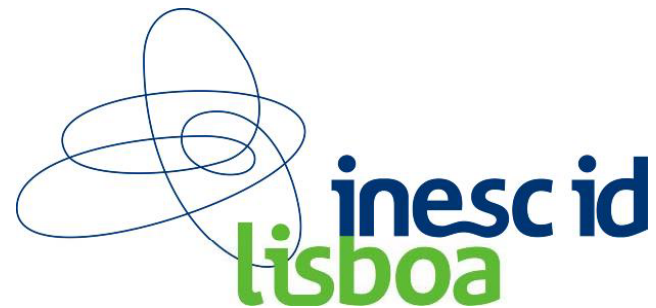# Cooperative Execution on Heterogeneous Multi-core Systems

**Leonel Sousa** and Aleksandar Ilić

**INESC-ID/IST, TU Lisbon**

COMMODITY COMPUTERS = HETEROGENEOUS SYSTEMS

– Multi-core General-Purpose Processors (CPUs)

– Many-core Graphic Processing Units (GPUs)

– …

– Special accelerators, co-processors, FPGAs

=> HUGE COMPUTING POWER

– Not yet completely explored for **COLLABORATIVE COMPUTING**

HETEROGENEITY MAKES PROBLEMS MUCH MORE COMPLEX!

– Performance modeling and load balancing

– Different programming models and languages

COLLABORATIVE ENVIRONMENT FOR HETEROGENEOUS COMPUTERS
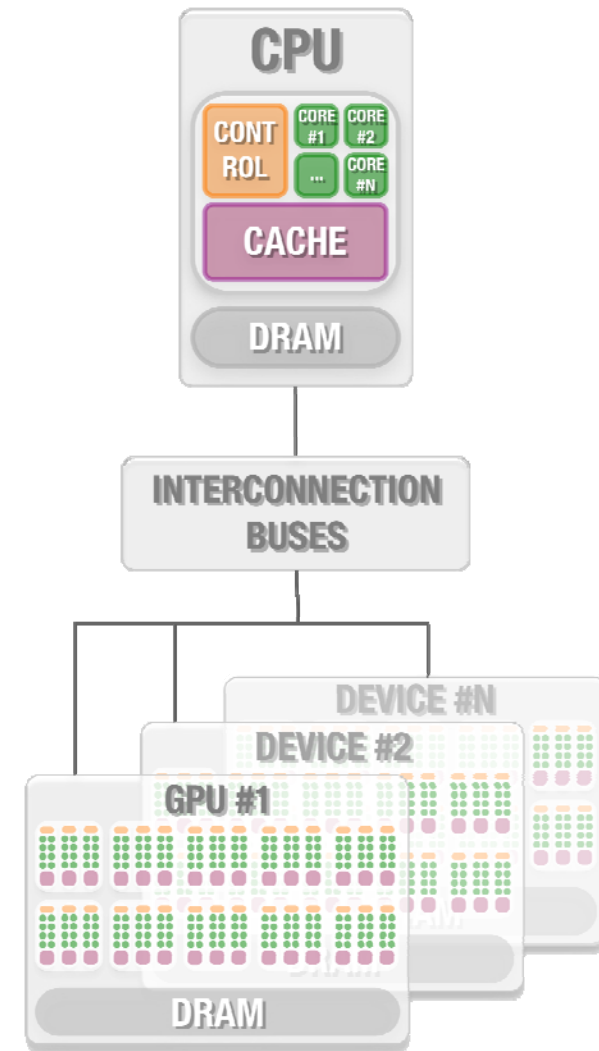
PERFORMANCE MODELING AND LOAD BALANCING

– State of the art for heterogeneous systems
– for CPU+GPU

CASE STUDY: 2D BATCH FAST FOURIER TRANSFORM
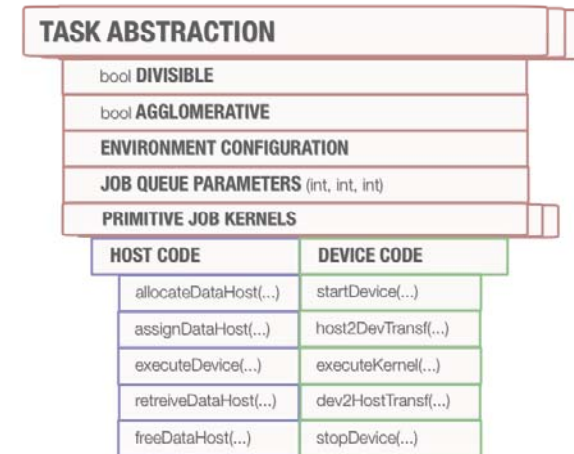
CONCLUSIONS AND FUTURE WORK

**MASTER-SLAVE** paradigm

- CPU (Master)

  – Global execution controller

  – Access the whole global memory

- INTERCONNECTION BUSSES

  – Limited and asymmetric communication bandwidth
  – Potential execution bottleneck

- UNDERLYING DEVICES (Slaves)

  – Different architectures and programming models

  – Computation performed using local memories

# Redefining Tasks and Primitive Jobs

**TASK –** basic programming unit (coarser-grained)

- CONFIGURATION PARAMETERS
  - Task: application and task dependency information
  - Environment: device type, number of devices…

- PRIMITIVE JOB WRAPPER
  - DIVISIBLE TASK – comprise several finer-grained Primitive Jobs
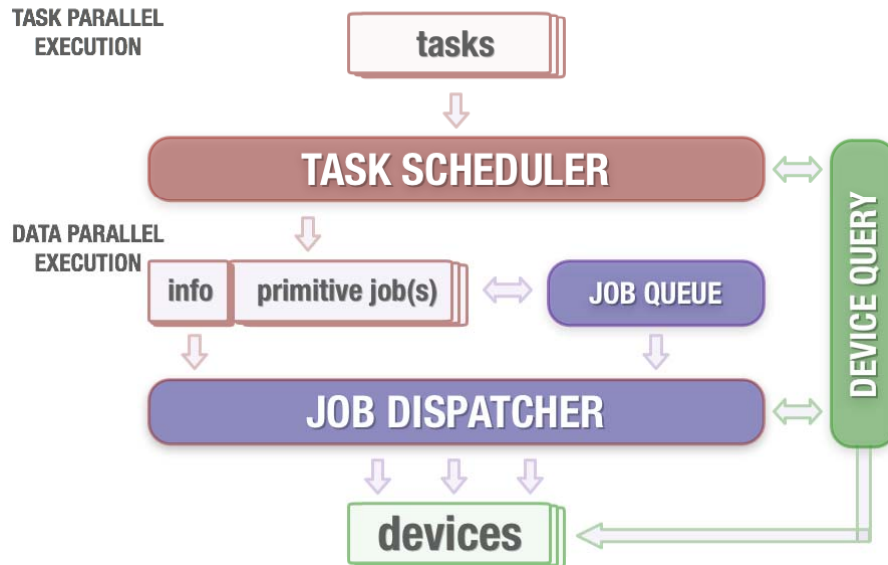  - AGGLOMERATIVE TASK – allows grouping of Primitive Jobs

**PRIMITIVE JOB –** minimal program portion for parallel execution

- CONFIGURATION PARAMETERS
  - I/O and performance specifics, …

- CARRIES PER-DEVICE-TYPE IMPLEMENTATIONS
  - Vendor-specific programming models and tools
  - Specific optimization techniques

| TASK ABSTRACTION | |
|---|---|
| bool DIVISIBLE | |
| bool AGGLOMERATIVE | |
| ENVIRONMENT CONFIGURATION | |
| JOB QUEUE PARAMETERS (int, int, int) | |
| PRIMITIVE JOB KERNELS | |
| HOST CODE | DEVICE CODE |
| allocateDataHost(...) | startDevice(...) |
| assignDataHost(...) | host2DevTransf(...) |
| executeDevice(...) | executeKernel(...) |
| retreiveDataHost(...) | dev2HostTransf(...) |
| freeDataHost(...) | stopDevice(...) |

| Primitive Job Granularity | Task Type | |
|---|---|---|
| | Divisible | Agglomerative |
| Coarser-grained | NO | -- |
| Balanced | YES | NO |
| Finer/Balanced | YES | YES |

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

# Collaborative Execution Environment for Heterogeneous Systems*



TASK PARALLEL EXECUTION

DATA PARALLEL EXECUTION

tasks → TASK SCHEDULER → info | primitive job(s) ⇔ JOB QUEUE → JOB DISPATCHER → devices — DEVICE QUERY

| Task Type | |
| --- | --- |
| Divisible | Agglomerative |
| NO | -- |
| YES | NO |
| YES | YES |

## Task Level Parallelism

– TASK SCHEDULER submits independent tasks to JOB DISPATCHER in respect to task and environment configuration parameters and current platform state from DEVICE QUERY structure

## Data Level Parallelism

– PRIMITIVE JOBS may be arranged into JOB QUEUES (currently, 1D-3D grid organization) for DIVISIBLE (AGGLOMERATIVE) TASKS

– JOB DISPATCHER uses DEVICE QUERY and JOB QUEUE information to map (agglomerated) PRIMITIVE JOBS to the requested devices; then initiates and controls further execution;

## Nested Parallelism

– If provided, JOB DISPATCHER can be configured to perceive certain number of cores of a multi-core device as a single device
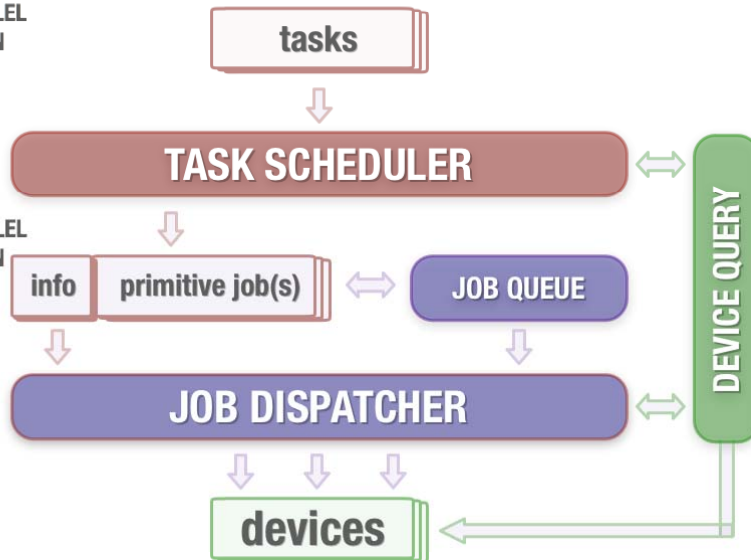
*Aleksandar Ilic and Leonel Sousa. "Collaborative Execution Environment for Heterogeneous Parallel Systems", In 12th Workshop on Advances in Parallel and Distributed Computational Models (APDCM/IPDPS 2010), April 2010.

# Collaborative Execution Environment for Heterogeneous Systems*

TASK PARALLEL EXECUTION

tasks

TASK SCHEDULER

DATA PARALLEL EXECUTION

info | primitive job(s) | JOB QUEUE

DEVICE QUERY

JOB DISPATCHER

devices

| Task Type | |
|-----------|---------------|
| Divisible | Agglomerative |
| NO | -- |
| YES | NO |
| YES | YES |

## PROBLEM

How to make good DYNAMIC LOAD BALANCING decisions using PERFORMANCE MODELS of the devices aware of :
- application demands
- implementation specifics
- platform / device heterogeneity
- complex memory hierarchies
- limited asymmetric communication bandwidth
- …

# CONSTANT PERFORMANCE MODELS (CPM)

- DEVICE PERFORMANCE (SPEED) : constant positive number
  - Typically represents relative speed when executing a serial benchmark of a given size
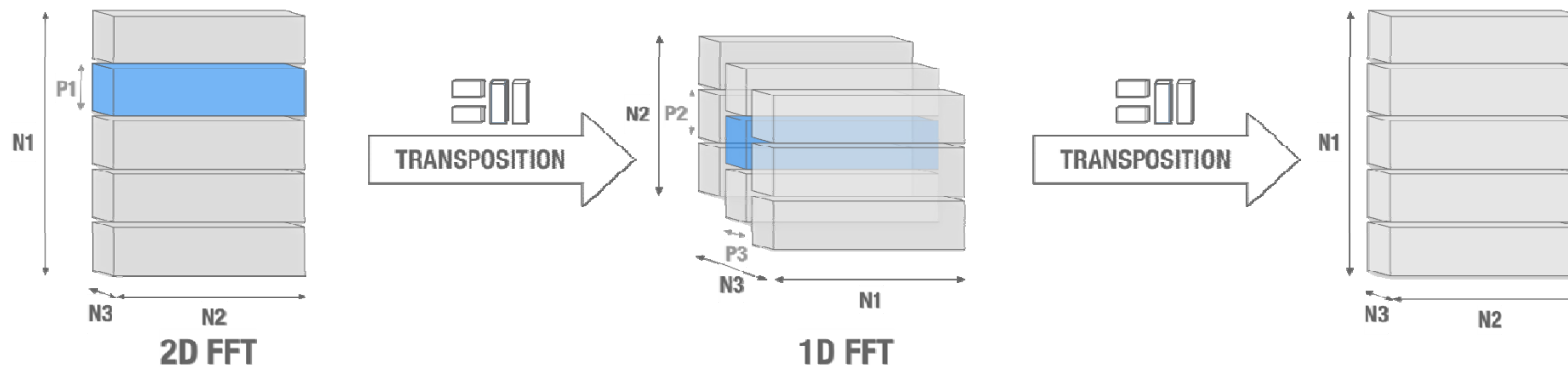- COMPUTATION DISTRIBUTION : proportional to the speed of device

# FUNCTIONAL PERFORMANCE MODELS (FPM)

- DEVICE PERFORMANCE (SPEED) : continuous function of the problem size
  - Typically require several benchmark runs and significant amount of time for building it
- COMPUTATION DISTRIBUTION : relies on the "functional speed" of the processor

# FPM VS. CPM

- MORE REALISTIC : integrates features of heterogeneous processor
  - Processor heterogeneity, the heterogeneity of memory structure, and the other effects (such as paging)
- MORE ACCURATE DISTRIBUTION of computation across heterogeneous devices
- APPLICATION-CENTRIC approach characterize speed for different applications with different functions

# Case Study : 2D FFT Batches



Part of a PARALLEL 3D FFT PROCEDURE : $H = FFT_{1D}(FFT_{2D}(h))$
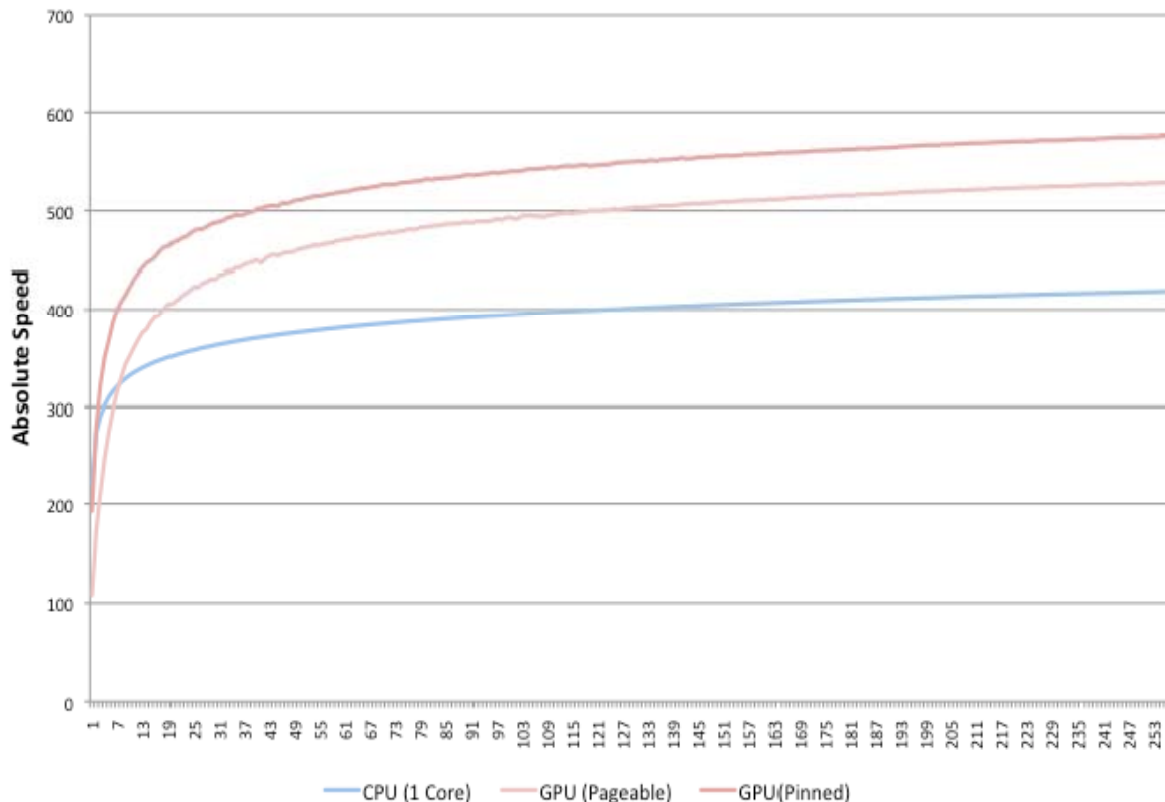
- Very HIGH COMMUNICATION-TO-COMPUTATION ratio

PROBLEM DEFINITION

- $N_{in} = N_{out} = N_1N_2N_3 * sizeof(data)$

- Performance [FLOPS] : $N_1N_2N_3 * log(N_1N_2N_3)$

- $N_2 = const; N_3 = const;$

- $N_1$ – total number of computational chunks (PRIMITIVE JOBS)

| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | Iteration |
|---|---|---|---|---|---|---|



CPU (1 Core) — GPU (Pageable) — GPU(Pinned)

## METRIC : ABSOLUTE SPEED

– p devices: $P_1$, $P_2$,…, $P_p$
– $N_1$ total #Primitive Jobs (chunks)
– Device Load [chunks]: $n_1$, $n_2$,…, $n_p$

– Absolute speed:
$$s_i(n_i) = n_i/t_i(n_i), \ 1 \le i \le p$$

## SOLUTION: OPTIMAL LOAD BALANCING

Lies on the straight line that passes through the origin of coordinate system, such that:

$$x_1/s_1(x_1) = x_2/s_2(x_2) =…= x_p/s_p(x_p)$$
$$x_1 + x_2 + … + x_p = N_1$$
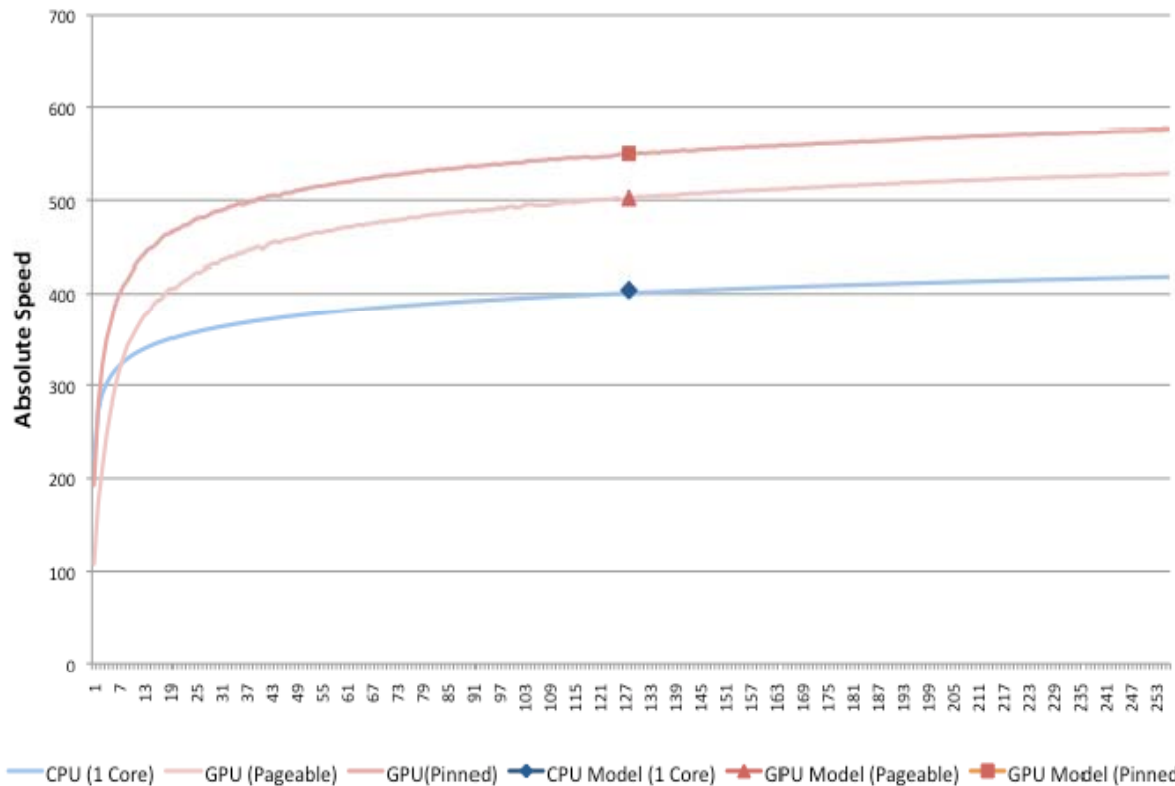
| Performance Metric | ⇨ | **Initialization** | ⇨ | Approximation | ⇨ | Iteration |



① All the $P$ computational units execute $N_1/p$ 2D FFT Batches **in parallel**

$$n_i = N_1/p, \quad 1 \le i \le p$$

② Record execution times: $t_i(N_1/p)$

③ IF $\max_{1 \le i, j \le p}\{((t_i(N_1/p) - t_j(N_1/p))/t_i(N_1/p)\} \le \varepsilon$
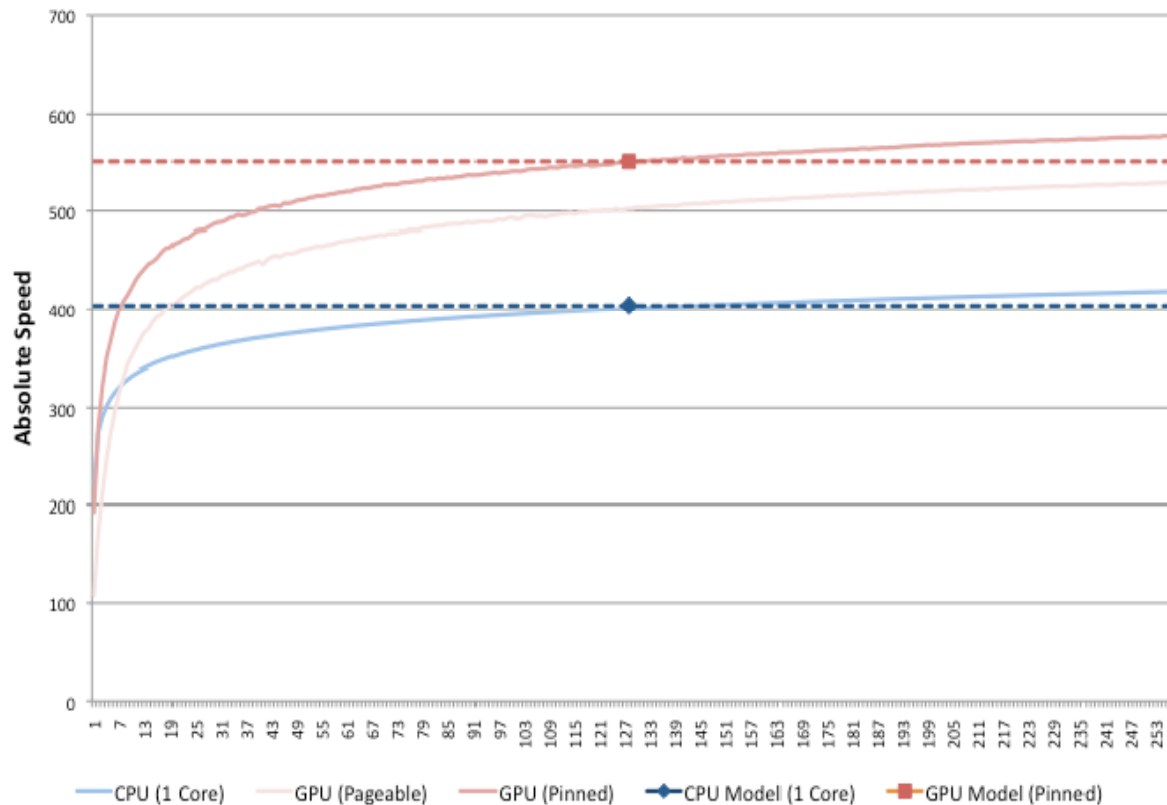THEN even distribution solves the problem and the algorithm stops;
ELSE absolute speeds of devices are calculated, such that:

$$s_i(N_1/p) = (N_1/p)/t_i(N_1/p), \quad 1 \le i \le p$$

Legend: CPU (1 Core) — GPU (Pageable) — GPU(Pinned) — CPU Model (1 Core) — GPU Model (Pageable) — GPU Model (Pinned)

*Lastovetsky, A., and R. Reddy, "Distributed Data Partitioning for Heterogeneous Processors Based on Partial Estimation of their Functional Performance Models", HeteroPar 2009, Netherlands, Lecture Notes in Computer Science, vol. 6043, Springer, pp. 91-101, 25/9/2009, 2010.*

| Performance Metric | ⟹ | Initialization | ⟹ | Approximation | ⟹ | Iteration |



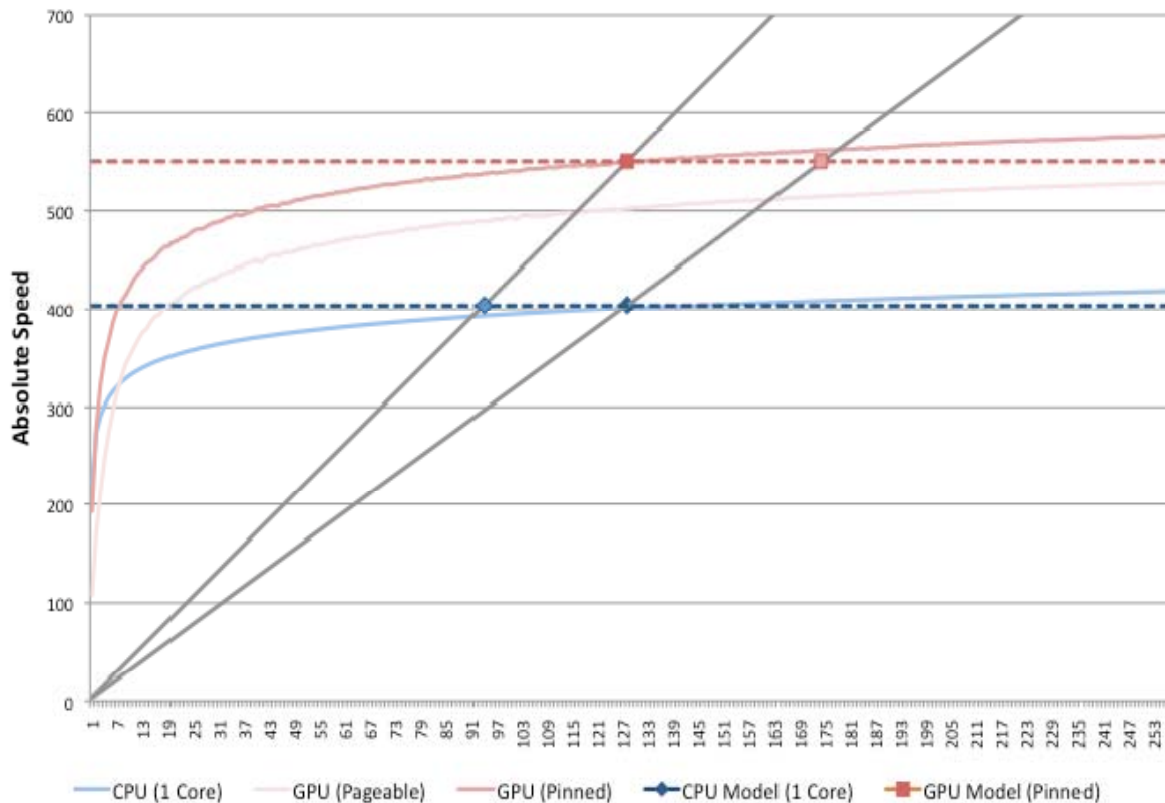① **Performance** of each device is **model**ed as a **constant**

$$s_i(x) = s_i(N_1/p), \quad 1 \le i \le p$$

*Lastovetsky, A., and R. Reddy, "Distributed Data Partitioning for Heterogeneous Processors Based on Partial Estimation of their Functional Performance Models", HeteroPar 2009, Netherlands, Lecture Notes in Computer Science, vol. 6043, Springer, pp. 91-101, 25/9/2009, 2010.*

# Case Study: Performance Modeling
# What we already know …** (4)

| Performance Metric | → | Initialization | → | Approximation | → | Iteration |
|---|---|---|---|---|---|---|



① Draw Upper $U$ and Lower $L$ lines through the following points:

`(0,0), (N₁/p, max_i{s_i(N_1/p)})`

$$(0,0), \quad (N_1/p, \ max_i\{s_i(N_1/p)\})$$

$$(0,0), \quad (N_1/p, \ min_i\{s_i(N_1/p)\})$$

② Let $x_i^{(U)}$, $x_i^{(L)}$ be the intersections with $s_i(x)$ IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$
THEN go to 3
ELSE go to 5

—CPU (1 Core) — GPU (Pageable) — GPU (Pinned) —◆—CPU Model (1 Core) —■—GPU Model (Pinned)

**Absolute Speed**

# Case Study: Performance Modeling
# What we already know …** (5)

| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | Iteration |



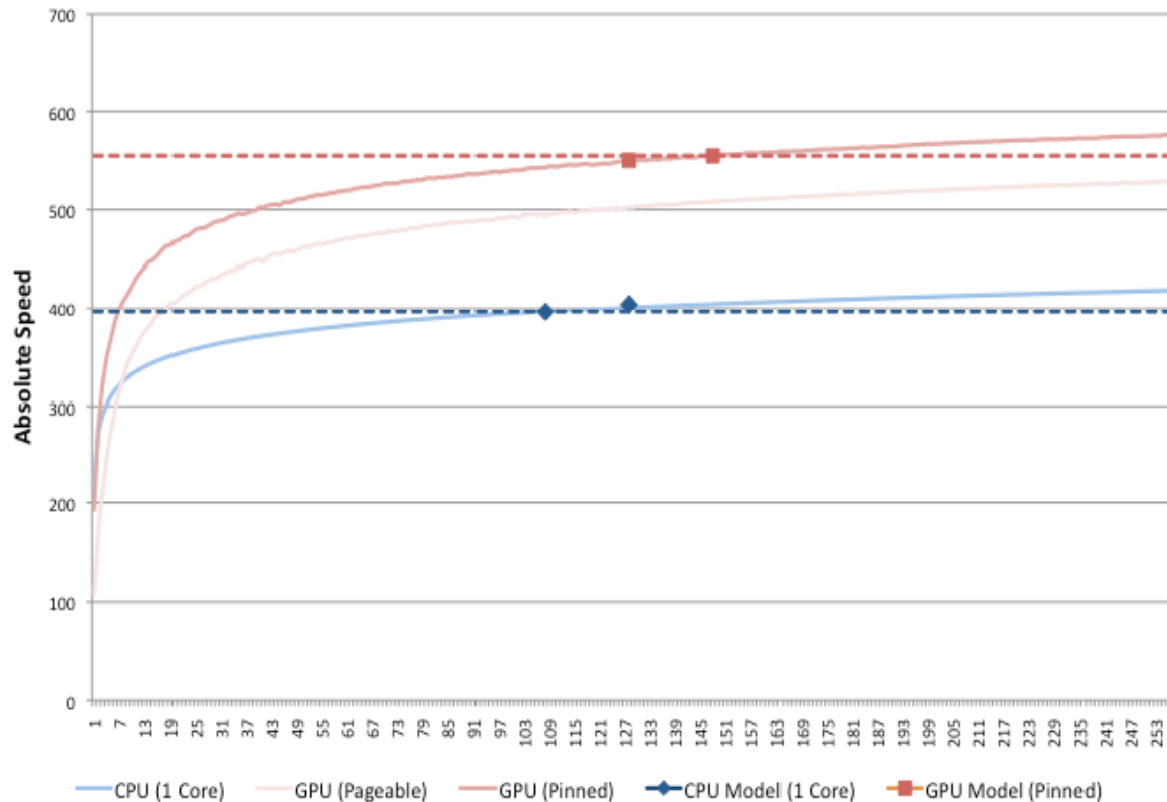CPU (1 Core) — GPU (Pageable) — GPU (Pinned) — CPU Model (1 Core) — GPU Model (Pinned)

1  Draw Upper $U$ and Lower $L$ lines through the following points:

$(0,0)$, $(N_1/p, \max_i\{s_i(N_1/p)\})$

$(0,0)$, $(N_1/p, \min_i\{s_i(N_1/p)\})$

2  Let $x_i^{(U)}$, $x_i^{(L)}$ be the intersections with $s_i(x)$ IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$
THEN go to 3
ELSE go to 5

③  Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$

**Lastovetsky, A., and R. Reddy, "Data Partitioning with a Functional Performance Model of Heterogeneous Processors", International Journal of High Performance Computing Applications, vol. 21, issue 1: Sage, pp. 76-90, 2007

# Case Study: Performance Modeling
## What we already know …** (6)

| Performance Metric | | Initialization | | Approximation | | Iteration |



1. Draw Upper $U$ and Lower $L$ lines through the following points:

   $(0,0), (N_1/p, \max_i\{s_i(N_1/p)\})$

   $(0,0), (N_1/p, \min_i\{s_i(N_1/p)\})$

2. Let $x_i^{(U)}, x_i^{(L)}$ be the intersections with $s_i(x)$ IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$
   THEN go to 3
   ELSE go to 5

3. Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$

4. IF $\Sigma_i\, x_i^{(M)} \leq N_1$
   THEN $U=M$
   ELSE $L=M$
   REPEAT 2

Chart legend: CPU (1 Core) — GPU (Pageable) — GPU (Pinned) — CPU Mode (1 Core) — GPU Model (Pinned)

**Lastovetsky, A., and R. Reddy, "Data Partitioning with a Functional Performance Model of Heterogeneous Processors", International Journal of High Performance Computing Applications, vol. 21, issue 1: Sage, pp. 76-90, 2007*

INSTITUTO SUPERIOR TÉCNICO

# Case Study: Performance Modeling
# What we already know … (7)

| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | Iteration |



A. Speeds of each device are modeled as **constants** equal to the observed values [1]

$$s_i(x) = s_i(n_i), \quad 1 \leq i \leq p$$

# Case Study: Performance Modeling
## What we already know … (8)

technology
from seed

inesc id
lisboa

Performance Metric ⟹ Initialization ⟹ **Approximation** ⟹ Iteration



A. Speeds of each device are modeled as **constants** equal to the observed values [1]

$$s_i(x) = s_i(n_i), \quad 1 \leq i \leq p$$

B. Functional performance modeling using **piecewise linear approximation** [2]

— CPU (1 Core)  — GPU (Pageable)  — GPU (Pinned)  ◆ CPU Model (1 Core)  ■ GPU Model (Pinned)

[1] *Galindo I., Almeida F., and Badía-Contelles J.M., "Dynamic Load Balancing on Dedicated Heterogeneous Systems", In EuroPVM/MPI 2008, Springer, pp. 64-74, 2008.*
[2] *Lastovetsky, A., and R. Reddy, "Distributed Data Partitioning for Heterogeneous Processors Based on Partial Estimation of their Functional Performance Models", HeteroPar 2009, vol. 6043, Springer, pp. 91-101, 25/9/2009, 2010*

# Goals in the CPU + GPU Environment

## 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
  – Modeling of the overall CPU and GPU performance for different problem sizes (+ kernel-only GPU performance)

## DYNAMIC LOAD BALANCING

– OPTIMAL DISTRIBUTION OF COMPUTATIONS (PRIMITIVE JOBS)
  – Partial estimations of the performance should be built and used to decide on optimal mapping
  – Returned solution should provide load balancing within a given accuracy

## COMMUNICATION AWARENESS

– MODELING THE BANDWIDTH for interconnection busses DURING THE EXECUTION
  – To select problem sizes that maximize the interconnection bandwidth
  – The algorithm should be aware of asymmetric bandwidth for Host-To-Device and Device-To-Host transfers

## CPU+GPU ARCHITECTURAL SPECIFICS

– Make use of ENVIRONMENT-SPECIFIC FUNCTIONS to ease performance modeling
  – Asynchronous transfers and CUDA streams to overlap communication with computation
  – Be aware of diverse capabilities of different devices, but also for devices of the same type (e.g. GT200 vs. Fermi)

# Case Study:
# Building Full Performance Models



## FULL PERFORMANCE MODELS: PER-DEVICE REAL PERFORMANCE

– Experimentally obtained using CPU+ GPU platform specifics (*Pageable/Pinned Memory*)

– Exhaustive search on the full range of problem sizes

– High cost of building it !!!

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

| 2D FFT Batch | CPU | GPU |
|---|---|---|
| FFT Type | Complex; Double Precision | |
| Total Size ($N_1$x$N_2$x$N_3$) | 256x256x256 | |
| **High Performance Software** | | |
| FFT | Intel MKL 10.2 | CUFFT 3.1 |

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

Performance Metric ⟹ Initialization ⟹ Approximation ⟹ Iteration



CPU (1 Core)    GPU (Pageable)    GPU(Pinned)

## PROBLEM DEFINITION: PRIMITIVE JOB

– $n^{in}$ – input data requirements
– $n^{out}$ – output data requirements
– $n^{size}$ – problem size
– $n^{flops}$ – #float-point operations

## METRIC : FLOPS

– p devices: $P_1$, $P_2$,…, $P_p$
– $N_1$ total #Primitive Jobs (chunks)
– Device Load [chunks]: $n_1$, $n_2$,…, $n_p$
– $n_i = \{n_i^{in}, n_i^{out}, n_i^{size}, n_i^{flops}\}$

– FLOPS:
   $Perf_i(n_i) = n_i^{flops}/t_i(n_i), \quad 1 \leq i \leq p$

## SOLUTION: OPTIMAL LOAD BALANCING

Lies on the straight line that passes through the origin of coordinate system, such that:

   $x_i^{flops} = nf(x_i), \quad x_i$ as $n_i^{size}$

   $nf(x_1)/Perf_1(x_1) = \ldots = nf(x_p)/Perf_p(x_p)$

   $x_1 + x_2 + \ldots + x_p = N_1$
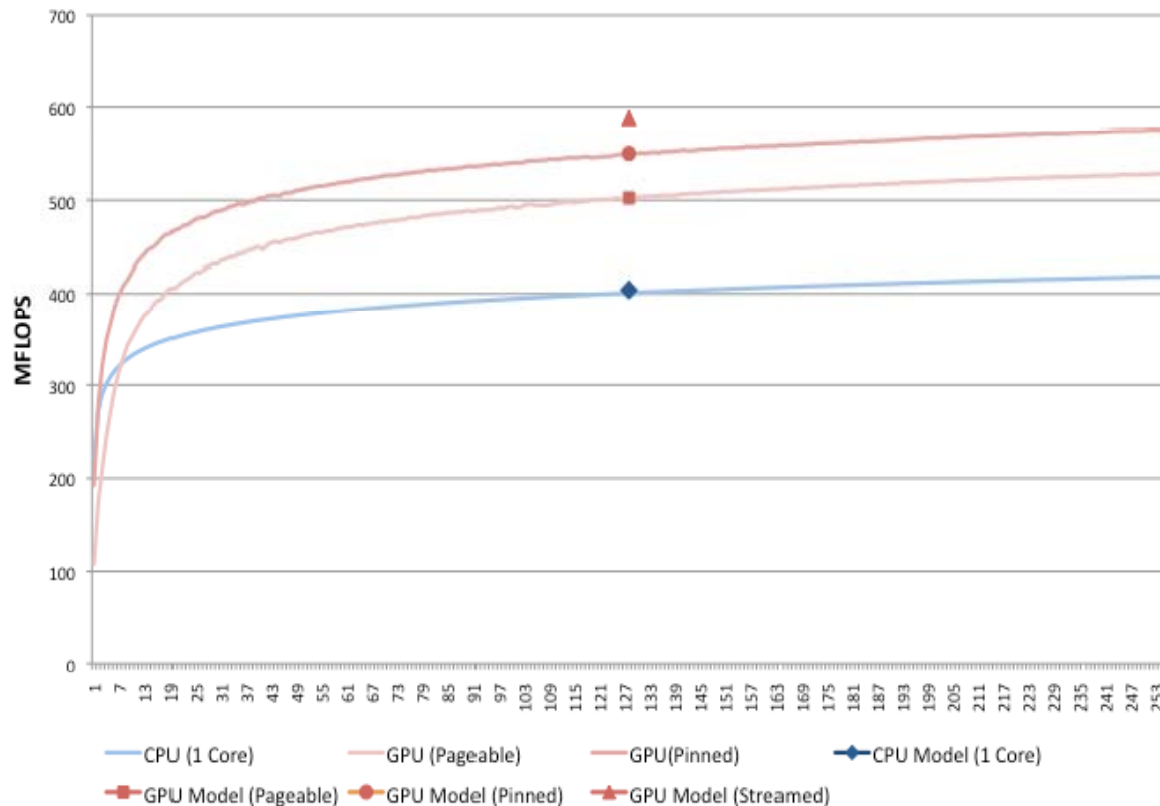
| Performance Metric | → | **Initialization** | → | Approximation | → | Iteration |

① All the $P$ computational units execute $N_1/p$ 2D FFT Batches **in parallel**

$$n_i = N_1/p, \quad 1 \le i \le p$$

② **IF** (device is GPU) AND (task is Divisible and Agglomerative)
**THEN** go to 3
**ELSE** go to 4

③ Split given computational load into streams and use asynchronous transfers to overlap communication with computation



Legend: CPU (1 Core) — GPU (Pageable) — GPU(Pinned)

| Performance Metric | $\Rightarrow$ | Initialization | $\Rightarrow$ | Approximation | $\Rightarrow$ | Iteration |
|---|---|---|---|---|---|---|

Streaming

- SUBDIVIDE $n_i$ computational chunks using DIV2 STRATEGY
  - No prior knowledge on the performance of an application!
  - The next stream has half the load of a previous stream
  - Algorithm may continue splitting the workload until the last stream is assigned with load equal to 1

- BANDWIDTH-AWARE DIV2 STRATEGY
  - Interconnection bandwidth is a subject to the amount of data that should be transferred and not to the application-specific demands
  - Run small pre-calibration tests for HOST-TO-DEVICE AND DEVICE-TO-HOST transfers
  - Tests can be stopped when saturation points are detected, or when transfers reach desired value (e.g. 60% of its theoretical)
  - CASE STUDY: $n^{min\_size} = 4$

1. All the $P$ computational units execute $N_1/p$ 2D FFT Batches **in parallel**

$$n_i = N_1/p, \quad 1 \le i \le p$$

2. IF (device is GPU) AND (task is Divisible and Agglomerative)
   THEN go to 3
   ELSE go to 4

③ Split given computational load into streams and use asynchronous transfers to overlap communication with computation

| Stream 0 |
|---|

| Stream 1 |
|---|

| Stream 2 |
|---|

Performance Metric ⟹ **Initialization** ⟹ Approximation ⟹ Iteration



1  All the $P$ computational units execute $N_1/p$ 2D FFT Batches **in parallel**

$$n_i = N_1/p, \quad 1 \le i \le p$$

2  IF `(device is GPU) AND (task is Divisible and Agglomerative)`
THEN  go to 3
ELSE  go to 4

3  Split given computational load into streams and use asynchronous transfers to overlap communication with computation
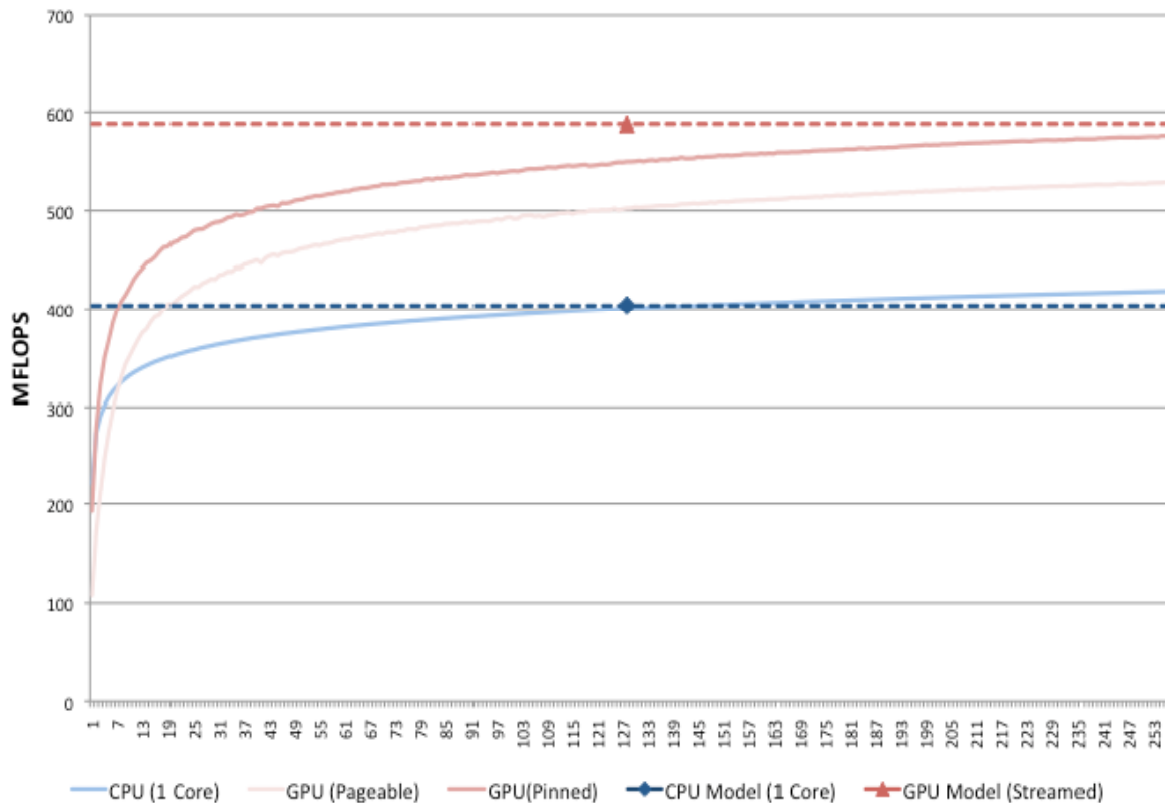
④  Execute & record execution times: $t_i(N_1/p)$

⑤  IF $\max_{1 \le i,j \le p}\{((t_i(N_1/p)-t_j(N_1/p))/t_i(N_1/p)\} \le \varepsilon$
THEN even distribution solves the problem and the algorithm stops;
ELSE performance of devices is calculated, such that:

$$\text{Perf}_i(N_1/p) = nf(N_1/p)/t_i(N_1/p), \quad 1 \le i \le p$$

| Performance Metric | | Initialization | | Approximation | | Iteration |



① Traditional approach: **Performance** of each device is **model**ed as a **constant**

$$\texttt{Perf}_i(\texttt{x}) = \texttt{Perf}_i(\texttt{N}_1/\texttt{p}), \quad 1 \leq i \leq p$$

Legend: CPU (1 Core) — GPU (Pageable) — GPU(Pinned) — CPU Model (1 Core) — GPU Model (Streamed)

# Case Study: Performance Modeling in CPU + GPU Environment (6)

| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | Iteration |

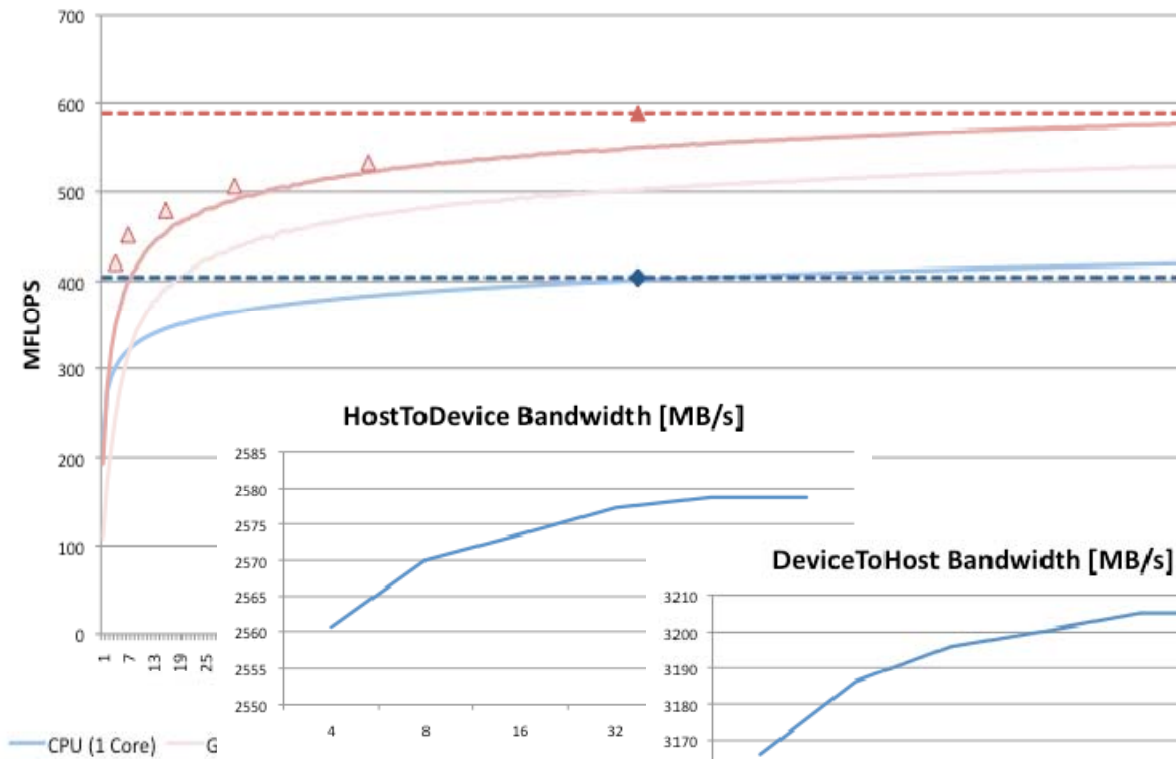

1. Traditional approach: **Performance** of each device is **model**ed as a **constant**

$$\texttt{Perf}_i(\texttt{x}) = \texttt{Perf}_i(\texttt{N}_1/\texttt{p}), \quad 1 \le i \le p$$

② GPU-specific Modeling: Using the obtained values from streaming execution

   – *HostToDevice* Bandwidth

| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | Iteration |
|---|---|---|---|---|---|---|



HostToDevice Bandwidth [MB/s]

DeviceToHost Bandwidth [MB/s]

CPU (1 Core)   GPU (Pageable)   GPU(Pinned)

1 Traditional approach: **Performance** of each device is **model**ed as a **constant**

$$\texttt{Perf}_i(x) = \texttt{Perf}_i(N_1/p), \quad 1 \leq i \leq p$$

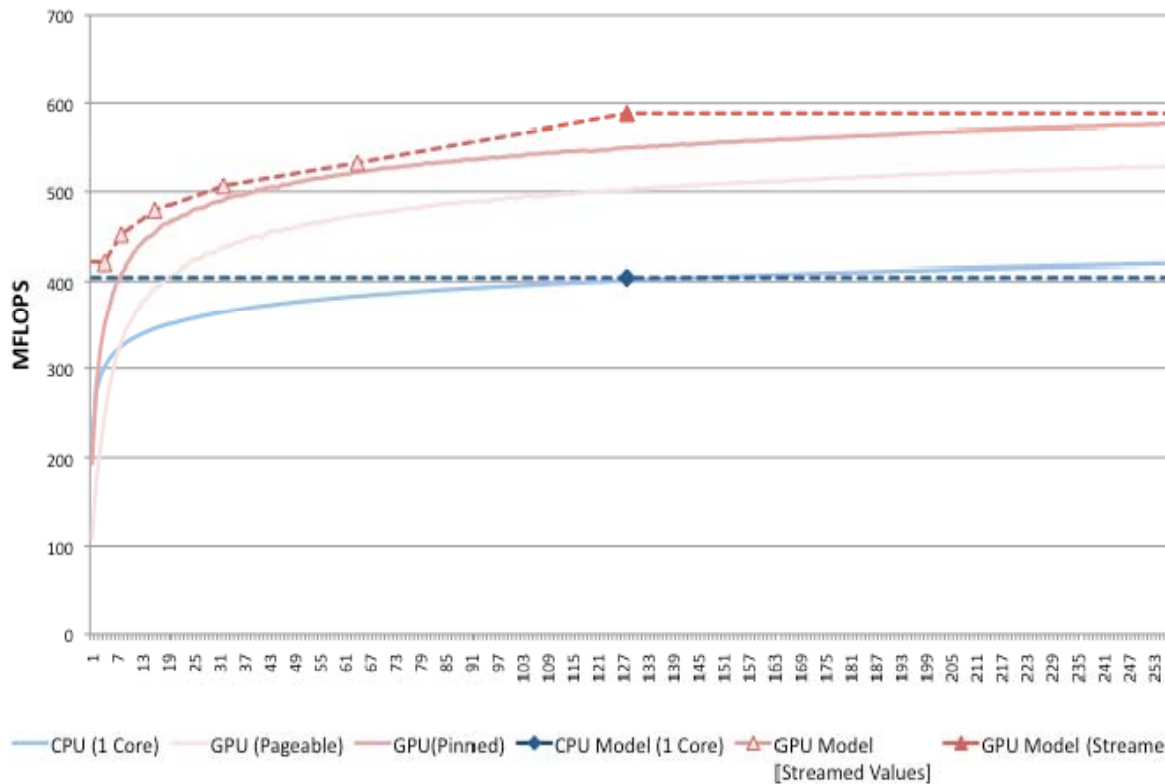② GPU-specific Modeling: Using the obtained values from streaming execution

- *HostToDevice* Bandwidth
- *DeviceToHost* Bandwidth

| Performance Metric | ⟹ | Initialization | ⟹ | **Approximation** | ⟹ | Iteration |
|---|---|---|---|---|---|---|



1 Traditional approach: **Performance** of each device is **model**ed as a **constant**

$$\texttt{Perf}_i(x) = \texttt{Perf}_i(N_1/p),\ 1 \le i \le p$$

② GPU-specific Modeling: Using the obtained values from streaming execution

  – *HostToDevice* Bandwidth
  – *DeviceToHost* Bandwidth
  – GPU *Kernel* Performance

# Case Study: Performance Modeling in CPU + GPU Environment (9)

| Performance Metric | $\Rightarrow$ | Initialization | $\Rightarrow$ | Approximation | $\Rightarrow$ | Iteration |
|---|---|---|---|---|---|---|



1  Traditional approach: **Performance** of each device is **model**ed as a **constant**

$$\texttt{Perf}_i(\texttt{x}) = \texttt{Perf}_i(\texttt{N}_1/\texttt{p}), \ 1 \leq i \leq p$$

② GPU-specific Modeling: Using the obtained values from streaming execution

– *HostToDevice* Bandwidth
– *DeviceToHost* Bandwidth
– GPU *Kernel* Performance

③ Incorporate streaming results

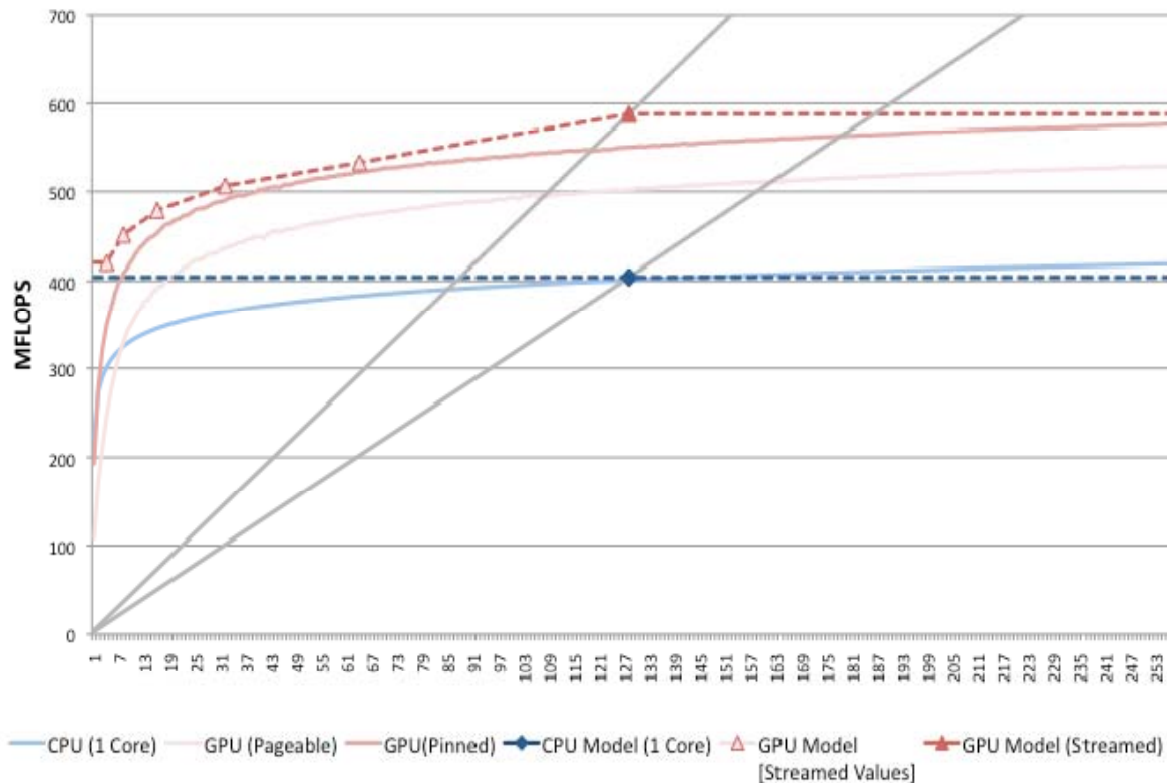# Case Study: Performance Modeling in CPU + GPU Environment (10)

| Performance Metric | → | Initialization | → | **Approximation** | → | Iteration |



CPU (1 Core) — GPU (Pageable) — GPU(Pinned) — CPU Model (1 Core) — GPU Model [Streamed Values] — GPU Model (Streamed)

1. Traditional approach: **Performance** of each device is **model**ed as a **constant**

   $$\mathtt{Perf_i(x) = Perf_i(N_1/p),\ 1{\le}i{\le}p}$$

2. GPU-specific Modeling: Using the obtained values from streaming execution

   - *HostToDevice* Bandwidth
   - *DeviceToHost* Bandwidth
   - GPU *Kernel* Performance

3. Incorporate streaming results

# Case Study: Performance Modeling in CPU + GPU Environment (11)
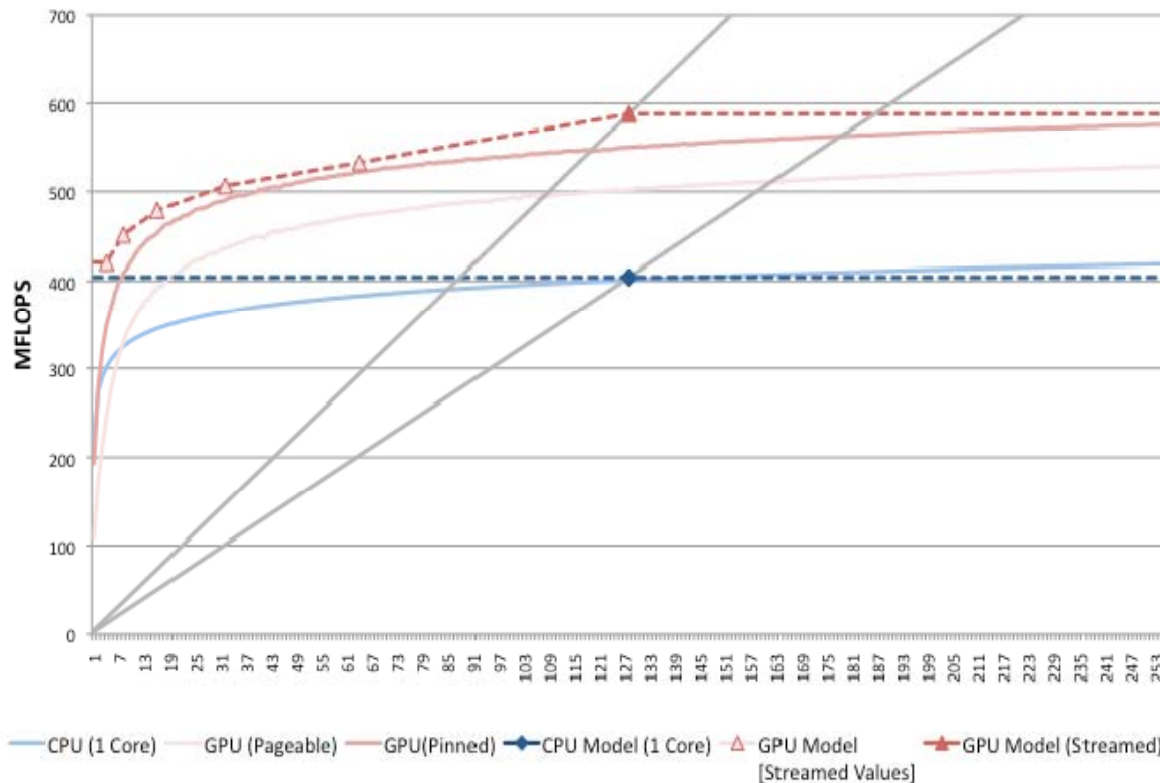
| Performance Metric | → | Initialization | → | Approximation | → | Iteration |



① Draw Upper $U$ and Lower $L$ lines through the following points:

$(0,0)$, $(N_1/p, \max_i\{Perf_i(N_1/p)\})$

$(0,0)$, $(N_1/p, \min_i\{Perf_i(N_1/p)\})$

② Let $x_i^{(U)}$ and $x_i^{(L)}$ be the intersections with $Perf_i(x)$

IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$

THEN go to 3

ELSE go to 5

③ Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$

④ IF $\Sigma_i x_i^{(M)} \leq N_1$

THEN $U=M$

ELSE $L=M$

REPEAT 2

Legend: CPU (1 Core) · GPU (Pageable) · GPU(Pinned) · CPU Model (1 Core) · GPU Model [Streamed Values] · GPU Model (Streamed)

| Performance Metric | ⟹ | Initialization | ⟹ | Approximation | ⟹ | Iteration |



1 Draw Upper $U$ and Lower $L$ lines through the following points:

$(0,0)$, $(N_1/p$, $\max_i\{Perf_i(N_1/p)\})$

$(0,0)$, $(N_1/p$, $\min_i\{Perf_i(N_1/p)\})$

2 Let $x_i^{(U)}$ and $x_i^{(L)}$ be the intersections with $Perf_i(x)$

IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$

THEN go to 3

ELSE go to 5

3 Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$

4 IF $\Sigma_i x_i^{(M)} \leq N_1$

THEN $U = M$

ELSE $L = M$

REPEAT 2

⑤ Employ **streaming strategy** on the calculated workload value

Performance Metric $\Longrightarrow$ Initialization $\Longrightarrow$ Approximation $\Longrightarrow$ Iteration / Streaming

- STREAMING STRATEGY
  - Results obtained using DIV2 STRATEGY consider application characterization (e.g. communication-to-computation ratio)
  - Workload size for the next stream should be chosen in order to OVERLAP TRANSFERS WITH COMPUTATION in the previous stream

- BANDWIDTH-AWARE STREAMING STRATEGY
  - Reuses the MINIMAL WORKLOAD SIZE FROM DIV2 STRATEGY (obtained via HOST-TO-DEVICE and DEVICE-TO-HOST tests)

    ```
    IF (n^curr ≥ n^min_size)
    THEN use strategy (cont. overlapping)
    ELSE restart strategy on n^curr load
    ```

  - In case that load drops under $n^{min\_size}$, strategy is restarted on remaining load

1. Draw Upper $U$ and Lower $L$ lines through the following points:

   $(0,0)$, $(N_1/p$, $\max_i\{Perf_i(N_1/p)\})$

   $(0,0)$, $(N_1/p$, $\min_i\{Perf_i(N_1/p)\})$

2. Let $x_i^{(U)}$ and $x_i^{(L)}$ be the intersections with $Perf_i(x)$
   IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$
   THEN go to 3
   ELSE go to 5

3. Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$

4. IF $\sum_i x_i^{(M)} \leq N_1$
   THEN $U=M$
   ELSE $L=M$
   REPEAT 2

⑤ Employ **streaming strategy** on the calculated workload value

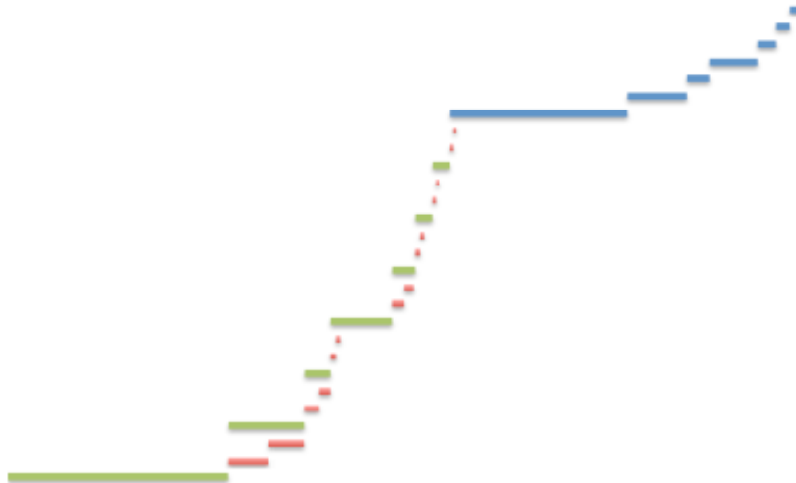| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | **Iteration** |
|---|---|---|---|---|---|---|
| | | | | | | **Streaming** |

– CASE STUDY: **2D FFT BATCH**
  – About 83% of total execution time goes on data transfers
    – **HOSTTODEVICE** Transfers – 46%
    – **KERNEL** Execution – 17%
    – **DEVICETOHOST** Transfers – 37%

– BANDWIDTH-AWARE STREAMING STRATEGY
  – $n^{min\_size} = 4$, overlap HOSTTODEVICE transfers and KERNEL execution of two streams



1  Draw Upper $U$ and Lower $L$ lines through the following points:

$(0,0)$, $(N_1/p, \max_i\{Perf_i(N_1/p)\})$

$(0,0)$, $(N_1/p, \min_i\{Perf_i(N_1/p)\})$

2  Let $x_i^{(U)}$ and $x_i^{(L)}$ be the intersections with $Perf_i(x)$
IF  exists $x_i^{(L)} - x_i^{(U)} \geq 1$
THEN  go to 3
ELSE  go to 5

3  Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$
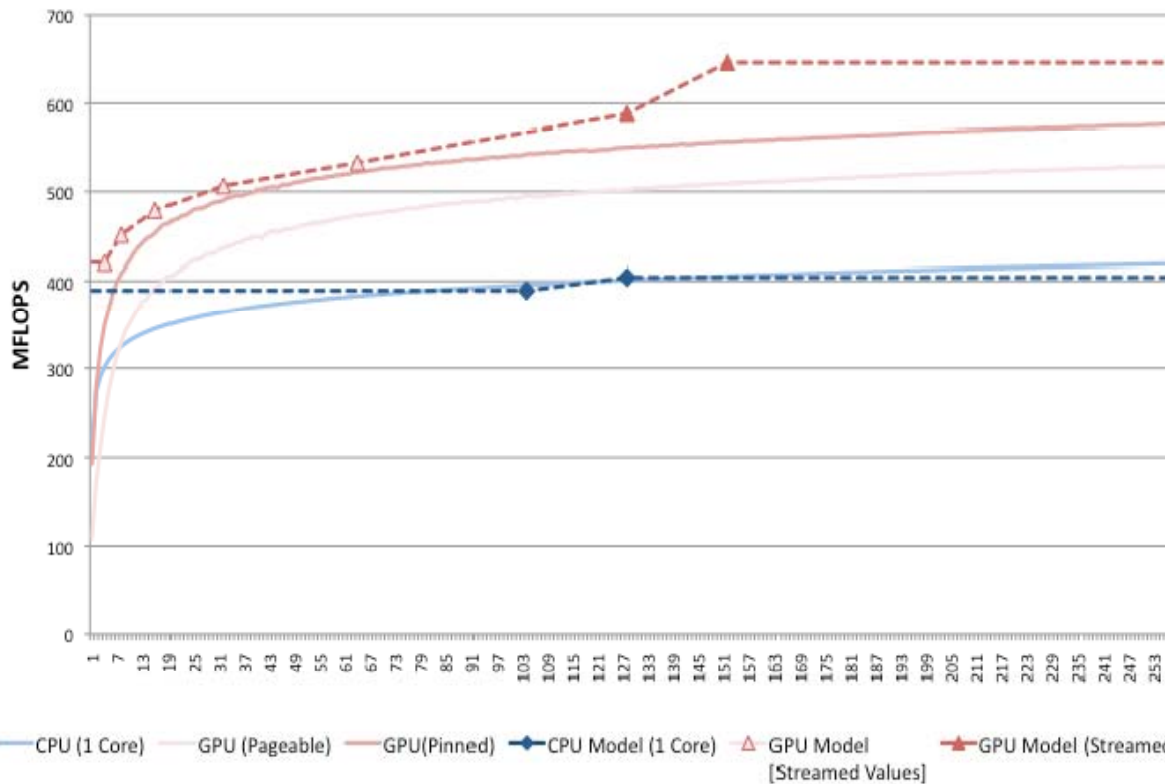
4  IF  $\Sigma_i x_i^{(M)} \leq N_1$
THEN  $U = M$
ELSE  $L = M$
REPEAT 2

⑤ Employ **streaming strategy** on the calculated workload value

| Performance Metric | ⟹ | Initialization | ⟹ | Approximation | ⟹ | Iteration |



1 Draw Upper $U$ and Lower $L$ lines through the following points:

$(0,0)$, $(N_1/p, \max_i\{Perf_i(N_1/p)\})$

$(0,0)$, $(N_1/p, \min_i\{Perf_i(N_1/p)\})$

2 Let $x_i^{(U)}$ and $x_i^{(L)}$ be the intersections with $Perf_i(x)$
IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$
THEN go to 3
ELSE go to 5

3 Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$

4 IF $\Sigma_i\ x_i^{(M)} \leq N_1$
THEN $U=M$
ELSE $L=M$
REPEAT 2

5 Employ streaming strategy on the calculated workload value
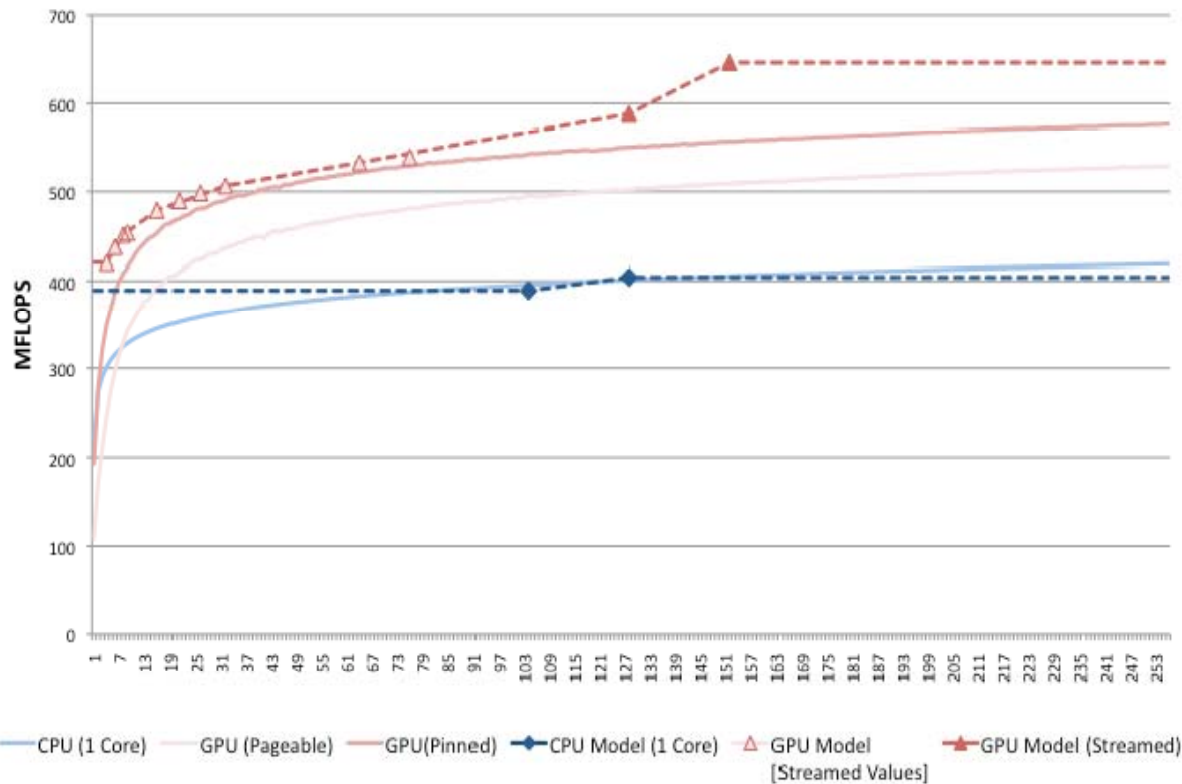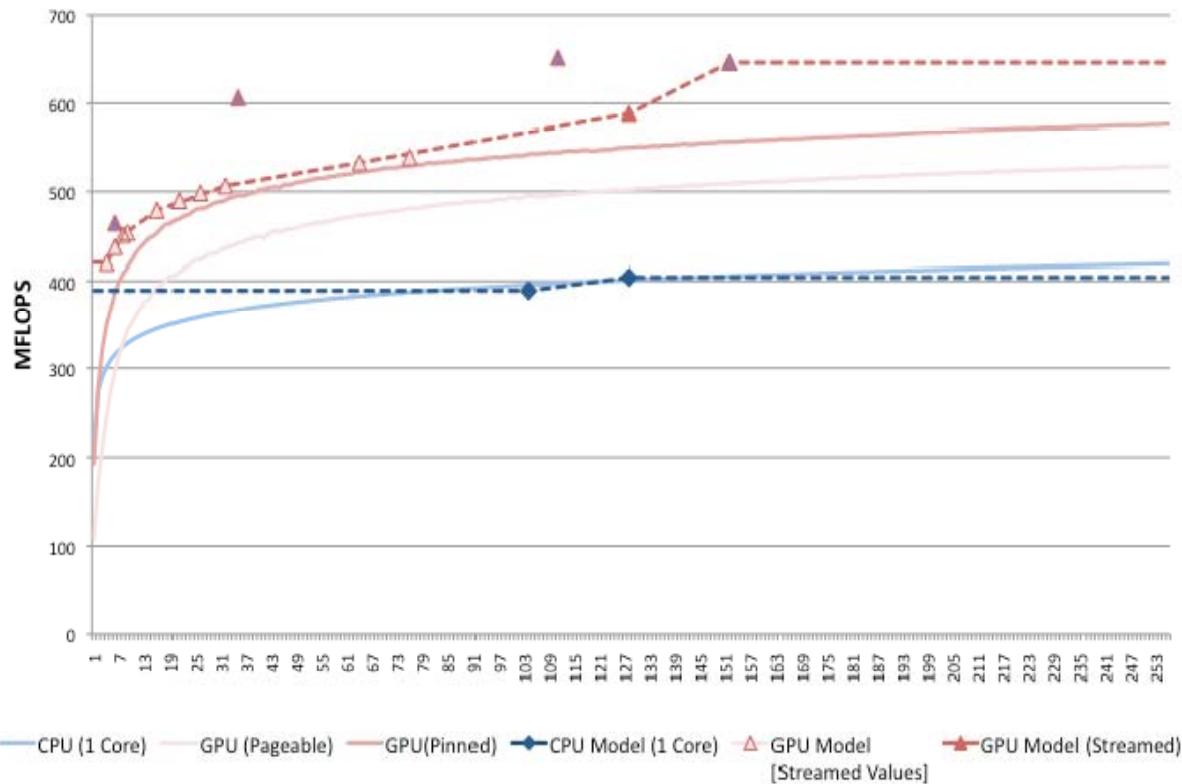
Performance Metric ⟹ Initialization ⟹ **Approximation** ⟹ Iteration



① Refine performance models with the newly obtained results

2 GPU-specific Modeling: Using the obtained values from streaming execution

- *HostToDevice* Bandwidth
- *DeviceToHost* Bandwidth
- GPU *Kernel* Performance

3 Incorporate streaming results
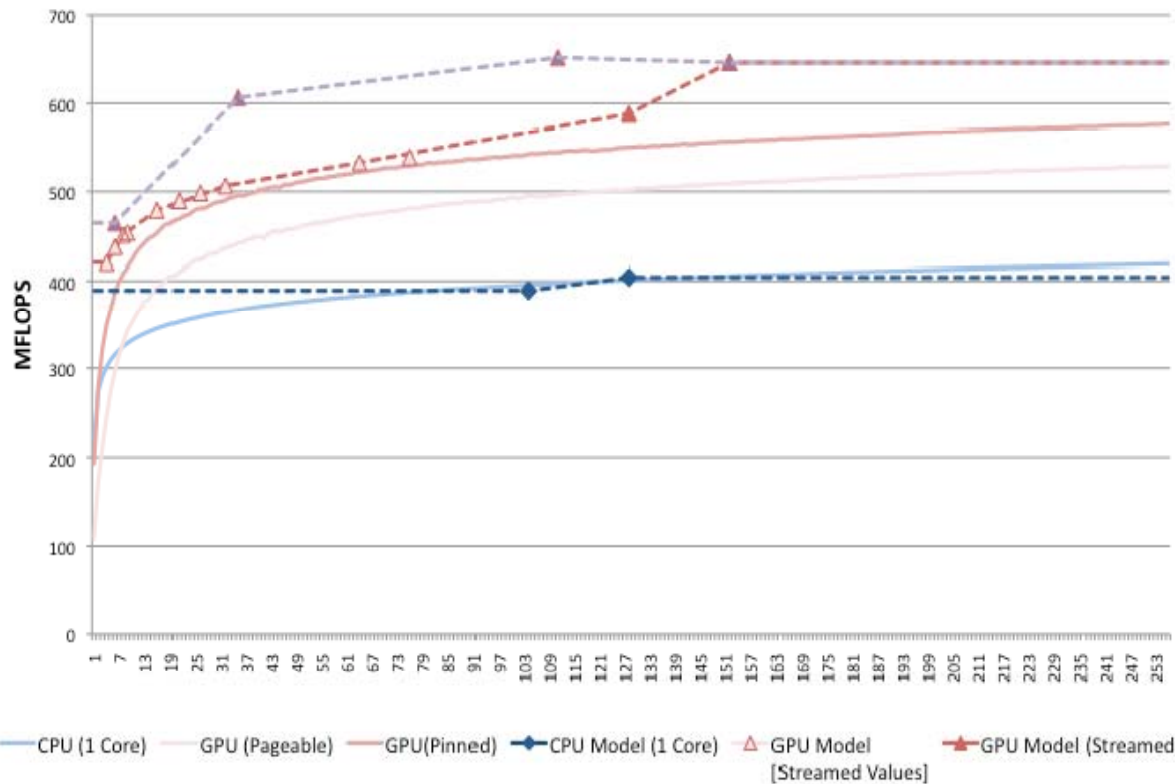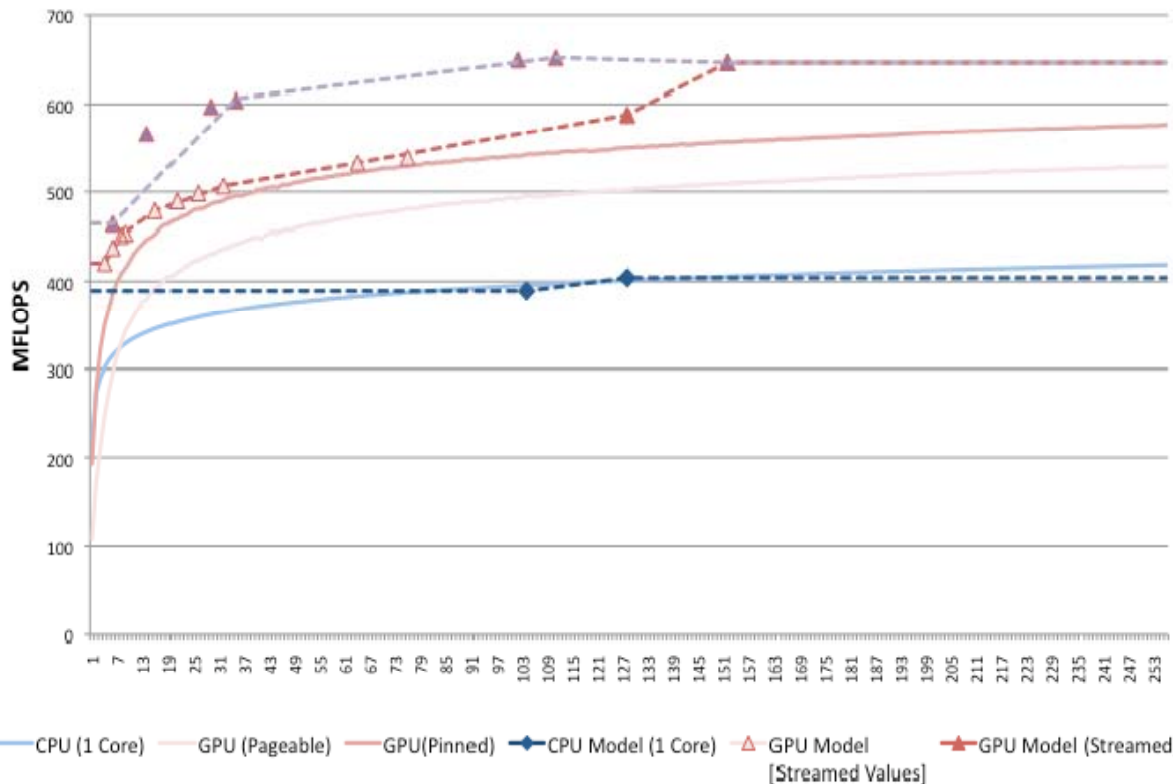
Performance Metric ⟹ Initialization ⟹ **Approximation** ⟹ Iteration



1  Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution

– *HostToDevice* Bandwidth
– *DeviceToHost* Bandwidth
– GPU *Kernel* Performance

③ Incorporate streaming results

– for each stream

# Case Study: Performance Modeling in CPU + GPU Environment (18)

| Performance Metric | ⇒ | Initialization | ⇒ | **Approximation** | ⇒ | Iteration |



1  Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution

  – *HostToDevice* Bandwidth
  – *DeviceToHost* Bandwidth
  – GPU *Kernel* Performance

③ Incorporate streaming results

  – for each stream
  – for each stream restart

# Case Study: Performance Modeling in CPU + GPU Environment (19)

| Performance Metric | ⇒ | Initialization | ⇒ | **Approximation** | ⇒ | Iteration |



1  Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution
   - *HostToDevice* Bandwidth
   - *DeviceToHost* Bandwidth
   - GPU *Kernel* Performance

③ Incorporate streaming results
   - for each stream
   - for each stream restart

# Case Study: Performance Modeling in CPU + GPU Environment (20)
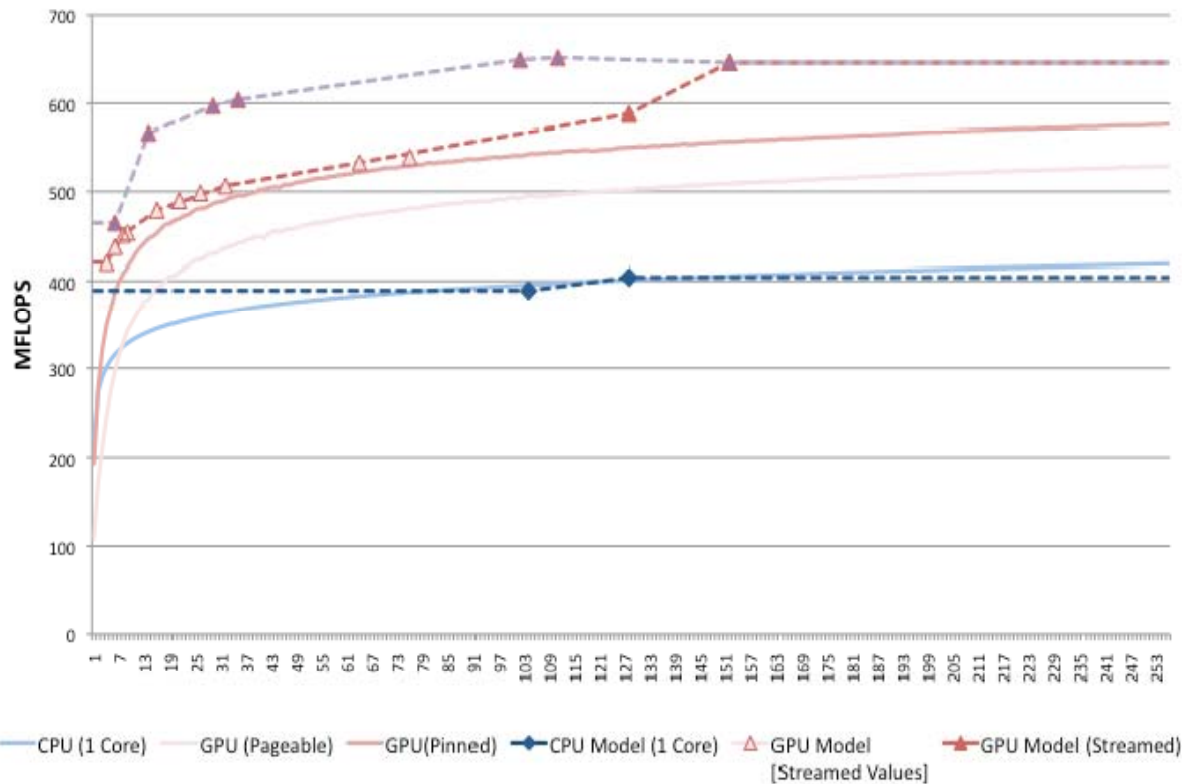
Performance Metric ⟹ Initialization ⟹ **Approximation** ⟹ Iteration



1 Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution

- *HostToDevice* Bandwidth
- *DeviceToHost* Bandwidth
- GPU *Kernel* Performance

③ Incorporate streaming results

- for each stream
- for each stream restart
- for every stream combination

| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | Iteration |
|---|---|---|---|---|---|---|



1   Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution

  – *HostToDevice* Bandwidth
  – *DeviceToHost* Bandwidth
  – GPU *Kernel* Performance

③ Incorporate streaming results

  – for each stream
  – for each stream restart
  – for every stream combination

Legend: CPU (1 Core) — GPU (Pageable) — GPU(Pinned) — CPU Model (1 Core) — GPU Model [Streamed Values] — GPU Model (Streamed)

| Performance Metric | $\Rightarrow$ | Initialization | $\Rightarrow$ | Approximation | $\Rightarrow$ | Iteration |



① Draw Upper $U$ and Lower $L$ lines through the following points:

$(0,0)$, $(N_1/p$, $\max_i\{Perf_i(N_1/p)\})$

$(0,0)$, $(N_1/p$, $\min_i\{Perf_i(N_1/p)\})$

② Let $x_i^{(U)}$ and $x_i^{(L)}$ be the intersections with $Perf_i(x)$

IF exists $x_i^{(L)} - x_i^{(U)} \geq 1$

THEN go to 3

ELSE go to 5

③ Bisect the angle between $U$ and $L$ by the line $M$, and calculate intersections $x_i^{(M)}$

④ IF $\sum_i x_i^{(M)} \leq N_1$

THEN $U=M$

ELSE $L=M$

REPEAT 2

⑤ Employ streaming strategy on the calculated workload value
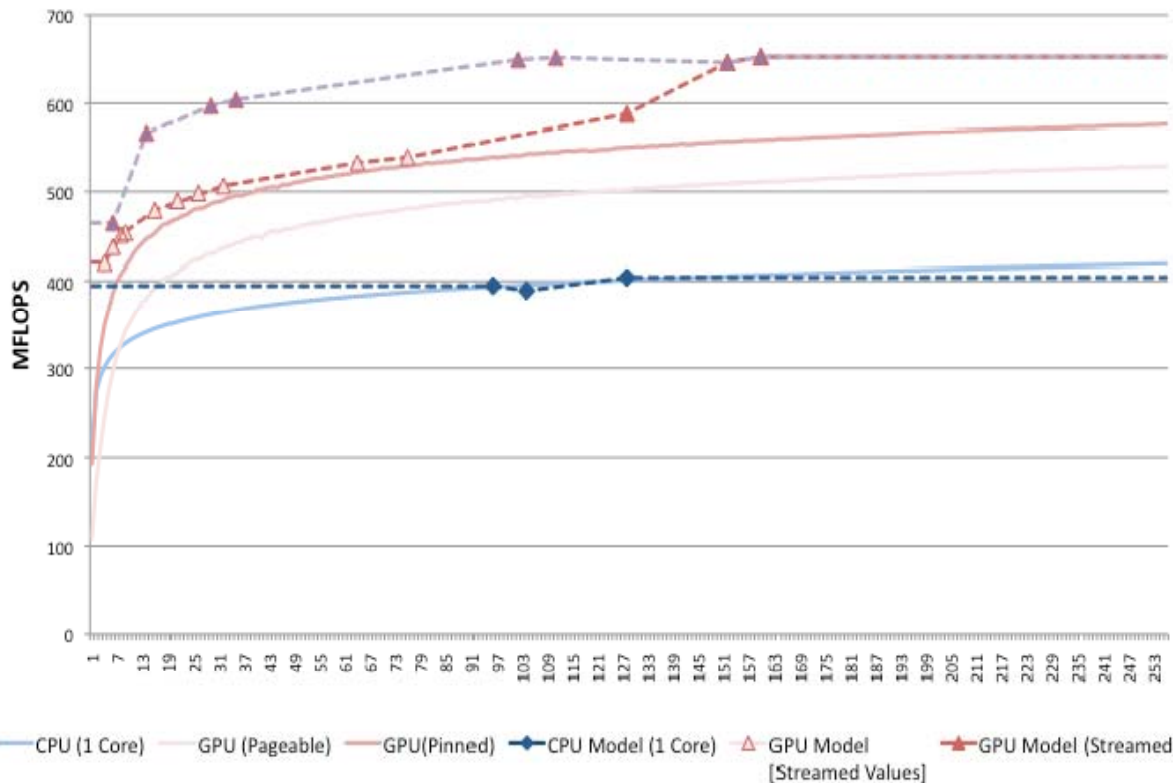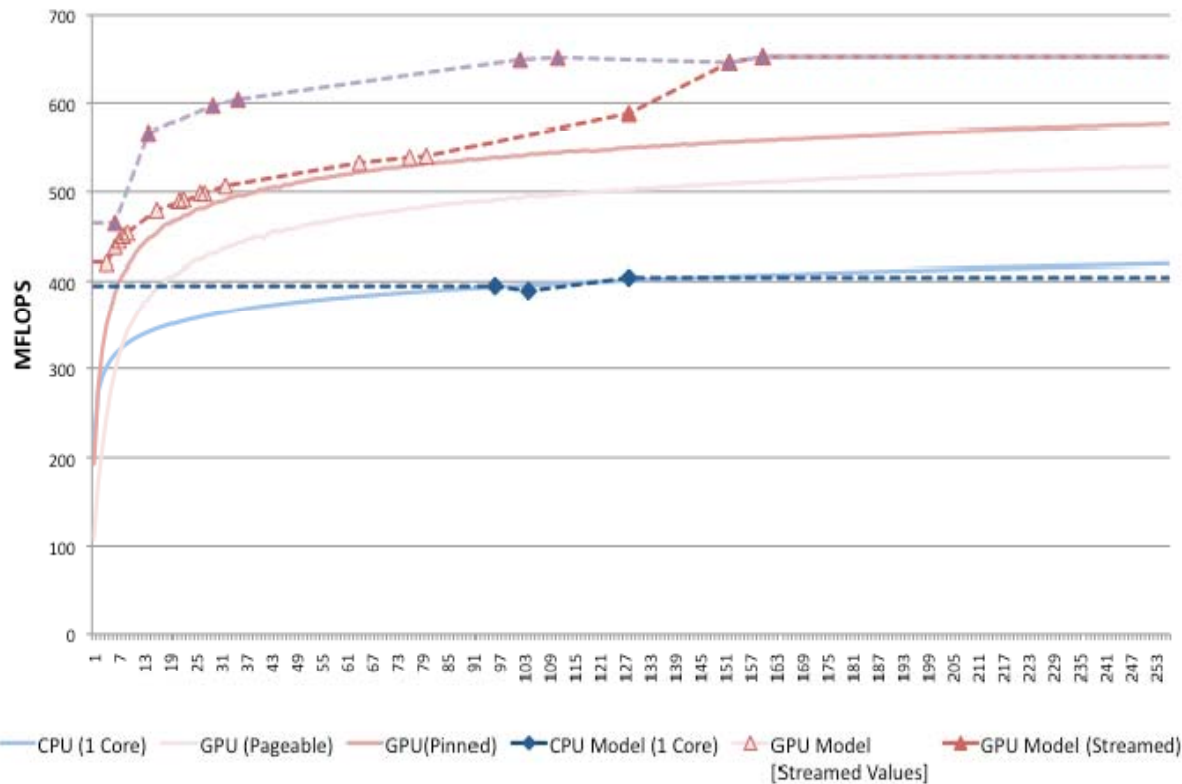
Performance Metric ⟹ Initialization ⟹ Approximation ⟹ Iteration



① Refine performance models with the newly obtained results

2   GPU-specific Modeling: Using the obtained values from streaming execution

– *HostToDevice* Bandwidth
– *DeviceToHost* Bandwidth
– GPU *Kernel* Performance

3   Incorporate streaming results

– for each stream
– for each stream restart
– for every stream combination

Legend: CPU (1 Core) — GPU (Pageable) — GPU(Pinned) — CPU Model (1 Core) — GPU Model [Streamed Values] — GPU Model (Streamed)

# Case Study: Performance Modeling in CPU + GPU Environment (24)

| Performance Metric | Initialization | **Approximation** | Iteration |



1 Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution

 – *HostToDevice* Bandwidth
 – *DeviceToHost* Bandwidth
 – GPU *Kernel* Performance

③ Incorporate streaming results

 – for each stream
 – for each stream restart
 – for every stream combination
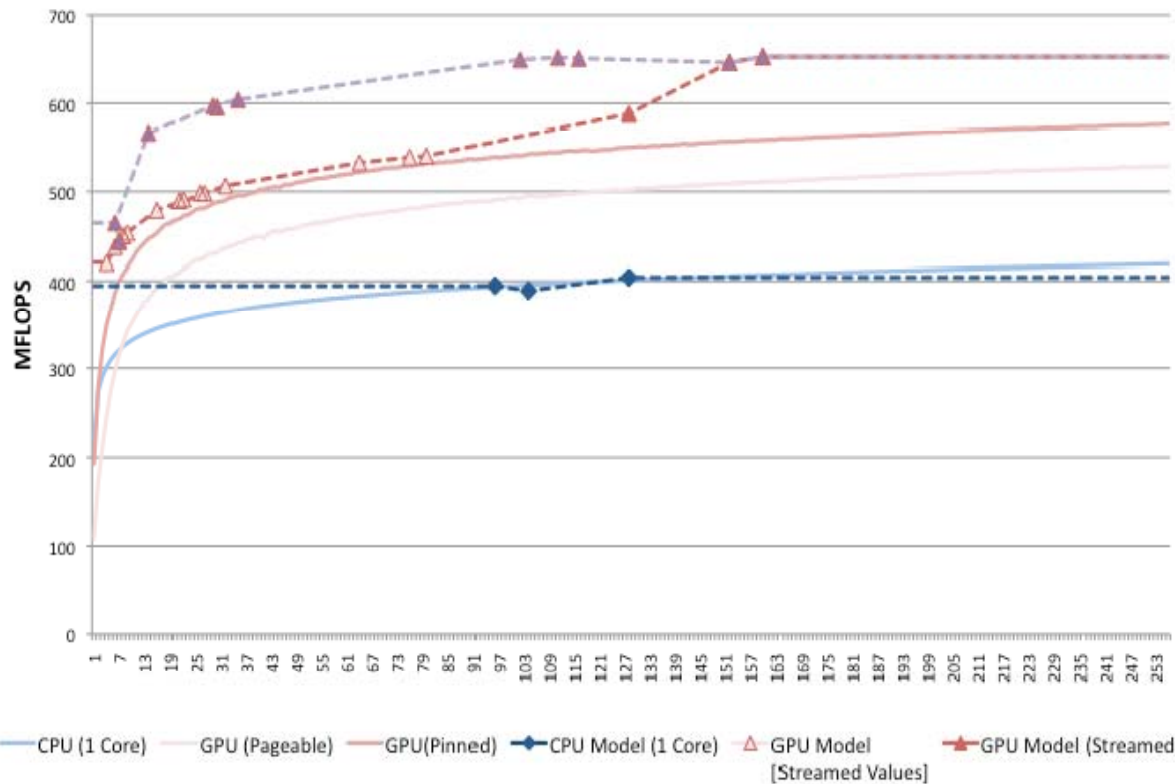
| Performance Metric | ⇒ | Initialization | ⇒ | Approximation | ⇒ | Iteration |



1  Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution

– *HostToDevice* Bandwidth
– *DeviceToHost* Bandwidth
– GPU *Kernel* Performance

③ Incorporate streaming results

– for each stream
– for each stream restart
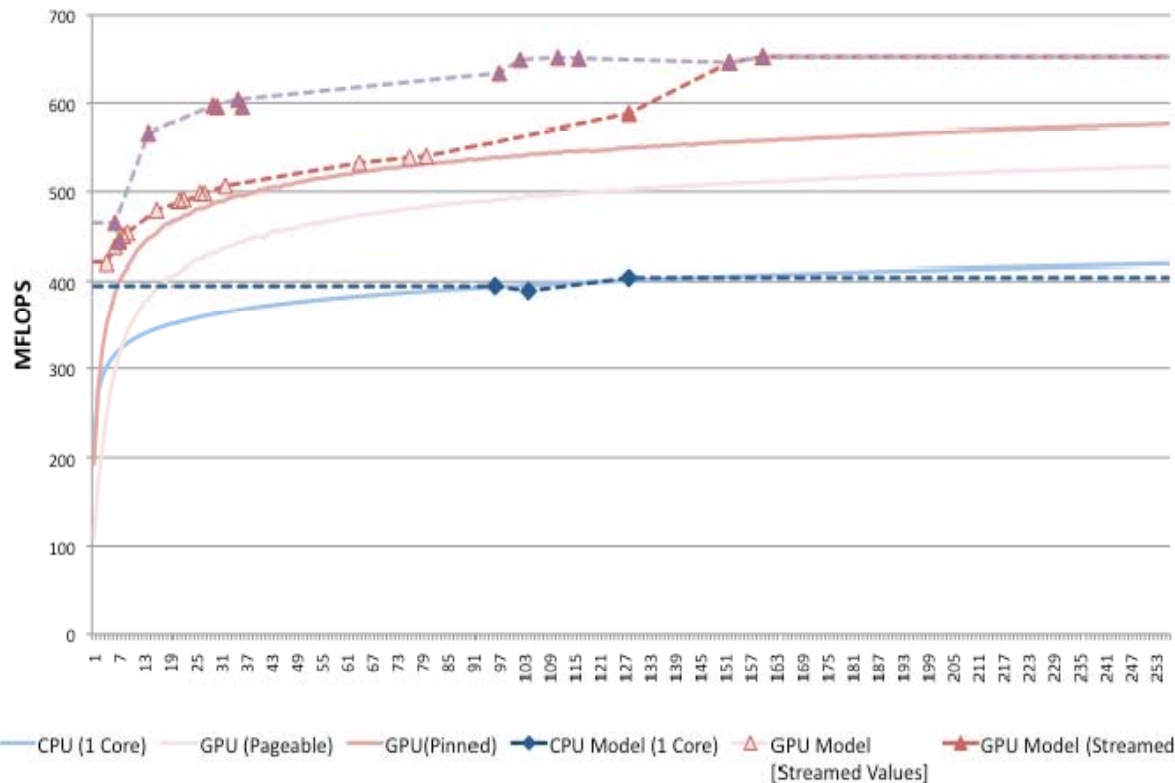– for every stream combination

Performance Metric ⟹ Initialization ⟹ **Approximation** ⟹ Iteration



1  Refine performance models with the newly obtained results

② GPU-specific Modeling: Using the obtained values from streaming execution
   – *HostToDevice* Bandwidth
   – *DeviceToHost* Bandwidth
   – GPU *Kernel* Performance

③ Incorporate streaming results
   – for each stream
   – for each stream restart
   – for every stream combination

# COLLABORATIVE ENVIRONMENT FOR HETEROGENEOUS COMPUTERS

## TRADITIONAL APPROACHES FOR PERFORMANCE MODELING

– Approximate the performance using number of points equal to the number of iterations
– In this case, **3 POINTS** per each device

## PRESENTED APPROACH FOR PERFORMANCE MODELING

– Models the performance using **MORE THAN 30 POINTS**, in this case
– **COMMUNICATION-AWARE** – schedules in respect to limited and asymmetric interconnection bandwidth
– Employs **STREAMING STRATEGIES** to overlap communication with computation across devices
– **BUILDS SEVERAL** PER-DEVICE **MODELS** AT THE SAME TIME
  – OVERALL PERFORMANCE for each device + STREAMING GPU PERFORMANCE
  – HOSTTODEVICE BANDWIDTH Modeling
  – DEVICETOHOST BANDWIDTH Modeling
  – GPU KERNEL PERFORMANCE Modeling