

From Pixels to Tokens: Byte-Pair Encoding on Quantized Visual Modalities

ICLR 2025

Wanpeng Zhang¹, Zilong Xie², Yicheng Feng¹, Yijiang Li³, Xingrun Xing^{4,5}, Sipeng Zheng⁴, Zongqing Lu^{1,6}

¹PKU ²CUHK ³UCSD ⁴BAAI ⁵CASIA ⁶BeingBeyond



Background

- MLLMs struggle with effectively aligning visual and textual modalities
- Current approaches:
 - Late-fusion with specialized encoders: Complex alignment challenges
 - Early-fusion token-based: Lacks explicit structural information
- Text-only LLMs benefit from BPE tokenization but visual modalities lack equivalent

Theoretical Foundation

Proposition 1. For data generating processes described in either Scenario 1 or Scenario 2, as $m \rightarrow \infty$, the optimal cross-entropy loss among unigram model family $\mathcal{Q}_{1\text{-gram}}$ satisfies

$$\liminf_{m \rightarrow \infty} \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}_m(Q) \geq H(\pi) = \sum_{a \in \mathcal{C}} \pi(a) \log(\pi(a)). \quad (3)$$

In contrast, the optimal unconstrained cross entropy loss satisfies

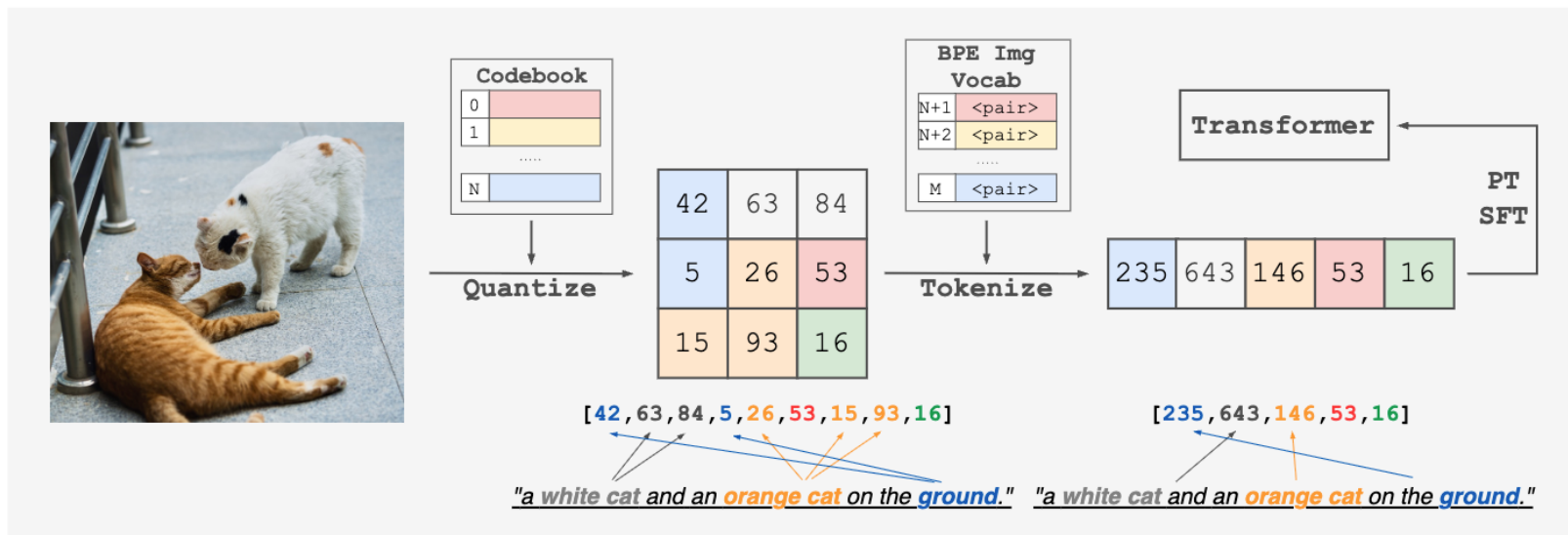
$$\lim_{m \rightarrow \infty} \min_Q \mathcal{L}_m(Q) = H_\infty \triangleq - \sum_{a \in \mathcal{C}} \sum_{a' \in \mathcal{C}} \pi(a) P(a' | a) \log(P(a' | a)). \quad (4)$$

Proposition 2. For data generating processes described in either Scenario 1 or Scenario 2, assume that $\delta \triangleq \min_{a, a' \in \mathcal{C}} P(a' | a) > 0$. Then there exists a tokenizer with a dictionary containing at most D tokens, along with an encoding function $\text{enc}(\cdot)$ applied to \mathbf{X} , such that

$$\limsup_{m \rightarrow \infty} \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) \leq \frac{1}{1 - \varepsilon} H_\infty, \quad (5)$$

where $\varepsilon = \log(1/\delta)/(0.99 \log(D))$ and $D \in \mathbb{N}$ is an arbitrary constant that is sufficiently large.

Our Method



- **BPE Image Tokenizer:**

- Step 1: Quantize image into initial token IDs using VQ-GAN (codebook size 8192)
- Step 2: Learn merged tokens based on frequency patterns (like text BPE)
- Step 3: Combine tokens into semantically meaningful units with structural information

Our Method

- Information loss theoretical upper bound:

$$L_{bpe} \leq (|D_{bpe}| - |D_{vq}|) \times (-p_{\min} \log(p_{\min})).$$

- For typical configuration: ~0.35% information loss, acceptable trade-off.

Our Method

Algorithm 4.1 BPE Image Tokenizer training procedure.

```
1: Input  $v_0, m, D$ .                                ▷  $v_0$ : initial vocab size,  $m$ : new vocab size,  $D$ : training data
2:  $v \leftarrow v_0$                                     ▷  $v$ : current vocab size
3:  $A \leftarrow \text{zeros}(v \times v)$                         ▷  $A$ : adjacency matrix
4:  $V \leftarrow \emptyset$                                 ▷  $V$ : extended vocabulary
5: for  $i \leftarrow 1$  to  $m$  do
6:    $A \leftarrow \text{UpdateMatrix}(D)$ 
7:    $(p, f) \leftarrow \text{MaxFreqPair}(A)$                 ▷  $p$ : best pair,  $f$ : frequency
8:   if  $f = 0$  then break
9:   end if
10:   $V \leftarrow V \cup \{(p, v)\}$ 
11:   $D' \leftarrow \emptyset$ 
12:  for each  $d \in D$  do
13:     $d' \leftarrow \text{Replace } p \text{ with } v \text{ in } d$ 
14:     $D' \leftarrow D' \cup \{d'\}$ 
15:  end for
16:   $D \leftarrow D'$ 
17:   $v \leftarrow v + 1$                                 ▷ set next id for new token
18: end for
19: return  $V$ 
```

Training Pipeline

- **Base Model:** Llama-3.1-8B
- **Two-Stage Training Process:**
 - **PT (Pretraining):**
 - Freeze original text embeddings
 - Train only visual embeddings
 - 595K images (CC-3M) + 558K (LCS)
 - **SFT (Supervised Fine-Tuning):**
 - Unfreeze all weights
 - 1.27M entries from LLaVA-OneVision Dataset
- **Key Difference:** Direct fusion of image modality without separate encoders

Experiment Results

	Training type	VQAv2	MMBench	MME ^p	MME ^c	POPE	VizWiz
LLM+VQ	SFT	51.1	35.9	972.3	231.8	73.8	43.1
	PT(full)+SFT	53.7	37.0	1037.2	261.4	75.3	44.2
	PT(freeze)+SFT	55.4	37.6	1054.5	277.0	76.0	45.3
LLM+VQ+BPE (Being-VL-0)	SFT	52.2	35.4	1029.7	269.6	76.3	45.3
	PT(full)+SFT	56.5	38.6	1144.6	284.3	77.3	45.8
	PT(freeze)+SFT	57.1	40.9	1223.5	307.1	79.0	46.0
Additional scaling (PT)	+RefCOCO(50.6K)	58.6	42.3	1257.4	314.3	79.8	47.1
	+AOKVQA (66.2K)	59.6	43.1	1288.1	321.4	80.4	47.5
Additional scaling (SFT)	+ShareGPT4o (57.3K)	60.2	43.7	1304.5	327.7	80.9	47.8
	+ALLaVA Inst (70K)	60.6	44.0	1316.2	331.0	81.3	48.2

- BPE Image Tokenizer consistently outperforms direct VQ approach
- Performance improves with vocabulary size up to 8K (balanced utilization)
- Model shows strong category-specific improvements (Existence: 145.0 vs 113.3)

Summary & Future Work

- First explicit tokenization of multimodal data like text-only LLMs
- Theoretical analysis validates benefits of tokenization for 2D data
- Significant performance gains with limited training data (0.1% of typical CLIP training)
- Future: Scale up data, apply to video, explore more sophisticated tokenization

Thanks!