



Outline

- Power and energy consumption basics
- Power consumption in processors
- Multicore: power and utilization walls
- Energy advantages of hardware accelerators
- Playing with accuracy for reducing energy

3

- Towards heterogeneous manycores
 - 3D stacking
 - Optical interconnect



















Minimum Energy per Operation

• Putting all together 4.5 4 3.5 **Energy per Operation [bJ]** - Dynami -Leakage 0.5 Total 0 0.3 0.2 0.5 0.7 0.8 0.4 Supply Voltage [V] 13

On-Chip Interconnect?

- Gate delay decreases but... wire delay increases
- Crossing chip in 5-10 clock cycles
- Also affected by noise...



- Metal layers to reduce wire delay
- Repeaters
- Towards networkon-chip

14

Conclusion: Power in CMOS

$\mathbf{P} = \sum_{\mathbf{i}} \left[\alpha_{\mathbf{i}}.\mathbf{f_{i}}.\mathbf{C_{i}}.\mathbf{Vdd^{2}} + \mathbf{I_{leak_{i}}}.\mathbf{Vdd} \right]$

- Dynamic power
 - 40-70% today
 - Decreasing relatively
 - DVFS becomes more and more difficult
- Leakage power
 - 20-50 % today
 - Increasing rapidly
 - number of transistors
 - Vdd/Vt scaling
 - Critical for memory

 $P = \frac{energy}{operation} \times rate + static \ power$

15

<section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item>

Power Consumption in Processors

• A typical (yet simple) processor pipeline





Energy Cost in a Processor

• Fetching operands costs more than computing





The Energy Cost of Data Movement

- Future processor up to 3 Tera-op/sec
- At minimum requires 64b x 9 Tera-operands to be moved each second
- If on average 1mm (10% of die size) then





Dynamic Power Management

- Dynamic Voltage and Frequency Scaling (DVFS)
- Reduce speed (clock freq.) and Vdd depending on processor activity



<section-header> Outline Power and energy consumption basics Power consumption in processors Multicore: power and utilization walls Muscalerations Playing with accuracy for reducing energy Powards heterogeneous manycores Can 3D stacking help? Optical interconnect





Shared-Memory Multiprocessor



- Processors communicate with shared address space by memory read/write
- Hardware-managed, implicitly-addressed, coherent caches
- Bandwidth depends on
 - Cache size, associativity
 - Replacement policy, coherence protocol
 - Application requirements

27



IBM Power 8

- Across 12 core chip – 4 TB/sec L2 BW
 - 3 TB/sec L3 BW
- 230 GB/s sustained external memory bandwidth

GB/sec shown assuming 4 GHz











