

Quantum Self-Attention

Dominic Pasquali

March, 23 2022

2022 CERN openlab Technical Workshop



QUANTUM
TECHNOLOGY
INITIATIVE



Agenda

- Motivation for Quantum GANs & Quantum Self-Attention
- Classical Self-Attention
- Quantum Self-Attention
- Multi-Head Attention

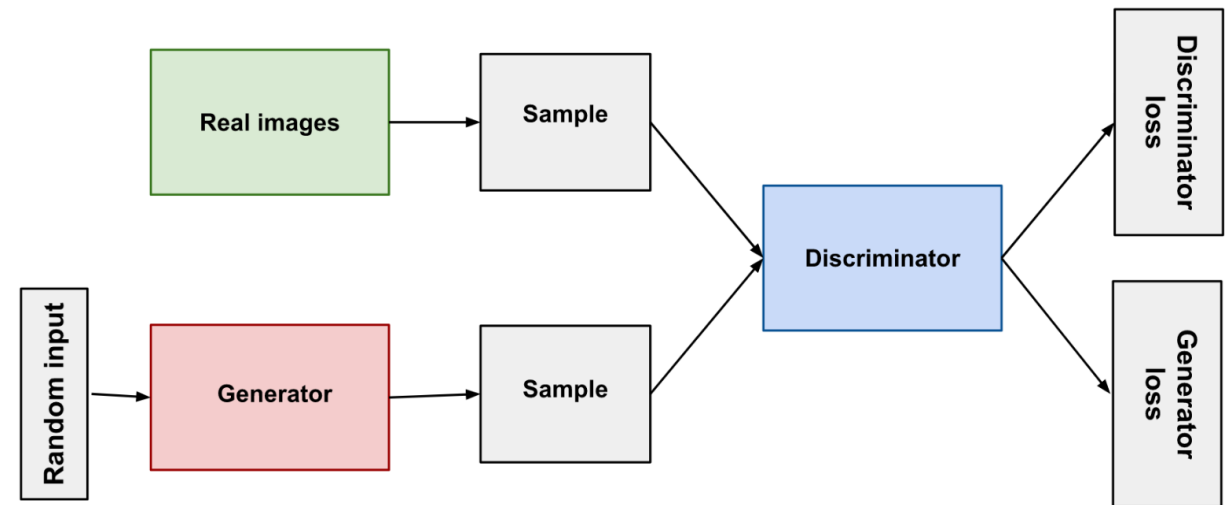
Motivation For Quantum GANs & Quantum Self-Attention

Why GANs:

- Simulation of particle transport through matter is fundamental for understanding the physics of High Energy Physics (HEP) experiments
- Most of LHC CPU budget (~ 1M CPU-years!!!) is dedicated to Monte Carlo simulation
- Faster approach: Replace Monte Carlo simulation with deep learning algorithms (e.g. GANs)

Why QGANs:

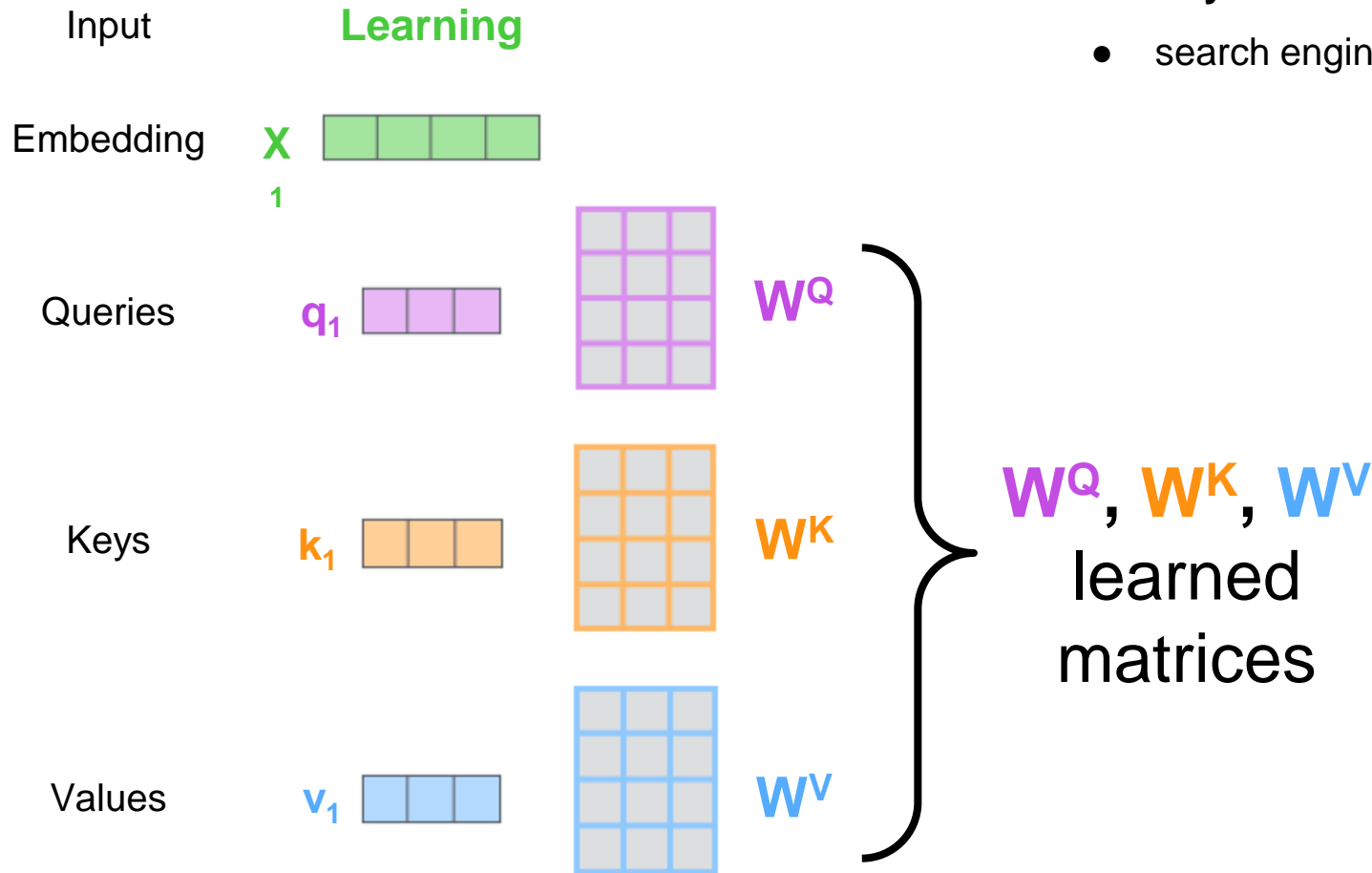
- compressed data representation in quantum states
- expect faster training with less number of parameters
- potential advantage of Quantum GAN^[1]



Explore different prototypes of Quantum GANs to improve model

- **Quantum Self-Attention** in **Classical GANs** to boost performance in hybrid architecture

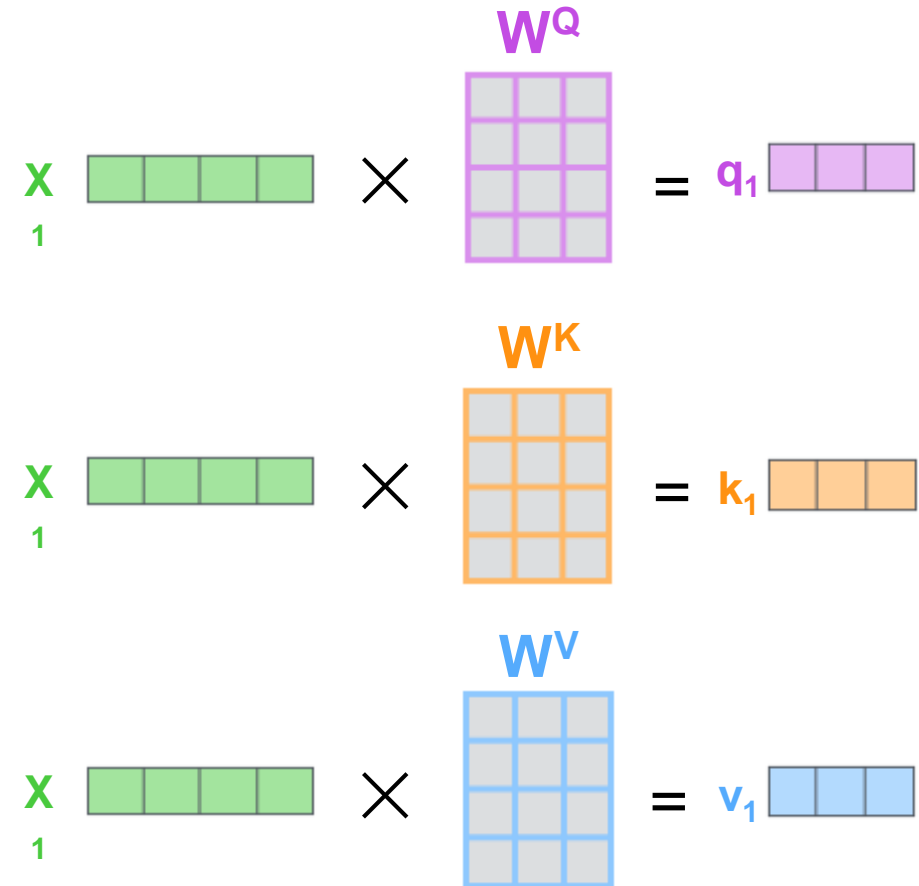
Classical Self-Attention



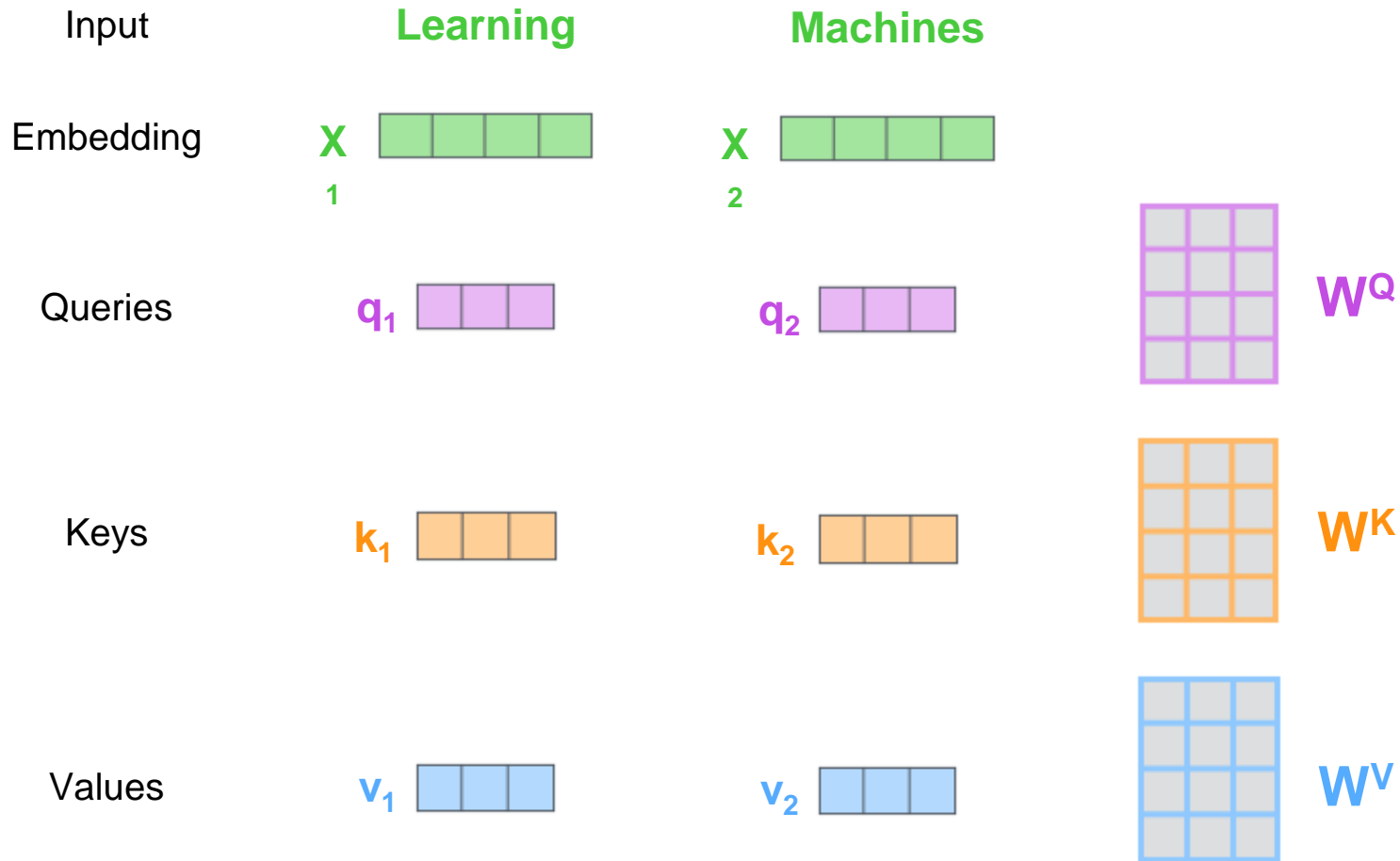
Query, Key, Value concept \square analogous to retrieval systems

Example: When you searching for videos on YouTube's search engine

- search engine maps the **Query (text in the search bar)** against Keys
- **Keys: descriptors** (video title, description, etc.) of YouTube videos
- search engine returns the **best matched videos (Values)**



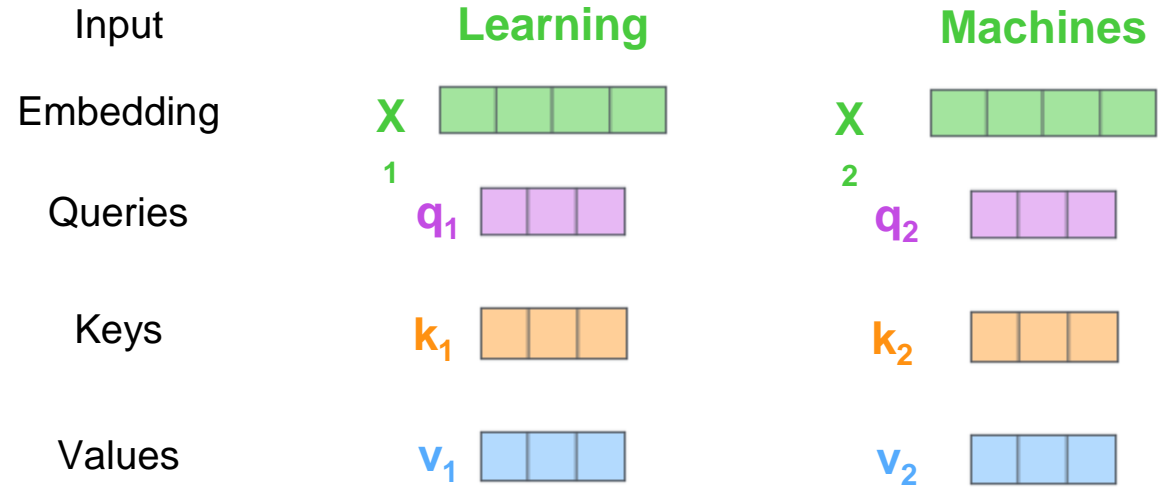
Classical Self-Attention



Adding more words adds more resulting vectors (while using the same learned matrices)

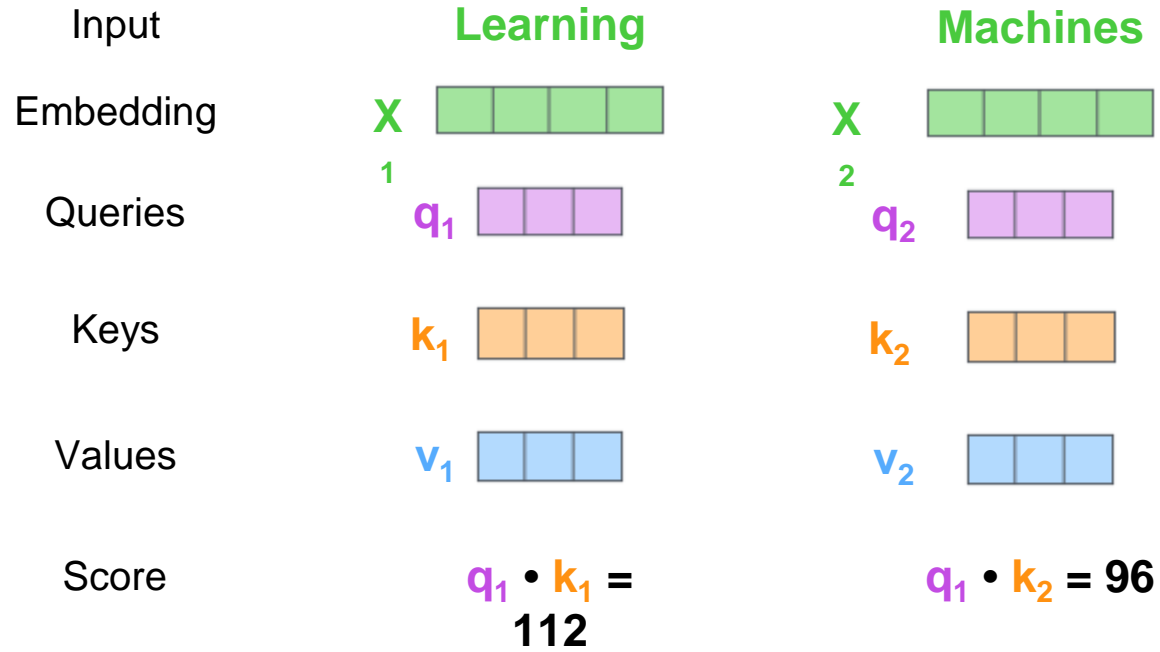
For now, let's only consider 2 words as our input: "Learning Machines"

Classical Self-Attention



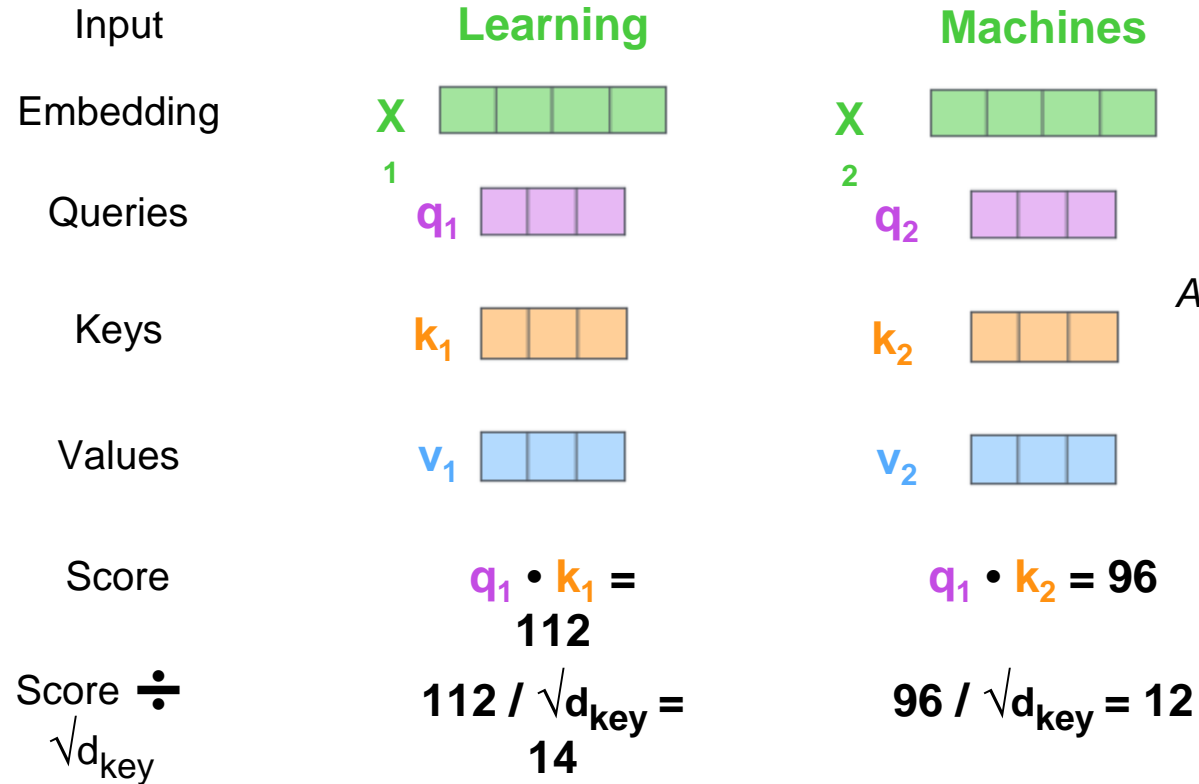
Classical Self-Attention

Computing
output for the
word *Learning*

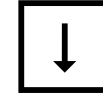


Classical Self-Attention

Computing output for the word *Learning*



- d_{key} = dimension of key vector
- Leads to more stable gradients
 - Hyperparameter (!!!)
 - Other values may be used

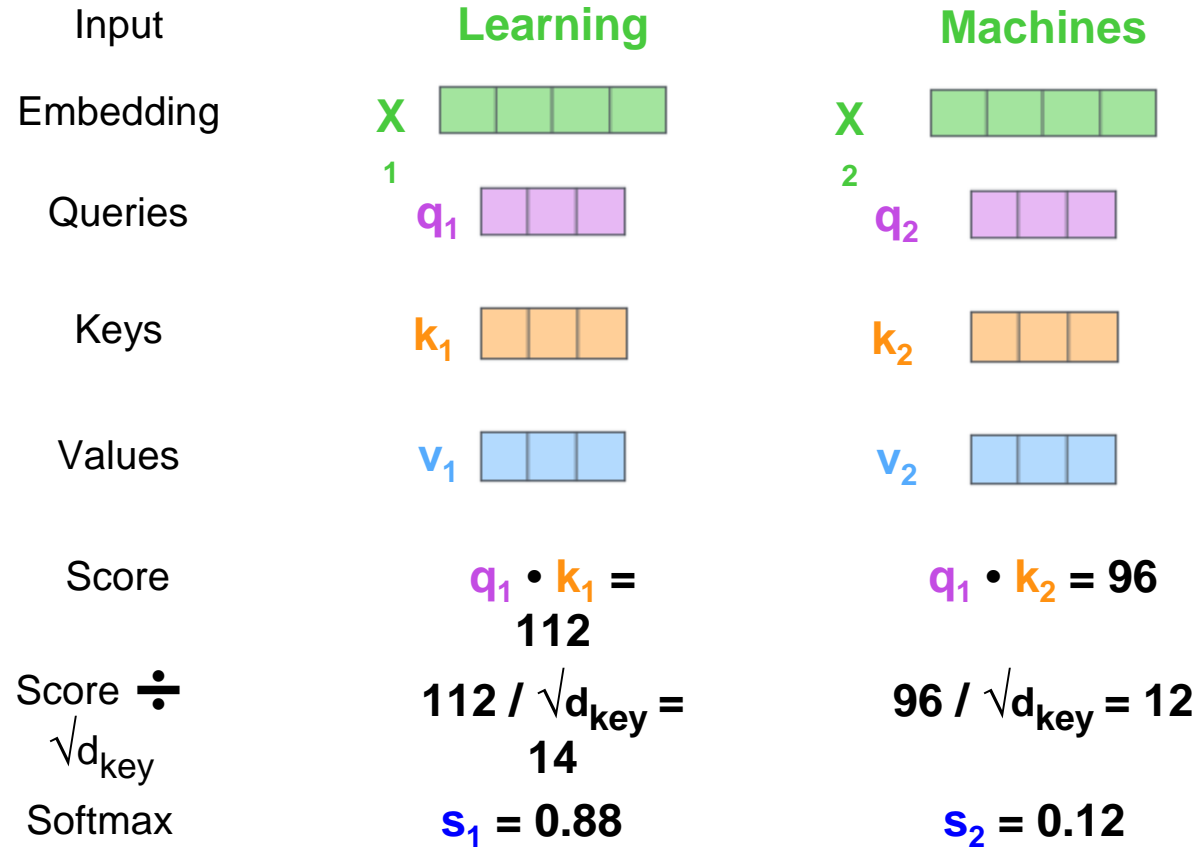


$$\sqrt{d_{\text{key}}} = \sqrt{64} = 8$$

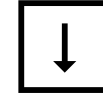
(8 is the value used in *Attention Is All You Need* (2017))

Classical Self-Attention

Computing output for the word *Learning*



- d_{key} = dimension of key vector
- Leads to more stable gradients
 - Hyperparameter (!!!)
 - Other values may be used

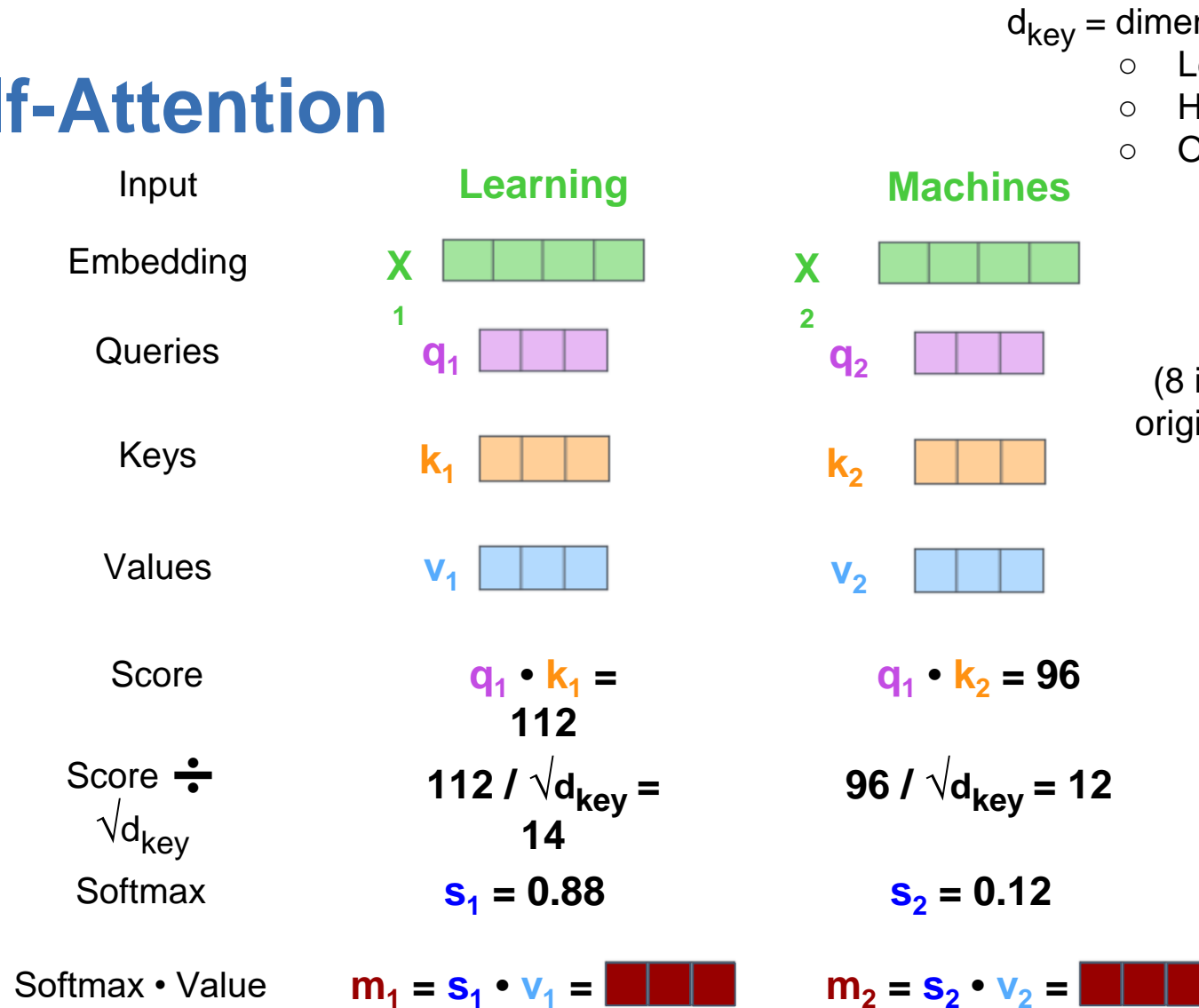


$$\sqrt{d_{\text{key}}} = \sqrt{64} = 8$$

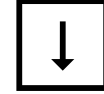
(8 is the value used in the original self-attention paper)

Classical Self-Attention

Computing output for the word *Learning*



- d_{key} = dimension of key vector
- Leads to more stable gradients
 - Hyperparameter (!!!)
 - Other values may be used



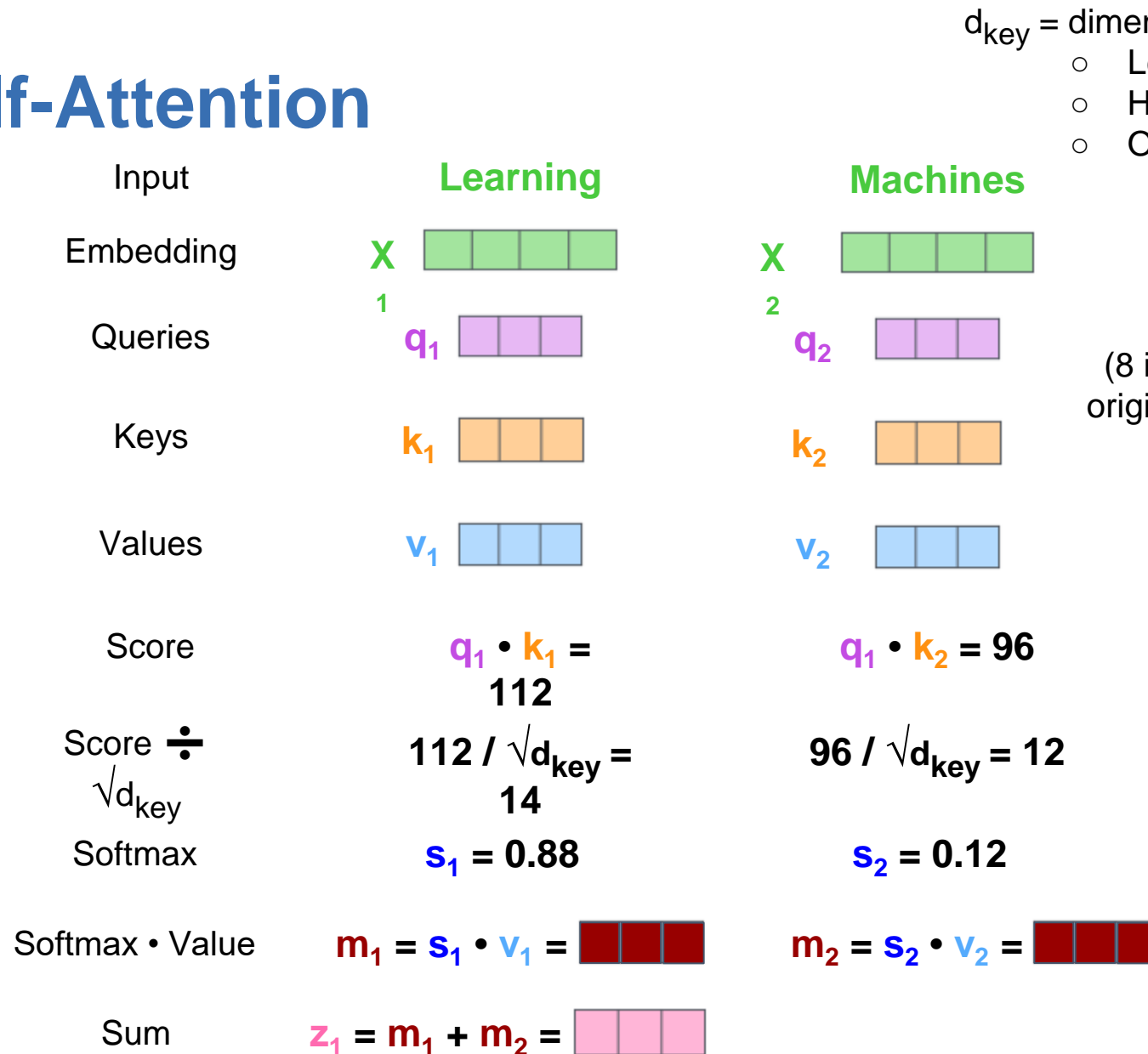
$$\sqrt{d_{\text{key}}} = \sqrt{64} = 8$$

(8 is the value used in the original self-attention paper)

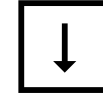


Classical Self-Attention

Computing output for the word *Learning*



- d_{key} = dimension of key vector
- Leads to more stable gradients
 - Hyperparameter (!!!)
 - Other values may be used

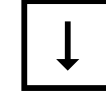


$$\sqrt{d_{\text{key}}} = \sqrt{64} = 8$$

(8 is the value used in the original self-attention paper)

Classical Self-Attention

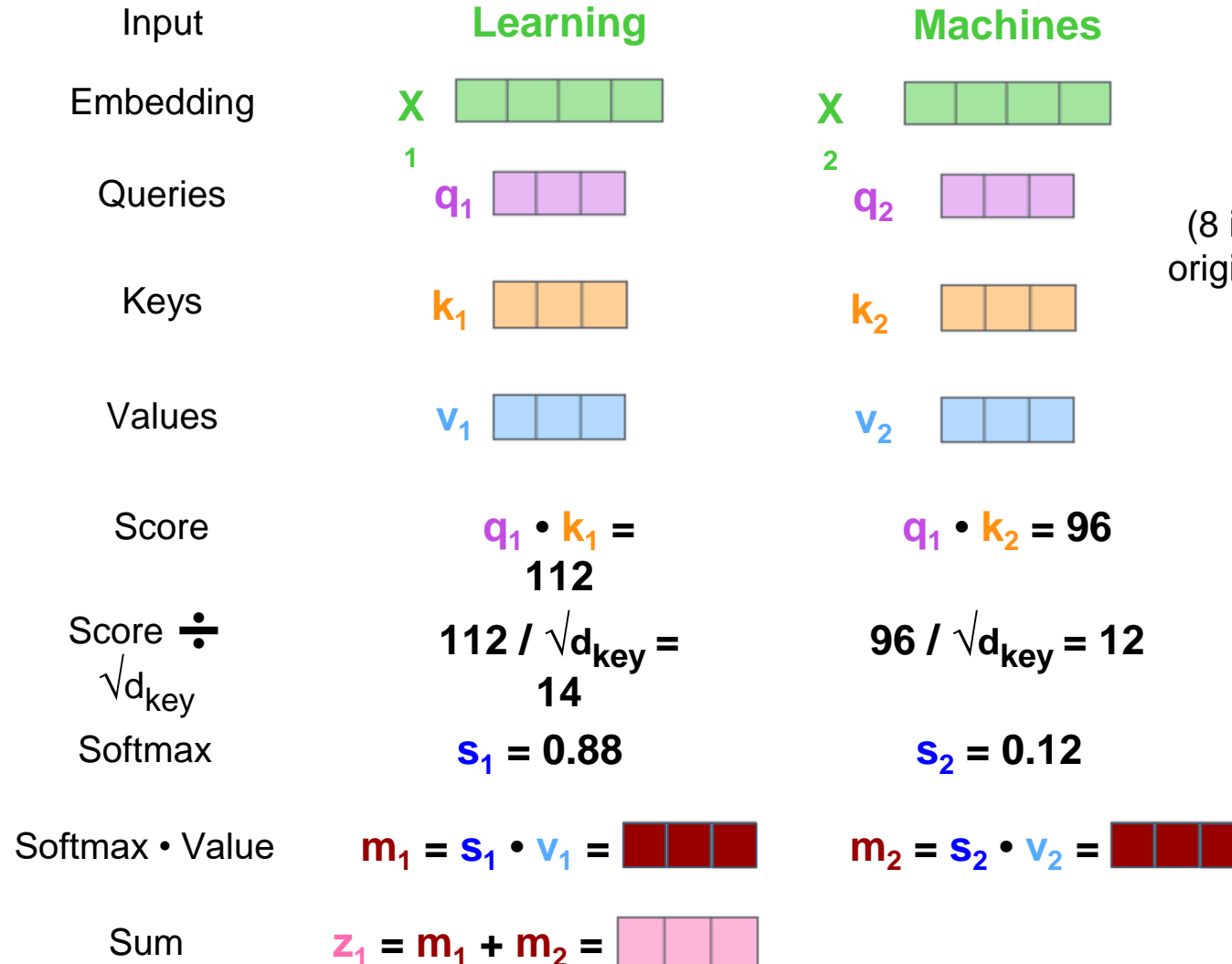
- d_{key} = dimension of key vector
- Leads to more stable gradients
 - Hyperparameter (!!!)
 - Other values may be used



$$\sqrt{d_{\text{key}}} = \sqrt{64} = 8$$

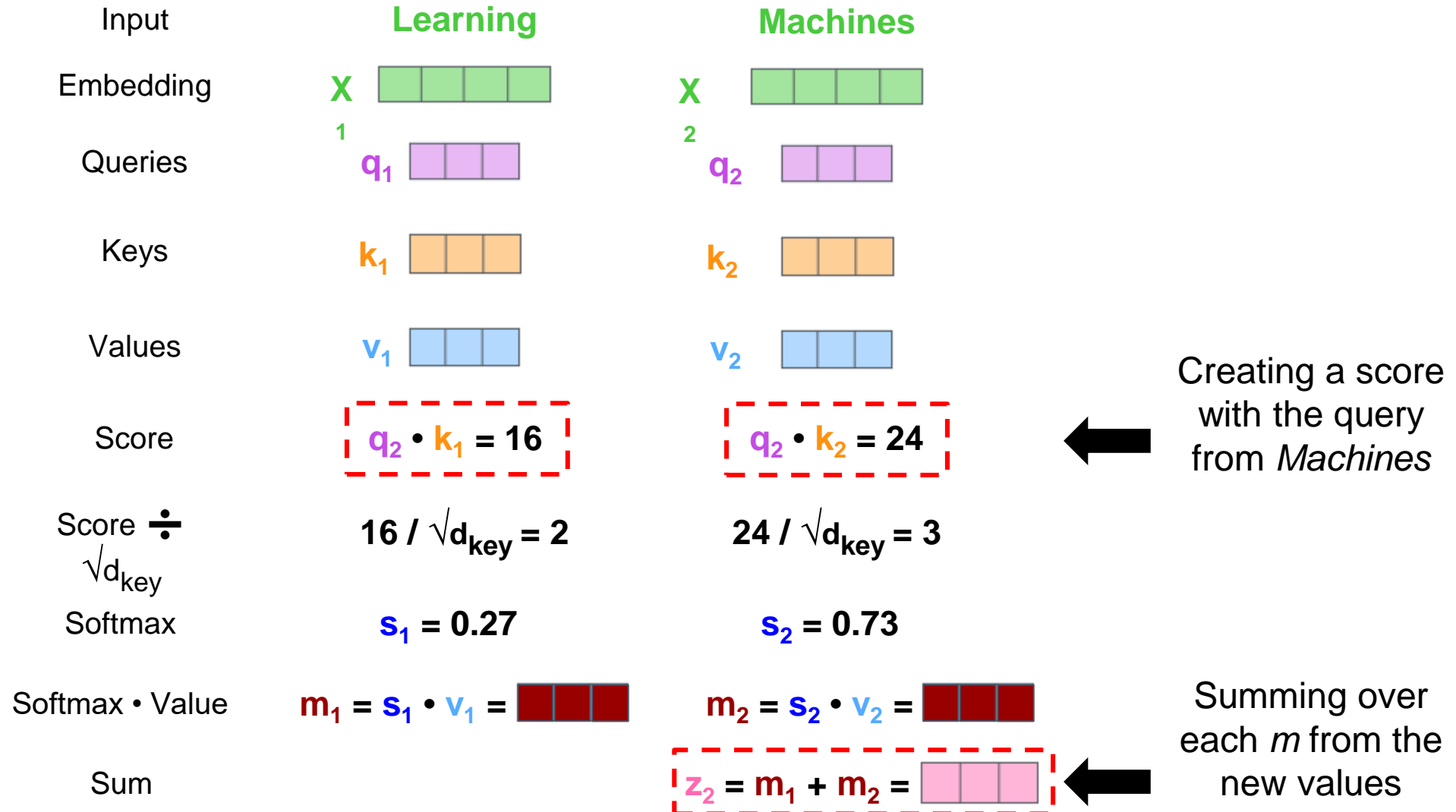
(8 is the value used in the original self-attention paper)

This is only for
the word
Learning!

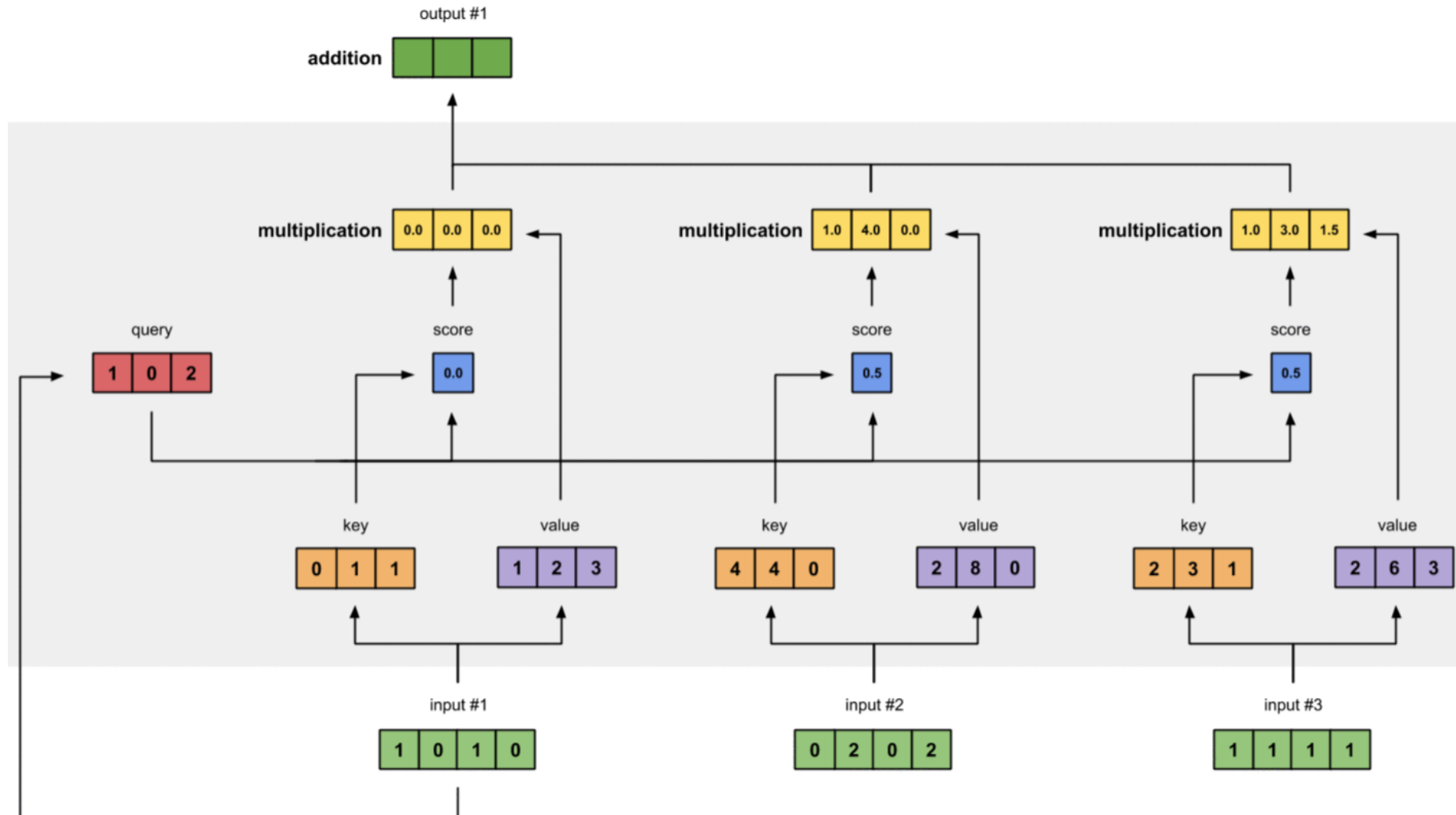


Classical Self-Attention

Doing this for
the word
Machines is
just as easy



Classical Self-Attention: Another Look

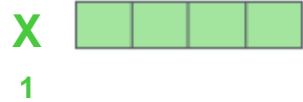


Quantum Self-Attention

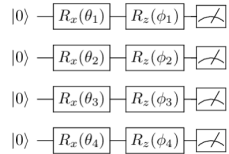
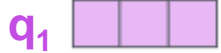
Input

Learning

Embedding



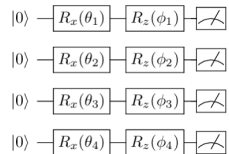
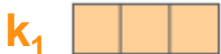
Queries



W^Q

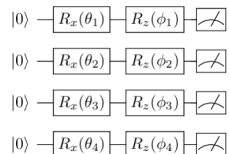
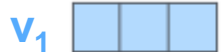
W^Q, W^K, W^V
learned
matrices
quantum
circuits!

Keys



W^K

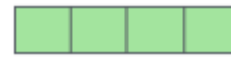
Values



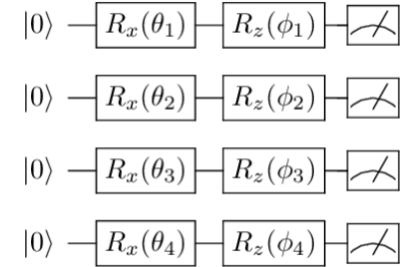
W^V

W^Q

X_1



\times

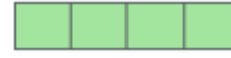


$=$

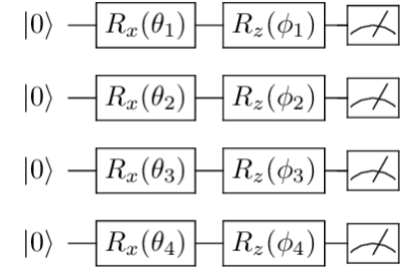


W^K

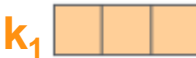
X_1



\times

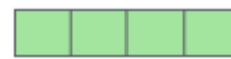


$=$

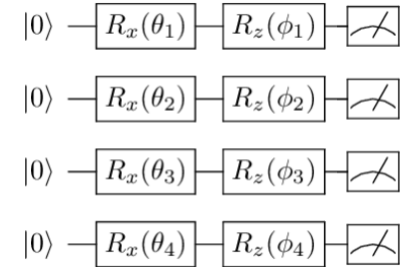


W^V

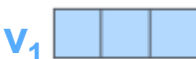
X_1



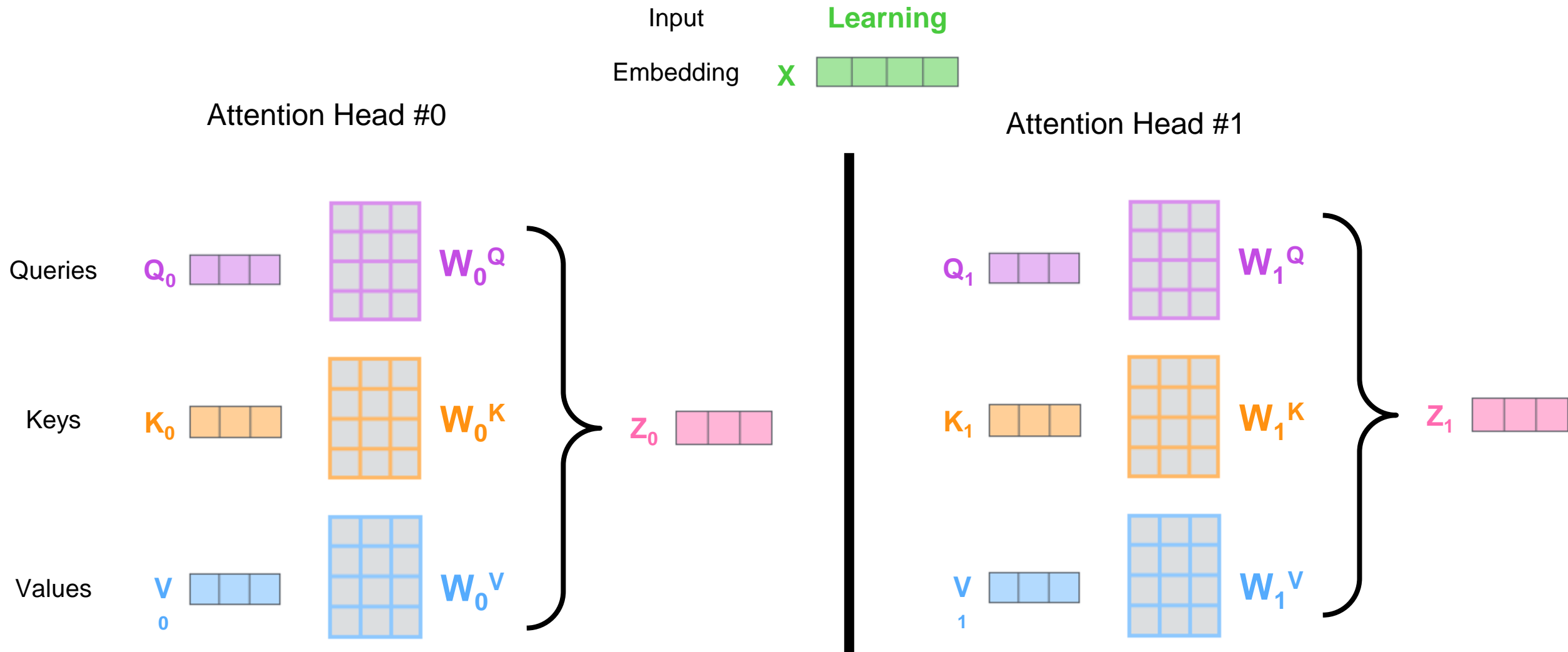
\times



$=$

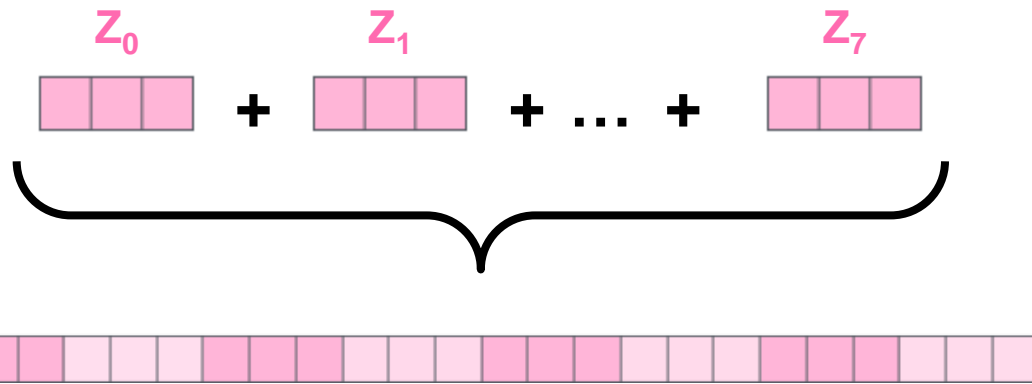


Multi-Head Attention

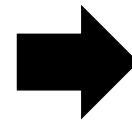


Multi-Head Attention

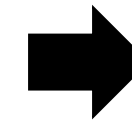
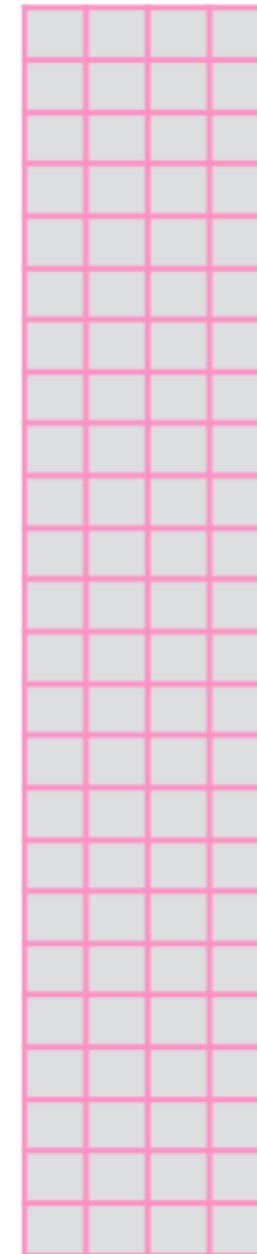
Concatenate Z from each attention head:



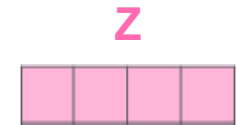
Multiply concatenated Z 's with learned matrix W^0



W^0



Resulting matrix Z contains information from all attention heads

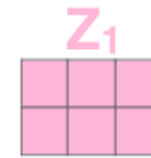
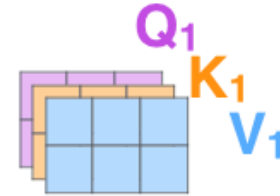
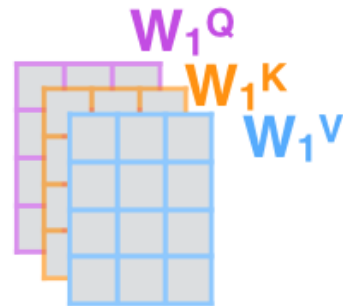
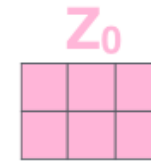
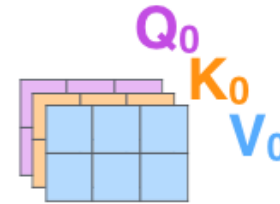
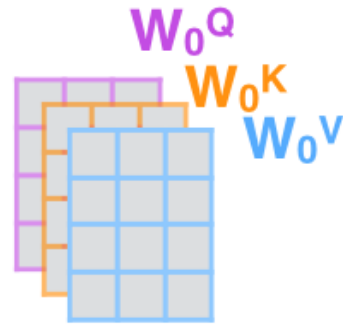


Efficient Multi-Head Attention

Stacking words results in a larger matrix

Allows for representing each input as a (larger) matrix

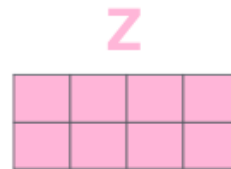
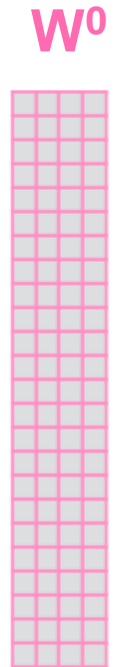
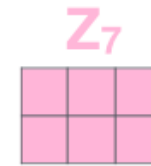
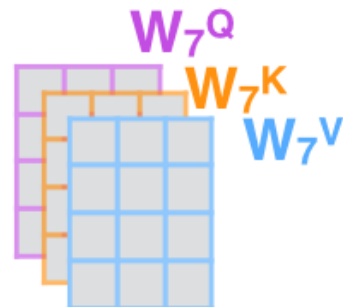
Learning Machines



...

...

...



CERN Quantum Technology Initiative

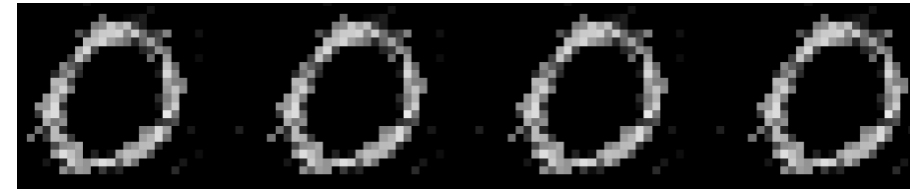
Accelerating Quantum Technology Research and Applications

Thank You!

In < 50 epochs, we get
(preliminary) results
with MNIST dataset!

Dr. Sofia Vallecorsa

Dr. Michele Grossi



Computational Resources:



Questions?



QUANTUM
TECHNOLOGY
INITIATIVE

