

XLNet

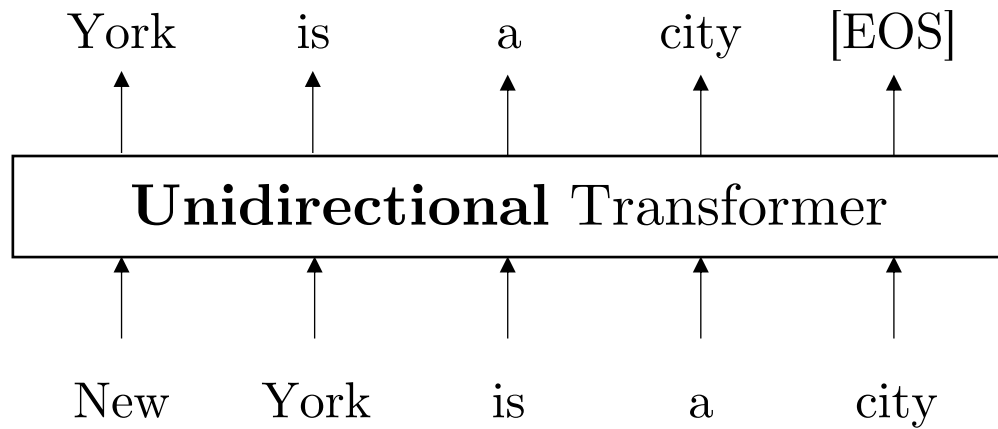
Zhilin Yang*, Zihang Dai*, Yiming Yang, Jaime Carbonell,
Ruslan Salakhutdinov, Quoc V. Le
(*: equal contribution)

Language Pretraining: Related Work

- RBMs (Salakhutdinov et al 2007), Autoencoders (Vincent et al 2008), Jigsaw (Noroozi and Favaro 2016), GANs (Donahue and Simonyan 2019)...
- word2vec (Mikolov et al 2013), GloVe (Pennington et al 2014)
- Semi-supervised sequence learning (Dai and Le 2015), ELMo (Peters et al 2017), CoVe (McCann et al 2017), GPT (Radford et al 2018), BERT (Devlin et al 2018)...

Two Notable Objectives for Language Pretraining

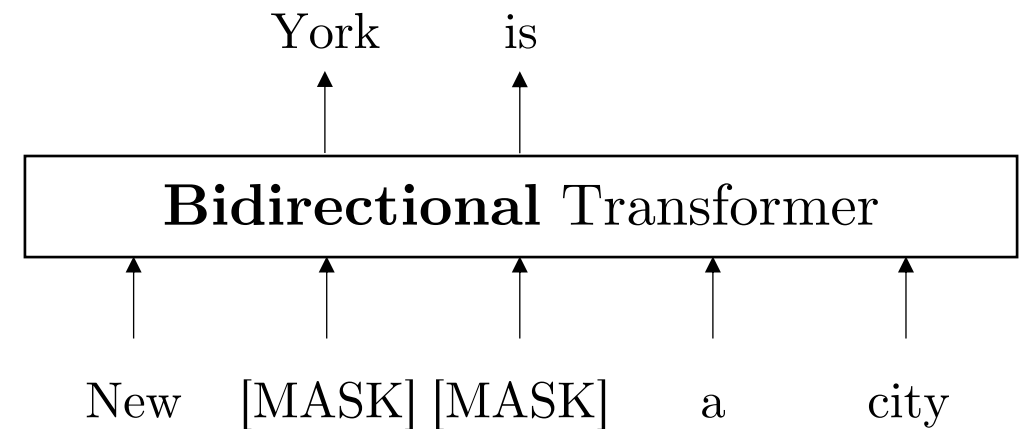
Auto-regressive Language Modeling



$$\log p(\mathbf{x}) = \sum_{t=1}^T \log p(x_t | \mathbf{x}_{<t})$$

- Next-token prediction

Denoising Auto-encoding (BERT)

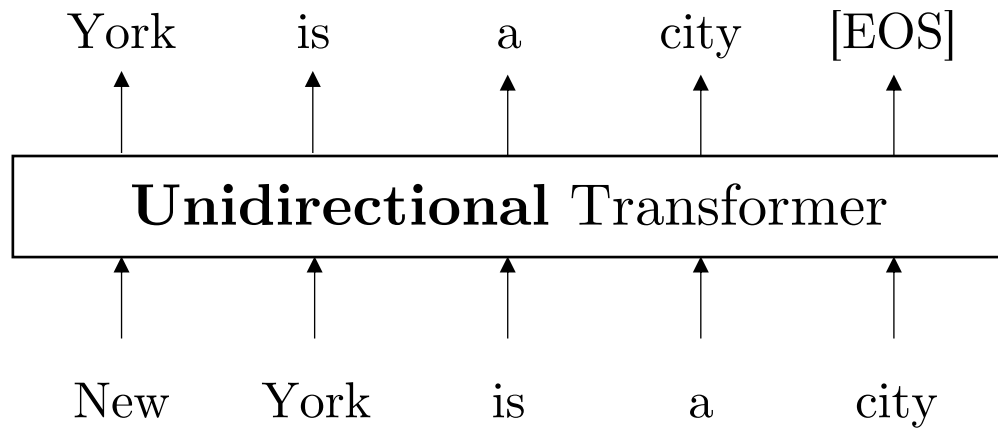


$$\log p(\bar{\mathbf{x}} | \hat{\mathbf{x}}) = \sum_{t=1}^T \text{mask}_t \log p(x_t | \hat{\mathbf{x}})$$

- Reconstruct masked tokens

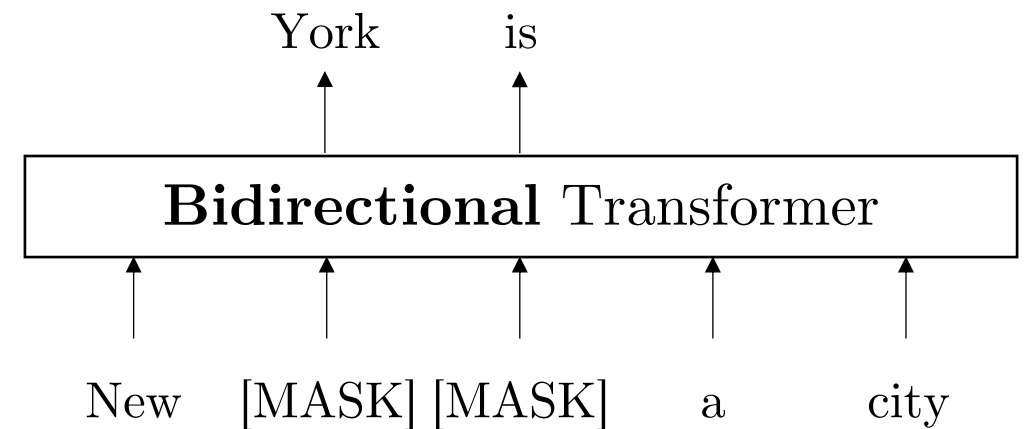
Two Notable Objectives for Language Pretraining

Auto-regressive Language Modeling



 No **Bidirectional** Context

Denoising Auto-encoding (BERT)

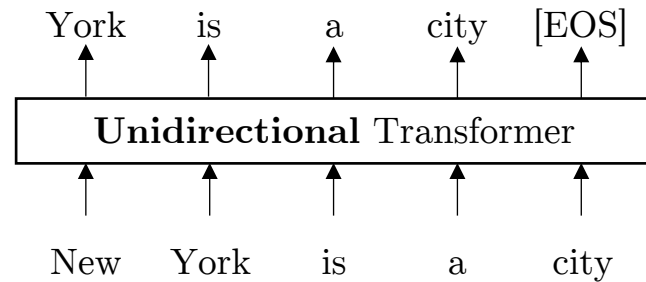


 Independent Predictions

 Artificial **Noise**: [MASK]

Two Notable Objectives for Language Pretraining

Auto-regressive Language Modeling

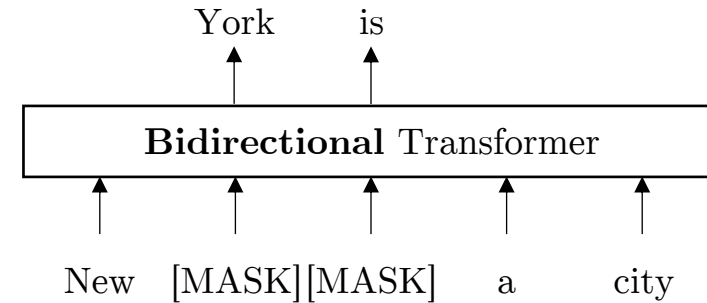


😊 Full Auto-regressive **Dependence**

😊 Free from artificial **Noise**

😞 No **Bidirectional Context**

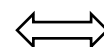
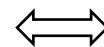
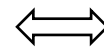
Denoising Auto-encoding (BERT)



😞 **Independent** Predictions

😞 Artificial **Noise**: [MASK]

😊 Natural **Bidirectional Context**



Desire: Combine the Pros and Remove the Cons

😊 Full Auto-regressive **Dependence**

😊 Free from **Noise**

😊 Natural **Bidirectional Context**

Desire: Combine the Pros and Remove the Cons



Full Auto-regressive Dependence

XLNet

- An **auto-regressive** model that captures **bidirectional context**



Natural Bidirectional Context

Context Depends on the Factorization Order

- **Standard LM:** Left-to-right factorization $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$

$$P(\mathbf{x}) = P(x_1)P(x_2 \mid \mathbf{x}_1)P(x_3 \mid \mathbf{x}_{1,2})P(x_4 \mid \mathbf{x}_{1,2,3}) \cdots$$

x_1

x_2

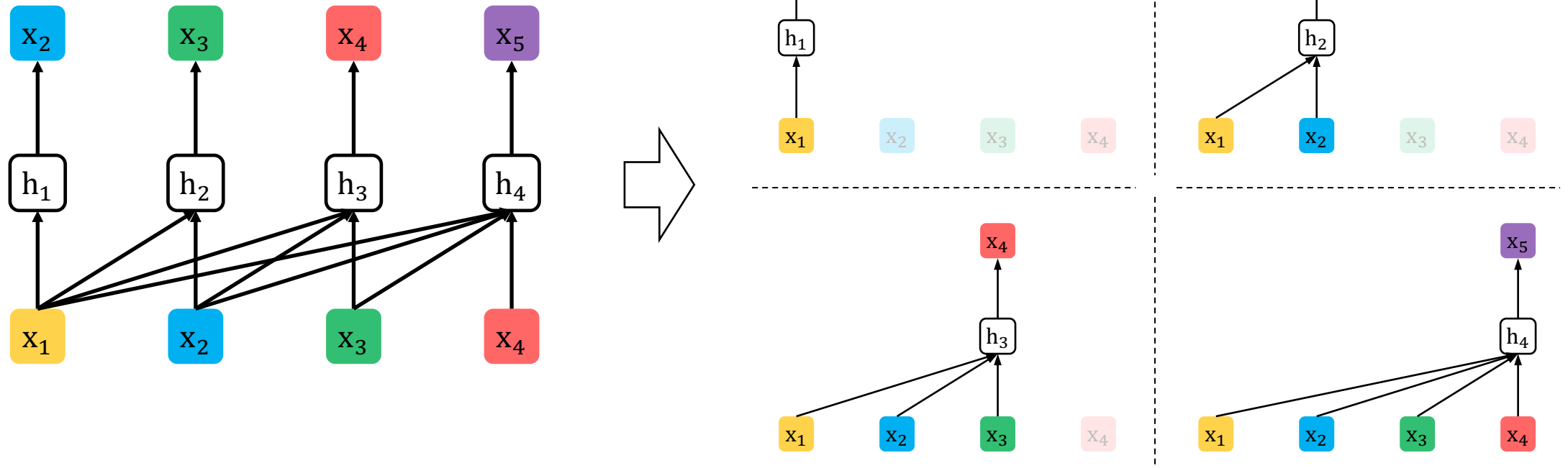
x_3

x_4

Context Depends on the Factorization Order

- **Standard LM:** Left-to-right factorization $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$

$$P(\mathbf{x}) = P(x_1)P(x_2 \mid \mathbf{x}_1)P(x_3 \mid \mathbf{x}_{1,2})P(x_4 \mid \mathbf{x}_{1,2,3}) \cdots$$



Context Depends on the Factorization Order

- Change the Factorization order to: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

$$P(\mathbf{x}) = P(x_4)P(x_1 \mid \mathbf{x}_4)P(x_3 \mid \mathbf{x}_{1,4})P(x_2 \mid \mathbf{x}_{1,2,4}) \cdots$$

x_1

x_2

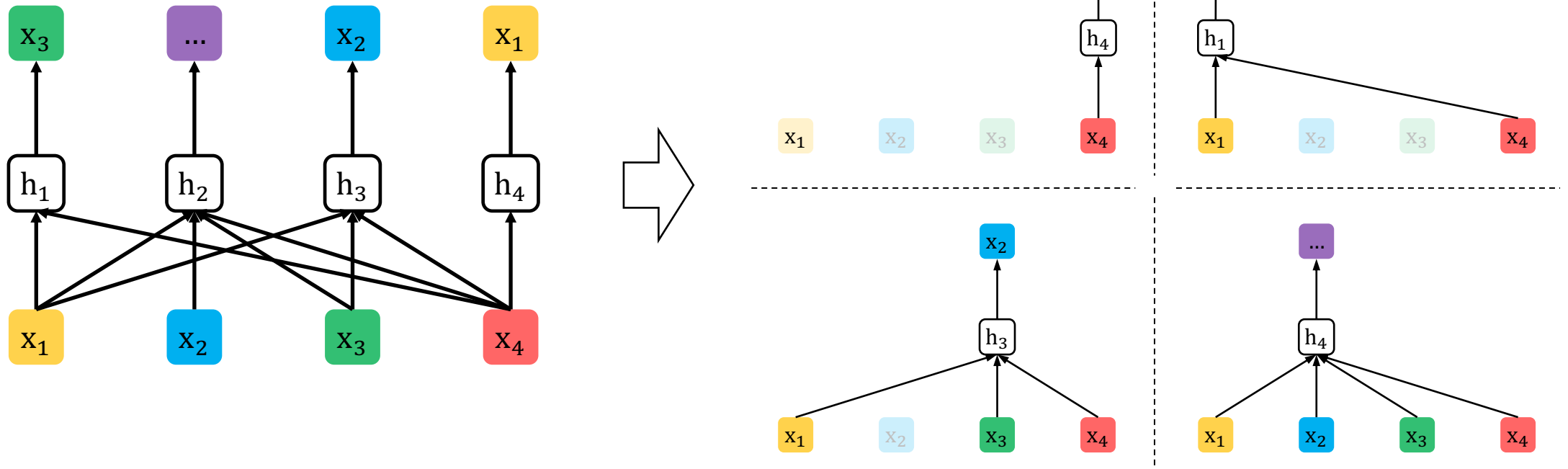
x_3

x_4

Context Depends on the Factorization Order

- Change the Factorization order to: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

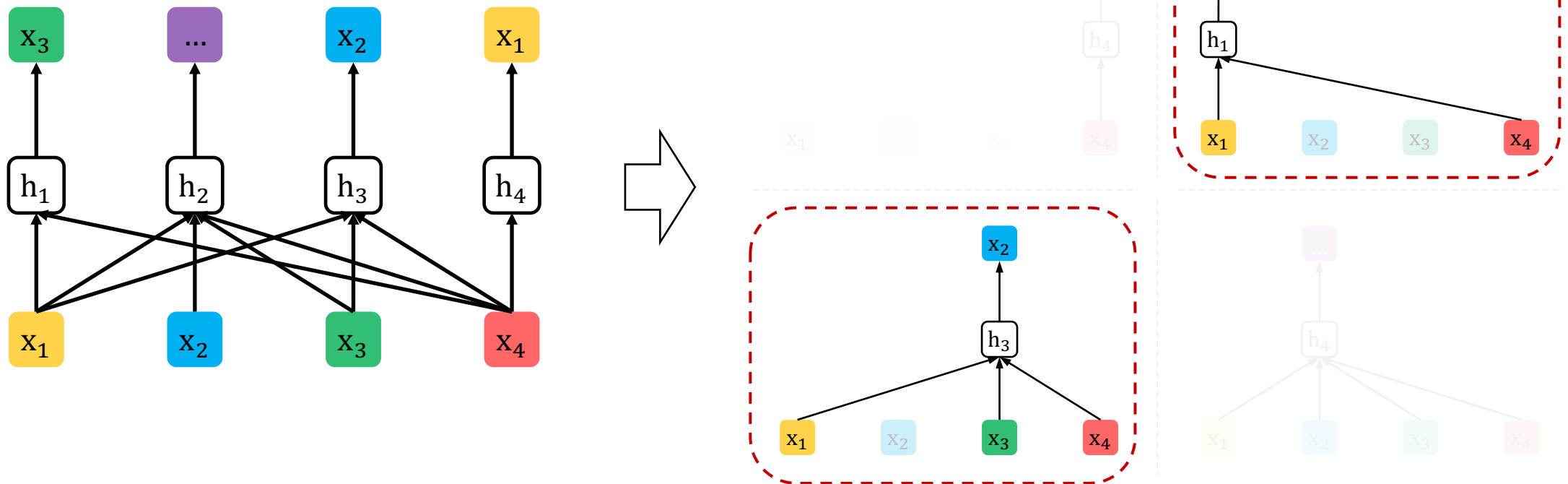
$$P(\mathbf{x}) = P(x_4)P(x_1 \mid \mathbf{x}_4)P(x_3 \mid \mathbf{x}_{1,4})P(x_2 \mid \mathbf{x}_{1,2,4}) \cdots$$



Context Depends on the Factorization Order

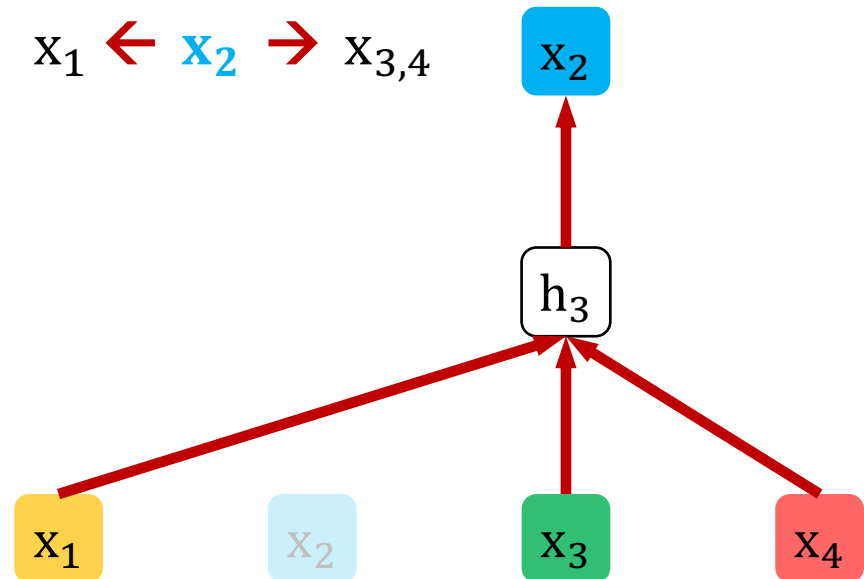
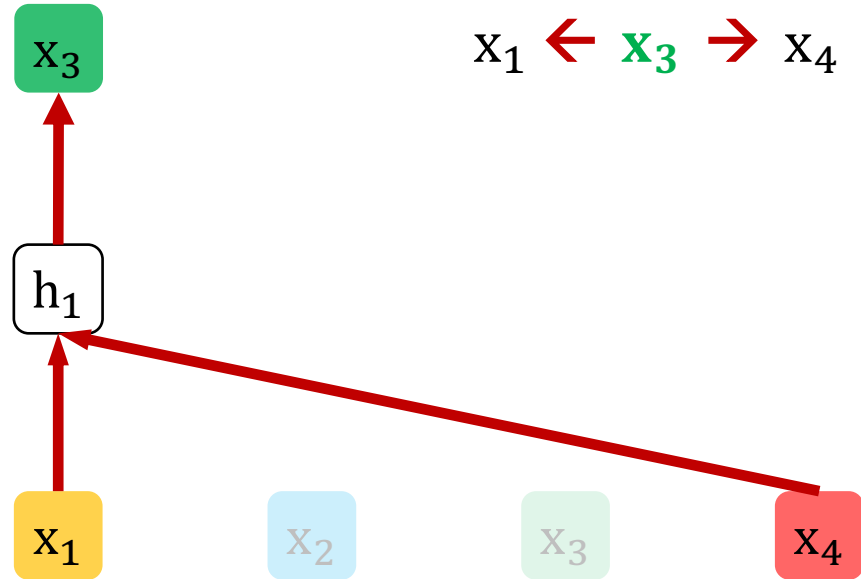
- Change the Factorization order to: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

$$P(\mathbf{x}) = P(x_4)P(x_1 \mid \mathbf{x}_4)P(x_3 \mid \mathbf{x}_{1,4})P(x_2 \mid \mathbf{x}_{1,2,4}) \cdots$$



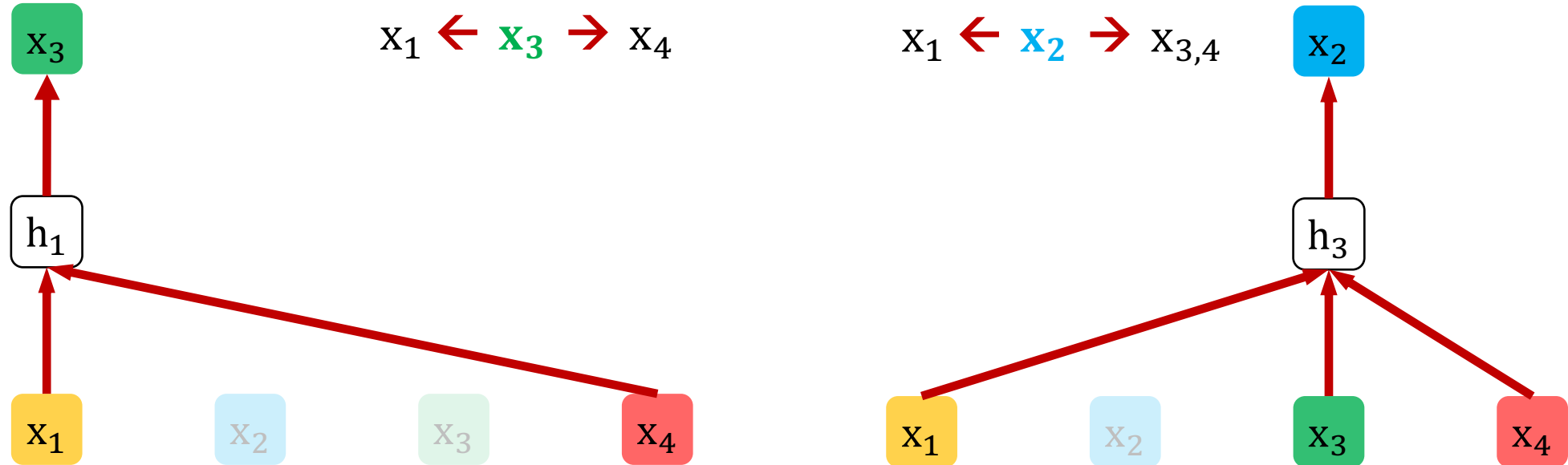
Bidirectional Context via Factorization Order

- Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$



Bidirectional Context via Factorization Order

- Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$



Bidirectional Context

Permutation Language Modeling

- Given a sequence \mathbf{x} of length T
- Uniformly sample a factorization order \mathbf{z} from all possible permutations
- Maximize the permuted log-likelihood

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{Z}_T} [\log P(\mathbf{x} \mid \mathbf{z})] = \mathbb{E}_{\mathbf{z} \sim \mathcal{Z}_T} \left[\sum_{t=1}^T P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) \right]$$

More examples

Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

x_1

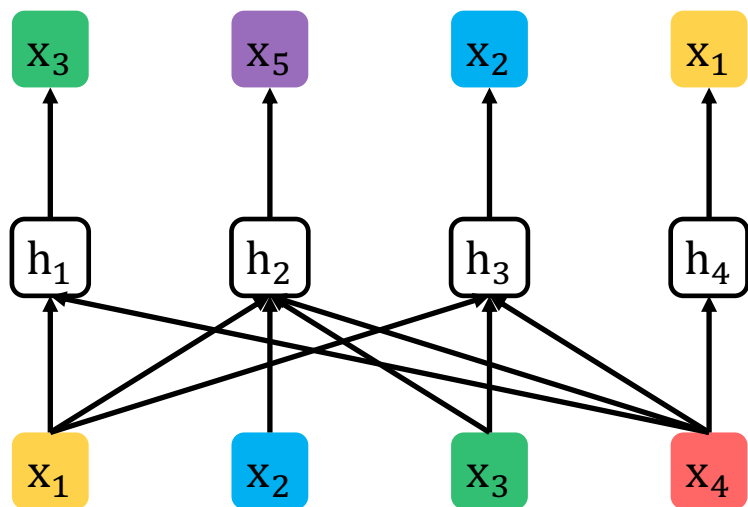
x_2

x_3

x_4

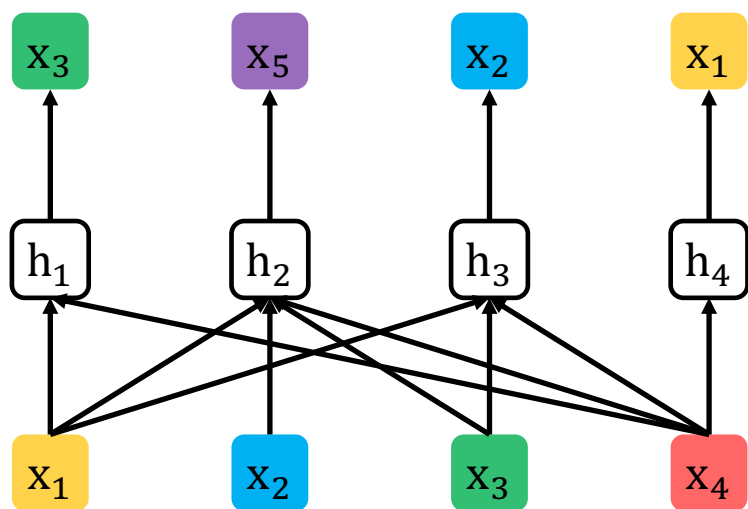
More examples

Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$



More examples

Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

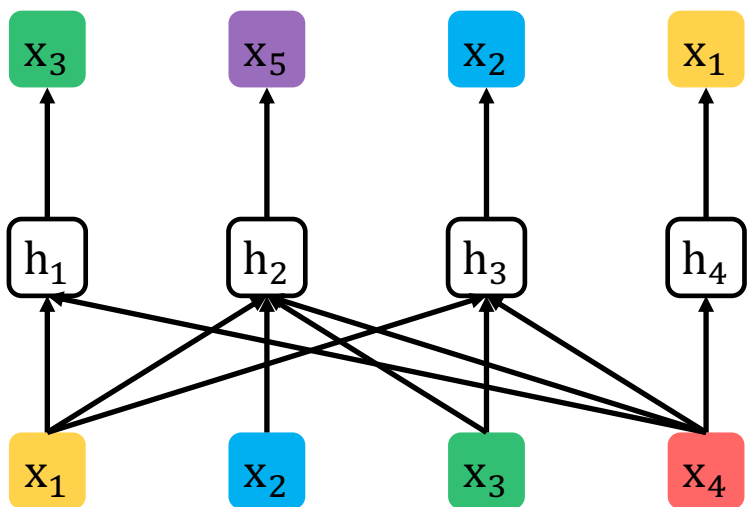


Factorization order: $2 \rightarrow 4 \rightarrow 1 \rightarrow 3$

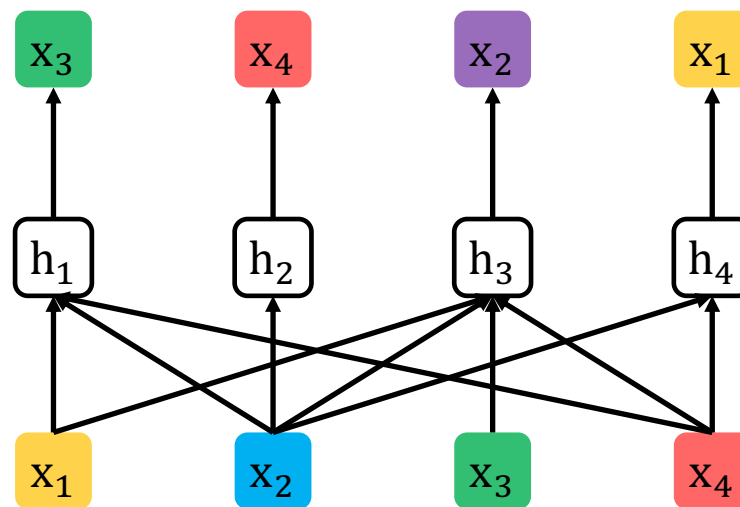


More examples

Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

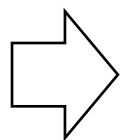


Factorization order: $2 \rightarrow 4 \rightarrow 1 \rightarrow 3$



Target-position-aware Distribution

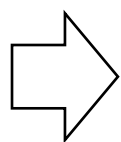
$$\mathbb{E}_{z_t \sim \mathcal{Z}_T} \left[\sum_{t=1}^T \log P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) \right]$$



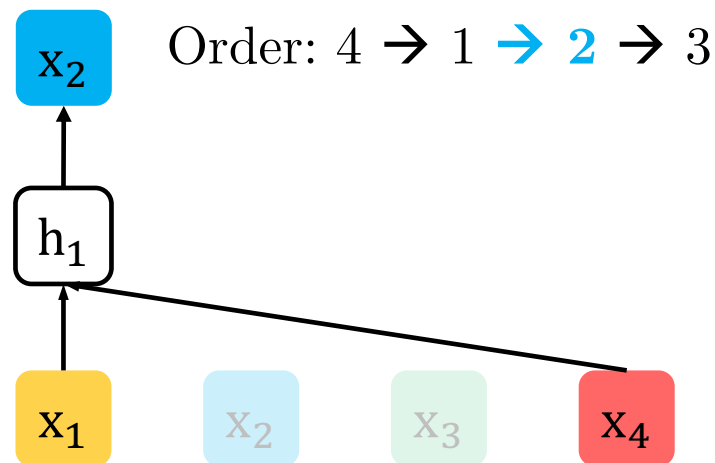
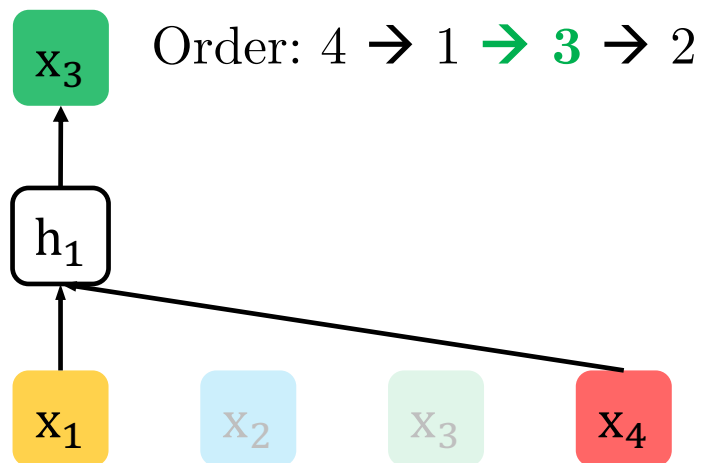
The distribution $P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t)$ must
condition on the target position z_t

Target-position-aware Distribution

$$\mathbb{E}_{z_t \sim \mathcal{Z}_T} \left[\sum_{t=1}^T \log P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) \right]$$

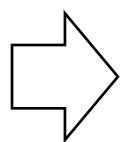


The distribution $P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t)$ must
condition on the target position z_t

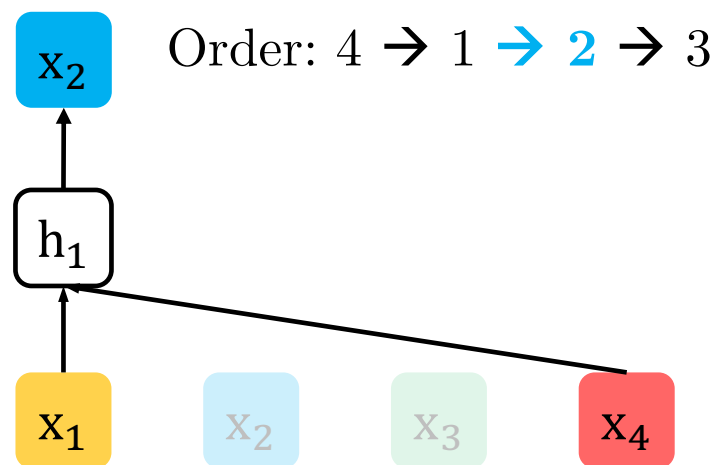
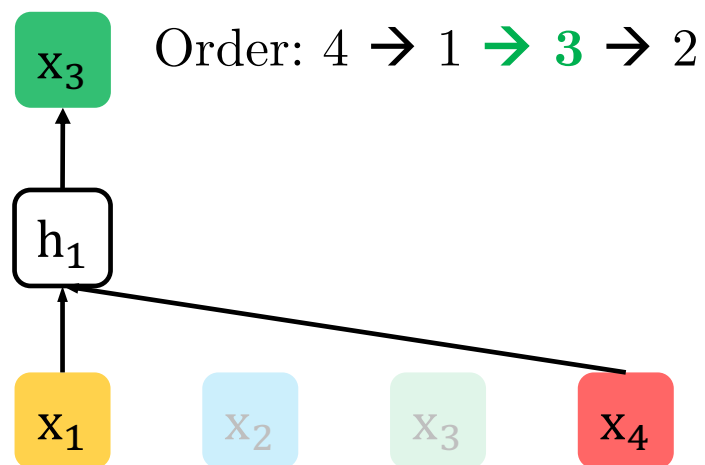


Target-position-aware Distribution

$$\mathbb{E}_{z_t \sim \mathcal{Z}_T} \left[\sum_{t=1}^T \log P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) \right]$$



The distribution $P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t)$ must **condition on the target position z_t**

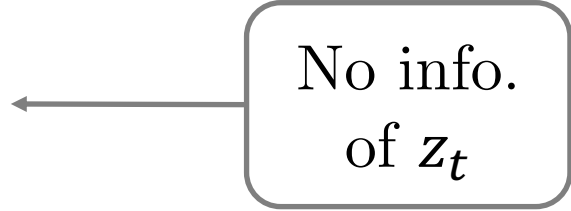


- Predicting **position 3** and **position 2** requires different prediction distributions
- The prediction distribution should **change according to the target position**

Reparameterization

- Standard Softmax does **NOT** work

$$P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) = \frac{\exp(e(x_{z_t})^\top h(\mathbf{x}_{\mathbf{z}_{<t}}))}{\sum_{x'} \exp(e(x')^\top h(\mathbf{x}_{\mathbf{z}_{<t}}))}$$




No info.
of z_t

Reparameterization

- Standard Softmax does **NOT** work

$$P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) = \frac{\exp(e(x_{z_t})^\top h(\mathbf{x}_{\mathbf{z}_{<t}}))}{\sum_{x'} \exp(e(x')^\top h(\mathbf{x}_{\mathbf{z}_{<t}}))}$$

No info.
of z_t



- Proposed** solution: incorporate z_t into **hidden states**

$$P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) = \frac{\exp(e(x_{z_t})^\top g(z_t, \mathbf{x}_{\mathbf{z}_{<t}}))}{\sum_{x'} \exp(e(x')^\top g(z_t, \mathbf{x}_{\mathbf{z}_{<t}}))}$$

Deep Net



Reparameterization

- Standard Softmax does **NOT** work

$$P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) = \frac{\exp(e(x_{z_t})^\top h(\mathbf{x}_{\mathbf{z}_{<t}}))}{\sum_{x'} \exp(e(x')^\top h(\mathbf{x}_{\mathbf{z}_{<t}}))}$$

No info.
of z_t

- Proposed** solution: incorporate \mathbf{z}_t into **hidden states**

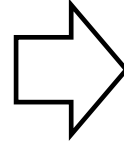
$$P(x_{z_t} \mid \mathbf{x}_{\mathbf{z}_{<t}}, z_t) = \frac{\exp(e(x_{z_t})^\top g(\mathbf{z}_t, \mathbf{x}_{\mathbf{z}_{<t}}))}{\sum_{x'} \exp(e(x')^\top g(\mathbf{z}_t, \mathbf{x}_{\mathbf{z}_{<t}}))}$$

Deep Net

Question: how to implement $g(\mathbf{z}_t, \mathbf{x}_{\mathbf{z}_{<t}})$?

Target Position Aware Representation: $g(\mathbf{z}_t, \mathbf{x}_{\mathbf{z}_{<t}})$

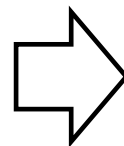
Reuse the Idea of Attention



- Stand at the target position \mathbf{z}_t
- Gather information from $\mathbf{x}_{\mathbf{z}_{<t}}$

Target Position Aware Representation: $g(\mathbf{z}_t, \mathbf{x}_{\mathbf{z}_{<t}})$

Reuse the Idea of Attention

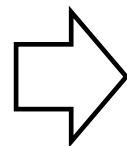


- Stand at the target position \mathbf{z}_t
- Gather information from $\mathbf{x}_{\mathbf{z}_{<t}}$

$$g(z_t, \mathbf{x}_{\mathbf{z}_{<t}}) = \text{Attn}_\theta \left(\underbrace{Q = \text{Enc}(\mathbf{z}_t)}_{\text{Stand at } \mathbf{z}_t}, \underbrace{KV = \mathbf{h}(\mathbf{x}_{\mathbf{z}_{<t}})}_{\text{Gather info. from } \mathbf{x}_{\mathbf{z}_{<t}}} \right)$$

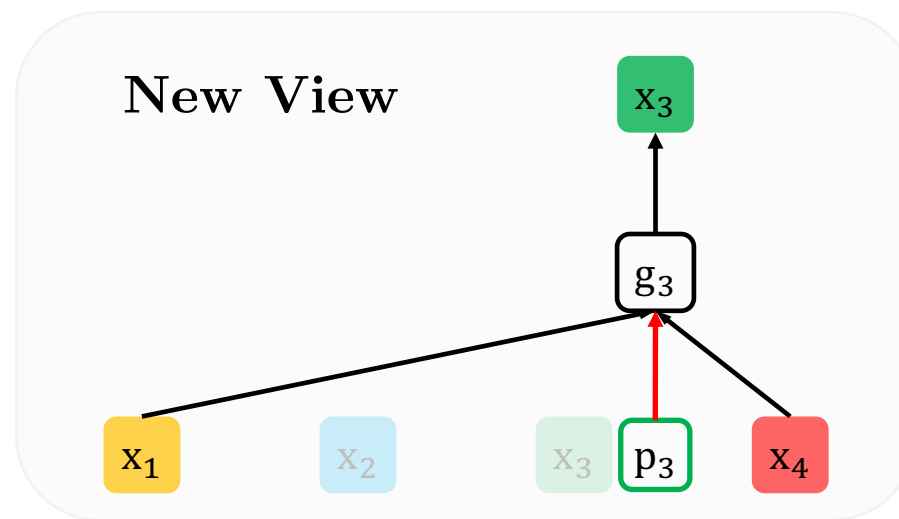
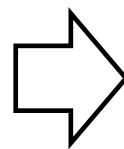
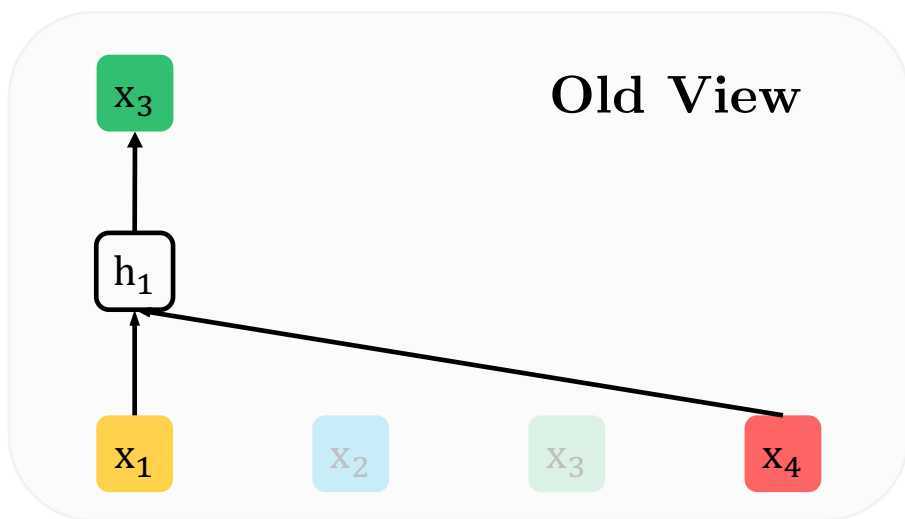
Target Position Aware Representation: $g(\mathbf{z}_t, \mathbf{x}_{\mathbf{z}_{<t}})$

Reuse the Idea of Attention



- Stand at the target position \mathbf{z}_t
- Gather information from $\mathbf{x}_{\mathbf{z}_{<t}}$

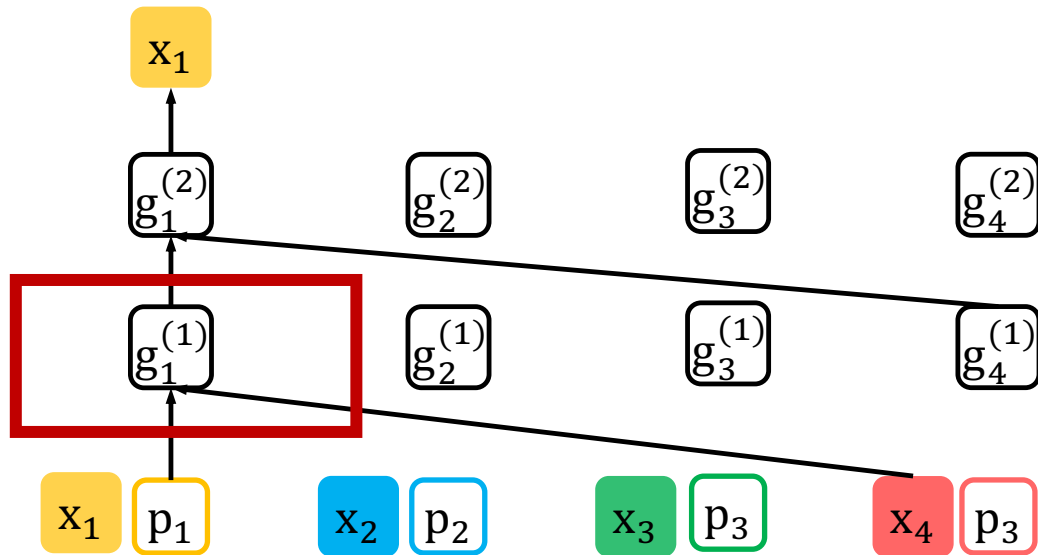
$$g(\mathbf{z}_t, \mathbf{x}_{\mathbf{z}_{<t}}) = \text{Attn}_\theta \left(\underbrace{Q = \text{Enc}(\mathbf{z}_t)}_{\text{Stand at } \mathbf{z}_t}, \underbrace{KV = \mathbf{h}(\mathbf{x}_{\mathbf{z}_{<t}})}_{\text{Gather info. from } \mathbf{x}_{\mathbf{z}_{<t}}} \right)$$



Contradiction: Predicting Self and Others

- Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

Use $g_1^{(1)}$ to predict \mathbf{x}_1 (self)

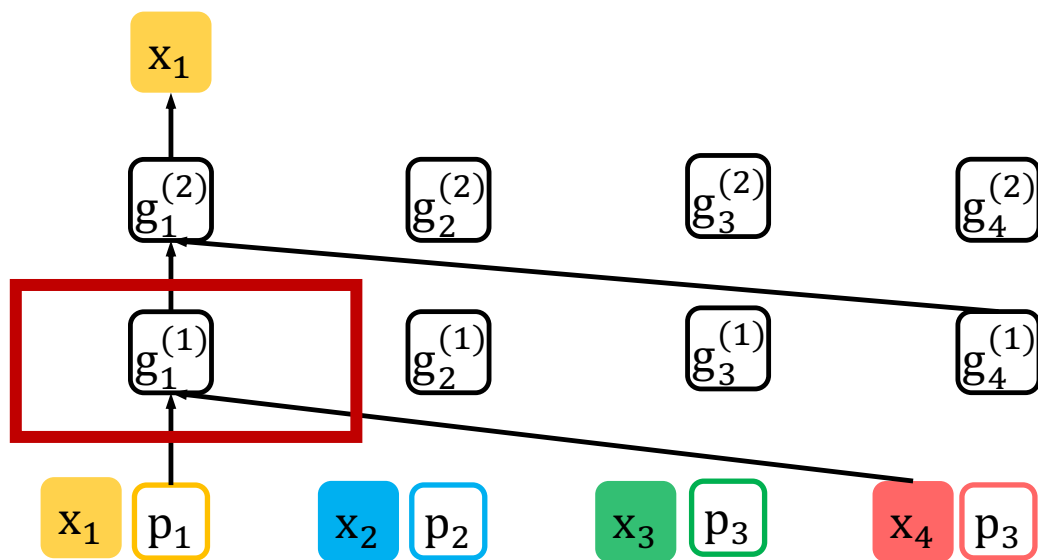


Should not encode \mathbf{x}_1

Contradiction: Predicting Self and Others

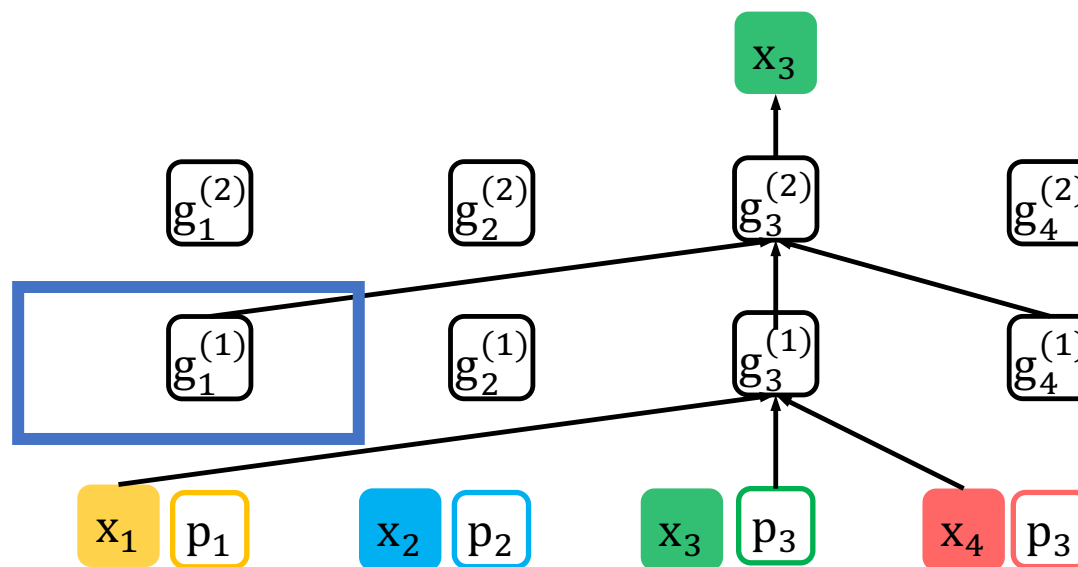
- Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

Use $g_1^{(1)}$ to predict x_1 (self)



Should not encode x_1

Use $g_1^{(1)}$ to predict x_3 (other)

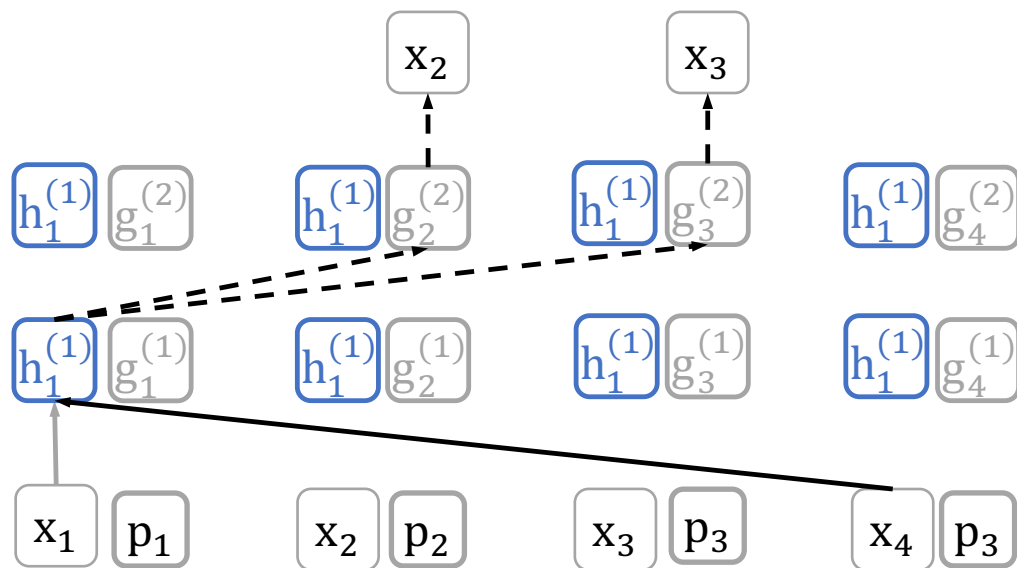


Should encode x_1

Two-Stream Attention

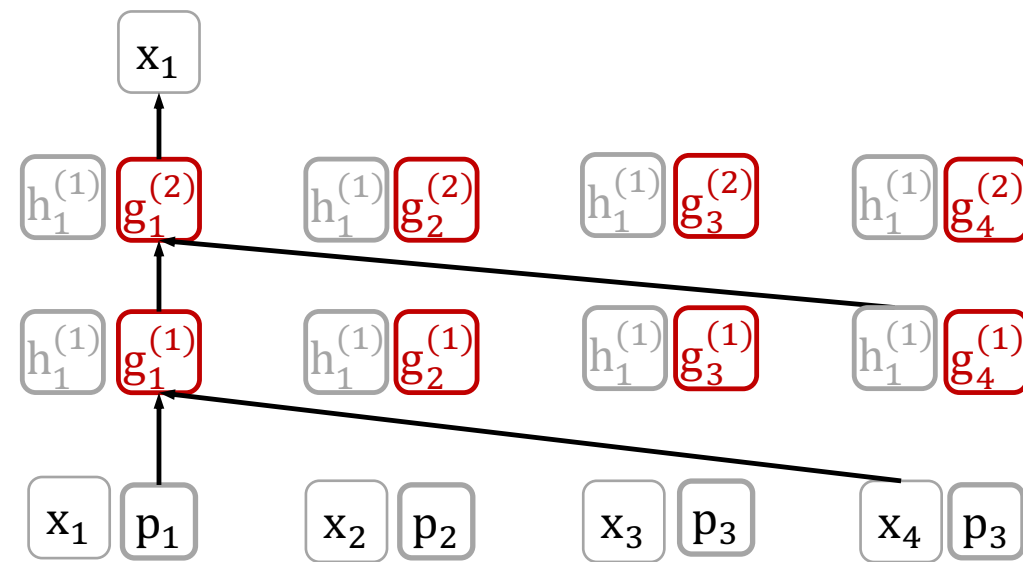
- Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

Encoding. Predicting x_2 and x_3 (others).



h_1 encodes x_1

Decoding. Predicting x_1 (self).



g_1 does not encode x_1

Two-Stream Attention

- Factorization order: $4 \rightarrow 1 \rightarrow 3 \rightarrow 2$

Encoding. P

- Allow Transformers to work with permutation LM.
- No additional parameters.
- No additional computational costs during finetuning.
- Additional costs during pretraining are minimal.

$h_1^{(1)}$ $g_1^{(2)}$

$h_1^{(1)}$ $g_1^{(1)}$

x_1 p_1

x_2 p_2

x_3 p_3

x_4 p_3

x_1 p_1

x_2 p_2

x_3 p_3

x_4 p_3

$g_3^{(2)}$

$h_1^{(1)}$ $g_4^{(2)}$

$g_3^{(1)}$

$h_1^{(1)}$ $g_4^{(1)}$

h_1 encodes x_1

g_1 does not encode x_1

Summarizing XLNet

Challenges



Solutions



Summarizing XLNet

Challenges

Independence assumption and
distribution discrepancy in BERT

Solutions

Summarizing XLNet

Challenges

Independence assumption and
distribution discrepancy in BERT

Solutions

Permutation language modeling



Summarizing XLNet

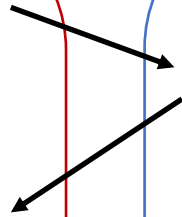
Challenges

Independence assumption and
distribution discrepancy in BERT

Standard parameterization is reduced
to bag-of-words

Solutions

Permutation language modeling



Summarizing XLNet

Challenges

Independence assumption and
distribution discrepancy in BERT

Standard parameterization is reduced
to bag-of-words

Solutions

Permutation language modeling

Reparameterization with positions

Summarizing XLNet

Challenges

Independence assumption and distribution discrepancy in BERT

Standard parameterization is reduced to bag-of-words

Contradiction for predicting both self and others

Solutions

Permutation language modeling

Reparameterization with positions

Summarizing XLNet

Challenges

Independence assumption and distribution discrepancy in BERT

Standard parameterization is reduced to bag-of-words

Contradiction for predicting both self and others

Solutions

Permutation language modeling

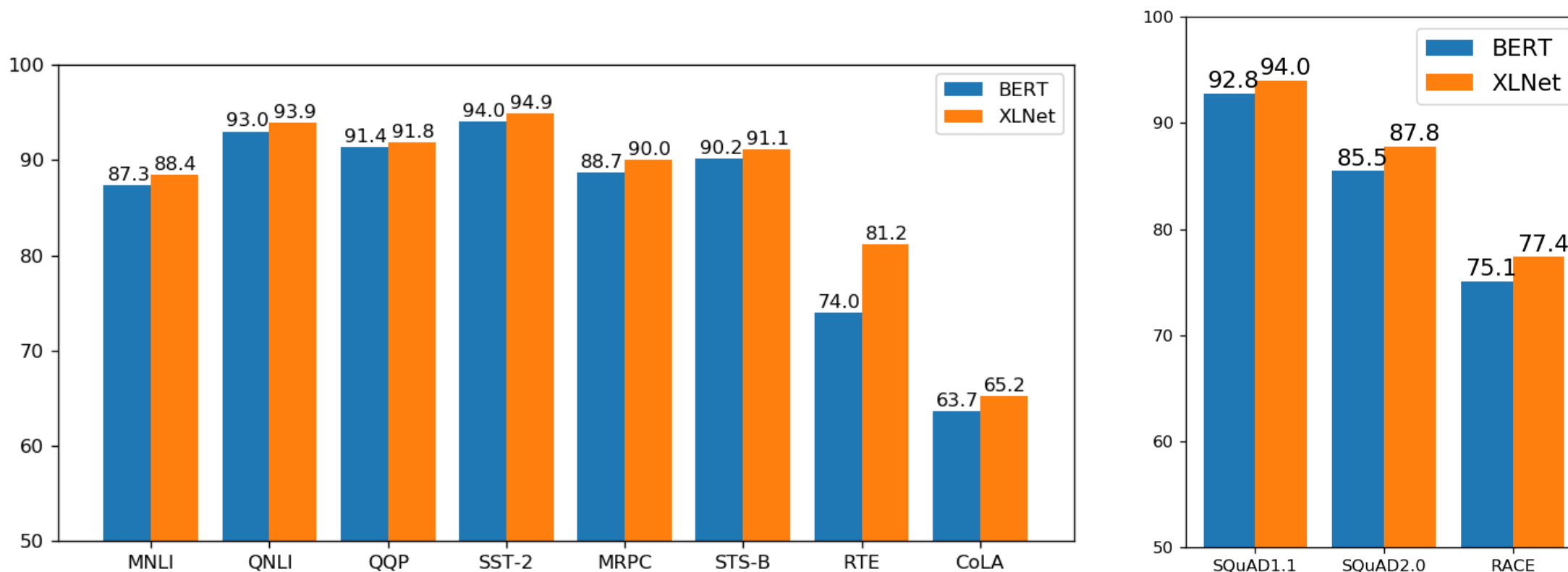
Reparameterization with positions

Two-stream attention

Experiment 1: Comparison with BERT

- **Same** training data **as in BERT**: Wikipedia + BooksCorpus
- **Same** hyperparameters for pretraining **as in BERT**
 - Model size: L=24, H=1024, A=16
 - Batch size: 256
 - Number of steps: 1M
 - ...
- **Same** hyperparameter search space for finetuning **as in BERT**

XLNet outperforms BERT on 20 tasks

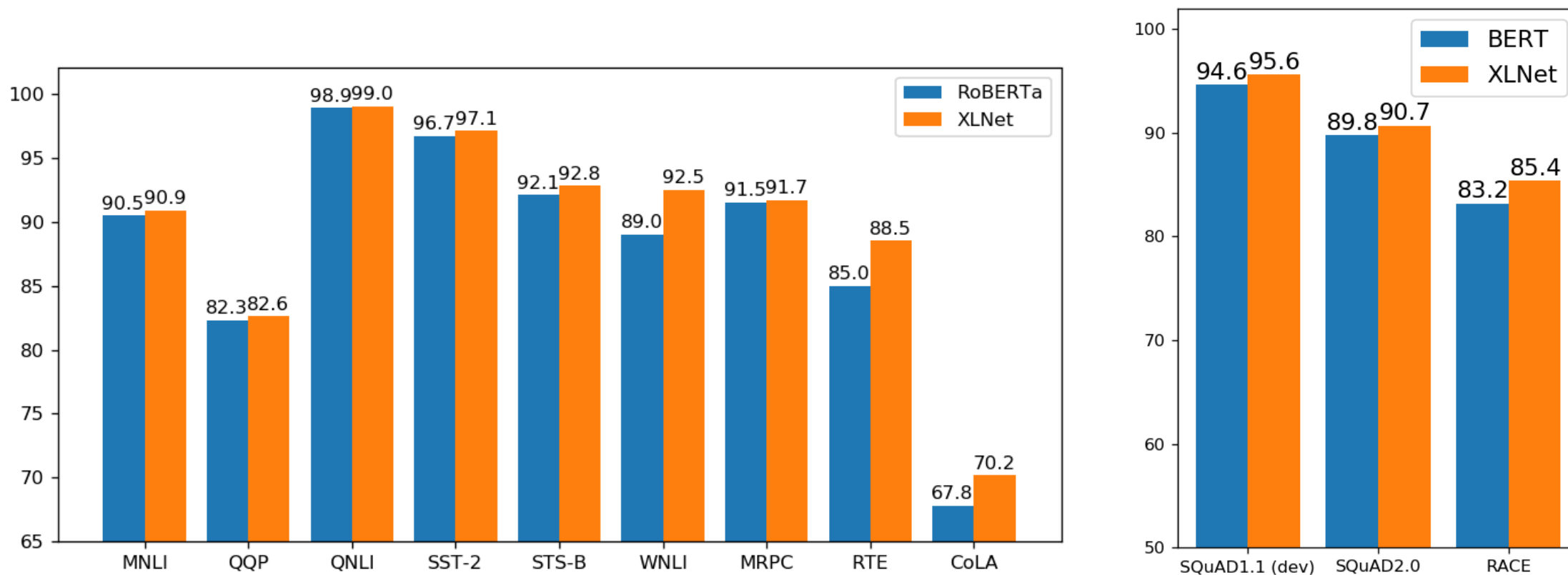


We report the **best of 3** BERT variants.
Almost **identical** training recipes.

Experiment 2: Comparison with RoBERTa

- Less training data for XLNet: 126GB vs 160GB
- **Same** hyperparameters for pretraining **as in RoBERTa**
 - Model size: L=24, H=1024, A=16
 - Batch size: 8192
 - Number of steps: 500K
 - ...
- **Same** hyperparameter search space for finetuning **as in RoBERTa**

XLNet outperforms RoBERTa on all considered tasks

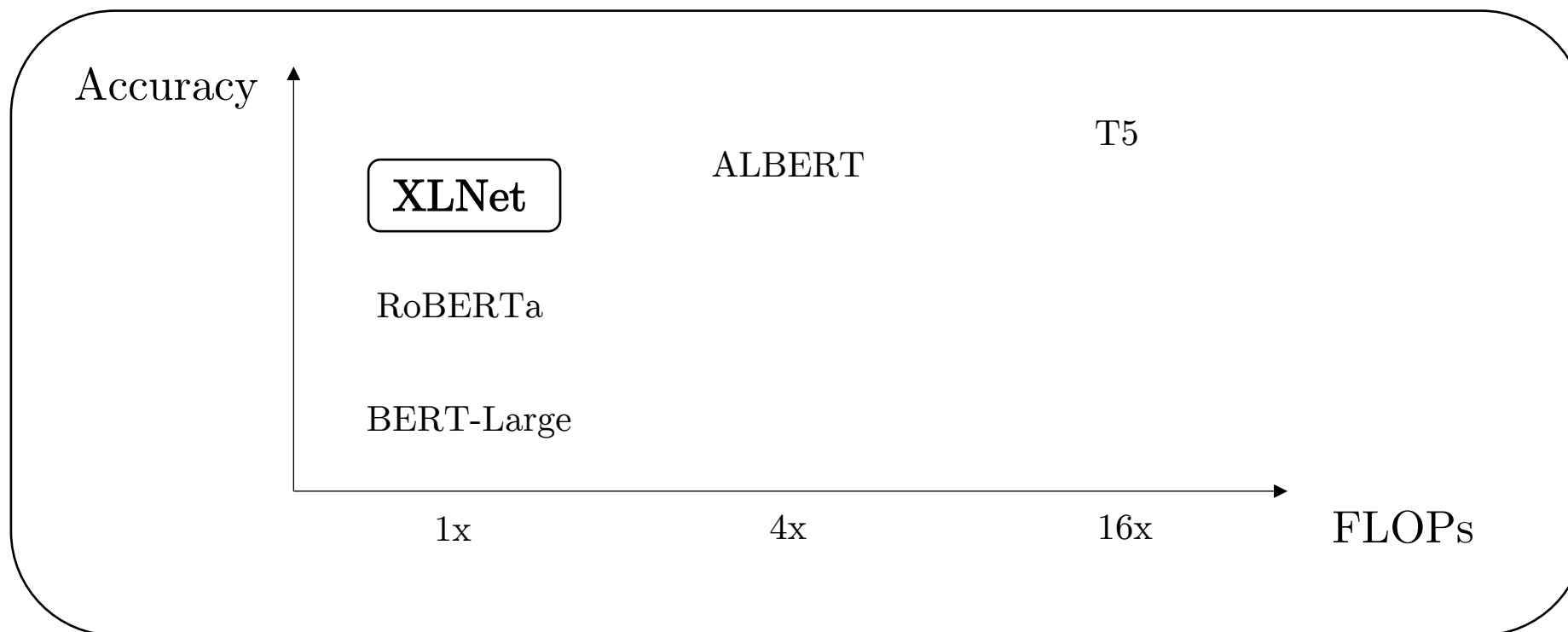


Almost **identical** training recipes.

XLNet is

The **best pretrained model today**

Given standard FLOPs.



XLNet-2 Coming Soon!

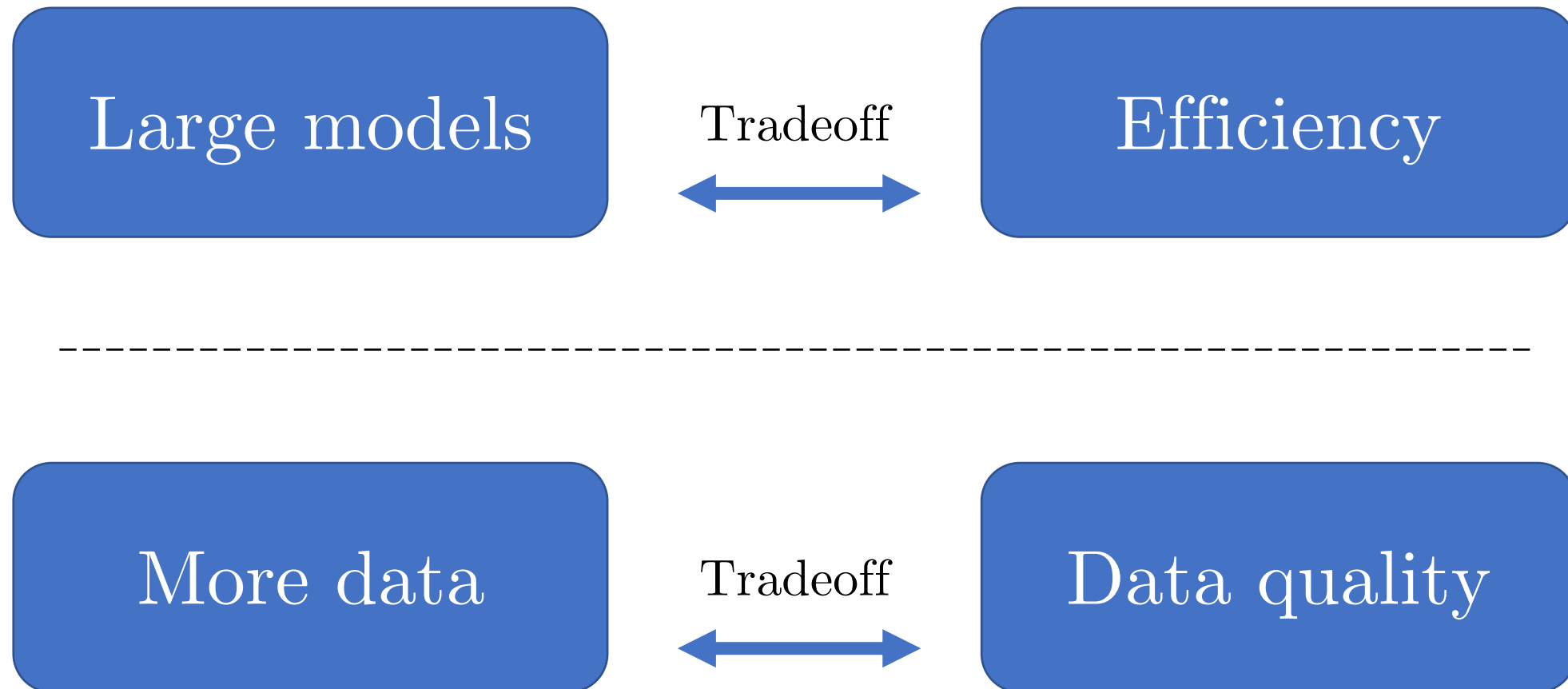
Optimized data
processing

Optimized model
implementation

- Only about 10% slower than BERT during pretraining
- Finetuning speed and memory are identical to BERT
- Outperforms BERT (and larger BERT-like models) consistently under all considered settings

To be release at <https://github.com/zihangdai/xlnet>

Future Work



Thanks!

Zhilin Yang*, Zihang Dai*, Yiming Yang, Jaime Carbonell,
Ruslan Salakhutdinov, Quoc V. Le
(*: equal contribution)