Stanford, ML Sys

Video 67 — flash Attention (F.A)

fast and Memory Efficient Exact Attention with IO Awareness

Why this ?

3 main reasons for Modelling Longer Sequences

① NLP :- large context is required to understand books, plays, .... etc

② CV :- higher resolution → better insights

③ Time series, Audio, Video, medical Images → data understanding

↓

there are Sequences of Million steps

Challenge :-

Scaling Transformers to longer Sequences :-

* Content length of GPT 3 → 2048 (Seq)

F.A → helps to train things faster and with longer content.

① Tiling

② Recomputation

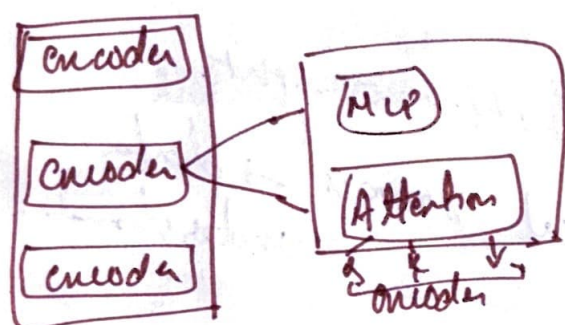] — 2 main things in F.A to reduce GPU memory IOs.

Speed → 3 times faster

Memory Efficient → 10-20 x for Exact attention

Seq length → up to 16k (Content)

# Table of Content

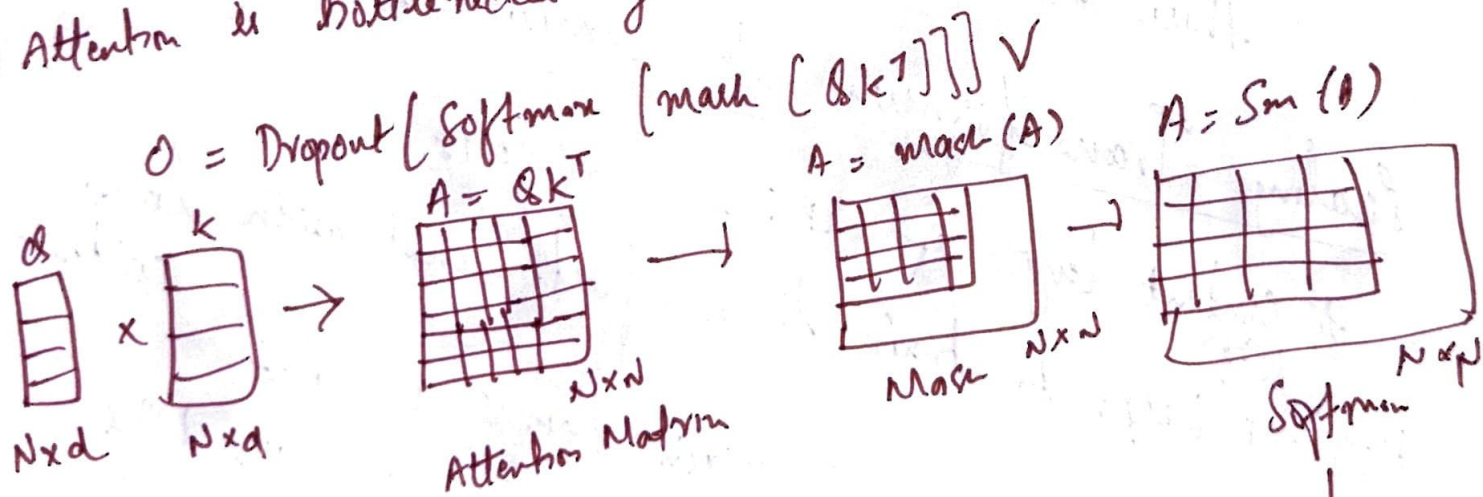## 1. Background



Transformer
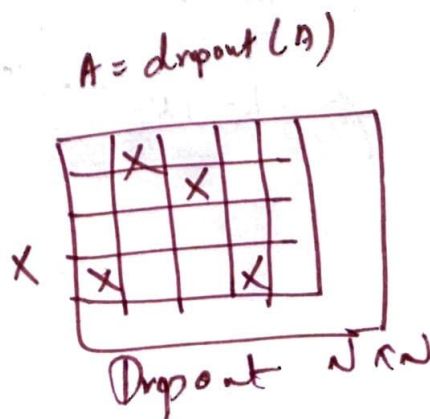
$Q$ = Query
$k$ = Key
$V$ = Value
$N$ = Seq length $(2k, 4k, 16k)$
$d$ = head dimension $(64, 128 \dots)$

→ Attention is heart of Transformer
→ Attention is bottlenecked by Memory Read/writes

$$O = Dropout \left[ softmax \left[ mask \left[ Qk^T \right] \right] \right] V$$



$A = Qk^T$

$N \times d$    $N \times d$

Attention Matrix $N \times N$

$A = mask(A)$

Mask $N \times N$

$A = Sm(A)$

Softmax $N \times N$

$A = dropout(A)$
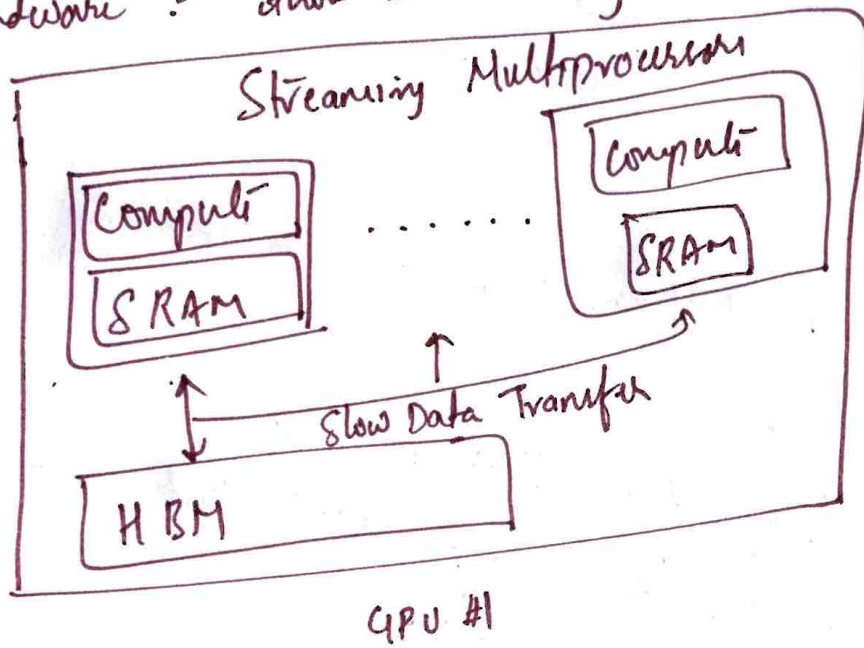
$O = AV$

$N \times d$    $N \times d$    Dropout $N \times N$

→ The problem is, lot of time is spent on Reading and writing in GPU instead of calculations.

→ Attention in GPT 2 : Masking, softmax and Dropout takes lot of time than matrix multiplication.

→ hardware :- there is memory hierarchy in hardware

→ SRAM is much much smaller but much much faster



Streaming Multiprocessors

| Compute |
| SRAM |

. . . . . . .

| Compute |
| SRAM |

↕ ↑ ↗ Slow Data Transfer

| H BM |

GPU #1

While any computation :-
inputs → HBM → SRAM ⇄ Compute → HBM
                              ↓
                        (Softmax, Matmul)
                           ...es

→ data moves back and forth between HBM and from SRAM

② Method . F.A

Challenges :- How to reduce HBM Read/write : compute by blocks :-

① compute softmax reduction without access to full input

② backward without the large attention matrix from forward.

Techniques used to address these issues are:-

① Tiling : Restructure algorithm to load
block by block from HBM to SRAM
to compute attention

② Recomputation : Dont store attention matrix from
forward, recompute it in the
backward

Tiling :-
decomposing large softmax into smaller ones by scaling

Steps :-
① load inputs by block from HBM to SRAM
② On chip, compute attention outputs wrt that block
③ Update output in HBM by Scaling

Recompute :-
by storing softmax normalization factors from fwd
(size N) quickly recompute attention is hand from inputs
in SRAM