# Streaming Automatic Speech Recognition with the Transformer Model

Niko Moritz, Takaaki Hori, Jonathan Le Roux

ICASSP
May, 2020

# Motivation

- End-to-end automatic speech recognition (ASR) has greatly simplified the pipeline for buidling and applying ASR systems.

- Offline end-to-end ASR systems have shown to surpass the performance of traditional hybrid DNN-HMM solutions.

- Streaming end-to-end architectures are still lacking behind this success.

- Encoder-decoder based architectures have demonstrated to achieve the best end-to-end ASR results but are difficult to apply in a streaming fashion.
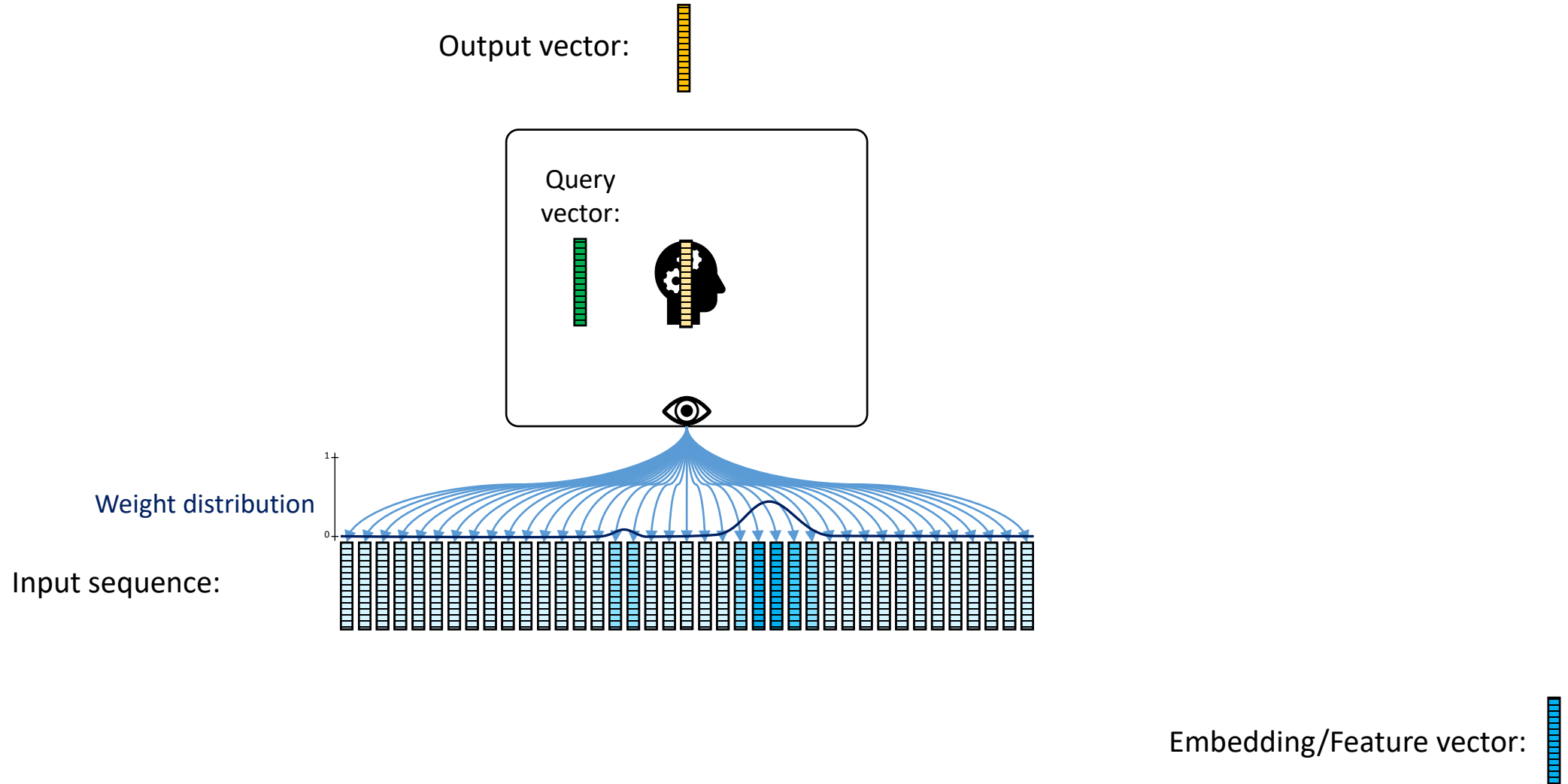
# This work

- Our proposed triggered attention (TA) concept is used to overcome these difficulties.

- The TA concept is applied to the transformer architecture, achieving SOTA streaming end-to-end ASR results.
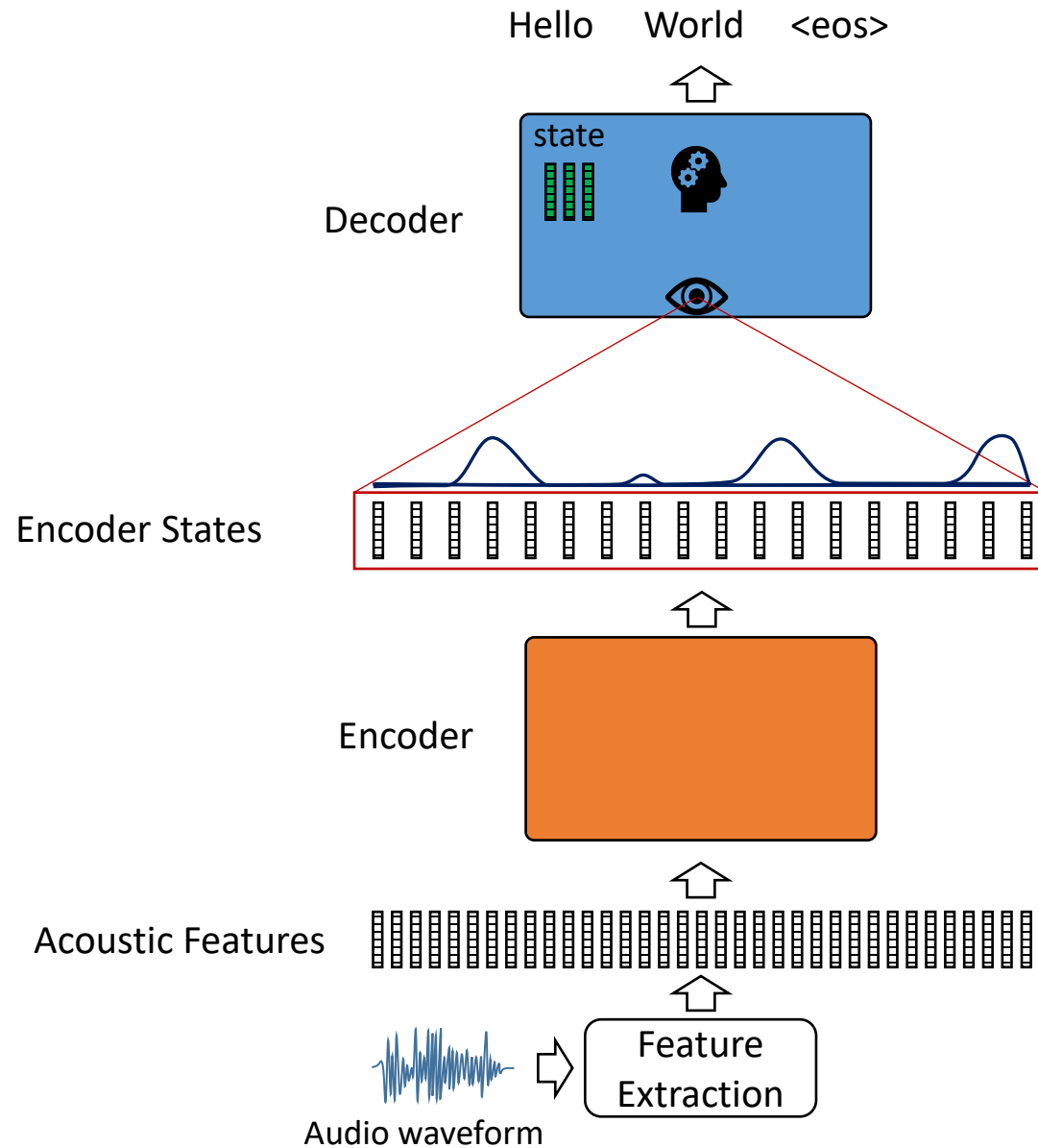
# Outline

- Encoder-Decoder Neural Networks
  - Attention
  - Transformer
  - Self-attention
  - Time-Restricted Self-Attention
  - Streaming Encoder-Decoder Attention (prior work)
- Triggered Attention
  - Architecture
  - Frame-Synchronous Decoding Algorithm
- LibriSpeech Results

Output vector:

Query vector:

Weight distribution

1

0

Input sequence:

Embedding/Feature vector:

# Encoder-Decoder Attention

Hello    World    <eos>

state

Decoder

Encoder States

Encoder

Acoustic Features

Feature Extraction

Audio waveform

# Transformer Architecture

Output sequence:

Query
vector:

Input sequence:

current frame

Input sequence:



past frames

current frame

future frames

time

# Time-Restricted Self-Attention

Output sequence:

Algorithmic delay: $\#layers \cdot \varepsilon^{\text{enc}} = 4$ frames

Output sequence: Layer 2

look-ahead $\varepsilon^{\text{enc}} = 2$ frames

current frame

Input sequence: Layer 1

look-ahead $\varepsilon^{\text{enc}} = 2$ frames

current frame

# Streaming Encoder-Decoder Attention (prior work)

## Adaptive Chunking based on Selection Probability



Example:
- Monotonic Chunkwise Attention (MoChA) [1]

Problems:
- Backpropagation with discrete decisions is not possible.
- No frame-synchronous decoding algorithm.
- Detecting word or word-piece positions is a good part of the ASR job that defines insertion and deletion errors.

[1] C. Chiu and C. Raffel, "Monotonic chunkwise attention," in Proc. ICLR, Apr. 2018.

# Triggered Attention (TA) Architecture

Decoding output: $Y = (\boldsymbol{y}_1, \ldots, \boldsymbol{y}_l) \rightarrow$

$\boxed{\boldsymbol{y}_1}\ \boxed{\boldsymbol{y}_2}\ \boxed{\boldsymbol{y}_3}\ \ldots\ \boxed{\boldsymbol{y}_l}$

shallow fusion → **Joint Decoding** $\quad \boldsymbol{y}_{1:l-1}$

**Language Model**    **CTC**    **Triggered Attention Decoder**

Transformer model
#encoder layers $E = 12$
#decoder layers $D = 6$
attention dimension $= 512$
#attention heads $= 8$

Encoder output: $X_E = (\boldsymbol{x}_1^E, \ldots, \boldsymbol{x}_N^E) \rightarrow$   $\boxed{\boldsymbol{x}_1^E}\ \boxed{\boldsymbol{x}_2^E}\ \boxed{\boldsymbol{x}_3^E}\ \boxed{\boldsymbol{x}_4^E}\ \boxed{\boldsymbol{x}_5^E}\ \cdots\ \boxed{\boldsymbol{x}_{n'}^E}\ \cdots\ \boxed{\boldsymbol{x}_{v}^E}\ \cdots\ \boxed{\boldsymbol{x}_N^E}$

Encoder    **Encoder**

Acoustic features: $X = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_T) \rightarrow$

N. Moritz, T. Hori, and J. Le Roux, "Triggered attention for end-to-end speech recognition," in Proc. ICASSP, May 2019, pp. 5666–5670.

# Frame-Synchronous Decoding

Frame-synchronous CTC prefix beam search [1]:

$\ell_7 = (\text{<sos>, Hey, Word, World})$  $\ell_8 = (\text{<sos>, Hello, World})$

Set of prefix sequences after pruning:



$$\log p(\ell|X_{1:n}^E) = \log p_{\text{prfx}}(\ell|X_{1:n}^E) + \alpha \log p_{\text{LM}}(\ell) + \beta |\ell|$$

$\ell$: prefix sequence
$X_{1:n}^E$: Encoder state sequence for frame $(1, \dots, n)$
$p_{\text{prfx}}$: CTC prefix probability
$p_{\text{LM}}$: Language model (LM) probability
$\alpha$: LM weight
$\beta$: insertion bonus weight
$|\ell|$: prefix sequence length

CTC probability

Encoder output $X_E$:

[1] A. L. Maas, A. Y. Hannun, D. Jurafsky, and A. Y. Ng, "Firstpass large vocabulary continuous speech recognition using bidirectional recurrent DNNs," arXiv preprint arXiv:1408.2873, 2014.

# Frame-Synchronous Decoding

Frame-synchronous one-pass TA decoding [1]:

$$\log p_{\text{joint}}(\ell|X_{1:n}^E) = \lambda \log p_{\text{prfx}}(\ell|X_{1:n}^E) + (1-\lambda) \log p_{\text{ta}}(\ell|X_{1:\nu}^E) + \alpha \log p_{\text{LM}}(\ell) + \beta|\ell|$$
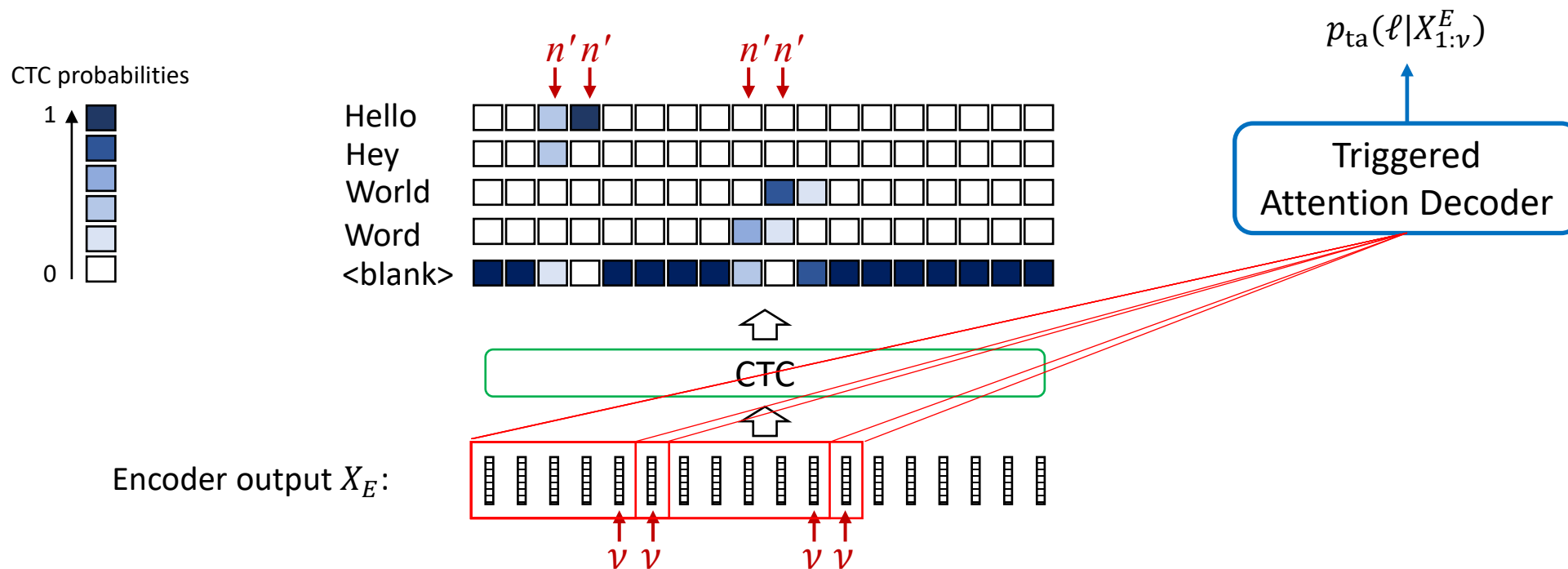
$p_{\text{ta}}$: Triggered attention probability
$\lambda$: CTC weight
$\nu = \text{n}' + \varepsilon^{\text{dec}}$
$\text{n}'$: trigger frame
$\varepsilon^{\text{dec}}$: decoder look-ahead



[1] N. Moritz, T. Hori, and J. Le Roux, "Streaming end-to-end speech recognition with joint CTC-attention based models," in Proc. ASRU, Dec. 2019, pp. 936–943.

# LibriSpeech Word Error Rates (WERs) [%]

| Encoder | Full-sequence CTC-attention decoding [1,2] | | | |
|---|---|---|---|---|
| | Clean | | Other | |
| | Dev | Test | Dev | Test |
| Full-sequence | 2.4 | 2.7 | 6.0 | 6.1 |

| Time-restricted encoder | Frame-synchronous CTC prefix beam search | | | | TA: $\varepsilon^{dec} = 18$, delay: $\varepsilon^{dec} \cdot 40$ ms $= 720$ ms | | | |
|---|---|---|---|---|---|---|---|---|
| | Clean | | Other | | Clean | | Other | |
| $\varepsilon^{enc}$ / delay* | Dev | Test | Dev | Test | Dev | Test | Dev | Test |
| 0 / 30 ms | 3.3 | 3.7 | 9.4 | 9.4 | 2.9 | 3.2 | 8.1 | 8.0 |
| 1 / 510 ms | 3.0 | 3.3 | 8.4 | 8.6 | 2.8 | 3.0 | 7.5 | 7.8 |
| 2 / 990 ms | 2.9 | 3.1 | 8.0 | 8.2 | 2.7 | 2.9 | 7.3 | 7.4 |
| 3 / 1470 ms | 2.8 | 2.9 | 7.8 | 8.1 | 2.7 | 2.8 | 7.1 | 7.2 |
| Full-sequence | 2.5 | 2.8 | 6.9 | 7.0 | 2.4 | 2.6 | 6.1 | 6.3 |

1.23 seconds

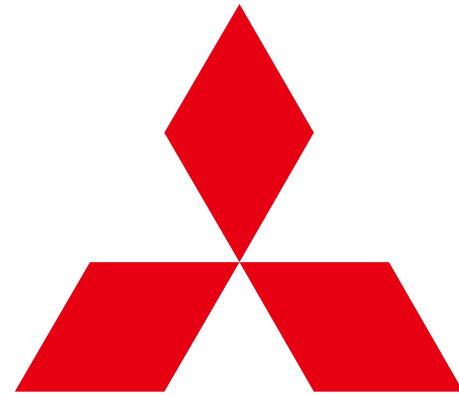* Algorithmic encoder delay: $E \cdot \varepsilon^{enc} \cdot$ frame−rate + CNN−delay

$E = 12$, frame-rate $= 40$ ms, CNN-delay $= 30$ ms

[1] S. Watanabe, T. Hori, S. Kim, J. R. Hershey, and T. Hayashi, "Hybrid CTC/attention architecture for end-to-end speech recognition," J. Sel. Topics Signal Processing, vol. 11, no. 8, pp. 1240–1253, 2017.
[2] S. Karita, N. Yalta, S. Watanabe, M. Delcroix, A. Ogawa, and T. Nakatani, "Improving transformer-based end-to-end speech recognition with connectionist temporal classification and language model integration," in Proc. ISCA Interspeech, Sep. 2019, pp. 1408–1412.

# Conclusions

- The triggered attention (TA) concept enables frame-synchronous decoding with an encoder-decoder based model for the first time.

- The TA concept enables joint scoring of an CTC and attention-based decoder model in a streaming fashion.

- The proposed system achieves state-of-the-art results for streaming end-to-end ASR on the LibriSpeech corpus.