

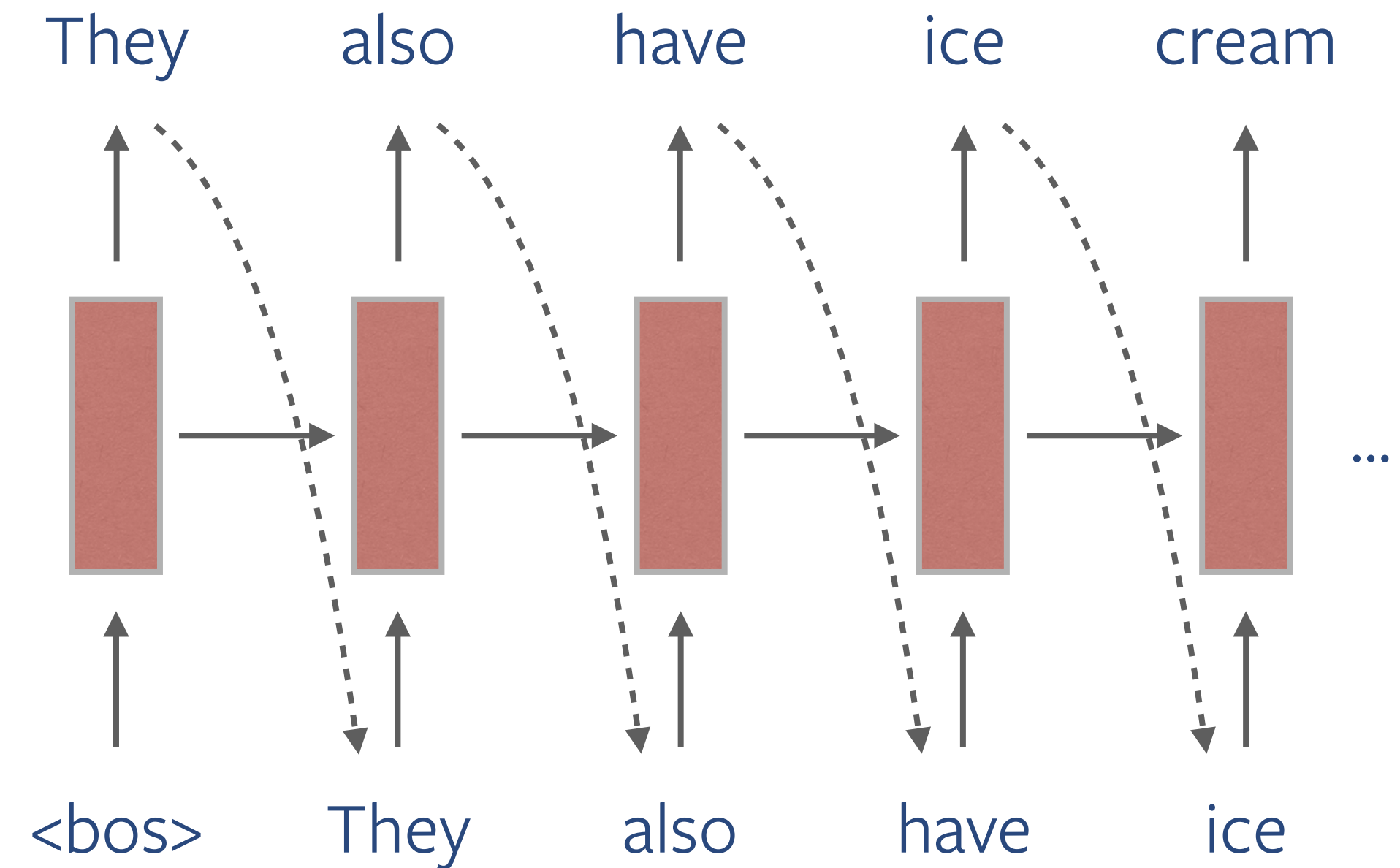
Blank Language Models

Tianxiao Shen* Victor Quach* Regina Barzilay Tommi Jaakkola (*: equal contribution)
tianxiao@mit.edu



Left-to-Right Language Model

- ✓ Generate from scratch
- ✗ Start with partially specified text
 - text editing
 - template filling
 - text restoration
 - ...



Blank Language Model (BLM)

Input: They also have _____ which _____.

Output: They also have ice cream which is really good.

- ✓ Fine-grained control over generation location
- ✓ Respect preceding and following context
- ✓ Variable number of missing tokens

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “ w ”, “___ w ”, “ w ___”, or “___ w ___”
- Stop when there is no “___”

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “ w ”, “___ w ”, “ w ___”, or “___ w ___”
- Stop when there is no “___”

They also have _____ which _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “w”, “___ w”, “w ___”, or “___ w ___”
- Stop when there is no “___”

_____ really _____
↓
They also have _____ which _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “ w ”, “___ w ”, “ w ___”, or “___ w ___”
- Stop when there is no “___”

They also have _____ which _____ really _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “w”, “___ w”, “w ___”, or “___ w ___”
- Stop when there is no “___”

ice _____
↓
They also have _____ which _____ really _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “w”, “___ w”, “w ___”, or “___ w ___”
- Stop when there is no “___”

They also have ice _____ which _____ really _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “w”, “___ w”, “w ___”, or “___ w ___”
- Stop when there is no “___”

They also have ice _____ which _____ really _____.

is
↓

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “ w ”, “___ w ”, “ w ___”, or “___ w ___”
- Stop when there is no “___”

They also have ice _____ which **is** really _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “w”, “___ w”, “w ___”, or “___ w ___”
- Stop when there is no “___”

cream



They also have ice _____ which is really _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “w”, “___ w”, “w ___”, or “___ w ___”
- Stop when there is no “___”

good



They also have ice cream which is really _____.

Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “ w ”, “___ w ”, “ w ___”, or “___ w ___”
- Stop when there is no “___”

They also have ice cream which is really good.

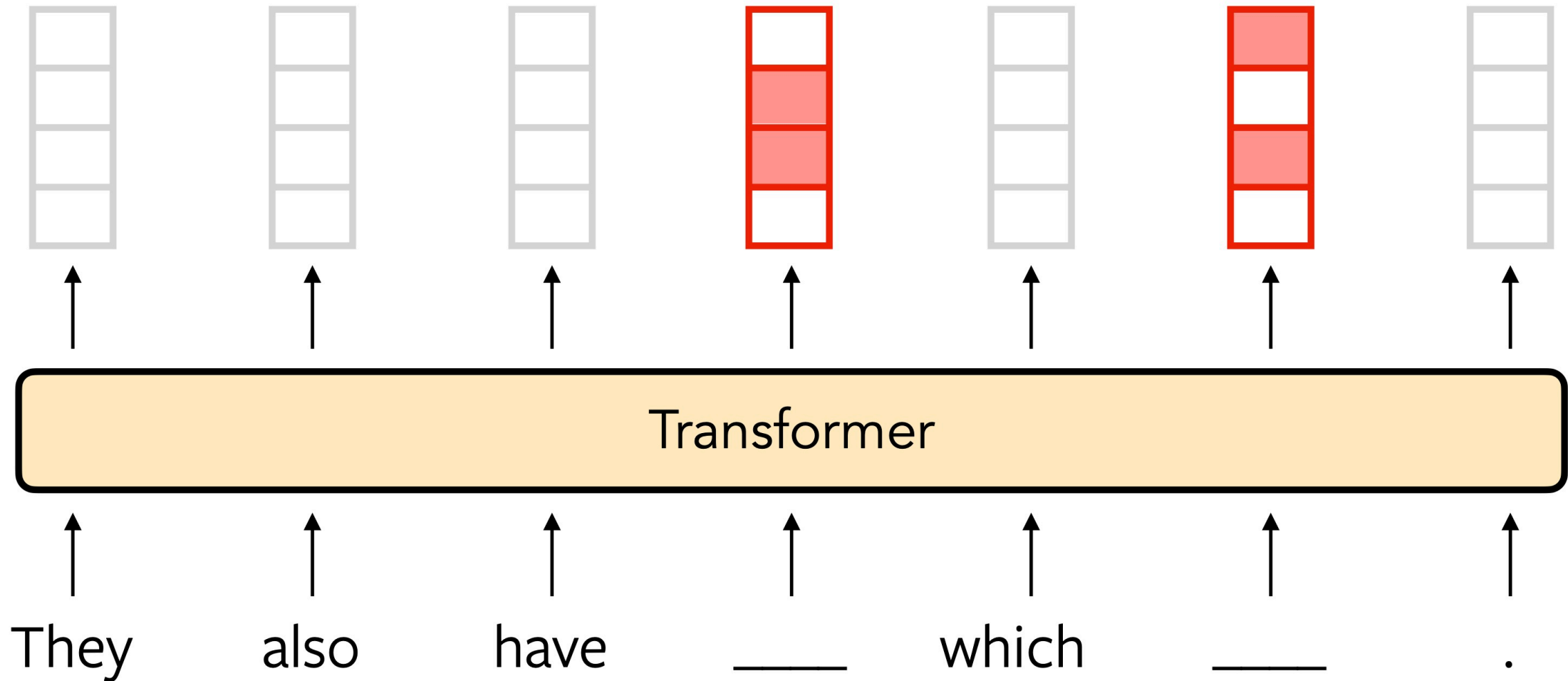
Blank Language Model — Overview

- Dynamic canvas where “___” controls where tokens can be placed
- At each step,
 1. select a “___”
 2. predict a word w
 3. replace that blank with “ w ”, “___ w ”, “ w ___”, or “___ w ___”
- Stop when there is no “___”

Grammar

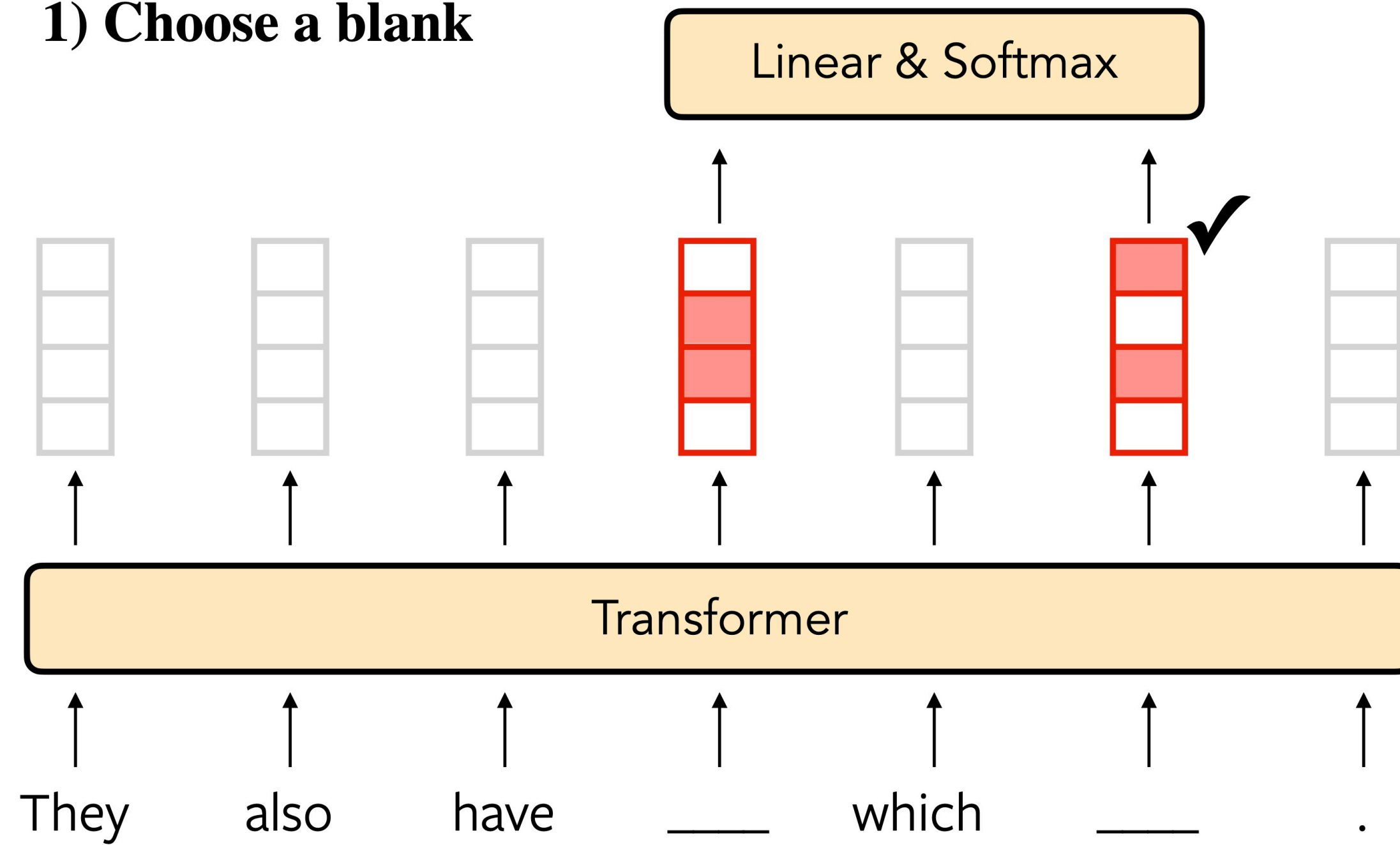
- Nonterminal: ___
- Terminals: $w \in V$
- Production rules: ___ \rightarrow ___? w ___?
(dist. depends on model and context)

Blank Language Model — Architecture

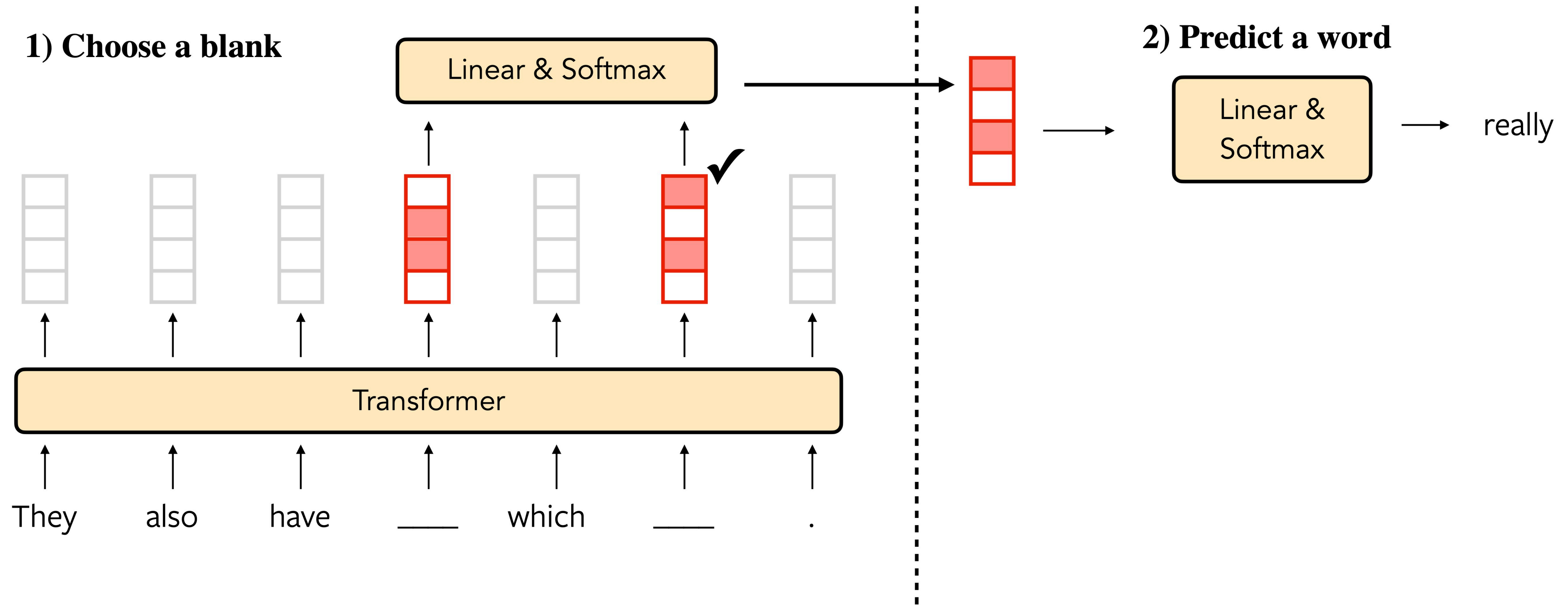


Blank Language Model — Architecture

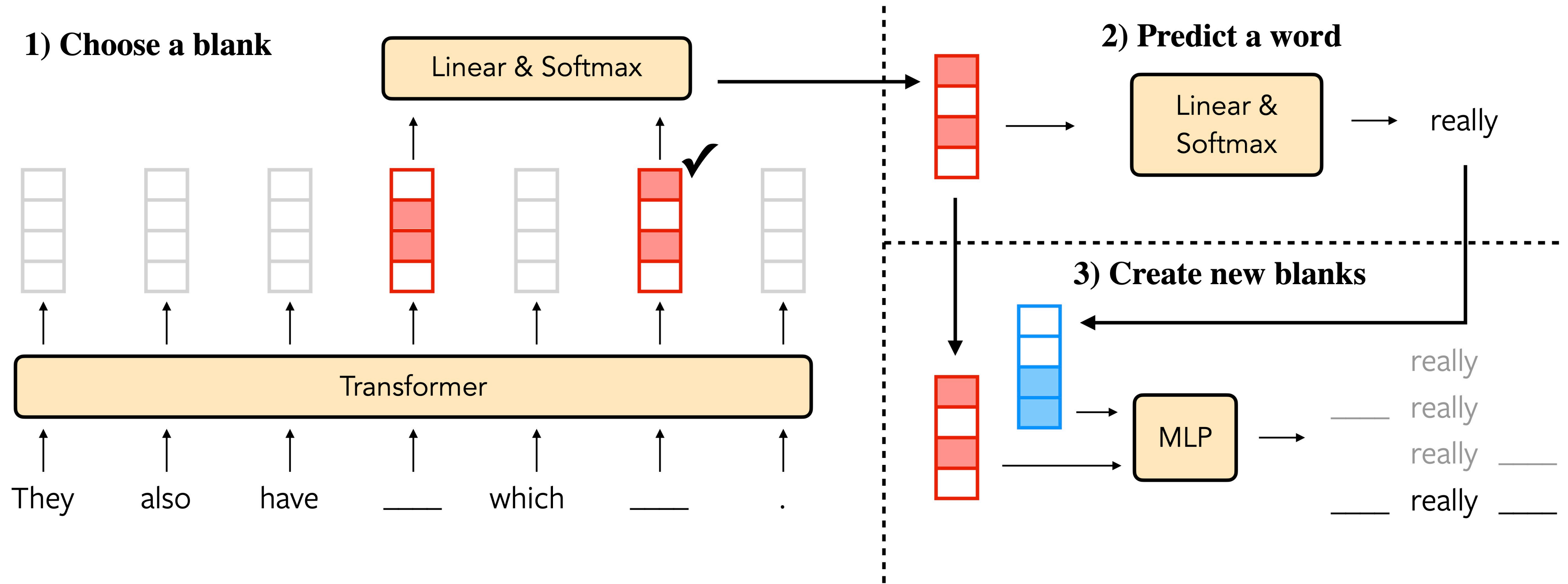
1) Choose a blank



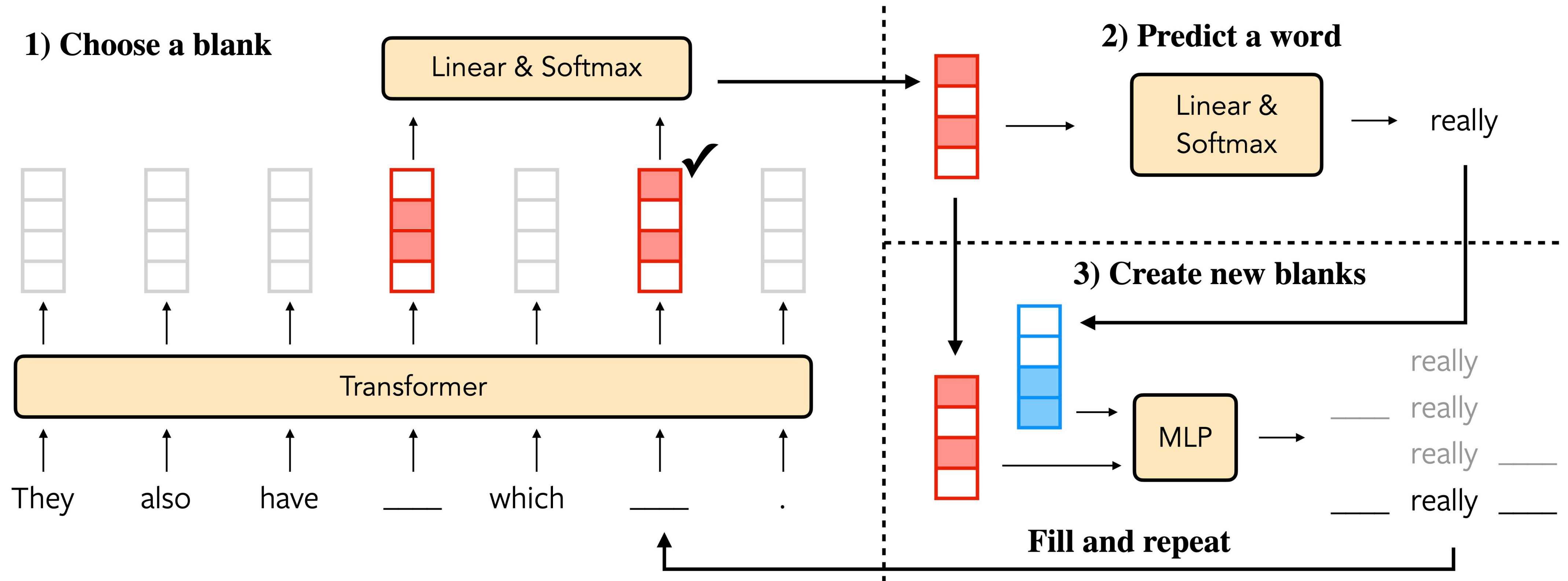
Blank Language Model — Architecture



Blank Language Model — Architecture



Blank Language Model — Architecture



Blank Language Model — Likelihood

Step t	Canvas c	Action a				
		Location b	Word w	(Left blank l ,	Right blank r)	
0.	<u> #1 </u>	#1	3 have	Yes	Yes	
1.	<u> #1 </u> have <u> #2 </u>	#1	1 They	No	Yes	
2.	They <u> #1 </u> have <u> #2 </u>	#2	10 .	Yes	No	
3.	They <u> #1 </u> have <u> #2 </u> .	#2	6 which	Yes	Yes	
4.	They <u> #1 </u> have <u> #2 </u> which <u> #3 </u> .	#1	2 also	No	No	
5.	They also have <u> #1 </u> which <u> #2 </u> .	#2	8 really	Yes	Yes	
6.	They also have <u> #1 </u> which <u> #2 </u> really <u> #3 </u> .	#1	4 ice	No	Yes	
7.	They also have ice <u> #1 </u> which <u> #2 </u> really <u> #3 </u> .	#2	7 is	No	No	
8.	They also have ice <u> #1 </u> which is really <u> #2 </u> .	#1	5 cream	No	No	
9.	They also have ice cream which is really <u> #1 </u> .	#1	9 good	No	No	
10.	They also have ice cream which is really good .			-End-		
						1 2 3 4 5 6 7 8 9 10

trajectory

A sentence x with n words can be realized by $n!$ trajectories, each corresponds to a different word insertion order

$$p(x; \theta) = \sum_{\sigma \in S_n} p(x, \sigma; \theta) = \sum_{\sigma \in S_n} \prod_{t=0}^{n-1} p(a_t^{x, \sigma} | c_t^{x, \sigma}; \theta)$$

↑
↑

order
action, canvas at step t

Blank Language Model — Training

$$\log p(x; \theta) = \log \sum_{\sigma \in S_n} \prod_{t=0}^{n-1} p(a_t^{x, \sigma} | c_t^{x, \sigma}; \theta) \quad \text{intractable}$$

$$\begin{aligned} \downarrow \log \left(\frac{1}{m} \sum_{i=1}^m b_i \right) &\geq \frac{1}{m} \sum_{i=1}^m \log b_i \\ &\geq \log(n!) + \frac{1}{n!} \sum_{\sigma \in S_n} \sum_{t=0}^{n-1} \log p(a_t^{x, \sigma} | c_t^{x, \sigma}; \theta) \end{aligned}$$

Blank Language Model — Training

$$\log p(x; \theta) = \log \sum_{\sigma \in S_n} \prod_{t=0}^{n-1} p(a_t^{x, \sigma} | c_t^{x, \sigma}; \theta) \quad \text{intractable}$$

$$\begin{aligned} &\downarrow \log \left(\frac{1}{m} \sum_{i=1}^m b_i \right) \geq \frac{1}{m} \sum_{i=1}^m \log b_i \\ &\geq \log(n!) + \frac{1}{n!} \sum_{\sigma \in S_n} \sum_{t=0}^{n-1} \log p(a_t^{x, \sigma} | c_t^{x, \sigma}; \theta) \end{aligned}$$

1. Uniformly sample σ from S_n
2. Uniformly sample t from 0 to $n - 1$
3. Construct canvas $c_t^{x, \sigma}$
4. Compute estimated loss $-\log(n!) - n \cdot \log p(a_t^{x, \sigma} | c_t^{x, \sigma}; \theta)$

one action loss per pass :(

Blank Language Model — Training

$c_t^{x,\sigma}$ only depends on $\sigma_{1:t}$

→ combine losses of trajectories with the same first t steps and different $(t + 1)$ -th step

$$\geq \log(n!) + \frac{1}{n!} \sum_{\sigma \in S_n} \sum_{t=0}^{n-1} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)$$

Blank Language Model — Training

$c_t^{x,\sigma}$ only depends on $\sigma_{1:t}$

→ combine losses of trajectories with the same first t steps and different $(t + 1)$ -th step

$$\begin{aligned} &\geq \log(n!) + \sum_{t=0}^{n-1} \frac{1}{n!} \sum_{\sigma \in S_n} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta) \\ &= \log(n!) + n \cdot \mathbb{E}_t \mathbb{E}_{\sigma_{1:t}} \mathbb{E}_{\sigma_{t+1}} \overline{\mathbb{E}_{\sigma_{t+2:n}}} [\log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)] \\ &= \log(n!) + \mathbb{E}_t \mathbb{E}_{\sigma_{1:t}} \left[\frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta) \right] \end{aligned}$$

Blank Language Model — Training

1. Uniformly sample t from 0 to $n - 1$

2. Uniformly sample $\sigma_{1:t}$

3. Construct canvas $c_t^{x,\sigma}$

4. Compute estimated loss $-\log(n!) - \frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)$

$n/2$ action losses per pass :)

$$= \log(n!) + \mathbb{E}_t \mathbb{E}_{\sigma_{1:t}} \left[\frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta) \right]$$

Blank Language Model — Training

1. Uniformly sample t from 0 to $n - 1$

2. Uniformly sample $\sigma_{1:t}$

3. Construct canvas $c_t^{x,\sigma}$

4. Compute estimated loss $-\log(n!) - \frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)$

$n/2$ action losses per pass :)

1 2 3 4 5 6 7 8 9 10

$x =$ They also have ice cream which is really good .

$n = 10$

Blank Language Model — Training

→ 1. Uniformly sample t from 0 to $n - 1$

2. Uniformly sample $\sigma_{1:t}$

3. Construct canvas $c_t^{x,\sigma}$

4. Compute estimated loss $-\log(n!) - \frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)$

$n/2$ action losses per pass :)

1 2 3 4 5 6 7 8 9 10

$x =$ They also have ice cream which is really good .

$n = 10$

$t = 5$

Blank Language Model — Training

1. Uniformly sample t from 0 to $n - 1$

$n/2$ action losses per pass :)

→ 2. Uniformly sample $\sigma_{1:t}$

3. Construct canvas $c_t^{x,\sigma}$

4. Compute estimated loss $-\log(n!) - \frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)$

1 2 3 4 5 6 7 8 9 10

$x =$ They also have ice cream which is really good .

$n = 10$

$t = 5$

$\sigma_{1:t} = (6, 2, 1, 3, 10)$

Blank Language Model — Training

1. Uniformly sample t from 0 to $n - 1$

2. Uniformly sample $\sigma_{1:t}$

→ 3. Construct canvas $c_t^{x,\sigma}$

4. Compute estimated loss $-\log(n!) - \frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)$

n/2 action losses per pass :)

1 2 3 4 5 6 7 8 9 10

$x =$ They also have ice cream which is really good .

$n = 10$

$t = 5$

$\sigma_{1:t} = (6, 2, 1, 3, 10)$

$c_t^{x,\sigma}$ They also have ____ which ____ .

Blank Language Model — Training

1. Uniformly sample t from 0 to $n - 1$

2. Uniformly sample $\sigma_{1:t}$

3. Construct canvas $c_t^{x,\sigma}$

→ 4. Compute estimated loss $-\log(n!) - \frac{n}{n-t} \sum_{\sigma_{t+1}} \log p(a_t^{x,\sigma} | c_t^{x,\sigma}; \theta)$

$n/2$ action losses per pass :)

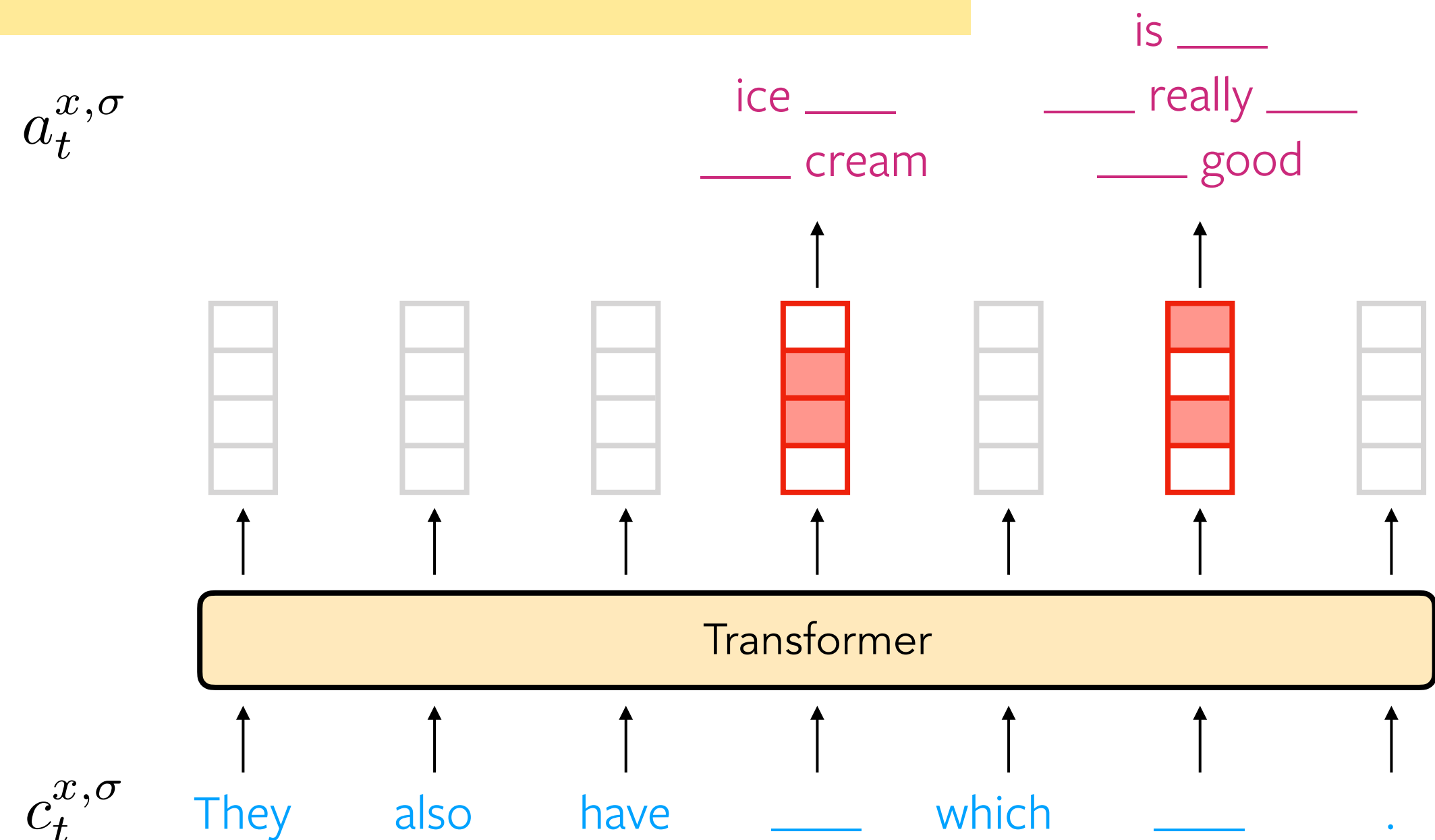
1 2 3 4 5 6 7 8 9 10

$x =$ They also have ice cream which is really good .

$n = 10$

$t = 5$

$\sigma_{1:t} = (6, 2, 1, 3, 10)$



Blank Language Model — Inference

- ✓ Simple greedy decoding or beam search to fill in the blanks in input

Experiments – Overview

Text Infilling

Input: They also have _____ which _____.

Output: They also have ice cream which is really good.

Ancient Text Restoration

Input: τε εγγονον εισαι???????σοφιαι

Output: τε εγγονον εισαιου τουσοφιαι

Sentiment Transfer

Input: The employees were **super nice** and **efficient** !

Output: The employees were rude and unprofessional !

Language Modeling

Output: They also have ice cream which is really good.

Text Infilling — Dataset

- Yahoo Answers dataset (100K documents, max length 200 words)
- Randomly mask tokens with different ratios
- Contiguous masked tokens → “___”

Mask Ratio	when time flies, where does it go? to the center of the universe to be recycled and made into new time.
10%	when time flies, ___ does it go? ___ the center of the ___ to be recycled ___ made into new time.

Mask Ratio	when time flies , where does it go? to the center of the universe to be recycled and made into new time .
50%	when time ___, where ___? ___ the ___ of ___ universe to ___ recycled ___ made into ___.

Text Infilling — Metrics

- Accuracy: BLEU score against original document
- Fluency: perplexity evaluated by a pre-trained left-to-right LM

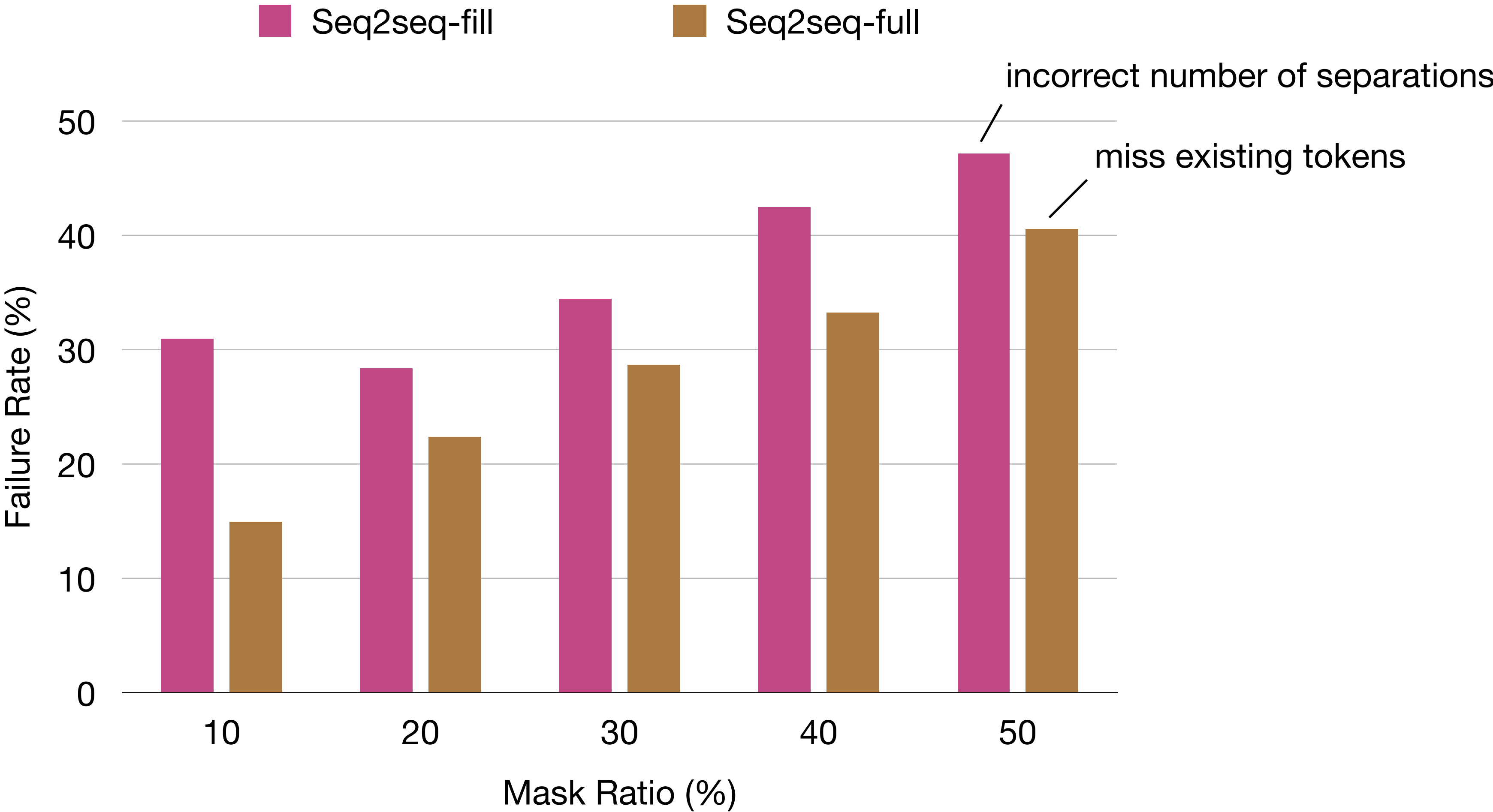
Mask Ratio	when time flies, where does it go? to the center of the universe to be recycled and made into new time.
10%	when time flies, _____ does it go? _____ the center of the _____ to be recycled _____ made into new time.

Mask Ratio	when time flies , where does it go ? to the center of the universe to be recycled and made into new time .
50%	when time _____, where _____? _____ the _____ of _____ universe to _____ recycled _____ made into _____.

Text Infilling — Baselines

- Seq2seq-fill [Donahue et al., 2020]
 - output tokens to fill in the blanks, separated by “|”
- Seq2seq-full [Donahue et al., 2020]
 - output full document from input

Text Infilling — Results



Text Infilling — Baselines

- BERT+LM
 - feed BERT representation of each blank to left-to-right LM that learns to generate tokens in that blank
 - at test time, fill in the blanks one by one

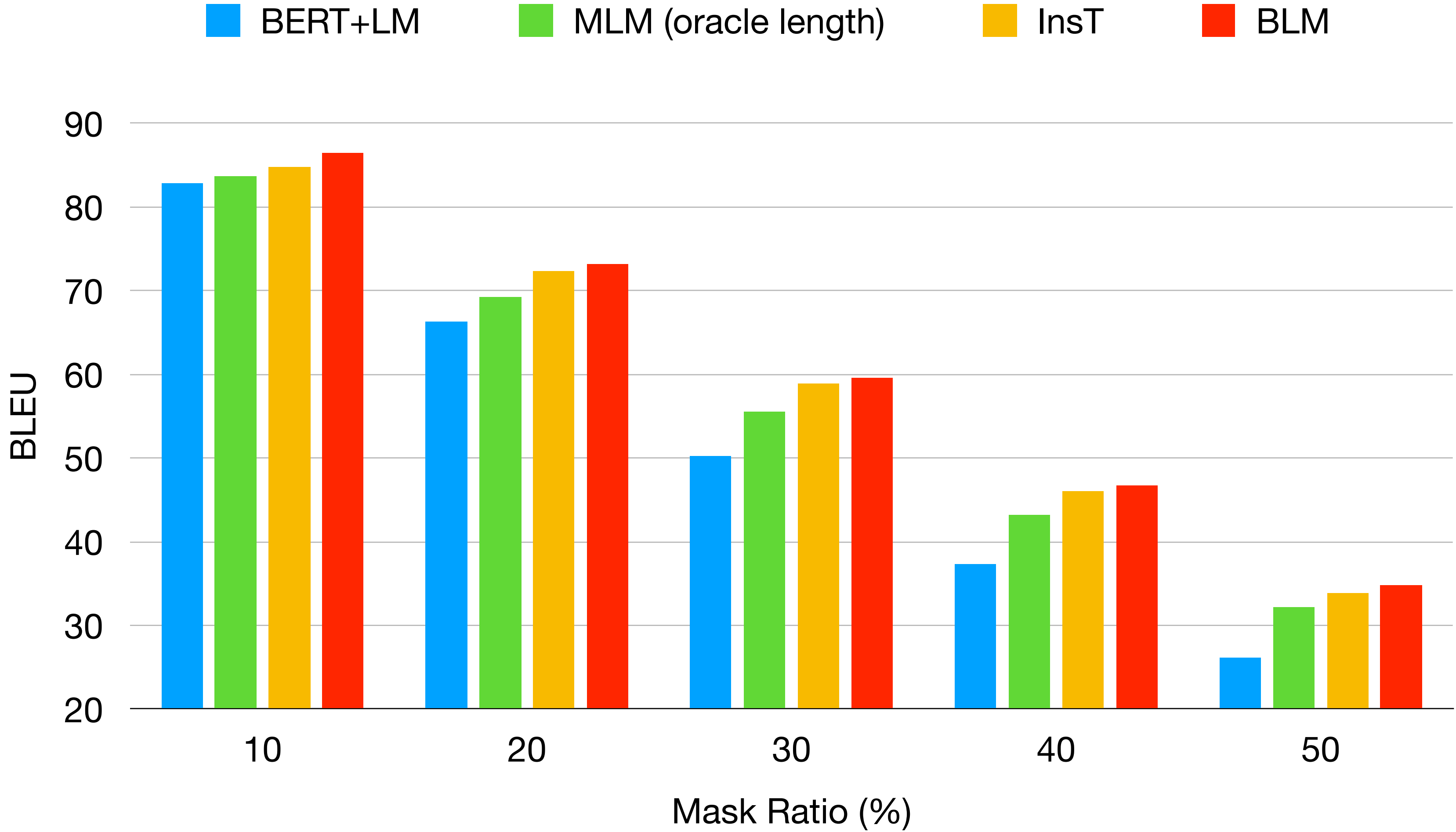
Text Infilling — Baselines

- BERT+LM
- Masked Language Model (MLM) with oracle length
 - replace blanks with the target number of masks
 - fill the masks autoregressively by most-confident-first heuristic

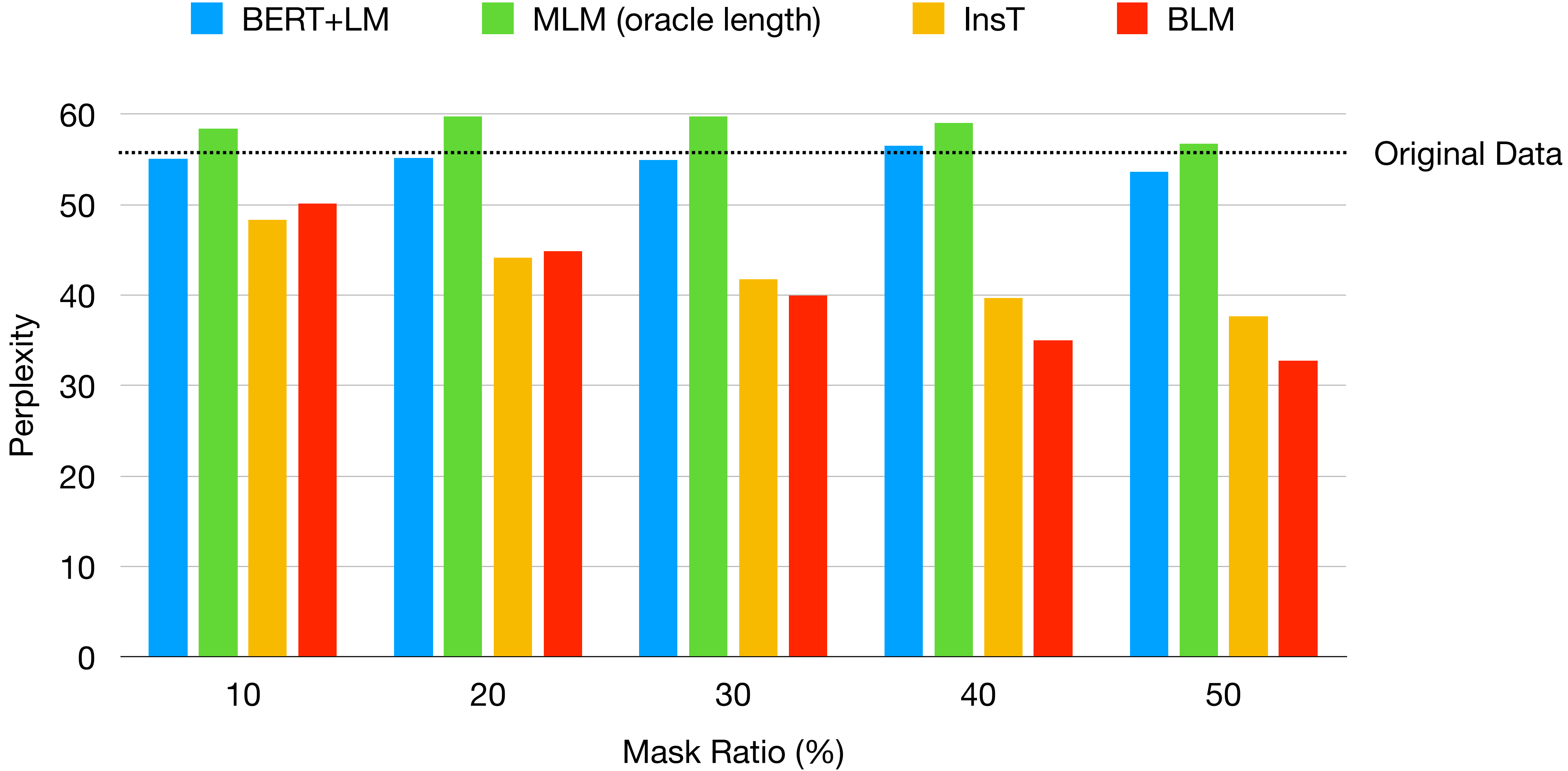
Text Infilling — Baselines

- BERT+LM
- Masked Language Model (MLM) with oracle length
- Insertion Transformer [Stern et al., 2019]
 - cannot specify insertion position
 - force it to generate at valid locations

Text Infilling — Results



Text Infilling — Results



Text Infilling — Examples

Original when time flies, **where** does it go? **to** the center of the **universe** to be recycled **and** made into new time.

Blanked when time flies, _____ does it go? _____ the center of the _____ to be recycled _____ made into new time.

BERT+LM when time flies, where does it go? to the center of the earth to be recycled came made into new time.

MLM (oracle len) when time flies, where does it go? from the center of the earth to be recycled converted made into new time.

InsT when time flies, where does it go? for the center of the earth has to be recycled and made into new time.

BLM when time flies, where does it go? for the center of the earth to be recycled and made into new time.

Mask Ratio 10%

Text Infilling — Examples

Original	when time flies , where does it go? to the center of the universe to be recycled and made into new time .
Blanked	when time _____, where _____? _____ the _____ of _____ universe to _____ recycled _____ made into _____.

BERT+LM	when time <u>is</u> , where <u>to</u> ? <u>i need to find the way of the</u> universe to <u>be</u> recycled <u>and</u> made into <u>a lot</u> .
MLM (oracle len)	when time <u>is</u> , where <u>is the universe</u> ? <u>from the creation of the</u> universe to <u>be</u> recycled <u>and</u> made into <u>the universe</u> .
InsT	when time <u>was created</u> , where <u>was it</u> ? <u>what was the name of the</u> universe to <u>be</u> recycled <u>and</u> made into <u>space</u> .
BLM	when time <u>was created</u> , where <u>did it come from</u> ? <u>it was the first part of the</u> universe to <u>be</u> recycled <u>and</u> made into <u>space</u> .

Mask Ratio 50%

Ancient Text Restoration — Setup

Ancient Greek Inscriptions dataset (18M characters / 3M words) [Assael et al., 2019]

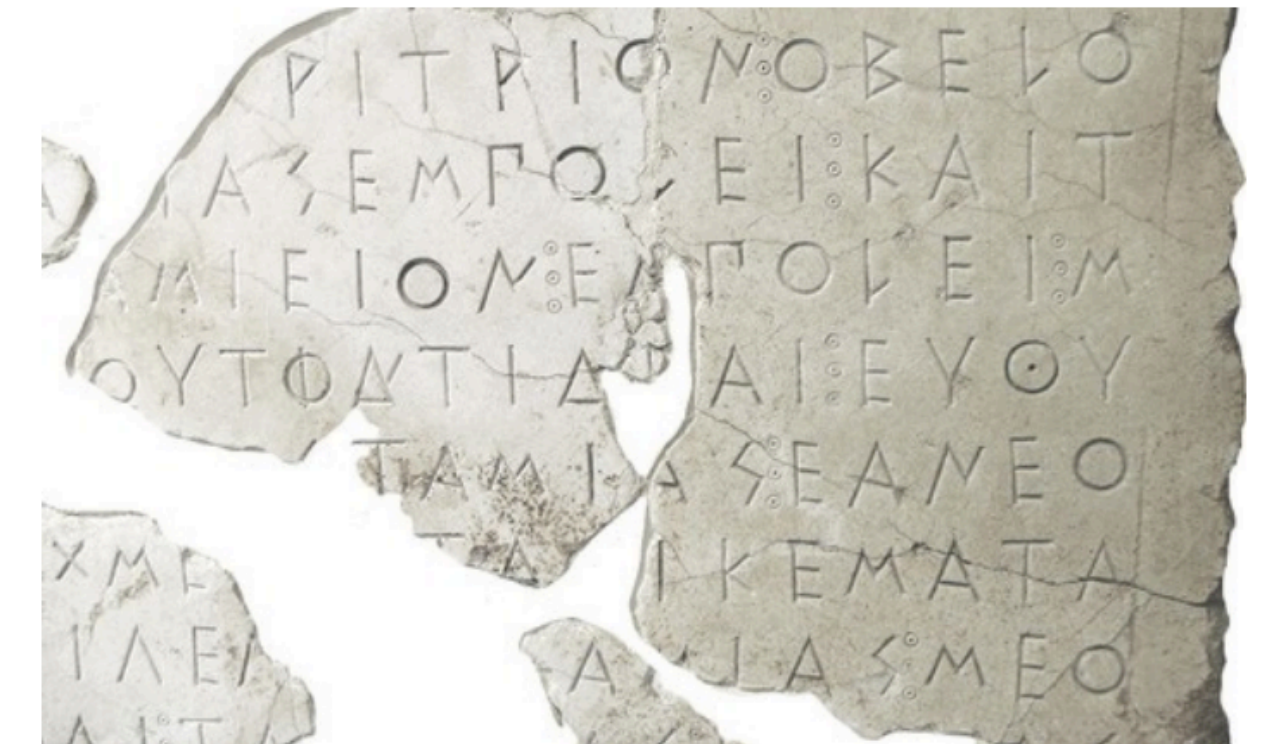
- number of characters to recover is assumed to be known

Length-aware BLM (L-BLM)

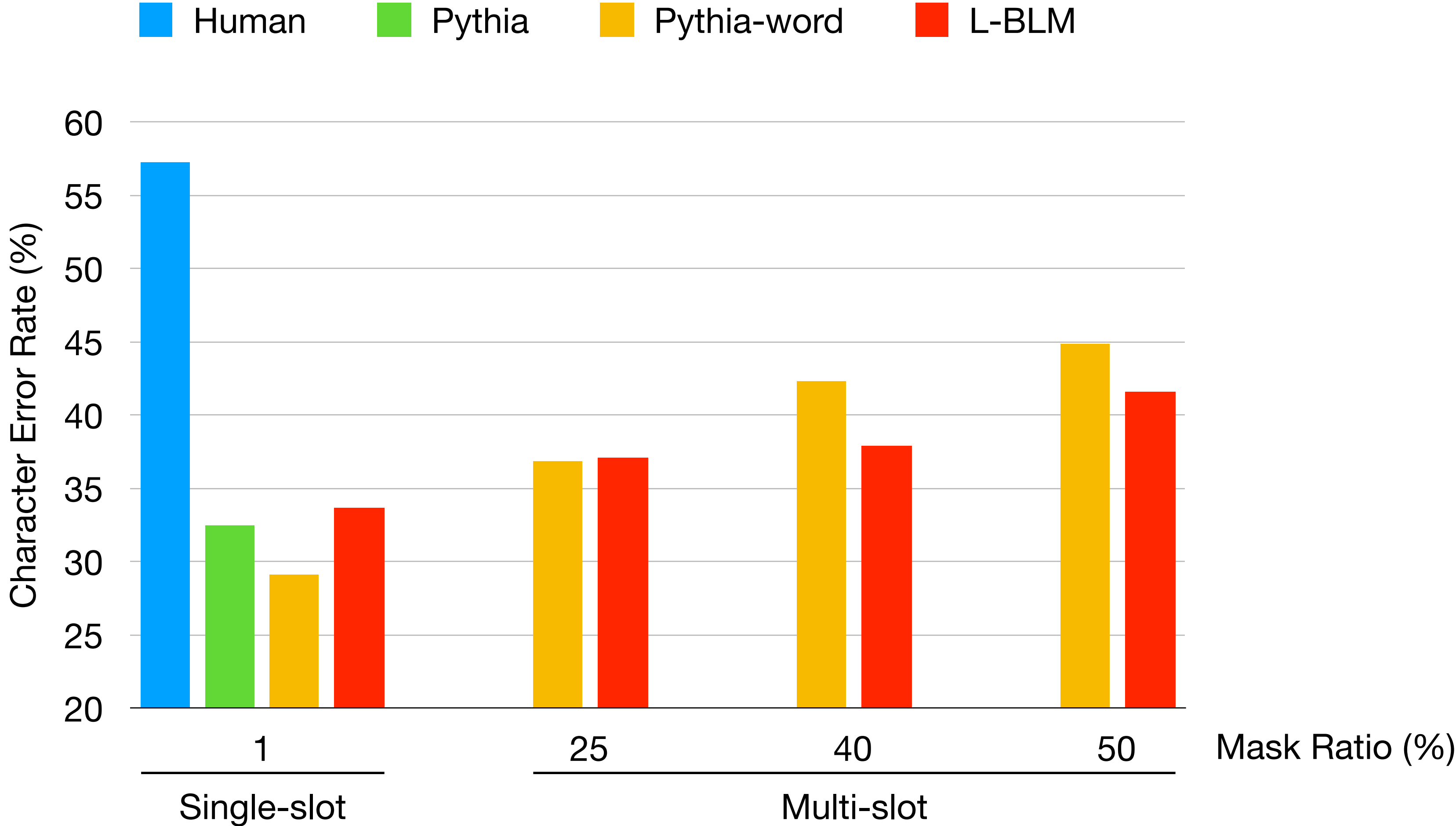
- [t] → [k] w [t-1-k]

Baselines [Assael et al., 2019]:

- Pythia: character-level seq2seq model to fill in one slot at a time
- Pythia-word: use both character and word representations



Ancient Text Restoration — Results



Sentiment Transfer — Approach

1. Remove expressions of high polarity
 - train a sentiment classifier and mask words with attention weight above average
2. Complete the partial sentence with expressions of the target sentiment
 - train two instances of BLM, one for each sentiment

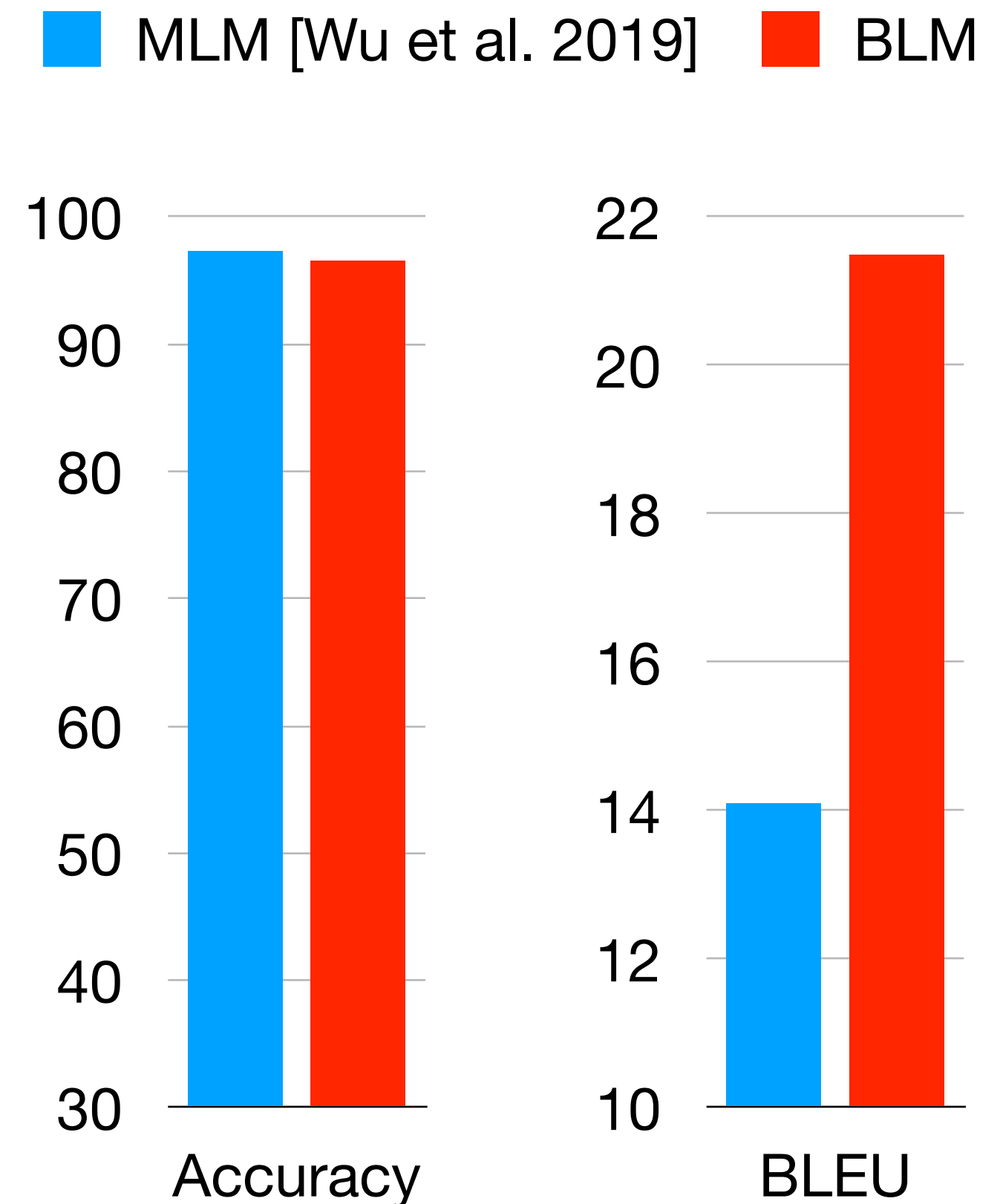
Sentiment Transfer — Results

Input everyone that i spoke with was **very helpful** and **kind** .
BLM everyone that i spoke with was rude and unprofessional .

Input there is **definitely not** enough **room** in that part of the venue .
BLM there is always enough parking in that part of the venue .

Input it is n't **terrible** , but it is **n't** very good either .
BLM it is n't fancy , but it is still very good either .

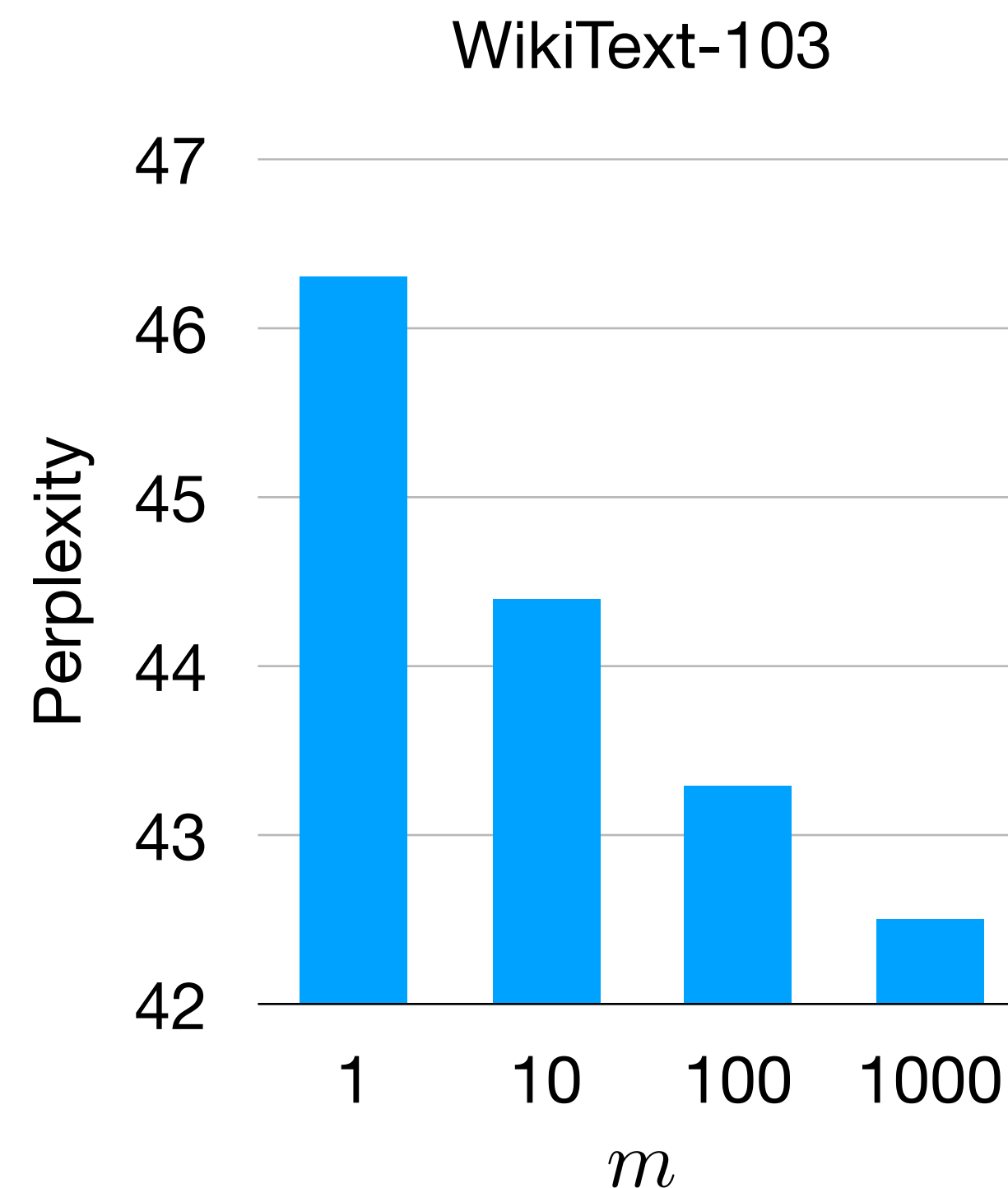
Yelp Reviews Dataset



Language Modeling — Estimation

Monte-Carlo sampling $p(x; \theta) = \sum_{\sigma \in S_n} p(x, \sigma; \theta) \leftarrow \frac{n!}{m} \sum_{i=1}^m p(x, \sigma_i; \theta)$

- estimated perplexity is likely to be higher than actual perplexity
- as m increases, it converges to actual perplexity



Language Modeling — Results

Datasets: Penn Treebank (1M tokens), WikiText-2 (2M), WikiText-103 (103M)

	PTB	WT2	WT103
LSTM (Grave et al., 2016)	82.3	99.3	48.7
TCN (Bai et al., 2018)	88.7	-	45.2
AWD-LSTM (Merity et al., 2017)	57.3	65.8	-
Transformer (Dai et al., 2019)	-	-	30.1
Adaptive (Baeovski and Auli, 2018)	-	-	18.7
Transformer-XL (Dai et al., 2019)	54.5	-	18.3
InsT (our implementation)	77.3	91.4	39.4
BLM	69.2	81.2	42.5

Room for improvements!

Summary

Input: They also have _____ which _____.

Output: They also have ice cream which is really good.

https://github.com/Varal7/blank_language_model

Thank you!

- Dynamically create and fill in blanks
- Effective on text infilling, ancient text restoration, style transfer

More Applications

- Template filling, information fusion, assisting human writing...
- Rewrite to mitigate toxicity and bias
- Representation learning

Extensions

- Add representation for blanks
- Conditional BLM: edit and refine machine translation, dialogue system...