

Deep Learning for Computer Vision (Discriminative way)

Andrii Liubonko
2024

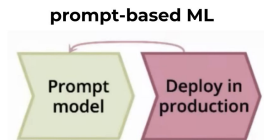
multimodal LLMs

CV foundational models

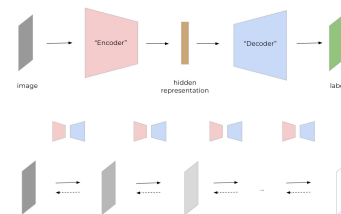
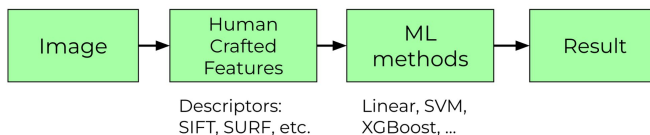
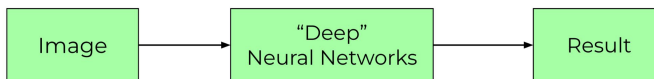
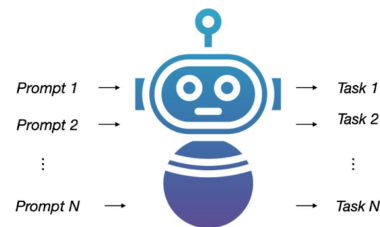
Deep Learning based CV

“classical” CV

math



CLIP, DINO, SAM, ...



- **discriminative** model $\Pr(\mathbf{w}|\mathbf{x})$
- **generative** model $\Pr(\mathbf{x}|\mathbf{w})$
- $\Pr(\mathbf{w}|\mathbf{x}) = \Pr(\mathbf{x}|\mathbf{w}) * \Pr(\mathbf{w}) / \int [\Pr(\mathbf{x} | \mathbf{w}) * \Pr(\mathbf{w})] d\mathbf{w}$

multimodal LLMs

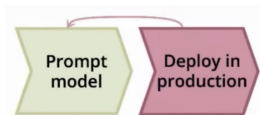
CV foundational models

Deep Learning based CV

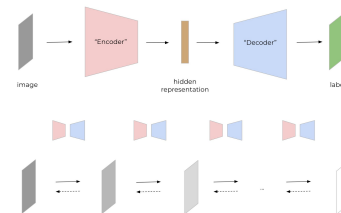
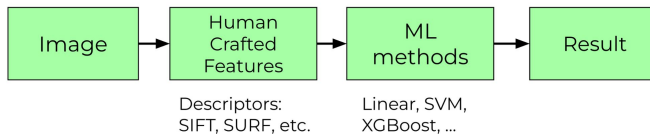
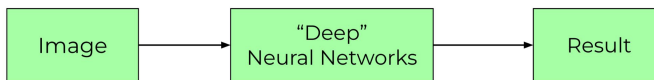
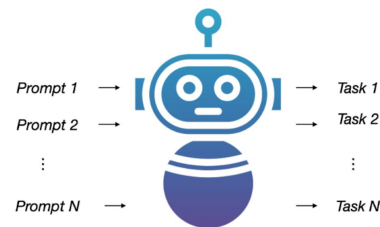
“classical” CV

math

prompt-based ML

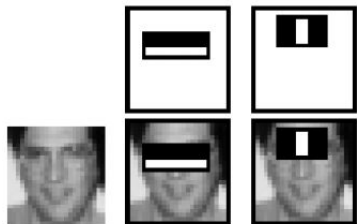


CLIP, DINO, SAM, ...



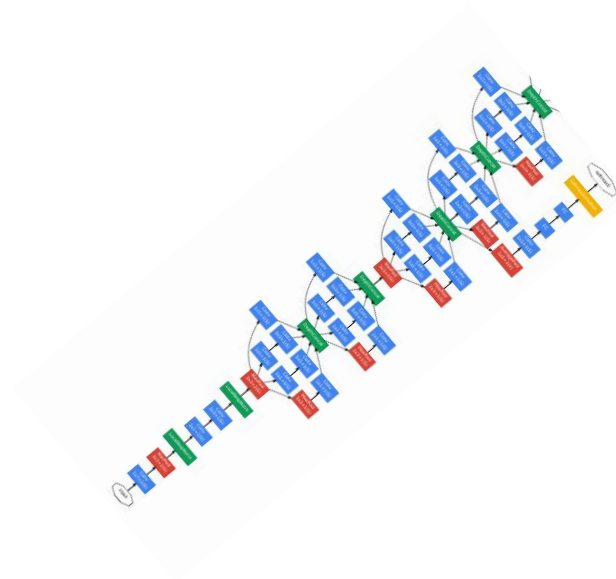
- **discriminative** model $\Pr(\mathbf{w}|\mathbf{x})$
- **generative** model $\Pr(\mathbf{x}|\mathbf{w})$
- $\Pr(\mathbf{w}|\mathbf{x}) = \Pr(\mathbf{x}|\mathbf{w}) * \Pr(\mathbf{w}) / \int [\Pr(\mathbf{x} | \mathbf{w}) * \Pr(\mathbf{w})] d\mathbf{w}$

Intro



Era of
Human-Crafter
Features

2012
AlexNet



Era of
Deep Learning

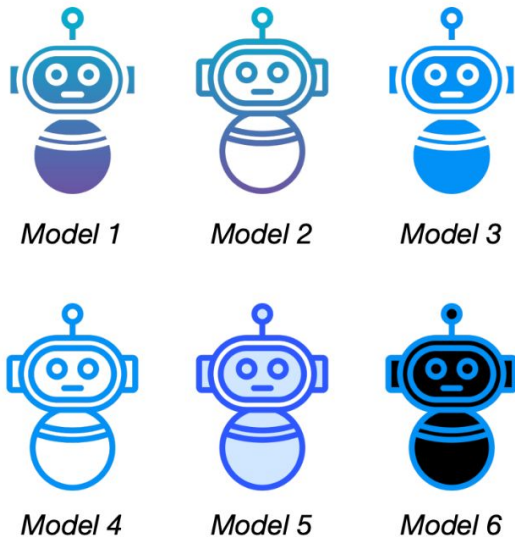
2022
ChatGPT



Era of
LLMs
(FoundM)

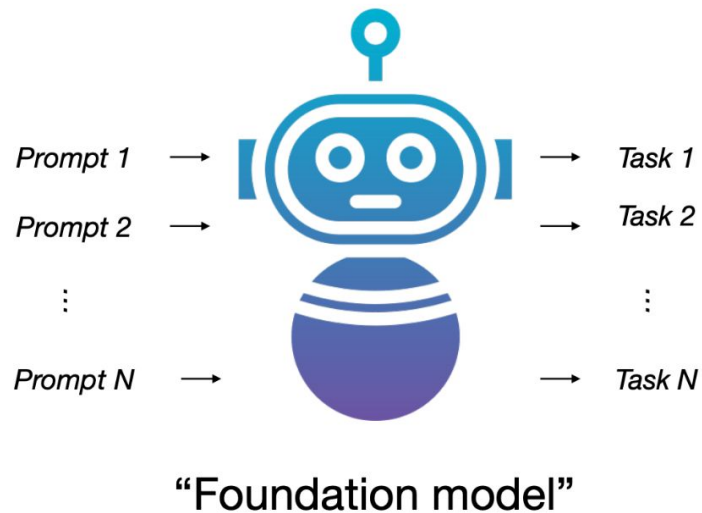
traditional ML

Old days: one model for one purpose

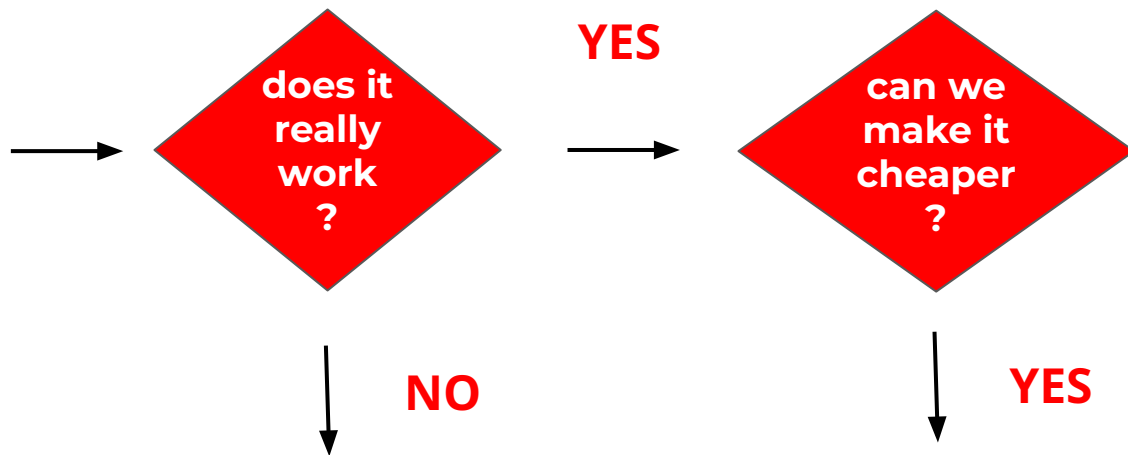
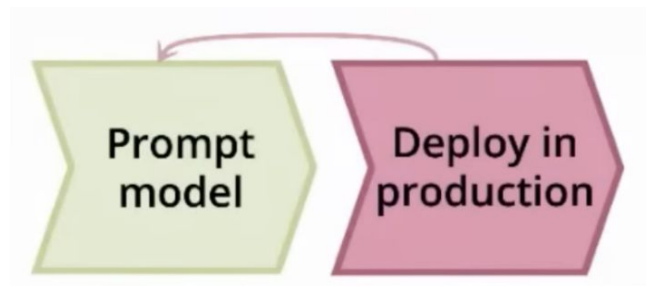


prompt-based ML

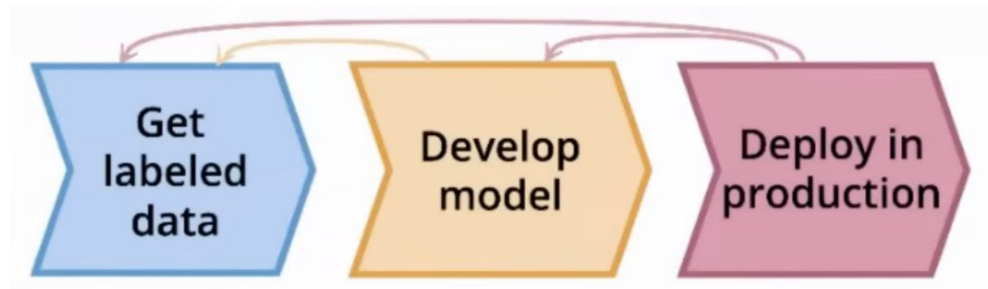
Now: one model for multiple purposes



prompt-based ML



traditional ML



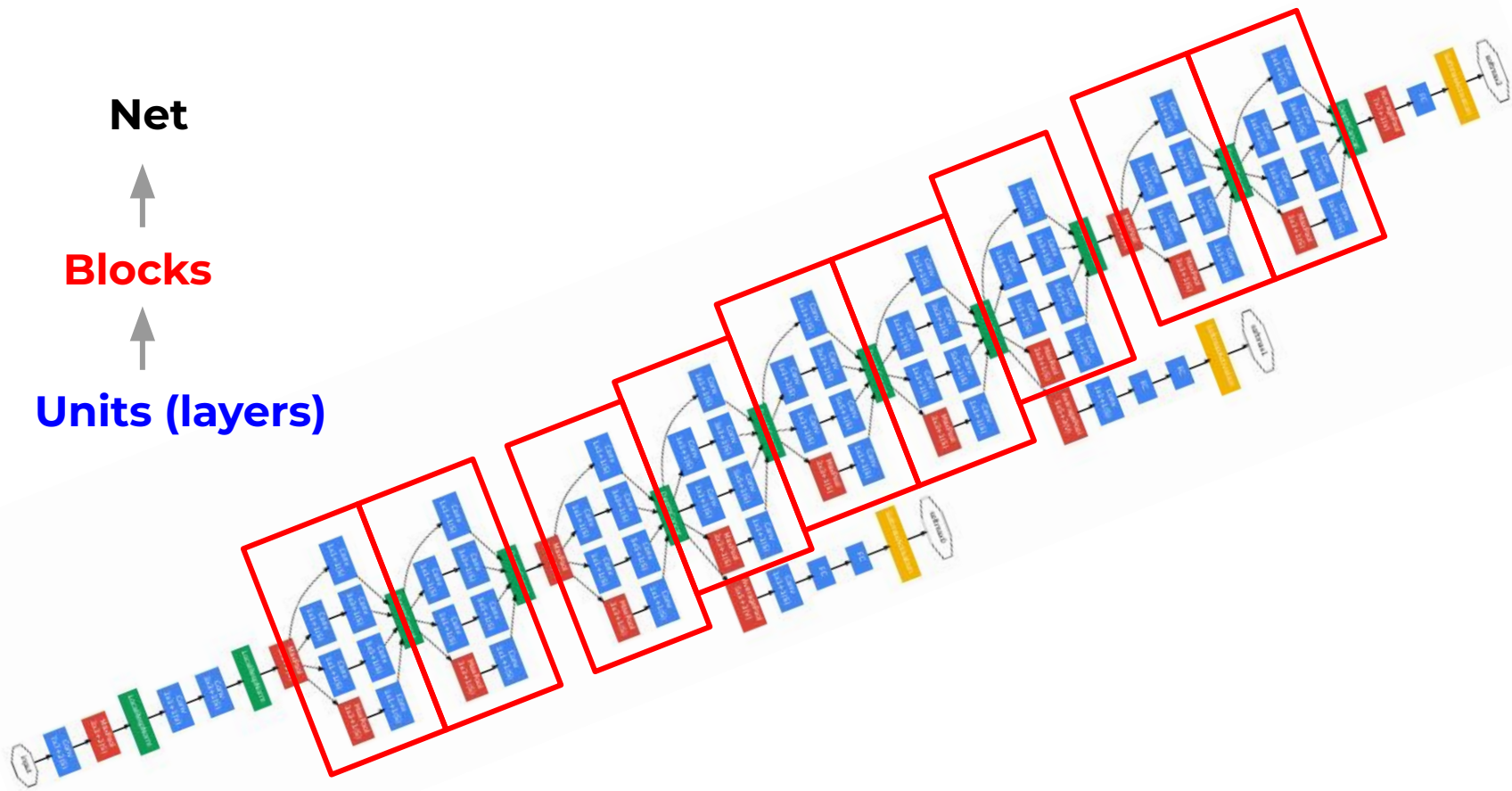
Net



Blocks



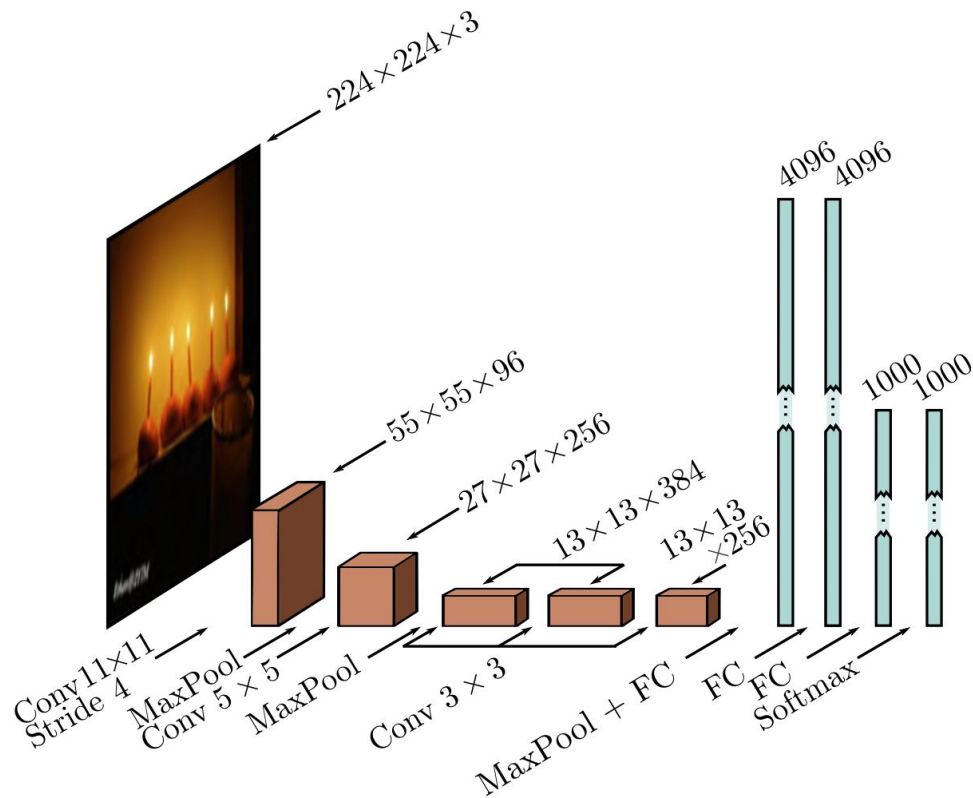
Units (layers)



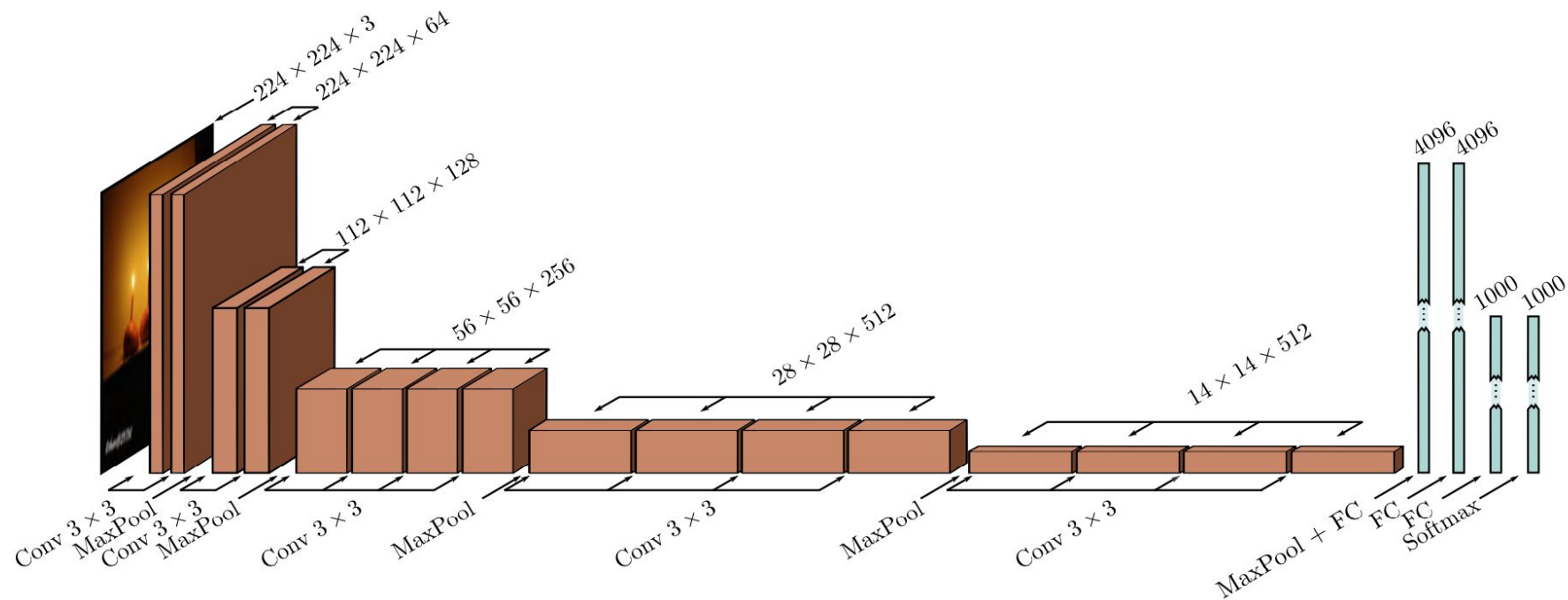
- Layers
 - Layers [Convolution] [Receptive field]
 - Layers [Dilated Convolution, Deformable Convolution]
 - Layers [Group Convolutions and its variants]
 - Layers [Upsampling, Learnable Upsampling]
 - Layers [Normalization, Batch Norm, Dropout]
- Blocks
 - VGG, Inception
 - ResNet*
 - MobileNet*
- Architectures

AlexNet, VGG, Inception,
ResNet, MobileNet, EfficientNet

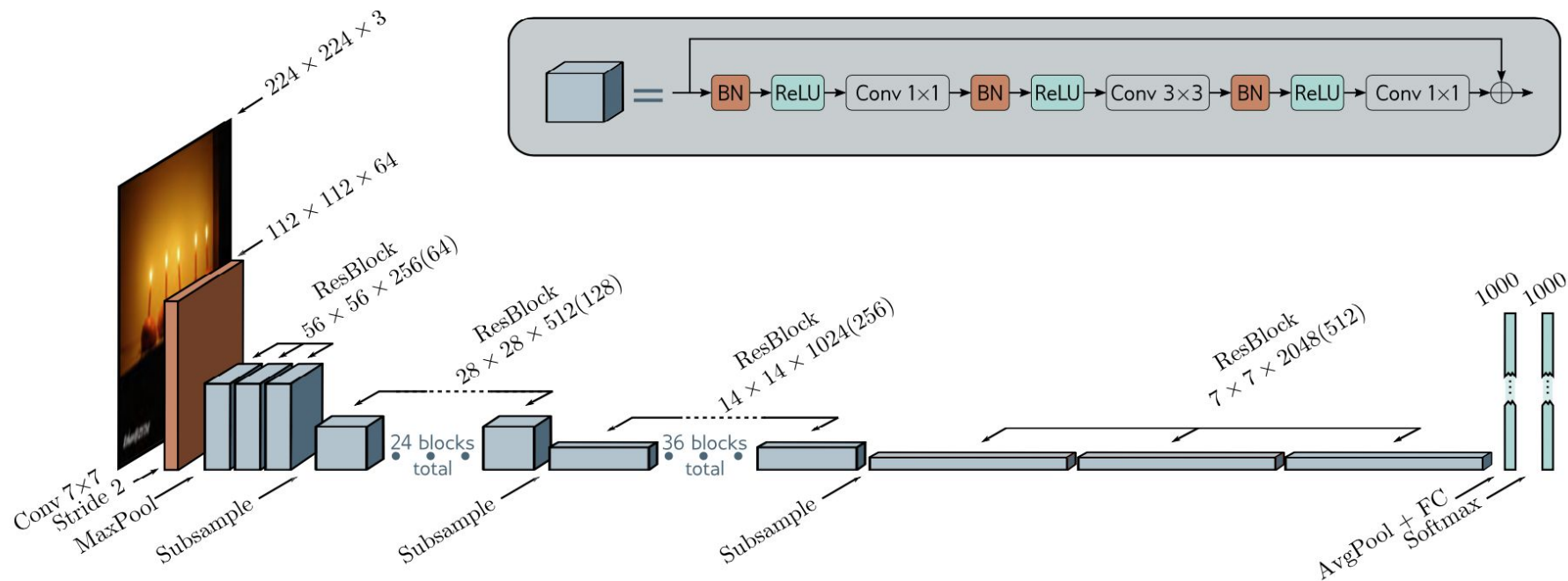
Net [AlexNet]



Net [VGG16]



Net [ResNet]



Content of today's lecture

- Attention Intro
 - Brief History
 - Seq2seq Attention
 - Self-Attention
 - Transformer
- Attention in CV
 - SE
 - ViT
 - latest developments

Attention Timeline

Attention Intro

Seq2Seq

- [Cho et al. \(2014\)](#)
- [Sutskever et al. \(2014\)](#)

Align & Translate

- [Bahdanau et al. \(2015\)](#)
- [Luong et al. \(2015\)](#)

Self attention / Transformer

- [Vaswani et al \(2017\)](#)

BERT

- [Devlin et al \(2018\)](#)

GPT-3

- [Brown et al \(2020\)](#)

GPT-4

- [OpenAI \(2023\)](#)

2014 2015 2016 2017 2018 2019 2020 2021 2022 2023

Attention in CV

Visual attention

- [Xu et al. \(2015\)](#)

SE networks

- [Hu et al. \(2018\)](#)

DETR

- [Carion et al.\(2020\)](#)

ViT

- [Dosovitskiy et al.\(2020\)](#)

MAE

- [He et al.](#)

SAM

- [Meta](#)

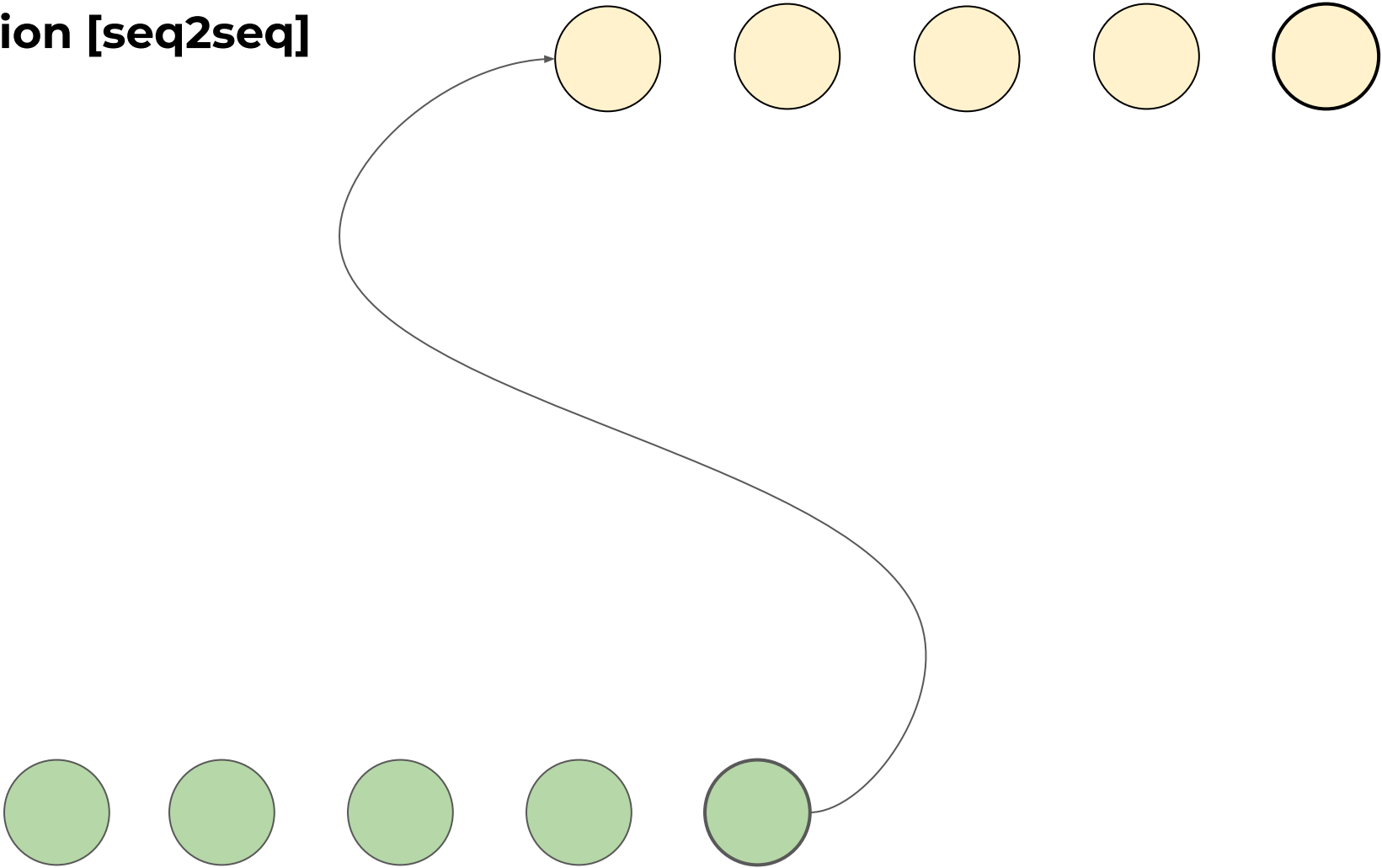
DINO

- [Caron et al.](#)

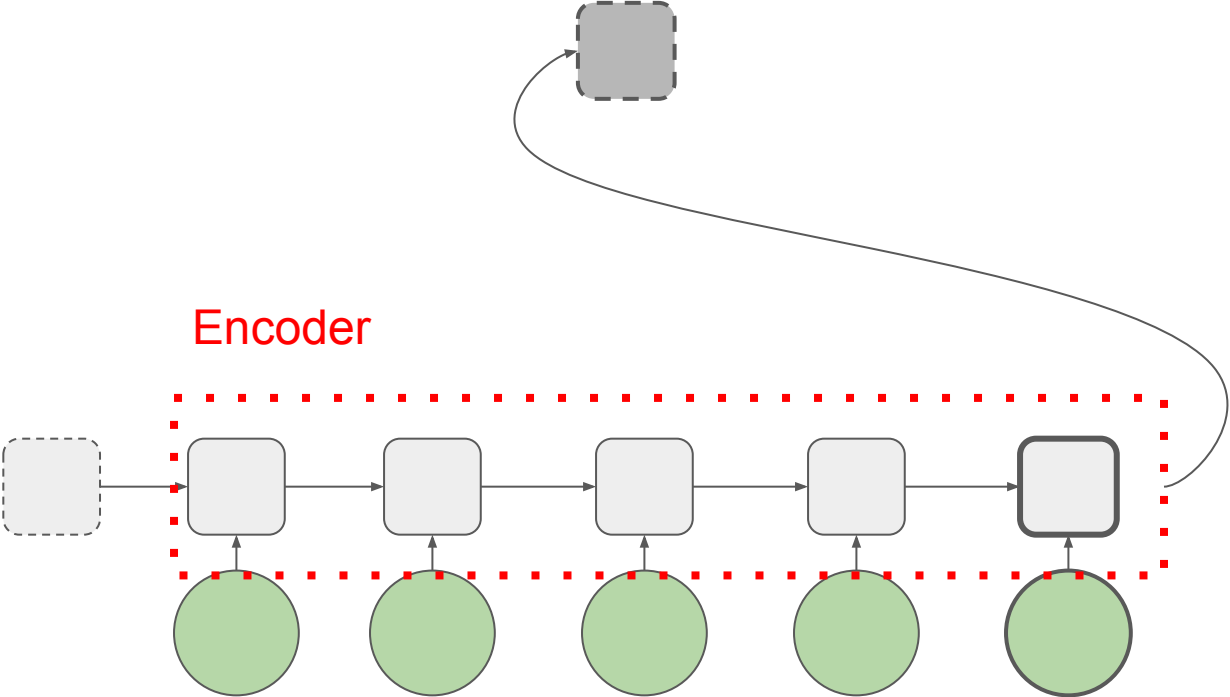
Attention [intro]



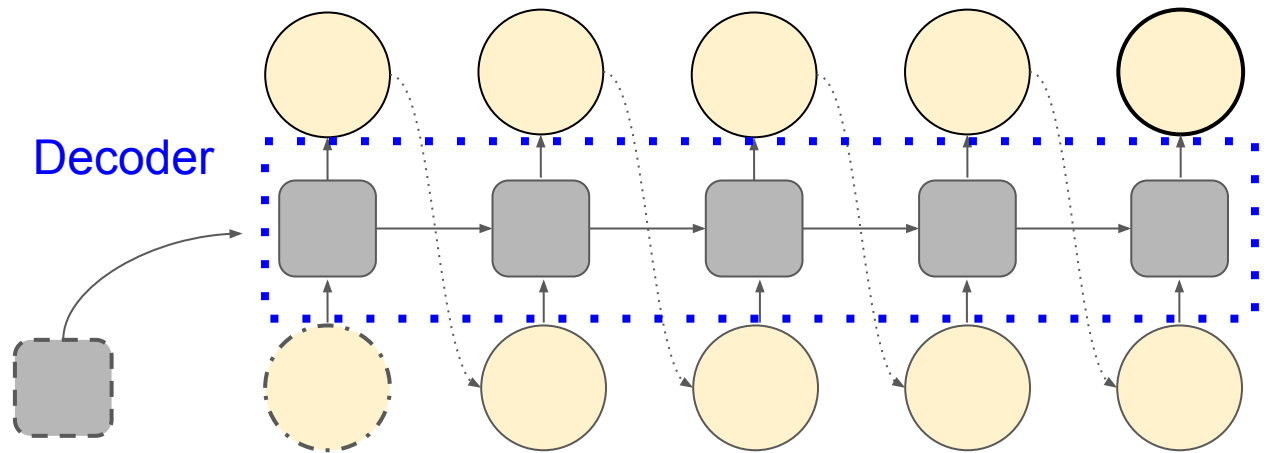
Attention [seq2seq]



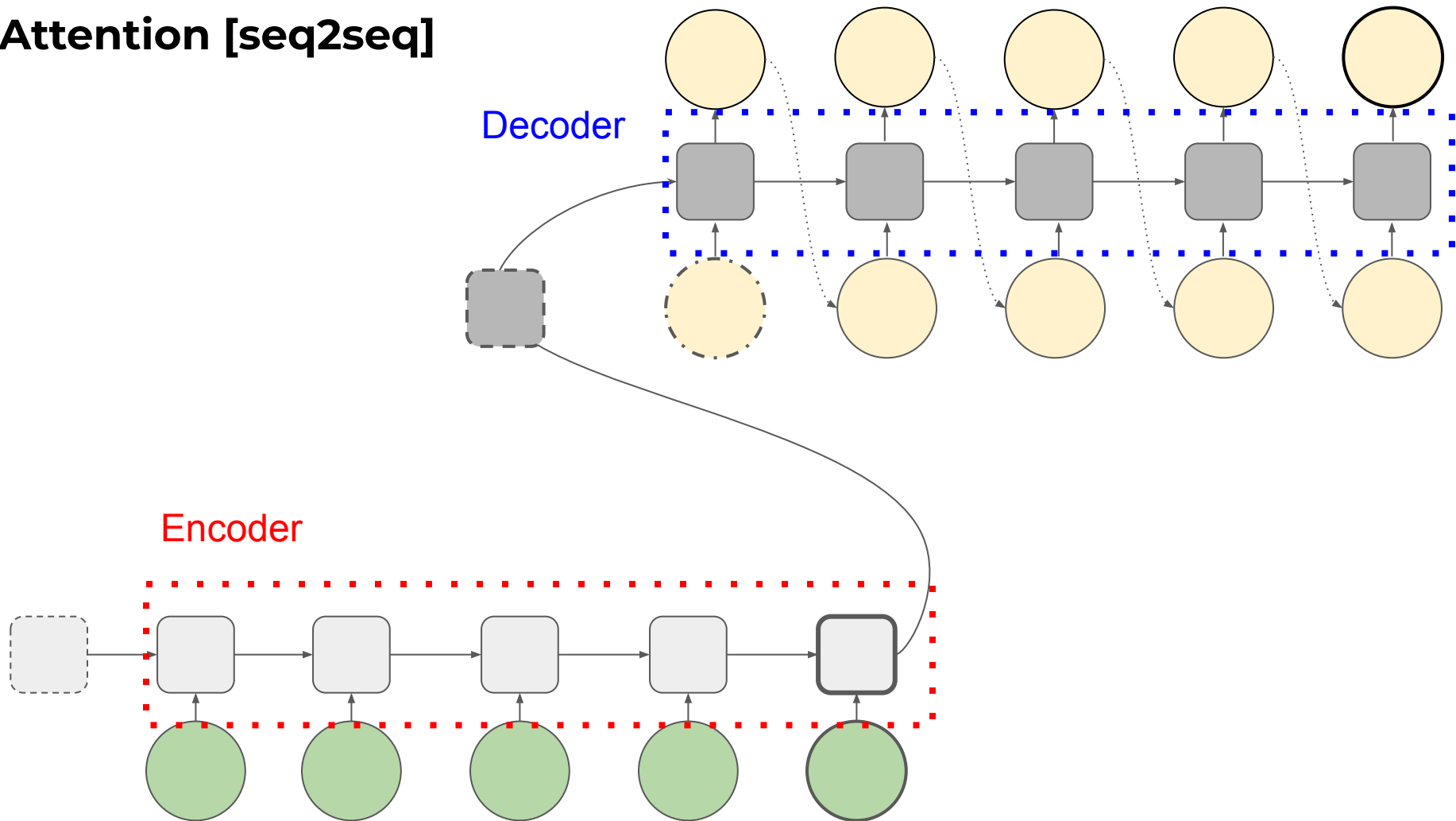
Attention [seq2seq]



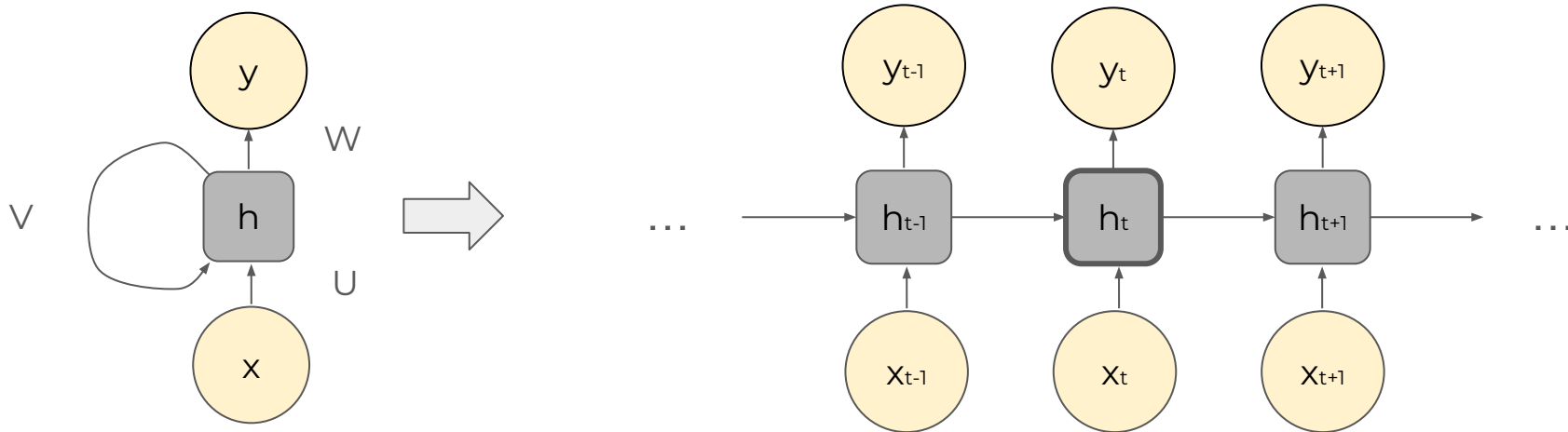
Attention [seq2seq]



Attention [seq2seq]



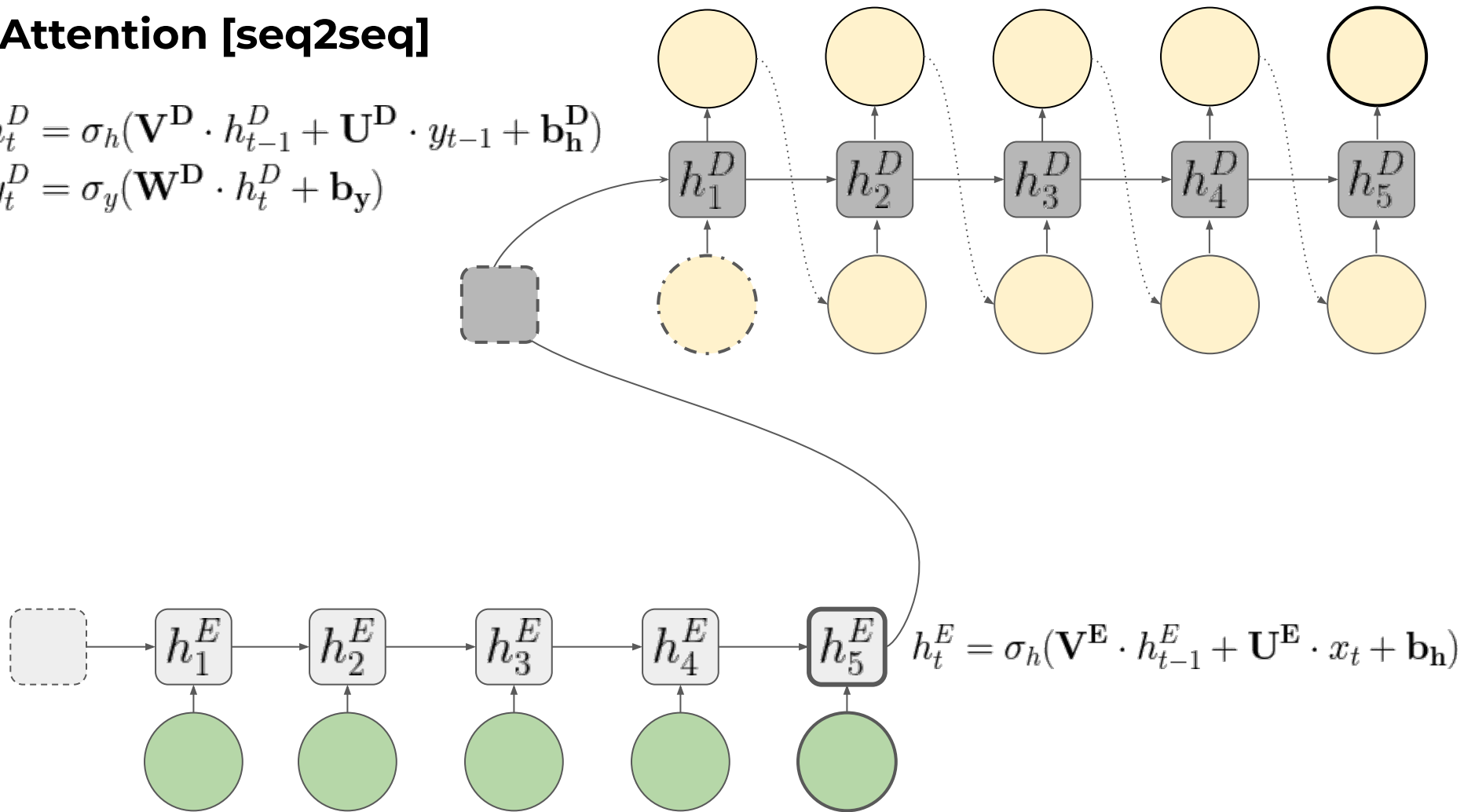
Attention [seq2seq]



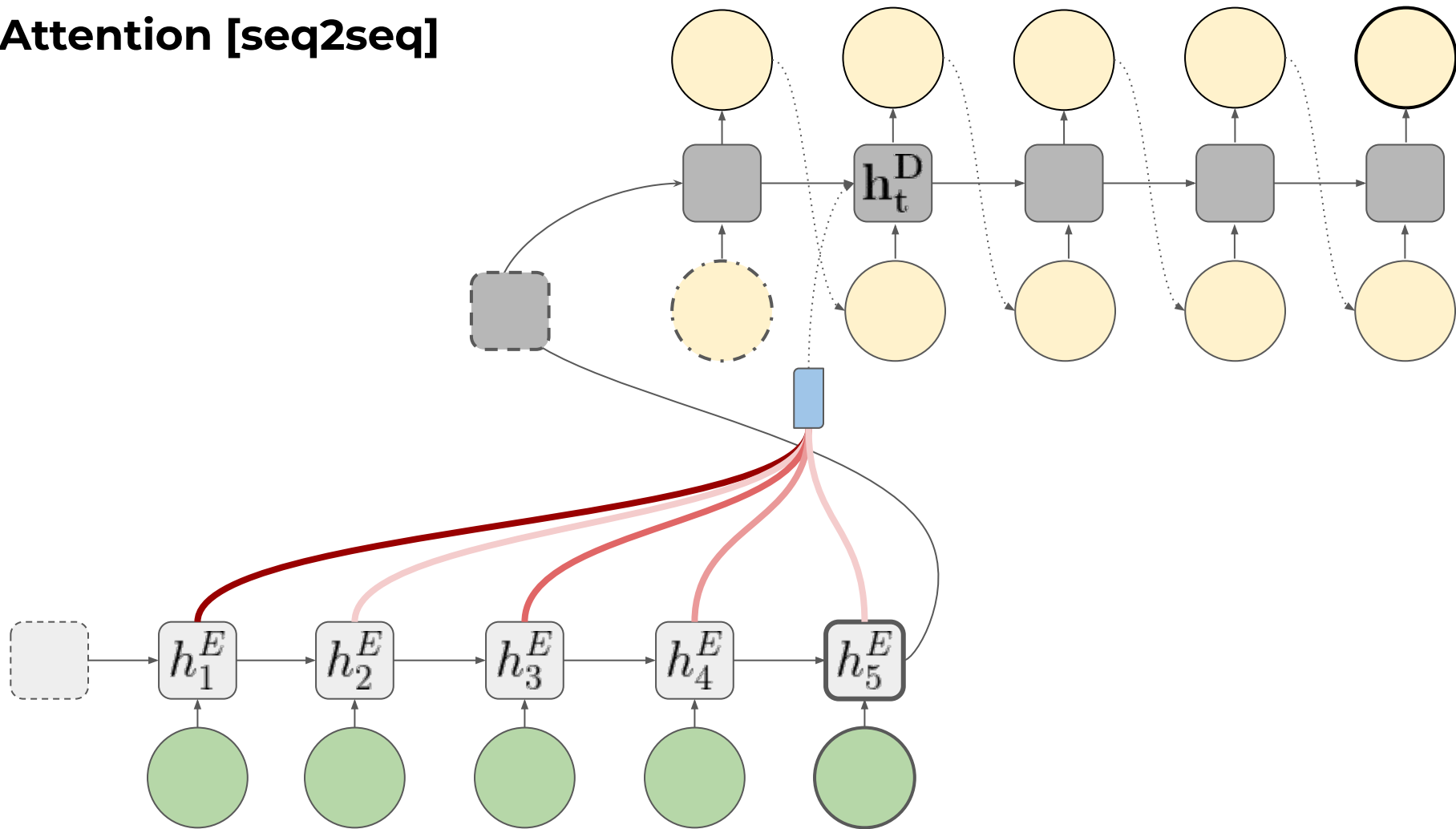
Attention [seq2seq]

$$h_t^D = \sigma_h(\mathbf{V}^D \cdot h_{t-1}^D + \mathbf{U}^D \cdot y_{t-1} + \mathbf{b}_h^D)$$

$$y_t^D = \sigma_y(\mathbf{W}^D \cdot h_t^D + \mathbf{b}_y)$$



Attention [seq2seq]



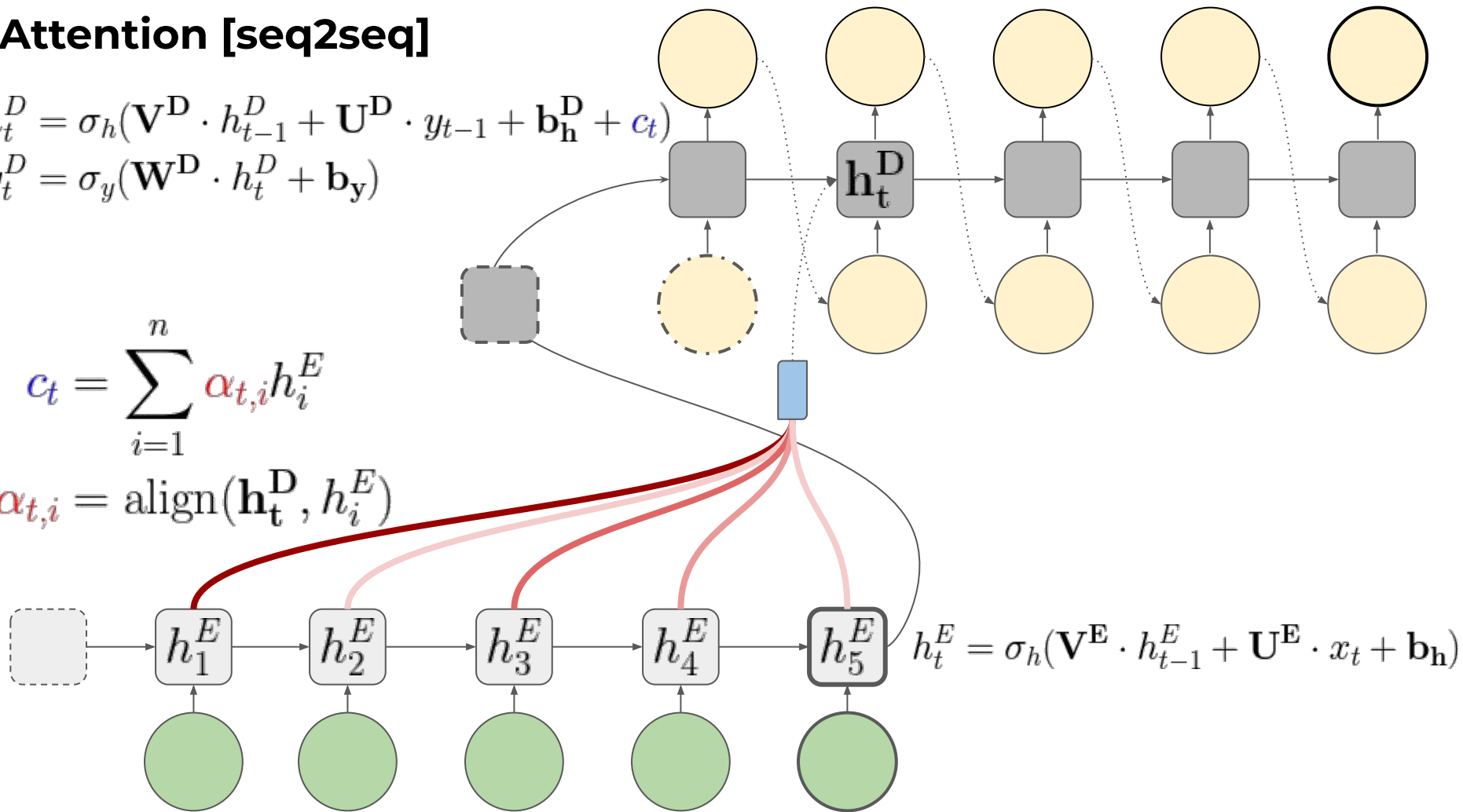
Attention [seq2seq]

$$h_t^D = \sigma_h(\mathbf{V}^D \cdot h_{t-1}^D + \mathbf{U}^D \cdot y_{t-1} + \mathbf{b}_h^D + \mathbf{c}_t)$$

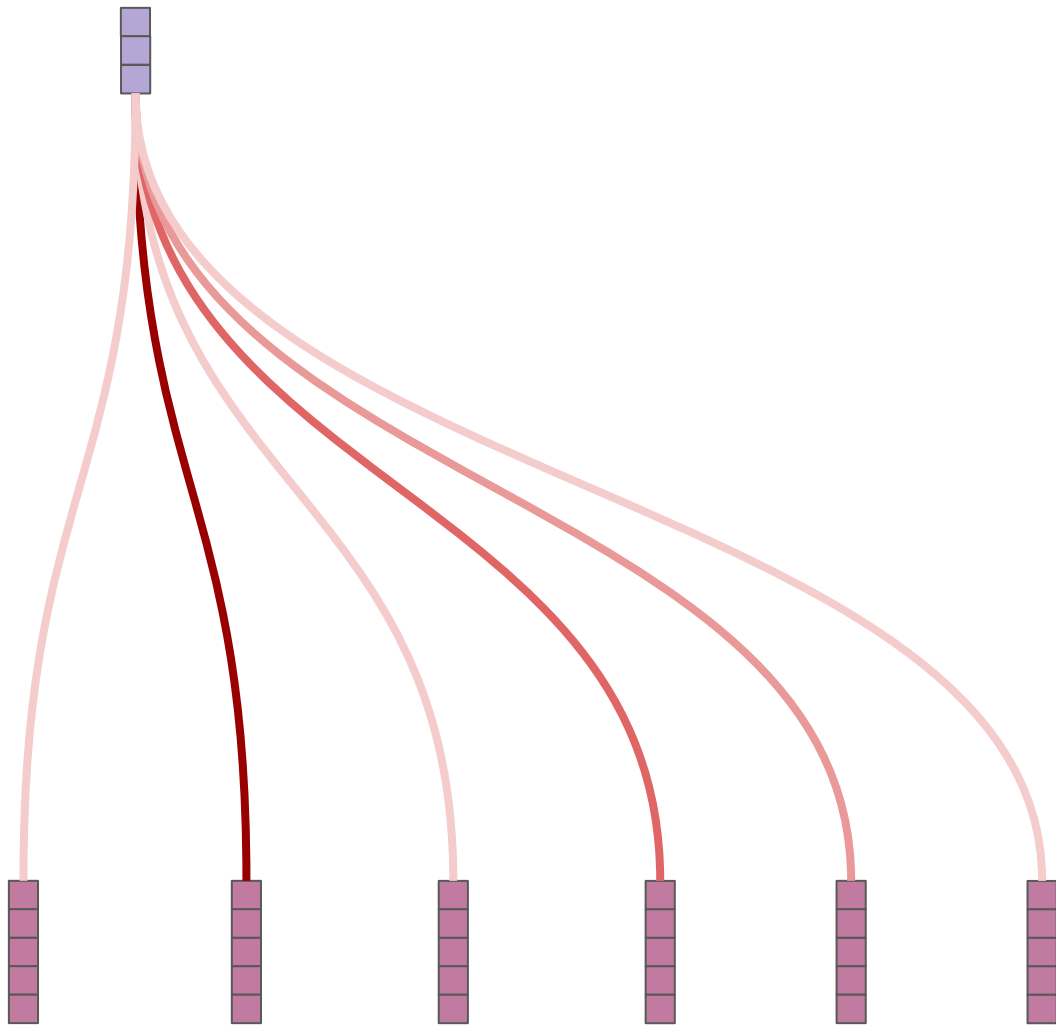
$$y_t^D = \sigma_y(\mathbf{W}^D \cdot h_t^D + \mathbf{b}_y)$$

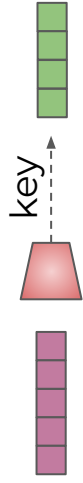
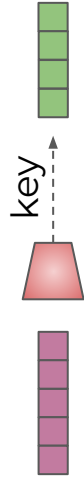
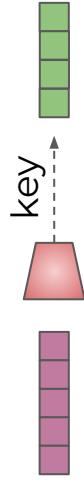
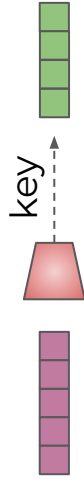
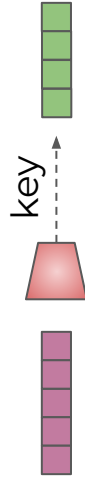
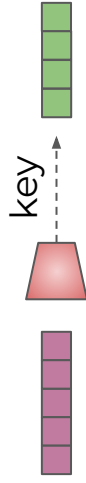
$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} h_i^E$$

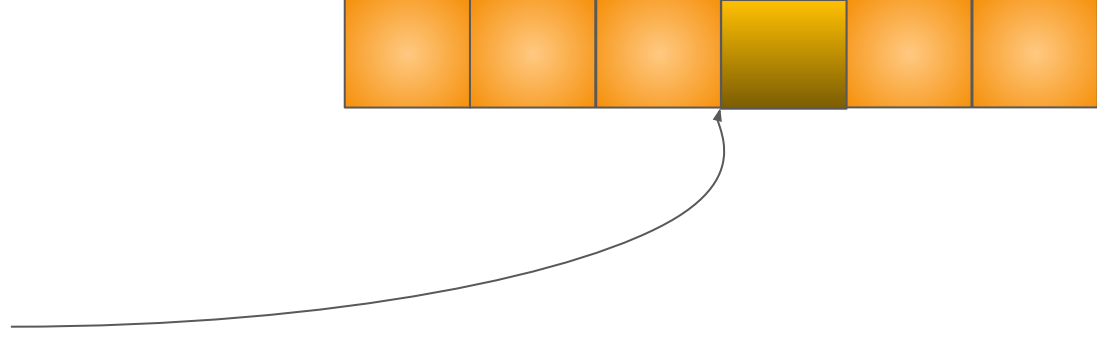
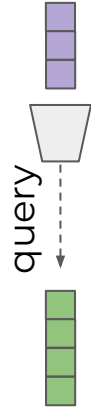
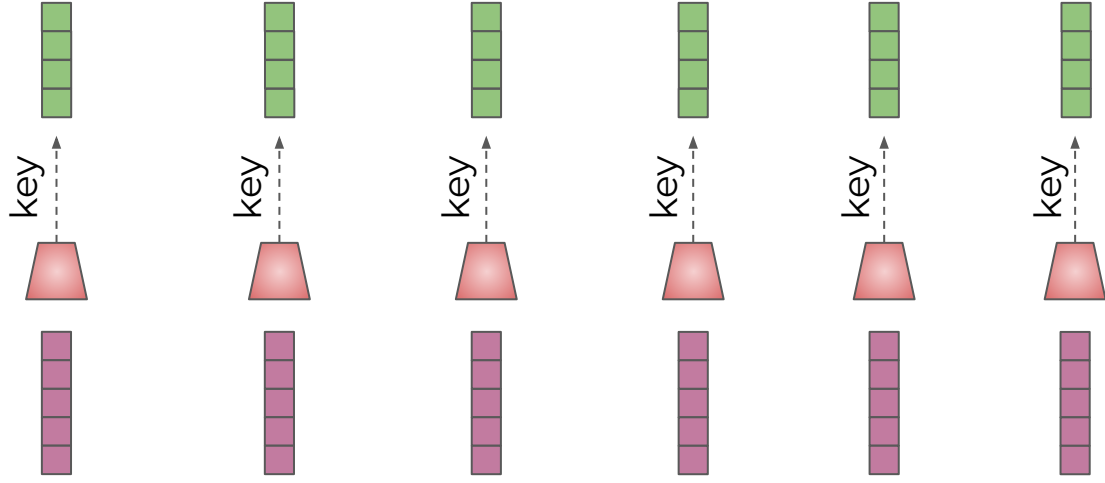
$$\alpha_{t,i} = \text{align}(\mathbf{h}_t^D, h_i^E)$$

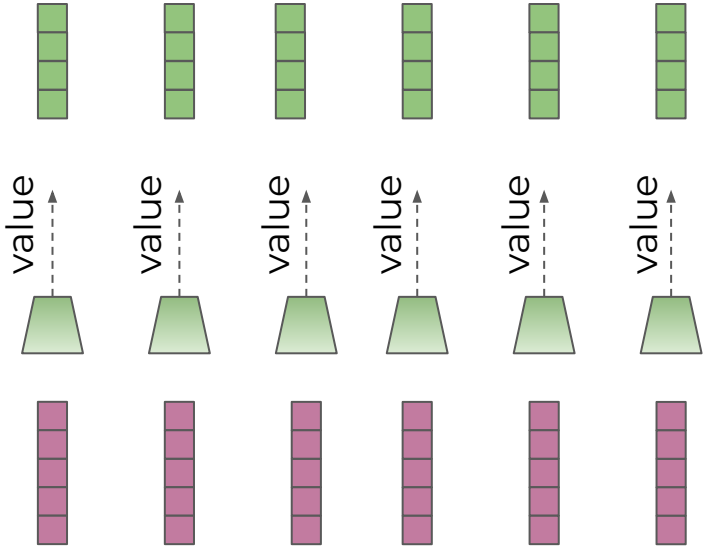


query-key-value abstraction

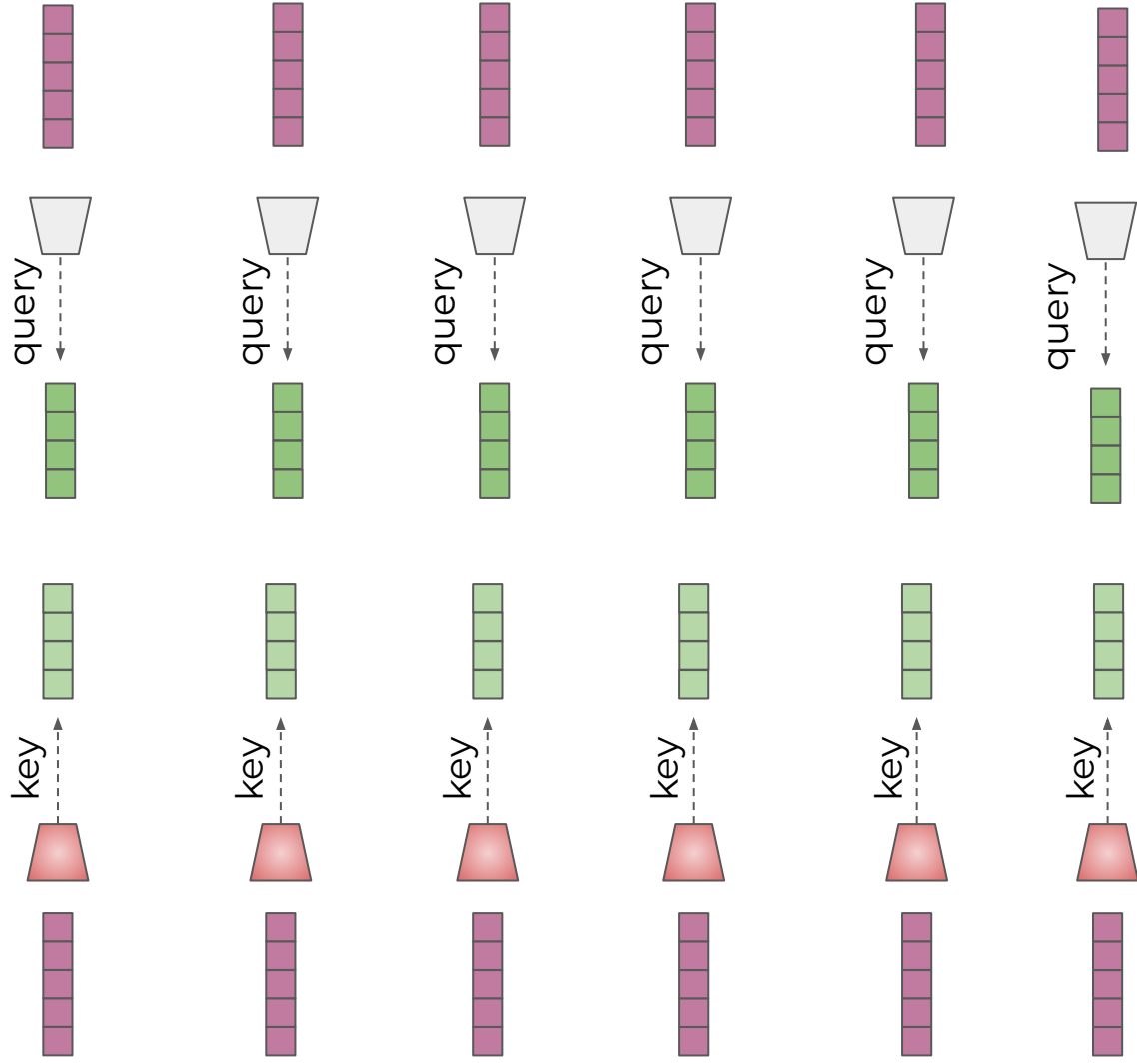


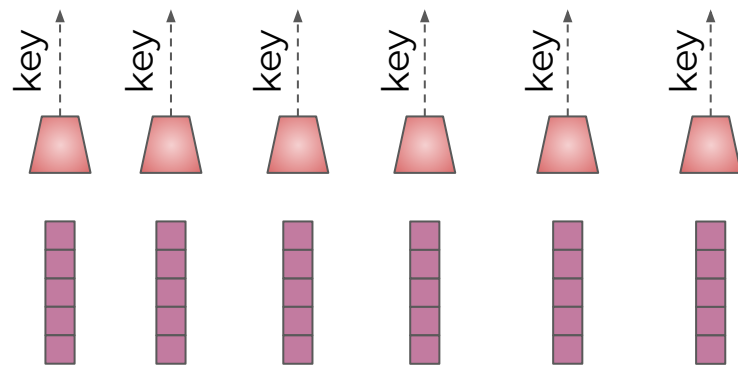
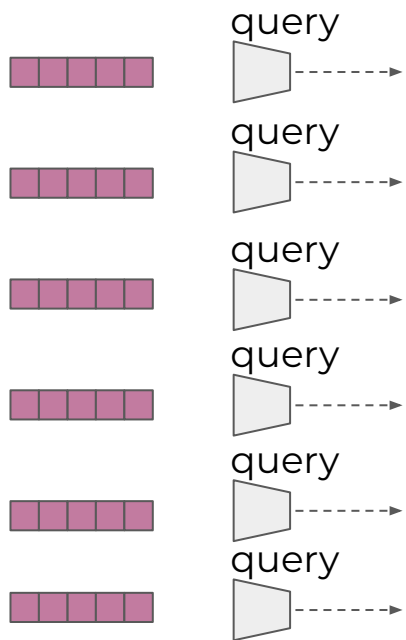






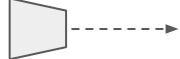
Self-attention



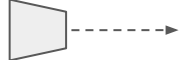




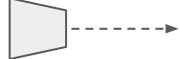
query



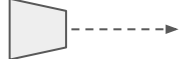
query



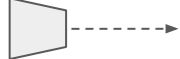
query



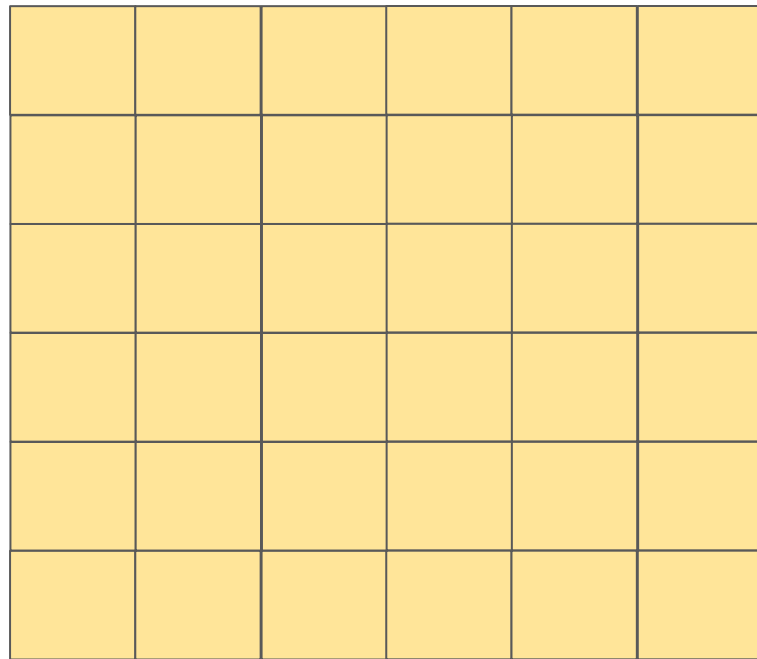
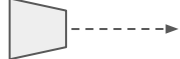
query



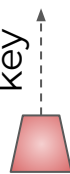
query



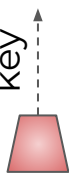
query



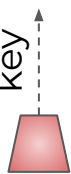
key



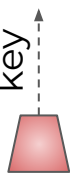
key



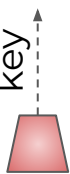
key



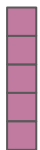
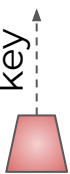
key

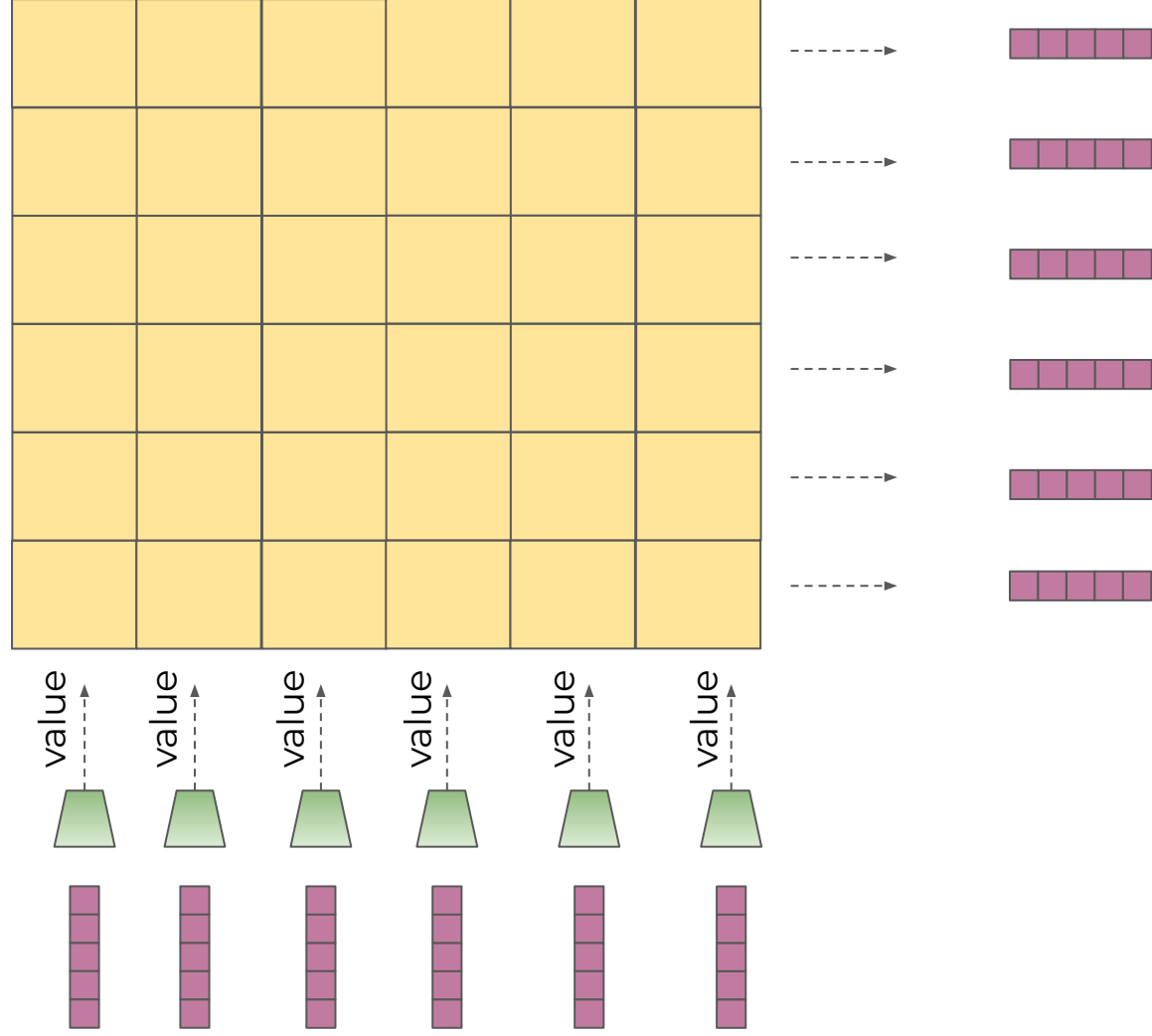


key

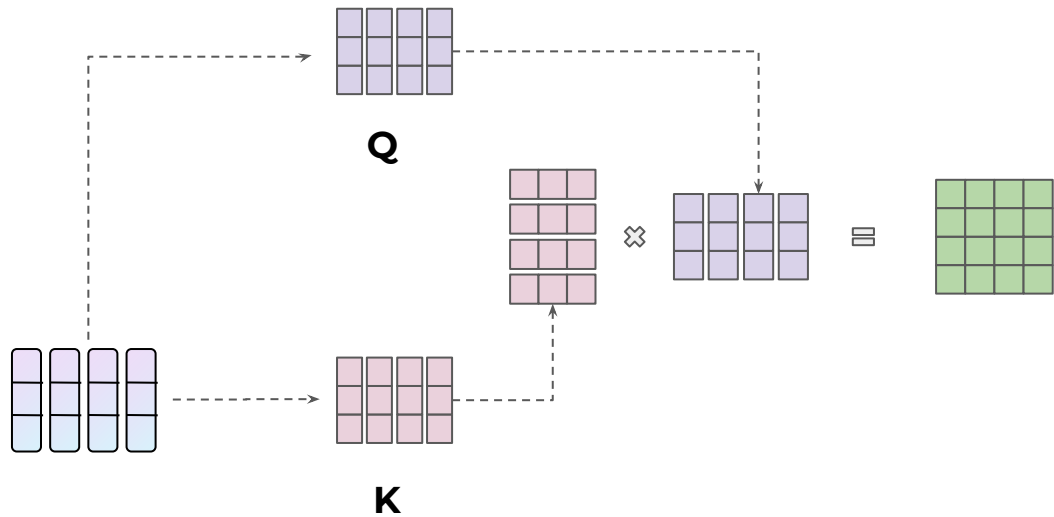


key

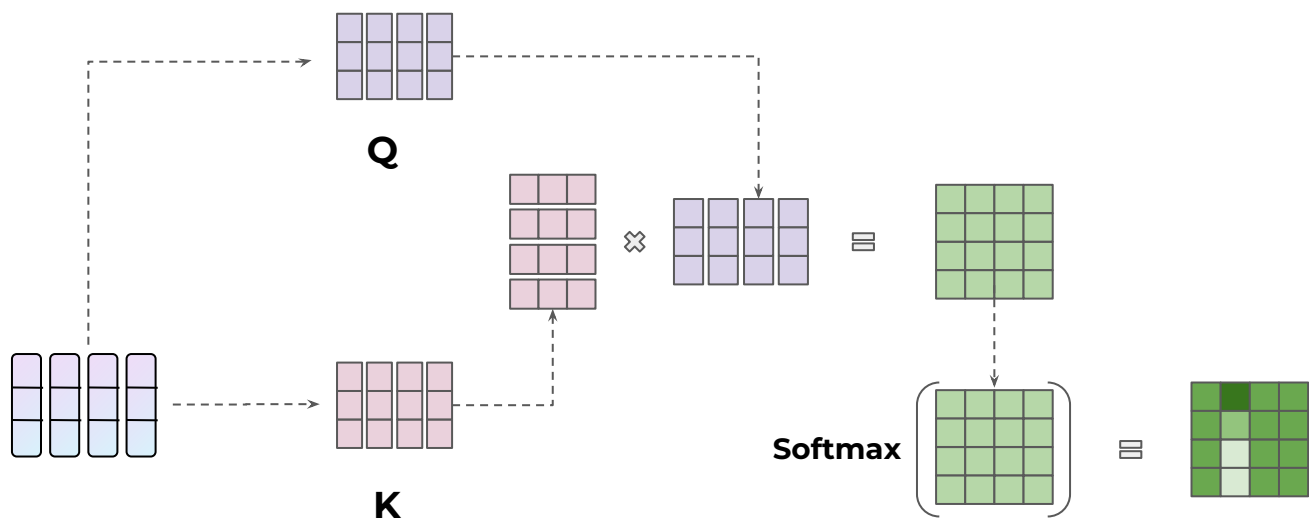




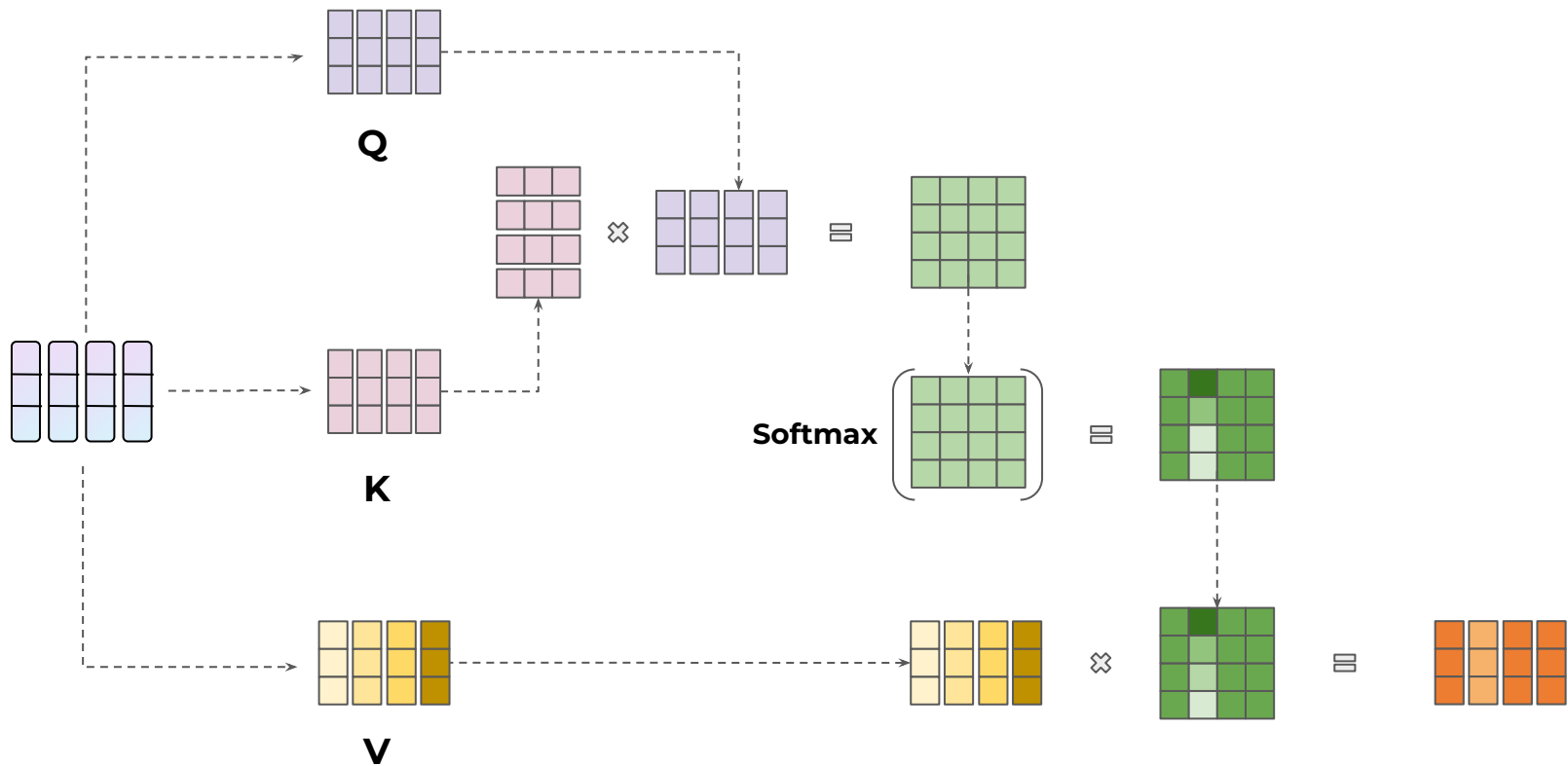
Self-Attention



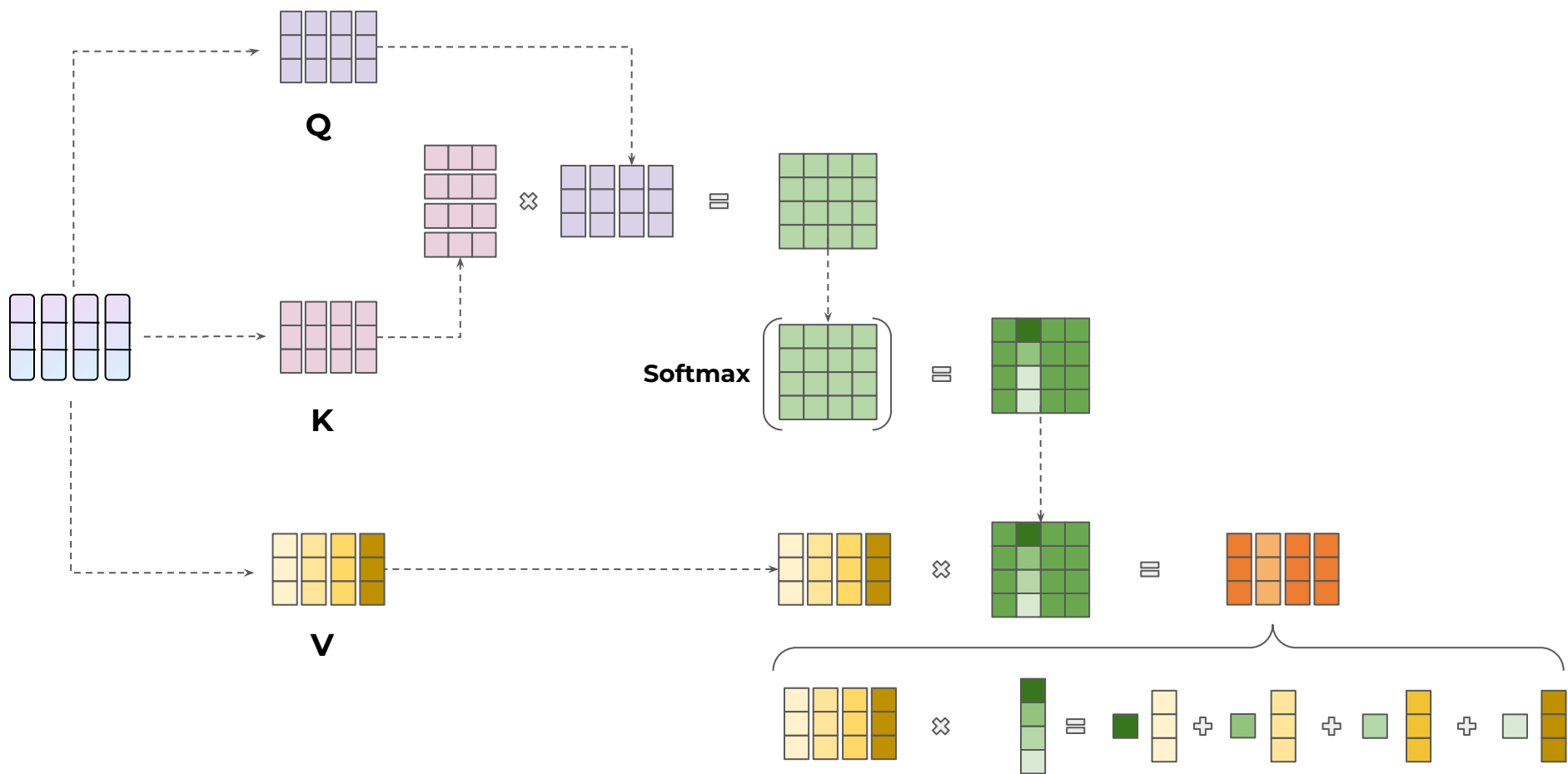
Self-Attention



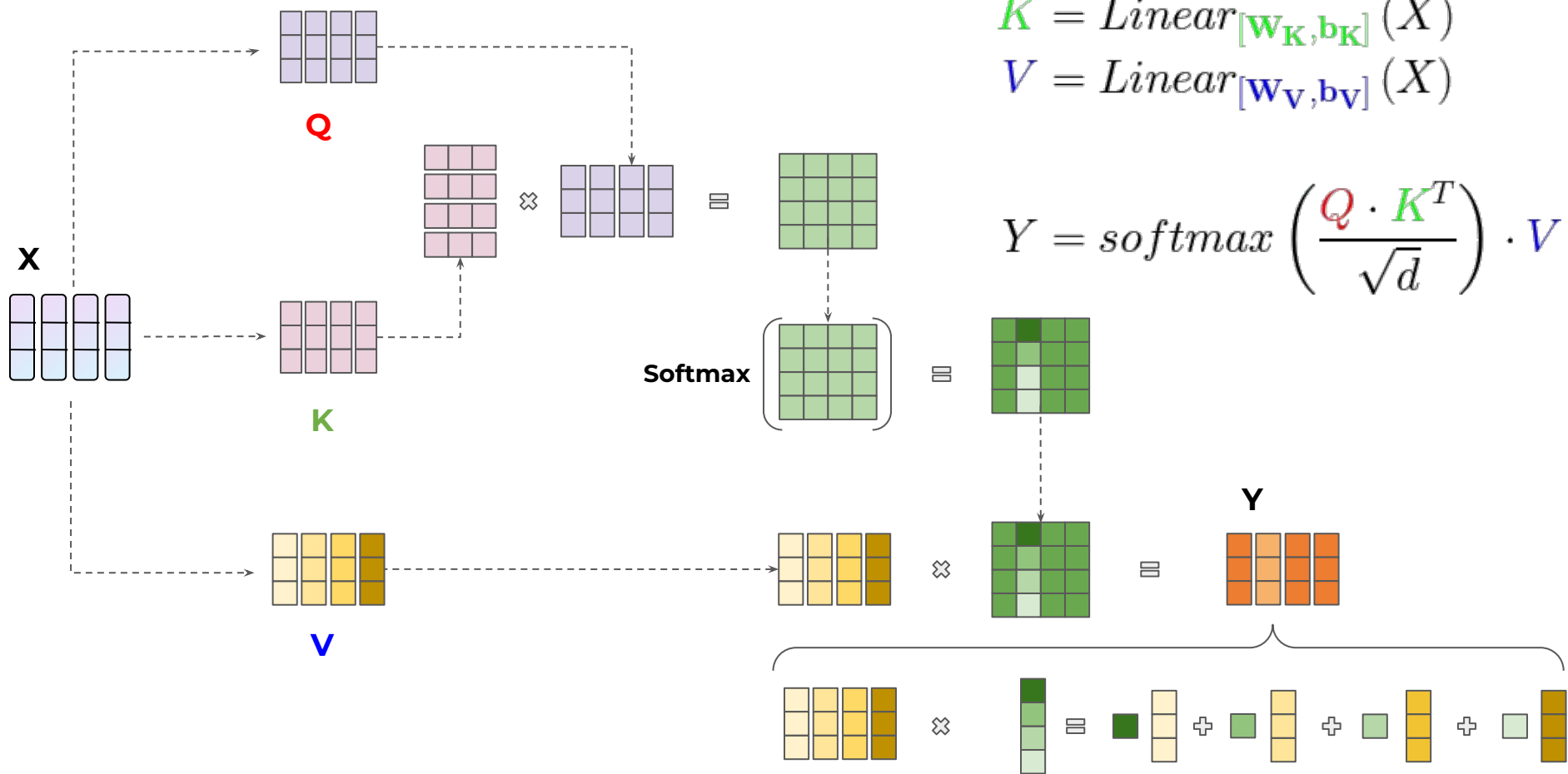
Self-Attention

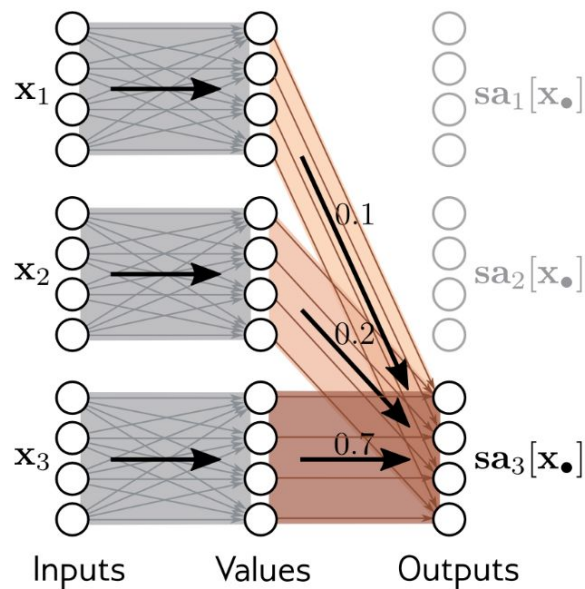
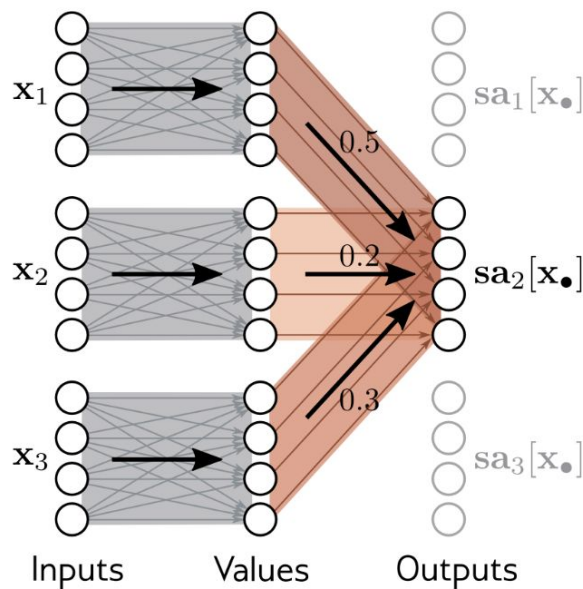
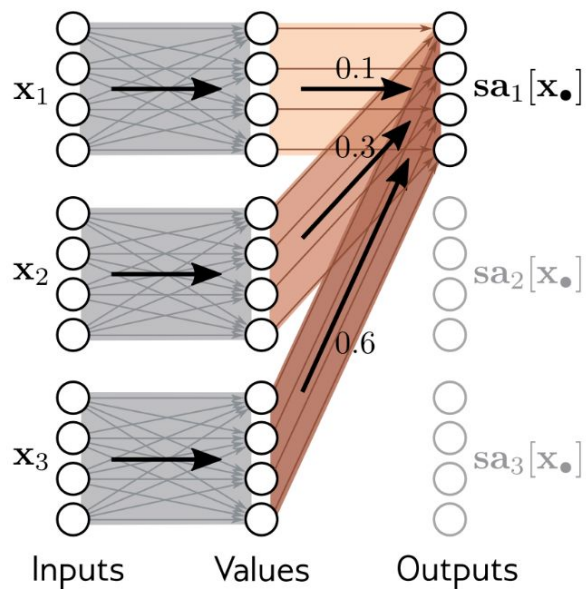


Self-Attention



Self-Attention

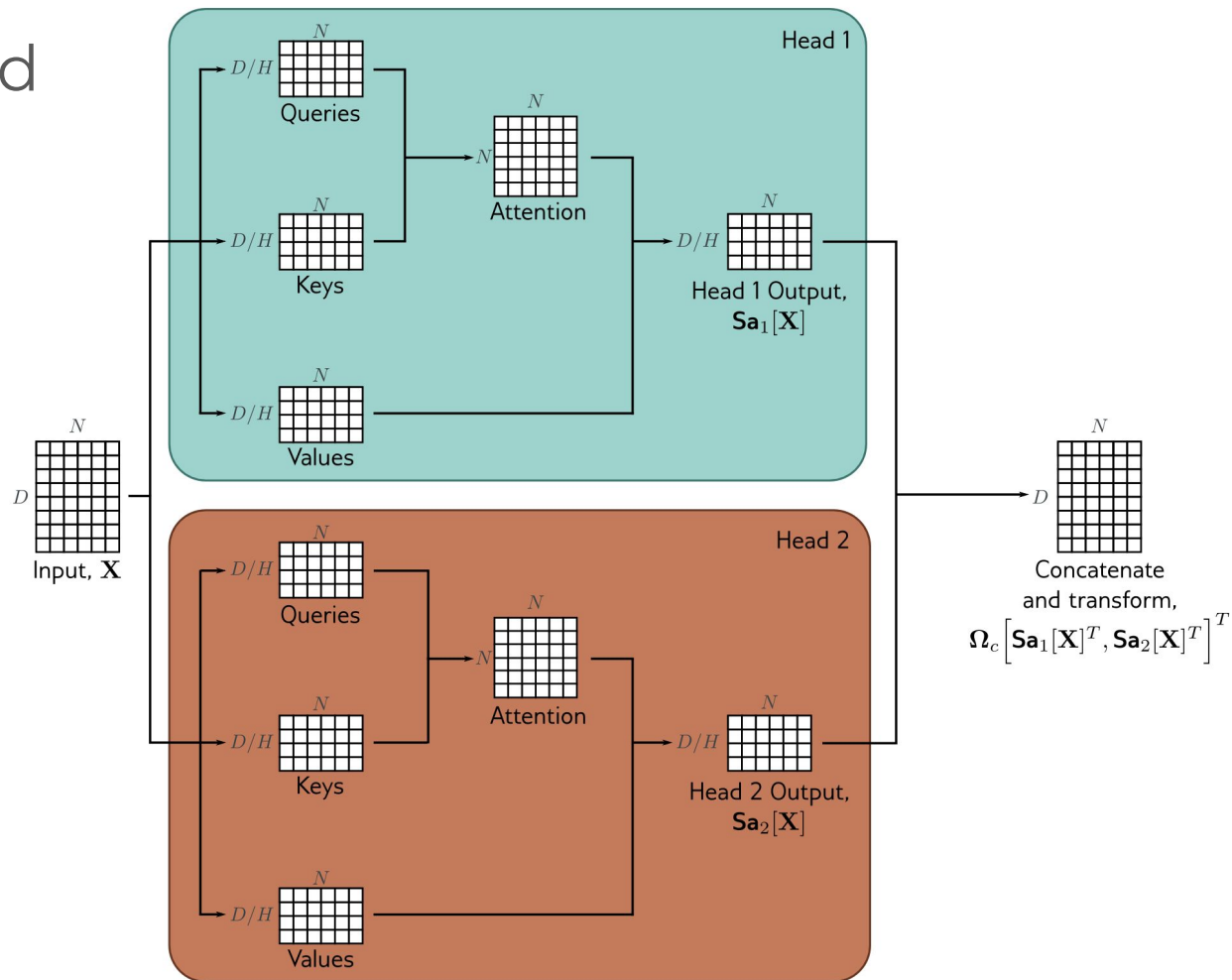




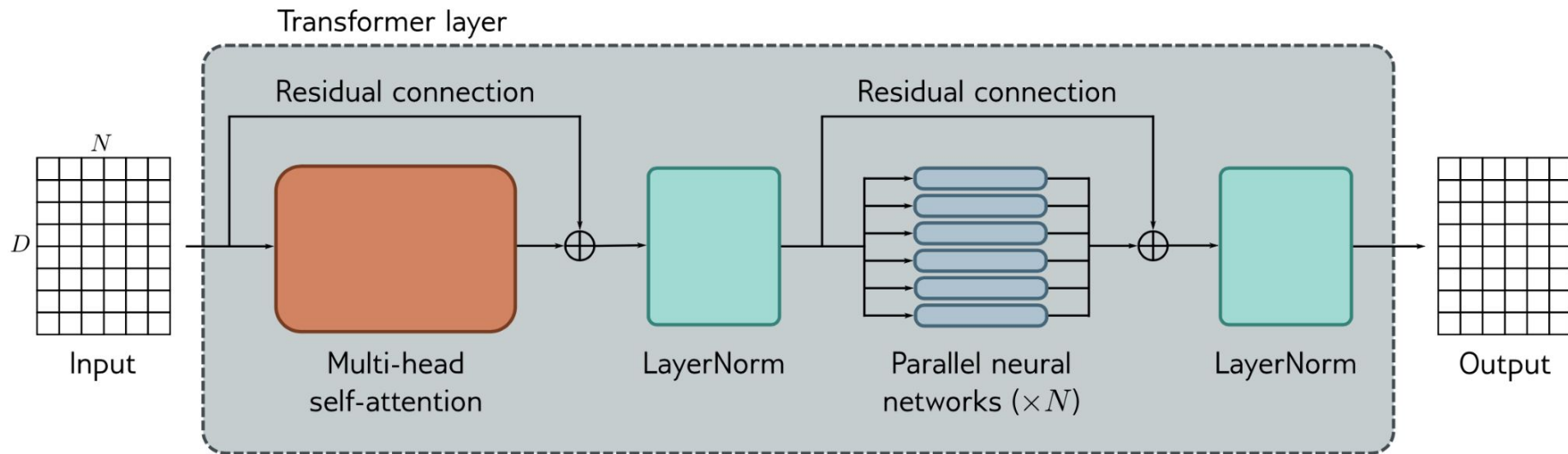
* image from [Understanding Deep Learning](#), book by Simon J.D. Prince, 2023

Transformer layer

Multihead

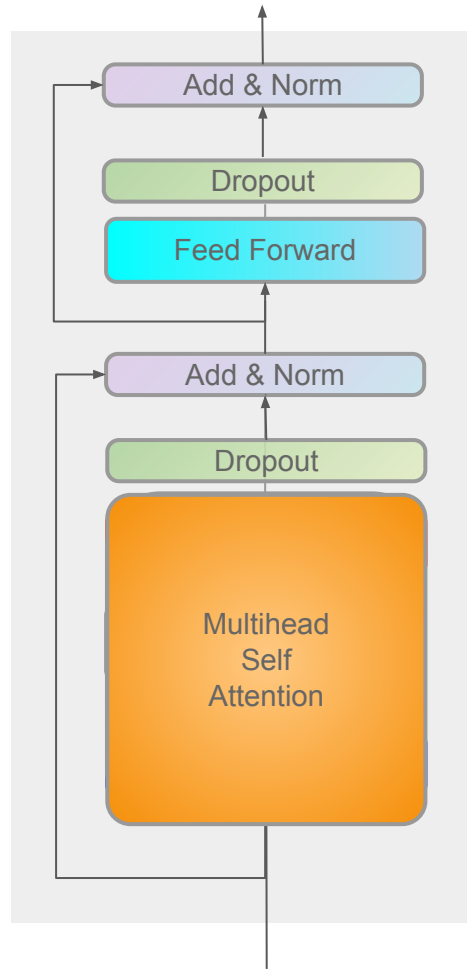


Transformer layer

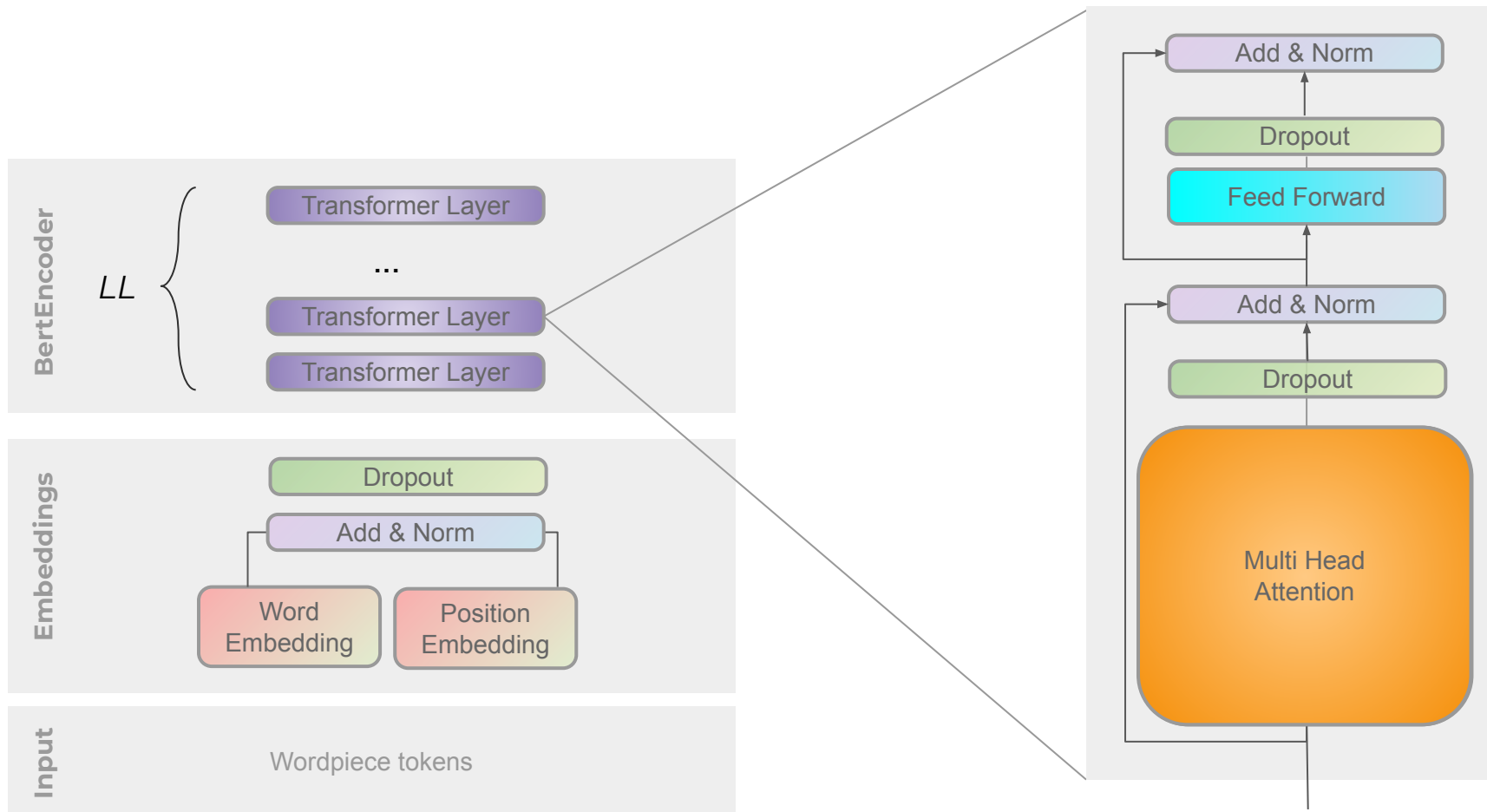


* image from [Understanding Deep Learning](#), book by Simon J.D. Prince, 2023

Transformer layer



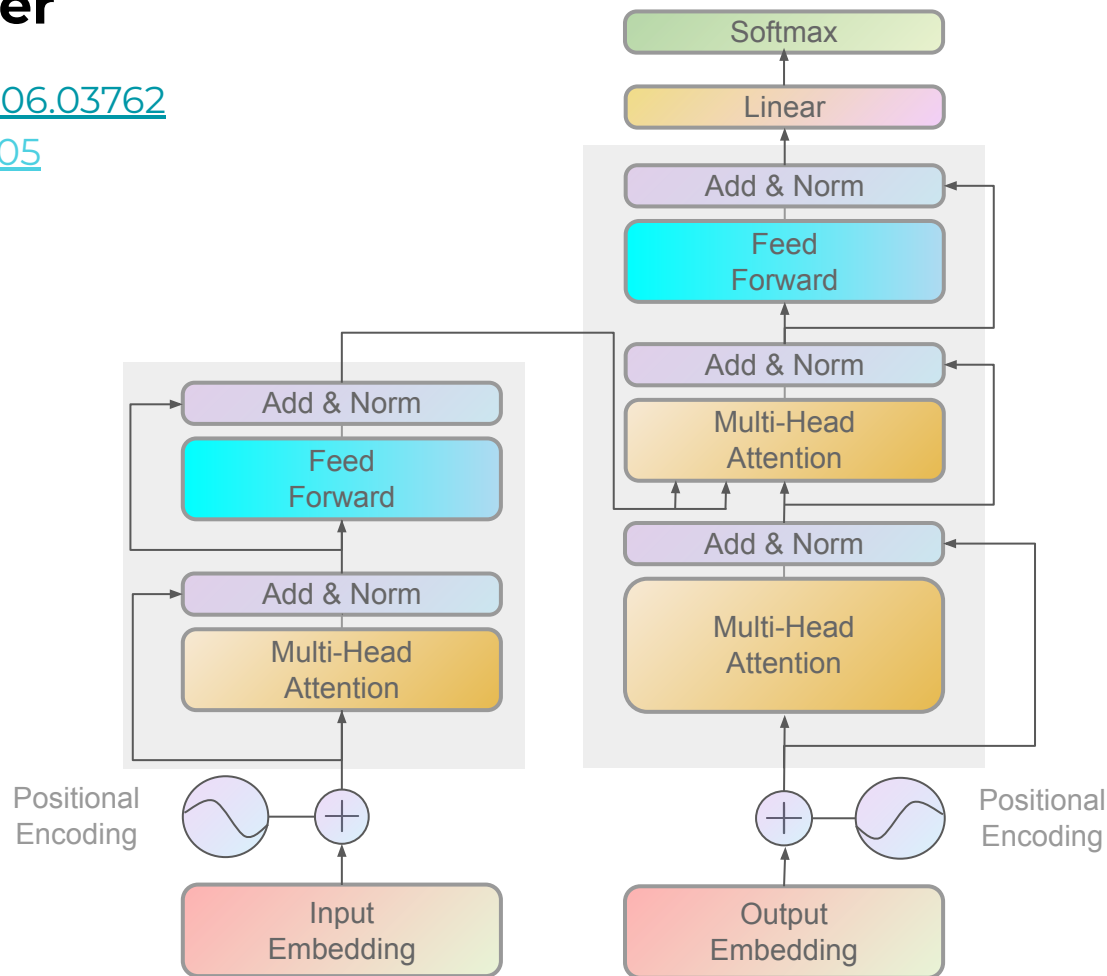
Transformer [Encoder]



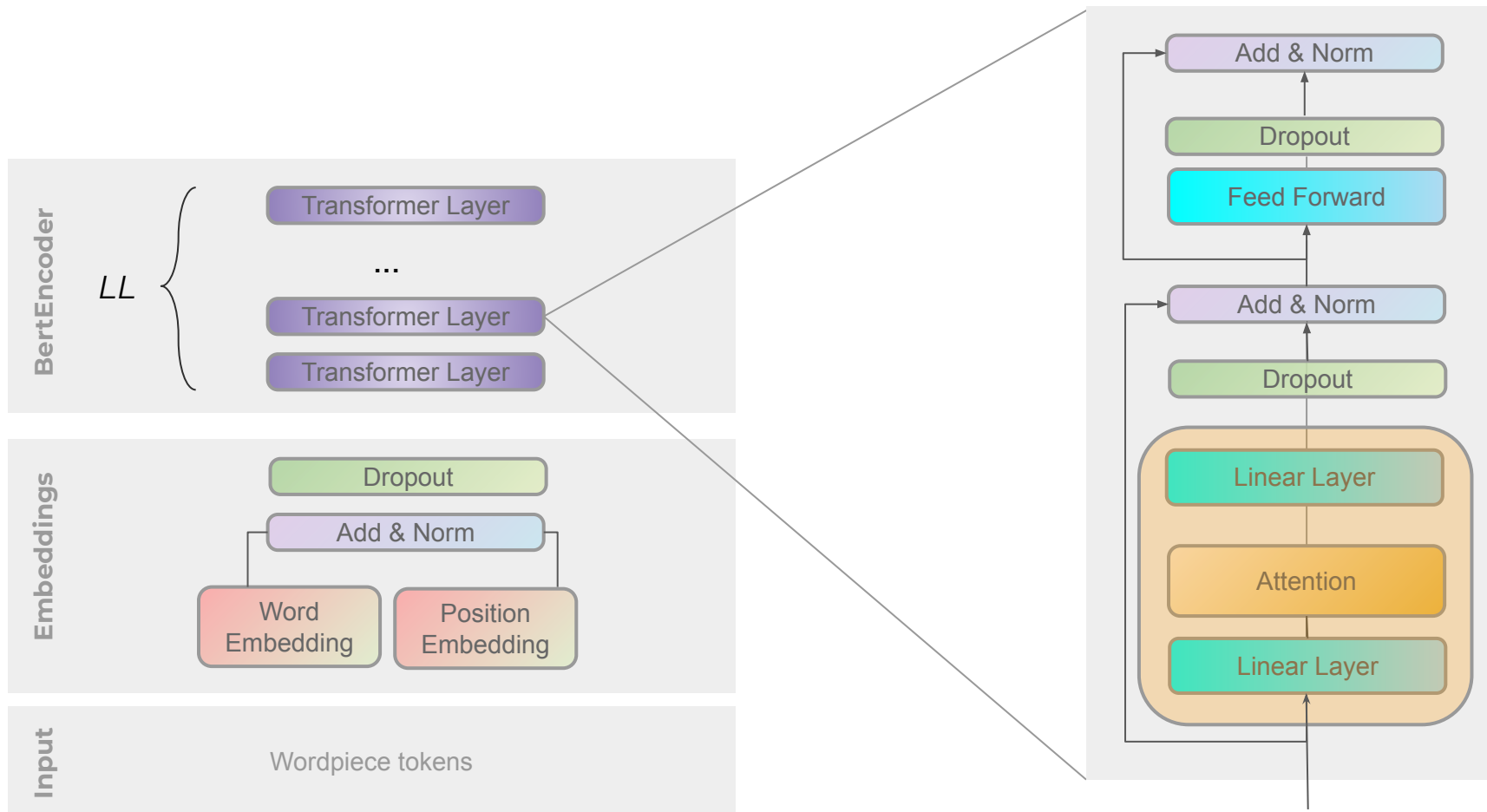
Transformer

Transformer: [1706.03762](#)

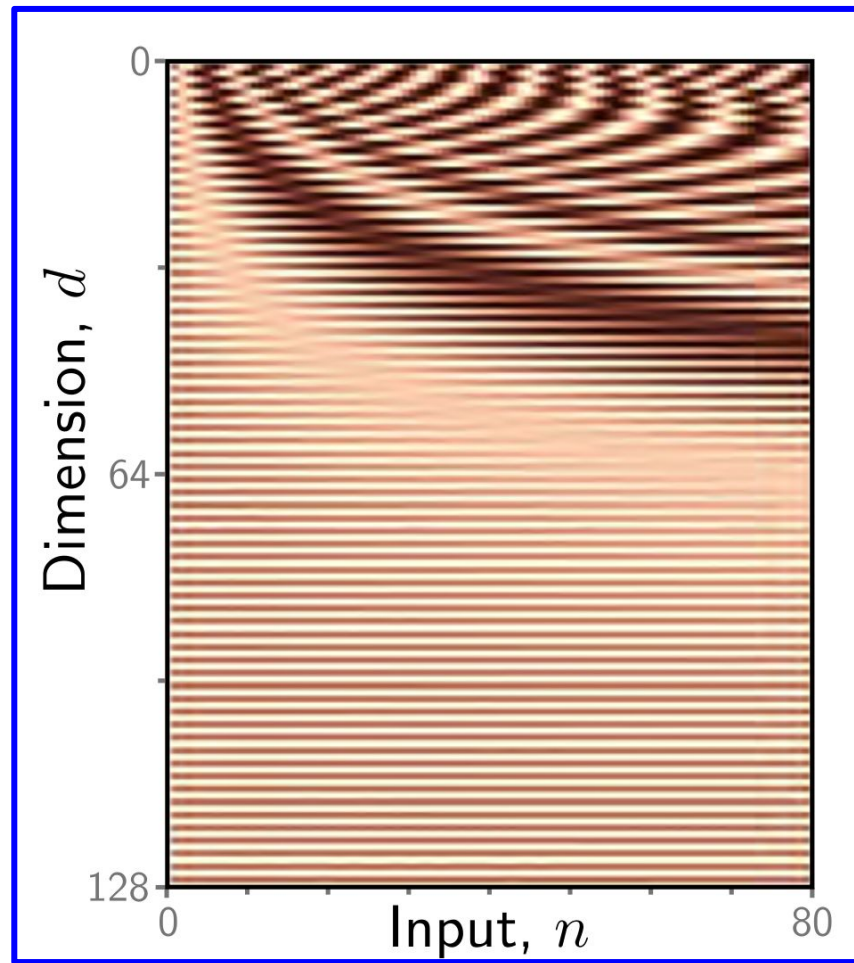
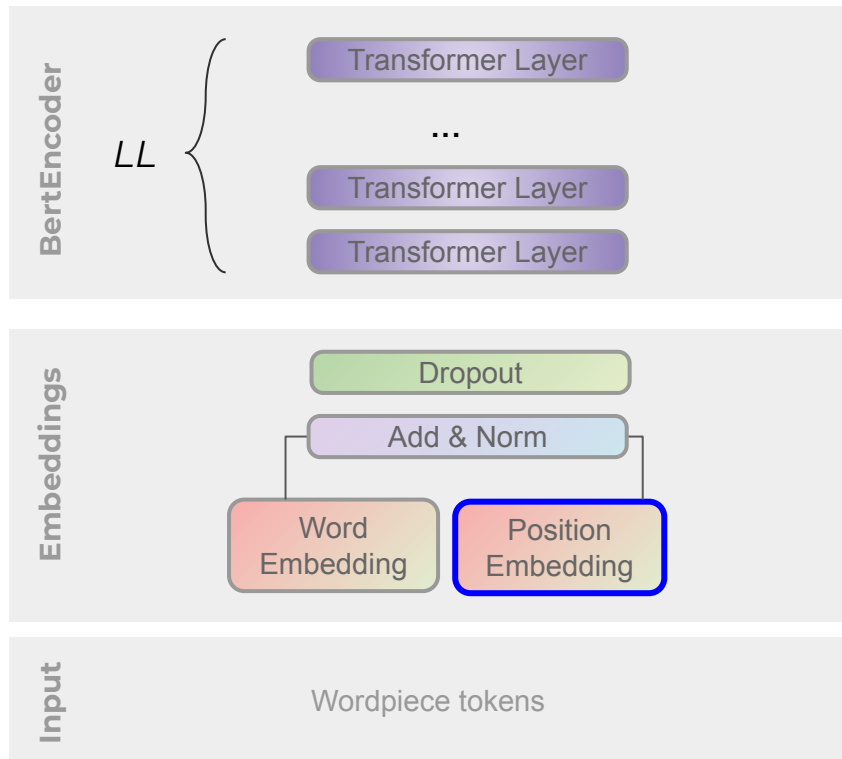
BERT: [1810.04805](#)



Transformer [Encoder]



Transformer [Encoder]



Memory & Compute

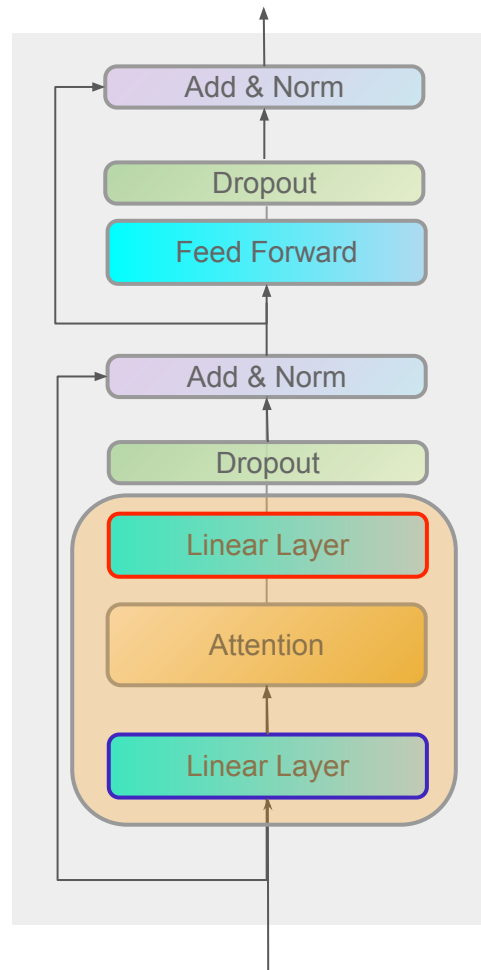
Transformer [Multi head attention]

$$FC(X) = Linear_{[W_{out}, b_{out}]}(\sigma(Linear_{[W_{in}, b_{in}]}(X)))$$

$$MultiHead(\mathbf{X}, \mathbf{X}) = Linear_{[W^o, b^o]} \left(concat_{i \in [A]} [\mathbf{H}^i] \right)$$

$$Attention(Q, K, V) = softmax \left(\frac{Q \cdot K^T}{\sqrt{d}} \right) \cdot V$$

$$\begin{aligned} \mathbf{H}^i = & Attention(Linear_{[W_Q^{(i)}, b_Q^{(i)}]}(X), \\ & Linear_{[W_K^{(i)}, b_K^{(i)}]}(X), \\ & Linear_{[W_V^{(i)}, b_V^{(i)}]}(X)) \end{aligned}$$



Transformer [Embeddings]

Memory

Computations

$$2 * [H]$$

$$B * L * H$$

$$[V \times H]$$

-

$$[M \times H]$$

-

$$[T \times H]$$

-

B - batch_size

H - hidden

V - vocabulary size

M - max token length

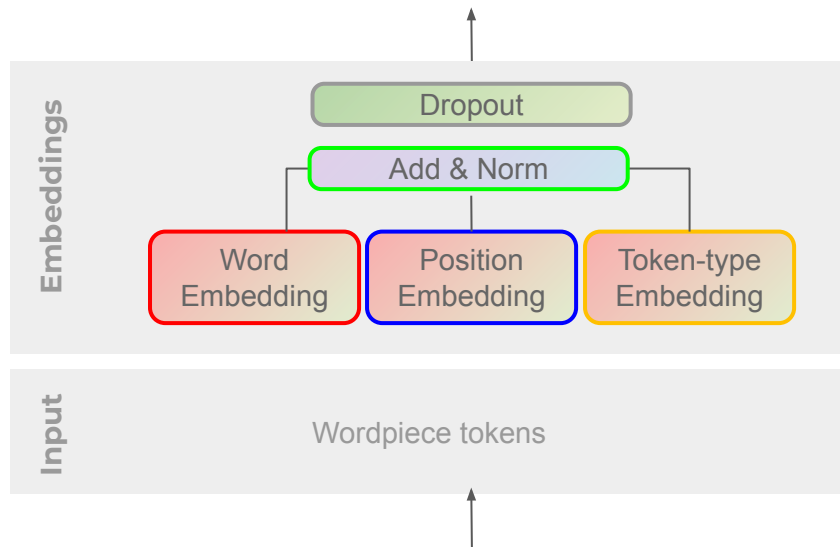
T - number of types

[bert-base: 768]

[bert-base: 50265]

[bert-base: 514]

[bert-base: 1]



Transformer [MHA]

Memory

$$2 * [H]$$

$$[A * S + 1 \times H]$$

-

$$3 * [H + 1 \times A * S]$$

Computations

$$B * L * H$$

$$B * L * (A * S) * H$$

$$B * A * L * L * S * 2$$

$$B * L * H * (A * S) * 3$$

B - batch_size

L - batch_length

H - hidden

A - num_heads

S - head_size (= H/A)

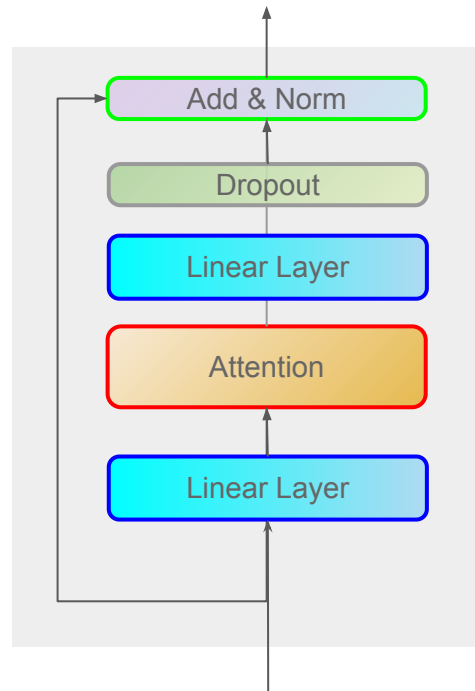
I - intermediate FC size

[bert-base: 768]

[bert-base: 12]

[bert-base: 64]

[bert-base: 3072]



Transformer [FC]

Memory

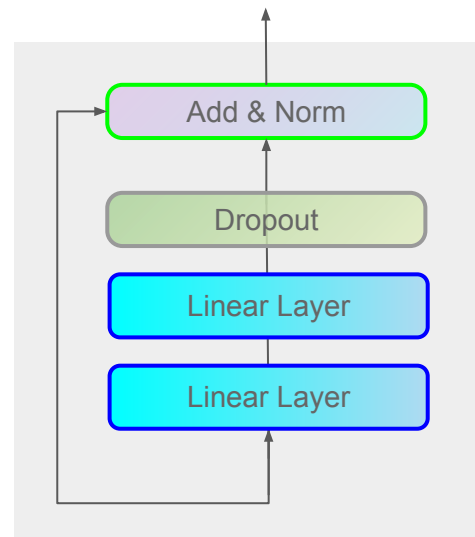
$$2 * [H]$$

$$[H + 1 \times I] + [I + 1 \times H]$$

Computations

$$B * L * H$$

$$2 * B * L * H * I$$



B - batch_size

L - batch_length

H - hidden

A - num_heads

S - head_size (= H/A)

I - intermediate FC size

[bert-base: 768]

[bert-base: 12]

[bert-base: 64]

[bert-base: 3072]

Transformer [summary]

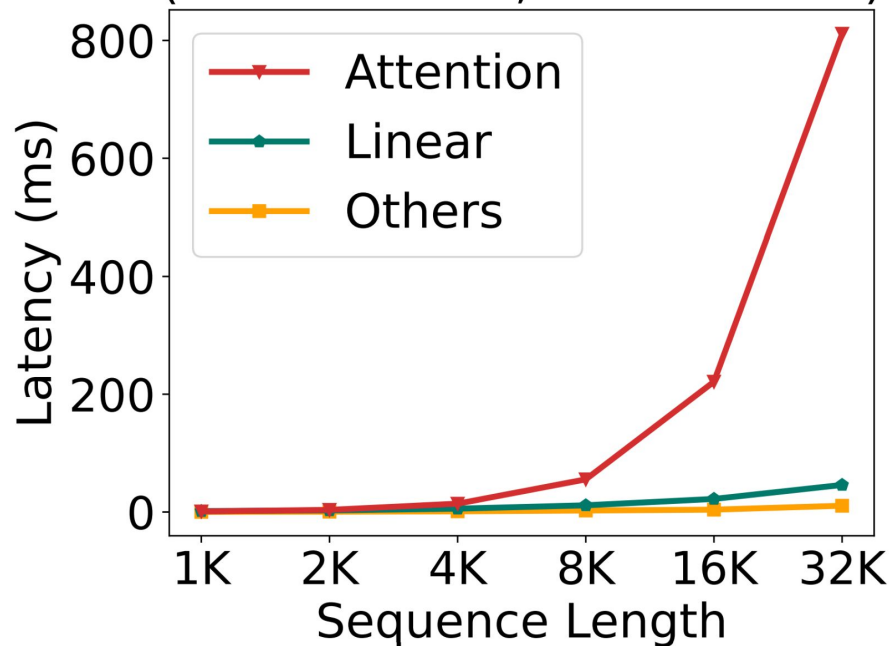
		<u>#params [M]</u>	<u>mocs [B]</u>
LL	FC-LayerNorm	$2 * H$	$B * L * H$
	FC-Out	$(1 + H) * I$	$B * L * H * I$
	FC-In	$(1 + I) * H$	$B * L * H * I$
	Atten-LayerNorm	$2 * H$	$B * L * H$
	Atten-LinearOutput	$(1 + H) * H$	$B * L * H * H$
	Atten-AttenScoreValue	0	$B * A * L * L * S$
	Atten-AttenScore	0	$B * A * L * L * S$
	Atten-LinearInput	$(1 + H) * H * 3$	$B * L * H * H * 3$
	Embedding-LayerNorm	$2 * H$	$B * L * H$
	Embedding-Word	$V * H$	0
	Embedding-Position	$M * H$	0
	Embedding-Token-Type	$T * H$	0

Transformer [summary]

		<u>#params [M]</u>	<u>mocs [B]</u>
LL	FC-LayerNorm	$2 * H$	$B * L * H$
	FC-Out	$(1 + H) * I$	$B * L * H * I$
	FC-In	$(1 + I) * H$	$B * L * H * I$
	Atten-LayerNorm	$2 * H$	$B * L * H$
	Atten-LinearOutput	$(1 + H) * H$	$B * L * H * H$
	Atten-AttenScoreValue	0	$B * A * \mathbf{L} * \mathbf{L} * S$
	Atten-AttenScore	0	$B * A * \mathbf{L} * \mathbf{L} * S$
	Atten-LinearInput	$(1 + H) * H * 3$	$B * L * H * H * 3$
	Embedding-LayerNorm	$2 * H$	$B * L * H$
	Embedding-Word	$V * H$	0
	Embedding-Position	$M * H$	0
	Embedding-Token-Type	$T * H$	0

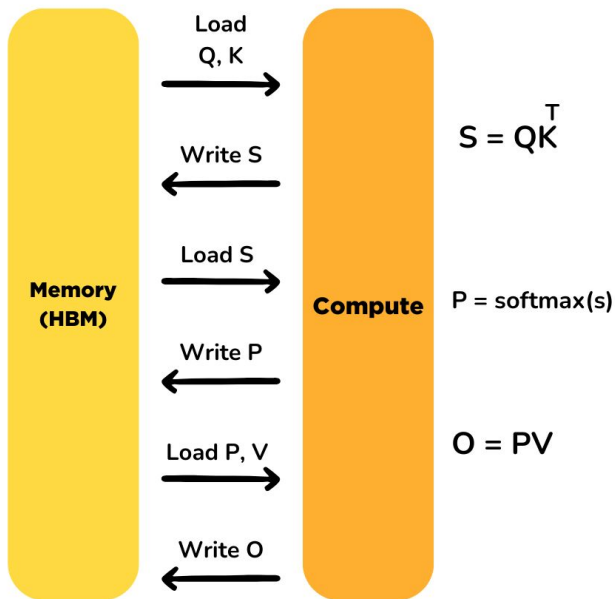
Transformer [Memory & Runtime]

Latency of operators in a transformer model
(NumHead=32, HeadDim=64)

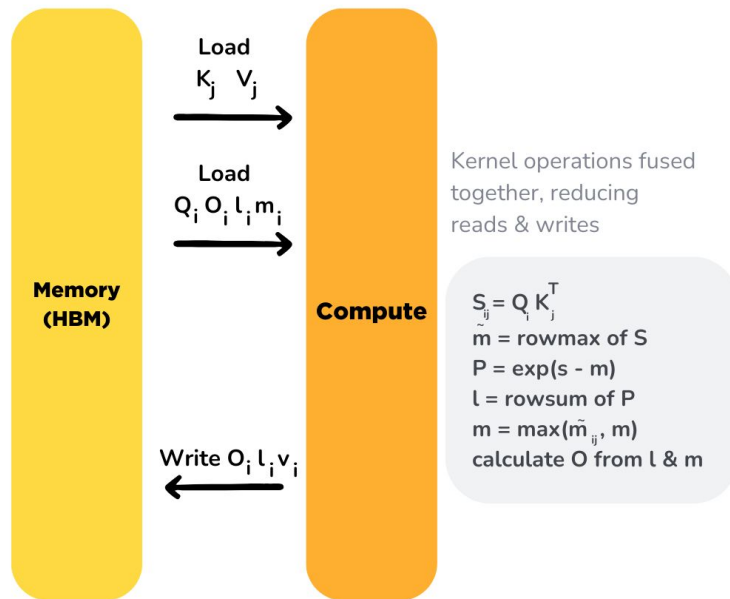


Transformer [Flash Attention]

Standard Attention Implementation

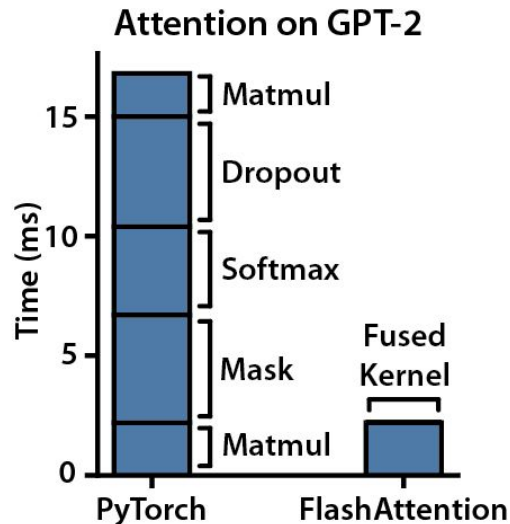
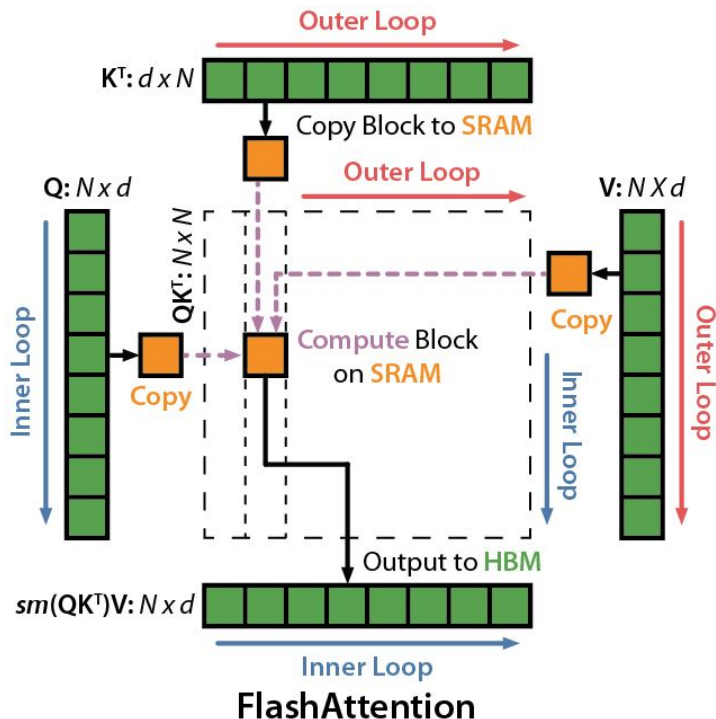
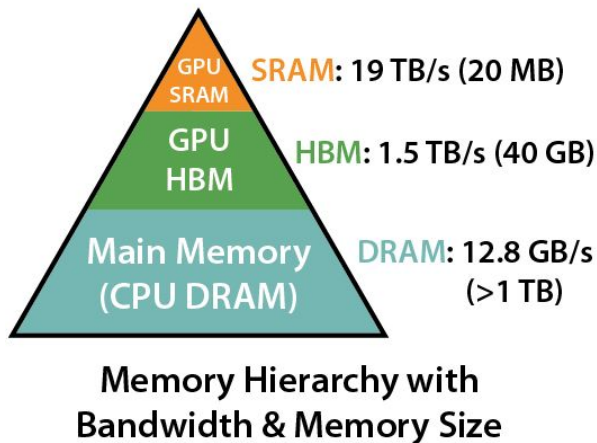


Flash Attention

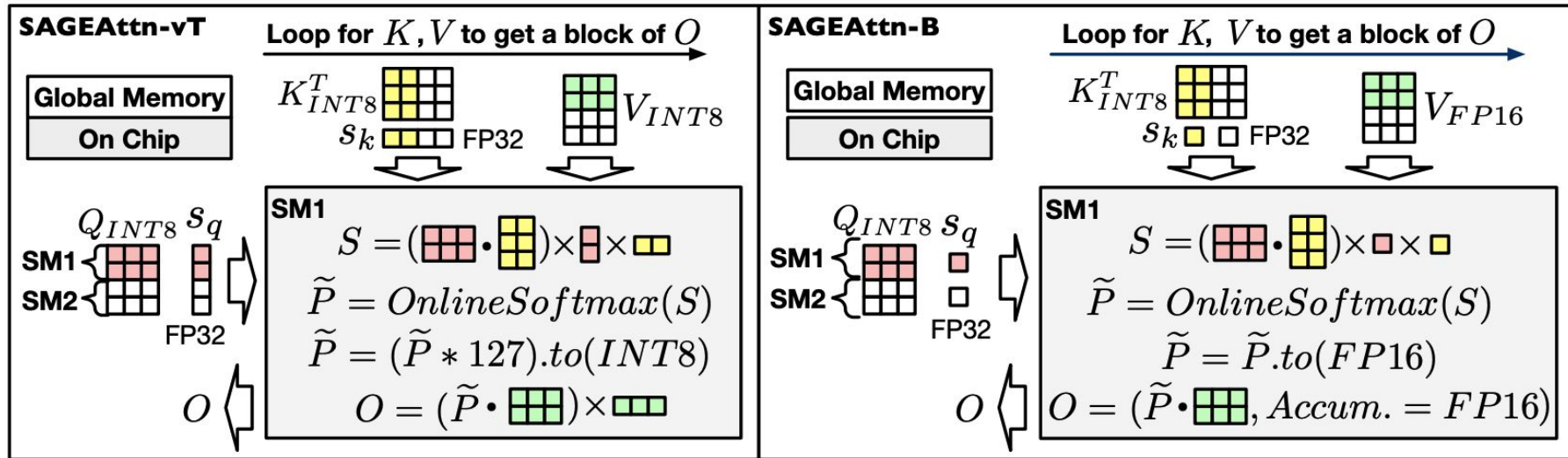


Initialize O, l and m matrices with zeroes. m and l are used to calculate cumulative softmax. Divide Q, K, V into blocks (due to SRAM's memory limits) and iterate over them, for i is row & j is column.

Transformer [Flash Attention]



Transformer [Sage Attention]



(a) SageAttention (per-token quantize Q,K; INT8 V)

(b) SageAttention (per-block quantize Q,K; FP16 V)

Transformer [another view]

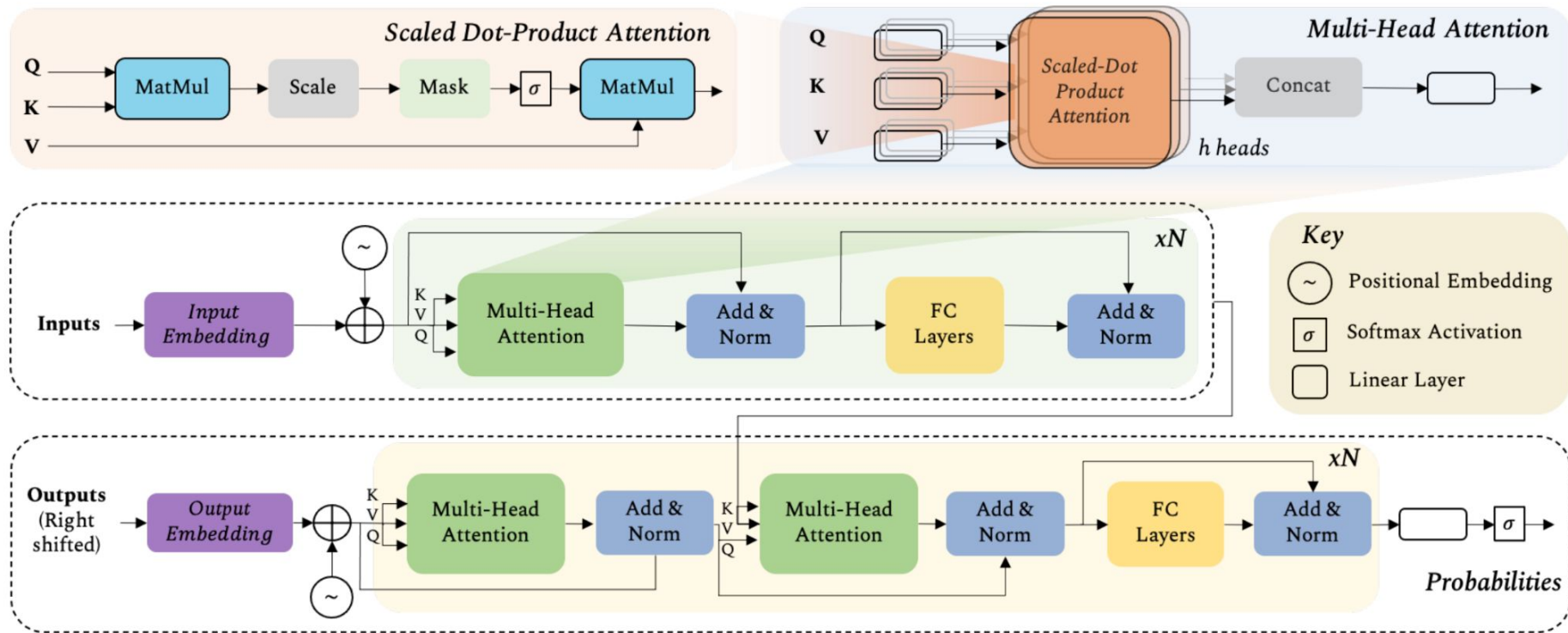
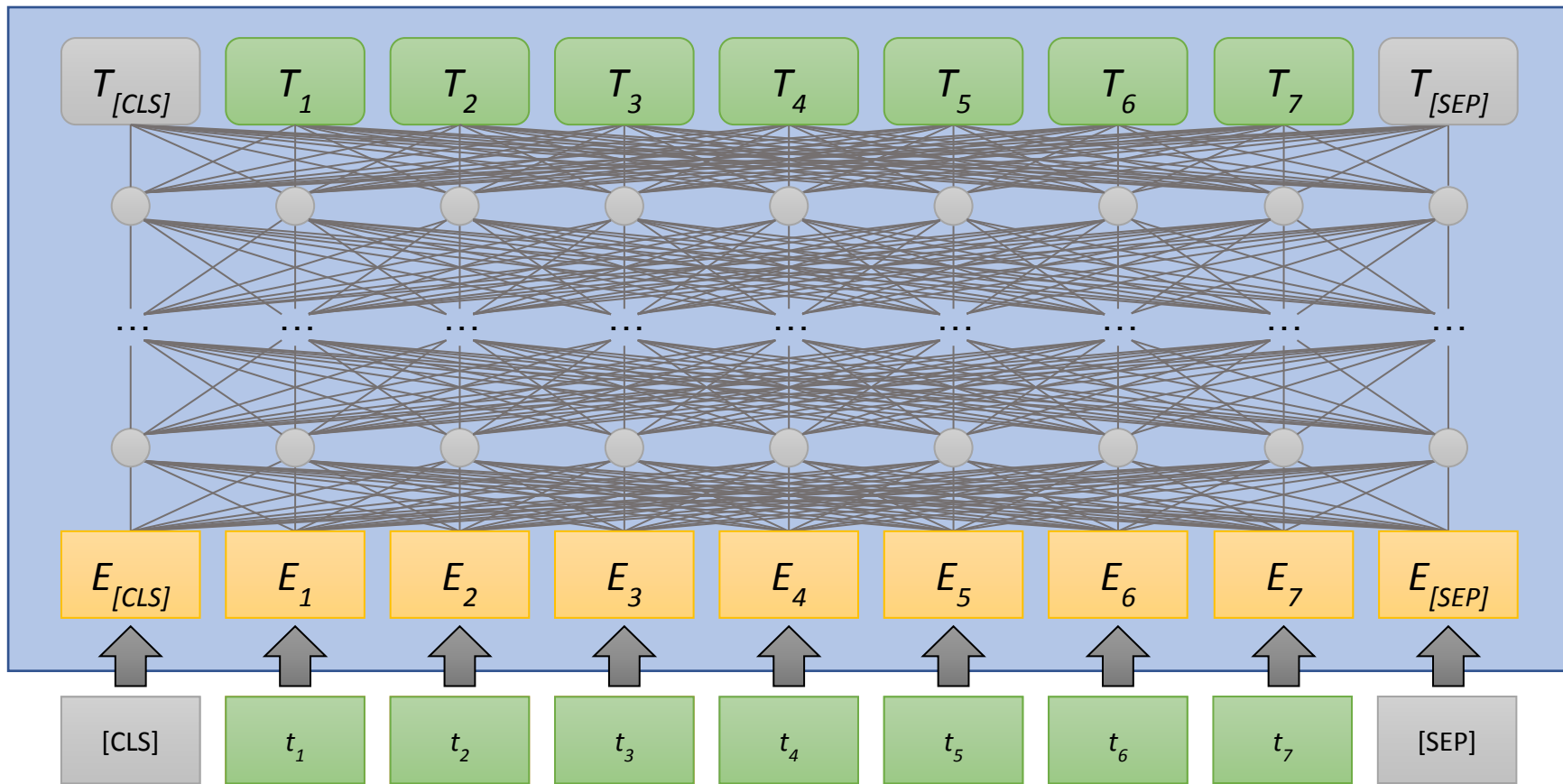
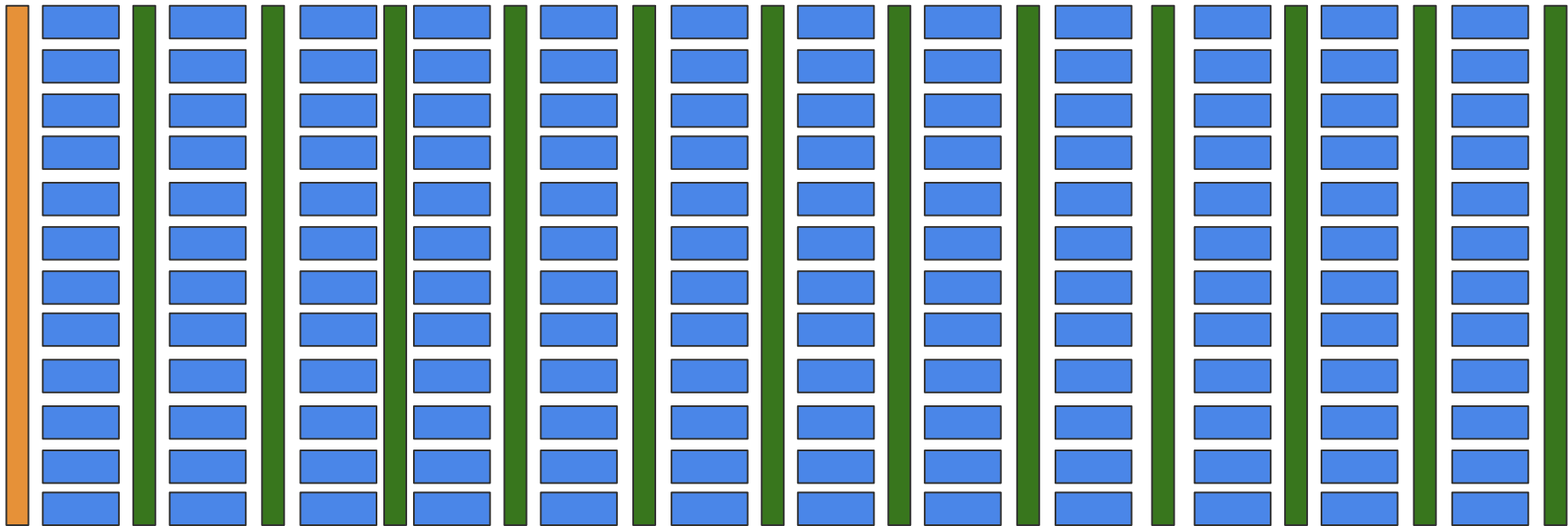


image source: 2101.01169

Transformer [another view]



Transformer [another view]



$L = 12$ (bert-base), $12 \times 12 = 144$ heads