# The Intel® Pentium® 4 Processor

**Doug Carmean**
**Sr. Principal Architect**
**Intel Architecture Group**

**April 26, 2005**
**Cornell**

Intel®

PDX

# The Intel® Pentium® 4 Processor

**Doug Carmean**
**Sr. Principal Architect**
**Intel Architecture Group**

**April 26, 2005**
**Cornell**

Intel®
PDX

# Agenda

- Review
- Pipeline Comparison
- Frequency Scaling
- Microarchitecture Description
- Analyzing Performance
- A Retrospective
- Summary

Intel®
PDX

# History



Performance

29K T
8086
1978

134K T
80286
1982

275K T
386
1986

1.2M T
486
1989

3.1M T
Pentium®*
1992

5.5M T
Pentium® Pro*
1995

42M T
Pentium® 4*
2000

1980    1985    1990    1995    2000

Intel®
PDX

# Intel® Netburst™ Micro-architecture vs P6

## Basic P6 Pipeline

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | | | 0 |
|---|---|---|---|---|---|---|---|---|---|
| Fetch | Fetch | Decode | Decode | Decode | Rename | ROB Rd | Rdy/Sc | | Exec |

**Intro at 733MHz .18µ**

## Basic Pentium® 4 Processor Pir

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 1. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TC Nxt IP | | TC Fetch | | Drive | Alloc | Rename | | Que | Sch | Sch | Sch | Disp | Disp | R. |

**Intro at 1.5GHz .18µ**

## Deeper Pipelines enable higher frequency and performance

intel™
PDX
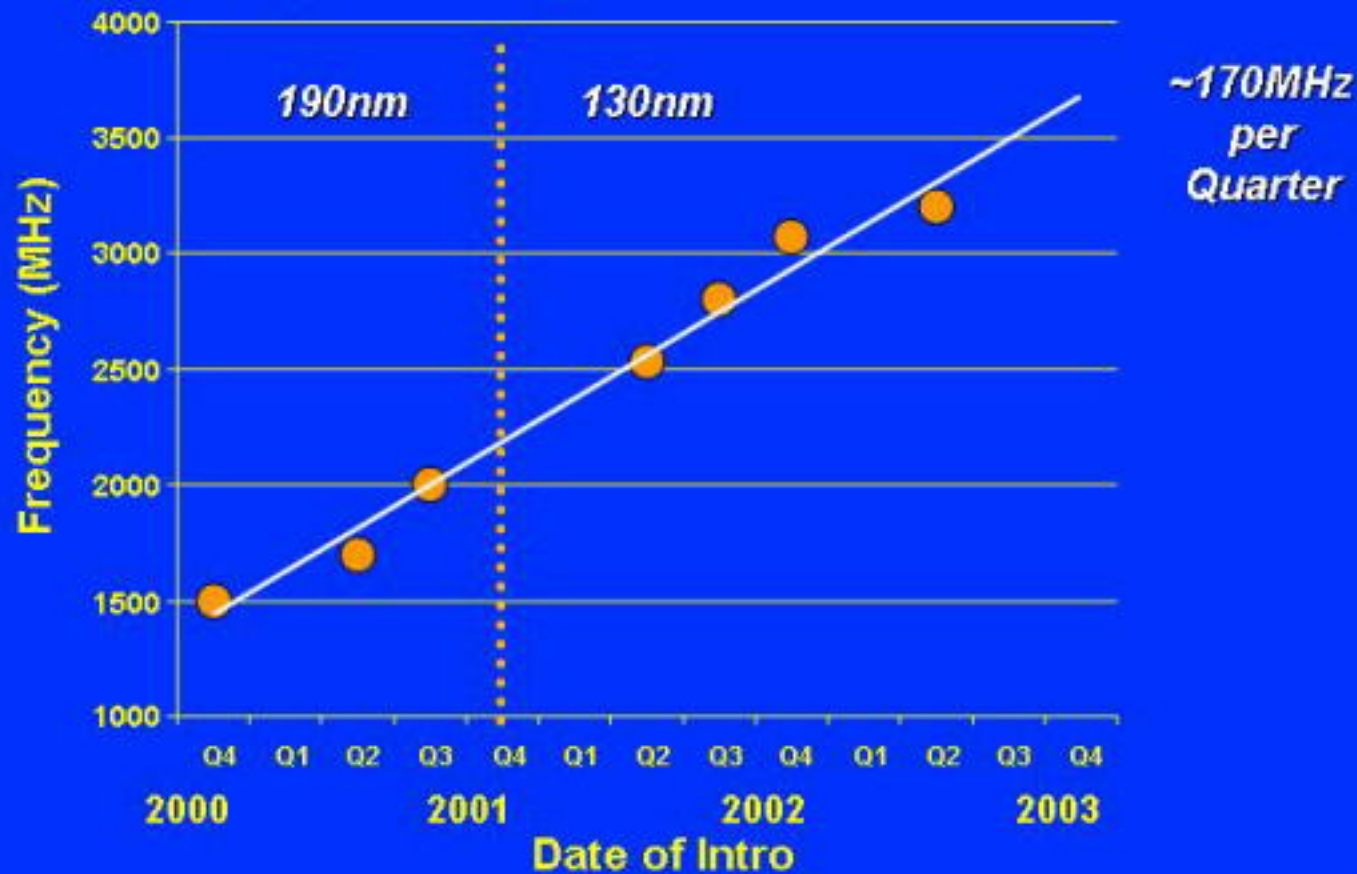
# Pipeline Comparison

- **P6:**
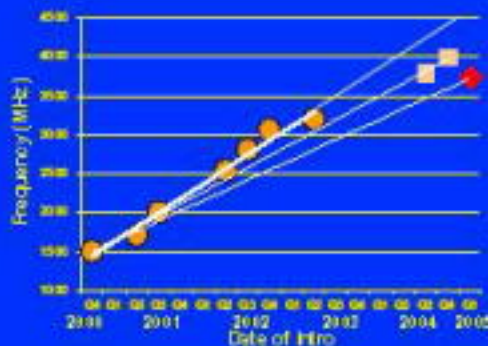  - **10 clocks @ 1.7GHz**

- **P4P:**
  - **20 clocks @ 3.2GHz**

6%

5.88nS

6.25nS

**Branch misprediction penalty remains similar**

Intel®

PDX

# Frequency Scaling

190nm | 130nm

~170MHz per Quarter

Frequency (MHz)

Date of Intro

**Solid Frequency Scaling Through 2003 (130nm process)**

# Frequency Scaling



**Original Ramp Rate**
- 170MHz per Q
- >4.5GHz in 2005
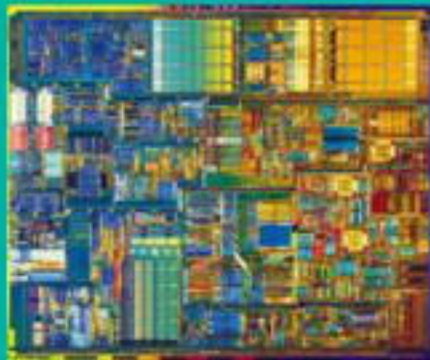- Circuit Tweaking
- Process Polishing

**Revised in 2003**
- 150MHz per Q
- 4GHz in 2005
- Design Sched Slip

**Revised in 2004**
- 130MHz per Q
- 3.7GHz in 2005
- More Sched Slip

## Schedule Slips Result in Significantly Lower Performance
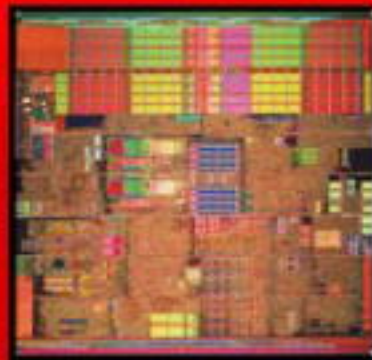
# The Family Tree



## Willamette
- 190nm
- 8KB L1
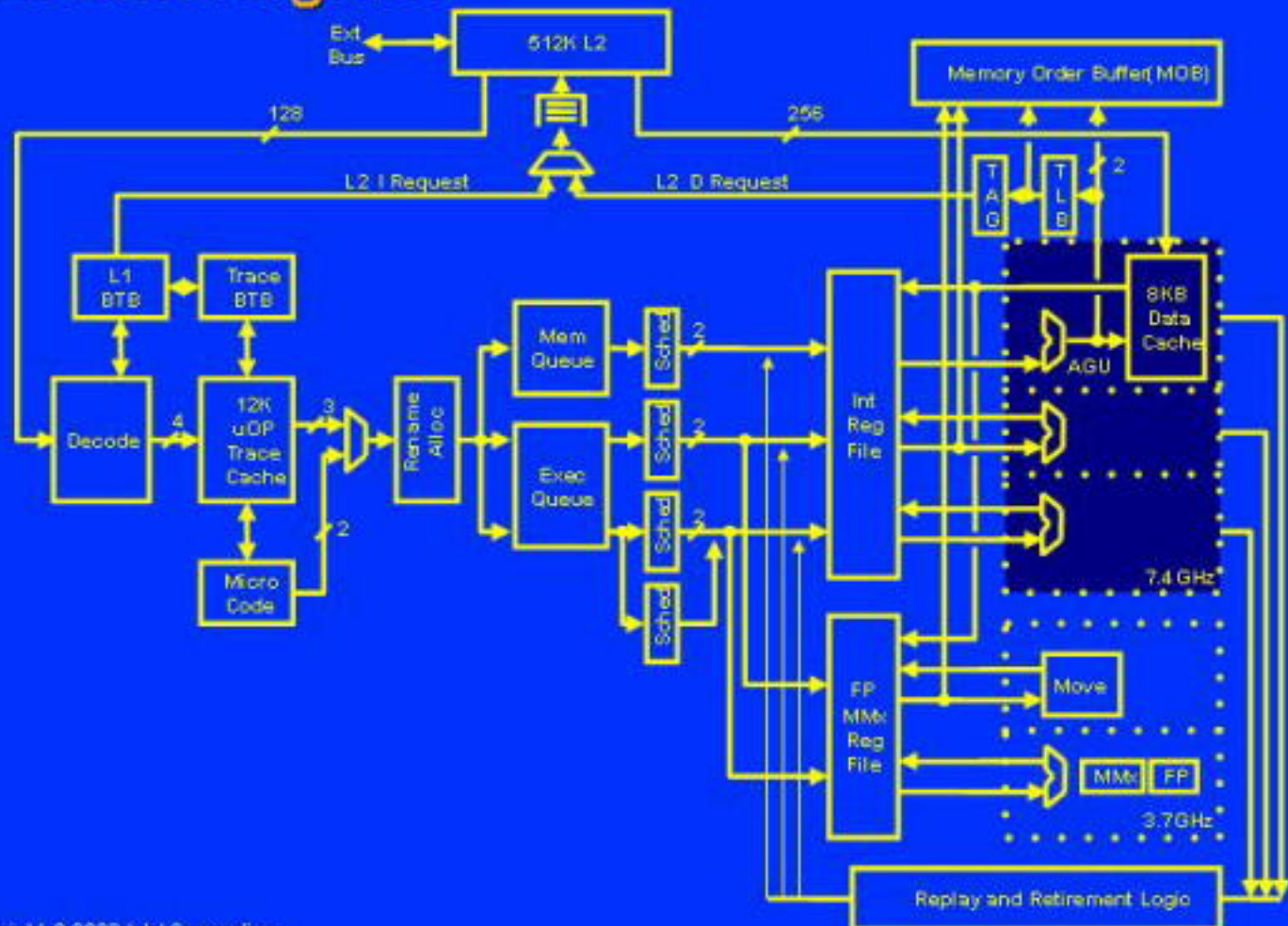- 256KB L2
- HT Capable

## Northwood
- 130nm
- 8KB L1
- 512KB L2
- HT Production
- Minor perf fixes

## Prescott
- 90nm
- 16KB L1
- 1MB L2
- 64-bit
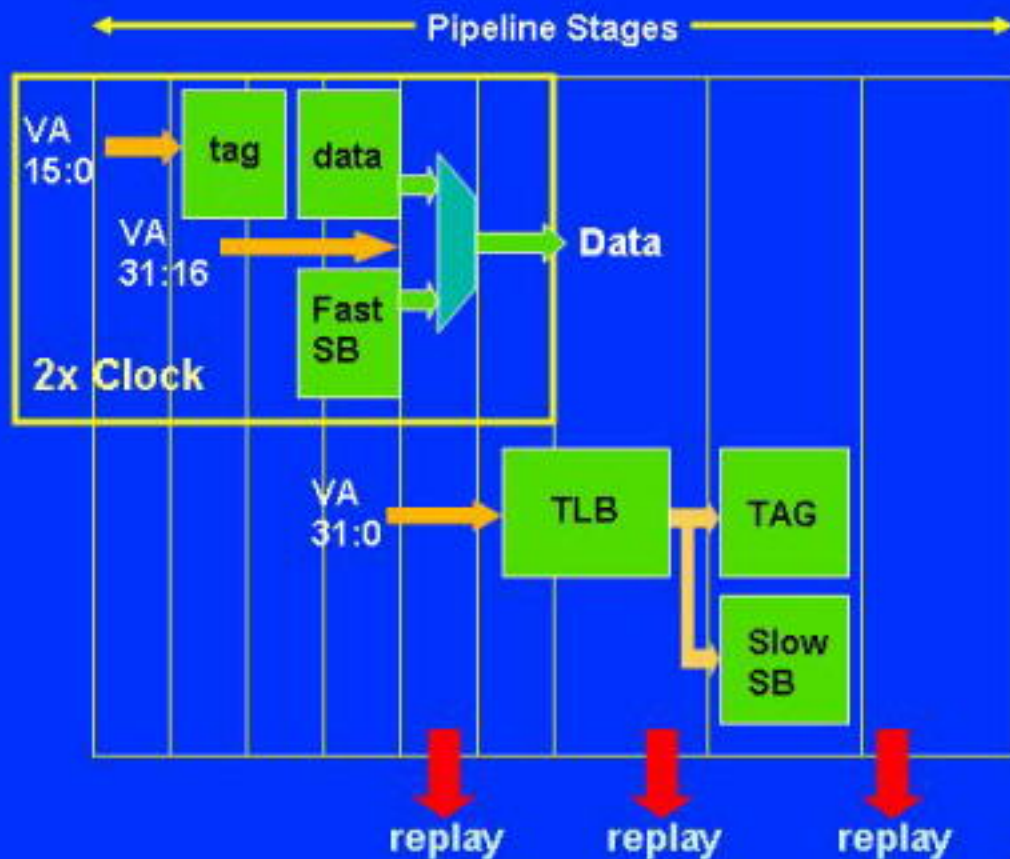- Shift/Mul fixes
- Thread instrs

# Block Diagram

512K L2

Ext Bus

Memory Order Buffer (MOB)

128

256

L2 I Request

L2 D Request

TAG

TLB

2

L1 BTB

Trace BTB

8KB Data Cache

Decode

12K uOP Trace Cache

Rename Alloc

Mem Queue

Sched

2

Int Reg File

AGU

4

3

Exec Queue

Sched

2

Sched

2

7.4 GHz

Micro Code

2

Sched

FP MMx Reg File

Move

MMx

FP

3.7GHz

Replay and Retirement Logic

# Agenda

- Review
- Pipeline Comparison
- Frequency Scaling
- <u>Microarchitecture Description</u>
- Analyzing Performance
- A Retrospective
- Summary

Intel® PDX

# L1 Data Cache

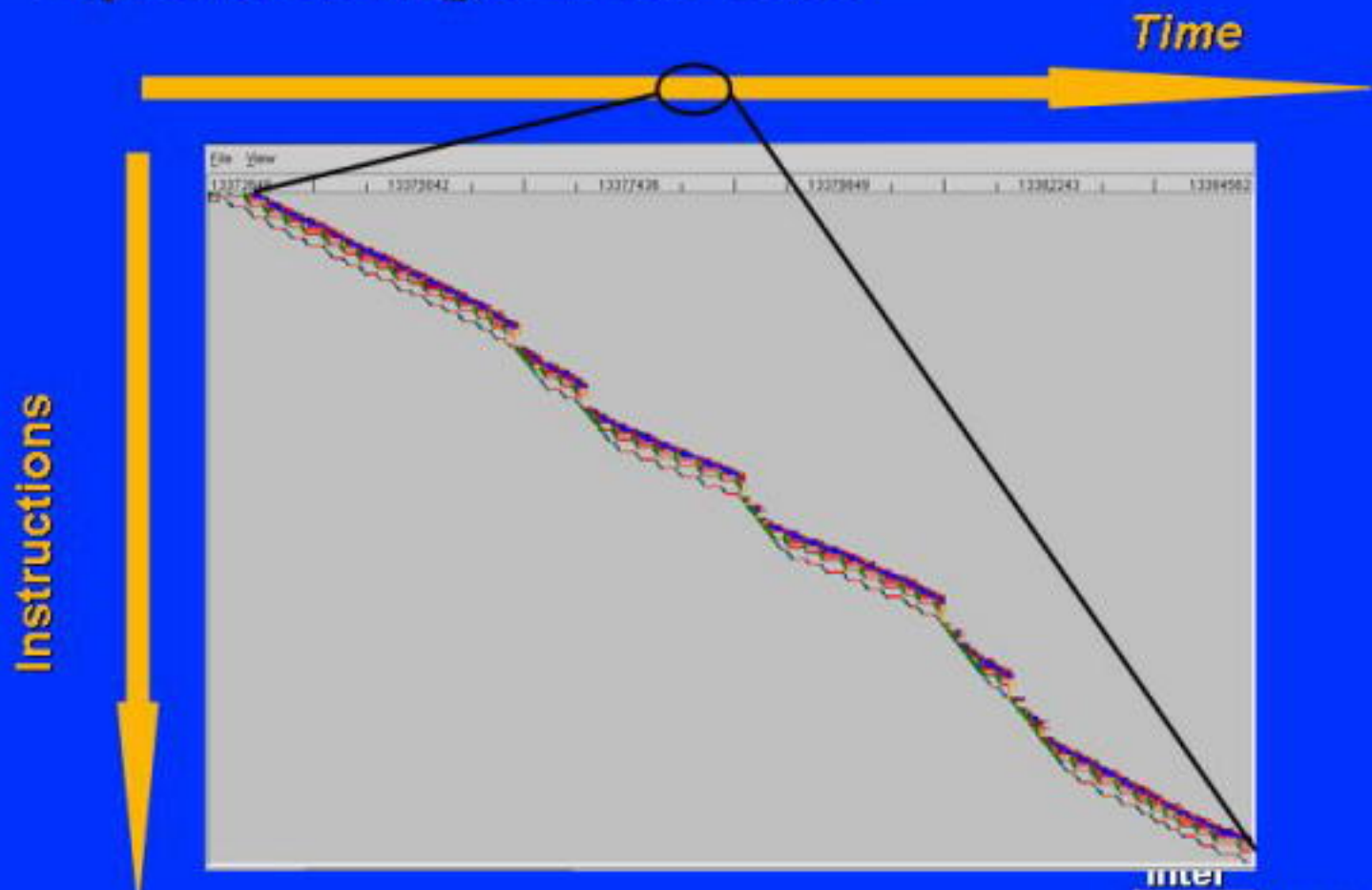# Agenda

- Review
- Pipeline Comparison
- Frequency Scaling
- Microarchitecture Description
- Analyzing Performance
- A Retrospective
- Summary

Intel®

PDX

# Analyzing Performance

- **Performance simulator**
  - Models low level microarchitecture details
  - Correlated to silicon
  - Tons of knobs (parameters) for various configurations
- **Application "Traces"**
  - Architectural state plus memory image
  - Multiple 30M instruction traces per application
  - A typical study runs on ~500 traces
  - An important study will run on ~1000 traces
- **Pipeline events**
  - Generates pipeline events (Fetch, Decode, ….)
  - Events are consistent between RTL & perf. Simulator
  - Graphical tool (ptv) used to analyze traces

Intel®

PDX

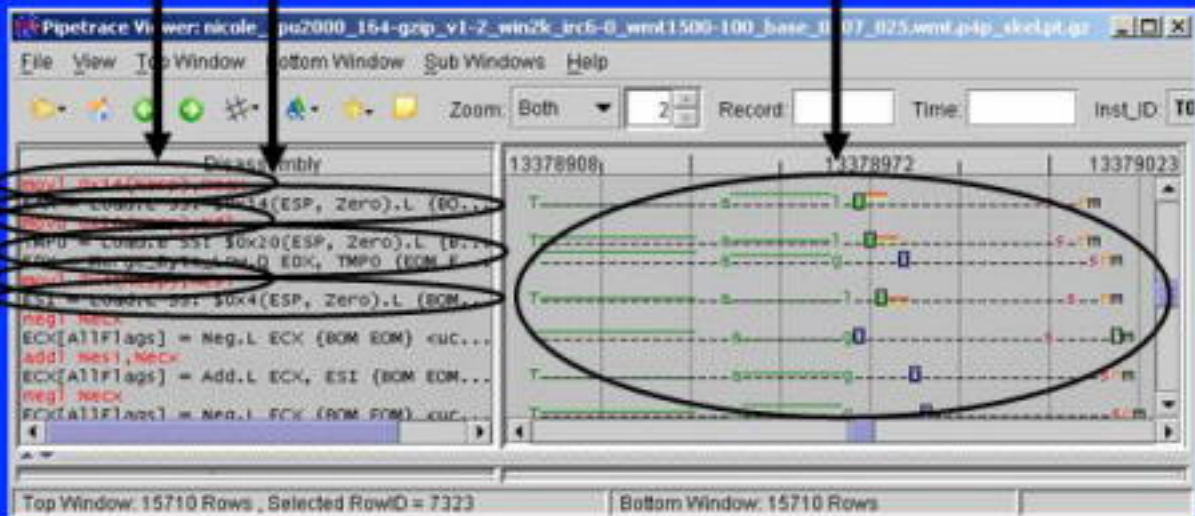# Pipetrace: High Level View



*Time*

Instructions

intel

PDX

# Pipetrace: uOP View Window



Copyright © 2003 Intel Corporation.

# Simple Example

# Replay loop

# A Well Behaved Application



TC Hitting Max Bandwidth per clock

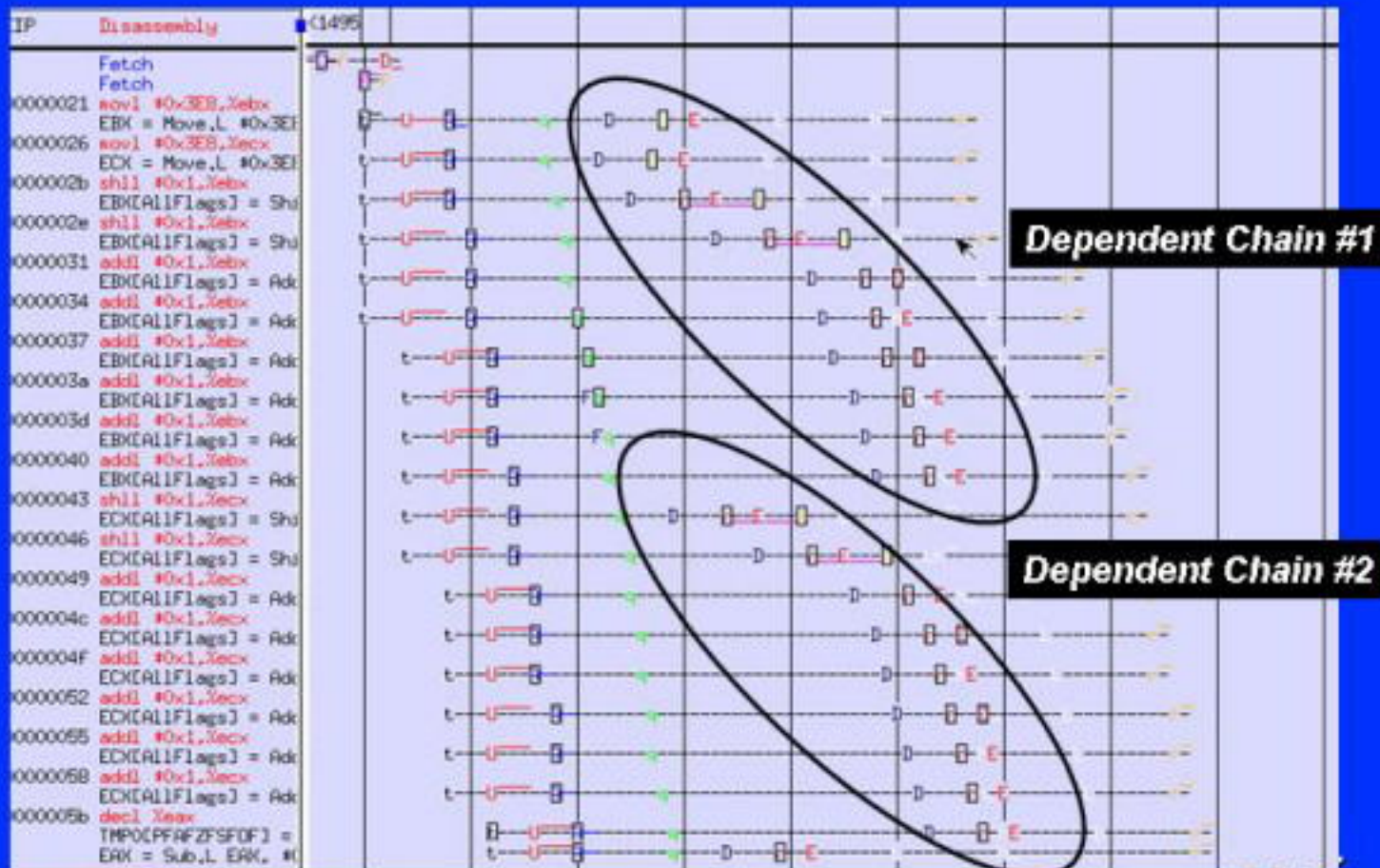Allocating 3 uOPs per clock

Good Execution

Retiring 3 uOPs per clock

Intel®

PDX

# Parallelism with dependencies

# Agenda

- Review
- Pipeline Comparison
- Frequency Scaling
- Microarchitecture Description
- Analyzing Performance
- A Retrospective
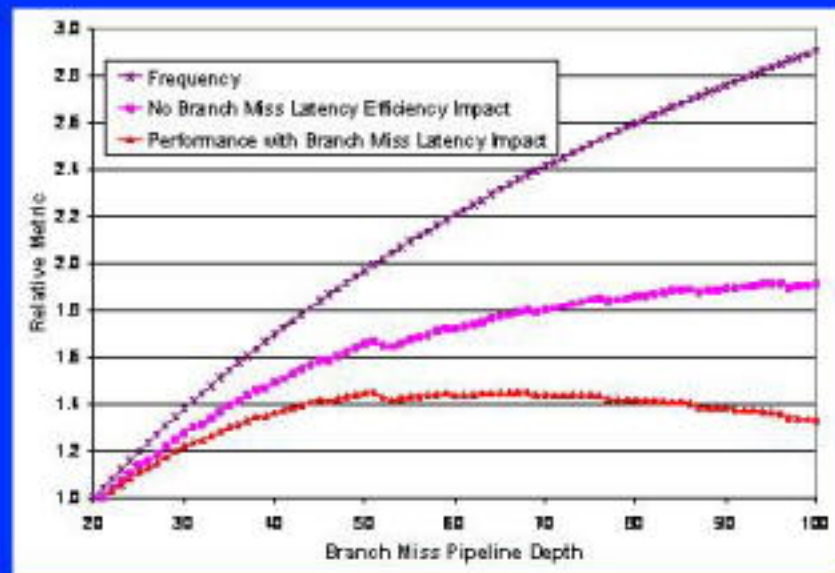- Summary

Intel®

PDX

# What *Were* We Thinking?

- **Need to improve traditional performance**
- **Start working on threaded performance**
- **The P6 family is out of gas**
- **Performance ~ Frequency**
- **Frequency Sells**
- **We can pipeline *anything***

*Improve Performance: Push IPC, Push Frequency,*
*Invent the new uArch for a family of processors,*
*Start down the SMT/CMP path*

# Increasing Performance by Implementing Deeper Pipelines and Larger Caches*
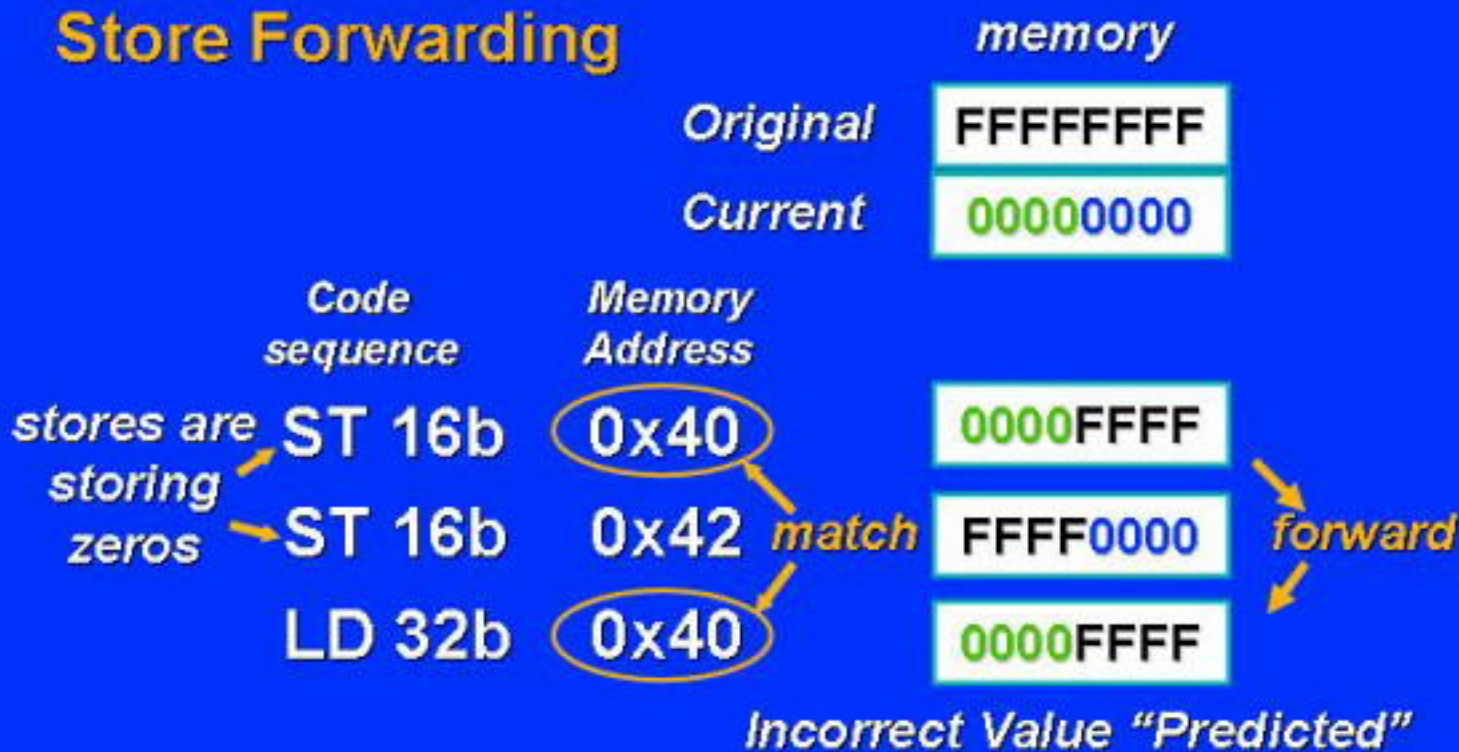
- **Performance vs. pipeline depth**
  - Cycle time = overhead time + useful time
  - Overhead time (latch + skew + jitter) is constant vs. pipeline depth
  - Useful time can be sliced up arbitrarily and is constant vs. pipeline depth
  - Use a simulator to calculate IPC vs. pipeline depth
  - The impact of increasing multiple pipelines is the product of impact of individual pipeline increases
  - Performance = Freq * IPC



Chart legend:
- Frequency
- No Branch Miss Latency Efficiency Impact
- Performance with Branch Miss Latency Impact

Y-axis: Relative Metric (1.0, 1.2, 1.4, 1.6, 1.8, 2.0, 2.2, 2.4, 2.6, 2.8, 3.0)
X-axis: Branch Miss Pipeline Depth (20, 30, 40, 50, 60, 70, 80, 90, 100)

**Intel®**  **PDX**

# But. . . as Pipelines get Deeper, Algorithms Change Fundamentally

- **Goal: Hold critical latencies to a constant number of *cycles* as frequency increases**
  - ALUs become staggered (16-bits per clock)
    - P4P ALUs are ½ cycle and running 2x frequency
    - Non critical functions moved - shifts take longer
  - Scheduler algorithms completely change
  - L1 data cache becomes smaller, virtually indexed
  - Data consumed on partial address matches
    - L1 data cache uses way predictors
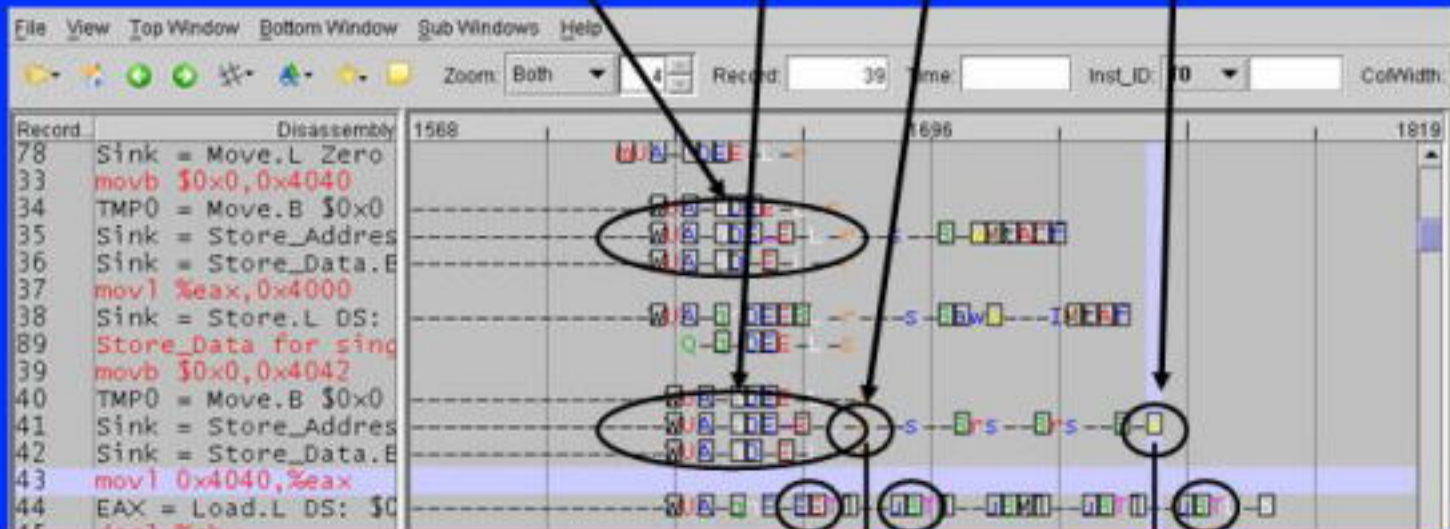    - Store forwarding uses subset of virtual addresses

Intel[®]

PDX

# Store Forwarding

memory

| | |
|---|---|
| Original | **FFFFFFFF** |
| Current | **00000000** |

| stores are storing zeros | Code sequence | Memory Address | | memory |
|---|---|---|---|---|
| → | ST 16b | (0x40) | | **0000FFFF** |
| → | ST 16b | 0x42 | *match* | **FFFF0000** → *forward* |
| | LD 32b | (0x40) | | **0000FFFF** |

Incorrect Value "Predicted"

Both stores have to commit
the data to the cache before
the load gets the right data

# Store forwarding



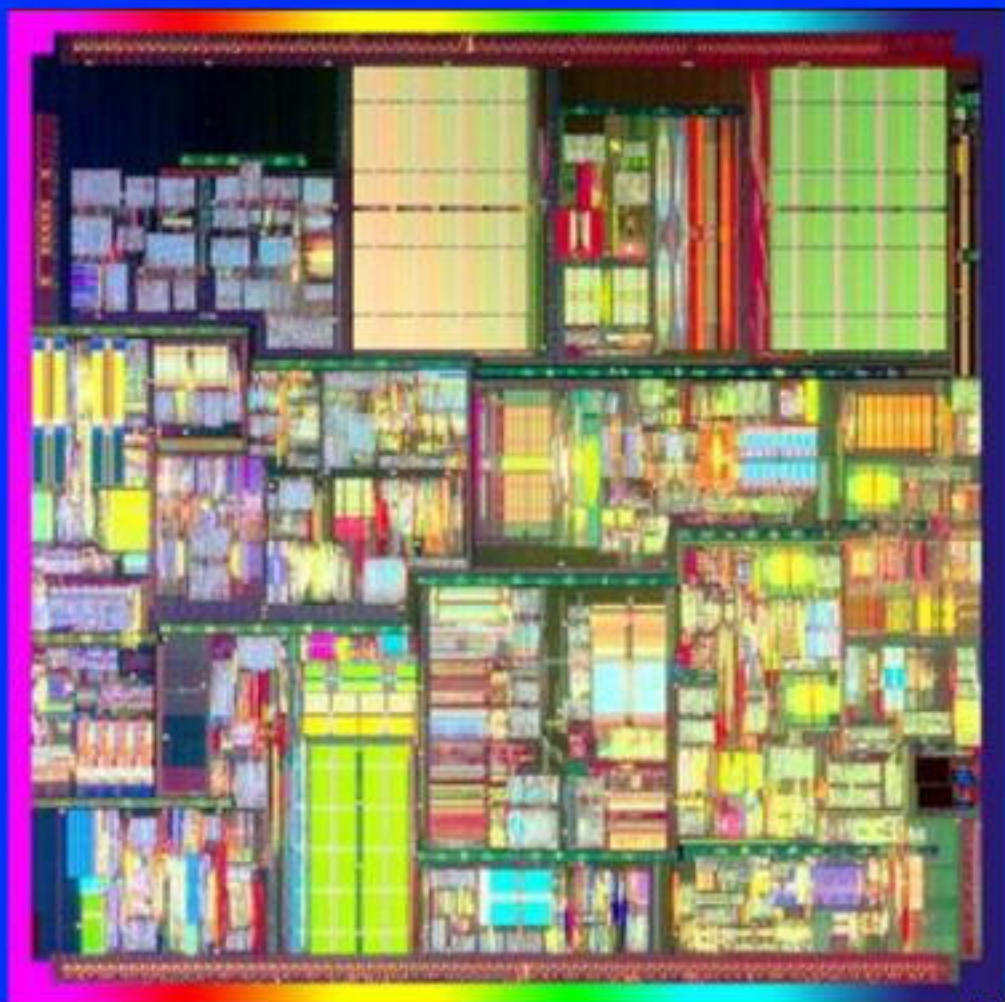Copyright © 2003 Intel Corporation.

**Pipeline Evolution**

Intel® PDX

# Debunking Deeper Pipelines

- **Deeper pipelines enable higher frequency, yielding higher performance**

- **But, additional algorithms introduce complexity which reduce potential performance gain**

- *should(?)*
- **But, engineering effort ~~can~~ overcome challenges represented by complexity increases**

# Summary

- **The Pentium® 4 Processor's deep pipelines provide high performance by enabling high frequency**

- **Introduced with, and maintained a 2x frequency advantage over the P6 pipeline**

- **Innovative algorithms introduced significant complexity, requiring huge effort to tune**

- **Evolution vs. Revolution – Evolution won this chapter**

Intel®

PDX

...tel®

PDX