Date 10/10/2024

→ TPU, NPU

→ CPU and Memory

→ MAC units

$$c = a + b$$
$$d = c + f$$

→ ultimate bottleneck is memory.

→ Harvard Arch. → seperate data and instruction.

→ Dataflow Arch.

→ 10 - 100 ms

→ Main Memory is based on DRAM.

L1 → 1 cycle
L2 → few ll
L3 → few 10 cycle

$$\boxed{\begin{array}{l} c = a + b \\ d = c + f \end{array}}$$

temporal
locality

→ Memory locality

Spatial

$\{ x_{00} \ x_{01} \ x_{02} \cdots \longrightarrow \}$ $\begin{cases} y_{00} \\ y_{01} \\ y_{02} \\ \vdots \\ \downarrow \end{cases}$ $\begin{bmatrix} y_{00} \\ y_{01} & \downarrow \\ y_{02} & \downarrow \\ \vdots & \downarrow \\ \downarrow \end{bmatrix}$

→ Pre-fetching

L1 is subset of L2
L2 is " " L3

→ DRAM
→ DIMM Module

Channel, ranks, banks

rank — 64B data
↳ bus width

→ 64 ms refresh rate of DRAM

→ Sense Amplifier ⟶ first sense then amplify
↳ Peripheral Circuitry
↳ A buffer
↳ row buffer locality

NVMAin

→ DRAMSIM

→ Ramulator

→ closed row policy
reads one destructive (then restore dat)

→ open row policy
to avoid precharging

→ order of commands → typical flow for every
command
ACT
RD
PRE

24/10/2024

Take   1K × 1K   matrix.

INT   4B

4B×1K

L2 → 1MB

Then   note   the   execution   time.

By   tiling   we   can   exploit   the   temporal   locality

for ( i = i - 100 ) {
    for ( j = 1 - 100 ) {

for ( i = i - 50 )
    for ( j =

-0   and   -3

Stency → 

0 - 1   Black & White

0 - 255   Grayscale

mccol

HBM → High Bandwidth Channel

64B → 4B

1 w/s     16 w/s

→ Decompression is on critical path

→ Latency & Throughput is same for non-pipelining system.

→ OOO processor

→ In order, keep Pipeline in order.

→ Asm A/S

① NVSim
② CACTi

SISD , SIMD

$$\begin{bmatrix} A_{00} & A_{10} \end{bmatrix} \begin{bmatrix} b_{00} \\ b_{10} \end{bmatrix}$$

GPUs are SIMD

→ GPT 4 training resources

DSL → Domain Specific language.

→ Scratchpad memory.
   └ equivalent of cache but is software managed.

HMC → multiple stack layers

Assignment #1
   DRAM timing parameters

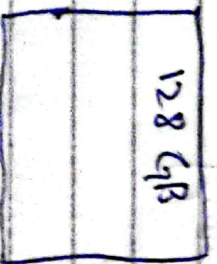   TACT
   upto 20 different commands.

→ GP

⇒ Compute bound

→ Memory ↑↑    (bounded by load store)

    Ops/byte ↑↓

→ SPEC 2017 → benchmark for compute bound/
                                          Memory ↑↓

→ Near Memory Computation

→ Two Categories ① CIM ② CNM

```
┌─────────┐
│ 128 GB  │
└─────────┘
```

→ Cost Model.

→ must read  " The landscape  "

→ french company upmem

⇒ DRAM near Computation (1970s)
⇒ MRAM (Main RAM)
⇒ WRAM (working ↑↑)

→ Software Emulated
  ↳ Software toolchain converts it int.

→ Synchronous / Asynchronous.

↳ Gcc LLVM has own IR

→ MLIR → Multilevel IR

→ Cinnamon flow

→ BLAS → baseline for optimization.

→ Command for broadcast, scatter and gather

→ Alveo also uses HBM memory.

↳ AiM → 2 neighbor banks can communicate

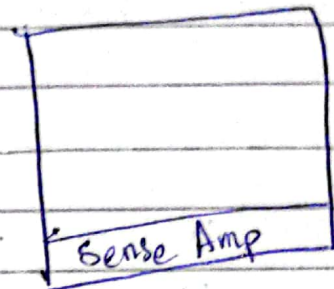↳ GDDR & compressed b/w HBM and HMC)

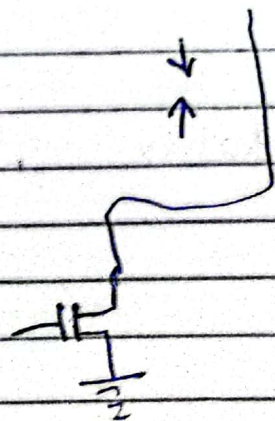↳ In Memory computing Next class.
    DRAM, SRAM.

Date 7/11/2024

Crossbar

→ $O(n^3)$ — $O(1)$

word size  512 - 4KB

CAMS → Good for searching.

→ SRAM are fastest

| 800 K | 100 k |
|-------|-------|
| LRS | HRS |
| 1 | 0 |

Sense Amp

| 800k | 400k | 200k | 100k |
|------|------|------|------|
| 11 | 10 | 01 | 00 |

→ Low Resistance State

High    "    "

→ Racetrack Memory → Similar to HDD

→ writing is extremely slow.

→ Endurance how many times a cell can be read/written.

→ Mem resistor

$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} g \qquad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} V$$

writing entire Matrix

→ Kcl

→ ISAAC
    Shift and hold.

→ PUMA Sim
→ Sim-P