



# Natural Language Processing

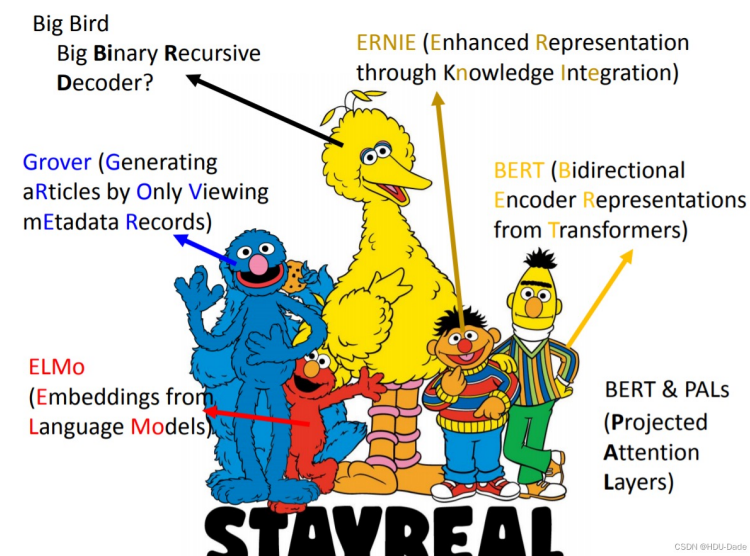
ELMo, BERT, GPT, and T5 (BERT and its Family)



# Outline

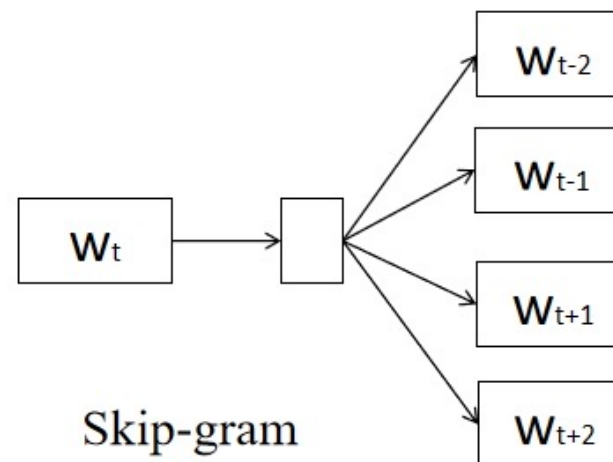
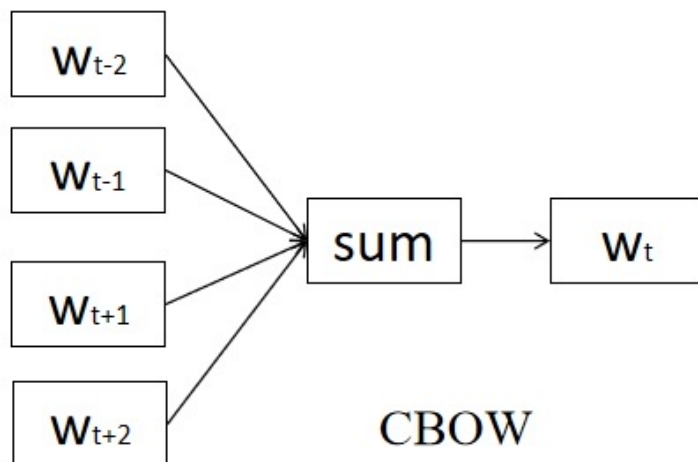
---

- Recap: Word Embeddings, RNN
- Pretraining Language Model from Word Embeddings
  - ELMo
- Pretraining by Transformer
  - Encoders (BERT)
  - Encoder-Decoders (T5)
  - Decoders (GPT)
- GPT-3 with In-Context Learning and Large Language Models



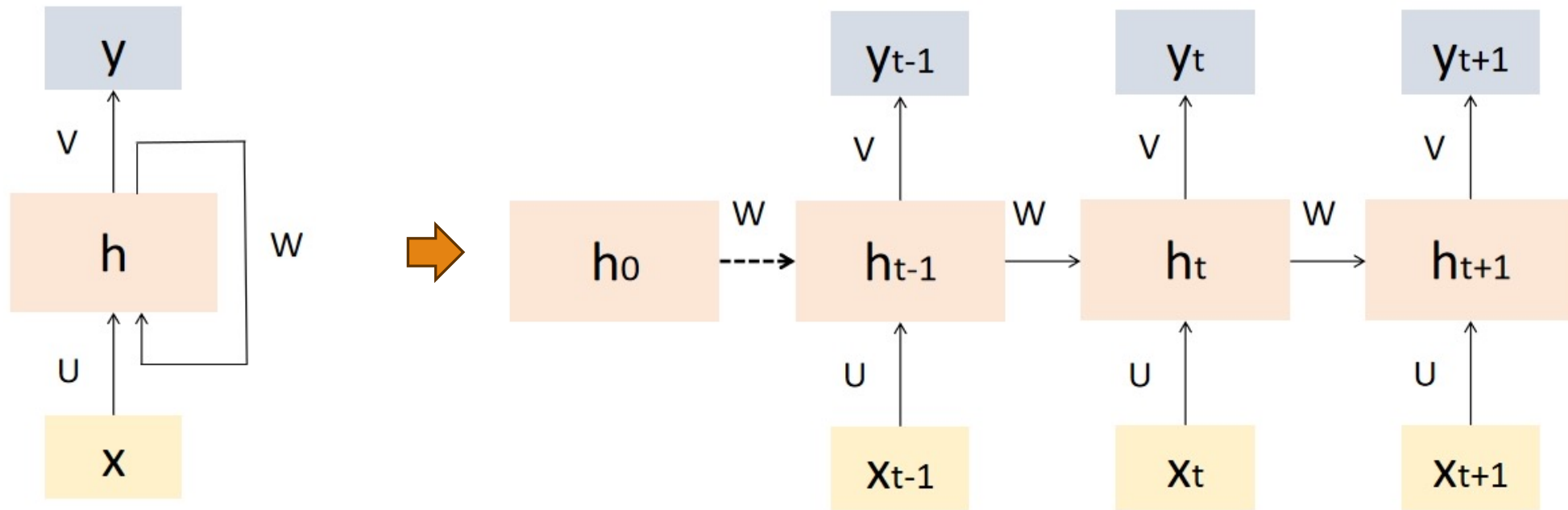
# Recap: Word Embeddings (Word2Vec)

- Skip-gram
  - Use words to predict their contexts.
- CBOW
  - Use the context to predict the target word.



# Recap: RNN

- Recurrent Neural Networks (RNNs) are a type of artificial neural network designed to handle sequential data by **capturing temporal dependencies**.
  - The same set of weights and biases are used across all time steps



# Pretrained Word Embeddings

---

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

-0.17	-0.89	0.06	-0.41	-1.08	-0.74	1.41	-0.94	-0.32	0.12
-------	-------	------	-------	-------	-------	------	-------	-------	------

Word2Vec or GloVe word embeddings of “**record**”

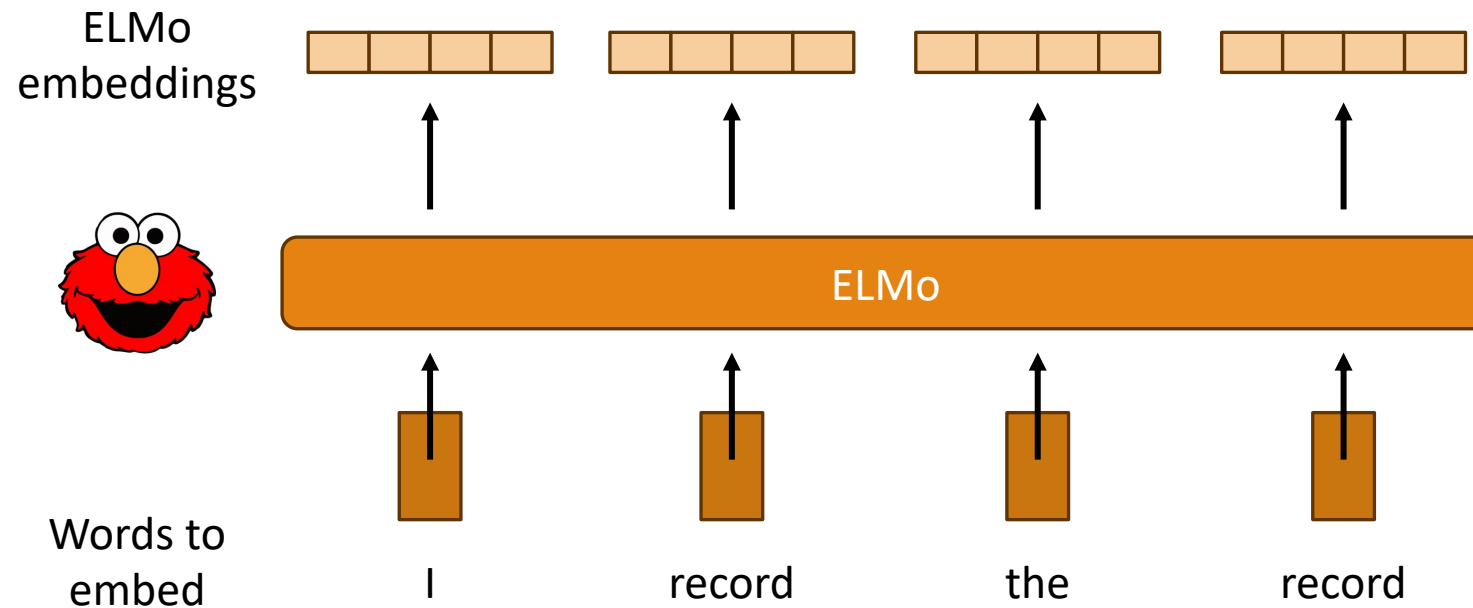


Vector representation

Issue: Both GloVe and Word2Vec represent the word “record” as the same vector regardless of the context.

# Context Matters: ELMo (Peters et al., 2018)

Rather than using a fixed representation for each word, **ELMo (Embeddings from Language Models)** considers the entire sentence before assigning each word in it an embedding.



# Reconstruct the Sentences

---

National Hsing Hua University is in \_\_\_\_\_. (Hsinchu)

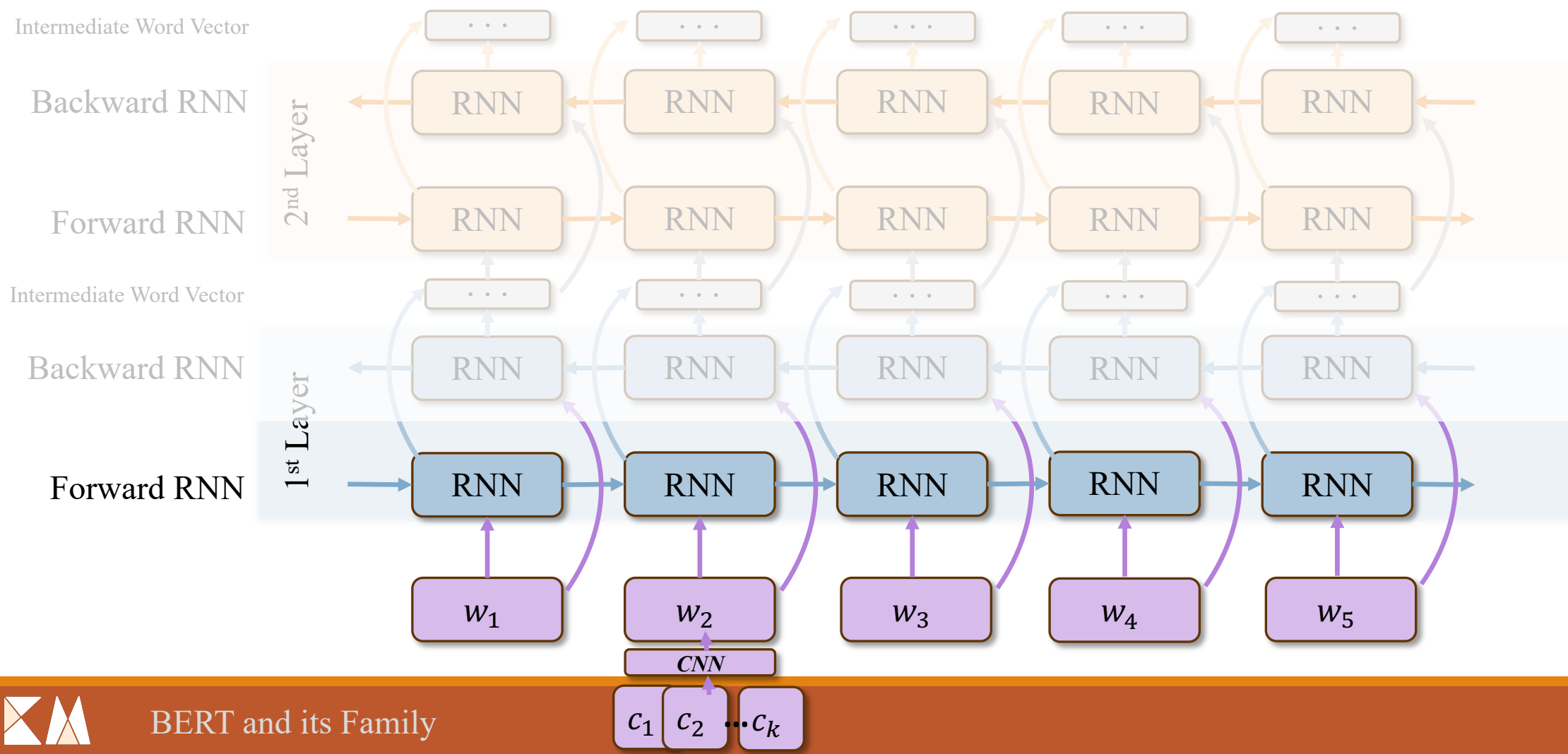
I put \_\_\_\_ bag on the table. (a)

I went to the zoo to see the lions, monkeys, elephants, and \_\_\_\_\_. (.)

The overall value I got from the two hours of watching it was the total sum of popcorn and the drink. The movie was \_\_\_\_\_. (bad)

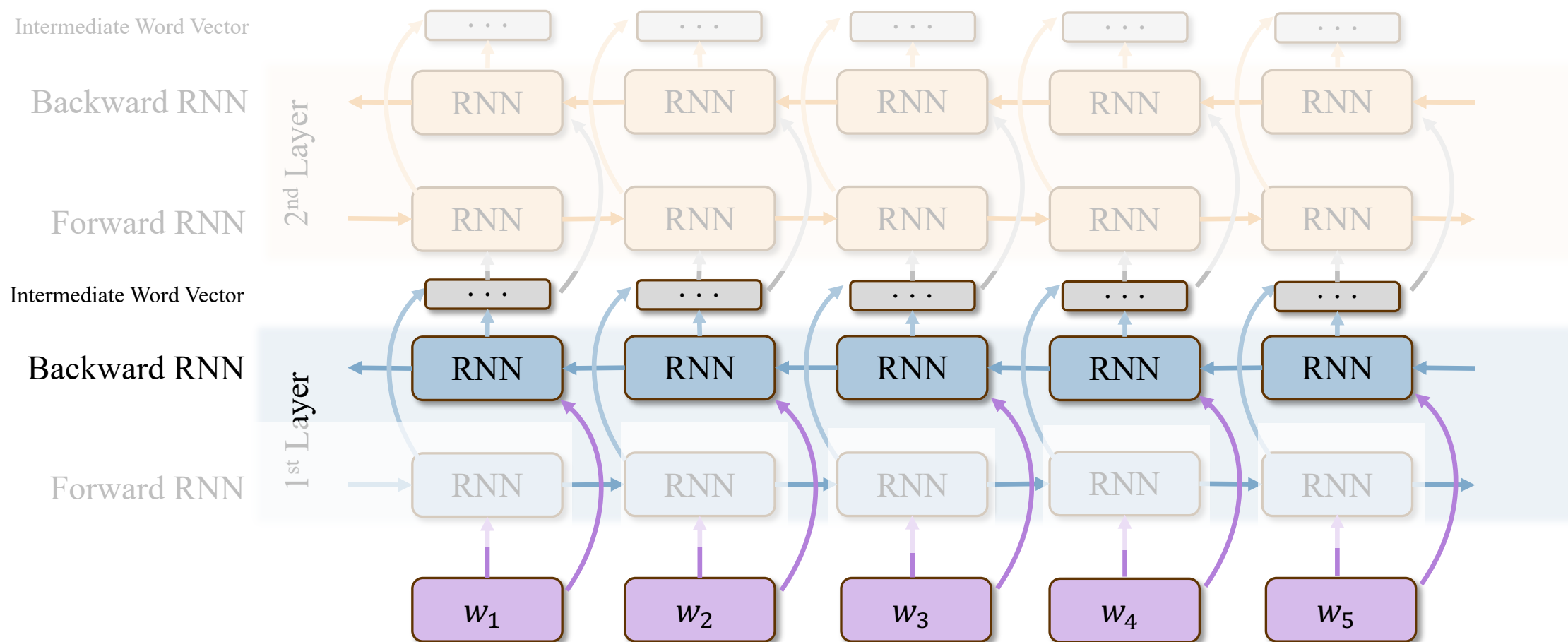
I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_(34)

# Bidirectional Language Model (biLM)

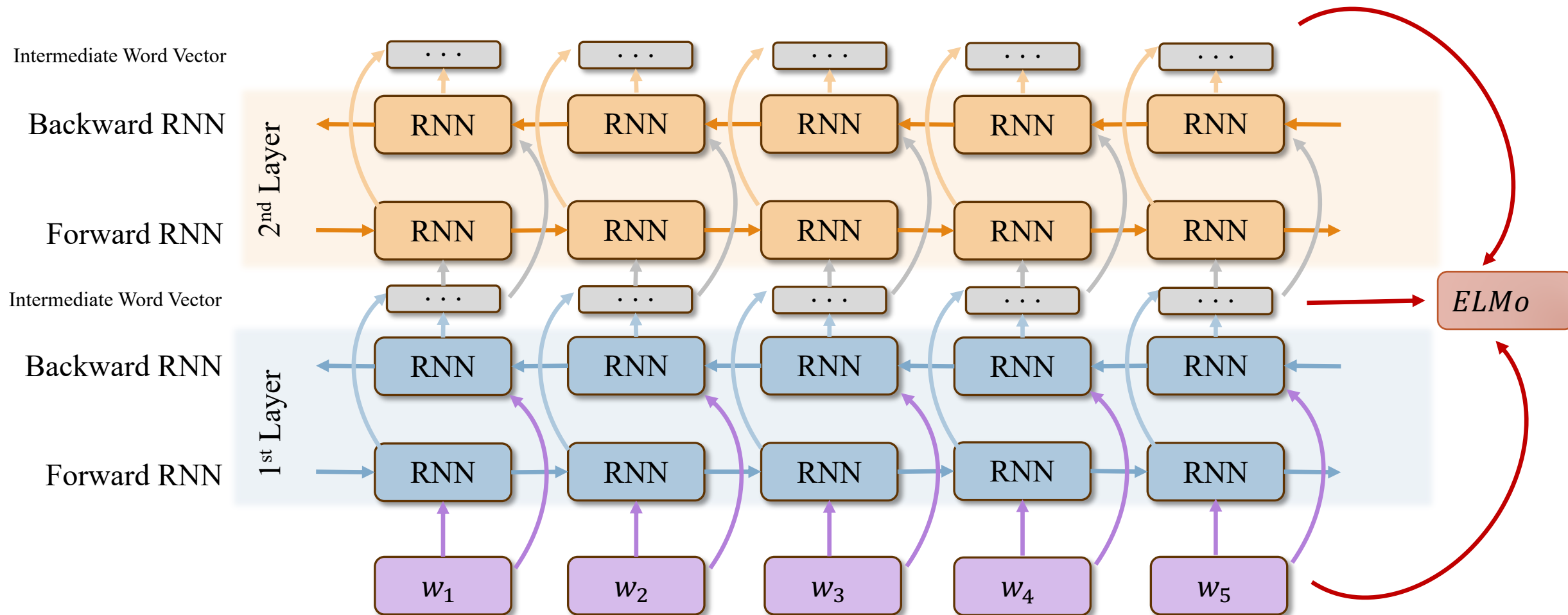




# Bidirectional Language Model (biLM)



# Bidirectional Language Model (biLM)



# Pretraining through Language Modeling

## Step 1. Concatenate hidden layers.

Concatenate forward and backward LM of each layer together (for token layer, we concatenate it with itself)

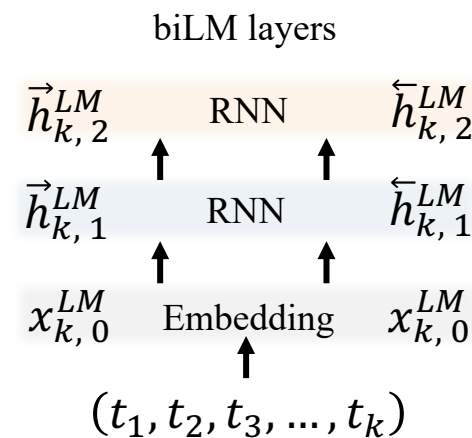
$\vec{h}_{k,0}^{LM}$  : Token Layer

$\vec{h}_{k,j}^{LM}$  : Forward LM at j layer

$\overleftarrow{h}_{k,j}^{LM}$  : Backward LM at j layer

$$\text{ELMo}_k^{\text{task}} = \gamma^{\text{task}} \sum_{j=0}^L \left\{ \begin{array}{l} s_2^{\text{task}} \times \overline{h}_{k,2}^{LM} (j=2) \\ s_1^{\text{task}} \times \overline{h}_{k,1}^{LM} (j=1) \\ s_0^{\text{task}} \times \overline{h}_{k,0}^{LM} (j=0) \end{array} \right\}, h_{k,j}^{LM} = [\vec{h}_{k,j}^{LM}; \overleftarrow{h}_{k,j}^{LM}]$$

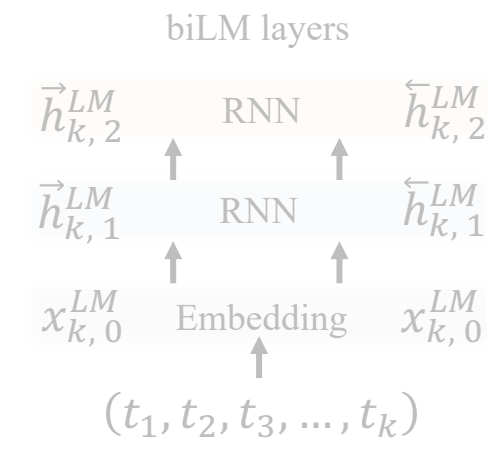
scalar



# Pretraining through Language Modeling

## Step 2. Multiply each vector by a weight

We assign each concatenated layer with softmax-normalized weight,  $s^{task}$

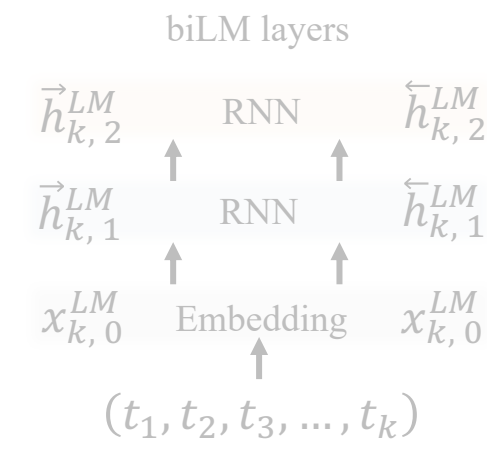
$$\text{ELMo}_k^{task} = \underbrace{\gamma^{task}}_{\text{scalar}} \sum_{j=0}^L \left\{ \begin{array}{l} s_2^{task} \times \overline{h_{k,2}^{LM}} \quad (j=2) \\ s_1^{task} \times \overline{h_{k,1}^{LM}} \quad (j=1) \\ s_0^{task} \times \overline{h_{k,0}^{LM}} \quad (j=0) \end{array} \right. , h_{k,j}^{LM} = [\vec{h}_{k,j}^{LM}; \overleftarrow{h}_{k,j}^{LM}]$$


The diagram illustrates the biLM layers architecture. At the bottom, a sequence of tokens  $(t_1, t_2, t_3, \dots, t_k)$  is processed by an 'Embedding' layer to produce forward vectors  $\vec{x}_{k,0}^{LM}$  and backward vectors  $\overleftarrow{x}_{k,0}^{LM}$ . These are then passed through two 'RNN' layers. The first RNN layer produces forward hidden states  $\vec{h}_{k,1}^{LM}$  and backward hidden states  $\overleftarrow{h}_{k,1}^{LM}$ . The second RNN layer produces forward hidden states  $\vec{h}_{k,2}^{LM}$  and backward hidden states  $\overleftarrow{h}_{k,2}^{LM}$ . The entire structure is labeled 'biLM layers'.

# Pretraining through Language Modeling

Step 3. Sum all weighted representations.

Once weighted vectors are added together, we use  $\gamma^{task}$  to allow the task model to scale the entire ELMo vector, which is crucial for optimization purposes.

$$\mathbf{ELMo}_k^{task} = \underbrace{\gamma^{task}}_{\text{scalar}} \sum_{j=0}^L \left\{ \begin{array}{l} s_2^{task} \times h_{k,2}^{LM} \quad (j=2) \\ s_1^{task} \times h_{k,1}^{LM} \quad (j=1) \\ s_0^{task} \times h_{k,0}^{LM} \quad (j=0) \end{array} \right. , h_{k,j}^{LM} = [\vec{h}_{k,j}^{LM}; \overleftarrow{h}_{k,j}^{LM}]$$


The diagram illustrates the biLM layers. At the bottom, the input sequence  $(t_1, t_2, t_3, \dots, t_k)$  is processed by an Embedding layer to produce  $x_{k,0}^{LM}$ . This input is fed into two parallel RNN layers. The forward RNN (left) produces hidden states  $\vec{h}_{k,1}^{LM}$  and  $\vec{h}_{k,2}^{LM}$ . The backward RNN (right) produces hidden states  $\overleftarrow{h}_{k,1}^{LM}$  and  $\overleftarrow{h}_{k,2}^{LM}$ . The hidden states for each layer  $j$  are combined into a vector  $h_{k,j}^{LM} = [\vec{h}_{k,j}^{LM}; \overleftarrow{h}_{k,j}^{LM}]$ .

# Pretraining through Language Modeling

---

ELMo is a type of deep contextualized word representation that are created through a learning process based on **hidden layers** of a deep bidirectional language model.

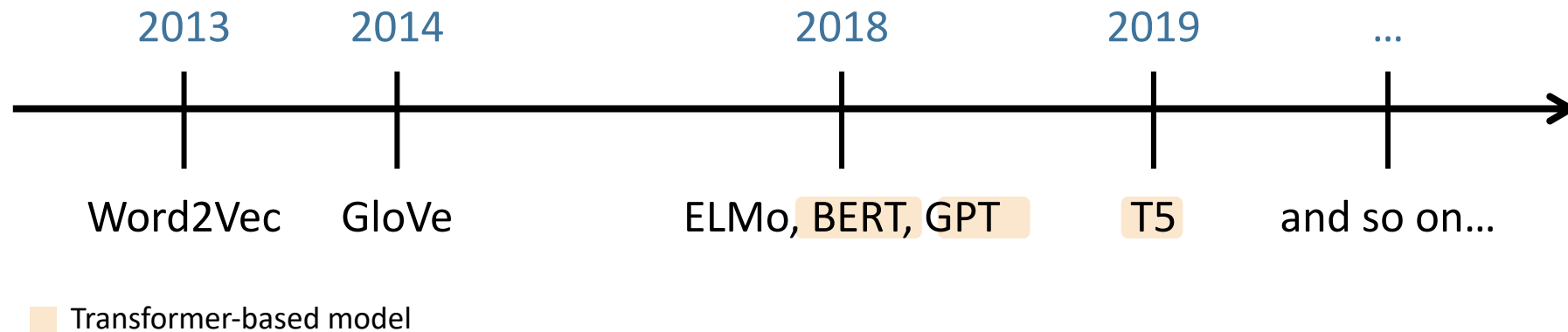
ELMo is a task specific combination of the intermediate layer representations in the biLM.

ELMo generates the contextualized embedding by concatenating the hidden states (and initial embedding) together followed by weighted summation.

# The Transformer: Going beyond RNNs

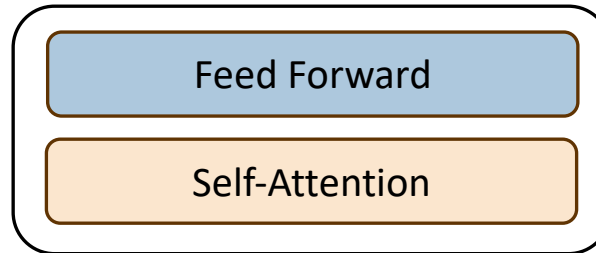
---

The release of the **Transformer** along with the results it achieved on tasks like machine translation led researchers to view it as a superior alternative to RNNs, given its improved handling long-term dependencies compared to RNNs.



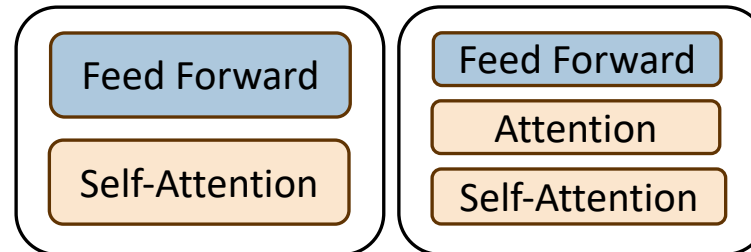
# Three Types of Pretraining Ways

Encoders



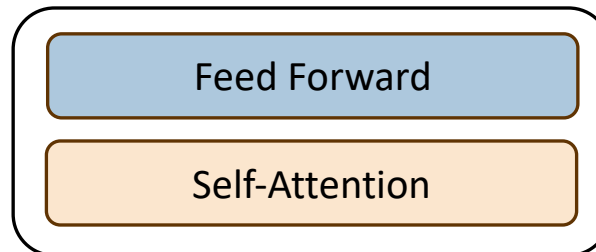
- Bidirectional – capable of looking into the future
- Suitable for downstream task

Encoder-Decoders



- Pros of decoders and encoders?
- Cons of having both during pretraining

Decoders

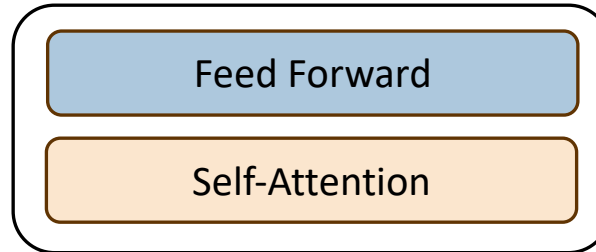


- What LMs we have seen so far
- Suitable for generation task
- Not bidirectional



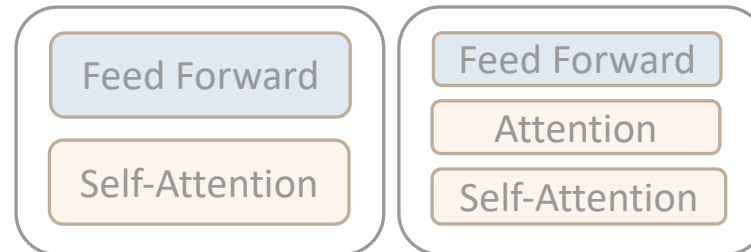
# Three Types of Pretraining Ways

## Encoders



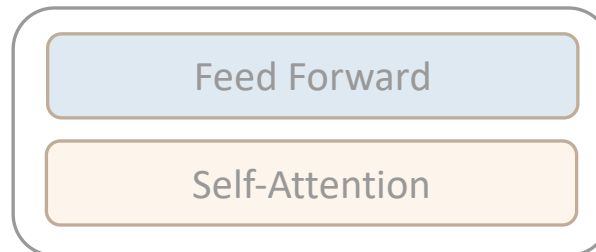
- Bidirectional – capable of looking into the future
- Suitable for downstream task

## Encoder-Decoders



- Pros of decoders and encoders?
- Cons of having both during pretraining

## Decoders



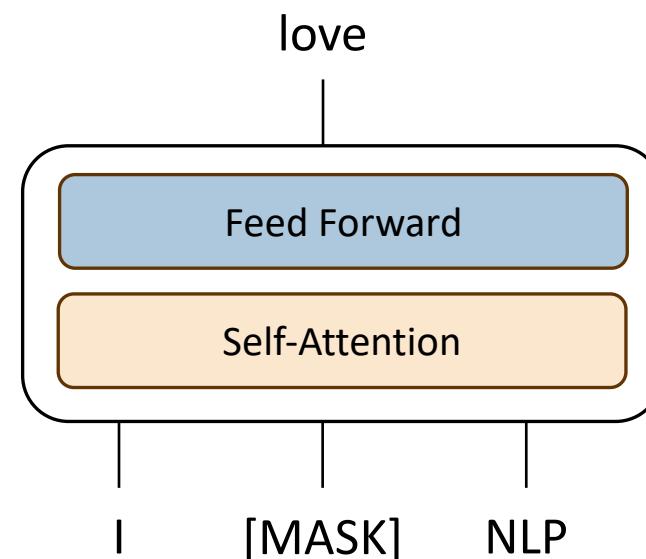
- What LMs we have seen so far
- Suitable for generation task
- Not bidirectional

# Pretraining Encoder:

---

Idea: Encoder works bidirectionally which is more powerful than either a left-to-right model or the shallow concatenation of a left-to-right and a right-to-left model.

- Pretraining task
  - #1 MLM (Masked Language Models):
    - Train the model to reconstruct the sentence by replacing some fraction of words in the input with a special token, [MASK]
    - The model will look at the sentence bidirectionally to try to predict the removed words.
  - #2 NSP (Next Sentence Prediction):
    - In order to train a model that understands sentence relationships, we pre-train for a binarized next sentence prediction task that can be trivially generated from any monolingual corpus



# BERT: Bidirectional Encoder Representations from Transformers

(Devlin et al., 2018)

---

## Task 1. “Masked Language Models”

The objective is proposed alongside the weights of a pretrained Transformer, a model they refer to as BERT.

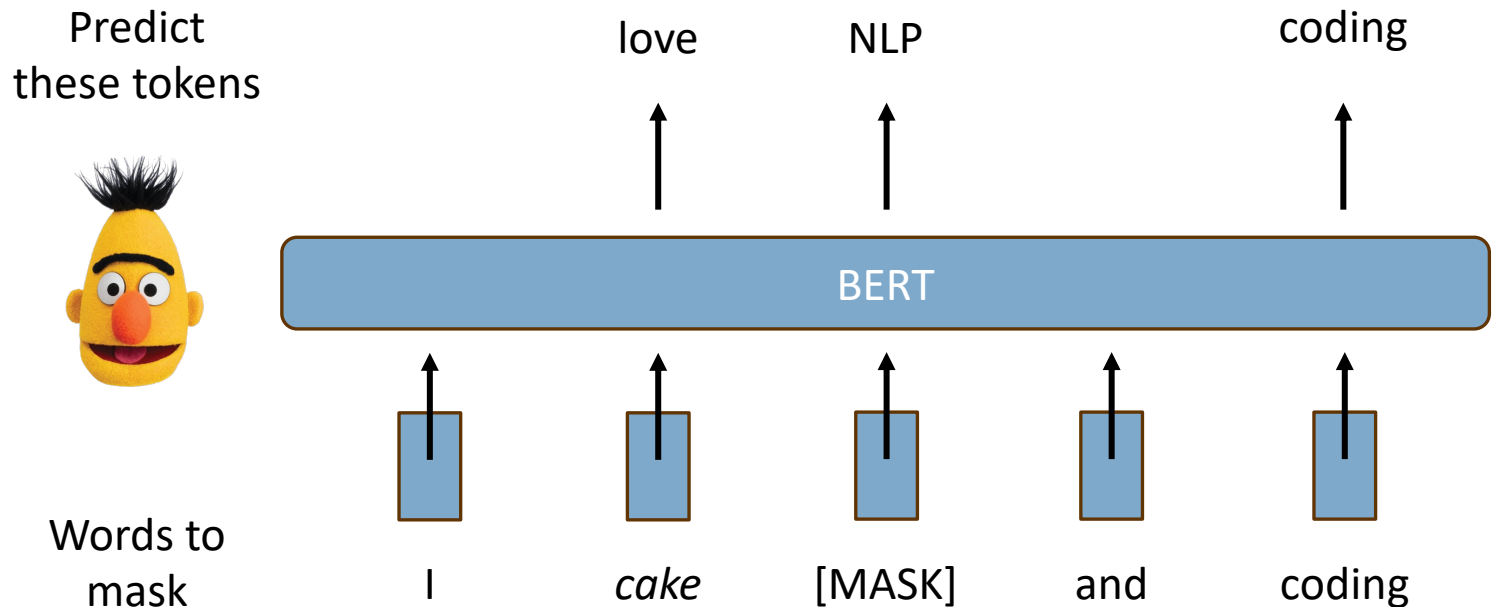
Replace the **15% of the tokens** with

- 1) [MASK] 80% of the time
- 2) Random 10% of the time
- 3) Unchanged 10 % of the time

Predict these tokens

### Why?

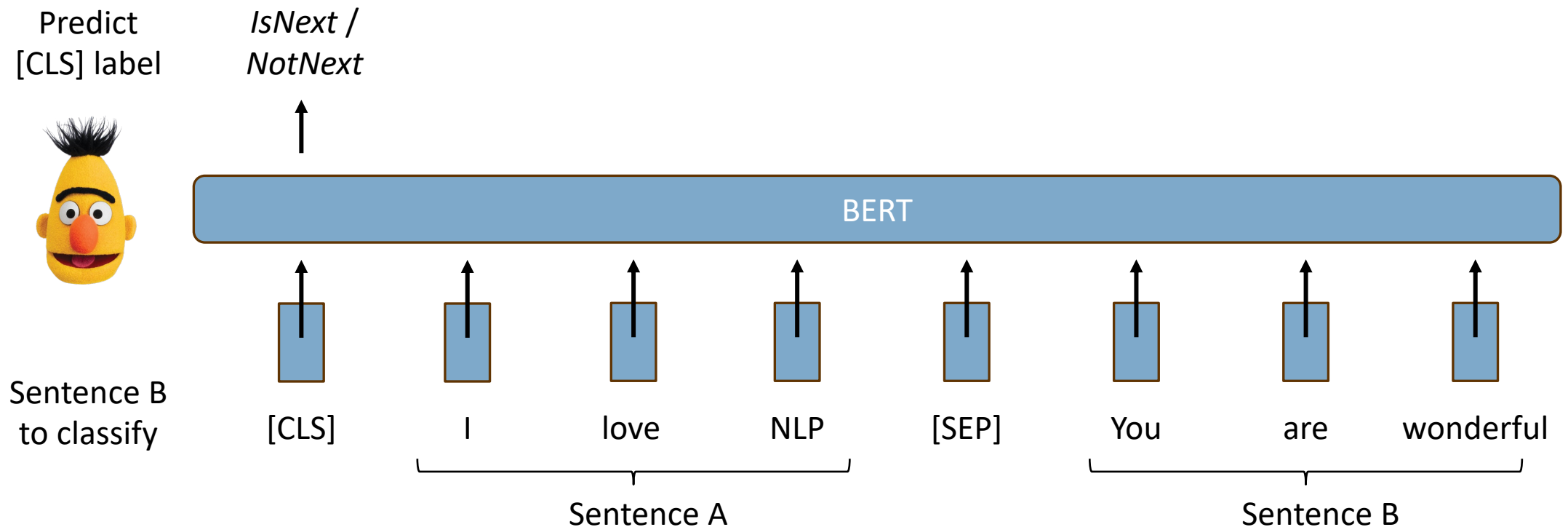
Prevent the model from becoming complacent and failing to construct robust representations for unmasked words. (No masks are seen at fine-tuning time!)



# BERT (Devlin et al., 2018)

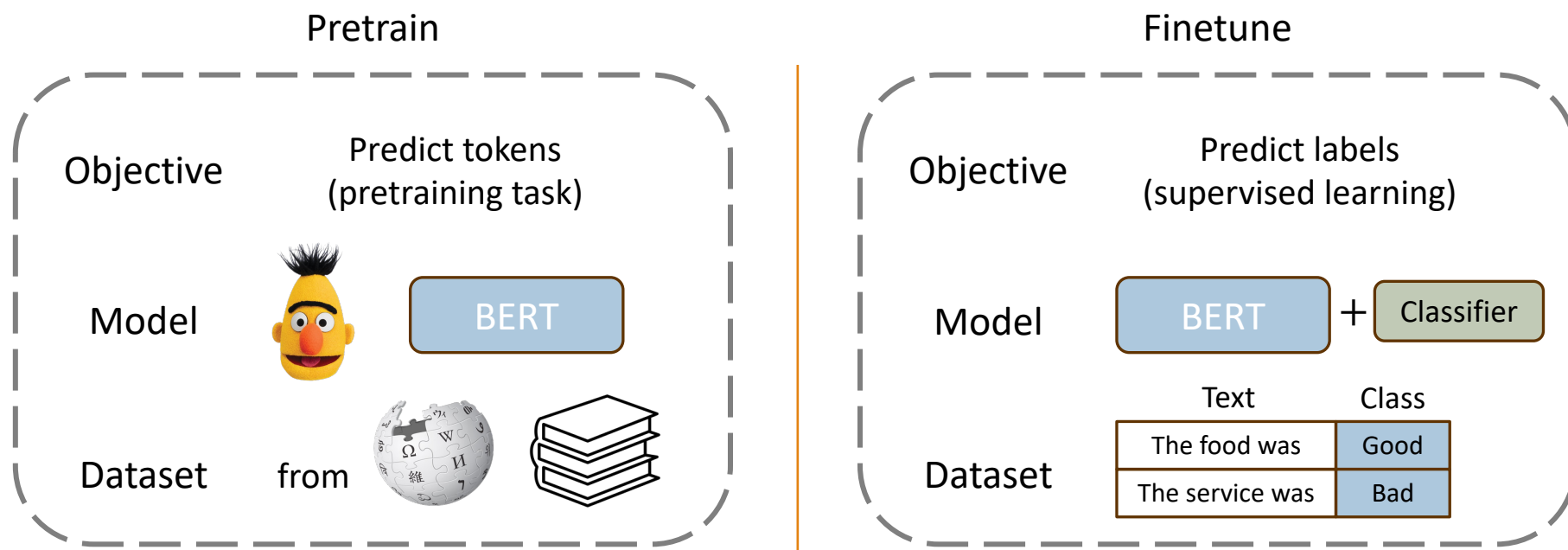
## Task 2. “Next Sentence Prediction”

To develop a model that understands sentence relationships, we pre-train for a binarized next sentence prediction task that can be easily generated from any monolingual corpus.



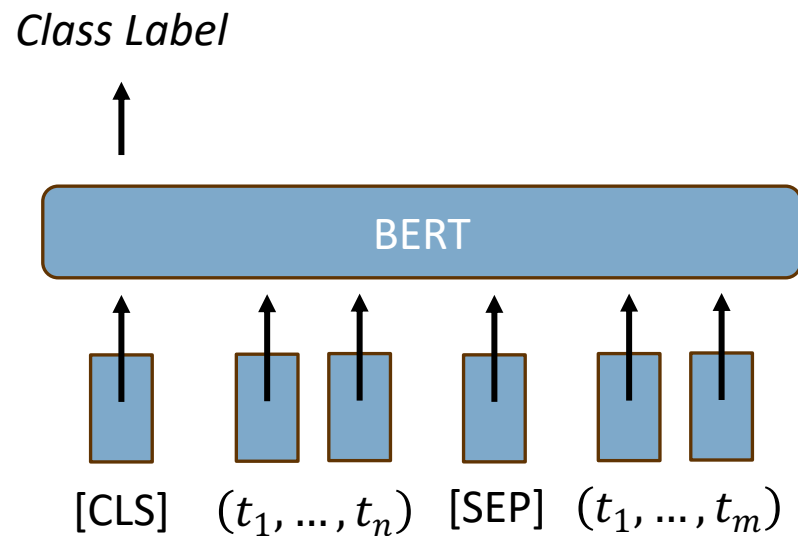
# Transfer Learning: Pretrain and Finetune

Modern-day standard approach is to **first pre-train then fine-tune**. During pretraining, a large amount of unlabeled text is used to build a general model of language understanding before fine-tuned on various specific NLP tasks

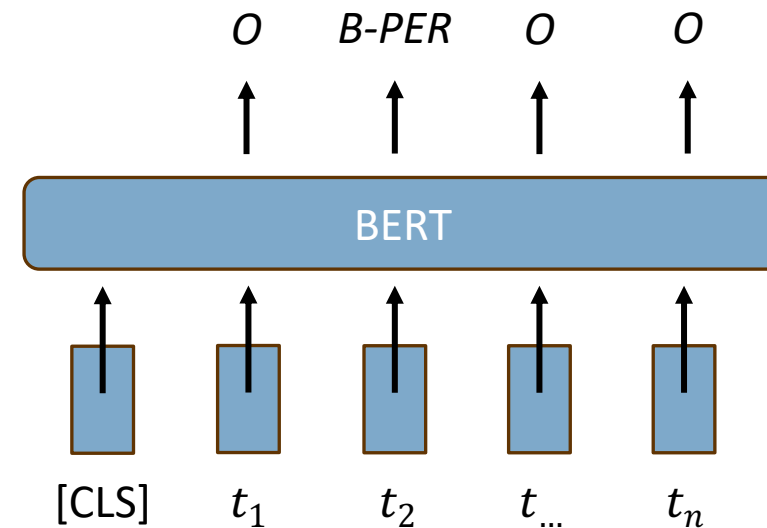


# Finetune BERT

Task-specific models are used to minimize the number of parameters need to be trained from scratch. This is done by integrating BERT with one additional output layer.



(A) Sentence Pair  
Classification



(B) Single Sentence  
Tagging Task

# Details about BERT

---

Google release two models at first:

- BERT-base (110M parameters): 12 layers, 768-dim hidden size, 12 attention heads.
- BERT-large (340M parameters): 24 layers, 1024-dim hidden size, 16 attention heads.

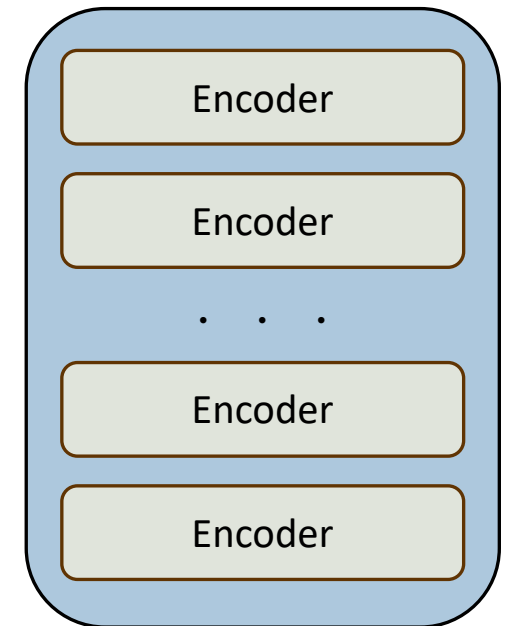
Trained on two datasets:

- BooksCorpus (800M words)
- English Wikipedia (2,500M 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**

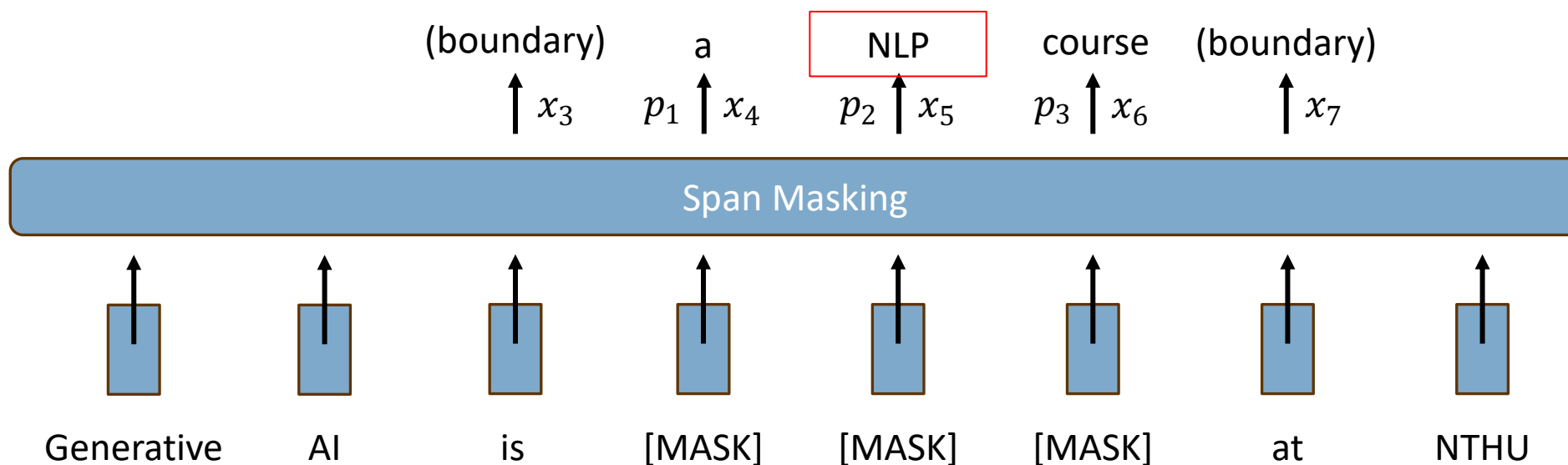
BERT



# Extensions of BERT

Span Masking: Span Boundary Objective (SBO)

$$\mathcal{L}(NLP) = \mathcal{L}_{MLM}(NLP) + \mathcal{L}_{SBO}(NLP) = -\log P(NLP \mid x_5) - \log P(NLP \mid x_3, x_7, p_2)$$





# Extensions of BERT

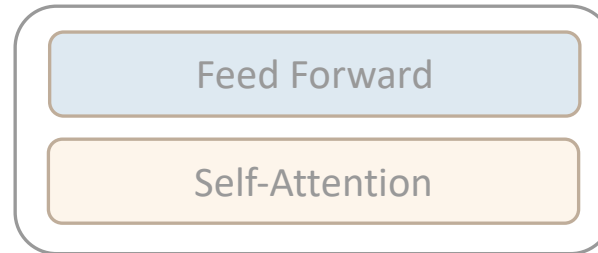
---

Model	MLM task	NSP task	Release
BERT	Static masking	Sentence relationship	2018/10
RoBERTa <a href="https://arxiv.org/pdf/1907.11692.pdf">https://arxiv.org/pdf/1907.11692.pdf</a>	Dynamic masking	Remove NSP task	2019/07
SpanBERT	Span masking	Remove NSP task	2019/07
Albert	N-gram masking	Sentence order	2019/09



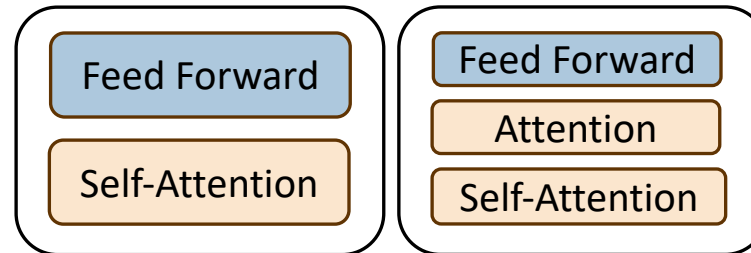
# Three Types of Pretraining Ways

Encoders



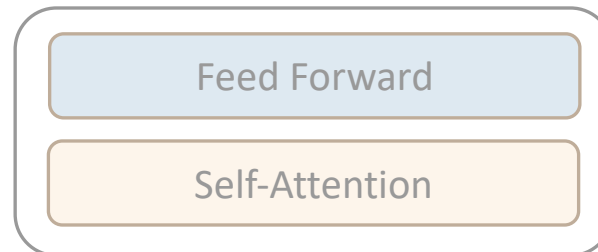
- Bidirectional – capable of looking into the future
- Suitable for downstream task

Encoder-Decoders



- Pros of decoders and encoders?
- Cons of having both during pretraining

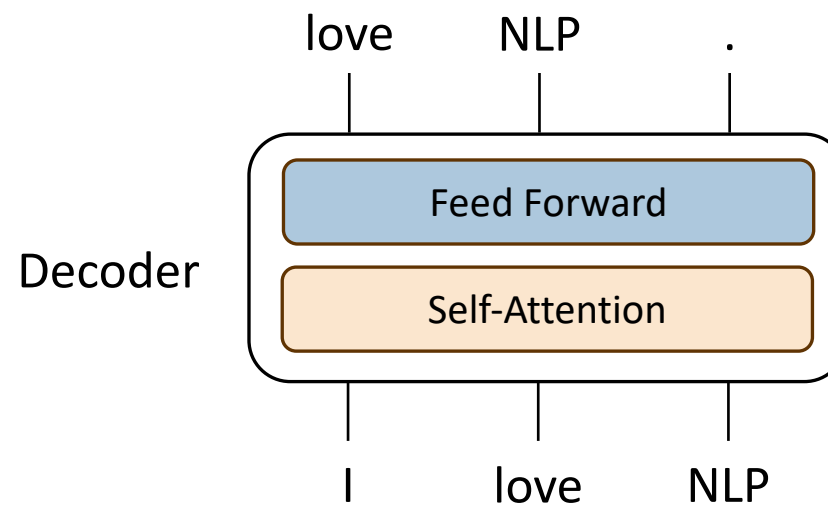
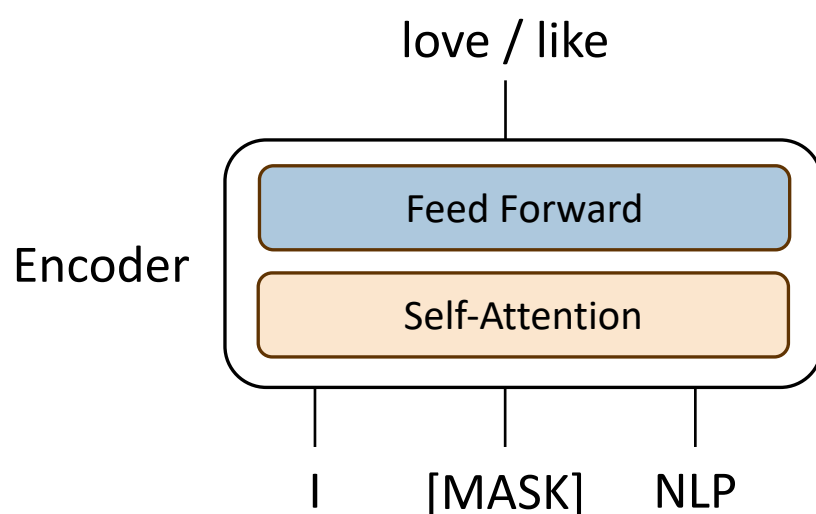
Decoders



- What LMs we have seen so far
- Suitable for generation task
- Not bidirectional

# Limitation of Encoders

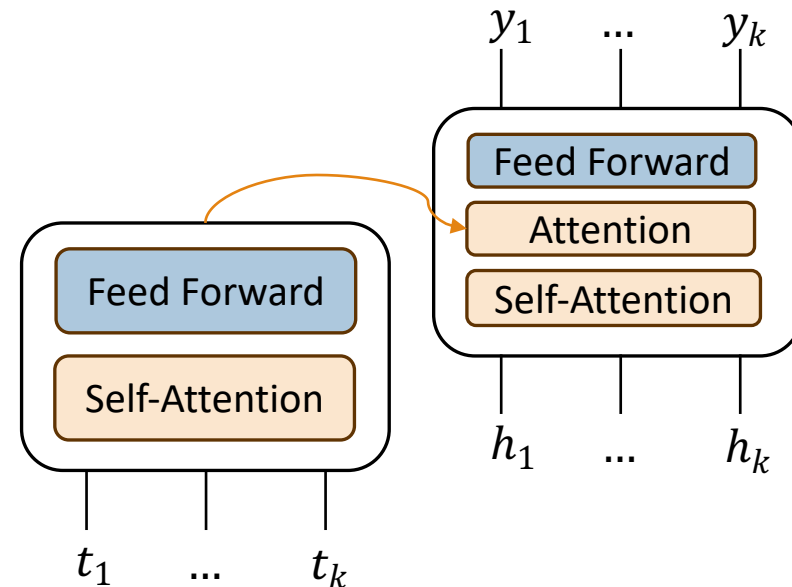
Despite its strong performance on various tasks, it struggles to perform well on generative task due to the lack of bidirectional context, therefore, when your task involves autoregressive (1-word-at-a-time), consider using a pretrained decoder.



# Pretraining Encoder-Decoder

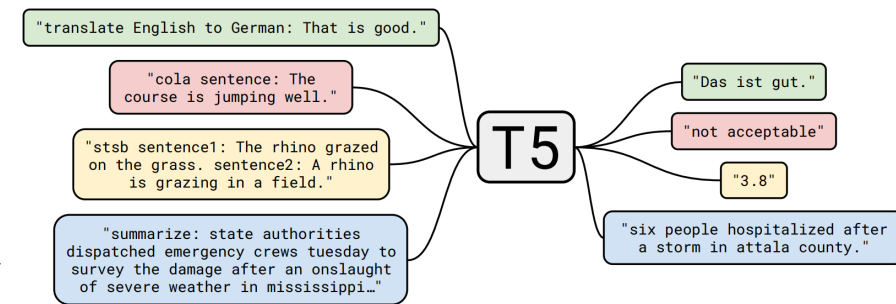
Idea: Capture the context bidirectionally with encoder then use decoder to train the whole model through language modeling.

- Pretraining task
  - **High-level approaches**
    - Language modeling
    - BERT-style
    - Deshuffling
  - **Corruption Strategy**
    - Replace by mask
    - Replace by spans
    - Drop the token



# T5 (Text-to-Text Transfer Transformer)

(Raffel et al., 2019) <https://arxiv.org/pdf/1910.10683.pdf>



Idea: Replace various spans with different length from the input with unique placeholders; then, reconstruct the removed spans during decoding.

Original text: *Thank you for inviting me to your party last week.*

BERT-style masking

Targets Thank you for inviting me to your party last week

T5

Inputs Thank you <M> <M> me to your party *apple* week

Replace spans

Target <X> for inviting <Y> last <Z>

T5

Input Thank you <X> me to your party <Y> week

# T5 (Raffel et al., 2018)

---

Later, BERT-style objective, a method originally proposed as a pre-training technique for an encoder-only model trained for classification and span prediction was found to have better performance over the alternative approach.

## BERT-style masking

Targets Thank you for inviting me to your party last week

T5

Inputs Thank you <M> <M> me to your party *apple* week

## Replace spans

Target <X> for inviting <Y> last <Z>

T5

Input Thank you <X> me to your party <Y> week

# T5 (Raffel et al., 2018)

---

“Span” refers to multiple consecutive tokens that have been corrupted. A single unique mask token is then used to replace the entire span, resulting in unlabeled text data being processed into shorter sequences.

## BERT-style masking

Targets Thank you for inviting me to your party last week

T5

Inputs Thank you <M> <M> me to your party *apple* week

## Replace spans

Target <X> for inviting <Y> last <Z>

T5

Input Thank you <X> me to your party <Y> week

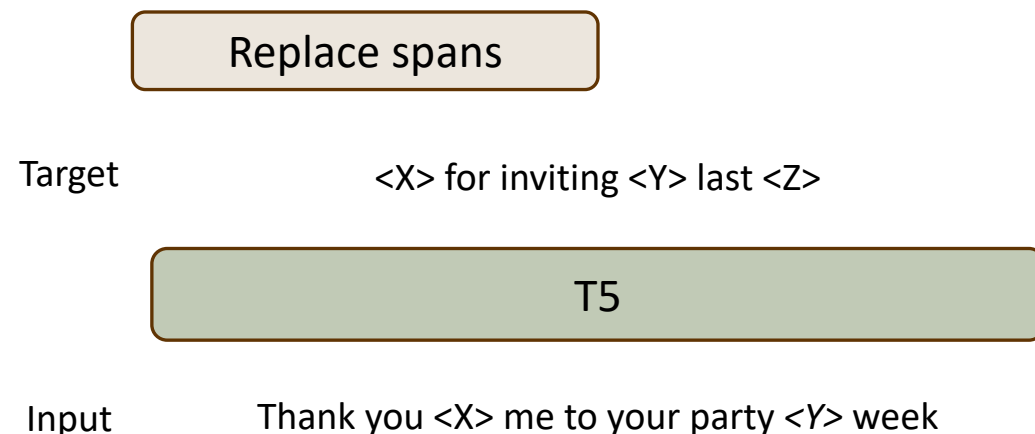
# Details about Replace Spans

Our objective of replacing spans can be parameterized by the proportion of tokens to be corrupted and the total number of corrupted. It is also shown that this setting is not sensitive to performance due to slight difference in performance with low corruption rate (below 50%).

- Corruption objective
  - Corruption rate : 15%
  - Corrupted average span length: 3

Both parameters can be used to specified total number of span by using this simple calculation.

$$avg\_span\_length = \frac{T_{token} \times rate_{corrupt}}{\#_{span}}$$





# Details about T5

---

There are five variants for T5 models:

- T5-small (60M parameters): 6 layers, 512-dim hidden size, 8 attention heads.
- T5-base (220M parameters): 12 layers, 768-dim hidden size, 12 attention heads.
- T5-large (770M parameters): 24 layers, 1024-dim hidden size, 16 attention heads.
- T5-3B: 24 layers, 1024-dim hidden size, 32 attention heads.
- T5-11B: 24 layers, 1024-dim hidden size, 128 attention heads.

Trained on:

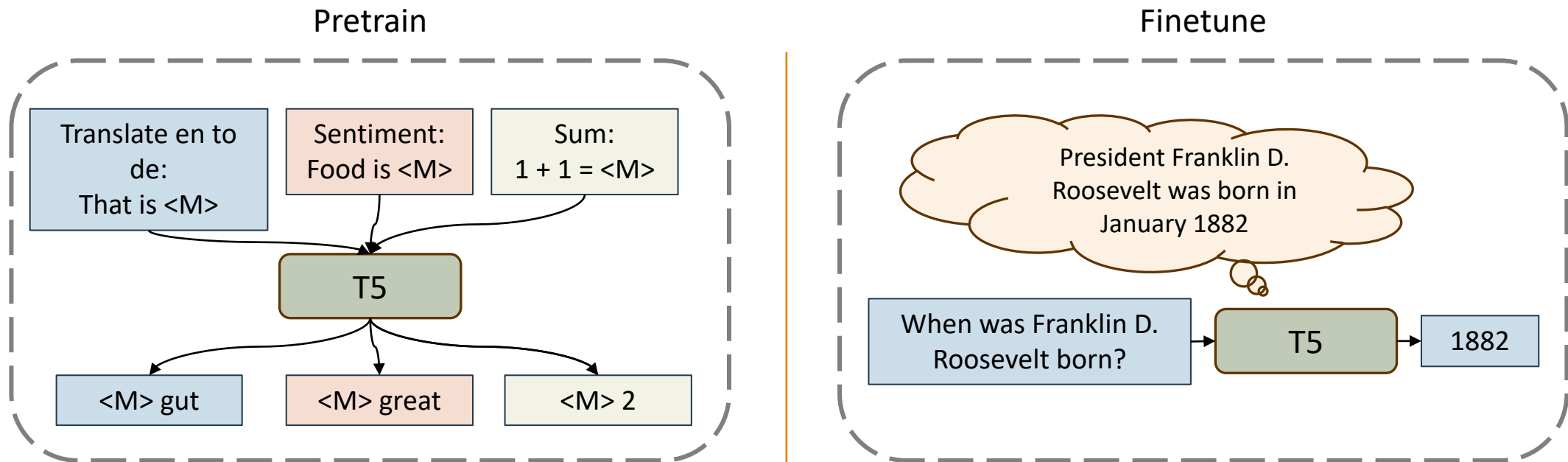
- C4 (34B words)

**Scaling** can make a significant boost in performance.

Converts all text-based language problems into a **text-to-text format**.

# Training Strategy and Finetune T5

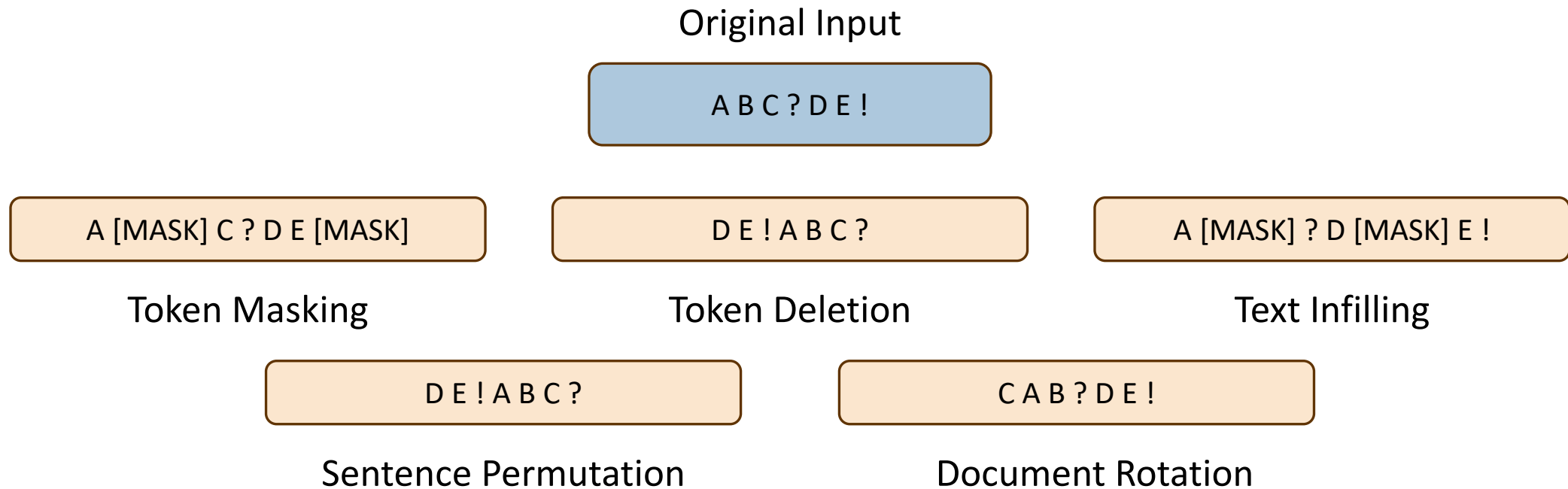
We pre-train T5 by having it learned to fill in dropped-out spans of text (denoted by <M>). To apply T5 to closed-book question answer, we fine-tuned it to answer questions without giving any additional input information or context. This forces T5 to answer questions based on “knowledge” that it learnt during pre-training.



# Extensions of T5

---

BART serves as a **denoising autoencoder** designed to pretrain **sequence-to-sequence** models. BART undergoes two training processes. First, corrupting text using arbitrary noising function. Second, learning a model to reconstruct the original text.



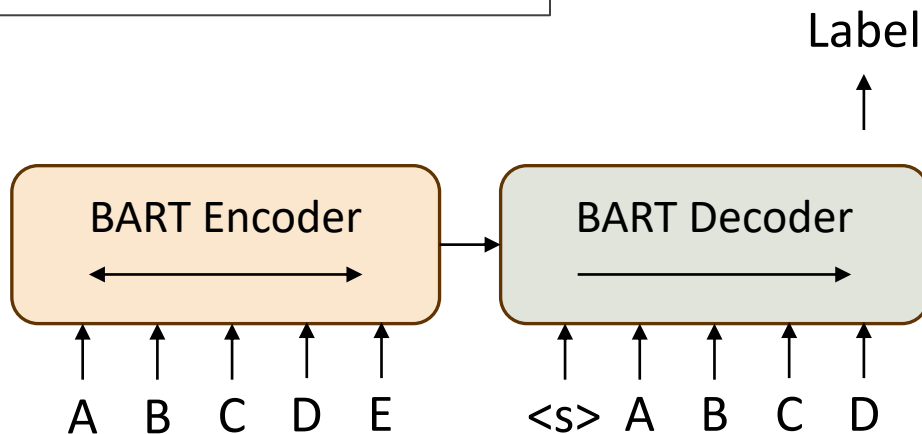
# Pretraining tasks of BART

Noise	Method	Goal
Token Masking	Random tokens are sampled and replaced with [MASK] elements.	Teach the model to predict masked token.
Token Deletion	Random tokens are deleted from the input.	Teach the model to predict masked token and its position.
Text Infilling	Like SpanBERT, but 0-length spans correspond to the insertion of [MASK] tokens.	Teach the model to predict how many tokens are missing from a span.
Sentence Permutation	A document is divided into sentences based on full stops, and these sentences are shuffled in a random order.	Teach the model to clarify the relationship between two sentences.
Document Rotation	A token is chosen uniformly at random, and the document is rotated so that it begins with that token.	Teach the model to identify the start of the document.

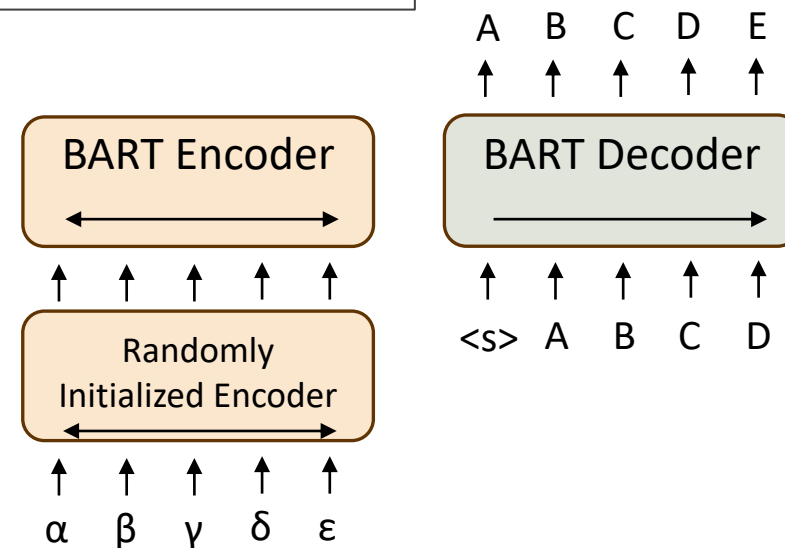
# Finetuning BART

For **classification** problems, both the encoder and decoder receive the same input, and the representation from the final output is utilized; For **generation**, BART's encoder is replaced with newly trained encoder the word embeddings in BART. The new encoder can use a separate vocabulary from original BART model.

## Classification Problems

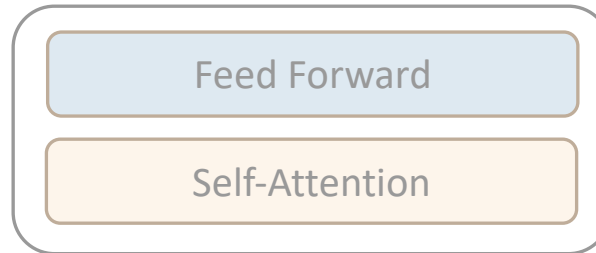


## Generation Problems



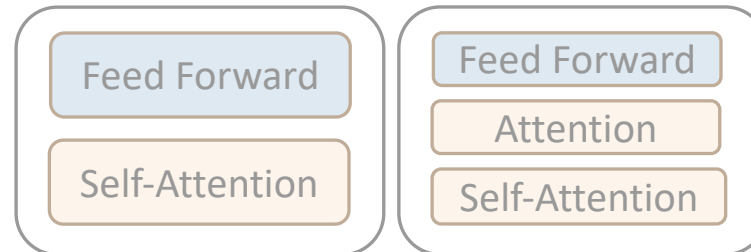
# Three Types of Pretraining Ways

Encoders



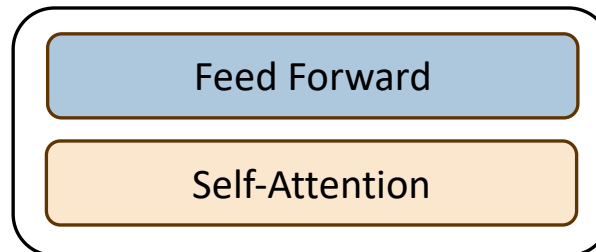
- Bidirectional – capable of looking into the future
- Suitable for downstream task

Encoder-Decoders



- Pros of decoders and encoders?
- Cons of having both during pretraining

Decoders



- What LMs we have seen so far
- Suitable for generation task
- Not bidirectional

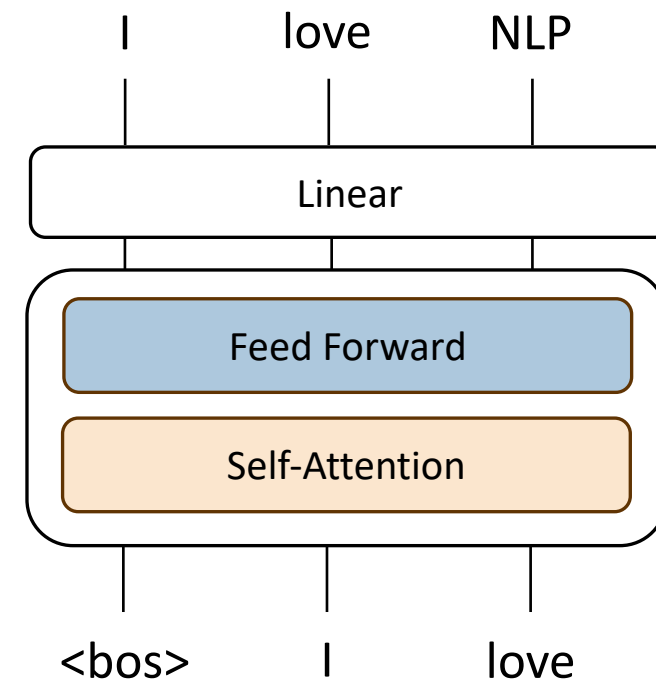
# Pretraining Decoder

Idea: Pretrain decoders on language modeling, then use them as generators. This is helpful in tasks where the output is a sequence with a vocabulary like that at pretraining time.

- Pretraining task
  - Next Token Prediction

$$p(x) = \prod_{i=1}^n p(s_n | s_1, \dots, s_{n-1})$$

Since language follows a specific rule and order, therefore, we leverage conditional probabilities to help us determine the most probable next word based on the preceding sequence of texts.



# GPT (Radford et al., 2018)

---

## Details

- Transformer decoder with 12 layers, 117M parameters
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers
- Byte-pair encoding with 40,000 merges

## Dataset

- Trained on BooksCorpus: over 7000 unique books
  - Contains long spans of contiguous text, for learning long-distance dependencies

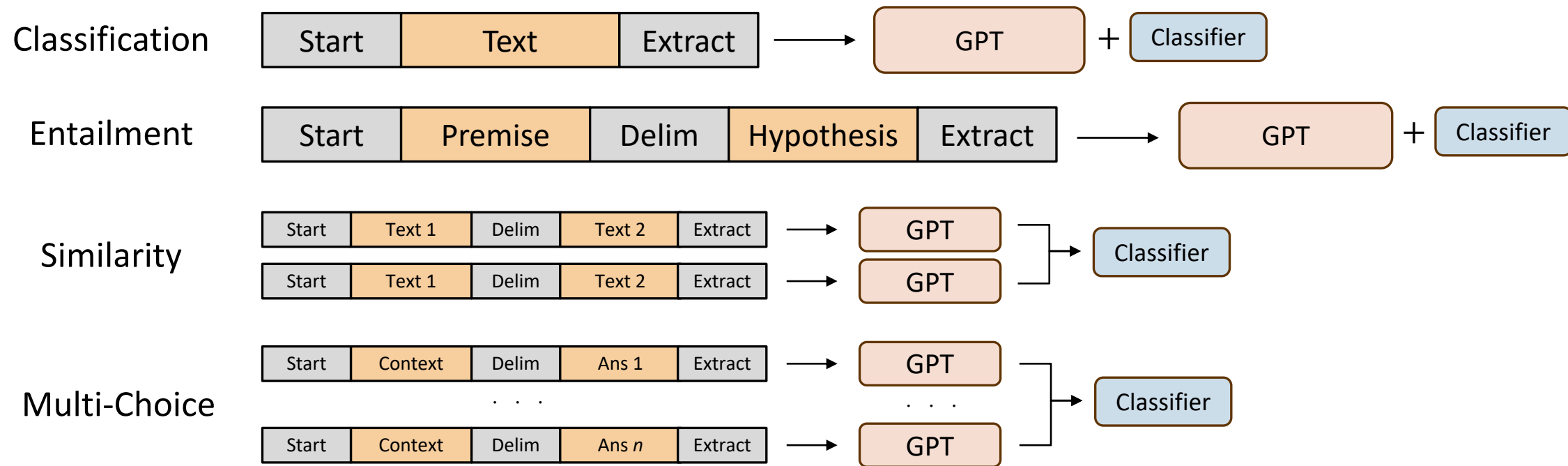
## Additional Insights

- The acronym “GPT” never showed up in the original paper
  - Could stand for “Generative PreTraining” or “Generative Pretrained Transformer”



# Finetune GPT

Idea: Convert all structured inputs into sequence of tokens to be processed by our pre-trained model, followed by a linear classifier



# GPT Family

---

Model	Pretraining method	Parameters	Dataset	Release
GPT-1	Transformer decoder followed by linear-softmax	117 M	BooksCorpus 4.5 GB	2018.06.11
GPT-2	Same as GPT-1 but different normalized layer	1.5 B	WebText 40 GB	2019.02.14
GPT-3	Larger version of GPT-2	175 B	CommonCrawl 45 TB	2020.05.15
GPT-4	Trained with RLHF	> 1.5 T	Undisclosed	2023.03.14

# GPT-3, In-Context learning, and LLMs

---

GPT-3 (Brown et al., 2020)

So far, there are two ways we interact with pre - trained models

1. Sample from the distributions they define
2. Fine-tune them on a task we care about

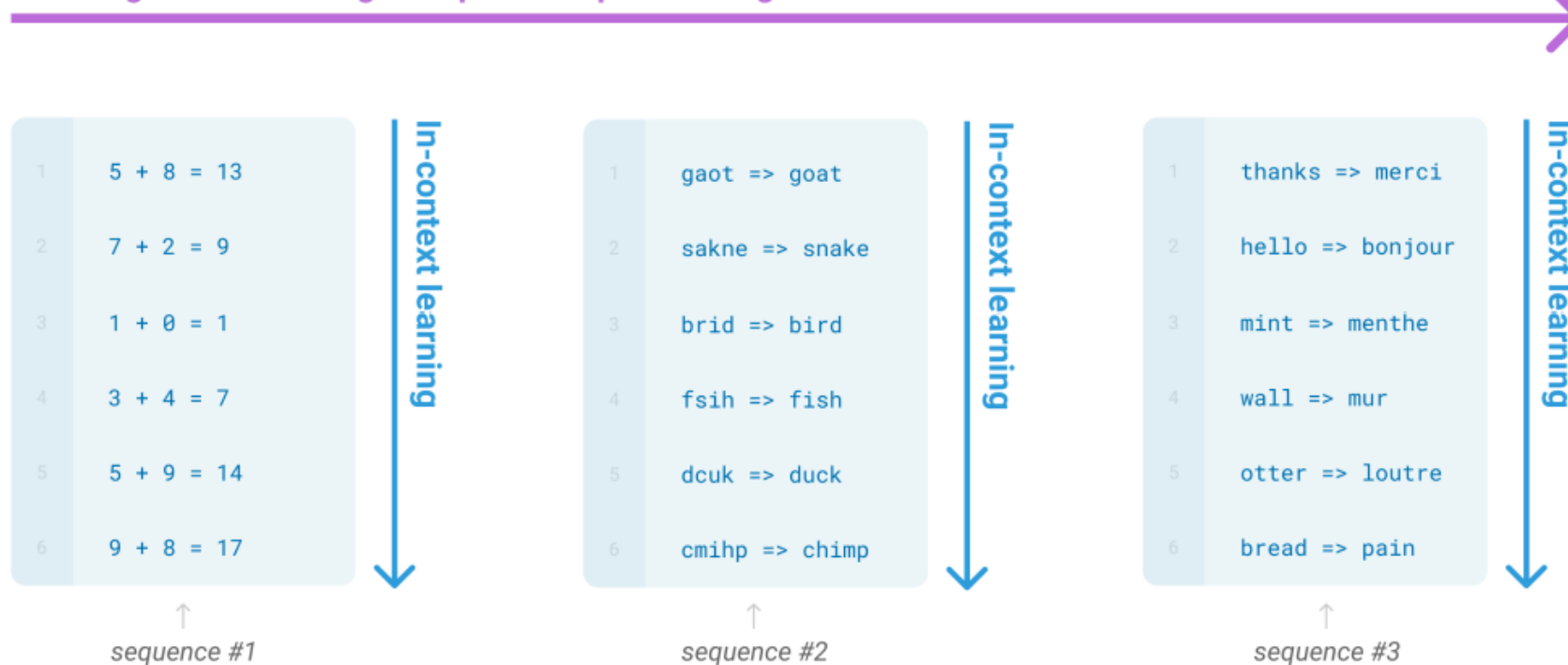
However, large models are shown to have this **emergence ability** which allows them to perform some kind of learning without gradient steps.

GPT-3, with 175 billion parameters, is a good example of this.

# GPT-3, In-Context learning, and LLMs

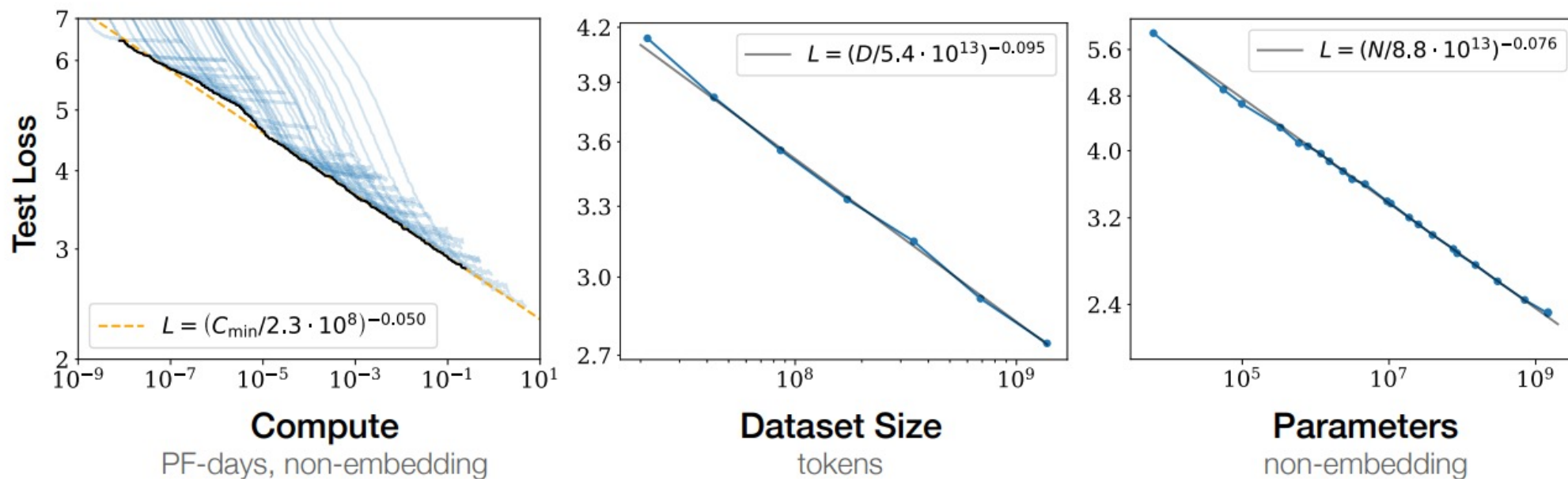
Very large language models have ability to **learn through examples** provided within their contexts **without gradient steps**. This is referred to as in-context learning ability.

Learning via SGD during unsupervised pre-training



# Scaling Laws (Kaplan et al., 2020)

Scaling laws are simple, predictive rules for model performance. Increasing the model size, dataset size, and the amount of computation during training lead to better performance which increase smoothly.



# Very Large Language Models

---

Model	Parameters	Publisher	Release
GPT-3	175B	OpenAI	2020.05
BLOOM	176B	BigScience	2022.07
Flan-PaLM	540B	Google	2022.10

# Large Language Models for Community

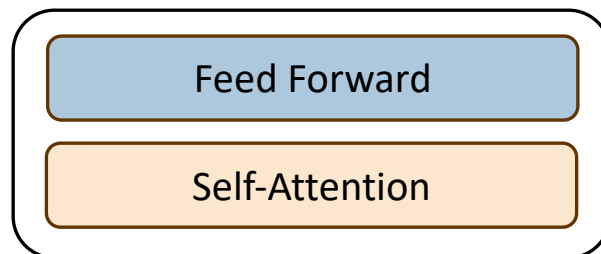
---

Model	Architecture	Parameters	Publisher	Release
Llama 2	Decoder-only	7B, 13B, 70B	Meta	2023.07.18
Mistral	Decoder-only	7B, 8×7B	MistralAI	2023.09.27
Phi 2	Decoder-only	2.7B	Microsoft	2023.12.12
Gemma	Decoder-only	2B, 7B	Google	2024.02.21

# Takeaways

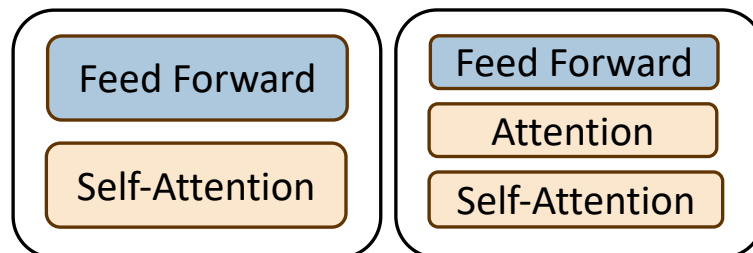
---

## Encoders



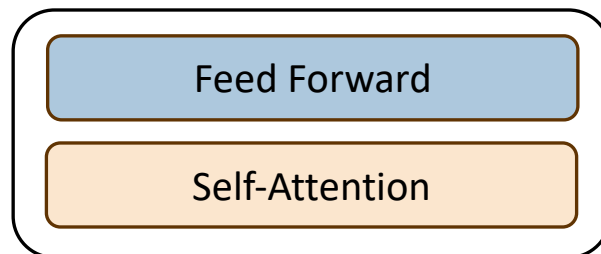
- BERT family
- MLM and NSP task for pretraining
- Not good at performing generation

## Encoder-Decoders



- T5, BART
- Span corruption for boosting MLM
- Targeting text-to-text format model

## Decoders



- GPT family
- Pure language modeling
- Most popular backbone nowadays