# DATA 8005 Advanced Natural Language Processing

## Lecture 2: Introduction to LLMs

Fall 2024

# Announcements

- Sign up for in-class presentations and final projects
  - In-class presentation: by Sep 22
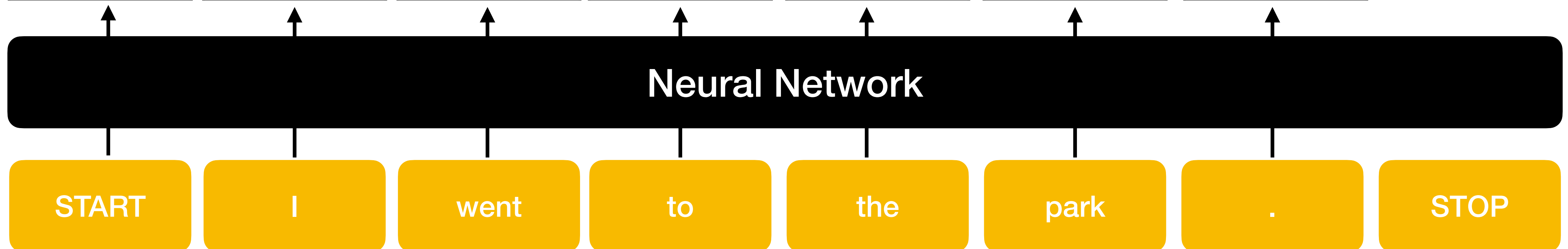  - Final projects: by Oct 4

# Neural language models: overview

# Neural language models: inputs/outputs

- **Input:** sequences of words (or tokens)

- **Output:** probability distribution over the next word (token)

| $p(x|\mathrm{START})$ | | $p(x|\mathrm{START\ I})$ | | $p(x|\cdots\mathrm{went})$ | | $p(x|\cdots\mathrm{to})$ | | $p(x|\cdots\mathrm{the})$ | | $p(x|\cdots\mathrm{park})$ | | $p(x|\mathrm{START\ I\ went\ to\ the\ park.})$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| The | 3 | think | 11% | to | 35% | the | 29% | bathroo | 3% | and | 14% | I | 21% |
| When | 2.5% | was | 5% | back | 8% | a | 9% | doctor | 2% | with | 9 | It | 6 |
| They | 2% | went | 2% | into | 5% | see | 5% | hospita | 2% | , | 8% | The | 3% |
| … | … | am | 1% | through | 4% | my | 3% | store | 1.5% | to | 7% | There | 3% |
| I | 1% | will | 1% | out | 3% | bed | 2% | … | … | … | … | … | … |
| … | … | like | 0.5% | on | 2% | school | 1% | park | 0.5% | . | 6% | STOP | 1% |
| Banana | 0.1% | … | … | … | …% | … | … | … | … | … | … | … | … |

**Neural Network**

START   I   went   to   the   park   .   STOP

# Usage example: GPT-3

# Usage example: GPT-3

# Neural language models: tokenization

# Why is tokenization important?

# Neural language models: inputs/outputs

- **Input:** sequences of words (or tokens)

- **Output:** probability distribution over the next word (token)

| $p(x\mid\text{START})$ | | $p(x\mid\text{START I})$ | | $p(x\mid\cdots\text{went})$ | | $p(x\mid\cdots\text{to})$ | | $p(x\mid\cdots\text{the})$ | | $p(x\mid\cdots\text{park})$ | | $p(x\mid\text{START I went to the park.})$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| The | 3 | think | 11% | to | 35% | the | 29% | bathroo | 3% | and | 14% | I | 21% |
| When | 2.5% | was | 5% | back | 8% | a | 9% | doctor | 2% | with | 9 | It | 6 |
| They | 2% | went | 2% | into | 5% | see | 5% | hospita | 2% | , | 8% | The | 3% |
| … | … | am | 1% | through | 4% | my | 3% | store | 1.5% | to | 7% | There | 3% |
| I | 1% | will | 1% | out | 3% | bed | 2% | … | … | … | … | … | … |
| … | … | like | 0.5% | on | 2% | school | 1% | park | 0.5% | . | 6% | STOP | 1% |
| Banana | 0.1% | … | … | … | …% | … | … | … | … | … | … | … | … |

**Neural Network**

START    I    went    to    the    park    .    STOP

# Neural language models: input vectors

*But neural networks take in real-valued vectors, not words…*

- Use one-hot or learned embeddings to map from words to vectors!
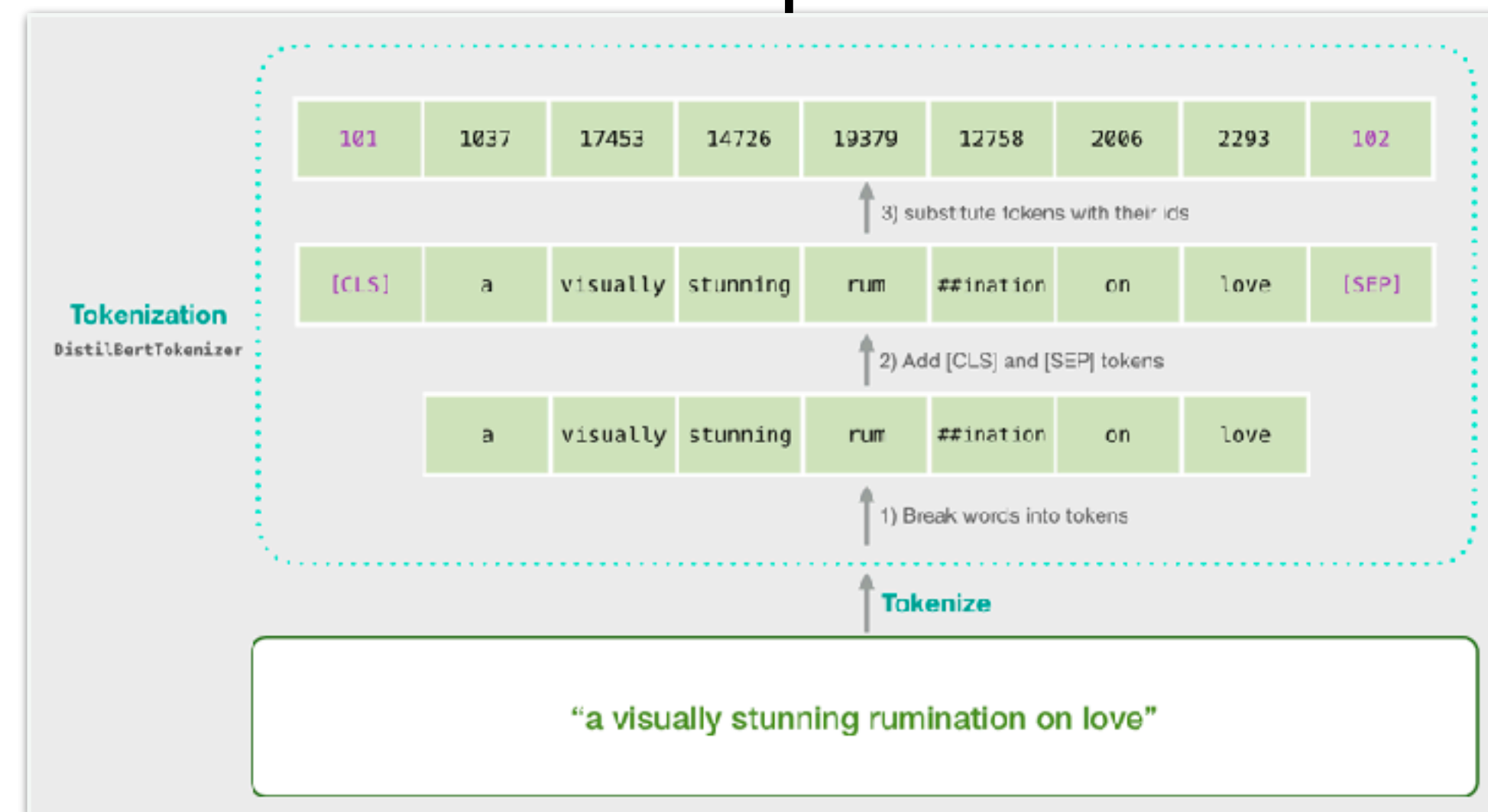  - Learned embeddings become part of parameters $\theta$

# Tokenization to input vectors

$p(x|\text{START})$ $p(x|\text{START I})$ $p(x|\cdots\text{went})$ $\quad p(x|\cdots\text{to})$ $\quad p(x|\cdots\text{the})$ $\quad p(x|\cdots\text{park})$ $\quad p(x|\text{START I went to the park.})$

## Neural Network

Mapping each tokenized id into its corresponding embeddings

Tokenization:



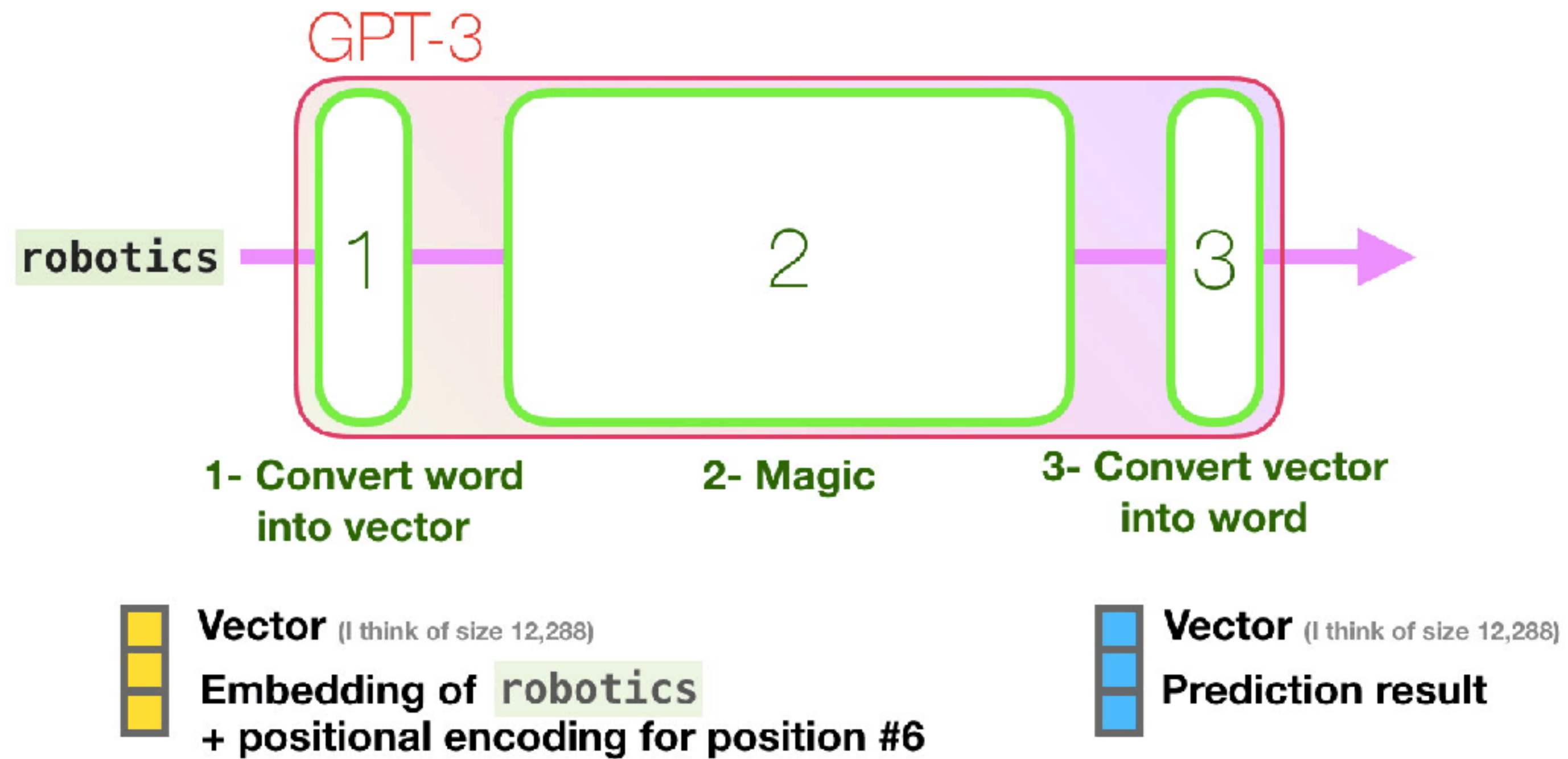| START | I | went | to | the | park | . | STOP |

# Usage example: GPT-3

# ChatGPT tokenization example



Call me Ishmael. Some years ago—never mind how long precisely—having little or no money in my purse, and nothing particular to interest me on shore, I thought I would sail about a little and see the watery part of the world. It is a way I have of driving off the spleen and regulating the circulation. Whenever I find myself growing grim about the mouth; whenever it is a damp, drizzly November in my soul; whenever I find myself involuntarily pausing before coffin warehouses, and bringing up the rear of every funeral I meet; and especially whenever my hypos get such an upper hand of me, that it requires a strong moral principle to prevent me from deliberately stepping into the street, and methodically knocking people's hats off—then, I account it high time tozz get to sea as soon as I can. This is my substitute for pistol and ball. With a philosophical flourish Cato throws himself upon his sword; I quietly take to the ship. There is nothing surprising in this. If they but knew it, almost all men in their degree, some time or other, cherish very nearly the same feelings towards the ocean with me.

TEXT     TOKEN IDS

**Tokens**
239

**Characters**
1109

[7368, 757, 57704, 1764, 301, 13, 4427, 1667, 4227, 2345, 37593, 4059, 1268, 1317, 24559, 2345, 69666, 2697, 477, 912, 3300, 304, 856, 53101, 11, 323, 4400, 4040, 311, 2802, 757, 389, 31284, 11, 358, 3463, 358, 1053, 30503, 922, 264, 2697, 323, 1518, 279, 30125, 727, 961, 315, 279, 1917, 13, 1102, 374, 264, 1648, 358, 617, 315, 10043, 1022, 279, 87450, 268, 323, 58499, 279, 35855, 13, 43633, 358, 1505, 7182, 7982, 44517, 922, 279, 11013, 26, 15716, 433, 374, 264, 41369, 11, 1377, 73825, 6841, 304, 856, 13836, 26, 15716, 358, 1505, 7182, 4457, 3935, 6751, 7251, 985, 1603, 78766, 83273, 11, 323, 12967, 709, 279, 14981, 315, 1475, 32079, 358, 3449, 26, 323, 5423, 15716, 856, 6409, 981, 636, 1778, 459, 8582, 1450, 315, 757, 11, 430, 433, 7612, 264, 3831, 16033, 17966, 311, 5471, 757, 505, 36192, 36567, 1139, 279, 8761, 11, 323, 1749, 2740, 50244, 1274, 753, 45526, 1022, 2345, 3473, 11, 358, 2759, 433, 1579, 892, 311, 10616, 636, 311, 9581, 439, 5246, 439, 358, 649, 13, 1115, 374, 856, 28779, 369, 40536, 323, 5041, 13, 3161, 264, 41903, 67784, 356, 4428, 3872, 5678, 5304, 813, 20827, 26, 358, 30666, 1935, 311, 279, 8448, 13, 2684, 374, 4400, 15206, 304, 420, 13, 1442, 814, 719, 7020, 433, 11, 4661, 682, 3026, 304, 872, 8547, 11, 1063, 892, 477, 1023, 11, 87785, 1890, 16024, 7119, 279, 18435, 449, 757, 13]

TEXT     **TOKEN IDS**

# Vocabulary: word-level

- For the n-gram model, our vocabulary $\mathcal{V}$ was comprised of all of the words in a language
- Some problems with this:
  - $|\mathcal{V}|$ **can be quite large** - ~470,000 words Webster's English Dictionary (3rd edition)
  - **Language is changing all of the time** - 690 words were added to Merriam Webster's in September 2023 ("rizz", "goated", "mid")
  - **Long tail of infrequent words**. Many words just occur a few times
  - **Some words may not appear** in a training set of documents
  - **No modeled relationship between words** - e.g., "run", "ran", "runs", "runner" are all separate entries despite being linked in meaning

# Character-level?

**What about representing text with characters?**

- $V = \{a, b, c, \ldots, z\}$

  - (Maybe add capital letters, punctuation, spaces, …)

- Pros:

  - Small vocabulary size ($|V| = 26$ for English)

  - Complete coverage (unseen words are represented by letters)

- Cons:

  - Encoding becomes very long - # chars instead of # words

  - Poor inductive bias for learning

# ~~Word~~ ~~Character~~ <u>Subword</u> tokenization!

How can we combine the high coverage of character-level representation with the efficiency of word-level representation?

**Subword tokenization! (e.g., Byte-Pair Encoding)**
- Start with character-level representations
- Build up representations from there

Original BPE Paper (Sennrich et al., 2016)
https://arxiv.org/abs/1508.07909

# Byte-pair encoding: ChatGPT example

# Byte-pair encoding: usage

- Basically state of the art in tokenization
- Used in all modern left-to-right large language models (LLMs), including ChatGPT

| Model/Tokenizer | Vocabulary Size |
|---|---|
| GPT-3.5/GPT-4/ChatGPT | 100k |
| GPT-2/GPT-3 | 50k |
| Llama2 | 32k |
| Falcon | 65k |

# Byte-pair encoding (BPE): algorithm

**Required:**

- Documents $\mathcal{D}$

- Desired vocabulary size $N$ (greater than characters in $\mathcal{D}$)

**Algorithm:**

- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation)

- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$

- Convert $\mathcal{D}$ into a list of tokens (characters)

- While $|\mathcal{V}| < N$:

  - Let $n := |\mathcal{V}| + 1$

  - Get counts of all bigrams in $\mathcal{D}$

  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)

    - Let $v_n := \mathrm{concat}(v_i, v_j)$

    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

# Byte-pair encoding: example

**Required:**

- Documents $\mathcal{D}$ $\longrightarrow$ $\mathcal{D} = \{$"i hug pugs", "hugging pugs is fun", "i make puns"$\}$

- Desired vocabulary size $N$ (greater than chars in $\mathcal{D}$) $\longrightarrow$ $N = 20$

**Algorithm:**

- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation) $\longrightarrow$

$\mathcal{D} = \{$"i", " hug", " pugs", "hugging", " pugs", " is", " fun", "i", " make", " puns"$\}$

- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$

- Convert $\mathcal{D}$ into a list of tokens (characters)

$\mathcal{V} = \{$' ', 'a', 'e', 'f', 'g', 'h', 'i', 'k', 'm', 'n', 'p', 's', 'u'$\}, |\mathcal{V}| = 13$

- While $|\mathcal{V}| < N$:

  - Let $n := |\mathcal{V}| + 1$

  - Get counts of all bigrams in $\mathcal{D}$

  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)

    - Let $v_n := \text{concat}(v_i, v_j)$

    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

$\mathcal{D} = \{$ ['i'], [' ', 'h', 'u', 'g'], [' ', 'p', 'u', 'g', 's'],
['h', 'u', 'g', 'g', 'i', 'n', 'g'], [' ', 'p', 'u', 'g', 's'],
[' ', 'i', 's'], [' ', 'f', 'u', 'n'], ['i'],
[' ', 'm', 'a', 'k', 'e'], [' ', 'p', 'u', 'n', 's']$\}$

Example inspired by: https://huggingface.co/docs/transformers/tokenizer_summary

# Byte-pair encoding: example

**Required:**

- Documents $\mathcal{D}$
- Desired vocabulary size $N$ (greater than chars in $\mathcal{D}$)

**Algorithm:**

- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation)
- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$
- Convert $\mathcal{D}$ into a list of tokens (characters)
- While $|\mathcal{V}| < N$:
  - Let $n := |\mathcal{V}| + 1$
  - Get counts of all bigrams in $\mathcal{D}$
  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)
    - Let $v_n := \mathrm{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

$$\mathcal{V} = \{1 : \text{` '}, 2 : \text{`a'}, 3 : \text{`e'}, 4 : \text{`f'}, 5 : \text{`g'}, 6 : \text{`h'}, 7 : \text{`i'},$$
$$8 : \text{`k'}, 9 : \text{`m'}, 10 : \text{`n'}, 11 : \text{`p'}, 12 : \text{`s'}, 13 : \text{`u'}\}$$

*Implementation aside: We normally store $\mathcal{D}$ with the token indices instead of the text itself!*

$$\mathcal{D} = \{\, [7]\,, [1, 6, 13, 5]\,, [1, 11, 13, 5, 12]\,,$$
$$[6, 13, 5, 5, 7, 10, 5]\,, [1, 11, 13, 5, 12]\,, [1, 7, 12]\,,$$
$$[1, 4, 13, 10]\,, [7]\,, [1, 9, 2, 8, 3]\,, [1, 11, 13, 10, 12]\}$$

*For legibility of the example, we will show the text corresponding to each token*

# Byte-pair encoding: example

**Required:**

- Documents $\mathcal{D}$
- Desired vocabulary size $N$ (greater than chars in $\mathcal{D}$)

**Algorithm:**

- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation)
- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$
- Convert $\mathcal{D}$ into a list of tokens (characters)
- While $|\mathcal{V}| < N$:
  - Let $n := |\mathcal{V}| + 1$
  - Get counts of all bigrams in $\mathcal{D}$
  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)
    - Let $v_n := \mathrm{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

$$\mathcal{D} = \{\ [\text{'i'}], [\text{' '}, \text{'h'}, \text{'u'}, \text{'g'}], [\text{' '}, \text{'p'}, \text{'u'}, \text{'g'}, \text{'s'}],$$
$$[\text{'h'}, \text{'u'}, \text{'g'}, \text{'g'}, \text{'i'}, \text{'n'}, \text{'g'}], [\text{' '}, \text{'p'}, \text{'u'}, \text{'g'}, \text{'s'}],$$
$$[\text{' '}, \text{'i'}, \text{'s'}], [\text{' '}, \text{'f'}, \text{'u'}, \text{'n'}], [\text{'i'}],$$
$$[\text{' '}, \text{'m'}, \text{'a'}, \text{'k'}, \text{'e'}], [\text{' '}, \text{'p'}, \text{'u'}, \text{'n'}, \text{'s'}]\}$$

| Bigram | Count |
|---|---|
| 'u','g' | 4 |
| 'p', 'u' | 3 |
| ' ', 'p' | 3 |
| 'h', 'u' | 2 |
| ... | ... |

$$v_{14} := \mathrm{concat}(\text{'u'}, \text{'g'}) = \text{'ug'}$$

# Byte-pair encoding: example

**Required:**

- Documents $\mathcal{D}$
- Desired vocabulary size $N$ (greater than chars in $\mathcal{D}$)

**Algorithm:**

- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation)
- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$
- Convert $\mathcal{D}$ into a list of tokens (characters)
- While $|\mathcal{V}| < N$:
  - Let $n := |\mathcal{V}| + 1$
  - Get counts of all bigrams in $\mathcal{D}$
  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)
    - Let $v_n := \mathrm{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

$\mathcal{D} = \{\ [\text{‘i’}], [\text{‘ ’, ‘h’, ‘u’, ‘g’}], [\text{‘ ’, ‘p’, ‘u’, ‘g’, ‘s’}],$
$[\text{‘h’, ‘u’, ‘g’, ‘g’, ‘i’, ‘n’, ‘g’}], [\text{‘ ’, ‘p’, ‘u’, ‘g’, ‘s’}],$
$[\text{‘ ’, ‘i’, ‘s’}], [\text{‘ ’, ‘f’, ‘u’, ‘n’}], [\text{‘i’}],$
$[\text{‘ ’, ‘m’, ‘a’, ‘k’, ‘e’}], [\text{‘ ’, ‘p’, ‘u’, ‘n’, ‘s’}]\}$

$v_{14} := \mathrm{concat}(\text{‘u’, ‘g’}) = \text{‘ug’}$

$\mathcal{D} = \{\ [\text{‘i’}], [\text{‘ ’, ‘h’, ‘ug’}], [\text{‘ ’, ‘p’, ‘ug’, ‘s’}],$
$[\text{‘h’, ‘ug’, ‘g’, ‘i’, ‘n’, ‘g’}], [\text{‘ ’, ‘p’, ‘ug’, ‘s’}],$
$[\text{‘ ’, ‘i’, ‘s’}], [\text{‘ ’, ‘f’, ‘u’, ‘n’}], [\text{‘i’}],$
$[\text{‘ ’, ‘m’, ‘a’, ‘k’, ‘e’}], [\text{‘ ’, ‘p’, ‘u’, ‘n’, ‘s’}]\}$

$\mathcal{V} = \{\text{‘ ’, ‘a’, ‘e’, ‘f’, ‘g’, ‘h’, ‘i’, ‘k’, ‘m’,}$
$\text{‘n’, ‘p’, ‘s’, ‘u’, ‘ug’}\}, |\mathcal{V}| = 14$

# Byte-pair encoding: example

**Required:**
- Documents $\mathcal{D}$
- Desired vocabulary size $N$ (greater than chars in $\mathcal{D}$)

**Algorithm:**
- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation)
- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$
- Convert $\mathcal{D}$ into a list of tokens (characters)
- While $|\mathcal{V}| < N$:
  - Let $n := |\mathcal{V}| + 1$
  - Get counts of all bigrams in $\mathcal{D}$
  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)
    - Let $v_n := \mathrm{concat}(v_i, v_j)$
  - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

$$\mathcal{D} = \{\ [\text{'i'}], [\text{' '}, \text{'h'}, \text{'ug'}], [\text{' '}, \text{'p'}, \text{'ug'}, \text{'s'}],$$
$$[\text{'h'}, \text{'ug'}, \text{'g'}, \text{'i'}, \text{'n'}, \text{'g'}], [\text{' '}, \text{'p'}, \text{'ug'}, \text{'s'}],$$
$$[\text{' '}, \text{'i'}, \text{'s'}], [\text{' '}, \text{'f'}, \text{'u'}, \text{'n'}], [\text{'i'}],$$
$$[\text{' '}, \text{'m'}, \text{'a'}, \text{'k'}, \text{'e'}], [\text{' '}, \text{'p'}, \text{'u'}, \text{'n'}, \text{'s'}]\}$$

| Bigram | Count |
|--------|-------|
| ' ', 'p' | 3 |
| 'p', 'ug' | 2 |
| 'ug', 's' | 2 |
| 'u', 'n' | 2 |
| ... | ... |

$$v_{15} := \mathrm{concat}(\text{' '}, \text{'p'}) = \text{' p'}$$

# Byte-pair encoding: example

**Required:**
- Documents $\mathcal{D}$
- Desired vocabulary size $N$ (greater than chars in $\mathcal{D}$)

**Algorithm:**
- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation)
- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$
- Convert $\mathcal{D}$ into a list of tokens (characters)
- While $|\mathcal{V}| < N$:
  - Let $n := |\mathcal{V}| + 1$
  - Get counts of all bigrams in $\mathcal{D}$
  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)
    - Let $v_n := \text{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

$$\mathcal{D} = \{\,['i'], ['\,', 'h', 'ug'], ['\,', 'p', 'ug', 's'],$$
$$['h', 'ug', 'g', 'i', 'n', 'g'], ['\,', 'p', 'ug', 's'],$$
$$['\,', 'i', 's'], ['\,', 'f', 'u', 'n'], ['i'],$$
$$['\,', 'm', 'a', 'k', 'e'], ['\,', 'p', 'u', 'n', 's']\}$$

$$v_{15} := \text{concat}('\,', 'p') = '\,p'$$

$$\mathcal{D} = \{\,['i'], ['\,', 'h', 'ug'], ['\,p', 'ug', 's'],$$
$$['h', 'ug', 'g', 'i', 'n', 'g'], ['\,p', 'ug', 's'],$$
$$['\,', 'i', 's'], ['\,', 'f', 'u', 'n'], ['i'],$$
$$['\,', 'm', 'a', 'k', 'e'], ['\,p', 'u', 'n', 's']\}$$

$$\mathcal{V} = \{\,'\,', 'a', 'e', 'f', 'g', 'h', 'i', 'k', 'm',$$
$$'n', 'p', 's', 'u', 'ug', '\,p'\}, |\mathcal{V}| = 15$$

# Byte-pair encoding: example

**Required:**

- Documents $\mathcal{D}$
- Desired vocabulary size $N$ (greater than chars in $\mathcal{D}$)

**Algorithm:**

- Pre-tokenize $\mathcal{D}$ by splitting into words (split before whitespace/punctuation)
- Initialize $\mathcal{V}$ as the set of characters in $\mathcal{D}$
- Convert $\mathcal{D}$ into a list of tokens (characters)
- While $|\mathcal{V}| < N$:
  - Let $n := |\mathcal{V}| + 1$
  - Get counts of all bigrams in $\mathcal{D}$
  - For the most frequent bigram $v_i, v_j$ (breaking ties arbitrarily)
    - Let $v_n := \text{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

*Repeat until* $|\mathcal{V}| = N$...

$$\mathcal{D} = \{\ ['i']\ ,\ ['\ hug']\ ,\ ['\ pugs']\ ,$$
$$['hug',\ 'g',\ 'i',\ 'n',\ 'g']\ ,\ ['\ pugs']\ ,$$
$$['\ ',\ 'i',\ 's']\ ,\ ['\ ',\ 'f',\ 'un']\ ,\ ['i']\ ,$$
$$['\ ',\ 'm',\ 'a',\ 'k',\ 'e']\ ,\ ['\ p',\ 'un',\ 's']\}$$

$$\mathcal{V} = \{'\ ',\ 'a',\ 'e',\ 'f',\ 'g',\ 'h',\ 'i',\ 'k',\ 'm',\ 'n',\ 'p',\ 's',\ 'u',$$
$$'ug',\ '\ p',\ '\ hug',\ '\ pug',\ '\ pugs',\ 'un',\ '\ hug'\},$$
$$|\mathcal{V}| = 20$$

CHANGES FROM START

# Byte-pair encoding: example

**Questions to think about:**

- Is every token we made used in the corpus? Why or why not?

- How much memory (#tokens) have we saved for each document?

- What would happen if you kept adding vocabulary until you couldn't anymore?

$\mathcal{D} = \{$ ['i'], [' hug'], [' pugs'],

['hug', 'g', 'i', 'n', 'g'], [' pugs'],

[' ', 'i', 's'], [' ', 'f', 'un'], ['i'],

[' ', 'm', 'a', 'k', 'e'], [' p', 'un', 's']$\}$

$\mathcal{D} = \{$ [7], [20], [18],

[16, 5, 7, 10, 5], [18],

[1, 7, 12], [1, 4, 19], [7],

[1, 9, 2, 8, 3], [15, 19, 12]$\}$

*(as tokens indices)*

$\mathcal{V} = \{1 : ' ', 2 : 'a', 3 : 'e', 4 : 'f', 5 : 'g', 6 : 'h', 7 : 'i',$

$8 : 'k', 9 : 'm', 10 : 'n', 11 : 'p', 12 : 's', 13 : 'u',$

$14 : 'ug', 15 : ' p', 16 : 'hug', 17 : ' pug', 18 : ' pugs',$

$19 : 'un', 20 : ' hug'\}$

# Byte-pair encoding: tokenization/encoding

**With this vocabulary, can you represent (or, tokenize/encode):**

- "apple"?
  - No, there is no 'l' in the vocabulary

- "huge"?
  - Yes - [16, 4]

- " huge"?
  - Yes - [20, 4]

- " hugest"?
  - No, there is no 't' in the vocabulary

- "unassumingness"?
  - Yes - [19, 2, 12, 12, 13, 9, 7, 10, 5, 10, 3, 12, 12]

$$\mathcal{V} = \{1 : `\ ', 2 : `a', 3 : `e', 4 : `f', 5 : `g', 6 : `h', 7 : `i',$$
$$8 : `k', 9 : `m', 10 : `n', 11 : `p', 12 : `s', 13 : `u',$$
$$14 : `ug', 15 : `\ p', 16 : `hug', 17 : `\ pug', 18 : `\ pugs',$$
$$19 : `un', 20 : `\ hug'\}$$

# Byte-pair encoding: tokenization/encoding

$$\mathcal{V} = \{1 : \text{` '}, 2 : \text{`a'}, 3 : \text{`e'}, 4 : \text{`f'}, 5 : \text{`g'}, 6 : \text{`h'}, 7 : \text{`i'},$$

$$8 : \text{`k'}, 9 : \text{`m'}, 10 : \text{`n'}, 11 : \text{`p'}, 12 : \text{`s'}, 13 : \text{`u'},$$

$$14 : \text{`ug'}, 15 : \text{` p'}, 16 : \text{`hug'}, 17 : \text{` pug'}, 18 : \text{` pugs'},$$

$$19 : \text{`un'}, 20 : \text{` hug'}\}$$

- Sometimes, there may be more than one way to represent a word with the vocabulary…
  - E.g., " hugs" = [20, 12] = [1, 16, 12] = [1, 6, 14, 12] = [1, 6, 13, 5, 13]
    - Which is the best representation? Why?

# Byte-pair encoding: tokenization/encoding

$$\mathcal{V} = \{1 : ` \text{'}, 2 : `a\text{'}, 3 : `e\text{'}, 4 : `f\text{'}, 5 : `g\text{'}, 6 : `h\text{'}, 7 : `i\text{'},$$

$$8 : `k\text{'}, 9 : `m\text{'}, 10 : `n\text{'}, 11 : `p\text{'}, 12 : `s\text{'}, 13 : `u\text{'},$$

$$14 : `ug\text{'}, 15 : ` p\text{'}, 16 : `hug\text{'}, 17 : ` pug\text{'}, 18 : ` pugs\text{'},$$

$$19 : `un\text{'}, 20 : ` hug\text{'}\}$$

**Encoding algorithm**

Given string $\mathcal{S}$ and (ordered) vocab $\mathcal{V}$,

- Pretokenize $\mathcal{D}$ in same way as before

- Tokenize $\mathcal{D}$ into characters

- Perform merge rules in same order as in training until no more merges may be done

# Byte-pair encoding: tokenization/encoding

$$\mathcal{V} = \{1 : \text{` '}, 2 : \text{`a'}, 3 : \text{`e'}, 4 : \text{`f'}, 5 : \text{`g'}, 6 : \text{`h'}, 7 : \text{`i'},$$
$$8 : \text{`k'}, 9 : \text{`m'}, 10 : \text{`n'}, 11 : \text{`p'}, 12 : \text{`s'}, 13 : \text{`u'},$$
$$14 : \text{`ug'}, 15 : \text{` p'}, 16 : \text{`hug'}, 17 : \text{` pug'}, 18 : \text{` pugs'},$$
$$19 : \text{`un'}, 20 : \text{` hug'}\}$$

**Encoding algorithm**

Given string $\mathcal{S}$ and (ordered) vocab $\mathcal{V}$,

- Pretokenize $\mathcal{D}$ in same way as before

- Tokenize $\mathcal{D}$ into characters

- Perform merge rules in same order as in training until no more merges may be done

$$\text{Encode}(\text{`` hugs''}) = [20, 12]$$
$$\text{Encode}(\text{``misshapenness''}) = [9, 7, 12, 12, 6, 2,$$
$$11, 3, 10, 10, 3, 12, 12]$$

# Byte-pair encoding: decoding

$$\mathcal{V} = \{1 : `\;', 2 : `a', 3 : `e', 4 : `f', 5 : `g', 6 : `h', 7 : `i',$$

$$8 : `k', 9 : `m', 10 : `n', 11 : `p', 12 : `s', 13 : `u',$$

$$14 : `ug', 15 : `\;p', 16 : `hug', 17 : `\;pug', 18 : `\;pugs',$$

$$19 : `un', 20 : `\;hug'\}$$

**Decoding algorithm**

Given list of tokens $T$:

- Initialize string $s := `'$

- Keep popping off tokens from the front of $T$ and appending the corresponding string to $s$

$$\text{Encode}(`\;hugs") = [20, 12]$$

$$\text{Encode}(``misshapenness") = [9, 7, 12, 12, 6, 2,$$

$$11, 3, 10, 10, 3, 12, 12]$$

$$\text{Decode}([20, 12]) = `\;hugs"$$

$$\text{Decode}([9, 7, 12, 12, 6, 2, 11, 3, 10, 10, 3, 12, 12])$$

$$= ``misshapenness"$$

# Byte-pair encoding: properties

- Efficient to run (greedy vs. global optimization)

- Lossless compression

- Potentially some shared representations - e.g., the token "hug" could be used both in "hug" and "hugging"

# Weird properties of tokenizers

- Token != word
- Spaces are part of token
  - "run" is a different token than " run"
- Not invariant to case changes
  - "Run" is a different token than "run"

run run RunRun

[6236, 1629, 6588, 6869]

TEXT TOKEN IDS

TEXT TOKEN IDS

# Weird properties of tokenizers

- Token != word
- Spaces are part of token
  - "run" is a different token than " run"
- Not invariant to case changes
  - "Run" is a different token than "run"
- Tokenization fits statistics of your data
  - e.g., while these words are multiple tokens…
  - These words are all 1 token in GPT-3's tokenizer!
  - *Why?*
    - Reddit usernames and certain code attributes appeared enough in the corpus to surface as its own token!

tokenization
NLP
don't
victory
lose

attRot
EStreamFrame
SolidGoldMagikarp
PsyNetMessage
embedreportprint
Adinida
oreAndOnline
StreamerBot
GoldMagikarp
externalToEVA
TheNitrome
TheNitromeFan
RandomRedditorWithNo
InstoreAndOnline

TEXT    TOKEN IDS

Example from https://www.lesswrong.com/posts/aPeJE8bSo6rAFoLqg/solidgoldmagikarp-plus-prompt-generation

# Other tokenization variants

| Model/Tokenizer | Objective | Spaces part of token? | Pre-tokenization | Smallest unit |
|---|---|---|---|---|
| GPT | BPE | No | Yes | Character-level |
| GPT-2/3/4, ChatGPT, Llama(2), Falcon, … | BPE | Yes | Yes | Byte-level |
| Jurassic | BPE | Yes | No. "SentencePiece" - treat whitespace like char | Byte-level |
| Bert, DistilBert, Electra | WordPiece | No | Yes | Character-level |
| T5, ALBERT, XLNet, Marian | Unigram | Yes | No. "SentencePiece" - treat whitespace like char* | Character-level |

*For non-English languages

# Neural language models: transformers

# Neural language models: neural networks

RNNs

LSTM

Transformers

# Transformers

## Attention Is All You Need

**Ashish Vaswani*** 
Google Brain 
avaswani@google.com

**Noam Shazeer*** 
Google Brain 
noam@google.com

**Niki Parmar*** 
Google Research 
nikip@google.com

**Jakob Uszkoreit*** 
Google Research 
usz@google.com

**Llion Jones*** 
Google Research 
llion@google.com

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

**Łukasz Kaiser*** 
Google Brain 
lukaszkaiser@google.com

**Illia Polosukhin*** [‡] 
illia.polosukhin@gmail.com

(Vaswani et al., 2017)

# Transformer encoder-decoder



- Transformer encoder + Transformer decoder
- First designed and experimented on NMT

# Transformer encoder-decoder



- Transformer encoder = a stack of **encoder layers**

- Transformer decoder = a stack of **decoder layers**

> **Transformer encoder**:  BERT, RoBERTa, ELECTRA
>
> **Transformer decoder**:  GPT-3, ChatGPT, Palm
>
> **Transformer encoder-decoder**: T5, BART

- Key innovation: **multi-head, self-attention**

- Transformers don't have any recurrence structures!

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t) \in \mathbb{R}^h$$

# Transformers: roadmap



- From attention to self-attention
- From self-attention to multi-head self-attention
- Feedforward layers
- Positional encoding
- Residual connections + layer normalization
- Transformer encoder vs Transformer decoder

# Neural language models: pretraining

# Traditional learning paradigm

- **Supervised training/fine-tuning only, NO pre-training**
  - Collect (x, y) task training pairs

Data:

| sentence | label |
|----------|-------|
| a stirring , funny and finally transporting re imagining of beauty and the beast and 1930s horror films | 1 |
| apparently reassembled from the cutting room floor of any given daytime soap | 0 |
| they presume their audience won't sit still for a sociology lesson | 0 |
| this is a visually stunning rumination on love , memory , history and the war between art and commerce | 1 |
| jonathan parker 's bartleby should have been the be all end all of the modern office anomie films | 1 |

# Traditional learning paradigm

- **Supervised training/fine-tuning only, NO pre-training**
  - Collect (x, y) task training pairs
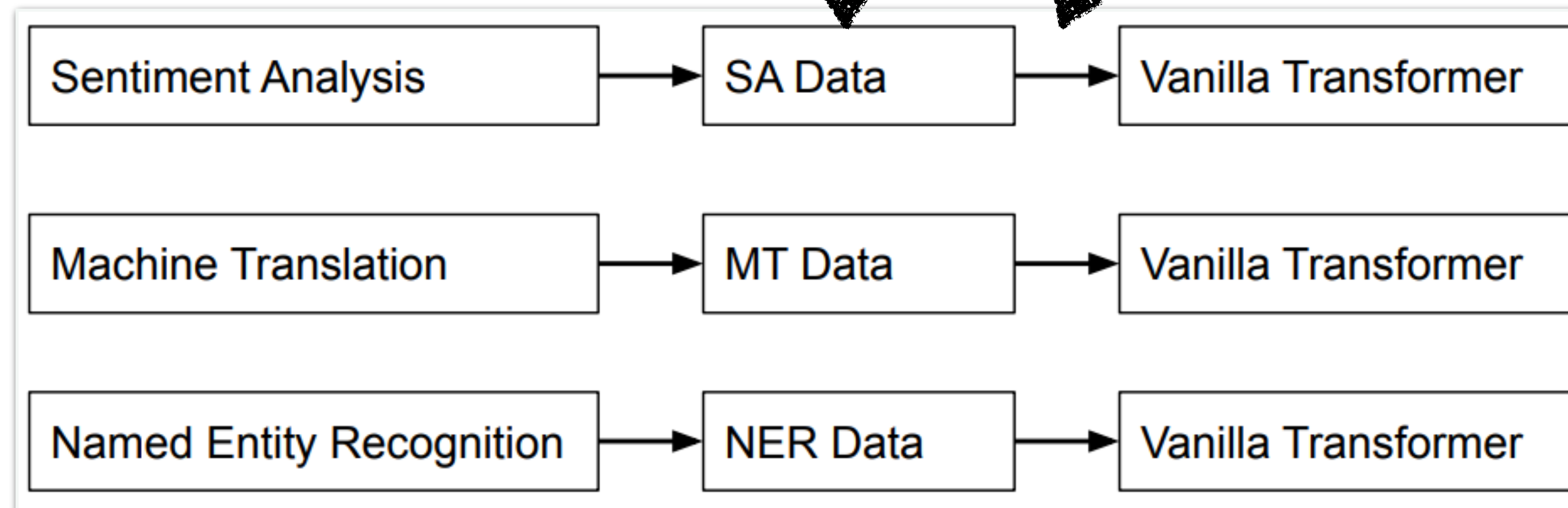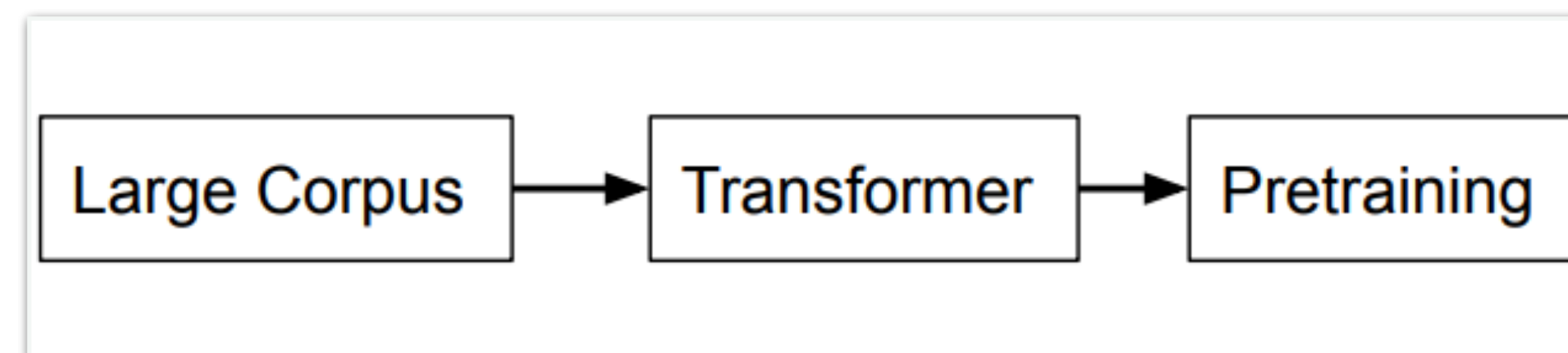  - Randomly initialize your models f(x) (e.g., vanilla Transformers)



Randomly initialized Transformers
**NO pretrained** parameters used

# Traditional learning paradigm

- **Supervised training/fine-tuning only, NO pre-training**
  - Collect (x, y) task training pairs
  - Randomly initialize your models f(x) (e.g., vanilla Transformers)
  - Train f(x) on (x, y) pairs



Train Transformers only on task labeled data

# Traditional learning paradigm

- **Supervised training/fine-tuning only, NO pre-training**
  - Collect (x, y) task training pairs
  - Randomly initialize your models f(x) (e.g., vanilla Transformers)
  - Train f(x) on (x, y) pairs



Then you get a trained Transformers **ONLY** for sentiment analysis
The model can be: NB, LR, RNNs, LSTM too

# Traditional learning paradigm

- **Supervised training/fine-tuning only, NO pre-training**
  - Train Transformer or other models separately for each task

| sentence | label |
|---|---|
| a stirring , funny and finally transporting re imagining of beauty and the beast and 1930s horror films | 1 |
| apparently reassembled from the cutting room floor of any given daytime soap | 0 |
| they presume their audience won't sit still for a sociology lesson | 0 |
| this is a visually stunning rumination on love , memory , history and the war between art and commerce | 1 |
| jonathan parker 's bartleby should have been the be all end all of the modern office anomie films | 1 |

Supervised fine-tuning only
NO pre-training

| Sentiment Analysis | → | SA Data | → | Vanilla Transformer |
| Machine Translation | → | MT Data | → | Vanilla Transformer |
| Named Entity Recognition | → | NER Data | → | Vanilla Transformer |

# Modern learning paradigm

- **Pre-training + supervised training/fine-tuning**
  - First train Transformer using a lot of general text using unsupervised learning. This is called **pretraining**.
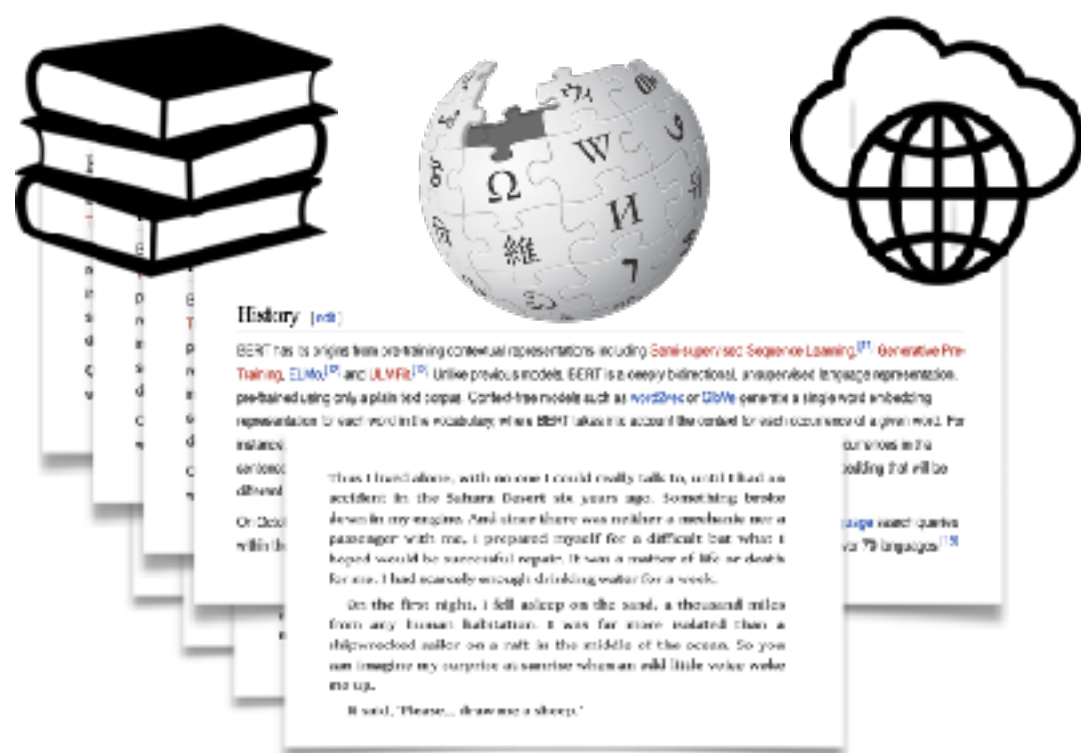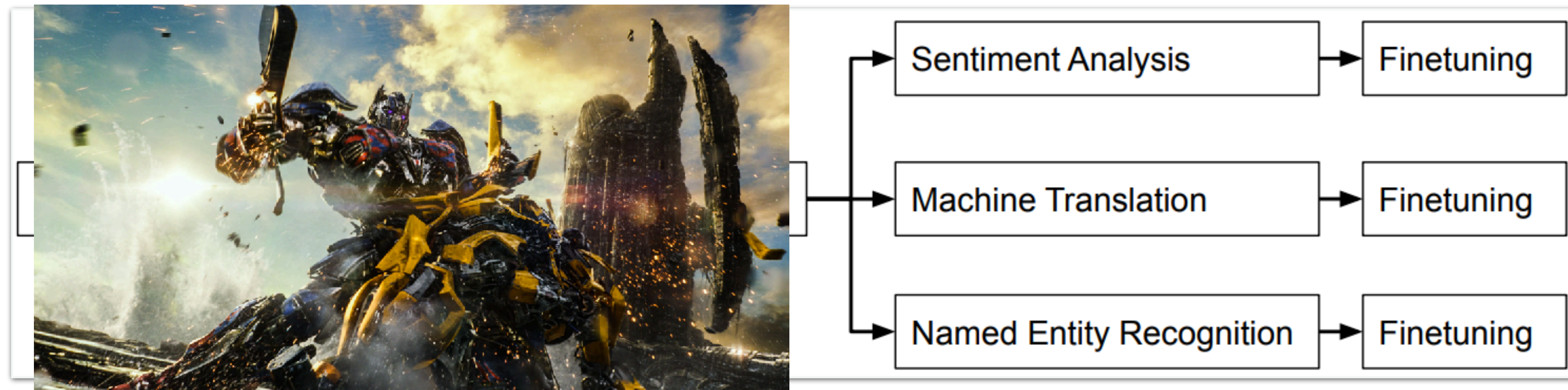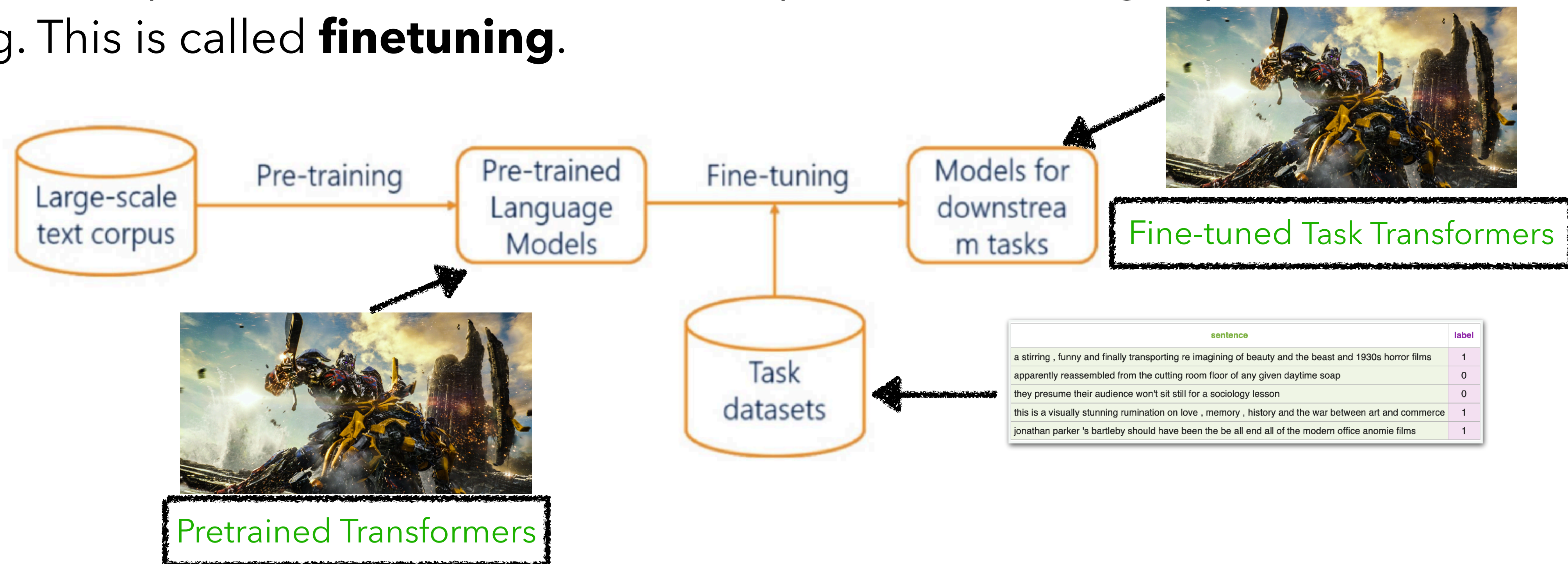
# Modern learning paradigm

- **Pre-training + supervised training/fine-tuning**
  - First train Transformer using a lot of general text using unsupervised learning. This is called **pretraining**.



1. Randomly initialized Transformers



Large Corpus → Transformer → Pretraining

# Modern learning paradigm

- **Pre-training + supervised training/fine-tuning**
  - First train Transformer using a lot of general text using unsupervised learning. This is called **pretraining**.



**general training objectives?**

2. Pre-train Transformers on large text data with **general objectives**

Large Corpus → Transformer → Pretraining

>2M GPU hours

# Modern learning paradigm

- **Pre-training + supervised training/fine-tuning**
  - First train Transformer using a lot of general text using unsupervised learning. This is called **pretraining**.
  - Then train the pretrained Transformer for a specific task using supervised learning. This is called **finetuning**.
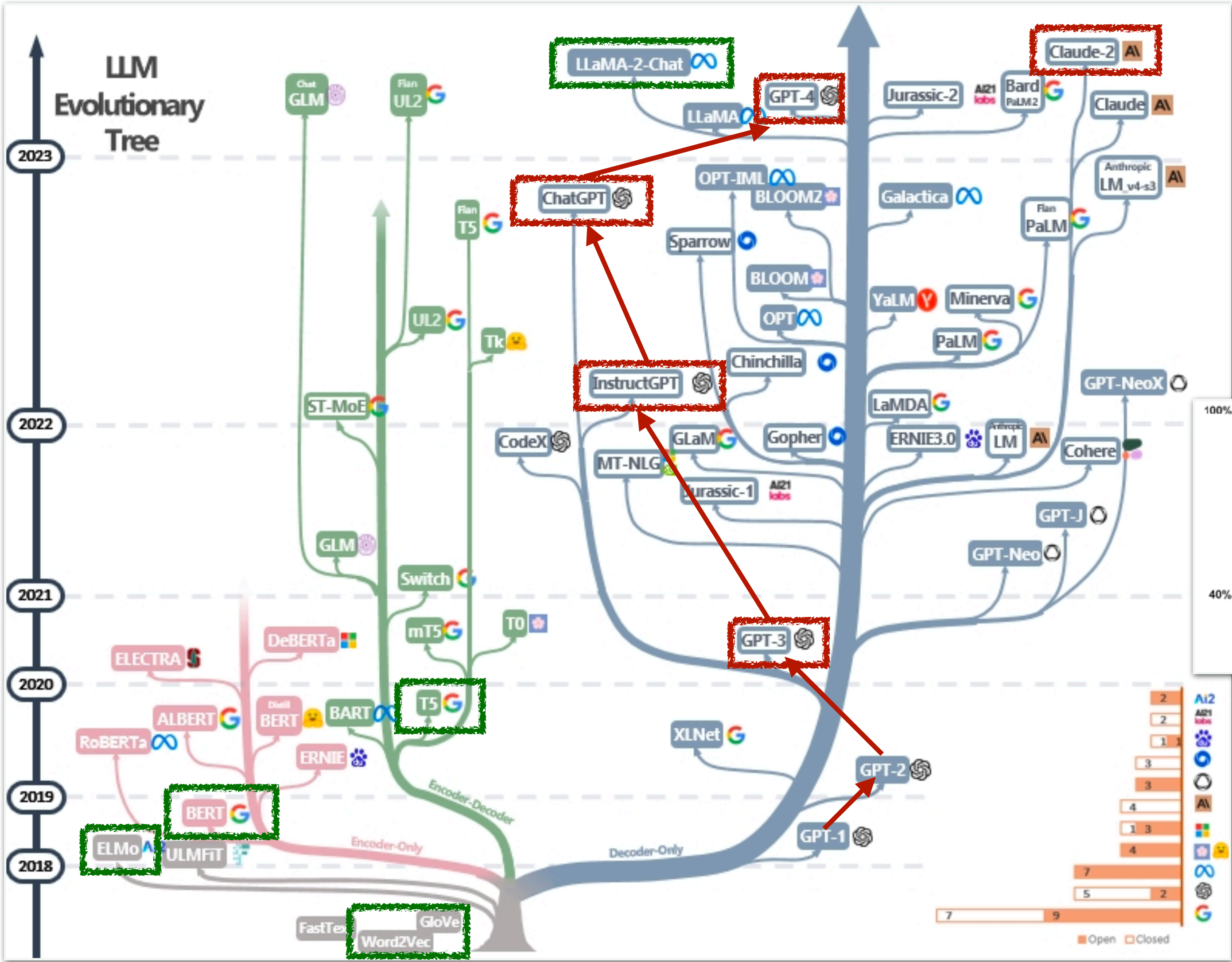
Pretrained Transformers

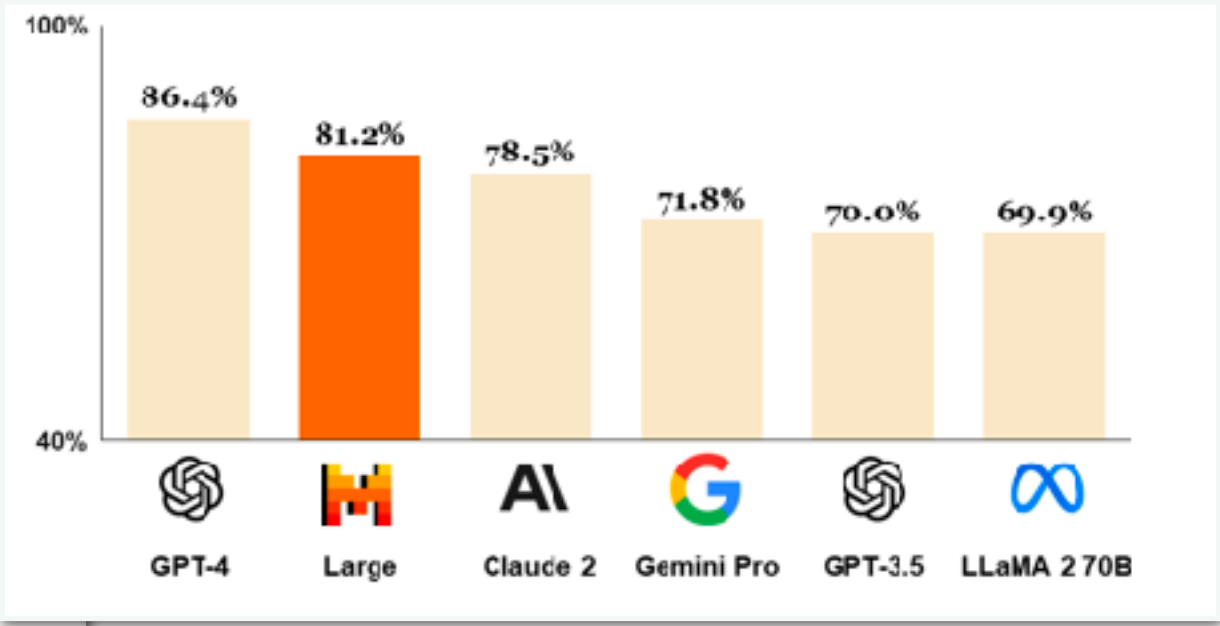# Modern learning paradigm

- **Pre-training + supervised training/fine-tuning**
  - First train Transformer using a lot of general text using unsupervised learning. This is called **pretraining**.
  - Then train the pretrained Transformer for a specific task using supervised learning. This is called **finetuning**.



Fine-tuned Task Transformers

Pretrained Transformers

# Evolution tree of pretrained LMs

# Latest learning paradigm with LLMs

- **Pre-training + prompting/in-context learning (no training this step)**

  - First train a **large (>7~175B)** Transformer using a lot of general text using unsupervised learning. This is called **large** language model **pretraining**.

# Latest learning paradigm with LLMs

- **Pre-training + prompting/in-context learning (no training this step)**

  - First train a **large (>7~175B)** Transformer using a lot of general text using unsupervised learning. This is called **large** language model **pretraining**.

  - Then **directly use** the pretrained large Transformer (**no further finetuning/training**) for any different task given only a natural language description of the task or a few task (x, y) examples. This is called **prompting/in-context learning**.



Zero-shot prompting



Few-shot prompting/in-context learning

# Example: Prompting ChatGPT for sentiment analysis

- **Pre-training + prompting/in-context learning (no training this step)**



**You**

what is the sentiment of "predictable with no fun"? just tell me: positive, negative, or neutral.
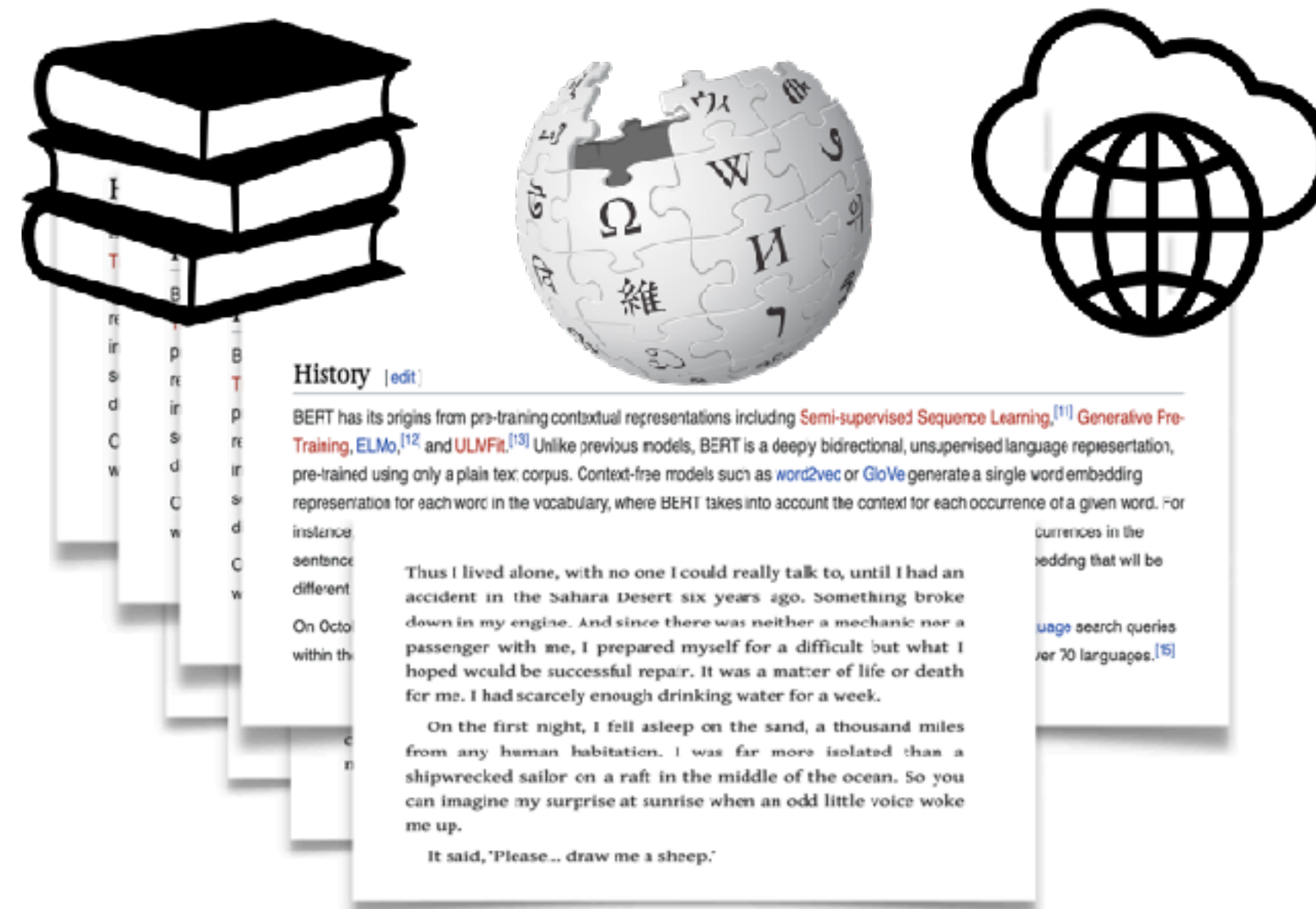
**ChatGPT**

Negative.

Already pretrained ChatGPT
No further training for sentiment analysis
Just prompting to conduct the task!

# Pretraining: training objectives?
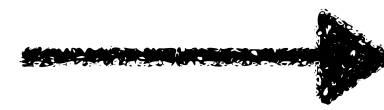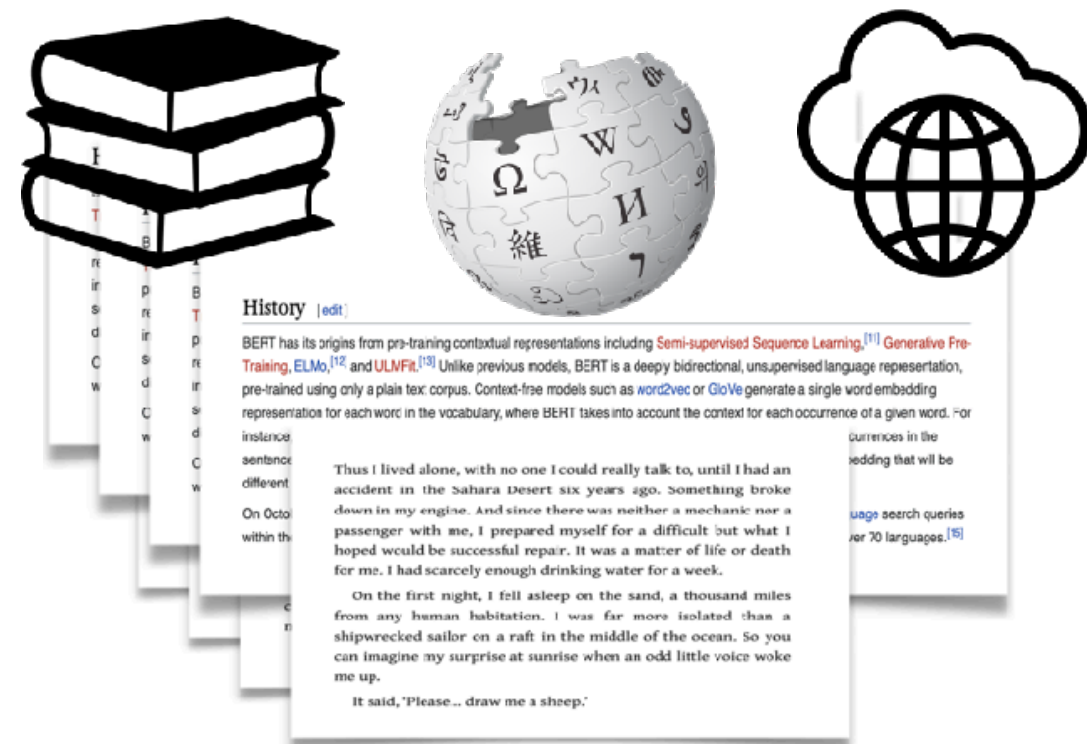
- During pretraining, we have a large text corpus (**no task labels**)
  - **Key question: what labels or objectives used to train the vanilla Transformers?**

# Pretraining: training objectives?

- During pretraining, we have a large text corpus (**no task labels**)
  - **Key question: what labels or objectives used to train the vanilla Transformers?**



Training labels/objectives?

# Pretraining: training objectives?



BERT
Devlin et al., 2018

T5
Raffel et al., 2019

OpenAI
GPT - 4

Masked token prediction    Denoising span-mask prediction    Next token prediction

# Advantages of pre-training

- **Leveraging rich underlying information** from abundant raw texts.
- **Reducing the reliance of task-specific labeled data** that is difficult or costly to obtain.
- **Initializing model parameters** for more **generalizable** NLP applications.
- **Saving training cost** by providing a reusable model checkpoints.
- **Providing robust representation** of language contexts.

# Pre-training architectures

**Encoder**

- E.g., BERT, RoBERTa, DeBERTa, …
- **Autoencoder** model
- **Masked** language modeling

**Encoder-Decoder**
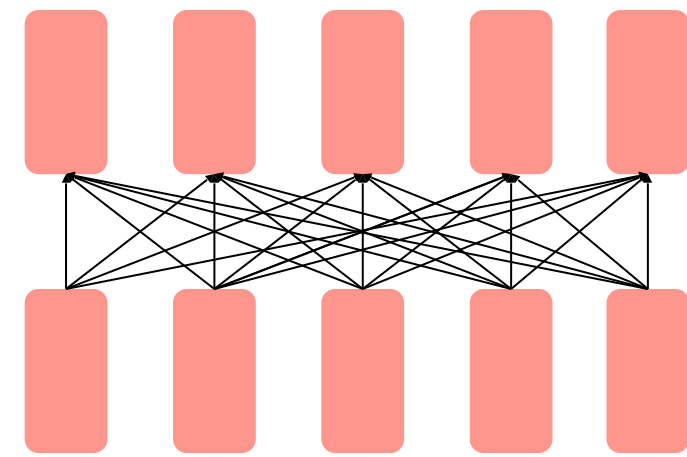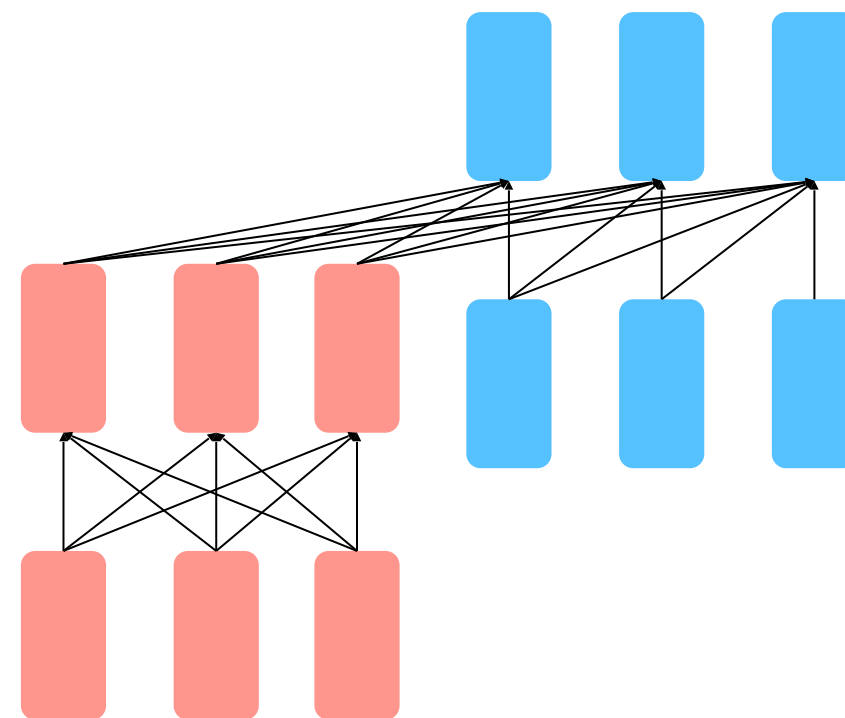
- E.g., T5, BART, …
- **seq2seq** model

**Decoder**

- E.g., GPT, GPT2, GPT3, …
- **Autoregressive** model
- **Left-to-right** language modeling
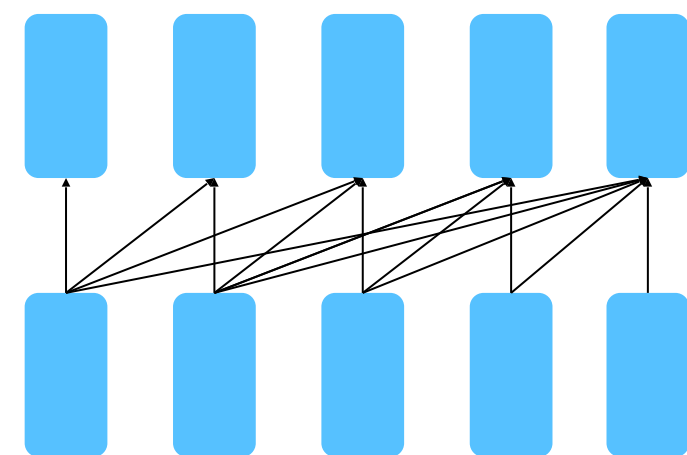
# Pre-training architectures

**Encoder**

- Bidirectional; can condition on the future context

**Encoder-Decoder**

- Map two sequences of different length together

**Decoder**

- Language modeling; can only condition on the past context

# BERT: Bidirectional Encoder Representations from Transformers



- It is a fine-tuning approach based on a deep **bidirectional Transformer encoder** instead of a Transformer decoder

- The key: learn representations based on **bidirectional contexts**

    Example #1: we went to the river bank.

    Example #2: I need to go to bank to make a deposit.

- Two new pre-training objectives:

    - **Masked language modeling (MLM)**

    - Next sentence prediction (NSP) - Later work shows that NSP hurts performance though..

(Devlin et al, 2019): BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

# Masked Language Modeling (MLM)

- Q: Why we can't do language modeling with bidirectional models?



- Solution: Mask out k% of the input words, and then predict the masked words

store                    gallon                    k = 15% in practice

↑                        ↑

the man went to [MASK] to buy a [MASK] of milk

# Masked Language Modeling (MLM)



Use the output of the masked word's position to predict the masked word

Possible classes: All English words

| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

1  2  3  4  5  6  7  8  ...  512

BERT

Randomly mask 15% of tokens

1  2  3  4  5  6  7  8  ...  512

[CLS]  Let's  stick  to  [MASK]  in  this  skit

Input

[CLS]  Let's  stick  to improvisation in  this  skit

# MLM: 80-10-10 corruption

For the 15% predicted words,

- 80% of the time, they replace it with [MASK] token

    went to the store ⟶ went to the [MASK]

- 10% of the time, they replace it with a random word in the vocabulary

    went to the store ⟶ went to the running

- 10% of the time, they keep it unchanged

    went to the store ⟶ went to the store

Why? Because [MASK] tokens are never seen during fine-tuning

(See Table 8 of the paper for an ablation study)

# Next Sentence Prediction (NSP)

- Motivation: many NLP downstream tasks require understanding the relationship between two sentences (natural language inference, paraphrase detection, QA)

- NSP is designed to reduce the gap between pre-training and fine-tuning

[CLS]: a special token always at the beginning

[SEP]: a special token used to separate two segments

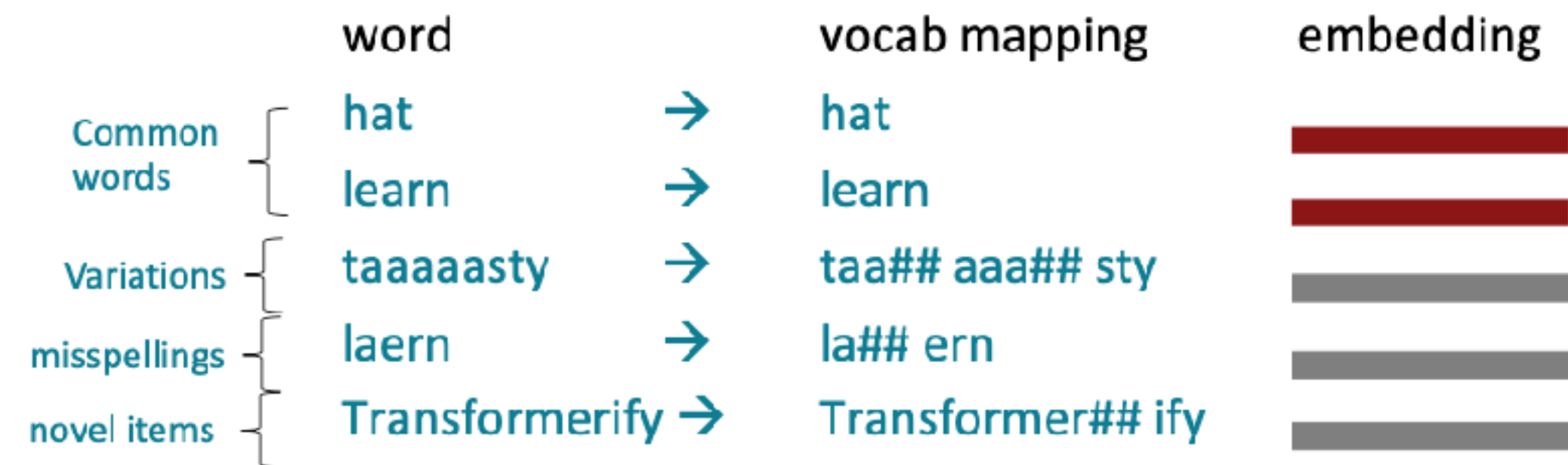Input = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

Label = IsNext

They sample two contiguous segments for 50% of the time and another random segment from the corpus for 50% of the time

Input = [CLS] the man [MASK] to the store [SEP]

penguin [MASK] are flight ##less birds [SEP]

Label = NotNext

This actually hurts model learning based on later work!

# BERT pre-training

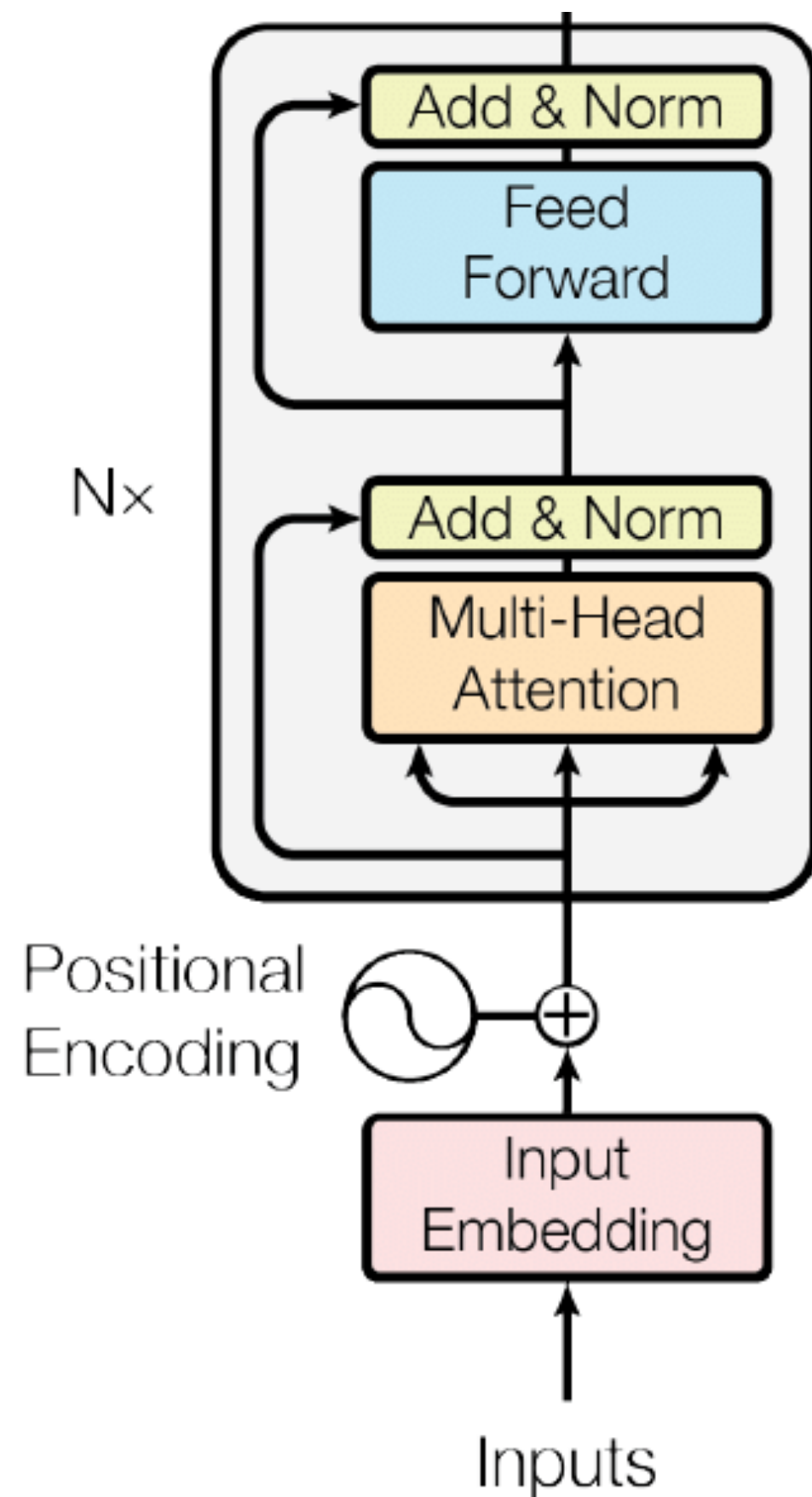- Vocabulary size: 30,000 wordpieces (common sub-word units) (Wu et al., 2016)

| | word | | vocab mapping | embedding |
|---|---|---|---|---|
| Common words | hat | → | hat | |
| | learn | → | learn | |
| Variations | taaaaasty | → | taa## aaa## sty | |
| misspellings | laern | → | la## ern | |
| novel items | Transformerify → | | Transformer## ify | |

(Image: Stanford CS224N)

- Input embeddings:

Special token added to the beginning of each input sequence

Special token to separate sentence A/B

Separate two segments

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

# BERT pre-training



- BERT-base: 12 layers, 768 hidden size, 12 attention heads, 110M parameters

- BERT-large: 24 layers, 1024 hidden size, 16 attention heads, 340M parameters

- Training corpus: Wikipedia (2.5B) + BooksCorpus (0.8B)

- Max sequence size: 512 wordpiece tokens (roughly 256 and 256 for two non-contiguous sequences)

- Trained for 1M steps, batch size 128k

# Encoder: other variations of BERT

- **ALBERT [Lan et al., 2020]**: incorporates two parameter reduction techniques that lift the major obstacles in scaling pre-trained models
- **DeBERTa [He et al., 2021]:** decoding-enhanced BERT with disentangled attention
- **SpanBERT [Joshi et al., 2019]:** masking contiguous spans of words makes a harder, more useful pre-training task
- **ELECTRA [Clark et al., 2020]:** corrupts texts by replacing some tokens with plausible alternatives sampled from a small generator network, then train a discriminative model that predicts whether each token in the corrupted input was replaced by a generator sample or not.
- **DistilBERT [Sanh et al., 2019]:** distilled version of BERT that's 40% smaller
- **TinyBERT [Jiao et al., 2019]:** distill BERT for both pre-training & fine-tuning
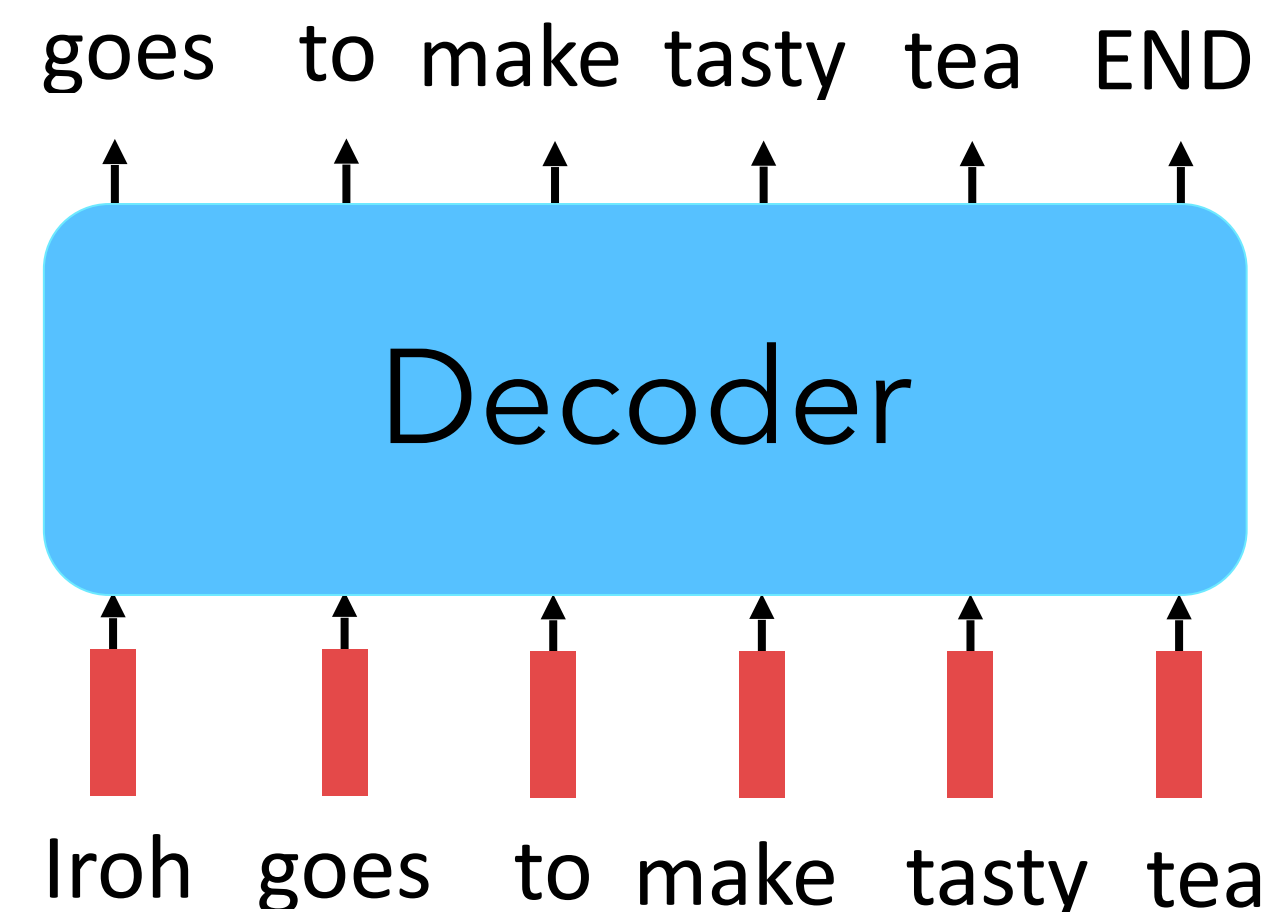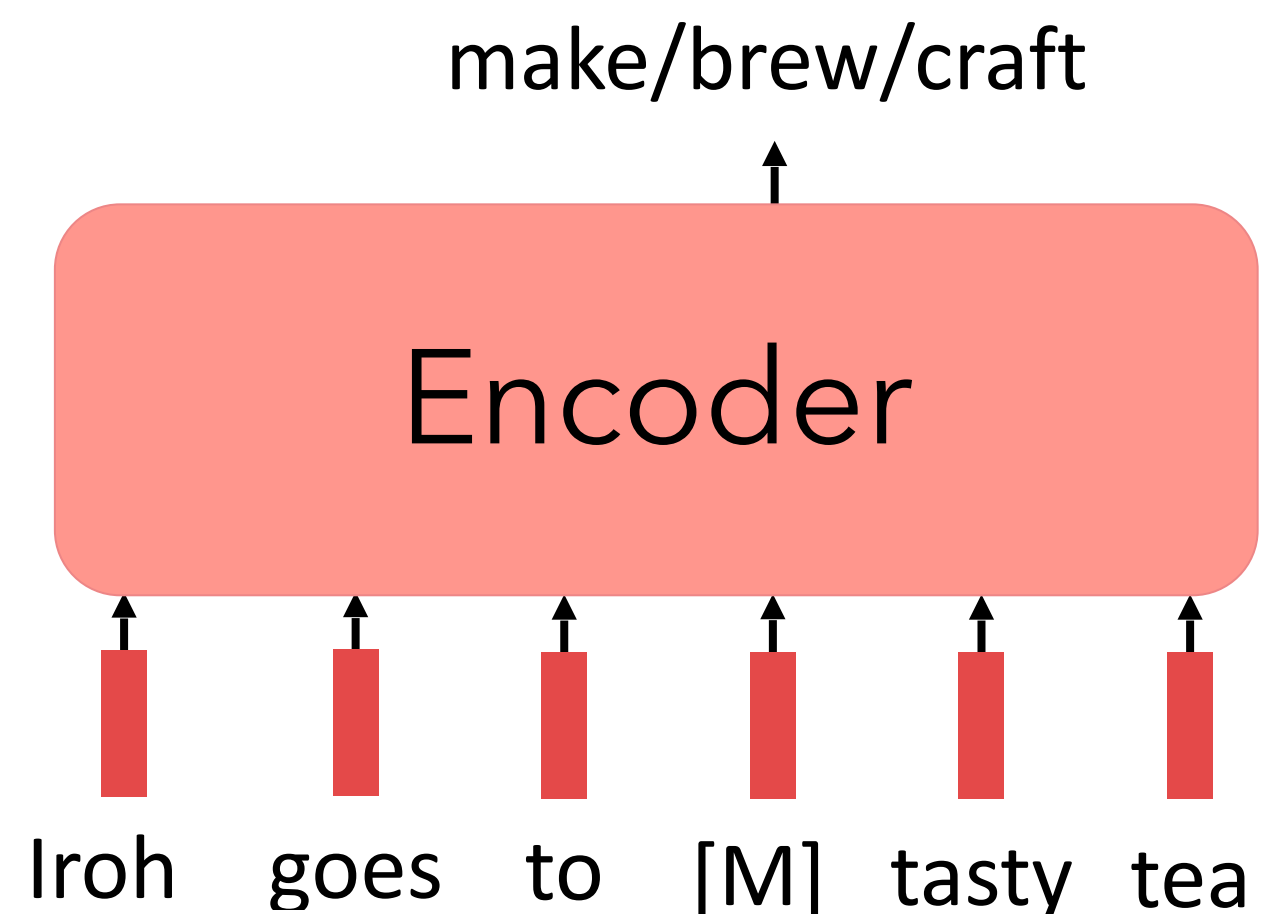- ...

# Encoder: pros & cons

- Consider both left and right context
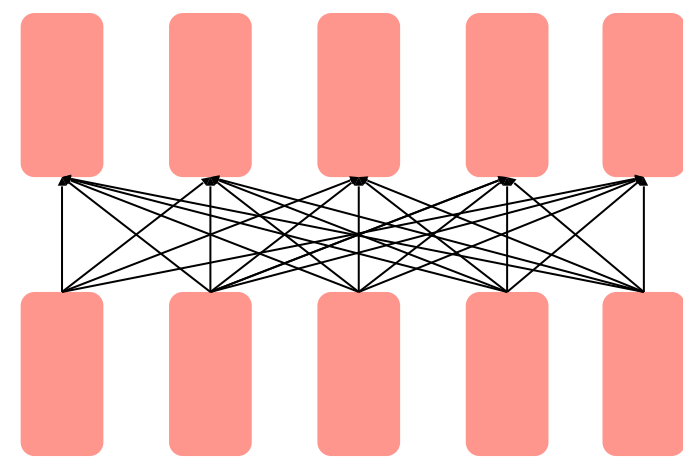- Capture intricate contextual relationships

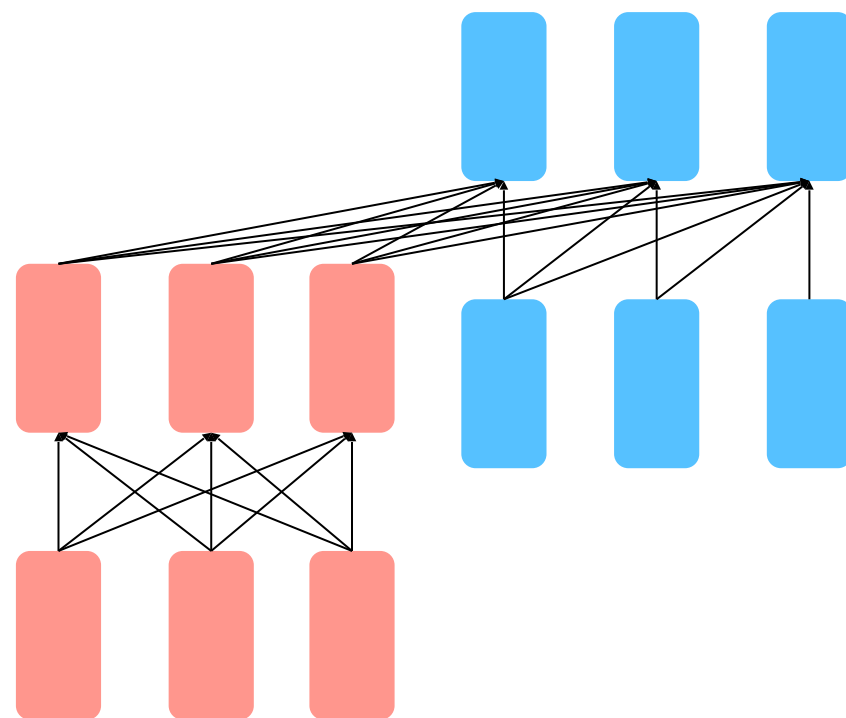- Not good at generating open-text from left-to-right, one token at a time
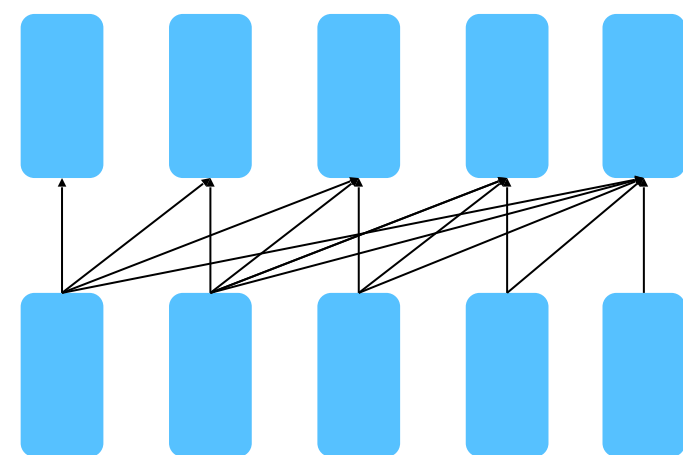
# Pre-training architectures



**Encoder**
- Bidirectional; can condition on the future context

**Encoder-Decoder**
- Map two sequences of different length together

**Decoder**
- Language modeling; can only condition on the past context

# Text-to-text models: the best of both worlds

- So bar, **encoder-only models (e.g., BERT)** enjoy the benefits of **bidirectionality** but they can't be used to generate text

- **Decoder-only models (e.g., GPT)** can do generation but they are left-to-right LMs..
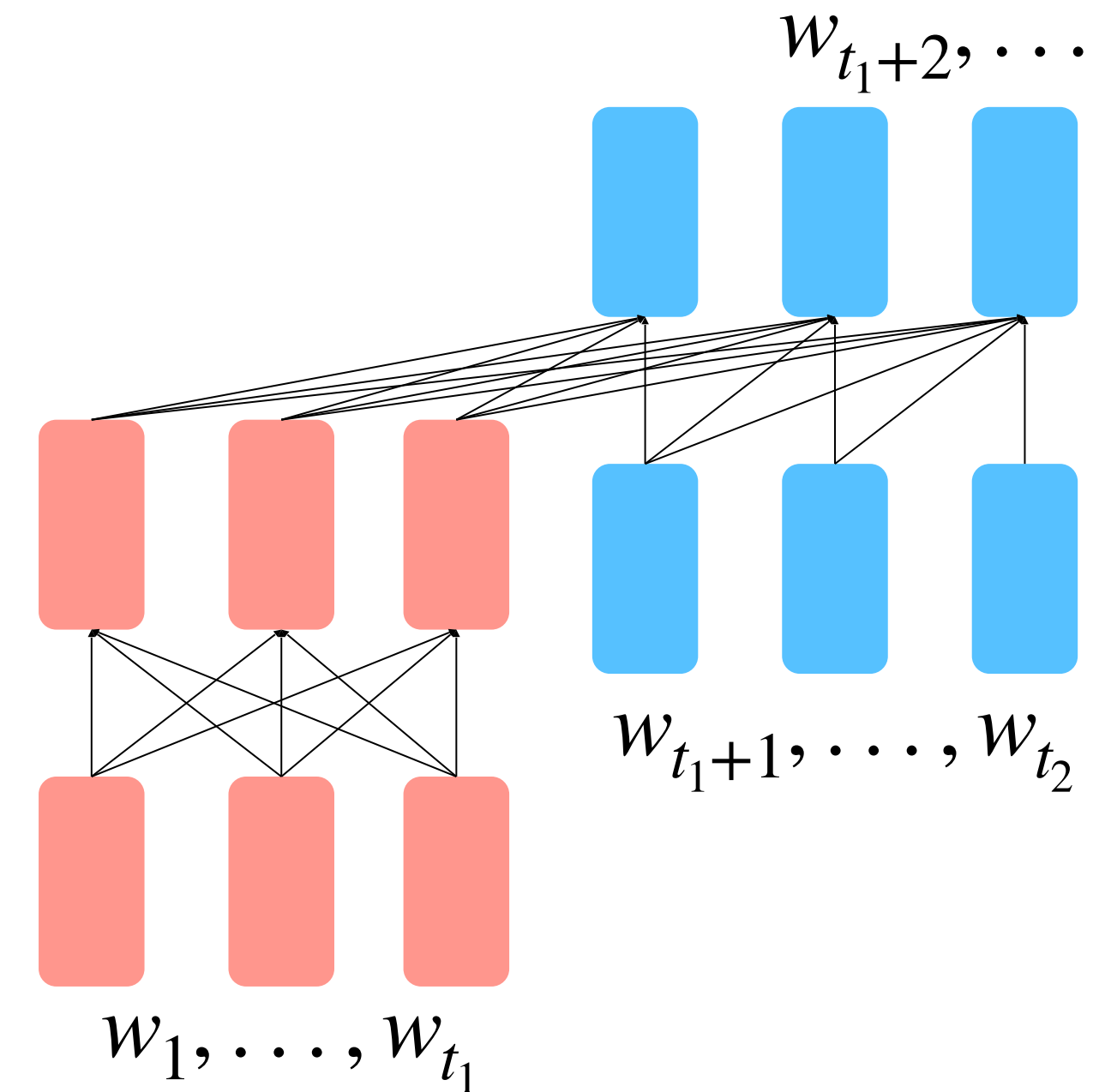
- Text-to-text models combine the best of both worlds!

T5 = **T**ext-**t**o-**T**ext **T**ransfer **T**ransformer



(Raffel et al., 2020): Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer

# Encoder-decoder: architecture

- Moving towards **open-text generation**…

- **Encoder** builds a representation of the source and gives it to the **decoder**

- **Decoder** uses the source representation to generate the target sentence

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

$$w_{t_1+2}, \ldots$$

$$w_{t_1+1}, \ldots, w_{t_2}$$

$$w_1, \ldots, w_{t_1}$$

$$h_1, \ldots, h_{t_1} = \text{Encoder}(w_1, \ldots, w_{t_1})$$

$$h_{t_1+1}, \ldots, h_{t_2} = \text{Decoder}(w_{t_1+1}, \ldots, w_{t_2}, h_1, \ldots, h_{t_1})$$

$$y_i \sim Ah_i + b, i > t$$

[Raffel et al., 2018]

# Encoder-decoder: machine translation example



P( * |Я видел котю на мате <eos>)

get probability distribution for the next token

Encoder → Decoder

Я видел котю на мате <eos>
"I" "saw" "cat" "on" "mat"

source

<bos> I saw a cat on a mat

previous history

process source and previous history

# Encoder-decoder: training objective

- **T5 [Raffel et al., 2018]**

- **Text span corruption (denoising):** Replace different-length spans from the input with unique placeholders (e.g., <extra_id_0>); decode out the masked spans.

  - Done during **text preprocessing**: training uses **language modeling** objective at the decoder side
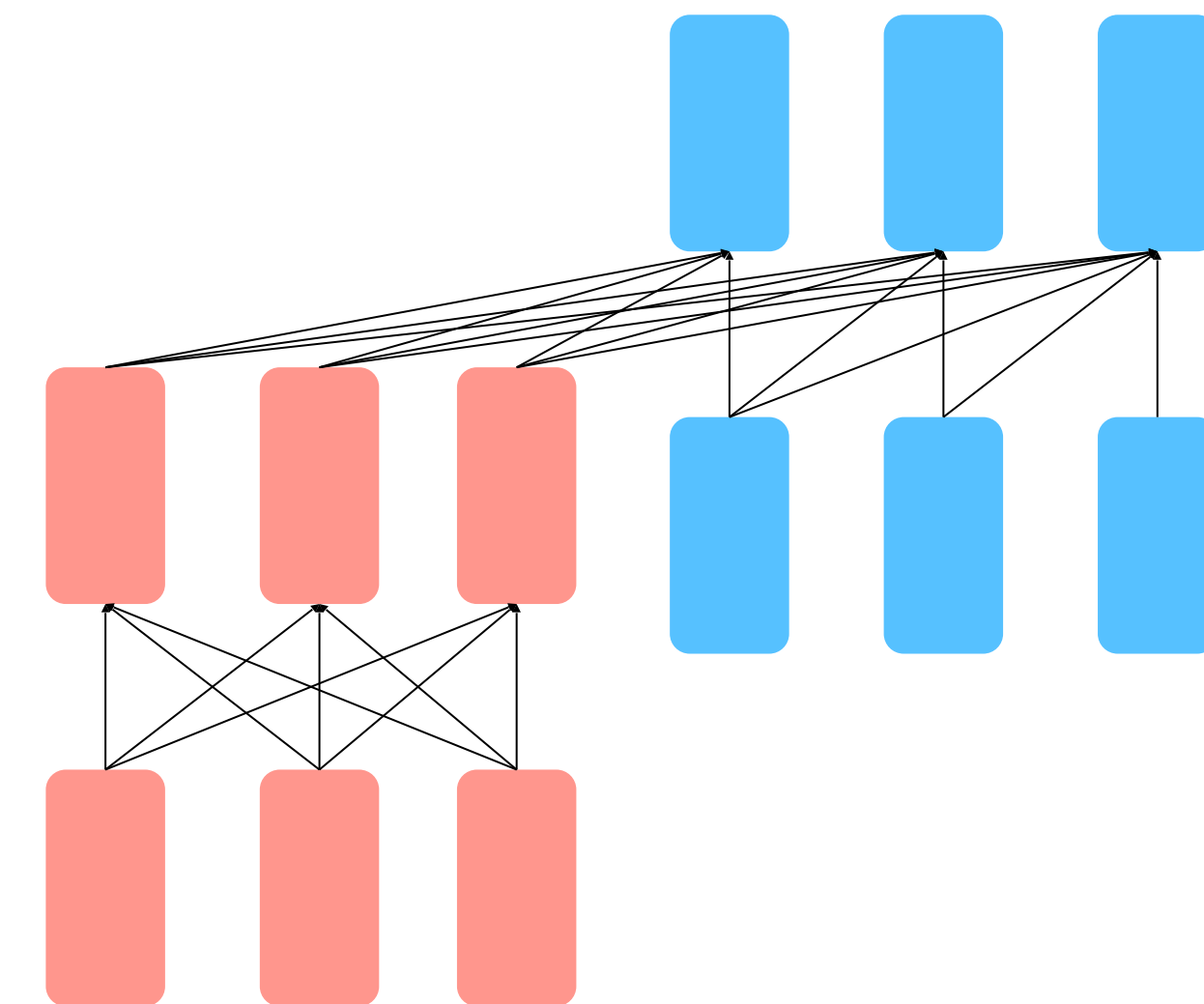


Targets

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

Inputs

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

Original text

Thank you for inviting me to your party last week.
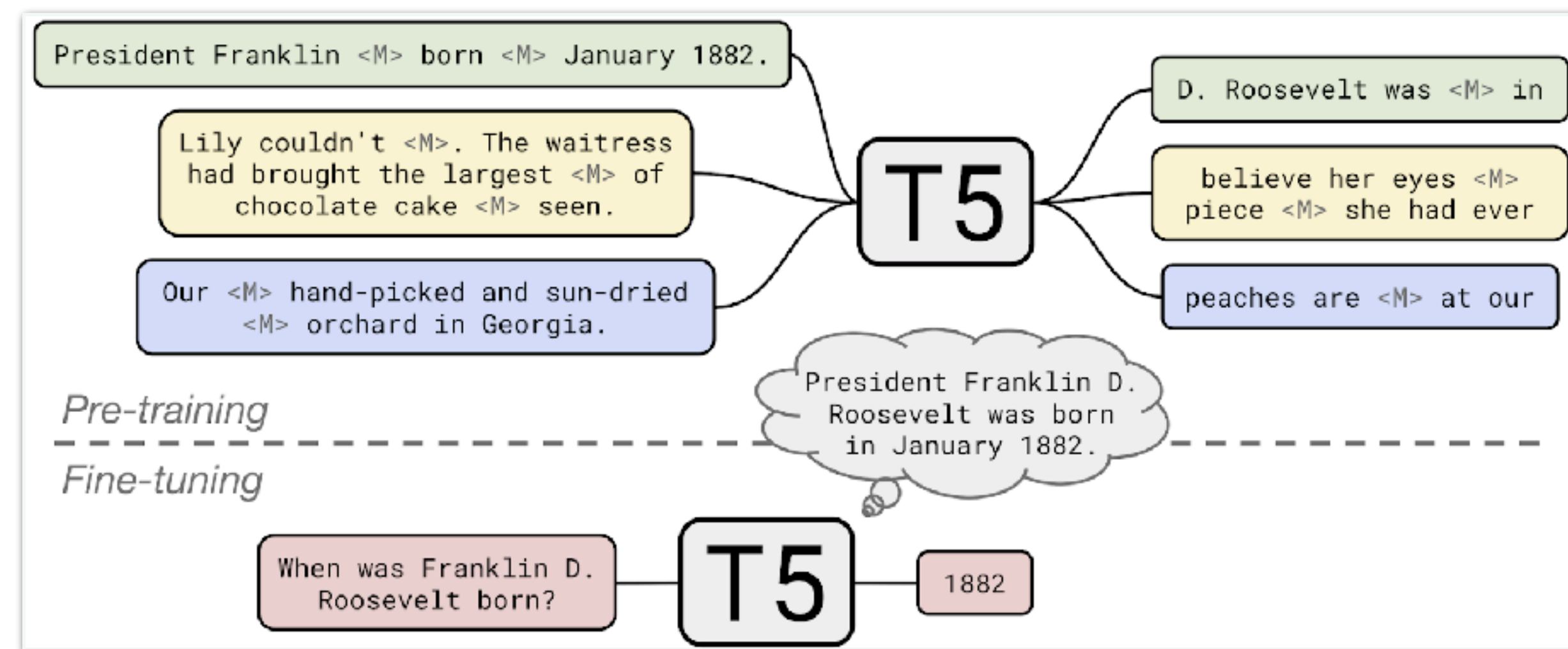
# Encoder-decoder: T5

[Raffel et al., 2018]

- **Encoder-decoders** works better than decoders
- **Span corruption (denoising)** objective works better than language modeling

| Architecture | Objective | Params | Cost | GLUE | CNNDM | SQuAD | SGLUE | EnDe | EnFr | EnRo |
|---|---|---|---|---|---|---|---|---|---|---|
| ★ Encoder-decoder | Denoising | $2P$ | $M$ | **83.28** | **19.24** | **80.88** | **71.36** | **26.98** | **39.82** | **27.65** |
| Enc-dec, shared | Denoising | $P$ | $M$ | 82.81 | 18.78 | **80.63** | **70.73** | 26.72 | 39.03 | **27.46** |
| Enc-dec, 6 layers | Denoising | $P$ | $M/2$ | 80.88 | 18.97 | 77.59 | 68.42 | 26.38 | 38.40 | 26.95 |
| Language model | Denoising | $P$ | $M$ | 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 |

# Encoder-decoder: T5

[Raffel et al., 2018]

- **Text-to-Text:** convert NLP tasks into input/ output text sequences

- **Dataset:** Colossal Clean Crawled Corpus (C4), 750G text data!

- **Various Sized Models:**
  - Base (222M)
  - Small (60M)
  - Large (770M)
  - 3B
  - 11B

- **Achieved SOTA with scaling & purity of data**



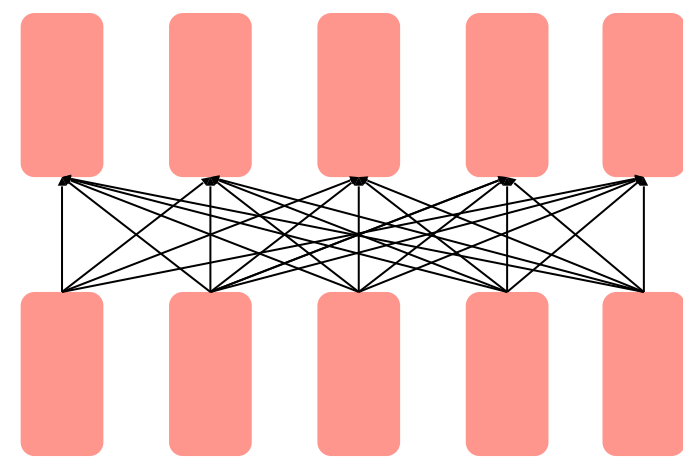[Google Blog]

# Encoder-decoder: pros & cons

- A nice middle ground between leveraging **bidirectional** contexts and **open-text** generation
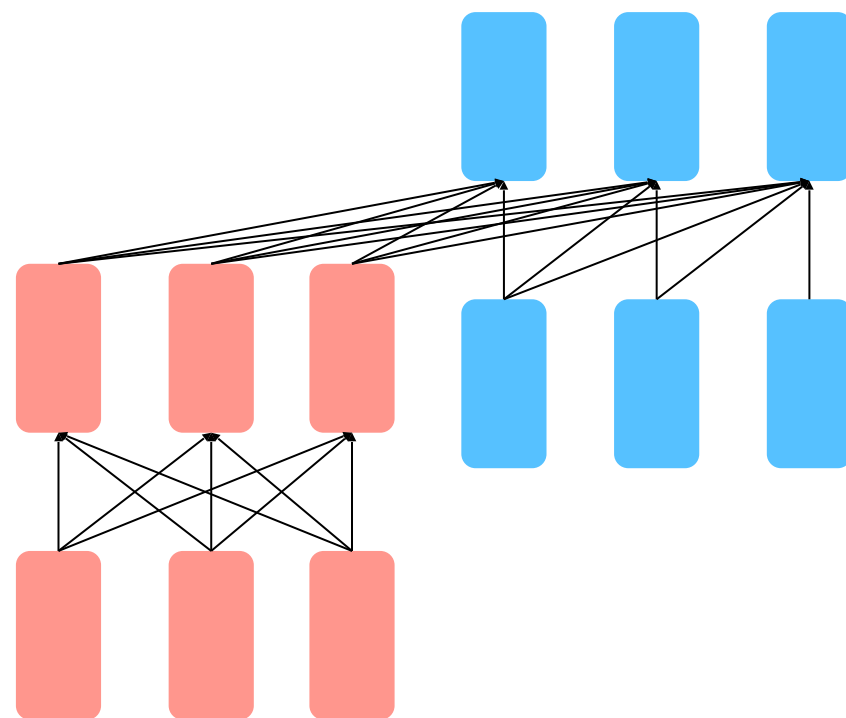- Good for **multi-task** fine-tuning

- Require more **text wrangling**
- **Harder to train**
- **Less flexible** for natural language generation
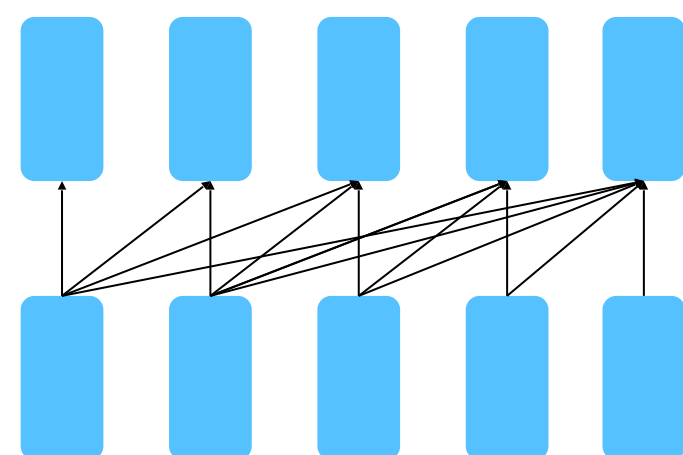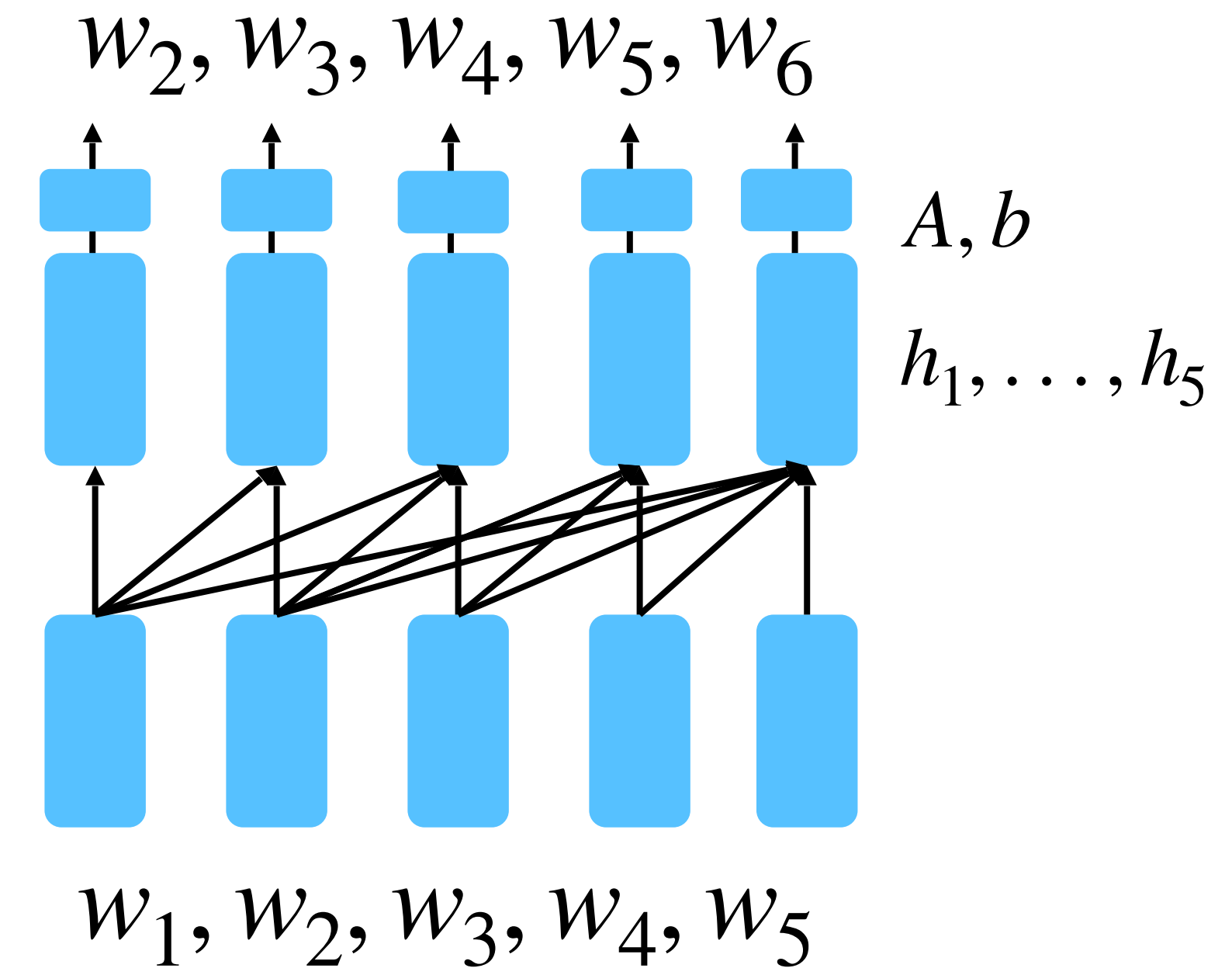
# Pre-training architectures

**Encoder**

- Bidirectional; can condition on the future context

**Encoder-Decoder**

- Map two sequences of different length together

**Decoder**

- Language modeling; can only condition on the past context

# Decoder: training objective

- Many most famous generative LLMs are **decoder-only**
  - e.g., GPT1/2/3/4, Llama1/2
- **Language modeling!** Natural to be used for **open-text generation**
- **Conditional LM:** $p(w_t | w_1, \ldots, w_{t-1}, x)$
  - Conditioned on a source context $x$ to generate from left-to-right
- Can be fine-tuned for **natural language generation (NLG)** tasks, e.g., dialogue, summarization.

$w_2, w_3, w_4, w_5, w_6$

$A, b$

$h_1, \ldots, h_5$

$w_1, w_2, w_3, w_4, w_5$

# Decoder: training objective

- Customizing the pre-trained model for downstream tasks:
  - Add a **linear layer** on top of the last hidden layer to make it a classifier!
  - During fine-tuning, trained the randomly **initialized linear layer**, along with **all parameters** in the neural net.



$A, b$

$h_1, \ldots, h_5$

Is Santa Claus real figure?

# Neural language models: scaling laws
# From GPT1 to GPT4

# Evolution tree of pretrained LMs

# From GPT1 to GPT-2 to GPT-3

- All **decoder-only Transformer-based language models**

- Model size ↑,  training corpora ↑

GPT-2



Context size = 1024



| GPT-2 SMALL | GPT-2 MEDIUM | GPT-2 LARGE | GPT-2 EXTRA LARGE |
|---|---|---|---|
| 117M Parameters | 345M Parameters | 762M Parameters | 1,542M Parameters |

.. trained on 40Gb of Internet text ..

(Radford et al., 2019): Language Models are Unsupervised Multitask Learners

# GPT-3: language models are few-shot learners

- GPT-2 → GPT-3: 1.5B → 175B (# of parameters), ~14B → 300B (# of tokens)



Total Compute Used During Training

Context size = 2048

Training computation is measured using floating-point operations or "FLOP".

One FLOP represents a single arithmetic operation involving floating-point numbers, such as addition, subtraction, multiplication, or division.

(Brown et al., 2020): Language Models are Few-Shot Learners

# Before GPT3: Modern learning paradigm

- **Pre-training + supervised training/fine-tuning**
  - First train Transformer using a lot of general text using unsupervised learning. This is called **pretraining**.
  - Then train the pretrained Transformer for a specific task using supervised learning. This is called **finetuning**.

# Paradigm shift since GPT-3

- Before GPT-3, **Pre-training + supervised training/fine-tuning** is the default way of doing learning in models like BERT/T5/GPT-2
  - SST-2 has 67k examples, SQuAD has 88k (passage, answer, question) triples

- Fine-tuning requires computing the gradient and applying a parameter update on every example (or every K examples in a mini-batch)

- However, this is very expensive for the 175B GPT-3 model



**Fine-tuning**

The model is trained via repeated gradient updates using a large corpus of example tasks.

```
1   sea otter => loutre de mer          ← example #1
```
↓
**gradient update**
↓
```
1   peppermint => menthe poivrée        ← example #2
```
↓
**gradient update**
↓
• • •
↓
```
1   plush giraffe => girafe peluche     ← example #N
```
**gradient update**
```
1   cheese => ........................  ← prompt
```

# Latest learning paradigm shift since GPT-3

- **Pre-training + prompting/in-context learning (no training this step)**

  - First train a **large (>7~175B)** Transformer using a lot of general text using unsupervised learning. This is called **large** language model **pretraining**.

  - Then **directly use** the pretrained large Transformer (**no further finetuning/training**) for any different task given only a natural language description of the task or a few task (x, y) examples. This is called **prompting/in-context learning**.
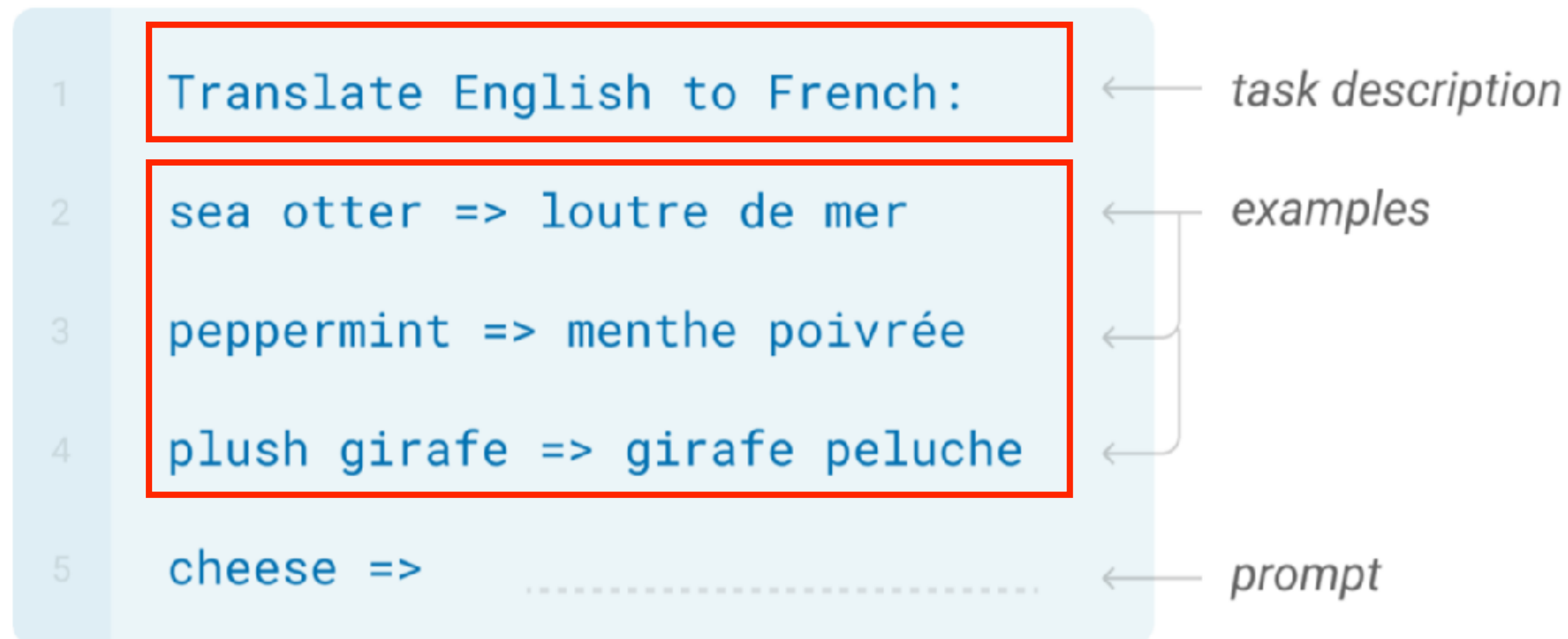
# GPT-3: few-shot in-context learning

- GPT-3 proposes an alternative: **in-context learning**

**Few-shot**

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

```
1   Translate English to French:        ←  task description

2   sea otter => loutre de mer          ←
                                            } examples
3   peppermint => menthe poivrée        ←

4   plush girafe => girafe peluche      ←

5   cheese =>  ........................  ←  prompt
```

- This is just a forward pass, **no gradient update at all**!

- You only need to feed a small number of examples (e.g., 32)

  (On the other hand, you can't feed many examples at once too as it is bounded by context size)

# GPT-3: task specifications

```
Context →   Passage:  Saint Jean de Brébeuf was a French Jesuit missionary who
            travelled to New France in 1625.  There he worked primarily with the Huron
            for the rest of his life, except for a few years in France from 1629 to
            1633.  He learned their language and culture, writing extensively about
            each to aid other missionaries.  In 1649, Brébeuf and another missionary
            were captured when an Iroquois raid took over a Huron village .  Together
            with Huron captives, the missionaries were ritually tortured and killed
            on March 16, 1649.  Brébeuf was beatified in 1925 and among eight Jesuit
            missionaries canonized as saints in the Roman Catholic Church in 1930.
            Question:  How many years did Saint Jean de Brébeuf stay in New France
            before he went back to France for a few years?
            Answer:
Target Completion →   4
```

DROP
(a reading comprehension task)

```
Context →   Please unscramble the letters into a word, and write that word:
            skicts =
Target Completion →   sticks
```

Unscrambling words

```
Context →   An outfitter provided everything needed for the safari.
            Before his first walking holiday, he went to a specialist outfitter to buy
            some boots.
            question:  Is the word 'outfitter' used in the same way in the two
            sentences above?
            answer:
Target Completion →   no
```
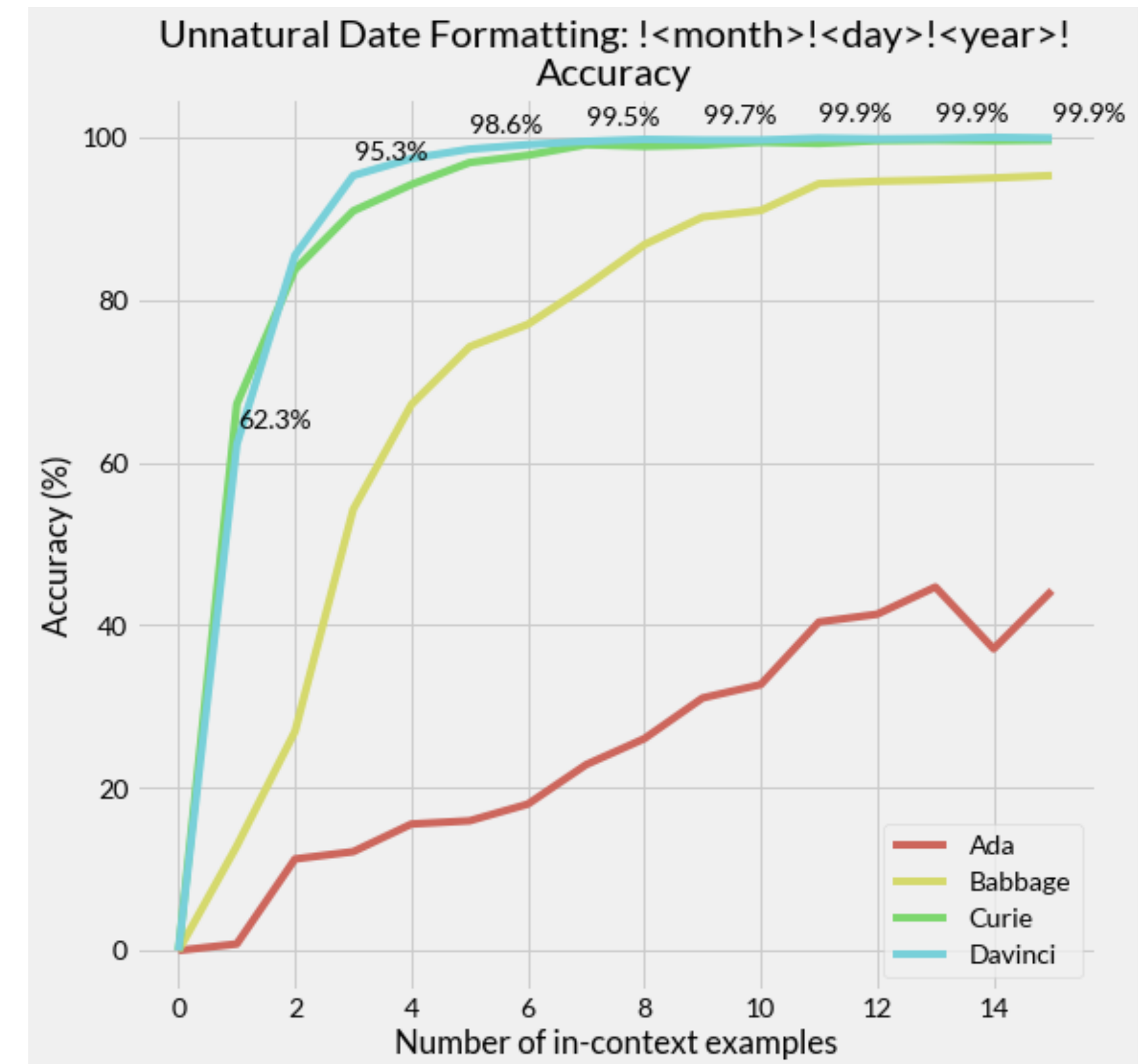
Word in context (WiC)

# GPT-3's in-context learning



```
Input: 2014-06-01
Output: !06!01!2014!
Input: 2007-12-13
Output: !12!13!2007!
Input: 2010-09-23
Output: !09!23!2010!
```
in-context examples

**Input: 2005-07-23**
test example

Output: !07!23!2005!
model completion

http://ai.stanford.edu/blog/in-context-learning/



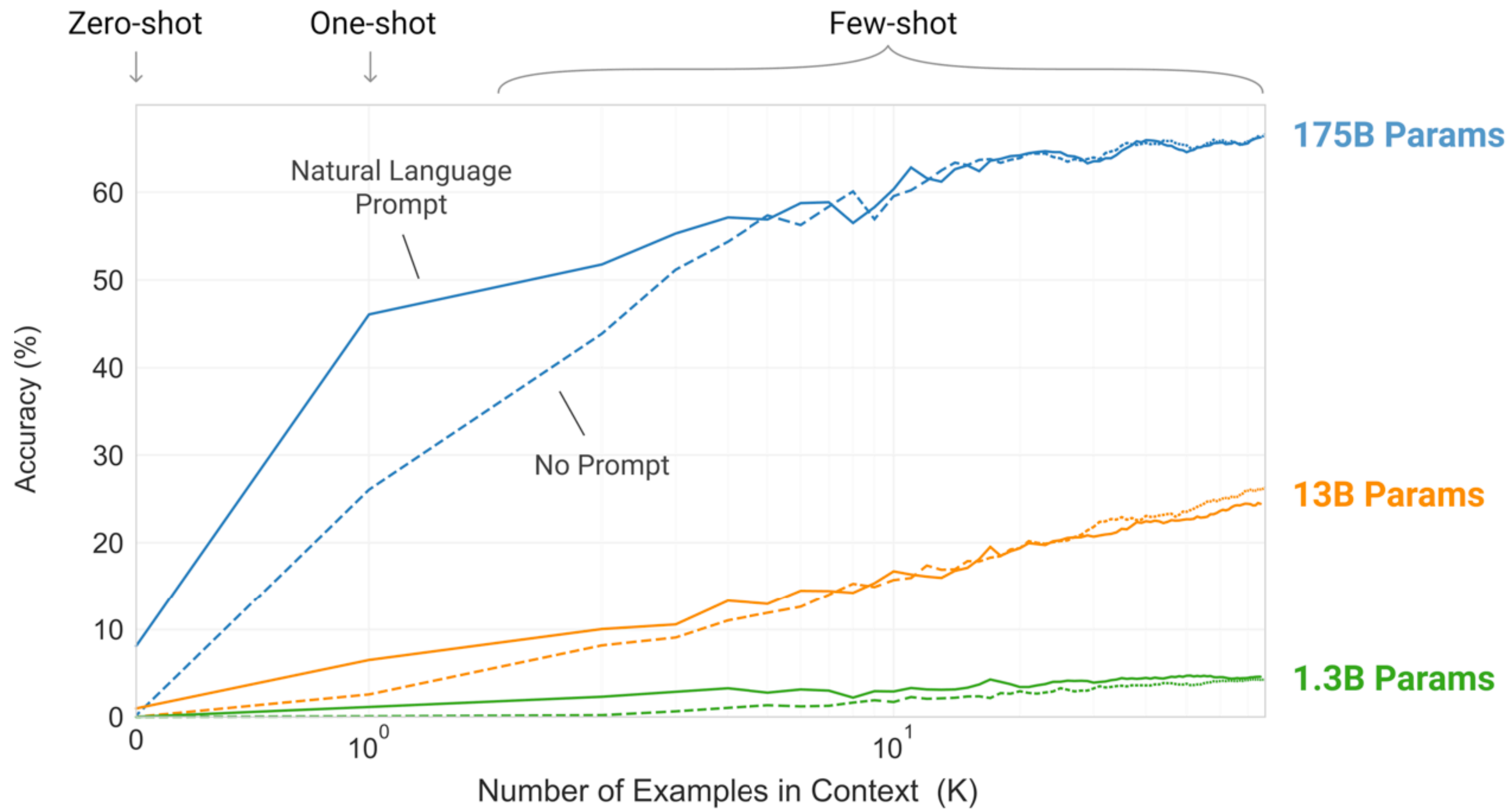Unnatural Date Formatting: !<month>!<day>!<year>!
Accuracy

- Ada
- Babbage
- Curie
- Davinci

# GPT-3's scaling laws in performance

# Chain-of-thought (CoT) prompting

**Standard Prompting**

**Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

**Chain of Thought Prompting**

**Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✅

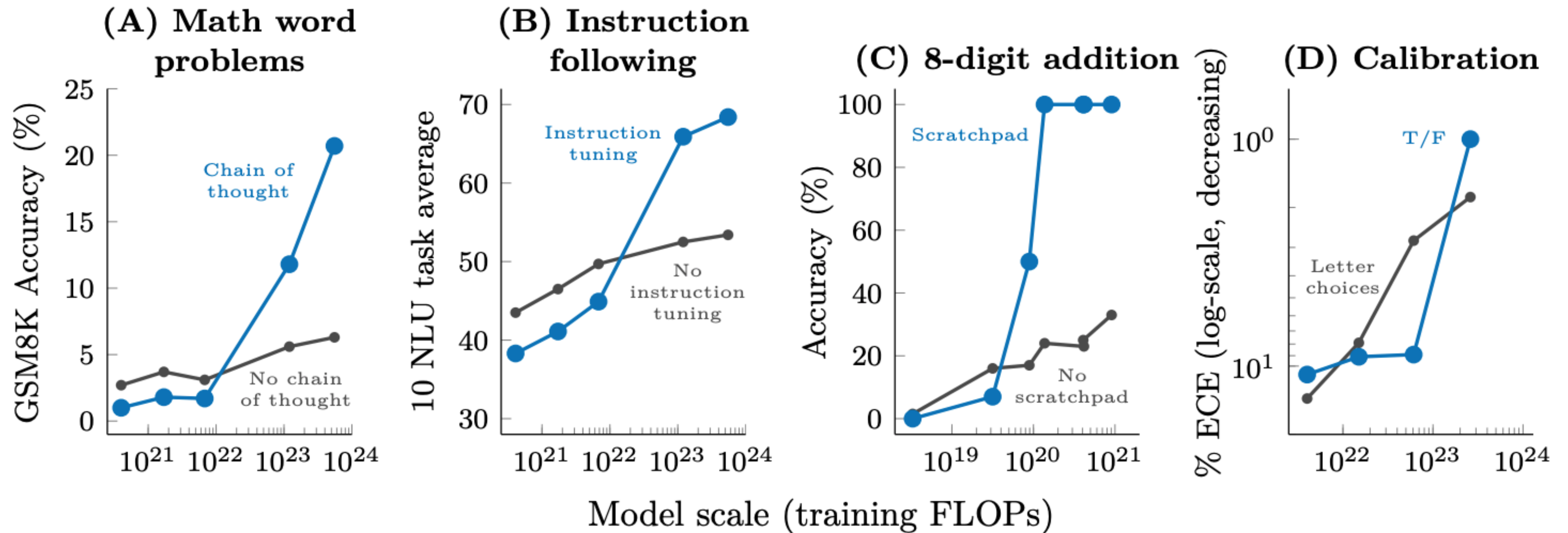(Wei et al., 2022): Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

# Why in-context learning with LLMs?

- Amazing zero/few-shot performance
  - Save a lot of annotation! 🎉
- Easy to use without training
  - Just talk to them! 👍
- One model for many NLP applications 😄
  - No need to annotate and fine-tune for different tasks


But, again, they are sensitive to prompts! Need to design a good prompt or train a good example retriever! 😂

# Emergent capabilities of LLMs?

# Emergent properties of LLMs



(Wei et al., 2022) Emergent Abilities of Large Language Models

# Emergent capabilities a mirage?

- (Schaeffer et al., 2023) take issue with the characterization of "emergent capabilities"
- Most metrics used in (Wei et al., 2022) were "hard" metrics which don't give partial credit like accuracy

---

**Are Emergent Abilities of Large Language Models a Mirage?**

Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo

Computer Science, Stanford University

Hard Accuracy:

A) 123 + 456 = 579 ✅

B) 123 + 456 = 578 ❌

C) 123 + 456 = 42 ❌

In (Wei et al., 2022), B and C are both wrong, even though B is much closer to correct than C

# Emergent capabilities a mirage?

- (Schaeffer et al., 2023) measure soft metrics (e.g., how many digits are correct, probability of the right answer) for "emergent abilities"

- Find much more predictable scaling

- Different metric choices lead to different appearances of "emergent" or not emergent

- "Emergent abilities" are a mirage(?)

Hard Accuracy:

A) 123 + 456 = 579 ✅

B) 123 + 456 = 578 ❌

C) 123 + 456 = 42 ❌

Soft Accuracy (# correct digits):

A) 123 + 456 = 579 3/3 ✅

B) 123 + 456 = 578 2/3 🤷‍♀️

C) 123 + 456 = 42 0/3 ❌