

# CMU's IWSLT 2022 Dialect Speech Translation System

**Brian Yan<sup>1</sup> Patrick Fernandes<sup>1,2</sup> Siddharth Dalmia<sup>1</sup> Jiatong Shi<sup>1</sup>  
Yifan Peng<sup>3</sup> Dan Berrebbi<sup>1</sup> Xinyi Wang<sup>1</sup> Graham Neubig<sup>1</sup> Shinji Watanabe<sup>1,4</sup>**

<sup>1</sup>Language Technologies Institute, Carnegie Mellon University, USA

<sup>2</sup>Instituto Superior Técnico & LUM LIS (Lisbon ELLIS Unit), Portugal

<sup>3</sup>Electrical and Computer Engineering, Carnegie Mellon University, USA

<sup>4</sup>Human Language Technology Center of Excellence, Johns Hopkins University, USA

{byan, pfernand, sdalmia, jiatongs}@cs.cmu.edu

Day1	May 26, 2022
------	--------------

Time (Dublin)	Session
---------------	---------

<a href="#">15:30-17:30</a>	<a href="#">Poster Session: System Papers</a>
-----------------------------	---

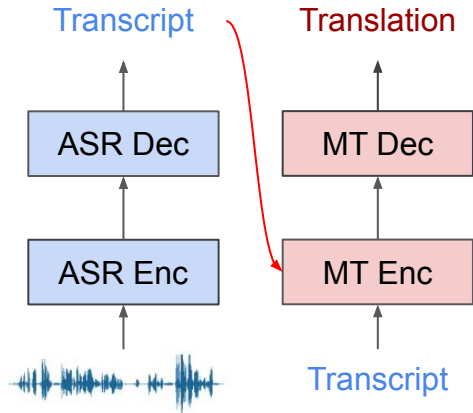


Carnegie Mellon University  
Language Technologies Institute



# The End-to-End Fallacy

## Fully Cascaded



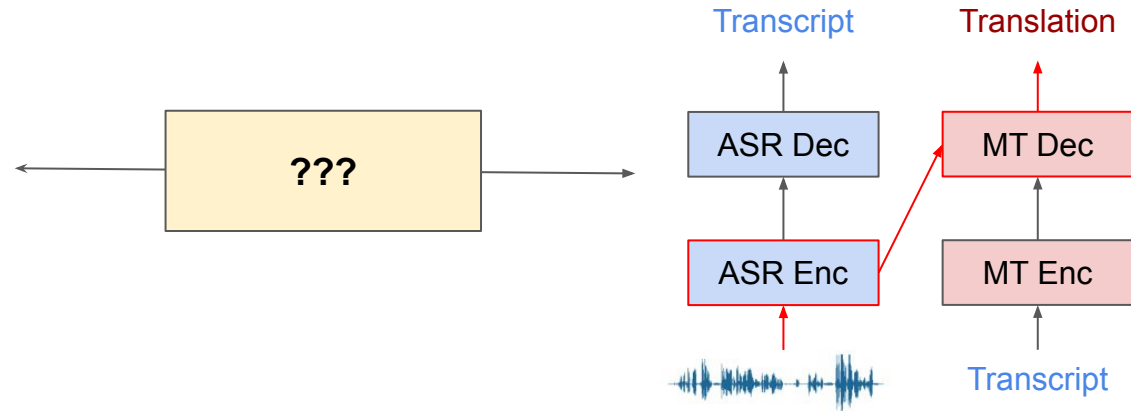
**Con:** early error propagation

**Con:** sensitive to noisy ASR transcription

**Pro:** relatively tons of ASR and MT data

**Pro:** post-processing (e.g. ROVER) / external models for ASR

## Fully Direct (w/ multi-tasking)



**Pro:** no intermediate representation

**Pro:** no intermediate representation

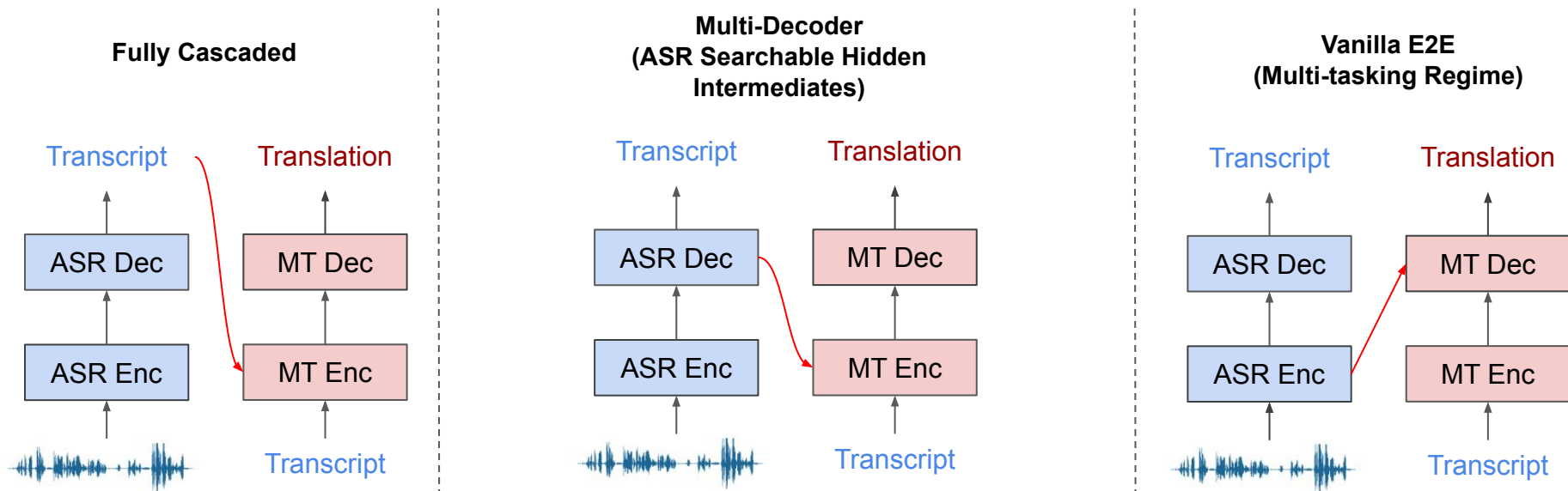
**Con:** less ST data + less data efficient with ASR & MT pre-training multi-task (wasted subnets)

**Con:** single retrieval stage

# Hybrid Approaches to Speech Translation

1. Multi-Decoder with Searchable Hidden Intermediates
2. Minimum Bayes-Risk Decoding

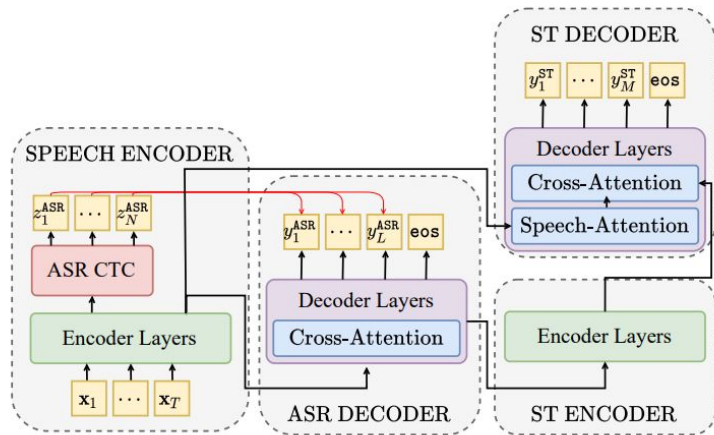
# Multi-Decoder vs. Cascade vs. Vanilla E2E



- ❑ Beam search over ASR output w/ use of external models (e.g. LM, CTC)
- ❑ Use of post-processing to improve ASR output (e.g. ROVER)

# Multi-Decoder with Searchable Hidden Intermediates

$$\mathcal{L} = \lambda_1 \mathcal{L}_{CE}^{ASR} + \lambda_2 \mathcal{L}_{CTC}^{ASR} + \lambda_3 \mathcal{L}_{CE}^{ST}$$



(a) Multi-Decoder

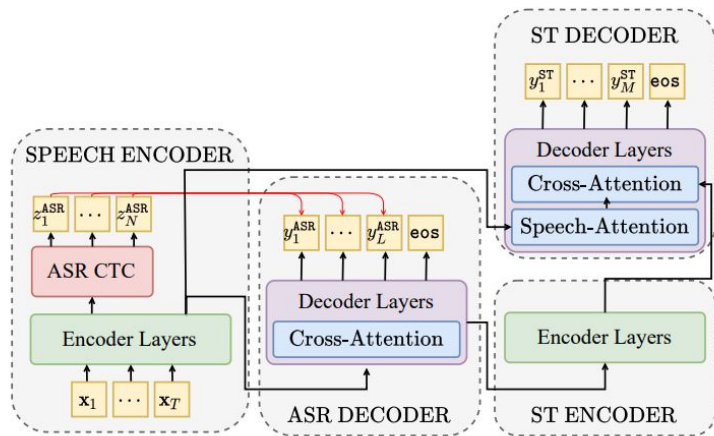
**+Searchable Hidden Intermediates:** ASR decoder representations are retrieved (e.g. via beam search) and passed to the ST Encoder

**Algorithm 1** Beam Search for Hidden Intermediates: We perform beam search to approximate the most likely sequence for the sub-task  $\mathcal{A} \rightarrow \mathcal{B}$ ,  $\mathbf{y}_{BEAM}^B$ , while collecting the corresponding  $DECODER_B$  hidden representations,  $\mathbf{h}_{BEAM}^{D_B}$ . The output  $\mathbf{h}_{BEAM}^{D_B}$  is passed to the final sub-network to predict final output  $\mathcal{C}$  and  $\mathbf{y}_{BEAM}^B$  is used for monitoring performance on predicting  $\mathcal{B}$ .

- 1: **Initialize:** BEAM  $\leftarrow$  {sos}; k  $\leftarrow$  beam size;
- 2:  $\mathbf{h}^{E_A} \leftarrow ENCODER_A(\mathbf{x})$
- 3: **for**  $l=1$  **to**  $max_{STEPS}$  **do**
- 4:   **for**  $\mathbf{y}_{l-1}^B \in BEAM$  **do**
- 5:      $\mathbf{h}_l^{D_B} \leftarrow DECODER_B(\mathbf{h}^{E_A}, \mathbf{y}_{l-1}^B)$
- 6:     **for**  $\mathbf{y}_l^B \in \mathbf{y}_{l-1}^B + \{\mathcal{V}\}$  **do**
- 7:        $s_l \leftarrow P_{\mathcal{A} \rightarrow \mathcal{B}}(\mathbf{y}_l^B | \mathbf{x})^{1-\lambda} P_{EXT}(\mathbf{y}_l^B)^\lambda$
- 8:        $\mathcal{H} \leftarrow (s_l, \mathbf{y}_l^B, \mathbf{h}_l^{D_B})$
- 9:     **end for**
- 10:   **end for**
- 11: BEAM  $\leftarrow \arg^{kmax}(\mathcal{H})$
- 12: **end for**
- 13:  $(s^B, \mathbf{y}_{BEAM}^B, \mathbf{h}_{BEAM}^{D_B}) \leftarrow \argmax(BEAM)$
- 14: **Return**  $\mathbf{y}_{BEAM}^B \rightarrow SUB_{\mathcal{A} \rightarrow \mathcal{B}}NET$  Monitoring
- 15: **Return**  $\mathbf{h}_{BEAM}^{D_B} \rightarrow$  Final  $SUB_{\mathcal{B} \rightarrow \mathcal{C}}NET$

# Multi-Decoder with Searchable Hidden Intermediates

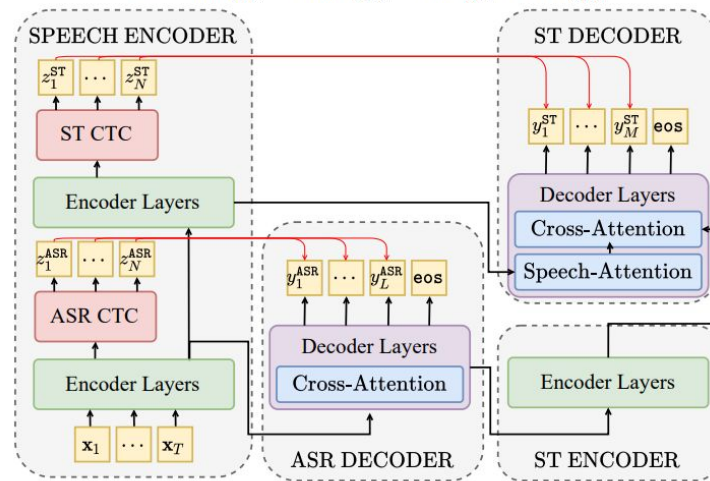
$$\mathcal{L} = \lambda_1 \mathcal{L}_{CE}^{ASR} + \lambda_2 \mathcal{L}_{CTC}^{ASR} + \lambda_3 \mathcal{L}_{CE}^{ST}$$



(a) Multi-Decoder

**+Searchable Hidden Intermediates:** ASR decoder representations are retrieved (e.g. via beam search) and passed to the ST Encoder

$$\mathcal{L} = \lambda_1 \mathcal{L}_{CE}^{ASR} + \lambda_2 \mathcal{L}_{CTC}^{ASR} + \lambda_3 \mathcal{L}_{CE}^{ST} + \lambda_4 \mathcal{L}_{CTC}^{ST}$$



(b) Multi-Decoder w/ Hierarchical Encoder + CTC/Attn ST Decoding

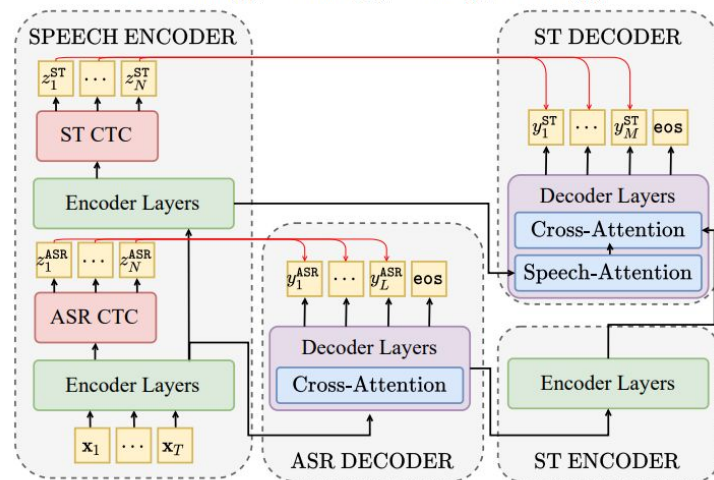
**+Hierarchical Encoder:** re-orders speech encoder  
**+Joint CTC/Attn ST Decoding:** length normalization  
**+ASR CTC Sampling:** simulates ASR errors in training

# Multi-Decoder with Searchable Hidden Intermediates

	test1
Model Name	BLEU(↑)
Encoder-Decoder	16.0
Multi-Decoder	17.1
+ ASR CTC Sampling	17.6
+ Hierarchical Encoder	17.9
+ Joint CTC/Attn ST Decoding (D4)	18.2
+ ASR CTC Sampling	<b>18.4</b>

+2.4 BLEU

$$\mathcal{L} = \lambda_1 \mathcal{L}_{CE}^{ASR} + \lambda_2 \mathcal{L}_{CTC}^{ASR} + \lambda_3 \mathcal{L}_{CE}^{ST} + \lambda_4 \mathcal{L}_{CTC}^{ST}$$



(b) Multi-Decoder w/ Hierarchical Encoder + CTC/Attn ST Decoding

- +Hierarchical Encoder:** re-orders speech encoder
- +Joint CTC/Attn ST Decoding:** length normalization
- +ASR CTC Sampling:** simulates ASR errors in training

# Guiding Multi-Decoder Representations

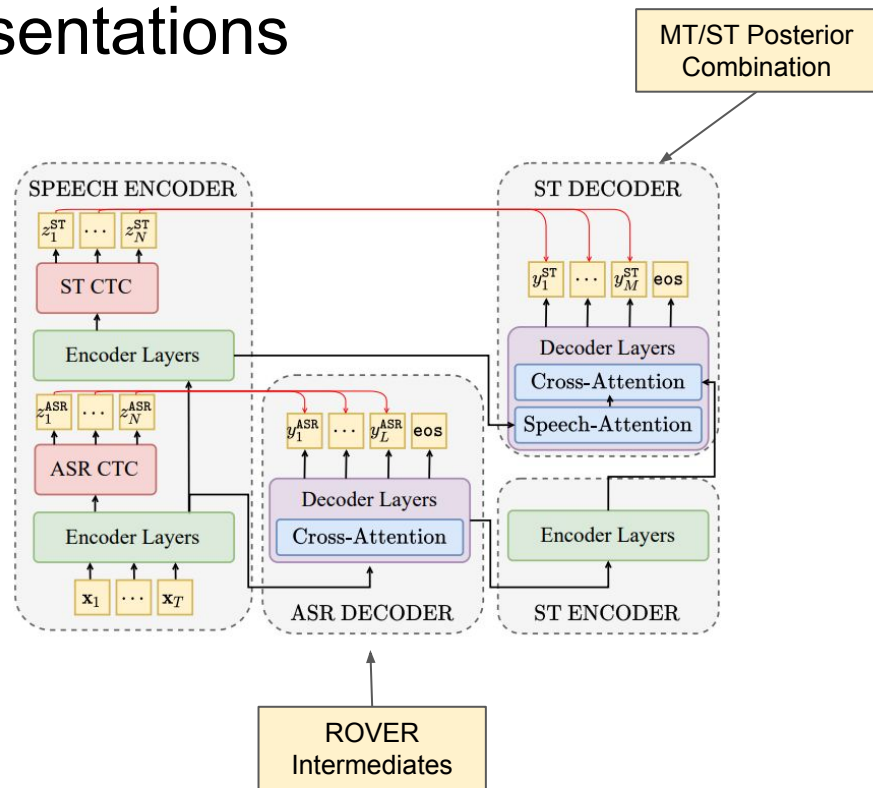
## ASR Decoder:

- We retrieve the hidden representations for ASR outputs generated by ROVER combination
- No intermediate beam search is required

## ST Decoder:

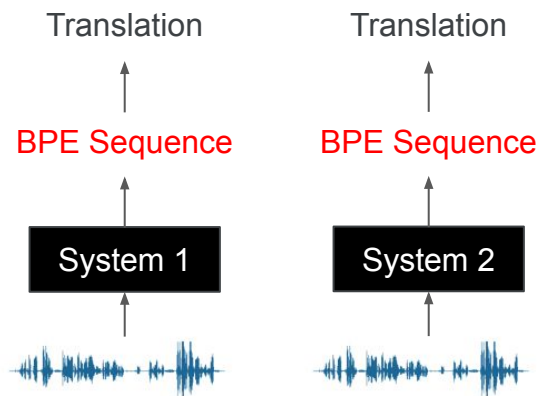
- External MT/ST models are used for posterior combination, along with joint CTC/Attn. ST decoding

\*we do not consider these Multi-Decoder variants to be purely E2E



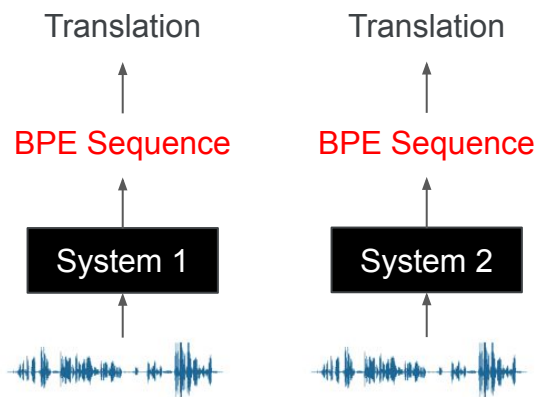


# Minimum Bayes-Risk Decoding



- Posterior comb. distinct BPE vocabularies?
- ROVER-like align-then-vote - fit for translation?

# Minimum Bayes-Risk Decoding



- Posterior comb. distinct BPE vocabularies?
- ROVER-like align-then-vote - fit for translation?

$$\hat{y}_{\text{MBR}} = \arg \max_{y \in \bar{\mathcal{Y}}_{\text{cands}}} \underbrace{\mathbb{E}_{Y \sim p_{\theta}(y|x)} [u(Y, y)]}_{\text{BLEU metric}},$$

Candidate set

$\approx \frac{1}{M} \sum_{j=1}^M u(y^{(j)}, y)$

BLEU metric      Sample set

- Both systems generate candidates and samples
- Risk of each candidate is measured as the avg. BLEU against all of the samples as references
- Systems can be black boxes

# Minimum Bayes-Risk Decoding

ID	Type	Model Name	Child	Dialect	test1	test2
			System(s)	Transfer	BLEU(↑)	BLEU(↑)
C1	Cascade	ASR Mixing Cascade	A1, B1	✗	16.4	-
C2	Cascade	+ ASR Rover Comb.	A2, B1	✗	16.7	-
C3	Cascade	+ MT Posterior Comb.	A2, B2	✗	17.5	18.6
C4	Cascade	ASR Mixing Cascade	A3, B3	✓	17.3	-
C5	Cascade	+ ASR Rover Comb.	A4, B3	✓	17.4	-
C6	Cascade	+ MT Posterior Comb.	A4, B4	✓	<b>17.9</b>	<b>19.4</b>
D1	E2E ST	Hybrid Multi-Decoder	-	✗	17.7	-
D2	Mix	+ ROVER Intermediates	A2	✗	18.1	19.1
D3	Mix	+ ST/MT Posterior Comb.	A2, B5	✗	18.7	19.7
D4	E2E ST	Hybrid Multi-Decoder	-	✓	18.2	-
D5	Mix	+ ROVER Intermediates	A4	✓	18.3	19.5
D6	Mix	+ ST/MT Posterior Comb.	A4, B5	✓	<b>18.9</b>	<b>19.8</b>
E1	Mix	Min. Bayes-Risk Ensemble	C3, D3	✗	19.2	20.4
E2	Mix	Min. Bayes-Risk Ensemble	C6, D6	✓	<b>19.5</b>	<b>20.8</b>

+1.3 BLEU

+0.6/+1.0 BLEU

# Thanks!

## **CMU's IWSLT 2022 Dialect Speech Translation System**

**Brian Yan<sup>1</sup> Patrick Fernandes<sup>1,2</sup> Siddharth Dalmia<sup>1</sup> Jiatong Shi<sup>1</sup>  
Yifan Peng<sup>3</sup> Dan Berrebbi<sup>1</sup> Xinyi Wang<sup>1</sup> Graham Neubig<sup>1</sup> Shinji Watanabe<sup>1,4</sup>**

<sup>1</sup>Language Technologies Institute, Carnegie Mellon University, USA

<sup>2</sup>Instituto Superior Técnico & LUM LIS (Lisbon ELLIS Unit), Portugal

<sup>3</sup>Electrical and Computer Engineering, Carnegie Mellon University, USA

<sup>4</sup>Human Language Technology Center of Excellence, Johns Hopkins University, USA

{byan, pfernand, sdalmia, jiatongs}@cs.cmu.edu