# Techniques for Mitigating Memory Latency Effects in the PA-8500 Processor
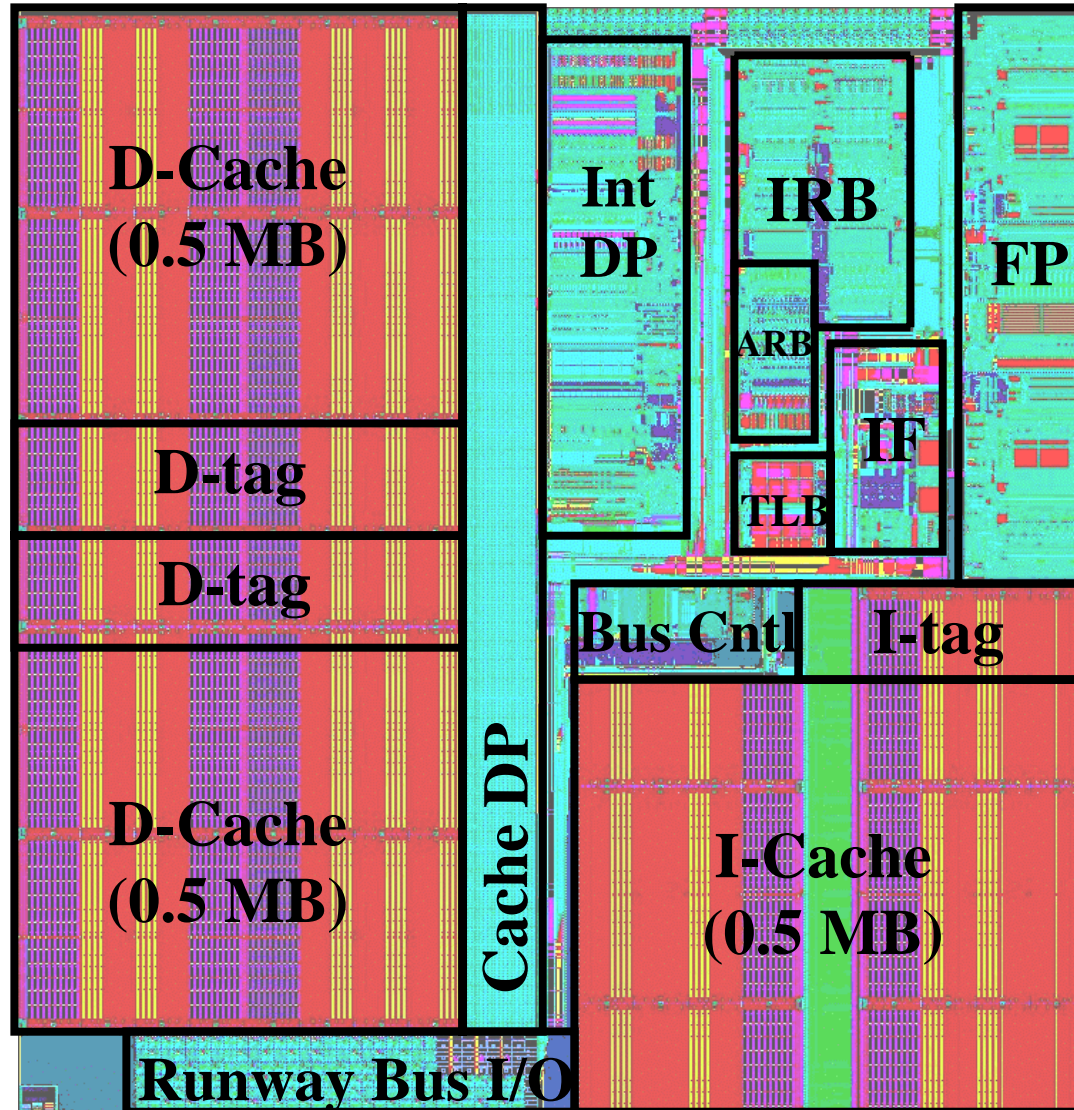
David Johnson
Systems Technology Division
Hewlett-Packard Company
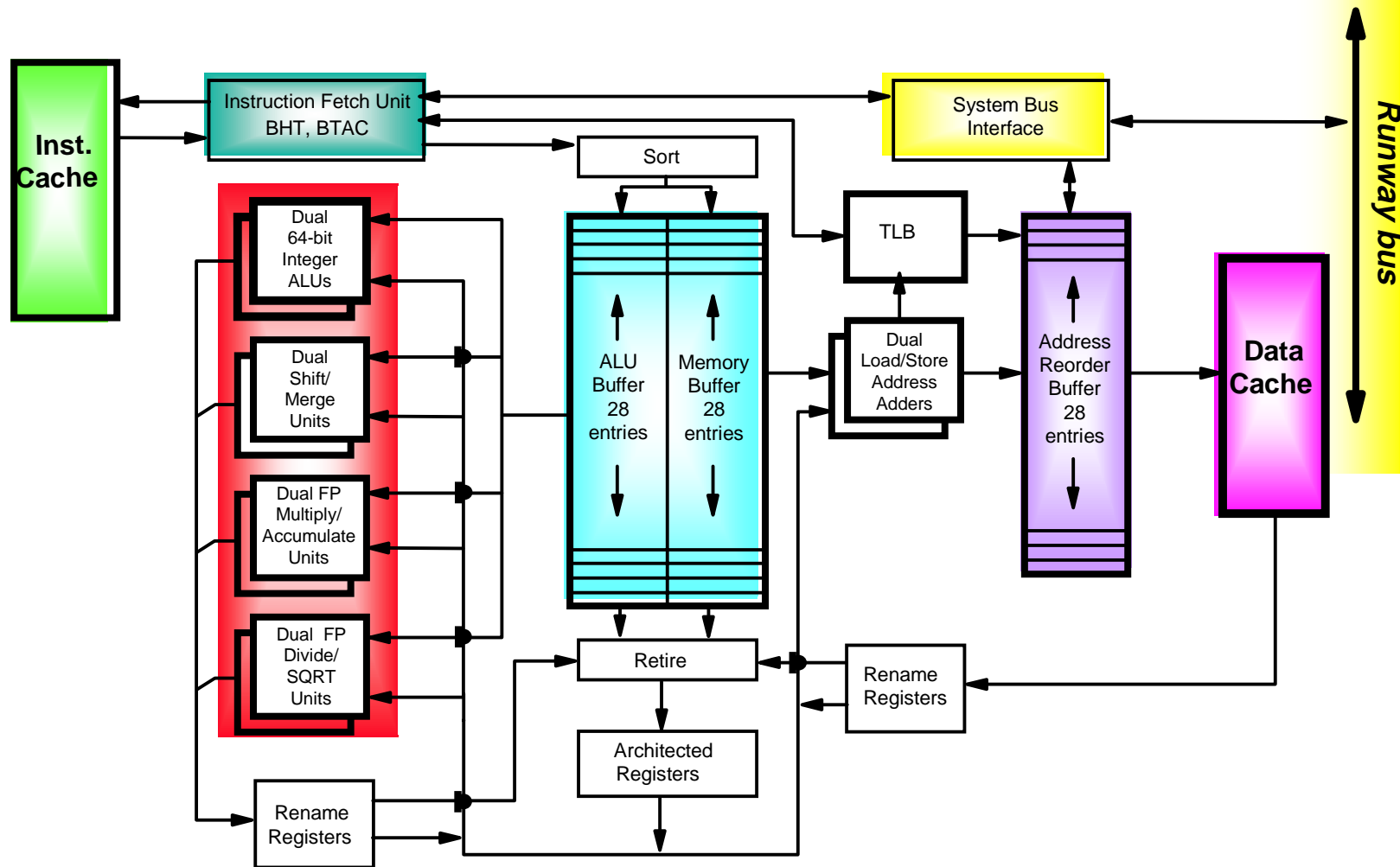
HEWLETT PACKARD

# Presentation Overview

- PA-8500 Overview

- Instruction Fetch Capabilities

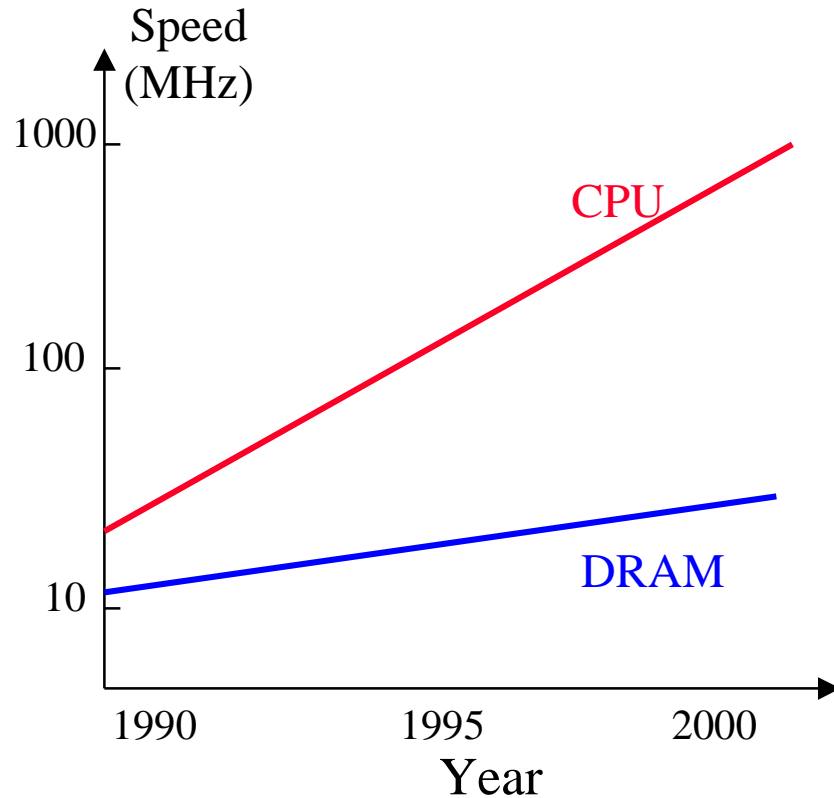- Reorder Buffers ("The Queue")

- Data Cache

- System Bus

HEWLETT PACKARD

# PA-8500



D-Cache (0.5 MB)

D-tag

D-tag

D-Cache (0.5 MB)

Cache DP

Int DP

IRB

ARB

IF

TLB

FP

Bus Cntl

I-tag

I-Cache (0.5 MB)

Runway Bus I/O

HEWLETT PACKARD

# PA-8500 Processor Core

# Memory Latency

# Instruction Fetch Features

- Instruction Cache
  - 0.5 MB on-chip cache
  - 4-way set associative
  - Pipelined 2-cycle access
  - Provides 4 instructions per cycle to CPU core
  - Supports 32-byte and 64-byte line sizes
- Instruction Prefetching
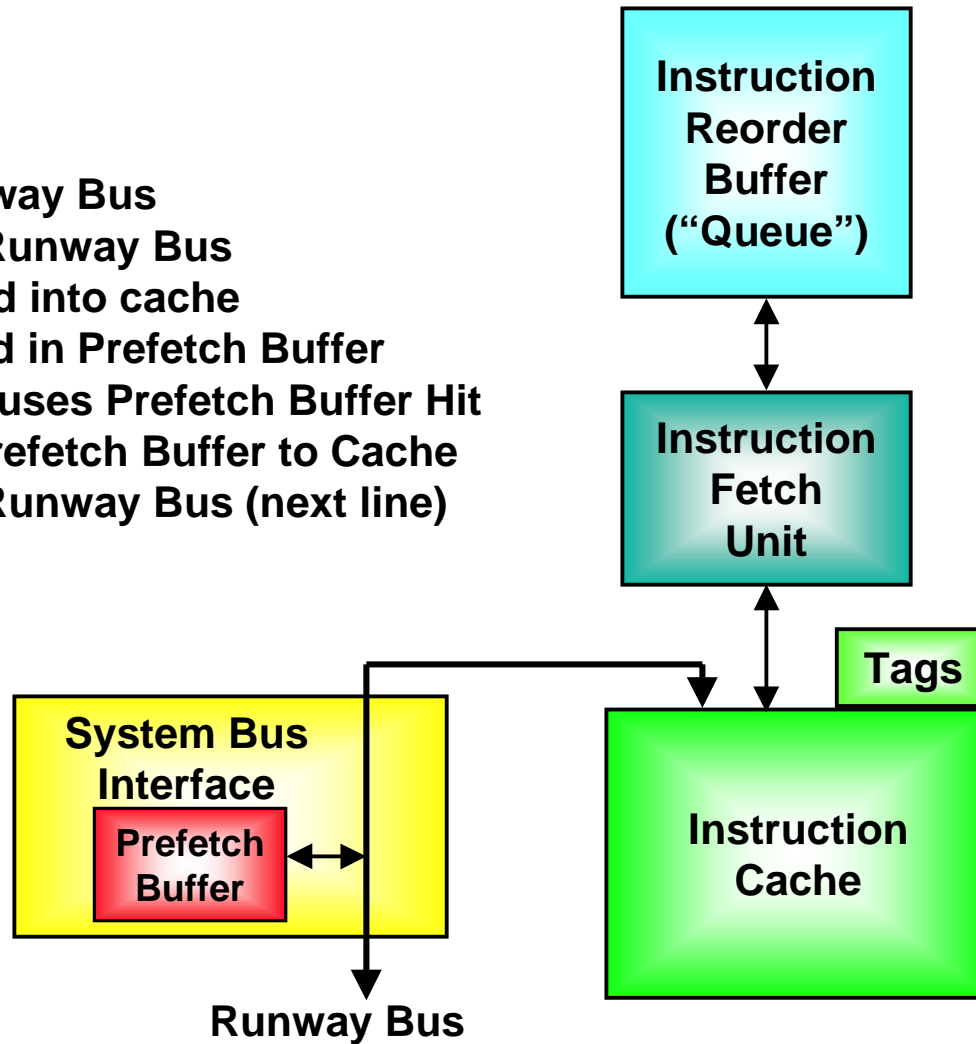
# PA-8500 I-Cache Composition

*4 Instructions per cycle
to Queue
from a 0.5 MB cache*

**Instruction Reorder Buffer ("Queue")**

**I-Fetch**

**mux**

**=**

**TAGS**

**I-Cache RAM**

**I-Cache RAM**

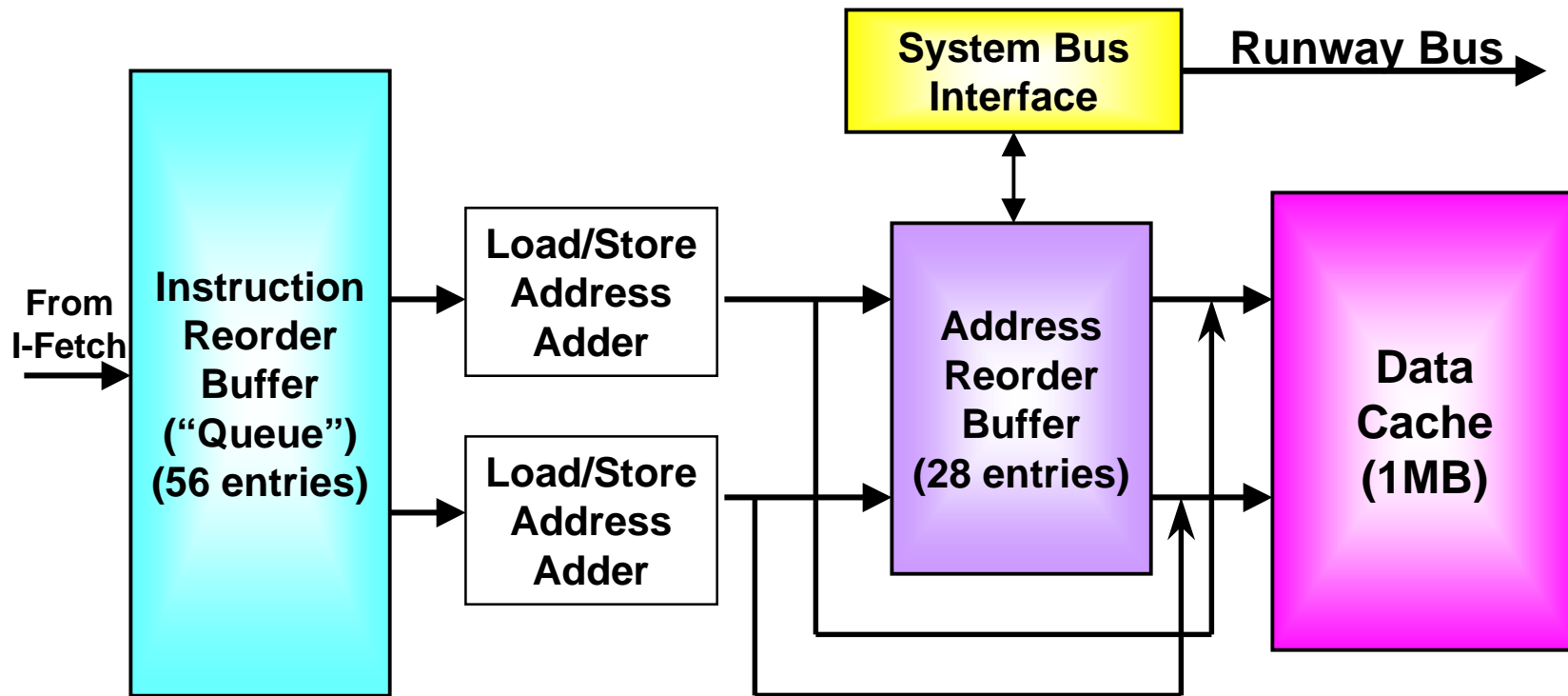**I-Cache RAM**

**I-Cache RAM**

Address

# PA-8500 Instruction Prefetching

1. I-Miss from cache
2. I-Miss issued to Runway Bus
3. I-Prefetch issued to Runway Bus
4. I-Miss Return inserted into cache
5. I-Prefetch Return held in Prefetch Buffer
6. I-Miss from Cache causes Prefetch Buffer Hit
7. I-Miss moved from Prefetch Buffer to Cache
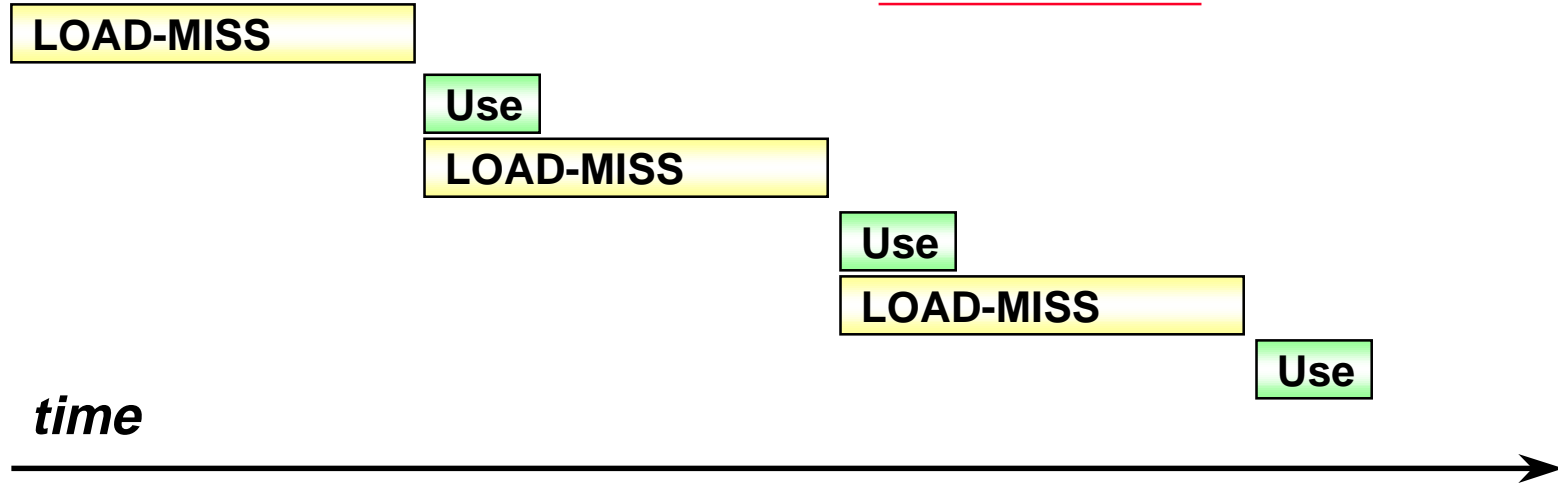8. I-Prefetch issued to Runway Bus (next line)

**Instruction Reorder Buffer ("Queue")**

**Instruction Fetch Unit**

**Tags**

**System Bus Interface**

**Prefetch Buffer**

**Instruction Cache**

**Runway Bus**

HEWLETT PACKARD

# Reorder Buffers

**System Bus Interface**

**Runway Bus** →

**Instruction Reorder Buffer ("Queue") (56 entries)**

**From I-Fetch** →

**Load/Store Address Adder**

**Load/Store Address Adder**

**Address Reorder Buffer (28 entries)**

**Data Cache (1MB)**

## Cycle by cycle progression of a load instruction

| Insert | Launch | Address | Cache | Cache | RR | Retire |
|--------|--------|---------|-------|-------|-----|--------|

# LOAD-MISS Overlapping

*The Problem*

| LOAD-MISS |

| Use |

| LOAD-MISS |

| Use |

| LOAD-MISS |

| Use |

*time* →

*PA-8500 Solution*

| LOAD-MISS |

| LOAD-MISS |

| LOAD-MISS |

| Use |

| Use |

| Use |

# Address Reorder Buffer: High-Speed Custom Circuitry

# Data Prefetching

LOAD-to-GR0

Instr

Instr

Instr

Instr

Instr

Instr

Instr

Instr

Instr

Instr

Instr

*time*

LOAD-HIT

*The Problem*
*Avoid the LOAD-MISS latency*

*Solution*
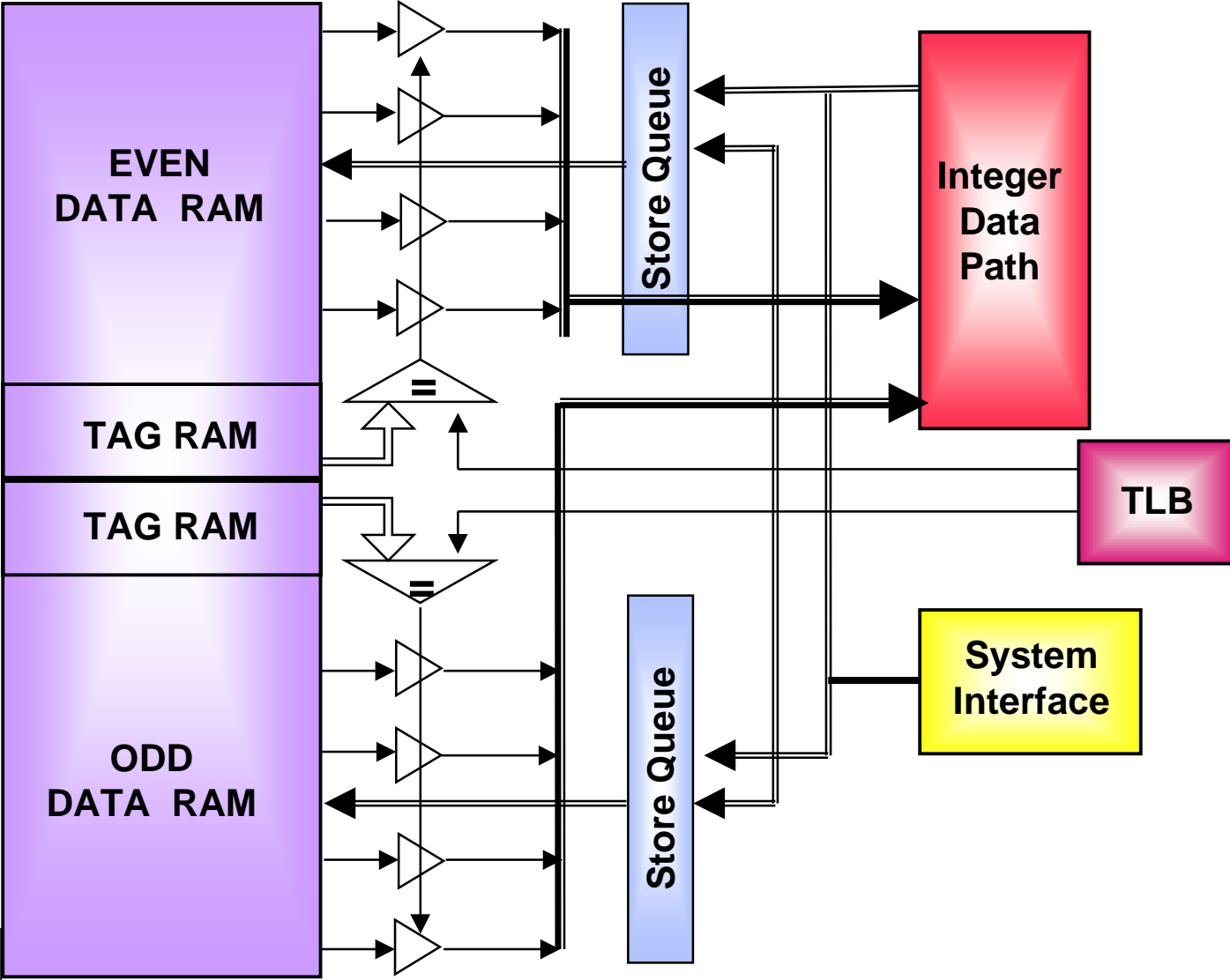*Compiler inserts Prefetch instruction*
*(LOAD to GR0)*
*Independent instructions executed*
*("Instr")*
*Data is resident in cache for LOAD*
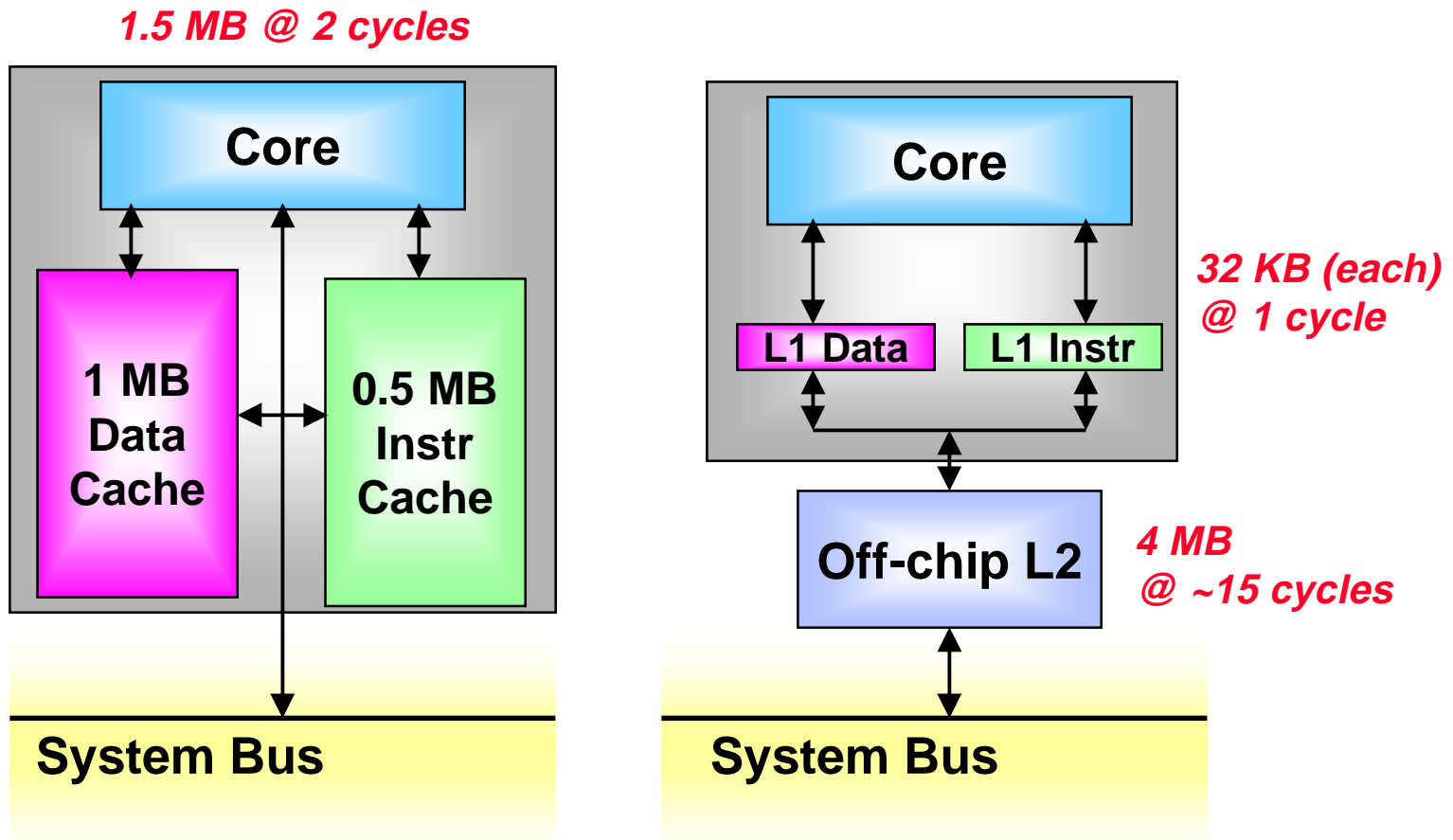*(LOAD-HIT)*

HEWLETT PACKARD

# Data Cache Features

- 1.0 MB on-chip cache
- 4-way set associative
- 2-cycle pipelined access
- Two accesses per cycle
- Supports 32-byte and 64-byte line sizes
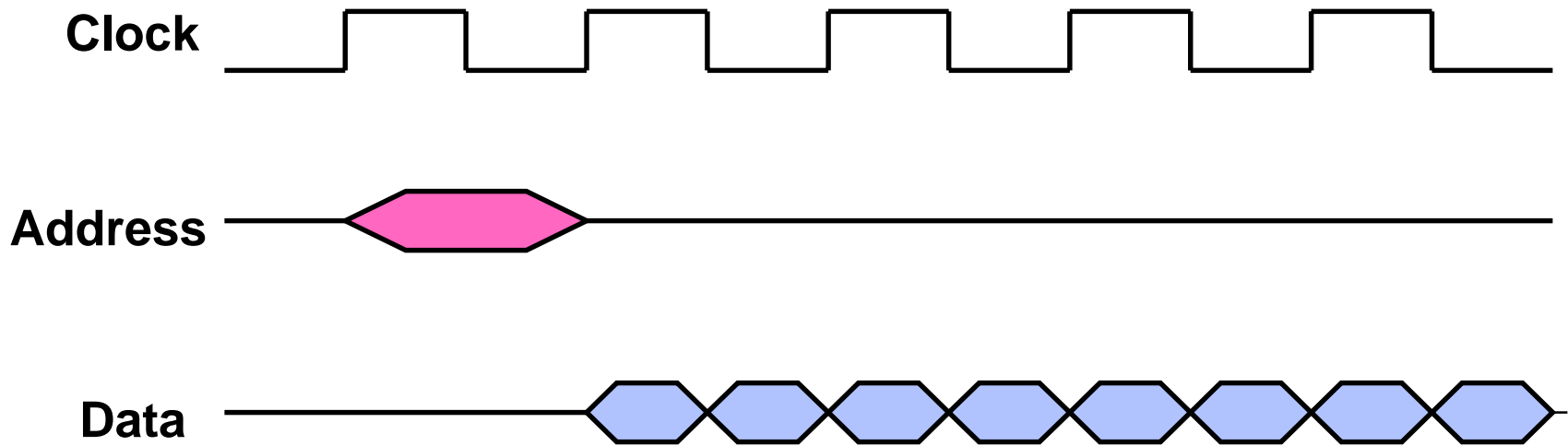- Sophisticated Store Queue

# Data Cache

# Single-Level vs. Multi-Level Cache Designs



*1.5 MB @ 2 cycles*

Core

1 MB Data Cache

0.5 MB Instr Cache

System Bus

Core

*32 KB (each) @ 1 cycle*

L1 Data

L1 Instr

Off-chip L2

*4 MB @ ~15 cycles*

System Bus

HEWLETT PACKARD

# System Bus Interface

- Split-transaction bus with out-of-order returns

- Multiple transactions in flight simultaneously

- Priority given to latency-sensitive transactions

- Asynchronous Interface

- Turbo Mode

# Turbo Mode

**Clock**

**Address**

**Data**

*High-Speed Data Transfer between Memory and CPU*

# Mitigating Memory Latency Effects

- Large Caches

- Out-of-Order Queue

- Flexible System Interface

- Custom Circuit Design

*The PA-8500 Achieves Superb Performance !*

HEWLETT PACKARD