Computer Architecture

Unit 7: Superscalar Pipelines

Slides developed by Milo Martin & Amir Roth at the University of Pennsylvania with sources that included University of Wisconsin slides by Mark Hill, Guri Sohi, Jim Smith, and David Wood

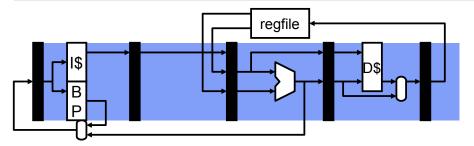
Computer Architecture | Prof. Milo Martin | Superscalar

1

A Key Theme: Parallelism

- Previously: pipeline-level parallelism
 - Work on execute of one instruction in parallel with decode of next
- Next: instruction-level parallelism (ILP)
 - Execute multiple independent instructions fully in parallel
- Then:
 - Static & dynamic scheduling
 - Extract much more ILP
 - Data-level parallelism (DLP)
 - Single-instruction, multiple data (one insn., four 64-bit adds)
 - Thread-level parallelism (TLP)
 - Multiple software threads running on multiple cores

"Scalar" Pipeline & the Flynn Bottleneck



- So far we have looked at scalar pipelines
 - One instruction per stage
 - With control speculation, bypassing, etc.
 - Performance limit (aka "Flynn Bottleneck") is CPI = IPC = 1
 - Limit is never even achieved (hazards)
 - Diminishing returns from "super-pipelining" (hazards + overhead)

Computer Architecture | Prof. Milo Martin | Superscalar

3

An Opportunity...

• But consider:

- Why not execute them **at the same time**? (We can!)
- What about:

- In this case, *dependences* prevent parallel execution
- What about three instructions at a time?
 - Or four instructions at a time?

What Checking Is Required?

For two instructions: 2 checks

```
ADD src1_1, src2_1 \rightarrow dest_1
ADD src1_2, src2_2 \rightarrow dest_2 (2 checks)
```

For three instructions: 6 checks

```
ADD src1_1, src2_1 \rightarrow dest_1

ADD src1_2, src2_2 \rightarrow dest_2 (2 checks)

ADD src1_3, src2_3 \rightarrow dest_3 (4 checks)
```

For four instructions: 12 checks

```
ADD \operatorname{srcl}_1, \operatorname{src2}_1 -> \operatorname{dest}_1

ADD \operatorname{srcl}_2, \operatorname{src2}_2 -> \operatorname{dest}_2 (2 checks)

ADD \operatorname{srcl}_3, \operatorname{src2}_3 -> \operatorname{dest}_3 (4 checks)

ADD \operatorname{srcl}_4, \operatorname{src2}_4 -> \operatorname{dest}_4 (6 checks)
```

Plus checking for load-to-use stalls from prior n loads

```
Computer Architecture | Prof. Milo Martin | Superscalar
```

5

What Checking Is Required?

For two instructions: 2 checks

```
ADD src1_1, src2_1 \rightarrow dest_1
ADD src1_2 src2_2 \rightarrow dest_2 (2 checks)
```

For three instructions: 6 checks

```
ADD src1_1, src2_1 \rightarrow dest_1

ADD src1_2 src2_2 \rightarrow dest_2 (2 checks)

ADD src1_3, src2_3 \rightarrow dest_3 (4 checks)
```

For four instructions: 12 checks

```
ADD \operatorname{src1}_1, \operatorname{src2}_1 \rightarrow \operatorname{dest}_1

ADD \operatorname{src1}_2, \operatorname{sre2}_2 \rightarrow \operatorname{dest}_2 (2 checks)

ADD \operatorname{src1}_3 \operatorname{src2}_4 \rightarrow \operatorname{dest}_4 (4 checks)

ADD \operatorname{src1}_4 \operatorname{src2}_4 \rightarrow \operatorname{dest}_4 (6 checks)
```

Plus checking for load-to-use stalls from prior n loads

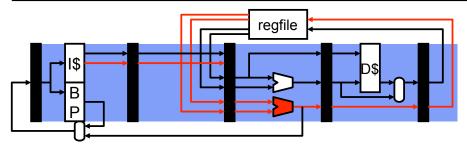
```
Computer Architecture | Prof. Milo Martin | Superscalar
```

How do we build such "superscalar" hardware?

Computer Architecture | Prof. Milo Martin | Superscalar

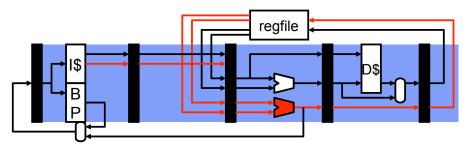
7

Multiple-Issue or "Superscalar" Pipeline



- Overcome this limit using multiple issue
 - Also called **superscalar**
 - Two instructions per stage at once, or three, or four, or eight...
 - "Instruction-Level Parallelism (ILP)" [Fisher, IEEE TC'81]
- Today, typically "4-wide" (Intel Core i7, AMD Opteron)
 - Some more (Power5 is 5-issue; Itanium is 6-issue)
 - Some less (dual-issue is common for simple cores)

A Typical Dual-Issue Pipeline (1 of 2)

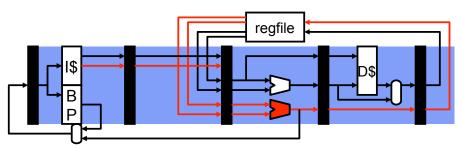


- Fetch an entire 16B or 32B cache block
 - 4 to 8 instructions (assuming 4-byte average instruction length)
 - Predict a single branch per cycle
- Parallel decode
 - Need to check for conflicting instructions
 - Is output register of I₁ is an input register to I₂?
 - Other stalls, too (for example, load-use delay)

Computer Architecture | Prof. Milo Martin | Superscalar

9

A Typical Dual-Issue Pipeline (2 of 2)



- Multi-ported register file
 - Larger area, latency, power, cost, complexity
- Multiple execution units
 - Simple adders are easy, but bypass paths are expensive
- Memory unit
 - Single load per cycle (stall at decode) probably okay for dual issue
 - Alternative: add a read port to data cache
 - Larger area, latency, power, cost, complexity

Computer Architecture | Prof. Milo Martin | Superscalar

How Much ILP is There?

- The compiler tries to "schedule" code to avoid stalls
 - Even for scalar machines (to fill load-use delay slot)
 - Even harder to schedule multiple-issue (superscalar)
- How much ILP is common?
 - · Greatly depends on the application
 - Consider memory copy
 - Unroll loop, lots of independent operations
 - Other programs, less so
- Even given unbounded ILP, superscalar has implementation limits
 - IPC (or CPI) vs clock frequency trade-off
 - Given these challenges, what is reasonable today?
 - ~4 instruction per cycle maximum

Computer Architecture | Prof. Milo Martin | Superscalar

11

Superscalar Implementation Challenges

Superscalar Challenges - Front End

Superscalar instruction fetch

- Modest: fetch multiple instructions per cycle
- Aggressive: buffer instructions and/or predict multiple branches

• Superscalar instruction decode

Replicate decoders

Superscalar instruction issue

- Determine when instructions can proceed in parallel
- More complex stall logic order N² for N-wide machine
- Not all combinations of types of instructions possible

Superscalar register read

- Port for each register read (4-wide superscalar → 8 read "ports")
- Each port needs its own set of address and data wires
 - Latency & area ∝ #ports²

Computer Architecture | Prof. Milo Martin | Superscalar

13

Superscalar Challenges - Back End

• Superscalar instruction execution

- Replicate arithmetic units (but not all, say, integer divider)
- Perhaps multiple cache ports (slower access, higher energy)
 - Only for 4-wide or larger (why? only ~35% are load/store insn)

Superscalar bypass paths

- More possible sources for data values
- Order (N² * P) for *N*-wide machine with execute pipeline depth *P*

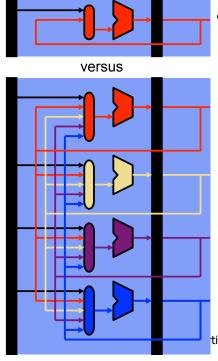
Superscalar instruction register writeback

- One write port per instruction that writes a register
- Example, 4-wide superscalar → 4 write ports

Fundamental challenge:

- Amount of ILP (instruction-level parallelism) in the program
- Compiler must schedule code and extract parallelism

Superscalar Bypass



N² bypass network

- N+1 input muxes at each ALU input
- N² point-to-point connections
- Routing lengthens wires
- Heavy capacitive load
- And this is just one bypass stage (MX)!
 - There is also WX bypassing
 - Even more for deeper pipelines
- One of the big problems of superscalar
 - Why? On the critical path of single-cycle "bypass & execute" loop

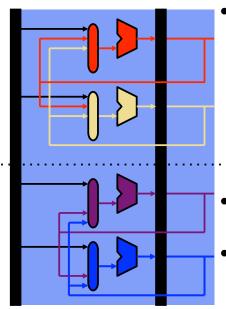
tin | Superscalar

15

Not All N² Created Equal

- N² bypass vs. N² stall logic & dependence cross-check
 - Which is the bigger problem?
- N² bypass ... by far
 - 64- bit quantities (vs. 5-bit)
 - Multiple levels (MX, WX) of bypass (vs. 1 level of stall logic)
 - Must fit in one clock period with ALU (vs. not)
- Dependence cross-check not even 2nd biggest N² problem
 - Regfile is also an N² problem (think latency where N is #ports)
 - And also more serious than cross-check

Mitigating N² Bypass & Register File



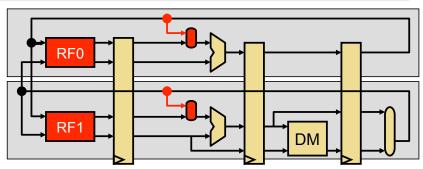
- **Clustering**: mitigates N² bypass
 - Group ALUs into K clusters
 - Full bypassing within a cluster
 - Limited bypassing between clusters
 - With 1 or 2 cycle delay
 - Can hurt IPC, but faster clock
 - (N/K) + 1 inputs at each mux
- (N/K)² bypass paths in each cluster
- **Steering**: key to performance
 - Steer dependent insns to same cluster
- Cluster register file, too
 - Replica a register file per cluster
 - All register writes update all replicas
 - Fewer read ports; only for cluster

Computer Architecture | Prof. Milo Martin | Superscalar

17

Mitigating N² RegFile: Clustering++





- Clustering: split N-wide execution pipeline into K clusters
 - With centralized register file, 2N read ports and N write ports
- Clustered register file: extend clustering to register file
 - Replicate the register file (one replica per cluster)
 - Register file supplies register operands to just its cluster
 - All register writes go to all register files (keep them in sync)
 - Advantage: fewer read ports per register!
 - K register files, each with 2N/K read ports and N write ports

Computer Architecture | Prof. Milo Martin | Superscalar

18

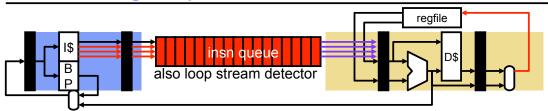
Another Challenge: Superscalar Fetch

- What is involved in fetching multiple instructions per cycle?
- In same cache block? → no problem
 - 64-byte cache block is 16 instructions (~4 bytes per instruction)
 - Favors larger block size (independent of hit rate)
- What if next instruction is last instruction in a block?
 - Fetch only one instruction that cycle
 - Or, some processors may allow fetching from 2 consecutive blocks
- What about taken branches?
 - How many instructions can be fetched on average?
 - Average number of instructions per taken branch?
 - Assume: 20% branches, 50% taken → ~10 instructions
- Consider a 5-instruction loop with an 4-issue processor
 - Without smarter fetch, ILP is limited to 2.5 (not 4, which is bad)

Computer Architecture | Prof. Milo Martin | Superscalar

19

Increasing Superscalar Fetch Rate



- Option #1: over-fetch and buffer
 - Add a queue between fetch and decode (18 entries in Intel Core2)
 - Compensates for cycles that fetch less than maximum instructions
 - "decouples" the "front end" (fetch) from the "back end" (execute)
- Option #2: "loop stream detector" (Core 2, Core i7)
 - Put entire loop body into a small cache
 - Core2: 18 macro-ops, up to four taken branches
 - Core i7: 28 micro-ops (avoids re-decoding macro-ops!)
 - Any branch mis-prediction requires normal re-fetch
- Other options: next-next-block prediction, "trace cache"

Computer Architecture | Prof. Milo Martin | Superscalar

20

Multiple-Issue Implementations

Statically-scheduled (in-order) superscalar

- What we've talked about thus far
- + Executes unmodified sequential programs
- Hardware must figure out what can be done in parallel
- E.g., Pentium (2-wide), UltraSPARC (4-wide), Alpha 21164 (4-wide)

Very Long Instruction Word (VLIW)

- Compiler identifies independent instructions, new ISA
- + Hardware can be simple and perhaps lower power
- E.g., TransMeta Crusoe (4-wide)
- Variant: Explicitly Parallel Instruction Computing (EPIC)
 - A bit more flexible encoding & some hardware to help compiler
 - E.g., Intel Itanium (6-wide)

Dynamically-scheduled superscalar (next topic)

- · Hardware extracts more ILP by on-the-fly reordering
- Core 2, Core i7 (4-wide), Alpha 21264 (4-wide)

Computer Architecture | Prof. Milo Martin | Superscalar

21

Multiple Issue Redux

- Multiple issue
 - Exploits insn level parallelism (ILP) beyond pipelining
 - Improves IPC, but perhaps at some clock & energy penalty
 - 4-6 way issue is about the peak issue width currently justifiable
 - Low-power implementations today typically 2-wide superscalar
- Problem spots
 - N² bypass & register file → clustering
 - Fetch + branch prediction → buffering, loop streaming, trace cache
 - N² dependency check → VLIW/EPIC (but unclear how key this is)
- Implementations
 - Superscalar vs. VLIW/EPIC

[spacer]