# Intro to Neural Nets

Attention & Transformers

# Today's Agenda

**Attention Mechanism in RNNs for Translation (2014)**
- Sequence to sequence models with attention in translation tasks.

**More General, Approach (No RNNs Required)**
- We can hard-code a measure of semantic similarity between a focal token, and all other tokens in the sequence

**QKV Framework & Multi-head Attention**
- We can incorporate trainable weights around the attention mechanism.

# TF-IDF is "Static Attention"

**Not all phrases are of equal importance...**

- E.g., David less important than Beckham
- If a term occurs all the time, observing its presence is less informative

**Inverse-document frequency (IDF) helps address this.**

$$\text{IDF} = \log(N/n_j)$$

- Term 'weighting' is then calculated as Term Frequency (TF) x IDF
- $n_j$= # of docs containing the term, N = total # of docs
- A term is deemed important if it has a high TF and/or a high IDF.
- As TF goes up, the word is more common generally. As IDF goes up, it means very few documents contain this term.
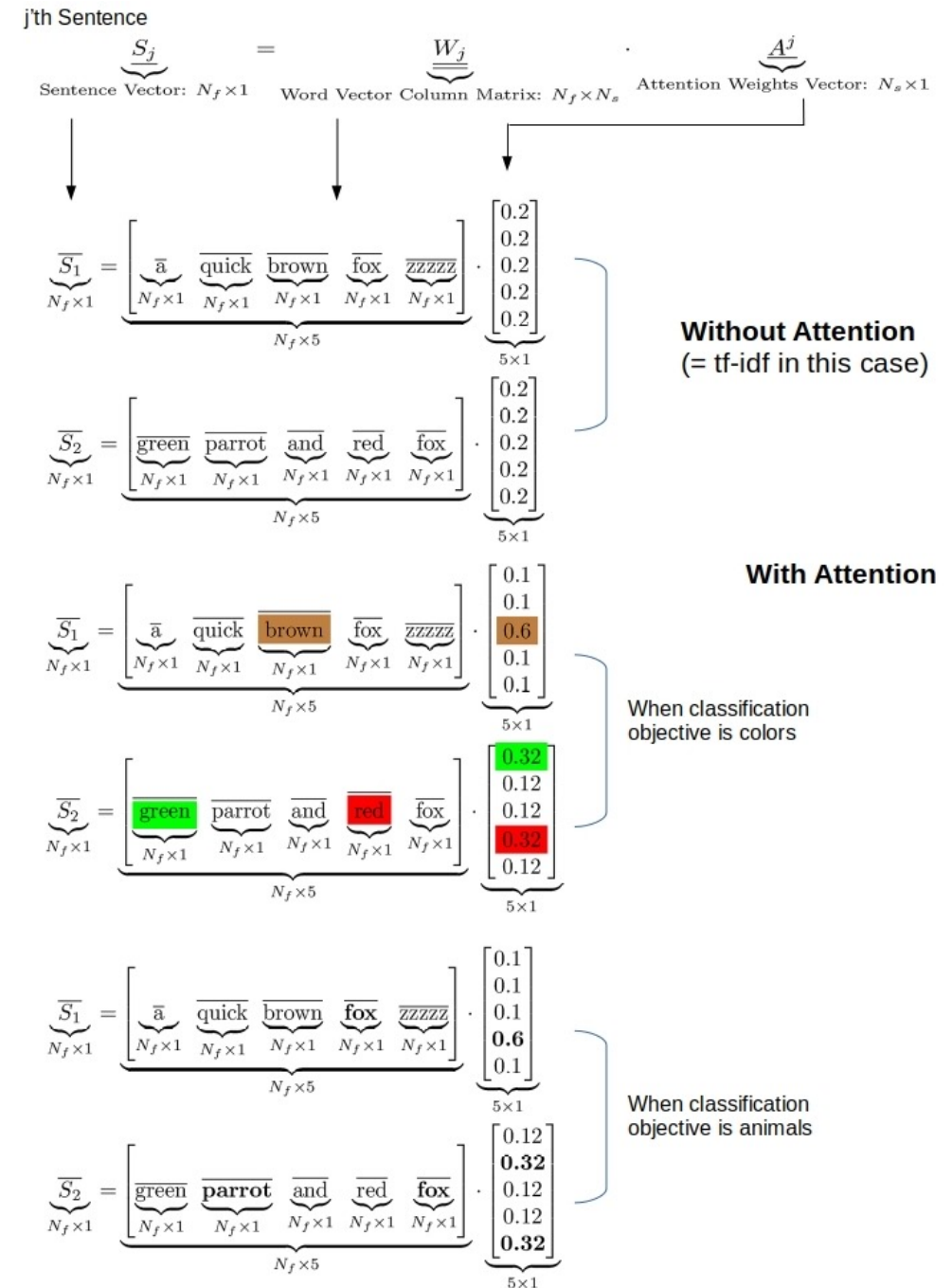
# What Would Adaptive Weights Look Like?

**Figure Out Weights Useful for Prediction**

- Rather than telling the network the weights, which is what TF-IDF effectively does, we provide the network a weight matrix to update!
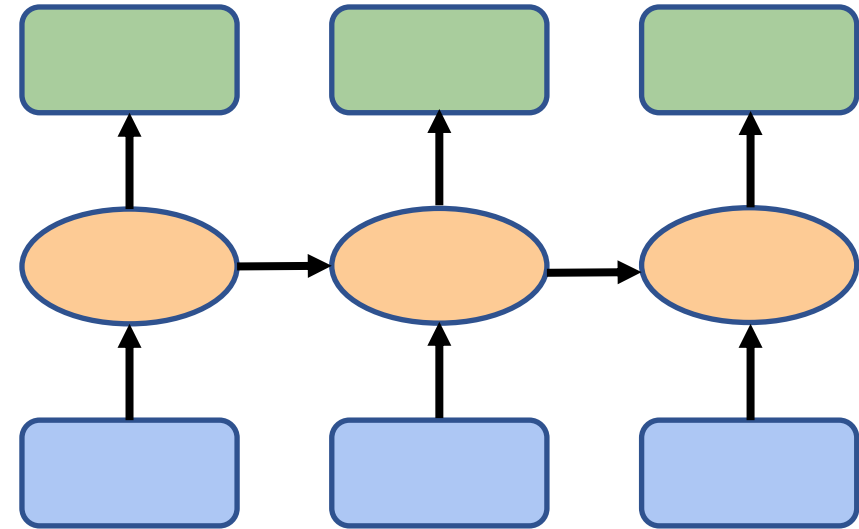
**This Yields Task-specific Token Weights**

- See some simple examples on the right.
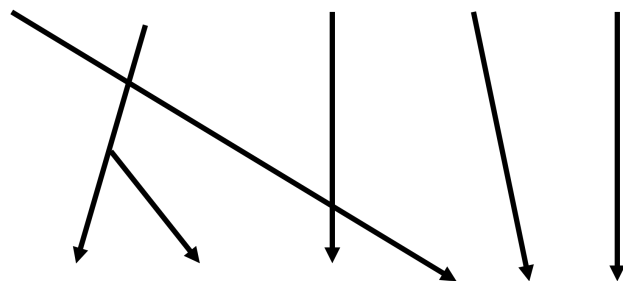- Our network may learn that certain tokens are more important for predicting a given label.

# Where Attention Came From: Sequence-to-Sequence Models
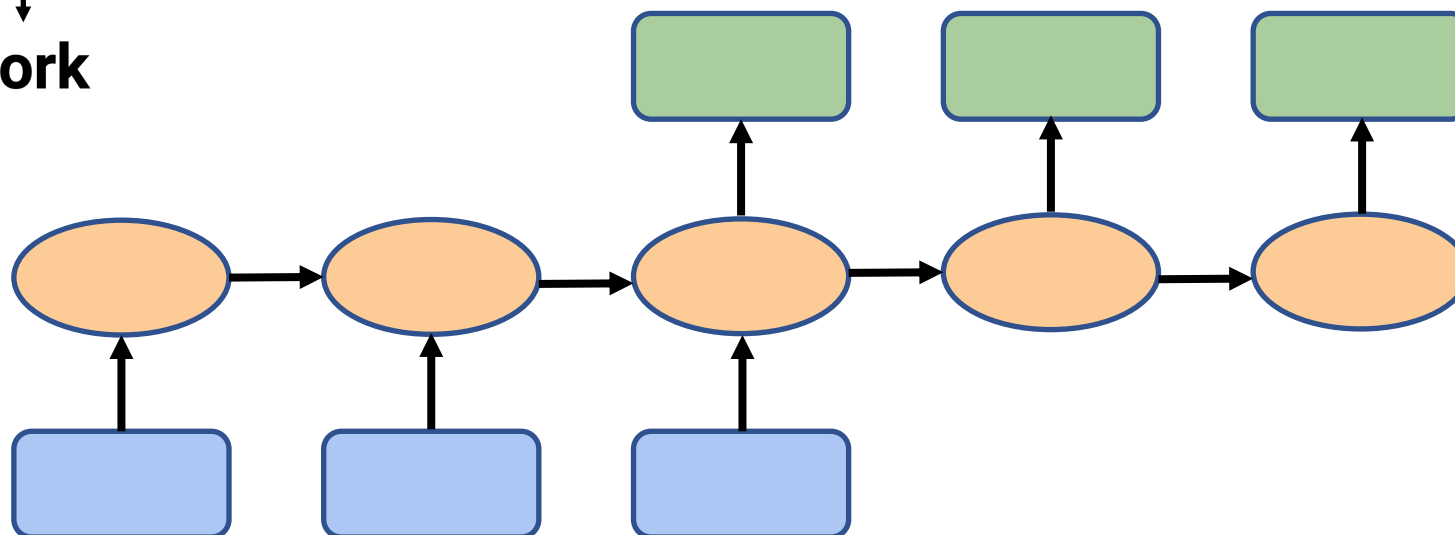
**Yo soy de Canada**

**I am from Canada**

# Struggling to Remember Long Sequences

# The First (RNN) Attention Model

## Neural Machine Translation by Jointly Learning to Align and Translate

**Dzmitry Bahdanau**
Jacobs University Bremen, Germany

**KyungHyun Cho**    **Yoshua Bengio**[*]
Université de Montréal

*" … allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word …"*
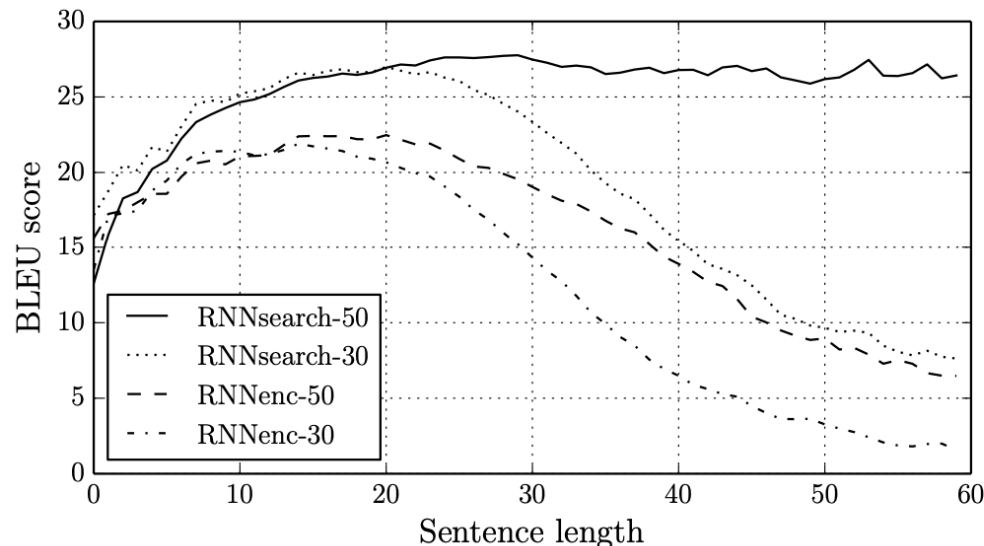


Figure 2: The BLEU scores of the generated translations on the test set with respect to the lengths of the sentences. The results are on the full test set which includes sentences having unknown words to the models.

# The First (RNN) Attention Model

**Two RNNs**
- A bidirectional RNN handling the input sentence (encoder).
- A regular RNN handling the output (decoder).

**Two RNNs**
- Bidirectional RNN reads sentence in both directions.
- Produces a hidden state output for each token of input.
- Attention weights are calculated too, per token, based on the entire input sequence.
- Finally, the attention-weighted output for each token is passed into the regular RNN (which does the final translation).
- Each output from that last RNN is based on attention-weighted input token and output hidden state of last translated token (i.e., s_t-1)
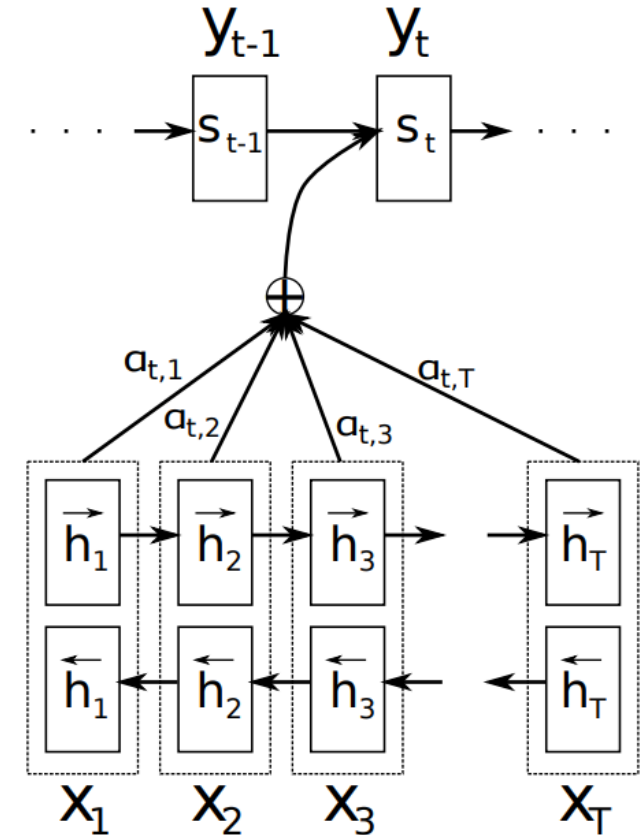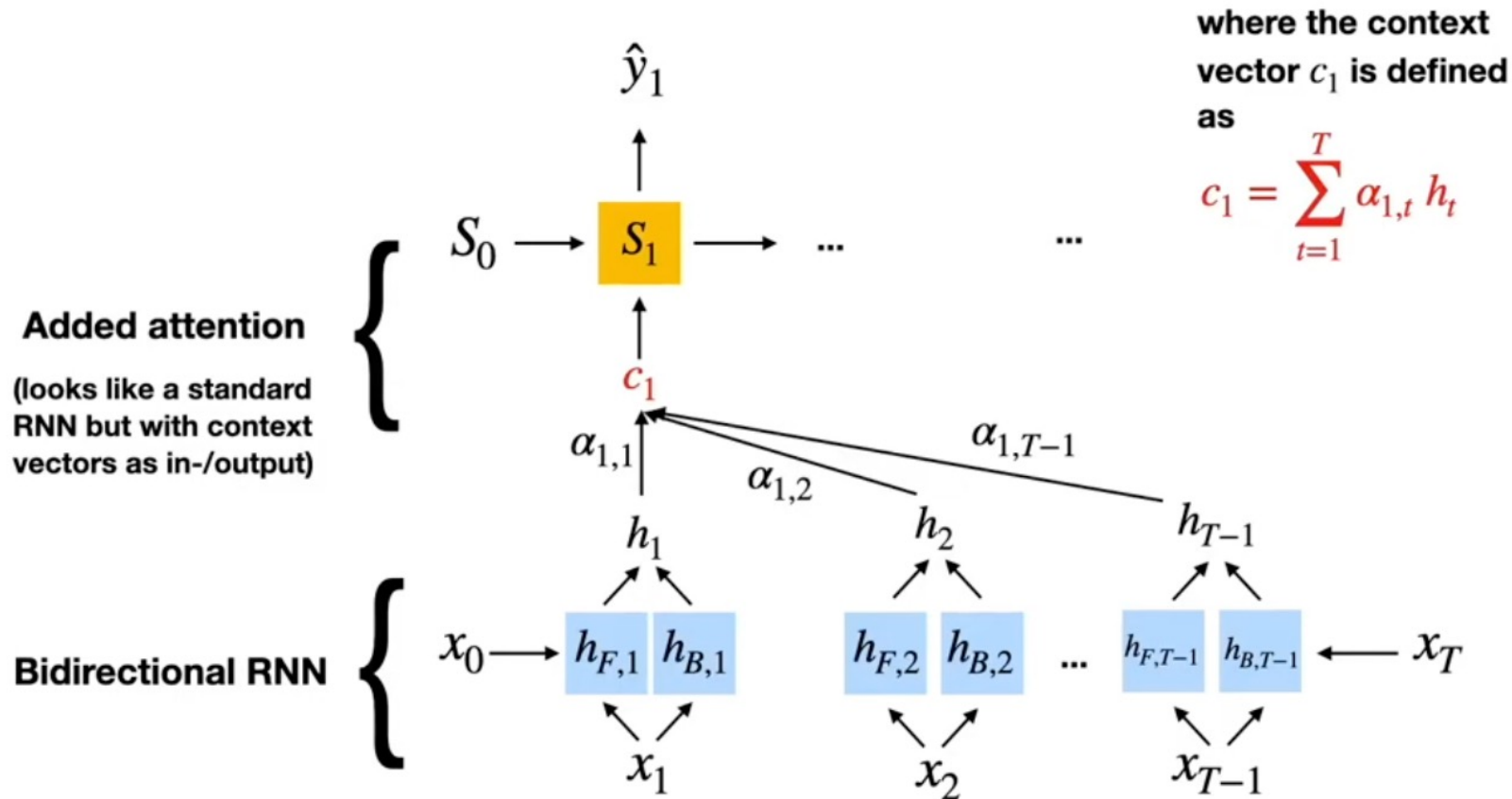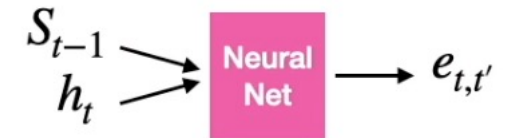


Figure 1: The graphical illustration of the proposed model trying to generate the $t$-th target word $y_t$ given a source sentence $(x_1, x_2, \ldots, x_T)$.

# The First (RNN) Attention Model



where the context vector $c_1$ is defined as

$$c_1 = \sum_{t=1}^{T} \alpha_{1,t}\, h_t$$

**Computing attention weights**

$$\alpha_{t,t'} = \frac{\exp(e_{t,t'})}{\sum_{t'=1}^{T} \exp(e_{t,t'})}$$

**Added attention**

(looks like a standard RNN but with context vectors as in-/output)

**Bidirectional RNN**

# The First (RNN) Attention Model

# We Don't Need an RNN to Do This

**We Can Drop the RNN and Process the Whole Sequence in Parallel**

- Note that transformers work with a similar autoencoder architecture, but no LSTMs needed (we use "stacked attention layers").

- Transformers are actually more computationally expensive (more calculations required), but because we can parallelize them they can be trained faster than RNNs.

- By moving away from RNN specific implementation we can also think about more general forms of attention (e.g.,. For CNNs).

---

**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

# A Simpler Attention Mechanism

**Main procedure:**

1) Derive attention <u>weights</u>: similarity between current input and all other inputs (next slide)

2) Normalize weights via softmax (next slide)

3) Compute attention value from normalized weights and corresponding inputs (below)

**Self-attention as weighted sum:**

$$A_i = \sum_{j=1}^{T} a_{ij} x_j$$

output corresponding to the i-th input          weight based on similarity between current input $x_i$ and all other inputs

# A Simple Attention Mechanism

**Self-attention as weighted sum:**

$$A_i = \sum_{j=0}^{T} a_{ij} x_j$$

output corresponding to the i-th input

weight based on similarity between current input $x_i$ and all other inputs

**How to compute the attention weights?**

here as simple dot product:

$$e_{ij} = x_i^\top x_j$$

repeat this for all inputs $j \in \{1...T\}$, then normalize

$$a_{ij} = \frac{\exp\left(e_{ij}\right)}{\sum_{j=1}^{T} \exp\left(e_{ij}\right)} = \text{softmax}\left(\left[e_{ij}\right]_{j=1....T}\right)$$
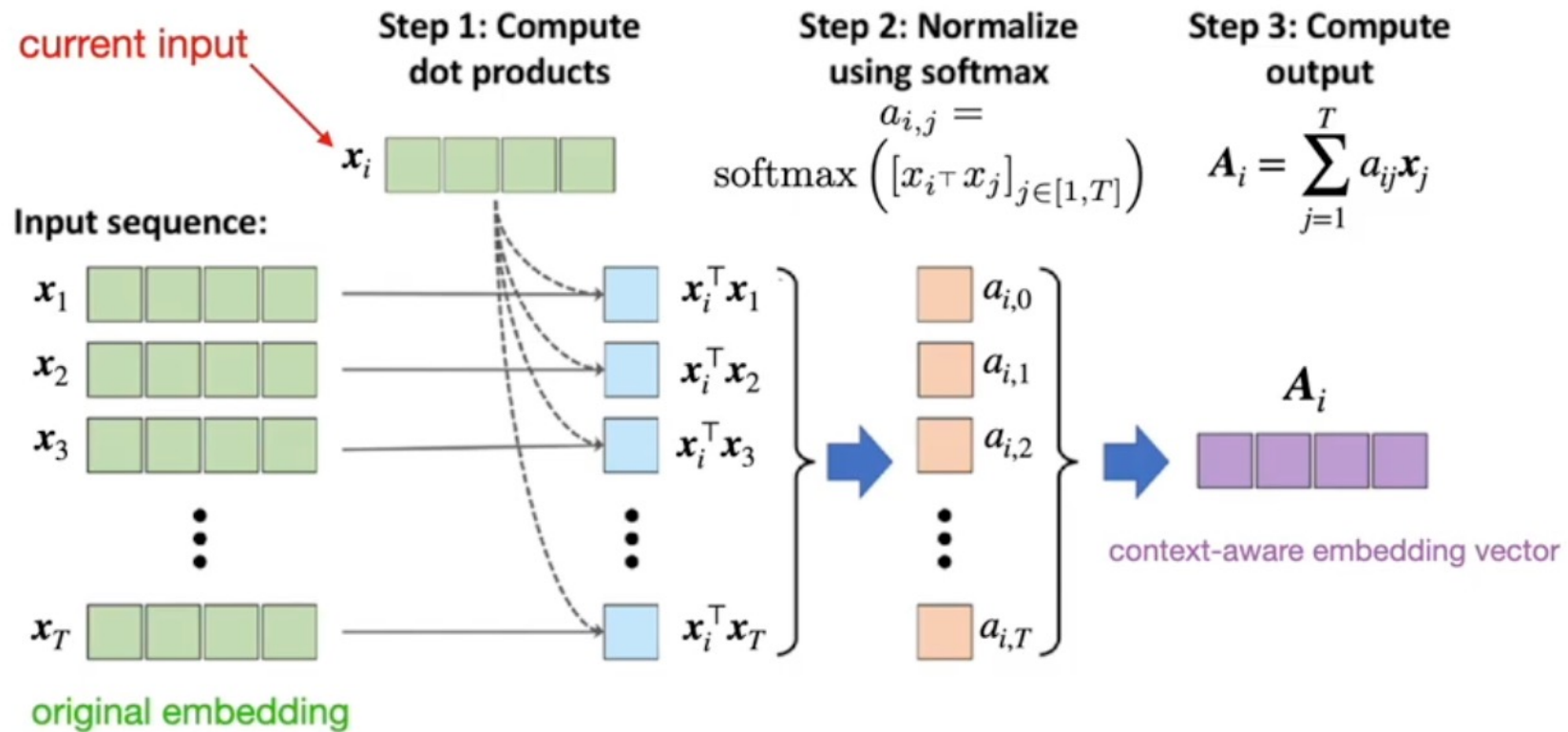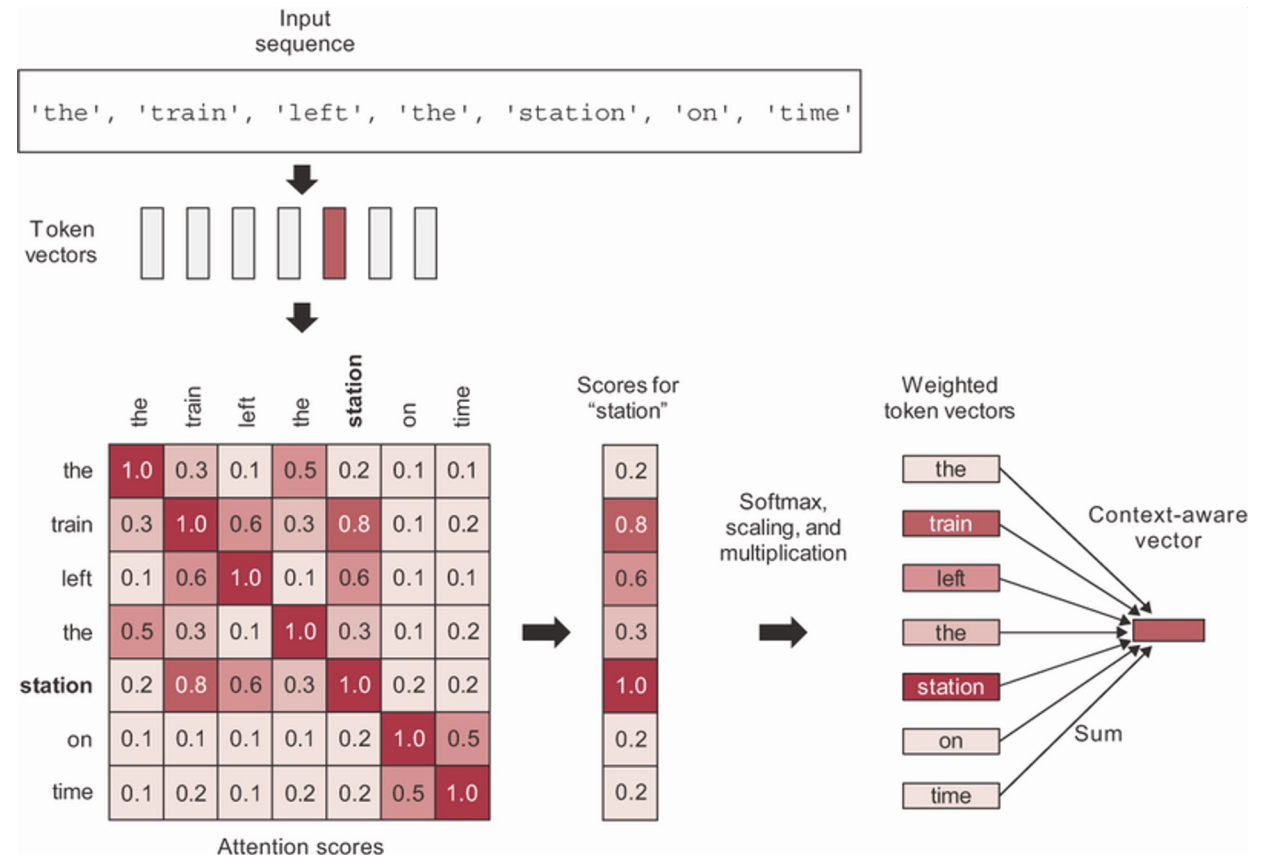
# A Simple Attention Mechanism

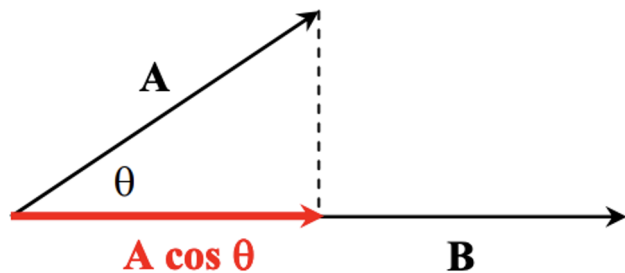

Image source: Raschka & Mirjalili 2019. Python Machine Learning, 3rd edition

# A Basic Self-Attention Mechanism

**Most Basic Form of Self-Attention**

- This is what the textbook introduces.
- We magnify vector representations of tokens based on how parallel (semantically related) they are to the other tokens in the input sequence.
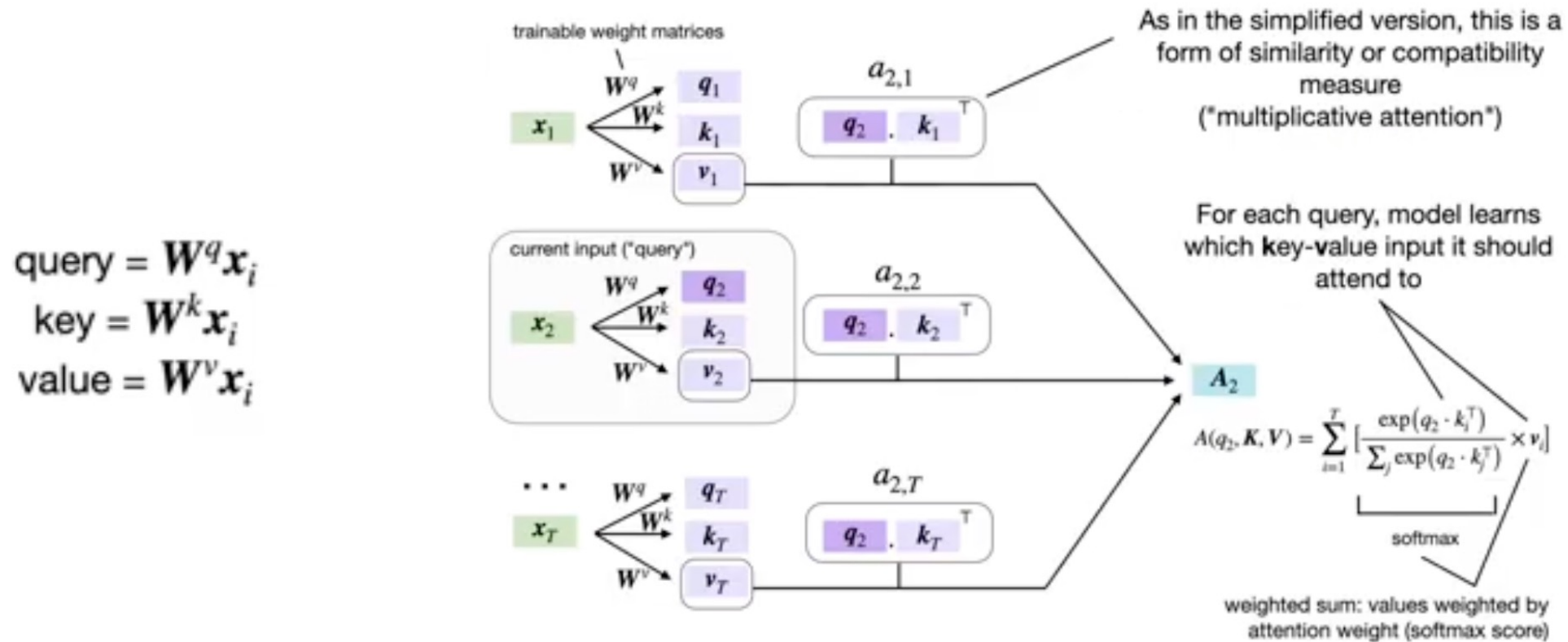
# The Basic Attention Mechanism is Rigid

**The Approach Considers the Input Context in a Very 'Fixed' Manner**

- We get context-specific attention weights, but the notion of context is pre-defined / hard-coded in the dot-product-summation procedure.
- The pre-defined approach to attention does not generalize outside of language.
- The pre-defined approach to attention assumes that shared semantic similarity is the best heuristic for attention, but perhaps it is not in some cases. Most obviously, perhaps in some scenarios it is ideal to focus attention on words that are semantically dissimilar.
- How can we enable a procedure of this sort, but using trainable parameters, so the network can figure out how to allocate attention by itself?

# Scaled Dot-Product Attention Mechanism



$$query = W^q x_i$$
$$key = W^k x_i$$
$$value = W^v x_i$$

trainable weight matrices

As in the simplified version, this is a form of similarity or compatibility measure ("multiplicative attention")

current input ("query")

For each query, model learns which key-value input it should attend to

$$A(q_2, K, V) = \sum_{i=1}^{T} \left[ \frac{\exp(q_2 \cdot k_i^T)}{\sum_j \exp(q_2 \cdot k_j^T)} \times v_i \right]$$

softmax

weighted sum: values weighted by attention weight (softmax score)

# Flexible Contextual Attention Mechanism

**Scaled Dot-Product Attention...**

- One other nuance is that we scale the Softmax output.
- This is done to prevent the output being too close to 1 or too close to 0, which can lead to vanishing gradient issues.
- Here, sqrt(d_k) is square root of the dimensionality of the keys (in the original paper this is the same as the dimension of the word embeddings).
- We write this in big-letter notation because we can parallelize everything as very large matrix / tensor calculations.
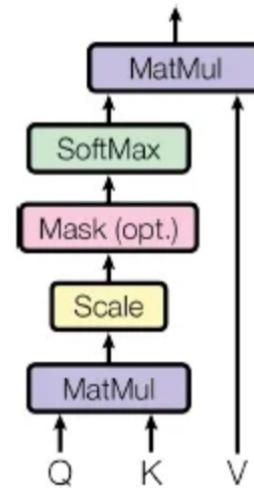
$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$
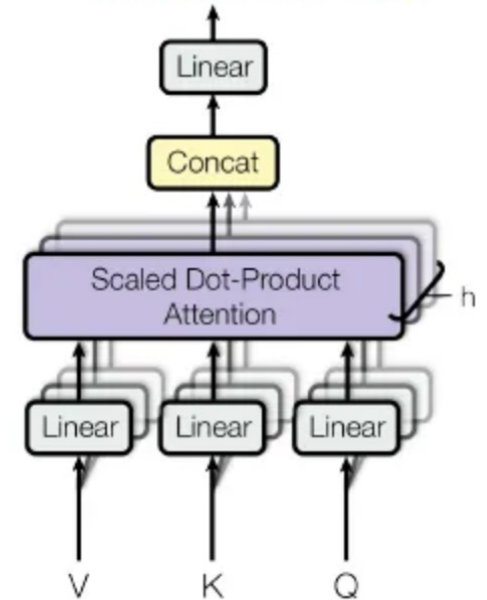
# Multi-Head Attention

**We have Multiple 'Heads' Paying Attention**

- Each head can be thinking about something very different.
- By using multiple sets of query, key and value weight matrices, the network can learn to consider different combinations of input features.
- This lets it capture more complex patterns in the data.



From "Attention is all you need" paper by Vaswani, et al., 2017 [1]

# Questions?