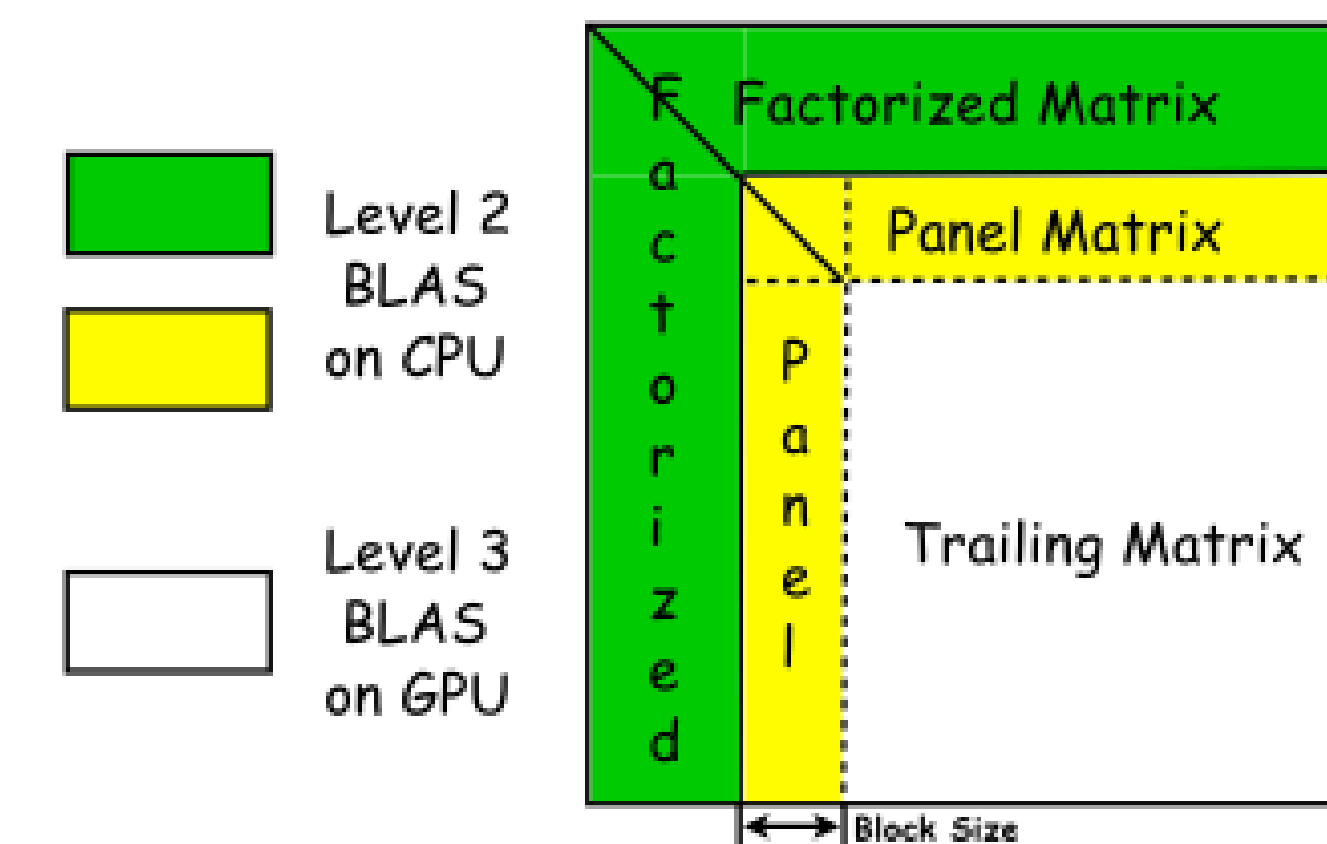# Optimizing Energy Efficiency for Distributed Dense Matrix Factorizations via Utilizing Algorithmic Characteristics

*Li Tan and Zizhong Chen, UC Riverside*
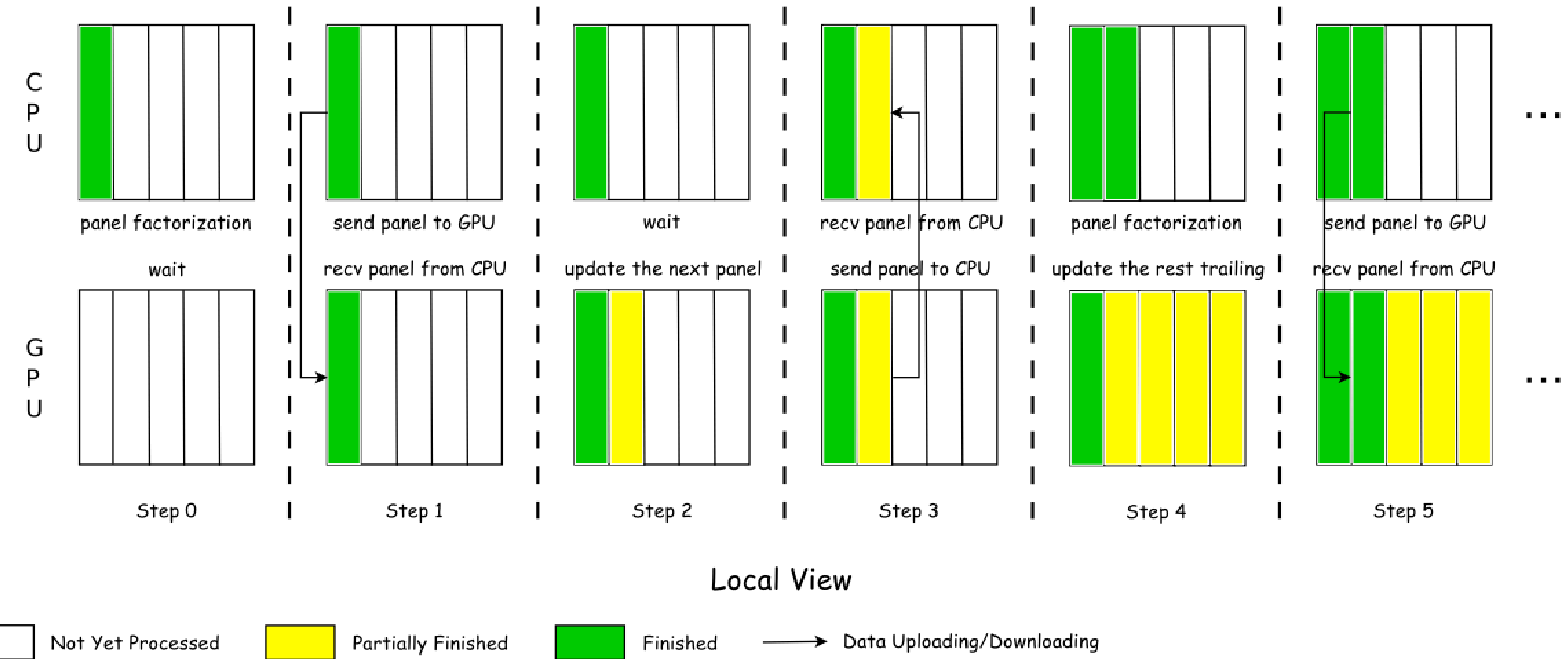
UNIVERSITY OF CALIFORNIA, RIVERSIDE

UC RIVERSIDE — UNIVERSITY OF CALIFORNIA

## Summary

<u>Goal</u>: Achieve the optimal energy efficiency for linear algebra operations by using algorithmic characteristics

- ❑ Pressing demand of improving energy efficiency for high performance scientific computing
  - o Scientific apps widely applied on supercomputers
  - o Costs of powering supercomputers are increasing

- ❑ Strategic DVFS energy efficient scheduling
  - o Switch ↓ hardware power if peak perf. not needed
  - o CPU, GPU, memory → Handy APIs for DVFS

- ❑ Numerical linear algebra algorithms: LU, QR, Chol.
  - o Dense Linear Algebra vs. Sparse Linear Algebra
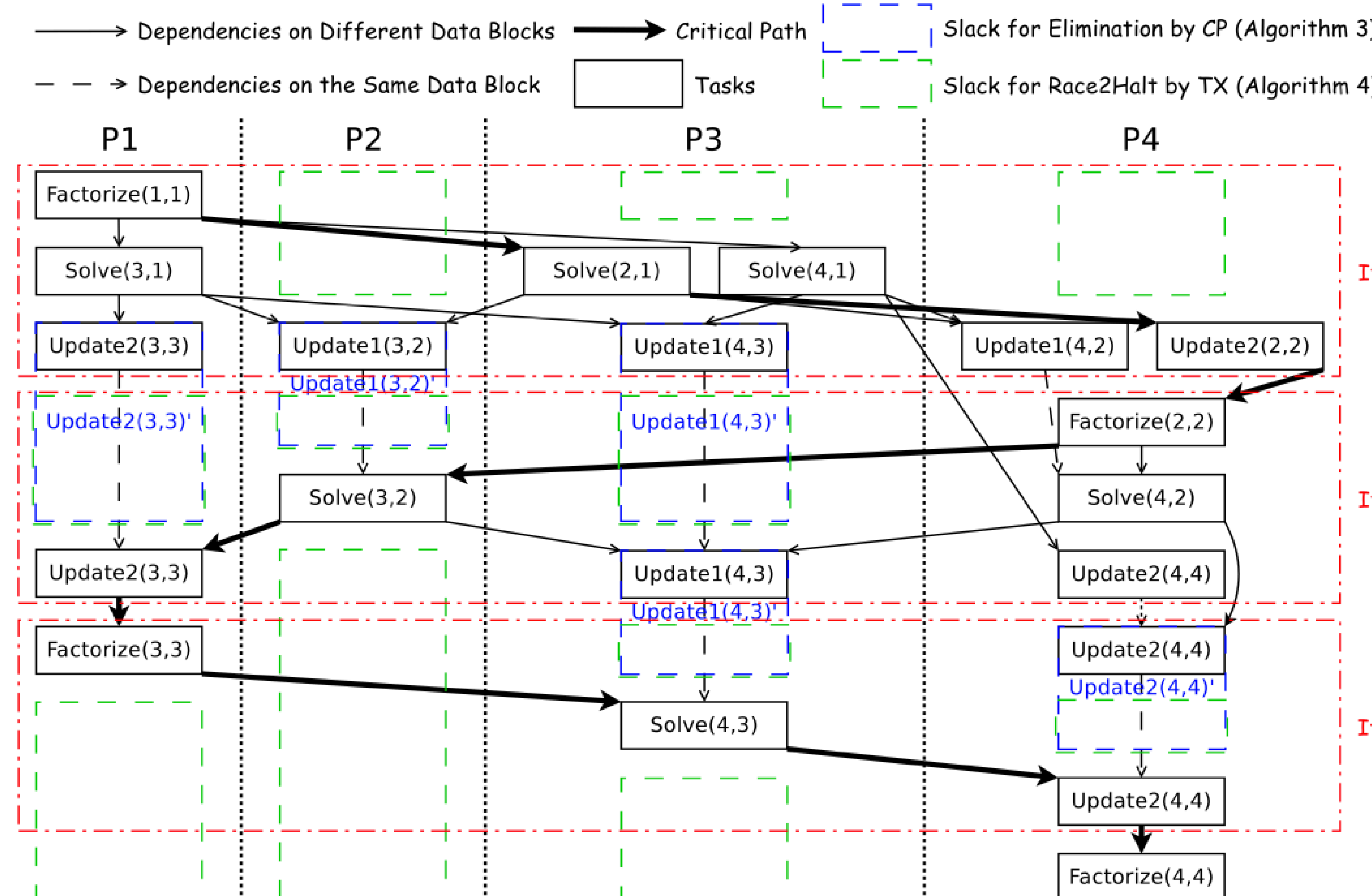  - o Homogeneous Systems vs. Heterogeneous Sys.

Level 2 BLAS on CPU · Level 3 BLAS on GPU

Factorized Matrix / Panel Matrix / Trailing Matrix / Block Size / Global View

C P U: panel factorization / send panel to GPU / wait / recv panel from CPU / panel factorization / send panel to GPU
wait / recv panel from CPU / update the next panel / send panel to CPU / update the rest trailing / recv panel from CPU

G P U

Step 0 / Step 1 / Step 2 / Step 3 / Step 4 / Step 5

Local View

Not Yet Processed · Partially Finished · Finished · Data Uploading/Downloading

## Background: Energy Saving Strategies

- ❑ Critical Path (CP) and Slack Analysis
  - o Slack: a period when a hardware waits for another
  - o Slack Examples: load imbalance, network latency, communication delay, memory and disk access stall
  - o Critical Path: a particular sequence of tasks where the total slack amounts to zero in task-parallel apps.

- ❑ CP-free Race-to-halt
  - o Enforce hardware to run at the highest F/V when workloads are ready and at the lowest F/V otherwise

- ❑ CP-aware Slack Reclamation
  - o Tasks on the CP: Run at highest F/V (peak perf.)
  - o Tasks off the CP: Run at appropriately scaled F/V

- ❑ Detailed Comparison
  - o Race-to-halt
    - ▪ Easy to implement
    - ▪ Hardware independent
    - ▪ Save additional energy due to load imbalance
  - o CP-aware
    - ▪ Require CP detection
    - ▪ Hardware dependent
    - ▪ Generally the optimal
  - o Hardware Utilizat. Polling
    - ▪ OS level monitoring
    - ▪ Coarse-grained
    - ▪ Runtime overhead

→ Dependencies on Different Data Blocks · ➡ Critical Path · Slack for Elimination by CP (Algorithm 3)
--→ Dependencies on the Same Data Block · Tasks · Slack for Race2Halt by TX (Algorithm 4)

P1 / P2 / P3 / P4

Factorize(1,1) / Solve(3,1) / Solve(2,1) / Solve(4,1) / Update2(3,3) / Update1(3,2) / Update1(4,3) / Update1(4,2) / Update2(2,2) — Iter1
Update1(3,2) / Update2(3,3)' / Update1(4,3)' / Factorize(2,2) / Solve(3,2) / Solve(4,2) / Update2(3,3) / Update1(4,3) / Update2(4,4) — Iter2
Update1(4,3)' / Factorize(3,3) / Update2(4,4) / Solve(4,3) / Update2(4,4)' / Update2(4,4) / Factorize(4,4) — Iter3

## Overall Comparison: Existing Approaches

### Operating Layer

- ❑ OS Level Approaches
  - o Work aside app. at runtime
  - o Make online decisions
  - o No need of app-specific knowledge
  - o General, no source modification

- ❑ Library Level Approaches
  - o Customize energy saving decisions
  - o Utilize app-specific knowledge
  - o Library source modif. & recompilation
  - o Optimize potential energy savings

### Prediction Mechanism

- ❑ Online Slack Prediction
  - o Online training of prediction model using execution info. in the same execution
  - o Dynamic training overhead
  - o Unexploited energy savings during training and prediction
  - o Accurate prediction requires sufficient training: overhead ↑

- ❑ Offline Slack Prediction
  - o Profile logged execut. history statically: no runtime overhead
  - o Impractical for running large-scale app. in HPC: too costly
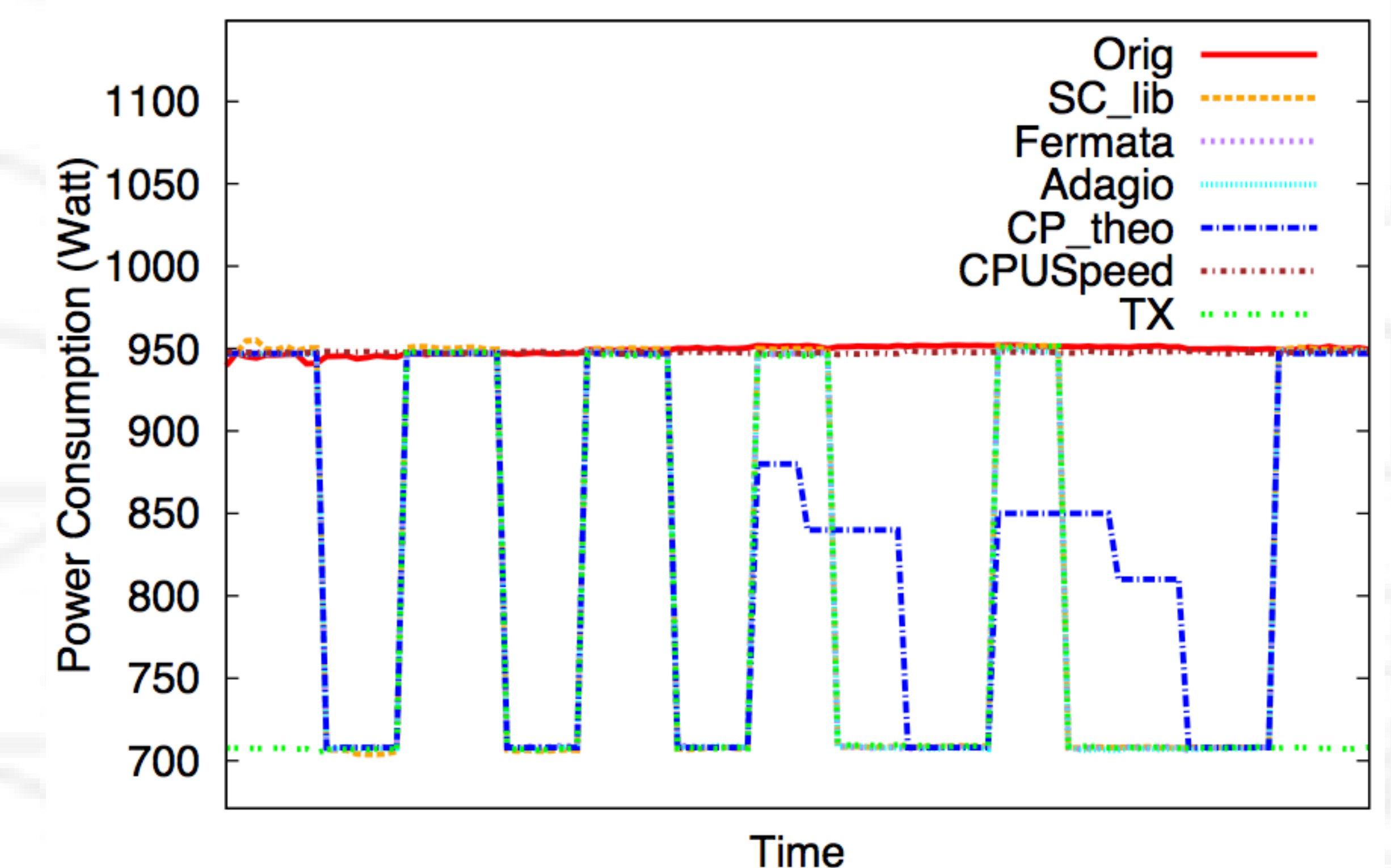
### Architecture

- ❑ Homogeneous Systems
  - o Slack arises among different processes through loop iterations

- ❑ Heterogeneous Systems
  - o Slack arises between CPU and accelerators (e.g., GPU)
  - o Similar techniques to HomoSys
  - o Sleep state scheduling for unused/idle CPU cores

## Algorithmic Energy Saving Strategy

- ❑ Utilize Algorithmic Characteristics of Target Applications
  - o Dense Cholesky, LU, and QR matrix factorizations
  - o Algorithmic Task Dependency Set (TDS) analysis
  - o Algorithmic online slack prediction

$$\frac{T_{i+1}^U}{T_i^U} = \frac{O(N_{i+1}'^3)}{O(N_i'^3)} = \frac{(N-(i+2) \times N_{proc} \times nb)^3}{(N-(i+1) \times N_{proc} \times nb)^3}$$

$$slack_{i+1} = \texttt{func}(T_i^U, T_{i+1}^U, param\_list)$$

## Preliminary Results: Power Savings

Matrix Size: 160000 x 160000, Power Costs of Three Nodes



Legend: Orig · SC_lib · Fermata · Adagio · CP_theo · CPUSpeed · TX

Power Consumption (Watt) vs Time

- ❑ Evaluated Power Saving Capability of Six Approaches
  - o Distributed Chol. factorization with 16 x 16 process grid
  - o OS level solutions failed to achieve the optimal savings
  - o Race-to-halt can be comparable to CP-aware solutions