

(TL;DR) Does the signal come from only one source, i.e. the so-called **canonical** (default) **tokenization**?

The answer is **no!**

By looking at **non-canonical tokenizations**, we get **consistent improvement** in downstream performance! 📈

(1)

Given a **string**, what's $p(\mathbf{x})$ under the LLM?

string \mathbf{x} = Caterpillar
tokenization \mathbf{v} = [C,ater,p,ill,ar]

Common assumption:

$$p(\mathbf{x}) = p(\mathbf{v}) \quad \times$$

What about *other* tokenizations?

[C,ater,pi,l,lar], [Cat,er,pi,lla,r],
..., [C,at,e,r,p,i,l,l,a,r]

are all valid tokenizations!

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

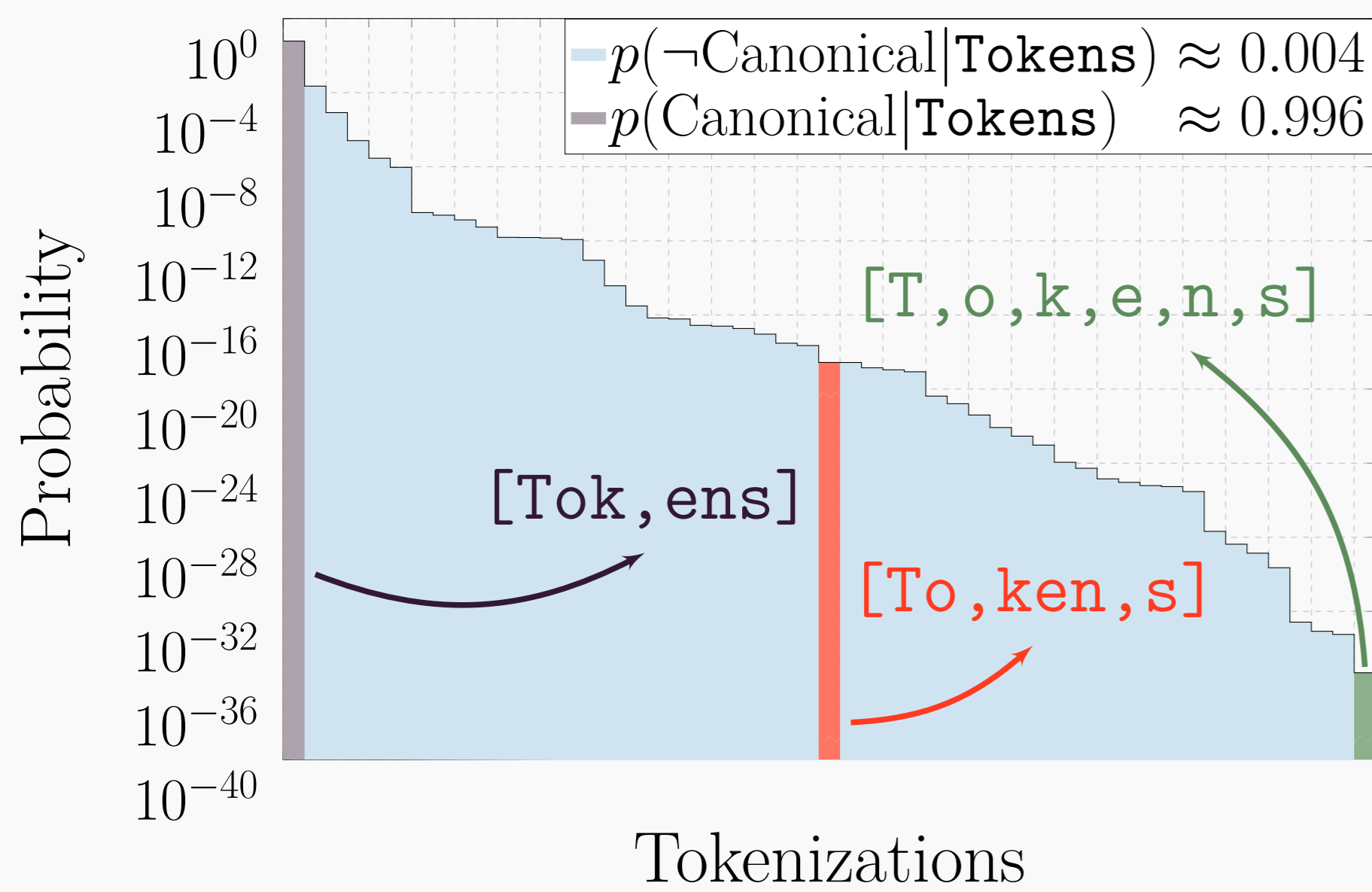
We should not neglect other tokenizations...

(4)

E.g., exactly compute the **most likely tokenization** for autoregressive models?

Answer: **No!** 😞

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



Yet experimentally, **canonical is often much more likely** in English.

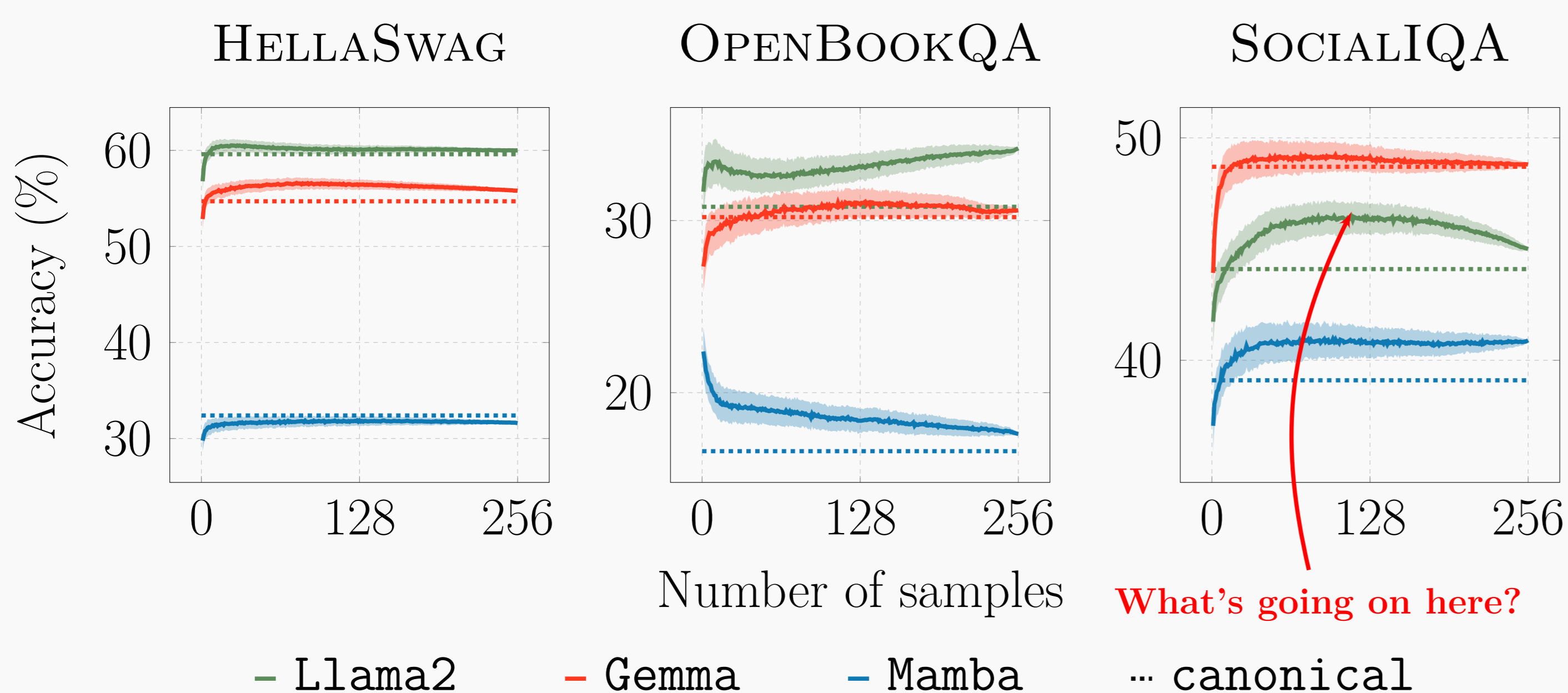
(7)

So can we aggregate over all tokenizations and exactly compute the marginal $p(\mathbf{x}) = \sum_{\mathbf{v}} p(\mathbf{v})$?

Answer: **No!** 😞

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

But we *can* approximate through **sequential importance sampling!**



That's evidence of signal in non-canonical tokenizations!

(2)

After all, some non-canonical tokenizations can have **non-negligible mass!**

string \mathbf{x} = Hypnopaturist
canonical \mathbf{v} = [Hyp,nop,atu,rist]
most likely \mathbf{v} = [Hyp,no,patu,rist]
canonical prob $p(\mathbf{v}|\mathbf{x}) \approx 0.0004$
most likely prob $p(\mathbf{v}|\mathbf{x}) \approx 0.9948$

Here, Gemma's **canonical tokenization** \mathbf{v} of \mathbf{x} is *much less likely* compared to the **most likely** (non-canonical) **tokenization** \mathbf{v} .

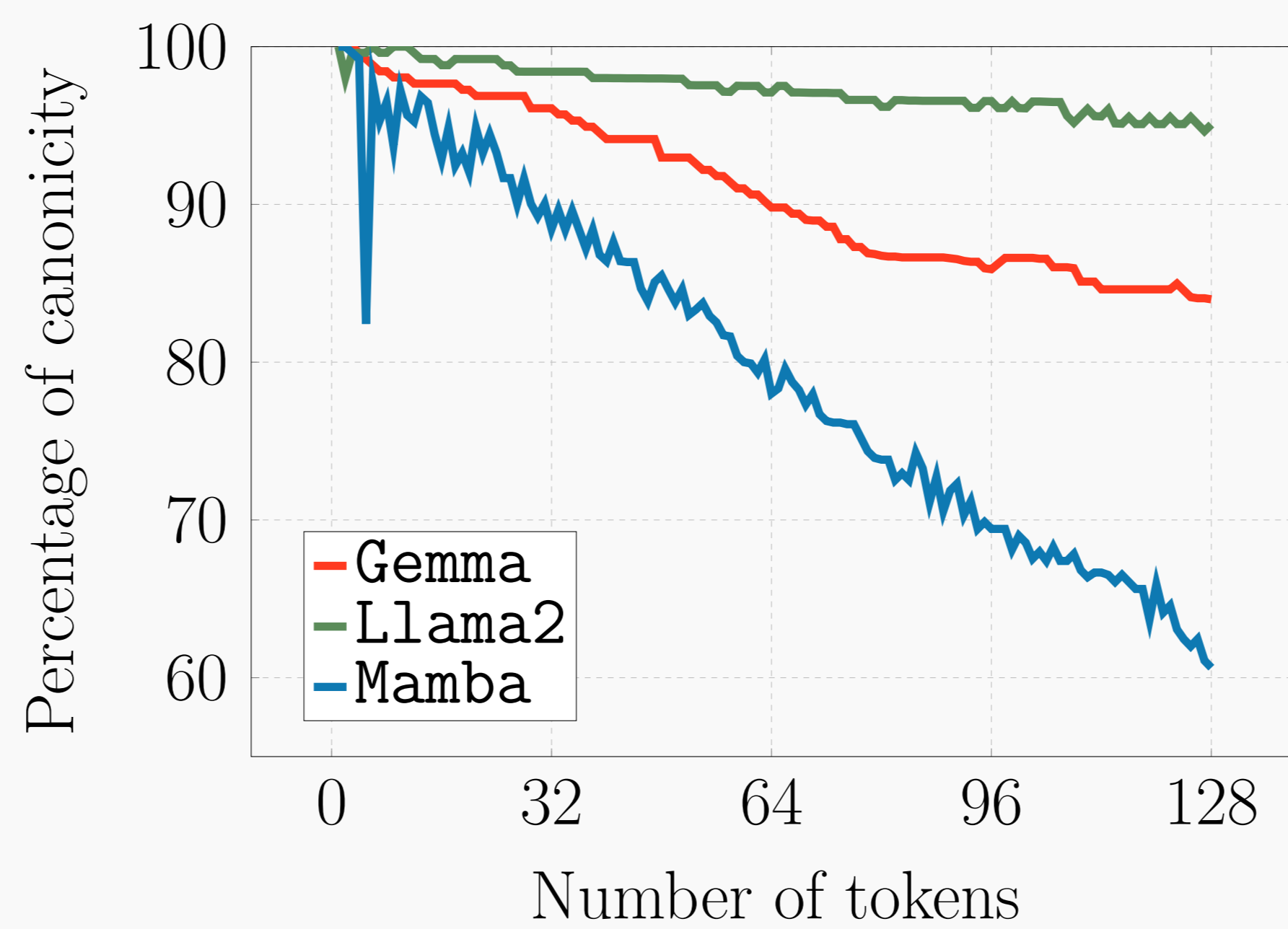
Read our paper!



<https://arxiv.org/abs/2408.08541>

(5)

Despite this, **sampling unconditionally** from the LLM reveals that...

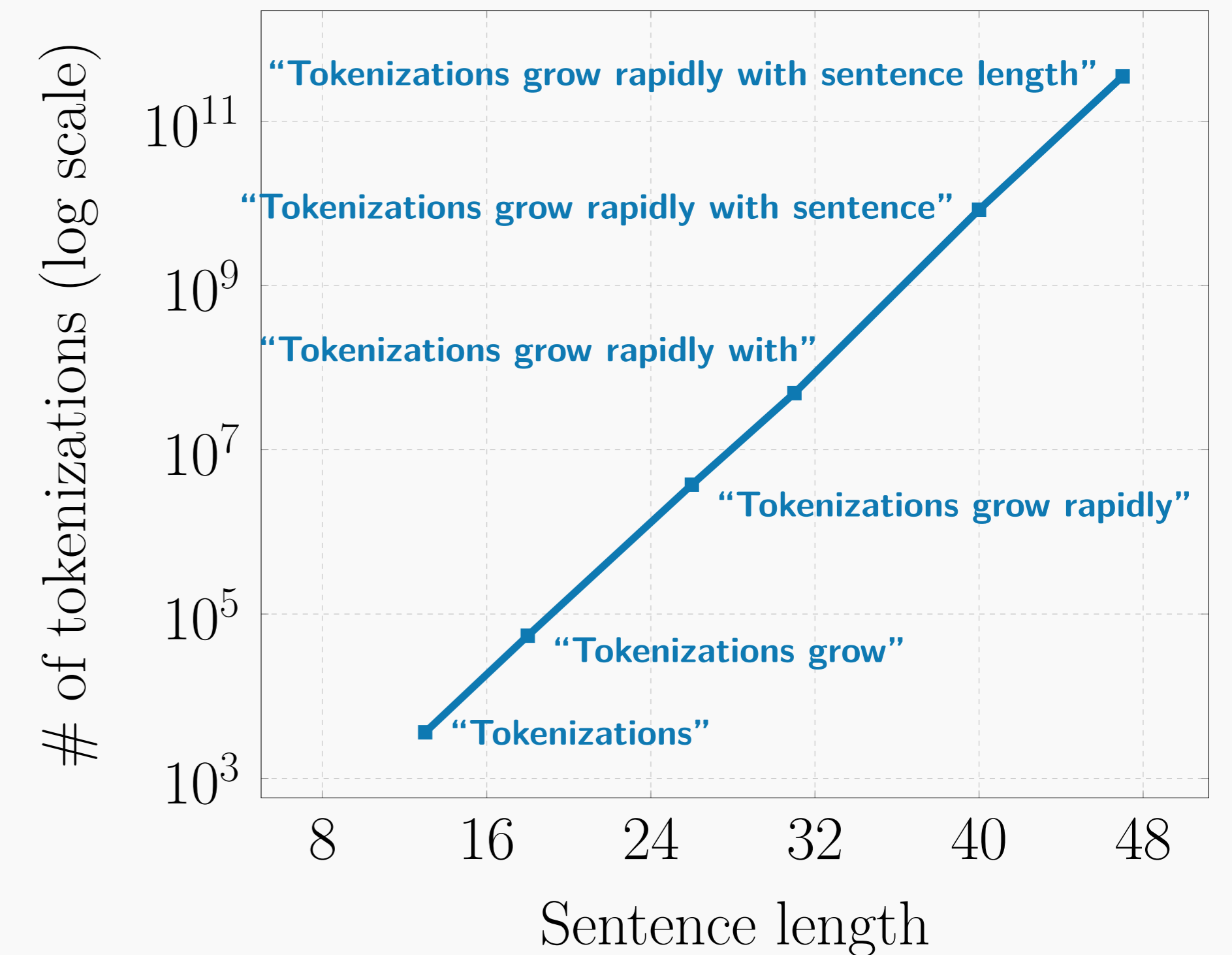


...they are generating **non-canonically for long texts!**

This is especially true for generated non-English (e.g. code, unicode characters, etc.)

(3)

...but the space of tokenizations is **exponential!**



So can we **reason probabilistically** about this tokenization space?

(6)

These generated **non-canonical tokenizations** can be more likely than **canonical!**

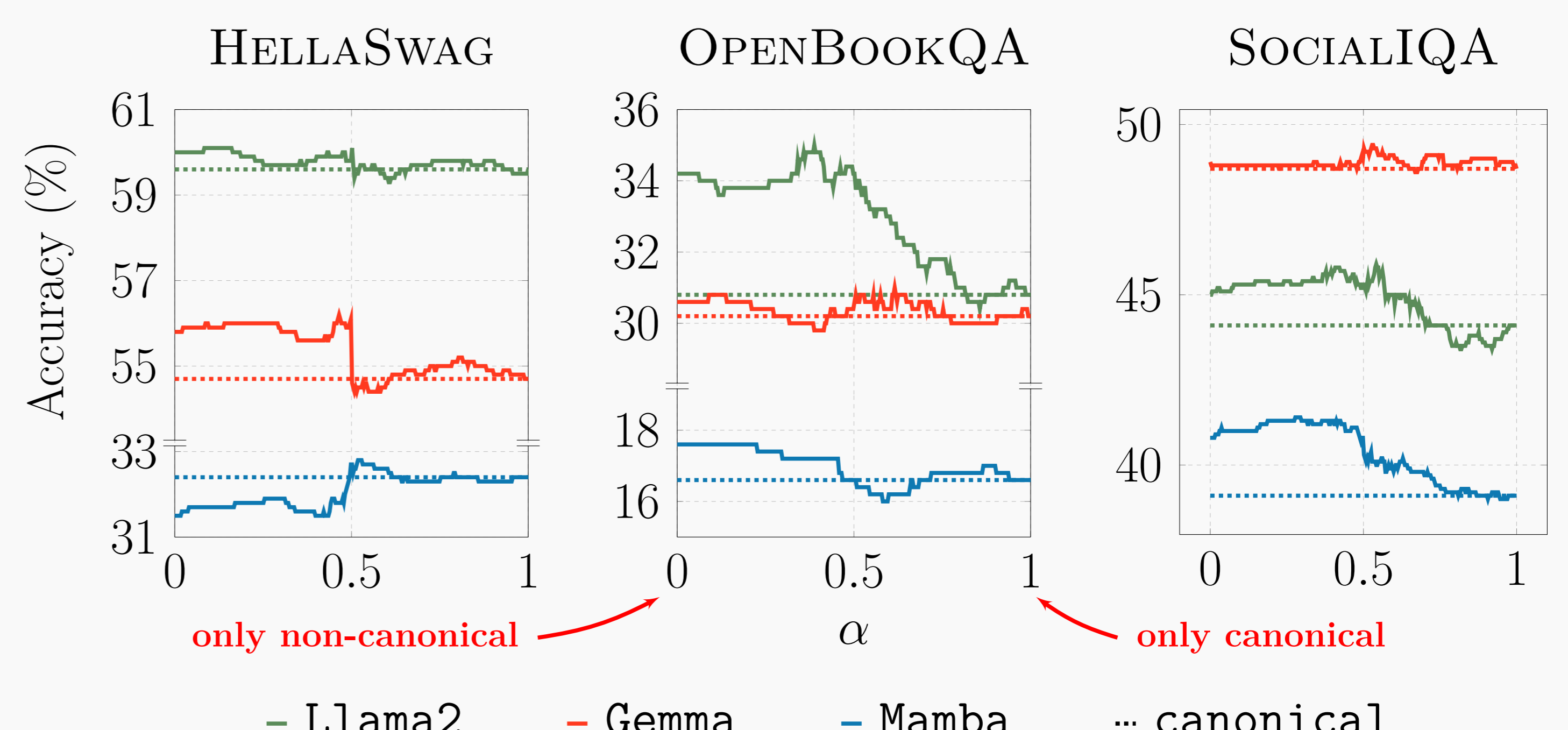
\mathbf{x} = ーラク
 $p(\mathbf{v} = [ー,ラ,ク]|\mathbf{x}) = 0.586$
 $p(\mathbf{v} = [ー,ラ,ク]|\mathbf{x}) = 0.012$
 $p(\mathbf{v} = [ー,ラク]|\mathbf{x}) = 0.402$
 \mathbf{x} = 舌less
 $p(\mathbf{v} = [舌,less]|\mathbf{x}) = 0.518$
 $p(\mathbf{v} = [舌,less]|\mathbf{x}) = 0.004$
 $p(\mathbf{v} = [舌,less]|\mathbf{x}) = 0.474$
 \mathbf{x} = _HEADER_DELIMITER
 $p(\mathbf{v} = [_ ,HEADER, _ ,DELIM, ITER]|\mathbf{x}) = 0.412$
 $p(\mathbf{v} = [_ ,HEADER, _ ,DELIM, ITER]|\mathbf{x}) = 0.330$
 $p(\mathbf{v} = [_ ,HEADER, _ ,DELIM, ITER]|\mathbf{x}) = 0.010$

Meaning there is possibly signal in non-canonical tokenizations!

(8)

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

$$\arg \max_{\mathbf{x}} \alpha \cdot \underbrace{p(\mathbf{v}, \mathbf{x} | \mathbf{v}_q)}_{\text{canonical}} + (1 - \alpha) \cdot \underbrace{p(-\mathbf{v}, \mathbf{x} | \mathbf{v}_q)}_{\text{non-canonical}}$$



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