

Pre-Training Strategies

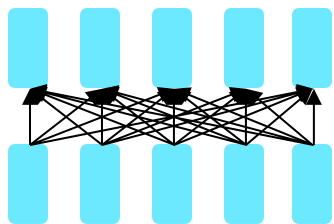
Encoder-Decoder and Decoder-only Models

Tanmoy Chakraborty

Associate Professor, IIT Delhi

<https://tanmoychak.com/>

Pre-Training for Different Types of Architectures

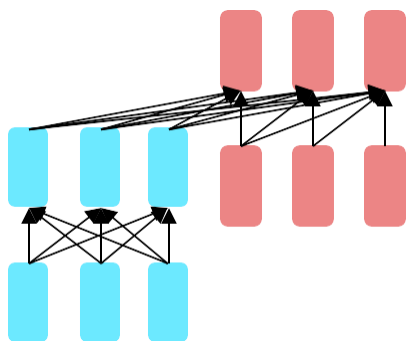


**Encoder-
only**

BERT

(already discussed)

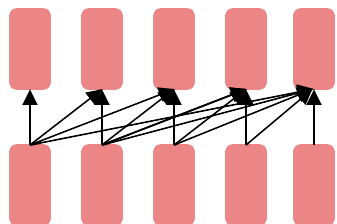
- Gets bi-directional context – can condition on future!
- How do we train them to build strong representations?



**Encoder-
Decoder**

BART, T5

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



**Decoder-
only**

GPT, Llama

- Language models!
- Nice to generate from; can't condition on future words

Pre-Training Encoder-Decoder Models

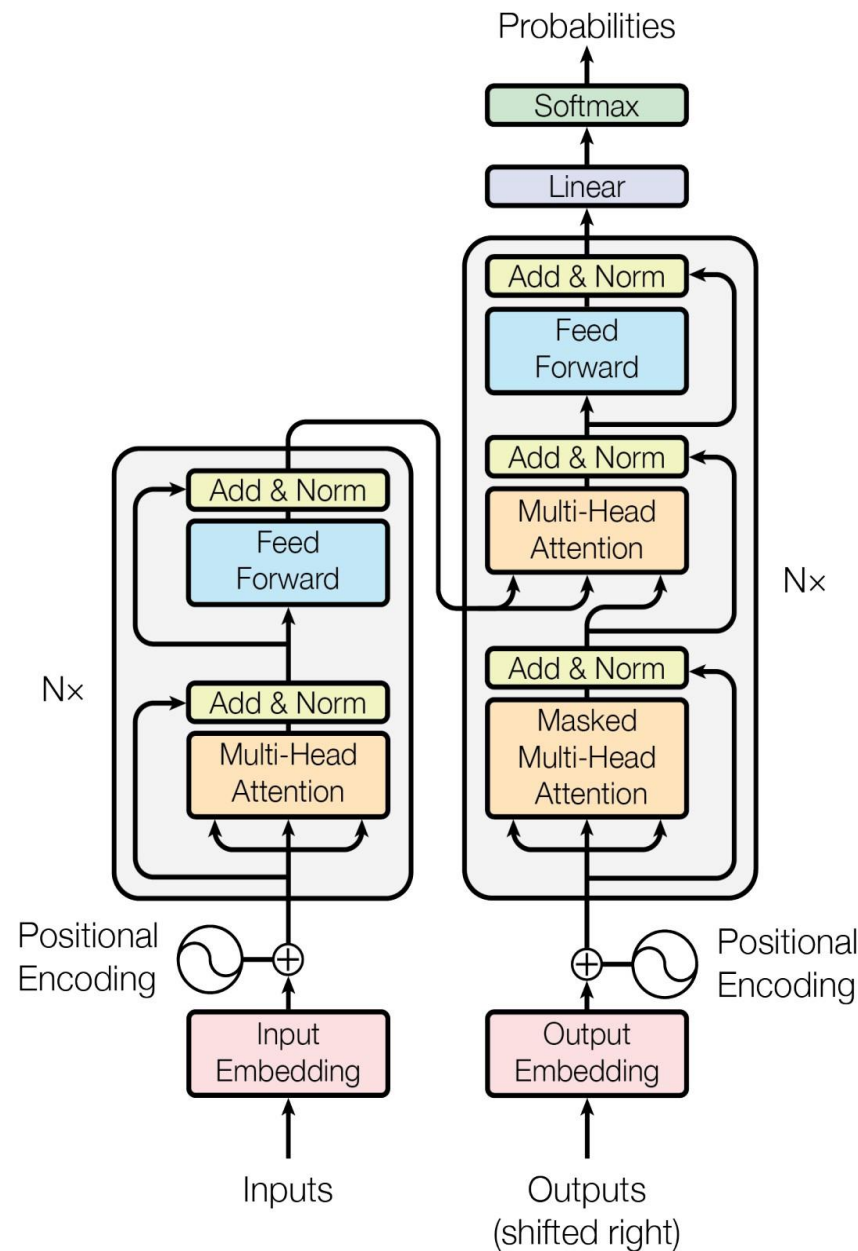
BART and T5

Pre-Training Encoder-Decoder Models

- Masked LMs: trained bidirectionally but with masking
- How can we pre-train a model for $P(\mathbf{y} \mid \mathbf{x})$?
- Why was BERT effective?
 - Predicting a mask requires some kind of text “understanding”.
- What would it take to do the same for sequence prediction?

Recall: Encoder-Decoder Architecture

- Standard Transformer Architecture
- Decoder attends back to the input. But the input doesn't change, so this just needs to be encoded once.



Pre-Training Encoder-Decoder Models

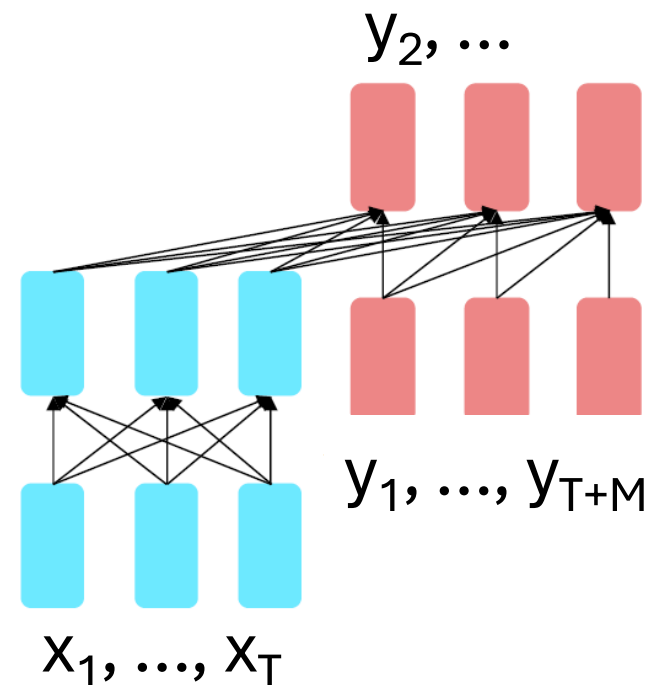
- For **encoder-decoders**, we could do something like **language modeling**, but where a prefix of every input is provided to the encoder and is not predicted.

$$h_1, \dots, h_T = \text{Encoder}(x_1, \dots, x_T)$$

$$h_{T+1}, \dots, h_{T+M} = \text{Decoder}(y_1, \dots, y_{i-1}, h_1, \dots, h_T)$$

$$P(y_i | y_{<i}, h_{1:T}) = \text{Softmax}(Wh_i + b)$$

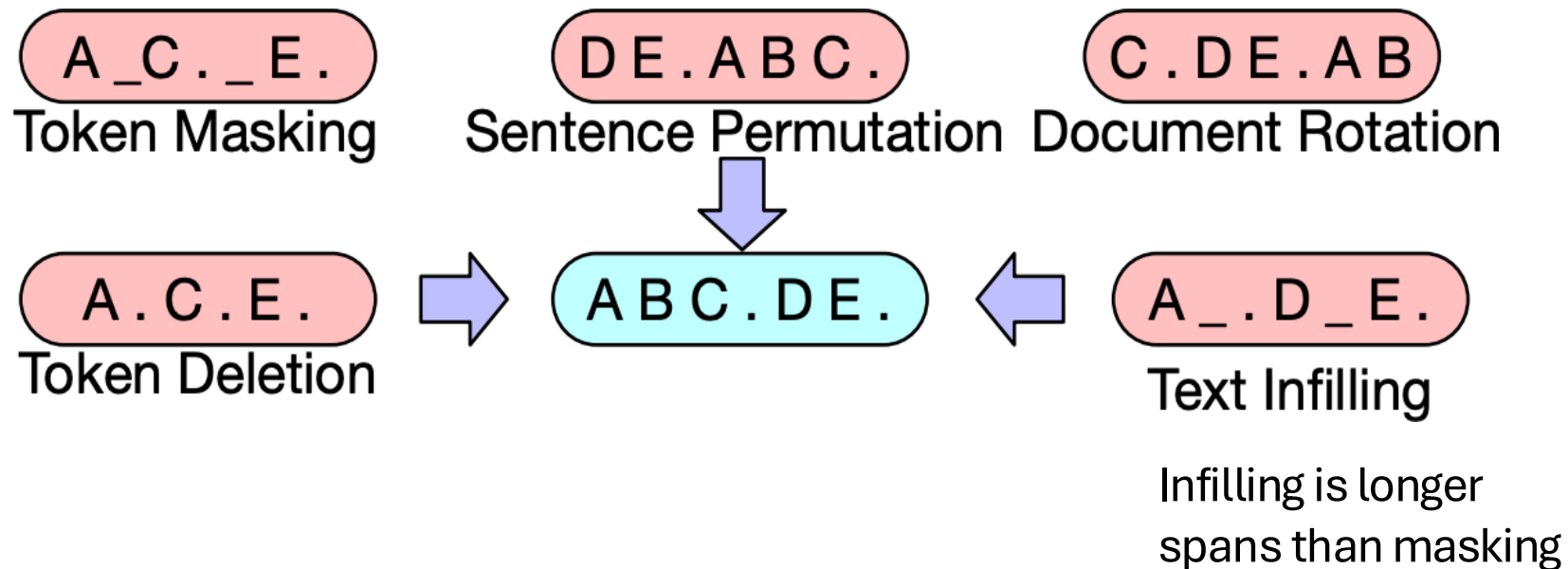
The **encoder** portion benefits from bidirectional context; the **decoder** portion is used to train the whole model through language modeling.



Pre-Training Encoder-Decoder Models

- How can we pre-train a model for $P(\mathbf{y} \mid \mathbf{x})$?
- **Requirements:**
 1. should use unlabeled data
 2. should force a model to attend from \mathbf{y} back to \mathbf{x}

Pre-Training BART (Bidirectional and Auto-Regressive Transformers)

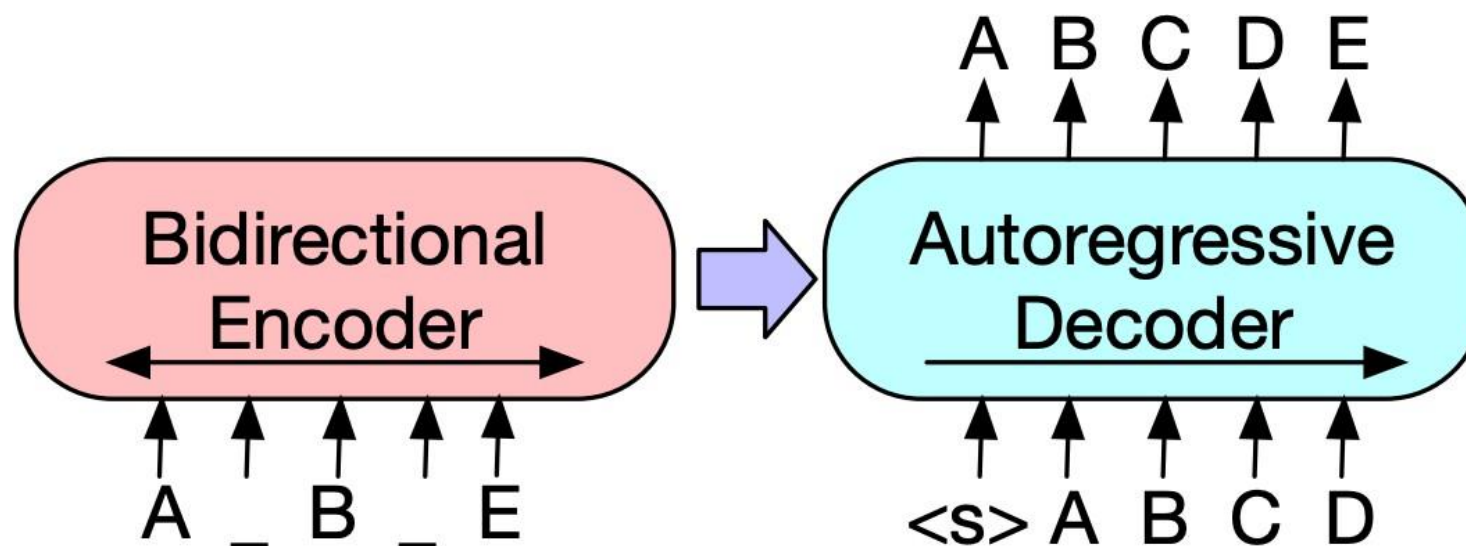


- Several possible strategies for corrupting a sequence are explored in the BART paper.

Lewis et al. (2019), "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension"

Pre-Training BART

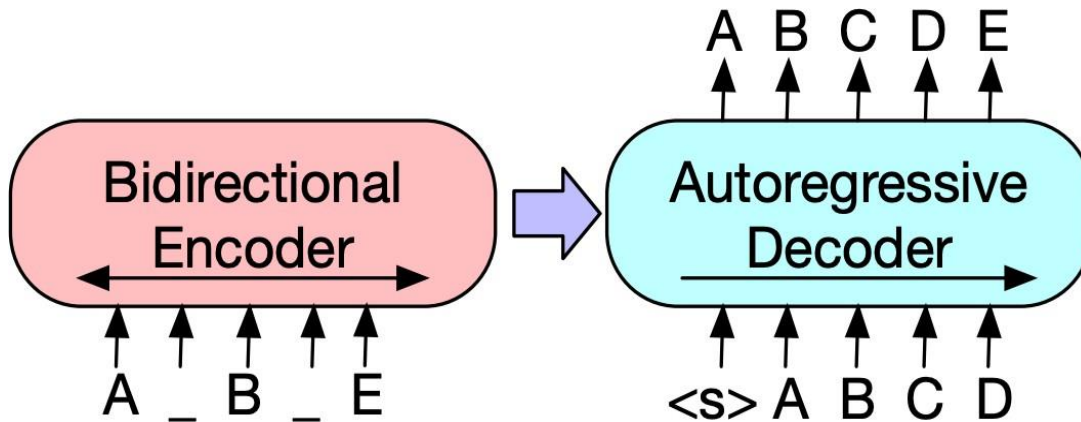
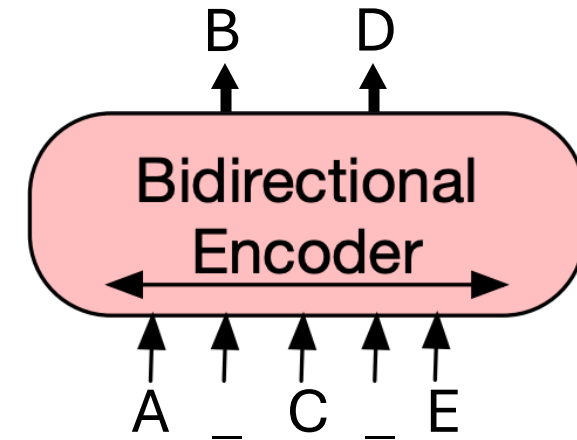
- Sequence-to-sequence Transformer trained on this data: permute/make/delete tokens, then predict full sequence autoregressively.



Lewis et al. (2019), "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension"

BERT vs. BART

- **BERT:** only an encoder, trained with masked language modeling objective. Cannot generate text or do Seq2Seq tasks (in standard form).



- **BART:** consists of both an encoder and a decoder. Can also use just the encoder wherever we would use BERT.

BART for Summarization

- **Pre-train** on the BART task: take random chunks of text, noise them according to the schemes described, and try to “decode” the clean text
- **Fine-tune** on a summarization dataset: a news article is the input and a summary of that article is the output (usually 1-3 sentences depending on the dataset)
- Can achieve good results even with **few summaries to fine-tune on**, compared to basic seq2seq models which require 100k+ examples to do well

BART for Summarization: Output Example

PG&E stated it scheduled the blackouts in response to forecasts for high winds amid dry conditions. The aim is to reduce the risk of wildfires. Nearly 800 thousand customers were scheduled to be affected by the shutoffs which were expected to last through at least midday tomorrow.



Power has been turned off to millions of customers in California as part of a power shutoff plan.



March 25, 2025

[Google OpenAI](#)

GEMINI 2.5

Gemini 2.5 models are **thinking models**, capable of reasoning through their thoughts before responding

Combines a significantly base model with **improved post-training**

Gemini 2.5 Pro **tops the LMArena leaderboard** by a significant margin.

Gemini 2.5 Pro scores **18.8% on Humanity's last exam** and leads in math and science benchmarks



OpenAI 4o Image Generation

Image generation that is not only beautiful, but **useful**.

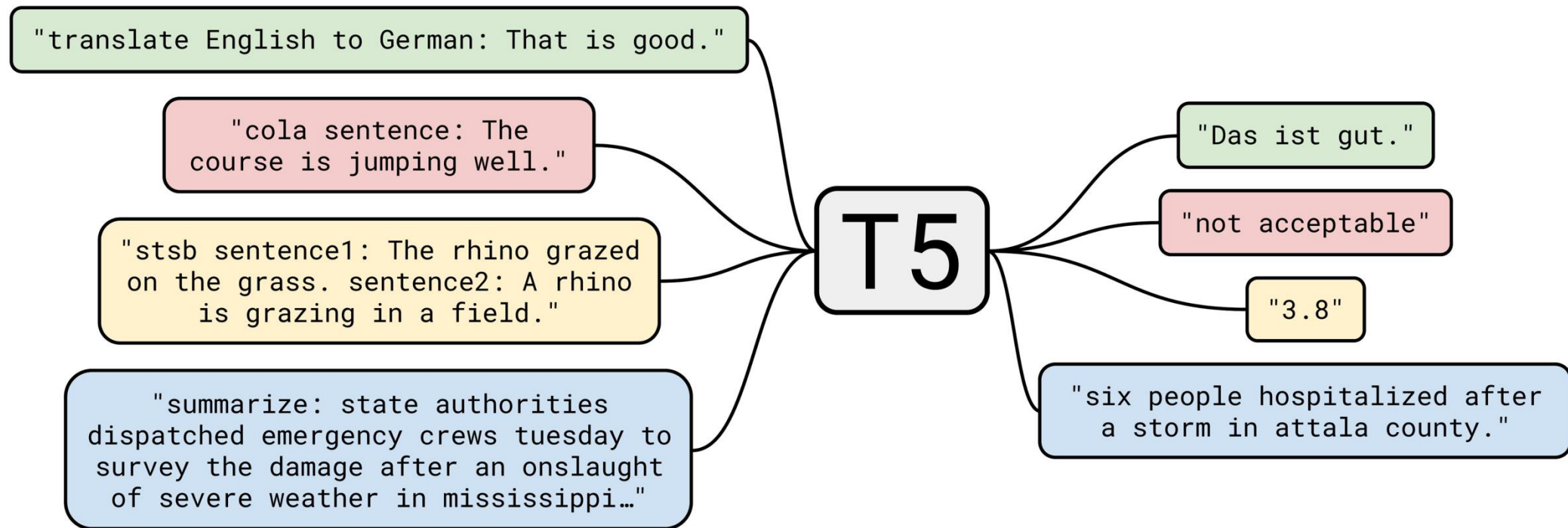
GPT-4o image generation excels at **accurately rendering text**, precisely **following prompts**, and **leveraging 4o's inherent knowledge base** and chat context

Trained models on the **joint distribution of online images and text**, and learnt how images relate to language and to each other

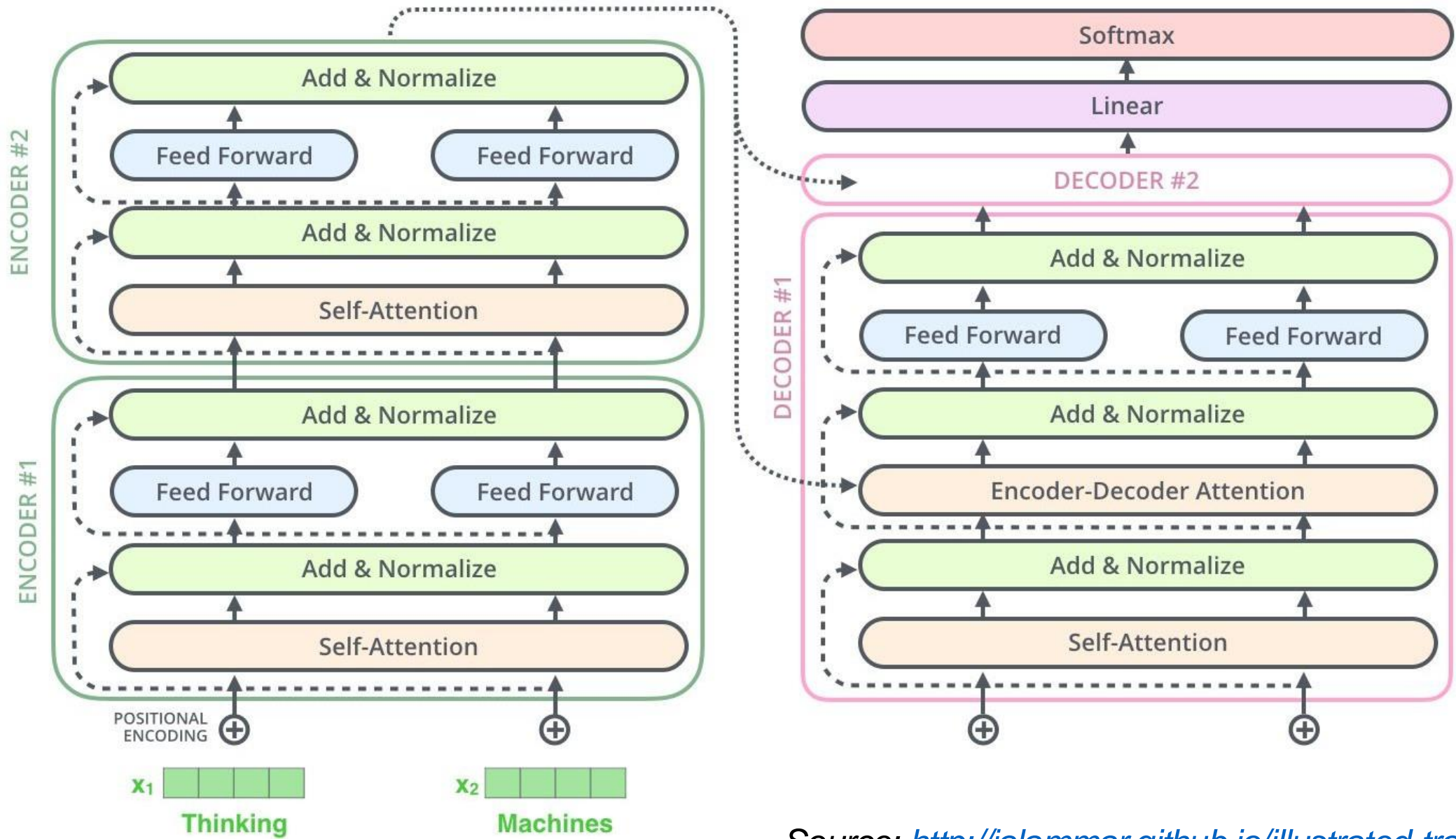
- Two major announcements:
- 1) Gemini 2.5 by Google - Advanced Reasoning
 - 2) OpenAI 4o Image Generation - Mind-blowing Image generation



T5: Text-to-Text Transfer Transformer



Raffel et al. (2019), "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer"



Source: <http://jalammar.github.io/illustrated-transformer/>

Pre-Training T5

- **Pre-training:** similar denoising scheme to BART (they were released within a week of each other in fall 2019)
- **Input:** text with gaps ; **Output:** a series of phrases to fill those gaps.

Original text

Thank you ~~for inviting~~ me to your party ~~last~~ week.

Inputs

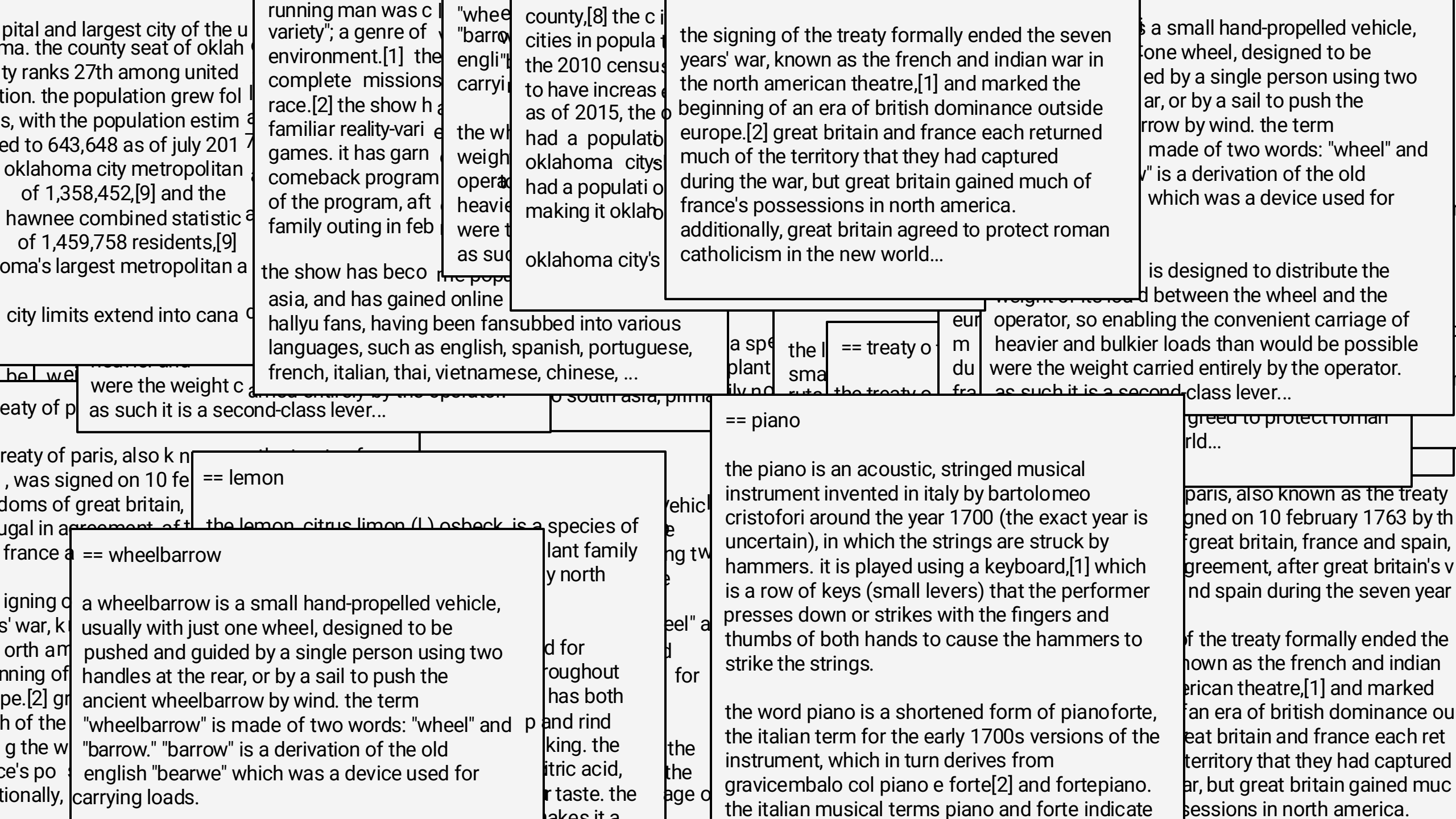
Thank you <X> me to your party <Y> week.

Targets

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

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

Raffel et al. (2019)



Common Crawl Web Extracted Text

Menu

Lemon

Introduction

The lemon, Citrus Limon (L.) Osbeck, is a species of small evergreen tree in the flowering plant family rutaceae. The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily as a flavoring which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a pH of around 2.2, giving it a sour taste.

Article

The origin of the lemon is unknown, though lemons are thought to have first grown in Assam (a region in northeast India), northern Burma or China. A genomic study of the lemon indicated it was a hybrid between bitter orange (sour orange) and citron.

Please enable JavaScript to use our site.

Home
Products
Shipping
Contact

Lorem ipsum dolor sit amet, consectetur adipiscing elit.
Curabitur in tempus quam. In mollis et ante at consectetur.
Aliquam erat volutpat.
Deneque et lacinia est.

- Removed lines that didn't end in a terminal punctuation mark.
- Language classifier to retain only English text
- Removed texts which look like placeholder texts
- Removed anything which look like code
- Removed duplicated texts

The tree's ellipsoidal yellow fruit is used for culinary and non-culinary purposes throughout the world, primarily for its juice, which has both culinary and cleaning uses. The juice of the lemon is about 5% to 6% citric acid, with a pH of around 2.2, giving it a sour taste.

```
function Ball(r) {  
  this.radius = r;  
  this.area = pi * r ** 2;  
  this.show = function(){  
    drawCircle(r);  
  }  
}
```

Datasets v1.3.2

[Overview](#)
[Catalog](#)
[Guide](#)
[API](#)
[Overview](#)

▶ [Audio](#)

▶ [Image](#)

▶ [Object_detection](#)

▶ [Structured](#)

▶ [Summarization](#)

▼ [Text](#)
[c4 \(manual\)](#)
[civil_comments](#)
[definite_pronoun_resolution](#)
[esnli](#)
[gap](#)
[glue](#)
[imdb_reviews](#)
[TensorFlow](#) > [Resources](#) > [Datasets v1.3.2](#) > [Catalog](#)


c4 (Manual download)

Contents ▼

[c4/en](#)
[Statistics](#)
[Features](#)
[Homepage](#)

...

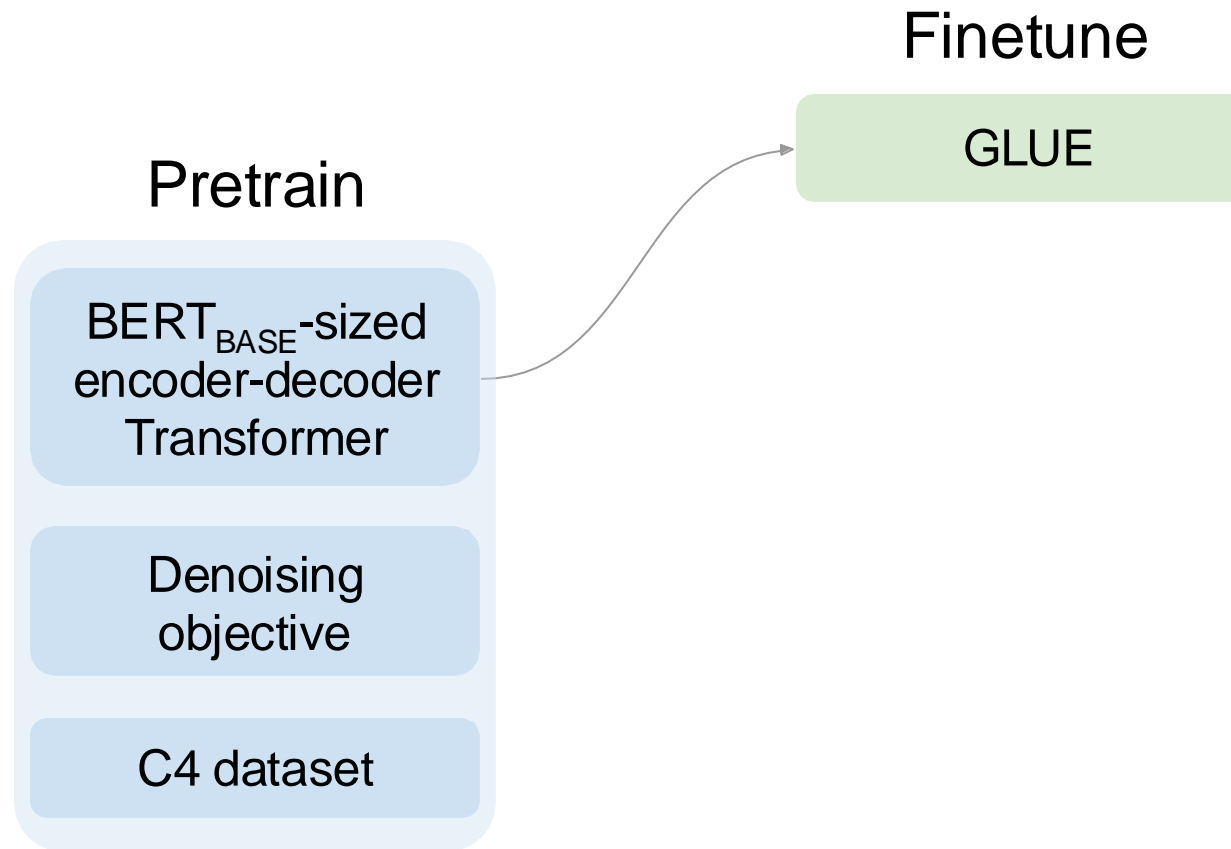
A colossal, cleaned version of Common Crawl's web crawl corpus.

Pretrain

BERT_{BASE}-sized
encoder-decoder
Transformer

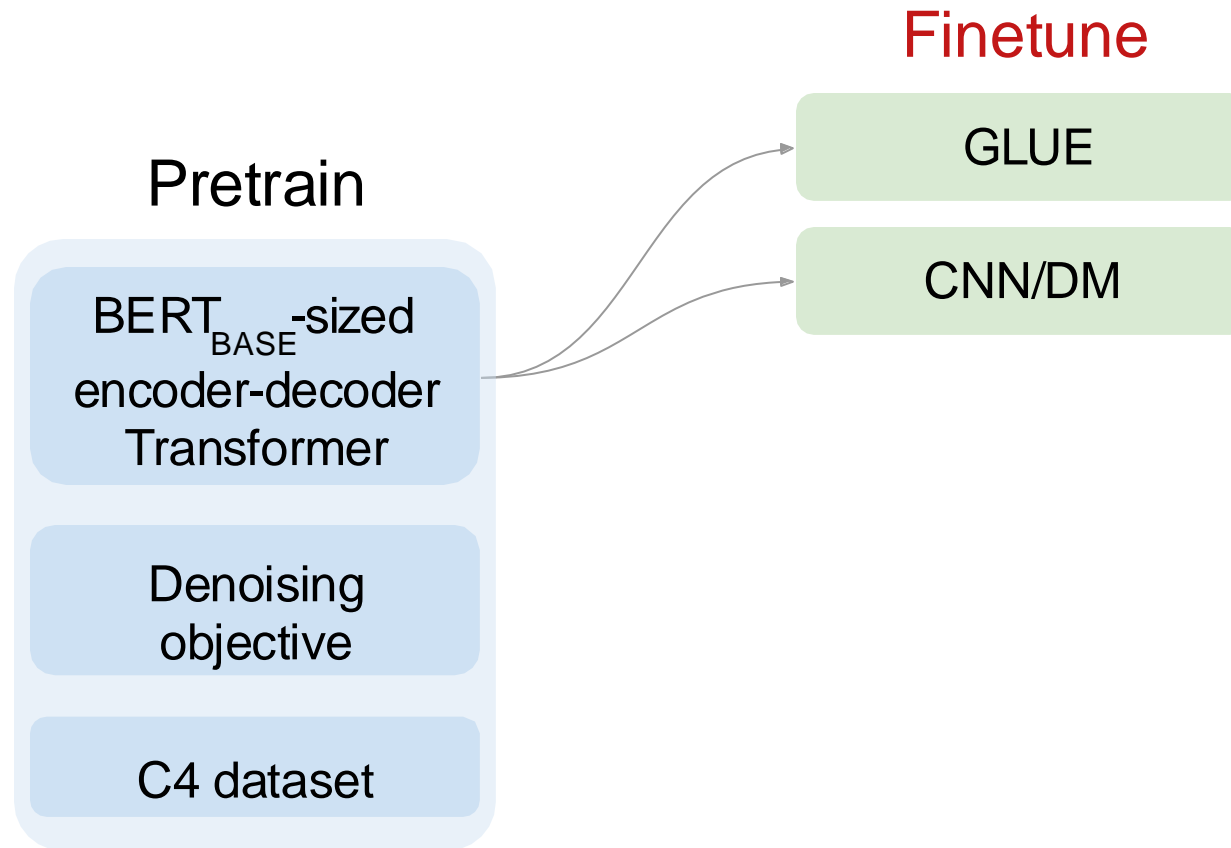
Denoising
objective

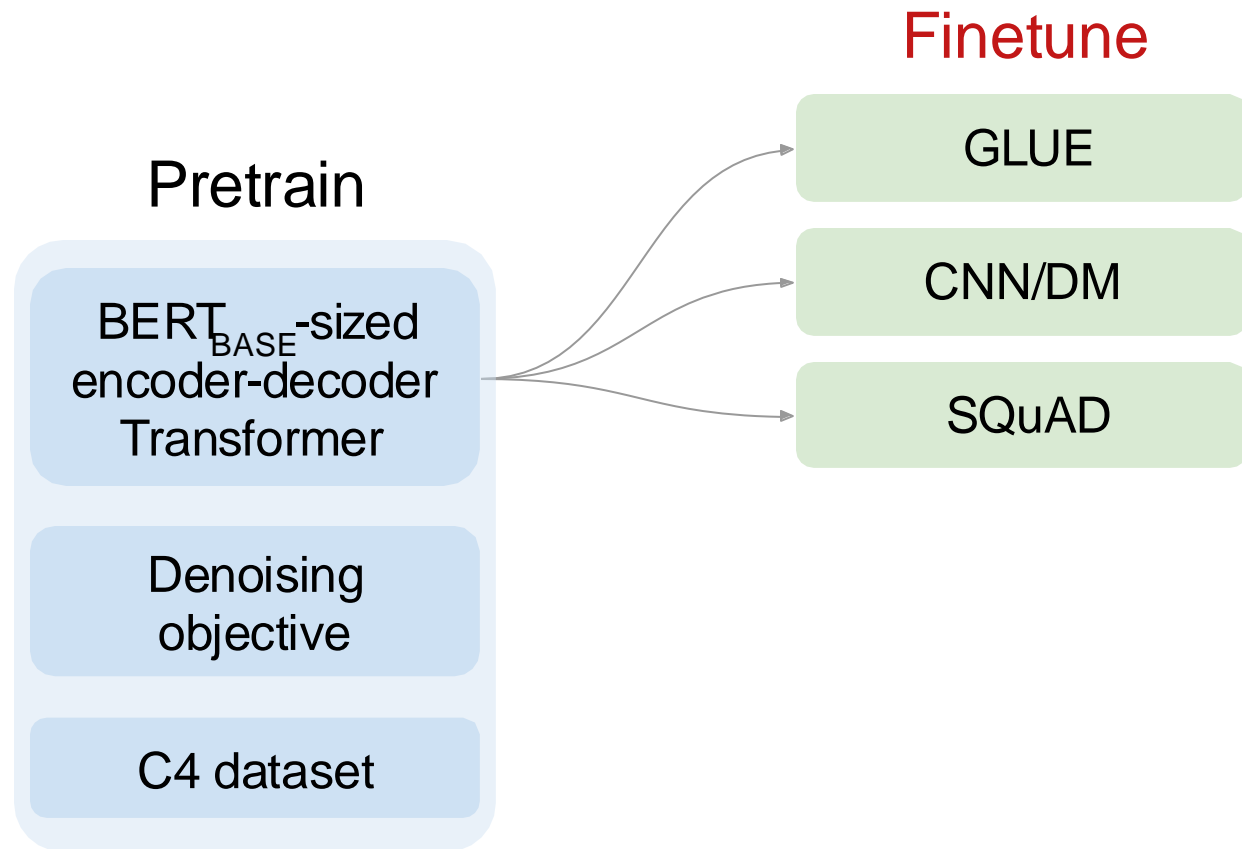
C4 dataset

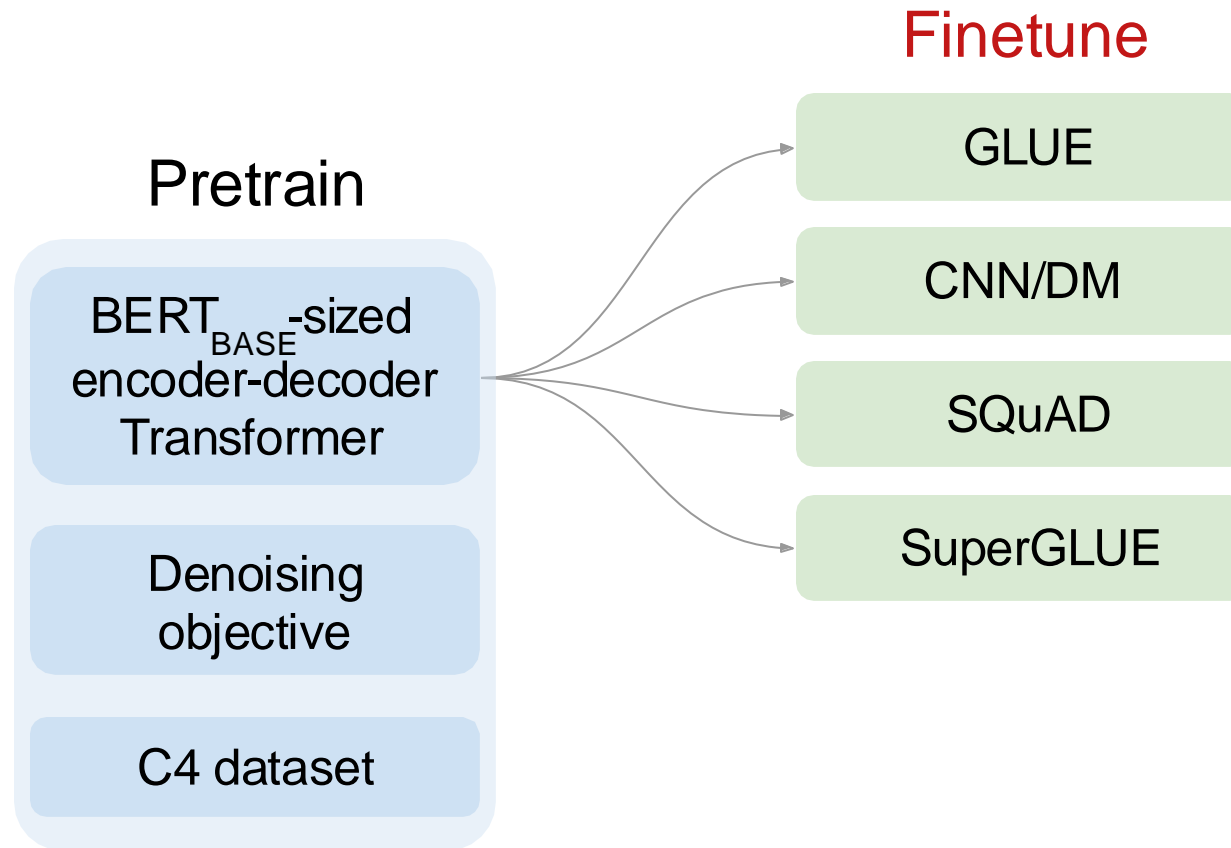


GLUE Benchmark





















Dataset	Description	Data example	Metric
CoLA	Is the sentence grammatical or ungrammatical?	"This building is than that one." = Ungrammatical	Matthews
SST-2	Is the movie review positive, negative, or neutral?	"The movie is funny , smart , visually inventive , and most of all , alive ." = .93056 (Very Positive)	Accuracy
MRPC	Is the sentence B a paraphrase of sentence A?	A) "Yesterday , Taiwan reported 35 new infections , bringing the total number of cases to 418 ." B) "The island reported another 35 probable cases yesterday , taking its total to 418 ." = A Paraphrase	Accuracy / F1
STS-B	How similar are sentences A and B?	A) "Elephants are walking down a trail." B) "A herd of elephants are walking along a trail." = 4.6 (Very Similar)	Pearson / Spearman
QQP	Are the two questions similar?	A) "How can I increase the speed of my internet connection while using a VPN?" B) "How can Internet speed be increased by hacking through DNS?" = Not Similar	Accuracy / F1
MNLI-mm	Does sentence A entail or contradict sentence B?	A) "Tourist Information offices can be very helpful." B) "Tourist Information offices are never of any help." = Contradiction	Accuracy
QNLI	Does sentence B contain the answer to the question in sentence A?	A) "What is essential for the mating of the elements that create radio waves?" B) "Antennas are required by any radio receiver or transmitter to couple its electrical connection to the electromagnetic field." = Answerable	Accuracy
RTE	Does sentence A entail sentence B?	A) "In 2003, Yunus brought the microcredit revolution to the streets of Bangladesh to support more than 50,000 beggars, whom the Grameen Bank respectfully calls Struggling Members." B) "Yunus supported more than 50,000 Struggling Members." = Entailed	Accuracy
WNLI	Sentence B replaces sentence A's ambiguous pronoun with one of the nouns - is this the correct noun?	A) "Lily spoke to Donna, breaking her concentration." B) "Lily spoke to Donna, breaking Lily's concentration." = Incorrect Referent	Accuracy

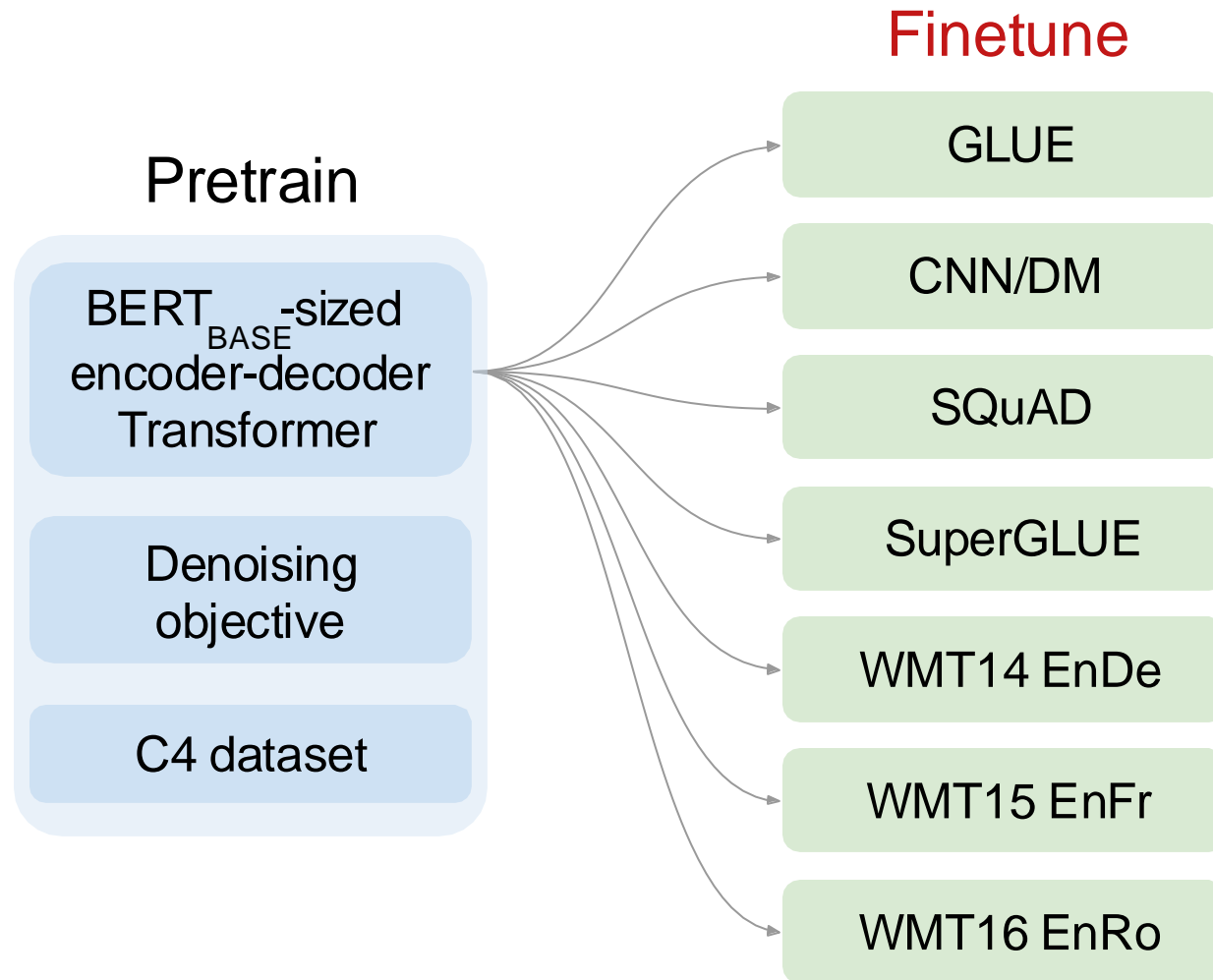


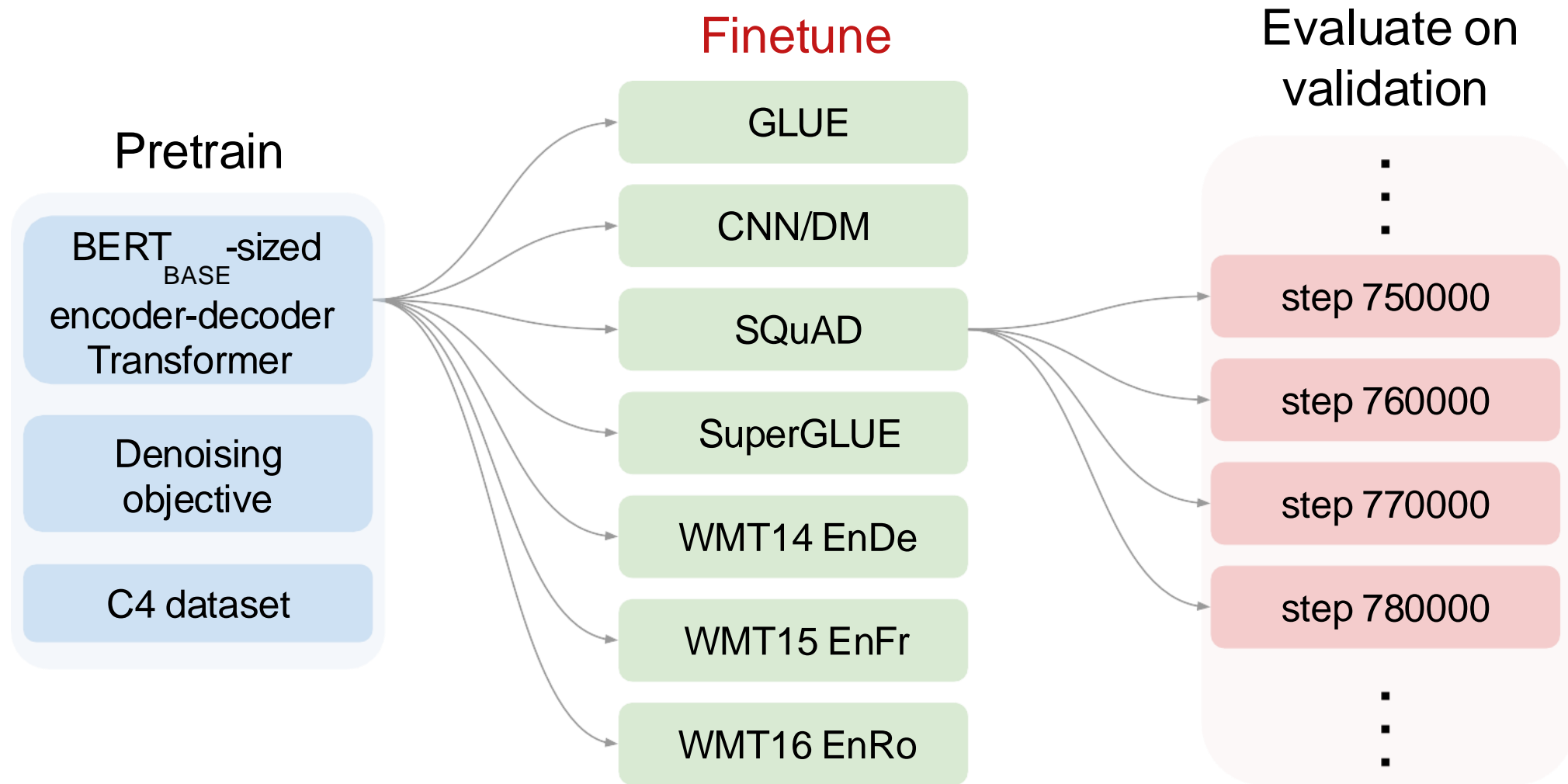




SuperGLUE Tasks

Name	Identifier	Download	More Info	Metric
Broadcoverage Diagnostics	AX-b			Matthew's Corr
CommitmentBank	CB			Avg. F1 / Accuracy
Choice of Plausible Alternatives	COPA			Accuracy
Multi-Sentence Reading Comprehension	MultiRC			F1a / EM
Recognizing Textual Entailment	RTE			Accuracy
Words in Context	WiC			Accuracy
The Winograd Schema Challenge	WSC			Accuracy
BoolQ	BoolQ			Accuracy
Reading Comprehension with Commonsense Reasoning	ReCoRD			F1 / Accuracy
Winogender Schema Diagnostics	AX-g			Gender Parity / Accuracy





	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline average	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Baseline standard deviation	0.235	0.065	0.343	0.416	0.112	0.090	0.108
No pre-training	66.22	17.60	50.31	53.04	25.86	39.77	24.04

Star denotes baseline

Comparable to BERT

Bold = 1 std. dev. of max

	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline average	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Baseline standard deviation	0.235	0.065	0.343	0.416	0.112	0.090	0.108
No pre-training	66.22	17.60	50.31	53.04	25.86	39.77	24.04

Big training set

No pre-training is dramatically worse, except EnFr!

C4: The Data

- C4: Colossal Clean Crawled Corpus
 - Web-extracted text
 - English language only
 - 750GB

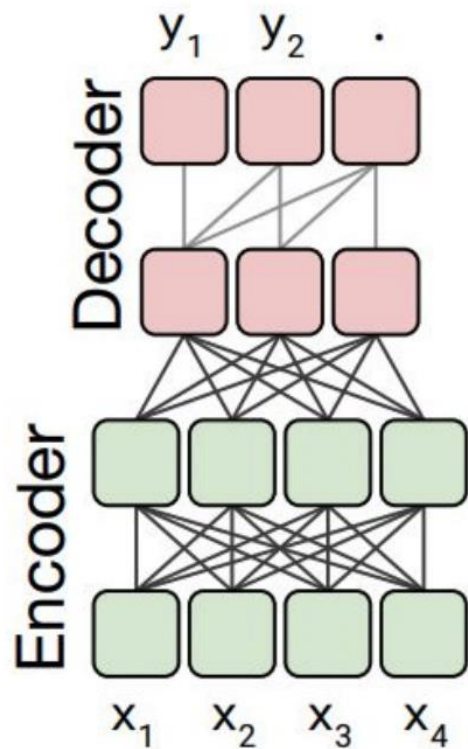
Data set	Size
★ C4	745GB
C4, unfiltered	6.1TB

Pre-Training Data: Experiment

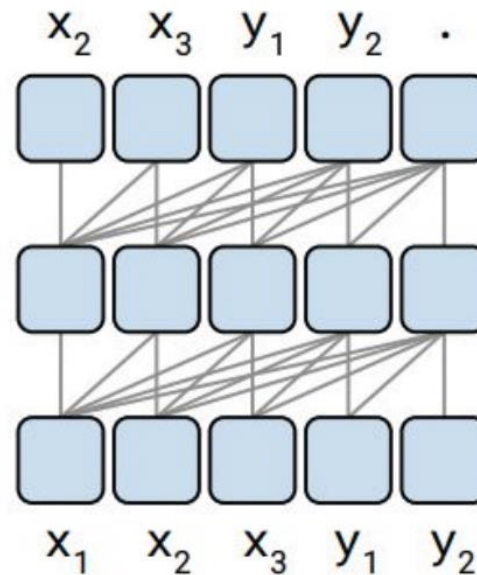
- Takeaway:
 - Clean and compact data is better than large, but noisy data.
 - Pre-training on in-domain data helps.

Data set	Size	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ C4	745GB	83.28	19.24	80.88	71.36	26.98	39.82	27.65
C4, unfiltered	6.1TB	81.46	19.14	78.78	68.04	26.55	39.34	27.21

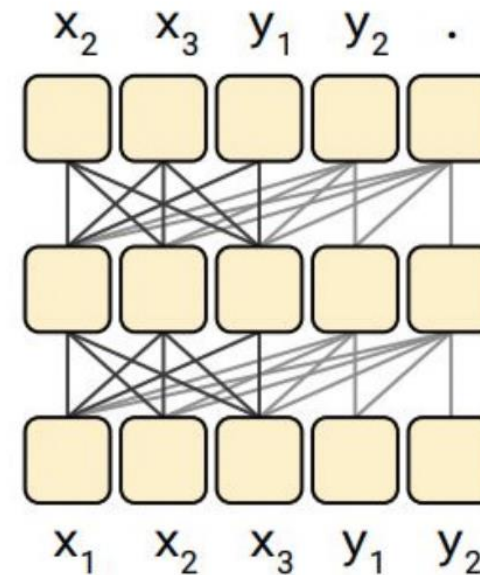
Architectures: Different Choices



Language model

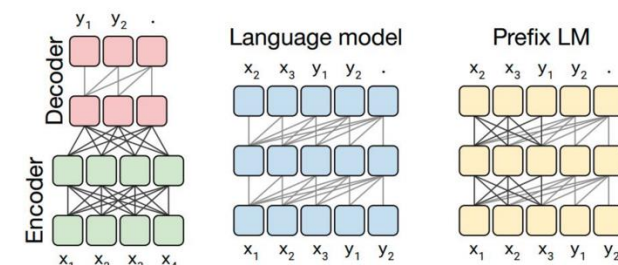
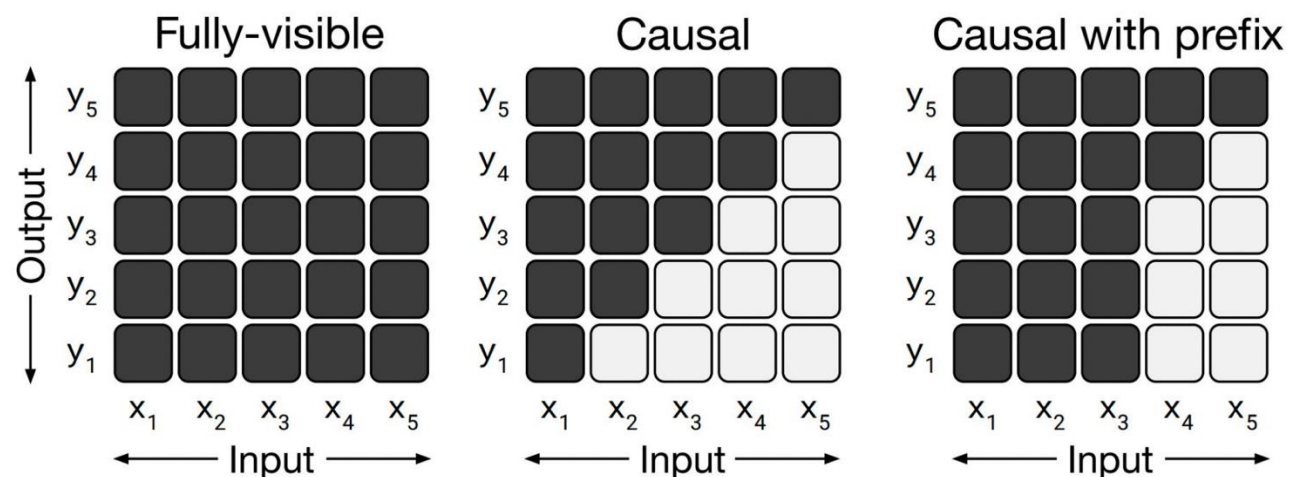


Prefix LM



Architectures: Different Attention Masks

- **Fully visible mask** allows the self attention mechanism to attend to the full input.
- A **causal mask** doesn't allow output elements to look into the future.
- **Causal mask with prefix** allows to fully-visible masking on a portion of input.



Architectural Variants: Experiments

Evaluated for classification tasks.

Architecture	Objective	Params	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
Encoder-decoder	Denoising	$2P$	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Enc-dec, shared	Denoising	P	82.81	18.78	80.63	70.73	26.72	39.03	27.46
Enc-dec, 6 layers	Denoising	P	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	P	74.70	17.93	61.14	55.02	25.09	35.28	25.86
Prefix LM	Denoising	P	81.82	18.61	78.94	68.11	26.43	37.98	27.39

Takeaways:

1. Halving the number of layers in encoder and decoder hurts the performance.
2. Performance of Encoder-Decoder with shared params is almost on-par with prefix LM.

T5: Pre-Training Objectives

- **Prefix language modeling**
 - **Input:** Thank you for inviting
 - **Output:** me to your party last week
- **BERT-style denoising**
 - **Input:** Thank you <M> <M> me to your party
apple week
 - **Output:** Thank you for inviting me to your party
last week
- **De-shuffling**
 - **Input:** party me for your to. last fun you inviting
week Thanks.
 - **Output:** Thank you for inviting me to your party
last week
- **Replace spans**
 - **Input:** Thank you <X> me to your party <X>
week
 - **Output:** <X> for inviting <Y> last <Z>
- **Drop tokens**
 - **Input:** Thank you me to your party week .
 - **Output:** for inviting last

Pre-Training Objectives: Experiments

- All the variants perform similarly
- “Replace corrupted spans” and “Drop corrupted tokens” are more appealing because **target sequences are shorter, speeding up training.**

Assuming Enc-Dec architecture.
Evaluated for classification tasks.

Objective	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
Prefix language modeling	80.69	18.94	77.99	65.27	26.86	39.73	27.49
Deshuffling	73.17	18.59	67.61	58.47	26.11	39.30	25.62
BERT-style (Devlin et al., 2018)	82.96	19.17	80.65	69.85	26.78	40.03	27.41
★ Replace corrupted spans	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Drop corrupted tokens	84.44	19.31	80.52	68.67	27.07	39.76	27.82

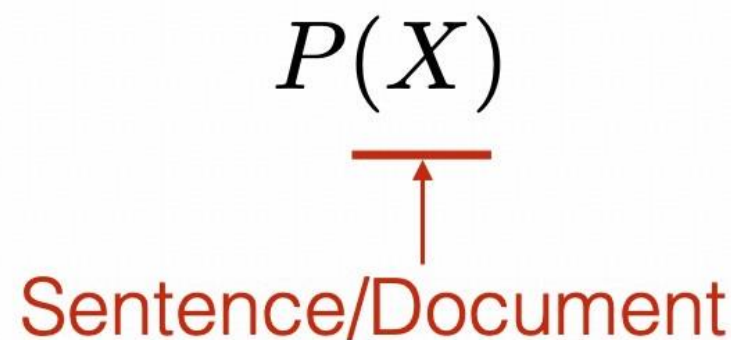
Pre-Training Decoder-only Models

GPT and Llama

Recall: Probabilistic Language Models


$$P(\underline{X})$$

Sentence/Document



A generative model that calculates the probability of language

Auto-regressive Language Models

$$P(X) = \prod_{i=1}^I P(x_i \mid x_1, \dots, x_{i-1})$$


The diagram illustrates the components of the probability formula. A red horizontal line is positioned under the term x_i in the numerator, with a red arrow pointing from the text "Next Token" below it to this line. A blue horizontal line is positioned under the denominator x_1, \dots, x_{i-1} , with a blue arrow pointing from the text "Context" below it to this line.

Next Token Prediction

- This is essentially **classification**!
 - We can think of neural language models as neural classifiers. They classify prefix of a text into $|V|$ classes, where the classes are vocabulary tokens.

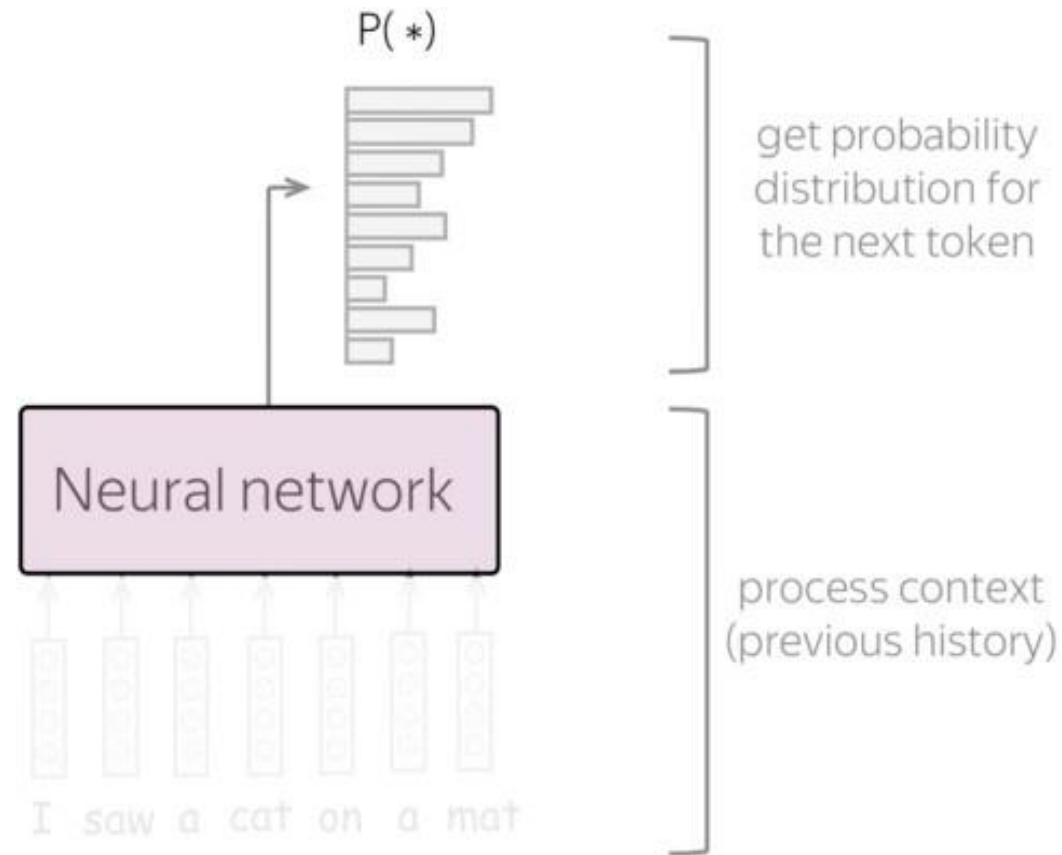
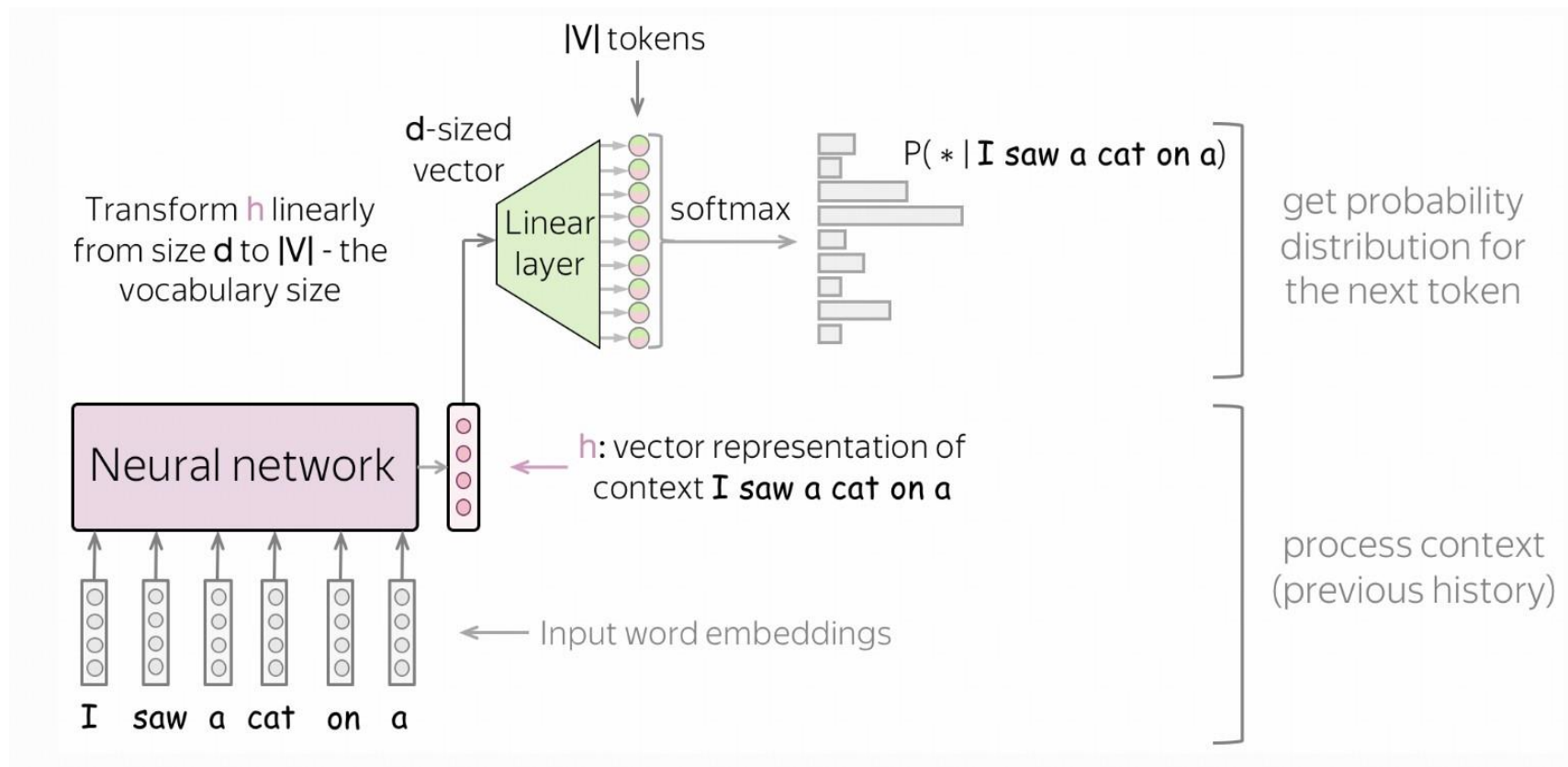


Image Credit: https://lena-voita.github.io/nlp_course/language_modeling.html

Next Token Prediction



- Feed word embedding for previous (context) words into a network.
- Get vector representation of context from the network.
- From this vector representation, predict a probability distribution for the next token.

Image Credit: https://lena-voita.github.io/nlp_course/language_modeling.html

Encoders vs. Decoders

- BERT is a Transformer **Encoder**: **bidirectional attention**, trained with masked language modelling.

$$P(x_i \mid x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

- GPT and many other Transformer language models (e.g., LLaMA) are **Decoders**: **unidirectional attention**, trained to predict the next token.

$$P(x_i \mid x_1, \dots, x_{i-1})$$

Generative Pre-trained Transformer (GPT)

- 2018's GPT was a big success in pretraining a decoder!
- **Transformer decoder with 12 layers, 117M parameters.**
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers.
- Byte-pair encoding with 40,000 merges
 - Trained on BooksCorpus: over 7000 unique books.
- Contains long spans of contiguous text, for learning long-distance dependencies.

Radford et al. (2018), "Improving Language Understanding by Generative Pre-Training"

GPT-2

GPT-2 is identical to GPT-1, but:

- Has Layer normalization in between each sub-block
- Vocab extended to 50,257 tokens and context size increased from 512 to 1024
- Data: 8 million docs from the web (Common Crawl), minus Wikipedia

Language Models are Unsupervised Multitask Learners

Alec Radford ^{* 1} Jeffrey Wu ^{* 1} Rewon Child ¹ David Luan ¹ Dario Amodei ^{** 1} Ilya Sutskever ^{** 1}

Increasingly Convincing Generations by GPT-2

- We discussed how we can sample sentences from auto-regressive LMs for text generation.
 - This is how pre-trained decoders are used **in their capacities as language models**.
- **GPT-2**, a larger version (1.5B) of GPT trained on more data, was shown to produce relatively convincing samples of natural language.

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

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

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

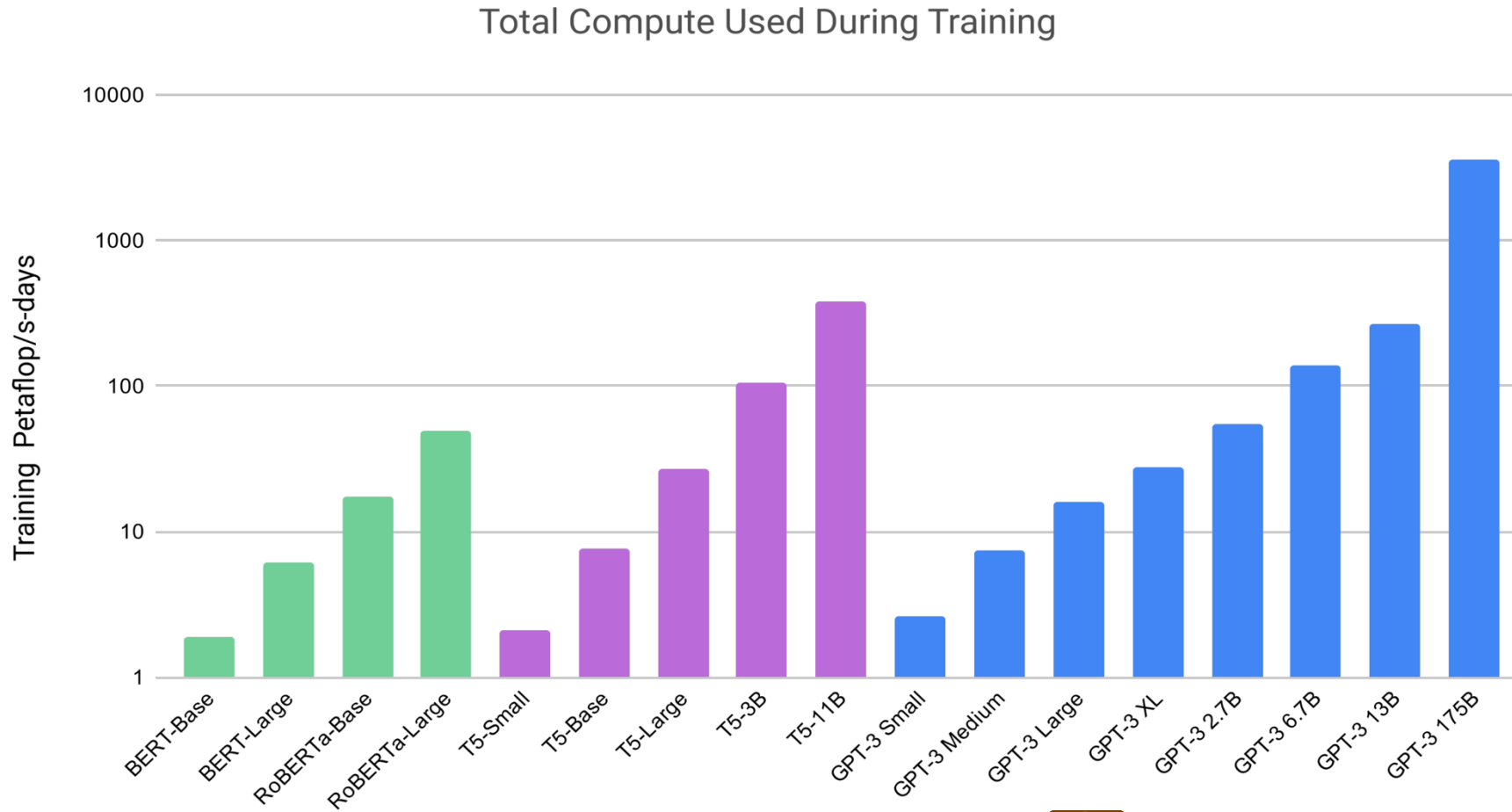
Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock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

Pre-Training Cost (with GCP/AWS)

- **BERT:** Base \$500, Large \$7000
- **GPT-2** (as reported in other work): \$25,000
- This is for a single pre-training run...developing new pre-training techniques may require many runs.
- Fine-tuning these models can typically be done with a single GPU (but may take 1-3 days for medium-sized datasets).

<https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/>

GPT-3



- **175B parameter model**
 - 96 layers, 96 heads, 12k-dim vectors
- Trained on Microsoft Azure, estimated to cost roughly **\$10M**

Comparison: GPT-1, 2, 3

Model	Parameters	Layers	Training Data	Key Advancement
GPT-1	117M	12	BooksCorpus	First large-scale Transformer for NLP
GPT-2	1.5B	48	WebText	Zero-shot learning, larger training data
GPT-3	175B	96	Common Crawl + others	In-context learning, emergent behaviors

GPT-4

- Transformer-based

- The rest is mystery!
- If we're going based on costs, GPT-4 is ~15-30 times costlier than GPT3. That should give you an idea how its likely size!

- Note, these language models involve more than just pre-training.

- Pre-training provides the foundation based on which we build the model.
- We will discuss the later stages next week.

Model	Usage
davinci-002	\$0.0020 / 1K tokens

Model	Input	Output
gpt-4	\$0.03 / 1K tokens	\$0.06 / 1K tokens

Llama: A Family of Open-Source LLMs from Meta AI

- **Llama-1 + Llama-2**

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.**

- Models have mostly gotten smaller since GPT-3, but haven't changed much.
- **Tokenizer: Byte-Pair Encoding (BPE)** [Recall: we have already discussed this algorithm in lecture on 'Tokenization Strategies']
- **Rotary positional encodings**, a few other small architecture changes
- **Optimized mix of pre-training data:** Common Crawl, GitHub, Wikipedia, Books, etc.

Next Week: How to Make Pre-Trained LMs Work?

- Instruction Tuning
 - Finetune the pre-trained model to follow instructions
- Prompting and In-context Learning
 - Give few examples of the task that you want the model to solve
- Reinforcement Learning from Human Feedback (RLHF)
 - Train the LMs to align their outputs with human preferences
 - Also called 'Preference Optimization'

