technology
from seed

# Stream-based concurrent computational models and programming tools

## Leonel Sousa

with

Shinichi Yamagiwa

inesc id lisboa

INSTITUTO SUPERIOR TÉCNICO

1. **Many-core platforms based on GPU's**

2. **GPGPU: Computation Models and Programming tools**
   1. Stream based computing
   2. Massively parallelism based on Multithreading
   3. APIs and Programming tools

3. **Caravela Project**
   1. Flow-Model and Caravela Platform
   2. Caravela Tools for programming GPUs (locally and remotely)
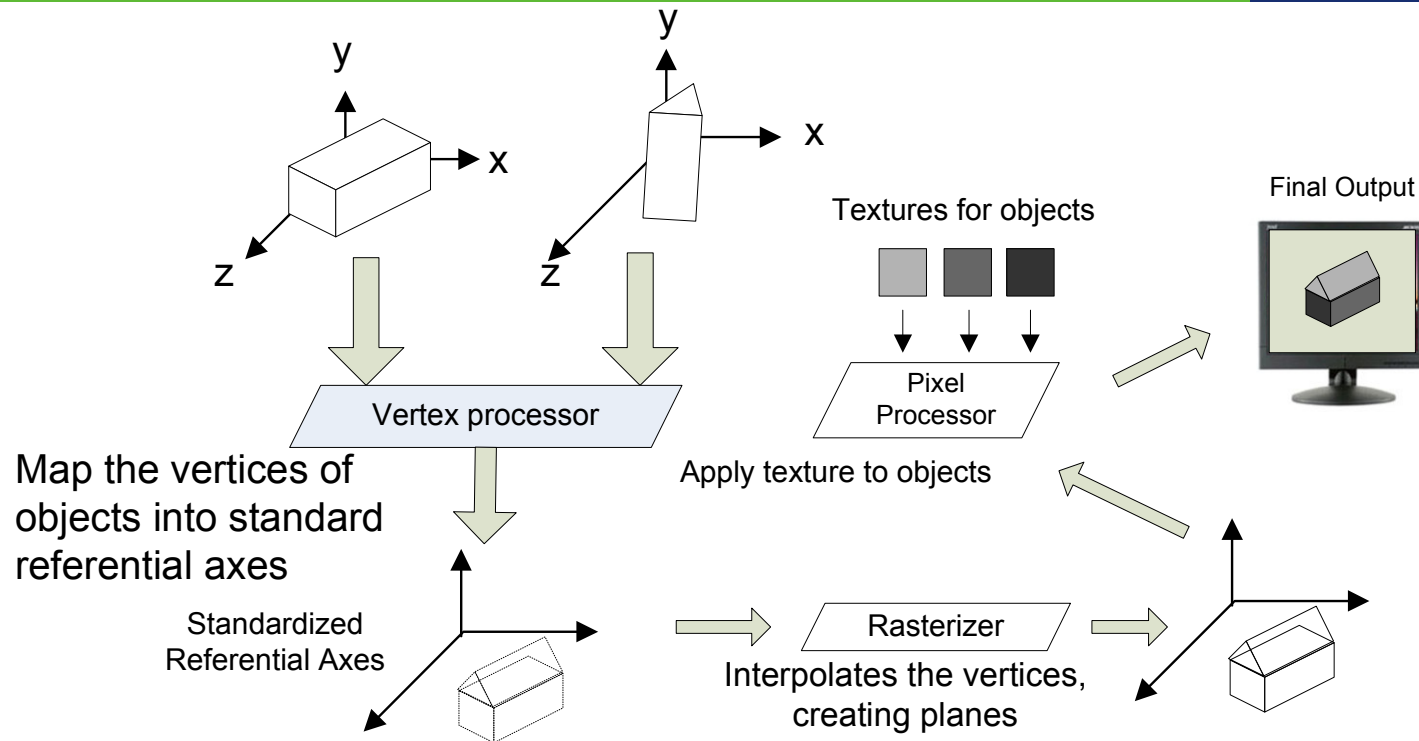   3. Optimizations for current GPUs/Systems
   4. Future Work

technology
from seed

- Graphics Processing Units (GPUs)

  – Available in all computers
  – Unused high computational capacity
  – **Manycore** processing systems

## GPGPU - General-Purpose computation on GPUs

- Usage of GPUs for GPGPU

  – Graphics APIs are not tuned for general-purpose applications
  – Programmer has to learn irrelevant graphics concepts
  – Data copy from main memory to video memory is slow
    - PCI-E  system bus

technology
from seed

inesc id
lisboa

y

y

x

x

Final Output

Textures for objects

z

z

Vertex processor

Pixel
Processor

Map the vertices of
objects into standard
referential axes

Apply texture to objects

Standardized
Referential Axes

Rasterizer

Interpolates the vertices,
creating planes

- Stream-based processing with four elements
  - Vertex processor: x, y, z, w
  - Pixel processor: operates on pixel data in a vector approach , issuing instructions to operate concurrently on the multiple color components of a pixel -R(ed), G(reen), B(lue) and A(lpha)
- Vertex and Pixel processors are programmable
  - DirectX assembly language and HLSL
  - OpenGL Shading Language (GLSL)

University of Murcia

INSTITUTO
SUPERIOR
TÉCNICO

```
ps_2_0        ←——— DirectX assembly language
                    Pixel Shader Model 2.0
def c0, 0.5,0.5,0.5,0  ←——— α
def c1, 1,1,1,1

dcl_2d s0
dcl_2d s1

dcl t0.xy     ←———— Coordinates of textures
dcl t1.xy

texld r2, t0, s0  ←———  Da Vinci
texld r3, t1, s1  ←———  Mona Lisa

mov r5, c1
sub r5, r5, c0
mul r2, r2, r5   ←———  Pa(1-α)
mad r4, r3, c0, r2  ←—— Pa(1-α)+Pbα

mov oC0, r4   ←——— Output of results
```
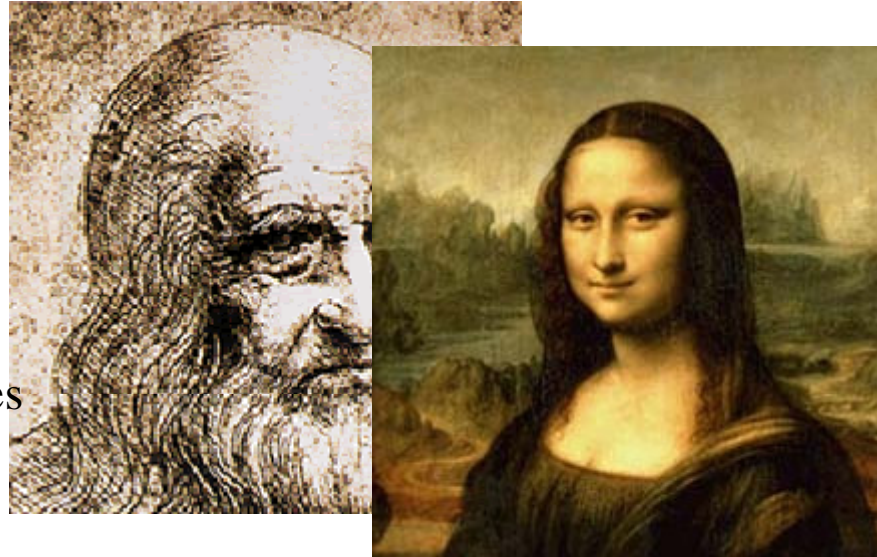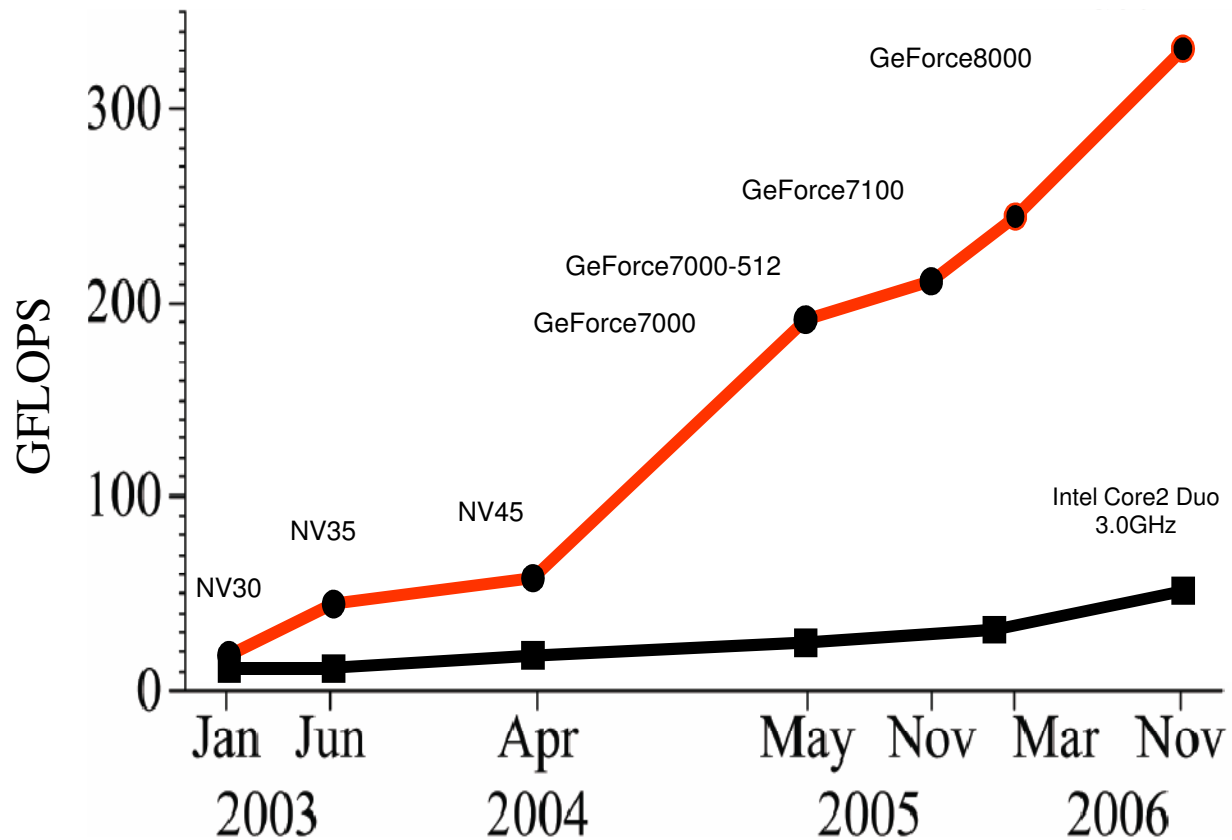
$$P' = Pa(1-\alpha)+Pb\alpha$$
Alpha blending
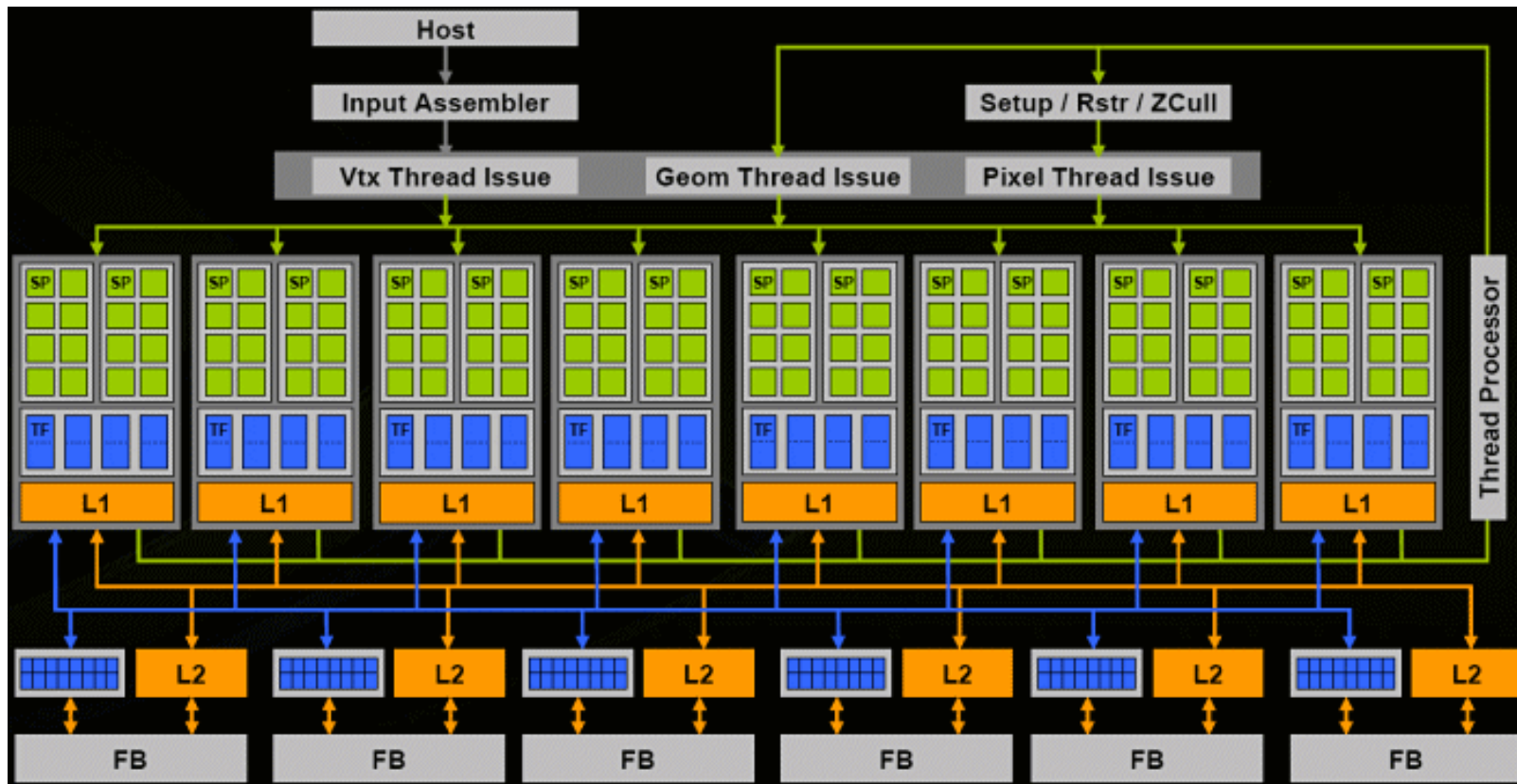
- **GPU supports general purpose processing (data-parallelism)**
    - **with high number of arithmetic calculations per memory access**

- **Examples (**www.**gpgpu**.org**)**
    - **Physics simulation**
    - **Signal processing**
    - **Computational geometry**
    - **Database management**
    - **Computational biology**
    - **Computational finance**
    - **Computer vision**
    - **…..**

- GPU drastically improves performance in the last 5 years

GeForce 8800  [source: NVIDIA]

## GPU architecture
## Case study: GeForce 8800
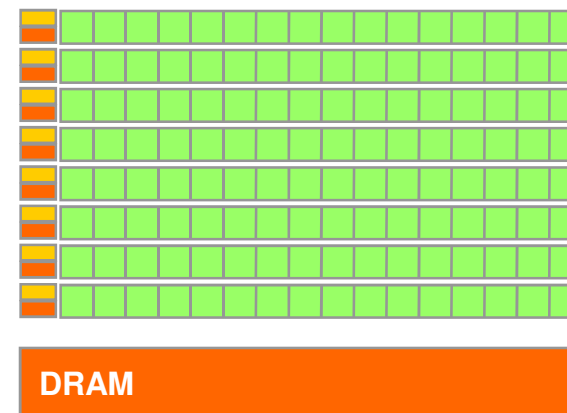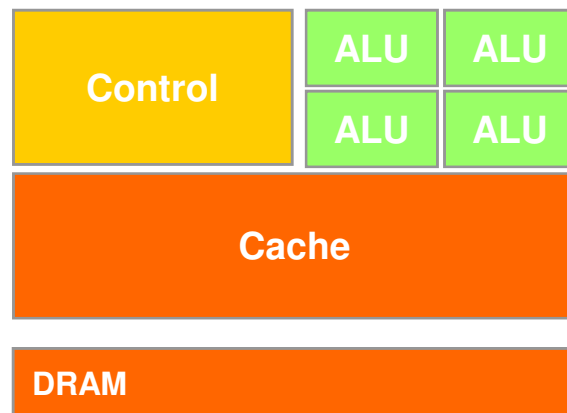
technology
from seed

inesc id
lisboa

330 Gflop/s (issue rate for MAC), 86.4 GB/s peak mem. bandwidth

- 128 stream processors: 8 clusters of 16 SPs
- SPs aren't vertex or pixel shaders: generalized floating-point processors capable of operating on vertices, pixels, or any data
  - most GPUs operate on pixel data in a way (R,G,B,A) but the G80's SP is scalar
- SPs are clocked at a relatively speedy 1.35GHz, while most of the rest of the chip is clocked independently at 575MHz
  - GeForce 8800: a tremendous amount of raw floating-point processing power
- The cores in a cluster share:
  - local memory (L1)
  - banks of specialized hardware (TF) for implementing texture fetch operations
- High performance access to the frame buffer memory (FB)
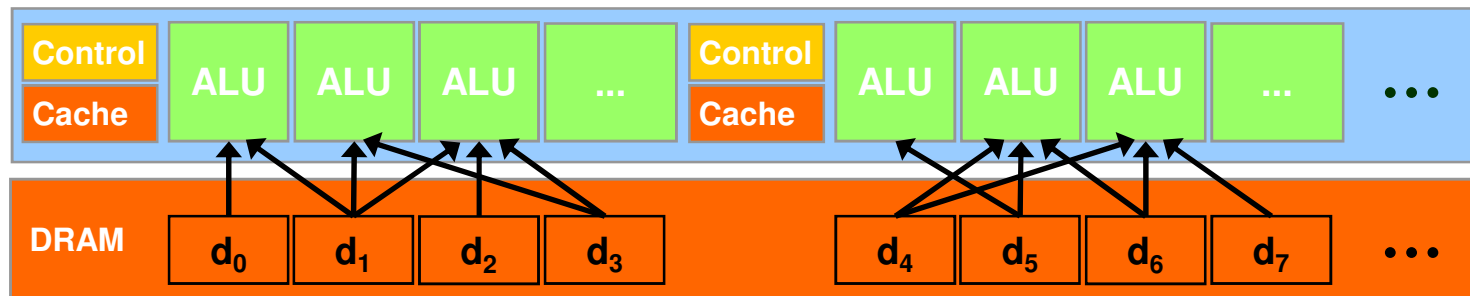  - to store both texture data and rendered images

- Input data is streamed in from one or more input arrays, processed by a **stream kernel**, and then streamed out to one or more output arrays

- A **stream kernel** can be thought of as:
  - function that is applied in parallel to every element of one or more input arrays and produces one or more output arrays
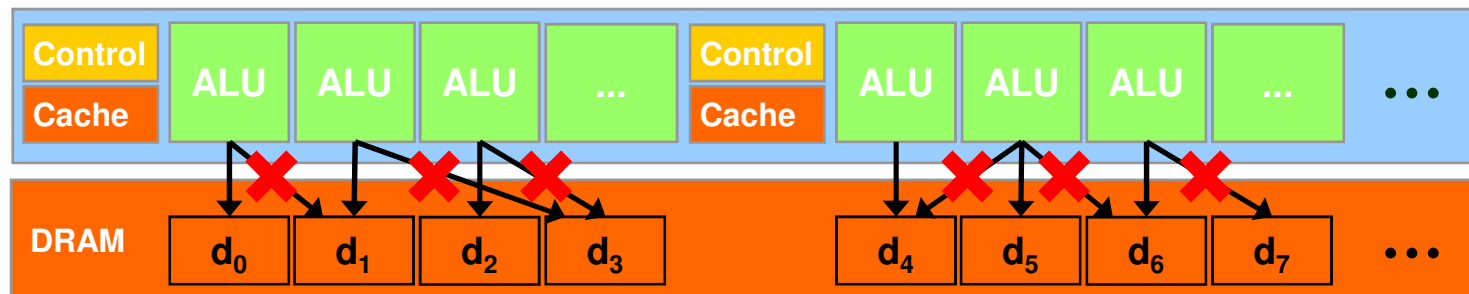
| Control | ALU | ALU |
|---------|-----|-----|
|         | ALU | ALU |

**Cache**

**DRAM**

**DRAM**

- ## Applications can easily be limited by memory bandwidth
  - Restrictions: memory accesses oriented to pixel processing
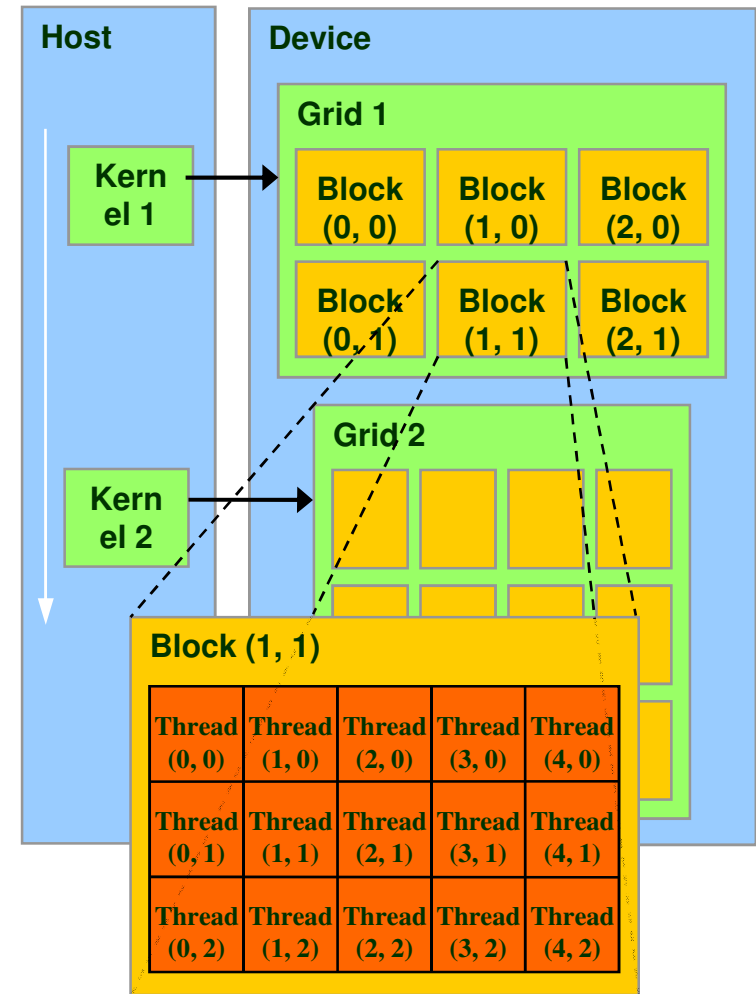  - Only gather: can read data from other pixels
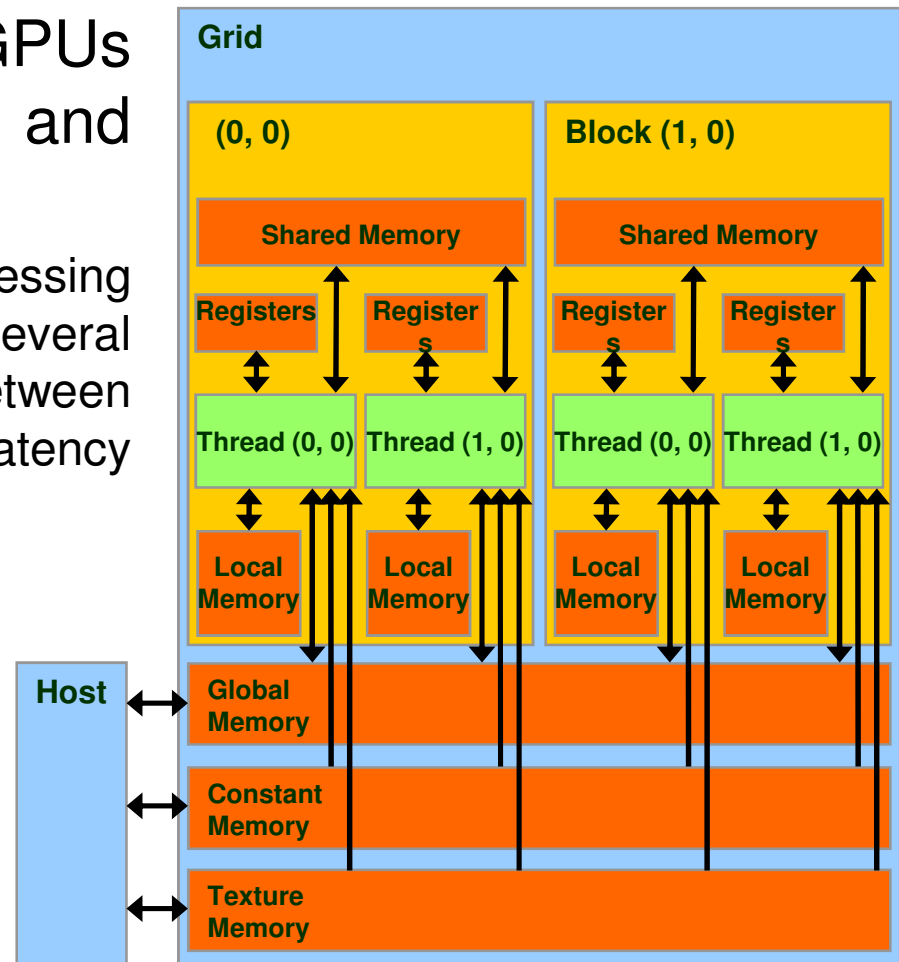


  - No scatter: (Can only write to one pixel)

**Computation Models:**
**Multithreading**

technology
from seed

inesc id
lisboa

**Compute Unified Device Architecture (CUDA): NVIDIA proprietary**

- **SPMD + SIMD** Model (SIMT)
  - Data-parallel portions of an application are executed as **kernels** which run in parallel on many threads

- A kernel is executed as a grid of thread blocks
  - A thread block is a batch of threads that can cooperate with each other through shared memory

- Two threads from two different blocks cannot cooperate

- ## Massive parallelism for GPUs to hide memory access and pipeline latencies

  - For instance, a single processing element in a GPU might run several threads at once and switch between them whenever a high-latency operation is encountered.

- ## Read/write per-thread

  - registers, local memory

- ## Read/write per-block

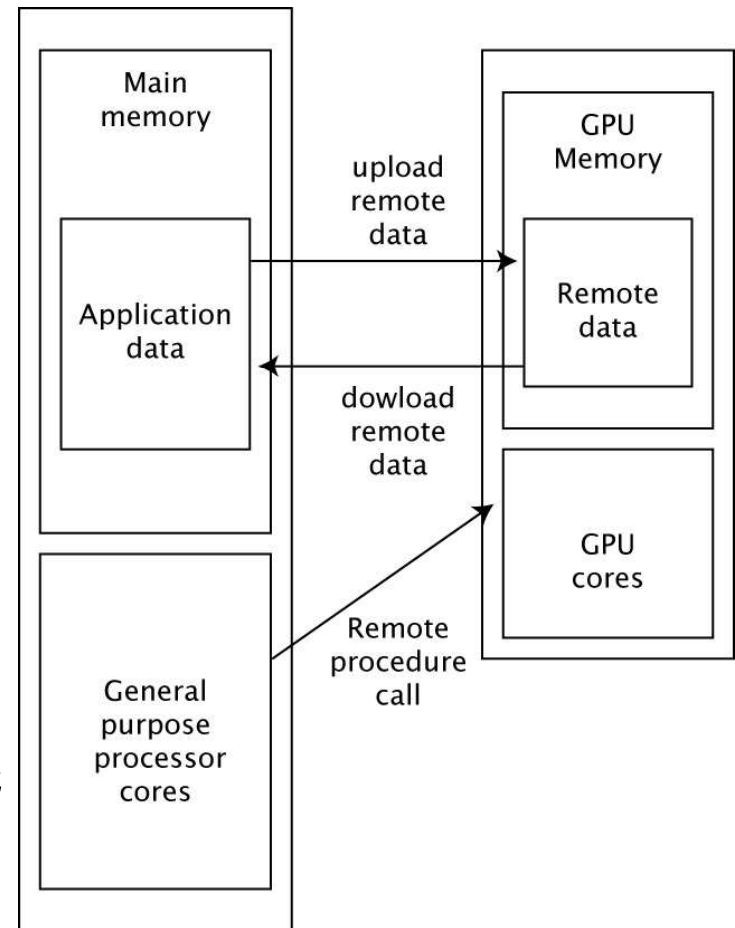  - shared memory

- ## Read/write per-grid

  - global memory

- ## CUDA API is an extension to the C language
  - extensions to target portions of the code for execution on the device
  - a runtime library split into
    - a common component providing built-in vector types and a subset of the C runtime library supported in both host and device codes
    - A host component to control and access one or more devices from the host
    - A device component providing device-specific functions

- ## The GPU is always:
  - a coprocessor of the CPU or host
  - has its own DRAM (device memory)

- ## Approach similar to OpenMP, but designed to handle hardware accelerators
  - application source code portable
    - sequential binary -> traditional compiler

- ## CAPS HMPP is:
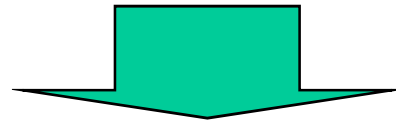  - a set of compiler directives and runtime software for multicore programming in C

Main memory

Application data

General purpose processor cores

upload remote data

dowload remote data

Remote procedure call

GPU Memory

Remote data

GPU cores

technology
from seed

inesc id
lisboa
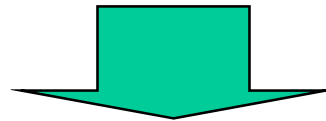
- A new execution model for local and remote computation is required
- Stream computing is expected to become a paradigm in high performance computing
- GPU never touches resources on host machine using stream-based computation, so security can be guaranteed

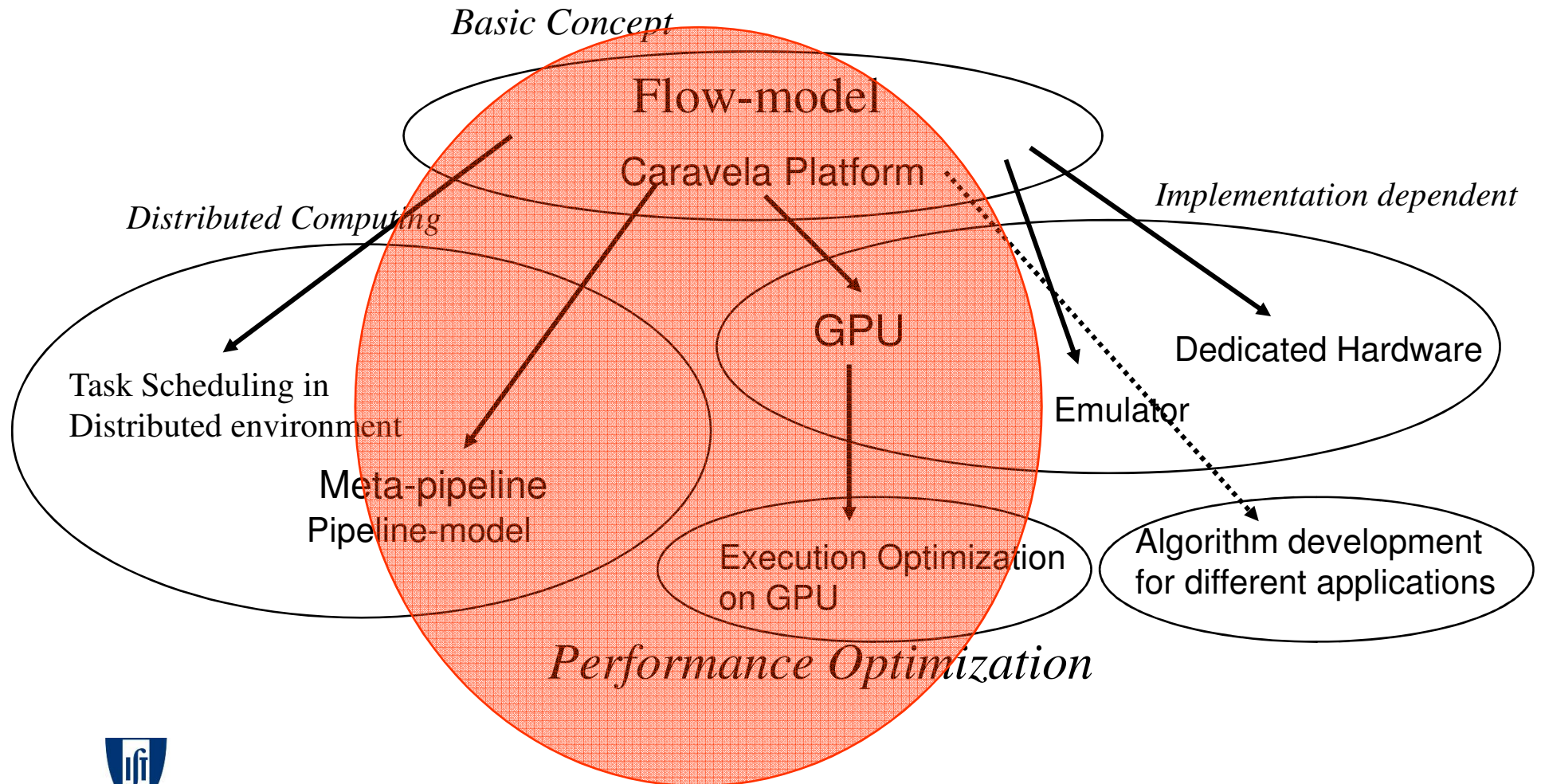Stream-based computation on GPU can be applied to distributed computing
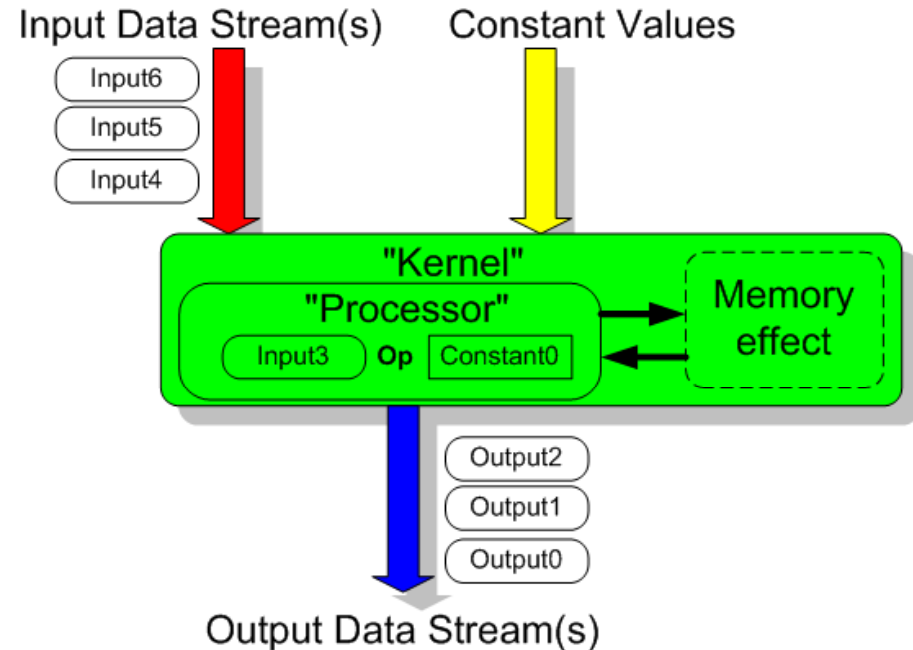
Caravela: A new platform for distributed computing

**Caravela Project:**
**Project Roadmap**

technology
from seed

inesc id
lisboa

*Basic Concept*

Flow-model

Caravela Platform

*Distributed Computing*

*Implementation dependent*

GPU

Task Scheduling in
Distributed environment

Dedicated Hardware

Emulator

Meta-pipeline

Pipeline-model

Execution Optimization
on GPU

Algorithm development
for different applications

*Performance Optimization*

**Caravela Platform:**
**Flow-model**

technology
from seed

inesc id
lisboa

- Memory effect by introducing feedback
- Program does not touch other resources beyond I/O streams
- Flow-model encapsulates a task object
- Flow-model unit can be fetched from remote site.



Caravela provides a set of tools
for executing a flow-model unit.

FlowModelCreator and Caravela Library

- Resource definition in Caravela library
    - Machine: has Adapter(s)
    - Adapter: has Shader(s)
    - Shader: Pixel Processor(s)
- Programming steps in application
    1. Acquire shaders
    2. Define flow-models
    3. Map flow-models to shaders
    4. Setup input streams
    5. Fire flow-models
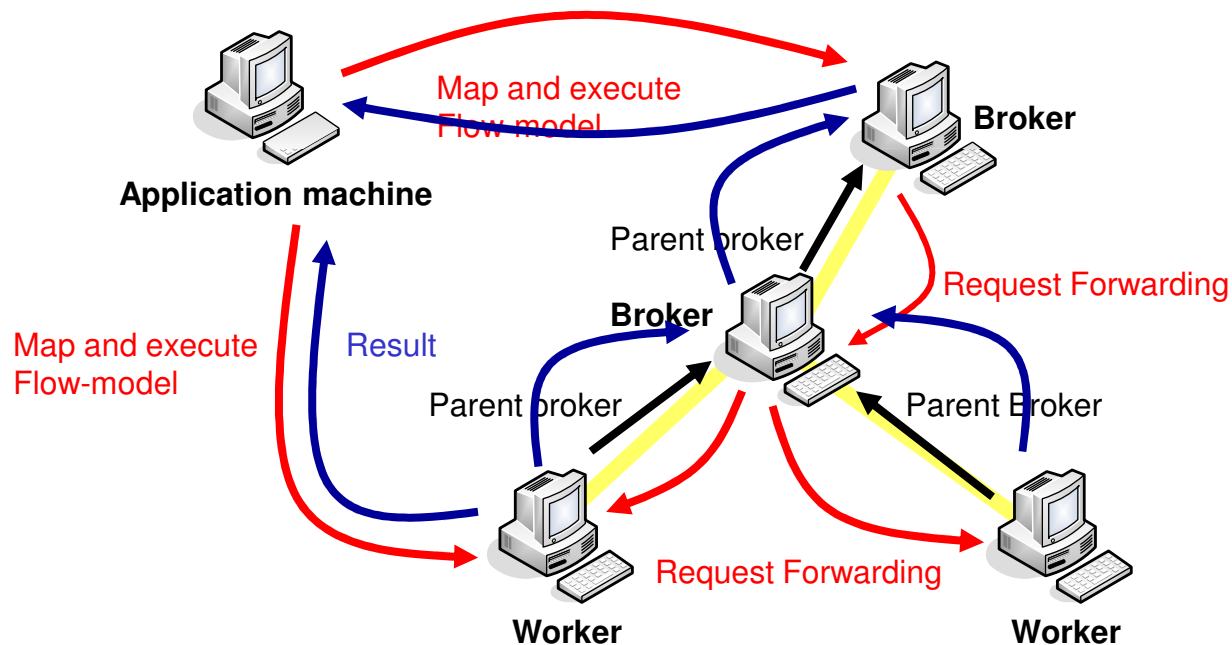    6. Get output data streams

Shader

Adapter

Machine

**Caravela Platform ：**
**Runtime for remote execution**

technology
from seed

inesc id
lisboa

- Remote execution runtime supports:
  - Worker server: executes flow-models.
  - Broker server: maintains routing information to worker servers.

**Caravela Platform:**
**Caravela library**

technology
from seed

inesc id
lisboa

- Initialization and Finalization
  CARAVELA_Initialize(RUNTIME), CARAVELA_Finalize(RUNTIME)
- Flow-model creation
  flow-model ← CARAVELA_CreateFlowModelFromFile(filename)
- Machine creation
  machine ← CARAVELA_CreateMachine(machine_type)
- Getting Shader
  shader ← CARAVELA_QueryShader(machine)
- Mapping Flow-model into Shader
  fuse ← CARAVELA_MapFlowModelIntoShader(shader, flow-model)
- Initialization for input data stream
  input data stream buffer ← CARAVELA_GetInputData(flow-model)
- Execution of Flow-model
  CARAVELA_FireFlowModel(fuse)
- Getting output data stream
  output data stream buffer ← CARAVELA_GetOutputData()

machine_type is "REMOTE" for remote execution.

INSTITUTO
SUPERIOR
TÉCNICO

inesc id
lisboa

```
void main(){
  int i,j;
  float inv = 1.0/Const4.x;
  vec4 res = vec4(0.0,0.0,0.0,0.0);

  vec2 coord = gl_TexCoord[0].xy;
  vec4 data0 = texture2D(CaravelaTex0, coord);
  coord.x+=inv;
  vec4 data1 = texture2D(CaravelaTex0, coord);
  coord.x+=inv;
  vec4 data2 = texture2D(CaravelaTex0, coord);

  // for x value
  for( j=0; j<4; j++ ){
     res.x += data0[j] * Const0[j];
     res.x += data1[j] * Const1[j];
  }

 // for y value
  for(j=1;j<4;j++)
     res.y += data0[j] * Const0[j-1];
  res.y += data1[0] * Const0[3];
  for( j=1; j<4; j++ )
     res.y += data1[j] * Const1[j-1];
  res.y += data2[0] * Const1[3];

...
  gl_FragData[0] = res;
}
```
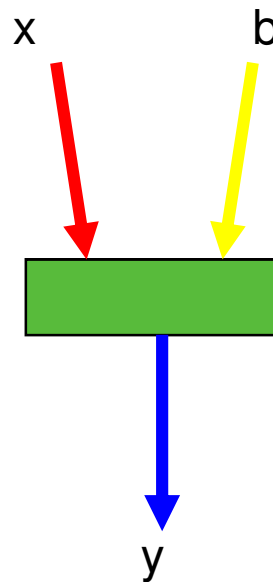
## OpenGL (GLSL)

$$y_n = \sum_{i=0}^{15} b_i * x_{n-i}$$

x          b

y

```
void main( in float2 t0:  TEXCOORD0,
           out float4 oC0: COLOR0){
   int j;
   float inv = 1.0/Const4.x;
   float4 res = 0;
   float2 coord = t0;

   float4 data0 = tex2D(CaravelaTex0, coord);
   coord.x += inv;
   float4 data1 = tex2D(CaravelaTex0, coord);
   coord.x += inv;
   float4 data2 = tex2D(CaravelaTex0, coord);

  // for x value
   for( j=0; j<4; j++ )
      res.x += data0[j] * taps[j][0];
   for( j=0; j<4; j++ )
      res.x += data1[j] * taps[j][1];
 // for y value
   for(j=1;j<4;j++)
      res.y += data0[j] * taps[j-1][0];
  res.y += data1[0] * taps[3][0];
  for( j=1; j<4; j++ )
     res.y += data1[j] * taps[j-1][1];
  res.y += data2[0] * taps[3][1];
  oC0 = res;
}
```

## DirectX (HLSL)

INSTITUTO
SUPERIOR
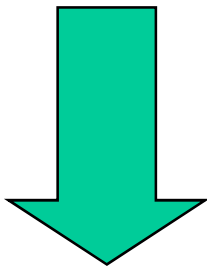TÉCNICO

**Caravela Platform:
Experimental Results**

technology
from seed

inesc id
lisboa

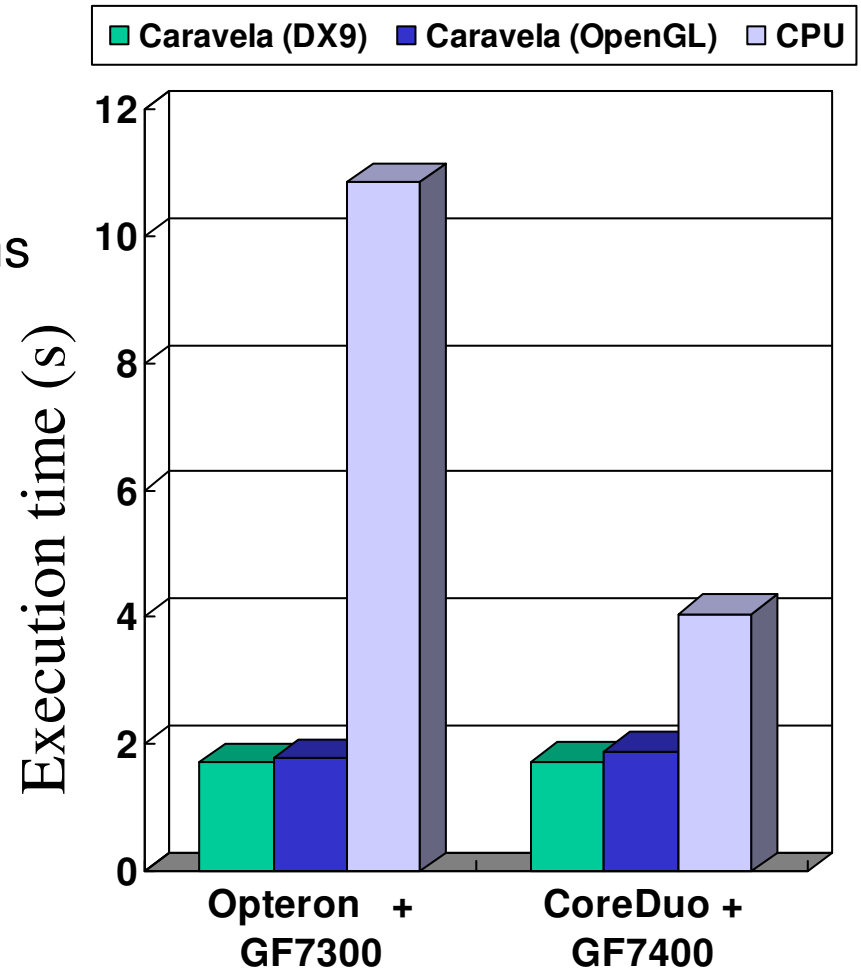| | Machine1 | Machine2 |
|---|---|---|
| CPU | AMD Opteron 2GHz 2GB DDR400 | Intel CoreDuo 1.66GHz 1GB DDR2 |
| GPU | NVIDIA GeForce 7300GS 256MB DDR | NVIDIA GeForce Go 7400 128MB DDR2 |

## Caravela Platform:
## Experimental Results

technology
from seed

inesc id
lisboa

- 1D FIR Filter

  - Input：1M samples ×30 iterations

- S= 4-10 times regarding to CPU

Caravela platform speeds up
local processing



Legend: ■ Caravela (DX9)　■ Caravela (OpenGL)　□ CPU

Y-axis: Execution time (s)

X-axis categories: Opteron + GF7300, CoreDuo + GF7400

**Local Optimizations:**
**Recursive processing**

technology
from seed

inesc id
lisboa

- Recursive processing with flow-model
  - Output streams must be copied to input streams
    → performance degrades due to the copy overhead
- Example: IIR Filter
  - Output "y" is feed-forwarded to input recursively.

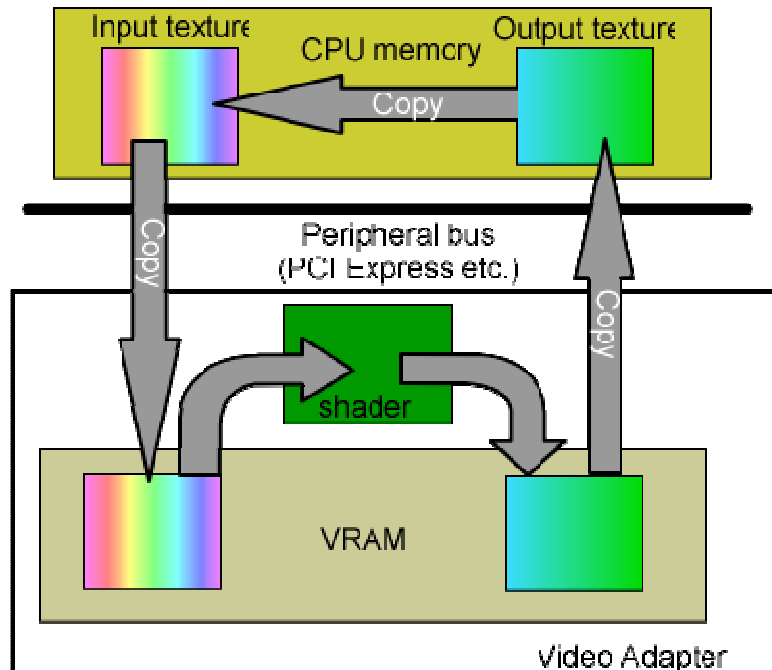$$y_n = \sum_{i=0}^{7} b_i * x_{n-i} + \sum_{k=1}^{8} b_k * y_{n-k}$$

x          b

Recursive  y

y

IIR filter

- ## Swap mechanism: Optimization for recursive I/O
  - *Pair* ← CARAVELA_CreateSwapIoPair(input_index, output_index)
  - CARAVELA_SwapFlowmodelIo(*Pair*)

Pair = CARAVELA_CreateSwapIoPair(0,1)

While(…){

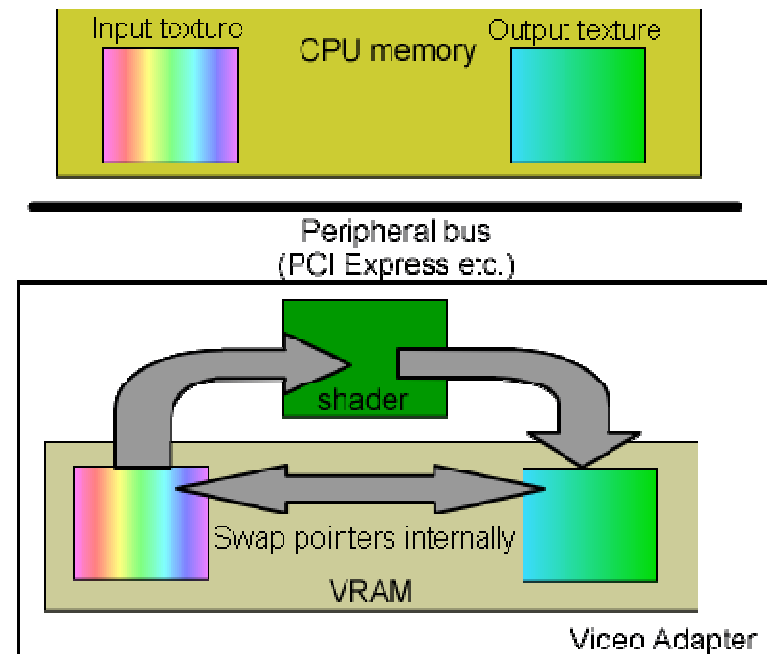    Fire flow-model…

    CARAVELA_SwapFlowmodelIo(pair)

}

## Conventional method



**Copy method (DirectX)**

Output stream copied
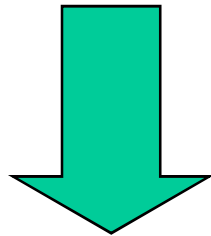VRAM →CPU memory and
CPU memory → VRAM

## Swap mechanism



**Swap method (OpenGL)**
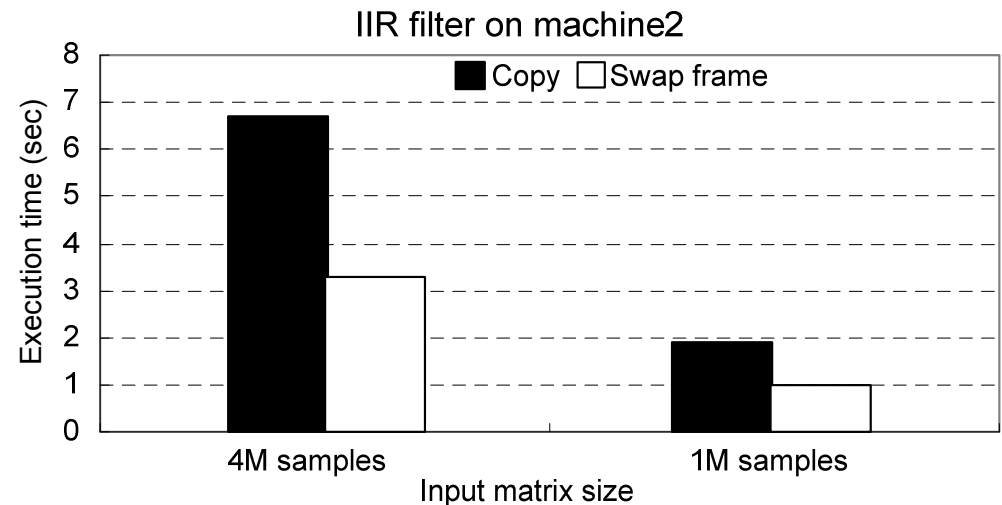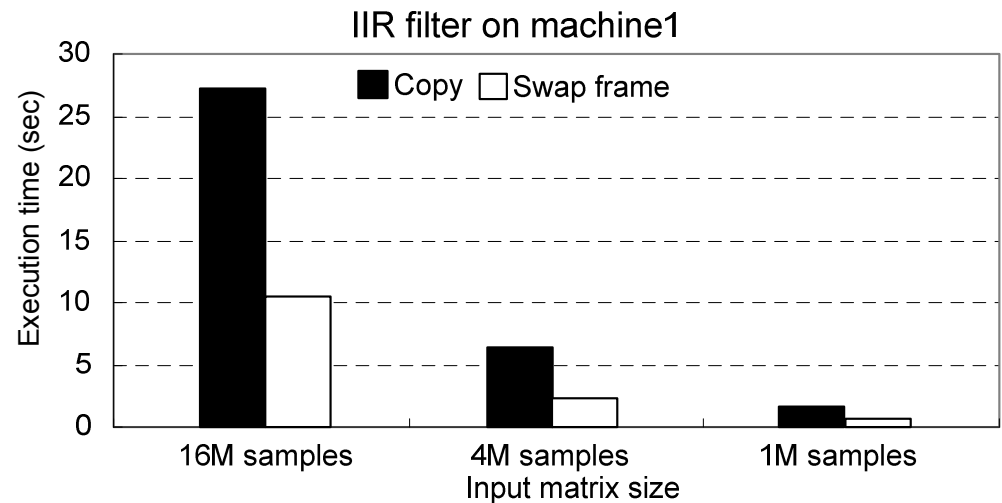
Exchanges pointers of I/O
buffers in the GPU side.

University of Murcia

- **OpenGL is used as the graphics runtime:**
  - CARAVELA_SwapFlowmodelIO() for swap mechanism
- **Swap:**

**Improves performance 55-60%**
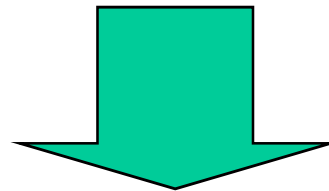
Swap mechanism is an effective optimization technique.

IIR filter on machine1

Execution time (sec)

Copy ☐ Swap frame

16M samples    4M samples    1M samples
Input matrix size

IIR filter on machine2

Execution time (sec)

Copy ☐ Swap frame

4M samples    1M samples
Input matrix size

**Local optimizations:**
**Remap method**

technology
from seed

inesc id
lisboa

- # I/O overhead of GPGPU application
  - Copy operation among CPU memory-VRAM

  - Overhead in GPU at writing output stream to VRAM

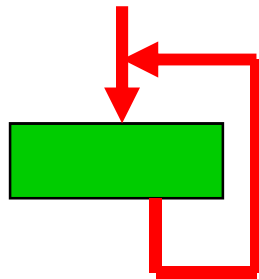  - Overhead in Pixel processor reading textures
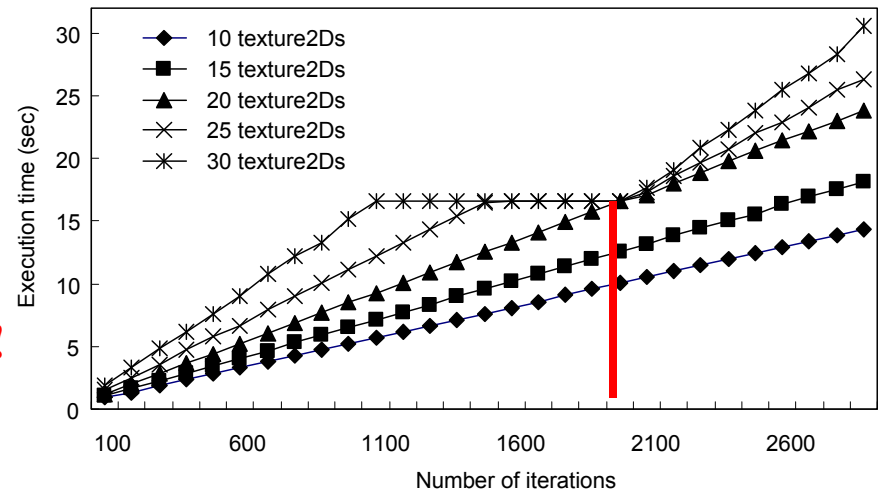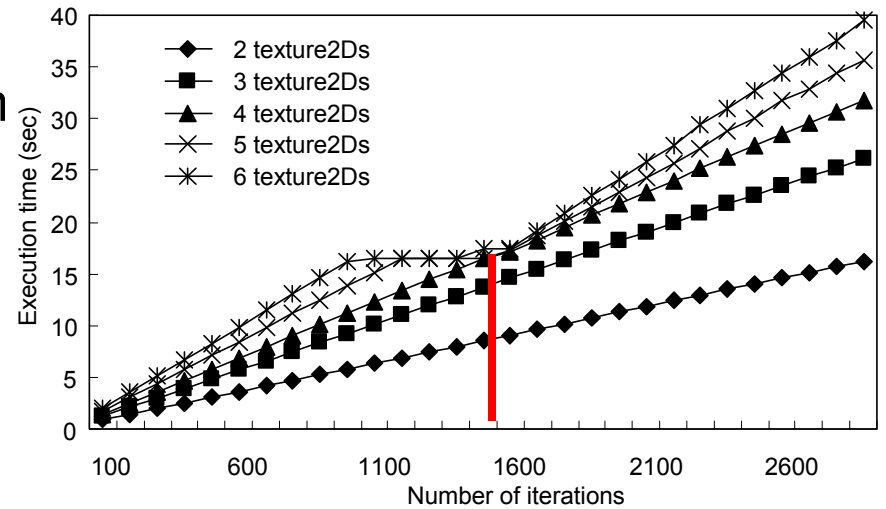
Smaller texture size may result in better performance.

## Local optimizations:
## Remap method

technology
from seed

inesc id
lisboa

- Iterating with 3000x3000 texture input and applying Swap mechanism
  - Spot depending on the number of iterations of Swap mechanism
  - GeForce7300: 1500 iterations

  - GeForce7900: 2000 iterations

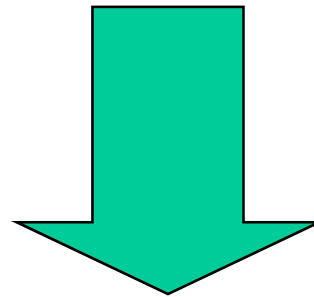**Swap iteration should be reset at the spot!**

- For applications where calculation size decreases:
  - Flow-model should be mapped again after the input texture sizes are reduced
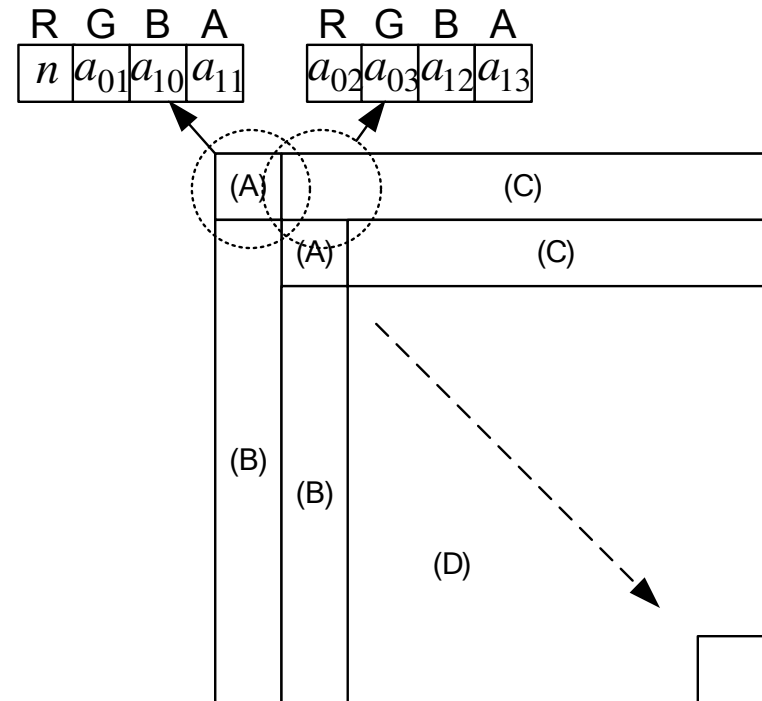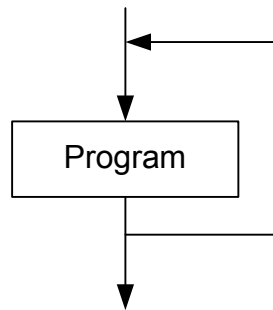  - Applying a threshold number of iterations for Swap, flow-model is mapped again at the spot

# Remap method

(A) Normalization of diagonal elements
(B) Orthogonalization
(C) Normalization



Elements previously calculated are forwarded to the
output data stream without any calculation
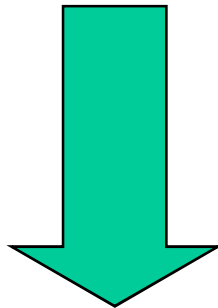
**Caravela Platform** ：
**Applying remap method to LU decomposition**

technology
from seed

inesc id
lisboa
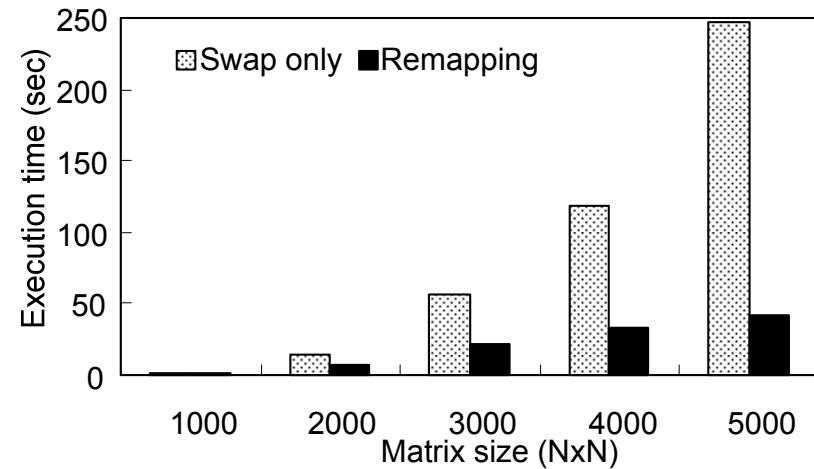
- GeForce7300
  - Remap flow-model every 1500 Swap iterations
- GeForce7900
  - Remap flow-model every 2000 Swap iterations

- Reduction of 80% in execution time
  - Remap method further improves performance in the top of the swap mechanism



(a) GeForce7300



(b) GeForce7900

## Remote execution:
## Meta-pipeline

technology
from seed

inesc id
lisboa

- Executing the flow-model in a remote machine:
  - Sending input data to the remote machine,
  - receiving output data from the remote machine,
  - scheduling the execution

Pipeline-model

*Meta-pipeline*

Map to worker machine network

Caravela worker machine network

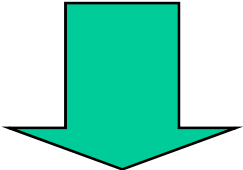**Remote execution:**
**Pipeline model**

technology
from seed

inesc id
lisboa

- ## I/O ports of the Pipeline-model
  – *ENTRANCE* port
  – *EXIT* port
  – *INTERMEDIATE* port
- ## When all input streams are ready, flow-model is executed
- ## Deadlock might occur if feedback edges exist
  – *INITONCE* port

ENTRANCE ports

INITONCE port

flowmodel0

flowmodel1    flowmodel2

EXIT ports

- Extended functions for Caravela library
  - CARAVELA_CreatePipeline()
  - CARAVELA_AddShaderToPipeline()
  - CARAVELA_AttachFlowModelToShader()
  - CARAVELA_ConnectIO()
  - CARAVELA_Specify[InitOnce | Exit | Intermediate]Port()
  - CARAVELA_ImplementPipelineModel()
  - CARAVELA_SendInputDataToPipeline()
  - CARAVELA_ReceiveOutputDataFromPipeline()

  During local execution: it promotes pipeline execution.
  During remote execution: communication with worker servers.

# Remote execution:
## 2D DWT

technology
from seed

inesc id
lisboa

- 2D Discrete Wavelet Transform
  - Image compression (JPEG2000), denoising, edge detection, enlarge…

$$LL_n = \sum_{k=0}^{K-1}\sum_{m=0}^{M-1} LL_{n-1}(2i+k, 2j+m)l(m)l(k)$$

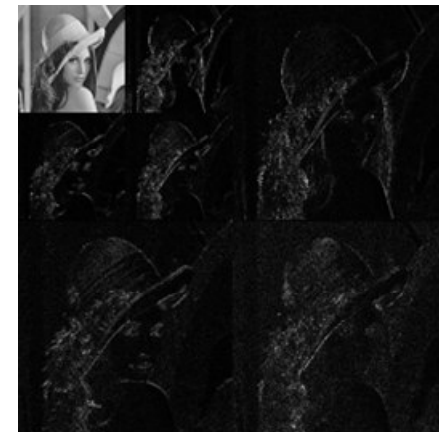$$HL_n = \sum_{k=0}^{K-1}\sum_{m=0}^{M-1} LL_{n-1}(2i+k, 2j+m)h(m)l(k)$$

$$LH_n = \sum_{k=0}^{K-1}\sum_{m=0}^{M-1} LL_{n-1}(2i+k, 2j+m)l(m)h(k)$$

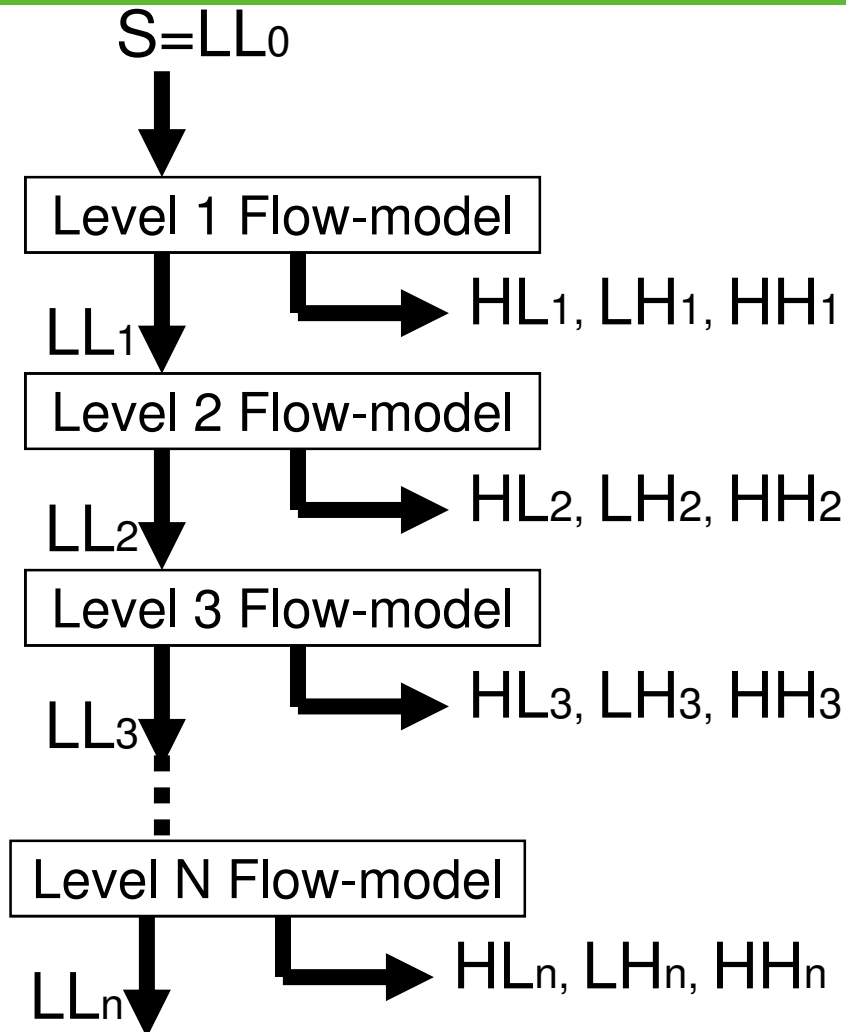$$HH_n = \sum_{k=0}^{K-1}\sum_{m=0}^{M-1} LL_{n-1}(2i+k, 2j+m)h(m)h(k)$$

| $LL_2$ | $HL_2$ | $HL_1$ |
| $LH_2$ | $HH_2$ | |
| $LH_1$ | | $HH_1$ |

2 decomposition level

# Remote execution:
## 2D DWT

technology
from seed

inesc id
lisboa

$S = LL_0$

$\downarrow$

Level 1 Flow-model

$\rightarrow$ $HL_1, LH_1, HH_1$

$LL_1 \downarrow$

Level 2 Flow-model

$\rightarrow$ $HL_2, LH_2, HH_2$

$LL_2 \downarrow$

Level 3 Flow-model

$\rightarrow$ $HL_3, LH_3, HH_3$

$LL_3 \downarrow$

$\vdots$

Level N Flow-model

$\rightarrow$ $HL_n, LH_n, HH_n$

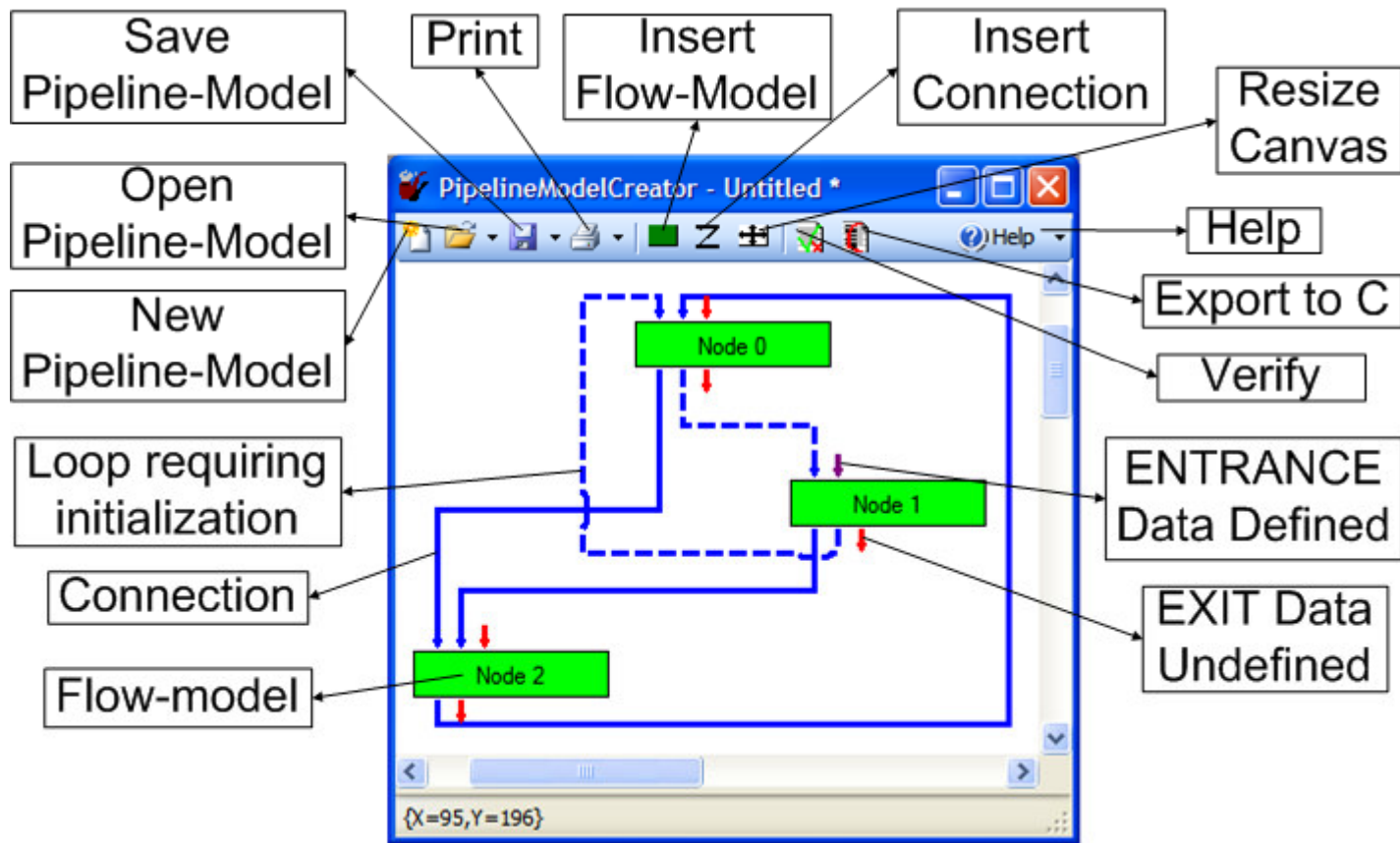$LL_n \downarrow$

```
void main()
{
  float delta = 1/NUMDATA;
  vec4 tmp, tmp0, tmp1, result;
  vec2 coord = gl_TexCoord[0].xy;
  vec2 caux;  int i;
  coord += coord;
  caux = coord;
  // horizontal direction
  for (i=0; i<4; i++, coord.y += delta){
    tmp.x = texture2D(CaravelaTex0, coord).x;
    coord.x += delta;
    tmp.y = texture2D(CaravelaTex0, coord).x;
    coord.x += delta;
    tmp.z = texture2D(CaravelaTex0, coord).x;
    coord.x += delta;
    tmp.w = texture2D(CaravelaTex0, coord).x;
    tmp0[i] += dot(tmp, const0);
    tmp1[i] += dot(tmp, const1);
    coord.x = caux.x;
  }
  // vertical direction
  result.x = dot(tmp0, const0);
  result.y = dot(tmp0, const1);
  result.z = dot(tmp1, const0);
  result.w = dot(tmp1, const1);
  // LL sub-band stream
  gl_FragData[0] = result;
  // LH, HL and HH sub-bands stream
  gl_FragData[1] = result;
}
```

**University of Murcia**

**Remote execution:**
**PipelineModelCreator tool**

technology
from seed

inesc id
lisboa

- MPI + flow-model = CaravelaMPI

- Caravela platform operated in command line mode (operating system)

- Attach other hardware co-processors to the Caravela platform (co-processors on FPGAs,….)  and generate hardware from the OpenGL programs

- Test Meta-Pipeline with large real problems
  – Japan-Cyprus-Portugal

- Virtualization -> trying to abstract programming  from architectures
  – programming kernels -> also needed to CUDA

- **Papers**

1. Shinichi Yamagiwa, Leonel Sousa, "Caravela: A Novel Environment for stream-based distributed computing", IEEE Computer Magazine, May 2007, pp.76-83

2. Shinichi Yamagiwa, Leonel Sousa, "Design and implementation of a stream-based distributed computing platform using graphics processing units", ACM International Conference on Computing Frontier, May 2007

3. Shinichi Yamagiwa, Leonel Sousa, Diogo Antão, "Data buffering optimization methods toward a uniform programming interface for GPU-based applications", ACM International conference of Computing Frontier, May 2007

4. Shinichi Yamagiwa, Leonel Sousa, Tomas Brandao, "Meta-Pipeline: A new execution mechanism for distributed pipeline processing", 6th International Symposium on Parallel and Distributed Computing (ISPDC 2007), August 2007

5. Shinichi Yamagiwa and Diogo Ricardo Cardoso Antao and Leonel Sousa, Design and Implementation of a Graphical User Interface for Stream-based Distributed Computing, the IASTED International Conference on Parallel and Distributed Computing and Networks (PDCN 2008), Feb. 2008

- **Book chapter**

1. Concurrent and Parallel Computing: Theory, Implementation and Applications, chapter 1, NOVA Publishers, May 2008

- **Patent**

1. "Program execution method applied to data streaming in distributed heterogeneous computing environment", Portuguese national patent

For more detailed information,
please visit:
http://www.caravela-gpu.org

**technology**
**from seed**

inesc id
lisboa