
Self-Supervised Learning

For Speech

Andy T. Liu
2020/06/03

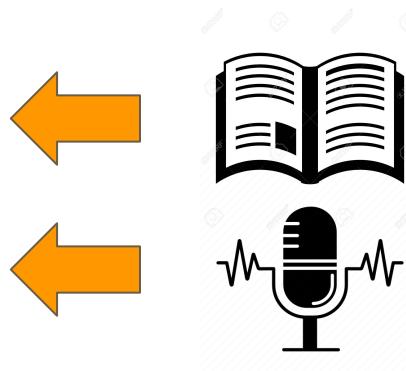
Overview

What is Self-Supervised Learning?

An analogy

How do Infants Learn?

Can Machine do the same?



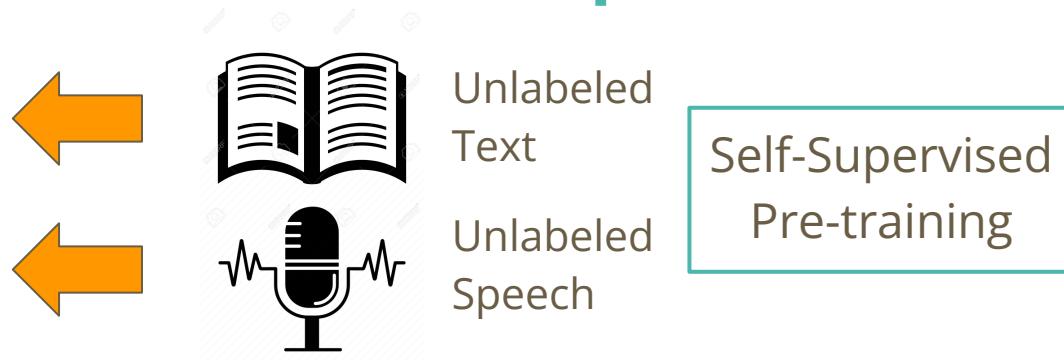
Looks at a lot of books!



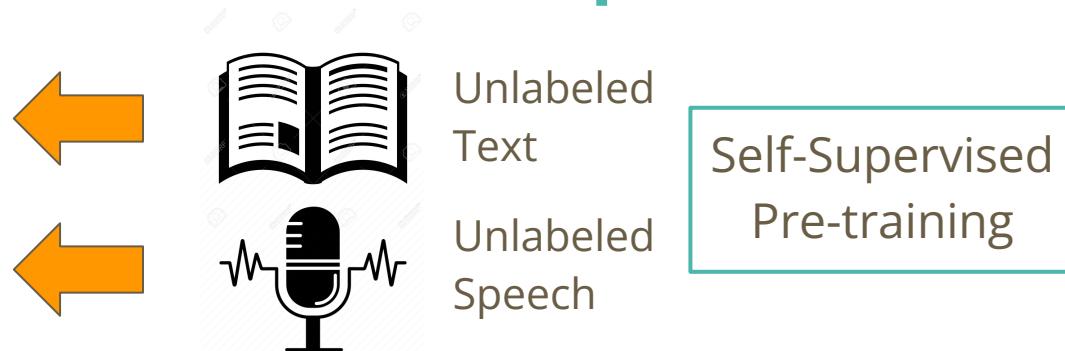
Listens to a lot of conversations!

Then their parents
teach them.

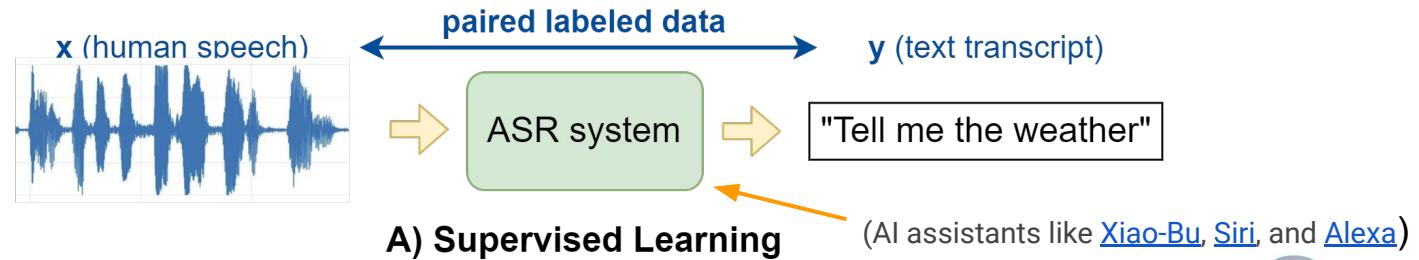
How do Infants Learn? Yes! Self-Supervised Learning



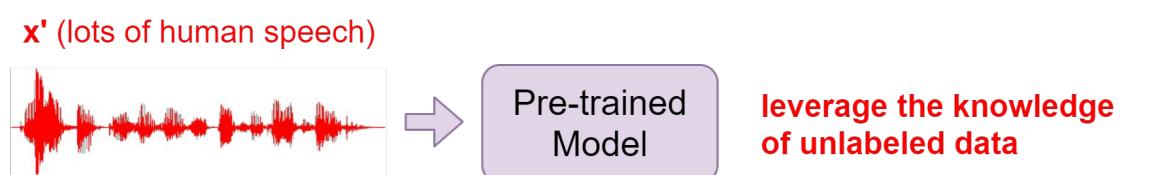
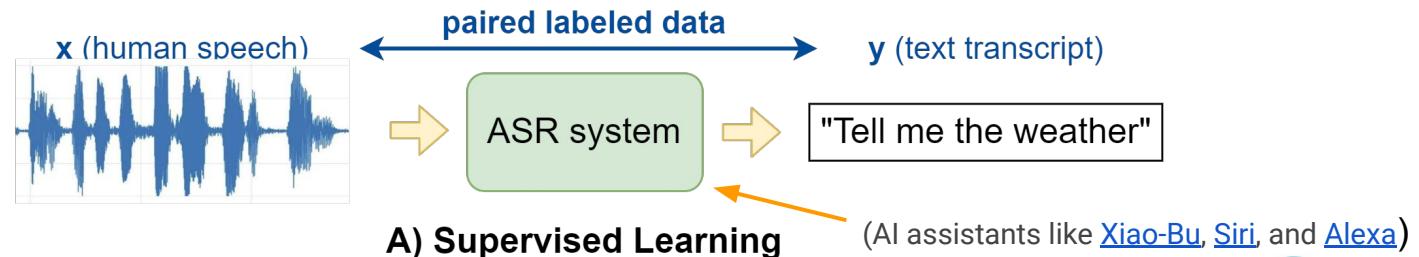
How do Infants Learn? Yes! Self-Supervised Learning



Self-Supervised Learning for Speech



Self-Supervised Learning for Speech



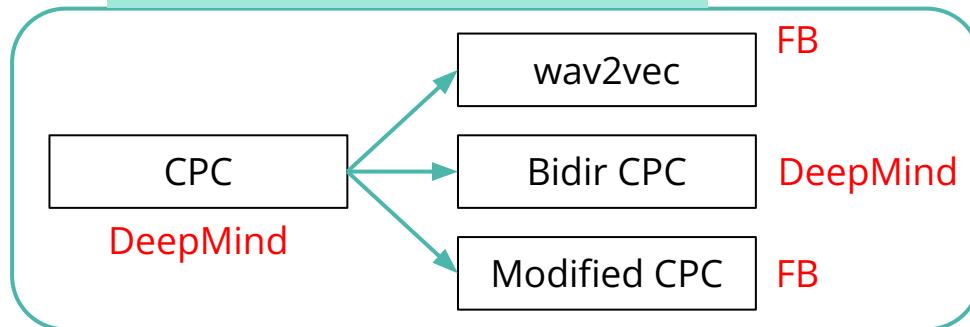
Pre-trained models
are evaluated on
downstream tasks.

B) Self-Supervised Learning for Improving Supervised Systems

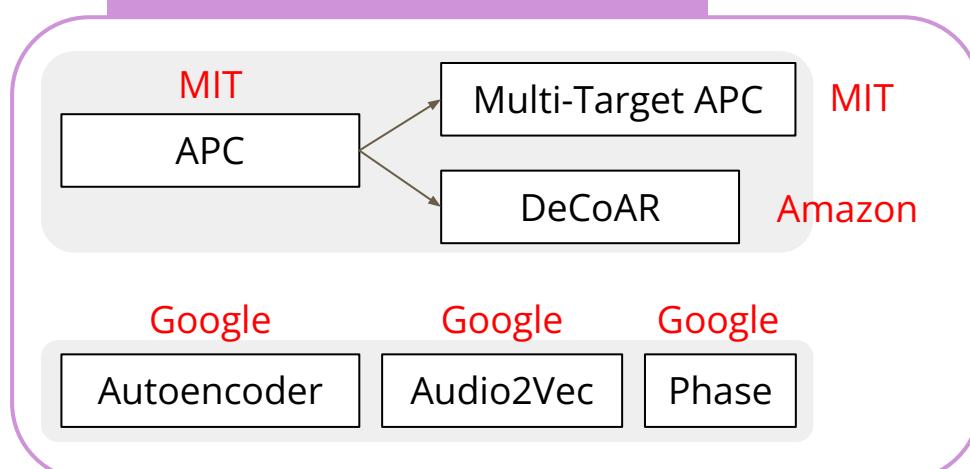
A Broad Introduction of All Recent Methods

A Broad Introduction

Contrastive Predictive Losses

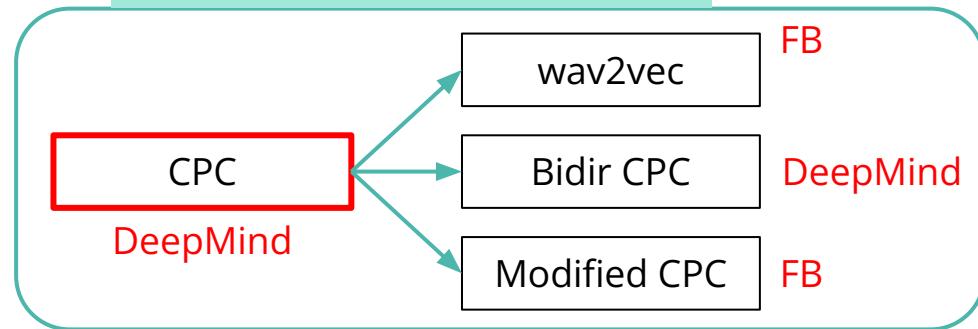


Reconstruction Losses



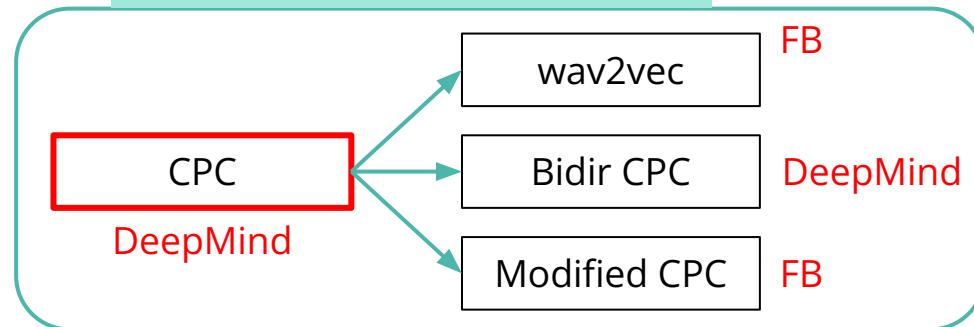
A Broad Introduction

Contrastive Predictive Losses



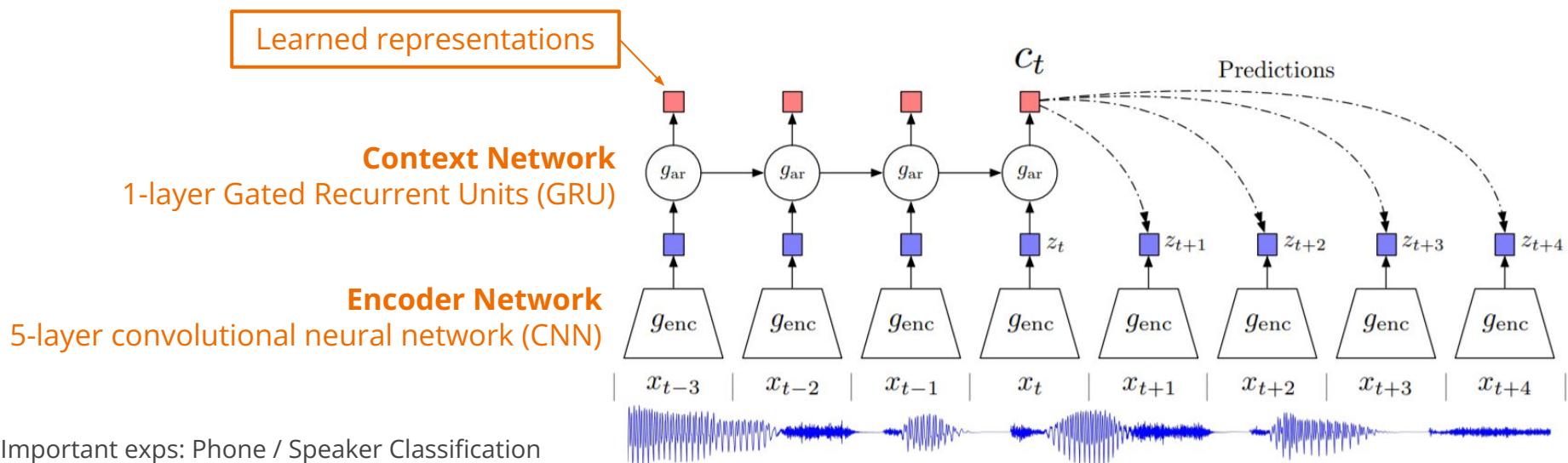
A Broad Introduction

Contrastive Predictive Losses



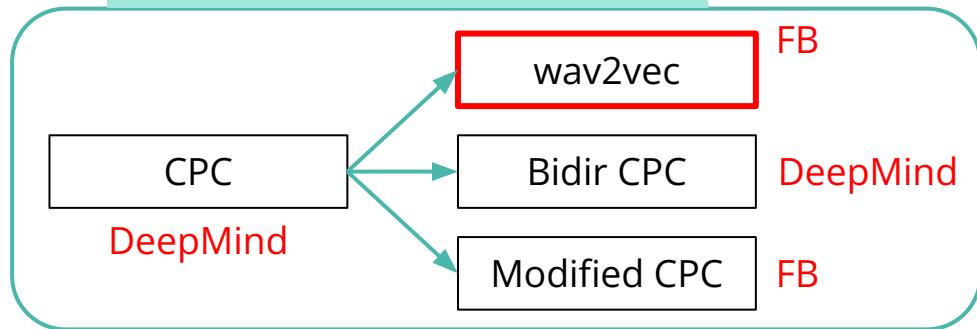
Intuition:

Pulls temporally nearby representations closer and pushes temporally distant ones further.



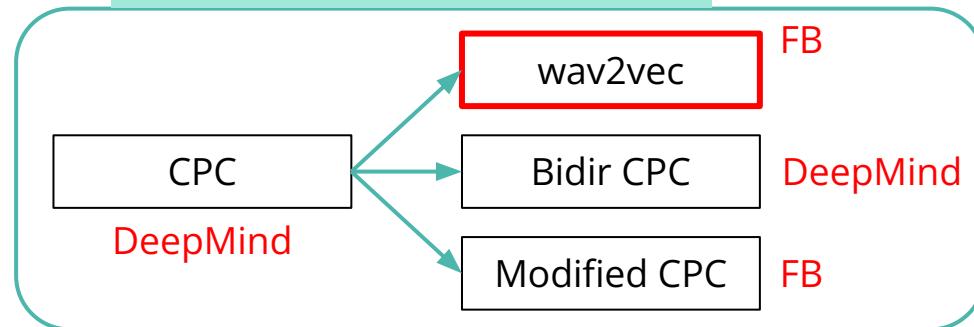
A Broad Introduction

Contrastive Predictive Losses



A Broad Introduction

Contrastive Predictive Losses

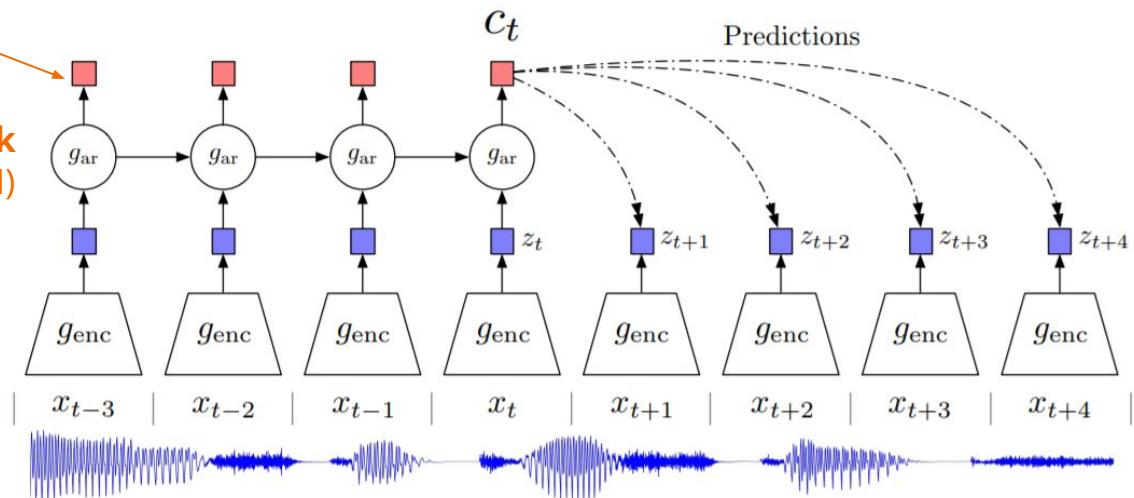


Contribution:
self-supervised pre-training is shown to improve supervised **ASR**

Used as input for ASR models,
replace acoustic features

Context Network
9-layer convolutional neural network (CNN)

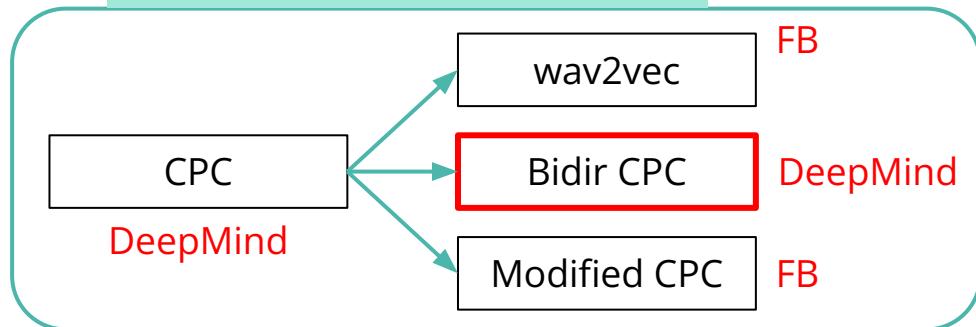
Encoder Network
5-layer convolutional neural network (CNN)



Important exps: ASR on WSJ / TIMIT

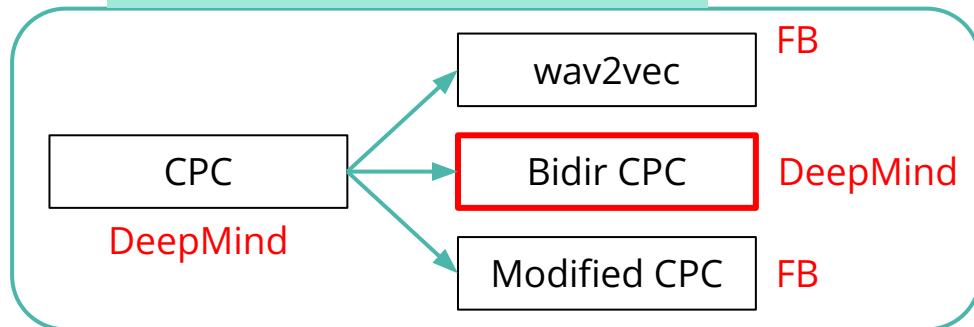
A Broad Introduction

Contrastive Predictive Losses



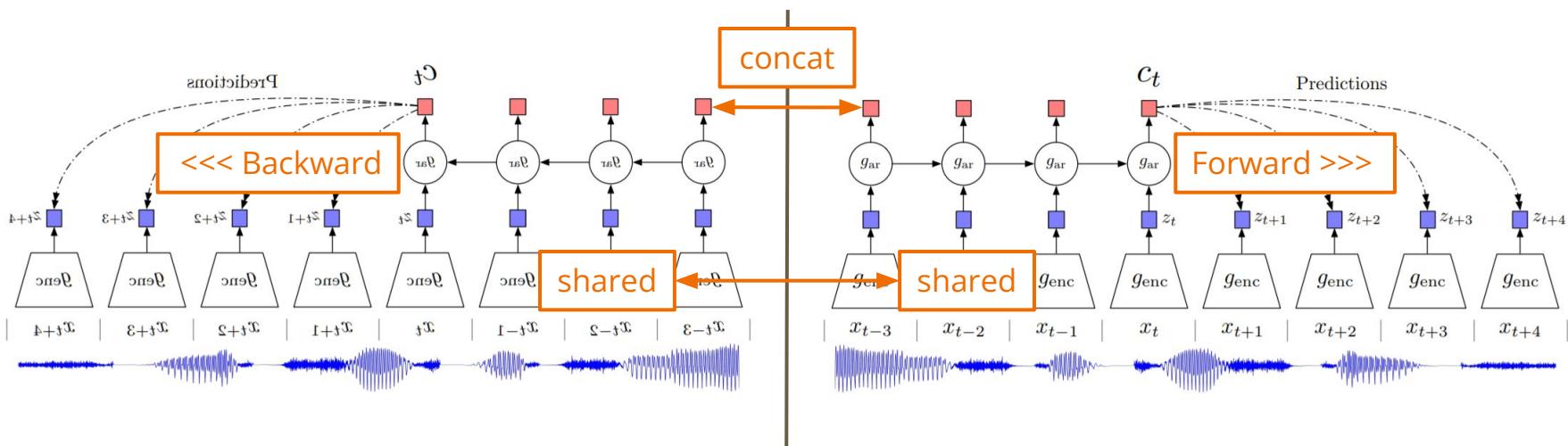
A Broad Introduction

Contrastive Predictive Losses



Contribution:
bidirectional context + ASR

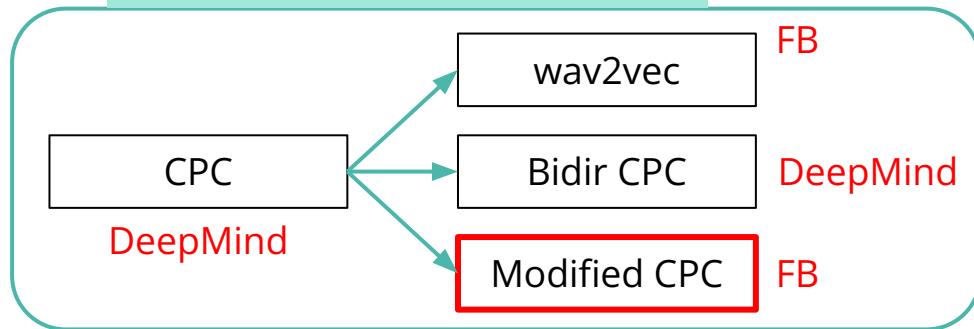
learning representations from large amount of unlabeled data (8000 hrs) can provide improvements for out-of-domain transfer (different datasets / cross-lingual).



Important exps: ASR on LibriSpeech

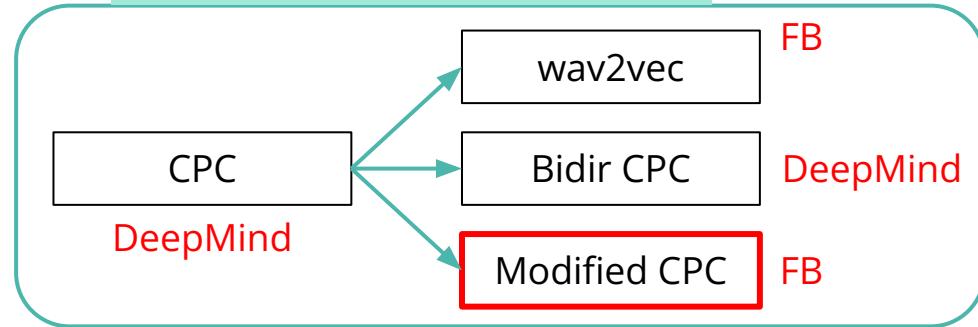
A Broad Introduction

Contrastive Predictive Losses



A Broad Introduction

Contrastive Predictive Losses

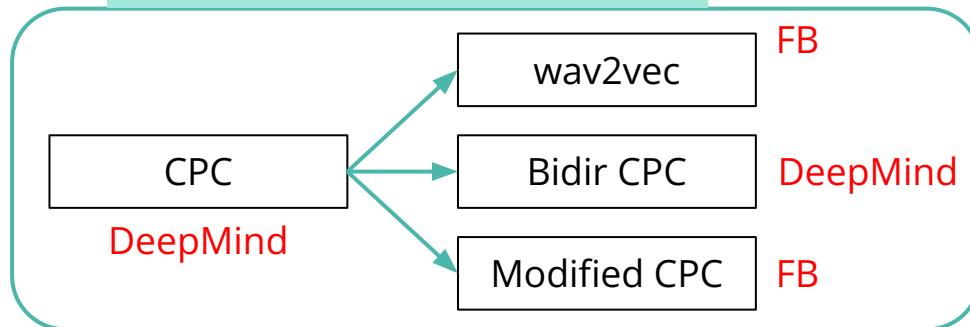


Contributions:

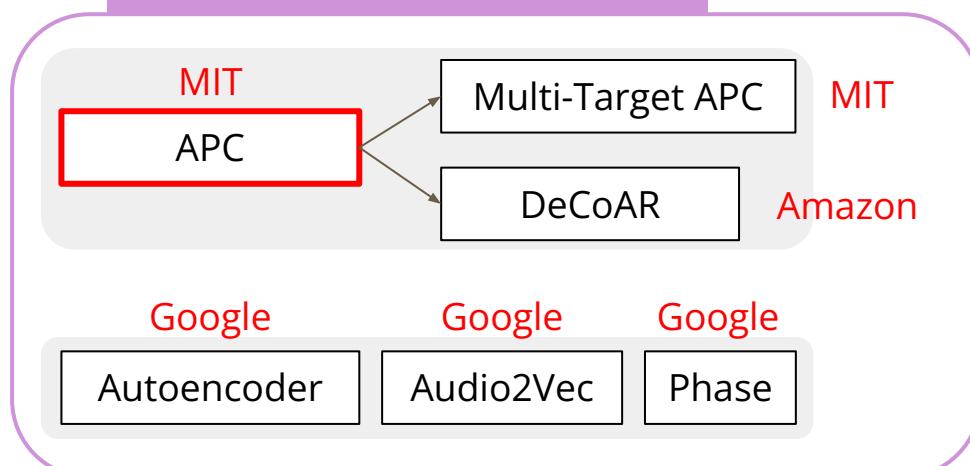
- 1) changing the batch normalization to channel-wise normalization
- 2) replace the linear prediction layer to a Transformer layer
- 3) and replacing the context network of GRUs with Long Short-Term Memory (LSTM) cells

A Broad Introduction

Contrastive Predictive Losses



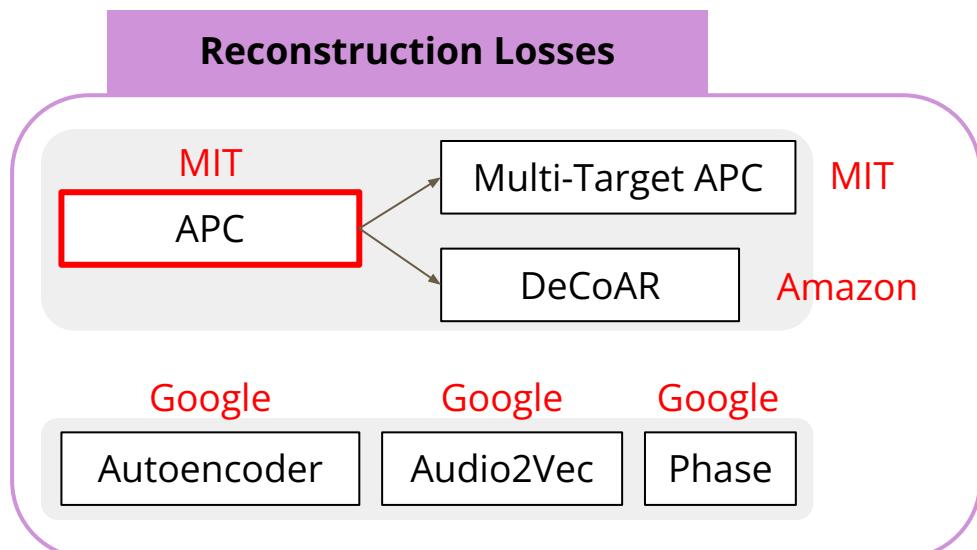
Reconstruction Losses



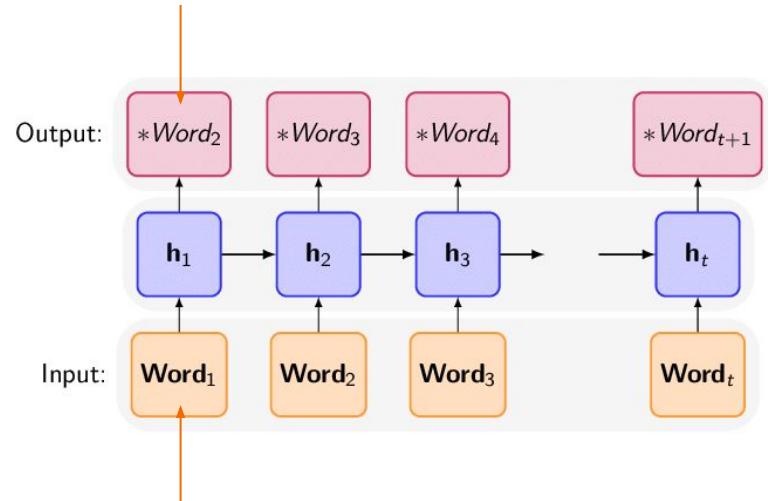
A Broad Introduction

Intuition:

Speech Version of a RNN Language Model

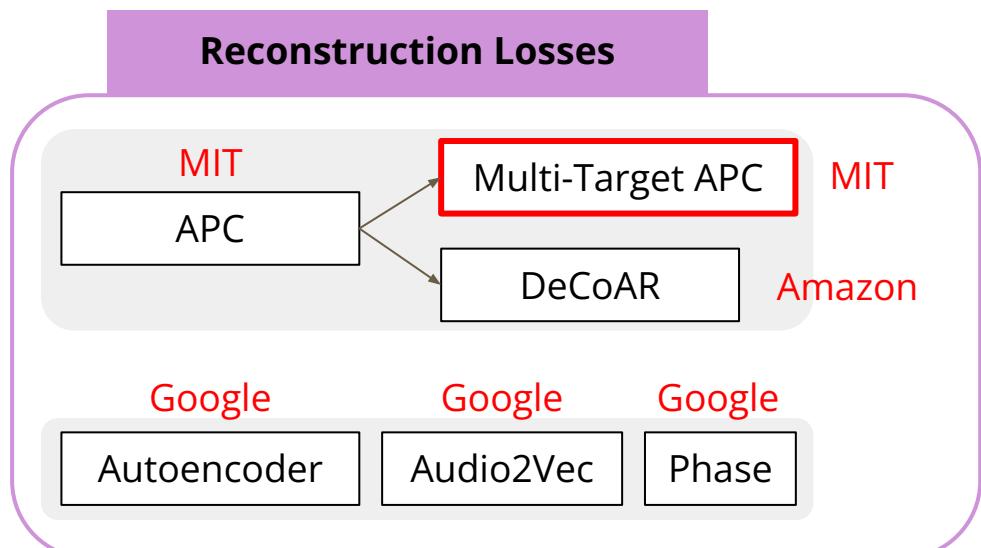


Change the softmax layer to regression layer for reconstruction



Instead of operating on word tokens, change them to acoustic frames

A Broad Introduction



Important exps: Phone Classification, ASR on WSJ
They use settings that are not conventional.

A Broad Introduction

Intuition:

The APC objective is extended to bidirectional.

An auxiliary RNN is used to refresh current hidden states with the knowledge learned in the past, allowing the model to remember more from the past.

Reconstruction Losses

MIT

APC

Multi-Target APC

MIT

Amazon

Google

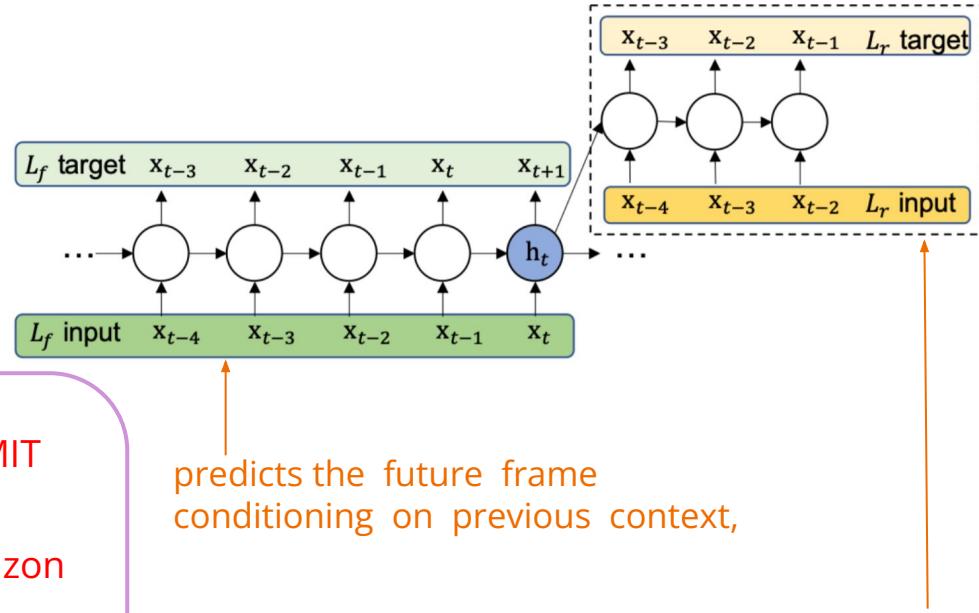
Google

Google

Autoencoder

Audio2Vec

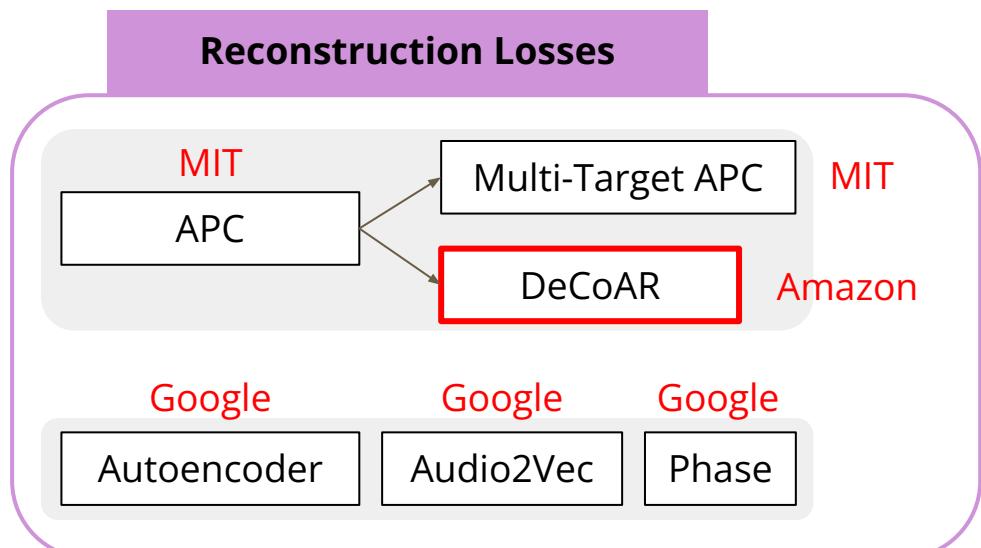
Phase



predicts the future frame
conditioning on previous context,

but also predicts the past
memory through reconstruction.

A Broad Introduction

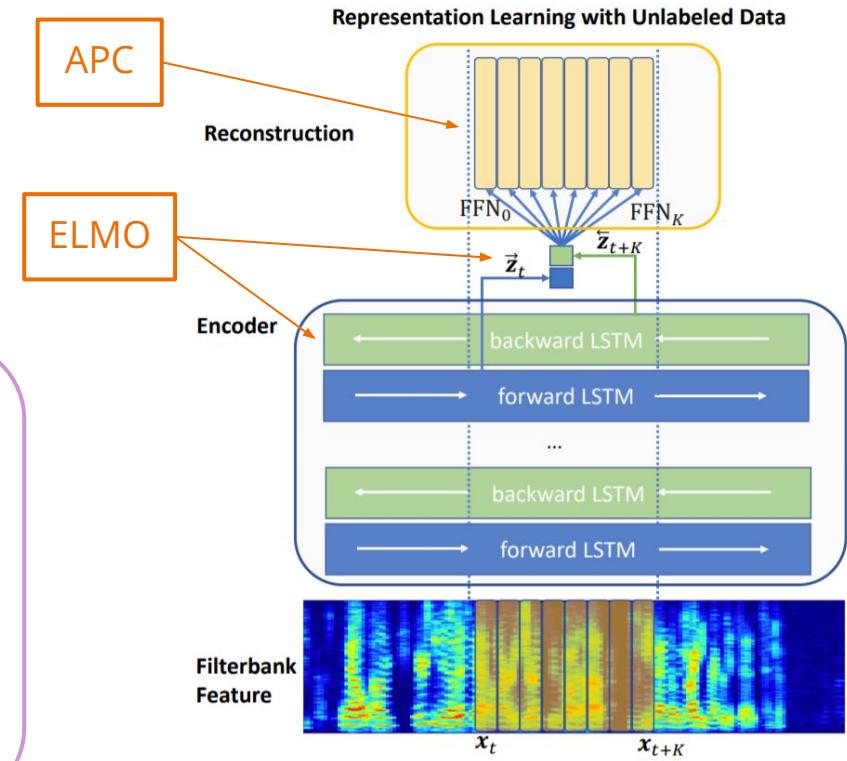
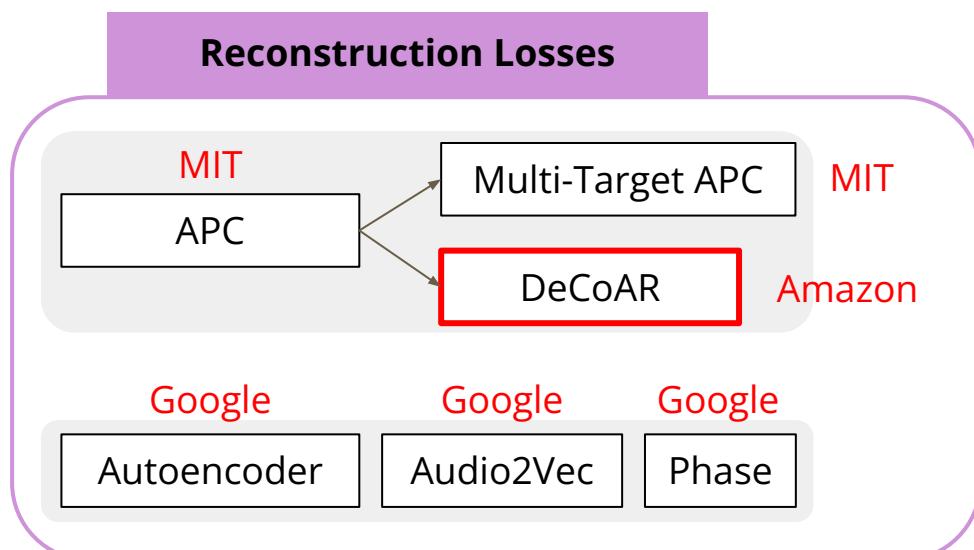


A Broad Introduction

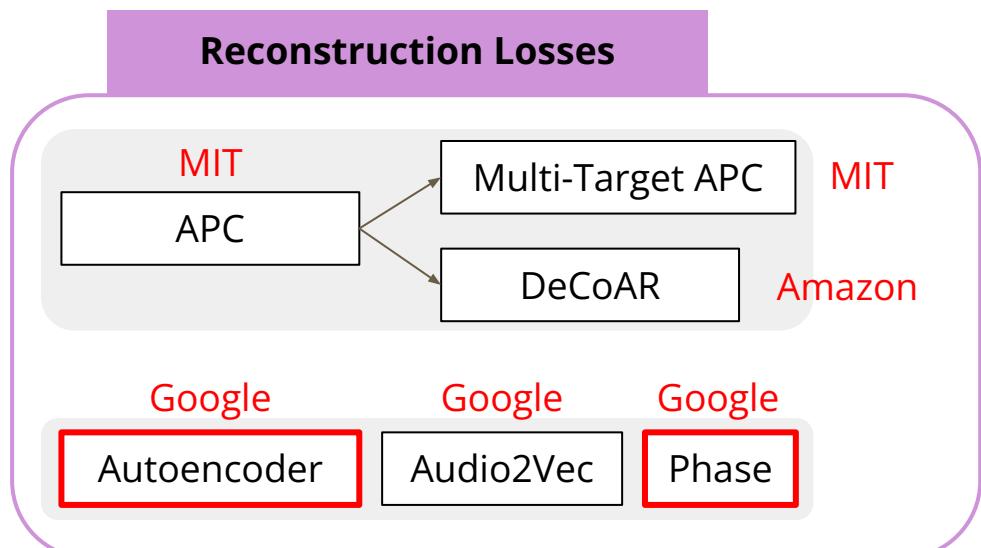
Intuition:

Deep Contextualized Acoustic Representations

Combining the bidirectionality of ELMo and the reconstruction objective of APC. Reconstruction loss is summed over all possible slices in the entire sequence.

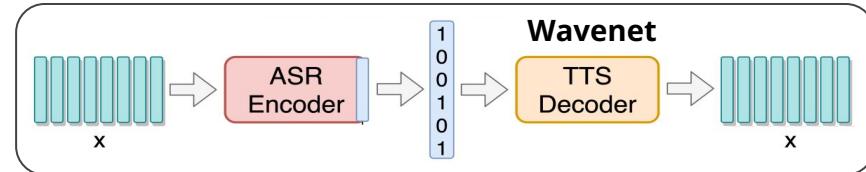


A Broad Introduction

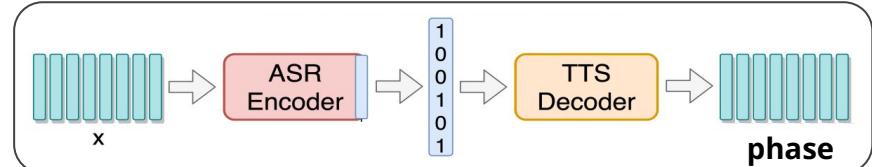


A Broad Introduction

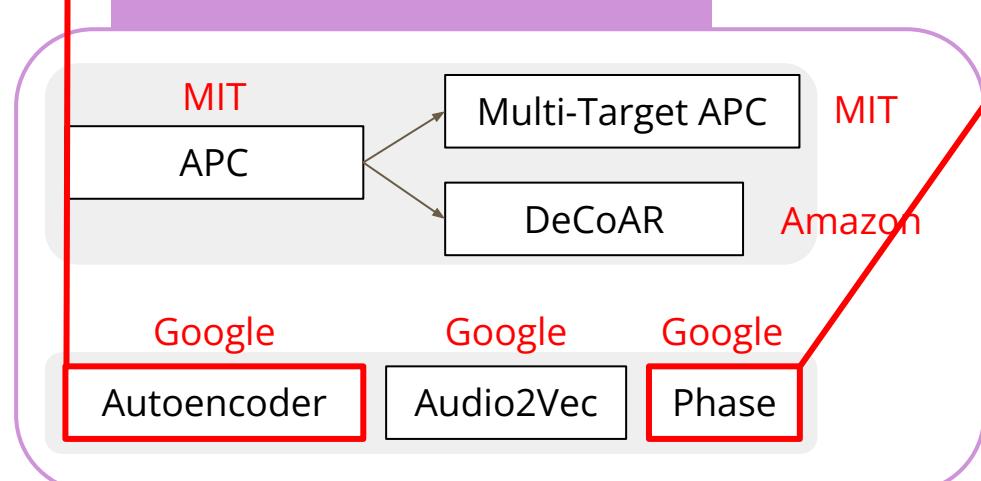
Important exps: Phone Classification



Important exps: Linear Classifications

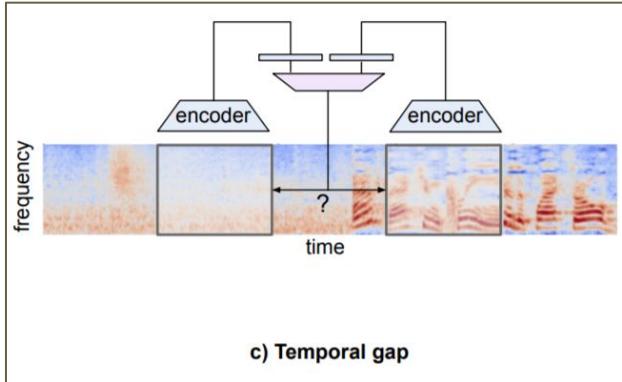


Reconstruction Losses



Intuition:
Very similar to the VC structure, learn through autoencoder bottleneck and reconstruction

Important exps: Linear Classifications



Intuition:
Audio version of
Word2Vec

Reconstruction Losses

MIT

APC

Multi-Target APC

MIT

DeCoAR

Amazon

Google

Autoencoder

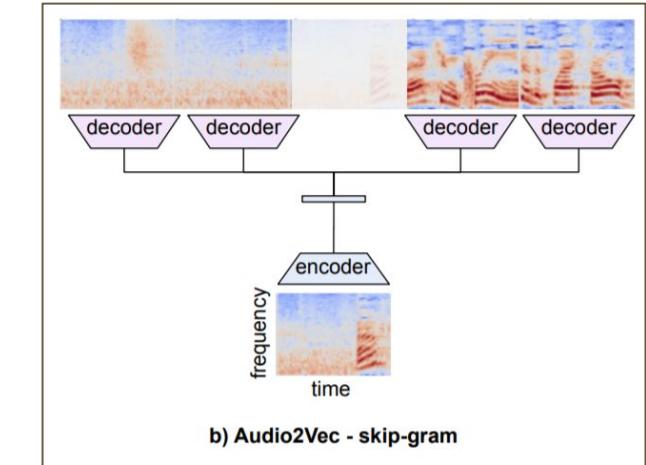
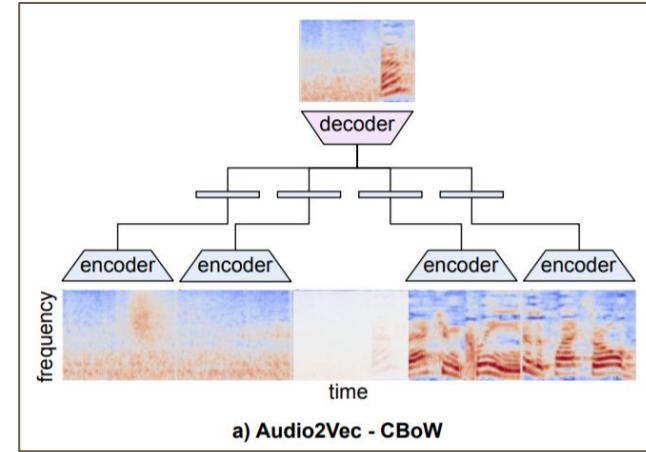
Google

Audio2Vec

Google

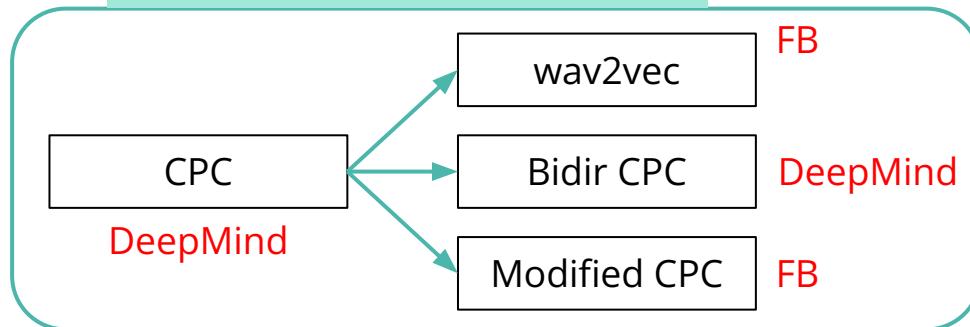
Phase

A Broad Introduction

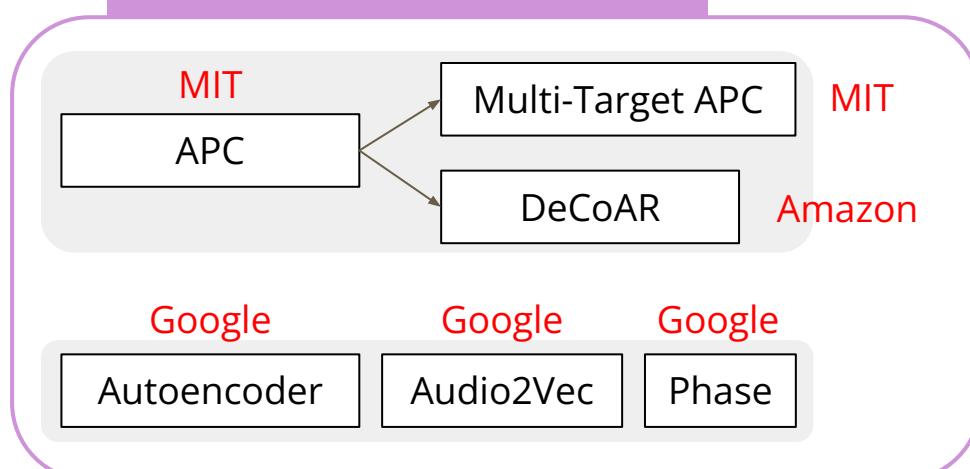


A Broad Introduction

Contrastive Predictive Losses

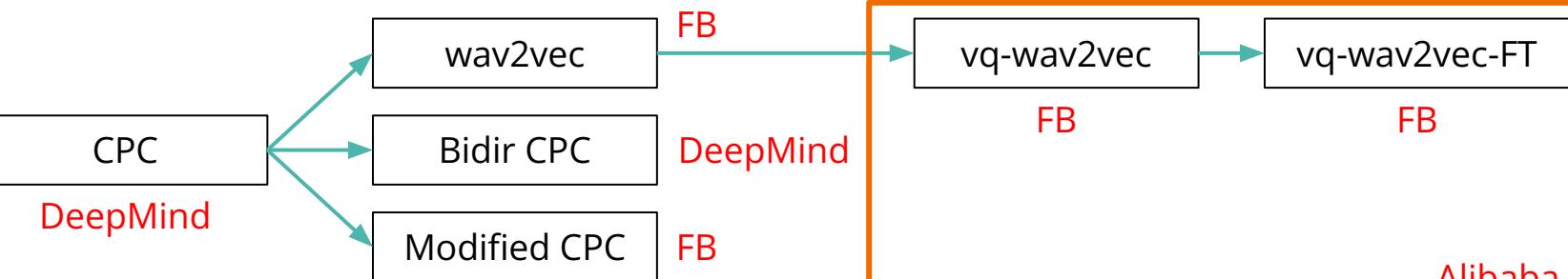


Reconstruction Losses



A Broad Introduction

Contrastive Predictive Losses



Reconstruction Losses

MIT

APC

Multi-Target APC

MIT

DeCoAR

Amazon

Google

Autoencoder

Google

Audio2Vec

Google

Phase

BERT-Style

SLU BERT

Speech XLNet

Tencent

Amazon

Speech Encoder

Didi Chuxing

MPC

Ours

Mockingjay

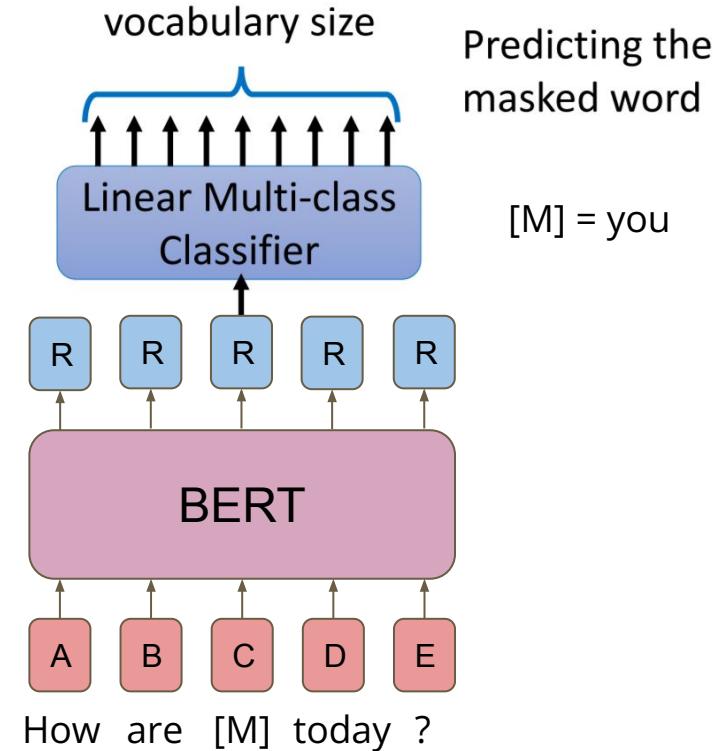
BERT (Bidirectional Encoder Representations from Transformers)

- Achieved **11 SOTA** when published.
- A technique for NLP pre-training developed by Google.



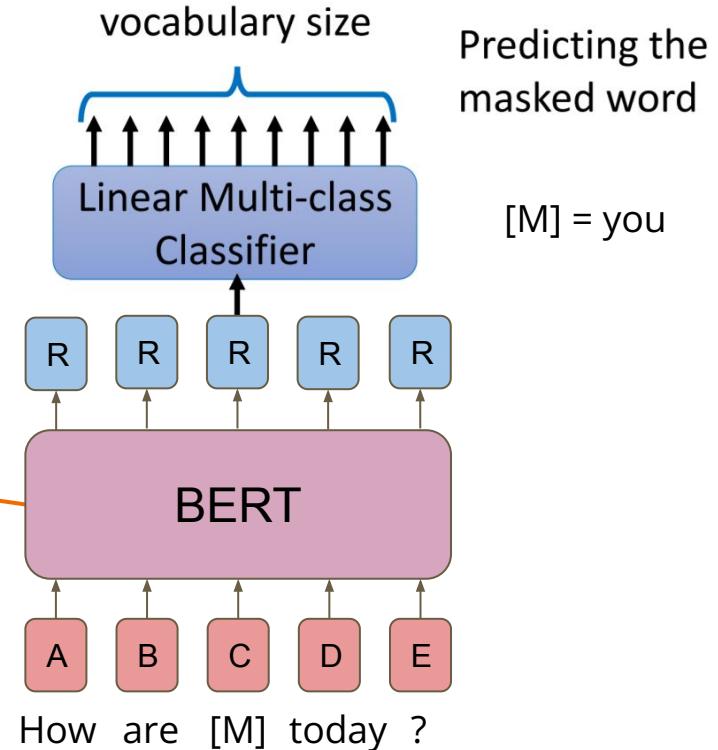
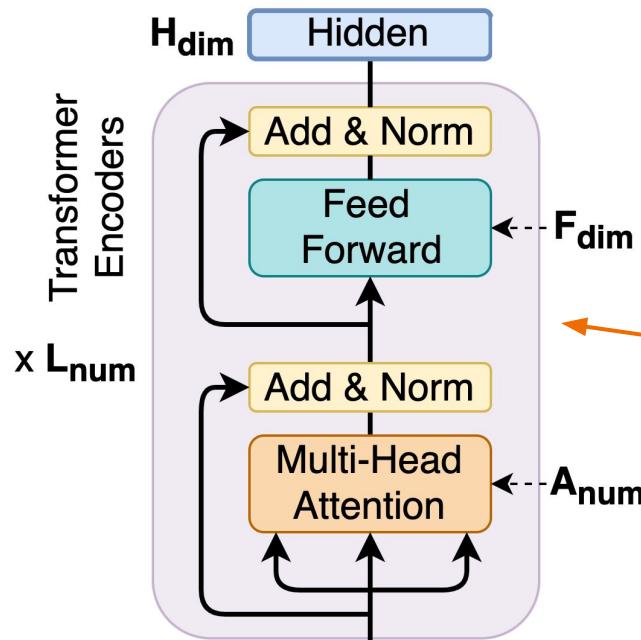
BERT (Bidirectional Encoder Representations from Transformers)

- Achieved **11 SOTA** when published.
- A technique for NLP pre-training developed by Google.
- Learn contextualized repr trough Masked LM:

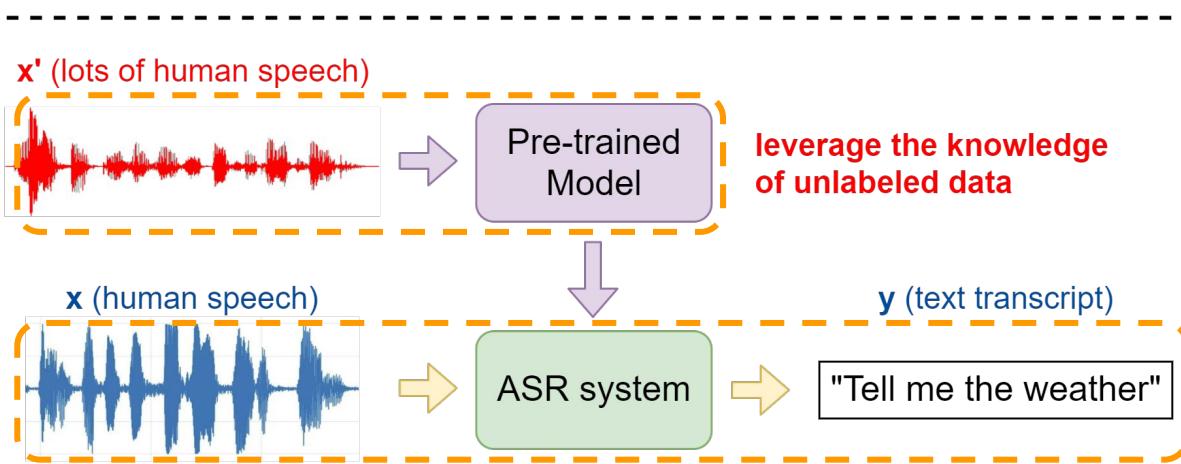
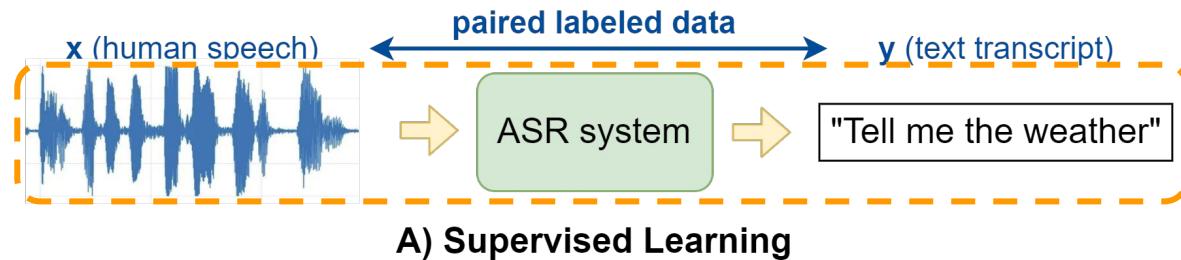


BERT (Bidirectional Encoder Representations from Transformers)

- Achieved **11 SOTA** when published.
- A technique for NLP pre-training developed by Google.
- Learn contextualized repr trough Masked LM:



Recall: Self-Supervised Learning for Speech



B) Self-Supervised Learning for Improving Supervised Systems

Self-Supervised Learning: BERT

P (passage)
Q (question)



paired labeled data

A (answer)

QA Model

"Covid-19 outbreaks
in year 2020"

A) Supervised Learning

x' (lots of raw text)



BERT

leverage the knowledge
of unlabeled data

P (passage)
Q (question)



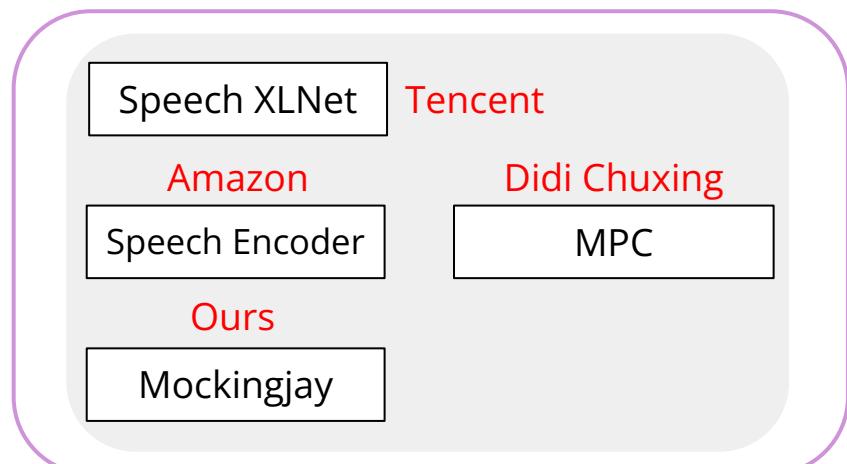
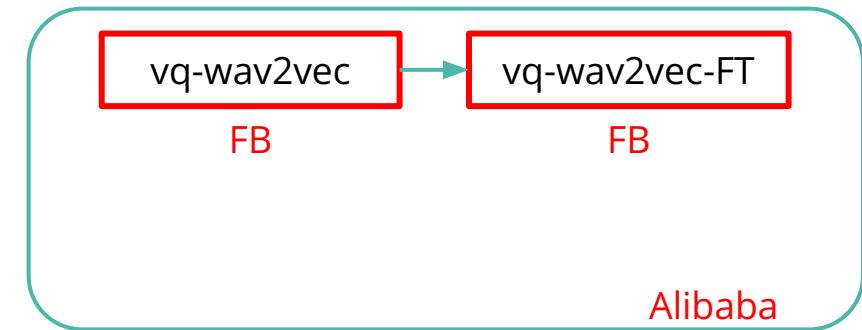
QA Model

A (answer)

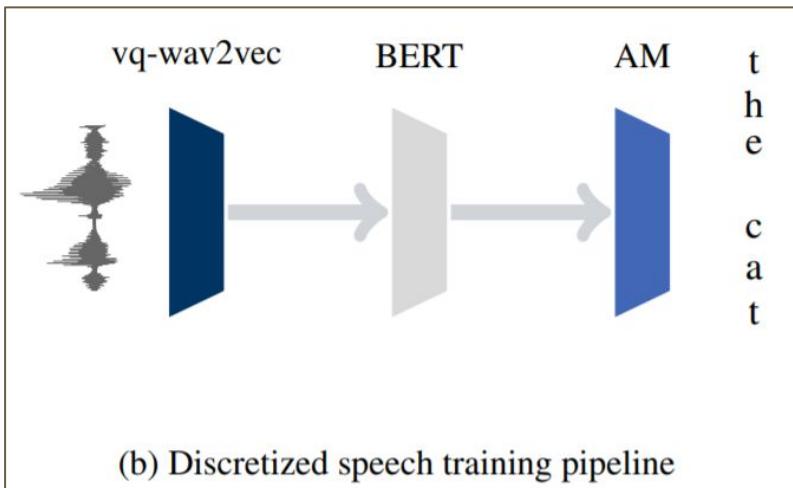
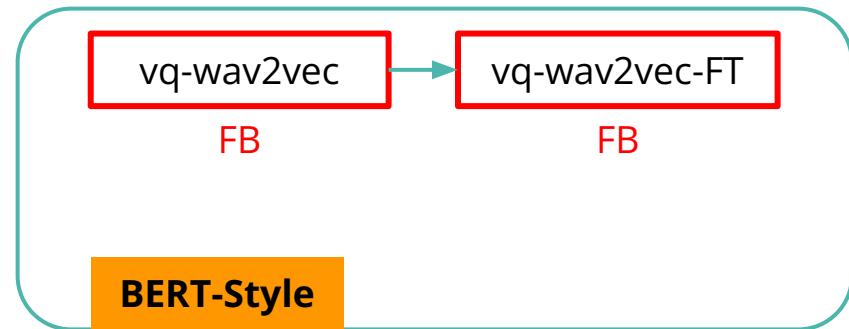
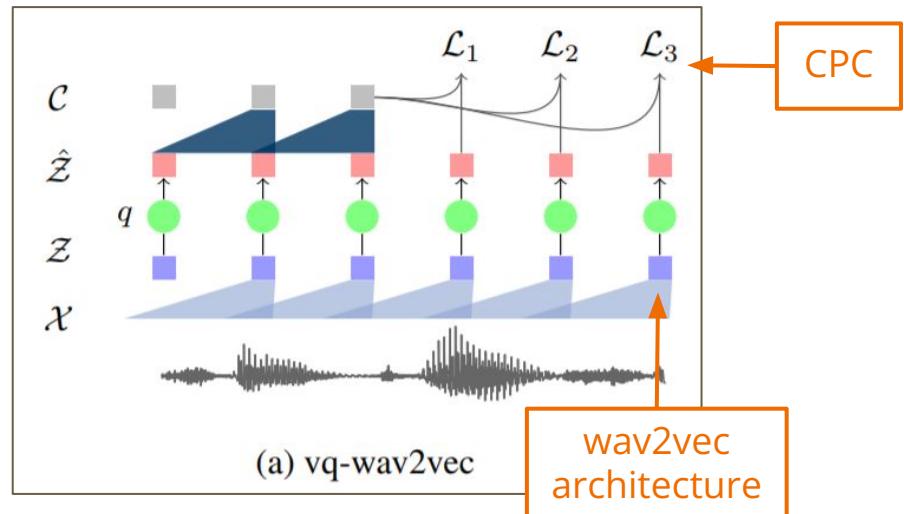
"Covid-19 outbreaks
in year 2020"

B) Self-Supervised Learning for Improving Supervised Systems

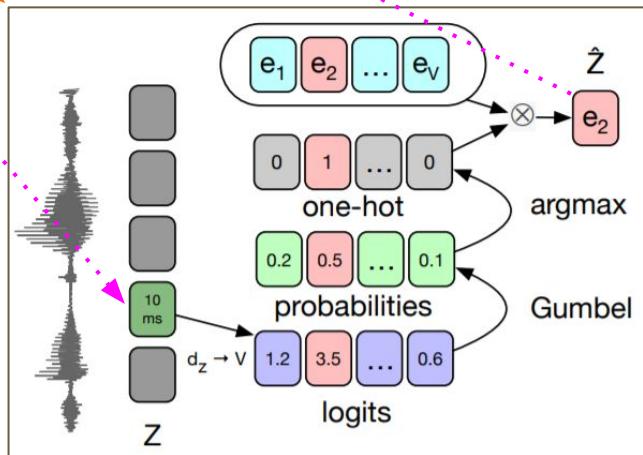
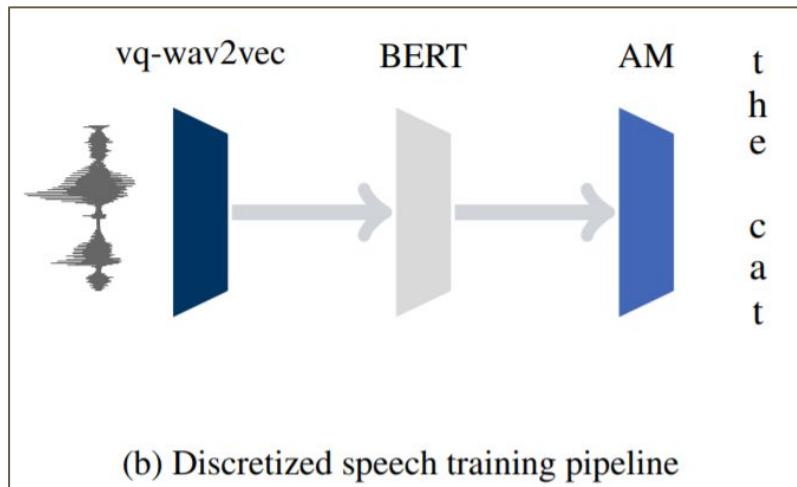
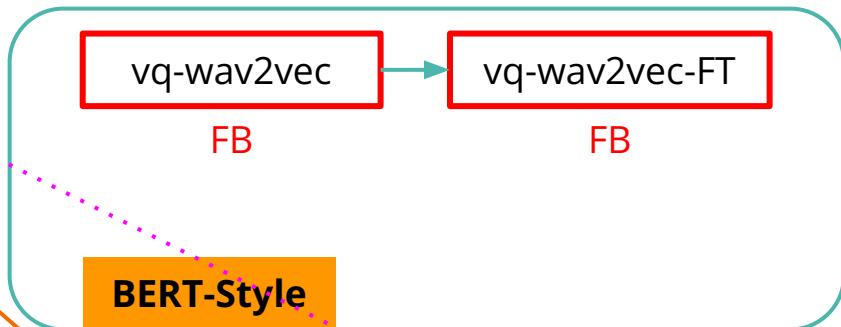
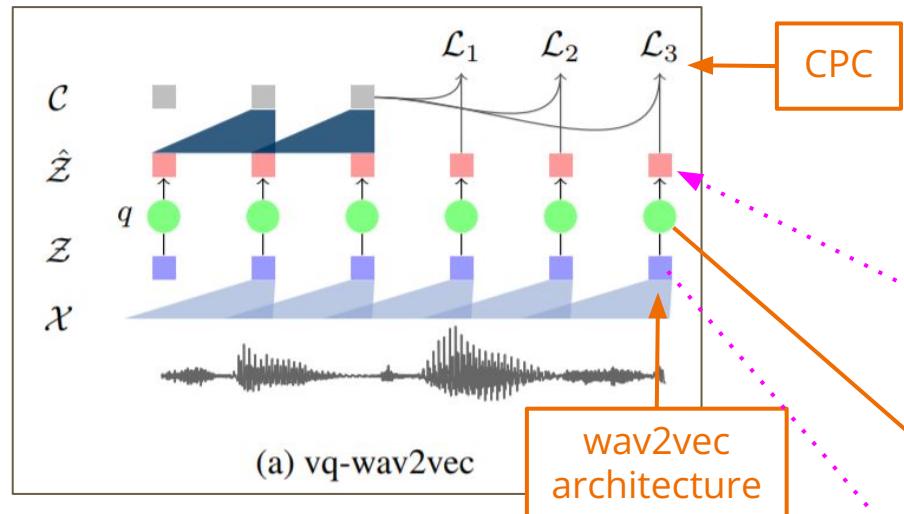
A Broad Introduction



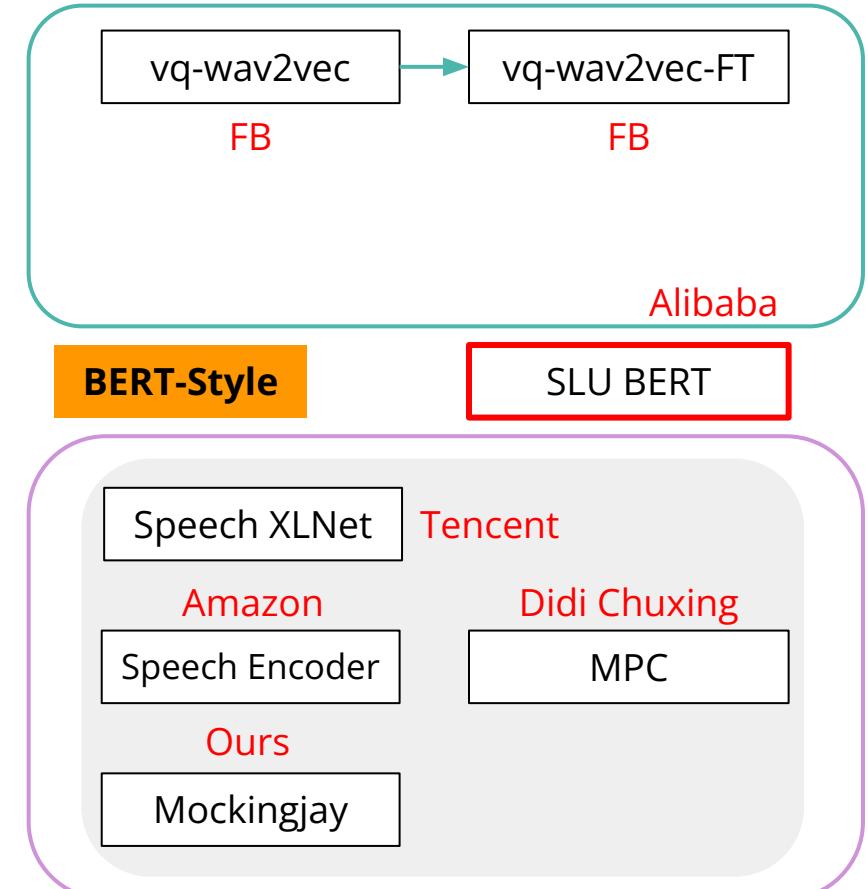
A Broad Introduction



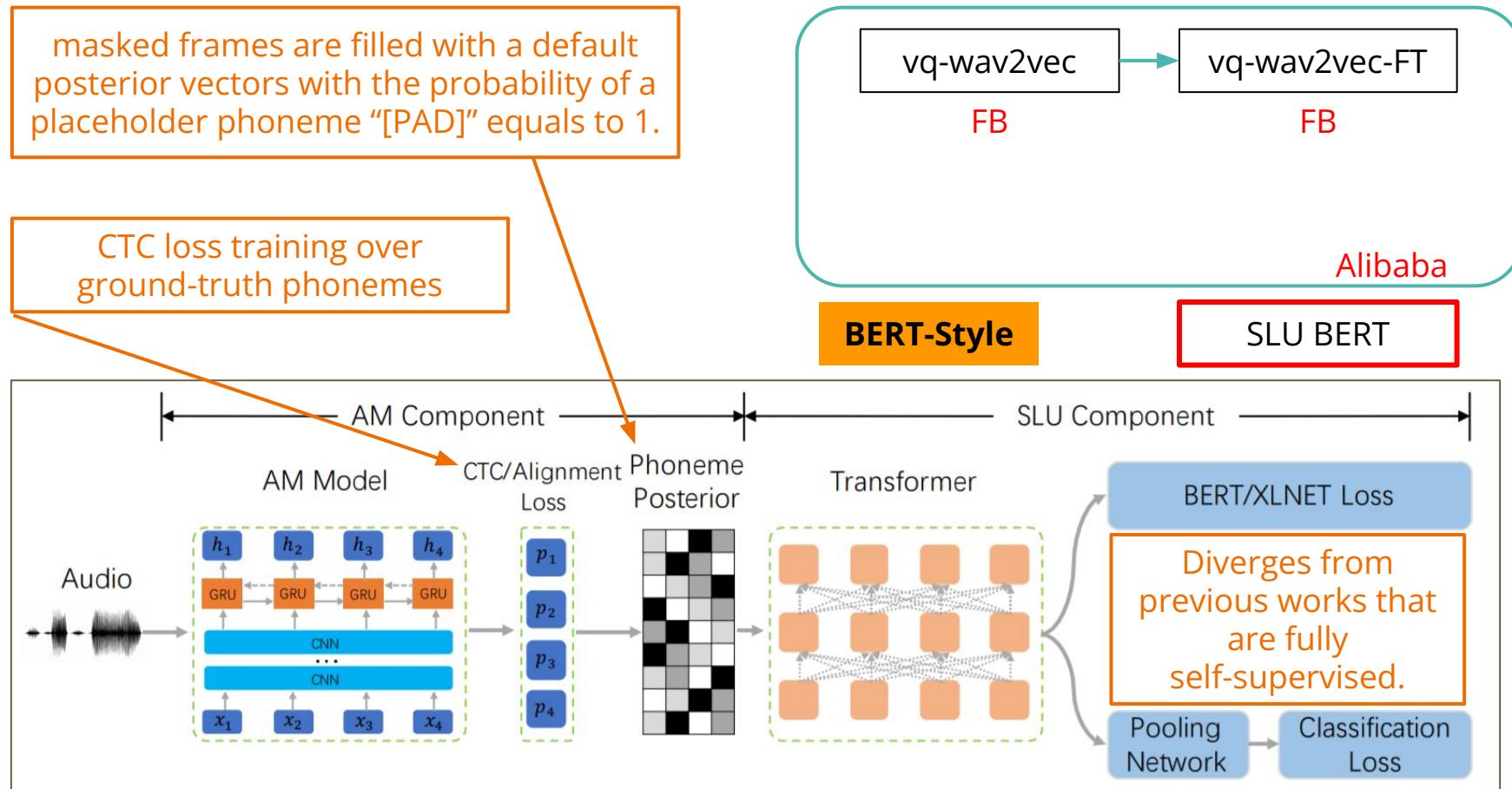
A Broad Introduction



A Broad Introduction



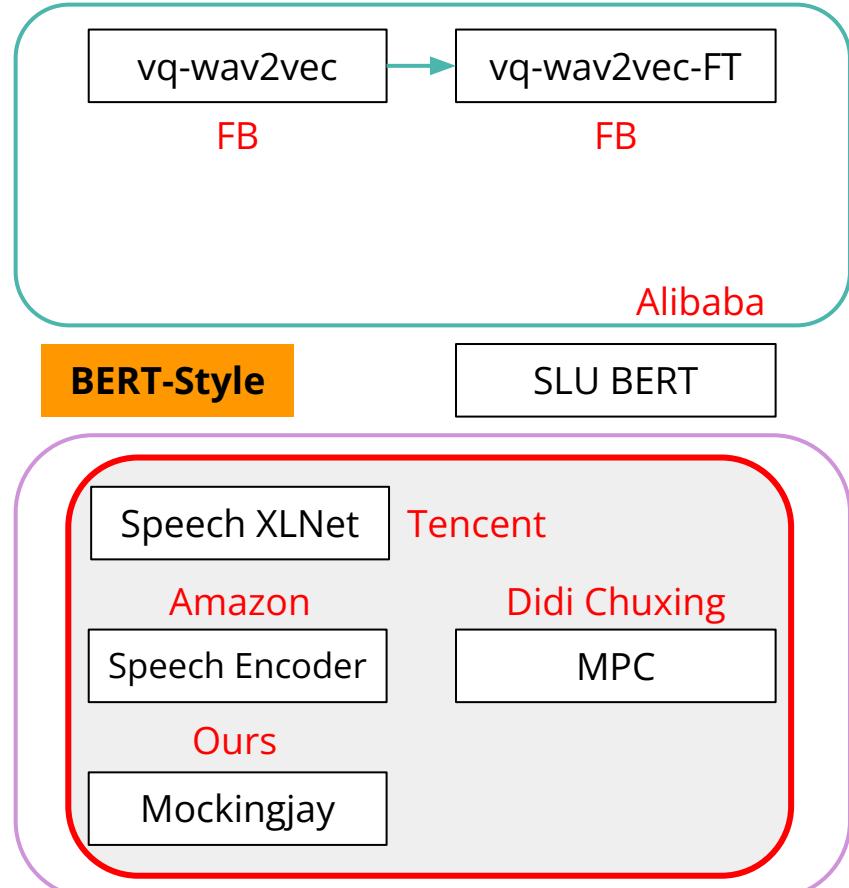
A Broad Introduction



A Broad Introduction

The Trend:
All of these works emerges around October, 2019.
All submitted to ICASSP 2020

(Speech XLNet and MPC did not make it)



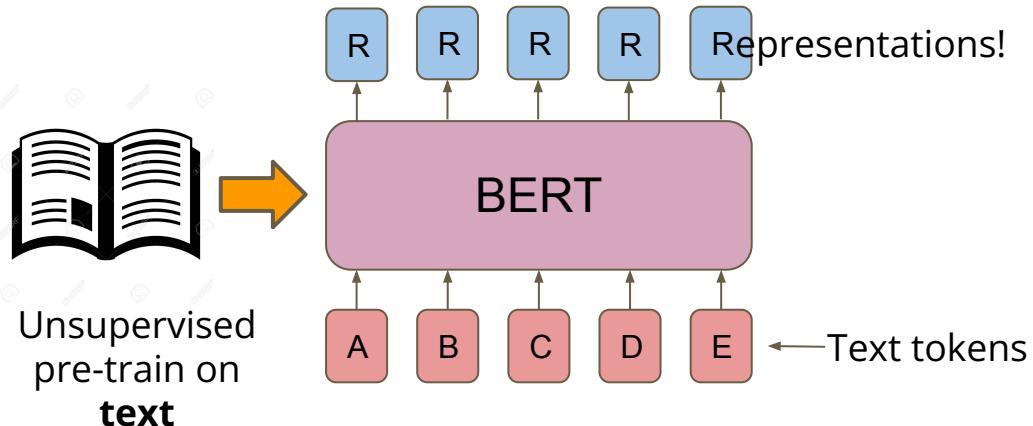
Mockingjay

From BERT to Speech BERT

From BERT to Speech BERT

NLP BERT:

Language Representation Learning



From BERT to Speech BERT

NLP BERT:

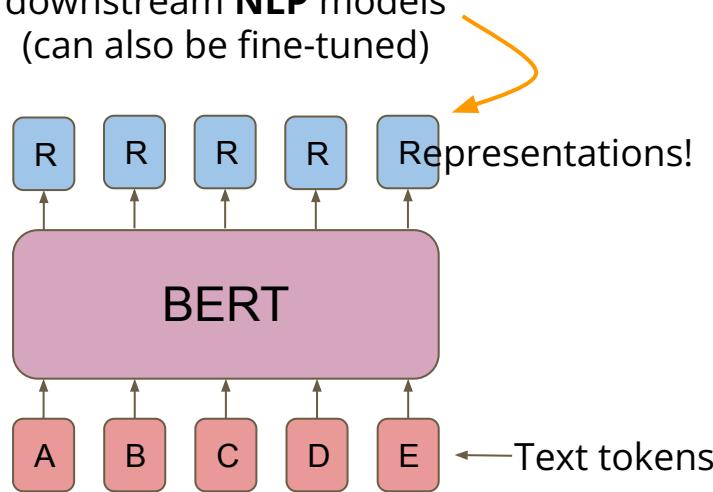
Language Representation Learning

Usage:

Extracts features for
downstream **NLP** models
(can also be fine-tuned)



Unsupervised
pre-train on
text



From BERT to Speech BERT

NLP BERT:

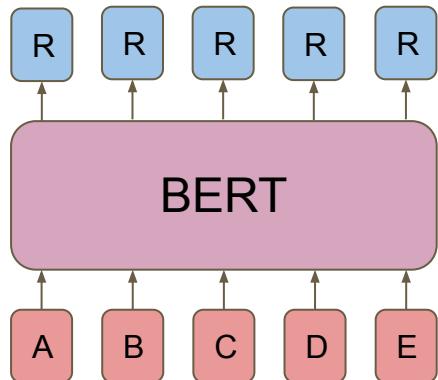
Language Representation Learning

Usage:

Extracts features for downstream **NLP** models
(can also be fine-tuned)



Unsupervised pre-train on
text

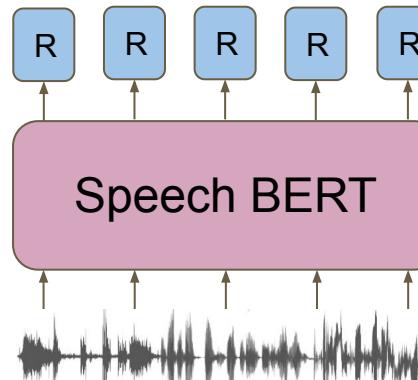


Speech BERT:

Speech Representation Learning

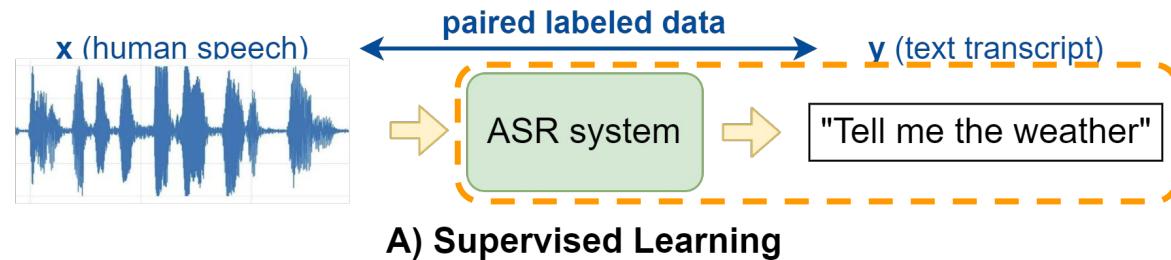
Usage:

Extracts features for downstream **SLP** models
(can also be fine-tuned)

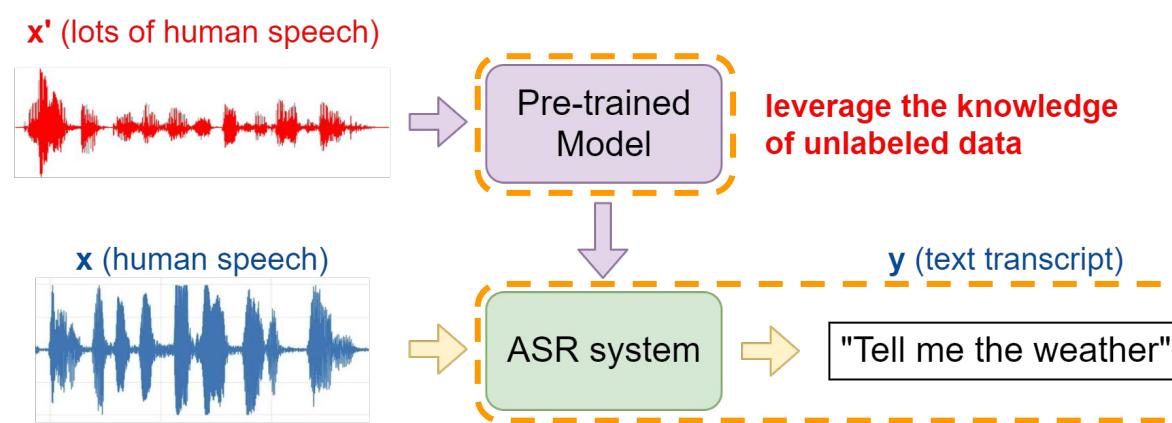


Unsupervised pre-train on
speech

Recall: Self-Supervised Learning for Speech



A) Supervised Learning

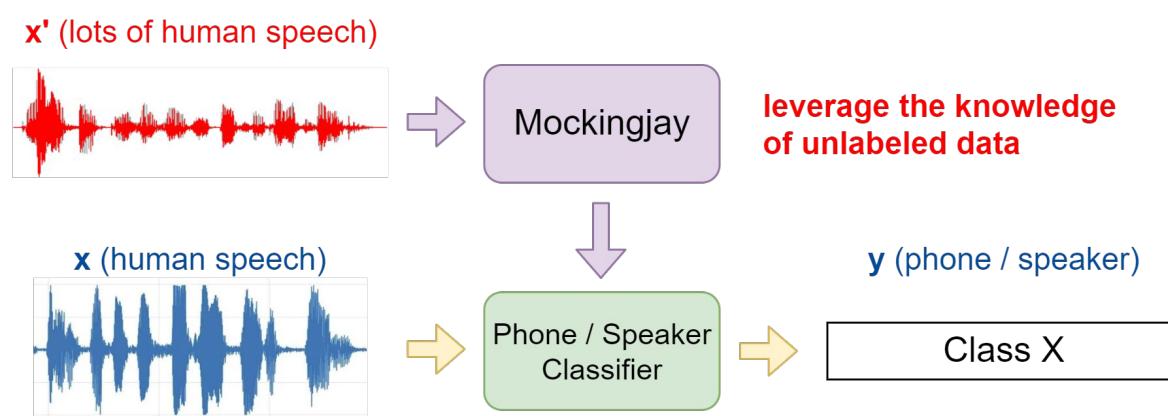


B) Self-Supervised Learning for Improving Supervised Systems

Self-Supervised Learning: Mockingjay

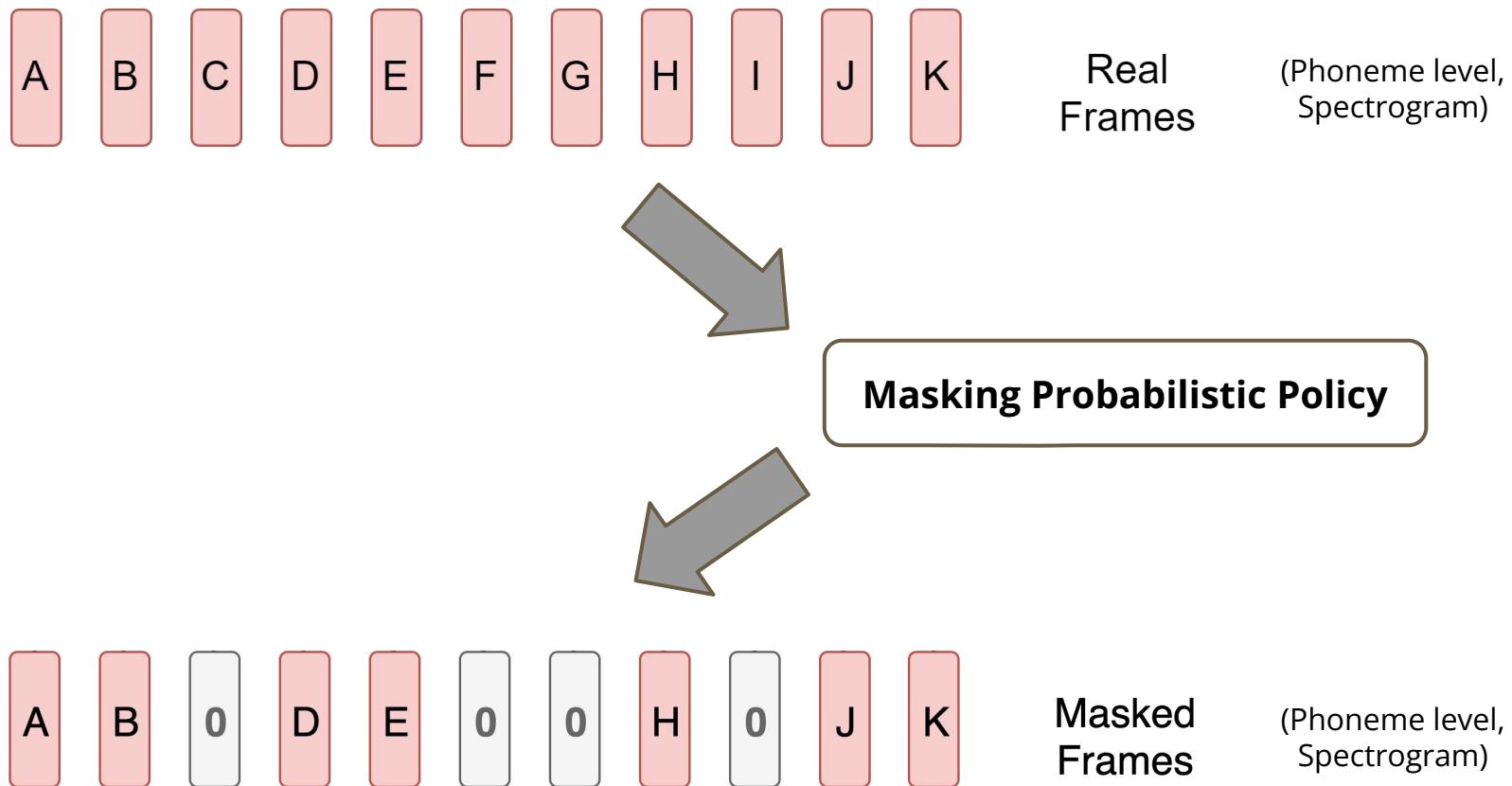


A) Supervised Learning

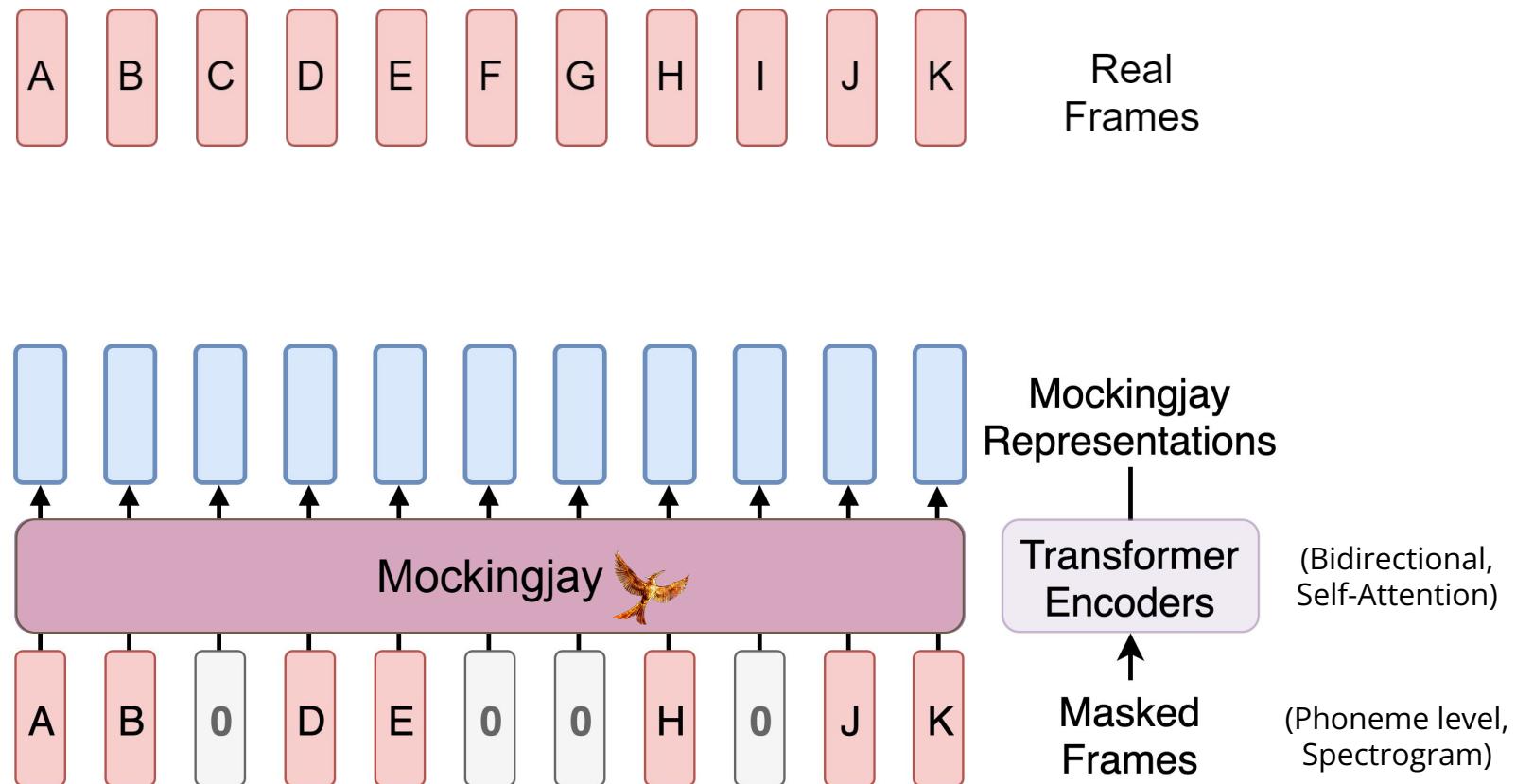


B) Self-Supervised Learning for Improving Supervised Systems

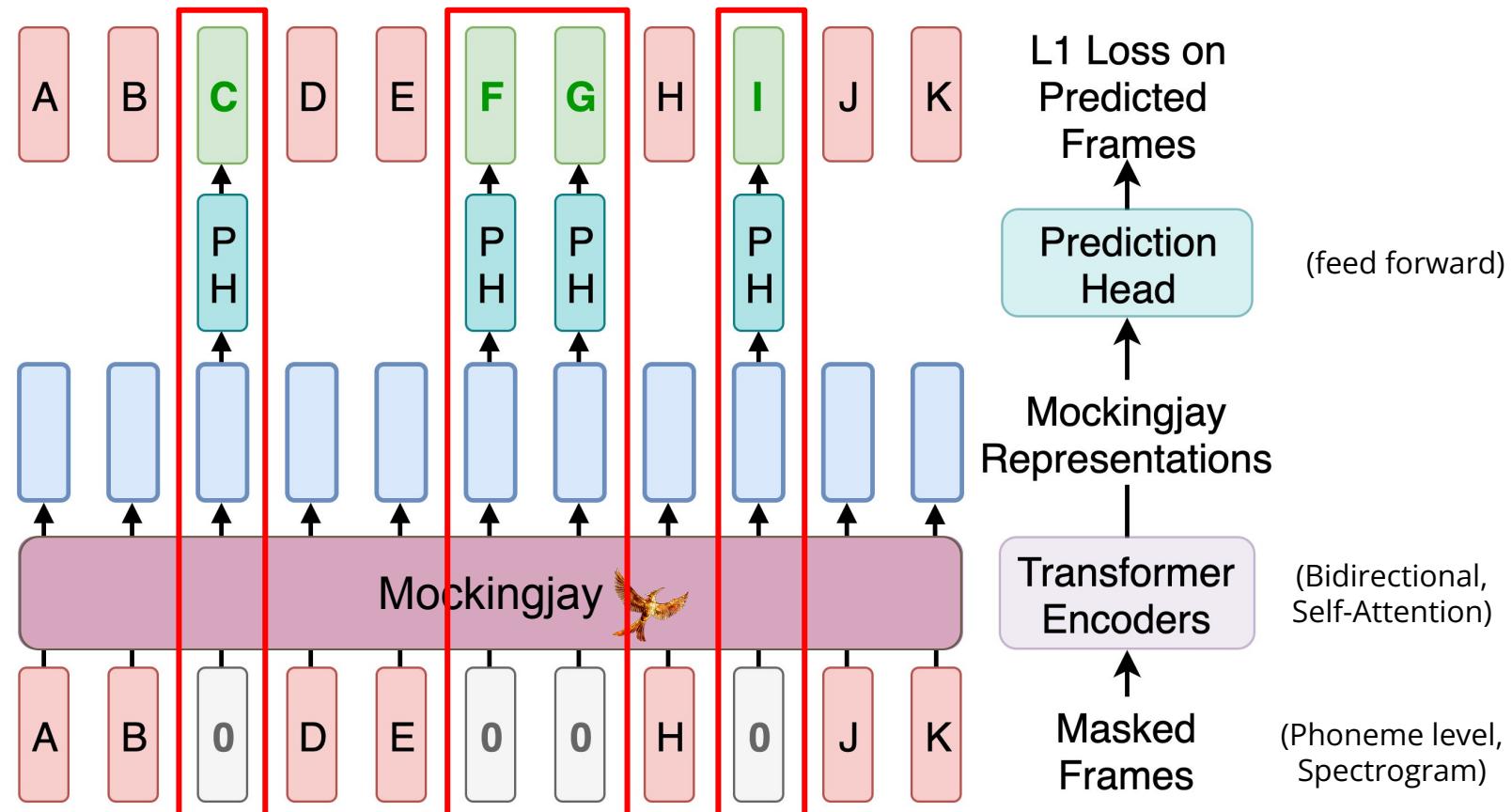
Pre-Training Task: Masked Acoustic Model



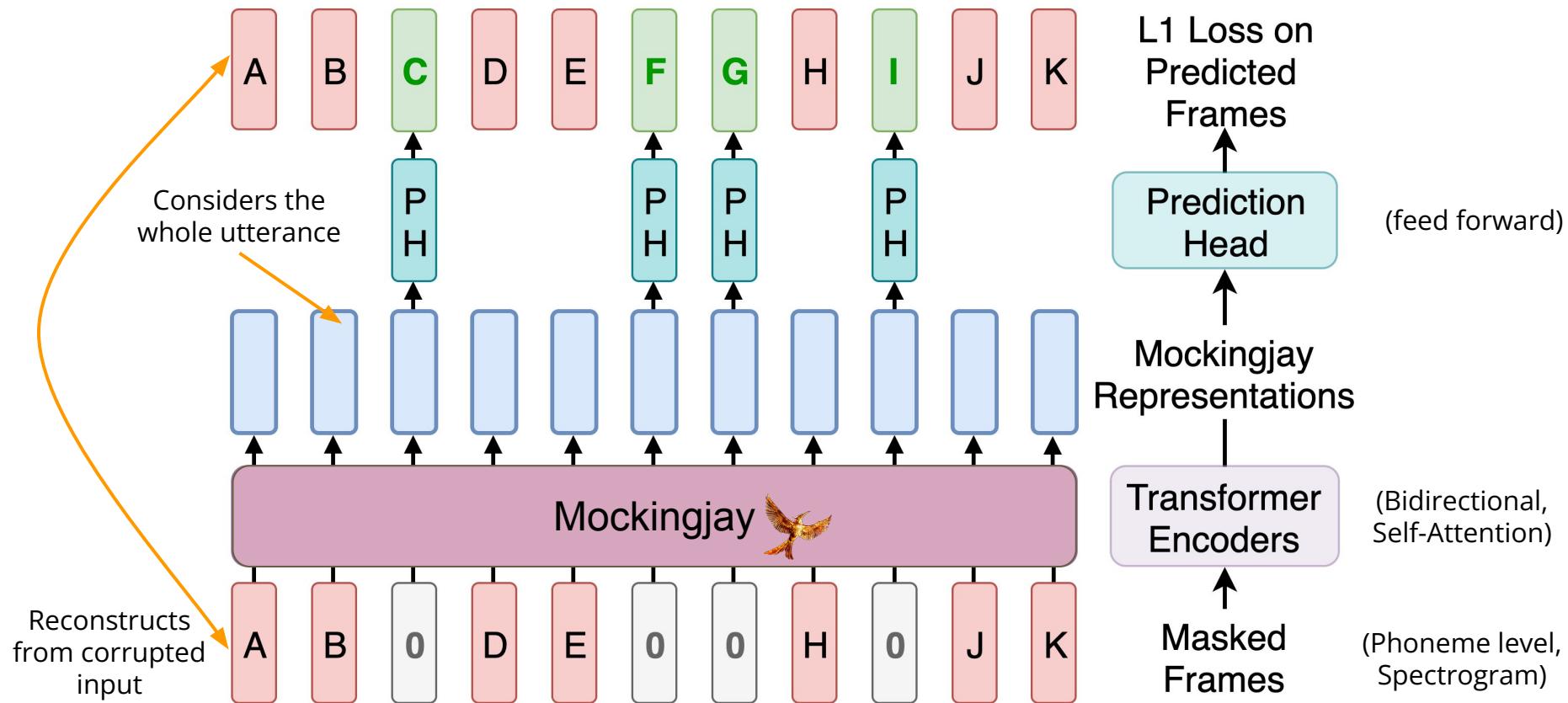
Pre-Training Task: Masked Acoustic Model



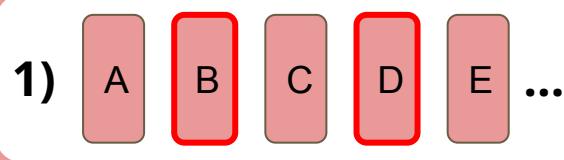
Pre-Training Task: Masked Acoustic Model



Pre-Training Task: Masked Acoustic Model

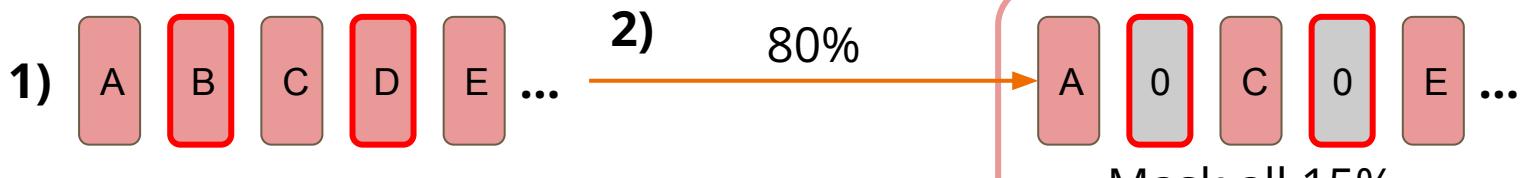


Probabilistic Policy for Masking Frames



1) Select **15%** of the frames for prediction (highlighted in green).

Probabilistic Policy for Masking Frames

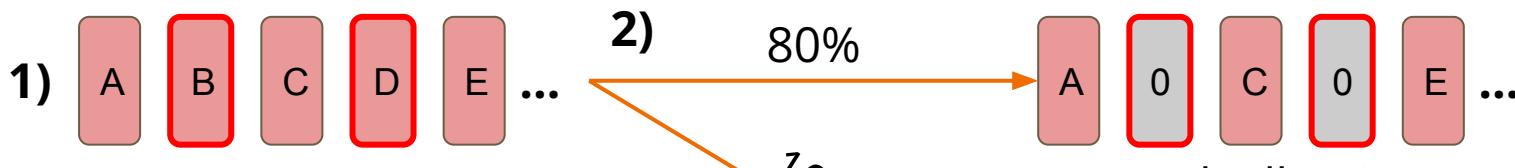


1) Select **15%** of the frames for prediction (highlighted in green).

2) For all selected frames:

- mask to zero **80%** of the time
- replace randomly **10%** of the time
- leave untouched **10%** of the time

Probabilistic Policy for Masking Frames



1) Select **15%** of the frames for prediction (highlighted in green).

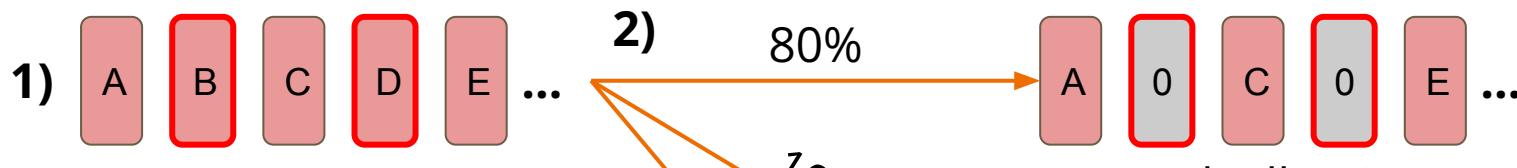
2) For all selected frames:

- mask to zero **80%** of the time
- replace randomly **10%** of the time
- leave untouched **10%** of the time

A G C Y E ...

Replace all 15%

Probabilistic Policy for Masking Frames



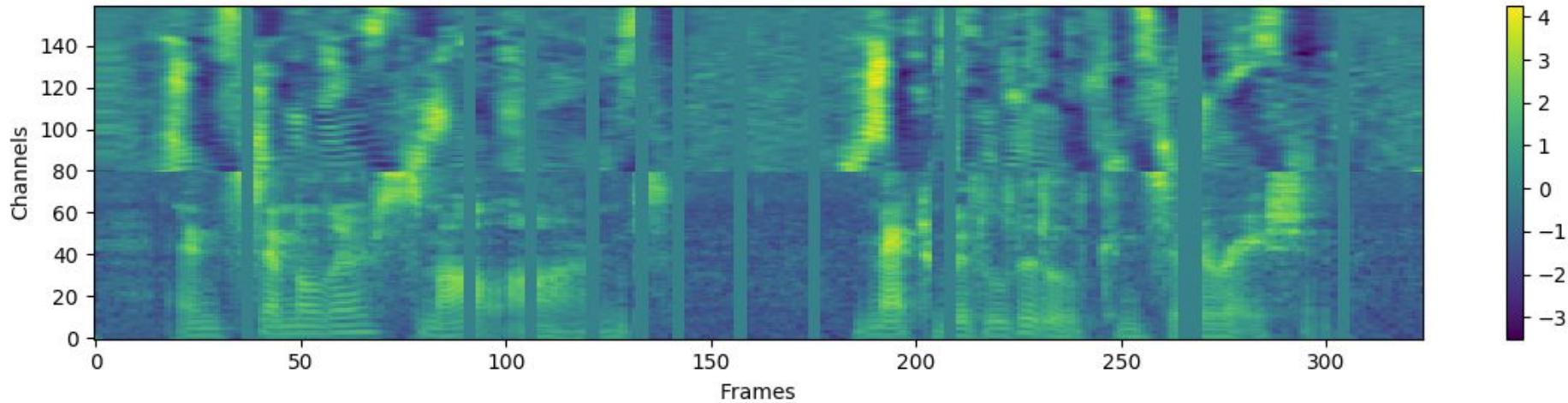
1) Select **15%** of the frames for prediction (highlighted in green).

2) For all selected frames:

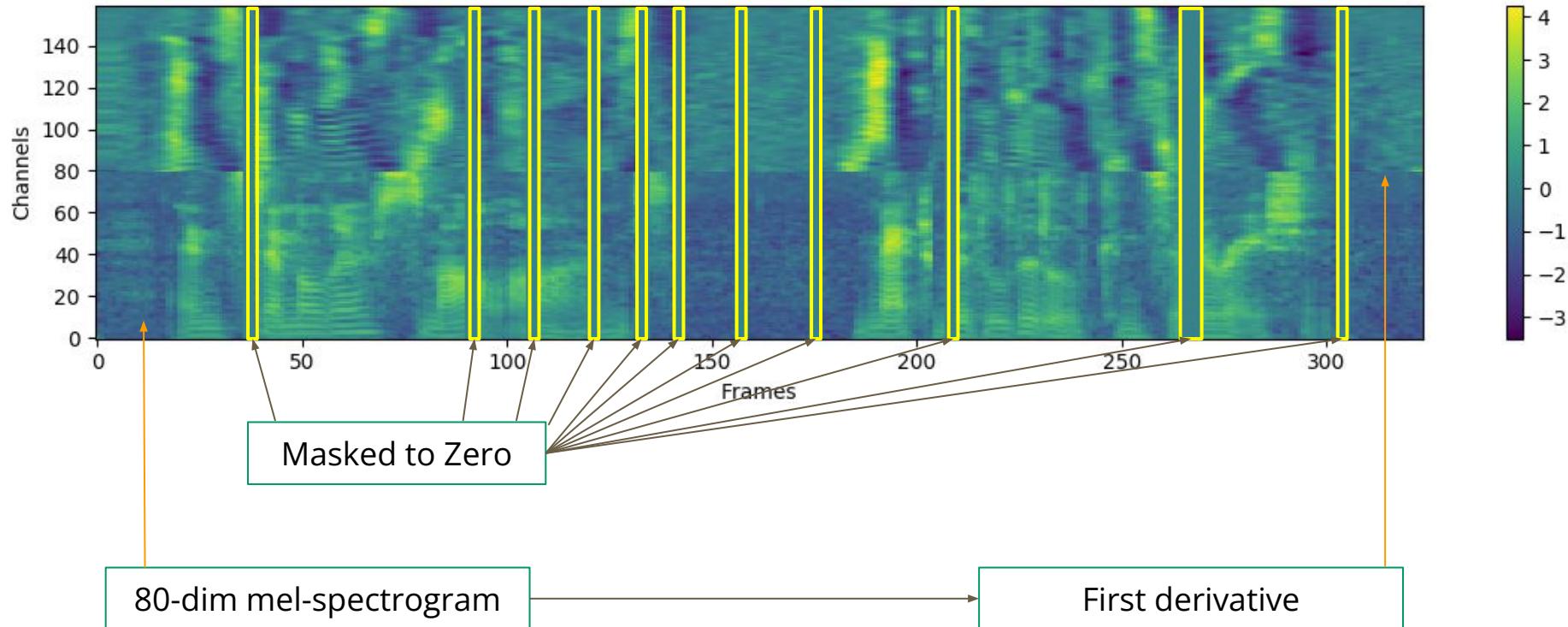
- mask to zero **80%** of the time
- replace randomly **10%** of the time
- leave untouched **10%** of the time

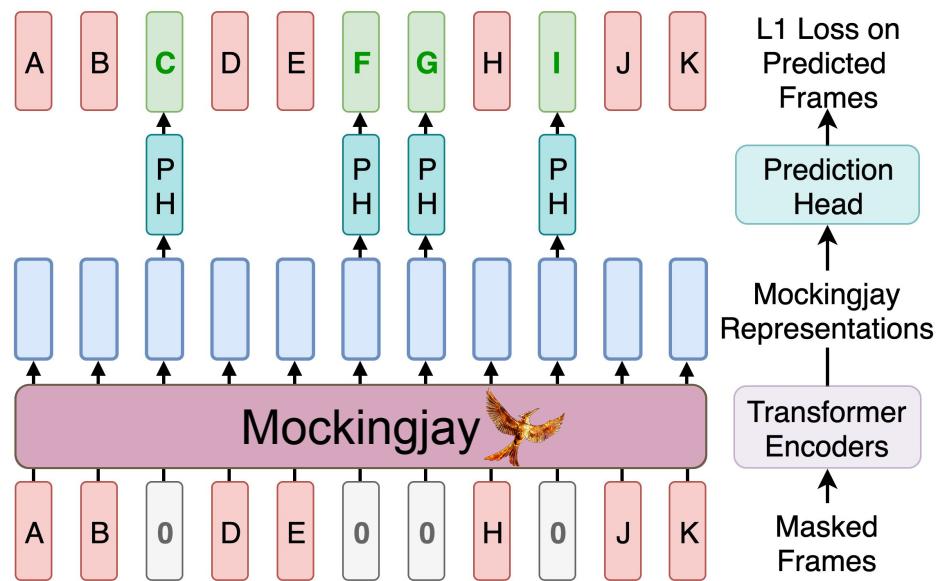
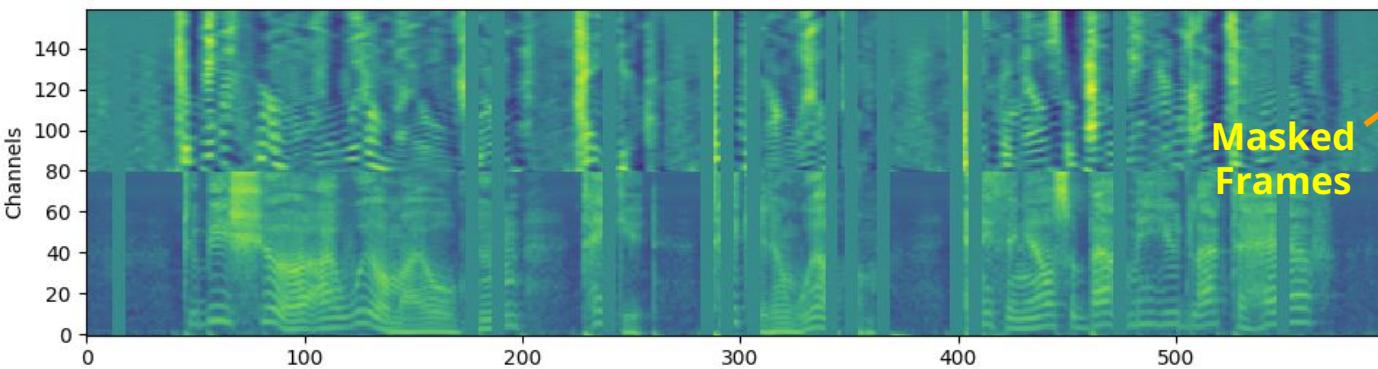
Do nothing, frames remain the same

Input Feature: Masked Spectrogram

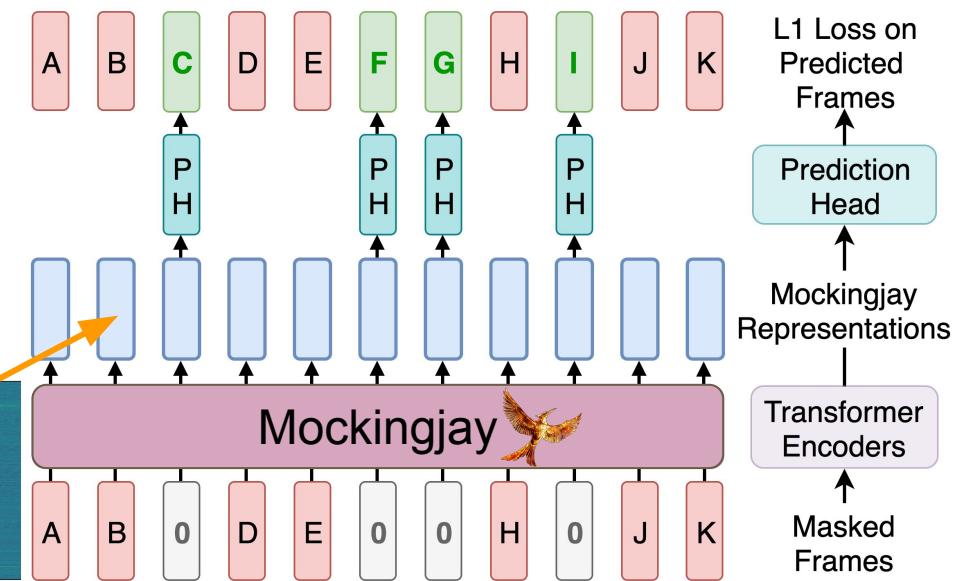
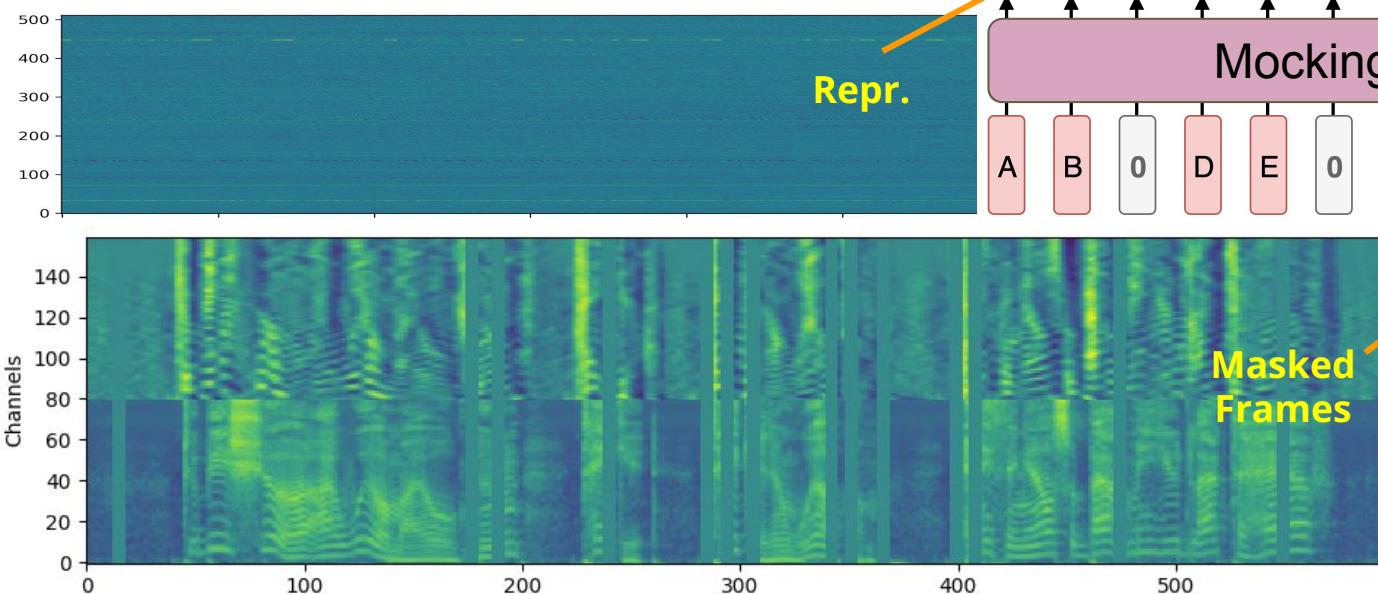


Input Feature: Masked Spectrogram

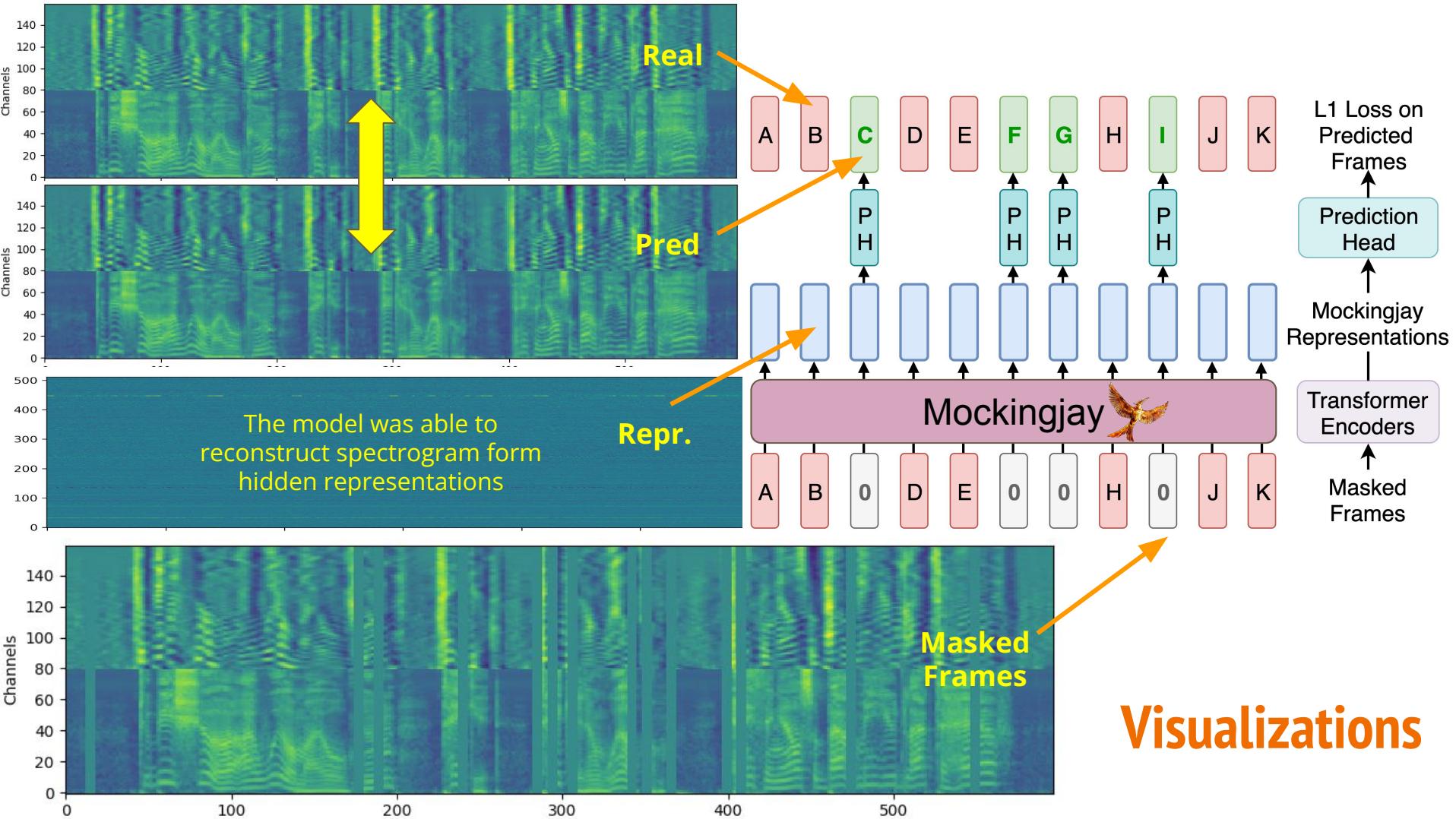




Visualizations



Visualizations



Migrating from text to speech

Acoustic Features: long and locally smooth in nature,

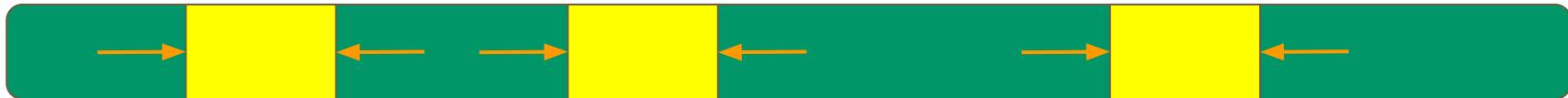
need to 1) shorten the sequence and 2) mask over a longer span



Migrating from text to speech

Acoustic Features: long and locally smooth in nature,

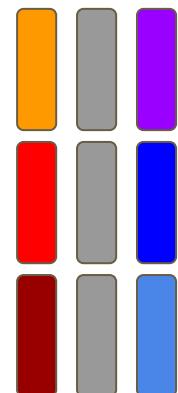
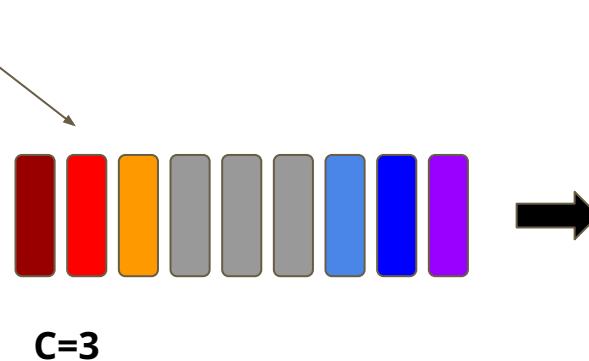
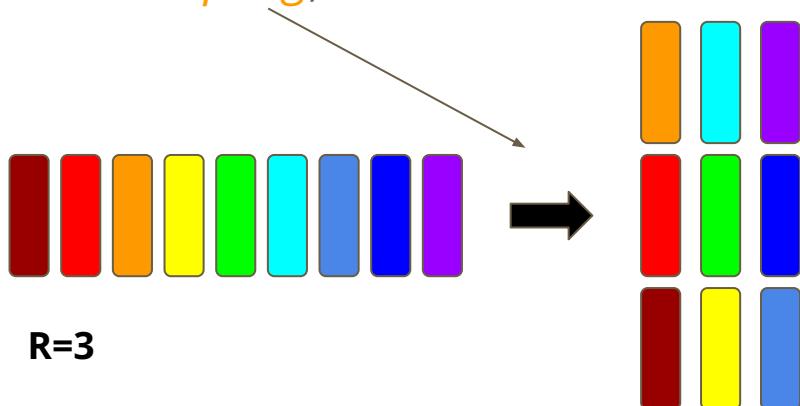
need to 1) shorten the sequence and 2) mask over a longer span



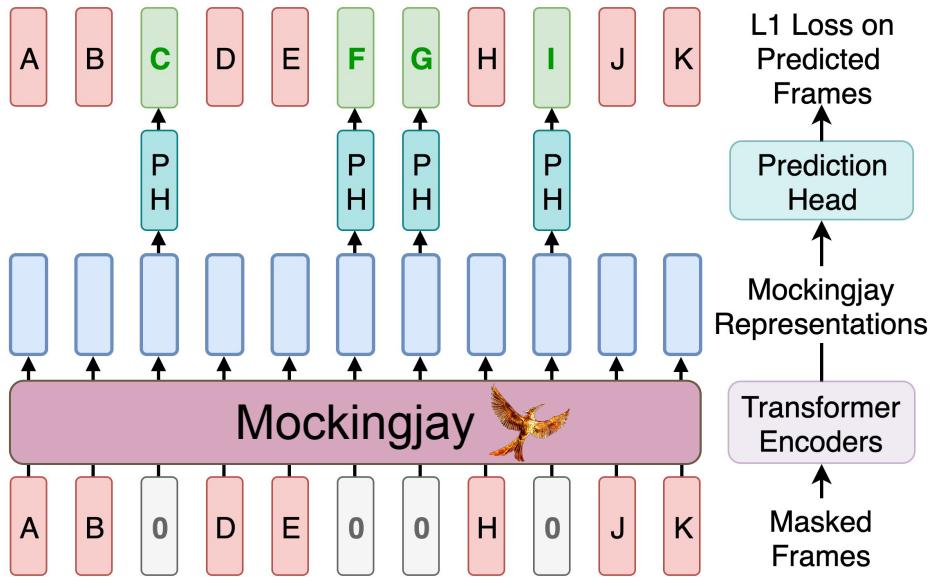
Address the long and smooth problem with:

Downsampling, and *consecutive masking*

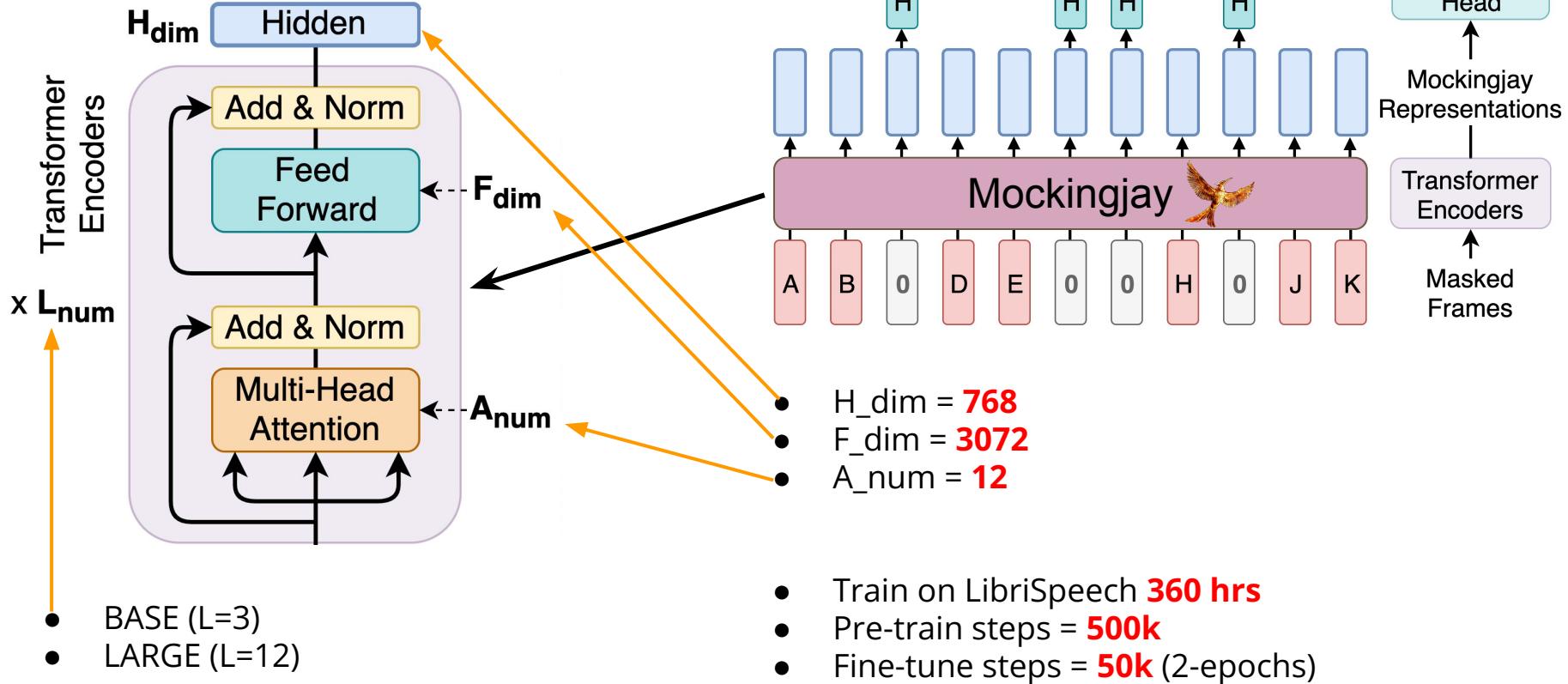
$R=3, C=3$



Model Architecture

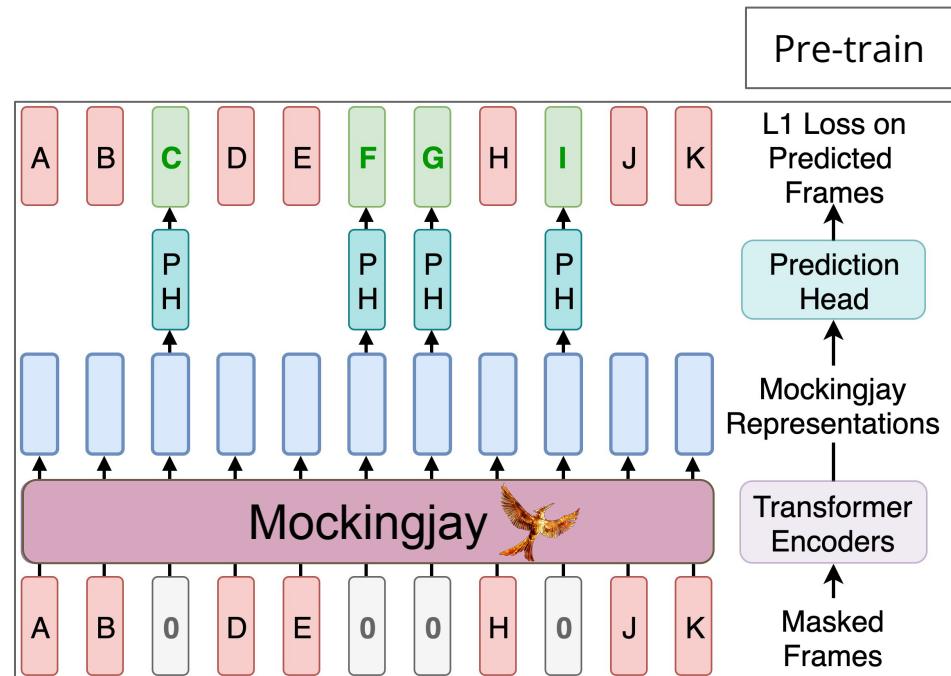


Model Architecture



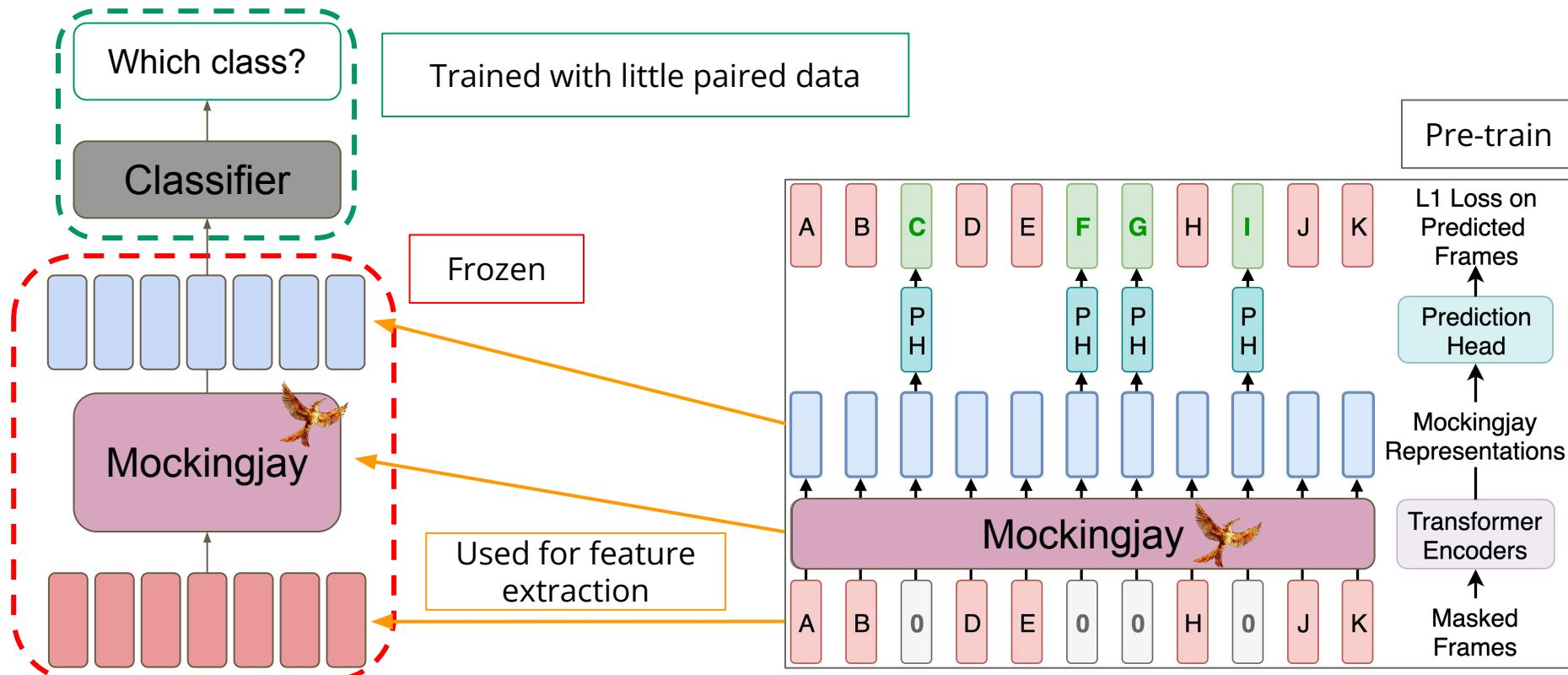
Incorporating with Downstream Tasks

1) Feature Extraction



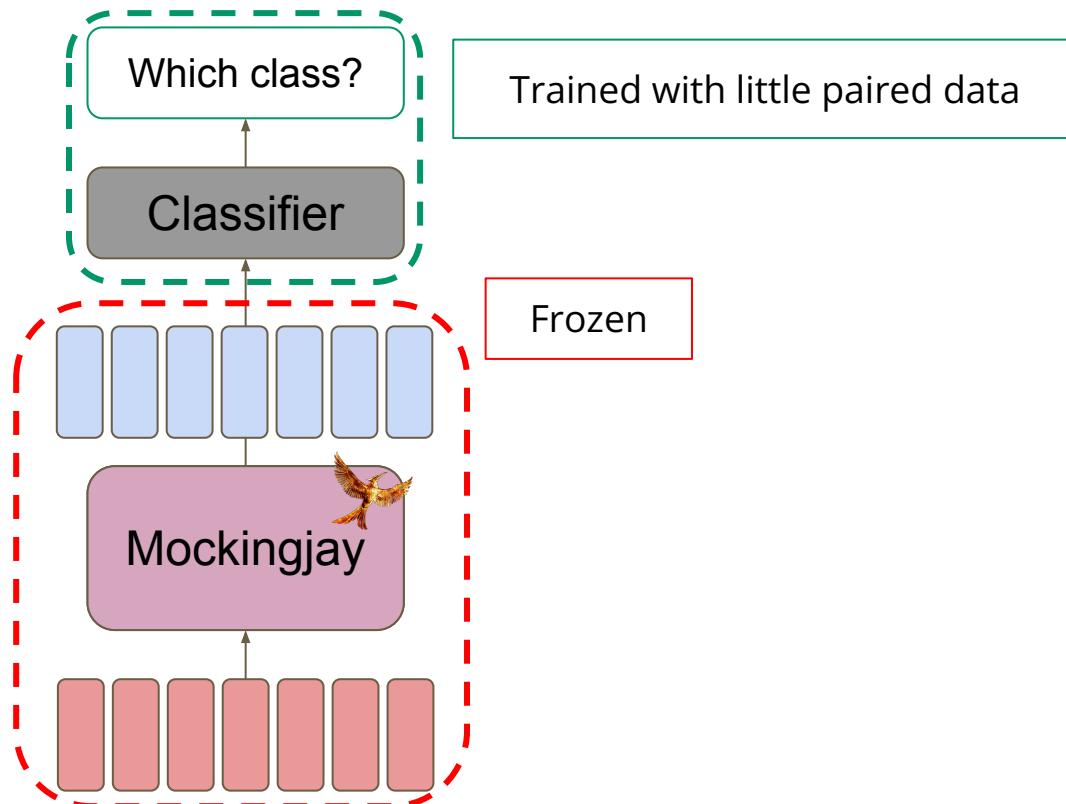
Incorporating with Downstream Tasks

1) Feature Extraction



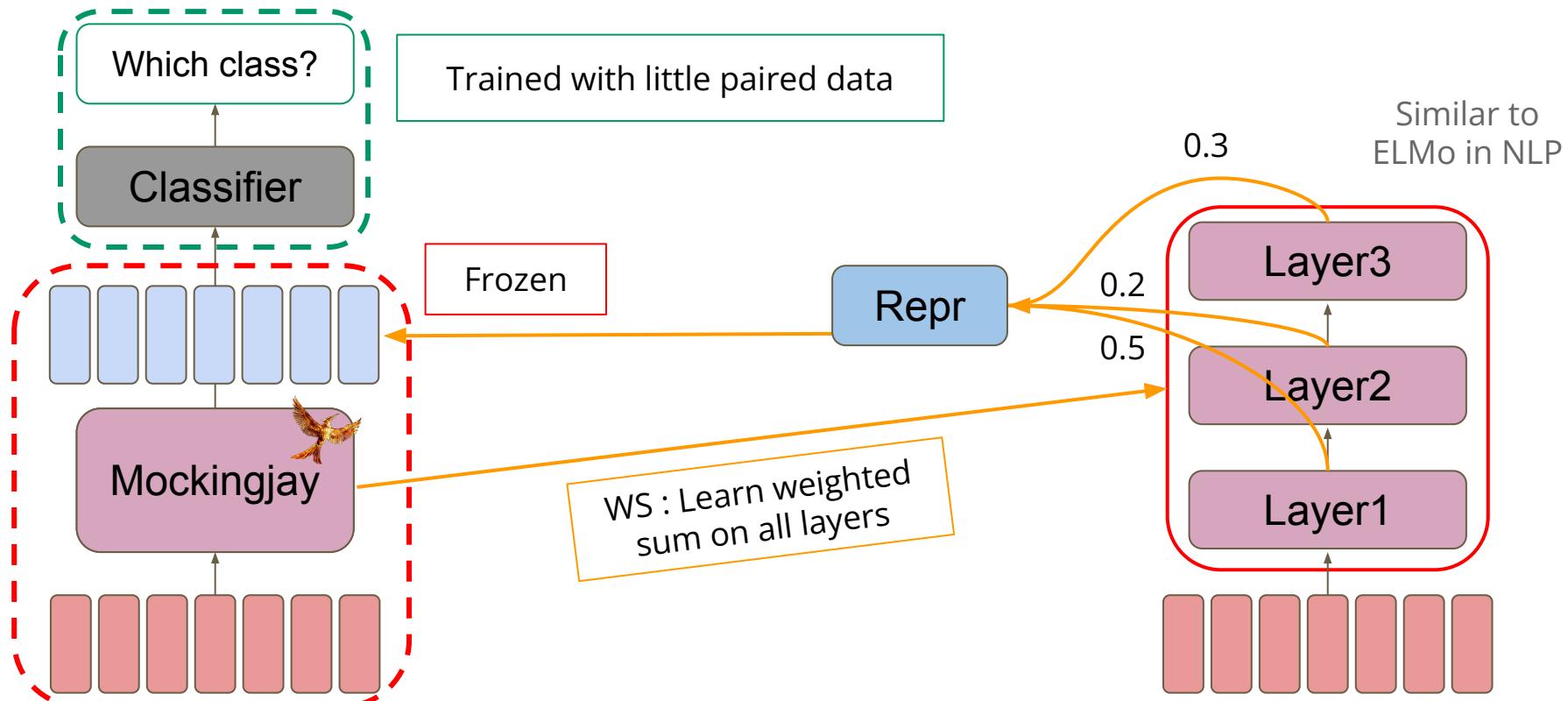
Incorporating with Downstream Tasks

2) Weighted Sum from All Layers (WS)



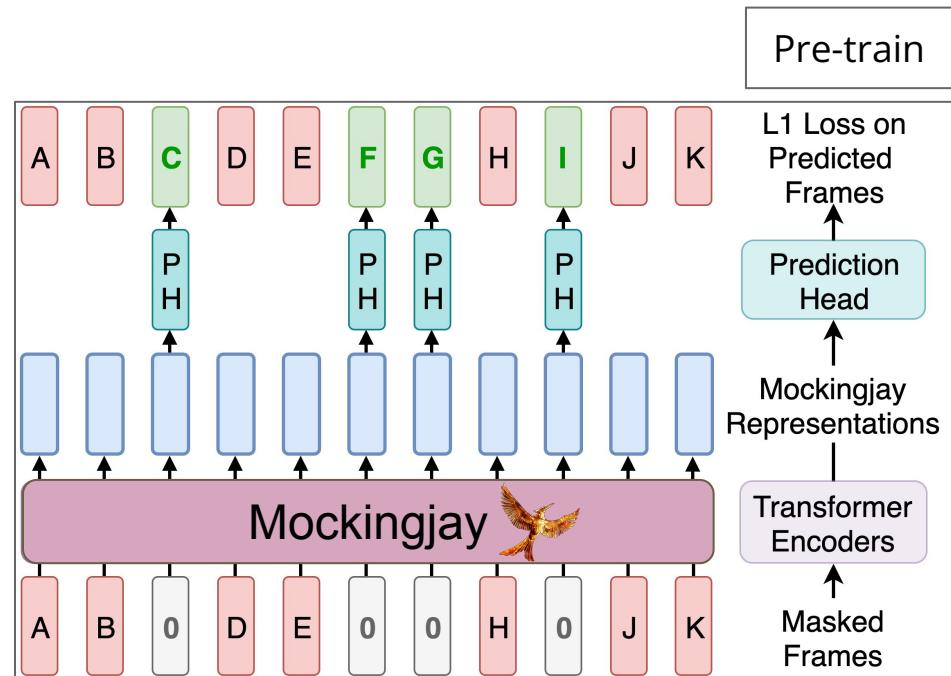
Incorporating with Downstream Tasks

2) Weighted Sum from All Layers (WS)



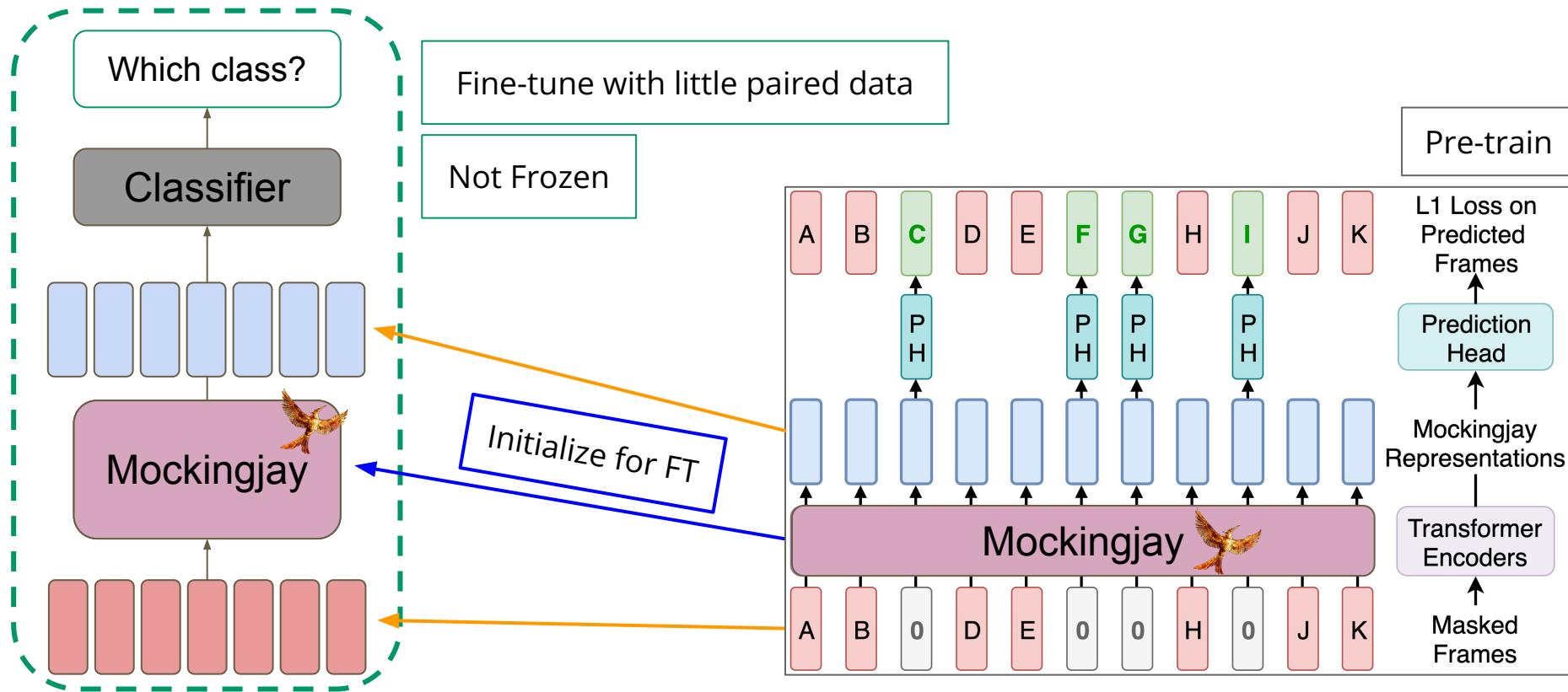
Incorporating with Downstream Tasks

3) Fine-tune (FT2)



Incorporating with Downstream Tasks

3) Fine-tune (FT2)



Experiments - 1/3

Acoustic Features	Phoneme Classification	Speaker Recognition	Sentiment Classification
Mel Features	49.1	70.1	64.6
BASE	60.9	94.5	67.4
LARGE	64.3	96.3	70.1

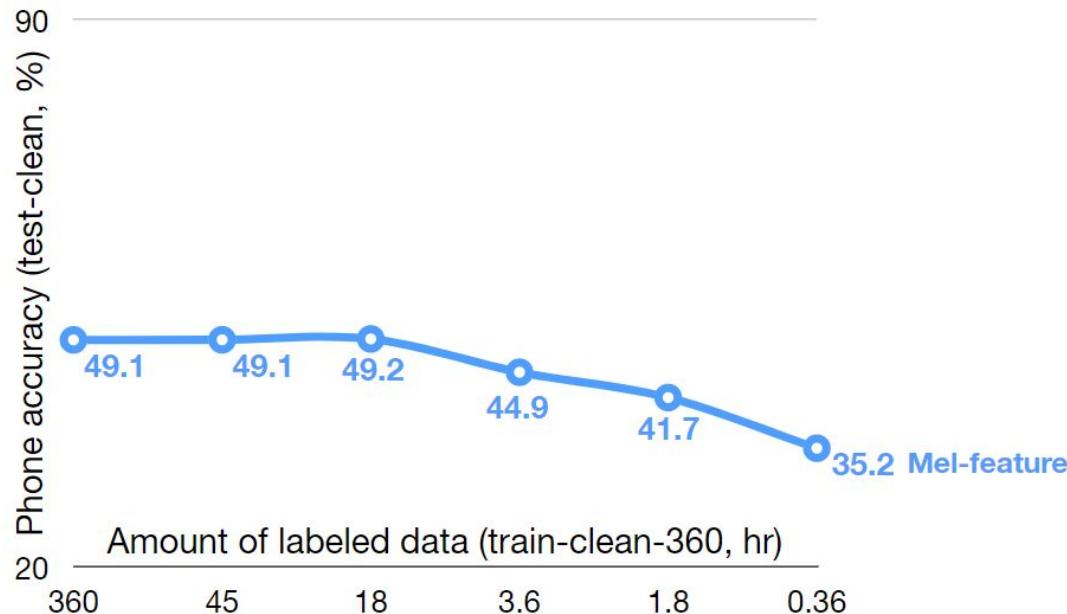
Consistent results over all three tasks:
Mel < BASE < LARGE

Experiments - 2/3

Acoustic Features	Phoneme Classification	Speaker Recognition	Sentiment Classification
Mel Features	49.1	70.1	64.6
BASE	60.9	94.5	67.4
LARGE	64.3	96.3	70.1
LARGE-WS	69.9	96.4	71.1

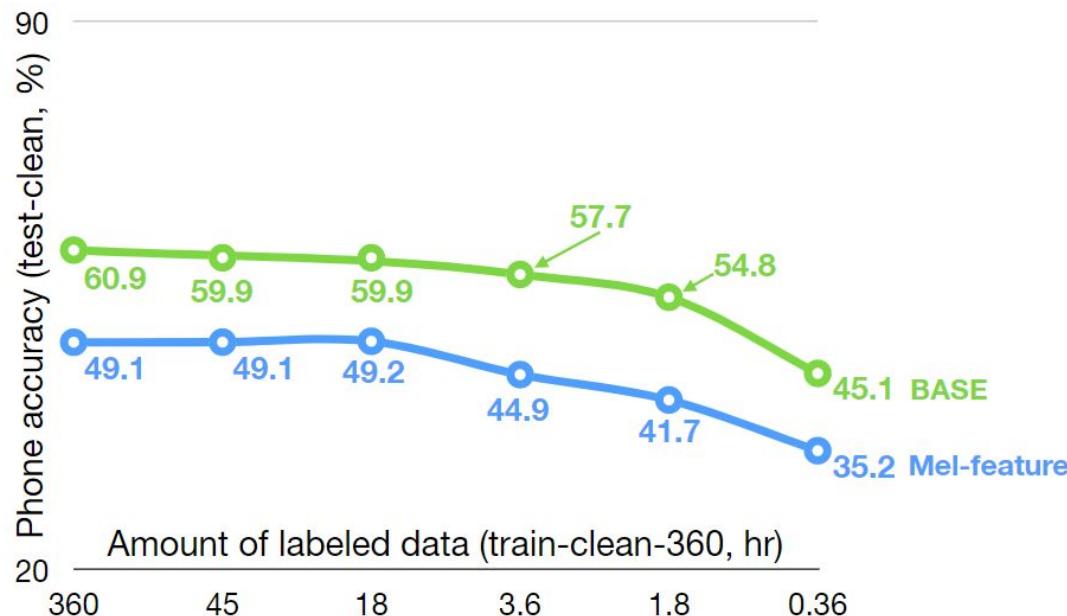
Consistent results over all three tasks:
LARGE < LARGE-WS

Low-Resource Experiments - 1/6



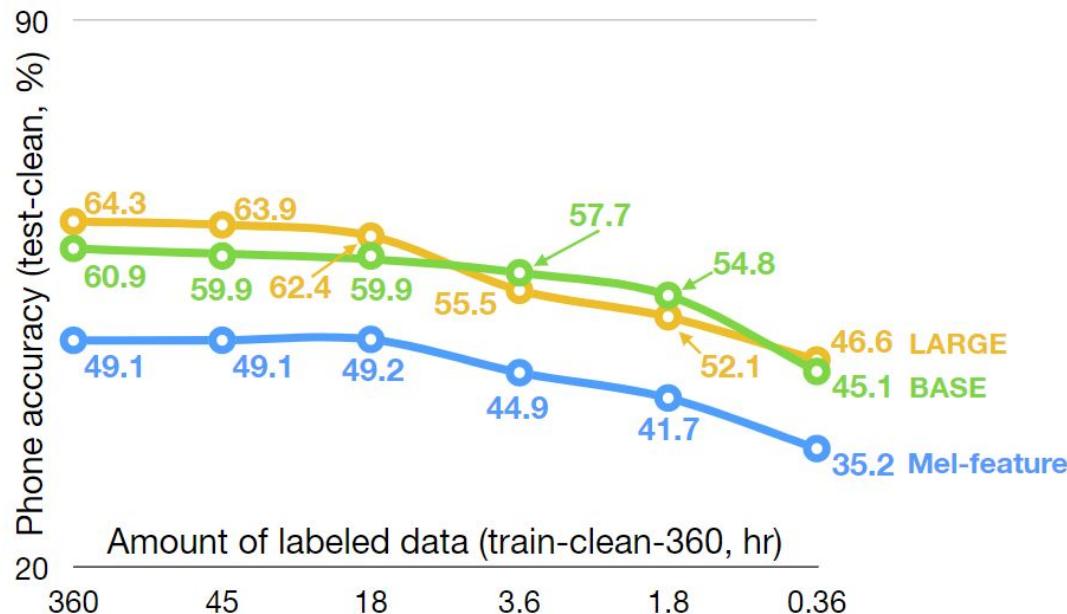
We demonstrate how pre-training on speech can improve supervised training in low resource scenarios, we train with reduced amount of labels.

Low-Resource Experiments - 2/6



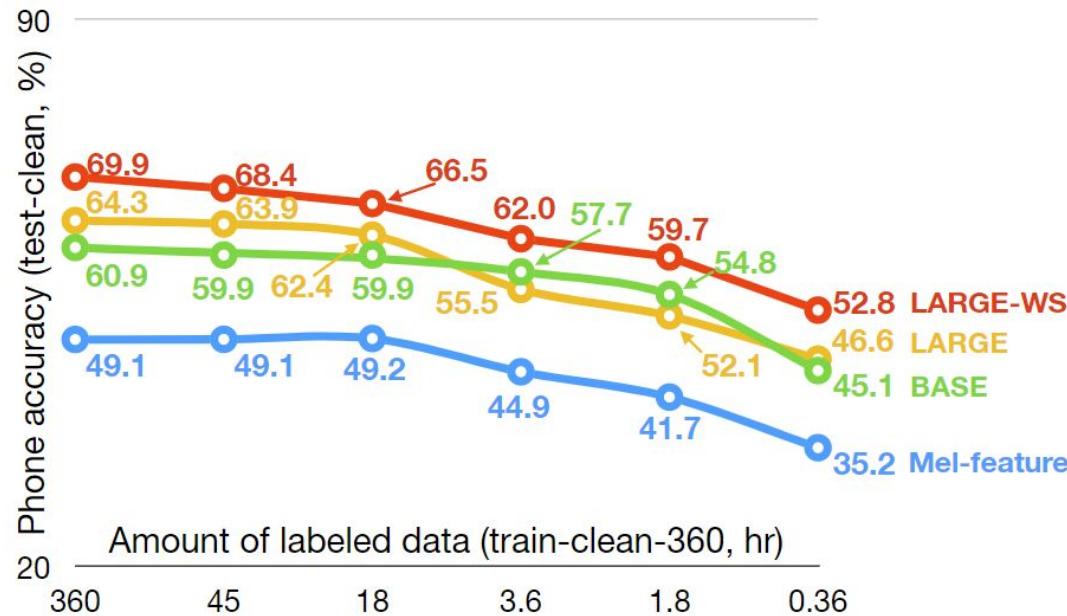
Boxed Text: Mel < BASE

Low-Resource Experiments - 3/6



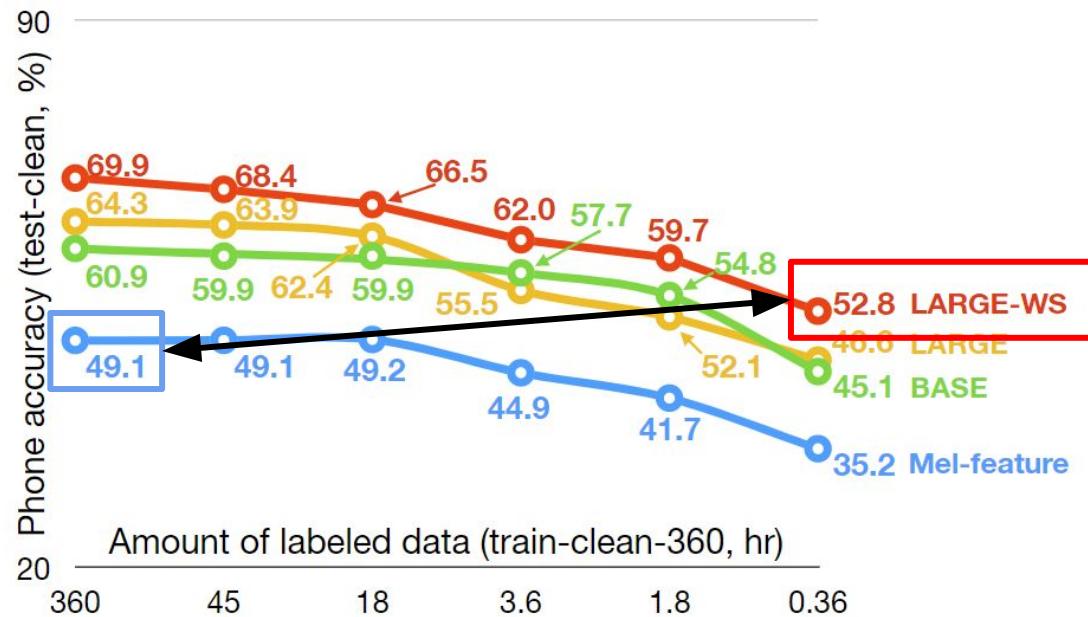
Mel < BASE < LARGE

Low-Resource Experiments - 4/6



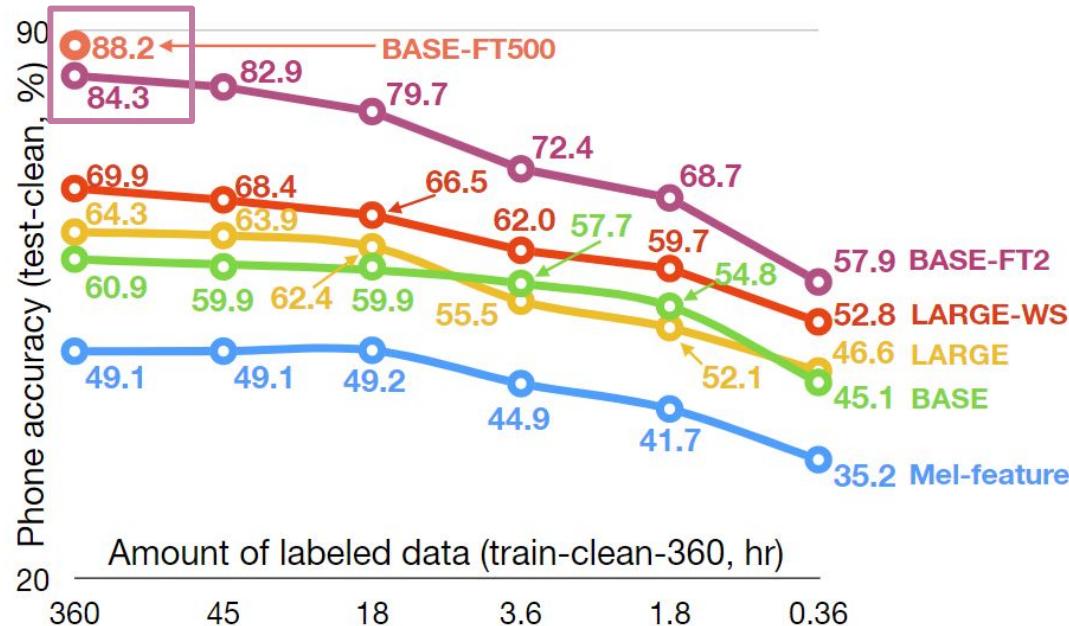
LARGE < LARGE-WS
with an avg 5.75% improvement

Low-Resource Experiments - 4/6



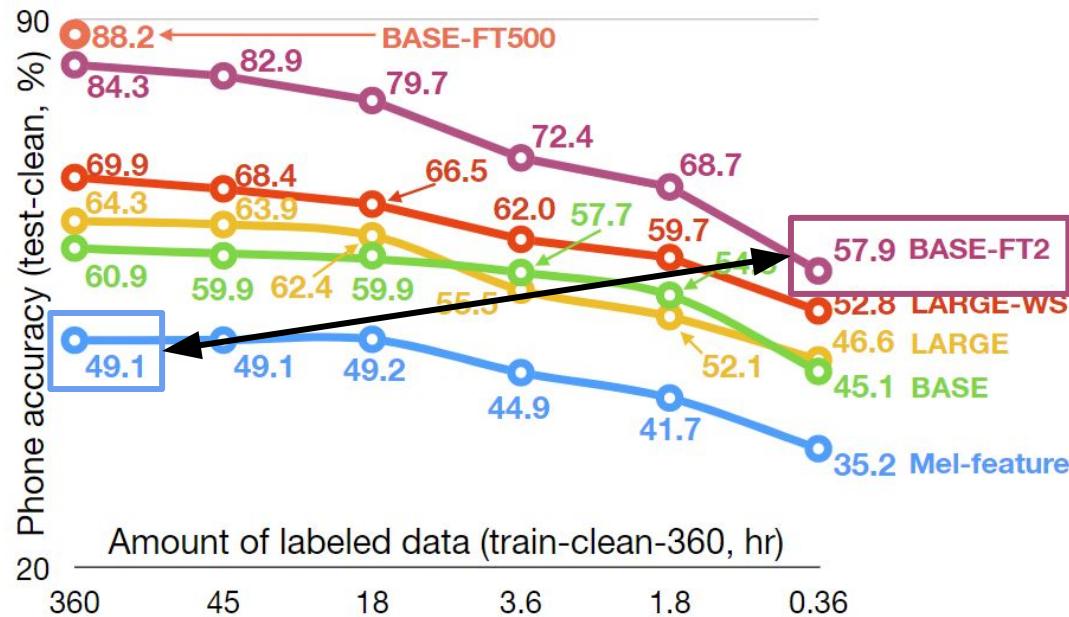
With 0.1% of labels,
LARGE-WS (52.8%) outperformed Mel (49.1%) that uses all 100% hours of labeled data.

Low-Resource Experiments - 5/6



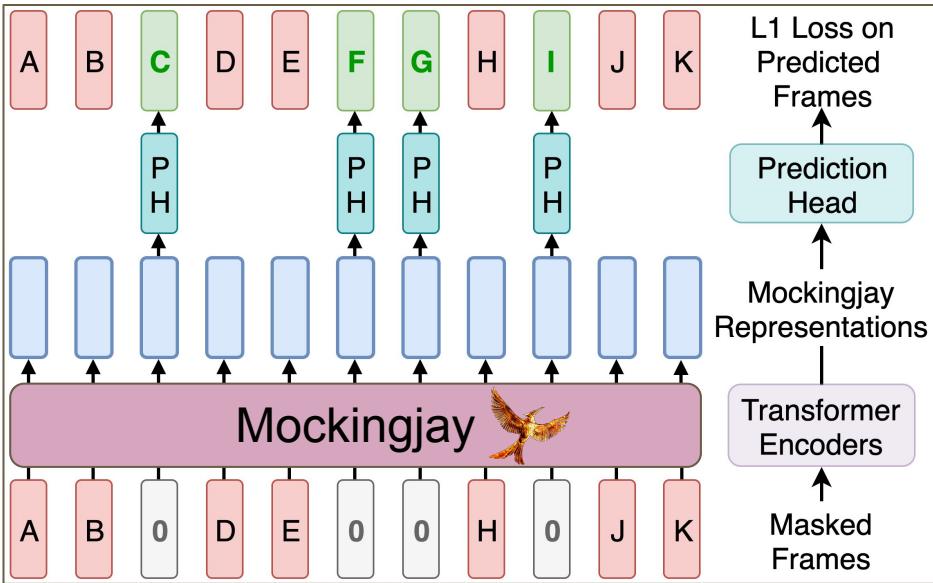
All < BASE-FT2

Low-Resource Experiments - 6/6

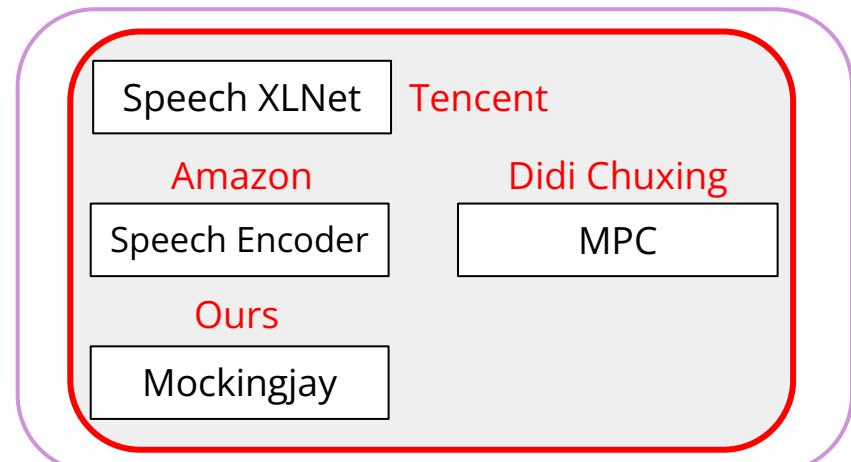
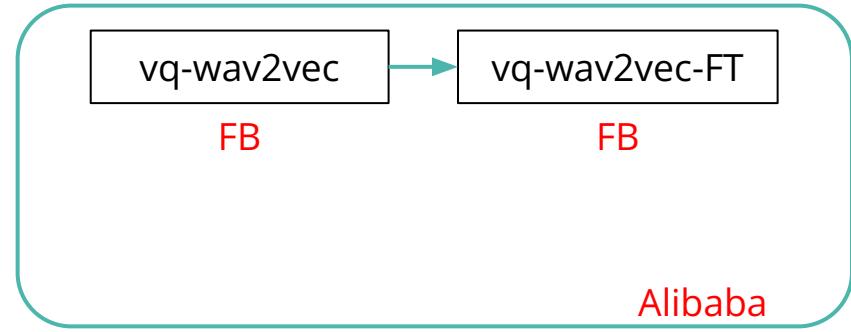


With 0.1% of labels,
BASE-FT2 (57.9%) outperformed Mel (49.1%) that uses all 100% hours of labeled data.

A Broad Introduction

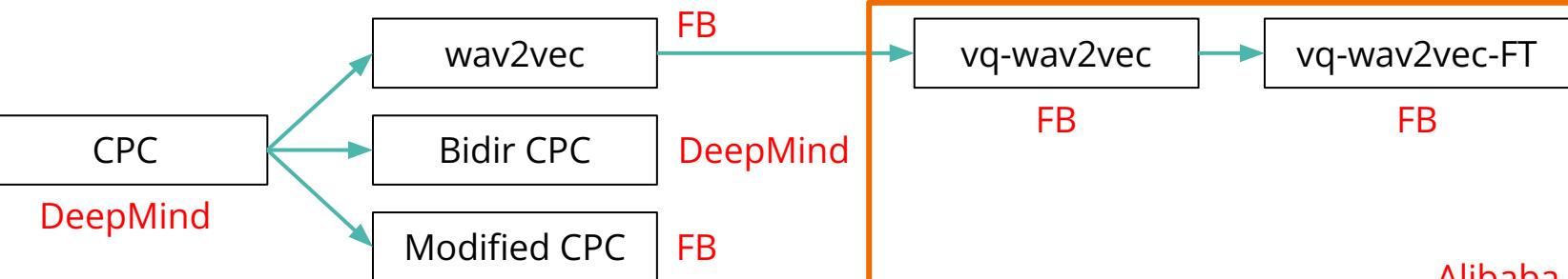


Intuition:
A model that can predict the partial loss
of small segments of speech,
should provide a contextualized
understanding of previous and later content.



A Broad Introduction

Contrastive Predictive Losses



Reconstruction Losses

MIT

APC

Multi-Target APC

MIT

DeCoAR

Amazon

Google

Autoencoder

Google

Audio2Vec

Google

Phase

BERT-Style

SLU BERT

Speech XLNet

Tencent

Amazon

Speech Encoder

Didi Chuxing

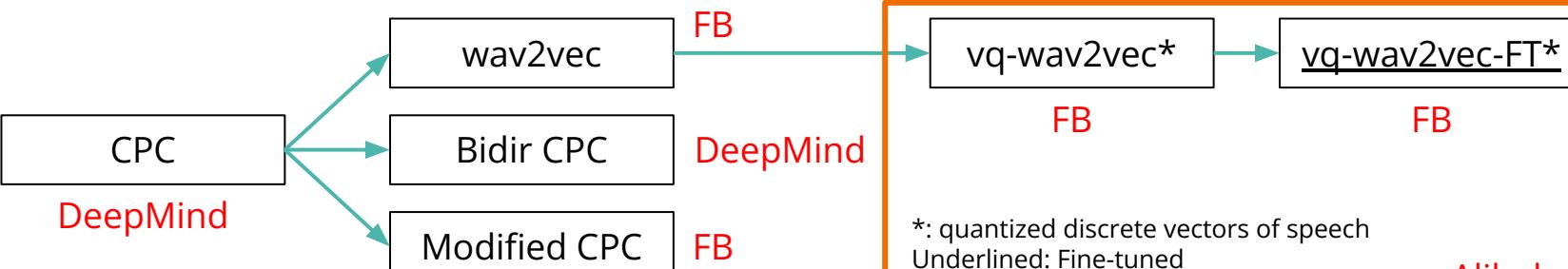
MPC

Ours

Mockingjay

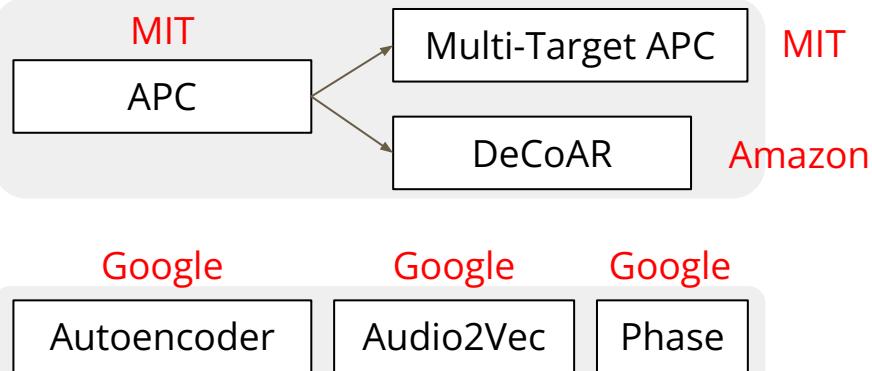
A Broad Introduction

Contrastive Predictive Losses

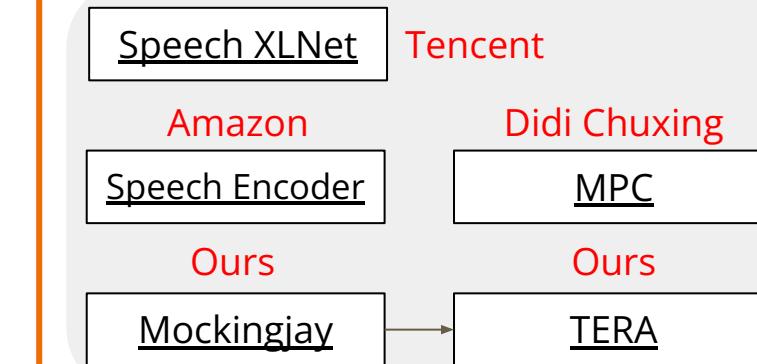


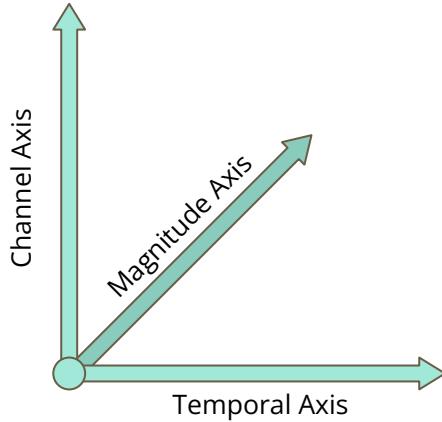
Alibaba

Reconstruction Losses



BERT-Style



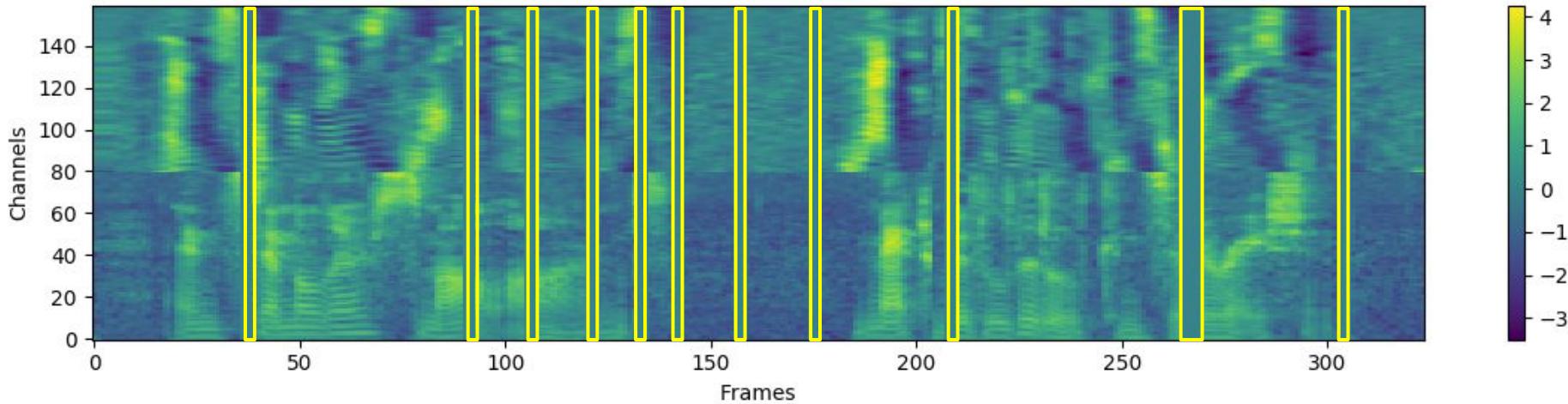


TERA

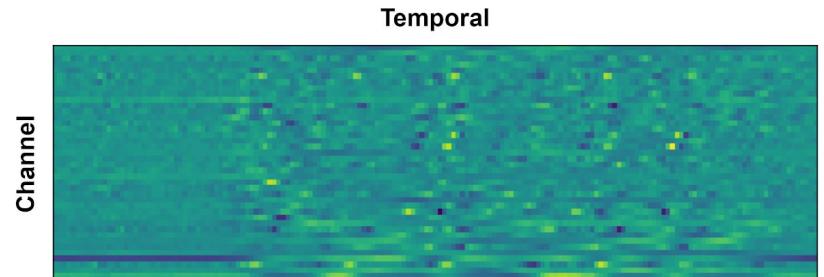
Transformer Encoder Representations from Alteration

Extending Mockingjay to multi-target learning on three dimensions

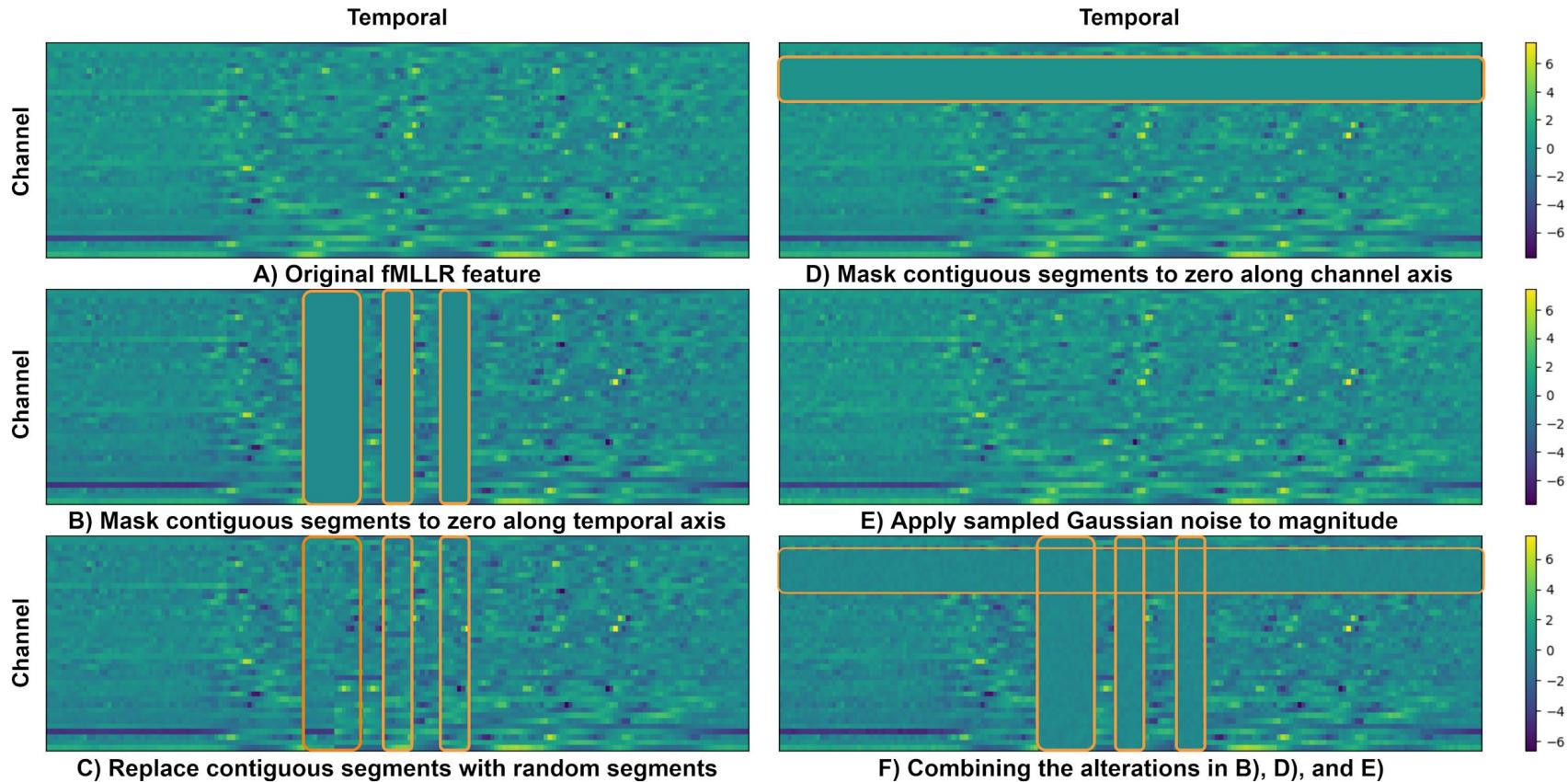
Recall: we mask mel spectrogram on time axis



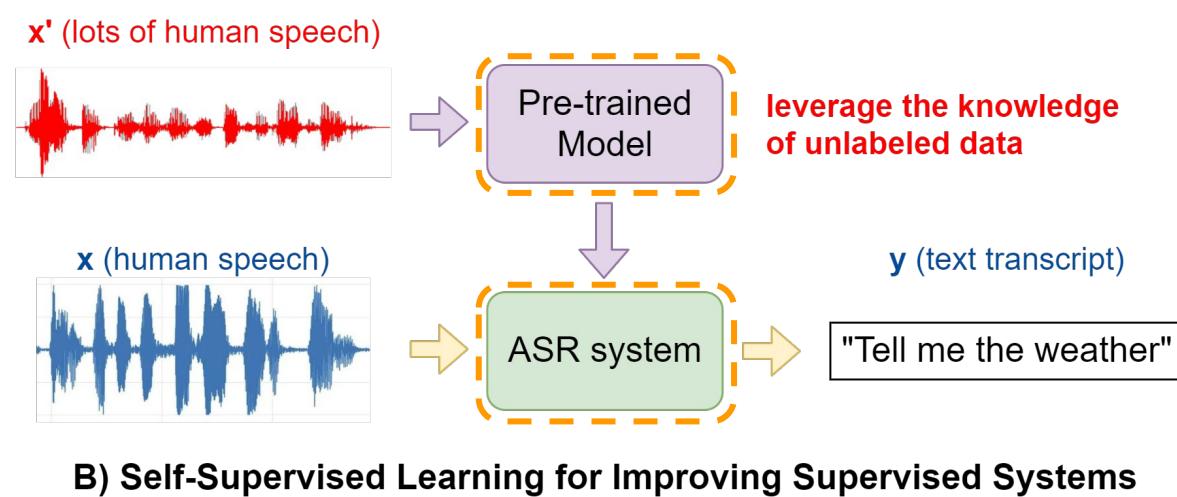
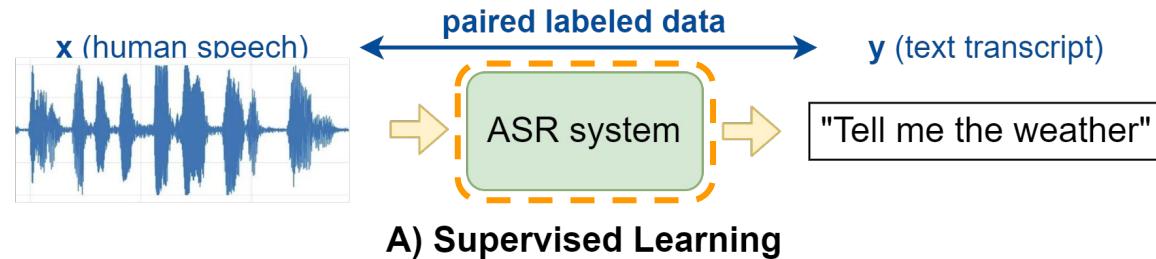
Consider fMLLR on 3 Axis:



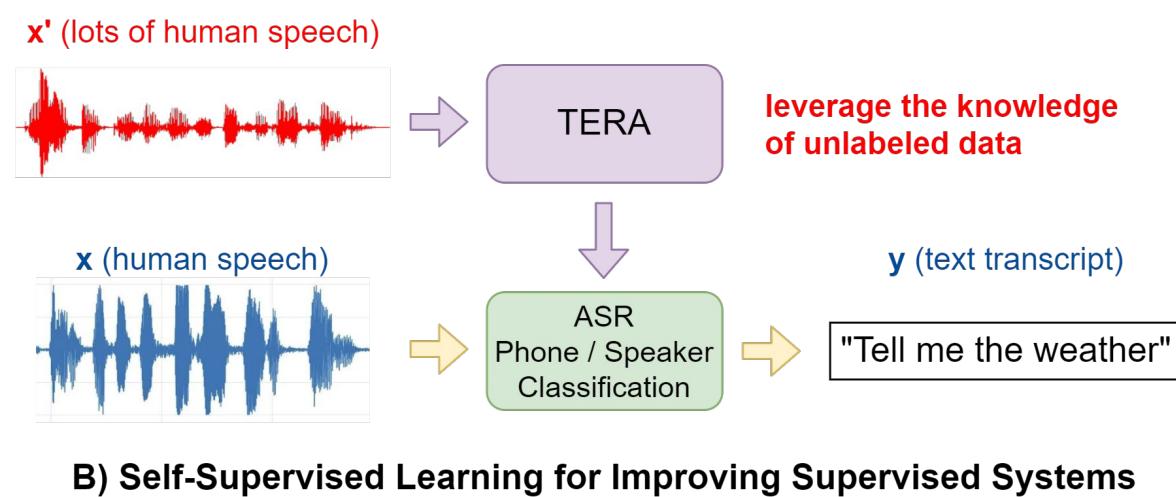
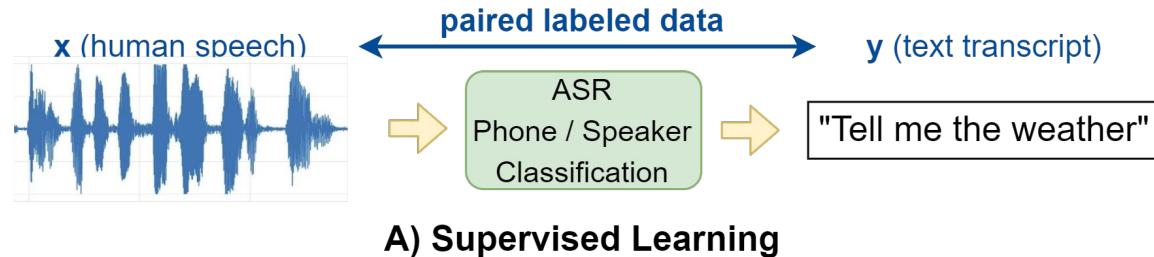
Multi-target Pre-training



Recall: Self-Supervised Learning for Speech

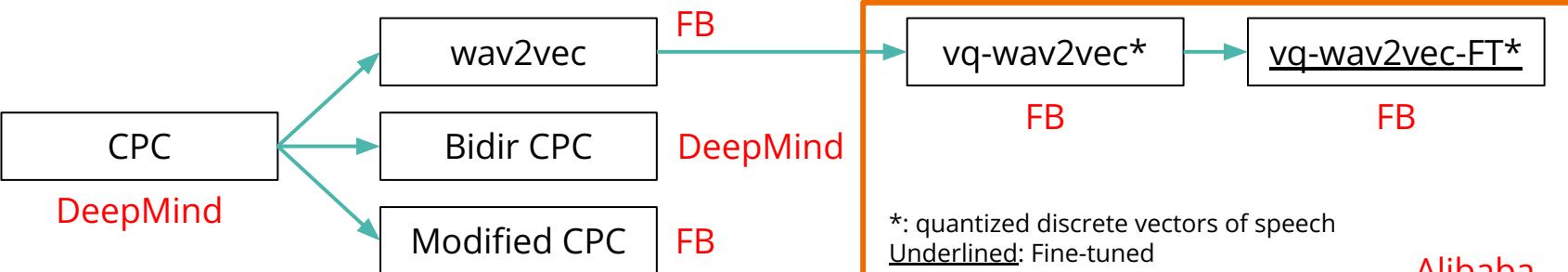


Self-Supervised Learning: TERA

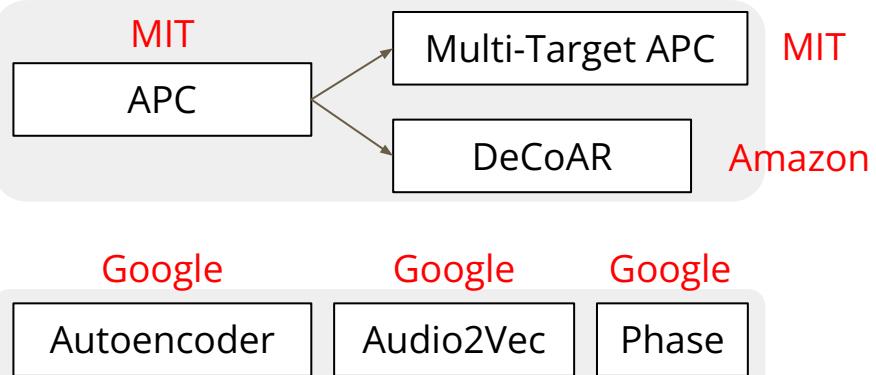


A Broad Introduction

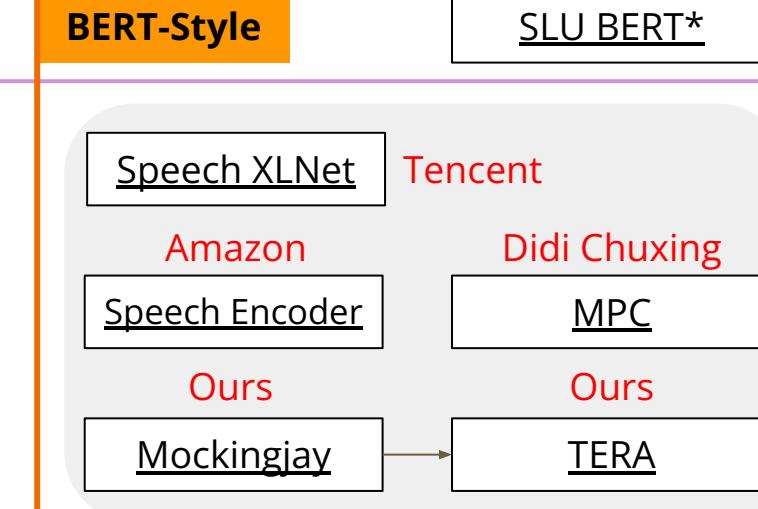
Contrastive Predictive Losses



Reconstruction Losses

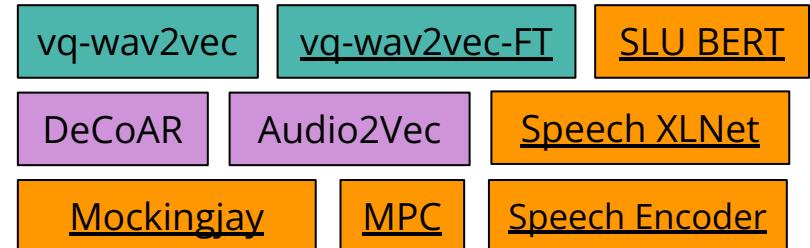
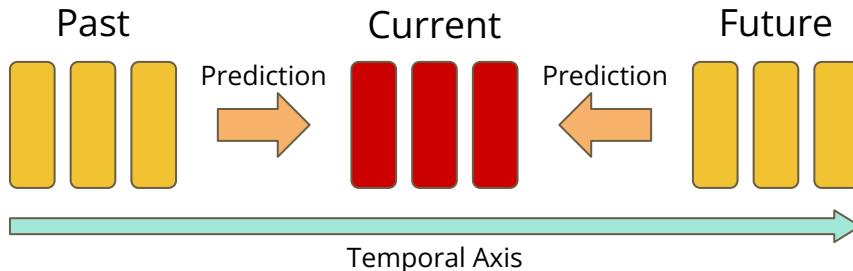
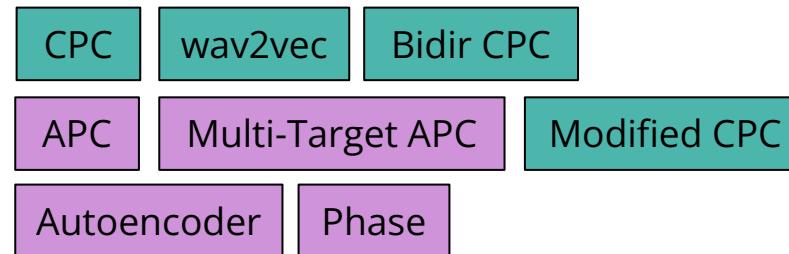
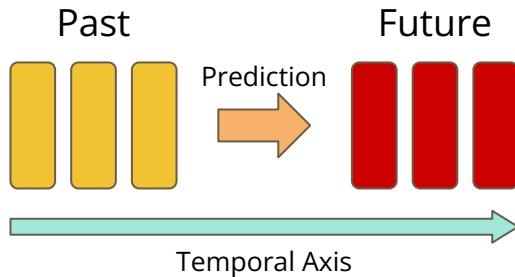


BERT-Style



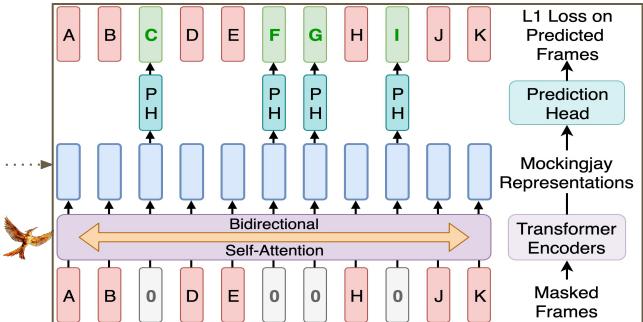
Quick Summary

The design of auxiliary task fundamentally defines what the model learns!



From here to beyond

- Mockingjay (ICASSP 2020)



Submitted to InterSpeech 2020 (5/15)

1. Mockingjay for Adversarial Defence
2. Improving Mockingjay: Speech ALBERT
3. Understanding Self-Attention

What else?

What else can we do with Mockingjay?

1. Adversarial Defense

Employ Mockingjay to protect models against adversarial attacks

1. Adversarial Defense

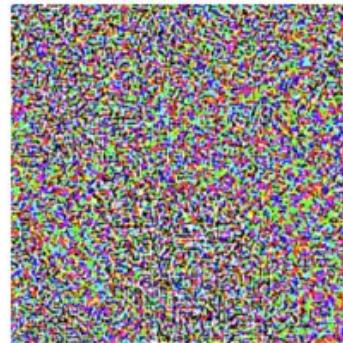
What is Adversarial Attack?



“panda”

57.7% confidence

+ .007 ×



noise

=



“gibbon”

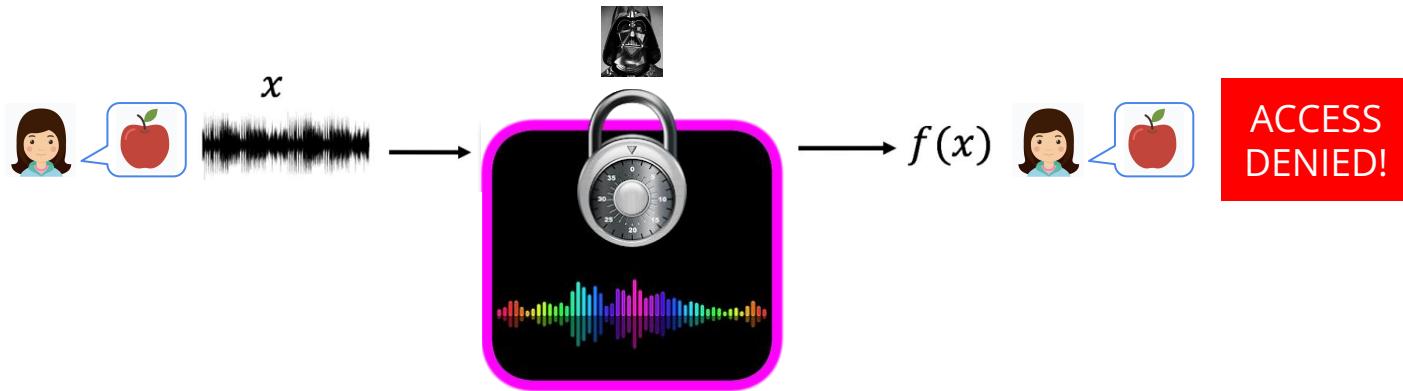
99.3% confidence



Hacking AI security systems: Face ID / Voice ID

1. Adversarial Defense

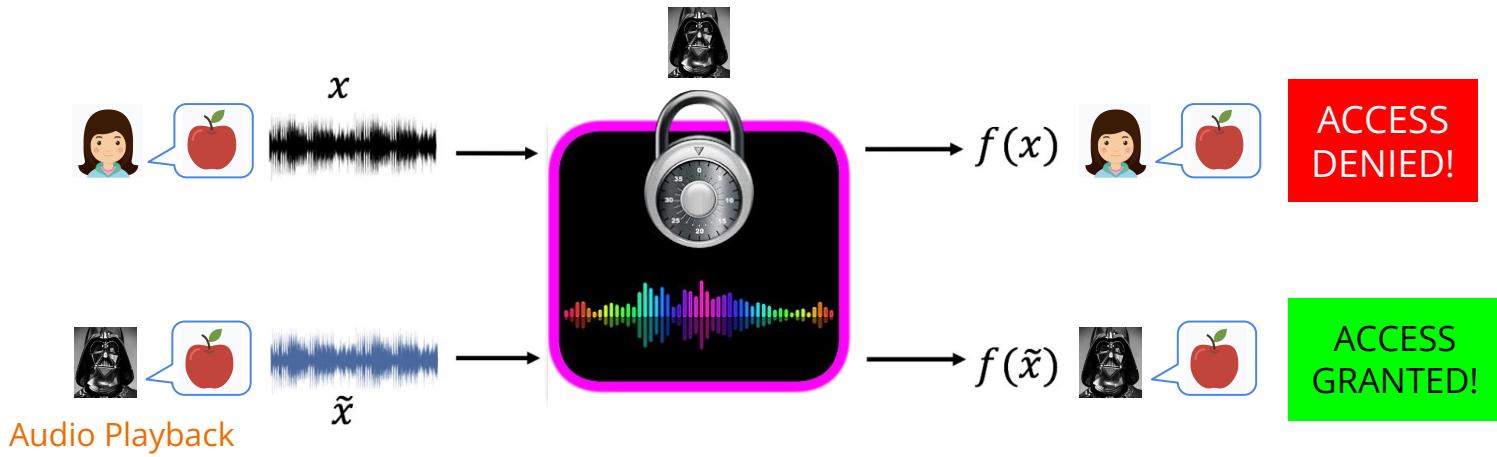
What is Adversarial Attack?



Hacking AI security systems: Face ID / Voice ID

1. Adversarial Defense

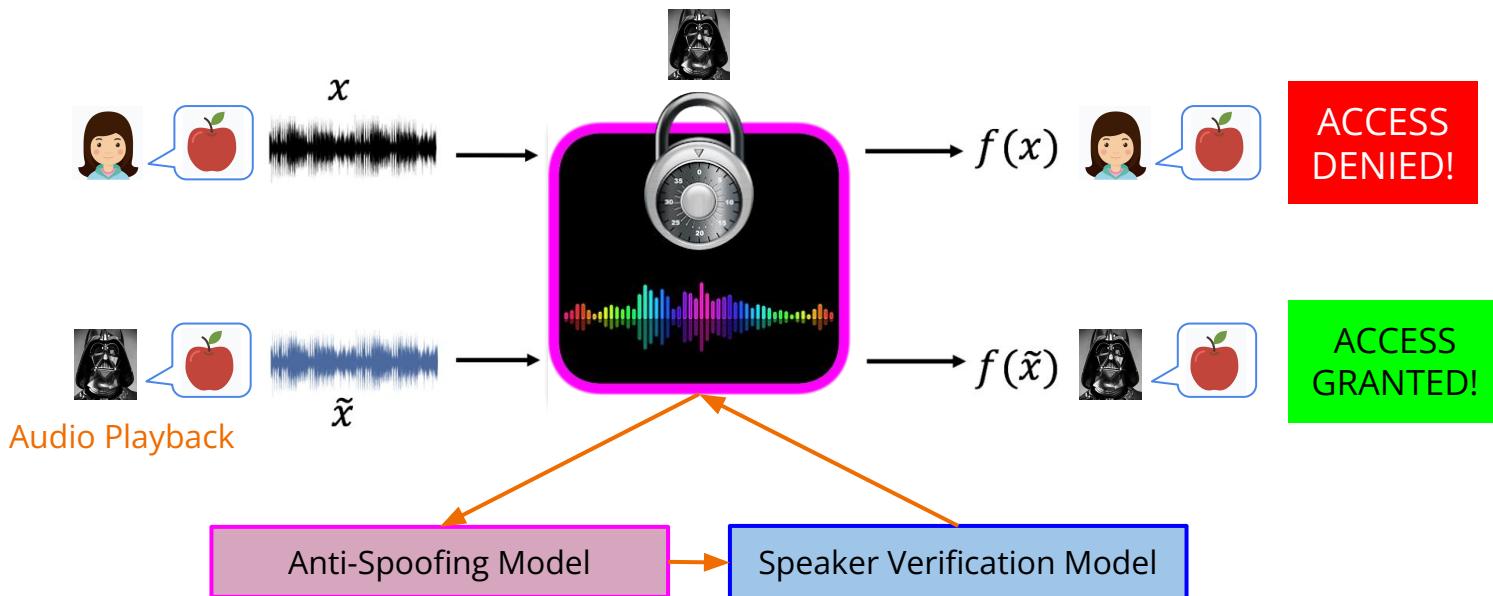
What is Adversarial Attack?



Hacking AI security systems: Face ID / Voice ID

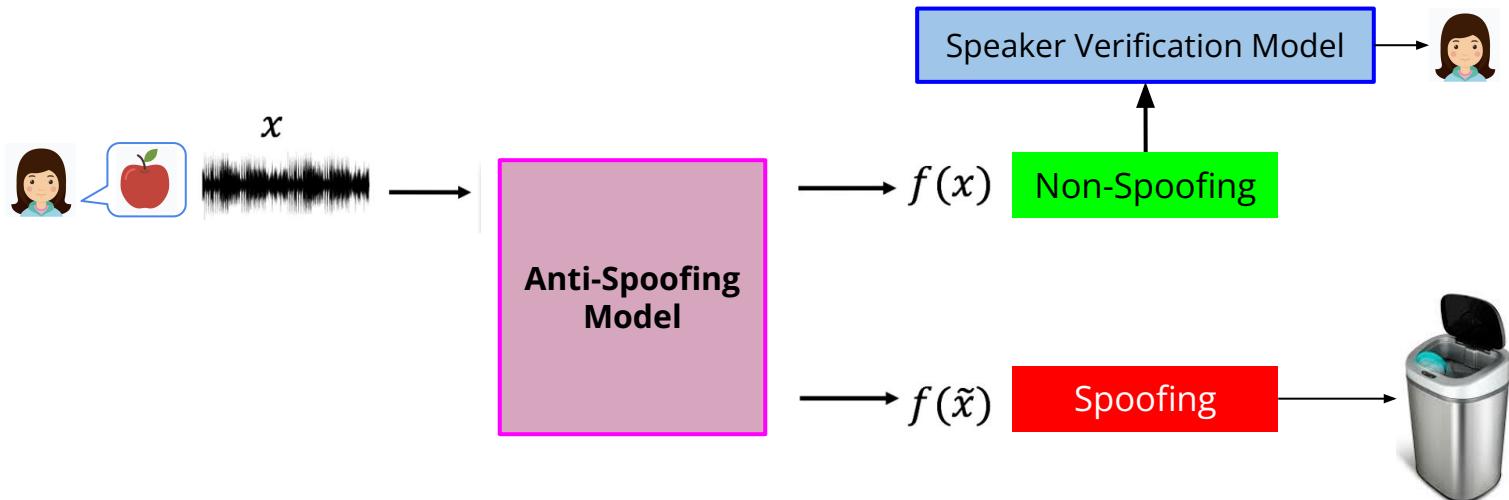
1. Adversarial Defense

Anti-Spoofing Model



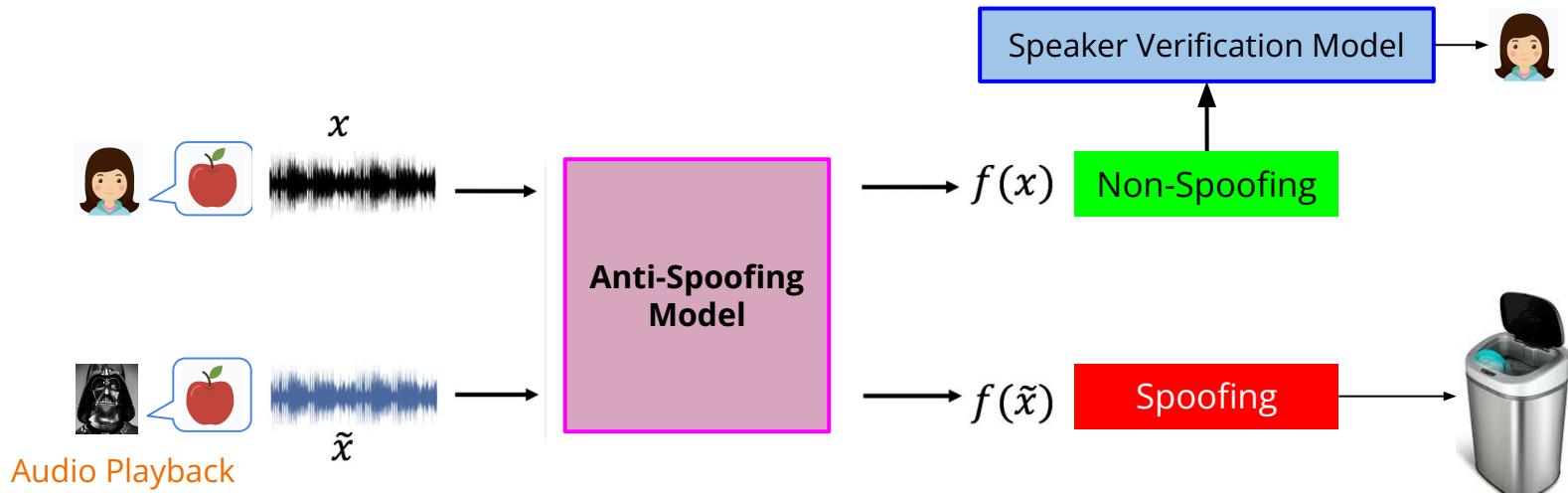
1. Adversarial Defense

Anti-Spoofing Model



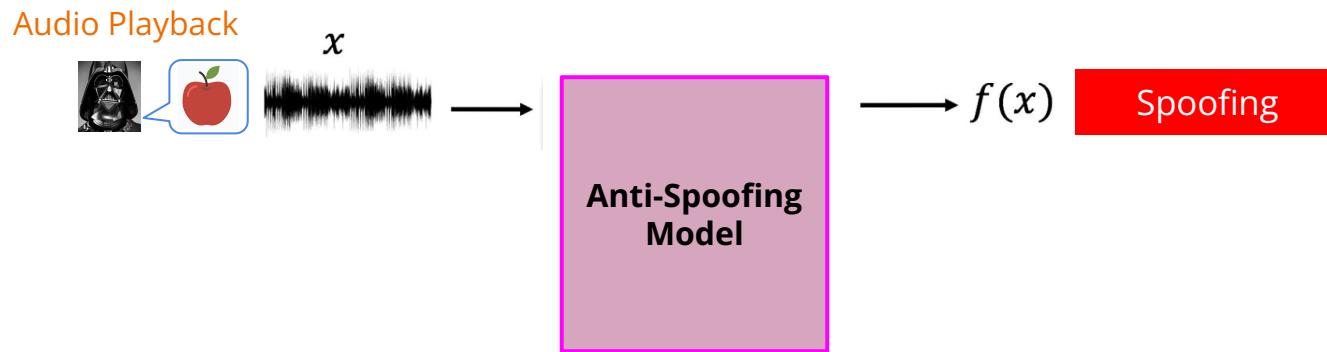
1. Adversarial Defense

Anti-Spoofing Model



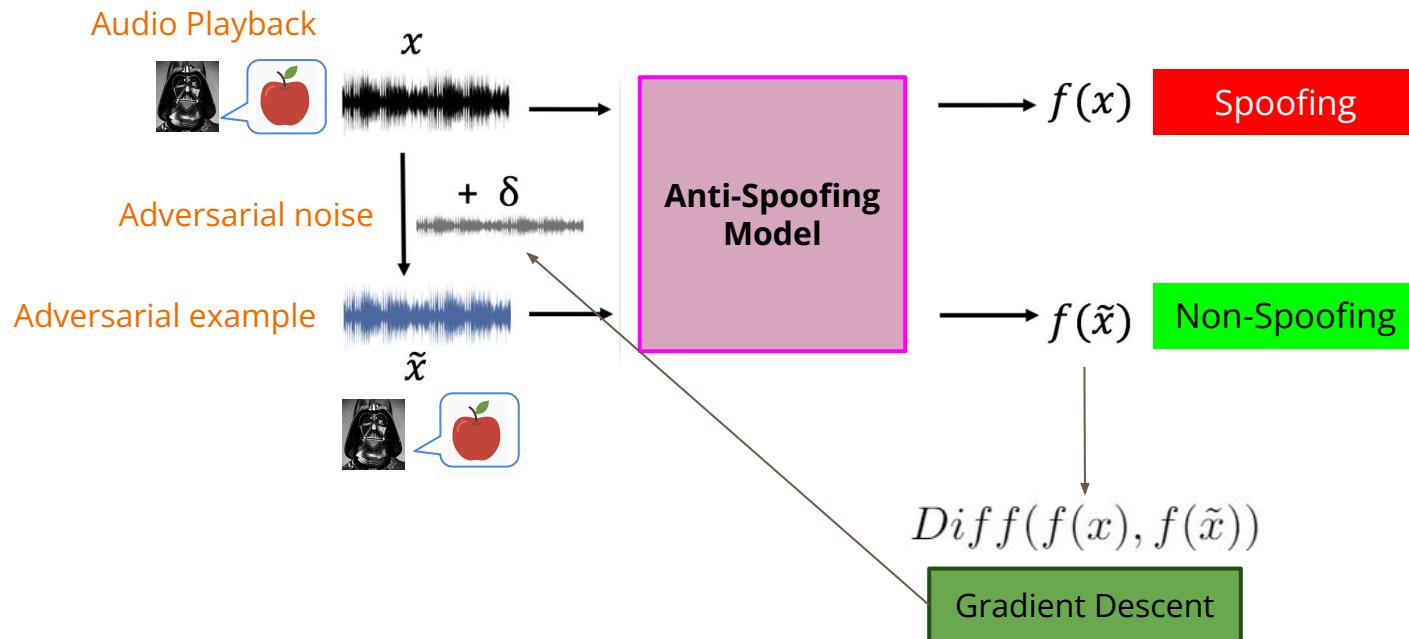
1. Adversarial Defense

How to Attack?



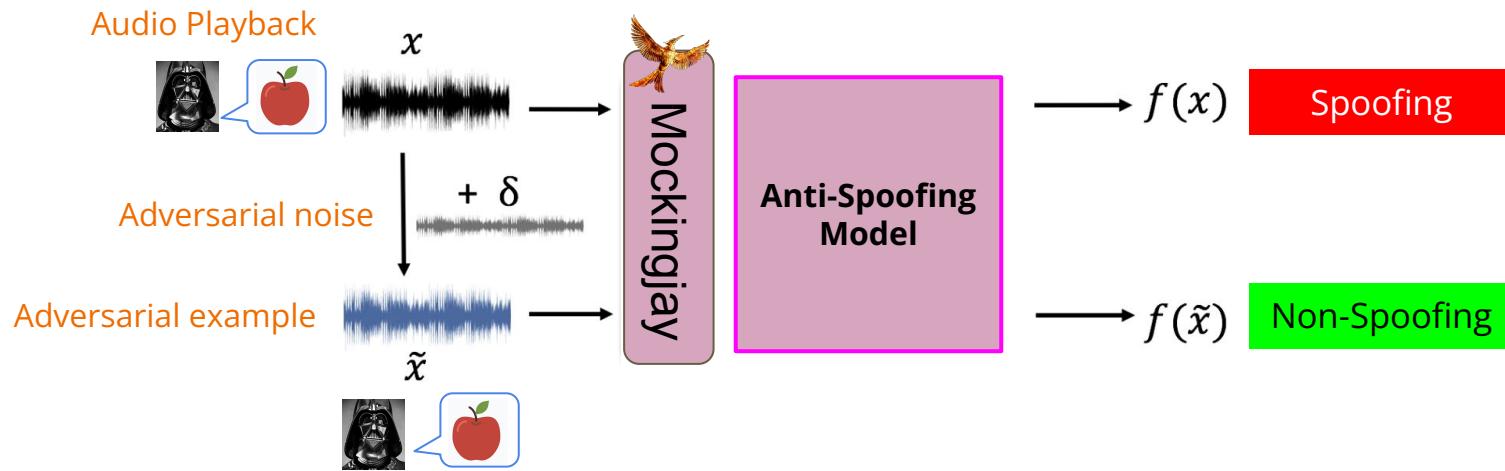
1. Adversarial Defense

How to Attack?



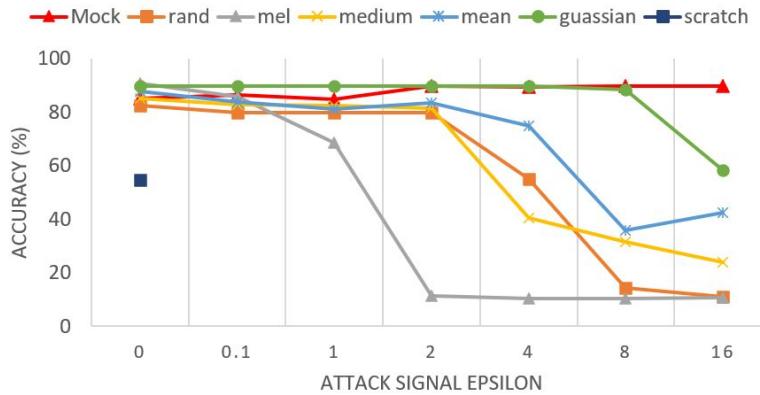
1. Adversarial Defense

Employing Mockingjay

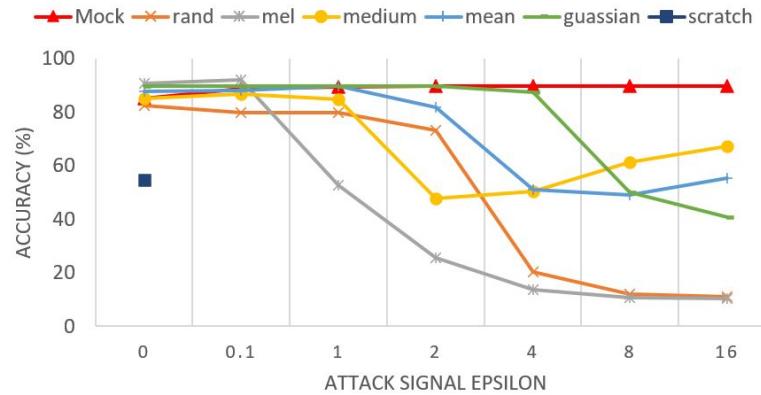


1. Adversarial Defense - Experiments

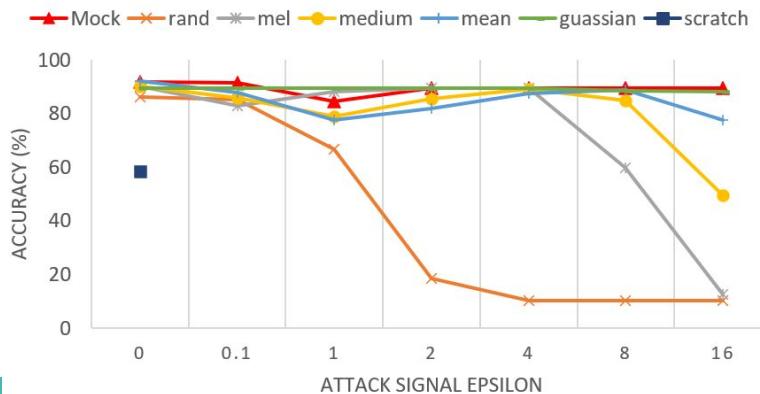
A) ATTACKING LCNN WITH PGD



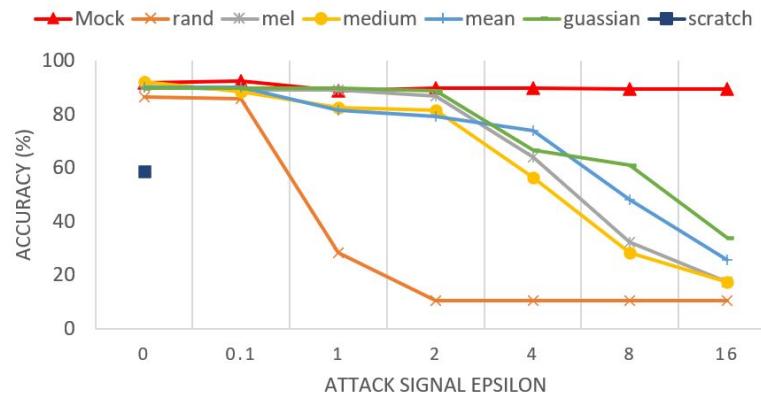
B) ATTACKING LCNN WITH FGSM



C) ATTACKING SENET WITH PGD

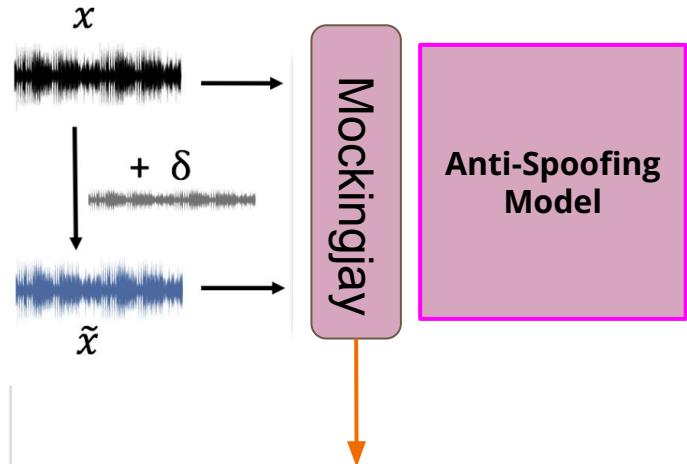
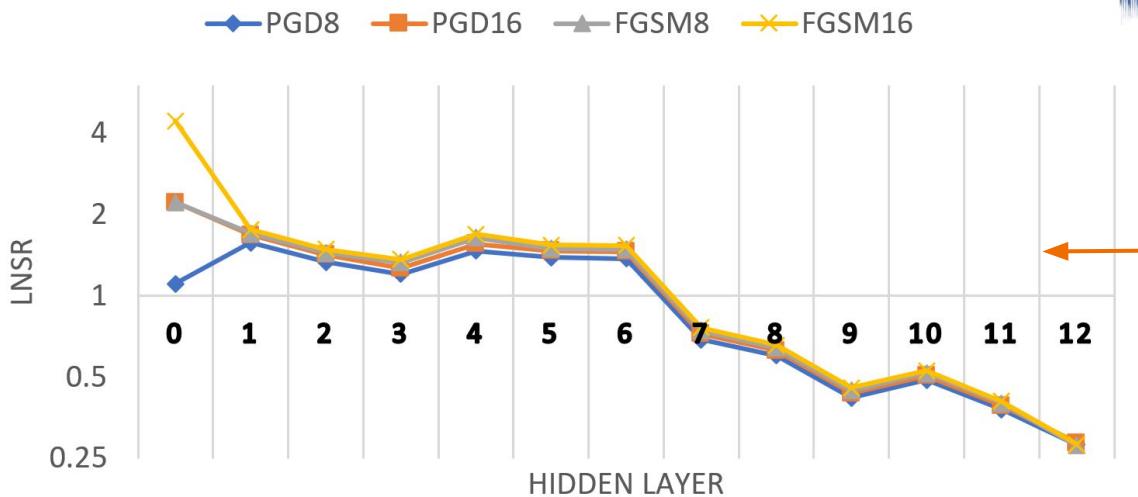


D) ATTACKING SENET WITH FGSM



1. Adversarial Defense - Experiments

HIDDEN DIFFERENCES OVER ALL LAYERS



Intuition:
LNSR- Measure the amount of
adversarial signal through the layers

$$LNSR_i = \sum_{n=1}^N \frac{\|\hat{h}_i^n - h_i^n\|_2}{\|h_i^n\|_2}$$

What else?

What else can we do with Mockingjay?

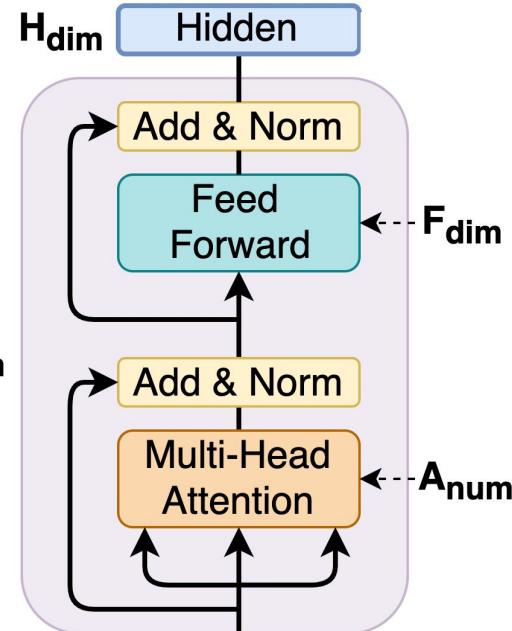
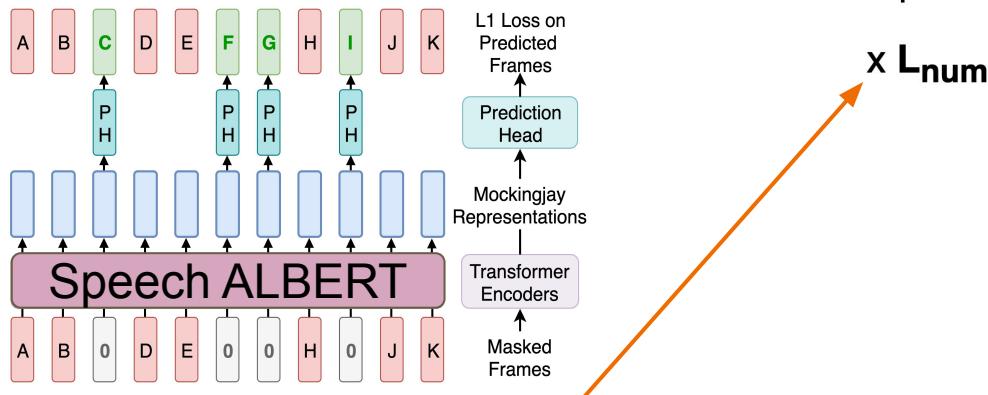
Talk 2: Audio ALBERT

Recall that:

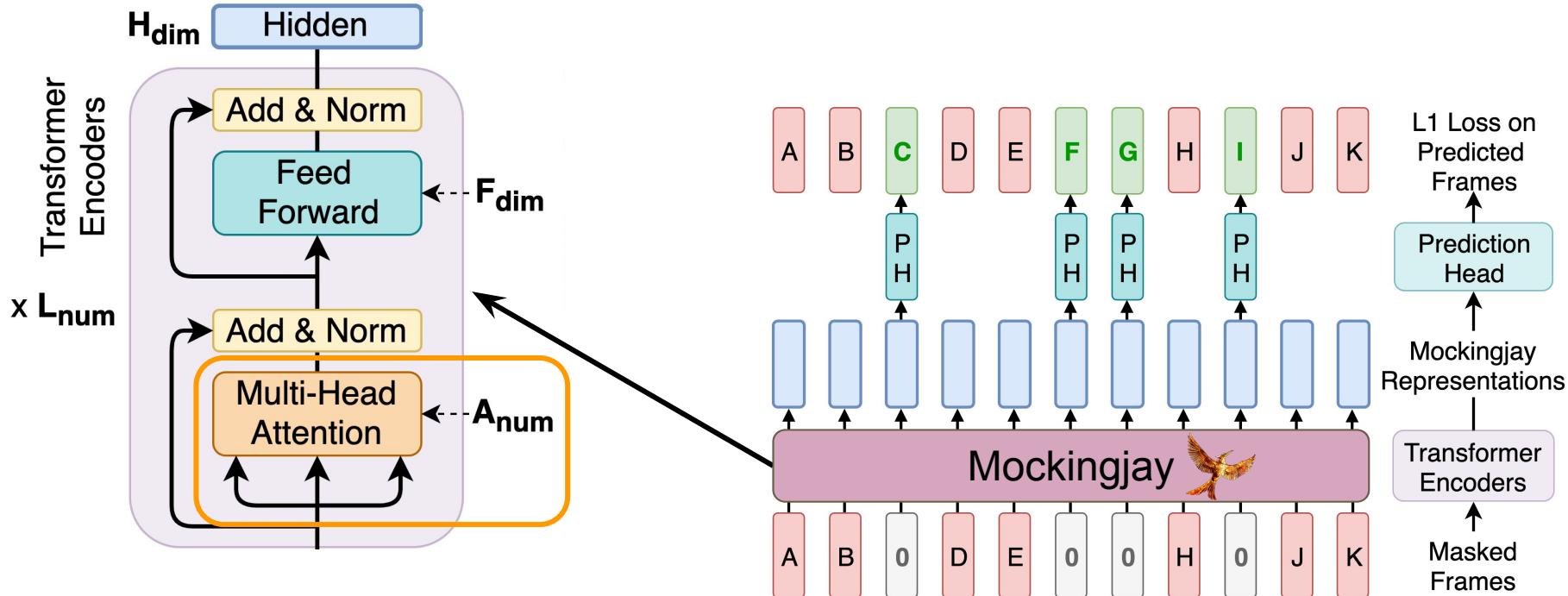
Audio BERT >>> Mockingjay

Now we also have:

ALBERT >>> Audio ALBERT



Talk 3: Understanding Self-Attention



Conclusion

**Self-supervised learning,
a brand new topic with lots of ideas that we can work on!**

Thank You

Q&A