, a method that returns human-interpretable insights by leveraging text to identify the best matching frames SeqSim relies on for retrieving cross-lingual semantic equivalent utterances. Specifically, building on the finding that the center frame of a spoken utterance retains word identity information (Pasad et al., 2024; Choi et al., 2024), for each word in an utterance we select the center frame $\mathbf{x}_{c(i)} \in X$ (Figure 1), which is determined by word-level timestamps inferred through applying dynamic time warping on cross-attention weights (Zusag et al., 2024). We then identify the frame in $Y$ that has the highest similarity to this center frame. The word whose time boundaries contain the maximally similar frame $\mathbf{y}_{j*}$ is selected as the aligned target word. Similar to SeqSim, this procedure is applied bidirectionally to obtain mutual

alignments. An overview of the full SeqSimInterp pipeline is shown in Figure 2. Assuming that the best matching frames should be semantically similar, if indeed SeqSim is relying on semantic knowledge for retrieval, mutually aligned words should be semantically equivalent. To quantify the semantic alignment captured by SeqSimInterp, we leverage text-based multilingual language model embeddings. For each mutually selected word pair, we compute the cosine similarity between their contextual embeddings produced by LaBSE (Feng et al., 2022), a BERT-based multilingual sentence embedding model. In our experiments, we compare the resulting scores against a random baseline, where we sample random word pairs from the same utterance pairs. We assume that if SeqSim relies on semantic cues for retrieval, the aligned word pairs should consistently score higher than randomly sampled pairs. As such, we perform a paired t-test comparing the similarity scores of word pairs produced by SeqSimInterp against the random baseline. In addition, to account for the fine-grainedness of Whisper's word-level timestamps on Mandarin and Japanese, which do not directly correspond with the word-level granularity of English, French, and German, we merge the timestamp boundaries of the fine-grained word tokens on the basis of language-specific tokenizers for the two languages[1]. This helps to perform the

---

[1]For Mandarin, we use jieba: `github.com/fxsjy/jieba`. For Japanese, we use Sudachi (Takaoka et al., 2018).

word matching process with word-level tokens on a more comparable basis.

To sum up, SeqSimInterp leverages the word-level timestamps inferred through cross-attention weights to produce insights as to what cues SeqSim relies on for retrieving spoken translations. Such interpretability helps lend confidence to findings pertaining to the cross-lingual aligmment of speech foundation models.

## 4.2 Logit Lens

To obtain insights as to how the encoder layers process the input to produce speech translation output in unseen directions, we follow nostalgebraist (2020) in viewing layers of a transformer model as performing incremental updates to latent predictions of the next token. This assumption implies that the hidden states can be decoded to gain insights as to how the input is being processed in said layer. This method has gained widespread usage in text-based decoder-only language models, where it is used to illustrate what the model captures in terms of vocabulary space at each internal layer (Belrose et al., 2023; Din et al., 2024). Langedijk et al. (2024) extend this method to encoder-decoder models, where the intermediate hidden states of the encoder are passed directly to the decoder. In this paper, we apply the method of Langedijk et al. (2024) to speech translation in unseen directions (Peng et al., 2023), with the goal of understanding how the representation changes as the layers move towards a language-agnostic space.

## 4.3 Word-Level Analyses

In addition to utterance-level analyses, we follow Choi et al. (2024) in measuring the cosine similarity between cross-lingual synonyms and near homophones for a given word. We include lower and upper similarity estimates, i.e., a random sample of a word is employed as a lower bound of similarity, and an occurrence of the same word in a different audio sample is used as an upper bound of similarity. For instance, if indeed the encoder of a speech foundation model encodes semantics, then the cosine similarity of *dog* in English should be high when measured against *chien* in French; and if the encoder encodes phonetic similarity, the word *dog* in English should have similarity to the near-homophone French word *dague*, meaning *dagger* in English, where the upper and lower bounds help to contextualize the degree of similarity observed. In addition, we follow Choi et al. (2024) in repeat-

| Full FLEURS test | | | | | |
|---|---|---|---|---|---|
| | **en** | **fr** | **de** | **zh** | **ja** |
| en | — | 80.80% | 78.45% | 42.15% | 47.54% |
| fr | 74.47% | — | 65.57% | 43.79% | 48.71% |
| de | 71.90% | 60.89% | — | 45.20% | 52.93% |
| zh | 29.98% | 23.19% | 19.91% | — | 49.41% |
| ja | 18.97% | 18.50% | 14.99% | 34.66% | — |
| FLEURS challenge subset | | | | | |
| | **en** | **fr** | **de** | **zh** | **ja** |
| en | — | — | — | 41.58% | 22.08% |
| fr | — | — | — | 40.00% | 23.61% |
| de | — | — | — | 42.55% | 46.27% |
| zh | 24.75% | 25.45% | 10.64% | — | — |
| ja | 7.79% | 18.06% | 7.46% | — | — |

Table 3: Speech retrieval R@1 for Whisper-large-v2 on the full FLEURS test set and on the FLEURS challenge subset (no pronunciation-similar pairs). Rows are source languages, columns are target languages.

| Full FLEURS test | | | | | |
|---|---|---|---|---|---|
| | **en** | **fr** | **de** | **zh** | **ja** |
| en | — | 89.5 | 83.25 | 51.79 | 44.28 |
| fr | 86.1 | — | 66.59 | 42.82 | 36.58 |
| de | 88.0 | 72.26 | — | 47.43 | 39.19 |
| zh | 52.29 | 44.38 | 46.1 | — | 54.21 |
| ja | 42.94 | 36.6 | 36.6 | 55.69 | — |
| FLEURS challenge subset | | | | | |
| | **en** | **fr** | **de** | **zh** | **ja** |
| en | — | — | — | 24.16 | 14.92 |
| fr | — | — | — | 21.13 | 13.32 |
| de | — | — | — | 21.24 | 12.5 |
| zh | 25.07 | 23.57 | 21.8 | — | — |
| ja | 15.51 | 12.27 | 11.86 | — | — |

Table 4: t-statistics of the paired t-tests. All results are statistically significant with $p < 0.001$.

ing the process five times, and show the 95% confidence interval. We also follow Choi et al. (2024) in subtracting the lower bound random values from all other values computed to help readability of the plot.

## 5 Results

### 5.1 Speech Retrieval

Table 3 shows our main results for speech retrieval on the full test set of FLEURS and our decon-founded challenge set. In all experiments, retrieval performance is measured using Recall@1 (R@1).

**Spoken retrieval in Whisper uses semantics.** We observe that in all language-pair setups, there is a noticable drop in performance on the challenge set, although retaining still a high accuracy
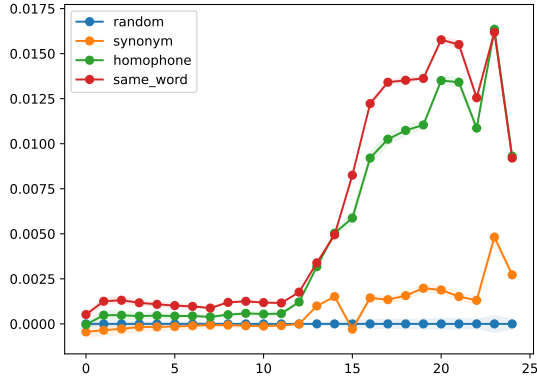
Figure 3: Word-level analyses on whisper-medium. x-axis is layer count, y-axis is cosine similarity.

overall. For instance, Mandarin->English drops from 29.98% to 24.75%, and Japanese->English drops from 18.97% to 7.79%. This indicates that although phonetic cues are indeed useful for retrieval, the models examined are able to rely mostly on semantic cues.

To confirm the validity of our speech retrieval findings, we look into whether the speech retrieval metric indeed relies on semantic cues. Table 4 shows the paired t-test results of applying SeqSim-Interp on the obtained word pairs and the random word pairs for all language pairs. To control for the family-wise error rate due to multiple comparisons, we applied the Holm–Bonferroni method (Holm, 1979) to adjust the p-values. We find that for all language pairs examined, the paired t-tests are statistically significant, with $p < 0.001$. This suggests spoken translation retrieval to indeed rely on semantic cues, as the best matching center frames between utterances of two given languages are also semantically similar.

**Word-Level Results** Figure 3 shows our results for our word-level experiments. Following Choi et al. (2024), the cosine similarity with the random samples is subtracted from all other lines to improve clearness of the plot. We observe that, corresponding to what Choi et al. (2024) observe for self-supervised speech models, Whisper[2] also exhibits higher similarity between cross-lingual near-homophones than between cross-lingual synonyms. We hypothesize that the weaker semantic

capacity than what we observe on utterance-level data to be due to word-level data being too short in duration for a multilingual model to obtain the semantics of the given input. For instance, without prior context, the sound *one* has many possible meanings across all the languages in the training data, but the increased context of *one man in Paris* constrains the possible meaning of the sound *one* considerably. We note that our procedure in SeqSimInterp also extracts word-level correspondences, but the word-level features obtained from our procedure come from *utterance-level context*, which provides the model more information to obtain the semantics of the given word. Correspondingly, this raises the question of whether the weaker degree of cross-lingual semantic knowledge in self-supervised speech models observed in Choi et al. (2024) is due to word-level supervision. Our results suggest it to be worthwhile to revisit the degree of semantic knowledge in self-supervised speech models with utterance-level data. We leave such an investigation to future work.

| | |
|---|---|
| **Layer 32** | bambini sviluppano un'conoscenza di stereotipi di razza e di **racia** abbastanza giovani e questi stereotipi di razza **affettano** il comportamento |
| **Layer 30** | I bambini sviluppano **un'awareness** di stereotipi **raci** e **raciali** abbastanza giovani e questi stereotipi **raciali affettano** il comportamento |
| **Layer 28** | **Cildren** develop an awareness of racial stereotypes quite young and these racial stereotypes affect behavior |
| **Source** | Children develop an awareness of race and racial stereotypes quite young and these racial stereotypes affect behavior |

Table 5: Early exiting outputs for English to Italian speech translation, where 32 is the final layer. Words in **blue** indicate **phonetic** error and words in **red** indicate **semantic** error.

## 5.2 Logit Lens

### 5.2.1 Speech Translation

Table 5 show an example on early exiting the encoder for EN-IT speech translation, showing how translation emerges throughout the decoder layers. Although the translations produced are largely coherent, we observe that translations produced in earlier layers exhibit errors that are *phonetically similar to the source language*. This is visible in the words highlighted in blue and red, which indicate respectively phonetic and semantic errors. In our example, the translation *affettano* for instance bears phonetic similarity to the source word *affect*, yet

---

[2]The figure shows Whisper-medium, although all model sizes show similar trends of stronger similarity to near-homophones than to synonyms. We show our results for other model sizes in the appendix.

its meaning (*affettare*) in Italian means *to cut*, and constitutes a translational error that arguably stems from a sound-based translation. A similar pattern is observed in the translation of the words *race* and *racial*. For instance, *race* is mistranslated as *racia*, the latter being a non-existent or malformed Italian word that nonetheless closely resembles the English pronunciation. Likewise, *racial* becomes *raciali*, which is morphologically plausible in Italian, but earlier layers also produce variants like *raci* that lack semantic grounding. These outputs illustrate a breakdown in meaning preservation, where the model opts for phonetically close forms, even if they are invalid or inappropriate in context, rather than true semantic equivalents.

### 5.2.2 Unseen Source Languages

Building on the early exiting observations from speech translation in the previous section, we hypothesize that intermediate layer predictions may be particularly valuable for low-resource languages, especially those with phonemically transparent orthographies and typological proximity to high-resource languages already supported by Whisper.

To test this hypothesis, our setup allows us to quantitatively test whether transcriptions from earlier layers retain more phonetic information and are thus better suited for zero-shot adaptation. To evaluate this, we adopt the six low-resource languages examined by Ma et al. (2025) for assessing Whisper's adaptability to unseen source languages—matching their target language selection except for Cebuano, where we replace English with Tagalog due to stronger typological relatedness. We compute Word Error Rate (WER) and Character Error Rate (CER) at different layer depths using early exit decoding. Additionally, we include Javanese, a language with a transparent orthography, to test whether early-layer outputs can yield usable transcriptions without explicit training data.

The results on the development set (Table 9) indicate a consistent trend: early layer predictions yield lower WER and CER across all evaluated languages, supporting our hypothesis that intermediate representations encode more phonetic information. Based on these findings, we selected the best-performing layer per language for subsequent evaluation on the test set (Table 6). The results support our hypothesis. Taking Javanese as example, while the WER fluctuates across layers, CER consistently decreases in earlier layers, reaching its minimum at layer 3 (34.9%). This 10-point

**CER results**

| Language | Best Layer | Last Layer | CER Decrease |
|---|---|---|---|
| ky | **100.2**% (24) | 157.0% | 56.8% |
| ceb | **16.2**% (29) | 36.7% | 20.5% |
| ga | **75.2**% (29) | 89.2% | 14% |
| jv | **24.7**% (29) | 35.0% | 10.3% |
| kea | **35.0**% (29) | 35.7% | 0.7% |
| ast | **17.5**% (31) | 17.6% | 0.1% |
| ckb | **50.9**% (32) | **50.9**% | 0.0% |

**WER results**

| Language | Best Layer | Last Layer | WER Decrease |
|---|---|---|---|
| ky | **112.0**% (24) | 210.0% | 98% |
| ga | **99.9**% (24) | 121.4% | 21.5% |
| ceb | **51.7**% (29) | 67.4% | 15.7% |
| jv | **75.3**% (29) | 82.5% | 7.2% |
| ckb | **1.05**% (24) | 1.13% | 0.8% |
| kea | **92.2**% (29) | 92.6% | 0.4% |
| ast | **63.7**% (31) | 64.0% | 0.3% |

Table 6: Early exiting results on FLEURS test set of low-resource languages. Numbers in parentheses indicate best layer as selected on dev set.

**OWSM v3.1 Small**

| Lang | en | fr | de | zh | ja |
|---|---|---|---|---|---|
| en | — | 85.71% | 82.9% | 27.17% | 38.64% |
| fr | 74.47% | — | 48.48% | 11.71% | 15.2% |
| de | 68.38% | 44.5% | — | 13.35% | 24.12% |
| zh | 12.65% | 1.41% | 4.92% | — | 30.91% |
| ja | 7.96% | 1.64% | 3.51% | 9.37% | — |

**OWSM v3.1 Small Low-Restriction**

| Lang | en | fr | de | zh | ja |
|---|---|---|---|---|---|
| en | — | 67.21% | 64.4% | 1.64% | 4.22% |
| fr | 64.64% | — | 44.5% | 2.34% | 0.7% |
| de | 54.8% | 41.45% | — | 2.11% | 0.7% |
| zh | 1.17% | 1.41% | 2.11% | — | 0.94% |
| ja | 2.81% | 2.11% | 1.87% | 1.64% | — |

Table 7: Speech retrieval scores for OWSM v3.1 small models on the whole FLEURS test set.

absolute drop in CER between the early and third layers reinforces the claim that intermediate representations encode more segmental-level phonetic information than deeper layers.

## 6 Discussion

Based on our results, we revisit below our research question of whether spoken translation retrieval relies on semantic features in the models we examine. In addition, prompted by our results, we ask the follow-up question of whether such an alignment is also observed to the same degree in a model trained only on the task of speech recognition, OWSM v3.1 small low-restriction (Peng et al., 2024), with no speech translation supervision.

| OWSM v3.1 Small | | | | | |
| --- | --- | --- | --- | --- | --- |
| Lang | en | fr | de | zh | ja |
| en | — | — | — | 26.73% | 7.79% |
| fr | — | — | — | 5.45% | 1.39% |
| de | | — | — | 8.51% | 1.49% |
| zh | 7.92% | 1.82% | 5.32% | — | — |
| ja | 2.6% | 0.0% | 1.49% | — | — |

| OWSM v3.1 Small Low-Restriction | | | | | |
| --- | --- | --- | --- | --- | --- |
| Lang | en | fr | de | zh | ja |
| en | — | — | — | 0.0% | 0.0% |
| fr | — | — | — | 0.91% | 1.39% |
| de | — | — | — | 0.0% | 0.0% |
| zh | 0.99% | 0.0% | 0.0% | — | — |
| ja | 1.3% | 0.0% | 0.0% | — | — |

Table 8: Speech retrieval scores for OWSM v3.1 models on the FLEURS challenge subset.

## Do speech foundation models make use of semantic cues for retrieving spoken translations?

While prior work has leveraged high retrieval scores to suggest shared semantic representations (Ma et al., 2025), our controlled challenge set demonstrates that such scores can be inflated by pronunciation-level cues like cognates and proper nouns. Nevertheless, as observed in Table 3, even after systematically removing these cues, retrieval accuracy—though reduced—remains substantially above random, suggesting that semantic alignment does indeed exist, but it coexists with and is partly entangled in phonetic similarity. SeqSimInterp provides further support for this interpretation. By isolating the token-level alignments contributing to retrieval scores, we show that many of the matched words across languages are semantically equivalent, confirmed via cross-lingual embedding similarity. This suggests cross-lingual alignment in speech foundation models such as Whisper to align both phonetically and semantically.

## Is cross-lingual alignment in speech foundation models induced by speech translation?

To answer the question of whether the cross-lingual alignment we observe is induced by the Any->En speech translation task in Whisper, or whether the alignment arises even only with multilingual speech recognition, we repeat our spoken translation retrieval experiments on Open Whisper-Style Speech Models (OWSMs) (Peng et al., 2024). OWSMs provide an ideal testing ground for this hypothesis due to the open source nature of its train-

ing data; their adherence to Whisper-style training specifications; and provide models of various sizes and setups. OWSM includes two models of comparable size and architecture, OWSM v3.1 Small and Small Low-Restriction, of which one the latter is not trained on speech translation. As such, we repeat our spoken translation retrieval experiments on the two models. Table 7 and Table 8 detail our results. We observe that on both the full test set and the challenge set, the model not trained on speech translation (OWSM Small Low-Restriction) exhibits significantly lower retrieval capabilities particularly on typologically distant language pairs, which deconfounds for pronunciation-level shortcuts by way of cognates, dropping to 0.0% in many cases. This suggests that speech translation is indeed a strong contributor for semantic cross-lingual alignment capabilities in speech models, which arises to a much less degree when trained only on multilingual speech recognition.

## 7 Conclusion

In this work, we revisit the question of whether cross-lingual alignment in multilingual speech foundation models is truly semantic. Through a series of controlled experiments, we demonstrated that although phonetic cues—such as cognates and proper nouns—can facilitate speech-to-speech retrieval, they are not solely responsible for the observed cross-lingual alignment. To probe the nature of this alignment, we introduce SeqSimInterp, a method for interpreting retrieval decisions at the word level. Our analysis revealed that semantically equivalent words contribute significantly to retrieval scores when used at utterance level. Additionally, we extended our evaluation across models with and without speech translation training. We found that speech foundation models trained *additionally* on speech translation tasks significantly enhances semantic cross-lingual alignment in speech, underscoring the role of supervised multilingual multi-task signals in pre-training to aid shaping useful semantic representations. Finally, by early exiting the encoder, we observed that earlier encoder layers preserve more phonetic detail—information that can be harnessed for zero-shot adaptation to low-resource, phonetically transparent languages.

## 8 Limitations

A limitation of our work is that SeqSimInterp is only able to capture word-level matches between

languages, whereas empirically cross-lingual semantic mappings may be one-to-many. In addition, our study focuses only on Whisper and Open Whisper-style Speech Models, but it may be interesting compare whether our findings hold also for speech foundation models of different architectures, such as SeamlessM4T (Communication et al., 2023).

# 9 Ethical Considerations

The authors acknowledge the usage of ChatGPT as an assistant tool in part of the source code's development, in assisting the creation of Figure 1 and Figure 2, and in enhancing the coherence of parts of the manuscript.

# 10 Acknowledgements

# References

Badr M. Abdullah, Mohammed Maqsood Shaik, and Dietrich Klakow. 2024. Wave to interlingua: Analyzing representations of multilingual speech transformers for spoken language translation. In *Interspeech 2024*, pages 362–366.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus. *Preprint*, arXiv:1912.06670.

Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. Eliciting latent predictions from transformers with the tuned lens. *Preprint*, arXiv:2303.08112.

Francis Bond and Ryan Foster. 2013. Linking and extending an open multilingual Wordnet. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1352–1362, Sofia, Bulgaria. Association for Computational Linguistics.

Kwanghee Choi, Ankita Pasad, Tomohiko Nakamura, Satoru Fukayama, Karen Livescu, and Shinji Watanabe. 2024. Self-supervised speech representations are more phonetic than semantic. In *Interspeech 2024*, pages 4578–4582.

Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, Christopher Klaiber, Pengwei Li, Daniel Licht, Jean Maillard, Alice Rakotoarison, Kaushik Ram Sadagopan, Guillaume Wenzek, Ethan Ye, and 49 others. 2023. Seamlessm4t: Massively multilingual & multimodal machine translation. *Preprint*, arXiv:2308.11596.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

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 *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 798–805.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Alexander Yom Din, Taelin Karidi, Leshem Choshen, and Mor Geva. 2024. Jump to conclusions: Shortcutting transformers with linear transformations. *Preprint*, arXiv:2303.09435.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Danny Halawi, Jean-Stanislas Denain, and Jacob Steinhardt. 2024. Overthinking the truth: Understanding how language models process false demonstrations. *Preprint*, arXiv:2307.09476.

Sture Holm. 1979. A simple sequentially rejective multiple test procedure. *Scandinavian journal of statistics*, pages 65–70.

Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python.

Yigitcan Kaya, Sanghyun Hong, and Tudor Dumitras. 2019. Shallow-deep networks: Understanding and mitigating network overthinking. *Preprint*, arXiv:1810.07052.

Guillaume Lample and Alexis Conneau. 2019. Cross-lingual language model pretraining. *Preprint*, arXiv:1901.07291.

Anna Langedijk, Hosein Mohebbi, Gabriele Sarti, Willem Zuidema, and Jaap Jumelet. 2024. DecoderLens: Layerwise interpretation of encoder-decoder transformers. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4764–4780, Mexico City, Mexico. Association for Computational Linguistics.

Rao Ma, Mengjie Qian, Yassir Fathullah, Siyuan Tang, Mark Gales, and Kate Knill. 2025. Cross-lingual transfer learning for speech translation. *Preprint*, arXiv:2407.01130.

Mark Mazumder, Sharad Chitlangia, Colby Banbury, Yiping Kang, Juan Ciro, Keith Achorn, Daniel Galvez, Mark Sabini, Peter Mattson, David Kanter, Greg Diamos, Pete Warden, Josh Meyer, and Vijay Janapa Reddi. 2021. Multilingual spoken words corpus. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.

Hosein Mohebbi, Grzegorz Chrupała, Willem Zuidema, Afra Alishahi, and Ivan Titov. 2024. Disentangling textual and acoustic features of neural speech representations. *Preprint*, arXiv:2410.03037.

David R. Mortensen, Siddharth Dalmia, and Patrick Littell. 2018. Epitran: Precision G2P for many languages. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

nostalgebraist. 2020. Interpreting gpt: The logit lens. https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens. Accessed: 2025-05-19.

Ankita Pasad, Chung-Ming Chien, Shane Settle, and Karen Livescu. 2024. What do self-supervised speech models know about words? *Transactions of the Association for Computational Linguistics*, 12:372–391.

Puyuan Peng, Brian Yan, Shinji Watanabe, and David Harwath. 2023. Prompting the hidden talent of web-scale speech models for zero-shot task generalization. In *Interspeech 2023*, pages 396–400.

Yifan Peng, Jinchuan Tian, William Chen, Siddhant Arora, Brian Yan, Yui Sudo, Muhammad Shakeel, Kwanghee Choi, Jiatong Shi, Xuankai Chang, Jee weon Jung, and Shinji Watanabe. 2024. Owsm v3.1: Better and faster open whisper-style speech models based on e-branchformer. *Preprint*, arXiv:2401.16658.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. *arXiv preprint*.

Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. 2025. A practical review of mechanistic interpretability for transformer-based language models. *Preprint*, arXiv:2407.02646.

Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Q. Tran, Yi Tay, and Donald Metzler. 2022. Confident adaptive language modeling. *Preprint*, arXiv:2207.07061.

Kazuma Takaoka, Sorami Hisamoto, Noriko Kawahara, Miho Sakamoto, Yoshitaka Uchida, and Yuji Matsumoto. 2018. Sudachi: a Japanese tokenizer for business. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Chih-Kai Yang, Kuan-Po Huang, Ke-Han Lu, Chun-Yi Kuan, Chi-Yuan Hsiao, and Hung-Yi Lee. 2024. Investigating zero-shot generalizability on mandarin-english code-switched asr and speech-to-text translation of recent foundation models with self-supervision and weak supervision. In *2024 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW)*, pages 540–544.

Mario Zusag, Laurin Wagner, and Bernhad Thallinger. 2024. Crisperwhisper: Accurate timestamps on verbatim speech transcriptions. In *Interspeech 2024*, pages 1265–1269.
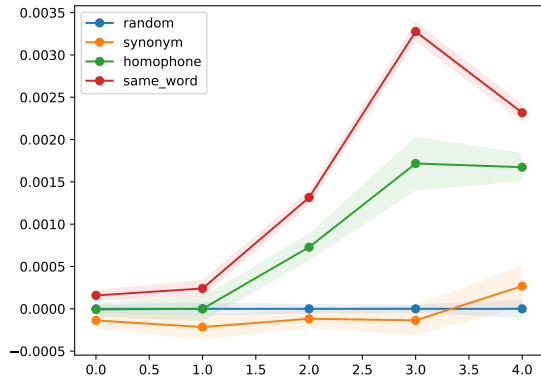
# A  Appendix



Figure 4: Word-level analyses on whisper-tiny. x-axis is layer count, y-axis is cosine similarity.
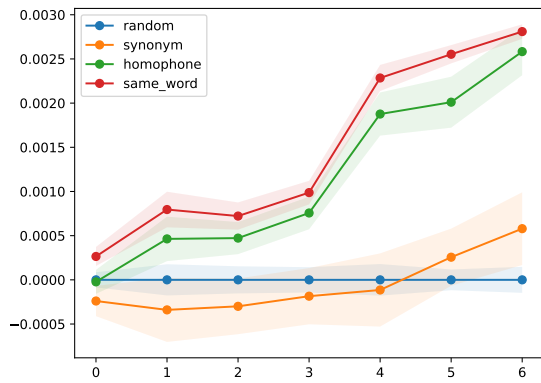


Figure 7: Word-level analyses on whisper-large. x-axis is layer count, y-axis is cosine similarity.



Figure 5: Word-level analyses on whisper-base. x-axis is layer count, y-axis is cosine similarity.
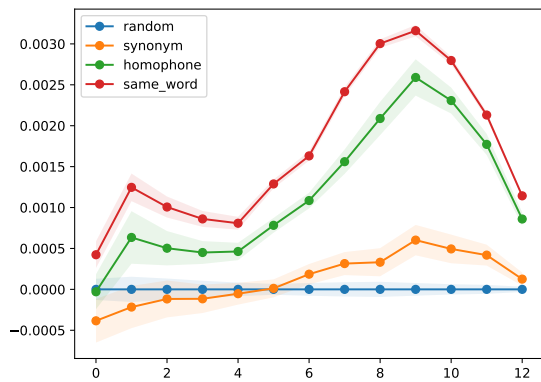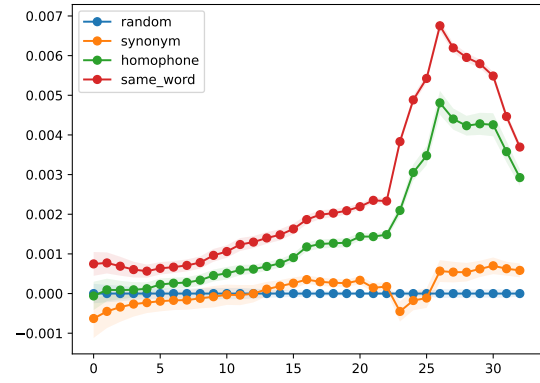


Figure 8: Word-level analyses on whisper-large-v2. x-axis is layer count, y-axis is cosine similarity.



Figure 6: Word-level analyses on whisper-small. x-axis is layer count, y-axis is cosine similarity.



Figure 9: Word-level analyses on whisper-large-v3. x-axis is layer count, y-axis is cosine similarity.
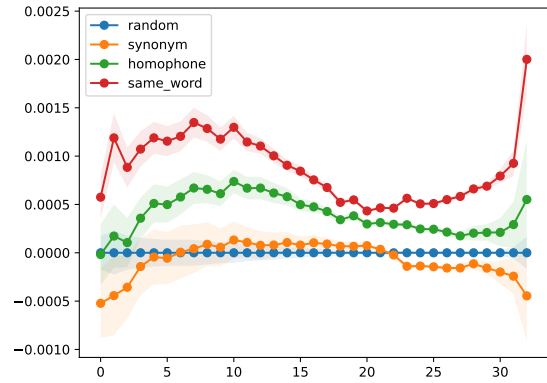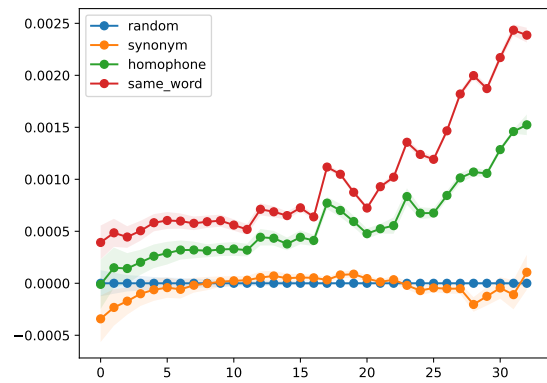
| Layer | Asturian | Cebuano | Irish | Kabuverdianu | Kyrgyz | Sorani Kurdish | Javanese |
|---|---|---|---|---|---|---|---|
| 0 | 17.1 / 61.5 | 12.6 / **46.1** | 78.5 / 115.6 | 34.1 / 91.4 | 122.8 / 176.6 | 60.1 / 129.3 | 46.4 / 90.7 |
| 1 | **17.0 / 61.3** | **12.5** / 46.7 | 71.8 / 108.2 | 34.2 / 91.7 | 127.7 / 204.9 | 63.5 / 139.6 | 45.2 / 94.7 |
| 2 | 17.8 / 61.7 | 12.6 / 48.1 | **70.3** / 109.2 | 34.3 / 91.2 | 146.3 / 232.1 | **56.1** / 127.5 | 55.3 / 97.5 |
| 3 | 17.5 / 62.2 | 12.6 / 48.0 | 72.6 / 107.6 | **32.9 / 89.5** | 103.2 / 150.5 | 63.9 / 136.7 | **34.9 / 86.5** |
| 4 | 20.1 / 62.9 | 22.3 / 51.7 | 116.0 / 152.6 | 35.8 / **88.7** | 159.1 / 205.5 | 114.3 / 179.1 | 42.1 / 90.4 |
| 5 | 37.0 / 77.4 | 46.5 / 69.4 | 128.4 / 138.4 | 55.4 / 92.7 | 145.9 / 189.8 | 115.7 / 139.0 | 73.8 / 98.5 |
| 6 | 84.5 / 103.4 | 82.3 / 98.2 | 100.9 / 100.2 | 83.5 / 99.1 | 127.3 / 137.3 | 96.9 / **101.0** | 109.8 / 143.3 |
| 7 | 99.6 / 103.5 | 94.6 / 98.5 | 104.3 / 108.9 | 106.7 / 119.5 | 103.2 / 120.1 | 97.5 / 102.5 | 106.6 / 143.8 |
| 8 | 101.4 / 107.0 | 99.2 / 99.8 | 99.2 / **99.9** | 99.4 / 102.1 | **99.3 / 104.6** | 101.1 / 103.9 | 98.2 / 100.4 |

Table 9: Layer-wise CER and WER (in %) for FLEURS dev set low-resource languages