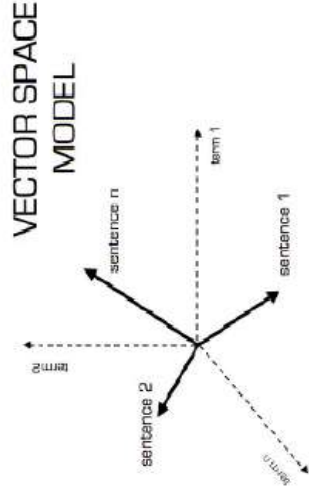


NLP

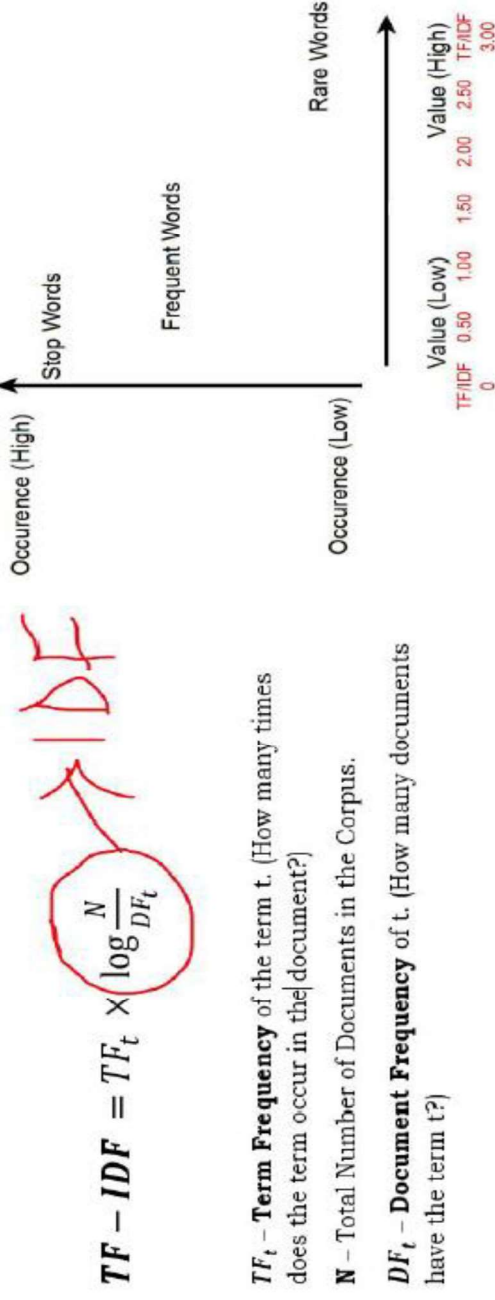
- Text representation

		Document 1		Document 2	
		Document 1		Document 2	
		Term			
Document 1	The quick brown fox jumped over the lazy dog's back.	aid		0 1	
		all		0 1	
		back		1 0	
		brown		1 0	
		come		0 1	
		dog		1 0	
		fox		1 0	
		good		0 1	
		jump		1 0	
		lazy		1 0	
		men		0 1	
		now		0 1	
		over		1 0	
		party		0 1	
Document 2	Now is the time for all good men to come to the aid of their party.	aid		1 0	
		all		1 0	
		back		0 1	
		brown		0 1	
		come		1 0	
		dog		0 1	
		fox		0 1	
		good		1 0	
		jump		0 1	
		lazy		0 1	
		men		1 0	
		now		0 1	
		over		1 0	
		party		1 0	
		quick		1 0	
		their		0 1	
		time		0 1	



Text representation

- TF-IDF - term frequency - inverse document frequency

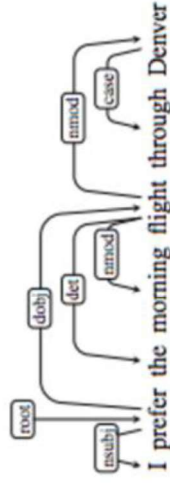
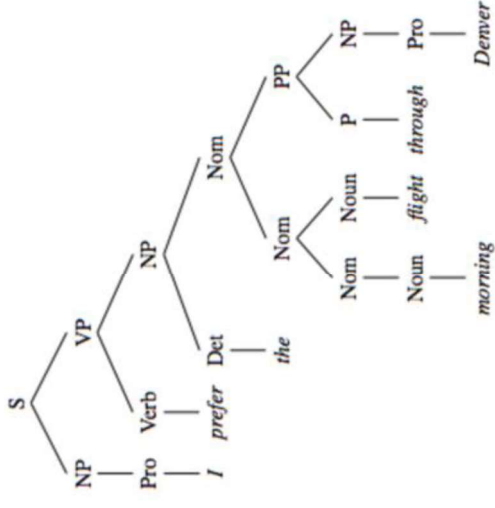
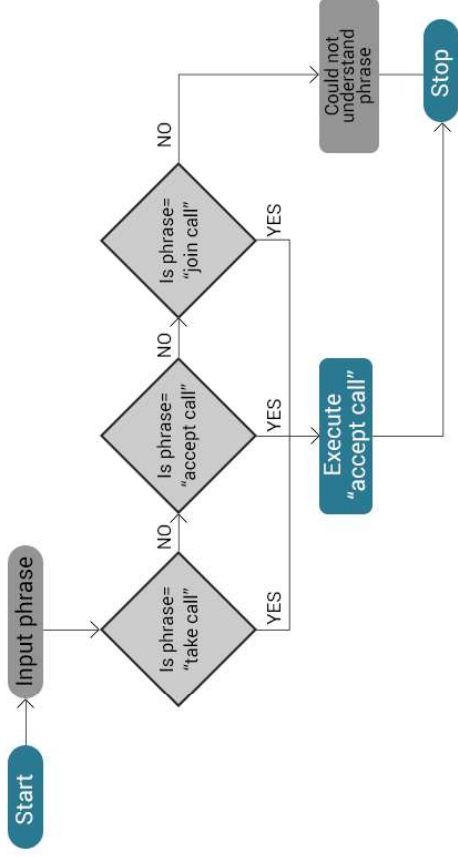


Linguistic Foundations

85

- ▶ Rule-based approaches
- ▶ Semantic parsing
- ▶ Analyzing linguistic structure and grammars of text

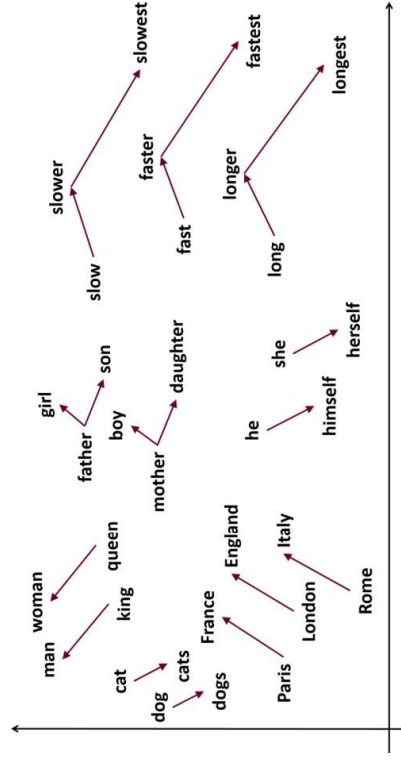
SIMPLE RULE BASED RULE



Word Embeddings

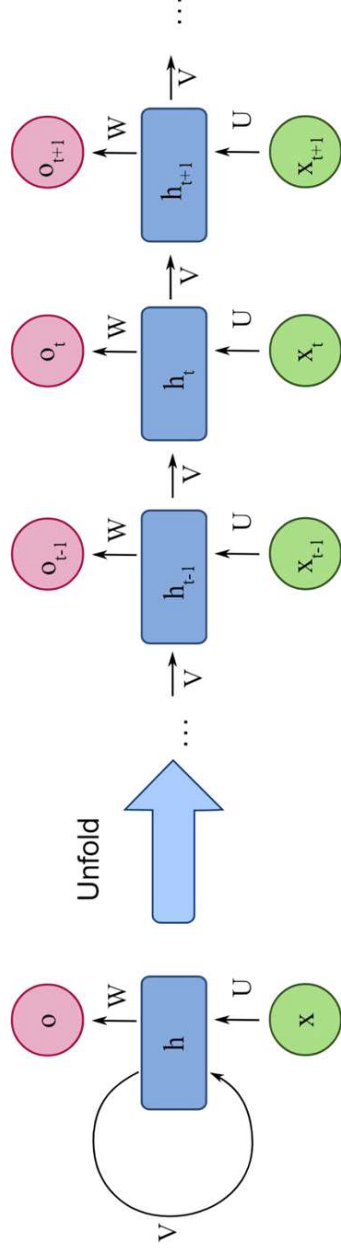
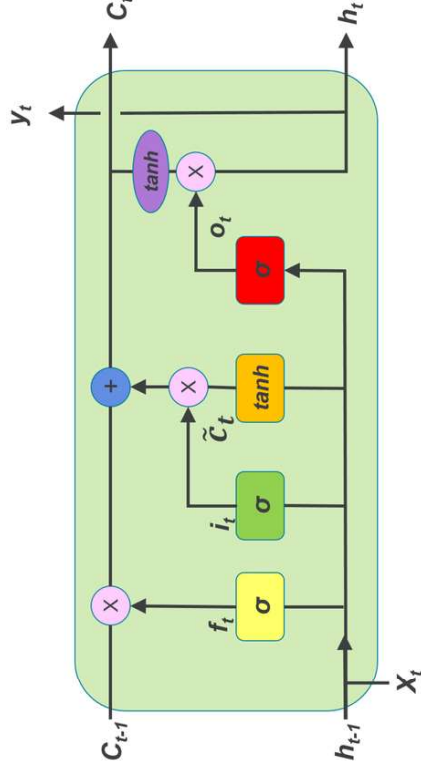
86

- ▶ Represent each word as a “vector” of numbers
- ▶ Converts a “discrete” representation to “continuous”, allowing for:
 - ▶ More “fine-grained” representations of words
 - ▶ Useful computations such as cosine/eucl distance
 - ▶ Visualization and mapping of words onto a semantic space
- ▶ Examples:
 - ▶ Word2Vec (2013), GloVe, BERT, ELMo



Seq2seq Models

- ▶ Recurrent Neural Networks (RNNs)
- ▶ Long Short-Term Memory Networks (LSTMs)
- ▶ “Dependency” and info between tokens
- ▶ Gates to “control memory” and flow of information

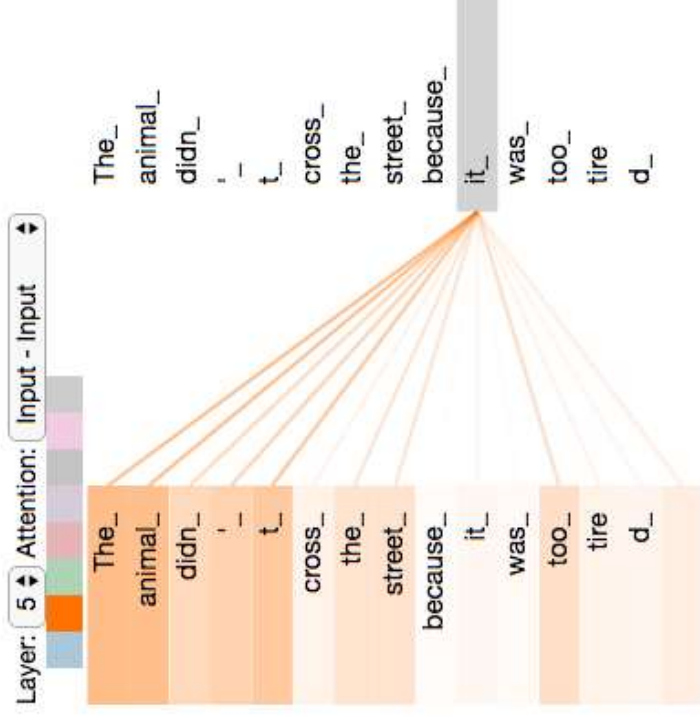


Attention and Transformers

- ▶ Allows to “focus attention” on particular aspects of the input text
- ▶ Done by using a set of parameters, called “weights,” that determine how much attention should be paid to each input at each time step
- ▶ These weights are computed using a combination of the input and the current hidden state of the model
- ▶ Attention weights are computed (dot product of the query, key and value matrix), then a softmax function is applied to the dot product

$$\text{attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

88



<https://arxiv.org/abs/1706.03762>

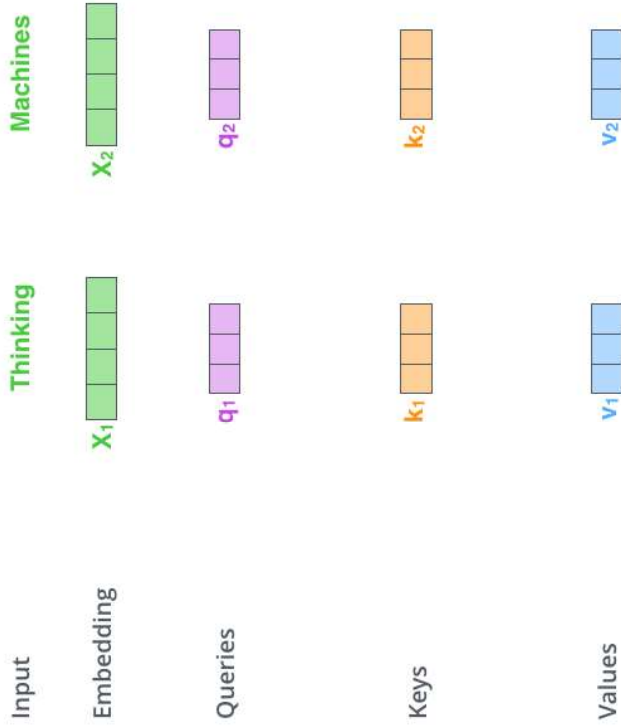
<https://jalamar.github.io/illustrated-transformer/>

Analogy for Q, K, V

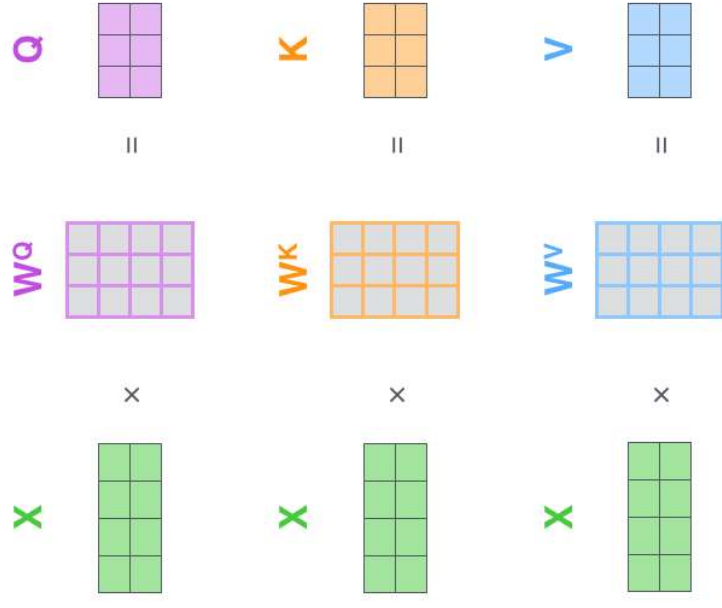
- ▶ Library system
- ▶ Imagine you're looking for information on a specific topic (query)
- ▶ Each book in the library has a summary (key) that helps identify if it contains the information you're looking for
- ▶ Once you find a match between your query and a summary, you access the book to get the detailed information (value) you need
- ▶ Here, in Attention, we do a "soft match" across multiple values, e.g. get info from multiple books ("book 1 is most relevant, then book 2, then book 3, etc.")

$$\text{attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Self-Attention



91



<https://ialammar.github.io/illustrated-transformer/>

Transformer & Multi-Head Attention

92

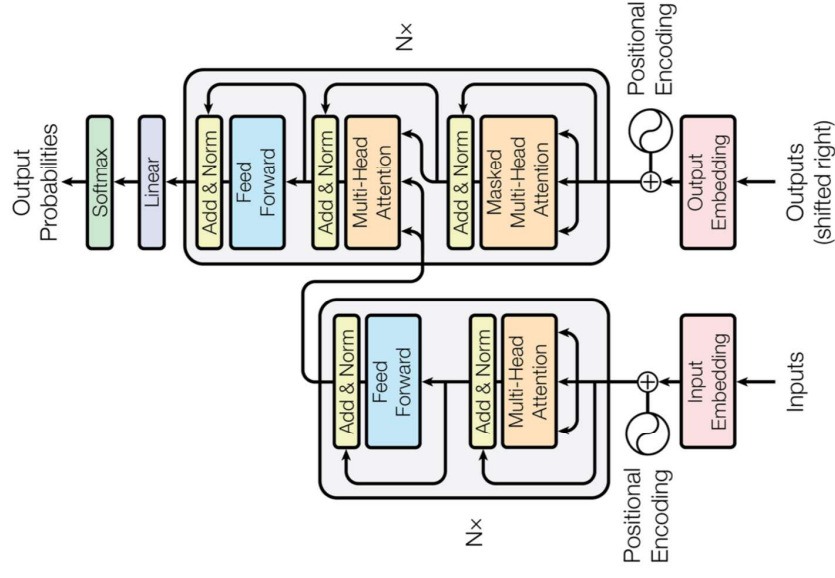
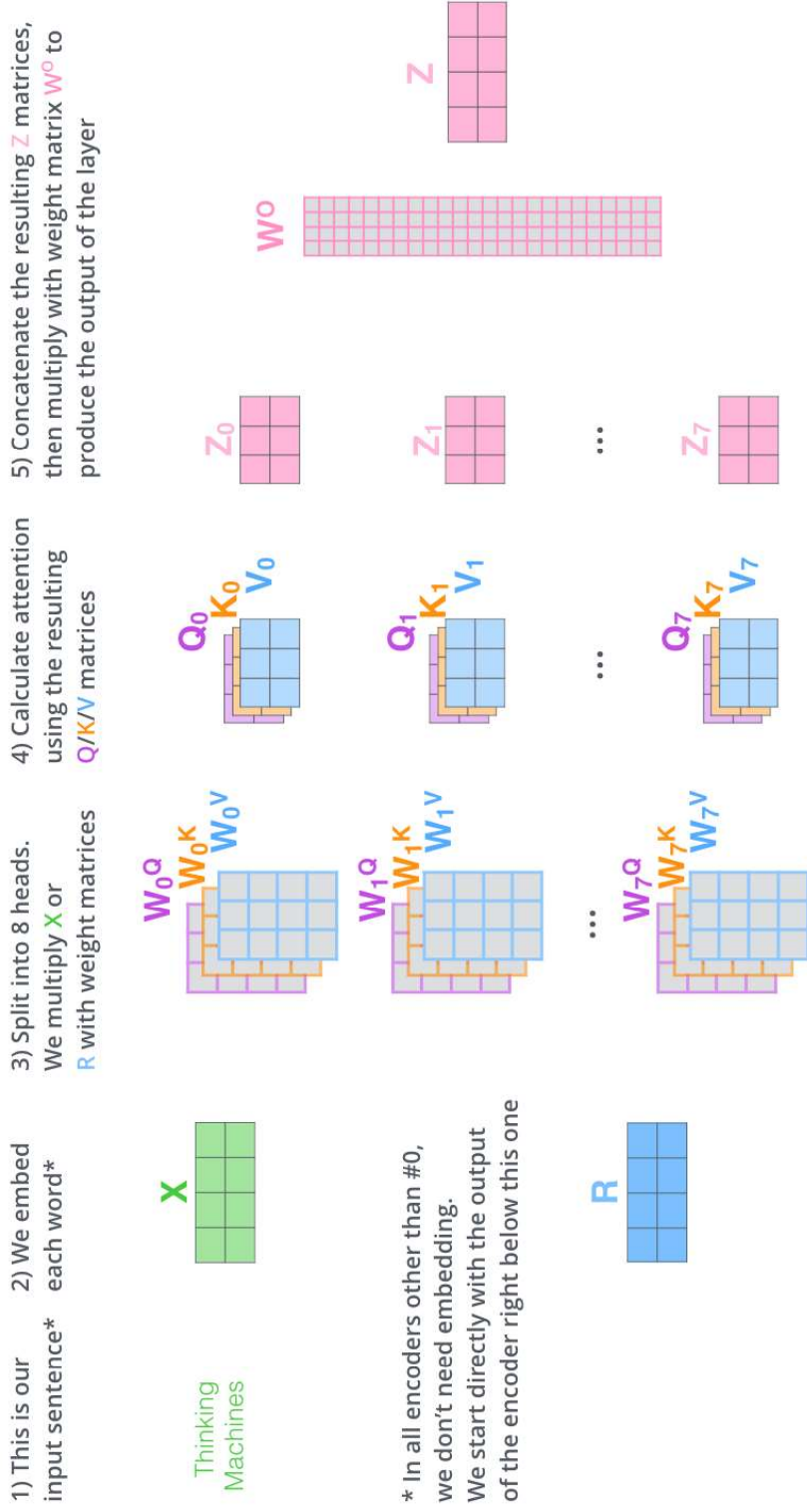


Figure 1: The Transformer - model architecture.

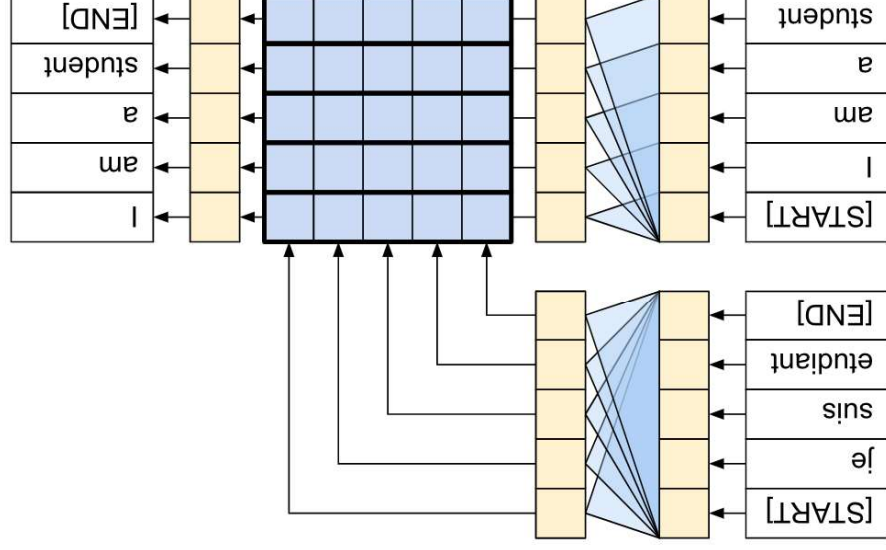
“Attention Is All You Need”
<https://arxiv.org/abs/1706.03762>

Multi-Head Attention

93



Cross-Attention (e.g. Machine Translation)⁹⁴



Transformers vs. RNNs

95

Challenges with RNNs	Transformers
<ul style="list-style-type: none">• Long range dependencies• Gradient vanishing and explosion• Large # of training steps• Sequential/recurrence → can't parallelize• Complexity per layer: $O(n \cdot d^2)$	<ul style="list-style-type: none">• Can model long-range dependencies• No gradient vanishing and explosion• Fewer training steps• Can parallelize computation!• Complexity per layer: $O(n^2 \cdot d)$

Large Language Models

96

- ▶ Scaled up versions of Transformer architecture, e.g. millions/billions of parameters
- ▶ Typically trained on massive amounts of “general” textual data (e.g. web corpus)
- ▶ Training objective is typically “next token prediction”: $P(W_{t+1} | W_t, W_{t-1}, \dots, W_1)$
- ▶ Emergent abilities as they scale up (e.g. chain-of-thought reasoning)
- ▶ Heavy computational cost (time, money, GPUs)
- ▶ Larger general ones: “plug-and-play” with few or zero-shot learning
 - ▶ Train once, then adapt to other tasks without needing to retrain
 - ▶ E.g. in-context learning and prompting

Emergent Abilities of Large Language Models⁹⁷

- ▶ Why do LLMs work so well? What happens as you scale up?
- ▶ Potential explanation: emergent abilities!
- ▶ An ability is emergent if it is present in larger but not smaller models
- ▶ Not have been directly predicted by extrapolating from smaller models
- ▶ Performance is near-random until a certain critical threshold, then improves heavily
 - ▶ Known as a “phase transition” and would not have been extrapolated

Wei et al., 2022. <https://arxiv.org/abs/2206.07682>