

Write a report on each of following speech transformers:

1. HUBERT:

Hubert is a speech model that accepts a float array corresponding to the raw waveform of the speech signal. HUBERT stands for Hidden Unit BERT. Hubert model was fine-tuned using connectionist temporal classification (CTC) so the model output has to be decoded using [Wav2Vec2CTCTokenizer](#).

Why we can't use NLP models on audio

There are 3 main problems when trying to apply BERT or other NLP models on speech data:

1. There are multiple sound units in each input expression
2. There is no lexicon of discrete sound units
3. Sound units have variable length and no explicit segmentation

Problem 1 prevents the usage of techniques such as instance classification, which are used in Computer Vision for pre-training.

Problem 2 hinders the usage of predictive losses, due to not having a reliable target to compare the prediction with. Finally, problem 3 complicates masked prediction pre-training due to unknown borders between sound units.

To solve these problems the authors propose **HuBERT**. HUBERT utilizes an offline clustering step to provide aligned target labels for a BERT-like prediction loss. A key ingredient of this approach is applying the prediction loss over the masked regions only, which forces the model to learn a combined acoustic and language model over the continuous inputs.

The HuBERT model architecture follows the wav2vec 2.0 architecture consisting of:

- Convolutional encoder
- BERT encoder
- Projection layer
- Code embedding layer

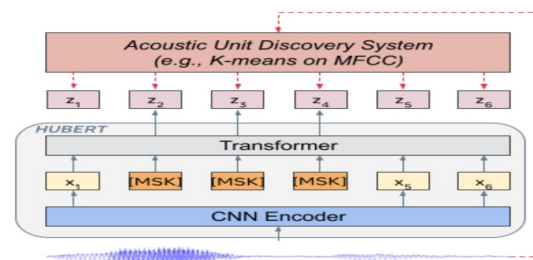
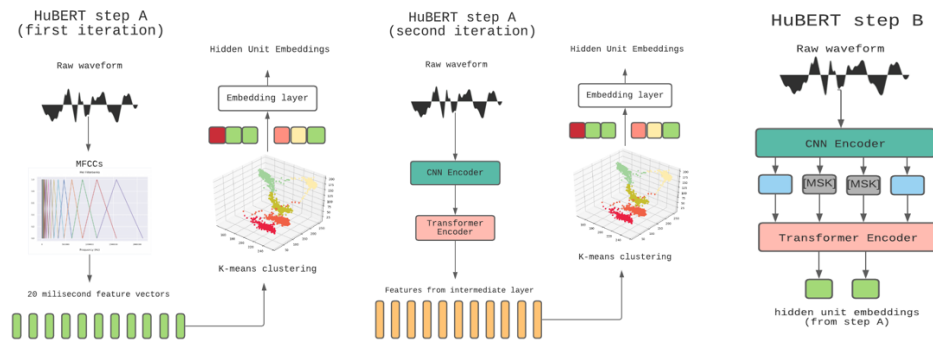


Fig. 1: The HuBERT approach predicts hidden cluster assignments of the masked frames (y_2, y_3, y_4 in the figure) generated by one or more iterations of k-means clustering.



HuBERT initial clustering step , HuBERT subsequent clustering step , HuBERT Prediction step

The first training step consists of discovering the hidden units, and the process begins with extracting [MFCCs\(Mel frequency cepstrum\)](#) from the audio waveform.

These are raw acoustic features useful for representing speech.

Each segment of audio is then passed to the [K-means clustering algorithm](#), and assigned to one of K clusters.

All audio frames will then be labeled according to which cluster they belong to, and these are the **hidden units**.

Afterward, these units are converted into embedding vectors to be used in step B of training.

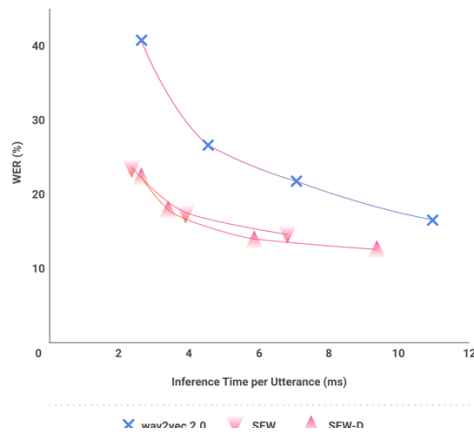
HuBERT relies primarily on the consistency of the unsupervised clustering step rather than the intrinsic quality of the assigned cluster labels.

Starting with a simple k-means teacher of 100 clusters, and using two iterations of clustering, the HuBERT model either matches or improves upon the state-of-the-art wav2vec 2.0 performance on the Librispeech (960h) and Libri-light (60,000h) benchmarks with 10min, 1h, 10h, 100h, and 960h fine-tuning subsets. Using a 1B parameter model, HuBERT shows up to 19% and 13% relative WER reduction on the more challenging dev-other and test-other evaluation subsets.

The second step is analogous to the training of the original BERT model, using masked language modeling. The CNN is responsible for generating features from the raw audio, which are then randomly masked and fed into the BERT encoder. The BERT encoder outputs a feature sequence, filling in the masked tokens. This output is then projected into a lower dimension to match the labels and the cosine similarity is computed between these outputs and each hidden unit embedding generated in step A. The cross-entropy loss is then used on the logits to penalize wrong predictions.

2. SEW

SEW stands for Squeezed and Efficient Wav2Vec. This paper is a study of performance-efficiency trade-offs in pre-trained models for automatic speech recognition (ASR). We focus on wav2vec 2.0, and formalize several architecture designs that influence both the model performance and its efficiency. Putting together all our observations, we introduce SEW (Squeezed and Efficient Wav2vec), a pre-trained model architecture with significant improvements along both performance and efficiency dimensions across a variety of training setups. For example, under the 100h-960h semi-supervised setup on LibriSpeech, SEW achieves a 1.9x inference speedup compared to wav2vec 2.0, with a 13.5% relative reduction in word error rate. With a similar inference time, SEW reduces word error rate by 25-50% across different model sizes. The paper focuses on performance efficiency trade-offs in pre-trained models for automatic speech recognition (ASR). The focus is on wav2vec 2.0 and formalizing several architectural designs that influence both the model performance and its efficiency. SEW is pre-trained model architecture with significant improvement in both performance and efficiency dimensions across a variety of training setups.



SEW (Squeezed and Efficient Wav2vec)

LibriSpeech 960h is used as training data for unsupervised pre-training, leaving 1% as a validation set for pre-training. Hyperparameters used are W2V2-base2. To speed up and reduce the cost of the experiments, pre-training all models for 100K updates similar to Hsu et al. (2020). All experiments use an AWS p3.16xlarge instance with 8 NVIDIA V100 GPUs and 64 Intel Xeon 2.30GHz CPU cores. Because Baevski et al. (2020b) use 64 GPUs, we set gradient accumulation steps to 8 to simulate their 64-GPU pre-training with 8 GPUs.

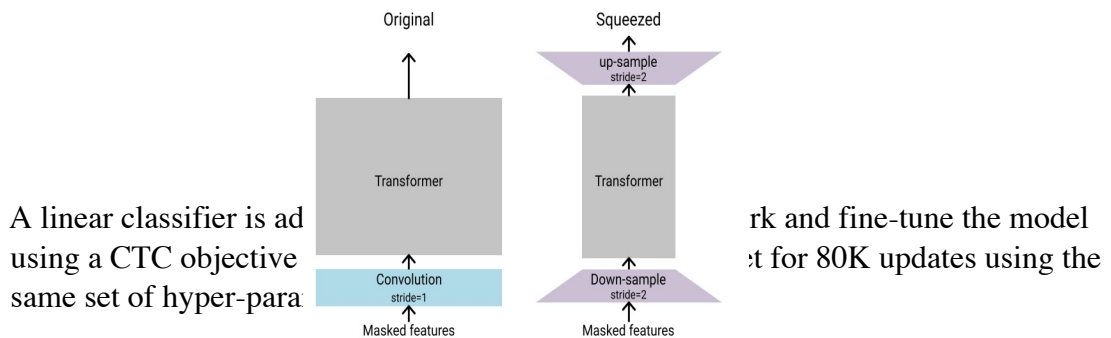


Figure 4: Original vs. squeezed context network. The sequence length is halved by the down-sampling layer.

The use of CTC is used for greedy decoding for all experiments because it is faster than Viterbi decoding (Viterbi, 1967) and we do not find any WER differences between the two using baseline W2V2 models. We use LibriSpeech dev-other for validation, and hold out test-clean and test-other as test sets. We consider three metrics to evaluate model efficiency and performance: pre-training time, inference time, and WER (word error rate). All evaluation is done on a NVIDIA V100 GPU with FP32 operations, unless specified otherwise. When decoding with a language model (LM), we use the official 4gram LM4 and wav2letter (Collobert et al., 2016) decoder5 with the default LM weight 2, word score -1, and beam size 50. Reducing the inference time with LM is an important direction for future work, as the wav2letter decoder is the bottleneck and is at least 3× slower than W2V2-base

Without an LM, compared with the W2V2-tiny, SEW-tiny reduces the WER by 53.5% (22.8% to 10.6%) and 43.7% (41.1% to 23.7%) on test-clean and test-other, while being faster. With an LM, WER improves by 38.6% and 43.4% on test-clean and test-other. Compared with the W2V2-mid, SEW-mid reduces WER by 30.2% (9.6% to 6.7%) and 32.9% (22.2% to 14.9%) with similar inference times. SEW does incur slight increase in training time compared to W2V2 with similar inference times. However, SEW has lower WER even compared to a slower W2V2 which takes longer to train (e.g., SEW-small vs. W2V2-mid or SEW-mid vs. W2V2-base). SEW-D has lower WER compared to SEW even with smaller width and half of the parameters. With large models, SEW-D is also more efficient. However, SEW-D-tiny is slower than SEW-tiny, due to the implementation difference

3. FAIRSEQ S2T: Fast Speech-to-Text Modeling with FAIRSEQ

Fairseq S2T, a fairseq extension for speech-to-text (S2T) modeling tasks such as end-to-end speech recognition and speech-to-text translation. It follows fairseq’s careful design for scalability and extensibility. We provide end-to-end workflows from data pre-processing, model training to offline (online) inference. We implement state-of-the-art RNN-based as well as Transformer-based models and open-source detailed training recipes. Fairseq’s machine translation models and language models can be seamlessly integrated into S2T workflows for multi-task learning or transfer learning. FAIRSEQ S2T is an extension for Speech2Text modelling task such as Speech recognition and Speech2Text translation.

Fairseq Models FAIRSEQ provides a collection of MT models and LMs that demonstrate state-of-the-art performance on standard benchmarks. They are open sourced with pretrained models. FAIRSEQ also supports other tasks such as text summarization, story generation and self-supervised speech pre-training.

S2T extension FAIRSEQ S2T adds attention-based RNN, Transformer models, as well as the latest Conformer models for ASR and ST. It also supports the CTC criterion for ASR. For the simultaneous ST setting, it includes online models with widely used policies: monotonic attention, wait-k (Ma monotonic infinite lookback attention, and monotonic multi-head attention

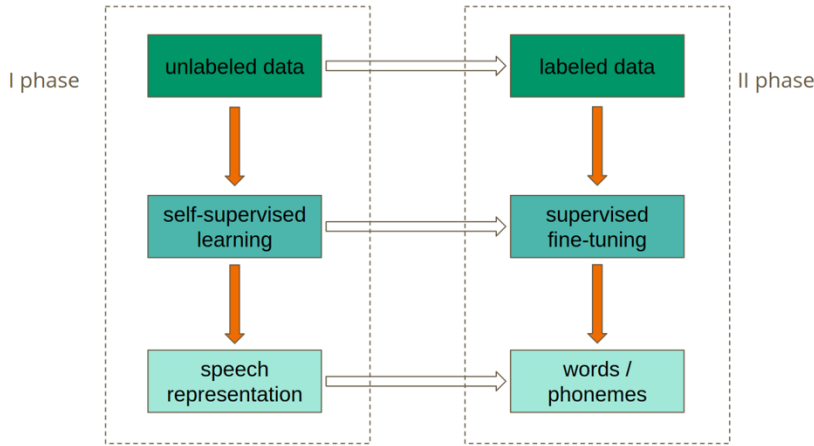
	Type	Config.	Params
B-Base	RNN ¹	512d, 3L enc./2L dec.	31M
B-Big		512d, 5L enc./3L dec.	52M
T-Sm	Trans- former ²	256d, 12L enc./6L dec.	31M
T-Md		512d, 12L enc./6L dec.	72M
T-Lg		1024d, 12L enc./6L dec.	263M
W-Lg	wav2vec	1024d, 24L	315M
CW-Lg	2.0 ³	1024d, 24L, Conformer ⁴	618M

Table 3: FAIRSEQ S2T models for benchmarking. For simplicity, we use the same (default) model hyperparameters and learning rate schedule across all experiments. ¹ Bérard et al. (2018). ² Vaswani et al. (2017). ³ Baevski et al. (2020). ⁴ Gulati et al. (2020).

Evaluation Metrics FAIRSEQ S2T provides common automatic metrics for ASR, ST, and MT, including WER (word error rate), BLEU, and chrF It also integrates SIMULEVAL for simultaneous ST/MT metrics such as AL (average lagging) and DAL (differentiable average Lagging)

4. Wav2Vec

Wav2Vec 2.0 is one of the current state-of-the-art models for Automatic Speech Recognition due to a self-supervised training which is quite a new concept in this field. This way of training allows us to pre-train a model on unlabeled data which is always more accessible. Then, the model can be fine-tuned on a particular dataset for a specific purpose. As the previous works show this way of training is very powerful.



Training phases of Wav2Vec 2.0

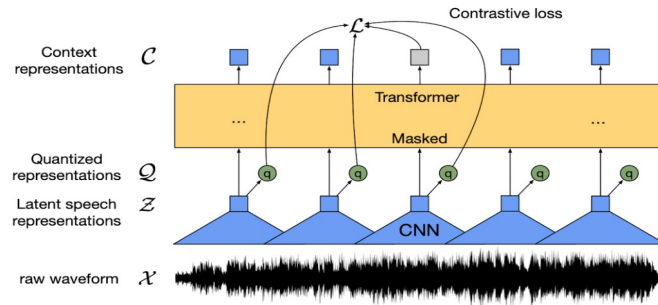
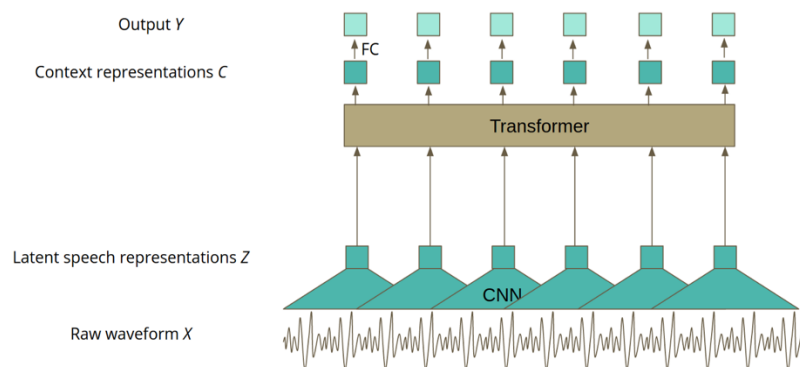


Illustration of the Wav2vec2 framework

Wav2Vec 2.0 Model Architecture

The architecture of the final model used for prediction consists of three main parts:

- convolutional layers that process the raw waveform input to get latent representation - Z ,
- transformer layers, creating contextualised representation - C ,
- linear projection to output - Y .



Fine-tuned Wav2Vec 2.0 model architecture (Image by author, based on wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations)

A major advantage of this approach is that we end up training a generic audio model that could be used for multiple downstream tasks! And because of the self-supervised learning, we don't need access to huge amount of labelled data. In the paper, after pre-training on unlabelled speech, the model is fine-tuned on small labelled data with a Connectionist Temporal Classification (CTC) loss for speech recognition task.

The complete architecture of the framework can be divided into 3 components, they are

Feature encoder: This is the encoder part of the model. It takes the raw audio data as input and outputs feature vectors. Input size is limited to 400 samples which is 20ms for 16kHz sample rate. The raw audio is first standardized to have zero mean and unit variance. Then it is passed to 1D convolutional neural network (temporal convolution) followed by layer normalization and GELU activation function. There could be 7 such convolution blocks with constant channel size (512), decreasing kernel width (10, 3x4, 2x2) and stride (5, 2x6). The output is list of feature vectors each with 512 dimensions.

Transformers: The output of the feature encoder is passed on to a transformer layer. One differentiator is use of relative positional embedding by using convolution layers, rather than using fixed positional encoding as done in original Transformers paper. The block size differs, as 12 transformers block with model dimension of 768 is used in BASE model but 24 blocks with 1024 dimension in LARGE version.

Quantization module: For self-supervised learning, we need to work with discrete outputs. For this, there is a quantization module that converts the continuous vector output to discrete representations, and on top of it, it automatically learns the discrete speech units.

This is done by maintaining multiple codebooks/groups (320 in size) and the units are sampled from each codebook are later concatenated ($320 \times 320 = 102400$ possible speech units). The sampling is done using Gumbel-SoftMax which is like argmax but differentiable.

Training

- To pre-train the model, Wav2Vec2 masks certain portions of time steps in the feature encoder which is similar to the masked language model. The aim is to teach the model to predict the correct quantized latent audio representation in a set of distractors for each time step.
- The overall training objective is to minimize contrastive loss (L_m) and diversity loss (L_d) in $L = L_m + \alpha L_d$. Contrastive loss is the performance on the self-supervised task. Diversity loss is designed to increase the use of the complete quantized codebook representations and not only a select subset.
- For pretraining, the datasets used were (1) Librispeech corpus with 960 hours of audio data, (2) LibriVox 60k hours of audio data that was later subset to 53.2k hours. Only unlabelled data was used for pretraining.
- To make the model more robust to different tasks, we can finetune the model on different task-specific modifications and datasets. Here, the paper finetuned for ASR by adding a randomly initialized classification layer on top on the Transformer layer with a class size equal to the size of vocab. The model is optimized by minimizing the CTC loss.
- Adam was used as an optimization algorithm and the learning rate is warmed up till 10% of the training duration, then kept constant for the next 40%, and finally linearly decayed for the remaining duration. Also, for the first 60k updates only the output classifier was trained after which the Transformer is also updated. The feature encoder is kept frozen (not trained at all).

Results

There are two interesting points to note from the results of the Wav2Vec2 model,

- The model is able to learn ASR with as minimum as 10 mins of labelled data! As shown below, LARGE model pre-trained on LV-60k and finetuned on Librispeech with CTC loss is giving 4.6/7.9 WER! This is a very good news in case you want to finetune the model for your domain or accent!

- The choice of decoder can lead to improvement in performance. As obvious from the results, Transformer decoder is giving best performance, followed by n-gram and then CTC decoding. But also note that the CTC decoding will gives the best inference speed. The suggested decoder could be 4-gram, as it provides huge improvement in performance by fixing the spelling mistakes and grammar issues of CTC and is still faster than Transformer decoders.

Model	Unlabeled data	LM	dev		test	
			clean	other	clean	other
10 min labeled						
BASE	LS-960	None	46.1	51.5	46.9	50.9
		4-gram	8.9	15.7	9.1	15.6
		Transf.	6.6	13.2	6.9	12.9
LARGE	LS-960	None	43.0	46.3	43.5	45.3
		4-gram	8.6	12.9	8.9	13.1
		Transf.	6.6	10.6	6.8	10.8
LARGE	LV-60k	None	38.3	41.0	40.2	38.7
		4-gram	6.3	9.8	6.6	10.3
		Transf.	4.6	7.9	4.8	8.2
1h labeled						
BASE	LS-960	None	24.1	29.6	24.5	29.7
		4-gram	5.0	10.8	5.5	11.3
		Transf.	3.8	9.0	4.0	9.3
LARGE	LS-960	None	21.6	25.3	22.1	25.3
		4-gram	4.8	8.5	5.1	9.4
		Transf.	3.8	7.1	3.9	7.6
LARGE	LV-60k	None	17.3	20.6	17.2	20.3
		4-gram	3.6	6.5	3.8	7.1
		Transf.	2.9	5.4	2.9	5.8
10h labeled						
BASE	LS-960	None	10.9	17.4	11.1	17.6
		4-gram	3.8	9.1	4.3	9.5
		Transf.	2.9	7.4	3.2	7.8
LARGE	LS-960	None	8.1	12.0	8.0	12.1
		4-gram	3.4	6.9	3.8	7.3
		Transf.	2.9	5.7	3.2	6.1
LARGE	LV-60k	None	6.3	9.8	6.3	10.0
		4-gram	2.6	5.5	3.0	5.8
		Transf.	2.4	4.8	2.6	4.9
100h labeled						
BASE	LS-960	None	6.1	13.5	6.1	13.3
		4-gram	2.7	7.9	3.4	8.0
		Transf.	2.2	6.3	2.6	6.3
LARGE	LS-960	None	4.6	9.3	4.7	9.0
		4-gram	2.3	5.7	2.8	6.0
		Transf.	2.1	4.8	2.3	5.0
LARGE	LV-60k	None	3.3	6.5	3.1	6.3
		4-gram	1.8	4.5	2.3	4.6
		Transf.	1.9	4.0	2.0	4.0

Wav2Vec 2.0 uses a self-supervised training approach for Automatic Speech Recognition, which is based on the idea of *contrastive learning*. Learning speech representation on a huge, raw (unlabeled) dataset reduces the amount of labeled data required for getting satisfying results.

Key takeaways from this article:

- Wav2Vec 2.0 takes advantage of self-supervised training,
- it uses convolutional layers to preprocess raw waveform and then it applies transformer to enhance the speech representation with context,
- its objective is a weighted sum of two loss functions:
- contrastive loss,
- diversity loss