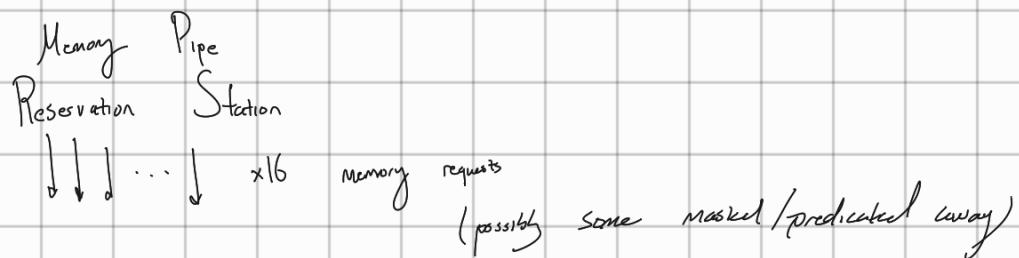
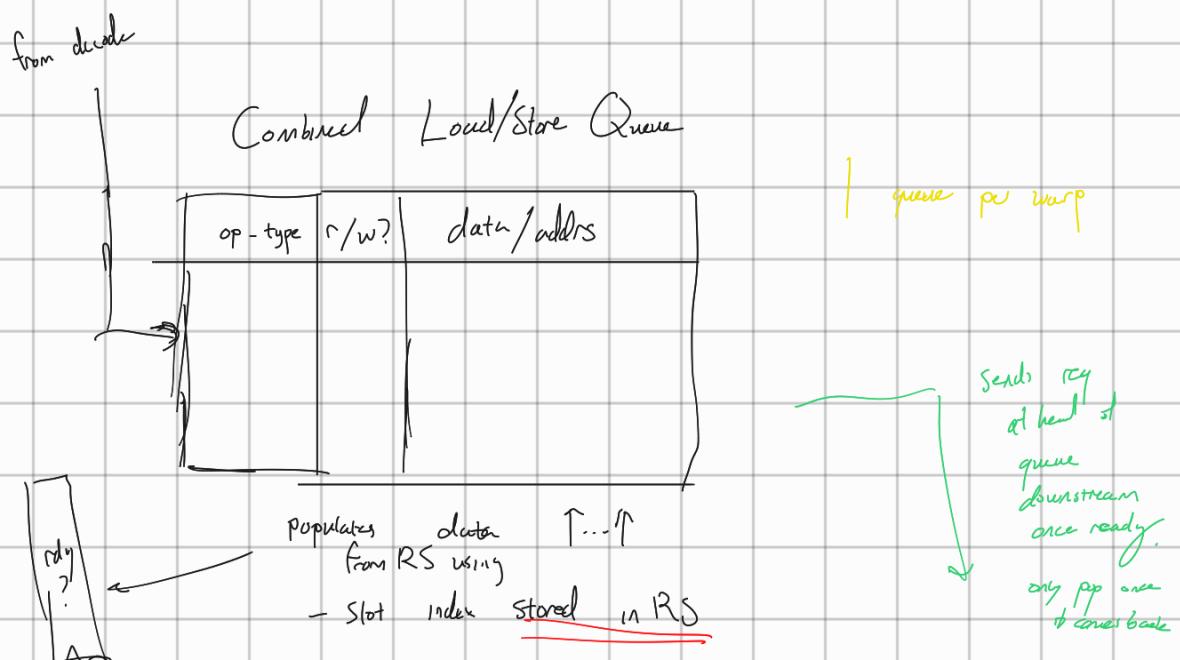


General Notes:

- don't support unaligned accesses
(i.e. for N -byte load, address % N must be 0)
thus is not a serious limitation \rightarrow CUDA has same
- Memory ordering is quite relaxed \rightarrow loads and stores by one thread can appear in any order to another
- Coherence is also very relaxed:
use write-through L1;
L2 is coherent. However L1 can contain stale data



D_{Config} 1 : keep loads/stores in program order



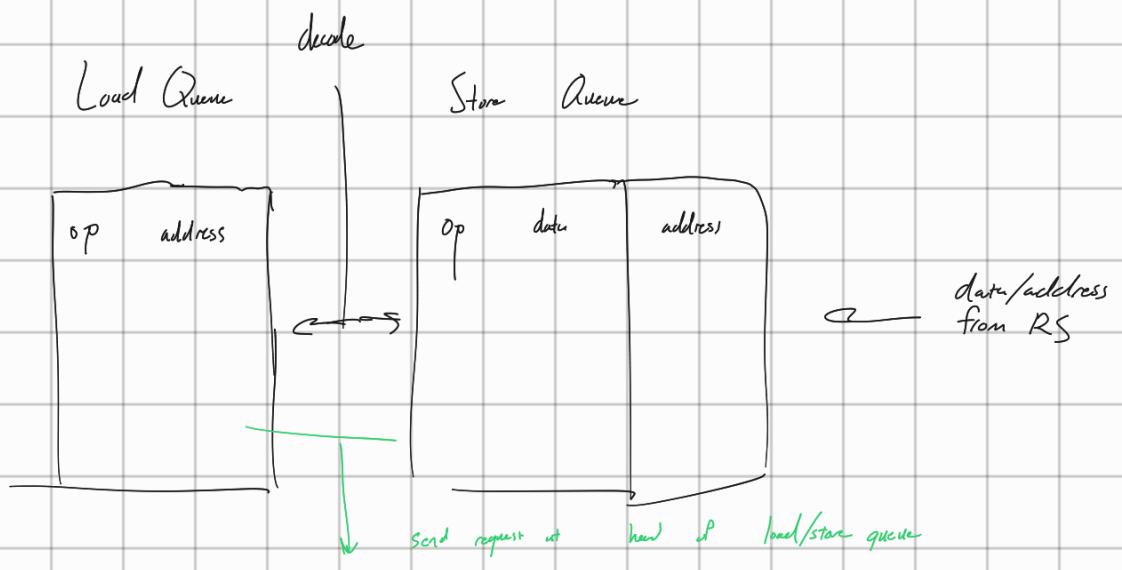
- Slots allocated in-order from decode stage
- only 1 outstanding memory transaction in flight per warp (too limiting?)

Cost is essentially 1 SRAM to act as queue

LRZW would be optimal but probably not feasible

Width dominated by $16 \times 32 = 512$ needed for store data and address
 $+ 16 \times 32 = 512$
(a lot is wasted for reads...)

Design 2: keep loads/stores in program order, but separate them physically



- as before, slots allocated in program order from decode,
slot index comes from RS

- when allocating a Load Queue entry,
write tail pointer of Store Queue

- Load not allowed to issue until
head of store queue reaches that stored value
(i.e. all older stores retired)

- Store not allowed until its at head of queue
and load queue head \geq stored tail ptr

- Also forces program order but allows sizing

the queues differently

Design 2.5: keep stores in program order, allows reordering consecutive loads

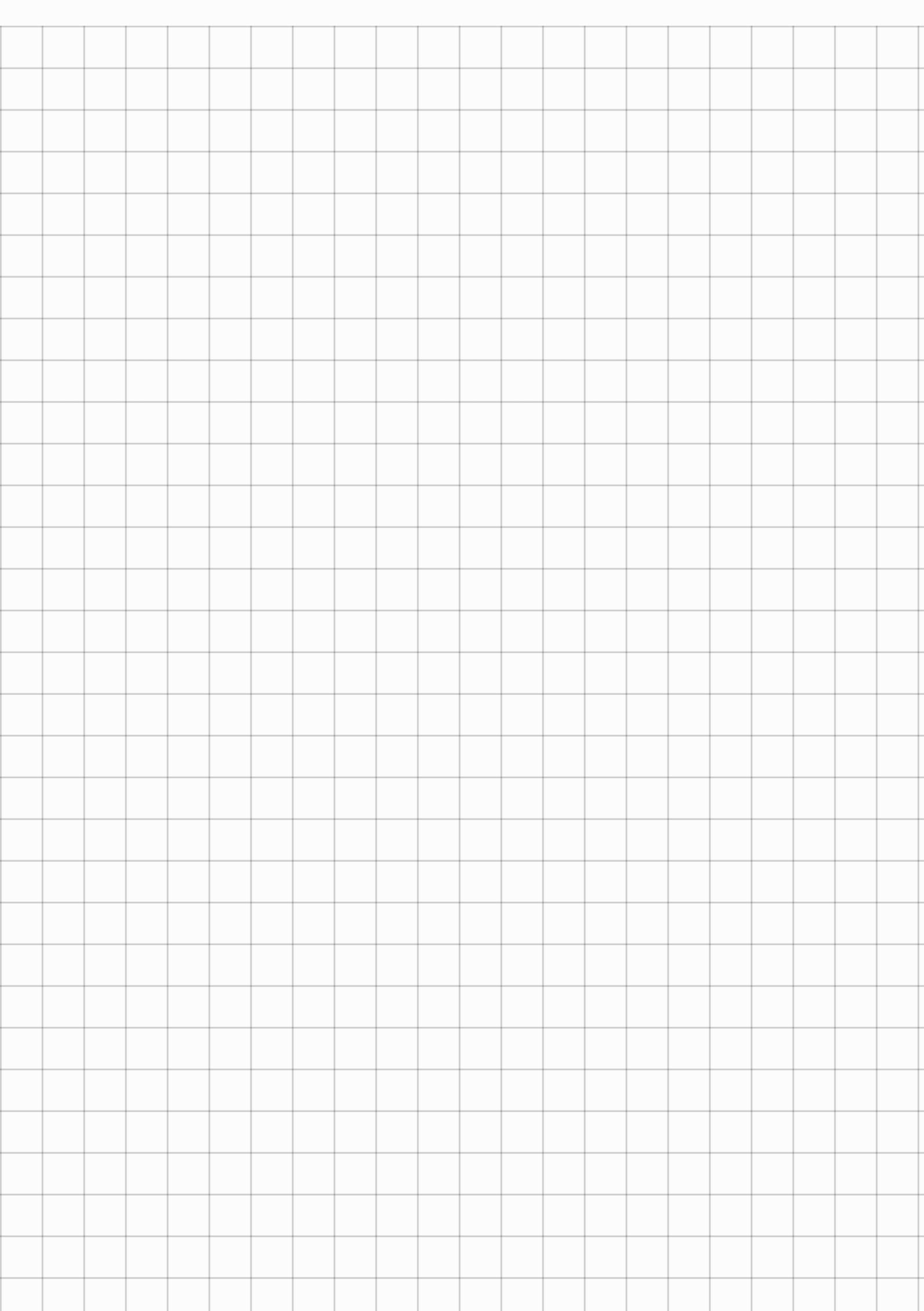
- Same as above, but instead of only issuing from head,
issue from anything which is ready
- actually this seems much more sensible since it's really
not much additional overhead over Design Z

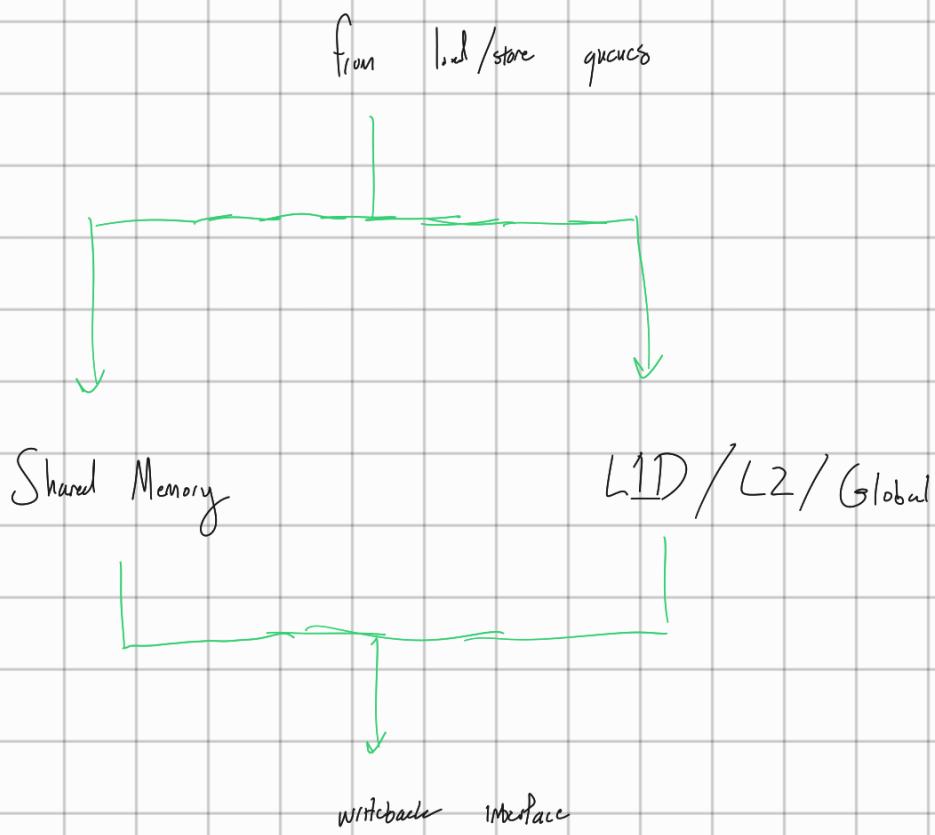
D_{CSIG} 3 : allows loads/stores to happen out of program order (full reordering)

- essentially a CPU load/store queue setup
- loads can be reordered w.r.t any other load, but not older stores to same address
 - (before starting a load, ensure no pending older store to same address)
- stores cannot be reordered w.r.t older stores to same address
 - or
 - (before starting a store, ensure no pending store to same address / load to same address)

* but we only enforce ordering within a single lane - if multiple threads non-atomically access same memory, it is very undefined ...

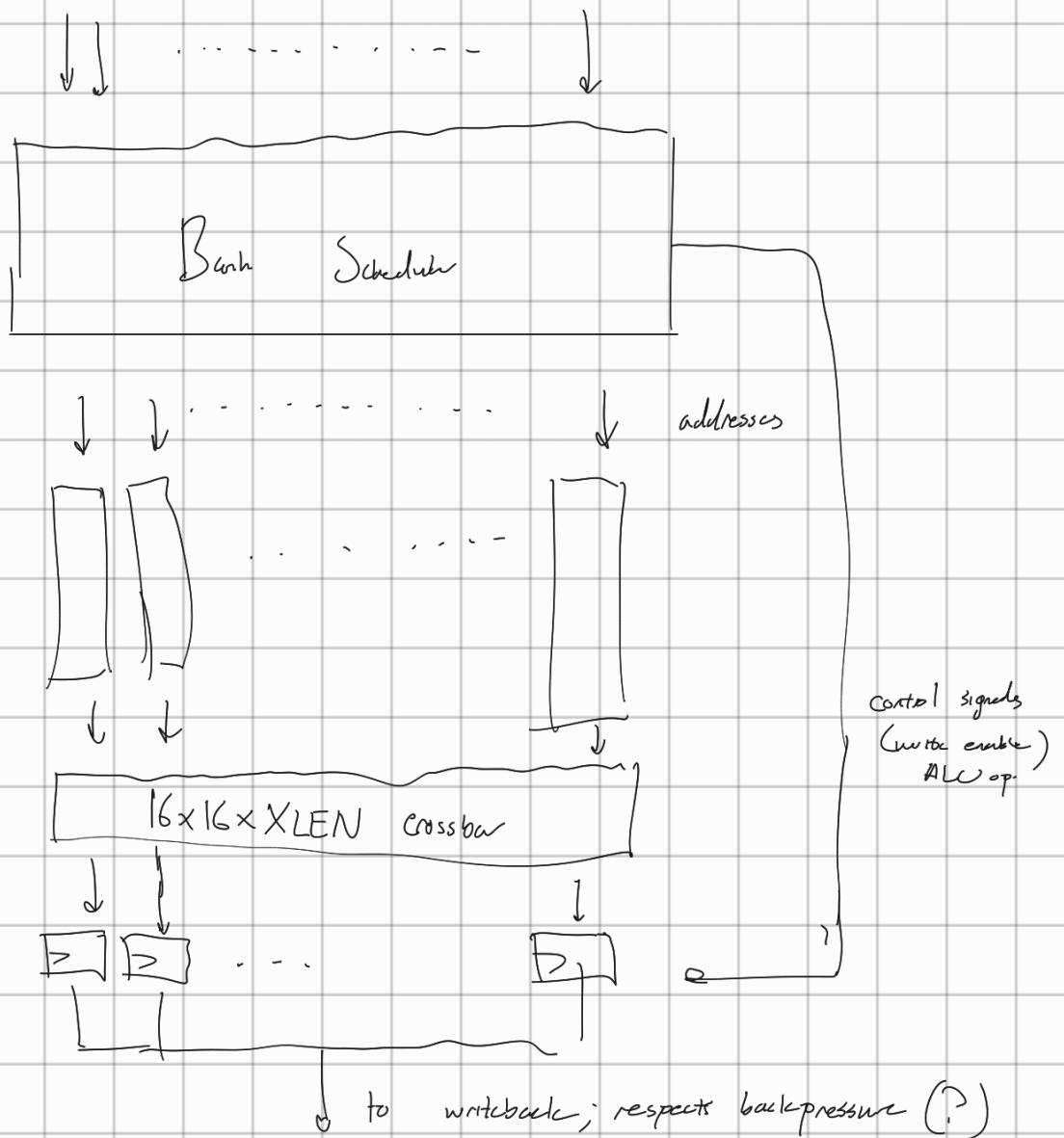
- i honestly think this is too expensive to do
 - alias detection for 16 threads or speculation/replay
 - both seem super expensive ...





- Note: not possible to do a mixed shared mem / global mem operation
 - ↳ they will be part of the same address space due to ISA limitation
(can we change this?)
- fault or undefined behavior if you try something weird
- Shared Mem lines in certain address range — switch between SMEM/GLOBAL w/ Z comparators
- Both can apply backpressure (e.g. scrubbing bank conflicts in SMEM)
 - to issuer, but busy shared MEM shouldn't affect global memory accesses

Shared Memory



If banks, word interleaved

(word 1 ∈ bank 1, word 2 ∈ bank 2,
word 16 ∈ bank 6, word 17 ∈ bank 1, ...)

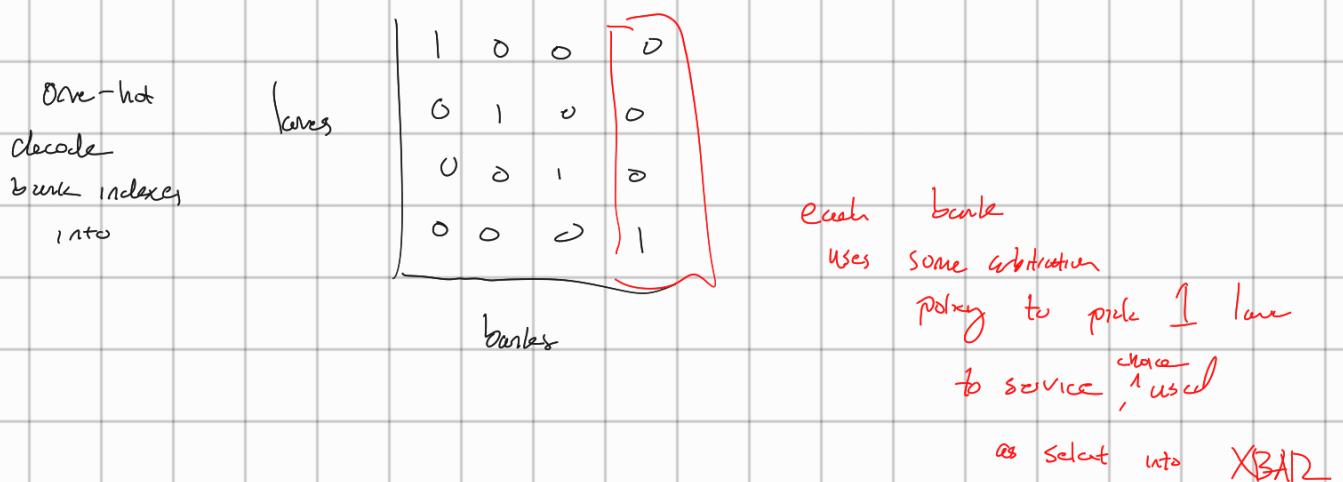
- Bank Scheduler responsible for serializing
bank conflicts

Stage 1. Decode + Broadcast Detection + Conflict Detection

LSBs = bank index (word-interleaving)

MSBs = address within bank (bank address)

- For reads, check if it's a broadcast
 $(address_1 = address_2 = \dots = address_{16})$
 - we can special case this to avoid very bad serialization
- Conflict detection / serialization



- = Serviced lanes tracked, continue serializing until compare equal to thread mask

Stage 2 . Read / Write Data

Atomic Support

TODO: how to efficiently support lr, sc?

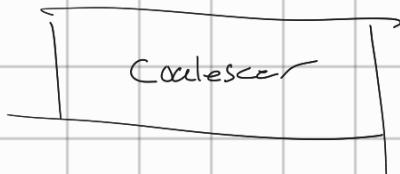
way too many logical threads to track
 $(\# \text{ nops} \times \# \text{ threads})$

do we actually need lr, sc?

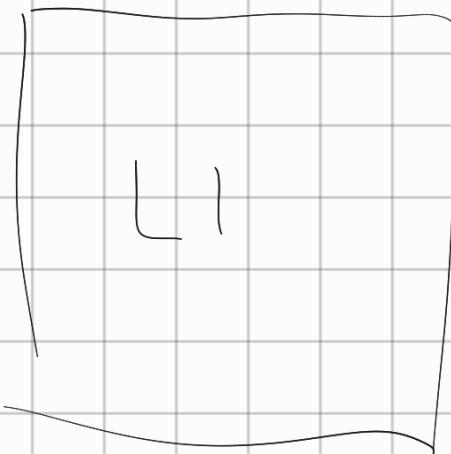
- AMO *
- $(+, -, \min, \max)$

is easy on the other hand - just pass in a
per-bank ALU and
have bank scheduler take some
extra cycles to read, modify, then write.

L1D / L2 / Global

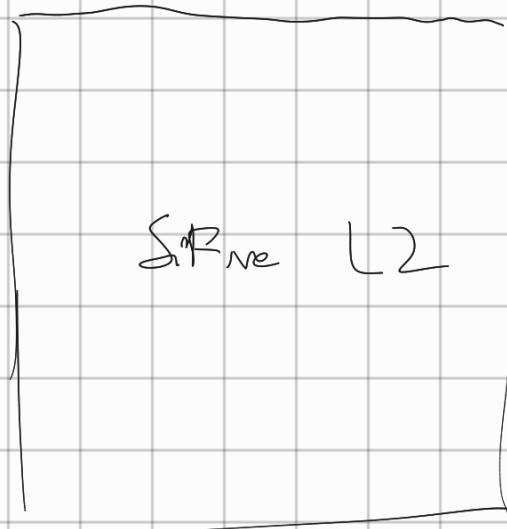


(?)
Coalescer: Can safely do
Spatial and temporal
coalescing, but
probably spatial
only for this
design



TODO: Can we
reuse Rocket Chip's L1?
Writing Cache is
hard

- probably want write-through
cache for
column reasons?



- need flush mechanism
for kernel boundary

L2