# Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory

**Aleksandar Ilić** and Leonel Sousa

INESC-ID/
IST, TU Lisbon, Portugal

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology
from seed

# COMMODITY COMPUTERS = HETEROGENEOUS SYSTEMS

- Multi-core General-Purpose Processors (CPUs)
- Many-core Graphic Processing Units (GPUs)
- Special accelerators, co-processors, FPGAs,…

# => SIGNIFICANT COMPUTING POWER

- Not yet completely explored for **COLLABORATIVE COMPUTING**
- TO USE THE AVAILABLE RESOURCES AND IMPROVE PERFORMANCE/WATT

# HETEROGENEITY MAKES PROBLEMS MUCH MORE COMPLEX!

- Scheduling, performance modeling and load balancing
- Different programming models, languages and implementations

# Outline

**DISCRETELY DIVISIBLE LOAD (DDL) PROCESSING**

**HETEROGENEOUS (CPU+GPU) DESKTOP SYSTEMS**

**PERFORMANCE MODELING AND DDL SCHEDULING ALGORITHM**

**CASE STUDY: 2D FFT Batch Execution**

**CONCLUSIONS AND FUTURE WORK**

# Discretely Divisible Load Processing

- **DISCRETELY DIVISIBLE LOAD (DDL) APPLICATIONS**

  – Computations divisible into pieces of arbitrary sizes (integers)

  – Fractions independently processed in parallel with no precedence constraints

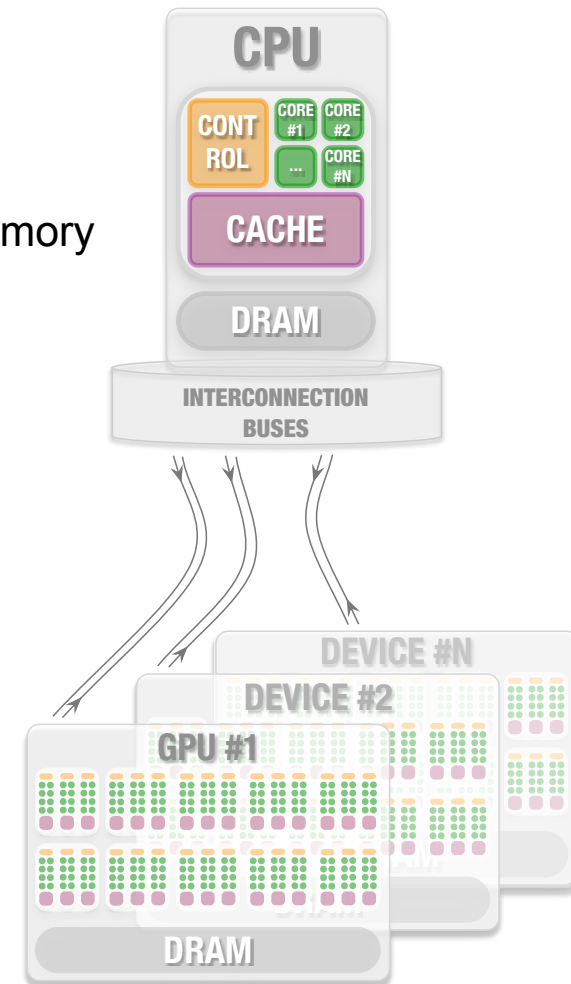- **APPLICABLE TO A WIDE RANGE OF SCIENTIFIC PROBLEMS**

  – Linear algebra, digital signal and image processing, database applications …

- **STATE OF THE ART DDL APPROACHES IN HETEROGENEOUS DISTRIBUTED COMPUTING**

  – Assume symmetric bandwidth and an one-port model for communication links

  – Limited memory: only the size of input load is considered; the exceeding load is simply redistributed among the nodes with available memory

  – Unrealistic: Computation/communication time is a linear/affine function of the #chunks

# HETEROGENEOUS STAR NETWORK (MASTER-WORKER)

- **MULTI-CORE CPU** (Master)

  - Global execution controller; access the whole global memory

  - All cores employed for execution

- **INTERCONNECTION BUSES**

  - Bidirectional full-duplex asymmetric communication
  - Different concurrency levels
  - Potential execution bottleneck

- **DEVICES** (Distant workers)

  - Different architectures and programming models
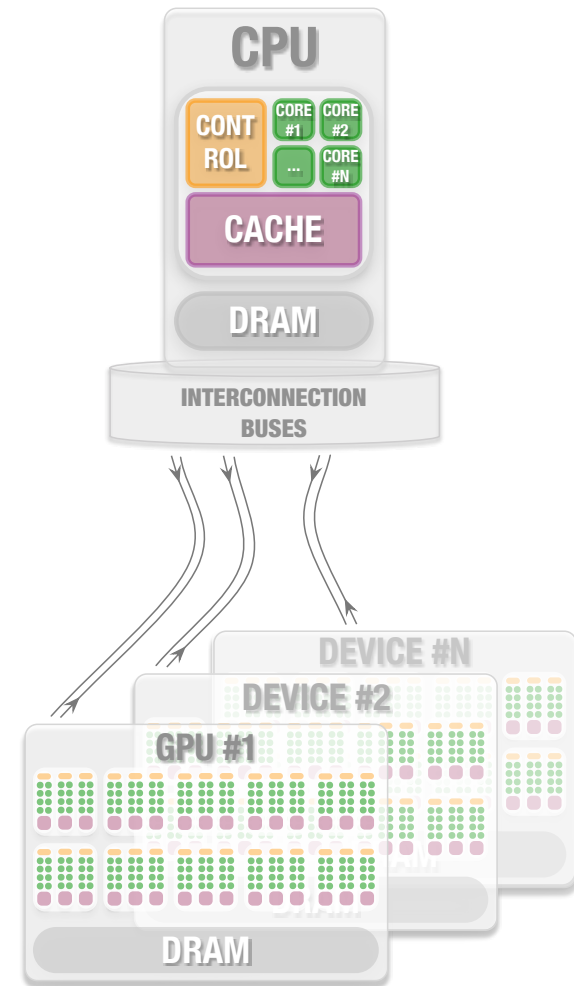
  - Computation performed using local memories

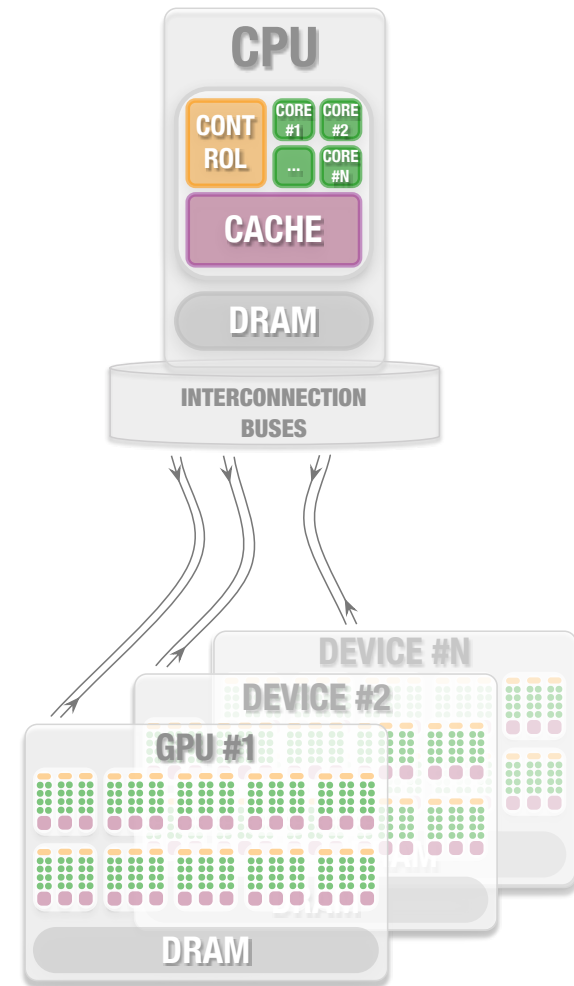technology
from seed

## 2-STEP DIVISIBLE LOAD SCHEDULING

DDL

## 2-STEP DIVISIBLE LOAD SCHEDULING

- ## STEP 1 – SYSTEM-LEVEL LOAD BALANCING

  - How many load units to send to each device?

technology
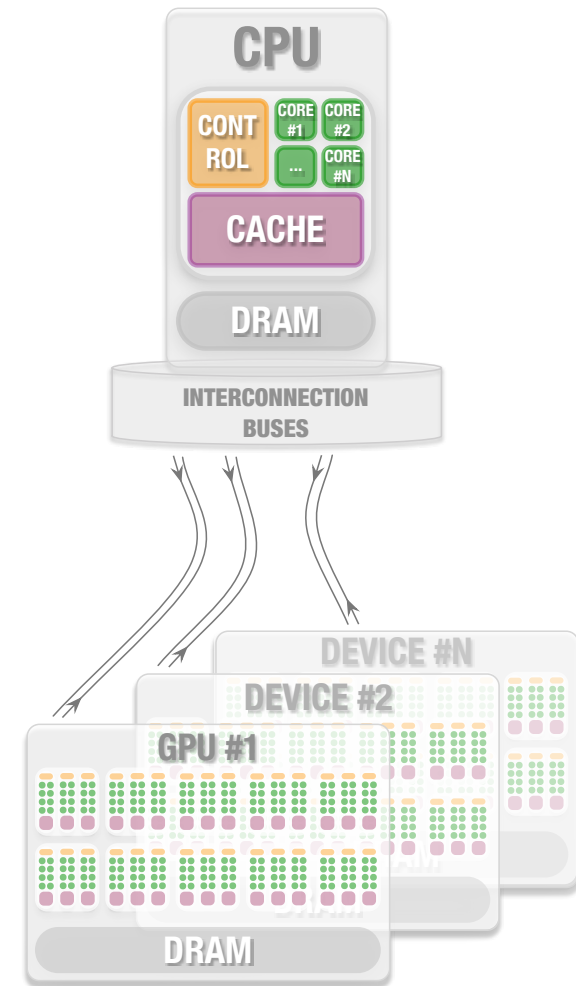from seed

inesc id
lisboa

## 2-STEP DIVISIBLE LOAD SCHEDULING

- ### STEP 1 – SYSTEM-LEVEL LOAD BALANCING

  – How many load units to send to each device?

## 2-STEP DIVISIBLE LOAD SCHEDULING

- **STEP 1 – SYSTEM-LEVEL LOAD BALANCING**

    - How many load units to send to each device?

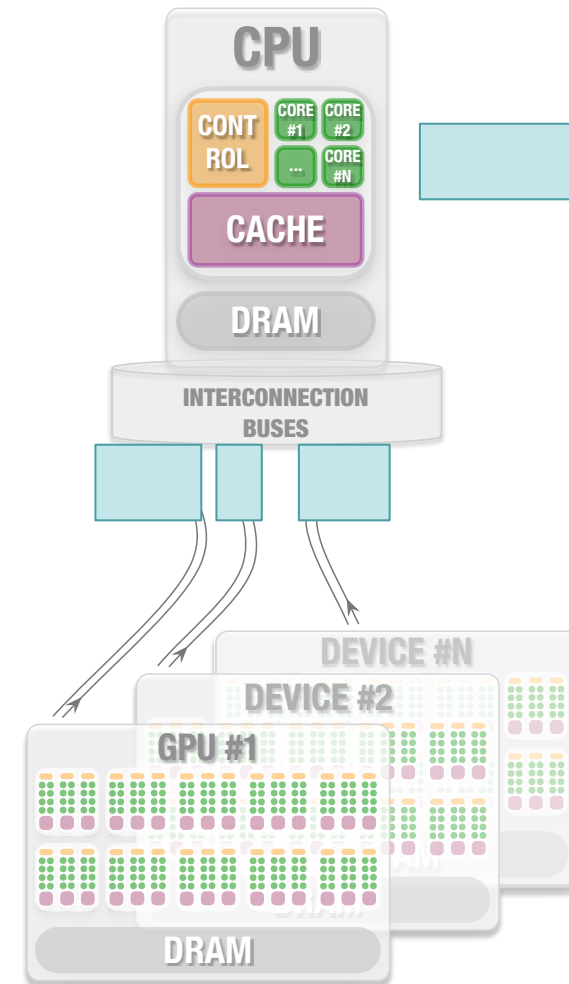- **STEP 2 – DEVICE-LEVEL LOAD SCHEDULING**

    - How to sub-partition the device load to:

        - Reduce delays when distributing and retrieving

        - Overlap computation and communication

        - Efficiently use the bidirectional asymmetric bandwidth of buses

        - Respect the amount of supported concurrency

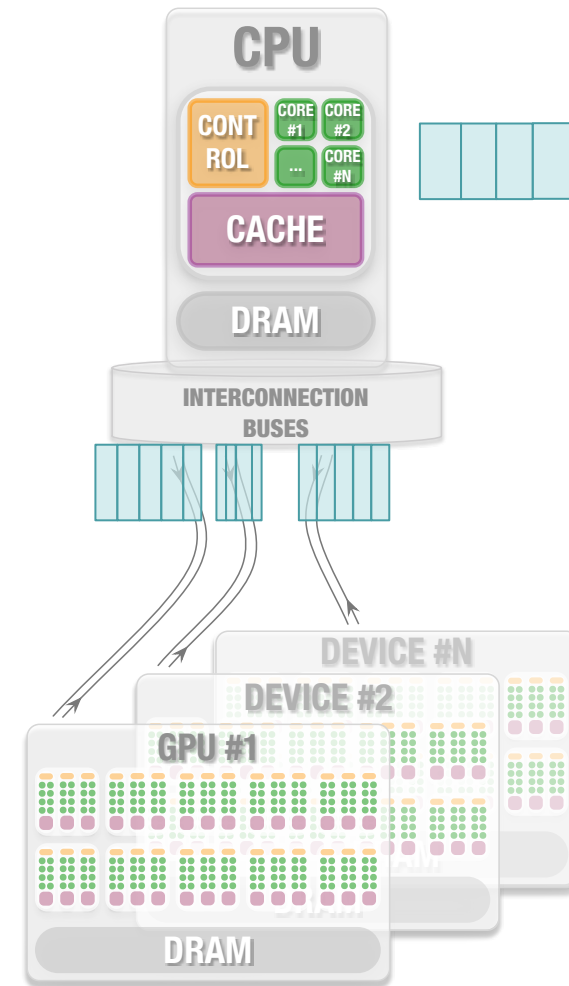        - Fit into device limited memory

**technology** from seed

## 2-STEP DIVISIBLE LOAD SCHEDULING

- ### STEP 1 – SYSTEM-LEVEL LOAD BALANCING

  – How many load units to send to each device?

- ### STEP 2 – DEVICE-LEVEL LOAD SCHEDULING

  – How to sub-partition the device load to:

    - Reduce delays when distributing and retrieving

    - Overlap computation and communication

    - Efficiently use the bidirectional asymmetric bandwidth of buses

    - Respect the amount of supported concurrency

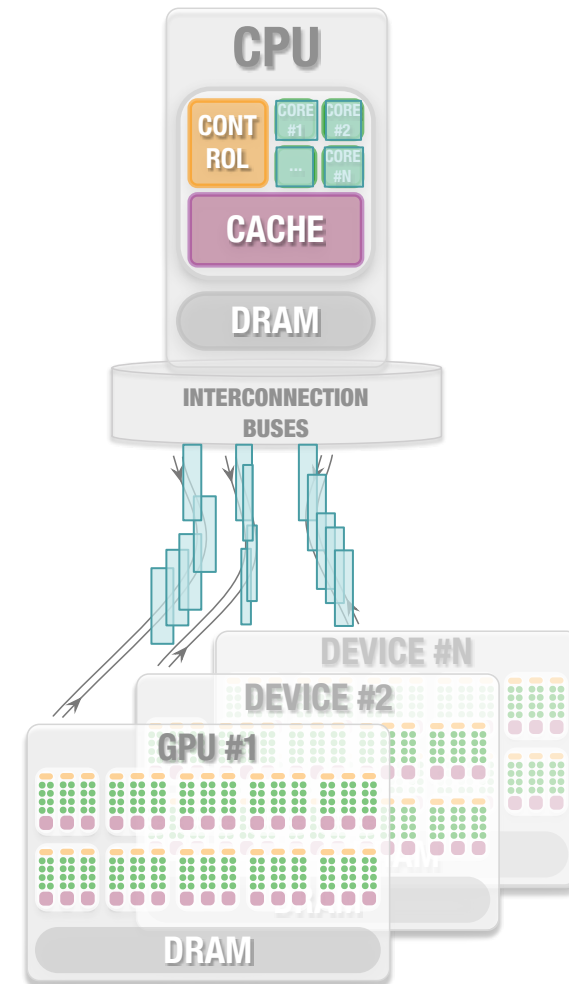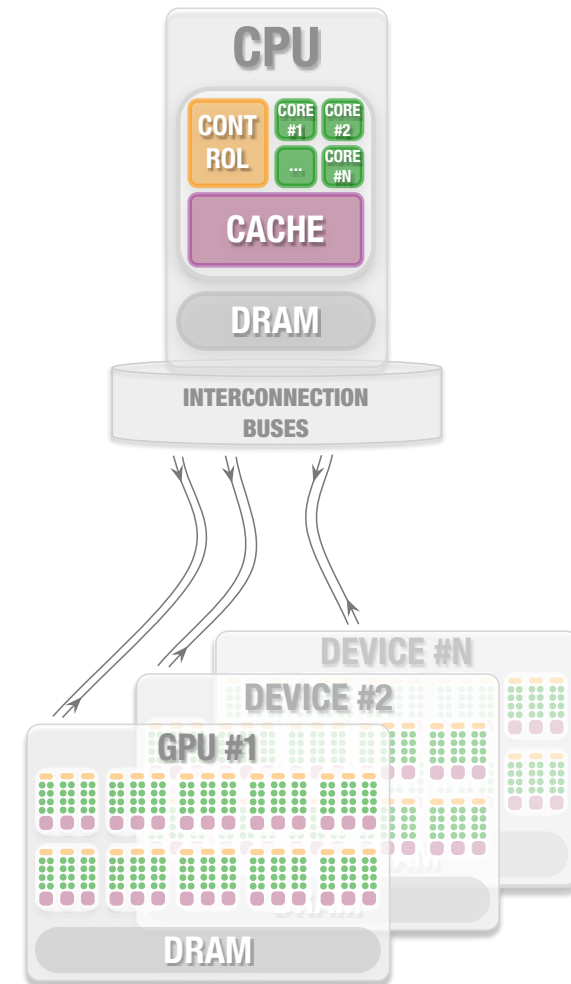    - Fit into device limited memory

## 2-STEP DIVISIBLE LOAD SCHEDULING

- **STEP 1 – SYSTEM-LEVEL LOAD BALANCING**

  – How many load units to send to each device?

- **STEP 2 – DEVICE-LEVEL LOAD SCHEDULING**

  – How to sub-partition the device load to:

    - Reduce delays when distributing and retrieving

    - Overlap computation and communication

    - Efficiently use the bidirectional asymmetric bandwidth of buses

    - Respect the amount of supported concurrency

    - Fit into device limited memory
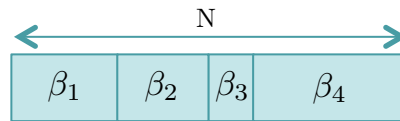
# FUNCTIONAL PERFORMANCE MODELS

- Continuous performance functions (#chunks/time)
- Built from the real application execution
  - *No assumptions being made to ease modeling!*

- ## COMPUTATION PERFORMANCE MODELS - $\psi_w$
  - For each master core and distant worker

- ## FULL-DUPLEX COMMUNICATION BANDWIDTH - $\sigma_\iota, \sigma_o$
  - Bidirectional and asymmetric for each link

- ## TOTAL PERFORMANCE $\psi_\tau$ OF EACH DEVICE
  - Including computation and communication

- ## STEP 1 – SYSTEM-LEVEL LOAD BALANCING

$$\xleftarrow{\qquad\qquad N \qquad\qquad}\rightarrow$$

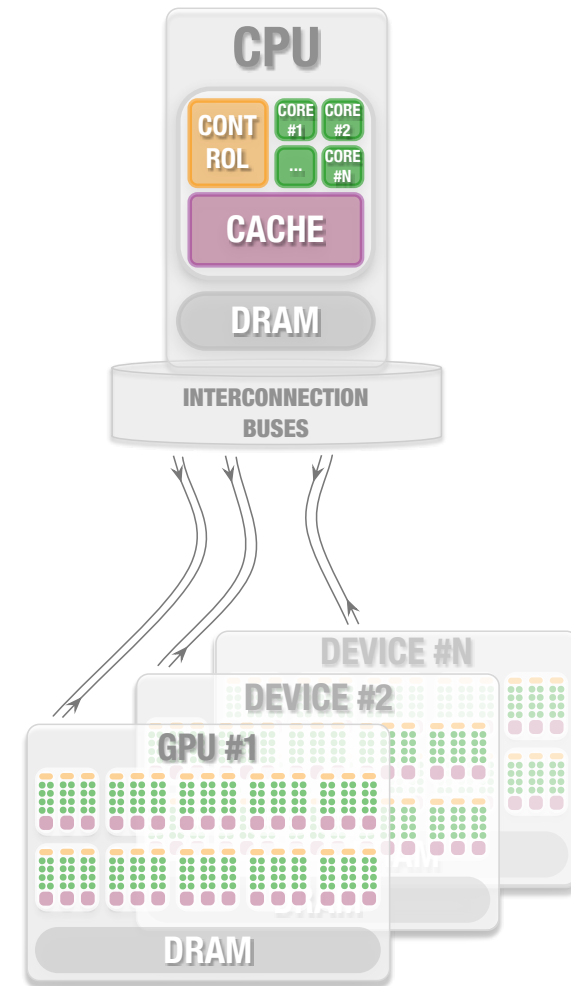| $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |

- The total load $N$ is partitioned between devices

- The optimal distribution lies on a straight line passing through the origin of coordinate system and intersecting communication-aware total performance curves ($\psi_\tau$), such that*:

$$\frac{\beta_1}{\psi_{\tau_1}(\beta_1)} = \cdots = \frac{\beta_k}{\psi_{\tau_k}(\beta_k)} = \cdots = \frac{\beta_n}{\psi_{\tau_n}(\beta_n)}$$
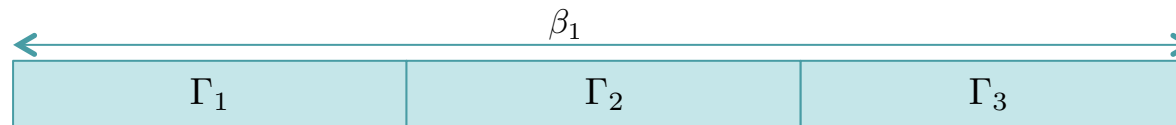
$$\sum_{j=1}^{n} \beta_j = N$$



*Lastovetsky, A., and R. Reddy, "Distributed Data Partitioning for Heterogeneous Processors Based on Partial Estimation of their Functional Performance Models", HeteroPar 2009, LNCS, vol. 6043, Springer, pp. 91-101, 2010.*

- **STEP 2 – DEVICE-LEVEL SCHEDULING**

  - Per-device distributions $\beta_j$ are allowed to exceed the device memory limits, $b_j$

    - Device-level **multi-installment** processing with **multi-distributions**

      - $\Gamma_k$ sub-distributions with $\gamma_{k,l}$ sub-load fractions

- **STEP 2 – DEVICE-LEVEL SCHEDULING**

  – Per-device distributions $\beta_j$ are allowed to exceed the device memory limits, $b_j$

    • Device-level **multi-installment** processing with **multi-distributions**

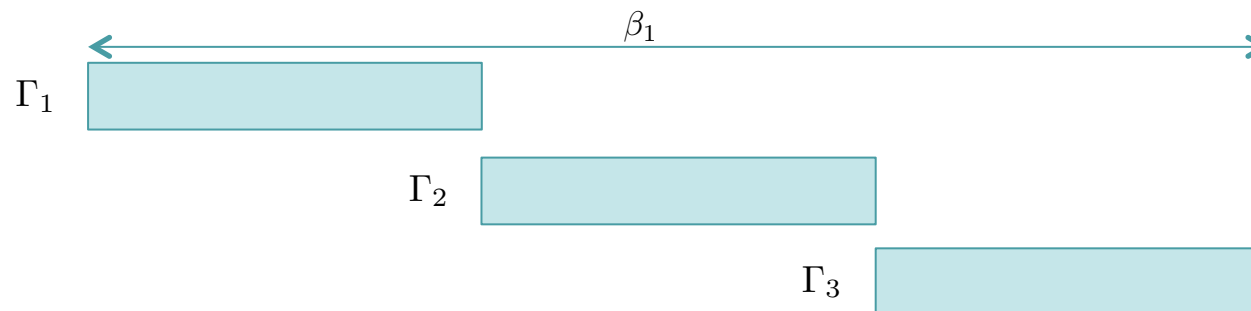      – $\Gamma_k$ sub-distributions with $\gamma_{k,l}$ sub-load fractions
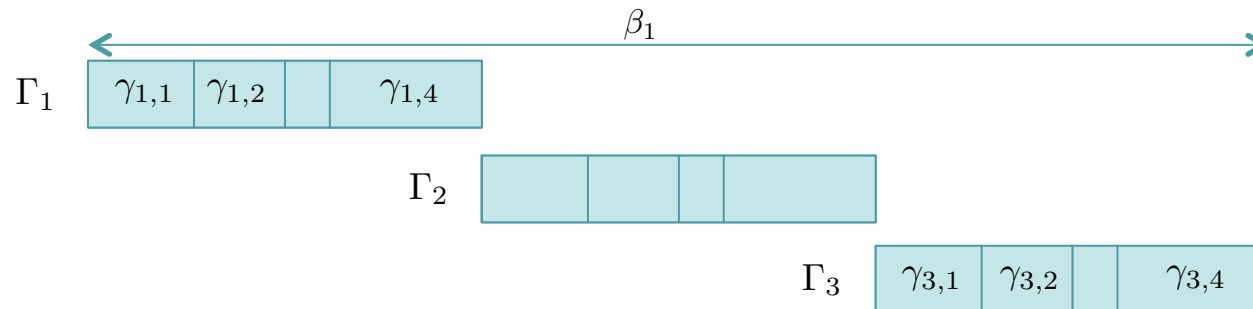
## Determination of Load Fractions (4)

- **STEP 2 – DEVICE-LEVEL SCHEDULING**
  - Per-device distributions $\beta_j$ are allowed to exceed the device memory limits, $b_j$
    - Device-level **multi-installment** processing with **multi-distributions**
      - $\Gamma_k$ sub-distributions with $\gamma_{k,l}$ sub-load fractions

- **STEP 2 – DEVICE-LEVEL SCHEDULING**

  - Per-device distributions $\beta_j$ are allowed to exceed the device memory limits, $b_j$

    - Device-level **multi-installment** processing with **multi-distributions**

      - $\Gamma_k$ sub-distributions with $\gamma_{k,l}$ sub-load fractions

  - Application **memory requirements** are modeled with three functions of load size:

    - *Input memory requirements*, $\mu_\iota(x)$

    - *Output memory requirements*, $\mu_o(x)$

    - **Execution memory requirements**, $\mu_w(x, P)$

      - Different implementations of the same problem might have different memory requirements!

$\Rightarrow$ Each $\Gamma_k$ *sub-distribution* may request *the* whole amount of memory, such that:

$$\sum_{l=1}^{|\Gamma_k|} \left( \mu_\iota(\gamma_{k,l}) + \mu_w(\gamma_{k,l}, p_j) + \mu_o(\gamma_{k,l}) \right) \leq b_j$$
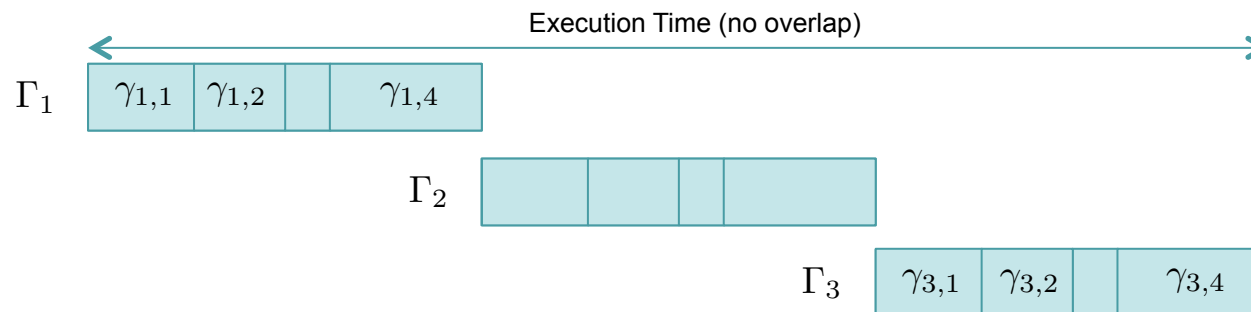
$$\sum_{k=1}^{|\Gamma|} \sum_{l=1}^{|\Gamma_k|} \gamma_{k,l} = \beta_j$$

**Determination of Load Fractions (6)**
**- Computation/Communication Overlapping -**

technology
from seed

inesc id
lisboa

- **STEP 2 – DEVICE-LEVEL SCHEDULING**

  – For each $\Gamma_k$ sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible *overlapping of computation and communication* between subsequent sub-fractions
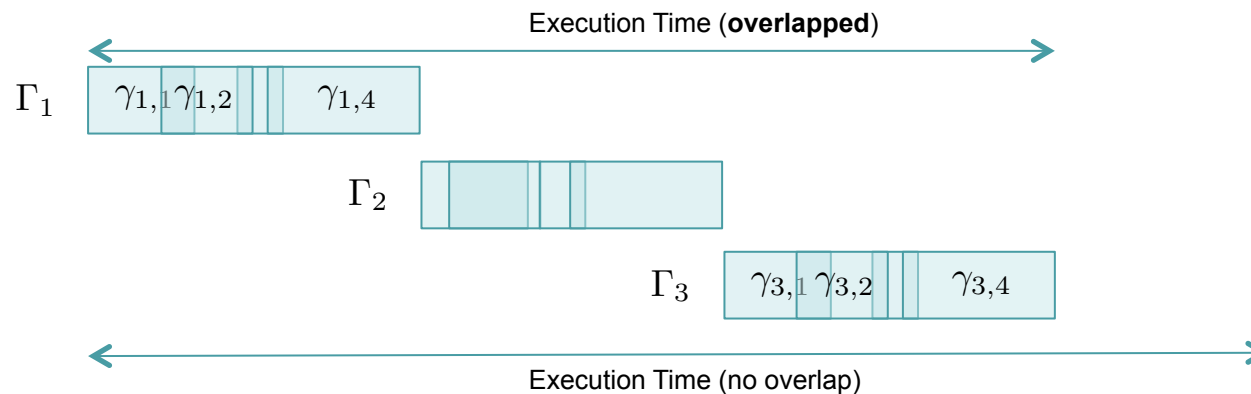
Execution Time (no overlap)

$\Gamma_1$ — $\gamma_{1,1}$ | $\gamma_{1,2}$ | | $\gamma_{1,4}$

$\Gamma_2$

$\Gamma_3$ — $\gamma_{3,1}$ | $\gamma_{3,2}$ | | $\gamma_{3,4}$

**Determination of Load Fractions (7)**
**- Computation/Communication Overlapping -**

technology
from seed

inesc id
lisboa

- ## STEP 2 – DEVICE-LEVEL SCHEDULING

  - For each $\Gamma_k$ sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible *overlapping of computation and communication* between subsequent sub-fractions
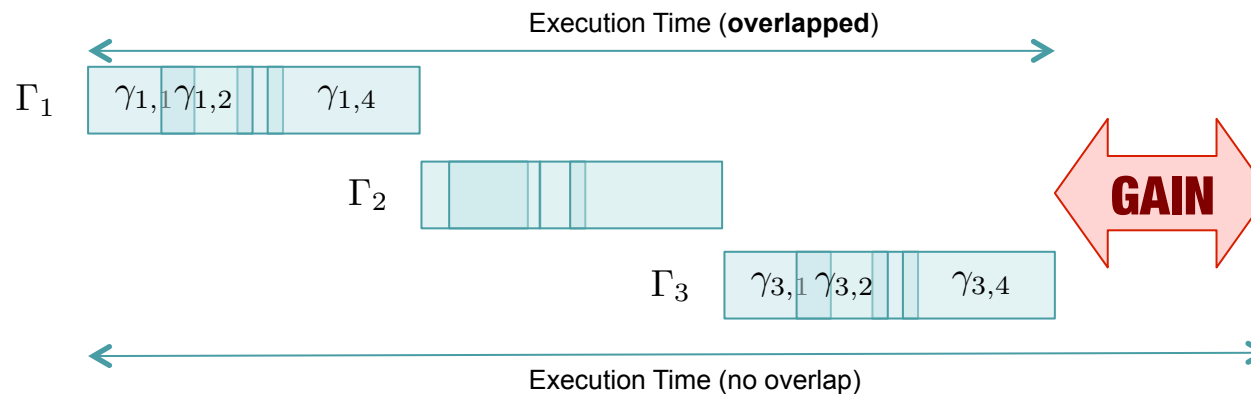
- **STEP 2 – DEVICE-LEVEL SCHEDULING**

  - For each $\Gamma_k$ sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible *overlapping of computation and communication* between subsequent sub-fractions
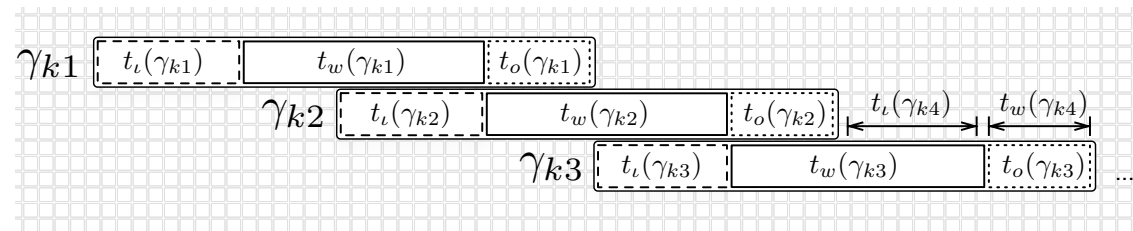


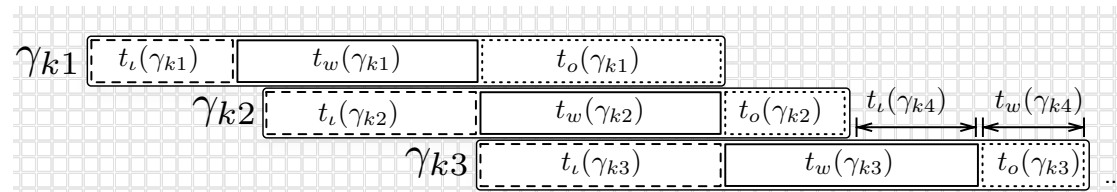  - Gain size depends on how well the chunks are overlapped!

- **STEP 2 – DEVICE-LEVEL SCHEDULING**

  - For each $\Gamma_k$ sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible *overlapping of computation and communication* between subsequent sub-fractions

  - The decisions are made according to the **amount of overlapping concurrency** supported by the device:



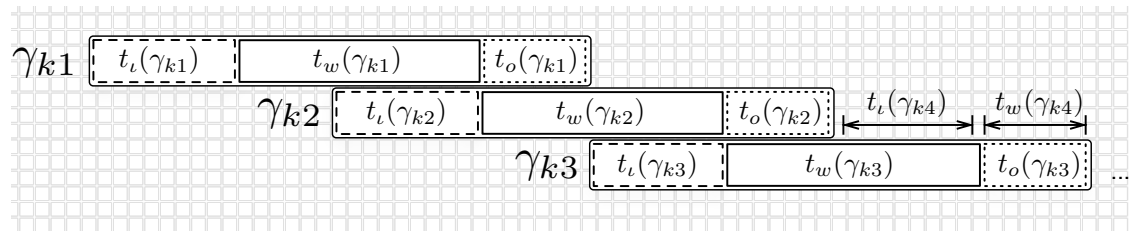(a) Overlap of a single communication with computation at the time



(b) Complete concurrency between communication and computation

- **STEP 2 – DEVICE-LEVEL SCHEDULING ALGORITHM\***



  – According to the performance models for device computation ($\psi_w$) and bidirectional asymmetric full-duplex communication links ($\sigma_\iota, \sigma_o$)
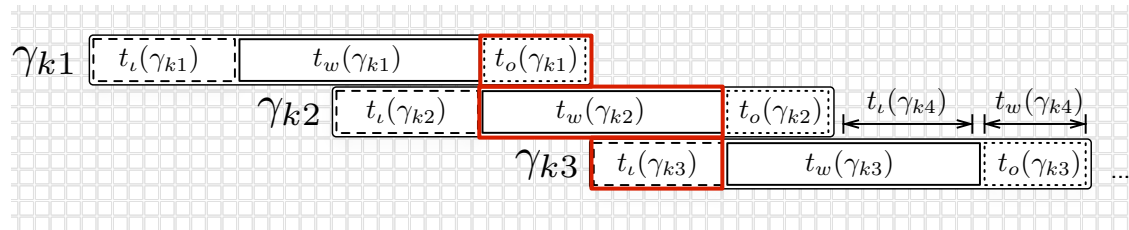
\* Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)

- **STEP 2 – DEVICE-LEVEL SCHEDULING ALGORITHM***



  - According to the performance models for device computation ($\psi_w$) and bidirectional asymmetric full-duplex communication links ($\sigma_\iota, \sigma_o$)

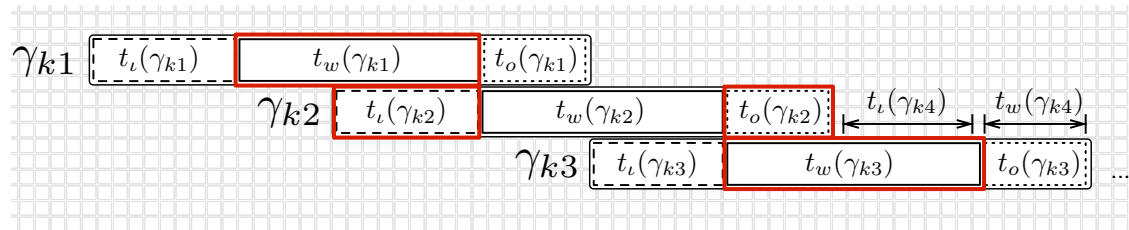    - **Step 2-I.** Determination of the initial optimal distribution with three load fractions.

*\* Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)*

- **STEP 2 – DEVICE-LEVEL SCHEDULING ALGORITHM\***



  – According to the performance models for device computation ($\psi_w$) and bidirectional asymmetric full-duplex communication links ($\sigma_\iota, \sigma_o$)

    - **Step 2-I.** Determination of the initial optimal distribution with three load fractions.

    - **Step 2-II.** Generate additional three-fraction distributions.

\* Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)

technology
from seed

- **STEP 2 – DEVICE-LEVEL SCHEDULING ALGORITHM\***



  - According to the performance models for device computation ($\psi_w$) and bidirectional asymmetric full-duplex communication links ($\sigma_\iota, \sigma_o$)

    - **Step 2-I.** Determination of the initial optimal distribution with three load fractions.

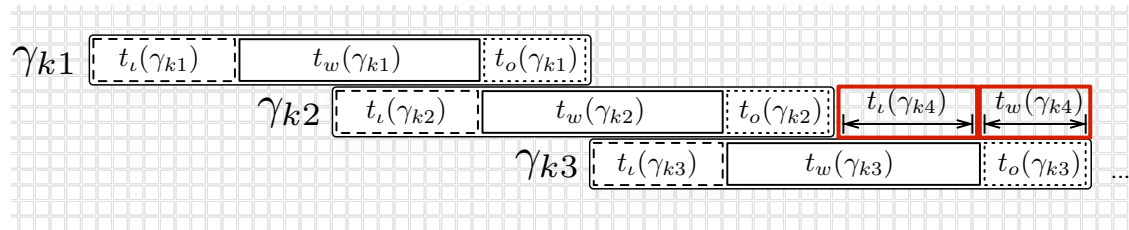    - **Step 2-II.** Generate additional three-fraction distributions.

    - **Step 2-III.** Insert additional load fractions into existing sub-distributions (iterative).

\* Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)

- **STEP 2 – DEVICE-LEVEL SCHEDULING ALGORITHM\***



  - According to the performance models for device computation ($\psi_w$) and bidirectional asymmetric full-duplex communication links ($\sigma_\iota, \sigma_o$)

    - **Step 2-I.** Determination of the initial optimal distribution with three load fractions.

    - **Step 2-II.** Generate additional three-fraction distributions.

    - **Step 2-III.** Insert additional load fractions into existing sub-distributions (iterative).
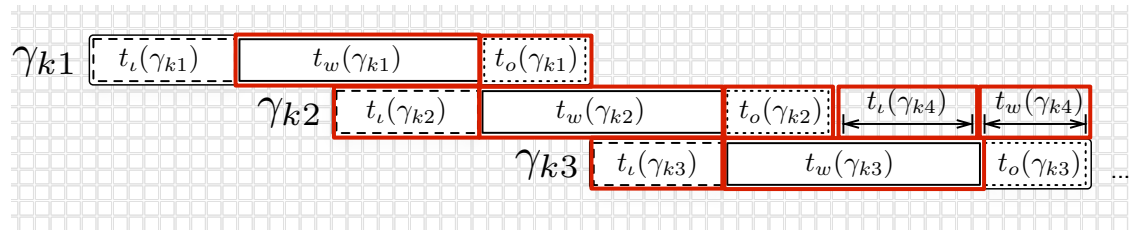
    - **Step 2-IV.** Generate new sub-distributions by restarting.

\* Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)

- ## STEP 2 – DEVICE-LEVEL SCHEDULING ALGORITHM*



- – According to the performance models for device computation ($\psi_w$) and bidirectional asymmetric full-duplex communication links ($\sigma_\iota, \sigma_o$)

  - **Step 2-I.** Determination of the initial optimal distribution with three load fractions.

  - **Step 2-II.** Generate additional three-fraction distributions.

  - **Step 2-III.** Insert additional load fractions into existing sub-distributions (iterative).

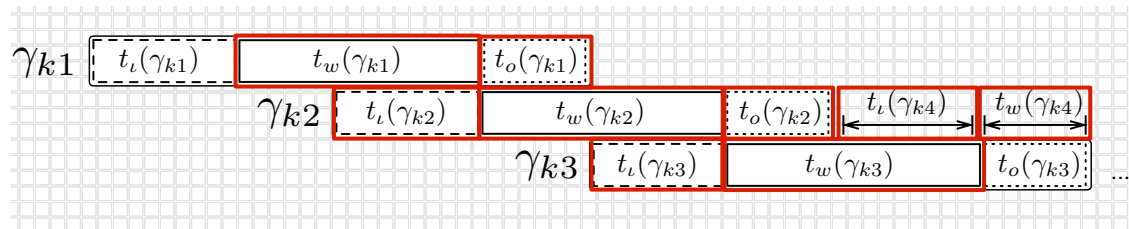  - **Step 2-IV.** Generate new sub-distributions by restarting.
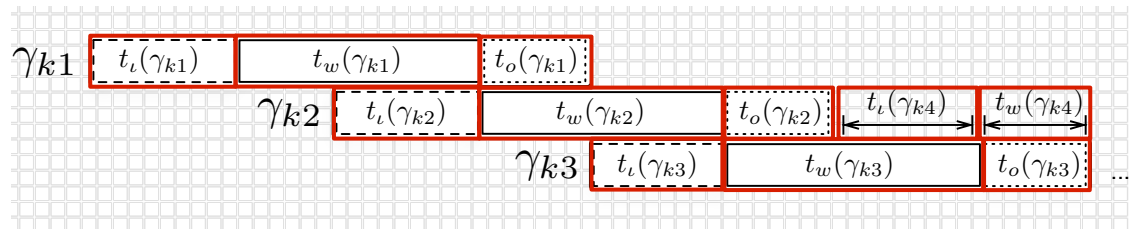
  - **Step 2-V.** Expand all sub-distributions.

\* Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)

- ## STEP 2 – DEVICE-LEVEL SCHEDULING ALGORITHM*



- According to the performance models for device computation ($\psi_w$) and bidirectional asymmetric full-duplex communication links ($\sigma_\iota, \sigma_o$)

  - **Step 2-I.** Determination of the initial optimal distribution with three load fractions.

  - **Step 2-II.** Generate additional three-fraction distributions.

  - **Step 2-III.** Insert additional load fractions into existing sub-distributions (iterative).

  - **Step 2-IV.** Generate new sub-distributions by restarting.

  - **Step 2-V.** Expand all sub-distributions.

  - **Step 2-VI.** Select the distribution with maximum relative performance.

*Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)*

technology
from seed

## DOUBLE FLOATING POINT COMPLEX 2D FFT BATCH EXECUTION

– Size: 256 times 512 × 512; divisible in the first dimension
– The optimal vendor-provided FFT implementations are used
  – NVIDIA's CUFFT 3.2 for the GPU and Intel MKL 10.3 for the CPU

## HETEROGENEOUS CPU+GPU DESKTOP SYSTEM

| Experimental Setup | CPU | GPU |
|---|---|---|
| | Intel Core 2 Quad | nVIDIA GeForce 285GTX |
| Speed/Core (GHz) | 2.83 | 1.476 |
| Global Memory (MB) | 4096 | 1024 |

## ITERATIVE PROCEDURE FOR ONLINE PERFORMANCE MODELING

– PERFORMANCE ESTIMATION of all heterogeneous devices DURING THE EXECUTION
  – No prior knowledge on the performance of an application is available on any of the devices
– **INITIALLY,** the load is distributed among devices using FACTORING-BY-TWO STRATEGY
  – Limited Memory: Factoring-by-two partitioning of the largest loads into new sub-distributions until satisfying the memory limitations
– **IN EACH FOLLOWING ITERATION,** the load is distributed using the PRESENTED APPROACH
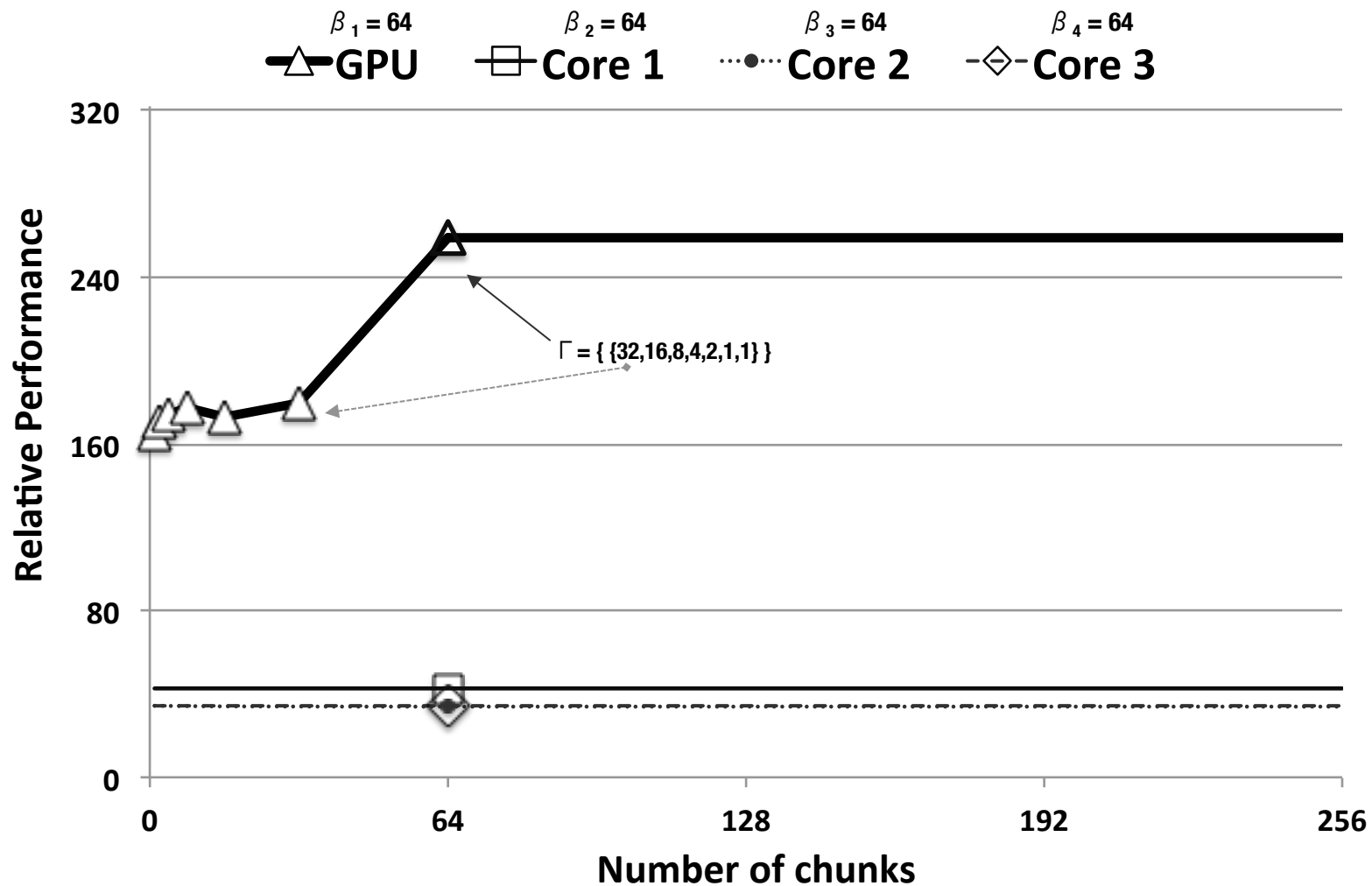
$\beta_1 = 64$  △—GPU  $\beta_2 = 64$  ⊟—Core 1  $\beta_3 = 64$  •—Core 2  $\beta_4 = 64$  ◇—Core 3

$\Gamma = \{ \{32,16,8,4,2,1,1\} \}$

**8 points for GPU modeling
(7 actual, 1 accuracy)**

Relative Performance

Number of chunks

Case Study: 2D FFT Batch
ITERATION 1: PROPOSED ALGORITHM

$\beta_1 = 181$ △ GPU   $\beta_2 = 29$ ⊟ Core 1   $\beta_3 = 23$ ⋯●⋯ Core 2   $\beta_4 = 23$ ◇ Core 3

$\Gamma = \{ \{8,14,18,20,9,1\}$
$\{8,14,8,20,21,21,9\} \}$

Relative Performance
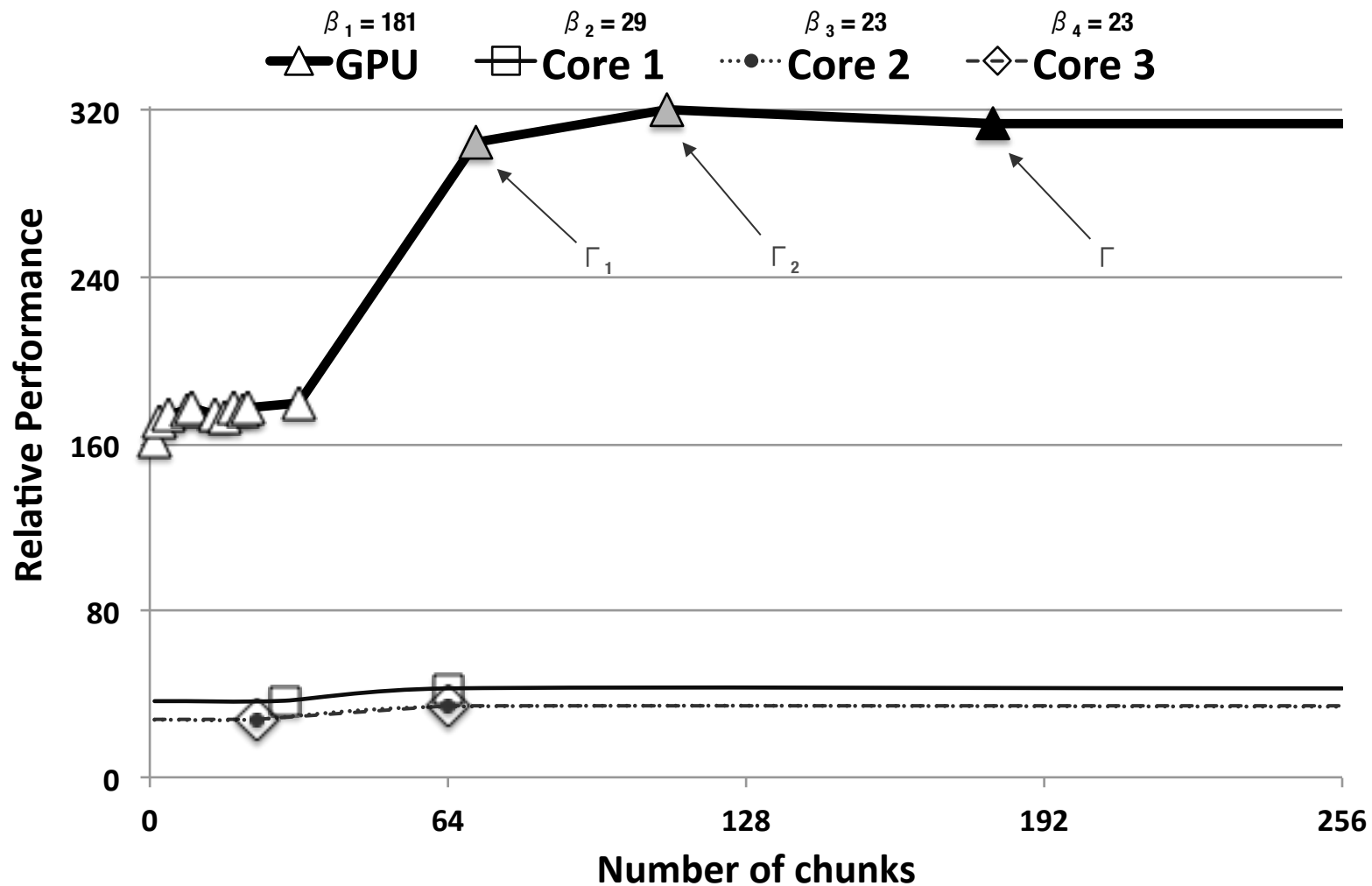
Number of chunks

Recalculate Distributions (from obtained models)

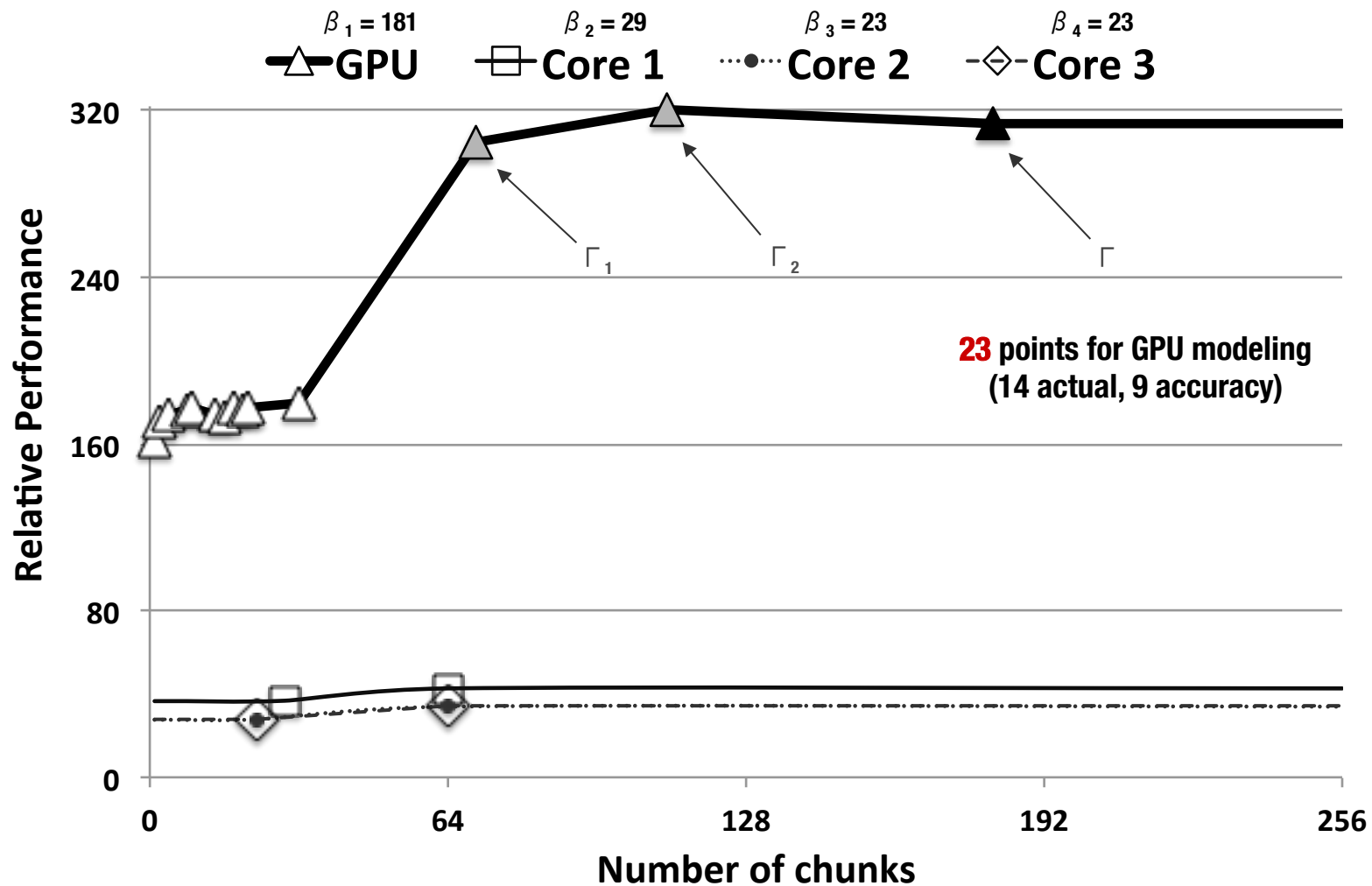technology from seed

Case Study: 2D FFT Batch
ITERATION 2: PERFORMANCE MODELS

Case Study: 2D FFT Batch
ITERATION 2: MODELING EFFICIENCY

## Case Study: 2D FFT Batch
### ITERATION 2: ALGORITHM EFFICIENCY

technology
from seed

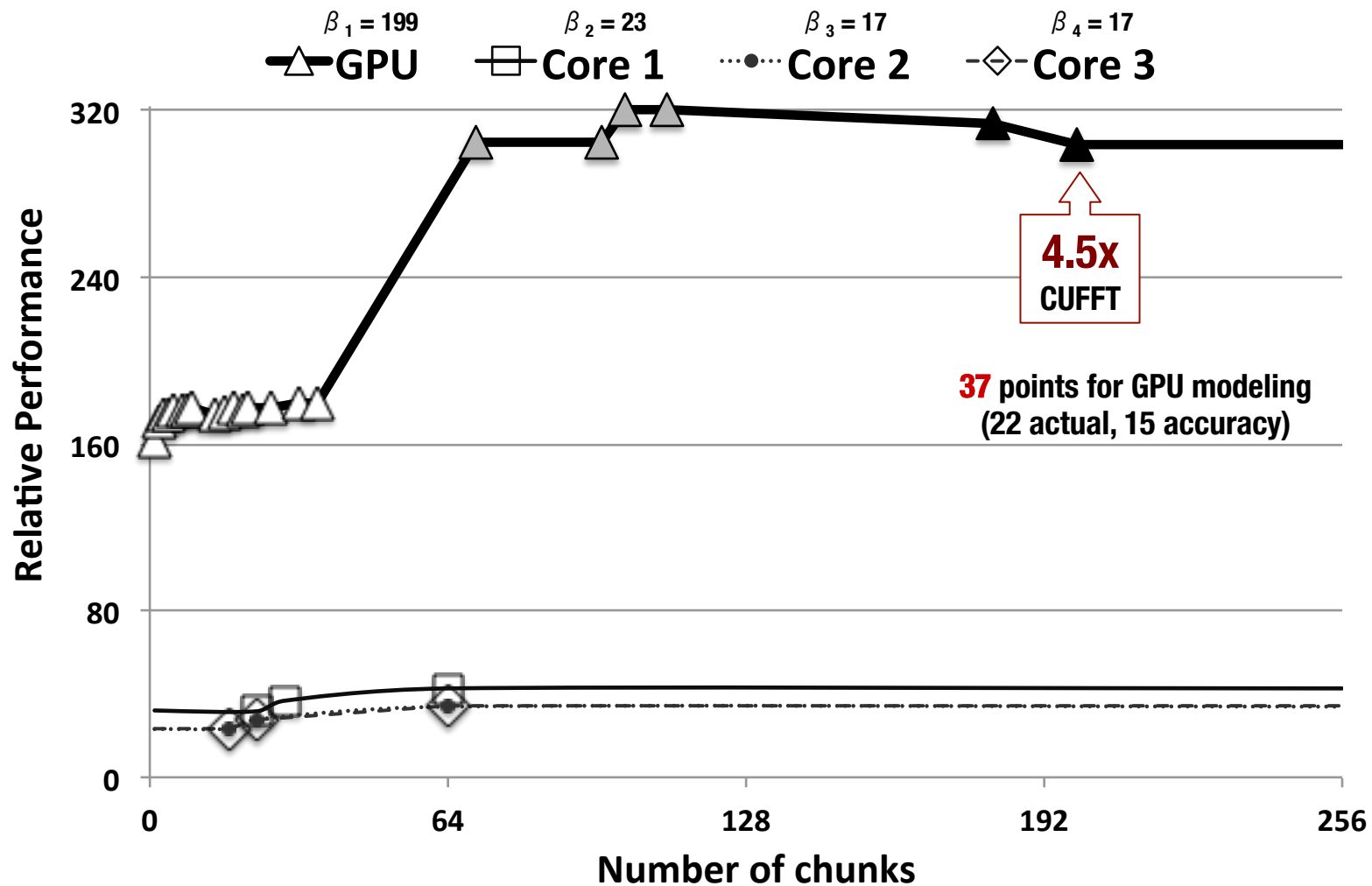inesc id
lisboa

$\beta_1 = 181$ △ GPU  $\beta_2 = 29$ Core 1  $\beta_3 = 23$ Core 2  $\beta_4 = 23$ Core 3

4.6x
CUFFT

**23** points for GPU modeling
(14 actual, 9 accuracy)

Case Study: 2D FFT Batch
ITERATION 3

# Case Study: 2D FFT Batch
## COMPARISON WITH STATE OF THE ART APPROACHES



|  | Proposed | [10] | [6] |
|---|---|---|---|
| #Iterations to converge | **4** | 10 | 8 |
| Load Balancing | **YES** | **NO** | **NO** |
| #Points for GPU Modeling | **53** (4 iters) | 10 (10 iters) | 8 (8 iters) |
| #Points per iteration | **13.25** | 1 | 1 |

- **4.3x faster** in determining the steady-state distribution (complete algorithm time) comparing to the approach from [10], and **3.2x faster** comparing to [6] (ping-pong state)

- Load balanced distribution achieves **2x better performance** comparing to the steady-state distribution from [10], and **2.2x better performance** than [6]

[6] Galindo, I., Almeida, F., Badıa-Contelles, J.M.: Dynamic load balancing on dedicated heterogeneous systems. In: PVM/MPI. pp. 64–74 (2008)

[10] Lastovetsky, A., Reddy, R.: Data partitioning with a functional performance model of heterogeneous processors. Int. J. High Perform. Comput. Appl. 21, 76–90 (2007)

## DYNAMIC LOAD BALANCING

– OBTAINED IN **4 ITERATIONS** AND **4.1 SECONDS** (IN TOTAL)

## TRADITIONAL APPROACHES FOR PERFORMANCE MODELING

– Approximate the performance using number of points equal to the number of iterations
– In this case, 4 **POINTS** in total for GPU performance modeling (4 iterations)

## PRESENTED DDL SCHEDULING APPROACH

– Models the GPU performance using **53 POINTS**, in this case ~13x more than with traditional modeling
– Load balancing solution is 2x faster than the current state of the art approach
– **IN THIS CASE, OBTAINED GPU PERFORMANCE IS AT LEAST 4.3X BETTER THAN THE "OPTIMAL" CUFFT EXECUTION**

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa