

### **CS4248: Natural Language Processing**

Lecture 10 — Transformers & LLMs

### **Course Logistics**

- Assignments
  - Submission deadline for A3: Tue, Apr 2, 11.59 pm
- Project
  - Grades and comments for Intermediate Update posted
  - Optional consultation session you can register <u>here</u>
  - Submission deadline: Thu, Apr 18, 11:59 pm
  - Considering participating in STePS

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

## Supervised Training (RNN)

Task A: Learning a Language Model

[...] Precipitation forms as smaller droplets coalesce via collision with other raindrops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".[...]



[...] Precipitation forms as smaller droplets coalesce via collision with other raindrops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".[...]

Quick Quiz: Which model is easier to build? Why?

#### Task B: Learning a QA Systems



[...] Precipitation forms as smaller droplets coalesce via collision with other raindrops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".[...]

### **Transfer Learning for NLP Models**



### Transfer Learning with Word2Vec (or GloVe)

- Word2Vec: (almost) context-independent
  - BoW model → no consideration of word order
  - Limited window size → no consideration of whole sentence
  - Combining all the senses of a word into a single vector









Problem: Same word vector for all occurrences of "light"!

### **Goal: Contextualized Word Embeddings**

- What we want
  - Word representations should vary depending on context
  - Context = whole sentence + word order

#### "A light wind will make the traffic light collapse and light up in flames."



#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### ELMo — Embeddings from Language Model

- ELMo = RNN-based Language model, but...
  - LSTM instead of Vanilla RNN (better handling of long dependencies)
  - Bi-LSTM Bidirectional LSTM (forward and backward processing of sequence)
  - Two Bi-LSTM layers

(output of 1st layer = input of 2nd layer)

P("his" (57, I, (1/2)

Recall: Vanilla RNN Language Model









### ELMo — Final Embeddings



### ELMo — Evaluation

#### • Improvement of NLP downstream tasks

TASK	PREVIOUS SOTA		OUR BASELIN	ELMO + E BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	$88.7\pm0.17$	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2/17.2%
Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2/9.8%
NER	Peters et al. (2017)	$91.93\pm0.19$	90.15	$92.22\pm0.10$	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	$54.7\pm0.5$	3.3/6.8%

### ELMo — Evaluation

• Qualitative understanding what ELMo learns

~ Wond Wec				
	Source	Nearest Neighbors		
GloVe	play	playing, game, games, played, players, plays, player, Play, football, multiplayer		
biLM	Chico Ruiz made a spec- tacular play on Alusik 's grounder {} Olivia De Havilland signed to do a Broadway play for Garson {}	Kieffer , the only junior in the group , was commended for his ability to hit in the clutch , as well as his all-round excellent play .) $\rightarrow$ (Leore {} they were actors who had been handed fat roles in a successful play , and had talent enough to fill the roles competently , with nice understatement .		

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **RNN** — **Problem: (Very) Long Sequences**

- Training
  - Vanishing & Exploding Gradients problem (not detailed here)
- Information capture
  - Hidden state  $h_t$  must capture all information from  $h_0, h_1, ..., h_{t-1}$
  - Information dilutes over time → bottleneck

#### • Performance

- Processing is intrinsically sequential → no parallelization
- GPU-based performance gain depends on parallelization



### Transformer — Architecture

- Encoder-decoder architecture without recurrences
  - No long-range dependencies → no bottleneck
  - No sequential processing → easy to parallelize (note: this does not mean transformers are easier/faster to train!)
- Core concept: Attention
  - Alignment scores between all word pairs
- Important: Positional Embeddings
  - Preserve order of words in sequence



#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **Positional Encodings**

- Recall: RNNs process words sequentially
  - Considers order of words
  - Considers distance between words

- Transformers
  - Process all words all at once
  - No in-built mechanism to consider word order and word distances

#### Can we somehow encode the position of words?

(as part of preprocessing the input for the transformer)



### In-Lecture Activity (5 mins) — Positional Encodings

• Basic idea: Add "some" position embeddings p to initial word embeddings e



### Positional Encodings — Naive Approach 1

• Set position embedding values to actual position



→ **Problem:** positional encodings quickly start "dominating" word embeddings

• Magnitude of positional embedding values depends on sequence length N

### Positional Encodings — Naive Approach 2

• Set position embedding values to  $\frac{pos}{N-1}$ 





- → Problem: positional encodings depend on the length of the sequence length
  - encoding of the same position will differ for sequences with different lengths

### **Positional Encodings** — **Proposed Approach**



### **Positional Encodings** — Visualized



#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers: Attention
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **RNN Attention (revisited)**

0.02

 $h_{s}^{(2)}$ 

ging

Ct

0.06

 $h_{s}^{(3)}$ 

nach

0.07

 $h_s^{(\exists}$ 

Hause

**Attention Layer** 

Encoder RNN

0.85

 $h_s^{(1)}$ 

Ich



Step 2: Calculation of Attention Weights





 $h_t$ 

<s>





Decoder RNN

### **RNN Attention (revisited)**





#### **Scaled Dot-Product Attention**

- Intuition: queries Q, keys K, values V
- $k \in K$ ,  $q \in Q$  are vector of size  $d_k$
- scaling by  $\sqrt{d_k}$  leads to more stable gradients

### **Scaled Dot-Product Attention**



Source: Attention is all You Need

# Attention Head size of P. "Goden stude" of tras

- Attention Head
  - Maps model size  $d_{model}$  to size of queries, keys, and values (by default: same size)

Proposed: 
$$d_q = d_k = d_v = (d_{model})$$
  
Number of heads;

see next slide

Quick Quiz: What do you think is the reason for dividing by the number of heads?

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers: Multi-Head Attention
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **Multi-Head Attention (MHA)**

- Multi-Head Attention purpose / intuition
  - A word may relate to multiple other words in a sentence
  - A single Attention Head considers only one instance of relationship between pairs of words
  - MHA allows to capture different relationships
    (note that each Attention Head comes with its own weight matrices!)
  - Parameter: number of heads  $\rightarrow h$

![](_page_30_Picture_6.jpeg)

### **Multi-Head Attention**

![](_page_31_Figure_1.jpeg)

Output Probabilities

### **Multi-Head Attention**

![](_page_32_Figure_1.jpeg)

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers: Feed-Forward Layer
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **Feed Forward Layer**

- Feed Forward Layer purpose
  - The original paper doesn't say
  - ...uh, increase capacity / complexity

Feed-forward layers constitute two-thirds of a transformer model's parameters, yet their role in the network remains under-explored.

![](_page_34_Figure_5.jpeg)

Source: Transformer Feed-Forward Layers Are Key-Value Memories (2021)

### **Feed Forward Layer**

![](_page_35_Figure_1.jpeg)

Output Probabilities

Softmax
#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges



# Encoder Layer / Lluck





### **Encoder** — Self-Attention

• Example: German-to-English machine translation





### **Decoder Layer**

• The same components as Encoder Layer

- Multi-Head Attention but 2 MHA blocks (one for output, once for input/output)
- Feed Forward Layer
- The same additional concepts (residual connections, dropout, layer normalization)
- Multiple layers for complete decoder



# **Decoder Layer**



Output

### **Decoder** — Attentions

• Example: German-to-English machine translation



# **Complete Decoder**



Output

# **Complete Transformer**





#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

# Masking — Purpose



- Masking: Ignore attention between "invalid" words most commonly
  - Padding in batches with sequences of different lengths
  - "Hidden" words in models for Language Modeling
  - "Future" words in models for text generation
    - "surve" special
- Masking padded words

	-		Qèr	
best	movie	ever	<pad></pad>	<pad></pad>
i	really	liked	only	the
top	movie	<pad></pad>	<pad></pad>	<pad></pad>
such	а	dumb	and	silly
could	have	been	much	worse
the	story	was	not	that



# Masking for Language Modeling

- Masked Language Model basic idea
  - Mask a random number of word in a inputs sequence (e.g., BERT: 15%)
  - Train model transformer encoder to predict masked words



# $\alpha_{ij} = O(1) \left( \text{opened}, \left[ MASK_{j} \right] \right)$ $C_{ij} \neq (-\alpha_{i}) = -20 \implies 0$





### **Masking for Text Generation**

- Decoder is autoregressive
  - Output is generated word-by-word
  - During training, decoder gets complete output sequence (i.e., the decoder could "cheat" and look at subsequent words)
  - Ignore attention between a word and "future" words
  - Only affects self-attention MHA block
- Example
  - German-to-English machine translation



#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **Attention** — **Performance Considerations**

- Attentions is all you need...but it doesn't come for free
  - Pro: no sequential processing required → easy parallelize
  - Cons: number of operations for attention:  $N^2$  (*N* = sequence length)



### **Attention** — **Performance Considerations**

- Alternative: "restricted" attention
  - Does not compute attention between all pairs of words
  - $\blacksquare$  Main goal: make number of operations to be in  ${\it O}(N)$
- Example:







#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **Architectures**



(individually or as a group; include all group members' names in the post)

# The LLM Craze

Observation: Decoder-only dominates!

- Simpler architecture & setup
- More cheaply to train (relatively)
- More suitable for text generation
- Good zero-shot generalization



Source: Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond (2023)

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **BERT** (Bidirectional Encoder Representations from Transformers)

### • BERT

- Uses only the Transformer Encoder
- Self-supervised training
- Train on 2 learning objectives
   MLM: Masked Language Model

(predicted the words masked in the input sentences)

NSP: Next Sentence Prediction
 (predict if 2nd sentence was indeed followed 1st sentence)



### **BERT** (Bidirectional Encoder Representations from Transformers)

#### Pretraining

**Fine-Tuning** for specific task



### **RoBERTa** (A Robustly Optimized Bidirectional Encoder Representations from Transformers)

### RoBERTa ≈ BERT scaled up

- Same architecture, similar training setup (MLM only) but longer training using more data
- Dynamic masking: masking done during training time (BERT uses "static" masking: masking done during preprocessing)

### • Other BERT variants

DistilBERT

ALBERT

BERT October 11, 2018	BERT         RoBERTa         DistilBERT           October 11, 2018         July 26, 2019         October 2, 2019		ALBERT September 26, 2019	
Base: 110M Large: 340M	Base: 125 Large: 355	<b>Base:</b> 66	Base: 12M Large: 18M	
Base: 12 / 768 / 12 Large: 24 / 1024 / 16	Base: 12 / 768 / 12 Large: 24 / 1024 / 16	Base: 6 / 768 / 12	Base: 12 / 768 / 12 Large: 24 / 1024 / 16	
Base: 8 x V100 x 12d Large: 280 x V100 x 1d	1024 x V100 x 1 day (4-5x more than BERT)	Base: 8 x V100 x 3.5d (4 times less than BERT)	[not given] Large: 1.7x faster	
Outperforming SOTA in Oct 2018	88.5 on GLUE	97% of BERT-base's performance on GLUE	89.4 on GLUE	
BooksCorpus + English Wikipedia = 16 GB	BERT + CCNews + OpenWebText + Stories = 160 GB	BooksCorpus + English Wikipedia = 16 GB	BooksCorpus + English Wikipedia = 16 GB	
Bidirectional Trans- former, MLM & NSP	BERT without NSP, Using Dynamic Masking	BERT Distillation	BERT with reduced para- meters & SOP (not NSP)	
	BERT October 11, 2018 Base: 110M Large: 340M Base: 12 / 768 / 12 Large: 24 / 1024 / 16 Base: 8 × V100 × 12d Large: 280 × V100 × 1d Outperforming SOTA in Oct 2018 BooksCorpus + English Wikipedia = 16 GB Bidirectional Trans- former, MLM & NSP	BERTRoBERTaOctober 11, 2018July 26, 2019Base: 110MBase: 125Large: 340MLarge: 355Base: 12 / 768 / 12Large: 355Base: 24 / 1024 / 16Base: 12 / 768 / 12Large: 24 / 1024 / 16IO24 x V100 x 1 dayBase: 8 x V100 x 12d1024 x V100 x 1 dayLarge: 280 x V100 x 1d1024 x V100 x 1 dayCoutperforming SOTA in Oct 201888.5 on GLUEBooksCorpus + English Wikipedia = 16 GBBERT + CCNews + OpenWebText + Stories = 160 GBBidirectional Trans- former, MLM & NSPBERT without NSP, Using Dynamic Masking	BERT October11, 2018RoBERTa July 26, 2019DistilBERT October 2, 2019Base: 110M Large: 340MBase: 125 Large: 355Base: 66Base: 12 / 768 / 12 Large: 24 / 1024 / 16Base: 12 / 768 / 12 Large: 24 / 1024 / 16Base: 6 / 768 / 12Base: 8 × V100 × 12d Large: 280 × V100 × 1 day Oct 20181024 × V100 × 1 day (4-5x more than BERT)Base: 8 × V100 × 3.5d (4 times less than BERT)Outperforming SOTA in Oct 201888.5 on GLUE97% of BERT-base's performance on GLUEBooksCorpus + English Wikipedia = 16 GBBERT + CCNews + OpenWebText + Stories = 160 GBBooksCorpus + English Wikipedia = 16 GBBidirectional Trans- former, MLM & NSPBERT without NSP, Using Dynamic MaskingBERT Distillation	

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### • T5 — core concepts

- Basic encoder-decoder Transformer architecture
- Multi-task learning: training of model on multiple tasks at the same time

   (e.g., machine translation, coreference resolution, text summarization, sentence acceptability judgment, sentiment analysis)
- Each task is (re-)formulated as text-to-text task to match encoder-decoder architecture (including task-specific prefixes)



### T5 (Text-to-Text Transfer Transformer)

- T5 evaluation
  - The authors evaluated multi-task learning approach for different architectures
  - Best results: encoder-decoder architecture



Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
$\star$ 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	$\dot{M}$	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

### **BART** (Bidirectional and Auto-Regressive Transformers)

- BART core concepts
  - Basic encoder-decoder Transformer architecture
  - Trained by corrupting documents and then optimizing a reconstruction loss → denoising (the cross-entropy between the decoder's output and the original document)
  - Various transformation techniques to corrupt input documents



### BART ≈ BERT + GPT

#### BERT

- Random tokens are replaced with masks (e.g., [MASK])
- Input is encoded bidirectionally (not suitable for text generation)



#### GPT

- Auto-regressively word prediction (suitable for text generation)
- Words can only condition on leftward context (cannot learn bidirectional interactions)

#### BART

- Arbitrary noise transformation (not just BERT-like masking)
- Bidirectional encoding + auto-regression word prediction



#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### **GPT** (Generative Pretrained Transformer)

### • GPT

- Uses only the Transformer Decoder without the encoder attention block (alternatively: encoder with "do not look ahead" masking)
- Self-supervised training
- Learning objectives
  - Auto-regressive Language Model
- (Very) oversimplified history of GPT
  - GPT-1/2/3: text only, "just" making it larger; GPT-4: multimodal
  - GPT-3+: reinforcement learning from human feedback (RLHF)

Shext-word prodiction



### **GPT** (Generative Pretrained Transformer)

### • GPT-3 models

Model Name	$n_{ m params}$	$n_{ m layers}$	$d_{ m model}$	$n_{ m heads}$	$d_{ m head}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0  imes 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5  imes 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0  imes 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1 <b>M</b>	$1.6  imes 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	$0.6  imes 10^{-4}$

### **GPT** — **RLHF** (Reinforcement Learning from Human Feedback)

- RLHF two common setups
  - Use human-generated responses to prompts to fine-tune the pretrained model
  - Generate multiple response for same prompt; human ranks response; use ranking for fine-tuning



### LLaMA (Large Language Model Meta AI)

- Excoder-only architecture very similar to GPT (any many others!) main tweaks
  - Pre-normalization: layer normalization is put inside the residual blocks
  - SwiGLU (Swish-Gated Linear Unit) activation: non-monotonic, parameterized activation function
  - Rotary Positional Embeddings: encode word positions by rotating word embedding vectors
- Open LLM
  - Trained exclusively on publicly available data



### LLaMA — Pre-Normalization

- Post vs. pre-normalization
  - Post: layer normalization between residual blocks (original transformer)
  - Pre: layer normalization inside residual blocks (LLaMA, etc.)
  - Observed benefit of pre-normalization:
    - Well-behaved gradients at initialization
    - Significantly faster training



### LLaMA — SwiGLU (Swish-Gated Linear Unit)

#### GLU – Gated Linear Unit (paper)

- Gating proposed in LSTM paper (1997!)
- Parameterized activation function
- Many other variants proposed

 $GLU(x) = (xW+b)\otimes \sigma(xV+c)$ 

- Simple parameterized activation function
- Approach: "try and see what works best"

$$Swish(x) = x \otimes \sigma(eta x)$$



1

$$SwiGLU(x) = (xW+b)\otimes Swish_eta(xV+c)$$
### LLaMA — SwiGLU (Swish-Gated Linear Unit)

#### **ReLU (Linear Rectified Unit)**







### LLaMA — Rotary Positional Embeddings



### Outline

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

#### • Transformer-based LLMs

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

# The Future of Large Language Models — Opportunities

- Language models are an old idea What changed?
  - New architectures (here: Transformers)
  - More computing power
  - More and diverse data
  - More resources (i.e., money, manpower)



LLMs show Emergent Abilities



Abilities that were not explicitly programmed into the model but emerge from the training process

# The Future of Large Language Models — Opportunities

#### • Emergent abilities

- Language Generation (coherent and fluent text in a variety of styles and genres, from news articles to poetry)
- Question Answering (answering complex questions by extracting information from large amounts of text data)
- **Translation** (translating text between different languages with high accuracy)
- Summarization (generate concise summaries of long documents, allowing for efficient information extraction and consumption)
- Dialogue Generation (engage in natural and coherent conversations with humans)
- Common Sense Reasoning (basic degree of common sense reasoning; predicting outcome of simple scenarios)

#### ➔ Question: Can a language model <u>really</u> do these tasks?

EXPLAINER: What is ChatGPT and why are schools blocking it?

Will ChatGPT take my job? Here are 20 professions that could be replaced by AI

Hallucinations, Plagiarism, and ChatGPT

Letters | How universities can start to grapple with ChatGPT's capabilities

Hollywood: Writers Guild considering ChatGPT-written scripts, no AI credits

ChatGPT

The impact of Large Language Models on Law Enforcement

Criminals will soon use ChatGPT to make scams more convincing, experts warn; only 'a matter of time' before S'pore hit

**ChatGPT Poses Dangers for Online Dating Apps** 

Cybercriminals are using ChatGPT to create malware

A fake news frenzy: why ChatGPT could be disastrous for truth in journalism

Pause Giant AI Experiments: An Open Letter

ChatGPT invented a sexual harassment scandal and named a real law prof as the accused

Italy orders ChatGPT blocked citing data protection concerns

1,100+ notable signatories just signed an open letter asking 'all AI labs to immediately pause for at least 6 months'

AI can be racist, sexist and creepy. What should we do about it?

GPT-4 kicks AI security risks into higher gear

Europol sounds alarm as crooks tap into ChatGPT-4

GPT-5 expected this year, could make ChatGPT indistinguishable from a human

Experts Warn of Nightmare Internet Filling With Infinite AI-Generated Propaganda

What Have Humans Just Unleashed?

Call it tech's optical-illusion era: Not even the experts know exactly what will come next in the AI revolution.

Did a Robot Write This? We Need Watermarks to Spot Al

Exclusive: OpenAI Used Kenyan Workers on Less Than \$2 Per Hour to Make ChatGPT Less Toxic Australian Mayor Threatens to Sue OpenAl for Defamation by Chatbot

Artists sue AI company for billions, alleging "parasite" app used their work for free

### ChatGPT banned on Q&A site over 'substantially harmful' answers

\$120bn wiped off Google after Bard AI chatbot gives wrong answer

Microsoft tries to justify A.I.'s tendency to give wrong answers by saying they're 'usefully wrong'

Chat-GPT Pretended to Be Blind and Tricked a Human Into Solving a CAPTCHA

ChatGPT lies about scientific results, needs open-source alternatives, say researchers

Al isn't close to becoming sentient – the real danger lies in how easily we're prone to anthropomorphize it

#### ...and the biggest questions: Why does this all seem to work?

We have extended the GLU family of layers and proposed their use in Transformer. In a transfer-learning setup, the new variants seem to produce better perplexities for the de-noising objective used in pre-training, as well as better results on many downstream language-understanding tasks. These architectures are simple to implement, and have no apparent computational drawbacks. We offer no explanation as to why these architectures seem to work; we attribute their success, as all else, to divine benevolence.

### Outline

#### • Contextual Word Embeddings

- Motivation
- ELMo

#### • Transformers

- Positional Encoding
- Core Layers
- Encoder & Decoder

#### • Extended Concepts

- Masking
- Restricted Attention

#### • Transformer-based LLMs

- Overview
- Encoder-only: BERT, RoBERTa
- Encoder-Decoder: T5, BART
- Decoder-only: GPT, LLaMA
- Opportunities & Challenges

### Summary

- Transformer architecture
  - Encoder-decoder architecture
  - Core concept: attention (self-attention + cross attention)
  - Additional concepts: positional encoding, masking
- Large Language Models (LLMs)
  - Currently dominated by transformer architecture
  - Main categorization: encoder-only, encoder-decoder, decoder-only (with decoder-only models right now dominating the field)
  - Still continuously growing model zoo of LLMs

→ Last lecture: LLMs – problems, challenges, strategies

## **Pre-Lecture Activity for Next Week**

- Assigned Task
  - Do a web search and for the question stated below
  - Post you answer(s) to the question into the Discussion on Canvas (please cite or quote your sources)

*"What is the relationship between information retrieval and natural language processing?"* 

#### Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers but his won't help you learn better

# **Solutions to Quick Quizzes**

- Slide 4
  - Learning a language model is arguably easier since it is a self-supervised task
  - The annotations / labels are (almost) the same as the inputs → (relatively) easier to set up
  - Note: this does not mean that it's also easier to get good results for Task A
- Slide 29
  - This make the total number of trainable parameters independent from the number of heads