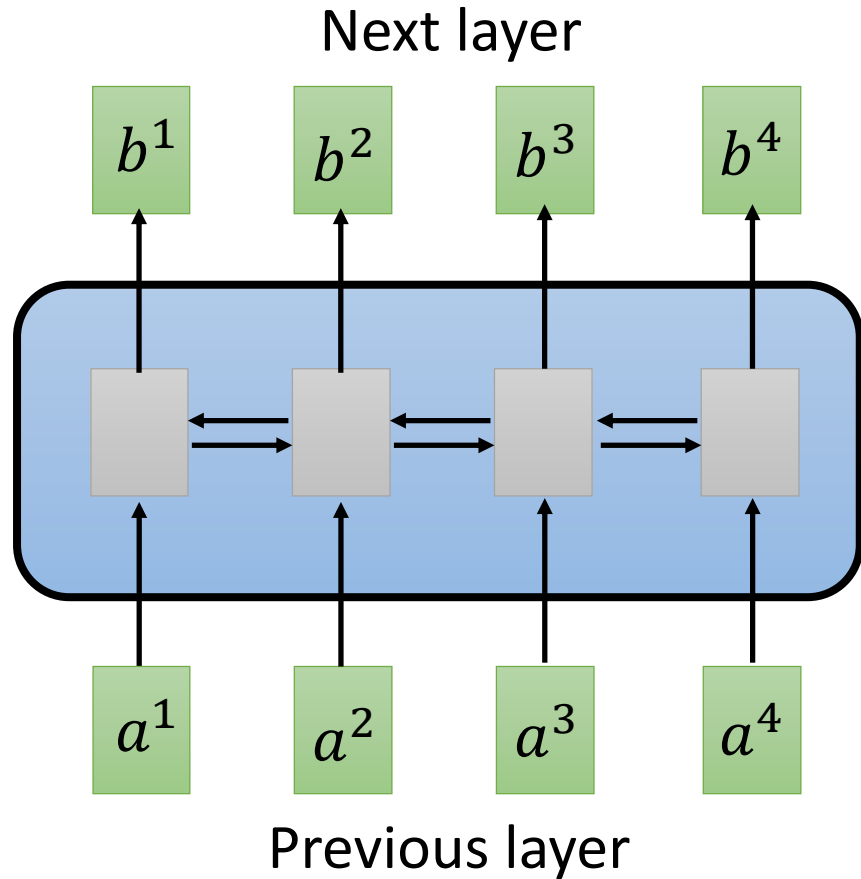


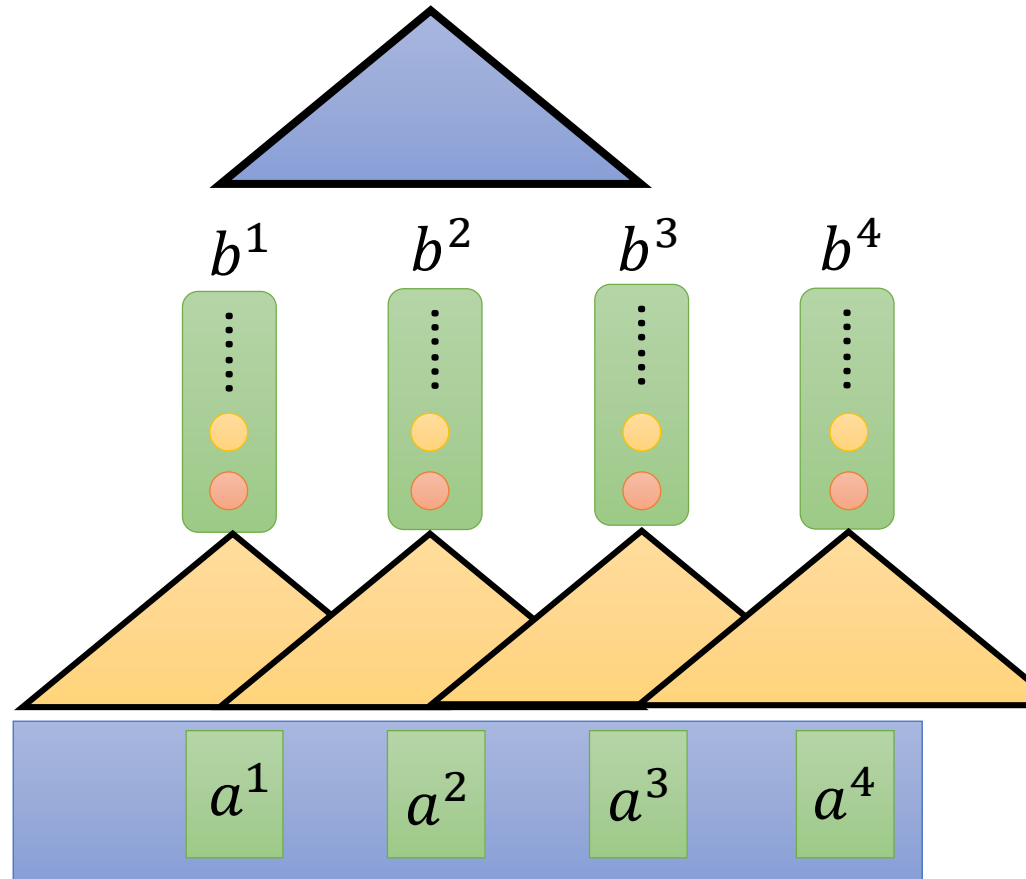
Natural Language Pretraining

Transformer



Hard to parallel

Filters in higher layer can consider longer sequence

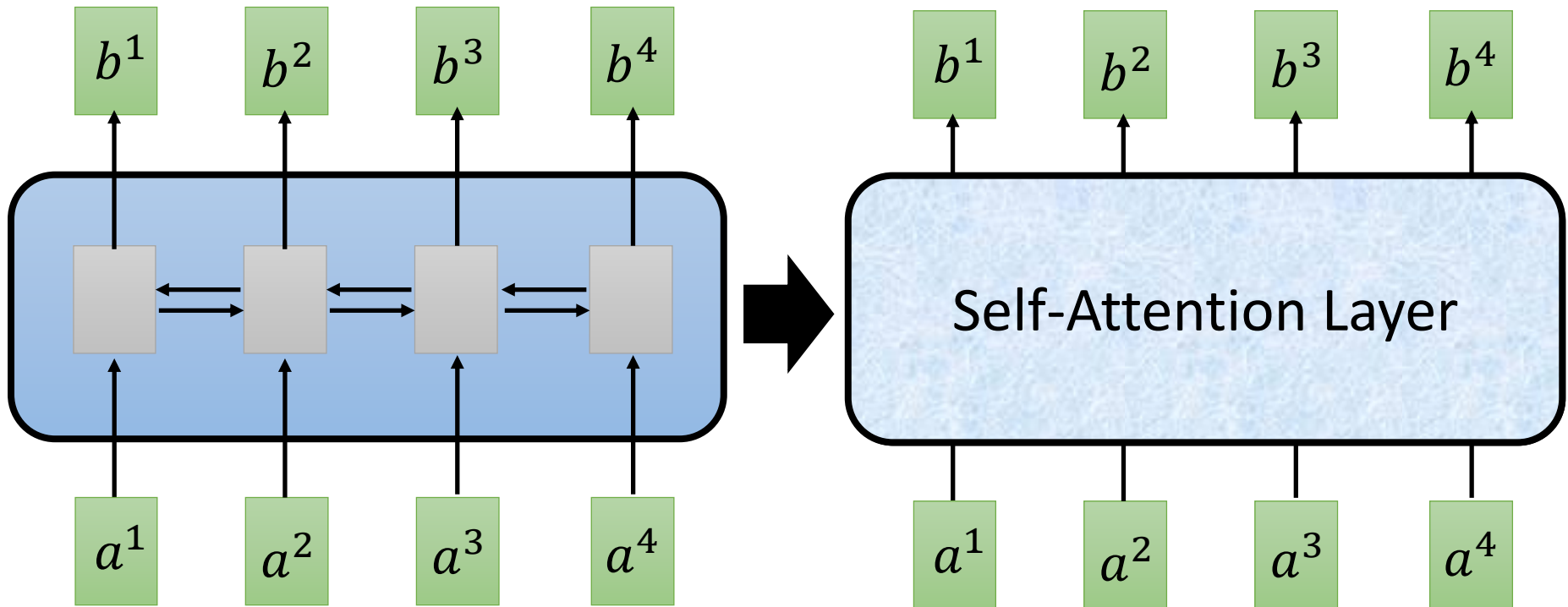


Using CNN to replace RNN
(CNN can parallel)

Self-Attention

b^i is obtained based on the whole input sequence.

b^1, b^2, b^3, b^4 can be parallelly computed.



You can try to replace any thing that has been done by RNN with self-attention.

Self-attention

<https://arxiv.org/abs/1706.03762>



q : query (to match others)

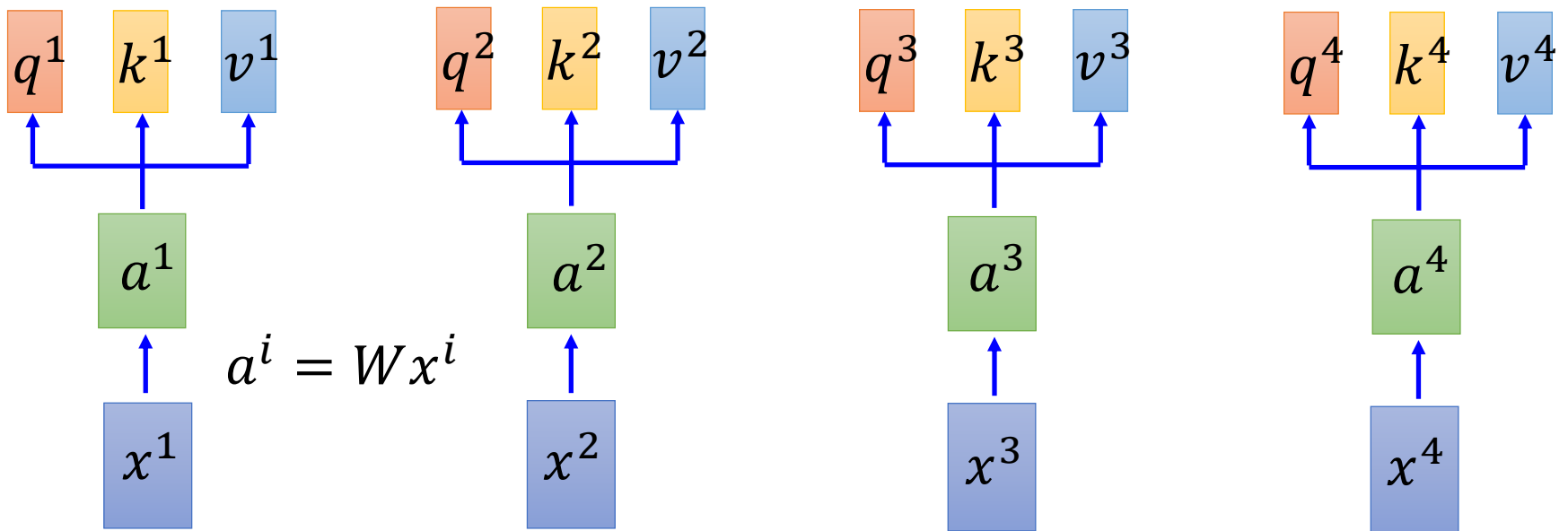
$$q^i = W^q a^i$$

k : key (to be matched)

$$k^i = W^k a^i$$

v : information to be extracted

$$v^i = W^v a^i$$



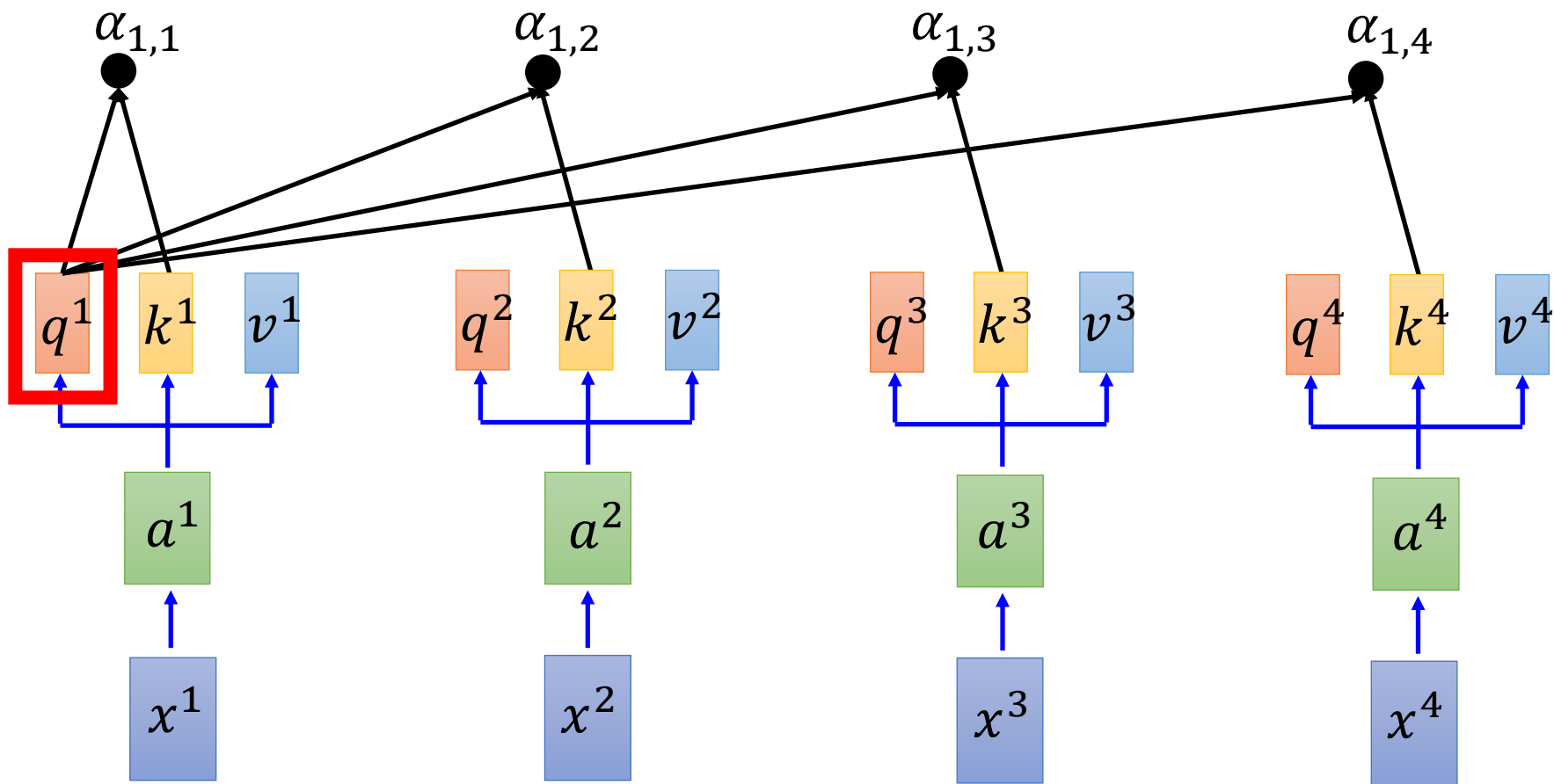
Self-attention

拿每个 query q 去对每个 key k 做 attention

d is the dim of q and k

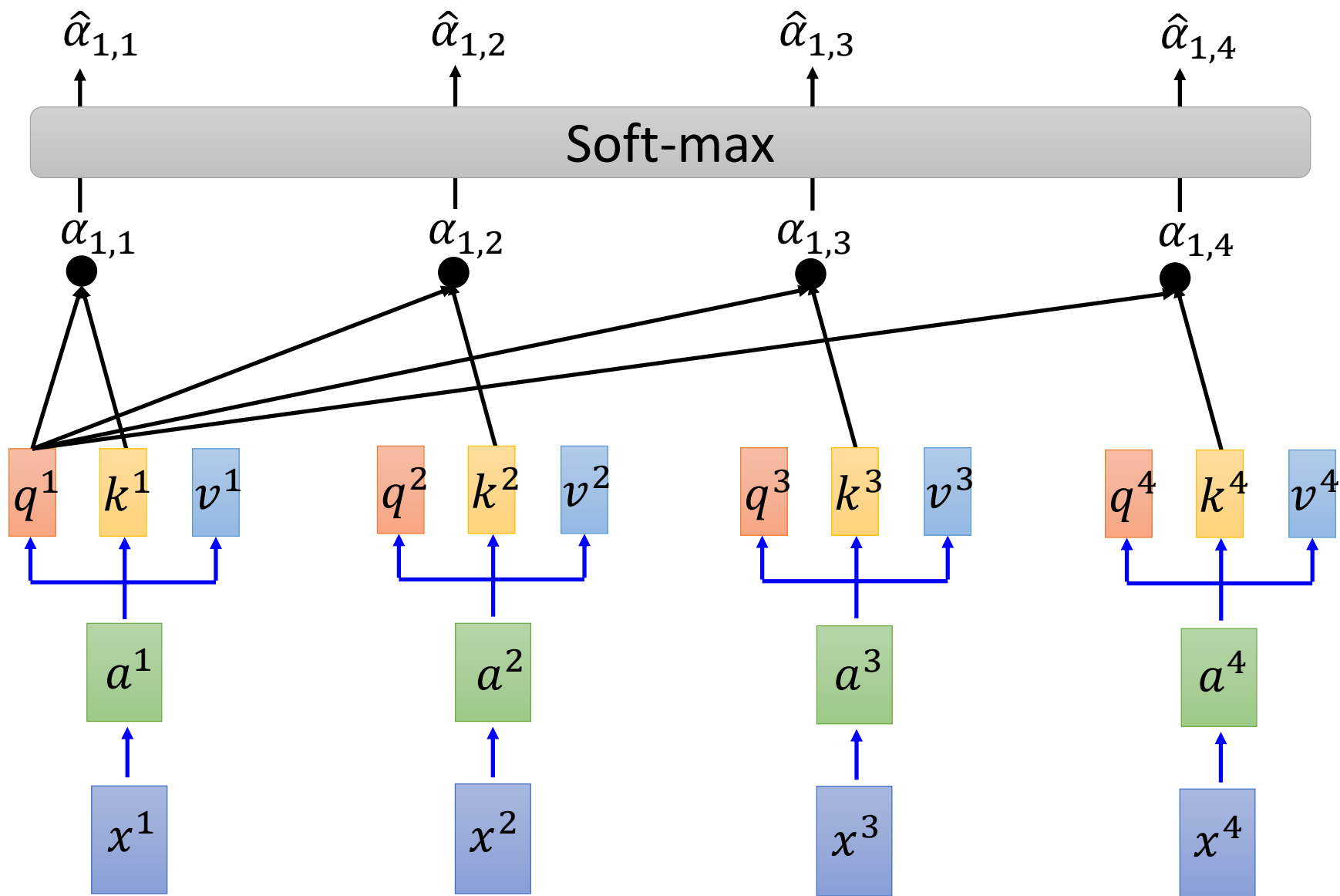
Scaled Dot-Product Attention: $\alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$

dot product



Self-attention

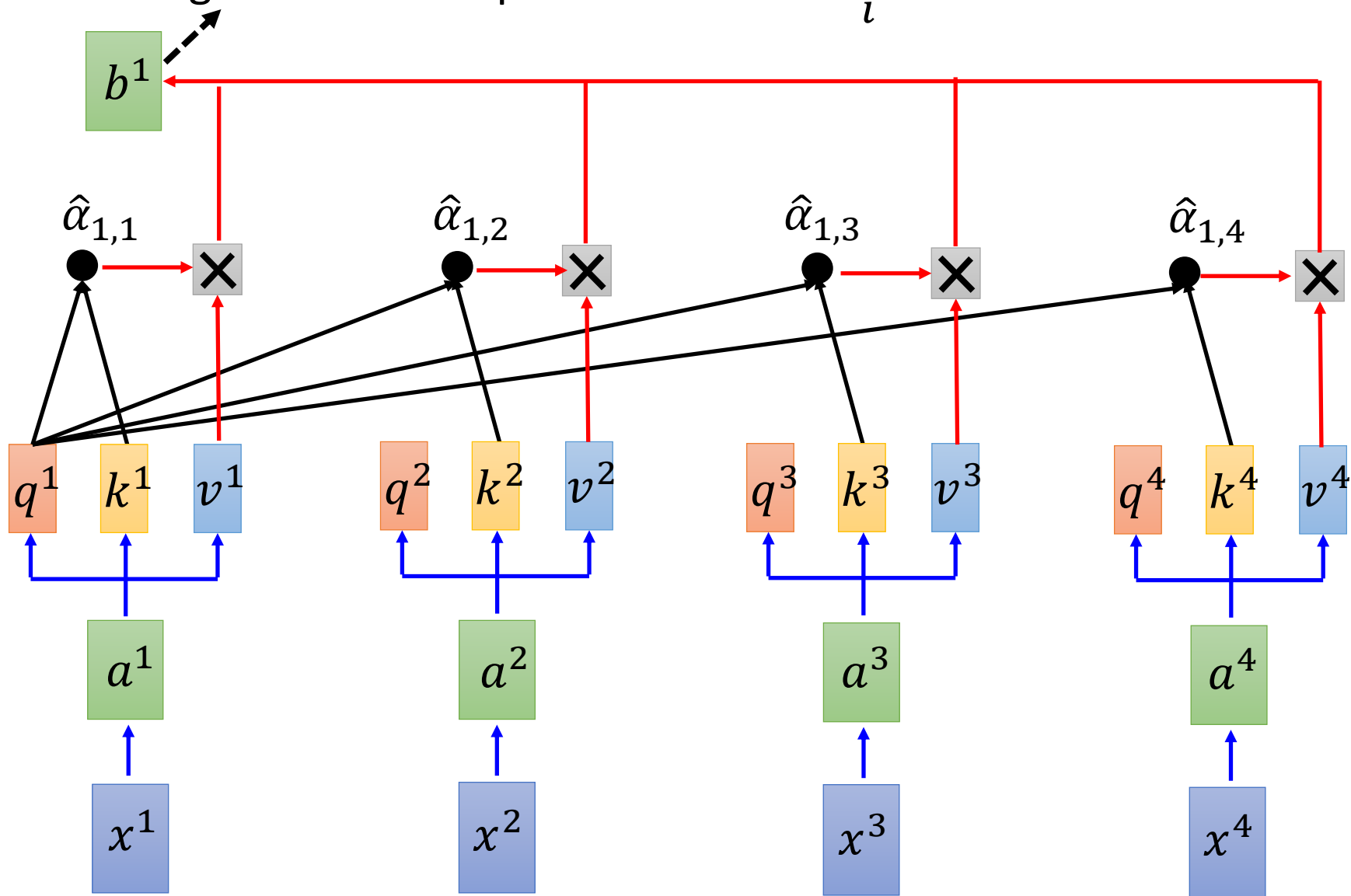
$$\hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



Self-attention

Considering the whole sequence

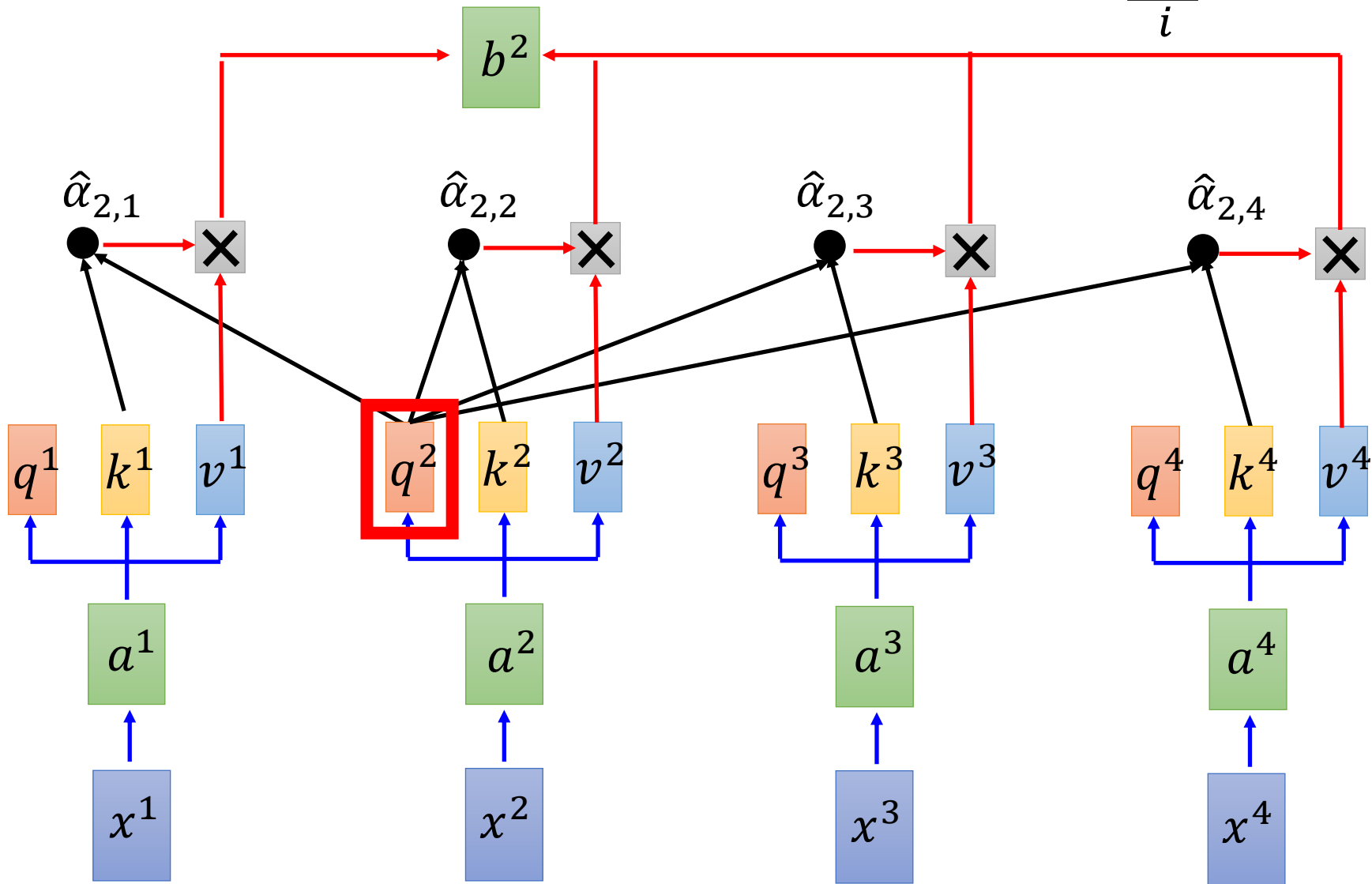
$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$



Self-attention

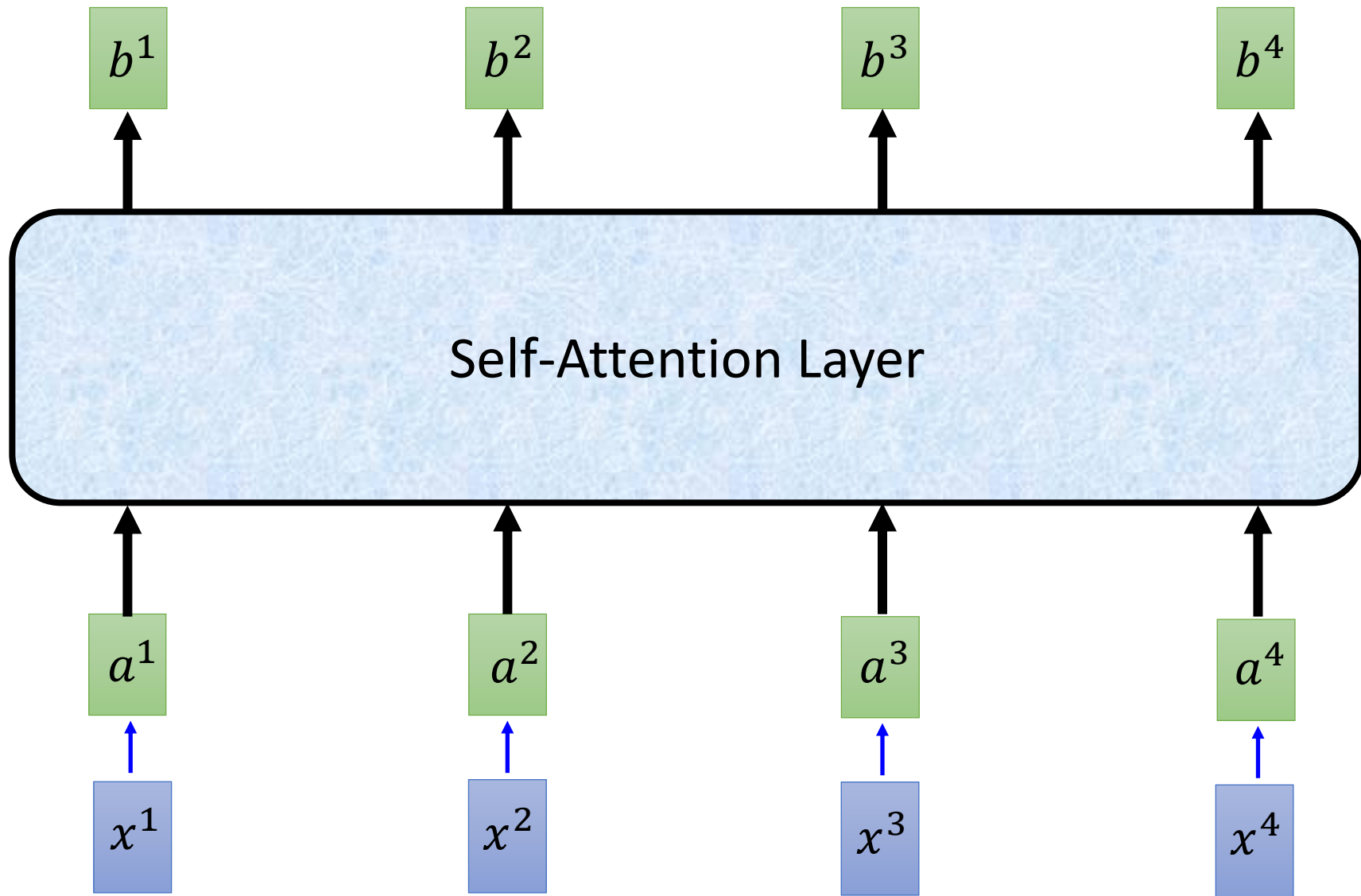
拿每个 query q 去对每个 key k 做 attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



Self-attention

b^1, b^2, b^3, b^4 can be parallelly computed.



Self-attention

$$\begin{matrix} q^1 & q^2 & q^3 & q^4 \\ \hline Q \end{matrix} = \begin{matrix} W^q & \\ \hline a^1 & a^2 & a^3 & a^4 \\ \hline I \end{matrix}$$

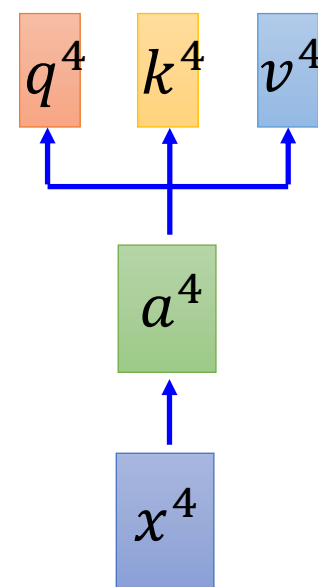
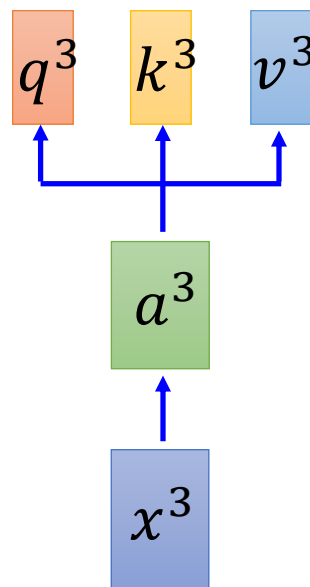
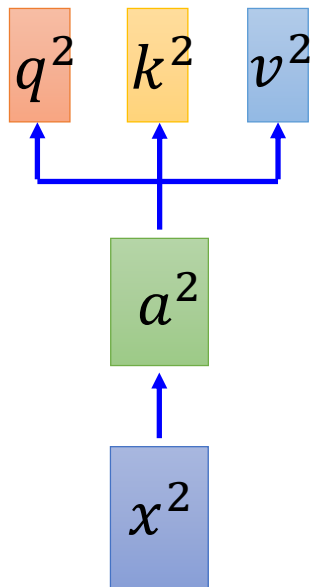
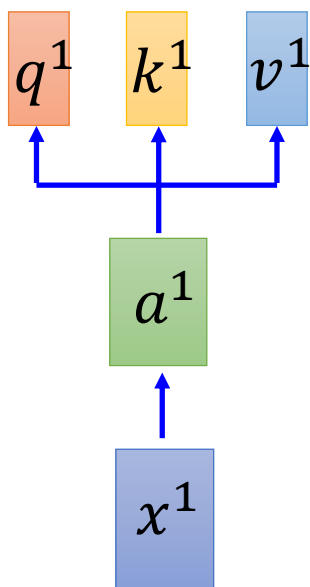
$$\begin{matrix} k^1 & k^2 & k^3 & k^4 \\ \hline K \end{matrix} = \begin{matrix} W^k & \\ \hline a^1 & a^2 & a^3 & a^4 \\ \hline I \end{matrix}$$

$$\begin{matrix} v^1 & v^2 & v^3 & v^4 \\ \hline V \end{matrix} = \begin{matrix} W^v & \\ \hline a^1 & a^2 & a^3 & a^4 \\ \hline I \end{matrix}$$

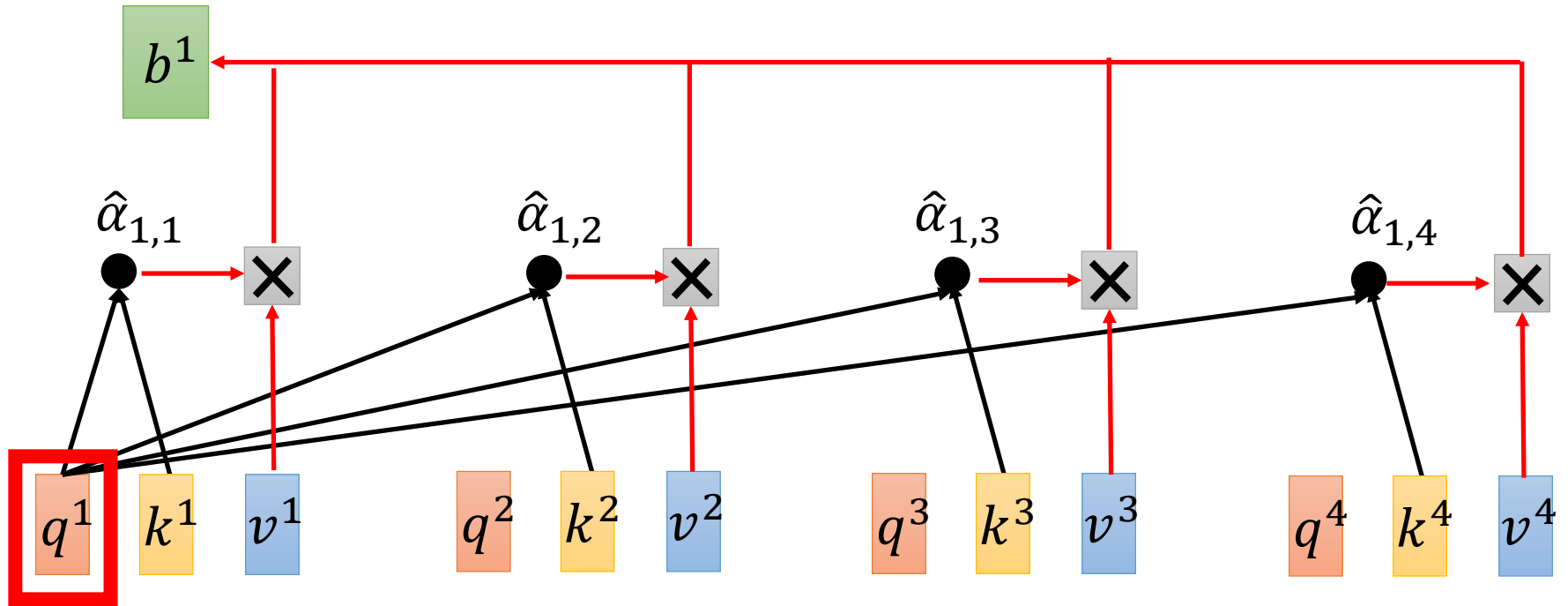
$$q^i = W^q a^i$$

$$k^i = W^k a^i$$

$$v^i = W^v a^i$$



Self-attention



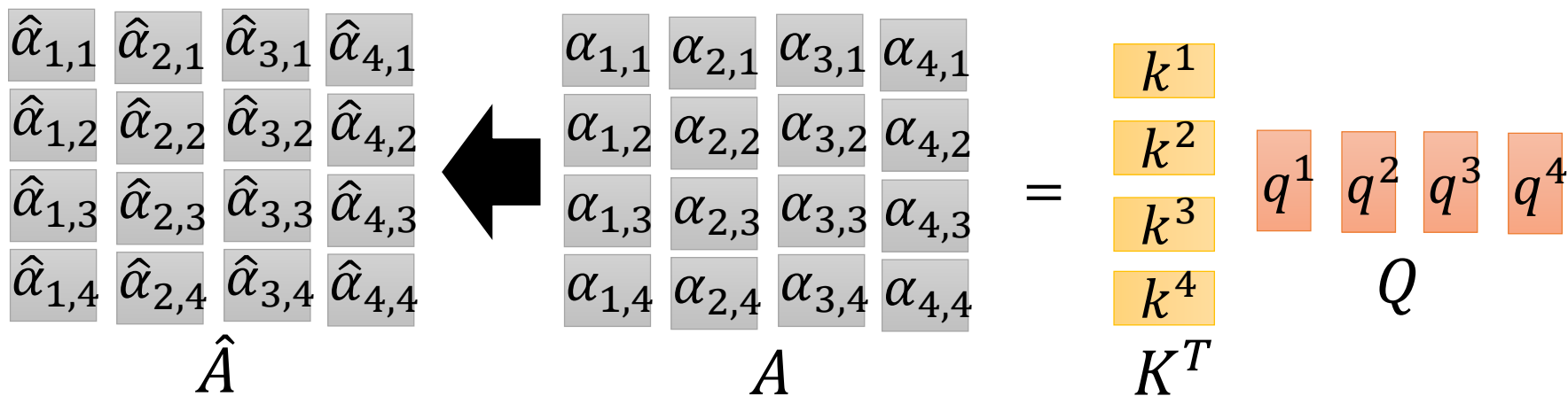
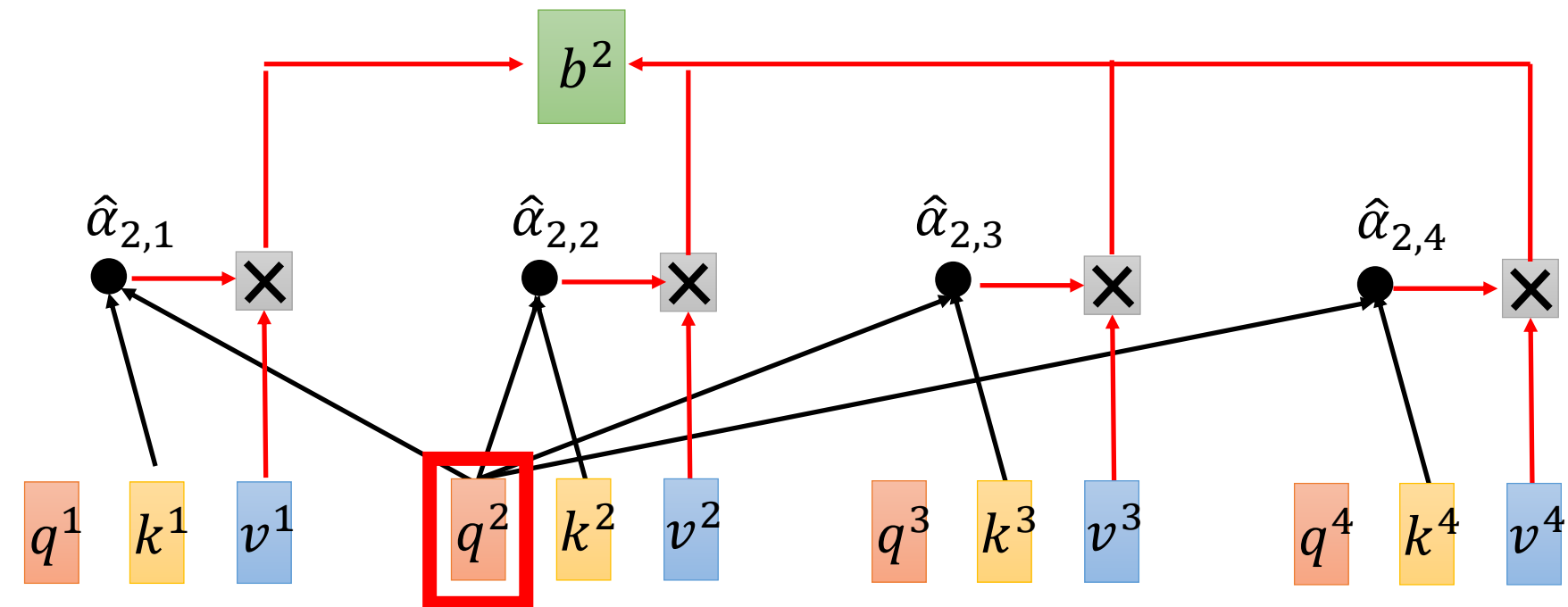
$$\alpha_{1,1} = k^1 q^1 \quad \alpha_{1,2} = k^2 q^1$$
$$\alpha_{1,3} = k^3 q^1 \quad \alpha_{1,4} = k^4 q^1$$

$$\begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} q^1$$

(ignore \sqrt{d} for simplicity)

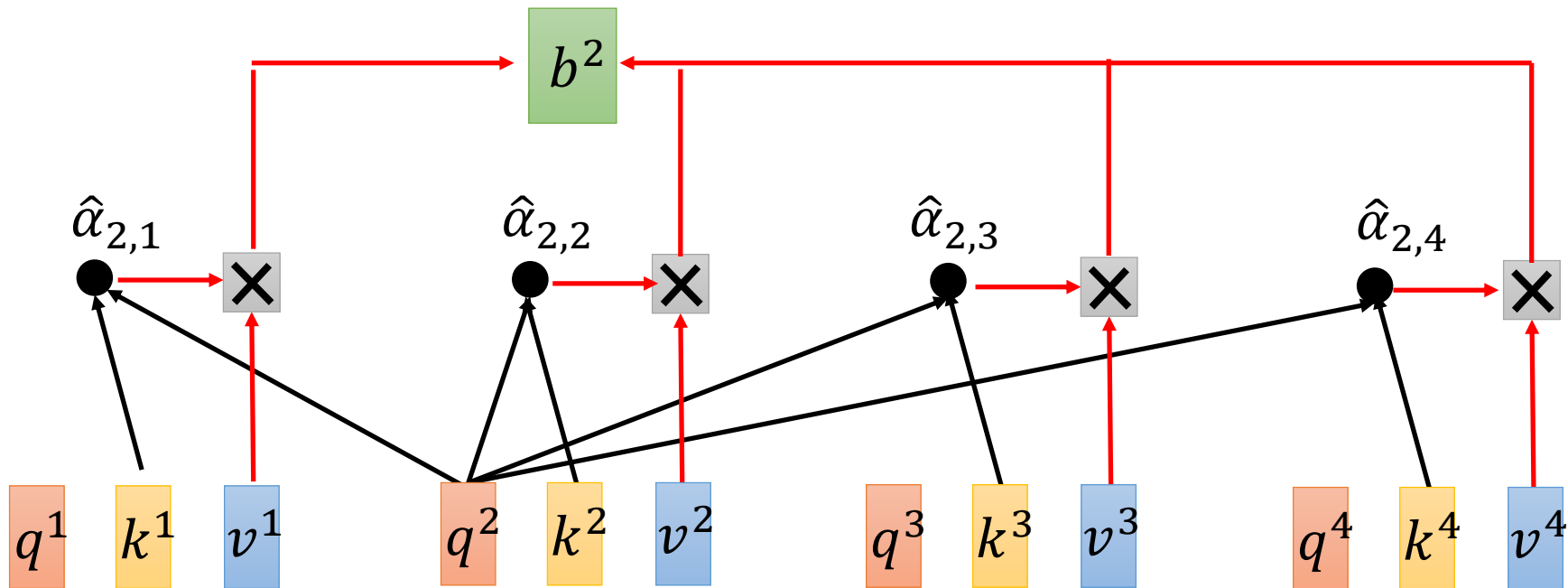
Self-attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



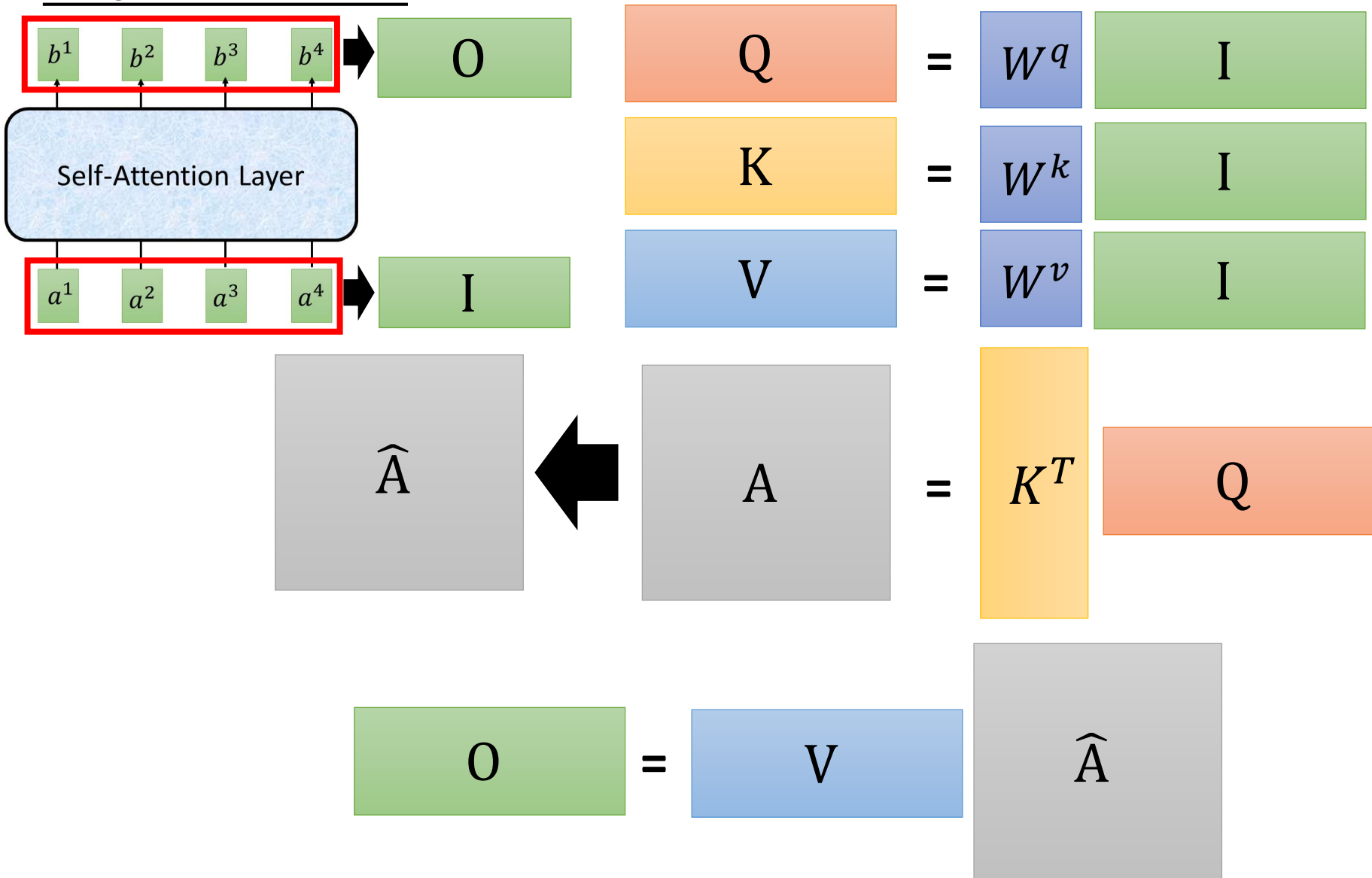
Self-attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



$$\begin{matrix} b^1 & b^2 & b^3 & b^4 \\ \hline \end{matrix} \quad = \quad \begin{matrix} v^1 & v^2 & v^3 & v^4 \\ \hline \end{matrix} \quad \begin{matrix} \hat{\alpha}_{1,1} & \hat{\alpha}_{2,1} & \hat{\alpha}_{3,1} & \hat{\alpha}_{4,1} \\ \hat{\alpha}_{1,2} & \hat{\alpha}_{2,2} & \hat{\alpha}_{3,2} & \hat{\alpha}_{4,2} \\ \hat{\alpha}_{1,3} & \hat{\alpha}_{2,3} & \hat{\alpha}_{3,3} & \hat{\alpha}_{4,3} \\ \hat{\alpha}_{1,4} & \hat{\alpha}_{2,4} & \hat{\alpha}_{3,4} & \hat{\alpha}_{4,4} \\ \hline \hat{A} \end{matrix}$$

Self-attention



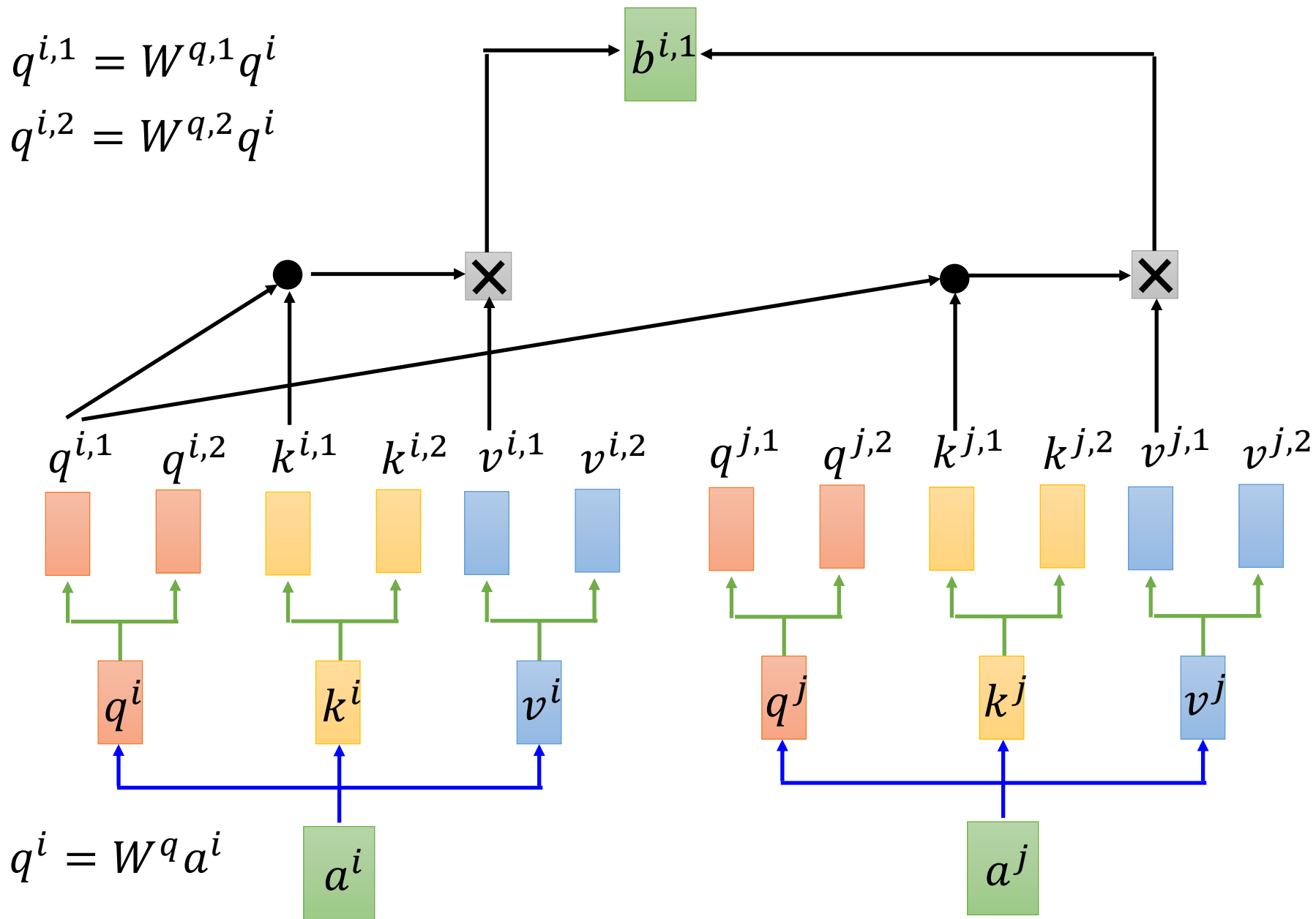
反正就是一堆矩阵乘法，用 GPU 可以加速

Multi-head Self-attention

(2 heads as example)

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$

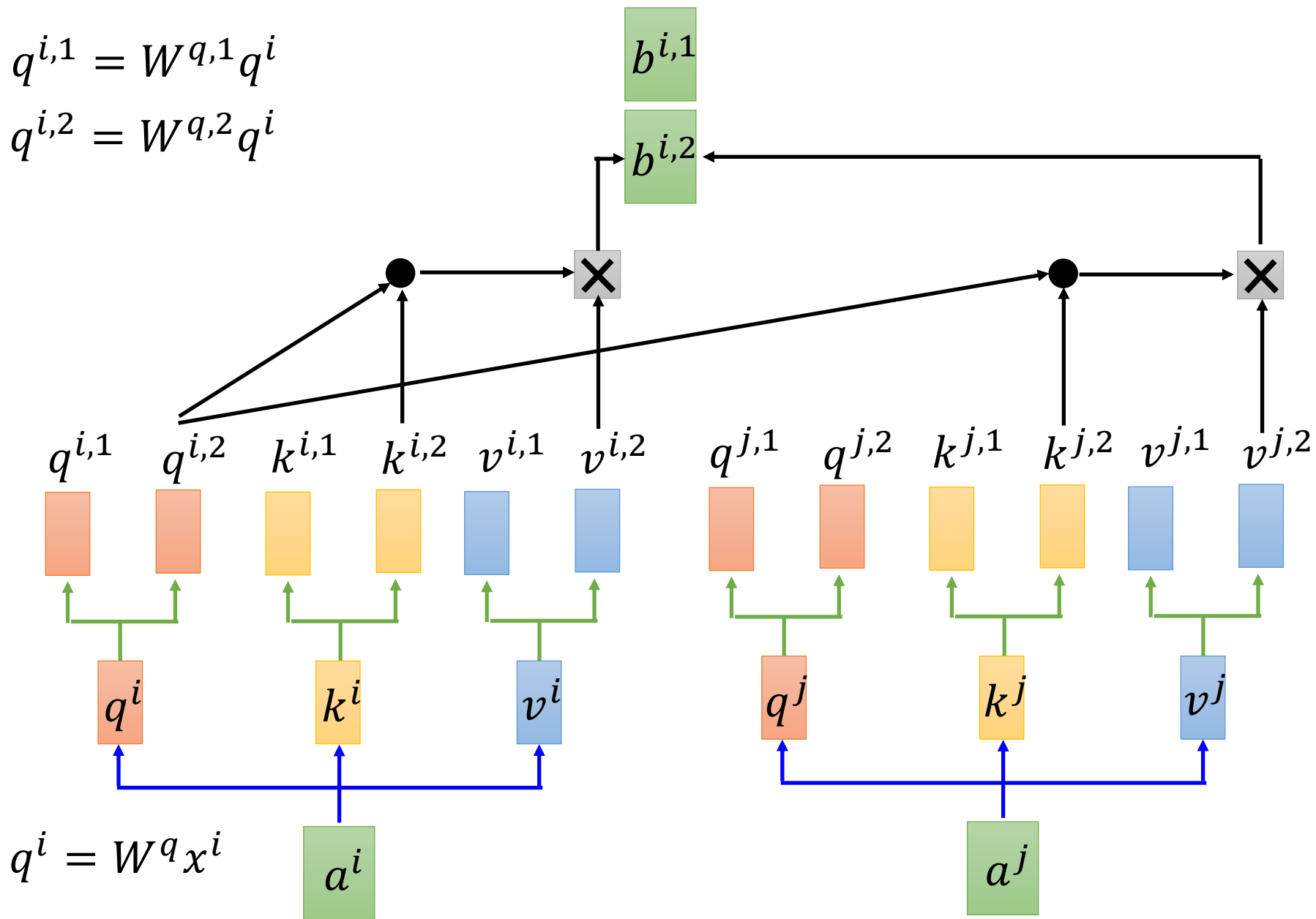


Multi-head Self-attention

(2 heads as example)

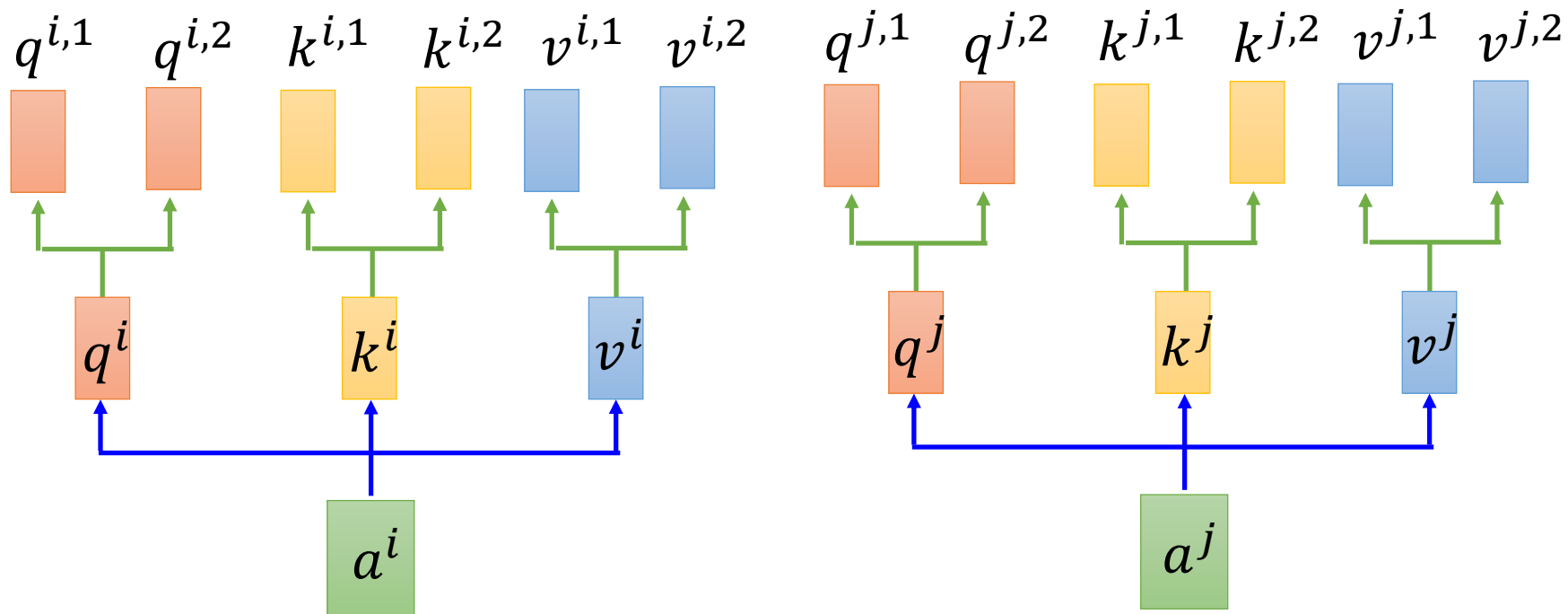
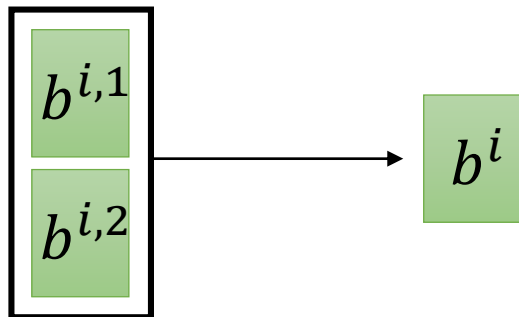
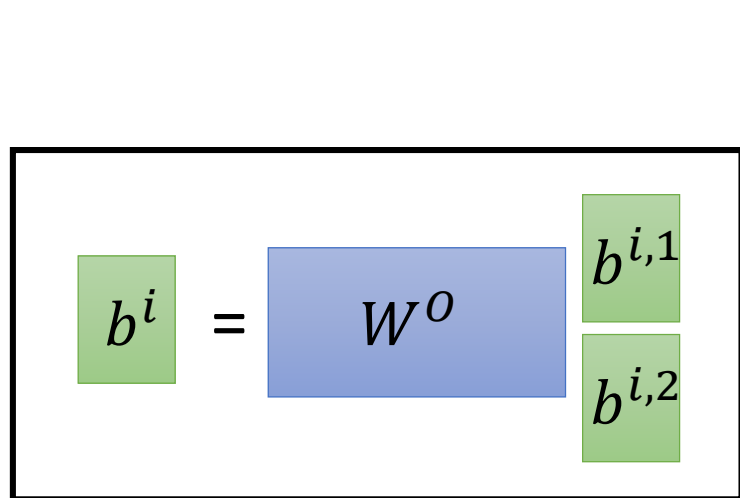
$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$



Multi-head Self-attention

(2 heads as example)



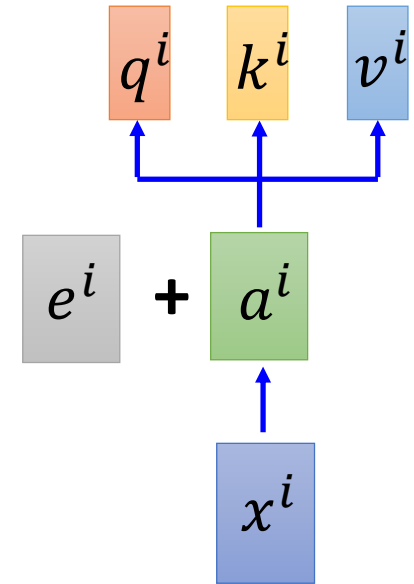
Positional Encoding

- No position information in self-attention.
- Original paper: each position has a unique positional vector e^i (not learned from data)

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

$$\omega_k = \frac{1}{10000^{2k/d}}$$

t is the desired position in an input sentence



$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

0:	0	0	0	0	8:	1	0	0	0
1:	0	0	0	1	9:	1	0	0	1
2:	0	0	1	0	10:	1	0	1	0
3:	0	0	1	1	11:	1	0	1	1
4:	0	1	0	0	12:	1	1	0	0
5:	0	1	0	1	13:	1	1	0	1
6:	0	1	1	0	14:	1	1	1	0
7:	0	1	1	1	15:	1	1	1	1

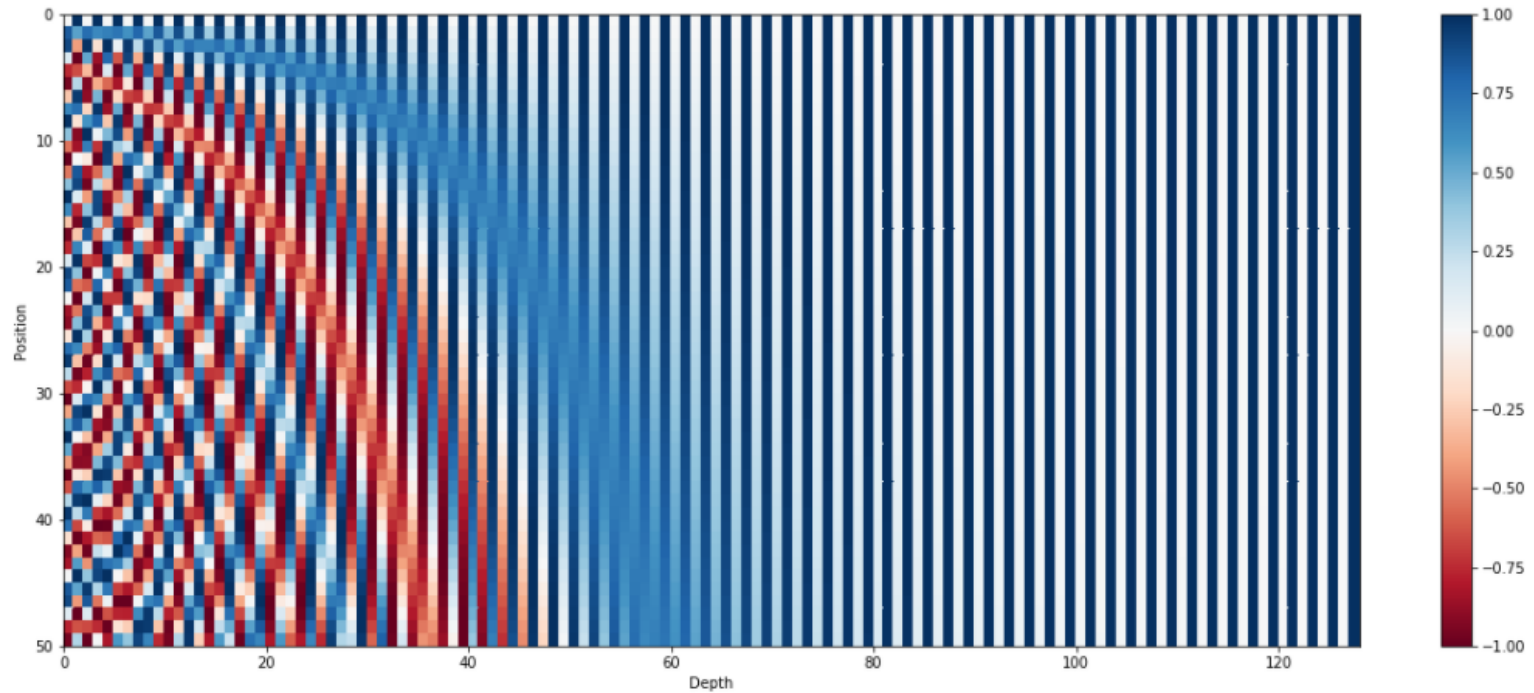
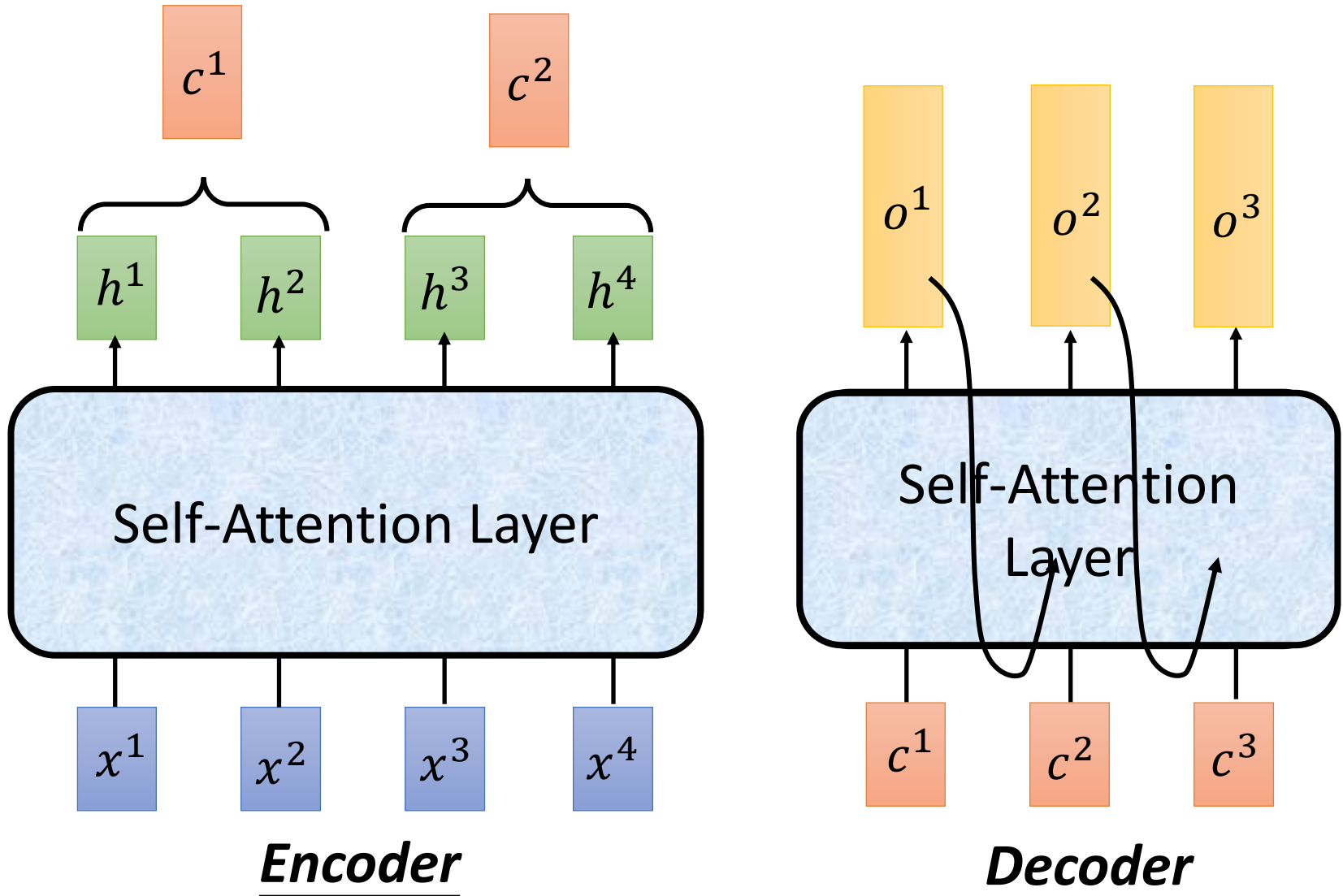


Figure 2 - The 128-dimensional positional encoding for a sentence with the maximum length of 50.

Each row represents the embedding vector \vec{p}_t

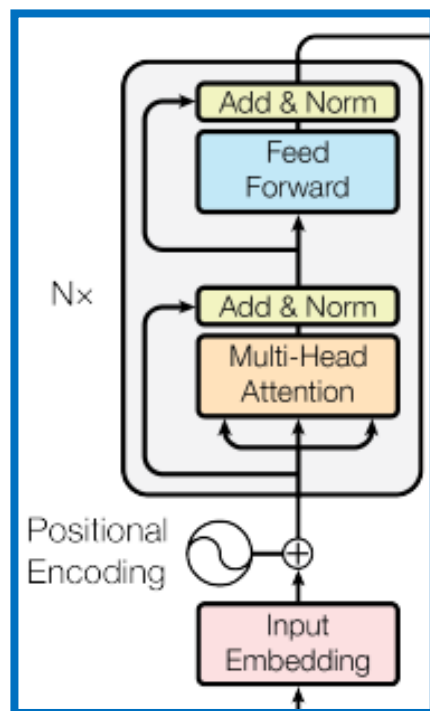
Seq2seq with Attention



Transformer

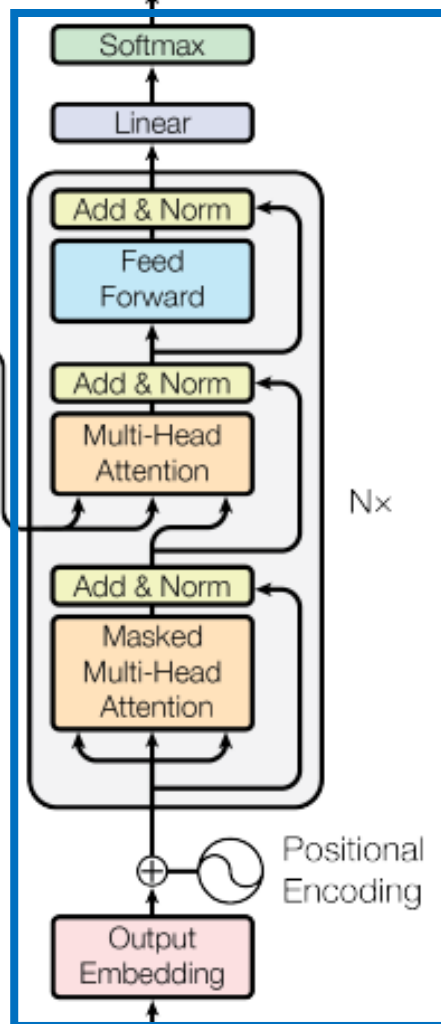
Using Chinese to English translation as example

Encoder



机器学习

machine learning
Output Probabilities



Decoder

Outputs
(shifted right)

<BOS>

machine






Outline

1. A brief note on subword modeling
2. Motivating model pretraining from word embeddings
3. Model pretraining three ways
 1. Decoders
 2. Encoders
 3. Encoder–Decoders
4. Interlude: what do we think pretraining is teaching?
5. Very large models and in-context learning

Word structure and subword models

Let's take a look at the assumptions we've made about a language's vocabulary.

We assume a fixed vocab of tens of thousands of words, built from the training set. All *novel* words seen at test time are mapped to a single UNK.

	word	vocab	embedding
Common words	hat	mapping	
	lear	pizza	
Variations	n	(index)	
	taaaaas	tasty	
misspellings novel items	transformerify	(index) UNK	
	laern	(index) UNK	
		(index) UNK	
		(index)	

Word structure and subword models

Finite vocabulary assumptions make even *less* sense in many languages.

- Many languages exhibit complex **morphology**, or word structure.
 - The effect is more word types, each occurring fewer times.

Example: Swahili verbs can have hundreds of conjugations, each encoding a wide variety of information. (Tense, mood, definiteness, negation, information about the object, ++)

Conjugation of <i>-ambia</i>																									
Form		Non-finite forms												Negative											
Infinitive		Positive												kutoambia											
Imperative		Simple finite forms												Plural											
Habitual		Singular												ambieni											
		huambia																							
		Complex finite forms																							
Polarity	Persons						Classes				Classes														
	Sg.	1st	Pl.	Sg.	2nd	Pl.	Sg./1	Pl./2	3	M-mi	4	5	Ma	6	7	KI-vi	8	9	N	10	U	Ku	Pa	Mu	
Positive	niliambia	tulambia	ulambia	mlambia	aliambia	valiambia	ulambia	ilambia	lilambia	yaliambia	kilambia	vilambia	ilambia	ziliambia	ulambia	kulambia	paliambia	mulambia							
Negative	sikuambia	hatuambia	hukuambia	hamkuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia	hakuambia
Positive	ninaambia	tunaambia	unaambia	mnaambia	anaambia	wanaambia	unaambia	inaambia	inaambia	yanaambia	kinaambia	vinaambia	inaambia	zinaambia	unaambia	kunaambia	panaambia	munaambia							
Negative	sinaambia	hatuambia	huambia	hamambia	hambia	havaambia	hauambia	halambia	halambia	hayaambia	hakiambia	haviambia	halambia	hazambia	hauambia	hakuambia	hapambia	hamuambia							
Positive	nitaambia	tutaambia	utaambia	mtaambia	ataambia	vataambia	utaambia	itaambia	itaambia	yataambia	kitaambia	vitaambia	itaambia	zitaambia	utaambia	kutaambia	pataambia	mutaambia							
Negative	sitaambia	hatuambia	hutaambia	hamtaambia	hataambia	havaambia	hutaambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia	hataambia
Positive	niambia	tumbia	umbia	mambia	aambia	waambia	umbia	lambia	lambia	yaambia	kiambia	viambia	lambia	ziambia	umbia	kuambia	paambia	muambia							
Negative	nisambia	tusambia	usambia	msambia	asambia	wasambia	usambia	lisambia	lisambia	yasambia	kisambia	visambia	lisambia	zisambia	usambia	kusambia	pasambia	musambia							
Positive	ningeambia	tungeambia	ungeambia	mngeambia	angeambia	wangeambia	ungeambia	ingambia	ingambia	yangeambia	kingambia	vingambia	ingambia	zingambia	ungeambia	kungeambia	pangeambia	mungeambia							
Negative	nisingeambia	hatingeambia	usingeambia	msingeambia	asingeambia	wasingeambia	usingeambia	isingeambia	isingeambia	hayingeambia	kingeambia	vingeambia	isingeambia	zingeambia	usingeambia	kusingeambia	pingeambia	musingeambia							
Positive	ningaliambia	tungaliambia	ungaliambia	mngaliambia	angaliambia	wangaliambia	ungaliambia	ingaliambia	ingaliambia	yangaliambia	kingaliambia	vingaliambia	ingaliambia	zingaliambia	ungaliambia	kungaliambia	pangaliambia	mungaliambia							
Negative	nisingaliambia	hatingaliambia	usingaliambia	msingaliambia	asingaliambia	wasingaliambia	usingaliambia	isingaliambia	isingaliambia	hayingaliambia	kingaliambia	vingaliambia	isingaliambia	zingaliambia	usingaliambia	kusingaliambia	pingaliambia	musingaliambia							
Positive	ningeliambia	tungeliambia	ungeliambia	mngeliambia	angeliambia	wangeliambia	ungeliambia	ingeliambia	ingeliambia	yangeliambia	kingeliambia	vingeliambia	ingeliambia	zingeliambia	ungeliambia	kungeliambia	pangeliambia	mungeliambia							
Negative	nisingeliambia	hatingeliambia	usingeliambia	msingeliambia	asingeliambia	wasingeliambia	usingeliambia	isingeliambia	isingeliambia	hayingeliambia	kingeliambia	vingeliambia	isingeliambia	zingeliambia	usingeliambia	kusingeliambia	pingeliambia	musingeliambia							
Positive	ningaliambia	tungaliambia	ungaliambia	mngaliambia	angaliambia	wangaliambia	ungaliambia	ingaliambia	ingaliambia	yangaliambia	kingaliambia	vingaliambia	ingaliambia	zingaliambia	ungaliambia	kungaliambia	pangaliambia	mungaliambia							
Negative	nisingaliambia	hatingaliambia	usingaliambia	msingaliambia	asingaliambia	wasingaliambia	usingaliambia	isingaliambia	isingaliambia	hayingaliambia	kingaliambia	vingaliambia	isingaliambia	zingaliambia	usingaliambia	kusingaliambia	pingaliambia	musingaliambia							
Positive	naambia	tvaambia	waambia	mvaambia	aambia	waambia	yaambia	yaambia	yaambia	yaambia	chaambia	yaambia	yaambia	zaambia	vaambia	kvaambia	paambia	meaambia							

[Wiktionary]

Here's a small fraction of the conjugations for *ambia* - to tell.

The byte-pair encoding algorithm

Subword modeling in NLP encompasses a wide range of methods for reasoning about structure below the word level. (Parts of words, characters, bytes.)

- The dominant modern paradigm is to learn a vocabulary of **parts of words (subword tokens)**.
- At training and testing time, each word is split into a sequence of known subwords.

Byte-Pair Encoding (BPE) is a simple, effective strategy for defining a subword vocabulary.

1. Start with a vocabulary containing only characters and an “end-of-word” symbol.
2. Using a corpus of text, find the most common adjacent characters “a,b” ; add “ab” as a subword.
3. Replace instances of the character pair with the new subword; repeat until desired vocab size.






[[Senrich et al., 2016](#), [Wu et al., 2016](#)]

Originally used in NLP for machine translation: now a similar method

Word structure and subword models

Common words end up being a part of the subword vocabulary, while rarer words are split into (sometimes intuitive, sometimes not) components.

In the worst case, words are split into as many subwords as they have characters.

	word	vocab	embedding
Common words	hat	mapping hat	
	lear	learn	
Variations	n	taa## aaa##	
	taaaaas	sty la##	
misspellings	Transformerify	ern##	
	laern	Transformer##	
novel items		ify	

Outline

1. A brief note on subword modeling
2. Motivating model pretraining from word embeddings
3. Model pretraining three ways
 1. Decoders
 2. Encoders
 3. Encoder–Decoders
4. Interlude: what do we think pretraining is teaching?
5. Very large models and in-context learning

Motivating word meaning and context

Recall the adage we mentioned at the beginning of the course:

“You shall know a word by the company it keeps” (J. R. Firth 1957: 11)

This quote is a summary of **distributional semantics**, and motivated **word2vec**. But:

“... the complete meaning of a word is always contextual, and no study of meaning apart from a complete context can be taken seriously.” (J. R. Firth 1935)

Consider *I **record** the **record***: the two instances of ***record*** mean different things.

Where we were: pretrained word embeddings

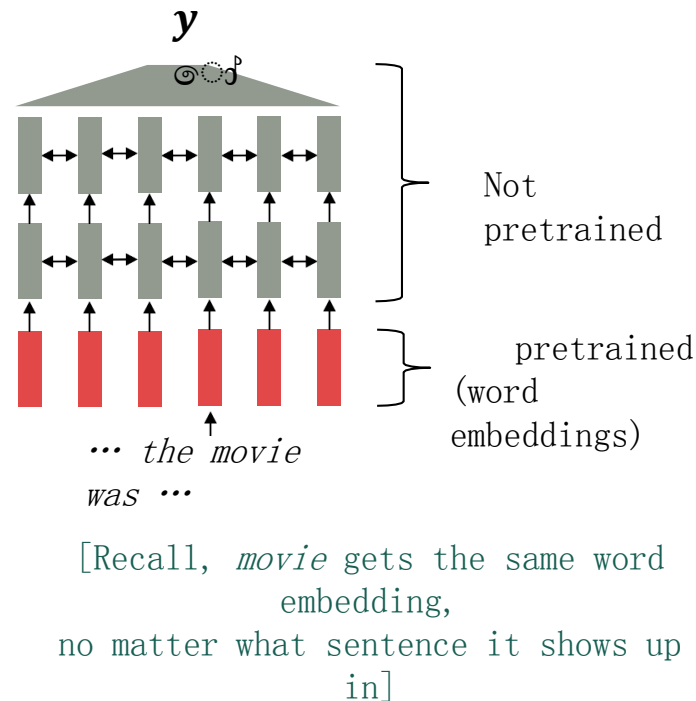
Circa

2017: Start with pretrained word embeddings (no context!)

- Learn how to incorporate context in an LSTM or Transformer while training on the task.

Some issues to think

- **about:** the training data we have for our **downstream task** (like question answering) must be sufficient to teach all contextual aspects of language.
- Most of the parameters in our network are randomly initialized!



Where we're going: pretraining whole models

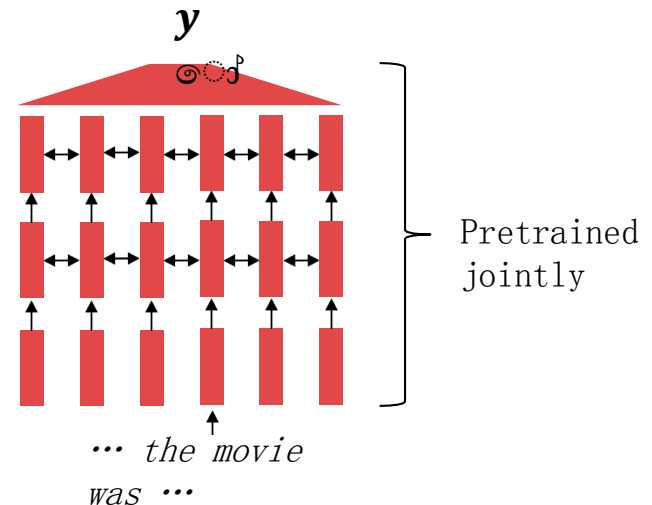
In modern

• **NLP** All (or almost all) parameters in NLP

networks are initialized via **pretraining**.

• Pretraining methods hide parts of the input from the model, and train the model to reconstruct those parts.

- This has been exceptionally effective at building strong:
 - **representations of language**
 - **parameter initializations** for strong NLP models.
 - **Probability distributions** over language that we can sample from



[This model has learned how to represent entire sentences through pretraining]

What can we learn from reconstructing
the input?

Shandong University is located in____, China.

Jinan
Qingdao
Weihai

What can we learn from reconstructing
the input?

I put__fork down on the
table.

the

What can we learn from reconstructing the input?

The woman walked across the street, checking for traffic over
_____her should
r.

What can we learn from reconstructing
the input?

I went to the ocean to see the fish, turtles,
seals, and

anything in the ocean

.

What can we learn from reconstructing the input?

Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink.

The movie was____.
suck

What can we learn from reconstructing
the input?

Iroh went into the kitchen to make some
tea.

Standing next to Iroh, Zuko pondered his
destiny.

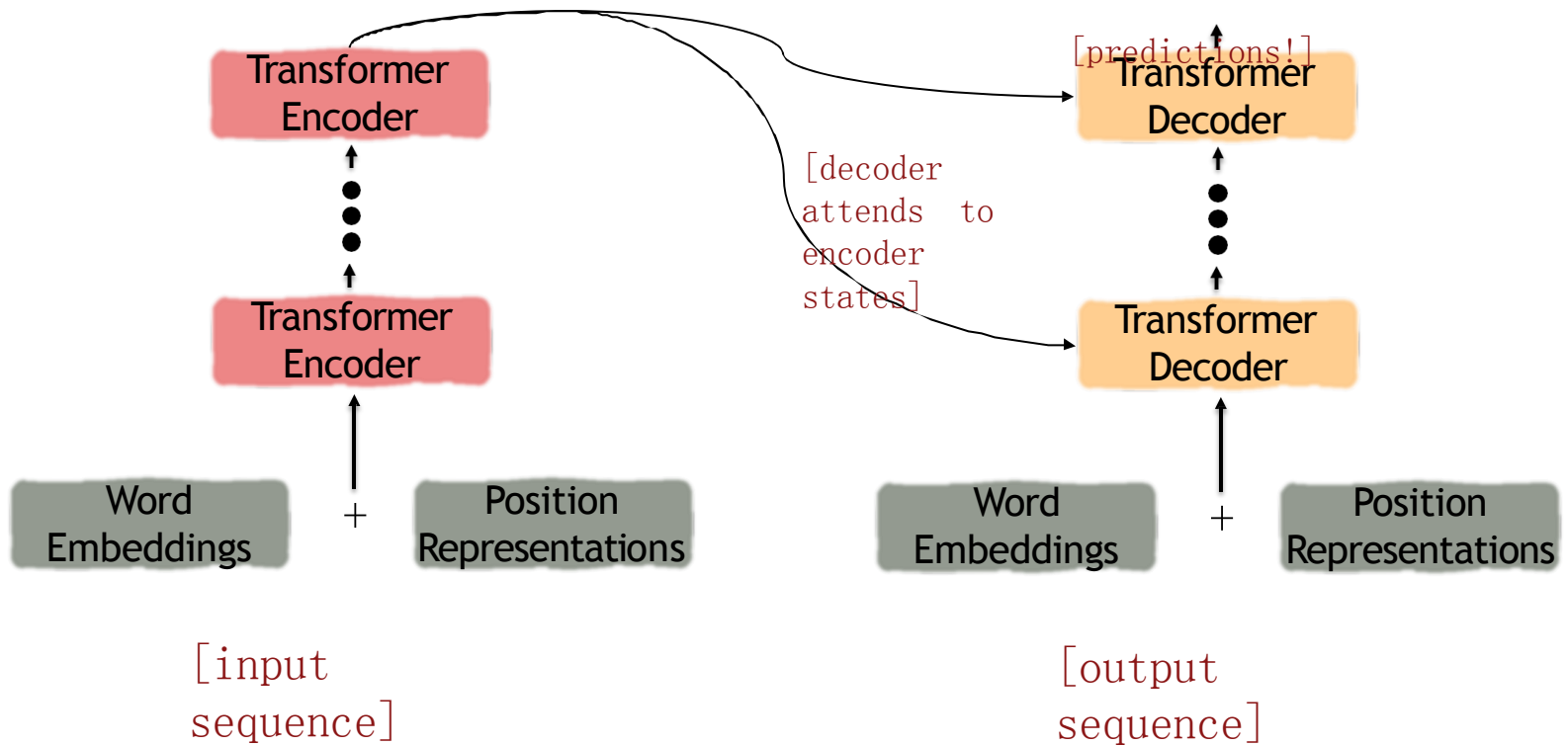
Zuko left the ~~kitchen~~.

What can we learn from reconstructing
the input?

I was thinking about the sequence
that goes 1, 1, 2, 3, 5, 8,
13, 21, _____

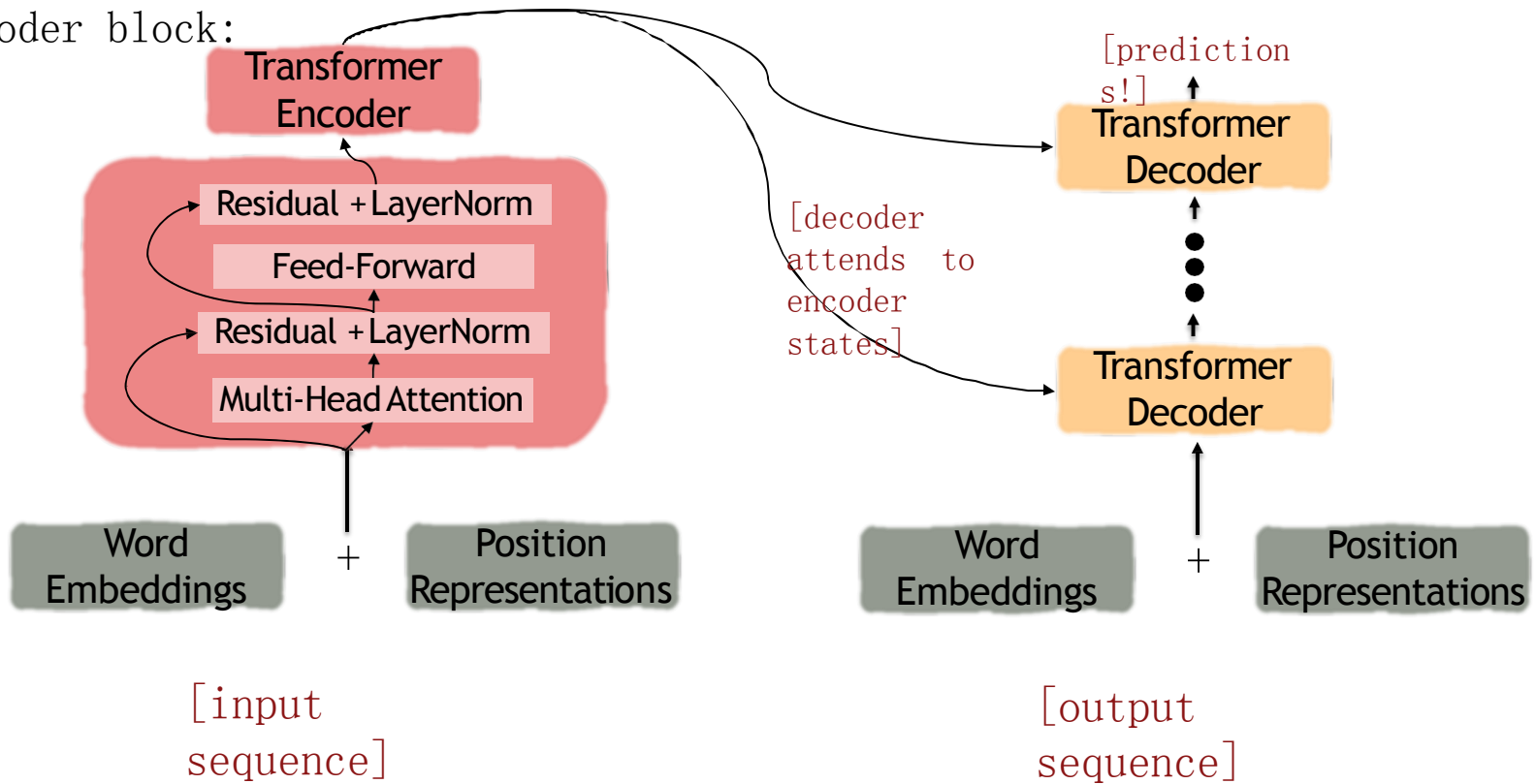
The Transformer Encoder-Decoder [\[Vaswani et al., 2017\]](#)

Looking back at the whole model, zooming in on an Encoder block:



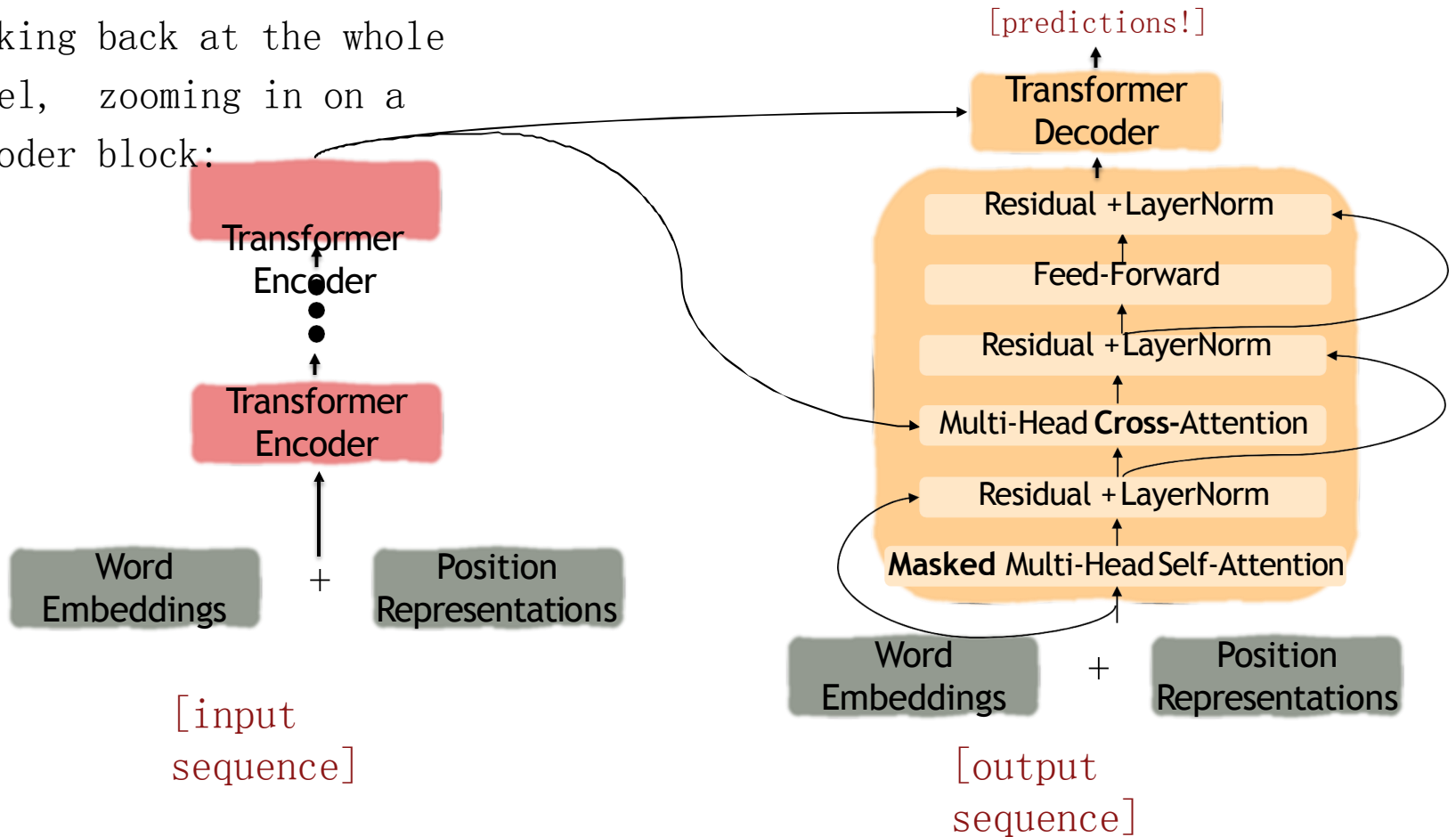
The Transformer Encoder-Decoder [\[Vaswani et al., 2017\]](#)

Looking back at the whole model, zooming in on an Encoder block:



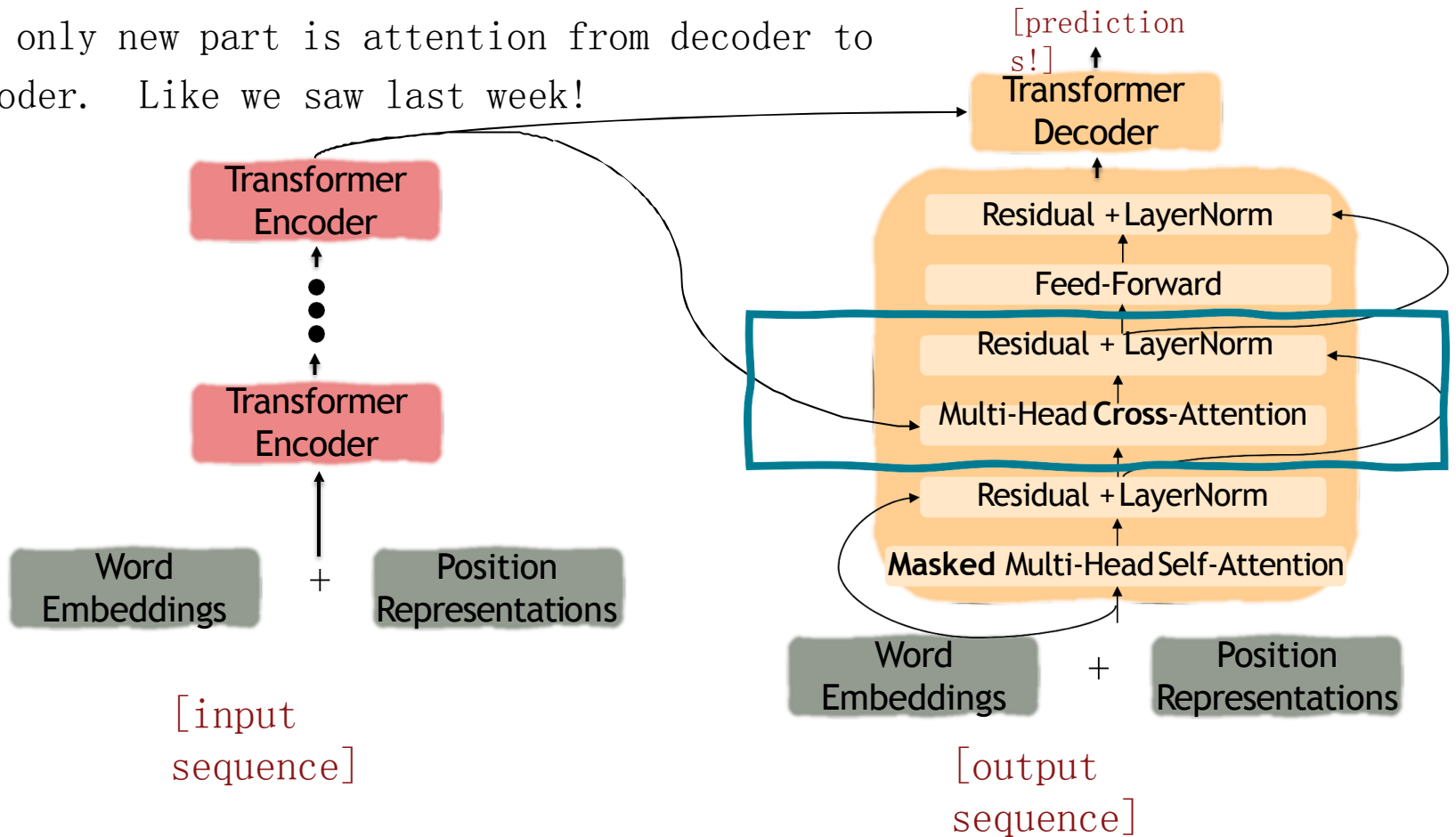
The Transformer Encoder-Decoder [Vaswani et al., 2017]

Looking back at the whole model, zooming in on a Decoder block:



The Transformer Encoder-Decoder [Vaswani et al., 2017]

The only new part is attention from decoder to encoder. Like we saw last week!



Pretraining through language modeling [\[Dai and Le, 2015\]](#)

Recall the **language modeling**

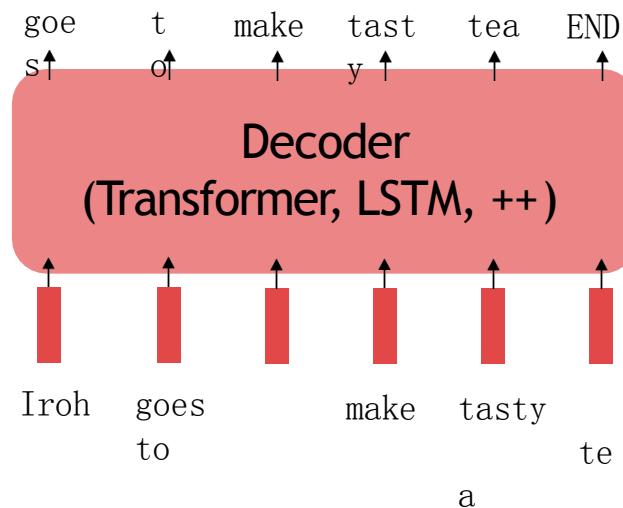
task: Model $p_{\theta}(w_t|w_{1:t-1})$, the probability distribution over words given their past contexts.

- There's lots of data for this!

Pretraining through language

modeling: train a neural network to perform language modeling on a large amount of text.

- Save the network parameters.



The Pretraining / Finetuning Paradigm

Pretraining can improve NLP applications by serving as parameter initialization.

Step 1: Pretrain (on language modeling)

Lots of text; learn general

goes t make tasty tea END



Decoder

(Transformer, LSTM, ++)

Iro goe t make tast te
h s o y a

Step 2: Finetune (on your task)

Not many labels; adapt to

the task! ☺



Decoder

(Transformer, LSTM, ++)

... the movie
was ...

Stochastic gradient descent and pretrain/finetune

Why should pretraining and finetuning help, from a “training neural nets” perspective?

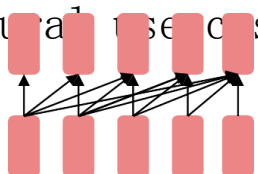
- Consider, provides parameters $\hat{\theta}$ by approximating $\min_{\theta} \mathcal{L}_{\text{pretrain}}(\theta)$
 - (The pretraining loss.)
- Then, finetuning approximates $\min_{\theta} \mathcal{L}_{\text{fine}}(\theta)$, starting at $\hat{\theta}$.
 - (The finetuning loss)
- The pretraining may matter because stochastic gradient descent sticks (relatively) close to $\hat{\theta}$ during finetuning.
 - So, maybe the finetuning local minima near $\hat{\theta}$ tend to generalize well!
 - And/or, maybe the gradients of finetuning loss near $\hat{\theta}$ propagate nicely!

Outline

1. A brief note on subword modeling
2. Motivating model pretraining from word embeddings
3. Model pretraining three ways
 1. Decoders
 2. Encoders
 3. Encoder–Decoders
4. Interlude: what do we think pretraining is teaching?
5. Very large models and in-context learning

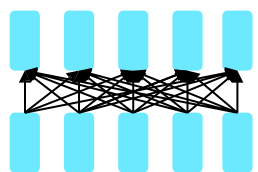
Pretraining for three types of architectures

The neural architecture influences the type of pretraining, and natural use cases.



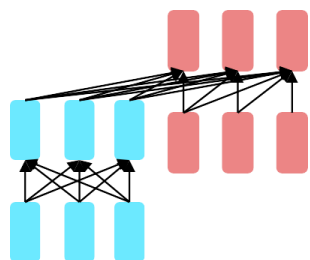
Decoder
s

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



Encoder
s

- Gets bidirectional context - can condition on future!
- Wait, how do we pretrain them?

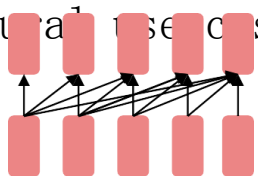


Encoder
-
Decoder
s

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

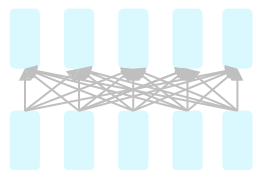
Pretraining for three types of architectures

The neural architecture influences the type of pretraining, and natural use cases.



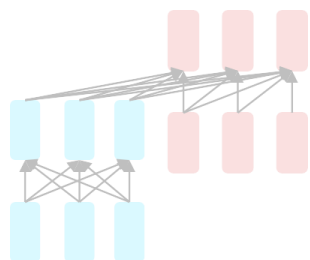
Decoder
s

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



Encoder
s

- Gets bidirectional context - can condition on future!
- Wait, how do we pretrain them?



Encoder
-
Decoder
s

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Pretraining decoders

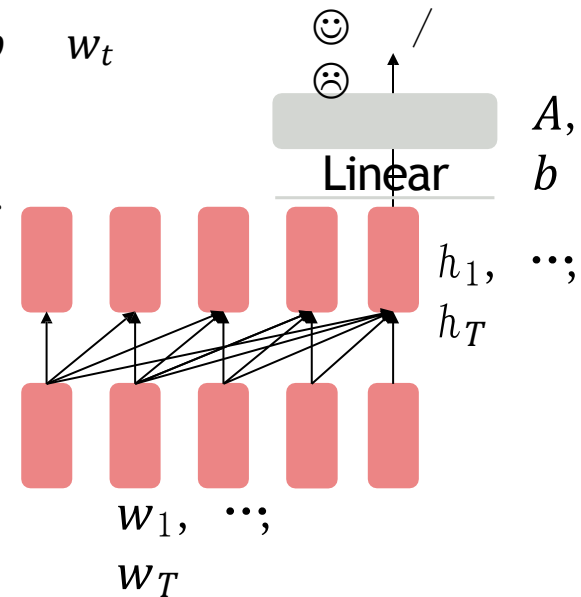
When using language model pretrained decoders, we can ignore that they were trained to model $p(w_t | w_{1:t-1})$.

We can finetune them by training a classifier on the last word's hidden state.

$$y \sim A w_T + b$$

Decoder w_1, \dots, w_T

Where A and b are randomly initialized and specified by the downstream task. Gradients backpropagate through the whole network.



[Note how the linear layer hasn't been pre-trained and must be learned from scratch.]

Pretraining decoders

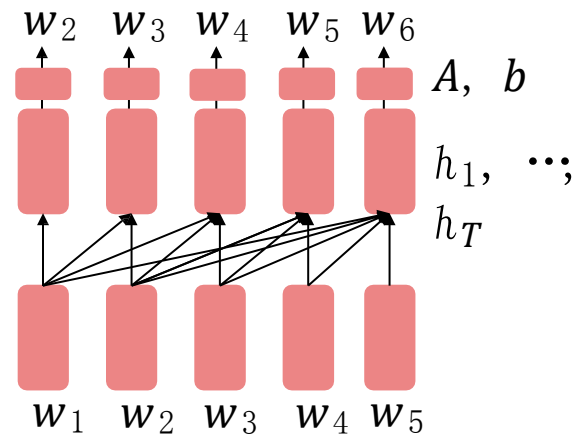
It's natural to pretrain decoders as language models and then use them as generators, finetuning their p_θ

This is helpful in tasks where the output is a sequence with a vocabulary like that at pretraining time!

- Dialogue (context=dialogue history)
- Summarization (context=document)

Where A, b were pre-trained in the language model!

$$(w_t | w_{1:t-1})!$$



[Note how the linear layer has been pre-trained.]

Generative Pretrained Transformer (GPT) [[Radford et al., 2018](#)]

2018' s GPT was a big success in pretraining a decoder!

- Transformer decoder with 12 layers.
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers.
- Byte-pair encoding with 40,000 merges
- Trained on BooksCorpus: over 7000 unique books.
 - Contains long spans of contiguous text, for learning long-distance dependencies.
- The acronym “GPT” never showed up in the original paper; it could stand for “Generative PreTraining” or “Generative Pretrained Transformer”

[[Devlin et al., 2018](#)]

Generative Pretrained Transformer (GPT) [[Radford et al., 2018](#)]

How do we format inputs to our decoder for **finetuning tasks**?

Natural Language Inference: Label pairs of sentences as

Premise: The man is in the doorway } **entailment**

Hypothesis: The person is near the door

Radford et al., 2018 evaluate on natural language inference.

Here's roughly how the input was formatted, as a sequence of tokens for the decoder.

[START] *The man is in the doorway* [DELIM] *The person is near the door* [EXTRACT]

The linear classifier is applied to the representation of the [EXTRACT] token.

Generative Pretrained Transformer (GPT) [[Radford et al., 2018](#)]

GPT results on various *natural language inference* datasets.

Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	<u>89.3</u>	-	-	-
CAFE [58] (5x)	80.2	79.0	<u>89.3</u>	-	-	-
Stochastic Answer Network [35] (3x)	<u>80.6</u>	<u>80.1</u>	-	-	-	-
CAFE [58]	78.7	77.9	88.5	<u>83.3</u>		
GenSen [64]	71.4	71.3	-	-	<u>82.3</u>	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	82.1	61.7
Finetuned Transformer LM (ours)	82.1	81.4	89.9	88.3	88.1	56.0

Increasingly convincing generations (GPT2) [[Radford et al., 2018](#)]

We mentioned how pretrained decoders can be used **in their capacities as language models**. GPT-2, a larger version of GPT trained on more data, was shown to produce relatively convincing samples of natural language.

Context (human-written): In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

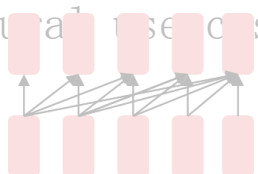
GPT-2: The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

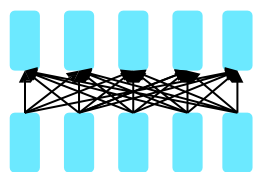
Pretraining for three types of architectures

The neural architecture influences the type of pretraining, and natural use cases.



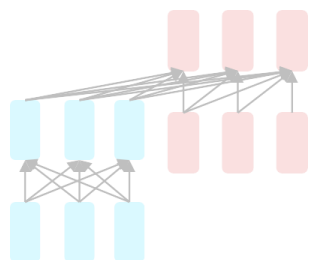
Decoder
s

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



Encoder
s

- Gets bidirectional context - can condition on future!
- Wait, how do we pretrain them?



Encoder
-
Decoder
s

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

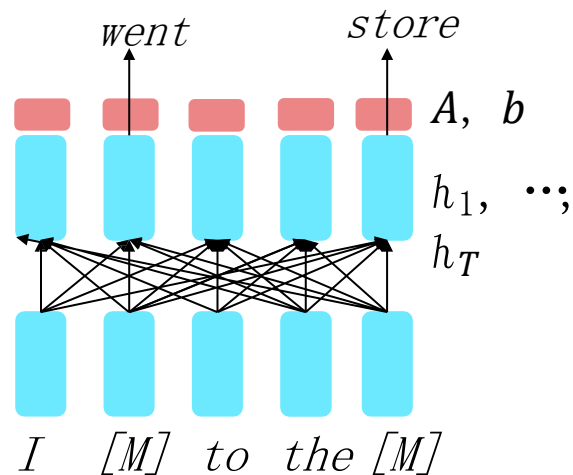
Pretraining encoders: what pretraining objective to use?

So far, we've looked at language model pretraining. But **encoders get bidirectional context**, so we can't do language modeling!

Idea: replace some fraction of words in the input with a special [MASK] token; predict these words.

$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$
$$y_i \sim Aw_i + b$$

Only add loss terms from words that are "masked out." If \tilde{x} is the masked version of x , we're learning $p_\theta(x | \tilde{x})$. Called **Masked LM**.



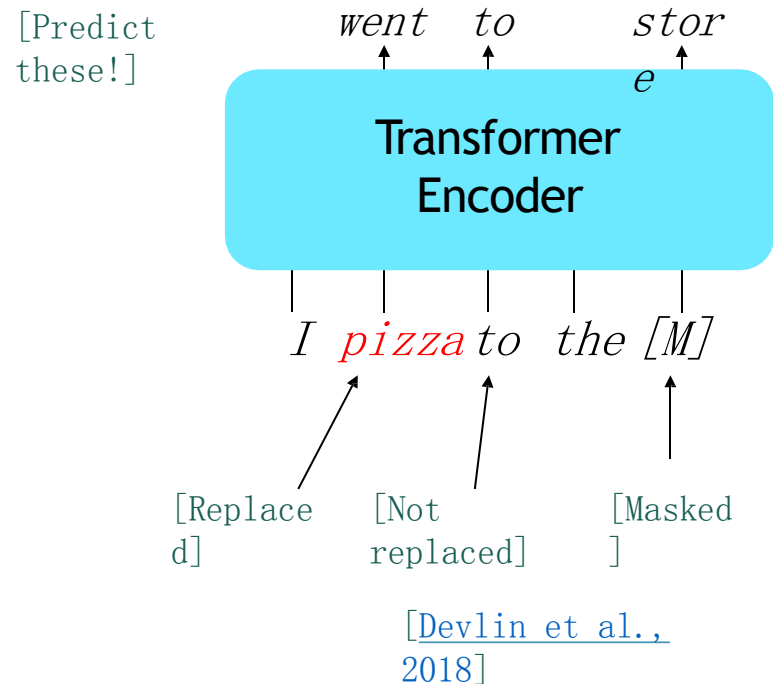
[[Devlin et al., 2018](#)]

BERT: Bidirectional Encoder Representations from Transformers

Devlin et al., 2018 proposed the “Masked LM” objective and **released the weights of a pretrained Transformer**, a model they labeled BERT.

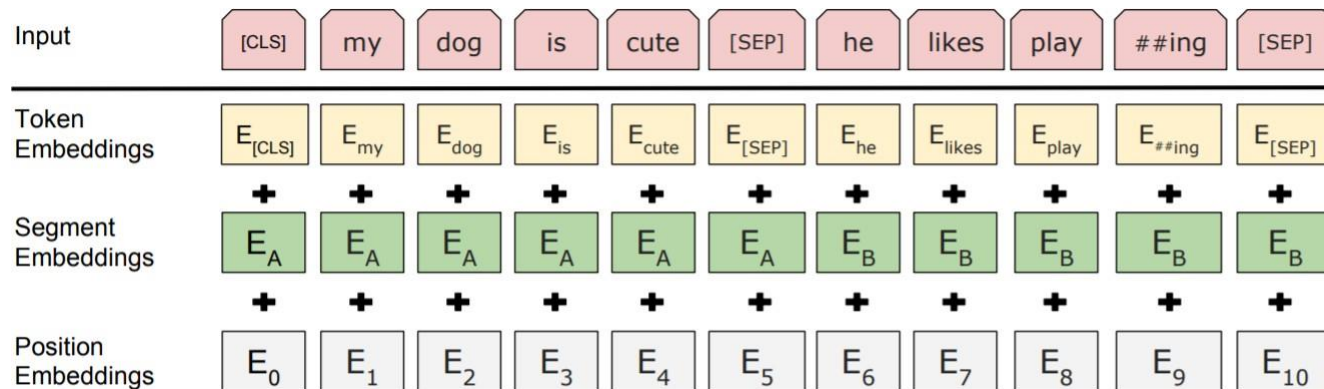
Some more details about Masked LM for BERT:

- Predict a random 15% of (sub)word tokens.
 - Replace input word with [MASK] 80% of the time
 - Replace input word with a random token 10% of the time
 - Leave input word unchanged 10% of the time (but still predict it!)
- Why? Doesn't let the model get complacent and not build strong representations of non-masked words. (No masks are seen at fine-tuning time!)



BERT: Bidirectional Encoder Representations from Transformers

- The pretraining input to BERT was two separate contiguous chunks of text:



- BERT was trained to predict whether one chunk follows the other or is randomly sampled.
 - Later work has argued this “next sentence prediction” is not necessary.

[\[Devlin et al., 2018, Liu et al., 2019\]](#)

BERT: Bidirectional Encoder Representations from Transformers

Details about BERT

- Two models were released:
 - BERT-base: 12 layers, 768-dim hidden states, 12 attention heads, 110 million params.
 - BERT-large: 24 layers, 1024-dim hidden states, 16 attention heads, 340 million params.
- Trained on:
 - BooksCorpus (800 million words)
 - English Wikipedia (2,500 million words)
- Pretraining is expensive and impractical on a single GPU.
 - BERT was pretrained with 64 TPU chips for a total of 4 days.
 - (TPUs are special tensor operation acceleration hardware)
- Finetuning is practical and common on a single GPU
 - “Pretrain once, finetune many times.”

[\[Devlin et al., 2018\]](#)

BERT: Bidirectional Encoder Representations from Transformers

BERT was massively popular and hugely versatile; finetuning BERT led to new state-of-the-art results on a broad range of tasks.

- **QQP**: Quora Question Pairs (detect paraphrase questions)
- **QNLI**: Question Natural Language Inference over question
- **SST-2**: Sentiment analysis
- **CoLA**: corpus of linguistic acceptability (detect whether sentences are grammatical.)
- **STS-B**: semantic textual similarity
- **MRPC**: microsoft paraphrase corpus
- **RTE**: a small natural language inference

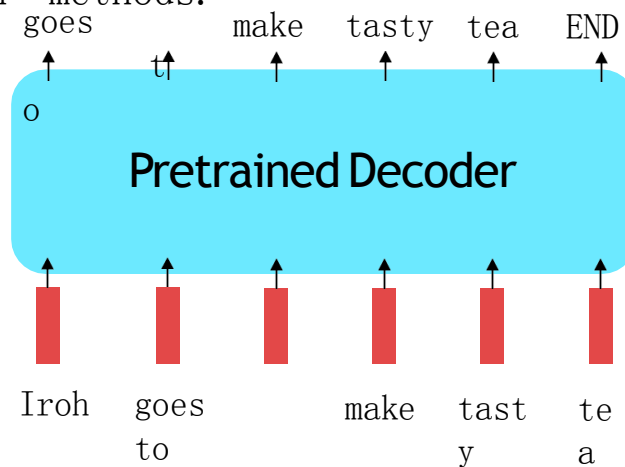
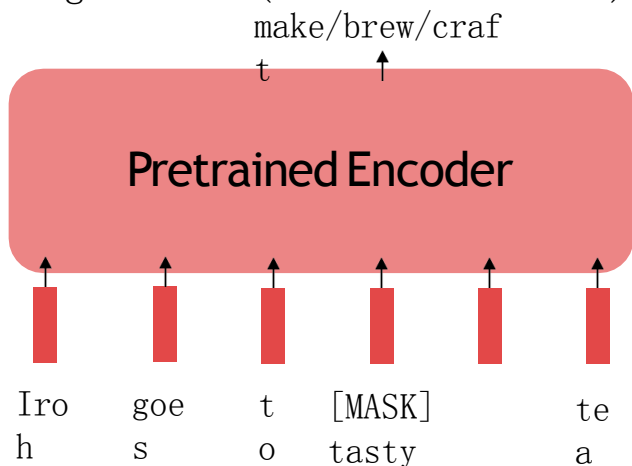
System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

[[Devlin et al., 2018](#)]

Limitations of pretrained encoders

Those results looked great! Why not used pretrained encoders for everything?

If your task involves generating sequences, consider using a pretrained decoder; BERT and other pretrained encoders don't naturally lead to nice autoregressive (1-word-at-a-time) generation methods.

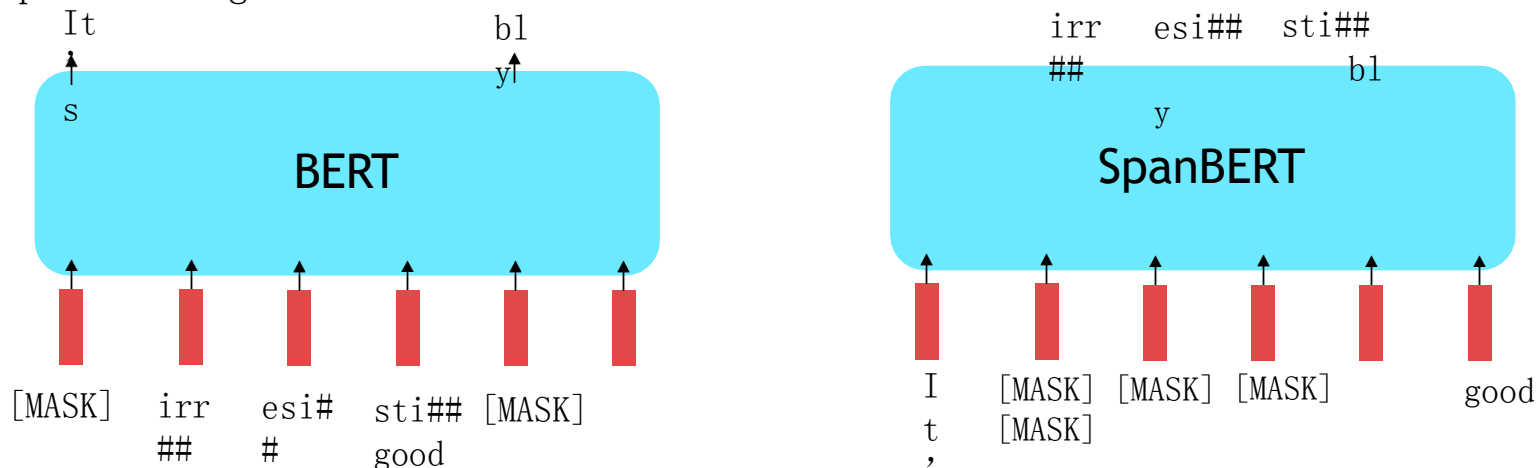


Extensions of BERT

You'll see a lot of BERT variants like RoBERTa, SpanBERT, +++

Some generally accepted improvements to the BERT pretraining formula:

- RoBERTa: mainly just train BERT for longer and remove next sentence prediction!
- SpanBERT: masking contiguous spans of words makes a harder, more useful pretraining task



[[Liu et al., 2019](#); [Joshi et al., 2020](#)]

Extensions of BERT

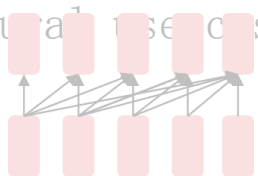
A takeaway from the RoBERTa paper: more compute, more data can improve pretraining even when not changing the underlying Transformer encoder.

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	94.6/89.4	90.2	96.4
BERT _{LARGE}						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7

[[Liu et al., 2019](#); [Joshi et al., 2020](#)]

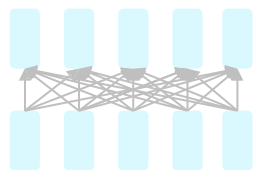
Pretraining for three types of architectures

The neural architecture influences the type of pretraining, and natural use cases.



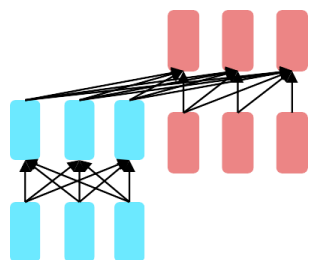
Decoder
s

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



Encoder
s

- Gets bidirectional context - can condition on future!
- Wait, how do we pretrain them?



Encoder
-
Decoder
s

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

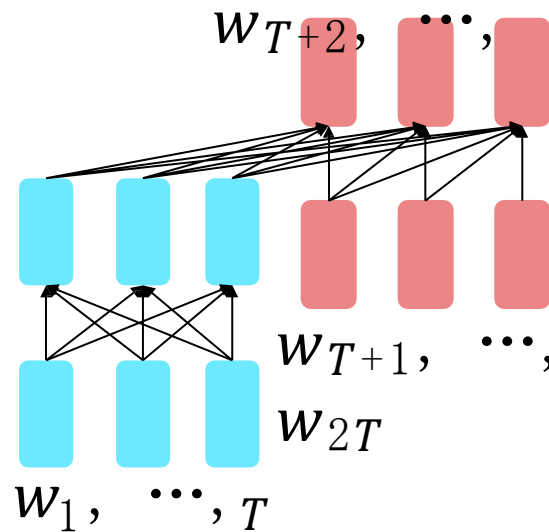
Pretraining encoder-decoders: what pretraining objective to use?

For **encoder-decoders**, we could do something like **language modeling**, but where a prefix of every input is provided to the encoder and is not predicted.

$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

$$h_{T+1}, \dots, h_{2T} = \text{Decoder}(w_1, \dots, w_T, h_1, \dots, h_T)$$

The **encoder** portion benefits from bidirectional context; the **decoder** portion is used to train the whole model through language modeling.



w [Raffel et al., 2018]

Pretraining encoder-decoders: what pretraining objective to use?

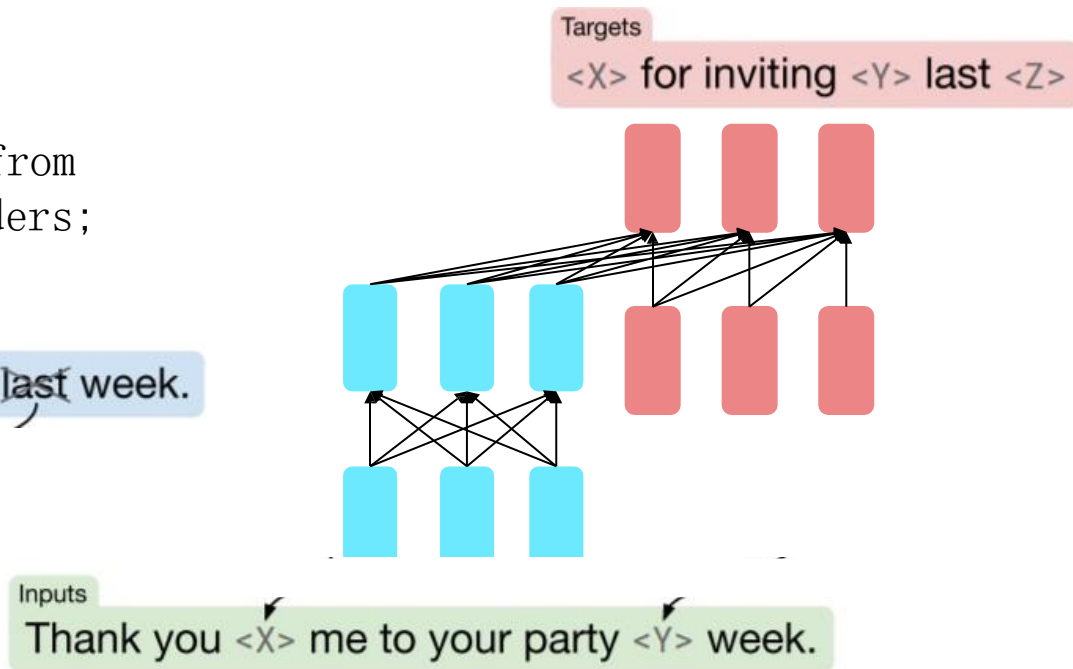
What [Raffel et al., 2018](#) found to work best was **span corruption**.
Their model: **T5**.

Replace different-length spans from the input with unique placeholders; decode out the spans that were

Original text

Thank you for inviting me to your party last week.

This is implemented in text preprocessing: it's still an objective that looks like **language modeling** at the decoder side.



Pretraining encoder-decoders: what pretraining objective to use?

[Raffel et al., 2018](#) found encoder-decoders to work better than decoders for their tasks, and span corruption (denoising) to work better than language modeling.

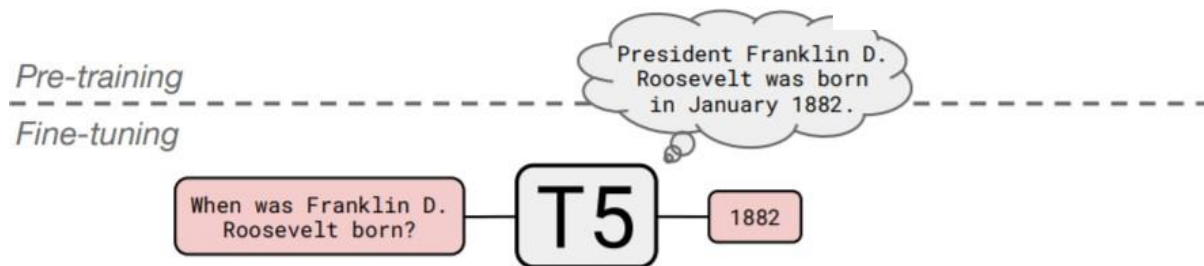
Architecture	Objective	Params	Cost	GLUE	CNN3M	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	M	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Enc-dec, shared	Denoising	P	M	82.81	18.78	80.63	70.73	26.72	39.03	27.46
Enc-dec, 6 layers	Denoising	P	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	P	M	74.70	17.93	61.14	55.02	25.09	35.28	25.86
Prefix LM	Denoising	P	M	81.82	18.61	78.94	68.11	26.43	37.98	27.39
Encoder-decoder	LM	$2P$	M	79.56	18.59	76.02	64.29	26.27	39.17	26.86
Enc-dec, shared	LM	P	M	79.60	18.13	76.35	63.50	26.62	39.17	27.05
Enc-dec, 6 layers	LM	P	$M/2$	78.67	18.26	75.32	64.06	26.13	38.42	26.89
Language model	LM	P	M	73.78	17.54	53.81	56.51	25.23	34.31	25.38
Prefix LM	LM	P	M	79.68	17.84	76.87	64.86	26.28	37.51	26.76

Pretraining encoder-decoders: what pretraining objective to use?

A fascinating property of T5: it can be finetuned to answer a wide range of questions, retrieving knowledge from its parameters.

NQ: Natural Questions
 WQ: WebQuestions
 TQA: Trivia QA

All “open-domain” ns



	NQ	WQ	TQA		
			dev	test	
<u>Karpukhin et al. (2020)</u>	41.5	42.4	57.9	–	
T5.1.1-Base	25.7	28.2	24.2	30.6	220 million
T5.1.1-Large	27.3	29.5	28.5	37.2	params
T5.1.1-XL	29.5	32.4	36.0	45.1	770 million
T5.1.1-XXL	32.8	35.6	42.9	52.5	params
<u>T5.1.1-XXL + SSM</u>	35.2	42.8	51.9	61.6	3 billion params
					11 billion
					params

[Raffel et al., 2018]

Outline

1. A brief note on subword modeling
2. Motivating model pretraining from word embeddings
3. Model pretraining three ways
 1. Decoders
 2. Encoders
 3. Encoder–Decoders
4. Interlude: what do we think pretraining is teaching?
5. Very large models and in-context learning

What kinds of things does pretraining learn?

There's increasing evidence that pretrained models learn a wide variety of things about the statistical properties of language. Taking our examples from the start of class:

- *Shandong University is located in___, China.* [Trivia]
- *I put___fork down on the table.* [syntax]
- *The woman walked across the street, checking for traffic over _____shoulder.*
[coreference]
- *I went to the ocean to see the fish, turtles, seals, and_____.* [lexical semantics/topic]
- *Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was_____.* [sentiment]
- *Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the_____.* [some reasoning - this is harder]
- *I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21,*
[some basic

Outline

1. A brief note on subword modeling
2. Motivating model pretraining from word embeddings
3. Model pretraining three ways
 1. Decoders
 2. Encoders
 3. Encoder–Decoders
4. Interlude: what do we think pretraining is teaching?
5. Very large models and in-context learning

GPT-3, In-context learning, and very large models

So far, we've interacted with pretrained models in two ways:

- Sample from the distributions they define (maybe providing a prompt)
- Fine-tune them on a task we care about, and take their predictions.

Very large language models seem to perform some kind of learning **without gradient steps** simply from examples you provide within their contexts.

GPT-3 is the canonical example of this. The largest T5 model had 11 billion parameters.

GPT-3 has 175 billion parameters.

GPT-3, In-context learning, and very large models

- Very large language models seem to perform some kind of learning **without gradient steps** simply from examples you provide within their contexts.

- The in-context examples seem to specify the task to be performed, and the conditional distribution mocks performing the task to a certain extent.

- **Input (prefix within a single Transformer decoder context):**

Output (conditional generations):

thanks ->

merci hello

-> bonjour

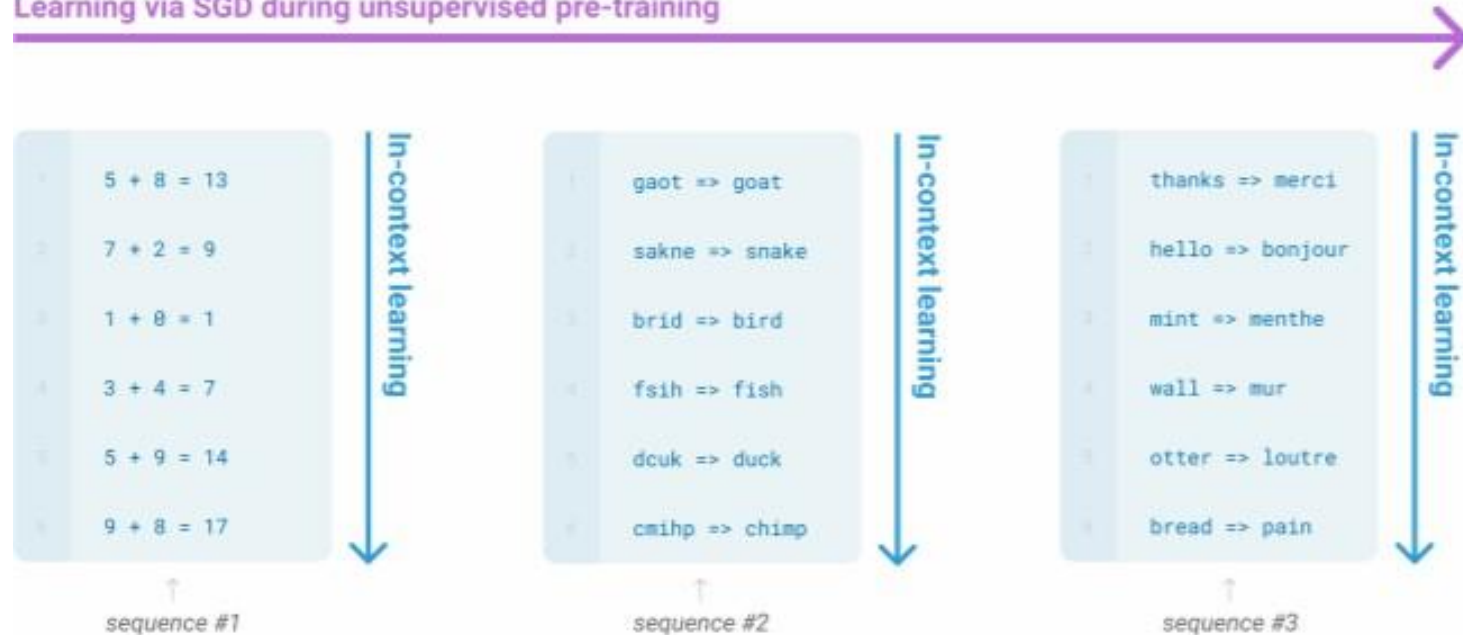
mint ->

menthe otter

GPT-3, In-context learning, and very large models

Very large language models seem to perform some kind of learning **without gradient steps** simply from examples you provide within their contexts.

Learning via SGD during unsupervised pre-training



Summary

1. A brief note on subword modeling
2. Motivating model pretraining from word embeddings
3. Model pretraining three ways
 1. Decoders
 2. Encoders
 3. Encoder–Decoders
4. Interlude: what do we think pretraining is teaching?
5. Very large models and in-context learning