



# Kuliah Umum Matematika: Pengenalan Komputasi Berbasis GPU

7 Oktober 2016 – FMIPA, Universitas Andalas

*Muhammad Teguh Satria, MSc*



**NovaGlobal Pte Ltd**  
Green & Scientific Computing Solutions



**PREFERRED  
SOLUTION  
PROVIDER**

## Muhammad Teguh Satria

*Application Systems Manager – NovaGlobal Pte Ltd, Singapore*

Email: [mteguhsat@gmail.com](mailto:mteguhsat@gmail.com) , [teguh@novaglobal.com.sg](mailto:teguh@novaglobal.com.sg)

### Pendidikan:

- ❑ S2, Computer Science, National Taipei University of Technology (Taipei Tech), Taiwan.
- ❑ S1, Matematika, Institut Teknologi Bandung, Indonesia.

### Pengalaman Penelitian:

- ❑ *GPU Performance Model* – Advanced Digital Sciences Center, Singapore.
- ❑ Implementasi GPU untuk:
  - *Attentive Teleconferencing Robot* – Advanced Digital Sciences Center, Singapore.
  - *Perkiraan Cuaca (WRF)* – SSEC, University of Wisconsin-Madison, US.
  - *Tsunami Model* – SSEC, University of Wisconsin-Madison, US.
  - *Geothermal* – TerraMath Indonesia
  - *Tsunami Simulation* – Taipei Tech, Taiwan.

# Pengenalan Komputasi Berbasis GPU



**Konsep Komputasi Paralel**

**Konsep Komputasi Berbasis GPU**

**Metode Komputasi Berbasis GPU**

**Contoh Kasus**



**Konsep Komputasi Paralel**

**Konsep Komputasi Berbasis GPU**

**Metode Komputasi Berbasis GPU**

**Contoh Kasus**

*Teknik mempersingkat komputasi dengan cara mengeksekusi instruksi-instruksi di dalamnya secara simultan.*

# Proses Komputasi di *Operating System*

## □ Vector addition

$$a + b = c$$

$$\begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \\ b_6 \\ b_7 \end{bmatrix} = \begin{bmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \\ c_6 \\ c_7 \end{bmatrix}$$

```
N = 8;
for( int i = 0; i < N; i++ )
{
    c[i] = a[i] + b[i];
}
```

Aplikasi: vector\_addition

*process\_0*

*start*

*thread\_0:*

- ambil elemen  $a_0$  di data  $a$
- ambil elemen  $b_0$  di data  $b$
- tambahkan  $a_0$  dan  $b_0$
- simpan hasilnya sebagai elemen  $c_0$  di data  $c$
- ambil elemen  $a_1$  di data  $a$
- ambil elemen  $b_1$  di data  $b$
- tambahkan  $a_1$  dan  $b_1$
- ...
- ...
- ...
- *exit*

*end*

# Definisi Dalam *Operating System*

## ❑ *process*

- ❑ = operasi yang dijalankan oleh *operating system*
- ❑ memiliki *virtual address space*
- ❑ memiliki unique id
- ❑ sedikitnya memiliki sebuah *thread*

## ❑ *thread*

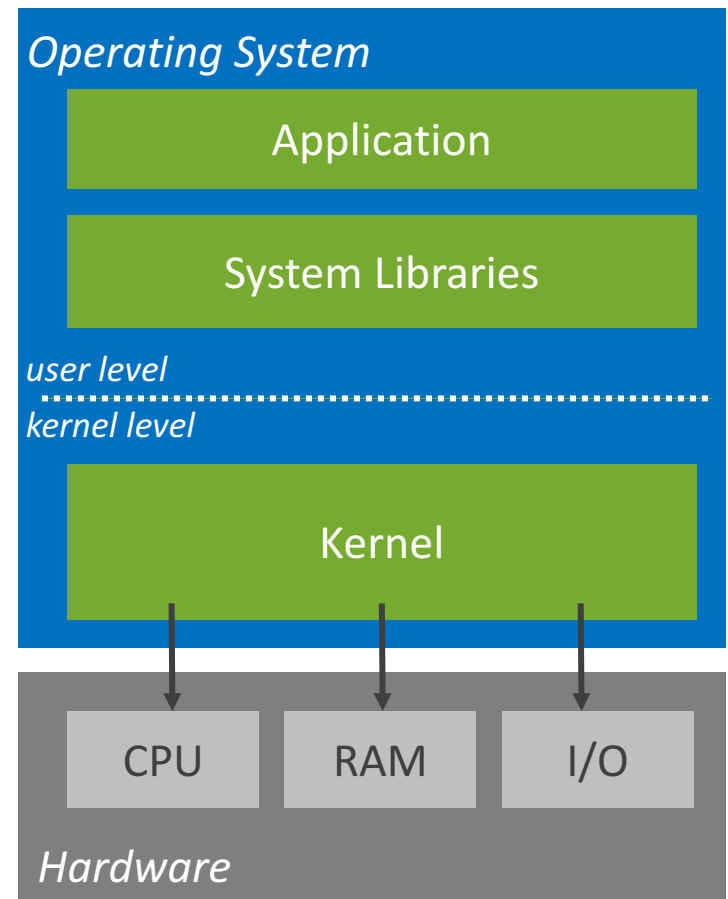
- ❑ = kumpulan perintah yang dieksekusi di processor
- ❑ memiliki unique id
- ❑ bagian dari *process*
- ❑ *thread-thread* dalam sebuah process berbagi *virtual address space*

## ❑ Reference

- ❑ [https://msdn.microsoft.com/en-us/library/windows/desktop/ms681917\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/desktop/ms681917(v=vs.85).aspx)

# Thread

- ❑ = kumpulan perintah yang dieksekusi di processor
- ❑ thread dibedakan dalam 2 tipe: *user thread* and *kernel thread*
- ❑ *user thread*
  - ❑ User level
  - ❑ Manageable
  - ❑ Handled by system library
  - ❑ Unlimited (depends on memory)
- ❑ *kernel thread*
  - ❑ kernel level
  - ❑ Unmanageable
  - ❑ Handled by kernel
  - ❑ Depends on number of cores or features in processor (hyperthreading)





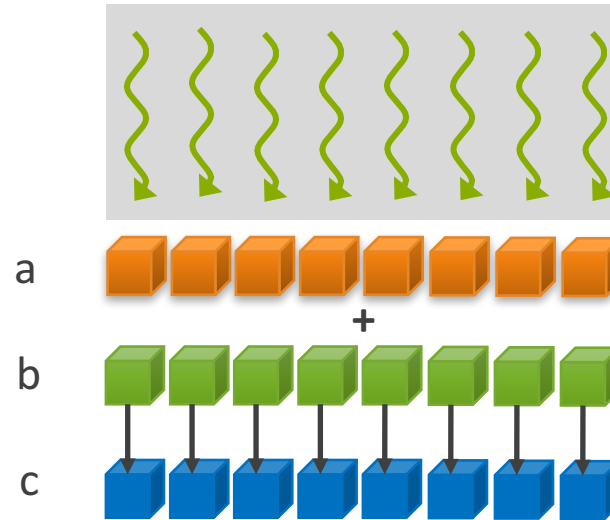
# Komputasi Paralel

## □ Vector addition

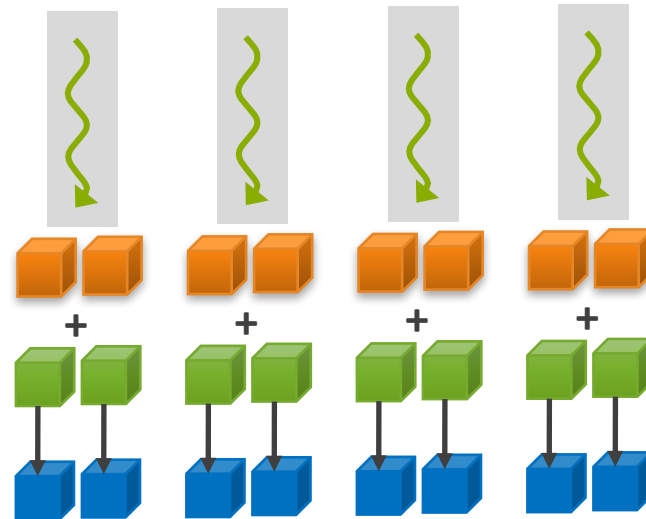
$$a + b = c$$

$$\begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \\ b_6 \\ b_7 \end{bmatrix} = \begin{bmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \\ c_6 \\ c_7 \end{bmatrix}$$

## □ Multithreading



## □ Multiprocessing



# Komputasi Parallel - Multithreading

*process\_0: vector addition  
start*

*thread\_0: create 7 more threads*

<i>thread_0:</i>	<i>thread_1:</i>	<i>thread_2:</i>	<i>thread_3:</i>	<i>thread_4:</i>	<i>thread_5:</i>	<i>thread_6:</i>	<i>thread_7:</i>
- ambil elemen $a_0$ di data $a$	- ambil elemen $a_1$ di data $a$	- ambil elemen $a_2$ di data $a$	- ambil elemen $a_3$ di data $a$	- ambil elemen $a_4$ di data $a$	- ambil elemen $a_5$ di data $a$	- ambil elemen $a_6$ di data $a$	- ambil elemen $a_7$ di data $a$
- ambil elemen $b_0$ di data $b$	- ambil elemen $b_1$ di data $b$	- ambil elemen $b_2$ di data $b$	- ambil elemen $b_3$ di data $b$	- ambil elemen $b_4$ di data $b$	- ambil elemen $b_5$ di data $b$	- ambil elemen $b_6$ di data $b$	- ambil elemen $b_7$ di data $b$
- tambahkan $a_0$ dan $b_0$	- tambahkan $a_1$ dan $b_1$	- tambahkan $a_2$ dan $b_2$	- tambahkan $a_3$ dan $b_3$	- tambahkan $a_4$ dan $b_4$	- tambahkan $a_5$ dan $b_5$	- tambahkan $a_6$ dan $b_6$	- tambahkan $a_7$ dan $b_7$
- simpan hasilnya sebagai elemen $c_0$ di data $c$	- simpan hasilnya sebagai elemen $c_1$ di data $c$	- simpan hasilnya sebagai elemen $c_2$ di data $c$	- simpan hasilnya sebagai elemen $c_3$ di data $c$	- simpan hasilnya sebagai elemen $c_4$ di data $c$	- simpan hasilnya sebagai elemen $c_5$ di data $c$	- simpan hasilnya sebagai elemen $c_6$ di data $c$	- simpan hasilnya sebagai elemen $c_7$ di data $c$
- exit	- exit	- exit	- exit	- exit	- exit	- exit	- exit

*end*

- 8 thread dieksekusi dalam waktu yang sama
- Tools: OpenMP, POSIX Thread, Win32 Thread

# Komputasi Parallel - Multiprocessing

Aplikasi: vector\_addition

<i>process_0</i> start	<i>process_1</i> start	<i>process_2</i> start	<i>process_3</i> start
<i>thread_0:</i> - ambil elemen $a_0$ di data $a$ - ambil elemen $b_0$ di data $b$ - tambahkan $a_0$ dan $b_0$ - simpan hasilnya sebagai elemen $c_0$ di data $c$ - ambil elemen $a_1$ di data $a$ - ambil elemen $b_1$ di data $b$ - tambahkan $a_1$ dan $b_1$ - simpan hasilnya sebagai elemen $c_1$ di data $c$ - exit	<i>thread_0:</i> - ambil elemen $a_0$ di data $a$ - ambil elemen $b_0$ di data $b$ - tambahkan $a_0$ dan $b_0$ - simpan hasilnya sebagai elemen $c_0$ di data $c$ - ambil elemen $a_1$ di data $a$ - ambil elemen $b_1$ di data $b$ - tambahkan $a_1$ dan $b_1$ - simpan hasilnya sebagai elemen $c_1$ di data $c$ - exit	<i>thread_0:</i> - ambil elemen $a_0$ di data $a$ - ambil elemen $b_0$ di data $b$ - tambahkan $a_0$ dan $b_0$ - simpan hasilnya sebagai elemen $c_0$ di data $c$ - ambil elemen $a_1$ di data $a$ - ambil elemen $b_1$ di data $b$ - tambahkan $a_1$ dan $b_1$ - simpan hasilnya sebagai elemen $c_1$ di data $c$ - exit	<i>thread_0:</i> - ambil elemen $a_0$ di data $a$ - ambil elemen $b_0$ di data $b$ - tambahkan $a_0$ dan $b_0$ - simpan hasilnya sebagai elemen $c_0$ di data $c$ - ambil elemen $a_1$ di data $a$ - ambil elemen $b_1$ di data $b$ - tambahkan $a_1$ dan $b_1$ - simpan hasilnya sebagai elemen $c_1$ di data $c$ - exit
end	end	end	end
address space	address space	address space	address space

- Dikenal dengan *Message Passing Interface (MPI)*
- Tools: OpenMPI, MPICH



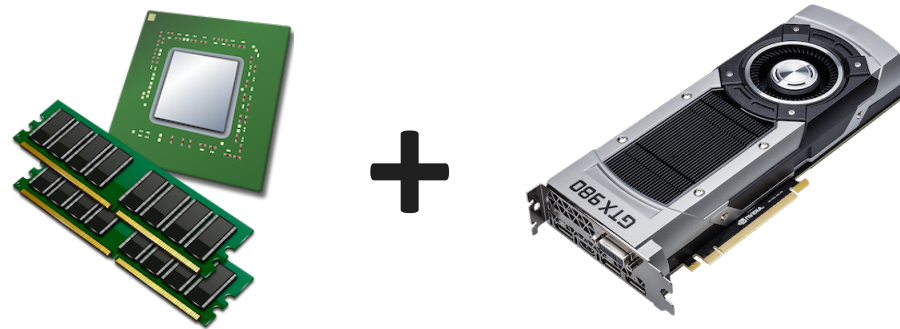
**Konsep Komputasi Paralel**

**Konsep Komputasi Berbasis GPU**

**Metode Komputasi Berbasis GPU**

**Contoh Kasus**

# Komputasi Berbasis GPU



- ❑ GPU sebagai CPU co-processor
- ❑ Traditional GPU Computing: Shader Programming
  - ❑ Cg, developed by NVIDIA
  - ❑ GLSL, developed by Khronos Group.
  - ❑ BrookGPU, developed by Stanford University
- ❑ Modern GPU Computing Platform
  - ❑ CUDA, proposed and developed by NVIDIA
  - ❑ OpenCL, proposed by Apple, developed and maintained by Khronos Group.
  - ❑ DirectCompute, developed by Microsoft, bundled in DirectX SDK.
- ❑ *Catatan: materi selanjutnya merujuk pada arsitektur di NVIDIA GPU dan CUDA platform.*

\*image source: [geforce.com](http://geforce.com), [freeiconspng.com](http://freeiconspng.com), [pixabay.com](http://pixabay.com)

## CPU

Optimized for  
Serial Tasks

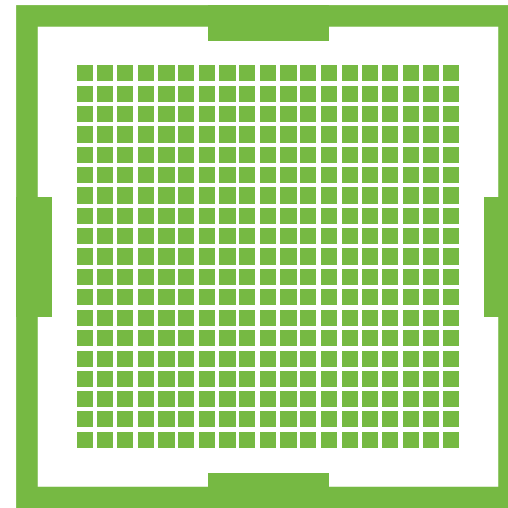


Low Latency

+

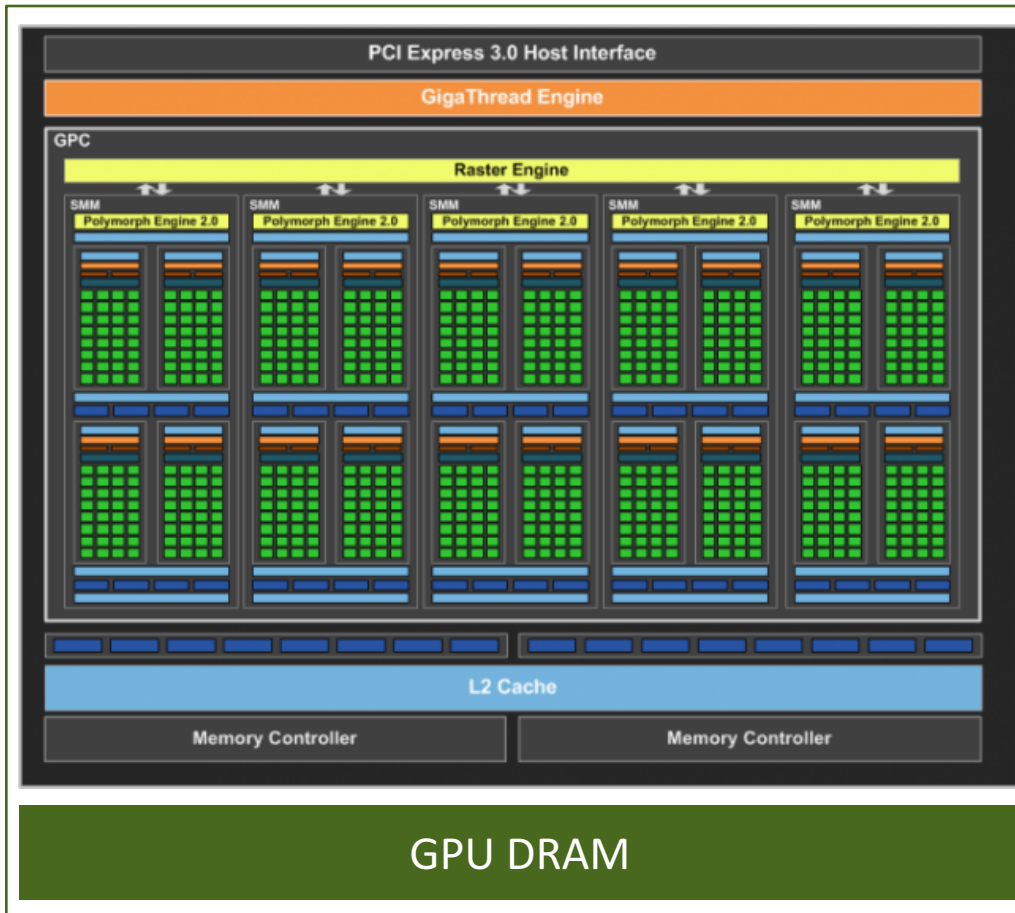
## GPU Accelerator

Optimized for Many  
Parallel Tasks



High Throughput

# NVIDIA GPU Architecture – Maxwell



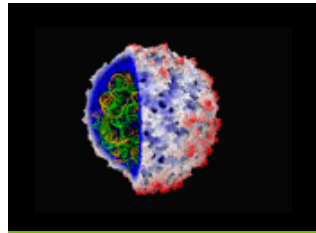
- ❑ Streaming Multiprocessor (SMM)
- ❑ L2 Cache for GPU Memory caching
- ❑ GPU DRAM (Global Memory)
- ❑ On-chip Memory

## GPUs ACCELERATE SCIENCE



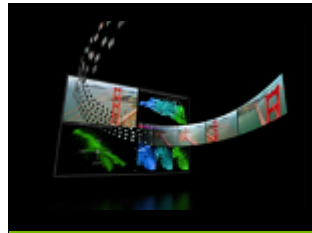
**146X**

Medical imaging  
U of Utah



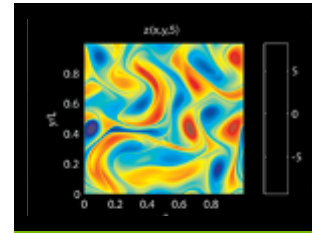
**36X**

Molecular Dynamics  
U of Illinois, Urbana



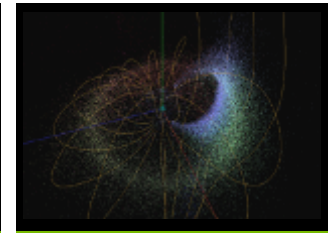
**18X**

Video Transcoding  
Elemental Tech



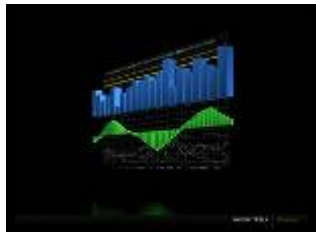
**50X**

Matlab Computing  
AccelerEyes



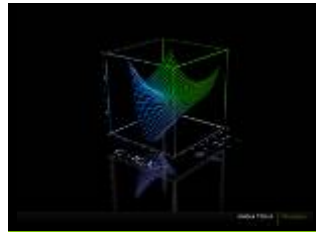
**100X**

Astrophysics  
RIKEN



**149X**

Financial Simulation  
Oxford



**47X**

Linear Algebra  
Universidad Jaime



**20X**

3D Ultrasound  
Techniscan



**130X**

Quantum Chemistry  
U of Illinois, Urbana

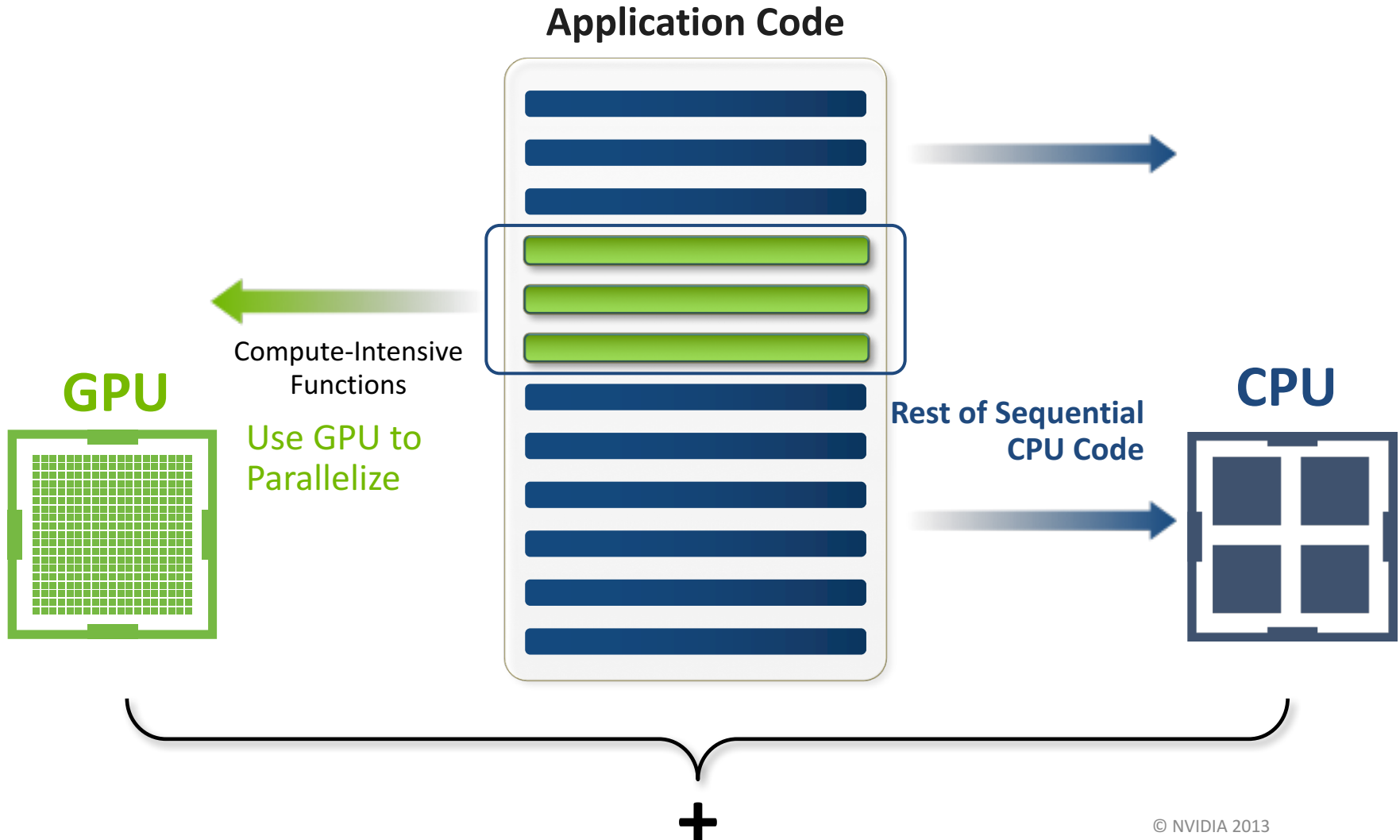


**30X**

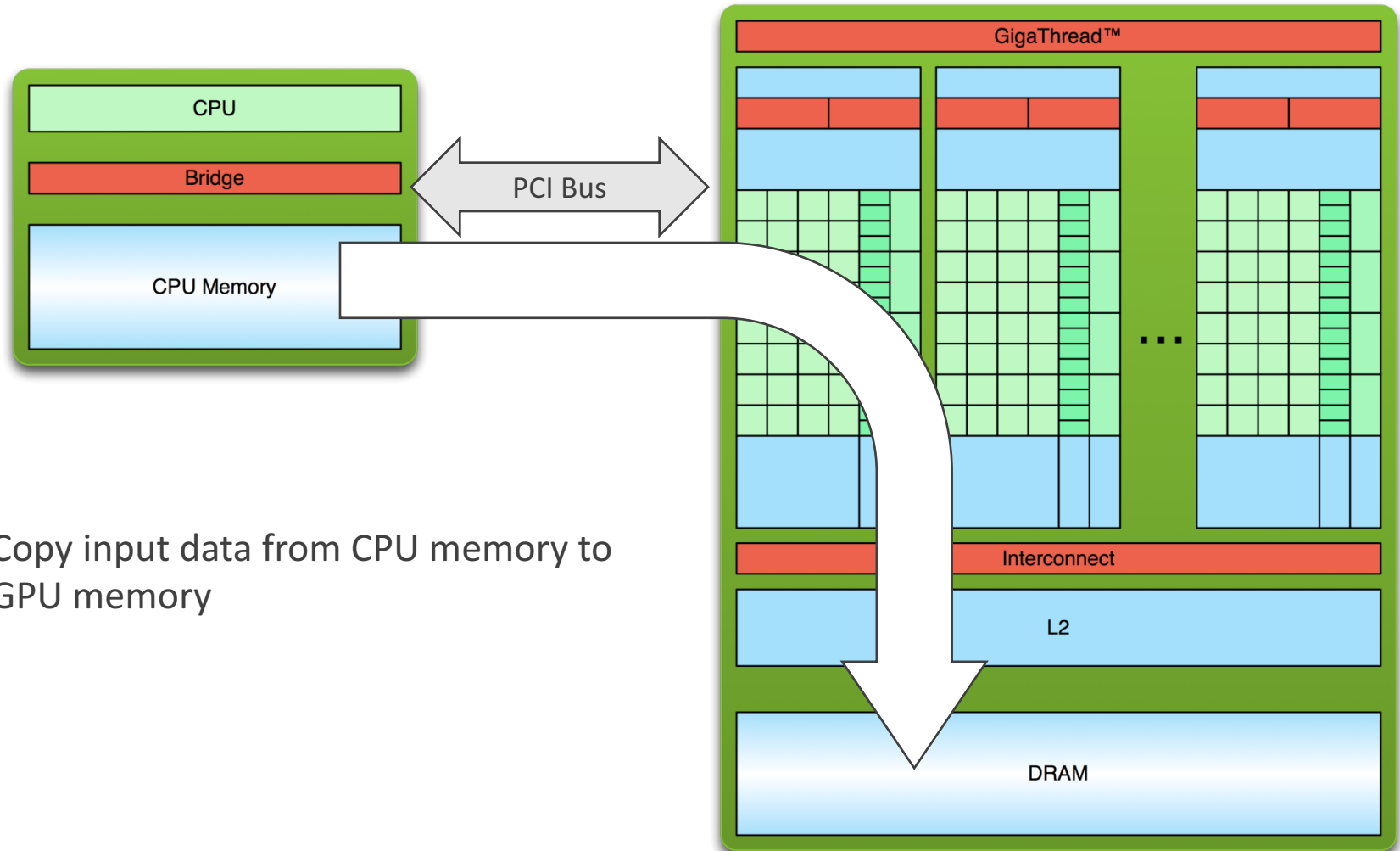
Gene Sequencing  
U of Maryland



# Konsep Implementasi *GPU Computing* – 1



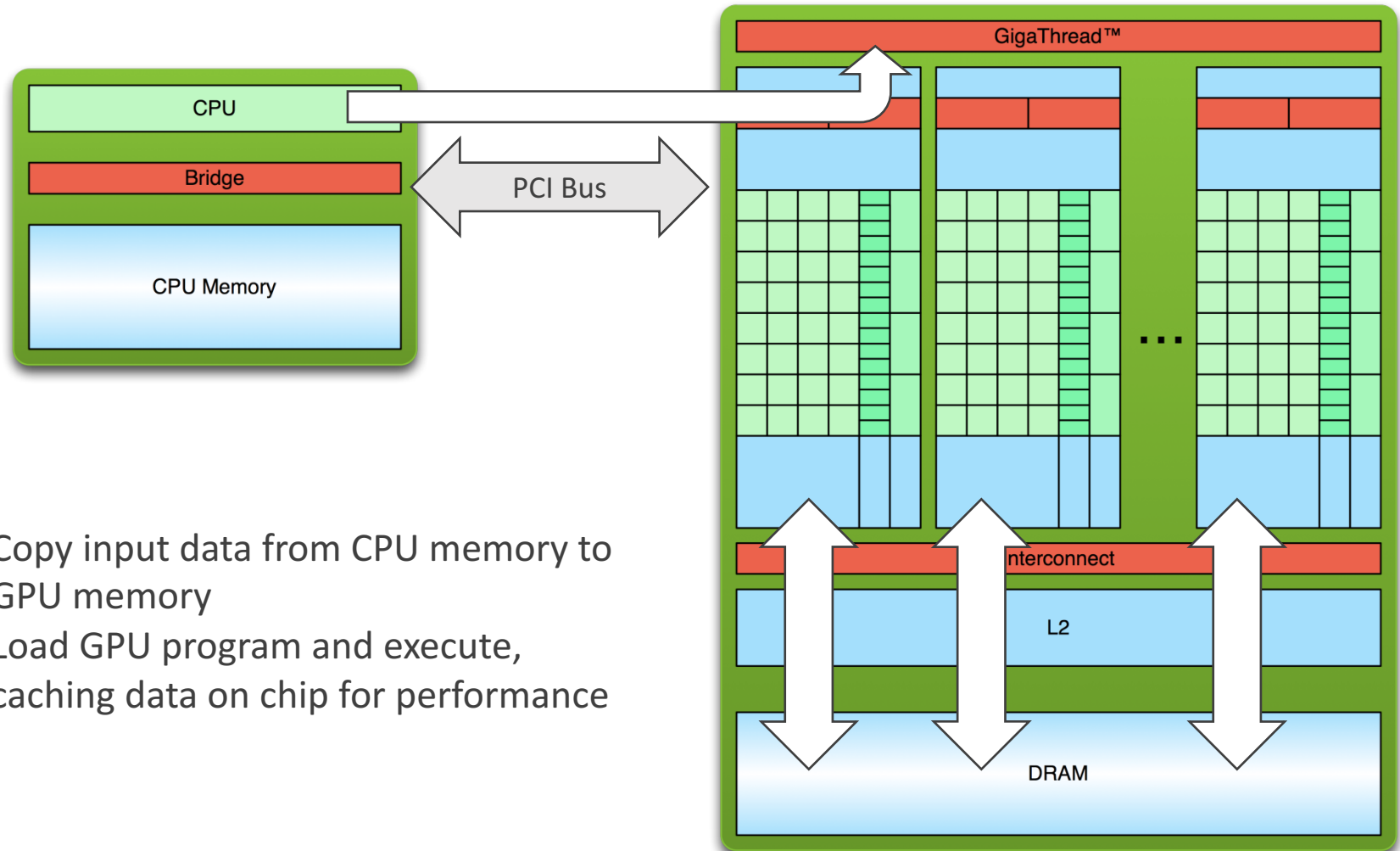
# Konsep Implementasi *GPU Computing* – 2



1. Copy input data from CPU memory to GPU memory

*Contents credit to NVIDIA*

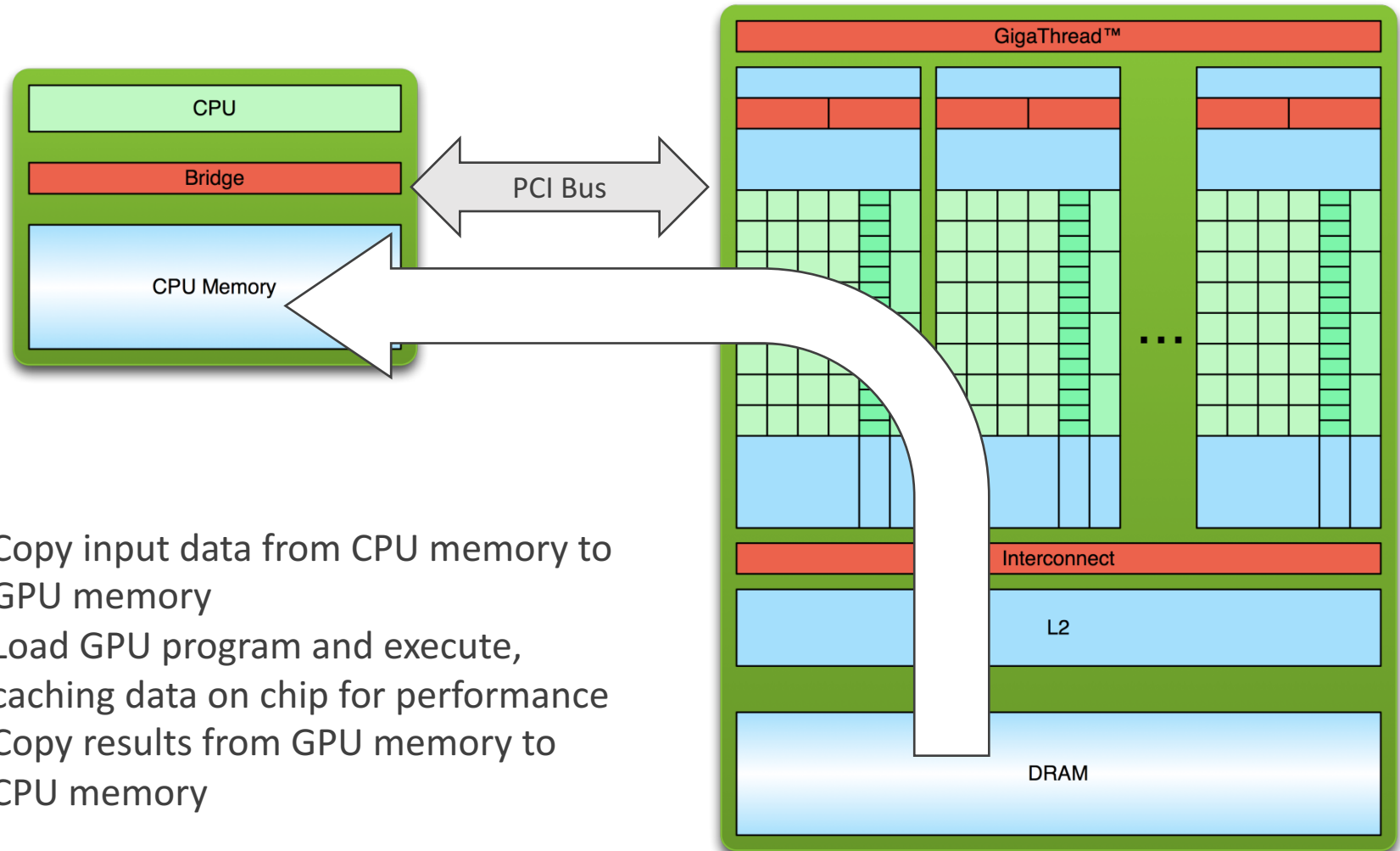
# Konsep Implementasi *GPU Computing* – 3



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

Contents credit to NVIDIA

# Konsep Implementasi *GPU Computing* – 4



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
3. Copy results from GPU memory to CPU memory

*Contents credit to NVIDIA*

# Konsep Implementasi *GPU Computing* – 5

## CPU

### *CPU process*

#### *CPU thread*

*Allocate space in GPU memory  
Copy data from CPU memory to  
GPU memory*

*Create GPU process with X GPU  
threads and execute it*

*Copy output from GPU memory to  
CPU memory*

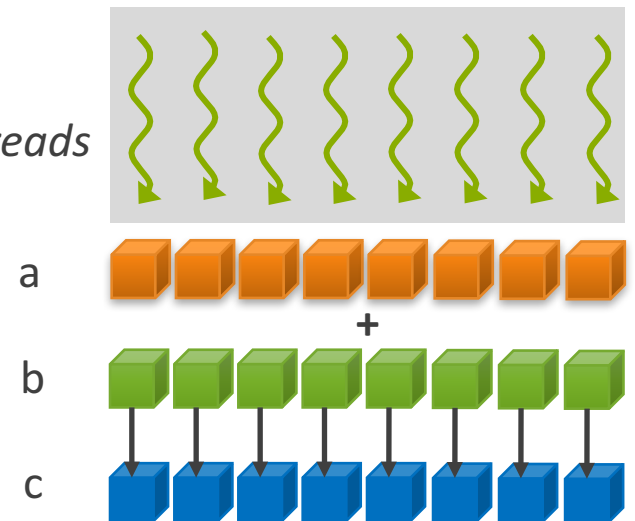
*Clear GPU memory*



## GPU

### *GPU process*

#### *GPU threads*





- ..... **Konsep Komputasi Paralel**

- ..... **Konsep Komputasi Berbasis GPU**

- ..... **Metode Komputasi Berbasis GPU**

- ..... **Contoh Kasus**

# Hal-hal Yang Dibutuhkan

- ❑ GPU
- ❑ GPU Driver
- ❑ CUDA Toolkit

# Cara Memanfaatkan *GPU Computing*

Owned Application

**Libraries**

**“Drop-in”  
Acceleration**

**OpenACC  
Directives**

**Easily  
Accelerate  
Applications**

**Programming  
Languages**

**Maximum  
Flexibility**

**GPU-ready Applications**

**MATLAB, Mathematica, dll**

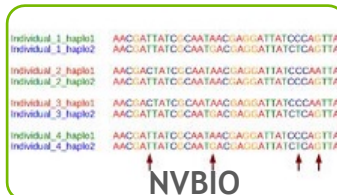


## GPU ACCELERATED LIBRARIES

“Drop-in” Acceleration for Your Applications

### Domain-specific

Deep Learning, GIS, EDA,  
Bioinformatics, Fluids



### Visual Processing

Image & Video



### Linear Algebra

Dense, Sparse, Matrix



### Math Algorithms

AMG, Templates, Solvers



[developer.nvidia.com/gpu-accelerated-libraries](http://developer.nvidia.com/gpu-accelerated-libraries)

# Directive-based Programming Model

## OPENACC

More Science, Less Programming

```
main()
{
    <serial code>
    #pragma acc kernels
    // automatically runs
    // on GPUs
    {
        <parallel code>
    }
}
```

**SIMPLE**

Minimum efforts  
Small code modifications

**POWERFUL**

Up to 10x faster application  
performance

**PORTABLE**

Optimize once,  
run on GPUs and CPUs

Free workstation license for academia

[developer.nvidia.com/openacc](http://developer.nvidia.com/openacc)

## PROGRAMMING LANGUAGES

### Numerical analytics

▶ MATLAB, Mathematica, LabVIEW

### Fortran

▶ OpenACC, CUDA Fortran

### C

▶ OpenACC, CUDA C

### C++

▶ Thrust, CUDA C++

### Python

▶ Anaconda, PyCUDA

### C#

▶ Alea, Hybridizer

[developer.nvidia.com/language-solutions](http://developer.nvidia.com/language-solutions)



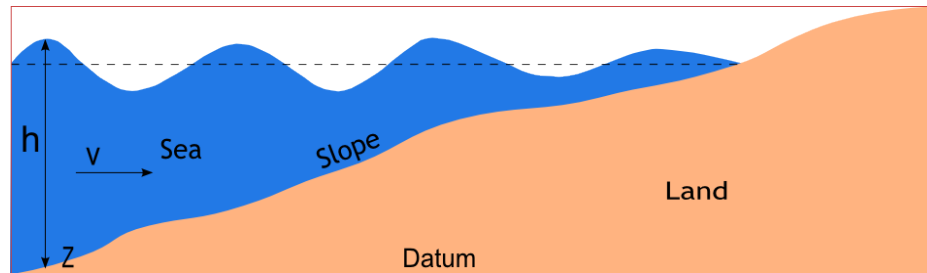
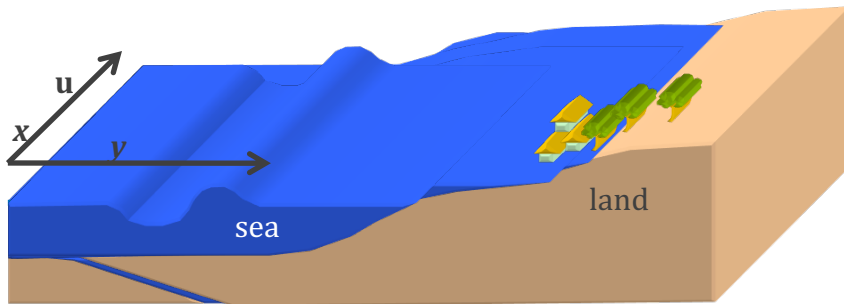
**Konsep Komputasi Paralel**

**Konsep Komputasi Berbasis GPU**

**Bagaimana Menggunakannya?**

**Contoh Kasus**

# Tsunami Propagation Model – 1



$$h_t + (hu)_x + (hv)_y = 0$$

$$(uh)_t + \left(u^2h + \frac{1}{2}gh^2\right)_x + (uhv)_y = -gh(So_x - Sf_x)$$

$$(vh)_t + (uhv)_x + \left(v^2h + \frac{1}{2}gh^2\right)_y = -gh(So_y - Sf_y)$$

$$Sf_x = \frac{\eta^2 u \sqrt{u^2 + v^2}}{h^{\frac{4}{3}}} \quad Sf_y = \frac{\eta^2 v \sqrt{u^2 + v^2}}{h^{\frac{4}{3}}}$$

$\eta$  is equal to 0.05, as a representation of light brush surface over floodplain

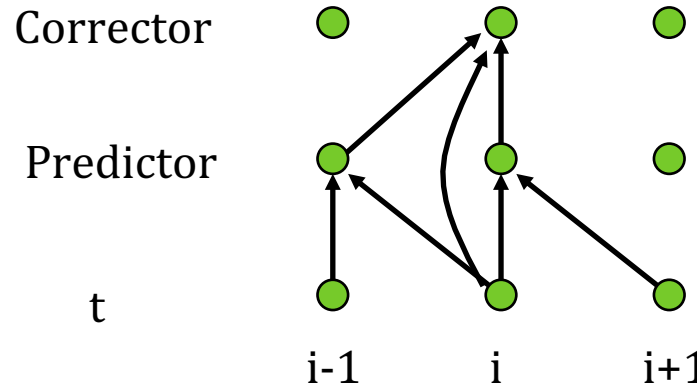
$$\frac{\partial h}{\partial t} + \frac{\partial(hu)}{\partial x} + \frac{\partial(hv)}{\partial y} = 0$$

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} + g \frac{\partial h}{\partial x} = -g(So_x - Sf_x)$$

$$\frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} + g \frac{\partial h}{\partial y} = -g(So_y - Sf_y)$$

# Tsunami Propagation Model – 2

## ❖ Scheme



## ❖ Predictor

$$h_{i,j}^{t+1} = h_{i,j}^t - \frac{\Delta t}{\Delta x} (u_{i+1,j}^t h_{i+1,j}^t - u_{i,j}^t h_{i,j}^t) - \frac{\Delta t}{\Delta y} (v_{i,j+1}^t h_{i,j+1}^t - v_{i,j}^t h_{i,j}^t)$$

$$u_{i,j}^{t+1} = u_{i,j}^t - u_{i,j}^t \frac{\Delta t}{\Delta x} (u_{i+1,j}^t - u_{i,j}^t) - v_{i,j}^t \frac{\Delta t}{\Delta y} (u_{i,j+1}^t - u_{i,j}^t) - g \frac{\Delta t}{\Delta x} (h_{i+1,j}^t - h_{i,j}^t) - g \Delta t (S_{o_x} - S_{f_x})$$

$$v_{i,j}^{t+1} = v_{i,j}^t - u_{i,j}^t \frac{\Delta t}{\Delta x} (v_{i+1,j}^t - v_{i,j}^t) - v_{i,j}^t \frac{\Delta t}{\Delta y} (v_{i,j+1}^t - v_{i,j}^t) - g \frac{\Delta t}{\Delta y} (h_{i,j+1}^t - h_{i,j}^t) - g \Delta t (S_{o_y} - S_{f_y})$$

## ❖ Corrector

$$h_{i,j}^{t+1} = \frac{(h_{i,j}^t + h_{i,j}^{t+1})}{2} - \frac{\Delta t}{2\Delta x} (u_{i,j}^{t+1} h_{i,j}^{t+1} - u_{i-1,j}^{t+1} h_{i-1,j}^{t+1}) \\ - \frac{\Delta t}{2\Delta y} (v_{i,j}^{t+1} h_{i,j}^{t+1} - v_{i,j-1}^{t+1} h_{i,j-1}^{t+1})$$

$$u_{i,j}^{t+1} = \frac{(u_{i,j}^t + u_{i,j}^{t+1})}{2} - u_{i,j}^{t+1} \frac{\Delta t}{2\Delta x} (u_{i,j}^{t+1} - u_{i-1,j}^{t+1}) - v_{i,j}^{t+1} \frac{\Delta t}{2\Delta y} (u_{i,j}^{t+1} - u_{i,j-1}^{t+1}) \\ - g \frac{\Delta t}{2\Delta x} (h_{i,j}^{t+1} - h_{i-1,j}^{t+1}) - g \frac{\Delta t}{2} (So_x - Sf_x)$$

$$v_{i,j}^{t+1} = \frac{(v_{i,j}^t + v_{i,j}^{t+1})}{2} - u_{i,j}^{t+1} \frac{\Delta t}{2\Delta x} (v_{i,j}^{t+1} - v_{i-1,j}^{t+1}) - v_{i,j}^{t+1} \frac{\Delta t}{2\Delta y} (v_{i,j}^{t+1} - v_{i,j-1}^{t+1}) \\ - g \frac{\Delta t}{2\Delta y} (h_{i,j}^{t+1} - h_{i,j-1}^{t+1}) - g \frac{\Delta t}{2} (So_y - Sf_y)$$

# Tsunami Propagation Model – 4

## ❖ Boundary Condition

- ❑ Neumann boundary → for  $h$  and velocity @corresponding axis.
- ❑ Copy boundary → for the velocity perpendicular to the axis.
- ❑ Wet-dry condition → for the wave front → by giving a tolerance value.  
If  $h$  is lower than the tolerance, then it set to be dry.

## ❖ Hansen filtering method

$$F_{i,j} = cF_{i,j} + (1 - c)(0.25)(F_{i-1,j} + F_{i+1,j} + F_{i,j-1} + F_{i,j+1})$$

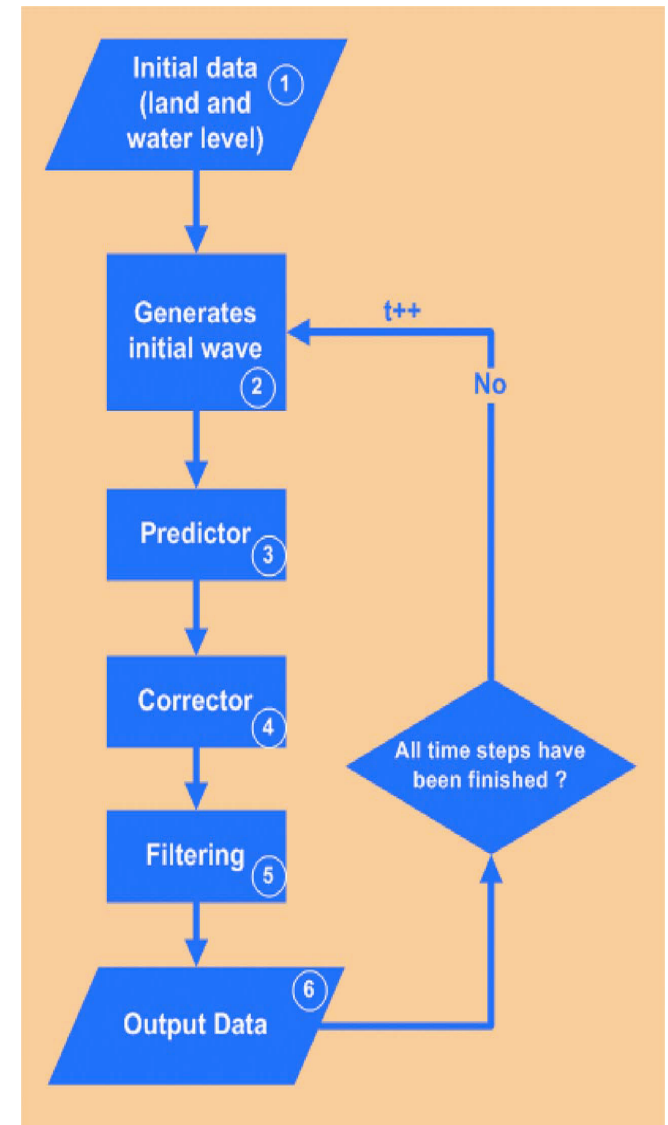
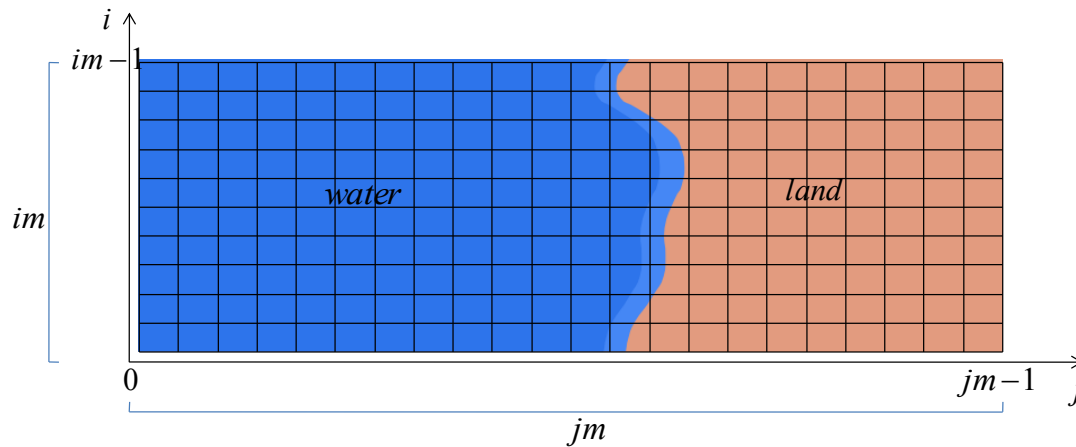
where  $F$  represents the values of  $h$ ,  $u$ , and  $v$ , and  $c$  is equal to 0.99.

## ❖ Wave Generation

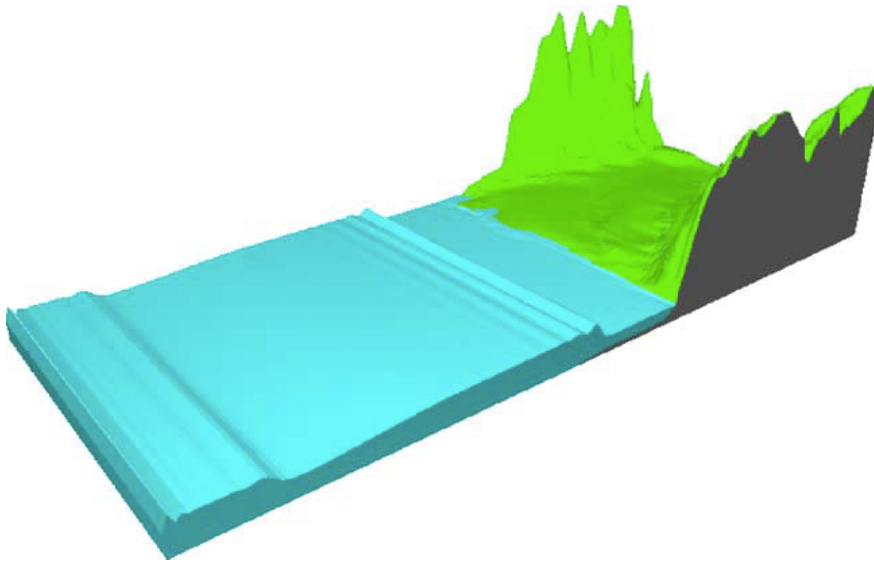
- ❖  $h$  is generated by an interpolation of water level data → *spline interpolation*.
- ❖  $u = 0$ .
- ❖  $v = \sqrt{gh}$



# Simulation Flow



# Performance



Performance on NVIDIA Fermi C2050 GPU with 448 cores

<b>Number of Time Steps</b>	<b>Execution Time (ms)</b>	<b>Speedup *</b>
100	0.80	422x
500	3.68	473x
1000	11.47	223x

\* Compared to single core on Intel Xeon E5620 2.4 GHz



**NovaGlobal Pte Ltd**  
Green & Scientific Computing Solutions

<http://novaglobal.com.sg>

- ❑ Established in 1996
- ❑ Operating in ASEAN region
  - ❑ Primarily in Singapore, Malaysia, Indonesia, Thailand and Vietnam
- ❑ Solution provider for
  - ❑ High Performance Computing
  - ❑ Accelerated GPU Computing
  - ❑ Cloud/Virtualization
- ❑ Platform-Independent
- ❑ Partnering with Technology Leaders



Intel® Solutions for Lustre\* Reseller

\* Other names and brands may be claimed as the property of others

