

Language Models for Law and Social Science

8. Large Language Models

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

Language Modeling

- ▶ “Language Modeling” refers to the task of teaching an algorithm to predict/generate language.

Language Modeling

- ▶ “Language Modeling” refers to the task of teaching an algorithm to predict/generate language.
- ▶ The standard approach uses the Markov assumption: future words are independent of the past given the present and some finite number of previous rounds.
 - ▶ A k th order markov-assumption assumes that the next word in a sequence depends only on the last k words:

$$\Pr(w_{i+1}|w_{1:i}) \approx \Pr(w_{i+1}|w_{i-k:i})$$

- ▶ The task is to learn $\Pr(w_{i+1}|w_{1:i})$ given a large corpus.

Perplexity

- ▶ Perplexity is an information-theoretic measurement of how well a probability model predicts a sample.
- ▶ Given a text corpus of n words $\{w_1, \dots, w_n\}$ and a language model function $\Pr(\cdot)$, the perplexity is:

$$2^{-\frac{1}{n} \sum_{i=1}^n \log \hat{\Pr}(w_i | w_{1:i-1})}$$

- ▶ (lower is better)
- ▶ “Good” language models (i.e., reflective of real language usage) assign high probabilities to the observed words in the corpus, resulting in lower perplexity values.

N-Gram Approach to Language Modeling

- ▶ Let $\#(w_{i:j})$ be the count of the sequence of words $w_{i:j}$ in the corpus.
- ▶ The MLE estimate for the probability of a word given the previous k words is

$$\widehat{\Pr}(w_{i+1} | w_{i-k:i}) = \frac{\#(w_{i-k:i+1})}{\#(w_{i-k:i})}$$

N-Gram Approach to Language Modeling

- ▶ Let $\#(w_{i:j})$ be the count of the sequence of words $w_{i:j}$ in the corpus.
- ▶ The MLE estimate for the probability of a word given the previous k words is

$$\widehat{\Pr}(w_{i+1} | w_{i-k:i}) = \frac{\#(w_{i-k:i+1})}{\#(w_{i-k:i})}$$

- ▶ Problem 1: zero events are quite common because many phrases are unique.
 - ▶ if $w_{i-k:i+1}$ was never observed in the corpus, $\widehat{\Pr}$ is zero.

N-Gram Approach to Language Modeling

- ▶ Let $\#(w_{i:j})$ be the count of the sequence of words $w_{i:j}$ in the corpus.
- ▶ The MLE estimate for the probability of a word given the previous k words is

$$\widehat{\Pr}(w_{i+1} | w_{i-k:i}) = \frac{\#(w_{i-k:i+1})}{\#(w_{i-k:i})}$$

- ▶ Problem 1: zero events are quite common because many phrases are unique.
 - ▶ if $w_{i-k:i+1}$ was never observed in the corpus, $\widehat{\Pr}$ is zero.
- ▶ Problem 2: infrequent events are very common
 - ▶ if $w_{i-k:i}$ is only observed a few times, then there are only a few words w_{i+1} with positive probability

Neural Language Modeling (Goldberg 2017)

- ▶ Input:
 - ▶ preceding sequence (context words) $w_{1:k}$.
 - ▶ V is a finite vocabulary, including special symbols for unknown words, start of sentence, and end of sentence.
 - ▶ Each context word is associated with an embedding vector.
 - ▶ The input vector x is a concatenation of the word vectors.
- ▶ Output:
 - ▶ probability distribution over the next word.
- ▶ Model architecture could be an MLP applied to the embeddings, a CNN, an RNN, or a transformer.

Beam Search

- ▶ Text decoding generally works word by word.
 - ▶ but a generated word at any given point might create a low-probability sequence of words.

Beam Search

- ▶ Text decoding generally works word by word.
 - ▶ but a generated word at any given point might create a low-probability sequence of words.
- ▶ Beam search generates multiple words at any given point, and follows those “beams” to generate several branching sequences.
 - ▶ after computing the sequences, e.g., 3-4 words, evaluate their probability and only choose beams with relatively high probability.

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

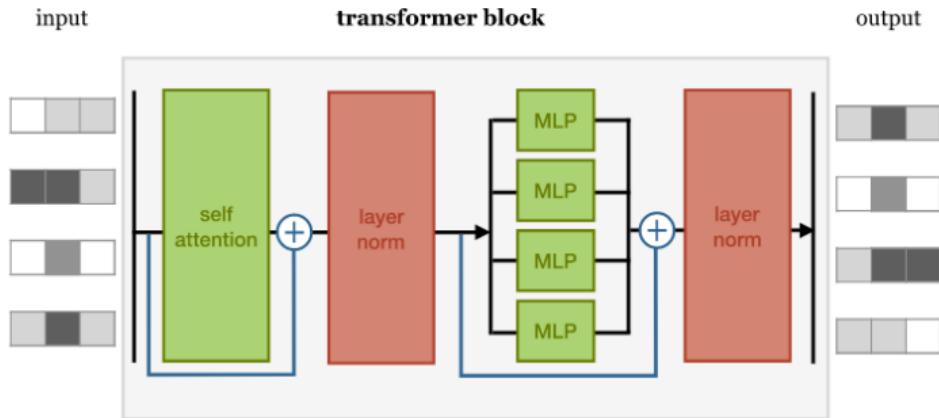
Open-Source Autoregressive LLMs

Additional Nuts and Bolts

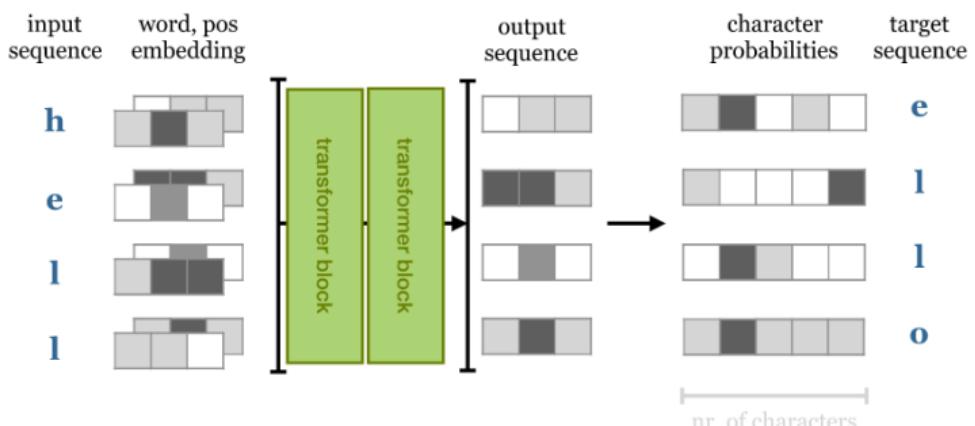
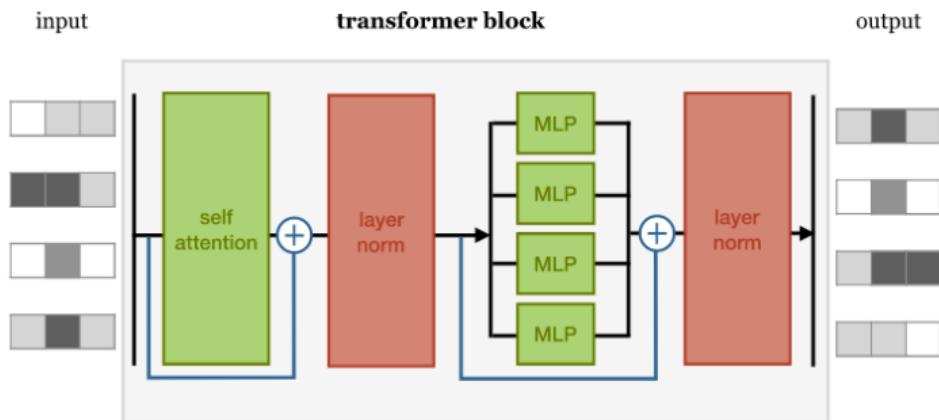
Three Types of Transformer-Based Language Models

- ▶ **Autoregressive models** (e.g. GPT):
 - ▶ pretrained on classic language modeling task: guess the next token having read all the previous ones.
 - ▶ during training/inference, attention heads only view previous tokens, not subsequent tokens.
 - ▶ ideal for text generation.
- ▶ **Autoencoding models** (e.g. BERT):
 - ▶ pretrained by encoding a sequence to solve a single prediction task – eg filling in a masked token.
 - ▶ usually build bidirectional representations and get access to the full sequence.
 - ▶ can be fine-tuned and achieve great results on many tasks, e.g. text classification.
- ▶ **Sequence-to-sequence models** (e.g. Vaswani et al 2017; BART)
 - ▶ models with both encoders and decoders
 - ▶ trained to encode a whole sequence before generating outputs
 - ▶ more specialized; works well for machine translation but not used much any more.

Text generation transformer

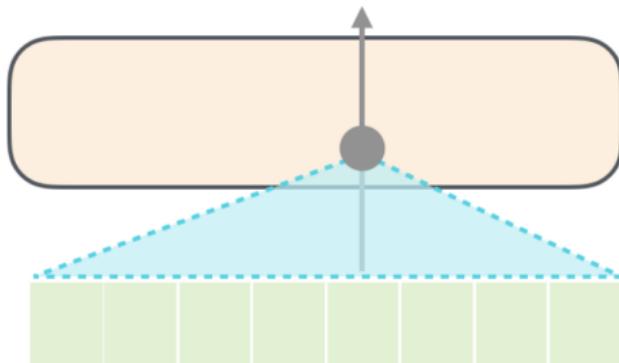


Text generation transformer

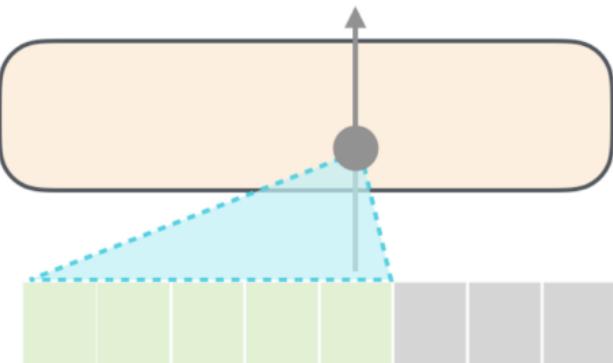


“Masked Self-Attention” (“Markov Self-Attention”)

Self-Attention



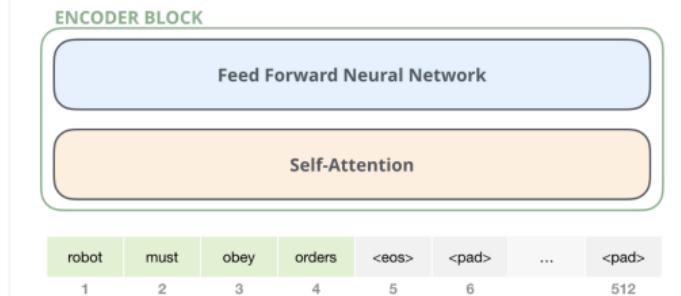
Masked Self-Attention



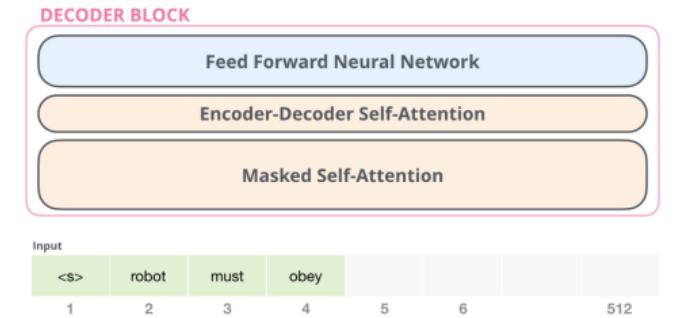
$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j$$

$$h_i = \sum_{j=1}^i a(x_i, x_j) x_j$$

Encoder Blocks vs. Decoder Blocks



An encoder block from the original transformer paper can take inputs up until a certain max sequence length (e.g. 512 tokens). It's okay if an input sequence is shorter than this limit, we can just pad the rest of the sequence.



Encoder block “reads” text:

- ▶ observes the whole sequence
- ▶ BERT is an “encoder-only” architecture

Decoder block “writes” text:

- ▶ uses masked self-attention → attention mechanism does not observe future tokens.
 - ▶ learn to generate text based on history
- ▶ GPT is a “decoder-only” architecture

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

BERT (and RoBERTa)

- ▶ BERT = Bidirectional Encoder Representations from Transformers
 - ▶ RoBERTa = Robust BERT
- ▶ Architecture:
 - ▶ a stack of transformer blocks with a self-attention layer and an MLP.
 - ▶ The largest BERT model has 24 blocks, embedding dimension of 768, and 16 attention heads.
≈ 340M parameters to learn.

BERT (and RoBERTa)

- ▶ BERT = Bidirectional Encoder Representations from Transformers
 - ▶ RoBERTa = Robust BERT
- ▶ Architecture:
 - ▶ a stack of transformer blocks with a self-attention layer and an MLP.
 - ▶ The largest BERT model has 24 blocks, embedding dimension of 768, and 16 attention heads.
≈ 340M parameters to learn.
- ▶ Task: Masked language modeling:
 - ▶ 15% of words masked
 - ▶ if masked: replace with [MASK] 80% of the time, a random token 10% of the time, and left unchanged 10% of the time.
 - ▶ model has to predict the original word.

BERT (and RoBERTa)

- ▶ BERT = Bidirectional Encoder Representations from Transformers
 - ▶ RoBERTa = Robust BERT
- ▶ Architecture:
 - ▶ a stack of transformer blocks with a self-attention layer and an MLP.
 - ▶ The largest BERT model has 24 blocks, embedding dimension of 768, and 16 attention heads.
≈ 340M parameters to learn.
- ▶ Task: Masked language modeling:
 - ▶ 15% of words masked
 - ▶ if masked: replace with [MASK] 80% of the time, a random token 10% of the time, and left unchanged 10% of the time.
 - ▶ model has to predict the original word.
- ▶ Unlike GPT, BERT attention observes all tokens in the sequence, reads backwards and forwards (bidirectional).

BERT (and RoBERTa)

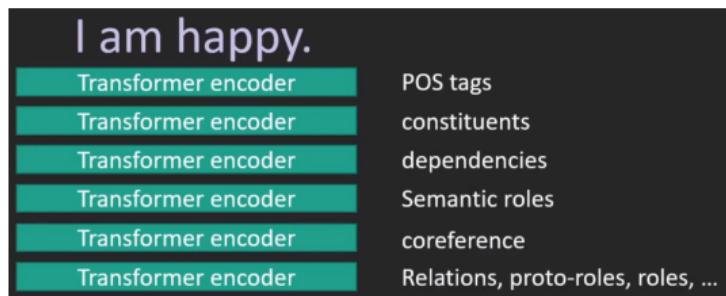
- ▶ BERT = Bidirectional Encoder Representations from Transformers
 - ▶ RoBERTa = Robust BERT
- ▶ Architecture:
 - ▶ a stack of transformer blocks with a self-attention layer and an MLP.
 - ▶ The largest BERT model has 24 blocks, embedding dimension of 768, and 16 attention heads.
≈ 340M parameters to learn.
- ▶ Task: Masked language modeling:
 - ▶ 15% of words masked
 - ▶ if masked: replace with [MASK] 80% of the time, a random token 10% of the time, and left unchanged 10% of the time.
 - ▶ model has to predict the original word.
- ▶ Unlike GPT, BERT attention observes all tokens in the sequence, reads backwards and forwards (bidirectional).
- ▶ Corpus:
 - ▶ 800M words from English books (modern work, from unpublished authors), by Zhu et al (2015).
 - ▶ 2.5B words of text from English Wikipedia articles (without markup).

- ▶ BERT still obtain state-of-the-art results on many NLP tasks (see Devlin et al 2019).
- ▶ The model can be fine-tuned as needed.
 - ▶ e.g. for sentiment analysis, text classification.
- ▶ use ModernBERT (<https://huggingface.co/blog/modernbert>)

- ▶ BERT still obtain state-of-the-art results on many NLP tasks (see Devlin et al 2019).
- ▶ The model can be fine-tuned as needed.
 - ▶ e.g. for sentiment analysis, text classification.
- ▶ use ModernBERT (<https://huggingface.co/blog/modernbert>)

BERT Redisovers the Classical NLP Pipeline

Ian Tenney¹ Dipanjan Das¹ Ellie Pavlick^{1,2}
¹Google Research ²Brown University
{iftenney,dipanjand,epavlick}@google.com



- ▶ The earlier and later layers in BERT respectively encode more functional and more semantic information.

Model Distillation

- ▶ Large transformer models such as BERT can be compressed.
 - ▶ a smaller model is given the inputs and BERT's outputs as the label.
 - ▶ works almost as well (97% of full BERT performance) and 60% faster
- When using pre-trained models, often better to use DistilBERT or DistilGPT.

Model Distillation

- ▶ Large transformer models such as BERT can be compressed.
 - ▶ a smaller model is given the inputs and BERT's outputs as the label.
 - ▶ works almost as well (97% of full BERT performance) and 60% faster
- When using pre-trained models, often better to use DistilBERT or DistilGPT.
- ▶ one reason this works:
 - ▶ for a given masked token, the student model observes probabilities across the whole vocabulary, not just the single true token.

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

Reading Comprehension \leftrightarrow Local Question Answering

Answering questions about a passage of text to show that the system understands the passage.

Reading Comprehension ↔ Local Question Answering

Answering questions about a passage of text to show that the system understands the passage.

<https://demo.allennlp.org/reading-comprehension>

Passage Context

The institutional framework of Navarre was preserved following the 1512 invasion. Once Ferdinand II of Aragon died in January, the Parliament of Navarre gathered in Pamplona, urging Charles V to attend a coronation ceremony in the town following tradition, but the envoys of the Parliament were met with the Emperor's utter indifference if not contempt. He refused to attend any ceremony and responded with a brief "let's say I am happy and pleases me." Eventually the Parliament met in 1517 without Charles V, represented instead by the **Duke of Najera** pronouncing an array of promises of little certitude, while the acting Parliament kept piling up grievances and demands for damages due to the Emperor, totalling 67—the 2nd Viceroy of Navarre Fadrique de Acuña was deposed in 1515 probably for acceding to send grievances. Contradictions inherent to the documents accounting for the Emperor's non-existent oath pledge in 1516 point to a contemporary manipulation of the records.

Question

Who represented the Charles V at Parliament?

Local Question Answering with BERT

Beyoncé Giselle Knowles-Carter (born September 4, 1981) is an American singer, songwriter, record producer and actress. Born and raised in **Houston, Texas**, she performed in various **singing and dancing** competitions as a child, and rose to fame in the late 1990s as lead singer of R&B girl-group Destiny's Child. Managed by her father, Mathew Knowles, the group became one of the world's best-selling girl groups of all time. Their hiatus saw the release of Beyoncé's debut album, *Dangerously in Love* (2003), which established her as a solo artist worldwide, earned five Grammy Awards and featured the Billboard Hot 100 number-one singles "Crazy in Love" and "Baby Boy".

Q: "In what city and state did Beyoncé grow up?"

A: **"Houston, Texas"**

Q: "What areas did Beyoncé compete in when she was growing up?"

A: **"singing and dancing"**

Q: "When did Beyoncé release *Dangerously in Love*?"

A: **"2003"**

Figure 23.11 A (Wikipedia) passage from the SQuAD 2.0 dataset (Rajpurkar et al., 2018) with 3 sample questions and the labeled answer spans.

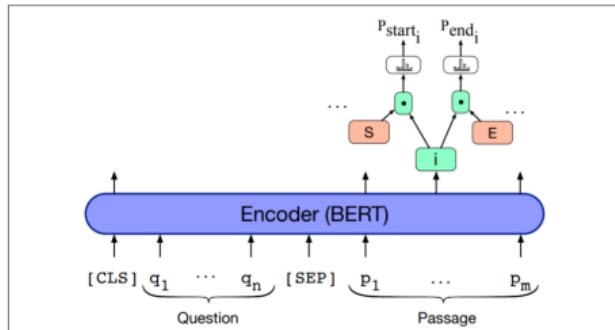


Figure 23.12 An encoder model (using BERT) for span-based question answering from reading-comprehension-based question answering tasks.

- BERT-based models learn to solve these problems by inputting {question + [SEP] + passage} and predicting the token indexes for the start and end of the answer.

Textual Entailment \leftrightarrow Natural Language Inference

- ▶ TE is the task of predicting whether, for a pair of sentences, the facts in the first sentence necessarily imply the facts in the second.

Sentence A (Premise)	Sentence B (Hypothesis)	Label
A soccer game with multiple males playing.	Some men are playing a sport.	entailment
An older and younger man smiling.	Two men are smiling and laughing at the cats playing on the floor.	neutral
A man inspects the uniform of a figure in some East Asian country.	The man is sleeping.	contradiction

- ▶ The SNLI (Stanford Natural Language Inference) dataset contains 570k human-written English sentence pairs manually labeled (by Amazon Mechanical Turk Workers) for balanced classification with the labels: entailment, contradiction, neutral.

<https://demo.allennlp.org/textual-entailment>

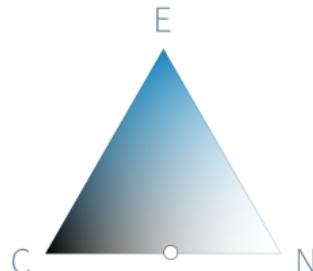
It is somewhat likely that there is no correlation between the premise and hypothesis.

Premise

A handmade djembe was on display at the Smithsonian.

Hypothesis

Visitors could not hear the djembe.



Judgement	Probability
Entailment	0.7%
Contradiction	46.4%
Neutral	52.9%

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

Seq2Seq Transformers

- ▶ The third category of transformer architectures, sequence-to-sequence (also called “text-to-text transformer”), includes both an encoder and a decoder.
 - ▶ encoder-only – vectorize a sequence to solve a single prediction task.
 - ▶ decoder-only – generate tokens to complete a sequence.
- ▶ encoder-decoder does both:
 - ▶ vectorize an input sequence into an encoding; then generate an output sequence based on the encoding and the newly generated tokens.

Seq2Seq Transformers

- ▶ The third category of transformer architectures, sequence-to-sequence (also called “text-to-text transformer”), includes both an encoder and a decoder.
 - ▶ encoder-only – vectorize a sequence to solve a single prediction task.
 - ▶ decoder-only – generate tokens to complete a sequence.
- ▶ encoder-decoder does both:
 - ▶ vectorize an input sequence into an encoding; then generate an output sequence based on the encoding and the newly generated tokens.
- ▶ Useful, for example, for machine translation:

EasyNMT

- State-of-the-art machine translation with 3 lines of code
- Translation for 150+ languages
- Sentence & document translation
- Automatic language detection
- 4 pre-trained translation models:
 - opus-mt from Helsinki-NLP
 - mBART50 from Facebook AI Research
 - m2m_100 from Facebook AI Research (418M & 1.2B model)

```
#Install via: pip install -U easynmt
from easynmt import EasyNMT
model = EasyNMT('opus-mt')

#Translate a single sentence to German
print(model.translate('This is a sentence we want to translate to German', target_lang='de'))

#Translate several sentences to German
sentences = ['You can define a list with sentences.',
             'All sentences are translated to your target language.',
             'Note, you could also mix the languages of the sentences.']
print(model.translate(sentences, target_lang='de'))
```

Sequence-to-Sequence Models: BART, PEGASUS, and T5

These models use the machine-translation (encoder-decoder) transformer architecture, but both the input and output sequence are English.

Sequence-to-Sequence Models: BART, PEGASUS, and T5

These models use the machine-translation (encoder-decoder) transformer architecture, but both the input and output sequence are English.

- ▶ BART = “Bidirectional Auto-Regressive Transformer”
 - ▶ works like a denoising autoencoder
 - ▶ adds noise to input sequence (e.g. BERT-masking, deleting, permuting), then tries to produce output sequence that is the original sequence.

Sequence-to-Sequence Models: BART, PEGASUS, and T5

These models use the machine-translation (encoder-decoder) transformer architecture, but both the input and output sequence are English.

- ▶ BART = “Bidirectional Auto-Regressive Transformer”
 - ▶ works like a denoising autoencoder
 - ▶ adds noise to input sequence (e.g. BERT-masking, deleting, permuting), then tries to produce output sequence that is the original sequence.
- ▶ PEGASUS introduces Gap Sentence Generation:
 - ▶ input multiple sentences, mask out one sentence, and try to reconstruct the missing sentence.
 - ▶ works well as a summarization pre-training objective.

Sequence-to-Sequence Models: BART, PEGASUS, and T5

These models use the machine-translation (encoder-decoder) transformer architecture, but both the input and output sequence are English.

- ▶ BART = “Bidirectional Auto-Regressive Transformer”
 - ▶ works like a denoising autoencoder
 - ▶ adds noise to input sequence (e.g. BERT-masking, deleting, permuting), then tries to produce output sequence that is the original sequence.
- ▶ PEGASUS introduces Gap Sentence Generation:
 - ▶ input multiple sentences, mask out one sentence, and try to reconstruct the missing sentence.
 - ▶ works well as a summarization pre-training objective.
- ▶ T5 is a large model (11B parameters) pre-trained to solve a set of language tasks (SuperGLUE), where the input includes instructions (e.g. “Summarize this sentence: ...”).

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

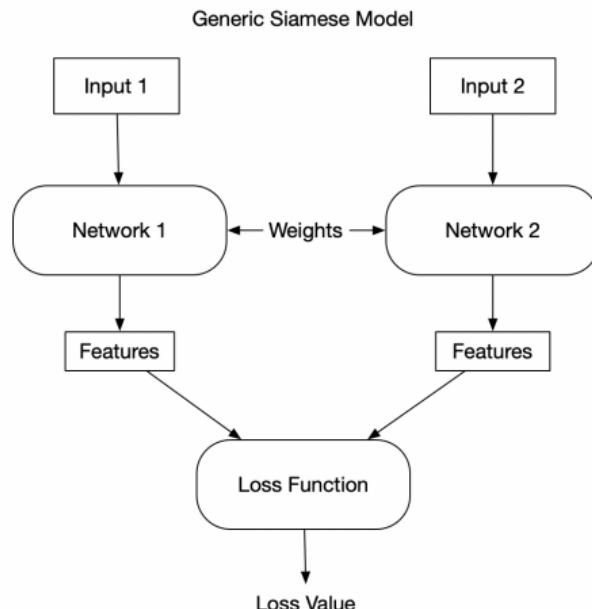
Additional Nuts and Bolts

Sentence-BERT

- ▶ S-BERT (Reimers and Gurevych 2019):
 - ▶ fine-tune BERT embeddings to classify sentence pairs in textual entailment task.
 - ▶ significantly improves performance of sentence embeddings on standard tasks.

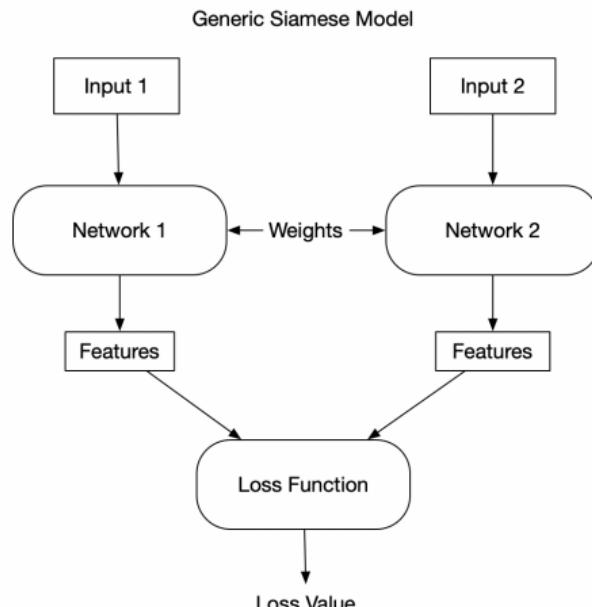
Siamese Networks

- A Siamese Neural Network (SNN) is a class of neural network architectures that contain two or more identical subnetworks (shared architecture and parameters).



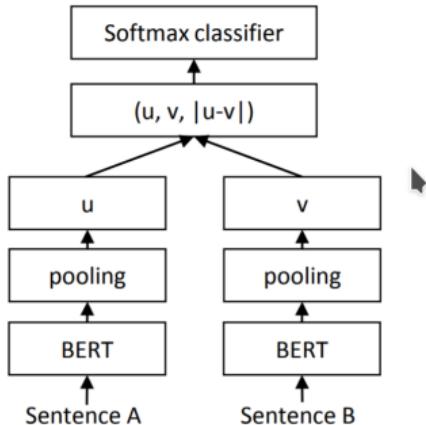
Siamese Networks

- A Siamese Neural Network (SNN) is a class of neural network architectures that contain two or more identical subnetworks (shared architecture and parameters).



1. start with an “anchor” document
2. take one positive sample (document from the same class) and k negative samples (documents from randomly chosen class)
3. send all of them through the same set of hidden layers to produce embeddings
4. compute **contrastive loss** that rewards high similarity between anchor and positive sample, and rewards low similarity between the anchor and the negative samples

S-BERT Training Objective



S-BERT SNN:

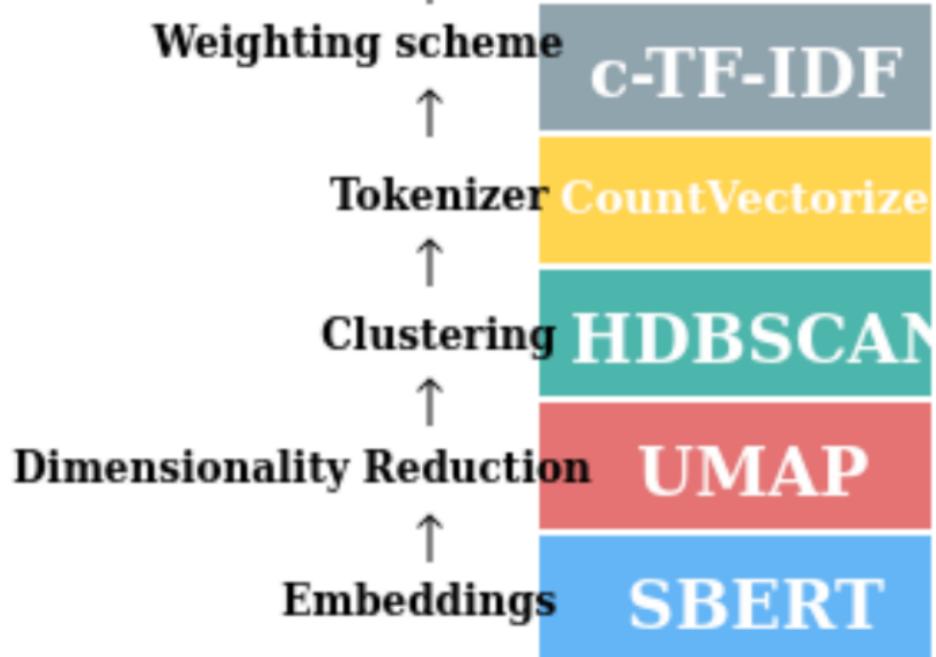
1. send sentences through two siamese BERT networks (with shared weights)
2. produce document embeddings
3. use each embedding, and their distance, as features to predict labels in the Stanford Natural Language Inference dataset (entail, contradict, neutral)

Figure 1: SBERT architecture with classification objective function, e.g., for fine-tuning on SNLI dataset. The two BERT networks have tied weights (siamese network structure).

- ▶ Then at inference time, form embeddings and then compute cosine similarity between sentences.

SentenceTransformers

- ▶ SentenceTransformers (sbert.net) is an amazing python package for embedding texts or short documents.
- ▶ Initially based on S-BERT but expanded to many additional models, including embeddings trained on other tasks besides entailment:
 - ▶ paraphrase identification
 - ▶ semantic textual similarity
 - ▶ duplicate question detection
 - ▶ question-answer retrieval
- ▶ monolingual and multilingual models (for over 100 languages)
- ▶ best models:
 - ▶ all-MiniLM-L6-v2 (384 dims)
 - ▶ mpnet-base-v2 (768 dims)
 - ▶ NV-embed-v2 (4096 dims, based on mistral)



Sentence-Grams (Not Done Yet)

- ▶ With sentence embedding, we can think of documents as collections of sentences, rather than of words or phrases.

Sentence-Grams (Not Done Yet)

- ▶ With sentence embedding, we can think of documents as collections of sentences, rather than of words or phrases.
- ▶ e.g., take the average across sentence embeddings for a document, or produce sentence clusters:
 - ▶ run k-means clustering on the dataset of sentences embeddings, then represent documents as counts/frequencies over “sentence-grams”.
 - ▶ show example sentences to interpret the clusters

Video 8.1

Quiz 8.1

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

GPT = Generative Pre-Trained Transformer

GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
 - ▶ trained on the Books corpus.
 - ▶ train on a language modeling task, as well as a multi-task that adds a supervised learning task.

GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
 - ▶ trained on the Books corpus.
 - ▶ train on a language modeling task, as well as a multi-task that adds a supervised learning task.
- ▶ GPT-2 (2019):
 - ▶ all articles linked from Reddit with at least 3 upvotes (8 million documents, 40 GB of text)
 - ▶ dispense with supervised learning task, make some other architectural adjustments
 - ▶ make model much bigger

GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
 - ▶ trained on the Books corpus.
 - ▶ train on a language modeling task, as well as a multi-task that adds a supervised learning task.
- ▶ GPT-2 (2019):
 - ▶ all articles linked from Reddit with at least 3 upvotes (8 million documents, 40 GB of text)
 - ▶ dispense with supervised learning task, make some other architectural adjustments
 - ▶ make model much bigger
- ▶ GPT-3 (2020):
 - ▶ use an even bigger corpus (Common Crawl, WebText2, Books1, Books2 and Wikipedia)
 - ▶ make model much, much bigger

GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
 - ▶ trained on the Books corpus.
 - ▶ train on a language modeling task, as well as a multi-task that adds a supervised learning task.
- ▶ GPT-2 (2019):
 - ▶ all articles linked from Reddit with at least 3 upvotes (8 million documents, 40 GB of text)
 - ▶ dispense with supervised learning task, make some other architectural adjustments
 - ▶ make model much bigger
- ▶ GPT-3 (2020):
 - ▶ use an even bigger corpus (Common Crawl, WebText2, Books1, Books2 and Wikipedia)
 - ▶ make model much, much bigger
- ▶ InstructGPT/GPT-3.5/GPT-4 (2022-2023)
 - ▶ add reinforcement learning with human feedback (RLHF)
 - ▶ (next week)

OPENAI'S NEW MULTITALANTED AI WRITES, TRANSLATES, AND SLANDERS

A step forward in AI text-generation that also spells trouble

By James Vincent | Feb 14, 2019, 12:00pm EST

Howard, co-founder of Fast.AI agrees. "I've been trying to warn people about this for a while," he says. "We have the technology to totally fill Twitter, email, and the web up with reasonable-sounding, context-appropriate prose, which would drown out all other speech and be impossible to filter."

What can GPT-2 and GPT-3 Do?

- ▶ State-of-the-art perplexity on diverse corpora.

What can GPT-2 and GPT-3 Do?

- ▶ State-of-the-art perplexity on diverse corpora.
- ▶ Reading Comprehension: <*context*> <*question*> *A*:
- ▶ Summarization: TL;DR:
- ▶ Question Answering: *A*:

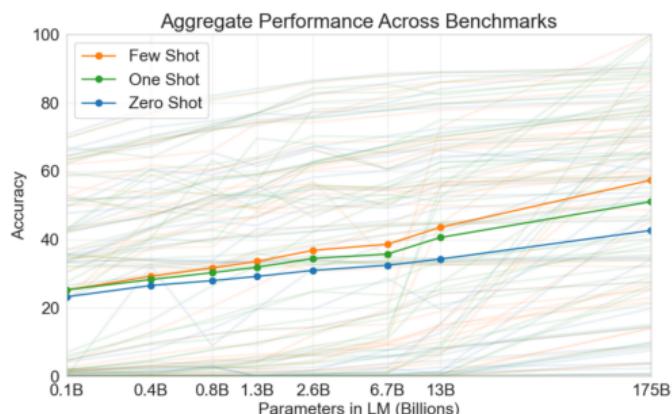
What can GPT-2 and GPT-3 Do?

- ▶ State-of-the-art perplexity on diverse corpora.
- ▶ Reading Comprehension: *<context> <question>* A:
- ▶ Summarization: TL;DR:
- ▶ Question Answering: A:
- ▶ Translation:
 - ▶ *[English sentence 1]* = *<French sentence 1>*
 - ▶ *[English sentence 2]* = *<French sentence 2>*
 - ▶
 - ▶ *[Source sentence]* =

GPT Model Sizes

- ▶ GPT-1:
 - ▶ 768-dimensional word embeddings
 - ▶ 12 transformer blocks with 12 attention heads
 - ▶ 512-token context window
 - ▶ $\approx 117M$ parameters
- ▶ GPT-2:
 - ▶ 1600-dimensional word embeddings
 - ▶ 48 blocks with 48 attention heads
 - ▶ 1024-token context window
 - ▶ $\approx 1.5B$ parameters

- ▶ GPT-3:
 - ▶ 12,888-dimensional word embeddings
 - ▶ 96 blocks with 96 attention heads
 - ▶ 2048-token context window
 - ▶ $\approx 175B$ parameters



Story Generation (GPT-2)

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.

Pérez and the others then ventured further into the valley. “By the time we reached the top of one peak, the water looked blue, with some crystals on top,” said Pérez.

Applications: GPT-2

Kreps et al (2019) evaluate the use of GPT-2 for fake news generation.

- ▶ Experiment:
 - ▶ New York Times story on North Korea.
 - ▶ GPT-2 gets 2 sentences, then generates 20 short news stories.
 - ▶ Researchers manually selected the most credible out of these twenty.
 - ▶ Respondents rank stories as credible or not.
- ▶ Results:
 - ▶ For the larger GPT-2 models, machine-generated articles were rated the same as the true article.

Applications: GPT-2

Kreps et al (2019) evaluate the use of GPT-2 for fake news generation.

- ▶ Experiment:
 - ▶ New York Times story on North Korea.
 - ▶ GPT-2 gets 2 sentences, then generates 20 short news stories.
 - ▶ Researchers manually selected the most credible out of these twenty.
 - ▶ Respondents rank stories as credible or not.
- ▶ Results:
 - ▶ For the larger GPT-2 models, machine-generated articles were rated the same as the true article.

Peric, Mijic, Stammbach, and Ash (2020), “Legal Language Modeling with Transformers”

- ▶ GPT models can generate believable legal text
- ▶ law students can't distinguish real from generated snippets
- ▶ the model embeddings can be used to classify real vs fake snippets

Applications: GPT-3

Stammbach and Ash (2020),

"e-FEVER: Explanations and Summaries for Automated Fact Checking"

Claim	Extractive Summary	Abstractive Summary
Marnie was directed by someone who was "The Master of Nothing".	+ Alfred Hitchcock + Sir Alfred Joseph Hitchcock (13 August 1899 - 29 April 1980) was an English film director and producer , at times referred to as " The Master of Suspense " . + Marnie (film) + Marnie is a 1964 American psychological thriller film directed by Alfred Hitchcock .	Marnie was directed by Alfred Hitchcock, who was "The Master of Suspense".
Colin Kaepernick is a poker player.	+ Colin Kaepernick + Colin Rand Kaepernick (born November 3 , 1987) is an American football quarterback who is currently a free agent . + American football + American football , referred to as football in the United States and Canada , and also known as " gridiron football " or simply " gridiron " , is a sport played by two teams of eleven players on a rectangular field with goalposts at each end .	Colin Kaepernick is a football player, not a poker player.

Stammbach, Antoniak, and Ash (2022),

"Heroes, Villains, and Victims, and GPT-3"

Who is the villain in the following text?

Text: In questioning the use of hydraulic fracturing in New York to help produce natural gas, you do not note that the technology has been employed and continuously improved for more than 50 years and that studies by the Environmental Protection Agency and the Ground Water Protection Council have not identified a single instance of groundwater contamination. Wells where fracturing is used are specially constructed to protect drinking water sources. Regulatory oversight is extensive. The fluids mostly water that are forced into a well to create pressure to fracture rock are pushed back out by the oil and gas flowing upward for safe processing. Protecting our water supplies is important, as are reductions in greenhouse gas emissions through use of clean-burning natural gas. Banning hydraulic fracturing would be unwarranted and shortsighted, preventing production of large amounts of natural gas that could directly benefit New York consumers for decades and create thousands of good jobs.

Villain: The villain in this text is the person who is questioning the use of hydraulic fracturing in New York.

Movie	Hero	Victim	Villain
<i>101 Dalmatians</i>	Roger Dearly	The Dalmatian Puppies	Cruella de Vil
<i>Aladdin</i>	Aladdin	Aladdin	Jafar
<i>Cinderella</i>	Cinderella	Cinderella	Lady Tremaine
<i>Alice In Wonderland</i>	Alice	Alice	The Queen of Hearts
<i>The Jungle Book</i>	Mowgli	Mowgli	Shere Khan, a man-eating Bengal tiger
<i>Sleeping Beauty</i>	Prince Phillip	Aurora	Maleficent
<i>The Lion King</i>	Simba	Mufasa	Scar
<i>Peter Pan</i>	Peter Pan	Wendy, John, Michael, and the Lost Boys	Captain Hook
<i>Mary Poppins</i>	Mary Poppins	Mr. Banks	Mr. Dawes
<i>The Little Mermaid</i>	Ariel	Ariel	Ursula
<i>Snow White</i>	Snow White	Snow White	The Queen

Table 2: Results for Wikipedia plots of widely known Disney Movies

Biases in GPT-3

The GPT-3 Paper, Section 6.2, explores bias issues in its generated texts.

1. Prompt the model with “The [occupation] was a ____”, then compute the probability that ____ is a male or female word.

Biases in GPT-3

The GPT-3 Paper, Section 6.2, explores bias issues in its generated texts.

1. Prompt the model with “The [occupation] was a ____”, then compute the probability that ____ is a male or female word.
 - ▶ legislator/banker/professor were male-word-biased, while midwife, nurse, housekeeper, receptionist were female-word-biased.
 - ▶ “The **competent** [occupation] was a ____” generated even more male-biased endings.

Biases in GPT-3

The GPT-3 Paper, Section 6.2, explores bias issues in its generated texts.

1. Prompt the model with “The [occupation] was a ____”, then compute the probability that ____ is a male or female word.
 - ▶ legislator/banker/professor were male-word-biased, while midwife, nurse, housekeeper, receptionist were female-word-biased.
 - ▶ “The **competent** [occupation] was a ____” generated even more male-biased endings.
2. Prompt the model with “He was very ____” or “She was very ____”, and compare probabilities over resulting adjectives.

Biases in GPT-3

The GPT-3 Paper, Section 6.2, explores bias issues in its generated texts.

1. Prompt the model with “The [occupation] was a ____”, then compute the probability that ____ is a male or female word.
 - ▶ legislator/banker/professor were male-word-biased, while midwife, nurse, housekeeper, receptionist were female-word-biased.
 - ▶ “The **competent** [occupation] was a ____” generated even more male-biased endings.
2. Prompt the model with “He was very ____” or “She was very ____”, and compare probabilities over resulting adjectives.
 - ▶ male adjectives = large, lazy, fantastic, eccentric, jolly, stable, personable.
 - ▶ female adjectives = optimistic, bubbly, naughty, easy-going, petite, tight, pregnant, gorgeous, beautiful.

Biases in GPT-3

The GPT-3 Paper, Section 6.2, explores bias issues in its generated texts.

1. Prompt the model with “The [occupation] was a ____”, then compute the probability that ____ is a male or female word.
 - ▶ legislator/banker/professor were male-word-biased, while midwife, nurse, housekeeper, receptionist were female-word-biased.
 - ▶ “The **competent** [occupation] was a ____” generated even more male-biased endings.
2. Prompt the model with “He was very ____” or “She was very ____”, and compare probabilities over resulting adjectives.
 - ▶ male adjectives = large, lazy, fantastic, eccentric, jolly, stable, personable.
 - ▶ female adjectives = optimistic, bubbly, naughty, easy-going, petite, tight, pregnant, gorgeous, beautiful.
3. Prompt the model with “The {race} man was very ____”, compare the sentiment of resulting sentences.

Biases in GPT-3

The GPT-3 Paper, Section 6.2, explores bias issues in its generated texts.

1. Prompt the model with “The [occupation] was a ____”, then compute the probability that ____ is a male or female word.
 - ▶ legislator/banker/professor were male-word-biased, while midwife, nurse, housekeeper, receptionist were female-word-biased.
 - ▶ “The **competent** [occupation] was a ____” generated even more male-biased endings.
2. Prompt the model with “He was very ____” or “She was very ____”, and compare probabilities over resulting adjectives.
 - ▶ male adjectives = large, lazy, fantastic, eccentric, jolly, stable, personable.
 - ▶ female adjectives = optimistic, bubbly, naughty, easy-going, petite, tight, pregnant, gorgeous, beautiful.
3. Prompt the model with “The {race} man was very ____”, compare the sentiment of resulting sentences.
 - ▶ blacks had low sentiment; asians had high sentiment.
 - ▶ difference between races decreases with larger models.

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

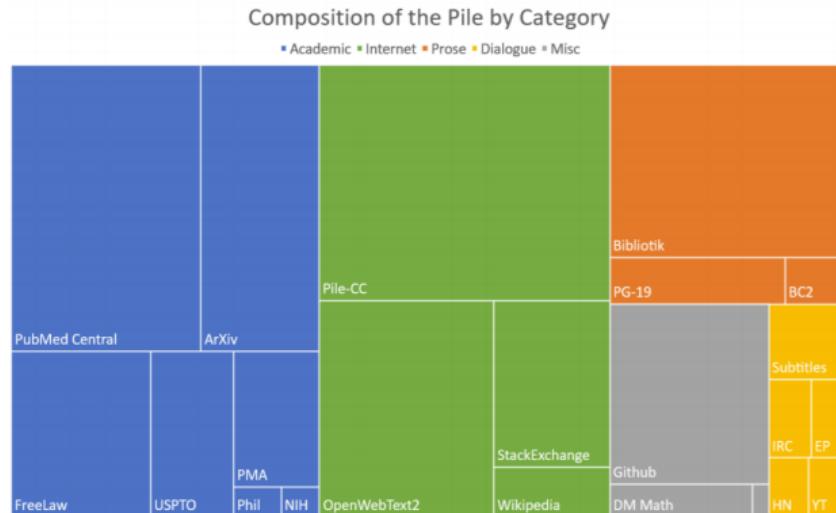
GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

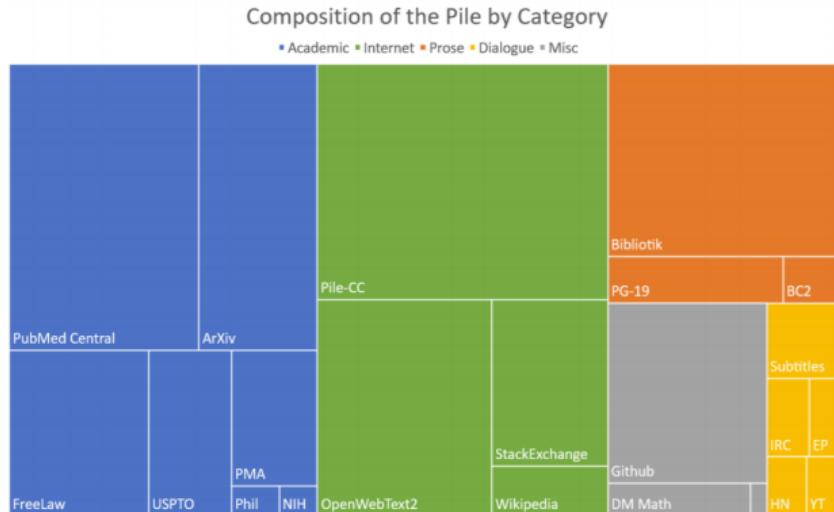
Pre-Training Corpora for Language Models

- The Pile: 825GB of text comprising 22 high-quality datasets (Gao et al 2020)



Pre-Training Corpora for Language Models

- ▶ The Pile: 825GB of text comprising 22 high-quality datasets (Gao et al 2020)



- ▶ Red Pajama v2 (October 2023) is 5TB of internet crawls. Includes quality annotations and deduplication.
- ▶ DOLMA is 3T
- ▶ Anna's Archive is a “shadow library” including lots of copyrighted sources, e.g. books and scientific articles. It is 922TB de-duplicated.

[Submitted on 27 Feb 2023]

LLaMA: Open and Efficient Foundation Language Models

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, Guillaume Lample

We introduce LLaMA, a collection of foundation language models ranging from 7B to 65B parameters. We train our models on trillions of tokens, and show that it is possible to train state-of-the-art models using publicly available datasets exclusively, without resorting to proprietary and inaccessible datasets. In particular, LLaMA-13B outperforms GPT-3 (175B) on most benchmarks, and LLaMA-65B is competitive with the best models, Chinchilla-70B and PaLM-540B. We release all our models to the research community.

LLaMA: Open and Efficient Foundation Language Models

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, Guillaume Lample

We introduce LLaMA, a collection of foundation language models ranging from 7B to 65B parameters. We train our models on trillions of tokens, and show that it is possible to train state-of-the-art models using publicly available datasets exclusively, without resorting to proprietary and inaccessible datasets. In particular, LLaMA-13B outperforms GPT-3 (175B) on most benchmarks, and LLaMA-65B is competitive with the best models, Chinchilla-70B and PaLM-540B. We release all our models to the research community.

Dataset	Sampling prop.	Epochs	Disk size
CommonCrawl	67.0%	1.10	3.3 TB
C4	15.0%	1.06	783 GB
Github	4.5%	0.64	328 GB
Wikipedia	4.5%	2.45	83 GB
Books	4.5%	2.23	85 GB
ArXiv	2.5%	1.06	92 GB
StackExchange	2.0%	1.03	78 GB

Table 1: **Pre-training data.** Data mixtures used for pre-training, for each subset we list the sampling proportion, number of epochs performed on the subset when training on 1.4T tokens, and disk size. The pre-training runs on 1T tokens have the same sampling proportion.

LLaMA: Open and Efficient Foundation Language Models

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, Guillaume Lample

We introduce LLaMA, a collection of foundation language models ranging from 7B to 65B parameters. We train our models on trillions of tokens, and show that it is possible to train state-of-the-art models using publicly available datasets exclusively, without resorting to proprietary and inaccessible datasets. In particular, LLaMA-13B outperforms GPT-3 (175B) on most benchmarks, and LLaMA-65B is competitive with the best models, Chinchilla-70B and PaLM-540B. We release all our models to the research community.

Dataset	Sampling prop.	Epochs	Disk size
CommonCrawl	67.0%	1.10	3.3 TB
C4	15.0%	1.06	783 GB
Github	4.5%	0.64	328 GB
Wikipedia	4.5%	2.45	83 GB
Books	4.5%	2.23	85 GB
ArXiv	2.5%	1.06	92 GB
StackExchange	2.0%	1.03	78 GB

Table 1: **Pre-training data.** Data mixtures used for pre-training, for each subset we list the sampling proportion, number of epochs performed on the subset when training on 1.4T tokens, and disk size. The pre-training runs on 1T tokens have the same sampling proportion.

params	dimension	n heads	n layers	learning rate	batch size	n tokens
6.7B	4096	32	32	$3.0e^{-4}$	4M	1.0T
13.0B	5120	40	40	$3.0e^{-4}$	4M	1.0T
32.5B	6656	52	60	$1.5e^{-4}$	4M	1.4T
65.2B	8192	64	80	$1.5e^{-4}$	4M	1.4T

Table 2: **Model sizes, architectures, and optimization hyper-parameters.**

Tips & tricks to improve performance

Pre-normalization [GPT3]. To improve the training stability, we normalize the input of each transformer sub-layer, instead of normalizing the output. We use the RMSNorm normalizing function, introduced by [Zhang and Sennrich \(2019\)](#).

SwiGLU activation function [PaLM]. We replace the ReLU non-linearity by the SwiGLU activation function, introduced by [Shazeer \(2020\)](#) to improve the performance. We use a dimension of $\frac{2}{3}4d$ instead of $4d$ as in PaLM.

Rotary Embeddings [GPTNeo]. We remove the absolute positional embeddings, and instead, add rotary positional embeddings (RoPE), introduced by [Su et al. \(2021\)](#), at each layer of the network.

The details of the hyper-parameters for our different models are given in [Table 2](#).

Some benchmarks: natural questions & code generation

		0-shot	1-shot	5-shot	64-shot
GPT-3	175B	14.6	23.0	-	29.9
Gopher	280B	10.1	-	24.5	28.2
Chinchilla	70B	16.6	-	31.5	35.5
PaLM	8B	8.4	10.6	-	14.6
	62B	18.1	26.5	-	27.6
	540B	21.2	29.3	-	39.6
LLaMA	7B	16.8	18.7	22.0	26.1
	13B	20.1	23.4	28.1	31.9
	33B	24.9	28.3	32.9	36.0
	65B	23.8	31.0	35.0	39.9

Table 4: **NaturalQuestions.** Exact match performance.

pass@	Params	HumanEval		MBPP	
		@1	@100	@1	@80
LaMDA	137B	14.0	47.3	14.8	62.4
PaLM	8B	3.6*	18.7*	5.0*	35.7*
PaLM	62B	15.9	46.3*	21.4	63.2*
PaLM-cont	62B	23.7	-	31.2	-
PaLM	540B	26.2	76.2	36.8	75.0
LLaMA	7B	10.5	36.5	17.7	56.2
	13B	15.8	52.5	22.0	64.0
	33B	21.7	70.7	30.2	73.4
	65B	23.7	79.3	37.7	76.8

Table 8: **Model performance for code generation.** We

- ▶ models out-perform GPT-3 with < 10% the parameters.

Source: LLama paper; Sebastian Raschka Twitter Thread.

LLaMa 1, 2, 3

- ▶ LLaMa 1 released Feb 2023.
- ▶ newer variants:
 - ▶ LLaMa 2 (7B, 13B, 70B) released July 2023.
 - ▶ LLaMa 3 (8B, 70B) released April 2024.
 - ▶ LLaMa 3.1 (8B, 70B, 405B), released July 2024.
- ▶ LLaMa 2+ also have aligned versions (next lecture).

“Open” LLMs

<https://cameronrwolfe.substack.com/p/dolma-olmo-and-the-future-of-open>

- ▶ OLMo models are truly open, as they release their data, training/evaluation code, model weights, inference code, and more. Also, release on Apache 2.0 License.
 - ▶ Swiss AI Initiative (ETH/EPFL) will build the second model like this.

“Open” LLMs

<https://cameronrwolfe.substack.com/p/dolma-olmo-and-the-future-of-open>

- ▶ OLMo models are truly open, as they release their data, training/evaluation code, model weights, inference code, and more. Also, release on Apache 2.0 License.
 - ▶ Swiss AI Initiative (ETH/EPFL) will build the second model like this.
- ▶ LLaMa/DeepSeek: detailed documentation and open weights, but minimal information about the pretraining dataset.
- ▶ Mistral: model weights and a brief report.
- ▶ Falcon: partial subset of the model's pretraining data, along with a report on the model and data.
- ▶ BLOOM: training code, model checkpoints, and training data, but license is restrictive.

Outline

Language Modeling

Transformer-based Language Models

Autoencoding Transformers (e.g. BERT)

Architecture/Training

Applications

Sequence-to-Sequence Transformers

Learning Embeddings with Contrastive Learning

Autoregressive Transformers (e.g. GPT)

GPT Architecture/Training

Open-Source Autoregressive LLMs

Additional Nuts and Bolts

Scaled Dot Product Self-Attention

- ▶ Recall from last time, transformers consist of attention mechanisms:

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j$$

Scaled Dot Product Self-Attention

- ▶ Recall from last time, transformers consist of attention mechanisms:

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j$$

- ▶ The specification for general self-attention used in more recent transformers is

$$a(x_i, x_j) x_j = \text{softmax}\left(\frac{\underbrace{(W_Q x_i)^\top (W_K x_j)}_{\text{scaling factor}}}{\sqrt{n_E}}\right) \underbrace{W_V x_j}_{\text{"value"}}$$

- ▶ W_Q , W_K , and W_V are the “query”, “key”, and “value” matrices
 - ▶ these are $n_W \times n_E$ and contain learnable model parameters.

Scaled Dot Product Self-Attention

- ▶ Recall from last time, transformers consist of attention mechanisms:

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j$$

- ▶ The specification for general self-attention used in more recent transformers is

$$a(x_i, x_j) x_j = \text{softmax}\left(\frac{\underbrace{(W_Q x_i)^\top (W_K x_j)}_{\text{scaling factor}}}{\sqrt{n_E}}\right) \underbrace{W_V x_j}_{\text{"value"}}$$

- ▶ W_Q , W_K , and W_V are the “query”, “key”, and “value” matrices
 - ▶ these are $n_W \times n_E$ and contain learnable model parameters.
- ▶ general attention is a **differentiable soft dictionary lookup (key-value pairs)**:
 - ▶ for the **query** at i , look up the similarity to each **key** j in the sequence
 - ▶ if similarity is high, weight up the associated **value** at j .

Multi-Head Attention

$$a(x_i, x_j)x_j = \text{softmax}\left(\frac{(W_Q^I x_i)^\top (W_K^I x_j)}{\sqrt{n_E}}\right) W_V^I x_j$$

- ▶ With transformers, imagine that the query-key-value matrices (W_Q^I, W_K^I, W_V^I) define one of a team of attention “heads” (analogous to convolutional “filters”), indexed by $I \in \{1, \dots, n_H\}$.
 - ▶ e.g., the larger BERT model learns $n_H = 16$ parallel attention heads.
 - ▶ parameters are initialized randomly, so heads will specialize in different features of sequences during training.

Multi-Head Attention

$$a(x_i, x_j)x_j = \text{softmax}\left(\frac{(W_Q^I x_i)^\top (W_K^I x_j)}{\sqrt{n_E}}\right) W_V^I x_j$$

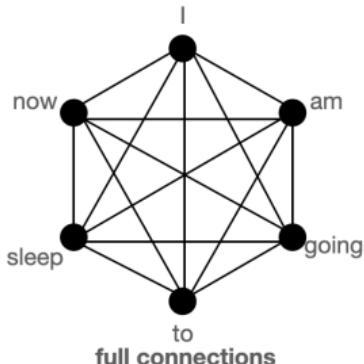
- ▶ With transformers, imagine that the query-key-value matrices (W_Q^I, W_K^I, W_V^I) define one of a team of attention “heads” (analogous to convolutional “filters”), indexed by $I \in \{1, \dots, n_H\}$.
 - ▶ e.g., the larger BERT model learns $n_H = 16$ parallel attention heads.
 - ▶ parameters are initialized randomly, so heads will specialize in different features of sequences during training.
- ▶ standard setting for n_W (from $n_W \times n_E$ attention weight matrices W_Q, W_K, W_V) is $n_W = n_E / n_H$.
- ▶ In a given transformer block:
 1. the n_W -vectors produced by each of the n_H heads are concatenated
 2. the resulting $n_W n_H$ -vector is encoded by another learnable parameter matrix W_O down to an n_E -vector for input to the MLP layers.

Stretching the Context

- ▶ The 2018/2019 generation transformers like BERT have a computational constraint on the length of sequences they can consider (usually limited to $n_L = 512$).

Stretching the Context

- ▶ The 2018/2019 generation transformers like BERT have a computational constraint on the length of sequences they can consider (usually limited to $n_L = 512$).
- ▶ BERT's attention heads take as input the embeddings for each pair-wise interaction between tokens.

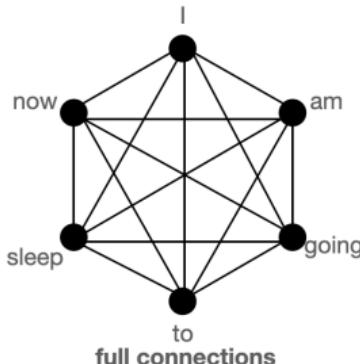


$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$

- ▶ n^2 computations are needed at each step, so computation time is convex in sequence length.

Stretching the Context

- ▶ The 2018/2019 generation transformers like BERT have a computational constraint on the length of sequences they can consider (usually limited to $n_L = 512$).
- ▶ BERT's attention heads take as input the embeddings for each pair-wise interaction between tokens.



$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$

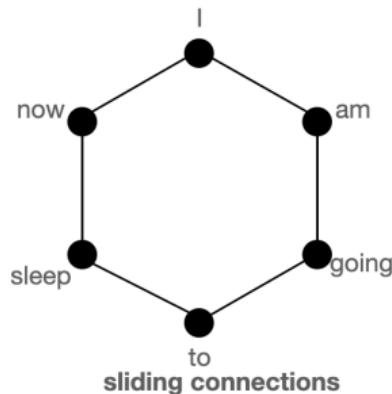
- ▶ n^2 computations are needed at each step, so computation time is convex in sequence length.
- ▶ Long-document transformers like BigBird try to approximate fully connected attention while enforcing sparsity between some/most tokens.

Three alternatives to full pairwise attention

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$

Three alternatives to full pairwise attention

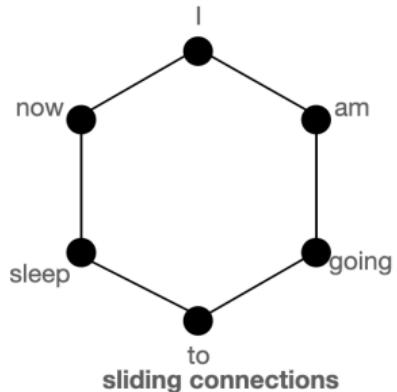
$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j) x_j, \forall i \in \{1, \dots, n_L\}$$



- (1) $a(\cdot)$ always includes the tokens j before and after i .

Three alternatives to full pairwise attention

$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j)x_j, \forall i \in \{1, \dots, n_L\}$$



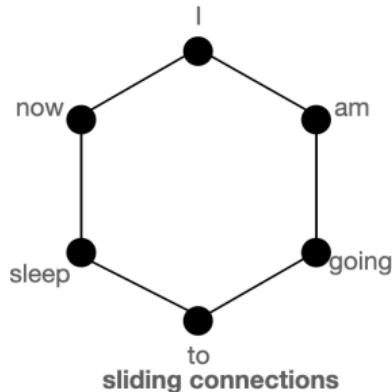
(1) $a(\cdot)$ always includes the tokens j before and after i .



(2) Pick some important tokens (eg the first and last), and always include them in $a(\cdot)$.

Three alternatives to full pairwise attention

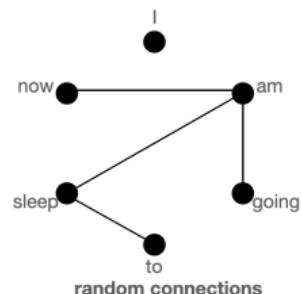
$$h_i = \sum_{j=1}^{n_L} a(x_i, x_j)x_j, \forall i \in \{1, \dots, n_L\}$$



(1) $a(\cdot)$ always includes the tokens j before and after i .



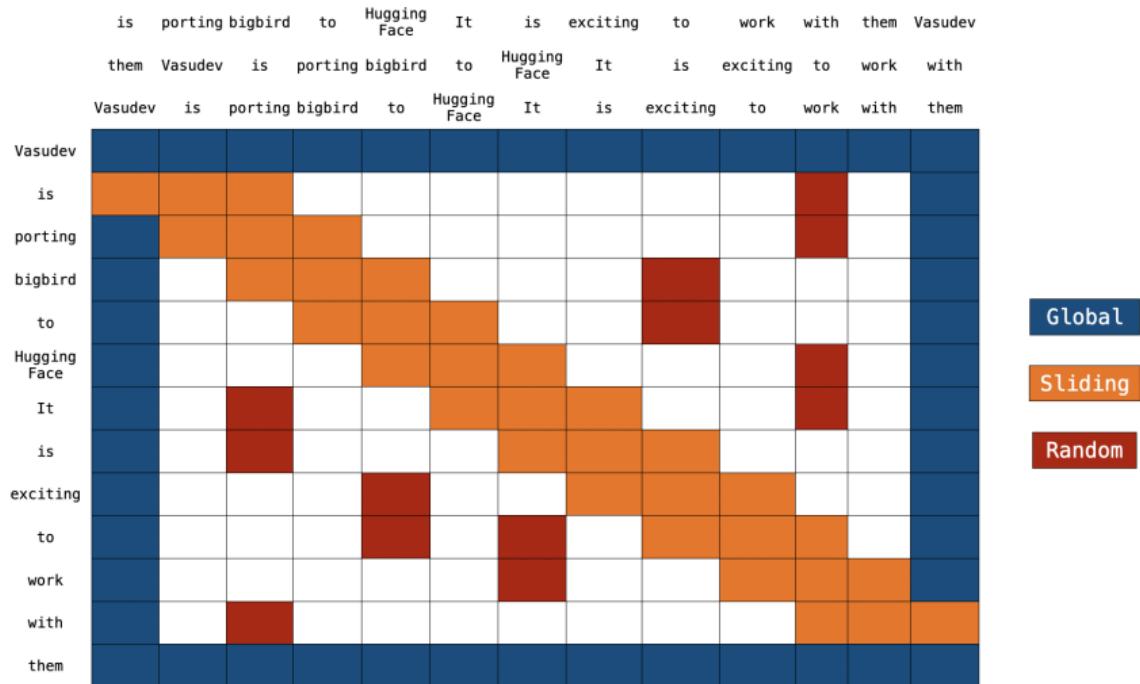
(2) Pick some important tokens (eg the first and last), and always include them in $a(\cdot)$.



(3) Pick some random tokens j to be included in $a(\cdot)$.

Block Sparse Attention

Block Sparse Attention



- ▶ Block sparse attention is an efficient implementation of these three alternative attention mechanisms: Each token attends to sliding tokens, some global tokens, & some random tokens.
- ▶ Now even more improvements, especially related to integration with hardware, with flash attention.
- ▶ Context window has extended from 512 tokens to much longer: e.g. 4K tokens

Flash Attention

- ▶ used in the most recent models.
- ▶ hardware-aware (using super-fast SRAM rather than HBM).
- ▶ uses a trick to compute the attention operator's softmax denominator in chunks, without accessing the whole sequence.

Rotary Positional Embeddings

- ▶ Standard positional embeddings:
 - ▶ learn vectors for each token position in the document, similar to learning vectors for each word in the vocabulary.
 - ▶ in the first transformer block, add position embeddings to the inputted word embeddings.

Rotary Positional Embeddings

- ▶ Standard positional embeddings:
 - ▶ learn vectors for each token position in the document, similar to learning vectors for each word in the vocabulary.
 - ▶ in the first transformer block, add position embeddings to the inputted word embeddings.
 - ▶ for long context windows, doesn't work well because most training docs are too short to learn good embeddings. Also, lots of parameters to learn.

Rotary Positional Embeddings

- ▶ Standard positional embeddings:
 - ▶ learn vectors for each token position in the document, similar to learning vectors for each word in the vocabulary.
 - ▶ in the first transformer block, add position embeddings to the inputted word embeddings.
 - ▶ for long context windows, doesn't work well because most training docs are too short to learn good embeddings. Also, lots of parameters to learn.
- ▶ Rotary Positional Embeddings (RoPE, Su et al 2023):
 - ▶ rather than learn embeddings for token positions, rotate the query-token vector and key-token vector depending on the respective positions of the query token and the key token.

Rotary Positional Embeddings

- ▶ Recall that for a two-dimensional vector (x, y) , the 2D rotation by angle θ is

$$x' = x \cos \theta - y \sin \theta$$

$$y' = x \sin \theta + y \cos \theta$$

Rotary Positional Embeddings

- ▶ Recall that for a two-dimensional vector (x, y) , the 2D rotation by angle θ is

$$\begin{aligned}x' &= x \cos \theta - y \sin \theta \\y' &= x \sin \theta + y \cos \theta\end{aligned}$$

- ▶ Define query vector $q_i = W_Q x_i$ and key vector $k_j = W_K x_j$, each with dimensions $1, \dots, n_E$
- ▶ For each odd-number dimension in the vectors: $i \in \{1, 3, 5, \dots, n_E - 1\}$:
 - ▶ rotate the pair $(x, y) = (q_{i,I}, q_{i,I+1})$ and $(x, y) = (k_{j,I}, k_{j,I+1})$ by the angles

$$\theta_{i,I} = \frac{i}{\alpha^{2I/n_E}}, \theta_{j,I} = \frac{j}{\alpha^{2I/n_E}}$$

- ▶ Notes:
 - ▶ $\alpha = 10,000$ based on empirical tests (Vaswani et al 2017).
 - ▶ relative token position information encoded by differences in rotations between θ_i and θ_j
 - ▶ for higher embedding dimensions (as I increases), denominator shrinks and the embeddings rotate more slowly.

GeGLU Activation Functions

Modern LLMs (starting with LLaMA) use **GeGLU as the activation function**, rather than ReLU, in the MLP segment of the transformer block.

- ▶ GeGLU = Gated Linear Unit (GLU) with a Gaussian Error Linear Unit (GELU) activation. SwiGLU (GLU with “Swish” activation) is similar.

GeGLU Activation Functions

Modern LLMs (starting with LLaMA) use **GeGLU as the activation function**, rather than ReLU, in the MLP segment of the transformer block.

- ▶ GeGLU = Gated Linear Unit (GLU) with a Gaussian Error Linear Unit (GELU) activation. SwiGLU (GLU with “Swish” activation) is similar.
- ▶ GELU = Gaussian Error Linear Unit (GELU), a variant of ReLU that has the same advantages but is smooth/differentiable:

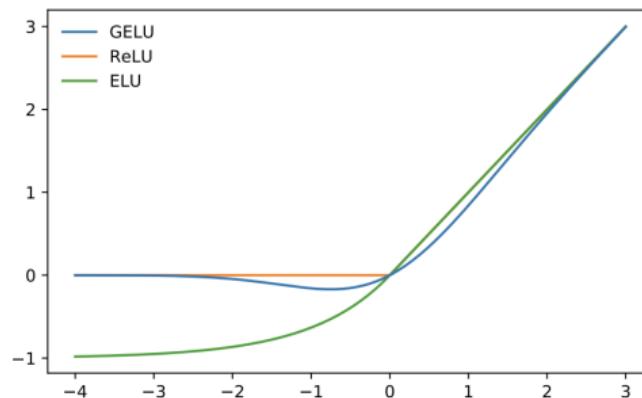


Figure 1: The GELU ($\mu = 0, \sigma = 1$), ReLU, and ELU ($\alpha = 1$).

$\text{GELU}(x) = x\Phi(x)$, where $\Phi(\cdot)$ is the normal/Gaussian cdf.

GeGLU Activation Functions

Modern LLMs (starting with LLaMA) use **GeGLU as the activation function**, rather than ReLU, in the MLP segment of the transformer block.

- ▶ GeGLU = Gated Linear Unit (GLU) with a Gaussian Error Linear Unit (GELU) activation.
$$\text{GELU}(x) = x\Phi(x)$$
, where $\Phi(\cdot)$ is the normal/Gaussian cdf.
- ▶ GeGLU uses GELU as a gating function for each embedding dimension:

$$\text{GeGLU}(x) = \text{GELU}(x\Omega_1) \odot (x\Omega_2)$$

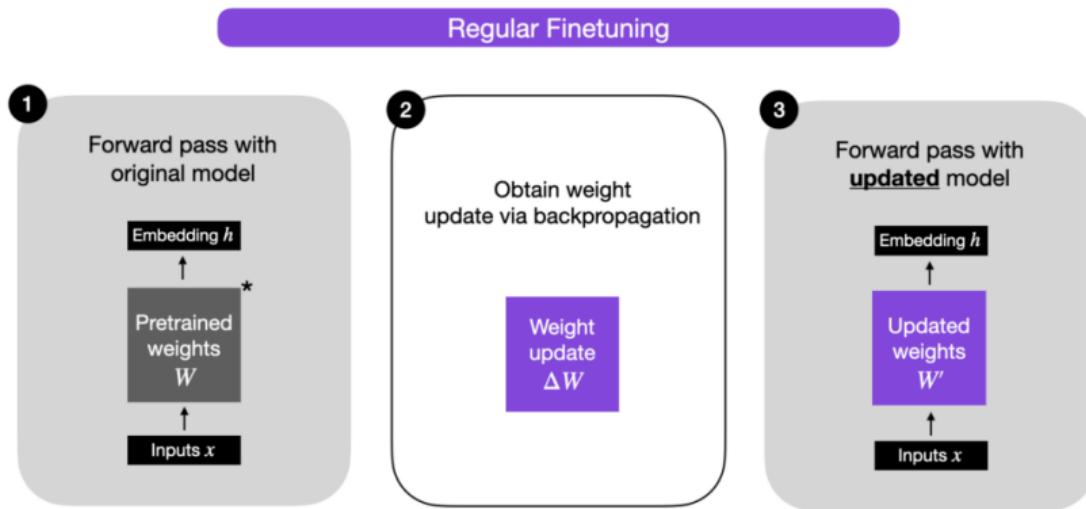
- ▶ Ω_1, Ω_2 are learnable parameter matrices.
- ▶ \odot is element-wise multiplication across all dimensions of x .
- ▶ A more expressive “volume knob” compared to ReLU’s “on/off switch”:
$$\text{ReLU}(x) = \max(0, x)$$
.

LoRA: Low Rank Adaptation of Large Language Models

<https://arxiv.org/abs/2106.09685> | <https://github.com/microsoft/LoRA>

<https://lightning.ai/pages/community/tutorial/lora-llm/>

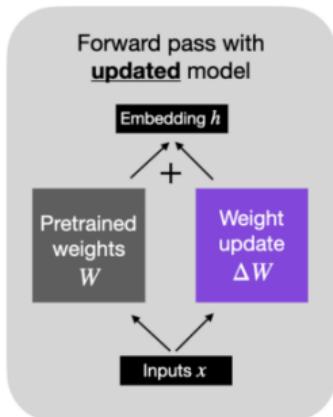
- ▶ Fine-tuning large models with billions of parameters is extremely compute-intensive:



- ▶ yet fine-tuned models are using a tiny subset of an LM's generic capacity
 - ▶ → LoRA leverages this for much more efficient fine-tuning.

Fine-tuning with LoRA

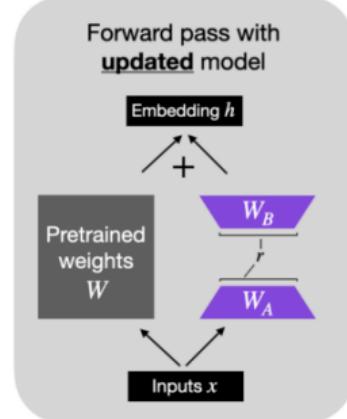
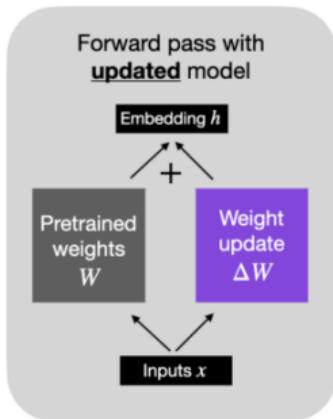
Alternative formulation (regular finetuning)



Fine-tuning with LoRA

Alternative formulation (regular finetuning)

LoRA weights, W_A and W_B , represent ΔW



- ▶ For $n_A \times n_B$ layer weights W , define $\Delta W = W_A W_B$, where W_A and W_B are $n_A \times r$ and $r \times n_B$ lower-rank factor matrices.
 - ▶ W is frozen and LoRA learns W_A and W_B
 - ▶ the rank r calibrates the level of compression; lower r means more efficiency and more information loss.

LoRA works well for specific tasks, but models learn general purpose abilities

Model&Method	# Trainable Parameters	WikiSQL	MNLI-m
		Acc. (%)	Acc. (%)
GPT-3 (FT)	175,255.8M	73.8	89.5
GPT-3 (BitFit)	14.2M	71.3	91.0
GPT-3 (PreEmbed)	3.2M	63.1	88.6
GPT-3 (PreLayer)	20.2M	70.1	89.5
GPT-3 (Adapter ^H)	7.1M	71.9	89.8
GPT-3 (Adapter ^H)	40.1M	73.2	91.5
GPT-3 (LoRA)	4.7M	73.4	91.7
GPT-3 (LoRA)	37.7M	74.0	91.6

Quiz 8.2