

Abstract

- 2 techniques for improving efficiency of transformers
 - a) Dot-product attention \rightarrow Locality Sensitive Hashing $(O(N^2)) \rightarrow (O(N \log N))$
 - b) Reversible Residual Layers - store activation once instead of storing every layer.
- performs on par with transformers while being memory & time efficient
- c) bringing power of transformers to other domains like time-series forecasting, music, image, video generation

I.) INTRODUCTION

"Generating Wikipedia articles" - 11k tokens of text in a single example input to a transformer

- models too large to be even fine-tuned on a single GPU.
- Are large transformers really efficient in their parameter usage?

$$0.5 \text{B params} = 2 \text{GB RAM}$$

$$64k \text{ input tokens} \times 1024 \text{ embedding size} \times 8 \text{ batch-size} = 0.5 \text{B floats} = 2 \text{GB RAM}$$

$$\Rightarrow 2 + 2 = 4 \text{ GB}$$

- If memory was per layer \Rightarrow fit a transformer easily
- whole corpus used to train BERT is only 17 GB. So, why we cannot even fine-tune on a single machine??

Reasons

- a) N layers = N time more memory
- b) d_f - depth of feed forward layers $\gg d_{\text{model}}$
- c) Attention Calculation is $O(L^2)$

Solution

- a) Reversible Layers \rightarrow only one copy of activations instead of N .
- b) Splitting activations inside f-f layers & processing in chunks
- c) Approximate attention computation using Locality Sensitive Hashing - $O(L^2) \rightarrow O(L \log L)$

is even numerically identical to original Transformer.

Tasks used:

↳ enwik8 (64k sequence)

↳ image generation - image net 64k-generation.

(12k sequence length).

2. LOCALITY-SENSITIVE HASHING ATTENTION - LSH Attention

$$\text{Att}(Q, K, V) = \text{Softmax} \left(\frac{QK^\top}{\sqrt{d_k}} \right) V$$

$$Q, K, V \rightarrow (\text{bsz}, \text{length}, d_{\text{model}}) \Rightarrow QK^\top = (\text{bsz}, \text{length}, \text{length}).$$

$64K \times 64K$ Matrix - 32 bit float

$$= 64,000 \times 64,000 \times 32 \\ = 16 \text{ GB}$$

Alternative : softmax $\left(\frac{q_i K^T}{\sqrt{d_{\text{model}}}} \right)$ \checkmark is linear in d_{model}

You have $A = (bsg, \text{length}, d_{\text{model}})$

Then 3 linear layers project from to

$$A \rightarrow L1 \rightarrow Q$$

$$A \rightarrow L2 \rightarrow K$$

$$A \rightarrow L3 \rightarrow V$$

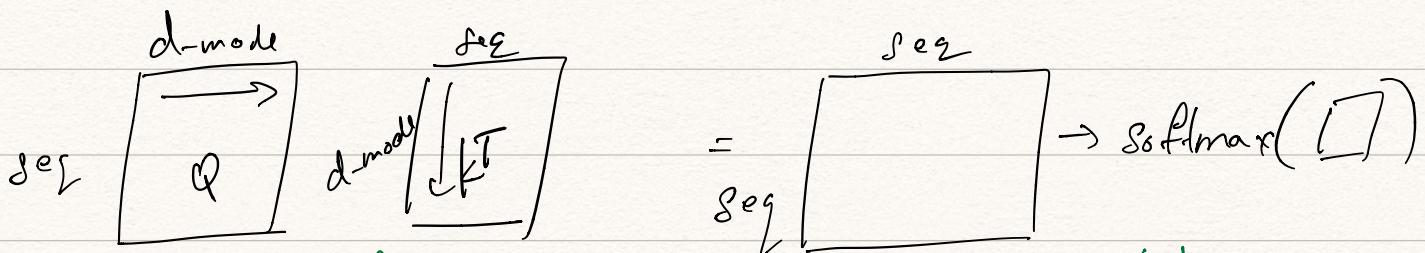
LSH $\Rightarrow L1 = L2$ - "Shared QK Trans for me"

Sharing Q & K does not affect the performance

Hashing Attention:

we need softmax (QK^T) \Rightarrow largest values & not all the values. ($\text{abs}(D) = \text{maximum value}$)

\Rightarrow We need values closer to q_i (may 32 of them).



① Even if you narrow it down, you still need all "seg x seg" of them right?

② Also, doing Nearest Neighbour is faster?

Locality Sensitive Hashing. (LSH)

(TUTORIAL TIME!!)

LSH \rightarrow hash similar items similarly. This is the opposite behavior of cryptographic hashing.
In LSH, we want to maximize the probability of collisions of similar items.

$P(h(a) == h(b)) = \text{low}$ if a near b in space.

LSH using Random Projection.

$$\begin{bmatrix} \text{Projected}(P) \\ k \times n \end{bmatrix} = \begin{bmatrix} \text{Random}(R) \\ k \times d \end{bmatrix} \begin{bmatrix} \text{Original}(D) \\ d \times n \end{bmatrix}$$

where $k < d$

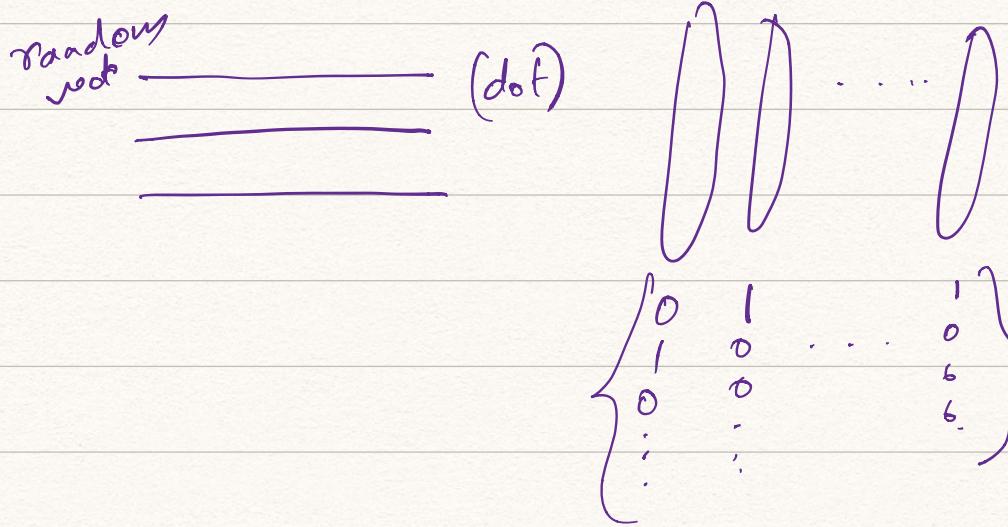
$k \left\{ \begin{array}{l} \text{---} \\ \text{---} \\ \text{---} \end{array} \right\}$ k -d-dimensional vectors sampled from $N(0,1)$

n-documents of d-dimension.

(.) D D D D

$P \rightarrow$ if $P > 0$, put binary 1

if $P \leq 0$, put binary 0



You will get
a vector of 0/1
for each document.

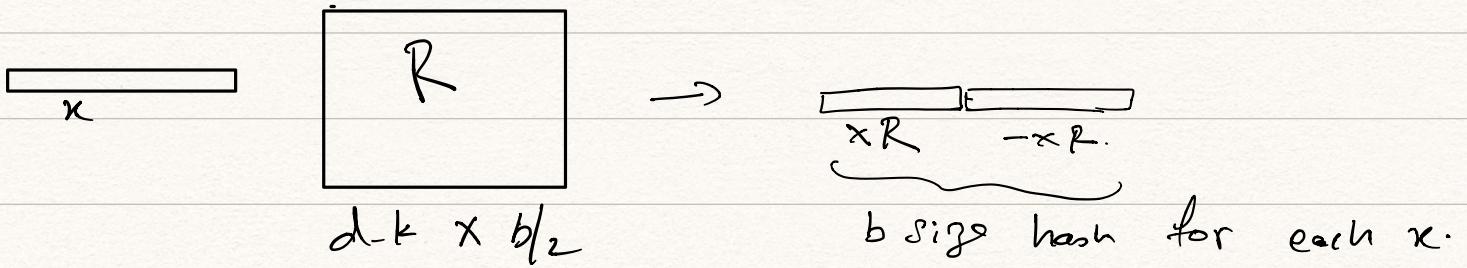
Hash Table \Rightarrow

If the binary representation is same for 2 documents then they are supposed to be similar.

- Due to the random nature we sometimes construct multiple hash tables (like multiple heads in a transformer).
- The documents that have similar binary representation are supposed to have higher cosine similarity.

B A C k !

Our requirement: Nearby vectors get the same hash with high probability.



$$\text{i.e., } h(x) = \arg\max([x^T; -x^T]). \quad \checkmark$$

gives axis of the maximum value.

Normal:

$$o_i = \text{softmax}\left(\sum_j q_i | k_j\right) \cdot V$$

$\lceil \log \rceil \left[\right]$
for d -model

LST:

$$o_i = \sum_{j \in P_i} \exp(q_i | k_j - \gamma(i, p_i)) \cdot V$$

$P_i \rightarrow$ set of keys matching the query

$\gamma \rightarrow$ Normalizing factor

$$P_i = \{j : h(q_i) == h(k_j)\}$$

	q_1	q_2	q_3	q_4	q_5	q_6
k_1
k_2		.			.	.
k_3				.		.
k_4				.		.
k_5				.		.
k_6	

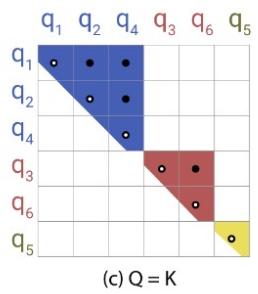
(a) Normal

\rightarrow Normal attn matrix which is sparse

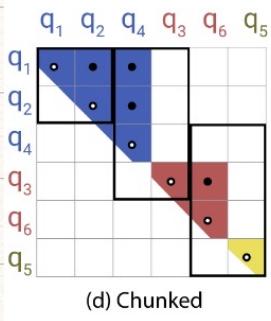
	q_1	q_2	q_4	q_3	q_6	q_5
k_1
k_2
k_3
k_4
k_5

(b) Bucketed

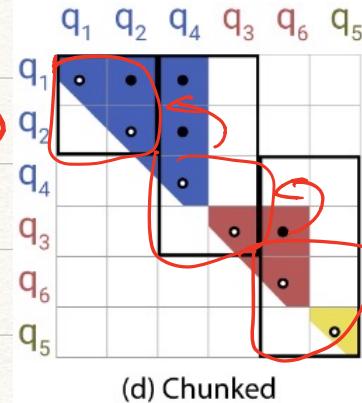
\rightarrow Find matching queries & keys. They could be uneven, also, number of query != number of keys. But we know nothing $Q = K$ word hamper performance



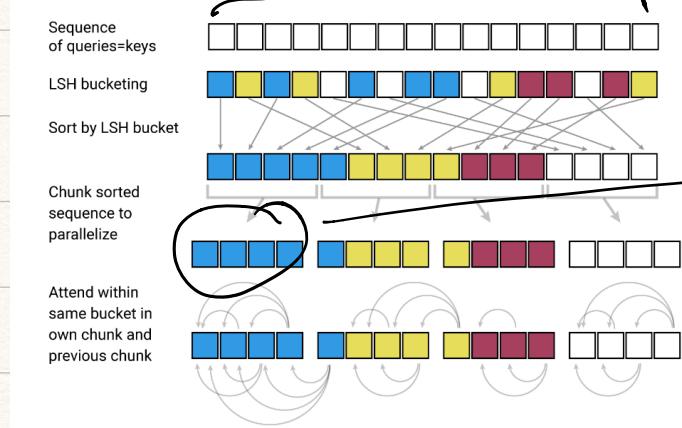
\rightarrow setting $\varphi = k$ ie, $k_j = \frac{e_j}{\|e_j\|}$



\rightarrow why not this ??



without chunking this is the attn matrix dim.



\rightarrow after chunking it is this -

\rightarrow cross chunk attention.

Train \ Eval	Full Attention	LSH-8	LSH-4	LSH-2	LSH-1
Full Attention	100%	94.8%	92.5%	76.9%	52.5%
LSH-4	0.8%	100%	99.9%	99.4%	91.9%
LSH-2	0.8%	100%	99.9%	98.1%	86.8%
LSH-1	0.8%	99.9%	99.6%	94.8%	77.9%

You can train
4 fast with
different
CSH models.

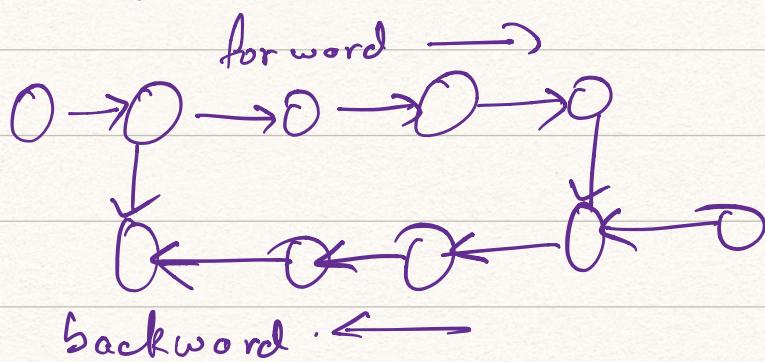
3) Reversible Transformer

Rev Nets:

Interesting paper: "The Reversible Residual Network: Back propagation without storing activations!"

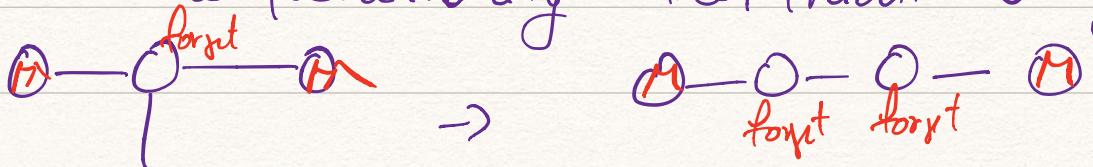
- instead of checkpointing for use in backward pass, you can reverse layer-by-layer during back prop - wow!

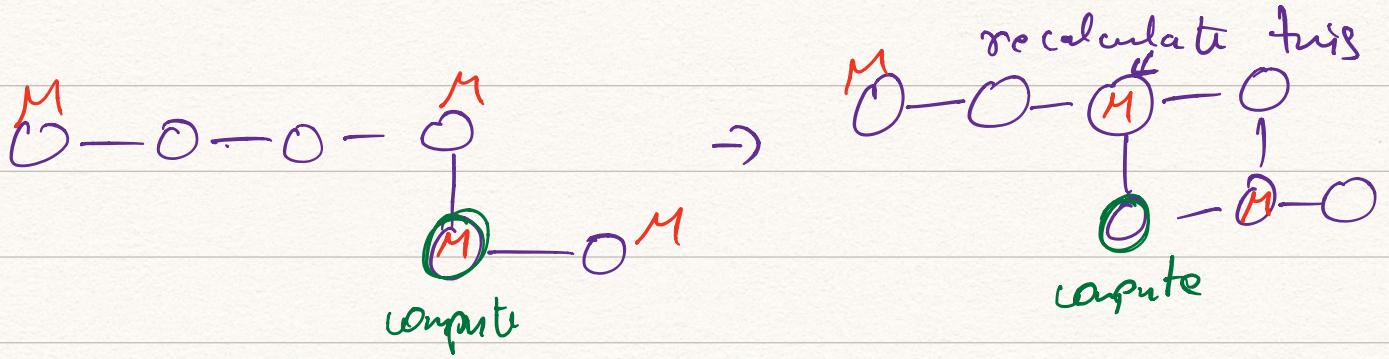
Checkpointing



You need to keep all the circles in memory to compute backward pass one layer at a time.

Instead there is a memory efficient - computationally inefficient way



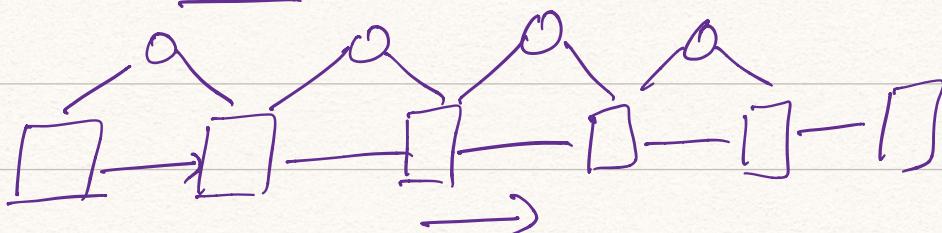


If, you recalculate the forward pass till you meet the required node & then compute backward pass on it:

You could pick some nodes & save their intermediate results.

e.g. For a chain of length 'n', place checkpoints at every $\sqrt[n]{n}$. This way you can be more memory & time efficient

Parents

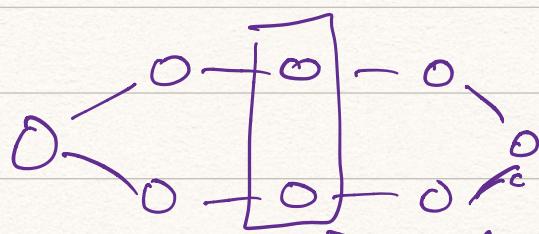


If makes no sense to checkpoint the circles.

Squares nodes are the best.

Knowing "on" would remove the need to recompute anything before - called - "graph separators"

e.g. example of graph separator of size 2

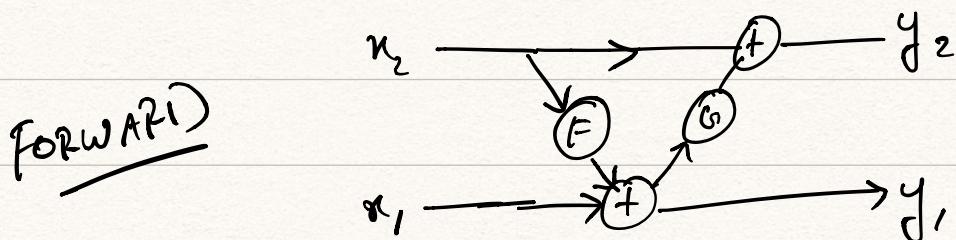


You need checkpoint pairs
like this

BACK!

Normal
Resnets : $x \rightarrow y$, $y = x + F(x)$

Rev Nets : Split your input(x) into (x_1) & (x_2)



$$y_1 = x_1 + F(x_2)$$

$$y_2 = x_2 + G(y_1)$$

BACKWARD we have y_2 & y_1 ,

$$x_2 = y_2 - G(y_1)$$

$$x_1 = y_1 - F(x_2)$$

This way we
can compute
the activation
of previous layer

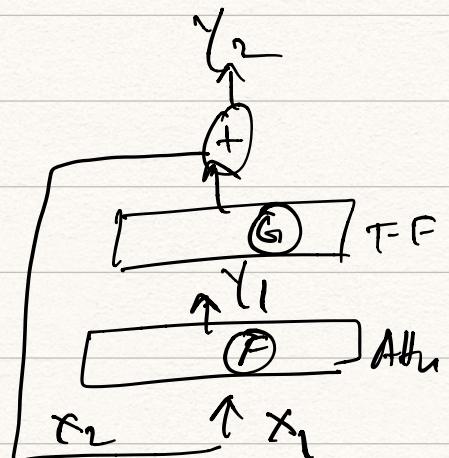
(x_1, x_2) from current
layer (y_1, y_2)

Reversible Transformer

$$y_1 = x_1 + \text{Attn}(x_2)$$

$$y_2 = x_2 + \text{FeedForward}(y_1)$$

$$x_1, x_2 \in \mathbb{R}^{d_{\text{model}}}$$



→ This can be chunked. & thereby can improve parallelism. Also operating on one chunk at a time can reduce memory.

* Now, due to chunking & reversible layers, the memory we use for the forward pass is independent of the number of layers

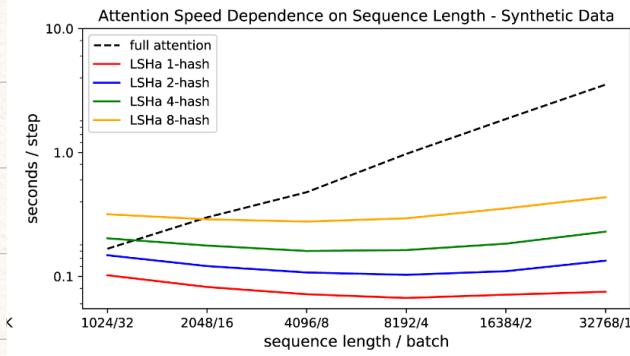
RELATED WORK

There is a lot of very interesting ideas. I have no idea if I'll be able to go through them.

5 EXPERIMENTS

- ① Reversible Layers & sharded QK does not hamper the performance.
- ② Sharded QK, $k_i = \frac{q_i}{\|q_i\|}$. Should prevent tokens attending themselves. why??
- ③ 20 layer Reformer models fit in the memory
- ④ LSH-2 almost matches full attention.

(5)



The speed of attention computation remains mostly constant for LSH when seq-length is increased.