# M-SPEECHCLIP: LEVERAGING LARGE-SCALE, PRE-TRAINED MODELS FOR MULTILINGUAL SPEECH TO IMAGE RETRIEVAL

*Layne Berry[1], Yi-Jen Shih[2], Hsuan-Fu Wang[2], Heng-Jui Chang[3], Hung-yi Lee[2], and David Harwath[1]*

[1]University of Texas at Austin
[2]National Taiwan University
[3]MIT CSAIL

## ABSTRACT

This work investigates the use of large-scale, English-only pre-trained models (CLIP and HuBERT) for multilingual image-speech retrieval. For non-English image-speech retrieval, we outperform the current state-of-the-art performance by a wide margin both when training separate models for each language, and with a single model which processes speech in all three languages. We identify key differences in model behavior and performance between English and non-English settings, attributable to the English-only pre-training of CLIP and HuBERT, and investigate how fine-tuning the pre-trained models impacts these differences. Finally, we show that our models can be used for mono- and cross-lingual speech-text retrieval and cross-lingual speech-speech retrieval, despite never having seen any parallel speech-text or speech-speech data during training.

***Index Terms***— visually-grounded speech, multimodal speech processing, multilingual speech processing, self-supervised learning

## 1. INTRODUCTION

Language is more than a probability distribution over words or phonemes–it is produced in context to express semantic information. The task of language-image retrieval targets semantic understanding by asking models to connect utterances to contexts in a different modality. This is especially important for low-resource and unwritten languages, for which limited data exists for training speech recognition and natural language understanding systems. The recently-proposed CLIP [1] model was pre-trained on an unprecedented amount of parallel English image-text data to encode each modality in a shared semantic embedding space. Speech-CLIP [2] was then proposed to map HuBERT representations of English speech into the same embedding space, achieving state-of-the-art performance on English image-speech retrieval and enabling zero-shot speech-text retrieval. Here, we investigate the value of English-only pre-training for non-English speech understanding by applying the SpeechCLIP [2] model to non-English image-speech retrieval. We find that our models beat the prior state-of-the-art for non-English image-speech retrieval by a wide margin.

We next train multilingual models which can take input text in any of the three languages investigated here. We experiment with scaling up these models, and achieve further gains. As with Speech-CLIP, we show that our models can perform zero-shot transfer to English speech-text retrieval, even outperforming image-text retrieval with the image embeddings M-SpeechCLIP used as labels during training. Finally, we consider the challenging settings of zero-shot

transfer to cross-lingual speech-text and speech-speech retrieval, setting strong baselines for the former and outperforming prior work for the latter.

This work demonstrates that large-scale pre-training is highly effective, even for tasks where both language and modality differ. We set a new state-of-the-art for non-English image-speech retrieval, and enable both speech-text and speech-speech retrieval in a cross-lingual setting without any parallel speech from different languages, parallel speech and text, or non-English text at all.

## 2. RELATED WORK

CLIP [1] is an image-text alignment model trained on 400 million web-scraped image and English caption pairs–a private dataset 4000x larger than the popular MS-COCO [3] and Visual Genome [4] and less noisy than the 4x smaller YFCC100M [5]. CLIP-Large consists of a ViT-L/14 [6] image encoder and a GPT-based [7] text encoder which map their respective input modalities into a shared embedding space learned with a contrastive loss. This embedding space has been shown to transfer well to tasks other than English image-text alignment. M-CLIP[8] uses the CLIP text encoder with English inputs as a teacher model to learn non-English and multilingual text encoders into the CLIP embedding space, and beat state-of-the-art image-text retrieval performance on the XTD [9] dataset of MS-COCO caption translations in 11 languages by an average of over 12% absolute R@10. SpeechCLIP [2] freezes the CLIP image encoder and learns to map English spoken captions into the CLIP embedding space, beating state-of-the-art image-speech retrieval [10] on the SpokenCOCO [11] and Flickr8k [12] datasets.

Non-English image-speech retrieval has most frequently been evaluated on the Places [13] dataset. [14] collected 400k spontaneous English spoken captions for Places images. [15] collected spontaneous Hindi spoken captions for 100k of these images, and [16] collected spontaneous Japanese spoken captions for the same subset of images. [17], [18], and [19] also investigate non-English image-speech retrieval, with [19] using transfer learning from English pre-training on HowTo100M to beat prior state-of-the-art by an average absolute R@10 of 5.7% for Japanese and 8.1% for Hindi.

[2] showed that encoding English speech into the CLIP embedding space allows speech-text retrieval to be learned without any speech-text pairs. Visual grounding has also been used to learn to retrieve speech given a text keyword without seeing any speech-text pairs during training. [20, 21] investigate English speech retrieval given an English query word, while [22] learns to retrieve English speech given one-word text query in German. [16] used a combination of cross-lingual and cross-modal loss to learn monolingual encoders that can be used for cross-lingual speech-to-speech retrieval.
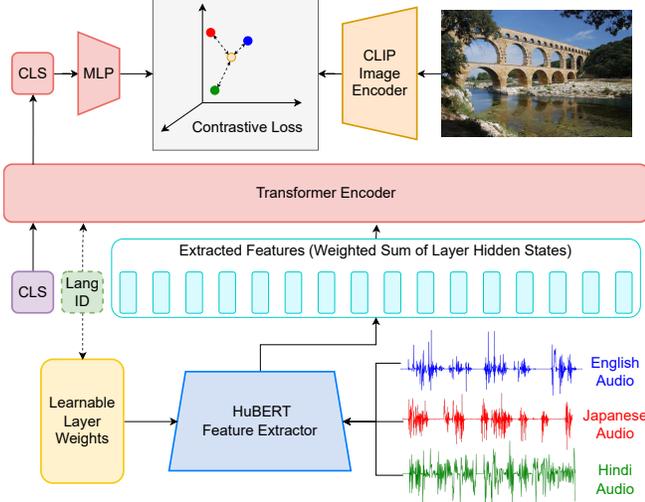
**Fig. 1**: Diagram of the Multilingual SpeechCLIP Model

## 3. MODEL DESIGN

### 3.1. Monolingual Models

Our model architecture is based on that of Parallel SpeechCLIP, introduced in [2]. Given a batch of $n$ image-caption pairs $(I_i, S_i)$ for $i = 1, ..., n$, we use a frozen CLIP [1] image encoder to generate an embedding $v_i^I$ for each image $I_i$ and train a speech encoder to produce vectors $v_j^S$ such that $v_i^I$ and $v_j^S$ are similar for $i = j$ and dissimilar for $i \neq j$. Our speech encoder uses a frozen pretrained HuBERT [23] model to extract frame-level audio features $f_1, ..., f_k$, where $k$ is the number of frames in the audio at 50Hz. We learn weights $\mathbf{w} = (w^1, ..., w^l)$ for each of the $l$ HuBERT layers, which we use to compute the representation for each audio frame $f_t$ from the hidden states $\mathbf{h}_t = (h_t^1, ..., h_t^l)$ at that frame as $f_t = \sum_{i=1}^{l} w^i h_t^i$. We then append a learnable [CLS] token to the beginning of the sequence and pass it to a trainable Transformer Encoder [24]. The hidden state of the [CLS] token at the last layer of the Transformer Encoder on input $S_i$ is projected to the same shape as the target $v_i^I$ and used as the output encoding $v_i^S$. We use Masked Margin Softmax [25] as our contrastive loss function in both retrieval directions, so that for each batch of $B$ speech encodings $\mathbf{v}^S = (v_1^S, ..., v_B^S)$ and $B$ image encodings $\mathbf{v}^I = (v_1^I, ..., v_B^I)$, our total loss is $\mathcal{L}(\mathbf{v}^S, \mathbf{v}^I) = L_{MMS}(\mathbf{v}^S, \mathbf{v}^I) + L_{MMS}(\mathbf{v}^I, \mathbf{v}^S)$, where $L_{MMS}$ with margin $\delta$ is defined as:

$$L_{MMS}(\mathbf{x}, \mathbf{y}) = -\frac{1}{B} \sum_{i=1}^{B} \left( \log \frac{e^{x_i \cdot y_i - \delta}}{e^{x_i \cdot y_i - \delta} + \sum_{j \neq i}^{B} e^{x_i \cdot y_j}} \right) \quad (1)$$

### 3.2. Multilingual Models

We consider two possible settings: a language-agnostic model identical to the monolingual models but trained on inputs from multiple languages, and a language-aware model which learns a [LangID] token for each language which is appended to the audio feature sequence after the [CLS] token. The language-aware model also learns a separate set of layer weights $\mathbf{w}^C$ for each language $C$.

When training multilingual models, individual batches may contain speech from just one language at a time or a mix of speech from

all languages. Since the other pairs in a batch are used as negatives in the contrastive loss, this choice determines whether distractors are all in the same language or may be in any language.

Finally, we consider introducing cross-lingual loss terms modeled after [16]. At each iteration, we randomly select two languages $C$ and $D$ and compute embeddings of the image ($\mathbf{v}^I$) and the spoken caption for that image in each selected language ($\mathbf{v}^C$ and $\mathbf{v}^D$). We then compute our total loss using six contrastive terms:

$$
\begin{aligned}
\mathcal{L}(\mathbf{v}^I, \mathbf{v}^C, \mathbf{v}^D) = & L_{MMS}(\mathbf{v}^I, \mathbf{v}^C) + L_{MMS}(\mathbf{v}^C, \mathbf{v}^I) + \\
& L_{MMS}(\mathbf{v}^I, \mathbf{v}^D) + L_{MMS}(\mathbf{v}^D, \mathbf{v}^I) + \\
& L_{MMS}(\mathbf{v}^C, \mathbf{v}^D) + L_{MMS}(\mathbf{v}^D, \mathbf{v}^C)
\end{aligned}
\quad (2)
$$

## 4. EXPERIMENTS

We train and evaluate our models on Places100k, the 100k image subset of Places for which spontaneous spoken captions are available in three languages (English, Hindi, and Japanese). All of our monolingual models are trained on 2 V100 GPUs, as are our four multilingual models ablating batch composition and the use of a [LangID] token. We empirically select 1 layer and 8 attention heads for our Transformer Encoders. Monolingual models are trained for 25 epochs and take 1 day to finish training, while multilingual models (which see 3x more training data per epoch) are trained for 15 epochs and take 1.5 days. The batch size used for these experiments is 256, and the feature extractor is HuBERT-Large.

We additionally conduct two experiments on 8 V100 GPUs with a batch size of 72, in order to test more computationally expensive ablations to our multilingual model. These experiments take longer to converge, so we run them for 25 epochs (∼3 days). The first, which we refer to as "Multi+TrainFeat", unfreezes the HuBERT-Large feature extractor, allowing it to be fine-tuned on our multilingual input data. For the second, which we call "Multi+TrainFeat+XLL", we include cross-lingual loss terms. Since this doubles the number of gradients which must be stored at each training step, we use a trainable HuBERT-Base as our feature extractor for this experiment. This allows the amount of compute required by the "Multi+TrainFeat" and the "Multi+TrainFeat+XLL" models to be comparable, and the same batch size, number of GPUs, and training time to be used for both experiments.

Our learning rates increase linearly from zero to a peak over the first 10% of iterations and decrease linearly back to zero over the remaining 90% of iterations. The peak learning rate was set empirically to 0.001 for English and Japanese monolingual models and 0.0005 for Hindi monolingual and all multilingual models. The MMS margin $\delta$ is set to 0.001 and spoken captions are zero-padded or truncated as necessary to 15s.

## 5. PERFORMANCE AND ANALYSIS

### 5.1. Monolingual Image-Speech Retrieval

We evaluate the quality of the learned speech encoders by measuring their retrieval performance on the Places100k test set of 1000 image-speech pairs. Following prior works [15, 16, 17, 19], we report Recall@$k$ for $k = 1, 5, 10$ in Table 1. Our models' improvement over prior work is significantly larger for English than other languages. Still, for Hindi and Japanese, our models achieve an average gain of over 7% absolute R@10 over the state-of-the-art [19]. Lower recall for Hindi than English and Japanese is in line with prior work; lower recall for Japanese than English, however, contrasts the

findings of [16] and [17]. This divergence from the trends of prior work indicates that English-only pre-training biases our models in favor of English speech.

## 5.2. Multilingual Image-Speech Retrieval

We evaluate our multilingual models on the same image-speech retrieval task, and report results in Table 1. When the feature extractors are frozen, monolingual models outperform multilingual ones in monolingual evaluation settings. Our language-aware and language-agnostic models perform similarly, indicating that providing an explicit language ID to the model is not necessary. On the other hand, we find that using mixed-language batches during training produces much better models than those which only consider distractors in the same language. This finding holds across both multilingual *and* monolingual tests. We therefore select the language-agnostic architecture and multilingual training batches for our two more computationally expensive experiments.

Allowing the feature extractor to be fine-tuned leads to sizable gains both in monolingual and in mixed-language evaluation settings, particularly for non-English inputs. Despite using only a HuBERT-Base feature extractor, the model trained with cross-lingual loss terms achieves additional gains on Hindi, Japanese, and mixed-language tests. Improvements from the cross-lingual loss terms are less pronounced in the monolingual English evaluation setting, making "Multi+TrainFeat+XLL" our only model to achieve better scores on the Places100k Japanese test set than the English.

For a multilingual baseline, we consider cascading ASR with the text-based multilingual M-CLIP [8] model. The English and Hindi Places datasets include ASR-generated transcriptions of the spoken captions; for the Japanese data, we use the XLSR-53 Large model fine-tuned on Japanese to generate transcriptions. We then use the XLM-R Large ViT-L/14 variant of M-CLIP to produce embeddings of the transcribed captions, which we use for image-speech retrieval. All of our models outperform this baseline, especially in the Japanese and multilingual evaluation settings.

## 5.3. Analysis of Learned Layer Weights

Fig. 2 summarizes the layer weights learned by our models with frozen feature extractors. Two separate patterns emerge for English and non-English weights, regardless of whether they are learned by a monolingual or a multilingual model. Layers 15 through 18 receive the highest weights when processing non-English inputs, whereas layers 17 through 21 receive the highest weights when processing English speech, suggesting that later HuBERT layers are more specialized to English than earlier ones. The maximum and minimum weights assigned to any layer by the language-aware multilingual model are more similar across languages than the weights learned by monolingual models. For the language-agnostic multilingual model, which uses the same weights for both English and non-English speech, layer weights are distributed according to the non-English pattern. This does not appear to hurt the model's retrieval performance on English-only test sets, which Table 1 shows is comparable for the language-agnostic and language-aware models.

## 5.4. Zero-Shot Speech-Text Retrieval

Since our targets $\mathbf{v}^I$ are image embeddings in the output space of CLIP-Large, we can compare our speech embeddings $\mathbf{v}^S$ to the output $\mathbf{v}^T$ of the CLIP-Large text encoder on the English captions for each image in the batch. We report the performance of our monolingual models and our two large multilingual models on this task.

| | Image→Speech | | | Speech→Image | | |
|---|---|---|---|---|---|---|
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| | | | Japanese | | | |
| Trilingual Embds. [16] | 20.0 | 46.8 | 62.3 | 20.3 | 52.0 | 66.7 |
| Pair Expansion [17] | 16.7 | 44.3 | 57.8 | 20.1 | 49.7 | 63.9 |
| AVLnet [19] | 24.3 | 56.6 | 70.0 | 23.5 | 57.3 | 70.4 |
| Japanese-SpeechCLIP | **32.1** | **66.6** | **78.7** | **32.9** | **66.4** | **77.5** |
| Monolingual Batches | 21.1 | 50.0 | 63.6 | 22.5 | 48.5 | 62.4 |
| Mono+LangID | 20.1 | 50.1 | 63.3 | 20.1 | 50.9 | 64.3 |
| Multilingual Batches | **26.0** | **56.3** | **69.7** | **26.5** | **55.1** | 68.0 |
| Multi+LangID | 24.5 | 55.4 | 69.0 | 24.6 | 54.5 | **69.8** |
| Multi+TrainFeat | 32.6 | 66.4 | 77.5 | 32.1 | 66.1 | 78.3 |
| Multi+TrainFeat+XLL | **51.0** | **83.3** | **91.2** | **48.8** | **80.0** | **90.0** |
| ASR→M-CLIP [8] | 12.0 | 35.2 | 46.3 | 22.9 | 47.4 | 58.5 |
| | | | Hindi | | | |
| Bilingual Embds. [15] | 7.4 | 23.5 | 35.4 | 8.0 | 25.0 | 35.6 |
| Trilingual Embds. [16] | 10.8 | 31.3 | 41.9 | 11.2 | 31.5 | 44.5 |
| Pair Expansion [17] | 9.3 | 29.5 | 38.2 | 9.4 | 29.8 | 41.8 |
| AVLnet [19] | 17.0 | 39.8 | 51.5 | 15.2 | 38.9 | 51.1 |
| Hindi-SpeechCLIP | **21.8** | **46.5** | **58.8** | **19.1** | **42.7** | **57.3** |
| Monolingual Batches | 16.5 | 35.6 | 46.9 | 12.3 | 32.8 | 44.5 |
| Mono+LangID | 17.1 | 37.1 | 49.4 | 12.5 | 35.7 | 47.5 |
| Multilingual Batches | 17.6 | 41.3 | 52.9 | **16.7** | **39.1** | **51.2** |
| Multi+LangID | **18.1** | **41.4** | **53.2** | 14.8 | 37.6 | 49.6 |
| Multi+TrainFeat | 23.5 | 52.5 | 64.7 | 20.7 | 48.4 | 62.2 |
| Multi+TrainFeat+XLL | **35.4** | **66.1** | **76.3** | **34.6** | **64.9** | **75.3** |
| ASR→M-CLIP [8] | 13.8 | 30.2 | 42.0 | 19.5 | 40.8 | 50.8 |
| | | | English | | | |
| Bilingual Embds. [15] | 8.0 | 25.2 | 36.5 | 8.3 | 28.2 | 42.4 |
| Trilingual Embds. [16] | 11.6 | 35.8 | 50.8 | 13.9 | 39.5 | 52.9 |
| Pair Expansion [17] | 12.3 | 35.3 | 47.7 | 13.8 | 40.2 | 51.6 |
| English-SpeechCLIP | **50.9** | **82.8** | **90.6** | **48.9** | **82.3** | **89.4** |
| Monolingual Batches | 37.5 | 73.2 | 84.7 | 36.1 | 70.9 | 83.1 |
| Mono+LangID | 37.8 | 72.9 | 84.0 | 36.1 | 71.3 | 83.9 |
| Multilingual Batches | **41.9** | 75.6 | 86.4 | **41.1** | 73.9 | 85.1 |
| Multi+LangID | 41.1 | **77.2** | **86.5** | 40.8 | **75.0** | **86.4** |
| Multi+TrainFeat | 44.6 | 79.8 | **89.1** | 44.2 | 78.3 | 87.3 |
| Multi+TrainFeat+XLL | **51.7** | **82.2** | 88.9 | **48.0** | **79.7** | **88.0** |
| ASR→M-CLIP [8] | 29.2 | 58.1 | 69.4 | 43.4 | 70.5 | 78.9 |
| | | | Multilingual | | | |
| Monolingual Batches | 25.0 | 55.3 | 66.8 | 29.7 | 53.5 | 65.4 |
| Mono+LangID | 25.0 | 52.0 | 65.3 | 28.3 | 52.4 | 66.2 |
| Multilingual Batches | **25.8** | 56.6 | 68.4 | **33.1** | **56.2** | 66.7 |
| Multi+LangID | 24.9 | **58.8** | **69.2** | 32.8 | 55.4 | **68.7** |
| Multi+TrainFeat | 34.0 | 63.1 | 75.9 | 31.0 | 61.5 | 74.1 |
| Multi+TrainFeat+XLL | **44.8** | **76.3** | **85.9** | **43.0** | **74.5** | **84.9** |
| ASR→M-CLIP [8] | 20.0 | 35.7 | 45.0 | 28.5 | 52.9 | 62.3 |

**Table 1**: Retrieval performance on the Places100k test set. For each language, we evaluate (from top to bottom): monolingual models; multilingual models trained on 2 V100 GPUs; multilingual models trained on 8 V100 GPUs; and monolingual ASR cascaded with multilingual image-text retrieval from prior work. All but the monolingual models are also evaluated on a multilingual test set formed by randomly selecting a language for each image's spoken caption.
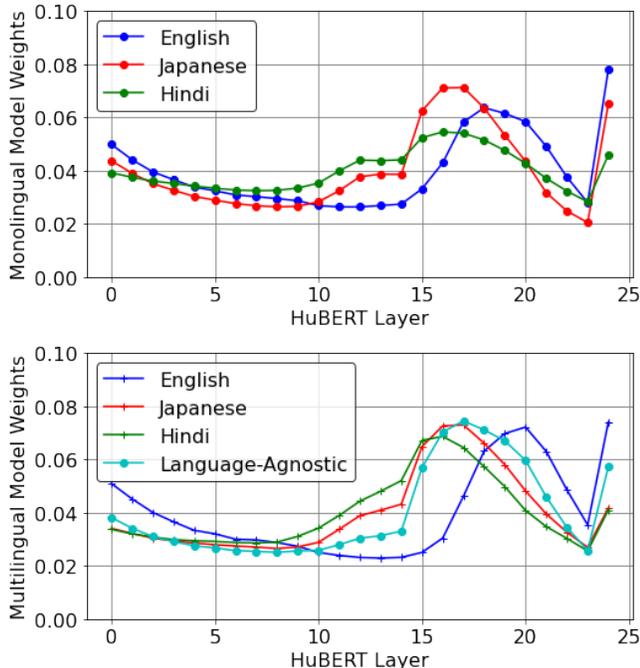
**Fig. 2**: Top: Layer weights learned by monolingual models. Bottom: Layer weights learned by multilingual models trained on multilingual batches with a frozen HuBERT-Large feature extractor.

| | Prompt→Target | | | Target←Prompt | | |
|---|---|---|---|---|---|---|
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| | *Image ↔ English Text* | | | | | |
| CLIP [1] | 26.1 | 51.3 | 63.8 | 46.4 | 71.8 | 78.4 |
| | *English Speech ↔ English Text* | | | | | |
| Monolingual* | **57.1** | **80.0** | **85.1** | 67.8 | 89.1 | 93.1 |
| Multi+TrainFeat* | 50.0 | 75.7 | 84.0 | 63.0 | 85.2 | 90.1 |
| Multi+TrainFeat+XLL* | 39.0 | 68.7 | 79.1 | **70.1** | **90.1** | **93.2** |
| | *Japanese Speech ↔ English Text* | | | | | |
| Monolingual* | **11.8** | **34.8** | **45.1** | 12.9 | 31.3 | 42.9 |
| Multi+TrainFeat* | 9.8 | 28.6 | 39.9 | 11.3 | 27.9 | 38.1 |
| Multi+TrainFeat+XLL* | 10.8 | 27.7 | 39.1 | **23.9** | **48.8** | **60.1** |
| | *Hindi Speech ↔ English Text* | | | | | |
| Monolingual* | 8.4 | 22.3 | 31.7 | 8.0 | 22.5 | 32.3 |
| Multi+TrainFeat* | 8.6 | 25.2 | 35.1 | 8.1 | 25.9 | 39.6 |
| Multi+TrainFeat+XLL* | **14.8** | **34.5** | **45.7** | **16.9** | **32.3** | **40.8** |
| | *English Speech ↔ Japanese Speech* | | | | | |
| Trilingual Embds. [16] | 10.5 | 31.2 | 43.7 | 10.6 | 31.7 | 44.1 |
| Multi+TrainFeat* | 10.2 | 28.9 | 41.1 | 10.7 | 30.0 | 42.4 |
| Multi+TrainFeat+XLL | **28.9** | **56.4** | **69.6** | **28.6** | **57.2** | **69.6** |
| | *English Speech ↔ Hindi Speech* | | | | | |
| Trilingual Embds. [16] | 7.6 | 22.5 | 31.3 | 7.6 | 22.5 | 31.3 |
| Multi+TrainFeat* | 9.3 | 27.6 | 40.2 | 10.1 | 28.6 | 40.6 |
| Multi+TrainFeat+XLL | **26.0** | **49.5** | **60.2** | **25.1** | **50.9** | **61.8** |
| | *Japanese Speech ↔ Hindi Speech* | | | | | |
| Trilingual Embds. [16] | 10.4 | 24.6 | 35.0 | 8.5 | 24.8 | 33.4 |
| Multi+TrainFeat* | 7.0 | 19.4 | 29.4 | 7.7 | 18.4 | 26.4 |
| Multi+TrainFeat+XLL | **22.5** | **46.1** | **56.3** | **22.0** | **45.8** | **59.0** |

**Table 2**: Image-text, speech-text, and speech-speech retrieval on the Places100k test set. * indicates a model is not trained on this task.

Inspired by [20], we consider image-text retrieval using the CLIP image and text encoders as a point of comparison. Since the CLIP image encoder outputs were used as our targets during training, we would expect this to function as a topline–the performance our model would achieve if it output precisely the training targets. Instead, as in [20], our models significantly outperform this experiment, with improvements of nearly 15% when English text is the query and over 20% when English text is being retrieved.

In the cross-lingual setting, despite the difficulty of the task, our models achieve reasonable performance. Our monolingual Japanese model performs the best in the Japanese Speech→English Text direction, while the multilingual model trained with cross-lingual loss performs best for all other directions. This is the first work to successfully perform zero-shot transfer from image-speech retrieval to cross-lingual text-speech retrieval with texts longer than a single word, and sets a strong baseline for the task. We achieve this with no text at all in the spoken source languages.

### 5.5. Cross-Lingual Speech-Speech Retrieval

Finally, we consider the task of cross-lingual speech-to-speech retrieval. We report results for our "Multi+TrainFeat" model, for which this is a zero-shot transfer task, as well as for our "Multi+TrainFeat+XLL" model, for which it is not. The monolingual models trained simultaneously with cross-lingual loss terms from [16] provide a baseline for this task. Our multilingual model trained with the cross-lingual loss terms from [16] outperforms prior work by 20-30% R@10 for all retrieval directions. Even our model trained without any cross-lingual supervision can perform zero-shot transferreasonably well. While it underperforms compared to [16] in the English ↔ Japanese and Japanese ↔ Hindi directions, it outperforms prior work for the English ↔ Hindi directions, despite the

significant disadvantage of having not been trained for cross-lingual speech-speech retrieval. This indicates that learning to align spoken captions with CLIP image embeddings can produce high-quality alignments between speech in different languages.

## 6. CONCLUSION

In this paper, we found that large-scale, English-only pre-training is effective not only for downstream tasks processing English speech, but can be used to achieve state-of-the-art performance even for non-English image-speech retrieval. The CLIP semantic embedding space can be used to represent not only multiple modalities, but multiple languages effectively. Examining the layer weights learned by our model revealed that pre-trained HuBERT speech encoders specialize for English in later layers, but that features extracted from their middle layers are useful for non-English downstream tasks even without any fine-tuning. We also trained a single model to encode speech in multiple languages into the CLIP semantic embedding space. This model can then be used to perform zero-shot cross-lingual speech-to-speech retrieval. Our models significantly outperform prior work when trained with a cross-lingual objective, and perform comparably to prior work even when trained without one. Finally, we showed that our learned speech encoders can perform zero-shot speech-text retrieval for English text even when the speech is not in English. In our future work, we plan to explore the use of these models to bootstrap speech-to-text and speech-to-speech translation for low-resource languages when neither parallel data nor non-English text is available.

# 7. REFERENCES

[1] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever, "Learning transferable visual models from natural language supervision," in *Proceedings of the 38th International Conference on Machine Learning*, 2021.

[2] Yi-Jen Shih, Hsuan-Fu Wang, Heng-Jui Chang, Layne Berry, Hung-yi Lee, and David Harwath, "Speechclip: Integrating speech with pre-trained vision and language model," in *IEEE Spoken Language Technology Workshop*, 2022.

[3] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision*, 2014, pp. 740–755.

[4] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," in *International Journal of Computer Vision*, 2016.

[5] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li, "Yfcc100m: The new data in multimedia research," *Communications of the Association for Computing Machinery*, vol. 59, 2016.

[6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *9th International Conference on Learning Representations*, 2021.

[7] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever, "Language models are unsupervised multitask learners," *OpenAI blog*, 2019.

[8] Fredrik Carlsson, Philipp Eisen, Faton Rekathati, and Magnus Sahlgren, "Cross-lingual and multilingual clip," in *Proceedings of the Language Resources and Evaluation Conference*, June 2022, pp. 6848–6854.

[9] Pranav Aggarwal, Ritiz Tambi, and Ajinkya Kale, "Towards zero-shot cross-lingual image retrieval and tagging," in *The Web Conference Workshop on Multilingual Search*, 2021.

[10] Puyuan Peng and David Harwath, "Fast-slow transformer for visually grounding speech," in *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2022.

[11] Wei-Ning Hsu, David Harwath, Christopher Song, and James R. Glass, "Text-free image-to-speech synthesis using learned segmental units," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 2021.

[12] David F. Harwath and James R. Glass, "Deep multimodal semantic embeddings for speech and images," *IEEE Workshop on Automatic Speech Recognition and Understanding*, 2015.

[13] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva, "Learning deep features for scene recognition using places database," in *Advances in Neural Information Processing Systems*, 2014, vol. 27.

[14] David Harwath, Antonio Torralba, and James Glass, "Unsupervised learning of spoken language with visual context," in *Advances in Neural Information Processing Systems*, 2016.

[15] David F. Harwath, Galen Chuang, and James R. Glass, "Vision as an interlingua: Learning multilingual semantic embeddings of untranscribed speech," in *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2018.

[16] Yasunori Ohishi, Akisato Kimura, Takahito Kawanishi, Kunio Kashino, David Harwath, and James Glass, "Trilingual semantic embeddings of visually grounded speech with self-attention mechanisms," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 4352–4356.

[17] Yasunori Ohishi, Akisato Kimura, Takahito Kawanishi, Kunio Kashino, David F. Harwath, and James R. Glass, "Pair expansion for learning multilingual semantic embeddings using disjoint visually-grounded speech audio datasets," in *21st Annual Conference of the International Speech Communication Association*, 2020.

[18] William Havard, Jean-Pierre Chevrot, and Laurent Besacier, "Models of visually grounded speech signal pay attention to nouns: a bilingual experiment on english and japanese," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019.

[19] Andrew Rouditchenko, Angie W. Boggust, David Harwath, Samuel Thomas, Hilde Kuehne, Brian Chen, Rameswar Panda, Rogério Feris, Brian Kingsbury, Michael Picheny, and James R. Glass, "Cascaded multilingual audio-visual learning from videos," in *22nd Annual Conference of the International Speech Communication Association*, 2021.

[20] Herman Kamper, Gregory Shakhnarovich, and Karen Livescu, "Semantic speech retrieval with a visually grounded model of untranscribed speech," 2018.

[21] Ankita Pasad, Bowen Shi, Herman Kamper, and Karen Livescu, "On the contributions of visual and textual supervision in low-resource semantic speech retrieval," in *20th Annual Conference of the International Speech Communication Association*, 2019.

[22] Herman Kamper and Michael Roth, "Visually grounded cross-lingual keyword spotting in speech," in *6th Workshop on Spoken Language Technologies for Under-Resourced Languages*, 2018.

[23] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," in *IEEE/ACM Transactions on Audio, Speech and Language Processing*, 2021, vol. 29.

[24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *31st Conference on Neural Information Processing Systems*, 2017.

[25] Gabriel Ilharco, Yuan Zhang, and Jason Baldridge, "Large-scale representation learning from visually grounded untranscribed speech," in *Proceedings of the 23rd Conference on Computational Natural Language Learning*, 2019.