

# GPU solution of optimized wavenumber model for steady incompressible Navier-Stokes equations

**Tony Wen-Hann Sheu**

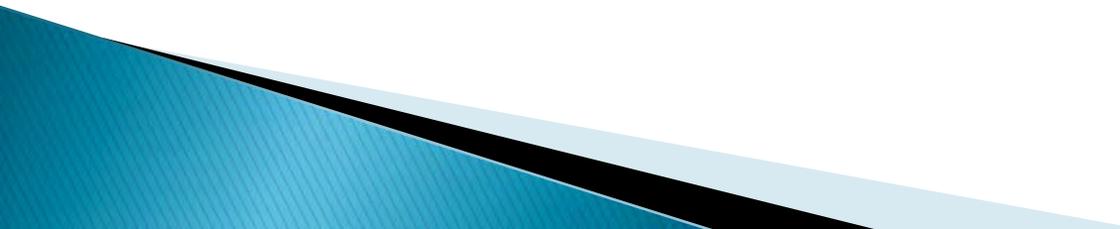
- Department of Engineering Science and Ocean Engineering, National Taiwan University
- Department of Mathematics, National Taiwan University
- Center for Advanced Studies on Theoretical Sciences (CASTS), National Taiwan University



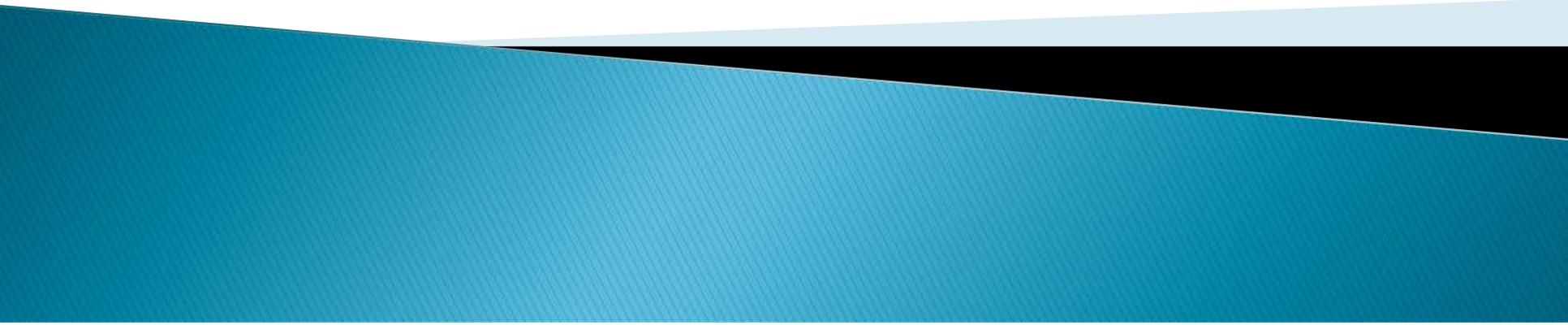
CASTS-LJLL workshop on Applied Mathematics and Mathematical Sciences, Taipei, May 26-29, 2014

Acknowledgement : Neo Shih-Chao Kao, Wei-Chung Wang

# Contents

- I. Computational challenges in solving 3D steady incompressible Navier-Stokes equations by finite element method
  - II. In-house developed finite element code
  - III. Code implementation in Kepler K20 GPU
  - IV. Verification and numerical results
  - V. Concluding remarks
- 

# **I. Computational Challenges in solving 3D steady incompressible Navier-Stokes equations by finite element method**



# Mathematical model

Three-dimensional steady incompressible **Navier-Stokes equations** in **primitive variables**  $(\underline{u}, p)$  are as follows :

Elliptic Incompressible Navier-Stokes equations

$$\left\{ \begin{array}{l} \boxed{(\underline{u} \cdot \nabla) \cdot \underline{u}} + \boxed{\nabla p} = \boxed{\frac{1}{Re} \nabla^2 \underline{u}} + \boxed{\underline{f}(\underline{x})} \\ \text{nonlinear convective term} \quad \text{Pressure gradient} \quad \text{Viscous term} \quad \text{Source term} \end{array} \right. \quad (1)$$

$$\boxed{\nabla \cdot \underline{u} = 0} \quad \text{Divergence-free constrain condition} \quad (2)$$

$$\underline{u}(\underline{x}, t) = \underline{u}_\Gamma \quad \text{on } \partial\Omega \quad (3)$$

The **Reynolds number**  $Re = \rho u_{ref} L_{ref} / \mu$  comes out as the result of normalization

- The **mixed** FEM formulation is adopted to get the solution  $(\underline{u}, p)$  simultaneously



Claude-Louis Navier (1785-1836)



George Gabriel Stokes, (1819-1903)

# Computational challenges in solving $(\underline{u}, p)$ solutions from (1-3)-Motivation

- C1 – Rigorous implementation of boundary condition and assurance of divergence-free constrain for the elliptic differential system of equations
- C2 - Accurate approximation of convection term  $\underline{u} \cdot \nabla \underline{u}$  in fixed grid to avoid oscillatory velocity field
- C3- Proper storage of primitive variables  $\underline{u}$  and  $p$  to avoid even-odd pressure solution
- C4- Effective calculation of  $(\underline{u}, p)$  from the unsymmetric and indefinite matrix equation using modern iterative solver
- C5- Parallel implementation of finite element program

# Resolving computational challenges C1~C5

## -Objectives

- O1- Adopt mixed formulation of equations cast in primitive variables
- O2- Derive a streamline upwinding Petrov-Galerkin formulation featuring the numerical wavenumber error reducing property for the convection term
- O3- Formulate finite element equations in weak sense in trilinear pressure /triquadratic element to satisfy the LBB compatibility condition
- O4- Apply preconditioner to the normalized symmetric and positive definite matrix equations to ensure fast unconditional convergence for high Reynolds number simulation
- O5- Implement CUDA programming finite element code in hybrid CPU-GPU (Graphic Process Unit ) platform to largely reduce the simulation time

## II. In-house developed finite element code

(Numerical Heat Transfer ; in press )



# Tri-quadratic wavenumber optimized finite element model

The proposed **Petrov-Galerkin** finite element model is as follows

$$\int_{\Omega} (\underline{\mathbf{u}} \cdot \nabla \underline{\mathbf{u}}) \underline{\mathbf{w}} d\Omega + \frac{1}{\text{Re}} \int_{\Omega} \nabla \underline{\mathbf{u}} : \nabla \underline{\mathbf{w}} - \int_{\Omega} \mathbf{p} \nabla \cdot \underline{\mathbf{w}} d\Omega = \int_{\Omega} \underline{\mathbf{f}} \underline{\mathbf{w}} d\Omega$$
$$\int_{\Omega} (\nabla \cdot \underline{\mathbf{u}}) \mathbf{q} d\Omega = 0$$

where  $\underline{\mathbf{w}}$  and  $\mathbf{q}$  are the **weighting functions** for the **momentum** and **continuity equations**, respectively

The weighted residual weak formulation for the convection term in one-dimensional

$$\sum_{e=1} \int_{\Omega^e} \Phi_x w_i d\Omega^e$$

- Galerkin model :  $w_i = N_i$

- Wavenumber optimized model :  $w_i = N_i + \mathbf{B}_i$

$\left. \begin{array}{l} \mathbf{B}_i = \tau u \frac{\partial N_i}{\partial x} \\ \tau = \frac{\delta u H}{2|u|^2} \end{array} \right\}$

biased part

Substituting the quadratic interpolation function into the above weak statement to derive the discrete equations at the center and corner nodes, respectively

The discretization equation for the convection term in a quadratic element is expressed as follows :

$$\Phi_x \approx \frac{1}{h} (a_1 \Phi_{i-1} + a_2 \Phi_i + a_3 \Phi_{i+1})$$

$$\begin{cases} a_1 = -\frac{1}{2} - \delta \\ a_2 = 2\delta \\ a_3 = \frac{1}{2} + \delta \end{cases}$$

To minimize the numerical **wavenumber error**, the **Fourier transform** and its inverse

$$\tilde{\Phi}(\alpha) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \Phi(x) \exp(-i\alpha x) dx$$

$$\Phi(x) = \int_{-\infty}^{\infty} \tilde{\Phi}(\alpha) \exp(i\alpha x) d\alpha$$

are applied to derive the **numerical wavenumber**

$$\alpha \approx \frac{-i}{h} \left[ a_1 \exp(-i\alpha h) + a_2 + a_3 \exp(i\alpha h) \right]$$

- ▶ The **actual wavenumber** can be regarded as the right hand side of the numerical wavenumber

$$\tilde{\alpha} = \frac{-\mathbf{i}}{h} \left[ a_1 \exp(-\mathbf{i}\alpha h) + a_2 + a_3 \exp(\mathbf{i}\alpha h) \right]$$

- ▶ We define the error function to minimize the discrepancy between  $\alpha$  and  $\tilde{\alpha}$

$$E = \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} |\tilde{\alpha} - \alpha|^2 d(\alpha h) = \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} |\tilde{\gamma} - \gamma|^2 d\gamma$$

- ▶ The limiting condition  $\frac{\partial E}{\partial a_3} = 0$  is enforced to get the expression

of  $\delta$  as  $\tau = \frac{\delta u H}{2|u|^2}$

$$\delta = \begin{cases} \frac{1}{2} & ; \text{ at end node} \\ \frac{8-3\pi}{-22+6\pi} & ; \text{ at center node} \end{cases}$$

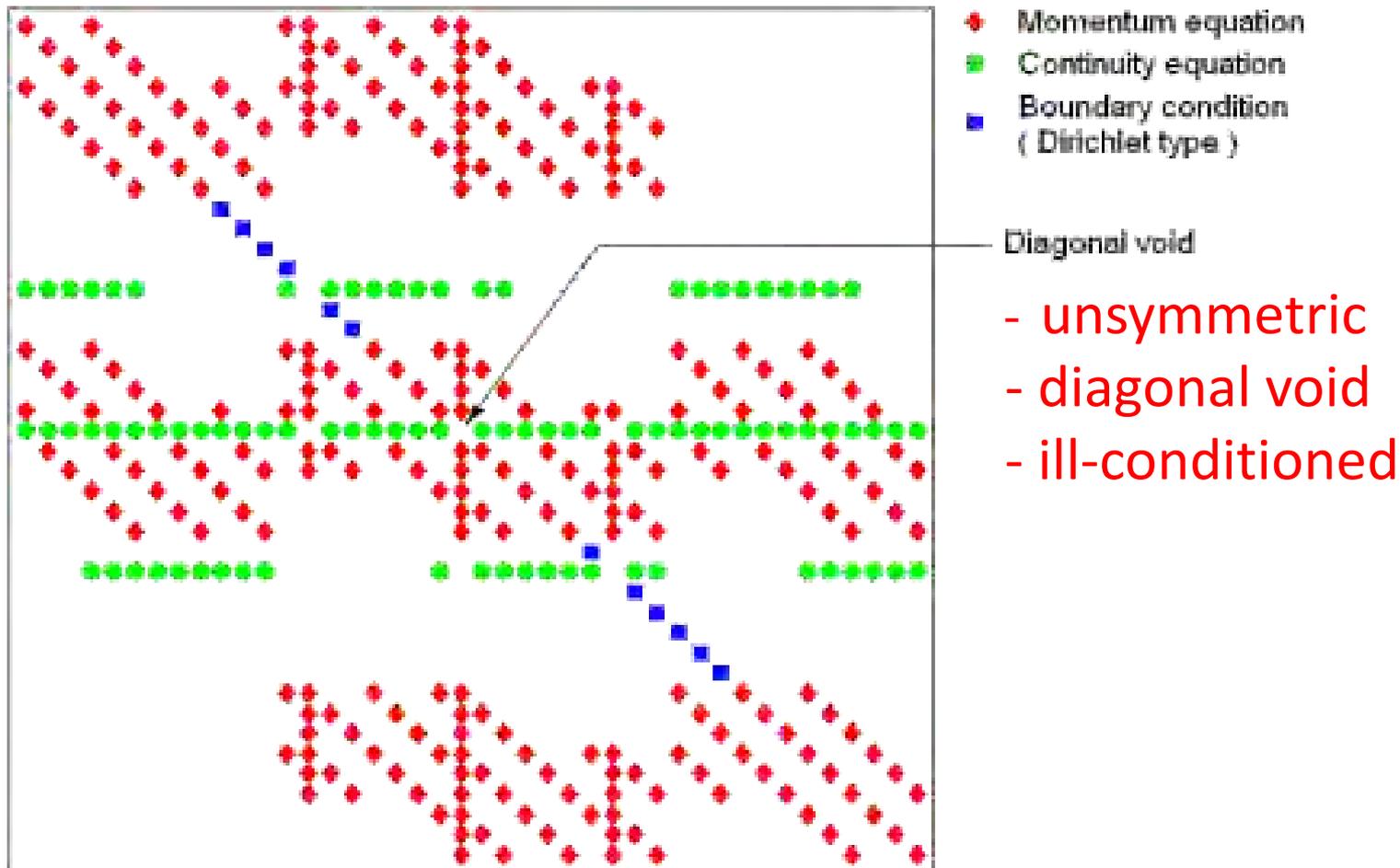
- ▶ In multi-dimensional problem, the method of **operator splitting** is adopted to get the stabilization term along the **streamline** direction

$$\tau = \frac{\delta_{\xi} V_{\xi} H_{\xi} + \delta_{\eta} V_{\eta} H_{\eta} + \delta_{\zeta} V_{\zeta} H_{\zeta}}{2V_i V_j}$$

$$V_{Y_i} = \hat{e}_{Y_i} \cdot \mathbf{u} \text{ , where } (Y_1, Y_2, Y_3) = (\xi, \eta, \zeta)$$

# The features of matrix equation

For the smallest problem ( $2^3$  elements), the indefinite matrix equations takes the following profile



- ▶ The linear system derived from the FEM is usually written in the form

$$\underline{\underline{\mathbf{A}}} \underline{\mathbf{X}} = \underline{\mathbf{b}} \quad \text{----- (*)}$$

where  $\underline{\underline{\mathbf{A}}}$  : Global matrix

$\underline{\mathbf{X}}$  : Solution vector

$\underline{\mathbf{b}}$  : Right hand side

- ▶ 90% of computation time of FEM is used to solve (\*) by **GMRES** or **BICGSTAB**
- ▶ It is a grand challenge to solve (\*) efficiently using iterative solution solver for three-dimensional large size problems  
=> our proposed iterative solver given below

- ▶ To avoid Lanczos or pivoting breakdown, the linear system has been normalized, leading to the following positive definite and symmetric matrix equation

$$\underline{\tilde{\mathbf{A}}} \underline{\mathbf{X}} = \underline{\tilde{\mathbf{b}}}, \text{ where } \underline{\tilde{\mathbf{A}}} = \underline{\mathbf{A}}^T \underline{\mathbf{A}}, \underline{\tilde{\mathbf{b}}} = \underline{\mathbf{A}}^T \underline{\mathbf{b}}$$

symmetric and positive definite

- ▶ Since the normalized linear system becomes **symmetric** and **positive definite**, the unconditionally convergent **Conjugate Gradient (CG)** iterative solver is applied
- ▶ Preconditioning of matrix is adopted to reduce the corresponding increase of condition number
- ▶ Two techniques will be adopted to accelerate matrix calculation
  - The mesh coloring technique (To avoid the race-condition)
  - The EBE technique (To avoid global matrix assembling)

# PCG algorithm

**Algorithm** The PCG iterative solver for  $\underline{\underline{A}}^T \underline{\underline{A}} \underline{x} = \underline{\underline{A}}^T \underline{b}$

Starting from an initial guess  $\underline{x}_0$

$\underline{\underline{B}}$ : Jacobi preconditioner

Compute  $\underline{x}'_0 = \sum_e \underline{\underline{A}}^e \underline{x}_0$ ,  $\underline{x}''_0 = \sum_e (\underline{\underline{A}}^e)^T \underline{x}'_0$ ,  $\underline{b}' = \sum_e (\underline{\underline{A}}^e)^T \underline{b}$   $\leftarrow$  **EBE step**

Compute the initial residual  $\underline{r}_0 = \underline{b}' - \underline{x}''_0$

$$\underline{z}_0 = \underline{\underline{B}}^{-1} \underline{r}_0$$

$$\underline{p}_0 = \underline{z}_0$$

For  $j = 1, 2, \dots$ ,

$$\underline{p}'_{j-1} = \sum_e \underline{\underline{A}}^e \underline{p}_{j-1} \leftarrow \text{EBE step}$$

$$\underline{p}''_{j-1} = \sum_e (\underline{\underline{A}}^e)^T \underline{p}'_{j-1} \leftarrow \text{EBE step}$$

$$\alpha_{j-1} = (\underline{p}_{j-1}, \underline{r}_{j-1}) / (\underline{p}_{j-1}, \underline{p}''_{j-1})$$

$$\underline{x}_j = \underline{x}_{j-1} + \alpha_{j-1} \underline{p}_{j-1} \leftarrow \text{Update step}$$

$$\underline{r}_j = \underline{r}_{j-1} - \alpha_{j-1} \underline{p}''_{j-1} \leftarrow \text{Update step}$$

**Convergence check**

$$\underline{z}_j = \underline{\underline{B}}^{-1} \underline{r}_j$$

$$\beta_{j-1} = (\underline{z}_j, \underline{r}_j) / (\underline{r}_{j-1}, \underline{r}_{j-1})$$

$$\underline{p}_j = \underline{z}_j + \beta_{j-1} \underline{p}_{j-1} \leftarrow \text{Update step}$$

**End**

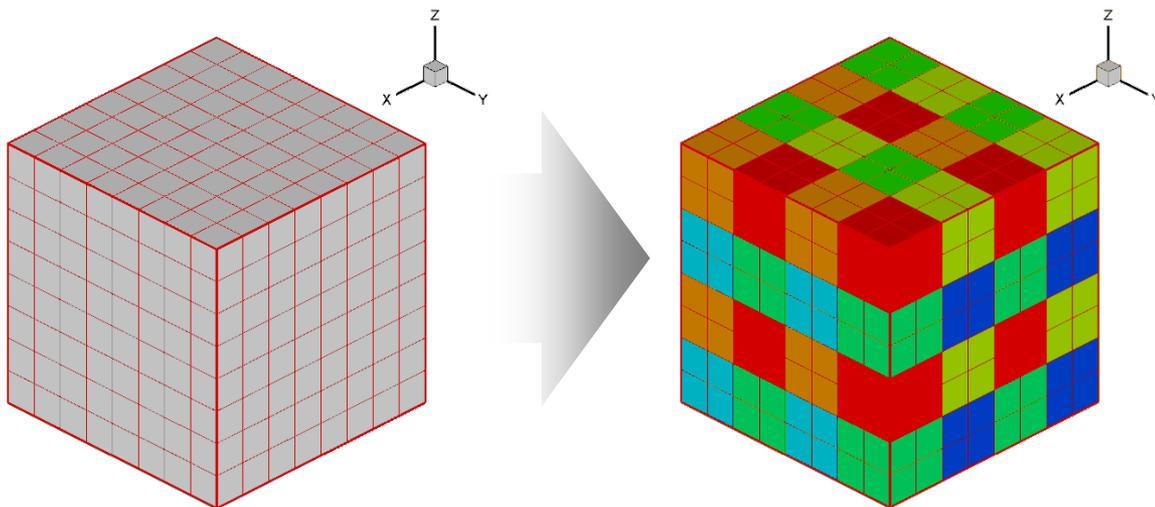
Mesh coloring  
technique

# Mesh coloring algorithm

- ▶ In the “**Update step**” of PCG algorithm, provided that any two threads update the same vector component **simultaneously**, it will result in the so-called **race-condition**

(If two men ride on a horse, one must ride behind) “一山不容二虎”

- ▶ To avoid the race-condition, we divide the all elements into a finite number of subsets such that any two elements in a subset are not allowed to share the same node



- \* Different colors will be proceeded in sequence
- \* Elements with the same color are proceeded simultaneously and in parallel

# EBE algorithm

In FEM, the global matrix  $\underline{\underline{\mathbf{A}}}$  can be formulated as :

$$\underline{\underline{\mathbf{A}}} = \sum_{e=1}^{Nel} (\underline{\underline{\mathbf{B}}}^e)^T \underline{\underline{\mathbf{A}}}^e \underline{\underline{\mathbf{B}}}^e, \quad \underline{\underline{\mathbf{B}}}^e : \text{ Boolean matrix}$$

Underlying this concept, we have

$$\underline{\underline{\mathbf{A}}} \underline{\underline{\mathbf{x}}} = \left( \sum_{e=1}^{Nel} (\underline{\underline{\mathbf{B}}}^e)^T \underline{\underline{\mathbf{A}}}^e \underline{\underline{\mathbf{B}}}^e \right) \underline{\underline{\mathbf{x}}} = \sum_{e=1}^{Nel} (\underline{\underline{\mathbf{B}}}^e)^T (\underline{\underline{\mathbf{A}}}^e \underline{\underline{\mathbf{x}}}^e)$$

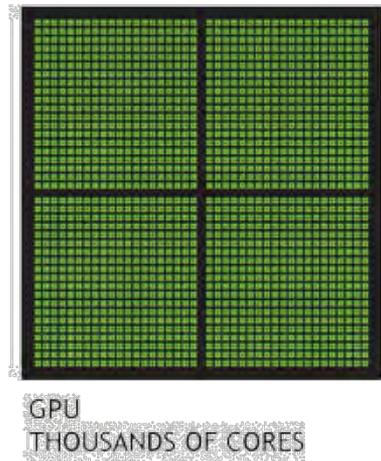
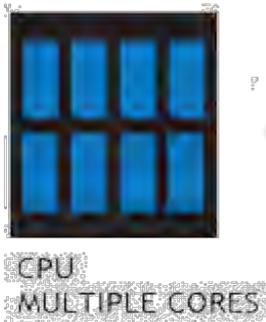
=> Time-consuming the **matrix-vector product** can be therefore implemented in an element-wise level

# III. Code implementation in Kepler K20 GPUs



# Introduction on CPU-GPU hybrid computing platform

- ▶ Computation is carried out on GPUs (**G**raphics **P**rocessing **U**nit) to get a good performance in scientific and engineering computing



## \* Parallelism

- CPU has 8 cores
- GPU has more than 1000 cores

## \* Memory

- CPU can reach 59 GB/s bandwidth
- GPU can reach 208 GB/s bandwidth

## \* Performance

- CPU can reach 59.2 GB (**double**)
- GPU can reach 3.52 TB/s (**single**)  
1.17 TB/s (**double**)



Intel i7 4820K  
Processor ~ 3.7GHz



NVIDIA Tesla K series  
Kepler K20

Processor clock rate = 705MHz  
Bandwidth = 208 GB/s  
Memory = 5 GB (DDR5)



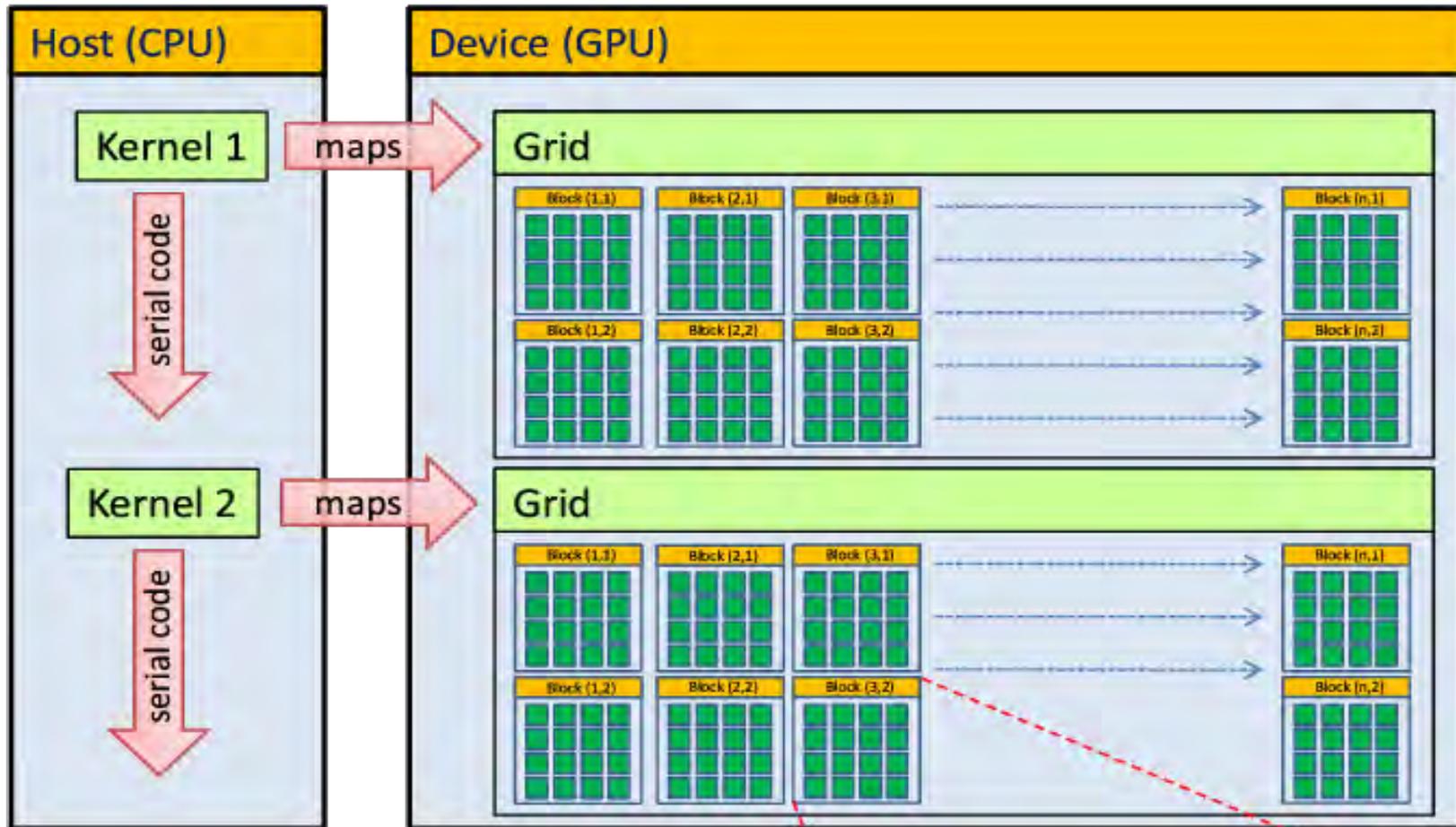
# CUDA programming

- ▶ (Compute Unified Device Architecture)
- ▶ CUDA is a general purpose parallel computational model released by NVIDIA in 2007
- ▶ CUDA only runs on NVIDIA GPUs
- ▶ Heterogeneous serial (CPU)-parallel (GPU) computing
- ▶ Scalable programming model
- ▶ CUDA permits Fortran (PGI) and C/C++ (NVIDIA NVCC) programming compilers

# CUDA programming

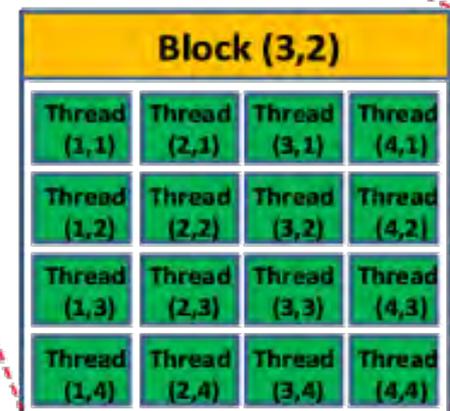
- ▶ CPU (**host**) plays the role of “**master**” and GPUs (**device**) play the role of “**workers**”
  - ◆ CPU maps computational tasks onto GPUs
  - ◆ CPU can carry out computations at the same time when code is run on GPUs
- ▶ Basic flow chart
  - CPU transfers the data to the GPU
  - Parallel code (**kernel**) are run on the GPU
  - GPU transfers the computational results back to CPU

Current CPU-GPU hardware architecture

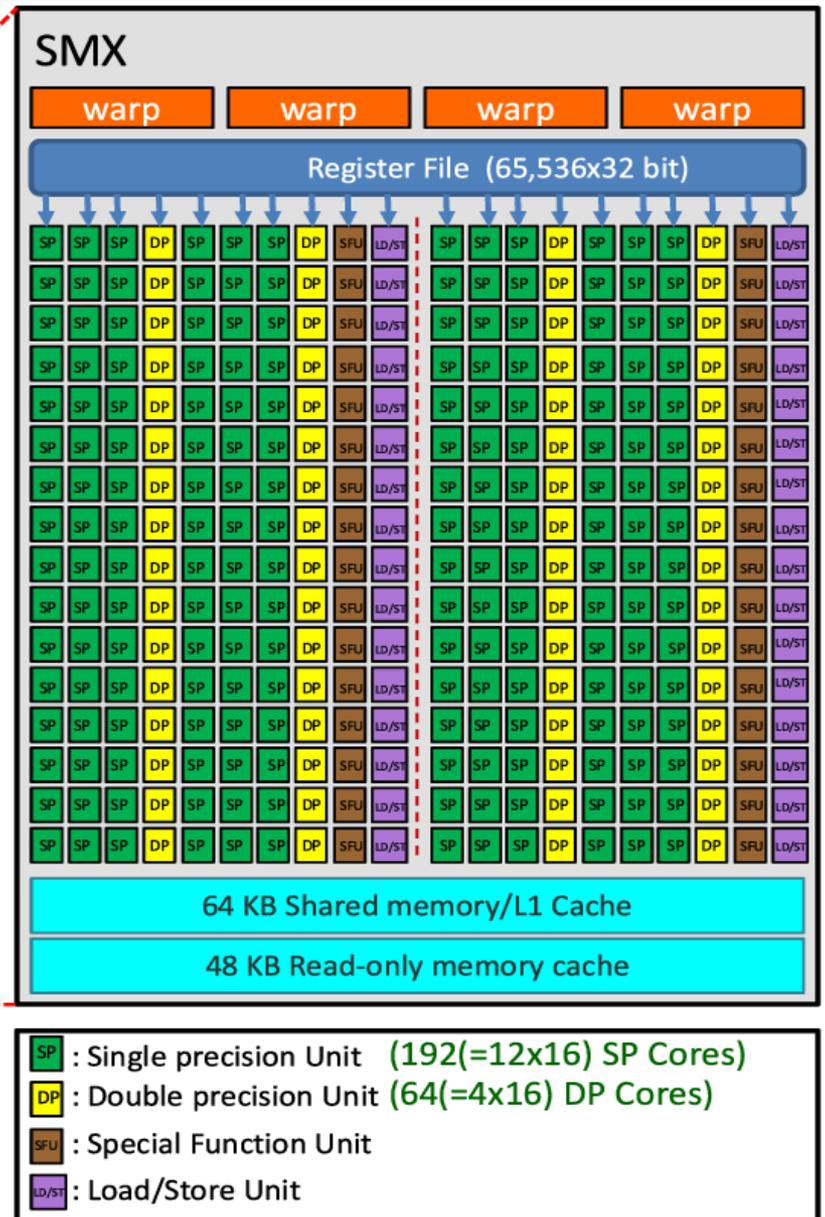


Note

- 1. Block is the smallest execution unit in GPU
- 2. All threads in a block are executed simultaneously



# GPU Kepler K20 architecture



# Classification of Memory in K20

1. Global memory (5GB DDR5) (Double Data Rate )
2. Shared memory (48KB/SMX)
3. Cache memory (1536KB L<sub>2</sub>)
4. Texture memory (48KB)

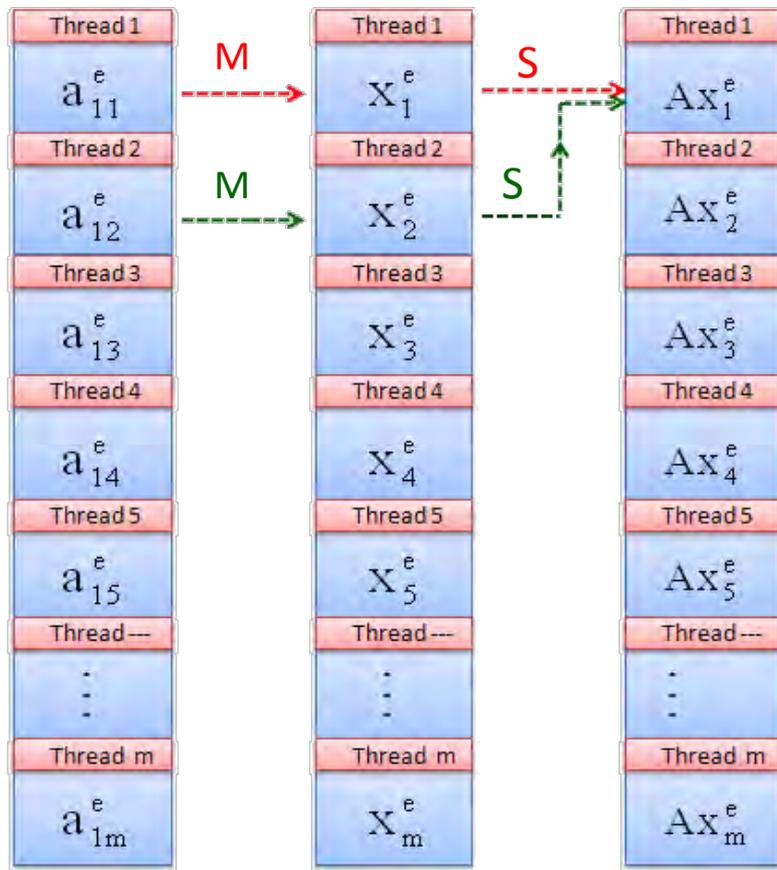
# EBE implementation in GPU (1)

- ▶ To implement the EBE on GPU, the new formulation will be adopted. Given an e-th element matrix, RHS and its product are as follows

$$\underline{\underline{A}}^e = \begin{pmatrix} a_{11}^e & a_{12}^e & \dots & \dots & a_{1m}^e \\ a_{21}^e & a_{22}^e & \dots & \dots & \dots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots \\ a_{m1}^e & & & & a_{mm}^e \end{pmatrix}, \underline{x}^e = \begin{pmatrix} x_1^e \\ x_2^e \\ \vdots \\ \vdots \\ x_m^e \end{pmatrix}, \underline{Ax}^e = \begin{pmatrix} Ax_1^e \\ Ax_2^e \\ \vdots \\ \vdots \\ Ax_m^e \end{pmatrix}$$

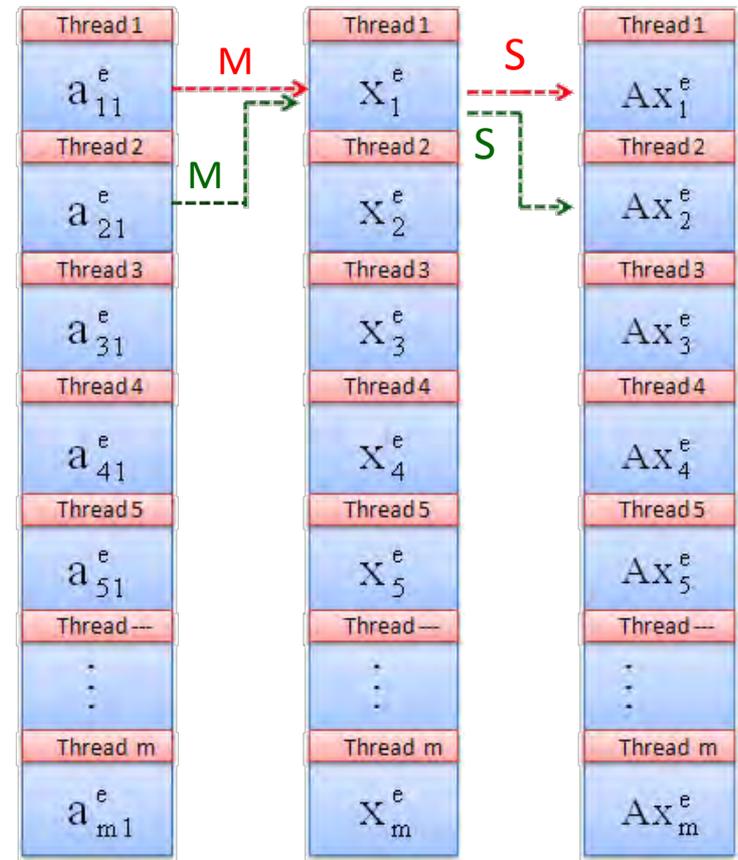
- ▶ One needs to compute the product  $\underline{Ax}_1^e, \dots, \underline{Ax}_m^e$ , The resulting matrix-vector product results in a race-condition

# EBE technique in GPU (2)



Race-condition algorithm

(Traditional)



Non-race-condition algorithm

(New operation)

Elements with the same color will be executed in parallel in a non-race-condition algorithm simultaneously

M : Multiplication operation  
S : Sum operation

# IV. Verification and numerical results



# Verification study

- ▶ The analytic solutions which satisfy the steady-state 3D Navier-Stokes equations are given as follows :

$$\mathbf{u} = \frac{1}{2}(y^2 + z^2), \mathbf{v} = -z, \mathbf{w} = y$$

The corresponding exact pressure is given below

$$p = \frac{1}{2}(y^2 + z^2) + \frac{2}{\text{Re}}x$$

- ▶ Problem setting
  - Re=1,000
  - Non uniform grid size : 11<sup>3</sup>,21<sup>3</sup>,41<sup>3</sup>,61<sup>3</sup>
  - Matrix Solver : PCG solver (iterative)

# Verification study

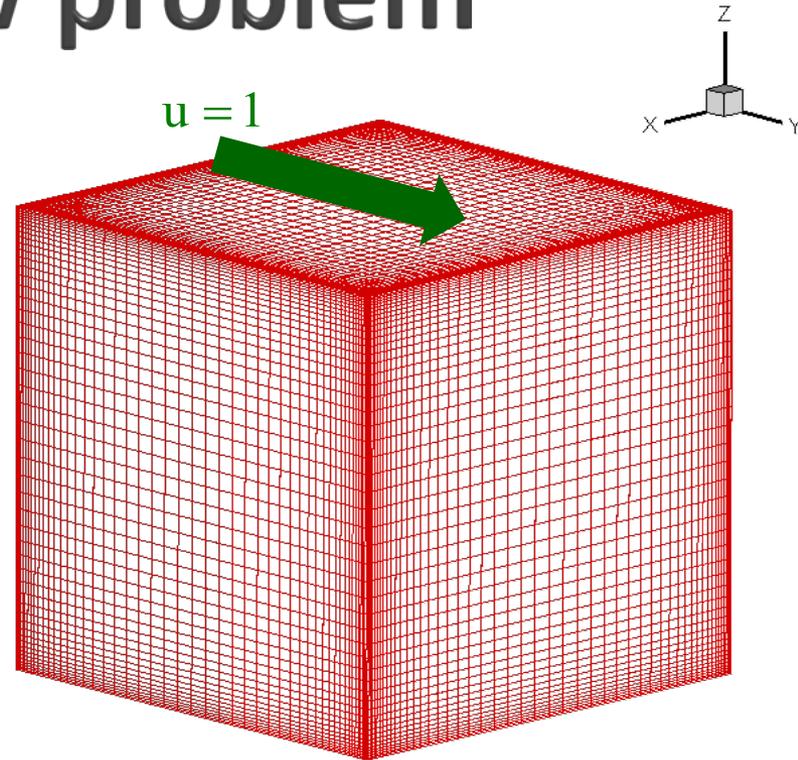
- ▶ The  $L_2$  error norms are listed as follows :

	L2U	L2V	L2W	L2P
$11^3$	$1.366 \times 10^{-3}$	$2.835 \times 10^{-3}$	$3.847 \times 10^{-3}$	$6.186 \times 10^{-3}$
$21^3$	$1.420 \times 10^{-4}$	$3.207 \times 10^{-4}$	$5.606 \times 10^{-4}$	$1.713 \times 10^{-3}$
$41^3$	$3.485 \times 10^{-5}$	$7.357 \times 10^{-5}$	$1.120 \times 10^{-4}$	$4.378 \times 10^{-4}$
$61^3$	$1.861 \times 10^{-5}$	$3.462 \times 10^{-5}$	$5.040 \times 10^{-5}$	$1.683 \times 10^{-4}$
R.O.C.	2.26	2.32	2.30	1.90

- ▶ R.O.C. (Rate Of Convergence )
- ▶ The predicted solutions show good agreement with the exact solutions

# Lid-driven cavity flow problem

- ▶ Problem domain :  $\Omega = [0,1] \times [0,1] \times [0,1]$
- ▶ Boundary condition
  - $u=1, v=w=0$  at upper plane
  - Non-slip at other planes
- ▶ Non-uniform grid number :  $51^3$
- ▶  $Re=100, 400, 1000, 2000$
- ▶ Jacobi preconditioner is adopted
- ▶ Newton-linearization is adopted



Schematic of the lid-driven cavity flow problem

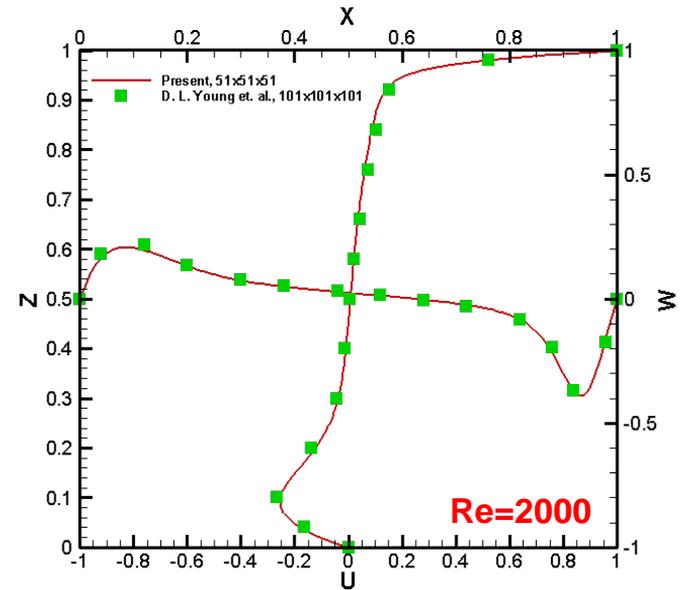
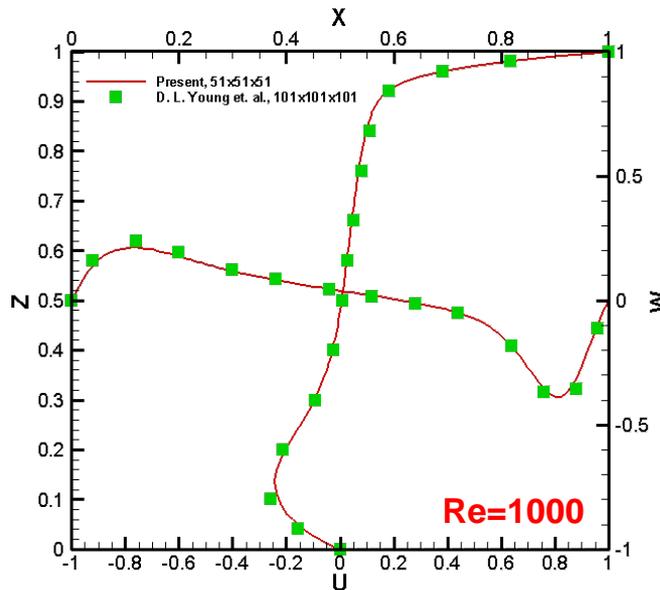
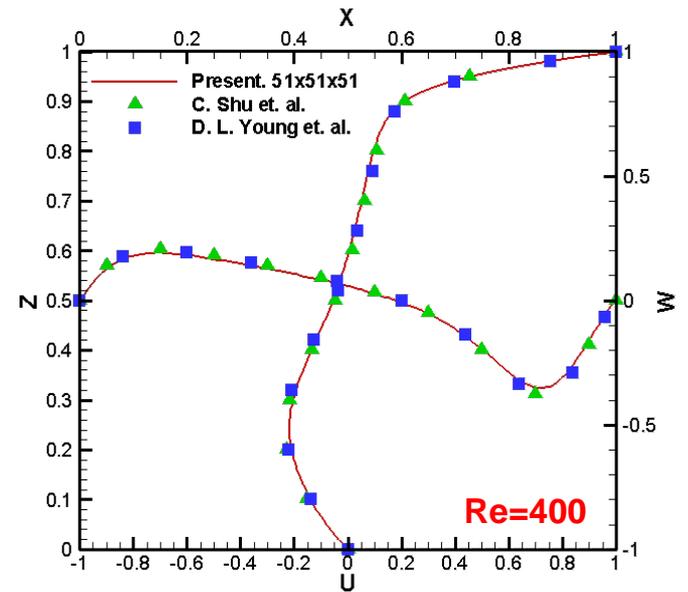
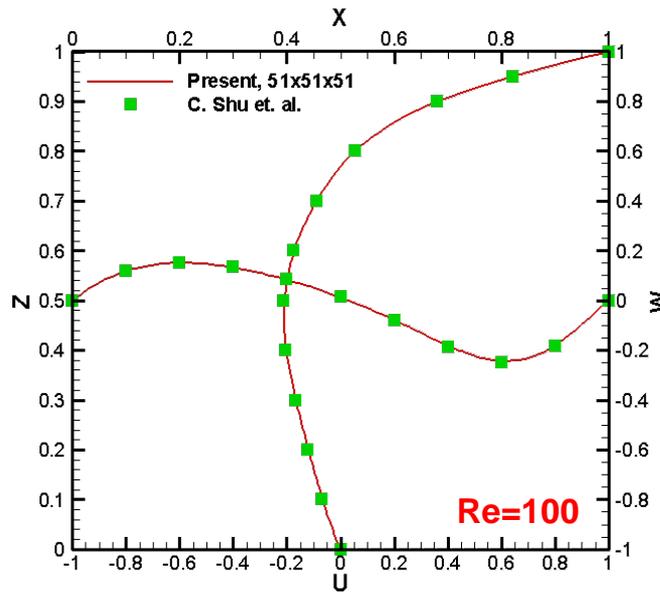
- D. C. Lo, K. Murugesan, D. L. Young, Numerical solution of three-dimensional velocity Navier-Stokes equations by finite difference method, *International Journal for numerical methods in Fluids*, **47**, 1469-1487, 2005
- C. Shu, L. Wang and Y. T. Chew Numerical computation of three-dimensional incompressible Navier-Stokes equations in primitive variable form by DQ method, *International Journal for Numerical methods in Fluids*, **43** :345-368, 2005.

# Velocity comparison

Compare the velocity profiles

$u(0.5, y, 0.5)$

$w(x, 0.5, 0.5)$



# Performance study

- ▶ Computations are implemented on the following two different platforms

Platform #1 specification	Platform #2 specification
Single CPU platform CPU : Intel i7-4820K GPU : Non	Hybrid CPU/GPU platform CPU : Intel i7-4820K GPU : NVIDIA Kepler K20

- ▶ The lid-driven cavity problem at different Reynolds and grid numbers is used to assess the performance of the current hybrid CPU/GPU calculation

# Performance study

► Grid number =  $31^3$

Platform	Re=100	Re=400	Re=600	Re=800	Re=1000
Platform # 1 (a)	416.5	860.4	1292.8	1769.5	2437.2
Platform # 2 (b)	4555.4	10087.9	15075.5	21515.7	29836.2
Speedup (b)/(a)	10.93	11.72	11.66	12.15	12.24

► Grid number =  $41^3$

Platform	Re=100	Re=400	Re=600	Re=800	Re=1000
Platform # 1 (a)	1377.1	2599.9	3424.1	4687.4	6502.3
Platform # 2 (b)	19823.2	36639.0	54635.3	65215.7	91780.6
Speedup (b)/(a)	14.3	14.0	15.9	13.9	14.11

# Performance study

► Grid number =  $51^3$

Platform	Re=100	Re=400	Re=600	Re=800	Re=1000
Platform # 1 (a)	3433.0	5983.7	7113.0	9606.1	12487.2
Platform # 2 (b)	50782.0	91012.0	113921.9	144358.3	196173.9
Speedup (b)/(a)	14.79	15.20	16.01	15.02	15.71

► Grid number =  $61^3$

Platform	Re=100	Re=400	Re=600	Re=800	Re=1000
Platform # 1 (a)	7702.8	11598.7	13939.7	17684.7	22222.5
Platform # 2 (b)	111338.8	168420.8	202867.8	256204.9	331559.7
Speedup (b)/(a)	14.45	14.52	14.55	14.48	14.92

# IV. Concluding remarks



# Concluding remarks

- ▶ A new streamline upwind FEM model accommodating the optimized numerical **wavenumber** along the **streamline** has been developed
- ▶ The **unconditionally convergent** finite element solution can be iteratively obtained from the **normalized** symmetric and positive definite matrix equation using the **PCG solver**
- ▶ Novel non-race-condition and **EBE** techniques have been successfully implemented on GPU architecture
- ▶ Computational time has been reduced more than ten times in a **hybrid CPU/GPU platform**

**Thank for your attention !**

**Q/A ?**