

VLSI, very large scale integration	ATA, Advanced Technology Attachment 高级技术附件	Flit, flow control unit
IC, integrated circuit 集成电路	SATA, Serial ATA, 串行高级技术附件	MIN, multistage interconnection network
HPC, high performance computer	SCSI, Small Computer System Interface	
IDC, Internet data center	SAS, Serial Attached SCSI	ICT, information and communication technology
UPS, uninterruptable power supply	FC, Fiber Channel, 光纤信道协议	WSC, warehouse-scale machine
PDU, power distribution units	JBOD, Just a bunch of disks	PUE, power usage effectiveness
TCO, total cost of ownership	RAID, redundant array of inexpensive drives	WUE, water usage effectiveness
CapEx, capital expenditure 资本支出	ECC, Error correcting coding, 错误校验码	CUE, carbon usage effectiveness
OpEx, operational expenditure 运营开销	DAS, Direct Access Storage, 直接存取存储器	ATS, automatic transfer switch
	NAS, Network Attached Storage 网络附属存储	STS, static transfer switch equipment
ISA, instruction set architecture	SAN, Storage Area Network, 存储区域网络	UPS, uninterruptable power supply
CISC, complex instruction set computer	NAND Flash cell	PDU, power distribution unit
RISC, reduced instruction set computer	SLC, single-level cell	PSU, power supply unit
PC, program count	MLC, multi-level cell	CRAC, computer room air conditioning
GPR, general-purpose register 通用寄存器	SSD, Solid-State Drive	COP系数, coefficient of performance
SPR, special-purpose register	WA, Write amplification	TOR, top of rack
ILP, instruction-level parallelism		EOR, end of rack
uOps, micro-ops	CPI, Clock cycle Per Instruction	MDA, main distribution area
	PCA, Principal Component Analysis 主成分分析	HAD, horizontal distribution area
RAW, read after write	NTT, normalized turnaround time	EDA, equipment distribution area
WAR, write after read	ANTT, average NTT	SLA, service level agreement
WAW, write after write		DFS,
FU, function unit	SISD, Single Instruction Single Data	DVFS, dynamic voltage and frequency scale
IF, instruction fetch	SIMD, Single Instruction Multiple Data	MDC, modular data center
ID, instruction decode	MISD, Multiple Instruction Single Data	
RO, read operands	MIMD, Multiple Instruction Multiple Data	ACPI, advanced configuration and power interface
EX, execution	UMA架构, Centralized Shared-Memory	TDP, thermal design power
WR, write result	SMP, symmetric multiprocessors, 对称多处理器	D2D, die to die
RS, reservation station	NUMA架构, DSM, Distributed Shared-Memory	C2C, core to core
CDB, common data bus	MSI, MESI	MPP, maximal power point
	TLP, thread-level parallelism	PTP, performance time product
IPC, instruction per cycle	CMP, Chip Multi-Processor, 片上多处理器	TPR, throughput power ratio
VLIW, very long instruction word	SMT, simultaneous multithreading, 超线程	ROI, return on investment
DB, dispatch buffer	LLC, Last Level Cache	
BHT, branch history tables	MIC, Intel Many Integrated Core 架构	MTTF, mean time to failure
BTB, branch target buffer	SCC, Intel Single-Chip Cloud Computer	MTBF, mean time between failure
BIA, branch instruction address		MTTR, mean time to repair
BTA, branch target address	DLP, Data-level parallelism, 数据级并行	ACE. Architecturally correct execution
OoO, out-of-order execution	MVL, maximum vector length	AVF, architectural vulnerability factor
ROB, reorder buffer	VLR, Vector Length Register	PSU,
EPIC, explicitly parallel instruction computing	GPU, Graphics processing unit	
	GPGPU, general-purpose computation on GPU GPU上的通用计算	
PTE, page table entry	SPMD, Single Program Multiple Data	
MMU, memory management unit	SIMT, Single Instruction Multiple Thread	
VPN, virtual page number	CUDA, Stands for Compute Unified Device Architecture	
PPN, physical page number	TPC, texture/processor clusters	
3C model, compulsory/capacity/conflict miss	SM, streaming multiprocessors, 流多处理器	
DIMM, dual in-line memory module	SP, streaming processor, 流处理器	
MC, memory controller	GPC, Graphics Processor Clusters, 图形处理器集群	
RAS, row access table	RF, Register File	
CAS, column access table		
SDRAM, synchronous DRAMS		
DDR, double data rate SDRAM		
IPM, instruction per miss		
MLP, memory-level parallelism		

#### Summary 1

- What is Computer Architecture
- History of IC(integrated circuit)
- Transistor basics
- Feature length
- HPC vs IDC
- Scale up/out
- Energy/power issues
- The trend of Computer Architecture research

#### Summary 2

- Architecture vs microarchitecture
- Evolution of instruction set
- CISC(IA32) vs RISC(MIPS)
- Machine interfaces: ISA
- User/System ISA
- MIPS instruction field
- Single-cycle MIPS
- Ideal pipeline
- Stage quantization
- Pipeline slot

#### Summary 3

- Pipeline stall and bubble
- Dependency and hazards
- RAW, WAR, WAW
- Forwarding and pipeline interlock
- Functional units
- Dynamic scheduling and OoO
- Scoreboarding
- Tomasulo's Algorithm

#### Summary 4

- Superscalar Pipeline
  - Limitations of scalar processor
  - Basic feature of superscalar pipeline
  - Multi-Issue Processor
  - Dispatch Buffer and Completion Buffer
  - Classification of ILP Machines
  - Rationale of branch prediction
  - 2-bit prediction
- Speculation
  - Precise exception: in-order completion
  - Reorder buffer (ROB)
  - Tomasulo's Algorithm with ROB
- VLIW and EPIC
  - CISC vs RISC vs VLIM
  - Loop unrolling
  - The concept of EPIC

#### Summary 5

- Memory hierarchy, uncore and off-chip
- Cache line, block, address
- Virtual memory, page table, PTE, TLB
- Locality principle, Inclusive and exclusive relationship
- Miss caching, victim caching, prefetching
- Cache write policies, write buffer/cache
- rank, bank, array, channel, MC, parallelism in DRAM
- 1T1C DRAM cell, data access, DRAM refresh
- DRAM access cost
- synchronous/asynchronous design
- Memory design challenges, memory wall
- MLP

#### Summary 6

- Disk concept; platter/track/sector
- Design good drive Interfaces
- Parallel/Serial ATA; Parallel/Serial SCSI
- RAID Organization
- DAS, NAS, SAN
- Flash memory cell, SLC/MLC
- SSD advantages, hybrid storage

#### Summary 7

- Amdahl's Law
- Calculating CPI (公式)
- Analyzing memory access time (公式)
- Little's Law
- Estimating server power
- Trace-/Execution- driven simulation
- Simulation acceleration
- Concepts of workload characterization
- Multi-programmed workload

#### Summary 8

- SIMD, MIMD, TLP
- Multiprocessors, UMA and NUMA
- Definition of cache coherency
- Cache coherency and memory consistency
- Basic facts of the snooping protocol
- A simple write-through invalidation protocol
- 3-state MSI protocol

#### Summary 9

- Thread, Multithreading, SMT
- CMP and multicore
- Benefits of multicore
- Multicore system architecture
- Heterogeneous multicore system
- Heterogeneous-ISA CMP
- Multicore and manycore
- Design challenge

#### Summary 10

- Throughput computing and data-level parallelism
- Vector processor and vector instruction
- VMIPS, vector registers, DAXPY, execution latency
- SIMD lanes, chaining, vector length register
- GPGPU, SIMT, CUDA programming model
- TPC, SM, SP, warp and warp scheduling
- Branch divergence
- GPU register file

#### Summary 11

- basic concepts, link/channel/buffer
- switch degree, average distance
- non-blocking network, direct/indirect network
- network performance, latency estimation
- network switch and switch strategy
- bus and crossbar
- array ring mesh torus tree butterfly hypercube

#### Summary 12

- what is a data center
- major metrics of data center design
- data center infrastructure
- the long tail concept
- data center capacity utilization
- types of power provisioning
- modular data center and cooling

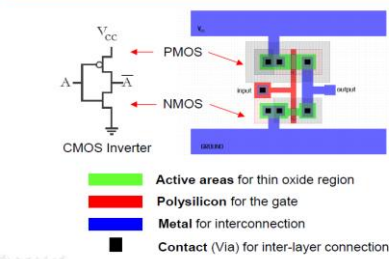
#### Summary 13

- G-States, S-States, C-States, P-States
- TDP, Turbo Boost
- Power management can be challenging

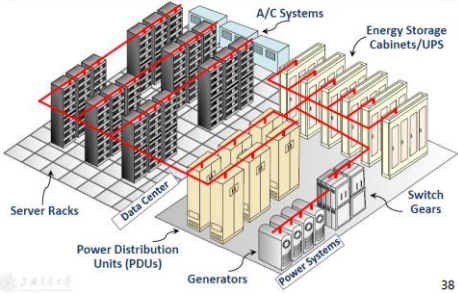
#### Summary 14

- Faults, error, and failure
- MTTF, MTBF, MTTR
- Availability, reliability
- ACE, AVF
- Redundancy

From Circuit to Layout

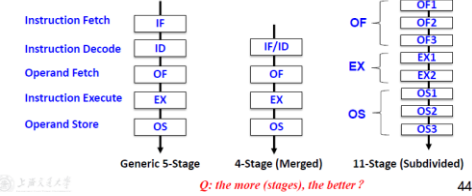


A Typical Data Center



Stage Quantization

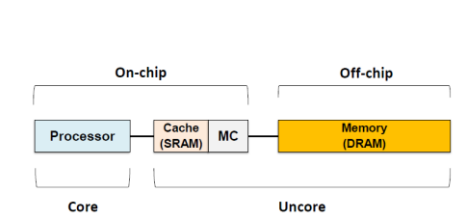
- Merge multiple sub-computations into one
  - Combining sub-computations with short latencies
- Subdivide a sub-computation into multiple
  - Fine-grained partition of sub-computations



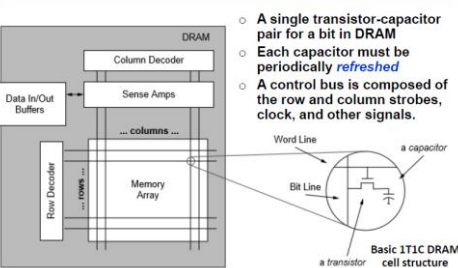
Architecture Comparison

Category	CISC	RISC	VLIW
Inst Size	Varies	Fixed (typically 32bits)	Fixed
Inst Format	Field placement varies	Regular, consistent placement of fields	
Inst Semantics	Complex; possibly many dependent ops per instr	Almost always one simple operation	Multiple, independent simple operations
Registers	Few, sometimes special	Many, general purpose	
Memory	Reg-Mem architecture	Load/Store architecture	
Hardware	Exploit microcode	Not microcoded implementation	
Example:			

Basic Concepts



1T1C DRAM Cell

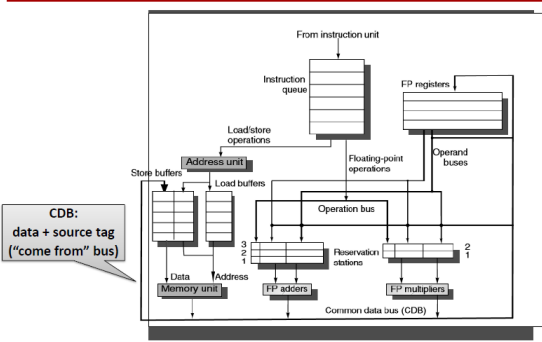


Simulation Approach Comparison

	Function Simulation	Trace-Driven Simulation	Execution-Driven Simulation
Development Time	Excellent	Poor	Very Poor
Evaluation Time	Good	Poor	Very Poor
Accuracy	Excellent	Very good	Excellent
Coverage	Poor	Excellent	Excellent

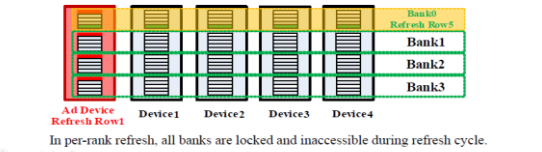
- Full system simulation
  - Trace-driven simulation and execution-driven simulation

Basic Structure Implementing Tomasulo's Alg.

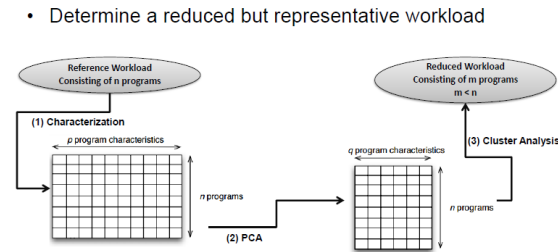


Refresh Mechanism

- A Row is the smallest refresh unit in bank.
- Typically, the retention time of data in DRAM cell is 64ms if the ambient temperature is less than 85 degree Celsius.
- Refresh operation can implement at rank level (per-rank refresh) or at bank level (per-bank refresh)

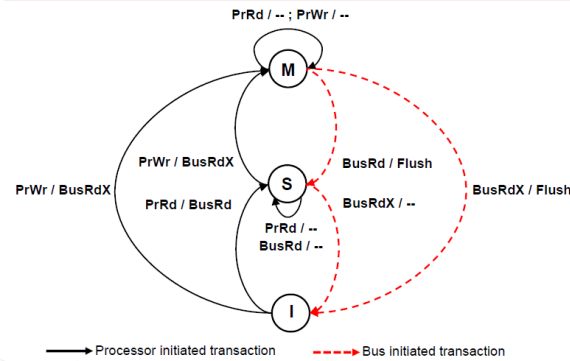


Workload Design (Cont'd)

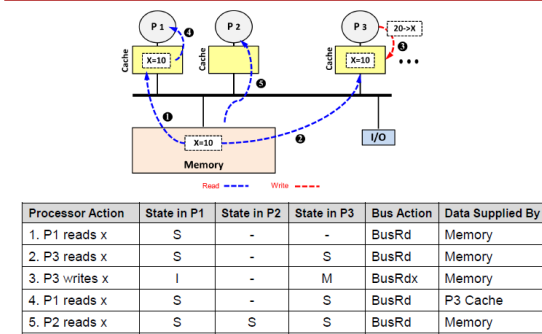


Schematic overview of the PCA-based workload reduction method

A 3-state (MSI) Write-Back Invalidation Protocol



A 3-state (MSI) Write-Back Invalidation Protocol



Processor Action	State in P1	State in P2	State in P3	Bus Action	Data Supplied By
1. P1 reads x	S	-	-	BusRd	Memory
2. P3 reads x	S	-	S	BusRd	Memory
3. P3 writes x	I	-	M	BusRdX	Memory
4. P1 reads x	S	-	S	BusRd	P3 Cache
5. P2 reads x	S	S	S	BusRd	Memory

Calculating CPI

$$CPI = \sum \frac{IC_i \times CPI_i}{\text{Instruction count}}$$

$$CPU \text{ time} = (\sum IC_i \times CPI_i) \times \text{Clock cycle time}$$

Example:

Suppose we have made the following measures:

- Frequency of FP operations = 25%
- Average CPI of FP operations = 4.0
- Average CPI of other instructions = 1.33

Then:

$$CPI \text{ original} = (4 \times 25\%) + (1.33 \times 75\%) = 2.0$$

Memory Performance Analysis

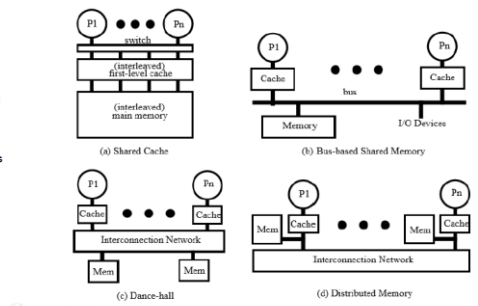
$$\frac{\text{Misses}}{\text{Instruction}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$

$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss Penalty}$$

$$\text{Average memory access time} = \text{Hit time}_{L1} + \text{Miss rate}_{L1} \times \text{Miss Penalty}_{L1}$$

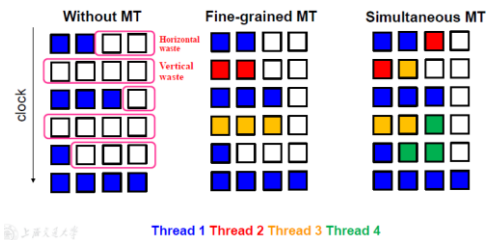
$$\text{Miss Penalty}_{L1} = \text{Hit time}_{L2} + \text{Miss rate}_{L2} \times \text{Miss Penalty}_{L2}$$

Common Memory Hierarchies in Multiprocessors



Impacts of SMT on Utilization

- Multithreaded processor improves hardware utilization in different dimensions

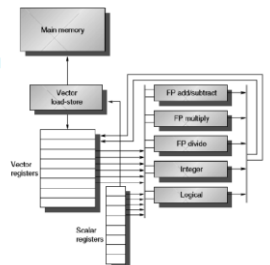


Classification of Heterogeneous Multicore

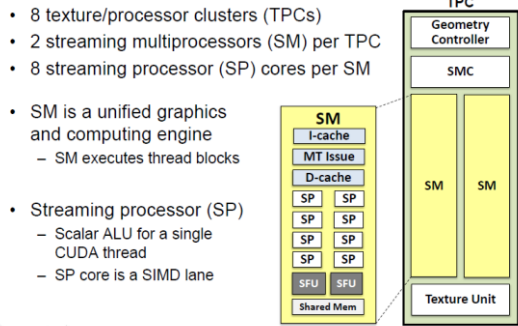
Different ISA	??	Heterogeneous-ISA Chip Multiprocessors
	Same ISA	Homogeneous Multi-/Many- Core
Same Cores	Same ISA	Cores of Different Capabilities
	Different Cores	

Case Study: VMIPS

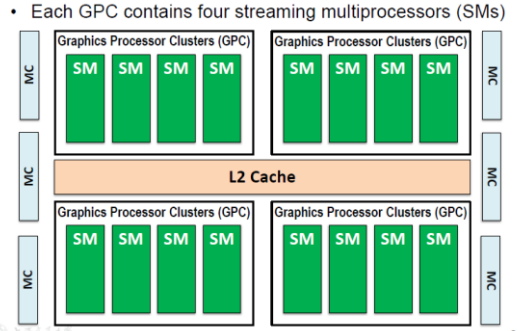
- Vector registers
  - 64-element register
  - Register file has 16 read ports and 8 write ports
- Vector functional units
  - 5 fully pipelined FUs
  - Hazards detection
- Vector load-store unit
  - Loads/stores a vector
  - 1 word per clock cycle



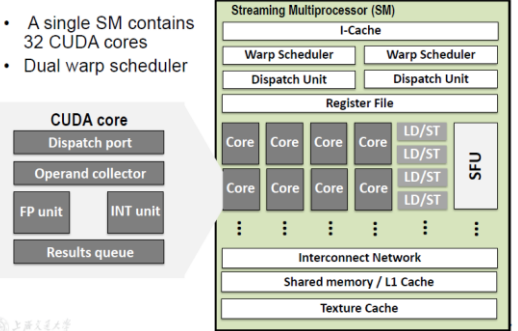
Case Study: Tesla GPU Architecture



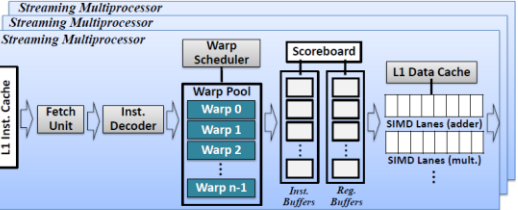
Case Study: Fermi GF100 Architectural Overview



Case Study: Fermi GF100 Architectural Overview



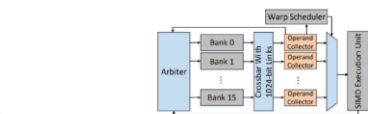
Warp Scheduling (Cont'd)



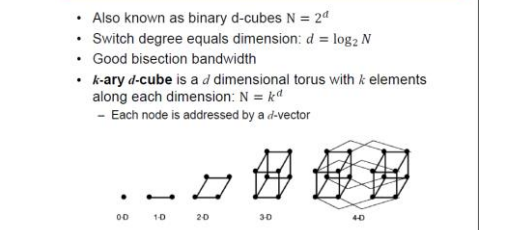
- Potential factors that can delay a warp's execution
  - Scheduling policies
  - Instruction/Data cache miss
  - Structural/Control/Data hazard
  - Synchronization primitives

Resource Limits

- The maximum parallelism in GPUs is often limited by the register file capacity
  - Applications with high TLP triggers more active warps
- Register file is a large SRAM structure
  - Fastest memory block available to the processor
  - Store intermediate results from units such as ALU
  - Power hungry structure



Hypercube (Optional)



Performance Evaluation

- Link width:  $w$
- Unit interval:  $\tau$
- Signaling rate:  $f = 1/\tau$
- Channel bandwidth:  $b = w \cdot f$
- Total bandwidth of all the channels (or links)
  - the number of channels times the bandwidth per channel

Latency (Lower Bound)

- Sending overhead:  $Overhead_s$
  - Receiving overhead:  $Overhead_r$
  - Total routing time:  $T_R$
  - Arbitration time:  $T_A$
  - Switching time:  $T_S$
  - Total time of flight of the packet  $T_{TotalProp}$
- Latency =  
 $Overhead_s + (T_{TotalProp} + T_R + T_A + T_S) + \frac{Packet\ size}{Bandwidth} + Overhead_r$

Crossbar

- Crossbar switch
  - A type of fully-connected network
  - Every node connected to all others
  - $O(N)$  bandwidth
  - Cost of interconnect:  $O(N^2)$
  - Good for small number of nodes

Multidimensional Topology: 2D Torus

- Reduces the diameter of a mesh network
- Torus: adding end-round direct connections
  - An extension of the 2D mesh
  - A regular torus has long wrap-around links
  - Slightly lower latency while higher cost

Trees (Optional)

- Trees features planar, hierarchical topology
- Employed as indirect networks with hosts as leaves
- Routing distance grows only logarithmically

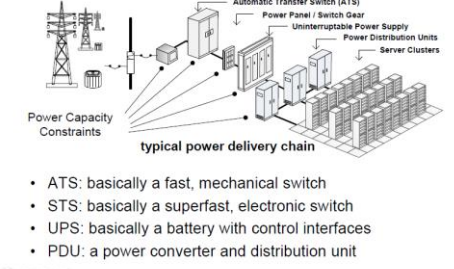
Multistage Interconnection Network (Optional)

- Indirect networks with multiple layers of switches
- Omega network:
  - $\log(N)$  number of stages and  $N/2$  switching units

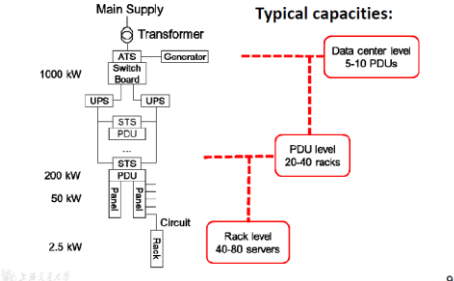
Butterfly Topology (Optional)

- Butterfly is an important logarithmic network
  - Can be viewed as a tree with multiple roots
- A d-dimensional indirect butterfly:
  - Connects  $N = 2^d$  nodes ( $d \geq 2$ )
  - $d = \log_2 N$  levels of switches

Datacenter Infrastructure - Power System



Datacenter Infrastructure - Power System



Global System States (G-states)

G-states are high-level description of the platform states

- G0 (working)
  - The working system state
- G1 (sleeping)
  - No computational task is performed
- G2 (soft off)
  - Powered down but can be restarted by interrupts
- G3 (hard off)
  - Mechanical off

Sleep States (S-states)

S-states are set in the BIOS and configured by the system

- G0-S0: normal operation
- G1
  - S1: processor clock is off
  - S2: processor is off
  - S3: suspend to RAM
  - S4: suspend to disk
- G2-S5: halt state

Processor Power States (C-states)

G0 and S0 together define a working platform state, at which a range of C-states are defined to save power

- C0 State (normal operating state)
  - code is being executed
- C1 State (auto halt):
  - The clock is gated, i.e., prevented from reaching the core
- C3 State (sleep):
  - Maintains architectural state but flushes data to shared LLC
  - Shut down the clock generators
- C6 State:
  - Architectural states are saved to a dedicated SRAM
  - Core voltage reduced to zero volts

Processor Performance State (P-states)

P-States talk about different operational modes (freq.)

- Multiple levels of clock frequency
  - From P0 (the highest performance) to Pn (the lowest performance)
- Sub states of C0
  - Defines dynamic voltage and frequency scaling (DVFS) steps
- Switching latencies are negligible for most purposes

Frequency	Voltage	P-State
1.6 GHz	1.484 V	P0
1.4 GHz	1.420 V	P1
1.2 GHz	1.276 V	P2
1.0 GHz	1.184 V	P3
800 MHz	1.036 V	P4
600 MHz	0.956 V	P5

Intel Pentium M at 1.6GHz

Tracking Coordination

