

STREAMING AUTOMATIC SPEECH RECOGNITION WITH THE TRANSFORMER MODEL

Niko Moritz, Takaaki Hori, Jonathan Le Roux

Mitsubishi Electric Research Laboratories (MERL), Cambridge,
MA, USA

ASR (Automatic Speech Recognition)

Traditional

HMM-DNN



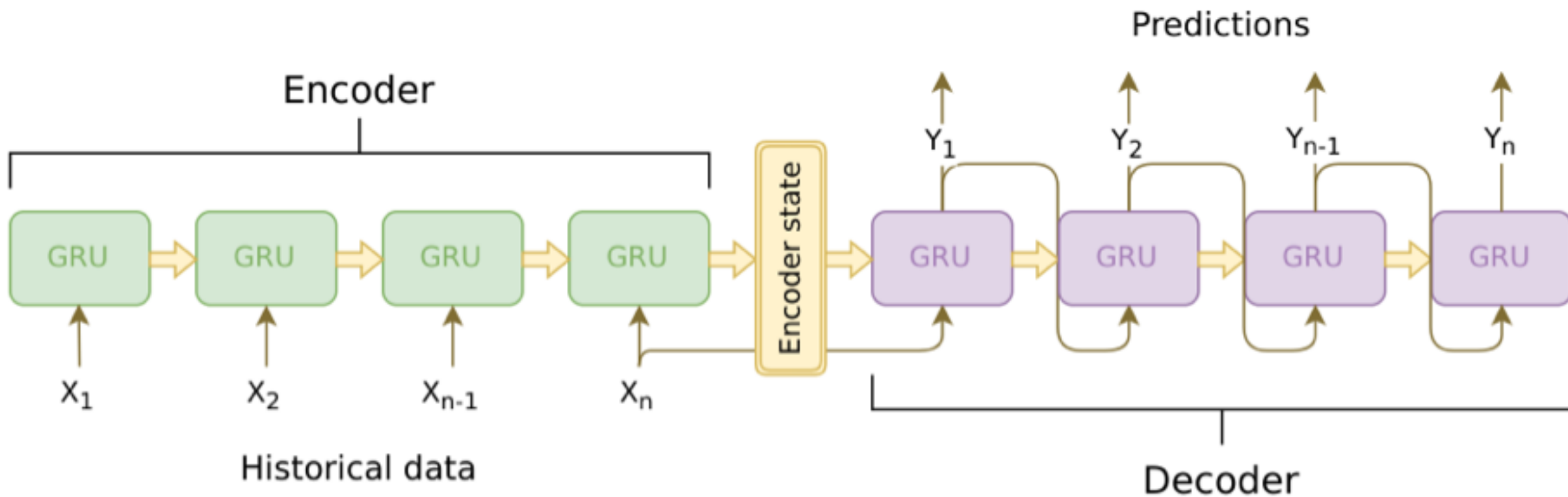
End-to-End

CTC (Connectionist Temporal Classification)

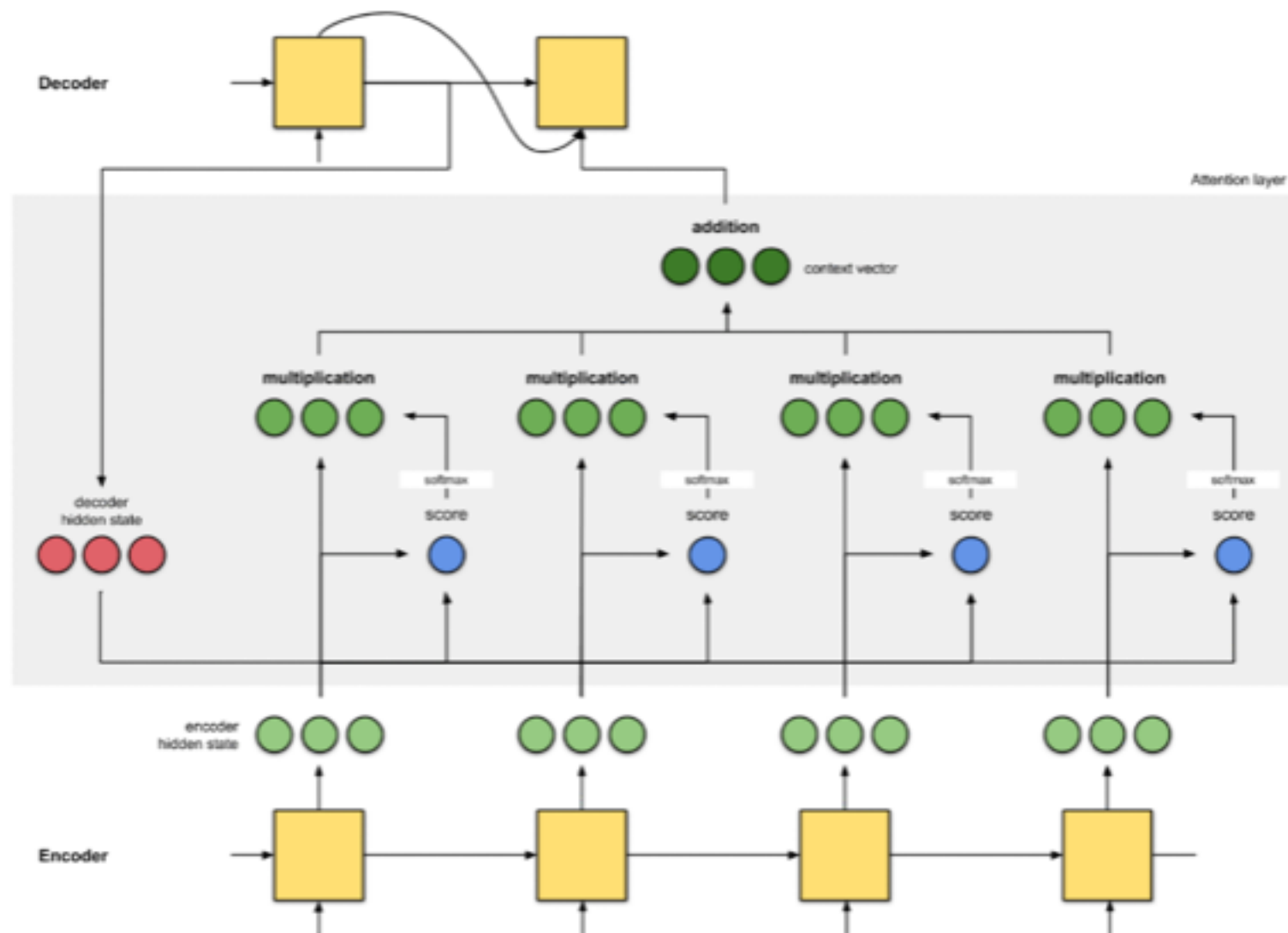
Attention Seq-to-Seq

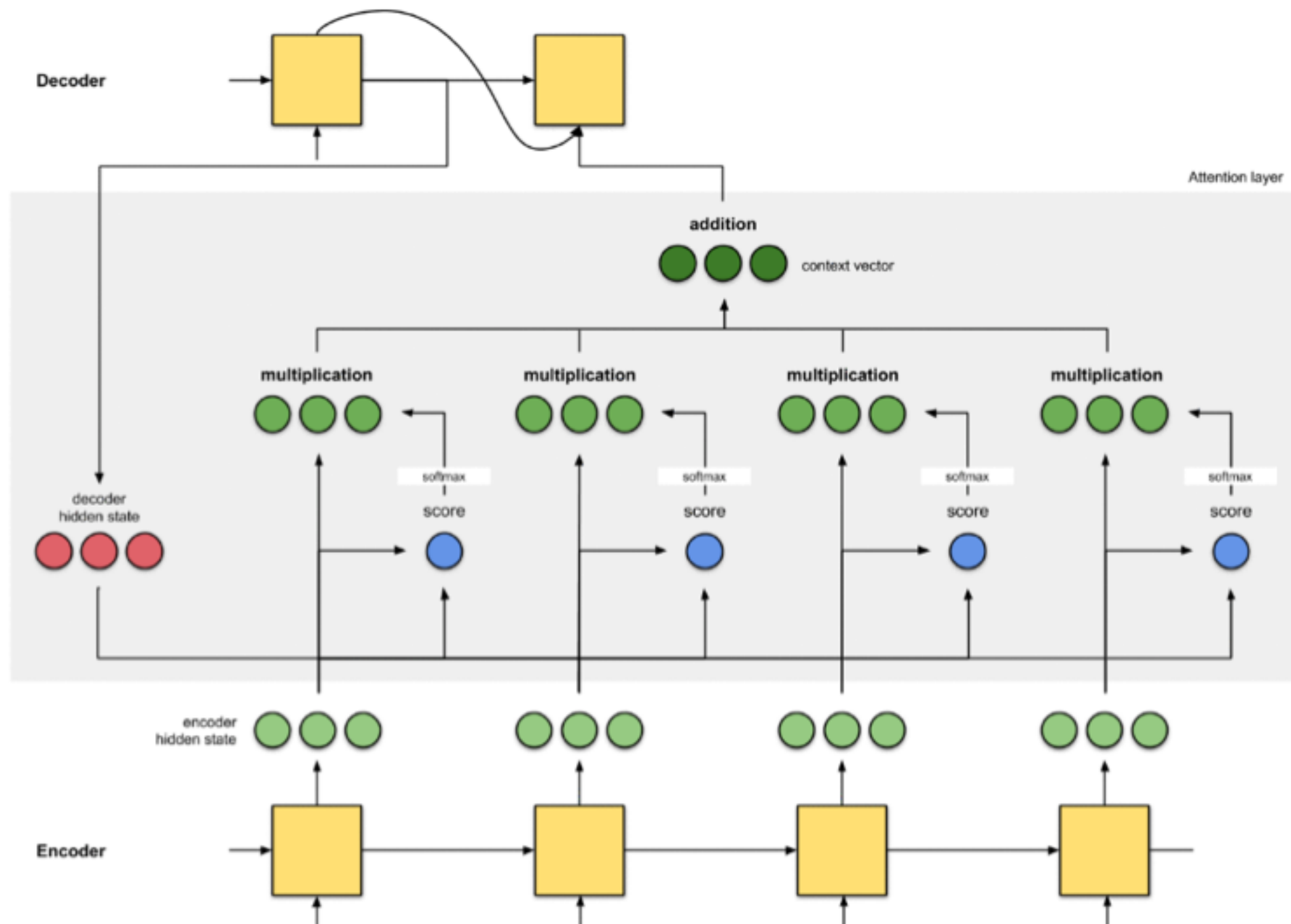
RNN-T

Seq-to-Seq

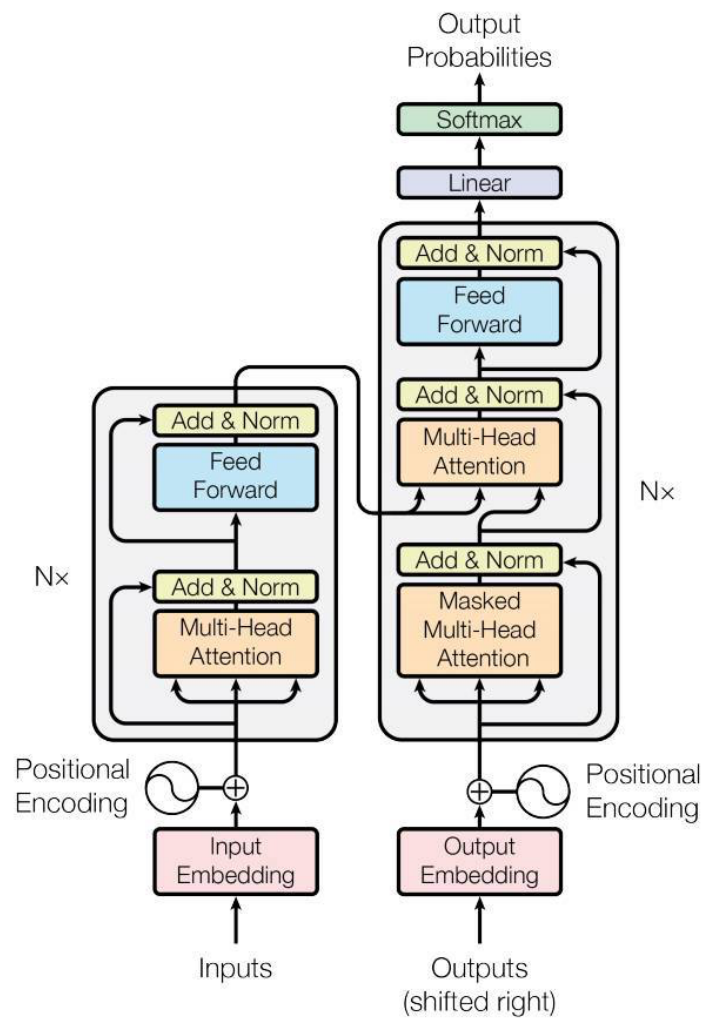


Attention Seq-to-Seq (RNN)

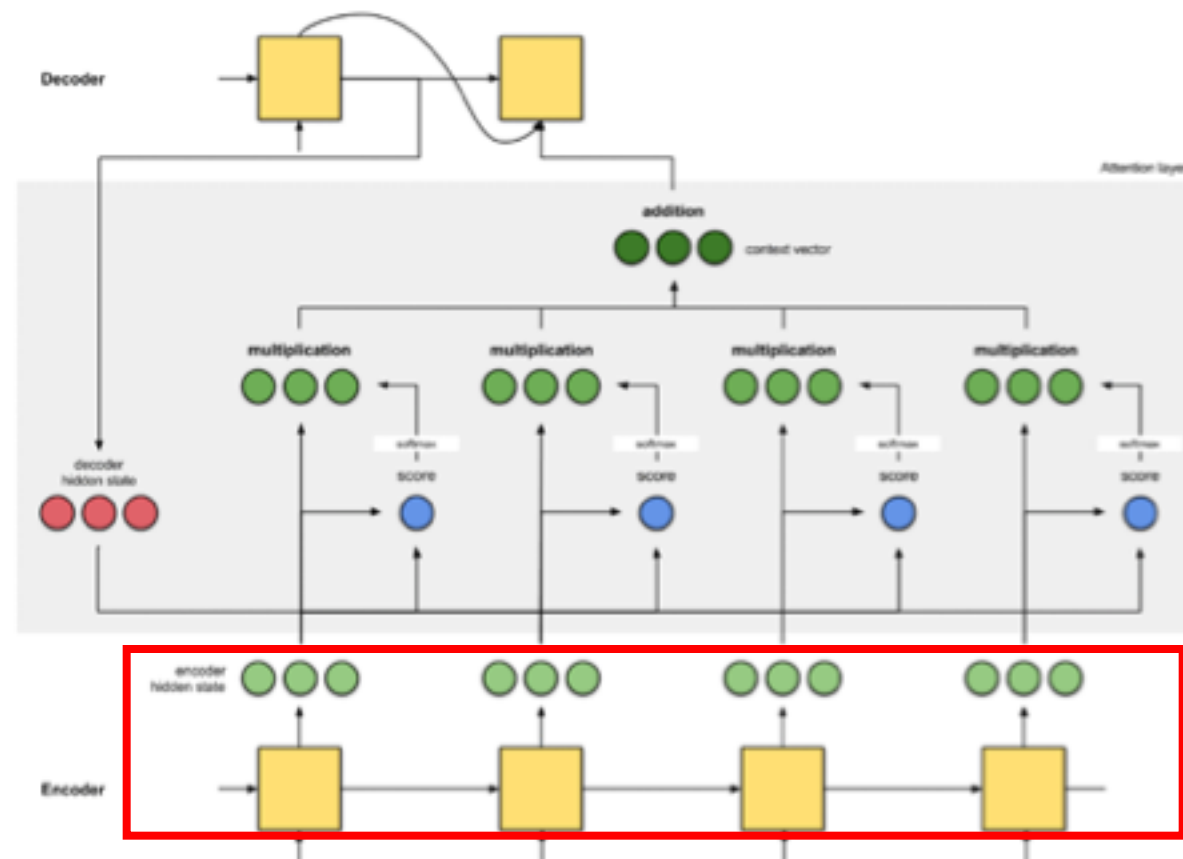
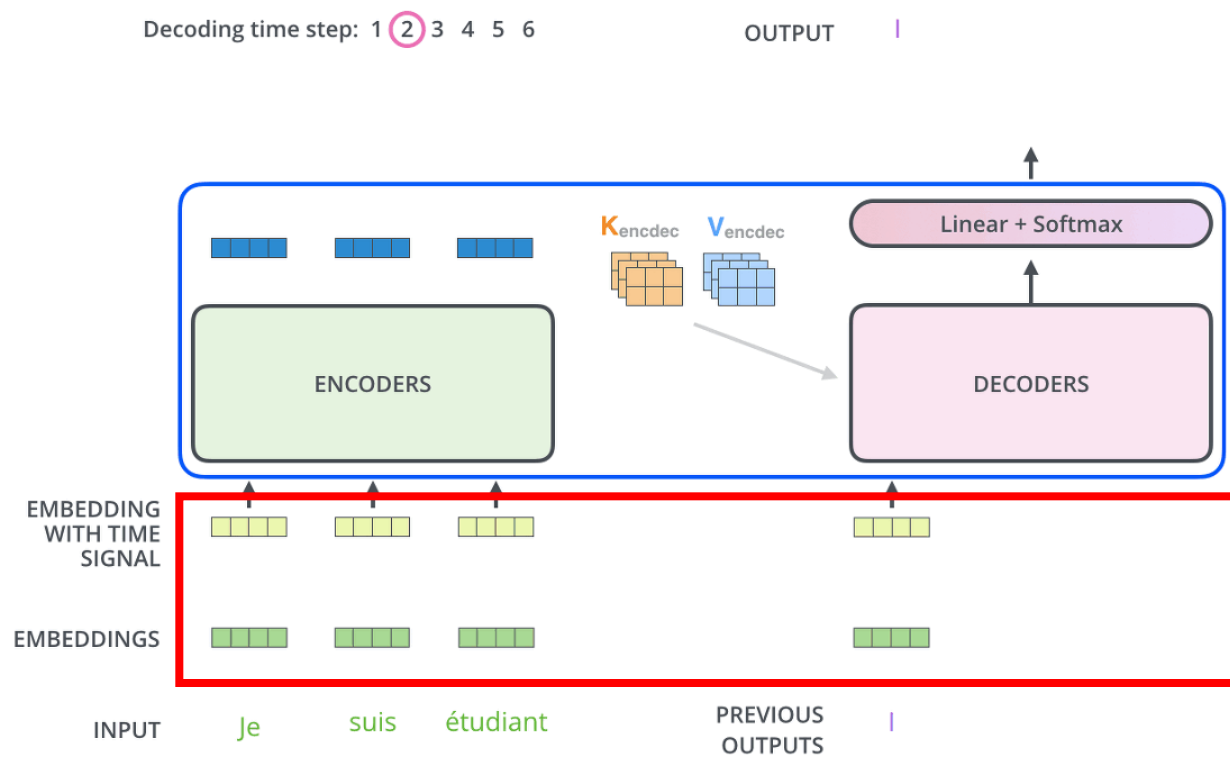




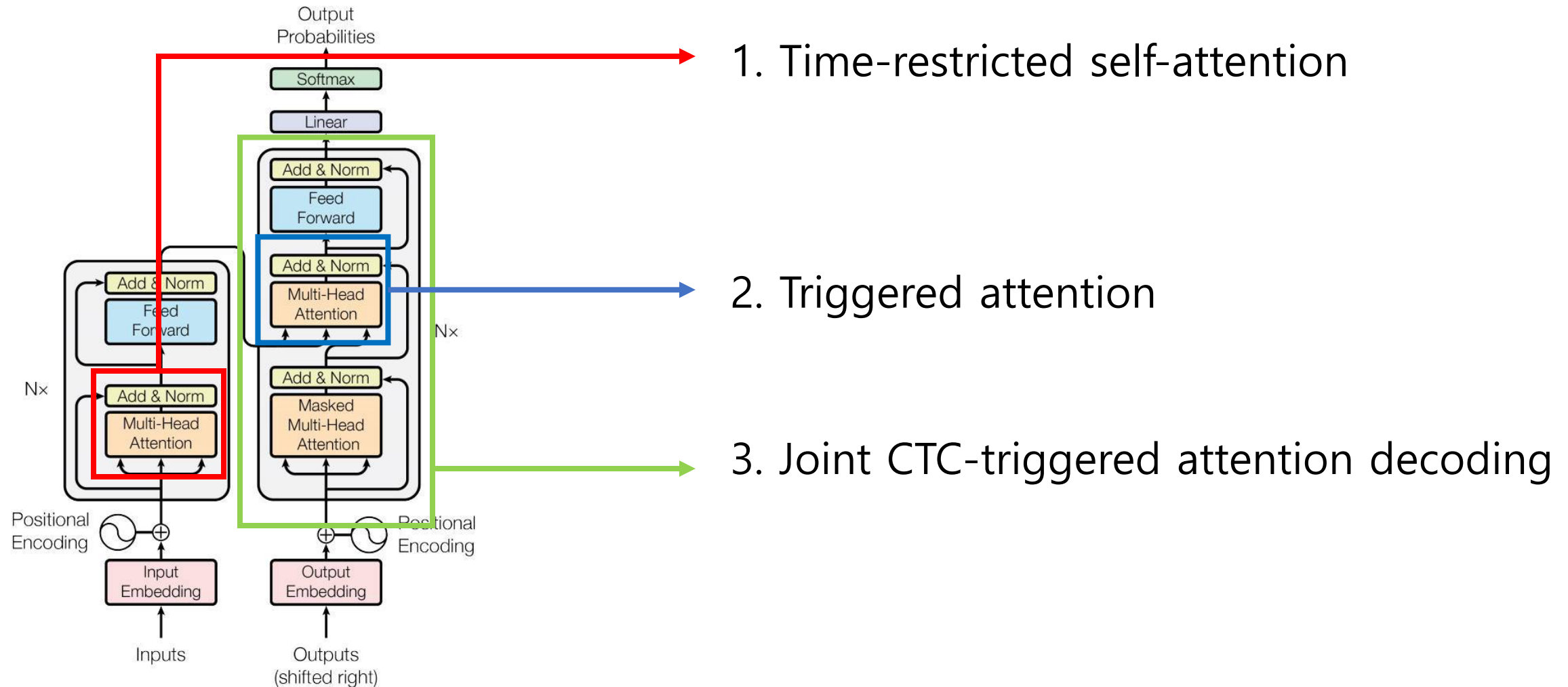
Attention Seq-to-Seq (Transformer)



Attention Seq-to-Seq



Transformer in streaming fashion



Time-restricted self-attention

$$X_0 = \text{ENCCNN}(X),$$

$$X_E = \text{ENCSA}(X_0),$$

$$X'_e = X_{e-1} + \text{MHA}_e(X_{e-1}, X_{e-1}, X_{e-1}),$$

Self-attention

$$X_e = X'_e + \text{FF}_e(X'_e),$$

Residual connection

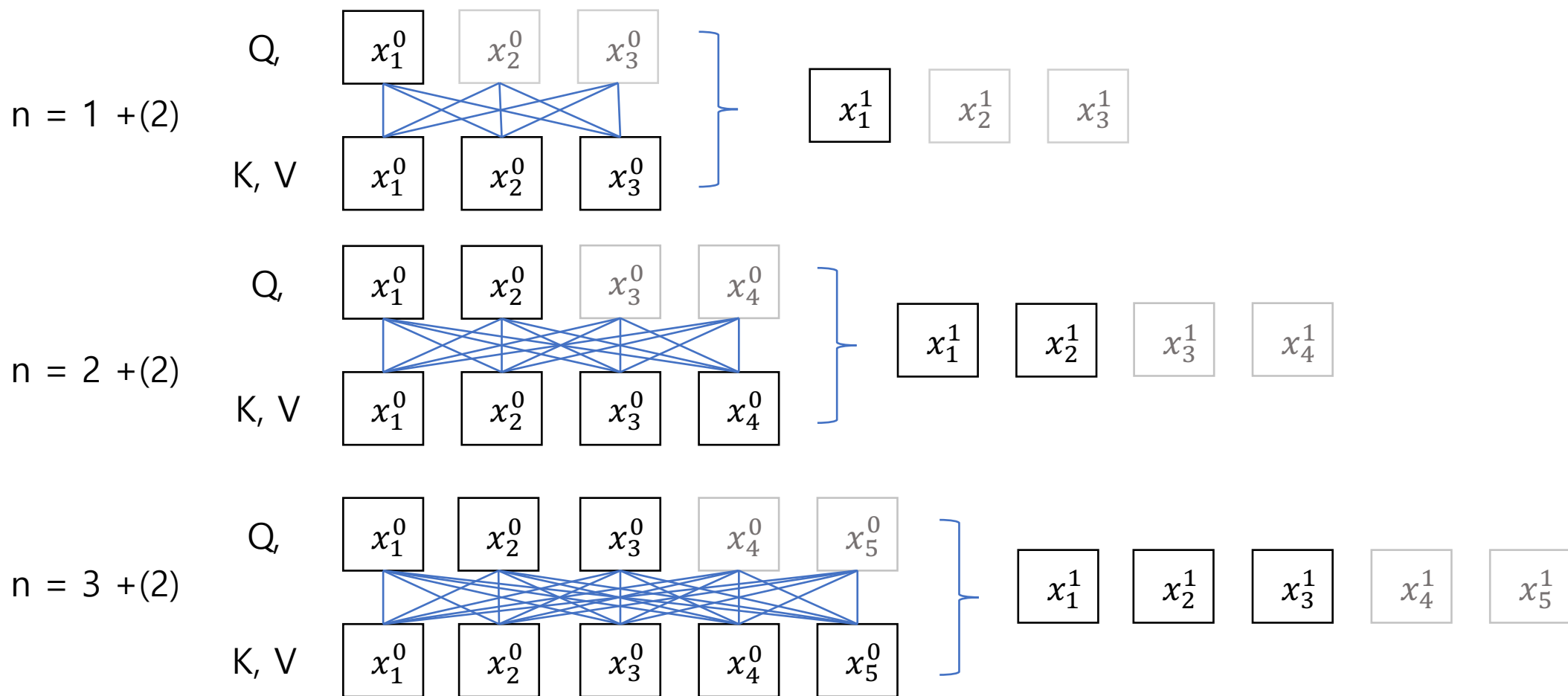
$$\text{with } \text{FF}_e(X'_e) = \text{ReLU}(X'_e W_{e,1}^{\text{ff}} + b_{e,1}^{\text{ff}}) W_{e,2}^{\text{ff}} + b_{e,2}^{\text{ff}},$$

$$\mathbf{x}_{1:n}^E = \text{ENCSA}^{\text{tr}}(\mathbf{x}_{1:n}^0 + \epsilon^{\text{enc}}),$$

Look-ahead frame

Time-restricted self-attention

$$\mathbf{x}_{1:n}^E = \text{ENC-SA}^{\text{tr}}(\mathbf{x}_{1:n+\varepsilon^{\text{enc}}}^0), \quad \text{If, } \varepsilon^{\text{enc}} = 2, E=1$$



Triggered attention

$$p_{\text{ta}}(Y|X_E) = \prod_{l=1}^L p(y_l | \mathbf{y}_{1:l-1}, \mathbf{x}_{1:\nu_l}^E)$$

$$p(y_l | \mathbf{y}_{1:l-1}, \mathbf{x}_{1:\nu_l}^E) = \text{DECTA}(\mathbf{x}_{1:\nu_l}^E, \mathbf{y}_{1:l-1}),$$

$$Y'_{d,l} = \mathbf{y}_{1:l-1}^{d-1} + \text{MHA}_d^{\text{self}}(\mathbf{y}_{1:l-1}^{d-1}, \mathbf{y}_{1:l-1}^{d-1}, \mathbf{y}_{1:l-1}^{d-1}),$$

$$Y''_{d,l} = Y'_{d,l} + \text{MHA}_d^{\text{dec}}(Y'_{d,l}, \mathbf{x}_{1:\nu_l}^E, \mathbf{x}_{1:\nu_l}^E),$$

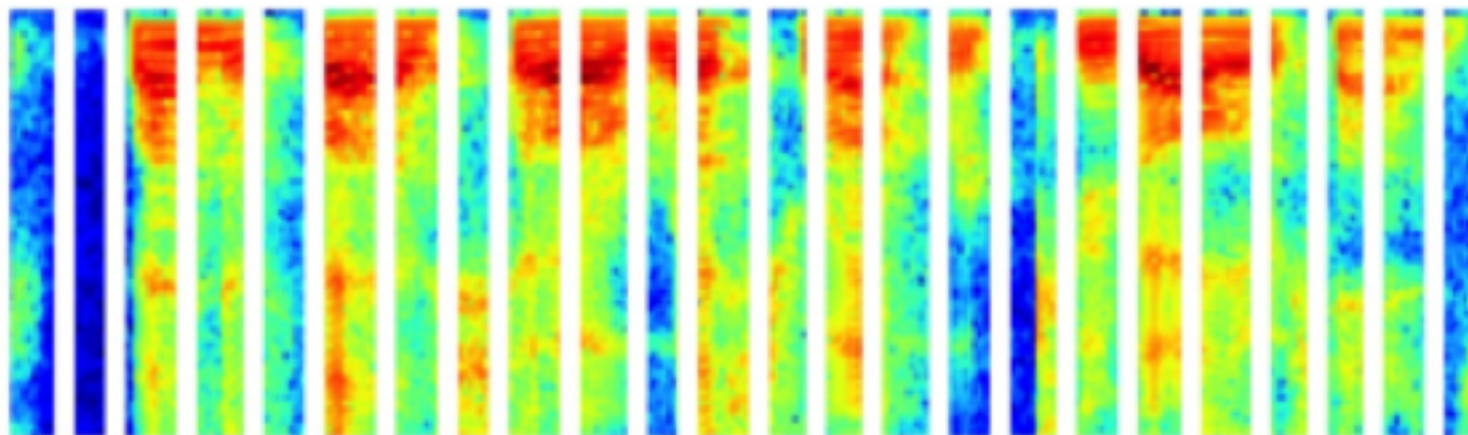
$$\mathbf{y}_{1:l-1}^d = Y''_{d,l} + \text{FF}_d(Y''_{d,l}),$$

$$\mathcal{L} = -\gamma \log p_{\text{ctc}} - (1 - \gamma) \log p_{\text{ta}},$$

Triggered attention

with $\nu_l = n'_l + \varepsilon^{\text{dec}}$, where n'_l denotes the position of the first occurrence of label y_l in the CTC forced alignment sequence [12, 14],

— T — — H — — EE — — — — C — — AA — — T — — —



Triggered attention

$$p_{\text{ta}}(Y|X_E) = \prod_{l=1}^L p(y_l | \mathbf{y}_{1:l-1}, \mathbf{x}_{1:\nu_l}^E)$$

$$p(y_l | \mathbf{y}_{1:l-1}, \mathbf{x}_{1:\nu_l}^E) = \text{DECTA}(\mathbf{x}_{1:\nu_l}^E, \mathbf{y}_{1:l-1}),$$

$$Y'_{d,l} = \mathbf{y}_{1:l-1}^{d-1} + \text{MHA}_d^{\text{self}}(\mathbf{y}_{1:l-1}^{d-1}, \mathbf{y}_{1:l-1}^{d-1}, \mathbf{y}_{1:l-1}^{d-1}),$$

$$Y''_{d,l} = Y'_{d,l} + \text{MHA}_d^{\text{dec}}(Y'_{d,l}, \mathbf{x}_{1:\nu_l}^E, \mathbf{x}_{1:\nu_l}^E),$$

$$\mathbf{y}_{1:l-1}^d = Y''_{d,l} + \text{FF}_d(Y''_{d,l}),$$

$$\mathcal{L} = -\gamma \log p_{\text{ctc}} - (1 - \gamma) \log p_{\text{ta}},$$

Joint CTC-triggered attention decoding

```

1: procedure DECODE( $X_E, p_{\text{ctc}}, \lambda, \alpha_0, \alpha, \beta, K, P, \theta_1, \theta_2$ )
2:    $\ell \leftarrow (\langle \text{sos} \rangle,)$ 
3:    $\Omega \leftarrow \{\ell\}, \Omega_{\text{ta}} \leftarrow \{\ell\}$ 
4:    $p_{\text{ab}}(\ell) \leftarrow 0, p_{\text{b}}(\ell) \leftarrow 1$ 
5:    $p_{\text{ta}}(\ell) \leftarrow 1$ 
6:   for  $n = 1, \dots, N$  do
7:      $\Omega_{\text{ctc}}, p_{\text{ab}}, p_{\text{b}} \leftarrow \text{CTCPREFIX}(p_{\text{ctc}}(n), \Omega, p_{\text{ab}}, p_{\text{b}})$ 
8:     for  $\ell$  in  $\Omega_{\text{ctc}}$  do ▷ Compute CTC prefix scores
9:        $p_{\text{prfx}}(\ell) \leftarrow p_{\text{ab}}(\ell) + p_{\text{b}}(\ell)$ 
10:       $\hat{p}_{\text{prfx}}(\ell) \leftarrow \log p_{\text{prfx}}(\ell) + \alpha_0 \log p_{\text{LM}}(\ell) + \beta|\ell|$ 
11:       $\hat{\Omega} \leftarrow \text{PRUNE}(\Omega_{\text{ctc}}, \hat{p}_{\text{prfx}}, K, \theta_1)$ 
12:      for  $\ell$  in  $\hat{\Omega}$  do ▷ Delete old prefixes in  $\Omega_{\text{ta}}$ 
13:        if  $\ell$  in  $\Omega_{\text{ta}}$  and  $\text{DCOND}(\ell, \hat{\Omega}, p_{\text{ctc}})$  then
14:          delete  $\ell$  in  $\Omega_{\text{ta}}$ 
15:      for  $\ell$  in  $\hat{\Omega}$  do ▷ Compute transformer scores
16:        if  $\ell$  not in  $\Omega_{\text{ta}}$  and  $\text{ACOND}(\ell, \hat{\Omega}, p_{\text{ctc}})$  then
17:           $p_{\text{ta}}(\ell) \leftarrow \text{DECTA}(\mathbf{x}_{1:n+\epsilon^{\text{dec}}}^E, \ell)$ 
18:          add  $\ell$  to  $\Omega_{\text{ta}}$ 
19:      for  $\ell$  in  $\hat{\Omega}$  do ▷ Compute joint scores
20:         $\hat{\ell} \leftarrow \ell$  if  $\ell$  in  $\Omega_{\text{ta}}$  else  $\ell_{:-1}$ 
21:         $p \leftarrow \lambda \log p_{\text{prfx}}(\ell) + (1 - \lambda) \log p_{\text{ta}}(\hat{\ell})$ 
22:         $p_{\text{joint}}(\ell) \leftarrow p + \alpha \log p_{\text{LM}}(\ell) + \beta|\ell|$ 
23:       $\Omega \leftarrow \text{MAX}(\hat{\Omega}, p_{\text{joint}}, P)$ 
24:       $\hat{\Omega} \leftarrow \text{PRUNE}(\hat{\Omega}, \hat{p}_{\text{prfx}}, P, \theta_2)$ 
25:       $\Omega \leftarrow \Omega + \hat{\Omega}$ 
26:      remove from  $\Omega_{\text{ta}}$  prefixes rejected due to pruning
27: return  $\text{MAX}(\hat{\Omega}, p_{\text{joint}}, 1)$ 

```

```

procedure PRUNE( $\Omega_{\text{in}}, \hat{p}_{\text{prfx}}, L, \theta$ )
   $\Omega_{\text{in}} \leftarrow \text{MAX}(\Omega_{\text{in}}, \hat{p}_{\text{prfx}}, L), \Omega_{\text{out}} \leftarrow \{\}$ 
  for  $\ell$  in  $\Omega_{\text{in}}$  do
    if  $\max(\hat{p}_{\text{prfx}}) - \theta < \hat{p}_{\text{prfx}}(\ell)$  then
      add  $\ell$  to  $\Omega_{\text{out}}$ 
  return  $\Omega_{\text{out}}$ 

procedure DCOND( $\ell, \hat{\Omega}, p_{\text{ctc}}$ )
   $c \leftarrow \ell_{-1}$  ▷ last element of  $\ell$ 
   $n_{\text{prev}} \leftarrow$  time index of  $\ell$  when added to  $\hat{\Omega}$ 
  if  $|\ell| > 1$  and  $n - n_{\text{prev}} > 2$  and  $p_{\text{ctc}}(n, c) > 0.01 >$ 
 $\max(p_{\text{ctc}}(n_{\text{prev}}, c), p_{\text{ctc}}(n_{\text{prev}} + 1, c))$  then
    return true
  else
    return false

procedure ACOND( $\ell, \hat{\Omega}, p_{\text{ctc}}$ )
   $c \leftarrow \ell_{-1}$  ▷ last element of  $\ell$ 
  if  $p_{\text{ctc}}(n, c) > \max(p_{\text{ctc}}(n + 1, c), p_{\text{ctc}}(n + 2, c))$  or
  any( $|\ell| > |\ell| + 1$  and  $\hat{\ell}$  starts with  $\ell$  for  $\hat{\ell}$  in  $\hat{\Omega}$ ) then
    return true
  else
    return false

```

Experiment

Dataset : LibriSpeech

Setting : transformer

small : $d_{\text{model}} = 256$, $d_{\text{ff}} = 2048$, $d_{\text{h}} = 4$, $E = 12$, $D = 6$

large : $d_{\text{model}} = 512$, $d_{\text{h}} = 8$

: joint CTC-TA decoding

CTC weight $\lambda = 0.5$, CTC LM weight $\alpha_0 = 0.7$, LM weight $\alpha = 0.5$,

pruning beam width $\theta_1 = 16.0$, pruning beam width $\theta_2 = 6.0$,

insertion bonus $\beta = 2.0$, pruning size $K = 300$, and pruning size $P = 30$.

Result

Theoretical delay

Encoder CNN : 30ms

Time restricted self-attention : $E \times \varepsilon^{enc} \times 40ms$

Triggered attention : $\varepsilon^{dec} \times 40ms$

WER : 2.8% / 7.3% $\rightarrow \varepsilon^{enc}=4, \varepsilon^{dec}=6 \rightarrow 2190ms$

Result

System	CTC-attention dec.				CTC beam search				Att. beam search			
	clean		other		clean		other		clean		other	
	dev	test	dev	test	dev	test	dev	test	dev	test	dev	test
baseline	4.7	4.9	13.0	12.9	6.1	6.1	15.7	15.9	6.0	7.8	14.5	14.9
+RNN-LM	2.9	3.1	8.0	8.4	3.1	3.4	9.3	9.6	4.7	7.2	10.7	11.5
+SpecAug.	2.4	2.8	6.4	6.7	2.9	3.2	7.6	7.9	4.2	5.2	8.3	8.6
+large	2.4	2.7	6.0	6.1	2.5	2.8	6.9	7.0	4.1	5.0	7.9	8.0

ϵ^{enc}	CTC beam search				TA: $\epsilon^{\text{dec}} = 6$				TA: $\epsilon^{\text{dec}} = 12$				TA: $\epsilon^{\text{dec}} = 18$			
	clean		other		clean		other		clean		other		clean		other	
	dev	test	dev	test	dev	test	dev	test	dev	test	dev	test	dev	test	dev	test
1	3.0	3.3	8.4	8.6	2.9	3.2	8.1	8.2	2.8	3.1	7.5	8.1	2.8	3.0	7.5	7.8
2	2.9	3.1	8.0	8.2	2.8	2.9	7.4	7.8	2.7	2.9	7.2	7.6	-	-	-	-
4	2.8	2.9	7.8	7.9	2.6	2.8	7.2	7.3	-	-	-	-	-	-	-	-
∞	2.5	2.8	6.9	7.0	2.5	2.7	6.3	6.5	2.5	2.7	6.3	6.4	2.4	2.6	6.1	6.3