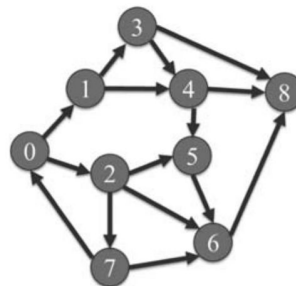# Graph Traversal

## Graph Search

# Graph and Matrix

- A graph is a data structure that represents the relationships between entities
  - A graph is a data structure that represents the relationships between entities
- Graphs are intrinsically related to sparse matrices
  - An intuitive representation of a graph is an adjacency matrix
  - Thus, graph computation can also be formulated in terms of sparse matrix operations
- Store them can be done using a different format
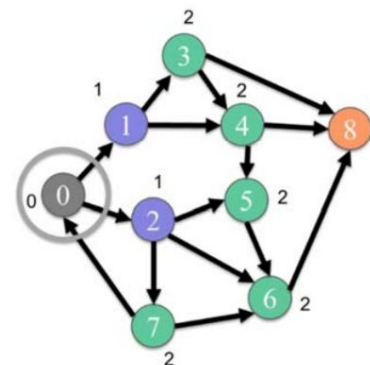  - CSR/CSC
  - COO

# Graph Usage Examples

- They are used to model different problems from different fields
  - Social media connection graphs
  - Driving directions
  - Telecommunication networks
  - Manufacturing process dependencies
  - Computation graph
  - 3D Meshes
  - Graphical models
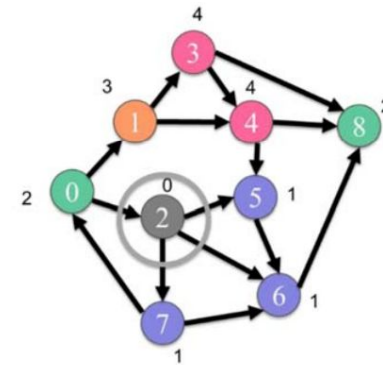- These graphs tend to be massive and sparse

POLITECNICO MILANO 1863

# Breadth-First Search (BFS)

- Given a source node $S$, find the number of steps required to reach each node $N$ in the graph
    - BFS is often used to discover the shortest number of edges that one needs to traverse to go from one vertex to another vertex of the graph
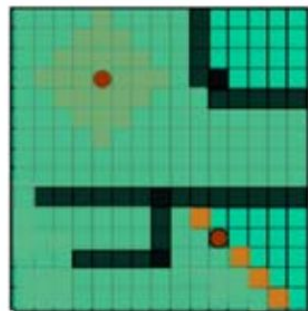


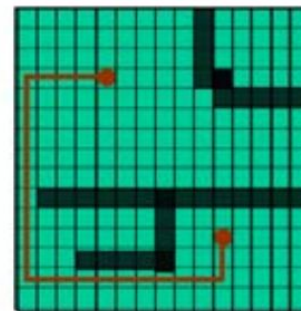(A) Vertex 0 is root          (B) Vertex 2 is root

# BSF Graph Computation Example

- While designing an integrated circuit chip, many electronic components need to be connected to complete the design
  - The routing software represents the chip as a grid of wiring blocks in which each block can potentially serve as a piece of a wire
- The maze routing application represents the chip as a graph
  - The routing blocks are vertices
  - An edge from vertex $i$ to vertex $j$ indicates that one can extend a wire from block $i$ to block $j$

● net terminal
■ blockage
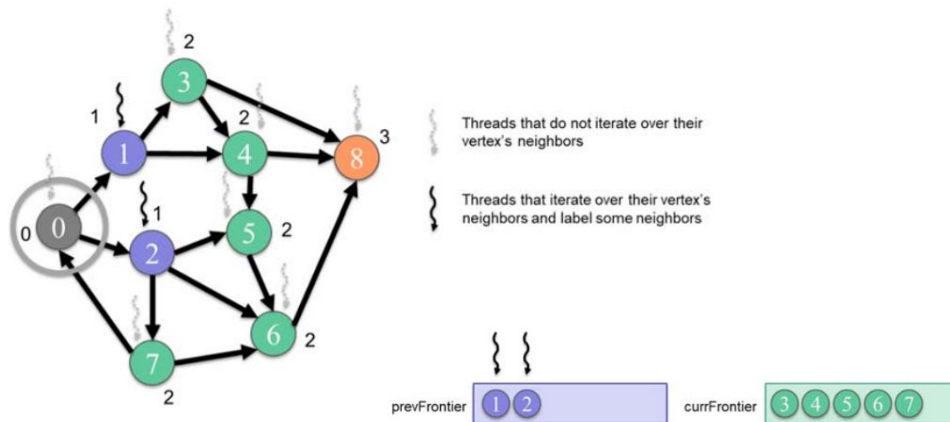
(A) Breadth-first search          (B) Identifying a routing path

# Vertix-Center BFS Computation

- Iterate or assign each thread to a vertex
  - For each iteration, check all incoming edges to see if the source vertex was just visited in the last iteration; if so, mark as visited in this iteration
  - Not very work efficient, complexity of $O(VL)$
    - $V$ number of vertices
    - $L$ length of the longest path
  - Difficult to detect stopping criterion

POLITECNICO MILANO 1863

# Frontier-Center Parallel BFS

- ## Assign threads to frontier vertices from the previous iteration
  - ### Add all the non-visited neighbors to the next frontier
  - ### The source will be the first element in the frontier
- ## Parallel execution between threads
  - ### A variable number of unnecessary threads are launched



Threads that do not iterate over their vertex's neighbors

Threads that iterate over their vertex's neighbors and label some neighbors

prevFrontier  1 2

currFrontier  3 4 5 6 7
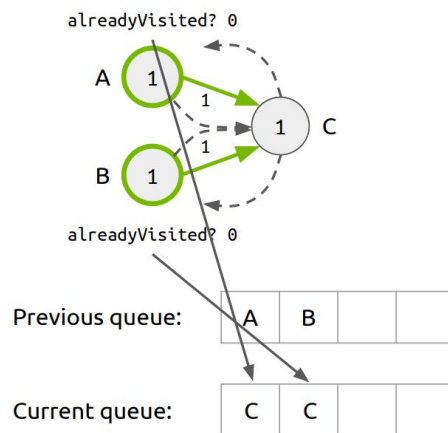
POLITECNICO MILANO 1863

7

# Problems and Optimizations

- Privatization
- Texture Memory
- Kernel launch overhead
- Sub Block-Level Queue
- Load Imbalance

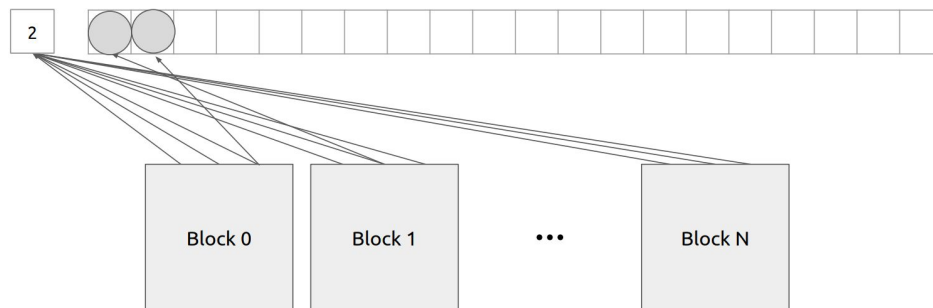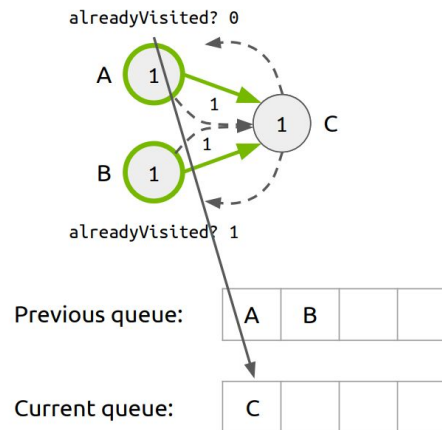| Optimization | Benefit to compute cores | Benefit to memory | Strategies |
|---|---|---|---|
| Maximizing occupancy | More work to hide pipeline latency | More parallel memory accesses to hide DRAM latency | Tuning usage of SM resources such as threads per block, shared memory per block, and registers per thread |
| Enabling coalesced global memory accesses | Fewer pipeline stalls waiting for global memory accesses | Less global memory traffic and better utilization of bursts/ cache lines | Transfer between global memory and shared memory in a coalesced manner and performing uncoalesced accesses in shared memory (e.g., corner turning) Rearranging the mapping of threads to data Rearranging the layout of the data |
| Minimizing control divergence | High SIMD efficiency (fewer idle cores during SIMD execution) | – | Rearranging the mapping of threads to work and/or data Rearranging the layout of the data |
| Tiling of reused data | Fewer pipeline stalls waiting for global memory accesses | Less global memory traffic | Placing data that is reused within a block in shared memory or registers so that it is transferred between global memory and the SM only once |
| Privatization (covered later) | Fewer pipeline stalls waiting for atomic updates | Less contention and serialization of atomic updates | Applying partial updates to a private copy of the data and then updating the universal copy when done |
| Thread coarsening | Less redundant work, divergence, or synchronization | Less redundant global memory traffic | Assigning multiple units of parallelism to each thread to reduce the price of parallelism when it is incurred unnecessarily |

# Frontier Output Interference
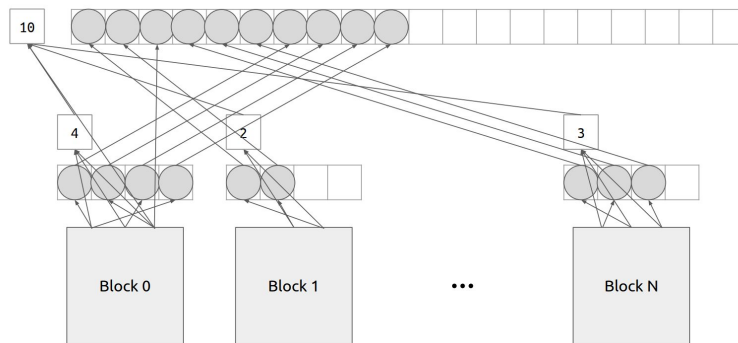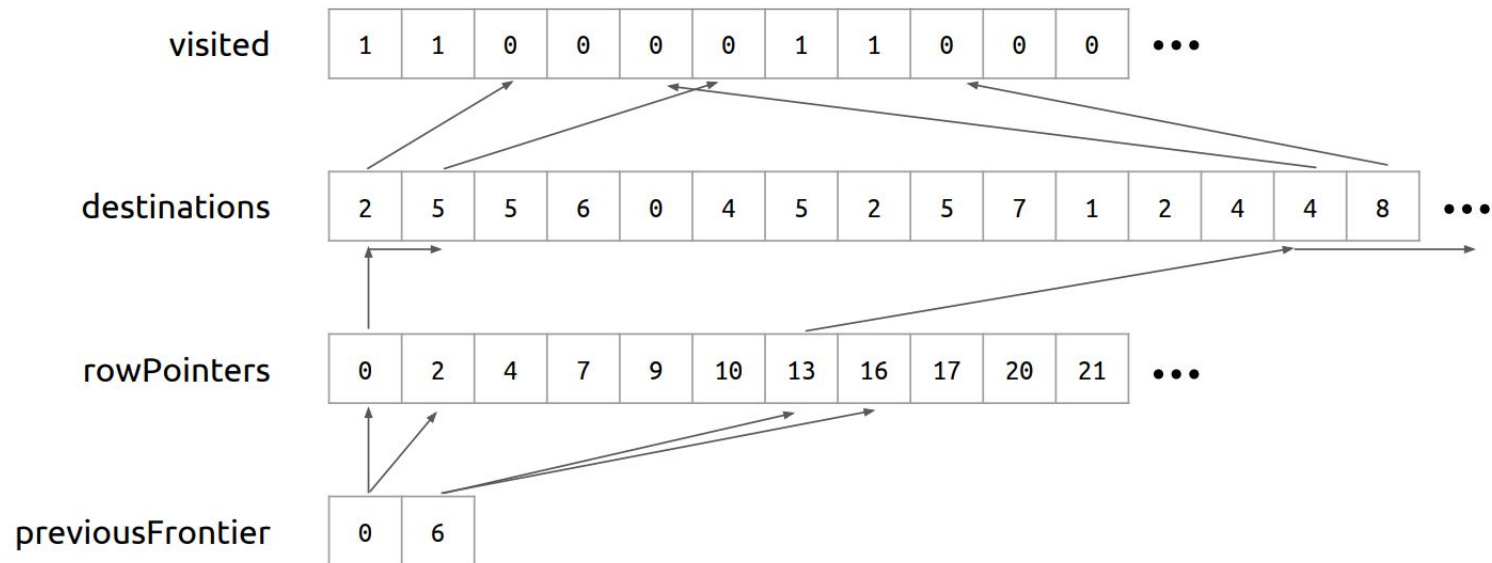
# Privatization

- Privatization reduces contention of atomics by applying partial updates to a private copy of the data, then updating the public copy when done
  - Threads will contend on the same data only with other threads in the same block
  - The local frontier and its counter can be stored in shared memory
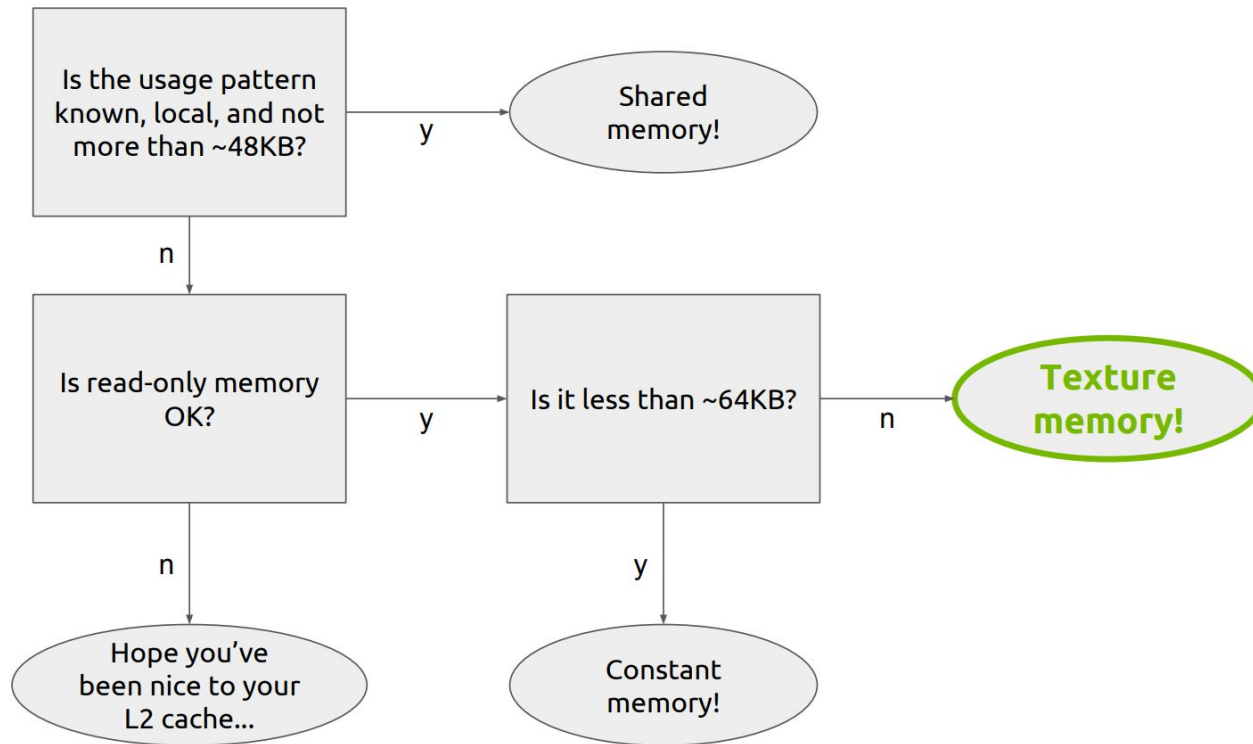    - Which enables lower-latency atomic operations

# Irregular global memory access

- Access patterns depend on graph structure and are unpredictable
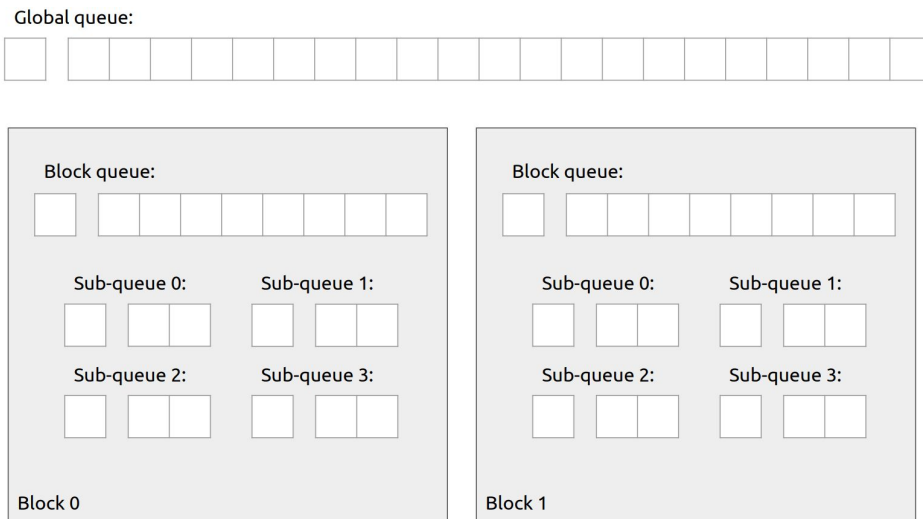
# Memory Optimization Flowchart

# Texture Memory

- Texture memory is another form of global memory
  - Like constant memory, it is aggressively cached for read-only access
- Originally developed and optimized for storing and reading textures for graphics applications
  - Has hardware-level support for 1-,2-, or3-D layouts and interpolated reads
  - The texture cache is also spatial layout-aware
- It can be useful for irregular access patterns with un-coalesced reads
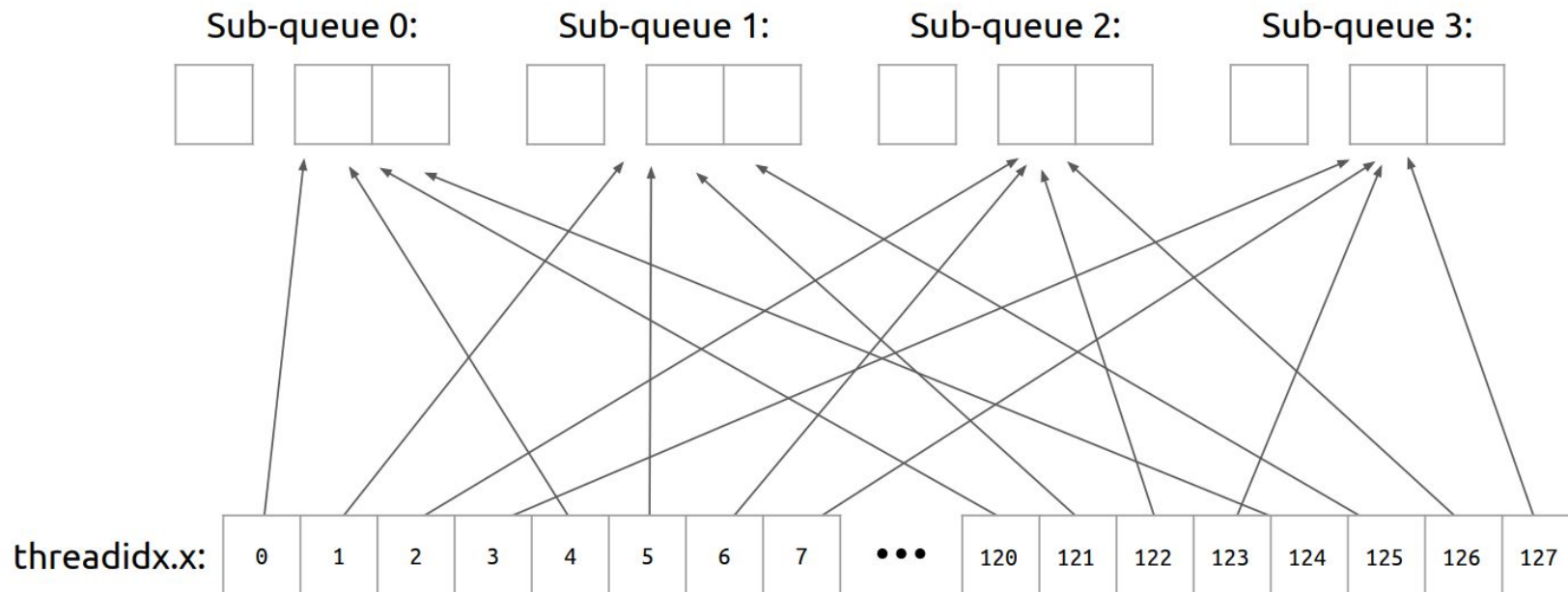
# Kernel Launch Overhead

- For some iterations of BFS (especially near the beginning), the <mark>frontier can be quite small</mark>
  - The benefits of parallelism only outweigh the kernel launch overhead when the frontier becomes large enough
- <mark>Use the CPU if the frontier size dips below some threshold</mark>

# Block-level queue length counter contention

- While the block-level queues reduced contention for global memory, the block-level counter is now the bottleneck
- We can extend the hierarchy by further splitting the block-level queue

Global queue:

Block queue:

Sub-queue 0:    Sub-queue 1:

Sub-queue 2:    Sub-queue 3:

Block 0

Block queue:

Sub-queue 0:    Sub-queue 1:

Sub-queue 2:    Sub-queue 3:

Block 1

# Sub-Queue

# Credits

- [CSE 599 I](#) Accelerated Computing - Programming GPUS - Parallel Pattern: Graph Traversal, Author *Tanner Schmidt*

# Thanks for your attention!

**Gianmarco Accordi**
*gianmarco.accordi@polimi.it*