

# Closing the Curious Case of Neural Text Degeneration

**Matthew Finlayson**<sup>1</sup> John Hewitt<sup>2</sup> Alexander Koller<sup>3</sup>  
Swabha Swayamdipta<sup>1</sup> Ashish Sabharwal<sup>4</sup>

<sup>1</sup>USC <sup>2</sup>Stanford <sup>3</sup>Saarland University <sup>4</sup>AI2

January 18, 2024

Accepted to ICLR 2024



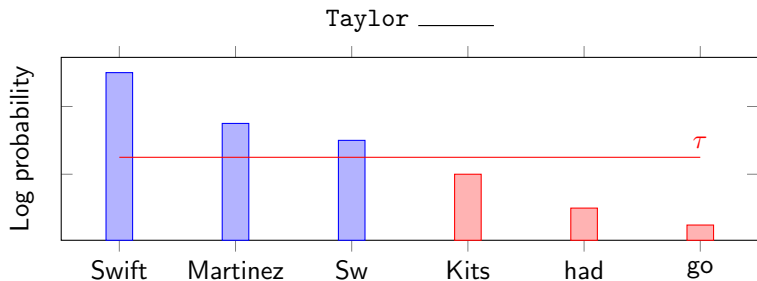
# Goal

- ▶ Improve threshold sampling methods (e.g., top- $k$ ) by directly addressing the source of errors.



# Threshold sampling

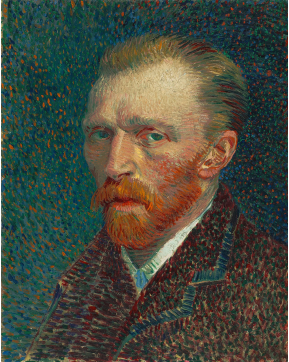
E.g., top- $k$ , top- $p$



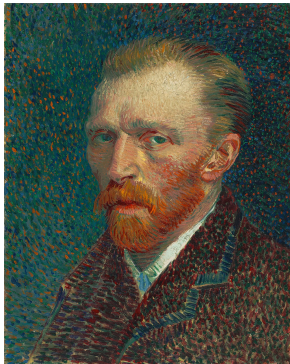
Choose a threshold  $\tau$  and only sample tokens with probability greater than  $\tau$ .



# Van Gogh



# Van Gogh





## The Art Institute of Chicago Recreates Van Gogh's Famous Bedroom to be Rented on Airbnb

FEBRUARY 9, 2016

KATE SIERZPUTOWSKI



# We don't know what Van Gogh's real bedroom looked like



What color is the floor? Is the towel brown? Is the pitcher glass?



# We don't know what Van Gogh's real bedroom looked like



What color is the floor? Is the towel brown? Is the pitcher glass?



*The Yellow House*, building destroyed June 25, 1944 by Allied bombing in France (WWII)



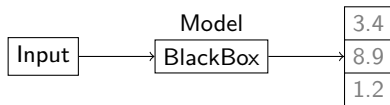






## What does this have to do with LM decoding?

A language model is like Van Gogh: it has a limited palette.

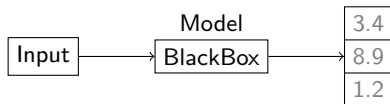


Contextualized embedding  $\mathbf{h} \in \mathbb{R}^d$



## What does this have to do with LM decoding?

A language model is like Van Gogh: it has a limited palette.



Contextualized embedding  $\mathbf{h} \in \mathbb{R}^d$

The size  $d$  of the embedding means the model paints with  $d$  colors.

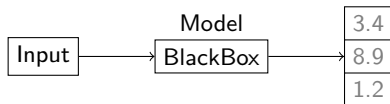
8.8	4.1	2.3	1.2	6.3	8.5
0.3	2.4	4.1	4.3	0.5	4.2
1.5	2.4	2.4	9.3	3.9	8.4

Softmax matrix  $\mathbf{W} \in \mathbb{R}^{v \times d}$



# What does this have to do with LM decoding?

A language model is like Van Gogh: it has a limited palette.



Contextualized embedding  $\mathbf{h} \in \mathbb{R}^d$

The size  $d$  of the embedding means the model paints with  $d$  colors.

8.8	4.1	2.3	1.2	6.3	8.5
0.3	2.4	4.1	4.3	0.5	4.2
1.5	2.4	2.4	9.3	3.9	8.4

Softmax matrix  $\mathbf{W} \in \mathbb{R}^{v \times d}$

The embedding  $\mathbf{h}$  specifies the combination of rows of  $\mathbf{W}$  to create the logits  $\mathbf{W}^T \mathbf{h}$ .

$\mathbf{W}^T$        $\mathbf{h}$       =       $\mathbf{W}^T \mathbf{h} \in \mathbb{R}^v$



LM logits define a distribution over vocab items

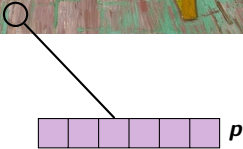
Martinez 0.2
Swift 0.1
had 0.1
Taylor 0.4
Sw 0.1
go 0.1

$$\mathbf{p} = \text{softmax}(\mathbf{W}^T \mathbf{h}) \in \mathbb{R}^v$$



# Models cannot output arbitrary distributions

Van Gogh's palette could not reproduce all the colors in the room:  
LMs cannot output arbitrary distributions over tokens.



This is known as the “softmax bottleneck”.





- ▶ The designers reverse engineered (kind of) the “true” colors.
- ▶ We want to reverse engineer the “true” distribution.



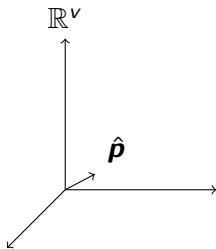


What are next-token distributions ( $\rho$ )?



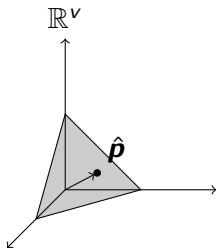
# What are next-token distributions ( $\boldsymbol{p}$ )?

- ▶ Distributions over  $v$  items are vectors in  $\mathbb{R}^v$ .



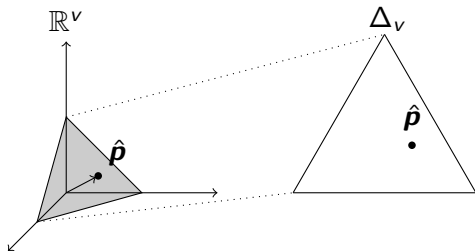
# What are next-token distributions ( $\rho$ )?

- ▶ Distributions over  $v$  items are vectors in  $\mathbb{R}^v$ .
- ▶ All distributions must sum to one.



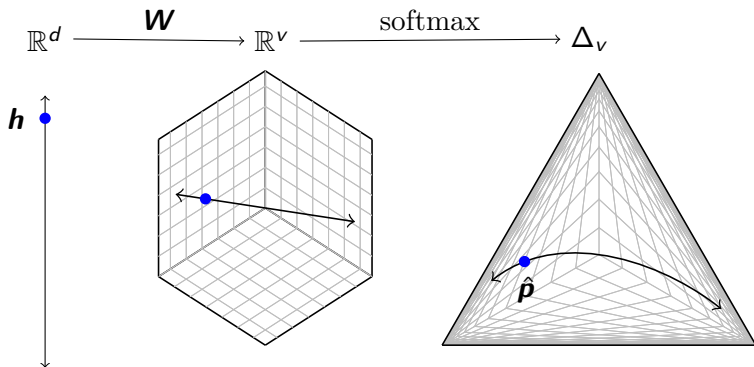
# What are next-token distributions ( $\rho$ )?

- ▶ Distributions over  $v$  items are vectors in  $\mathbb{R}^v$ .
- ▶ All distributions must sum to one.
- ▶ Distributions over  $v$  items are vectors in the  $v$ -simplex, or  $\Delta_v$ .



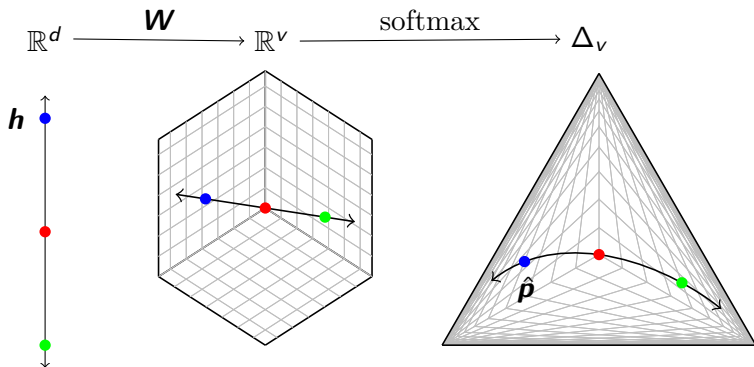
# The softmax bottleneck restricts model outputs

Example:  $d = 1$ ,  $v = 3$ .



# The softmax bottleneck restricts model outputs

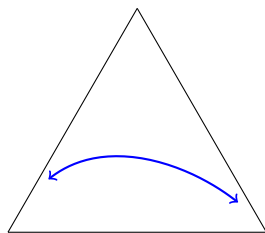
Example:  $d = 1$ ,  $v = 3$ .



Any choice of embedding  $h$  will result in a distribution  $\mathbf{p}$  within a low-dimensional subspace (curved line on right).



## Directly addressing the softmax bottleneck

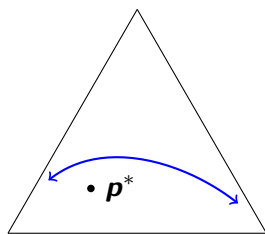


↔ Possible model outputs

- ▶ Model outputs are constrained.



## Directly addressing the softmax bottleneck



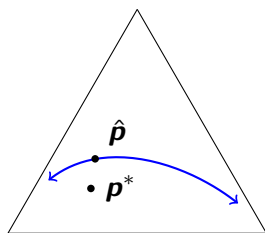
↔ Possible model outputs

- ▶ Model outputs are constrained.
- ▶ Model cannot output true distribution  $p^*$ .





## Directly addressing the softmax bottleneck

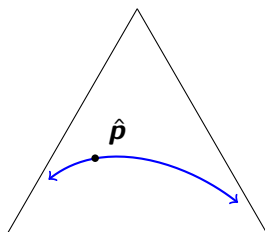


↔ Possible model outputs

- ▶ Model outputs are constrained.
- ▶ Model cannot output true distribution  $\mathbf{p}^*$ .
- ▶ The model outputs  $\hat{\mathbf{p}}$  that minimizes cross-entropy with  $\mathbf{p}^*$ .



## Directly addressing the softmax bottleneck

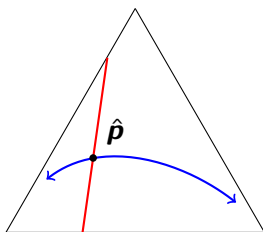


↔ Possible model outputs

- ▶ Reverse: given  $\hat{\rho}$ , what is  $\rho^*$ ?



## Directly addressing the softmax bottleneck

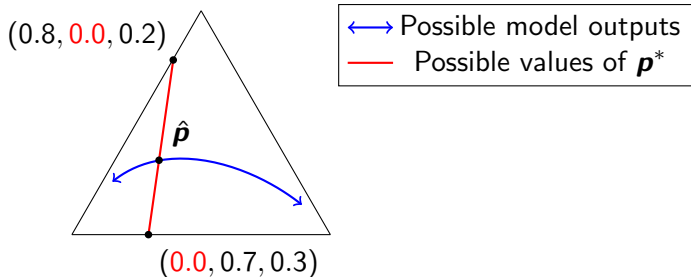


↔ Possible model outputs  
— Possible values of  $\mathbf{p}^*$

- ▶ Reverse: given  $\hat{\mathbf{p}}$ , what is  $\mathbf{p}^*$ ?
- ▶  $\mathbf{p}^*$  could be any distribution  $\hat{\mathbf{p}}$  minimizes cross-entropy with!



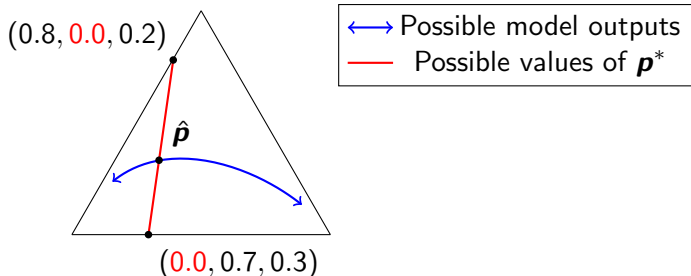
## Directly addressing the softmax bottleneck



- ▶ Reverse: given  $\hat{\mathbf{p}}$ , what is  $\mathbf{p}^*$ ?
- ▶  $\mathbf{p}^*$  could be any distribution  $\hat{\mathbf{p}}$  minimizes cross-entropy with!
- ▶ Don't sample tokens that could have 0 true probability.



## Directly addressing the softmax bottleneck



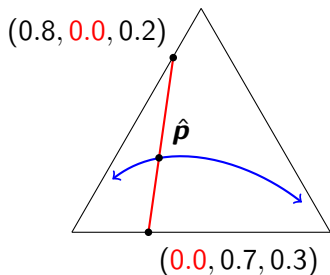
- ▶ Reverse: given  $\hat{\mathbf{p}}$ , what is  $\mathbf{p}^*$ ?
- ▶  $\mathbf{p}^*$  could be any distribution  $\hat{\mathbf{p}}$  minimizes cross-entropy with!
- ▶ Don't sample tokens that could have 0 true probability.
- ▶ The result: a more specific sampling rule!

Only sample tokens  $i$  if there is no solution  $\mathbf{p} \in \Delta_v$  to

$$p_i = 0, \quad \mathbf{W}^T \mathbf{p} = \mathbf{W}^T \hat{\mathbf{p}}.$$



## Directly addressing the softmax bottleneck

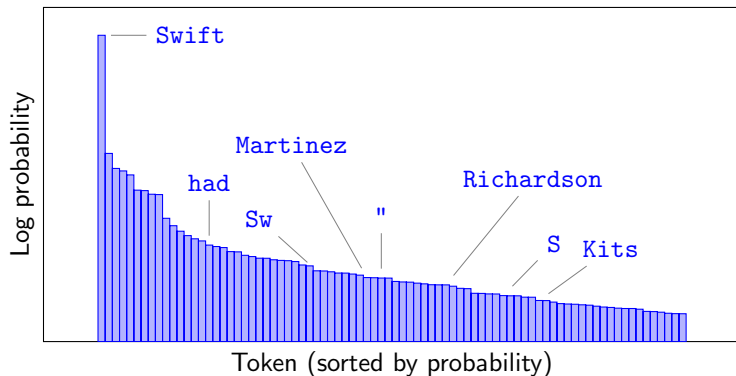


↔ Possible model outputs  
— Possible values of  $p^*$

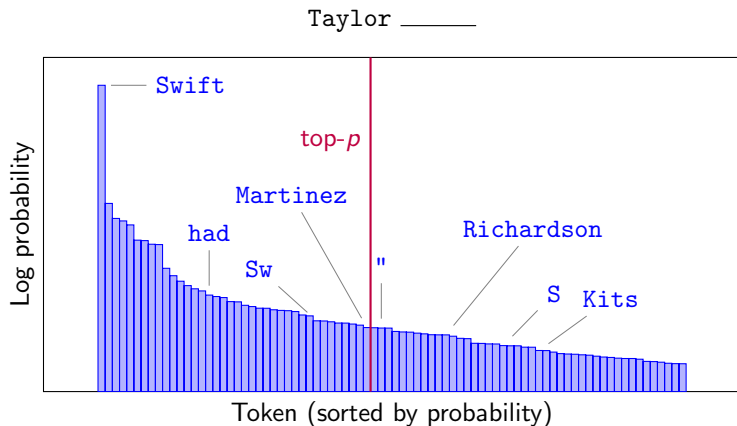


# Anecdotal evidence

Taylor \_\_\_\_\_

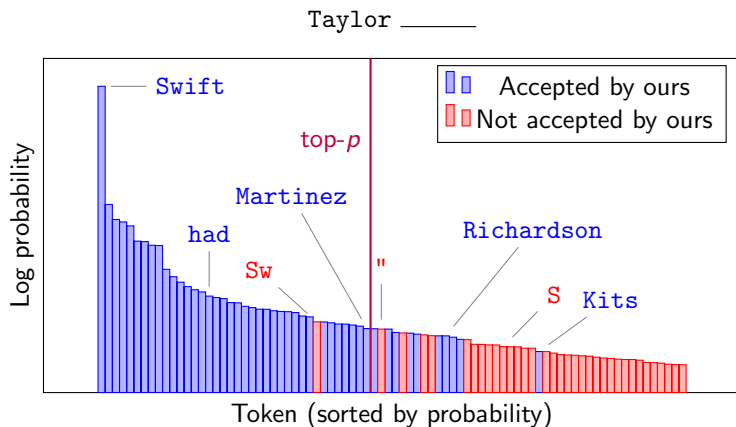


# Anecdotal evidence



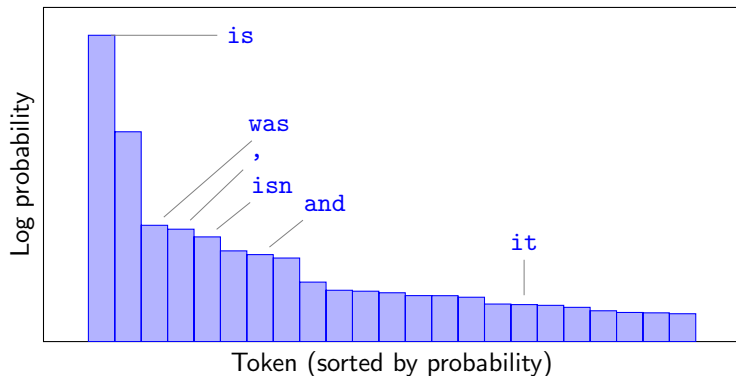


# Anecdotal evidence



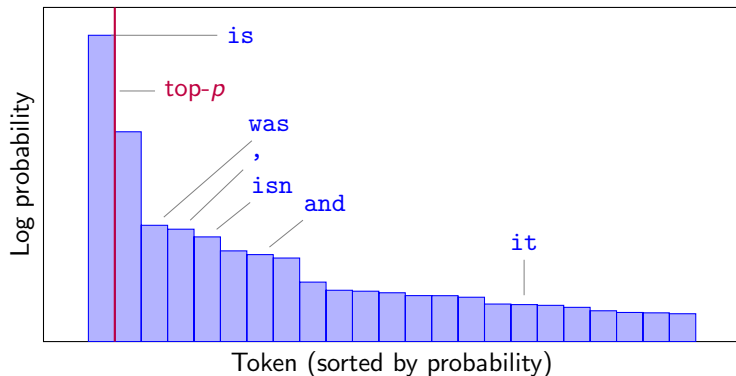
# Anecdotal evidence

My name \_\_\_\_\_



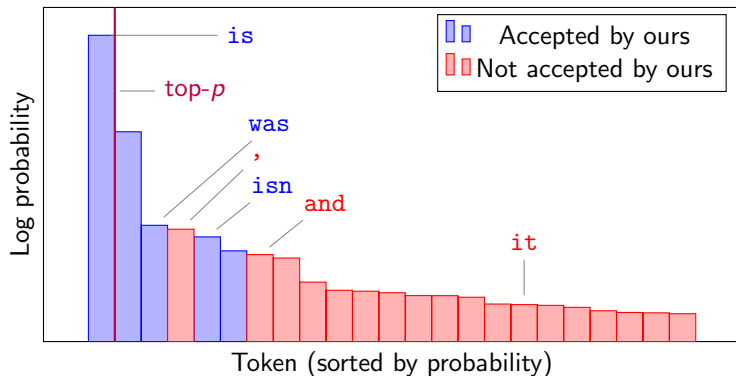
# Anecdotal evidence

My name \_\_\_\_\_



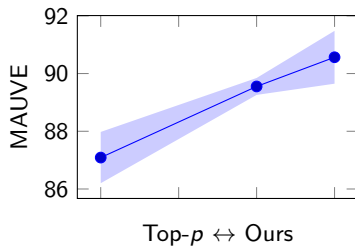
# Anecdotal evidence

My name \_\_\_\_\_



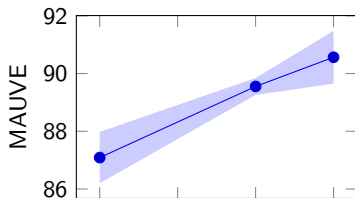
# Basis-aware threshold (BAT) sampling

- ▶ More BAT-like, higher MAUVE (similarity to human text).

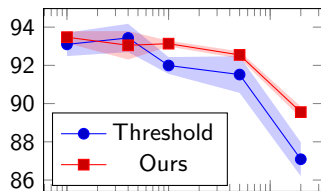


# Basis-aware threshold (BAT) sampling

- ▶ More BAT-like, higher MAUVE (similarity to human text).
- ▶ Low-entropy (closer to greedy) BAT > threshold sampling.



Top-p ↔ Ours



→ closer to greedy



# Basis-aware threshold sampling

The smoking gun

Our method outperforms greedy-like decoding across model sizes

Size	Small	Medium	Large	XL
Method				
Threshold	85.0 <sub>1.4</sub>	90.4 <sub>0.1</sub>	86.0 <sub>0.5</sub>	87.1 <sub>1.2</sub>
Ours	<b>87.8<sub>1.0</sub></b>	<b>92.2<sub>0.6</sub></b>	<b>88.4<sub>0.5</sub></b>	<b>89.6<sub>0.4</sub></b>



## Recap

- ▶ Threshold sampling (e.g., top- $p$ ) is a coarse heuristic to avoid model errors.
- ▶ Understanding a source of model errors (the softmax bottleneck) allows us to better recover the true distribution.
- ▶ We can leverage this method to sample more precisely.
- ▶ BAT outperforms threshold sampling in greedy-like settings.
- ▶ Lots of room for improvement :)



Thank you!

