

CS252
Graduate Computer Architecture
Lecture 2

Review of Cost, Integrated Circuits, Benchmarks,
Moore's Law, & Prerequisite Quiz

January 19, 2001
 Prof. David A. Patterson
 Computer Science 252
 Spring 2001

Review #1/3:
Pipelining & Performance

- Just overlap tasks; easy if tasks are independent
- Speed Up \leq Pipeline Depth; if ideal CPI is 1, then:

$$\text{Speedup} = \frac{\text{Pipeline depth}}{1 + \text{Pipeline stall CPI}} \cdot \frac{\text{Cycle Time}_{\text{unpipelined}}}{\text{Cycle Time}_{\text{pipelined}}}$$
- Hazards limit performance on computers:
 - Structural: need more HW resources
 - Data (RAW, WAR, WAW): need forwarding, compiler scheduling
 - Control: delayed branch, prediction
- Time is measure of performance: latency or throughput
- CPI Law:

$\text{CPU time} = \frac{\text{Seconds}}{\text{Program}} = \frac{\text{Instructions} \times \text{Cycles} \times \text{Seconds}}{\text{Program} \quad \text{Instruction} \quad \text{Cycle}}$

Review #2/3: Caches

- The Principle of Locality:
 - Program access a relatively small portion of the address space at any instant of time.
 - Temporal Locality: Locality in Time
 - Spatial Locality: Locality in Space
- Three Major Categories of Cache Misses:
 - Compulsory Misses: sad facts of life. Example: cold start misses.
 - Capacity Misses: increase cache size
 - Conflict Misses: increase cache size and/or associativity.
- Write Policy:
 - Write Through: needs a write buffer.
 - Write Back: control can be complex
- Today CPU time is a function of (ops, cache misses) vs. just f(ops): What does this mean to Compilers, Data structures, Algorithms?

Now, Review of Virtual Memory

Basic Issues in VM System Design

- size of information blocks that are transferred from secondary to main storage (M)
- block of information brought into M, and M is full, then some region of M must be released to make room for the new block --> replacement policy
- which region of M is to hold the new block --> placement policy
- missing item fetched from secondary memory only on the occurrence of a fault --> demand-based policy

Paging Organization
 virtual and physical address space partitioned into blocks of equal size
 page frames
 pages

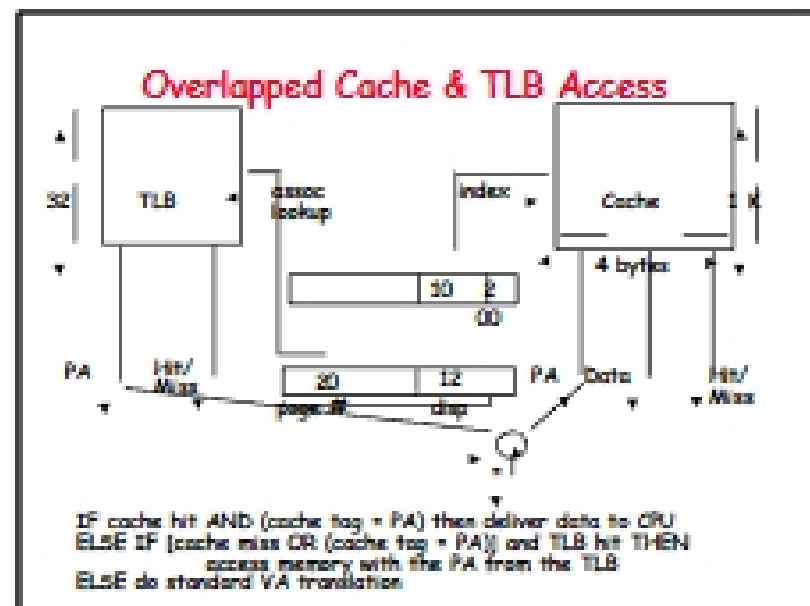
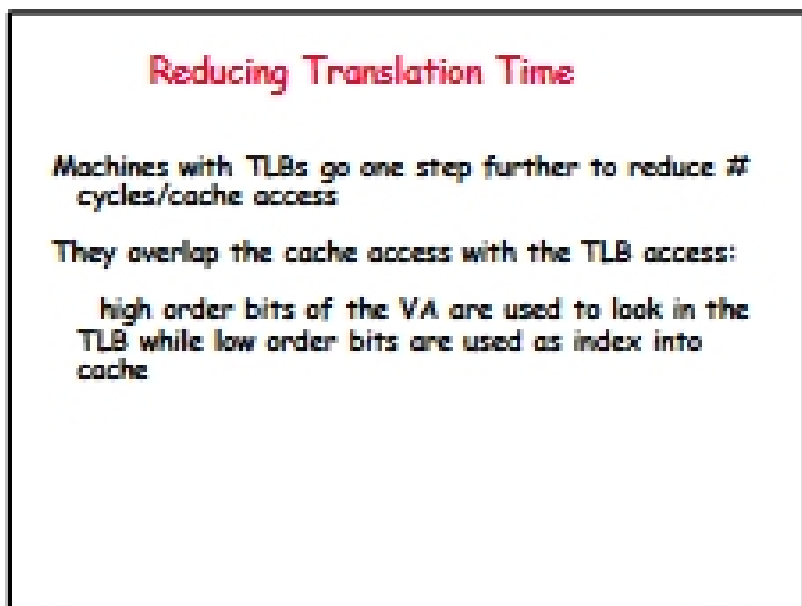
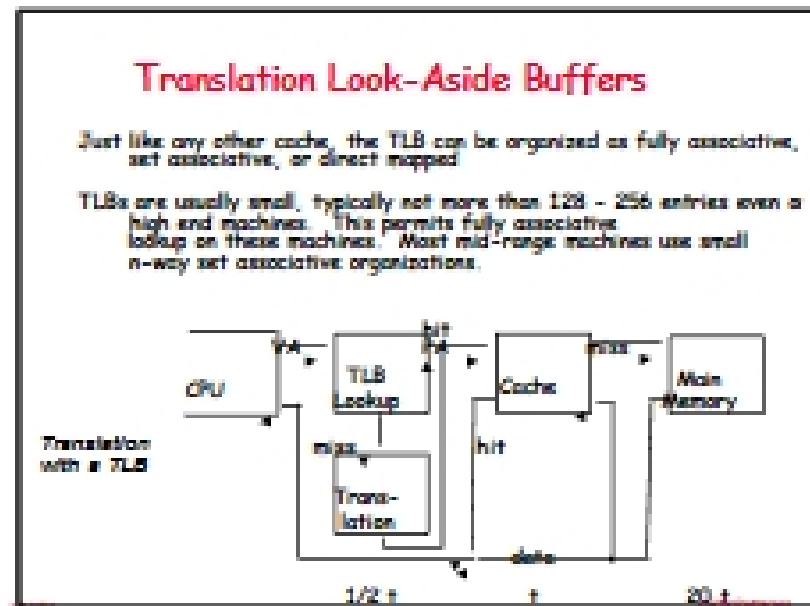
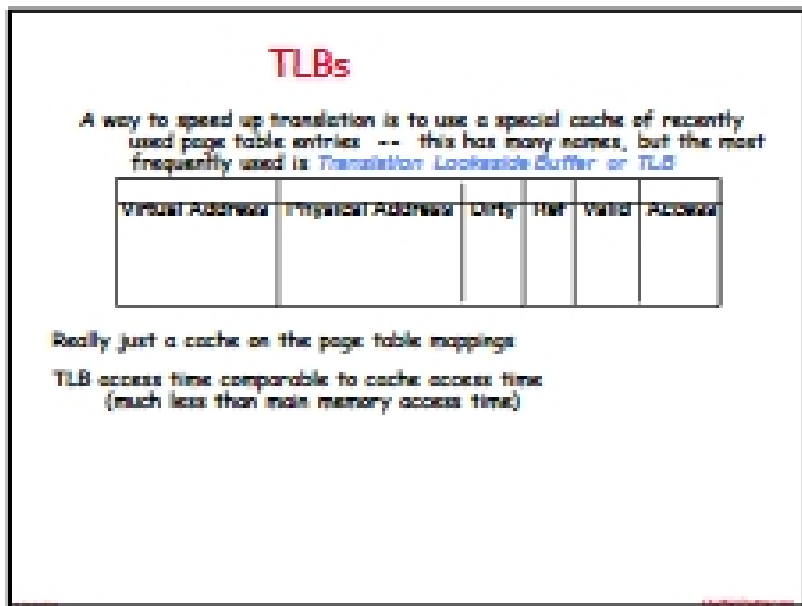
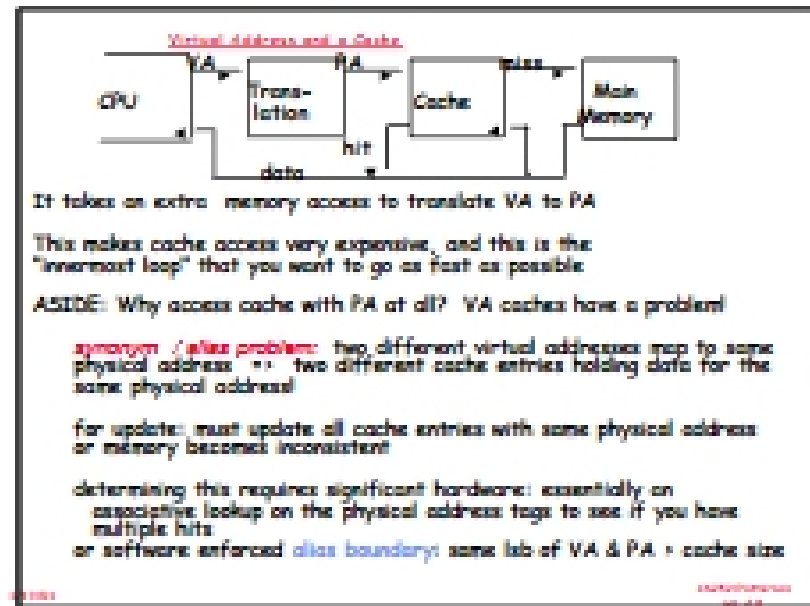
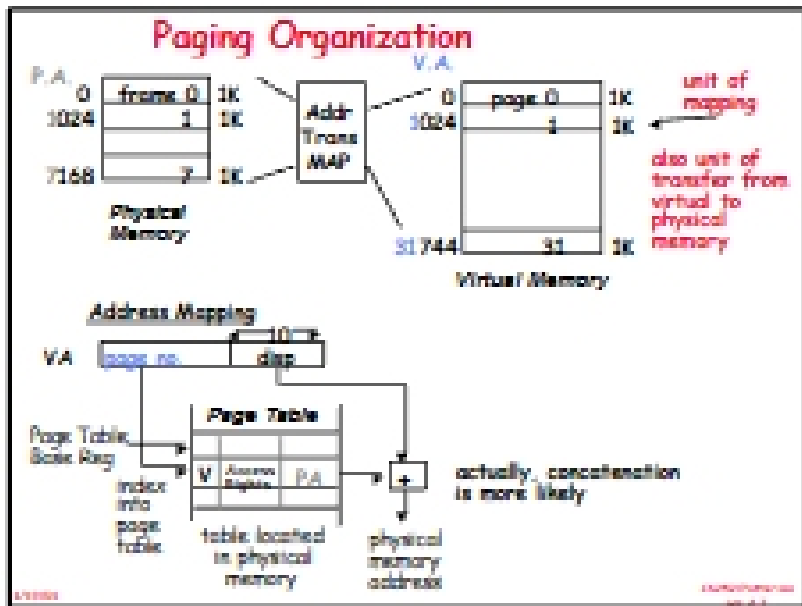
Address Map

$V = \{0, 1, \dots, n - 1\}$ virtual address space $n = m$
 $M = \{0, 1, \dots, m - 1\}$ physical address space

MAP: $V \rightarrow M \cup \{0\}$ address mapping function

$\text{MAP}(a) = a'$ if data at virtual address a is present in physical address a' and a' in M
 $= 0$ if data at virtual address a is not present in M

physical address
 OS performs this transfer

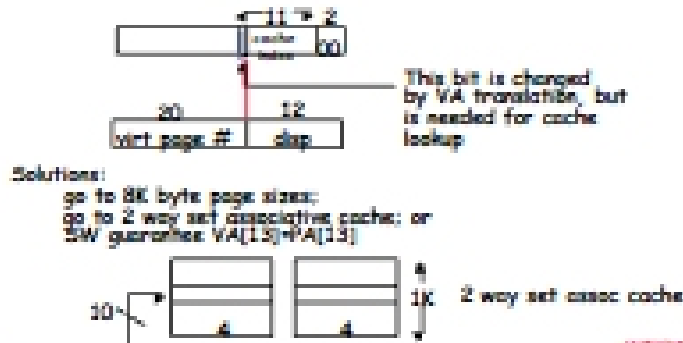


Problems With Overlapped TLB Access

Overlapped access only works as long as the address bits used to index into the cache **do not change** as the result of VA translation

This usually limits things to small caches, large page sizes, or high n-way set associative caches if you want a large cache

Example: suppose everything the same except that the cache is increased to 8 K bytes instead of 4 K:



SPEC: System Performance Evaluation Cooperative

- First Round 1989
 - 10 programs yielding a single number (SPECmark)
- Second Round 1992
 - SPECint92 (6 integer programs) and SPECfp92 (14 floating point programs)
 - Compiler Flags unlimited, March 93 of DEC 4000 Model 610:


```
apico: unix .cc/daf-(agnt,base_bcopy, "bcopy(a,b,c)-
                                memcpy(b,a,c) "
```
 - was5: /all-(all,docs-dac) /ag-a/cr-4/cr-200
 - was7: /nopec/ag-a/cr-4/cr-200/1c-bias
- Third Round 1995
 - new set of programs: SPECint95 (8 integer programs) and SPECfp95 (10 floating point)
 - "benchmarks useful for 3 years"
 - Single flag setting for all programs: SPECint_base95, SPECfp_base95

SPEC: System Performance Evaluation Cooperative

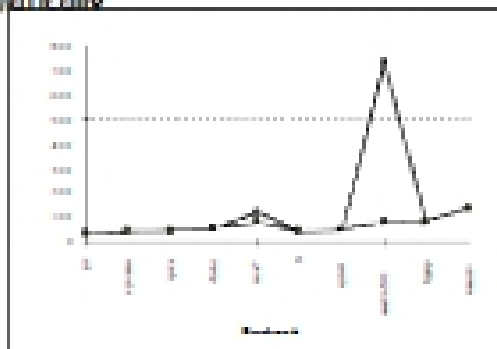
- Fourth Round 2000: SPEC CPU2000
 - 12 Integer
 - 14 Floating Point
 - 2 choices on compilation: "aggressive" (SPECint2000, SPECfp2000), "conservative" (SPECint_base2000, SPECfp_base); flags same for all programs, no more than 4 flags, some compiler for conservative, can change for aggressive
 - multiple data sets so that can train compiler if trying to collect data for input to compiler to improve optimization

How to Summarize Performance

- Arithmetic mean (weighted arithmetic mean) tracks execution time: $\sum(T_i)/n$ or $\sum(W_i * T_i)$
- Harmonic mean (weighted harmonic mean) of rates (e.g., MFLOPS) tracks execution time: $n/\sum(1/R_i)$ or $n/\sum(W_i/R_i)$
- Normalized execution time is handy for scaling performance (e.g., X times faster than SPARCstation 10)
- But do not take the arithmetic mean of normalized execution time, use the geometric mean: $(\prod T_i / N)^{1/n}$

SPEC First Round

- One program: 99% of time in single line of code
- New front-end compiler could improve dramatically



Impact of Means on SPECmark89 for IBM 550

Program	Ratio to VAX:		Time:		Weighted Time:	
	Before	After	Before	After	Before	After
gcc	30	23	49	51	6.91	9.22
expresso	35	34	65	67	7.64	7.88
apico	47	47	510	510	5.89	5.89
dotuc	46	49	41	38	5.81	5.45
was7	76	144	258	140	3.43	1.88
li	34	34	183	183	7.88	7.88
sqriott	40	40	28	28	6.88	6.88
matrix300	78	730	58	8	3.43	0.37
fpccc	50	87	34	33	2.97	3.07
tomcatv	35	138	20	13	2.81	1.94
Mean	54	72	124	108	54.42	40.59
	Geometric		Arithmetic		Weighted Arith	
	Ratio		Ratio		Ratio	
	1.33		1.18		1.09	