

10/08/2024

Introduction to CUDA Python with Numba

⇒ Objective - ~~perf~~ GPU accelerate python code that performs element-wise operations on Numpy arrays.

⇒ Numba: Numba is just-in-time, type-specializing, function compiler for accelerating numerically-focused Python either for CPU or GPU.

function-compiler: - Compiles Python function to a faster function
type-specializing: Speeds up function by generating a specialized implementation of data type specific data types you are using.

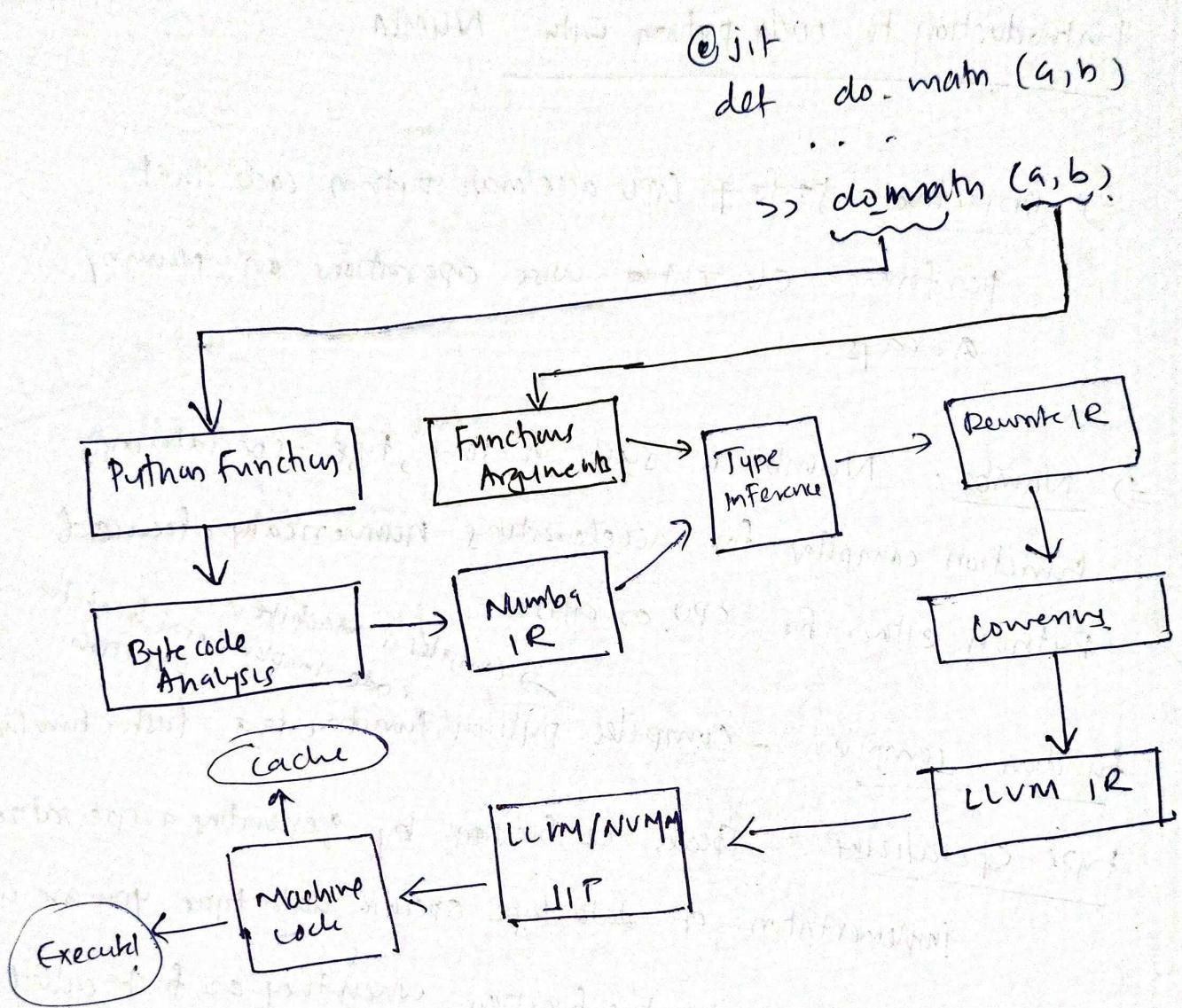
Just-in-time: Compiles the function when they are first called.

numerically-focused: Focused on numerical data types like int, float & complex.

Note

- Python built-in functions are faster than "Numba" because "numba" introduces some overhead to each function call that is larger than the overhead of function call overhead of python itself.

JIT compilation workflow



Byte code Analysis: Byte code is a lower-level, platform-independent representation of your python code.

Numba IR: Intermediate representation of python code by Numba. (Simplified low-level representation easier to analyze and optimize)

Type - Inference: Finding the type of the function arguments. Numba needs to know the data types to generate machine code.

Rewrite IR \Rightarrow Numba IR is rewritten with Type inference.

Lowering: \Rightarrow Numba IR is "lowered" to even a lower-level representation that is closer to machine code, known as LLVM IR.

LLVM IR \Rightarrow low-level Virtual Machine is a powerful compiler framework that Numba uses to generate highly optimized machine code.

LLVM/NVMM JIT: Compiler translates LLVM IR into machine code that can run on the CPU (using LLVM) or GPU (using NVMM).

Machine code: machine code is specific to the architecture of the machine where code is being run. The machine code is highly optimized for performance.

\hookrightarrow If function is again called with same argument type, Numba can use cache machine code from cache, avoiding need to go through the entire compilation process again.

Note

- * After @jit compilation the annotated source code could be accessed via "inspect-types."

No python and object mode

- * ~~@jit~~ @jit by default works in object mode, which does not do type-specialization.
- * @jit(nopython=True) is called the nopython mode, that do type-specialization and tells whether type inference step fails or not.

Note:

- Numba cannot compile every python data types.
- for e.g.: - Numba cannot compile, dictionaries
- @ jit(nopython=True) ≡ @njit

Numpy Ufunc

⇒ Numpy ufunc (or universal function) is a Numpy function that operates element-wise on arrays supporting array broadcasting, type casting & several other standard features. Such as vectorization and reduction operators.

Eg:- Arithmetic Ufunc \Rightarrow np.add, np.subtract, np.multiply etc,

Trigonometric Ufunc \Rightarrow np.sin, np.cos, np.tan etc,

Exponential and Logarithmic Ufunc \Rightarrow np.exp, np.log etc //

Comparison Operufunc \Rightarrow np.greater, np.less etc //

Bitwise operation ufunc \Rightarrow np.bitwise_and

np.bitwise_or etc //

⇒ One can also create custom ufunc using ~~"numpy.fromfunction"~~.

"numpy.fromfunction" function or by using the @vectorize decorator.

⇒ numpy ufunc are ideal to be carried out in GPU.

- Since numpy ufunc carry out same operation on different elements of the array, they are data parallel.
- GPU hardware is inherently data parallel. i.e., maximum throughput is achieved when ~~the~~ the GPU is computing the same operation on many different elements at once.

⇒ Using NumPy also one can create ufunc that can be run on GPU.

- this is by using the decorator `@vectorize`.

(NumPy vector works in same manner as NumPy ~~decor~~ vectors.)

- For running it on GPU (CUDA enabled device) we need to give an explicit type signature and target attribut.

~~a type signature is of this form:~~

- a type signature and target is of this form:

~~not~~ `@vectorize ('return_value_type(argument1_value-type,
argument2_value-type,...)', target = 'cuda')`.

⇒ GPUs have advantage over CPUs on certain cases:

- Inputs have to very large.

- Advanced calculation.

- Usage of GPU requires execution overhead which makes execution of simple function costly compared to CPU.

in GPU

CUDA Device Functions

- "ufuncs" is fantastic when you want to perform elementwise operations and can be accelerated CPU accelerated easily using numba.
- To compile functions that are not element wise, vectorized functions, we can use "numba.cuda.jit"
 - ↳ "numba.cuda.jit" can also be used in helperfunctions compile helperfunctions that ~~can be used inside~~ is used inside a GPU-accelerated "ufunc".

Syntax:

```
from cuda import jit
@cuda.jit(dence = True)
def helper_func():
    :
    return()

@vectorize(['float32(...)', target = 'cuda')
def -ufunc=():
    :
    return()
```

Note

- * The argument "device = True" ensures that the decorated function (decorated by `cuda.jit`) could be only used inside a function that is running on a CPU. (without CPU accelerated Ufunc)
and not on CPU hosted code.
- * During helper function calls, there is very little overhead since, CUDA compiler aggressively inlines device functions.
- * Allowed Python on the GPU:

1. if | elif | else
2. while and for loops
3. Basic math operations
4. Selected functions from the math and cmath modules
5. Tuples.

Managing GPU Memory

- ⇒ Numba transfers the data to GPU when it undergoes operation ~~from~~ from host (CPU)
- ⇒ After the operation the data gets ~~is~~ automatically transferred back to the host / CPU
- ⇒ Data is transferred to and from GPU while performing some computation and it is handled by Numba.
- ⇒ But ~~round-trip~~ data transfer is a time-intensive operation
- ⇒ But we can prevent the automatic data ~~from~~ transfer from GPU to host and store it GPU for additional process and copy it back to the host only after end of process in GPU.
- ⇒ The way to do it is to create "CUDA Device Arrays".
 - and pass them to our GPU functions.
 - Device arrays will not be automatically transferred back to the host after processing, and can be used again as we wish on the device before ultimately, and only if necessary, sending them, or parts of them, back to host.

Syntax

`x-device = cuda.to_device(x)`
`empty cuda device array: out-device = cuda.device_array(
shape=(n,), dtype='float32')`

`out_host = out_device.copy_to_host()`
↳ to copy the data back to the Host //

Custom CUDA kernels in python with Numba

⇒ custom CUDA kernels require more work to be implemented on GPU.

↳ we cannot use `@vectorize` as in the case of `ufunc`

- `Ufuncs` are limited to ~~scalar~~ case of scalar function operation that is operated on element wise data.
- There are more class of problems which cannot be solved by applying the same function to each element on dataset.
- examples of such cases : stencil algorithms or reduction algorithms
- These algorithms are inherently parallelizable but cannot be done using `ufunc`

CUDA kernels

- * Function written for GPU (while programming in CUDA) is called "kernel"
- * Execution of Function \Leftrightarrow Launching a kernel
- * Kernels are launched on the GPU's many cores with "parallel" threads. \Rightarrow "Execution configuration" or "launch configuration" is used to describe parallel execution configuration.

CPU - Accelerated v/s CPU - only application

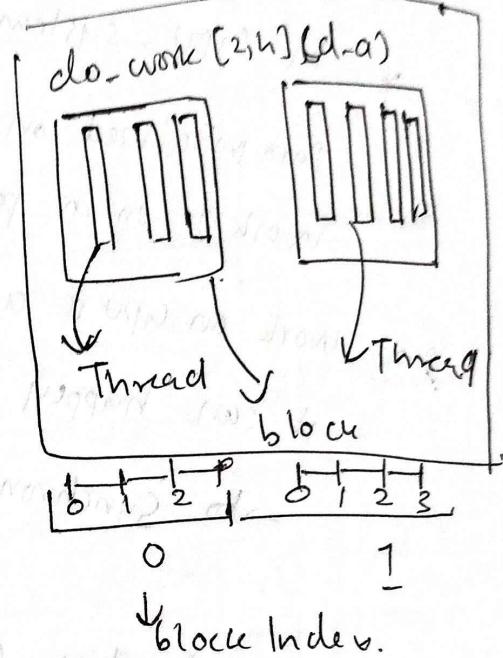
- * In CPU - only application data is allocated on the CPU and all the work is performed serially.
- * In GPU system there is both host and device memory. Data initialized on CPU can be copied to GPU where it can be worked on in parallel.
- * Work on GPU is asynchronous w.r.t host (CPU), thus it can happen in same time.
 - to synchronize GPU and CPU one could use: `cuda::synchronize()`
 - data then can be copied back to the CPU for verification.

Execution in GPU

- GPU do work in parallel and work done is done in "thread"
- and many threads run in parallel.
- A collection of threads is called a block
- A collection of blocks is called with a grid kernel
- Launch is called a grid.
- To launch a kernel we require a execution configuration which specifies defines the number of blocks in a grid and number of threads in each block.
- Every block in the grid contains the same number of threads.

-) Inside ~~__global__~~ or kernel definitions, CUDA - provided variables describe its executing thread, block and grid.

- ↳ $\text{gridDim} \cdot x$ → number of blocks in the grid. ↳ grid.
- ↳ $\text{blockIdx} \cdot x$ → index of the current block
- ↳ $\text{blockDim} \cdot x$ → number of threads in each block (all blocks will have same number of threads).
- ↳ $\text{threadIdx} \cdot x$ → describes the index of the thread within a block.



⇒ Accessing individual threads (for a parallel mapping of data)

" $\text{threadIdx} \cdot x + \text{blockIdx} \cdot x + \text{blockDim} \cdot x$ " will return the thread's unique index in the whole grid which can map to data elements.

↳ number of threads provide a unique way to access individual thread using `cudaGrid(1)`

It gives unique indexing for each thread.

⇒ '1' implies one-dimensional grid of threads.

- ⇒ CUDA allows to construct 1, 2, and 3D grids.
 - ⇒ While doing launching of kernel we pass execution configuration that contains information about the number of threads and blocks for the problem.
 - ⇒ To accelerate custom functions in GPU we use `"cuda.jit"` using "number"
 ↳ Note: While using `@cuda.jit` we pass done the function does not return an output value but rather it will be written to an array.
 - Note: Taking copies of arrays from device to host memory implies implicit synchronization and manual synchronization is not required in such cases.
 - ⇒ Deciding the very best size for the CUDA thread grid is a complex problem and depends upon both the algorithm and the specific GPU's compute capability.
- Rough Heuristics that we tend to follow are these:
- (1) Size of ~~thread~~ a block should be a multiple of 32 threads (since each streaming multiprocessor in GPUs are made up of 32 warps!)
- with typical block sizes b/w 128 to 512 threads per block.

- size of grid should ensure the full GPU is utilized where possible. Some things in range like 200-100 blocks is a good starting place.
- CUDA launch overhead increases with blocks size.
 - so when the input size is very large it is best not to launch a grid where number of threads equal to the number of input elements.
 - ↳ this would increase number of blocks
 - in such cases approaches like Grid Stride Looping is performed.

Grid stride looping

- ⇒ used when the number of threads is less than the data elements to be processed.
- ⇒ Grid stride looping is same as a ~~for loop~~ while setting step size while writing the loop in python:

Syntax:

For loop with step size:

for i in range (start, stop, stepsize)

Grid stride looping

for i in range (start, data.shape[0], stride),

\Rightarrow "start" is given by: `cuda::grid(1)`
stride is given by: `cuda::gridsize(1)`
 \hookrightarrow number of threads in
the grid

$$(\text{cuda::gridsize}(1) = \text{blockDim.x} + \text{gridDim.x})$$

\Rightarrow Without grid stride looping, all data elements will not
get processed by threads and some data elements will be
left unprocessed.
~~With grid stride looping, the device coalesces memory
reads/writes into as few transactions as possible for performance.~~

\Rightarrow Thus in grid stride looping:
Each thread processes multiple elements by iterating
over them in steps that are equal to the total
number of threads (grid size). This allows the entire
dataset to be processed even when the grid size is
smaller than the total number of elements.

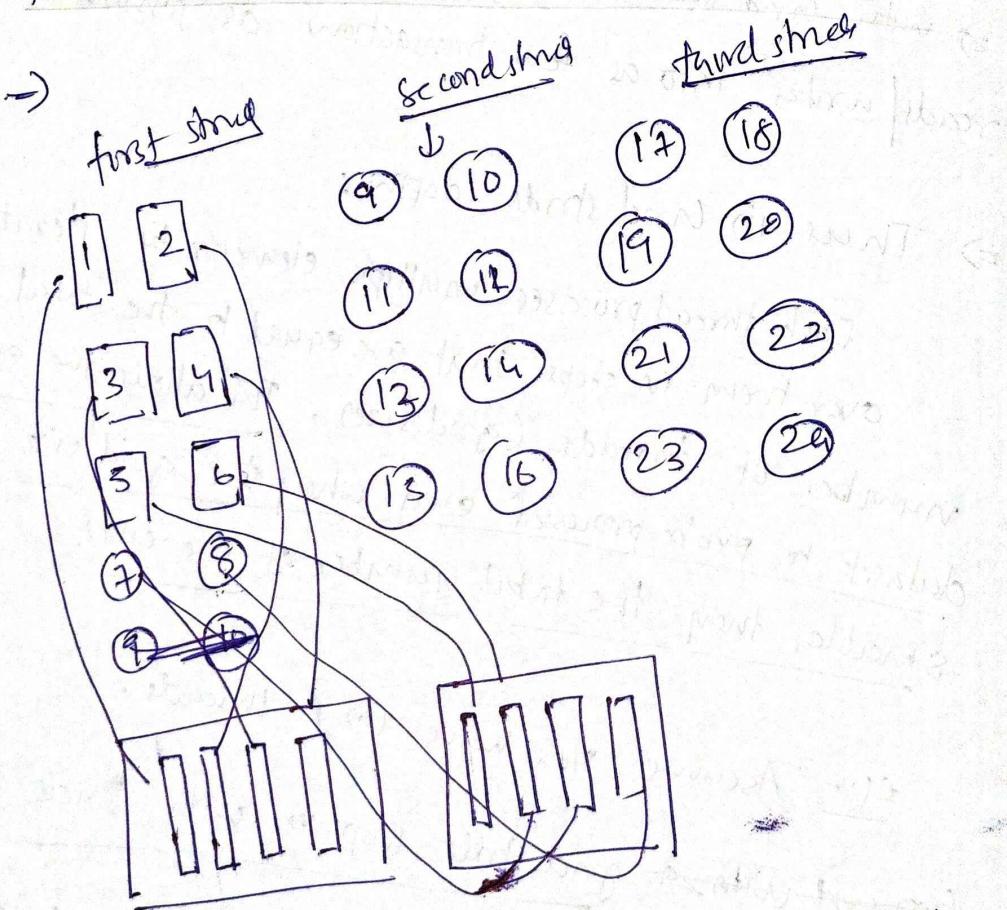
eg:- Assume you have 1024 threads.
~~Explain with a grid stride loop with stride = 1024.~~

1st thread access elements in 0, 1024, 2048
2nd thread access elements 1, 1025, 2049

goes like this

Memory coalescing

- ⇒ Memory coalescing occurs when threads in a warp (Subsystem of Streamline Multiprocessors containing 32 threads) access memory in a way that these access can be combined into a single or a few memory transactions.
- ⇒ Under stride loops facilitate memory coalescing and allows read/writes with few transactions.



- ⇒ End stride loops ensure that consecutive threads in a warp access consecutive memory locations even when they are looping over multiple elements.

NVIDIA - TASK 3.

Effective use of Memory subsystem

Global Memory Coalescing

Contiguous Memory location

Refers to a sequence of memory addresses that are next to each other without any gaps.

e.g.: - array = {10, 20, 30, 40}:

these integers in the array might be stored in contiguous memory locations like this:

Value 10	stored at address 0x1000
Value 20	stored at address 0x1004
Value 30	stored at address 0x1008
Value 40	stored at address 0x100C

⇒ Contiguous memory allows for efficient data access and transfer.

If a program reads data from contiguous memory, it can retrieve ~~over~~ a large block of data in one operation, rather than multiple smaller operations.

→ Memory Bus is the connection b/w CPU and RAM (Main memory) across the memory.
Now, data is transferred in blocks.

Modern computers are also equipped with cache lines (a type of small, fast memory close to the CPU) that stores data temporarily to speed up access.

↳ A cache block is a type of contiguous memory, of 64 bytes (or more) in size.

⇒ While reading contiguous memory, if the data is stored in contiguous locations, the CPU can load the data in one to cache in one operation.

↳ For instance, if the CPU needs the first element of an array, it might load the entire cache line (which could include several elements of the array) into cache. Thus subsequent access to other elements of the array can be satisfied directly from the cache, without needing additional trips to the slower main memory.

Memory coalescing in GPU

⇒ In GPU computing threads are organized in blocks and blocks are further divided in "warp" which consist of 32 threads.

⇒ Instructions are issued parallelly at warp level, i.e. each thread in a warp will carry out some instruction of different data elements.

⇒ Memory transfer b/w global device memory and CPU occur in segments, typically 32 byte segments.

↳ For data stored in L1 cache, it occurs in 128 byte segments → Cache memory closest to GPU/CPU

⇒ The objective of Memory coalescing is to reduce the number of memory lines (segments) needed to fulfill the read/write requirements of a warp.

↳ memory coalescing done properly leads to efficient memory transfer.

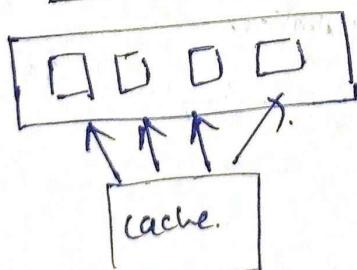
Contiguous Memory Access

⇒ If the data requested by the warps are stored at contiguous memory location, the memory access can be considered to be fully coalesced.

↳ If all the data in the memory ~~is~~ line will be used and data transfer occurs in few operations.

accessing contiguous location

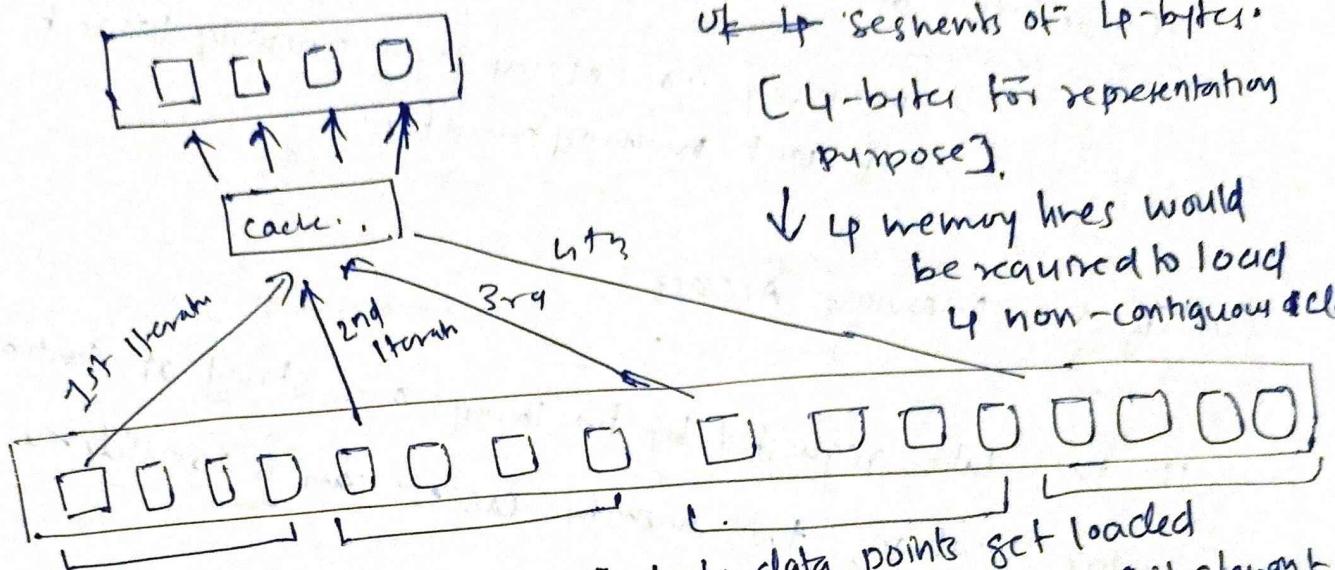
Warp



Entire memory lane will be loaded in one go and will be processed by Warp.

Data (global memory).

non-contiguous memory



Memory transfer occurs in ~~bytes~~
of 4 segments of 4-bits.

[4-bits for representation
purpose].

↓ 4 memory lines would
be required to load
4 non-contiguous data.

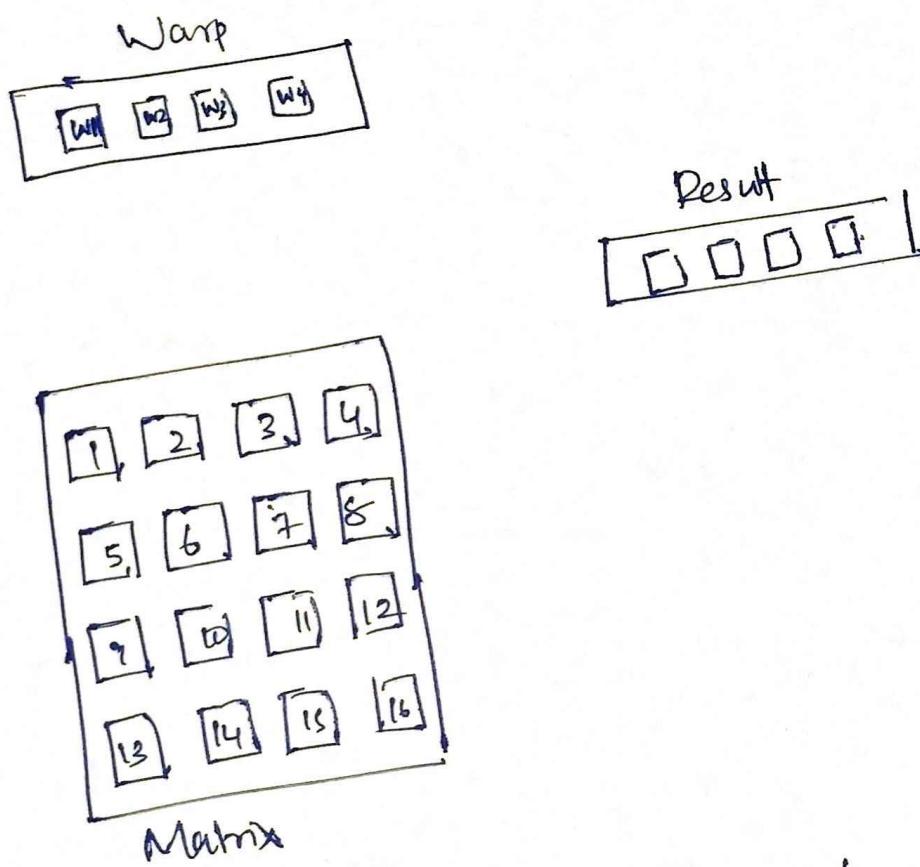
- ⇒ During 1st transfer ⇒ First 4 data points get loaded but GPU or warp only process 1st element and remaining 3 are discarded.
- During 2nd transfer ⇒ second 4 data points get loaded and warp only process the 1st element and remaining 3 are discarded.
- ⋮

It goes like this

- ⇒ Thus ~~accessing~~ non-contiguous memory location is inefficient and time consuming.
- ⇒ But we need to access non-contiguous data in cases like:
- (1) Sparse Matrix operations
 - (2) Graph Algorithms
 - (3) Datalog mining and database queries.
 - (4) Image & signal processing, etc.

Memory Coalescing while working with Matrix

Example (1): Consider a kernel that stores the sum of each row in a matrix (which is depicted by 4 contiguous data elements) in a result vector.



⇒ In this case there are two options available.

option -1: A single thread can iterate over a row summing it and store the result into a vector.

option -2: Parallelising the operation. ⇒ i.e.

~~like w_1 access~~

⇒ 1st row of matrix will be transferred to L1P0, and w_1 access the elements iteratively.

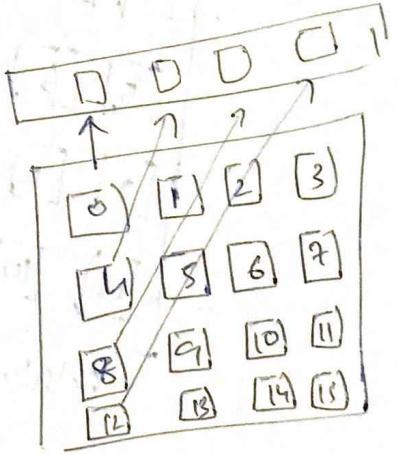
and now
they, ~~w₂~~ will be transferred to CPU ~~as~~ w₂ will be sums
the rows elements iteratively:
This process is repeated:

option (2): Accessing Row elements parallelly by a thread in warp

- (1) W[1] loads ~~the~~ first row to CPU access first element and
i.e.,
discard all other element.
- (2) W[1] loads second row to CPU access first element and
discard a

option (2): Parallelizing the process of summation: Each ~~one~~ thread
in warp ~~is~~ is requesting for data in different line of memory

- W[1] load first row access only first element and discard all other element
- W[2] load second row access only first element and discard all other element



⇒ In this process, in each iteration 4 data lines will be transferred to GPU but 75% data is discarded

- ⇒ Now this process is repeated to access all other elements of row in subsequent iterations
- ⇒ It also follows uncoalesced pattern. Thus it ~~will~~ load the column summation, there would be 16 memory line transfers in total and each use only 25% of data for each line transferred.

Example (2) : If each thread sums a column, the memory access pattern can be coalesced

Summing ~~for~~ each column and storing the resultant into a vector.

option (1) : Iterating over elements of column one by one

$w[1]$ loads ~~first~~ first row

$w[2]$ loads second row

$w[3]$ loads third row

$w[4]$ loads fourth row

access first element

access first element

access first element

access first element

This process is repeated ~~for all~~ by all threads. This process involves transfer of 16 memory lines and at each transfer ~~only~~ only 25% of data transferred ~~is~~ is used.

Option (2) : ~~Iterate~~ parallel execution:

First row is loaded to LPU and $w[1]$ access first element

~~now~~ $w[2] \rightarrow$ 2nd ~~row~~ element

$w[3] \rightarrow$ 3rd element

$w[4] \rightarrow$ 4th element

This process is repeated for other rows.

⇒ In this way memory access ~~will~~ be coalesced.

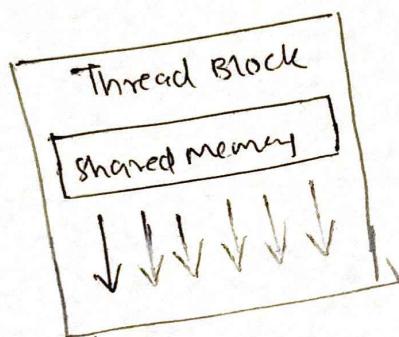
⇒ The computation requires only transfer of 4 memory lines and 100% of data ~~transferred~~ is used in each iteration.

Conclusion

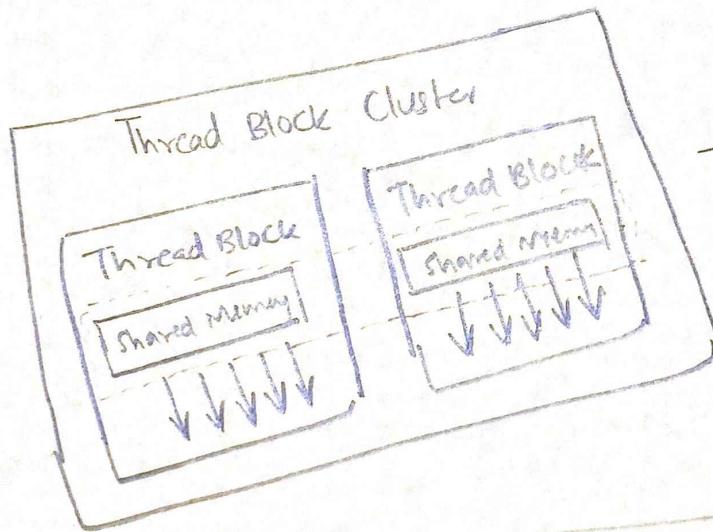
- ⇒ while performing row operations
- ⇒ If we assume the row elements forms a contiguous memory operations like summing elements in row is best to be performed iteratively rather than parallelly (in terms of memory coalescing) in GPU. For operations
- ⇒ Operations like summing elements in a column is best to be performed parallelly.
- ⇒ Note: ~~The~~ The way to perform memory coalescing depends upon the task in hand.
- ⇒ For more info: check the code write.

Memory Hierarchy of CUDA

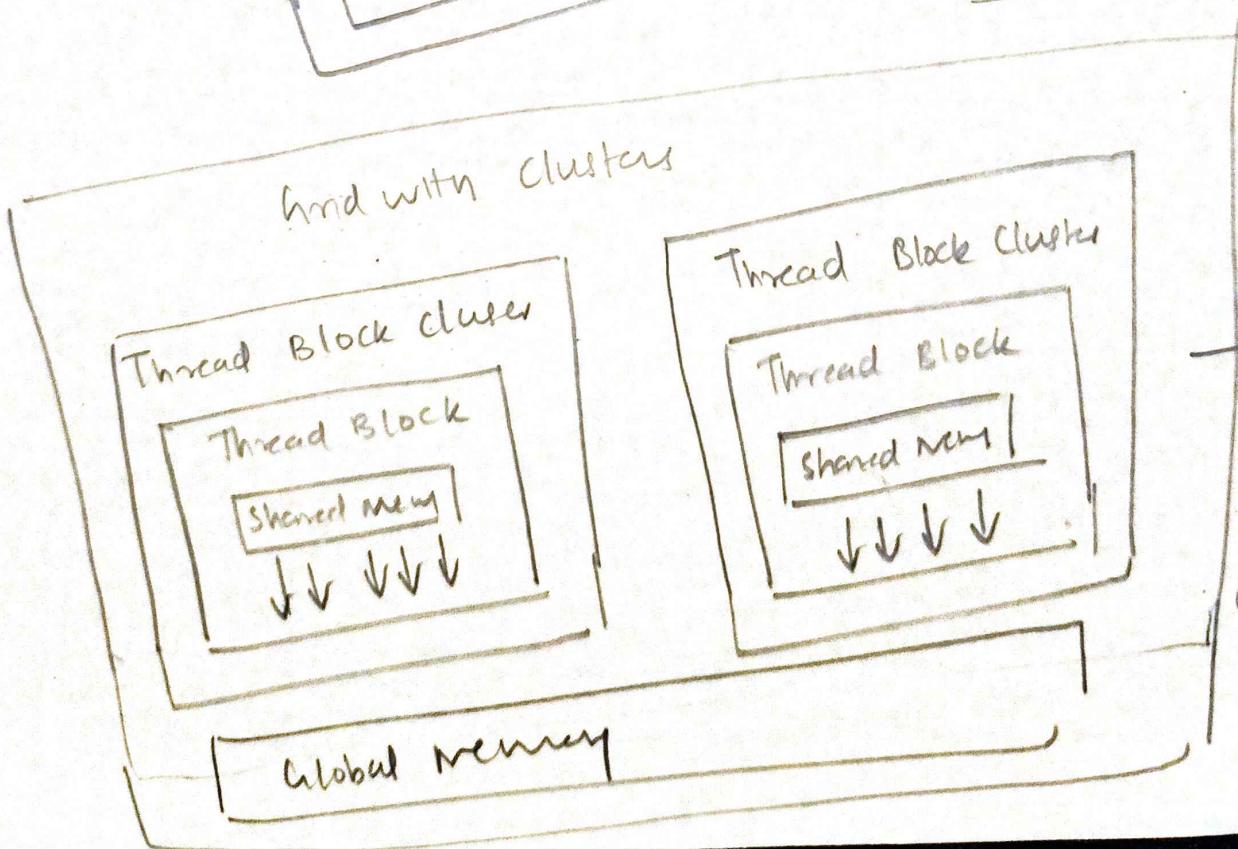
→ Each thread have has its private local memory.



→ Each thread block has a shared memory visible to all threads of the block with the same lifetime as the block.



→ Thread blocks in a Thread blocks cluster can perform, read, write and atomic operations on each other's shared memory.



→ All threads have access to Global memory.

Further, there are also two additional read-only memory spaces accessible by all threads: the constant & texture memory spaces.

→ The global, constant & texture memory spaces are optimized for different memory usage.

Global Memory vs Shared Memory

* Global memory is accessible to any thread or block on the device and can persist for the lifetime of the application, and it is relatively large memory space.

* Shared memory → A user-defined memory which has wider bandwidth than global memory that is shared between all threads in a block.

→ It does not persist after a kernel finishes executing.

→ It can be considered as a programmable cache.

→ Higher bandwidth implies that shared memory allows data to be transferred to and from the processing core at much faster rate than global memory.

* How to use shared memory

(1) caching memory: Read from global memory

→ Threads within a block may need to access same data multiple times during execution. If the data is stored in global memory, each access is slow and costly in terms of latency.

(2) Buffering output from threads for coalesced writes to Global memory

⇒ When the thread write results to global memory, writing one value at a time will be inefficient. That thus the results can be staged in the shared memory or in other words it could be used as temporary buffers and then be written back to the global memory in coalesced manner.

(3) staging Data for scatter/ data operations within a Block

⇒ In some algorithms different thread need to read/write data to non-contiguous regions of memory.

⇒ direct access of global memory will be inefficient.

⇒ direct access of global memory at a staging area.

⇒ Hence also, we can use shared memory as a staging area for loading or storing data efficiently.

+ setting up shared Memory

⇒ As said earlier, shared memory is ~~not~~ user defined and persists on till the execution of kernel is completed.

⇒ Numba provides function to define shared memory and also synchronize threads automatically which is often required while reading/writing to a shared memory.

⇒ While declaring a shared memory ⇒ we need to specify

- the size of shared array
- type of data in the array

(using numba type)

Note: While declaring the size/shape of array, we only can pass constant value as argument. Thus one cannot pass variables of a function and other variables provided by

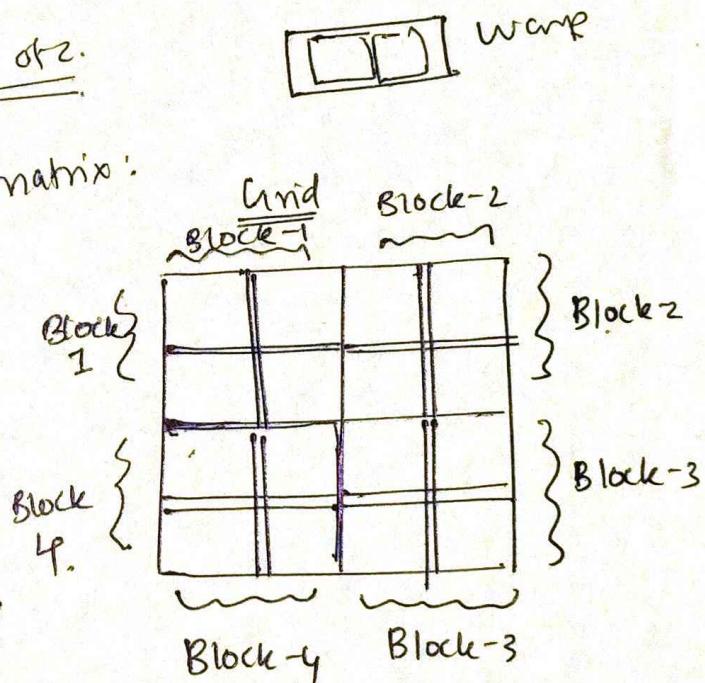
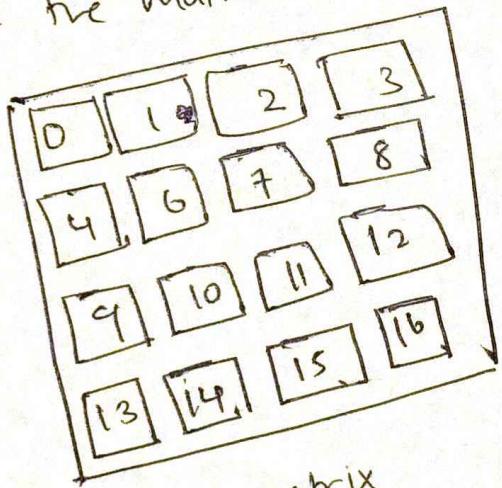
number (eg:- ~~cuda~~ numba.cuda.blockDim.x), or the calculation of ~~eg:-~~ ~~cu~~ cuda.griddim.

Using shared Memory for Memory Coalescing

\Rightarrow ~~Problem with~~ naive approach of not using shared memory

Let, the warp size be threads of 2.

Let the matrix be a 4×4 matrix:

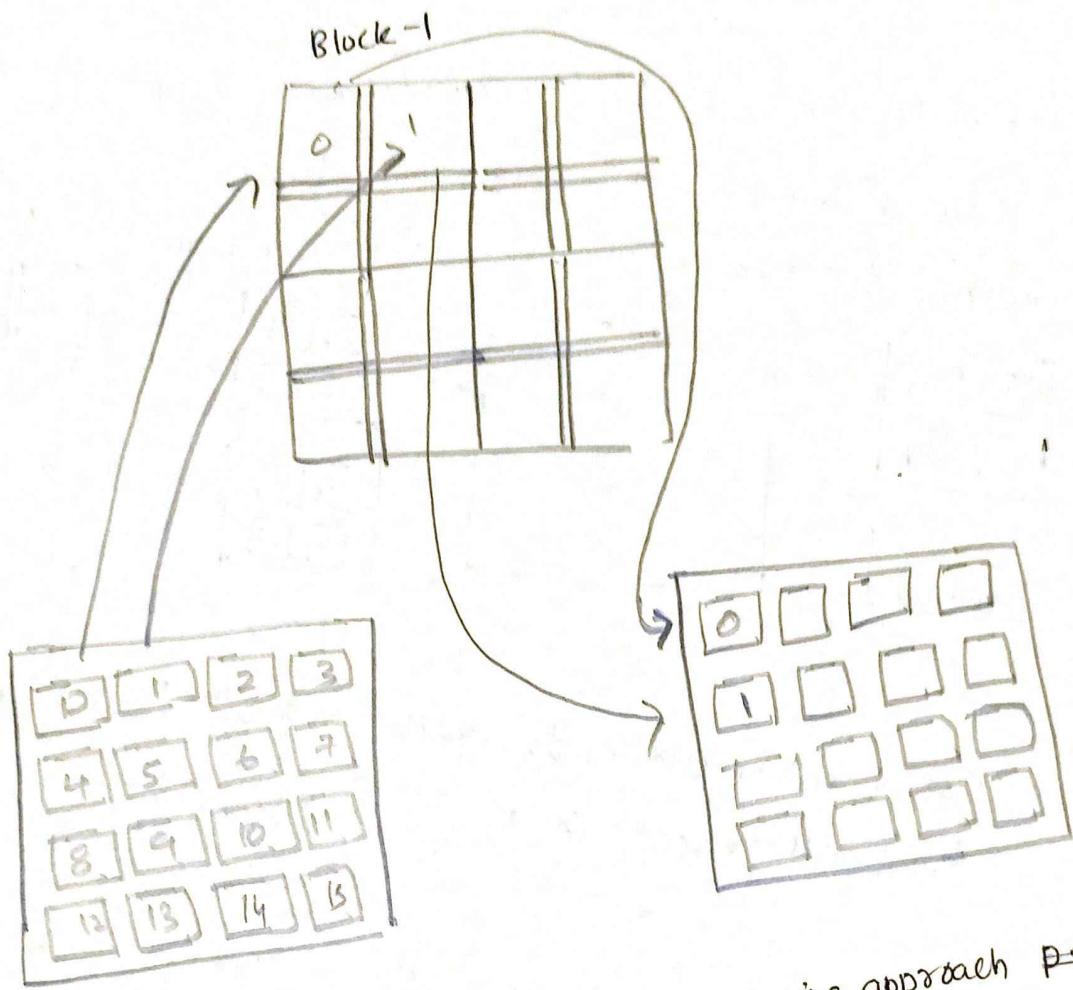


* Assume the task let task be taking transpose of a matrix
- naive approach is to make each thread read a element in

the matrix

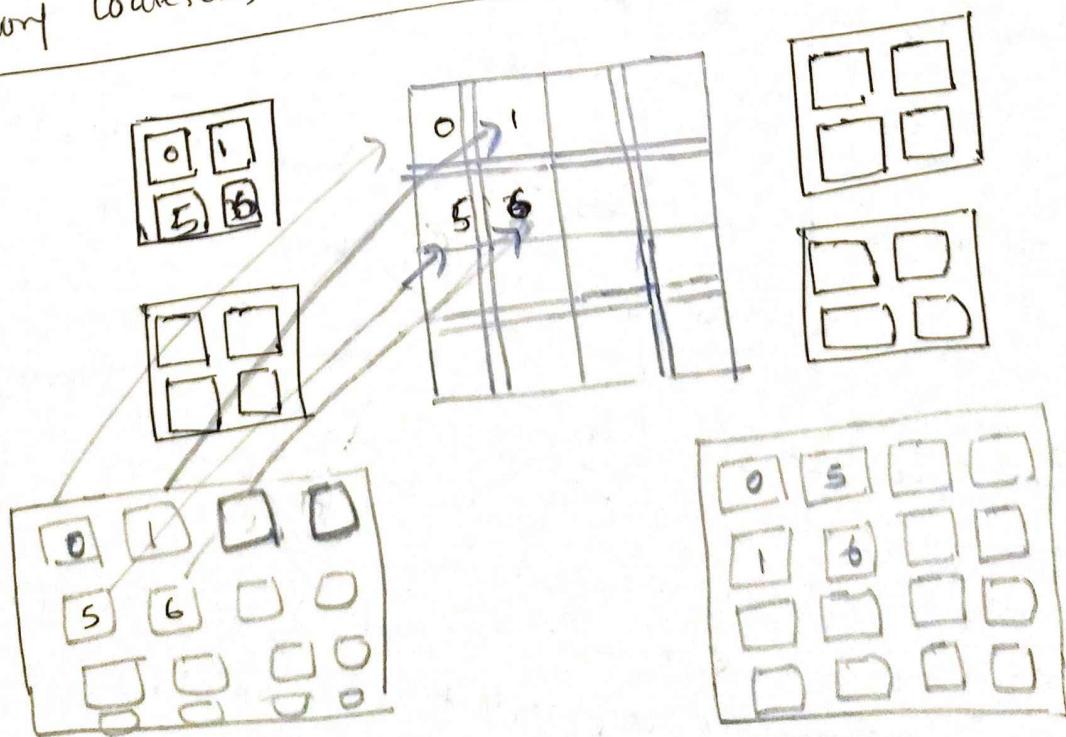
\Rightarrow Then write the element in the output ~~matrix~~ matrix in the transposed location.

- This scheme makes read in coalesced fashion and writing in uncoalesced fashion.
- Note: Accessing

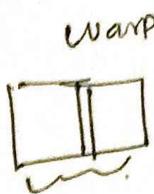
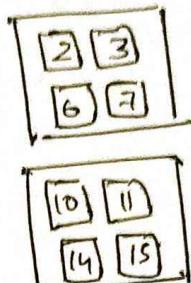
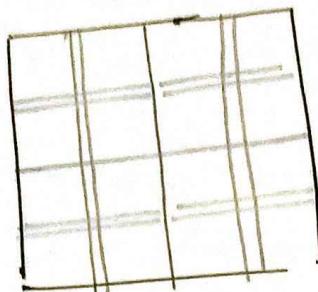
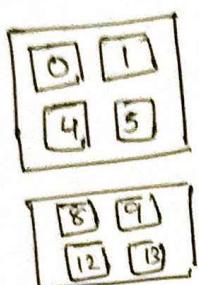


⇒ Here we can see, ~~thread per~~ naive approach ~~provide~~
performs coalesced reads while do not perform coalesced writes//.
⇒ This can be rectified by using shared memory.

Memory coalescing using shared Memory



- ⇒ Every thread in block makes coalesced reads and stored the data in shared memory of respective blocks.
After reading, the shared memory looks like this



Note: we need to use cuda::synchronize all threads ~~not~~ such that all threads in block is synchronized.

- ⇒ To perform each coalesced write, each warp reads column elements from shared memory and write ~~to~~ it row of ~~the~~ output array (or global memory in general)

↳ In a row major order convention followed by CUDA, writing to a ~~row~~ ~~row~~ is coalesced since it forms continuous memory.

- ⇒ This we can ~~perform from~~ take transpose of a matrix using coalesced reads & writes.

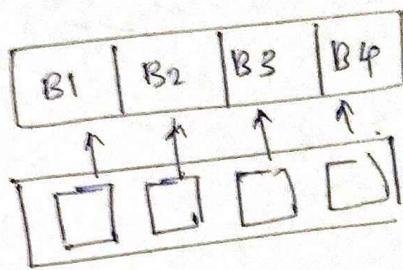
⇒ Here they performs column wise read from shared memory and write it to ~~the~~ appropriate row of output matrix (global memory).

⇒ I don't know how to column wise read is uncoalesced.

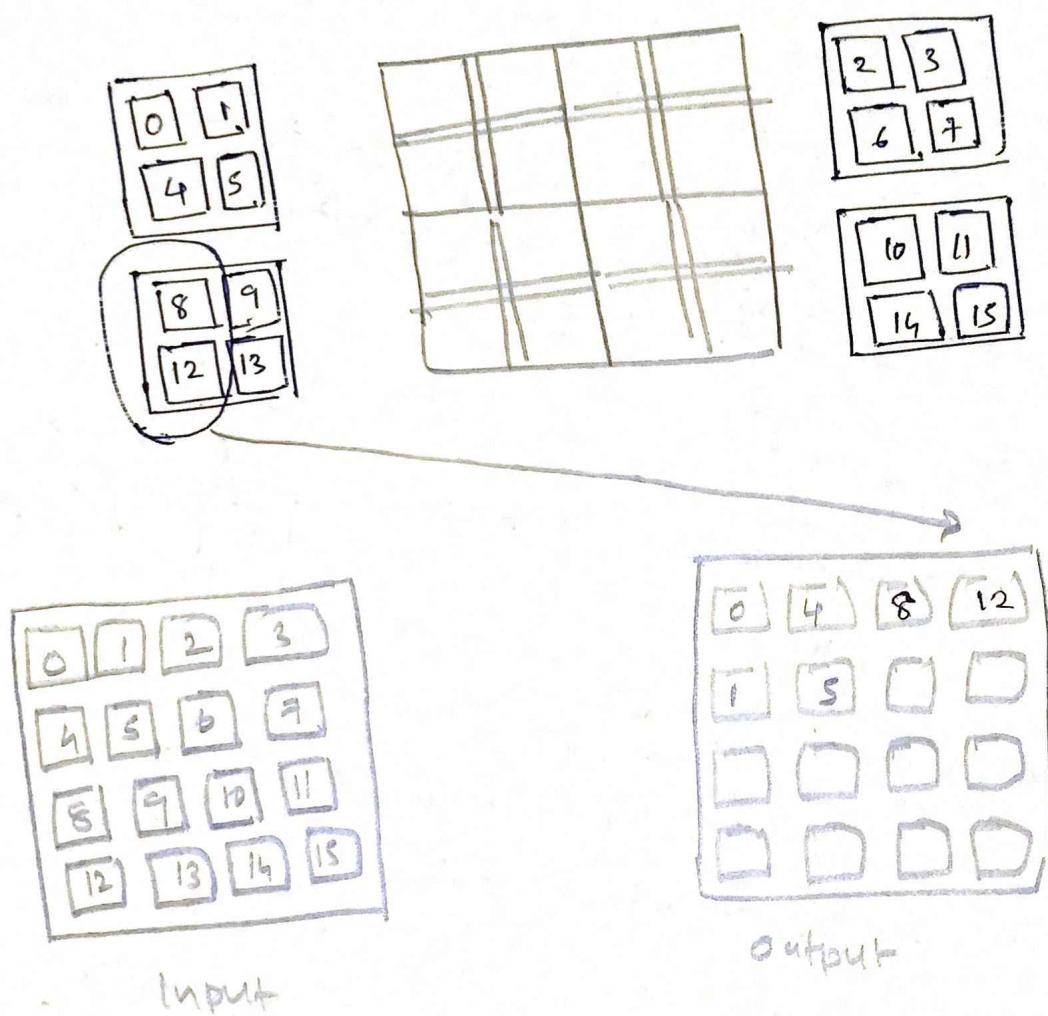
⇒ Further, columnwise read leads to something known as "Bank conflict".

Shared Memory Bank conflicts

- * Banks are the subunits of shared memory that allows parallel access
- * In most of hardware a shared memory will be divided to 32 banks
 - Since a warp has 32 threads, the shared memory is divided to 32 banks such that each thread can access one bank parallelly.
- * Banks in a shared memory form contiguous memory.
- * Each bank can serve one thread per cycle.
 - i.e., If you have 32 threads in a warp and each thread accessing a different bank, all 32 accesses will occur in parallel leading to high efficiency.
- * In CUDA, addresses are mapped to banks in a cyclic manner, if you declare a shared memory array, first element of array goes into bank 0, the second into bank 1 and so on. Once it reaches last bank, it wraps around and start at 0 again.
- * Bank conflict arise when ~~multiple~~ two or more threads tries to access same bank.
 - In such cases, the memory access becomes ~~semi~~ serialized, meaning the threads have to wait and access the bank one after the other, reducing the performance.



Blank conflict in the case of transpose of Matrix

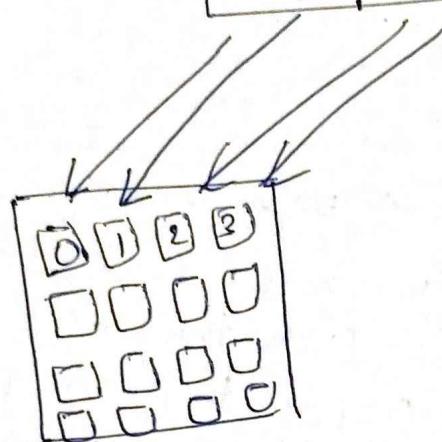


- ⇒ ~~In previous~~ In transpose of a matrix, we use shared memory to make reading and writing coalesced,
- we first read the data in coalesced manner (row wise) and store it in appropriate shared memory.
 - Then we read data column wise from shared memory and write it in the transposed location at output array.
 - It leads to Bank conflict.

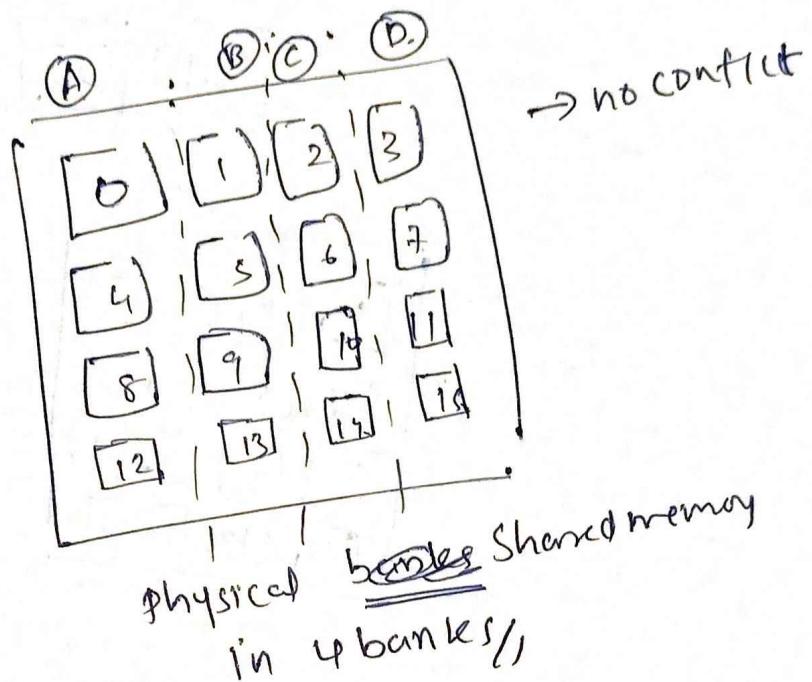
Warp + ~~Warp~~



(Here we assume warp to be composed of 4 threads)

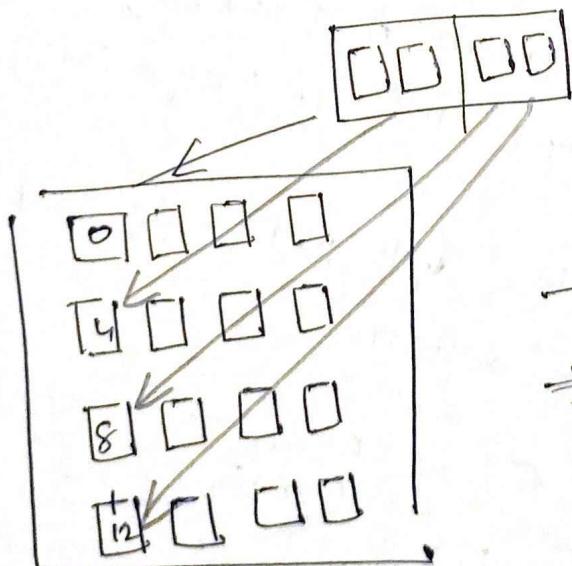


logical shared memory



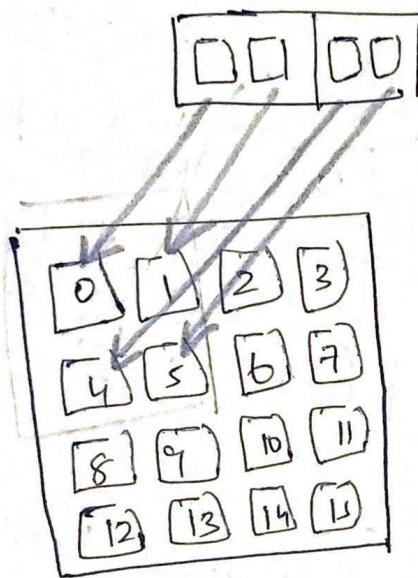
physical ~~banks~~ Shared memory
in 4 banks/

- ⇒ Elements in a column are stored in same bank.
- ⇒ Thus accessing elements in a column by more than one thread leads to Bank conflict.



→ leads to 4-way bank conflict.
⇒ the memory access will be serialized over ~4 cycles

⇒ Similarly consider the situation.



⇒ In this case, we have 2-way bank conflict that would require the memory access to be ~~serialized~~ serialized over 2 cycles.

⇒ So in the previous example, ~~there occurs~~ of transpose a matrix there occurs 2-way bank conflict.

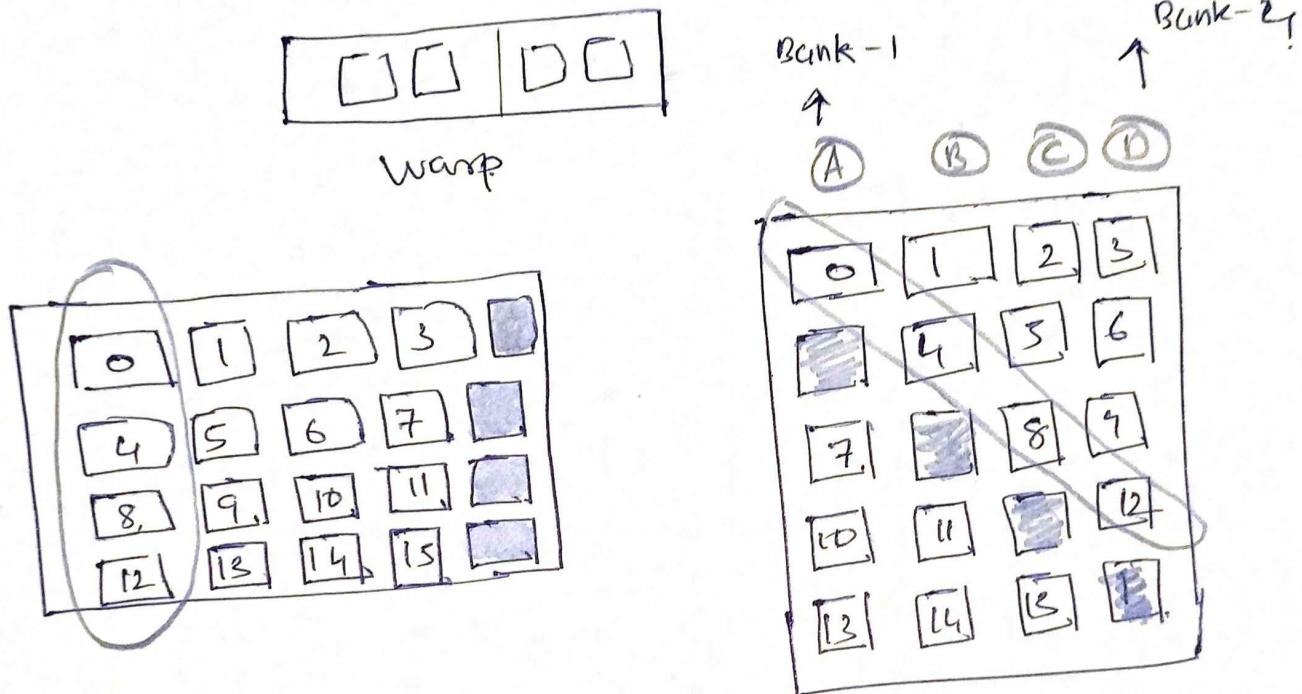
Solving ~~Bank~~ Bank conflict

⇒ In the case of transposing a matrix (or the cases where we need to access multiple elements in a column) we can use a simple trick of padding to prevent Bank conflict. For doing padding, we will pad the shared memory with an extra column.

i.e., if the shared memory array size = 4×4 (or tile), we will make it 4×5 .

- Then we write elements of input array to memory file to addresses in the range $[0:4][0:4]$
- ~~This when it~~

⇒ This padded memory file will be stored in the physical shared memory with 4 banks, in such a way that we can access multiple elements in a column of memory file without facing bank conflict.



⇒ Here one only here accessing the first column of memory file will be possible without memory bank conflict since the elements are stored in different banks.

⇒ Rowwise access also will be devoid of bank conflict.