



Transformer Language Models (Continued)

Slido: <https://app.sli.do/event/jbsxB9wcw7PojJ4ur9nM83>

Yu Meng
University of Virginia
yumeng5@virginia.edu

Oct 15, 2025

Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling & Recurrent Neural Networks (RNNs)
- Week 6: Language Modeling with Transformers
- **Week 8: Transformer and Pre-Training**
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

Reminder

- Guest lecture postponed (date to be announced later)
- Midterm report due next Monday (10/20) 11:59pm (guideline on Canvas)

(Recap) Transformer: Overview

- Transformer is a specific kind of sequence modeling architecture (based on DNNs)
 - Use attention to replace recurrent operations in RNNs
 - The most important architecture for language modeling (almost all LLMs are based on Transformers)!
-

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

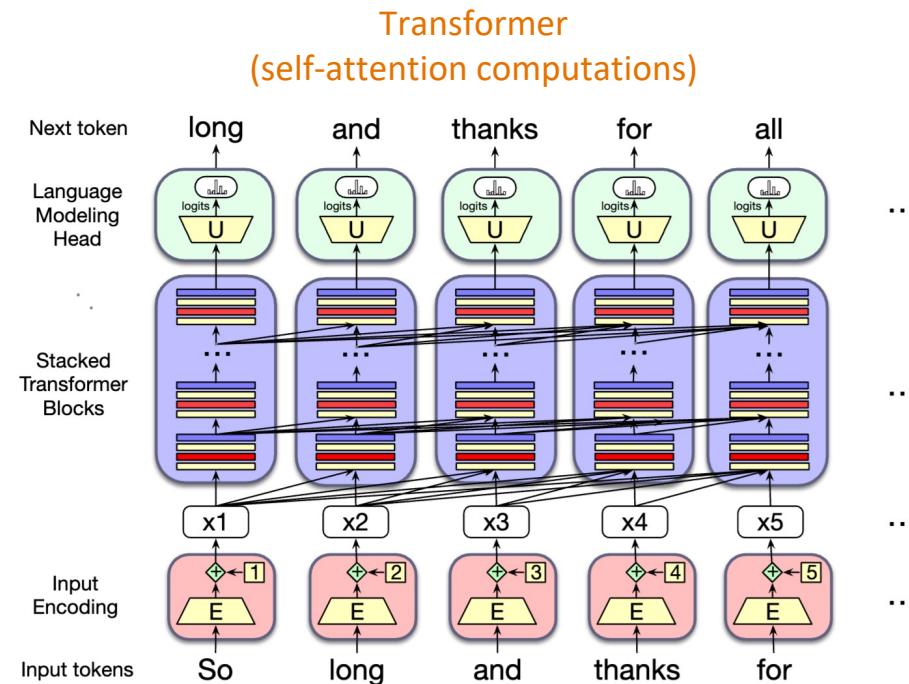
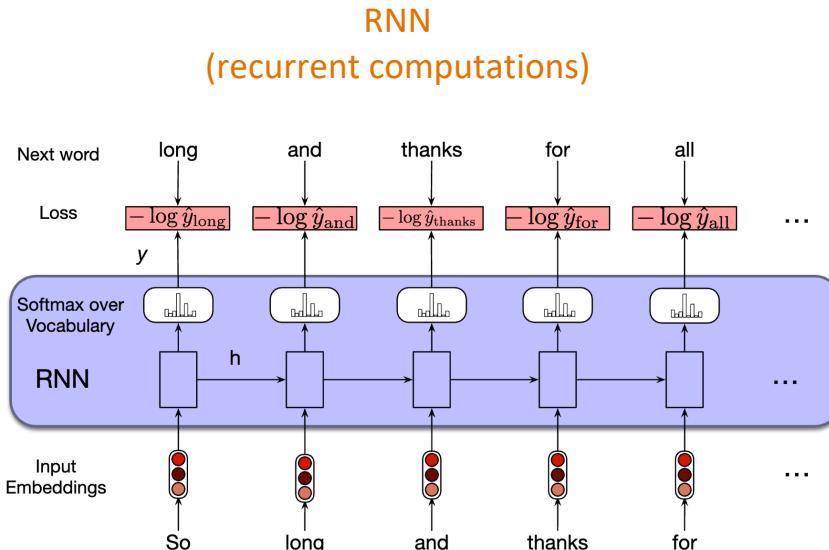
Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

(Recap) Transformer vs. RNN



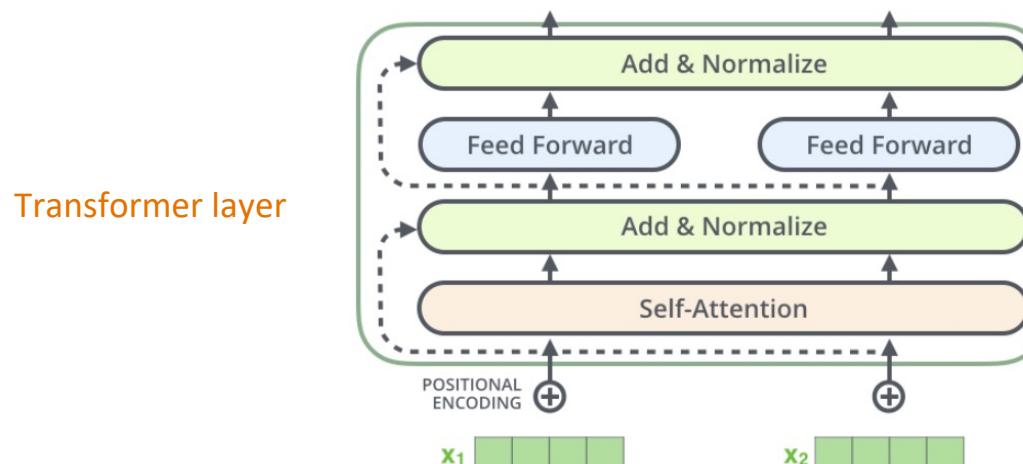
(Recap) Transformer: Motivation

- Parallel token processing
 - RNN: process one token at a time (computation for each token depends on previous ones)
 - Transformer: process all tokens in a sequence in parallel
- Long-term dependencies
 - RNN: bad at capturing distant relating tokens (vanishing gradients)
 - Transformer: directly access any token in the sequence, regardless of its position
- Bidirectionality
 - RNN: can only model sequences in one direction
 - Transformer: inherently allow bidirectional sequence modeling via attention

(Recap) Transformer Layer

Each Transformer layer contains the following important components:

- Self-attention
- Feedforward network
- Residual connections + layer norm



(Recap) Self-Attention: Intuition

- Attention: weigh the importance of different words in a sequence when processing a specific word
 - “When I’m looking at this word, which other words should I pay attention to in order to understand it better?”
- **Self-attention:** each word attends to other words in the **same** sequence
- Example: “The quick brown fox jumps over the lazy dog.”
 - Suppose we are learning attention for the word “**jumps**”
 - With self-attention, “**jumps**” can decide which other words in the sentence it should focus on to better understand its meaning
 - Might assign high attention to “fox” (the subject) & “over” (the preposition)
 - Might assign less attention to words like “the” or “lazy”

(Recap) Self-Attention: Example

Derive the center word representation as a weighted sum of context representations!

Center word representation

Context word representation

$$\mathbf{a}_i = \sum_{x_j \in \mathbf{x}} \alpha_{ij} \mathbf{x}_j, \quad \sum_{x_j \in \mathbf{x}} \alpha_{ij} = 1$$

Attention score $i \rightarrow j$, summed to 1

Context word (key)	Center word (query)
The	The
chicken	chicken
didn't	didn't
cross	cross
the	the
road	road
because	because
it	it
was	was
too	too
tired	tired

Current word = "it"

(Recap) Self-Attention: Attention Score Computation

- Attention score is given by the softmax function over vector dot product

$$\mathbf{a}_i = \sum_{x_j \in \mathbf{x}} \alpha_{ij} \mathbf{x}_j, \quad \sum_{x_j \in \mathbf{x}} \alpha_{ij} = 1$$

$$\alpha_{ij} = \text{Softmax}(\mathbf{x}_i \cdot \mathbf{x}_j)$$

Center word (query) representation Context word (key) representation

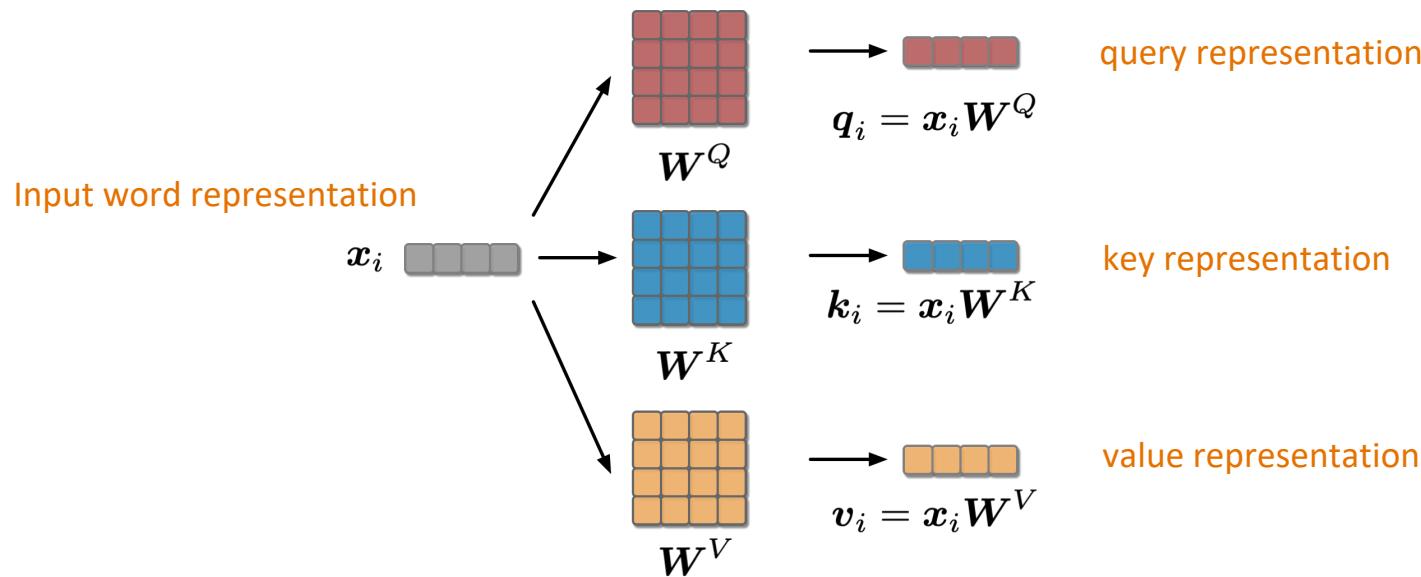
- Why use two copies of word representations for attention computation?
 - We want to reflect the different roles a word plays (as the target word being compared to others, or as the context word being compared to the target word)
 - If using the same copy of representations for attention calculation, a word will (almost) always attend to itself heavily due to high dot product with itself!

(Recap) Self-Attention: Query, Key, and Value

- Each word in self-attention is represented by three different vectors
 - Allow the model to flexibly capture different types of relationships between tokens
- **Query (Q):**
 - Represent the current word seeking information about
- **Key (K):**
 - Represent the reference (context) against which the query is compared
- **Value (V):**
 - Represent the actual content associated with each token to be aggregated as final output

(Recap) Self-Attention: Query, Key, and Value

Each self-attention module has three weight matrices applied to the input word vector to obtain the three copies of representations



(Recap) Self-Attention: Overall Computation

- Input: single word vector of each word \mathbf{x}_i
- Compute Q, K, V representations for each word:

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V$$

- Compute attention scores with Q and K
 - The dot product of two vectors usually has an expected magnitude proportional to \sqrt{d}
 - Divide the attention score by \sqrt{d} to avoid extremely large values in softmax function

$$\alpha_{ij} = \text{Softmax} \left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d}} \right)$$

.....

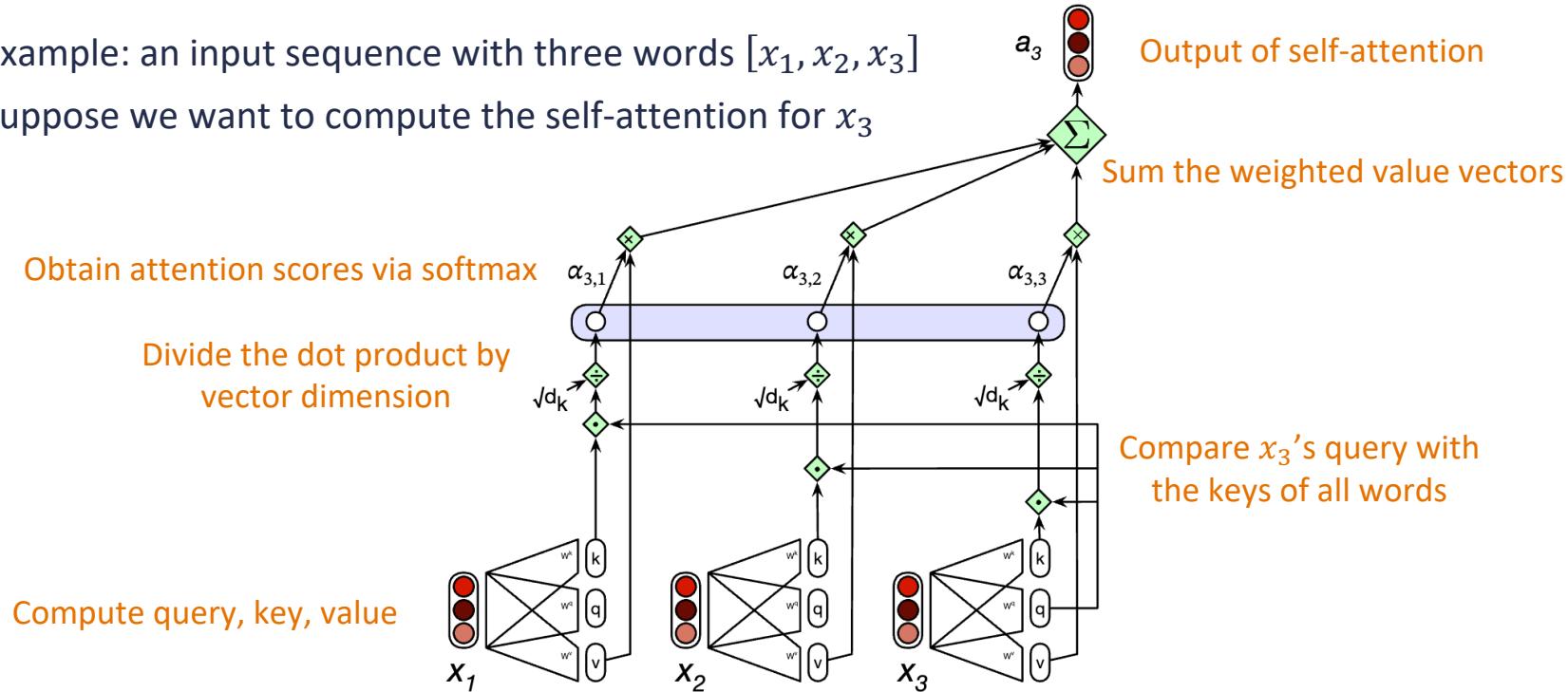
Dimensionality of q and k

- Sum the value vectors weighted by attention scores

$$\mathbf{a}_i = \sum_{x_j \in \mathbf{x}} \alpha_{ij} \mathbf{v}_j$$

(Recap) Self-Attention: Illustration

- Example: an input sequence with three words $[x_1, x_2, x_3]$
- Suppose we want to compute the self-attention for x_3



(Recap) Multi-Head Self-Attention

- Transformers use multiple attention heads for each self-attention module
- Intuition:
 - Each head might attend to the context for different purposes (e.g., particular kinds of patterns in the context)
 - Heads might be specialized to represent different linguistic relationships

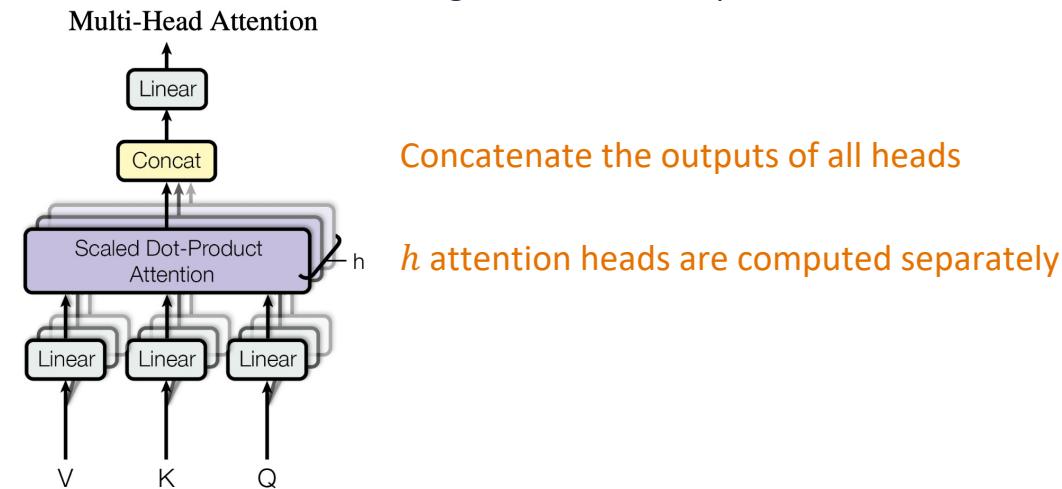


Figure source: <https://arxiv.org/pdf/1706.03762>

(Recap) Parallel Computation of QKV

- Self-attention computation performed for each token is independent of other tokens
- Easily parallelize the entire computation, taking advantage of the efficient matrix multiplication capability of GPUs
- Process an input sequence with N words in parallel

Compute QKV for one word: $\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V \in \mathbb{R}^d$

Stacking N input vectors: $\mathbf{Q} = \mathbf{X} \mathbf{W}^Q \quad \mathbf{K} = \mathbf{X} \mathbf{W}^K \quad \mathbf{V} = \mathbf{X} \mathbf{W}^V \in \mathbb{R}^{N \times d}$

$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}_1 & \text{---} \\ \text{---} & \mathbf{x}_2 & \text{---} \\ \dots & \dots & \dots \\ \text{---} & \mathbf{x}_N & \text{---} \end{bmatrix}$$

(Recap) Parallel Computation of Attention

Attention computation can also be written in matrix form

Compute attention for one word: $a_i = \text{Softmax} \left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d}} \right) \cdot \mathbf{v}_j$

Compute attention for one N words: $\mathbf{A} = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V}$

Attention matrix

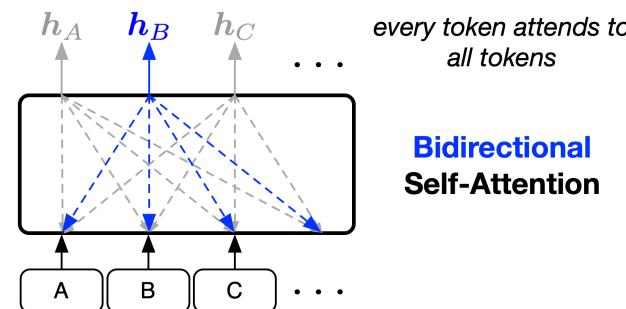
q1·k1	q1·k2	q1·k3	q1·k4
q2·k1	q2·k2	q2·k3	q2·k4
q3·k1	q3·k2	q3·k3	q3·k4
q4·k1	q4·k2	q4·k3	q4·k4

N

N

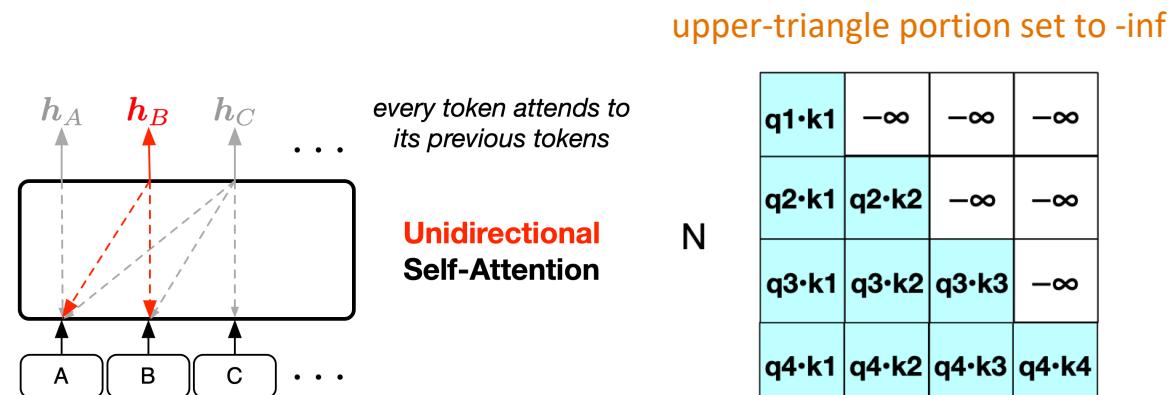
(Recap) Bidirectional vs. Unidirectional Self-Attention

- Self-attention can capture different context dependencies
- **Bidirectional** self-attention:
 - Each position can attend to all other positions in the input sequence
 - Transformers with bidirectional self-attention are called Transformer **encoders** (e.g., BERT)
 - Use case: natural language understanding (NLU) where the entire input is available at once, such as text classification & named entity recognition



(Recap) Bidirectional vs. Unidirectional Self-Attention

- Self-attention can capture different context dependencies
- **Unidirectional (or causal)** self-attention:
 - Each position can only attend to earlier positions in the sequence (including itself).
 - Transformers with unidirectional self-attention are called Transformer **decoders** (e.g., GPT)
 - Use case: natural language generation (NLG) where the model generates output sequentially



(Recap) Position Encoding

- Motivation: inject positional information to input vectors

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V \in \mathbb{R}^d$$

$$\mathbf{a}_i = \text{Softmax} \left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d}} \right) \cdot \mathbf{v}_j \quad \text{When } \mathbf{x} \text{ is word embedding, } \mathbf{q} \text{ and } \mathbf{k} \text{ do not have positional information!}$$

- How to know the word positions in the sequence? Use position encoding!

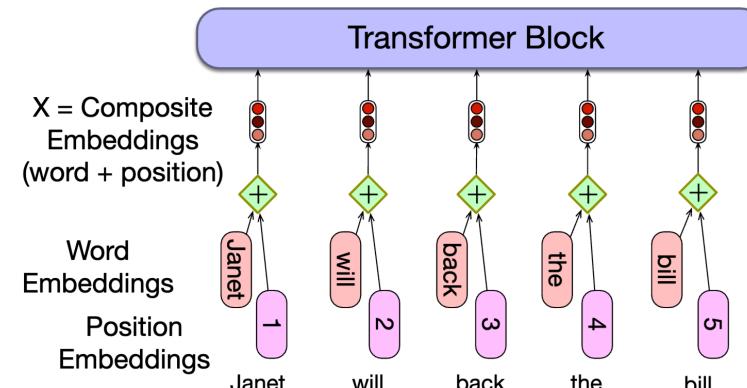


Figure source: <https://web.stanford.edu/~jurafsky/slp3/8.pdf>

(Recap) Position Encoding Methods

- Absolute position encoding (the original Transformer paper)
 - Learn position embeddings for each position
 - Not generalize well to sequences longer than those seen in training
- Relative position encoding ([Self-Attention with Relative Position Representations](#))
 - Encode the relative distance between words rather than their absolute positions
 - Generalize better to sequences of different lengths
- Rotary position embedding ([RoFormer: Enhanced Transformer with Rotary Position Embedding](#))
 - Apply a rotation matrix to the word embeddings based on their positions
 - Incorporate both absolute and relative positions
 - Generalize effectively to longer sequences
 - Widely-used in latest LLMs

Agenda

- Tokenization
- Other Transformer Modules
- Transformer Language Model Pretraining

Tokenization: Overview

- Tokenization: splitting a string into a sequence of tokens
- Simple approach: use whitespaces to segment the sequence
 - One token = one word
 - We have been using “tokens” and “words” interchangeably
- However, segmentation using whitespaces is not the approach used in modern large language models

Multiple models, each with different capabilities and price points. Prices can be viewed in units of either per 1M or 1K tokens. You can think of tokens as pieces of words, where 1,000 tokens is about 750 words.

Limitation of Word-Based Segmentation

- Out-of-vocabulary (OOV) issues:
 - Cannot handle words never seen in our training data
 - Reserving an [UNK] token for unseen words is a remedy
- Subword information:
 - Loses subword information valuable for understanding word meaning and structure
 - Example: “unhappiness” -> “un” + “happy” + “ness”
- Data sparsity and exploded vocabulary size:
 - Require a large vocabulary (vocabulary size = number of unique words)
 - The model sees fewer examples of each word (harder to generalize)

Single-Character Segmentation?

- How about segmenting sequences by character?
 - No OOV issue
 - Small vocabulary size
- Increased sequence length:
 - Significantly increases the length of input sequences
 - Transformer's self-attention has quadratic complexity w.r.t. sequence length!
- Loss of word-level semantics:
 - Characters alone often don't carry semantic meaning/linguistic patterns

Subword Tokenization

- Strike a balance between character-level and word-level tokenization
 - Capture meaningful subword semantics
 - Handle out-of-vocabulary words better
 - Efficient sequence modeling
- Three common algorithms:
 - Byte-Pair Encoding (BPE): [Sennrich et al. \(2016\)](#)
 - WordPiece: [Schuster and Nakajima \(2012\)](#)
 - SentencePiece: [Kudo and Richardson \(2018\)](#)
- Subword tokenization usually consists of two parts:
 - A token learner that takes a raw training corpus and induces a **vocabulary** (a set of tokens)
 - A token segmenter that takes a raw sentence and **tokenizes** it according to that vocabulary

Byte-Pair Encoding (BPE) Overview

- BPE is the most commonly used tokenization algorithm in modern LLMs
- Intuition: start with a character-level vocabulary and iteratively merge the most frequent pairs of tokens
- Initialization: let vocabulary be the set of all individual characters: {A, B, C, D, ..., a, b, c, d,}
- Frequency counting: count all adjacent symbol pairs (could be a single character or a previously merged pair) in the training corpus
- Pair merging: merge the most frequent pair of symbols (e.g. 't', 'h' => "th")
- Update corpus: replace all instances of the merged pair in the corpus with the new token & update the frequency of pairs
- Repeat: repeat the process of counting, merging, and updating until a predefined number of merges (or vocabulary size) is reached

BPE: Token Learner

Token learner of BPE

```
function BYTE-PAIR ENCODING(strings  $C$ , number of merges  $k$ ) returns vocab  $V$ 
     $V \leftarrow$  all unique characters in  $C$           # initial set of tokens is characters
    for  $i = 1$  to  $k$  do                      # merge tokens til k times
         $t_L, t_R \leftarrow$  Most frequent pair of adjacent tokens in  $C$ 
         $t_{NEW} \leftarrow t_L + t_R$                   # make new token by concatenating
         $V \leftarrow V + t_{NEW}$                       # update the vocabulary
        Replace each occurrence of  $t_L, t_R$  in  $C$  with  $t_{NEW}$       # and update the corpus
    return  $V$ 
```

BPE Example

Suppose we have the following corpus

low low low low lowest lowest newer newer newer
newer newer newer wider wider wider new new

corpus

5	l	o	w	_			
2	l	o	w	e	s	t	_
6	n	e	w	e	r	_	
3	w	i	d	e	r	_	
2	n	e	w	_			

Special “end-of-word” character
(distinguish between subword units
vs. whole word)

vocabulary

_, d, e, i, l, n, o, r, s, t, w

BPE: Counting & Merging

The adjacent symbol pair with the highest frequency is “er” (count = 9)

low low low low lowest lowest newer newer newer
newer newer newer wider wider wider new new

corpus

5	l	o	w	_			
2	l	o	w	e	s	t	_
6	n	e	w	e r	_		
3	w	i	d	e r	_		
2	n	e	w	_			

Merge “er”



corpus

5	l	o	w	_			
2	l	o	w	e	s	t	_
6	n	e	w	e r	_		
3	w	i	d	e r	_		
2	n	e	w	_			

vocabulary

_, d, e, i, l, n, o, r, s, t, w



vocabulary

_, d, e, i, l, n, o, r, s, t, w, er

BPE: Counting & Merging

The adjacent symbol pair with the highest frequency is “er_” (count = 9)

low low low low lowest lowest newer newer newer
newer newer newer wider wider wider new new

corpus

5	l o w _
2	l o w e s t _
6	n e w er _
3	w i d er _
2	n e w _

Merge “er_”

**corpus**

5	l o w _
2	l o w e s t _
6	n e w er _
3	w i d er _
2	n e w _

vocabulary

_, d, e, i, l, n, o, r, s, t, w, er

**vocabulary**

, d, e, i, l, n, o, r, s, t, w, er, er

BPE: Counting & Merging

The adjacent symbol pair with the highest frequency is “ne” (count = 8)

low low low low lowest lowest newer newer newer
newer newer newer wider wider wider new new

corpus

5	l o w _
2	l o w e s t _
6	n e w e r _
3	w i d e r _
2	n e w _

Merge “ne”

corpus

5	l o w _
2	l o w e s t _
6	ne w e r _
3	w i d e r _
2	ne w _

vocabulary

, d, e, i, l, n, o, r, s, t, w, er, er



vocabulary

, d, e, i, l, n, o, r, s, t, w, er, er, ne



BPE: Counting & Merging

Continue the process to merge more adjacent symbols

low low low low lowest lowest newer newer newer
newer newer newer wider wider wider new new

corpus

5	l o w _
2	l o w e s t _
6	n e w er_
3	w i d er_
2	n e w _

merge**current vocabulary**

(ne, w)	_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new
(l, o)	_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo
(lo, w)	_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low
(new, er_)	_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_
(low, __)	_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_, low_

BPE: Token Segmenter

- Once we learn our vocabulary, we need a token segmenter to tokenize an unseen sentence (from test set)
- Just run (greedily based on training data frequency) on the merge rules we have learned from the training data on the test data
- Example:
 - Assume the merge rules: [(e, r), (er, _), (n, e), (ne, w), (l, o), (lo, w), (new, er_), (low, _)]
 - First merge all adjacent “er”, then all adjacent “er_”, then all adjacent “ne”...
 - “newer_” from the test set will be tokenized as a whole word
 - “lower_” from the test set will be tokenized as “low” + “er_”

low low low low lowest lowest newer newer newer
newer newer newer wider wider wider new new

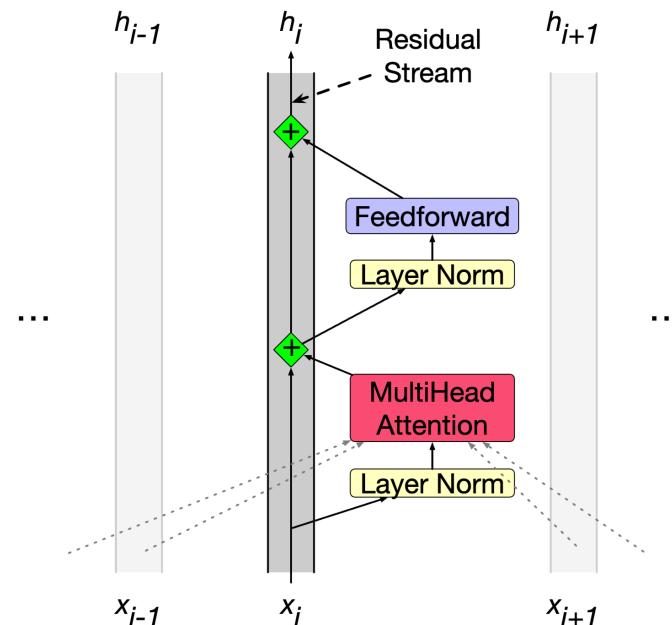
“lower_” is an unseen word from the training set

Agenda

- Tokenization
- Other Transformer Modules
- Transformer Language Model Pretraining

Transformer Block

- Modules in Transformer layers:
 - Multi-head attention
 - Layer normalization (LayerNorm)
 - Feedforward network (FFN)
 - Residual connection



Layer Normalization: Motivation

- Proposed in [Ba et al. \(2016\)](#)
- The distribution of inputs to DNN can change during training – “internal covariate shift”
- Slow down the training process: the model constantly adapts to changing distributions

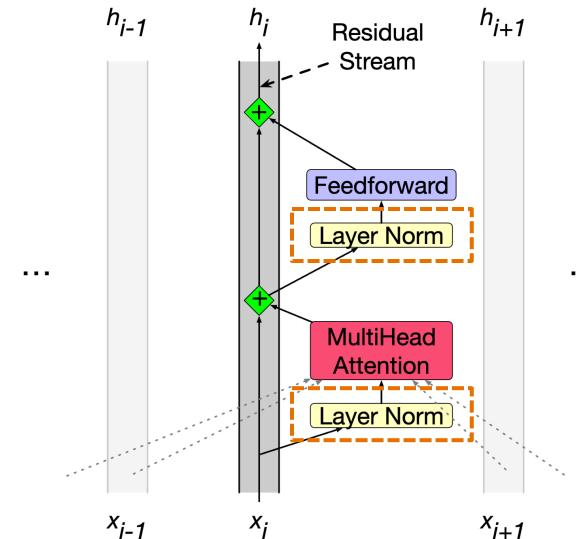


Figure source: <https://web.stanford.edu/~jurafsky/slp3/8.pdf>

Layer Normalization: Solution

- Normalize the input vector \mathbf{x}
 - Calculate the mean & standard deviation over the input vector dimensions

$$\mu = \frac{1}{d} \sum_{i=1}^d x_i \quad \sigma = \sqrt{\frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2}$$

- Apply normalization
- $$\hat{\mathbf{x}} = \frac{\mathbf{x} - \mu}{\sigma}$$
- Learn to scale and shift the normalized output with parameters

$$\text{LayerNorm}(\mathbf{x}) = \gamma \frac{\mathbf{x} - \mu}{\sigma} + \beta$$

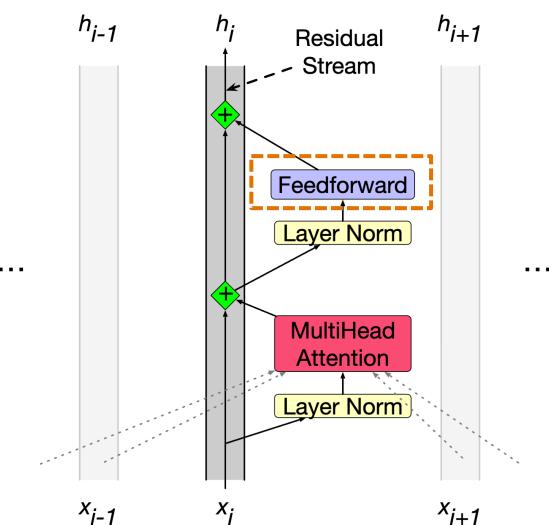
↑ ↑
Learnable parameters

Feedforward Network (FFN)

- FFN in Transformer is a 2-layer network (one hidden layer, two weight matrices)

$$\text{FFN}(\mathbf{x}_i) = \text{ReLU}(\mathbf{x}_i \mathbf{W}_1) \mathbf{W}_2$$

- Apply non-linear activation after the first layer
- Same weights applied to every token
- Weights are different across different Transformer layers

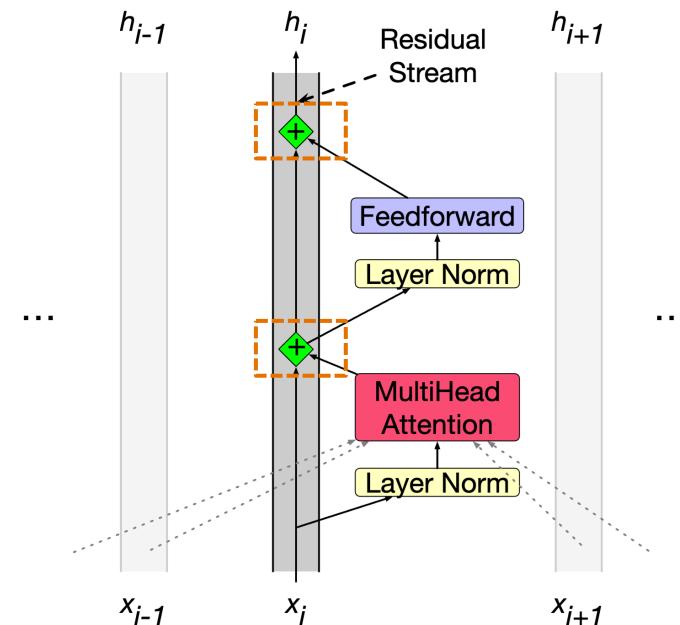


Residual Connections

- Add the original input to the output of a sublayer (e.g., attention/FFN)

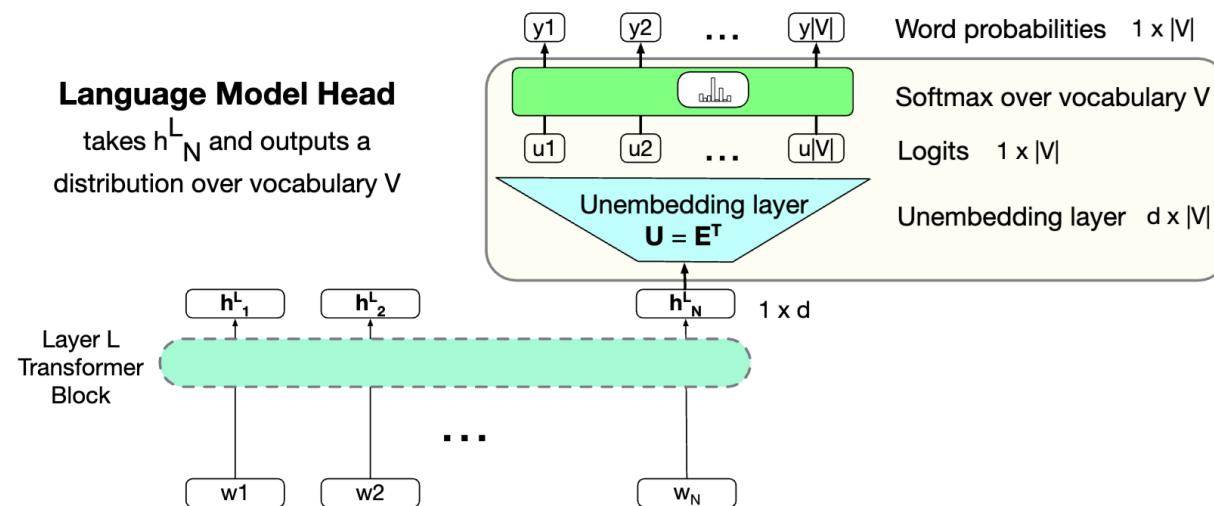
$$\mathbf{y} = \mathbf{x} + f(\mathbf{x})$$

- Benefits
 - Address the vanishing gradient problem
 - Facilitate information flow across the network
 - Help scale up model



Language Model Head

- Language model head is added to the final layer
- Usually apply the weight tying trick (share weights between input embeddings and the output embeddings)



Transformer Language Model: Overview

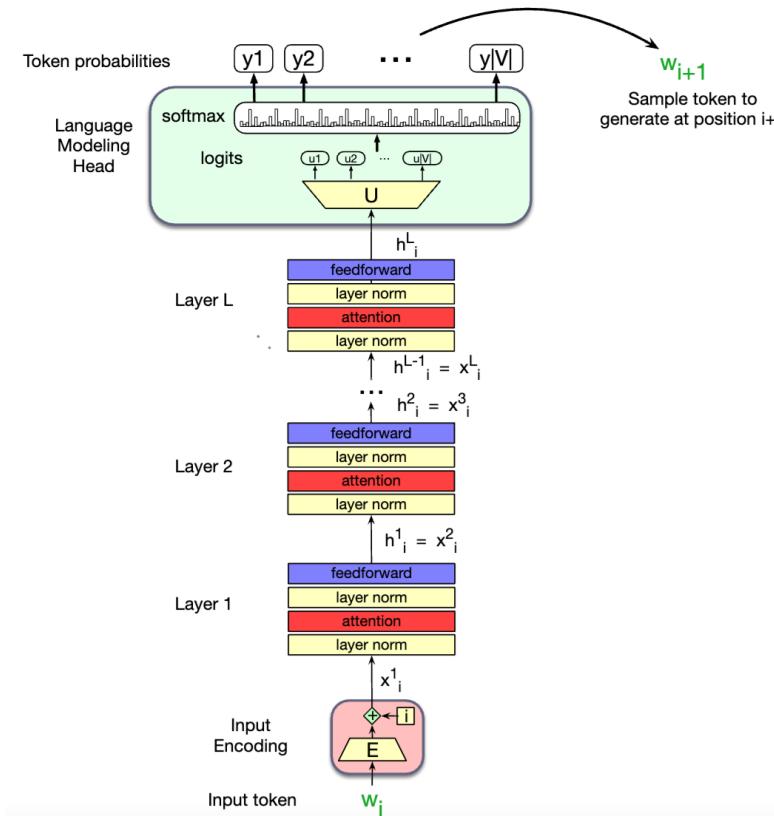
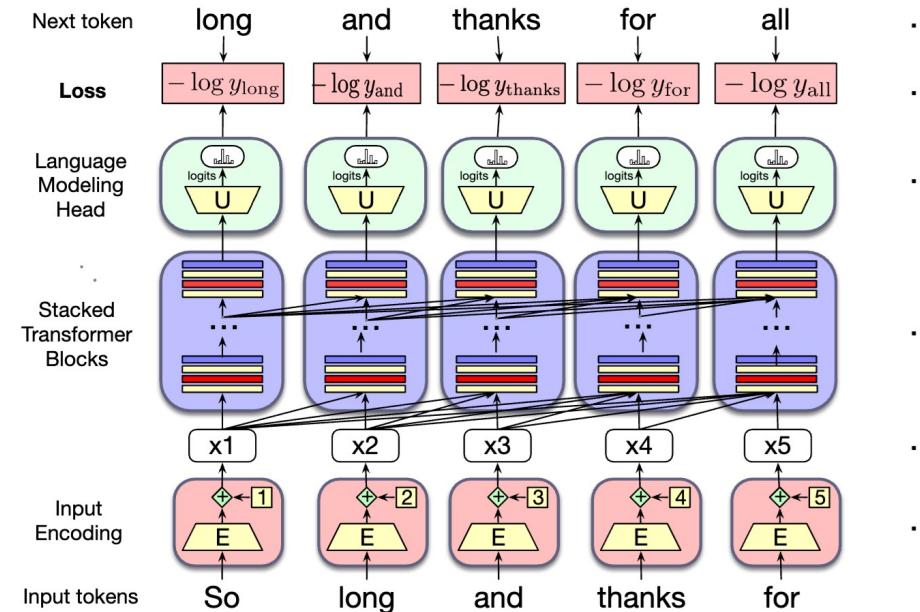


Figure source: <https://web.stanford.edu/~jurafsky/slp3/8.pdf>

Transformer Language Model Training

Use cross-entropy loss to train Transformers for language modeling (like RNN LMs)



Agenda

- Tokenization
- Other Transformer Modules
- Transformer Language Model Pretraining

Pretraining: Motivation

- Before pretraining became prevalent in NLP, most NLP models were trained from scratch on downstream task data
- **Data scarcity:** many NLP tasks do not have large labeled datasets available (costly to obtain)
- **Poor generalization:** models trained from scratch on specific tasks do not generalize well to unseen data or other tasks
- **Sensitivity to noise and randomness:** models are more likely to learn spurious correlations or be affected by annotation errors/randomness in training

Pretraining: Motivation

- There are abundant text data on the web, with rich information of linguistic features and knowledge about the world
- Learning from these easy-to-obtain data greatly benefits various downstream tasks



WIKIPEDIA
The Free Encyclopedia

The
New York
Times



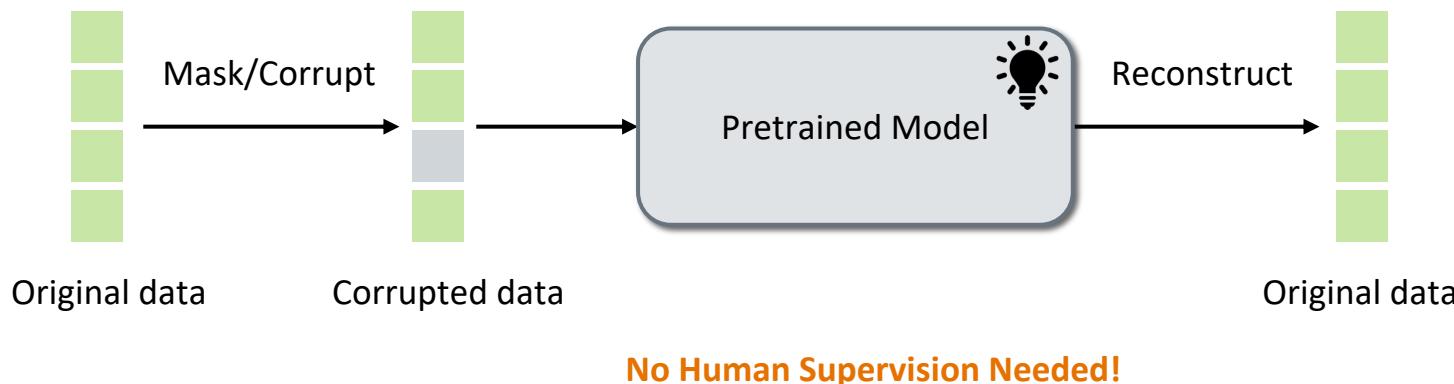
Pretraining: Multi-Task Learning

- In my free time, I like to **{run, banana}** (*Grammar*)
- I went to the zoo to see giraffes, lions, and **{zebras, spoon}** (*Lexical semantics*)
- The capital of Denmark is **{Copenhagen, London}** (*World knowledge*)
- I was engaged and on the edge of my seat the whole time. The movie was **{good, bad}** (*Sentiment analysis*)
- The word for “pretty” in Spanish is **{bonita, hola}** (*Translation*)
- $3 + 8 + 4 = \{15, 11\}$ (*Math*)
- ...

Examples from: https://docs.google.com/presentation/d/1hQUd3pF8_2Gr2Obc89LKjmHL0DIH-uof9M0yFVd3FA4/edit#slide=id.g28e2e9aa709_0_1

Pretraining: Self-Supervised Learning

- Pretraining is a form of **self-supervised** learning
- Make a part of the input unknown to the model
- Use other parts of the input to reconstruct/predict the unknown part



Pretraining + Fine-Tuning

- Pretraining: trained with pretext tasks on large-scale text corpora
- Fine-tuning (continue training): adjust the pretrained model's parameters with fine-tuning data
- Fine-tuning data can have different forms:
 - Task-specific labeled data (e.g., sentiment classification, named entity recognition)
 - (Multi-turn) dialogue data (i.e., instruction tuning)

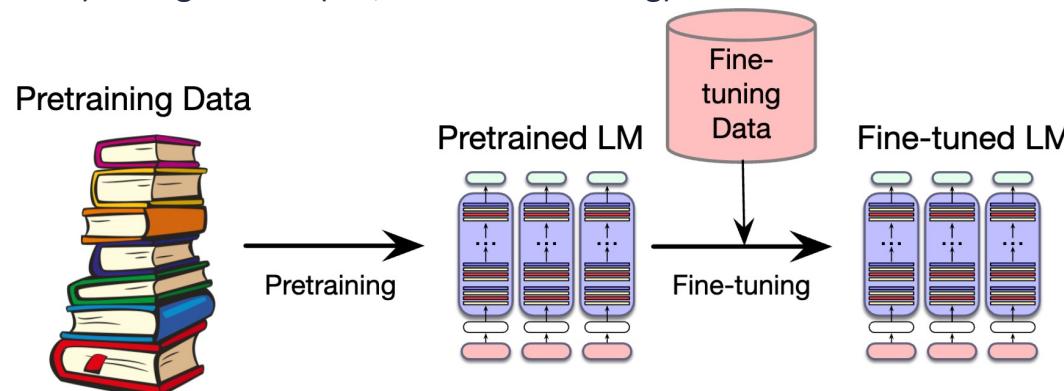


Figure source: <https://web.stanford.edu/~jurafsky/slp3/7.pdf>

Transformer for Pretraining

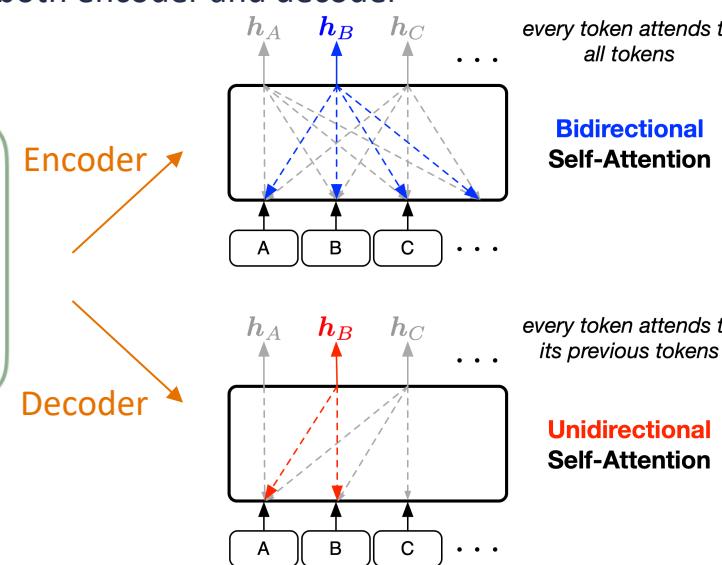
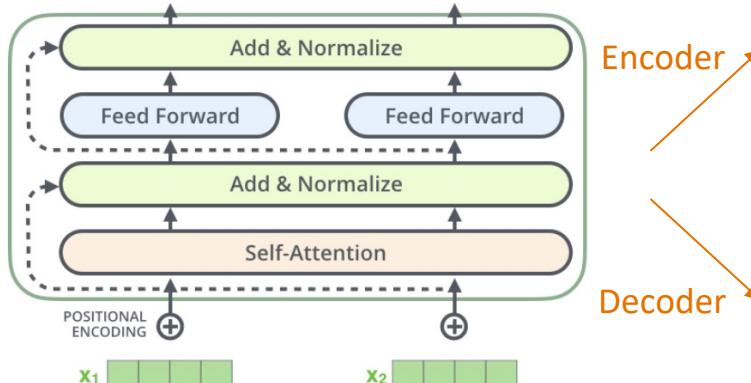
- Transformer is the common backbone architecture for language model pretraining
- **Efficiency:** Transformer processes all tokens in a sequence simultaneously – fast and efficient to train, especially on large datasets
- **Scalability:** Transformer architectures have shown impressive scaling properties, with performance improving as model size and training data increase (more on this later!)
- **Versatility:** Transformer can be adapted for various tasks and modalities beyond just text, including vision, audio, and other multimodal applications

Agenda

- Other Transformer Modules
- Language Model Pretraining: Overview
- Pretraining for Different Transformer Architectures

Transformer Architectures

- Based on the type of self-attention, Transformer can be instantiated as
 - Encoder: Bidirectional self-attention
 - Decoder: Unidirectional self-attention
 - Encoder-decoder: Use both encoder and decoder



N	q1·k1	q1·k2	q1·k3	q1·k4
	q2·k1	q2·k2	q2·k3	q2·k4
	q3·k1	q3·k2	q3·k3	q3·k4
	q4·k1	q4·k2	q4·k3	q4·k4

N	q1·k1	$-\infty$	$-\infty$	$-\infty$
	q2·k1	q2·k2	$-\infty$	$-\infty$
	q3·k1	q3·k2	q3·k3	$-\infty$
	q4·k1	q4·k2	q4·k3	q4·k4

Applications of Transformer Architectures

- Encoder (e.g., BERT):
 - Capture bidirectional context to learn each token representations
 - Suitable for natural language understanding (NLU) tasks
- Decoder (modern large language models, e.g., GPT):
 - Use prior context to predict the next token (conventional language modeling)
 - Suitable for natural language generation (NLG) tasks
 - Can also be used for NLU tasks by generating the class labels as tokens
- Encoder-decoder (e.g., BART, T5):
 - Use the encoder to process input, and use the decoder to generate outputs
 - Can conduct all tasks that encoders/decoders can do

NLU:

Text classification
Named entity recognition
Relation extraction
Sentiment analysis
...

NLG:

Text summarization
Machine translation
Dialogue system
Question answering
...

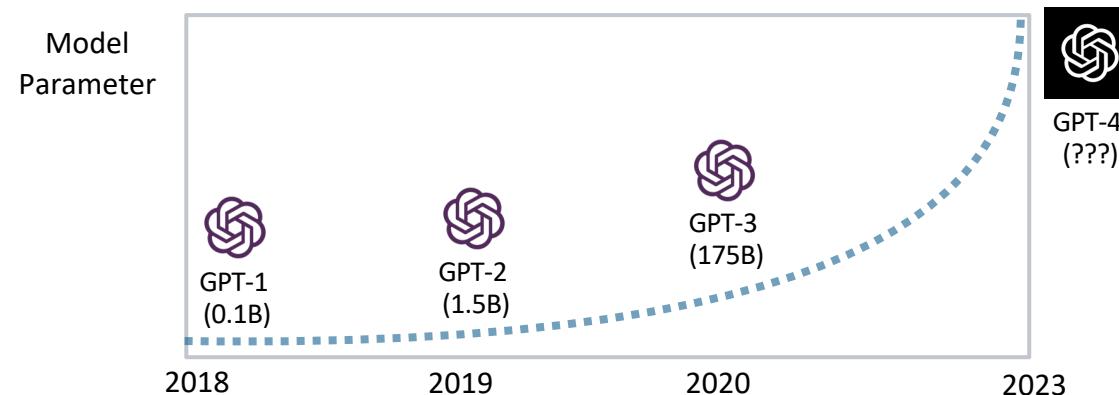
Decoder Pretraining

- Decoder architecture is the prominent choice in large language models
- Pretraining decoders is first introduced in GPT (generative pretraining) models
- Follow the standard language modeling (cross-entropy) objective

$$\mathcal{L}(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(x_i | x_1, x_2, \dots, x_{i-1})$$

GPT Series

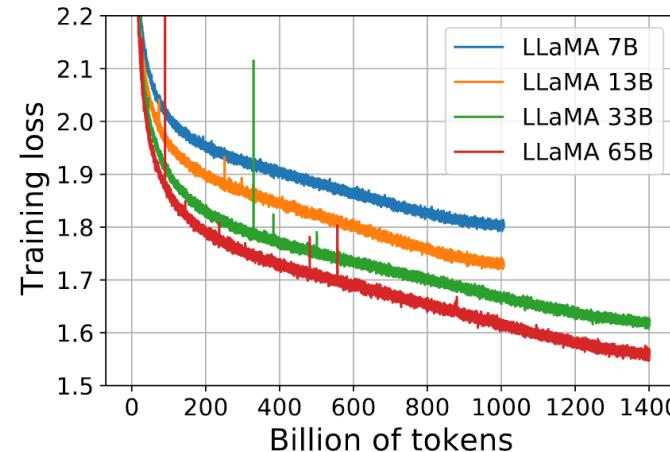
- GPT-1 (2018): 12 layers, 117M parameters, trained in ~1 week
- GPT-2 (2019): 48 layers, 1.5B parameters, trained in ~1 month
- GPT-3 (2020): 96 layers, 175B parameters, trained in several months



Papers: (GPT-1) https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf
(GPT-2) https://d4mucfpksywv.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf
(GPT-3) <https://arxiv.org/pdf/2005.14165.pdf>

Llama Series

- Llama-1 (2023/02): 7B/13B/33B/65B
- Llama-2 (2023/07): 7B/13B/70B
- Llama-3 (3.1 & 3.2) (2024/07): 1B/3B/8B/70B/405B w/ multi-modality



Larger models learn
pretraining data better

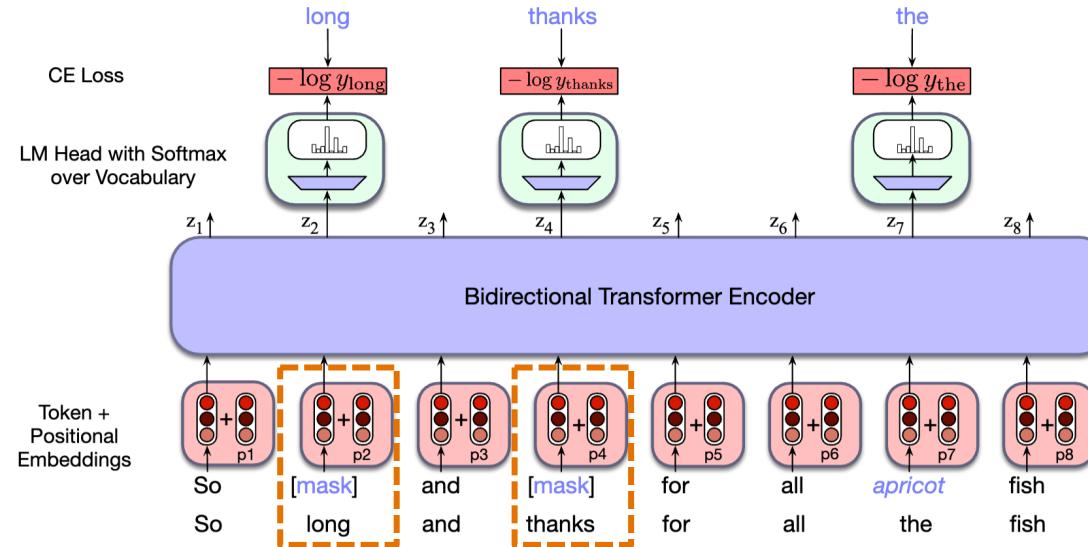
Papers: (Llama-1) <https://arxiv.org/pdf/2302.13971.pdf>
(Llama-2) <https://arxiv.org/pdf/2307.09288.pdf>
(Llama-3) <https://arxiv.org/pdf/2407.21783.pdf>

Further Reading on Decoder LMs

- [Mistral 7B](#) [Jiang et al., 2023]
- [Qwen Technical Report](#) [Bai et al., 2023]
- [GPT-4 Technical Report](#) [OpenAI, 2023]
- [Gemma: Open Models Based on Gemini Research and Technology](#) [Gemma Team, 2024]

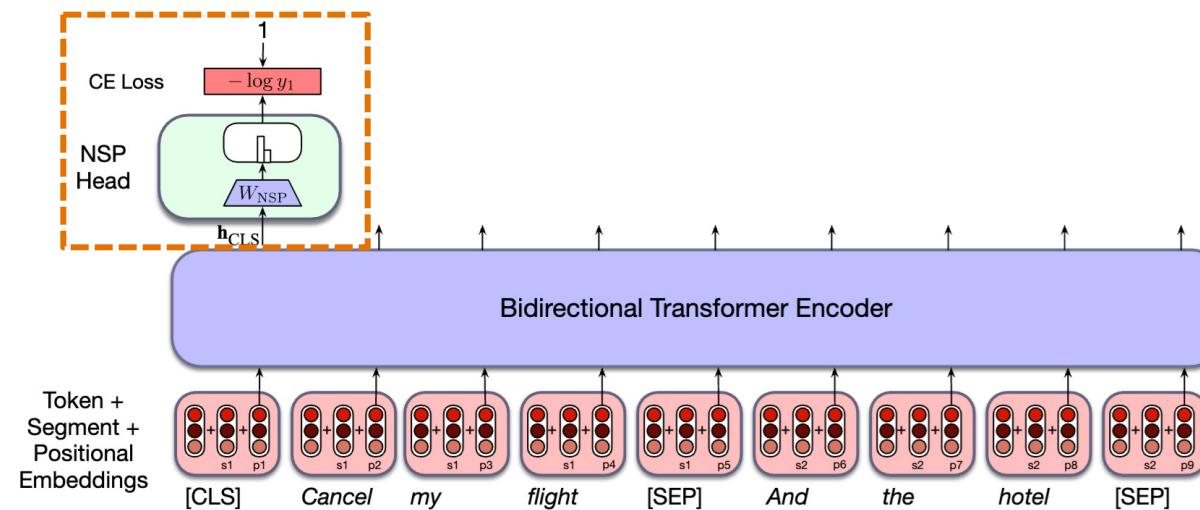
Encoder Pretraining: BERT

- BERT pretrains encoder models with bidirectionality
- **Masked language modeling (MLM):** With 15% words randomly masked, the model learns bidirectional contextual information to predict the masked words



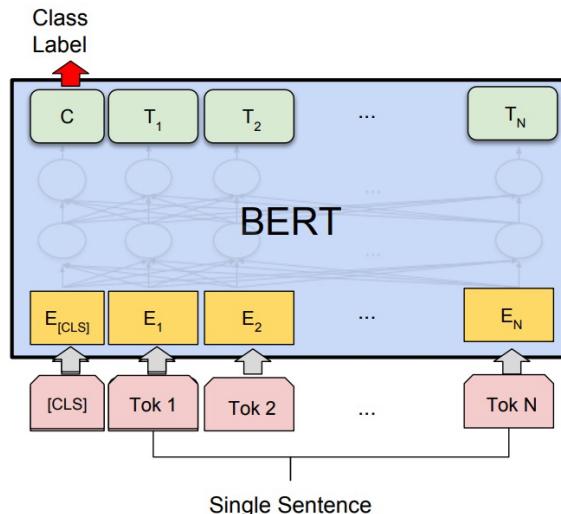
Encoder Pretraining: BERT

- **Next sentence prediction (NSP):** the model is presented with pairs of sentences
- The model is trained to predict whether each pair consists of an actual pair of adjacent sentences from the training corpus or a pair of unrelated sentence

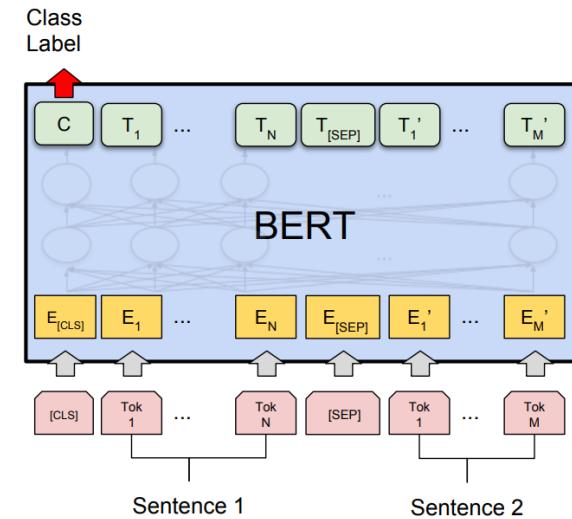


BERT Fine-Tuning

- Fine-tuning pretrained BERT models takes different forms depending on task types
- Usually replace the LM head with a linear layer fine-tuned on task-specific data



Single sequence classification



Sequence-pair classification

BERT vs. GPT on NLU tasks

- BERT outperforms GPT-1 on a set of NLU tasks
- Encoders capture **bidirectional** contexts – build a richer understanding of the text by looking at both preceding and following words
- Are encoder models still better than state-of-the-art (large) decoder models?
 - LLMs can be as good as (if not better than) encoders model on NLU: [Can ChatGPT Understand Too?](#)
 - The sheer model size + massive amount of pretraining data compensate for LLMs' unidirectional processing

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

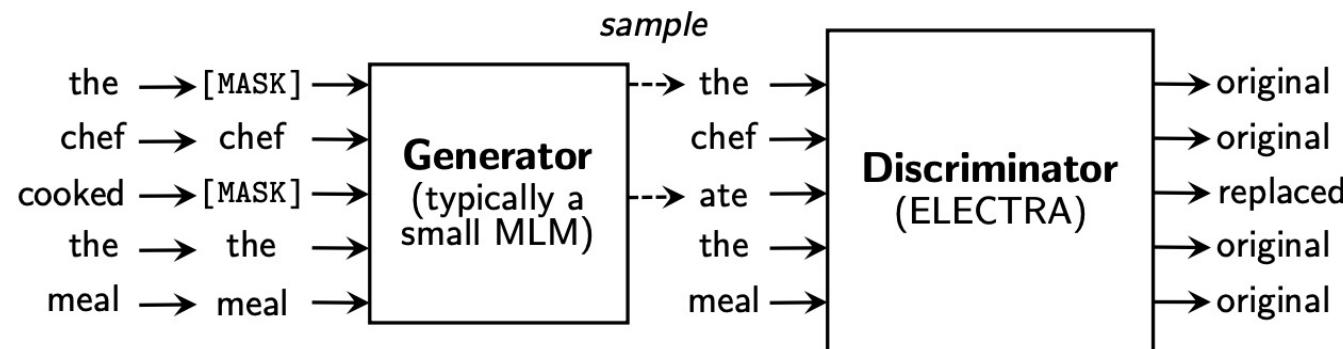
BERT Variant I: RoBERTa

- Pretrain the model for longer, with bigger batches over more data
- Pretrain on longer sequences
- Dynamically change the masking patterns applied to the training data in each epoch

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

BERT Variant II: ELECTRA

- Use a small MLM model as an auxiliary generator (discarded after pretraining)
- Pretrain the main model as a discriminator
- The small auxiliary MLM and the main discriminator are jointly trained
- The main model's pretraining task becomes more and more challenging in pretraining
- Major benefits: sample efficiency + learning curriculum



ELECTRA Performance

- ELECTRA pretraining incurs lower computation costs compared to MLM
- Better downstream task performance

Model	Train FLOPs	Params	CoLA	SST	MRPC	STS	QQP	MNLI	QNLI	RTE	Avg.
BERT	1.9e20 (0.27x)	335M	60.6	93.2	88.0	90.0	91.3	86.6	92.3	70.4	84.0
RoBERTa-100K	6.4e20 (0.90x)	356M	66.1	95.6	91.4	92.2	92.0	89.3	94.0	82.7	87.9
RoBERTa-500K	3.2e21 (4.5x)	356M	68.0	96.4	90.9	92.1	92.2	90.2	94.7	86.6	88.9
XLNet	3.9e21 (5.4x)	360M	69.0	97.0	90.8	92.2	92.3	90.8	94.9	85.9	89.1
BERT (ours)	7.1e20 (1x)	335M	67.0	95.9	89.1	91.2	91.5	89.6	93.5	79.5	87.2
ELECTRA-400K	7.1e20 (1x)	335M	69.3	96.0	90.6	92.1	92.4	90.5	94.5	86.8	89.0
ELECTRA-1.75M	3.1e21 (4.4x)	335M	69.1	96.9	90.8	92.6	92.4	90.9	95.0	88.0	89.5

Further Reading on Encoder LMs

- [XLNet: Generalized Autoregressive Pretraining for Language Understanding](#) [Yang et al., 2019]
- [ALBERT: A Lite BERT for Self-supervised Learning of Language Representations](#) [Lan et al., 2020]
- [DeBERTa: Decoding-enhanced BERT with Disentangled Attention](#) [He et al., 2020]
- [COCO-LM: Correcting and Contrasting Text Sequences for Language Model Pretraining](#) [Meng et al. 2021]

Encoder-Decoder Pretraining: BART

- Pretraining: Apply a series of noising schemes (e.g., masks, deletions, permutations...) to input sequences and train the model to recover the original sequences
- Fine-Tuning:
 - For NLU tasks: Feed the same input into the encoder and decoder, and use the final decoder token for classification
 - For NLG tasks: The encoder takes the input sequence, and the decoder generates outputs autoregressively



BART Performance

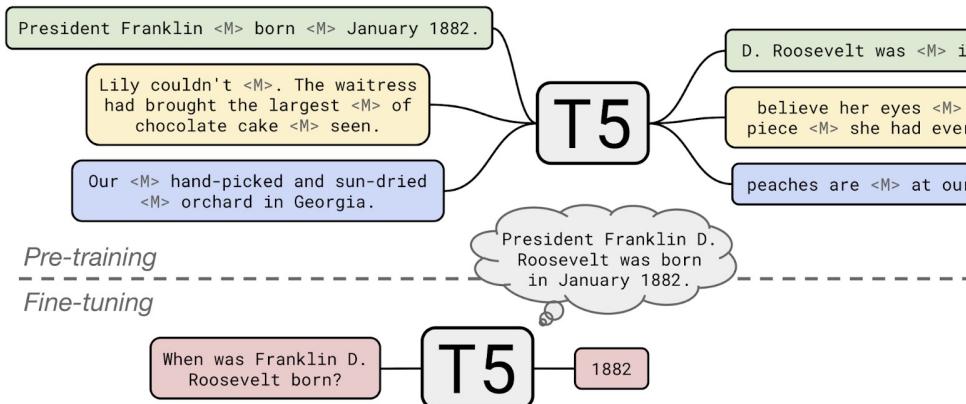
- Comparable to encoders on NLU tasks
- Good performance on NLG tasks

	SQuAD 1.1 EM/F1	SQuAD 2.0 EM/F1	MNLI m/mm	SST Acc	QQP Acc	QNLI Acc	STS-B Acc	RTE Acc	MRPC Acc	CoLA Mcc
BERT	84.1/90.9	79.0/81.8	86.6/-	93.2	91.3	92.3	90.0	70.4	88.0	60.6
UniLM	-/-	80.5/83.4	87.0/85.9	94.5	-	92.7	-	70.9	-	61.1
XLNet	89.0/94.5	86.1/88.8	89.8/-	95.6	91.8	93.9	91.8	83.8	89.2	63.6
RoBERTa	88.9/94.6	86.5/89.4	90.2/90.2	96.4	92.2	94.7	92.4	86.6	90.9	68.0
BART	88.8/94.6	86.1/89.2	89.9/90.1	96.6	92.5	94.9	91.2	87.0	90.4	62.8

	CNN/DailyMail			XSum		
	R1	R2	RL	R1	R2	RL
Lead-3	40.42	17.62	36.67	16.30	1.60	11.95
PTGEN (See et al., 2017)	36.44	15.66	33.42	29.70	9.21	23.24
PTGEN+COV (See et al., 2017)	39.53	17.28	36.38	28.10	8.02	21.72
UniLM	43.33	20.21	40.51	-	-	-
BERTSUMABS (Liu & Lapata, 2019)	41.72	19.39	38.76	38.76	16.33	31.15
BERTSUMEXTABS (Liu & Lapata, 2019)	42.13	19.60	39.18	38.81	16.50	31.27
BART	44.16	21.28	40.90	45.14	22.27	37.25

Encoder-Decoder Pretraining: T5

- T5: Text-to-Text Transfer Transformer
- Pretraining: Mask out spans of texts; generate the original spans
- Fine-Tuning: Convert every task into a sequence-to-sequence generation problem
- We'll see this model again in the instruction tuning lectures



T5 Performance

- Good performance across various tasks
- T5 vs. BART performance: unclear comparison due to difference in model sizes & training setups

Model	GLUE Average	CoLA Matthew's	SST-2 Accuracy	MRPC F1	MRPC Accuracy	STS-B Pearson	STS-B Spearman
Previous best	89.4 ^a	69.2 ^b	97.1 ^a	93.6^b	91.5^b	92.7 ^b	92.3 ^b
T5-Small	77.4	41.0	91.8	89.7	86.6	85.6	85.0
T5-Base	82.7	51.1	95.2	90.7	87.5	89.4	88.6
T5-Large	86.4	61.2	96.3	92.4	89.9	89.9	89.2
T5-3B	88.5	67.1	97.4	92.5	90.0	90.6	89.8
T5-11B	90.3	71.6	97.5	92.8	90.4	93.1	92.8
Model	QQP F1	QQP Accuracy	MNLI-m Accuracy	MNLI-mm Accuracy	QNLI Accuracy	RTE Accuracy	WNLI Accuracy
Previous best	74.8 ^c	90.7^b	91.3 ^a	91.0 ^a	99.2^a	89.2 ^a	91.8 ^a
T5-Small	70.0	88.0	82.4	82.3	90.3	69.9	69.2
T5-Base	72.6	89.4	87.1	86.2	93.7	80.1	78.8
T5-Large	73.9	89.9	89.9	89.6	94.8	87.2	85.6
T5-3B	74.4	89.7	91.4	91.2	96.3	91.1	89.7
T5-11B	75.1	90.6	92.2	91.9	96.9	92.8	94.5

Encoder-Decoder vs. Decoder-Only

- Modern LLMs are mostly based on the decoder-only Transformer architecture
- Simplicity:
 - Decoder-only models are simpler in structure (one Transformer model)
 - Encoder-decoder models require two Transformer models
- Efficiency:
 - Decoder-only models are more parameter-efficient for text generation
 - Encoder-decoder models' encoder part does not contribute to generation
- Scalability:
 - Decoder-only models scale very well with increased model size and data
 - Encoder-decoder models do not outperform decoder-only models at large model sizes



Thank You!

Yu Meng
University of Virginia
yumeng5@virginia.edu