# Where is the signal in tokenization space?

Renato Geh, Honghua Zhang, Kareem Ahmed,
Benjie Wang, Guy Van den Broeck

**University of California, Los Angeles**

UCLA

ST★R
AI
RESEARCH LAB
UCLA

# Tokenization

Most language models represent distributions over sequences of *tokens* (subwords), not strings.

$$\text{string} \quad \mathbf{x} = (x_1, x_2, \ldots, x_n)$$
$$\text{tokenization} \quad \mathbf{v} = (v_1, \ldots, v_m)$$

For example:

$$\text{string} \quad \mathbf{x} = \texttt{Caterpillar}$$
$$\text{tokenization} \quad \mathbf{v} = \texttt{[C,ater,p,ill,ar]} \equiv \texttt{[315,1008,29886,453,279]}$$

# Canonical Tokenization

How do we tokenize? There is usually a unique *canonical* tokenization:

$$\textbf{string} \quad \textbf{x} = \texttt{Caterpillar}$$
$$\textbf{canonical} \quad \textbf{v} = \texttt{[C,ater,p,ill,ar]}$$

(Llama 2)

Common assumption:

$$p(\textbf{x}) = p(\textbf{v}) \qquad \textcolor{red}{\times}$$

A string can be tokenized in an exponential number of ways (784 here!)

$$\texttt{[C,ater,pi,l,lar]}, \texttt{[Cat,er,pi,lla,r]}, \texttt{[Cat,er,pi,l,lar]},$$
$$\texttt{[Ca,ter,p,ill,ar]}, \texttt{[Ca,ter,p,illa,r]}, \texttt{[Cat,er,pi,ll,ar]}, \quad \text{(Llama 2)}$$
$$\ldots$$
$$\texttt{[Ca,t,e,r,p,i,l,l,a,r]}, \texttt{[C,a,t,e,r,p,i,l,l,a,r]}$$

# Tokenization

Why does this tokenization problem matter?

$$\textbf{string} \quad \mathbf{x} = \texttt{Hypnopaturist}$$

$$\textbf{canonical} \quad \mathbf{v} = \texttt{[Hyp,nop,atu,rist]}$$

$$\textbf{most likely} \quad \mathbf{v} = \texttt{[Hyp,no,patu,rist]}$$

$$\textbf{canonical prob} \quad p(\mathbf{v}|\mathbf{x}) \approx 0.0004$$

$$\textbf{most likely prob} \quad p(\mathbf{v}|\mathbf{x}) \approx 0.9948$$

(Gemma 2B)

**How canonical are unconditional samples?**

Less likely for non-English (code, unicode characters, etc)



**We're ignoring an exponential number of tokenizations!**

# Tokenization is a Neurosymbolic Problem!

↪ Tokens are symbols.
↪ A tokenization of a text is a constraint over these symbols.

$$p(\mathbf{v}, \mathbf{x}) = \begin{cases} p_{\mathrm{LLM}}(\mathbf{v}) & \text{if } \mathbf{v} \models \mathbf{x}; \\ 0 & \text{otherwise.} \end{cases}$$

$$\mathbf{v} = (v_1, v_2, \ldots, v_n) \models \mathbf{x} \Leftrightarrow v_1 \circ v_2 \circ \cdots \circ v_n = \mathbf{x}$$

concatenation

Example:

$$p(\mathbf{v} = [\text{一} \ \text{ラ}, \text{ク}] | \mathbf{x} = \text{一ラク}) = 0.586 \qquad p(\mathbf{v} = [\text{一}, \text{ラク}] | \mathbf{x} = \text{一ラク}) = 0.402$$

$$p(\mathbf{v} = [\text{一}, \text{ラ}, \text{ク}] | \mathbf{x} = \text{一ラク}) = 0.012 \qquad p(\mathbf{v} = [\text{Tok}, \text{ens}] | \mathbf{x} = \text{一ラク}) = 0$$

# Reasoning in Tokenization Space

Instead of the *canonical* tokenization, we might want to compute:

1. The most likely tokenization ❌

$$\arg\max_{\mathbf{v}\models\mathbf{x}} p(\mathbf{v}, \mathbf{x})$$

**Theorem.** *The most likely tokenization problem is NP-hard.*

For autoregressive models, e.g. transformers and state space models

2. The true probability of a text ❌

$$p(\mathbf{x}) = \sum_{\mathbf{v}\models\mathbf{x}} p(\mathbf{v}, \mathbf{x})$$

**Theorem.** *The marginal string probability problem is #P-hard.*

# (Approximate) Reasoning in Tokenization Space

## 1. The most likely tokenization

$$\arg\max_{\mathbf{v} \models \mathbf{x}} p(\mathbf{v}, \mathbf{x})$$

Branch-and-bound

↪ Lower bound: canonical likelihood

↪ Anytime: candidate at least as good as canonical

What did we learn?

↪ Runtime exponential on string length!

↪ Canonical best candidate for almost all cases…

…not always!

$p(\mathbf{v} = [\texttt{\_tongue,less}]|\mathbf{x} = \texttt{\_tongueless}) = 0.518$ → most likely tokenization

$p(\mathbf{v} = [\texttt{\_t,ong,uel,ess}]|\mathbf{x} = \texttt{\_tongueless}) = 0.004$

$p(\mathbf{v} = [\texttt{\_tong,uel,ess}]|\mathbf{x} = \texttt{\_tongueless}) = 0.474$

canonical tokenization

$p(\mathbf{v} = [\texttt{\_,HEADER,\_,DELIM,ITER}]|\mathbf{x} = \texttt{\_HEADER\_DELIMITER}) = 0.412$

$p(\mathbf{v} = [\texttt{\_HEAD,ER,\_,DELIM,ITER}]|\mathbf{x} = \texttt{\_HEADER\_DELIMITER}) = 0.330$

$p(\mathbf{v} = [\texttt{\_HEADER,\_,DELIM,ITER}]|\mathbf{x} = \texttt{\_HEADER\_DELIMITER}) = 0.010$

canonical tokenization

(Gemma 2B)

# (Approximate) Reasoning in Tokenization Space

2. The true probability of a text

$$p(\mathbf{x}) = \sum_{\mathbf{v} \models \mathbf{x}} p(\mathbf{v}, \mathbf{x})$$

Sequential importance sampling

Unbiased estimator converging to the true probability of text as #samples grows

$$p(\mathbf{x}) = \mathbb{E}_{\mathbf{v} \sim q(\mathbf{v}|\mathbf{x})} \left[ \frac{p(\mathbf{v}, \mathbf{x})}{q(\mathbf{v}|\mathbf{x})} \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^{N} \frac{p\left(\mathbf{v}^{(i)}, \mathbf{x}\right)}{q\left(\mathbf{v}^{(i)}|\mathbf{x}\right)}$$

proposal distribution

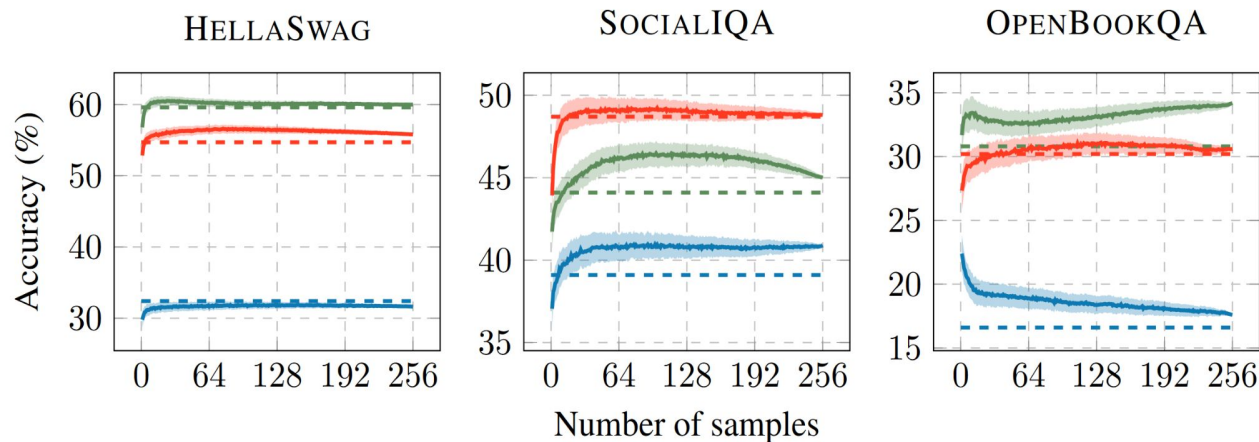$$q(v_j | \mathbf{v}_{1:j-1} = [\texttt{Tok,eni}], \mathbf{x} = \texttt{Tokenization}) = \begin{cases} 0.50 & , \text{if } v_j = \texttt{zation}; \\ 0.30 & , \text{if } v_j = \texttt{zat}; \\ 0.15 & , \text{if } v_j = \texttt{za}; \\ 0.05 & , \text{if } v_j = \texttt{z}; \\ 0.00 & , \text{if } v_j = \texttt{a}; \\ \vdots \\ 0.00 & , \text{if } v_j = \texttt{zzz}; \end{cases}$$

zero-out next tokens inconsistent with constraint

# Where is the signal in tokenization space?

$$\arg\max_{\text{answer}} \sum_{\mathbf{v} \models \text{answer}} p(\mathbf{v}, \text{answer} | \mathbf{v}_{\text{question}})$$



California experiences heavy earthquake activity due to
(a)   erosion
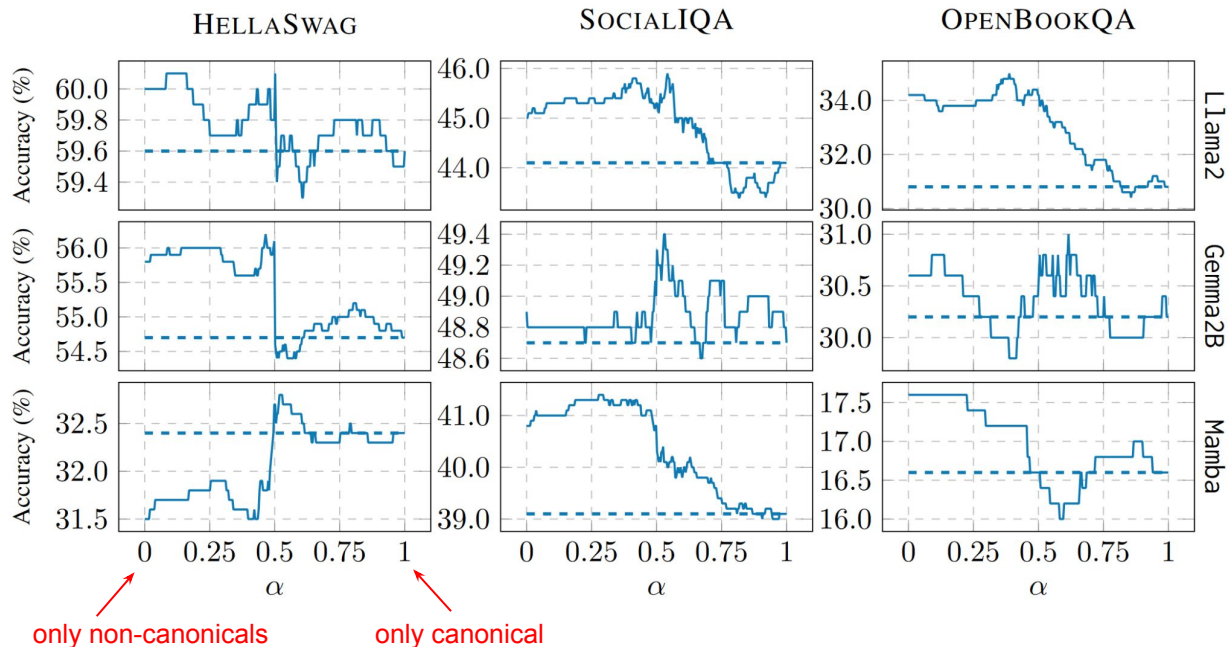(b)   techtonics
(c)   volcanic activity
(d)   fire

**There is signal in non-canonical tokenizations!**

# Mixtures of tokenizations can boost LLM accuracy!

Can we quantify how much signal is in non-canonical tokenizations?

$$\underset{\text{answer}}{\arg\max}\ \alpha \cdot \underbrace{p(\mathbf{v}_{\text{answer}}|\mathbf{v}_{\text{question}})}_{\text{canonical}} + (1 - \alpha) \cdot \underbrace{p(\text{noncanonical}|\mathbf{v}_{\text{question}})}_{\text{non-canonicals}}$$



Tune for α

|  | MIXTURE | CANONICAL |  |
|---|---|---|---|
|  | \multicolumn{2}{c}{Accuracy (%)} |  |  |
| Llama2 | **59.7** | 59.6 | HELLASWAG |
| Gemma | **55.8** | 54.7 | |
| Mamba | 31.6 | **32.4** | |
| Llama2 | **44.8** | 44.1 | SOCIALIQA |
| Gemma | **48.8** | 48.7 | |
| Mamba | **39.8** | 39.1 | |
| Llama2 | **34.0** | 30.8 | OPENBOOKQA |
| Gemma | **30.6** | 30.2 | |
| Mamba | **17.6** | 16.6 | |

**Consistent improvement!**

only non-canonicals    only canonical

# Main Takeaways

**Probabilistic reasoning is hard**
- ✗   Computing the most likely tokenization (exactly) is **hard**
- ✗   Computing the true text probability (exactly) is **hard**

**Non-canonical tokenizations appear in the wild**
- ✔   LLMs sample non-canonical tokenizations
- ✔   Non-canonical tokenizations can be more likely

**Non-canonical tokenizations matter**
- ✔   Mixtures of canonical and non-canonical boost performance
- ✔   More inference time compute, better performance

# Where is the signal in tokenization space?

Renato Geh, Honghua Zhang, Kareem Ahmed,
Benjie Wang, Guy Van den Broeck

**University of California, Los Angeles**

UCLA

ST★R
AI
RESEARCH LAB
UCLA