18-643 Lecture 10: FPGA Memory Architecture

James C. Hoe Department of ECE Carnegie Mellon University

18-643-F23-L10-S1, James C. Hoe, CMU/ECE/CALCM, ©2023

Housekeeping

Your goal today: think about memory from a spatial (non-monolithic) computing perspective

(I assume you have taken a comp arch course)

- Notices
 - Handout #5: lab 2, due noon, 10/9
 - Project status report due each Friday
- Readings: (see lecture schedule online)
 - Intel[®] FPGA SDK for OpenCL[™] Pro Edition: Best Practices Guide
 - Compute Express Link[™] 2.0 White Paper

Stored Program Architecture a.k.a. von Neumann



• Memory holds both program and data

paradigm since its invention

- instructions and data in a linear memory array
- - update some state (e.g. PC and memory) as a function of current state according to instruction
 - 4. repeat

18-643-F23-L10-S3

CMU/ECE/CALCM, ©2023

program (count	ter		٦		
0 1	L 2 3	4 5	 	↓	 _	



Necessary Awareness #1:

Fast means Small and Expensive

- Bigger is slower
 - SRAM 512 Bytes @ sub-nsec
 SRAM KByte~MByte @ nsec
 DRAM GByte @ ~50 nsec
 SSD TByte @ msec
 Hard Disk TByte @ ~10 msec
- Faster is more expensive (dollars and chip area)
 - SRAM ~\$10K per GByte
 - DRAM ~\$10 per GByte
 - "Drives" ~\$0.1 per GByte

Treat the values as ×/÷ 3x

not hov

man

times

think it mean

Necessary Awareness #2: Locality is Good

- Temporal: after accessing A, how many other distinct addresses before accessing A again
- Spatial: after accessing A, how many other distinct addresses before accessing a <u>near-by</u> B not what you
- Good locality implies
 - easier to cover more of working-set with a small (fast) memory
 - lower BW and/or more efficient use of BW to large (slow) memory

MMM::good; sparse MMM::bad; streaming::1 of 2



Necessary Awareness #3: Data Movement not Free

- Latency: transit time between source and dest
- Overhead: dead time spent in the act of sending or receiving <u>not overlapped with concurrent</u> <u>compute</u>
- Gap: wait time in between successive send's or receive's due to limited transfer BW

see LogP [Culler, et al., PPoPP93]



accounting for

communication

^{cost} with memory



Necessary Awareness #3.1: Memory Perf Matters

- An algorithm has a cost in terms of operation count
 - runtime_{compute-bound} = # operations / FLOPS
- An algorithm also has a cost in terms of number of bytes communicated (ld/st or send/receive)
 - runtime_{BW-bound} = # bytes / BW
- Which one dominates depends on
 - ratio of FLOPS and BW of platform
 - ratio of ops and bytes of algorithm
- Average Arithmetic Intensity (AI)
 - how many ops performed per byte accessed
 - # operations / # bytes

DRAM, SRAM and all that

On-Chip Fast/Small Memory



- LUT-RAM: 64x1b or 32x2b
 - dual-port: [sync write/async read] + async read
 - also good as one 32b or two 16b shift-regs
- BRAM: 18Kb, variable aspect ratio 1~36b wide
 - true dual-port: 2x[sync read/write], separate clocks
 - fast-enough to be double-pump'ed
- Become FF array if infeasible as hard macro



FPGA Memory Peculiarities

- Quantized, disjoint memory options
 - jumps between FF-based vs. LUT-RAM vs. BRAMs
 - choose from fixed menu of sizes and aspect ratios
 - subtle limitations (comb-read, multi-write)
- Large memory (BRAM) abnormally fast and "free" until your run-out
- Must manage RAM usage
 - don't waste BRAM on small buffers
 - tune buffer sizes to natural granularities, e.g., zero incremental cost to go from 1KB to 2KB
 - pack buffers to share same physical array



- Simple asynchronous request/reply queues
 - in-order or out-of-order (need tags)

er than SRAMs since coarser, further away from "action"

- multiple queues, separate read/write queues
- Need very wide data word for bandwidth

4B@200MHz is only 0.8GB/sec

• Long and variable <u>but predictable</u> latency

18-643-F23-L10-S11, James C. Hoe, CMU/ECE/CALCM, ©2023

Access Pattern affects Latency and BW

DRAM organization

- (multiple ranks per DIMM)
- (multiple chips per rank)
- multiple banks per chip
- Per bank
 - long delay to new row
 - very fast to same row
- Chip/package/board add to total latency
- Rows refreshed every 64ms
 - bank unavailable for
 30~40ns at a time
 - avg. ~1% unavailability



FPGAs gaining more SRAM capacity and DRAM bandwidth

- Interest in FPGA computing sets new "balance"
- 10 years ago
 - single-digit GB/sec to DRAM
 - single-digit Mbit SRAM
- New FPGAs (more than Moore increase)
 - GPU-level memory bandwidth (HBM, HMC, ...)
 - 10s MByte SRAM
 - also add GPU-level FP throughput

No free lunch in Watt/Perf though

• Will we see hardened memory hierarchy?

E.g., Stratix-10 MX



HBM2 on Stratix-10 MX

- Not a monolithic structure or abstraction
- 32 x "pseudo" channels (32B wide upto 500MHz)
 - each channel to a separate DRAM array
 - 16GB/sec/chnl if bursting consecutive addresses
- 512GB/sec if algorithm
 - partition data 32 ways
 - issue 32 concurrent,
 - burst accesses wo. conflict
 - digest 32 x 32B per cyc
 - run at 500MHz



No walk-in-the-park on GPUs either

18-643-F23-L10-S15, James C. Hoe, CMU/ECE/CALCM, ©2023

Programmable Fabric

Memory Controller

Hard NoC: making memory BW usable



Memory Controller

Memory Controller

Compute Express Link (CXL)



In PCI, host cannot cache device space; device can read host space coherently with host cache but cannot cache

Memory Organization as Design Optimization

Classic Memory Hierarchy 101 (pre-multicore)

- Memory hierarchy level i has access time of t_i
- Perceived access time T_i is longer than t_i
 - a chance (hit-rate h_i) you find what you want $\Rightarrow t_i$
 - a chance (miss-rate m_i) you don't find it $\Rightarrow t_i + T_{i+1}$
 - $-h_{i} + m_{i} = 1.0$
- In general

$$T_{i} = h_{i} \cdot t_{i} + m_{i} \cdot (t_{i} + T_{i+1})$$
 think this of as

$$T_{i} = t_{i} + m_{i} \cdot T_{i+1} \leftarrow m_{i} \cdot T_{i+1}$$
 "miss penalty"
ote: h_{i} and m_{i} are of the references missing at level i-
h_{bottom-most}=1.



FPGA hierarchy used differently

- 200MHz soft-logic cache
 - a miss to DRAM is not too many cycles away
 - if 4-byte access, 100%
 hit-rate only 0.8 GB/sec
- Remember to think **SPATIAL**
 - distributed concurrent
 ibandwidth! for spatial
 kernels
 - reduce off-chip memory bandwidth and **jpower!**



10s GB/sec per core 1 cycle latency @ GHZ+

> 10s GB/sec per channel 100s core clock latency



18-643-F23-L10-S20, James C. Hoe, CMU/ECE/CALCM, ©2023

Cache vs Scratchpad

- Manual scratchpad is easy for regular/structured locality
 - per-kernel scratchpad more opportunity and benefit in specialization
 - HW management does not lengthen code
 - prefetching can hide memory latency completely
 - 95% of the time: streaming or double-buffering

These easy cases actually against cache heuristics

• Cache is useful when locality is not predictable ahead of time

Customized the cache if you use one!!!

Banking

- not the same as DRAM chip banks • Partition storage onto multiple structures
- More BW for parallel, non-conflicting accesses



applicable to SRAM, DRAM or disk "storage"

word sel

lg₂W

Control over Data Layout

• An array of N words; index is lg₂N bits

lg₂N array index (sequential)

line

lg₂N/B/W

- N-word total storage Ig₂B
 - divided into B banks; bank number is lg₂B bits
 - each bank is W-word wide; word-select is lg₂W bits

bank #

- line index within bank is $lg_2(N/B/W)$ bits
- Assign bank #, word select and index to maximize





18-643-F23-L10-S24, James C. Hoe, CMU/ECE/CALCM, ©2023

A Small Concrete Example: N=16, W=4



18-643-F23-L10-S25, James C. Hoe, CMU/ECE/CALCM, ©2023

Blocked CNN Kernel on Local Memory

```
local cnndata_t BufI[_Tn][_Tr*_S_wts+_K_wts-1][_Tc*_S_wts+_K_wts-1];
local cnndata_t BufO[_Tm][_Tr][_Tc];
local cnndata_t BufW[_Tm][_Tn][_K_wts][_K_wts];
```

 Assuming sequential execution, what is the access sequence on BufW and Bufl?

}}}}}

ignore **BufO** for now



Bufl more challenging



What should happen?



How to Layout BufW

- # of bank, width of bank (height is derived)
- bank number, line index, word select



How to Layout Bufl

- # of bank, width of bank (height is derived)
- bank number, line index, word select



What should happen?



How many '*'/cycle possible?

What about BufO?





Some Hints on Lab 2

- Optimize tiled kernel throughput out-of-context
 - maximize arithmetic op/cyc
 - buffers need to match arithmetic concurrency
 - stay within Ultra96 resources (don't run out of LUT, DSP, or BRAM)

Don't worry about full layer size and data reuse

- Degrees of freedom
 - tile size
 - loop ordering (affects access pattern)
 - loopnest level to pipeline (higher perf but also cost)
 Buffer design and data layout tied to above choices

Start from the Shallow End

for(row b=0;row b< Tr;row b++){</pre> for(col_b=0;col_b<_Tc;col_b++){</pre> for(to_b=0;to_b<_Tm;to_b++){</pre> for(ti_b=0;ti_b<_Tn;ti_b++){</pre> for(i=0;i<_K_wts;i++){</pre> for(j=0;j<_K_wts;j++){
 #pragma HLS_PIPELINE</pre> Buf0[to b][row b][col b]+= BufW[to b][ti b][i][j]* BufI[ti b][S wts*row b+i] [S wts*col b+j];

}}}}}

Play with the Design Space

}}}}}

Parting Thoughts

- Memory architecture and memory performance important to computing
- Spatial FPGA computing needs adjusted intuition
 - wider, more concurrent access
 - slow clock ticks
 - amenable to extreme specialization
- Look forward to
 - more SRAM, faster DRAM
 - more data movement BW within fabric
 - cache-coherence, better integration
 - hardwired, native memory architecture?