# Introduction

## NLP BERT:
## Language Representation Learning

# Introduction

## NLP BERT:
## Language Representation Learning

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

| R | R | R | R | R | Representations! |

**BERT**

| A | B | C | D | E | ← Text tokens |

Unsupervised pre-train on **text**

# Introduction



**NLP BERT:**
**Language Representation Learning**

**Speech BERT:**
**Speech Representation Learning**

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

R R R R R

BERT

Unsupervised pre-train on **text**

A B C D E

R R R R R   Representations!

Speech BERT

Acoustic Frames

Unsupervised pre-train on **speech**

# Introduction



**NLP BERT:**
**Language Representation Learning**

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

R R R R R

BERT

A B C D E

Unsupervised pre-train on **text**

**Speech BERT:**
**Speech Representation Learning**

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

R R R R R Representations!

Speech BERT

Acoustic Frames

Unsupervised pre-train on **speech**
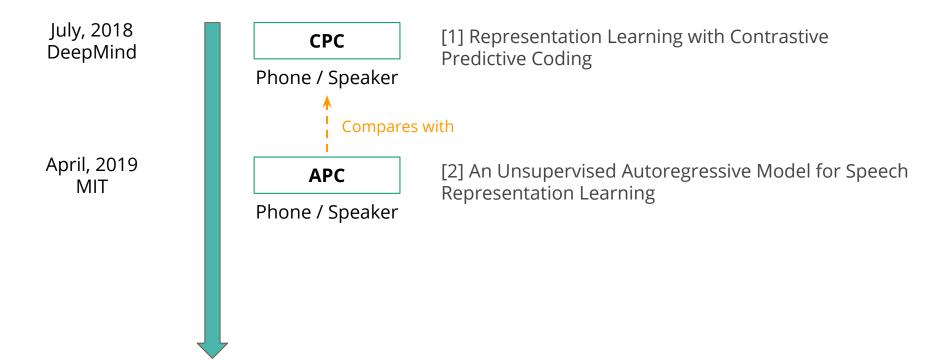
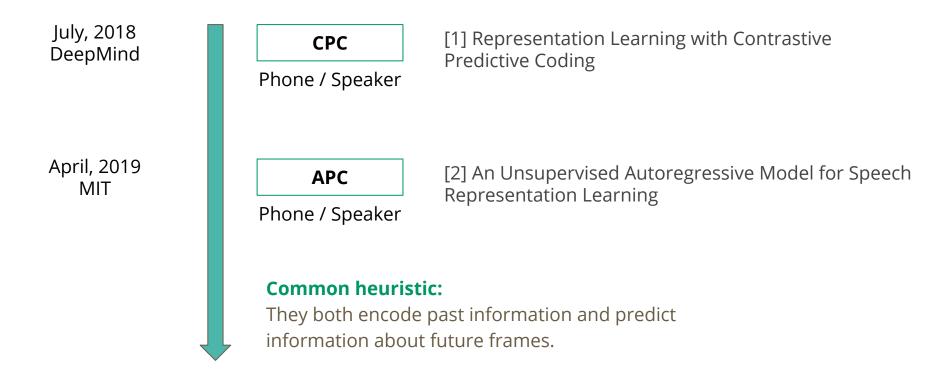# A View of Recent Unsupervised Speech Representation Learning Approaches
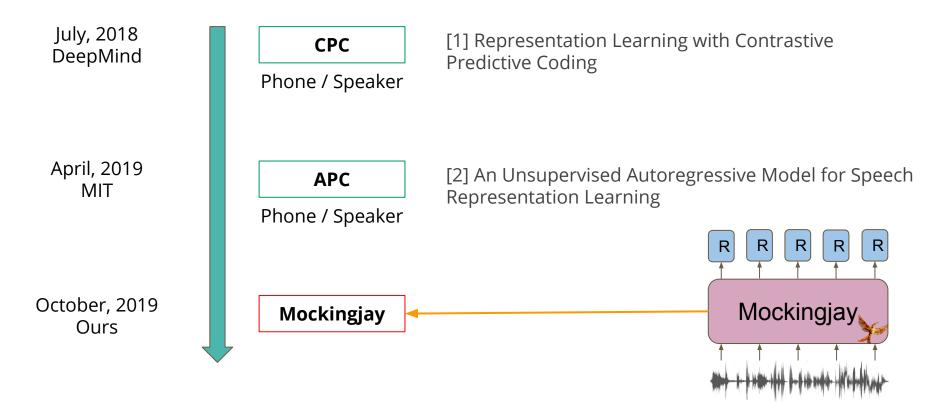
July, 2018
DeepMind

CPC

Phone / Speaker

[1] Representation Learning with Contrastive Predictive Coding

# A View of Recent Unsupervised Speech Representation Learning Approaches

July, 2018
DeepMind

CPC

Phone / Speaker

[1] Representation Learning with Contrastive Predictive Coding

Compares with

April, 2019
MIT

APC

Phone / Speaker

[2] An Unsupervised Autoregressive Model for Speech Representation Learning

# A View of Recent Unsupervised Speech Representation Learning Approaches

July, 2018
DeepMind

**CPC**

Phone / Speaker
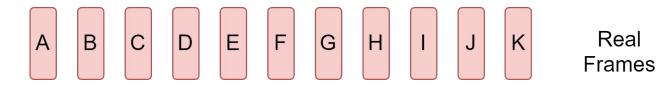
[1] Representation Learning with Contrastive Predictive Coding

April, 2019
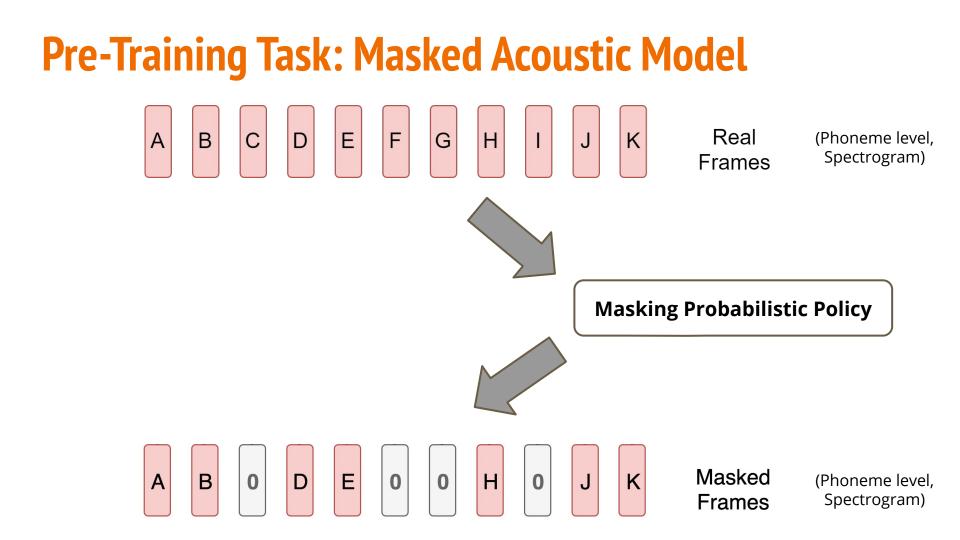MIT

**APC**

Phone / Speaker

[2] An Unsupervised Autoregressive Model for Speech Representation Learning

**Common heuristic:**
They both encode past information and predict information about future frames.

# A View of Recent Unsupervised Speech Representation Learning Approaches

July, 2018
DeepMind

**CPC**

Phone / Speaker

[1] Representation Learning with Contrastive Predictive Coding

April, 2019
MIT

**APC**

Phone / Speaker

[2] An Unsupervised Autoregressive Model for Speech Representation Learning

October, 2019
Ours

**Mockingjay**

| R | R | R | R | R |

Mockingjay

# A View of Recent Unsupervised Speech Representation Learning Approaches

**CPC**

Phone / Speaker

[1] Representation Learning with Contrastive Predictive Coding

Compares with

**APC**

Phone / Speaker

[2] An Unsupervised Autoregressive Model for Speech Representation Learning

Compares with

October, 2019
Ours

**Mockingjay**

Phone / Speaker / Sentiment

R  R  R  R  R

Mockingjay

# Pre-Training Task: Masked Acoustic Model

A B C D E F G H I J K    Real Frames    (Phoneme level, Spectrogram)

# Pre-Training Task: Masked Acoustic Model

# 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)** A B C D E ...

**2)** 80% → A 0 C 0 E ...

Mask all 15%

**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 untouch **10%** of the time

# Probabilistic Policy for Masking Frames

**1)** A B C D E ...

**2)**

80% → A 0 C 0 E ...
Mask all 15%

10% → A G C Y E ...
Replace all 15%

**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 untouch **10%** of the time

# Probabilistic Policy for Masking Frames



**1)** A B C D E ...

**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 untouch **10%** of the time

**2)**

80% → A 0 C 0 E ...
Mask all 15%

10% → A G C Y E ...
Replace all 15%

10% → A B C D E ...
Do nothing, frames remain the same

# Input Feature: Masked Spectrogram

# Input Feature: Masked Spectrogram



Masked to Zero

80-dim mel-spectrogram

First derivative

Visualizations

Visualizations

The model was able to reconstruct spectrogram form hidden representations

Real

Pred

Repr.

Masked Frames

L1 Loss on Predicted Frames

A B **C** D E **F** **G** H **I** J K

P H    P H P H    P H

Prediction Head

Mockingjay Representations

Mockingjay

Transformer Encoders

A B 0 D E 0 0 H 0 J K

Masked Frames

**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:
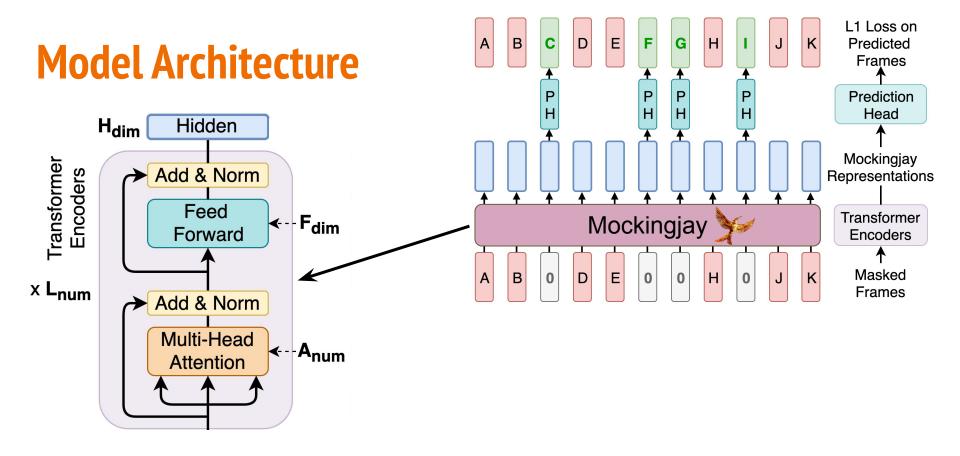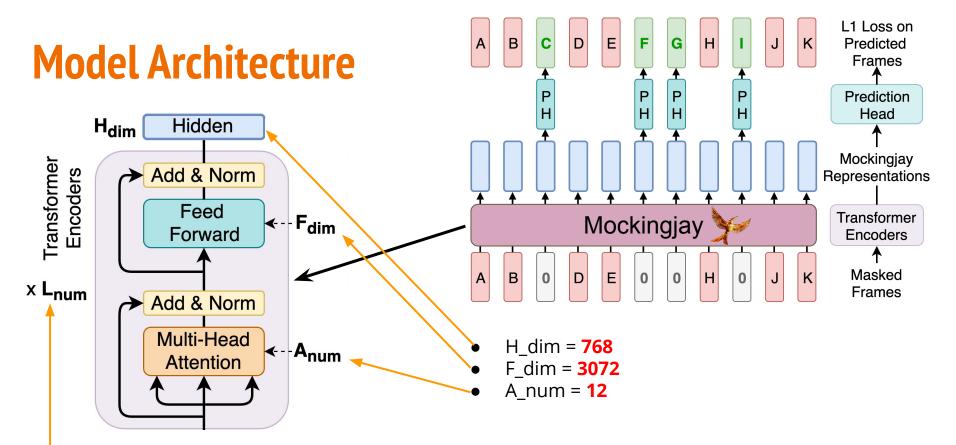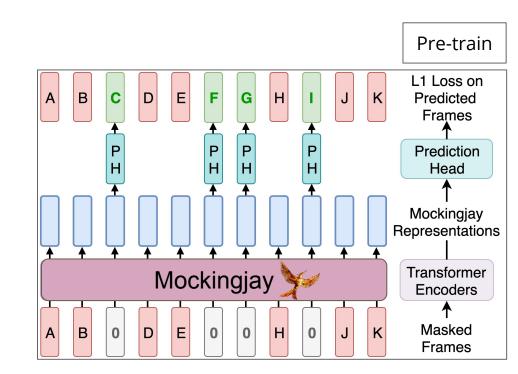*Downsampling*, and  *consecutive masking*

# 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**

# 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

R=3, C=3

# Model Architecture

# Model Architecture

# Model Architecture



- H_dim = **768**
- F_dim = **3072**
- A_num = **12**

- BASE (L=3)
- LARGE (L=12)

- Train on LibriSpeech **360 hrs**
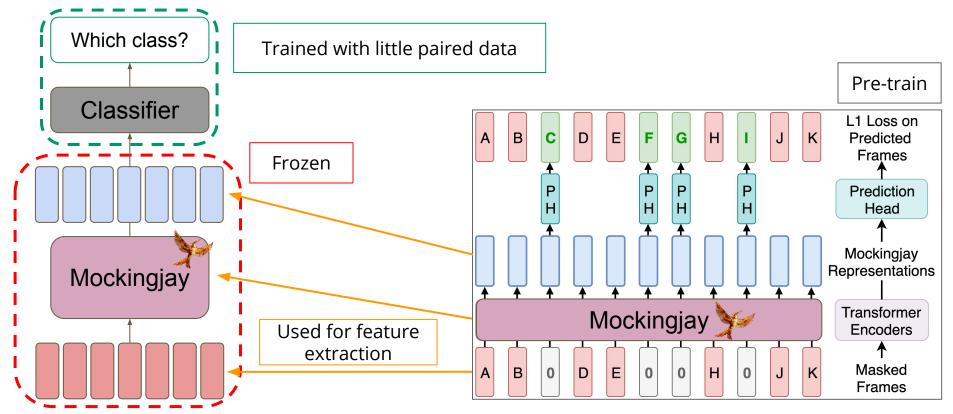- Pre-train steps = **500k**
- Fine-tune steps = **50k** (2-epochs)

# Incorporating with Downstream Tasks

## 1) Feature Extraction

# Incorporating with Downstream Tasks

## 1) Feature Extraction

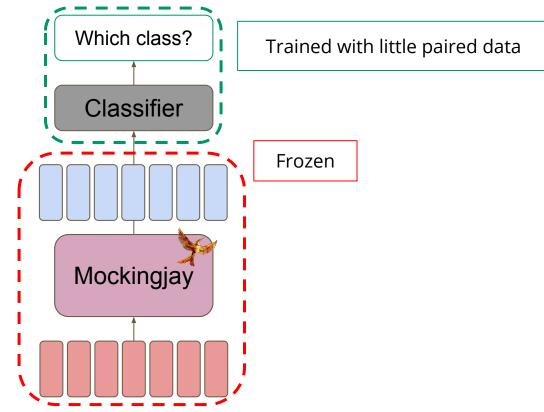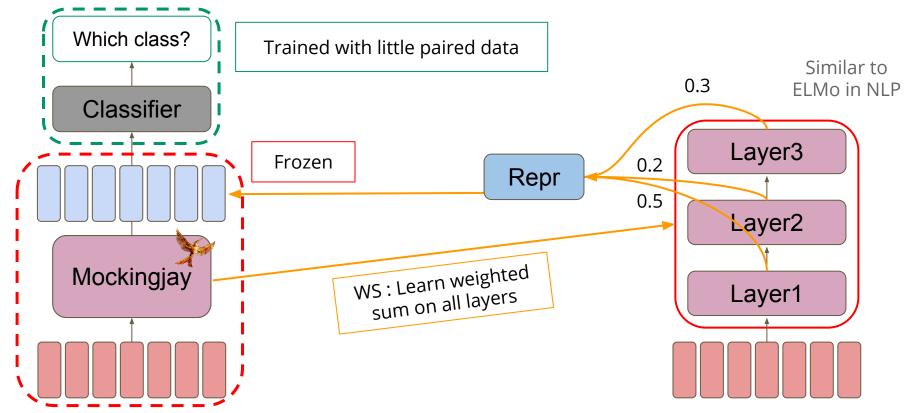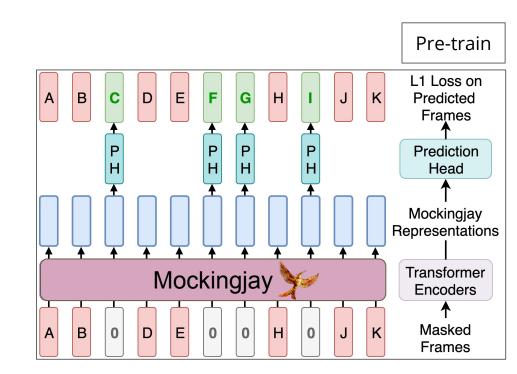# Incorporating with Downstream Tasks

## 1) Feature Extraction

# Incorporating with Downstream Tasks

## 1) Feature Extraction

# Incorporating with Downstream Tasks
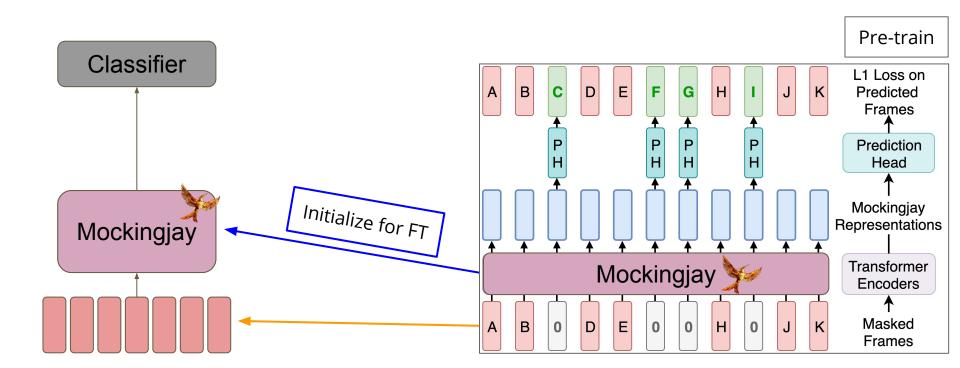
## 2) Weighted Sum from All Layers (WS)

Which class?

Trained with little paired data

Classifier

Frozen

Mockingjay

# Incorporating with Downstream Tasks

## 2) Weighted Sum from All Layers (WS)

Which class?

Trained with little paired data

Classifier

Frozen

Mockingjay

Repr

0.3

0.2

0.5

Similar to ELMo in NLP
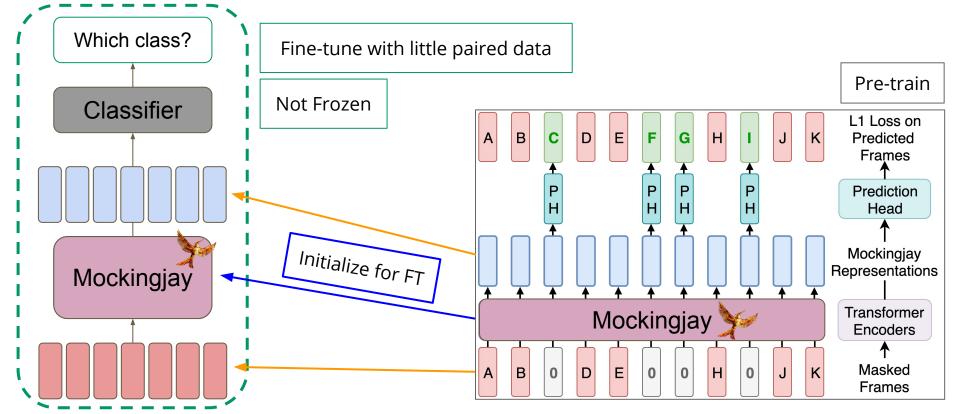
Layer3

Layer2

Layer1

WS : Learn weighted sum on all layers

# Incorporating with Downstream Tasks

## 3) Fine-tune (FT2)

# Incorporating with Downstream Tasks

## 3) Fine-tune (FT2)

# Incorporating with Downstream Tasks
## 3) Fine-tune (FT2)

# Experiments

We report results on 3 different downstream tasks:

- Phoneme Classification

- Speaker Recognition

- Sentiment Classification on spoken content

# Experiments

We report results on 3 different downstream tasks:

- ### Phoneme Classification (72 classes):
  Train: LibriSpeech 360 / Test: LibriSpeech test-clean

features     →     Feed-forward Classifier     →     0 0 0 3 3 8 8   phone

- ### Speaker Recognition

- ### Sentiment Classification on spoken content

# Experiments

We report results on 3 different downstream tasks:

- Phoneme Classification (72 classes):
  Train: LibriSpeech 360 / Test: LibriSpeech test-clean



- Speaker Recognition (63 classes):
  Train: 90% of LibriSpeech 100 / Test: 10% of LibriSpeech 100

- Sentiment Classification on spoken content (2 classes):
  To demonstrate domain invariant transferability, we use another dataset: MOSEI [3]

# Experiments - 1/3

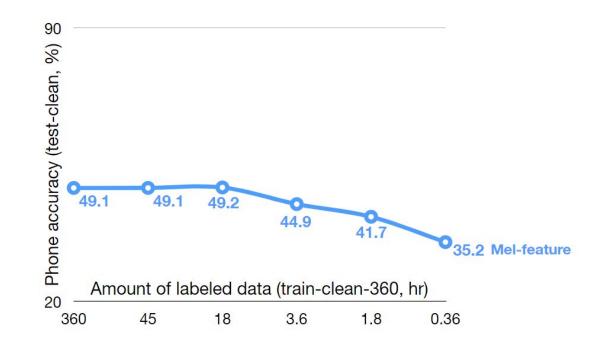| Acoustic Features | Phoneme Classification | Speaker Recognition | Sentiment Classification |
|---|---|---|---|
| Mel Features | 49.1 | 70.1 | 64.6 |
| BASE | 60.9 | 94.5 | 67.4 |
| LARGE | **64.3** | **96.3** | **70.1** |

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

# Experiments - 3/3

| 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** |
| BASE-FT2 | **84.3** | **98.1** | 68.5 |
| APC [2] | 74.1 | 85.9 | 66.0 |

[2] An Unsupervised Autoregressive Model for Speech Representation Learning

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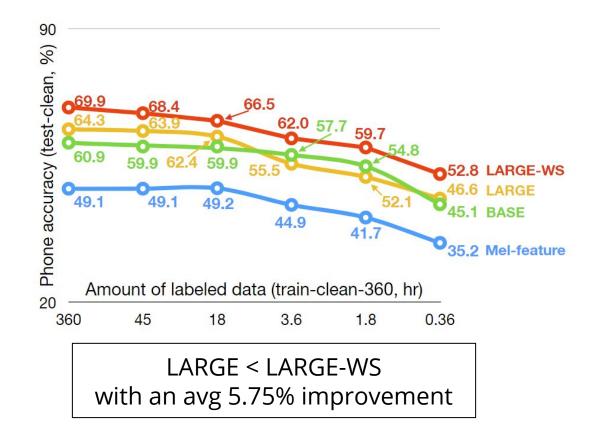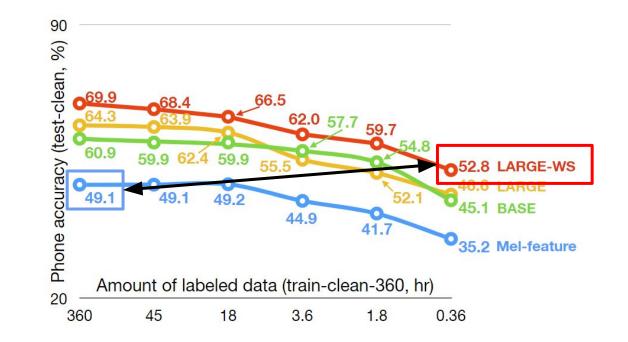
Mel < BASE

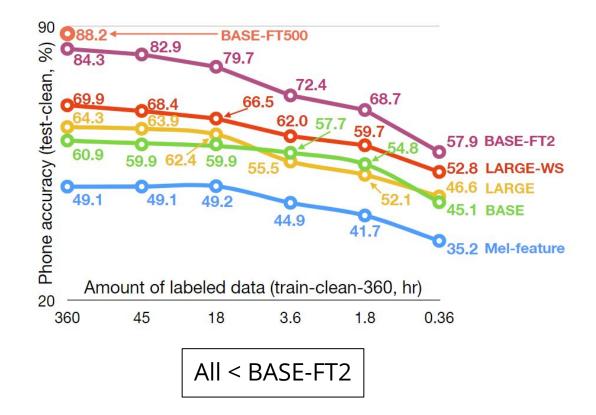# Low-Resource Experiments - 3/6



Mel < BASE < LARGE

LARGE < LARGE-WS
with an avg 5.75% improvement
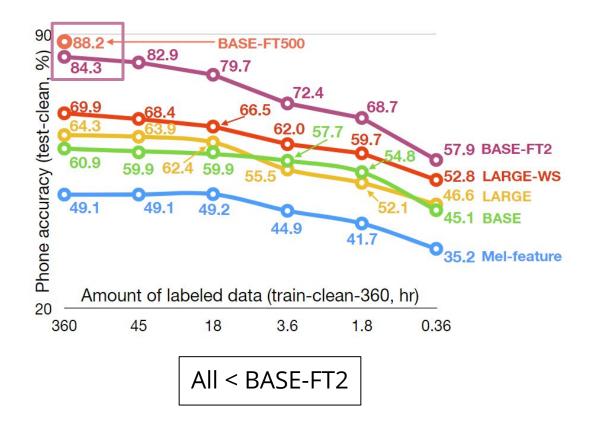
With 0.1% of labels,
LARGE-WS (52.8%) outperformed Mel (49.1%) that uses all 100% hours of labeled data.

# Low-Resource Experiments - 5/6



All < BASE-FT2

All < BASE-FT2

# Low-Resource Experiments - 5/6



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

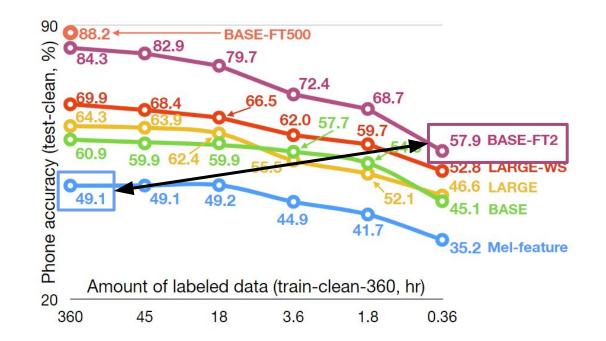APC works well on full resource but fails to generalize for limited labeled data.

# Conclusion

We conclude that unsupervised Mockingjay improves supervised training!

# Links

This slide (with speaker notes) can be found here:
https://bit.ly/icassp2020-mockingjay

Our code and implementation can be found here:
https://github.com/andi611/Mockingjay-Speech-Representation