# Pretext tasks
# 4. CPC

**|| AR / LPC / <span style="color:red">VAR / RNN / APC</span> / CPC ||**

| | | | |
|---|---|---|---|
| **1** **2** | SELF-PREDICTION | INNATE RELATIONSHIP (Context-based) | 1. ROTATION<br><br>2. RELATIVE POSITION | IMAGE |
| **3** | CONTRASTIVE LEARNING | INTER-SAMPLE CLASSIFICATION | 1. Instance Discrimination<br>2. SimCLR [Contrastive Loss]<br>3. Theory – Guarantees / Bounds | IMAGE |
| **4** | CONTRASTIVE LEARNING | INTER-SAMPLE CLASSIFICATION | Contrastive Predictive Coding (CPC), [NCE, InfoNCE Loss] | AUDIO/ SPEECH |
| **5** | SELF-PREDICTION | GENERATIVE (VAE) | 1. AE – Variational Bayes<br><br>2. VQ-VAE + AR | IMAGE<br><br>AUDIO/ SPEECH |
| **6** | SELF-PREDICTION | GENERATIVE (AR) | 1. AR-LM – GPT<br>2. Masked-LM – BERT | LANGUAGE |
| **7** | SELF-PREDICTION | MASKED-GEN (Masked LM for ASR) | 1. Wav2Vec / 2.0<br>2. HuBERT | AUDIO/ SPEECH |

# Learning with or without supervision – speech and audio

- Next frame prediction



- Masked prediction
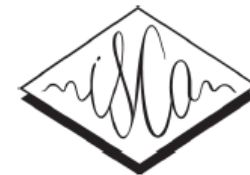


- Future prediction
  - To predict future audio features from the historical ones
    - Contrastive predictive coding (CPC) [Oord et al., 2018]
    - Autoregressive predictive coding (APC) [Chung et al., 2019]
    - wav2vec [Schneider et al., 2019]

- [Oord et al., 2018] Representation learning with contrastive predictive coding, arXiv
- [Chung et al., 2019] An unsupervised autoregressive model for speech representation learning, Interspeech
- [Schneider et al., 2019] wav2vec: Unsupervised pre-training for speech recognition, Interspeech

[Oord et al., 2018] Representation learning with contrastive predictive coding, arXiv

[Chung et al., 2019] An unsupervised autoregressive model for speech representation learning, Interspeech

[Schneider et al., 2019] wav2vec: Unsupervised pre-training for speech recognition, Interspeech

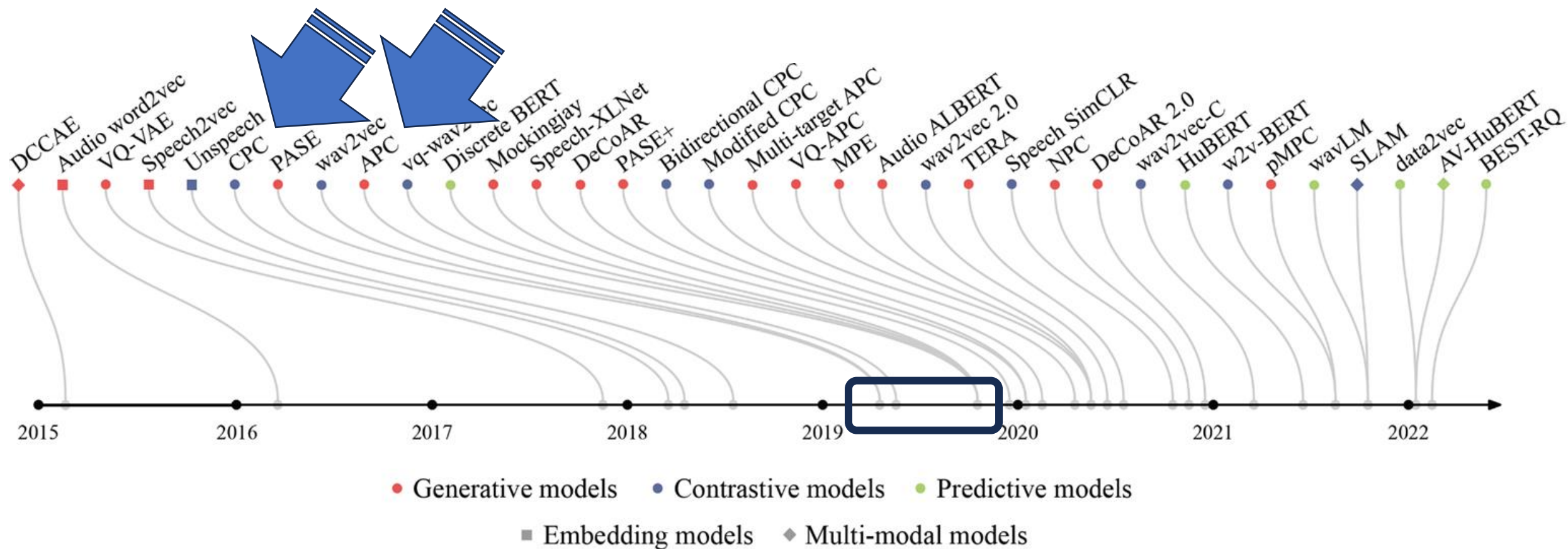# An Unsupervised Autoregressive Model for Speech Representation Learning

*Yu-An Chung, Wei-Ning Hsu, Hao Tang, James Glass*

Computer Science and Artificial Intelligence Laboratory
Massachusetts Institute of Technology
Cambridge, MA 02139, USA
{andyyuan,wnhsu,haotang,glass}@mit.edu

# Autoregressive Predictive Coding: A Comprehensive Study

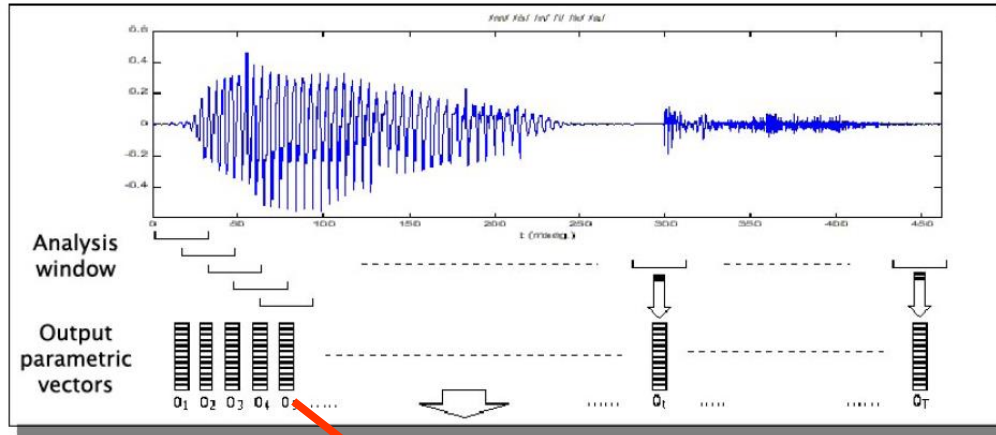Gene-Ping Yang , Sung-Lin Yeh, Yu-An Chung , James Glass , and Hao Tang

# Speech representation learning methods

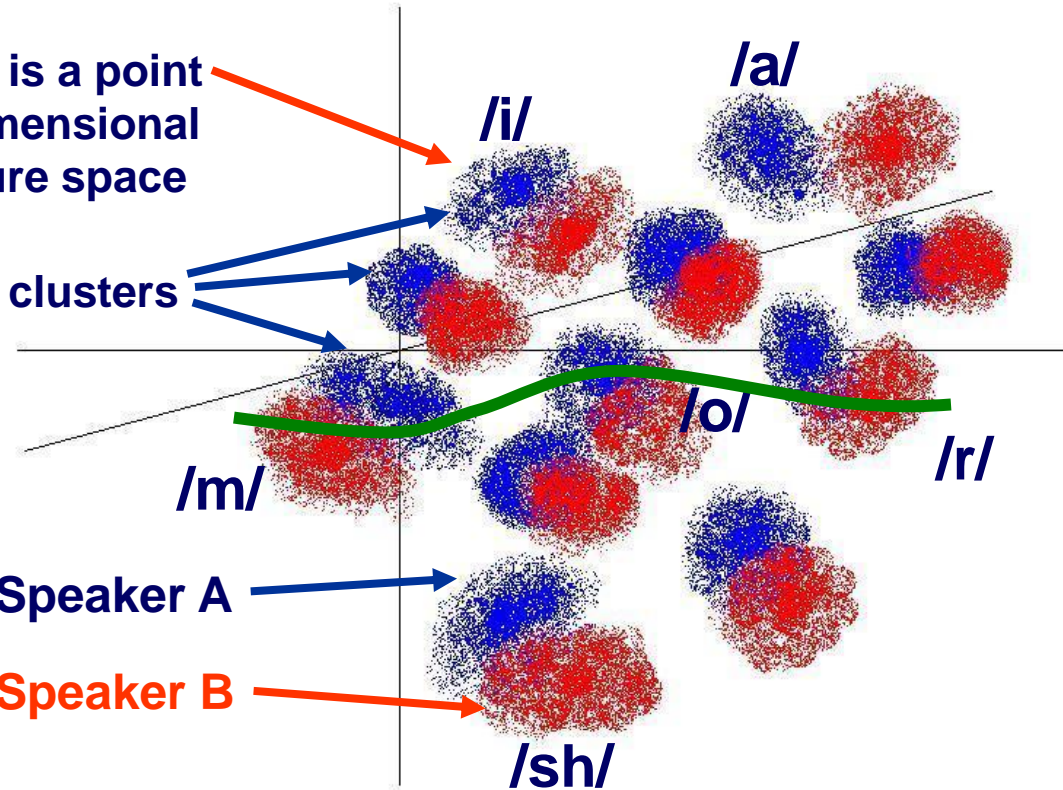## Short-time Analysis and Parameterization

**Feature Space**

Analysis window

Output parametric vectors

Bag of vectors representation of speaker's acoustic space

**Each vector is a point in the 13-dimensional MFCC feature space**

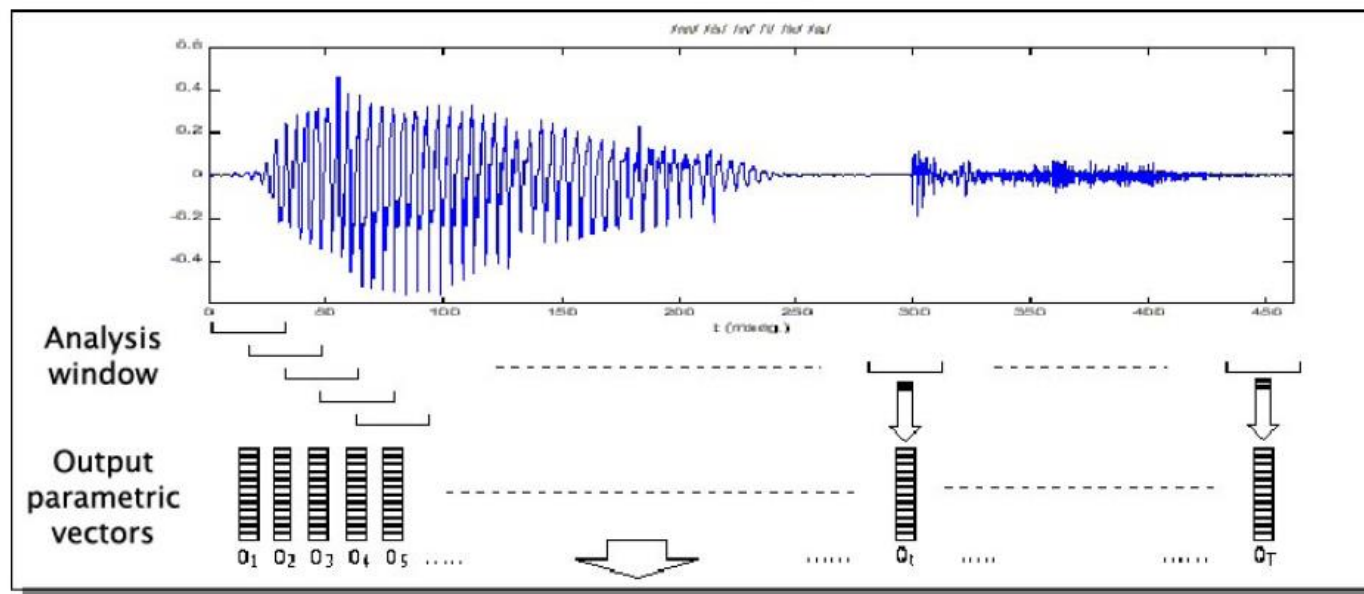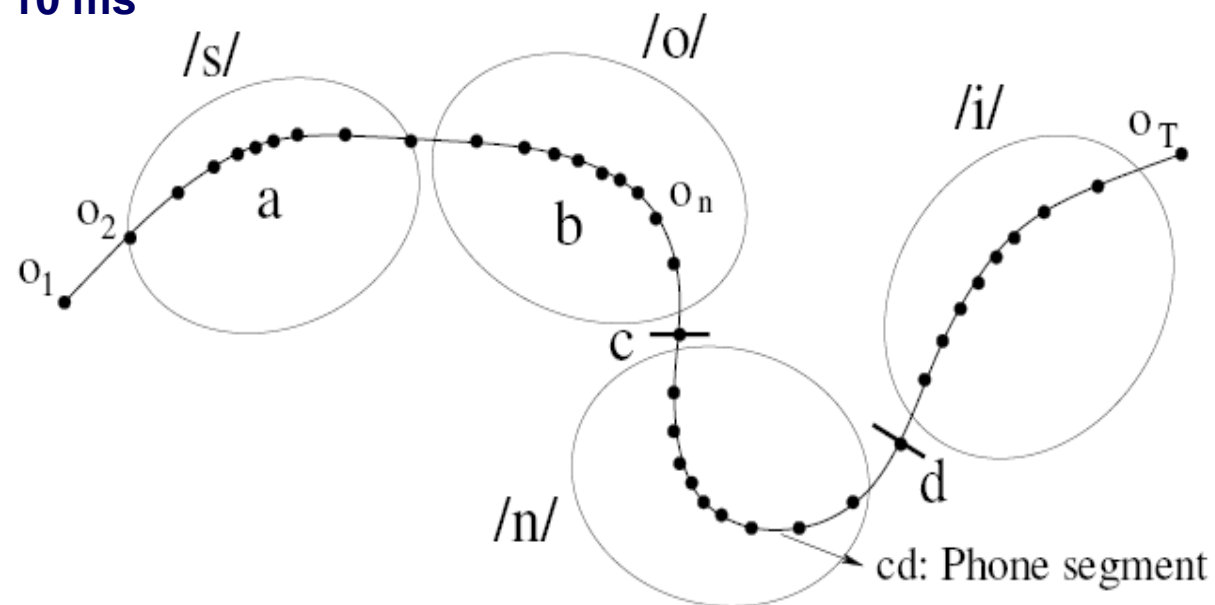**Phone clusters**

/a/
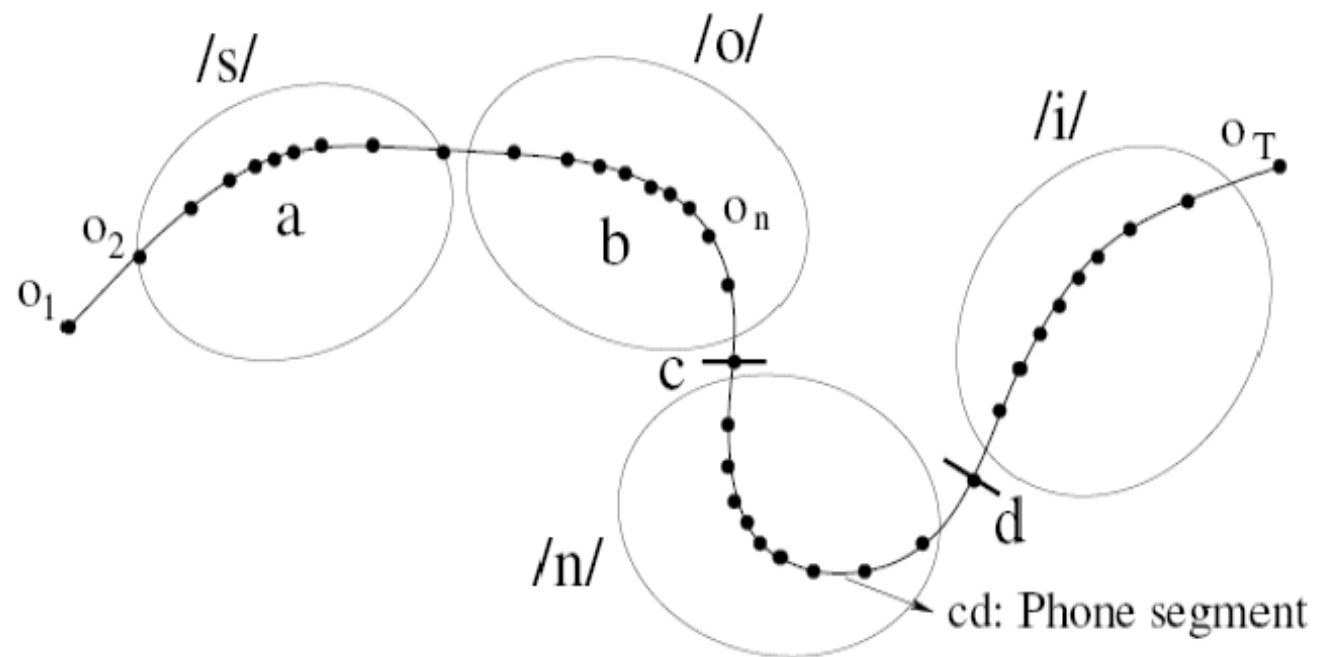
/i/

/o/

/r/

/m/

/sh/

**Speaker A**

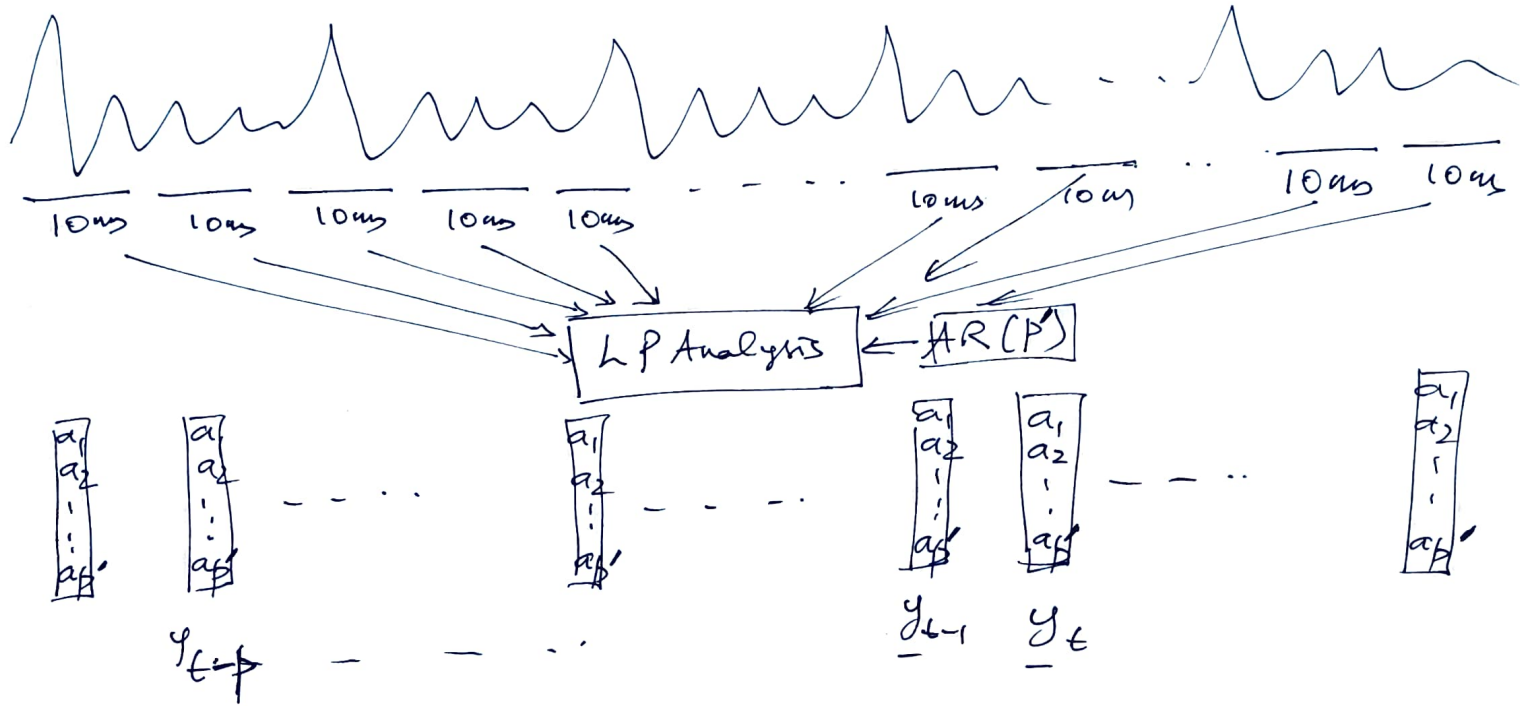**Speaker B**

**One feature vector every 10 ms**

**SPEECH RECOGNITION ALGORITHMS**
- ❑ **TAKE THIS FEATURE VECTOR SEQUENCE**
- ❑ **AS INPUT AND DETERMINE "WHAT HAS BEEN SAID"**
- ❑ **e.g. SEQUENCE OF PHONES / SEQUENCE OF WORDS etc.**

# VAR: Vector Auto regression



$$y_t = c + A_1 \underline{y}_{t-1} + A_2 \underline{y}_{t-2} + \cdots + A_i \underline{y}_{t-i} + \cdots + A_p \underline{y}_{t-p} + \underline{e}_t$$

$p \times 1$

$p \times 1$ intercept Constant

$A_i : p \times p$ time invariant matrix

error vector $p \times 1$

$p=1 \implies \text{VAR}(1)$, 2-d case

$$\underset{\underline{y}_t}{\begin{bmatrix} y_{1,t} \\ y_{2,t} \end{bmatrix}} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + \overset{A_1}{\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}} \overset{\underline{y}_{t-1}}{\begin{bmatrix} y_{1,t-1} \\ y_{2,t-1} \end{bmatrix}} + \begin{bmatrix} e_{1,t} \\ e_{2,t} \end{bmatrix}$$
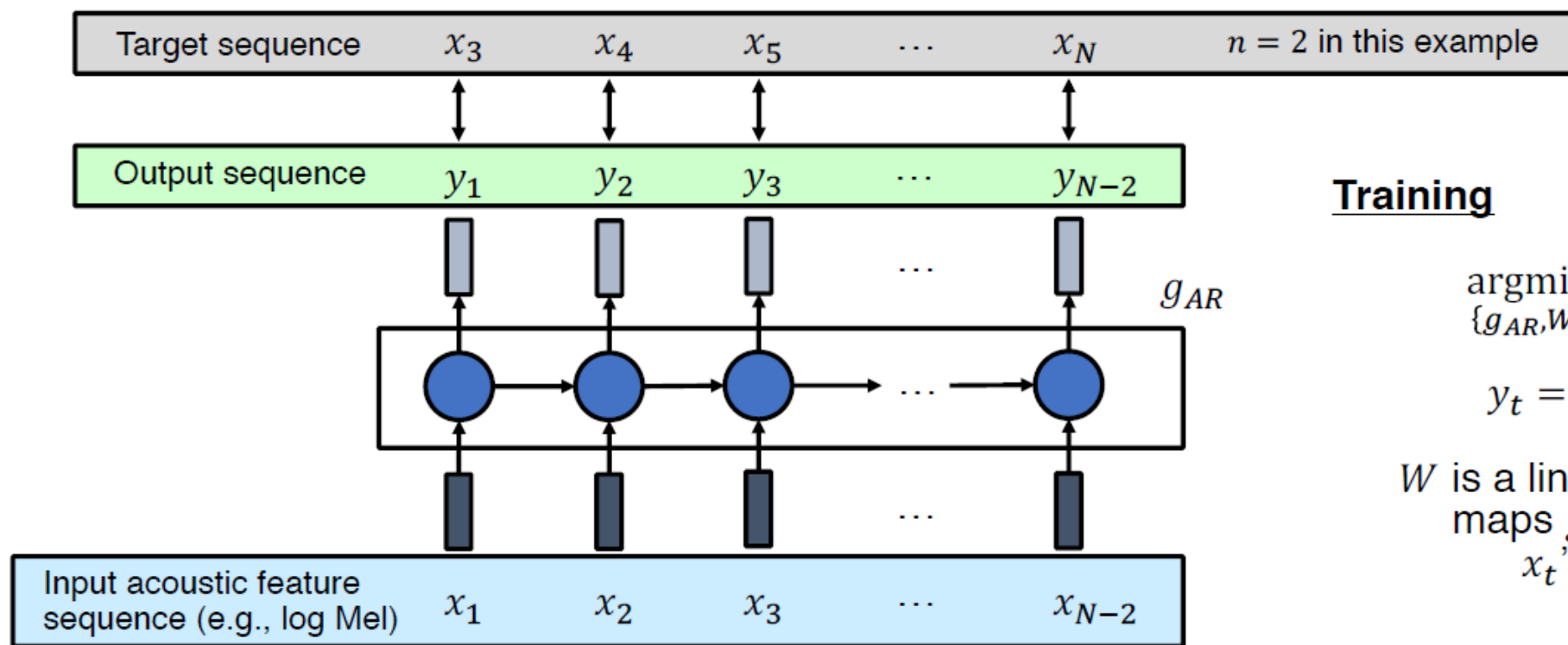
$$\implies \begin{array}{l} y_{1,t} = c_1 + a_{11}\, y_{1,t-1} + a_{1,2}\, y_{2,t-1} + e_{1,t} \\ y_{2,t} = c_2 + a_{21}\, y_{1,t-1} + a_{22}\, y_{2,t-1} + e_{2,t} \end{array}$$



$$t-p \qquad\qquad\qquad\qquad t-1 \quad t$$

$\underline{y}_{t-p}$ — — — $\cdots$ — — — $\underline{y}_{t-1}$ $\underline{y}_t$ — $p \times 1$

$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\text{History of } p\text{-vectors}}$

$\underset{p \times 1}{\phantom{x}} \underset{p \times p}{\left\{ \underline{c}, \underline{A}_i \Big|_{i=1}^{p} \right\}} \longrightarrow$

$\underline{\hat{y}}_t$ — $p \times 1$

$\underline{e}_t$ — $p \times 1$

# Autoregressive Predictive Coding (APC)

- Given a previous context $(x_1, x_2, \ldots, x_t)$, APC tries to predict a future audio feature $x_{t+n}$ that is $n$ steps ahead of $x_t$
  - Uses an autoregressive model $g_{AR}$ to summarize history and produce output
  - $n \geq 1$ encourages $g_{AR}$ to infer more global underlying structures of the data rather than simply exploiting local smoothness of speech signals



**Training**

$$\underset{\{g_{AR}, W\}}{\text{argmin}} \ \sum_{t=1}^{N-n} |x_{t+n} - y_t|,$$

$$y_t = g_{AR}(x_1, \ldots, x_t) \cdot W$$

$W$ is a linear transformation that maps $g_{AR}$'s output back to $x_t$'s dimensionality

# Types of autoregressive model $g_{AR}$

- $g_{AR}$
  - Input: $\mathbf{x} = (x_1, x_2, \ldots, x_N)$
  - Output: $\mathbf{y} = (y_1, y_2, \ldots, y_N)$

- $L$-layer Unidirectional RNN:
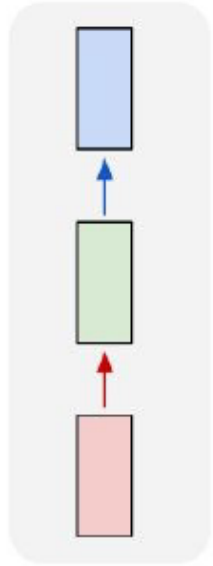
  $$h_0 = \mathbf{x}$$
  $$h_l = \text{RNN}^{(l)}(h_{l-1}), \forall l \in [1, L]$$
  $$\mathbf{y} = h_L \cdot W$$
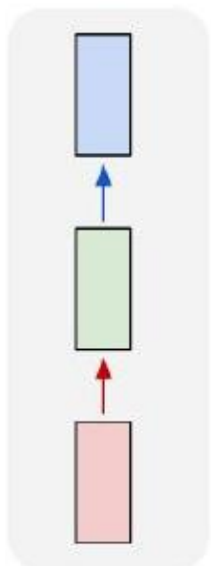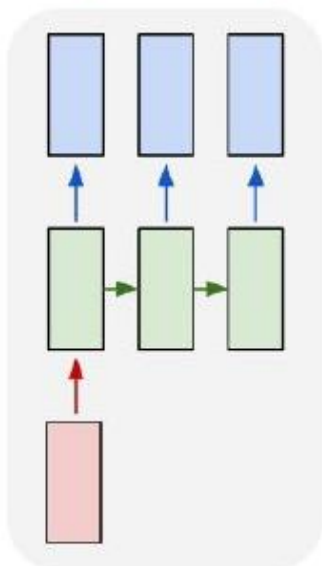


- Feature extraction: $\mathbf{h}_L$

one to one



**Vanilla Neural Networks**

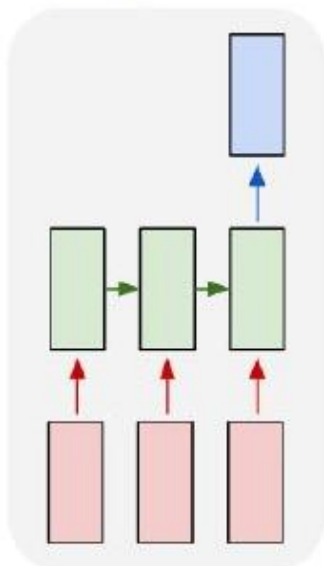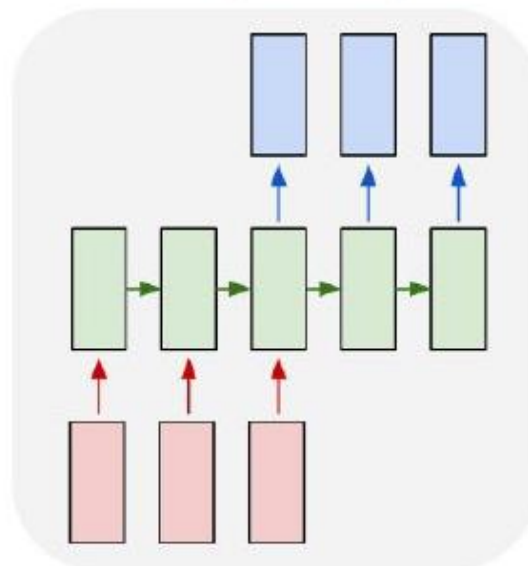# Recurrent Neural Networks: Process Sequences
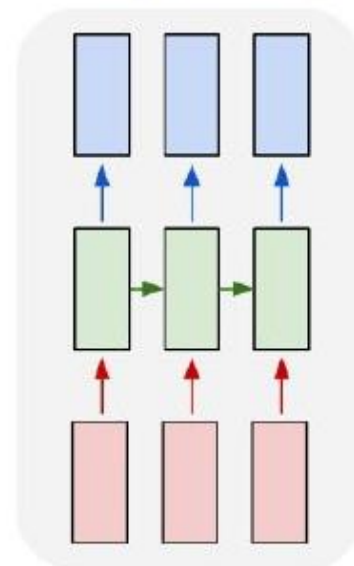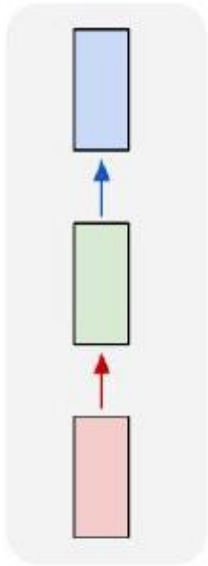


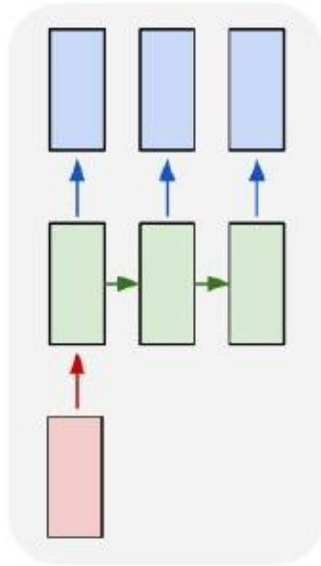one to one　　one to many　　many to one　　many to many　　many to many

e.g. **Image Captioning**
image -> sequence of words
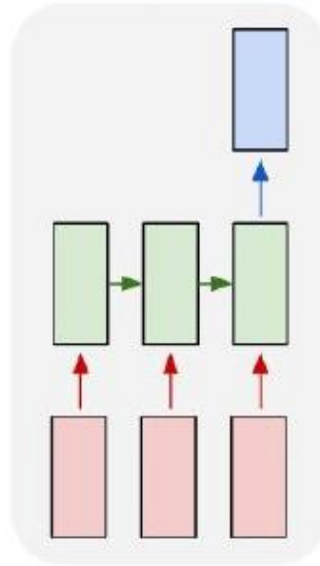
# Recurrent Neural Networks: Process Sequences



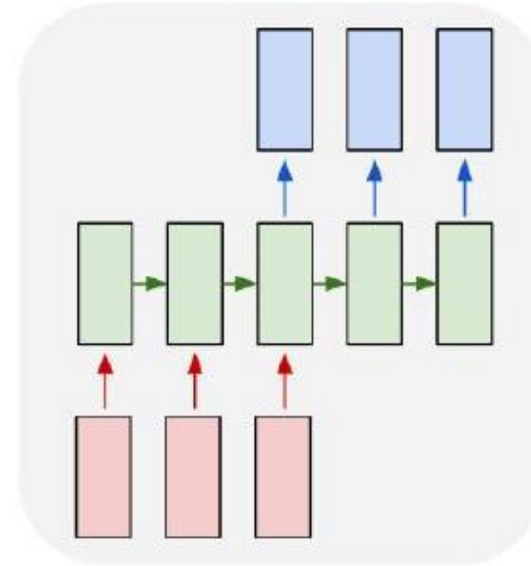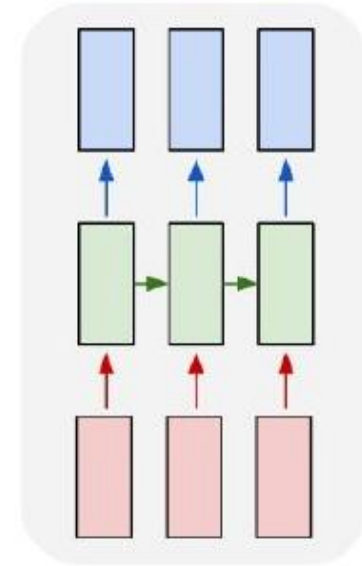one to one        one to many        many to one        many to many        many to many

e.g. **Sentiment Classification**
sequence of words -> sentiment

# Recurrent Neural Networks: Process Sequences

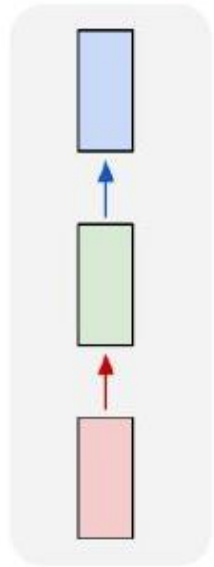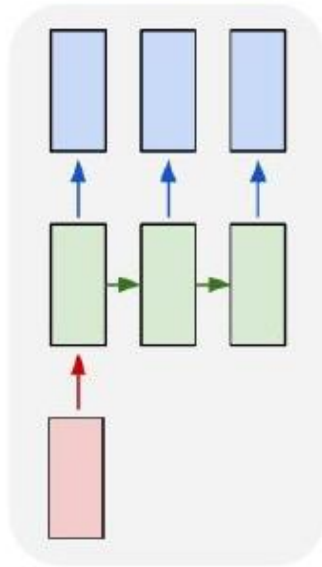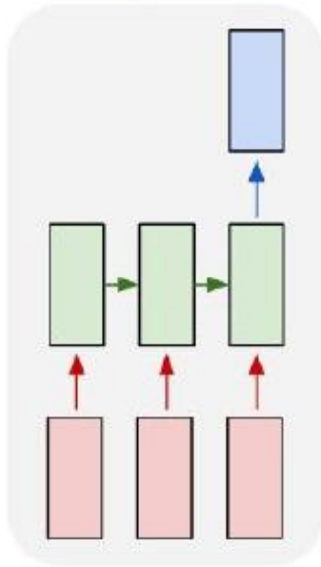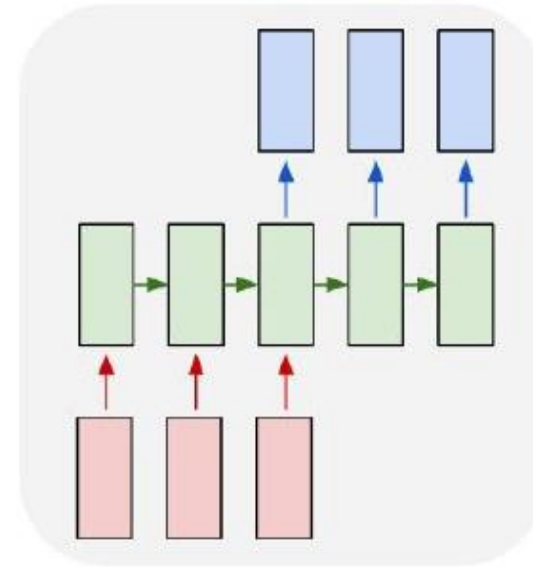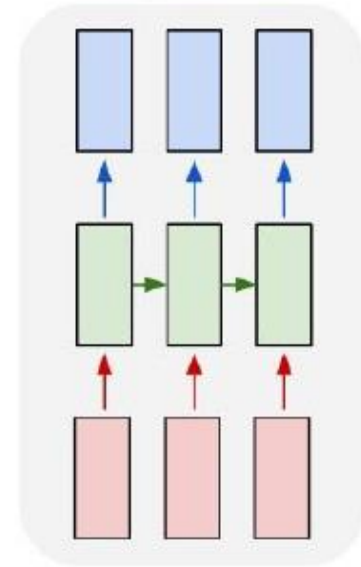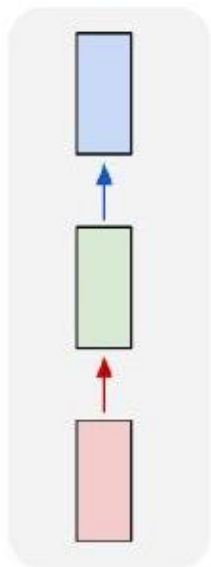one to one    one to many    many to one    many to many    many to many

e.g. **Machine Translation**
seq of words -> seq of words

# Recurrent Neural Networks: Process Sequences



e.g. **Video classification on frame level**

$\hat{\mathbf{y}}_t$

$V$

$W$

$\mathbf{h}_t$

$U$

$\mathbf{x}_t$

$\hat{\mathbf{y}}_t = \phi(V\mathbf{h}_t)$

$\mathbf{h}_t = \psi(U\mathbf{x}_t + W\mathbf{h}_{t-1})$

$\psi$ can be $\tanh$ and $\phi$ can be softmax

## Unrolling the Recurrence

$\hat{\mathbf{y}}_1 \qquad \hat{\mathbf{y}}_2 \qquad \hat{\mathbf{y}}_3 \qquad\qquad \hat{\mathbf{y}}_\tau$

$V \qquad\quad V \qquad\quad V \qquad\qquad\quad V$

$\mathbf{h}_1 \xrightarrow{W} \mathbf{h}_2 \xrightarrow{W} \mathbf{h}_3 \overset{W}{\dashrightarrow} \quad \overset{W}{\dashrightarrow} \mathbf{h}_\tau$

$U \qquad\quad U \qquad\quad U \qquad\qquad\quad U$

$\mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3 \qquad \cdots \qquad \mathbf{x}_\tau$

# Unrolling the Recurrence



$$\mathbf{a}_t = b + W\mathbf{h}_{t-1} + U\mathbf{x}_t$$

$$\mathbf{h}_t = \tanh \mathbf{a}_t$$

$$\mathbf{o}_t = \mathbf{c} + V\mathbf{h}_t$$

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{o}_t)$$

$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = \mathrm{F}(h_t)$$

$$C_t = \mathrm{Loss}(y_t, \mathrm{GT}_t)$$

------- indicates shared weights

We can process a sequence of vectors **x** by
applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state

some function
with parameters W

old state   input vector at
some time step

Notice: the same function and the same set
of parameters are used at every time step.

$$h_t = f_W(h_{t-1}, x_t)$$

$$\downarrow$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

Re-use the same weight matrix at every time-step

RNN: Computational Graph: Many to Many

# Sequence to Sequence: Many-to-one + one-to-many



**One to many**: Produce output sequence from single input vector

**Many to one**: Encode input sequence in a single vector

# Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

# Multilayer RNNs

Figure 4.2: **LSTM memory block with one cell.** The internal state of the cell is maintained with a recurrent connection of fixed weight 1.0. The three gates collect activations from inside and outside the block, and control the cell via multiplicative units (small circles). The input and output gates scale the input and output of the cell while the forget gate scales the internal state. The cell input and output activation functions (g and h) are applied at the indicated places.

$v_c^{(t)} = s_c^{(t)} \cdot o_c^{(t)}$

$o_c^{(t)}$

$\Pi$

$\sigma$

$s_c^{(t)} = g_c^{(t)} \cdot i_c^{(t)} + s_c^{(t-1)}$

$\Pi$

$\sigma$

$i_c^{(t)}$

$\sigma$

$g_c^{(t)}$

Edge to next time step

Edge from previous time step (and current input)

weight fixed at 1

output

recurrent

SRN unit

$g$

input          recurrent

**LSTM block**

block output

output

recurrent

**y**

output gate

**o**

$\sigma$

recurrent

input

$\odot$

$h$

peepholes

recurrent

$+$          $\sigma$     **f**     $\odot$     cell

forget gate

input

$+$     **c**

recurrent

$+$

input

$\odot$

**z**

**i**

$\sigma$

recurrent

input gate

input

$g$

block input

$+$

input          recurrent

**Legend**

——— unweighted connection

——— weighted connection

- - - connection with time-lag

● branching point

⊙ mutliplication

⊕ sum over all inputs

σ gate activation function (always sigmoid)

g input activation function (usually tanh)

h output activation function (usually tanh)

output

output gate

$\times$

peephole

$\Delta t$

forget gate

$\times$  1.0

CEC

$\Delta t$

$\times$

input gate

input

**IN**

tanh

mul.

1.0

mul.

**OUT**

$\mathbf{W}$

$\otimes$

$\otimes$

$\mathbf{y_{in}} \in [0,1]$

$\mathbf{y_{out}} \in [0,1]$

logsig

*input gate*

logsig

*output gate*

$\mathbf{W_{in}}$

$\mathbf{W_{out}}$

# Vanilla RNN vs LSTM

Vanilla RNN cell



LSTM cell

# Gated Recurrent Unit (GRU)

- A dramatic variant of LSTM
  - It combines the forget and input gates into a single update gate
  - It also merges the cell state and hidden state, and makes some other changes
  - The resulting model is simpler than LSTM models
  - Has become increasingly popular



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Autoregressive Predictive Coding (APC)

- Given a previous context $(x_1, x_2, \ldots, x_t)$, APC tries to predict a future audio feature $x_{t+n}$ that is $n$ steps ahead of $x_t$
  - Uses an autoregressive model $g_{AR}$ to summarize history and produce output
  - $n \geq 1$ encourages $g_{AR}$ to infer more global underlying structures of the data rather than simply exploiting local smoothness of speech signals



**Training**

$$\underset{\{g_{AR}, W\}}{\text{argmin}} \ \sum_{t=1}^{N-n} |x_{t+n} - y_t|,$$

$$y_t = g_{AR}(x_1, \ldots, x_t) \cdot W$$

$W$ is a linear transformation that maps $g_{AR}$'s output back to $x_t$'s dimensionality

# Types of autoregressive model $g_{AR}$

- $g_{AR}$
  - Input: $\mathbf{x} = (x_1, x_2, \ldots, x_N)$
  - Output: $\mathbf{y} = (y_1, y_2, \ldots, y_N)$

- $L$-layer Unidirectional RNN:

$$h_0 = x$$
$$h_l = \text{RNN}^{(l)}(h_{l-1}), \forall l \in [1, L]$$
$$y = h_L \cdot W$$



- Feature extraction: $\mathbf{h}_L$

# APC: Autoregressive Predictive Coding

▢ APCs exploit local smoothness of the speech signal
— to predict "__next__" frame.

▢ To encourage APCs to infer more "global structures"
— Rather than local information in the signal.

— ASK the model to predict a frame
"$n$" steps ahead of the current one.

$1 sec$
$= 8000\ Samples$

$25 ms$    $25 ms$    $25 ms$    $25 ms$

$25 ms$    $25 ms$    $25 ms$

$400 Samples$

$x_1$    $x_2$    —    $x_t$ —    $x_{t+n}$    $x_T$

| A P C |
| --- |

$y_1$    $y_2$    —    $y_t$

$L_1$ loss
between $y_t$ & $x_{t+n}$

Given an utterance represented as a sequence of acoustic feature vectors

$$x_1 \quad x_2 - - - \cdot x_f - - - - x_{f+n} - \cdot \cdot x_T$$

RNN/LSTM/GRU — processes $[x_1 \cdots x_t]$ to yield an output prediction $y_t$

$$x_t \in \mathbb{R}^d \quad, \quad y_t \in \mathbb{R}^d$$

History for $y_t$
$$x_1 \quad x_2 - - - - - x_f \mid x_{f+1} - - - x_{f+n} - \cdot \cdot x_T$$

$\boxed{\text{RNN}}$

$L_1$ loss.

$$y_1 \quad y_2 - - - - y_t$$

$$\to L_1 \text{ loss} = \sum_{t=1}^{T-n} |x_{f+n} - y_t|$$

Model is optimised by minimising $L_1$ loss

6|a
example

70 — 120 ms

$x_3$ . . . . . .
$x_2$ ○
$x_1$ ○ 10ms      10ms ○
○
10ms                    10ms ○

**10 ms**

|o|

Trivial RNN
Task.

n-step look-ahead.
. . . ○
○
○
○

$y_t$

$x_t$  ✗ - - - - - ✗ ○
10ms  ○ . . . ○ ○ ○ ○ $x_{t+n}$

Short-time
Spectral feature
[10 — 25 ms]

|o|

○ ○ ○ ○
○
○
○ ○
○ ○ ○ ○ ○ ○

||x||

$L_1$ loss between $y_t$ & $x_{t+n}$

— Within phone Spectral Continuity
  — Articulatory Constraints [slow moving articulators]
  — Coarticulation effects [across adjacent phons]

n - hyperparameter   ☐ n = 1 ⟹ objective ~ LM ~ Too simple a task
                      ☐ n↑ ⟹ More difficult the task.

$X_{1:T}$ : $\boxed{x_1 \; x_2 \; - \; - \; - \; \boxed{x_t}}$ $- \; - \; - \; \cdot \cdot$ $\boxed{x_{t+n}}$ $- \cdot \cdot \; x_T$

$\downarrow f \rightarrow$ RNN/ LSTM/ GRU

$h_T$

$H_{1:T}$ : $\boxed{h_1 \; h_2 \qquad \boxed{h_t}}$ $- \; - \; -$

$\downarrow g \rightarrow$ Linear Projection

[Regression layer that predicts $x_{t+n}$ as $y_{t+n}$ from $h_t$

$\boxed{y_{t+n}}$ $\Longleftarrow$

$L_t = \| y_{t+n} - x_{t+n} \|$

$$H_{1:T} = f(X_{1:T})$$
$$y_{t+n} = g(h_t) = \bar{W} h_t$$
$$L_t = \| y_{t+n} - x_{t+n} \|_1$$

Optimize $f$ & $g$ to minimize $L_t$ over $t = 1, \cdots, T-n$

Table 1: *Comparing APCs with a series of CPC models on phone classification. PERs are reported.*

| Method | #(step) | | | |
|---|---|---|---|---|
| | 2 | 5 | 10 | 20 |
| cpc-n9all | 51.3 | 48.8 | 50.8 | 54.6 |
| cpc-n9same | 47.5 | 48.2 | 50.0 | 53.0 |
| cpc-ctx-n9same | 42.1 | 46.1 | 48.8 | 53.8 |
| cpc-ctx-exhaust | 42.9 | 43.1 | 45.6 | 49.1 |
| apc (proposed) | 36.5 | 35.6 | 35.4 | 37.7 |

Table 2: *PERs on phone classification. All features are fed to a linear classifier unless otherwise stated. The number of steps to the target #(steps) is not relevant in the first four rows.*

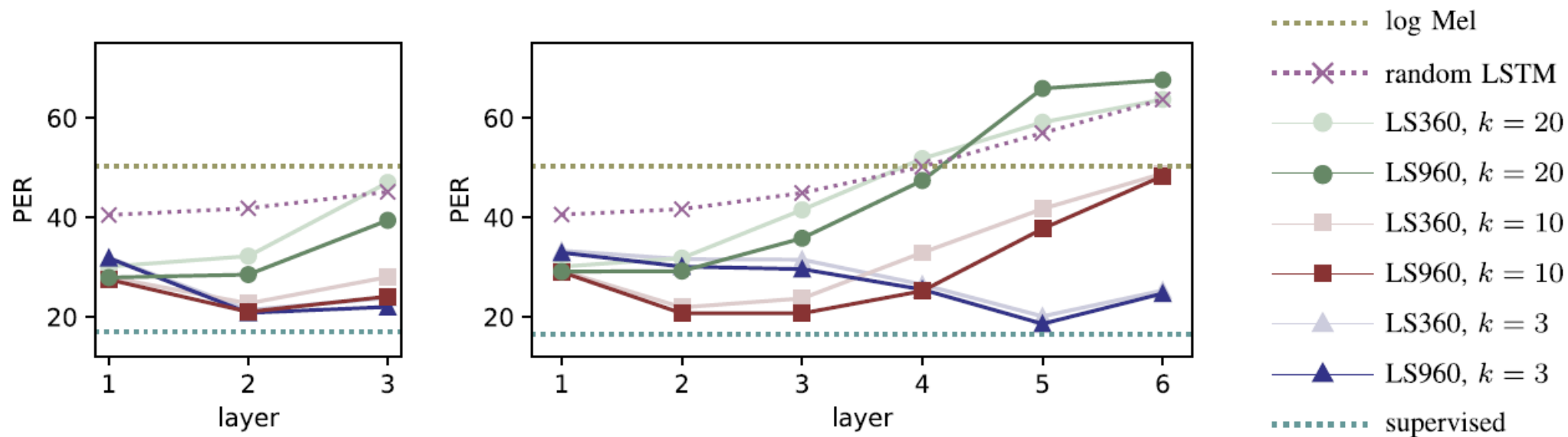| Method | #(step) | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 5 | 10 | 20 |
| Mel | | | | 50.0 | | |
| Mel + MLP-1 | | | | 43.4 | | |
| Mel + MLP-3 | | | | 41.3 | | |
| cpc best | | | | 42.1 | | |
| apc 1-layer | 39.4 | 36.5 | 35.4 | 35.6 | 35.4 | 37.7 |
| apc 2-layer | 38.5 | 34.6 | 35.9 | 35.7 | 34.6 | 38.8 |
| apc 3-layer | 37.2 | 36.7 | 33.5 | 36.1 | 37.1 | 38.8 |
| apc 4-layer | 36.2 | 34.4 | 34.5 | 35.3 | 36.9 | 39.6 |

Fig. 1. Phone error rates (PERs) of frame classification on dev93 with representations produced by 3-layer LSTMs (*left*) and 6-layer LSTMs (*right*). We use LS360 and LS960 to denote the LSTMs trained on the 360-hour subset and the 960 hours combined of LibriSpeech, respectively. We use $k$ to denote the number of time steps into the future in the APC objective.

# Transfer learning experiments

- Setup: pre-training + fine-tuning

- Pre-training data
  - Speech portion of the LibriSpeech 360 hours subset
  - 921 speakers
  - 80-dimensional log Mel spectrograms as input acoustic features (i.e., $x_t \in \mathbb{R}^{80}$)
  - Use extracted features to replace log Mel as new inputs to downstream models

- Considered downstream tasks
  - Speech recognition
  - Speech translation
  - Speaker identification (skipped in this talk, see paper!)

- Comparing methods
  - Contrastive predictive coding (CPC)
  - Problem-agnostic speech encoder (PASE)

# Speech Recognition

- Considered dataset: Wall Street Journal
  - Training: 90% of si284 (~ 72 hours of audio)
  - Validation: 10% of si284
  - Test: dev93

- APC $g_{AR}$
  - RNNs: 4-layer, 512-dim GRUs
  - Transformers: 4-layer, 512-dim Transformer decoder blocks

- Downstream ASR model
  - Seq2seq with attention [Chorowski et al., 2015]
  - Beam search with beam size = 5
  - No language model rescoring

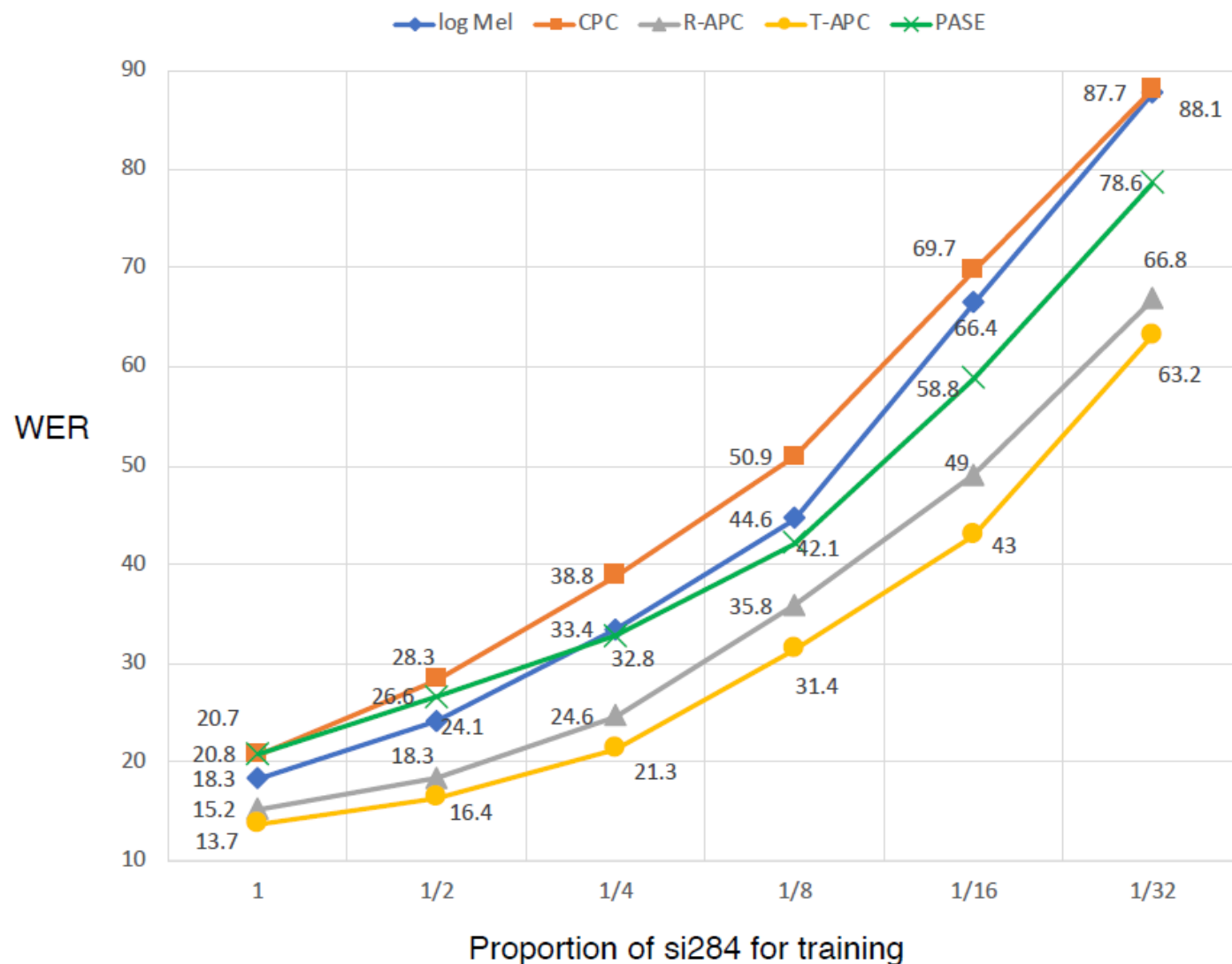# Choice of $n$, and whether to fine-tune $g_{AR}$



**Notations**

- R stands for RNN
- T stands for Transformer
- **Scratch**: $g_{AR}$ randomly initialized and concatenate with ASR model
- **Frozen**: keep $g_{AR}$ frozen when training ASR model
- **Finetuned**: fine-tune $g_{AR}$ along with ASR model

**Findings**

- Sweet spot exists for both Frozen and Finetuned when varying $n$
- Scratch performance is poor, even worse than log Mel baseline
- APC outperforms log Mel most of the time
- For both R and T, Frozen outperforms Finetuned
- Will use R-APC Frozen with $n = 3$ and T-APC Frozen with $n = 5$ for the rest

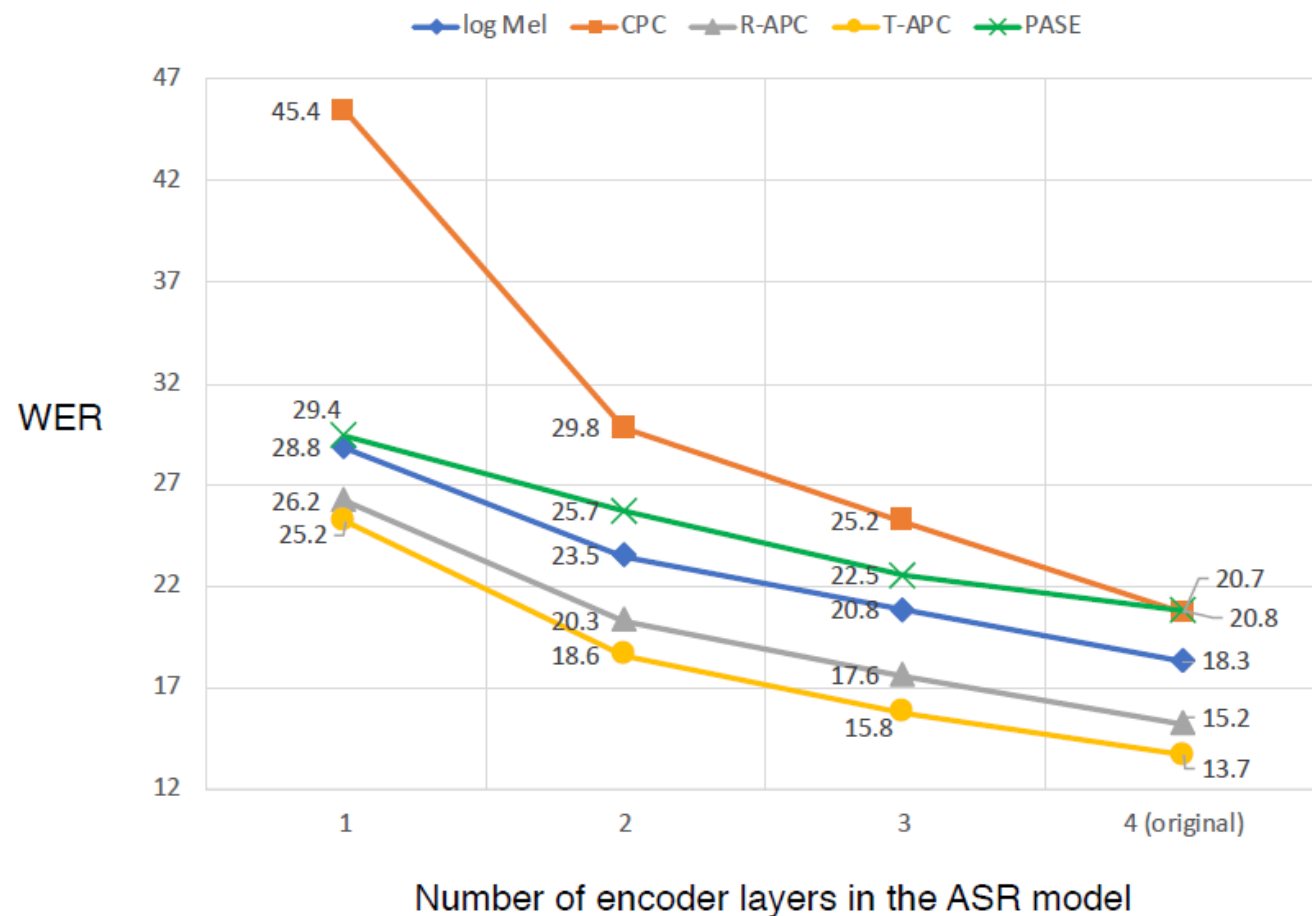# APC for reducing the amount of labeled training data



Recap: all feature extractors were pre-trained with 360 hours of LibriSpeech data; we did not fine-tune any feature extractor with the ASR model

**Findings**

- Full set:
  - 25% and 17% relative improvement for T-APC (13.7) and R-APC (15.2) over log Mel baseline (18.3), respectively

- As we decrease the amount of training data:
  - T-APC (yellow) and R-APC (gray) always outperform other methods
  - Gap between T-APC / R-APC and log Mel (blue) becomes larger
  - Using just half of si284, T-APC (16.4) already outperforms log Mel trained on full set (18.3)

- In the paper we also have the figure where all feature extractors were pre-trained on only 10 hrs of LibriSpeech data. **TLDR**: pre-training still helps even with just 10 hrs of pre-training data

# APC for reducing downstream model size



Note: all models trained on full si284

**Findings**

- T-APC (yellow) and R-APC (gray) always outperform other methods

- T-APC with just 2 layers (18.6) performs similar to log Mel with 4 layers (18.3)

Thank you !!