

Generative AI and Symbolic Knowledge Representations

LLMs, Knowledge and Reasoning

1

Damir Cavar & Billy Dickson
ESSLLI 2024

July 2024

Course Page

- https://damir.cavar.me/ESSLLI24_LLM_KG.github.io/



Resources and Contact

- [Website](#)
- [GitHub repo 1](#)
- [Damir's Python notebooks on GitHub](#)
- Get in touch: @iu.edu (during [ESSLLI 2024](#) also in person)
 - Damir Cavar: dcavar
 - Billy Dickson: dicksonb

Questions

- Who has never used Python and Jupyter Notebooks?
- Who has never experimented with Word Embeddings?
- Who has no background in Computational Linguistics or Natural Language Processing?
- Who has some understanding of how Large Language Models work?

Questions

- Who does not have experience with Knowledge Graphs?
- Who does not have experience with Description Logic, OWL, RDF, and Triple Stores?
- Who has never interacted with ChatGPT, Claude 3, Gemini, Llama 3, etc.?

Section 1

- Vector Models
- Words, Text, Semantic Space

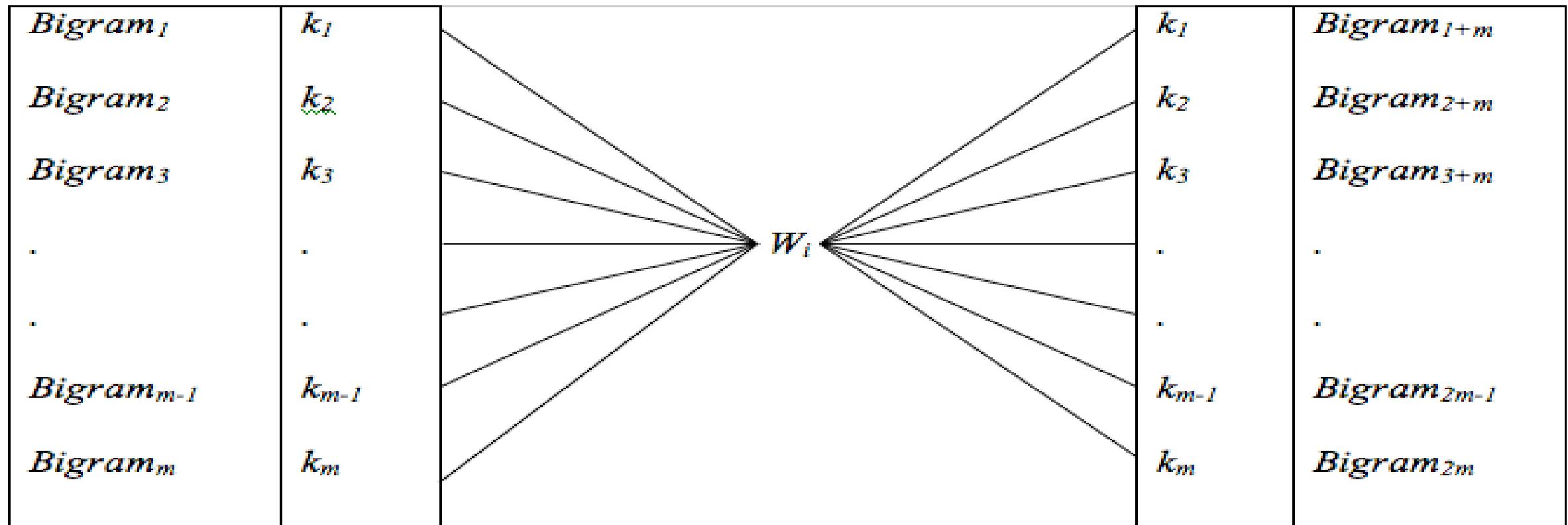
Words and Documents

- Document Vectorization (Jurafsky & Martin, [slp3](#))

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

Lexical Vector Space



Vector Space Model of Distributional Properties

- Count co-occurrence of words in a bag of words 5 words left and right:

	...	cheese	bark	purr	fly	swim	...
dog
cat
bird
mouse
...

Vector Space

Left and right context:

	left					right				
	...	the	a	big	the	a	big	...
dog	...	45	32	3	0	0	0	...
cat	...	39	27	1	0	0	0	...
walked	...	0	0	0	67	49	1	...

Vector Models

- Why vector models?
 - Architecture argument:
 - Neural Models for NLP have a fixed input-output architecture: vectors of a specific length
 - Mapping language units to mathematical representations that can be efficiently computed
- What kind of vector models?
 - Hand-crafted feature vectors
 - Distributional vectors representing words in context (also tf-idf)
 - Machine-learned distributional semantics dense vectors
 - Encoding semantic similarity and relatedness reflected in distance in semantic space

Distributional Semantics

- Distributional Properties correlate with semantic similarities

- Word similarities:

- Euclidean Distance

$$d = \sqrt{\sum_{i=1}^n (y_i - x_i)^2}$$

- Cosine Similarity

$$S_C(A, B) := \cos(\theta) = \frac{A \cdot B}{|A| |B|}$$

cosine similarity = $S_C(A, B) := \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$

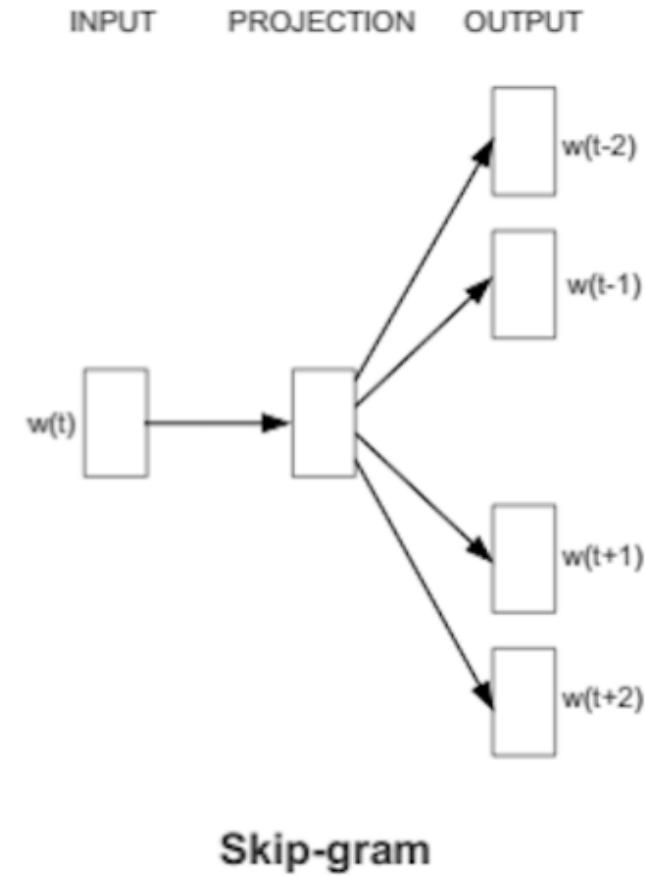
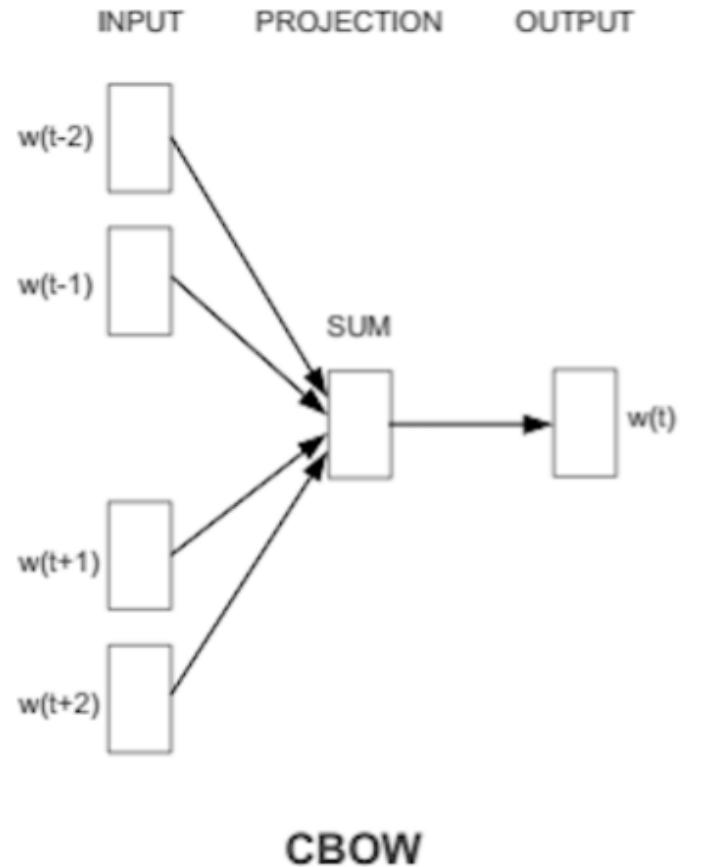
Vector Models

- Word (and text) embeddings as n-dimensional real vectors in classical Natural Language Processing (NLP)/AI applications
 - [Word2vec](#) (Mikolov et al., 2013)
 - fastText (Bojanowski et al., 2017; Joulin et al., 2018)
 - GloVe (Pennington et al. 2014)
 - Numberbatch (Speer et al., 2017)
 - BERT (Devlin et al., 2019)
 - Many derivatives

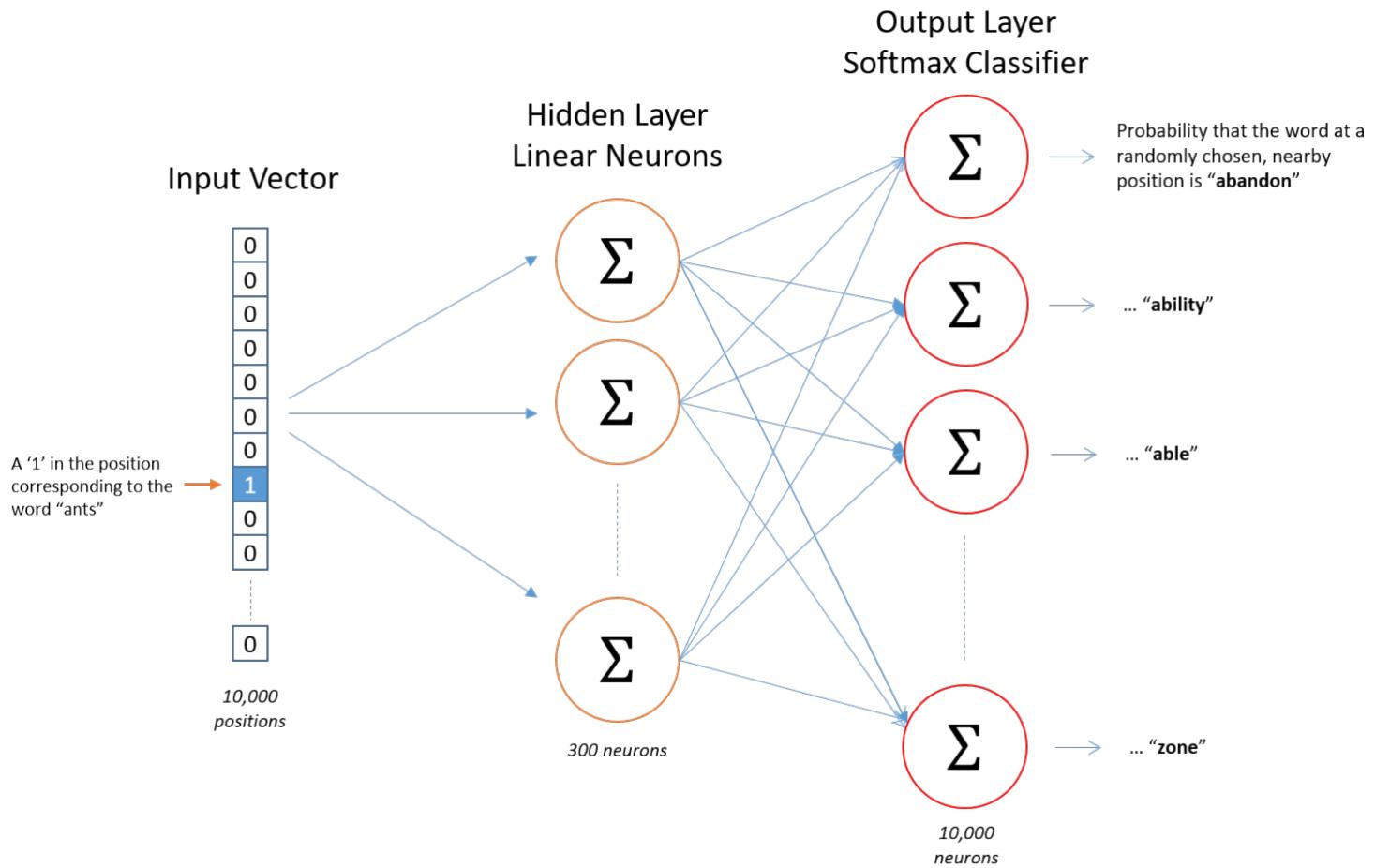
Training Vector Models

- **Continuous bag-of-words model:** This model predicts the middle word based on surrounding context words, i.e., n words before and after the current (middle) word, and the order of words in the context is not important.
- **Continuous skip-gram model:** predicts words within a certain range before and after the current word in the same sentence.
- **Predicting:** ([see tutorial](#))
 - Using a neural model, word vectors are rows in a hidden layer.

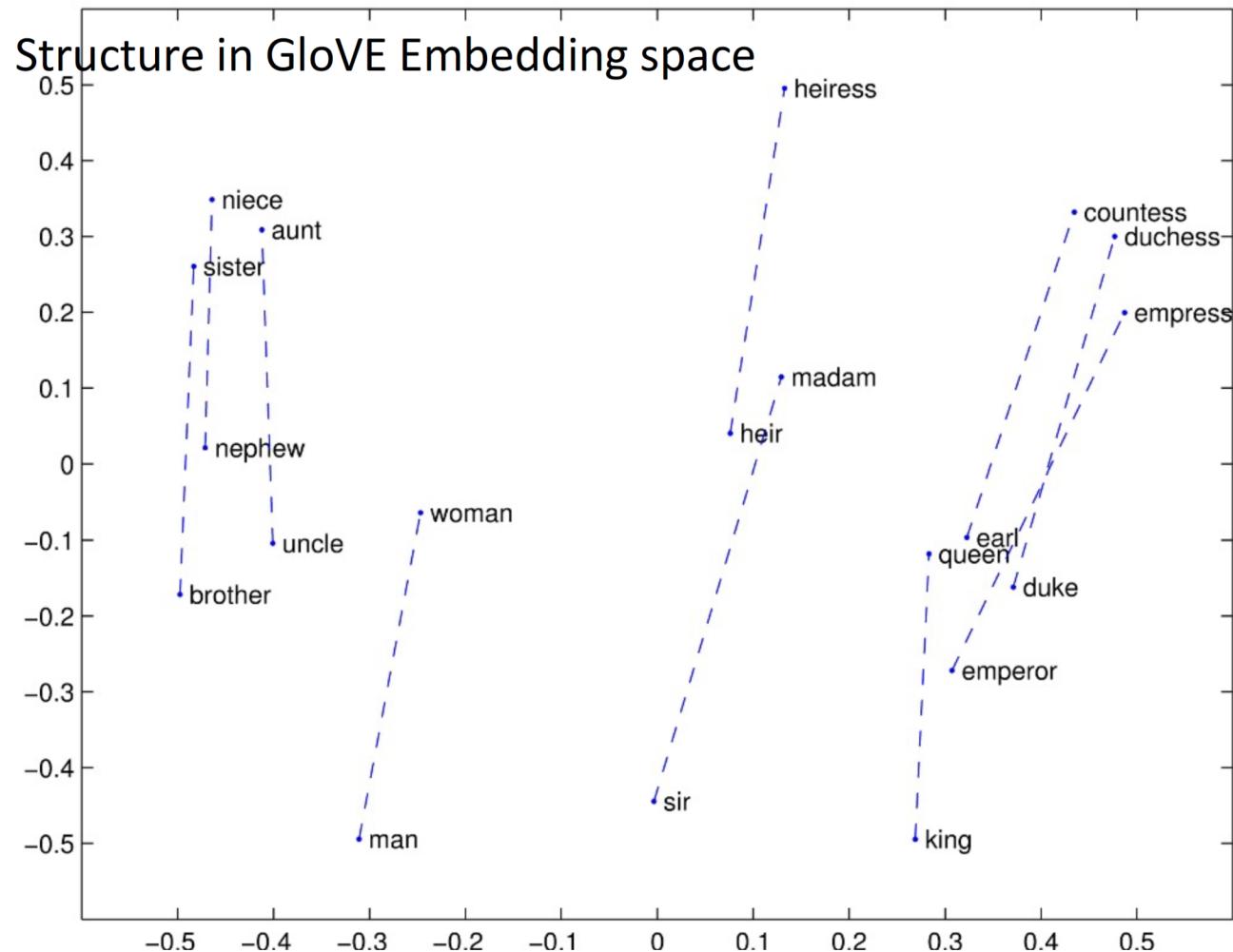
Word2vec Predictions



Word2vec as a Network



Similarity



Vector Models

- Issues with machine-learned vector models?
 - Large amount of data necessary
 - Costly: Time, Computational Effort (and carbon emission)
- Advantages:
 - Many models are freely and openly available online
 - GitHub, Hugging Face, NLP-pipeline providers

Training Vector Models

- **Cloze test:**
 1. Used in Machine Learning – Marked Word Prediction in BERT (LM)
 - *The house ___ I was born.* (a. *where* , b. *which*)
 2. Next word prediction as in Large Language Models (LLMs)
- Tokens = Words or word fragments in BERT
- Tokens ≠ Words in LLMs
 - Byte-pair encoding

Vector Models

- Word embeddings vs. Distributional Vector Models
 - [fastText](#): 300-dimensional word vectors, 2.5 mil. Words
 - [GloVe](#): 840 billion tokens, 300-dimensional word vectors, 2.1 mil. Words
 - [Numberbatch](#): 300-dimensional vectors, 516,783 words
 - [BERT](#): 768-dimensional word vectors
- Dimensions:
 - Real numbers
 - Example: $300 \times \text{real} = 300 \times 32 \text{ bit} = 9,600 \text{ bit}$
 - Real numbers could also be encoded using 64 bit

Vector Models

- Word2Vec
 - Lexical ambiguity is not captured, i.e., *date* (fruit, time expression)
 - Limited vocabulary (retraining for new words)
 - Unknown Word Problem
- fastText
 - Solves the Unknown Word Problem by using sub-word segments (not true morphemes)
- Numberbatch
 - Merges previous models and adds ConceptNet Knowledge Graph representations to improve semantic properties
 - Graph Embeddings (Concepts instead of strings)

Vector Models

- BERT
 - Adds sub-word tokens: WordPiece algorithm
 - *locking* could be tokenized as *lock* and <##*ing*
 - <##*ing* indicates that *ing* should be treated like a word continuation, i.e., a dependent morpheme in the linguistic sense
 - Processing of sequences of words
 - reduces the unknown word problem and increases the coverage at the lexical level
 - WordPiece in Schuster and Nakajima (2012) as a method to cope with out-of-vocabulary words and improve the general performance of machine learning and NLP tasks

Vector Models

- Orthographic / Linguistic tokenization
 - Too large a number of tokens (= number of embeddings)
- WordPiece and similar strategies:
 - Still a large number of tokens on multi-lingual content, programming code, data structures, etc.
- Simpler encoding and tokenization:
 - Byte-pair encoding
 - non-linguistic or orthographic
 - language independent
 - Predefined vocabulary size

Tokenization

- Tokenization methods:
 - Word tokenization
 - WordPiece
 - Byte-Pair Encoding (BPE)
 - Unigram
 - SentencePiece

Word Tokenization

- Training data might be missing words: the Unknown Word Problem
 - The training missing: football
 - But not missing: "foot" or "ball"
- Inflected forms in specific tense, case, number, gender agreement might have been never seen
- A lemmatizer might help, but: the concrete lemmatized token might be out of vocabulary

Word Tokenization

- Other problems:
 - slang
 - abbreviations
- Or language specific orthography:
 - no spaces between words (e.g., Japanese, Chinese, Arabic)

Tokenization

- Alternative: Character-based tokenization
- Compute embeddings for characters
- Compose words as embeddings using character embeddings
- Robust with input processing
 - Typos
 - Variation in orthography
 - Unknown words

Character Tokenization

- Problem:
 - each character is a token
 - this results in a lot more tokens
 - this results in a lot more input computation
- Example:
 - 5 word sentence might have 30 character tokens and 5 word tokens
- Word generation using characters can result in wrong spellings

Character Tokenization

- Distributional semantics might not work with such models
- Compare:
 - listen
 - silent
- We would need character embeddings and positional encoding to catch the difference between the two words.

Sub-word Tokenization

- Character-based models:
 - risk losing the semantic features of the word
- Word-based tokenization:
 - very large vocabulary needed to encompass all the possible variations of every word
- Goal: develop an algorithm that could
 - retain the semantic features of the token, that is information per token
 - tokenize without demanding a very large vocabulary with a finite set of words

Sub-word Tokenization

- Breaking down the words based on a set of prefixes and suffixes:
 - identify sub-words like: "##s", "##ing", "##ify", "un##"...
 - the position of the double hash denotes prefix and suffixes
- Tokenize: "unhappily"
 - "un##", "happ", "##ily"
- Result: few subwords and combinatorics to create words
- Compression: memory requirement and little effort required to create a large vocabulary

Sub-word Tokenization

- Problems:
 - some created sub-words might be very rare and occupy space
 - language specific rules and tokens:
 - for every language we need to define a different set of rules to create sub-words - every language has a different morphology
- Strategy:
 - most modern tokenizers have a training phase:
 - identification of recurring text in the input corpus and create new sub-word tokens
 - rare patterns: stick to word-based tokens

Sub-word Tokenization

- Important factor:
 - vital role takes the size of the vocabulary
 - large vocabulary size allows for more common words to be tokenized
 - smaller vocabulary requires more sub-words to be created to create every word in the text without using the <UNK> token

Byte-Pair Encoding

- Byte-Pair Encoding (BPE) used in language models:
 - GPT-2
 - RoBERTa
 - XLM
 - FlauBERT
 - ...
- Pre-tokenization:
 - space tokenization

Byte-Pair Encoding

- Simple form of data compression algorithm:
 - the most common pair of consecutive bytes of data is replaced with a byte that does not occur in that data
- first described in the article “A New Algorithm for Data Compression” published in 1994

Byte-Pair Encoding

- Use in NLP to find the best representation of text using the smallest number of tokens:
 - Add an identifier (</w>) at the end of each word to identify the end of a word and calculate the word frequency in the text.
 - Split the word into characters and then count character frequency.
 - From the character tokens, for a predefined number of iterations, count the frequency of the consecutive byte pairs and merge the most frequently occurring byte pairing.
 - Iterate until the iteration or token limit is reached.

Different Pre-tokenizers

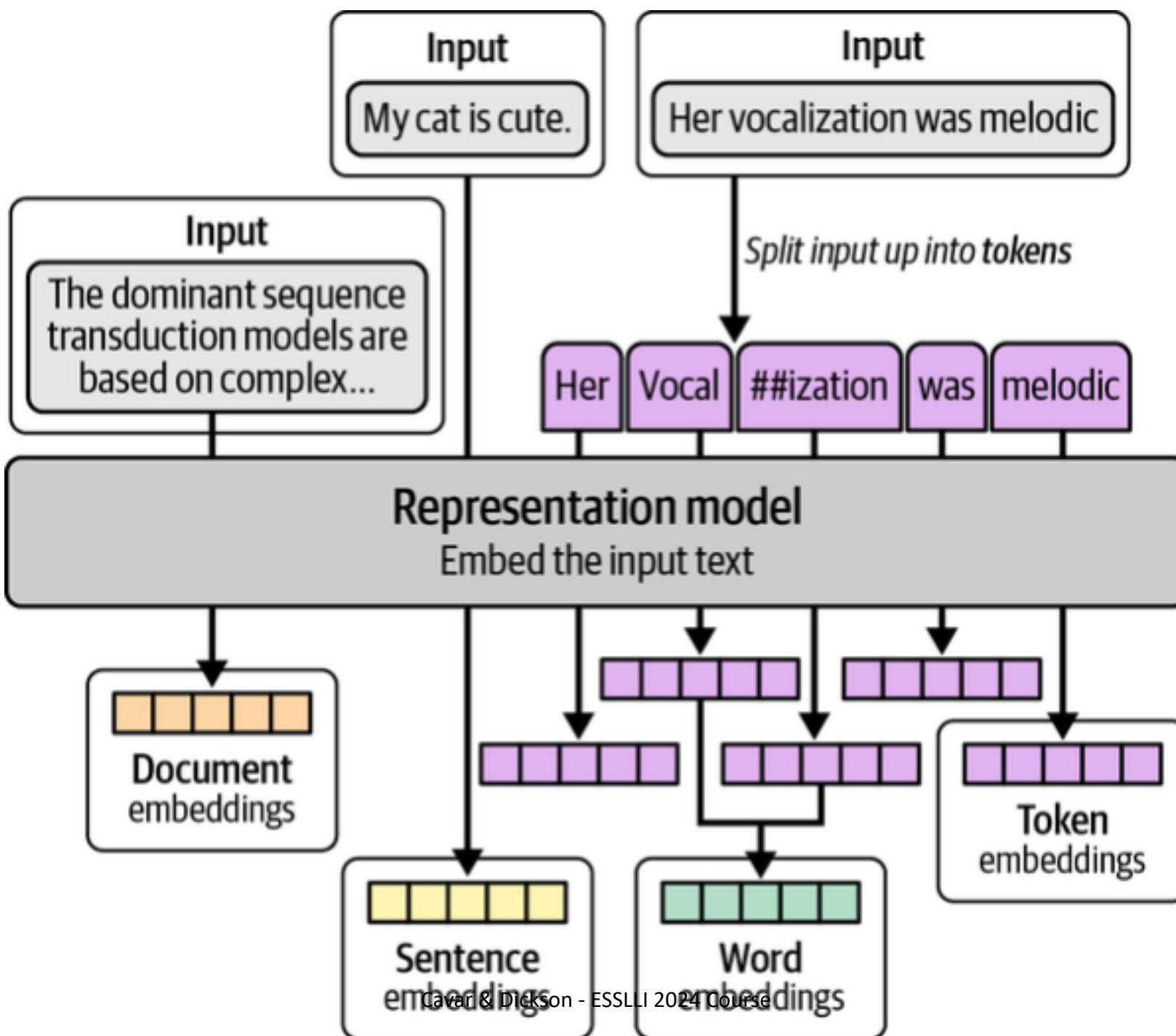
- Whitespace Pre-tokenizer:
 - "this", "sentence's", "content", "includes:", "characters", "spaces", "and", "punctuation.",
- BERT Pre-tokenizer:
 - "this", "sentence", "", "s", "content", "includes", ":", "characters", "", "spaces", "", "and", "punctuation", ".",
- GPT-2 Pre-tokenizer:
 - "this", "sentence", "s", "content", "includes", ":", "characters", "", "spaces", "", "and", "punctuation", ".",
- ALBERT Pre-tokenizer:
 - "_this", "_sentence's", "_content", "_includes:", "_characters", "_spaces", "_and", "_punctuation.",

Next Word Prediction

- Applications
 - Texting
 - Chatbots
- Examples:
 - Finite State Machines
 - N-gram Models

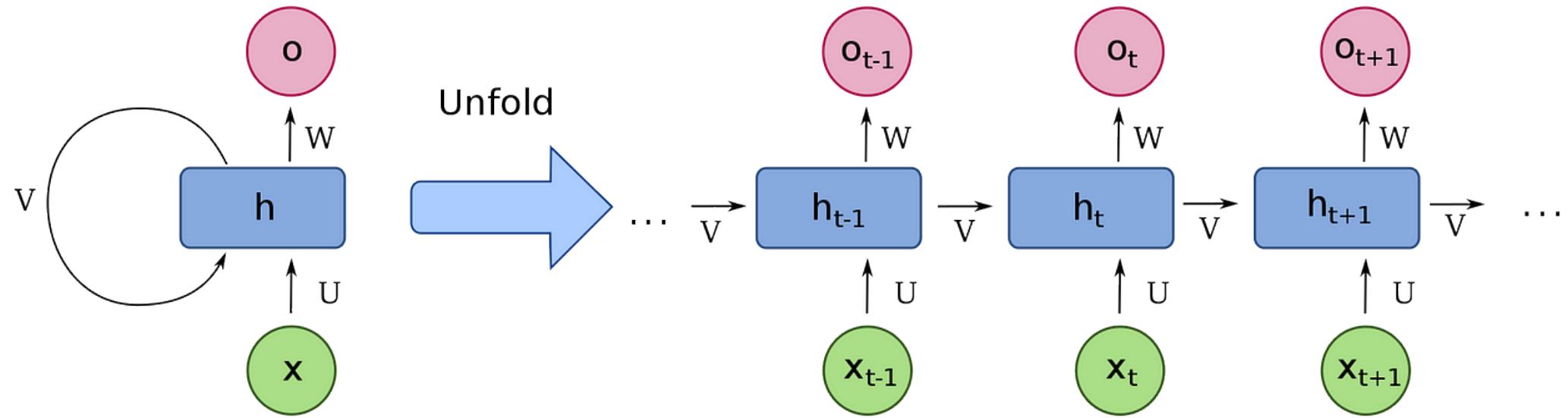
Section 2

- Recurrent Neural Networks
- Attention
- Transformers



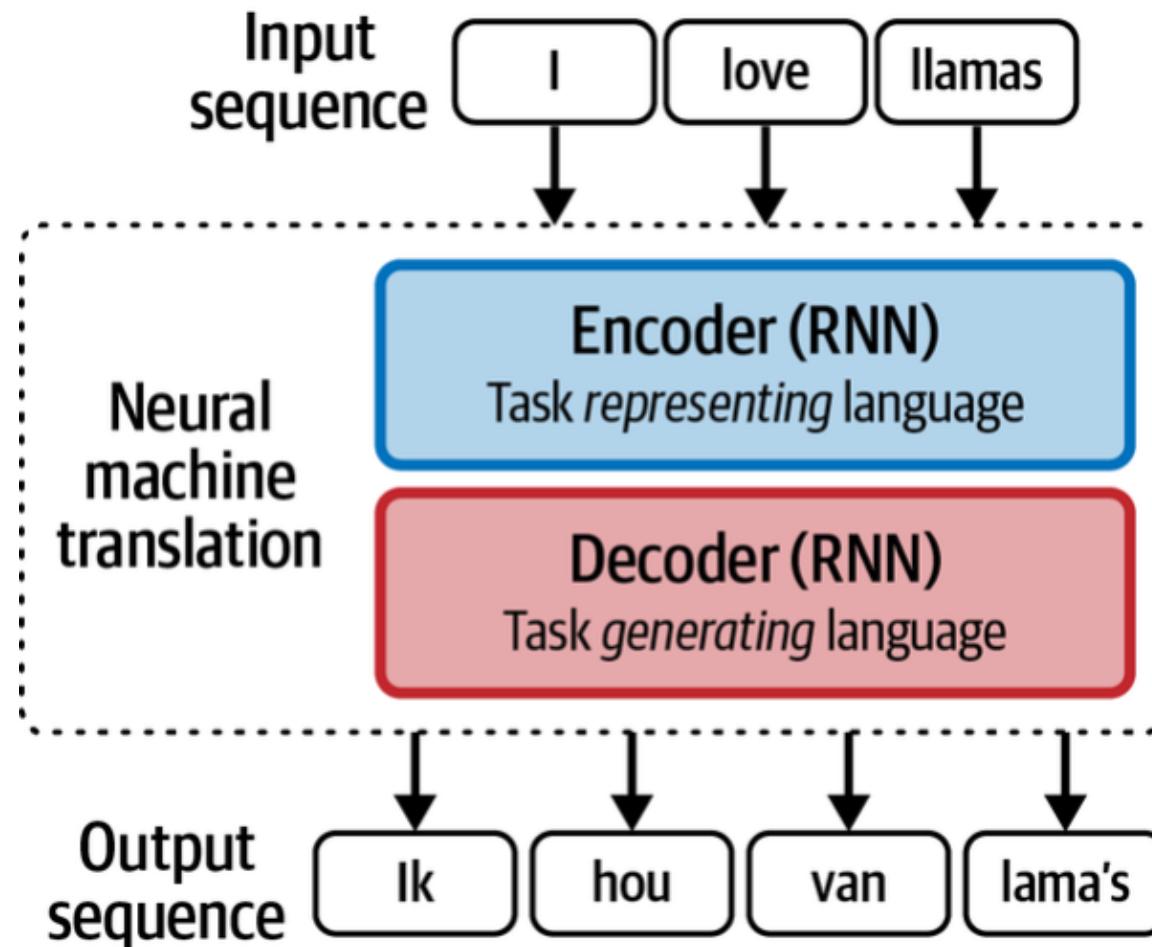
Recurrent Neural Networks (RNNs)

- **Sequence Data Processing:**
- RNNs are well-suited for tasks where the order of input data is important.
- They process inputs sequentially, maintaining a hidden state that captures information from previous steps.
- **Memory Capability:**
- RNNs can remember previous inputs due to their internal state, making them effective for tasks requiring context, like language translation or sentiment analysis.

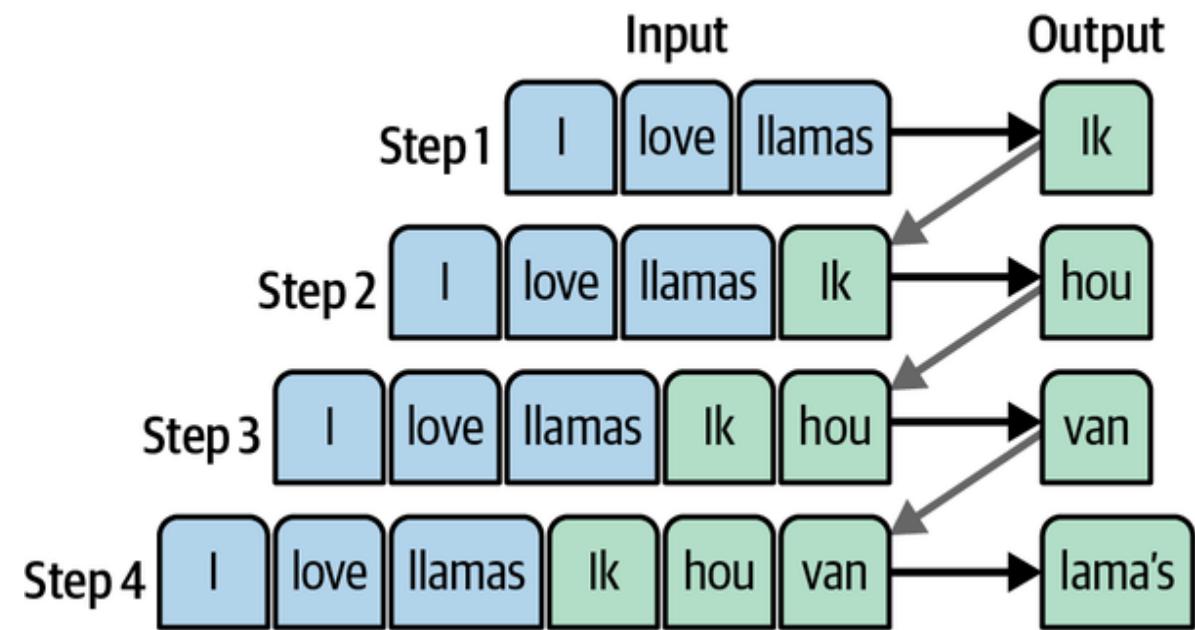


Word2vec provides a static vector embedding for each word. We can use RNNs to encode the **context sensitive** word embeddings.

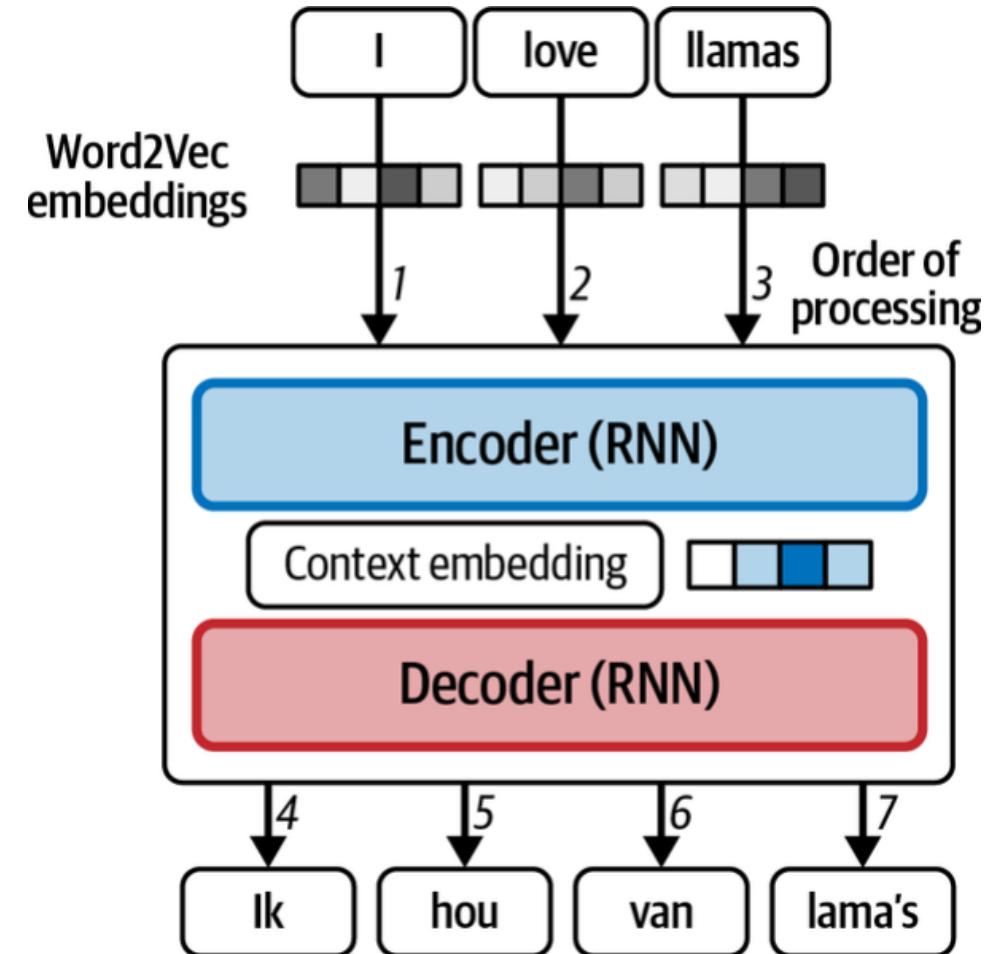
<eos> token signals that the model is finished generating



Word generation is **autoregressive**:
when predicting the next word,
the model must ingest all
previously generated words



We can use **word2vec** as input the initial starting point. Note the order that each token is processed and generated.



Limitations of using the context embedding

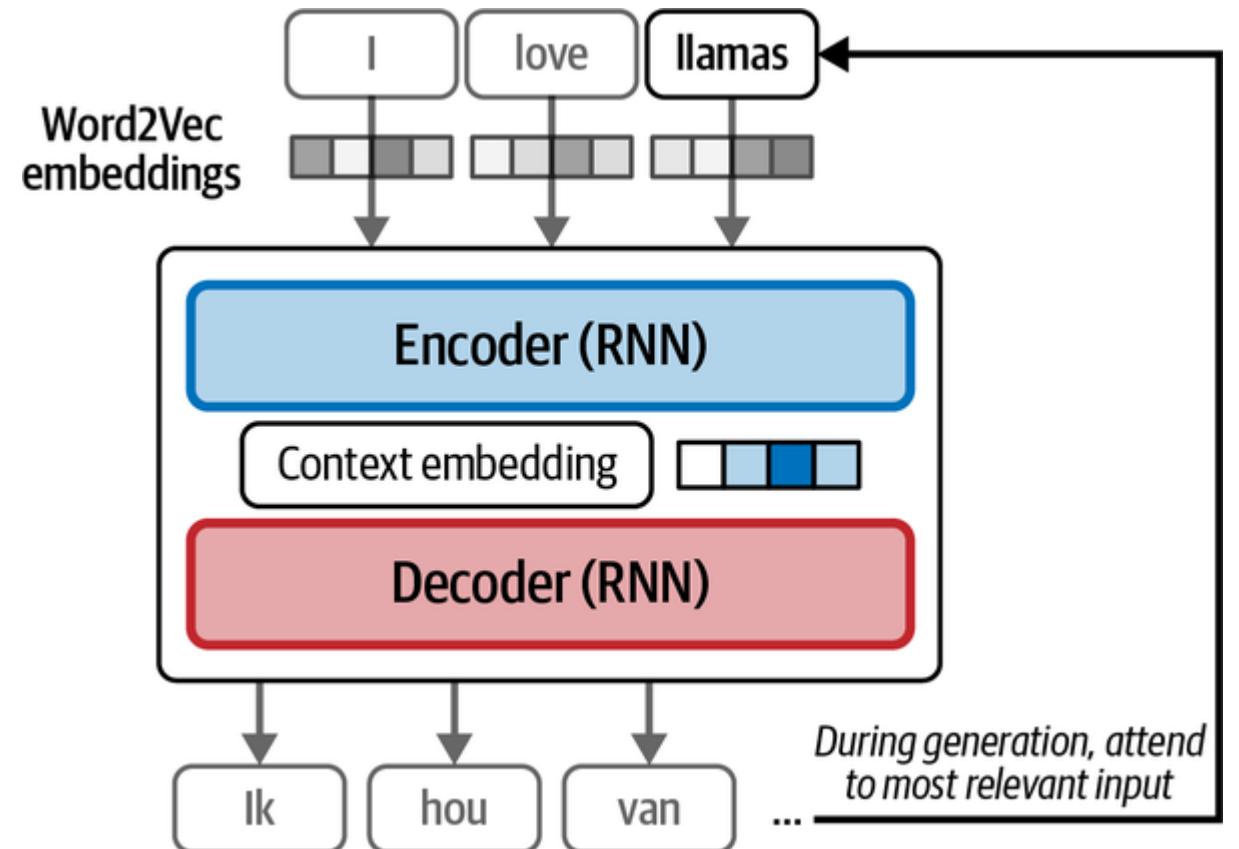
- **Information Bottleneck:**
 - Loss of important details due to compressing the entire sequence into a fixed-size vector.
- **Difficulty with Long Sequences:**
 - Struggles to retain information over long sequences, leading to the vanishing gradient problem.

Solution?

Solution?

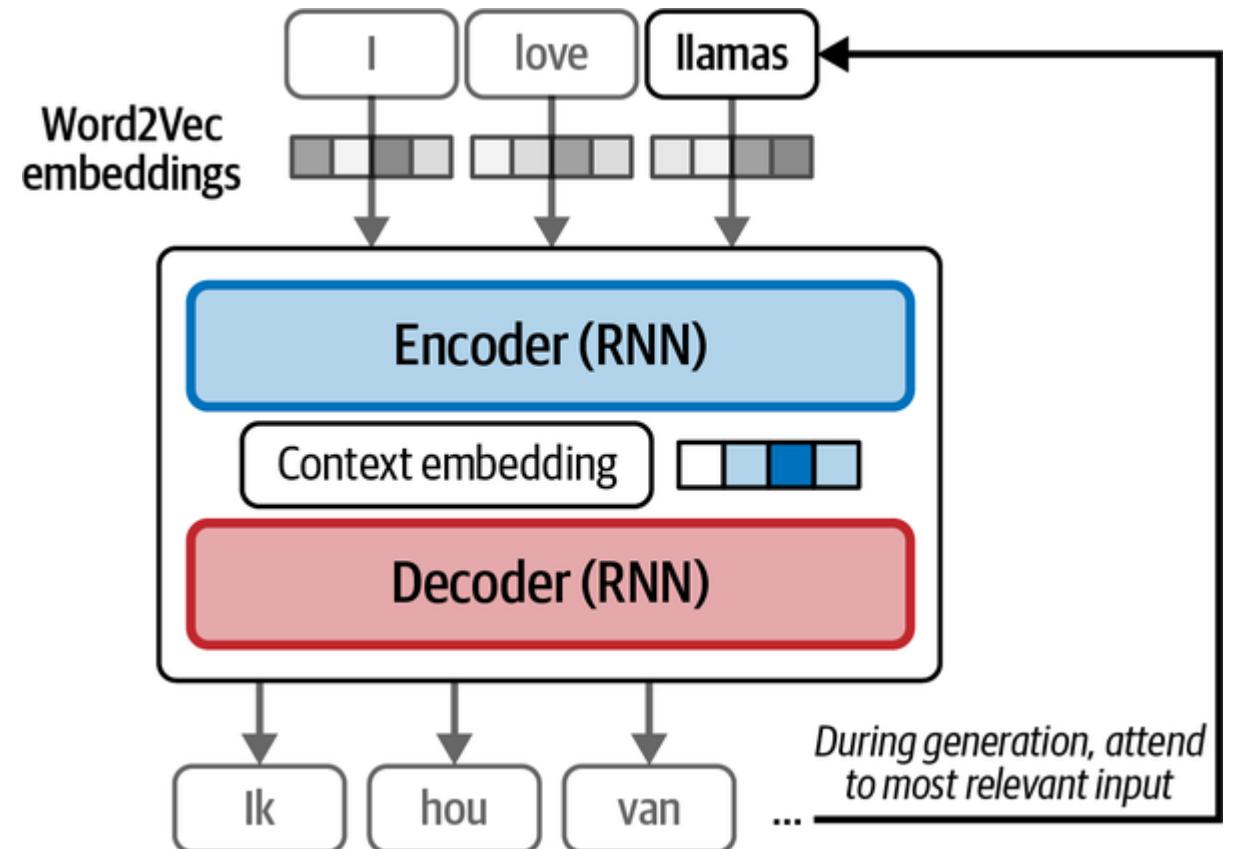
- Attention
 - The model pays attention to the most helpful part of the input for predicting the next word

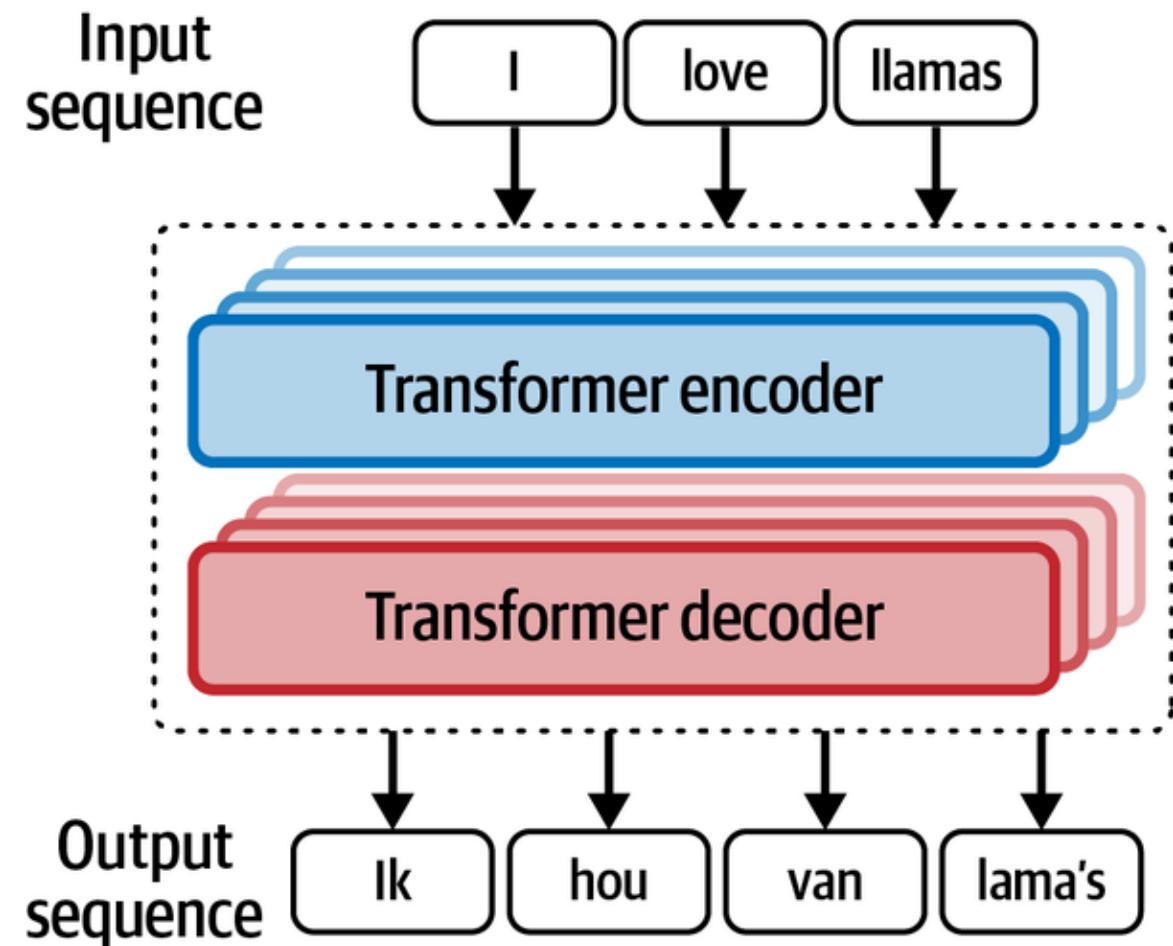
Model is trained to know how to distribute its attention



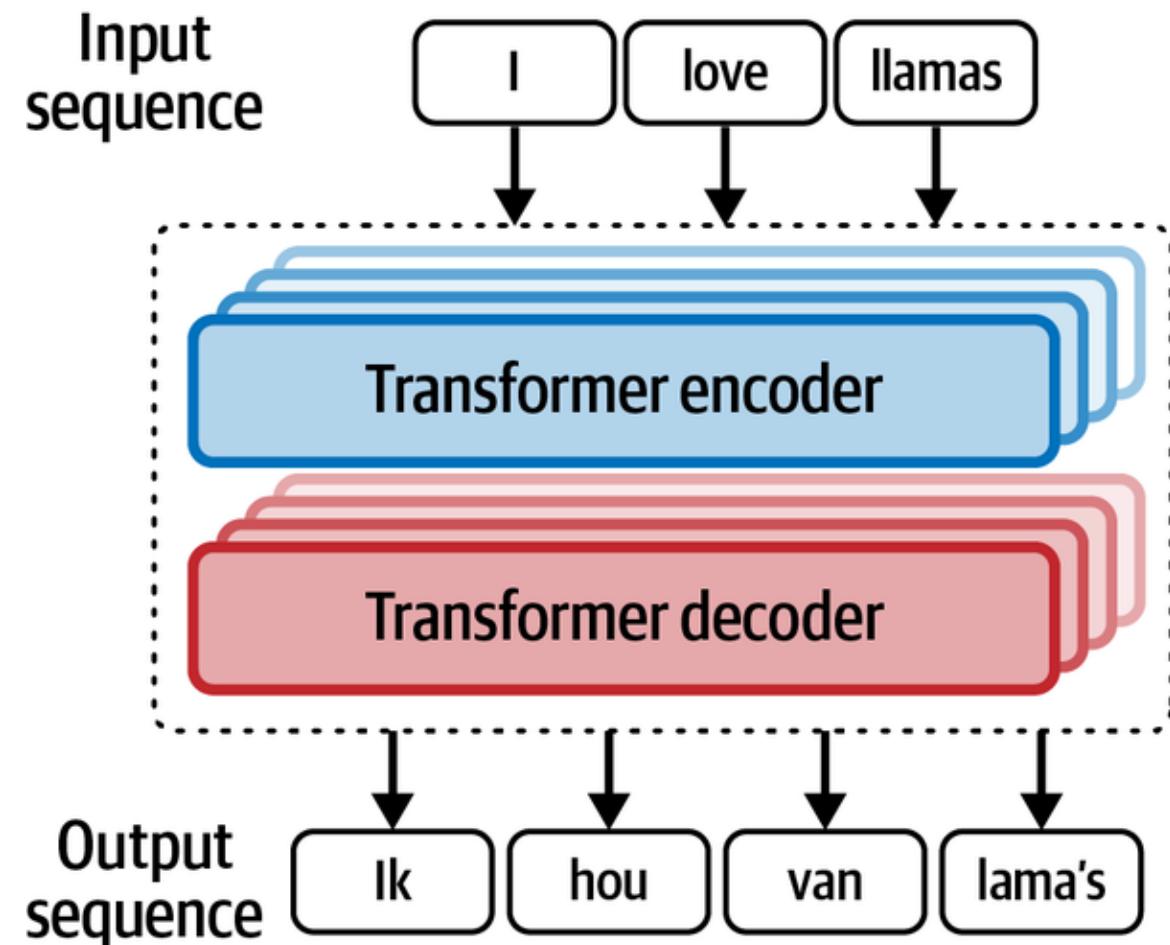
Model is trained to know how to distribute its attention

Still have problems with **catastrophic forgetting**



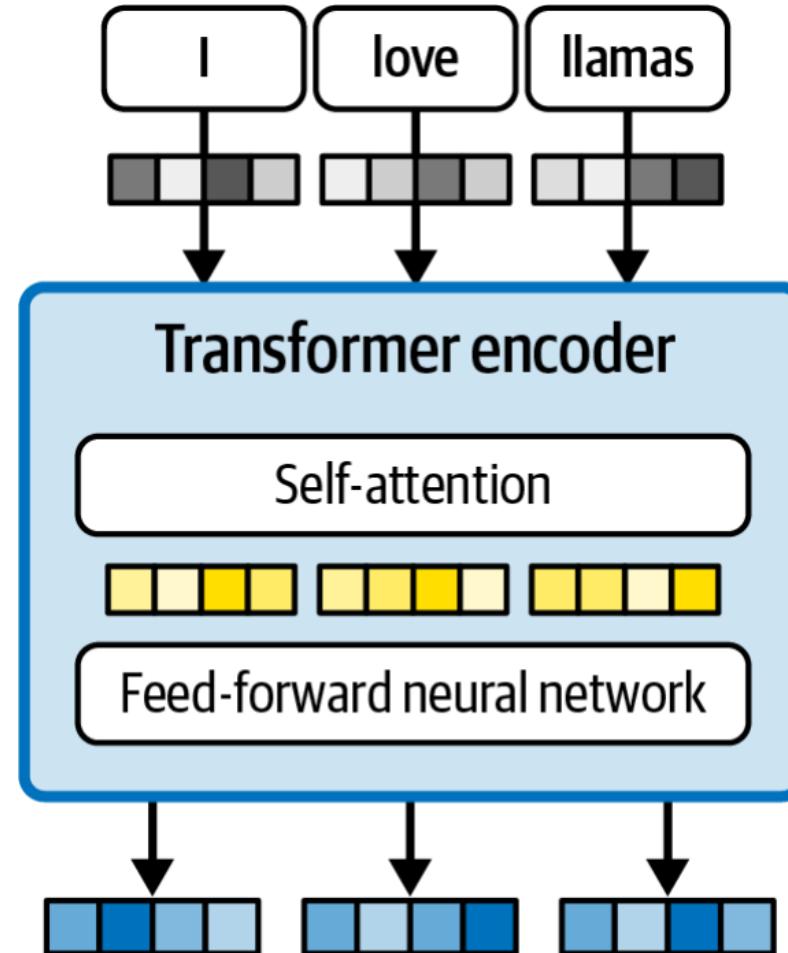


Attention is All You Need (Vaswani et al. 2017)

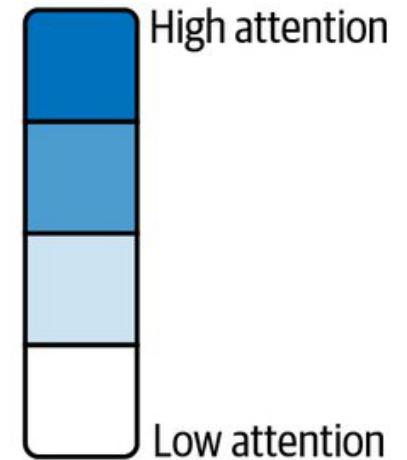
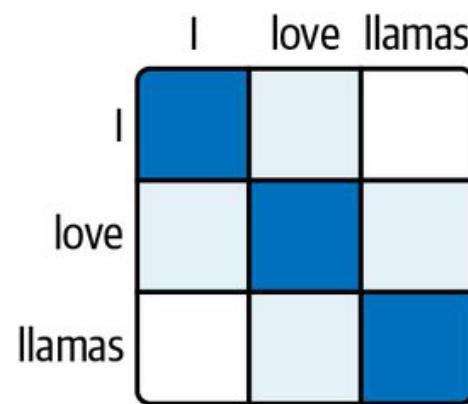


Input can be processed in parallel, all at once

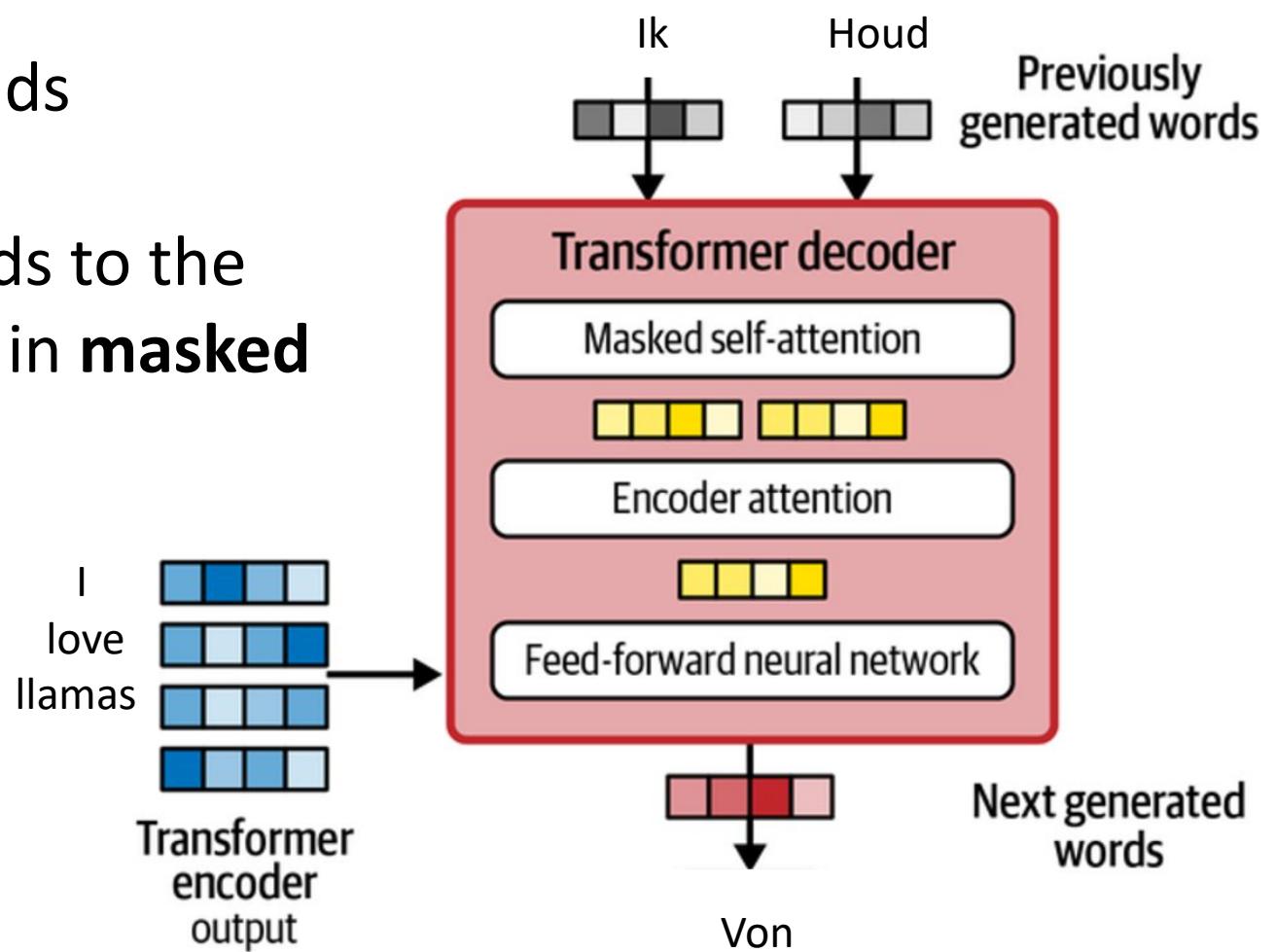
self-attention allows each token to distribute its attention to each other token including itself



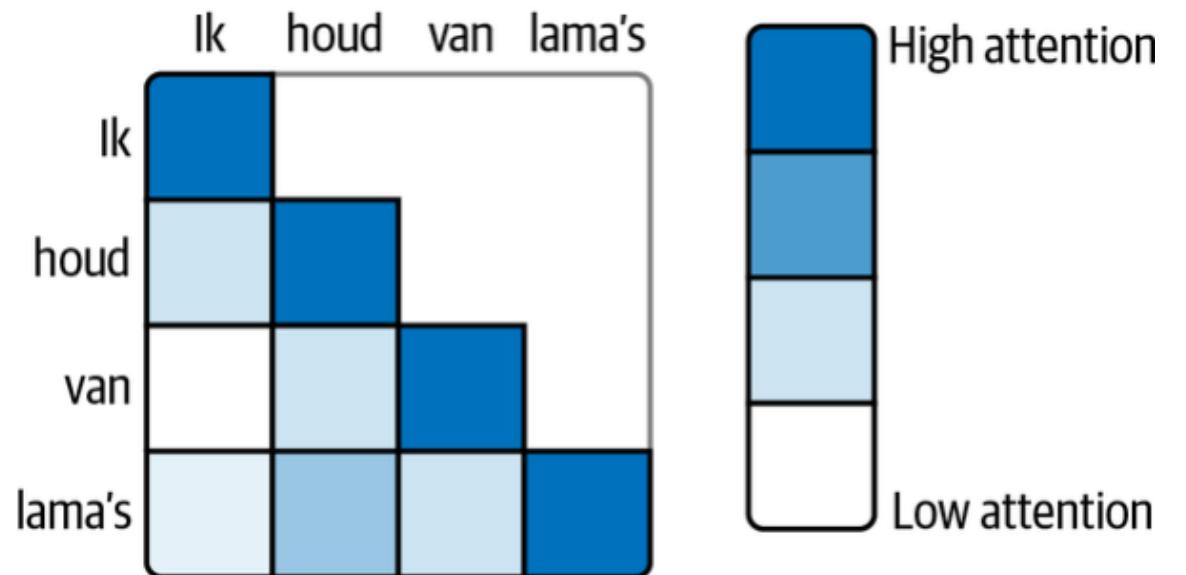
In the encoder, each token sees every other token of the input



In the decoder, the model attends to the output of the encoder in **encoder attention** and attends to the currently generated translation in **masked self-attention**



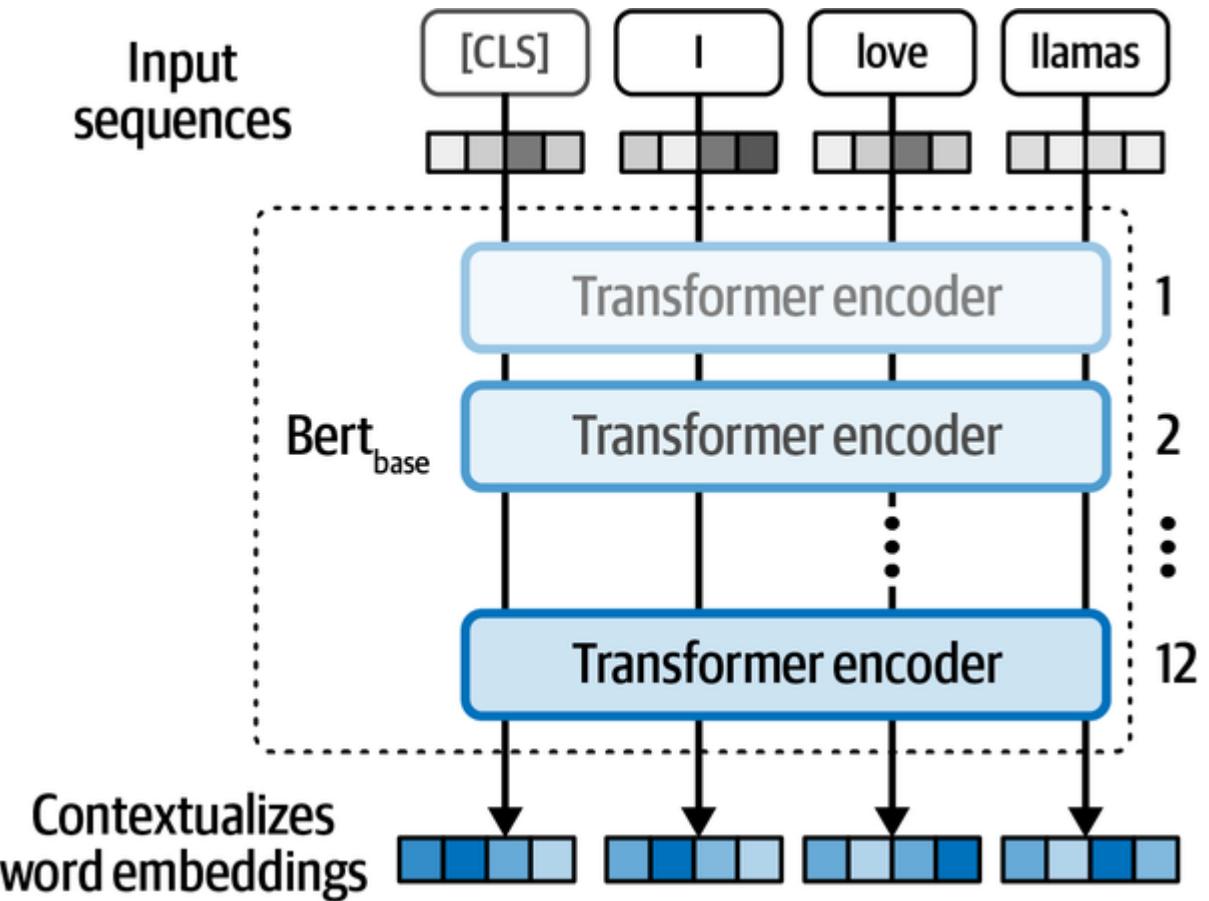
Masked self-attention
prevents the model
from cheating and
looking into the future
during training



Encoder Only Models

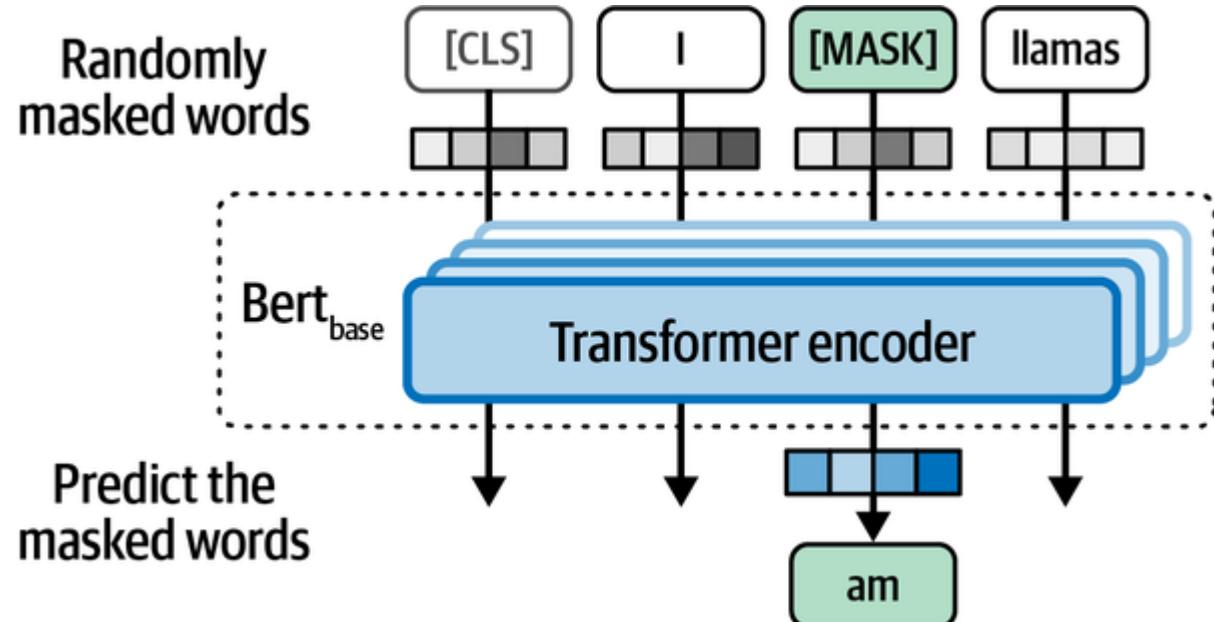
- **Purpose:**
 - Used for understanding and classification tasks.
 - Examples: Text classification, sentiment analysis, named entity recognition (NER).
- **Architecture:**
 - Multiple layers of encoders.
 - Processes input sequence to capture contextual information.
 - Outputs a fixed-size vector or sequence of vectors.
- **BERT (Bidirectional Encoder Representations from Transformers):**
 - Considers both past and future tokens for bidirectional context.

[CLS] or classification token, is used to represent the entire sequence

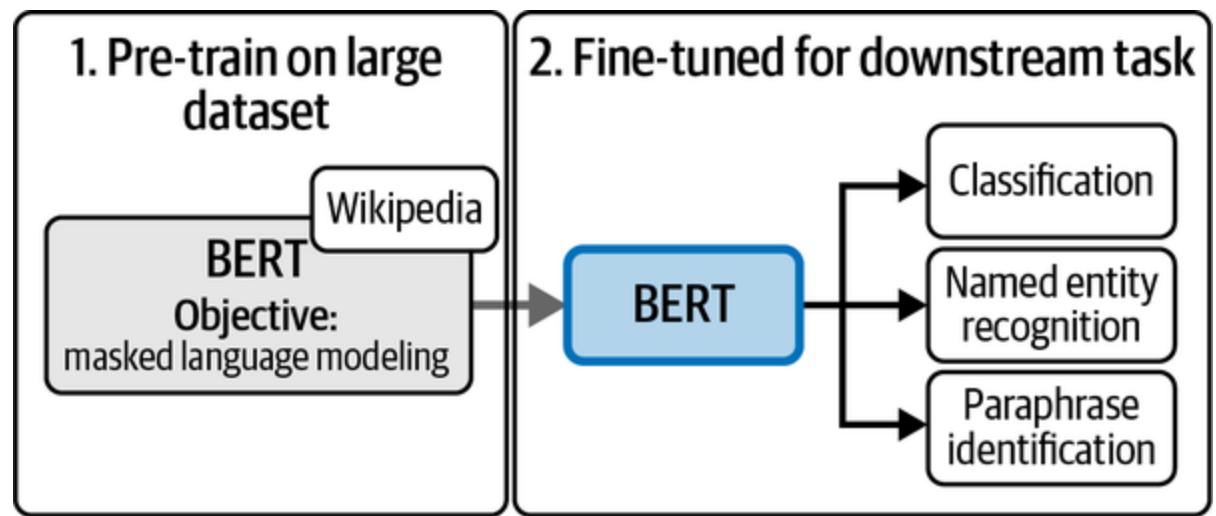


BERT implements **masked language modeling** which randomly blocks out tokens in the sequence.

The models must learn to fill in the missing gaps.

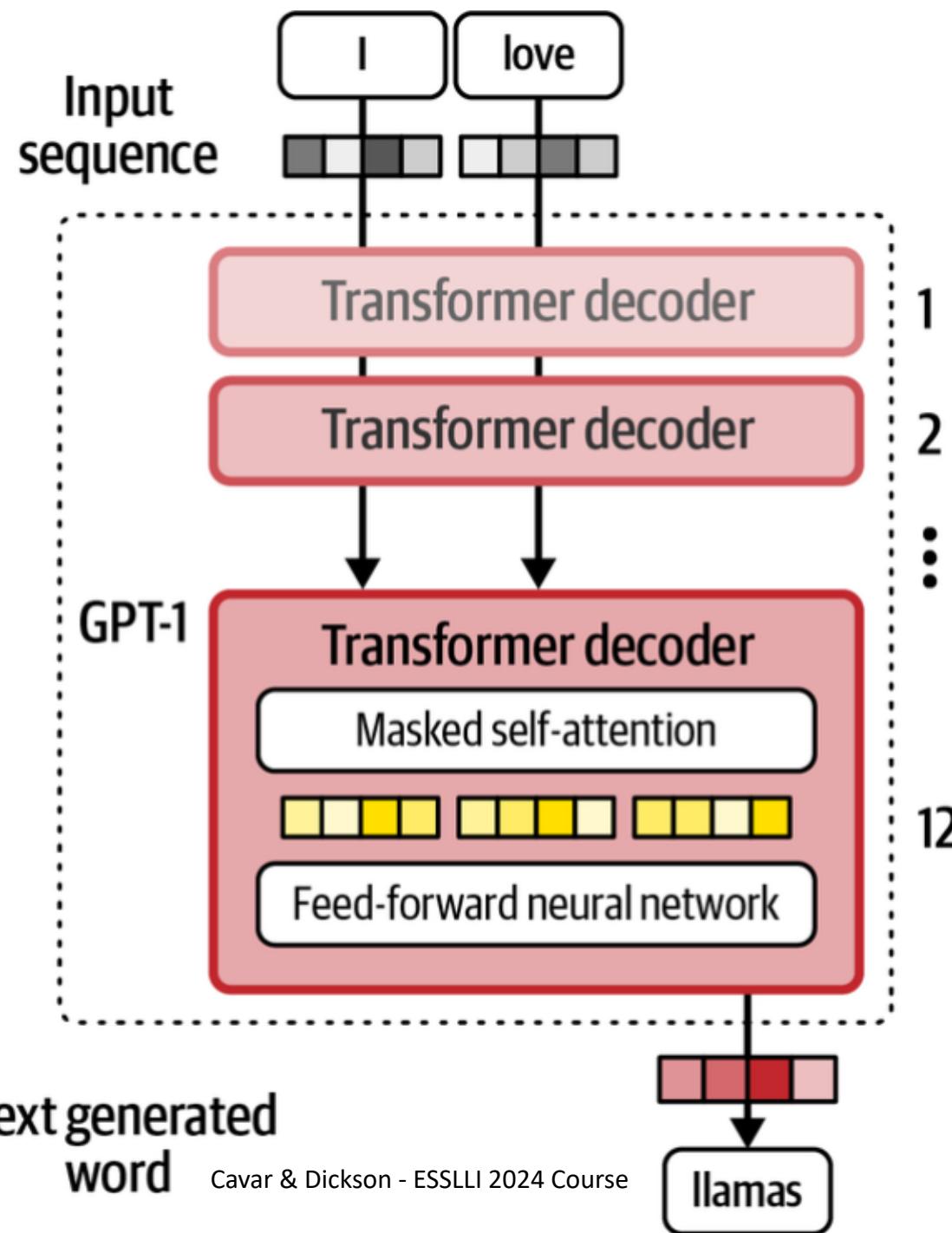


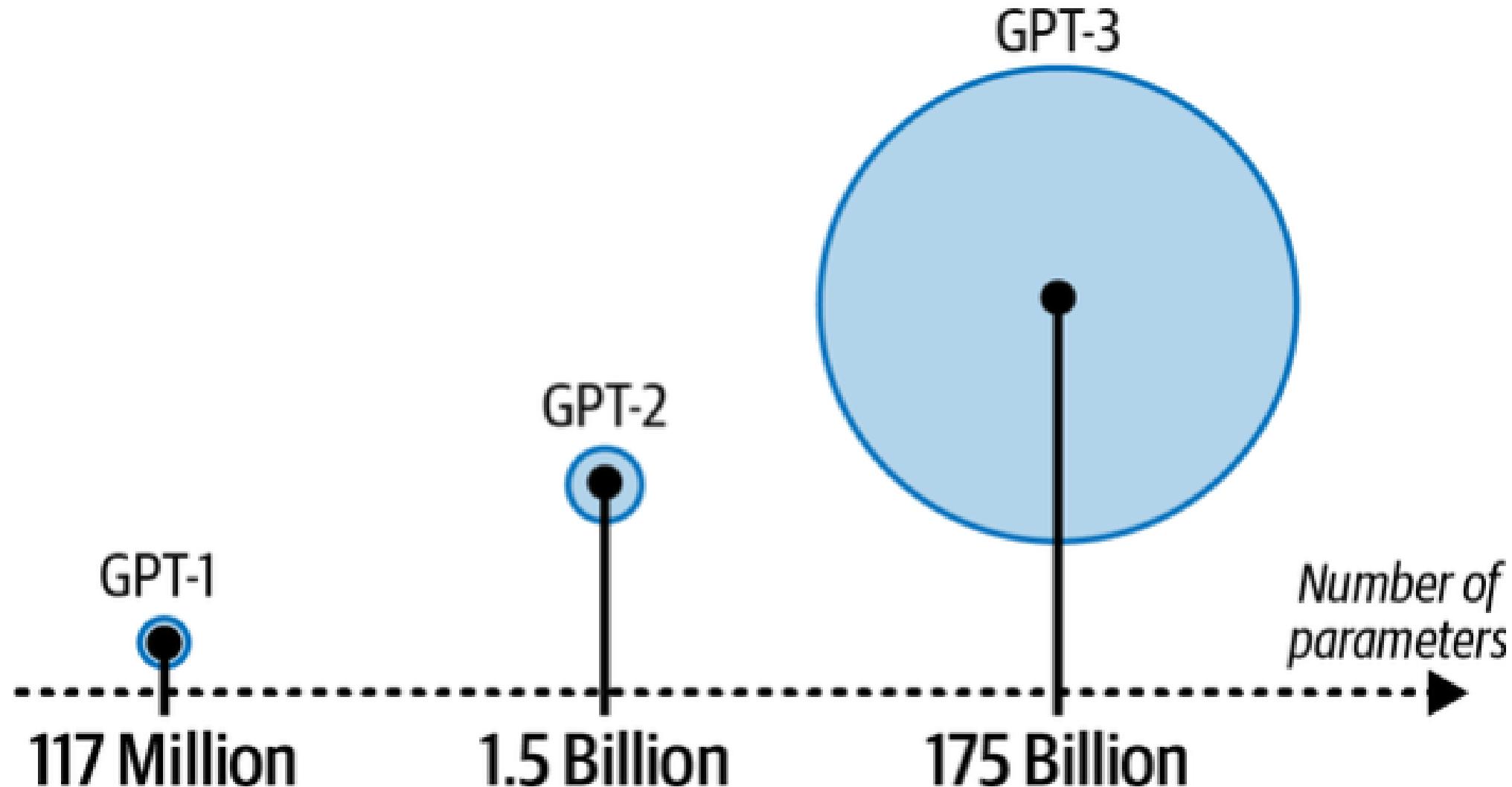
Fine-tuning involves initially pre-training on a large dataset, and then additionally training for a short amount of time on task specific data.



Decoder Only Models

- **Purpose:**
 - Used for sequence generation or prediction tasks.
 - Examples: Text generation, language modeling, autocomplete.
- **Architecture:**
 - Multiple layers of decoders.
 - Generates output tokens step-by-step, attending to previous tokens.
 - Predicts the next token based on the context of generated tokens.
- GPT, Llama, Mistral, Claude, etc.





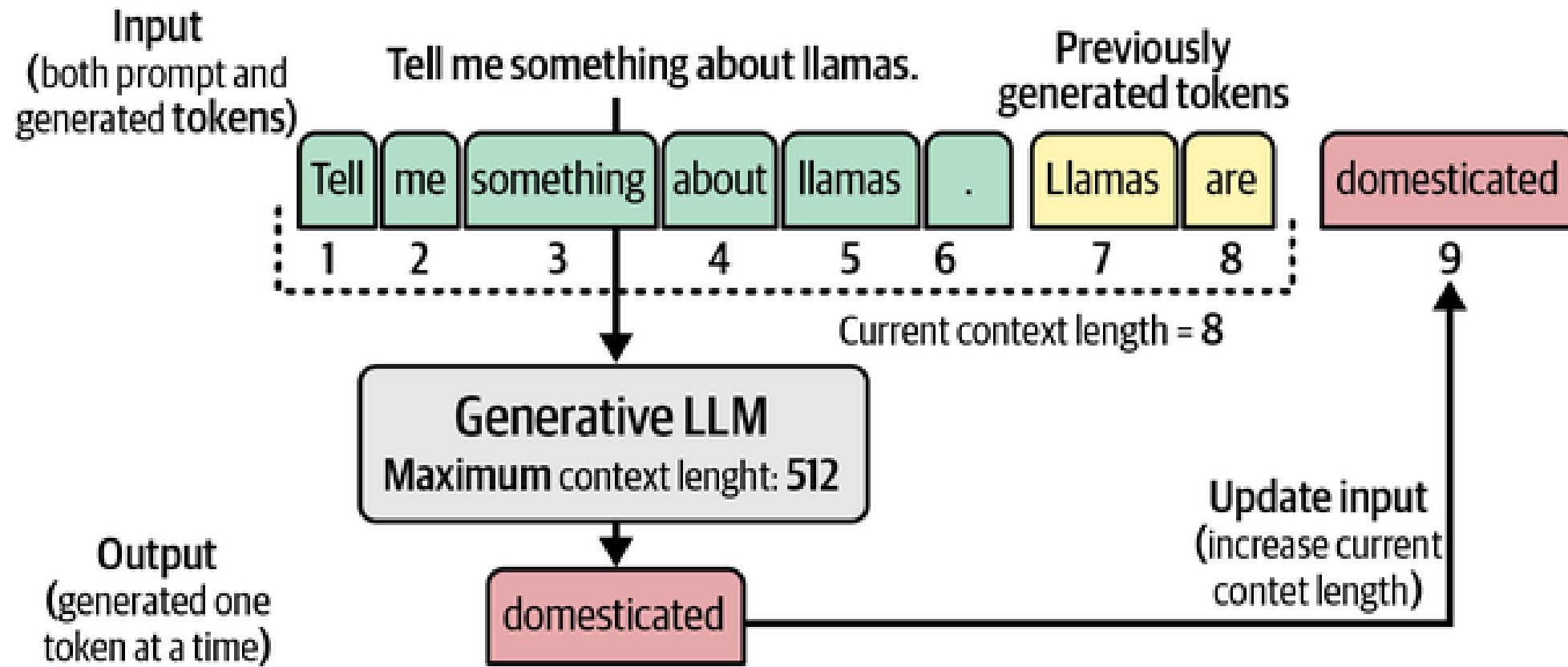
User query
(prompt)

Tell me something about llamas

Generative LLM
Task:
complete the input

Output
(completion)

Llamas are domesticated South American camelids, widely used as pack animals by Andean cultures since pre-Hispanic times. With their fluffy coat, long neck, and distinctive facial features...



Context length is the maximum number of tokens we can give the model at one time

How does the model choose the word?

Large Language Models

- Most common method for decoding in LLMs: Sampling
 - Sampling from a model's distribution over words means to choose random words according to their probability assigned by the model.
 - Iteratively choose a word to generate according to its probability in context as defined by the model.
 - More likely: generate words that the model predicts have a high probability in the context and less likely to generate words that the model predicts have a low probability.

Large Language Models

- Random sampling
- Problem:
 - mostly generating sensible, high-probable words
 - low-probability words in the tail of the distribution get chosen often enough to result in generating weird sentences.
- Instead of random sampling:
 - Use alternative sampling methods that avoid generating the very unlikely words.

Large Language Models

- Sampling methods: quality and diversity
- Most probable words tend to produce generations that are
 - more accurate
 - more coherent
 - more factual
 - but also more boring and more repetitive

Large Language Models

- Methods that give more weight to the middle-probability words:
 - are more creative
 - more diverse
 - less factual
 - more likely to be incoherent or otherwise low-quality

Large Language Models

- Top-k sampling (greedy decoding)
 - truncate the distribution of continuation words to the top k most likely words
 - renormalize to produce a legitimate probability distribution
 - randomly sample from within these k words according to their renormalized probabilities
- When $k = 1$, top-k sampling is identical to random sampling

Large Language Models

- Top-k sampling (greedy decoding)
- Setting $k > 1$ leads to
 - selected continuation word is not necessarily the most probable one
 - still probable enough
 - choice results in generating more diverse but still high-enough-quality text

Large Language Models

- Choose in advance a number of words k
- For each word in the vocabulary V , use the language model to compute the likelihood of this word given the context $p(w_t | w_{<t})$
- Sort the words by their likelihood, and throw away any word that is not one of the top k most probable words.
- Renormalize the scores of the k words to be a legitimate probability distribution.
- Randomly sample a word from within these remaining k most-probable words according to its probability.

Large Language Models

- Problems with: Top-k
 - k is fixed, but the shape of the probability distribution over words differs in different contexts
 - k = 10, sometimes the top 10 words will be very likely and include most of the probability mass, but other times the probability distribution will be flatter and the top 10 words will only include a small part of the probability mass

Large Language Models

- Nucleus or Top-p Sampling
 - keep not the top k words, but the top p percent of the probability mass
 - goal: truncate the distribution to remove very unlikely words
 - by measuring probability rather than the number of words
 - hope: the measure will be more robust in very different contexts, dynamically increasing and decreasing the pool of word candidates

Large Language Models

- Temperature Sampling
 - don't truncate the distribution, but instead reshape it
 - Intuition:
 - thermodynamics, where a system at a high temperature is very flexible and can explore many possible states, while a system at a lower temperature is likely to explore a subset of lower energy (better) states.
- Low-temperature sampling
 - smoothly increase the probability of the most probable words and decrease the probability of the rare words

Large Language Models

- Implementation of Temperature Sampling
 - dividing the logit by a temperature parameter τ before it is normalized by passing it through the softmax in low-temperature sampling, $\tau \in (0, 1]$
 - Instead of computing the probability distribution over the vocabulary directly from the logit, first divide the logits by τ , computing the probability vector y as:

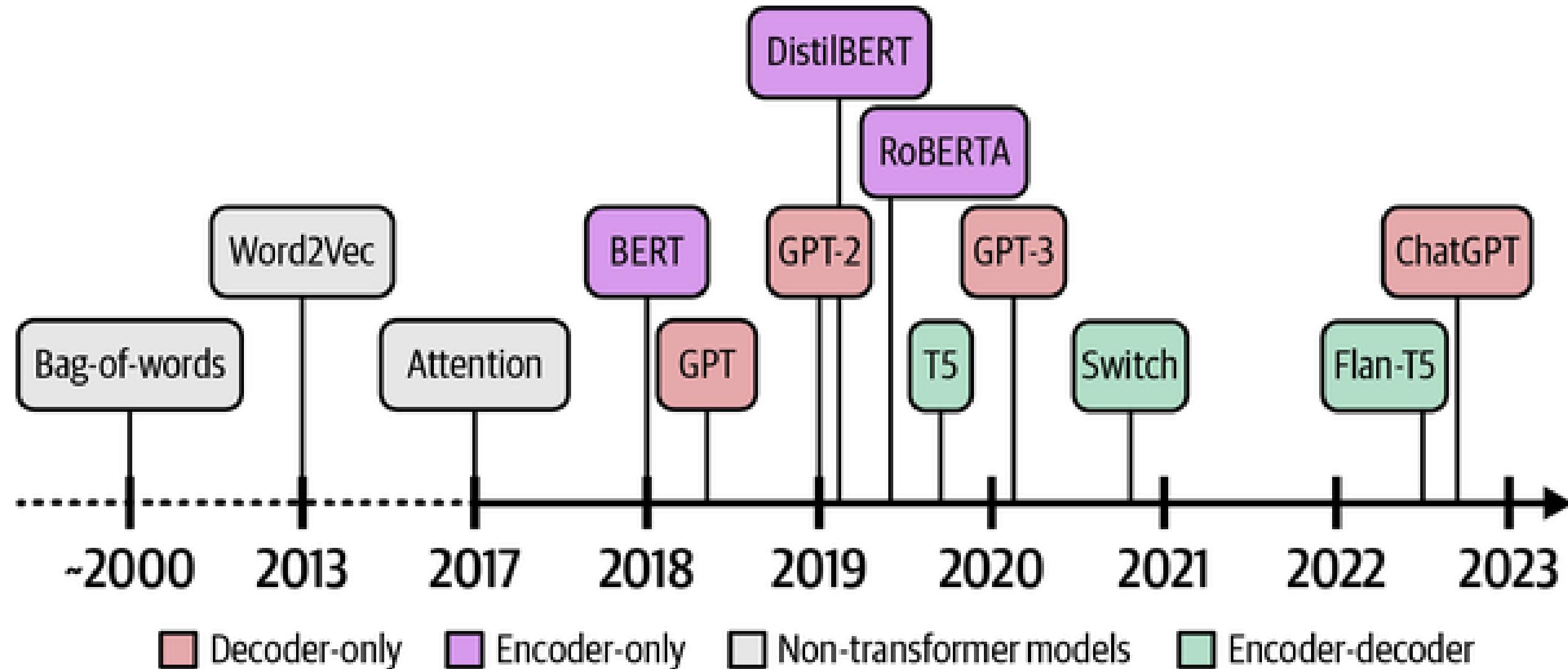
$$\mathbf{y} = \text{softmax}(u/\tau)$$

Large Language Models

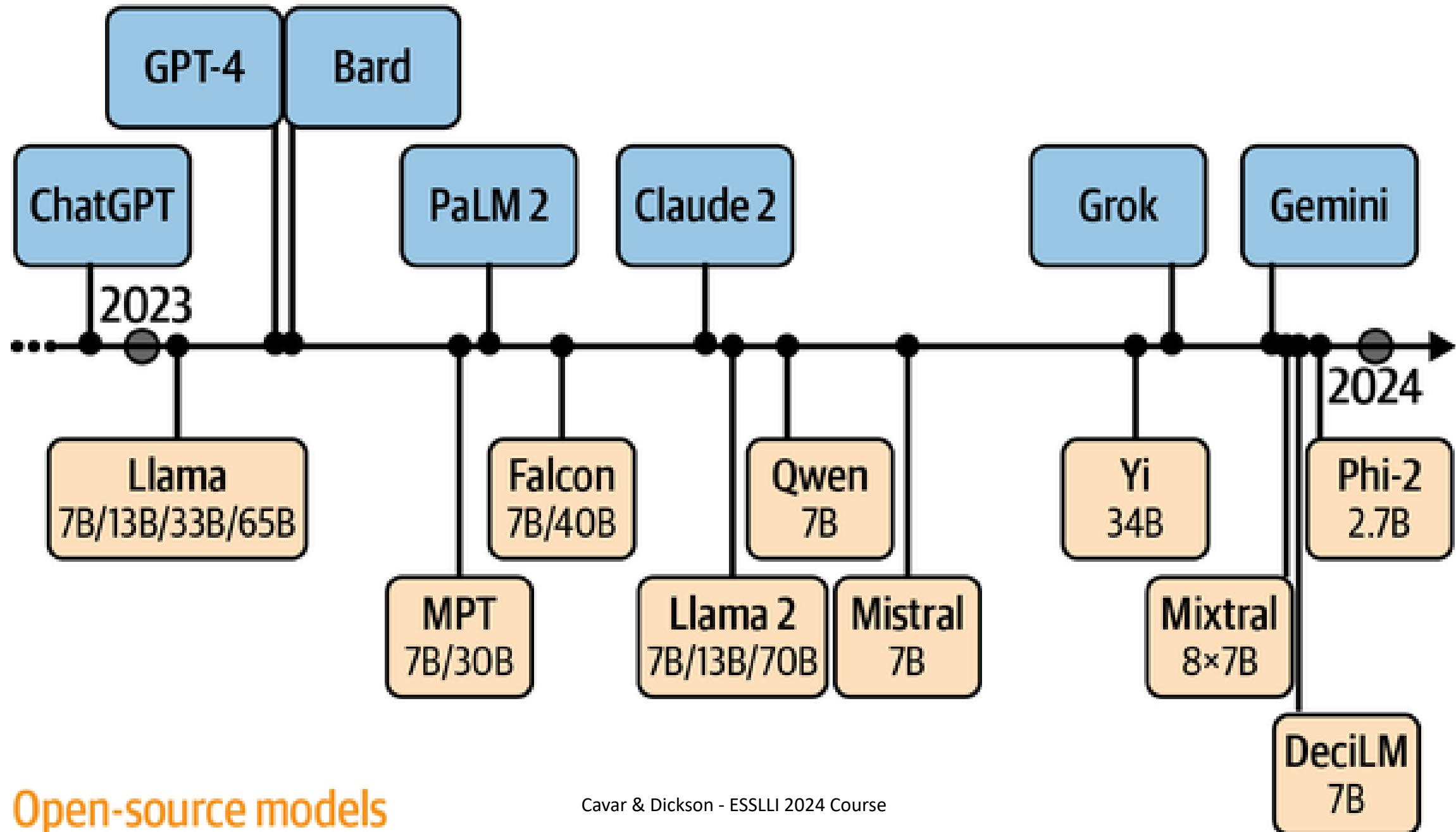
- When τ is close to 1 the distribution doesn't change significantly
- The lower τ the larger the scores being passed to softmax
 - dividing by a smaller fraction $\tau \leq 1$ results in larger scores
 - softmax pushes high values toward 1 and low values toward 0
 - Larger numbers passed to a softmax result in a distribution with increased probabilities of high-probability words and decreased probabilities of low probability words
- τ closer to 0 the probability of most likely word approaches 1

Large Language Models

- We may want to flatten the word probability distribution:
 - Temperature sampling can help
 - high-temperature sampling: use $\tau > 1$.



Proprietary models



Open-source models

@danielhanchen x.com

7/30/2024

Meta Llama 3.1 8b, 70b, 405b Base & Instruct fp16, fp8 unsplash
15.6T tokens, 128K context
Tool Multilingual

Benchmark	Llama 2 9B	Mistral 7B	Llama 70B	Mistral 22B	GPT-3.5 Turbo	Llama 3 405B	GPT-4o	Gemma-2.6b
MMLU (5-shot)	69.4	72.3	61.1	83.6	70.7	85.1	89.1	89.9
MMLU (whole, CvTF)	73.0	72.3 ^a	60.5	86.0	79.9	68.8	85.4	88.7
MMLU-Pro (5-shot, CvTF)	48.3	-	36.9	66.4	56.3	49.2	73.3	77.0
IPEval	80.4	73.6	57.6	87.5	72.7	69.9	88.6	84.5
HumanEval (5-shot)	72.6	54.1	41.2	80.5	75.6	68.0	80.9	86.4
MBPP EvalPlus (whole)	72.8	71.1	69.5	95.1	80.0	81.6	96.8	94.2
GSM8K (whole, CvTF)	84.5	76.7	53.2	95.1	80.0	81.6	96.9	96.1
MATH (whole, CvTF)	51.9	41.3	18.8	68.0	54.1	43.1	73.5	71.1

Llama3 RoPE Patch

```
def apply_scaling(freq: torch.Tensor):
    # Values obtained from grid search
    # "low_freq_factor": 1.0,
    # "high_freq_factor": 4.0,
    # "scale_factor" = 1
    low_freq_factor = freq * 1.0
    high_freq_factor = freq * 4.0
    original_max_headings = 8192,
    scaling = 1 / sqrt(query_pre_attn_scalar)
    query_pre_attn_scalar = 256 * 128
    head_dim = 256
    attn_logit_softcapping = 50
    final_logit_softcapping = 30
    Q *= ops.cast(scaling, dtype=Q.dtype)
    return Q
```

Chat Template

```
def __init__(self, scope_type="llama3", low_freq_factor=1.0, high_freq_factor=4.0, scale_factor=1, original_max_headings=8192, max_position_embeddings=4096, swiglu=False):
    self.scope_type = scope_type
    self.low_freq_factor = low_freq_factor
    self.high_freq_factor = high_freq_factor
    self.scale_factor = scale_factor
    self.original_max_headings = original_max_headings
    self.max_position_embeddings = max_position_embeddings
    self.swiglu = swiglu
```

Data Mixture Preprocessing

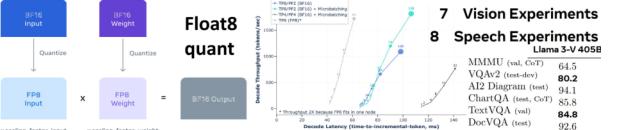
50% General knowledge
25% Math & Reasoning
17% Code Tasks
8% Multilingual

Pipeline Parallelism

38 to 43% bfloat16 MFU
6 staged long context expansion

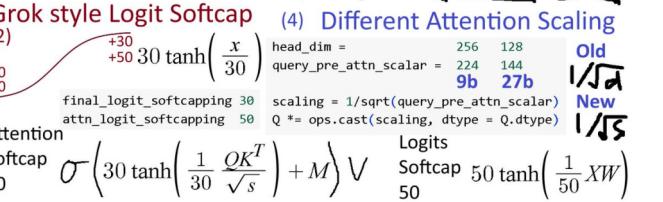
This long-context pre-training stage was performed using approximately 800B training tokens.
4.1.5 Model Averaging Languages supported: German, French, Italian, Portuguese, Hindi, Spanish, and Thai.

Finally, we average models obtained from experiments using various versions of data or hyperparameters at each RM, SFT, or DPO stage (Izmailov et al., 2019; Wortzman et al., 2022; Li et al., 2022).



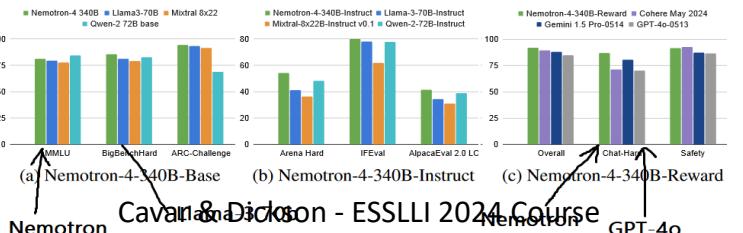
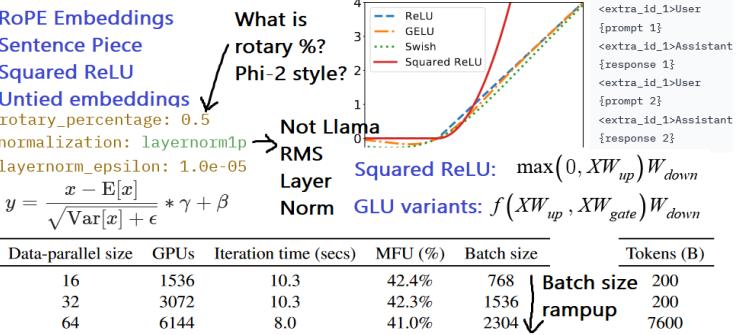
7 Vision Experiments 8 Speech Experiments Llama 3-405B

	Mistral	LLaMA-3	Gemma-2	Released
Benchmark	7B	8B	2.6B	9B
MMLU	62.5	66.6	51.3	71.3
MMLU Instruct	68.4	72.3	76.2	75.2
Base and Instruct released!				



Nemotron-4 340 billion 3 models: 8 trillion tokens

Number of transformer layers	Hidden dimension	Number of attention heads	Number of KV heads	Sequence length	Vocabulary size
96	18432	96	8	4096	256,000



Mistral Nemo 12b 2x faster 70% less VRAM	unslsotn
Mistral v3 Tokenizer = LlamaTokenizerFast	Uses GPT2
Mistral Nemo Tokenizer = GPT2TokenizerFast	Tokenizer
MT Bench	
Mistral 7B 6.48	
Llama 3 8B 6.85	
Mistral NeMo 7.84	

</unk>: 0, <s>: 1, </s>: 2, <INST>: 3, Tool Usage </INST>: 4, <AVAILABLE_TOOLS>: 5, </AVAILABLE_TOOLS>: 6, <TOOL_RESULTS>: 7, </TOOL_RESULTS>: 8, <TOOL_CALLS>: 9, <PREFIX>: 11, <MIDDLE>: 12, <SUFFIX>: 13, <SPECIAL_14>: 14, <SPECIAL_999>: 999

Base </s> amax diff to instruct (1.1270e-23, 1.3379e-01, 1.1219e-23, ([1.1064e-23, 1.2988e-01, 2.0294e-03, 2])

Pad token Fill in the Middle Padded Tokens

Auto appended EOS token (wrong) tokenizer.old("Hello!").input_ids[1, 22177, 1033, 2] tokenizer_new("Hello!").input_ids[1, 22177, 1033]

New HF (q_proj) (in_features=5120, out_features=4096) Old HF (k_proj) (in_features=5120, out_features=5120)

Larger hidden dim

Mistral Nemo 12b 1M possible max position embeddings More layers!

Uses 20GB! Now use dynamic RoPE

Phi-3 3.8b Mini

sliding_window: 2047, max_position_embeddings: 4096, SWA like Mistral?

use_sliding_windows = (_flash_supports_window_size and getattr(self.config, "sliding_window", None) is not None and kv_seq_len > self.config.sliding_window)

Upcasted RoPE like Gemma

```
q_embed = (q.float() * cos.float() + (rotate_half(q).float() * sin.float()))
k_embed = (k.float() * cos.float() + (rotate_half(k).float() * sin.float()))
return q_embed.to(q.dtype), k_embed.to(k.dtype)
```

Dynamic RoPE

```
class Phi3SuScaledRotaryEmbedding(_Phi3ScaledRotaryEmbedding):
    def __init__(self):
        super().__init__()
        self.qscaled = self.qscaled.to(torch.float32)
        self.kscaled = self.kscaled.to(torch.float32)
```

```
class Phi3YarnScaledRotaryEmbedding(_Phi3ScaledRotaryEmbedding):
    def __init__(self):
        super().__init__()
        self.qscaled = self.qscaled.to(torch.float32)
        self.kscaled = self.kscaled.to(torch.float32)
```

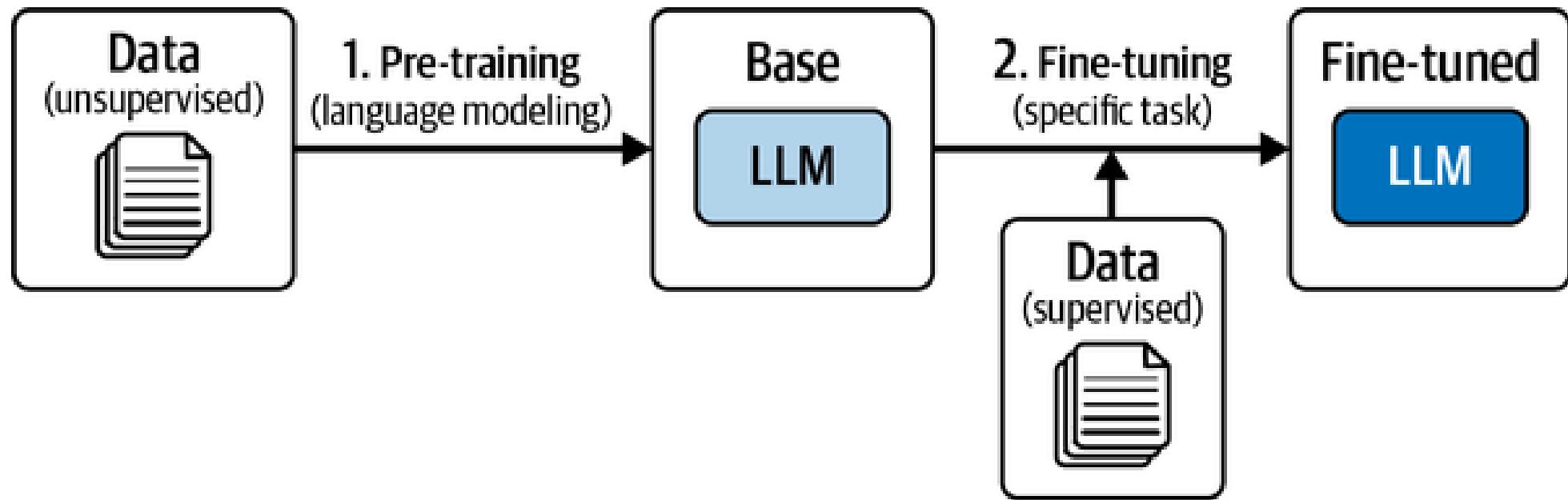
Fused MLP

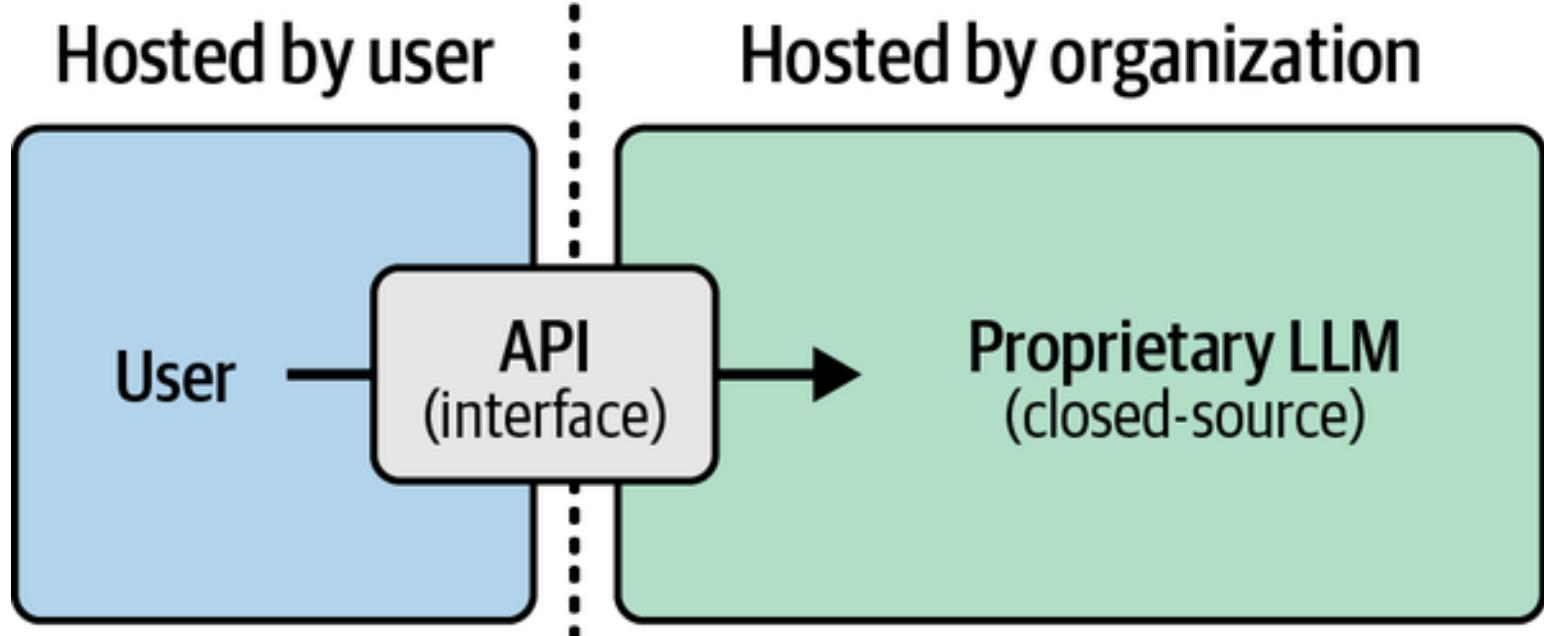
```
y = self.gate_up_proj(hidden_states)
gate, y = y.chunk(2, dim=-1)
y = y * self.activation_fn(gate)
```

Need to unfuse

Phi-3-mini	Mistral	Gemma	Llama-3-In	Mixtral	GPT-3.5
3.8b	7b	7b	8b	8x7b	version 1106
MMLU (5-Shot) 68.8	61.7	63.6	66.0	68.4	71.4

Llama-3





- Reinforcement Learning with Human Feedback (RLHF)
- Positional Encoding