# Communication-Avoiding Gang Scheduling of Resources in Tessellation OS
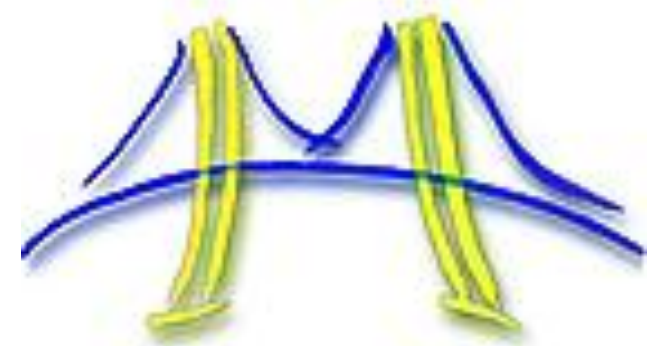
Juan A. Colmenares and John D. Kubiatowicz

Par Lab, CS Division, UC Berkeley
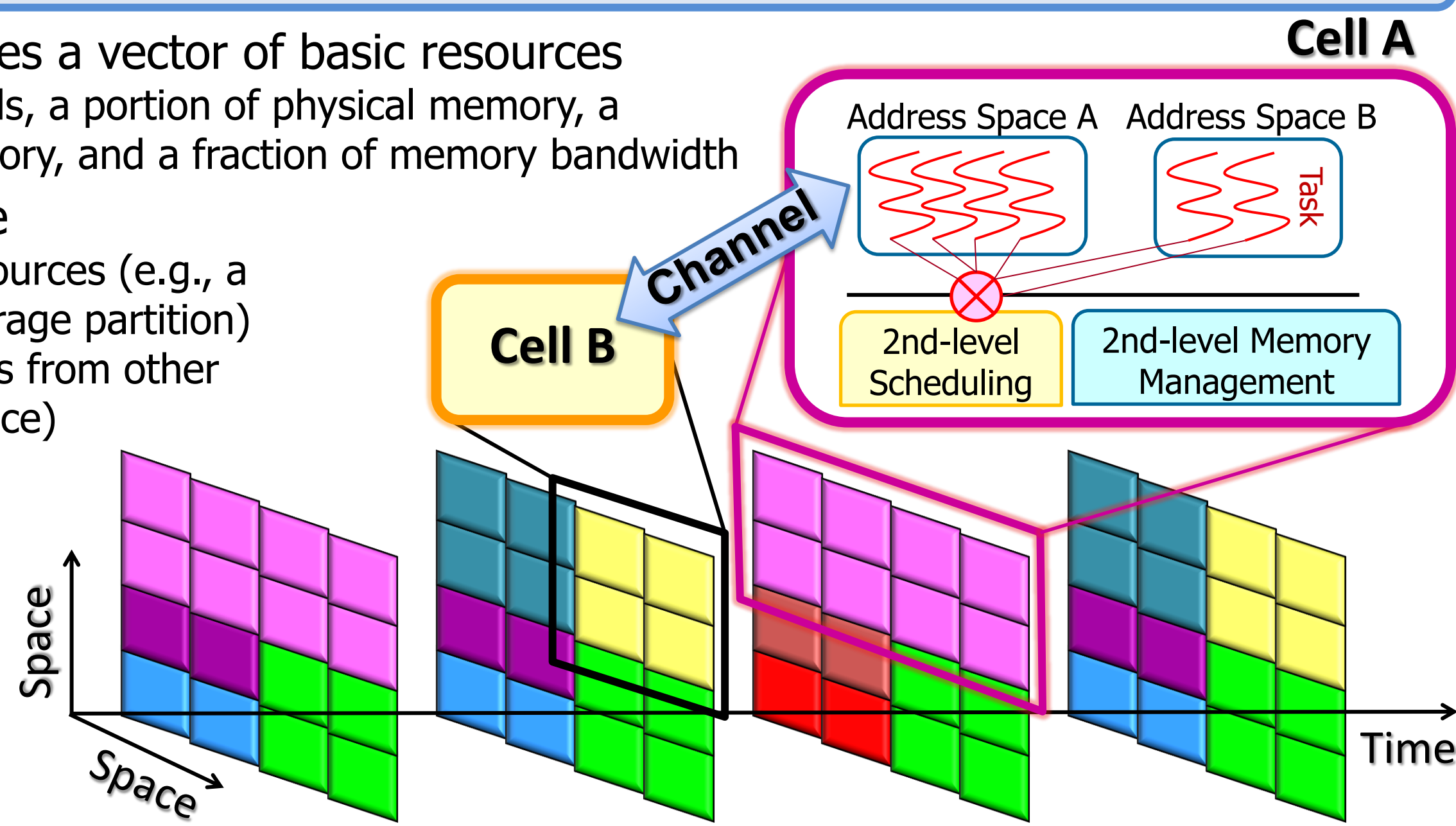
**Parallel Computing Laboratory**

## 1. Motivation

• Performance gap between computation and communication
• Communication cost increases with core count
• We want to determine the actual benefits of trading communication for computation in resource gang-scheduling, which is key to Tessellation OS?

## 2. Basic Goals in Tessellation OS

• Support a simultaneous mix of high-throughput parallel, interactive, and real-time applications
• Allow applications to consistently deliver performance

## 3. Space-time Partitioning and Two-level Scheduling

• A **Spatial Partition** receives a vector of basic resources
  – A number of hardware threads, a portion of physical memory, a portion of shared cache memory, and a fraction of memory bandwidth
• A Partition may also receive
  – Exclusive access to other resources (e.g., a hardware device and raw storage partition)
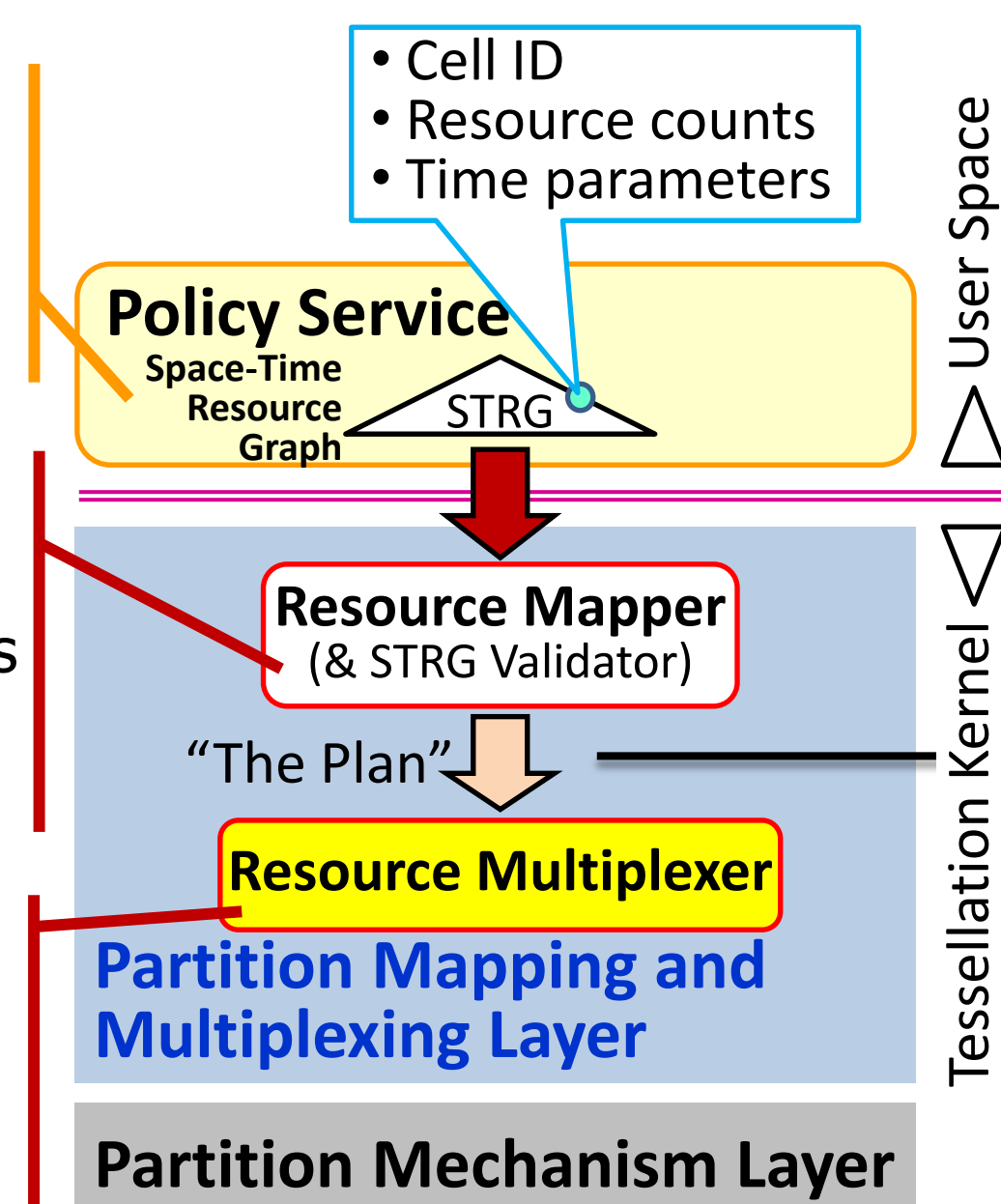  – Guaranteed fractional services from other partitions (e.g., network service)

• Spatial partitioning may vary over time
  – Partitions can be **time multiplexed**; resources are **gang-scheduled**
  – Partitioning adapts to needs of the system

**Cell A**

Address Space A   Address Space B   Task

**Cell B**   Channel

2nd-level Scheduling | 2nd-level Memory Management

Space / Space / Time

• **The Cell**: Our partitioning abstraction
  – User-level software container with guaranteed access to resources
• Basic properties of a cell
  – Full control over resources it owns when mapped to hardware
  – One or more address spaces
  – Communication channels

• **Scheduling at Level 1**: Coarse-grained resource allocation and distribution at the cell level
• **Scheduling at Level 2**: Fine-grained application-specific scheduling *within* a cell

## 4. Resource Allocation Architecture

*A Partial and Simplified View*

• Distributes resources among cells
• Establishes how cells should be time multiplexed

• Assigns specific resources to cells
• Produces only feasible mappings
  – Rejects invalid and infeasible STRGs

• Determines when cells should be activated and suspended
• Actually activates and suspends cells

• Cell ID
• Resource counts
• Time parameters

**Policy Service**
Space-Time Resource Graph   STRG

**Resource Mapper** (& STRG Validator)

"The Plan"

**Resource Multiplexer**

**Partition Mapping and Multiplexing Layer**

**Partition Mechanism Layer**

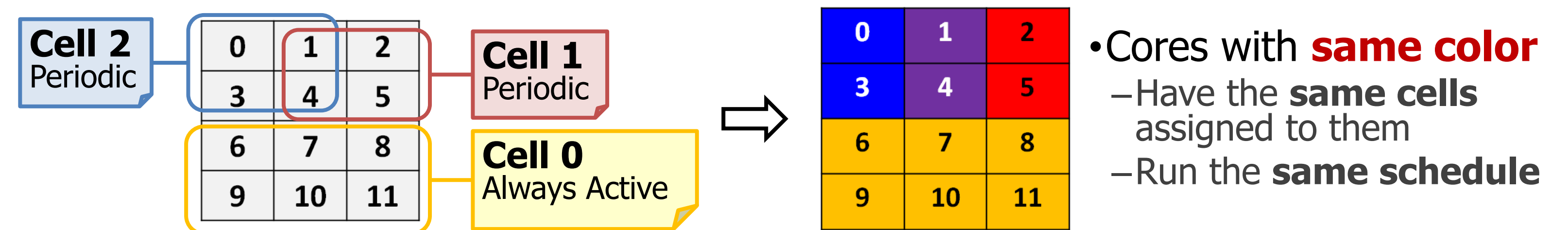User Space / Tessellation Kernel

**Separation between Mapper and Multiplexer**

• Decision Making Process
  – Mapping of multiple resources
  – Centralized because it requires global knowledge
  – Often expensive

• Execution
  – Relatively simple and fast
    • If the set of cells does not change
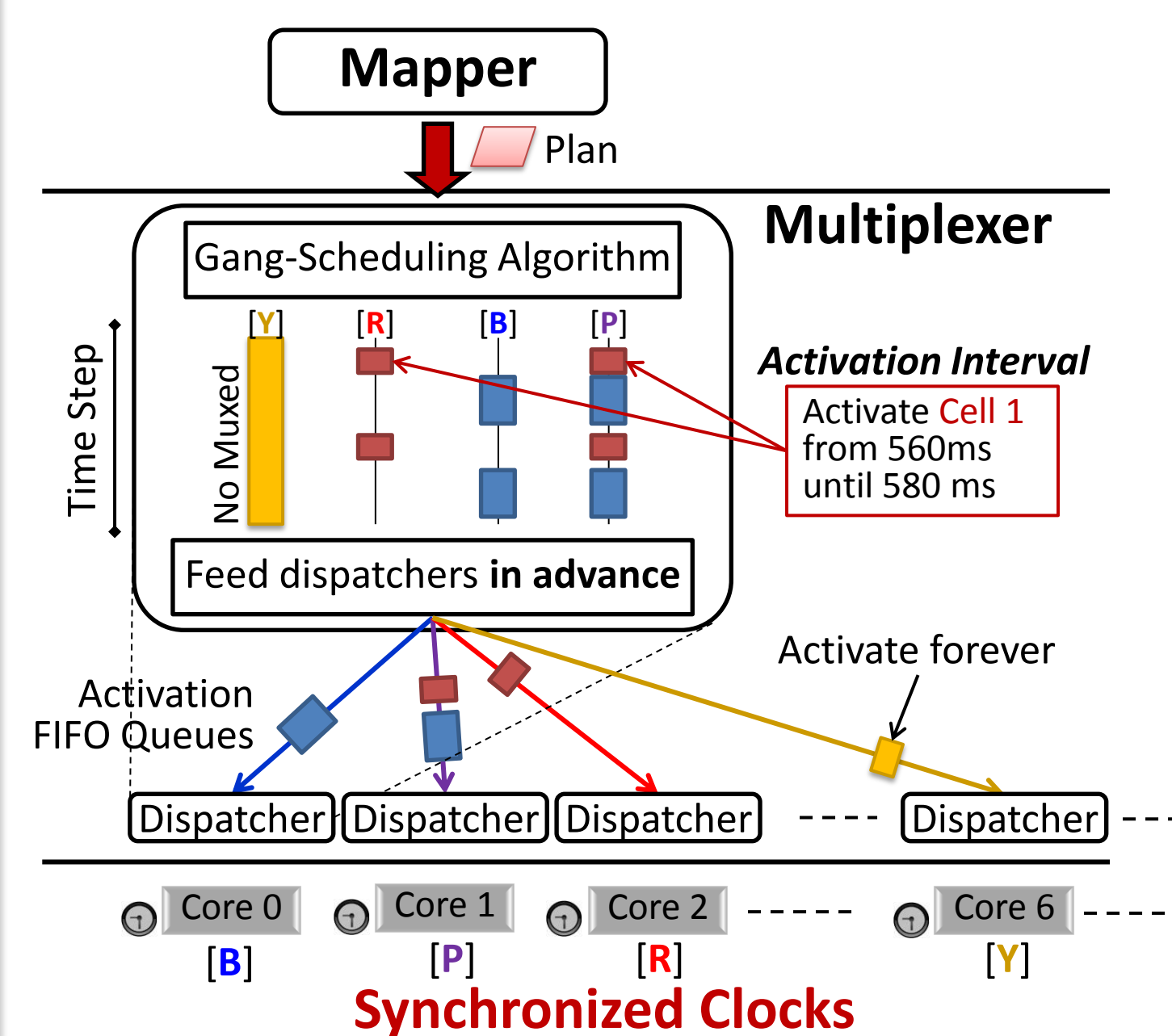  – Allows us to explore decentralized approaches

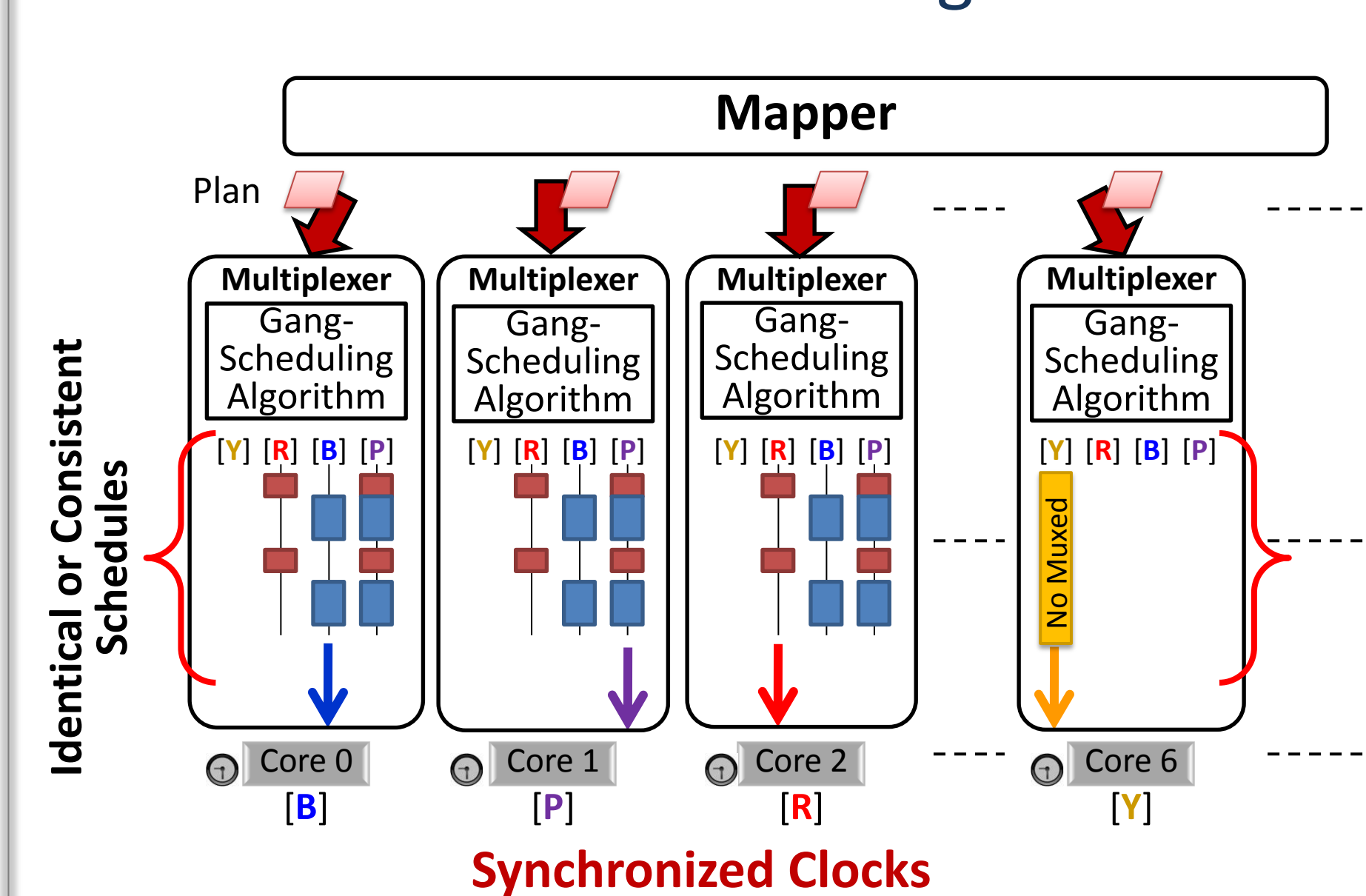## 5. Communication-Free and Centralized Multiplexers

**Sample Cell-Core Mapping**
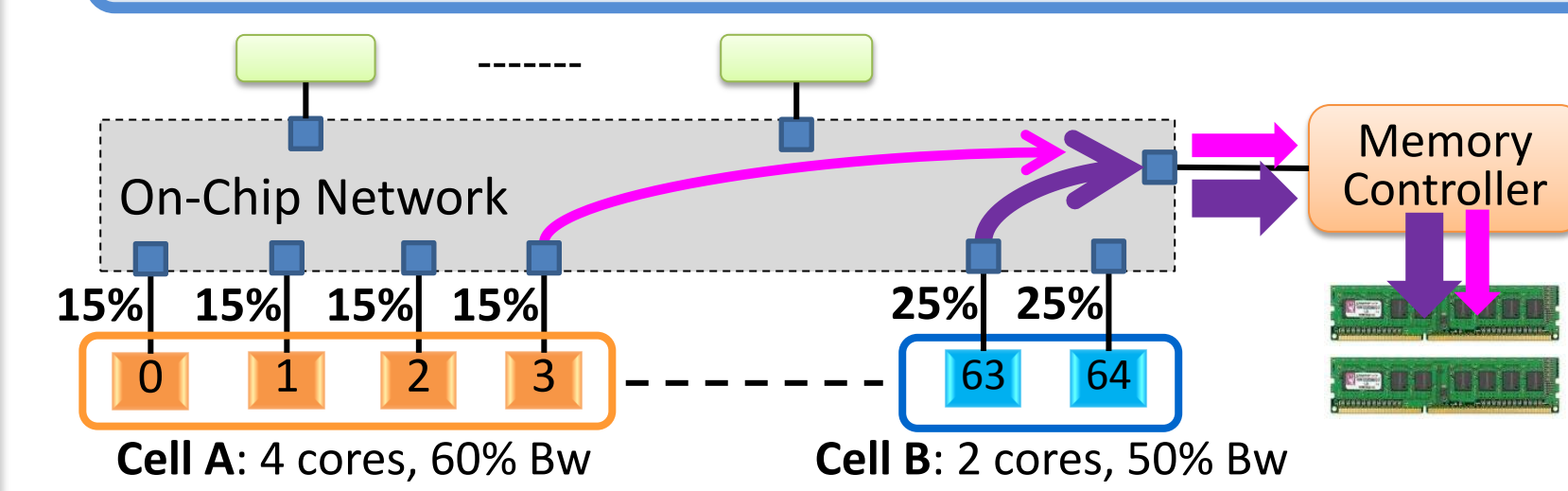Produced by the **Mapper** after checking feasibility

**Cell 2** Periodic   **Cell 1** Periodic
**Cell 0** Always Active

| 0 | 1 | 2 |
| 3 | 4 | 5 |
| 6 | 7 | 8 |
| 9 | 10 | 11 |

• Cores with **same color**
  – Have the **same cells** assigned to them
  – Run the **same schedule**

### A Centralized Version

**Mapper** → Plan

**Multiplexer**

Gang-Scheduling Algorithm

[Y] [R] [B] [P]   No Muxed

*Activation Interval*
Activate Cell 1 from 560ms until 580 ms

Feed dispatchers **in advance**

Activation FIFO Queues   Activate forever

Dispatcher | Dispatcher | Dispatcher | ---- | Dispatcher

Core 0 [B] | Core 1 [P] | Core 2 [R] | ---- | Core 6 [Y]

**Synchronized Clocks**

### Communication-avoiding Version

**Mapper**

Plan

Multiplexer | Multiplexer | Multiplexer | ---- | Multiplexer

Gang-Scheduling Algorithm (×4)

[Y] [R] [B] [P]   No Muxed

Identical or Consistent Schedules

Core 0 [B] | Core 1 [P] | Core 2 [R] | Core 6 [Y]

**Synchronized Clocks**

## 6. Gang-Scheduling of Cores and Memory Bandwidth

On-Chip Network

15% 15% 15% 15%   25% 25%

Memory Controller   DRAM

| 0 | 1 | 2 | 3 | ---- | 63 | 64 |

**Cell A**: 4 cores, 60% Bw   **Cell B**: 2 cores, 50% Bw
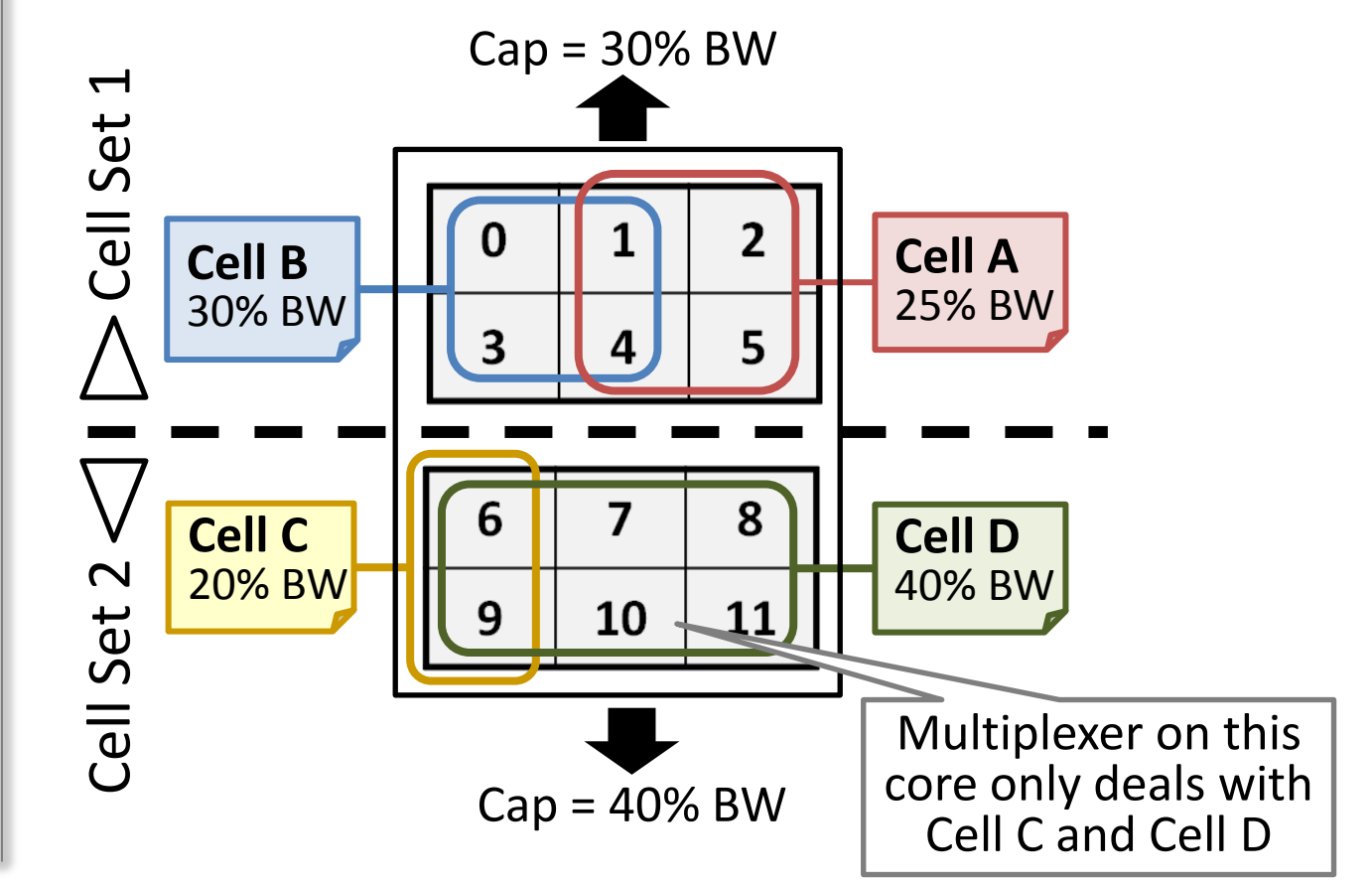
Memory bandwidth is a global **shared** resource

Assume that each core can request its own guaranteed minimum fraction of memory bandwidth

### Cell-Core Mapping

**Cell C** • Periodic • 30% BW
**Cell B** • Periodic • 50% BW
**Cell A** • Always Active • 40% BW

| 0 | 1 | 2 |
| 3 | 4 | 5 |
| 6 | 7 | 8 |
| 9 | 10 | 11 |

**BW** Guaranteed Minimum Memory Bandwidth

100%

**BW**: 50% | 30% | 50% / 40%

Need to start here due to bandwidth constraint

Cores: {2, 5}   Cell B | Cell B
Cores: {0,1,3,4}   Cell C
Cores: {6, 7,…, 11}   Cell A

Time

• Independent sets of cells with caps on guaranteed bandwidth
  — Less computation cost for each multiplexer

Cap = 30% BW

**Cell B** 30% BW | 0 1 2 / 3 4 5 | **Cell A** 25% BW

**Cell C** 20% BW | 6 7 8 / 9 10 11 | **Cell D** 40% BW

Cell Set 1 / Cell Set 2

Cap = 40% BW

Multiplexer on this core only deals with Cell C and Cell D

## 7. Status

• Initial versions of the gang-scheduling algorithm, centralized multiplexer, and communication-free multiplexer exist and they are being tested
• Implementation of the Mapper is underway