

A Non-monotonic Self-terminating Language Model

Eugene Choi
New York University

Kyunghyun Cho
New York University
Prescient Design, Genentech
CIFAR Fellow

Cheolhyoung Lee
New York University



ML² Machine Learning
for Language



Prescient
Design
A Genentech Accelerator

CIFAR

NLP Progress:

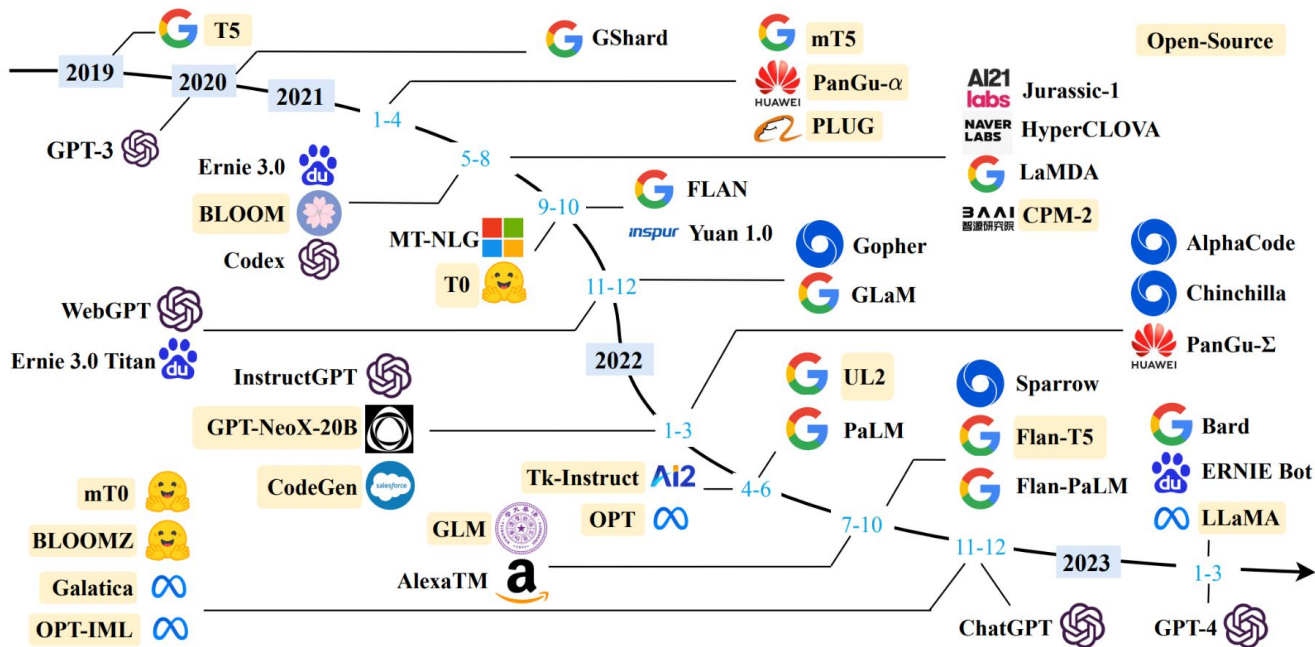


Fig. 1. A timeline of existing large language models (having a size larger than 10B) in recent years. We mark the open-source LLMs in yellow color.

(Fig. from Zhao et. al. 2023)

Identifying the Cause: **Incomplete Decoding**

- A decoding algorithm \mathcal{S} is *incomplete* if there exists $\emptyset \subsetneq \mathcal{V}_t \subsetneq \mathcal{V}$ such that $\sum_{v \in \mathcal{V}_t} q_{\mathcal{S}(p_\theta)}(y_t = v | \mathbf{y}_{<t}, \mathbf{x}) = 1$.
- Examples of incomplete decoding algorithms:
 - Greedy decoding
 - Top- k sampling
 - Nucleus sampling
 - Beam search

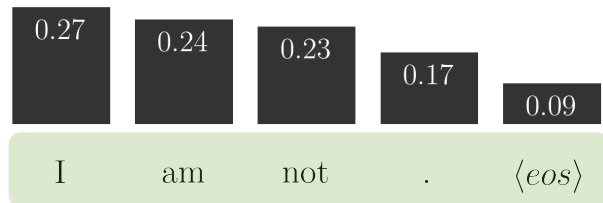
Identifying the Cause: **Incomplete Decoding**

- A decoding algorithm \mathcal{S} is *incomplete* if there exists $\emptyset \subsetneq \mathcal{V}_t \subsetneq \mathcal{V}$ such that $\sum_{v \in \mathcal{V}_t} q_{\mathcal{S}(p_\theta)}(y_t = v | \mathbf{y}_{<t}, \mathbf{x}) = 1$.
- Examples of incomplete decoding algorithms:
 - Greedy decoding
 - Top- k sampling
 - Nucleus sampling
 - Beam search

Identifying the Cause: Incomplete Decoding

- A decoding algorithm \mathcal{S} is *incomplete* if there exists $\emptyset \subsetneq \mathcal{V}_t \subsetneq \mathcal{V}$ such that $\sum_{v \in \mathcal{V}_t} q_{\mathcal{S}(p_\theta)}(y_t = v | \mathbf{y}_{<t}, \mathbf{x}) = 1$.
- Examples of incomplete decoding algorithms:
 - Greedy decoding
 - Top- k sampling
 - Nucleus sampling
 - Beam search

$p_\theta(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})$: model distribution

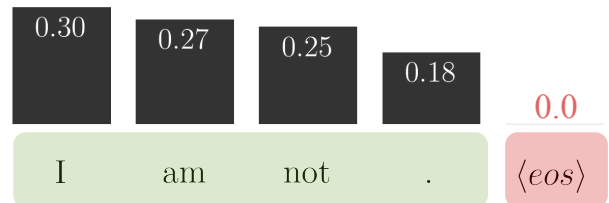


Properly terminating!

$\mathcal{S} : p_\theta \rightarrow q_{\mathcal{S}(p_\theta)}$

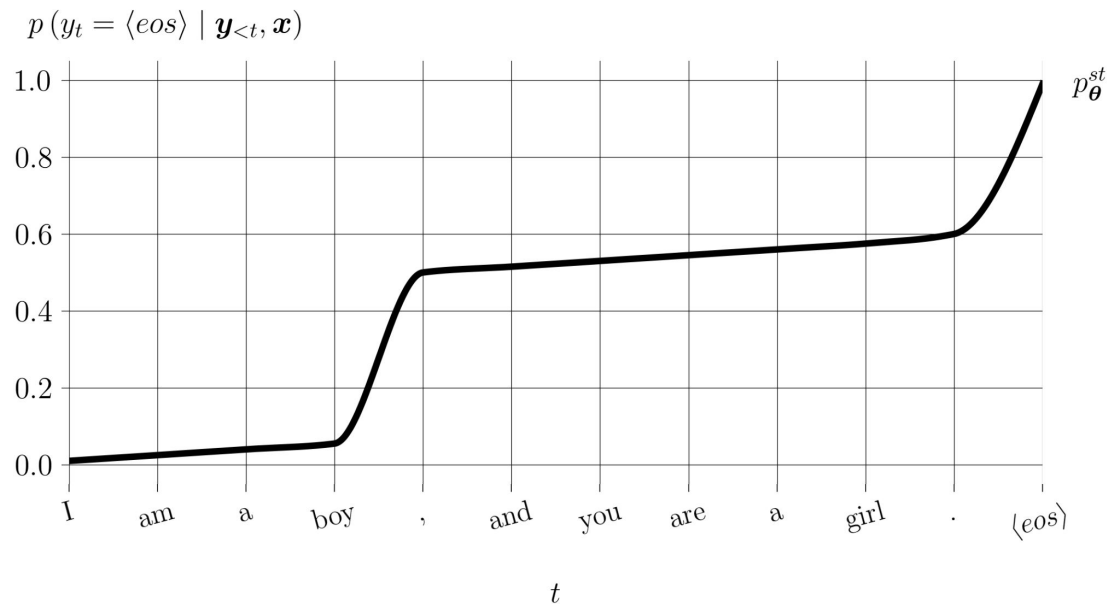
$\mathcal{S} = \text{Top-}k \text{ sampling,}$
with $k = 4$

$q_{\mathcal{S}(p_\theta)}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})$: induced distribution by \mathcal{S}



Non-terminating!

Self-terminating (ST+) Language Model (Welleck et al., 2020)



$$p_{\theta}^{st}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x}) \uparrow 1 \text{ as } t \rightarrow \infty$$

Is ST+ an Optimal LM Parametrization?

- Suppose there are two sequences in our dataset:

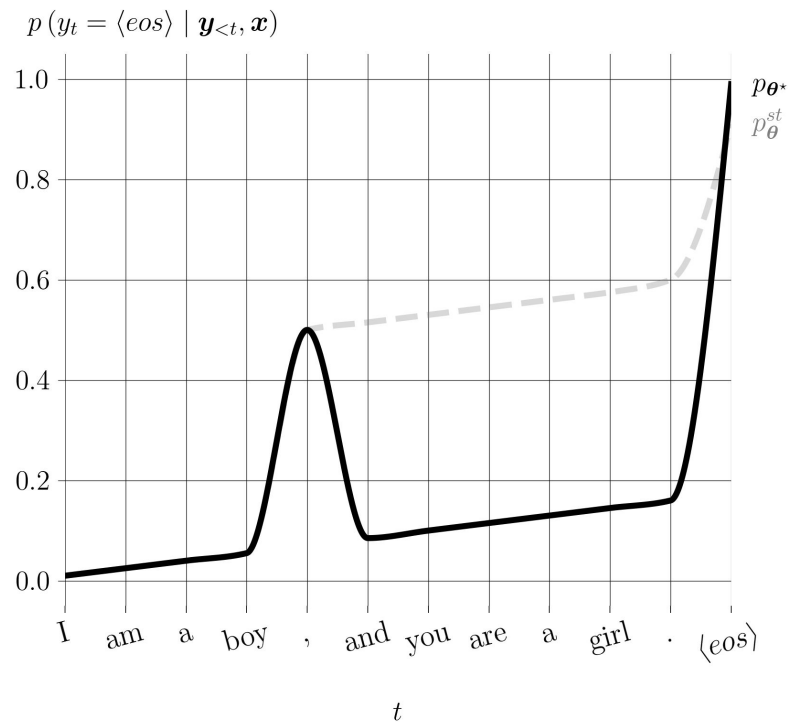
$$\mathcal{D} = \left\{ \begin{array}{l} \text{“I am a boy”} \\ \text{“I am a boy, and you are a girl.”} \end{array} \right\}$$

- Our language model trained on this dataset may or may not terminate after the former.
 - Once our model decides not to end, **it should dramatically reduce the termination probability** to continue:

$$p_{\theta^*}(y_t = \langle eos \rangle | \mathbf{y}_{<t} = \text{“I am a boy”}) \gg p_{\theta^*}(y_{t+1} = \langle eos \rangle | \mathbf{y}_{<t+1} = \text{“I am a boy,”})$$

- The ST+, which monotonically increase the termination probability, cannot capture such a case, where one sequence is a prefix of another.

Improving Self-terminating



$$p_{\theta}^{st}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x}) \uparrow 1 \text{ as } t \rightarrow \infty$$

Too strong!

$$p_{\theta^*}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x}) \rightarrow 1 \text{ as } t \rightarrow \infty$$

Enough for termination!

Non-monotonic Self-terminating Language Model (NMST)

⚠ *Avoiding non-termination requires $\lim_{t \rightarrow \infty} p_{\theta}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x}) = 1$.*

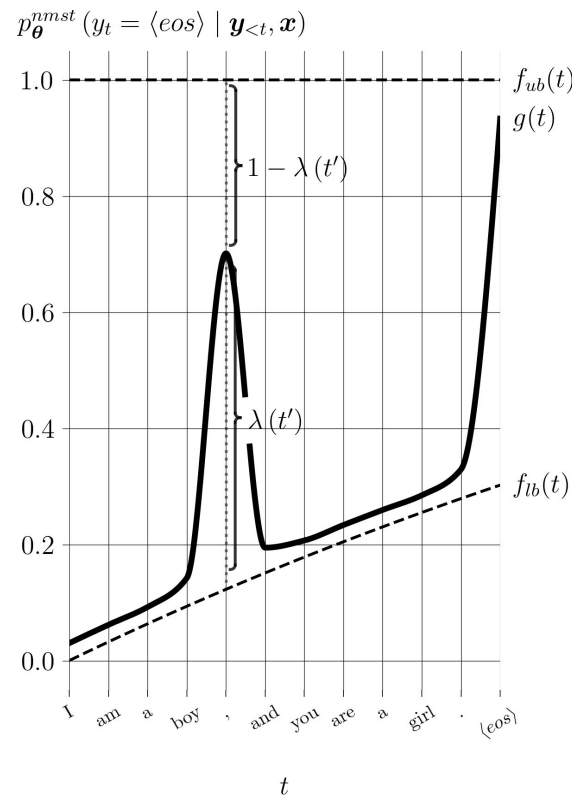
NMST permits a non-monotonic behavior, while preventing non-termination:

$$p_{\theta}^{nmst}(y_t | \mathbf{y}_{<t}, \mathbf{x}) = \begin{cases} \alpha_t, & \text{if } y_t = \langle eos \rangle, \\ (1 - \alpha_t) \cdot \frac{\exp(\mathbf{u}_v^{\top} \mathbf{h}_t)}{\sum_{v' \in \mathcal{V} \setminus \{\langle eos \rangle\}} \exp(\mathbf{u}_{v'}^{\top} \mathbf{h}_t)}, & \text{if } y_t = v \in \mathcal{V} \setminus \{\langle eos \rangle\}, \end{cases}$$

$$\text{where } \alpha_t = \underbrace{p_{\theta}^{nmst}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x})}_{g(t)} = \underbrace{(1 - \sigma(\mathbf{u}_{\langle eos \rangle}^{\top} \mathbf{h}_t))}_{(1 - \lambda(t))} \underbrace{(1 - (1 - \epsilon)^t)}_{f_{lb}(t)} + \underbrace{\sigma(\mathbf{u}_{\langle eos \rangle}^{\top} \mathbf{h}_t)}_{\lambda(t)}$$

and $\sigma(x)$ is a sigmoid function and $\epsilon \in (0, 1)$.

The figure on the right shows that $p_{\theta}^{nmst}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x})$ is modelled by taking a convex combination, $g(t) = (1 - \lambda(t))f_{lb}(t) + \lambda(t)f_{ub}(t)$, of a monotonically increasing lower bound-function, $f_{lb}(t)$, and a constant upper bound function, $f_{ub}(t) = 1$.



Non-monotonic Self-terminating Language Model (NMST)

⚠ *Avoiding non-termination requires* $\lim_{t \rightarrow \infty} p_{\theta}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x}) = 1$.

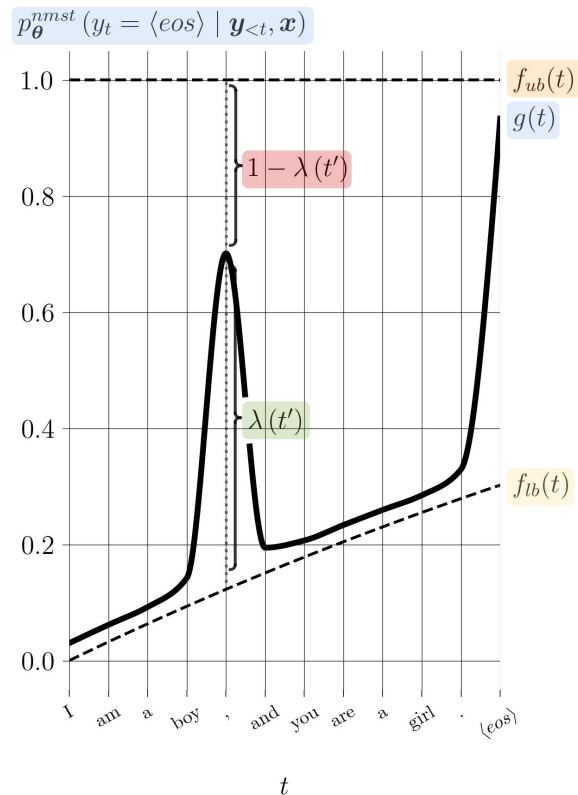
NMST permits a non-monotonic behavior, while preventing non-termination:

$$p_{\theta}^{nmst}(y_t | \mathbf{y}_{<t}, \mathbf{x}) = \begin{cases} \alpha_t, & \text{if } y_t = \langle eos \rangle, \\ (1 - \alpha_t) \cdot \frac{\exp(\mathbf{u}_v^{\top} \mathbf{h}_t)}{\sum_{v' \in \mathcal{V} \setminus \{\langle eos \rangle\}} \exp(\mathbf{u}_{v'}^{\top} \mathbf{h}_t)}, & \text{if } y_t = v \in \mathcal{V} \setminus \{\langle eos \rangle\}, \end{cases}$$

$$\text{where } \alpha_t = \underbrace{p_{\theta}^{nmst}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x})}_{g(t)} = \underbrace{(1 - \sigma(\mathbf{u}_{\langle eos \rangle}^{\top} \mathbf{h}_t))}_{(1-\lambda(t))} \underbrace{(1 - (1 - \epsilon)^t)}_{f_{lb}(t)} + \underbrace{\sigma(\mathbf{u}_{\langle eos \rangle}^{\top} \mathbf{h}_t)}_{\lambda(t)}$$

and $\sigma(x)$ is a sigmoid function and $\epsilon \in (0, 1)$.

The figure on the right shows that $p_{\theta}^{nmst}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x})$ is modelled by taking a convex combination, $g(t) = (1 - \lambda(t))f_{lb}(t) + \lambda(t)f_{ub}(t)$, of a monotonically increasing lower bound-function, $f_{lb}(t)$, and a constant upper bound function, $f_{ub}(t) = 1$.



Non-monotonic Self-terminating Language Model (NMST)

⚠ *Avoiding non-termination requires $\lim_{t \rightarrow \infty} p_{\theta}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x}) = 1$.*

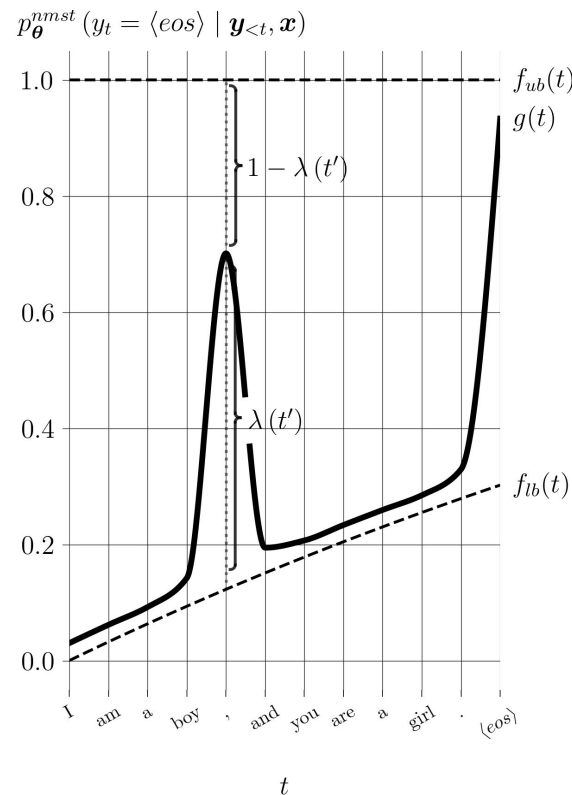
NMST permits a non-monotonic behavior, while preventing non-termination:

$$p_{\theta}^{nmst}(y_t | \mathbf{y}_{<t}, \mathbf{x}) = \begin{cases} \alpha_t, & \text{if } y_t = \langle eos \rangle, \\ (1 - \alpha_t) \cdot \frac{\exp(\mathbf{u}_v^{\top} \mathbf{h}_t)}{\sum_{v' \in \mathcal{V} \setminus \{\langle eos \rangle\}} \exp(\mathbf{u}_{v'}^{\top} \mathbf{h}_t)}, & \text{if } y_t = v \in \mathcal{V} \setminus \{\langle eos \rangle\}, \end{cases}$$

$$\text{where } \alpha_t = \underbrace{p_{\theta}^{nmst}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x})}_{g(t)} = \underbrace{(1 - \sigma(\mathbf{u}_{\langle eos \rangle}^{\top} \mathbf{h}_t))}_{(1-\lambda(t))} \underbrace{(1 - (1 - \epsilon)^t)}_{f_{lb}(t)} + \underbrace{\sigma(\mathbf{u}_{\langle eos \rangle}^{\top} \mathbf{h}_t)}_{\lambda(t)},$$

and $\sigma(x)$ is a sigmoid function and $\epsilon \in (0, 1)$.

The figure on the right shows that $p_{\theta}^{nmst}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x})$ is modelled by taking a convex combination, $g(t) = (1 - \lambda(t))f_{lb}(t) + \lambda(t)f_{ub}(t)$, of a monotonically increasing lower bound-function, $f_{lb}(t)$, and a constant upper bound function, $f_{ub}(t) = 1$.



Experiments

- Setup: Sequence completion

$$p_{\theta}(\underbrace{\mathbf{y} = \text{“and you are a girl. } \langle eos \rangle\text{”}}_{\text{continuation}} \mid \underbrace{\mathbf{x} = \text{“I am a boy,”}}_{\text{context}})$$

- Parameterizations:

- Vanilla Autoregressive LMs (VA+)
- (Monotonic) self-terminating recurrent LMs (ST+)
- Non-monotonic self-terminating LMs (NMST+)

- Architectures: RNN, LSTM, **GPT-2**

- Datasets: WikiText-2, **WikiText-103**

- Metrics:

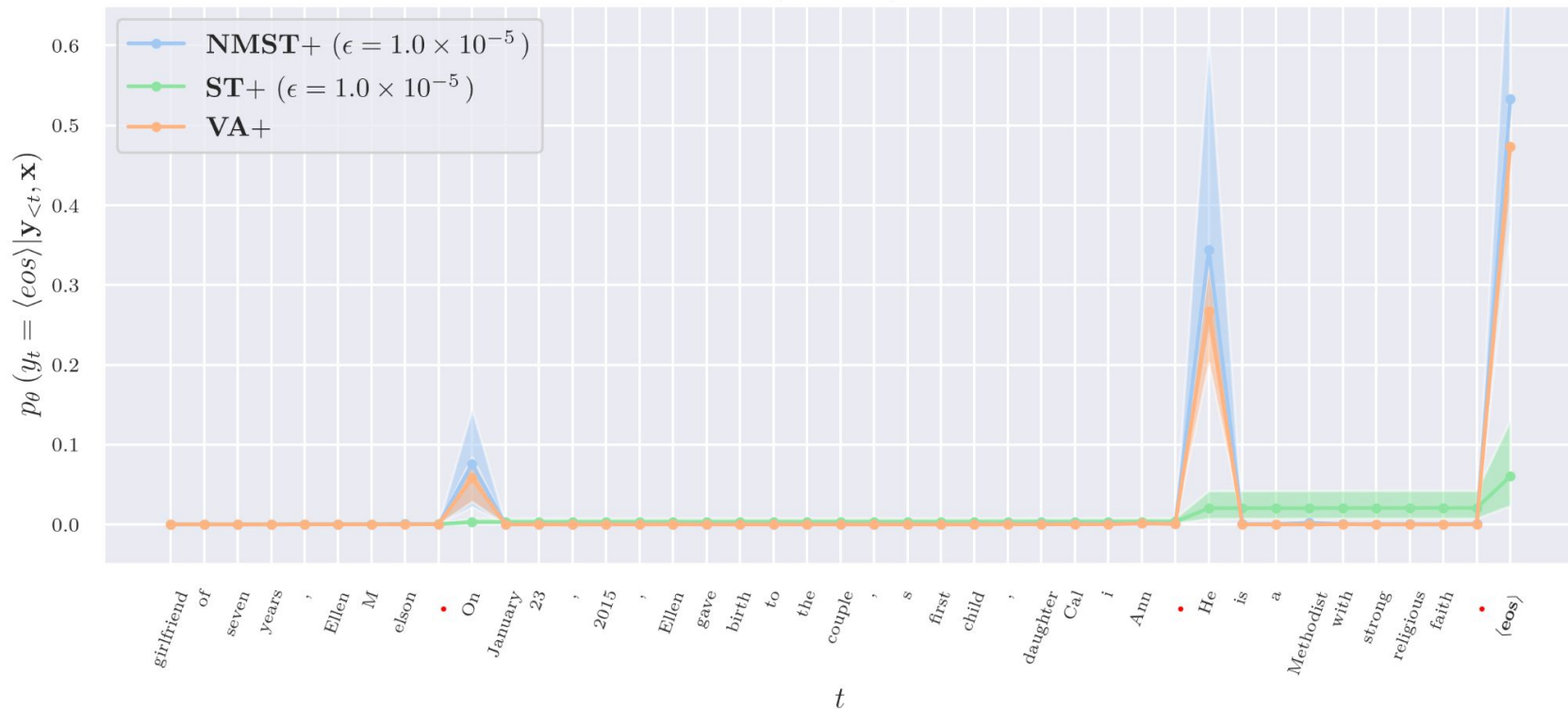
- Perplexity
- Non-termination Ratio (r_{nt}): To present the consistency of p_{θ} with respect to a given decoding algorithm \mathcal{S} , we need to compute $r_{nt} = q_{\mathcal{S}(p_{\theta})}(|\mathbf{y}| = \infty)$. Instead, based on

$$r_{nt} = q_{\mathcal{S}(p_{\theta})}(|\mathbf{y}| = \infty) = \lim_{L \rightarrow \infty} q_{\mathcal{S}(p_{\theta})}(|\mathbf{y}| > L), \quad (11)$$

we use $r_{nt}(L) = q_{\mathcal{S}(p_{\theta})}(|\mathbf{y}| > L)$ with a sufficiently large threshold L to estimate r_{nt} .

Result: $p_{\theta}(y_t = \langle eos \rangle | \mathbf{y}_{<t}, \mathbf{x})$ Plot (GPT-2 w/ WikiText-103)

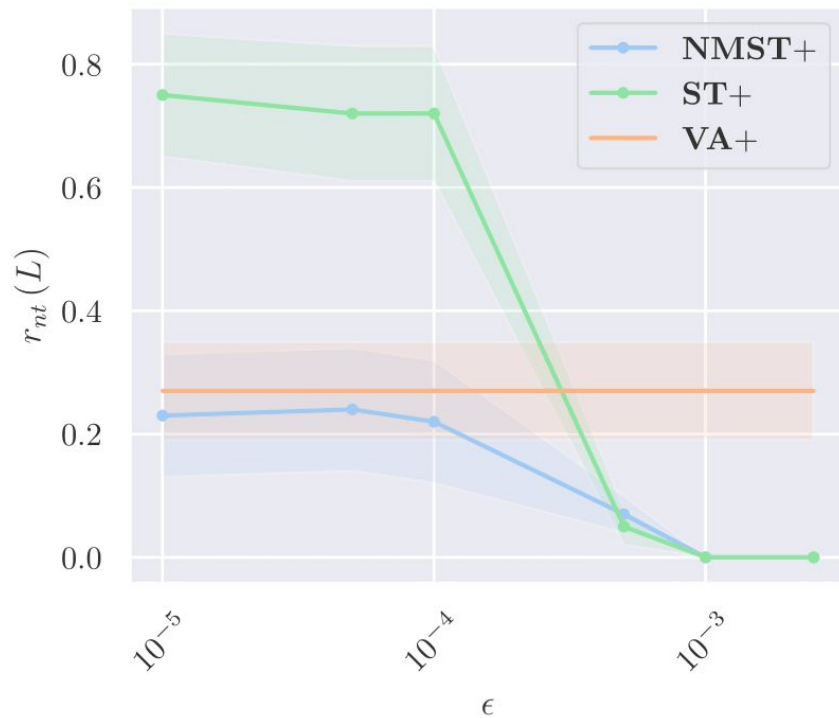
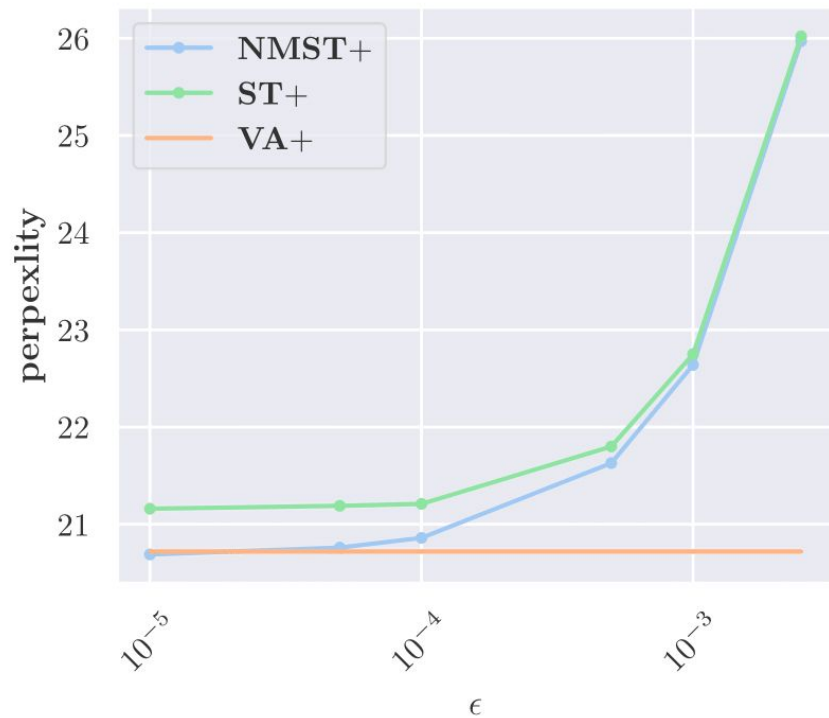
$\mathbf{x} =$ “ On December 4 , 2010 , Kershaw married his ”



Result: Completion Examples (GPT-2 w/ WikiText-103)

Context	<i>Made of concrete, steel, and wood, the</i>
VA+	building was built in the mid @-@ 19th century. It was the first building in the United States to be built in concrete, and the first to be built in wood. It was also the first building in the United States to be built in steel. It was the first building in the United States to be built in wood. It was ...
ST+	building is constructed of steel and concrete. The building's exterior is made of steel and concrete. The building's interior is made of wood, and the building's exterior is made of concrete. The building's exterior is made of concrete, and the building's interior is made of wood. The building's exterior is made of concrete, and the building's interior is made of wood. The building's ...
NMST+	building was designed by the architectural firm of Bowers & Wainwright, and was completed in 1892. The building is the largest of its kind in the United States. <eos>

Result: Perplexity-Consistency Tradeoff (GPT-2 w/ WikiText-103)



Conclusion:

We confirmed both theoretically and empirically that:

- Our use of incomplete decoding methods leads to generating a 0-probability sequence (e.g. ∞ -length sequence).
- Non-monotonic self-terminating language model can prevent non-termination, while maintaining generation quality.

Thank you for your attention!