# COMP 3361 Natural Language Processing

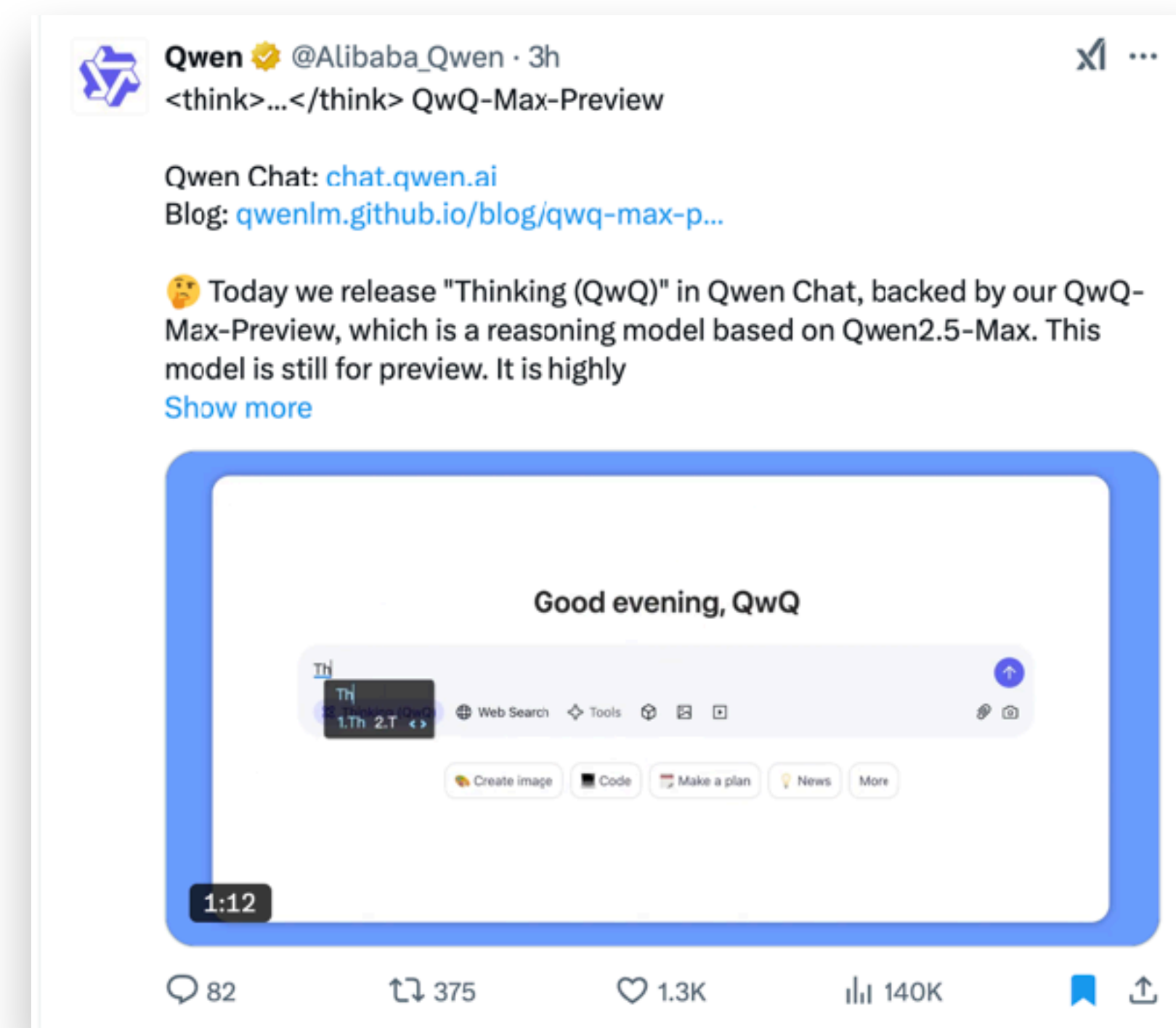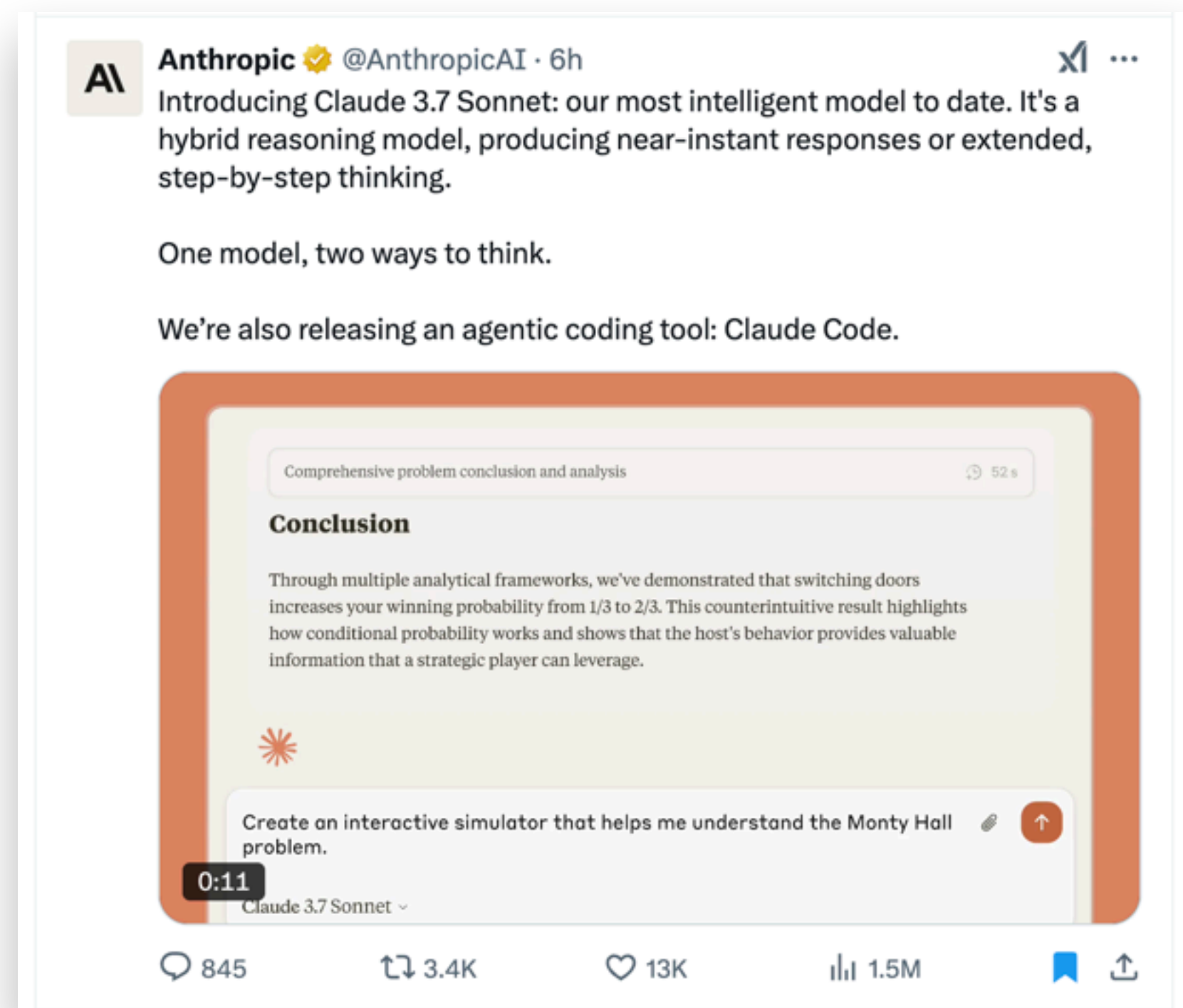## Lecture 8: Neural language models: Tokenization

Spring 2025

# Announcements

- Assignment 1 is due on Mar 4!
  - Will provide a short coding tutorial next Friday
  - Book a TA slot via the link on the course page
  - Also you can always ask questions on Slack

# Latest AI news

- Try Grok 3 for free (access w/o VPN in HK): https://x.com/i/grok
- OpenAI roadmap update for GPT-4.5 and GPT-5 (coming in May)
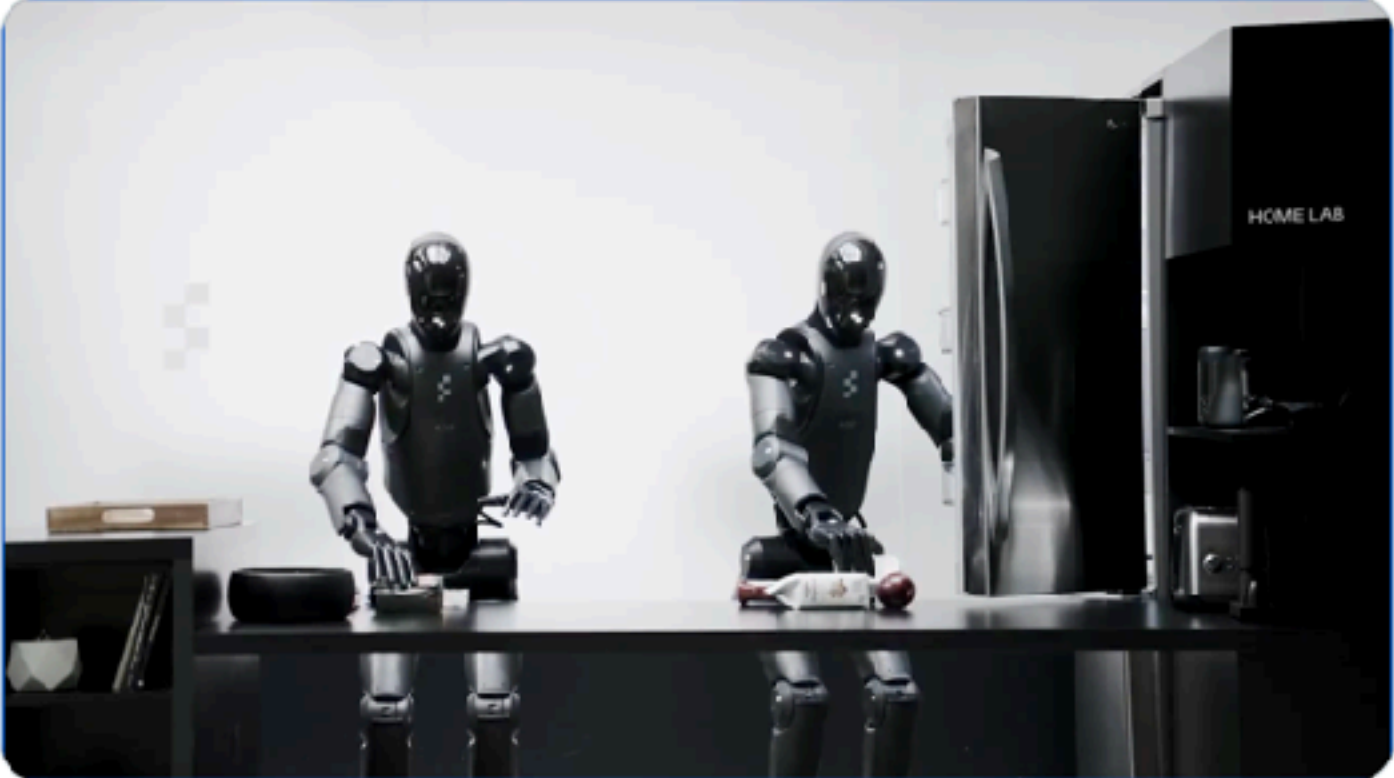- Anthropic Claude 3.7 Sonnet is out, also QwQ-Max-Preview

# Latest AI news

- <u>Helix: A Vision-Language-Action Model for Generalist Humanoid Control</u>

# Neural language models: tokenization

# Neural language models: inputs/outputs

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

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

$p(x|\text{START})$

| The | 3 |
|---|---|
| When | 2.5% |
| They | 2% |
| … | … |
| I | 1% |
| … | … |
| Banana | 0.1% |

$p(x|\text{START I})$

| think | 11% |
|---|---|
| was | 5% |
| went | 2% |
| am | 1% |
| will | 1% |
| like | 0.5% |
| … | … |

$p(x|\cdots\text{went})$

| to | 35% |
|---|---|
| back | 8% |
| into | 5% |
| through | 4% |
| out | 3% |
| on | 2% |
| … | …% |

$p(x|\cdots\text{to})$

| the | 29% |
|---|---|
| a | 9% |
| see | 5% |
| my | 3% |
| bed | 2% |
| school | 1% |
| … | … |

$p(x|\cdots\text{the})$

| bathroo | 3% |
|---|---|
| doctor | 2% |
| hospita | 2% |
| store | 1.5% |
| … | … |
| park | 0.5% |
| … | … |

$p(x|\cdots\text{park})$

| and | 14% |
|---|---|
| with | 9 |
| , | 8% |
| to | 7% |
| … | … |
| . | 6% |
| … | … |

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

| I | 21% |
|---|---|
| It | 6 |
| The | 3% |
| There | 3% |
| … | … |
| STOP | 1% |
| … | … |

**Neural Network**

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

# Tokenization to input vectors

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

## Neural Network

Mapping each tokenized id into its corresponding embeddings

Tokenization:



**Tokenization**
DistilBertTokenizer

| 101 | 1037 | 17453 | 14726 | 19379 | 12758 | 2006 | 2293 | 102 |

3) substitute tokens with their ids

| [CLS] | a | visually | stunning | rum | ##ination | on | love | [SEP] |

2) Add [CLS] and [SEP] tokens

| a | visually | stunning | rum | ##ination | on | love |

1) Break words into tokens

**Tokenize**

"a visually stunning rumination on love"

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

# ChatGPT tokenization example

# 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 $\mathcal{D}$

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

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

$N = 20$

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

$\mathcal{D} = \{$"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 $v_n := \mathrm{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

$\mathcal{D} = \{$ ['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 $v_n := \text{concat}(v_i, v_j)$

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

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

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

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

# 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 $v_n := \text{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

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

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

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

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

# 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 $v_n := \text{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

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

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

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

# 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 $v_n := \mathrm{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

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

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

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

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

# 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 $v_n := \mathrm{concat}(v_i, v_j)$
    - Change all instances in $\mathcal{D}$ of $v_i, v_j$ to $v_n$ and add $v_n$ to $\mathcal{V}$

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

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

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

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]

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

# 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

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

**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

# 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

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

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

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

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

**Decoding algorithm**

Given list of tokens $T$:

- Initialize string $s := \text{``''}$

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

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

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

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

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

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

$$= \text{``misshapenness''}$$

# 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