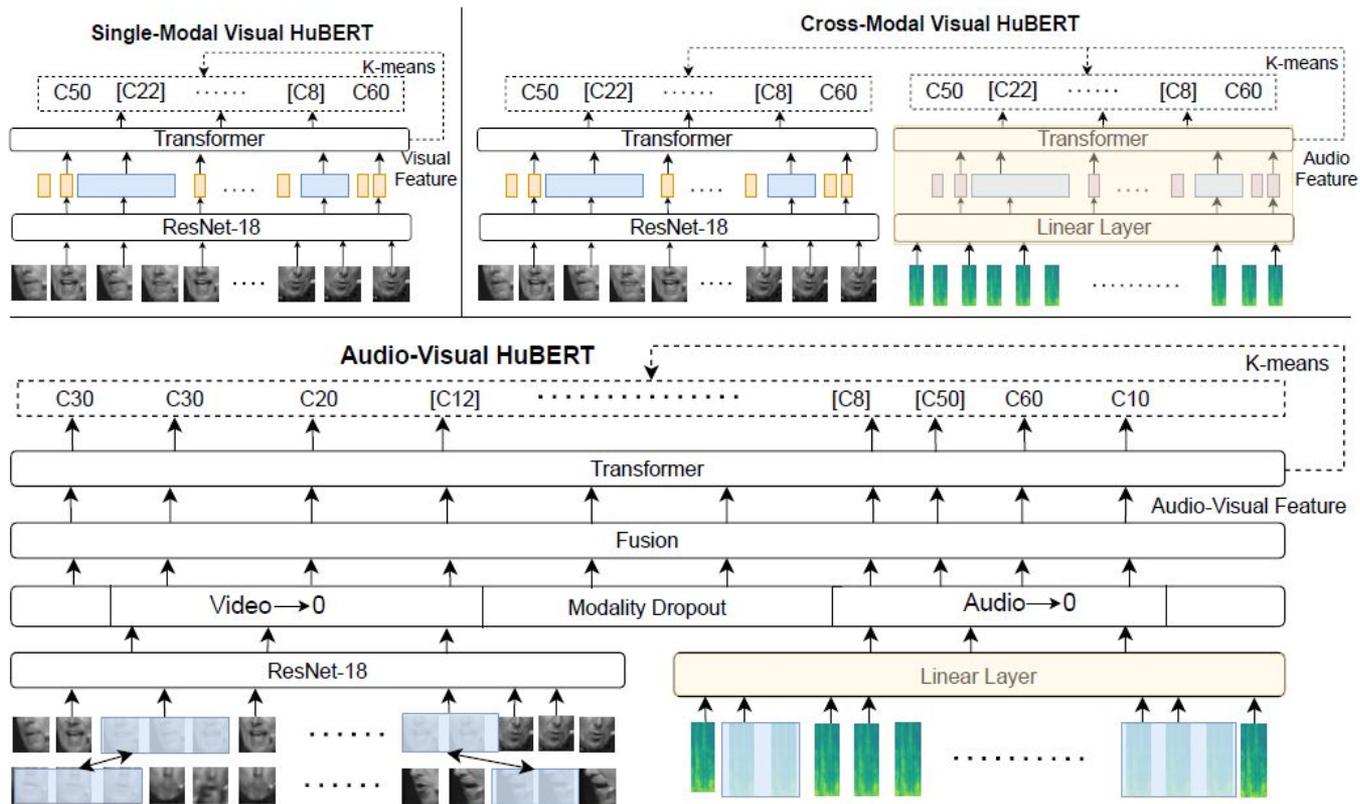# Streaming Audio-Visual Perception

Karthik Ganesan

# Why streaming ?

We need embodied agents to understand in real-time

# AV-HuBERT

# Drawbacks of AV-Hubert

1.  We need bidirectional context , thus we need to wait until entire input is provided
2.  Chunking is also non-trivial as end-point detection also needs supervision to train
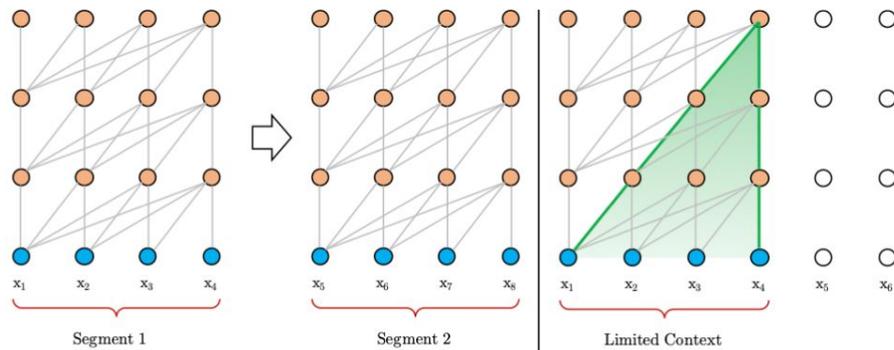
How did we solve this problem in the uni-modal speech recognition ?

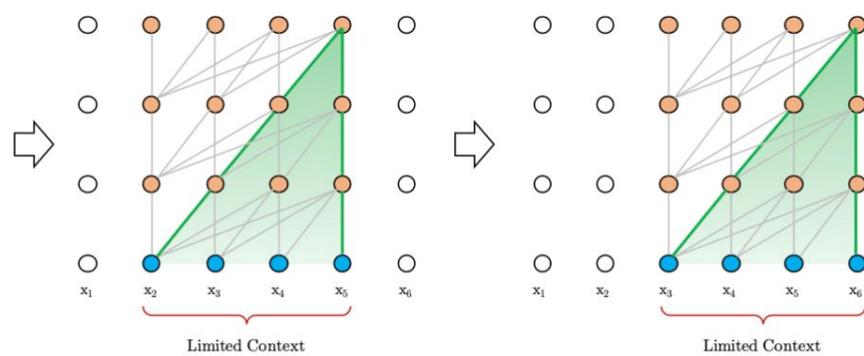# STREAMING TRANSFORMER ASR WITH BLOCKWISE SYNCHRONOUS BEAM SEARCH

*Emiru Tsunoo[1], Yosuke Kashiwagi[1], Shinji Watanabe[2]*

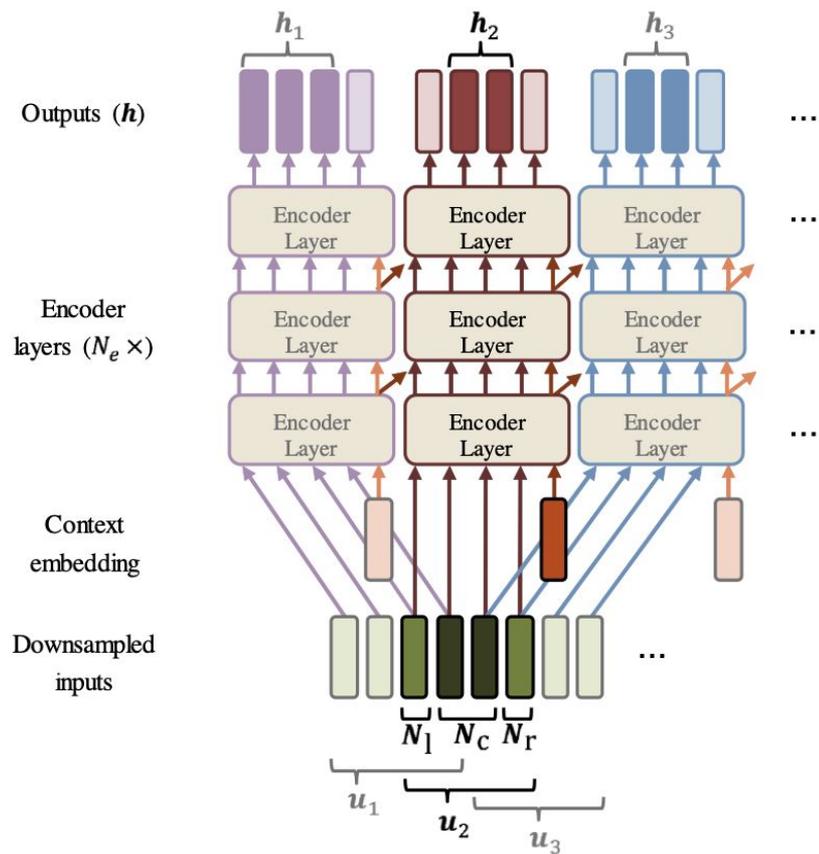[1]Sony Corporation, Japan
[2]Johns Hopkins University, USA
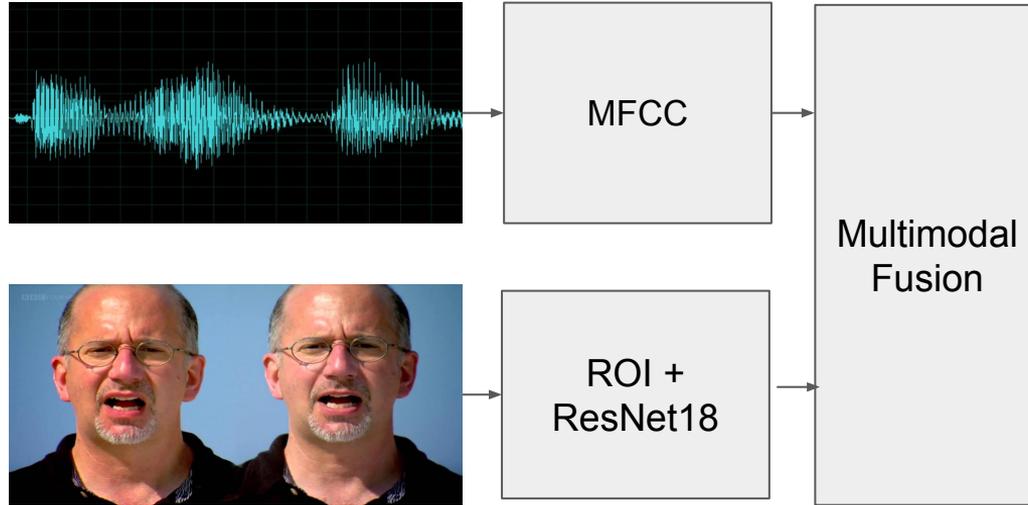
(a) Train phase.

(b) Evaluation phase.

# STREAMING TRANSFORMER ASR WITH BLOCK-WISE SYNCHRONOUS BEAM SEARCH

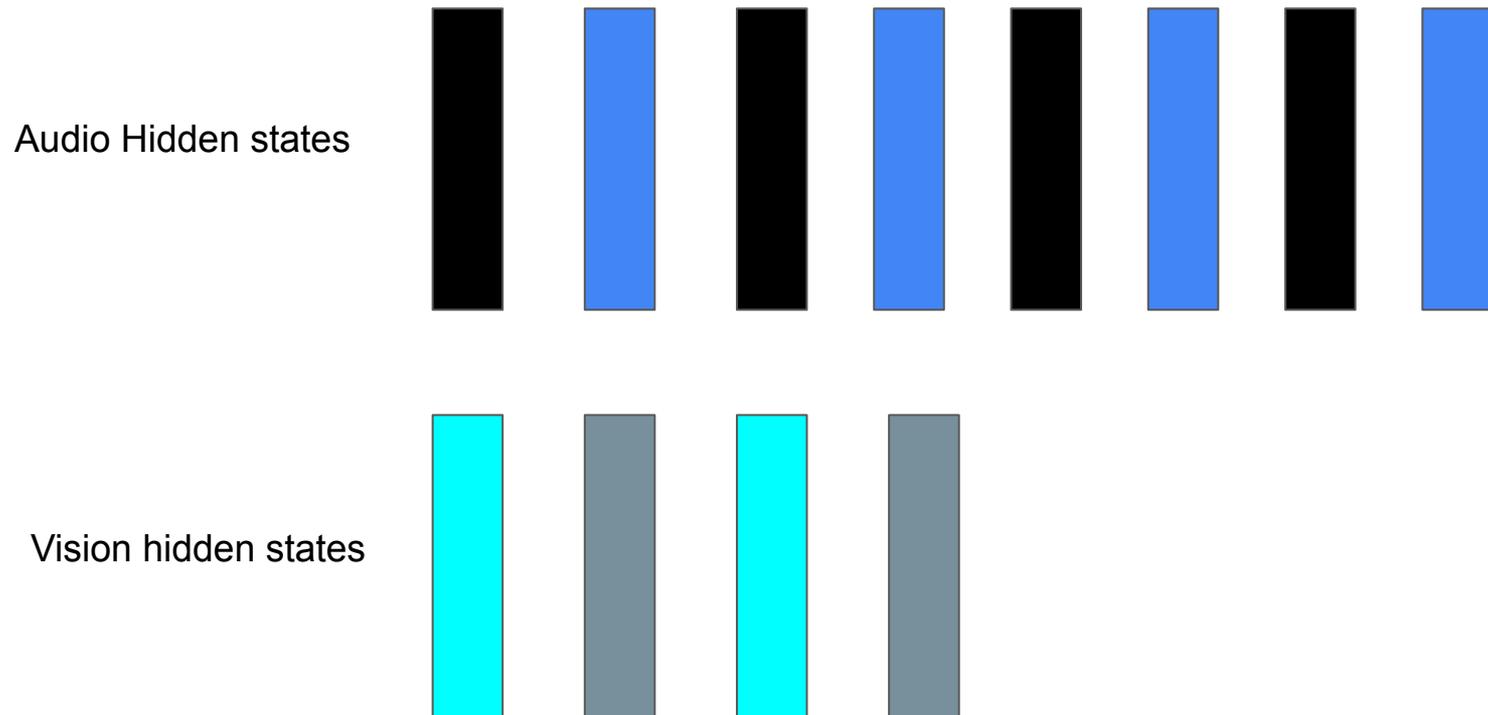How do we make this approach multimodal?

# Proposed Architecture

# Multimodal fusion

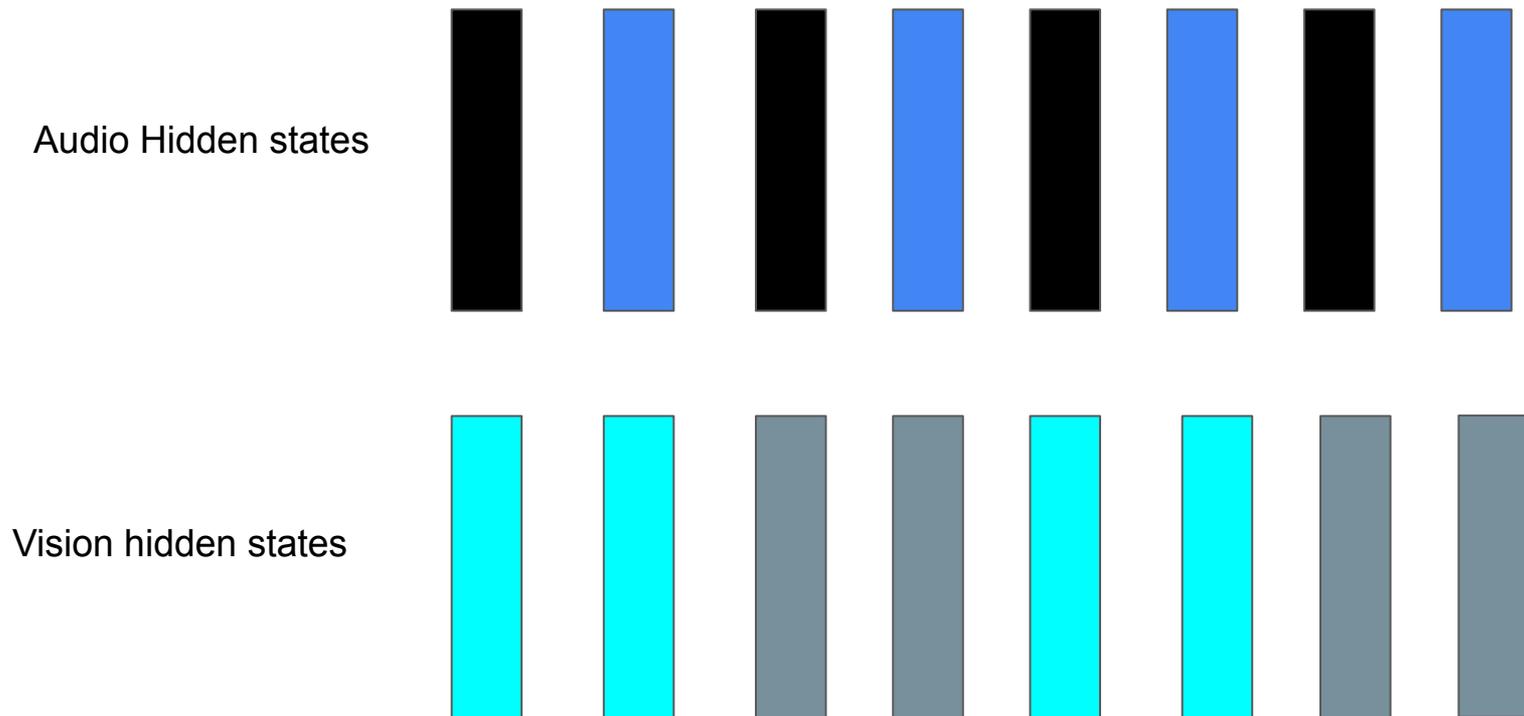Audio Hidden states

Vision hidden states

# TORCH.REPEAT_INTERLEAVE

```
>>> x = torch.tensor([1, 2, 3])

>>> x.repeat_interleave(2)
tensor([1, 1, 2, 2, 3, 3])
```

# Multimodal fusion

Audio Hidden states

Vision hidden states

# Proposed Architecture

# Results

Table 1: Evaluation Results

| Model Name | Model Type | WER | Latency (sec) |
|---|---|---|---|
| Av-HuBERT | Multimodal | 4.10 | 4.823 |
| Av-HuBERT | Speech Only | 4.75 | 4.786 |
| Av-HuBERT | Vision Only | 42.5 | 4.781 |
| Conf-trans | Speech Only | 11.8 | 3.517 |
| Stream | Speech Only | 17.8 | 2.434 |
| Conf-trans-ROI | Multimodal | 10.5 | 4.182 |
| Stream-ROI | Multimodal | 15.2 | 3.106 |

In progress research direction

# DISTILHUBERT: SPEECH REPRESENTATION LEARNING BY LAYER-WISE DISTILLATION OF HIDDEN-UNIT BERT
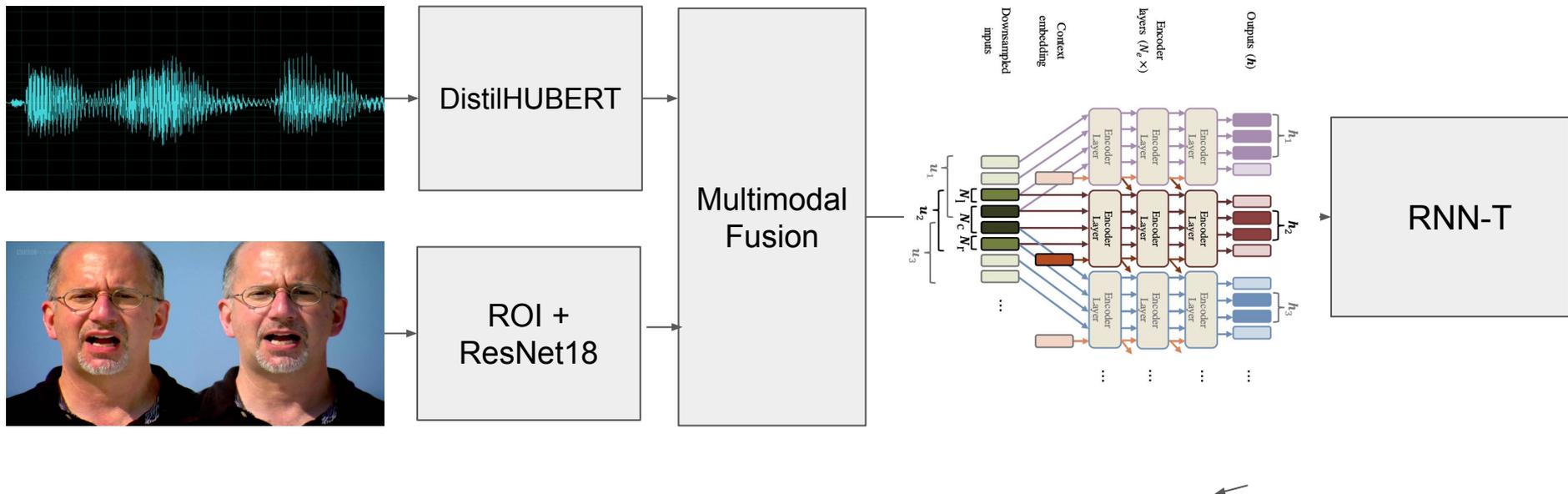
*Heng-Jui Chang, Shu-wen Yang, Hung-yi Lee*

College of Electrical Engineering and Computer Science, National Taiwan University

$$
\mathcal{L}^{(l)} = \mathcal{L}^{(l)}_{\ell 1} + \lambda \mathcal{L}^{(l)}_{\cos}
$$

$$
= \sum_{t=1}^{T} \left[ \frac{1}{D} \left\| \boldsymbol{h}_t^{(l)} - \hat{\boldsymbol{h}}_t^{(l)} \right\|_1 - \lambda \log \sigma \left( \cos \left( \boldsymbol{h}_t^{(l)}, \hat{\boldsymbol{h}}_t^{(l)} \right) \right) \right],
$$

$$(1)$$

(I) Pre-training

# Proposed Architecture

# TRANSFORMER BASED DELIBERATION FOR TWO-PASS SPEECH RECOGNITION

*Ke Hu, Ruoming Pang, Tara N. Sainath, Trevor Strohman*

Google, Inc., USA
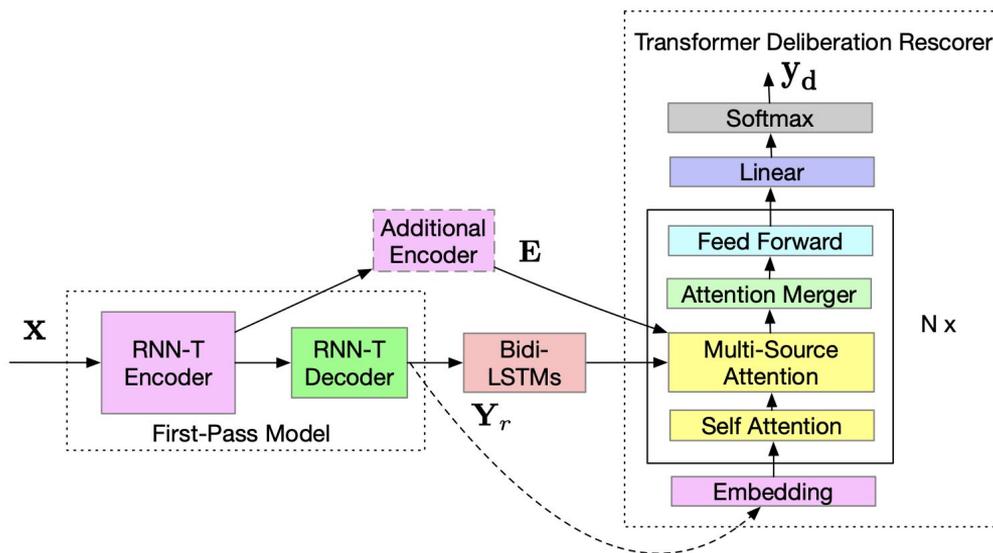
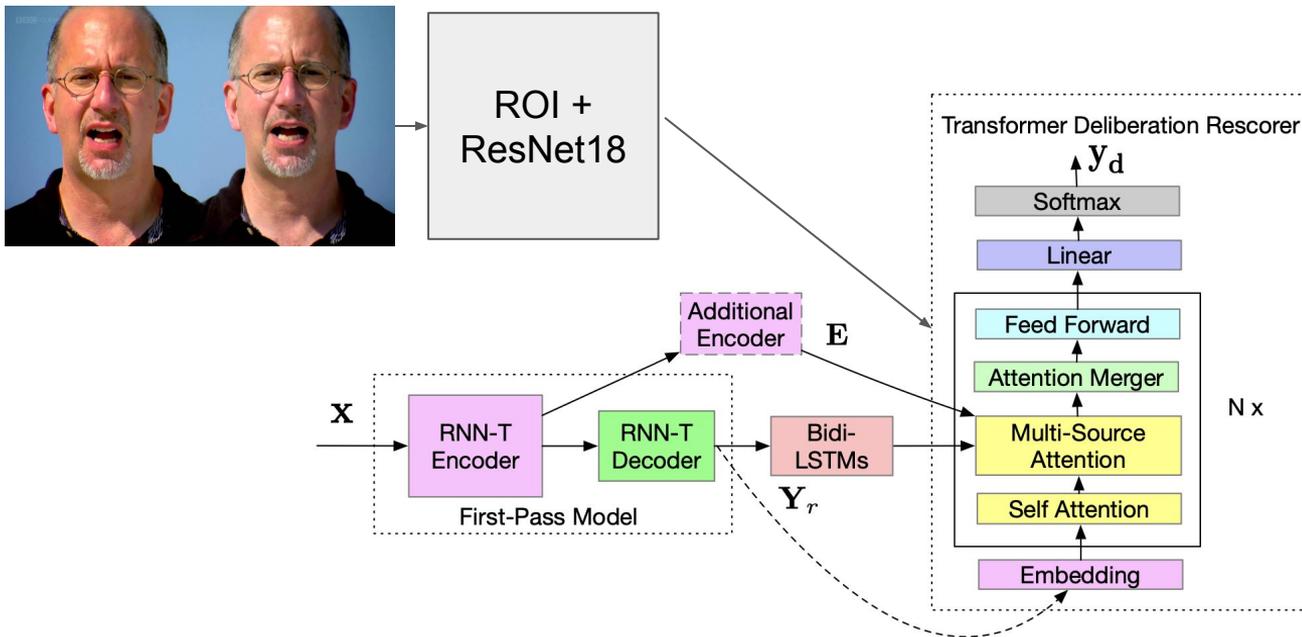{huk,rpang,tsainath,strohman}@google.com

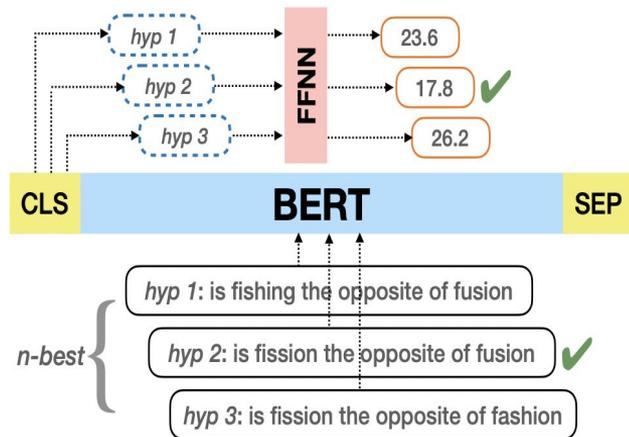# TRANSFORMER BASED DELIBERATION FOR TWO-PASS SPEECH RECOGNITION

*Ke Hu, Ruoming Pang, Tara N. Sainath, Trevor Strohman*

Google, Inc., USA

{huk, rpang, tsainath, strohman}@google.com

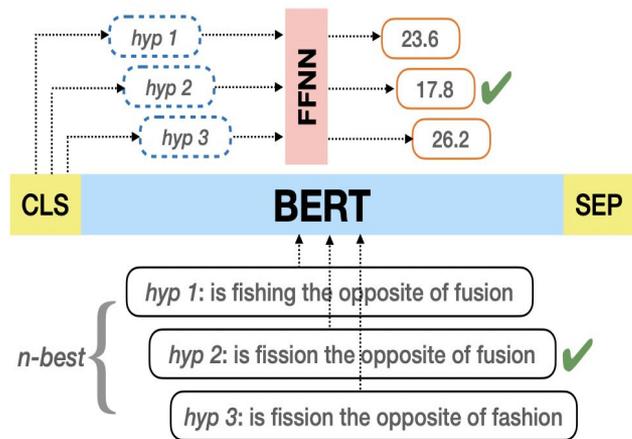# MWER loss based rescoring



$+$ Audio cross attention + some secret sauce

$$P_i = \frac{e^{-s_i}}{\sum_{j=1}^{n} e^{-s_j}}$$

$$\bar{\epsilon}_H = \frac{1}{n} \sum_{i=1}^{n} \epsilon_i$$

$$\mathcal{L}_{\text{MWER}} = \sum_{i=1}^{n} P_i \cdot (\epsilon_i - \bar{\epsilon}_H).$$

hyp 1

hyp 2 ✔

hyp 3

23.6

17.8 ✔

26.2

FFNN

CLS    **BERT**    SEP

*hyp 1*: is fishing the opposite of fusion

n-best { *hyp 2*: is fission the opposite of fusion ✔

*hyp 3*: is fission the opposite of fashion

+

Audio cross attention + some secret sauce

ROI + ResNet18

$$P_i = \frac{e^{-s_i}}{\sum_{j=1}^{n} e^{-s_j}}$$

$$\bar{\epsilon}_H = \frac{1}{n} \sum_{i=1}^{n} \epsilon_i$$

$$\mathcal{L}_{\text{MWER}} = \sum_{i=1}^{n} P_i \cdot (\epsilon_i - \bar{\epsilon}_H).$$

Thank you for your multimodal streaming synchronous attention :)