

Constructive Computer Architecture

Realistic Memories and Caches

Arvind
Computer Science & Artificial Intelligence Lab.
Massachusetts Institute of Technology

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-1

Contributors to the course material

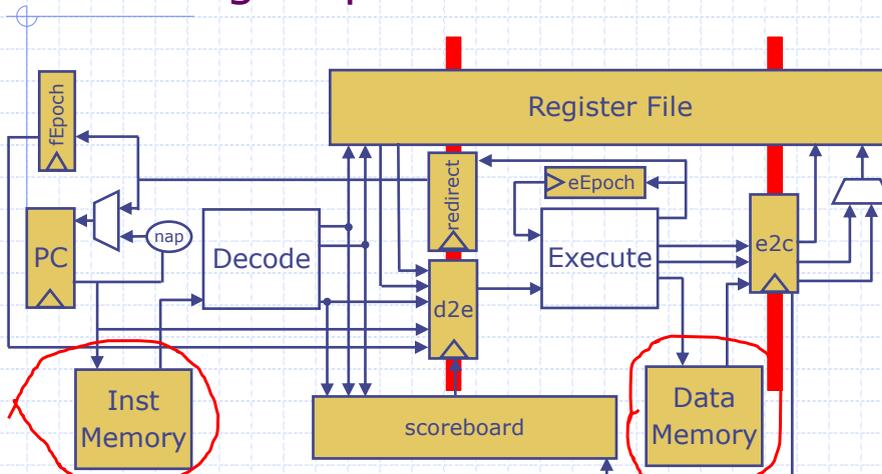
- ◆ Arvind, Rishiyur S. Nikhil, Joel Emer, Muralidaran Vijayaraghavan
- ◆ Staff and students in 6.375 (Spring 2013), 6.S195 (Fall 2012), 6.S078 (Spring 2012)
 - Asif Khan, Richard Ruhler, Sang Woo Jun, Abhinav Agarwal, Myron King, Kermin Fleming, Ming Liu, Li-Shiuan Peh
- ◆ External
 - Prof Amey Karkare & students at IIT Kanpur
 - Prof Jihong Kim & students at Seoul Nation University
 - Prof Derek Chiou, University of Texas at Austin
 - Prof Yoav Etsion & students at Technion

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-2

Multistage Pipeline



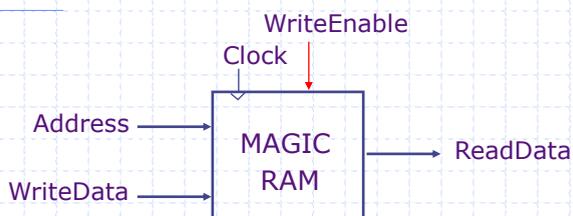
The use of magic memories (combinational reads) makes such design unrealistic

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-3

Magic Memory Model



- ◆ Reads and writes are always completed in one cycle
 - a Read can be done any time (i.e. combinational)
 - If enabled, a Write is performed at the rising clock edge
(the write address and data must be stable at the clock edge)

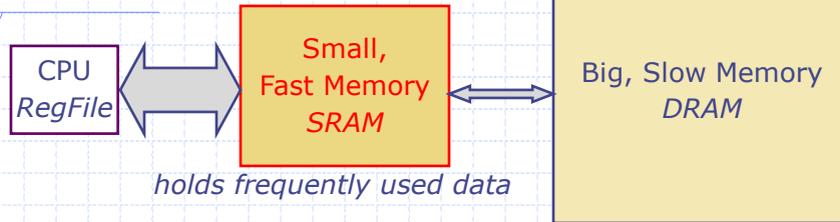
In a real DRAM the data will be available several cycles after the address is supplied

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-4

Memory Hierarchy



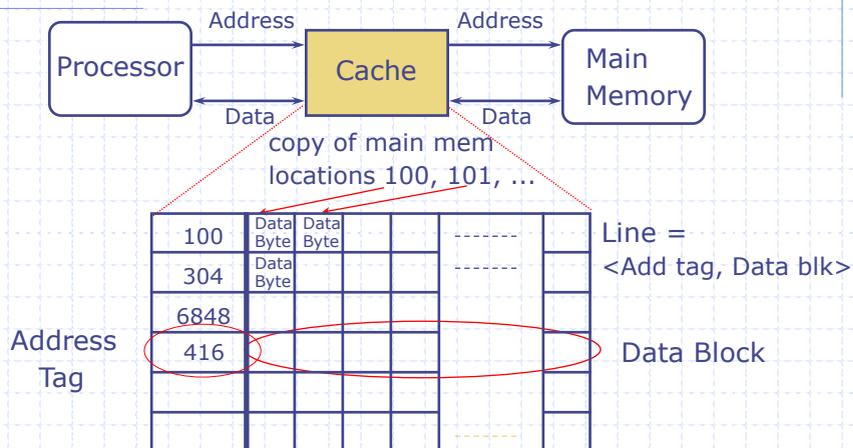
size: RegFile << SRAM << DRAM
 latency: RegFile << SRAM << DRAM
 bandwidth: on-chip >> off-chip

why?

On a data access:

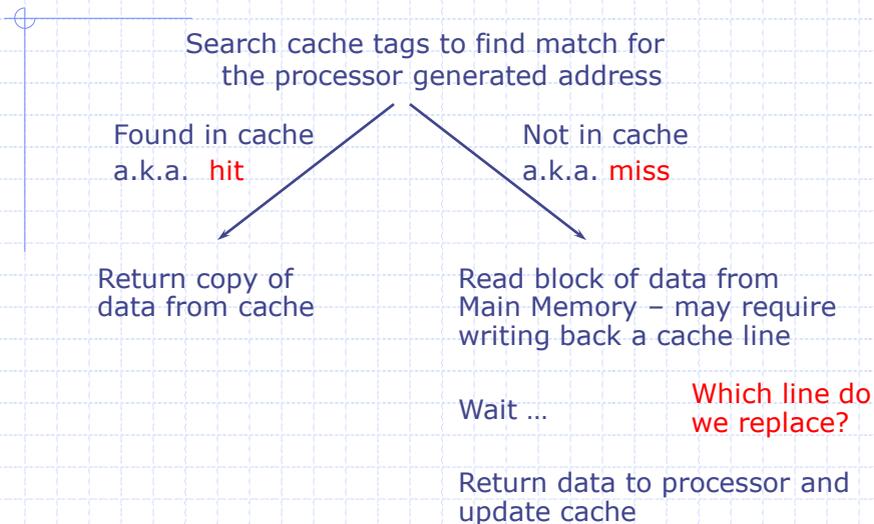
hit (data ∈ fast memory) ⇒ low latency access
miss (data ∉ fast memory) ⇒ long latency access (DRAM)

Inside a Cache



How many bits are needed for the tag?
 Enough to uniquely identify the block

Cache Read



October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-7

Write behavior

- ◆ On a write hit
 - Write-through: write to both cache and the next level memory
 - write-back: write only to cache and update the next level memory when line is evacuated
- ◆ On a write miss
 - Allocate – because of multi-word lines we first fetch the line, and then update a word in it
 - Not allocate – word modified in memory

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-8

Cache Line Size

- ◆ A cache line usually holds more than one word
 - Reduces the number of tags and the tag size needed to identify memory locations
 - Spatial locality: Experience shows that if address x is referenced then addresses $x+1$, $x+2$ etc. are very likely to be referenced in a short time window
 - ◆ consider instruction streams, array and record accesses
 - Communication systems (e.g., bus) are often more efficient in transporting larger data sets

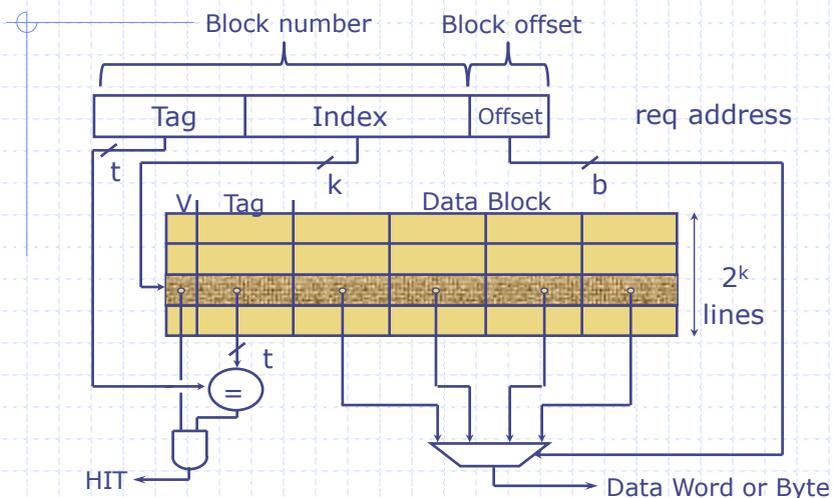
Types of misses

- ◆ Compulsory misses (cold start)
 - First time data is referenced
 - Run billions of instructions, become insignificant
- ◆ Capacity misses
 - Working set is larger than cache size
 - Solution: increase cache size
- ◆ Conflict misses
 - Usually multiple memory locations are mapped to the same cache location to simplify implementations
 - Thus it is possible that the designated cache location is full while there are empty locations in the cache.
 - Solution: Set-Associative Caches

Internal Cache Organization

- ◆ Cache designs restrict where in cache a particular address can reside
 - *Direct mapped:* An address can reside in exactly one location in the cache. The cache location is typically determined by the lowest order address bits
 - *n-way Set associative:* An address can reside in any of the a set of n locations in the cache. The set is typically determine by the lowest order address bits

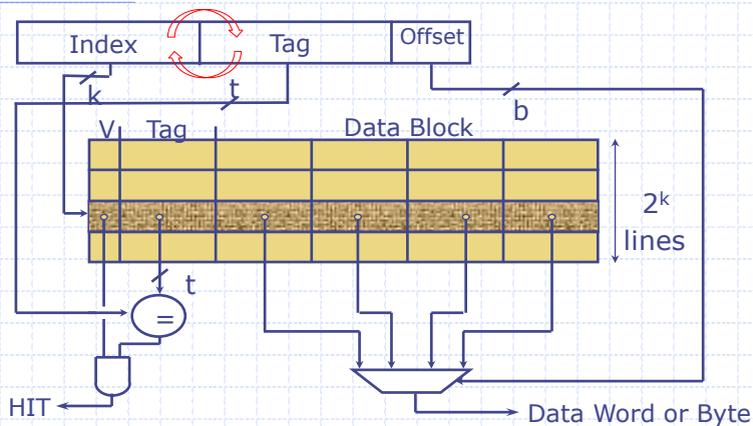
Direct-Mapped Cache



What is a bad reference pattern? Strided = size of cache

Direct Map Address Selection

higher-order vs. lower-order address bits



Why higher-order bits as tag may be undesirable?

Spatially local blocks conflict

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-13

Reduce Conflict Misses

Memory time =
Hit time + Prob(miss) * Miss penalty

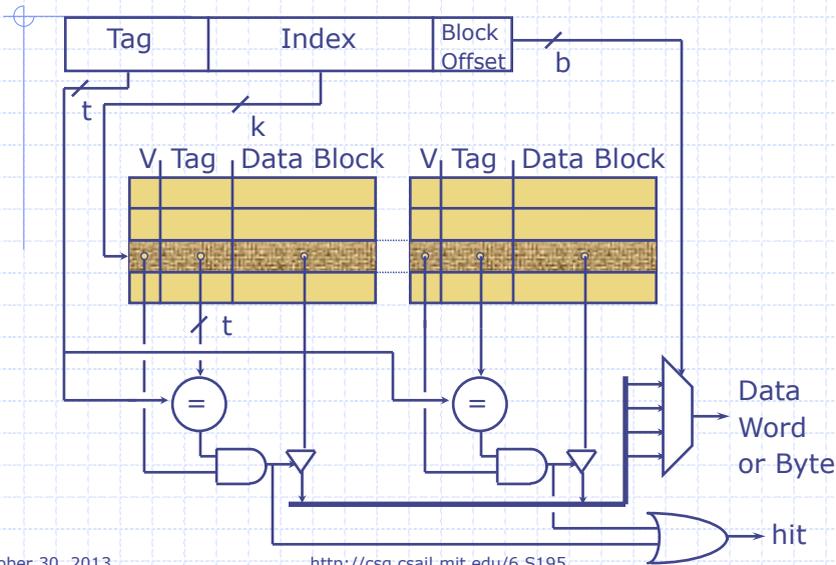
- ◆ **Associativity:** Reduce conflict misses by allowing blocks to go to several sets in cache
 - 2-way set associative: each block can be mapped to one of 2 cache sets
 - Fully associative: each block can be mapped to any cache frame

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-14

2-Way Set-Associative Cache



Replacement Policy

- ◆ In order to bring in a new cache line, usually another cache line has to be thrown out. Which one?
 - No choice in replacement if the cache is direct mapped
- ◆ Replacement policy for set-associative caches
 - One that is not dirty, i.e., has not been modified
 - ◆ In I-cache all lines are clean
 - ◆ In D-cache if a dirty line has to be thrown out then it must be written back first
 - Least recently used?
 - Most recently used?
 - Random?

How much is performance affected by the choice?

Difficult to know without benchmarks and quantitative measurements

Blocking vs. Non-Blocking cache

◆ Blocking cache:

- At most one outstanding miss
- Cache must wait for memory to respond
- Cache does not accept requests in the meantime

◆ Non-blocking cache:

- Multiple outstanding misses
- Cache can continue to process requests while waiting for memory to respond to misses

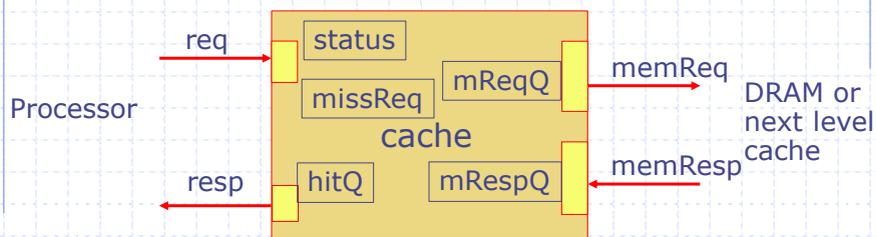
We will first design a write-back, No write-miss allocate, blocking cache

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-17

Blocking Cache Interface



```
interface Cache;  
  method Action req(MemReq r);  
  method ActionValue#(Data) resp;  
  
  method ActionValue#(MemReq) memReq;  
  method Action memResp(Line r);  
endinterface
```

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-18

Interface dynamics

- ◆ The cache either gets a hit and responds immediately, or it gets a miss, in which case it takes several steps to process the miss
- ◆ Reading the response dequeues it
- ◆ Requests and responses follow the FIFO order
- ◆ Methods are guarded, e.g., the cache may not be ready to accept a request because it is processing a miss
- ◆ A status register keeps track of the state of the cache while it is processing a miss

```
typedef enum {Ready, StartMiss, SendFillReq,  
             WaitFillResp} CacheStatus deriving (Bits, Eq);
```

Blocking Cache code structure

```
module mkCache (Cache);  
  RegFile# (CacheIndex, Line) dataArray <-  
    mkRegFileFull; ...  
  
  rule startMiss ... endrule;  
  
  method Action req (MemReq r) ...      endmethod;  
  method ActionValue# (Data) resp ...    endmethod;  
  
  method ActionValue# (MemReq) memReq ... endmethod;  
  method Action memResp (Line r) ...     endmethod;  
endmodule
```

- ◆ Internal communications is in line sizes but the processor interface, e.g., the response from the hitQ is word size

Blocking cache state elements

```
RegFile#(CacheIndex, Line) dataArray <- mkRegFileFull;
RegFile#(CacheIndex, Maybe#(CacheTag))
    tagArray <- mkRegFileFull;
RegFile#(CacheIndex, Bool) dirtyArray <- mkRegFileFull;

Fifo#(1, Data) hitQ <- mkBypassFifo;
Reg#(MemReq) missReq <- mkRegU;
Reg#(CacheStatus) status <- mkReg(Ready);

Fifo#(2, MemReq) memReqQ <- mkCFFifo;
Fifo#(2, Line) memRespQ <- mkCFFifo;

function CacheIndex getIdx(Addr addr) = truncate(addr>>2);
function CacheTag getTag(Addr addr) = truncateLSB(addr);
```

Tag and valid bits are kept together as a Maybe type

CF Fifos are preferable because they provide better decoupling. An extra cycle here may not affect the performance by much

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-21

Req method hit processing

```
method Action req(MemReq r) if(status == Ready);
    let idx = getIdx(r.addr); let tag = getTag(r.addr);
    let currTag = tagArray.sub(idx);
    let hit = isValid(currTag)?
        fromMaybe(?,currTag)==tag : False;
if(r.op == Ld) begin
    if(hit) hitQ.enq(dataArray.sub(idx));
    else begin missReq <= r; status <= StartMiss; end
    end
else begin // It is a store request
    if(hit) begin dataArray.upd(idx, r.data);
        dirtyArray.upd(idx, True); end
    else memReqQ.enq(r); // write-miss no allocate
    end
endmethod
```

It is straightforward to extend the cache interface to include a cacheline flush command

In case of multiword cache line, we only overwrite the appropriate word of the line

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-22

Rest of the methods

```
method ActionValue#(Data) resp;  
  hitQ.deq;  
  return hitQ.first;  
endmethod
```

```
method ActionValue#(MemReq) memReq;  
  memReqQ.deq;  
  return memReqQ.first;  
endmethod
```

```
method Action memResp(Line r);  
  memRespQ.enq(r);  
endmethod
```

Memory side
methods

Start miss rule

```
Ready -> StartMiss -> SendFillReq -> WaitFillResp -> Ready
```

```
rule startMiss(status == StartMiss);  
  let idx = getIdx(missReq.addr);  
  let tag = tagArray.sub(idx);  
  let dirty = dirtyArray.sub(idx);  
  if(isValid(tag) && dirty) begin // write-back  
    let addr = {fromMaybe(?, tag), idx, 2'b0};  
    let data = dataArray.sub(idx);  
    memReqQ.enq(MemReq{op: St, addr: addr, data: data});  
  end  
  status <= SendFillReq;  
endrule
```

Send-fill and Wait-fill rules

```
Ready -> StartMiss -> SendFillReq -> WaitFillResp -> Ready
```

```
rule sendFillReq (status == SendFillReq);  
    memReqQ.enq(missReq);    status <= WaitFillResp;  
endrule
```

```
Ready -> StartMiss -> SendFillReq -> WaitFillResp -> Ready
```

```
rule waitFillResp(status == WaitFillResp);  
    let idx = getIdx(missReq.addr);  
    let tag = getTag(missReq.addr);  
    let data = memRespQ.first;  
    dataArray.upd(idx, data);  
    tagArray.upd(idx, Valid (tag));  
    dirtyArray.upd(idx, False);  
    hitQ.enq(data); memRespQ.deq;  
    status <= Ready;
```

```
endrule
```

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-25

Hit and miss performance

◆ Hit

- Combinational read/write, i.e. 0-cycle response
- Requires req and resp methods to be concurrently schedulable, which in turn requires

$\text{hitQ.enq} < \{\text{hitQ.deq}, \text{hitQ.first}\}$

i.e., hitQ should be a bypass Fifo

◆ Miss

- No evacuation: memory load latency plus combinational read/write
- Evacuation: memory store followed by memory load latency plus combinational read/write

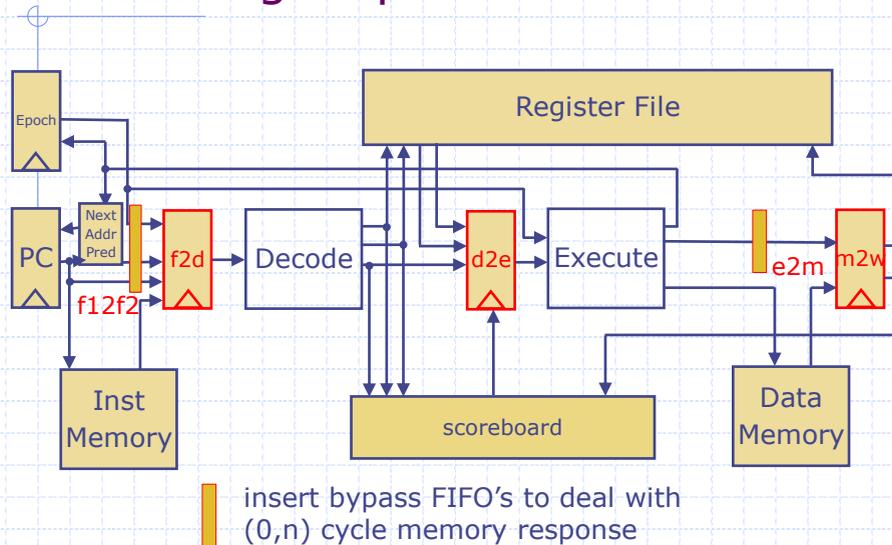
Adding an extra cycle here and there in the miss case should not have a big negative performance impact

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-26

Four-Stage Pipeline



October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-27

now some coding ...

- ◆ Integrating caches in to the 4-stage pipeline (F, D&R, E&M, W) from the last lecture
 - ◆ Direction predictor training is incompletely specified
- ◆ In L13 we discussed splitting a pipeline stage into two stages by inserting a bypass FIFO. We show it again here

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-28

4-Stage pipeline with BTB+BHT *without caches*

```
module mkProc(Proc);
  Reg#(Addr)      pc <- mkRegU;
  RFile          rf <- mkBypassRFile;
IMemory         iMem <- mkIMemory;
DMemory         dMem <- mkDMemory;
  Fifo#(1, Decode2Execute) d2e <- mkPipelineFifo;
  Fifo#(1, Exec2Commit)    e2c <- mkPipelineFifo;
  Scoreboard#(2) sb <- mkPipelineScoreboard;
  Reg#(Bool)      feEp <- mkReg(False);
  Reg#(Bool)      fdEp <- mkReg(False);
  Reg#(Bool)      dEp <- mkReg(False);
  Reg#(Bool)      deEp <- mkReg(False);
  Reg#(Bool)      eEp <- mkReg(False);
  Fifo#(ExecRedirect) redirect <- mkBypassFifo;
  Fifo#(DecRedirect) decRedirect <- mkBypassFifo;
  NextAddrPred#(16) nap <- mkBTB;
  DirPred#(1024) dirPred <- mkBHT;
```

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-29

4-Stage pipeline with BTB+BHT *with caches*

```
module mkProc(Proc);
  Reg#(Addr) pc <- mkRegU; Rfile rf <- mkBypassRFile;
  Cache#(ICacheSize) iCache <- mkCache;
  Cache#(DCacheSize) dCache <- mkCache;
  Fifo#(1, Fetch2Decode) f12f2 <- mkBypassFifo;
  Fifo#(1, Maybe#{Einst}) e2m <- mkBypassFifo;
  Fifo#(1, Decode2Execute) d2e <- mkPipelineFifo;
  Fifo#(1, Exec2Commit)    m2w <- mkPipelineFifo;
  Scoreboard#(2) sb <- mkPipelineScoreboard;
  Reg#(Bool)      feEp <- mkReg(False);
  Reg#(Bool)      fdEp <- mkReg(False);
  Reg#(Bool)      dEp <- mkReg(False);
  Reg#(Bool)      deEp <- mkReg(False);
  Reg#(Bool)      eEp <- mkReg(False);
  Fifo#(ExecRedirect) execRedirect <- mkBypassFifo;
  Fifo#(DecRedirect) decRedirect <- mkBypassFifo;
  AddrPred#(16) addrPred <- mkBTB;
  DirPred#(1024) dirPred <- mkBHT;
```

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-30

4-Stage pipeline with BTB+BHT *without caches*

```

rule doFetch;
  let inst = iMem.req(pc);
  if(redirect.notEmpty)
    begin redirect.deq; nap.update(redirect.first); end
  if(redirect.notEmpty && redirect.first.mispredict)
    begin pc <= redirect.first.nextPc; feEp <= !feEp; end
  else if(decRedirect.notEmpty) begin
    if(decRedirect.first.eEp == feEp) begin
      fdEp <= !fdEp; pc <= decRedirect.first.nextPc; end
    decRedirect.deq; end;
  else begin
    let ppc = nap.predPc(pc);
    f2d.enq(Fetch2Decode{pc: pc, ppc: ppc, inst: inst,
      eEp: feEp, dEp: fdEp});
  end
endrule

```

make a iMem request
and enqueue into f12f2

where?

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

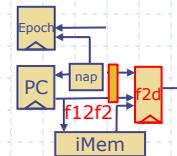
L17-31

4-Stage pipeline with BTB+BHT *with caches*

```

rule doFetch1;
  if(redirect.notEmpty)
    begin redirect.deq; nap.update(redirect.first); end
  if(redirect.notEmpty && redirect.first.mispredict)
    begin pc <= redirect.first.nextPc; feEp <= !feEp; end
  else if(decRedirect.notEmpty) begin
    if(decRedirect.first.eEp == feEp) begin
      fdEp <= !fdEp; pc <= decRedirect.first.nextPc; end
    decRedirect.deq; end;
  else begin
    let ppc = nap.predPc(pc);
    iCache.req(MemReq{op: Ld, addr: pc, data:??});
    f12f2.enq(Fetch2Decode{pc: pc, ppc: ppc, inst: ?,
      eEp: feEp, dEp: fdEp});
  end
endrule

```

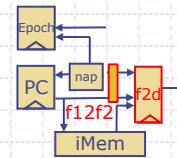


October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-32

4-Stage pipeline Fetch2 rule



```

rule doFetch2;
  let inst <- iCache.resp;
  let f2dVal = f12f2.first;
  f2dVal.inst = inst;
  f12f2.deq;
  f2d.enq(f2dVal);
endrule

```

Execute rule can be split in two rules in a similar manner to deal with dCache

next store buffers and non-blocking caches

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-33

4-Stage pipeline Execute rule

```

rule doExecute;
  let x = d2e.first;
  let dInst = x.dInst; let pc = x.pc;
  let ppc = x.ppc; let epoch = x.epoch;
  let rVal1 = x.rVal1; let rVal2 = x.rVal2;
  if(epoch == eEpoch) begin
    let eInst = exec(dInst, rVal1, rVal2, pc, ppc);
    if(eInst.iType == Ld) eInst.data <-
      dMem.req(MemReq{op:Ld, addr:eInst.addr, data:?});
    else if (eInst.iType == St) let d <-
      dMem.req(MemReq{op:St, addr:eInst.addr, data:eInst.data});
    e2c.enq(Exec2Commit{dst:eInst.dst, data:eInst.data});
    if(eInst.mispredict) eEpoch <= !eEpoch
    if(eInst.iType == J || eInst.iType == Jr || eInst.iType == Br)
      redirect.enq(Redirect{pc: pc, nextPc: eInst.addr,
        taken: eInst.brTaken, mispredict: eInst.mispredict,
        brType: eInst.iType}); end
    else e2c.enq(Exec2Commit{dst:Invalid, data:?});
    d2e.deq;
  endrule

```

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-34

4-Stage pipeline Execute rule with caches

```
rule doExecute1;
  let x = d2e.first;
  let dInst = x.dInst; let pc = x.pc;
  let ppc = x.ppc; let epoch = x.epoch;
  let rVal1 = x.rVal1; let rVal2 = x.rVal2;
  if(epoch == eEpoch) begin
    let eInst = exec(dInst, rVal1, rVal2, pc, ppc);
    if(eInst.iType == Ld) eInst.data <-
      dCache.req(MemReq{op:Ld, addr:eInst.addr, data:?});
    else if (eInst.iType == St) let d <-
      dCache.req(MemReq{op:St, addr:eInst.addr, data:eInst.data});
    e2m.eng(Valid (eInst));
    if(eInst.mispredict) eEpoch <= !eEpoch
    if(eInst.iType == J || eInst.iType == Jr || eInst.iType == Br)
      redirect.eng(Redirect{pc: pc, nextPc: eInst.addr,
        taken: eInst.brTaken, mispredict: eInst.mispredict,
        brType: eInst.iType}); end
    else e2m.eng(Invalid);
    d2e.deq;
  endrule
```

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-35

4-Stage pipeline Execute2 rule

```
rule doExecute2;
  let eInst = e2m.first;
  if(isValid(eInst)) begin
    let x = validValue(eInst);
    if(x.iType == Ld)
      x.data <- dCache.resp;
    m2w.eng(Exec2Commit{dst:x.dst, data:x.data});
  end
  else
    m2w.eng(Exec2Commit{dst:Invalid, data:?});
  e2m.deq;
endrule
```

October 30, 2013

<http://csg.csail.mit.edu/6.S195>

L17-36