

# Byte Latent Transformer: Patches Scale Better Than Tokens

Pagnoni et. al. from FAIR at Meta (Unreviewed Preprint)

Presenter: Marian Schneider  
Advisor: Frédéric Berdoz



# Demo



# Background

How do (large) language models work?

# Text as a Sequence

## Patches Scale Better Than Tokens

Characters

P a t c h e s \_ S c a l e \_ B e t t e r ...  
🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡 🟡

- Flexible
- Long sequence

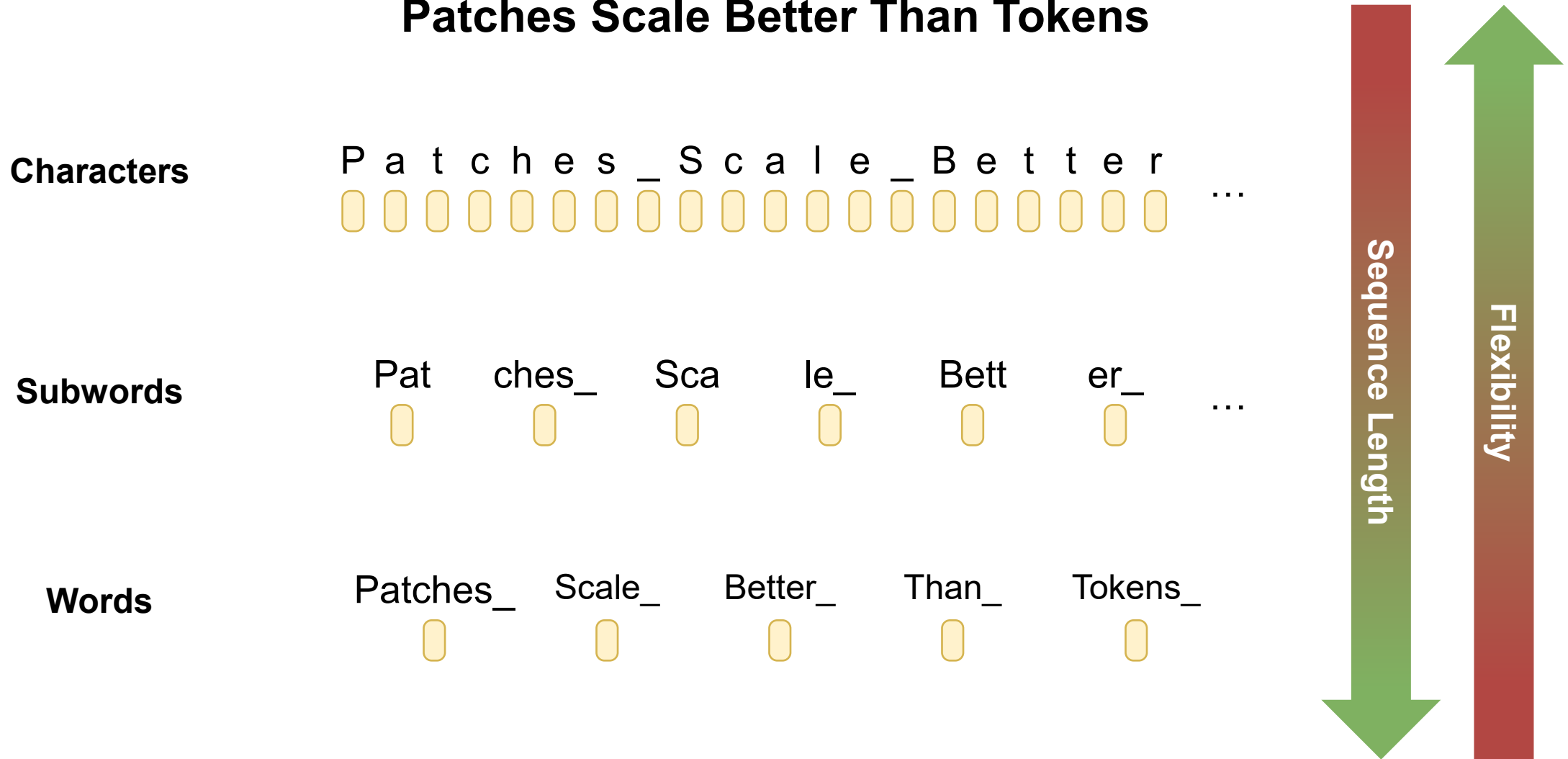
Words

Patches\_ Scale\_ Better\_ Than\_ Tokens\_  
🟡 🟡 🟡 🟡 🟡

- Short sequence
- Inflexible

# Text as a Sequence

## Patches Scale Better Than Tokens

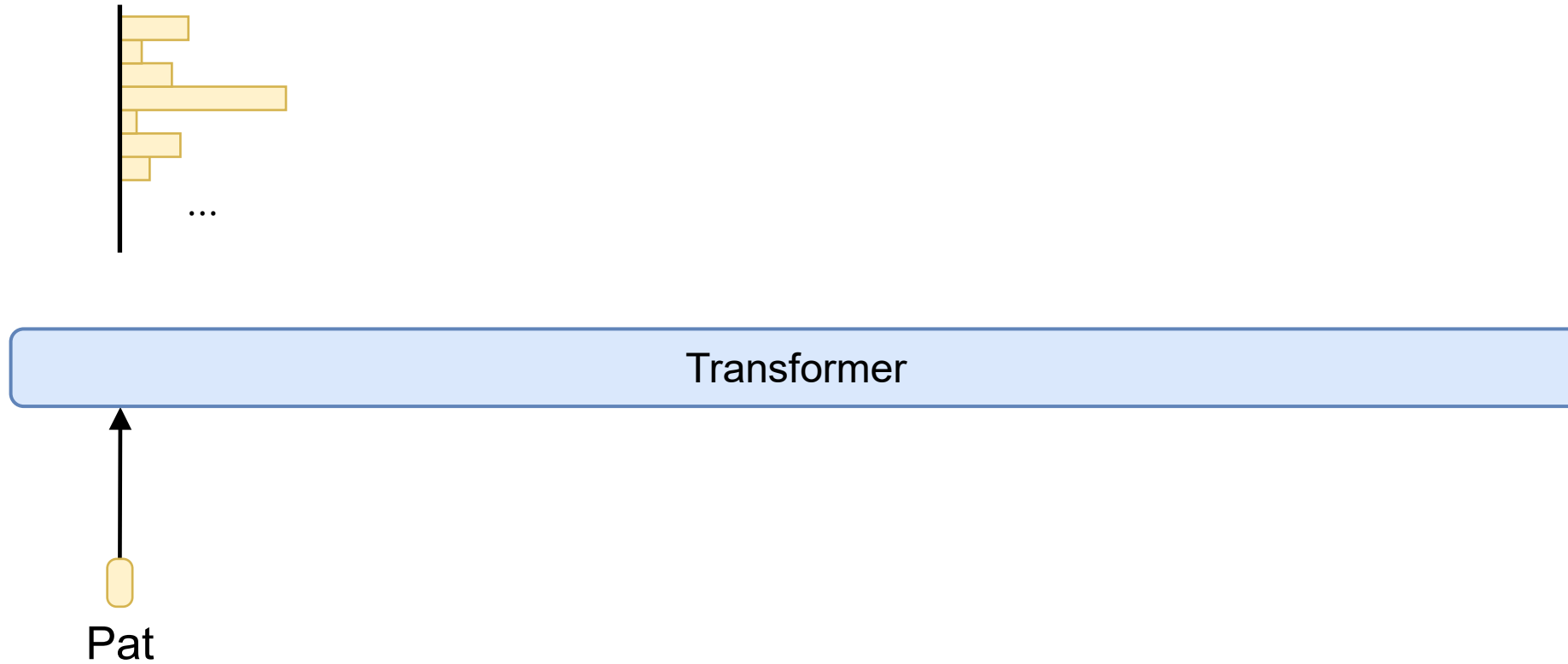


# Autoregressive Generation

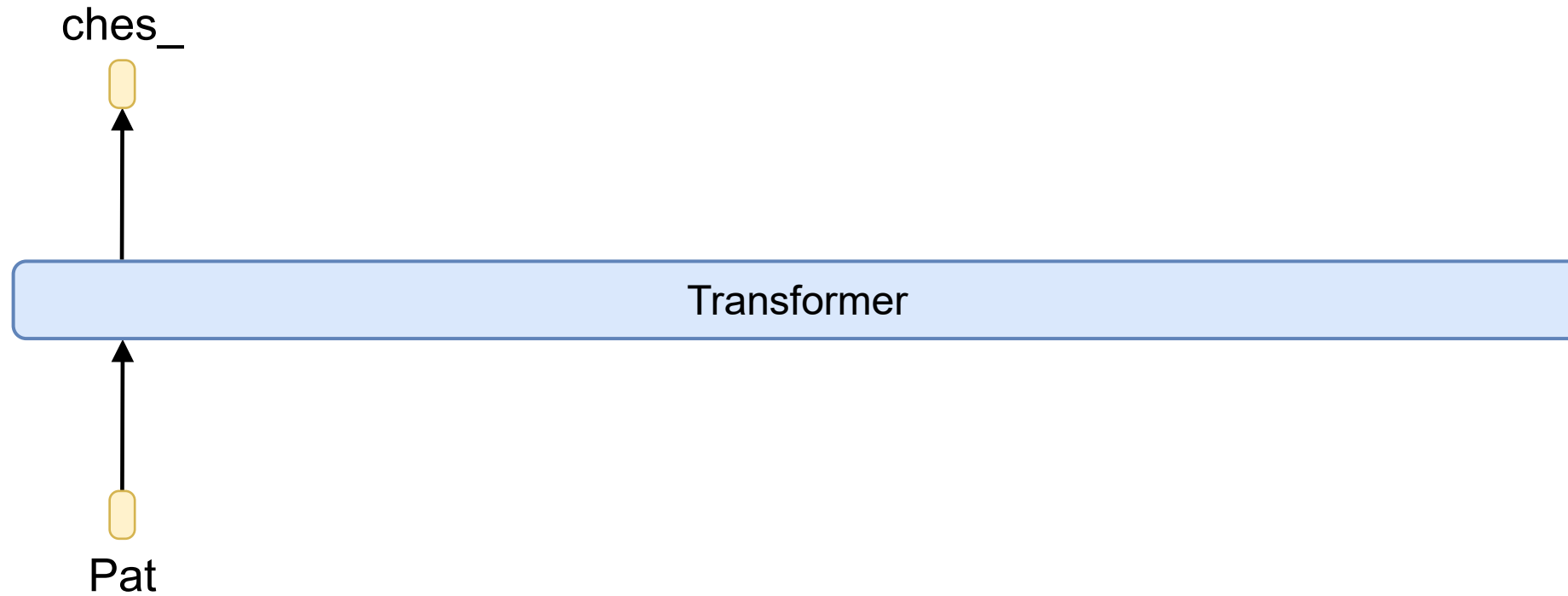


Transformer

# Autoregressive Generation

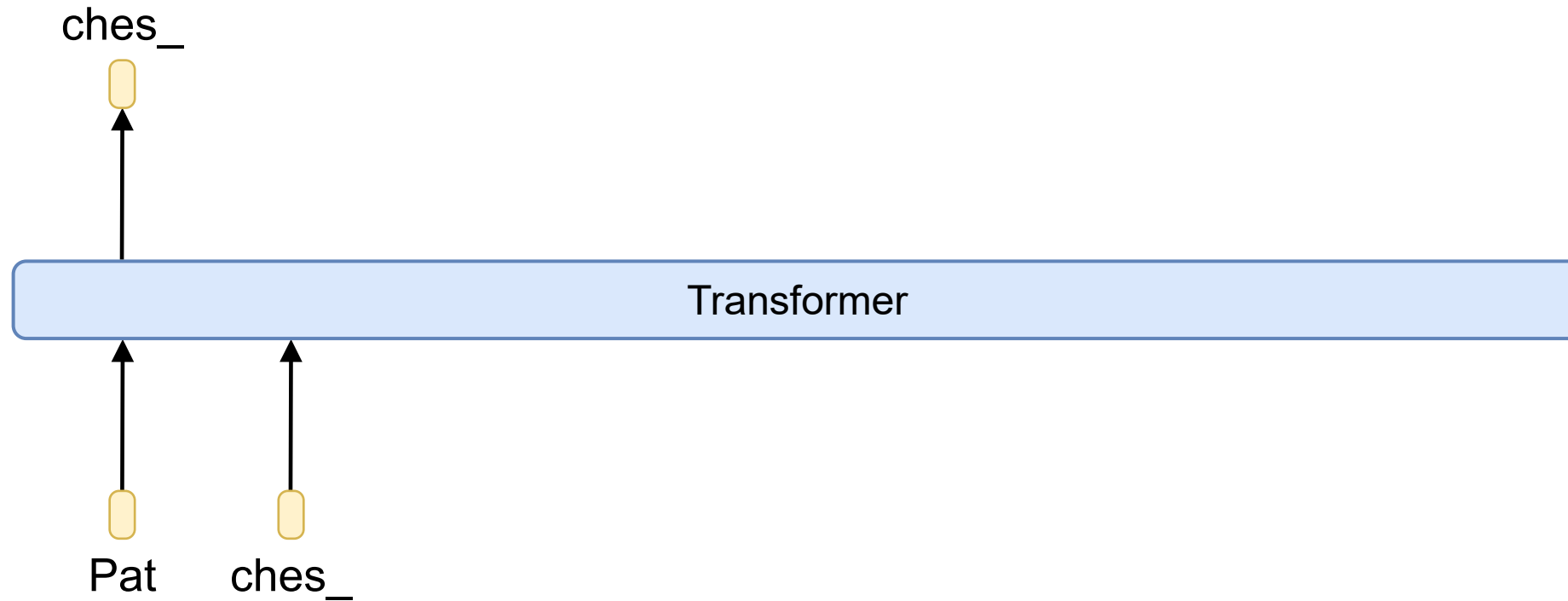


# Autoregressive Generation

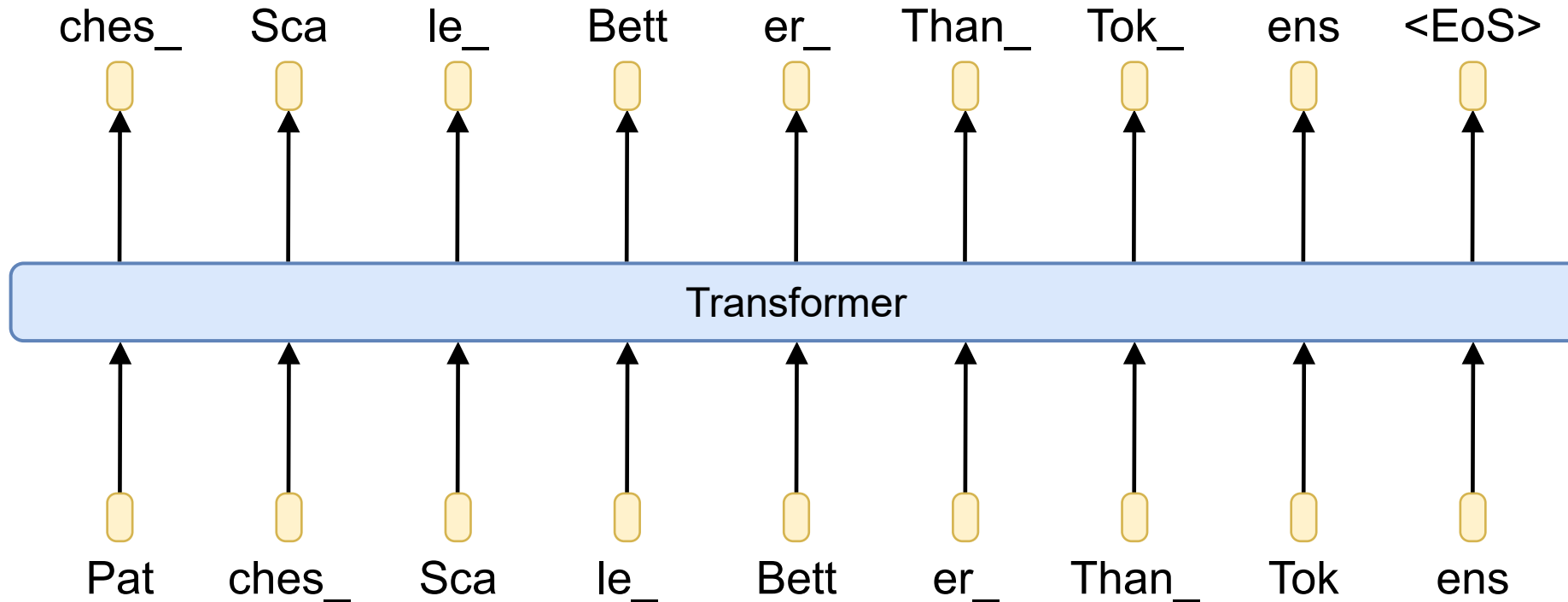




# Autoregressive Generation



# Autoregressive Generation

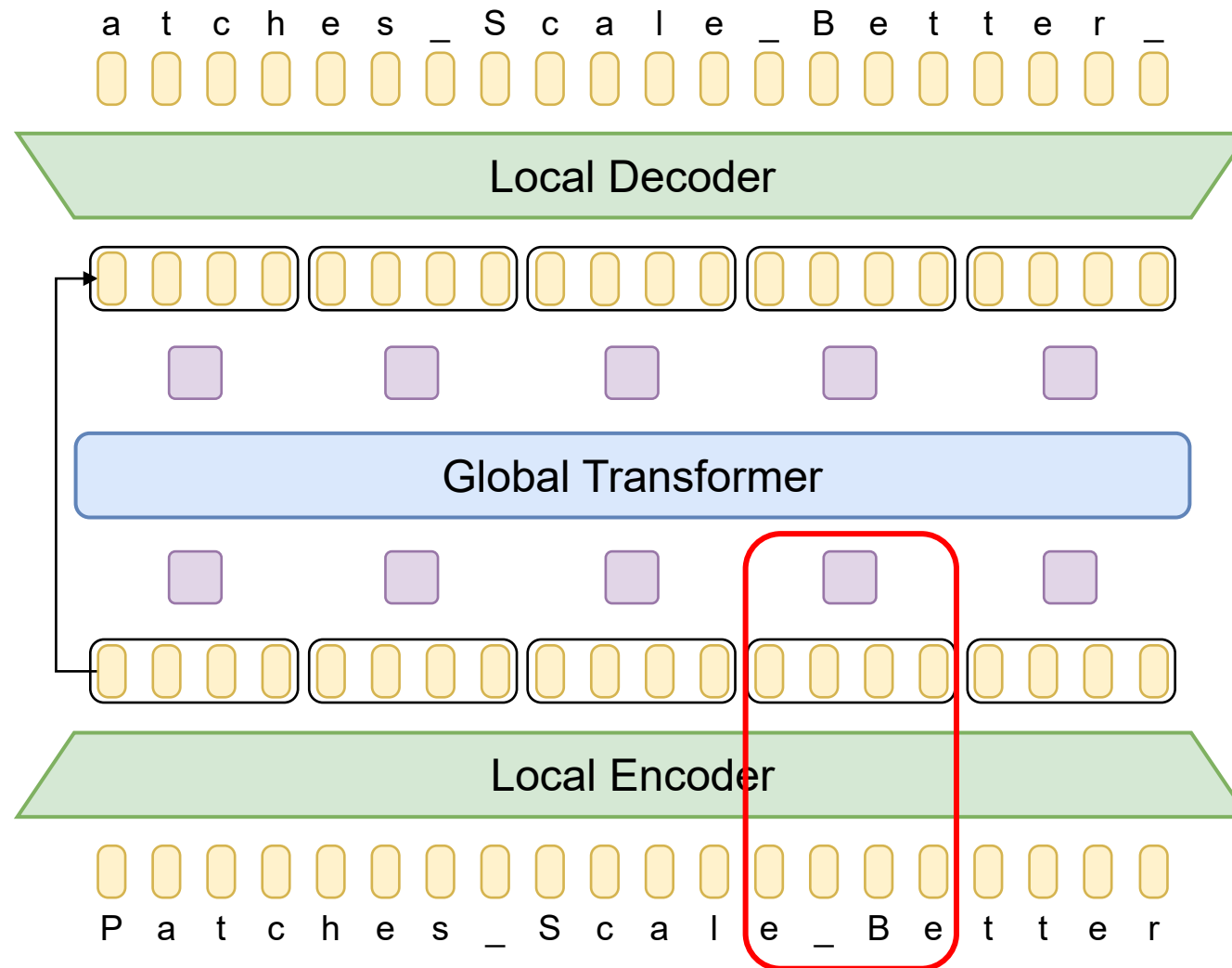


# Previous Work

How can we overcome the limitations of tokenization?

# MegaByte

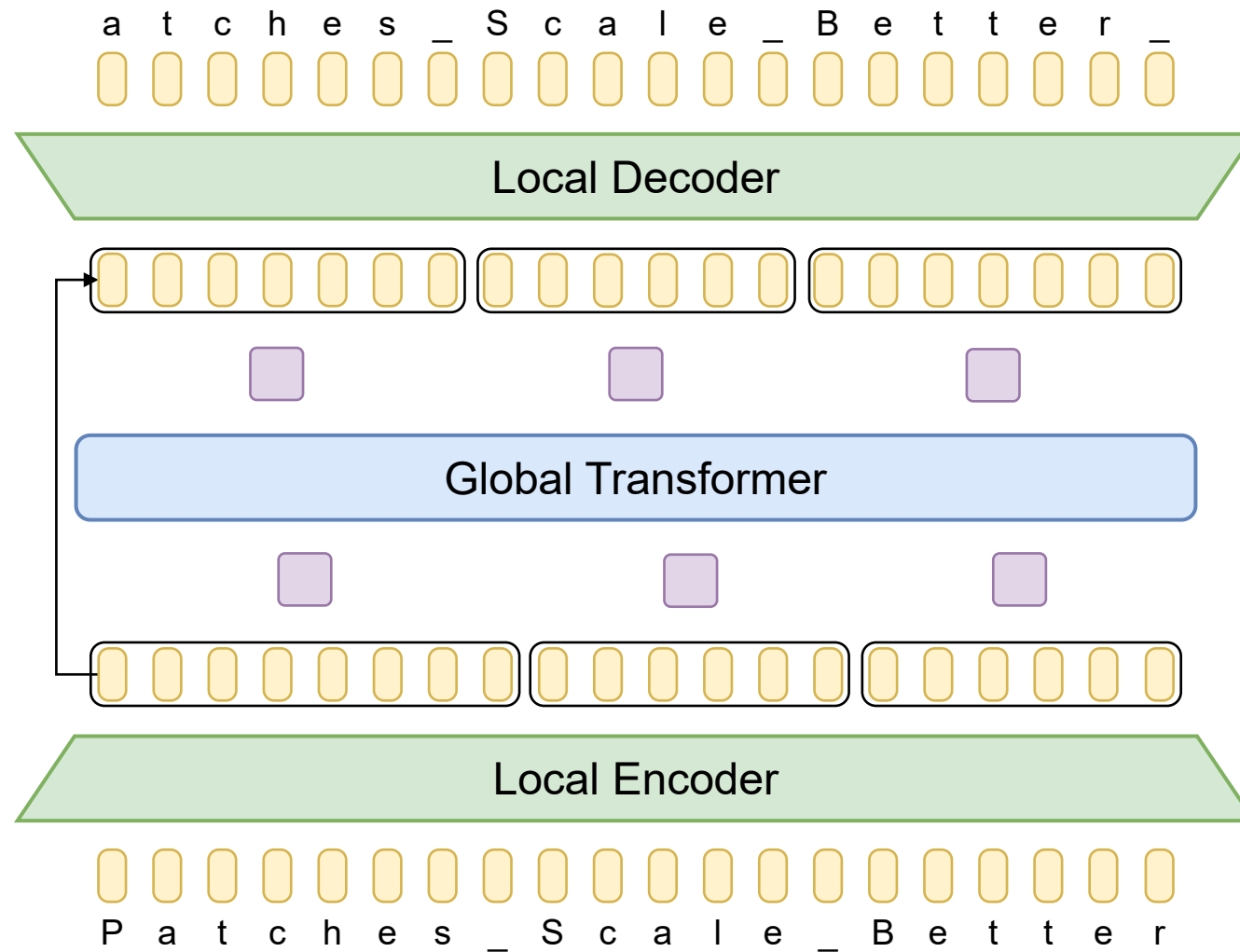
$k$  bytes per patch



- Inconsistent patching

# SpaceByte

1 word per patch



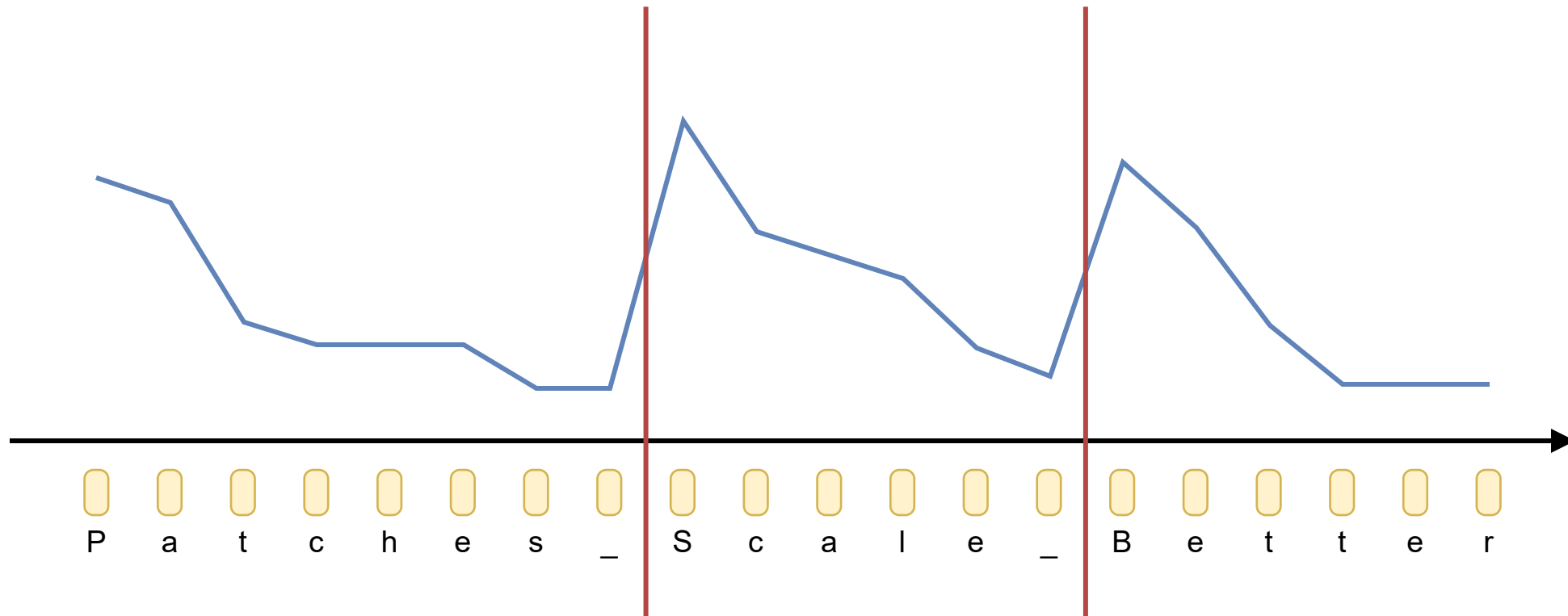
- Inflexible patching

# Dynamic Patching

No fixed patch size

Seminar in Deep → Neural Networks

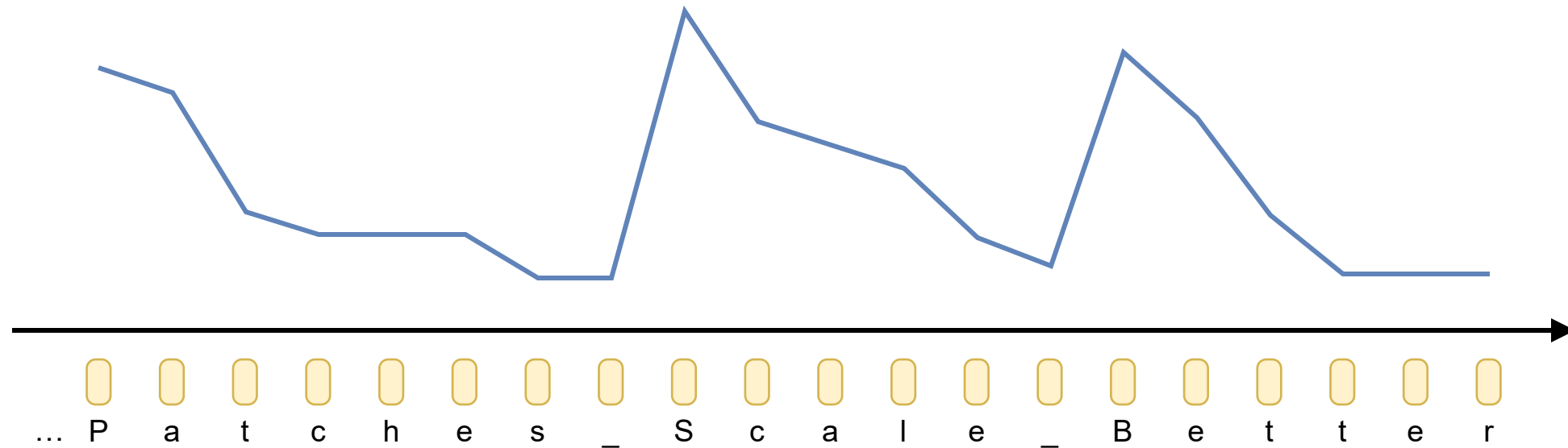
This is →  
→ a good presentation  
→ a beautiful day  
→ the guy she told you not to worry about



# Byte Latent Transformer

What is the contribution of the paper?

# BLT Dynamic Patching



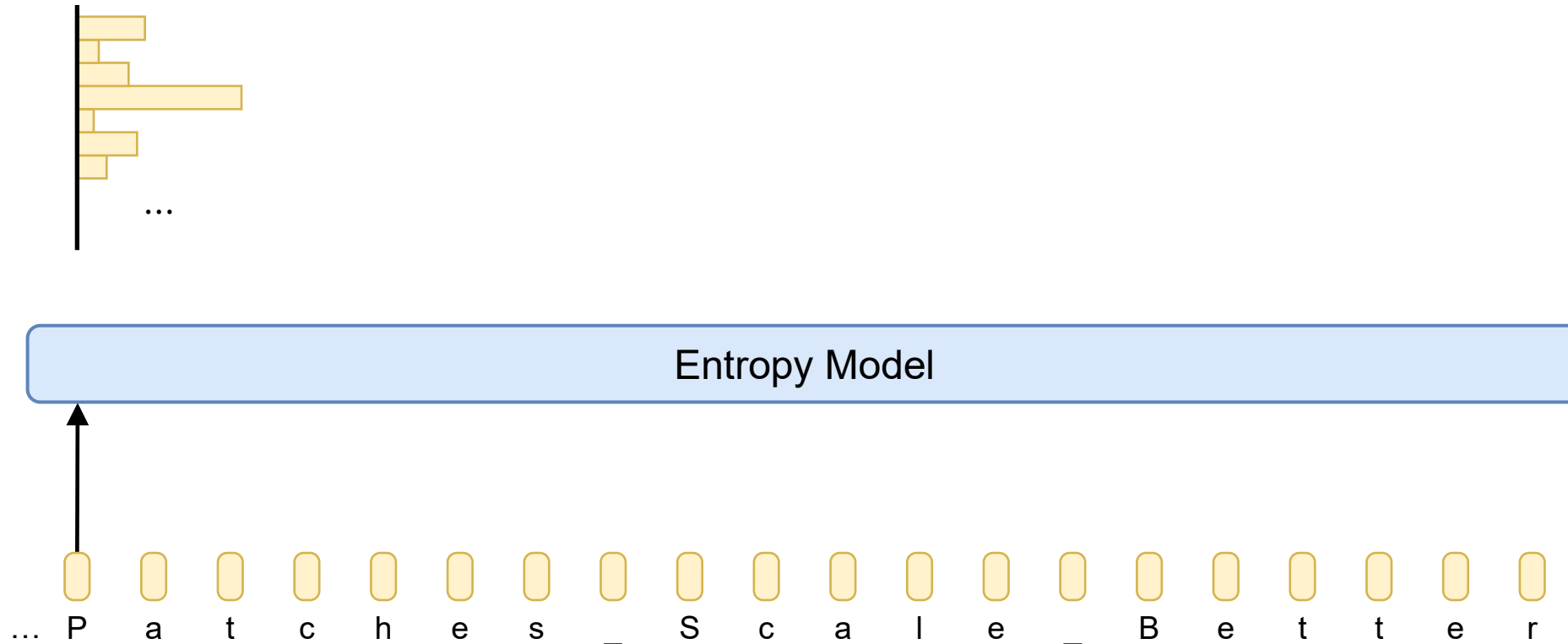


# BLT Dynamic Patching

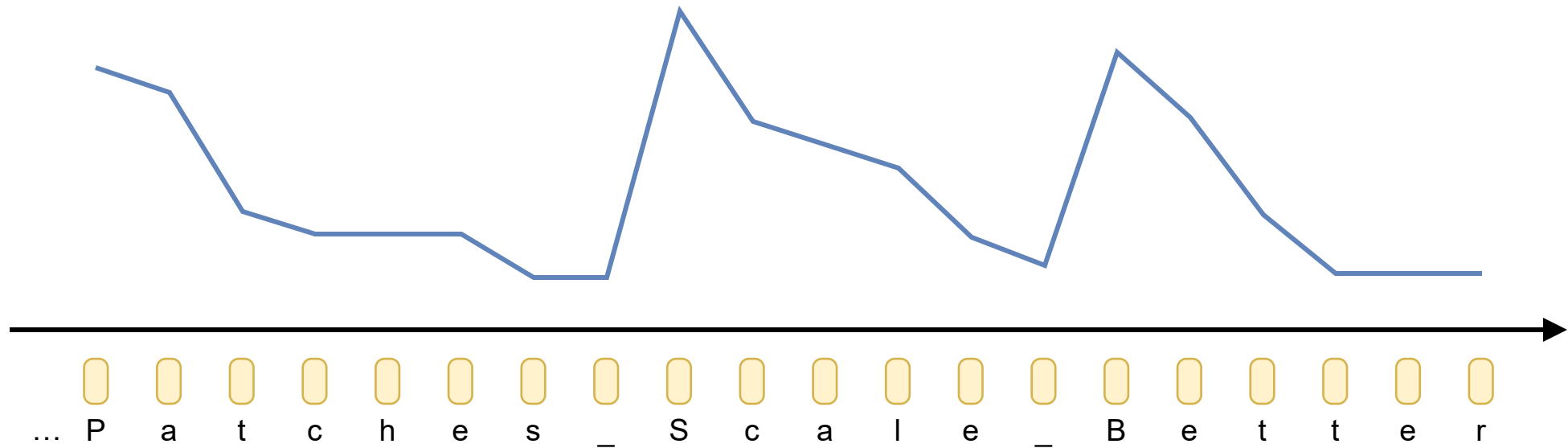
Entropy Model

... P a t c h e s \_ S c a l e \_ B e t t e r

# BLT Dynamic Patching



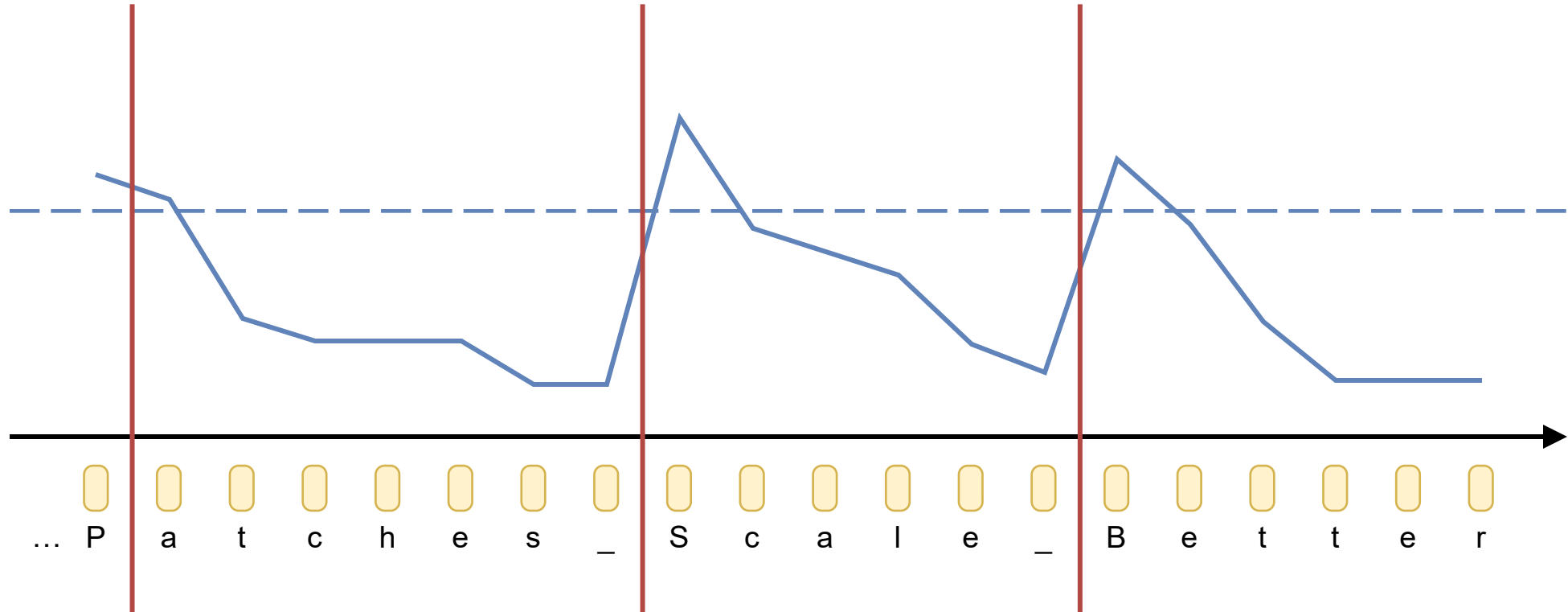
# BLT Dynamic Patching



# BLT Dynamic Patching

## Global Constraint

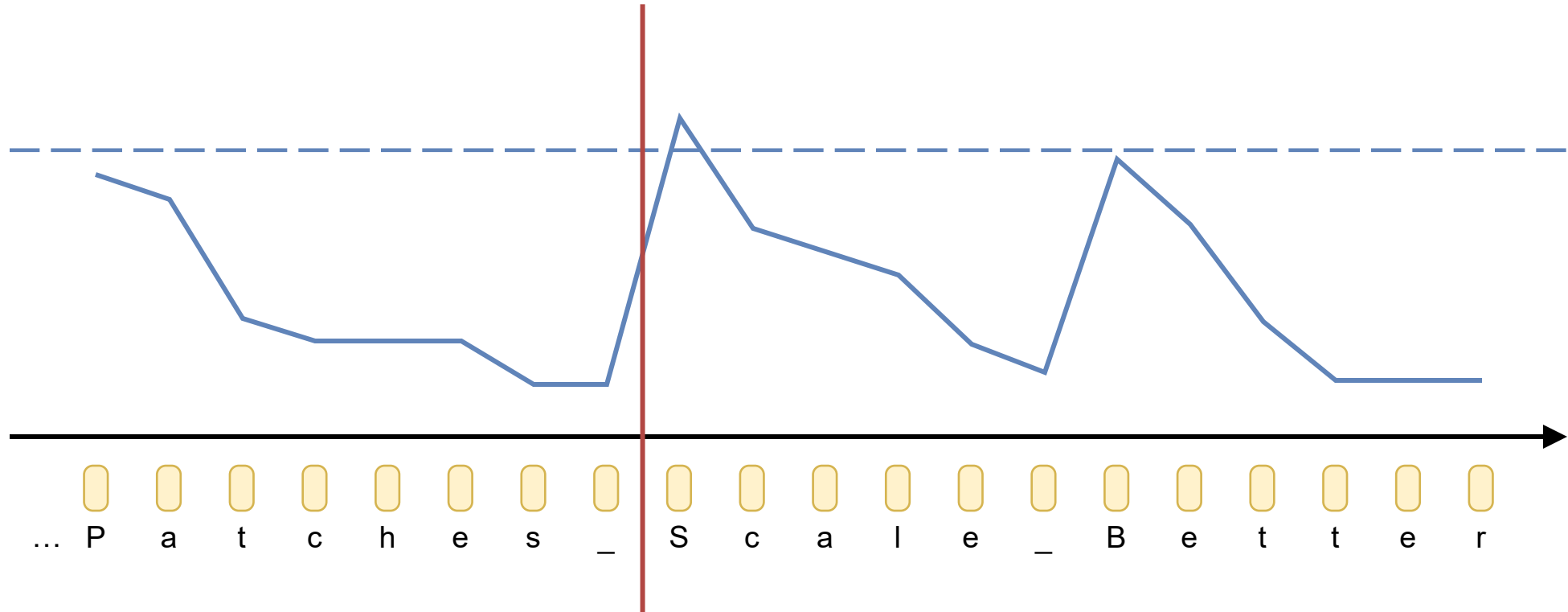
$$H(x_t) > \theta_g$$



# BLT Dynamic Patching

## Global Constraint

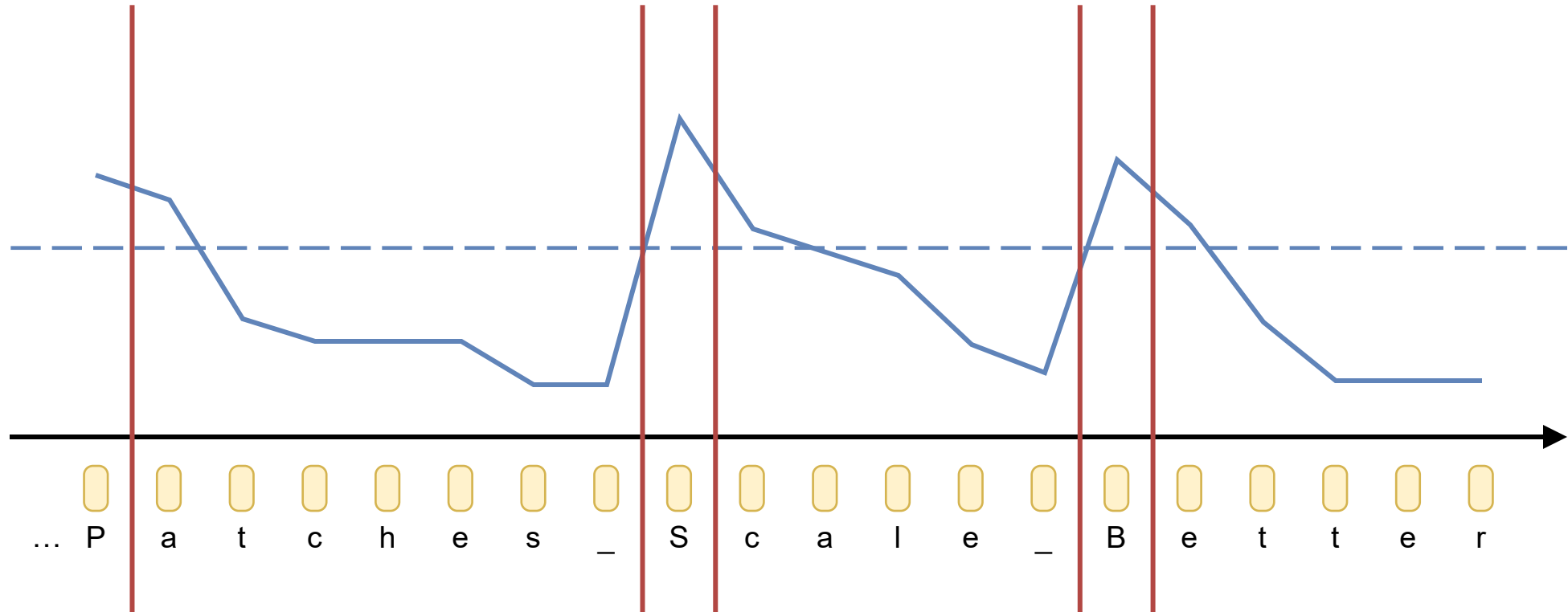
$$H(x_t) > \theta_g$$



# BLT Dynamic Patching

## Global Constraint

$$H(x_t) > \theta_g$$



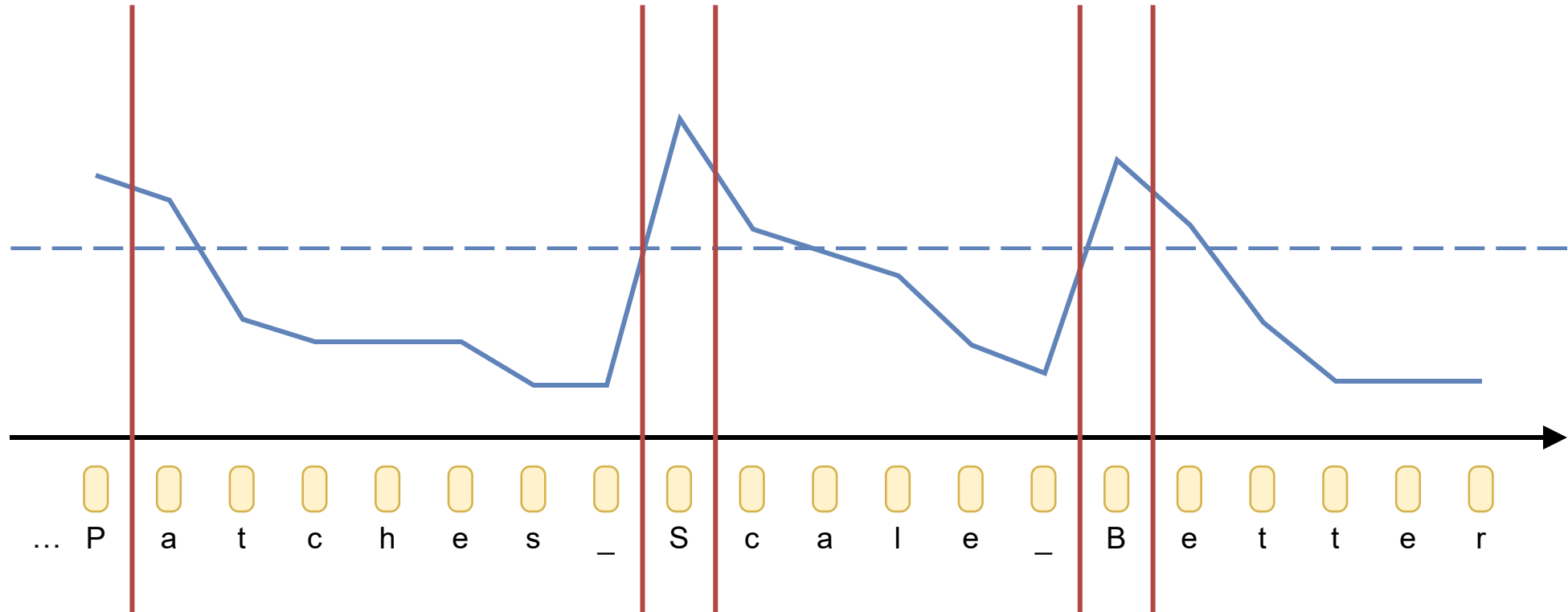
# BLT Dynamic Patching

## Global Constraint

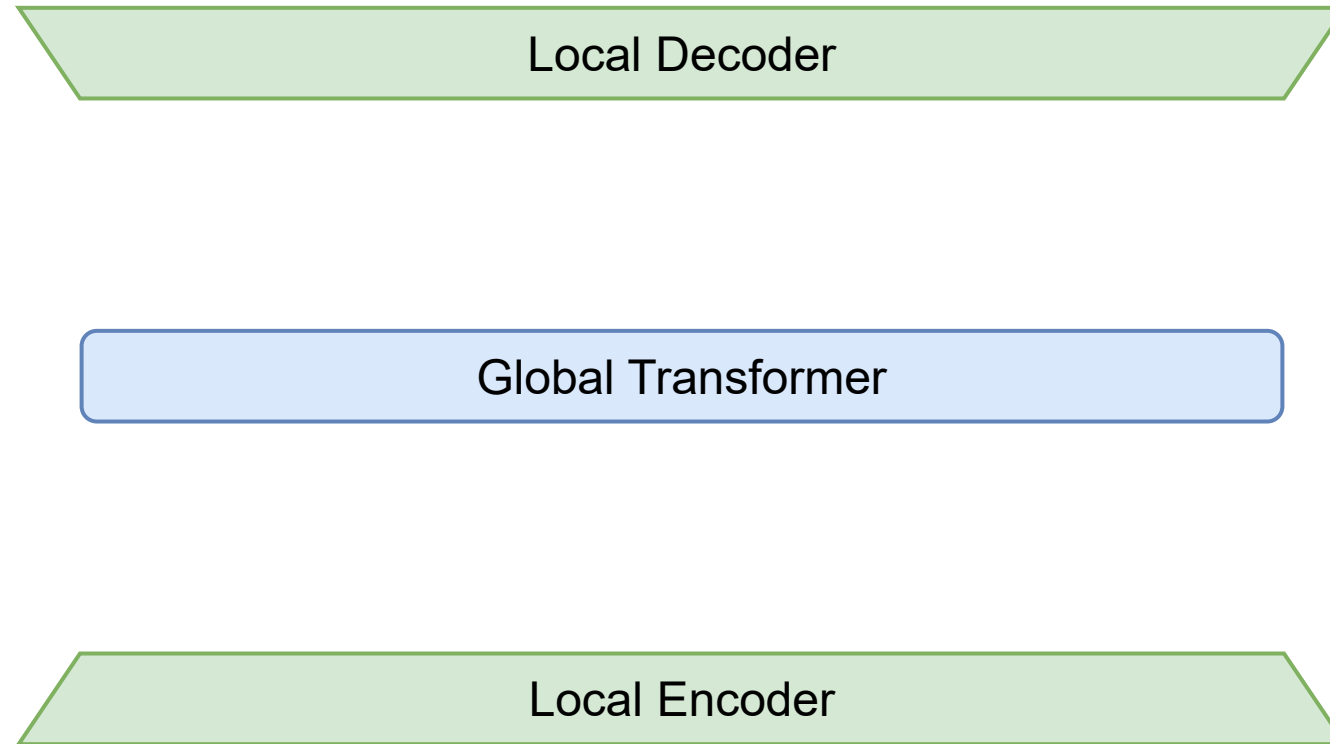
$$H(x_t) > \theta_g$$

## Approximate Monotonic Constraint

$$H(x_t) - H(x_{t-1}) > \theta_r$$

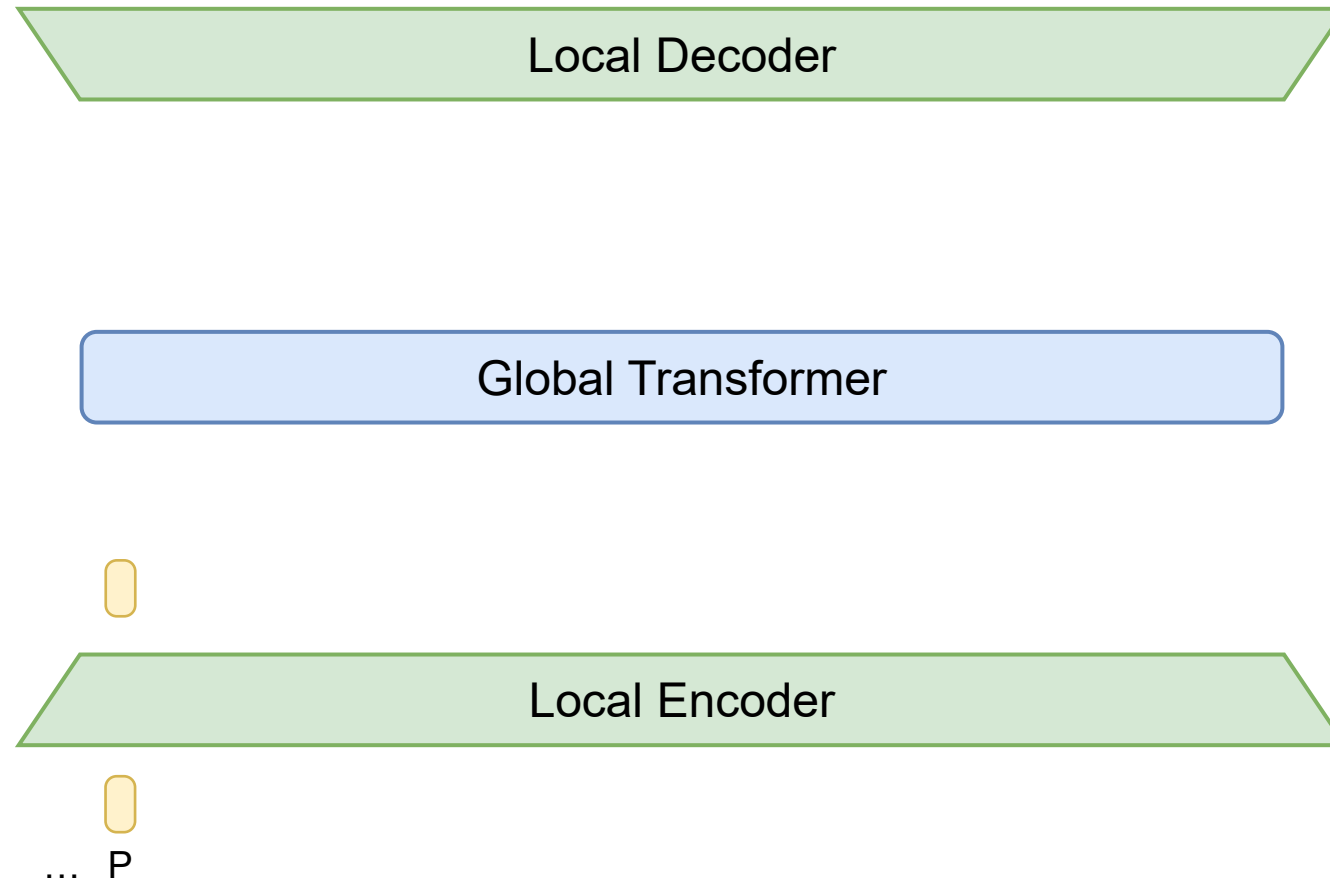


# BLT Generation

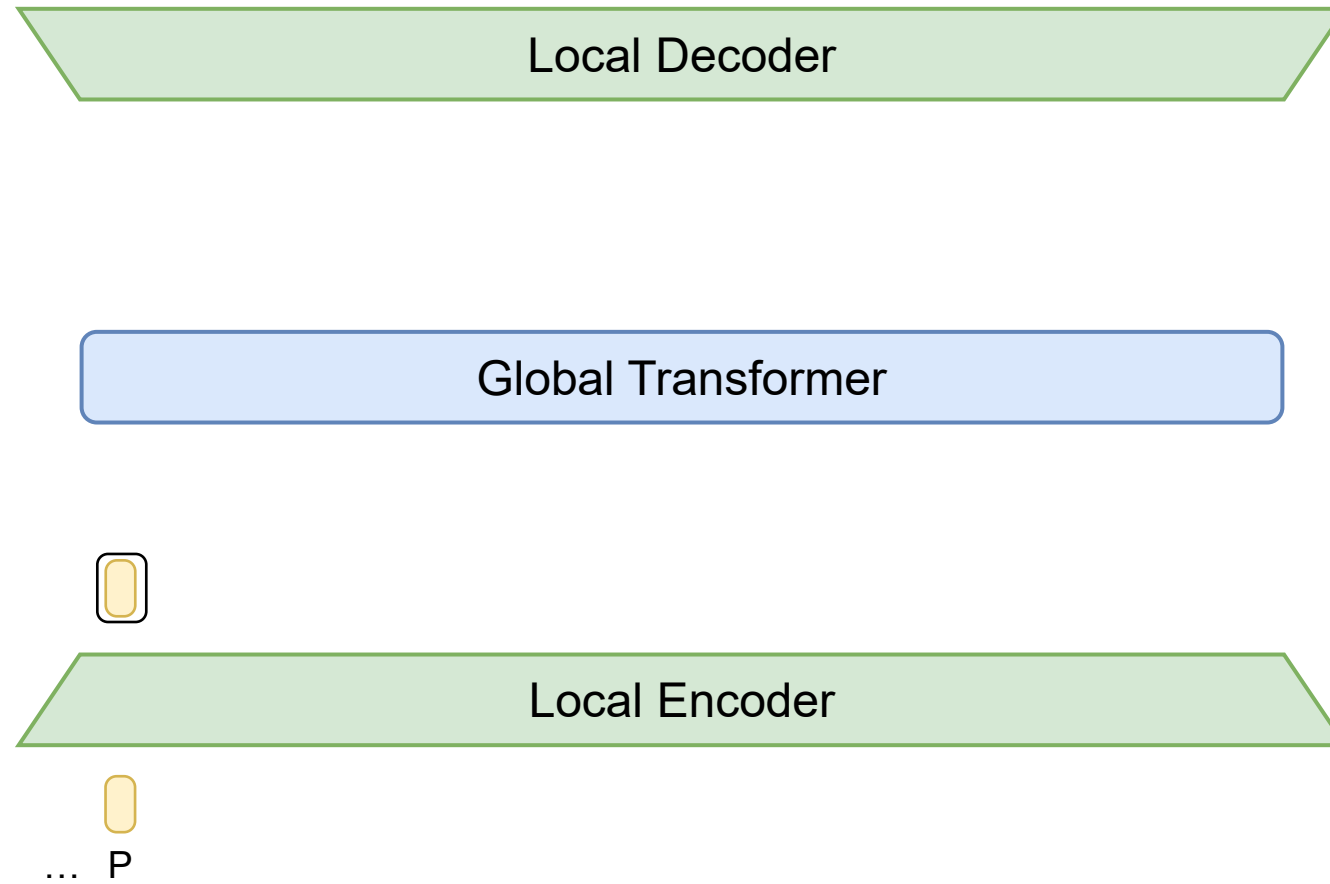




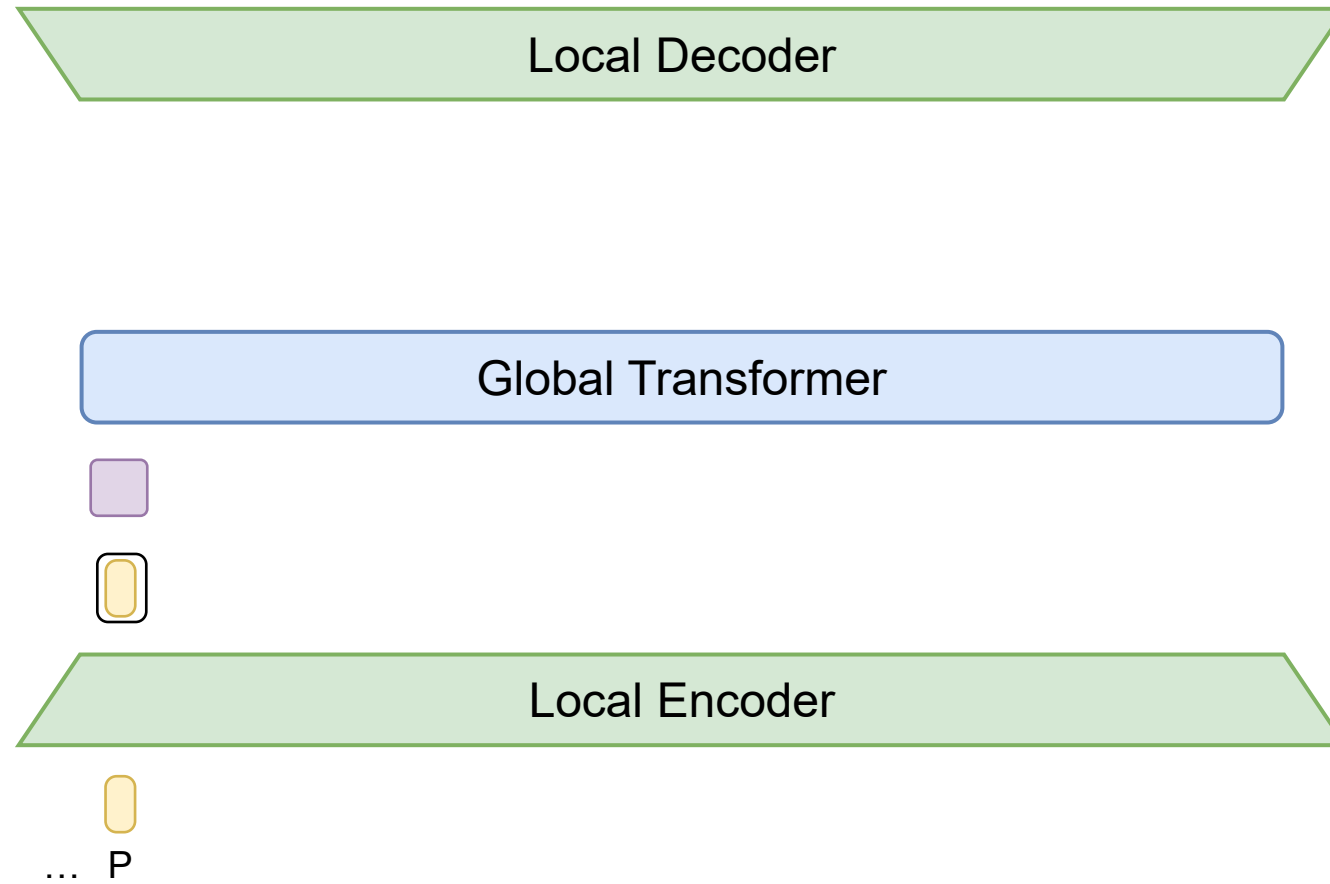
# BLT Generation



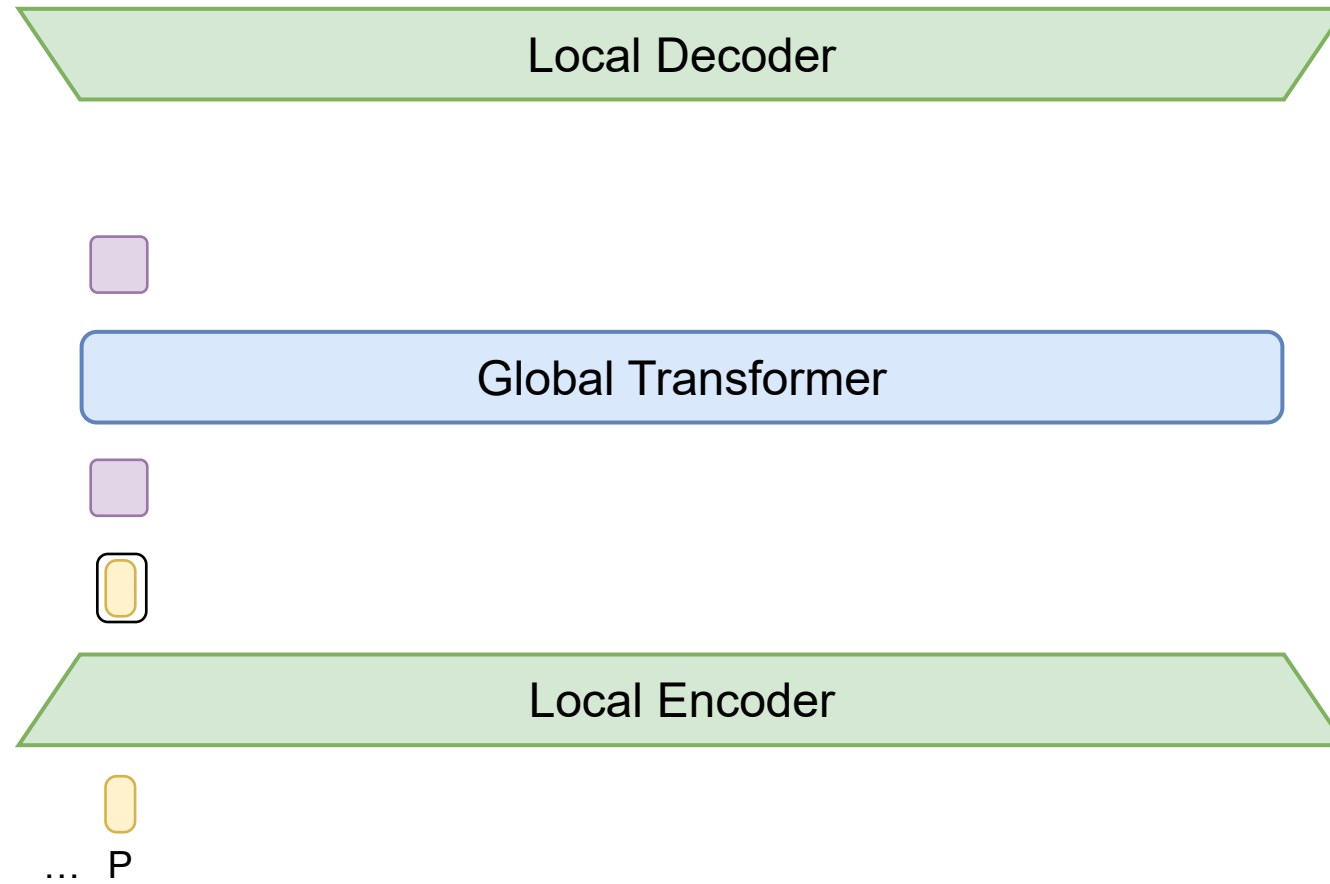
# BLT Generation



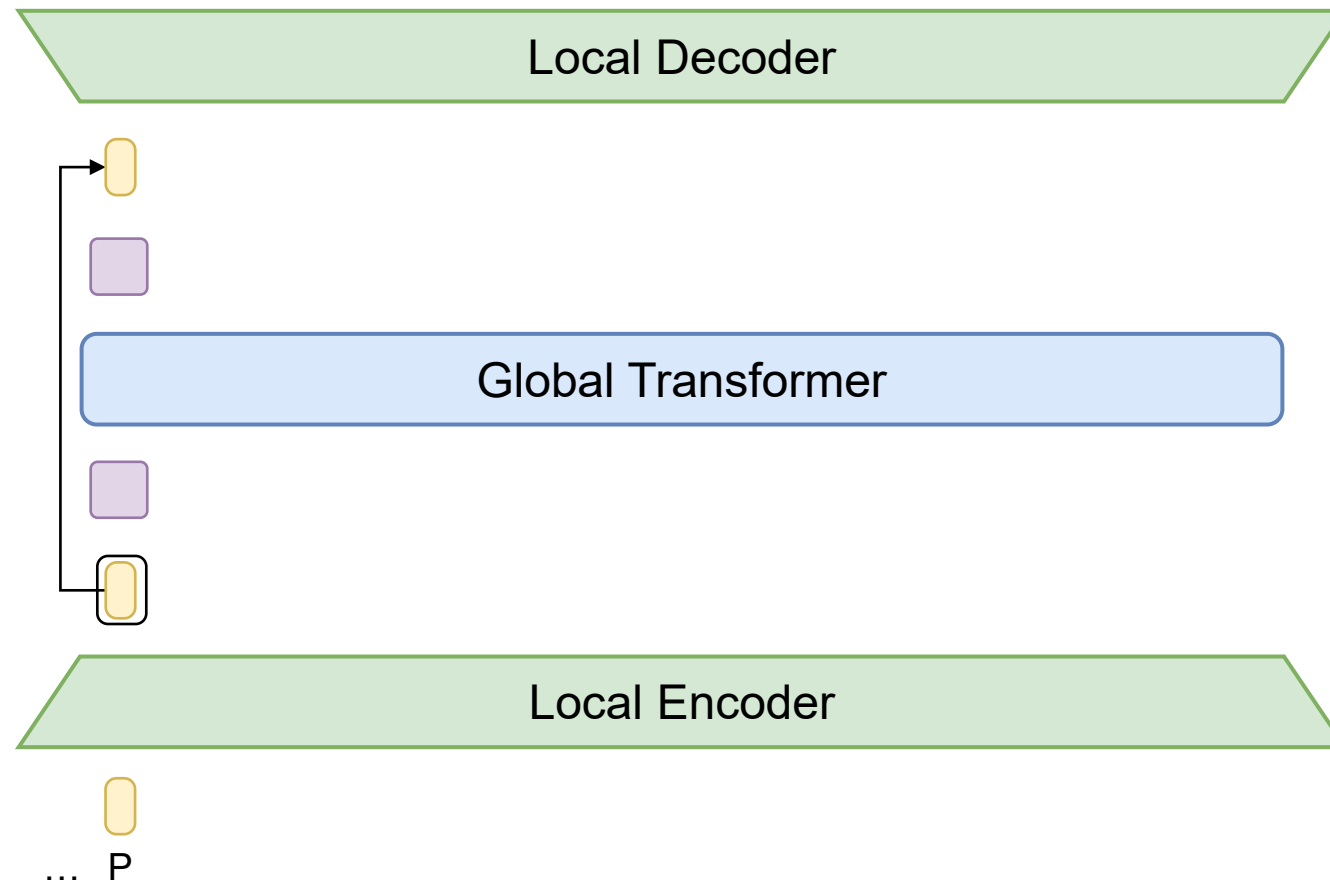
# BLT Generation



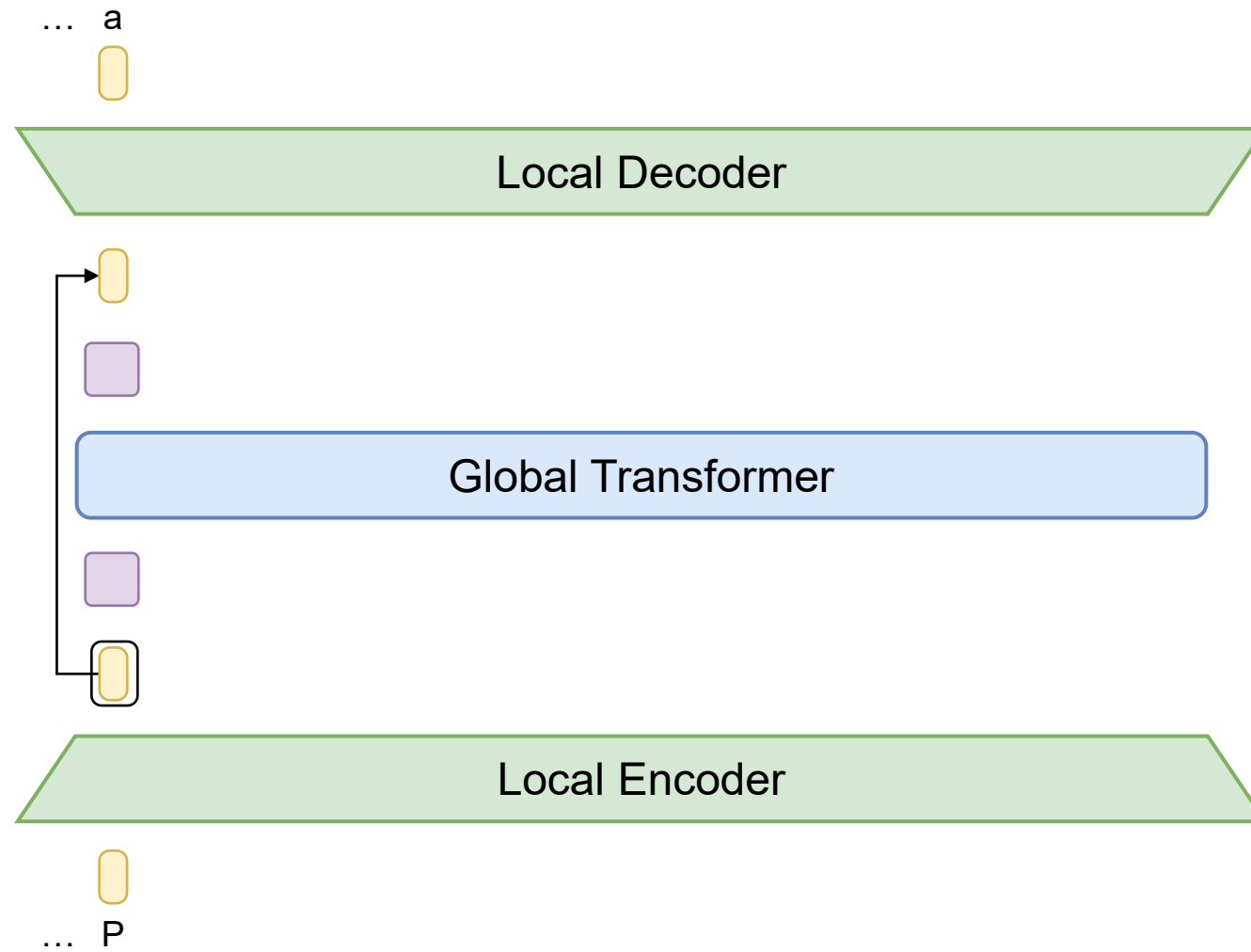
# BLT Generation



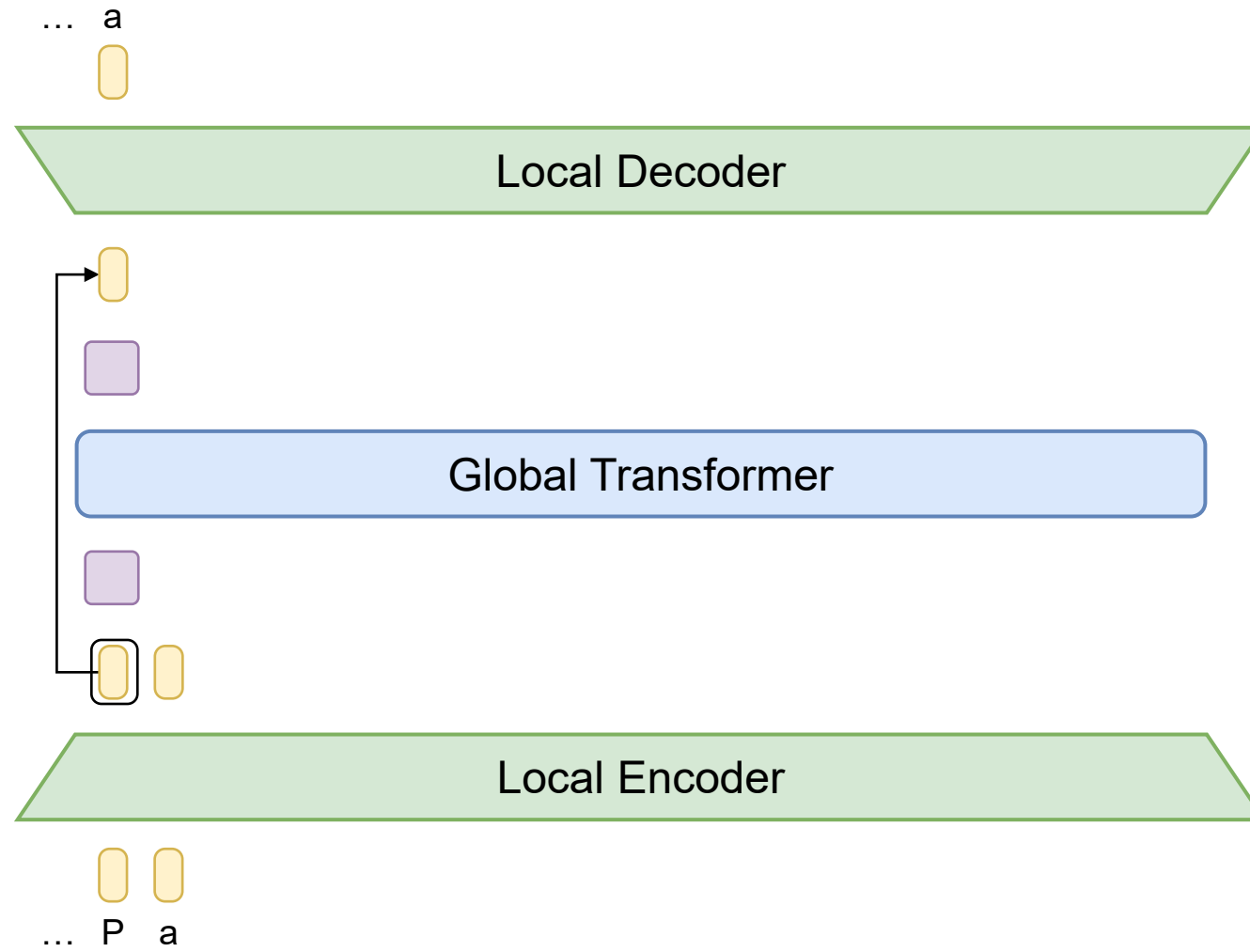
# BLT Generation



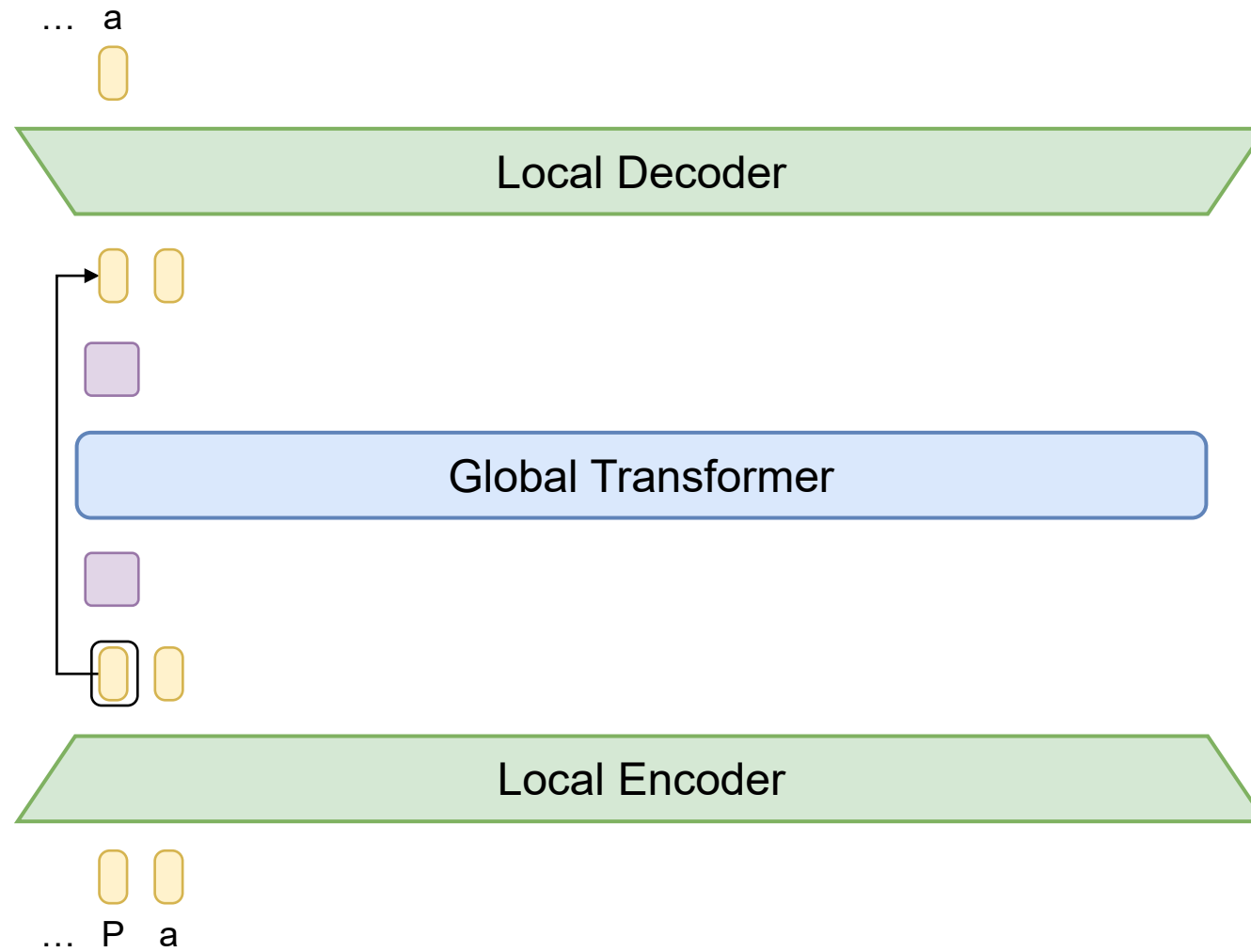
# BLT Generation



# BLT Generation

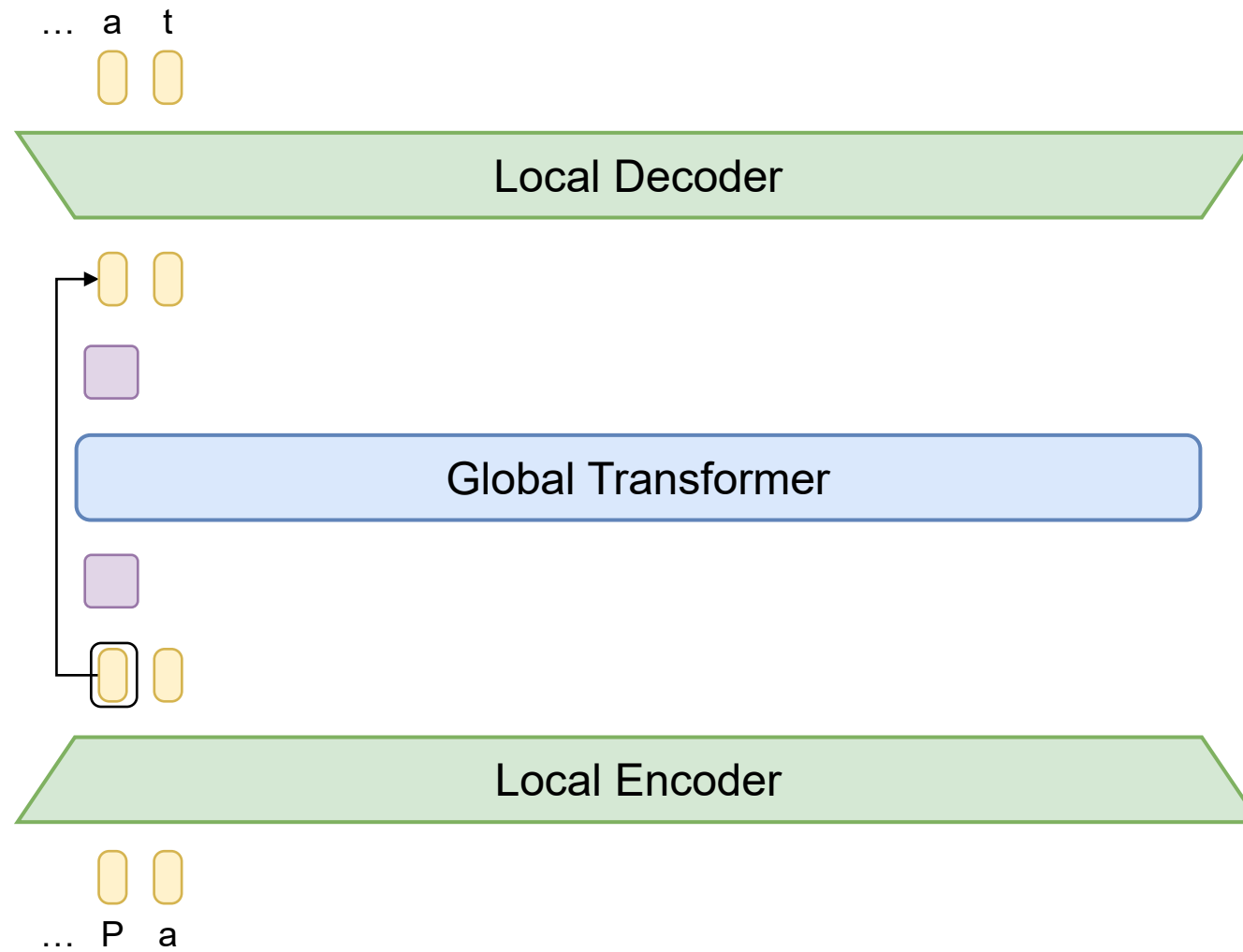


# BLT Generation

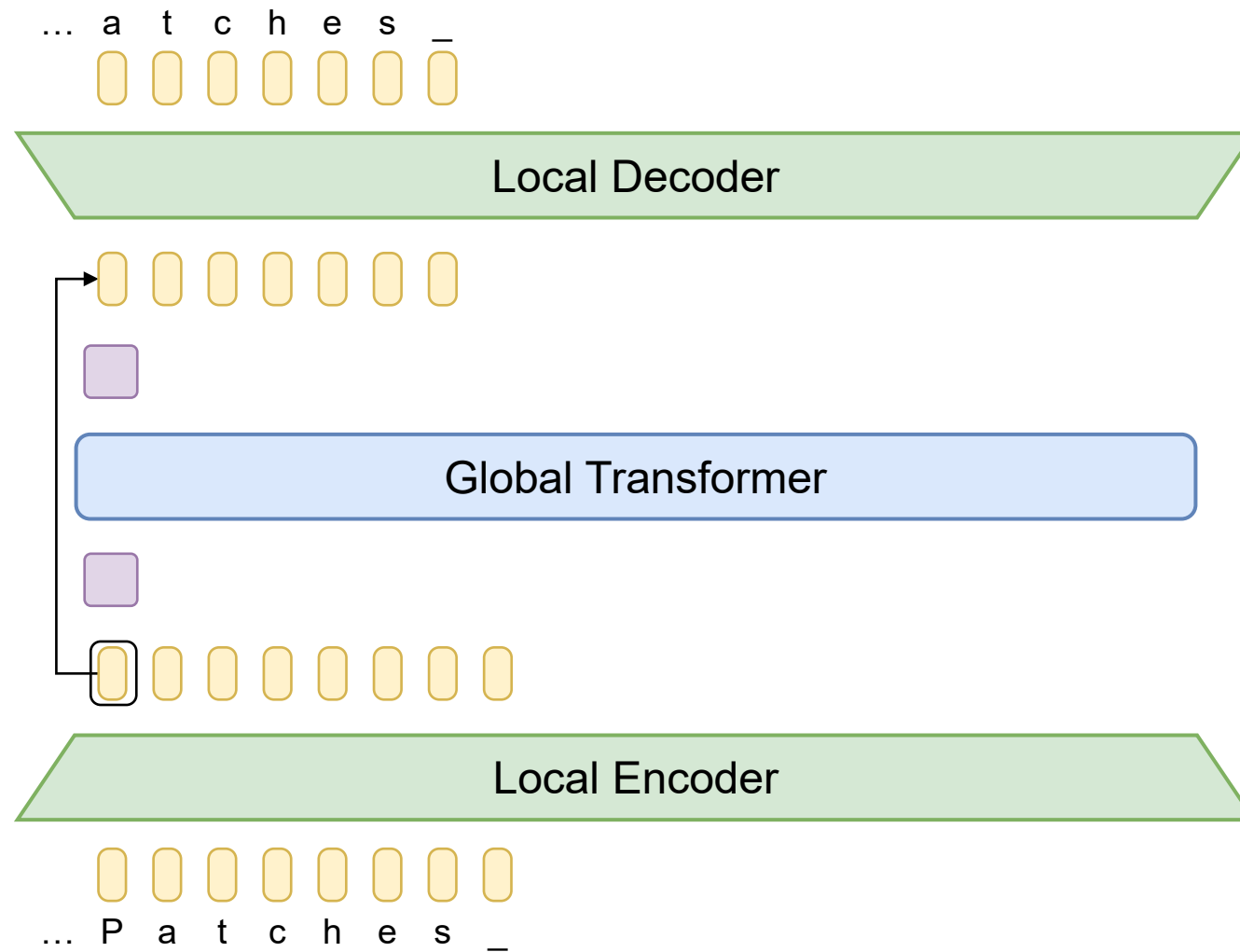




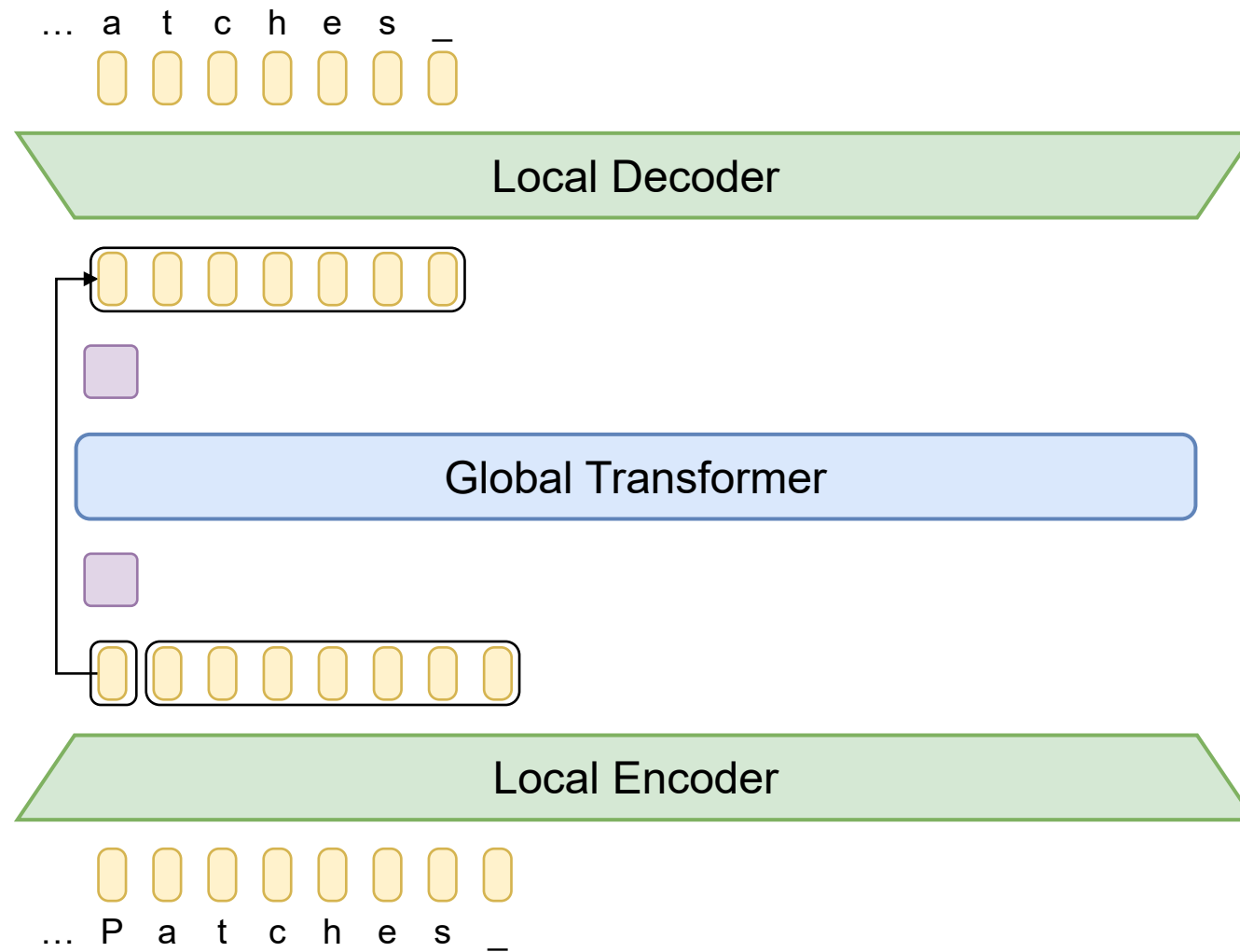
# BLT Generation



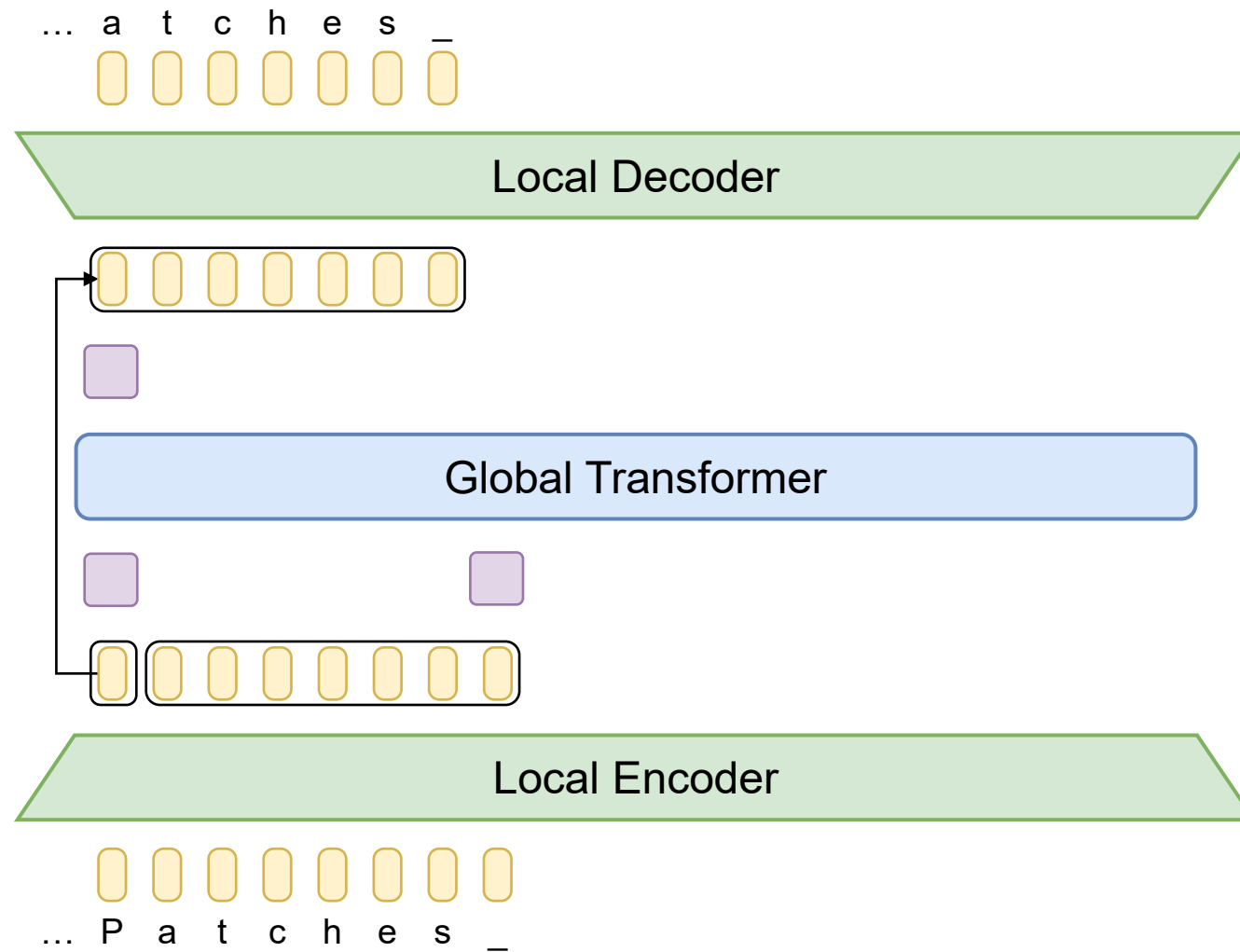
# BLT Generation



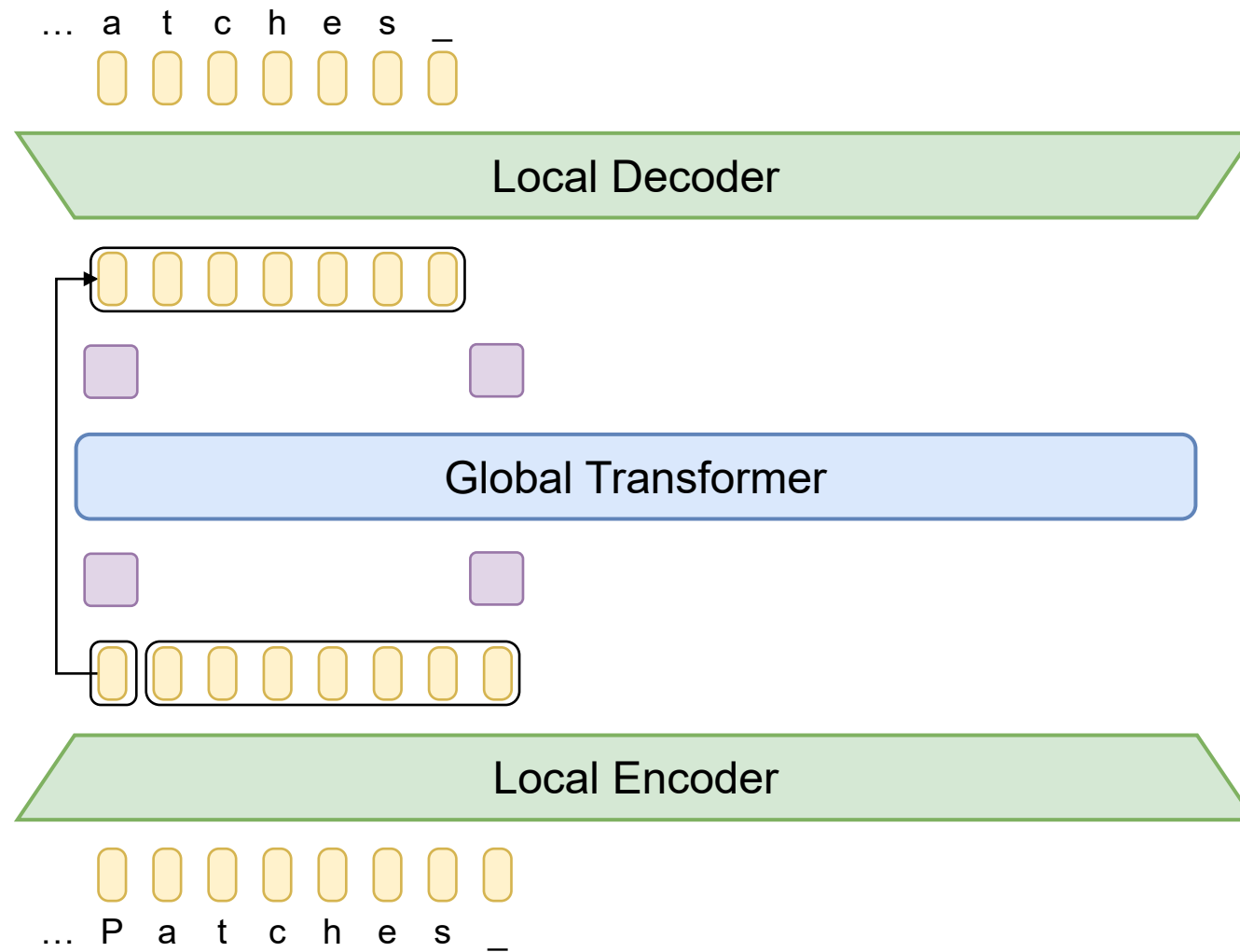
# BLT Generation



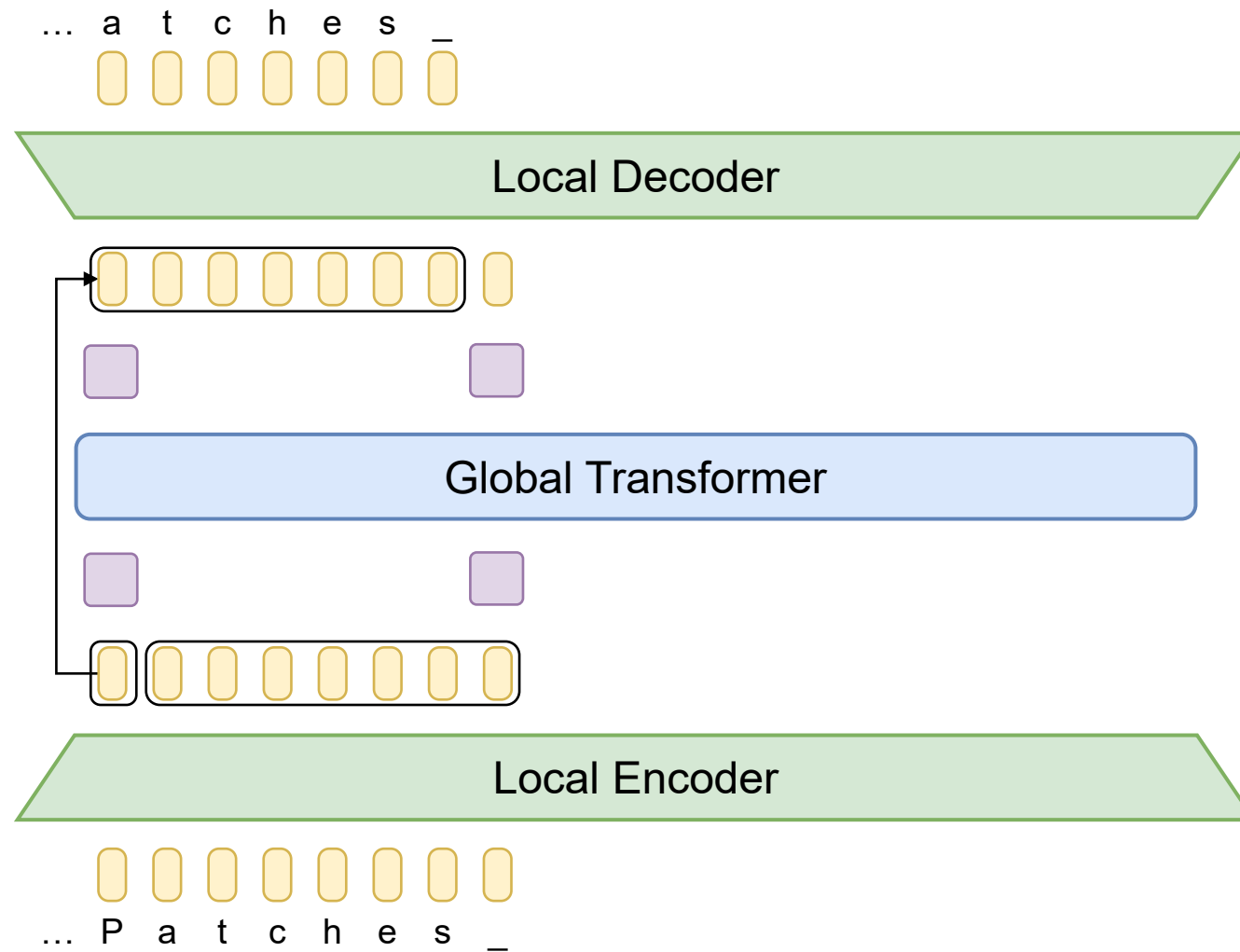
# BLT Generation



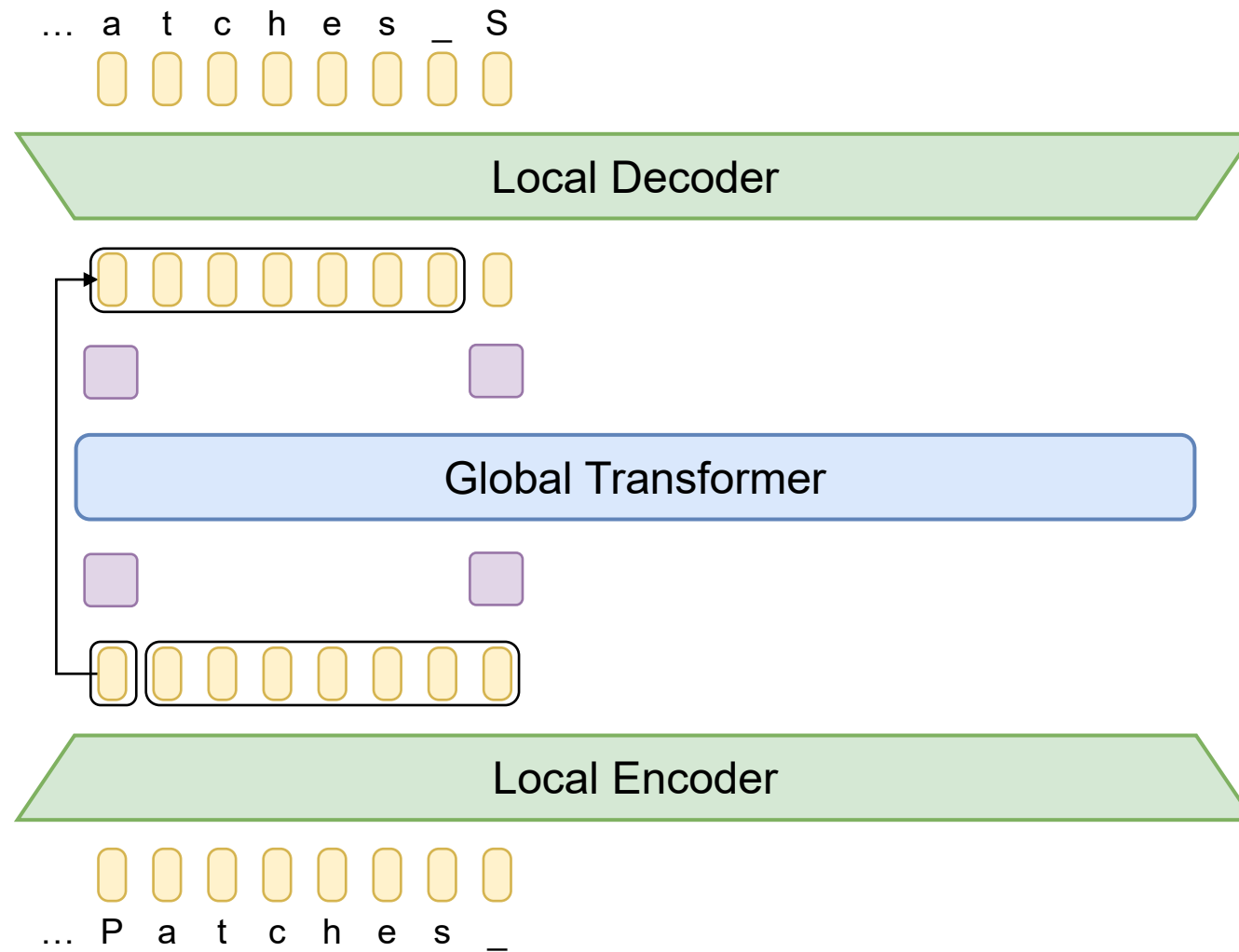
# BLT Generation



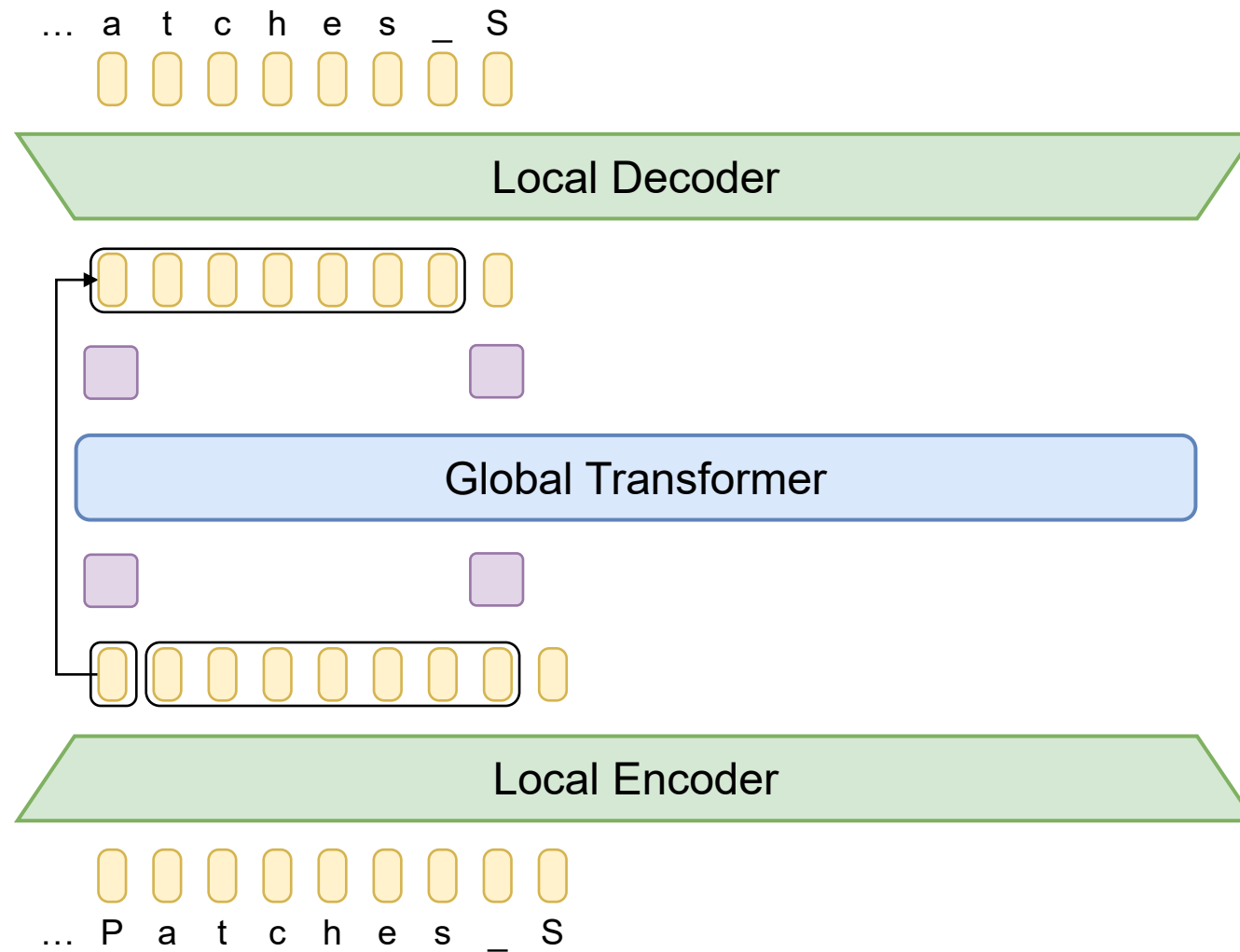
# BLT Generation



# BLT Generation

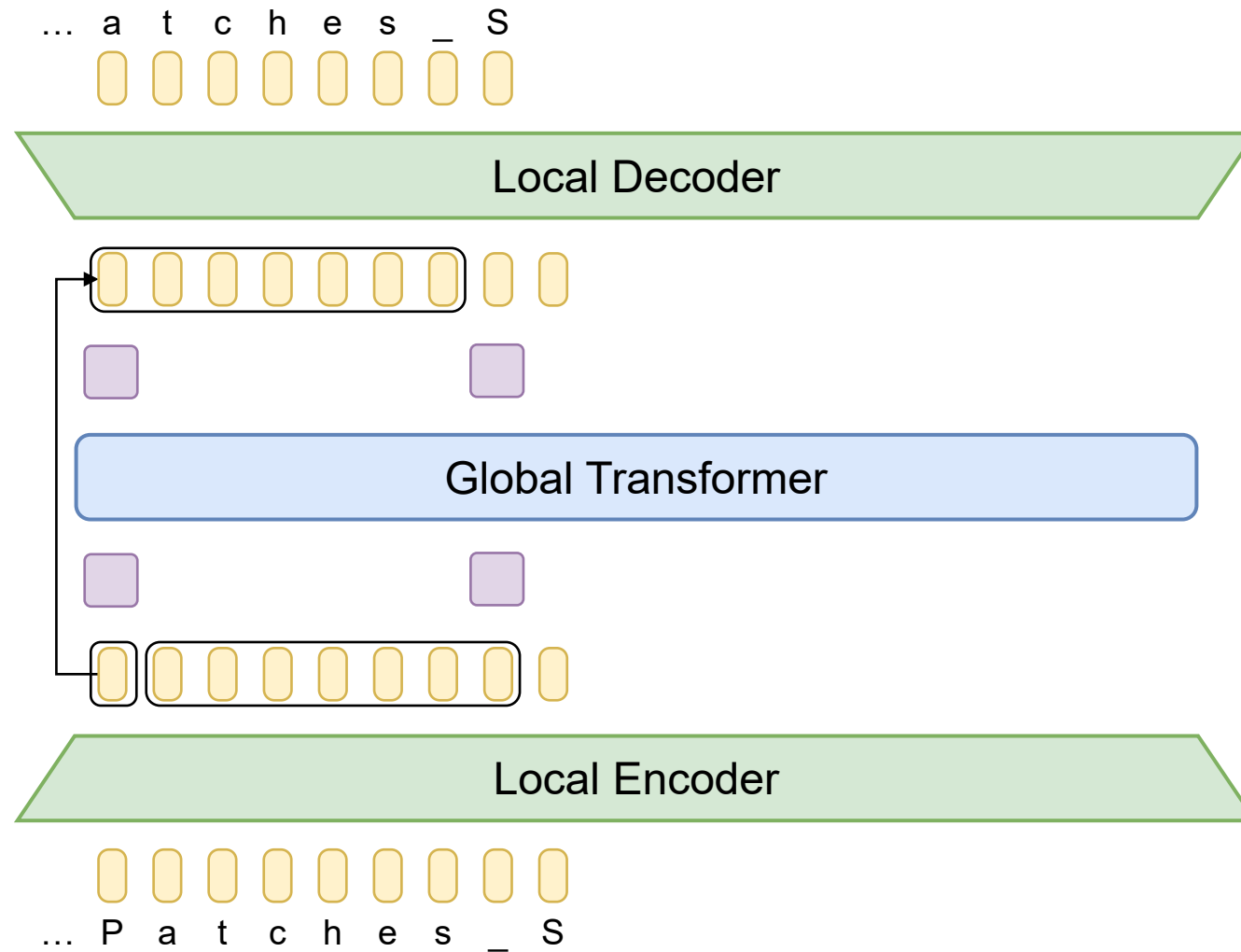


# BLT Generation

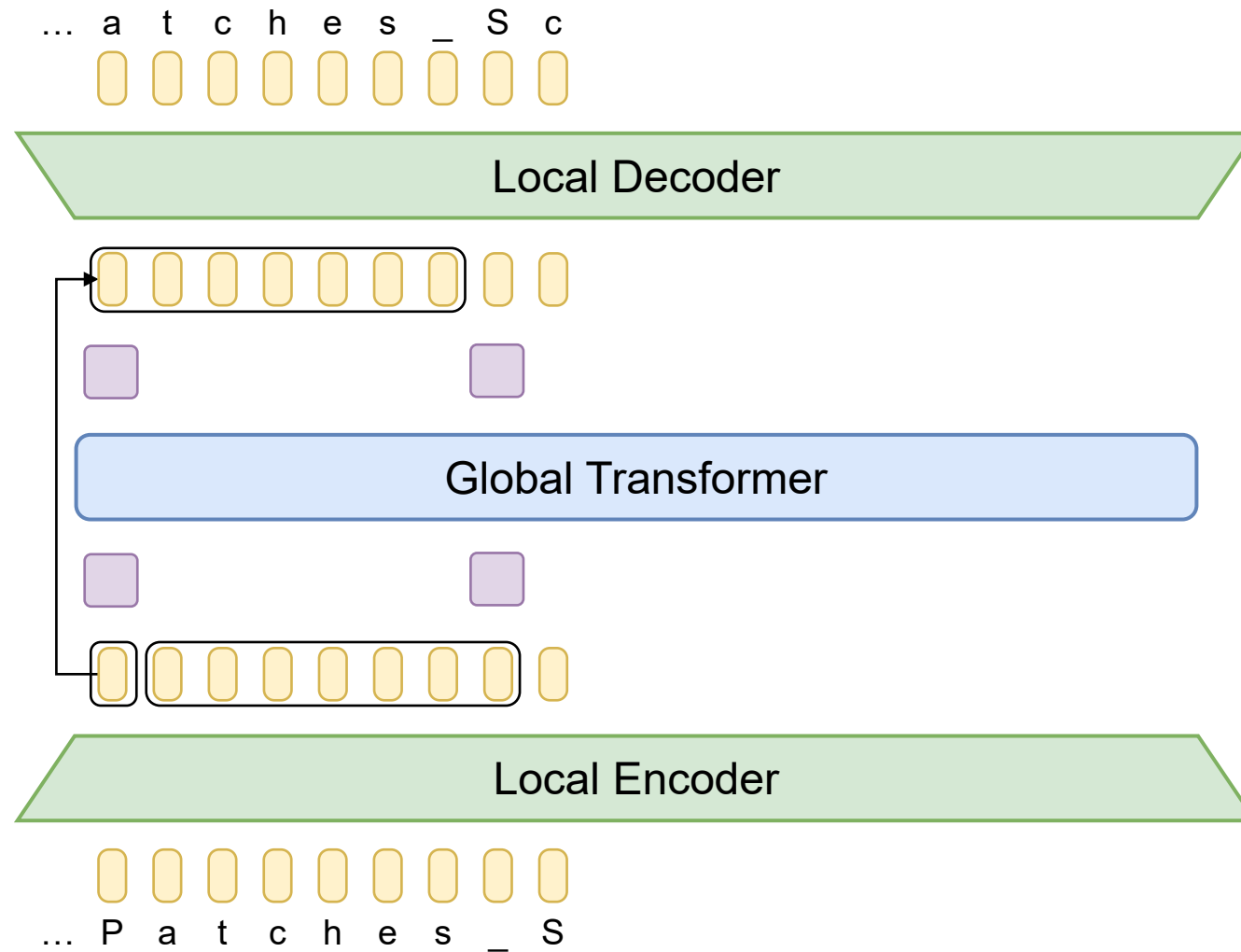




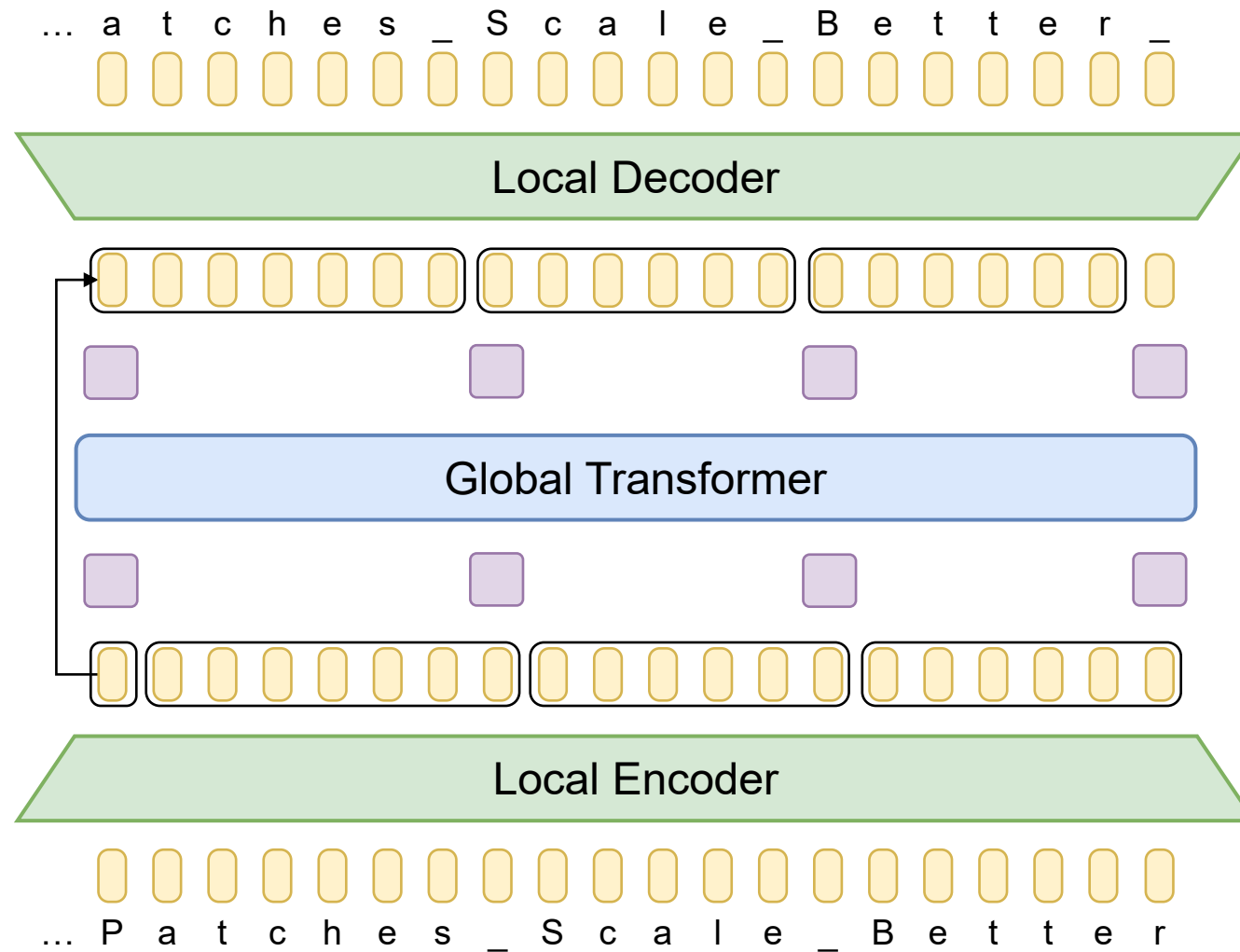
# BLT Generation



# BLT Generation



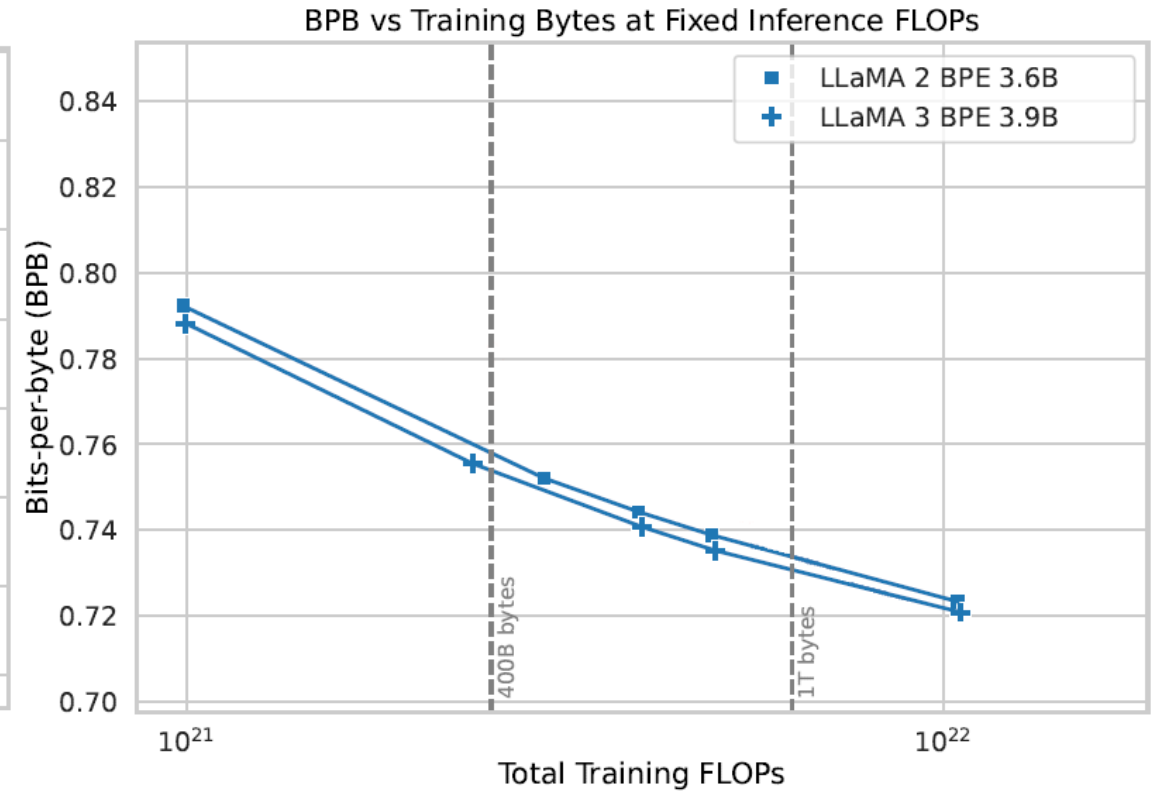
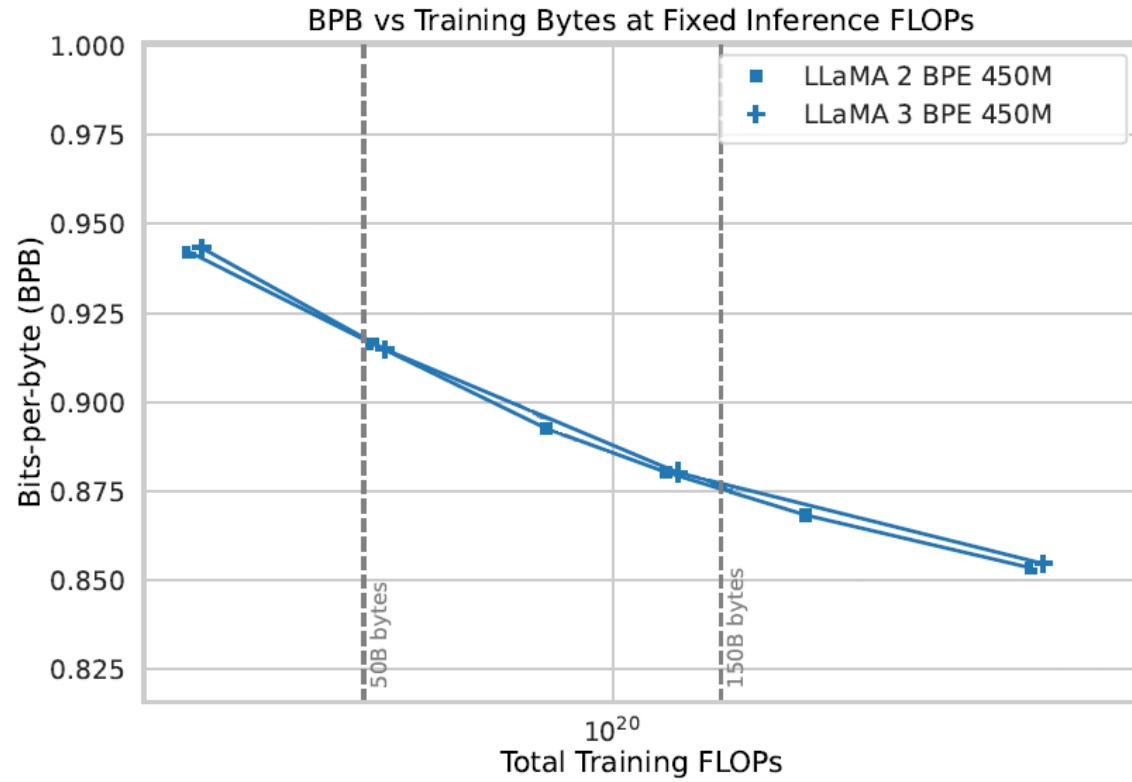
# BLT Generation



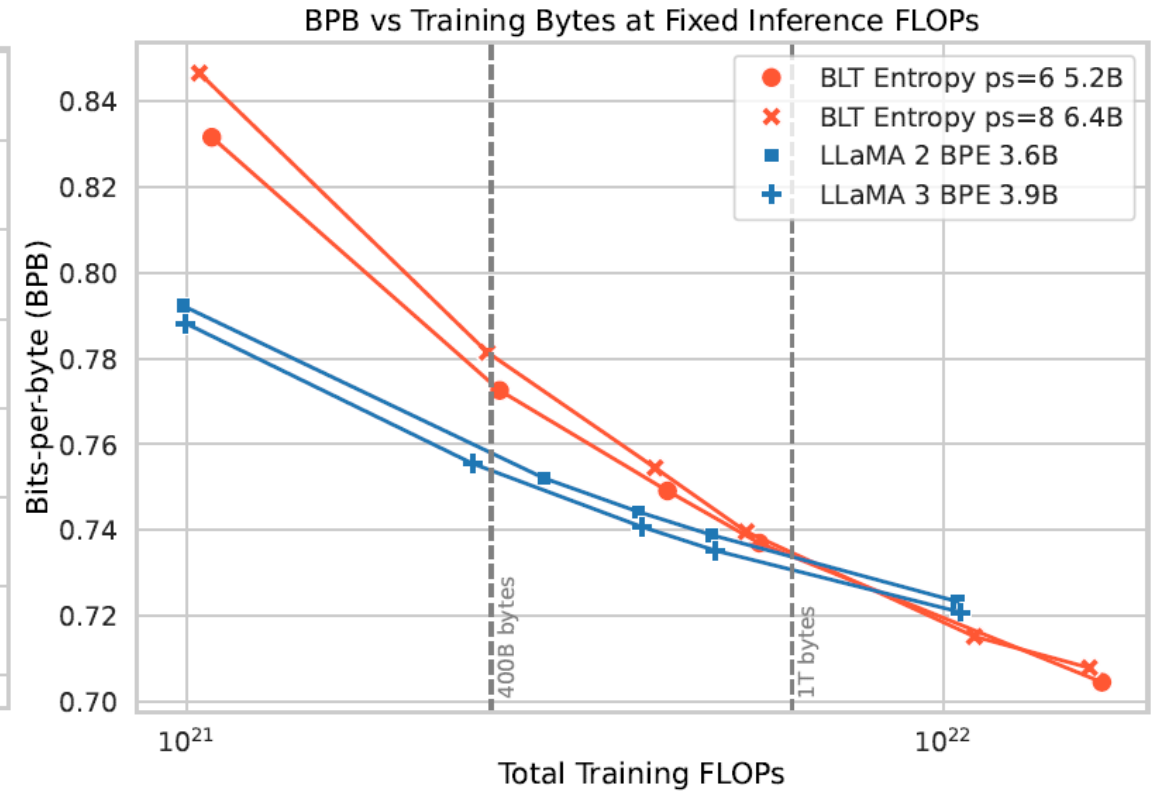
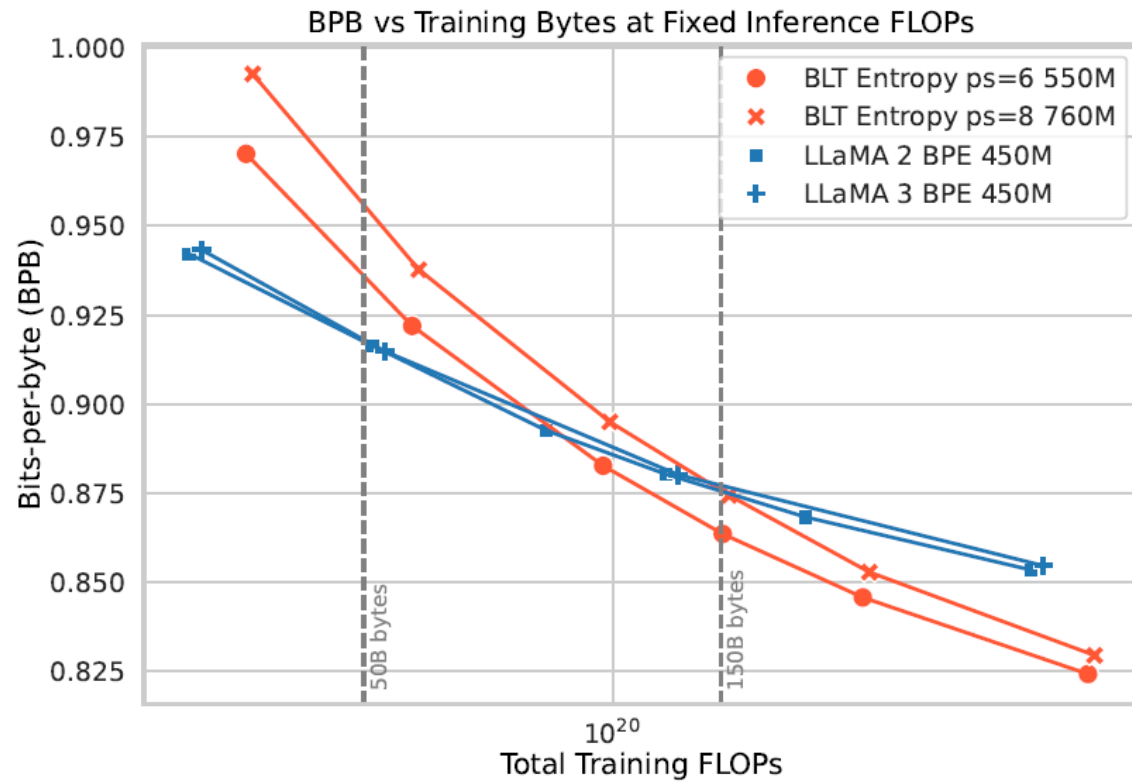
# Results

How does the BLT perform?

# Scaling Trends



# Scaling Trends



# From Token- to Patch-based

---

	Llama 3 8B (220B tokens)	BLT 8B (220B tokens)	BLT from Llama 3.1 8B (220B tokens)	Llama 3.1 8B (15T tokens)
<b>Arc-E</b>				
<b>Arc-C</b>				
<b>HellaSwag</b>				
<b>PIQA</b>				
<b>MMLU</b>				
<b>MBPP</b>				
<b>HumanEval</b>				

---

# From Token- to Patch-based

	Llama 3 8B (220B tokens)	BLT 8B (220B tokens)	BLT from Llama 3.1 8B (220B tokens)	Llama 3.1 8B (15T tokens)
<b>Arc-E</b>	<u>67.4</u>	66.8	66.6	83.4
<b>Arc-C</b>	40.4	38.8	<u>45.8</u>	55.2
<b>HellaSwag</b>	71.2	72.2	<u>76.1</u>	80.7
<b>PIQA</b>	77.0	<u>78.2</u>	77.4	80.7
<b>MMLU</b>	26.5	25.2	<u>63.7</u>	66.3
<b>MBPP</b>	11.8	10.0	<u>38.2</u>	47.2
<b>HumanEval</b>	9.2	7.3	<u>34.2</u>	37.2



# From Token- to Patch-based

	Llama 3 8B (220B tokens)	BLT 8B (220B tokens)	BLT from Llama 3.1 8B (220B tokens)	Llama 3.1 8B (15T tokens)
<b>Arc-E</b>	<u>67.4</u>	66.8	66.6	83.4
<b>Arc-C</b>	40.4	38.8	<u>45.8</u>	55.2
<b>HellaSwag</b>	71.2	72.2	<u>76.1</u>	80.7
<b>PIQA</b>	77.0	<u>78.2</u>	77.4	80.7
<b>MMLU</b>	26.5	25.2	<u>63.7</u>	66.3
<b>MBPP</b>	11.8	10.0	<u>38.2</u>	47.2
<b>HumanEval</b>	9.2	7.3	<u>34.2</u>	37.2

# From Token- to Patch-based

	Llama 3 8B (220B tokens)	BLT 8B (220B tokens)	BLT from Llama 3.1 8B (220B tokens)	Llama 3.1 8B (15T tokens)
<b>Arc-E</b>	<u>67.4</u>	66.8	66.6	83.4
<b>Arc-C</b>	40.4	38.8	<u>45.8</u>	55.2
<b>HellaSwag</b>	71.2	72.2	<u>76.1</u>	80.7
<b>PIQA</b>	77.0	78.2	77.4	80.7
<b>MMLU</b>	26.5	25.2	<u>63.7</u>	66.3
<b>MBPP</b>	11.8	10.0	<u>38.2</u>	47.2
<b>HumanEval</b>	9.2	7.3	<u>34.2</u>	37.2

# Personal Opinion & Potential Future Work

- Very exciting results
- Well executed evaluation
- Mostly empirically based
- Performance might not transfer to downstream tasks
- Change patching between benchmarks

# Personal Opinion & Potential Future Work

	Llama 3 8B (220B tokens)	BLT 8B (220B tokens)	BLT from Llama 3.1 8B (220B tokens)	Llama 3.1 8B (15T tokens)
<b>Arc-E</b>	67.4	66.8	66.6	83.4
<b>Arc-C</b>	40.4	38.8	45.8	55.2
<b>HellaSwag</b>	71.2	72.2	76.1	80.7
<b>PIQA</b>	77.0	78.2	77.4	80.7
<b>MMLU</b>	26.5	25.2	63.7	66.3
<b>MBPP</b>	11.8	10.0	38.2	47.2
<b>HumanEval</b>	9.2	7.3	34.2	37.2

**Table 5** Initializing the global transformer model of BLT from the non-embedding parameters of Llama 3 improves performance on several benchmark tasks. **First three models trained on the Llama 2 data for compute-optimal steps.**

	Llama 3 BPE	Space Patching BLT	Entropy BLT
<b>Arc-E</b>	67.4	67.2	68.9
<b>Arc-C</b>	40.5	37.6	38.3
<b>HellaSwag</b>	71.3	70.8	72.7
<b>PIQA</b>	77.0	76.5	77.6

**Table 6** Benchmark evaluations of two patching schemes for 8b BLT models and BPE Llama3 baseline. These models are **trained on the Llama 2 data for the optimal number of steps** as determined by [Dubey et al. \(2024\)](#).

# Personal Opinion & Potential Future Work

- Very exciting results
  - Well executed evaluation
  - Well written paper
  - Open-source
  - Open-weights (soon)
- 
- Apply to other domain
  - Multimodal models
- 
- Mostly empirically based
  - Performance might not transfer to downstream tasks
  - Change patching between benchmarks
  - Claim that the idea of dynamic patching is novel

# Q&A and Discussion

# Tokenization: Byte-Pair Encoding

b a n a n a

b a n d a n a

p a n

b a n d

a n

Pair	Count
an	7
na	3
ba	3
nd	2



b an a n a

b an d an a

p an

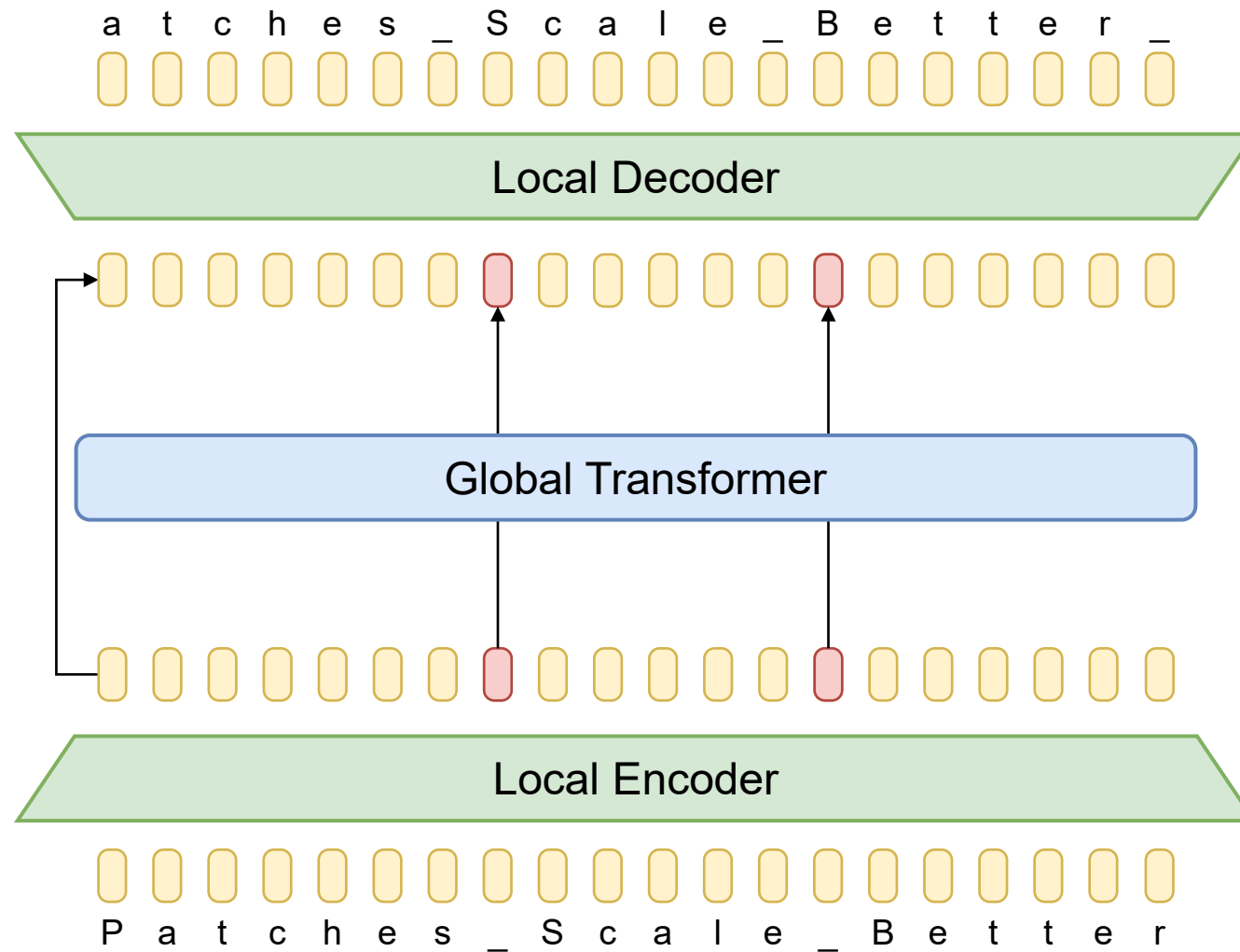
b an d

an

Vocabulary: {b, a, n, d, p}

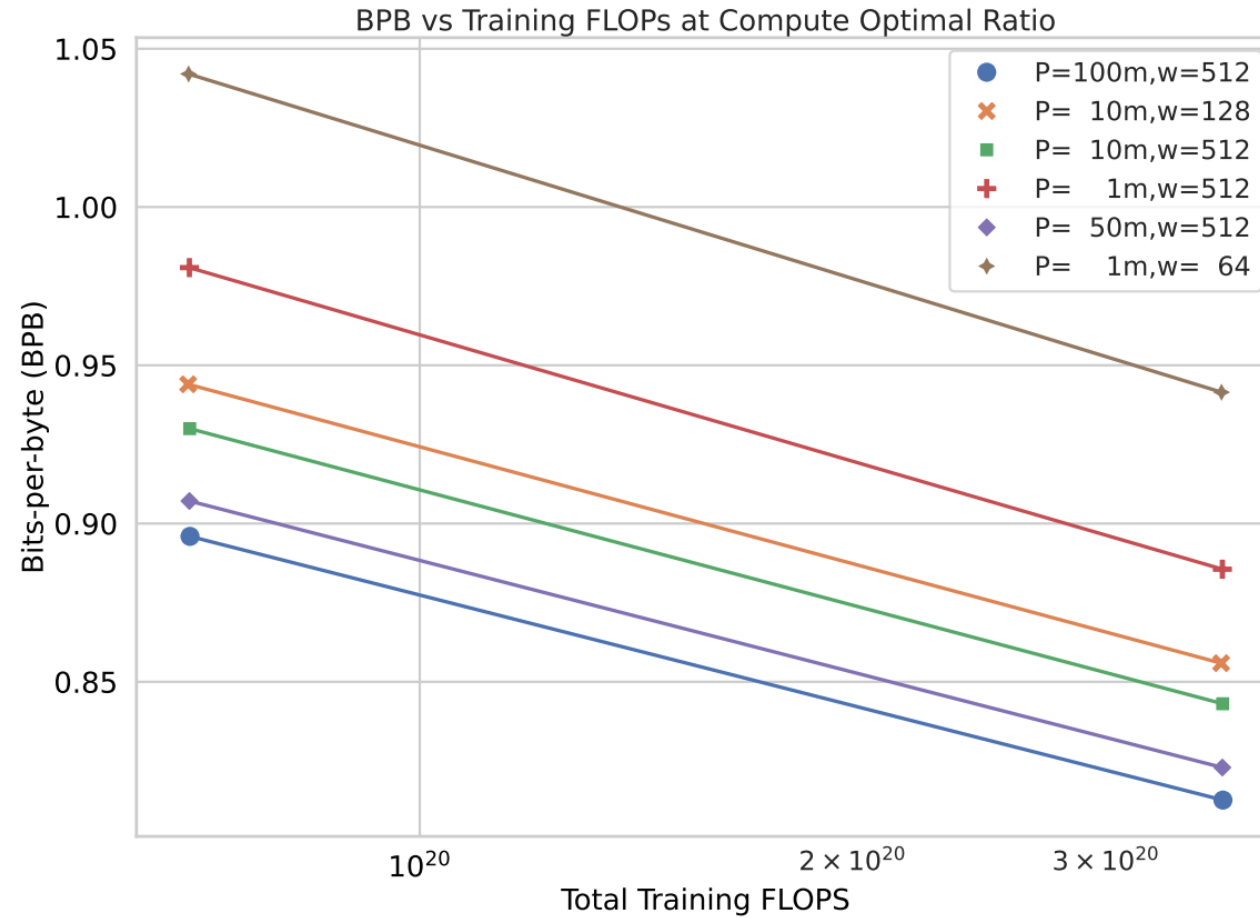
Vocabulary: {b, a, n, d, p, an}

# SpaceByte (Original Interpretation)

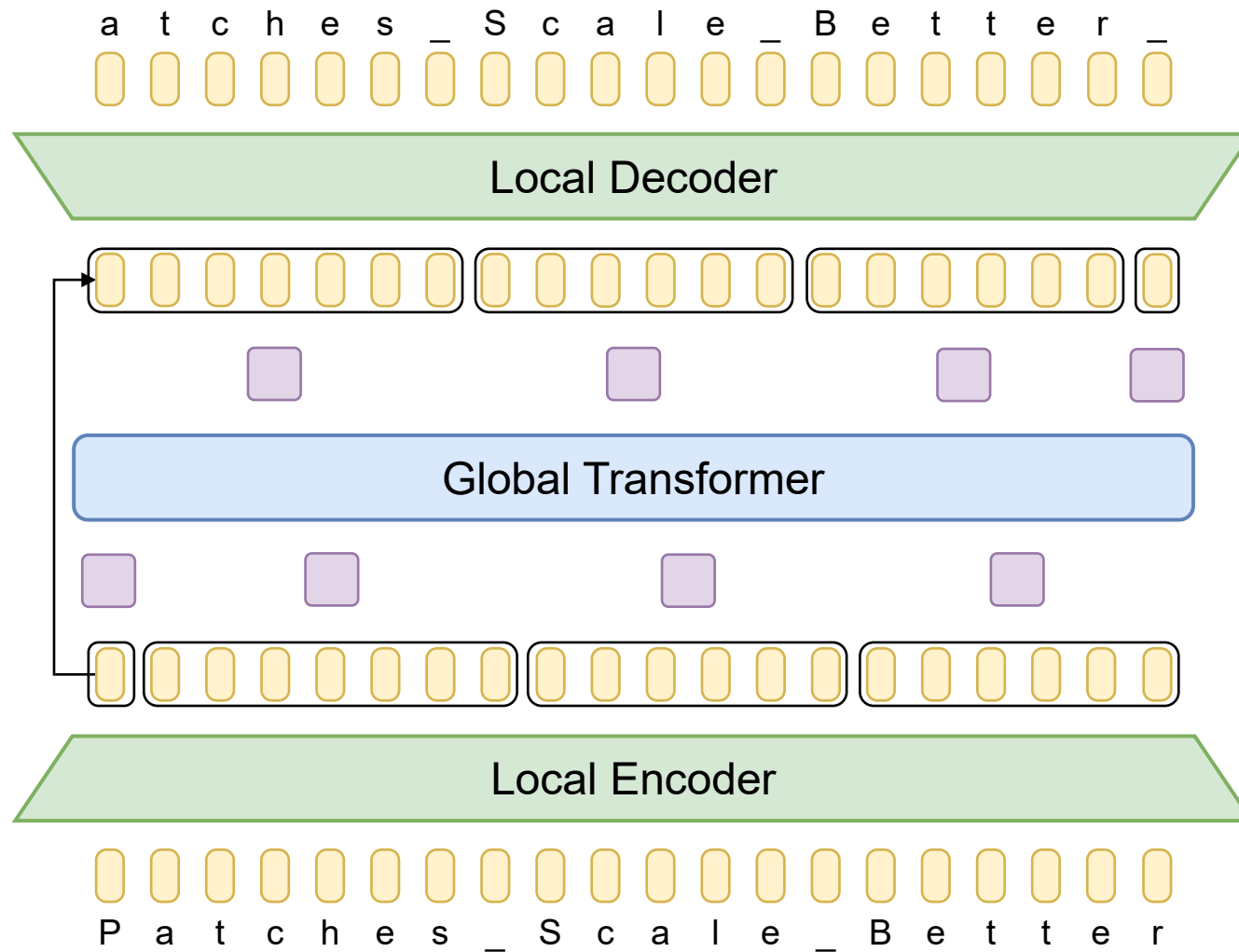




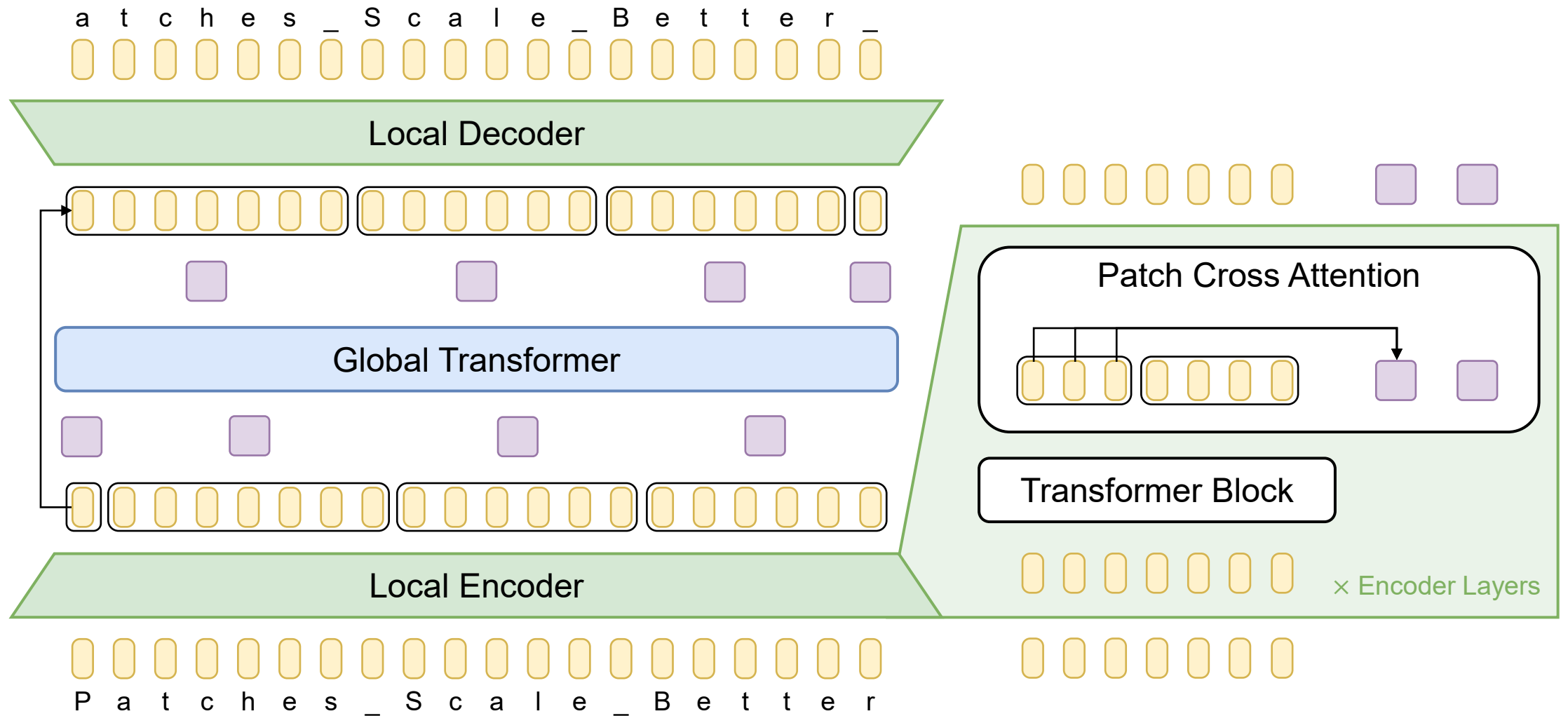
# BLT Entropy Model Scaling



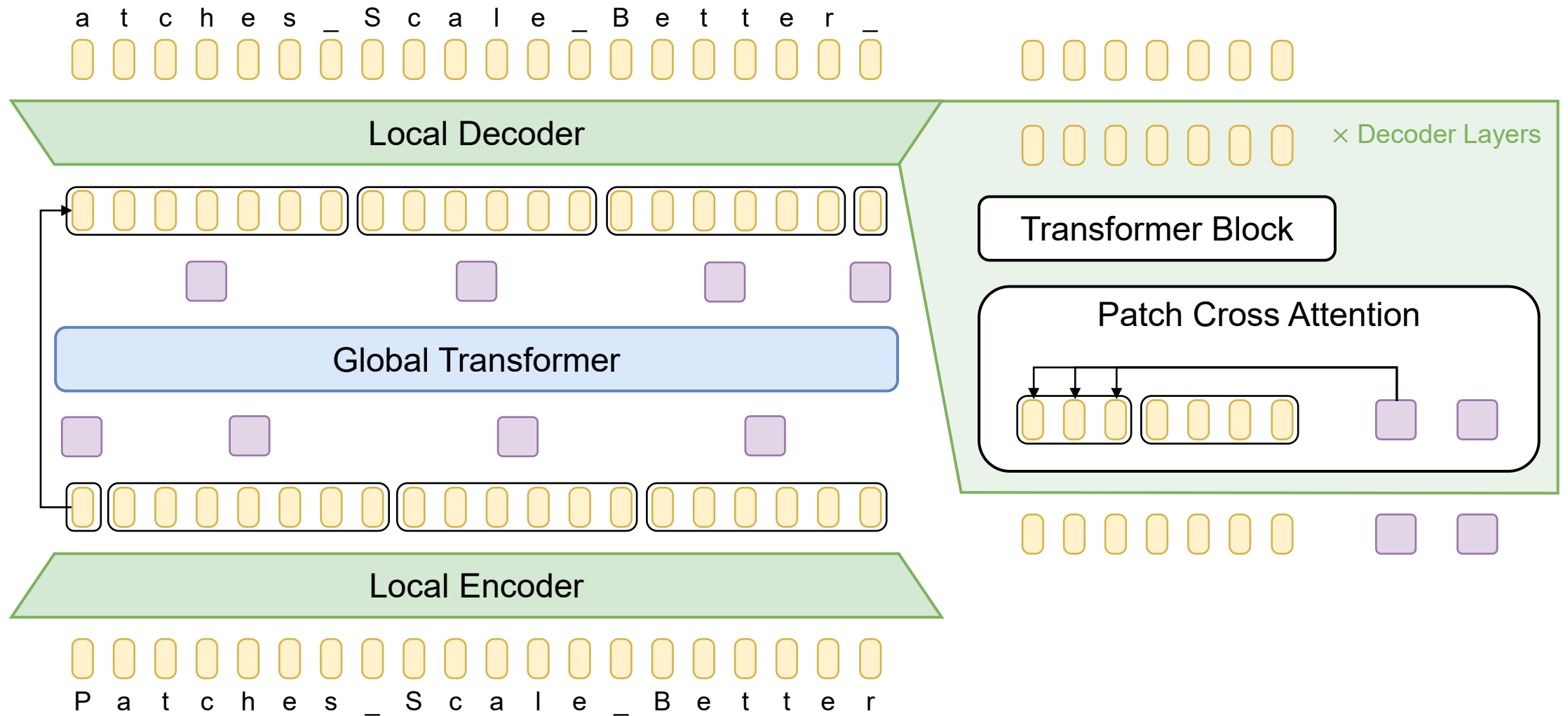
# BLT Architecture (Cross Attention)



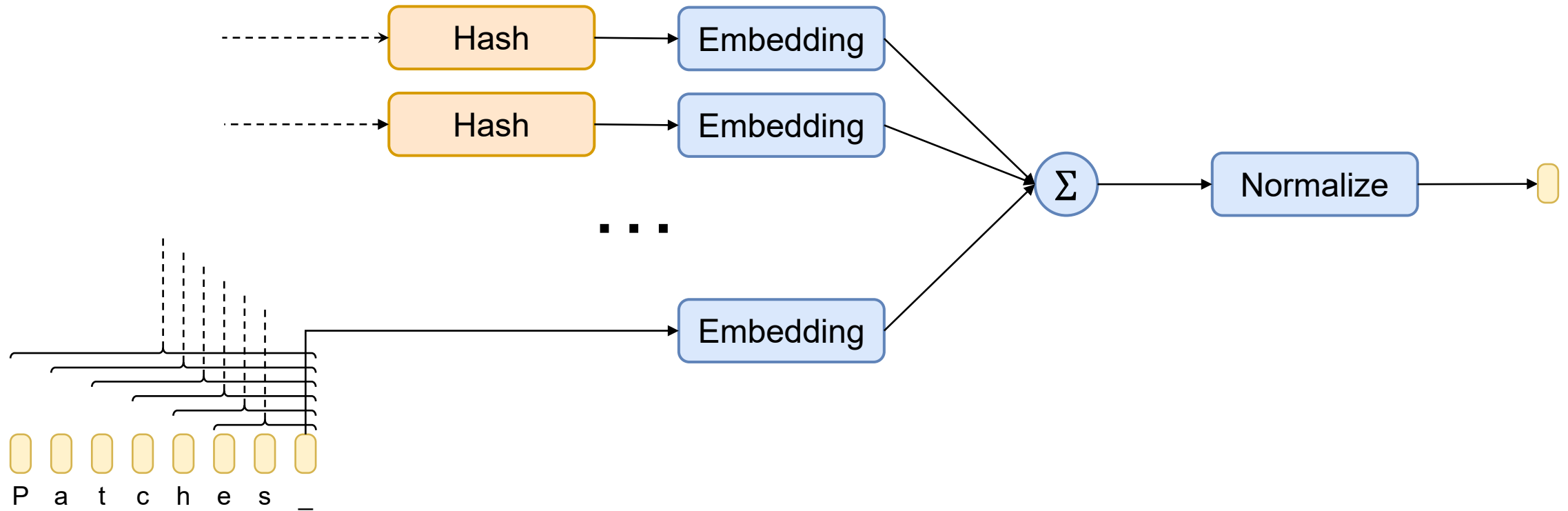
# BLT Architecture (Cross Attention)



# BLT Architecture (Cross Attention)



# BLT Byte Embedding

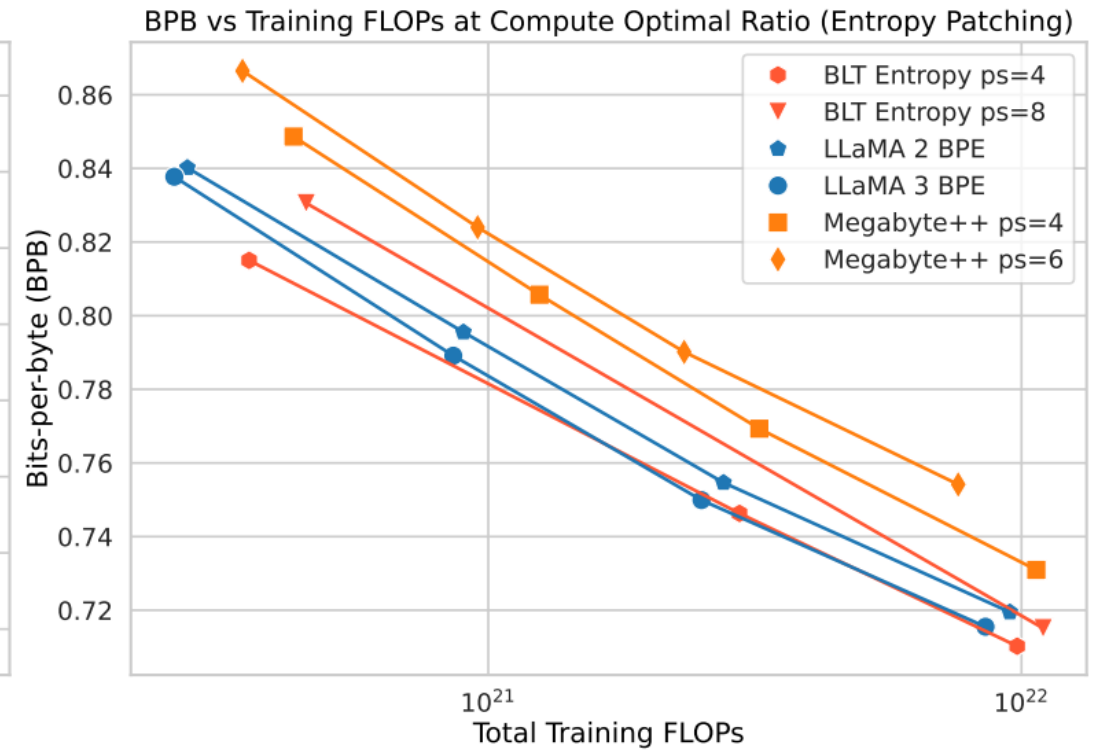
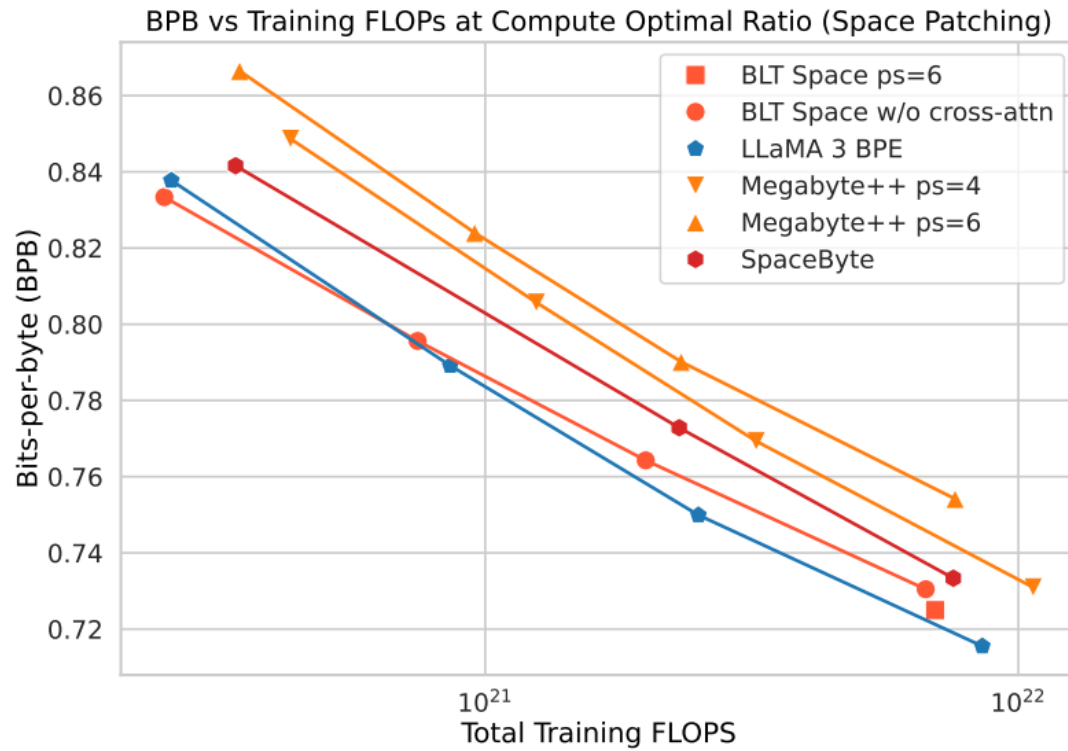


# Bits-per-Byte

- Tokenization independent Perplexity (Uncertainty)

$$BPB(x) = \frac{L_{CE}(x)}{\ln(2) \cdot n_{bytes}}$$

# MegaByte vs SpaceByte vs BLT vs LLaMA



# Downstream FLOP Matched Training

	Llama 3 1T Tokens	BLT-Space 6T Bytes	BLT-Entropy 4.5T Bytes
<b>Arc-E</b>	77.6	75.4	<b>79.6</b>
<b>Arc-C</b>	<b>53.3</b>	49.8	52.1
<b>HellaSwag</b>	79.1	79.6	<b>80.6</b>
<b>PIQA</b>	80.7	<b>81.1</b>	80.6
<b>MMLU</b>	<b>58.1</b>	54.8	57.4
<b>MBPP</b>	40.2	37.6	<b>41.8</b>
<b>HumanEval</b>	31.1	27.4	<b>35.4</b>
<b>Average</b>	60.0	58.0	<b>61.1</b>
<b>Bytes/Patch on Train Mix</b>	4.4	<b>6.1</b>	4.5



# Robustness

	Llama 3 (1T tokens)	Llama 3.1 (16T tokens)	BLT (1T tokens)
<b>HellaSwag Original</b>	79.1	<u>80.7</u>	<b>80.6</b>
<b>HellaSwag Noise Avg.</b>	56.9	<u>64.3</u>	<b>64.3</b>
- AntSpeak	45.6	<u>61.3</u>	<b>57.9</b>
- Drop	53.8	<u>57.3</u>	<b>58.2</b>
- RandomCase	55.3	<u>65.0</u>	<b>65.7</b>
- Repeat	57.0	<u>61.5</u>	<b>66.6</b>
- UpperCase	72.9	<u>76.5</u>	<b>77.3</b>
<b>Phonology-G2P</b>	11.8	<u>18.9</u>	<b>13.0</b>
<b>CUTE</b>	27.5	20.0	<b>54.1</b>
- Contains Char	0.0	0.0	<b>55.9</b>
- Contains Word	55.1	21.6	<b>73.5</b>
- Del Char	34.6	34.3	<b>35.9</b>
- Del Word	<b>75.5</b>	<u>84.5</u>	56.1
- Ins Char	7.5	0.0	<b>7.6</b>
- Ins Word	<b>33.5</b>	<u>63.3</u>	31.2
- Orthography	43.1	0.0	<b>52.4</b>
- Semantic	65	0.0	<b>90.5</b>
- Spelling	1.1	-	<b>99.9</b>
- Spelling Inverse	30.1	3.6	<b>99.9</b>
- Substitute Char	0.4	1.2	<b>48.7</b>
- Substitute Word	16.4	6.8	<b>72.8</b>
- Swap Char	2.6	2.4	<b>11.5</b>
- Swap Word	20.1	4.1	<b>21</b>

# Hyperparameters

Model	Encoder				Global Latent Transf.				Decoder				Cross-Attn.	
	$l_{\mathcal{E}}$	#heads	$h_{\mathcal{E}}$	#Params	$l_{\mathcal{G}}$	#heads	$h_{\mathcal{G}}$	#Params	$l_{\mathcal{D}}$	#heads	$h_{\mathcal{D}}$	#Params	#heads	k
<b>400M</b>	1	12	768	7M	24	10	1280	470M	7	12	768	50M	10	2
<b>1B</b>	1	16	1024	12M	25	16	2048	1B	9	16	1024	113M	16	2
<b>2B</b>	1	16	1024	12M	26	20	2560	2B	9	16	1024	113M	16	3
<b>4B</b>	1	16	1024	12M	36	24	3072	4.1B	9	16	1024	113M	16	3
<b>8B</b>	1	20	1280	20M	32	32	4096	6.4B	6	20	1280	120M	20	4

# (Authors) Future Work

- Compute optimal training for BLT models
- Scale larger
- Optimizations
- End-to-End patch learning

# Dynamic Patching on MMLU

The following are multiple choice questions (with answers) about college physics.

A refracting telescope consists of two converging lenses separated by 100 cm. The eyepiece lens has a focal length of 20 cm. The angular magnification of the telescope is

A. 4

B. 5

C. 6

D. 20

Answer: A

...

The muon decays with a characteristic lifetime of about  $10^{-6}$  s into an electron, a muon neutrino, and an electron anti neutrino. The muon is forbidden from decaying into an electron and just a single neutrino by the law of conservation of

A. charge

B. mass

C. energy and momentum

D. lepton number

Answer: D

The quantum efficiency of a photon detector is 0.1. If 100 photons are sent into the detector, one after the other, the detector will detect photons

A. an average of 10 times, with an rms deviation of about 4

B. an average of 10 times, with an rms deviation of about 3

C. an average of 10 times, with an rms deviation of about 1

D. an average of 10 times, with an rms deviation of about 0.1

Answer: A

# Patching Methods

4-Strided	Daen	erys	Tar	gary	en	i	s	in	Gam	e	of	Thr	ones	,	a	fant	asy	epic	by	Geor	ge	R	.	R.	Mart	in.
BPE	Da	enery	s	T	arg	ary	en	is	in	Game	of	Thrones	,	a	fantasy	epic	by	George	R	.	R.	.	Martin	.		
Entropy	D	a	e	nerys	Targaryen	is	in	G	ame	of	Thrones,	a	fa	ntasy	epic	by	G	eorge	R.R.	Martin.						
Entropy + Monotonicity	D	aenerys	Targar	yen	is	in	Game	of	Thrones	,	a	fantasy	epic	by	George	R.R.	Martin	.								
Space	Daenerys	Targaryen	is	in	Game	of	Thrones,	a	fantasy	epic	by	George	R.	R.	Martin.											
CNN	Daenerys	Targaryen	is	in	Game	of	Thrones,	a	fantasy	epic	by	George	R.R.	Martin.												