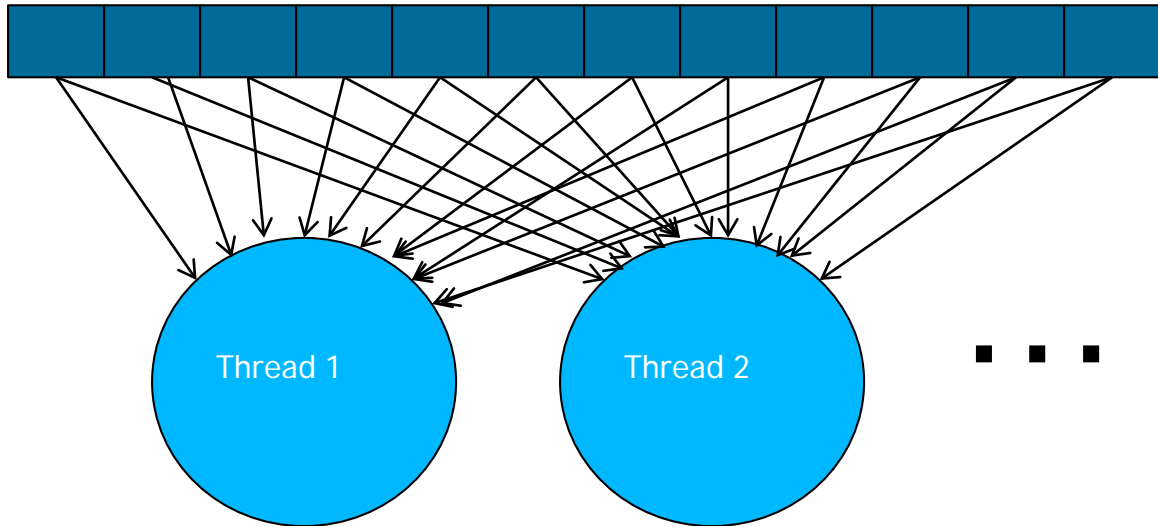# Objective

– To understand the motivation and ideas for tiled parallel algorithms
  – Reducing the limiting effect of memory bandwidth on parallel kernel performance
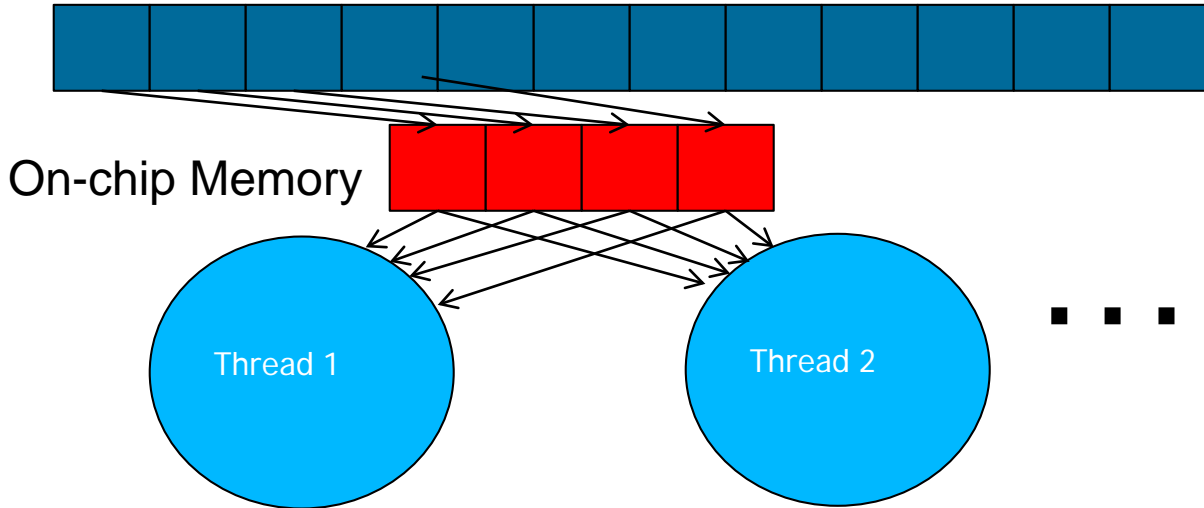  – Tiled algorithms and barrier synchronization

# Global Memory Access Pattern
# of the Basic Matrix Multiplication Kernel

Global Memory

# Tiling/Blocking - Basic Idea

Global Memory



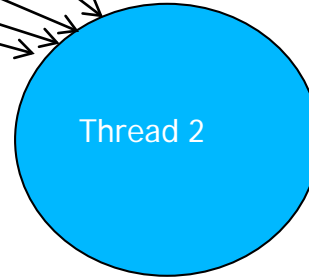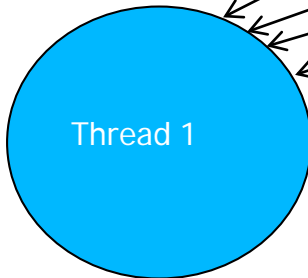Divide the global memory content into tiles

Focus the computation of threads on one or a small number
of tiles at each point in time

# Tiling/Blocking - Basic Idea

Global Memory

# Basic Concept of Tiling

– In a congested traffic system, significant reduction of vehicles can greatly improve the delay seen by all vehicles
  – Carpooling for commuters
  – Tiling for global memory accesses
    – drivers = threads accessing their memory data operands
    – cars = memory access requests

# Some Computations are More Challenging to Tile

– Some carpools may be easier than others
  – Car pool participants need to have similar work schedule
  – Some vehicles may be more suitable for carpooling
– Similar challenges exist in tiling
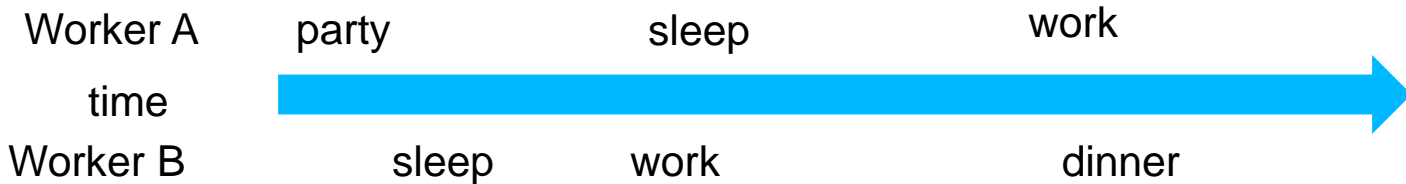
# Carpools need synchronization.

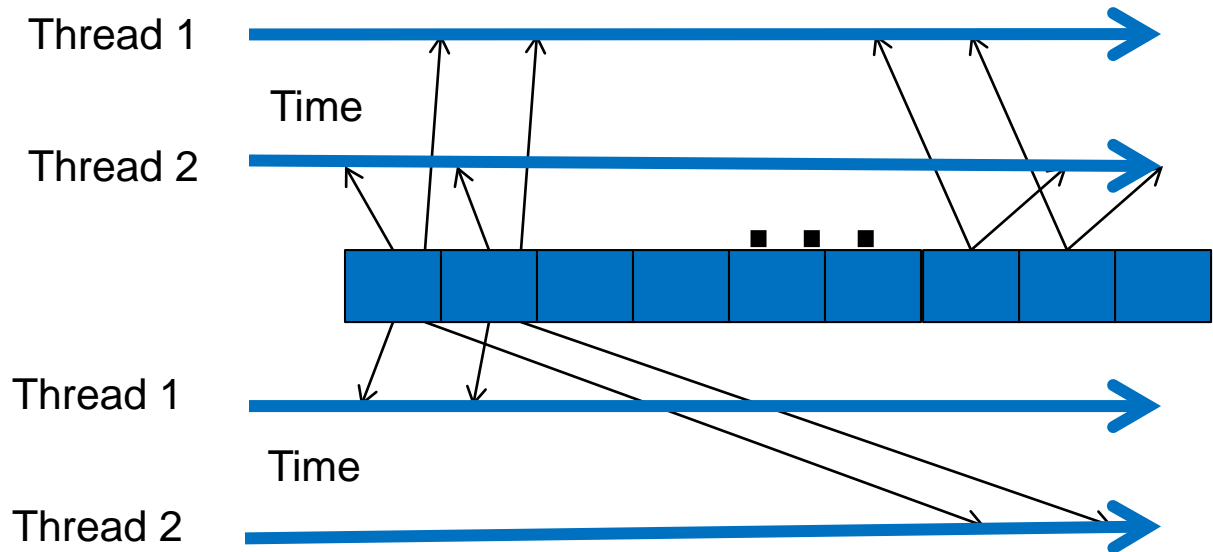– Good: when people have similar schedule

| Worker A | sleep | work | dinner |
|----------|-------|------|--------|
| Time | | | |
| Worker B | sleep | work | dinner |

# Carpools need synchronization.

– Bad: when people have very different schedule

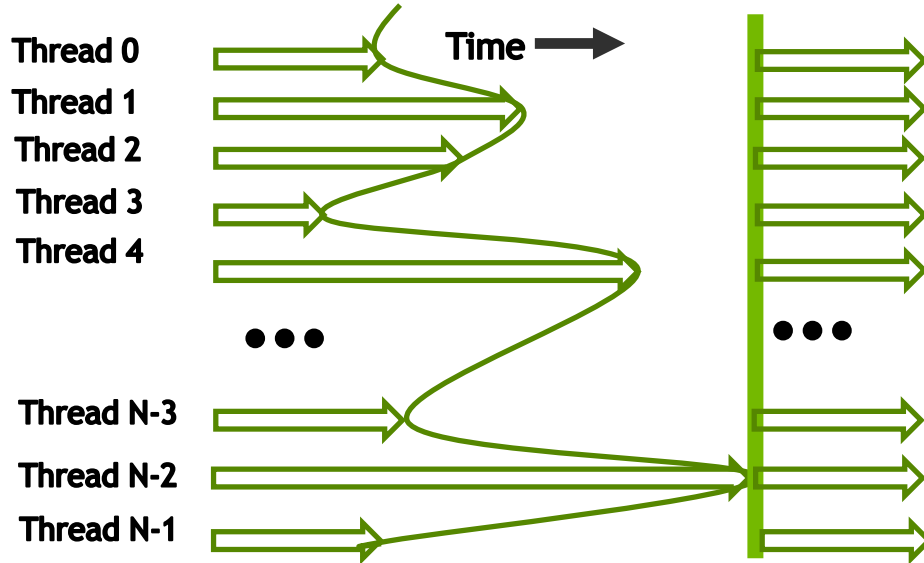| | | | |
|---|---|---|---|
| Worker A | party | sleep | work |
| time | | | |
| Worker B | sleep | work | dinner |

# Same with Tiling

– Good: when threads have similar access timing



– Bad: when threads have very different timing

# Barrier Synchronization for Tiling

# Outline of Tiling Technique

– Identify a tile of global memory contents that are accessed by multiple threads
– Load the tile from global memory into on-chip memory
– Use barrier synchronization to make sure that all threads are ready to start the phase
– Have the multiple threads to access their data from the on-chip memory
– Use barrier synchronization to make sure that all threads have completed the current phase
– Move on to the next tile

GPU Teaching Kit

Accelerated Computing