CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Contemporary DRAM Architectures and Beyond

**Bruce Jacob**

**Electrical & Computer Engineering**
**University of Maryland, College Park**
`http://www.ece.umd.edu/~blj/`

**OUTLINE:**

- **Motivation & Background**

- **Experiments**

- **Results**

- **More Recent Results**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Sources

**"A Performance Study of Contemporary DRAM Architectures," *Proc. ISCA '99.* V. Cuppu, B. Jacob, B. Davis, and T. Mudge**

**Recent experiments by Vinodh Cuppu, Ph.D. student at University of Maryland**

**Recent experiments by Brian Davis, Ph.D. student at University of Michigan**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Dilemma: THIS ...

**STATUS QUO in
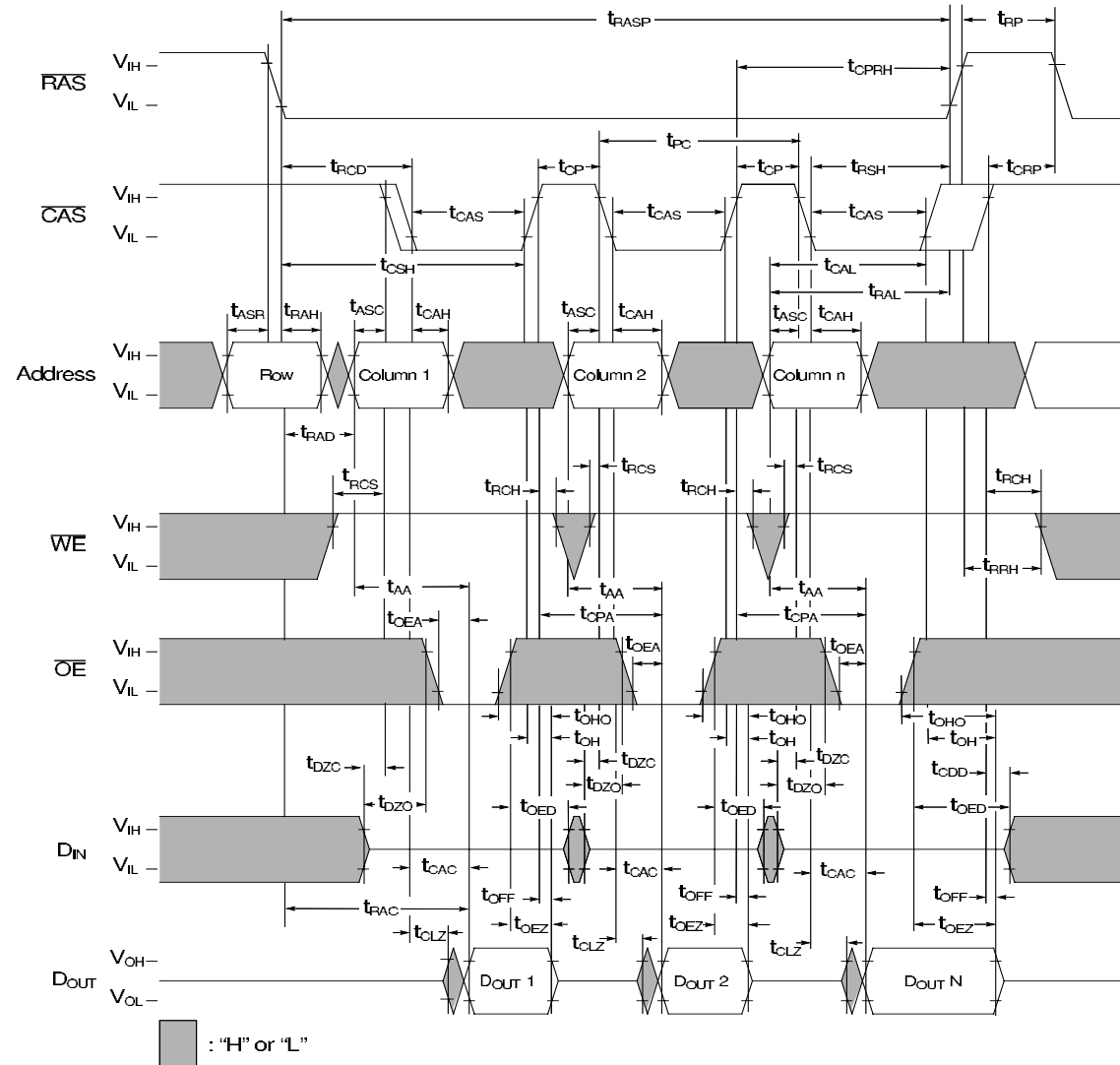MEMORY-SYSTEM RESEARCH:**

```
...

if (memory_instruction(INSTR)) {
    if (L1_cache_miss( data_addr(INSTR) ){
        if (L2_cache_miss( data_addr(INSTR) ){

            cycles += DRAM_LATENCY;

        }
    }
}

...
```
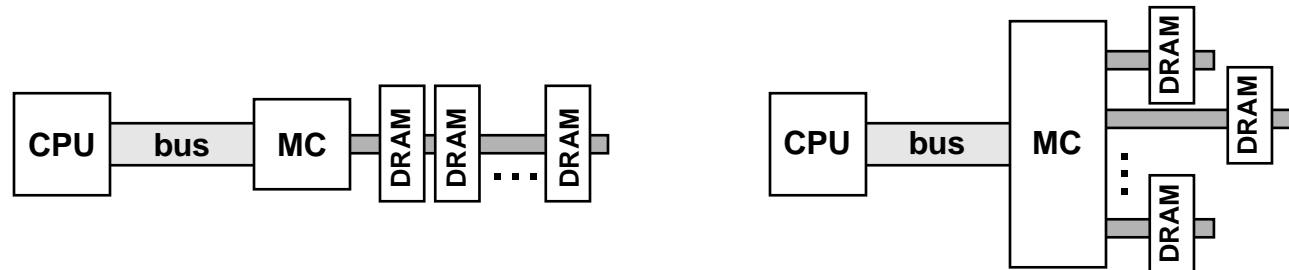
CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# ... or THIS

**Fast Page Mode Read Cycle**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Motivation

## HERE'S WHAT YOU MISS:



## DRAM LATENCY:



DATA TRANSFER

OVERLAP

COLUMN ACCESS

ROW ACCESS

BUS TRANSMISSION

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

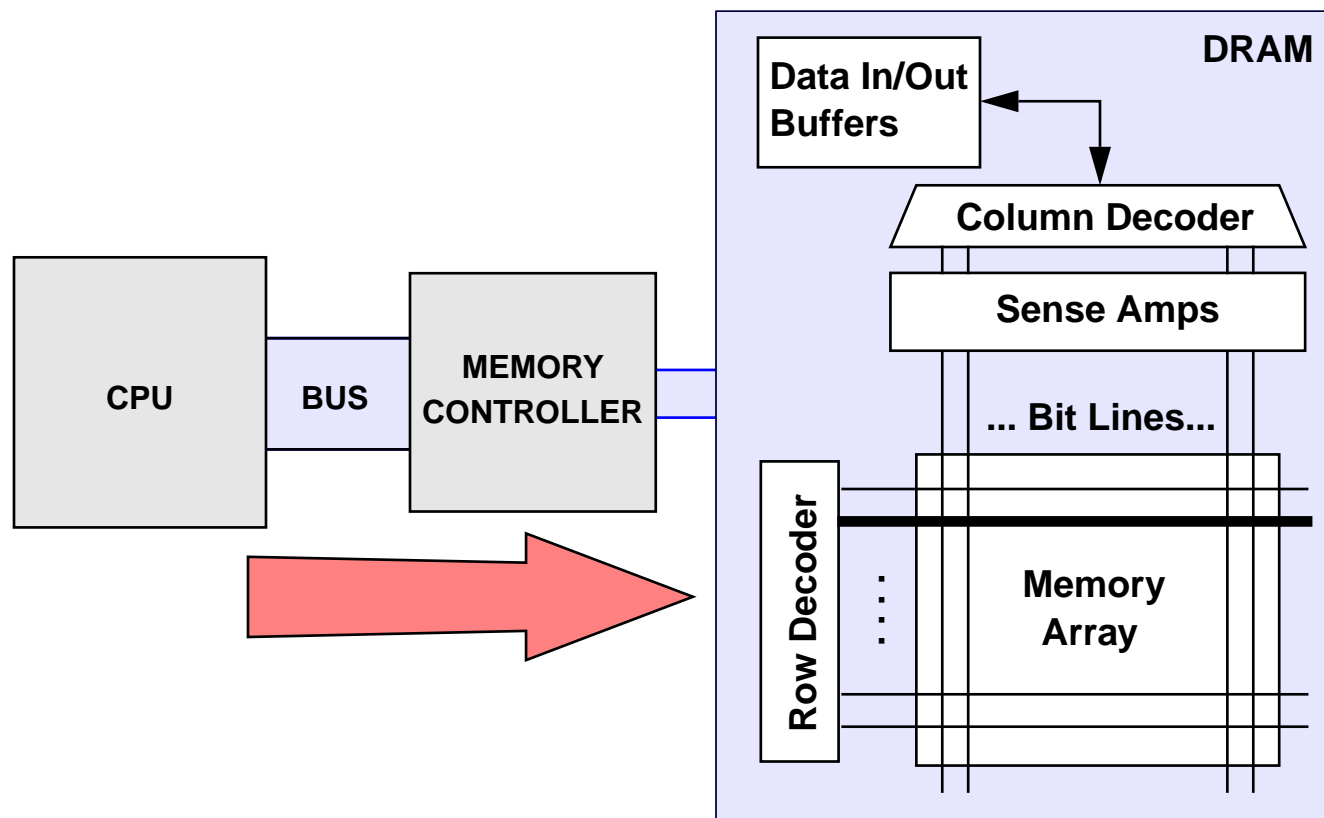Bruce Jacob

University of
Maryland

# Goal

## PRELIMINARY DRAM STUDY:

- **Bus Transmission**

- **Row Access**

- **Column Access**

- **Data Transfer**

- **Bus Wait/Synch Time**

- **Stalls Due to Refresh**

- **The OVERLAP of These Components
  (with each other)
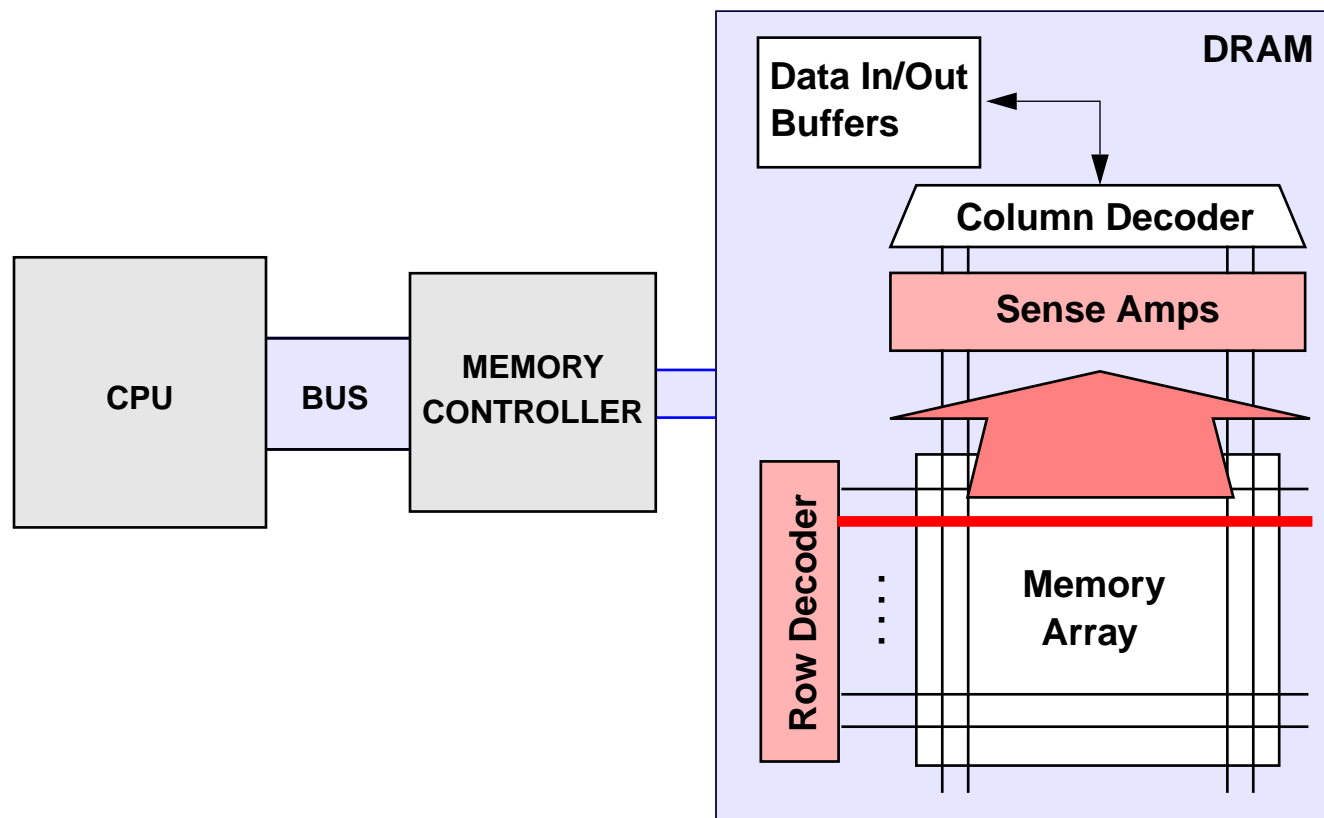  (with CPU execution)**

## MODEL EXISTING TECHNOLOGY

CONTEMPORARY
DRAM
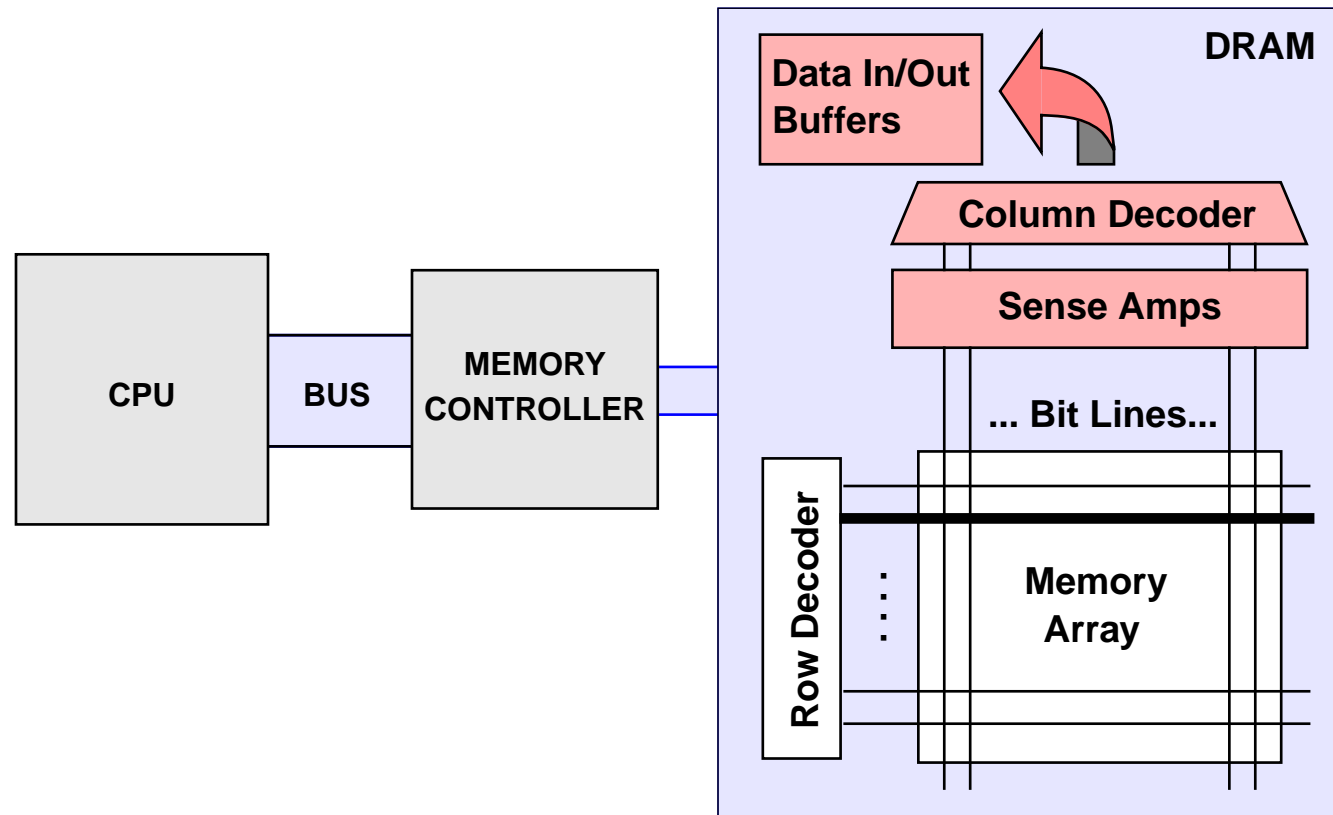ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## BUS TRANSMISSION

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## ROW ACCESS

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## COLUMN ACCESS

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## DATA TRANSFER

**DRAM**

**Data In/Out Buffers**

**Column Decoder**

**Sense Amps**

**... Bit Lines...**

**Row Decoder**

**Memory Array**

**CPU**

**BUS**

**MEMORY CONTROLLER**

## note: page mode enables overlap with COL

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## BUS TRANSMISSION

DRAM

Data In/Out Buffers

Column Decoder

Sense Amps

... Bit Lines...

Row Decoder

Memory Array

CPU

BUS

MEMORY CONTROLLER

## note: overlapped component not shown

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## Read Timing for Conventional DRAM

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## Read Timing for Fast Page Mode DRAM

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## Read Timing for Extended Data Out DRAM

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## Read Timing for Synchronous DRAM

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Primer

## Read Timing for Rambus DRAM

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Simulator Overview

## CPU: SimpleScalar v3.0a

- **8-way out-of-order**

- **L1 cache: split 64K/64K, lockup free x32**

- **L2 cache: unified 1MB, lockup free x1**
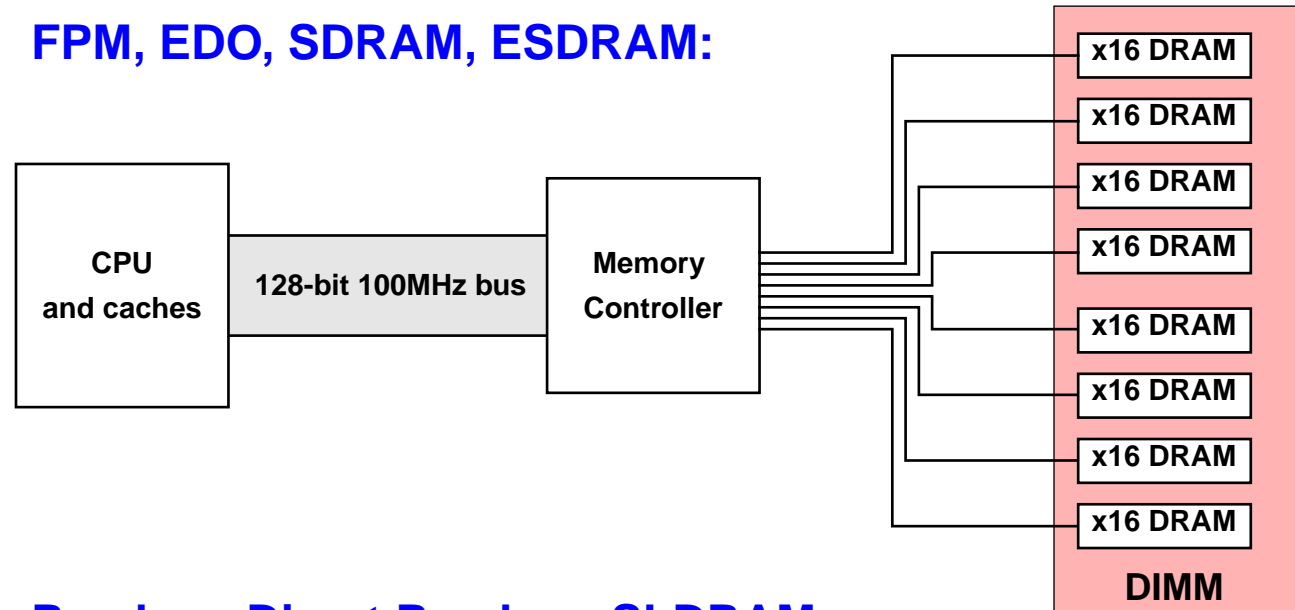
- **L2 blocksize: 128 bytes**

## Main Memory: 8 64Mb DRAMs

- **100MHz/128-bit memory bus**

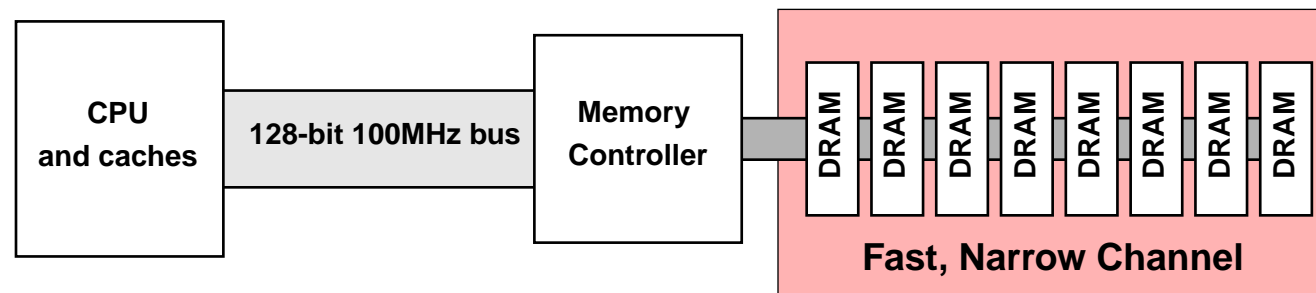- **Optimistic *open-page* policy (*close-immediately* can be calculated)**

## Represents a "typical" workstation

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Configurations

**FPM, EDO, SDRAM, ESDRAM:**

```
CPU
and caches  —[128-bit 100MHz bus]—  Memory
                                     Controller
```

x16 DRAM
x16 DRAM
x16 DRAM
x16 DRAM
x16 DRAM
x16 DRAM
x16 DRAM
x16 DRAM

**DIMM**

**Rambus, Direct Rambus, SLDRAM:**

```
CPU
and caches  —[128-bit 100MHz bus]—  Memory
                                     Controller
```

DRAM DRAM DRAM DRAM DRAM DRAM DRAM DRAM

**Fast, Narrow Channel**

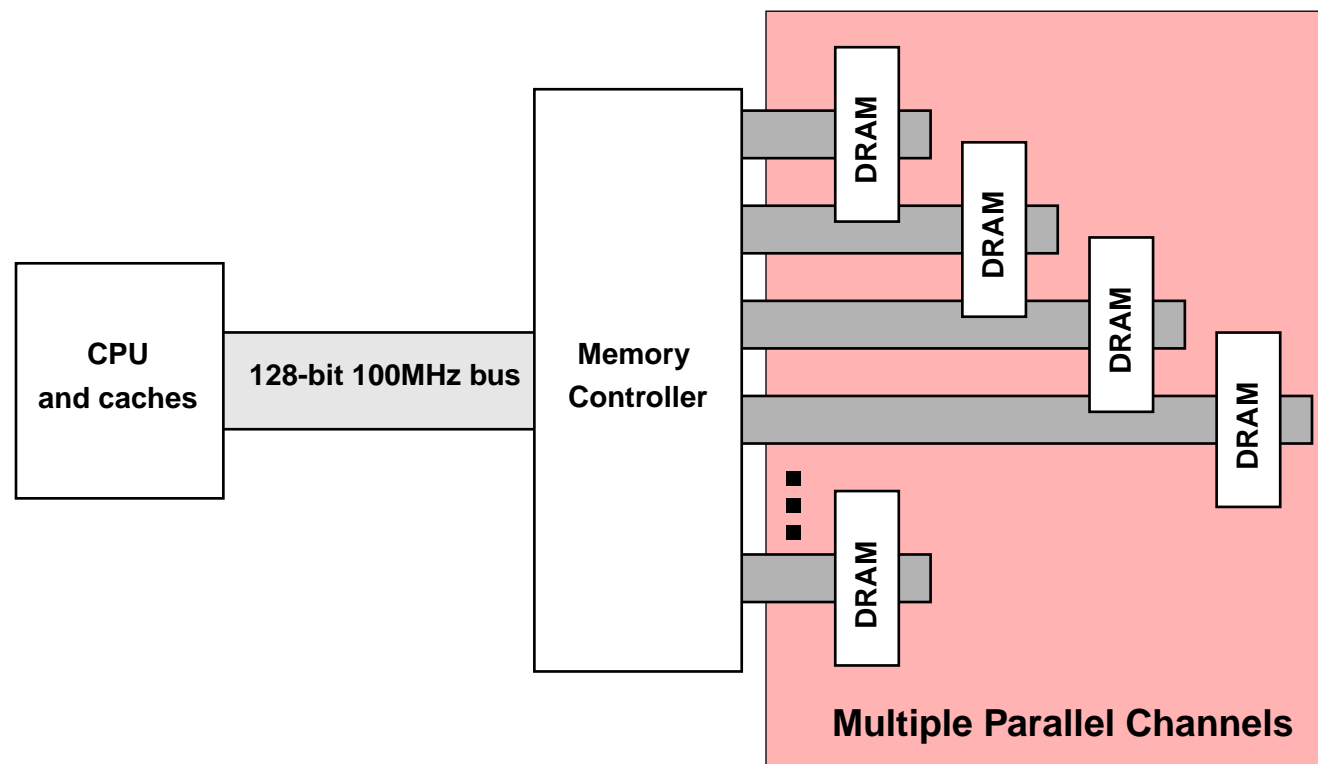**Note: TRANSFER WIDTH of Direct Rambus Channel**

- **equals that of ganged FPM, EDO, etc.**
- **is 2x that of Rambus & SLDRAM**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# DRAM Configurations

## Strawman: Rambus, etc.

CPU and caches

**128-bit 100MHz bus**

Memory Controller

DRAM

DRAM

DRAM

DRAM

DRAM

**Multiple Parallel Channels**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland



# Overhead: Memory vs. CPU

Variable: speed of processor & caches

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
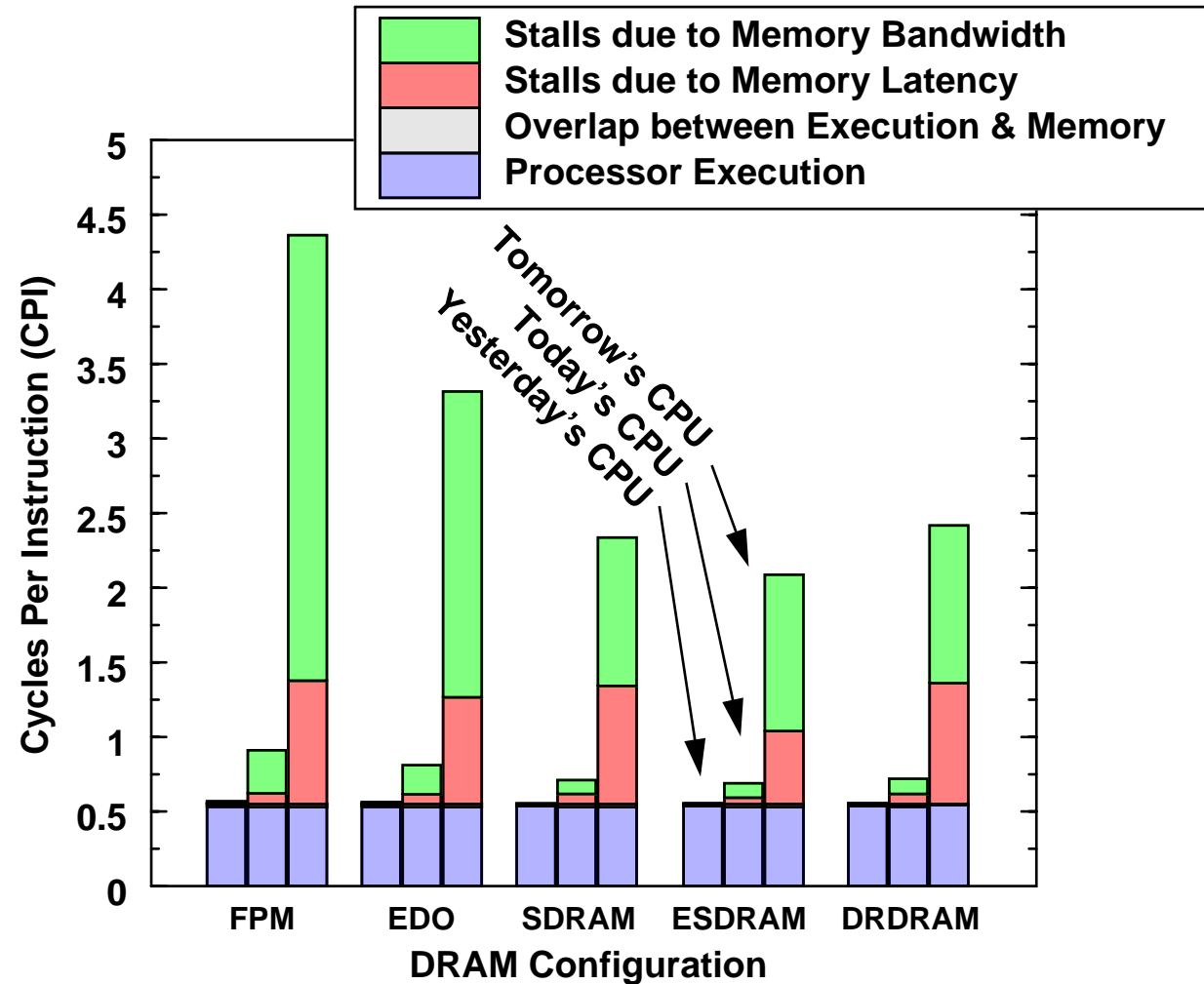Maryland

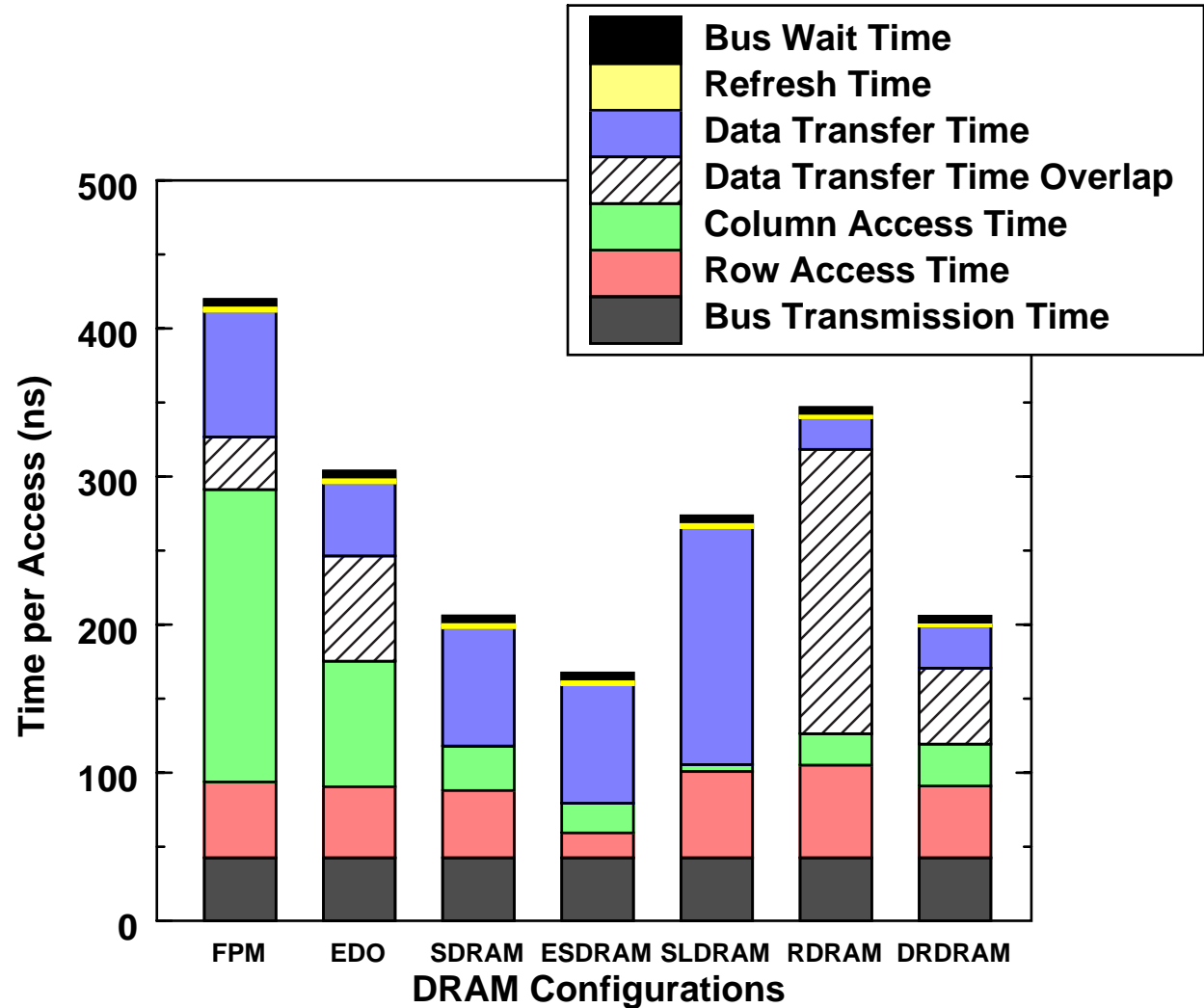# Definitions (var. on Burger, et al)

- $t_{PROC}$ — processor with perfect memory
- $t_{REAL}$ — realistic configuration
- $t_{BW}$ — CPU with wide memory paths
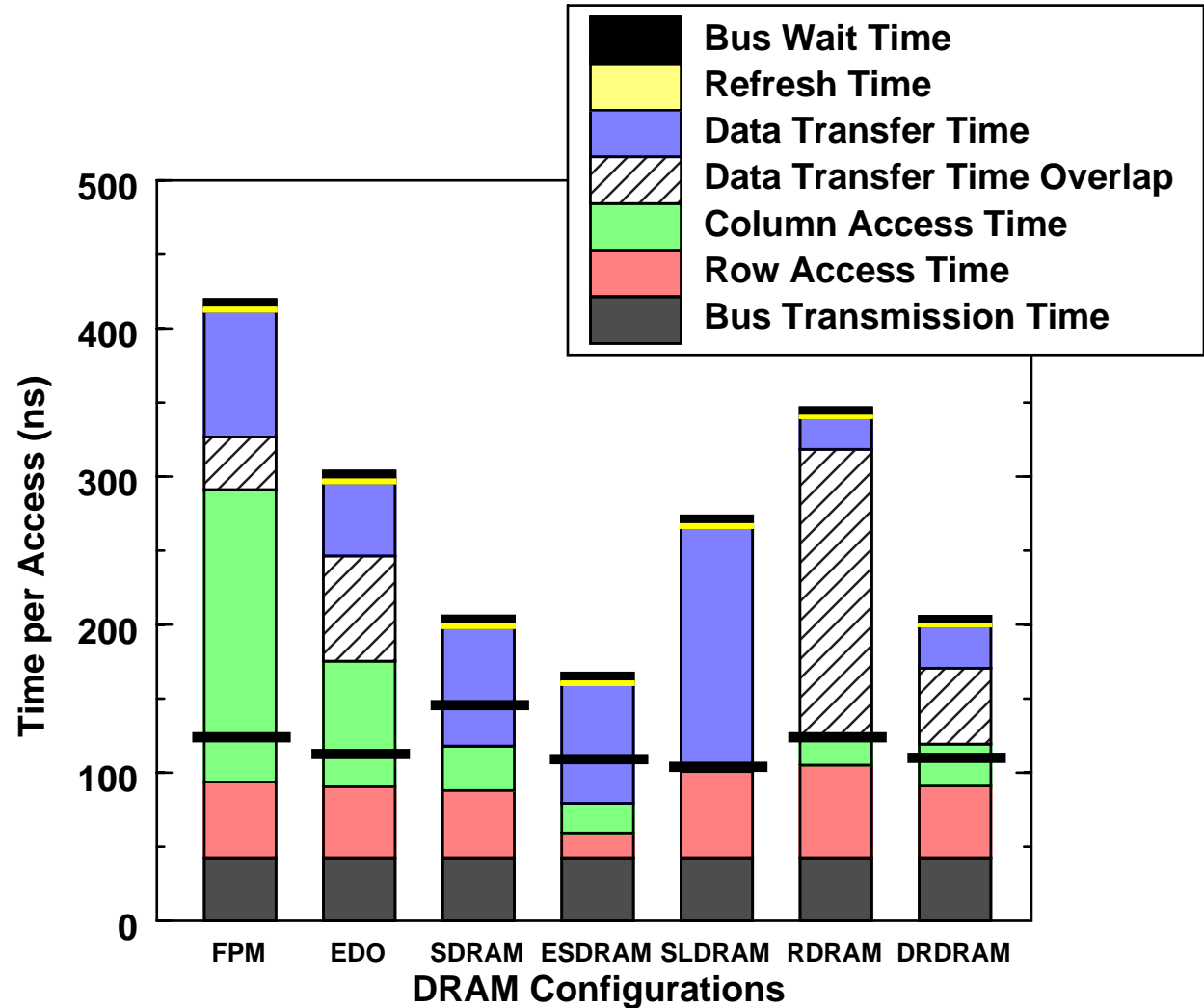- $t_{DRAM}$ — time seen by DRAM system

**Stalls Due to BANDWIDTH** — $t_{REAL} - t_{BW}$

**Stalls Due to LATENCY** — $t_{BW} - t_{PROC}$

**CPU-Memory OVERLAP** — $t_{PROC} - (t_{REAL} - t_{DRAM})$

**CPU+L1+L2 Execution** — $t_{REAL} - t_{DRAM}$

$t_{REAL}$    $t_{DRAM}$

$t_{BW}$

$t_{PROC}$

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Average Latency of DRAMs



note: SLDRAM & RDRAM 2x data transfers

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Average Latency of DRAMs



note: SLDRAM & RDRAM 2x data transfers

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Cost-Performance

**FPM, EDO, SDRAM, ESDRAM:**

- **Lower Latency => Wide/Fast Bus**

- **Increase Capacity => Decrease Latency**

- **Low System Cost**

**Rambus, Direct Rambus, SLDRAM:**

- **Lower Latency => Multiple Channels**

- **Increase Capacity => Increase Capacity**

- **High System Cost**

**1 DRDRAM = Multiple SDRAM**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Conclusions

**100MHz/128-bit Bus is Current Bottleneck**

- **Solution: Fast Bus/es & MC on CPU
(*e.g.* Compaq Alpha, Sony Emotion, ...)**

**Current DRAMs Solving Bandwidth Problem
(but not Latency Problem)**

**There is Locality in DRAM Accesses
(but how important is this?)**

**SPECint '95 Fits in 1MB Cache**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Recent (Unfinished) Work

**Investigation of Organization-Level Parameters:**

- **Channel widths & speeds, turnaround**

- **Independent vs. ganged channels**

- **Banks per channel, burst widths**

**Detailed Study of DRDRAM vs. SDRAM in Highly Concurrent Environment**
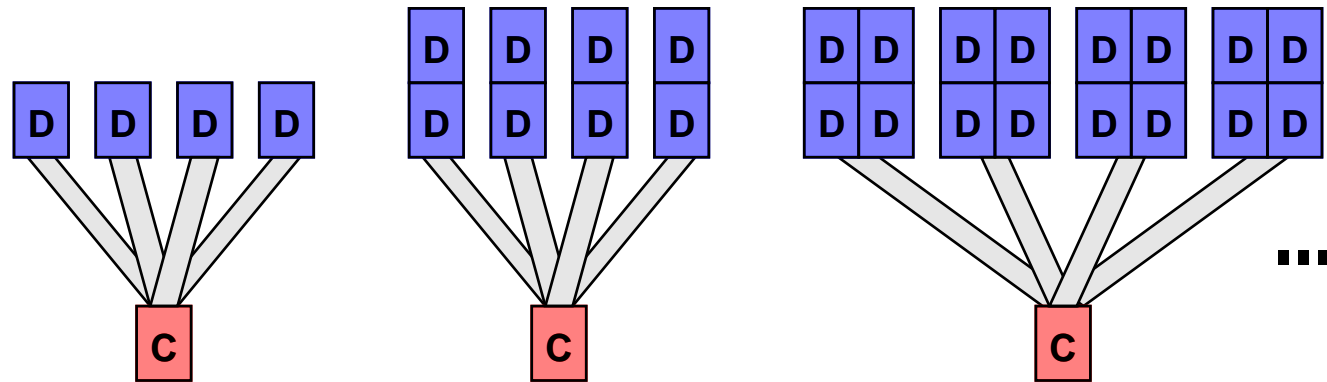
**Embedded DRAM+DSP Architectures**

**Detailed Study of Multiprocessor Buses**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Channel/Bank Model



**One independent channel**
**Banking degrees of 1, 2, 4, ...**

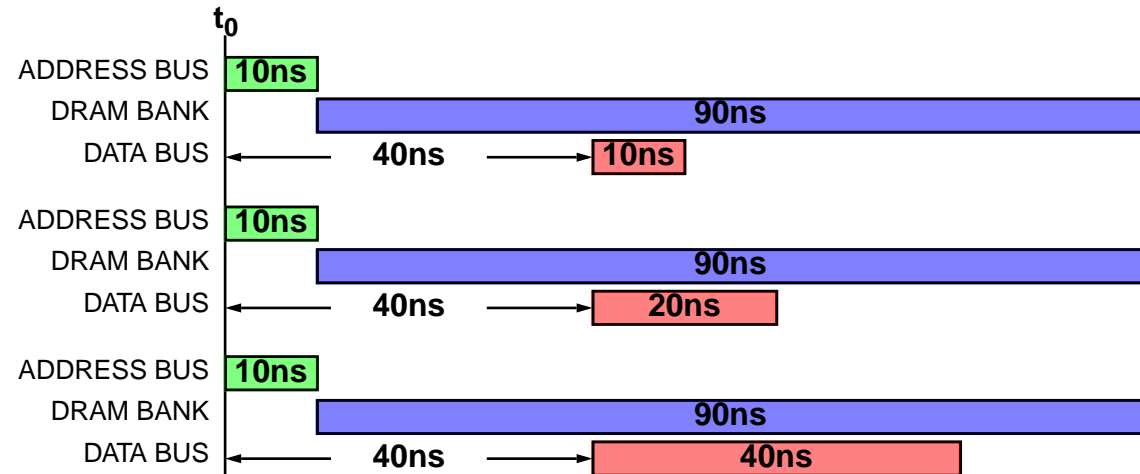**Two independent channels**
**Banking degrees of 1, 2, 4, ...**

**Four independent channels**
**Banking degrees of 1, 2, 4, ...**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Read/Write Request Model

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Bandwidth vs. Burst Width



**Legend:**
- 8-Byte Burst Width
- 16-Byte Burst Width
- 32-Byte Burst Width
- 64-Byte Burst Width
- 128-Byte Burst Width

Y-axis: Cycles per Instruction (0, 0.25, 0.5, 0.75, 1, 1.25)

X-axis: System Bandwidth (GB/s = Channels * Width * Speed) — 0.4, 0.8, 1.6, 3.2, 6.4

**PERL: 1 channel, 4 banks, 2GHz CPU**

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Exploiting Concurrency



PERL: 2 banks, 16-byte burst, 2GHz CPU

CONTEMPORARY
DRAM
ARCHITECTURES
AND BEYOND

Bruce Jacob

University of
Maryland

# Conclusions

**None yet ... preliminary data**

**CONTACT INFO:**

**Prof. Bruce Jacob**

**Electrical & Computer Engineering**
**University of Maryland, College Park**
`http://www.ece.umd.edu/~blj/`
`blj@eng.umd.edu`