

Intuition Behind Self-Attention

Attending to the most important parts of an input.



Intuition Behind Self-Attention

Attending to the most important parts of an input.



1. Identify which parts to attend to
2. Extract the features with high attention

Intuition Behind Self-Attention

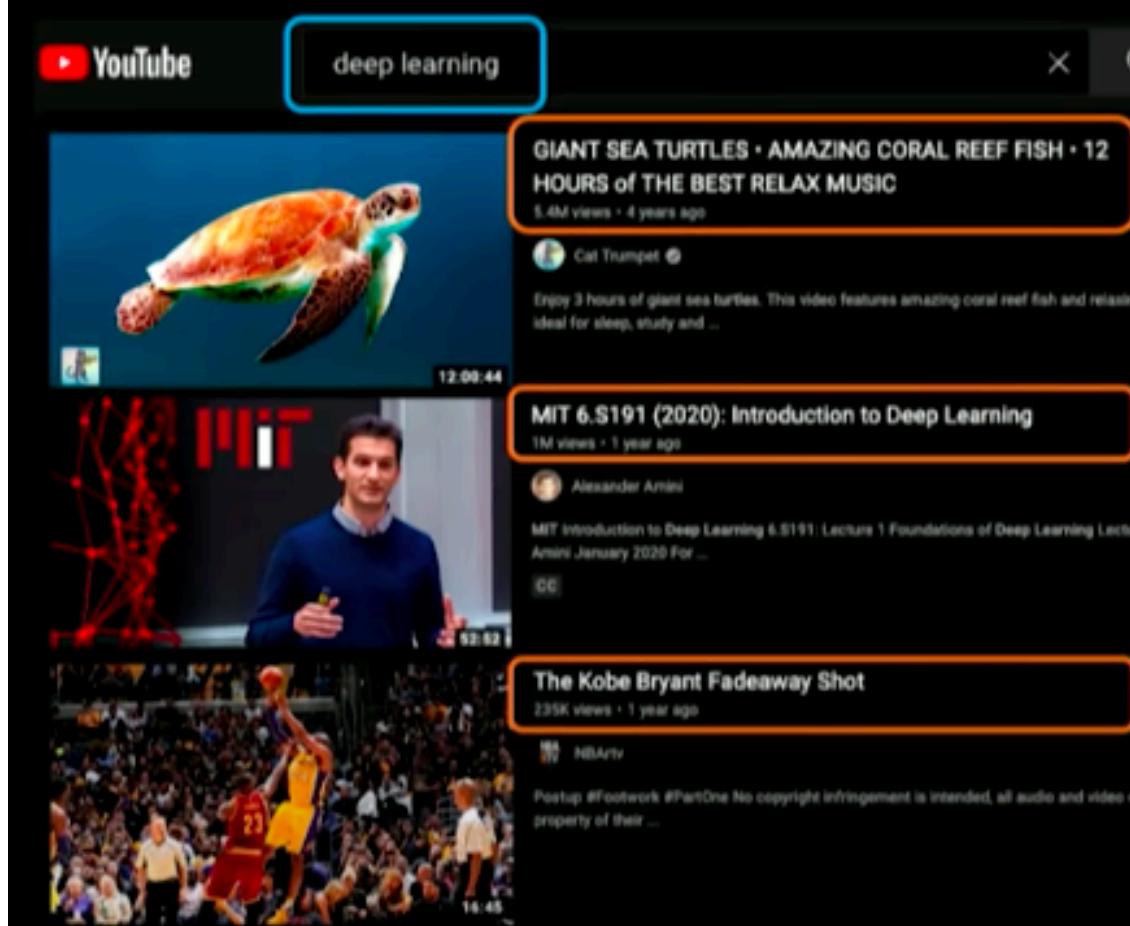
Attending to the most important parts of an input.



1. Identify which parts to attend to
2. Extract the features with high attention

Similar to a
search problem!

Understanding Attention with Search



Query (Q)

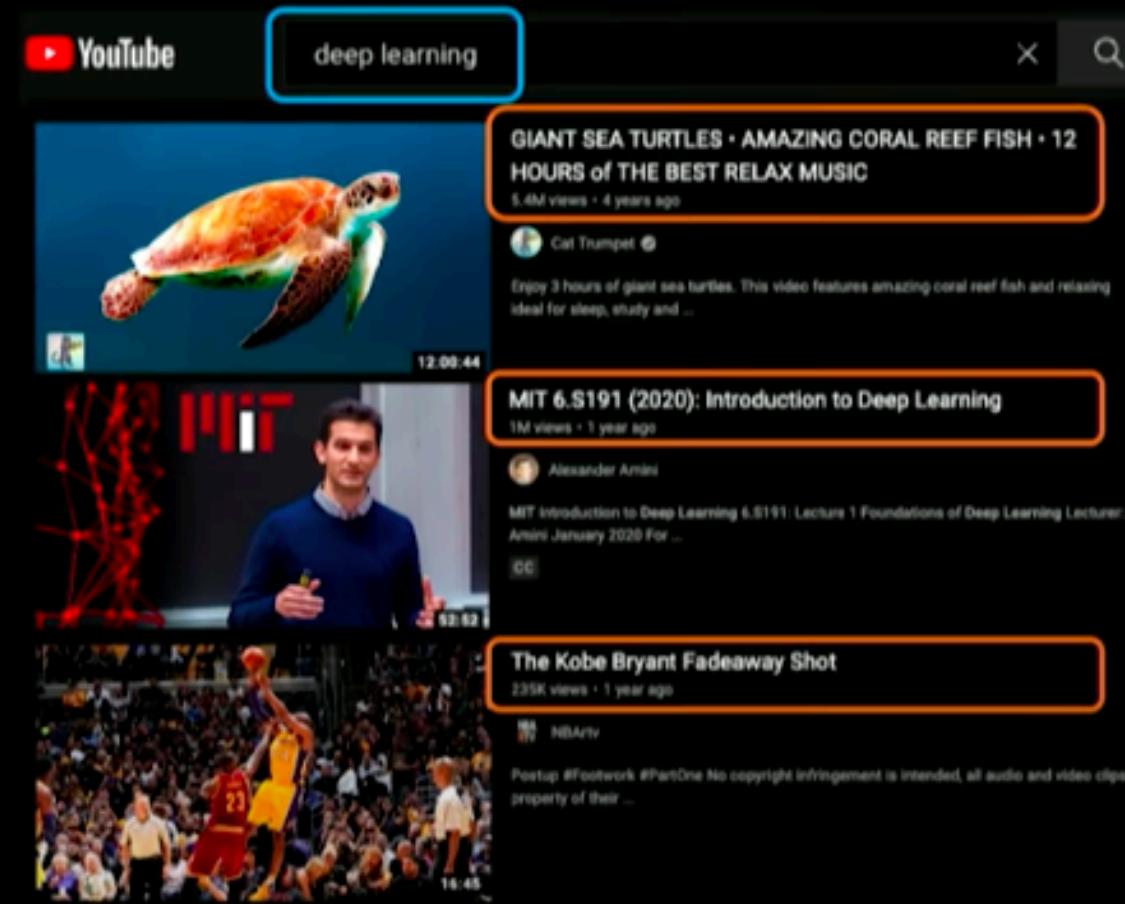
Key (K_1)

Key (K_2)

Key (K_3)

I. **Compute attention mask:** how similar is each key to the desired query?

Understanding Attention with Search



Query (Q)

Key (K_1)

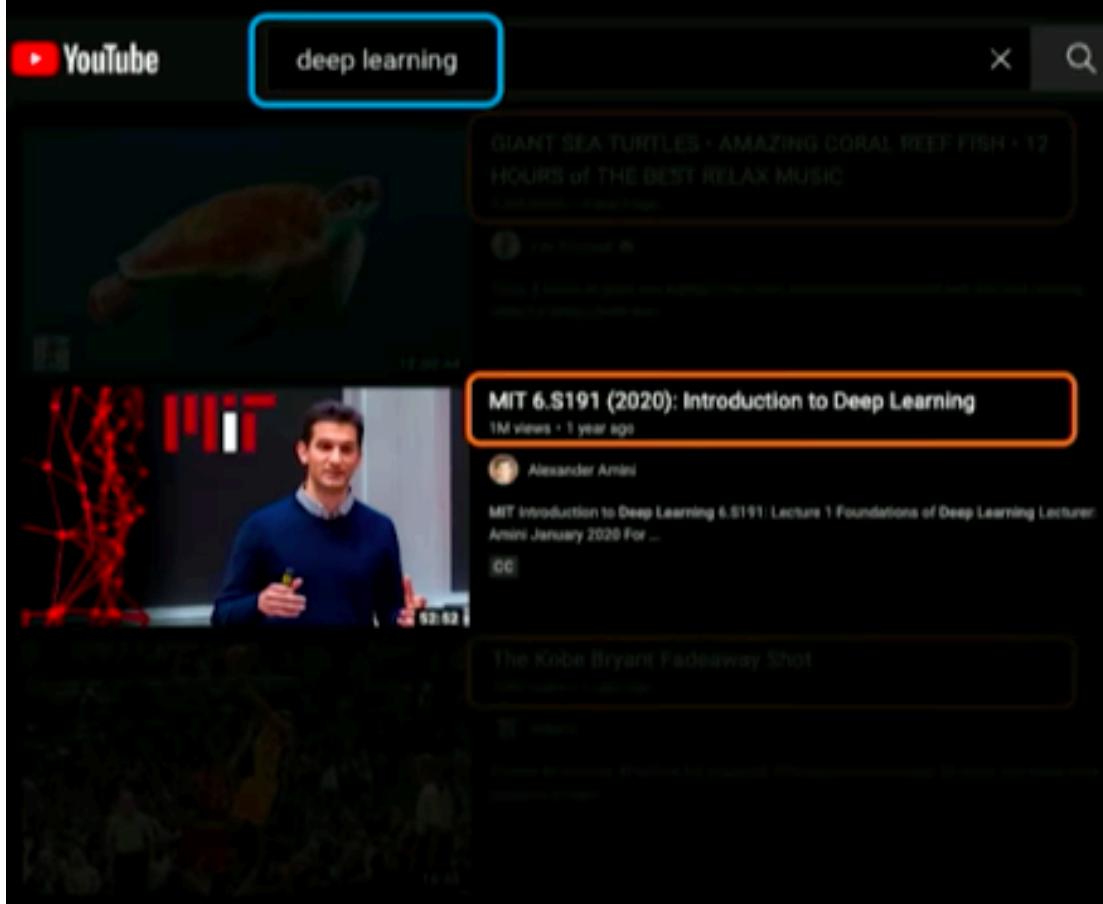
How similar is the key to the query?

Key (K_2)

Key (K_3)

I. **Compute attention mask:** how similar is each key to the desired query?

Understanding Attention with Search



Query (Q)

Key (K_1)

Key (K_2)

Key (K_3)

How similar is the key to the query?

I. **Compute attention mask:** how similar is each key to the desired query?

Understanding Attention with Search

YouTube deep learning X Query (Q)

GIANT SEA TURTLES • AMAZING CORAL REEF FISH • 12 HOURS of THE BEST RELAX MUSIC

MIT 6.S191 (2020): Introduction to Deep Learning

Alexander Amini

MIT Introduction to Deep Learning 6.S191: Lecture 1 Foundations of Deep Learning Lecture Amini January 2020 For ...

CC

MIT 6.S191 (2020): Introduction to Deep Learning

Key (K_1)

Key (K_2)

Value (V)

Key (K_3)

2. Extract values based on attention:
Return the values highest attention

A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

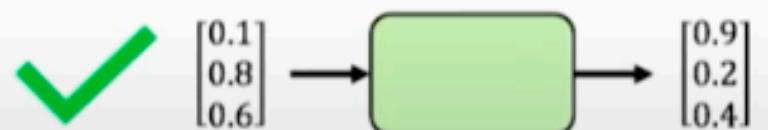
given these words

predict the
next word

Representing Language to a Neural Network



Neural networks cannot interpret words

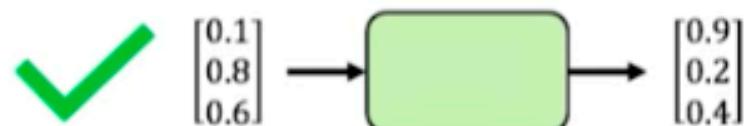


Neural networks require numerical inputs

Encoding Language for a Neural Network



Neural networks cannot interpret words



Neural networks require numerical inputs

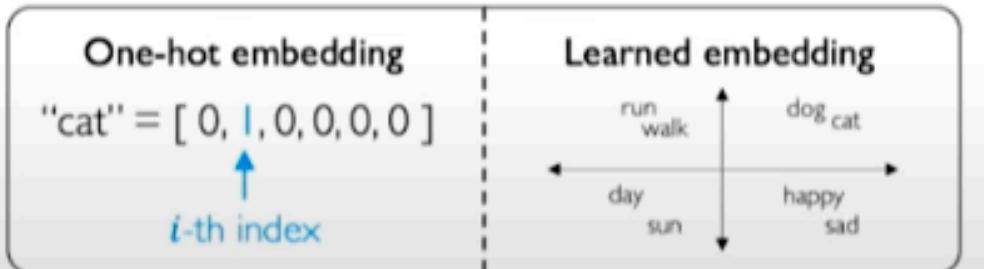
Embedding: transform indexes into a vector of fixed size.

this	cat	for
my	took	I
a	walk	morning

1. Vocabulary:
Corpus of words

a	→	1
cat	→	2
...
walk	→	N

2. Indexing:
Word to index



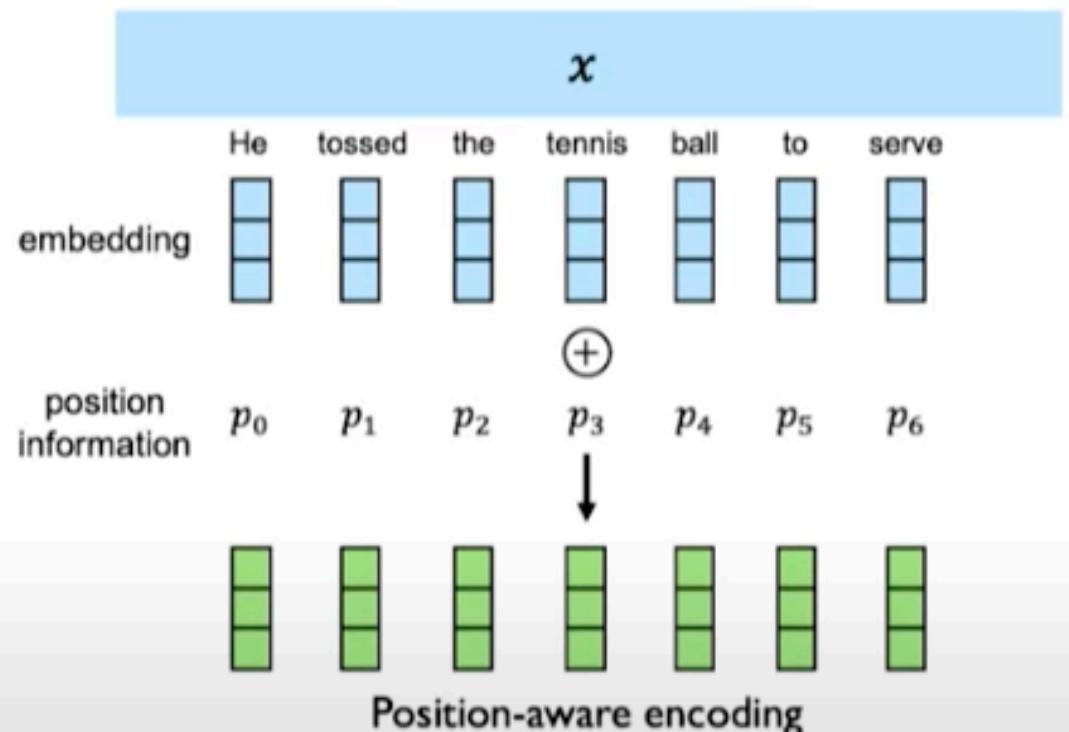
3. Embedding:
Index to fixed-sized vector

Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

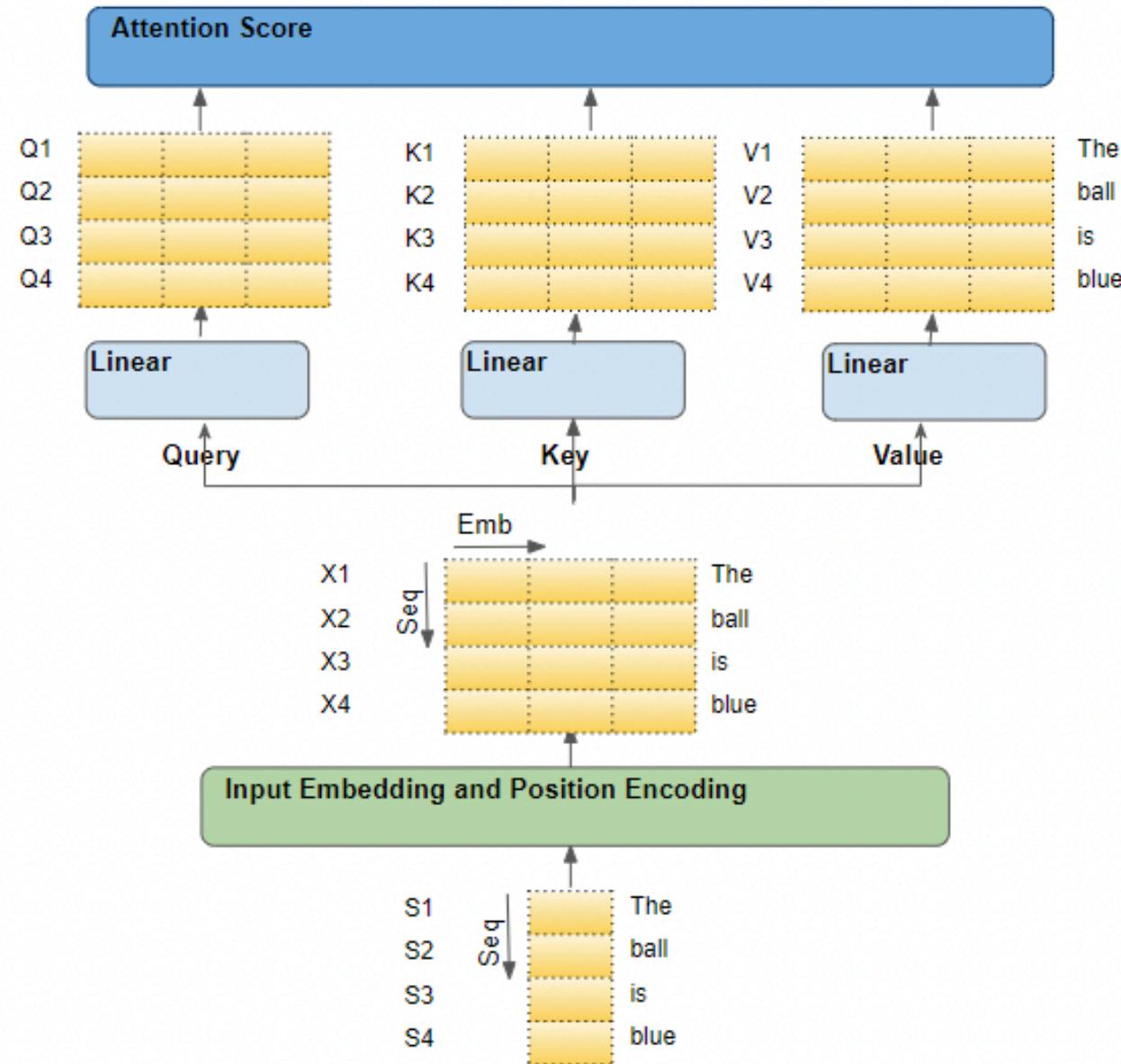
I. Encode position information

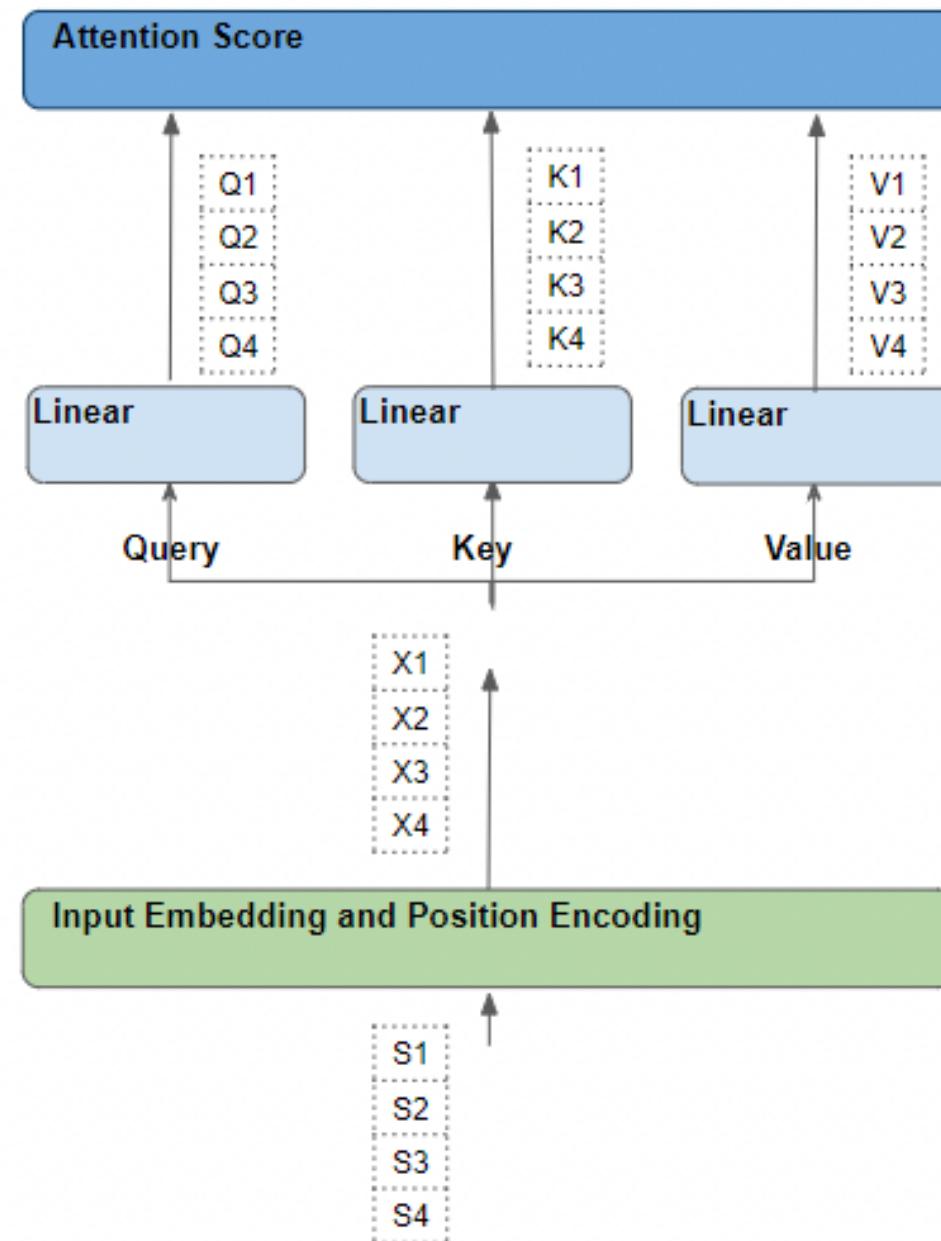
- |↓ learn query, key, value functions
- |↓ compute attention weighting
- |↓ learn features with high attention

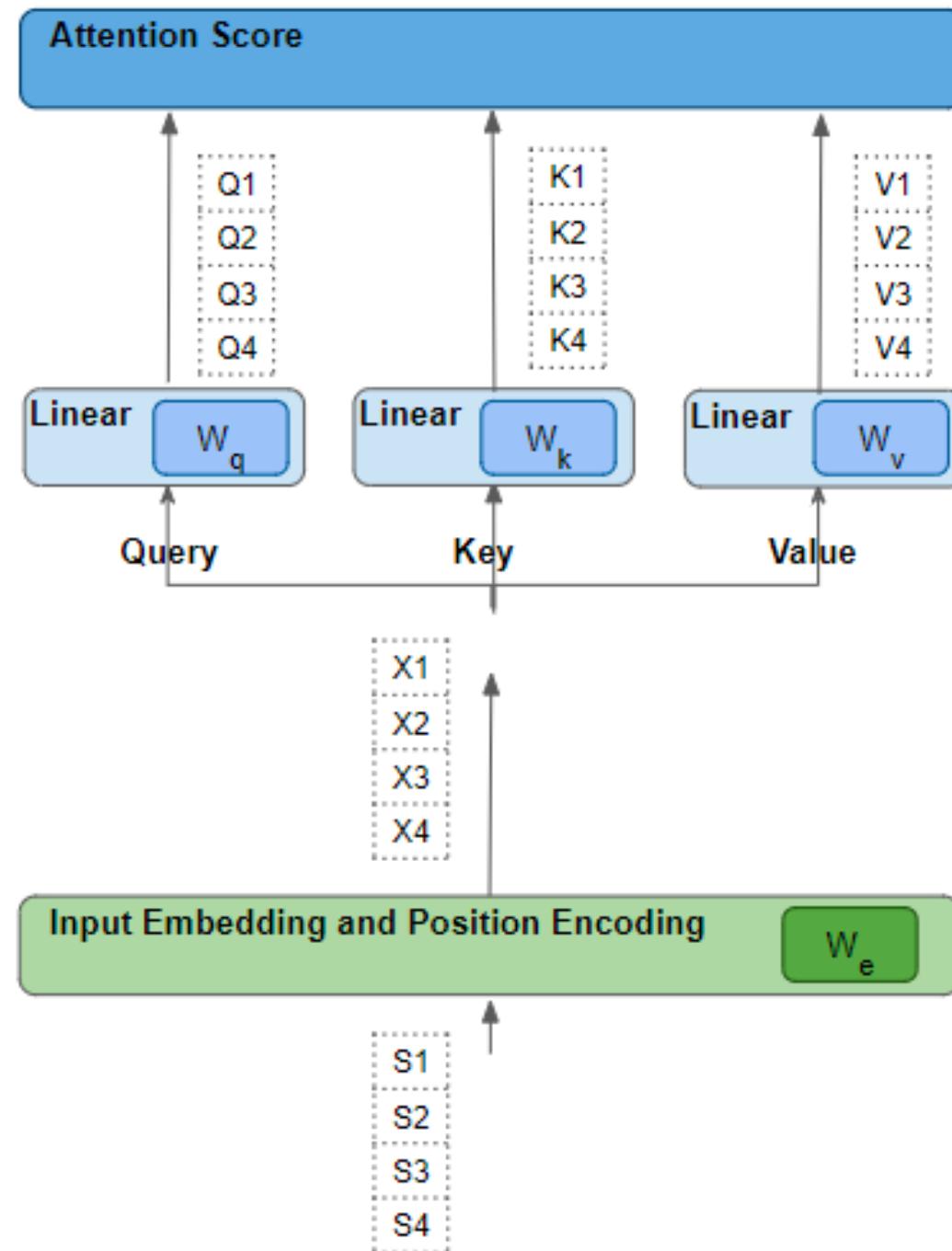


Data is fed in all at once! Need to encode position information to understand order.

The Attention Module

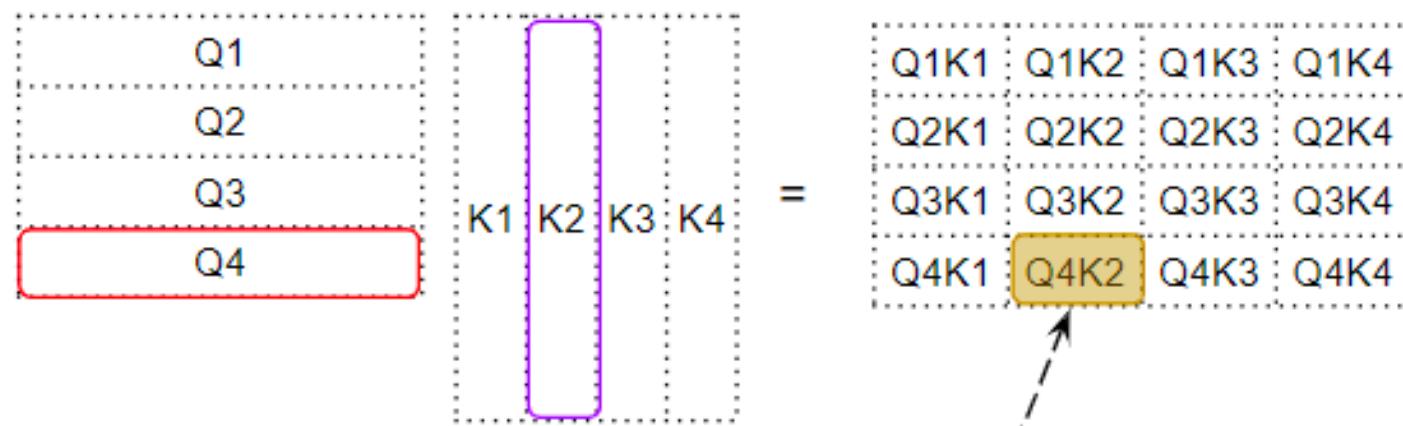




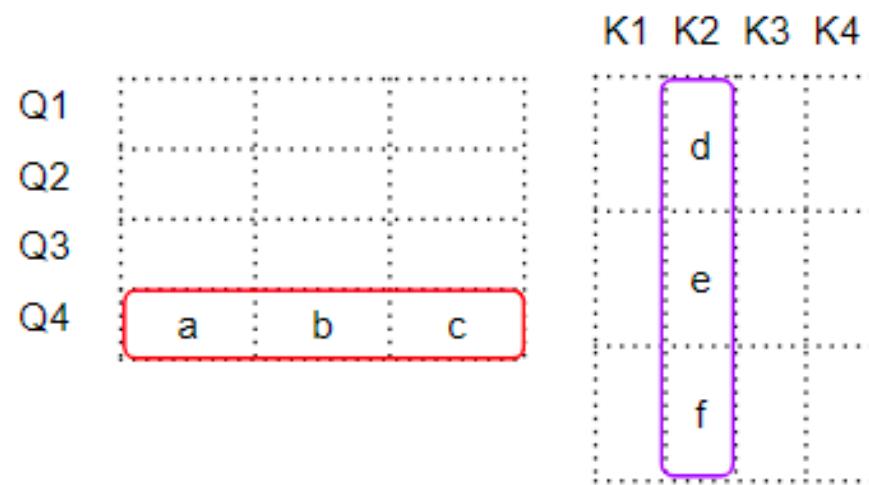


Q1					Q1K1	Q1K2	Q1K3	Q1K4
Q2					Q2K1	Q2K2	Q2K3	Q2K4
Q3	K1	K2	K3	K4	Q3K1	Q3K2	Q3K3	Q3K4
Q4					Q4K1	Q4K2	Q4K3	Q4K4

Q1					Q1K1	Q1K2	Q1K3	Q1K4
Q2					Q2K1	Q2K2	Q2K3	Q2K4
Q3	K1	K2	K3	K4	Q3K1	Q3K2	Q3K3	Q3K4
Q4					Q4K1	Q4K2	Q4K3	Q4K4



$$Q4K2 = a*d + b*e + c*f$$



Learning Self-Attention with Neural Networks

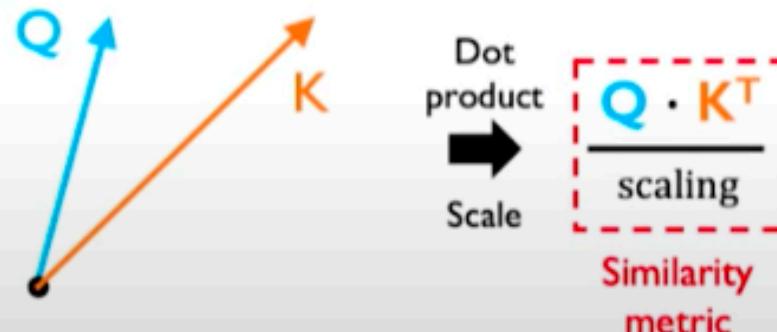
Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query**, **key**, **value** for search
3. Compute **attention weighting**

Identify features with high attention

Attention score: compute pairwise similarity between each **query** and **key**

How to compute similarity between two sets of features?



Also known as the "cosine similarity"

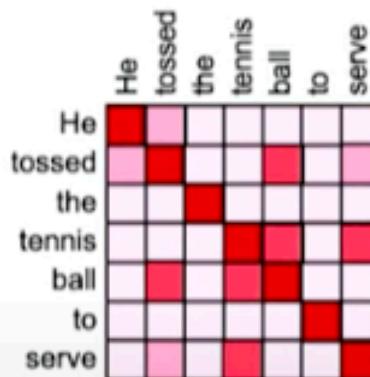
Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**

1. Identify features with high attention

Attention weighting: where to attend to!
How similar is the key to the query?



$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right)$$

Q1K1	Q1K2	Q1K3	Q1K4
Q2K1	Q2K2	Q2K3	Q2K4
Q3K1	Q3K2	Q3K3	Q3K4
Q4K1	Q4K2	Q4K3	Q4K4

V1
V2
V3
V4

$$= \begin{array}{l} Q1K1V1 + Q1K2V2 + Q1K3V3 + Q1K4V4 \\ Q2K1V1 + Q2K2V2 + Q2K3V3 + Q2K4V4 \\ Q3K1V1 + Q3K2V2 + Q3K3V3 + Q3K4V4 \\ Q4K1V1 + Q4K2V2 + Q4K3V3 + Q4K4V4 \end{array}$$

$$= \begin{array}{l} Z1 \\ Z2 \\ Z3 \\ Z4 \end{array}$$

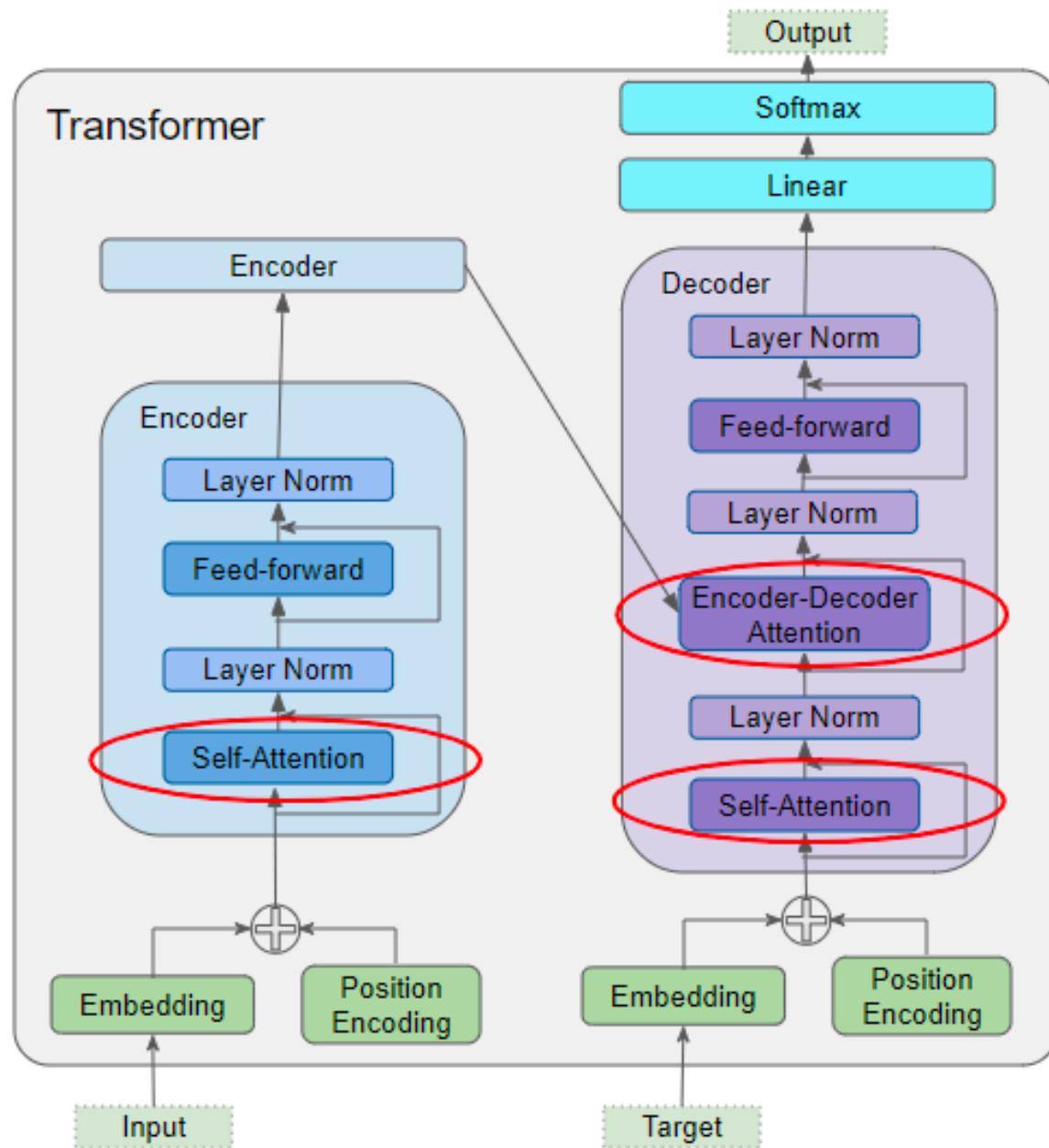
Fourth word Score

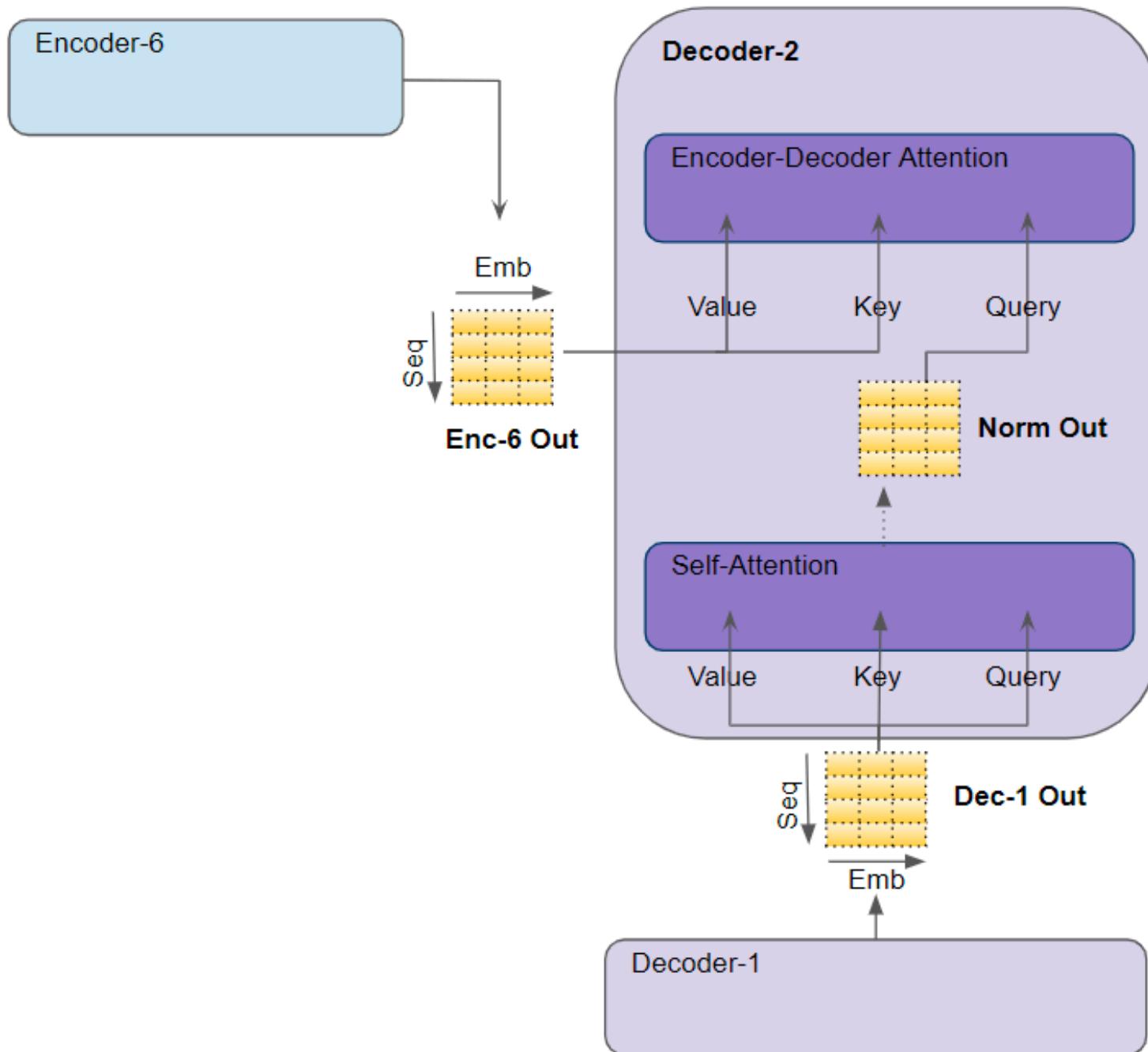
*Fourth Query word * first Key word*

$$Z_4 = (Q_4 K_1) V_1 + (Q_4 K_2) V_2 + (Q_4 K_3) V_3 + (Q_4 K_4) V_4$$

*Fourth Query word * second Key word*

$$Z = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



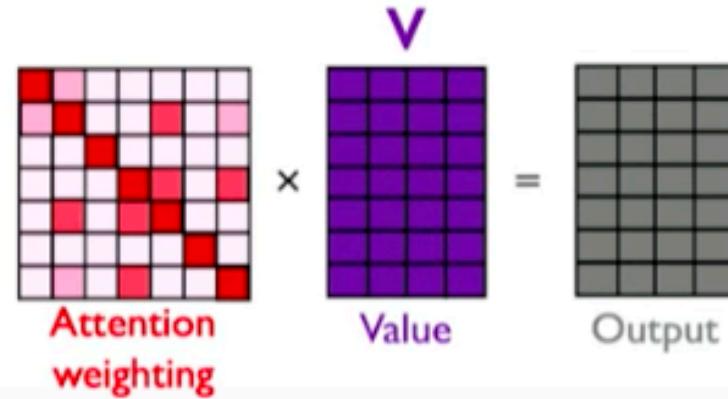


Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

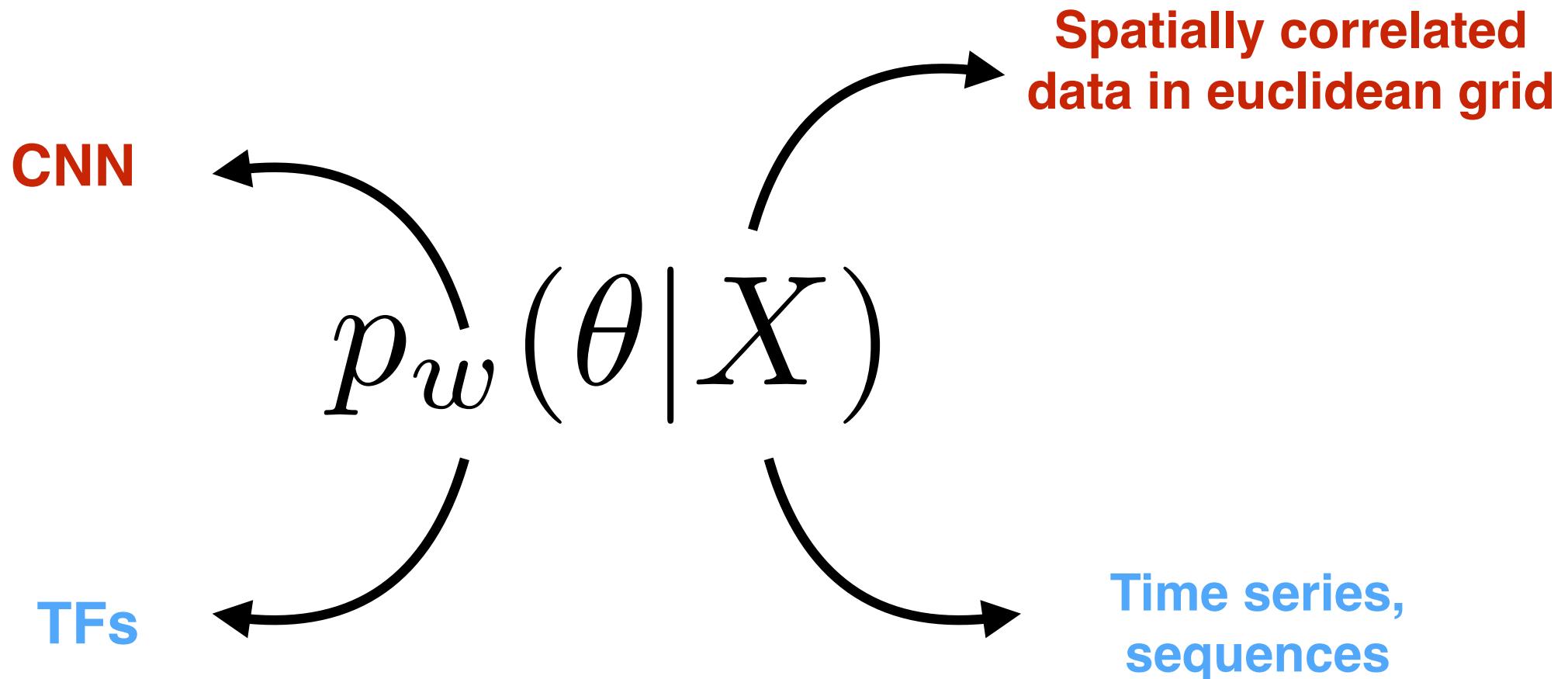
1. Encode position information
2. Extract **query**, **key**, **value** for search
3. Compute **attention weighting**
4. Extract features with high attention

Last step: self-attend to extract features



$$\frac{\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right) \cdot V}{\text{---}} = A(Q, K, V)$$

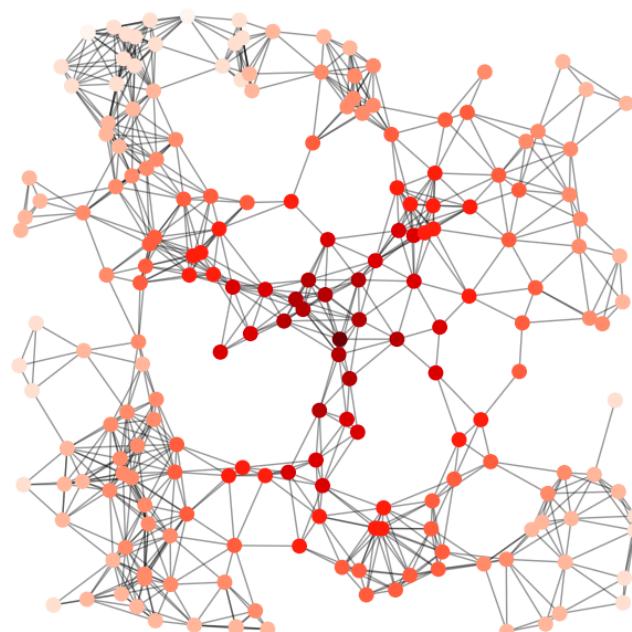
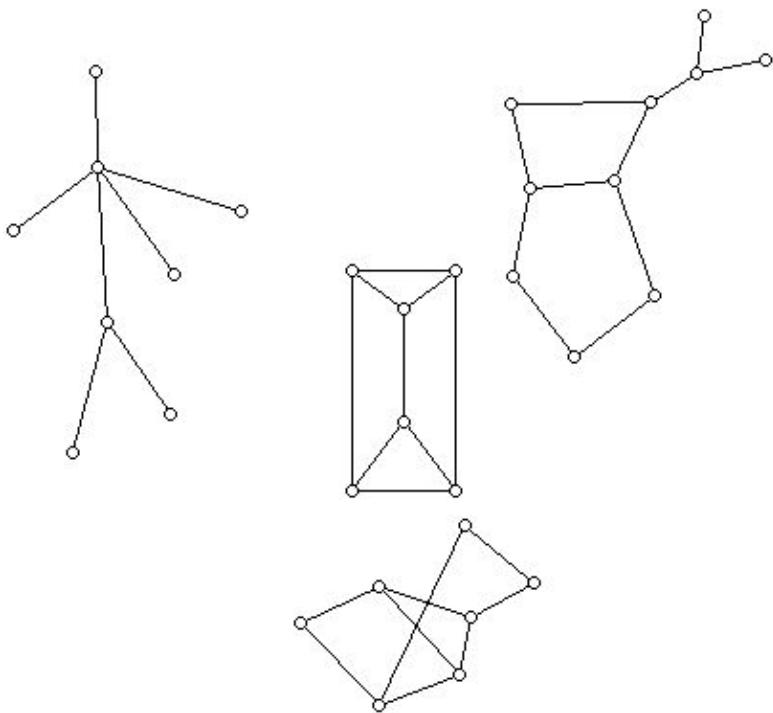
RECAP4:



Approach	Discriminative		Generative	
	Target Function	Method	Target Function	Method
Data type	$p(y x)$	All supervised networks with bottleneck	$p(x y)$	Conditional generative models
Labels	$p(\hat{y} x)$	Self-supervised learning	$p(x)$	Generative models, Autoencoders
No Labels				

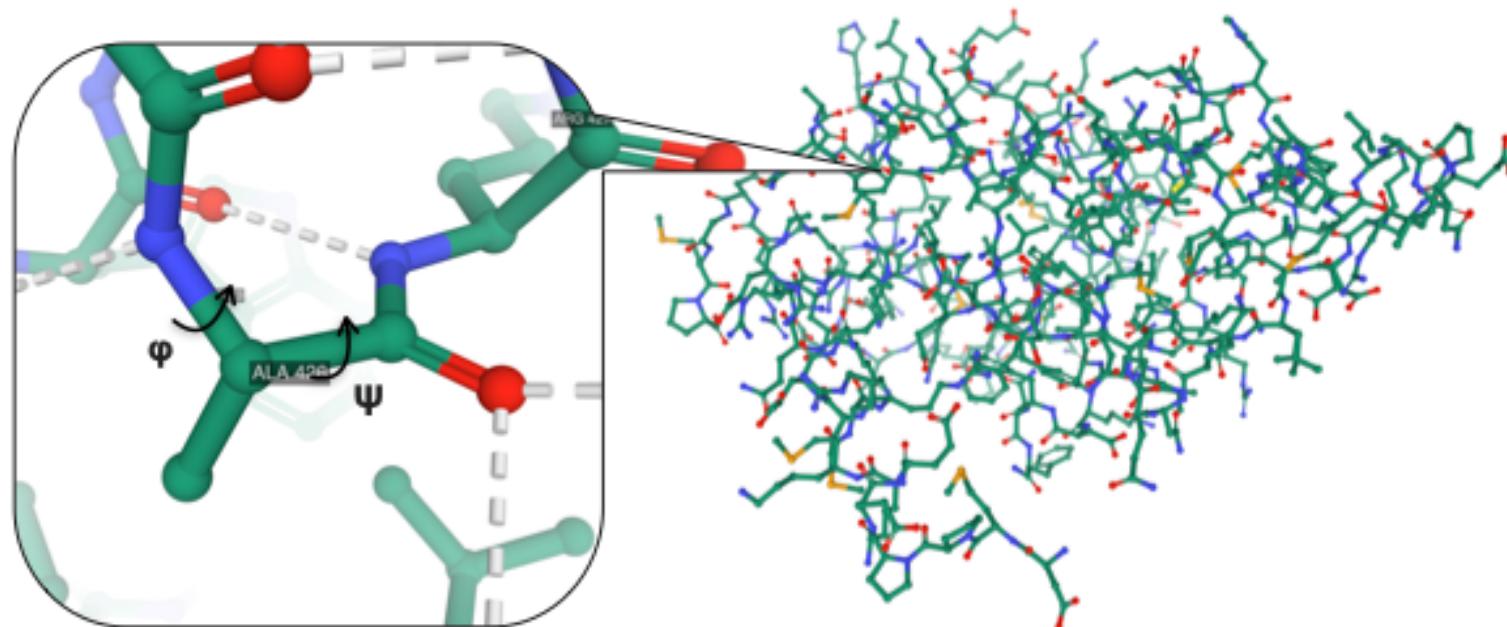
A very brief incursion into GNNs...

What is a graph?



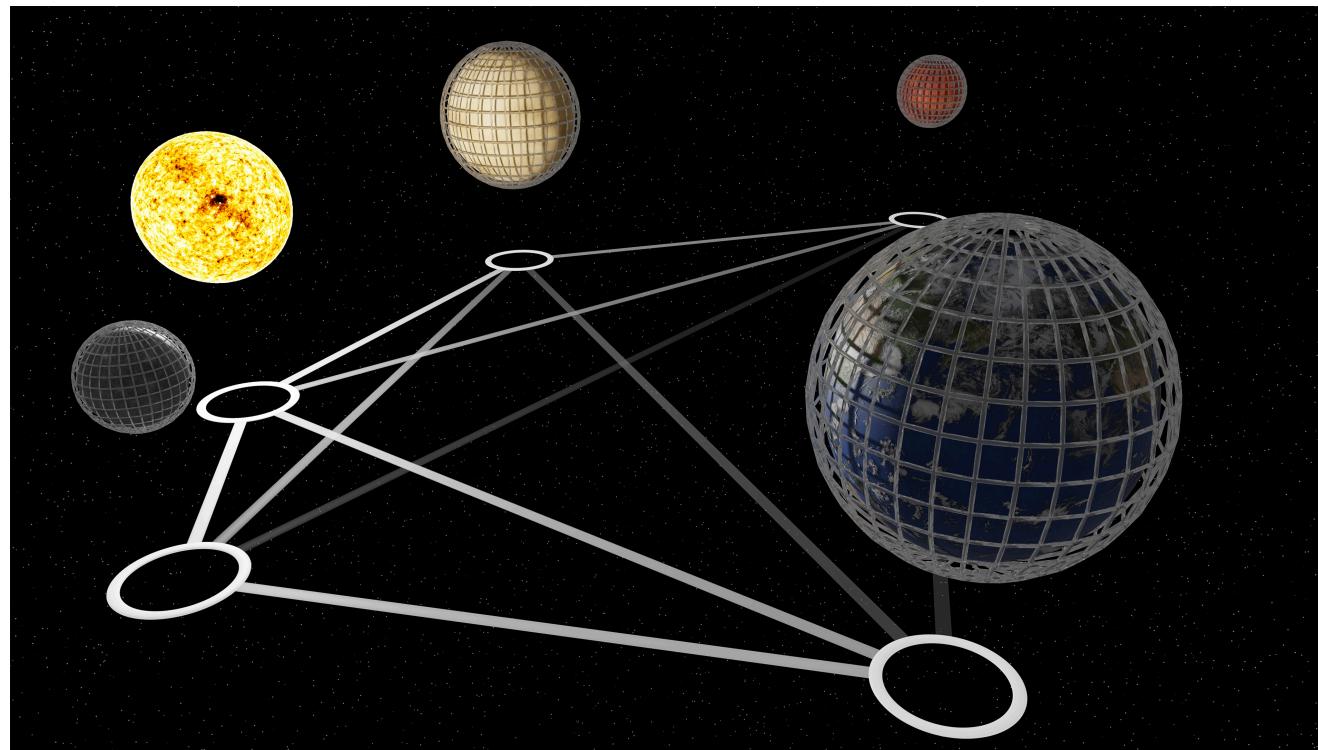
!!!

Molecules are Graphs

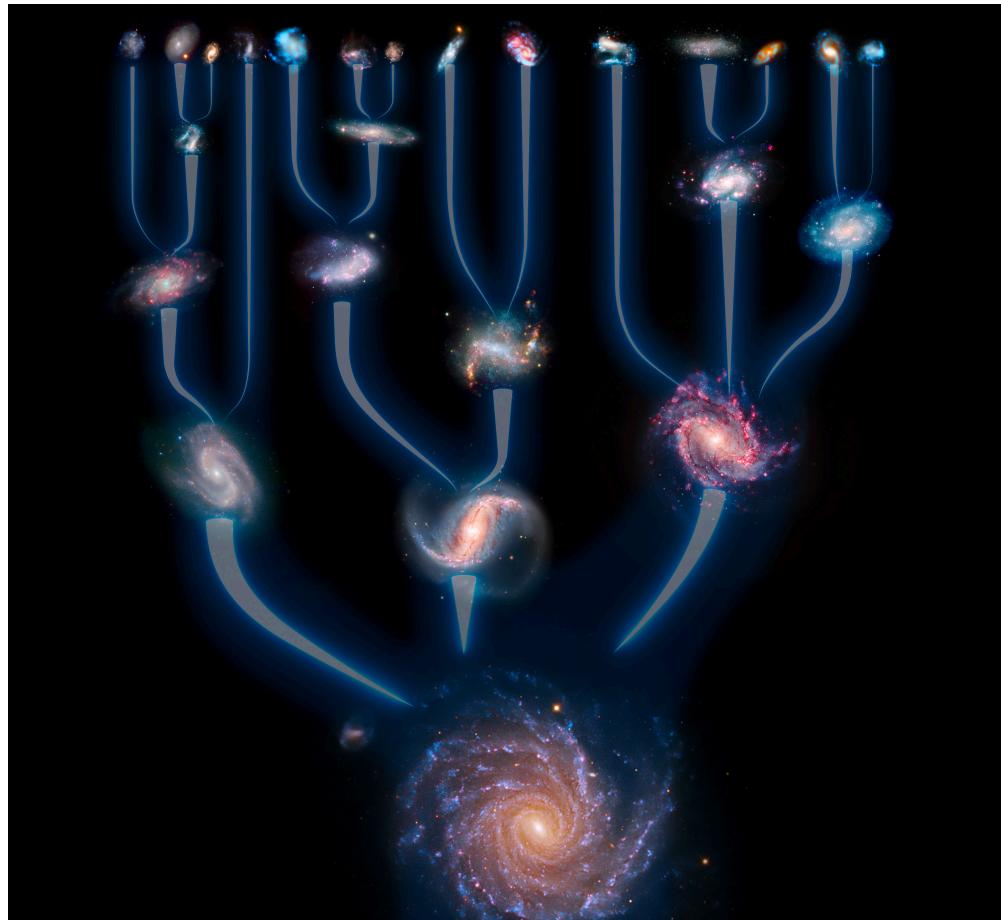


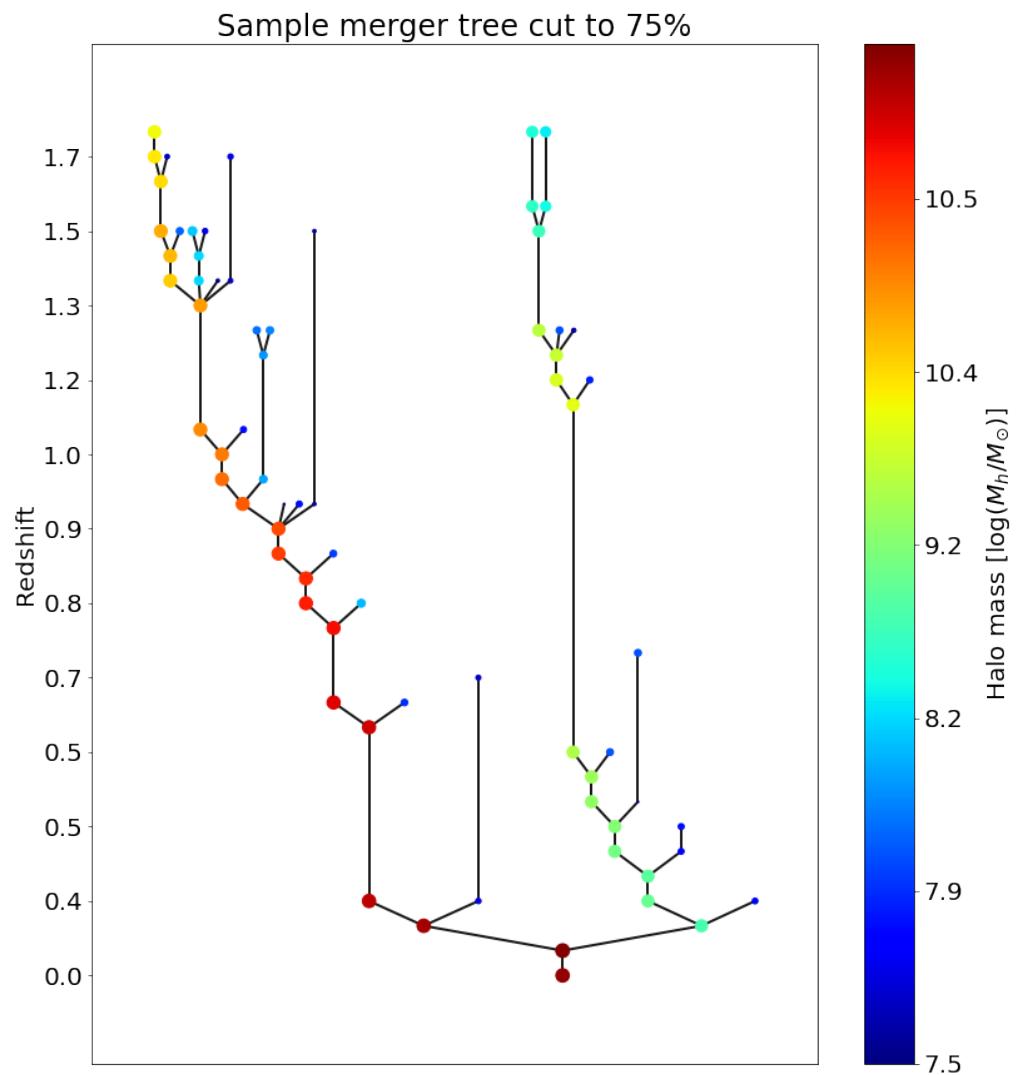
The Solar System is a Graph

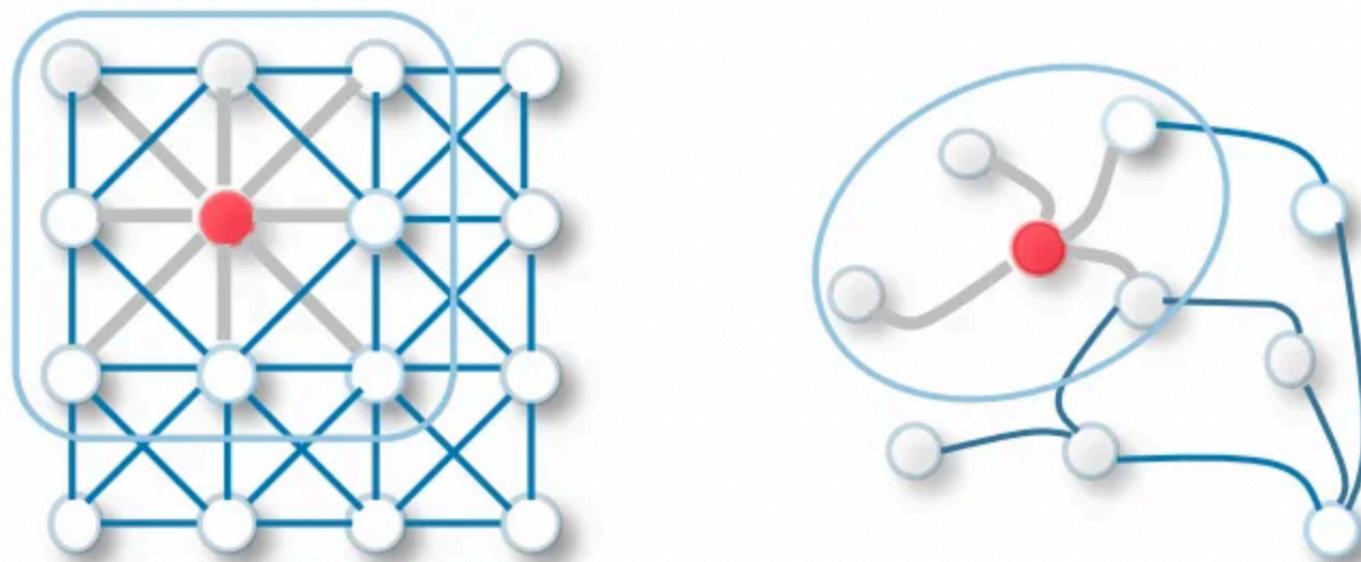
Or at least a natural abstraction



Galaxy Merger trees are Graphs



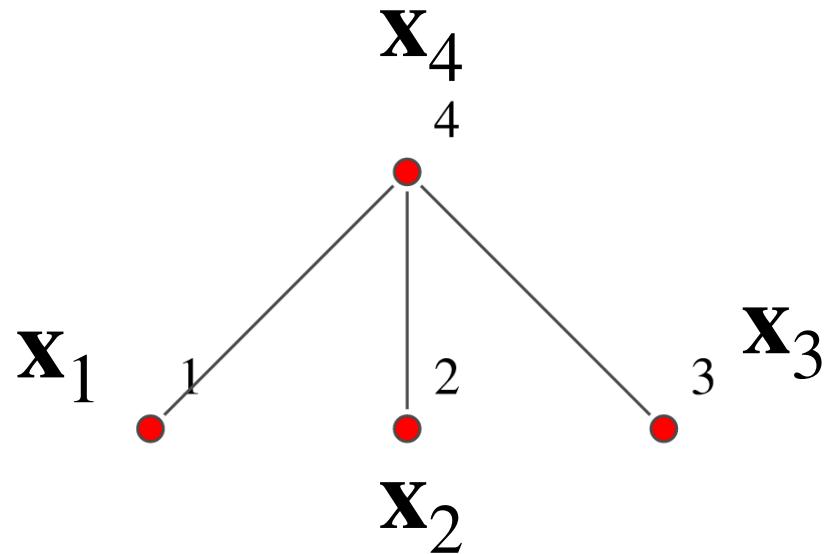




The main difference compared to images
is in the definition of neighbours

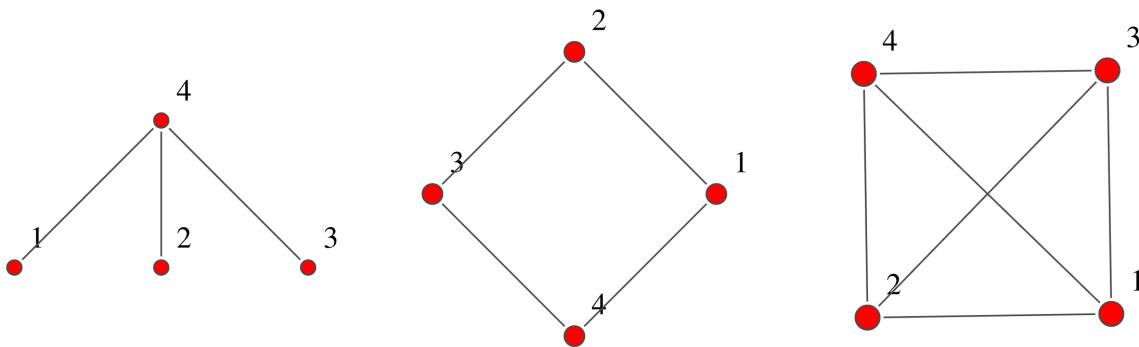
What makes up a graph?

Nodes and node features



What makes up a graph?

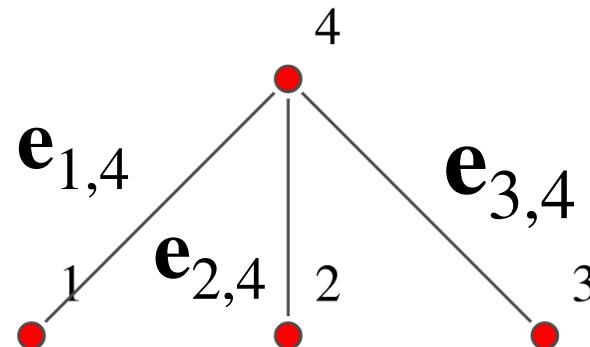
Nodes and node features



$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4]^T \rightarrow \mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{|\mathcal{V}|}]^T \rightarrow \mathbf{X} = \{\mathbf{x}_i \mid i \in N_v\} \in \mathbb{R}^{N_v \times d_v}$$

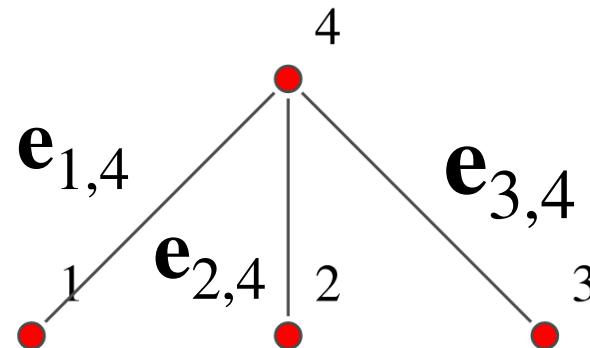
What makes up a graph?

Edges and edge features



What makes up a graph?

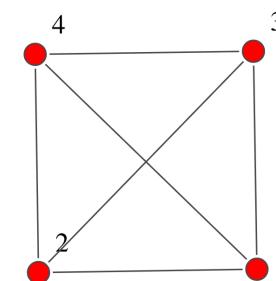
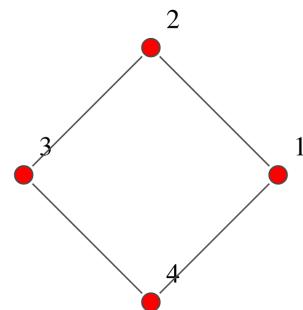
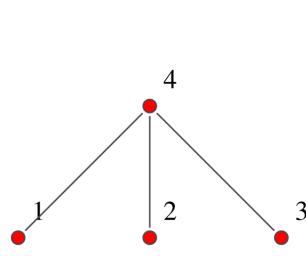
Edges and edge features



$$\mathbf{E} = [\mathbf{e}_{12}, \mathbf{e}_{13}, \mathbf{e}_{14}]^\top \rightarrow \mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{|\mathcal{E}|}]^\top \rightarrow \mathbf{E} = \{\mathbf{e}_i \mid i \in N_e\} \in \mathbb{R}^{N_e \times d_e}$$

What makes up a graph?

The Adjacency Matrix

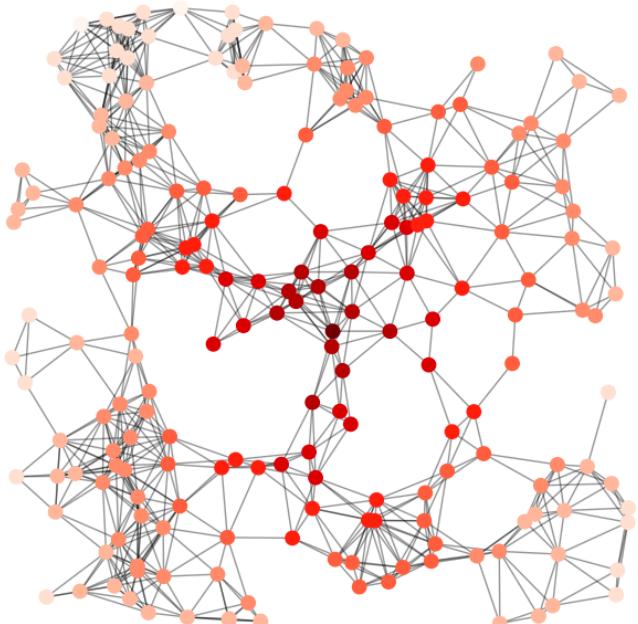


$$\begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

Global Features



Could be physical properties e.g.

$$\mathbf{U} = [\Omega_m, \sigma_8]$$

Statistical properties,
e.g. manually inserted

$$\mathbf{U} = [\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}, \max(\mathbf{x}), \min(\mathbf{x}) \dots P_n(\mathbf{x})]$$

Or completely learnable

Summarizing...

- Graph defined as $G(V, E, U)$ with G being the graph, V being the set of nodes, E being the set of edges and U being the global feature vector
- Adjacency matrix (A) encodes the structure of the graph and is a binary matrix

Locality constraints

The edges of a graph naturally define the concept of a local **neighborhood**

$$\mathcal{N}_u = \{v \mid (u, v) \in E \vee (v, u) \in E\}$$

Which allows us to define local functions, which work on a given node and its neighborhood.

$$\phi(\mathbf{x}_u, \mathbf{X}_{\mathcal{N}_u})$$

Increasing
complexity/
expressivity/
training time

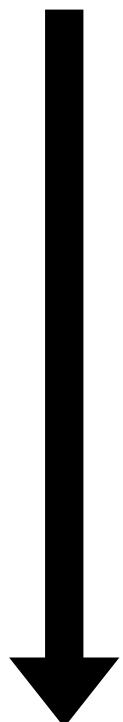
Graph Neural Network Layers

- Convolutional Layer

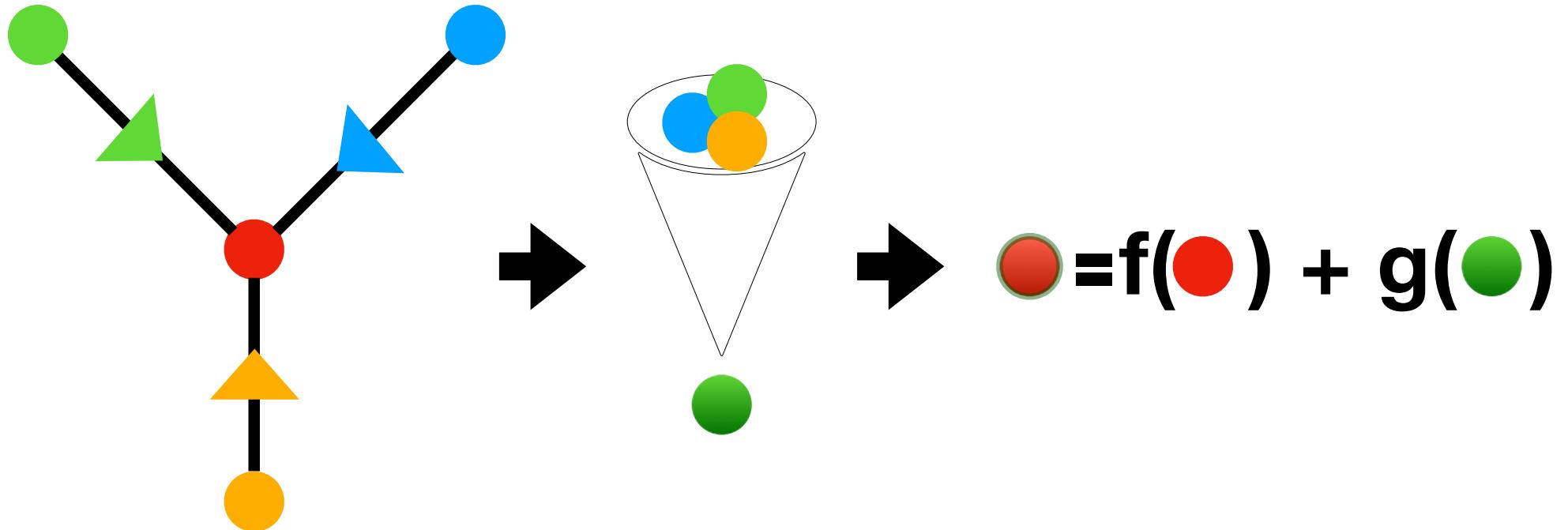
$$\mathbf{h}_u = \phi \left(\mathbf{x}_u, \bigoplus_{v \in \mathcal{N}_u} c_{vu} \psi(\mathbf{x}_v) \right)$$

- Attentional Layer

$$\mathbf{h}_u = \phi \left(\mathbf{x}_u, \bigoplus_{v \in \mathcal{N}_u} a(\mathbf{x}_u, \mathbf{x}_v, \mathbf{e}_{uv}) \psi(\mathbf{x}_v) \right)$$



Looks complicated, but!



Looks complicated, but!

Essentially, these are all variants of:

1. Define a node function f and a neighborhood function g (f and g can be the same function)
2. Act with f on the node and with g on the neighborhood

Graphs are good representations of physical systems

In modern physics ... a central theme will be a Geometric Principle: The laws of physics must all be expressible as geometric (coordinate-independent and reference frame-independent) relationships between geometric objects (scalars, vectors, tensors, ...) that represent physical entities.

Physics on Graphs

- A natural abstraction of physical systems and inherently local -> easier physical interpretations/separability
- Embeds inductive biases easily by restructuring the graph -> more efficient learning and no need to learn things we already know
- Can embed permutational, rotational, translational and reflectional symmetries

