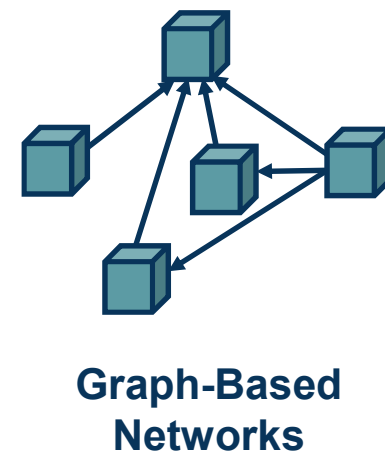
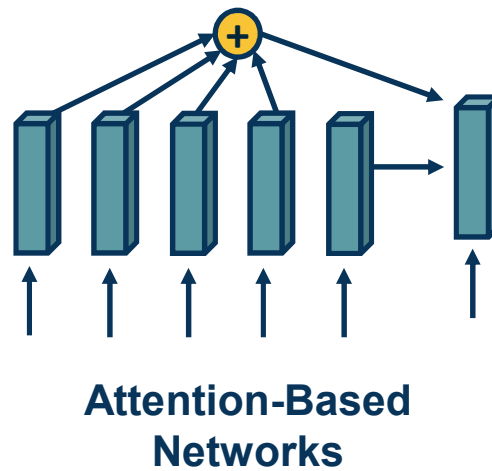
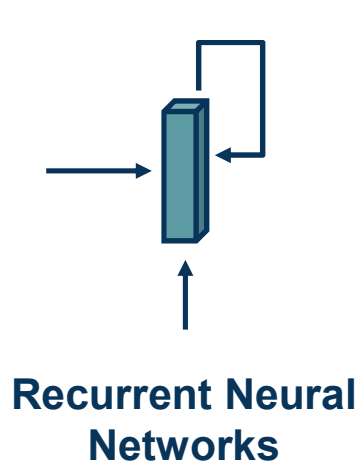
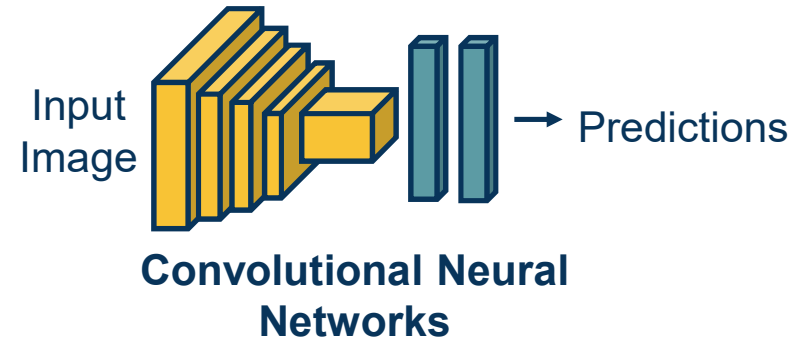
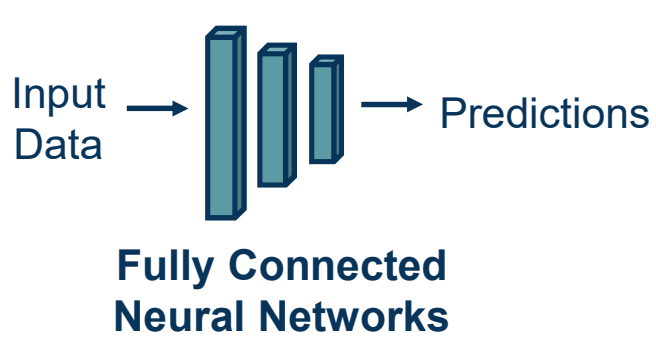
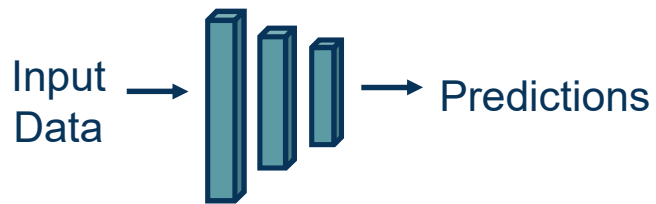


Module 3

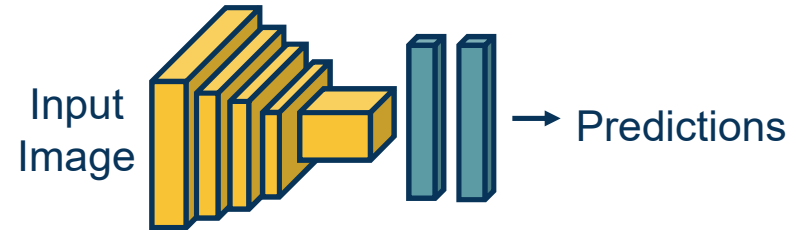
Introduction



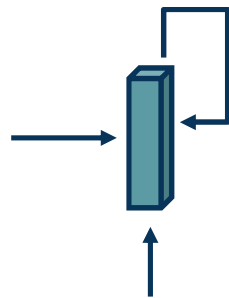
The Space of Architectures



**Fully Connected
Neural Networks**

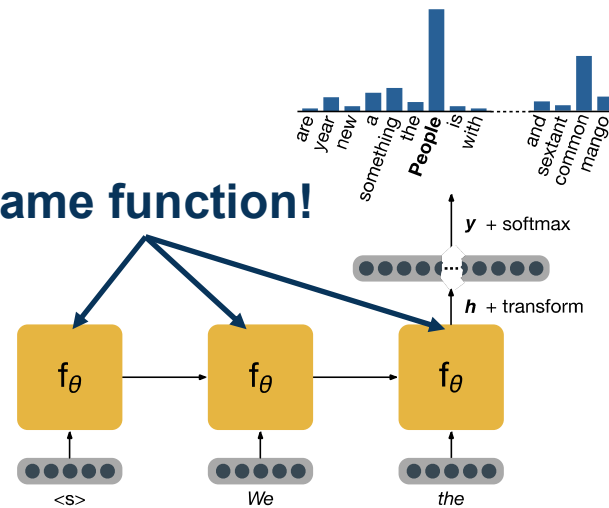


**Convolutional Neural
Networks**

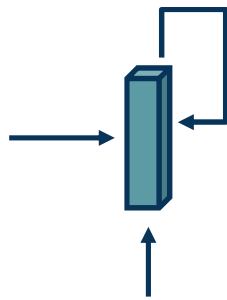


**Recurrent Neural
Networks**

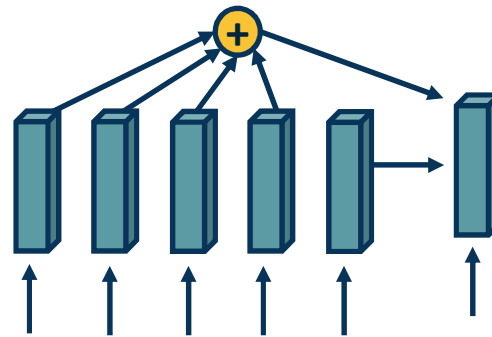
Same function!



Recurrent Neural Networks



**Recurrent Neural
Networks**



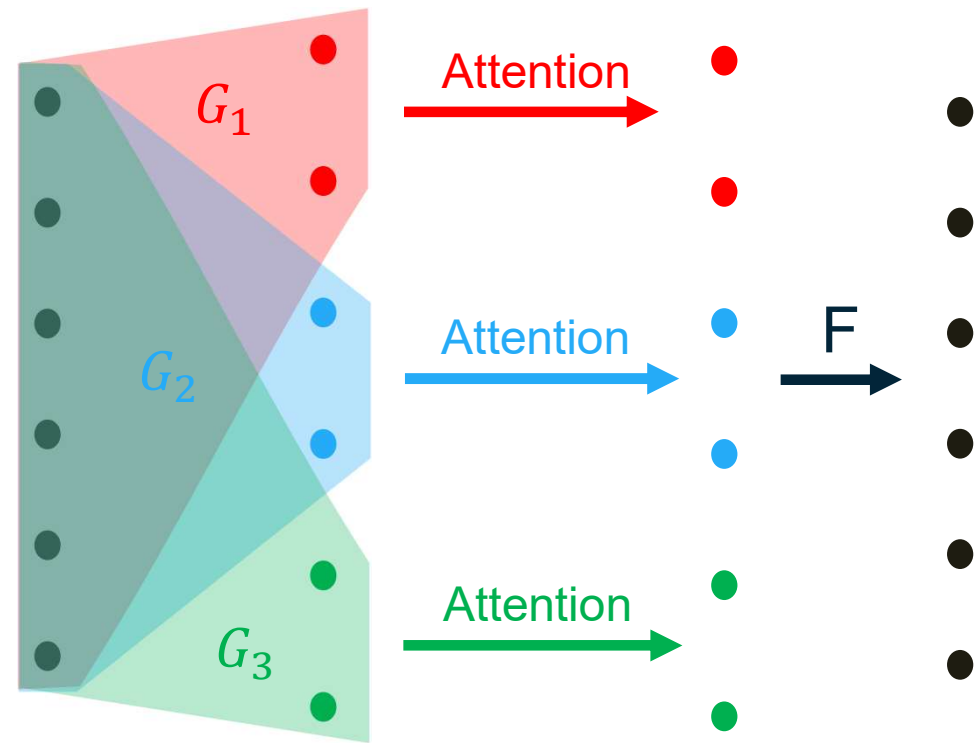
**Attention-Based
Networks**

The Space of Architectures

Transformer [Vaswani et. al. 2017] is a multi-layer attention model that is currently state of the art in most language tasks (and in many other things!)

Has superior performance compared to previous attention based architectures via

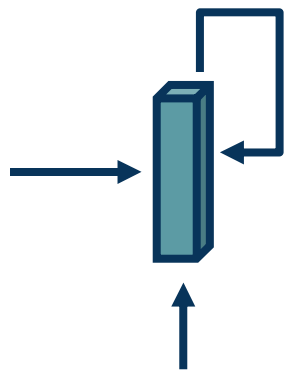
- Multi-query hidden-state propagation (“self-attention”)
- Multi-head attention**
- Residual connections, LayerNorm



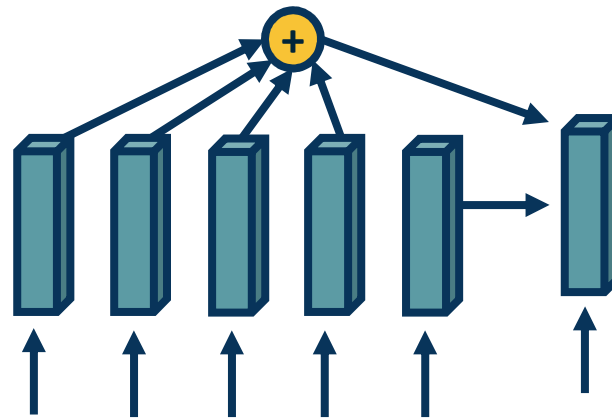
Enter the Transformer

FACEBOOK AI

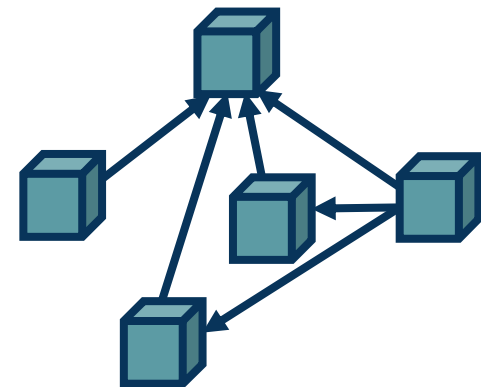




**Recurrent
Neural Networks**



**Attention-Based
Networks**



**Graph-Based
Networks**

The Space of Architectures

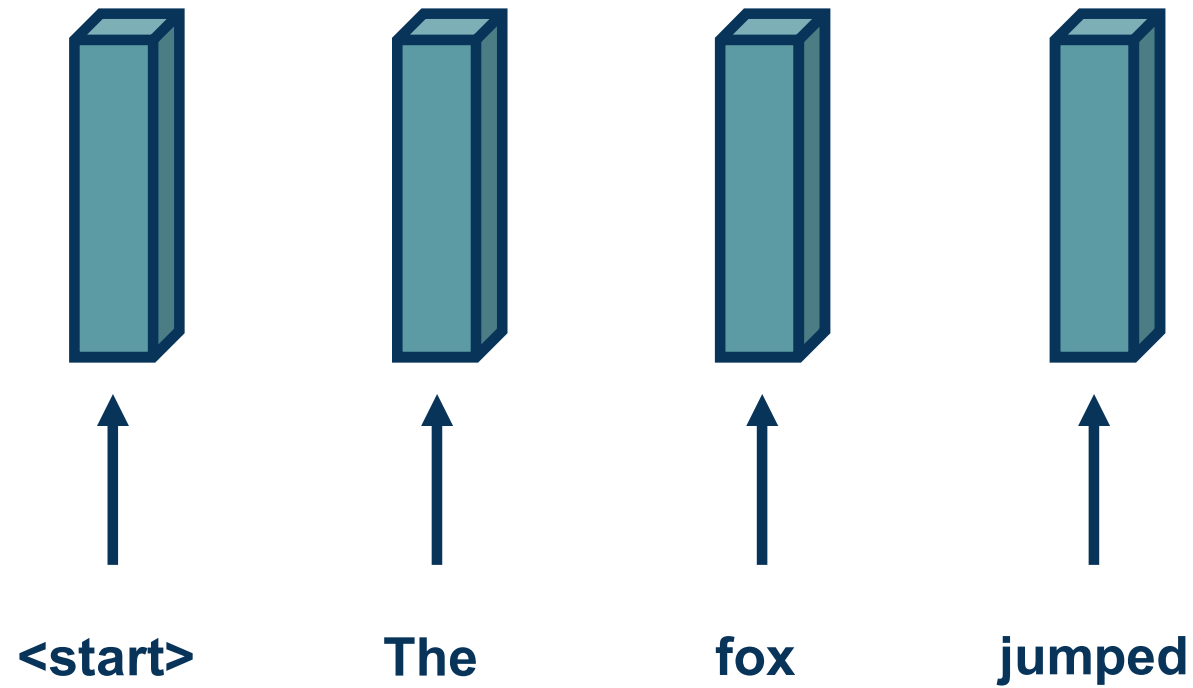
- ◆ **Many** → **many**: speech recognition, optical character recognition



- ◆ **Many** → **one**: sentiment analysis, topic classification



- ◆ Also consider: **one** → **many**, **one** → **one**.



Example Application: NLP

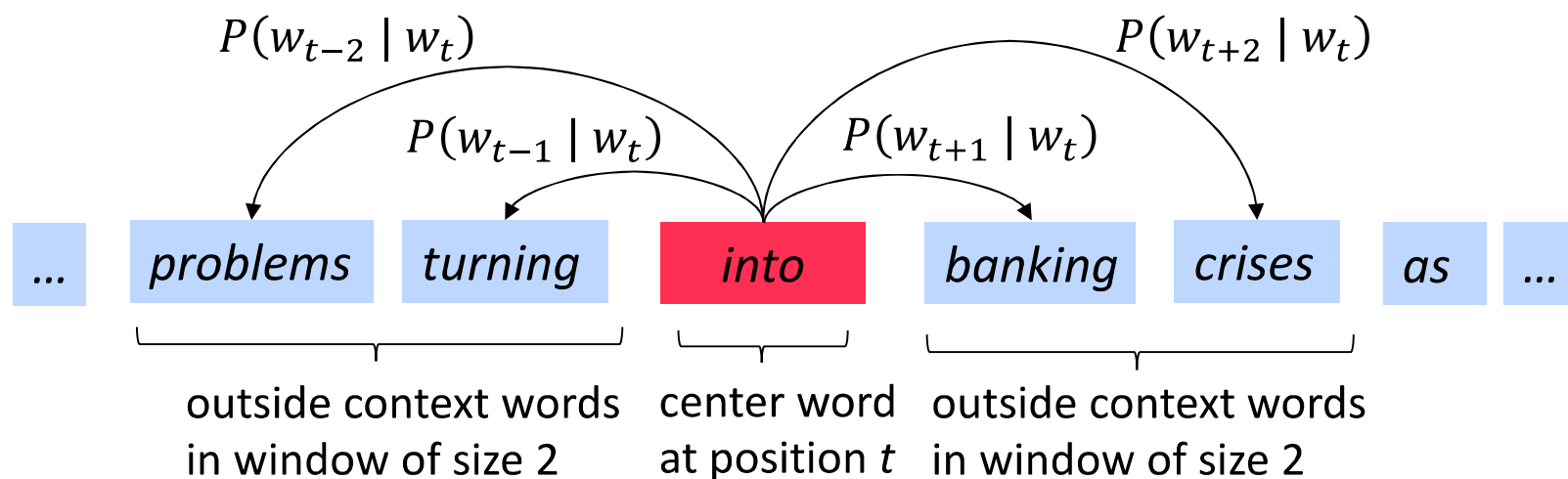
$$p(\mathbf{s}) = p(w_1, w_2, \dots, w_n)$$

$$= p(w_1) p(w_2 \mid w_1) p(w_3 \mid w_1, w_2) \cdots p(w_n \mid w_{n-1}, \dots, w_1)$$

$$= \prod_i p(\text{next word } w_i \mid \text{history } w_{i-1}, \dots, w_1)$$

Word2vec: the Skip-gram model

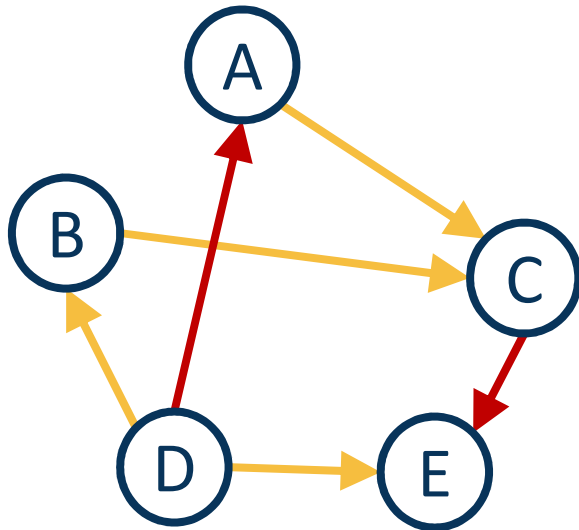
- ✦ The idea: use words to **predict** their context words
- ✦ Context: a fixed window of size **$2m$**



Slide Credit: Richard Socher, Christopher Manning



A Sequential Structure



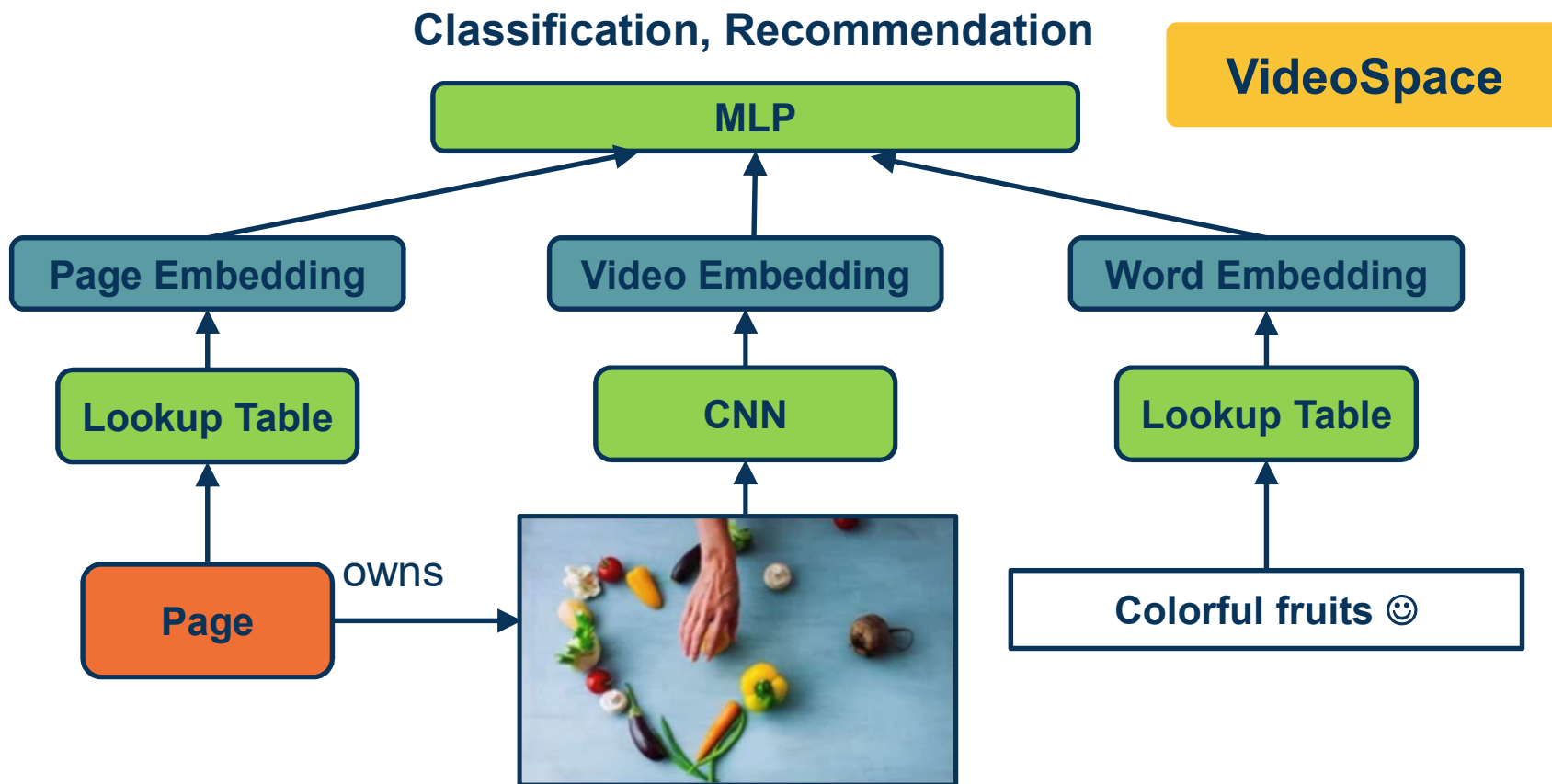
A Multi-Relation Graph

Embedding: A learned map from entities to vectors of numbers that encodes similarity

- Word embeddings: word → vector
- Graph embeddings: node → vector

Graph Embedding: Optimize the objective that **connected nodes have more similar embeddings** than unconnected nodes via gradient descent.

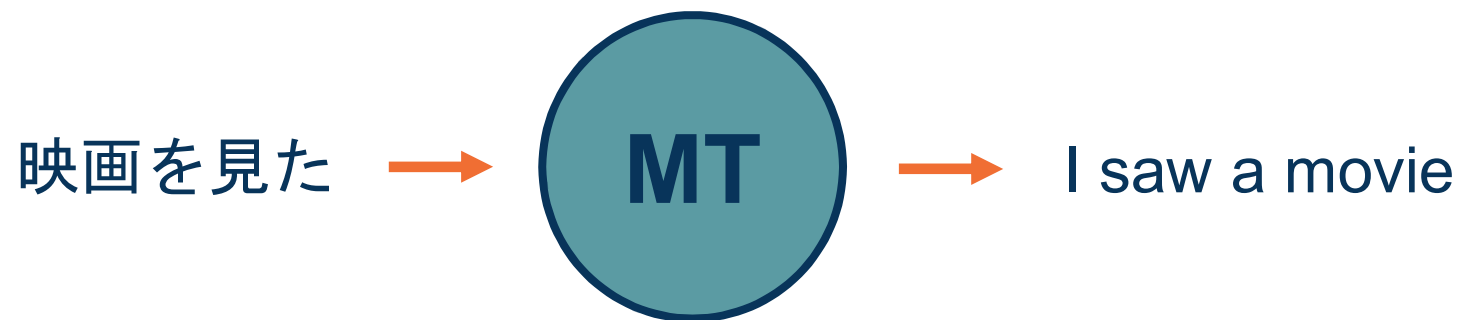
Slide Credit: Adam Lerer



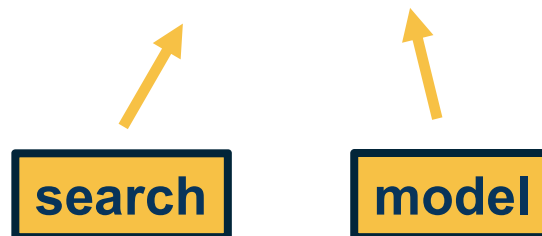
Application: VideoSpace

FACEBOOK AI





$$t = \operatorname{argmax}_t p(t|s)$$



Alignment in machine translation: for each word in the target, get a distribution over words in the source [Brown et. al. 1993], (lots more)

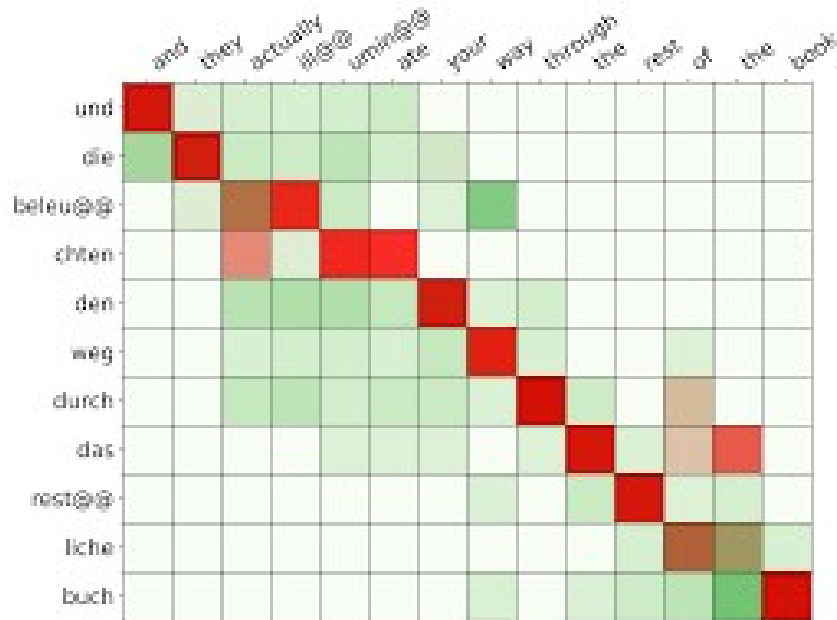
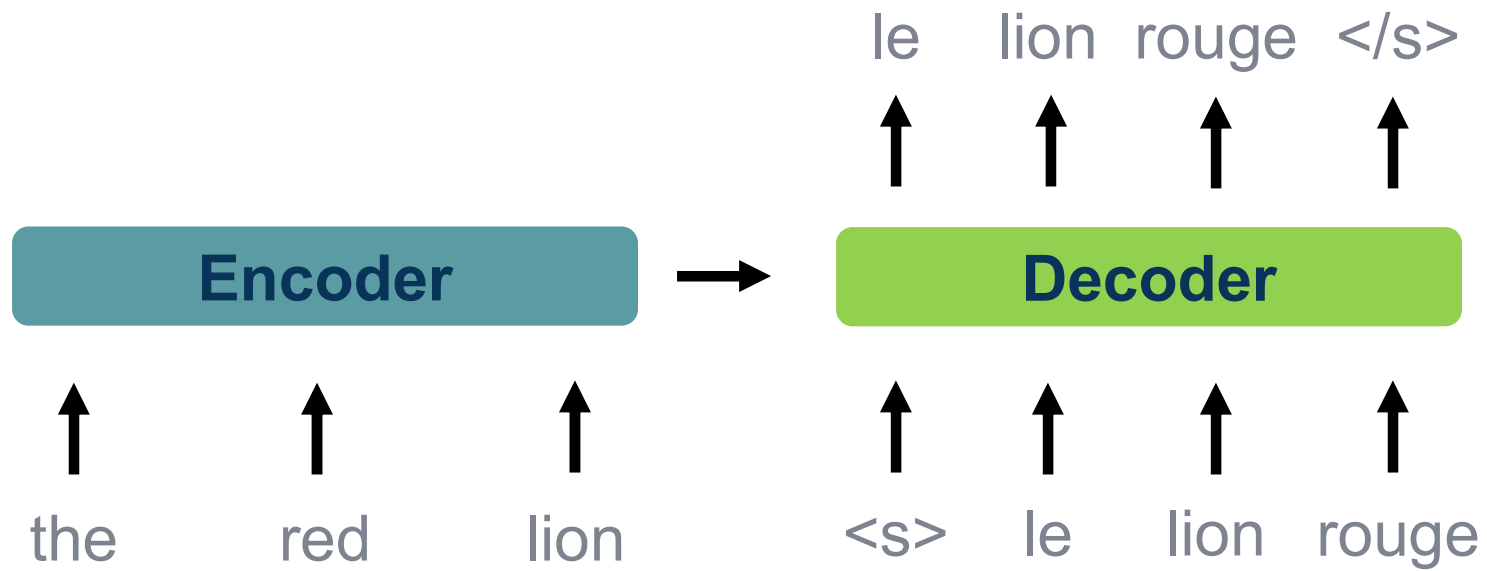


Figure from Latent Alignment and Variational Attention by Deng et. al.



Sequence-to-Sequence Model

FACEBOOK AI





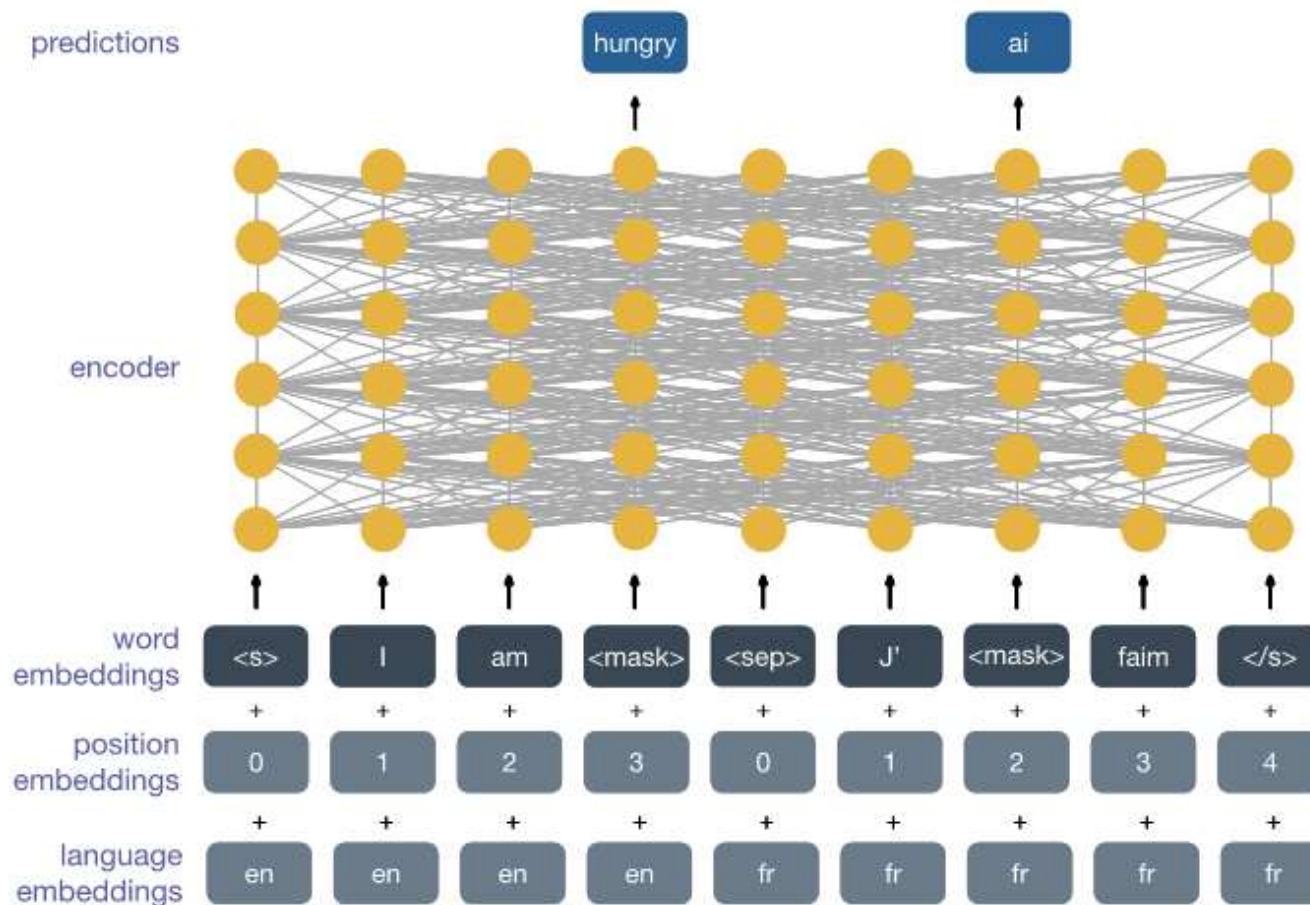
User Language
Prediction

English: 0.99
French : 0.95



Language
Identification

Turkish: 0.99

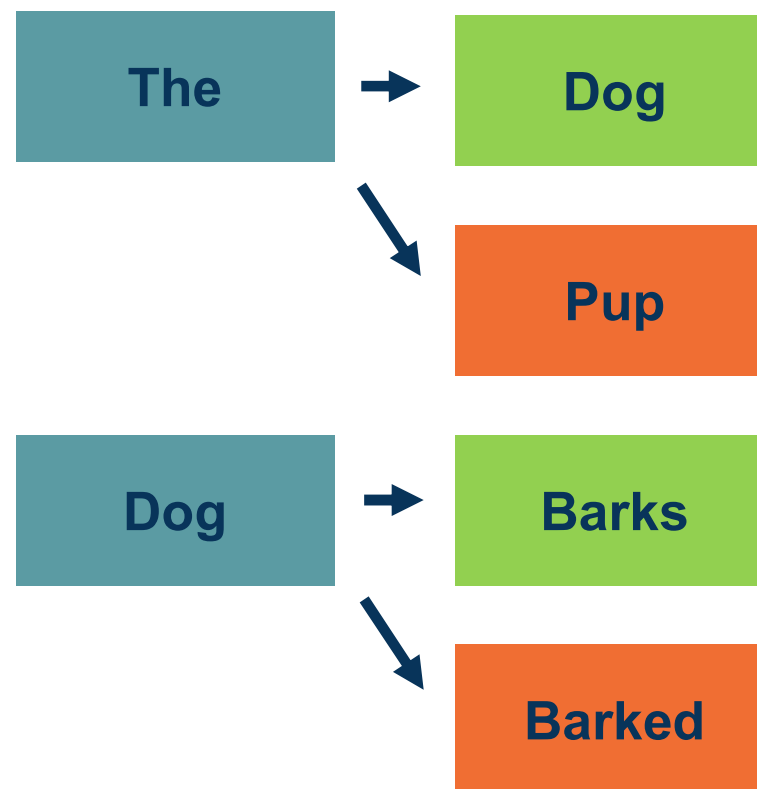


Cross-Lingual Masked Language Modeling

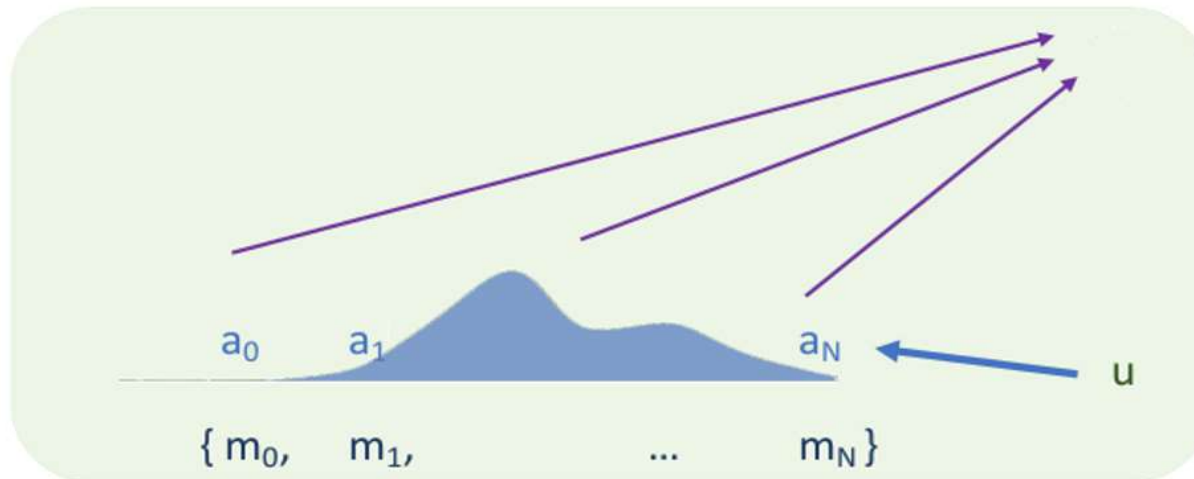
FACEBOOK AI



- ⬢ Search exponential space in linear time
- ⬢ Beam size k determines “width” of search
- ⬢ At each step, extend each of k elements by one token
- ⬢ Top k overall then become the hypotheses for next step



- Given a set of vectors $\{u_1, \dots, u_N\}$ and a “query” vector q
- We can select the most similar vector to q via $p = \text{Softmax}(Uq)$



$$a_j = \frac{e^{u_j \cdot q}}{\sum_k e^{u_k \cdot q}}$$

$$\text{output} = \sum_k a_k u_k$$