### Advanced Computer Architecture Chapter 7:

### **Multi-threading**



These lecture notes are partly based on the course text, Hennessy and Patterson's Computer Architecture, a quantitative approach (3<sup>rd</sup> and 4<sup>th</sup> eds), and on the lecture slides of David Patterson and John Kubiatowicz's Berkeley course

#### Instruction Issue

#### Chip Multiprocessor



Reduced function unit utilization due to dependencies Superscalar Issue



Superscalar leads to more performance, but lower utilization

#### **Predicated Issue**



Adds to function unit utilization, but results are thrown away



Limited utilization when only running one thread

#### Fine Grained Multithreading



Intra-thread dependencies still limit performance

#### Simultaneous Multithreading



Maximum utilization of function units by independent operations



Adds to function unit utilization, but results are thrown away

#### **Basic Out-of-order Pipeline**





#### 🕨 Alpha 21464

One CPU with 4 Thread Processing Units (TPUs)

"6% area overhead over singlethread 4issue CPU"

### **SMT** Pipeline



#### Multiprogrammed workload



#### ng\_technology.pdf

Application

# SMT performance

Description

crash simulation.

pharmaceuticals. Mechanical Design Analysis

Genetics

genes.

flows.

omputational Chemistry

Mesoscale Weather Modeling

regional-scale atmospheric circulation.

omputational Fluid Dynamics

geophysical engineering applications.

studies and computational fluid dynamics.

Finite Element Analysis

Finite Element Analysis

**Genetics** 

Mechanical Design Analysis (finite element method)

This application is used for metal-forming, drop testing, and

genetics application that correlates DNA samples from multiple animals to better understand congenital diseases.

This application uses the self-consistent field method to

compute chemical properties of molecules such as new

This application simulates the metal-stamping process.

This application simulates and predicts mesoscale and

This application is designed to generate Expressed Sequence Tags (EST) clusters, which are used to locate important

This application is used to model free-surface and confined

This finite element application is specifically targeted toward

This explicit time-stepping application is used for crash test

Ipha 21464 🕨



**FP/SIMD** execution cluster

http://www.tomshardware.co.uk/intel-atom-cpu,review-30931-5.html, adgets/2008/02/sm https://ar

### Each thread runs slow?

## **SMT** issues

The point of Simultaneous Multithreading is that resources are dynamically assigned, so if only one thread can run it can run faster

#### SMT threads contend for resources

- Possibly symbiotically?
  - One thread is memory-intensive, one arithmetic-intensive?
- Possibly destructively
  - thrashing the cache? Other shared resources.... (TLB?)
- Which resources should be partitioned per-thread, and which should be shared on-demand?

### SMT threads need to be scheduled *fairly*

- Can one thread monopolise the whole CPU?
  - Denial of service risk
  - Slow thread that suffers lots of cache misses fills RUU and blocks issue

### Side channels:

one thread may be able observe another's traffic and deduce what it's doing

# SMT – latency-hiding

### SMT threads exploit memory-system parallelism

- Easy way to get lots of memory accesses in-flight
- "Latency hiding" overlapping data access with compute

# What limits the number of threads we can have?

## SMT threads need a *lot* of registers

A lot of logical registers – but they share physical registers?

### In a machine without register renaming

- What about statically partitioning the register file based on the number of registers each thread actually needs?
- This is what many GPUs do
- Leads to tradeoff: lots of lightweight threads to maximise latency hiding? Or fewer heavyweight threads that benefit from lots of registers?
- Nvidia and AMD call this "occupancy"

# **Mapping threads into the register file**<sup>10</sup>



- If each thread needs few registers, we can have lots of them co-existing in the same physical register file
- Alternatively, we could have fewer, fatter threads
- More threads=higher "occupancy"
- Better latency hiding
- Tricky tradeoff!



#### 





ultiprocessor 32 Physical Max Warps/SM = 48 Occupancy = 48 / 48 = 100%

CUDA Occupancy Calculator	
Version:	11.1
Copyright and License	

Note: Occupancy limiter is shown in orange

20

# Chapter summary

### We have explored:

- Pipeline parallelism
- Dynamic instruction scheduling
- Static instruction scheduling
- Multiple instructions per cycle
- Very long instruction words (VLIW)
- Multi-threading
  - Coarse-grain
  - Fine-grain
  - Simultaneous multithreading (SMT)
  - Statically-partitioned multithreading

Vector instructions and SIMD – coming soon

SIMT and GPUs – coming soon

Multicore – coming soon

# **Extra slides for interest/fun**<sup>22</sup>

Is the "minimum" operator associative?

immin(X, Y) = if X<Y then X else Y</pre>

min(min(10, x), 100) = 100

# **Extra slides for interest/fun**<sup>23</sup>

Is the "minimum" operator associative?

```
>>min(X, Y) = if X<Y then X else Y</pre>
```

All comparisons on NaNs always fail.... min(min(10, NaN), 100) = 100

# **Extra slides for interest/fun**<sup>24</sup>

Is the "minimum" operator associative?

>>min(X, Y) = if X<Y then X else Y</pre>

All comparisons on NaNs always fail....

>>min(X, NaN) = NaN
>>>min(NaN, Y) = Y

min(min(10, NaN), 100) = 100