

Fast Attention with less memory

Advanced NLP: Summer 2023

Anoop Sarkar

REFORMER: THE EFFICIENT TRANSFORMER

Nikita Kitaev*

U.C. Berkeley & Google Research

kitaev@cs.berkeley.edu

Łukasz Kaiser*

Google Research

{lukaszkaiser,levskaya}@google.com

Anselm Levskaya

Google Research

<https://openreview.net/forum?id=rkgNKkHtvB>

https://iclr.cc/virtual_2020/poster_rkgNKkHtvB.html

Transformer

From the BERT documentation:

Using the default training scripts (`run_classifier.py` and `run_squad.py`), we benchmarked the maximum batch size on single Titan X GPU (12GB RAM) with TensorFlow 1.11.0:

System	Seq Length	Max Batch Size
BERT-Large	64	12
...	128	6
...	256	2
...	320	1
...	384	0
...	512	0

ZERO!



Outlook

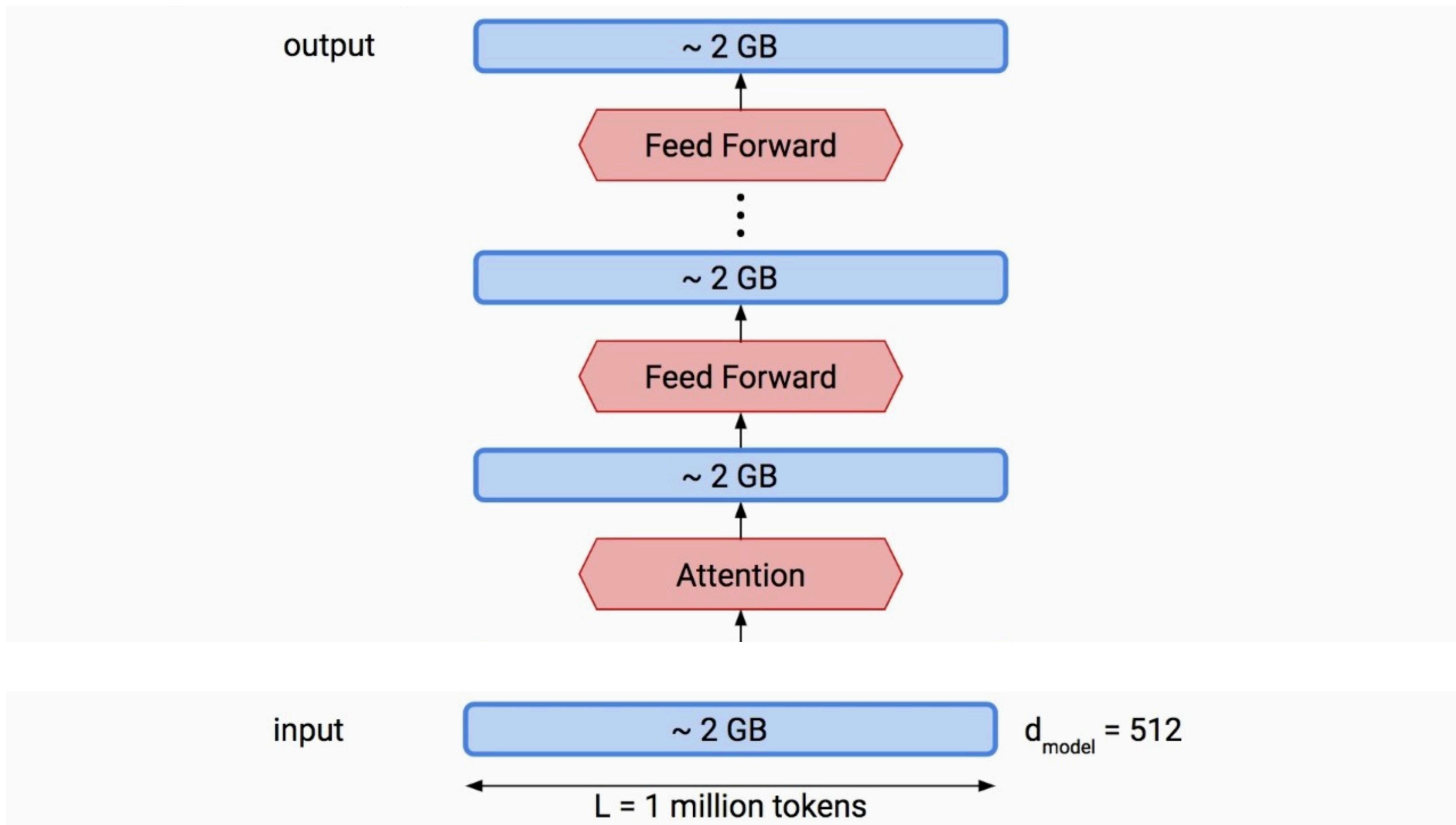
In the near future, it will be impossible to even fine tune state of the art models without datacenter scale hardware resources.

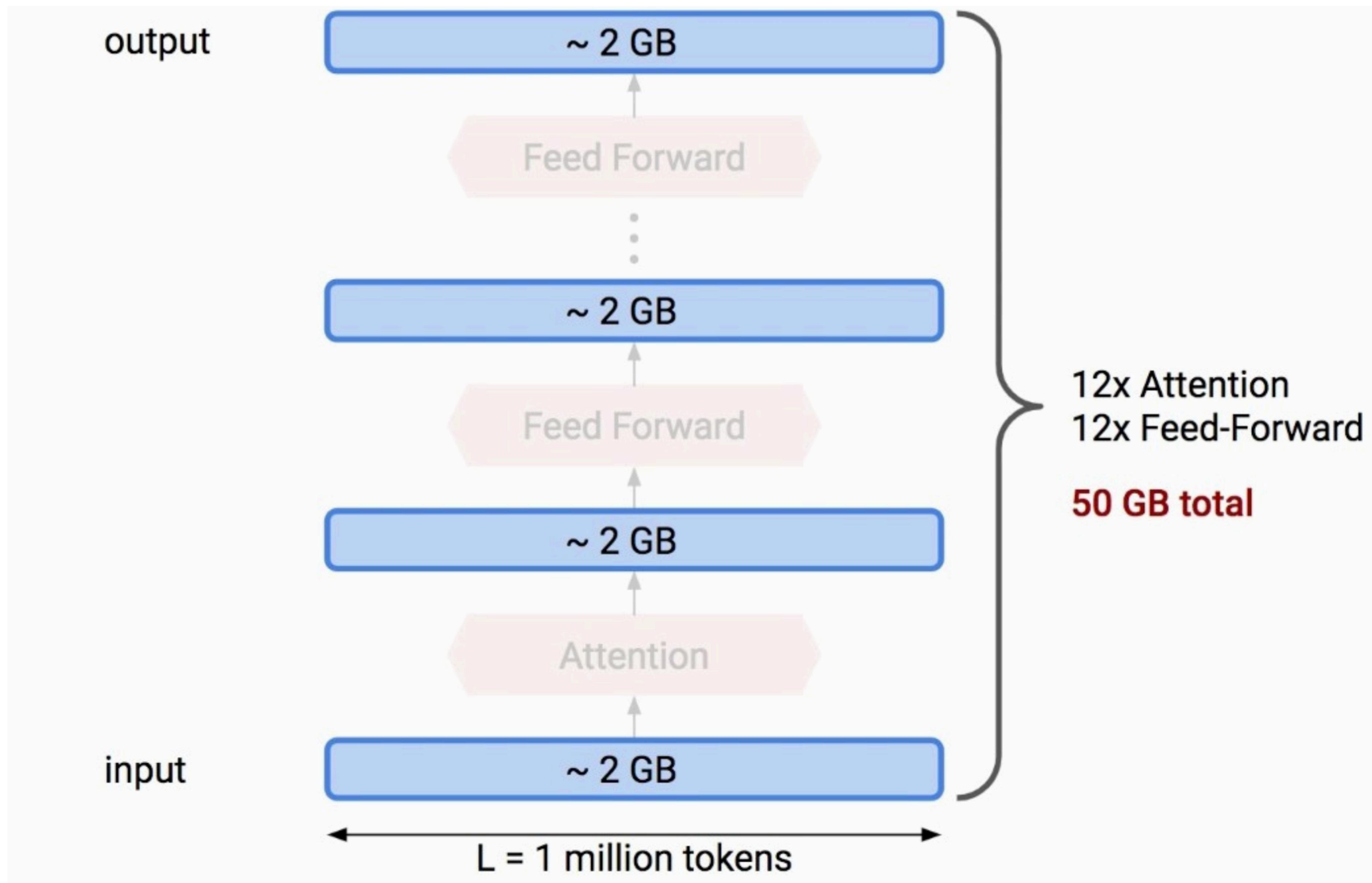
Transformers can be adapted to run on today's hardware over entire chapters or documents of text -- up to 1 million tokens at a time.

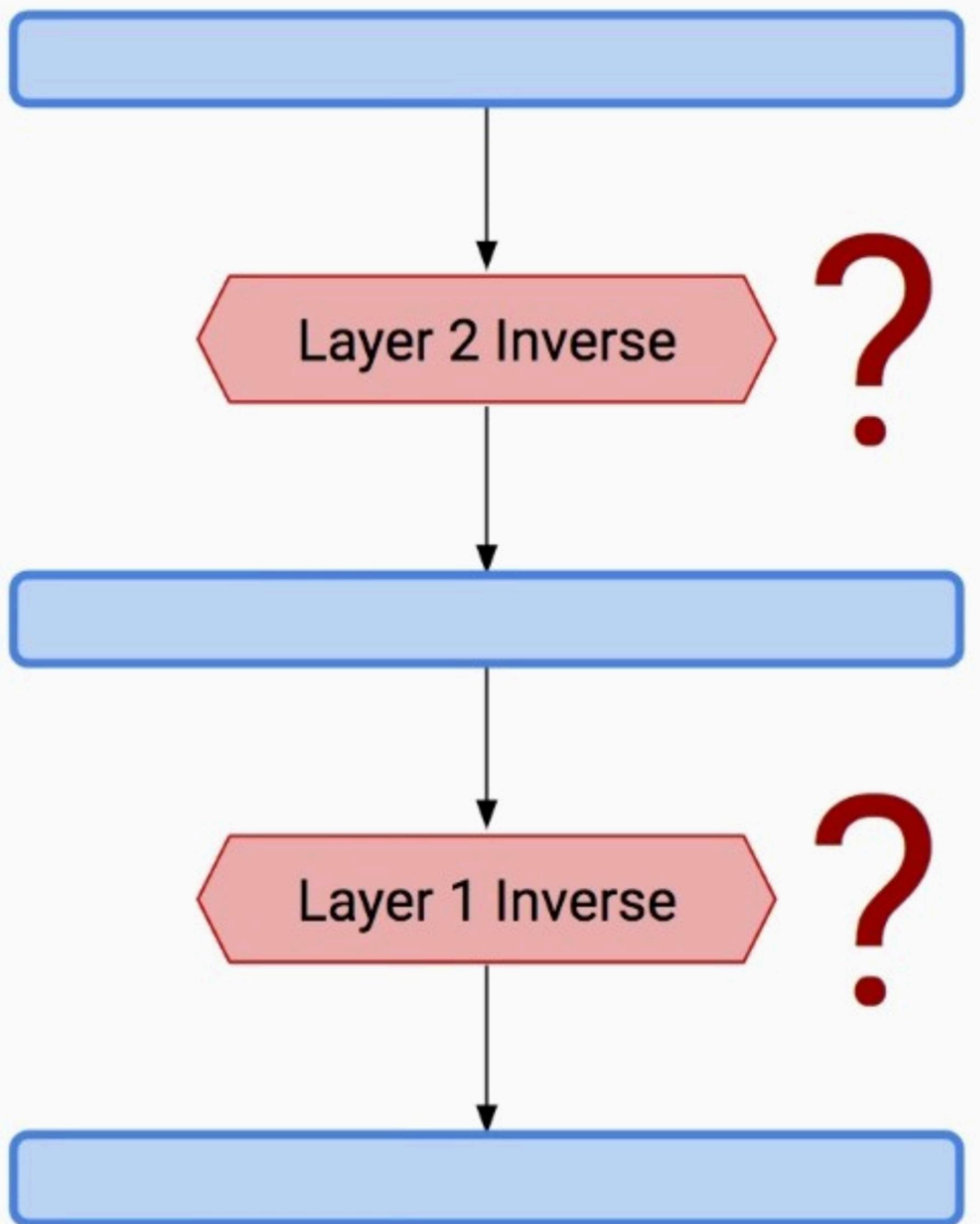
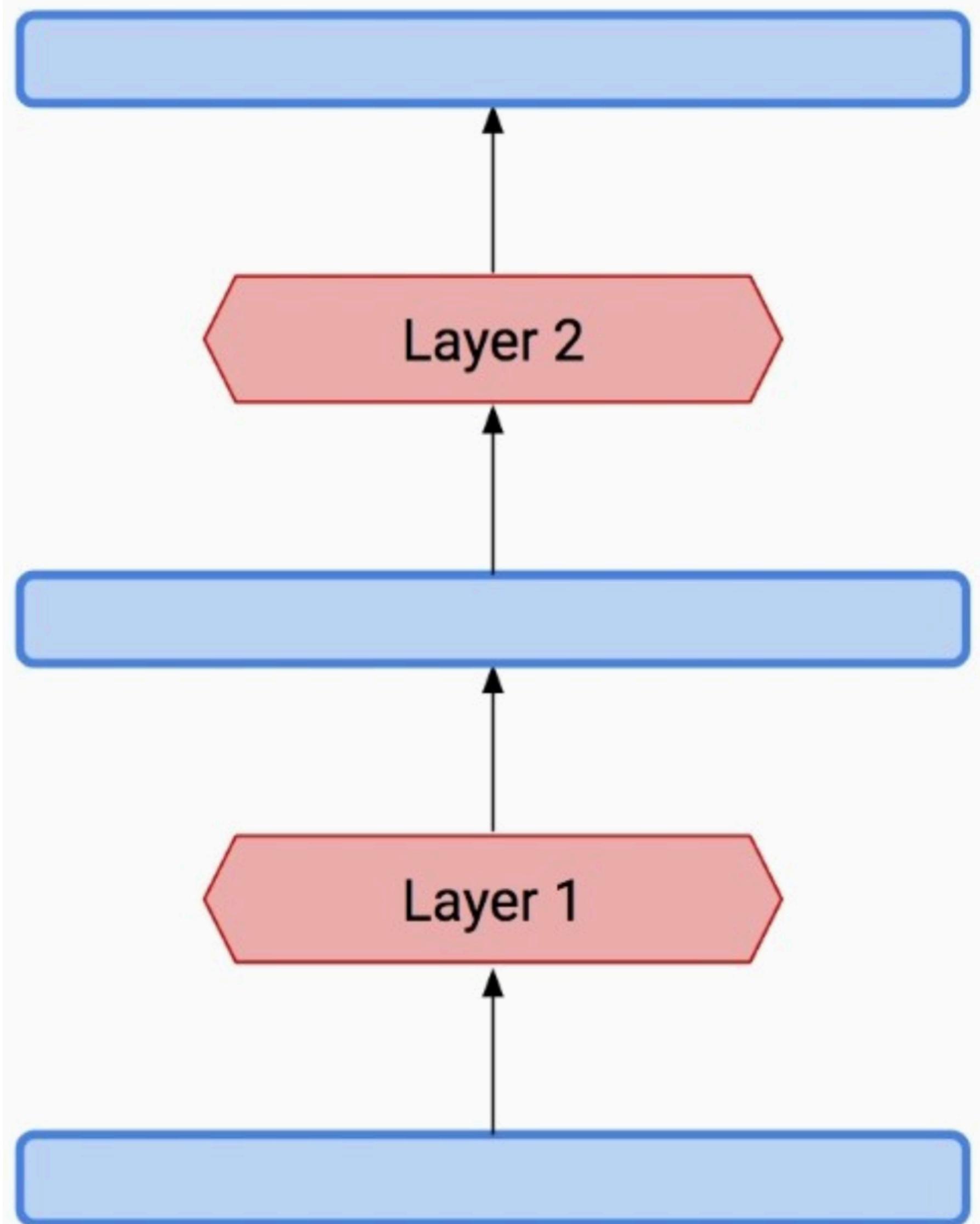
Moreover, the model should run on a single GPU or TPU device.

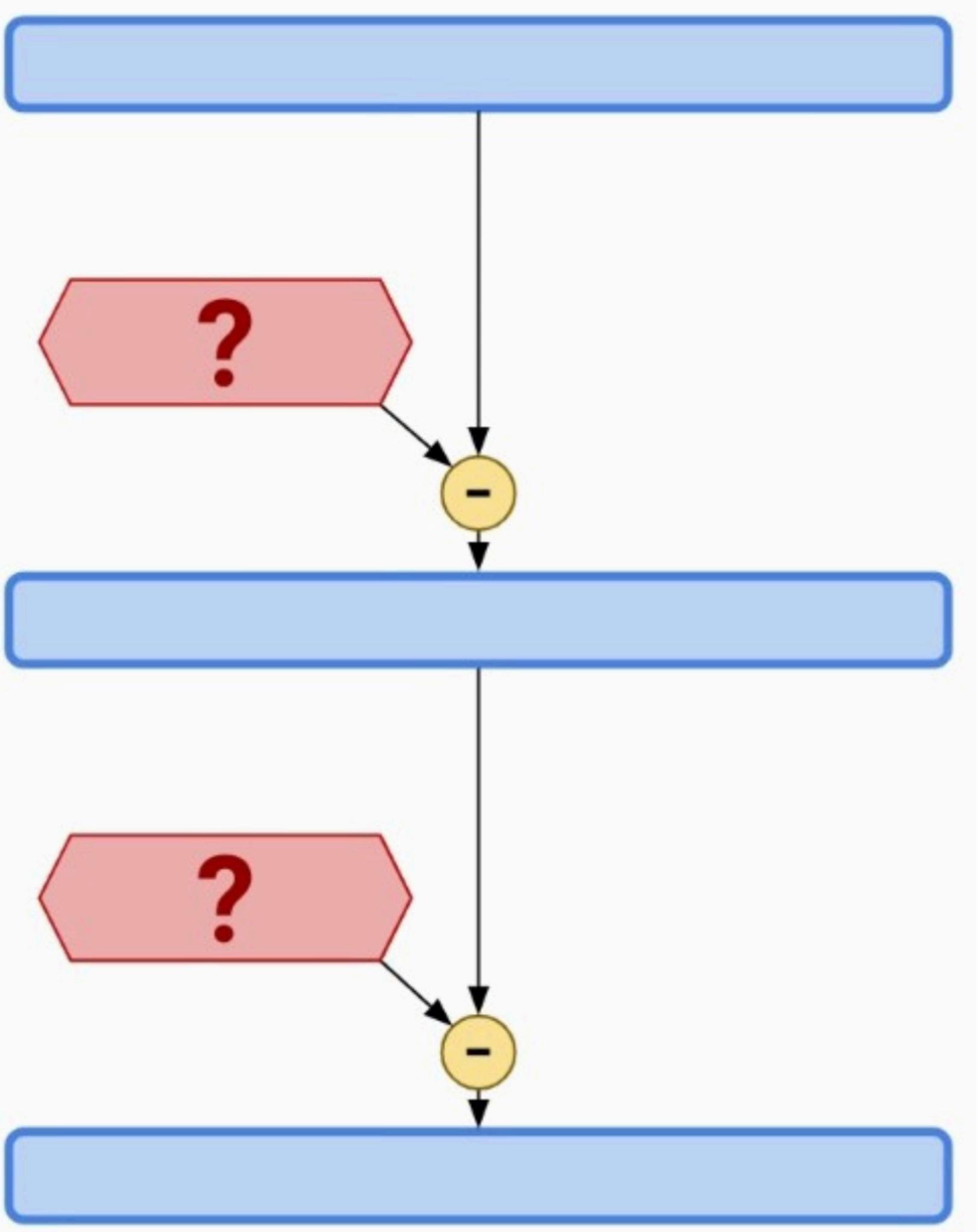
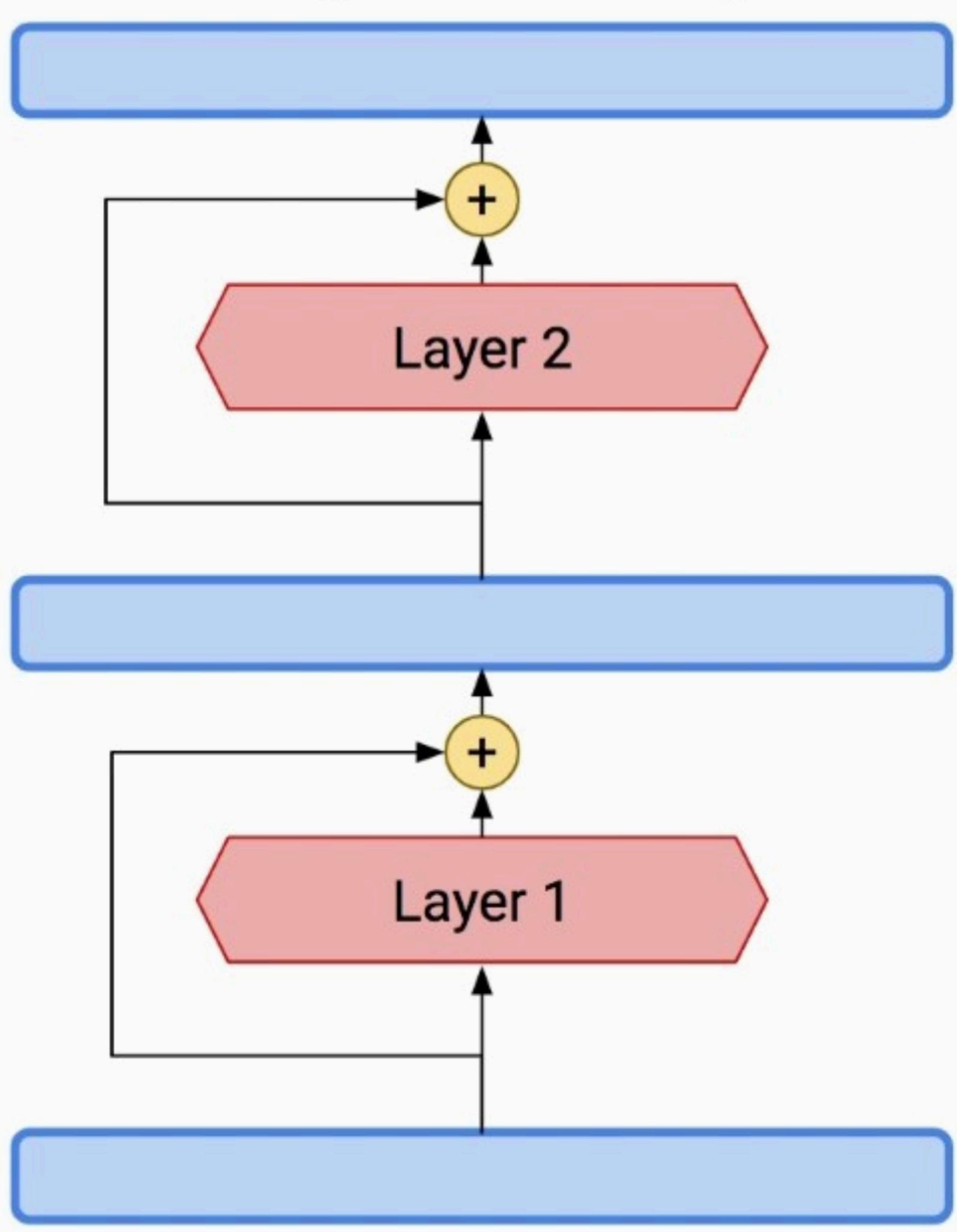
Efficiency Challenges

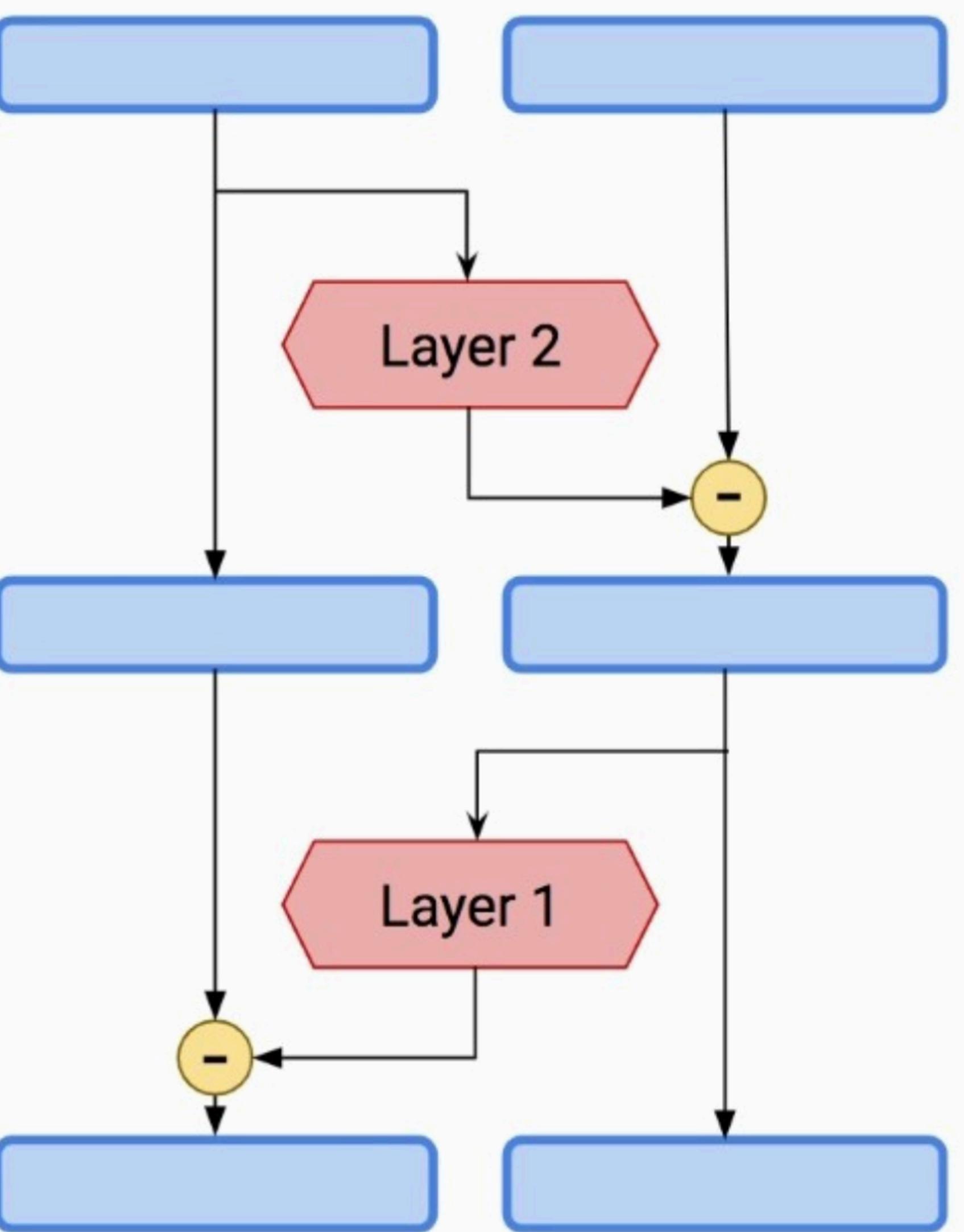
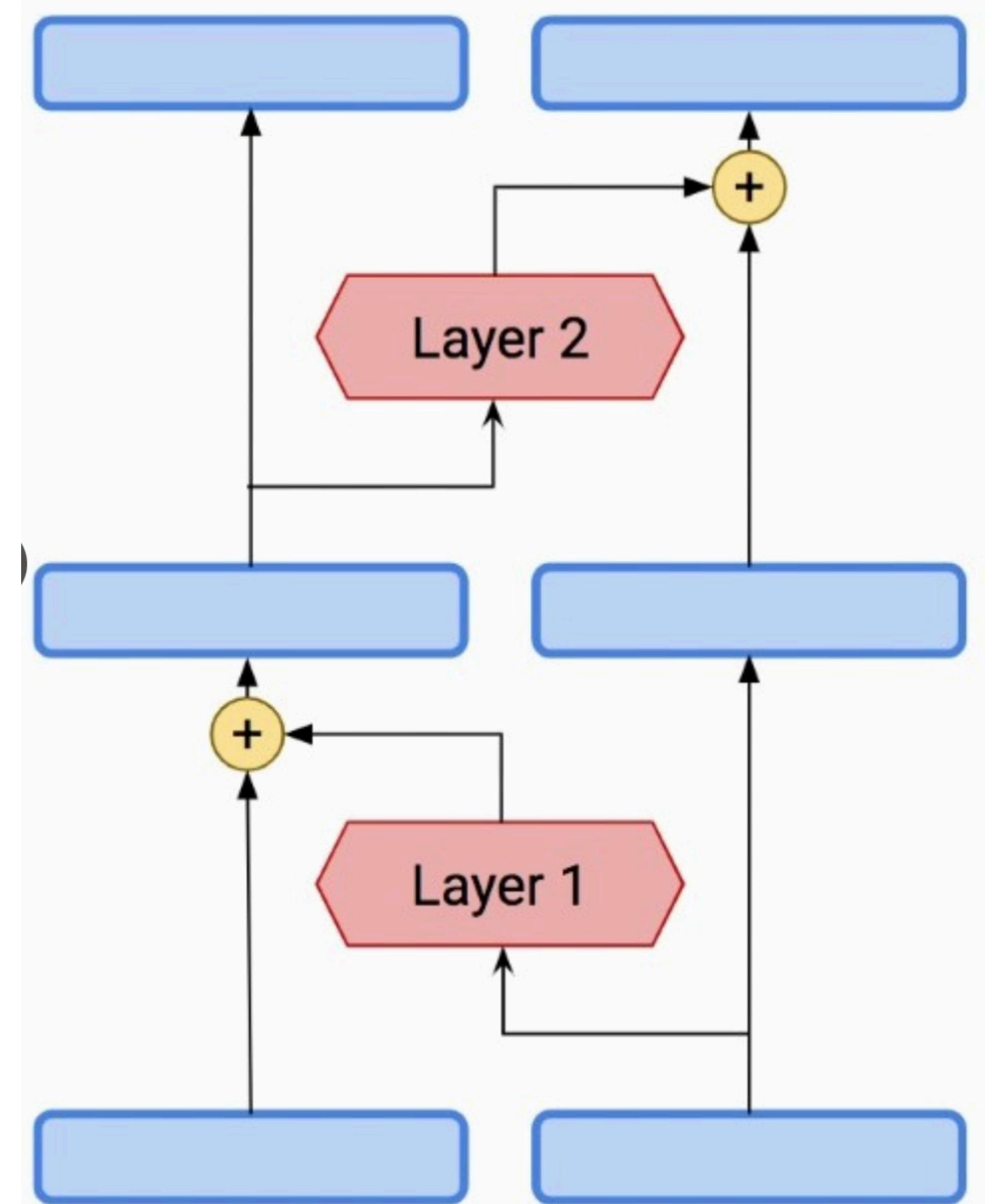
- Memory Efficiency
 - Reduce memory usage with reversible residual layers, as in RevNet [Gomez+ 17]
- Time Complexity
 - Introduce fast attention with locality sensitive hashing (LSH)











output

~ 4 GB

Feed Forward

⋮

~ 4 GB

Feed Forward

~ 4 GB

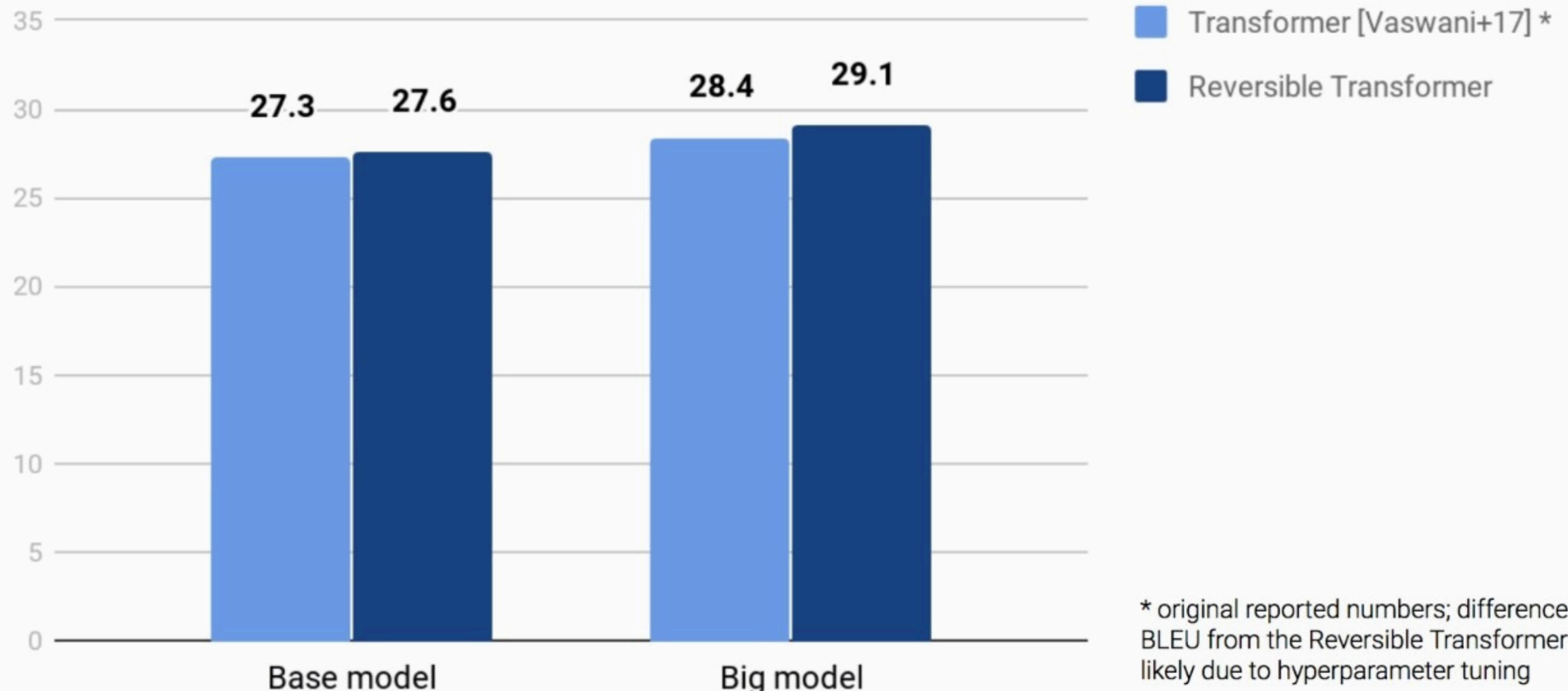
Attention

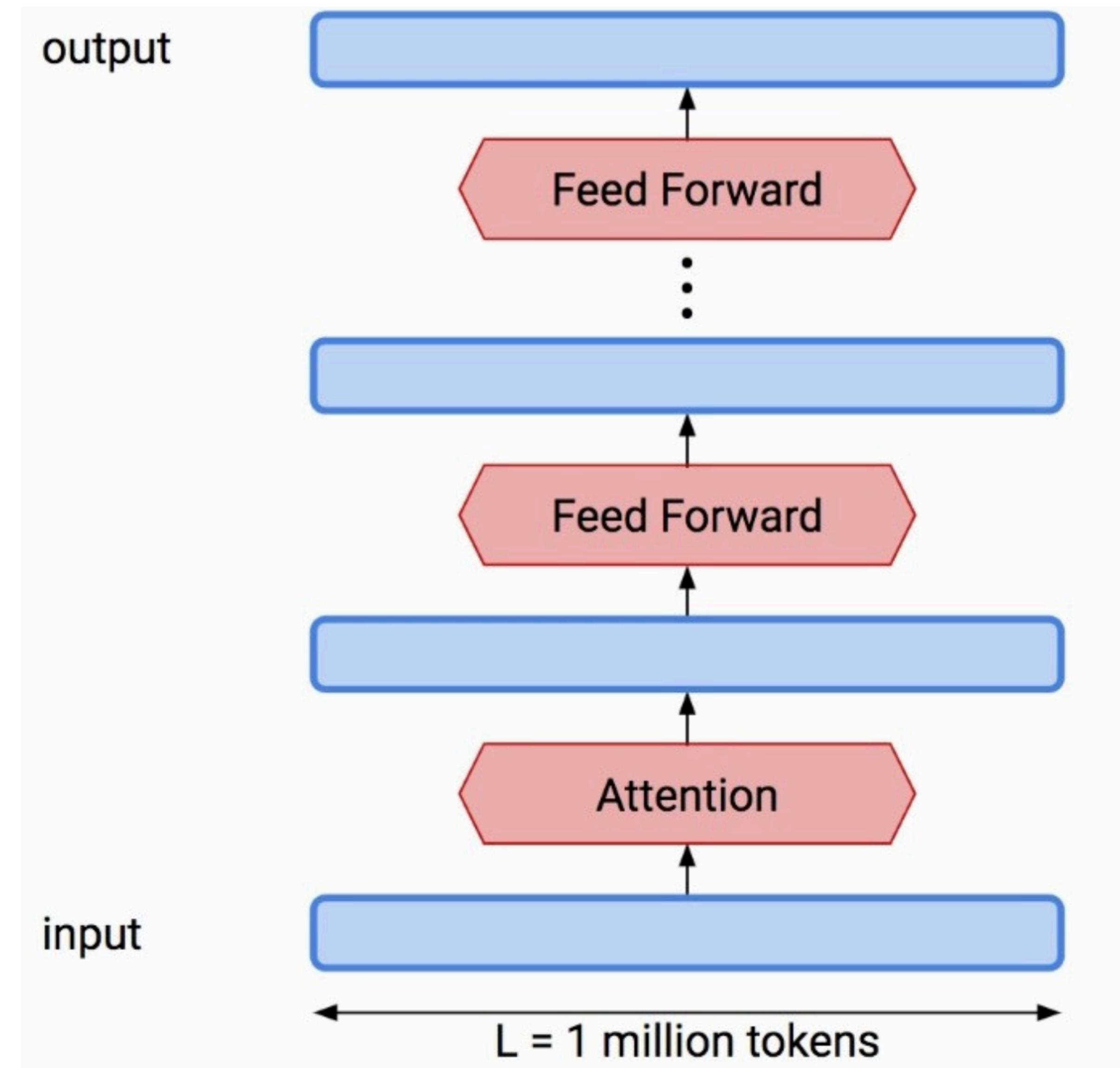
input

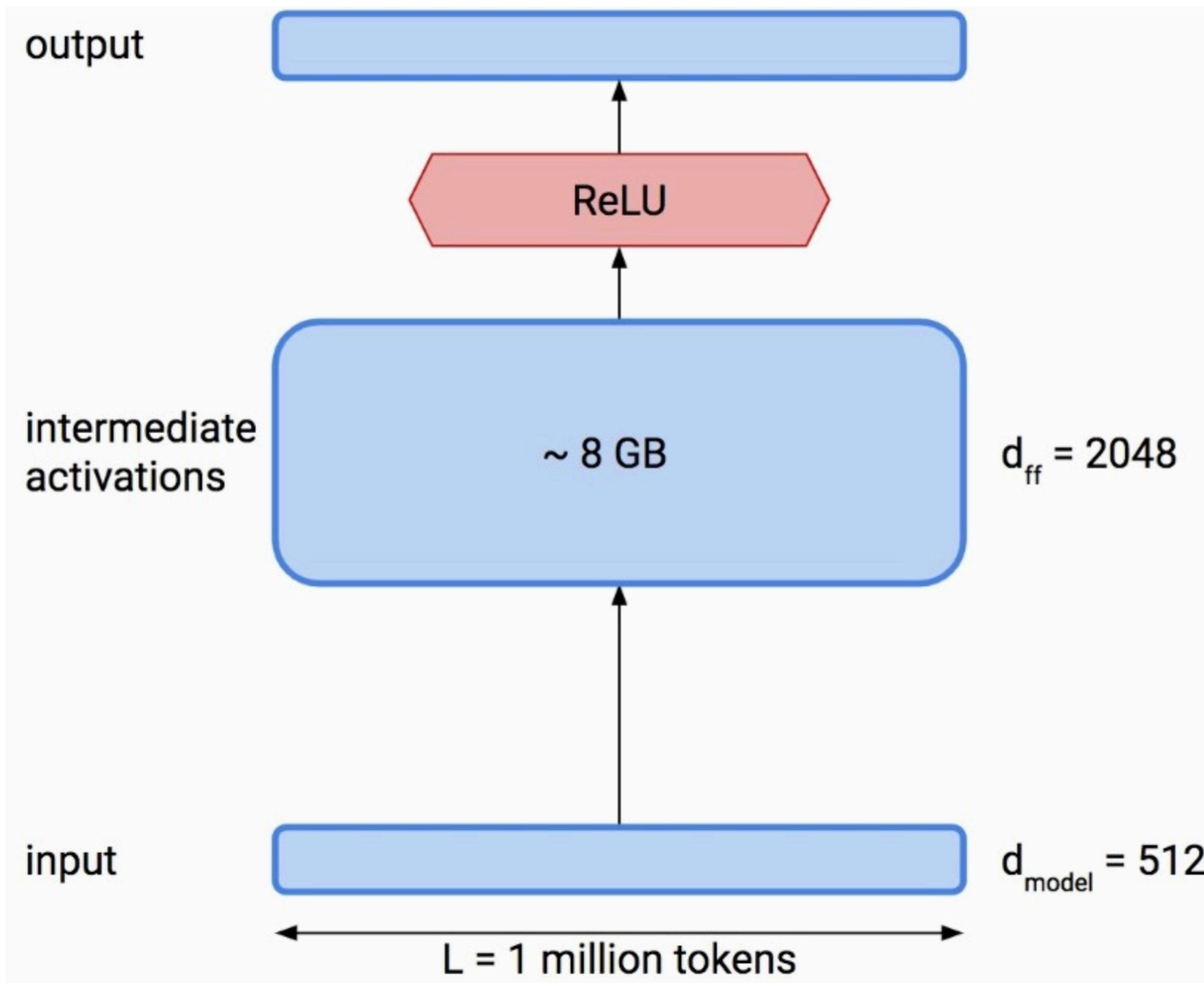
~ 4 GB

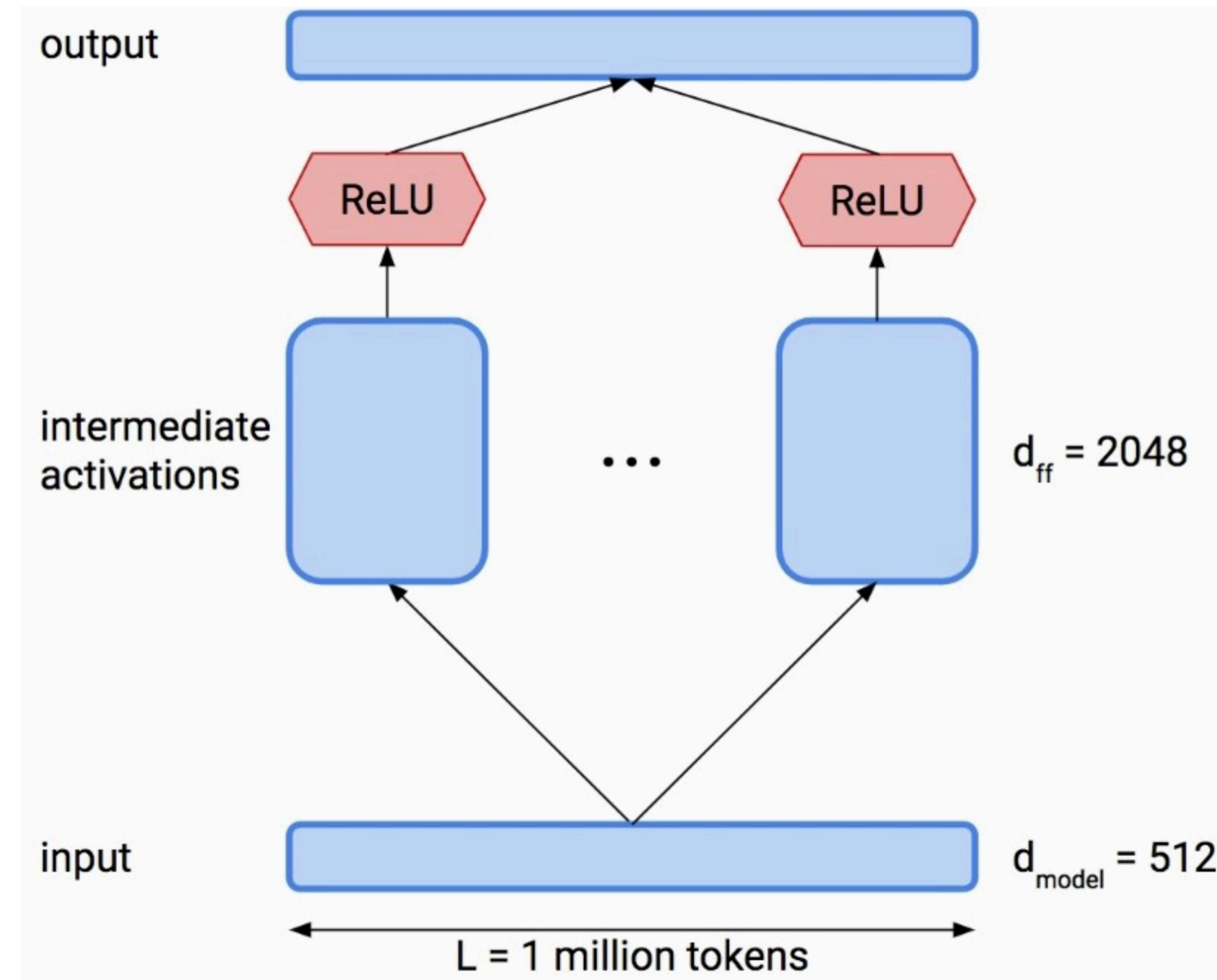
No caching needed when
using reversible layers

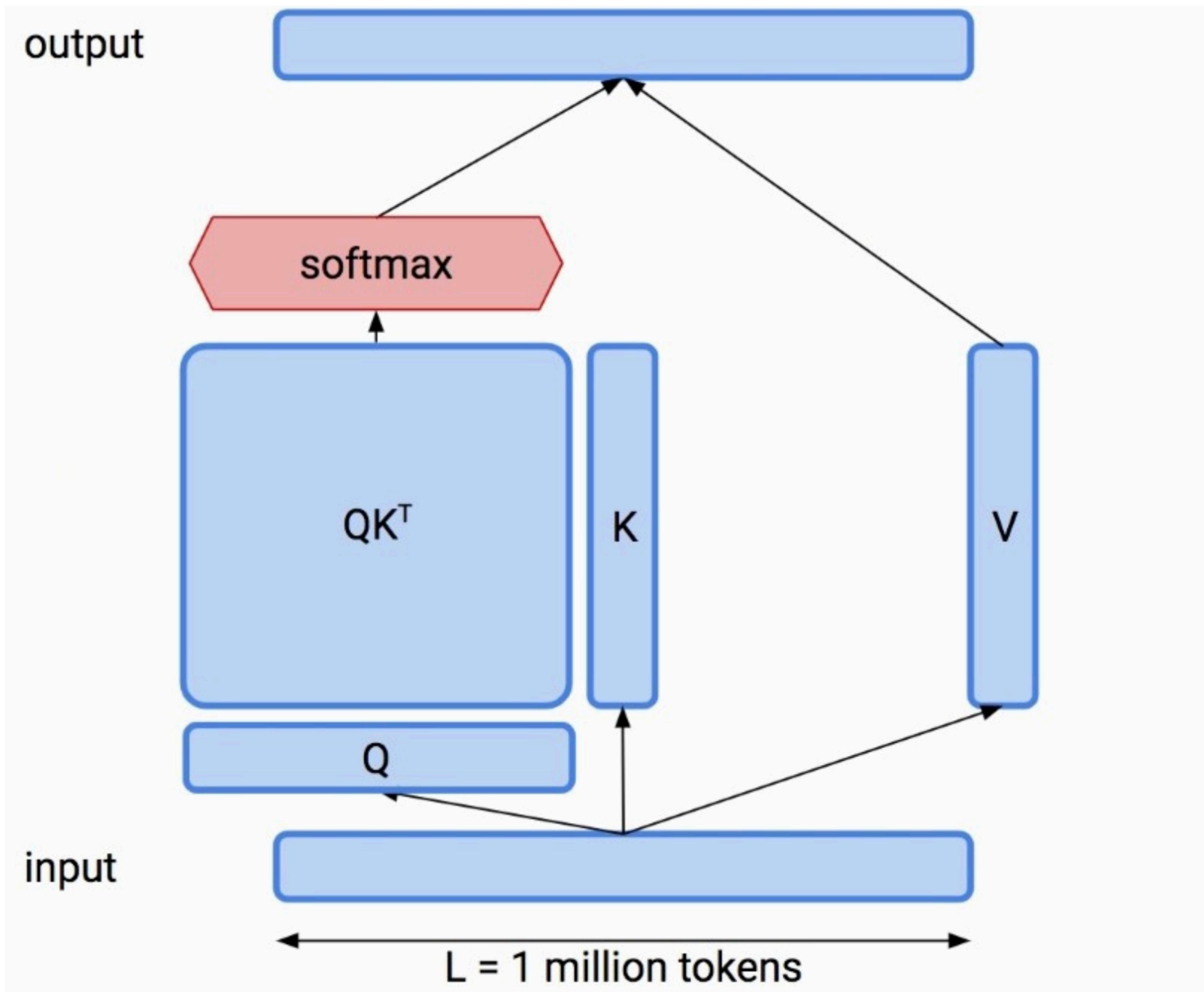
Reversible Transformer: BLEU Scores on WMT English-German

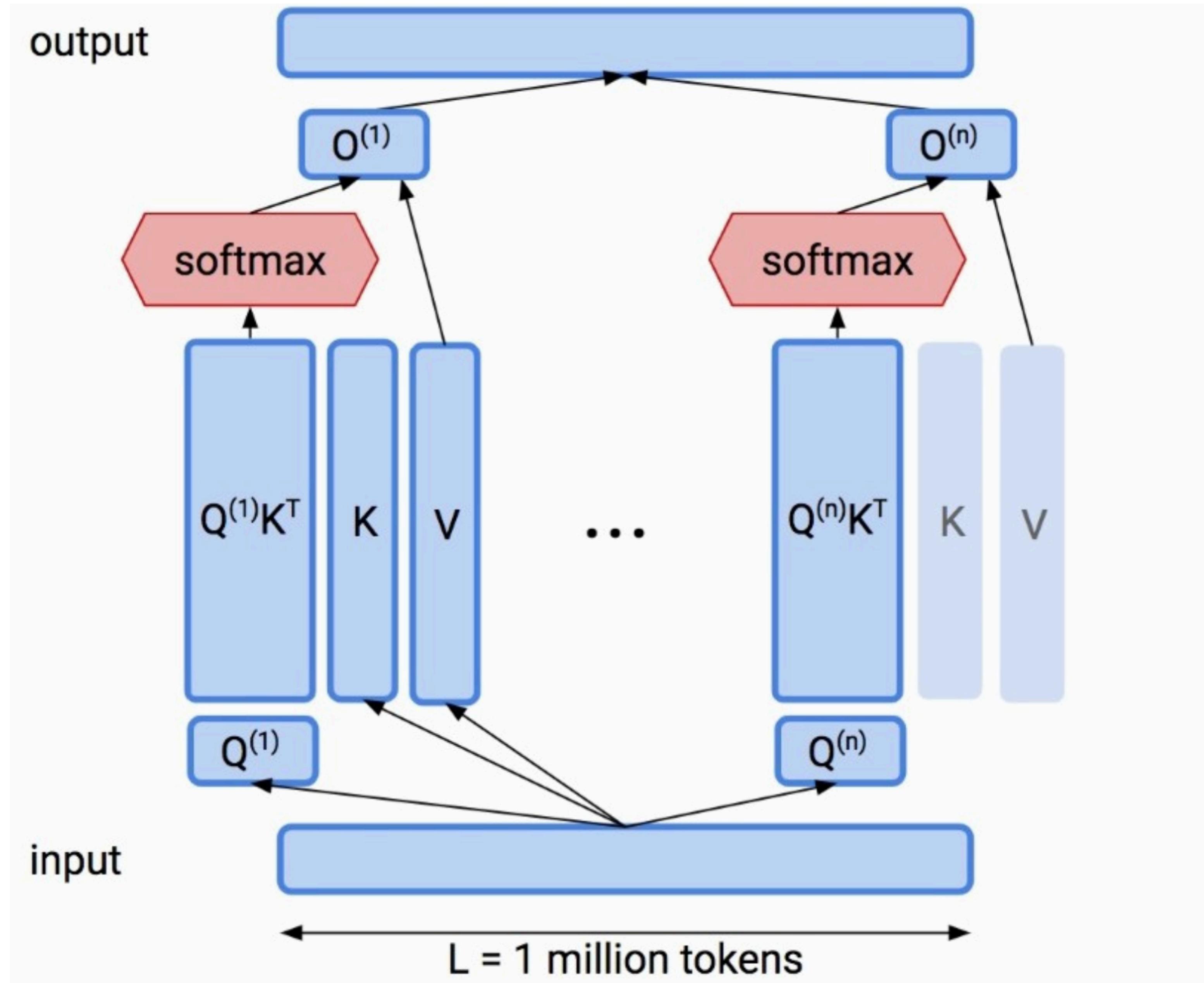


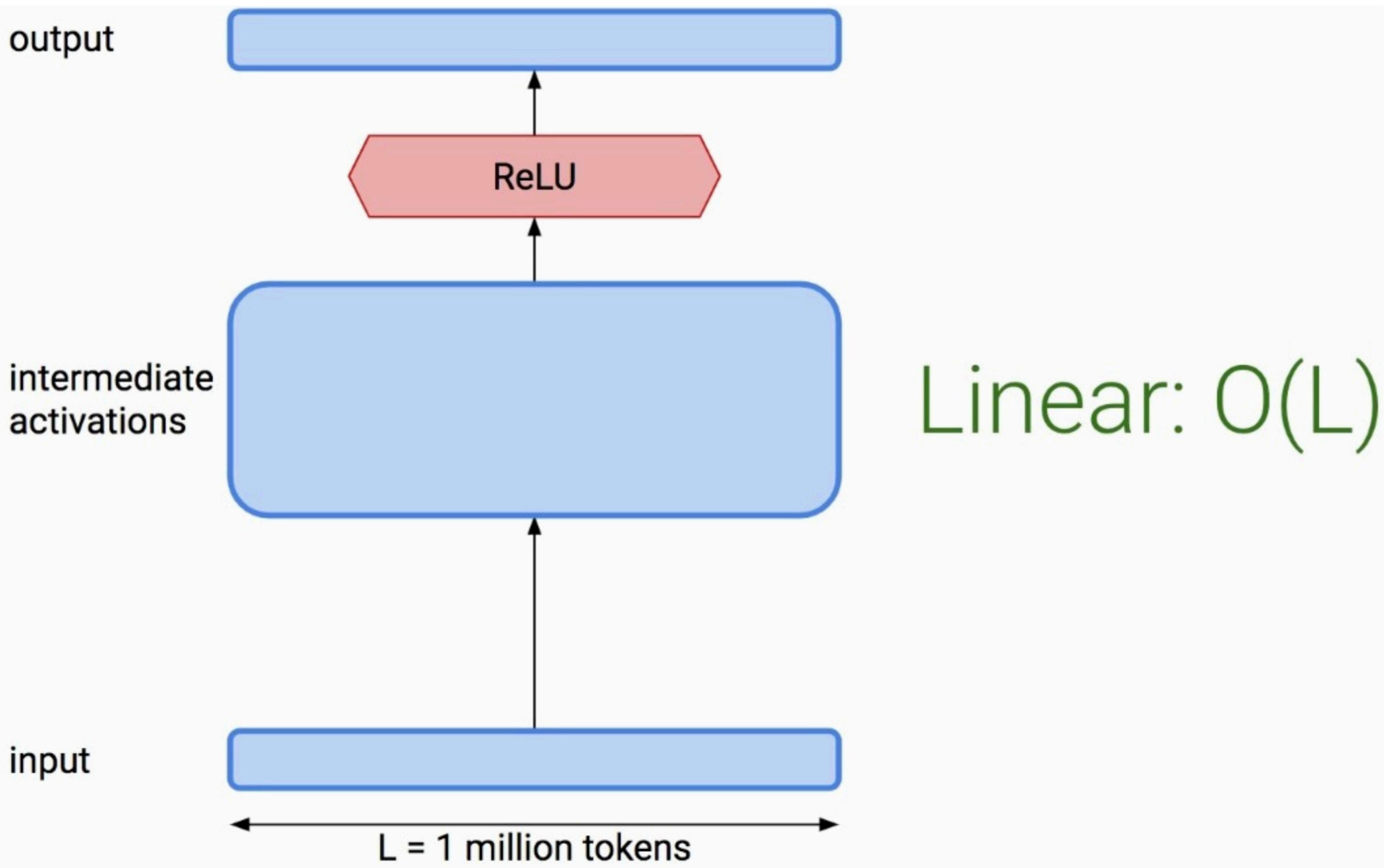


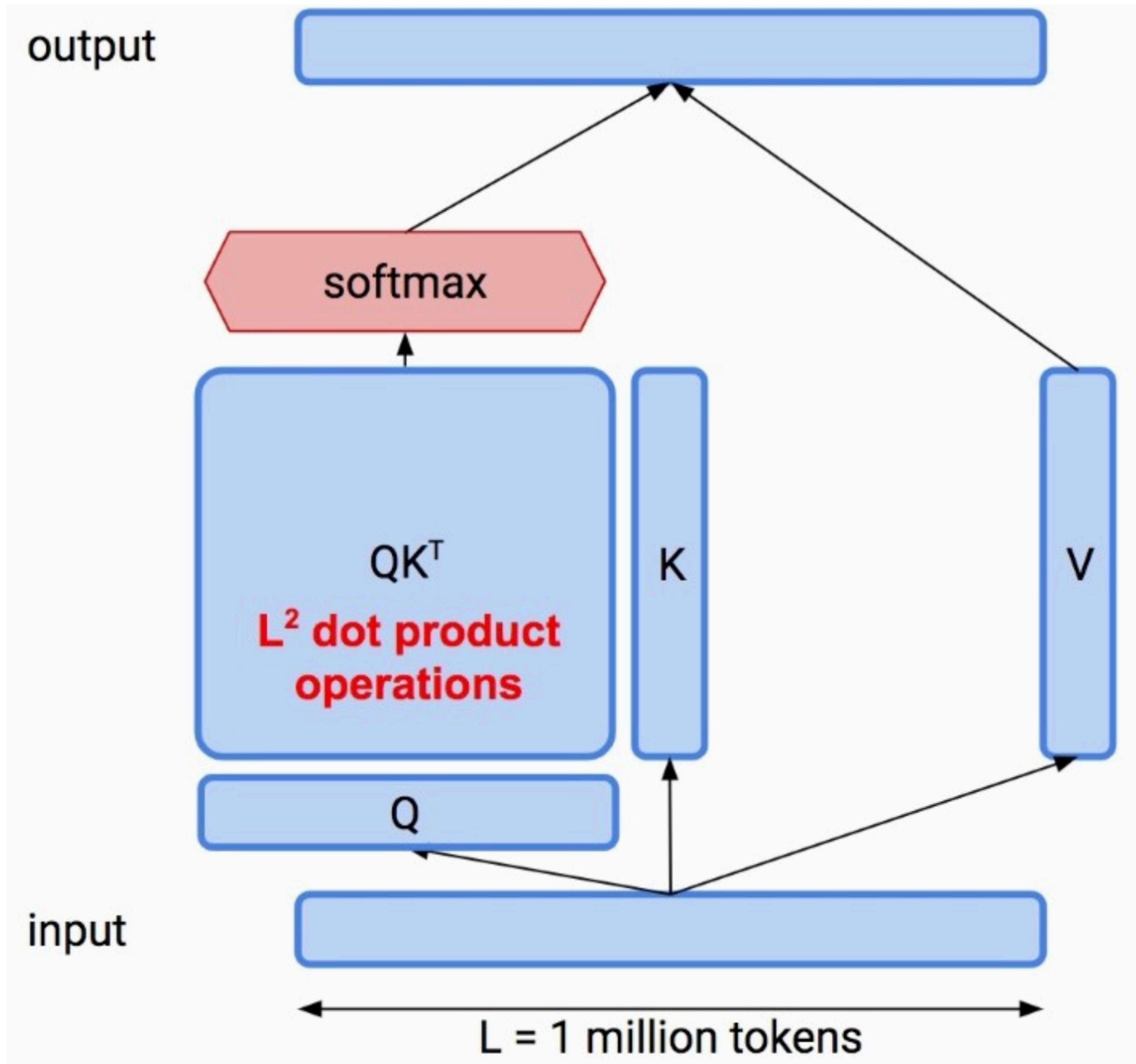


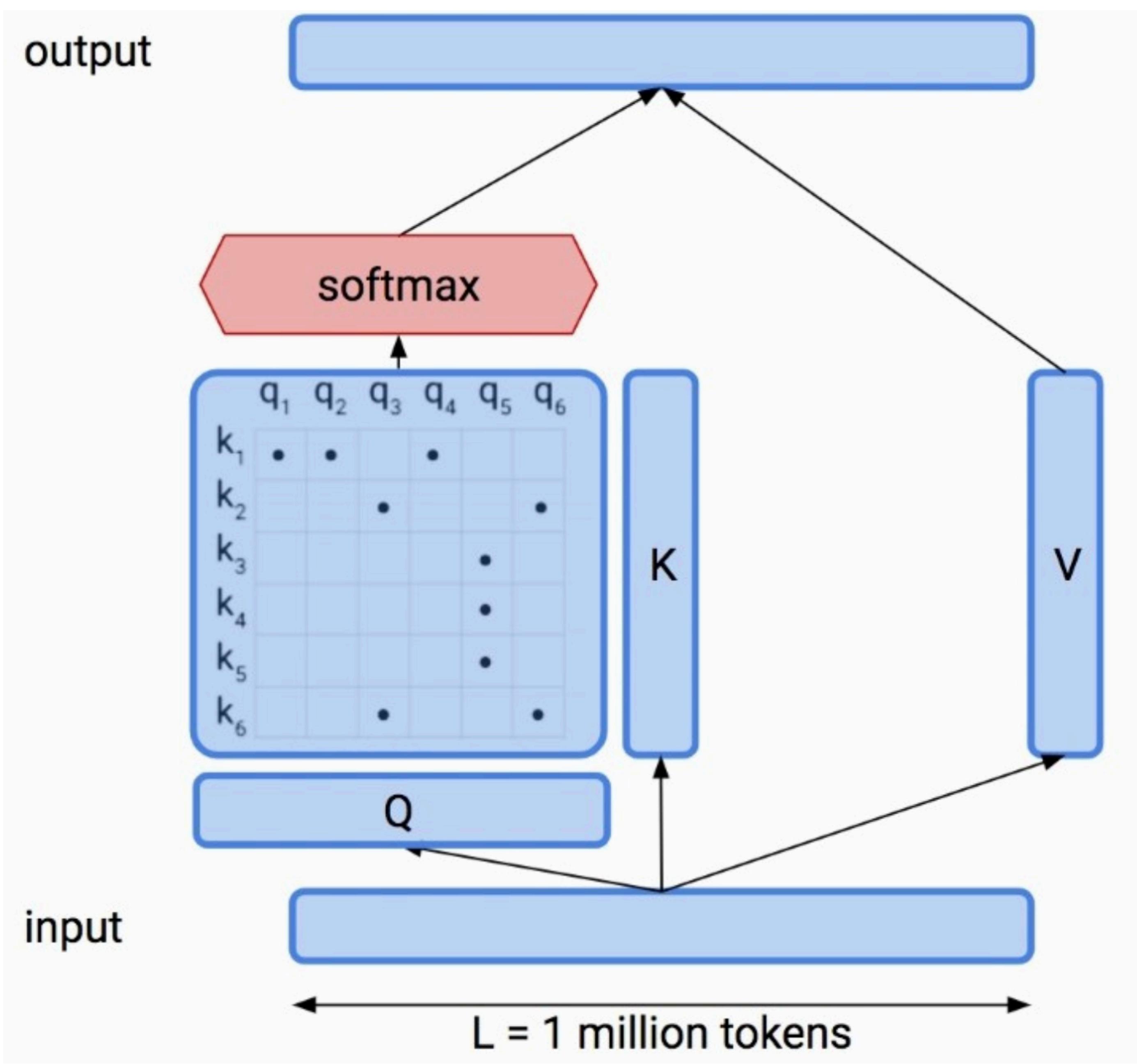


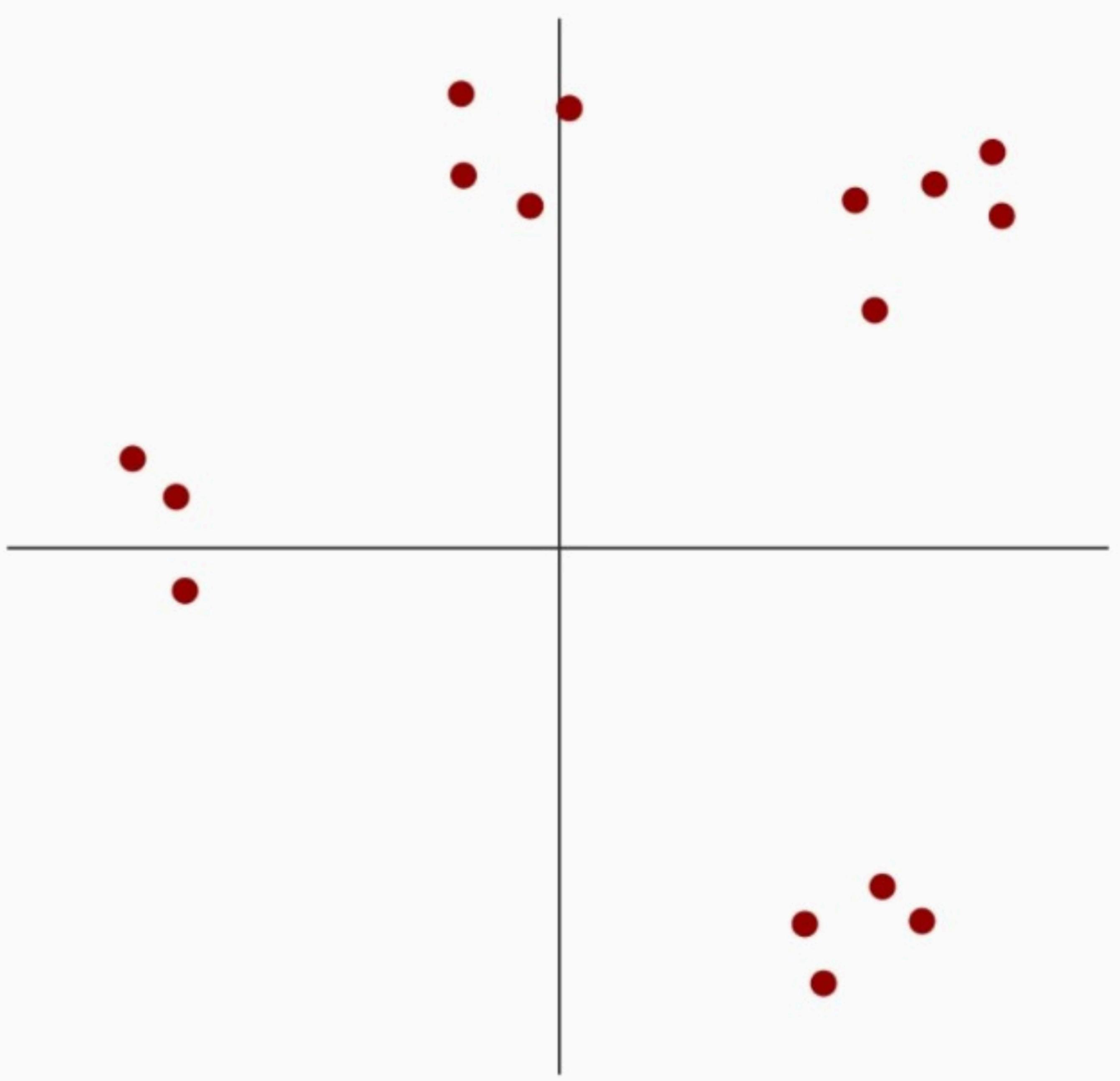


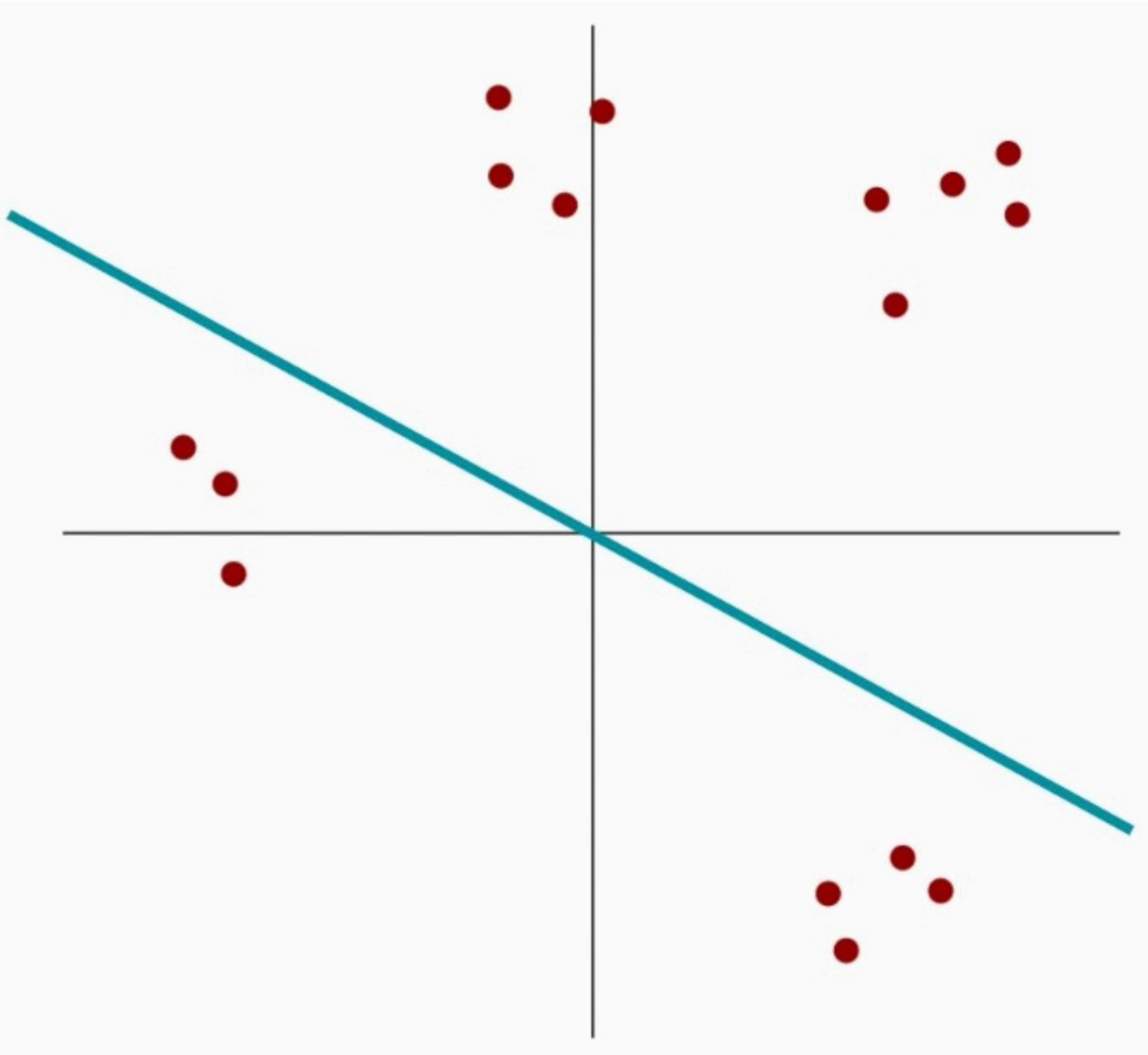


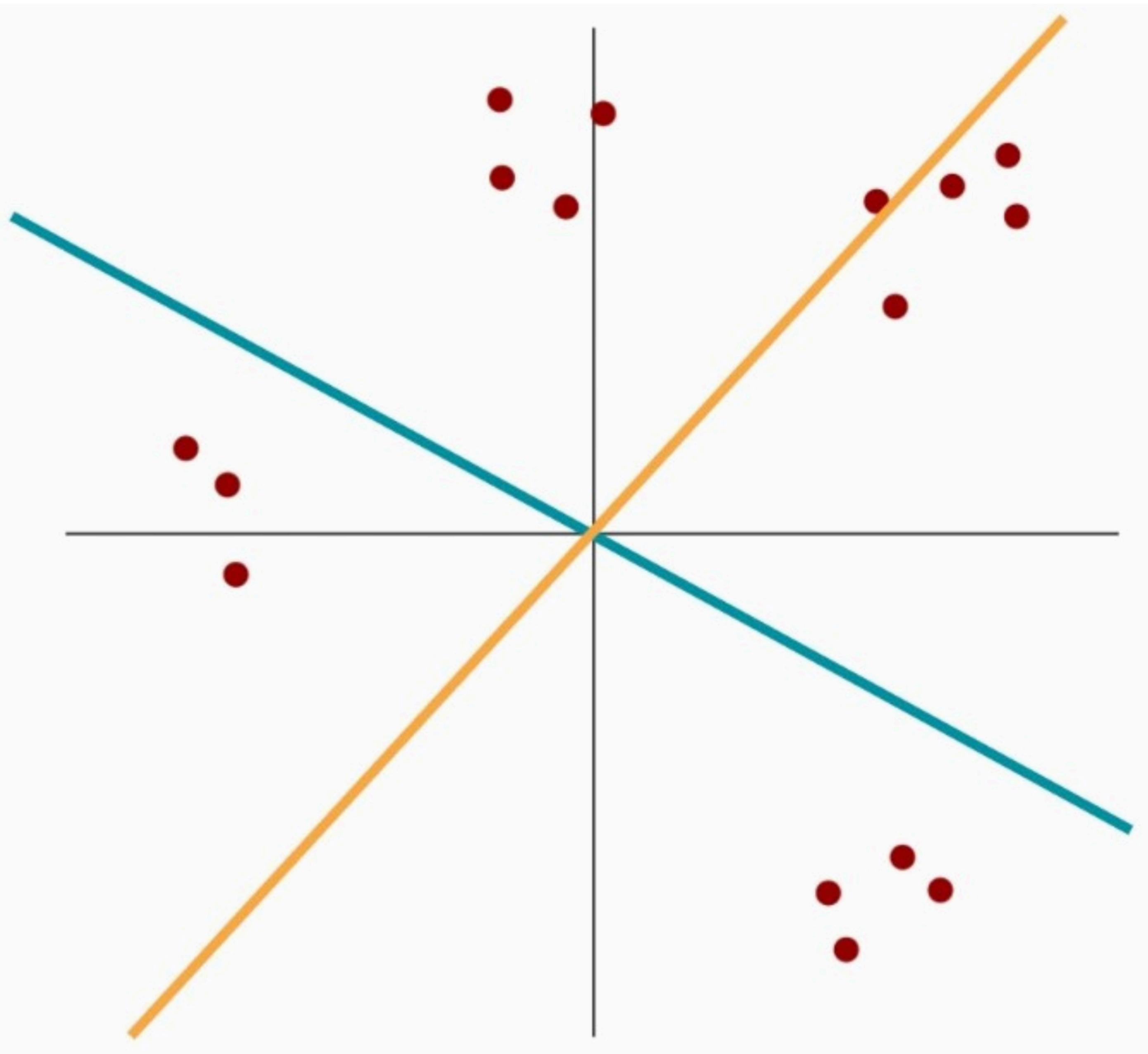


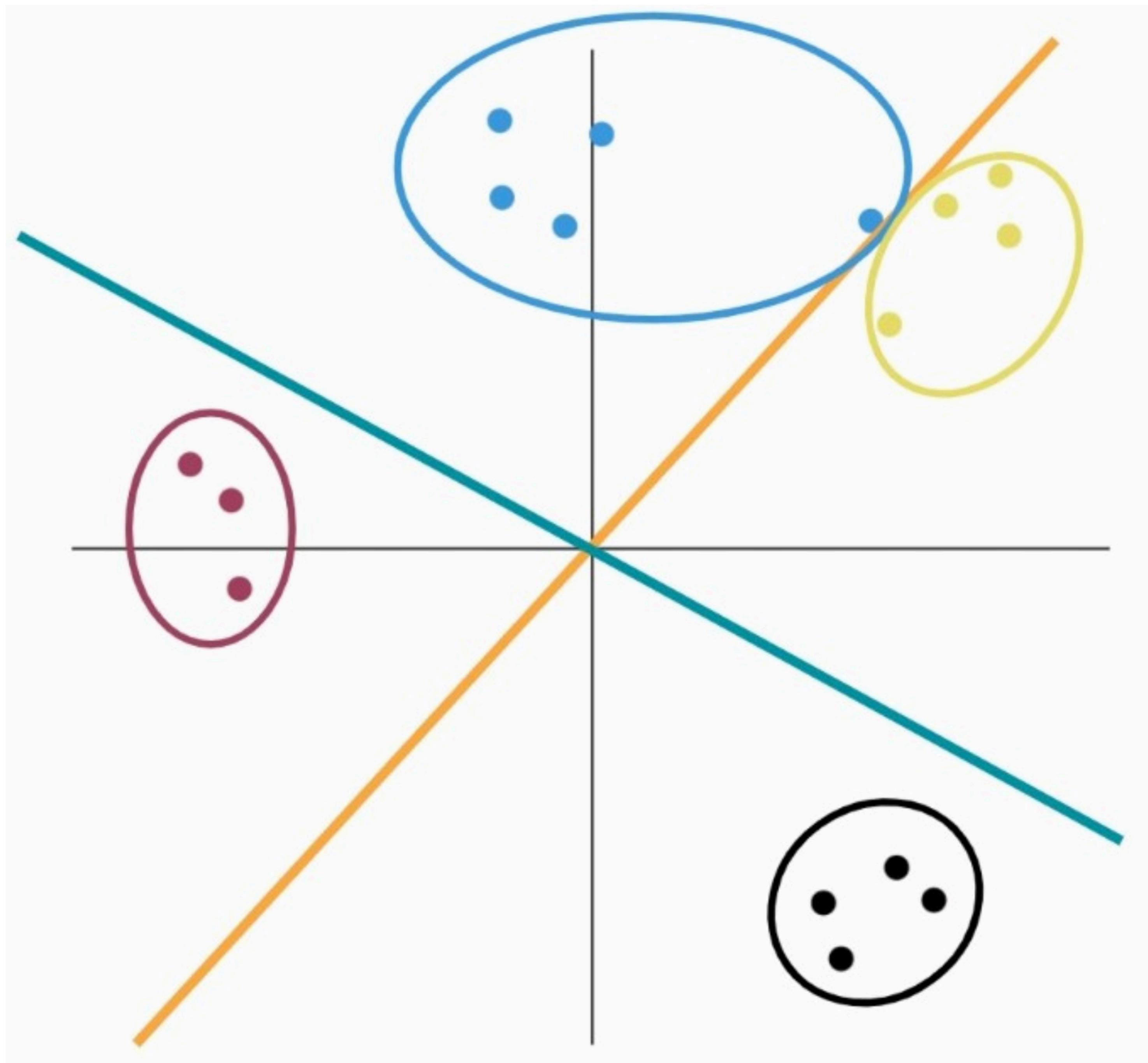


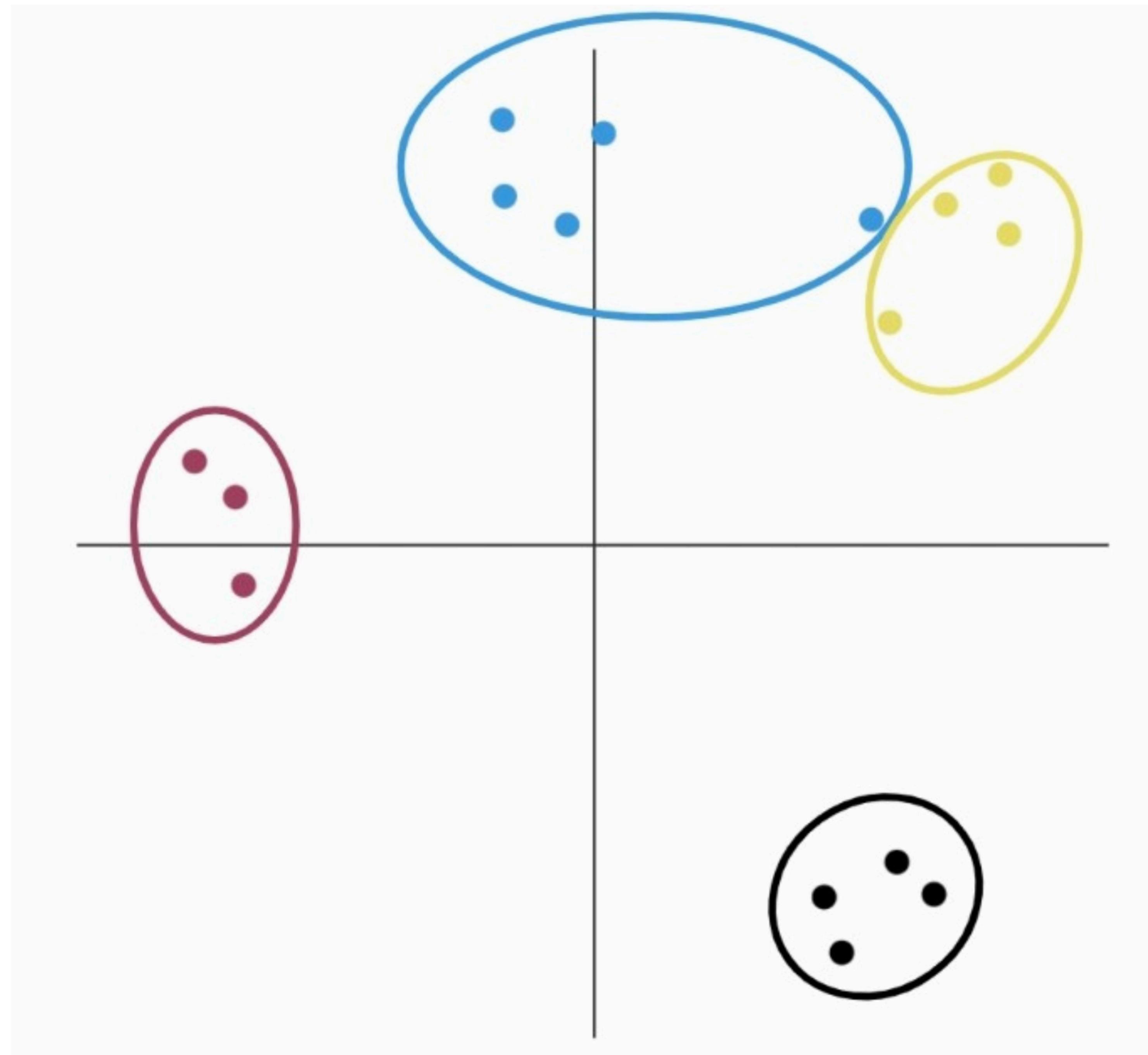












Sequence
of queries=keys



LSH bucketing

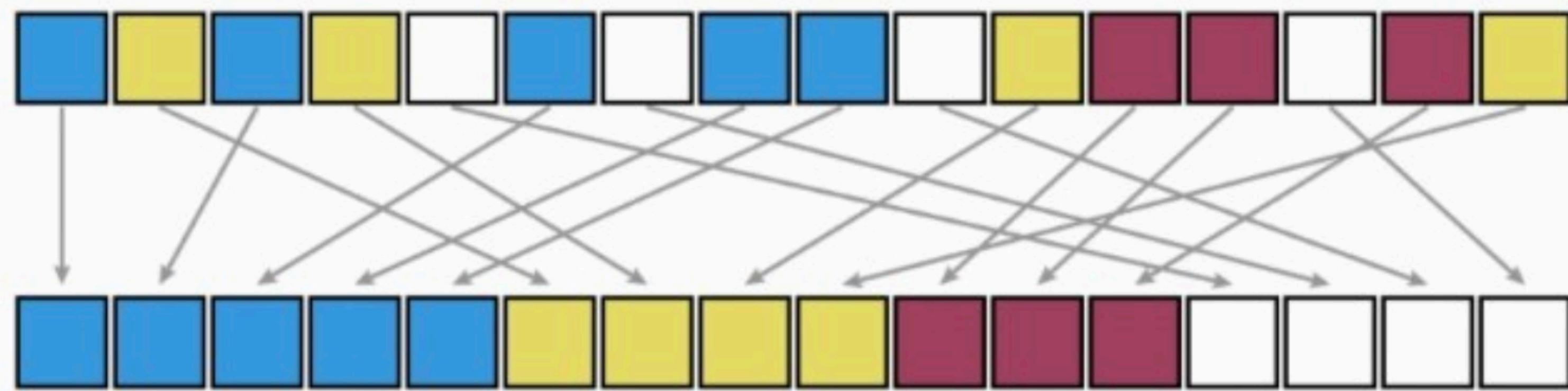


Sequence
of queries=keys



LSH bucketing

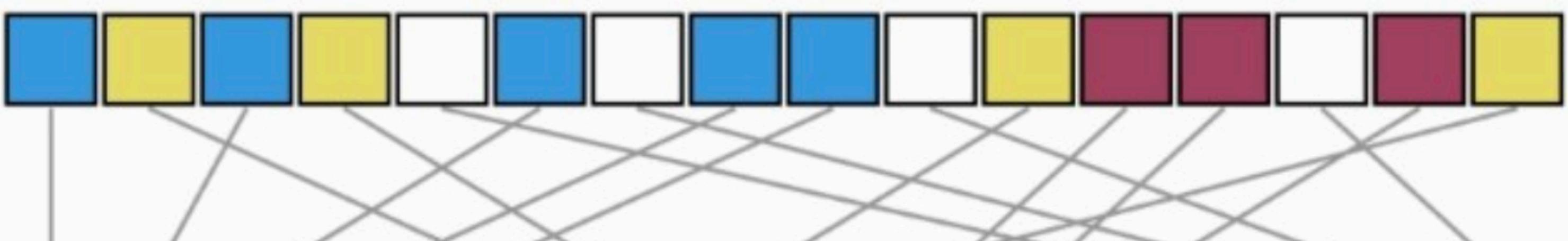
Sort by LSH bucket



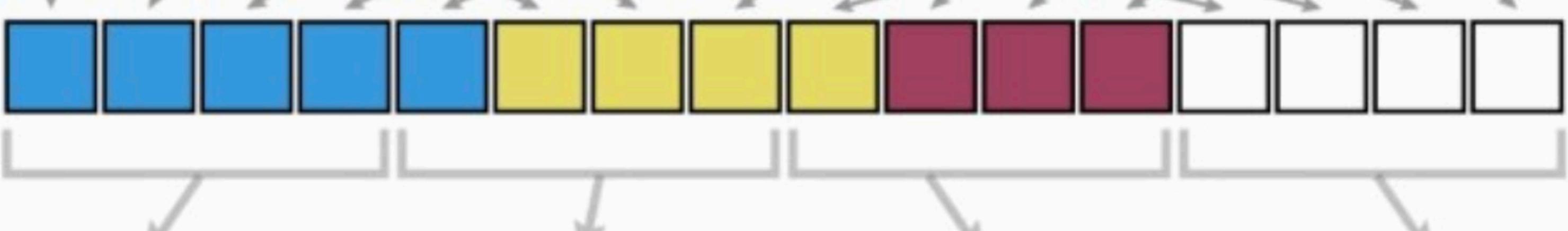
Sequence
of queries=keys



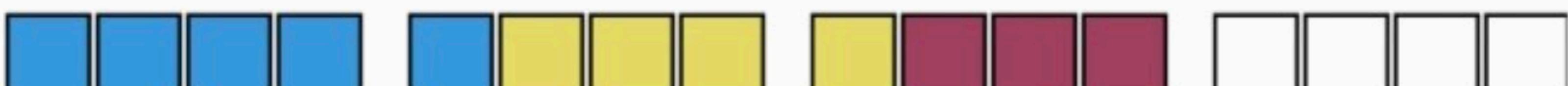
LSH bucketing



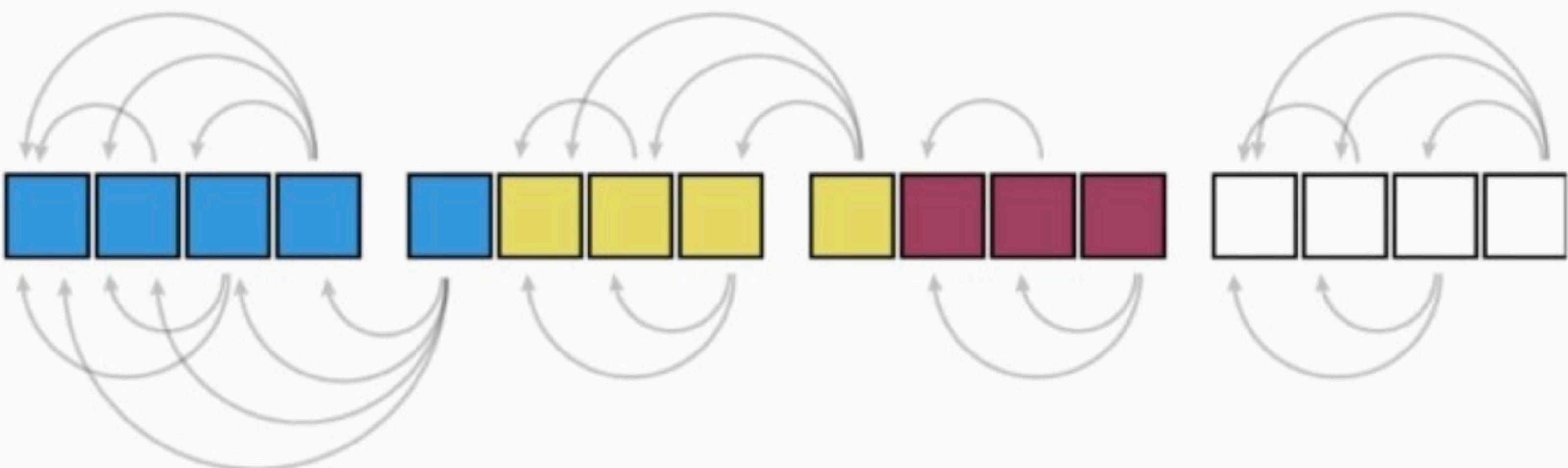
Sort by LSH bucket

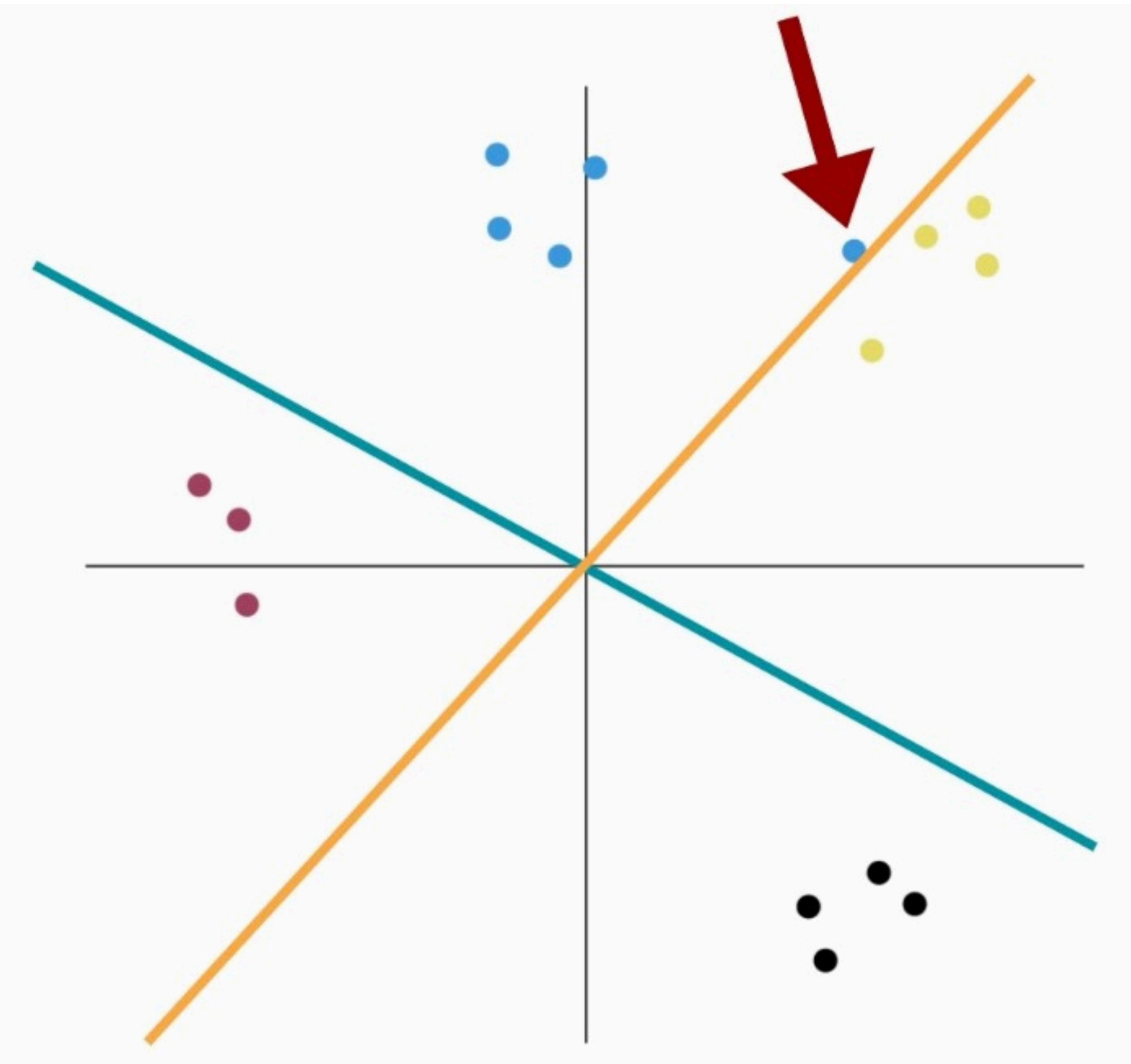


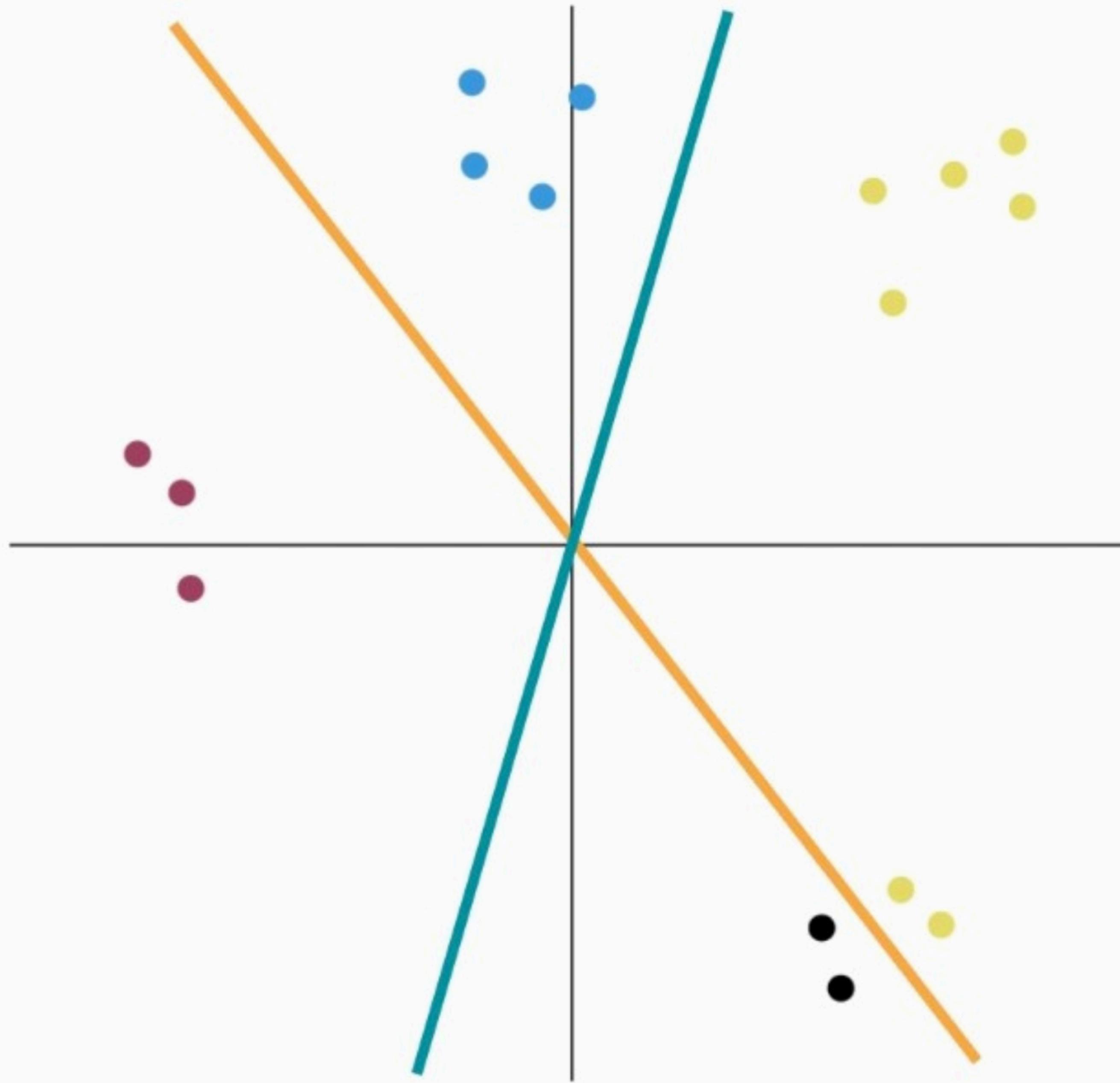
Chunk sorted
sequence to
parallelize



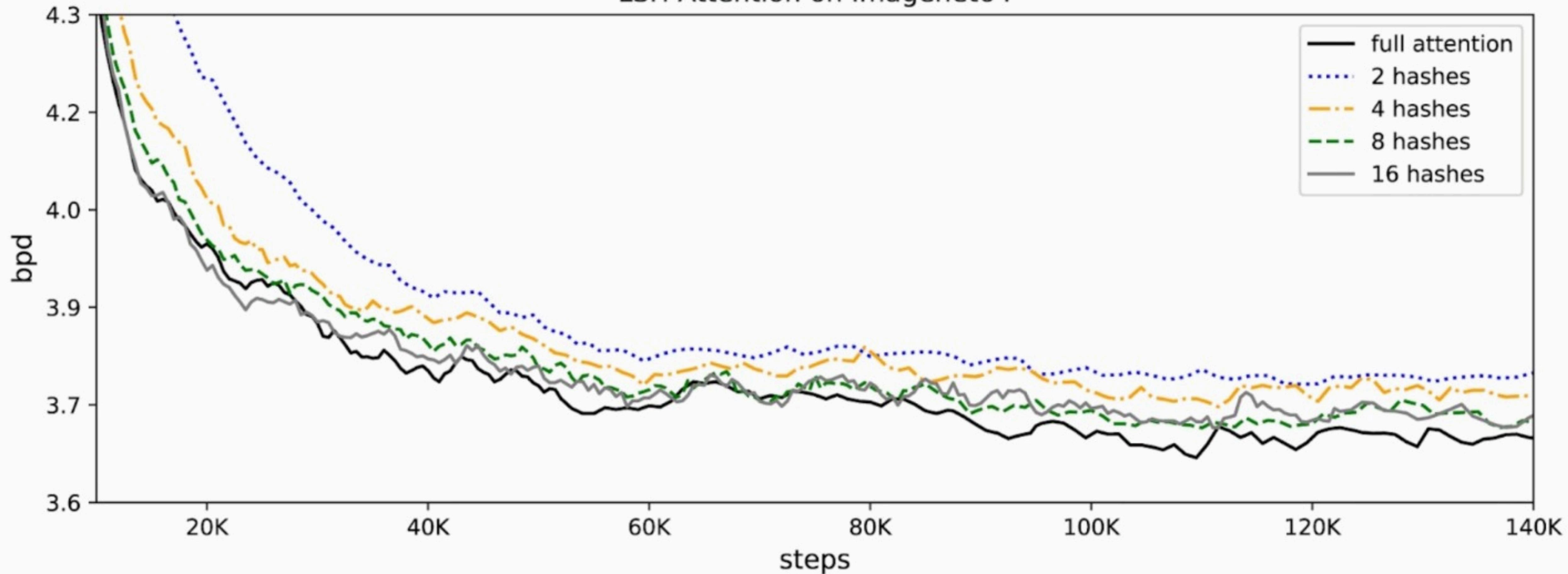
Attend within
same bucket in
own chunk and
previous chunk

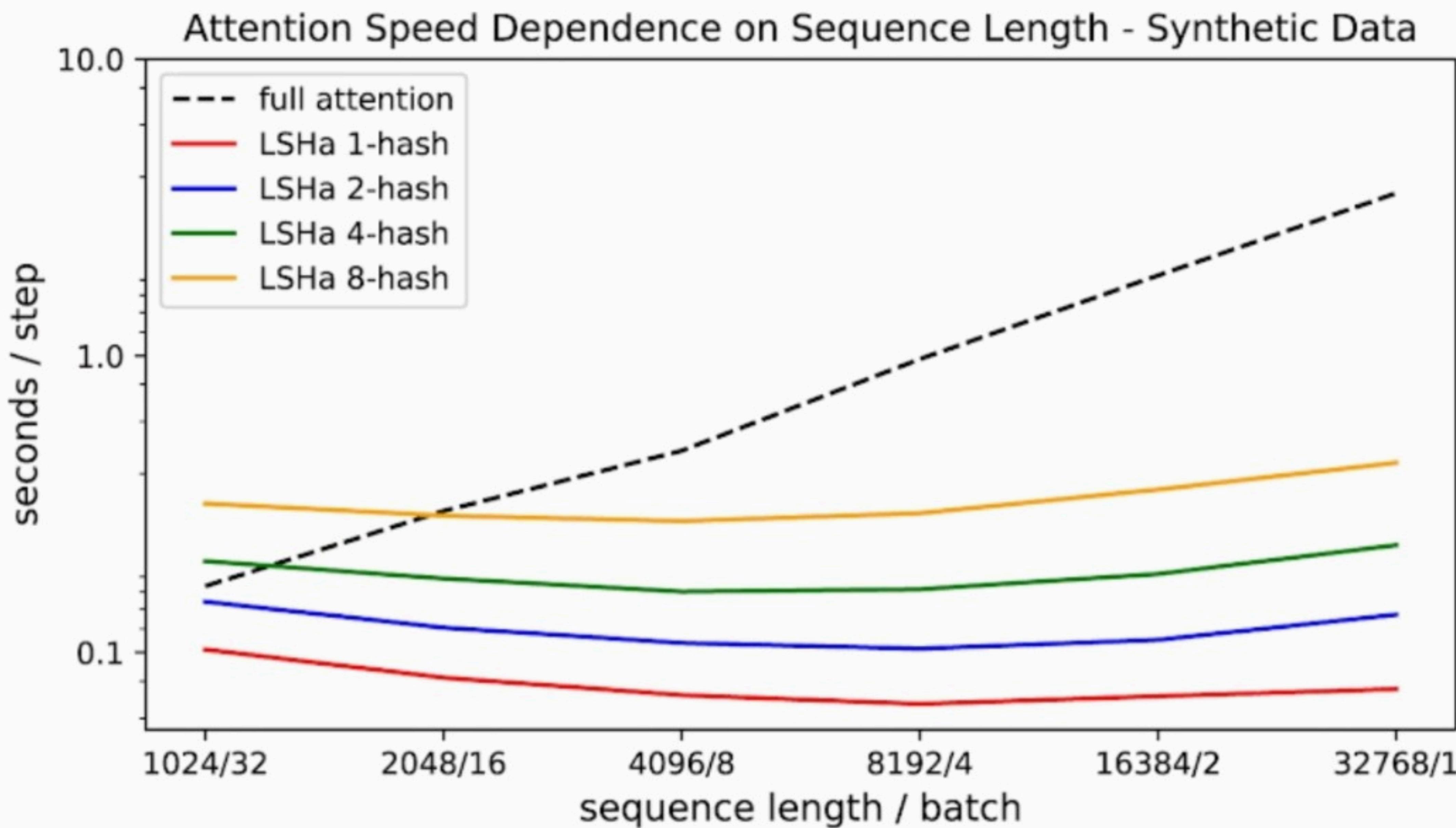






LSH Attention on Imagenet64





- Reversible Layers and Chunking make the Transformer memory-efficient without sacrificing accuracy
- LSH Attention approximates *global attention* with $O(L \log L)$ time complexity

Efficient Transformers: A Survey

Yi Tay

Google Research

YITAY@GOOGLE.COM

Mostafa Dehghani

Google Research, Brain team

DEHGHANI@GOOGLE.COM

Dara Bahri

Google Research

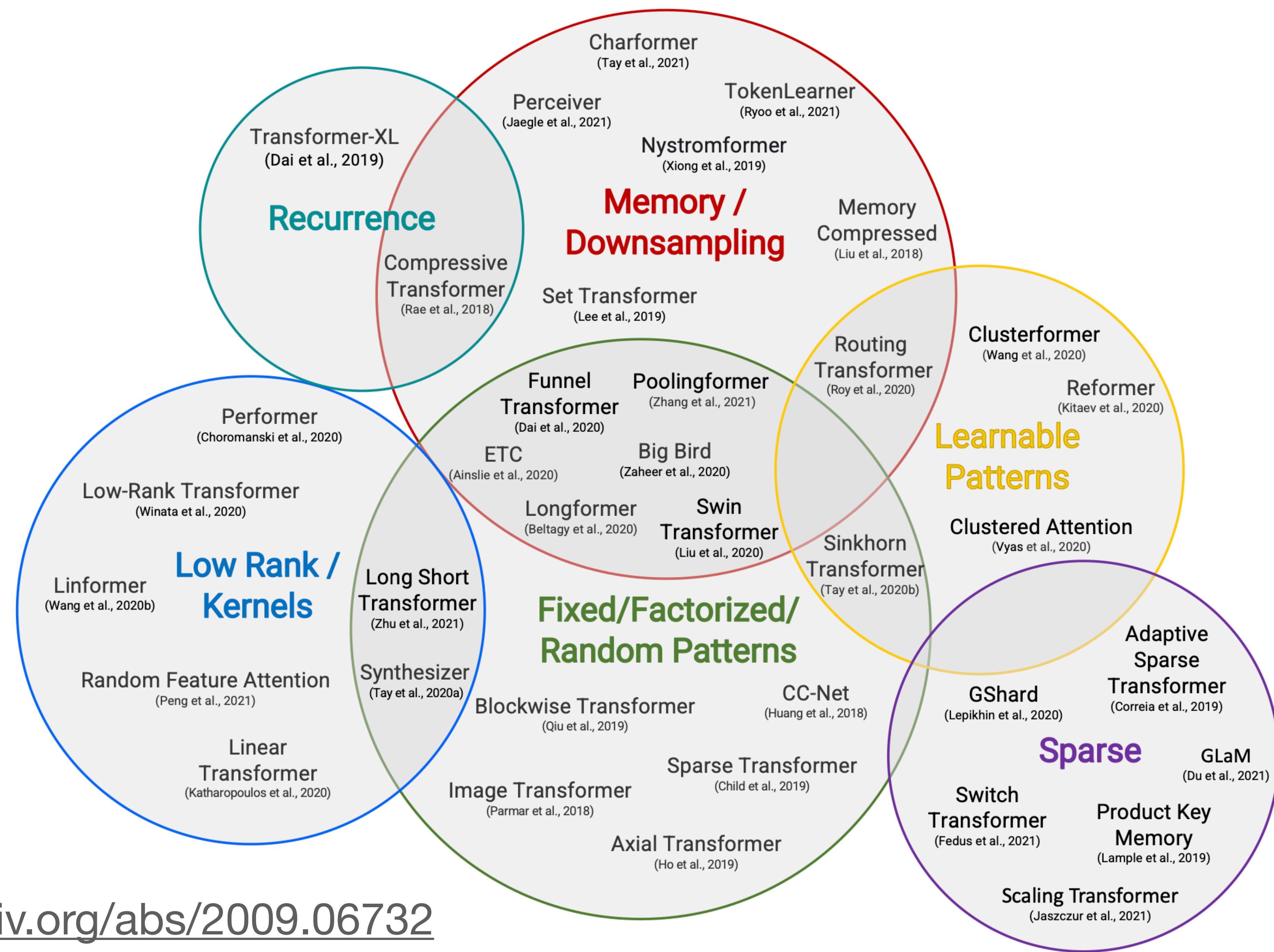
DBAHRI@GOOGLE.COM

Donald Metzler

Google Research

METZLER@GOOGLE.COM

Editor: Preprint, Version 2, Updated Mar 2022



Model / Paper	Complexity	Decode	Class
Memory Compressed (Liu et al., 2018)	$\mathcal{O}(N_c^2)$	✓	FP+M
Image Transformer (Parmar et al., 2018)	$\mathcal{O}(N.m)$	✓	FP
Set Transformer (Lee et al., 2019)	$\mathcal{O}(kN)$	✗	M
Transformer-XL (Dai et al., 2019)	$\mathcal{O}(N^2)$	✓	RC
Sparse Transformer (Child et al., 2019)	$\mathcal{O}(N\sqrt{N})$	✓	FP
Reformer (Kitaev et al., 2020)	$\mathcal{O}(N \log N)$	✓	LP
Routing Transformer (Roy et al., 2020)	$\mathcal{O}(N\sqrt{N})$	✓	LP
Axial Transformer (Ho et al., 2019)	$\mathcal{O}(N\sqrt{N})$	✓	FP
Compressive Transformer (Rae et al., 2020)	$\mathcal{O}(N^2)$	✓	RC
Sinkhorn Transformer (Tay et al., 2020b)	$\mathcal{O}(B^2)$	✓	LP
Longformer (Beltagy et al., 2020)	$\mathcal{O}(n(k + m))$	✓	FP+M
ETC (Ainslie et al., 2020)	$\mathcal{O}(N_g^2 + NN_g)$	✗	FP+M
Synthesizer (Tay et al., 2020a)	$\mathcal{O}(N^2)$	✓	LR+LP
Performer (Choromanski et al., 2020a)	$\mathcal{O}(N)$	✓	KR
Funnel Transformer (Dai et al., 2020)	$\mathcal{O}(N^2)$	✓	FP+DS

Table 1: Summary of Efficient Transformer Models. Models in the first section are mainly efficient attention methods. Models in the subsequent lower section generally refer to sparse models. Class abbreviations include: FP = Fixed Patterns or Combinations of Fixed Patterns, M = Memory, LP = Learnable Pattern, LR = Low-Rank, KR = Kernel RC = Recurrence, and DS = Downsampling. Furthermore, N generally refers to the sequence length and B is the local window (or block) size. N_g and N_c denote global model memory length and convolutionally-compressed sequence lengths respectively.

Linformer (Wang et al., 2020c)	$\mathcal{O}(N)$	\times	LR
Linear Transformers (Katharopoulos et al., 2020)	$\mathcal{O}(N)$	\checkmark	KR
Big Bird (Zaheer et al., 2020)	$\mathcal{O}(N)$	\times	FP+M
Random Feature Attention (Peng et al., 2021)	$\mathcal{O}(N)$	\checkmark	KR
Long Short Transformers (Zhu et al., 2021)	$\mathcal{O}(kN)$	\checkmark	FP + LR
Poolingformer (Zhang et al., 2021)	$\mathcal{O}(N)$	\times	FP+M
Nyströmformer (Xiong et al., 2021b)	$\mathcal{O}(kN)$	\times	M+DS
Perceiver (Jaegle et al., 2021)	$\mathcal{O}(kN)$	\checkmark	M+DS
Clusterformer (Wang et al., 2020b)	$\mathcal{O}(N \log N)$	\times	LP
Luna (Ma et al., 2021)	$\mathcal{O}(kN)$	\checkmark	M
TokenLearner (Ryoo et al., 2021)	$\mathcal{O}(k^2)$	\times	DS

Table 1: Summary of Efficient Transformer Models. Models in the first section are mainly efficient attention methods. Models in the subsequent lower section generally refer to sparse models. Class abbreviations include: FP = Fixed Patterns or Combinations of Fixed Patterns, M = Memory, LP = Learnable Pattern, LR = Low-Rank, KR = Kernel RC = Recurrence, and DS = Downsampling. Furthermore, N generally refers to the sequence length and B is the local window (or block) size. N_g and N_c denote global model memory length and convolutionally-compressed sequence lengths respectively.

Adaptive Sparse Transformer (Correia et al., 2019)	$\mathcal{O}(N^2)$	✓	Sparse
Product Key Memory (Lample et al., 2019)	$\mathcal{O}(N^2)$	✓	Sparse
Switch Transformer (Fedus et al., 2021)	$\mathcal{O}(N^2)$	✓	Sparse
ST-MoE (Zoph et al., 2022)	$\mathcal{O}(N^2)$	✓	Sparse
GShard (Lepikhin et al., 2020)	$\mathcal{O}(N^2)$	✓	Sparse
Scaling Transformers (Jaszcjur et al., 2021)	$\mathcal{O}(N^2)$	✓	Sparse
GLaM (Du et al., 2021)	$\mathcal{O}(N^2)$	✓	Sparse

- **Fixed Patterns (FP)** - The earliest modifications to self-attention simply sparsifies the attention matrix by limiting the field of view to fixed, predefined patterns such as local windows and block patterns of fixed strides.
 - **Blockwise Patterns** The simplest example of this technique in practice is the blockwise (or chunking) paradigm which considers blocks of local receptive fields by chunking input sequences into fixed blocks. Examples of models that do this include Blockwise (Qiu et al., 2019) and/or Local Attention (Parmar et al., 2018). Chunking input sequences into blocks reduces the complexity from N^2 to B^2 (block size) with $B \ll N$, significantly reducing the cost. These blockwise or chunking methods serve as a basis for many more complex models.
 - **Strided Patterns** Another approach is to consider strided attention patterns, i.e., only attending at fixed intervals. Models such as Sparse Transformer (Child et al., 2019) and/or Longformer (Beltagy et al., 2020) employ strided or “dilated” windows.
 - **Compressed Patterns** - Another line of attack here is to use some pooling operator to down-sample the sequence length to be a form of fixed pattern. For instance, Compressed Attention (Liu et al., 2018) uses strided convolution to effectively reduce the sequence length.

- **Recurrence** - A natural extension to the blockwise method is to connect these blocks via recurrence. Transformer-XL (Dai et al., 2019) proposed a segment-level recurrence mechanism that connects multiple segments and blocks. These models can, in some sense, be viewed as *fixed pattern* models. However, we decided to create its own category due to its deviation from other block / local approaches.

- **Combination of Patterns (CP)** - The key idea of combined² approaches is to improve coverage by combining two or more distinct access patterns. For example, the Sparse Transformer (Child et al., 2019) combines strided and local attention by assigning half of its heads to each pattern. Similarly, Axial Transformer (Ho et al., 2019) applies a sequence of self-attention computations given a high dimensional tensor as input, each along a single axis of the input tensor. In essence, the combination of patterns reduces memory complexity in the same way that fixed patterns does. The difference, however, is that the aggregation and combinaton of multiple patterns improves the overall coverage of the self-attention mechanism.

- **Learnable Patterns (LP)** - An extension to fixed, pre-determined pattern are *learnable* ones. Unsurprisingly, models using learnable patterns aim to learn the access pattern in a data-driven fashion. A key characteristic of learning patterns is to determine a notion of token relevance and then assign tokens to buckets or clusters (Vyas et al., 2020; Wang et al., 2020b). Notably, Reformer (Kitaev et al., 2020) introduces a hash-based similarity measure to efficiently cluster tokens into chunks. In a simlar vein, the Routing Transformer (Roy et al., 2020) employs online k -means clustering on the tokens. Meanwhile, the Sinkhorn Sorting Network (Tay et al., 2020b) exposes the sparsity in attention weights by learning to sort blocks of the input sequence. In all these models, the similarity function is trained end-to-end jointly with the rest of the network. The key idea of learnable patterns is still to exploit fixed patterns (chunked patterns). However, this class of methods learns to sort/cluster the input tokens - enabling a more optimal global view of the sequence while maintaining the efficiency benefits of fixed patterns approaches.

- **Low-Rank Methods** - Another emerging technique is to improve efficiency by leveraging low-rank approximations of the self-attention matrix. The key idea is to assume low-rank structure in the $N \times N$ matrix. The Linformer (Wang et al., 2020c) is a classic example of this technique, as it projects the length dimension of keys and values to a lower-dimensional representation ($N \rightarrow k$). It is easy to see that the low-rank method ameliorates the memory complexity problem of self-attention because the $N \times N$ matrix is now decomposed to $N \times k$.

- **Kernels** - Another recently popular method to improve the efficiency of Transformers is to view the attention mechanism through kernelization. The usage of kernels (Katharopoulos et al., 2020; Choromanski et al., 2020a) enable clever mathematical re-writing of the self-attention mechanism to avoid explicitly computing the $N \times N$ matrix. Since kernels are a form of approximation of the attention matrix, they can be also viewed as a type of low-rank approach (Choromanski et al., 2020a). Examples of recent work in this area include Performers, Linear Transformers and Random Feature Attention (RFA, (Peng et al., 2021))

Long-range Arena (LRA) dataset

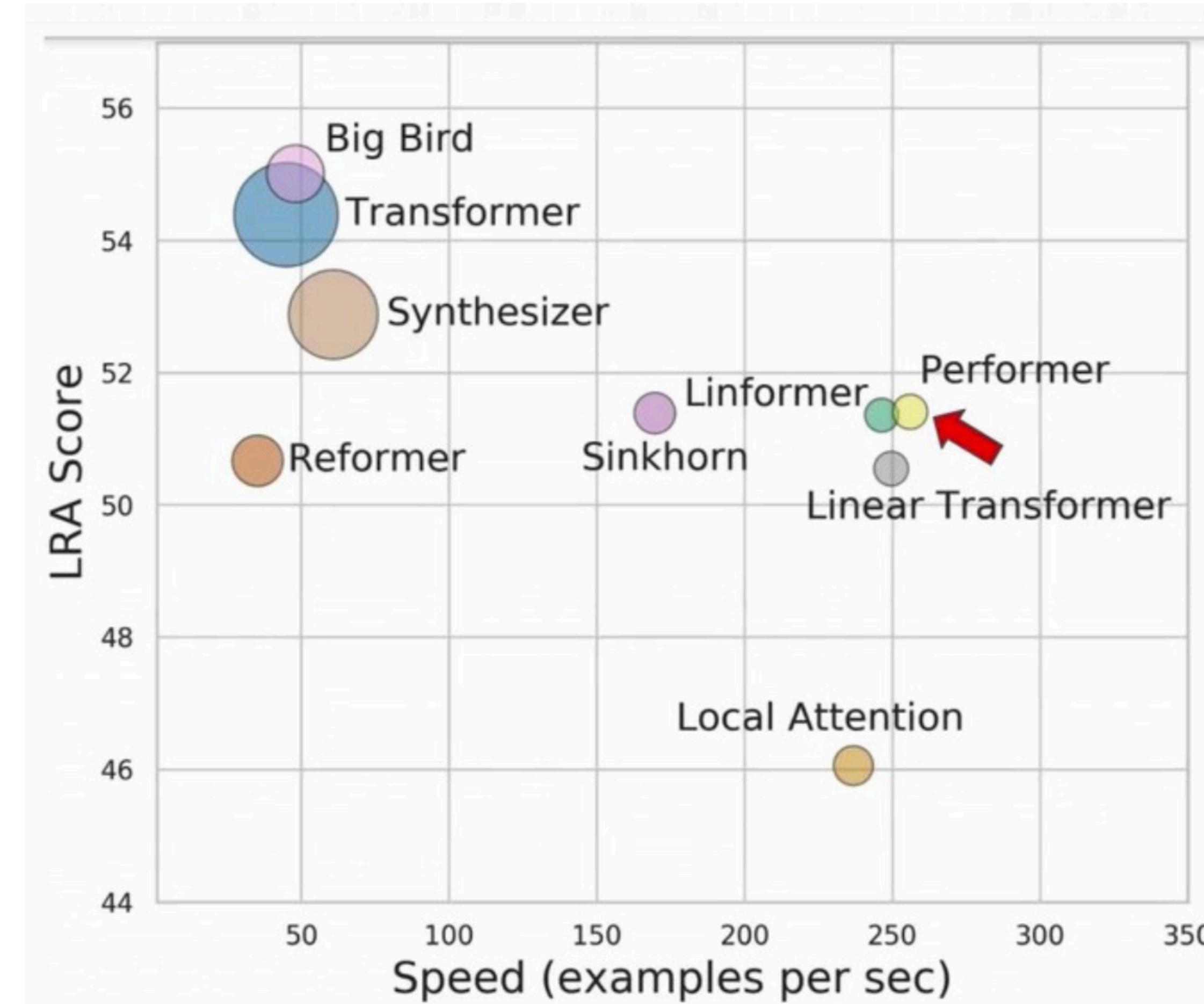
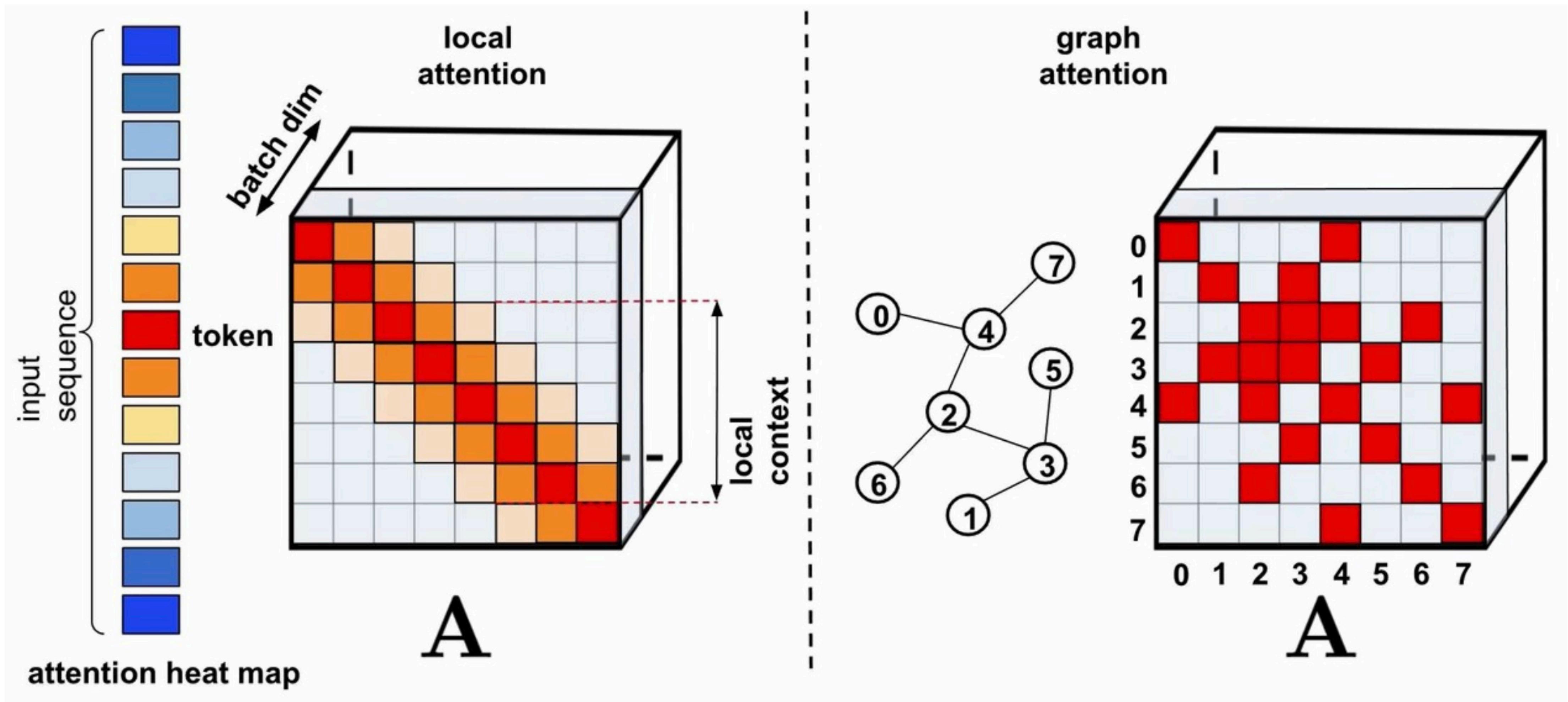
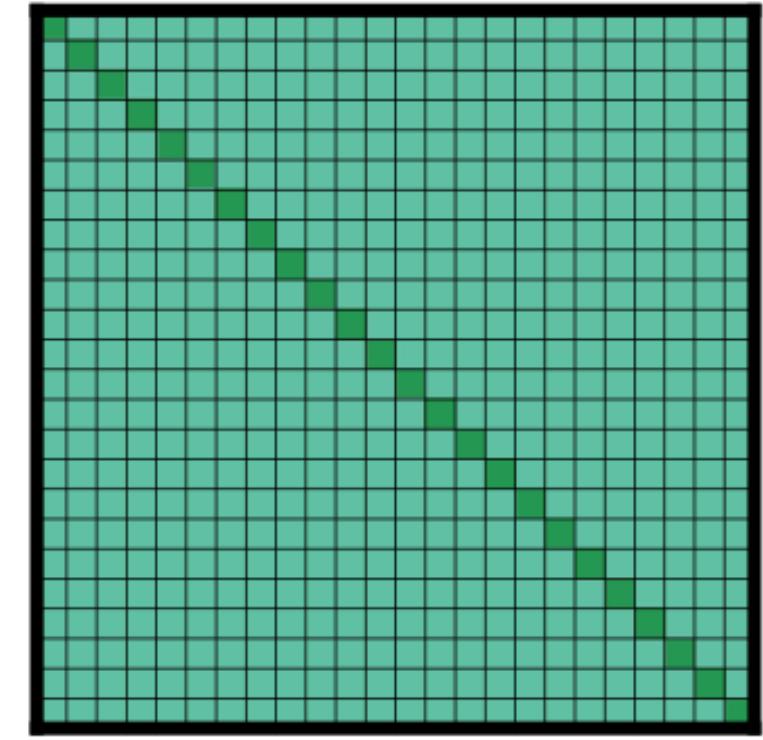


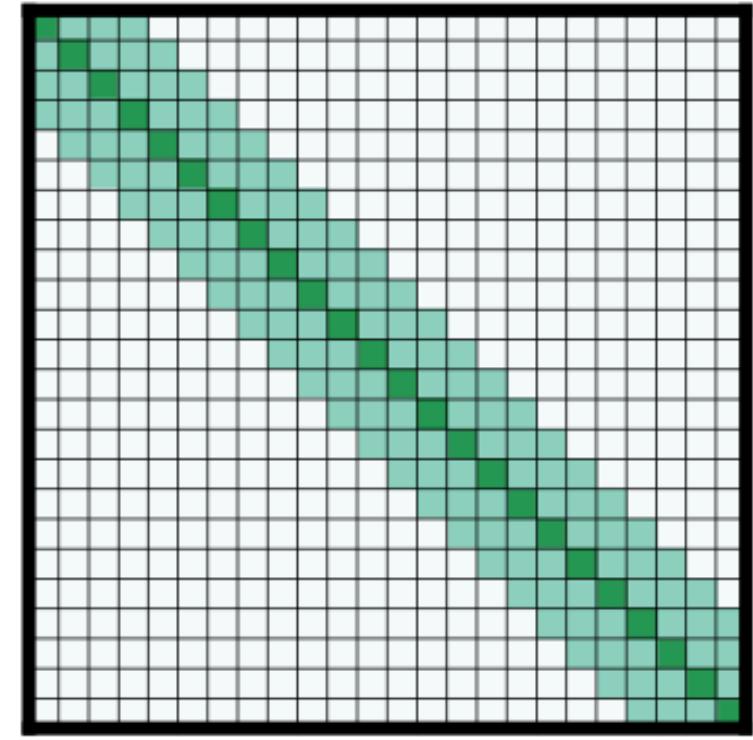
Figure 3: Performance (y axis), speed (x axis), and memory footprint (size of the circles) of different models.

Sparsifying Attention

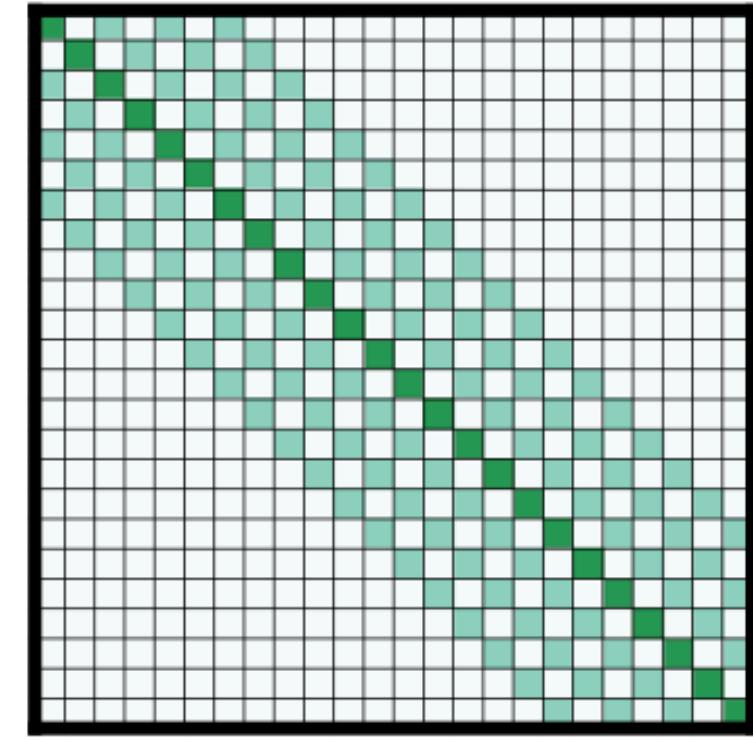




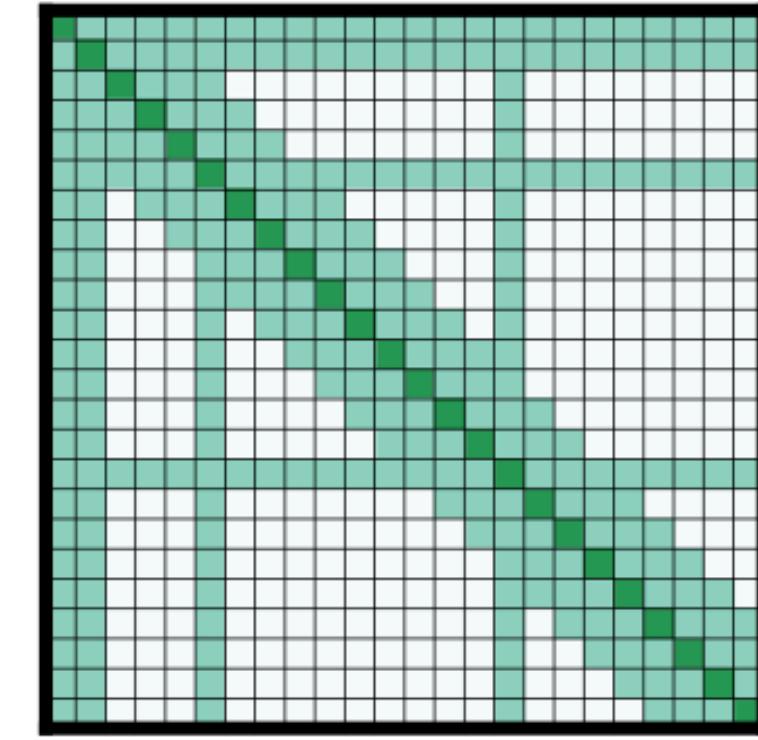
(a) Full n^2 attention



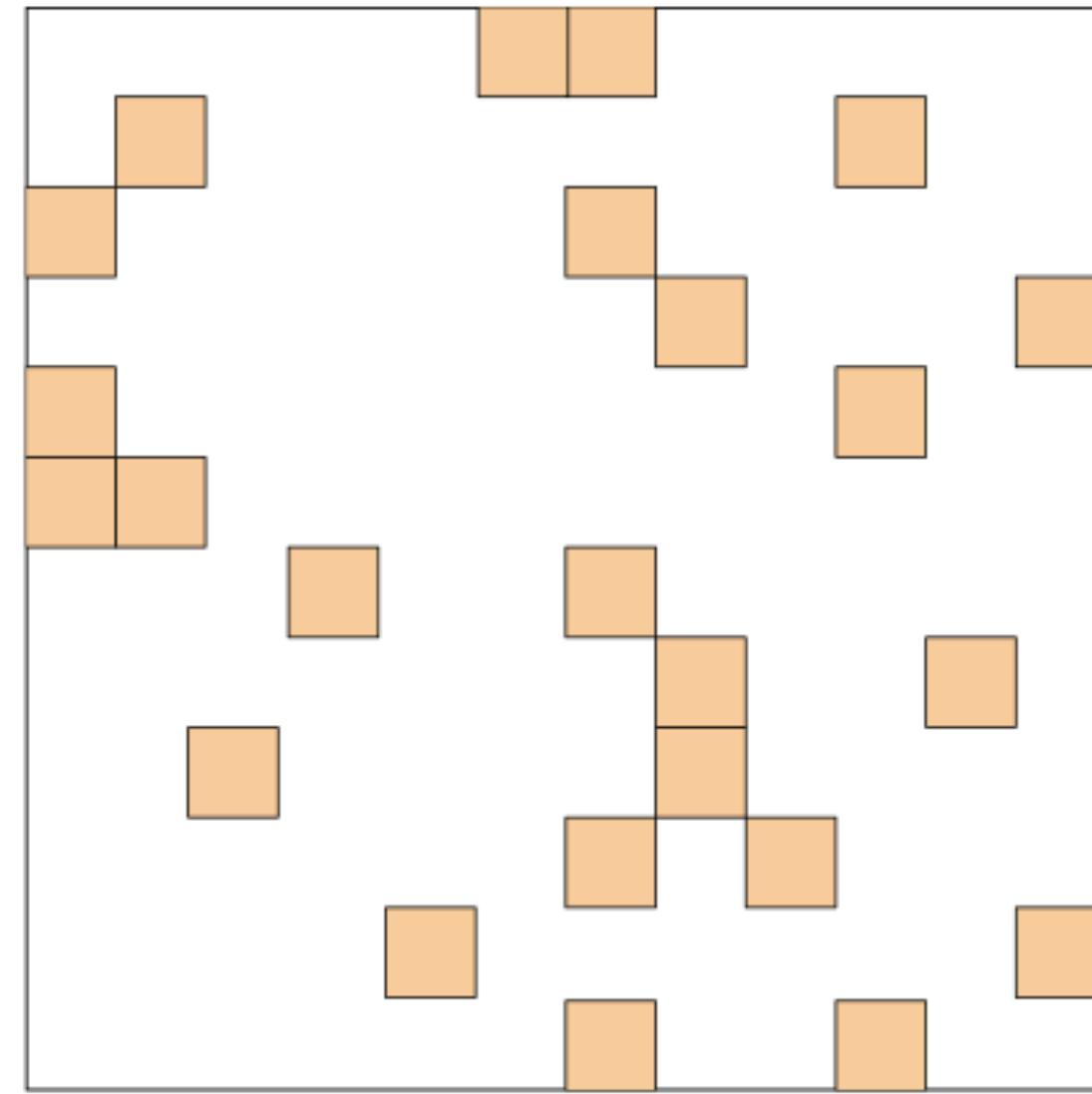
(b) Sliding window attention



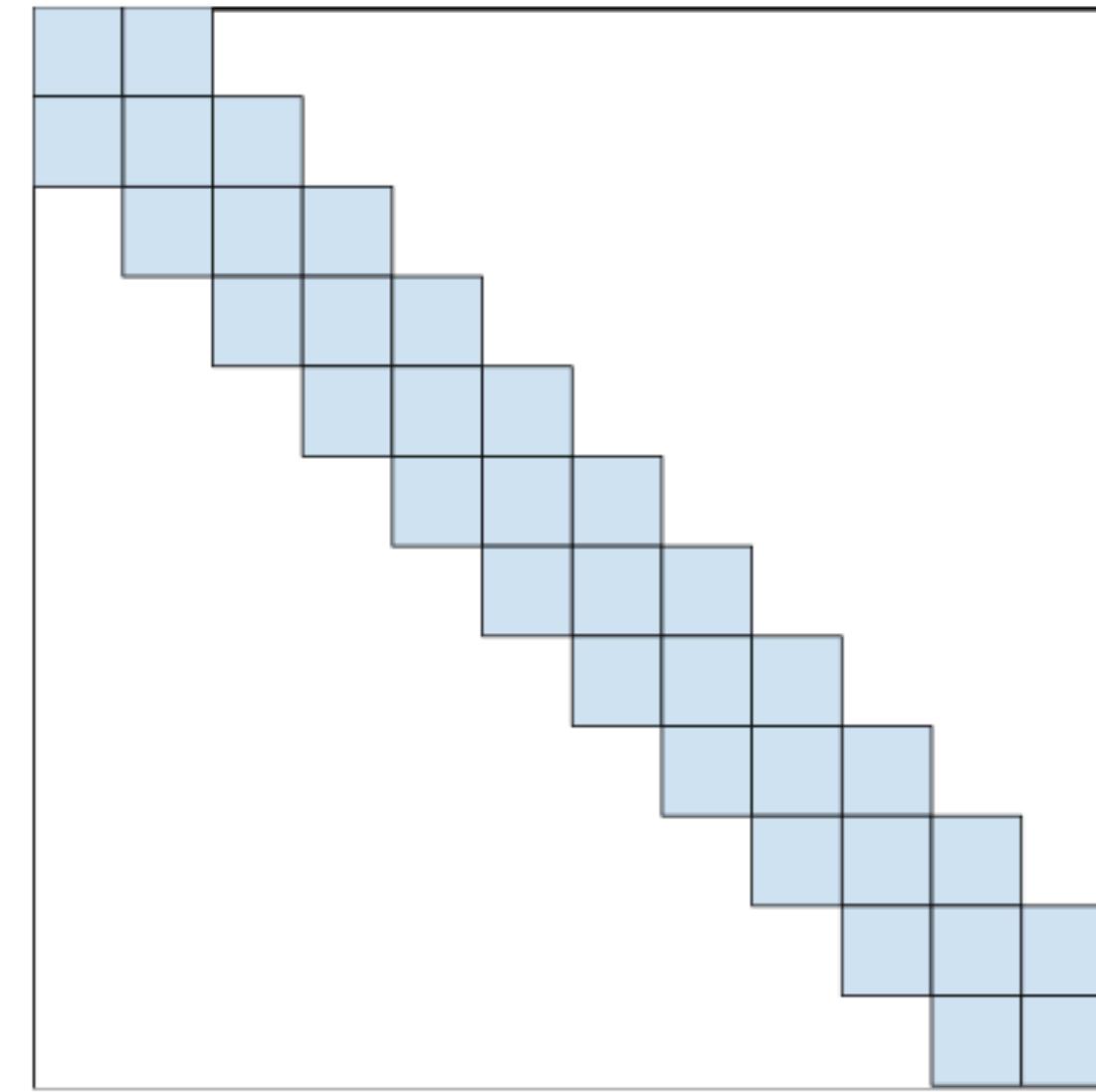
(c) Dilated sliding window



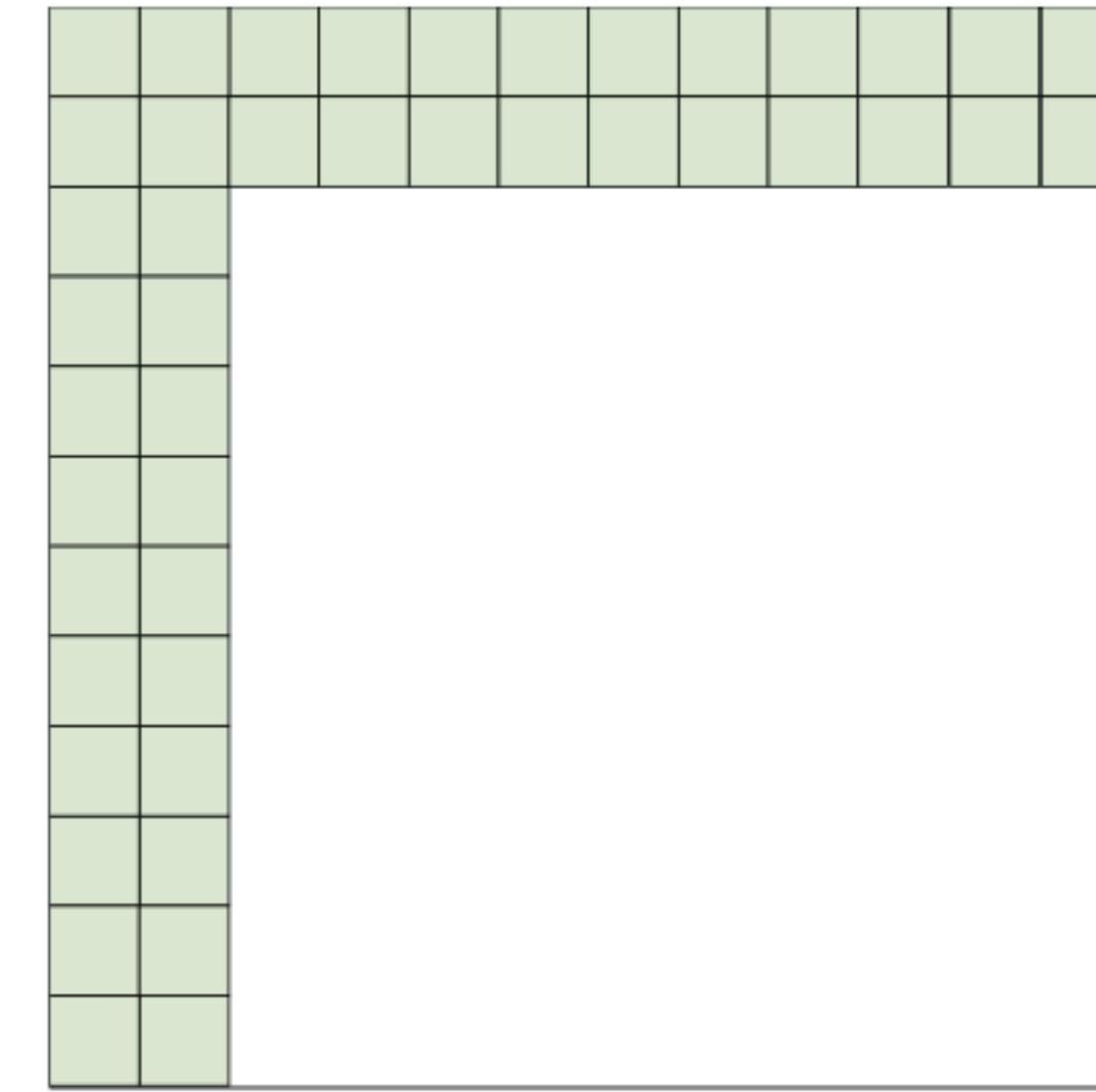
(d) Global+sliding window



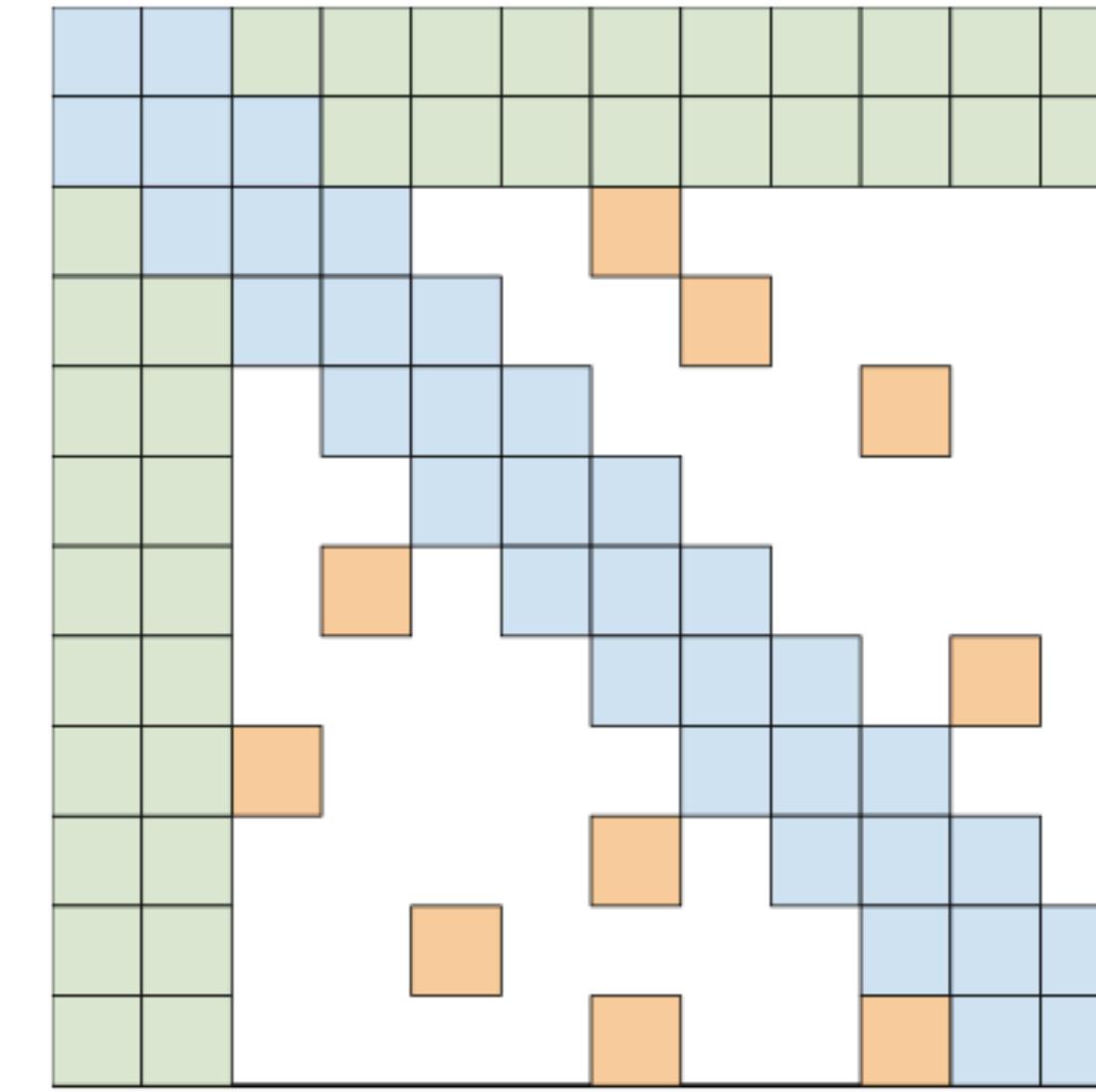
(a) Random attention



(b) Window attention



(c) Global Attention



(d) BIGBiRD