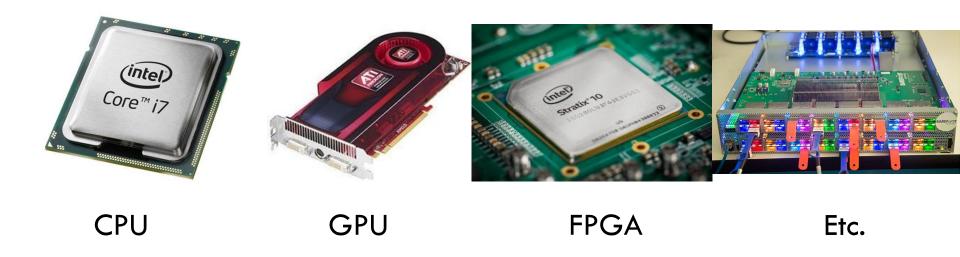
Lecture 2:

Instruction-Level Parallelism

15-418 Parallel Computer Architecture and Programming CMU 15-418/15-618, Spring 2021

Many kinds of processors



Why so many? What differentiates these processors?

Why so many kinds of processors?

Each processor is designed for different kinds of programs

- CPUs
 - "Sequential" code i.e., single / few threads
- GPUs
 - Programs with lots of independent work "Embarrassingly parallel"

Many others: Deep neural networks, Digital signal processing, Etc.

Parallelism pervades architecture

- Speeding up programs is all about parallelism
 - 1) Find independent work
 - 2) Execute it in parallel
 - 3) Profit

- Key questions:
 - Where is the parallelism?
 - Whose job is it to find parallelism?

Where is the parallelism?

Different processors take radically different approaches

- CPUs: Instruction-level parallelism
 - Implicit
 - Fine-grain
- GPUs: Thread- & data-level parallelism
 - Explicit
 - Coarse-grain

Whose job to find parallelism?

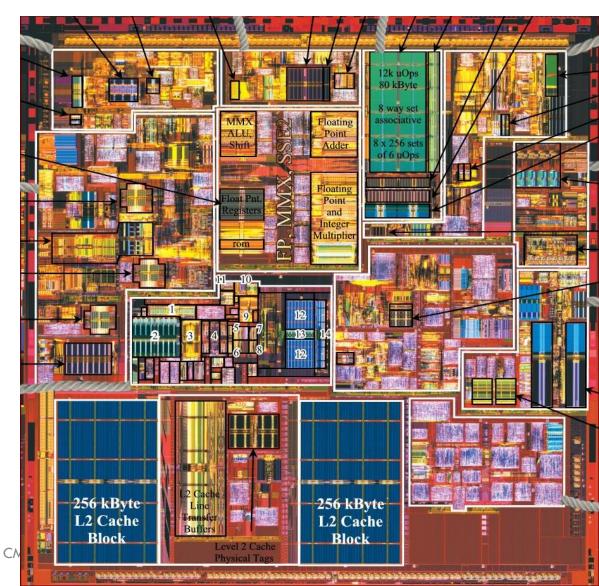
Different processors take radically different approaches

- CPUs: Hardware dynamically schedules instructions
 - Expensive, complex hardware → Few cores (tens)
 - (Relatively) Easy to write fast software



- GPUs: Software makes parallelism explicit
 - Simple, cheap hardware → Many cores (thousands)
 - (Often) Hard to write fast software

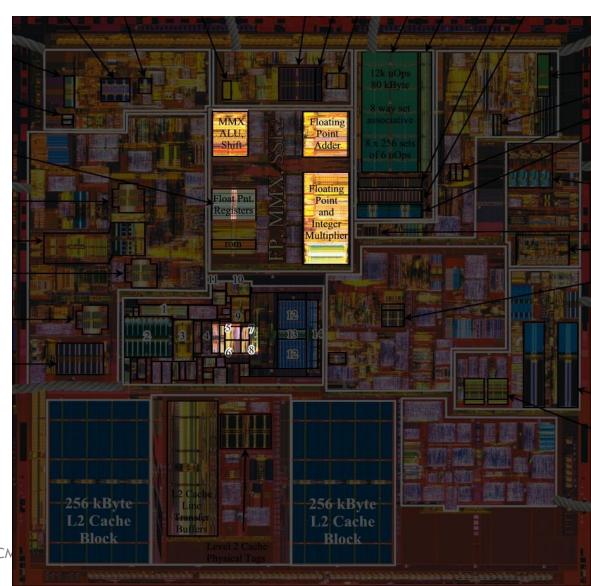
Pentium 4"Northwood" (2002)



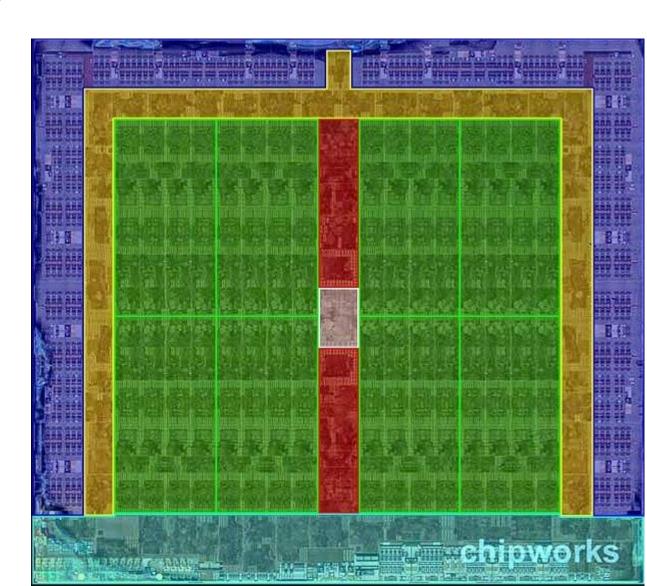
- Pentium 4"Northwood" (2002)
- Highlighted areas actually execute instructions



Most area spent on scheduling (not on executing the program)



■ AMD Fiji (2015)



- AMD Fiji (2015)
- Highlighted areas actually execute instructions
 - → Most area spent executing the program
 - (Rest is mostly I/O & memory, not scheduling)

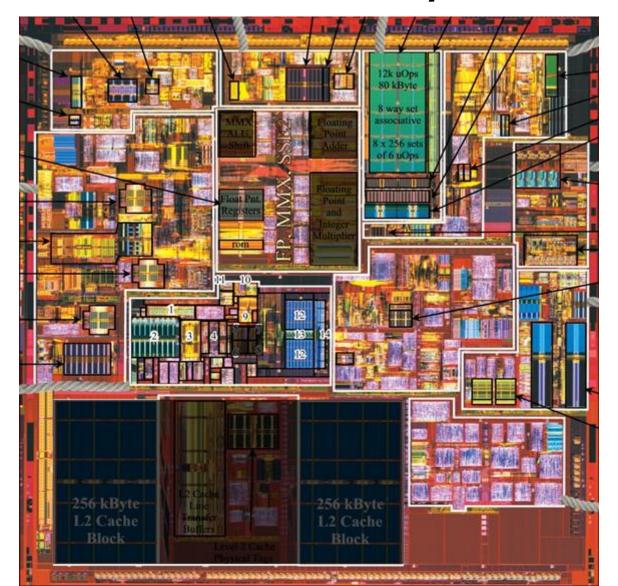


Today you will learn...

How CPUs exploit ILP to speed up sequential code

- Key ideas:
 - Pipelining & Superscalar: Work on multiple instructions at once
 - Out-of-order execution: Dynamically schedule instructions whenever they are "ready"
 - <u>Speculation</u>: Guess what the program will do next to discover more independent work, "rolling back" incorrect guesses
- CPUs must do all of this while preserving the <u>illusion</u> that instructions execute in-order, one-at-a-time

In other words... Today is about:



Buckle up!

...But please ask questions!

```
Compiling on ARM
                                         poly:
                                                    r1, #0
                                           cmp
                                           ble
                                                    .L4
int poly(int *coef,
                                           push
                                                    {r4, r5}
                                                    r3, r0
                                           mov
         int terms, int x) {
                                           add
                                                    r1, r0, r1, lsl #2
  int power = 1;
                                                    r4, #1
                                           movs
                                                    r0, #0
                                           movs
  int value = 0;
                                          .L3:
  for (int j = 0; j < terms; j++) {
                                           ldr
                                                    r5, [r3], #4
                                                    r1, r3
   value += coef[j] * power;
                                           cmp
                                           mla
                                                    r0, r4, r5, r0
    power *= x;
                                           mul
                                                    r4, r2, r4
                                           bne
                                                    .L3
                                                    {r4, r5}
                                           pop
  return value;
                                           bx
                                          .L4:
                                                    r0, #0
                                           movs
                                                    ٦r
                                           bx
                       CMU 15-418/15-618, Spring 2019
```

r0: value

r4: power

r5: coef[j]

r3: &coef[0]

r2: x

r1: &coef[terms]

Compiling on ARM

```
int poly(int *coef,
         int terms, int x) {
  int power = 1;
  int value = 0;
  for (int j = 0; j < terms; j++) {
   value += coef[j] * power;
    power *= x;
  return value;
```

```
r0: value
r1: &coef[terms]
r2: x
r3: &coef[j]
r4: power
r5: coef[j]
```

```
poly:
                                              Preamble
                            r1, #0
                    cmp
                    ble
                             . L4
                    push
                            {r4, r5}
                            r3, r0
                    mov
                            r1, r0, r1, lsl #2
                    add
                            r4, #1
                    movs
                            r0, #0
                    movs
                  .L3:
                            r5, [r3], #4
                    ldr
                                              teration
                            r1, r3
                    cmp
                    mla
                            r0, r4, r5, r0
                            r4, r2, r4
                    mul
                    bne
                            .L3
                            \{r4, r5\}
                    pop
                    bx
                  .L4:
                            r0, #0
                    movs
                             1r
                    bx
CMU 15-418/15-618, Spring 2019
```

r0: value r1: &coef[terms] r2: x r3: &coef[j] r4: power r5: coef[j]

Compiling on ARM

■ Executing poly(A, 3, x)

```
cmp r1, #0
ble .L4
push {r4, r5}
mov r3, r0
add r1, r0, r1, lsl #2
movs r4, #1
movs r0, #0
ldr r5, [r3], #4
cmp r1, r3
mla r0, r4, r5, r0
mul r4, r2, r4
bne .L3
```

CMU 15-418/15-618, Spring 2019

■ Executing poly(A, 3, x)



■ Executing poly(A, 3, x)

```
Preamble
    r1, #0
CMD
ble .L4
    {r4, r5}
push
mov r3, r0
add r1, r0, r1, lsl #2
    r4, #1
movs
    r0, #0
movs
ldr
     r5, [r3], #4
                       =0 Meration
    r1, r3
cmp
mla r0, r4, r5, r0
mul
     r4, r2, r4
bne
       . L3
```

```
ldr
        r5, [r3], #4
                         iteration
        r1, r3
cmp
mla
        r0, r4, r5, r0
mul
        r4, r2, r4
bne
        . L3
ldr
        r5, [r3], #4
                         iteration
        r1, r3
cmp
mla
        r0, r4, r5, r0
        r4, r2, r4
mul
bne
        .L3
        {r4, r5}
pop
bx
        1r
```

■ Executing poly(A, 3, x)

```
Preamble
    r1, #0
CMD
ble .L4
                                 1dr
                                         r5, [r3], #4
                                                          iteration
    \{r4, r5\}
push
                                         r1, r3
                                 cmp
mov r3, r0
                                 mla
                                         r0, r4, r5, r0
add r1, r0, r1, lsl #2
                                 mul
                                         r4, r2, r4
    r4, #1
                                 bne
                                         . L3
movs
    r0, #0
                                 ldr
                                         r5, [r3], #4
                                                          J=2 iteration
movs
ldr
     r5, [r3], #4
                                         r1, r3
                         =0 heration
                                 cmp
                                 mla
                                         r0, r4, r5, r0
cmp r1, r3
mla r0, r4, r5, r0
                                 mu l
                                         r4, r2, r4
mul
     r4, r2, r4
                                 bne
                                         .L3
                                         {r4, r5}
bne
        . L3
                                 pop
                                 bx
                                         1r
                    CMU 15-418/15-618, Spring 2019
```

The software-hardware boundary

- The instruction set architecture (ISA) is a <u>functional</u> <u>contract</u> between hardware and software
 - It says what each instruction does, but not how
 - Example: Ordered sequence of x86 instructions

A processor's microarchitecture is how the ISA is implemented

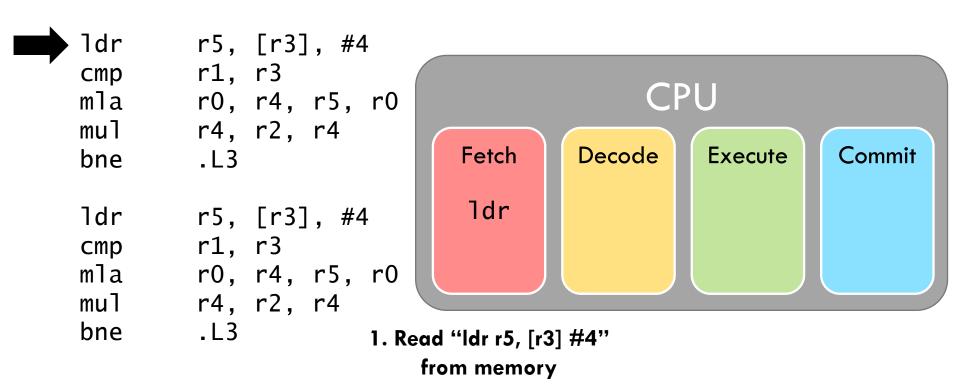
Arch: μ Arch:: Interface: Implementation

Simple CPU model

■ Execute instructions in program order

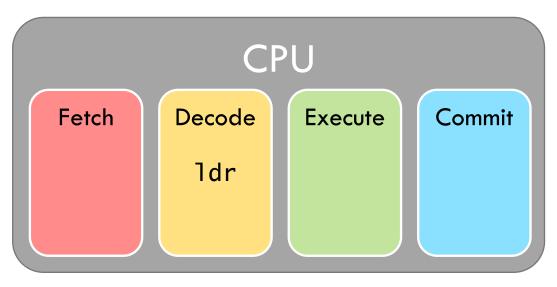
- Divide instruction execution into stages, e.g.:
 - 1. Fetch get the next instruction from memory
 - 2. Decode figure out what to do & read inputs
 - 3. Execute perform the necessary operations
 - 4. Commit write the results back to registers / memory
 - (Real processors have many more stages)

```
ldr
       r5, [r3], #4
      r1, r3
cmp
                                       CPU
mla
     r0, r4, r5, r0
     r4, r2, r4
mu1
                          Fetch
                                  Decode
                                           Execute
                                                    Commit
bne
        .L3
       r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mu1
bne
        .L3
```



CMU 15-418/15-618, Spring 2019

```
ldr
       r5, [r3], #4
      r1, r3
cmp
mla
       r0, r4, r5, r0
     r4, r2, r4
mul
bne
        .L3
ldr
       r5, [r3], #4
       r1, r3
cmp
       r0, r4, r5, r0
mla
        r4, r2, r4
mu1
bne
        .L3
```



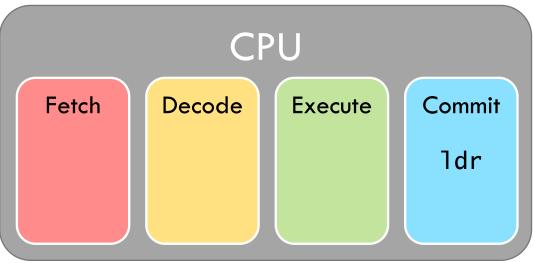
2. Decode "ldr r5, [r3] #4" and read input regs

```
ldr
        r5, [r3], #4
      r1, r3
cmp
                                          CPU
mla
        r0, r4, r5, r0
mul
        r4, r2, r4
                            Fetch
                                     Decode
                                              Execute
                                                        Commit
bne
         .L3
                                                ldr
        r5, [r3], #4
ldr
        r1, r3
cmp
        r0, r4, r5, r0
mla
        r4, r2, r4
mul
bne
         .L3
                                        3. Load memory at r3 and
```

. . .

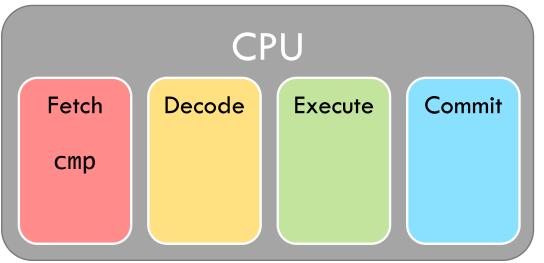
compute r3 + 4

```
ldr
       r5, [r3], #4
      r1, r3
cmp
mla
       r0, r4, r5, r0
mul
      r4, r2, r4
                          Fetch
bne
        .L3
ldr
        r5, [r3], #4
       r1, r3
cmp
       r0, r4, r5, r0
mla
        r4, r2, r4
mul
bne
        .L3
```



4. Write values into regs r5 and r3

```
ldr
       r5, [r3], #4
      r1, r3
\mathsf{cmp}
mla r0, r4, r5, r0
      r4, r2, r4
mul
                           Fetch
bne
        .L3
                           cmp
        r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mul
bne
        .L3
```

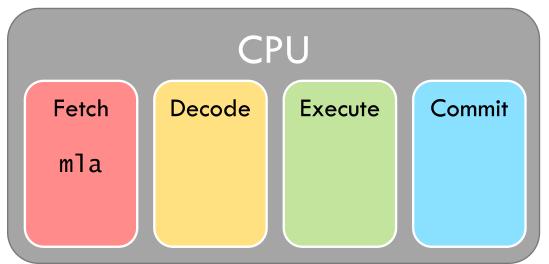


```
ldr
        r5, [r3], #4
      r1, r3
\mathsf{cmp}
                                         CPU
mla r0, r4, r5, r0
      r4, r2, r4
mul
                           Fetch
                                    Decode
                                             Execute
                                                      Commit
bne
        .L3
                                     cmp
        r5, [r3], #4
ldr
        r1, r3
cmp
        r0, r4, r5, r0
mla
        r4, r2, r4
mul
bne
        .L3
```

```
ldr
        r5, [r3], #4
      r1, r3
\mathsf{cmp}
                                        CPU
mla r0, r4, r5, r0
mul r4, r2, r4
                           Fetch
                                    Decode
                                             Execute
                                                      Commit
bne
        .L3
                                              cmp
        r5, [r3], #4
ldr
        r1, r3
cmp
        r0, r4, r5, r0
mla
        r4, r2, r4
mul
bne
        .L3
```

```
ldr
       r5, [r3], #4
      r1, r3
\mathsf{cmp}
                                        CPU
mla r0, r4, r5, r0
mul r4, r2, r4
                           Fetch
                                    Decode
                                             Execute
                                                      Commit
bne
        .L3
                                                        cmp
        r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mul
bne
        .L3
```

```
ldr
       r5, [r3], #4
     r1, r3
cmp
mla r0, r4, r5, r0
     r4, r2, r4
mu1
bne
       .L3
       r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mul
bne
       .L3
```



How fast is this processor? Latency? Throughput?

	1				_	_	_				
Fetch	ldr				cmp				mla		
Decode		ldr				cmp				mla	
Execute			ldr				стр				
Commit				ldr				стр			
							7				

Latency = 4 ns / instr

Throughput = 1 instr / 4 ns

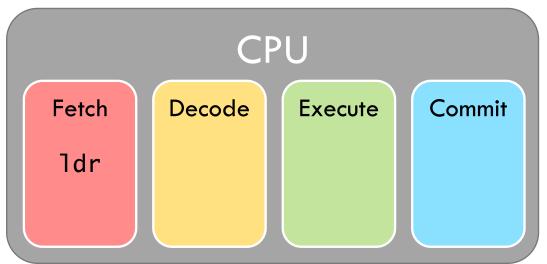
Simple CPU is very wasteful

	1 ns		TIME								
Fetch	ldr				стр				mla		
Decode		ldr		I	dle	стр				mla	
Execute			ldr	Har	dwai	e	стр				•••
Commit				ldr				стр			

Pipelining

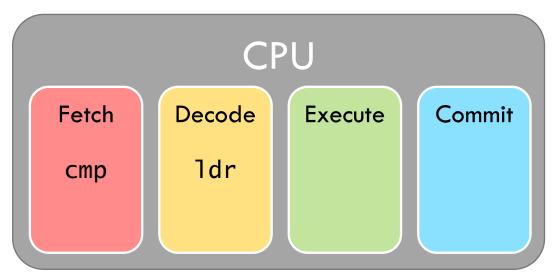
Idea: Start on the next instr'n immediately

```
ldr
       r5, [r3], #4
     r1, r3
cmp
mla r0, r4, r5, r0
     r4, r2, r4
mul
bne
       .L3
       r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mul
       .L3
bne
```



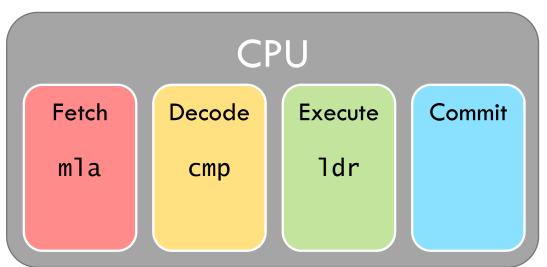
Idea: Start on the next instr'n immediately

```
ldr
       r5, [r3], #4
     r1, r3
cmp
mla r0, r4, r5, r0
     r4, r2, r4
mul
bne
       .L3
       r5, [r3], #4
ldr
       r1, r3
cmp
mla
       r0, r4, r5, r0
       r4, r2, r4
mul
       .L3
bne
```

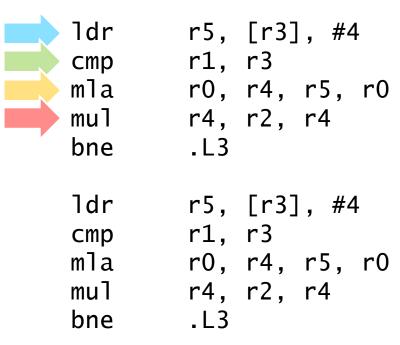


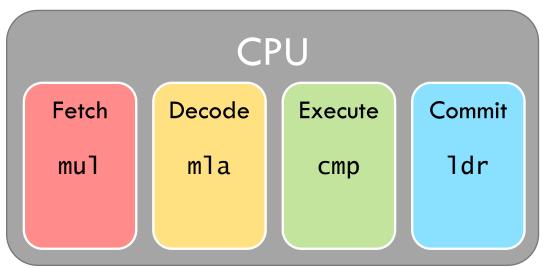
Idea: Start on the next instr'n immediately

```
ldr
        r5, [r3], #4
      r1, r3
\mathsf{cmp}
mla r0, r4, r5, r0
      r4, r2, r4
mul
bne
       .L3
        r5, [r3], #4
ldr
        r1, r3
cmp
mla
        r0, r4, r5, r0
        r4, r2, r4
mul
        .L3
bne
```

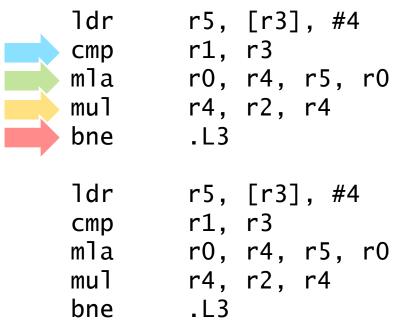


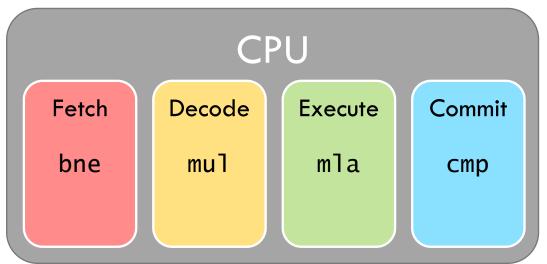
Idea: Start on the next instr'n immediately



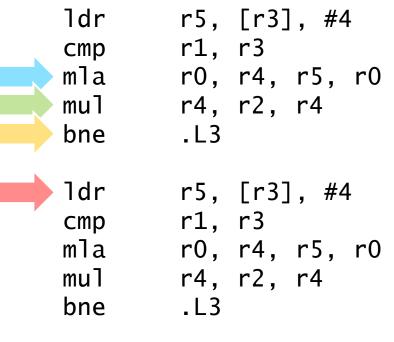


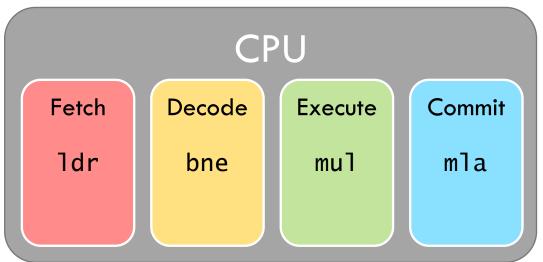
Idea: Start on the next instr'n immediately





Idea: Start on the next instr'n immediately





Evaluating polynomial on the pipelined CPU

How fast is this processor? Latency? Throughput?

1 ns		
	TIME	
1		

Fetch	ldr	стр	mla	mul	bne	ldr	cmp	mla	mul	bne
Decode		ldr	стр	mla	mul	bne	ldr	стр	mla	mul
Execute			ldr	стр	mla	mul	bne	ldr	стр	mla
Commit				ldr	стр	mla	mul	bne	ldr	стр

Latency = 4 ns / instr

Throughput = 1 instr / ns

4X speedup!

CMU 15-418/15-618, Spring 2019

Speedup achieved through pipeline parallelism

TIME

Processor works on 4 instructions at a time

Fetch	ldr	стр	mla	mul	bne	ldr	стр	mla	mul	bne
Decode		ldr	стр	mla	mul	bne	ldr	стр	mla	mul
Execute			ldr	стр	mla	mul	bne	ldr	стр	mla
Commit				ldr	стр	mla	mul	bne	ldr	стр

Limitations of pipelining

Parallelism requires <u>independent</u> work

Q: Are instructions independent?

A: No! Many possible hazards limit parallelism...

Data hazards

```
ldr ra, [rb], #4 // ra ← Memory[rb]; rb ← rb + 4
cmp rc, rd // rc ← rd + re
```

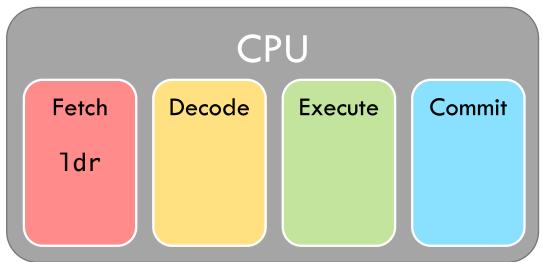
Q: When can the CPU pipeline the cmp behind 1dr?

Fetch	ldr	стр				
Decode		ldr	стр			
Execute			ldr	стр		
Commit				ldr	стр	

- A: When they use different registers
 - Specifically, when cmp does not read any data written by 1dr
 - E.g., rb != rd

■ Cannot pipeline cmp (1dr writes r3)

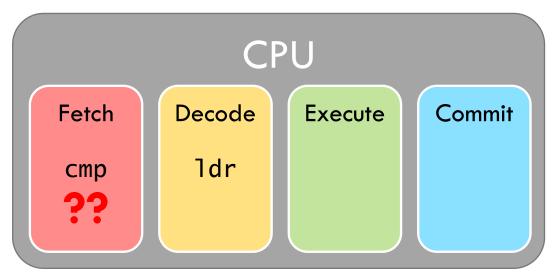
```
ldr
       r5, [r3], #4
      r1, *r3
\mathsf{cmp}
mla r0, r4, r5, r0
mul r4, r2, r4
     .L3
bne
       r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mul
bne
        .L3
```



■ Cannot pipeline cmp (1dr writes r3)

```
ldr
      r5, [r3], #4
     r1, 13
cmp
mla r0, r4, r5, r0
mul r4, r2, r4
     .L3
bne
      r5, [r3], #4
ldr
     r1, r3
cmp
     r0, r4, r5, r0
mla
     r4, r2, r4
mul
```

.L3

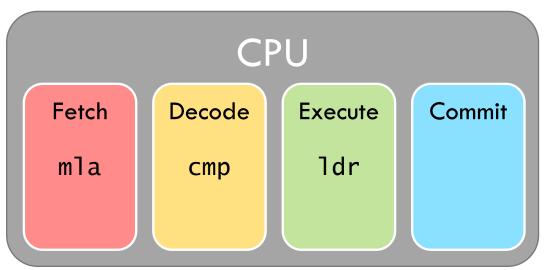


. . .

bne

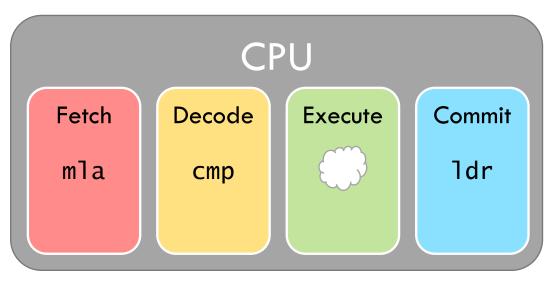
■ Cannot pipeline cmp (1dr writes r3)

```
ldr
       r5, [r3], #4
     r1, r3
cmp
mla r0, r4, r5, r0
mul r4, r2, r4
     .L3
bne
       r5, [r3], #4
ldr
      r1, r3
cmp
     r0, r4, r5, r0
mla
      r4, r2, r4
mul
bne
       .L3
```



■ Cannot pipeline cmp (1dr writes r3)

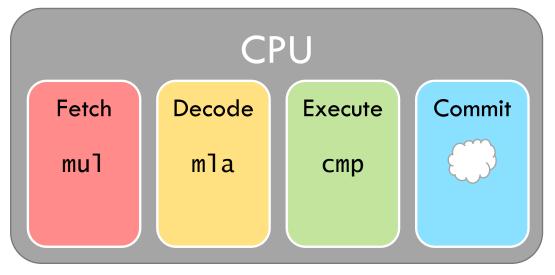
```
ldr
        r5, [r3], #4
      r1, *r3
\mathsf{cmp}
mla r0, r4, r5, r0
mul r4, r2, r4
bne
       .L3
        r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mul
        .L3
bne
```



Inject a "bubble" (NOP) into the pipeline

Cannot pipeline cmp (1dr writes r3)

```
ldr
       r5, [r3], #4
       r1, r3
cmp
mla r0, r4, r5, r0
mul r4, r2, r4
bne
     .L3
       r5, [r3], #4
ldr
       r1, r3
cmp
       r0, r4, r5, r0
mla
       r4, r2, r4
mul
       .L3
bne
```



cmp proceeds once 1dr
has committed

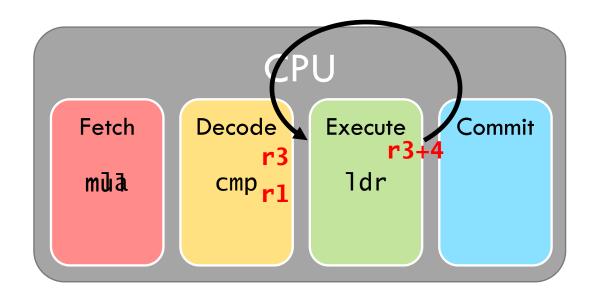
Stalling degrades performance

Processor works on 3 instructions at a time TIME mla mul bne | 1dr mla **Fetch** ldr cmp mul cmp cmp mla ldr mul bne mla mul ldr cmp Decode ldr mla ldr mu l bne Execute cmp ldr mla bne ldr Commit cmp mul

- But stalling is sometimes unavoidable
 - E.g., long-latency instructions (divide, cache miss)

Dealing with data hazards: Forwarding data

Wait a second... data is available after Execute!



Forwarding eliminates many (not all) pipeline stalls

Speedup achieved through pipeline parallelism

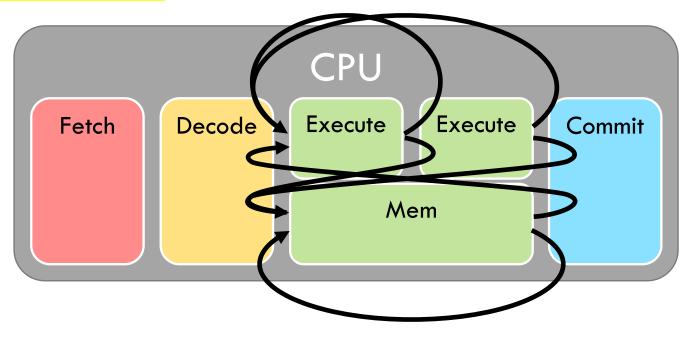
Processor works on 4 instructions at a time ©

Fetch	ldr	стр	mla	mul	bne	ldr	стр	mla	mul	bne
Decode		ldr	стр	mla	mul	bne	ldr	стр	mla	mul
Execute			ldr	стр	mla	mul	bne	ldr	стр	mla
Commit				ldr	стр	mla	mul	bne	ldr	стр

CMU 15-418/15-618, Spring 2019

Pipelining is not free!

- Q: How well does forwarding scale?
- A: Not well... many forwarding paths in deep & complex pipelines



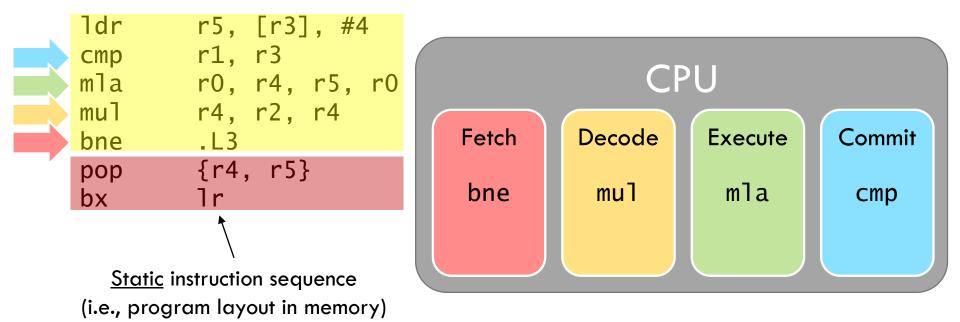
Control hazards + Speculation

- Programs must appear to execute in program order
 - → All instructions depend on earlier ones

- Most instructions implicitly continue at the next...
- But branches redirect execution to new location

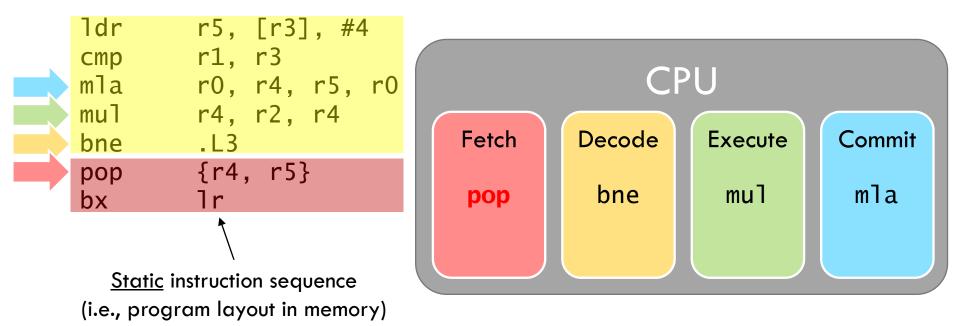
Dealing with control hazards: Flushing the pipeline

What if we always fetch the next instruction?



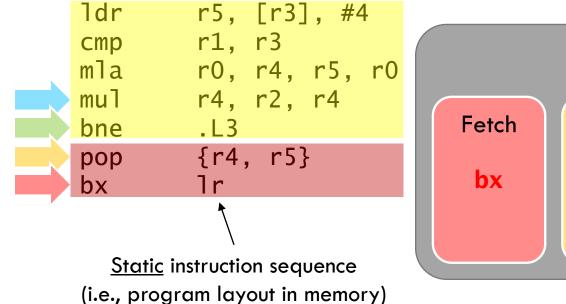
Dealing with control hazards: Flushing the pipeline

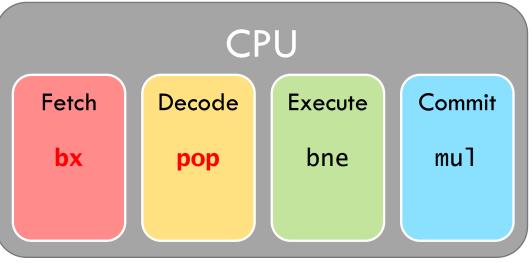
What if we always fetch the next instruction?



Dealing with control hazards: Flushing the pipeline

What if we always fetch the next instruction?



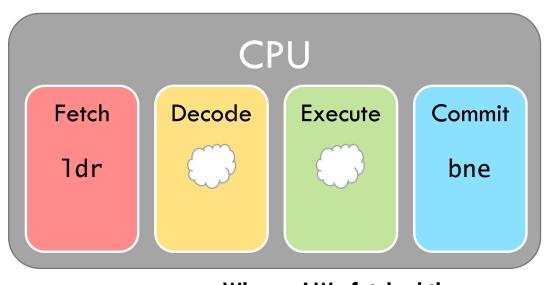


Whoops! We fetched the wrong instructions!
(Loop not finished)

Dealing with control hazards: Flushing the pipeline

What if we always fetch the next instruction?

Static instruction sequence (i.e., program layout in memory)



Whoops! We fetched the wrong instructions!
(Loop not finished)

Pipeline flushes destroy

TIME

performance

Processor works on <u>2 or 3</u> instructions at a time

Fetch	ldr	стр	mla	mul	bne			ldr	стр	mla
Decode		ldr	стр	mla	mul	bne			ldr	стр
Execute			ldr	стр	mla	mul	bne			1dr
Commit				ldr	стр	mla	mul	bne		

Penalty <u>increases</u> with deeper pipelines

Dealing with control hazards: Speculation!

- Processors do not wait for branches to execute
- Instead, they speculate (i.e., guess) where to go next
 + start fetching
- Modern processors use very sophisticated mechanisms
 - E.g., speculate in Fetch stage—before processor even knows instrn is a branch!
 - >95% prediction accuracy
 - Still, branch mis-speculation is major problem

Pipelining Summary

- Pipelining is a simple, effective way to improve throughput
 - lacktriangle N-stage pipeline gives up to N imes speedup



- Pipelining has limits
 - Hard to keep pipeline busy because of hazards
 - Forwarding is expensive in deep pipelines
 - Pipeline flushes are expensive in deep pipelines
- ightharpoonup Pipelining is ubiquitous, but tops out at N pprox 15

Software Takeaways

- Processors with a simple "in-order" pipeline are very sensitive to running "good code"
 - Compiler should target a specific model of CPU
 - Low-level assembly hacking
- ...But very few CPUs are in-order these days
 - E.g., embedded, ultra-low-power applications
- Instead, ≈all modern CPUs are "out-of-order"
 - Even in classic "low-power domains" (like mobile)

Out-of-Order Execution

Increasing parallelism via dataflow

 Parallelism limited by many false dependencies, particularly sequential program order

- <u>Dataflow</u> tracks how instructions actually depend on each other
 - True dependence: read-after-write

Dataflow increases parallelism by eliminating unnecessary dependences

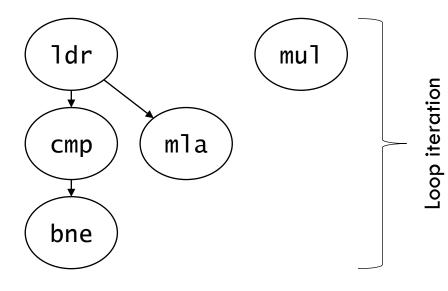


Example: Dataflow in polynomial

evaluation

ldr	r5, [r3], #4
cmp	r1, r3
mla	r0, r4, r5, r0
mul	r4, r2, r4
bne	.L3

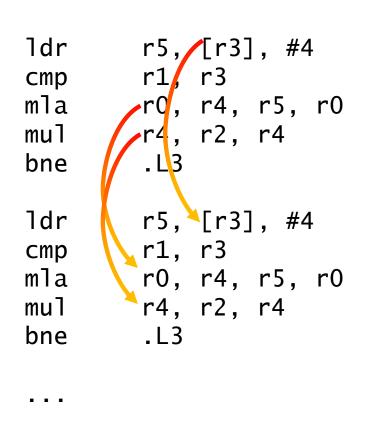
ldr	r5,	[r3]	, i	#4
cmp	r1,	r3		
mla	rO,	r4,	r5	, r0
mul	r4,	r2,	r4	
bne	.L3			





Example: Dataflow

evaluation



bne polynomial ldr mu1 mla cmp bne ldr mul mla cmp bne CMU 15-418/15-618 ldr mul

Example: Dataflow polynomial execution

- Execution only, with perfect scheduling & unlimited execution units
 - Idr, mul execute in 2 cycles
 - cmp, bne execute in 1 cycle
 - mla executes in 3 cycles

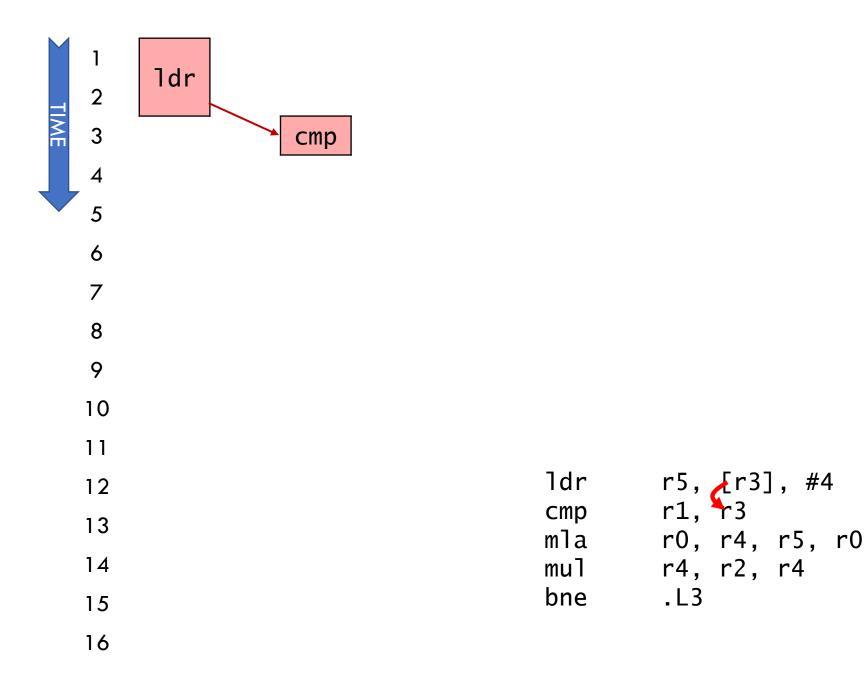


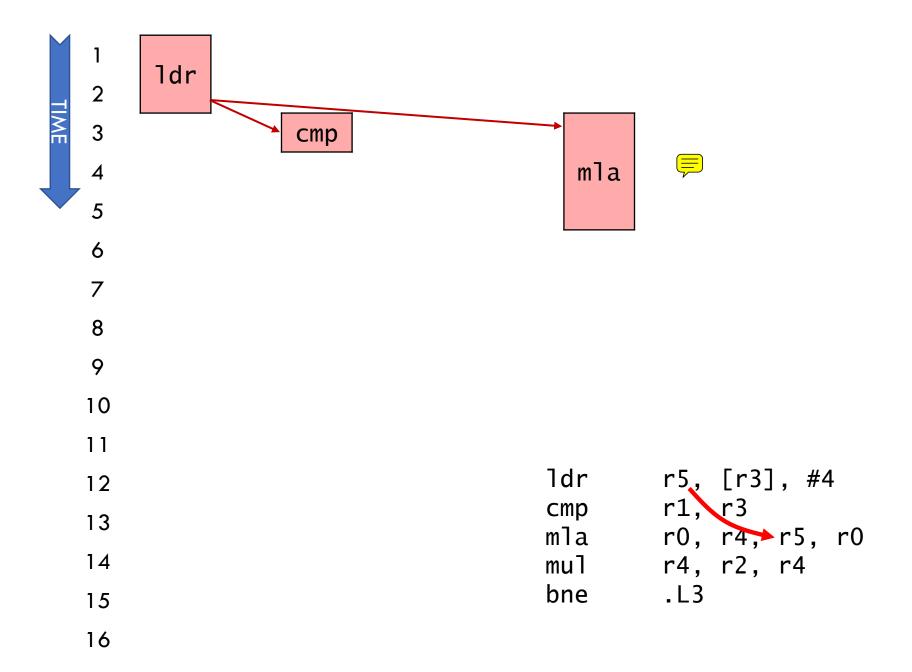
Q: Does dataflow speedup execution? By how much?

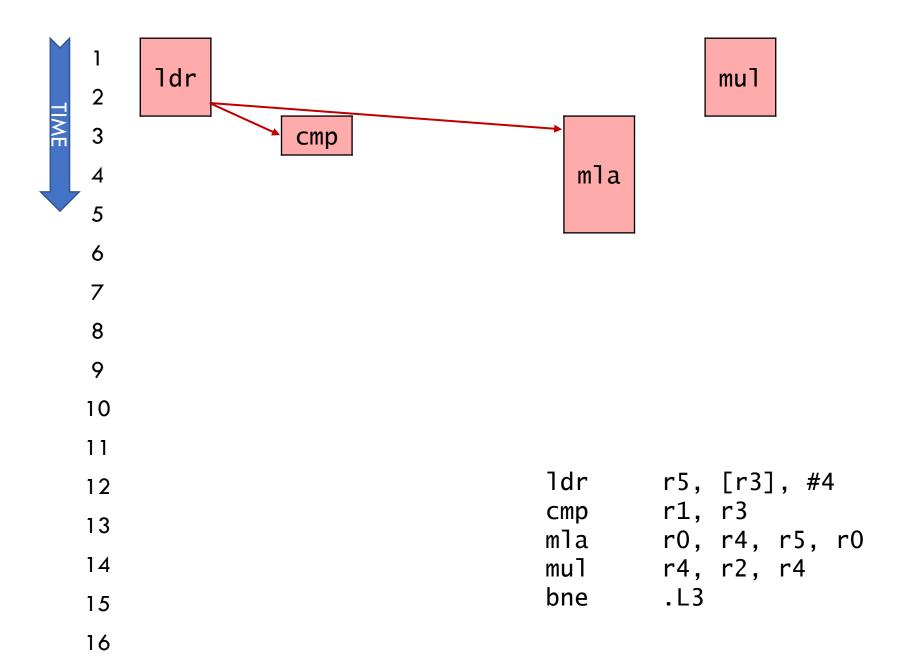
Q: What is the performance bottleneck?

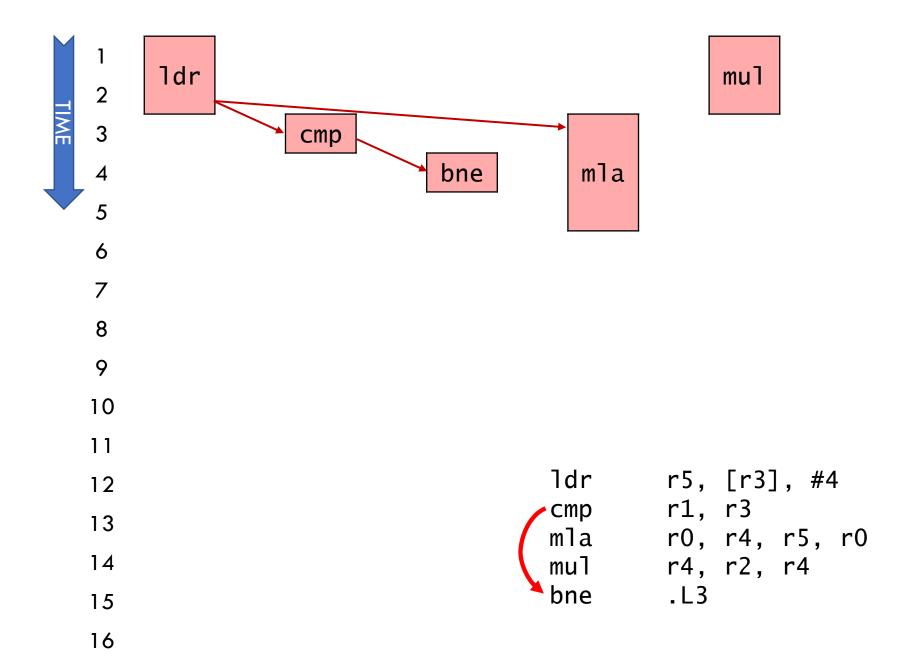
	ldr	
2		
3		
4		
5		
6		
7		
8		
9		
10		
11		
12		
13		
14		
15		
16		

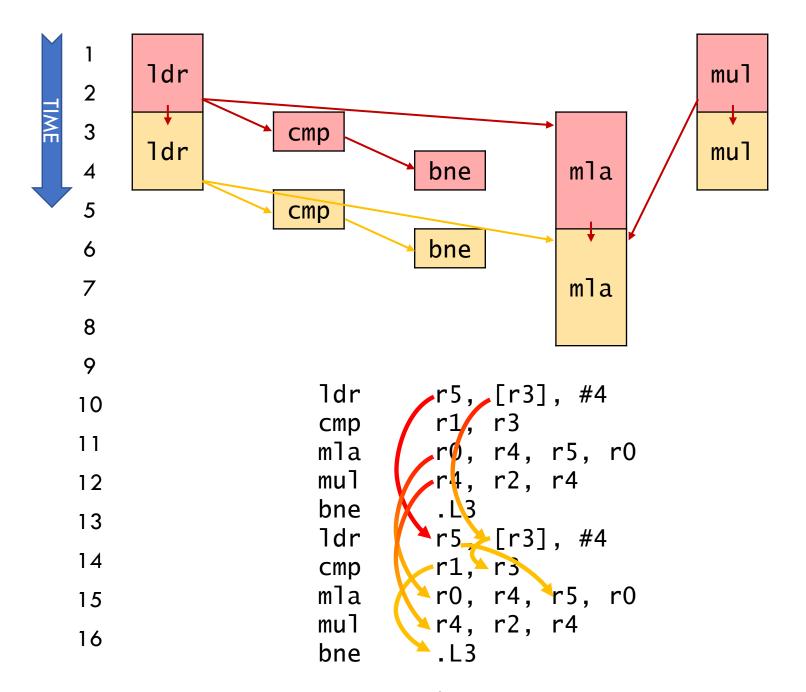
ldr r5, [r3], #4 cmp r1, r3 mla r0, r4, r5, r0 mul r4, r2, r4 bne .L3

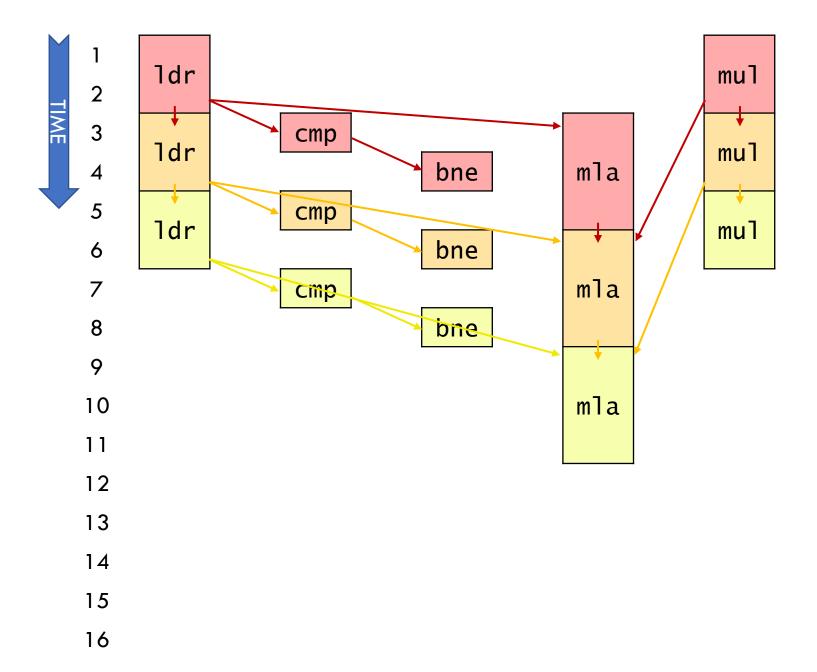


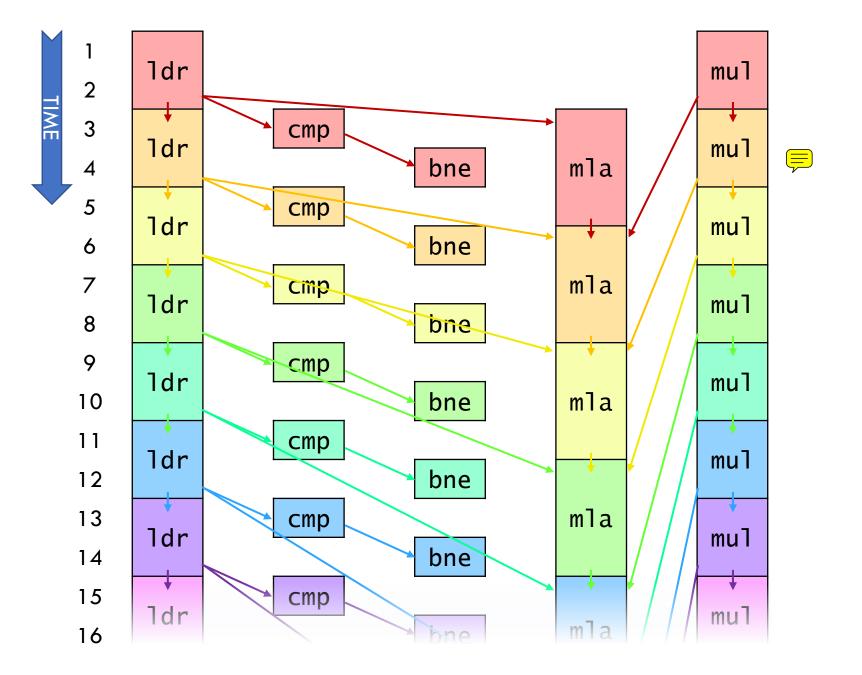












Example: Dataflow polynomial execution

- Q: Does dataflow speedup execution? By how much?
 - Yes! 3 cycles / loop iteration



■ Instructions per cycle (IPC) = $5/3 \approx 1.67$ (vs. 1 for perfect pipelining)

- Q: What is the performance bottleneck?
 - mla: Each mla depends on previous mla & takes 3 cycles
 - → This program is **latency-bound**

Latency Bound

- What is the "critical path" of the computation?
 - Longest path across iterations in dataflow graph
 - E.g., mla in last slide (but could be multiple ops)

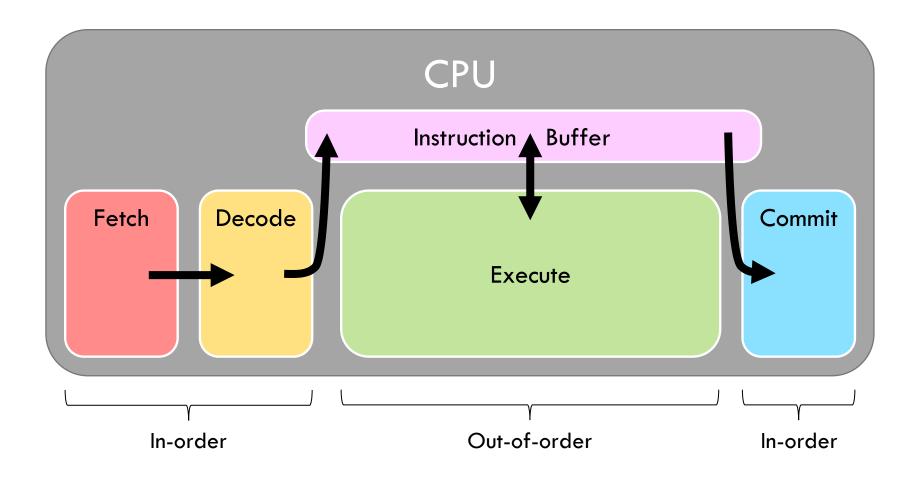
- Critical path limits maximum performance
- Real CPUs may not achieve latency bound, but useful mental model + tool for program analysis

Out-of-order (OoO) execution uses dataflow to increase parallelism

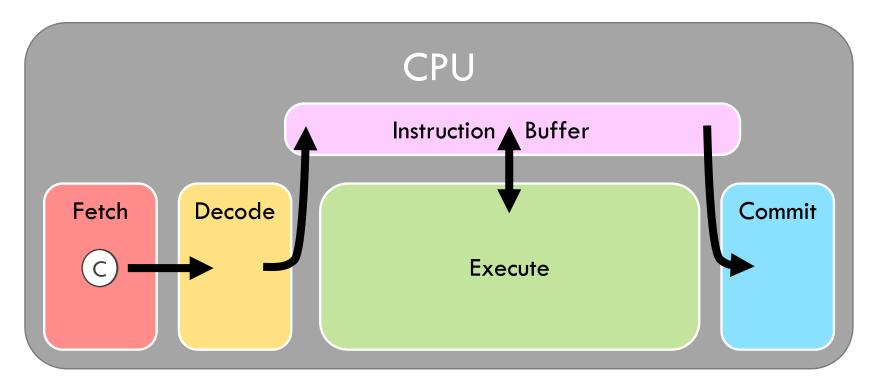
Idea: Execute programs in dataflow order, but give the illusion of sequential execution

- This is a "restricted dataflow" model
 - Restricted to instructions near those currently committing
 - (Pure dataflow processors also exist that expose dataflow to software)

High-level OoO microarchitecture



OoO is **hidden** behind in-order frontend & commit



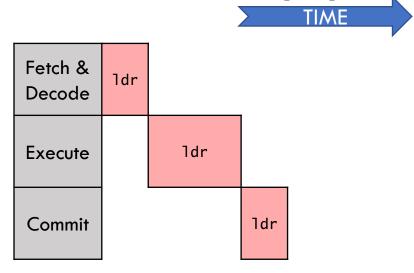
Instructions only enter & leave instruction buffer in program order; all bets are off in between!

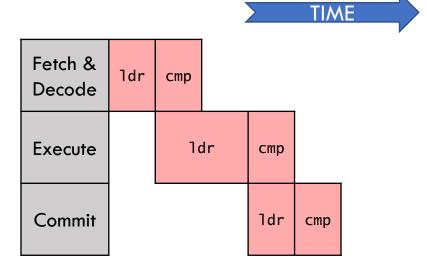
Example: OoO polynomial evaluation

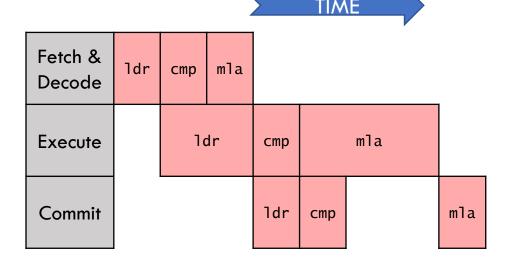
Q: Does OoO speedup execution? By how much?

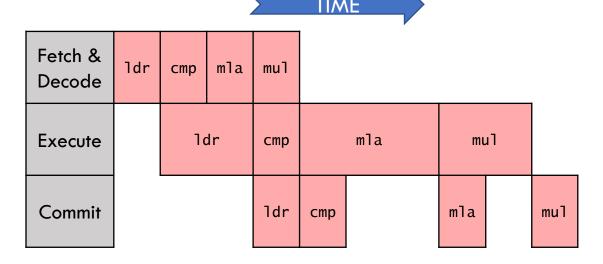
Q: What is the performance bottleneck?

Assume perfect forwarding & branch prediction









Fetch & Decode	ldr	стр	mla	mul	bne	,				
Execute		10	dr	стр		mla	mı	u]	bne	
Commit				ldr	стр		mla		mul	bne

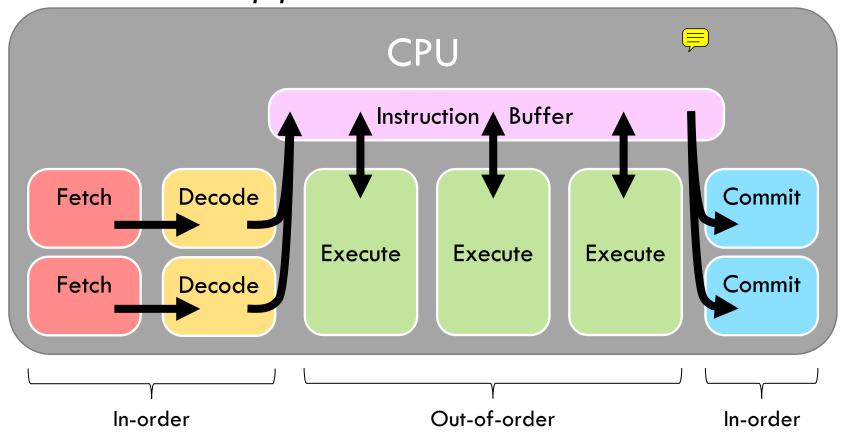
Fetch & Decode	ldr	стр	mla	mul	bne	ldr	стр	mla	mul	bne	ldr	стр	mla	mul	bne	ldr
Execute		10	dr	стр		mla		mul		bne	bne ldr		стр	mla		
Commit				ldr	стр			mla		mul	bne		ldr	стр		

Fetch & Decode	ldr	стр	mla	mul	bne	ldr	стр	mla	mul	bne	ldr	стр	mla	mul	bne	ldr
Execute		10	dr	стр		mla		mul		bne	bne ldr		стр	mla		
Commit				ldr	стр			mla		mul	bne		ldr	стр		

- Wait a minute... this isn't OoO... or even faster than a simple pipeline!
- Q: What went wrong?
- A: We're throughput-limited: can only exec 1 instrn

High-level **Superscalar** OoO microarchitecture

■ Must increase pipeline width to increase ILP > 1



Focus on Execution, not Fetch & Commit

- Goal of OoO design is to only be limited by dataflow execution
- Fetch and commit are over-provisioned so that they (usually) do not limit performance
 - → Programmers can (usually) ignore fetch/commit

- **Big Caveat:** Programs with *inherently unpredictable* control flow will often be limited by fetch stalls (branch misprediction)
 - E.g., branching based on random data

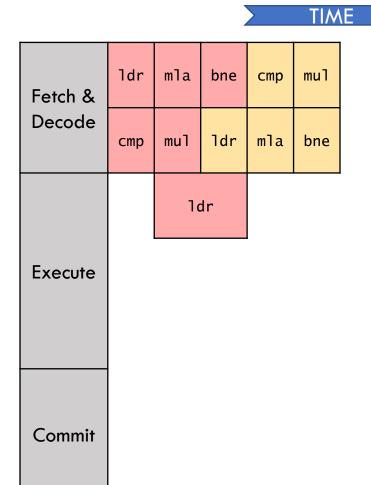
1dr Fetch & Decode cmp Execute Commit

ldr r5, [r3], #4 r1, r3 cmp mla r0, r4, r5, r0 mul r4, r2, r4 bne .L3 ldr r5, [r3], #4 r1, r3 cmp r0, r4, r5, r0 mla mu1 r4, r2, r4 .13 bne

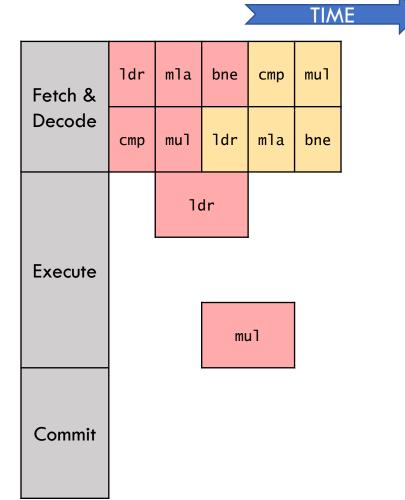
TIME	

Fetch &	ldr	mla	bne	стр	mul
Decode	стр	mul	ldr	mla	bne
Execute					
Commit					

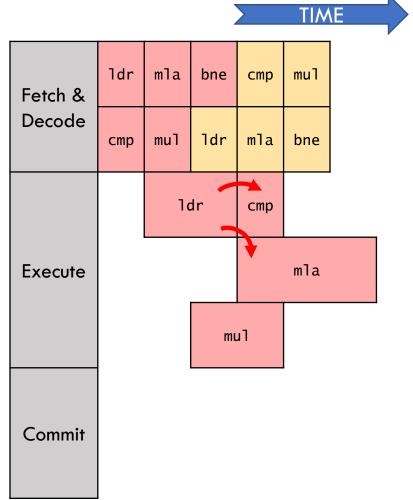
```
ldr
        r5, [r3], #4
        r1, r3
cmp
        r0, r4, r5, r0
mla
mul
        r4, r2, r4
bne
        .L3
ldr
        r5, [r3], #4
        r1, r3
cmp
        r0, r4, r5, r0
mla
mul
        r4, r2, r4
bne
        .13
```



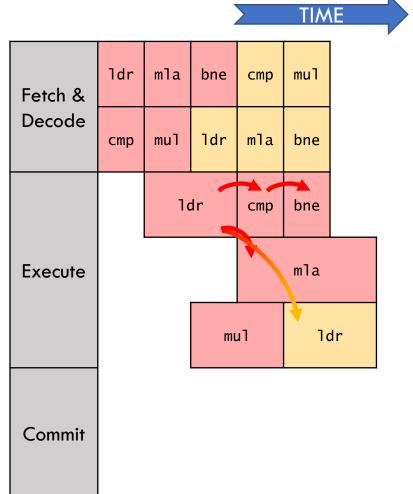
ldr	r5 <mark>,,[</mark> r3], #4
cmp	r1, r3
mla	r0, r4, r5, r0
mu1	r4, r2, r4
bne	.L3
ldr	r5, [r3], #4
cmp	r1, r3
mla	r0, r4, r5, r0
mul	r4, r2, r4
bne	.L3

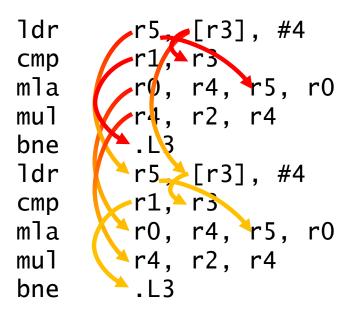


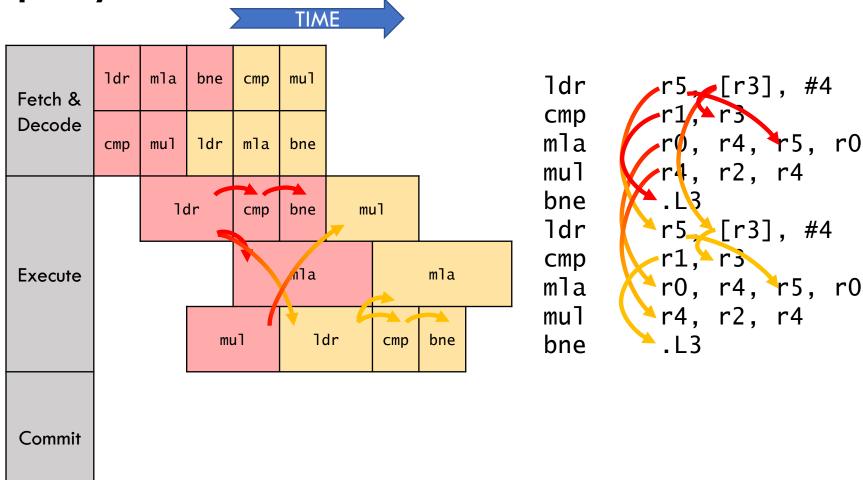
ldr	r5[r3], #4
cmp	r1, r3
mla	r0, r4, r5, r0
mul	r4, r2, r4
bne	•.L3
ldr	r5, [r3], #4
cmp	r1, r3
mla	r0, r4, r5, r0
mul	r4, r2, r4
bne	.L3

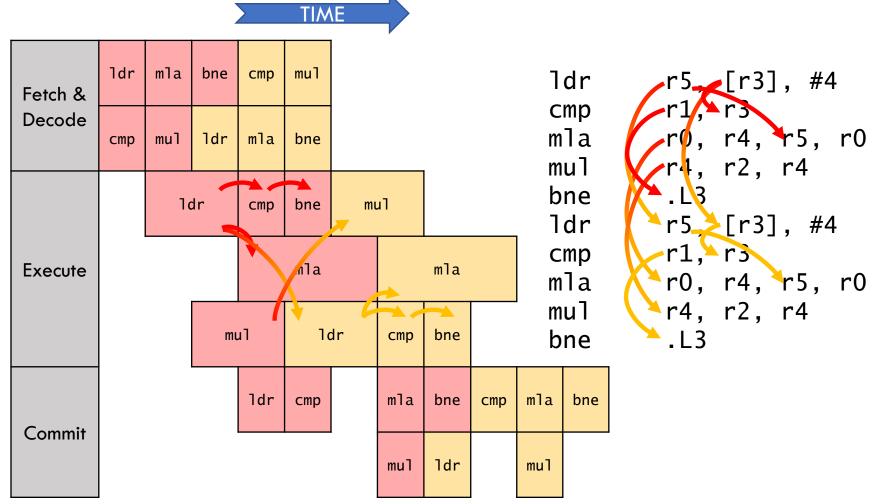


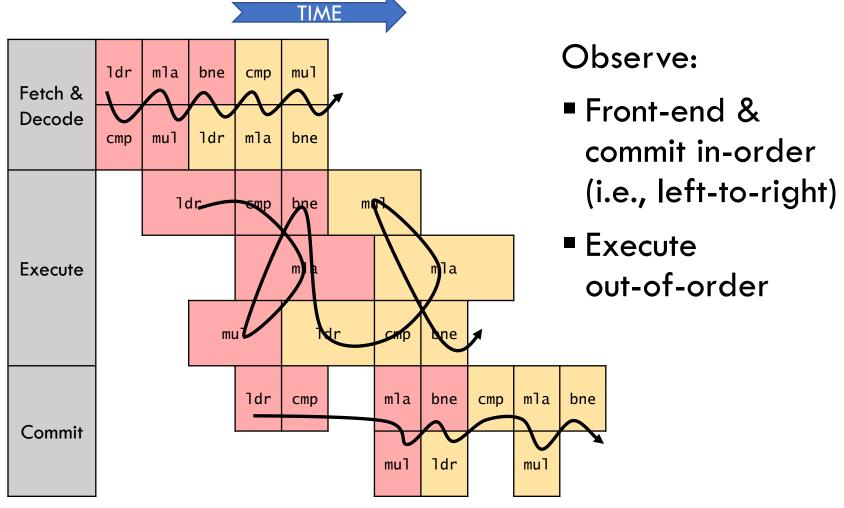
```
ldr
        r5 [r3], #4
cmp
        r0, r4, r5, r0
mla
         r4, r2, r4
mul
bne
ldr
        r5, [r3], #4
        r1, r3
cmp
        r0, r4, r5, r0
mla
        r4, r2, r4
mu1
         .13
bne
```











CMU 15-418/15-618, Spring 2019

1dr mla mla mla bne mu1 1dr bne cmp mul ldr mu1 ldr cmp bne cmp Fetch & Decode mu1 ldr mla 1dr mla. ldr mu1 mla cmp bne cmp bne cmp mu1 bne cmp 1dr 1dr 1dr cmp bne mu1 cmp bne mu1 cmp mla mla mla mla Execute mla ldr 1dr mu1 cmp bne mu1 bne mu1 cmp 1dr mla mla cmp mla. bne cmp bne 🔊 cmp bne cmp mla Commit mul ldr ldr mu1 1dr mu1

Structural hazards: Other throughput limitations

- Execution units are specialized
 - Floating-point (add/multiply)
 - Integer (add/multiply/compare)
 - Memory (load/store)
- Processor designers must choose which execution units to include and how many



 Structural hazard: Data is ready, but instr'n cannot issue because no hardware is available

Example: Structural hazards can severely limit performance

Fetch &	ldr	mla	bne	стр	mul	ldr	mla	bne	стр	mul	ldr	mla	bne	стр	mul	ldr
Decode	стр	mul	ldr	mla	bne	стр	mul	ldr	mla	bne	стр	mul	ldr	mla	bne	стр
Mem Execute	ldr		ldr		ldr		ldı		ldr		1		dr 10		dr	
Int Execute				стр	bne	стр	bne		стр	bne	стр	bne		стр	bne	стр
Mult Execute				mla	mla		ı]		mla		mı	ul		mla		mul
Commit				ldr	стр	mla		mul	ldr		mla		mul	ldr		mla
Commit								bne	emp			_/	bne	стр		

Throughput Bound

- Ingredients:
 - Number of operations to perform (of each type)
 - Number & issue rate of "execution ports"/"functional units" (of each type)
- Throughput bound = ops / issue rate
 - E.g., (1 mla + 1 mul) / (2 + 3 cycles)



Again, a real CPU might not exactly meet this bound

Software Takeaway

- OoO is much less sensitive to "good code"
 - Better performance portability
 - Of course, compiler still matters

- OoO makes performance analysis much simpler
 - Throughput bound: Availability of execution ports
 - **Latency bound:** "Critical path" latency
 - Slowest gives good approximation of program perf

Out-of-Order Execution: Under the Hood

OoO x86: Microcoding

- Each x86 instruction describes several operations
 - E.g., add [esp+4], 5 means:
 - 1. Load Mem[esp+4]
 - 2. Add 5 to it
 - 3. Store result to Mem[esp+4]
- This is too much for (fast) hardware

- Instead, hardware decodes instr'ns into micro-ops
 - Rest of pipeline uses micro-ops

Register Renaming

- "False dependences" can severely limit parallelism
 - Write-after-read
 - Write-after-write

- OoO processors eliminate false dependences by transparently renaming registers
 - CPU has many more "physical" than "architectural" registers
 - Each time register is written, it is allocated to a new physical register
 - Physical registers freed when instructions commit

Memory Disambiguation

- CPU must respect store → load ordering
 - E.g., a later instruction reads a value from memory written by an earlier instruction
- But what if the OoO CPU executes the load first?
 - Must "rollback" + execute the load again (next slide)

 Corollary: OoO CPU must track the order of all loads & stores, and only write memory when a store commits

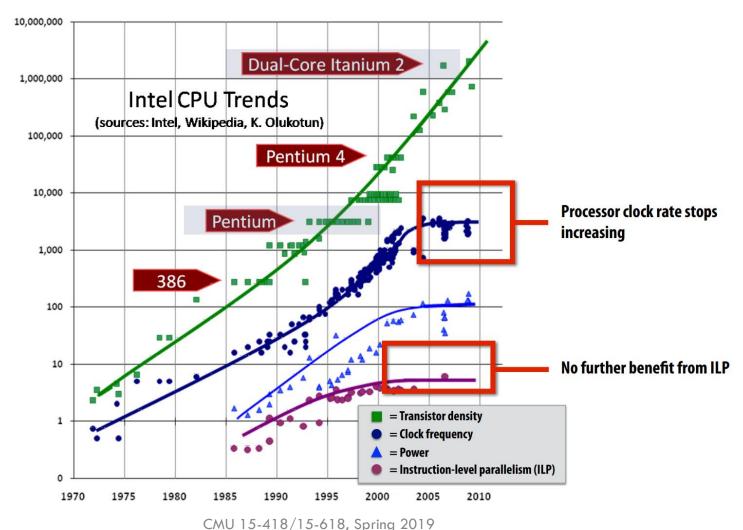
Rollback & Recovery

- OoO CPUs speculate constantly to improve performance
 - E.g., even guessing the results of a computation ("value prediction")

- Need mechanisms to "rollback" to an earlier point in execution when speculation goes awry
 - Complex: Need to recover old register names, flush pending memory operations, etc
- Very expensive: Up to hundreds of instrns of work lost!

Scaling Instruction-Level Parallelism

Recall from last time: ILP & pipelining tapped out... why?



Superscalar scheduling is complex & hard to scale

Q: When is it safe to issue two instructions?



- A: When they are independent
 - Must compare <u>all pairs</u> of input and output registers
- lacktriangle Scalability: $O(W^2)$ comparisons where W is "issue width" of processor
 - Not great!

Limitations of ILP

- 4-wide superscalar \times 20-stage pipeline = **80** instrns in flight
- High-performance OoO buffers hundreds of instructions
- Programs have limited ILP
 - Even with perfect scheduling, >8-wide issue doesn't help
- Pipelines can only go so deep
 - Branch misprediction penalty grows
 - Frequency (GHz) limited by power
- Dynamic scheduling overheads are significant
- Out-of-order scheduling is expensive

Limitations of ILP Multicore

- ILP works great! ...But is complex + hard to scale
- From hardware perspective, multicore is much more efficient, but...

Parallel software is hard!

- Industry resisted multicore for as long as possible
- When multicore finally happened, CPU μ arch simplified \rightarrow more cores
- Many program(mer)s still struggle to use multicore effectively