

Self-Supervised Learning for Speech Representations

Automatic Speech Processing, Master Cycle

Motivation

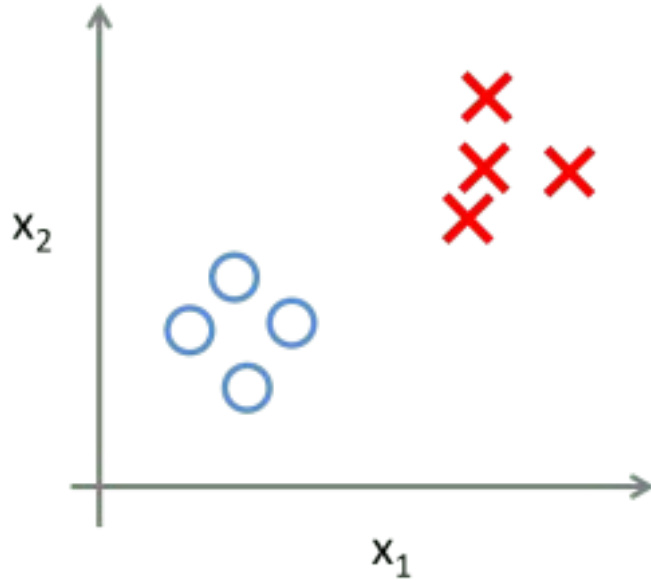
Speech datasets with transcriptions can be scarce / costly to create, depending on the task

Meanwhile, there is an abundance of unlabeled speech data

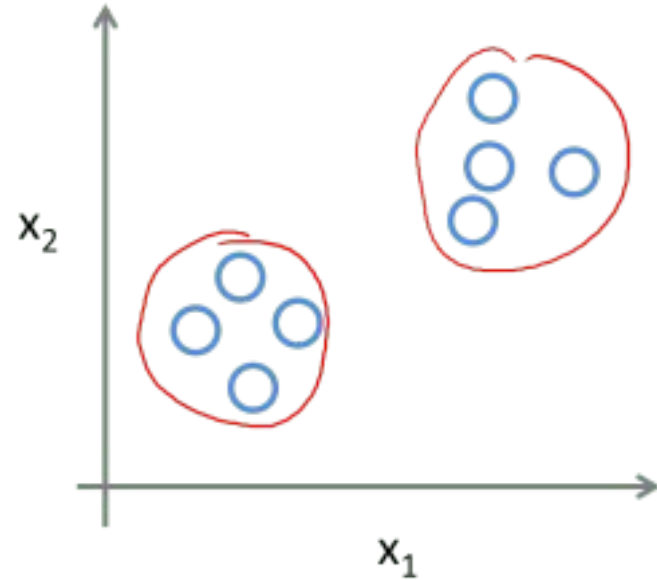
How can we benefit from unlabeled data?

→ Self-supervised Learning

Supervised vs. Unsupervised Learning



Supervised
Learn from labeled examples



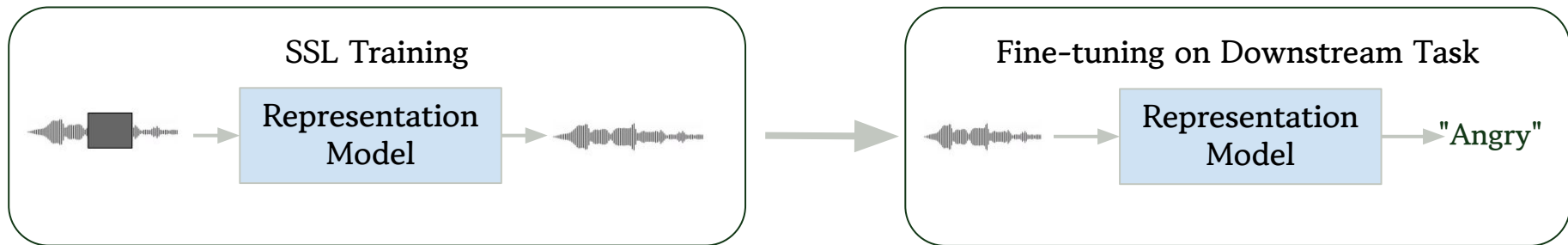
Unsupervised
Find structure in unlabeled data

Self-Supervised Learning (SSL)

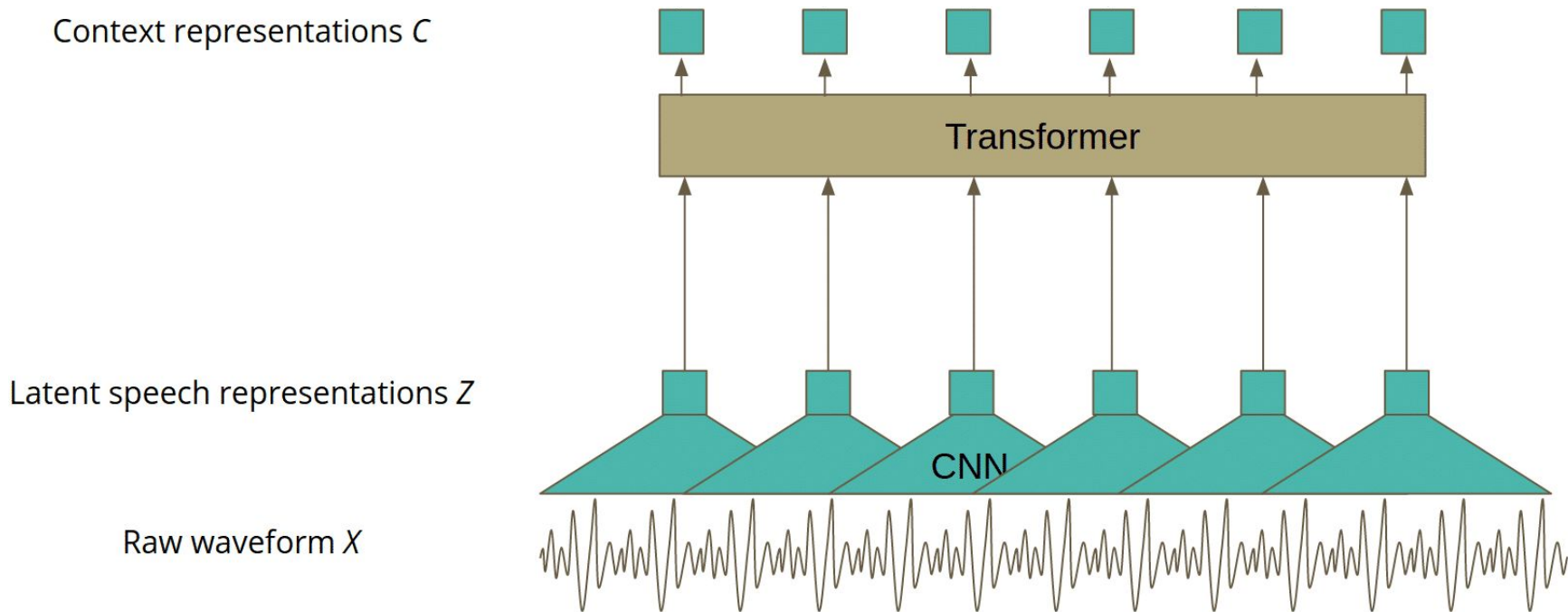
Data is still unlabeled but we design a pretext task to generate pseudo-labels.

A supervision signal is created from the unlabeled data itself

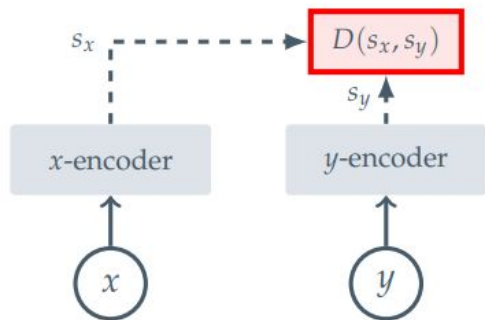
Goal: Learn useful speech representations which can perform well for a wide-range of downstream tasks



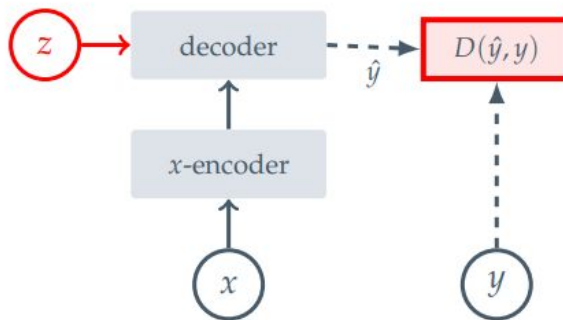
Typical Speech SSL Model Architecture



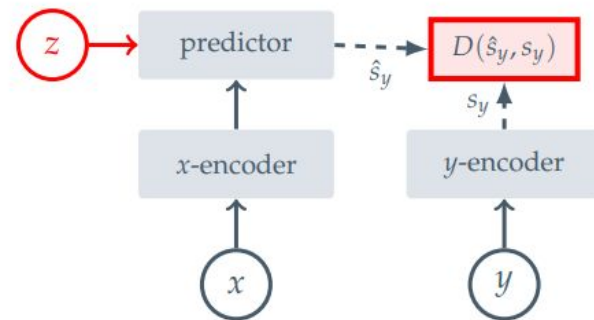
SSL Training Architectures



(a) Joint-Embedding Architecture



(b) Generative Architecture

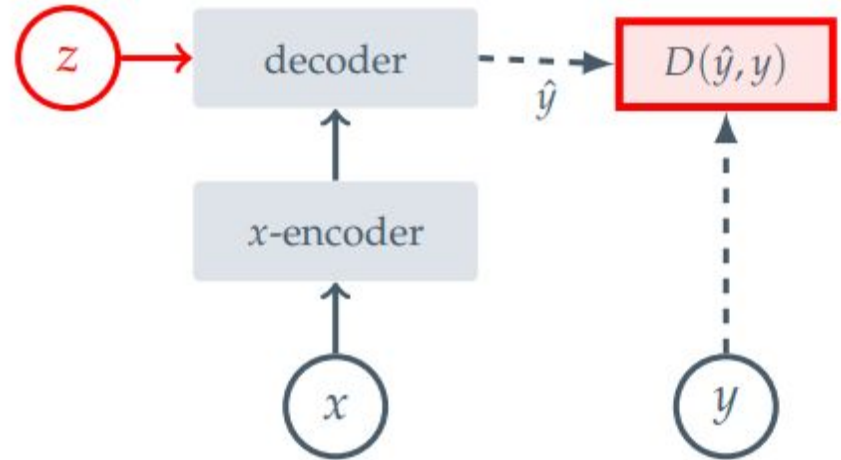


(c) Joint-Embedding Predictive Architecture

Generative Architecture

Learns to directly reconstruct a signal y from a compatible signal x

Decoder conditioned on additional variables z to facilitate reconstruction



(b) Generative Architecture

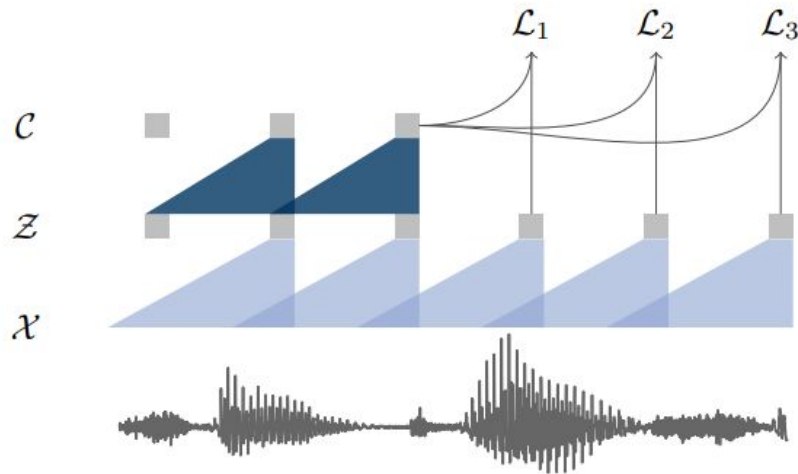
Generative Architecture

Wav2vec

Solves a next time-step prediction task

Predicts input features (raw audio encoded by the CNN encoder network)

Objective: Learns to distinguish a sample z_{i+k} that is k steps into the future from distractor samples drawn from a proposal distribution



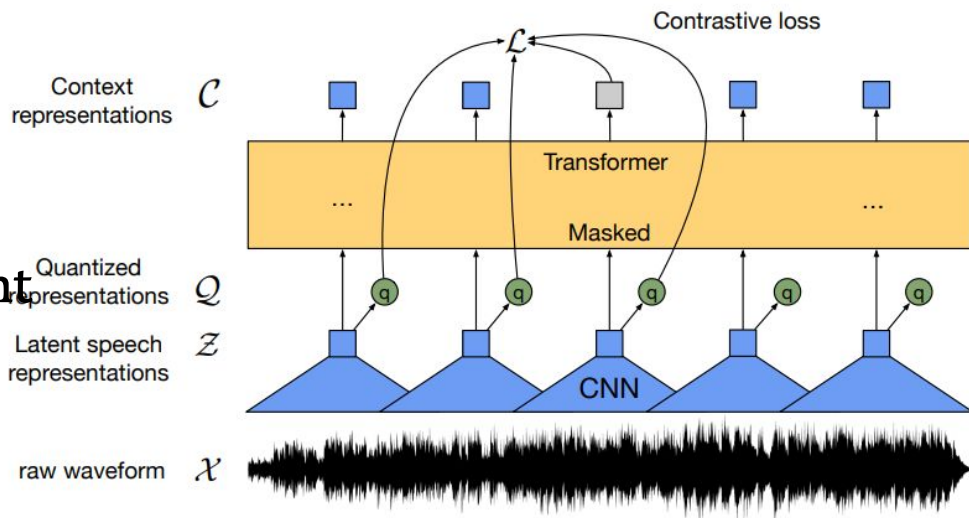
Generative Architecture

Wav2vec2.0

Solves a contrastive task defined over a quantization of the latent representations z

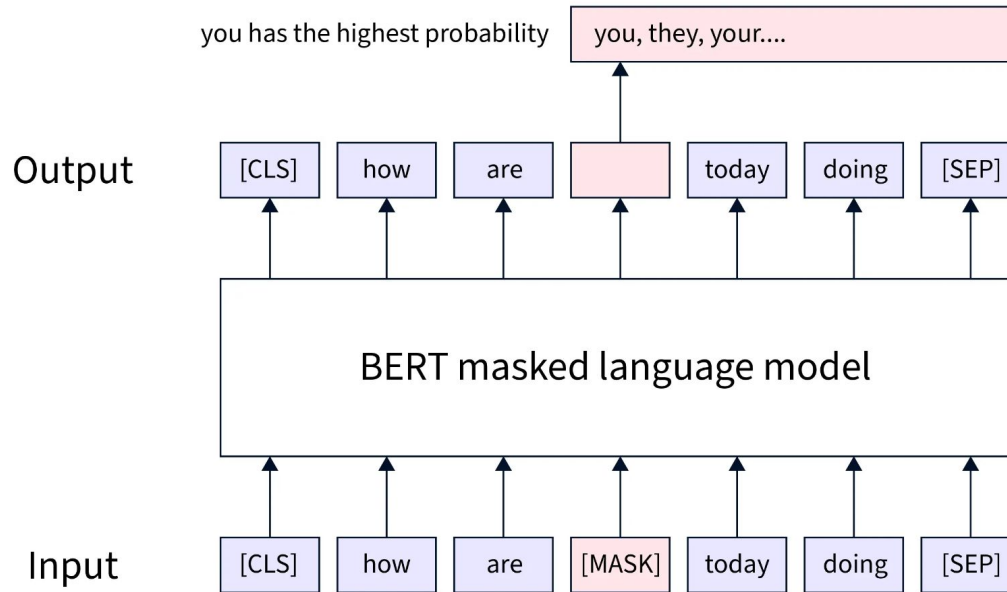
Masks the speech input in the latent space (CNN encoder output)

Objective: Learns to distinguish a true latent from distractors using quantized targets



Generative Architecture

Bert (Text)

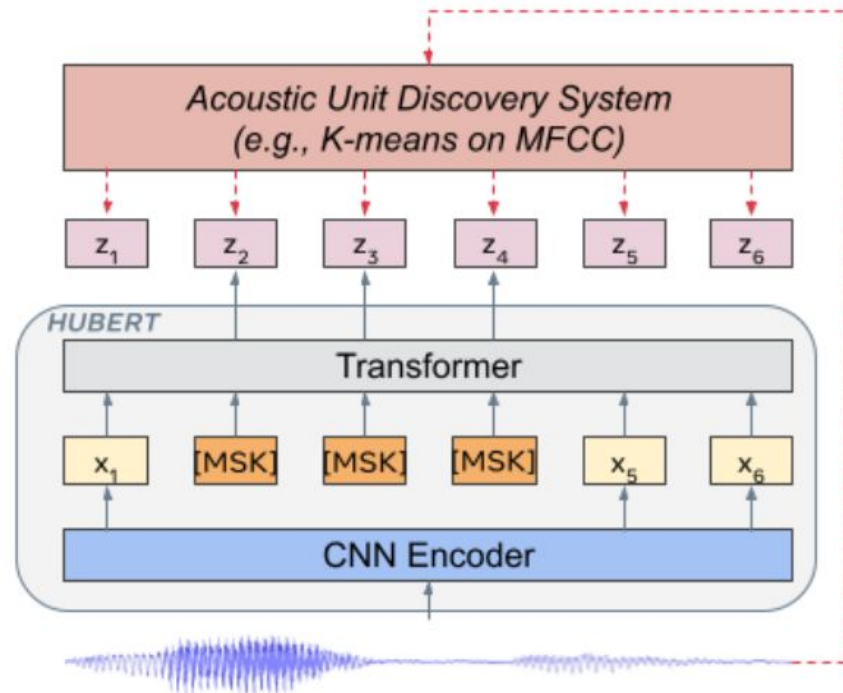


Generative Architecture

HuBERT

Uses offline clustering to provide aligned target labels for a BERT-like prediction loss

Prediction loss is applied over the masked regions only

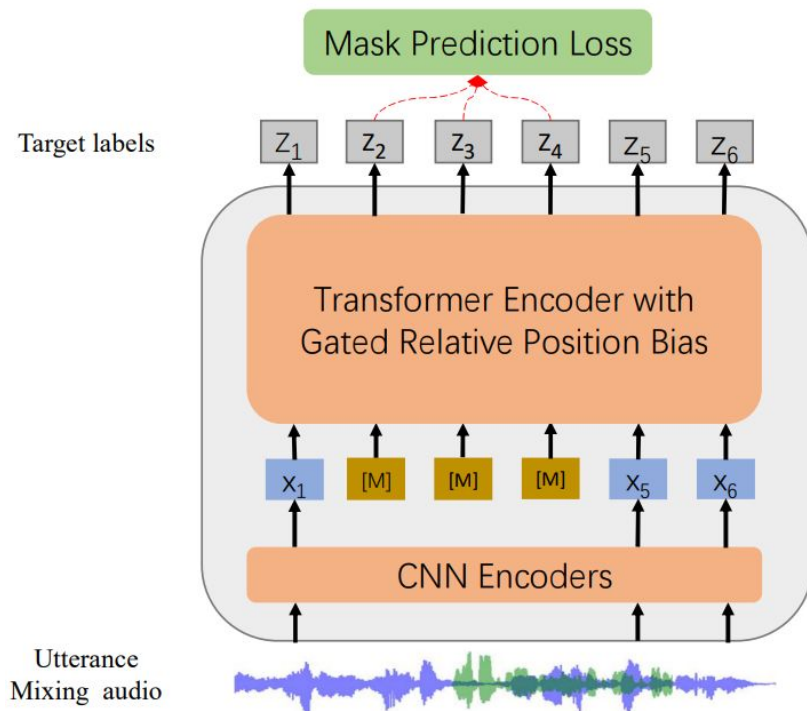


Generative Architecture

WavLM

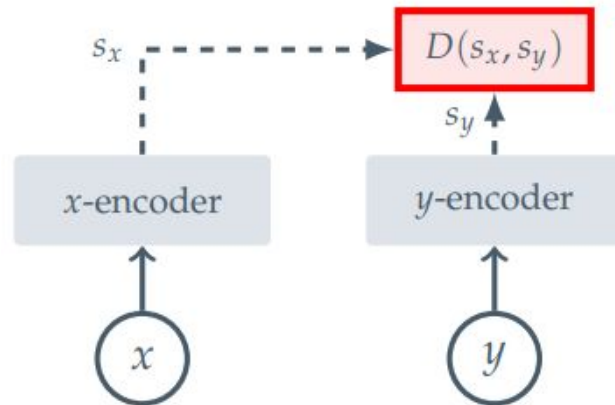
Jointly learns masked speech prediction and denoising

Diversifies data sources for better generalization (vs. other models which are trained on podcast data only)



Joint-Embedding Architecture

Learns to output similar embeddings for compatible inputs x , y and dissimilar embeddings for incompatible inputs

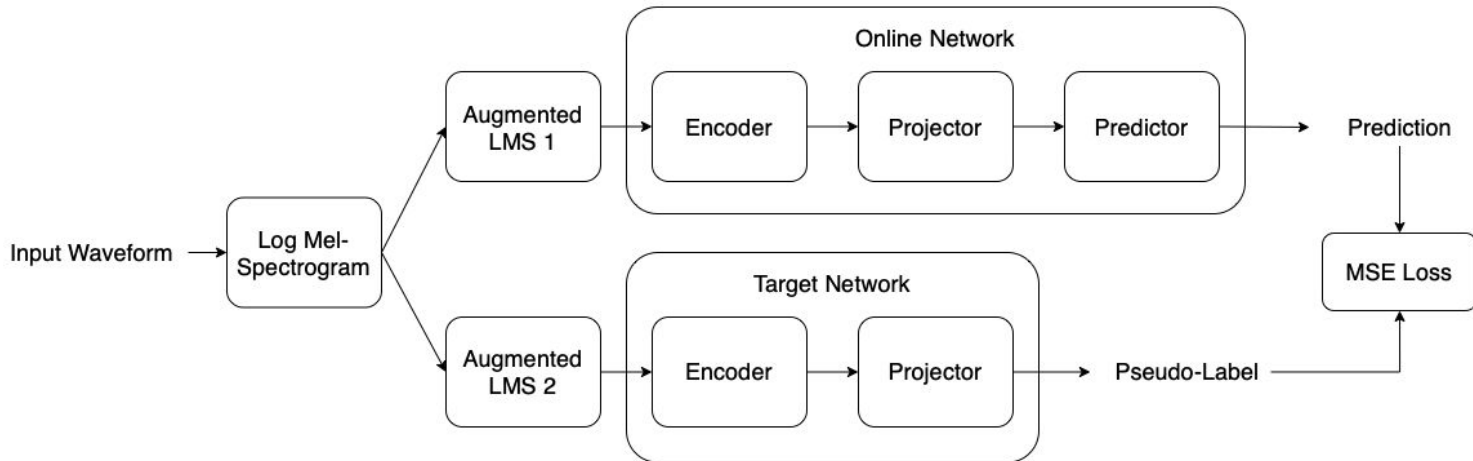


(a) Joint-Embedding Architecture

Joint-Embedding Architecture

BYOL

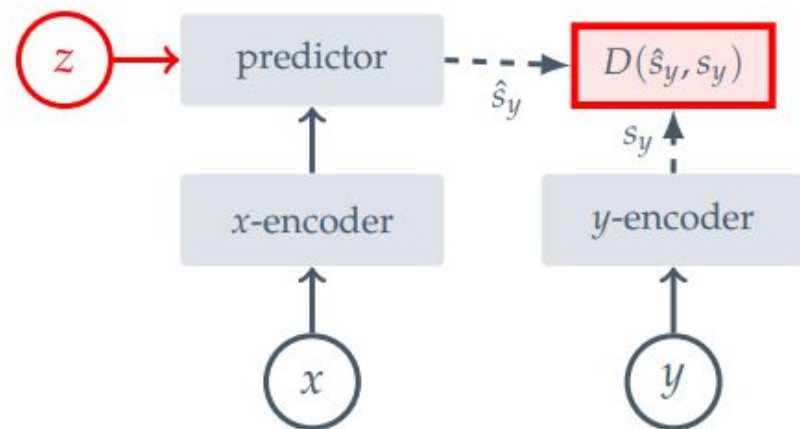
From an augmented view of an input waveform, we train the online network to predict the target network representation of the same image under a different augmented view



Joint-Embedding Predictive Architecture

Learns to predict the embeddings of a signal y from a compatible signal x

Predictor conditioned on additional variables z to facilitate prediction



(c) Joint-Embedding Predictive Architecture

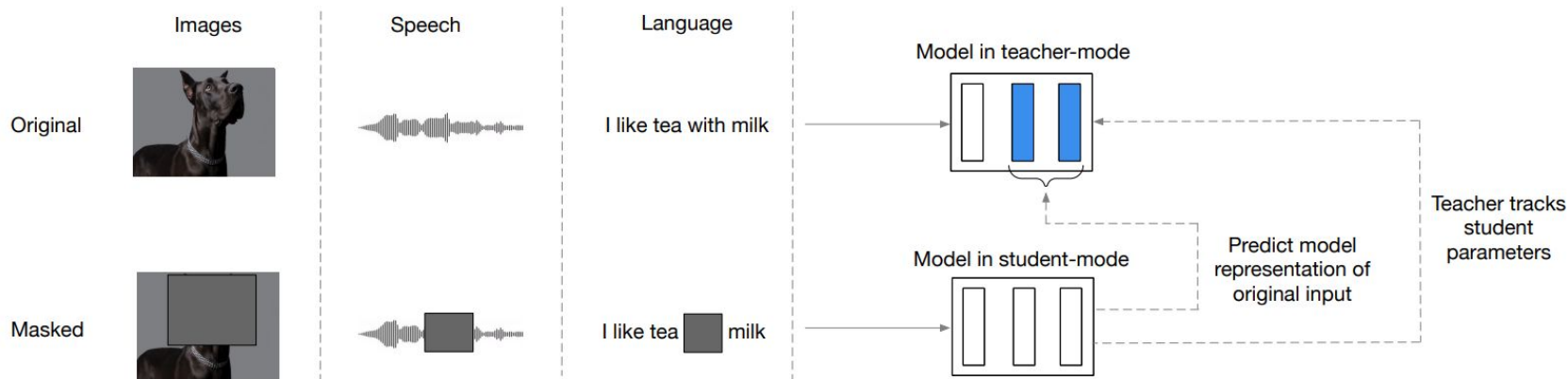
Joint-Embedding Predictive Architecture

Data2vec

Combines masked prediction with the learning of latent target representations

Teacher: builds representation of the full input data which serves as targets in the learning task

Student: encodes a masked version of the input sample with which we predict the full data representations



Why does SSL work?

Can make use of very large amounts of data

Transformer architecture is well suited to learn contextual information

Pre-text tasks allow learning representations that encode general information that is useful for a wide-range of tasks

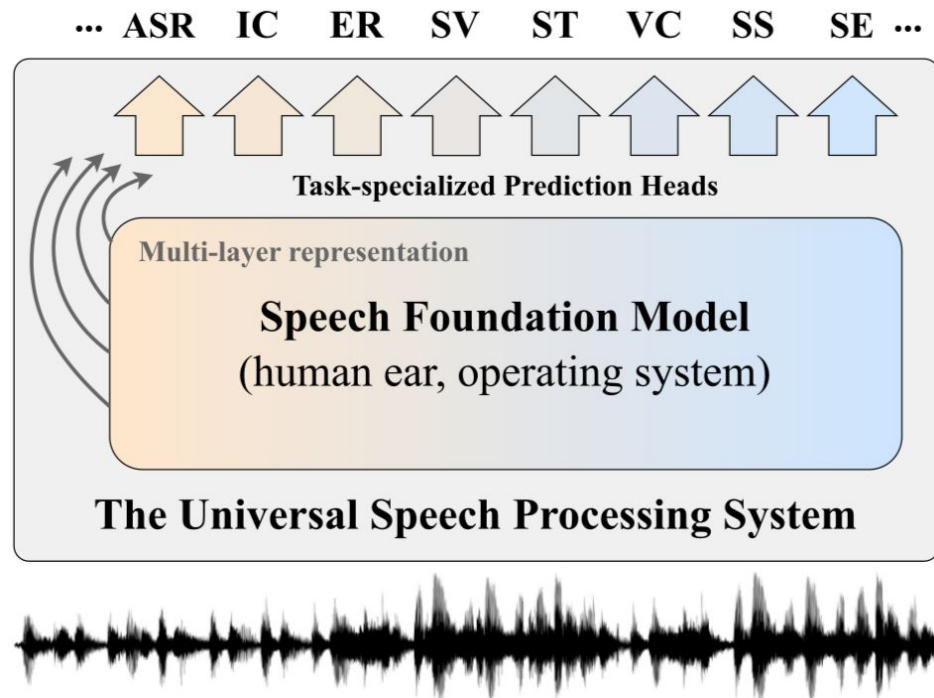
For example, to predict what's missing, the model must grasp the structure of the data

→ The model builds rich internal representations

SSL Model Evaluation

Models are compared to each other by evaluating the representations across a wide variety of tasks

Benchmarks have been set up to perform this evaluation



SUPERB Benchmark

Category	Tasks
Content	Phoneme Recognition (PR), Automatic Speech Recognition (ASR), Out-of-Domain ASR (OOD-ASR), Keyword Spotting (KS), Query-by-Example (QbE)
Speaker	Speaker Identification (SID), Speaker Verification (SV), Speaker Diarization (SD)
Prosody	Emotion Recognition (ER)
Semantics	Intent Classification (IC), Slot Filling (SF), Speech Translation (ST)
Generation	Voice Conversion (VC), Source Separation (SS), Speech Enhancement (SE)

HEAR Benchmark

Category	Tasks
Speech	Speech Command Classification; Emotion Recognition; Language ID; Speaker Count
Music	Pitch Classification; Music Genre Classification; Music/Speech Classification; Music Transcription; Percussion Classification; Instrument Classification
Environmental / Other Audio	Office Sound Detection; Environmental Sound Classification; General Audio Tagging; Beehive Condition Classification; Gunshot Triangulation; Imitated Sound Type Classification

Leaderboards

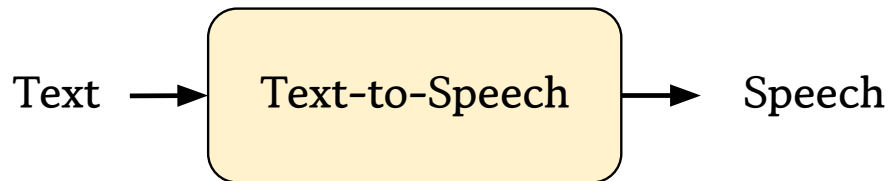
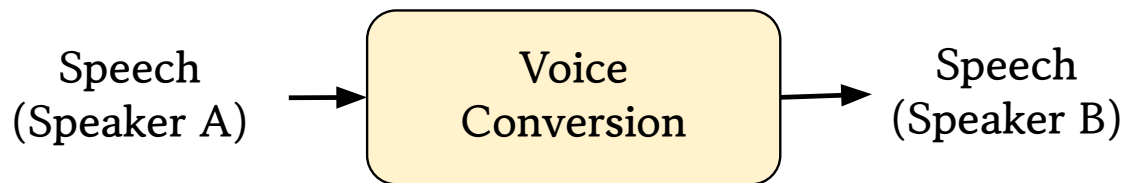
<https://superbenchmark.github.io/#/leaderboard>

<https://hearbenchmark.com/hear-leaderboard.html>

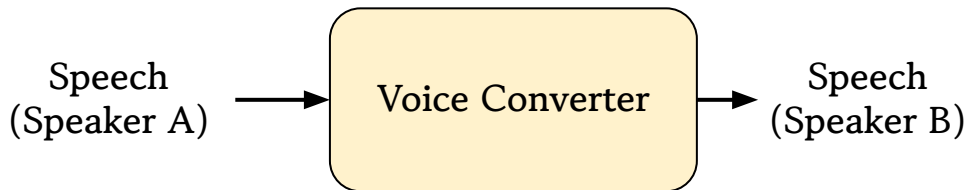
SSL for Synthesis

In a synthesis task, we generate new speech waveforms given inputs

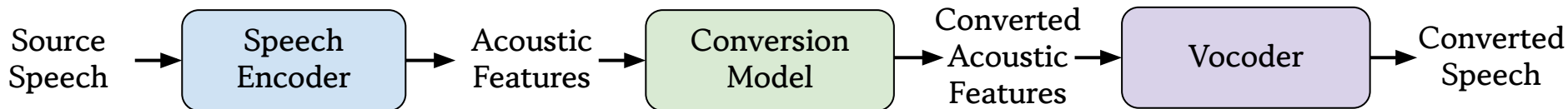
Examples:



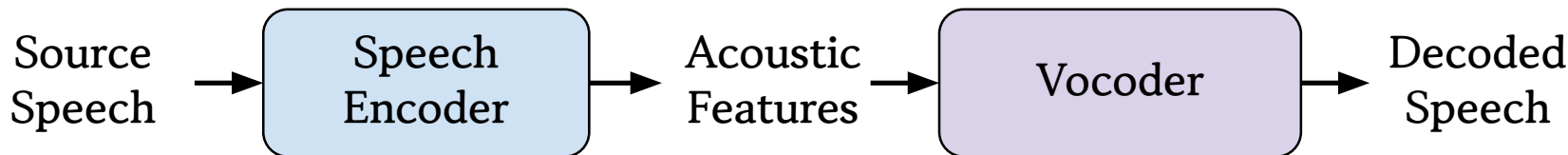
Voice Conversion





Typical Pipeline:



Copy-Synthesis



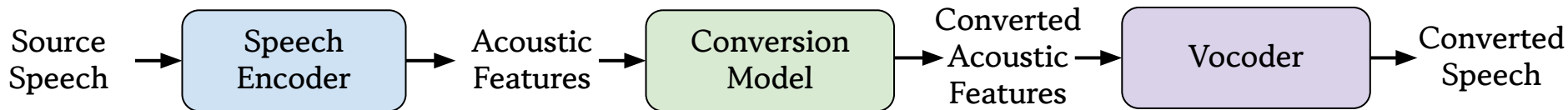
	Features	Vocoders
Signal Processing-based Examples	Mel-spectrogram, MFCCs	Griffin-Lim, WORLD 
Neural Examples	SSL Representations	HiFiGAN, BigVGAN 

Voice Conversion

How can we leverage the strengths of SSL representations for voice conversion?

SSL models encode speech into a sequence of frames, each corresponding to a window of 25ms of speech, with a 20ms hop between windows.

Property: frames are linearly close if they contain similar phonetic, linguistic information, even if they are spoken by very different voices





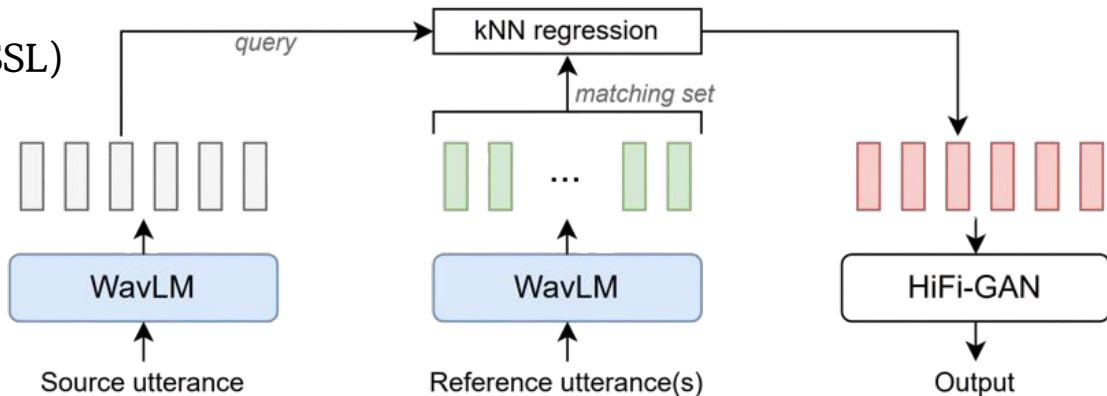
Voice Conversion: kNN-VC

Any-to-Any VC method

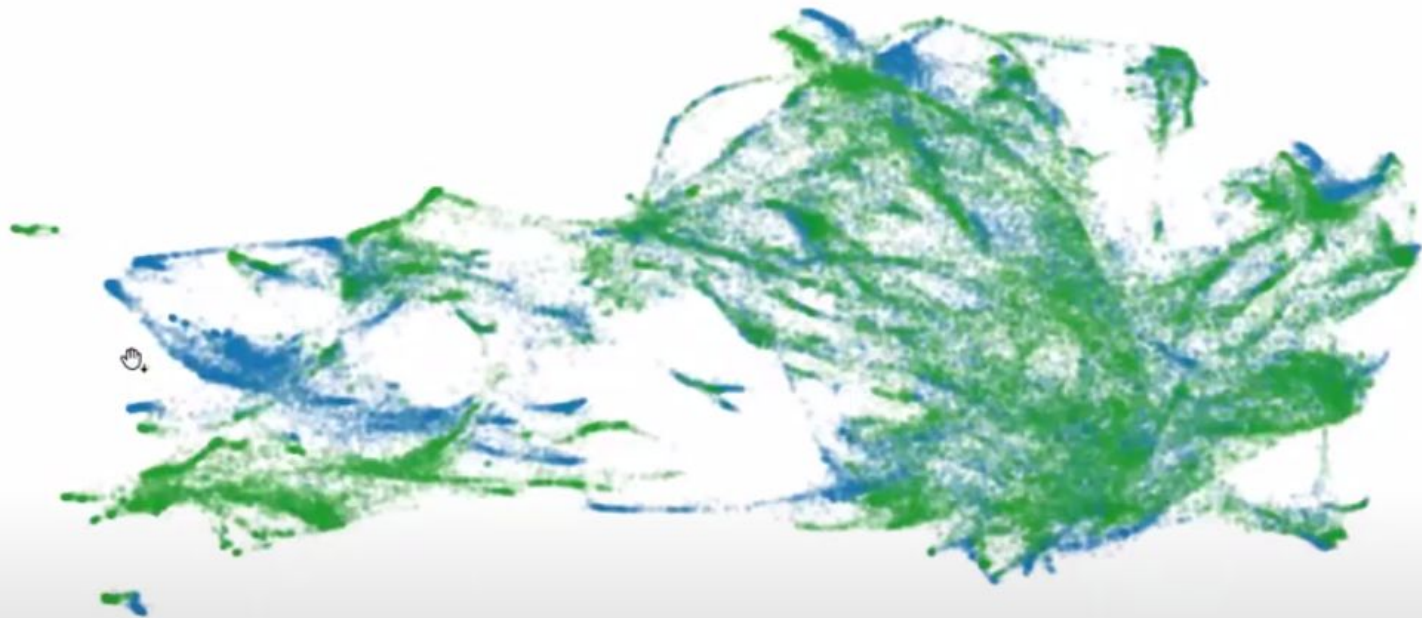
Leveraging self-supervised learning (SSL) features for zero-shot conversion

Requires no training

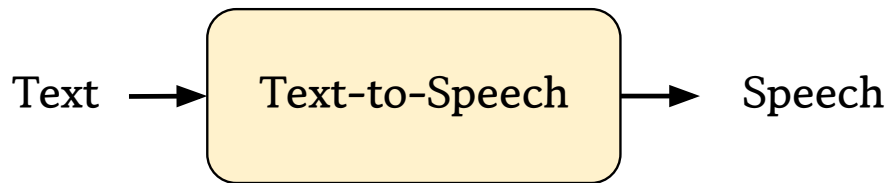
<i>Source</i>	<i>Target</i>	<i>Converted</i>
		



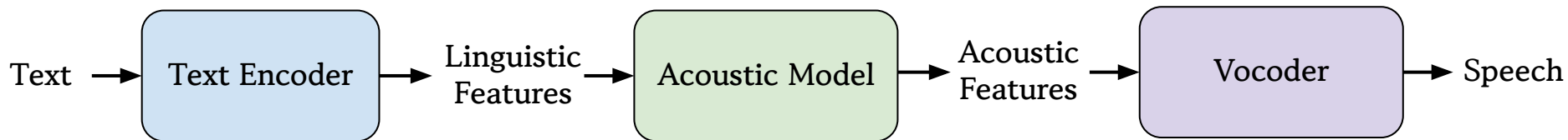
3D Projection of WavLM Features for Two Speakers



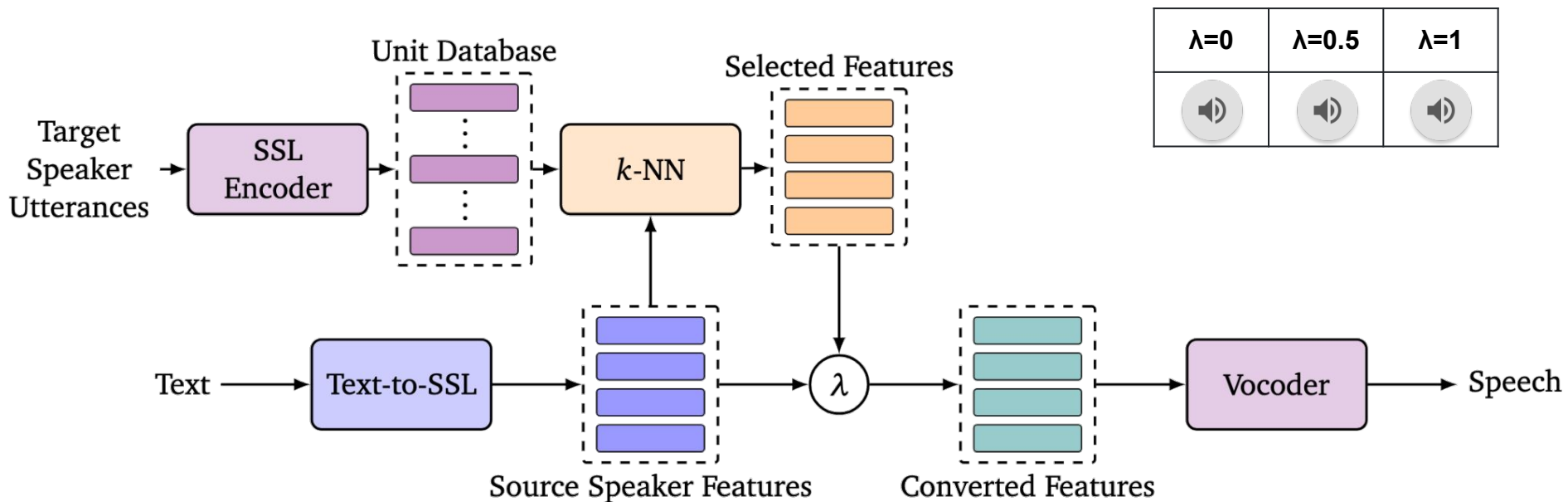
Text-to-Speech






Typical Pipeline:



kNN-TTS



$\lambda=0$	$\lambda=0.5$	$\lambda=1$
		

$$y_{\text{converted}} = \lambda y_{\text{selected}} + (1 - \lambda) y_{\text{source}}$$

Summary

Self-supervised learning enables learning useful general representations from unlabeled data

Evaluations show that SSL representations enable performance improvements for wide range of tasks;

→ Especially whenever labeled data is scarce

As seen in Voice Conversion for example, good representations can simplify task-specific models