

RELEVANT LINKS



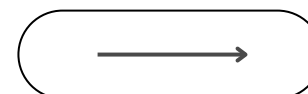
DATE

06/11/2025

# BLOCKDECODER: BOOSTING ASR DECODERS WITH **CONTEXT** AND **MERGER** MODULES



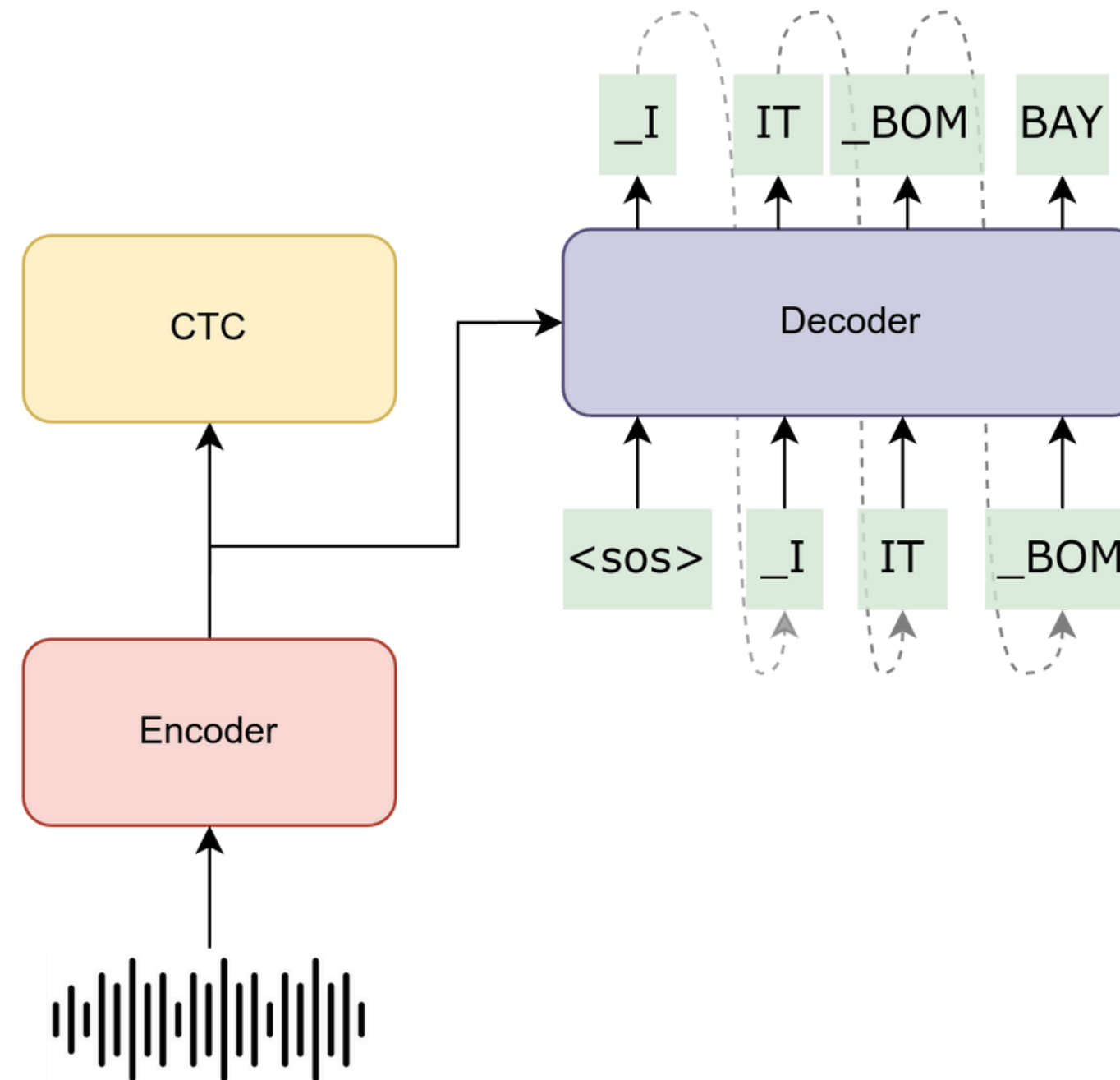
Darshan Prabhu Preethi Jyothi



PRESENTED BY  
Darshan Prabhu



# The Architecture of Focus





# Prior Literature

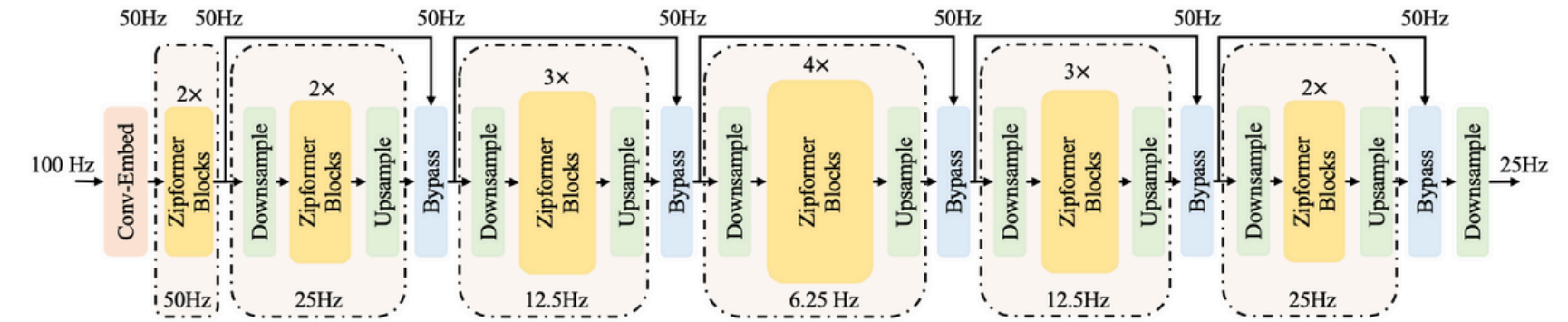
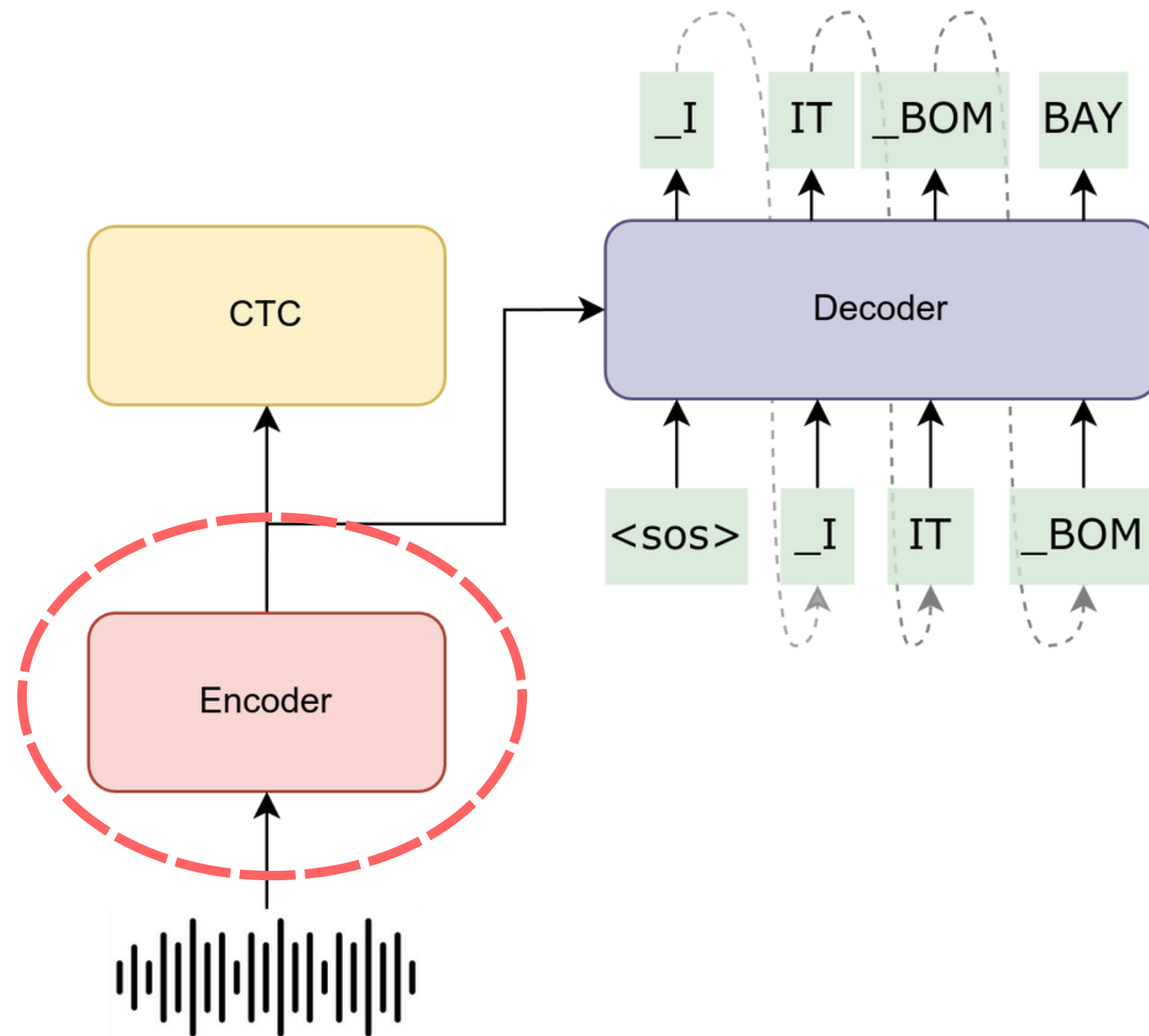
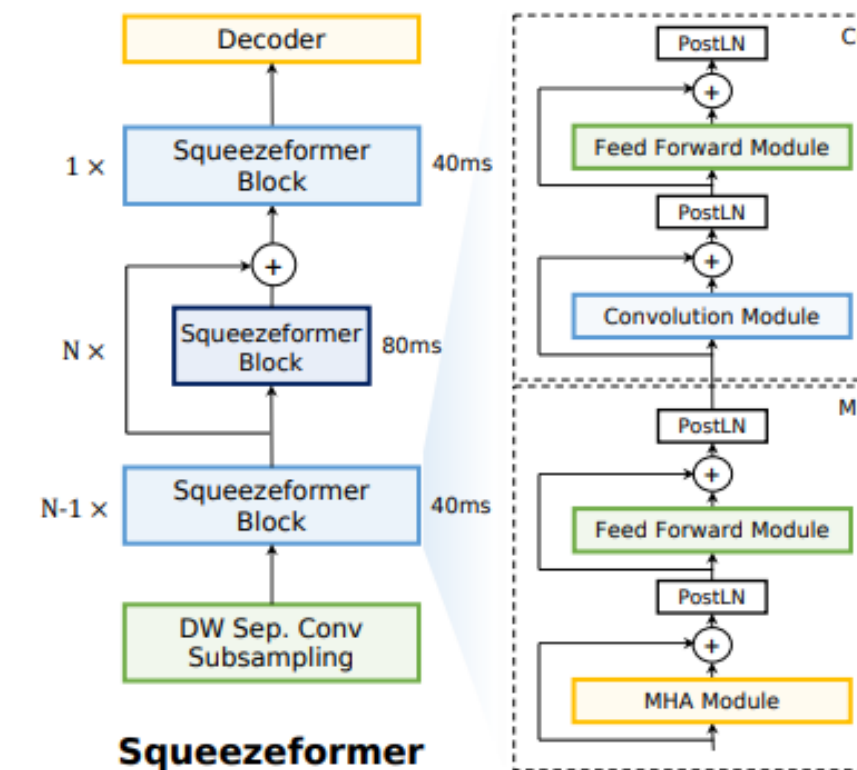


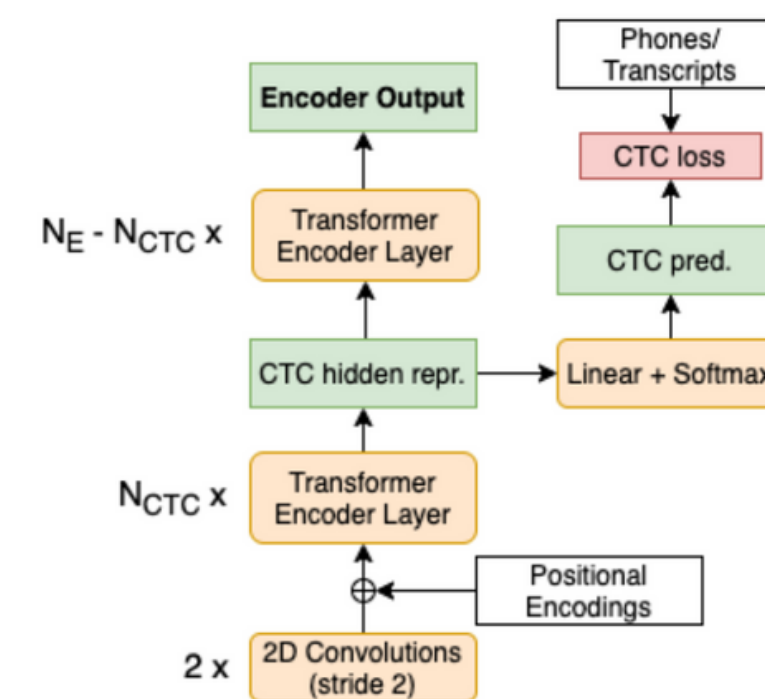
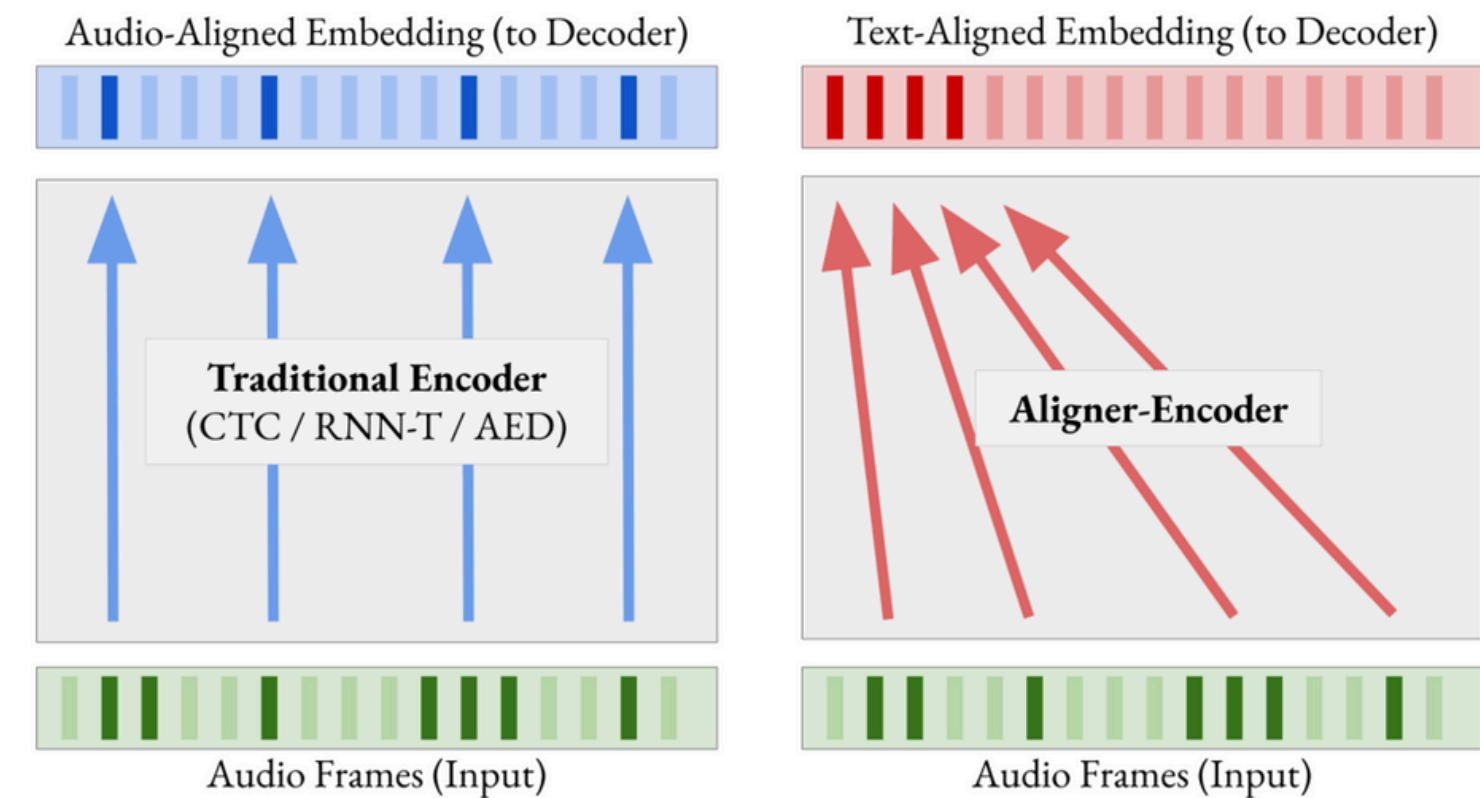
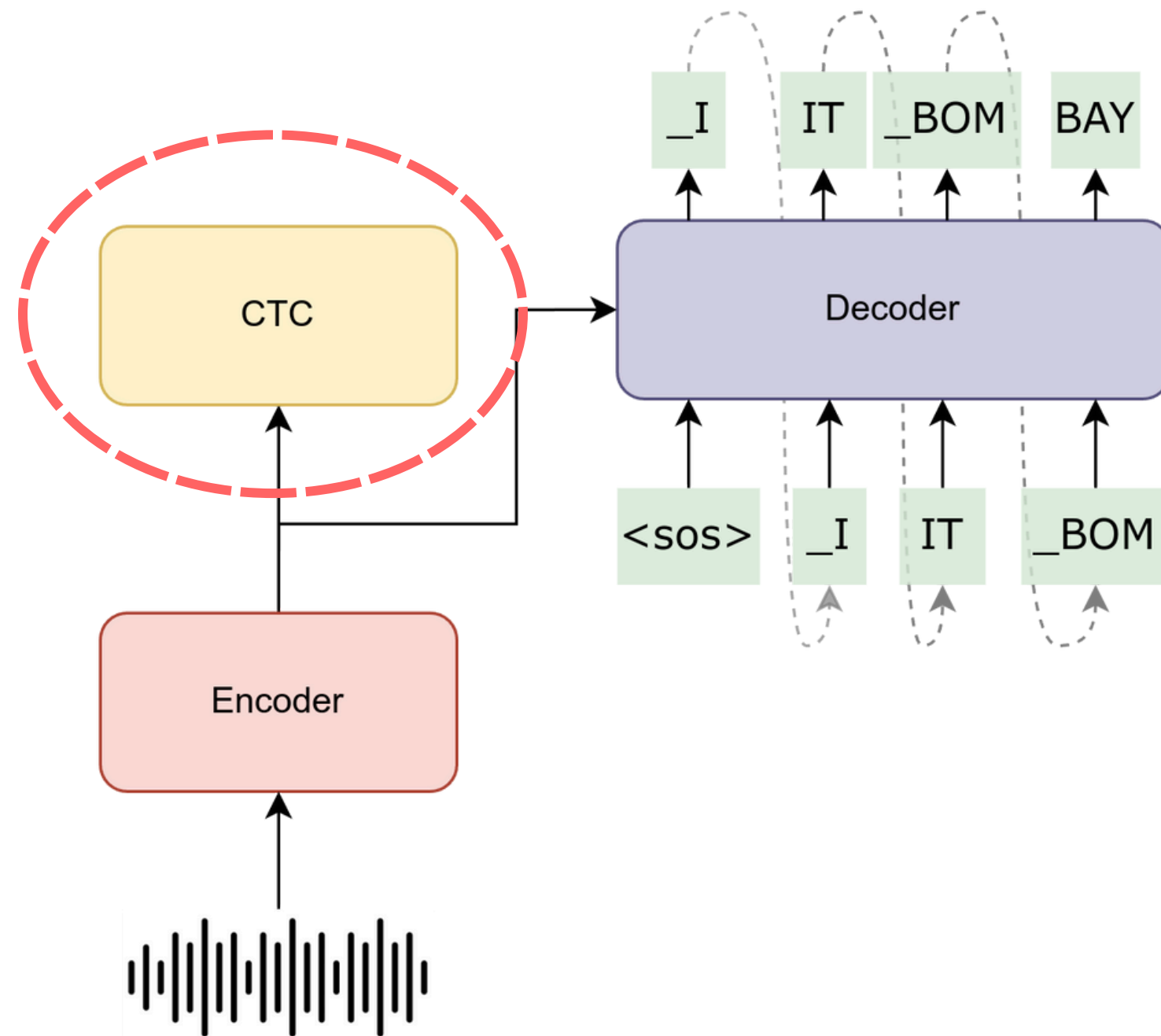
Figure 1: Overall architecture of Zipformer.



Reference: Yao, Z., Guo, L., Yang, X., Kang, W., Kuang, F., Yang, Y., & Povey, D. (2023). Zipformer: A faster and better encoder for automatic speech recognition. arXiv preprint arXiv:2310.11230. (ICLR 2024)



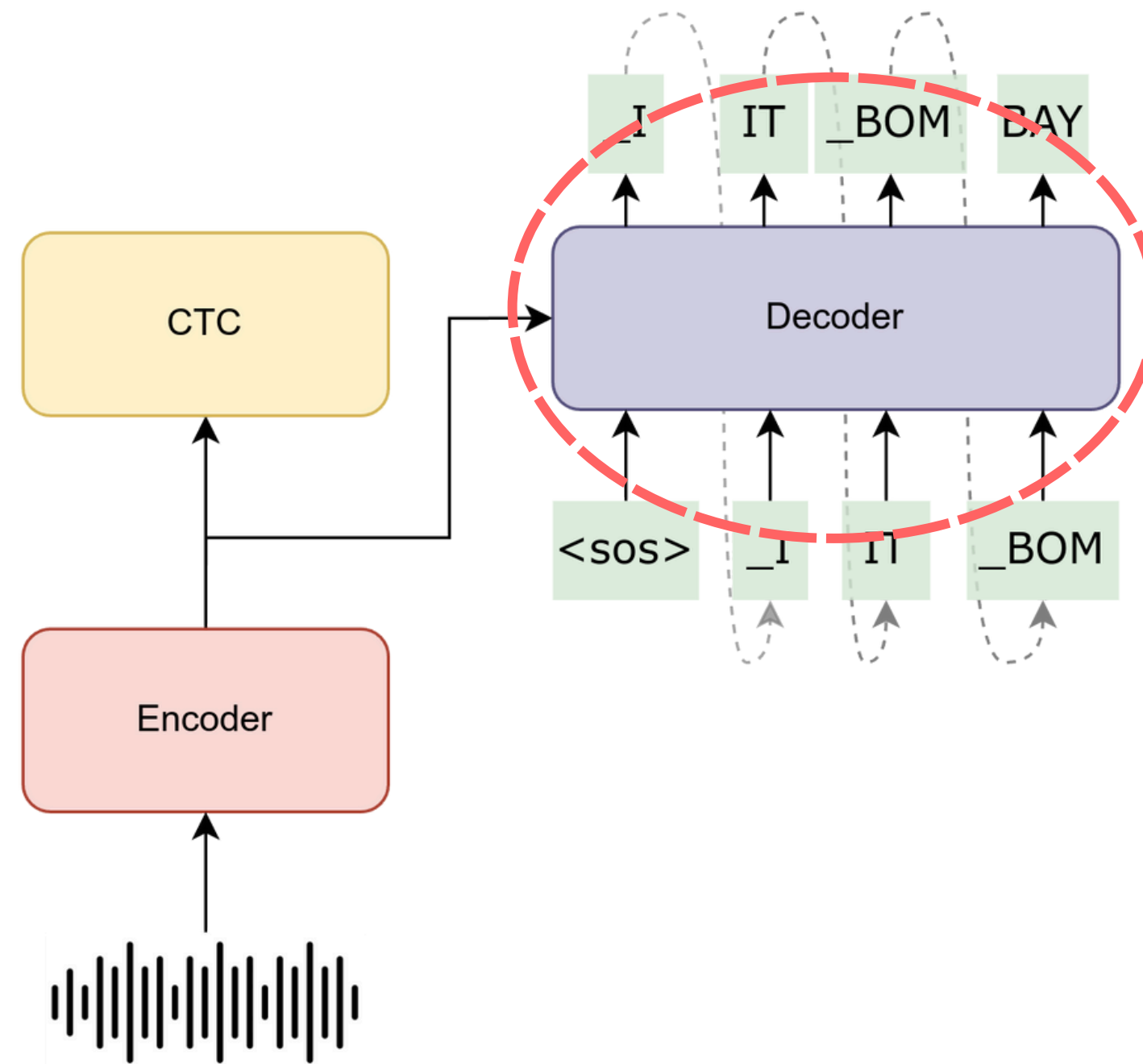
# Prior Literature (contd.)



Reference: Stooke, A., Prabhavalkar, R., Sim, K., & Moreno Mengibar, P. (2024). Aligner-Encoders: Self-Attention Transformers Can Be Self-Transducers. Advances in Neural Information Processing Systems, 37, 100318-100340. (NeurIPS 2024)

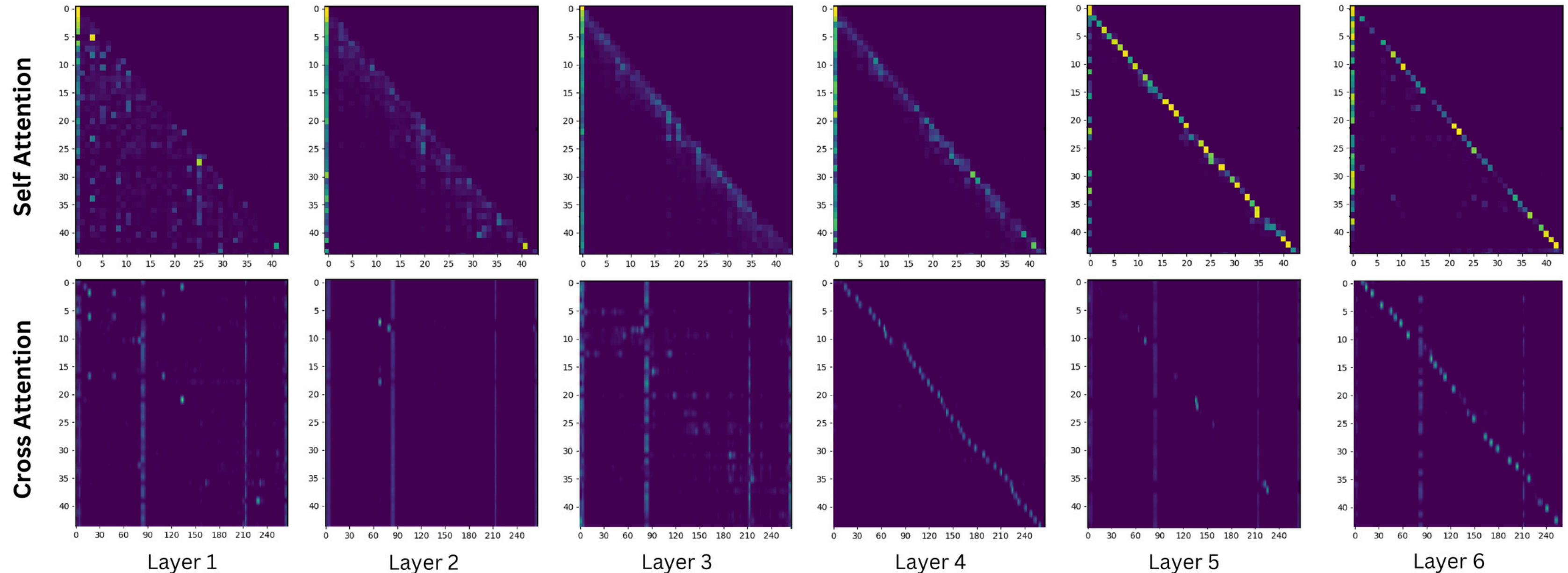


# What is missing ?





# Analysis of attention patterns

**1**

**Self-attention** increasingly focuses on **local context** as we progress **deeper into the decoder**.

**2**

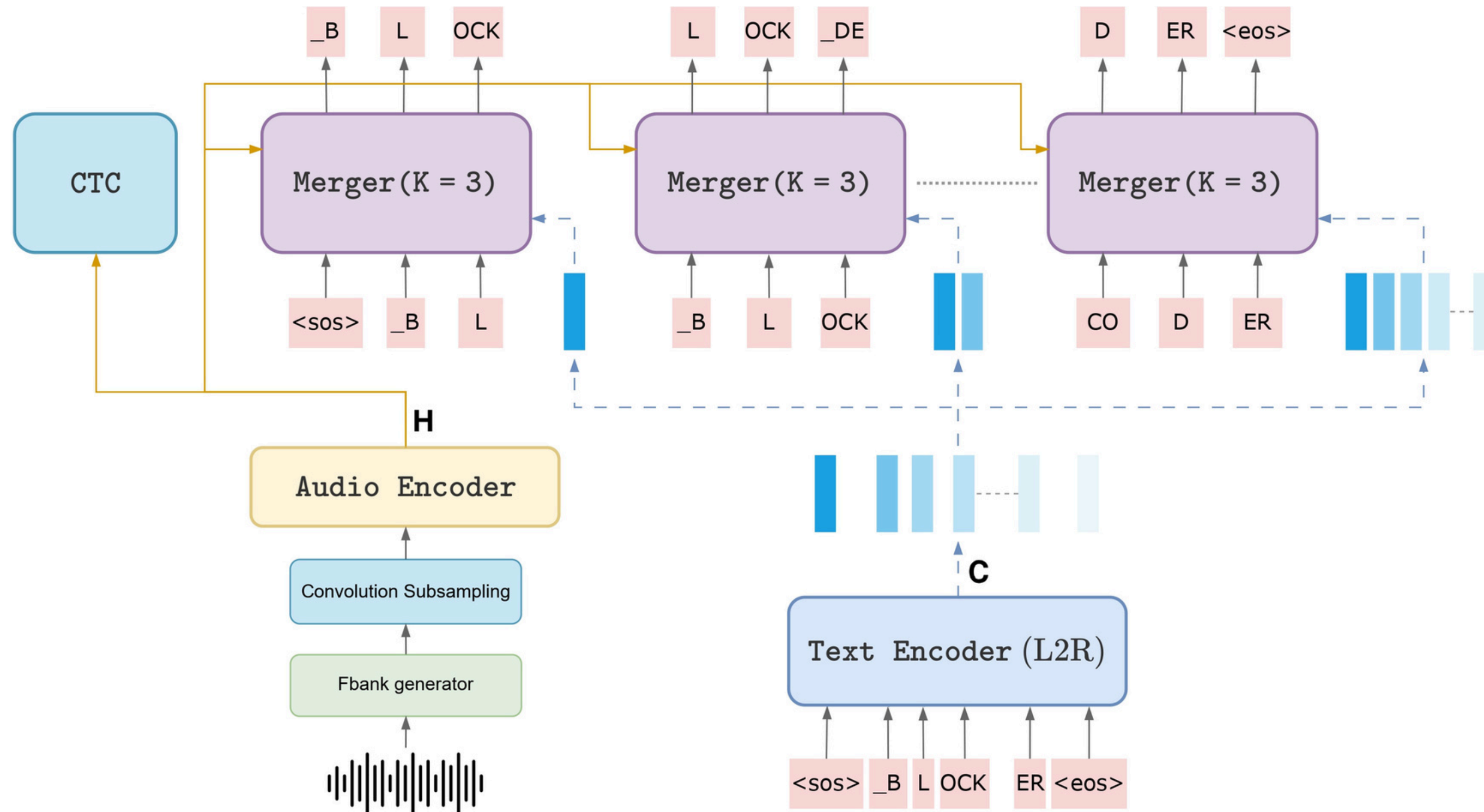
**Cross-attention** blocks appear to be **less effective** in the **initial decoder layers**.





# BlockDecoder

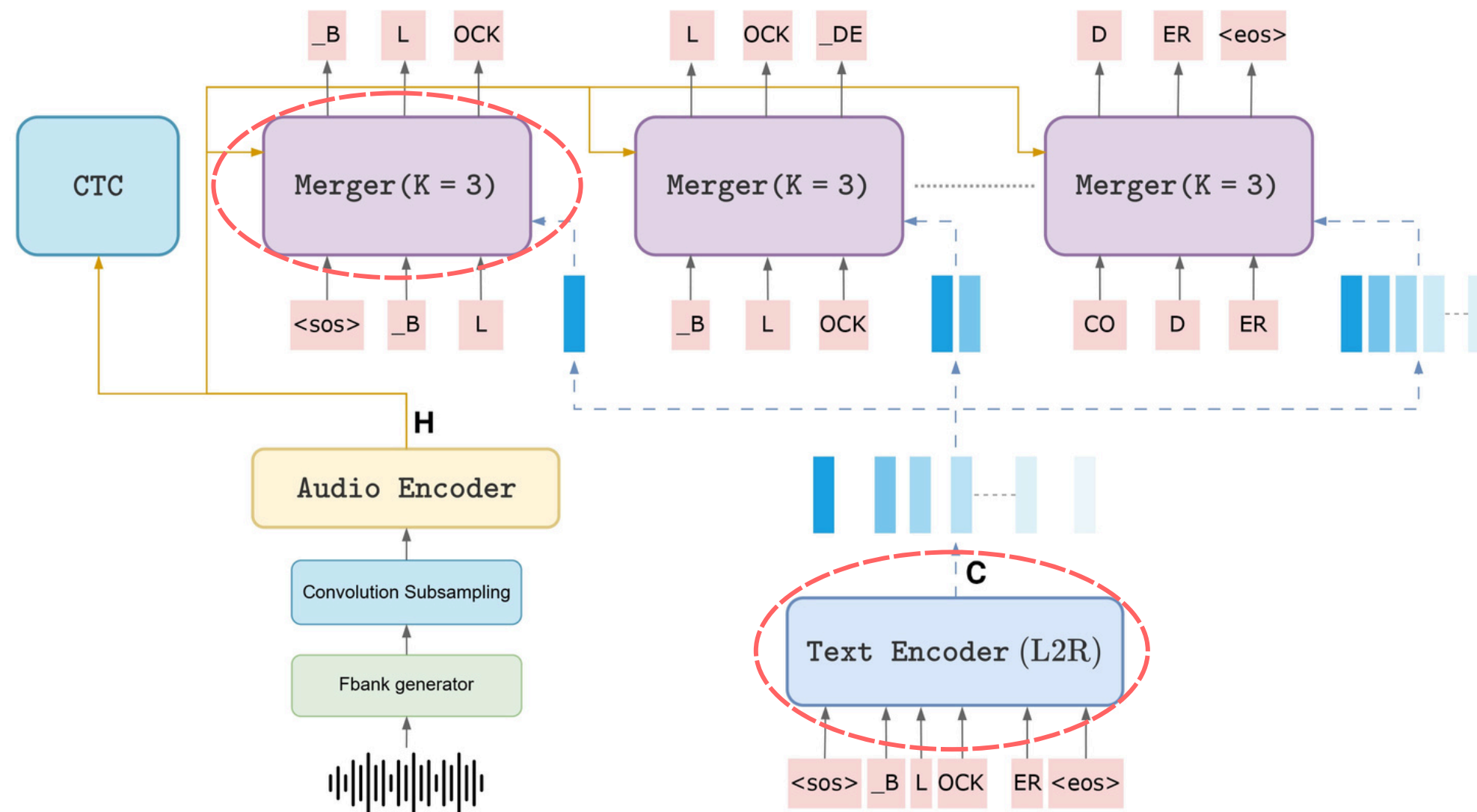
Boosting ASR Decoders with Context and Merger Modules





# BlockDecoder

Boosting ASR Decoders with Context and Merger Modules



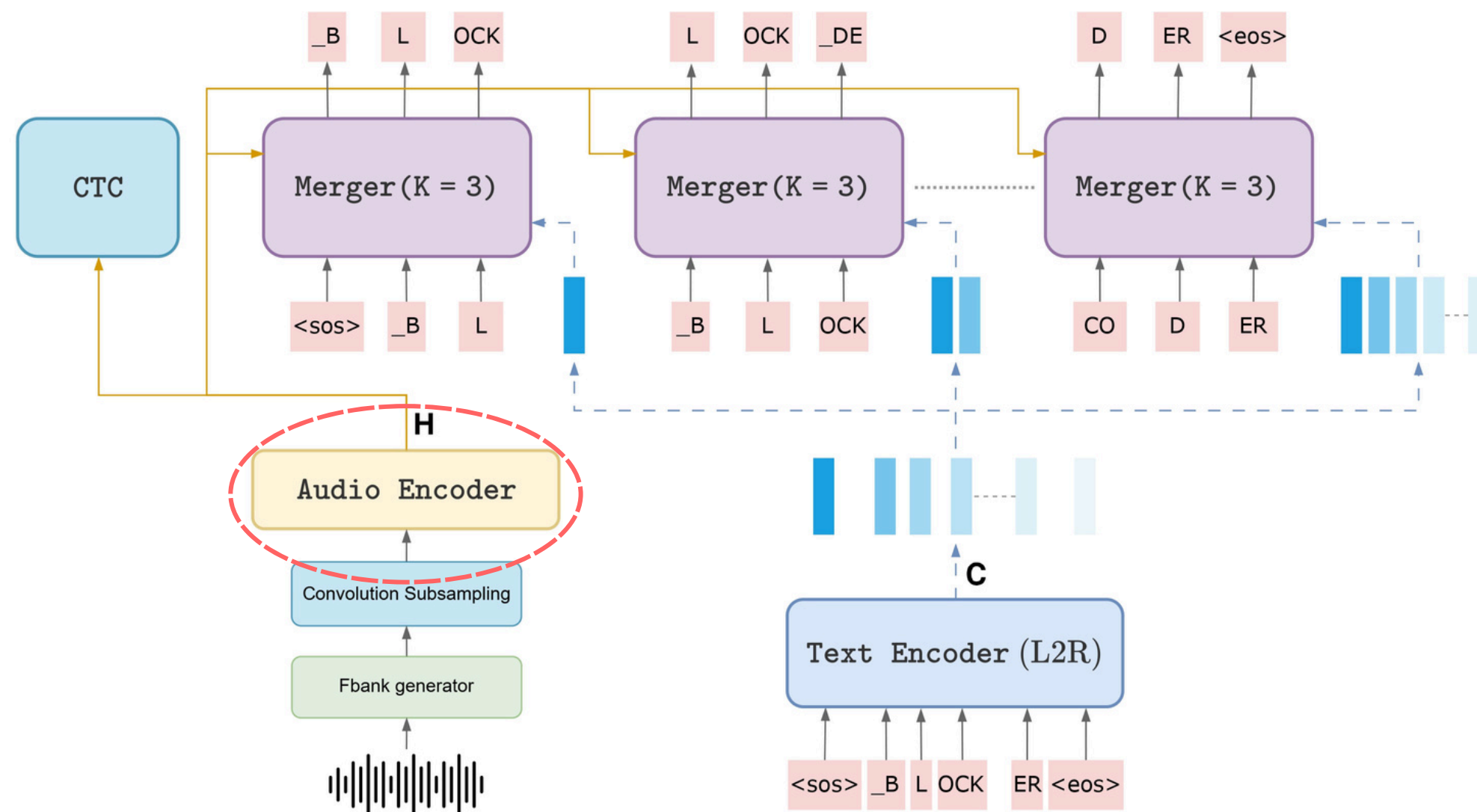
## Key Insights

- Replace the Traditional decoder with **Text Encoder** and **Merger**





## Boosting ASR Decoders with Context and Merger Modules



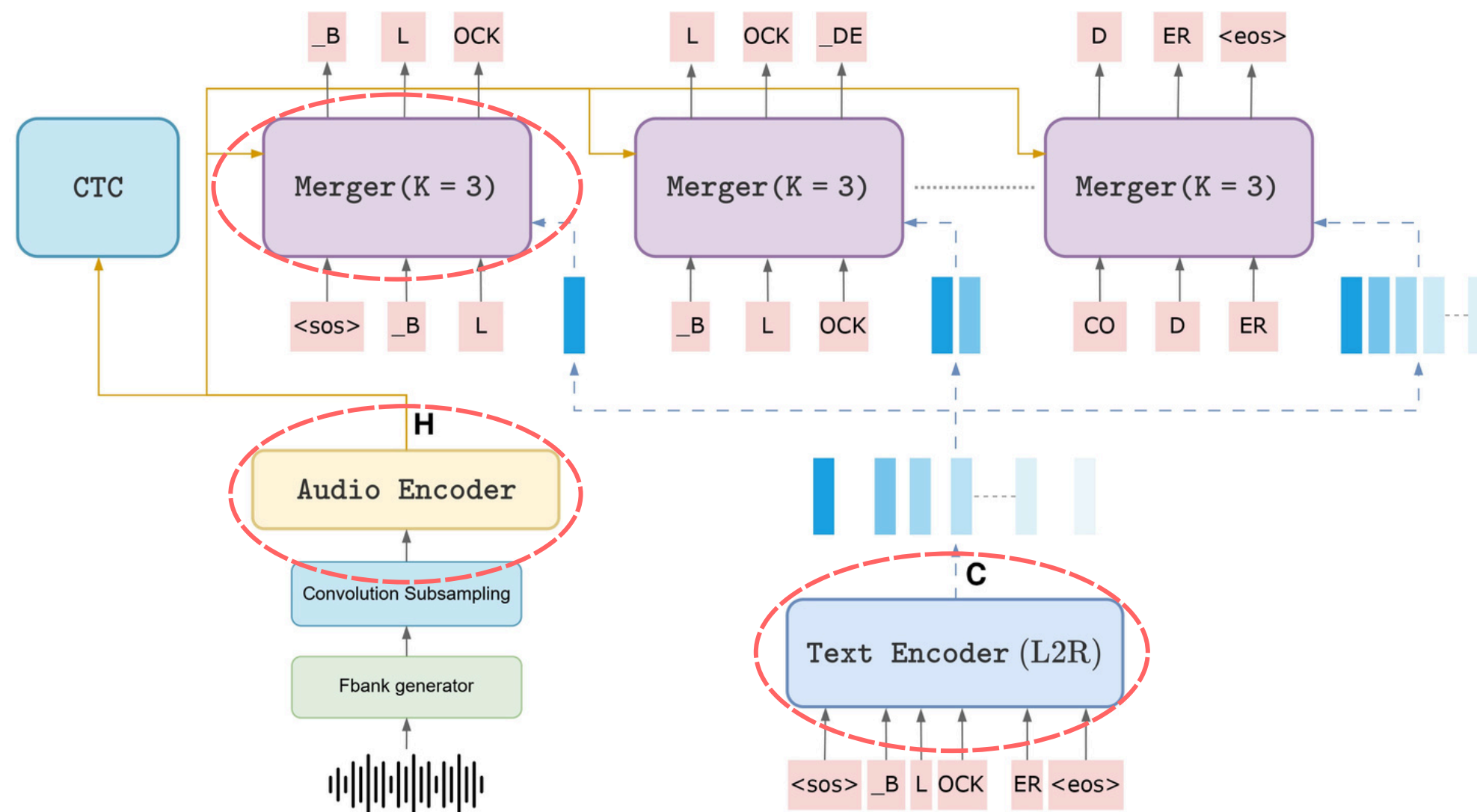
### Key Insights

- Replace the Traditional decoder with **Text Encoder** and **Merger**
- **Audio Encoder** → Builds rich audio context



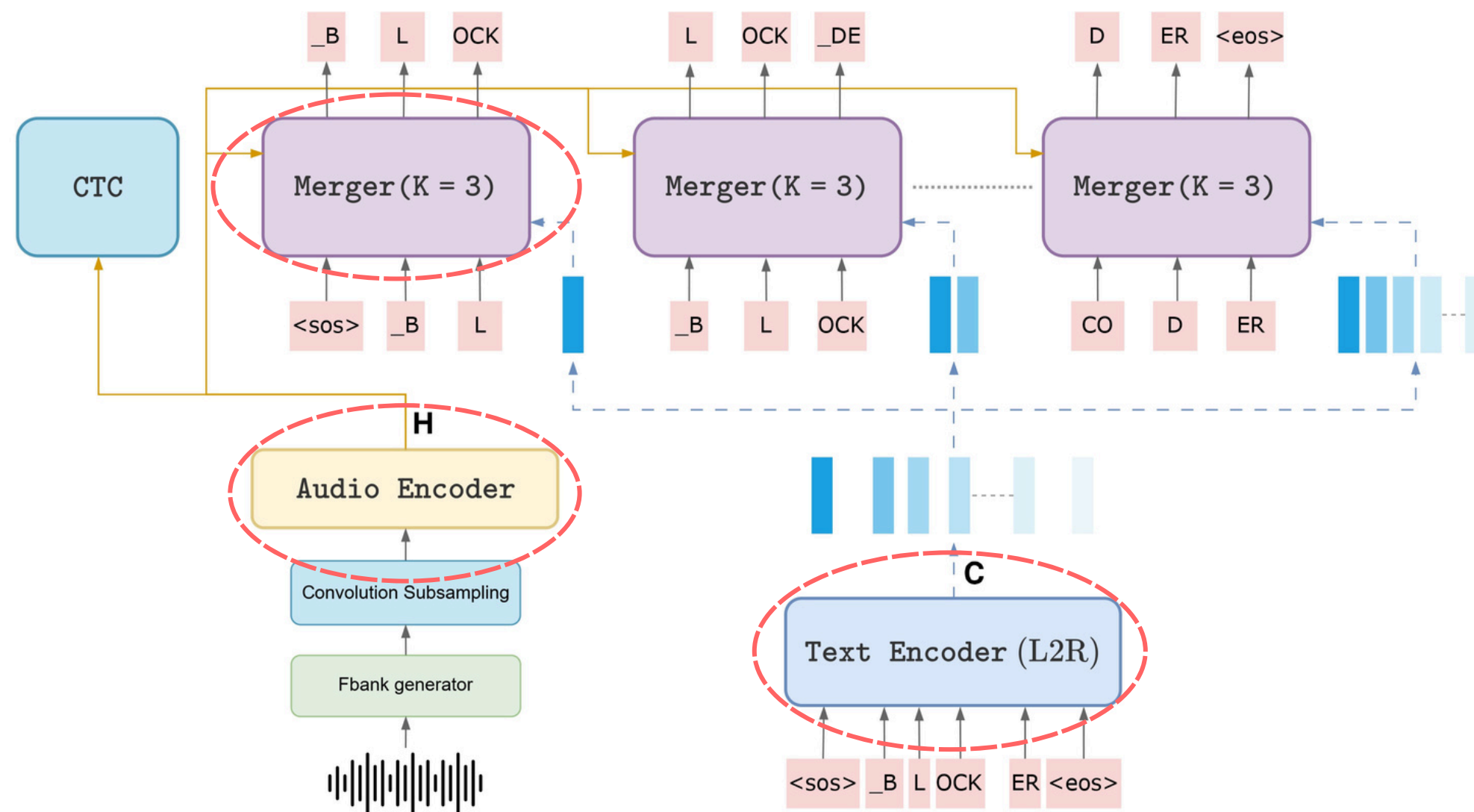
# BlockDecoder

Boosting ASR Decoders with Context and Merger Modules



## Key Insights

- Replace the Traditional decoder with **Text Encoder** and **Merger**
- **Audio Encoder** → Builds rich audio context
- **Text Encoder** → Builds rich text context, free from cross attention



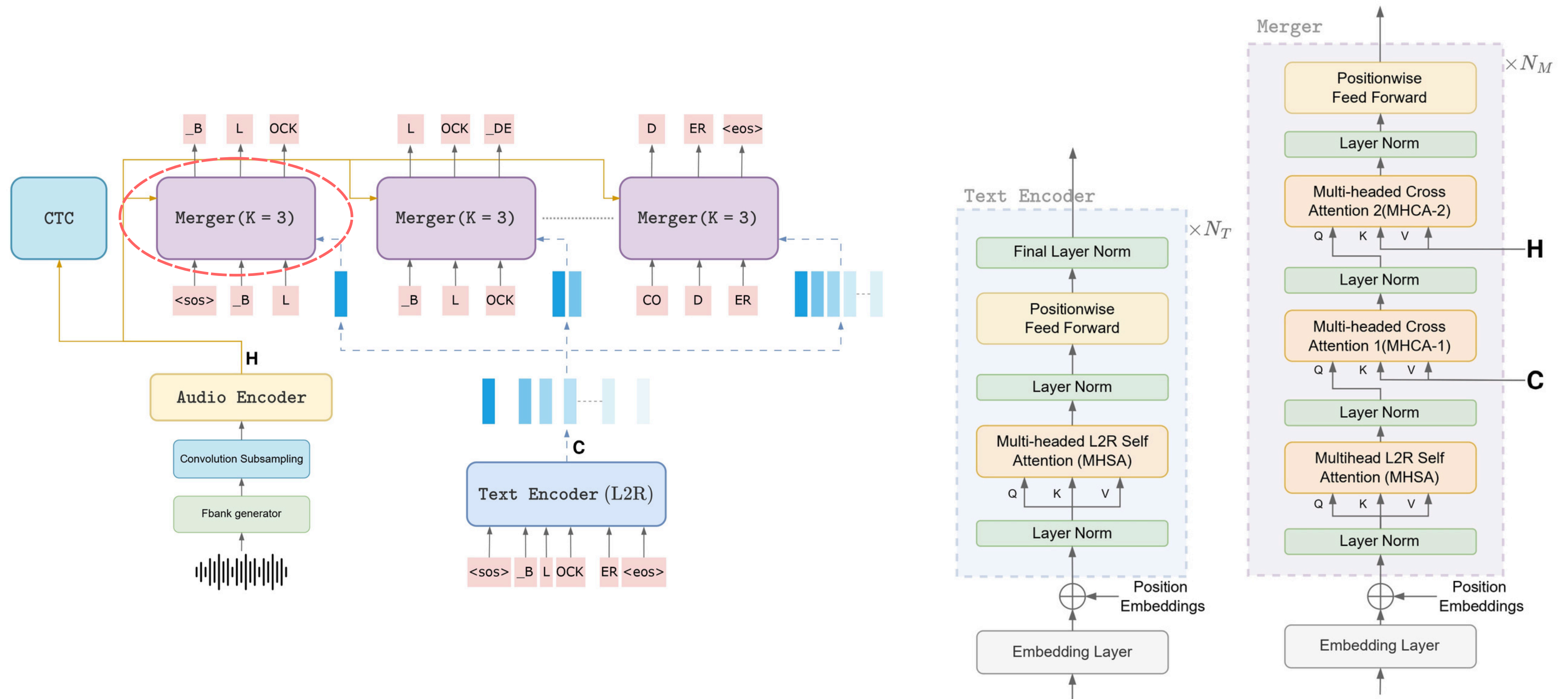
### Key Insights

- Replace the Traditional decoder with **Text Encoder** and **Merger**
- **Audio Encoder** → Builds rich audio context
- **Text Encoder** → Builds rich text context, free from cross attention
- **Merger** → Combines contexts and auto-regressively predicts a **block of K tokens**



# BlockDecoder

Boosting ASR Decoders with Context and Merger Modules

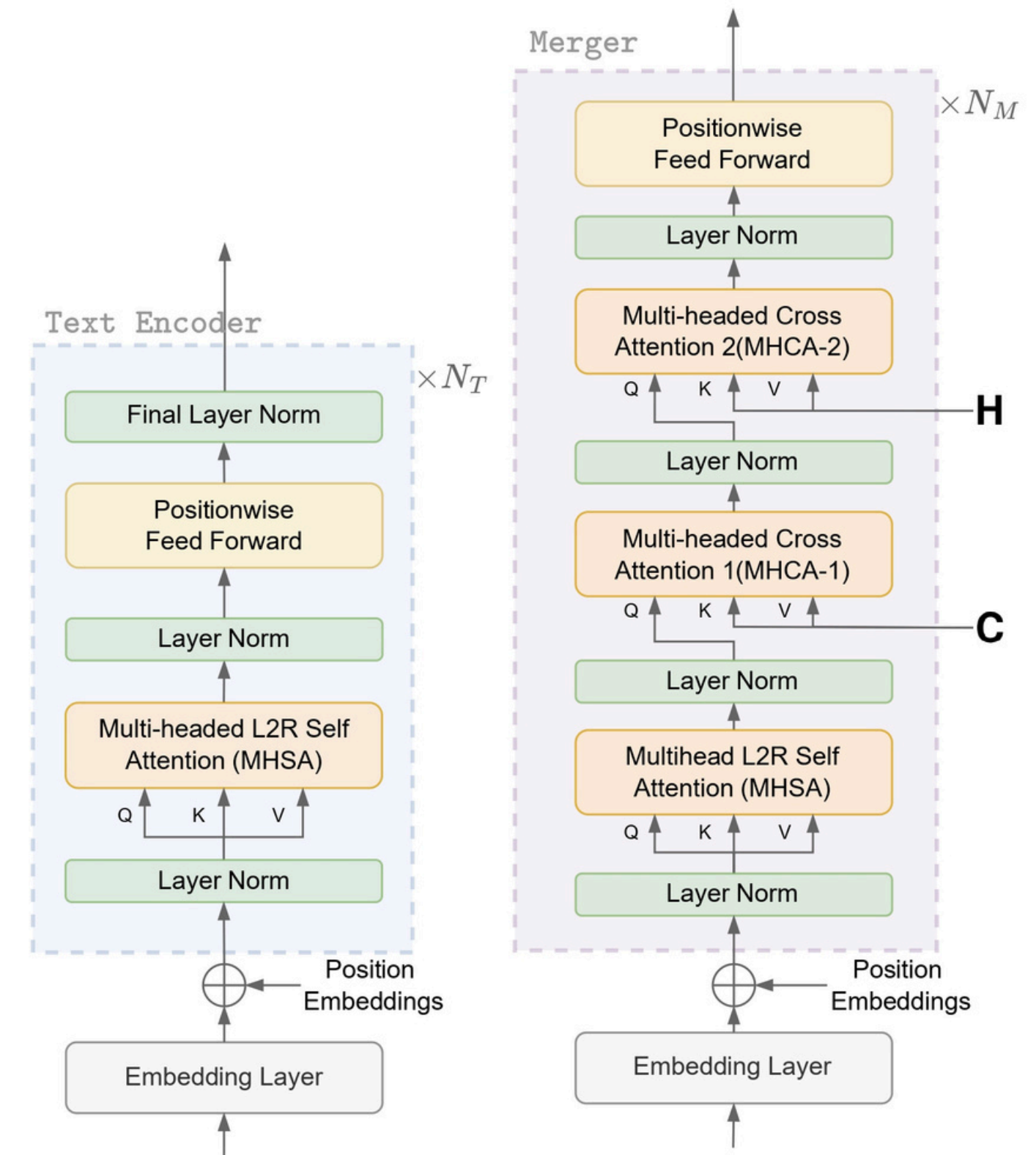
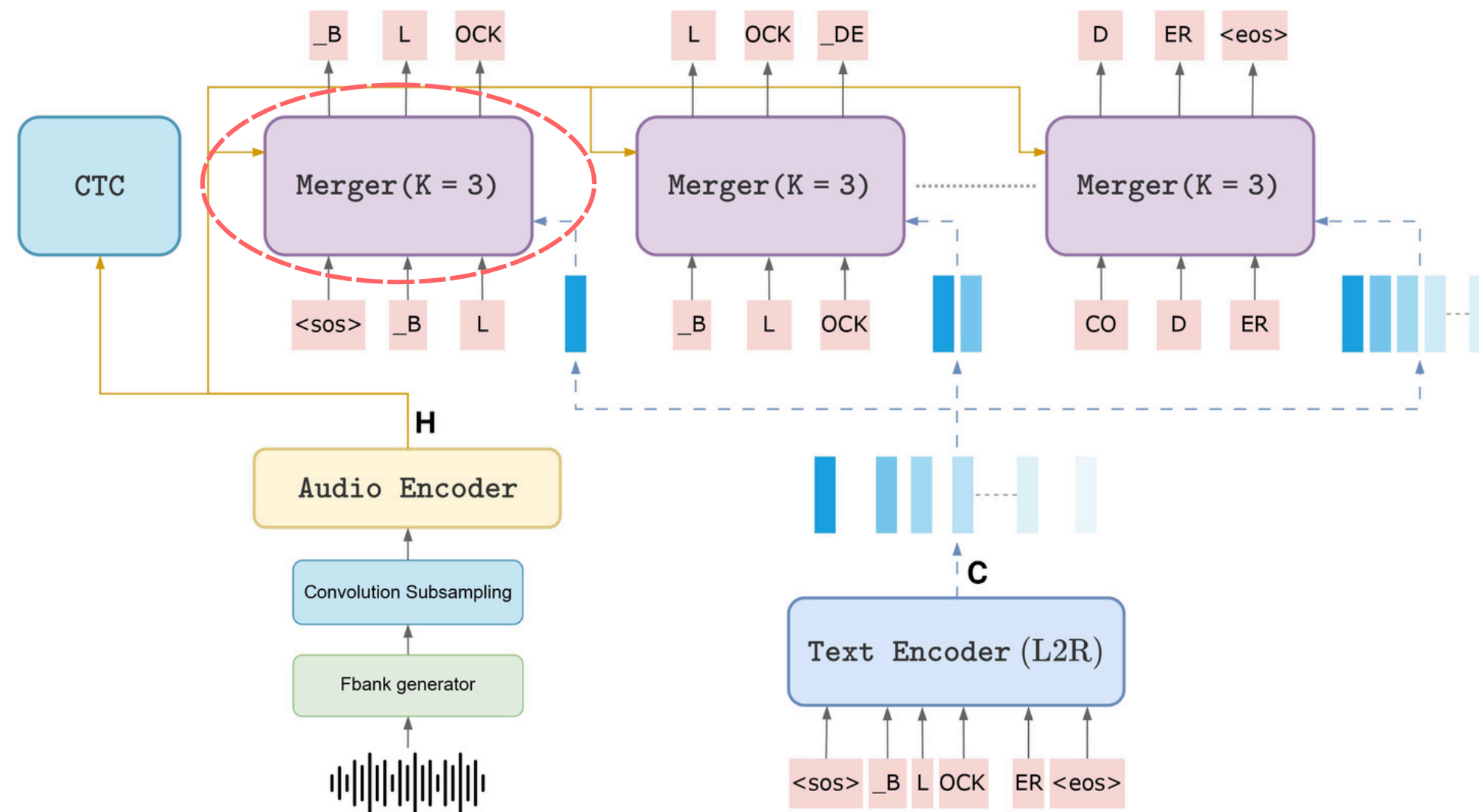






# BlockDecoder

Boosting ASR Decoders with Context and Merger Modules



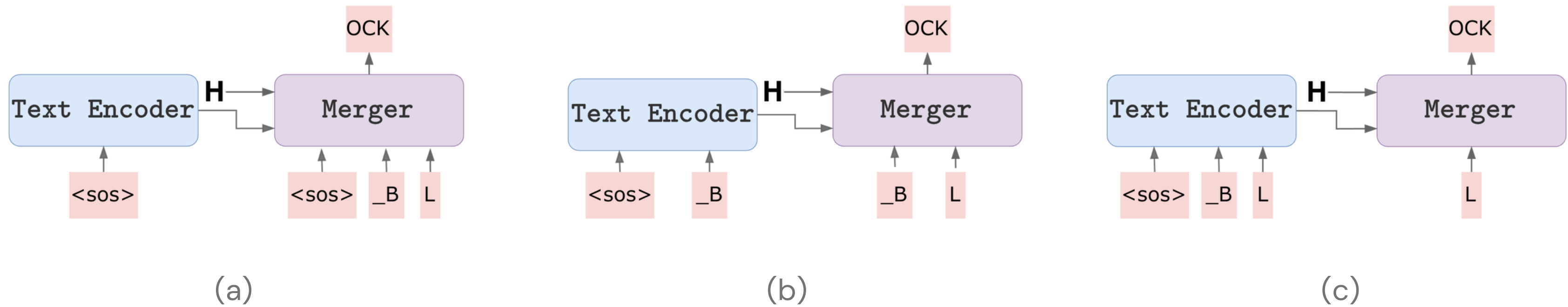
*With rich contexts, a very small number of merger layers are enough, thereby reducing parameters.*





# BlockDecoder

Boosting ASR Decoders with Context and Merger Modules

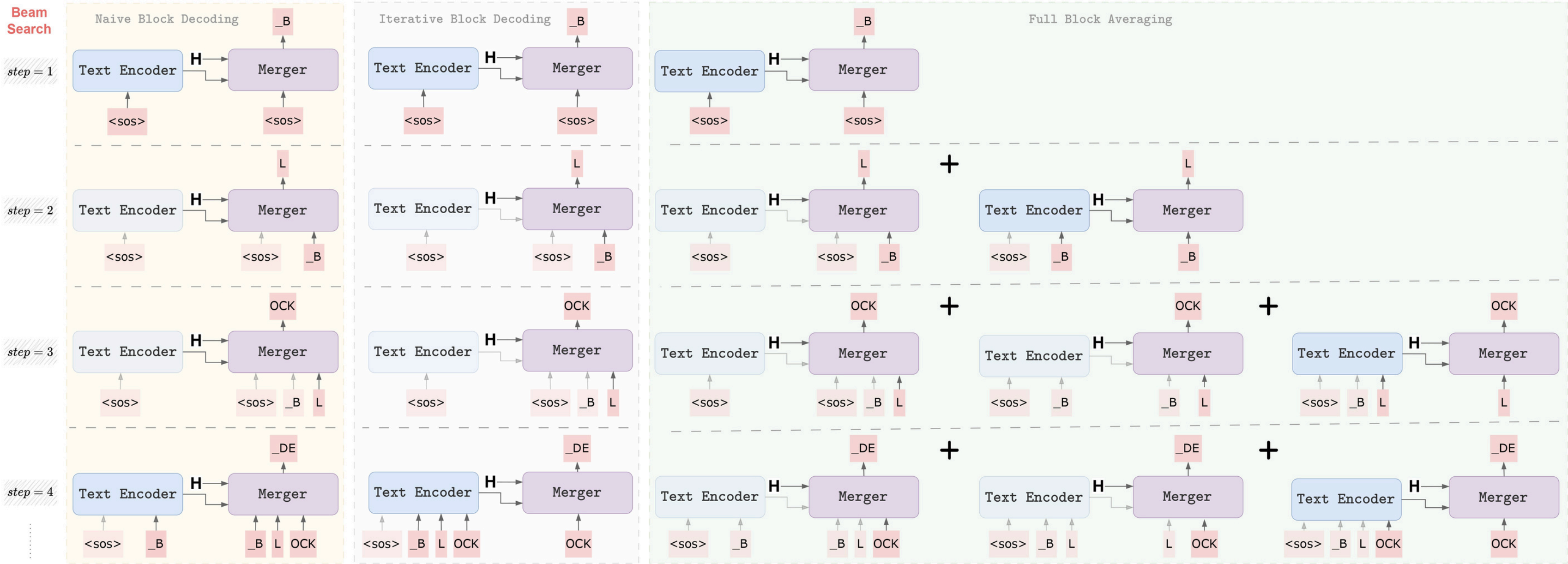


$$P(\text{OCK} / \langle \text{SOS} \rangle, \_B, L)$$



# BlockDecoder

Boosting ASR Decoders with Context and Merger Modules





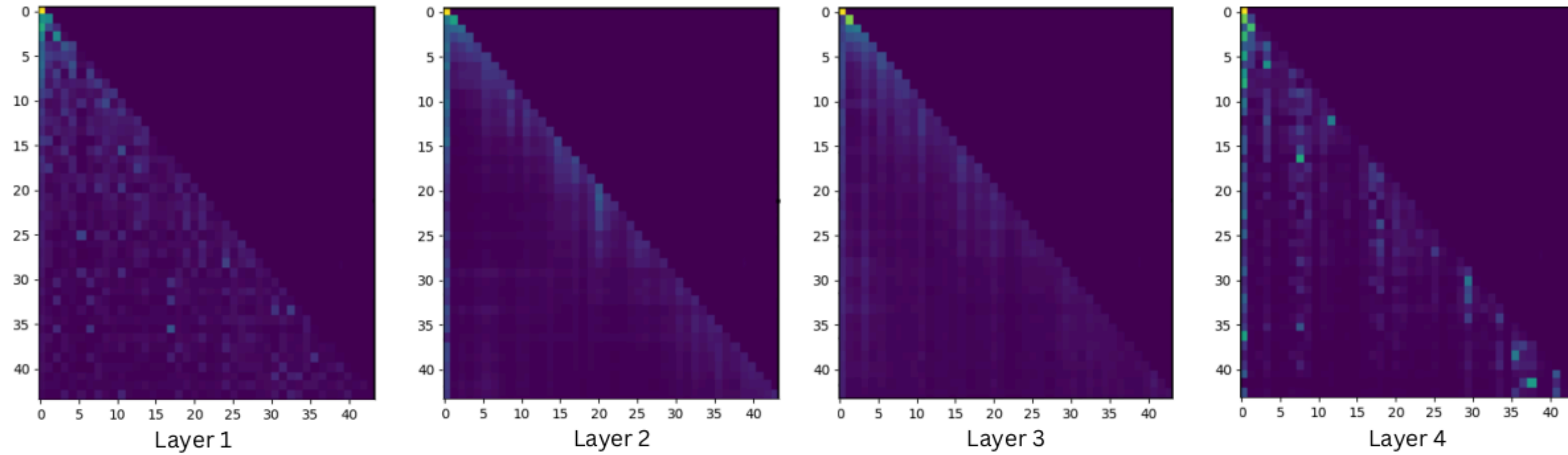
Boosting ASR Decoders with Context and Merger Modules

Method	Librispeech-100h (WER)				Tedlium2 (WER)			AISHELL (CER)		
	Params (M)	Test Clean ↓	Test Other ↓	RTF ↓	Params (M)	Test ↓	RTF ↓	Params (M)	Test ↓	RTF ↓
Conformer <a href="#">10</a>										
w/ Transformer Decoder	34.2	6.75	17.74	1.38	30.8	<b>7.69</b>	2.16	33.6	<b>4.58</b>	0.44
w/ BLOCKDECODER										
– Naive Block Decoding	33.7	6.63	17.63	0.73 <sub>(~1.9x)</sub>	30.2	7.81	1.02 <sub>(~2.1x)</sub>	33.1	4.76	0.28 <sub>(~1.6x)</sub>
– Iterative Block Decoding	33.7	6.72	17.70	<b>0.66</b> <sub>(~2.1x)</sub>	30.2	7.92	<b>0.90</b> <sub>(~2.4x)</sub>	33.1	4.75	<b>0.25</b> <sub>(~1.8x)</sub>
– Full Block Averaging	33.7	<b>6.58</b>	<b>17.62</b>	1.68 <sup>▷</sup>	30.2	7.79	2.39 <sup>▷</sup>	33.1	4.63	0.51 <sup>▷</sup>
E-Branchformer <a href="#">21</a>										
w/ Transformer Decoder	38.5	6.39	17.03	1.52	35.0	<b>7.44</b>	2.17	37.9	<b>4.50</b>	0.45
w/ BLOCKDECODER										
– Naive Block Decoding	37.9	6.15	<b>16.82</b>	0.77 <sub>(~2.0x)</sub>	34.5	7.61	1.04 <sub>(~2.1x)</sub>	37.4	4.61	0.30 <sub>(~1.5x)</sub>
– Iterative Block Decoding	37.9	6.19	16.94	<b>0.67</b> <sub>(~2.3x)</sub>	34.5	7.60	<b>0.91</b> <sub>(~2.4x)</sub>	37.4	4.61	<b>0.26</b> <sub>(~1.7x)</sub>
– Full Block Averaging	37.9	<b>6.14</b>	16.85	1.68 <sup>▷</sup>	34.5	7.61	2.40 <sup>▷</sup>	37.4	4.53	0.54 <sup>▷</sup>

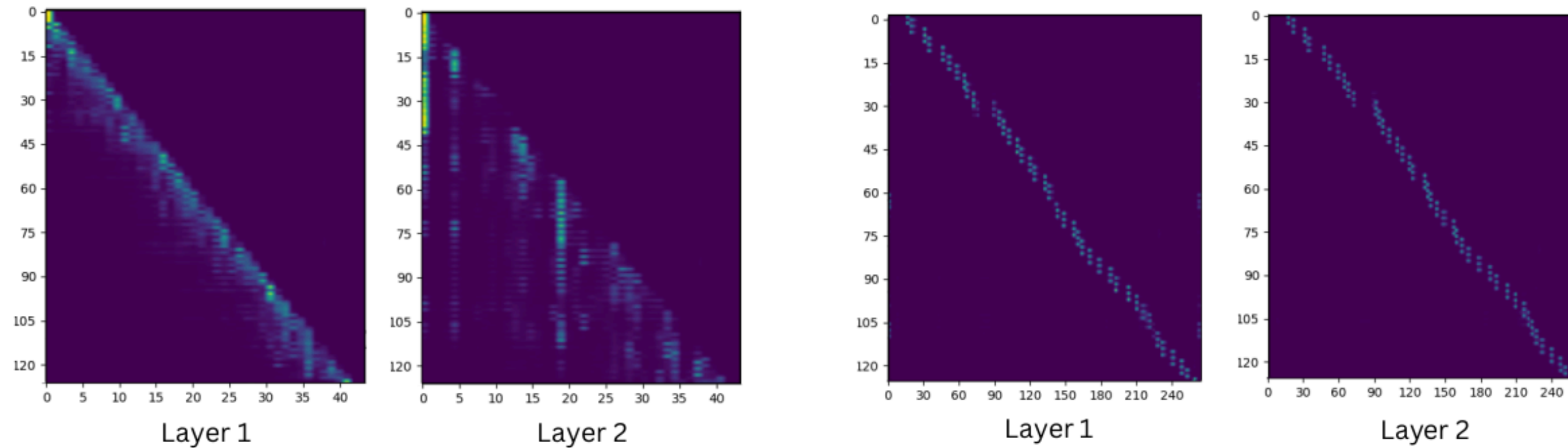


# BlockDecoder

Boosting ASR Decoders with Context and Merger Modules



(a)



RELEVANT LINKS



[https://openreview.net/forum?  
id=cGmcHJEFnY](https://openreview.net/forum?id=cGmcHJEFnY)



Email: [darshanp@cse.iitb.ac.in](mailto:darshanp@cse.iitb.ac.in)

# Want to know more?

DATE

06/11/2025

## POSTER SESSION

**Exhibit Hall C,D,E**  
**Wed 3 Dec 11 a.m. – 2 p.m. PST**

If you are interested, please stop by. 😊





RELEVANT LINKS



DATE

06/11/2025

[https://openreview.net/forum?  
id=cGmcHJEFnY](https://openreview.net/forum?id=cGmcHJEFnY)



Email: [darshanp@cse.iitb.ac.in](mailto:darshanp@cse.iitb.ac.in)

# Thank you!



Darshan Prabhu



Preethi Jyothi

