

# Softmax Bottleneck Makes Language Models Unable to Represent Multi-mode Word Distributions

Haw-Shiuan Chang

Andrew McCallum

UMassAmherst

Manning College of Information  
& Computer Sciences

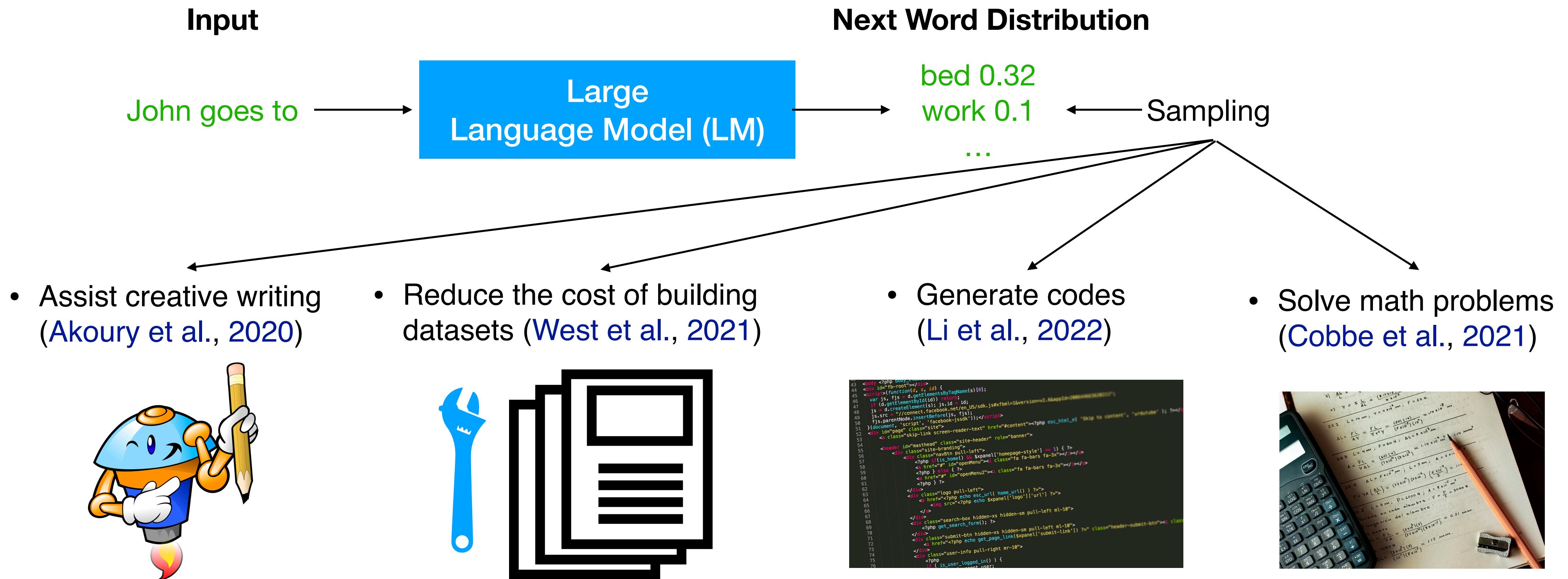
# Outline

- Introduction
- Theoretical Analysis
- Method
- Experiments
- Conclusion and Future Work

# Outline

- Introduction
- Theoretical Analysis
- Method
- Experiments
- Conclusion and Future Work

# Sampling Distributions from Large LMs



Nader Akoury, Shufan Wang, Josh Whiting, Stephen Hood, Nanyun Peng, and Mohit Iyyer. 2020. STORIUM: A Dataset and Evaluation Platform for Machine-in-the-Loop Story Generation. In EMNLP

Peter West, Chandra Bhagavatula, Jack Hessel, Jena D Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2021. Symbolic knowledge distillation: from general language models to commonsense models. arXiv preprint arXiv:2110.07178.

Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittweis, R mi Leblond, Tom Eccles, James Keeling, et al. 2022. Competition-level code generation with alphacode. arXiv preprint arXiv:2203.07814.

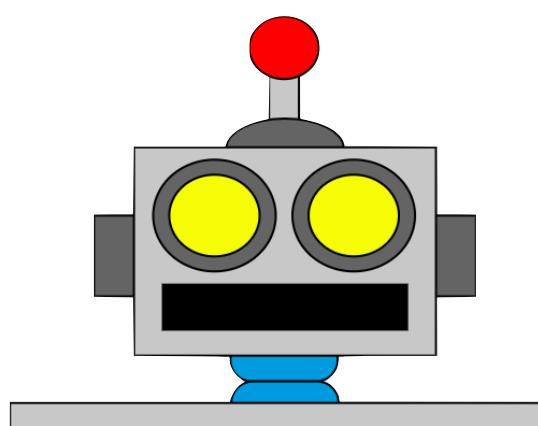
Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.

Can large LMs learn any  
distribution over the next word?

# An Ambiguous Context



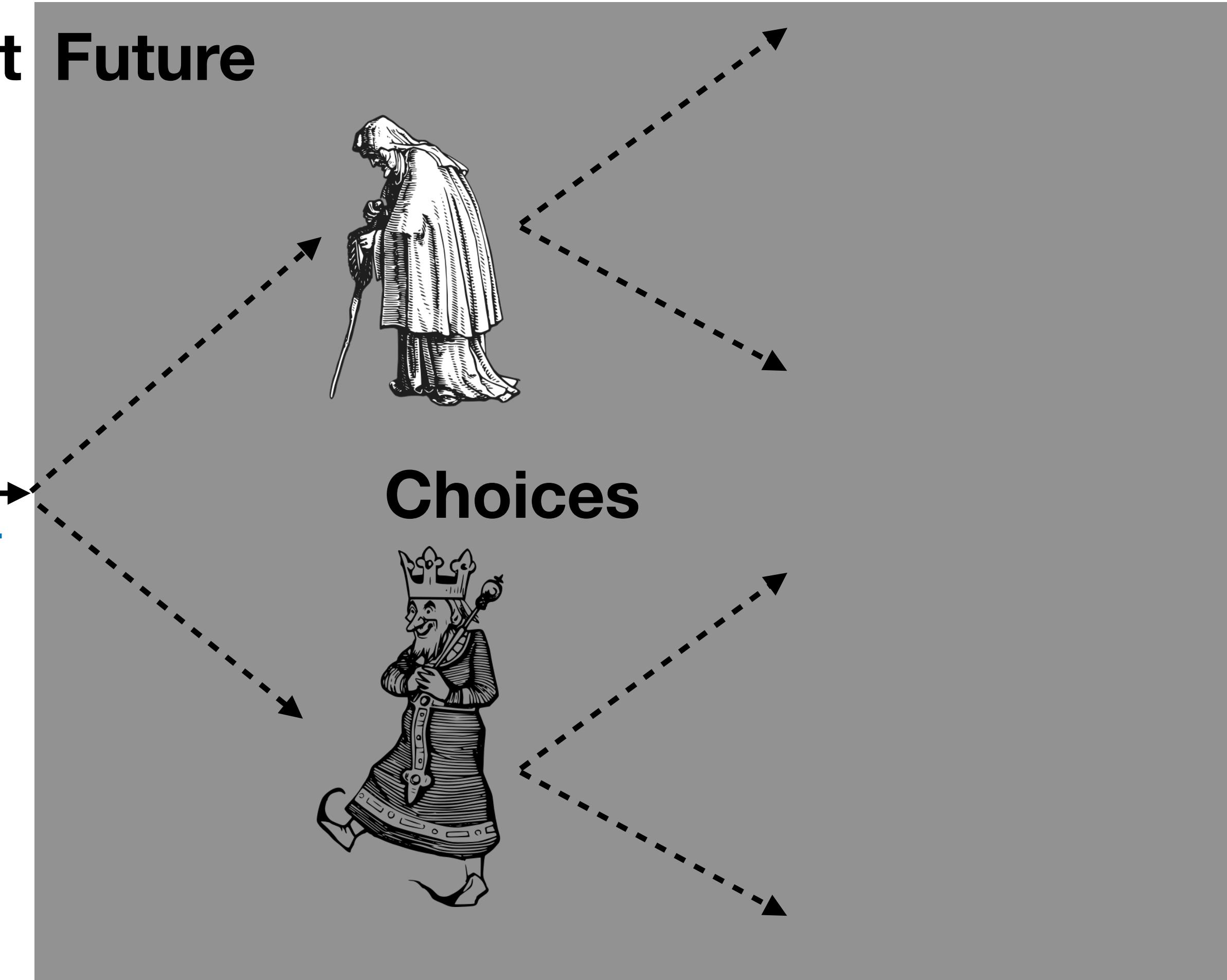
After debating whether to bow to the **king** or the **woman** first, the jester decided on the



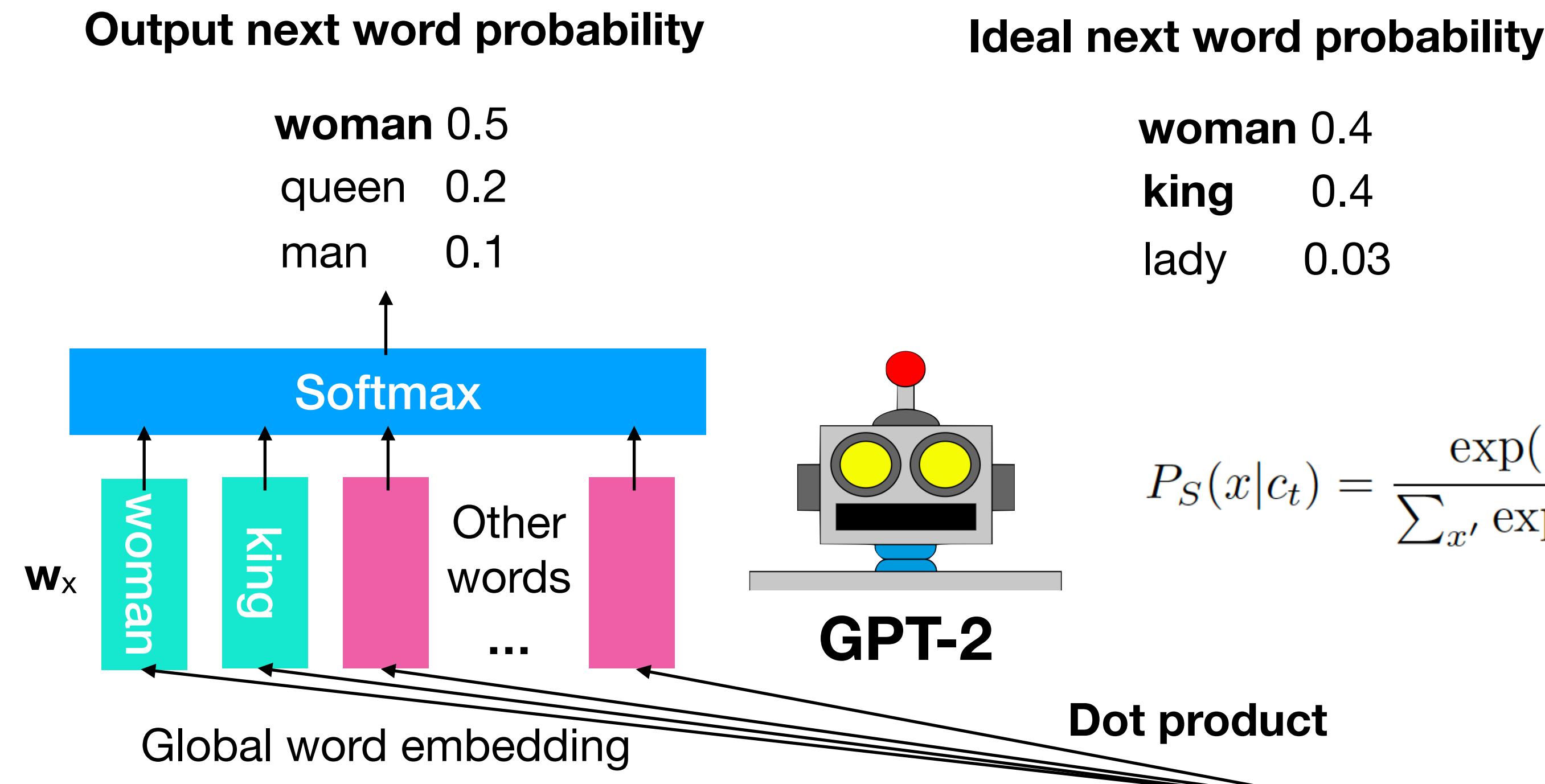
GPT-2

Context Future

Choices

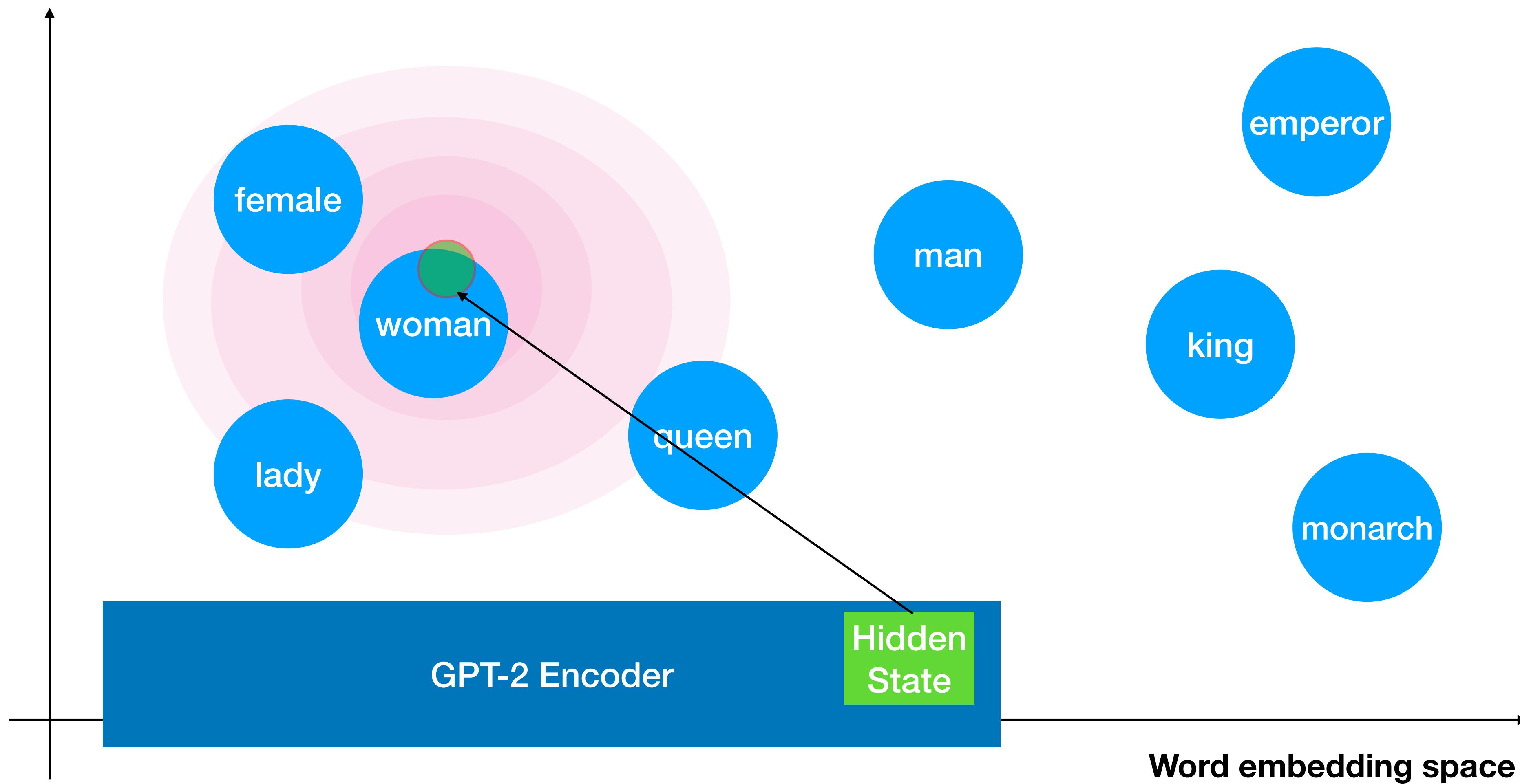


# Most of Existing Approaches

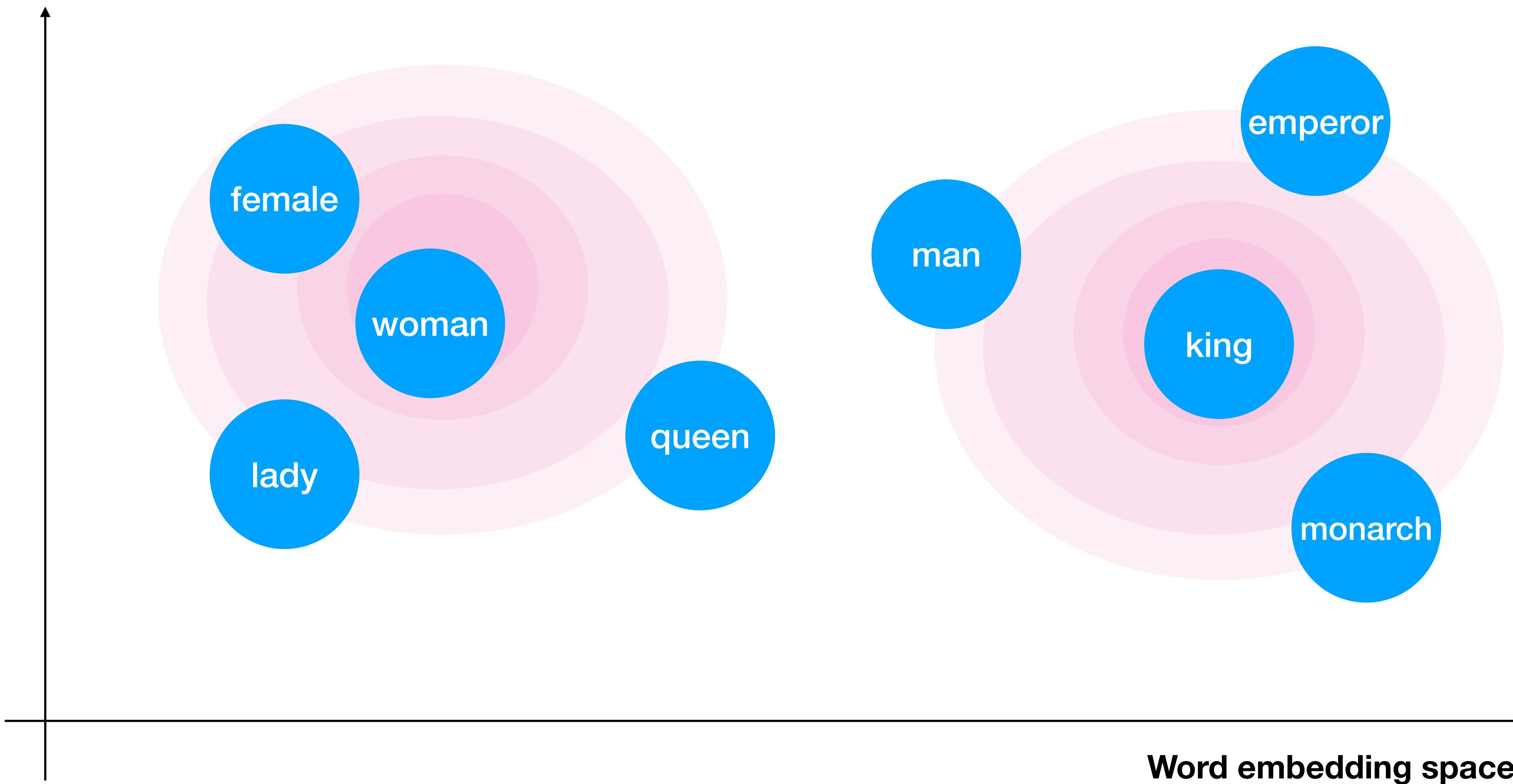


$c_t$  After debating whether to bow to the **king** or the **woman** first, the jester decided on the

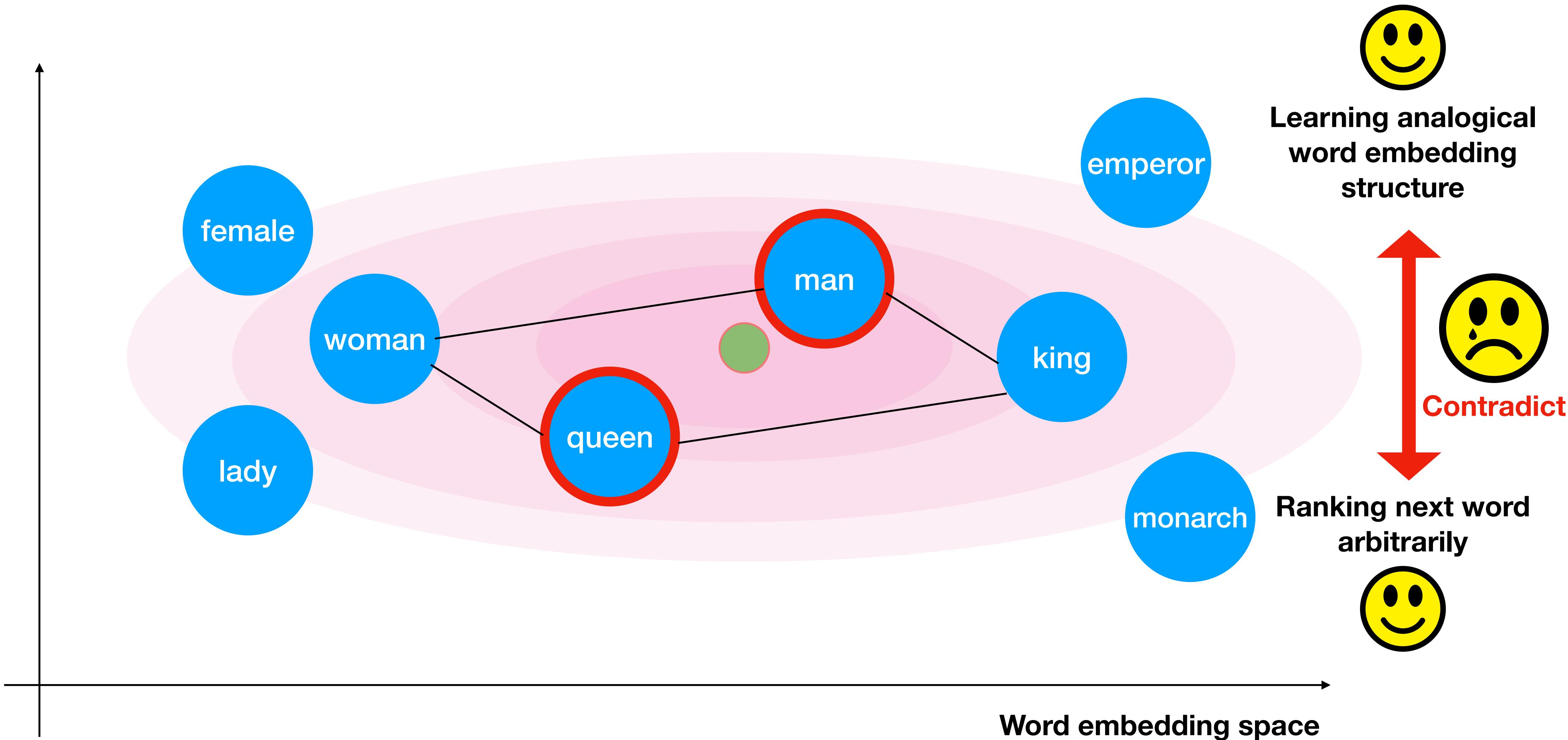
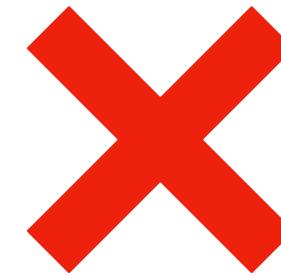
# Predicting “woman” as the Next Word ✓



# Could GPT-2 Predict Both “woman” and “king” as the Next Word? ?



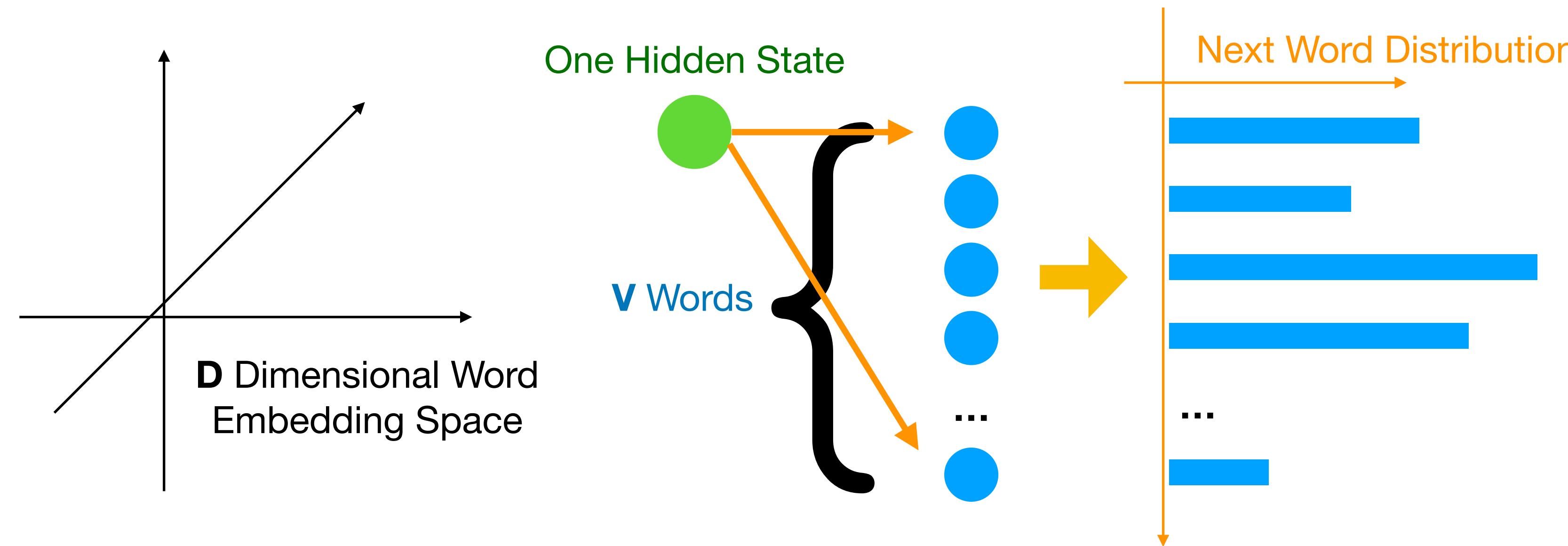
No, if there are some words between them  
and GPT-2 has only one hidden state



# Outline

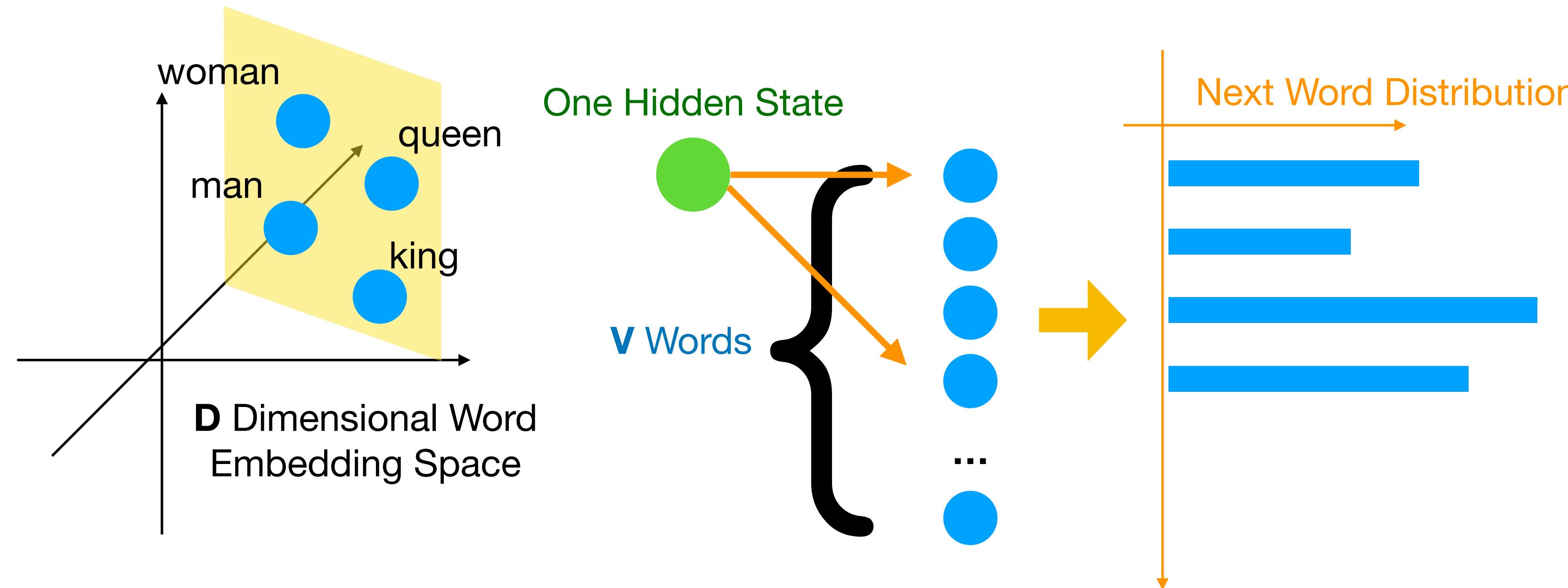
- Introduction
- Theoretical Analysis
- Method
- Experiments
- Conclusion and Future Work

# Softmax Bottleneck (Yang et al., 2018)



- If  $V > D$ , we cannot output arbitrary probabilities over  $V$  words
- Limitations
  - Serious among which words?
  - Affect the top words? If yes, when?
  - Disappears after making  $D > V$ ?

# Our Theoretical Improvements



- If  $N$  words are in a small subspace, we cannot rank  $N$  words arbitrarily
- Improvements
  - Serious among which words?      -> Among words in a small subspace
  - Affect the top words? If yes, when? -> Yes. When the ideal distribution is multi-mode
  - Disappears after making  $D > V$ ?      -> No, if some words are in a small subspace

# A Limitation of Single Embedding

**Theorem 1** (simplified): If many word embeddings are linearly dependent, the softmax in a LM cannot rank the words arbitrarily

**Example:** If “woman - man = queen - king”, GPT-2 cannot rank the word woman and king as the top 2 words

**Example:** If “UMass = 0.2 University + 0.2 Massachusetts”, GPT-2 cannot rank a rare word UMass on top of the similar popular words University and Massachusetts (Demeter et al., 2020).

**Linear Algebra Intuition:** N+1 words are linear dependent →

They are in subspace with  $d < N \rightarrow$  cannot have arbitrary probabilities

# Approximately Linearly Dependent

**Theorem 2** (simplified): If many word embeddings are approximately linearly dependent and the magnitude of the hidden state has an upperbound, the softmax in a LM cannot assign very small probabilities to some words

**Example:** If "woman + king = queen + man +  $\varepsilon$ ", GPT-2 cannot make the logits of queen and man much smaller than the logits of king and woman

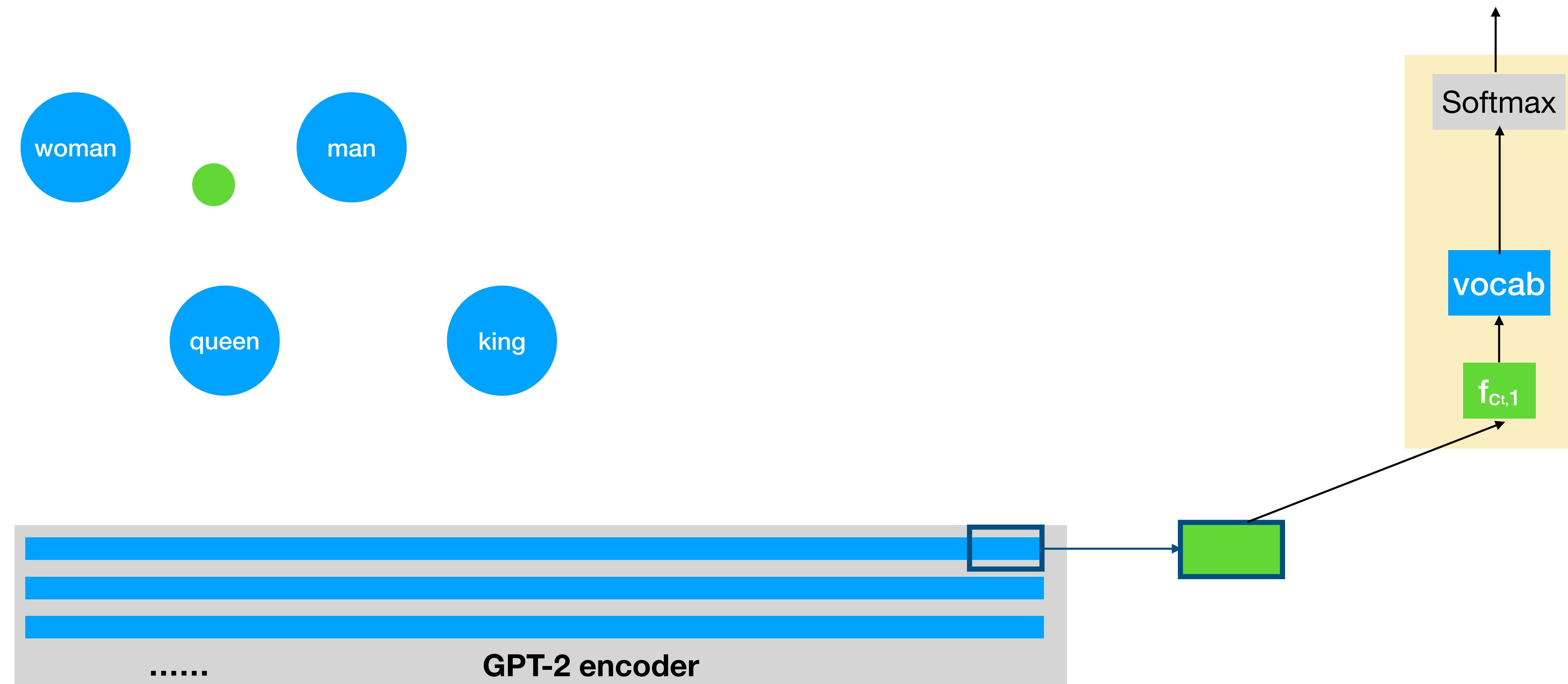
**Example:** If "woman = man +  $\varepsilon$ ", GPT-2 cannot make the logits of man much smaller than the logits of woman

**Intuition:**  $\underline{h}^T \underline{\text{king}} + \underline{h}^T \underline{\text{woman}} = \underline{h}^T \underline{\text{queen}} + \underline{h}^T \underline{\text{man}} + \underline{h}^T \underline{\varepsilon}$ , and we can ignore  $\underline{h}^T \underline{\varepsilon}$  if  $\|\underline{h}\|$  and  $\|\underline{\varepsilon}\|$  are both small

# Outline

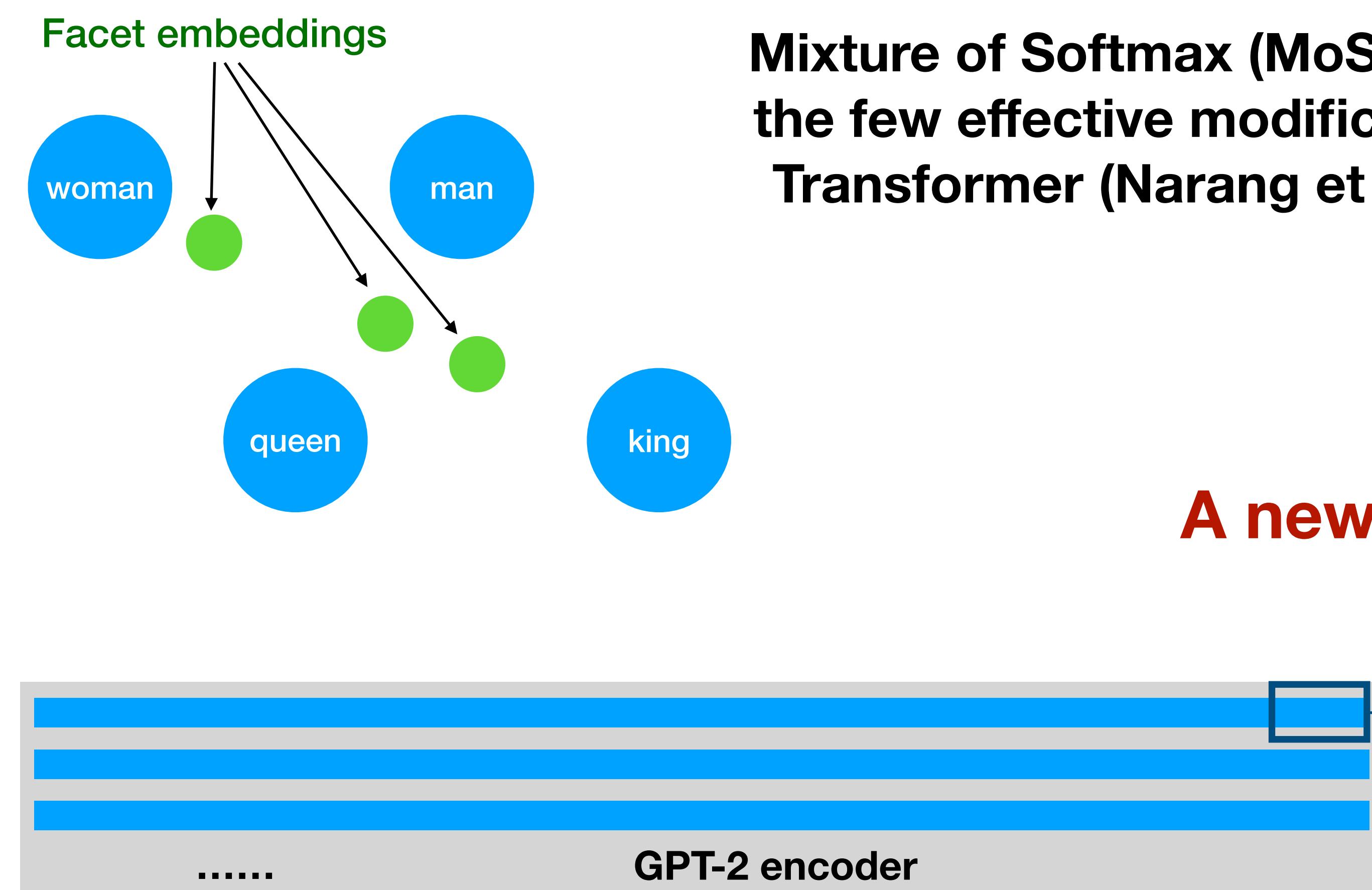
- Introduction
- Theoretical Analysis
- Method
- Experiments
- Conclusion and Future Work

# GPT-2 (Softmax)



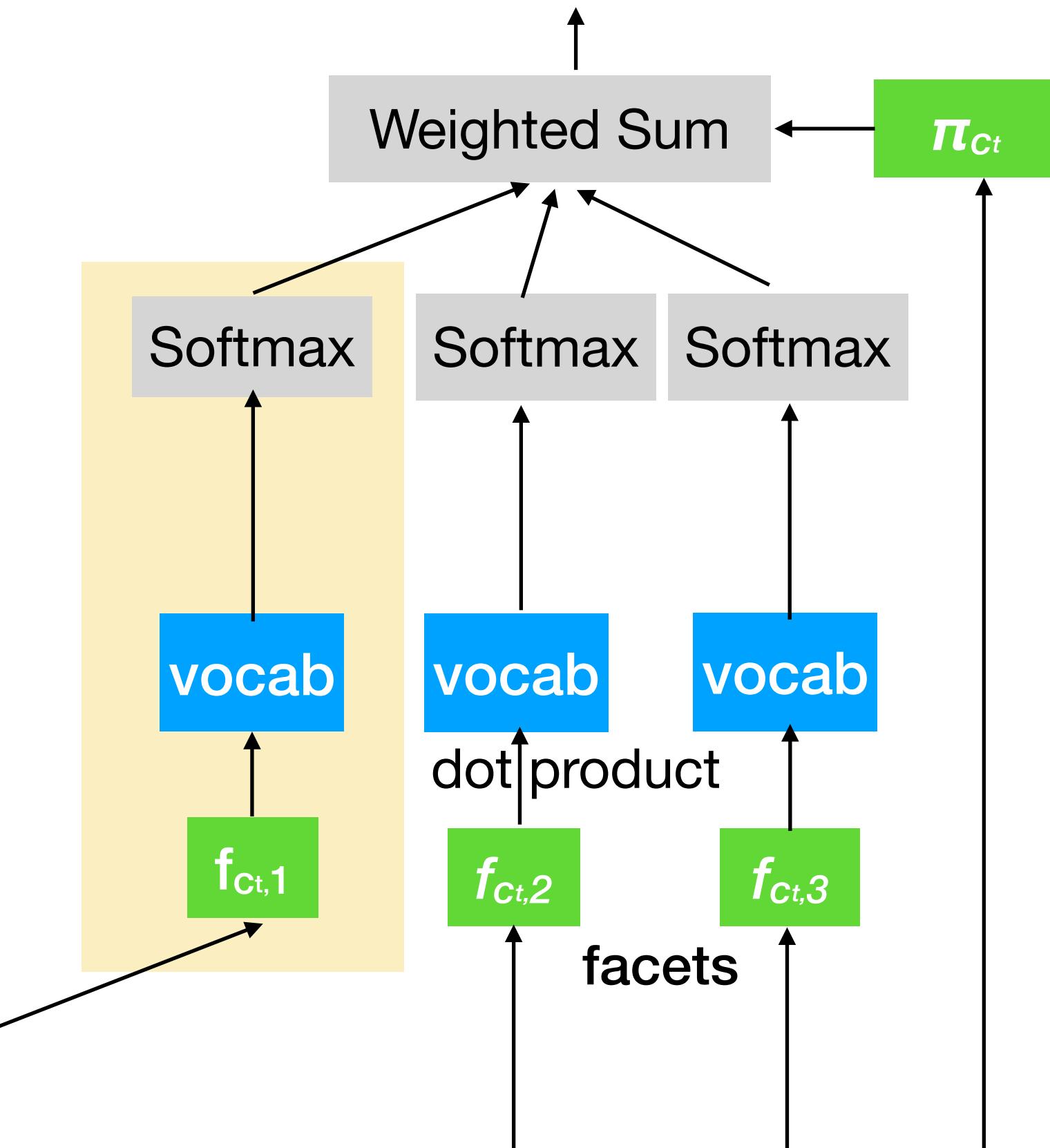
After debating whether to bow to the **king** or the **woman** first, the jester decided on the

# Mixture of Softmax (Yang et al., 2018)



**Mixture of Softmax (MoS) is one of the few effective modifications for Transformer (Narang et al., 2021)**

**A new bottleneck**

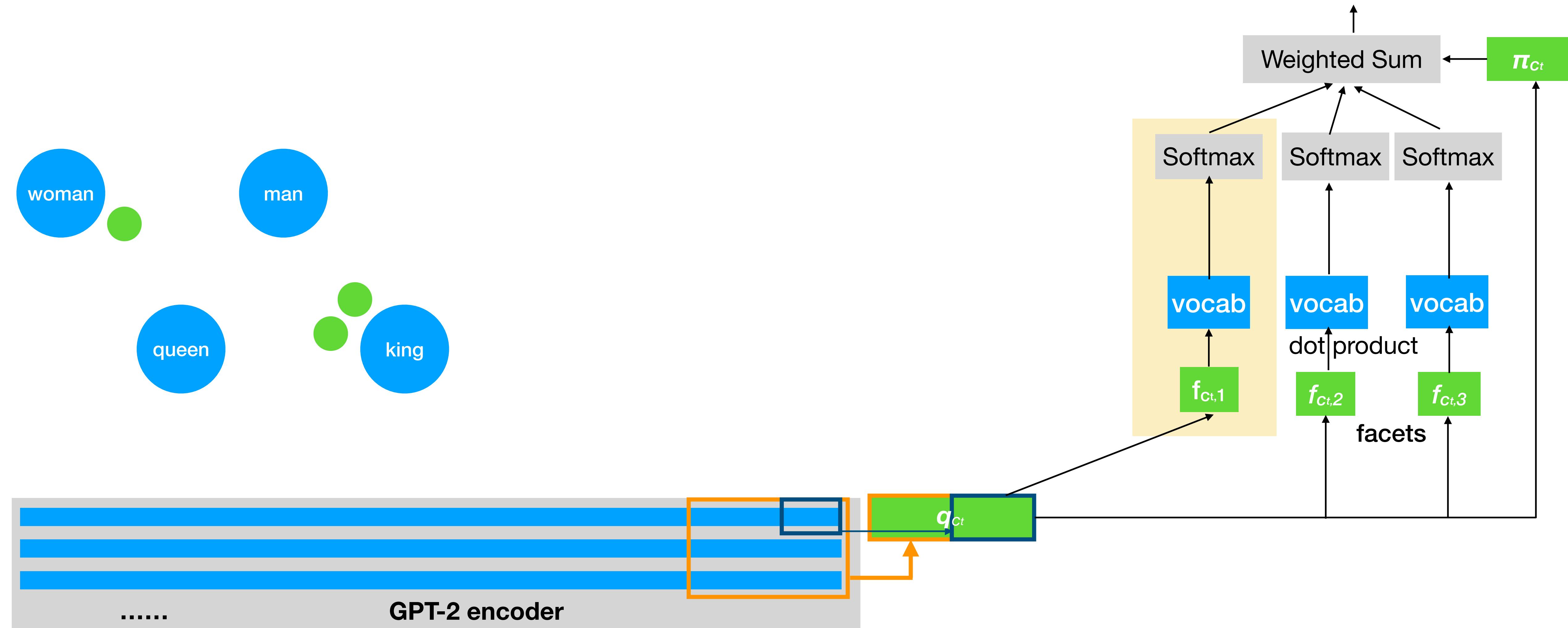


Yang, Zhilin, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen.  
"Breaking the Softmax Bottleneck: A High-Rank RNN Language Model."  
In *International Conference on Learning Representations*. 2018.

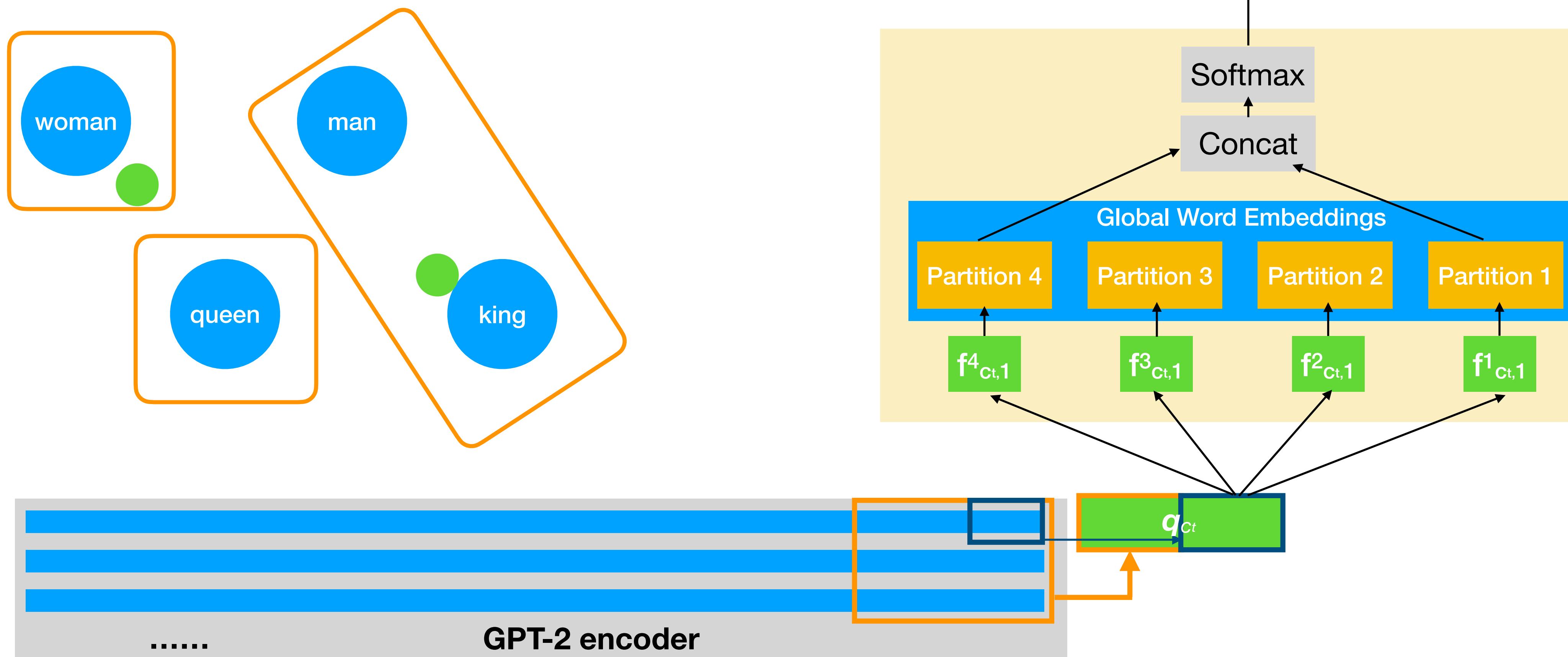
Sharan Narang, et al. Do transformer modifications transfer across implementations and applications? *EMNLP 2021*

After debating whether to bow to the **king** or the **woman** first, the jester decided on the

# MoS + Multi-input

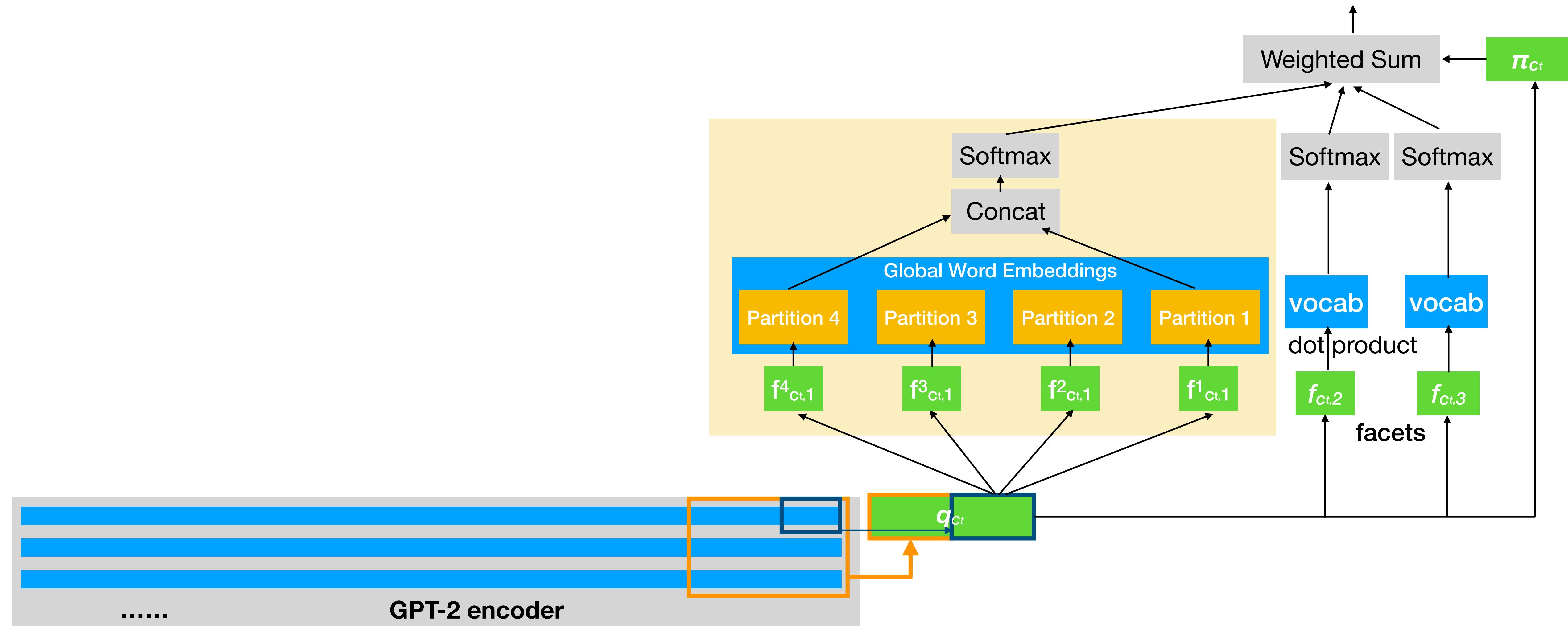


# Multiple Partitions



After debating whether to bow to the **king** or the **woman** first, the jester decided on the

# Multi-facet Softmax (MFS)



After debating whether to bow to the **king** or the **woman** first, the jester decided on the

# Outline

- Introduction
- Theoretical Analysis
- Method
- Experiments
- Conclusion and Future Work

# Multi-facet Softmax (MFS) Perplexity

Only adding nonlinearity is not enough (Parthiban et al., 2021)

Models ↓	Configuration			GPT-2 Small				GPT-2 Medium			
	#S	#I	#P	Size	Time	OWT	Wiki	Size	Time	OWT	Wiki
Softmax (GPT-2)	1	1	1	163.6M	84ms	18.72	24.06	407.3M	212ms	15.89	20.34
SigSoftmax (Kanai et al., 2018)	1	1	1	163.6M	91ms	18.63	24.06	407.3M	221ms	16.07	20.65
Softmax + Multi-input	1	9	1	169.5M	87ms	18.50	23.89	417.8M	219ms	15.76	20.29
Softmax + Multi-partition	1	1	4	165.4M	88ms	18.77	24.08	410.5M	218ms	15.89	20.30
MoS (Yang et al., 2018) (4)	4	1	1	165.4M	152ms	18.61	23.77	410.5M	299ms	15.75	20.08
MoS (Yang et al., 2018) (3)	3	1	1	164.8M	130ms	18.63	23.81	409.4M	270ms	15.79	20.11
DOC (Takase et al., 2018)	3	3	1	164.8M	130ms	18.69	24.02	409.4M	270ms	15.88	20.34
MFS w/o Multi-partition	3	9	1	171.9M	133ms	18.37	23.56	422.0M	276ms	15.65	20.06
MFS w/o Multi-input	3	1	4	166.6M	134ms	18.60	23.72	412.6M	275ms	15.71	20.08
MFS (Ours)	3	9	4	175.4M	138ms	18.29	23.45	428.3M	283ms	15.64	20.02

Multiple input hidden states help

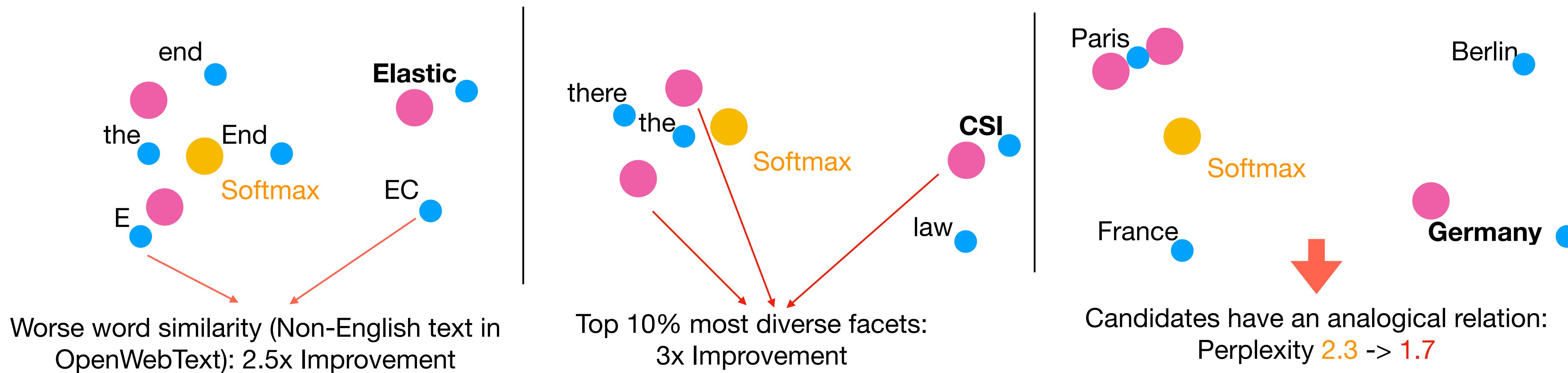
Multiple partitions help

Table 1: Perplexity comparison between MFS (Ours) and baselines. #S, #I, #P are the number of softmaxes (i.e.,  $K$ ), input hidden states, and partitions, respectively. The top four baselines use a single softmax. OWT and Wiki are the test set perplexity of OpenWebText and Wikipedia 2021, respectively. The standard errors of all models are smaller than 0.02 perplexity. We also compare the number of parameters and the inference time on one batch.

Improvement of MFS over Softmax is around 15% between GPT-2 Small and GPT-2 Medium (with 3x parameters)

# Examples

Corpus →	OpenWebText	Wikipedia 2021	Analogy in Templates (Section 5)
Input Context	... The Elastic Endpoint Security and Elastic SIEM solutions mentioned in this post are now referred to as <b>Elastic</b>	... law and chance working together cannot generate CSI, either. Moreover, he claims that <b>CSI</b>	I went to Paris and Germany before, and I love one of the places more, which is <b>Germany</b>
Softmax (GPT-2)	the 0.087, E 0.043, End 0.039	the 0.174, this 0.054, if 0.038	Paris 0.893, France 0.045, <b>Germany</b> 0.033
MFS (Ours)	<b>Elastic</b> 0.220, the 0.089, EC 0.033	CSI 0.186, the 0.140, there 0.033	Paris 0.544, <b>Germany</b> 0.389, France 0.064
MFS Softmax 1	end 0.051, the 0.043, security 0.023	the 0.191, law 0.127, if 0.053	Paris 0.979, France 0.013, <b>Germany</b> 0.007
MFS Softmax 2	<b>Elastic</b> 0.652, EC 0.080, ES 0.046	the 0.191, there 0.049, this 0.047	Paris 1.000 Berlin 0.000 ##Paris 0.000
MFS Softmax 3	the 0.193, E 0.040, a 0.014	<b>CSI</b> 0.677, law 0.029, laws 0.019	<b>Germany</b> 0.852, France 0.139, China 0.004



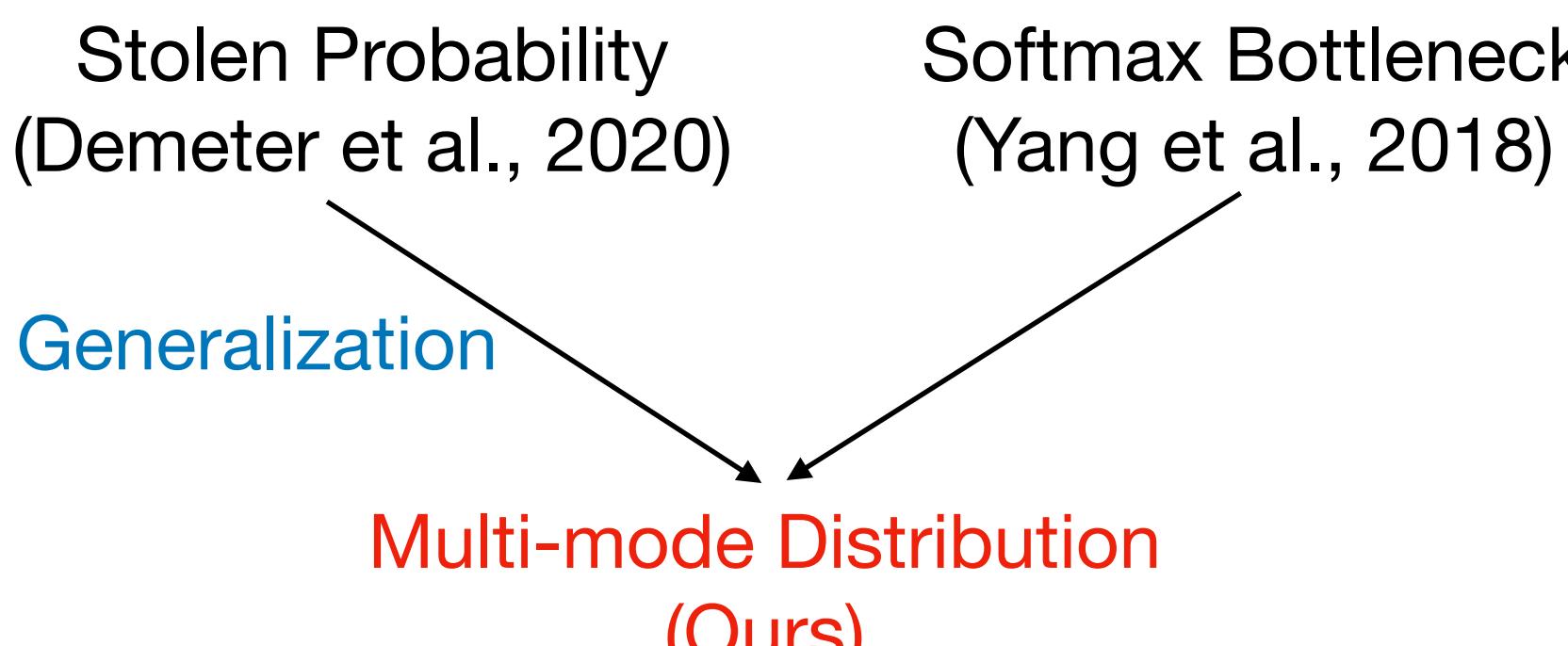
# Outline

- Introduction
- Theoretical Analysis
- Method
- Experiments
- Conclusion and Future Work

# Conclusion

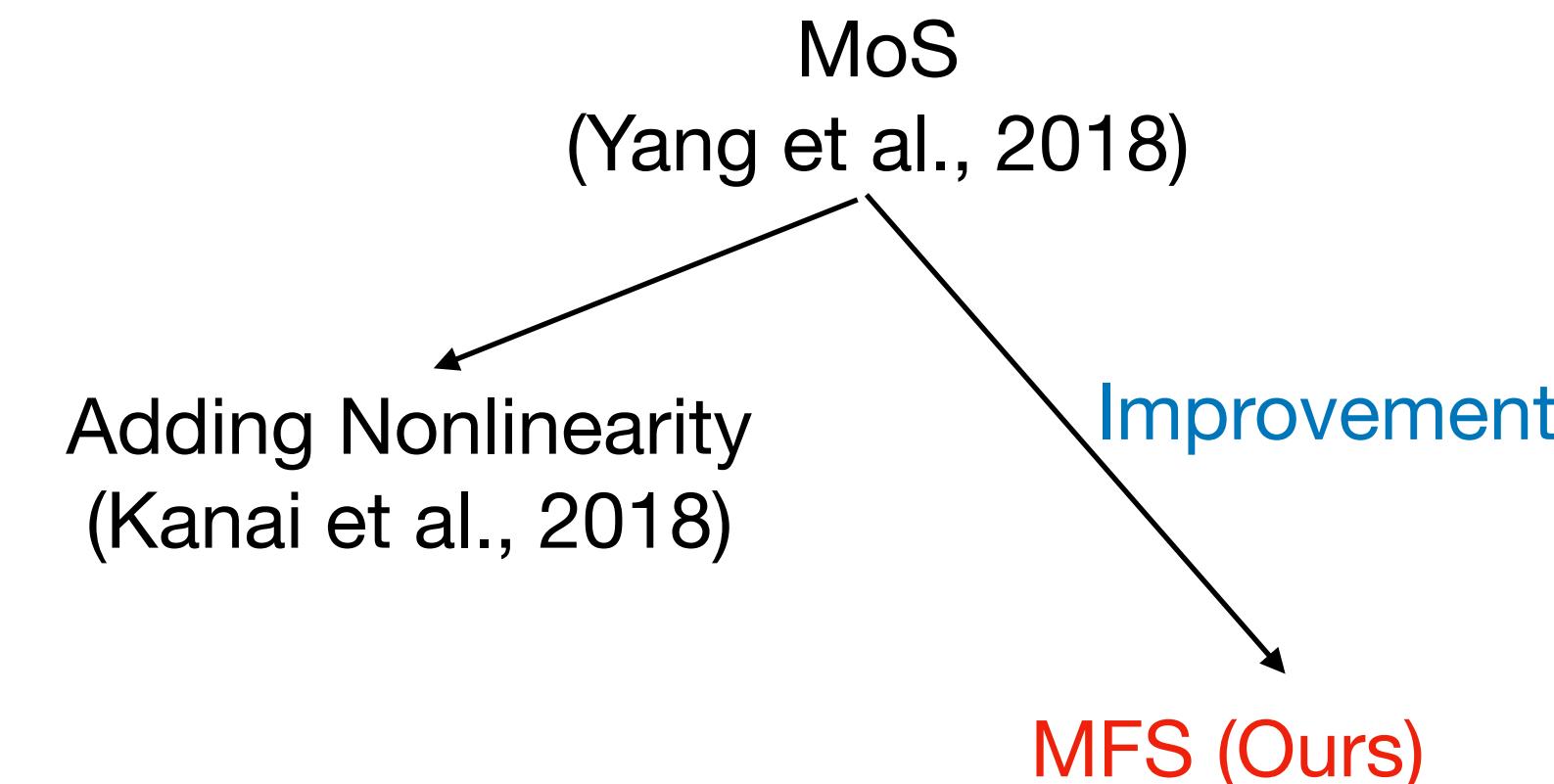
## Theory

- Multi-mode distribution must exist if some word embeddings are in a small subspace



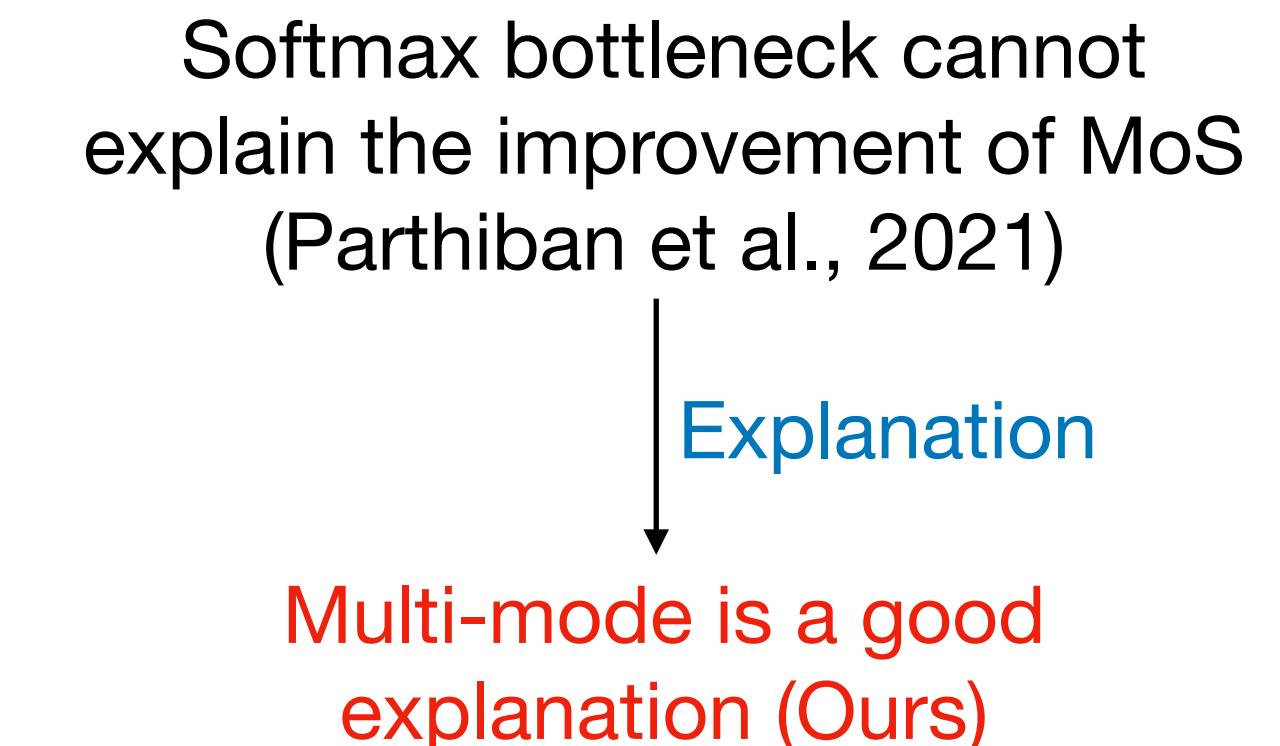
## Method

- We propose two enhancements for mixture of softmax (MoS)



## Analysis

- Our proposed method improves softmax layer in GPT-2 especially when the ideal next word distribution is multi-mode



Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen. "Breaking the Softmax Bottleneck: A High-Rank RNN Language Model." In ICLR. 2018.

David Demeter, Gregory Kimmel, and Doug Downey. Stolen probability: A structural weakness of neural language models. In ACL. 2020

Sekitoshi Kanai, Yasuhiro Fujiwara, Yuki Yamanaka, and Shuichi Adachi. Sigsoftmax: Reanalysis of the softmax bottleneck. In NeurIPS 2018

Dwarak Govind Parthiban, Yongyi Mao, and Diana Inkpen. On the softmax bottleneck of recurrent language models. In AAAI 2021

# Future Work

- How much MFS could help huge language models (e.g., GPT-3)
- Whether MFS could improve
  - NLU tasks
  - NLG tasks
  - Other extreme classification models using an output softmax layer
- The word similarity should be context dependent rather than globally fixed

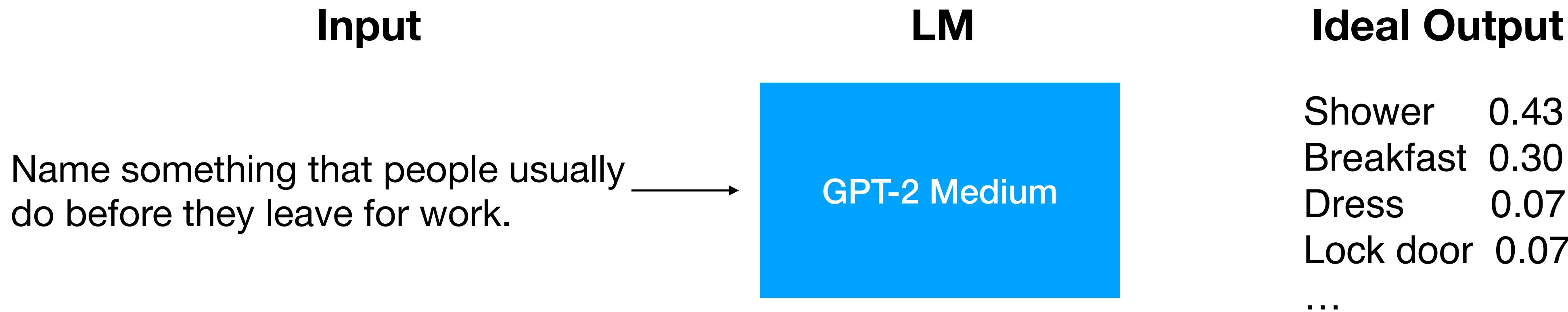
# **Appendix**

# BERT Experiment

Input	LM	Ideal Output
After debating whether to bow to the <b>king</b> or the <b>woman</b> first, the jester decided on the	GPT-2 Small	<b>king</b> or <b>woman</b>
After debating whether to bow to the <b>king</b> or the <b>woman</b> first, the jester decided on the [MASK], which makes him pleased.	BERT Base	<b>king</b>

- Perplexity improvement of MoS and MFS is much smaller compared to GPT-2
- MFS still doubles the improvement of MoS over Softmax

# ProtoQA



Models ↓	Perplexity on Scraped Development Set	Max Answers				Max Incorrect		
		Top 1	Top 3	Top 5	Top 10	Top 1	Top 3	Top 5
Softmax (GPT-2)	$1.5432 \pm 0.0003$	$34.1 \pm 0.8$	$35.2 \pm 0.5$	$37.8 \pm 0.4$	$45.0 \pm 0.5$	$18.3 \pm 0.4$	$30.7 \pm 0.5$	$38.5 \pm 0.6$
MoS (Yang et al., 2018) (3)	$1.5407 \pm 0.0004$	$33.9 \pm 0.8$	$36.0 \pm 0.6$	$37.7 \pm 0.6$	$44.9 \pm 0.4$	$18.3 \pm 0.4$	$31.7 \pm 0.6$	$38.2 \pm 0.6$
MFS w/o Multi-partition	$1.5411 \pm 0.0003$	$34.3 \pm 0.7$	$36.7 \pm 0.7$	$38.1 \pm 0.5$	$45.2 \pm 0.4$	$19.4 \pm 0.4$	$32.0 \pm 0.5$	$38.6 \pm 0.3$
MFS (Ours)	$\mathbf{1.5402 \pm 0.0005}$	$34.1 \pm 0.6$	$36.7 \pm 0.5$	$38.6 \pm 0.4$	$45.4 \pm 0.5$	$19.7 \pm 0.4$	$32.1 \pm 0.4$	$39.7 \pm 0.4$

Table 5: ProtoQA performances. All the numbers except perplexity are the percentages of the predictions that match the ground truth exactly on the crowdsourced development set. Max answers top k implies only evaluating the top k answers. Max incorrect top k indicates only evaluating the top answers that contain k errors. The best average performances are highlighted and the standard errors are reported as the confidence interval.

# Proof Sketch

- $1 \underline{w}_{king} - 1 \underline{w}_{queen} = 1 \underline{w}_{man} - 1 \underline{w}_{woman}$

- $1 \underline{w}_{king} + 1 \underline{w}_{woman} = 1 \underline{w}_{queen} + 1 \underline{w}_{man}$

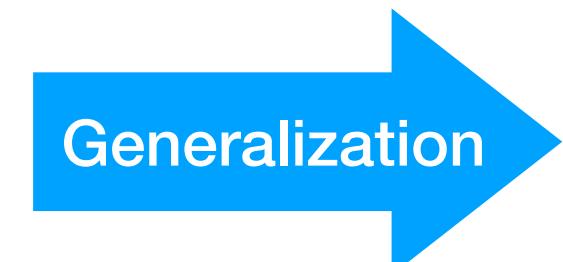
$\times \underline{h}$  (hidden state) on both side

- $1 \underline{h}^T \underline{w}_{king} + 1 \underline{h}^T \underline{w}_{woman} = 1 \underline{h}^T \underline{w}_{queen} + 1 \underline{h}^T \underline{w}_{man}$

- If  $\exists \underline{h}$ , s.t  $\min(\underline{h}^T \underline{w}_{king}, \underline{h}^T \underline{w}_{woman}) > \max(\underline{h}^T \underline{w}_{queen}, \underline{h}^T \underline{w}_{man})$

$\text{Large}$	$\text{Large}$
$\text{Large}$	$\text{Large}$
$\text{Small}$	$\text{Small}$
$1 \underline{h}^T \underline{w}_{king} + 1 \underline{h}^T \underline{w}_{woman} \geq$	
$2 \min(\underline{h}^T \underline{w}_{king}, \underline{h}^T \underline{w}_{woman}) >$	
$2 \max(\underline{h}^T \underline{w}_{queen}, \underline{h}^T \underline{w}_{man}) \geq$	
$1 \underline{h}^T \underline{w}_{queen} + 1 \underline{h}^T \underline{w}_{man} (\rightarrow \leftarrow)$	
$\text{Small}$	$\text{Small}$

- Thus, the logits of LM cannot rank both king and woman on top of queen and man



- Linearly dependent among  $\{\underline{w}_{l_1}, \dots, \underline{w}_{l_L}, \underline{w}_{r_1}, \dots, \underline{w}_{r_R}\}$

- $a_{l_1} \underline{w}_{l_1} + \dots + a_{l_L} \underline{w}_{l_L} = a_{r_1} \underline{w}_{r_1} + \dots + a_{r_R} \underline{w}_{r_R}$

- All coefficient  $a_{l_i} > 0, a_{r_j} > 0$

- WLOG  $a_{l_1} + \dots + a_{l_L} \geq a_{r_1} + \dots + a_{r_R}$

- $a_{l_1} \underline{h}^T \underline{w}_{l_1} + \dots + a_{l_L} \underline{h}^T \underline{w}_{l_L} = a_{r_1} \underline{h}^T \underline{w}_{r_1} + \dots + a_{r_R} \underline{h}^T \underline{w}_{r_R}$

- If  $\exists \underline{h}$ , s.t  $\min_i(\underline{h}^T \underline{w}_{l_i}) > \max_j(\underline{h}^T \underline{w}_{r_j})$

- $a_{l_1} \underline{h}^T \underline{w}_{l_1} + \dots + a_{l_L} \underline{h}^T \underline{w}_{l_L} \geq$   
 $(a_{l_1} + \dots + a_{l_L})\min_i(\underline{h}^T \underline{w}_{l_i}) >$   
 $(a_{r_1} + \dots + a_{r_R})\max_j(\underline{h}^T \underline{w}_{r_j}) \geq$   
 $a_{r_1} \underline{h}^T \underline{w}_{r_1} + \dots + a_{r_R} \underline{h}^T \underline{w}_{r_R} (\rightarrow \leftarrow)$

- Thus, the logits of LM cannot rank all the left words on top of the right words.