

Energy Efficient Big Data Processing at the Software Level

Da-Qi Ren, Zane Wei

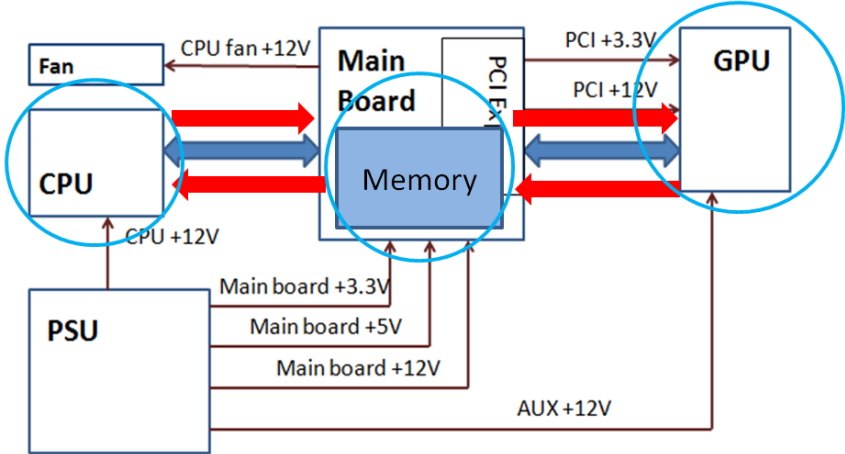
*Huawei US R&D Center
Santa Clara, CA 95050*



Power Measurement on Big Data Systems

1. If the System Under Test (SUT) is physically standalone, ----- measure from the input of the rack/cabinet/case.
2. If devices are separately powered --- separately measured.
3. The total power ----- summation of each component.

$$P = \sum_{1 \leq i \leq m} p_i$$



Rack 01

	Huawei S5700 Gigabit switches
	00 Master Node
...	01 Slave Nodes
	02 DN+TT
	03 DN+TT
...	04 DN+TT
	05 DN+TT
...	06 DN+TT
	07 DN+TT
...	08 DN+TT
	09 DN+TT
	10 DN+TT

RH2288 V2 Server
 CPU: E5-2680 2.7GHZ x 2
 Disk: 300GB 10Krpm SAS Interface x 4
 Memory: DDR3 16GB 1333Hz x 16

Component level power need to be measured

Energy Metrics

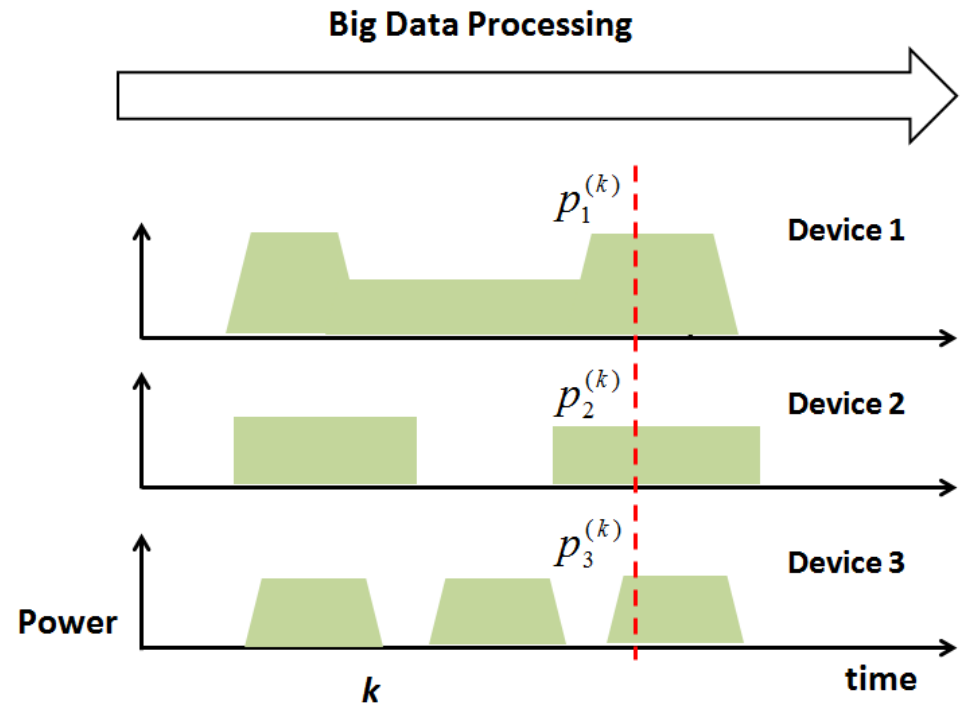
Energy Calculation Based on Measurement Results

For each device or subsystem, the calculation defined for Energy metrics is:

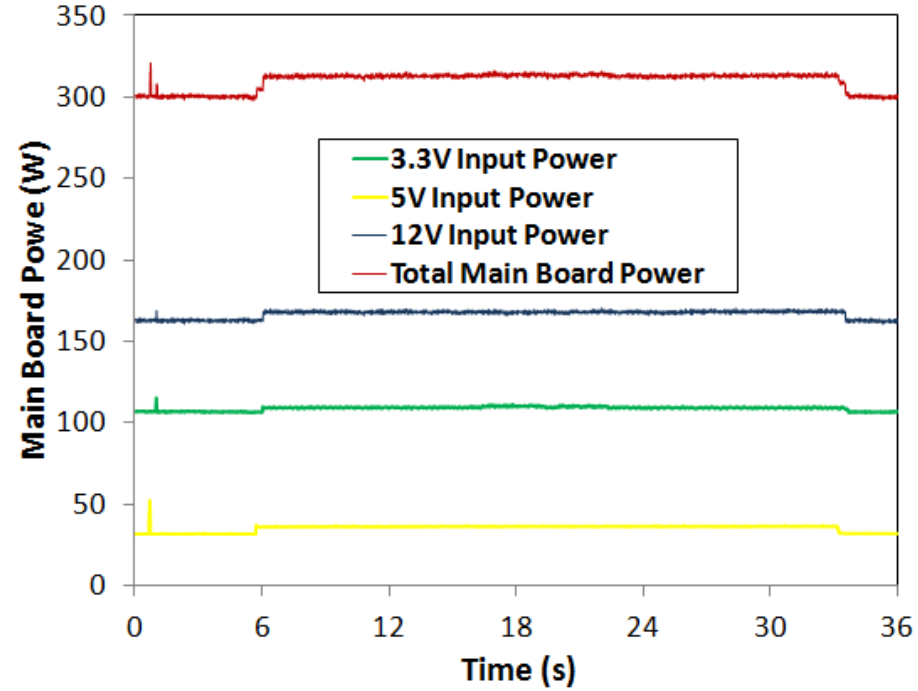
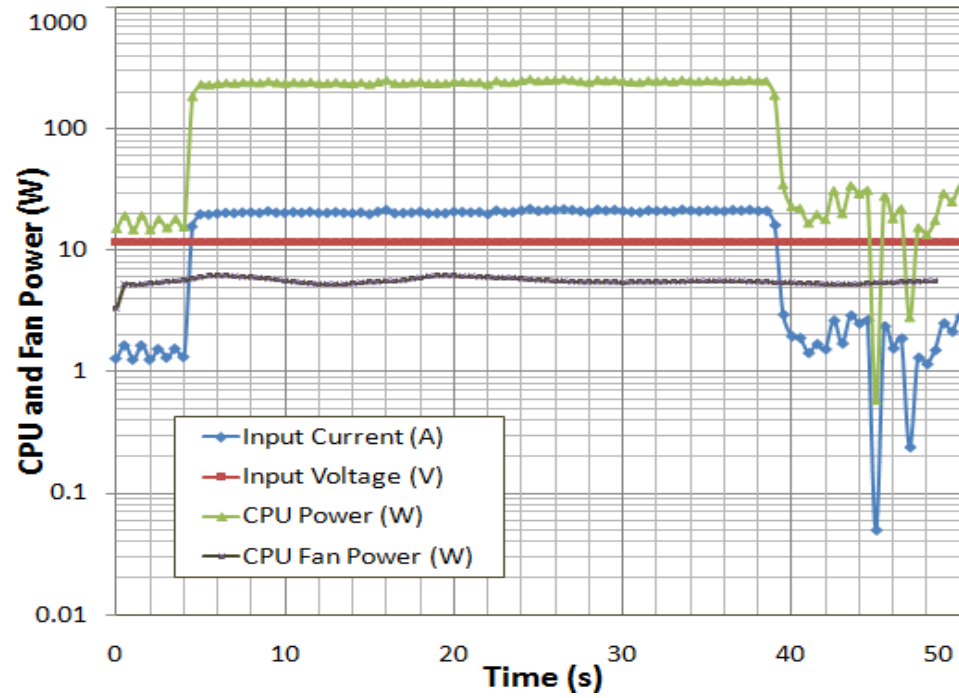
$$E = \int_0^T P(t)dt \quad E = \sum_0^T P^{(k)}k$$

Equation (1), *Energy* E is the elapsed time t times power P .

When sampling interval is small enough, the energy calculation of equation (1) becomes the format in equation (2). All real measurements can be done follow it by using power analyzers.



Sample Power Chart



CPU-GPU Processing Element Power Model

- 1. CPU power Measurement.** From CPU socket on main board ---- to measure the CPU input current and voltage at the 8-pin power plug. (Most of the onboard CPUs are powered only by this type of connector)
- 2. GPU power measurement.**
- 3. Memory and main board power estimation.** we can make an approximation on its power by measuring the power change on the main board.

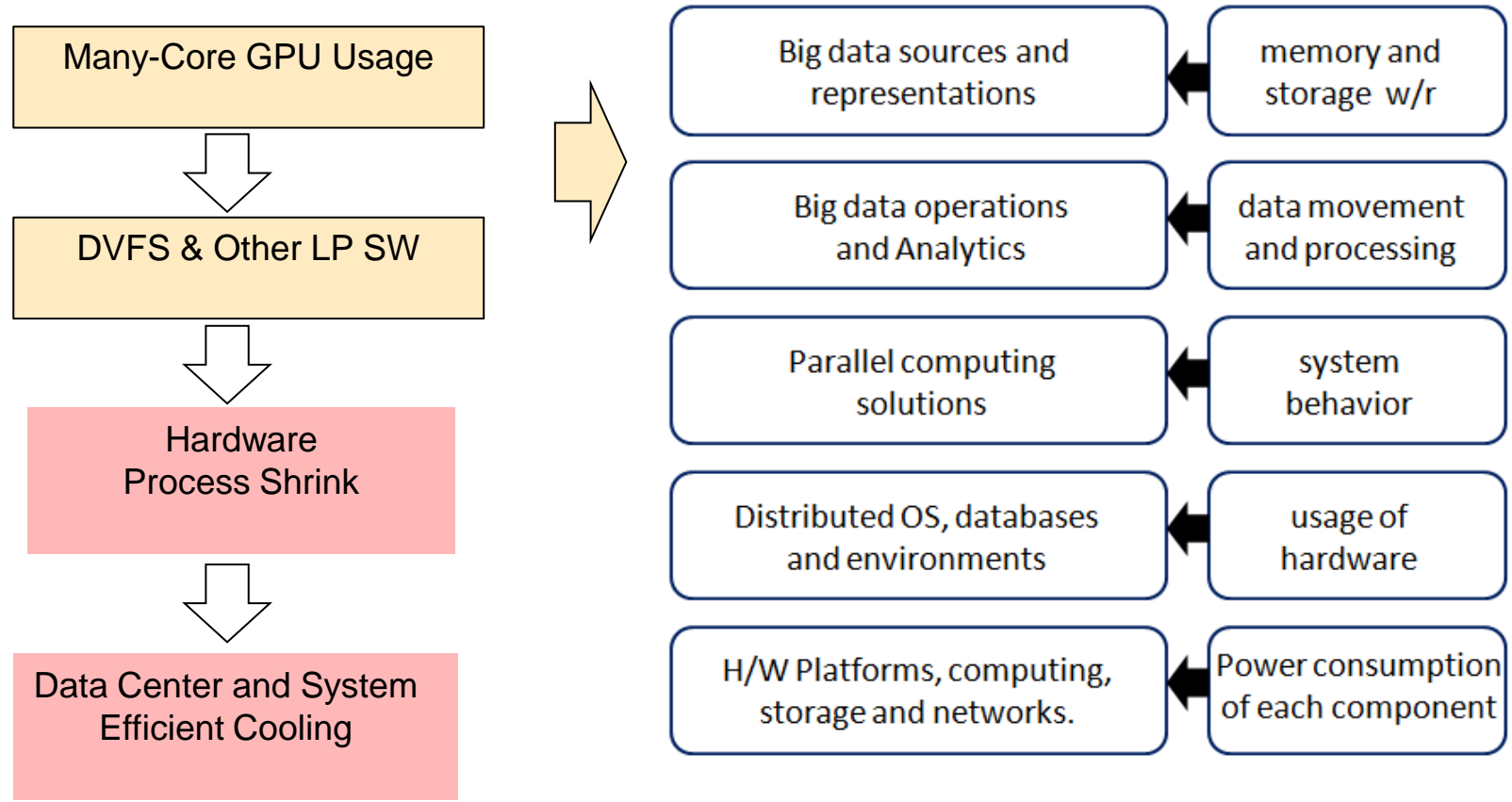
$$P_{total}(w) = \sum_{i=1}^N P_{GPU}^i(w^i) + \sum_j^M P_{CPU}(w^j) + P_{mainboard}(w)$$

CPU-GPU PE Power Feature Determination

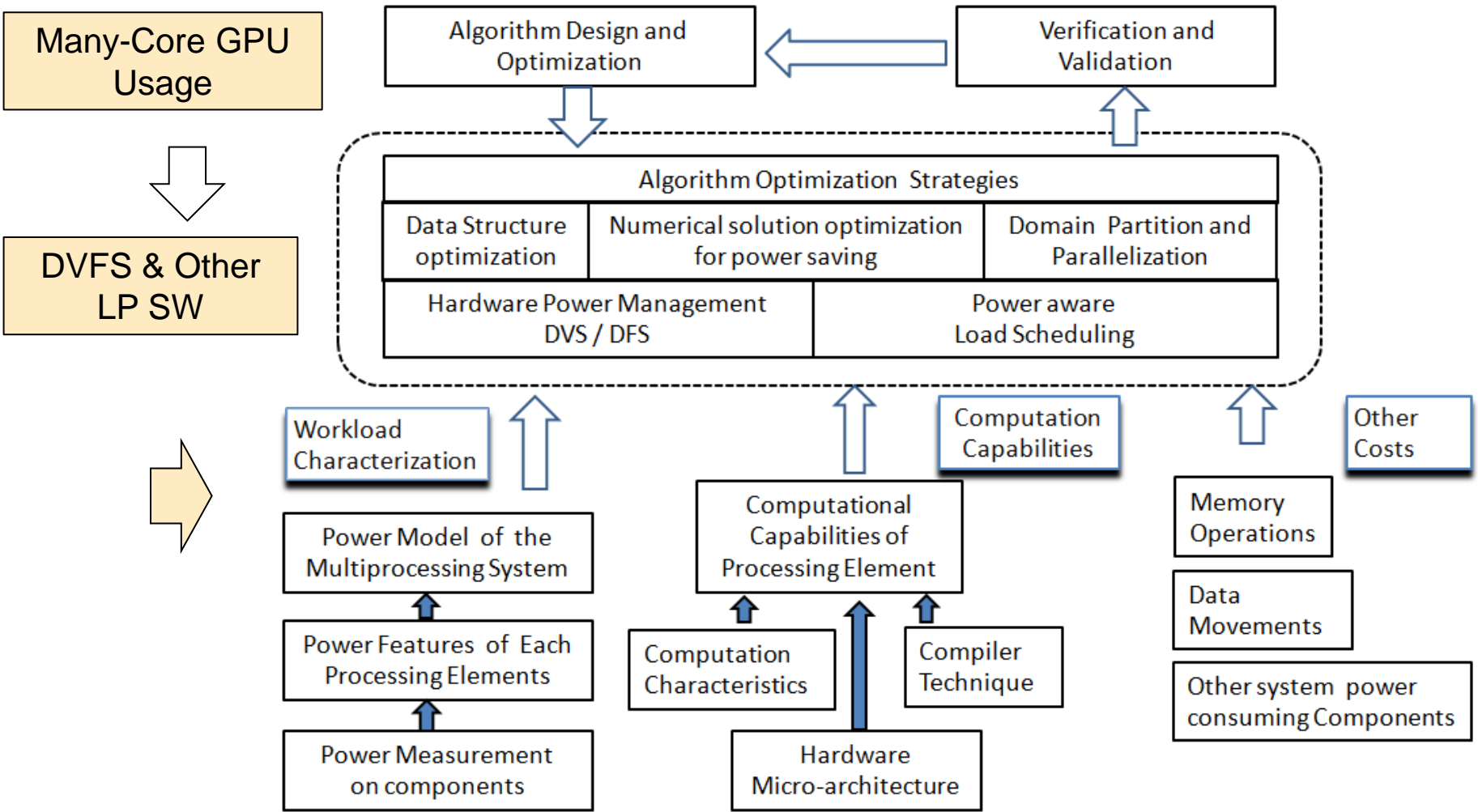
1. Measuring the power from each component of the PE;
2. Find out the energy consumption features
 - ❖ **FLOPS / Watt**
 - ❖ **Data Size / Watt**
 - ❖ **Workload / Watt**
 -to this application;
3. Estimated execution time --- the total workload FLOP to be computed divides by the computational speed ----- the CPU-GPU processing element can support;
4. Estimated energy consumption to complete the program ----- the summation of products of the component powers and the execution times.

Power issues in each layer of Big Data Computing

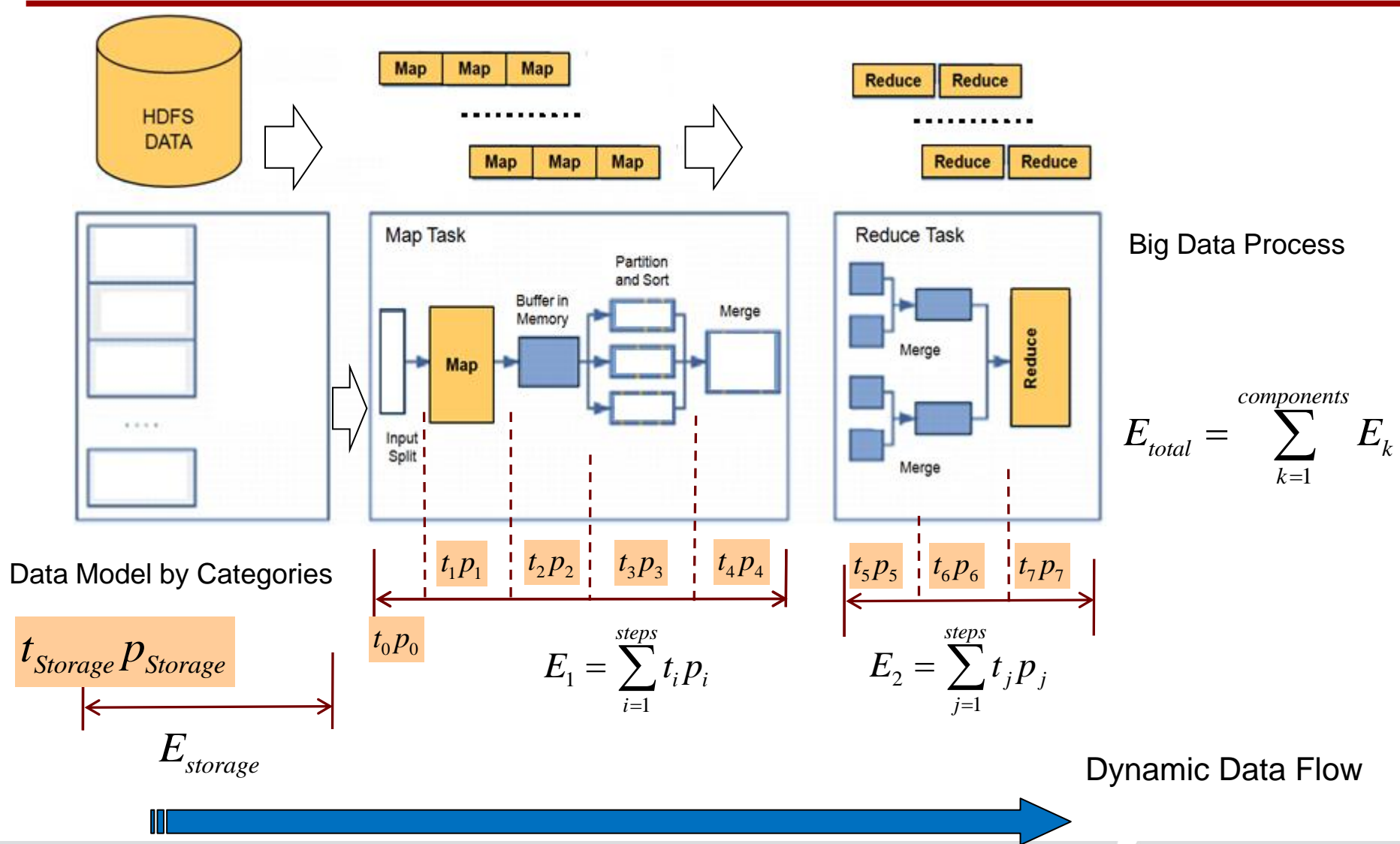
Power Effective Factors, when original power performance is x



Software Solutions for Saving the Energy



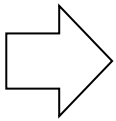
Power Analysis in Big Data Computing



$$E_{total} = \sum_{k=1}^{components} E_k$$

Power Efficient Implementation by Algorithm Optimization

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}, \mathbf{B} = \begin{bmatrix} \mathbf{B}_{1,1} & \mathbf{B}_{1,2} \\ \mathbf{B}_{2,1} & \mathbf{B}_{2,2} \end{bmatrix}, \mathbf{C} = \begin{bmatrix} \mathbf{C}_{1,1} & \mathbf{C}_{1,2} \\ \mathbf{C}_{2,1} & \mathbf{C}_{2,2} \end{bmatrix}$$



$$\begin{aligned} \mathbf{C}_{1,1} &= \mathbf{A}_{1,1}\mathbf{B}_{1,1} + \mathbf{A}_{1,2}\mathbf{B}_{2,1} \\ \mathbf{C}_{1,2} &= \mathbf{A}_{1,1}\mathbf{B}_{1,2} + \mathbf{A}_{1,2}\mathbf{B}_{2,2} \\ \mathbf{C}_{2,1} &= \mathbf{A}_{2,1}\mathbf{B}_{1,1} + \mathbf{A}_{2,2}\mathbf{B}_{2,1} \\ \mathbf{C}_{2,2} &= \mathbf{A}_{2,1}\mathbf{B}_{1,2} + \mathbf{A}_{2,2}\mathbf{B}_{2,2} \end{aligned}$$

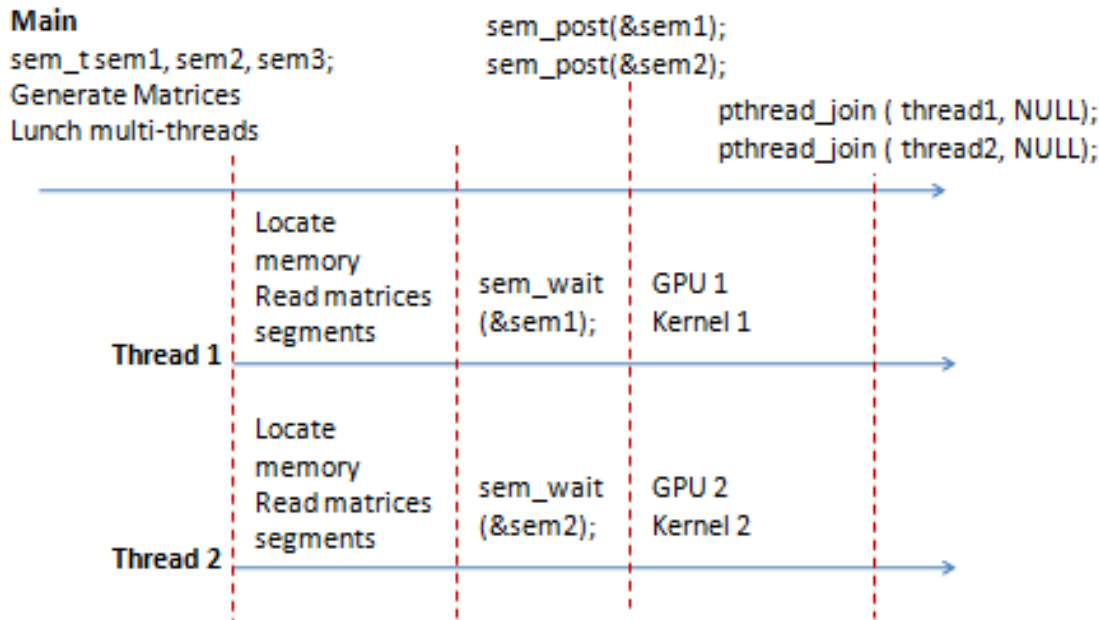
	Simple Algorithm	Enhance Algorithm
CPU Energy Consumption	0.24 wh	0.019 wh
GPU Energy Consumption	0.654 wh	0.055 wh
Main Board (Main Memory) Energy Consumption	0.47 wh	0.0402 wh
Overall Time Consumption	7.5 s	0.6s
Overall Energy Consumption	1.364 wh	0.1142 wh

Build up a power model incorporates physical power constrains of hardware;

Partition of smaller matrix-blocks to fill the shared memories of GPU.

Reduce the data transmission between GPU and main memory, will significantly enhance the GPU performance and power efficiency.

Parallel GPUs



Parallel GPU approach with signal synchronization.

Multithread kernel control to save the use of CPU cores.

Partition jobs and distribute them to each GPU device;

Design synchronization signal to synchronize each CUDA kernel.

Remove CUDA Overhead

- ❑ A GPU/CUDA overhead is on kernel initialization, memory copy and kernel launch, before start real kernel computation.
- ❑ Remove CUDA overhead by calling C function to compute small workload on CPU.
- ❑ Save the time and energy of CUDA kernel and GPU.

A threshold can be determined by experiment by analyzing the follows:

$$T_{CPU}^k \leq T_{GPU}^k = T_{GPU\text{overhead}}^k + T_{CUDA\text{kernel}}^k$$

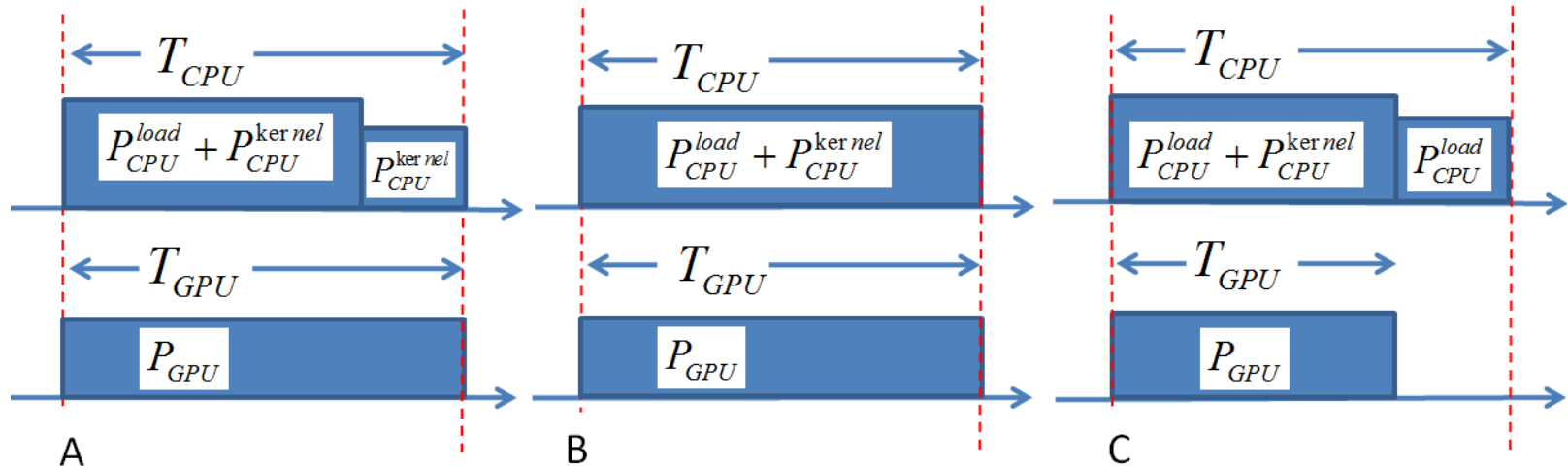
$$E_{CPU}^k = P_{CPU} \times T_{CPU}^k$$

$$E_{GPU}^k = P_{CPU-GPU-PE} \times T_{GPU}^k$$

C function will be selected when matrix size

less than k where $E_{CPU}^k \leq E_{GPU}^k$.

CPU Shares GPU's Workload



$$A) E_{total} = T_{total} (P_{GPU} + P_{CPU}^{kernel}) + T_{CPU}^{load} P_{CPU}^{load}$$

$$B) E_{total} = T_{total} (P_{GPU} + P_{CPU}^{load} + P_{CPU}^{kernel})$$

$$C) E_{total} = T_{GPU} (P_{GPU} + P_{CPU}^{load} + P_{CPU}^{kernel}) + (T_{total} - T_{GPU}) P_{CPU}^{load}$$

when $T_{GPU} > T_{CPU}$

when $T_{GPU} = T_{CPU}$

when $T_{GPU} < T_{CPU}$

CPU Frequency Scaling

- ❑ CPU frequency scaling can save CPU power, without decreasing the computation performance.
- ❑ CPU frequency should match the requirement of GPU/CUDA kernel calls.
- ❑ CPU frequency can be scaled down without compromising with the PE's performance however to save the CPU's power.

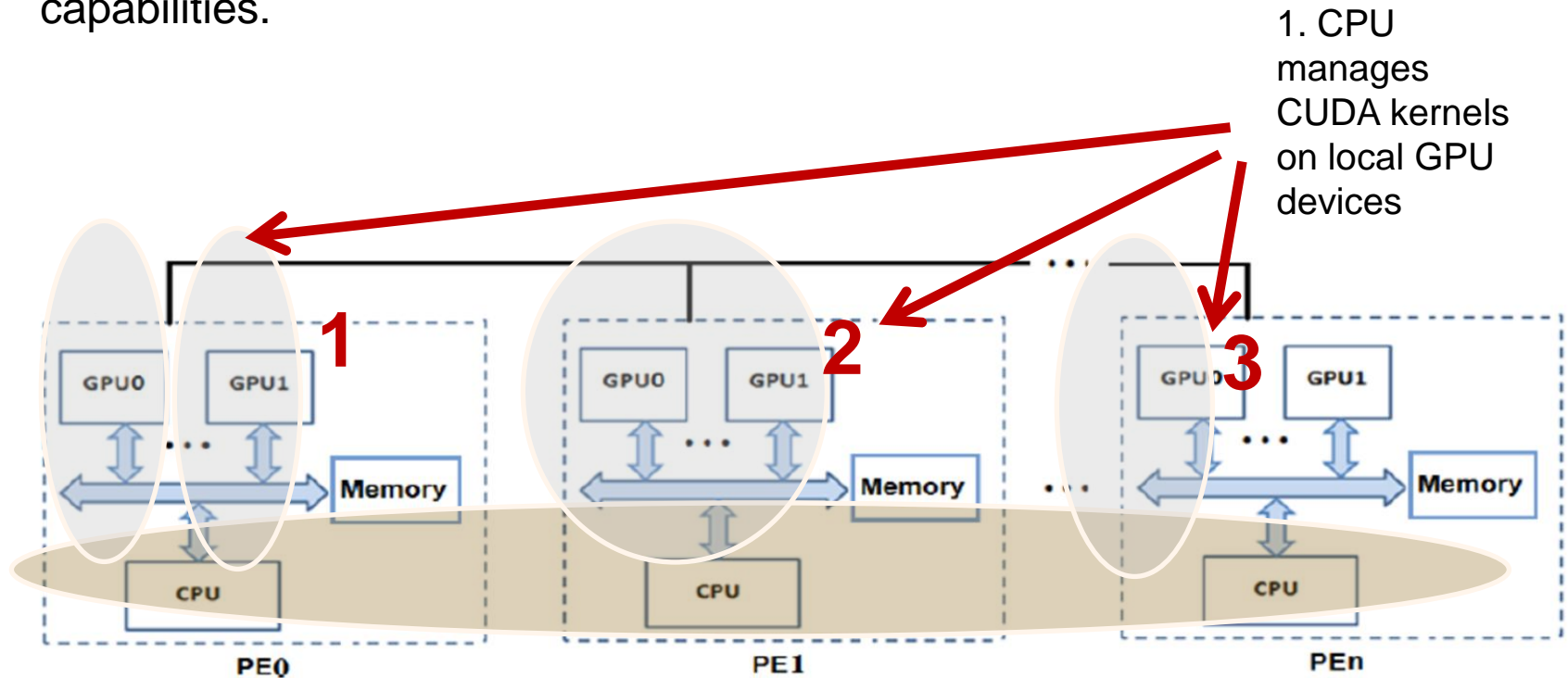
A rough estimation for the minimum CPU frequency value is :

$$F_{CPU} \geq F_{CudaCalls} \text{ (most of the cases)}$$

$$F_{CPU} \geq F_{GPUMemory} \text{ (if } F_{CudaCalls} \geq F_{GPUMemory} \text{)}$$

Energy Aware Load Scheduling with CUDA/MPI

- ❑ SIMD/SPMD (Single Instruction/Program Multiple Data) split up tasks to run simultaneously
- ❑ CUDA Processing Element (PE)
- ❑ CUDA/MPI An important choice in data intensive HPC in various applications.
- ❑ Scheduling the workload based on the PE's power feature and computing capabilities.



Conclusion

- Introducing energy efficient software design methodologies for big data processing, power performance metrics and measurements
- Modeling and evaluating large scale computer architectures with multi-core and GPU
- Global optimization for choosing the best power-efficient alternative based on data characteristics and quantitative performance analysis
- Validation of the energy efficiency improvements from the new designed approach

Thank you

www.huawei.com