# Matrix transpose
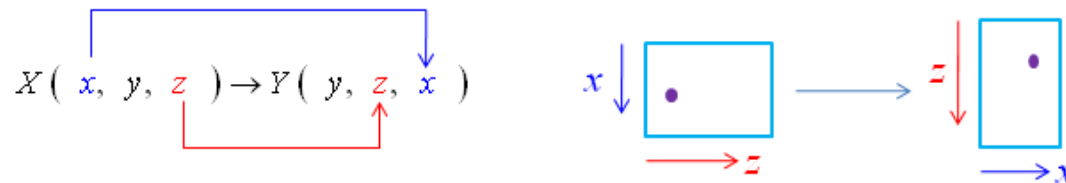
Author : Lung-Sheng Chien

National Tsing Hua university, R.O.C (Taiwan)

Mail:  d947207@oz.nthu.edutw

Abstract: we provide 2D transpose and 3D transpose in this document, source code of 2D transpose comes from SDK and then we use the same idea to build 3D transpose of the form
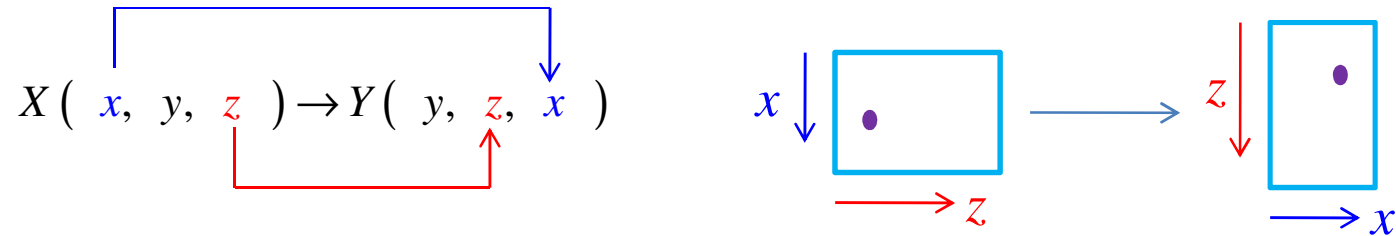
$$X(\ x,\ y,\ z\ ) \rightarrow Y(\ y,\ z,\ x\ )$$

Experiment shows that our 3D result has "good" performance comparable with 2D result (2D transpose is optimized version)

# Transpose on 3D data

Objective: given 3D data $X\left(1:n_1,1:n_2,1:n_3\right)$ , we want to do transpose operation under $\left(y,z,x\right)\leftarrow\left(x,y,z\right)$

such that $X\left(1:n_1,1:n_2,1:n_3\right)\rightarrow Y\left(1:n_2,1:n_3,1:n_1\right)$ with utilization of coalesce property.

Observation: it is similar to 2D transpose, if we only consider x-z slice when fixed y

$$X\left(\;x,\;y,\;z\;\right)\rightarrow Y\left(\;y,\;z,\;x\;\right)$$

The simplest way is (1) use 2D grid to represent (x, z) slice and do transpose operation along y

$for\;\; y=1:n_2$

 for each threads in $x$ - $z$ slice, do transpose

 $X\left(idx,y,idz\right)\rightarrow Y\left(y,idz,idx\right)$   ⟶   we must use share memory to decrease latency
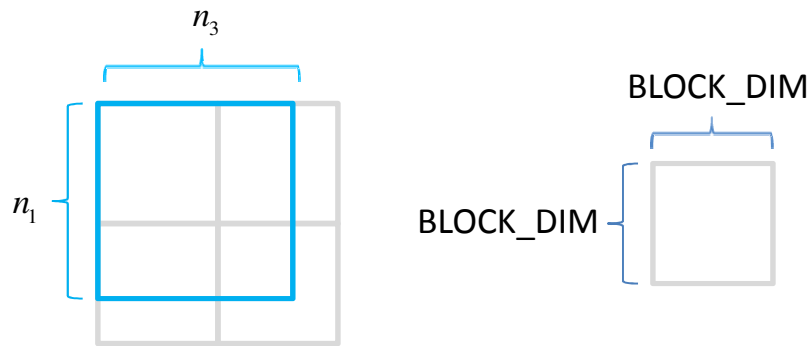
 $endfor$

Remark: motivation of transpose on 3D data comes from 3D FFT (see Sine_transform_3D.ppt)

# Transpose 3D: framework [1]

Objective: define a framework to do xyz2yzx transpose operation and use different kind of techniques

The number of grids in z-direction to cover x-z slice is (n3 + BLOCK_DIM-1)/ BLOCK_DIM

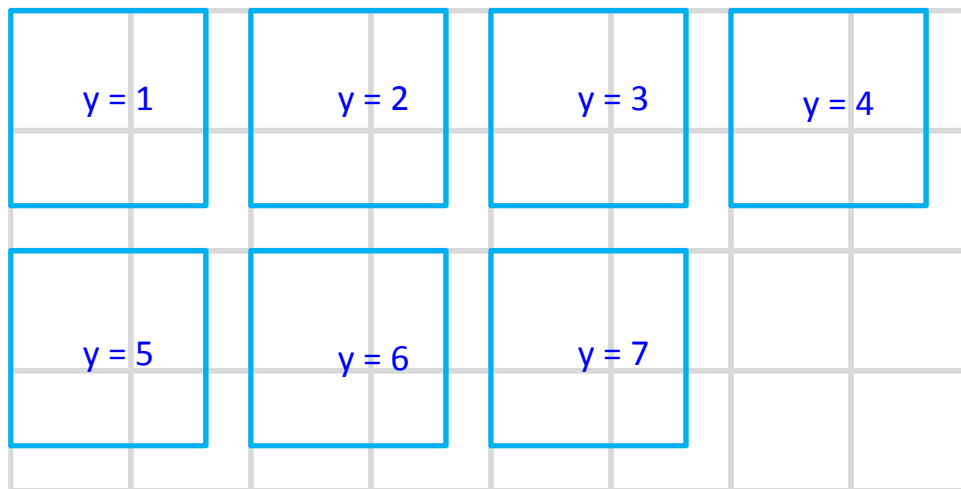The number of grids in x-direction to cover x-z slice is (n1 + BLOCK_DIM-1)/ BLOCK_DIM

$n_3$

$n_1$

BLOCK_DIM

BLOCK_DIM

$$Gx \equiv \frac{n_1 + BLOCK\_DIM - 1}{BLOCK\_DIM} = 2$$

$$Gz \equiv \frac{n_3 + BLOCK\_DIM - 1}{BLOCK\_DIM} = 2$$

Assume n2 = 7 ( 7 x-z slice), then what is configuration of grid?

| y = 1 | y = 2 | y = 3 | y = 4 |

| y = 5 | y = 6 | y = 7 |

$$k_1 \equiv floor\left(\sqrt{n_2}\right), \quad k_2 \equiv ceil\left(\frac{n_2}{k_1}\right)$$

$$then \quad grid = \left(k_2 \cdot Gz, k_1 \cdot Gx\right)$$

$$n_2 = 7$$

$$k_1 = floor\left(\sqrt{7}\right) = 2, \quad k_2 = ceil\left(\frac{7}{2}\right) = 4$$

**Trick:** we can do better grid configuration such that resource utilization is highest

$$for \ \ k_1 = floor\left(\sqrt{n_2}\right): -1:1$$

$$k_2 = ceil\left(\frac{n_2}{k_1}\right)$$

Then we waste one block at most

$$if \ \ k_1 k_2 - n_2 \leq 1, \ \ then \ \ grid = \left(k_2 \cdot Gz, k_1 \cdot Gx\right), \ break$$

$$end$$

---

Express  $\begin{array}{l} blockIdx.x = Gz \cdot s_1 + t_1 \\ blockIdx.y = Gx \cdot s_2 + t_2 \end{array}$  such that  $\begin{array}{l} \left(s_1, s_2\right): \text{index to y-direction} \\ \left(t_1, t_2\right): \text{index to } x\text{-}z \text{ slice} \end{array}$  where  $\begin{array}{l} s_1 = floorf\left(blockIdx.x \ / \ Gz\right) \\ t_1 = blockIdx.x - Gz \cdot s_1 \end{array}$

$(0,0)$　　　　$(1,0)$　　　　$(2,0)$　　　　$(3,0)$

$\left(s_1, s_2\right) = \left(0,0\right)$

| y = 0 | y = 1 | y = 2 | y = 3 |
|---|---|---|---|

y = 0

coarse grid

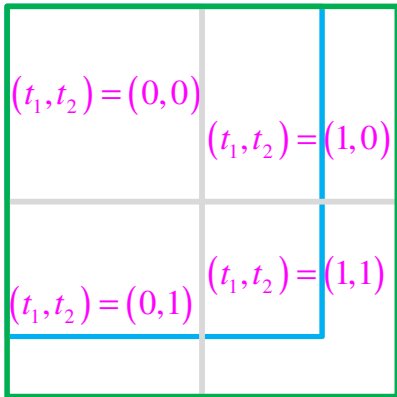| y = 4 | y = 5 | y = 6 | |
|---|---|---|---|

$(0,1)$　　　　$(1,1)$　　　　$(2,1)$　　　　$(3,1)$

$$yIndex = s_2 k_2 + s_1$$

where  $grid = \left(k_2 \cdot Gz, k_1 \cdot Gx\right)$
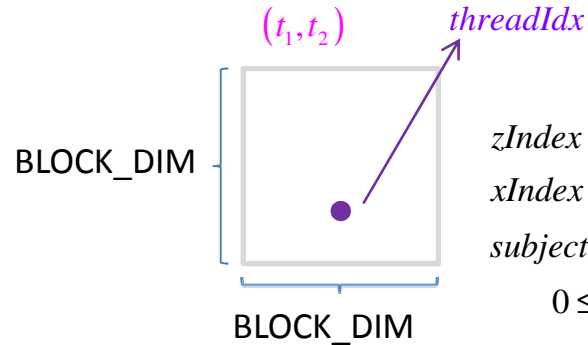
$$k_2 = \# \text{ of coarse grid along } z$$

$$= gridDim.x \ / \ Gz$$

# Transpose 3D: framework   [3]

$(s_1, s_2) \rightarrow yIndex$

$(t_1, t_2)$ may be regarded as local block ID

$(t_1, t_2) = (0,0)$
$(t_1, t_2) = (1,0)$
$(t_1, t_2) = (0,1)$
$(t_1, t_2) = (1,1)$

$(t_1, t_2)$          $threadIdx$

BLOCK_DIM

BLOCK_DIM

$zIndex = t_1 \cdot BLOCK\_DIM + threadIdx.x$

$xIndex = t_2 \cdot BLOCK\_DIM + threadIdx.y$

$subject \quad (zIndex, xIndex) \equiv (k-1, i-1) \quad since \; X(1:n_1, 1:n_2, 1:n_3)$

$0 \le xIndex < n_1 \quad , \quad 0 \le zIndex < n_3$
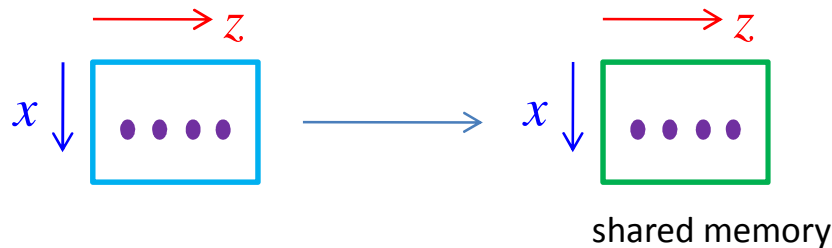
**1** Fixed y, real (x, y, z) to shared memory

$row - major(i, j, k) = (i-1)n_2 n_3 + (j-1)n_3 + (k-1)$

$for \; i = 1:n_1$
$\quad for \; k = 1:n_3$
$\qquad shared[i][k] \leftarrow X(i, j, k)$
$\quad endfor$
$endfor$

$X(1:n_1, 1:n_2, 1:n_3)$

$xIndex = i - 1$
$yIndex = j - 1$
$zIndex = k - 1$

$if \; (xIndex < n_1) \; and \; (zIndex < n_3)$
$\quad index\_in = xindex \cdot n_2 n_3 + yIndex \cdot n_3 + zIndex$
$\qquad shared[threadIdx.y][threadIdx.x] := X[index\_in]$
$endif$
$\_\_syncthreads(\;)$

$\xrightarrow{} z$

$x \downarrow$   • • • •

$\xrightarrow{} z$

$x \downarrow$   • • • •

shared memory

# Transpose 3D: framework    [4]

2  Transpose data in shared memory to $Y\left(1:n_2,1:n_3,1:n_1\right)$

$$for \ \ k=1:n_3$$
$$\quad for \ \ i=1:n_1$$
$$\quad\quad Y\left(j,k,i\right)\leftarrow shared\left[i\right]\left[k\right]$$
$$\quad endfor$$
$$endfor$$

shared memory

$(t_1,t_2)$

$z$

$x$

$X$    BLOCK_DIM

BLOCK_DIM

Transpose

$(t_2,t_1)$

$x$

$z$

$Y$    BLOCK_DIM

BLOCK_DIM

$zIndex = t_1 \cdot BLOCK\_DIM + threadIdx.x$
$xIndex = t_2 \cdot BLOCK\_DIM + threadIdx.y$

$xIndex = t_2 \cdot BLOCK\_DIM + threadIdx.x$
$zIndex = t_1 \cdot BLOCK\_DIM + threadIdx.y$

$$Y\left(1:n_2,1:n_3,1:n_1\right)$$

$$row-major\left(j,k,i\right)=\left(j-1\right)n_3n_1+\left(k-1\right)n_1+\left(i-1\right)$$

$$if \ \left(xIndex<n_1\right) \ \ and \ \left(zIndex<n_3\right)$$
$$\quad index\_out = yindex\cdot n_3n_1 + zIndex\cdot n_1 + xIndex$$
$$\quad Y\left[index\_out\right]:= shared\left[threadIdx.x\right]\left[threadIdx.y\right]$$
$$endif$$