# New Dimensions in Performance

## Harnessing 3D Integration Technologies

**Kerry Bernstein**

**IBM T.J. Watson Research Center**
**Yorktown Heights, NY**

**9 September, 2008**          **Fort Collins, CO**

IBM

# 3D Press Release

**Wall Street Journal: IBM Touts Breakthrough in 3-D Chips**

**By WILLIAM M. BULKELEY**

**April 12, 2007; Page B3**

**International Business Machines Corp. said it achieved a breakthrough in developing a three-dimensional semiconductor chip that can be stacked on top of another electronic device in a vertical configuration long sought by engineers to reduce size and power use.**

# Agenda

1) **Workloads**

2) **The Memory Wall, Bandwidth, and Latency**

3) **The Technologies of 3D Integration**

4) **"The future ain't what it used to be."**
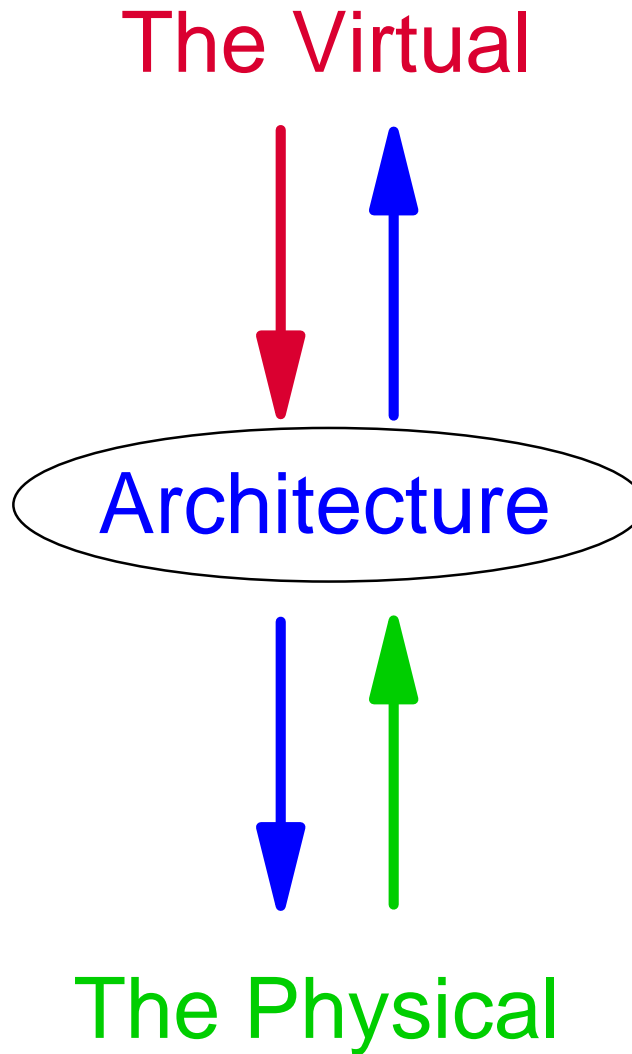
5) **Summary**

6) **Human Scaling**

# Computer Workloads & Thru-put

# Workloads – What do Computers do?

- **Scientific** (i.e. Lawrence Livermore National Labs)
  - Highly regular, predictable patters allows streaming data from cache to processor
  - Performance is directly proportional to bus bandwidth
  - High utilization, full data bus at all times

- **Commercial** (i.e. Starbucks)
  - Unpredictable irregular patterns
  - Miss rate follows Poisson process (random)
  - Requires low bus utilization to avoid clogs in the event of a burst of misses (usually 30% bus utilization)

- **Both application spaces need BW, but for different reasons**

**The Software Stack:**
- System
  - Hypervisor
  - Operating System
- Applications Layer
- Program
- Compiler
- Machine Language

The Virtual

**Architecture:**
A fully-specified unambiguous contract

Architecture

The Physical

**The Hardware Stack:**
- Logical Level Description
- Machine Organization
- Schematic Representation
- Circuit Design
- Physical Design
- Device Level (transistors)
- Atomic Level

# Processor Cores and Memory Subsystems
The New Units of Design

---

## (Systems/Thread) x (Threads/Core) x (Cores / Die)

**Puts pressure on Memory
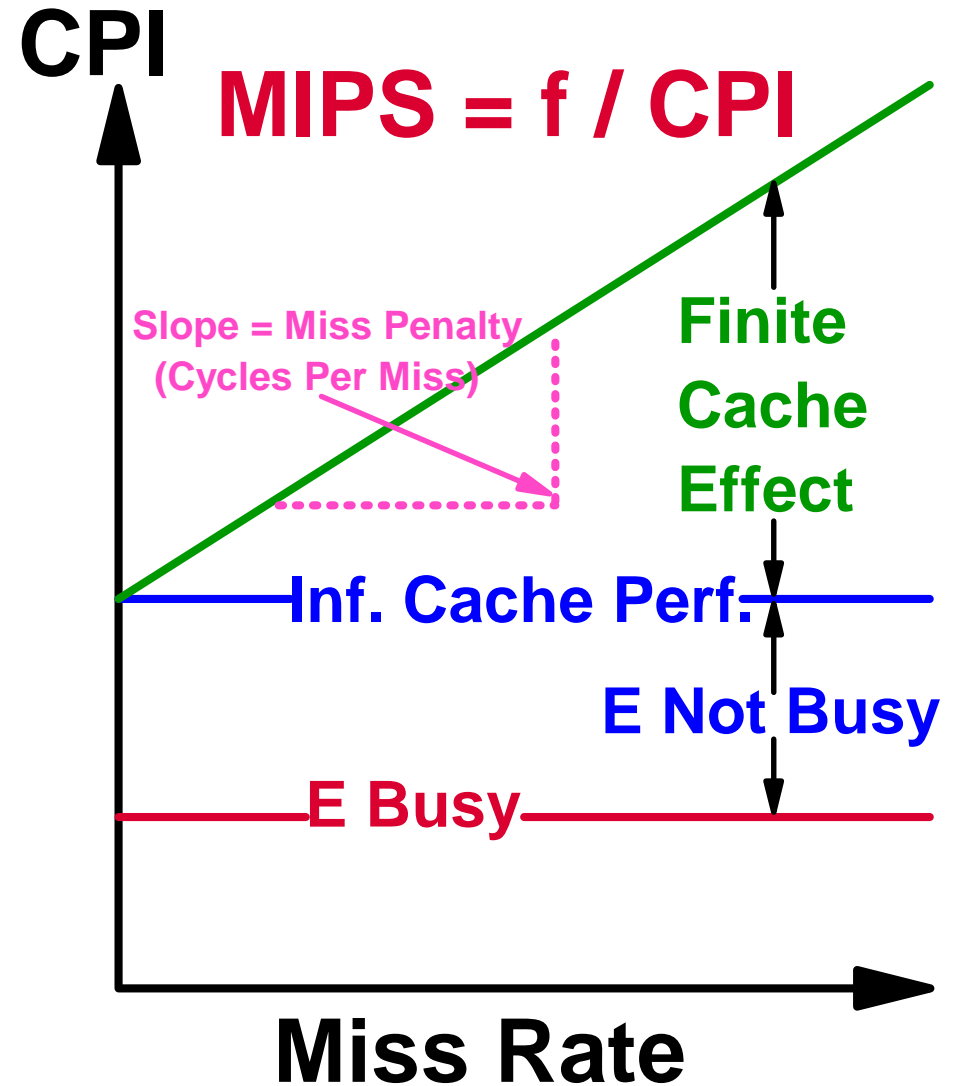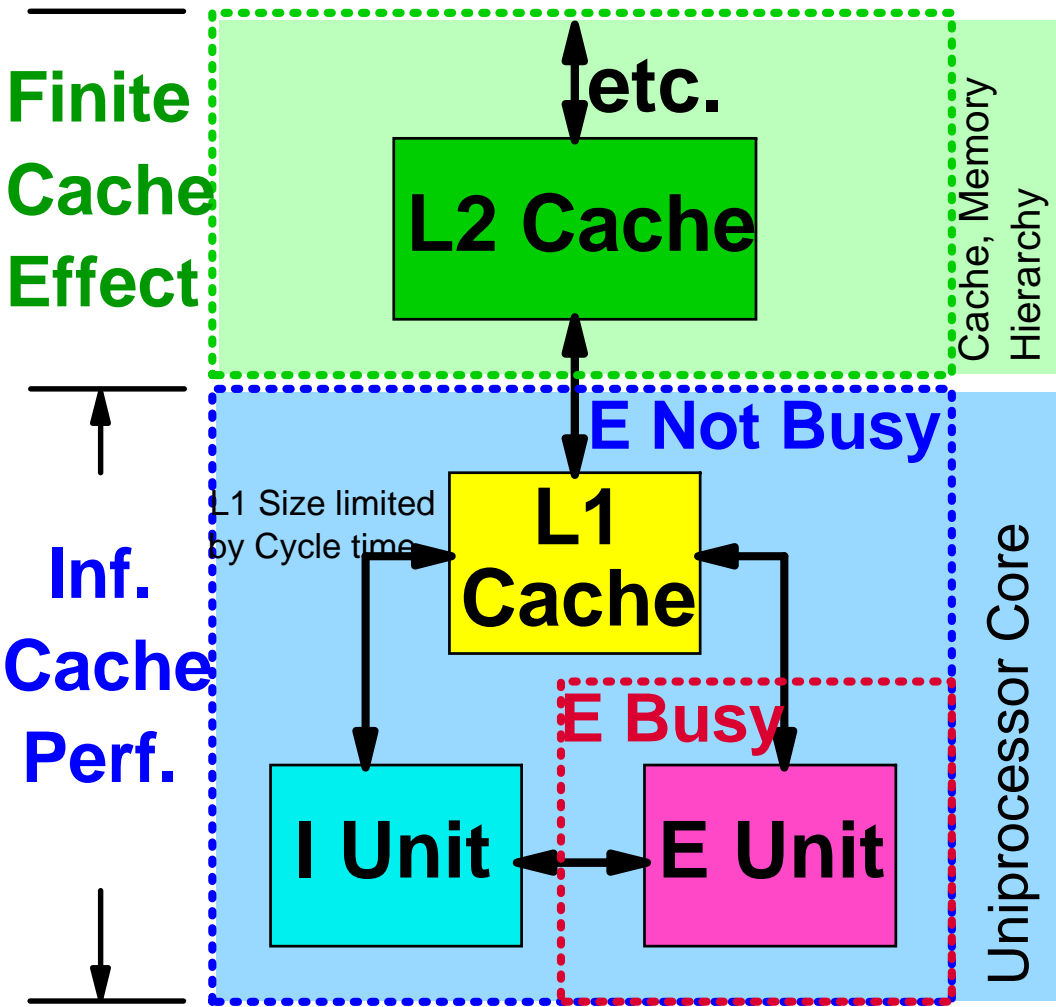Subsystem, Communication**

**On-chip content
(cache capacity)**

**Inter, Intra-chip BW**

Integration Focus moves from the device and circuit to core
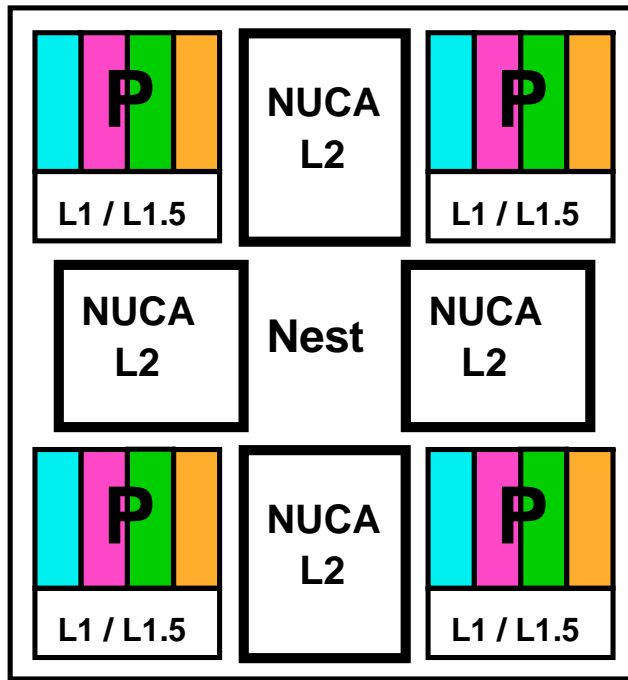
# Components of Processor Performance

**MIPS = f / CPI**

Delay is sequentially determined by a) ideal processor,
b) access to local cache, and c) refill of cache

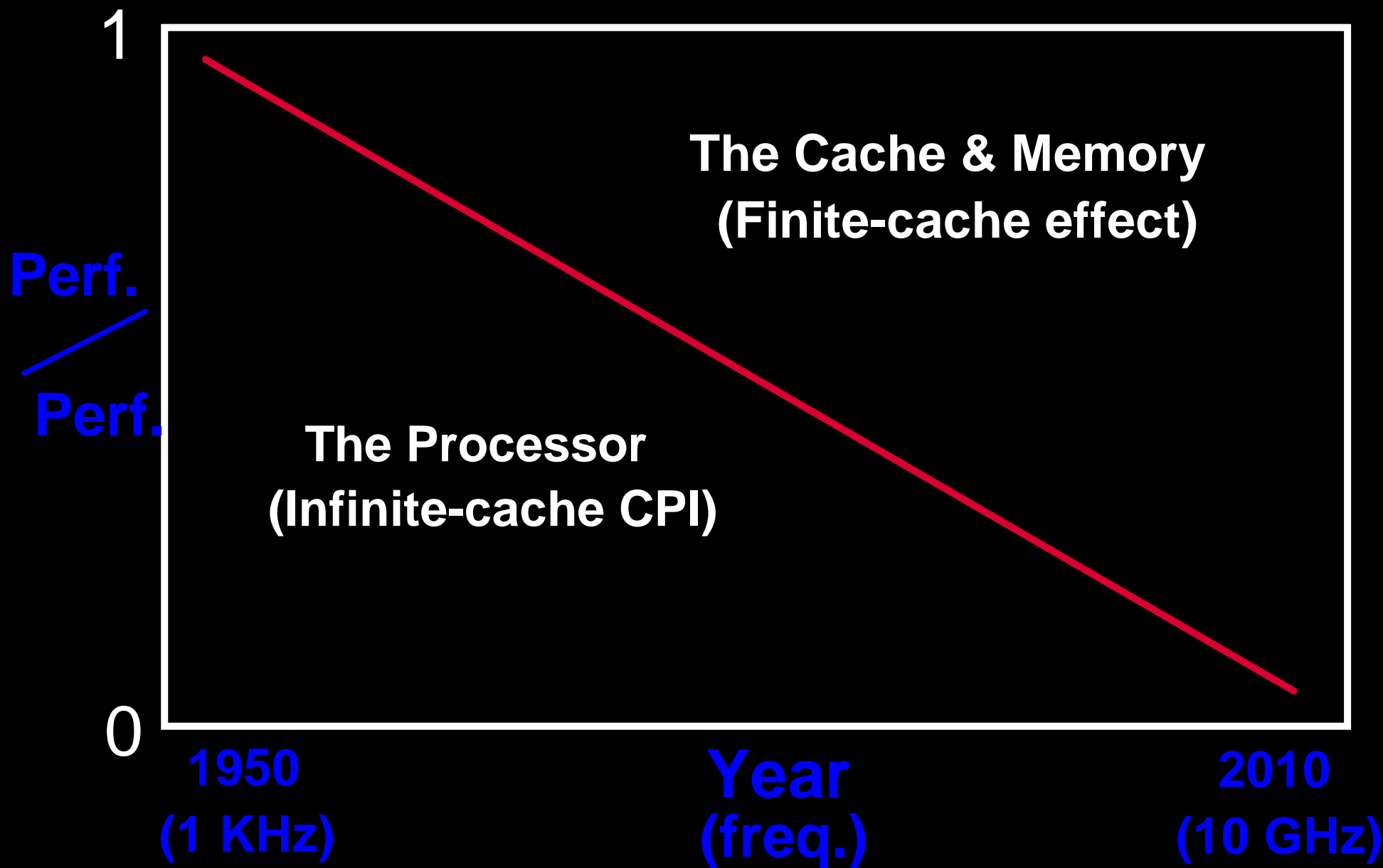# Virtualization
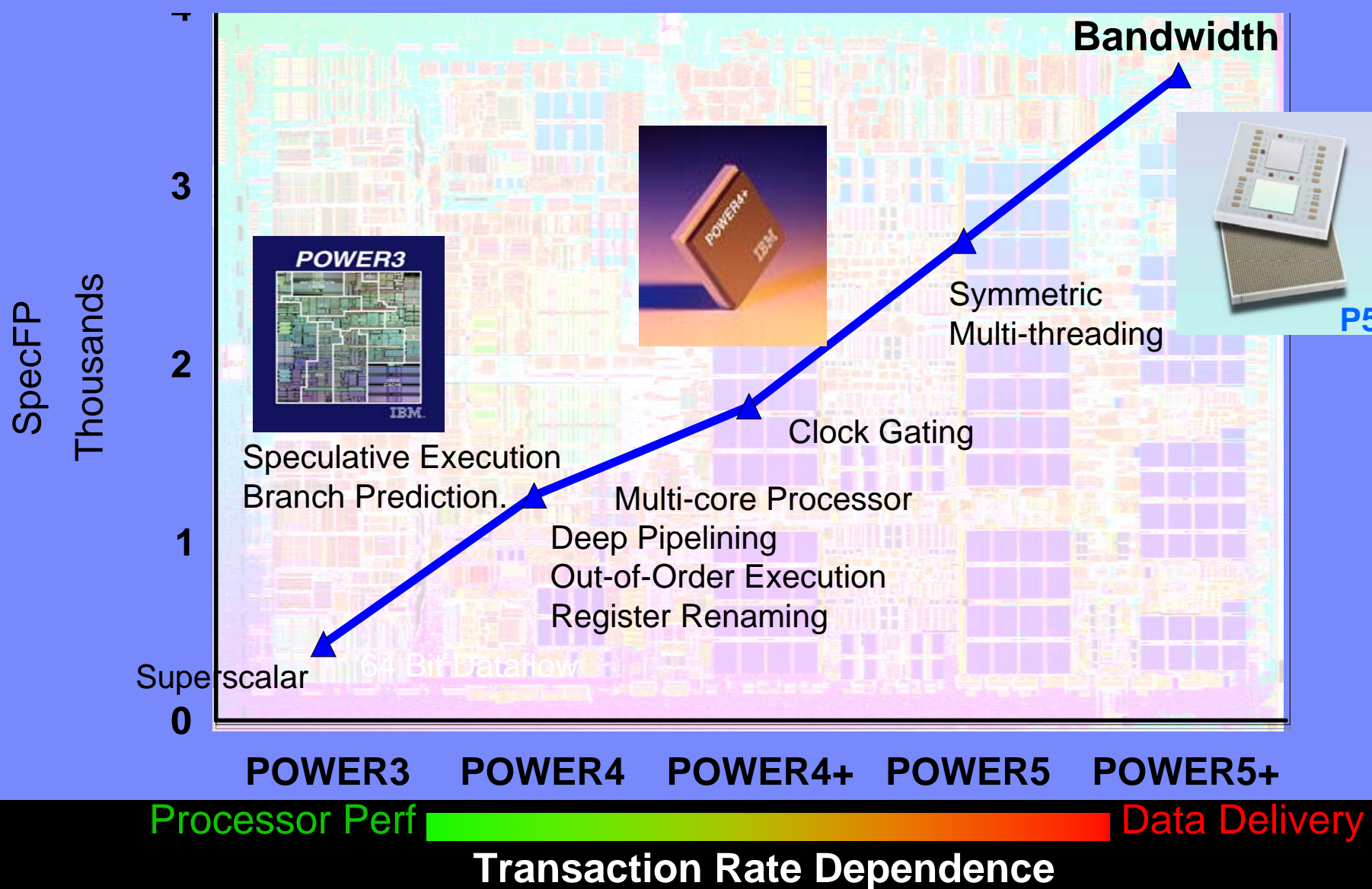
**Real System**

**Independent Virtual Systems**



**"Looks like" 4 independent systems, each with 16 cores!**

# What Dominates System Performance?
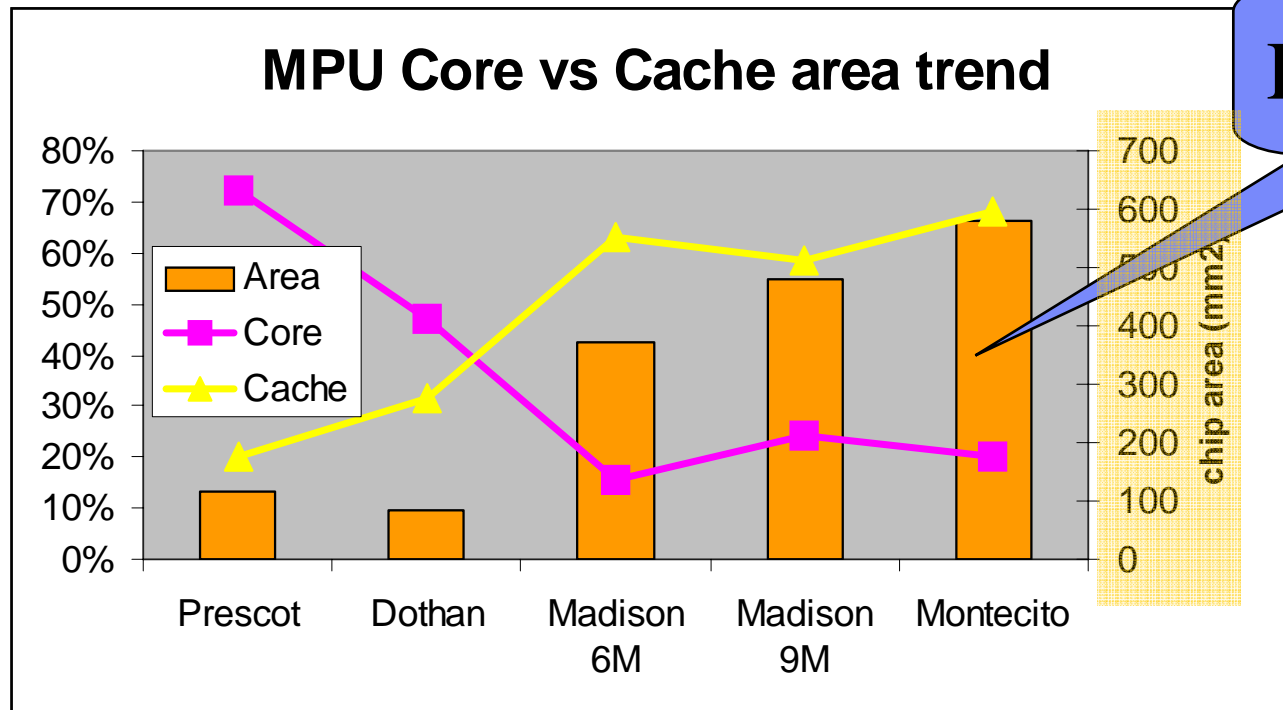
From ISCA '06
Keynote address by
Phil Emma, IBM



**The Cache & Memory
(Finite-cache effect)**

**The Processor
(Infinite-cache CPI)**

1

Perf.
/
Perf.

0

1950
(1 KHz)

Year
(freq.)

2010
(10 GHz)

# POWER Series Architectural Perf Contributions



SpecFP Thousands (y-axis: 0, 1, 2, 3)

**Bandwidth**

**Symmetric Multi-threading**

Clock Gating

Speculative Execution Branch Prediction.

Multi-core Processor
Deep Pipelining
Out-of-Order Execution
Register Renaming

Superscalar    64 Bit Dataflow

P5

**POWER3    POWER4    POWER4+    POWER5    POWER5+**

Processor Perf    Data Delivery
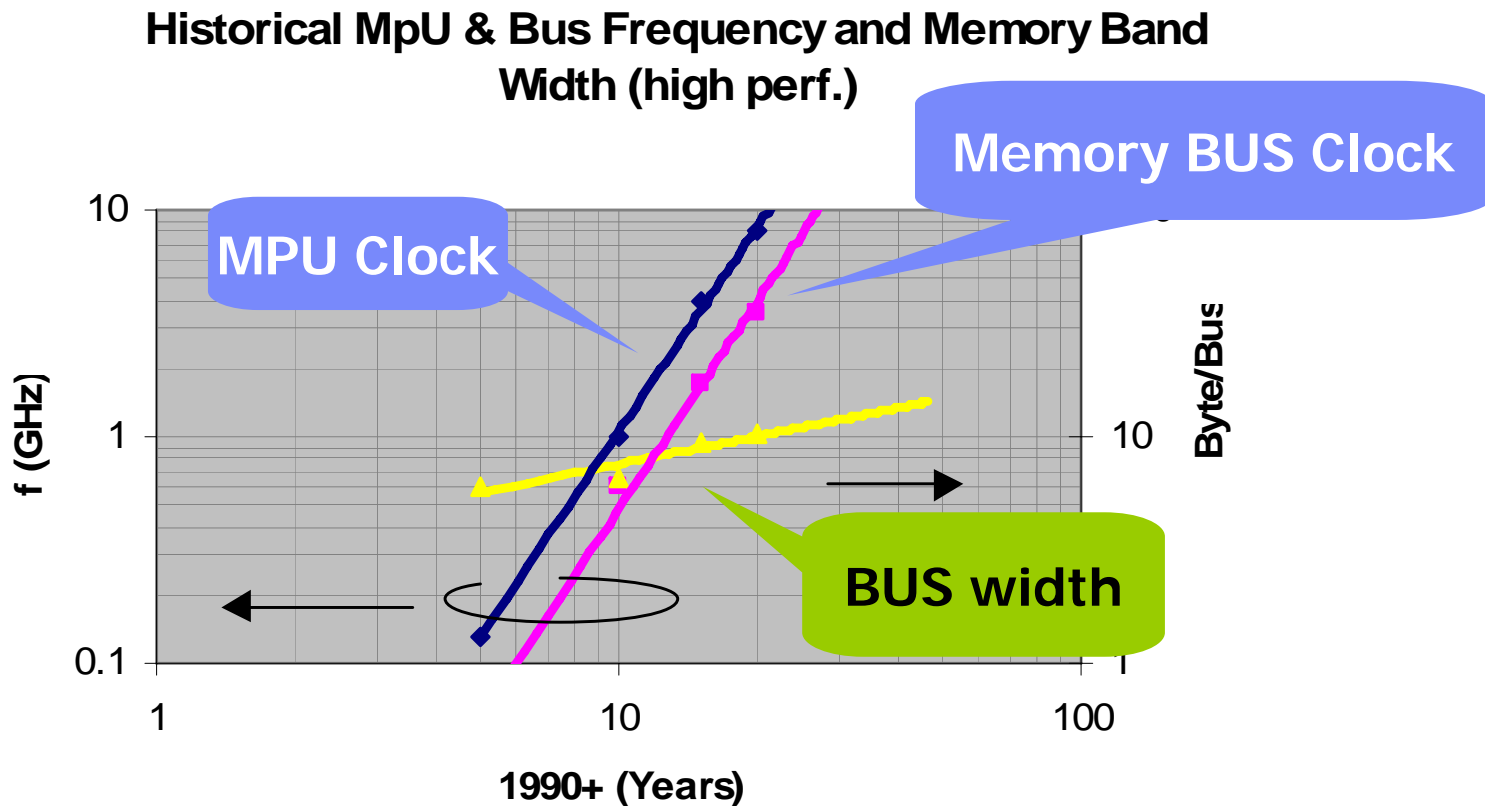
**Transaction Rate Dependence**

# Increasing Cache Size Drives Chip Size

Floorplans source: P. DeMone, "Sizing Up the
Super Heavyweights," *Real World Technologies*
Report, 9/17/2004



**MPU Core vs Cache area trend**

Dual Core

- Growing data sets will increasingly stress cache size
- Multi-core floor planning and SRAM concerns will halt cache size growth to maintain manageable chip size

# Frequency Drives Datarate



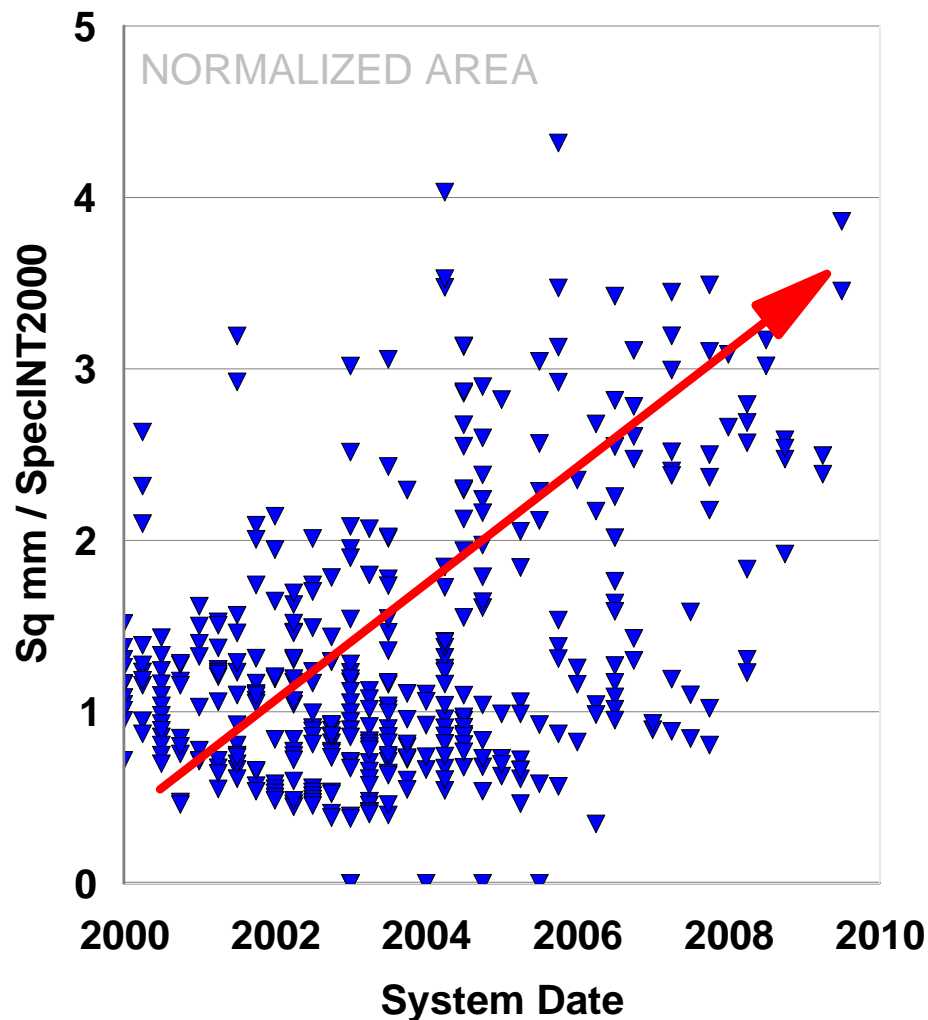**Historical MpU & Bus Frequency and Memory Band Width (high perf.)**

- Data bus frequency follows MPU frequency at a ratio 1:2 roughly doubling every 18 to 24 month
- Data bus band width shows only a moderate increase
  - Data bus transfer rate is basically scaled by bus frequency
- **When clock growth slows, BUS data rate growth will slow too!**

# Architecture Net

- **Growing the number of cores/chip increases demand for bandwidth**

- **Transaction retirement rate dependence on data delivery is *increasing***

- **Transaction retirement rate dependence on $\lambda$P performance is *decreasing***

# Die Area Increase



1) Architecture overhead increasing area of die

2) Accessible portion of chip over normalized cycle time is decreasing generation over generation
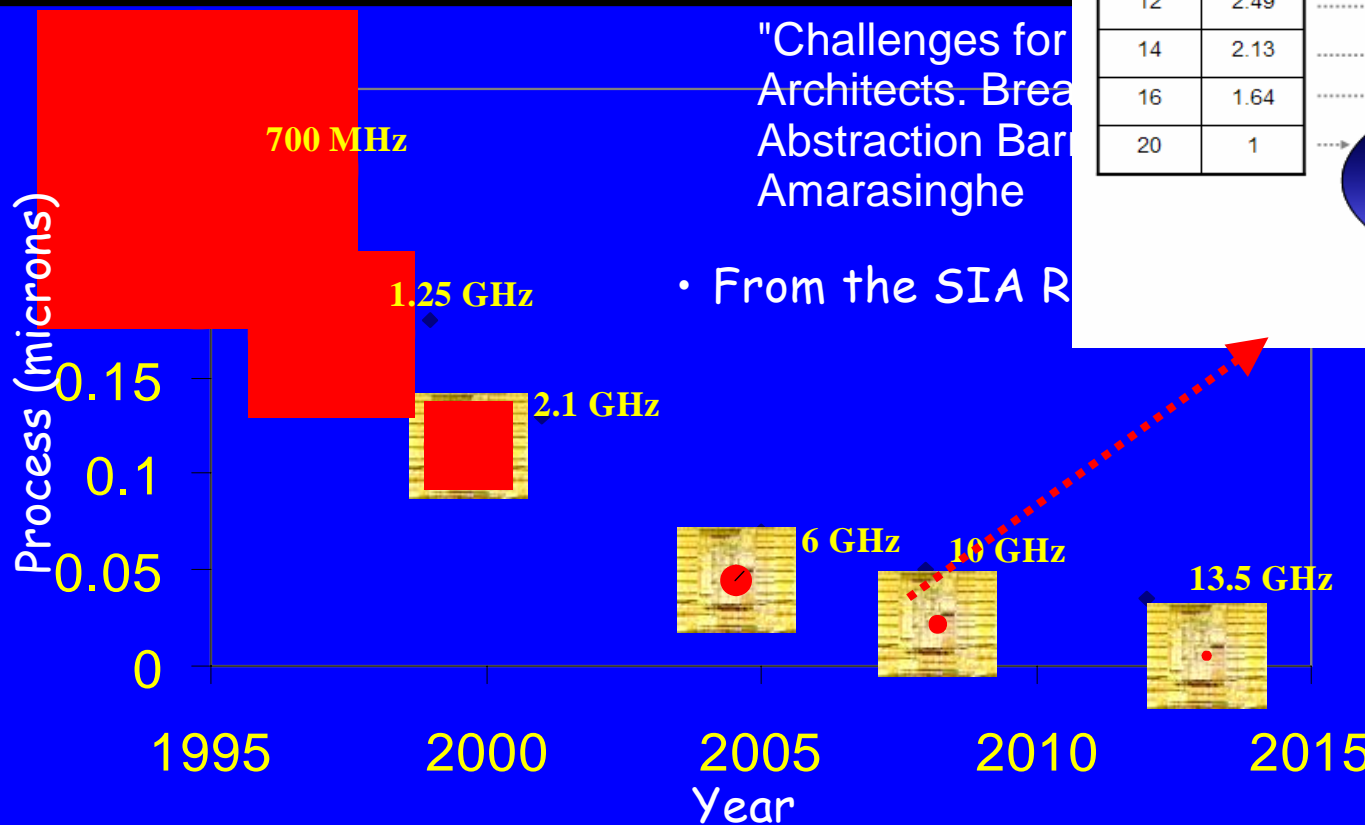
3) Deeper Pipes are decreasing delay per cycle

**Performance is expensive when left to architects!!**

# "Span of Control" with Scaling

Lack of wire delay improvement, die-size growth, and shorter relative cycle stage-depth together cause reduction in fan-out capability
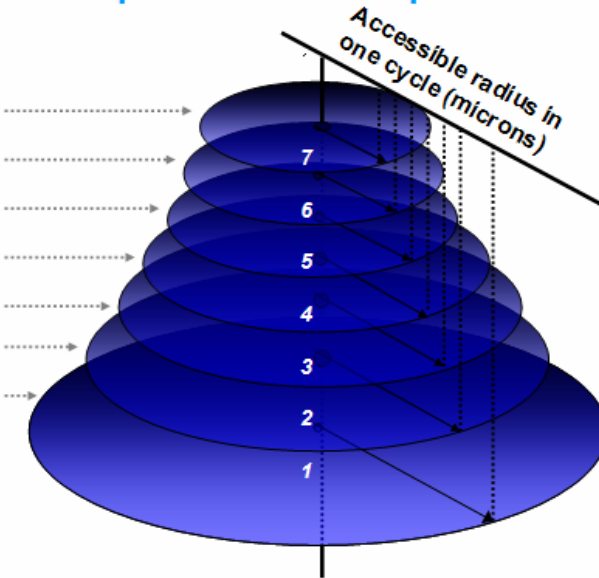
| No. of stage delays at this level | Span of control Imprvmt (Norm) |
|---|---|
| 8 | 3.10 |
| 9 | 2.94 |
| 10 | 2.74 |
| 12 | 2.49 |
| 14 | 2.13 |
| 16 | 1.64 |
| 20 | 1 |

**3DI Span of Control Improvement**

Accessible radius in one cycle (microns)

7
6
5
4
3
2
1

"Challenges for Architects. Brea Abstraction Bar Amarasinghe

• From the SIA R

700 MHz

1.25 GHz

2.1 GHz

6 GHz

10 GHz

13.5 GHz

Process (microns)

0.15
0.1
0.05
0

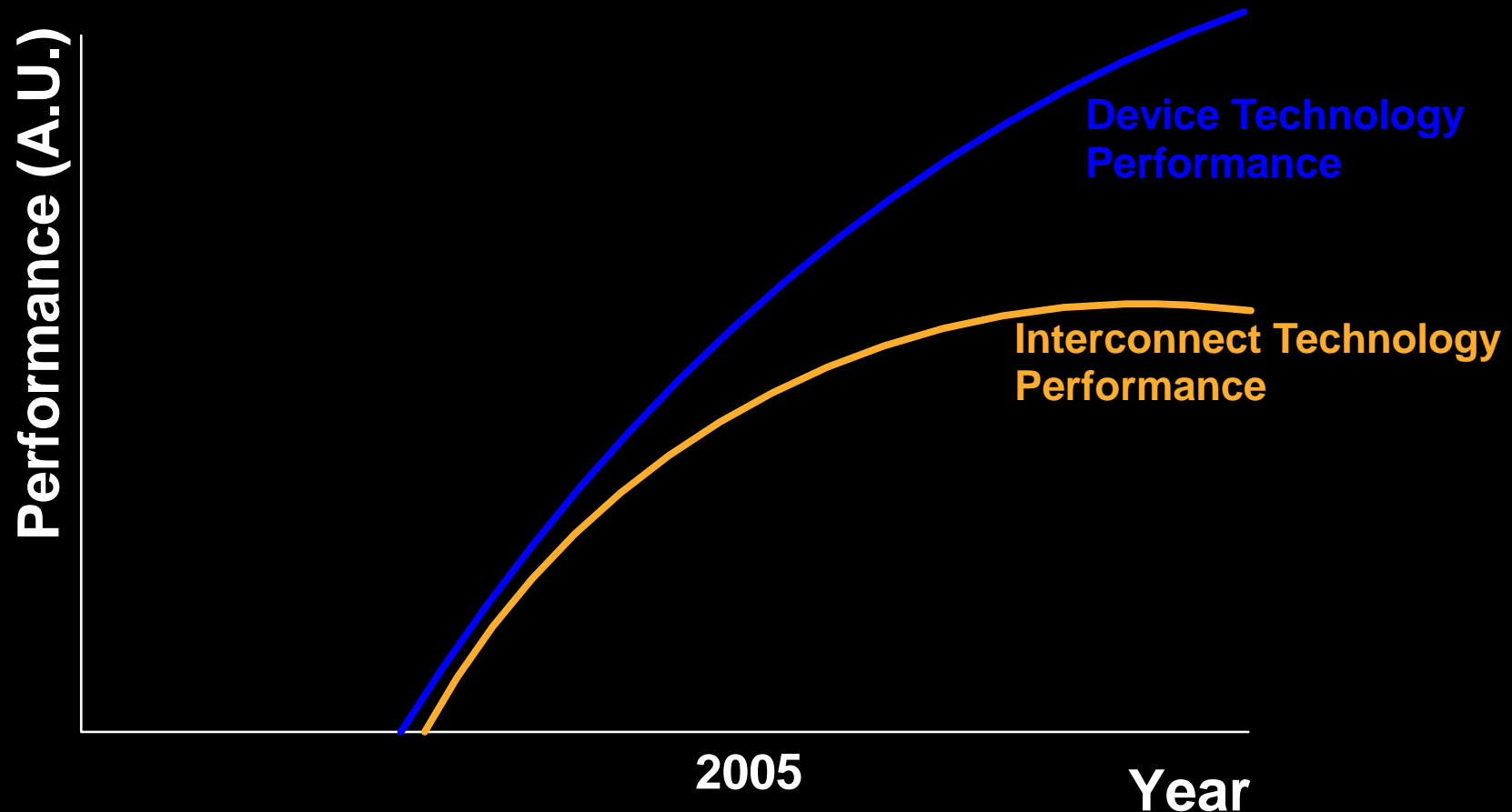1995    2000    2005    2010    2015

Year

- Perfect Storm: (a)Wire non-scaling; (b) die size growth; (c)Shorter FO4 stages
- Power Cost of Cross-Chip Latency increase

# First, a look at a "coincidence"……

Device performance (i.e. $I/C_G$) continues to improve, however at a decreasing rate…….



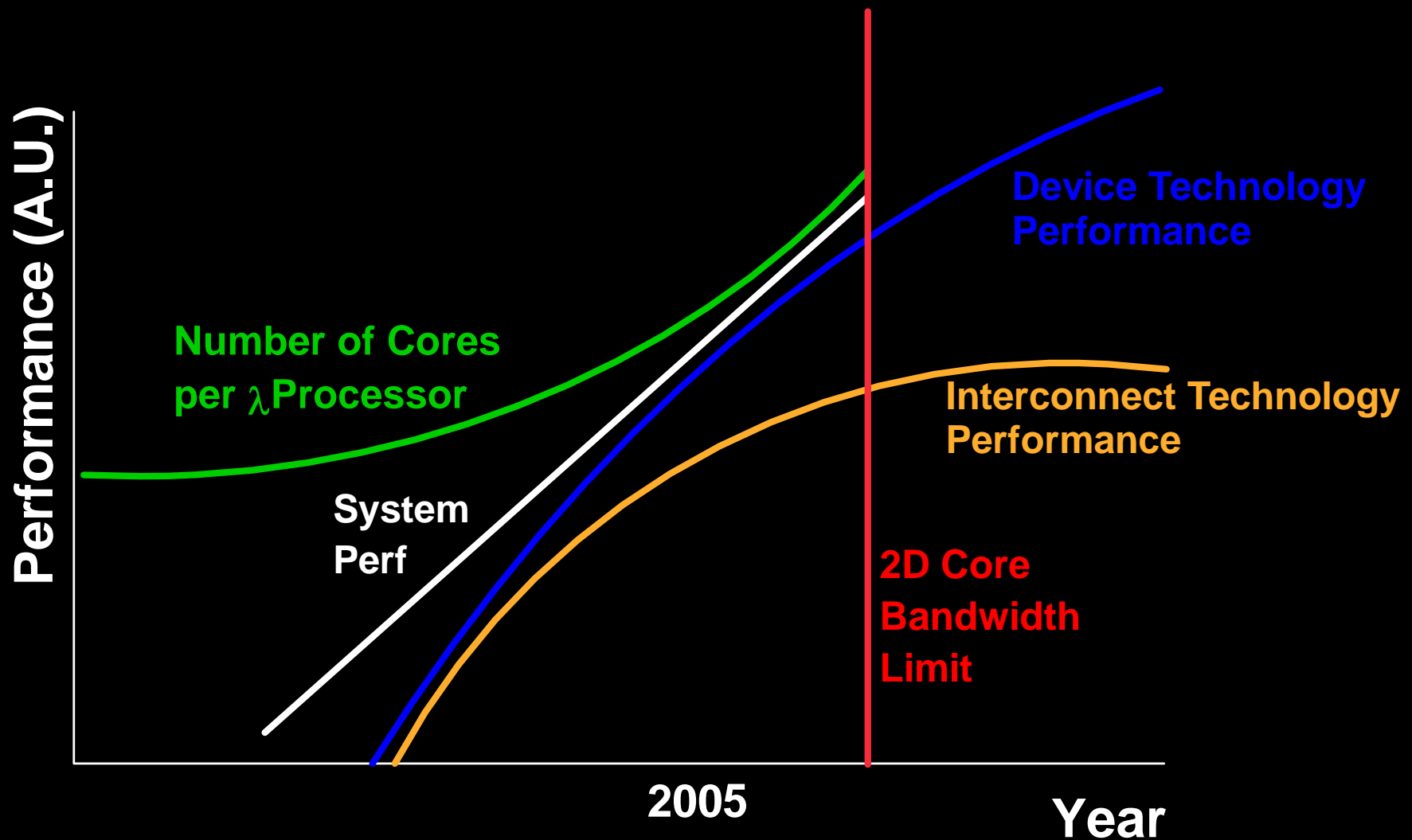Device Technology Performance

Performance (A.U.)

2005

Year

Despite constant infusion of new materials and processes however, interconnect technology performance has at best remained flat.
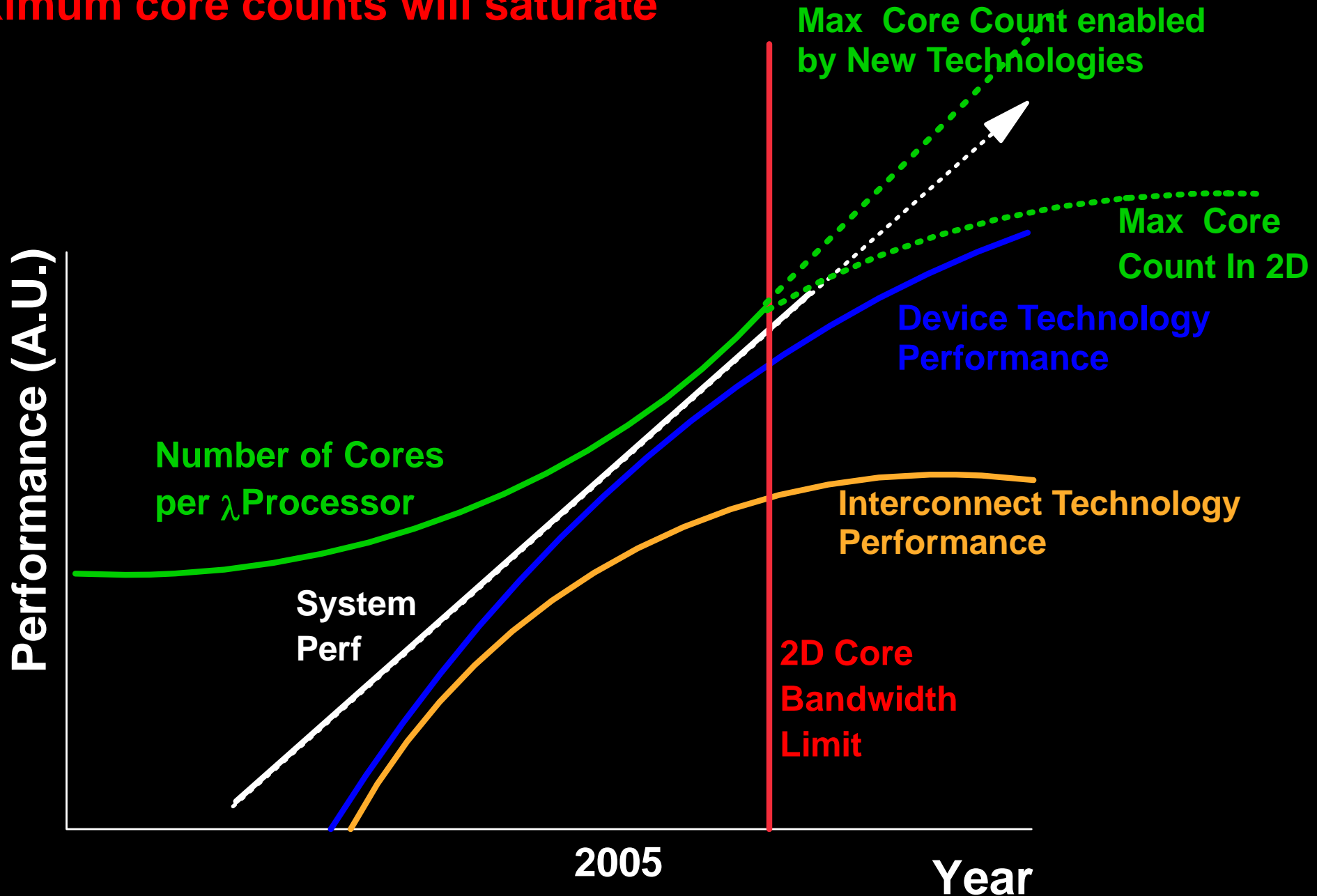


**Scaling has increased the divergence between FEOL and BEOL contributions to performance improvement**

As effective distances on chip increased due to interconnect, cores/chip has begun to climb. The bandwidth needed to feed these cores will ultimately **limit** number of cores



**Performance (A.U.)**

**Device Technology Performance**

**Number of Cores per λ Processor**

**Interconnect Technology Performance**

**System Perf**

**2D Core Bandwidth Limit**

**2005**

**Year**

**System Performance improvement is sustained more by the number of cores rather than by the performance of each core**

IBM

**Without more bandwidth at low latencies, maximum core counts will saturate**



Max Core Count enabled by New Technologies

Max Core Count In 2D

Device Technology Performance

Number of Cores per $\lambda$Processor

Interconnect Technology Performance

System Perf

2D Core Bandwidth Limit

Performance (A.U.)

2005

Year

**3D extends transfer of performance from the device to the core level**

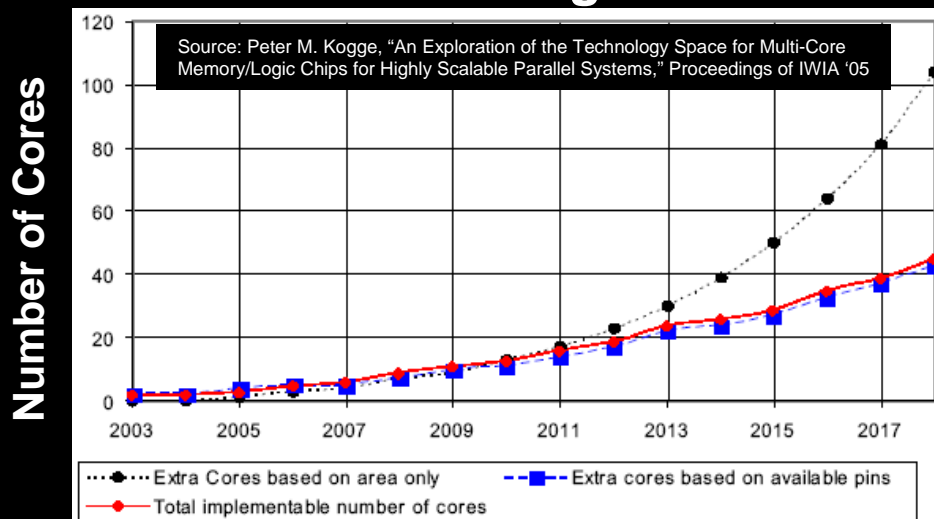# The Memory Wall, Bandwidth, and Latency

# Getting over the Memory Wall
## Microprocessor Architectures Fundamental Bus Limits

### Latency Challenge



### Bandwidth Challenge

**Number of Cores**



Source: Peter M. Kogge, "An Exploration of the Technology Space for Multi-Core Memory/Logic Chips for Highly Scalable Parallel Systems," Proceedings of IWIA '05

Extra Cores based on area only

Extra cores based on available pins

Total implementable number of cores

- **Processor speed has increased much quicker than memory access**

- **Result: $\lambda$P's data appetite has grown quicker than ability to feed it.**

  - What needs higher BW?

    - Multi-cores with limited cache

    - Multi-threading

    - Virtualization

- **Increasing "cores per chip" addresses memory latency. Core count Limit after 2010 will be from pins used to provide memory bandwidth**

  - The "Memory Wall" is back with a vengeance

# Architecture
## Cache Miss Penalty Calculation

L2 Memory

L1 Memory

λP

4. Cache Access
*"Your number please"*

3. L2 Cache Prioritization
*"Wait your turn"*

2. Address Time of Flight Up
*"L2, here's the address I need"*

1. Directory search + logic:
*"oops, a miss"*

5. Data Time of Flight Back
*"L1, here's your order"*

6. Error Correction Code
*"Are you sure?"*

7. L1 Cache Prioritization
*"Come on in"*

**Memory Latency is the delay encountered completing the loop above**

Bernstein

# What Is Bandwidth Used For?

**In a computer, it is mostly for handling cache misses:1**

**Miss**      **Access**      Processor Events

**Leading Edge**

Bus Events

**Rest of Cache Line**

**Trailing Edge**

Time

**First Data**      **Last Data**

**Miss Penalty = Leading Edge + Effects(Trailing Edge)**

**Where**

**Trailing Edge Effect** = (Line Size / Bus Width) x ($F_{(\lambda P)}$ / $F_{(Bus)}$)

**Bus Utilization** = (Trailing Edge / Intermiss Distance)

# Intermiss Distance Density

% of Misses

Intermiss Distance (# Instructions)

Queueing Effects vs. Log Miss Rate

# Server Trends are hard on Bandwidth

- **Frequency is no longer increasing**
  - Logic speed scaled faster than memory bus
  - (Processor clocks / Bus clock) consumes bandwidth

- **More speculation multipliers prefetch attempts**
  - Wrong guesses increase miss traffic

- **Reducing line length is limited by directory as cache grows**
  - But doubling line size doubles bus occupancy

- **Cores / die increasing each generation**
  - Multiplies off-chip bus transactions by N / 2*Sqrt(2)

- **More threads per core, and increase in virtualization**
  - Multiplies off-chip bus transactions by N

- **Total number of Processors / SMP increasing**
  - Aggravates queuing throughout the system

# 3D - Bandwidth and Latency

**Bandwidth and Latency Boundaries**
**General Purpose Processor Loads**

**Processor load trade-off between I/O Bandwidth, Bus Latency.**

**- For generic workloads, uni-processor perf saturates bandwidth benefit, becomes latency-limited.**

**- As core counts increase, I/O Bandwidth becomes increasingly important**



◆ Single Core
▲ Double Core
▲ Quad Core

Band width limited

Latency limited

Arch Perf, TpCC (Norm)

Bandwidth, GB/Sec (Norm)

**3D opportunity for improving High Perf Compute thru-put in sustaining a higher number of cores per chip**

# 3D Solution
## Hierarchical Memory Access



**2-D: Connections on the periphery**
- Long global connections
- CPU to off-chip main memory with latency and misses

**3-D: Connections across the area**
- Connections short + vertical
- Suitable for high-bandwidth and vector operations
- No pin cost, large block access of data

**Latency:** Important for random access (servers, e.g.), single core
**Bandwidth:** Multiple cores, multi-threads, graphics

S. Tiwari; "Potential, Characteristics, and Issues of 3D SOI; 3D SOI Opportunities"
Short Course, 2005 International SOI Conference

# The Technologies of 3D Integration
## (and their challenges)

# The 3D Integration Technology Spectrum

**Stacked Com.DRAM**

**Simple Chip Stack**

**3D Chip Support**

**Integrated eDRAM**

**Hierarchical Cache**

**3D Multicore uProc**

*Applications*

**Chip Stacks**

**3DI Integration**

| | Chip Stacks | | | | 3DI Integration | |
|---|---|---|---|---|---|---|
| **Via Density (pins/cm2)** | 1E2 | 1E3 | | 1E4 | 1E5 | 1E6 |
| **Via Size (um)** | 200 | 50 | | 10 | 1 | 0.100 |
| **Year (approx)** | 1990 | 1995 | | 2000 | 2010 | 2020 |
| **Supported Freq (Hz)** | 1E6 | | 1E7 | | 1E8 | 1E9 |

# Precedent for 3D Integration:
# When Real Estate Becomes Pricey



1900

1930

1970

2005

**NYC Office Inventory, Rent, and Skyscrapers**

Vertical Integration isn't new!

**Midtown Vacancy & Asking Rents**

**Manhattan Office Space**

Total SF
Vacant SF
Asking Rent

**Data courtesy of Richard Persichetti Grubb & Ellis, New York, NY**

# Chip-Package Technology Gap



- **Technology gap in the design rule between on-chip wiring and packaging interconnects**

# 2. Present Vertical Interconnect Schemes



Images used by permission, W.R. Davis, North Carolina State Univ,

Wire Bonding

Microbump

Coupled Virtual Connections

(a) Bulk

(b) SOI

Through-Via

# Evolution of 3D Integration

## Technology Investment in the Z-Dimension

- **3D Technologies continue the sequence of interconnect advances**

- **Return balance to device scaling**

- **Enable new capabilities not available in 2D**

*Increasing 3D Integration*

**CMOS 3D**

| Analog |
|---|
| Flash |
| DRAM |
| DRAM |
| CPU |

**3D Through Via Chip Stack**

Die 7
Die 6
Die 5
Die 4
Die 3
Die 2
Die 1

Pkg. Substrate
Metal Pad

**3D Pkg Chip Stack**

12345 12345-67
000 STAKTEK

**3D Flip Chip Package Stack**

- **3D Packaging R&D now pervasive in industry, academia**

- **Through-via technology emerging as predominant path**

- **3D has *always* been large volume, but now integrating higher technologies**

**Wire bonded chip stacked 3D**

# Key 3DI Processes

Images courtesy of Anna Topol,
IBM T.J. Watson Research Center

### Bonding

### Transfer/Alignment



### Electrical Contacting

### Release Process

# IBM 3D Process:
# SOI-Based 3DI Layer Transfer

- **Device layers stacked using wafer bonding**

- **Each layer fabricated by conventional processes**

- **Layers fabricated and tested simultaneously**

GLASS

SOI

BOX

Circuit Layer 1

GLASS

Circuit Layer 2

3D IC

- **Attach circuit to glass handle wafer**

- **Remove original substrate**

- **Align & bond top circuit to bottom circuit**

- **Remove handle wafer & adhesives**

- **Form vertical interconnects**

# Wafer Transfer / Thinning

K. Guarini, IEDM, 2002.

## Transparent Circuit
### 200 mm Wafer
### 130nm SOI Technology





**GLASS**



- SOI device layer + back-end metallization transferred onto glass

- Defect-free lamination over 200 mm wafers

# 3D Fly-Thru Movies of  IBM Assembly

# 3D Challenges
## Heat Dissipation and Natural Selection

*Why is area vs volume such a big deal?*

# Power/Energy Issues

- **It now takes more energy to move data than to generate it, even just across chip**

  – Compute: 50pJ / FLOP / bit

  – Read: 10 pJ / operand from Reg….but
  1 nJ / operand from cache

- **Worst power nets on chip are data, instruction nets: go from mm(2D) to λm(3D)**

Heatsink

Core

SRAM

DRAM

Heatsink

DRAM

SRAM

Core

Delta

Compliments of Sri Sri-Jayantha, IBM Research

Slice: Temperature [°C]   Arrow: Total heat flux [W/m²]   Streamline: Total heat flux [W/m²]

Scenario 1

Scenario 2

Slice: Temperature [°C]

# C. EDA and 3D Integration Trends



Interlayer Via Density (number of vias/mm²) — left axis: 100, 1000, 10000, 100000

Area efficiency ($F_0/(F_0+contact)$) — right axis: 12%, 15%, 18%, 21%

Sweet Spot

Chip Performance is limited by global paths at core/unit level. For significant performance improvement, 3D integration at core or unit level is desirable.

Chip labels: Tree, FPU1, PU0, L2, PU1, Torus, FPU0, Eth, JTAG, Perf, L3

d_rom_kmac

3D Design partitioning Level

| Transistor Level 0.25u via | Macro Level 0.75u via | Unit Level 2.0u via | Core Level 4.0u via |

# 7. Summary

- $\lambda$P architecture tricks to avoid atomistic, QM scaling boundaries overwhelm present interconnects

- Integration into Z-plane again postpones interconnect-related limitations to extending classic scaling.

- Transaction retirement rate dependence on data delivery is *increasing*: dependence on $\lambda$P performance and CMOS device speed is *decreasing*

- 3D Integration improves storage density & access to that storage

- 3D Integration will enable previously unattainable capabilities characterized by realtime access to massive amounts of storage.

# *Human* Scaling

Tomorrow's microprocessors will be improved with capabilities developed using today's machines

Tomorrow's engineers will design microprocessors with insights they learn from today's engineers and professors.

Engineers/professors today insure a bright tomorrow by transferring **ideas** as well as **technologies** to the next generation.