

From Instruction Level To Thread Level Parallelism

John Paul Shen

**Director of Microarchitecture Research
Intel Labs**

October 29, 2002
Microprocessor Research Forum

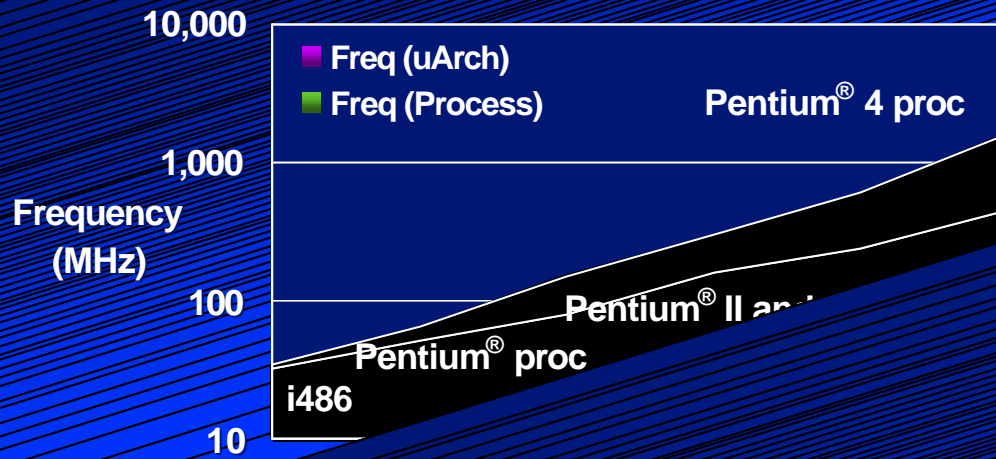


"Iron Law" of Microprocessor Performance

$$\begin{aligned} 1/\text{Processor Performance} &= \frac{\text{Time}}{\text{Program}} \\ &= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Time}}{\text{Cycle}} \\ &\quad \text{(inst. count)} \quad \quad \quad \text{(CPI)} \quad \quad \quad \text{(cycle time)} \end{aligned}$$

$$\text{Processor Performance} = \frac{\text{IPC} \times \text{GHz}}{\text{inst. count}}$$

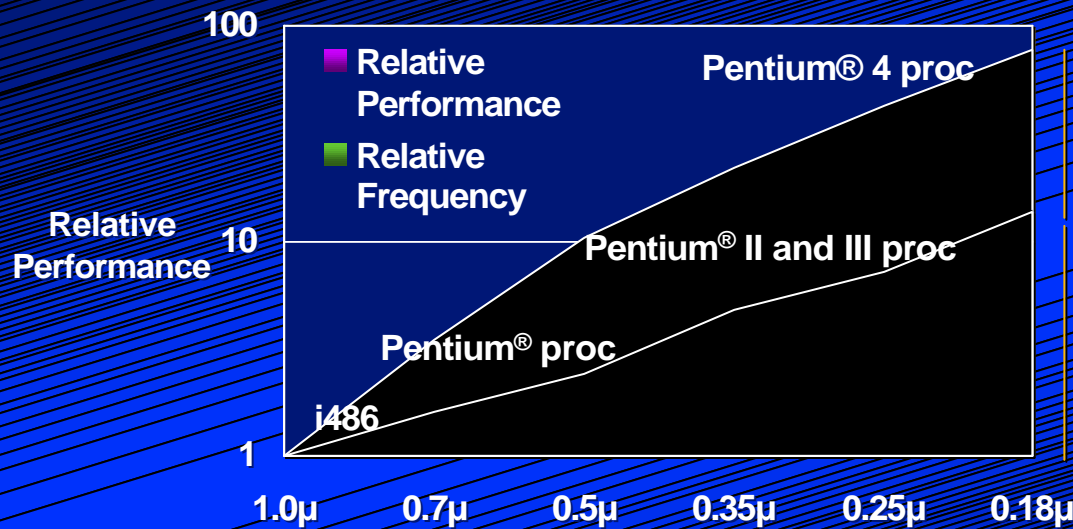
Frequency



Frequency Increased **50X**

4X

• 13X



Performance Increased **>75X**

6X

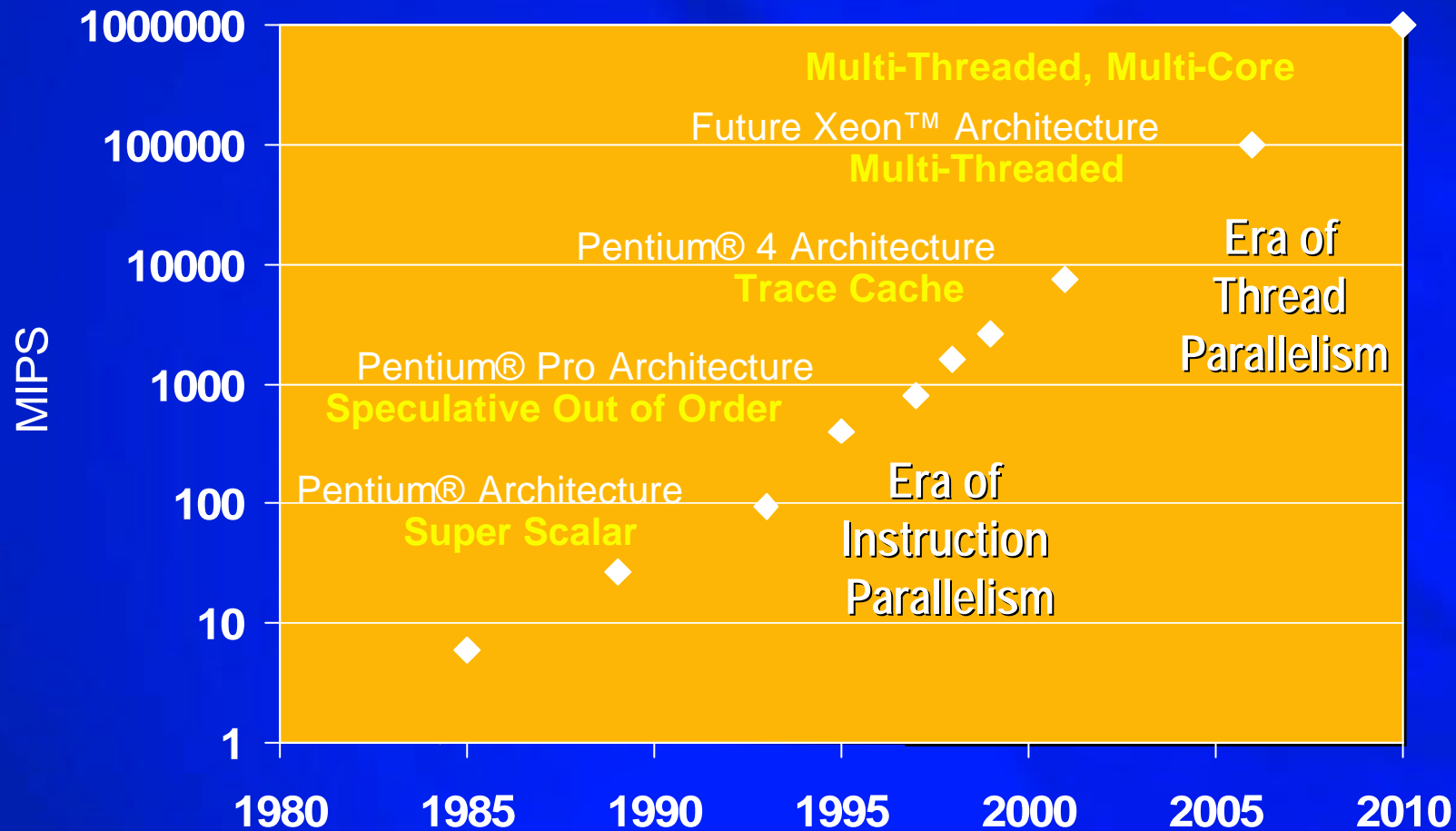
13X

- 13X due to process technology
- Additional >6X due to microarchitecture

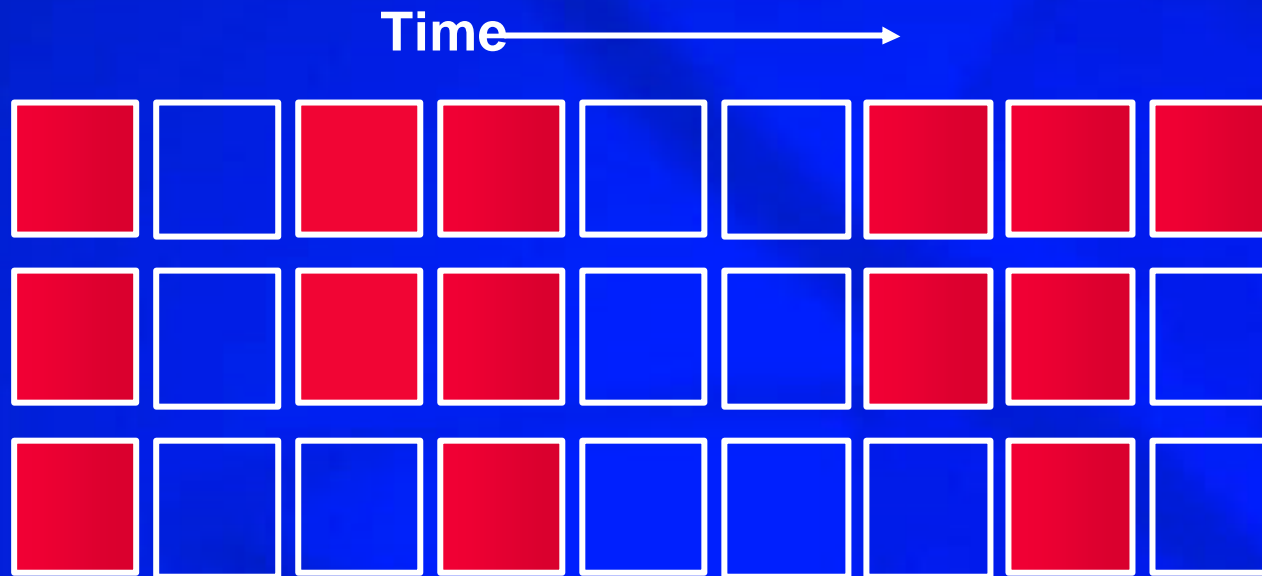
*Note: Performance measured using SpecINT and SpecFP

Source: Intel Corporation

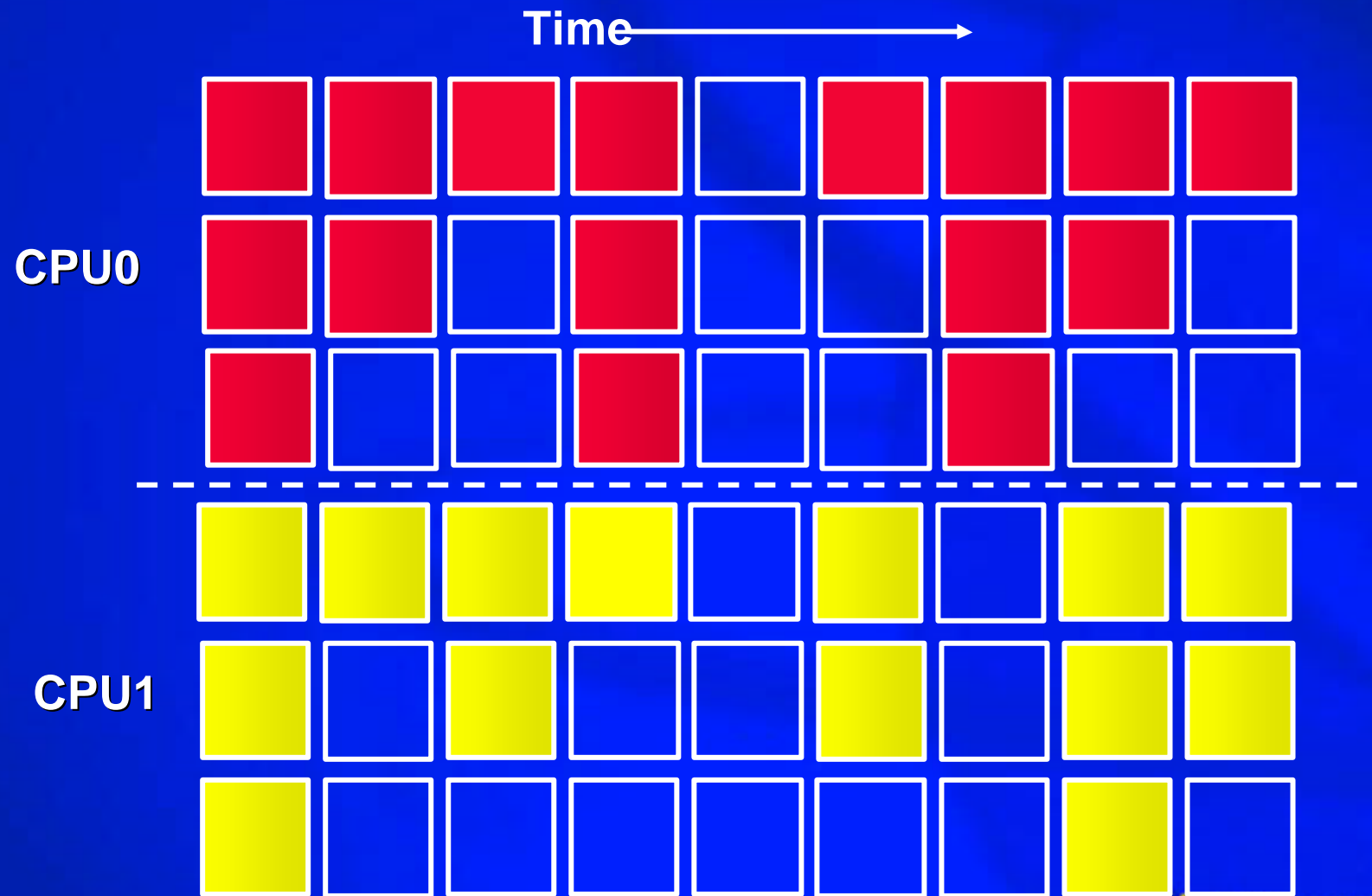
Parallelism in Transition



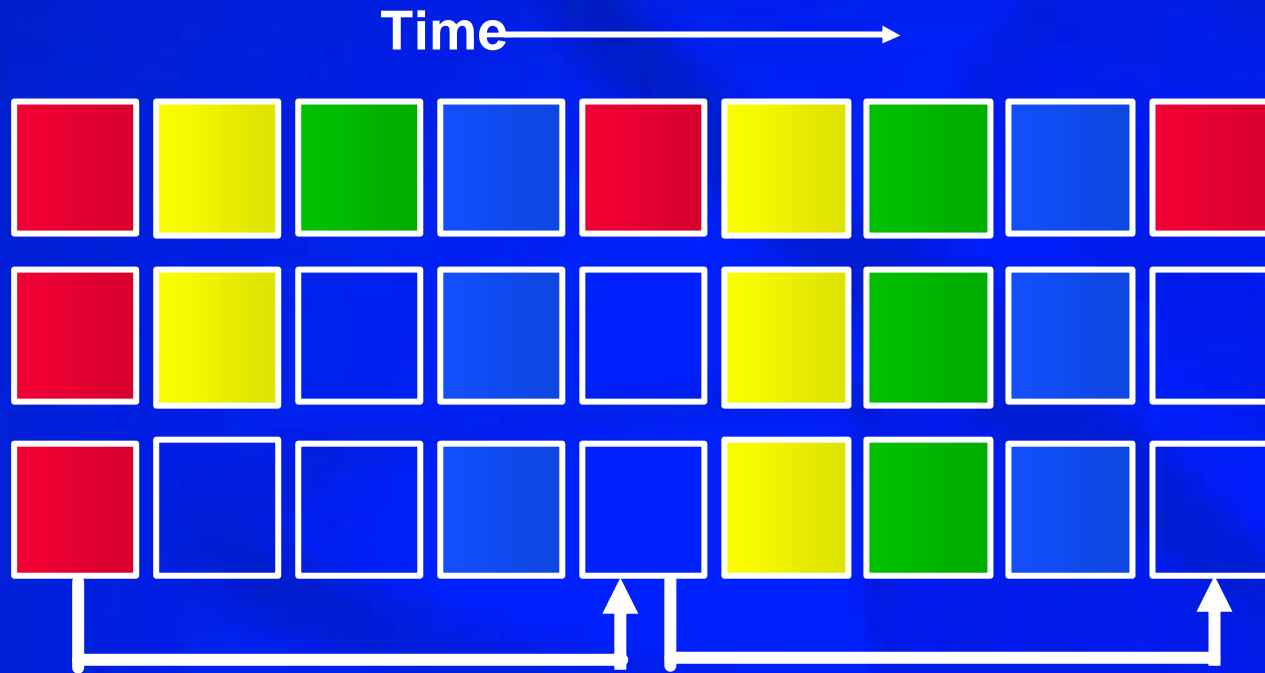
Superscalar Issue



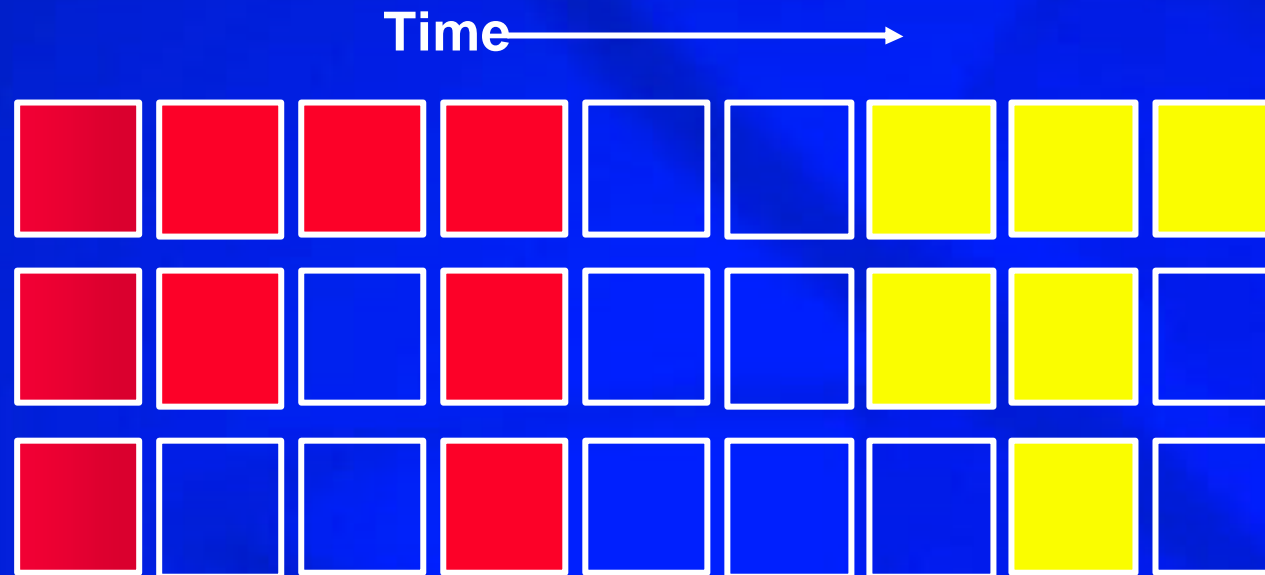
Chip Multiprocessor (CMP)



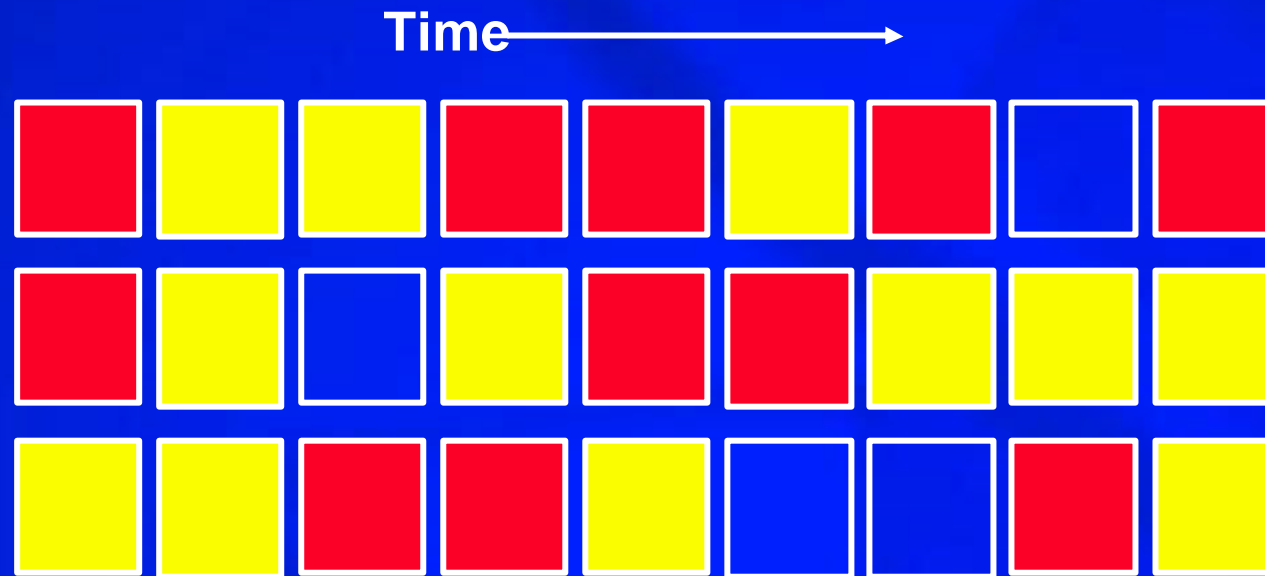
Time-slicing Multi-Threading



Switch-on-Event Multi-Threading

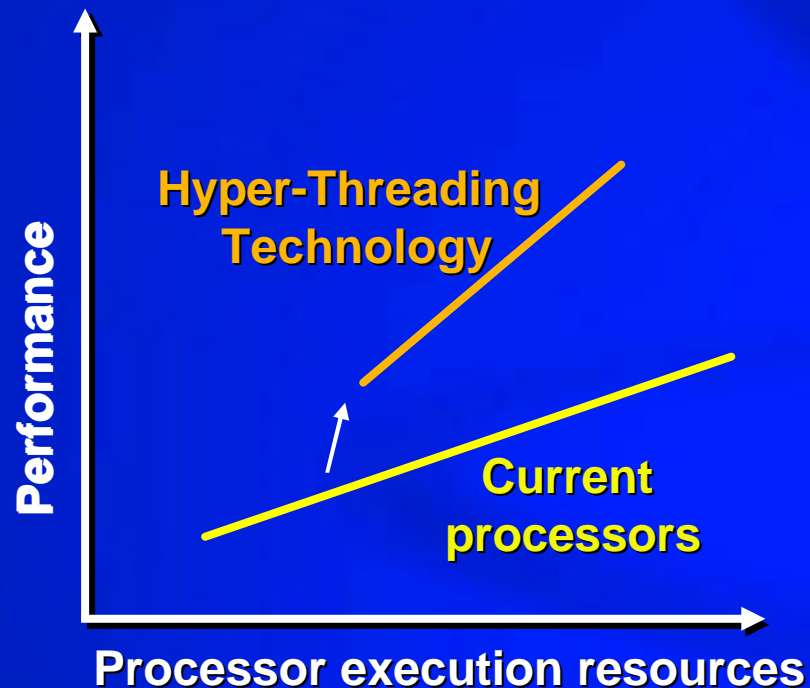


Simultaneous Multi-Threading (SMT)



Maximum utilization of function units by independent operations

Accelerate Performance of Threaded Applications



- ✍ SMT: most efficient/highest performance option
- ✍ More performance for CPU execution resources
 - Uses resources more fully
- ✍ Greater performance improvements from additional processor resources

Hyper-Threading Technology

✍ Executes two tasks simultaneously

- Two different applications
- Two threads of same application

✍ CPU maintains architecture state for two processors

- Two logical processors per physical processor

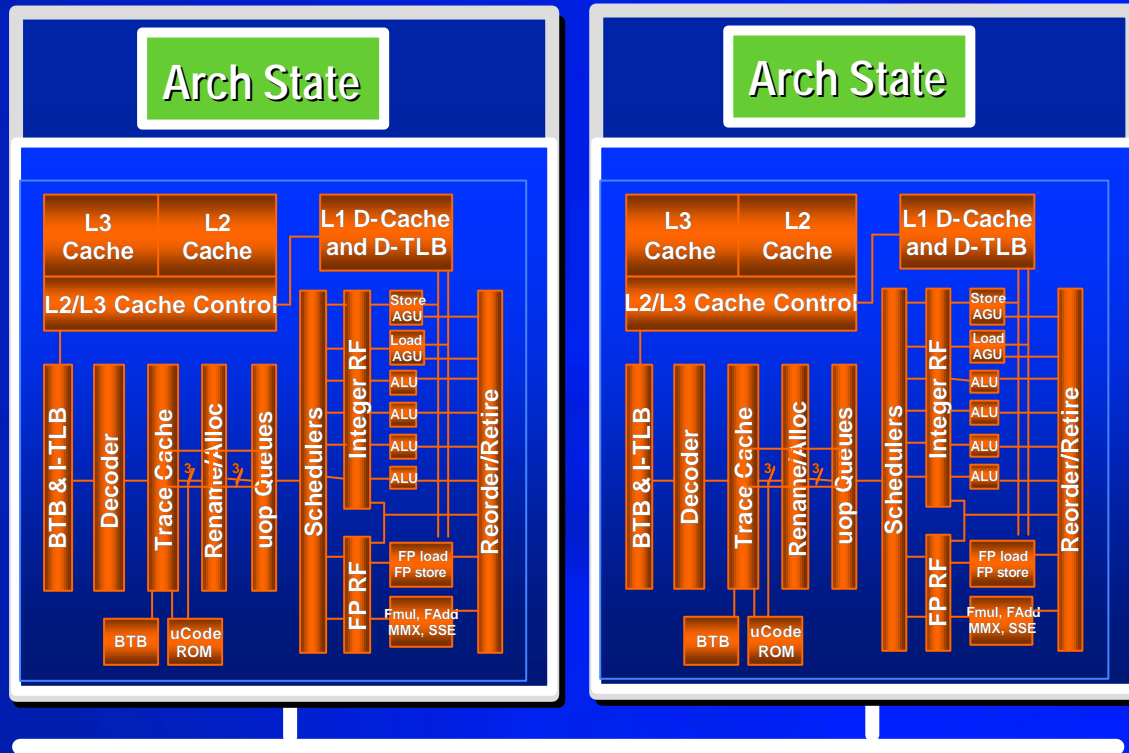
✍ Demonstrated on prototype Intel® Xeon™ Processor MP

- Two logical processors for < 5% additional die area
- Power efficient performance gain
- Result of significant research, design effort, and validation

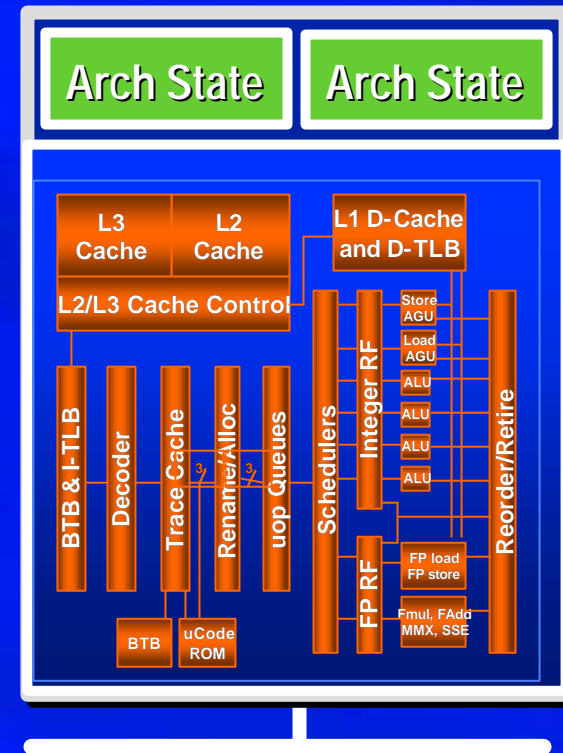


Replicated vs. Shared Resources

Multiprocessor



Hyper-Threading



Multi-Processors replicate execution resources
Hyper-Threading Technology shares resources

Changes for Hyper-Threading

Replicate resources

- All per-CPU architectural state
- Instruction Pointers, renaming logic
- Some smaller resources
 - E.g, return stack predictor, ITLB, etc

Partition resources

- Several buffers (Re-order buffer, load/store buffers, queues, etc)

Share most resources

- Out-of-Order execution engine
- Caches

What Was Added

Instruction Streaming Buffers

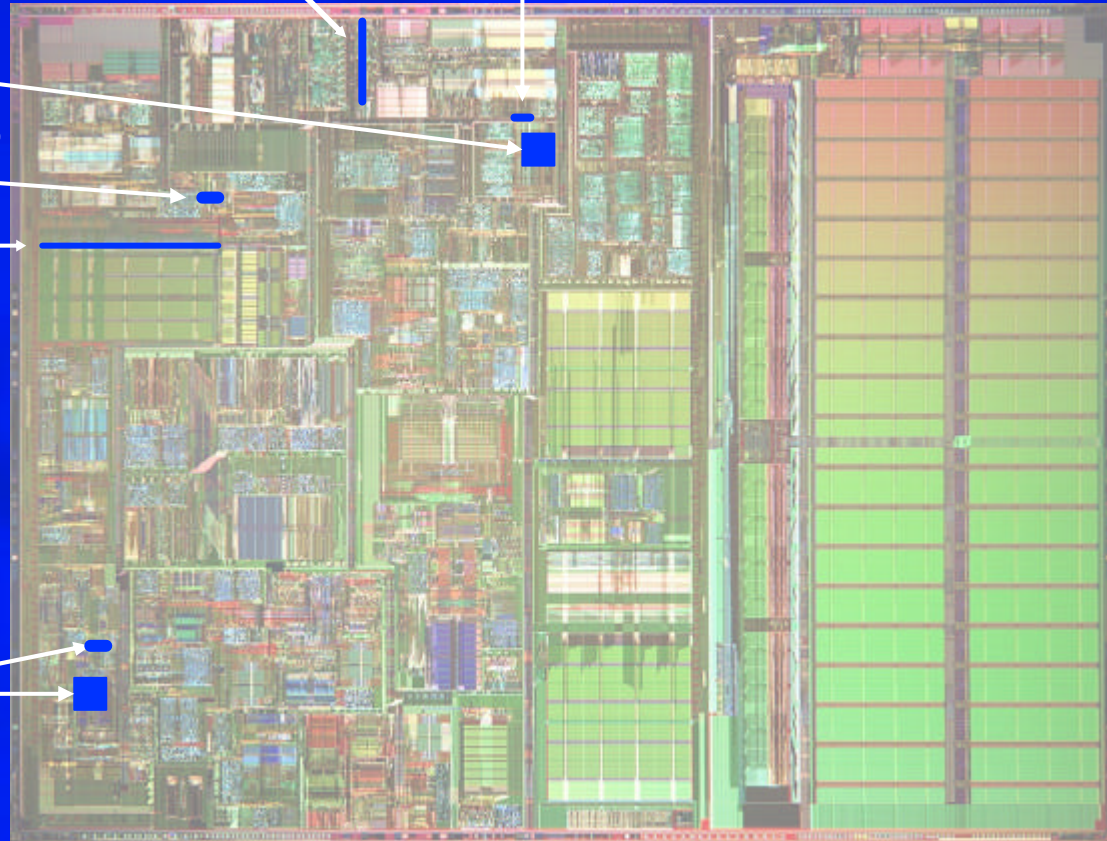
Next IP

Instruction TLB

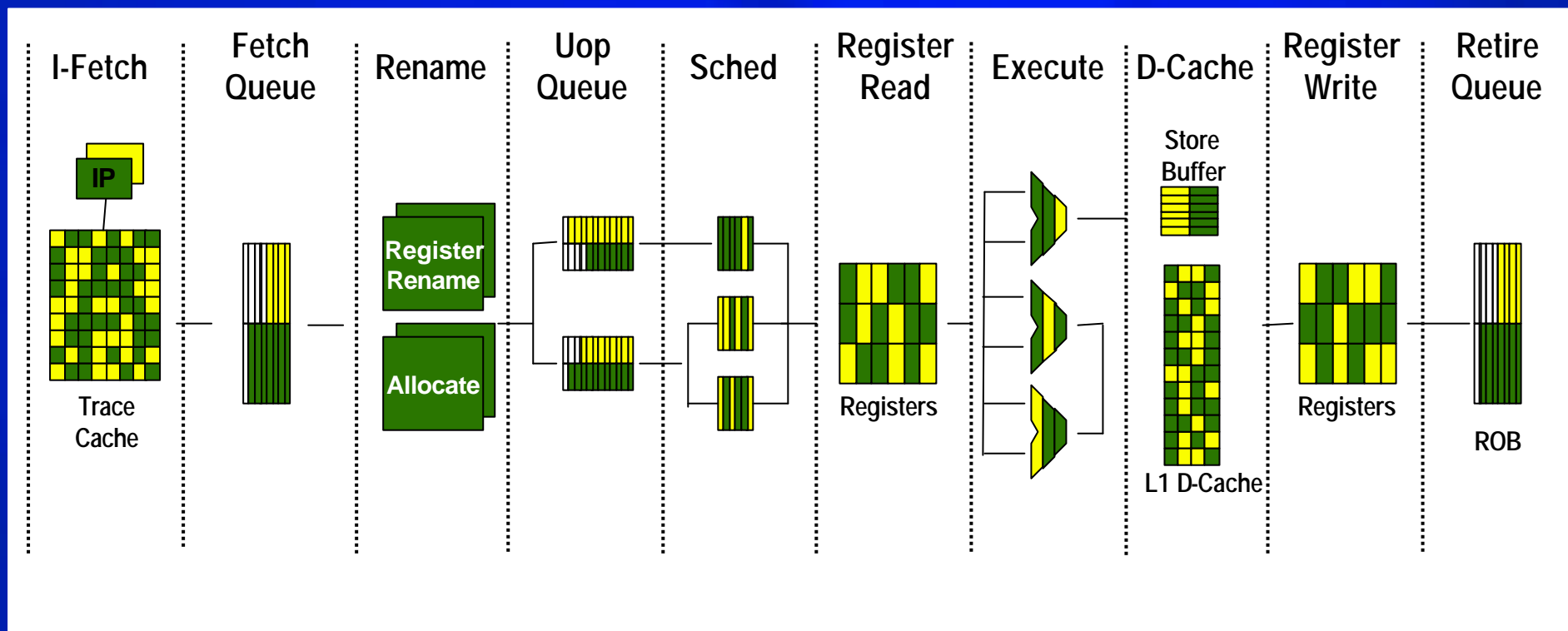
Trace Cache Next IP

Trace Cache
Fill Buffers

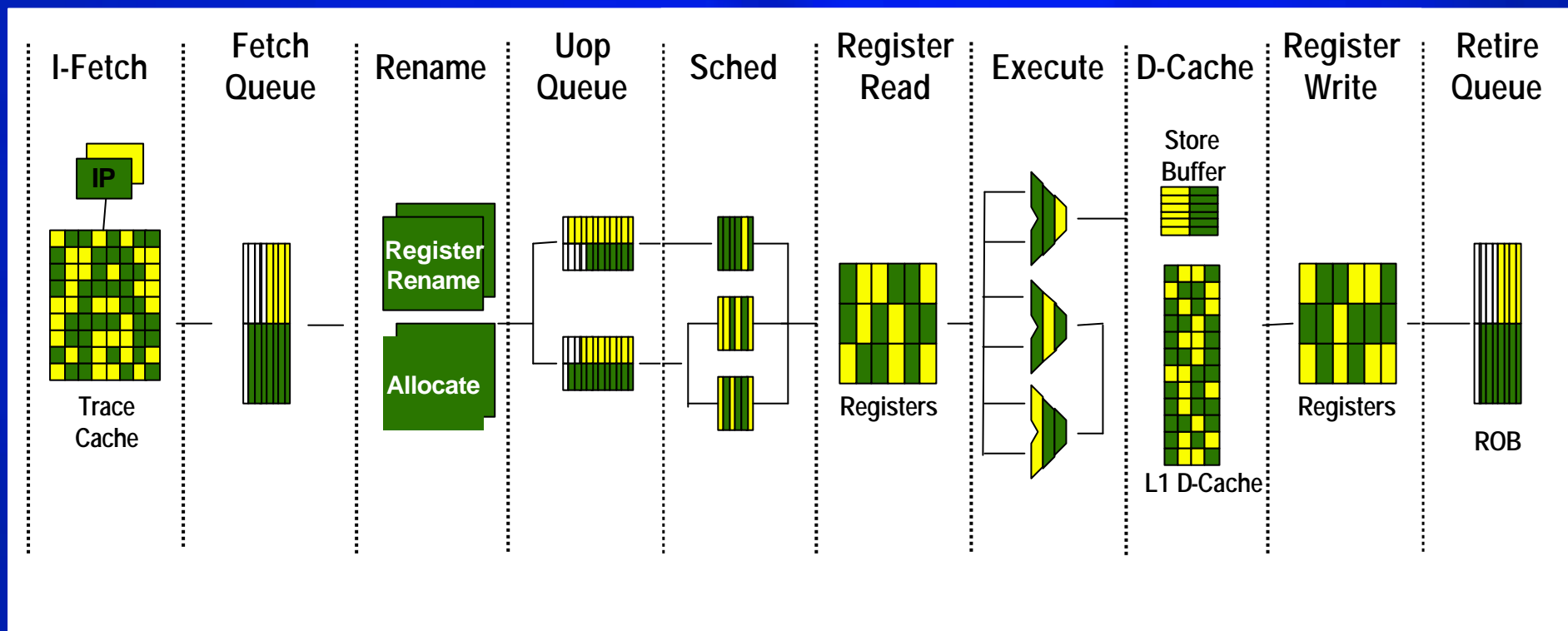
Register Alias
Tables



Execution Pipeline



Execution Pipeline



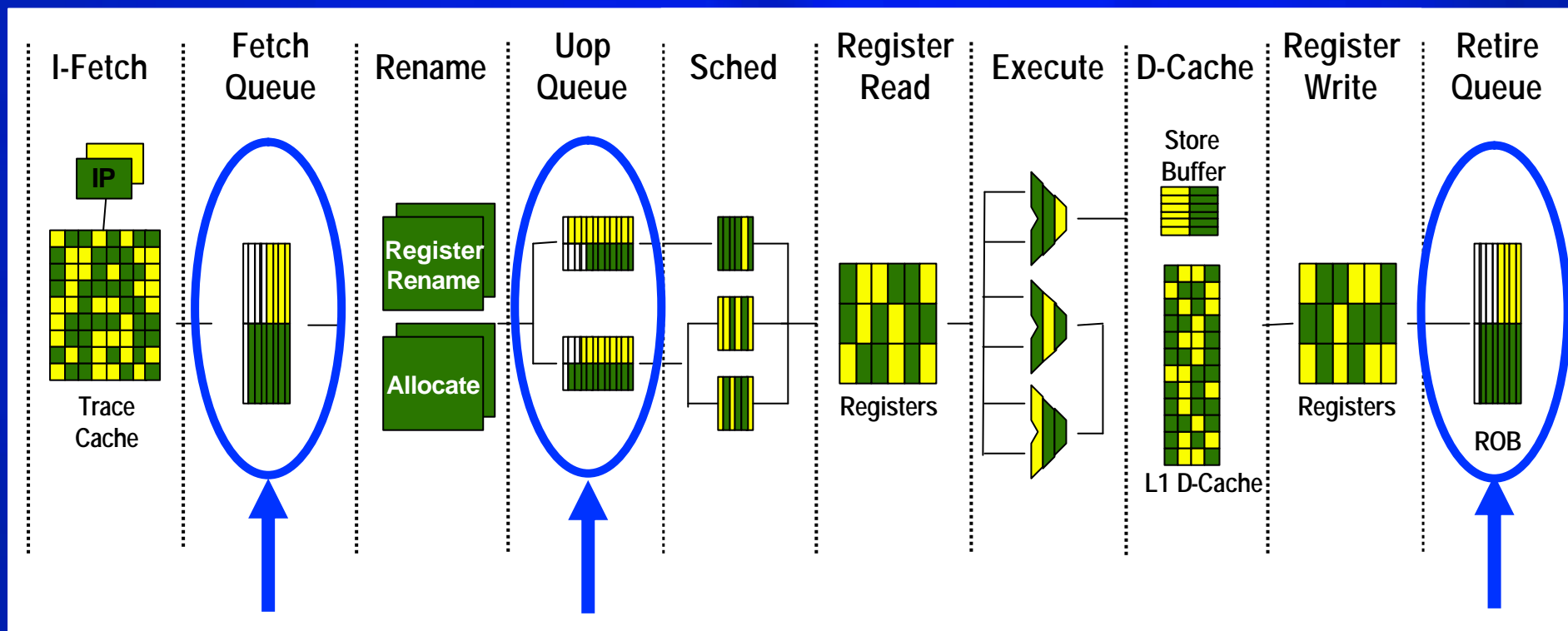
In-Order Pipeline

Out-of-Order Pipeline

In-Order

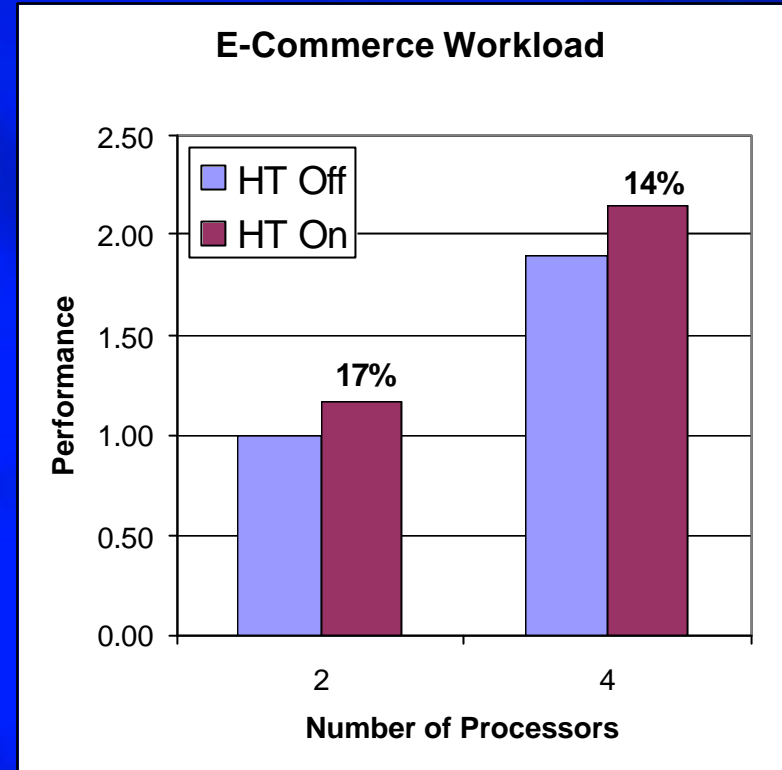
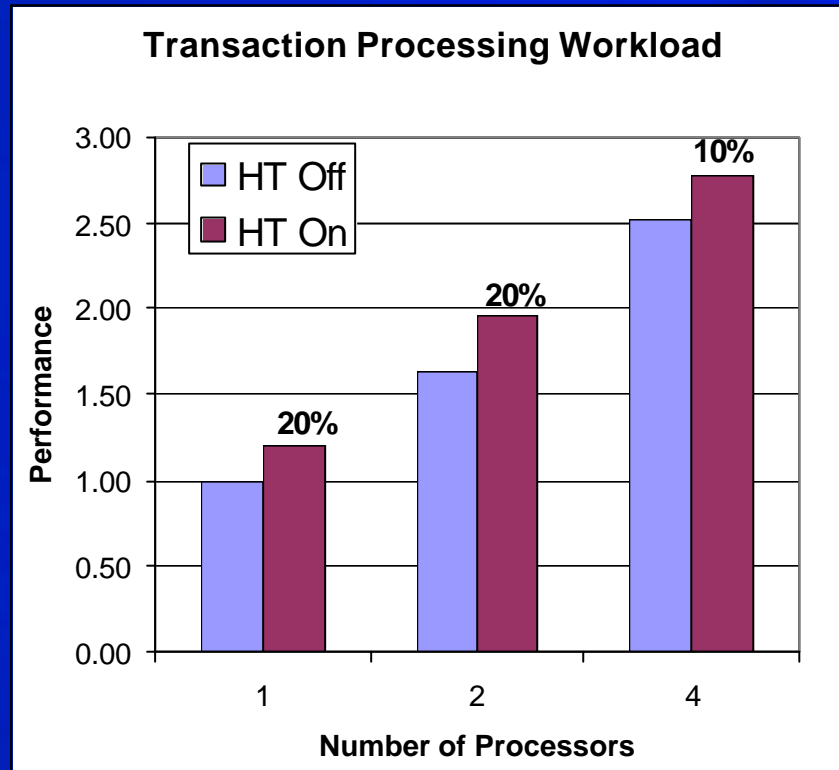


Execution Pipeline



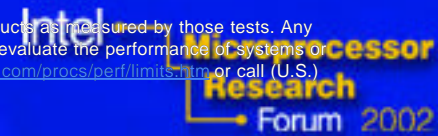
Partition queues between major pipestages of pipeline

Server Performance

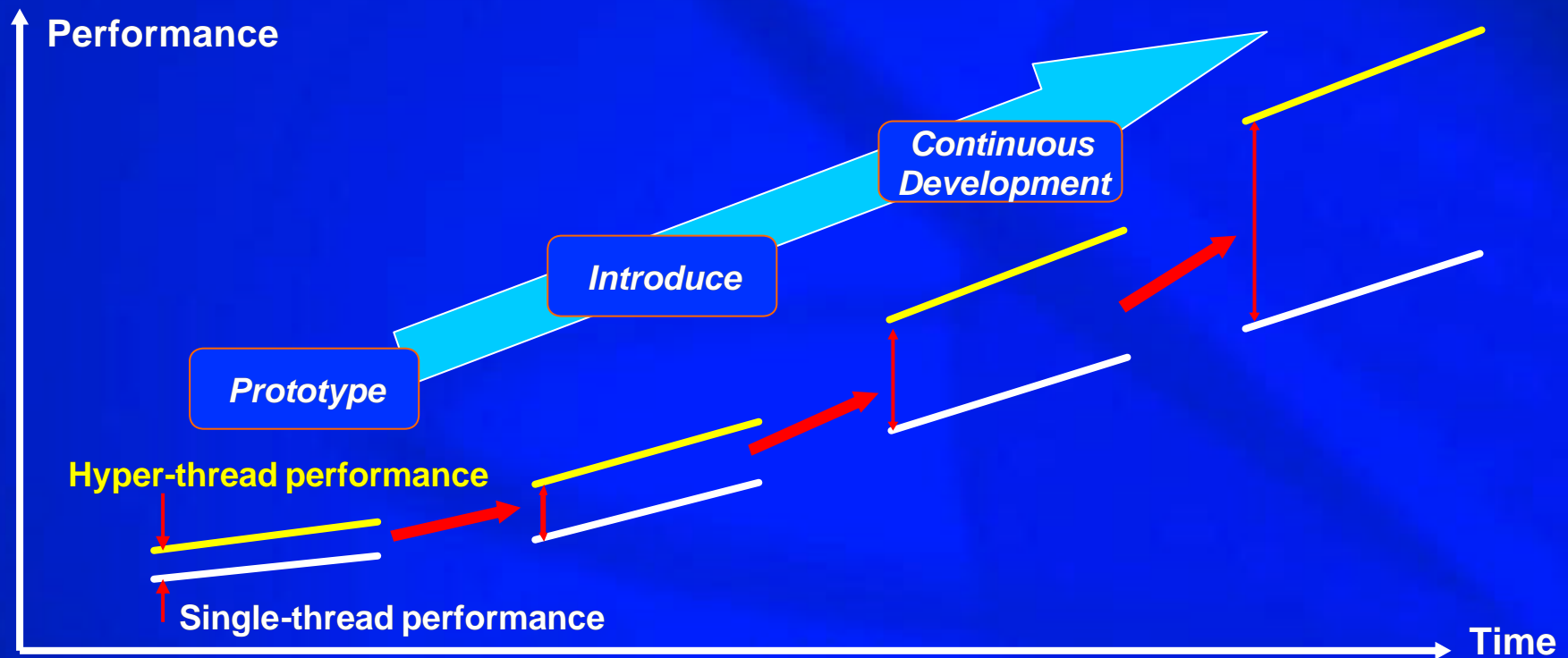


Good performance benefit from small die area investment

Performance tests and ratings are measured using specific computer systems and/or components and reflect the approximate performance of Intel products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance. Buyers should consult other sources of information to evaluate the performance of systems or components they are considering purchasing. For more information on performance tests and on the performance of Intel products, reference www.intel.com/procs/perf/limits.htm or call (U.S.) 1-800-628-8686 or 1-916-356-3104



Intel's Long-Term Hyper-Threading Strategy



Hyper-Threading Technology expected to deliver increasingly higher performance

Challenges for Multithreading

✍ Multithreaded Applications

- Multithread Programming
- Automatic Thread Partitioning

✍ Design Complexity

- Additional Validation Overhead
- Increasing Scalar Inefficiency
- SMT vs. CMP Tradeoffs ?

✍ Latency of Single-Threaded Applications

Multi-Threaded Processors

✍ Targeting:

- Throughput of Multi-tasking Workloads
- Latency of Multi-threaded Applications

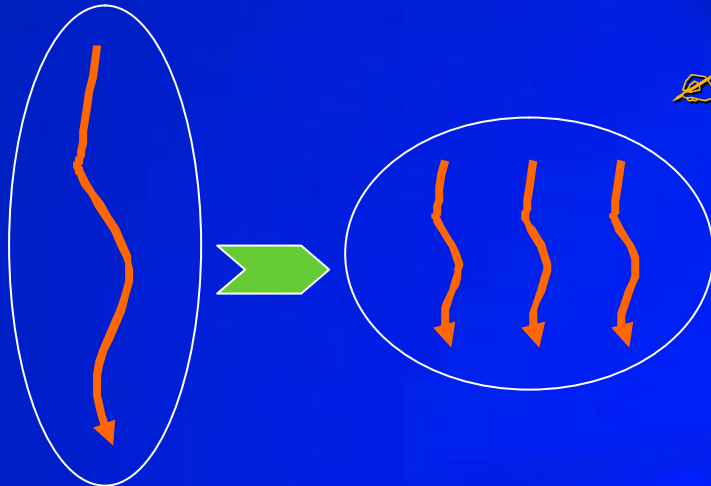
✍ Not Targeting:

- Latency of Single-threaded Applications

✍ Research Challenge:

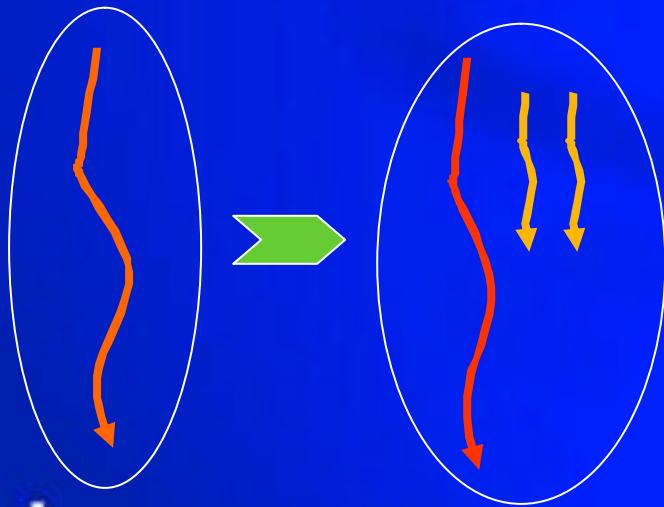
- Leverage Multi-threaded CPU to Improve Latency of Single-threaded Applications

Asymmetric Multi-Threading



Symmetric Multi-Threading

- Partition single thread into multiple threads
- Achieve performance by parallel execution of multiple threads



Asymmetric Multi-Threading

- Attach helper threads to original single-threaded code
- Achieve speedup via memory prefetching by helper threads

Asymmetric Helper Threads

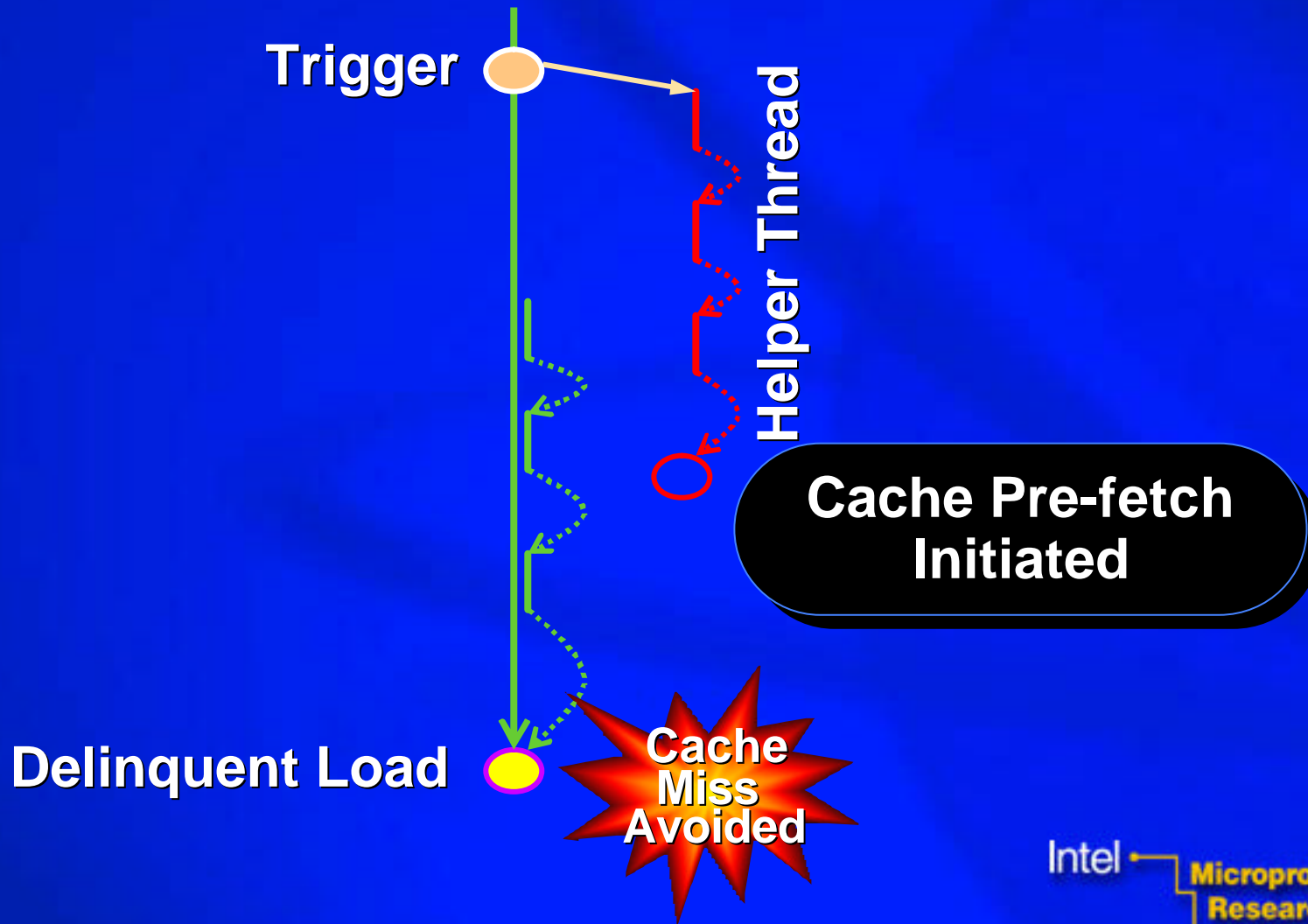
✎ Symmetric Multi-Threading Compiler

- Partition single thread into multiple threads
- Must ensure semantic correctness
- Difficult for common and legacy code

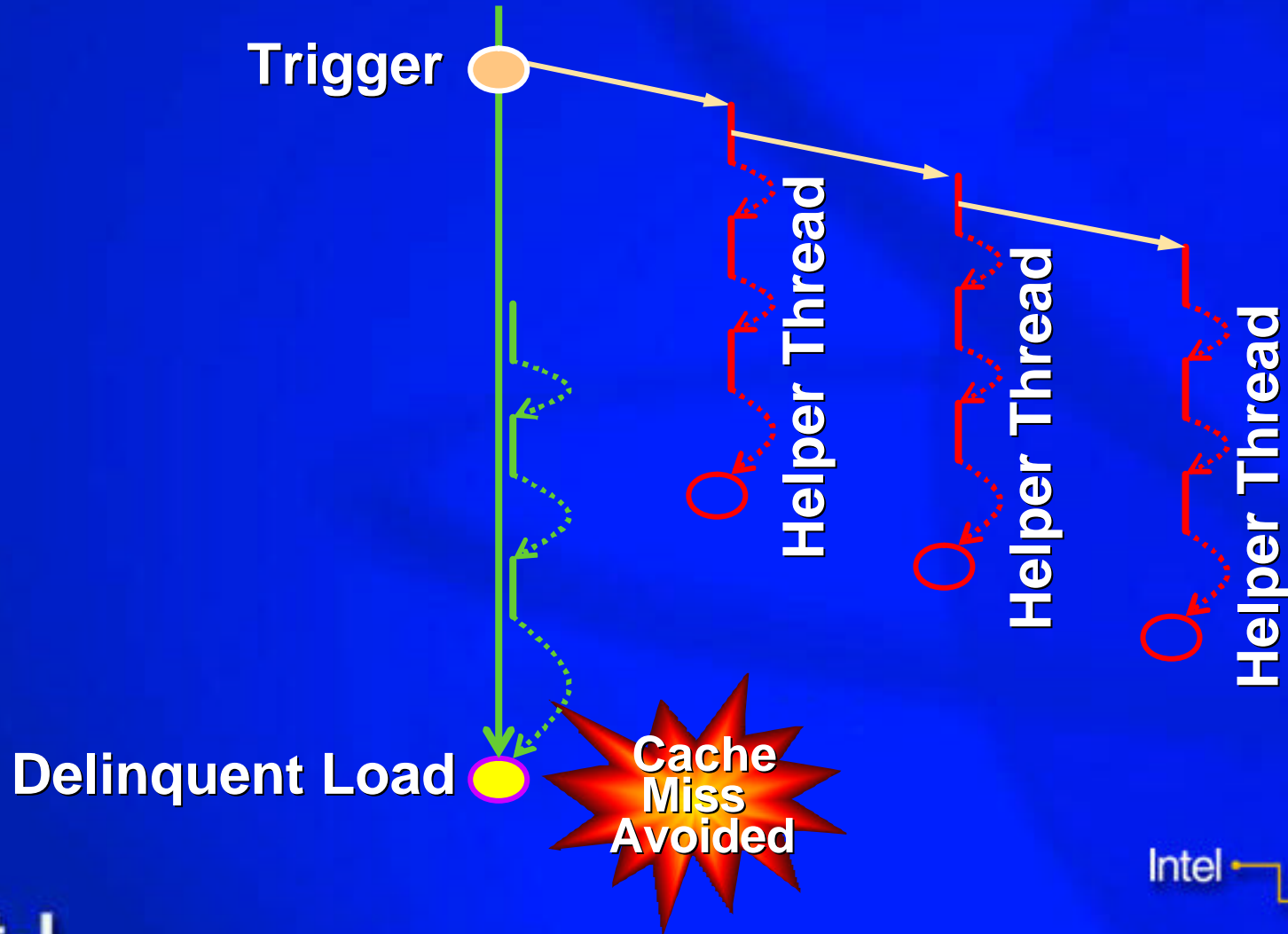
✎ Asymmetric Multi-Threading Compiler

- Attach helper threads to original code
- Leverage side effect of helper threads
- Can be dynamically invoked/controlled

Pre-fetch via Helper Threads



Chaining Triggers of Helper Threads



Chaining Trigger Advantages

✍ Low-cost Thread Spawning:

- Chaining triggers initiate helper threads without impacting main thread performance

✍ Long-range Prefetching:

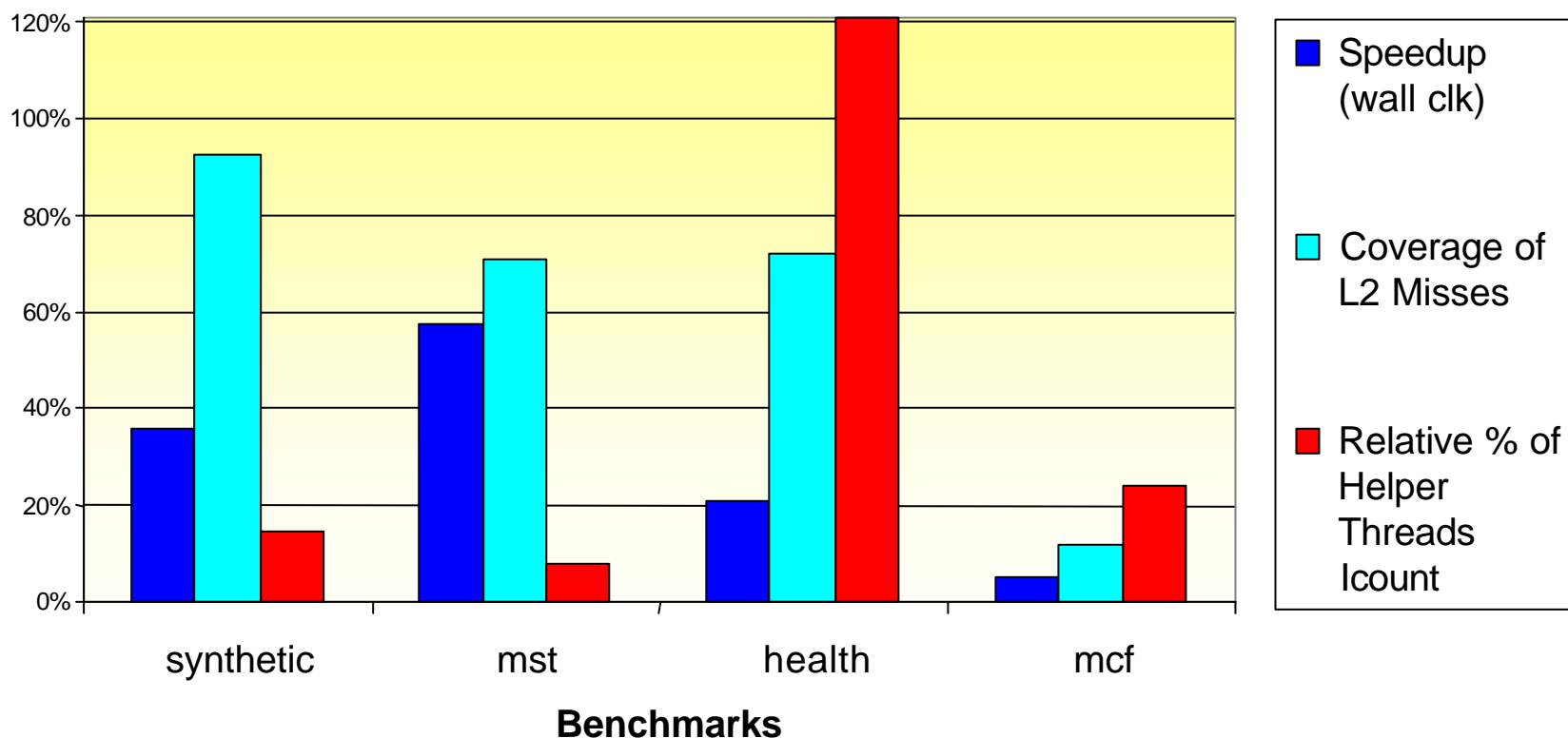
- Can target delinquent loads far ahead of the main thread
- Helper threads make progress independent of main thread's lack of progress

Prefetching Effectiveness of Helper Threads on HT Machine

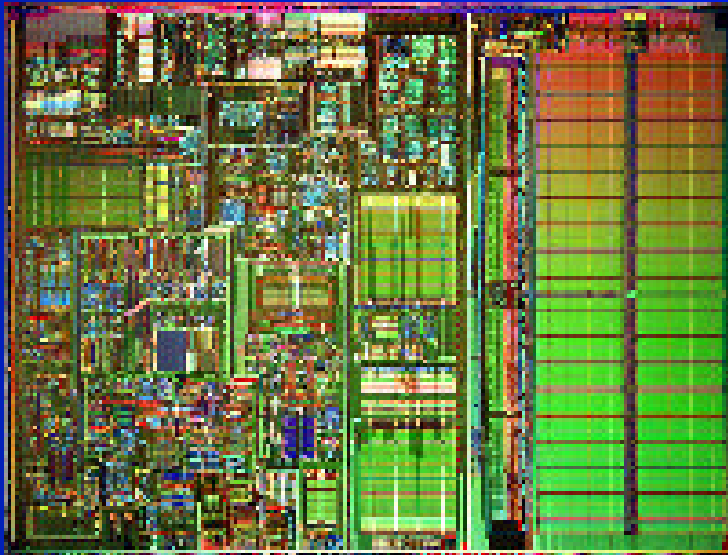
Benchmark	Description	Speedup
<i>Synthetic</i>	Graph traversal in large random graph simulating large database retrieval	22% - 45%
<i>MST</i> (Olden)	Minimal Spanning Tree algorithm used for Data Clustering	23% - 40%
<i>Health</i> (Olden)	Hierarchical database modeling health care system	11% - 24%
<i>MCF</i> (SPECint)	Integer programming algorithm used for bus scheduling	7.08%

Helper Threads on Hyper-Threading Machine

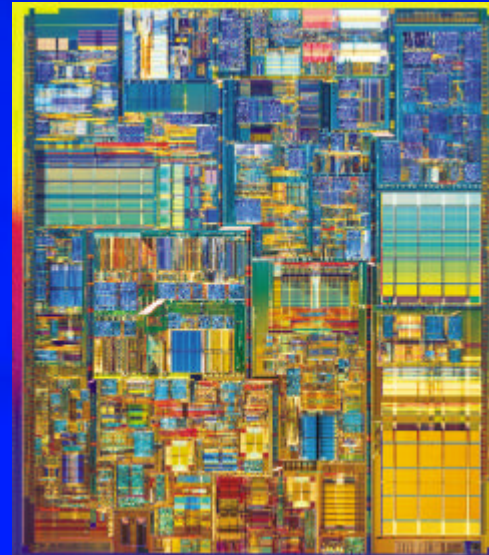
Asymmetric MT effectiveness:
speedup, L2 miss coverage; % helper thread lcount



Intel® Processors with Netburst™ Microarchitecture



Intel® Xeon™ MP Processor
256KB 2nd-Level Cache
1MB 3rd-Level Cache
.18u process



Intel Xeon Processor
256KB 2nd-Level Cache
.18u process



Intel Xeon Processor
512KB 2nd-Level Cache
13u process
Intel Microprocessor Research Forum 2002

