

# Modular Refinement

Arvind

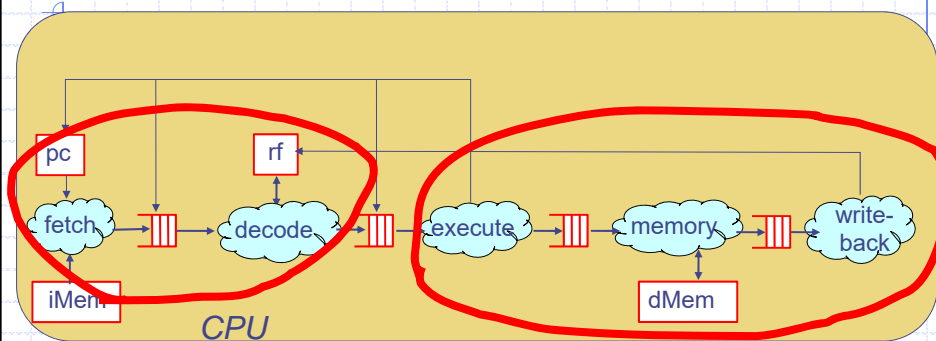
Computer Science & Artificial Intelligence Lab  
Massachusetts Institute of Technology

March 7, 2016

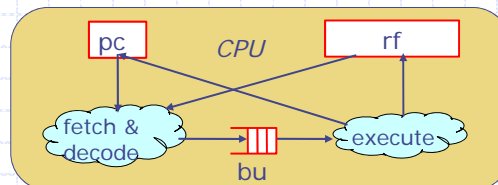
<http://csg.csail.mit.edu/6.375>

L11-1

## Successive refinement & Modular Structure



Can we derive a multi-stage pipeline by successive refinement of a 2-stage pipeline?



March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-2

## Architectural refinements

- ◆ Latency-insensitive communication between modules for better decoupling
- ◆ Refine each module independently
  - Separate Fetch and Decode
  - Replace magic memory by multicycle memory
  - Introduce Multicycle functional units (MUL, DIV)
  - ...
- ◆ A refined module should be able to work with other unrefined or refined modules

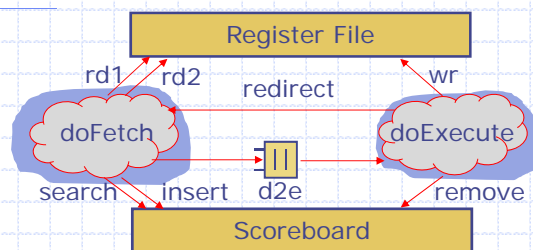
Nirav Dave, M.C. Ng, M. Pellauer, Arvind [Memocode 2010]  
A design flow based on modular refinement

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-3

## A 2-stage Processor Pipeline



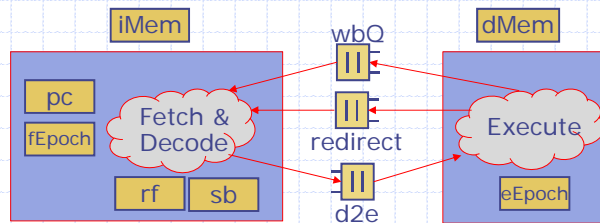
- ◆ Suppose we encapsulate each rule into its own module
- ◆ Communication from fetch to Execute via d2e is latency insensitive but not the redirection from Execute to Fetch
- ◆ Interactions between reads and writes to the register file or the interactions between search and remove to the scoreboard are also not latency insensitive

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-4

# Making the design latency Insensitive



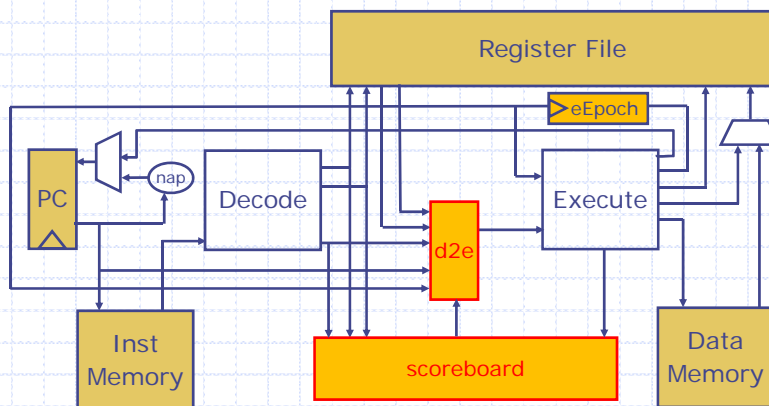
- ◆ Instead of directly resetting the pc and epoch in case of a redirection, let the Execute send the relevant information to Fetch via the redirect fifo
- ◆ Include rf and sb in the Fetch module and let the Execute send the information to update rf and sb to Fetch via the wbQ

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-5

# Tightly coupled 2-Stage-pipeline: Scoreboard and Stall logic



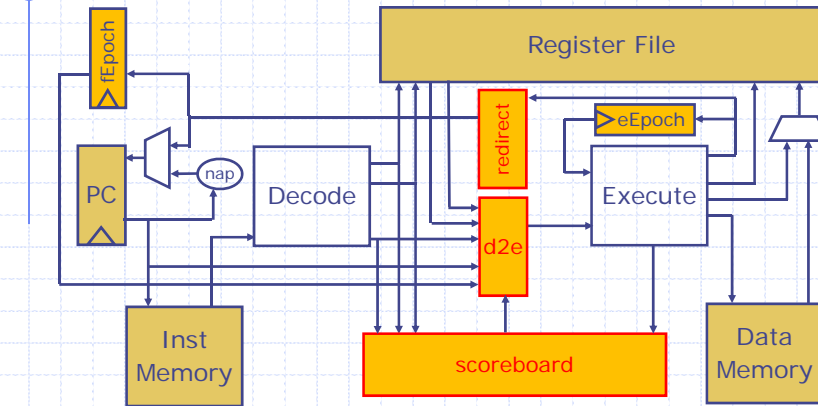
From lecture L10

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-6

## Making redirection latency-insensitive



- No direct access to PC from Execute;
- Fetch accesses its own epoch which is kept consistent with the epoch in the Execute lazily

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-7

## A decoupled solution using epochs

Fetch fEpoch | eEpoch Execute

- ◆ Add fEpoch and eEpoch registers to the processor state; initialize them to the same value
- ◆ The epoch changes whenever Execute detects the pc prediction to be wrong. This change is reflected immediately in eEpoch and eventually in fEpoch via a message from Execute to Fetch
- ◆ Associate fEpoch with every instruction when it is fetched
- ◆ In the execute stage, reject, i.e., kill, the instruction if its epoch does not match eEpoch

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-8

## Fixing the doExecute rule for latency-independent pc update

```
rule doExecute;
  let x = d2e.first;
  let dInstE = x.dInst; let pcE = x.pc;
  let ppcE = x.ppc; let inEpoch = x.epoch;
  let rVal1E = x.rVal1; let rVal2E = x.rVal2;
  if(inEpoch == epoch) begin
    let eInst = exec(dInstE, rVal1E, rVal2E, pcE, ppcE);
    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});
    if (isValid(eInst.dst))
      rf.wr(fromMaybe(?), eInst.dst, eInst.data);
    if(eInst.mispredict) begin
      pc[1] <= eInst.addr; epoch <= !epoch; end
    end
    d2e.deq; sb.remove;
  endrule
```

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-9

## Fixing the doFetch rule

```
rule doFetch;

  let instF = iMem.req(pc[0]);
  let ppcF = nap(pc[0]);
  let dInst = decode(instF);
  let stall = sb.search1(dInst.src1) || sb.search2(dInst.src2);
  if(!stall) begin
    let rVal1 = rf.rd1(fromMaybe(?), dInst.src1);
    let rVal2 = rf.rd2(fromMaybe(?), dInst.src2);
    d2e.enq(Decode2Execute{pc: pc[0], ppc: ppcF,
      dInst: dInst, epoch: epoch,
      rVal1: rVal1, rVal2: rVal2});
    sb.insert(dInst.rDst); pc[0] <= ppcF; end
  endrule
```

-pc is no longer an EHR

-redirected pc value must be read from the redirect FIFO

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-10

## Fetch rules for latency independent pc update,

```
rule redirectPC;
  fEpoch <= !fEpoch; pc <= redirect.first;
  redirect.deq;
endrule

rule doFetch (!redirect.notEmpty);
  let instF = iMem.req(pc);
  let ppcF = nap(pc); let dInst = decode(instF);
  let stall = sb.search1(dInst.src1) || sb.search2(dInst.src2);
  if(!stall) begin
    let rVal1 = rf.rd1(fromMaybe(? , dInst.src1));
    let rVal2 = rf.rd2(fromMaybe(? , dInst.src2));
    d2e.eng(Decode2Execute{pc: pc, ppc: ppcF,
      dInst: dInst, epoch: fEpoch,
      rVal1: rVal1, rVal2: rVal2});
    sb.insert(dInst.rDst); pc <= ppcF; end
  end
endrule
```

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-11

## Fixing the doExecute rule for latency-independent rf update

```
rule doExecute;
  let x = d2e.first;
  ...
  ...
  ...
  if(inEpoch == eEpoch) begin
    let eInst = exec(dInstE, rVal1E, rVal2E, pcE, ppcE);
    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});
    if (isValid(eInst.dst))
      rf.wr(fromMaybe(? , eInst.dst), eInst.data);
    if(eInst.mispredict) begin
      redirect.eng(eInst.addr); eEpoch <= !eEpoch; end
    end
    d2e.deq; sb.remove;
  end
endrule
```

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-12

## Latency-insensitive doExecute

```
rule doExecute;
  let x = d2e.first;
  let dInstE = x.dInst; let pcE      = x.pc;
  let ppcE      = x.ppc; let inEpoch = x.epoch;
  let rVal1E = x.rVal1; let rVal2E = x.rVal2;
  if(inEpoch == eEpoch) begin
    let eInst = exec(dInstE, rVal1E, rVal2E, pcE, ppcE);
    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});
    wbQ.enq(eInst.dst, eInst.data);
    if(eInst.mispredict) begin
      redirect.enq(eInst.addr); eEpoch <= !eEpoch; end
    end
  else wbQ.enq(tuple2(Invalid, ?));
  d2e.deq;
endrule
```

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-13

## We need to introduce a WB rule in the Fetch stage

```
rule writeback;
  match {.idx, .val} = wbQ.first;
  if(isValid(idx)) rf.write(fromMaybe(?, idx), val);
  wbQ.deq; sb.remove;
endrule
```

- ◆ Are data hazards being handled properly?
- ◆ That is, can stale values be read from the rf?

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-14

# Modular Processor

```
module mkModularProc(iMem, dMem);
  Fifo#(Decode2Execute) d2e <- mkPipelineFifo;
  Fifo#(Addr)  redirect <- mkBypassFifo;
  Fifo#(Tuple2#(Maybe#(Indx),Data) wbQ<-mkBypassFifo;
  Fetch      fetch <- mkFetch(iMem, d2e, redirect, wbQ);
  Execute execute <- mkExecute(dMem, d2e, redirect, wbQ);
endmodule
```

We only need to pass the enq interfaces to execute

Similarly we need to pass only deq and first interfaces to Fetch

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-15

# Fetch Module

```
module mkFetch(iMem, d2e, redirect, wbQ); types not shown
  Make pc, fEpoch, rf, sb;
rule writeback;
  match {.idx, .val} = wbQ.first;
  if(isValid(idx)) rf.write(fromMaybe(?, idx), val);
  wbQ.deq; sb.remove;
endrule
writeback and redirectPC are CF;
doFetch and redirectPC are ME

rule redirectPC ;
  fEpoch <= !fEpoch; pc <= redirect.first;
  redirect.deq;
endrule;
interaction between writeback and doFetch depends
upon what assumptions are made about rf and sb

rule doFetch
...
  bypass rf and pipeline sb ==> writeback < doFetch
  normal rf and CF sb ==> doFetch < writeback
endrule
```

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-16



## Execute Module

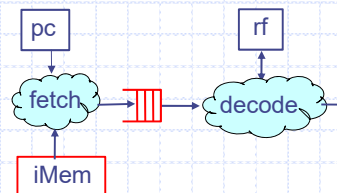
```
module mkExecute(dMem, d2e, redirect, wbQ) types not shown
  Make eEpoch;
rule doExecute;
  let x=d2e.first; let dInst=x.dInst; let pc=x.pc; let ppc=x.ppc;
  let epoch = x.epoch; let rVal1 = x.rVal1; let rVal2 = x.rVal2;
  if(epoch == eEpoch) begin
    let eInst = exec(dInst, rVal1, rVal2, pc, ppc);
    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});
    wbQ.enq(tuple2(eInst.dst, eInst.data));
    if(eInst.mispredict) begin
      execRedirect.enq(eInst.addr); eEpoch <= !eEpoch; end end
    else wbQ.enq(tuple2(Invalid, ?));
  d2e.deq; endrule
```

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-17

## Modular refinement: Separating Fetch and Decode



March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-18

## Fetch Module refinement

```
module mkFetch(iMem,d2e,redirect,wbQ);
  Reg#(Addr)      pc <- mkRegU;
  Reg#(Bool)      fEpoch <- mkReg(False);
  RFile           rf <- mkBypassRFile;
  Scoreboard#(1) sb <- mkPipelineScoreboard;
  Fifo#(Fetch2Decode) f2d <- mkPipelineFifo;

  rule writeback; ... endrule
  rule redirectPC; ... endrule

  rule fetch ; .... endrule
  rule decode ; .... endrule
```

doFetch rule will be split into fetch and decode rules

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-19

## Fetch Module: Fetch rule

```
rule fetch ;
  let inst = iMem.req(pc);
  let ppc = nextAddrPredictor(pc);
  f2d.enq(Fetch2Decode{pc: pc, ppc: ppc,
                      inst: inst, epoch: fEpoch});
  pc <= ppc
endrule
```

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-20

## Fetch Module: Decode rule

```
rule decode ;
  let x = f2d.first;
  let inst = x.inst; let inPC = x.pc; let ppc = x.ppc
  let inEp = x.epoch
  let dInst = decode(inst);
  let stall = sb.search1(dInst.src1) || sb.search2(dInst.src2);
  if(!stall) begin
    let rVal1 = rf.rd1(validRegValue(dInst.src1));
    let rVal2 = rf.rd2(validRegValue(dInst.src2));
    d2e.enq(Decode2Execute{pc: inPC, ppc: ppc,
      dInst: dInst, epoch: inEp;
      rVal1: rVal1, rVal2: rVal2});
    sb.insert(dInst.dst);
  f2d.deq end
endrule
```

with pipelined f2d decode < fetch

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-21

## Separate refinement

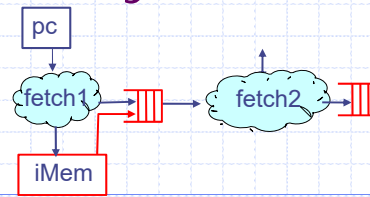
- ◆ Notice our refine Fetch&Decode module should work correctly with the old Execute module or its refinements
- ◆ This is a very important aspect of modular refinements

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-22

## Modular refinement: Replace magic memory by multicycle memory



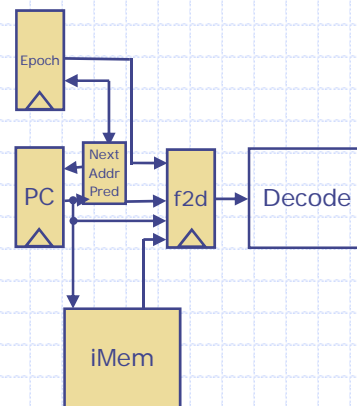
March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-23

## Memory and Caches

- ◆ Suppose iMem is replaced by a cache which takes one cycle in case of a hit and unknown number of variable cycles on a cache miss
- ◆ View iMem as a request/response system and split the fetch stage into two rules – to send a req and to receive a res



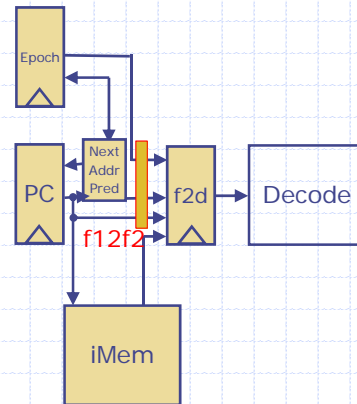
March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-24

## Splitting the fetch stage

- ◆ To split the fetch stage into two rules, insert a FIFO to deal with multicycle memory responses



March 7, 2016

<http://csg.csail.mit.edu/6.375>

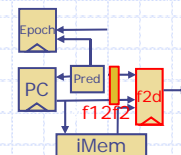
L11-25

## Fetch Module: 2nd refinement

```

module mkFetch(iMem,d2e,redirect,wbQ);
    Reg#(Addr)      pc <- mkRegU;
    Reg#(Bool)      fEpoch <- mkReg(False);
    RFile           rf <- mkBypassRFile;
    Scoreboard#(1) sb <- mkPipelineScoreboard;
    Fifo#(Fetch2Decode) f2d <- mkPipelineFifo;
    Fifo#(Fetch2Decode) f12f2 <- mkPipelineFifo;

    rule writeback; ... endrule
    rule redirectPC; ... endrule
    rule fetch1 ; .... endrule
    rule fetch2 ; ..... endrule
    rule decode ; ..... endrule
  
```

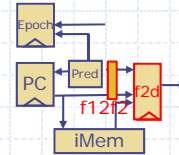


March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-26

## Fetch Rules



```
rule fetch1 ;
  let ppc = nap(pc); pc <= ppc;
  iCache.req(MemReq{op: Ld, addr: pc, data:?});
  f12f2.enq(Fetch2Decode{pc: pc, ppc: ppc,
                        inst: ?, epoch: fEpoch});
endrule
```

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

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-27

## Takeaway

- ◆ Modular refinement is a powerful idea; lets different teams work on different modules with only an early implementation of other modules
- ◆ BSV compiler currently does not permit separate compilation of modules with interface parameters
- ◆ Recursive call structure amongst modules is supported by the compiler in a limited way.
  - The syntax is complicated
  - Compiler detects and rejects recursive call structures

March 7, 2016

<http://csg.csail.mit.edu/6.375>

L11-28

## Interface issues

- ◆ For better safety only partial interfaces should to be passed to a module, e.g.,

- Fetch module needs only `deq` and `first` methods of `execRedirect` and `wbQ`, and `enq` method of `d2e`

```
interface FifoEnq#(t); method Action enq(t x); endinterface
interface FifoDeq#(t); method Action deq; method t first;
endinterface

function FifoEnq#(t) getFifoEnq(Fifo#(n, t) f); return
  interface FifoEnq#(t); method Action enq(t x) = f.enq(x);
endinterface endfunction

function FifoDeq#(t) getFifoDeq(Fifo#(n, t) f); return
  interface FifoDeq#(t); method Action deq=f.deq;
  method t first=f.first; endinterface endfunction

module mkModularProc(Proc);
  Fetch fetch <- mkFetch(iMem, getFifoEnq(d2e),
    getFifoDeq(execRedirect), getFifoDeq(wbQ));
```