# Memory and IO in deep learning

## Lecture 12 — CS 577 Deep Learning

Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology

November 6, 2024

Lec 13 is a seminar

# Memory and IO in deep learning

## Lecture 12 — CS 577 Deep Learning

So far...

- Backprop
- 2D convol
  lower to Linear Alg
- Parallelize seq model
  RNN → Transformer

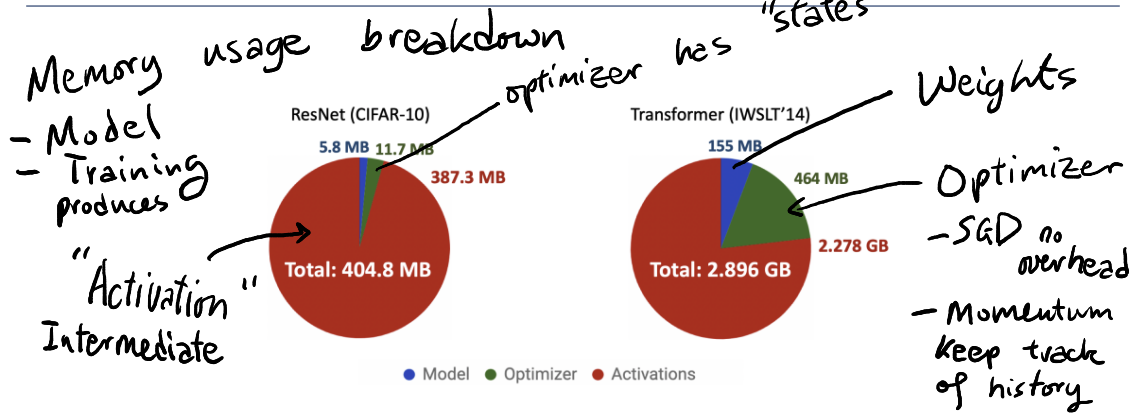Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology
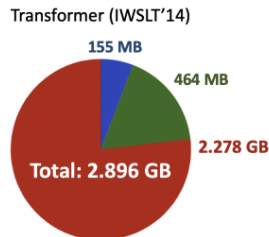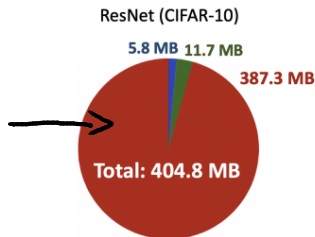
November 6, 2024

No discussion of memory so far

# The memory wall



*Handwritten annotations:*

Memory usage breakdown
- Model
- Training produces "Activation" Intermediate

optimizer has "states" - stateful

Weights

Optimizer
- SGD no overhead
- Momentum keep track of history

ResNet (CIFAR-10)

5.8 MB   11.7 MB
387.3 MB
Total: 404.8 MB

Transformer (IWSLT'14)

155 MB
464 MB
2.278 GB
Total: 2.896 GB

● Model  ● Optimizer  ● Activations

From Sohoni et al. 2019 "Low-memory neural network training: A technical report"
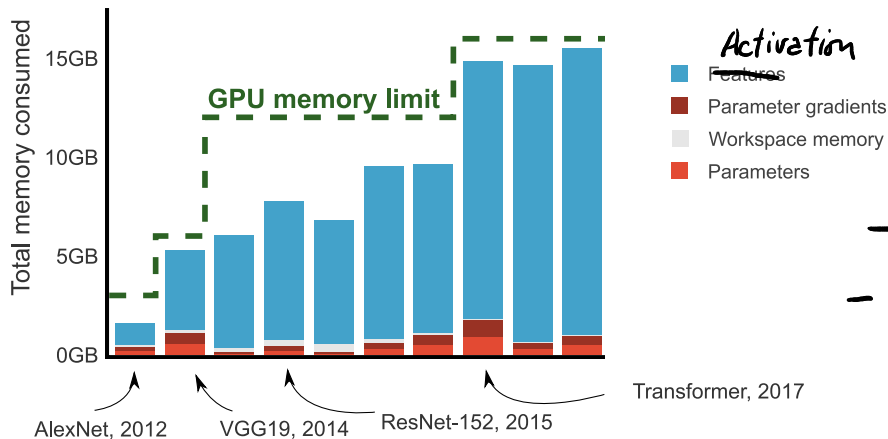
# The memory wall



Biggest piece of pie

ResNet (CIFAR-10)

5.8 MB  11.7 MB
387.3 MB

Total: 404.8 MB

Transformer (IWSLT'14)

155 MB
464 MB
2.278 GB

Total: 2.896 GB

● Model ● Optimizer ● Activations

From Sohoni et al. 2019 "Low-memory neural network training: A technical report"
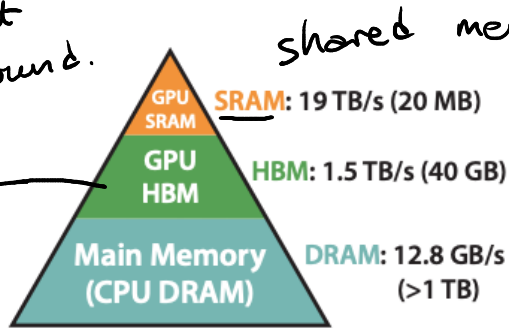
# The memory wall



From Jain et al. 2020 "Checkmate: Breaking the memory wall with optimal tensor rematerialization"

# Flash attention

IO input-output
moving data around.
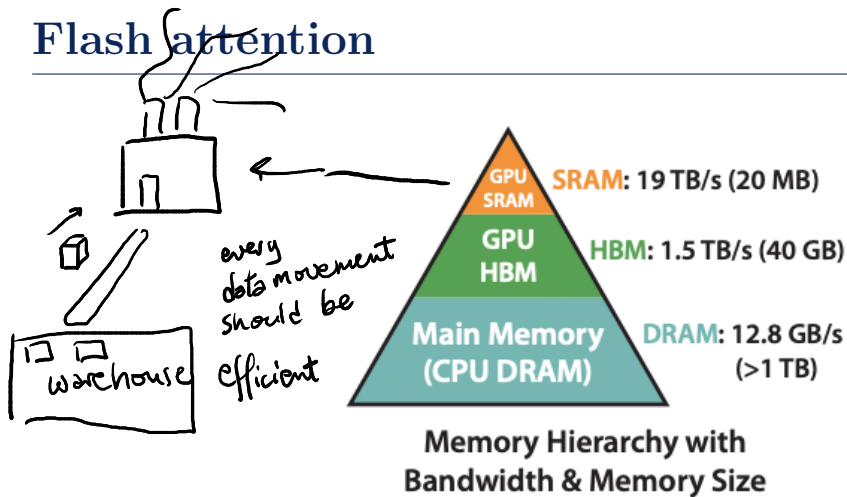
"off-chip"
more space
but slow

shared memory on chip
physical close to compute



**Memory Hierarchy with Bandwidth & Memory Size**

From Dao et al. 2022 "FlashAttention: Fast and memory-efficient exact attention with io-awareness"

# Flash attention



**Memory Hierarchy with Bandwidth & Memory Size**

From Dao et al. 2022 "FlashAttention: Fast and memory-efficient exact attention with io-awareness"

# Memory and IO

- DL workloads often memory limited
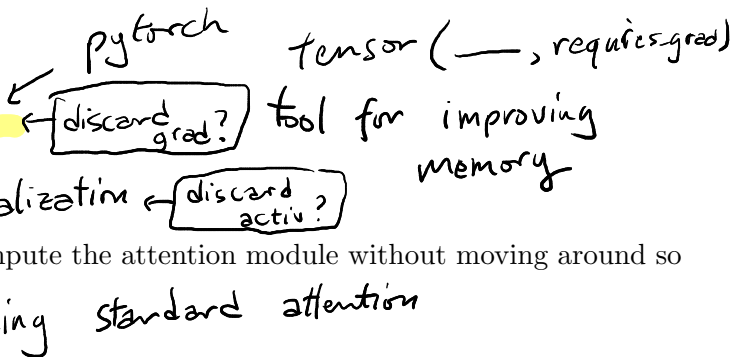- Moving data (IO) is costly

# Roadmap for today

*pytorch* &larr;

*tensor ( ___ , requires_grad)*

Memory:

- Warm-up: `requires_grad` &larr; *discard grad?* *tool for improving memory*
- Gradient checkpointing
  - *rematerialization* &larr; *discard activ?*

IO:

- FlashAttention: how to compute the attention module without moving around so much data &larr; *Computing standard attention*

**Next**: `requires_grad`

# Size of NumPy array

```
1 import numpy as np
2 import sys
3
4 n = 1000
5 X_np = np.random.randn(n,n)
6 sys.getsizeof(X_np)
7
8 # >>> 8000128      # <- unit in bytes
```

# Array in C

```
1 double *X;
2
3 X = (double *)malloc(n * n * sizeof(double));
```

```
1 n*n*8
2
3 # >>> 8000000
```

≈ 8 Mb

RAM

8 bytes
in 1 float 64
double

1000 × 1000 × 8

# Size of NumPy array

```python
import numpy as np
import sys

n = 1000
X_np = np.random.randn(n,n)
sys.getsizeof(X_np)

# >>> 8000128    # <--- extra bytes for metadata
```

```python
n*n*8      # <--- size_of_float64  times  num_of_elements_in_X_np

# >>> 8000000
```

# Aside: Size of PyTorch array

```
1 import torch
2
3 n = 1000
4 X_torch = torch.randn(n,n)          ← reference
5 sys.getsizeof(X_torch.untyped_storage())
6
7 # >>> 4000072
```

```
1 n*n*4       # <--- size_of_float32   times   num_of_elements_in_X_torch
2
3 # >>> 4000000
```

bfloat16 — google brain

— Quantization: save memory by making things less precise

# tracemalloc

Allows tracking only memory allocated by numpy

```
1  import tracemalloc          track memory allocated by a module
2                                                         numpy
3  # [...]
4  print("Simple example")
5
6  start_trace()
7
8  a = 1.0*np.zeros((n,n))      # n = 1000
9
10 print_trace_stats(snapshot_trace())
11 end_trace()
12 # >>> Simple example
13 # >>> memory allocated: 8 MB
```

$$\begin{bmatrix} 0 & \cdots & 0 \\ \vdots & & \vdots \\ 0 & \cdots & 0 \end{bmatrix}$$

$1000 \times 1000 \times 8$
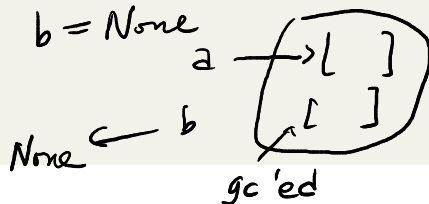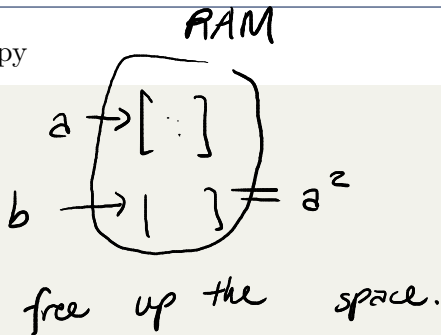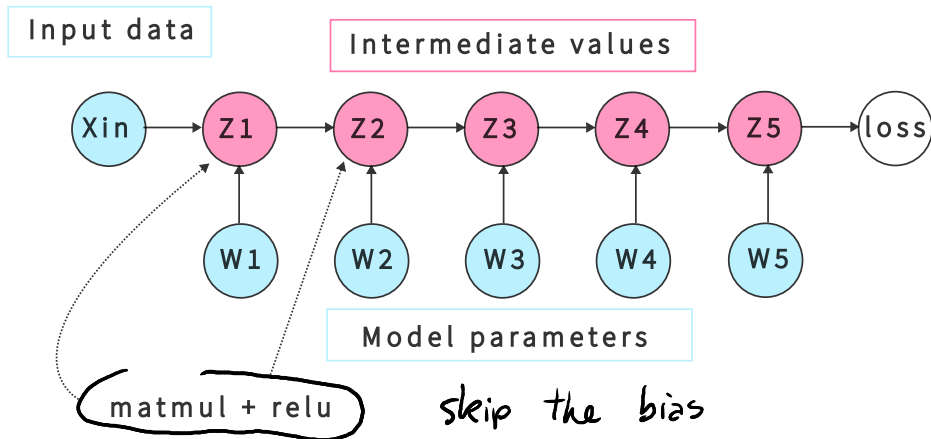
# tracemalloc

Allows tracking only memory allocated by numpy

```
1  import tracemalloc
2
3  # [...]
4  print("No discarding")
5  start_trace()
6
7  a = 1.0*np.zeros((n,n))    # n = 1000
8  b = a*a
9  c = np.sum(b,axis=1)
10
11 print_trace_stats(snapshot_trace())
12 end_trace()
13 # >>> No discarding
14 # >>> memory allocated: 16 MB
```

*(handwritten annotations):* 8 MB, 8 MB, 1000×8, 8 Kb, 16.008

# tracemalloc

Allows tracking only memory allocated by numpy

```python
import tracemalloc

# [...]
print("Discarding")
start_trace()

a = 1.0*np.zeros((n,n))      # n = 1000
b = a*a
c = np.sum(b,axis=1)
b = None

print_trace_stats(snapshot_trace())
end_trace()
# >>> With discarding
# >>> memory allocated: 8 MB
```

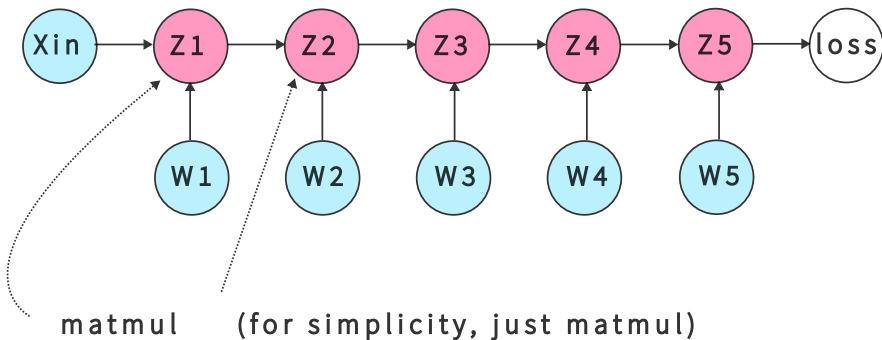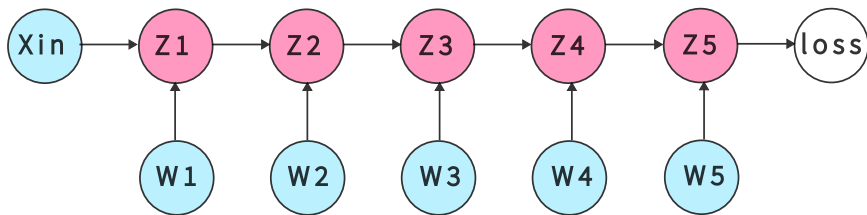# What is needed to calculate the backward?

# What is needed to calculate the backward?

Not a good model, but simple enough
for memory analysis



matmul    (for simplicity, just matmul)

# What is needed to calculate the backward?

# What is needed to calculate the backward?

$$Z1.shape = [4000, 1000]$$
batch_size

Z1 = ag.matmul(Xin, W1)

32 Mb



shape

(N, d)

(4000, 1000)

(1000, 1000)

W5.value has 8 MB

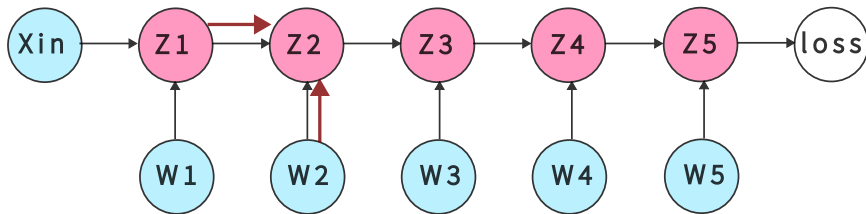# What is needed to calculate the backward?

Z2 = ag.matmul(Z1, W2)



Z2. value  — 32 MB
.grad  — 32 MB

total 64 MB
per
layer

# What is needed to calculate the backward?

(Inside constructor for Z2)...

```
self.value = 1.0*value
self.grad = np.zeros_like(self.value)
```
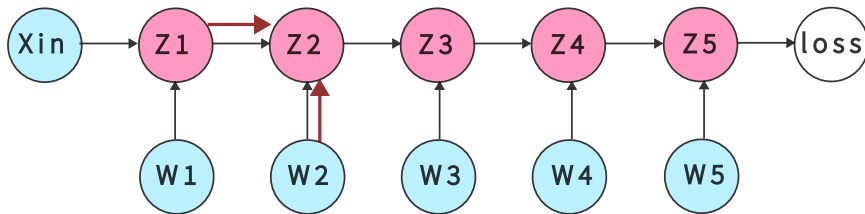
← wasted

# Array sizes!

(Inside constructor for Z2)...

self.value = 1.0*value
self.grad = np.zeros_like(self.value)

64 MB

# What is needed to calculate the backward?

```
1 def forward(x, weights):
2     for w in weights:
3         x = ag.matmul(x, w)
4     return ag.sum(x)
```

# What is needed to calculate the backward?

```
1  def forward_traced(x, weights):
2      start_trace()                          # tracing
3      mem_usage = []                         # tracing
4      for w in weights:
5          x = ag.matmul(x, w)
6          mem_usage.append(snapshot_trace()) # tracing
7      l = ag.sum(x)
8      mem_usage.append(snapshot_trace()) # tracing
9      end_trace() # tracing
10     return l, mem_usage
```
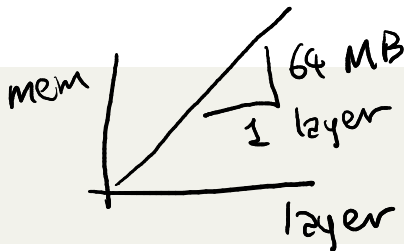
64 MB

128 MB

# What is needed to calculate the backward?

```python
for i, trace_stats in enumerate(mem_usage_forward):
    print(f"layer {i}")
    print_trace_stats(trace_stats)
```
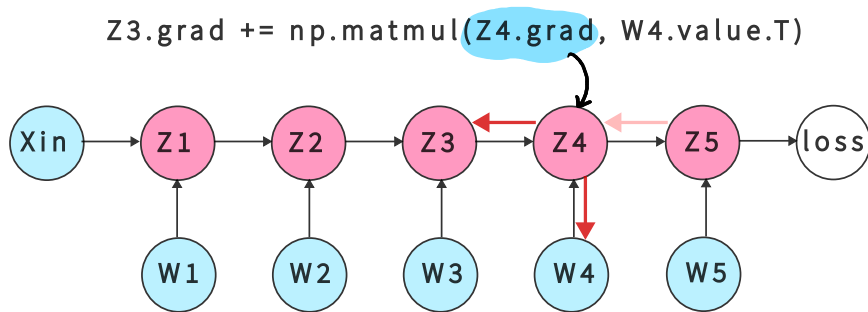
Output:

```
layer 0
memory allocated: 64 MB
layer 1
memory allocated: 128 MB
layer 2
memory allocated: 192 MB
```

**Observation**: double the amount of memory actually needed!

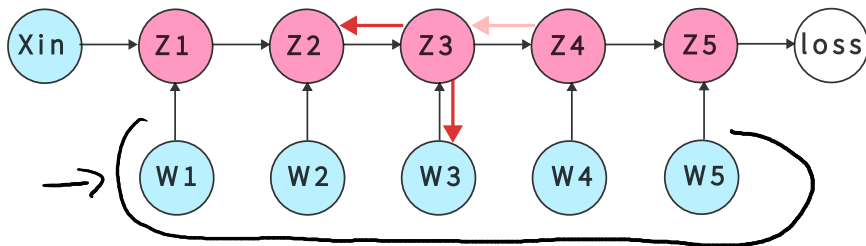# What is needed to calculate the backward?



Z3.grad += np.matmul(Z4.grad, W4.value.T)

W4.grad += np.matmul(Z3.value.T, Z4.grad)

# What is needed to calculate the backward?

`Z2.grad += np.matmul(Z3.grad, W3.value.T)`



only
care
about
grads

`W3.grad += np.matmul(Z2.value.T, Z3.grad)`

`Z4.grad never used again`

# Backward

```python
def backward(self):
    self.grad = np.array(1.0)

    topo_order = self.topological_sort()

    for node in reversed(topo_order):
        node._backward()
    return None
```

# Backward with tracing

```
1    def backward(self):
2        self.grad = np.array(1.0)
3
4        topo_order = self.topological_sort()
5
6        start_trace()  # tracing
7        mem_usage = [] # tracing
8
9        for node in reversed(topo_order):
10           node._backward()
11           mem_usage.append(snapshot_trace()) # tracing
12       end_trace() # tracing
13       return mem_usage
```

# Backward with tracing

```
1  for i, trace_stats in enumerate(mem_usage_backward):
2      print(f"backward step {i}")
3      print_trace_stats(trace_stats)
```

Output:

```
1  backward step 0
2  memory allocated: 0 MB
3  backward step 1
4  memory allocated: 0 MB
5  backward step 2
6  memory allocated: 0 MB
```

allocated
everything
already

**Observation**: After calling backward on a tensor, its grad is no longer needed

# Backward with tracing

```
1 for i, trace_stats in enumerate(mem_usage_backward):
2     print(f"backward step {i}")
3     print_trace_stats(trace_stats)
```

Output:

```
1 backward step 0
2 memory allocated: 0 MB
3 backward step 1
4 memory allocated: 0 MB
5 backward step 2
6 memory allocated: 0 MB
```

**Observation**: After calling backward on a tensor, its grad is no longer needed

**Next**: Shift the task of initialize the grad to the backward pass

# Keeping track of grad — is it needed?

*Constructor*

```
1    class Tensor: # Tensor with grads
2        def __init__(self,
3                     value,
4            # [...]
5
6            self.value = 1.0*value
7            self.grad = np.zeros_like(self.value) # <---- WASTEFUL!
```

# Keeping track of grad — is it needed?

```
1    class Tensor: # Tensor with grads
2        def __init__(self,
3                     value,
4                     requires_grad=False, # <-- Flag for keeping grad
5            # [...]
6
7            self.value = 1.0*value
8
9            self.grad = None
10           if self.requires_grad:
11               self.grad = np.zeros_like(self.value)
```

# Keeping track of grad — is it needed?

```
1    class Tensor: # Tensor with grads
2        def __init__(self,
3                     value,
4                     requires_grad=False, # <-- Flag for keeping grad
5            # [...]
6
7            self.value = 1.0*value
8
9            self.grad = None
10           if self.requires_grad:
11               self.grad = np.zeros_like(self.value)
12   # [...]
13
14   weights = [ag.Tensor(0.02*np.random.randn(dim_hidden, dim_hidden),
15                        requires_grad = True) for _ in range(num_layers)]
16   #
```

# Keeping track of `grad` — is it needed?

```
1    class Tensor: # Tensor with grads
2        def __init__(self,
3                     value,
4                     requires_grad=False, # <-- Flag for keeping grad
5            # [...]
6
7            self.value = 1.0*value
8
9            self.grad = None
10           if self.requires_grad:
11               self.grad = np.zeros_like(self.value)
12 # [...]
13
14   for w in weights:
15       x = ag.matmul(x, w)
16       # x.requires_grad == False by default
```

**Problem**: during backward, we might encounter `None`'s when we should expect `grad`.
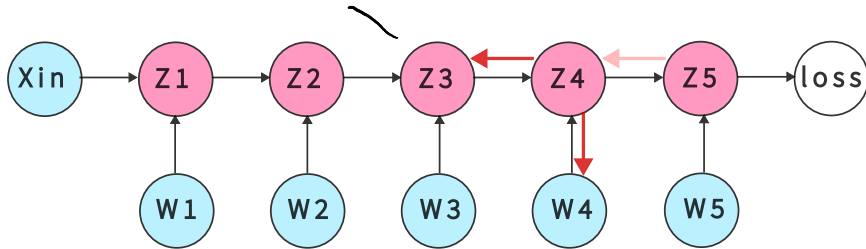
# What is needed to calculate the backward?

is None

↓

assume somehow
we initialized
this

↓

```
Z3.grad += np.matmul(Z4.grad, W4.value.T)
```
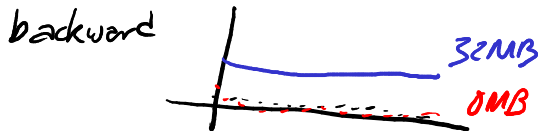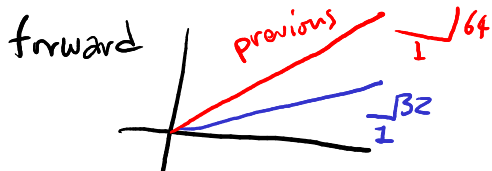


```
W4.grad += np.matmul(Z3.value.T, Z4.grad)
```

↳ this is fine

# Exercise 1: discarding the grad

```python
# inside ag.__matmul__
        def _backward():
            # YOUR CODE HERE FOR initializing the grad

            self.grad += np.matmul(output.grad, other.value.T)
            other.grad += np.matmul(self.value.T, output.grad)

            # YOUR CODE HERE FOR discarding the grad
            return None
```

*and   def sum(·):*
*        _backward*

If your answer is correct, then the "sanity checks" should have the expected output.



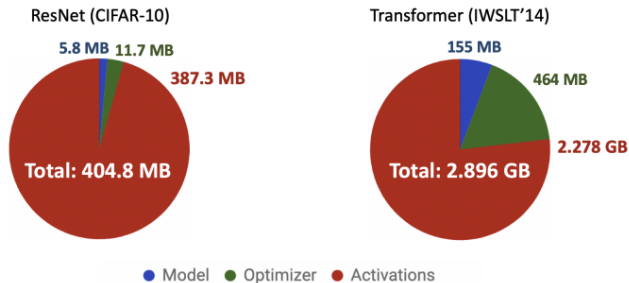forward — previous — 64/1, 32/1

backward — 32MB, 0MB

# The memory wall

- We've reduced the memory consumption of forward by a factor 2 without incurring much additional backward memory burden.

# The memory wall

- We've reduced the memory consumption of forward by a factor 2 without incurring much additional backward memory burden.
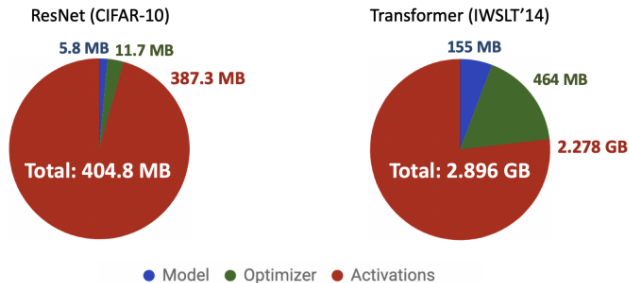- But recall that...



From Sohoni et al. 2019 "Low-memory neural network training: A technical report"

# The memory wall

- We've reduced the memory consumption of forward by a factor 2 without incurring much additional backward memory burden.

# The memory wall

- We've reduced the memory consumption of forward by a factor 2 without incurring much additional backward memory burden.
- But recall that...



From Sohoni et al. 2019 "Low-memory neural network training: A technical report"
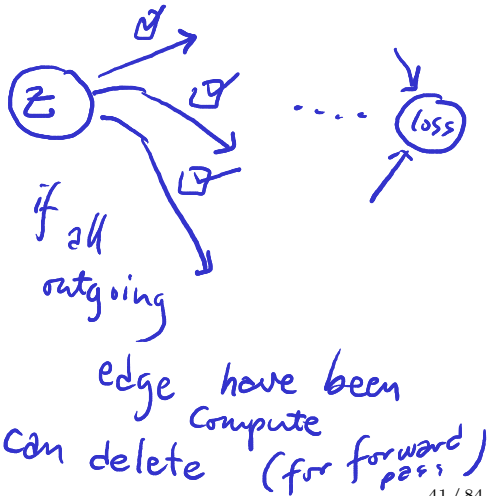
# The memory wall

Can we do something about this?

```
 1 layer 0
 2 memory allocated: 32 MB
 3 layer 1
 4 memory allocated: 65 MB
 5 layer 2
 6 memory allocated: 98 MB
 7 layer 3
 8 memory allocated: 131 MB
 9 layer 4
10 memory allocated: 163 MB
11 ...
```

# Tensor Rematerialization

- Discard the activations as soon as you can



$z$

$\nabla$
$\nabla$
$\nabla$

$\cdots$ loss

if all
outgoing
edge have been
compute
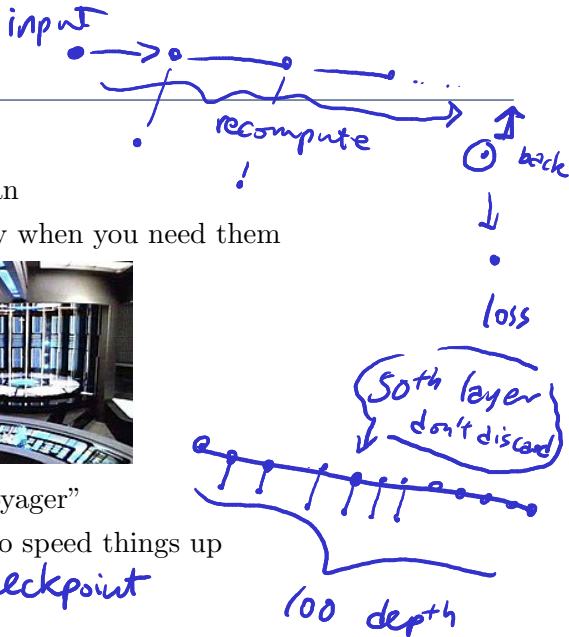can delete (for forward pass)

# Tensor Rematerialization

- Discard the activations as soon as you can
- "Beam" them back/recompute them only when you need them



From Berman et al. 1995 "Star Trek: Voyager"

# Tensor Rematerialization

- Discard the activations as soon as you can
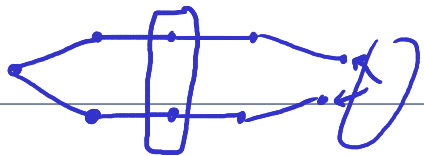- "Beam" them back/recompute them only when you need them



From Berman et al. 1995 "Star Trek: Voyager"

- Save *some* activations as "checkpoints" to speed things up

Only recompute from checkpoint

# Tensor Rematerialization

- Discard the activations as soon as you can
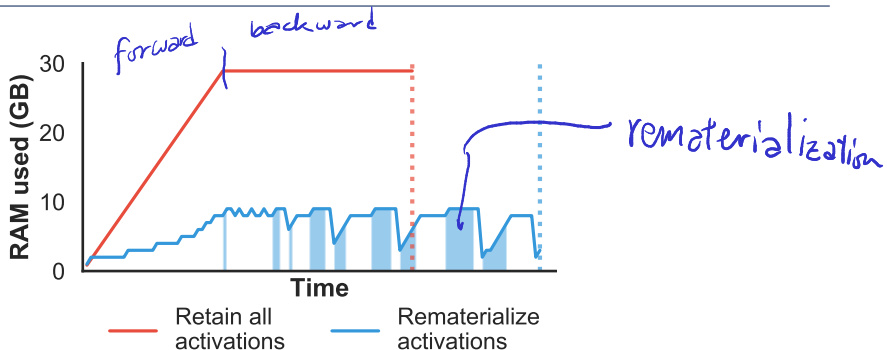- "Beam" them back/recompute them only when you need them



From Berman et al. 1995 "Star Trek: Voyager"

- Save *some* activations as "checkpoints" to speed things up

# Tensor Rematerialization



$h^{(1)} \to h^{(2)} \to h^{(3)}$

forward    backward

rematerialization

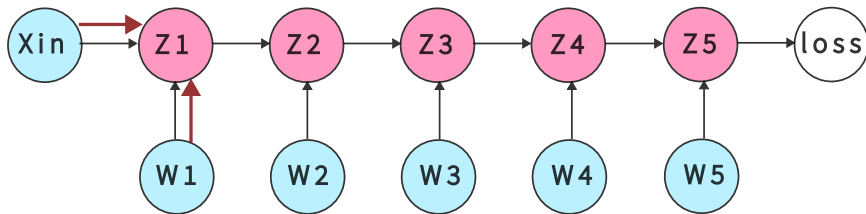Retain all activations

Rematerialize activations

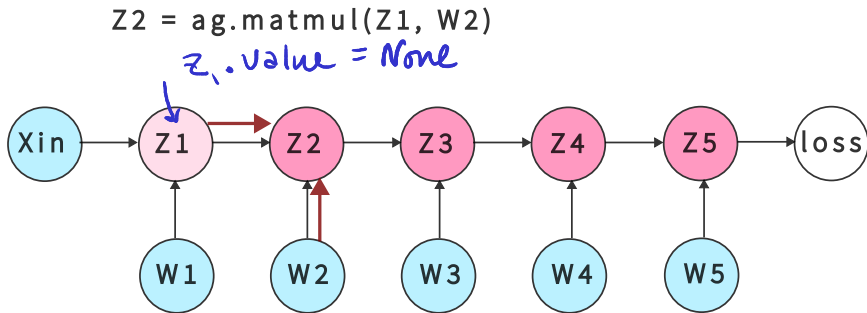From Jain et al. 2020 "Checkmate: Breaking the memory wall with optimal tensor rematerialization"

- Also known as **gradient checkpoint**, activation checkpointing, recomputation...
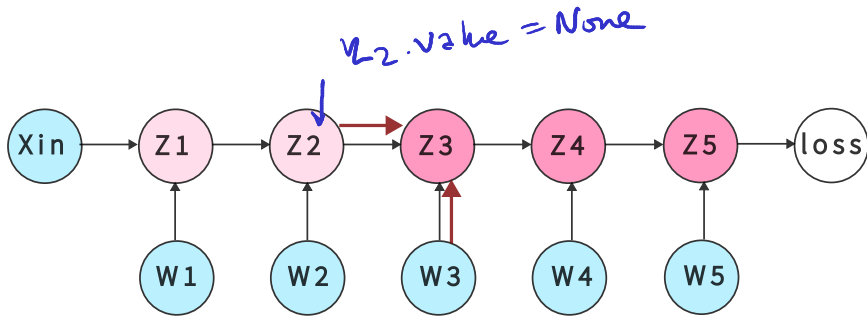
# Rematerialization: forward phase

```
Z1 = ag.matmul(Xin, W1)
```
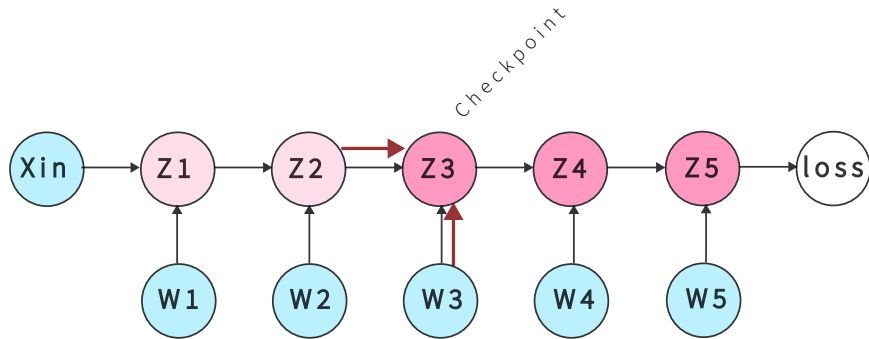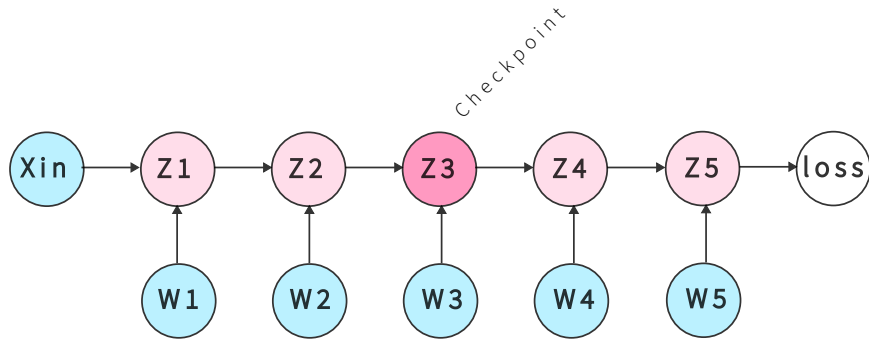
# Rematerialization: forward phase

# Rematerialization: forward phase

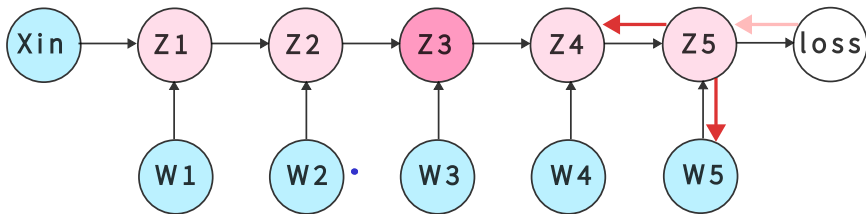# Rematerialization: forward phase

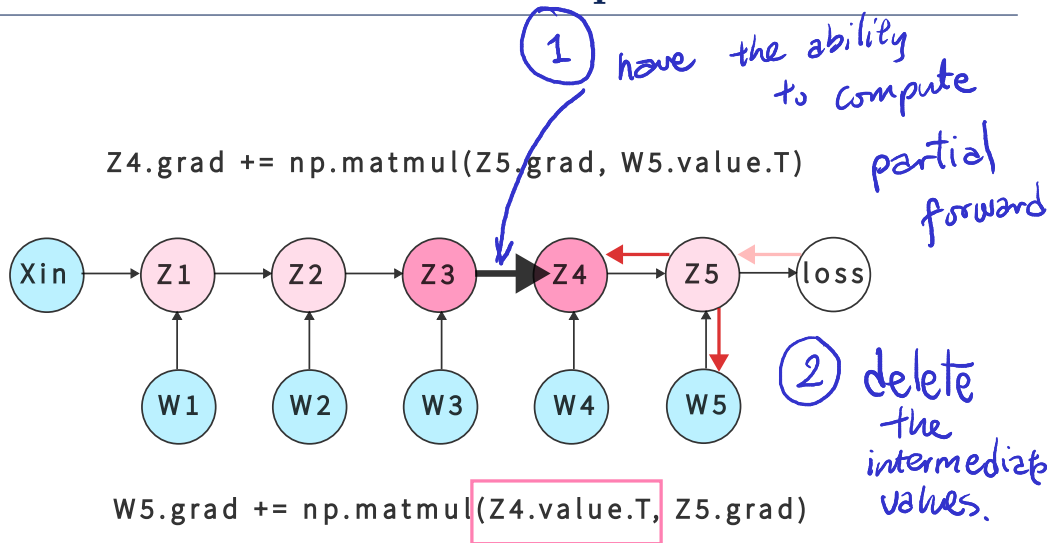# Rematerialization: forward phase

# Rematerialization: backward phase



`Z4.grad += np.matmul(Z5.grad, W5.value.T)`

`W5.grad += np.matmul(Z4.value.T, Z5.grad)`

# Rematerialization: backward phase



1) have the ability to compute partial forward
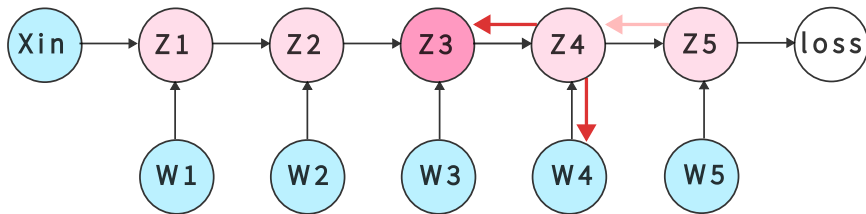
Z4.grad += np.matmul(Z5.grad, W5.value.T)

2) delete the intermediate values.

W5.grad += np.matmul(Z4.value.T, Z5.grad)

# Rematerialization: backward phase

Z3.grad += np.matmul(Z4.grad, W4.value.T)



W4.grad += np.matmul(Z3.value.T, Z4.grad)

# Exercise 2: rematerialization

```
1    class Tensor: # Tensor with grads
2        def __init__(self,
3                     value,
4                     requires_grad=False,
5                     rematerializer = None, # None means don't
    rematerialize, KEEP
6  # [...]
```

*None means checkpoint this tensor no discard.*

*partial forward pass information to each node a function to recompute itself.*

# Exercise 2: rematerialization

```
1  def forward_traced_with_rematerializer(x, weights):
2      start_trace()
3      mem_usage = [] # tracing
4
5      checkpoints = [5] # these layers are checkpoints
6
7      farthest_checkpoint = 0 # this is the input data
8      x_at_farthest_checkpoint = x
9
10     for i, w in enumerate(weights): # main training loop
11         if i in checkpoints:
12             x = ag.matmul(x, w)          ← rematerizlizer = None
13             farthest_checkpoint = i       ← update check
14             x_at_farthest_checkpoint = x  ← keep pointer alive
15 # [...]
```

checkpoint

Xin

# Exercise 2: rematerialization

*ag ↔ pytorch*

*numpy ↔ ATEN*

```
1   for i, w in enumerate(weights): # main training loop
2       if i in checkpoints:
3           x = ag.matmul(x, w)
4           farthest_checkpoint = i
5           x_at_farthest_checkpoint = x
6       else:
7           def _rematerializer():
8               xval = x_at_farthest_checkpoint.value
9               for w in weights[farthest_checkpoint:(i+1)]:
10                  xval = np.matmul(xval, w.value)
11              return xval
12          x = ag.matmul(x, w)
13          x.rematerializer = _rematerializer
```

*not checkpoint*

*do with plain numpy to avoid Autograd*

*recipe for reconstructing x*

# Exercise 2: task 1: rematerialize!

```
1    def backward(self):
2        # [...]
3        for node in reversed(topo_order):
4            # YOUR CODE HERE FOR rematerializing "input.value", if
   it is none
5
6        node._backward()
```
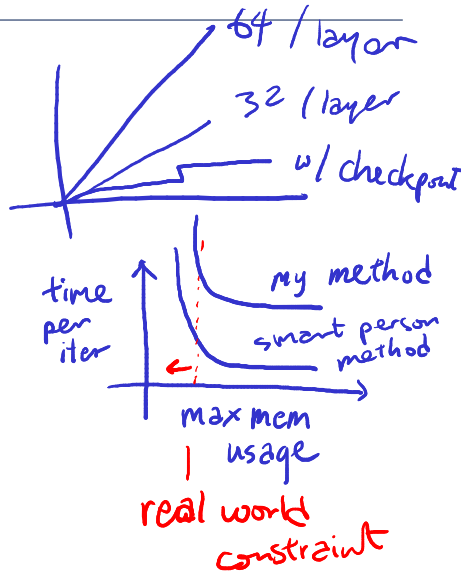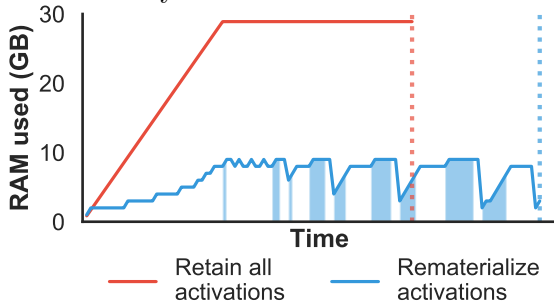
*should see inputs whose values are rematerialization*

# Exercise 2: task 2: discard!

```
1          def __matmul__(self, other):
2              # [...]
3              def _backward():
4                  # [...]
5                  self.grad += np.matmul(output.grad, other.value.T)
6                  other.grad += np.matmul(self.value.T, output.grad)
7
8                  # YOUR CODE HERE FOR discarding activations for "self"
     and "other"
9                  # hint: add a helper function to make it neater
10                 # hint: see "discard_value_if_has_rematerializer" below
11                 # [lines skipped ...]
12                 return None
13
14             output._backward = _backward
15             # YOUR CODE HERE FOR discarding activations for "self" and "
     other"
```

# Summary

- Less memory used in the forward



RAM used (GB) vs Time chart with legend:
- Retain all activations
- Rematerialize activations

# Summary

- Less memory used in the forward



- More computation (for rematerialization).
  (75 matmuls w/ rematerialization vs. 30 matmuls w/o rematerialization).

# Summary

- Less memory used in the forward



- More computation (for rematerialization).
  (75 matmuls w/ rematerialization vs. 30 matmuls w/o rematerialization).
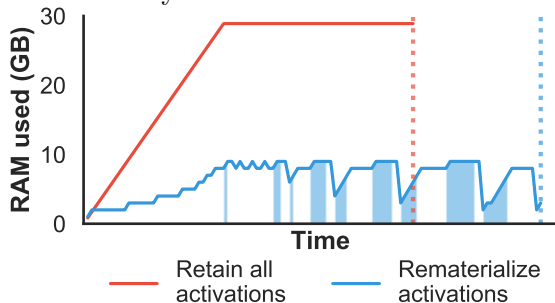- **Question**: what about more complicated computational graph?
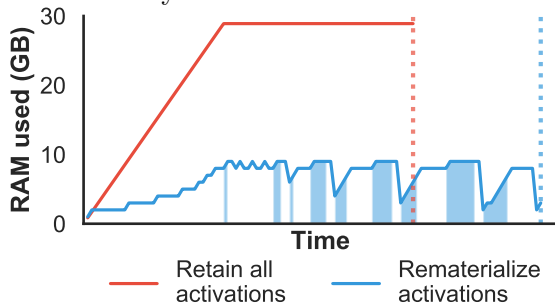
# Summary

- Less memory used in the forward



- More computation (for rematerialization).
  (75 matmuls w/ rematerialization vs. 30 matmuls w/o rematerialization).
- **Question**: what about more complicated computational graph?
- **Question**: optimal way of selecting the checkpoints?
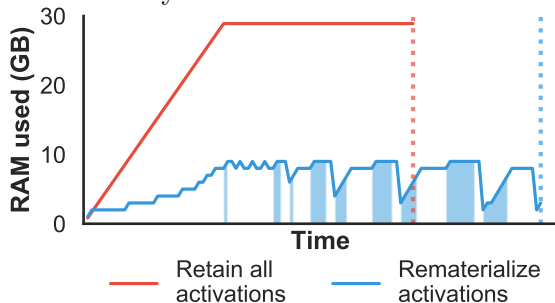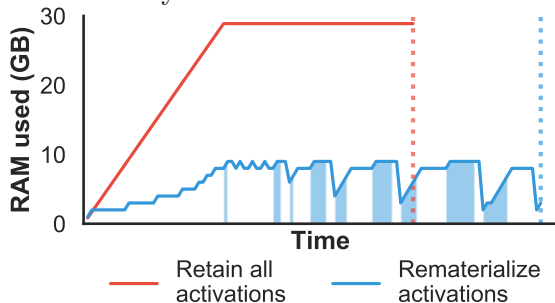
# Summary

- Less memory used in the forward



- More computation (for rematerialization).
  (75 matmuls w/ rematerialization vs. 30 matmuls w/o rematerialization).
- **Question**: what about more complicated computational graph?
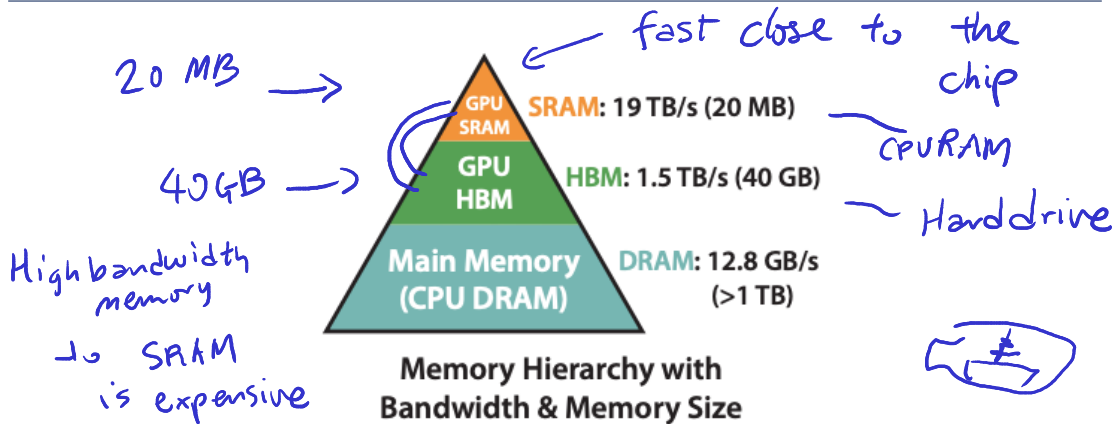- **Question**: optimal way of selecting the checkpoints?
- **Next**: IO and FlashAttention

# Flash attention



From Dao et al. 2022 "FlashAttention: Fast and memory-efficient exact attention with io-awareness"

# Self-attention

$X = $ (matrix) } C Context    Seq_len

d embedding    leftmost dimension

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\texttt{attention}(\mathbf{X}; \theta) := \text{softmax}\left(\mathbf{X}\mathbf{W}^{(Q)}\mathbf{W}^{(K)\top}\mathbf{X}^\top\right)\mathbf{X}\mathbf{W}^{(V)} \in \mathbb{R}^{d \times C}.$$

C = want this to be very large (more tokens)

d = (moderately large) ~ 700

Huge 1,000,000 = C ≫ d = 700

# Self-attention

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\texttt{attention}(\mathbf{X}; \theta) := \text{softmax}\left(\mathbf{X}\mathbf{W}^{(Q)}\mathbf{W}^{(K)\top}\mathbf{X}^\top\right)\mathbf{X}\mathbf{W}^{(V)} \in \mathbb{R}^{d \times C}.$$

- parameters

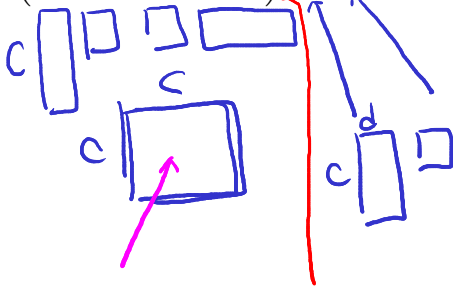$$\theta^{(\texttt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

# Self-attention

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\texttt{attention}(\mathbf{X}; \theta) := \text{softmax}\left(\mathbf{X}\mathbf{W}^{(Q)}\mathbf{W}^{(K)\top}\mathbf{X}^\top\right)\mathbf{X}\mathbf{W}^{(V)} \in \mathbb{R}^{d \times C}.$$

- parameters

$$\theta^{(\texttt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively)
  where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

# Self-attention

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\mathtt{attention}(\mathbf{X}; \theta) := \mathrm{softmax}\left(\mathbf{X}\mathbf{W}^{(Q)}\mathbf{W}^{(K)\top}\mathbf{X}^\top\right)\mathbf{X}\mathbf{W}^{(V)} \in \mathbb{R}^{d \times C}.$$

$\mathbb{R}^{C \times d}$

- parameters

$$\theta^{(\mathtt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

$\mathbb{R}^{C \times d}$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively) where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

- "**seq-2-seq**": maps the sequence $\mathbf{X}$ to another sequence $\mathtt{attention}(\mathbf{X}; \theta)$

# Self-attention

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\texttt{attention}(\mathbf{X}; \theta) := \text{softmax}\Big( \underbrace{\mathbf{X}\mathbf{W}^{(Q)}}_{\mathbf{Q}} \underbrace{\mathbf{W}^{(K)\top}\mathbf{X}^{\top}}_{\mathbf{K}^{\top}} \Big) \underbrace{\mathbf{X}\mathbf{W}^{(V)}}_{\mathbf{V}} \in \mathbb{R}^{d \times C}.$$

- parameters

$$\theta^{(\texttt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively)
  where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

- "**seq-2-seq**": maps the sequence $\mathbf{X}$ to another sequence $\texttt{attention}(\mathbf{X}; \theta)$

## Self-attention

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\mathtt{attention}(\mathbf{X}; \theta) := \mathrm{softmax}\Big(\underbrace{\mathbf{Q}\mathbf{K}^\top}_{\mathbf{S}}\Big)\mathbf{V} \in \mathbb{R}^{d \times C}.$$

- parameters

$$\theta^{(\mathtt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively)
  where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

- "**seq-2-seq**": maps the sequence $\mathbf{X}$ to another sequence $\mathtt{attention}(\mathbf{X}; \theta)$

# Self-attention

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\texttt{attention}(\mathbf{X}; \theta) := \text{softmax}\Big( \underbrace{\underbrace{\mathbf{Q}\mathbf{K}^\top}_{\mathbf{S}}}_{\mathbf{P} = \text{softmax}(\mathbf{S})} \Big) \mathbf{V} \in \mathbb{R}^{d \times C}.$$



row
Stochastic

- parameters

$$\theta^{(\texttt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively) where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

- "**seq-2-seq**": maps the sequence $\mathbf{X}$ to another sequence $\texttt{attention}(\mathbf{X}; \theta)$

## Self-attention

- Let $\mathbf{X} \in \mathbb{R}^{C \times d}$ (note everything is transposed compared to last time)

$$\texttt{attention}(\mathbf{X}; \theta) := \mathrm{softmax}\left(\mathbf{Q}\mathbf{K}^\top\right)\mathbf{V} \in \mathbb{R}^{d \times C}.$$
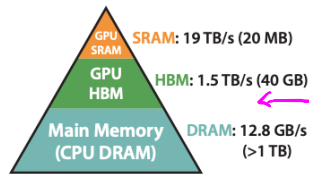
- parameters

$$\theta^{(\texttt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively)
  where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

- "**seq-2-seq**": maps the sequence $\mathbf{X}$ to another sequence $\texttt{attention}(\mathbf{X}; \theta)$

# Flash attention

*small enough*

*loop block over block* $P_{ij}$

$C = 1,000,000$

1: Compute $S = QK^T$

2: $P = \text{softmax}(S)$

3: $O = PV$

*Memory Hierarchy with Bandwidth & Memory Size*
- GPU SRAM: 19 TB/s (20 MB)
- GPU HBM: 1.5 TB/s (40 GB)
- Main Memory (CPU DRAM): 12.8 GB/s (>1 TB)

$Q, K, V$

*rather move $Cd$*

---

**Algorithm 0** Standard Attention Implementation

**Require:** Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.
1: Load $\mathbf{Q}, \mathbf{K}$ by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$, write $\mathbf{S}$ to HBM.
2: Read $\mathbf{S}$ from HBM, compute $\mathbf{P} = \text{softmax}(\mathbf{S})$, write $\mathbf{P}$ to HBM.
3: Load $\mathbf{P}$ and $\mathbf{V}$ by blocks from HBM, compute $\mathbf{O} = \mathbf{P}\mathbf{V}$, write $\mathbf{O}$ to HBM.
4: Return $\mathbf{O}$.

*Many read & write of $C^2$ sized matrix*

From Dao et al 2022 (Flash Attention)

# Tiling

$$d \ll C$$



**Q**  $d$  **K^T**  **S**

$C$

$=$

$$Q = X \, W^{(q)}$$

# Tiling

# Tiling

# Tiling

# Tiling



focus on top row

softmax( Q1@K1^T  Q1@K2^T ) @

V

# Online softmax

# Online softmax

```
1  C = 100
2
3  s = np.random.rand(C)
4
5  I1 = np.arange(0, C//2)
6  I2 = np.arange(C//2, C)
7  s1 = s[I1]
8  s2 = s[I2]
9
10 s - np.hstack([s1,s2]) # make sure we didn't make a silly error
```

= top row

left

right

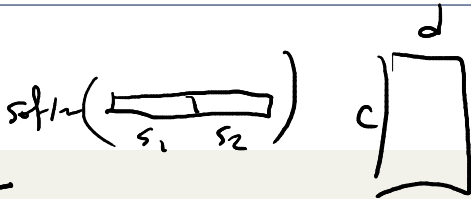# Online softmax

$$\sum \exp(s^{(\cdot)})$$

```python
def softmax(z):
    expz = np.exp(z)
    normalizer = np.sum(expz)
    return expz / normalizer, normalizer
```

probability vector

normalizer

# Online softmax



```
1 d = 3
2 np.random.seed(42)
3 V = np.random.randn(C, d)
4
5 p, n = softmax(s)
6 p @ V
7 # >>> array([ 0.13778369, -0.16034761,  0.04310764])
```

# Online softmax

softmax(s) V ✓

what if

```
1 # block 1
2 p1, n1 = softmax(s1)
3 O = p1 @ V[I1,:]
4
5 # block 2
6 p2, n2 = softmax(s2)
7 O = None # YOUR CODE HERE
8 O
```

$$\text{softmax}(s_1) @ \boxed{\phantom{xx}} V[I1,:] = O_1$$

$$\text{softmax}(s_2) @ \boxed{\phantom{xx}} V[I2,:] = O_2$$

How can we relate this back?

# Online softmax

$$s = \left[ \begin{array}{cc} s^1 & s^2 \end{array} , \begin{array}{cc} s^3 & s^4 \end{array} \right]$$

softmax

$$\underbrace{\qquad}_{s1} \qquad \underbrace{\qquad}_{s2}$$

$$s1 = \left[ s^1 \ s^2 \right]$$

```
1  # block 1
2  p1, n1 = softmax(s1)
3  O = p1 @ V[I1,:]
4
5  # block 2
6  p2, n2 = softmax(s2)
7  O = None # YOUR CODE HERE
8  O
```

$$Softmax(s1) = \left( \begin{array}{c} \dfrac{exp(s^1)}{exp(s^1) + exp(s^2)} \\[3ex] \dfrac{exp(s^2)}{exp(s^1) + exp(s^2)} \end{array} \right.$$

$$\frac{n_1}{n_1 + n_2} = exp(s^1) + \dots + exp(s^4)$$

# Online softmax

$$\left[ \frac{\exp(s^{(1)}), \ \exp(s^{(2)})}{\phantom{x}} \right]$$

$$n1$$

$$\text{softmax}(S) = \left[ \frac{n_1}{(n_1 + n_2)} \cdot \text{softmax}(S1) \right.$$

$$\left. \frac{n_2}{n_1 + n_2} \cdot \text{softmax}(S2) \right]$$

$$\exp(s^{(1)}) + \dots + \exp(s^{(\omega)})$$

$$\text{1-st half}$$
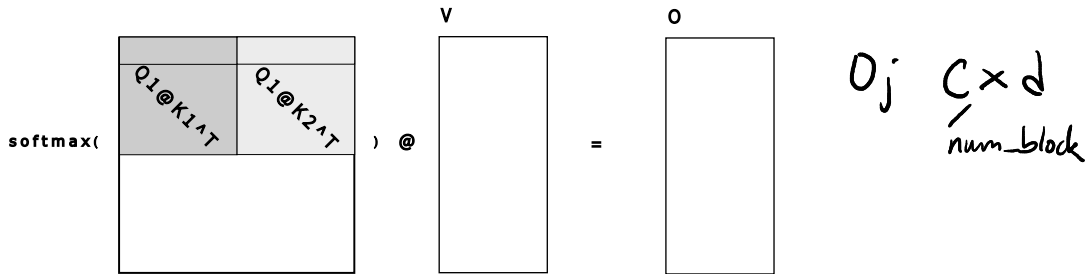
```
1  # block 1
2  p1, n1 = softmax(s1)
3  O = p1 @ V[I1,:]
4
5  # block 2
6  p2, n2 = softmax(s2)
7  O = None # YOUR CODE HERE
8  O
```

# Summary



softmax( $Q_1@K_1{}^T$ | $Q_1@K_2{}^T$ ) @ **V** = **O**

$O_j$ $C \times d$ num_block

- Load $\mathbf{Q}_i$, $\mathbf{K}_j$, $\mathbf{V}_j$ one block at a time
- Compute $\mathbf{S}_{ij} = \mathbf{Q}_i \mathbf{K}_j^\top$  (S = QK usually)
- Compute $\text{softmax}(\mathbf{S}_{ij})\mathbf{V}_j$
- Update $\mathbf{O}_j$ using online softmax

small enough you can compute on chip

# References I

📄 Dao, Tri, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré (2022).
"FlashAttention: Fast and memory-efficient exact attention with io-awareness".
In: *Advances in Neural Information Processing Systems* 35, pp. 16344–16359.

📄 Jain, Paras, Ajay Jain, Aniruddha Nrusimha, Amir Gholami, Pieter Abbeel,
Joseph Gonzalez, Kurt Keutzer, and Ion Stoica (2020). "Checkmate: Breaking the
memory wall with optimal tensor rematerialization". In: *Proceedings of Machine
Learning and Systems* 2, pp. 497–511.

📄 Sohoni, Nimit S, Christopher R Aberger, Megan Leszczynski, Jian Zhang, and
Christopher Ré (2019). "Low-memory neural network training: A technical
report". In: *arXiv preprint arXiv:1904.10631.*