# Toward Real-Time Simulation of Cardiac Dynamics

## Ezio Bartocci

## SUNY Stony Brook

**Joint work with**

**E. Cherry, J. Glimm, R. Grosu, S. A. Smolka, F. Fenton**
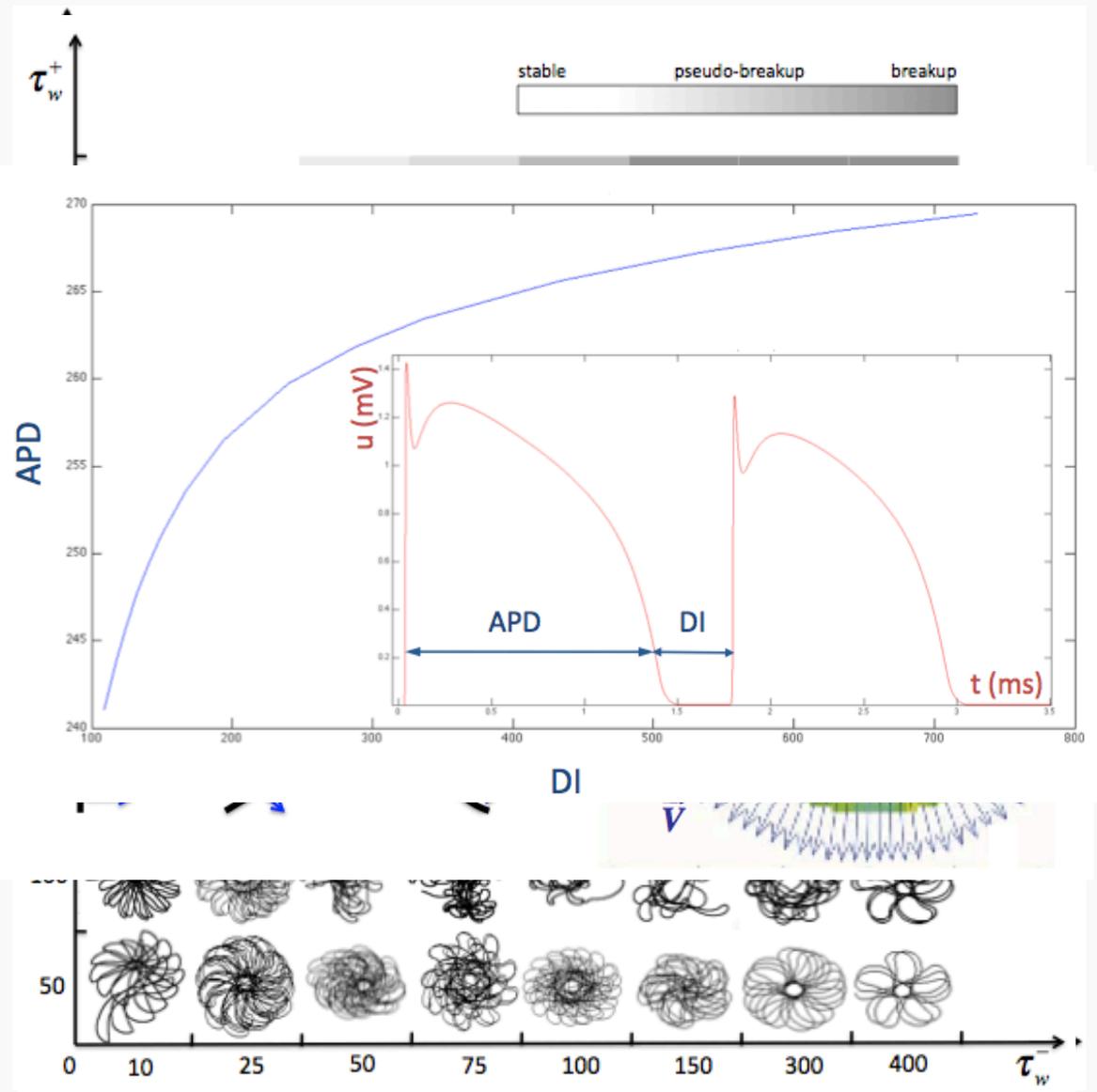
# Outline

- **Motivation**
- **Cardiac Models as Reaction Diffusion Systems**
- **CUDA Programming Model**
- **Reaction Diffusion in CUDA**
- **Case Studies**
- **Work in Progress**

# Motivation for our Work

Simulation-Based Analysis

- **Spiral Formation**
- **Spiral Breakup**
- **Tip Tracking**
- **Front Wave Tracking**
- **Curvature Analysis**
- **Conduction Velocity**
- **Restitution Analysis**

# Cardiac Models as Reaction Diffusion Systems

**Membrane's AP depends on:**

- **Stimulus** (voltage or current):
  - **External / Neighboring cells**
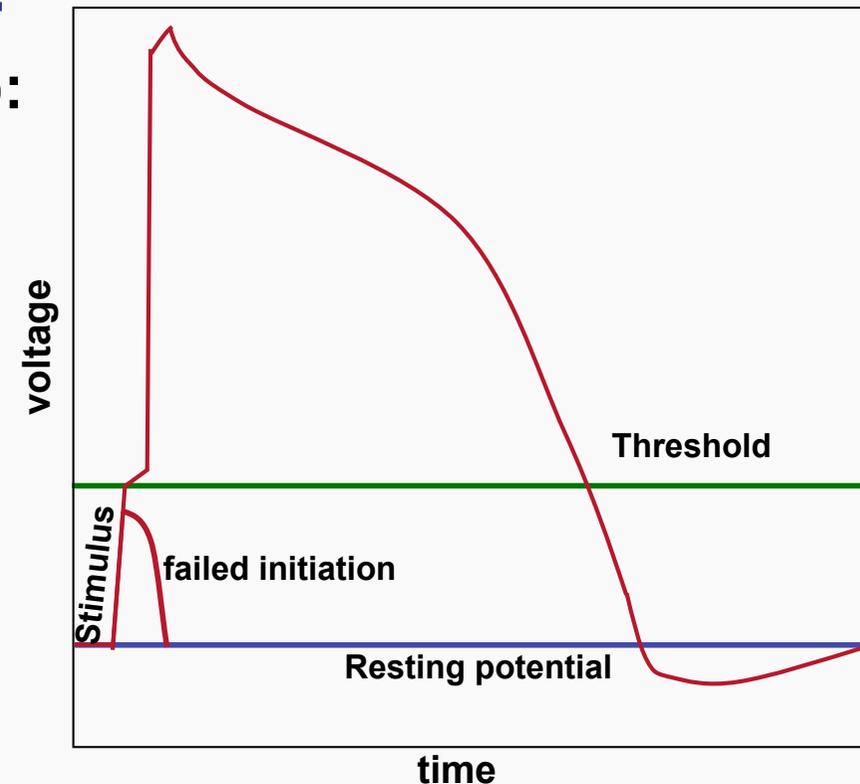
- **Cell's state**

**AP has nonlinear behavior!**

- **Reaction diffusion system:**

$$\frac{\partial \mathbf{u}}{\partial t} = R(\mathbf{u}) + \nabla(D\nabla\mathbf{u})$$

Behavior In time

Reaction

Diffusion

**Schematic Action Potential**

voltage

Stimulus
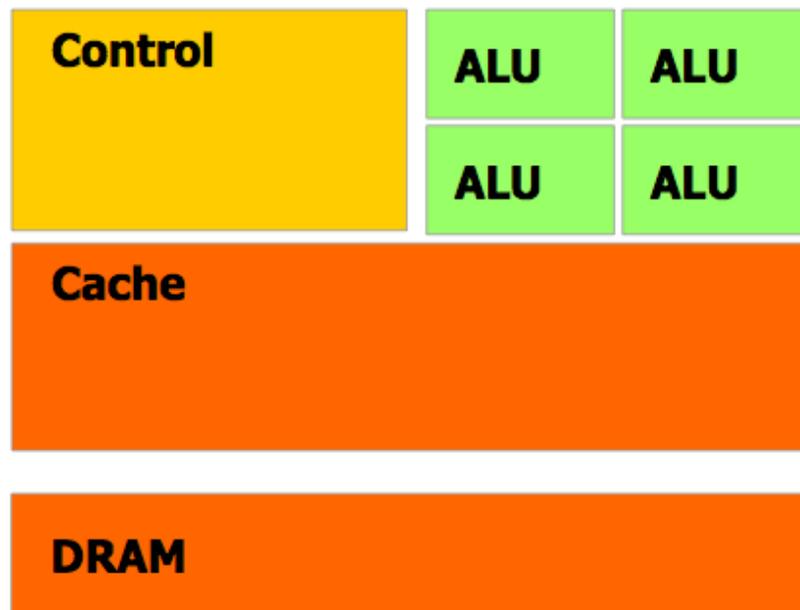
Threshold

failed initiation

Resting potential

time

# Cardiac Models

- **Minimal Model (Flavio-Cherry) (4 v) Human**
- **Beeler-Reuter (8 v) Canine**
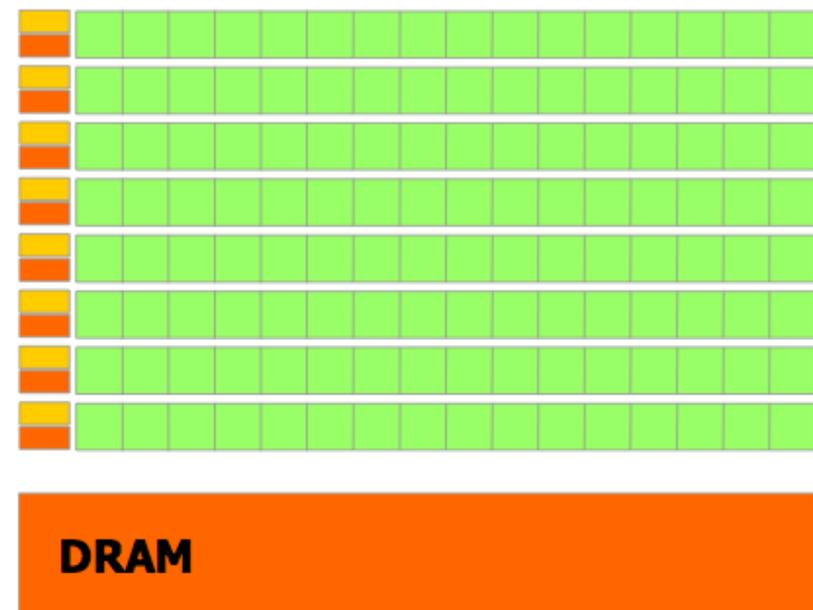- **Ten Tussher Panfilov (19 v) Human**
- **Iyer (67 v) Human**

# Available Technologies

## CPU based

| Control | ALU | ALU |
|---------|-----|-----|
|         | ALU | ALU |

Cache

DRAM

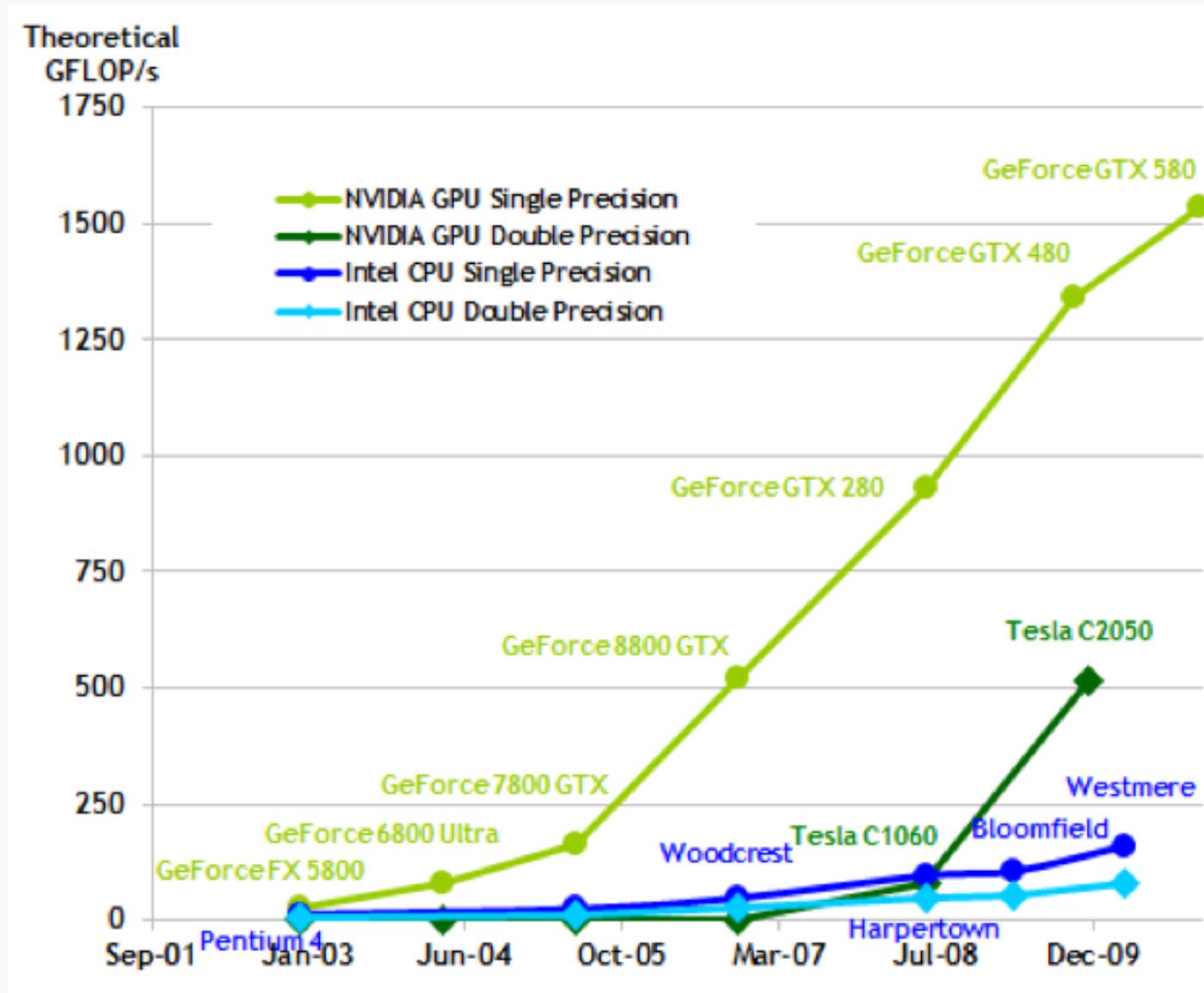**CPU**

## GPU based

DRAM

**GPU**

The GPU devotes more transistors to data processing
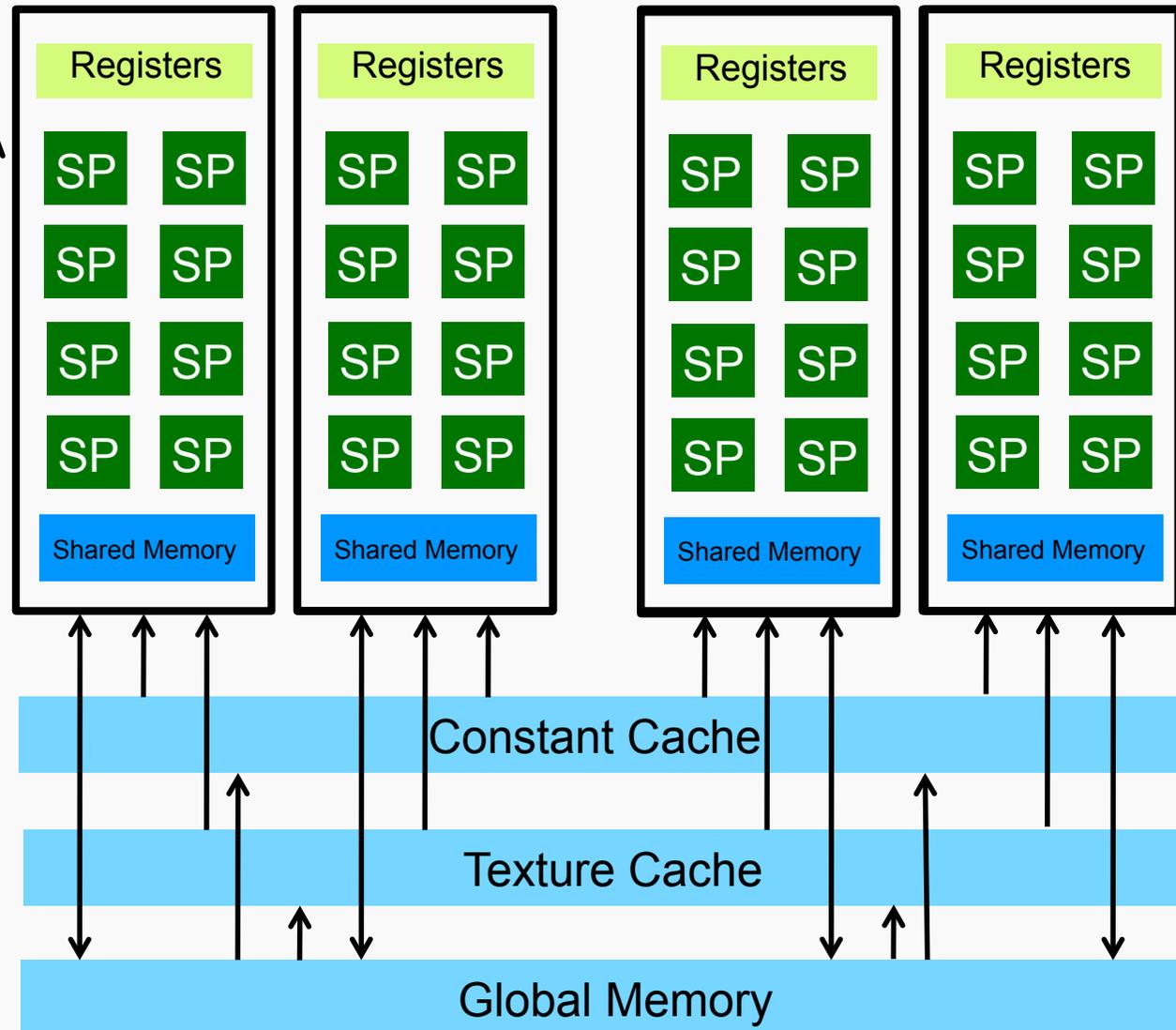
*This image is from CUDA programming guide*

# GPU vs CPU

# GPU Architecture

MULTIPROCESSORS

Each GPU consists of a Set of multiprocessors.

| Registers | | Registers | | Registers | | Registers |
|---|---|---|---|---|---|---|
| SP SP | | SP SP | | SP SP | | SP SP |
| SP SP | | SP SP | | SP SP | | SP SP |
| SP SP | | SP SP | | SP SP | | SP SP |
| SP SP | | SP SP | | SP SP | | SP SP |
| Shared Memory | | Shared Memory | | Shared Memory | | Shared Memory |

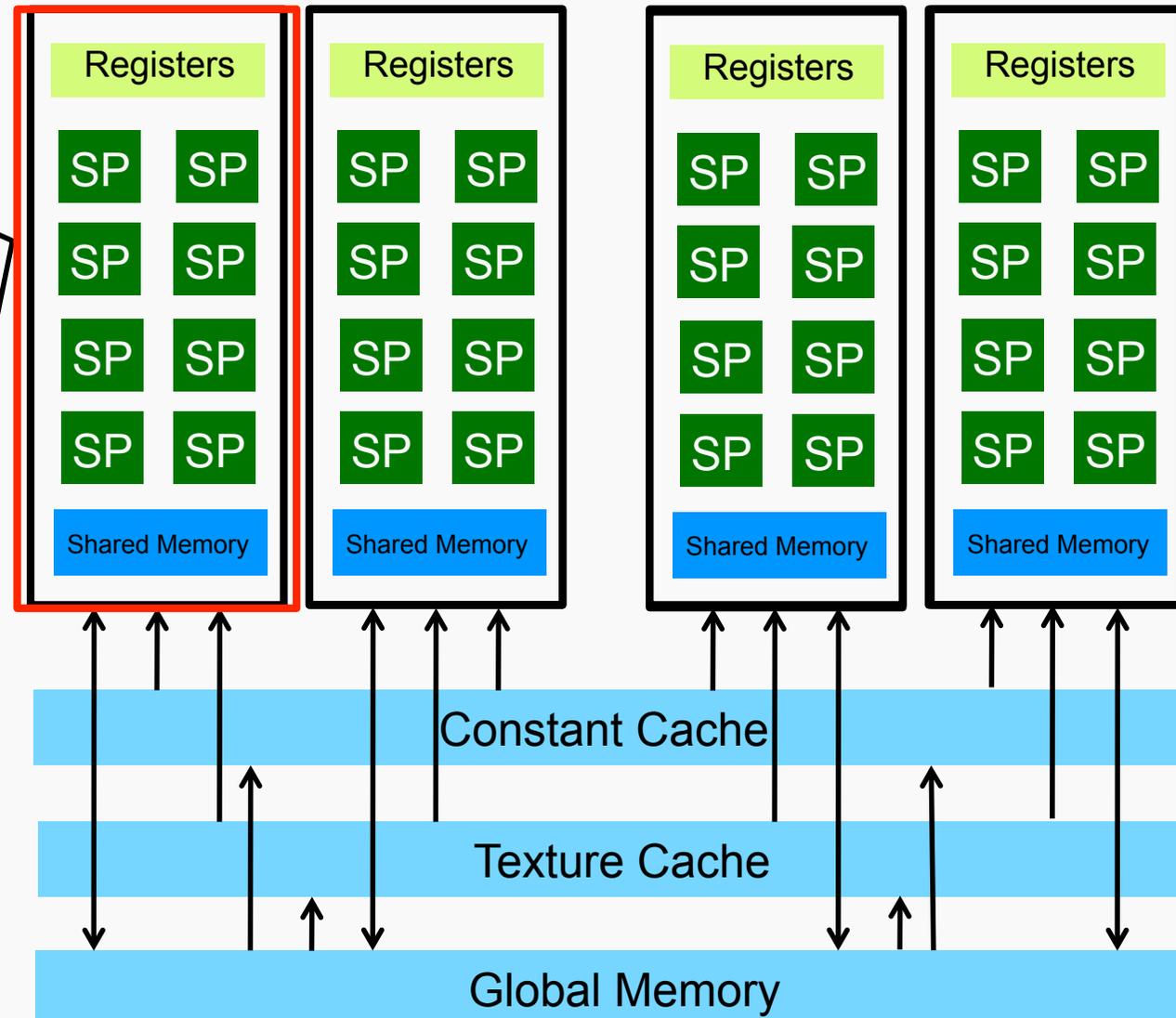Constant Cache

Texture Cache

Global Memory

# GPU Architecture

MULTIPROCESSORS

Each Multiprocesssor can have 8/32 Stream Processors (SP) (called by NVIDIA also cores) which share access to local memory.
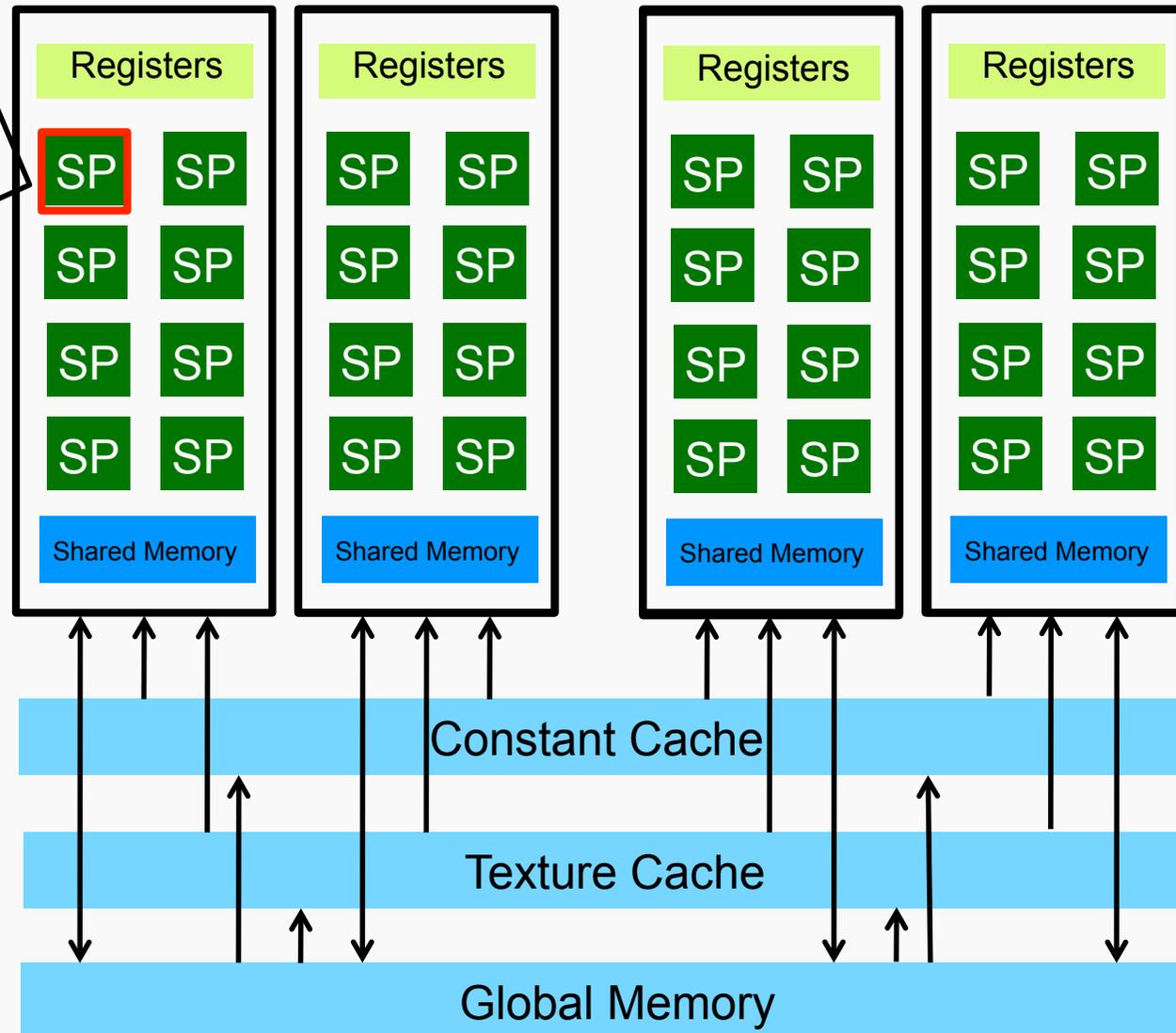
| Registers | Registers | Registers | Registers |
|---|---|---|---|
| SP  SP | SP  SP | SP  SP | SP  SP |
| SP  SP | SP  SP | SP  SP | SP  SP |
| SP  SP | SP  SP | SP  SP | SP  SP |
| SP  SP | SP  SP | SP  SP | SP  SP |
| Shared Memory | Shared Memory | Shared Memory | Shared Memory |

Constant Cache

Texture Cache

Global Memory

# GPU Architecture

MULTIPROCESSORS

Each Stream Processor (core) contains a fused multiply-adder capable of single precision arithmetic. It is capable of completing 3 floating point operations per cycle - a fused MADD and a MUL.
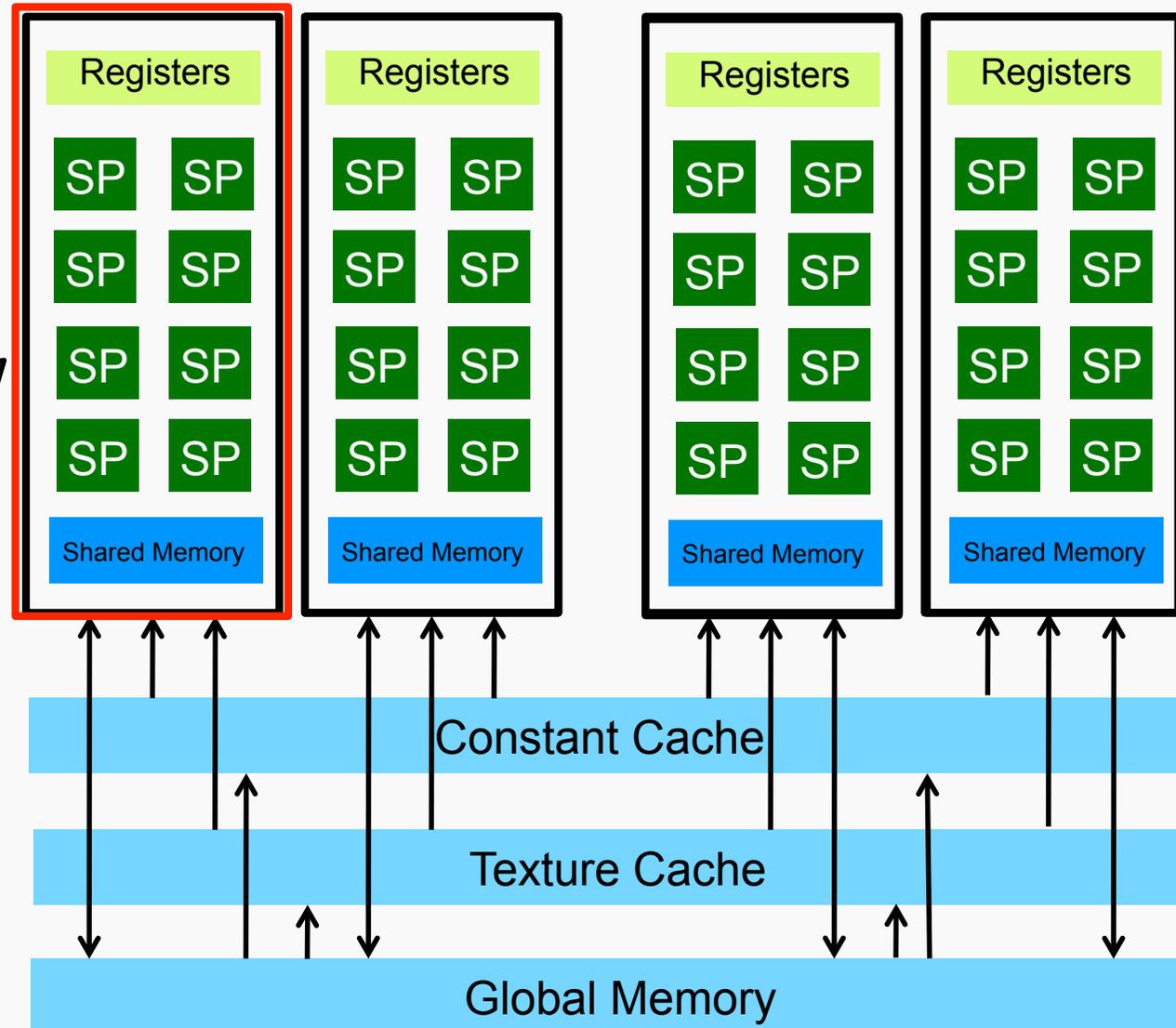
| Registers |
|---|
| SP SP |
| SP SP |
| SP SP |
| SP SP |
| Shared Memory |

| Registers |
|---|
| SP SP |
| SP SP |
| SP SP |
| SP SP |
| Shared Memory |

| Registers |
|---|
| SP SP |
| SP SP |
| SP SP |
| SP SP |
| Shared Memory |

| Registers |
|---|
| SP SP |
| SP SP |
| SP SP |
| SP SP |
| Shared Memory |

Constant Cache

Texture Cache

Global Memory

# GPU Architecture

MULTIPROCESSORS

Each multiprocessor can contain one or more 64-bit fused multiple adder for double precision operations.

| Registers | Registers | Registers | Registers |
|---|---|---|---|
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| Shared Memory | Shared Memory | Shared Memory | Shared Memory |

Constant Cache

Texture Cache

Global Memory

# Memory Hierarchy



MULTIPROCESSORS

The fastest available Memory for GPU computation is device registers. Each multiprocessor contains 16KB of registers. The registers are partitioned among the MP-resident threads

Registers

SP SP
SP SP
SP SP
SP SP

Shared Memory

Registers

SP SP
SP SP
SP SP
SP SP

Shared Memory

...

Registers

SP SP
SP SP
SP SP
SP SP

Shared Memory

Registers

SP SP
SP SP
SP SP
SP SP

Shared Memory

Constant Cache

Texture Cache

Global Memory

# Memory Hierarchy

# Memory Hierarchy

MULTIPROCESSORS

| Registers | Registers | Registers | Registers |
|---|---|---|---|
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| Shared Memory | Shared Memory | Shared Memory | Shared Memory |

...

Constant Cache

Texture Cache

Global Memory

GPUs have also DRAM
The latency is 150x is slower
then registers and
shared memory

# Memory Hierarchy

Constant memory, as the name implies, is a read-only region which also has a small cache.

# Memory Hierarchy

MULTIPROCESSORS

| Registers | Registers | Registers | Registers |
|---|---|---|---|
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| Shared Memory | Shared Memory | Shared Memory | Shared Memory |

...

Texture memory is read-only with a small cache optimized for manipulation of textures. It also provides built-in linear interpolation of the data. also provides built-in linear interpolation of the data.

Constant Cache

Texture Cache

Global Memory

# Memory Hierarchy

| Registers | Registers | Registers | Registers |
|---|---|---|---|
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| SP SP | SP SP | SP SP | SP SP |
| Shared Memory | Shared Memory | Shared Memory | Shared Memory |

...

Constant Cache

Texture Cache

Global memory is available to all threads and is persistent between GPU calls.

Global Memory

# CUDA Programming Model

**Single Instruction Multiple Threads (SIMT)**
**similar to**
**Single Instruction Multiple Data (SIMD)**

| THREAD ID 0 | THREAD ID 1 | THREAD ID 2 | THREAD ID 3 |
|---|---|---|---|
| A[ID]=ID | A[ID]=ID | A[ID]=ID | A[ID]=ID |
| if (ID%2) | if (ID%2) | if (ID%2) | if (ID%2) |
| A[ID]+=2; | A[ID]+=2; | A[ID]+=2; | A[ID]+=2; |
| A[ID]*=2; | A[ID]*=2; | A[ID]*=2; | A[ID]*=2; |
| else | else | else | else |
| A[ID]=0; | A[ID]=0; | A[ID]=0; | A[ID]=0; |
| endif | endif | endif | endif |
| A[ID]+=2; | A[ID]+=2; | A[ID]+=2; | A[ID]+=2; |

A vector

## CPU Host

Serial Code i

Kernel Invocation

Serial Code j

## GPU Device

Grid k

| Block (0,0) | Block (0,1) |
|---|---|
| Block (1,0) | Block (1,1) |
| Block (2,0) | Block (2,1) |
| Block (3,0) | Block (3,1) |
| Block (4,0) | Block (4,1) |
| Block (5,0) | Block (5,1) |

When branches occur in the code (e.g. due to if statements) the divergent threads will become inactive until the conforming threads complete their separate execution.

When execution merges, the threads can continue to operate in parallel.

# CUDA Programming Model

GRID

BLOCK

Different threads are multiplexed and executed by the same core in order to reduce the latency of memory access.

Each Thread block is executed by a multiprocessors

The max number of threads for a thread block is 512 and it depends on the amount of registers that each thread may need.

GPU DEVICE

THREAD BLOCK

Global Memory

REGISTERS

SHARED MEMORY

# Tesla C1060    Fermi C2070





30 Multiprocessors
240 Cores
Processor core clock: 1.296 GHz
933 Gigaflops (Single precision)
78 Gigaflops (Double Precision)
Max Bandwidth(102 Gigabytes/sec)
4 GB of DRAM

Cost: **1000 $**

14 Multiprocessors
448 Cores
Processor core clock: 1.15 GHz
1030 Gigaflops (Single precision)
515 Gigaflops (Double precision)
Max Bandwidth (144 GBytes/sec)
6 GB of DRAM

Cost: **3200 $**

# Reaction-Diffusion in CUDA

$$\frac{\partial \mathbf{u}}{\partial t} = R(\mathbf{u}) + \nabla(D\nabla\mathbf{u})$$

For each time step, a set of (ODEs) and Partial Differential Equations (PDEs) must be solved.

```
for (timestep=1; timestep < nend; i++){
        solveODEs <<grid, block>> (….);
        calcLaplacian <<grid, block>> (….);
}
```

Solving ODEs using different method depending on the model:

- Runge-Kutta 4th order
- Forward Euclid
- Backward Euclid
- Semi-Implicit
……

Solving PDEs (Calc the Laplacian)



$$\nabla(D\nabla\mathbf{u})_{i,j} = \frac{Ddt}{dx^2}\left(u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1} - 4u_{i,j}\right)$$

# Optimize the Reaction Term in CUDA

## Heaviside Simplification

- In order to avoid divergent branches among threads we substitute if then else condition of Heaviside functions in this way:

```
If (x > a){
    y = b;
}else{
    y = c;
}
```

$y = b + (x > a)*(b\_c)$

$a, b, b\_c$ (b-c) are constant

## Precomputing Lookup Tables using Texture

- We build tables where we precompute nonlinear part of the ODEs, we bind these table to textures and we exploit the built-in linear interpolation feature of the texture.

# Optimize the Reaction Term in CUDA

**Kernel splitting:**

- For complex models, we need to split the ODEs solving in many Kernels in order to have enough registers for thread to perform our calculation.

**Use –use_fast_math compiler option to substitute**

- In the models that use log, exp, sqrt, functions we substitute them with the GPU built-in functions.

**Using Costant memory for common parameters (dt, dx, ..)**

# Optimize the Diffusion Term in CUDA

Solving PDEs (Calc the Laplacian)

$$\nabla(D\nabla\mathbf{u})_{i,j} = \frac{Ddt}{dx^2}\left(u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1} - 4u_{i,j}\right)$$

Using texture we can reduce the latency
In texture the data is cached (optimize for 2D Locality)
Drawback: It supports only single precision

Each location is a float (4 bytes) The global memory latency is very slow. The memory is accessed in multiples of 64 bytes

# Optimize the Diffusion Term in CUDA

Another Technique is using SHARED MEMORY

| The yellow and red threads read the location from the global memory into the shared memory. | → SYNCH → | The red threads calculates the laplacian using the values in the shared memory. |

$$\nabla(D\nabla\mathbf{u})_{i,j} = \frac{Ddt}{dx^2}\left(u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1} - 4u_{i,j}\right)$$

THREAD BLOCK



REGISTERS

SHARED MEMORY

Step 1          Step 2          Step 3

This technique supports single and double precision

Drawback: The number of threads is greater than the number of elements

# Case Study 1: Minimal Model (4V)

# Perfomances



CPU Computation                    GPU Computation

# Naïve Implementation

# Reaction optimized (Diffusion with Shared Memory)



Minimal Model (4V) Diffusion Shared memory

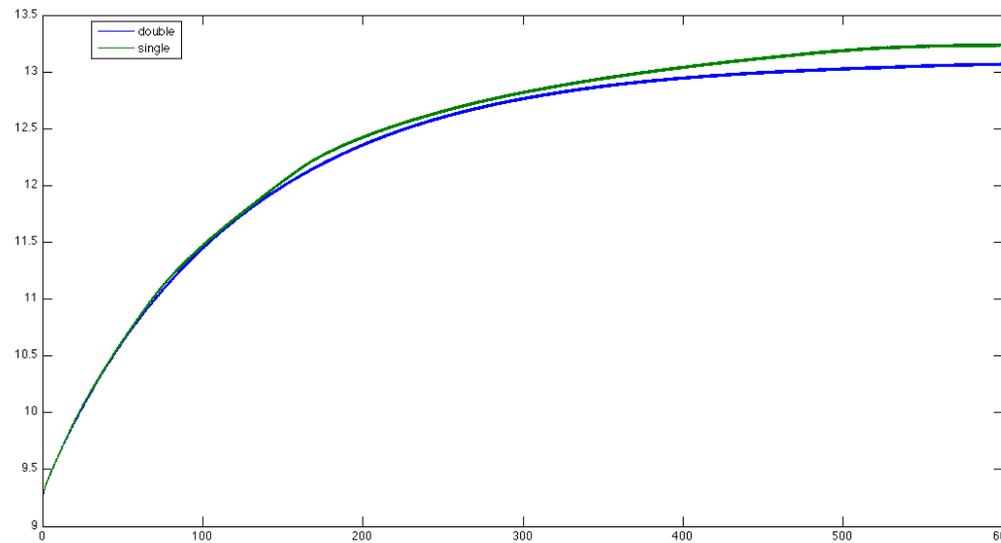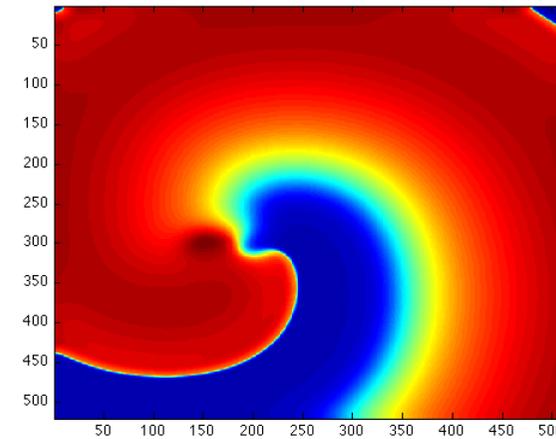# Reaction optimized (Diffusion with Texture)



Minimal Model (4V) Optimized Version Texture for Diffusion

# Beeler-Reuter Model (8V)
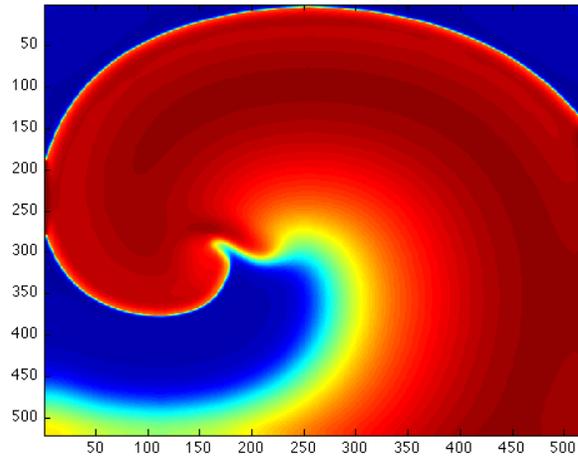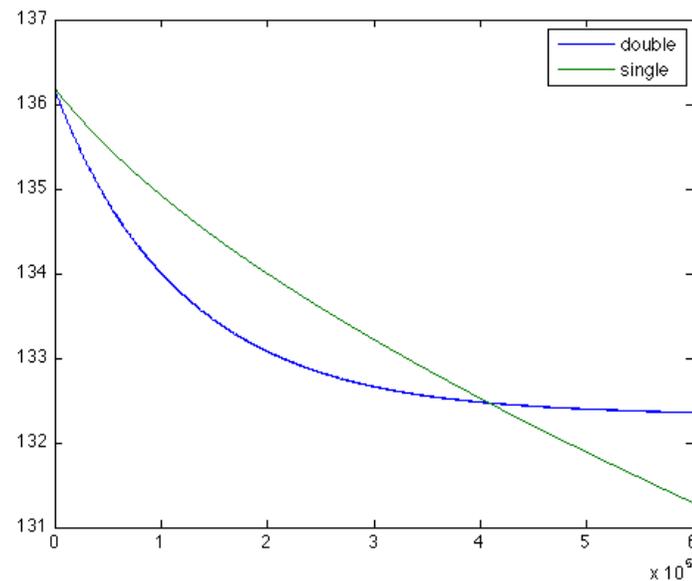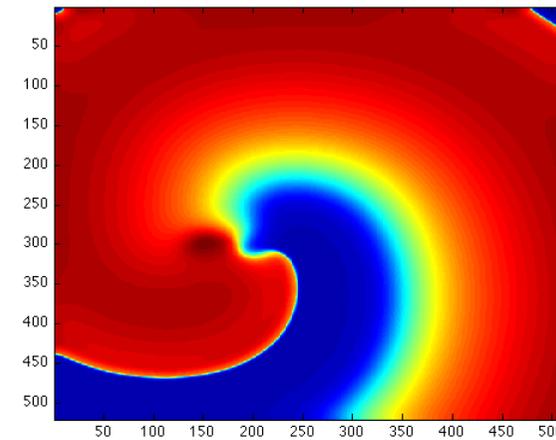
# Reaction optimized (Diffusion with Shared Memory)



Beeler Reuter Model (8 Variables) Implentation using Shared Memory for Diffusion

# Reaction optimized (Diffusion with Texture)



Beeler Reuter Model (8 Variables) Implentation using Texture for Diffusion

# Ten Tusscher Panfilov Model (19V)

# Reaction optimized (Diffusion with Shared Memory)
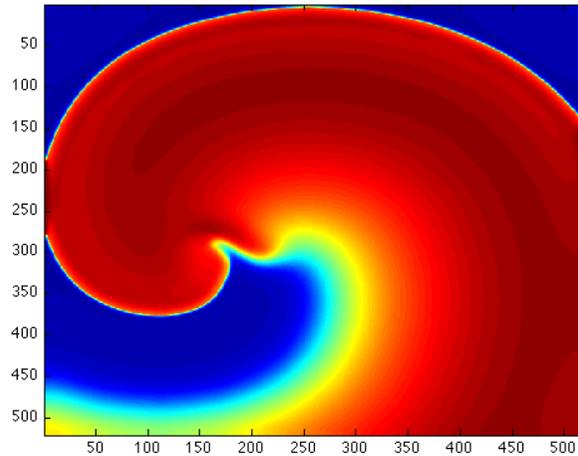
# Reaction optimized (Diffusion with Texture)

# Double vs Single


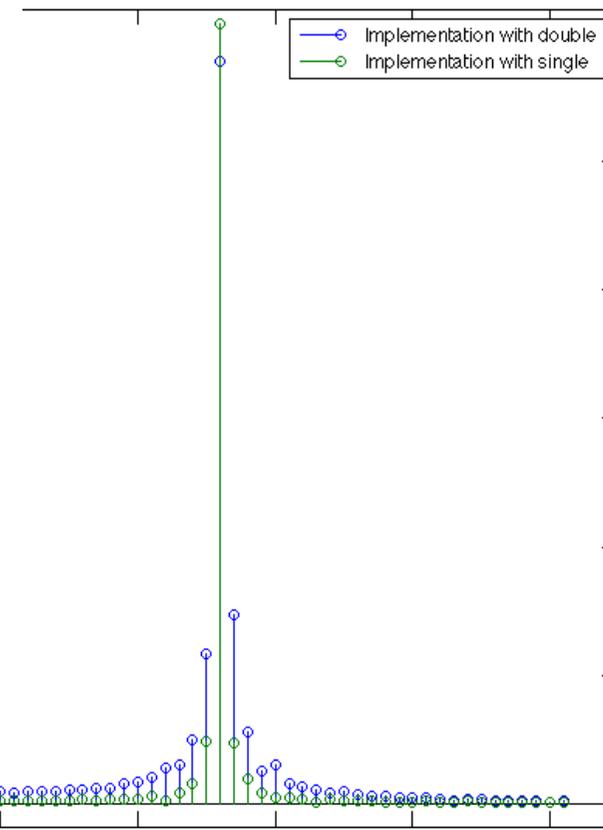
Nai
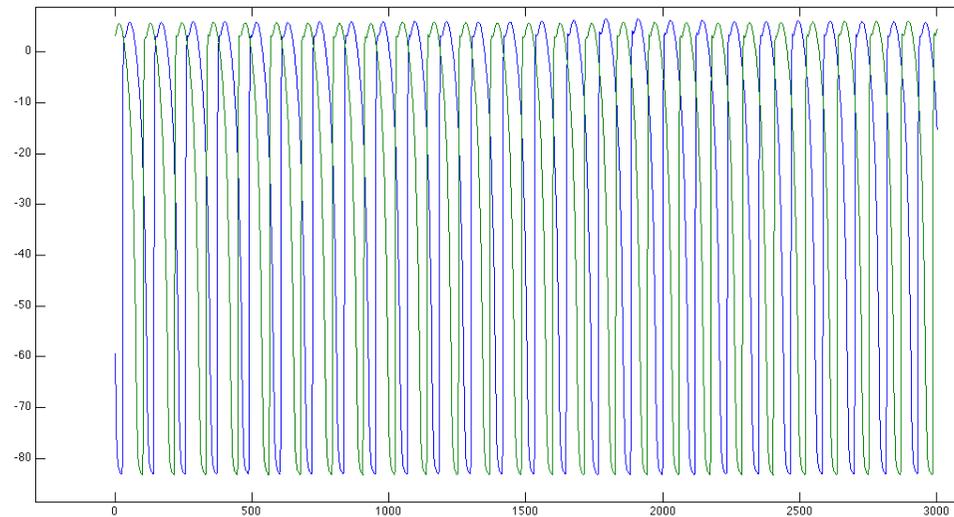
After 10 minutes
of simulation:

# Double vs Single



Ki

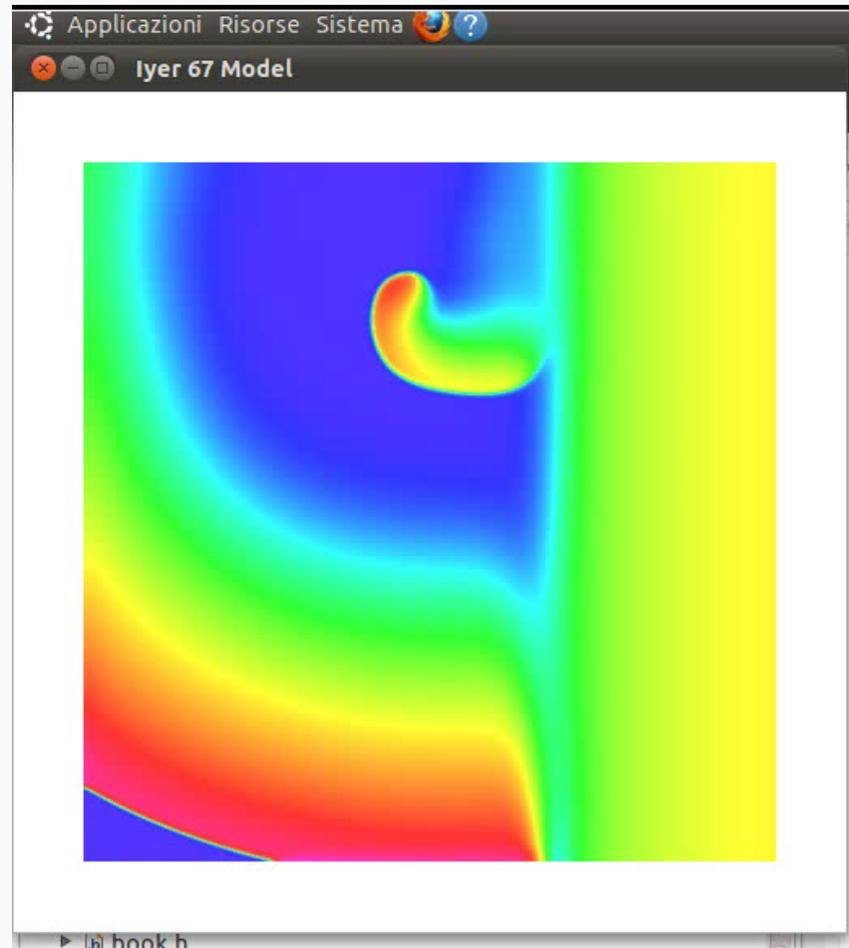After 10 minutes
of simulation:

# Double vs Single

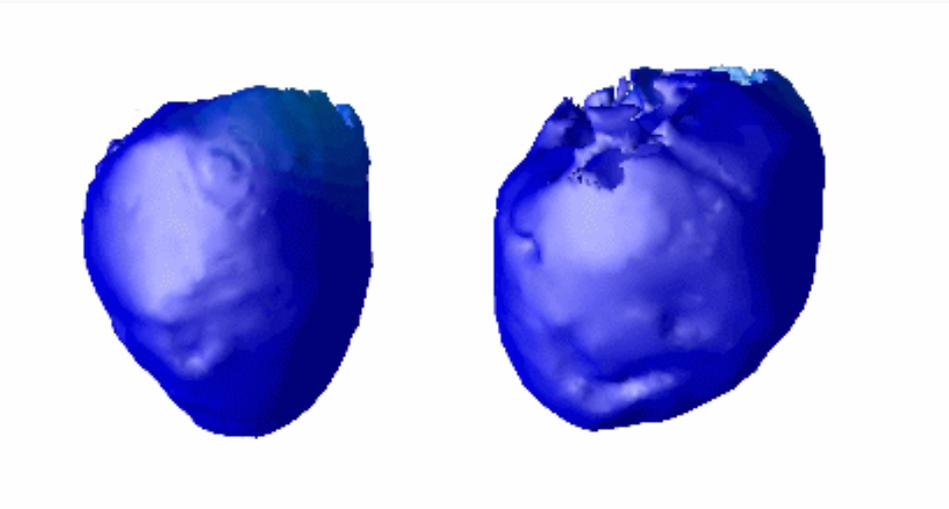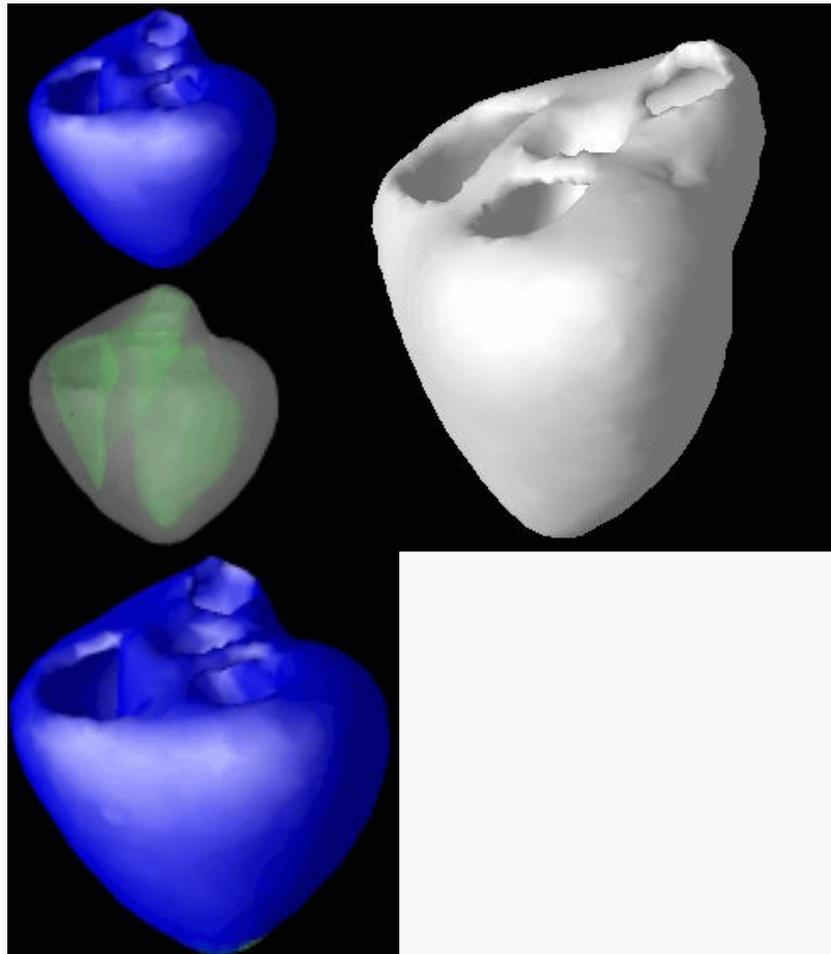After 10 minutes of simulation:

# Work in Progress Iyer Model (67V)

# Work in Progress 3D Models

# Conclusions

- **Many other challenge problems of CMACS expedition can take advantage of GPU technologies.**

- **The curve of developing of these technologies seems very promising for the future years.**

- **We are definitely interesting to collaborate with the other teams of the CMACS expedition in order to develop new revolutionary highly scalable GPU-based analysis tools for complex systems.**

# Thank you