6.5930/1
Hardware Architectures for Deep Learning

# Kernel Computation - Impact of Memory Hierarchy

February 21, 2024

Joel Emer and Vivienne Sze

Massachusetts Institute of Technology
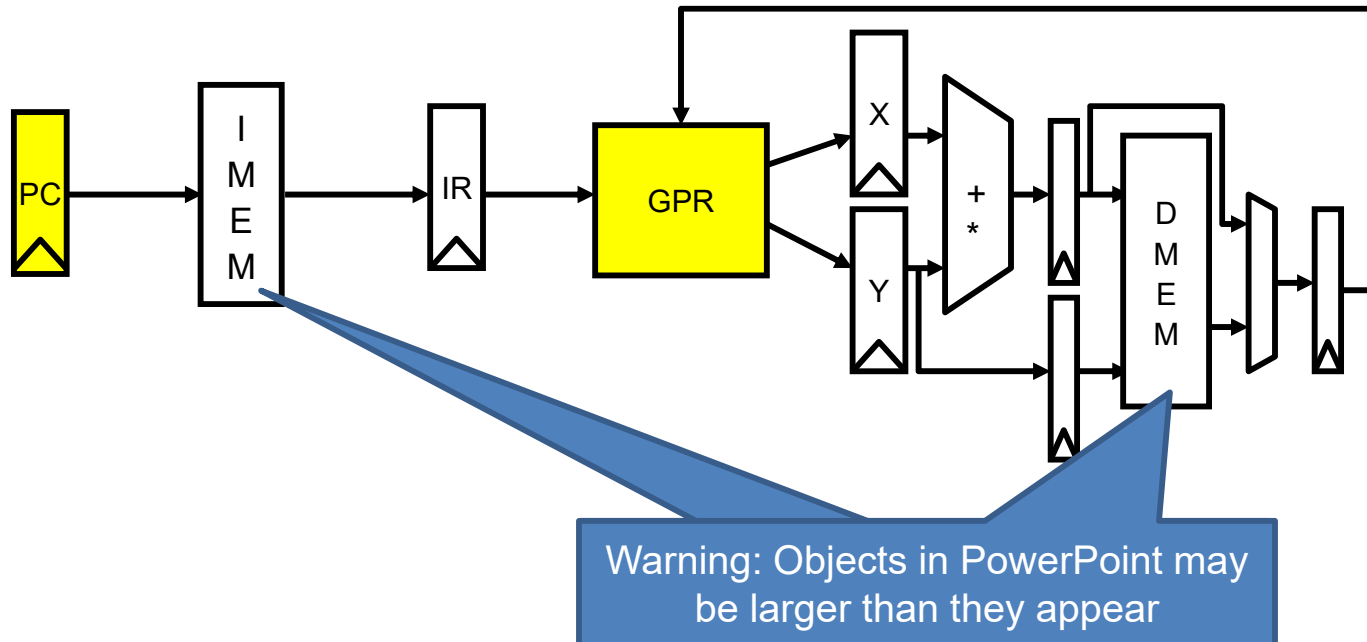Electrical Engineering & Computer Science

# Goals of Today's Lecture

- Understand impact of memory hierarchy

  – Overview of caches

  – Structuring algorithms to work well in caches using tiling
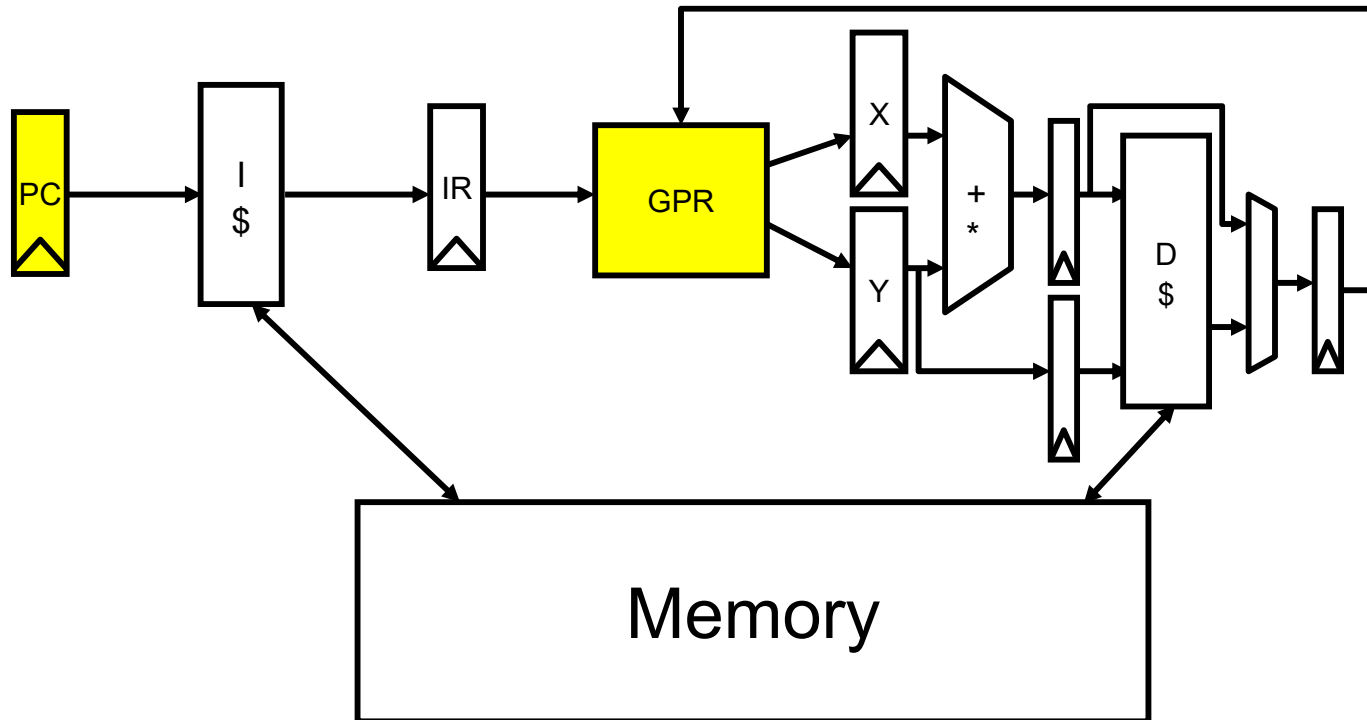
  – Storage technologies

Sze and Emer

# Readings for this Week

- Efficient Processing of Deep Neural Networks

  – Chapter 4 of https://doi.org/10.1007/978-3-031-01766-7

.

# Simple Pipelined µArchitecture



PC · IMEM · IR · GPR · X · Y · + * · DMEM

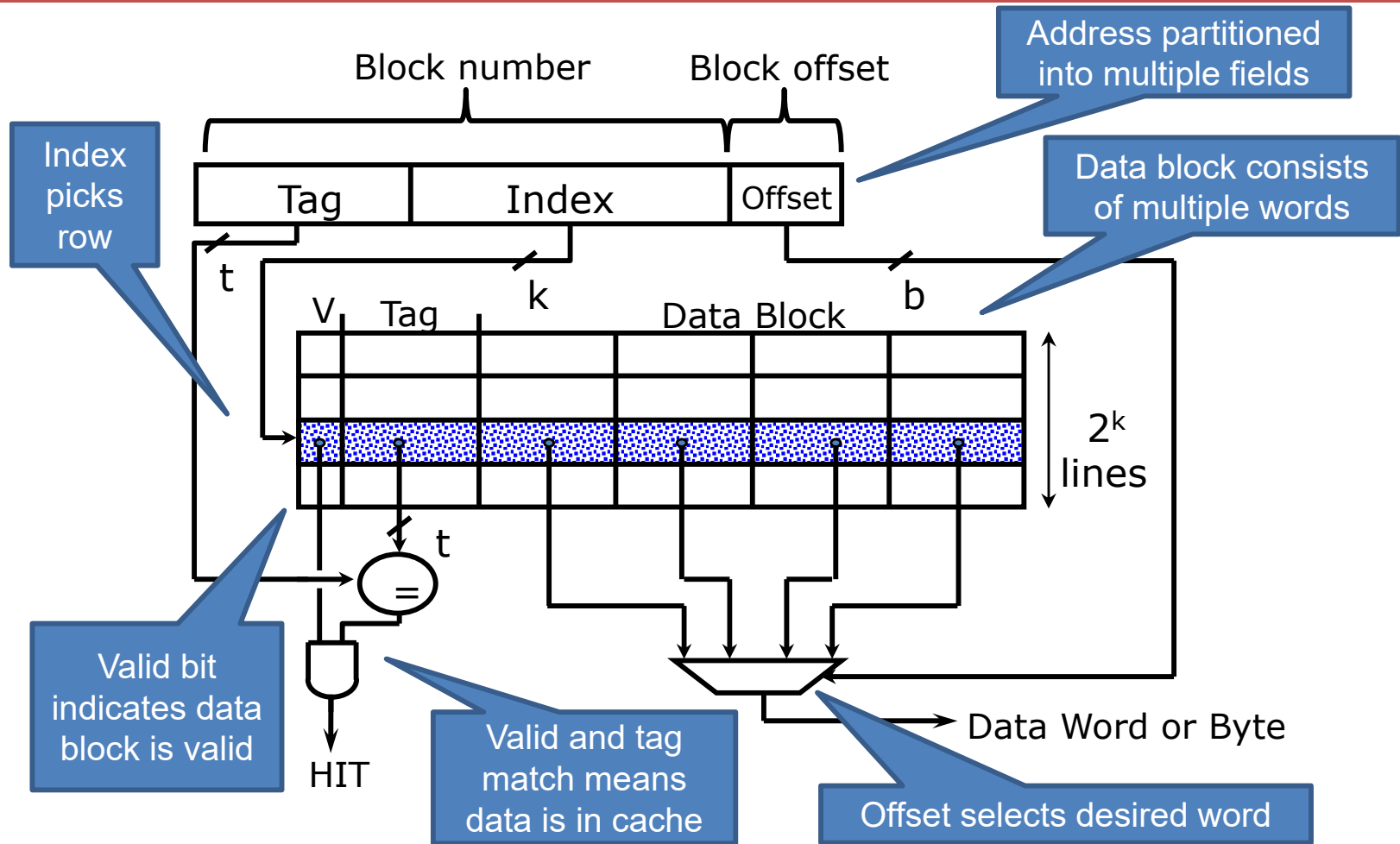Warning: Objects in PowerPoint may be larger than they appear

What are consequences of putting
large memory (e.g., megabytes)
directly in pipeline?

# Pipelined µArchitecture with Caches



Instruction cache (I$) and data cache (D$) hold memory data
for reuse in small energy efficient buffer

# Direct Mapped Cache

Block number    Block offset

Address partitioned into multiple fields

| Tag | Index | Offset |

Data block consists of multiple words

Index picks row

$t$    $k$    $b$

V    Tag    Data Block

$2^k$ lines

Valid bit indicates data block is valid

$t$

=

HIT

Valid and tag match means data is in cache

Data Word or Byte

Offset selects desired word

# Cache Operation

Look at data address, search cache tags to find match.
Then if…

Found in cache
a.k.a.  HIT

Not in cache
a.k.a. MISS

Return copy
of data from
cache

Read block of data from
Main Memory

Wait …

Return data to processor
and update cache

Fill

Metric: Hit Rate = #Hits / (#Hits + #Misses)

Sze and Emer

# Treatment of Writes

- Cache hit:
  - *write through:* write both cache & memory
    - generally higher traffic but simplifies cache in processor pipeline
  - *write back:*    write cache only
    (memory is written only when the entry is evicted)
    - a dirty bit per block can further reduce the traffic

- Cache miss:
  - *no write allocate:*  only write to main memory
  - *write allocate (aka fetch on write):*  fetch into cache

- Common combinations:
  -    write through and no write allocate
  -    write back with write allocate

Sze and Emer

# Cache Locality

Caches **implicitly** try to optimize data movement by trying to exploit two common properties of memory references:

– *Spatial Locality:* If a location is referenced it is likely that locations near it will be referenced in the near future.

  • Exploited by having block size larger than a word, which also amortizes fill overheads by getting more bytes with one access

– *Temporal Locality:* If a location is referenced it is likely to be referenced again in the near future.

  • Exploited by holding blocks for future access

Sze and Emer

# Fully Connected (FC) Computation

Predefined constant

```
int i[CHW];    # Input activations
int f[M*CHW];  # Filter Weights
int o[M];      # Output activations

CHWm = -CHW
for m in [0, M):
  o[m] = 0
  CHWm += CHW
  for chw in [0, CHW):
    o[m] += i[chw] * f[CHWm + chw]
```

M iterations

C*H*W iterations

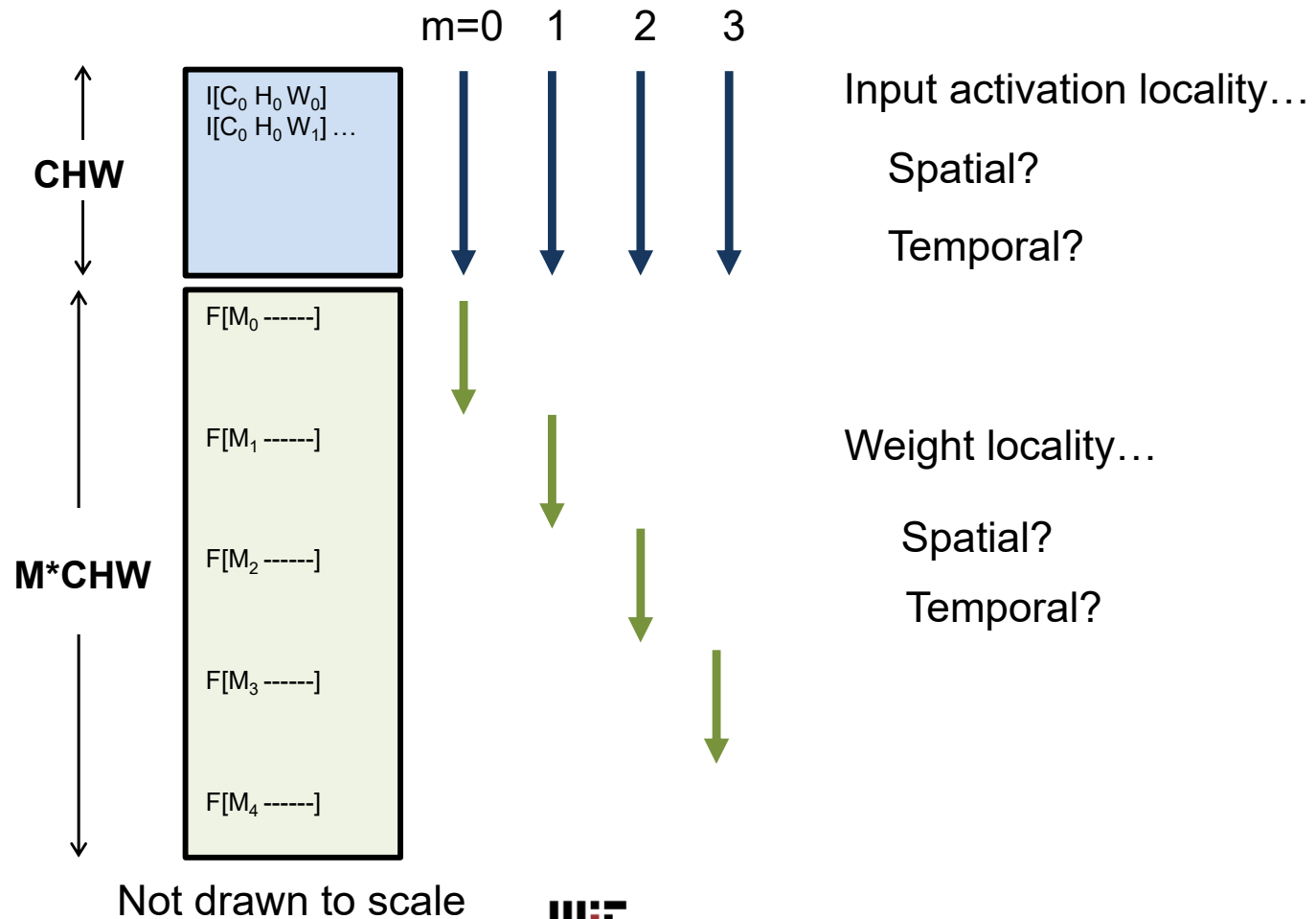M*C*H*W loads of each weight and input activation

February 21, 2024

MIT

# Impact of spatial locality

- Typical in-pipeline cache size
  - 64K bytes => 16K FP32 words
  - 64 byte blocks => 16 FP32 words/block

Hit rate of long sequential reference
streams due to spatial locality?

February 21, 2024

Sze and Emer

# FC – Data Reference Pattern

m=0    1    2    3

$I[C_0 H_0 W_0]$
$I[C_0 H_0 W_1] \dots$

**CHW**

Input activation locality…

Spatial?

Temporal?

$F[M_0 \text{ ------}]$

$F[M_1 \text{ ------}]$

**M*CHW**

$F[M_2 \text{ ------}]$

Weight locality…

Spatial?

Temporal?

$F[M_3 \text{ ------}]$

$F[M_4 \text{ ------}]$

Not drawn to scale

MIT

Sze and Emer

# FC – Data Reference Pattern

# Amount of temporal locality

- Typical in-pipeline cache size

  – 64K bytes => 16K FP32 words

  – 64 byte blocks => 16 FP32 words/block

- Typical layer size:

  – H, W = 256    C = 128

Size of input activations?

What does this imply for
Input activation temporary locality?

# Computational Intensity

Computational Intensity = $\dfrac{MACS}{Data\ Words}$

$$O_m = I_{chw} \times F_{m,chw}$$

Number MACS:

# Computational Intensity – Ideal FC

$$\text{Computational Intensity} = \frac{MACS}{Data\ Words}$$

$$O_m = I_{chw} \times F_{m,chw}$$

Number MACS:    M*C*H*W    Filter weight accesses:

Input activation accesses:

Output activation accesses:

Computational Intensity =

If in one cycle a processor can deliver one word from memory and perform one MAC how well will the machine perform running this code?

MIT

# Computational Intensity – Naïve FC

$$\text{Computational Intensity} = \frac{MACS}{Data\ Words}$$

```
CHWm = -CHW;
for m in [0, M):
  o[m] = 0;
  CHWm += CHW
  for chw in [0, CHW):
    o[m] += i[chw] * f[CHWm + chw]
```

Number MACS:        M*C*H*W        Filter weight accesses

                                   Input activation accesses

                                   Output activation accesses

Computational Intensity =

February 21, 2024

Sze and Emer

# Einsum for strip mined FC

$$O_m = I_{chw} \times F_{m,chw}$$

$$\underbrace{I_{chw/T}}_{chw1}\underbrace{_{,chw\%T}}_{chw0} = I_{chw} \qquad \underbrace{F_{m,chw/T}}_{chw1}\underbrace{_{,chw\%T}}_{chw0} = F_{m,chw}$$

$$O_m = I_{chw1,chw0} \times F_{m,chw1,chw0}$$

# Fully Connected – Strip Mined

```
for m in [0, M):
  for chw in [0, C*H*W):
    o[m] += i[chw] * f[CHW*m + chw]
```

```
// Level 1
for chw1 in [0, CHW1):
  for m in [0, M):
// Level 0
    for chw0 in [0, CHW0):
    chw = CHW0*chw1+chw0
    o[m] += i[chw] * f[CHW*m + chw]
```

CHW1*CHW0 = C*H*W

Inner loop working set = CHW0

Just considering input activations,
what value should CHW0 be?

Sze and Emer

# FC - Strip Mined Data Reference Pattern



CHW

M*CHW

I[$C_0$ $H_0$ $W_0$]
I[$C_0$ $H_0$ $W_1$]
…

F[$M_0$ ------]

F[$M_1$ ------]

F[$M_2$ ------]

F[$M_3$ ------]

F[$M_4$ ------]

Not drawn to scale

Untiled

Tiled

Cache Hits?

# Computational Intensity – Strip Mined

```
// Level 1
for chw1 in [0, CHW1):
  for m in [0, M):
// Level 0
    for chw0 in [0, CHW0):
      chw = CHW0*chw1+chw0
      o[m] += i[chw] * f[CHW*m + chw]
```

Number MACS:         M*C*H*W         Filter weight accesses:

Input activation accesses:

Output activation accesses

Computational Intensity =

# Matrix-Vector Multiply – Strip Mined



Tensor: f_MCHW[M, CHW]

Rank: CHW

Rank: M

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| 0 | 9 | 4 | 7 | 2 | 4 | 5 | 4 | 3 |
| 1 | 2 | 7 | 3 | 9 | 8 | 4 | 3 | 3 |
| 2 | 2 | 5 | 8 | 7 | 9 | 2 | 5 | 7 |
| 3 | 7 | 2 | 2 | 1 | 7 | 6 | 2 | 3 |

Tensor: i_CHW[CHW]

Rank: CHW

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 3 | 4 | 7 | 5 | 1 | 6 | 2 | 5 |

Tensor: unknown[M]

Rank: M

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |

February 21, 2024

Sze and Emer

# Associative Cache



Allows multiple streams to be resident at same time

Pick data from 'way' that 'hits'

Data Word or Byte

HIT

# Cache Miss Pipeline Diagram

HIT

| | | | | | | |
|---|---|---|---|---|---|---|
| ld r6, w(r5) | IF | ID | RF | EX | D$ | WB |
| mul r7,r4,r6 | | IF | ID | RF | stall | EX | D$ | WB |

MISS

| | | | | | | |
|---|---|---|---|---|---|---|
| ld r6, w(r5) | IF | ID | RF | EX | D$ MISS - MEM | WB |
| mul r7,r4,r6 | | IF | ID | RF | stall | EX | D$ | WB |

Time (cycles)

MIT

# Avoiding Cache Miss Stalls

- Reorganize code so loads are far ahead of use

  - Requires huge amount of unrolling

  - Consumes lots of registers


- Add 'prefetch' instructions that just load cache

  - Consumes instruction issue slots


- Add hardware that automatically loads cache

# Hardware Data Prefetching

- ## Prefetch-on-miss:
  - Prefetch b + 1 upon miss on b

- ## One Block Lookahead (OBL) scheme
  - Initiate prefetch for block b + 1 when block b is accessed
  - Can extend to N block lookahead

- ## Strided prefetch
  - If observe sequence of accesses to block b, b+N, b+2N, then prefetch b+3N etc.

Example: IBM Power 5 [2003] supports eight independent streams of strided prefetch per processor, prefetching 12 lines ahead of current access

MIT

# Multi-level Caches

- A memory cannot be large and fast

- Add level of cache to reduce miss penalty
  - Each level can have longer latency than level above
  - So, increase sizes of cache at each level



Metrics:

Local miss rate = misses in cache/ accesses to cache

Global miss rate = misses in cache / CPU memory accesses

Misses per instruction = misses in cache / number of instructions

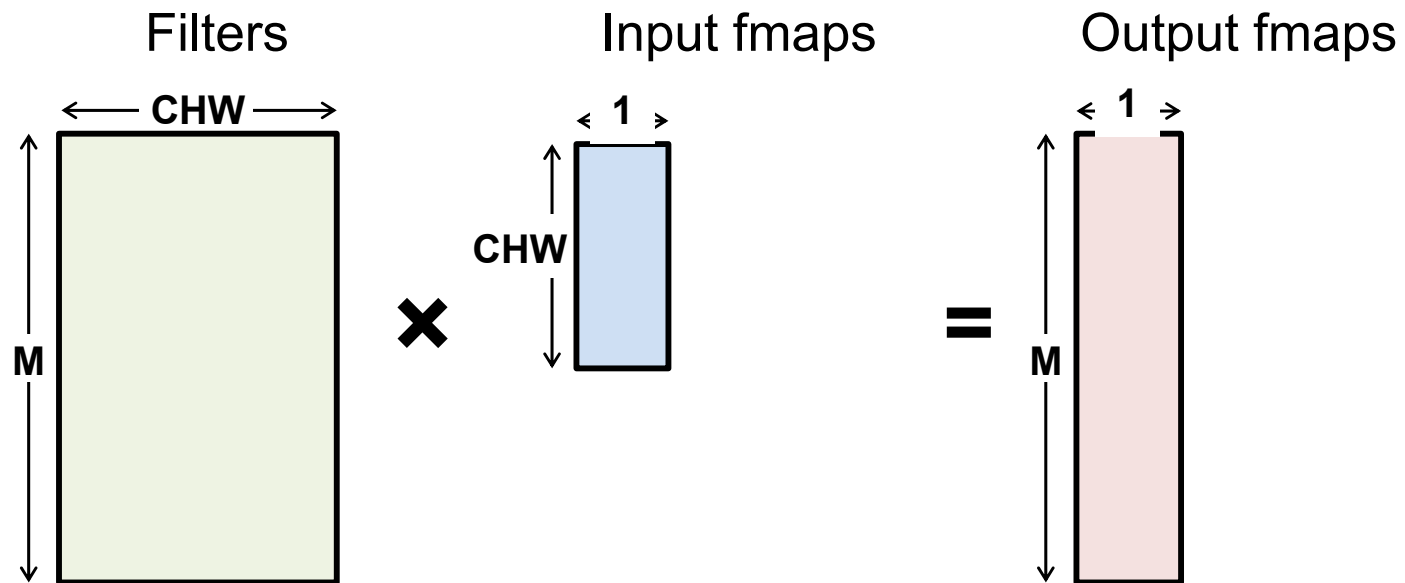# Contemporary CPU Cache Hierarchy



(a) Memory hierarchy for server

(b) Memory hierarchy for a personal mobile device

# FC Layer – Multichannel

filters

input fmaps

output fmaps

# Fully-Connected (FC) Layer

Filters $\times$ Input fmaps $=$ Output fmaps
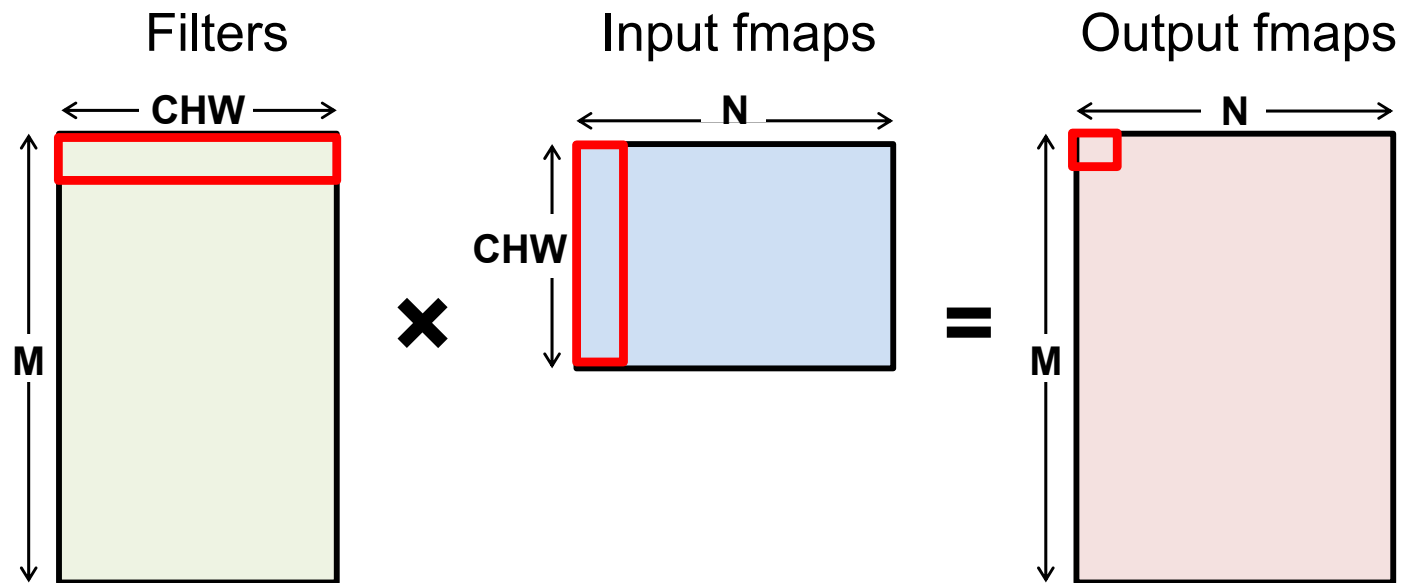
# Fully-Connected (FC) Layer

Filters          Input fmaps        Output fmaps



- After flattening, having a batch size of N turns the matrix-vector operation into a matrix-matrix multiply

# FC Einsum Notation
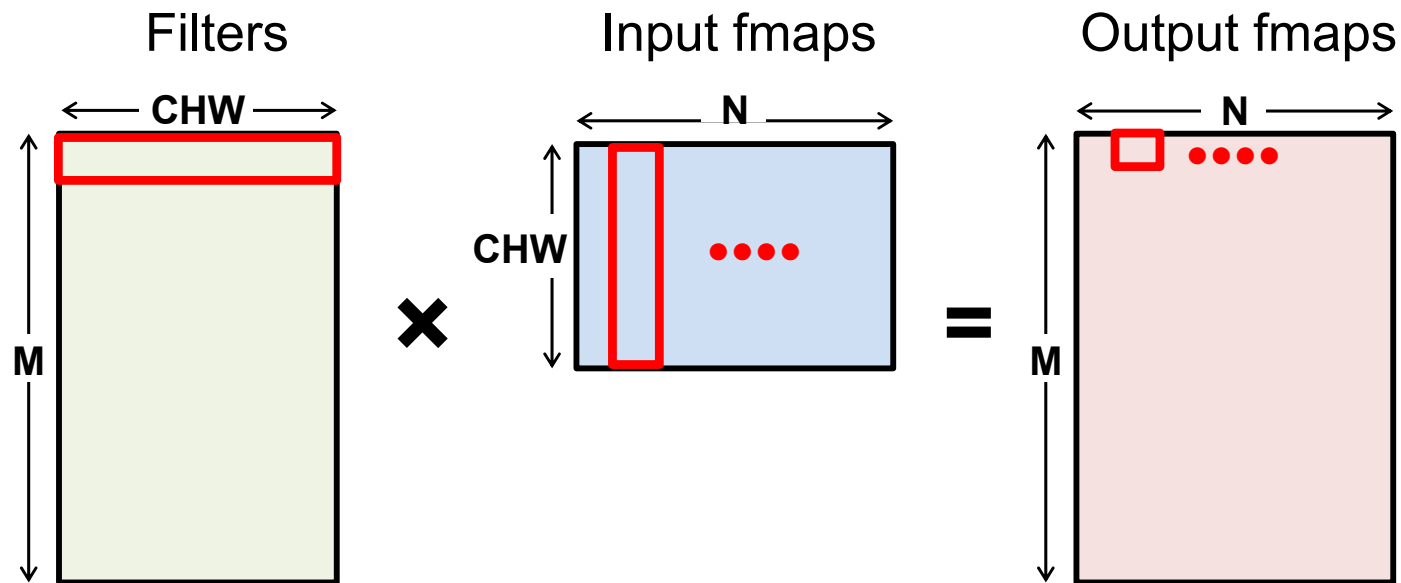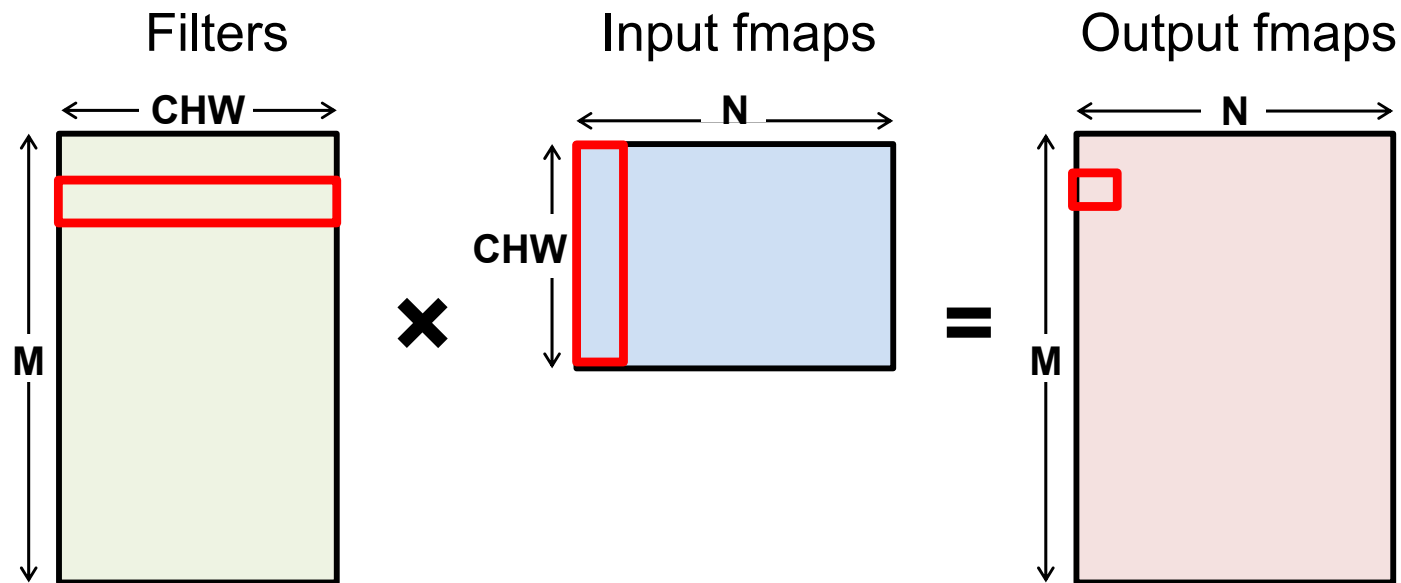
$$O_{n,m} = F_{m,chw} \times I_{n,chw}$$



Filters      Input fmaps      Output fmaps

# Fully-Connected (FC) Layer

Filters      Input fmaps      Output fmaps



- After flattening, having a batch size of N turns the matrix-vector operation into a matrix-matrix multiply

Sze and Emer

# Fully-Connected (FC) Layer

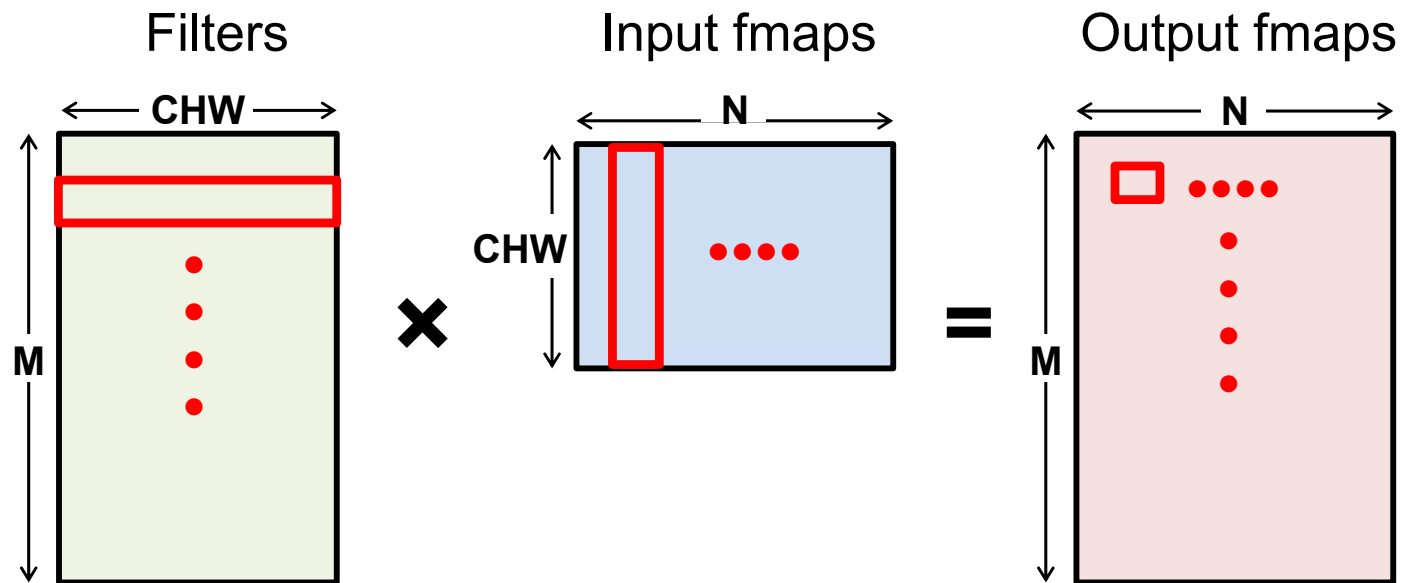Filters                    Input fmaps              Output fmaps



- After flattening, having a batch size of N turns the matrix-vector operation into a matrix-matrix multiply

# Fully-Connected (FC) Layer



Filters — Input fmaps — Output fmaps

- After flattening, having a batch size of N turns the matrix-vector operation into a matrix-matrix multiply

# Fully-Connected (FC) Layer

Filters      Input fmaps      Output fmaps



- After flattening, having a batch size of N turns the matrix-vector operation into a matrix-matrix multiply

How much temporal locality for naïve implementation?  None

# Matrix-Matrix Multiply

Tensor: f_MCHW[M, CHW]

Rank: CHW

| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|
| Rank: M | 0 | 9 | 4 | 7 | 2 | 4 | 5 | 4 | 3 |
| | 1 | 2 | 7 | 3 | 9 | 8 | 4 | 3 | 3 |
| | 2 | 2 | 5 | 8 | 7 | 9 | 2 | 5 | 7 |
| | 3 | 7 | 2 | 2 | 1 | 7 | 6 | 2 | 3 |

Tensor: i_NCHW[N, CHW]

Rank: CHW

| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|
| Rank: N | 0 | 6 | 7 | 5 | 9 | 7 | 1 | 5 | 1 |
| | 1 | 4 | 2 | 5 | 4 | 8 | 9 | 3 | 8 |
| | 2 | 9 | 2 | 9 | 7 | 7 | 7 | 8 | 9 |
| | 3 | 8 | 9 | 2 | 4 | 1 | 3 | 3 | 9 |

Tensor: unknown[M, N]

Rank: N

| | | 0 | 1 | 2 | 3 |
|---|---|---|---|---|---|
| Rank: M | 0 | 0 | 0 | 0 | 0 |
| | 1 | 0 | 0 | 0 | 0 |
| | 2 | 0 | 0 | 0 | 0 |
| | 3 | 0 | 0 | 0 | 0 |

February 21, 2024

Sze and Emer

# Matrix-Matrix Multiply Tiled

Tensor: f_MCHW[M, CHW]

Rank: CHW

|        |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|--------|---|---|---|---|---|---|---|---|---|
| Rank: M | 0 | 9 | 4 | 7 | 2 | 4 | 5 | 4 | 3 |
|        | 1 | 2 | 7 | 3 | 9 | 8 | 4 | 3 | 3 |
|        | 2 | 2 | 5 | 8 | 7 | 9 | 2 | 5 | 7 |
|        | 3 | 7 | 2 | 2 | 1 | 7 | 6 | 2 | 3 |

Tensor: i_NCHW[N, CHW]

Rank: CHW

|        |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|--------|---|---|---|---|---|---|---|---|---|
| Rank: N | 0 | 6 | 7 | 5 | 9 | 7 | 1 | 5 | 1 |
|        | 1 | 4 | 2 | 5 | 4 | 8 | 9 | 3 | 8 |
|        | 2 | 9 | 2 | 9 | 7 | 7 | 7 | 8 | 9 |
|        | 3 | 8 | 9 | 2 | 4 | 1 | 3 | 3 | 9 |

Tensor: unknown[M, N]

Rank: N

|        |   | 0 | 1 | 2 | 3 |
|--------|---|---|---|---|---|
| Rank: M | 0 | 0 | 0 | 0 | 0 |
|        | 1 | 0 | 0 | 0 | 0 |
|        | 2 | 0 | 0 | 0 | 0 |
|        | 3 | 0 | 0 | 0 | 0 |

February 21, 2024

Sze and Emer

# Tiled Fully-Connected (FC) Layer

### Filters

$\xleftarrow{\text{CHW}}\rightarrow$

| $F_{0,0}$ | $F_{0,1}$ |
|:---:|:---:|
| $F_{1,0}$ | $F_{1,1}$ |

M

$\times$

### Input fmaps

$\xleftarrow{\quad N \quad}\rightarrow$

CHW

| $I_{0,0}$ | $I_{0,1}$ |
|:---:|:---:|
| $I_{1,0}$ | $I_{1,1}$ |

$=$

### Output fmaps

$\xleftarrow{\quad N \quad}\rightarrow$

M

| $F_{0,0}I_{0,0}$ + $F_{0,1}I_{1,0}$ | $F_{0,0}I_{0,1}$ + $F_{0,1}I_{1,1}$ |
|:---:|:---:|
| $F_{1,0}I_{0,0}$ + $F_{1,1}I_{1,0}$ | $F_{1,0}I_{0,1}$ + $F_{1,1}I_{1,1}$ |

Matrix multiply tiled to fit in cache
and computation ordered to maximize reuse of data in cache

MIT

# Tiled Fully-Connected (FC) Layer

Filters

Input fmaps

Output fmaps



Matrix multiply tiled to fit in cache
and computation ordered to maximize reuse of data in cache

# Tiled Fully-Connected (FC) Layer



Filters — CHW, M

Input fmaps — N, CHW

Output fmaps — N, M

$$\begin{bmatrix} F_{0,0} & F_{0,1} \\ F_{1,0} & F_{1,1} \end{bmatrix} \times \begin{bmatrix} I_{0,0} & I_{0,1} \\ I_{1,0} & I_{1,1} \end{bmatrix} = \begin{bmatrix} F_{0,0}I_{0,0} + F_{0,1}I_{1,0} & F_{0,0}I_{0,1} + F_{0,1}I_{1,1} \\ F_{1,0}I_{0,0} + F_{1,1}I_{1,0} & F_{1,0}I_{0,1} + F_{1,1}I_{1,1} \end{bmatrix}$$

Matrix multiply tiled to fit in cache
and computation ordered to maximize reuse of data in cache

*Dotted line means partial result

February 21, 2024

Sze and Emer

# Tiled Fully-Connected (FC) Layer

Filters

Input fmaps

Output fmaps

$$
\begin{array}{c|c}
F_{0,0} & F_{0,1} \\
\hline
F_{1,0} & F_{1,1}
\end{array}
\quad \times \quad
\begin{array}{c|c}
I_{0,0} & I_{0,1} \\
\hline
I_{1,0} & I_{1,1}
\end{array}
\quad = \quad
\begin{array}{c|c}
F_{0,0}I_{0,0} + F_{0,1}I_{1,0} & F_{0,0}I_{0,1} + F_{0,1}I_{1,1} \\
\hline
F_{1,0}I_{0,0} + F_{1,1}I_{1,0} & F_{1,0}I_{0,1} + F_{1,1}I_{1,1}
\end{array}
$$

Matrix multiply tiled to fit in cache
and computation ordered to maximize reuse of data in cache

# Einsum for tiled FC

$$O_m = I_{n,chw} \times F_{m,chw}$$

$$I_{n,chw} \rightarrow I_{n1,chw1,n0,chw1}$$

$$F_{m,chw} \rightarrow F_{m1,chw1,m0,chw0}$$

$$O_{m1,m0} = I_{n1,chw1,n0,chw0} \times F_{m1,chw1,m0,chw0}$$

Sze and Emer

# Fully-Connected (FC) Layer

- Implementation: **Matrix Multiplication (GEMM)**

  - **CPU:** OpenBLAS, Intel MKL, etc
  - **GPU:** cuBLAS, cuDNN, etc

- Library will note shape of the matrix multiply and select implementation optimized for that shape.

- Optimization usually involves proper tiling to storage hierarchy

# Tradeoffs in Memories

# Overview of Memories

Memory consist of arrays of cells that hold a value.

- Types of Memories/Storage
    - Latches/Flip Flops (Registers)
    - SRAM (Register File, Caches)
    - DRAM (Main Memory)
    - Flash (Storage)

# Elements of Memory Operation

Implementations vary based on:

- How a memory cell holds a value?

- How is a value obtained from a memory cell?

- How is a value set in a memory cell?

- How is array constructed out of individual cells?

- Results in tradeoffs between cost, density, speed, energy and power consumption

MIT

Sze and Emer

# Latches/Flip Flops

D-flip flop

- Fast and low latency

- Located with logic

CLK

Image source: 6.111

*Example from CPU pipeline*

# Latches/Flip Flops (< 0.5 kB)

D-flip flop

- Fast and low latency

- Located with logic

- Not very dense

  - 10+ transistors per bit

  - Usually use for arrays smaller than 0.5kB

CLK

D

Image source: 6.111

*Array of Flip flops*

Read
Address
$[A_2:A_0]$

MIT

# Latches/Flip Flops (< 0.5 kB)

*Array of Flip flops*



Read
Address
$[A_2:A_0]$

Sze and Emer

# SRAM

- ## Higher density than register

  - Usually, 6 transistors per bit-cell

IC wafer



- ## Less robust and slower than latches/flip-flop

Bit cell size 0.75um$^2$ in 14nm

# SRAM (kB – MB)

# SRAM

# SRAM Power Dominated by Bit Line
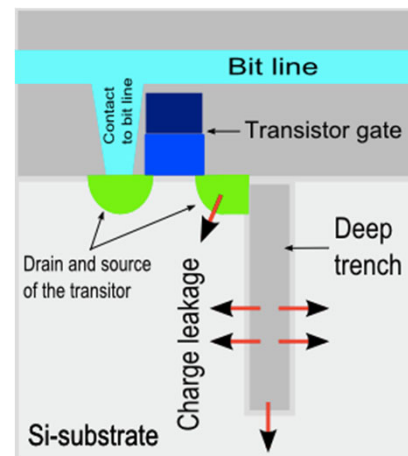
## Measured SRAM Power Breakdown



Legend:
- Bit-lines (BL)
- Word-line (WL)
- Sensing Ntwk.
- Other

Pie chart values: 56%, 22%, 15%, 6%

$@V_{DD}=0.6V$

Image Source: Mahmut Sinangil

Larger array → Longer bit-lines
→ Higher capacitance → Higher power

February 21, 2024

Sze and Emer

# DRAM

- ## Higher density than SRAM

  - 1 transistor per bit-cell

  - Needs periodic refresh
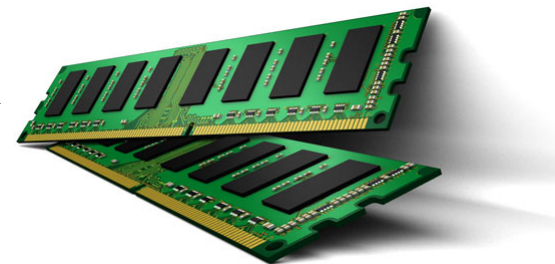
- ## Special device process

# DRAM (GB)

- ## Higher density than SRAM

  - 1 transistor per bit-cell

  - Needs periodic refresh

- ## Special device process

  - Usually off-chip (except eDRAM – which is pricey!)

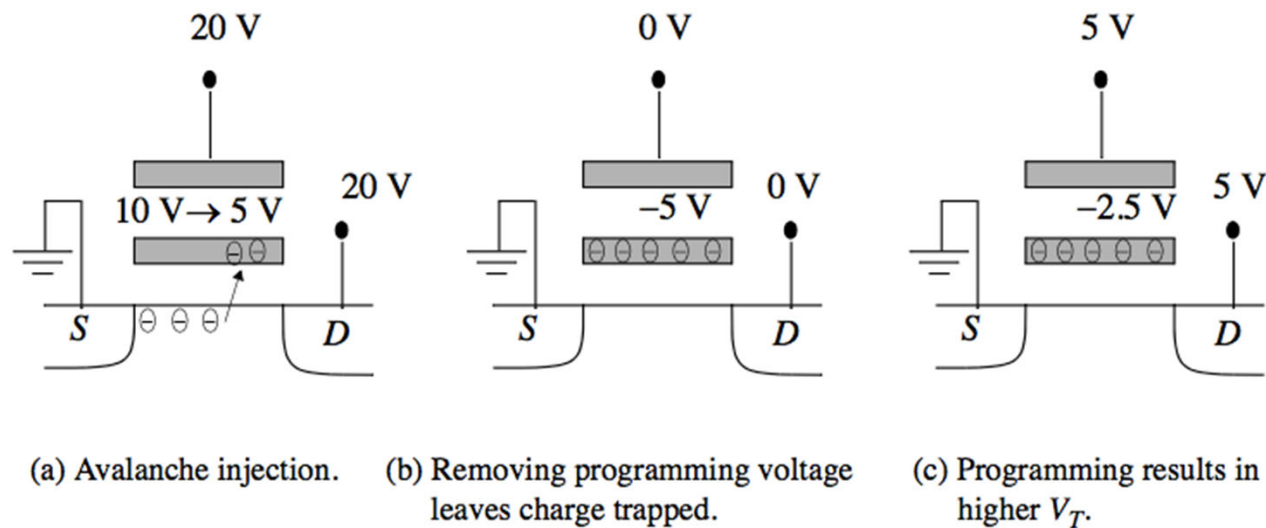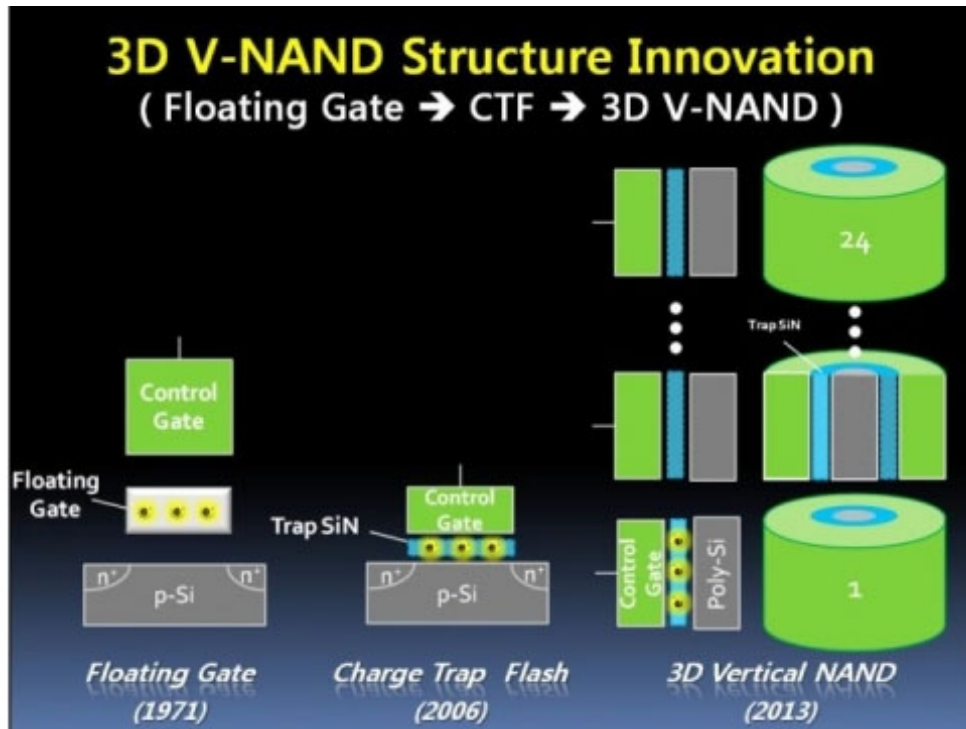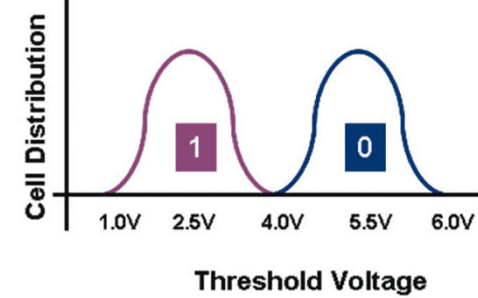  - Off-chip interconnect has much higher capacitance



pJ

nJ

Sze and Emer

# Flash (100GB to TB)

- More dense than DRAM

- Non-volatile

  – Needs high powered write (change $V_{TH}$ of transistor)



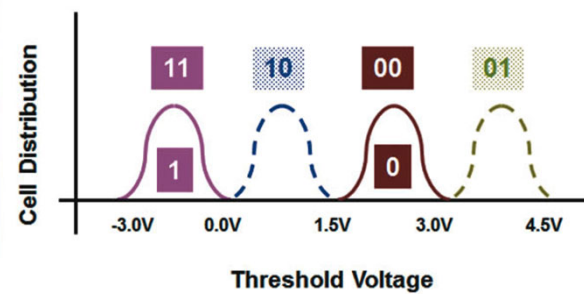(a) Avalanche injection.   (b) Removing programming voltage leaves charge trapped.   (c) Programming results in higher $V_T$.

Sze and Emer

# Flash Memory



Single Level Cell (SLC)

Multi-levels cell (MLC)

48 layer, Ternary level cell (TLC)
*Aug 2015*
*256 Gb per die (for SSD)*

# Memory Tradeoffs



**Cost/bit**
Function of circuit type (smaller → cheaper)

**Latency**
Function of circuit type
(smaller → slower)
and total capacity
(smaller → faster)

**Energy/access/bit**
Function of total capacity
(smaller → less energy)
and circuit type
(smaller → less energy)

**Bandwidth**
Increases with
parallelism

**Density**
Function of circuit type
(smaller → denser)

Most attributes tend to improve with technology scaling,
lower voltage and sometimes smaller capacitors

# Summary

- Reduce main memory access with caches

  - Main memory (i.e., DRAM) is slow and has high energy consumption

  - Exploits spatial and temporal locality

- Tiling to reduce cache misses

  - Possible since processing order does not affect result (MACs are commutative)

  - Add levels to loop nest to improve temporal locality

  - Size of tile depends on cache size and cache associativity

- Tradeoffs in storage technology

  - Various tradeoffs in cost, speed, energy, capacity…

  - Different technologies appropriate at different spots in the design

Sze and Emer

## Next Lecture: Vectorization

Thank you!