

# **LLMs and the future of NLP**

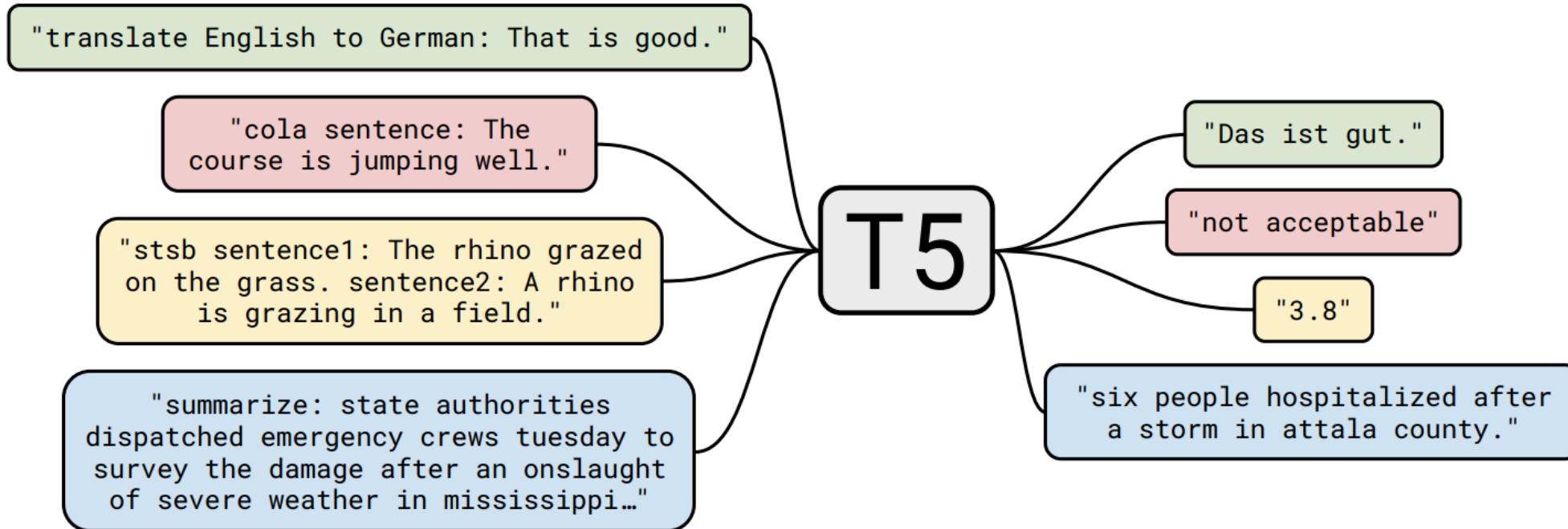
## **NLP Week 11**

**Thanks to Dan Jurafsky for most of the slides this week!**

# Plan for today

1. T5: everything is text-to-text
2. Large language models are few-shot learners
3. Sampling from language models
4. Prompting and in-context learning
5. Scaling up
6. Harms of large language models
7. *Group exercises*

# T5: seq-to-seq everywhere



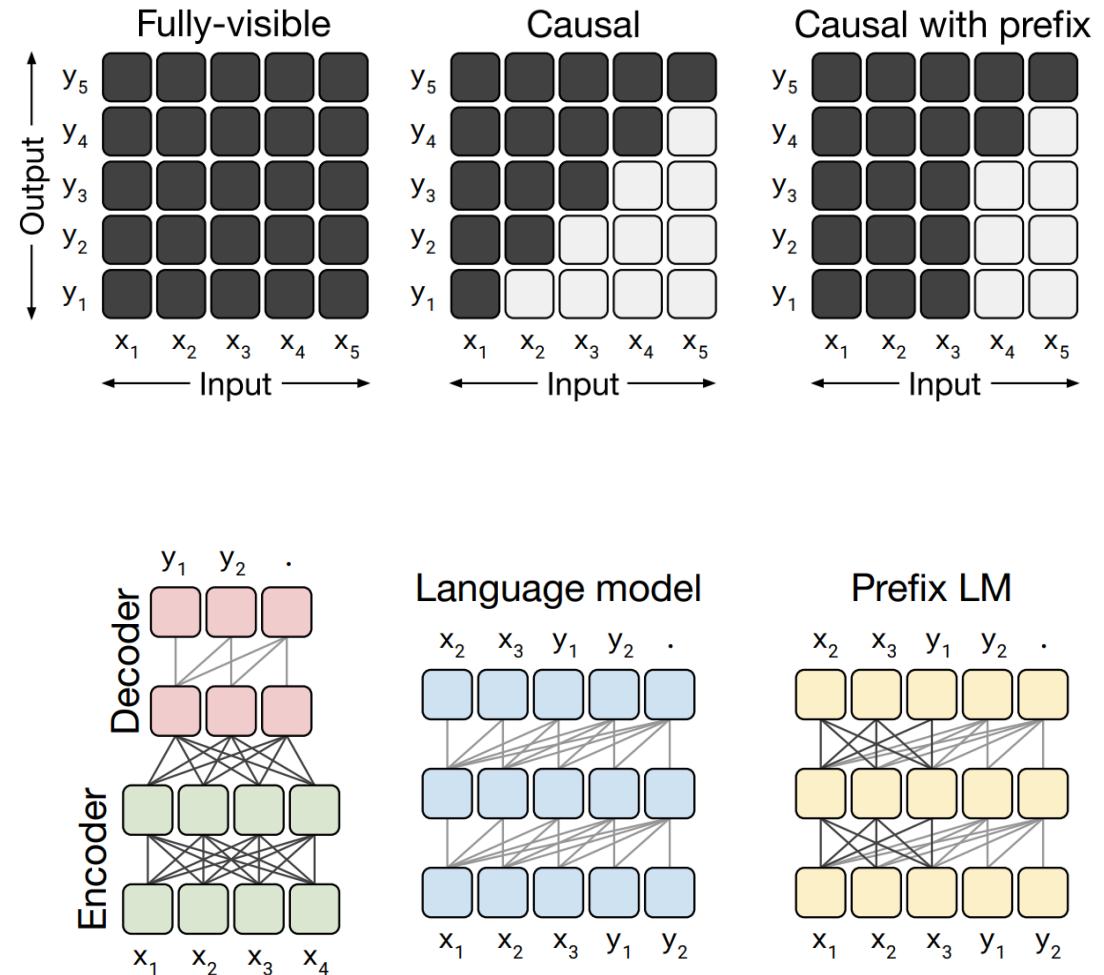
There are a few famous figures in our field, this is one of them :)

# T5: seq-to-seq everywhere

- General **text-to-text** format
- Machine translation
  - Translate this text to French: Montréal is great. Montréal est génial
- Summarization
  - Summarize this news article: [Some article] [Summary of the article]
- What other tasks can you think of?
- Can you think of a NLP task that **can't** be formatted as text-to-text?

# What's the best possible model for text-to-text leaning?

- **Model structure**
  - Which attention should we use?
    - Bi-directional
    - Causal
    - A mix of the two
- **Model architecture**
  - Encoder-decoder
  - Decoder-only
  - A mix of the two



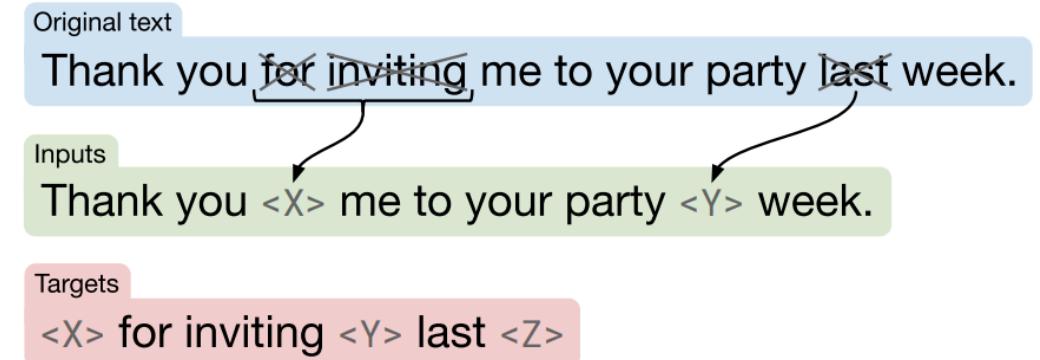
# What's the best possible model for text-to-text leaning?

Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	$M$	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
Enc-dec, shared	Denoising	$P$	$M$	82.81	18.78	<b>80.63</b>	<b>70.73</b>	26.72	39.03	<b>27.46</b>
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	$P$	$M$	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
Encoder-decoder	LM	$2P$	$M$	79.56	18.59	76.02	64.29	26.27	39.17	26.86
Enc-dec, shared	LM	$P$	$M$	79.60	18.13	76.35	63.50	26.62	39.17	27.05
Enc-dec, 6 layers	LM	$P$	$M/2$	78.67	18.26	75.32	64.06	26.13	38.42	26.89
Language model	LM	$P$	$M$	73.78	17.54	53.81	56.51	25.23	34.31	25.38
Prefix LM	LM	$P$	$M$	79.68	17.84	76.87	64.86	26.28	37.51	26.76

- Some combinations are better than others
- Denoising seems to be almost always better than LMing

# What's the best possible model for text-to-text learning?

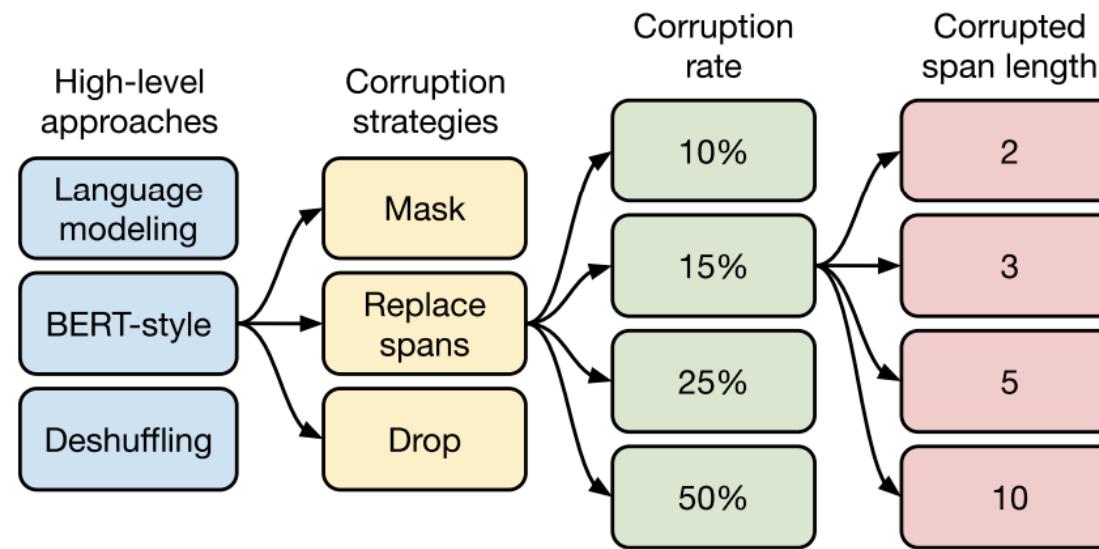
- Training objective
- What's the best way to encode useful information about and decode from a sequence?



Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style Devlin et al. (2018)	Thank you <M> <M> me to your party apple week .	(original text)
Deshuffling	party me for your to . last fun you inviting week Thank	(original text)
MASS-style Song et al. (2019)	Thank you <M> <M> me to your party <M> week .	(original text)
I.i.d. noise, replace spans	Thank you <X> me to your party <Y> week .	<X> for inviting <Y> last <Z>
I.i.d. noise, drop tokens	Thank you me to your party week .	for inviting last
Random spans	Thank you <X> to <Y> week .	<X> for inviting me <Y> your party last <Z>

# What's the best possible model for text-to-text leaning?

- Simply try all possible combinations



# What's the best possible model for text-to-text learning?

Objective	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
BERT-style (Devlin et al., 2018)	82.96	19.17	<b>80.65</b>	69.85	26.78	<b>40.03</b>	27.41
MASS-style (Song et al., 2019)	82.32	19.16	80.10	69.28	26.79	<b>39.89</b>	27.55
★ Replace corrupted spans	83.28	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	39.82	<b>27.65</b>
Drop corrupted tokens	<b>84.44</b>	<b>19.31</b>	<b>80.52</b>	68.67	<b>27.07</b>	39.76	<b>27.82</b>
Corruption rate	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
10%	<b>82.82</b>	19.00	<b>80.38</b>	69.55	<b>26.87</b>	39.28	<b>27.44</b>
★ 15%	<b>83.28</b>	19.24	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
25%	<b>83.00</b>	<b>19.54</b>	<b>80.96</b>	70.48	<b>27.04</b>	<b>39.83</b>	<b>27.47</b>
50%	81.27	19.32	79.80	70.33	<b>27.01</b>	<b>39.90</b>	<b>27.49</b>
Span length	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Baseline (i.i.d.)	<b>83.28</b>	19.24	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
2	<b>83.54</b>	19.39	<b>82.09</b>	<b>72.20</b>	<b>26.76</b>	<b>39.99</b>	<b>27.63</b>
3	<b>83.49</b>	<b>19.62</b>	<b>81.84</b>	<b>72.53</b>	<b>26.86</b>	39.65	<b>27.62</b>
5	<b>83.40</b>	19.24	<b>82.05</b>	<b>72.23</b>	<b>26.88</b>	39.40	<b>27.53</b>
10	82.85	19.33	<b>81.84</b>	70.44	<b>26.79</b>	39.49	<b>27.69</b>

- Each of these matters to some extent

# **What's the best possible model for text-to-text leaning?**

- One important ingredient is missing.
- Can you guess which one?

# What's the best possible model for text-to-text leaning?

- One important ingredient is missing
- Can you guess which one?

Data



# What's the best possible model for text-to-text learning?

- Data
  - C4: Colossal Clean Crawled Corpus

Data set	Size	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ C4	745GB	83.28	<b>19.24</b>	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
C4, unfiltered	6.1TB	81.46	19.14	78.78	68.04	26.55	39.34	27.21
RealNews-like	35GB	<b>83.83</b>	<b>19.23</b>	80.39	72.38	<b>26.75</b>	<b>39.90</b>	<b>27.48</b>
WebText-like	17GB	<b>84.03</b>	<b>19.31</b>	<b>81.42</b>	71.40	<b>26.80</b>	<b>39.74</b>	<b>27.59</b>
Wikipedia	16GB	81.85	<b>19.31</b>	81.29	68.01	<b>26.94</b>	39.69	<b>27.67</b>
Wikipedia + TBC	20GB	83.65	<b>19.28</b>	<b>82.08</b>	<b>73.24</b>	<b>26.77</b>	39.63	<b>27.57</b>

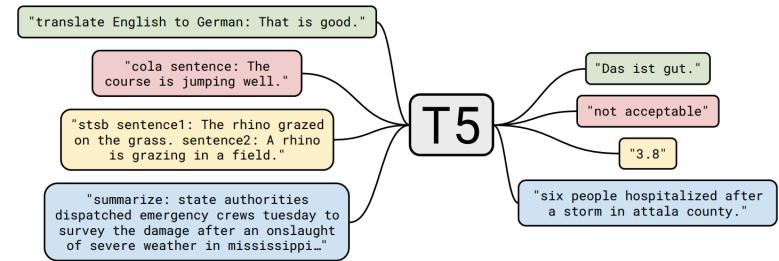
Number of tokens	Repeats	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Full data set	0	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
$2^{29}$	64	<b>82.87</b>	<b>19.19</b>	<b>80.97</b>	<b>72.03</b>	<b>26.83</b>	<b>39.74</b>	<b>27.63</b>
$2^{27}$	256	82.62	<b>19.20</b>	79.78	69.97	<b>27.02</b>	<b>39.71</b>	27.33
$2^{25}$	1,024	79.55	18.57	76.27	64.76	26.38	39.56	26.80
$2^{23}$	4,096	76.34	18.33	70.92	59.29	26.37	38.84	25.81

# What's the best possible model for text-to-text learning?

- **Fine-tuning**
- Recall from last week:
  - You can choose which weights to update
  - Updating all weights seems to work better for T5

Fine-tuning method	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ All parameters	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
Adapter layers, $d = 32$	80.52	15.08	79.32	60.40	13.84	17.88	15.54
Adapter layers, $d = 128$	81.51	16.62	79.47	63.03	19.83	27.50	22.63
Adapter layers, $d = 512$	81.54	17.78	79.18	64.30	23.45	33.98	25.81
Adapter layers, $d = 2048$	81.51	16.62	79.47	63.03	19.83	27.50	22.63
Gradual unfreezing	82.50	18.95	79.17	<b>70.79</b>	26.71	39.02	26.93

# Overall performance



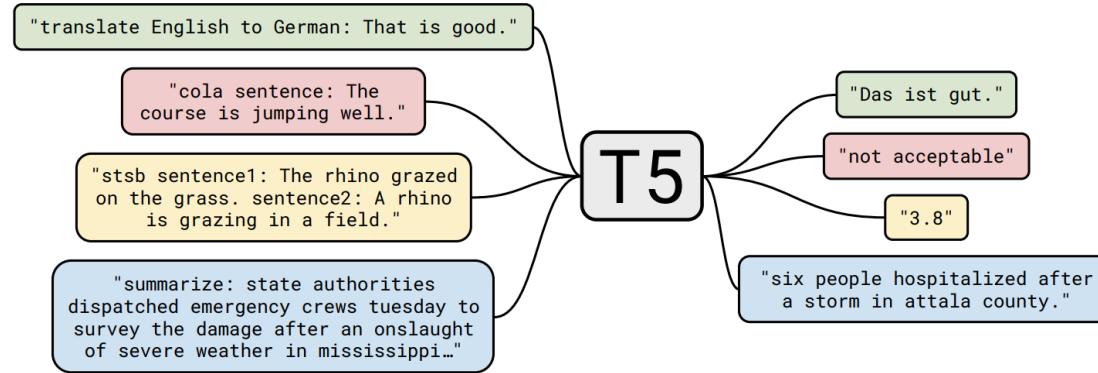
Model	GLUE	CoLA	SST-2	MRPC	MRPC	STS-B	STS-B
	Average	Matthew's	Accuracy	F1	Accuracy	Pearson	Spearman
Previous best	89.4 <sup>a</sup>	69.2 <sup>b</sup>	97.1 <sup>a</sup>	<b>93.6<sup>b</sup></b>	<b>91.5<sup>b</sup></b>	92.7 <sup>b</sup>	92.3 <sup>b</sup>
T5-Small	77.4	41.0	91.8	89.7	86.6	85.6	85.0
T5-Base	82.7	51.1	95.2	90.7	87.5	89.4	88.6
T5-Large	86.4	61.2	96.3	92.4	89.9	89.9	89.2
T5-3B	88.5	67.1	97.4	92.5	90.0	90.6	89.8
<b>T5-11B</b>	<b>90.3</b>	<b>71.6</b>	<b>97.5</b>	92.8	90.4	<b>93.1</b>	<b>92.8</b>

Model	QQP	QQP	MNLI-m	MNLI-mm	QNLI	RTE	WNLI
	F1	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Previous best	74.8 <sup>c</sup>	<b>90.7<sup>b</sup></b>	91.3 <sup>a</sup>	91.0 <sup>a</sup>	<b>99.2<sup>a</sup></b>	89.2 <sup>a</sup>	91.8 <sup>a</sup>
T5-Small	70.0	88.0	82.4	82.3	90.3	69.9	69.2
T5-Base	72.6	89.4	87.1	86.2	93.7	80.1	78.8
T5-Large	73.9	89.9	89.9	89.6	94.8	87.2	85.6
T5-3B	74.4	89.7	91.4	91.2	96.3	91.1	89.7
<b>T5-11B</b>	<b>75.1</b>	90.6	<b>92.2</b>	<b>91.9</b>	96.9	<b>92.8</b>	<b>94.5</b>

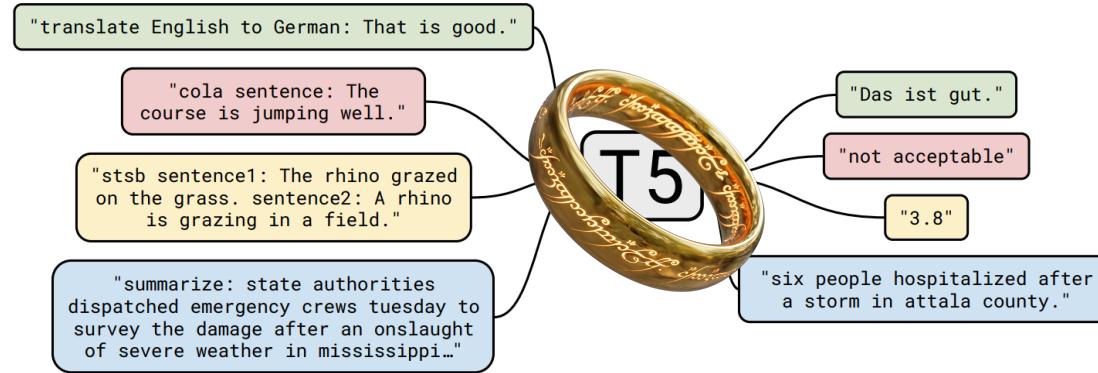
- Really good across the board. New SOTA on several tasks.
- More results in the paper

# T5: seq-to-seq everywhere



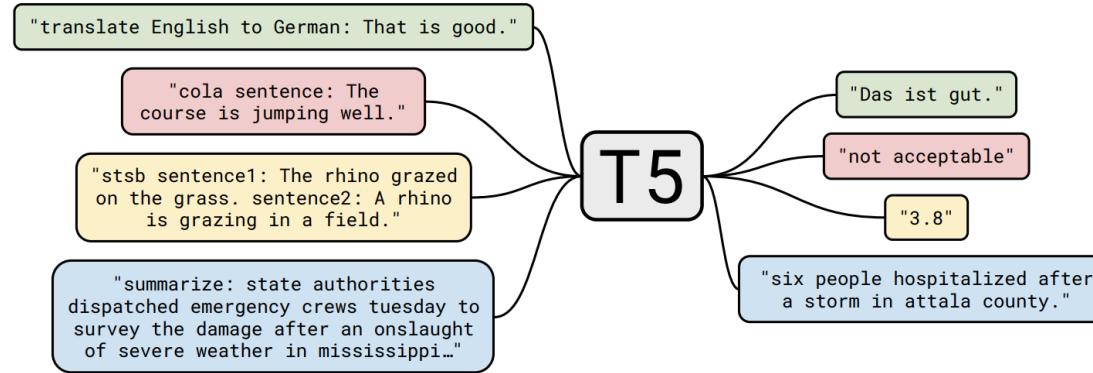
- Exhaustive search over the space of possible seq-to-seq models
- Very strong enc-dec model for NLP downstream tasks

# T5: seq-to-seq everywhere



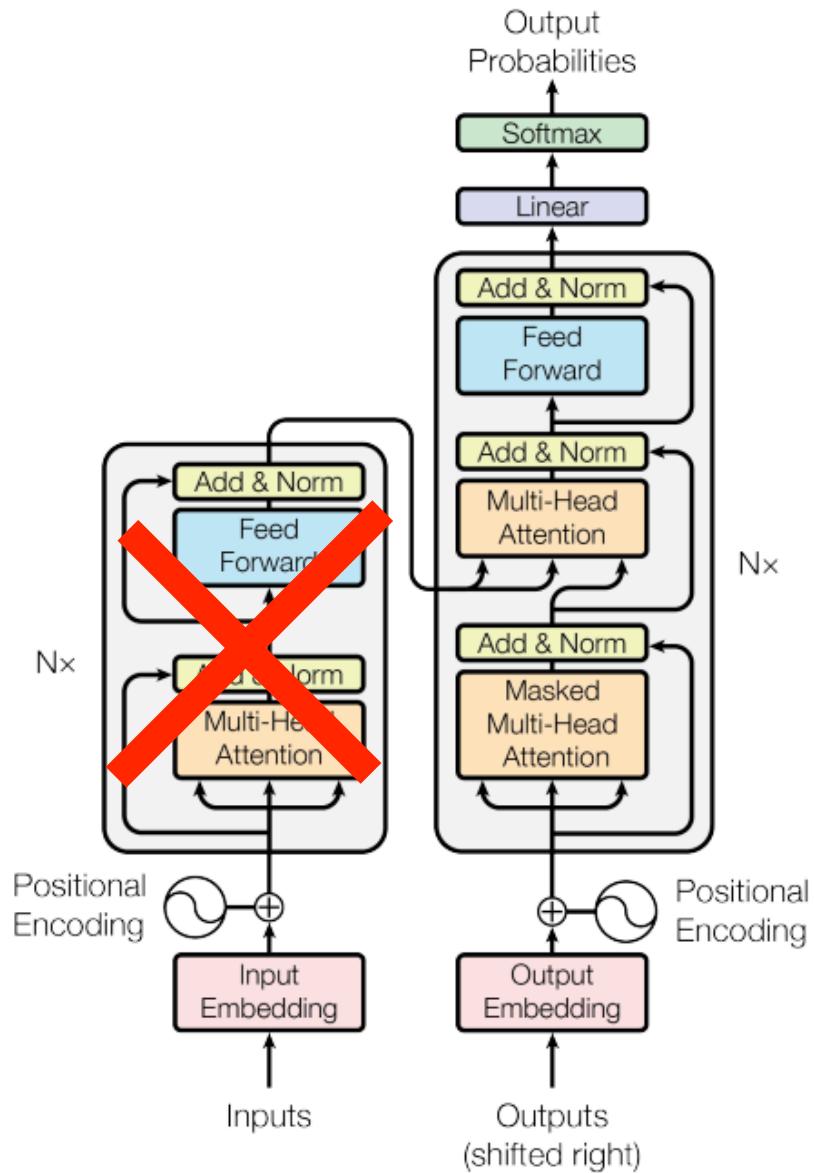
- Paradigm-shift: **One model to rule them all?**
- Preview for next week:
  - ~> This can be simplified and pushed even further ...

# Can we simplify this even more?

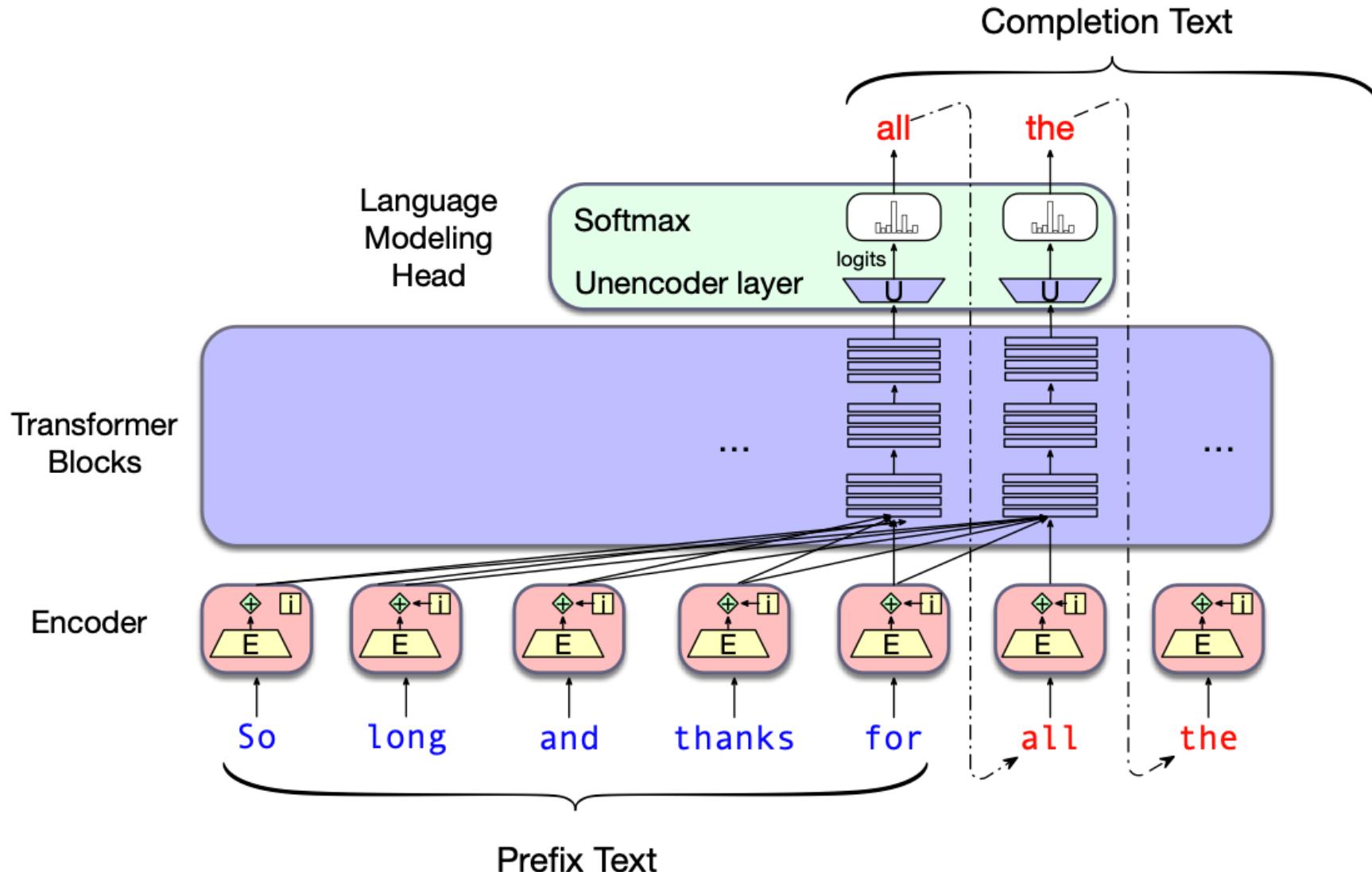


# Big idea

- Let's keep only the decoder
- Many tasks are text-to-text
- Learn them by predicting the next word



# Conditional Generation: Generating text conditioned on previous text!



# Many practical NLP tasks can be cast as word prediction!

Sentiment analysis: "I like Jackie Chan"

1. We give the language model this string:

The sentiment of the sentence "I like Jackie Chan" is:

2. And see what word it thinks comes next:

$P(\text{positive} | \text{The sentiment of the sentence } "I \text{ like Jackie Chan}" \text{ is:})$

$P(\text{negative} | \text{The sentiment of the sentence } "I \text{ like Jackie Chan}" \text{ is:})$

# Framing lots of tasks as conditional generation

QA: “Who wrote The Origin of Species”

1. We give the language model this string:

Q: Who wrote the book ‘‘The Origin of Species”? A:

2. And see what word it thinks comes next:

$P(w|Q: \text{Who wrote the book ‘‘The Origin of Species”? } A:)$

3. And iterate:

$P(w|Q: \text{Who wrote the book ‘‘The Origin of Species”? } A: \text{ Charles})$

# Summarization

Original

The only thing crazier than a guy in snowbound Massachusetts boxing up the powdery white stuff and offering it for sale online? People are actually buying it. For \$89, self-styled entrepreneur Kyle Waring will ship you 6 pounds of Boston-area snow in an insulated Styrofoam box – enough for 10 to 15 snowballs, he says.

But not if you live in New England or surrounding states. “We will not ship snow to any states in the northeast!” says Waring’s website, [ShipSnowYo.com](http://ShipSnowYo.com). “We’re in the business of expunging snow!”

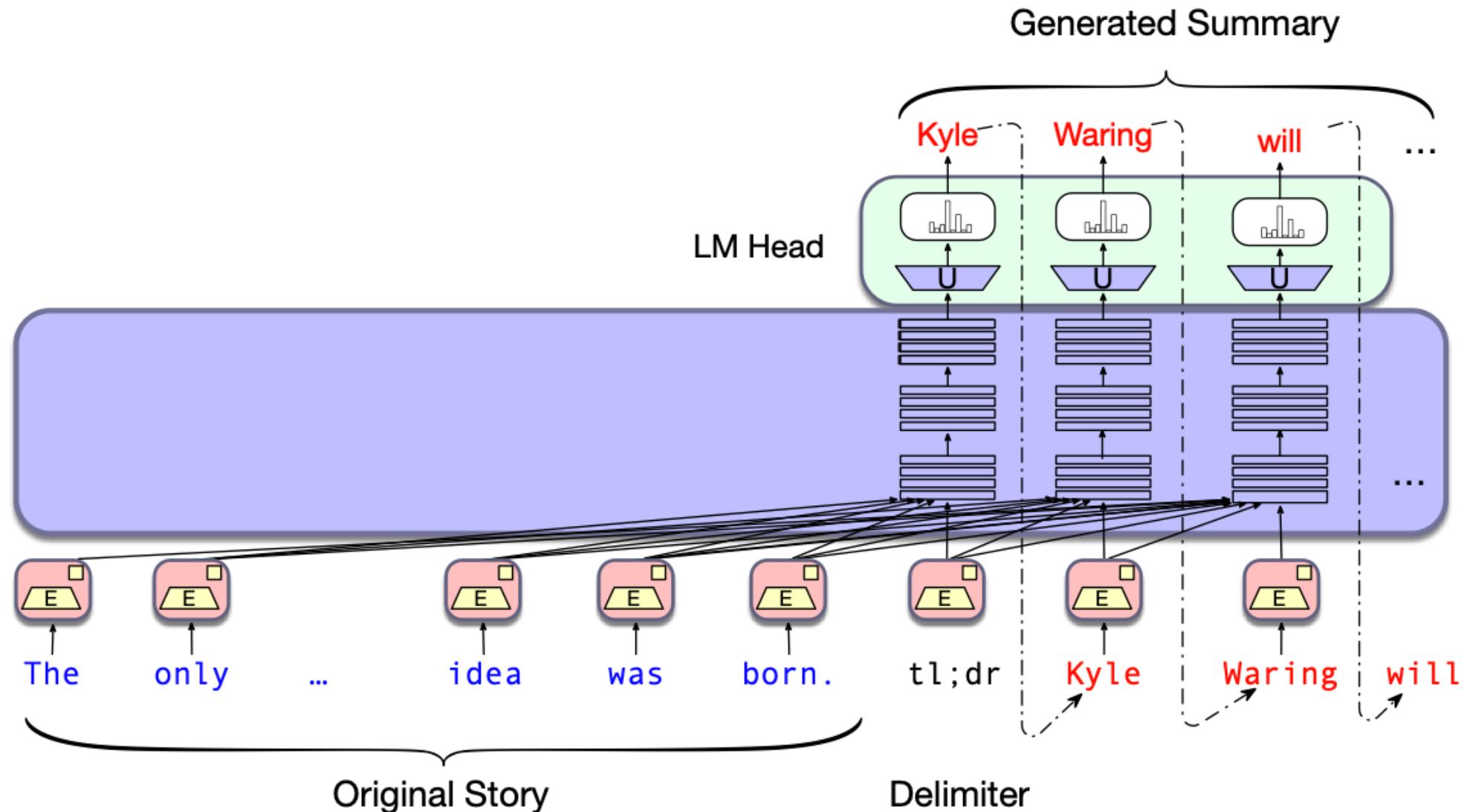
His website and social media accounts claim to have filled more than 133 orders for snow – more than 30 on Tuesday alone, his busiest day yet. With more than 45 total inches, Boston has set a record this winter for the snowiest month in its history. Most residents see the huge piles of snow choking their yards and sidewalks as a nuisance, but Waring saw an opportunity.

According to Boston.com, it all started a few weeks ago, when Waring and his wife were shoveling deep snow from their yard in Manchester-by-the-Sea, a coastal suburb north of Boston. He joked about shipping the stuff to friends and family in warmer states, and an idea was born. [...]

Summary

Kyle Waring will ship you 6 pounds of Boston-area snow in an insulated Styrofoam box – enough for 10 to 15 snowballs, he says. But not if you live in New England or surrounding states.

# LLMs for summarization (using tl;dr)



# **Reminder: Self-supervised training algorithm**

We just train them to predict the next word!

1. Take a corpus of text
2. At each time step  $t$ 
  - i. ask the model to predict the next word
  - ii. train the model using gradient descent to minimize the error in this prediction

**"Self-supervised"** because it just uses the next word as the label!

# Reminder: Cross-entropy loss for language modeling

**CE loss:** difference between the correct probability distribution and the predicted distribution

$$L_{CE} = - \sum_{w \in V} \mathbf{y}_t[w] \log \hat{\mathbf{y}}_t[w]$$

The correct distribution  $\mathbf{y}_t$  knows the next word, so is 1 for the actual next word and 0 for the others.

So in this sum, all terms get multiplied by zero except one: the logp the model assigns to the correct next word, so:

$$L_{CE}(\hat{\mathbf{y}}_t, \mathbf{y}_t) = -\log \hat{\mathbf{y}}_t[w_{t+1}]$$

# Reminder: Perplexity

Just as for n-gram grammars, we use perplexity to measure how well the LM predicts unseen text

The perplexity of a model  $\theta$  on an unseen test set is the **inverse probability that  $\theta$  assigns to the test set, normalized by the test set length.**

For a test set of  $n$  tokens  $w_{1:n}$  the perplexity is :

$$\begin{aligned}\text{Perplexity}_\theta(w_{1:n}) &= P_\theta(w_{1:n})^{-\frac{1}{n}} \\ &= \sqrt[n]{\frac{1}{P_\theta(w_{1:n})}}\end{aligned}$$

# Reminder: Perplexity

Just as for n-gram grammars, we use perplexity to measure how well the LM predicts unseen text

The perplexity of a model  $\theta$  on an unseen test set is the **inverse probability that  $\theta$  assigns to the test set, normalized by the test set length.**

For a test set of  $n$  tokens  $w_{1:n}$  the perplexity is :

$$\begin{aligned}\text{Perplexity}_\theta(w_{1:n}) &= P_\theta(w_{1:n})^{-\frac{1}{n}} \\ &= \sqrt[n]{\frac{1}{P_\theta(w_{1:n})}}\end{aligned}$$

# Why perplexity instead of raw probability of the test set?

- Probability depends on size of test set
  - Probability gets smaller the longer the text
  - Better: a metric that is **per-word**, normalized by length
- **Perplexity** is the inverse probability of the test set, normalized by the number of words
  - (The inverse comes from the original definition of perplexity from cross-entropy rate in information theory)
- Probability range is  $[0,1]$ , perplexity range is  $[1,\infty]$

# Reminder: Perplexity

Alternative view on perplexity

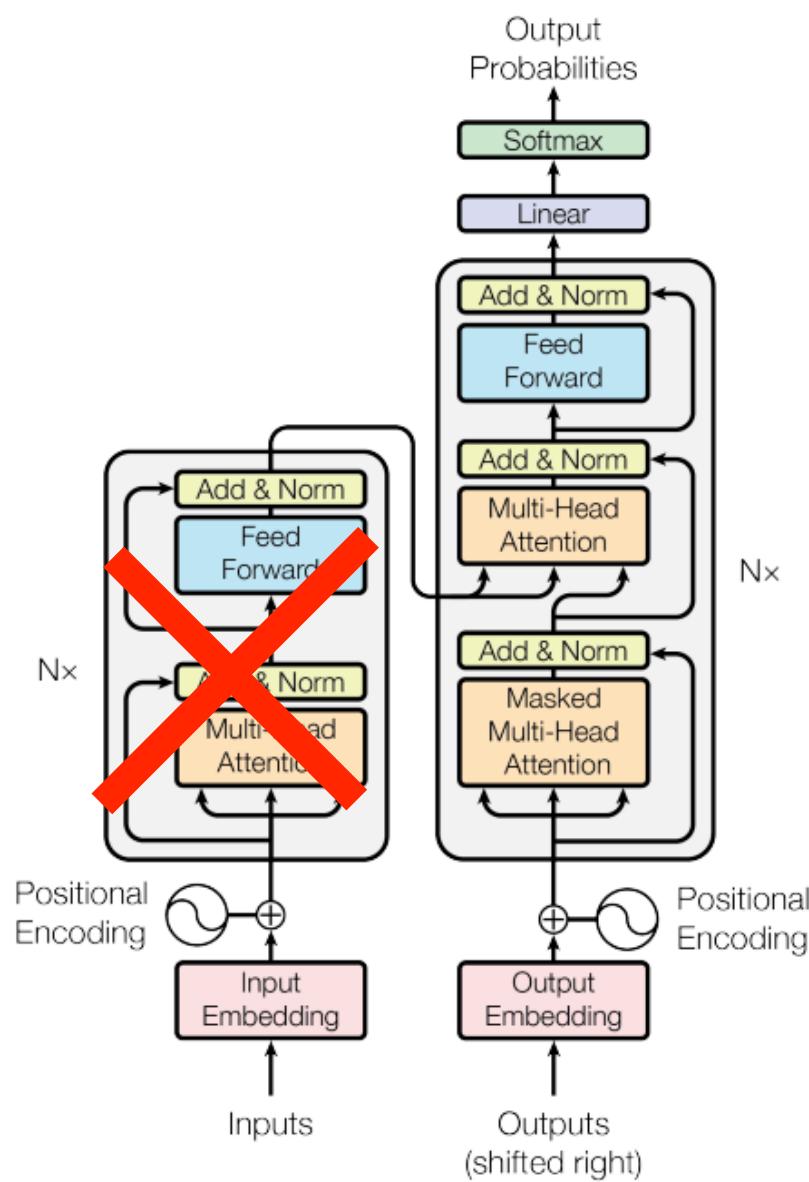
- Intuitively, ppl can be seen as a measure of among how many tokens a models has to choose when predicting the next tokens
- $\text{ppl} = 2^{H(x)}$ , where Cross-entropy is an upper-bound of the entropy:  
 $H(p) \leq H(p, m)$
- $$H(x) = -\frac{1}{N} \log P(x_1, \dots, x_n)$$
- Hence, in practice we simply exponentiate the loss of our model to get perplexity
- Hyper-parameter: how do you choose your context?

# Perplexity

- The higher the probability of the word sequence, the lower the perplexity.
- Thus the lower the perplexity of a model on the data, the better the model.
- **Minimizing perplexity is the same as maximizing probability**
- Also: perplexity is sensitive to length/tokenization so best used when comparing LMs that use the same tokenizer.

# Let's go back to the big idea

- Let's keep only the **decoder**
- Many tasks are text-to-text
- Learn them by predicting the next word



**Why would this work?**

# Pre-training

The technique that underlies (most of) the amazing performance of language models

First **pre-train** a transformer model on  
**enormous amounts** of text  
Then **apply** it to new tasks.

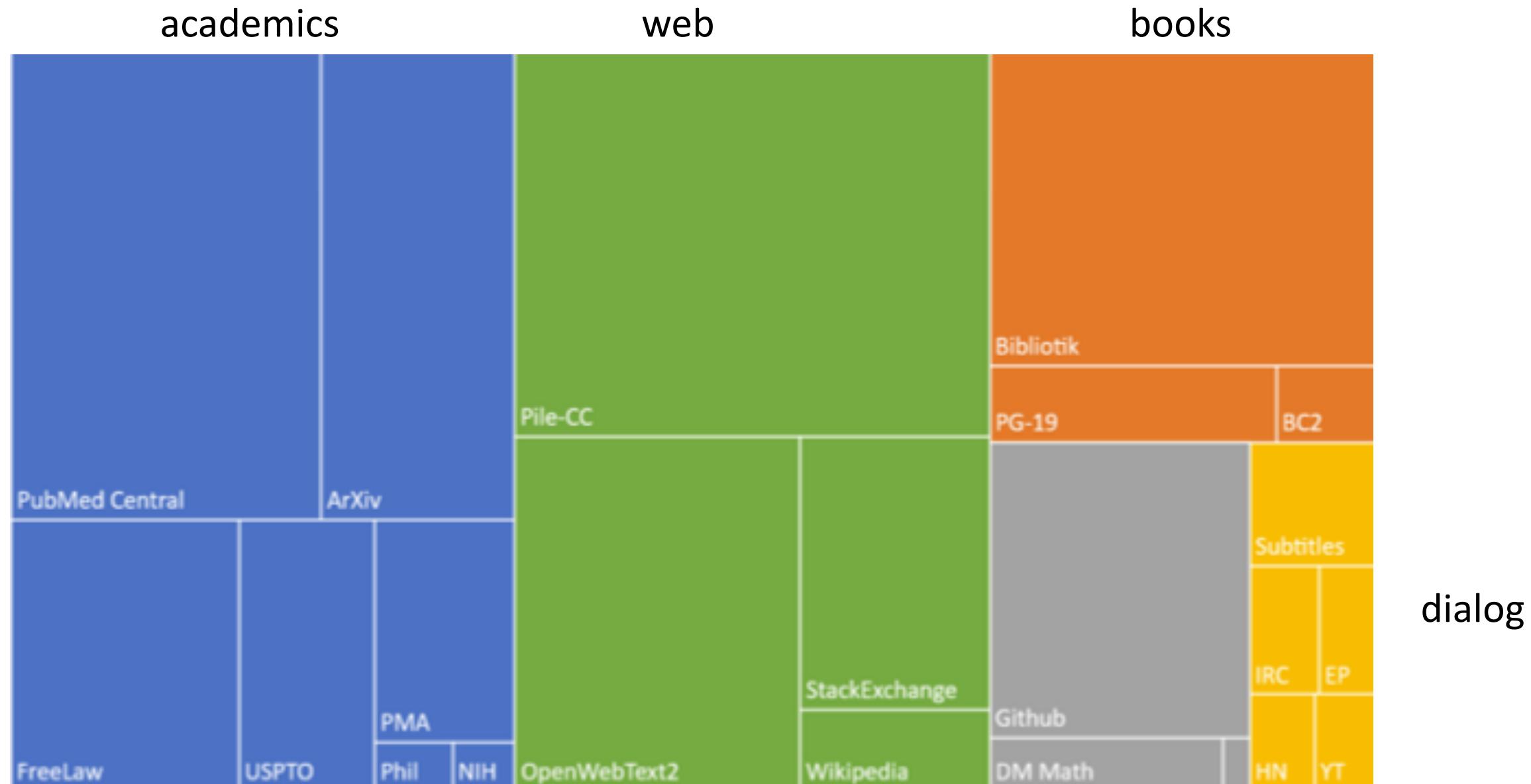
# LLMs are mainly trained on the web

Common crawl, snapshots of the entire web produced by the non- profit Common Crawl with billions of pages

Colossal Clean Crawled Corpus (C4; [Raffel et al. 2020](#)),  
156 billion tokens of English, filtered

What's in it? Mostly patent text documents, Wikipedia,  
and news sites

# The Pile: a pre-training corpus



# What does a model learn from pretraining?

- There are canines everywhere! One dog in the front room, and two dogs
- It wasn't just big it was enormous
- The author of "A Room of One's Own" is Virginia Woolf
- The doctor told me that he
- The square root of 4 is 2

# Filtering for quality and safety

Quality is subjective

- Many LLMs attempt to match Wikipedia, books, particular websites
- Need to remove boilerplate, adult content
- Deduplication at many levels (URLs, documents, even lines)

Safety also subjective

- Toxicity detection is important, although that has mixed results
- Can mistakenly flag data written in dialects like African American English

# **But there are problems with scraping from the web**

**Copyright:** much of the text in these datasets is copyrighted

- Not clear if fair use doctrine in US allows for this use
- This remains an open legal question

**Data consent**

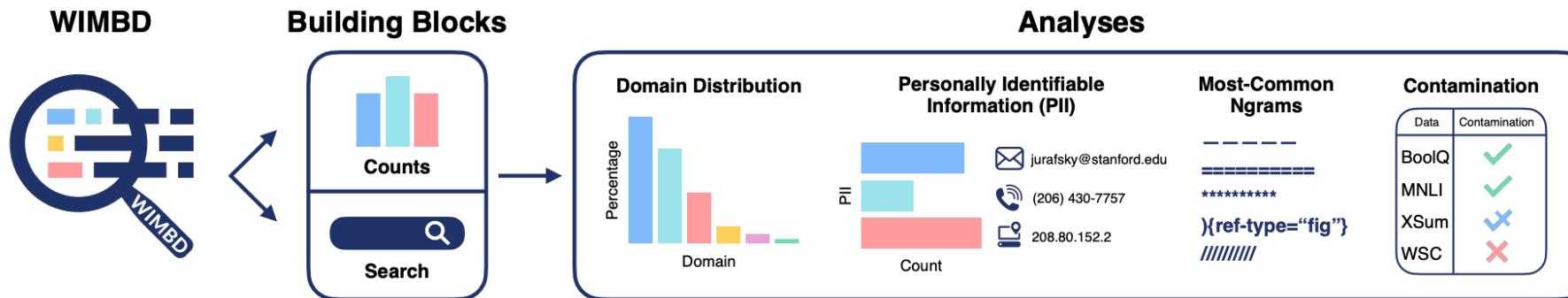
- Website owners can indicate they don't want their site crawled

**Privacy:**

- Websites can contain private IP addresses and phone numbers

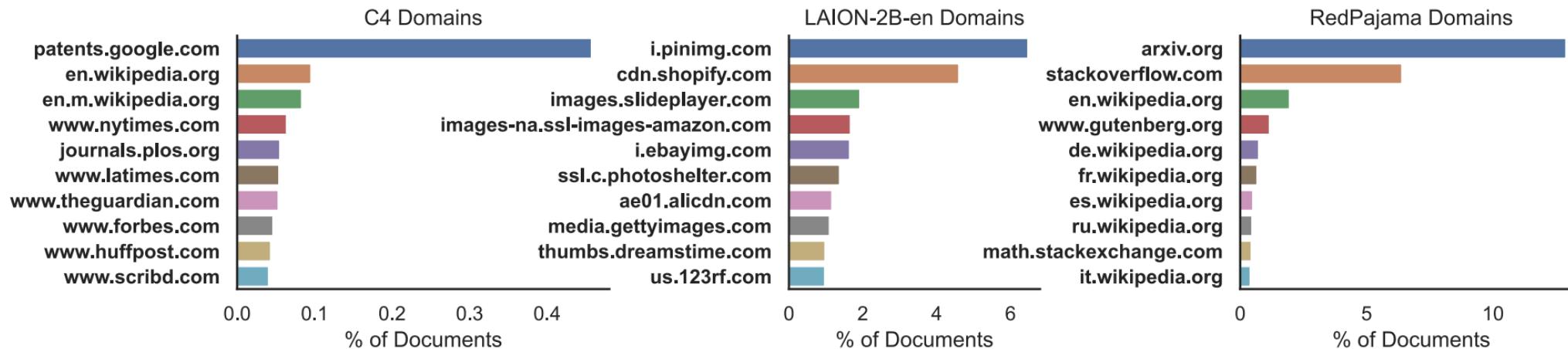
# What's in my big data?

- A tool to count and search in pre-training corpora
- Allows you to investigate the pre-training data of a model
- Alternatives: infinigram (Liu et al. 2024; <https://infini-gram.io/>)



# What's in my big data?

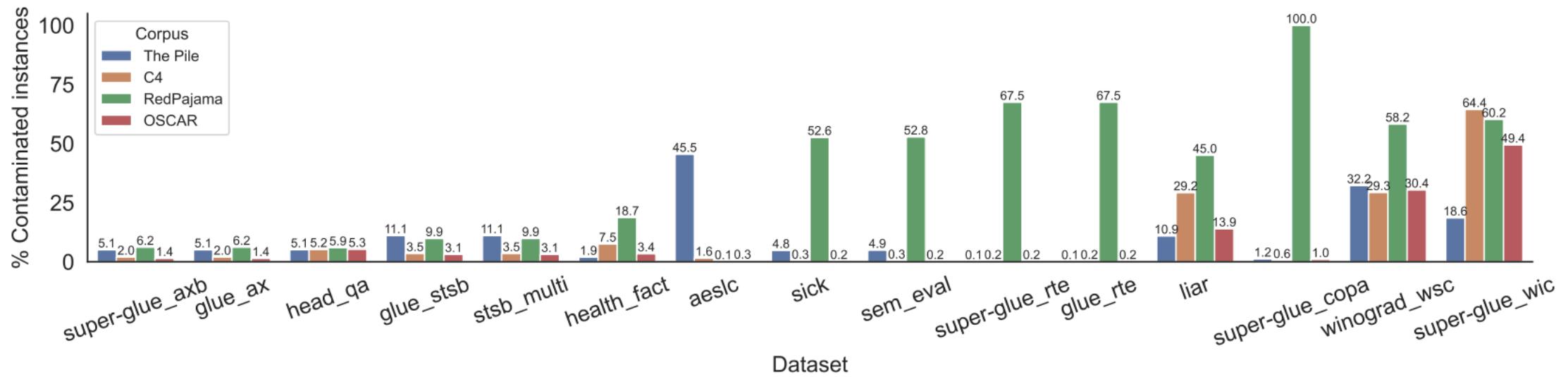
Elazar et al. (2024) - What's in my big data?



	OpenWebText	C4	mC4-en	OSCAR	The Pile
<i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram
-----	3.4M   ? ? ? ? ? ? ? ?	9M   -----	1.76B   \\\\\\\\\\\\\\\\\	773M   -----	3.64B   -----
.....	1.05M   .....	7.27M   -----	823M   -----	395M   =====	602M   =====
=====	830K   -----	4.41M   ♦♦♦♦♦♦♦♦♦♦♦♦♦♦♦♦	349M   -----	175M   * * * * * * * * *	188M   * * * * * * * * *
* * * * * * * * *	595K   * * * * * * * * *	3.87M   * * * * * * * * *	314M   -----	91.6M   ) { ref - type = " fig " }	59.1M   ) { ref - type = " fig " }
# ##### # ##### #	302K   ! ! ! ! ! ! ! !	1.91M   \ / s \ / files \ / 1 \ /	183M   * * * * * * * * *	34.9M   // / / / / / /	56.2M   // / / / / / /
amp ; amp ; amp ; amp ; amp ;	278K   . You can follow any responses to this entry through	784K   / s \ / files \ / 1 \ /	183M   =====	22.9M   .....	54.9M   .....
; amp ; amp ; amp ; amp ; amp	265K   ♦♦♦♦♦♦♦♦♦♦♦♦♦♦♦♦	753K   \ / \ / cdn.shopify.com \ / s \ /	182M   ( Opens in new window ) Click to share on	15.7M   ##### # ##### #	38.3M   ##### # ##### #
-----	249K   You can follow any responses to this entry through the	752K   cdn.shopify.com \ / s \ / files \ /	182M   Log Out / Change ) You are commenting using your	13.6M   -----	30.1M   -----
-----	88.1K   You can follow any responses to this entry through the RSS	752K   \ / cdn.shopify.com \ / s \ / files \ /	182M   ( Log Out / Change ) You are commenting using	13.6M   { ref - type = " fig " } )	28.9M   { ref - type = " fig " } )
-----	83.3K   follow any responses to this entry through the RSS 2.0	748K   \ / cdn.shopify.com \ / s \ / files	182M   . ( Log Out / Change ) You are commenting	13.6M   =====	21.8M   =====
	RedPajama	S2ORC	peS2o	LAION-2B-en	The Stack
<i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram	Count   <i>n</i> -gram
.....	670M   q q q q q q q q q	30.2M   .....	1.42M   -----	1.65M   -----	4.29B   -----
-----	507M   .....	5.49M   [ 1 ] [ 2 ] [ 3 ] [	457K   ♦♦♦♦♦♦♦♦♦♦♦♦♦♦♦♦	1.43M   * * * * * * * * *	3.87B   * * * * * * * * *
\\\\\\\\\\\\\\\\\	213M   + + + + + + + + +	3.03M   [ 1 ] [ 2 ] [ 3 ] [ 4 ]	453K   .....	1.15M   0 0 0 0 0 0 0 0	2.75B   0 0 0 0 0 0 0 0
* * * * * * * * *	195M   * * * * * * * * *	1.93M   [ 1 ] [ 2 ] [ 3 ] [ 4 ]	453K   \\\\\\\\\\\\\\\\\	809K   =====	2.62B   =====
=====	145M   oooooooo	1.73M   [ 5 ] [ 6 ] [ 7 ] [	450K   < br / > < br / > < br	797K   , " resolved " : " https : /	1.46B   , " resolved " : " https : /
// / / / / /	79.3M   .....	1.56M   [ 6 ] [ 7 ] [ 8 ] [	448K   / > < br / > < br / >	796K   " resolved " : " https : /	1.46B   " resolved " : " https : /
. / . / . / .	35.3M   .....	1.11M   [ 6 ] [ 7 ] [ 8 ]	448K   br / > < br / > < br /	796K   " resolved " : " https : / / registry.npmjs.org	1.42B   " resolved " : " https : / / registry.npmjs.org /
. / . / . / .	35.3M   [ 5 ] [ 6 ] [ 7 ] [	646K   [ 5 ] [ 6 ] [ 7 ] [ 8 ]	446K   > < br / > < br / > <	576K   resolved " : " https : / / registry.npmjs.org /	1.42B   resolved " : " https : / / registry.npmjs.org /
. / . / . / .	35.2M   [ 1 ] [ 2 ] [ 3 ] [	645K   [ 7 ] [ 8 ] [ 9 ]	446K   Price : 1 Credit ( USD \$ 1 )	437K   .....,	1B   .....,
##### # ##### #	33M   [ 6 ] [ 7 ] [ 8 ] [	644K   [ 6 ] [ 7 ] [ 8 ] [ 9 ]	444K   vector   Price : 1 Credit ( USD \$ 1 )	437K   .tgz " , " integrity " : " sha512	938M   .tgz " , " integrity " : " sha512

# What's in my big data?

Elazar et al. (2024) - What's in my big data?



**Once a model is trained, how can we generate text from it?**

# Decoding and Sampling

This task of choosing a word to generate based on the model's probabilities is called **decoding**.

The most common method for decoding in LLMs: **sampling**.

Sampling from a model's distribution over words:

- choose random words according to their probability assigned by the model.

After each token we'll sample words to generate according to their probability *conditioned on our previous choices*,

- A transformer language model will give the probability

# Random sampling

i  $\leftarrow$  1

$w_i \sim p(w)$

**while**  $w_i \neq \text{EOS}$

i  $\leftarrow$  i + 1

$w_i \sim p(w_i \mid w_{<i})$

# **Random sampling doesn't work very well**

Even though random sampling mostly generate sensible, high-probable words,

There are many odd, low- probability words in the tail of the distribution

Each one is low- probability but added up they constitute a large portion of the distribution

So they get picked enough to generate weird sentences

# **Factors in word sampling: quality and diversity**

**Emphasize high-probability words**

- + **quality**: more accurate, coherent, and factual,
- **diversity**: boring, repetitive.

**Emphasize middle-probability words**

- + **diversity**: more creative, diverse,
- **quality**: less factual, incoherent

# Top-k sampling:

1. Choose # of words  $k$
2. For each word in the vocabulary  $V$ , use the language model to compute the likelihood of this word given the context  $p(w_t | w_{<t})$
3. Sort the words by likelihood, keep only the top  $k$  most probable words.
4. **Renormalize** the scores of the  $k$  words to be a legitimate probability distribution.
5. Randomly sample a word from within these remaining  $k$  most-probable words according to its probability.

# Top-p sampling (= nucleus sampling)

Holtzman et al., 2020

Problem with top- $k$ :  $k$  is fixed so may cover very different amounts of probability mass in different situations

Idea: Instead, keep the top  $p$  percent of the probability mass

Given a distribution  $P(w_t | \mathbf{w}_{<t})$ , the top- $p$  vocabulary  $V^{(p)}$  is the smallest set of words such that

$$\sum_{w \in V^{(p)}} P(w | \mathbf{w}_{<t}) \geq p.$$

# Temperature sampling

Reshape the distribution instead of truncating it

Intuition from thermodynamics,

- a system at high temperature is flexible and can explore many possible states,
- a system at lower temperature is likely to explore a subset of lower energy (better) states.

In **low-temperature sampling**, ( $\tau \leq 1$ ) we smoothly

- increase the probability of the most probable words
- decrease the probability of the rare words.

# Temperature sampling

Divide the logit by a temperature parameter  $\tau$  before passing it through the softmax.

Instead of

$$\mathbf{y} = \text{softmax}(u)$$

We do

$$\mathbf{y} = \text{softmax}(u/\tau)$$

# Temperature sampling

$$0 \leq \tau \leq 1$$

$$\mathbf{y} = \text{softmax}(u/\tau)$$

Why does this work?

- When  $\tau$  is close to 1 the distribution doesn't change much.
- The lower  $\tau$  is, the larger the scores being passed to the softmax
- Softmax pushes high values toward 1 and low values toward 0.
- Large inputs pushes high-probability words higher and low probability word lower, making the distribution more greedy.
- As  $\tau$  approaches 0, the probability of most likely word approaches 1

Show this in VSCode

**How can we guide the generations of a model?**

# Prompting and in-context learning (ICL)

- A prompt refers to the context provided to the language model
- Various components of a prompt, e.g., instruction, description, etc.

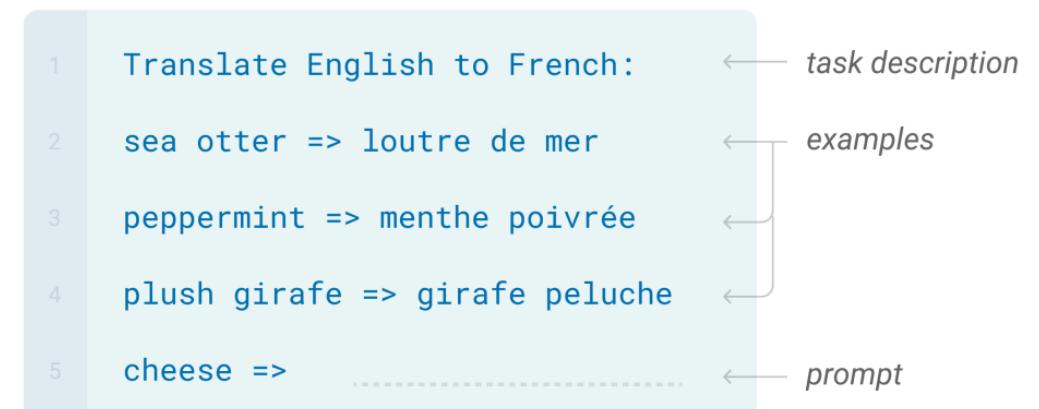
## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

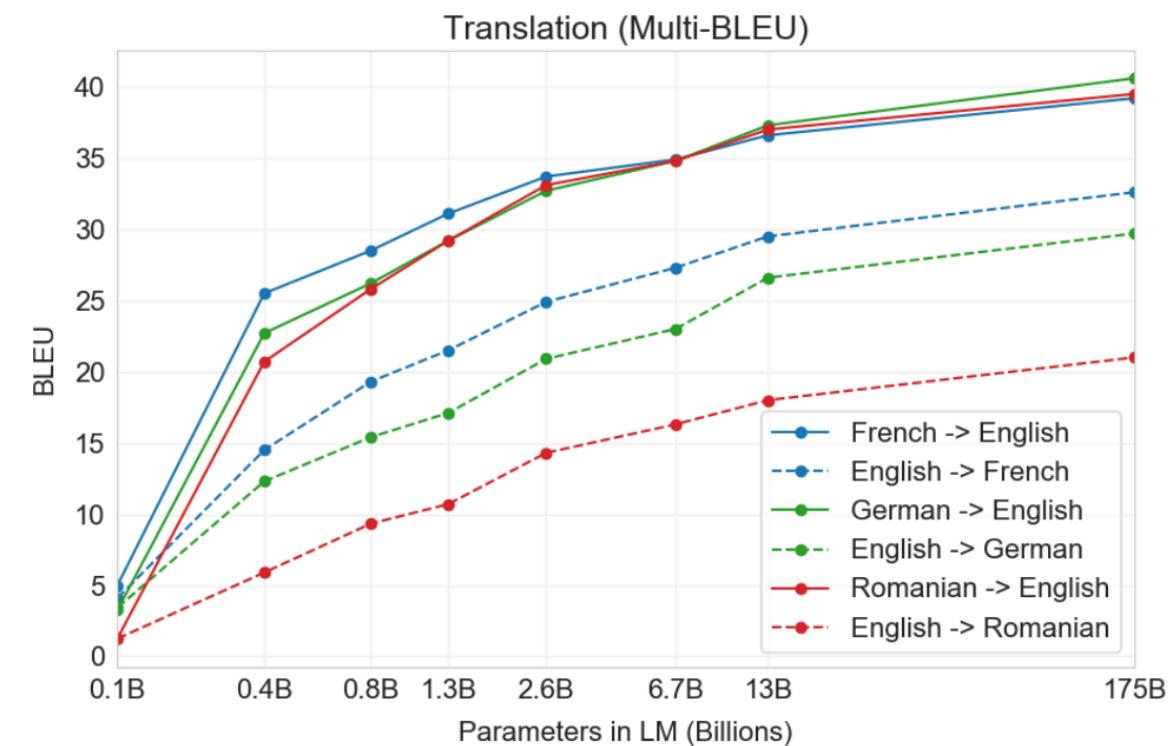
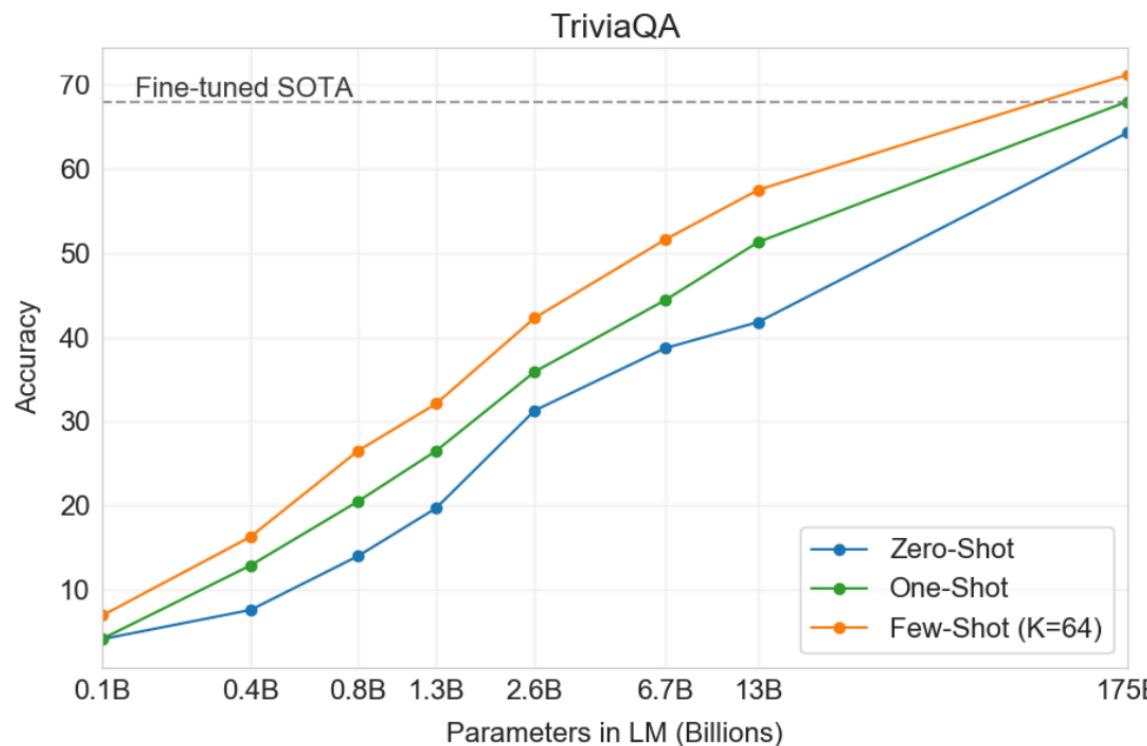


## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

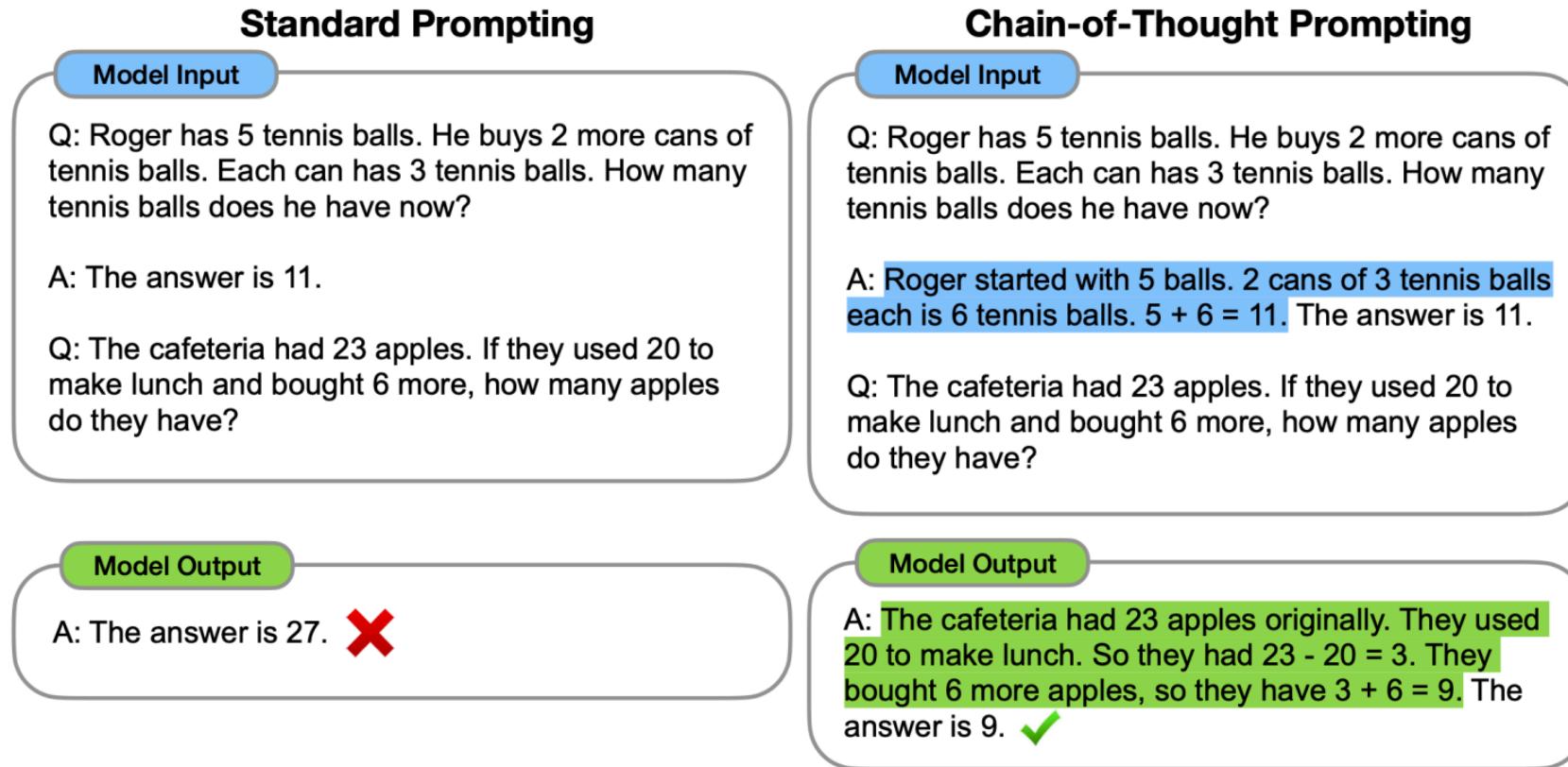


# How good are LLMs?



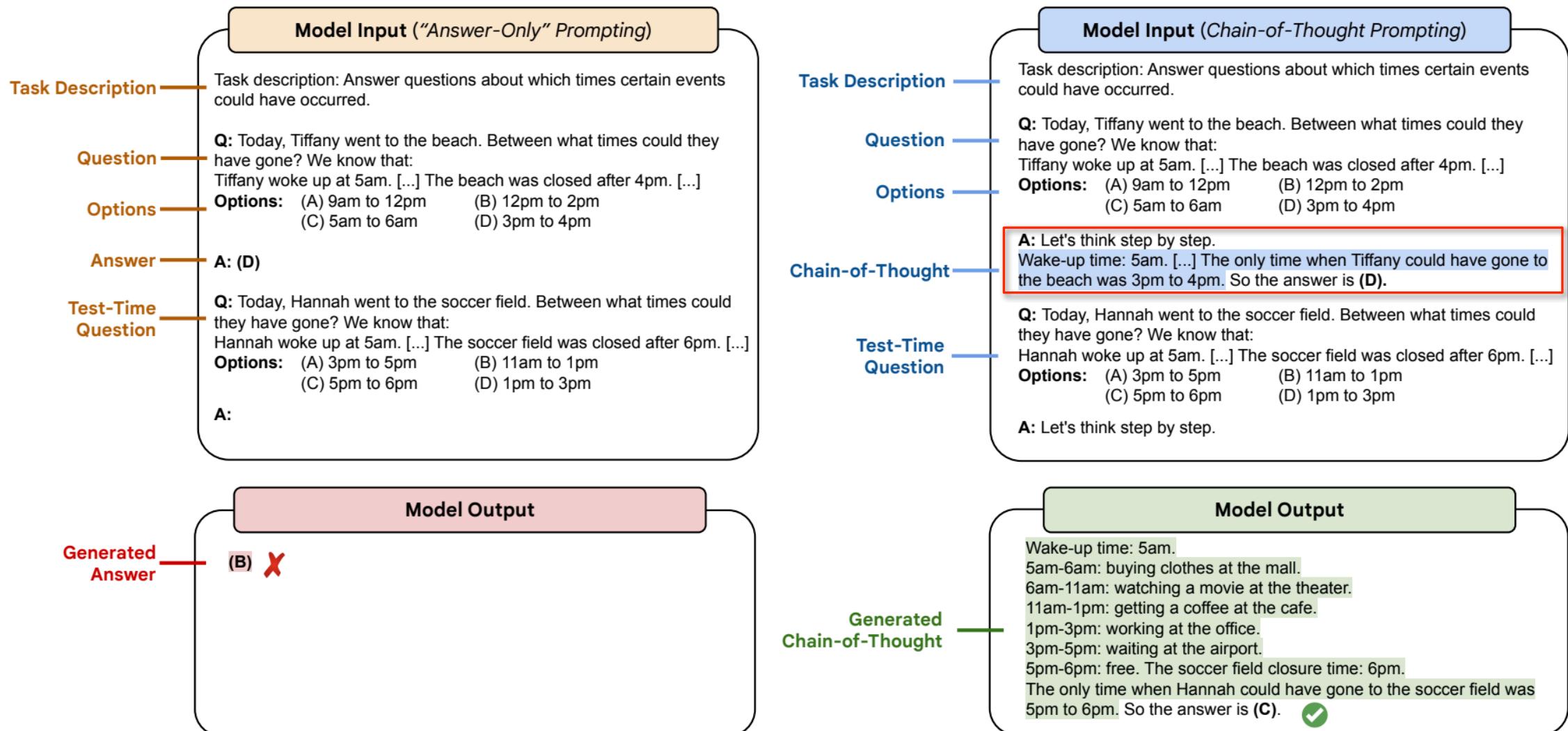
- As model increase in size, zero-shot and few-shot prompting performs better

# Chain-of-thought prompting (CoT)



Main idea: show a model how to decompose a problem into simpler sub-problems

# Chain-of-thought prompting (CoT)



# **Many other factors that we evaluate, like:**

## **Size**

Big models take lots of GPUs and time to train, memory to store

## **Energy usage**

Can measure kWh or kilograms of CO<sub>2</sub> emitted

## **Fairness**

Benchmarks measure gendered and racial stereotypes, or decreased performance for language from or about some groups.

# Scaling Laws

LLM performance depends on

- Model size: the number of parameters not counting embeddings
- Dataset size: the amount of training data
- Compute: Amount of compute (in FLOPS or etc)

Can improve a model by adding parameters (more layers, wider contexts), more data, or training for more iterations

The performance of a large language model (the loss) scales as a power-law with each of these three

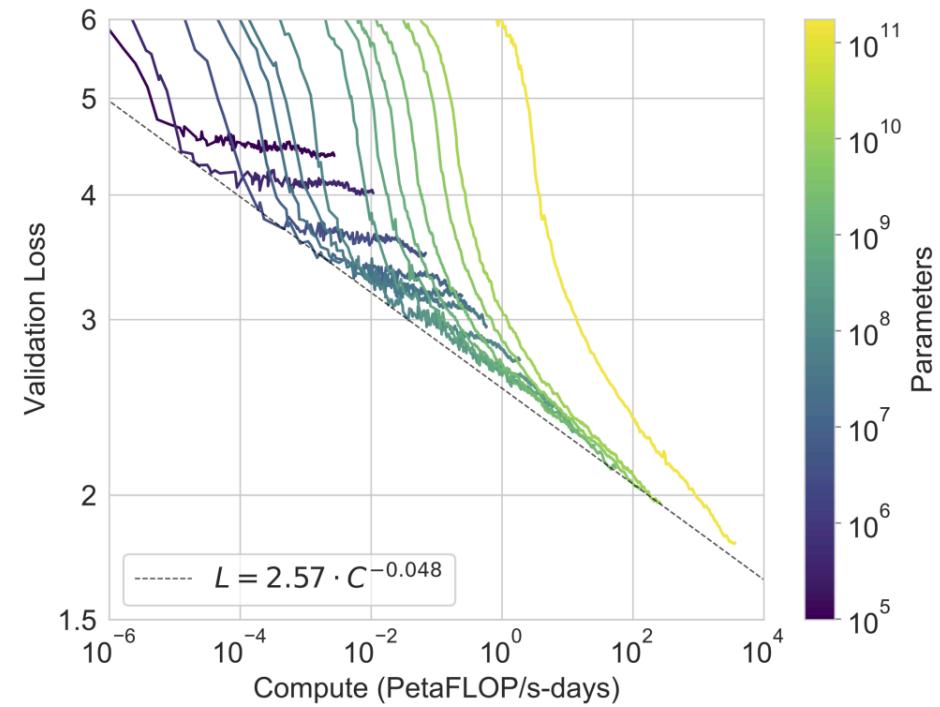
# Scaling Laws

Loss  $L$  as a function of # parameters  $N$ , dataset size  $D$ , compute budget  $C$  (if other two are held constant)

$$L(N) = \left(\frac{N_c}{N}\right)^{\alpha_N}$$

$$L(D) = \left(\frac{D_c}{D}\right)^{\alpha_D}$$

$$L(C) = \left(\frac{C_c}{C}\right)^{\alpha_C}$$



Scaling laws can be used early in training to **predict** what the loss would be if we were to add more data or increase model size.

# Number of non-embedding parameters $N$

$$\begin{aligned} N &\approx 2d n_{\text{layer}}(2d_{\text{attn}} + d_{\text{ff}}) \\ &\approx 12n_{\text{layer}} d^2 \\ &\quad (\text{assuming } d_{\text{attn}} = d_{\text{ff}}/4 = d) \end{aligned}$$

Thus GPT-3, with  $n = 96$  layers and dimensionality  $d = 12288$ , has  $12 \times 96 \times 12288^2 \approx 175$  billion parameters.

# Hallucination

*Chatbots May ‘Hallucinate’  
More Often Than Many Realize*

## *What Can You Do When A.I. Lies About You?*

People have little protection or recourse when the technology creates and spreads falsehoods about them.

### **Air Canada loses court case after its chatbot hallucinated fake policies to a customer**

The airline argued that the chatbot itself was liable. The court disagreed.

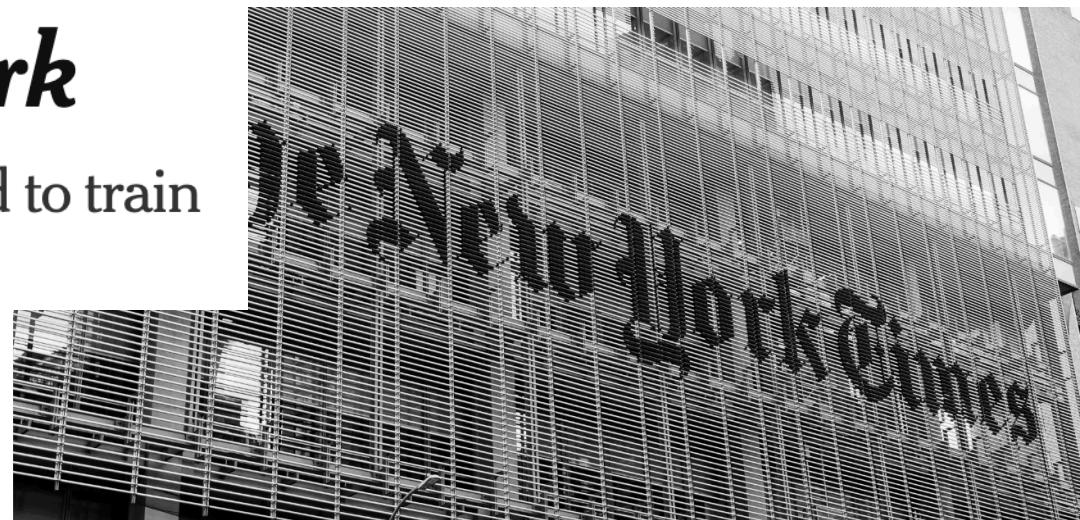
# Copyright



Authors Sue OpenAI Claiming Mass Copyright Infringement of Hundreds of Thousands of Novels

## *The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work*

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.



# Privacy

How Strangers Got My Email  
Address From ChatGPT's Model

# Toxicity and Abuse

The New AI-Powered Bing Is Threatening Users.

**Cleaning Up ChatGPT Takes Heavy Toll on Human Workers**

Contractors in Kenya say they were traumatized by effort to screen out descriptions of violence and sexual abuse during run-up to OpenAI's hit chatbot

# Misinformation

**Chatbots are generating false and misleading information about U.S. elections**

**[15 minute break]**