

Early Scheduling in Parallel State Machine Replication

Eduardo Alchieri, Fernando Dotti, and Fernando Pedone

Universidade de Brasilia, Pontifca Universidade Católica do Rio Grande do Sul, and University of Lugano



State Machine Replication (SMR)

🔹 Fundamental approach to fault tolerance

🔹 Google Spanner

🔹 Apache Zookeeper

🔹 Windows Azure Storage

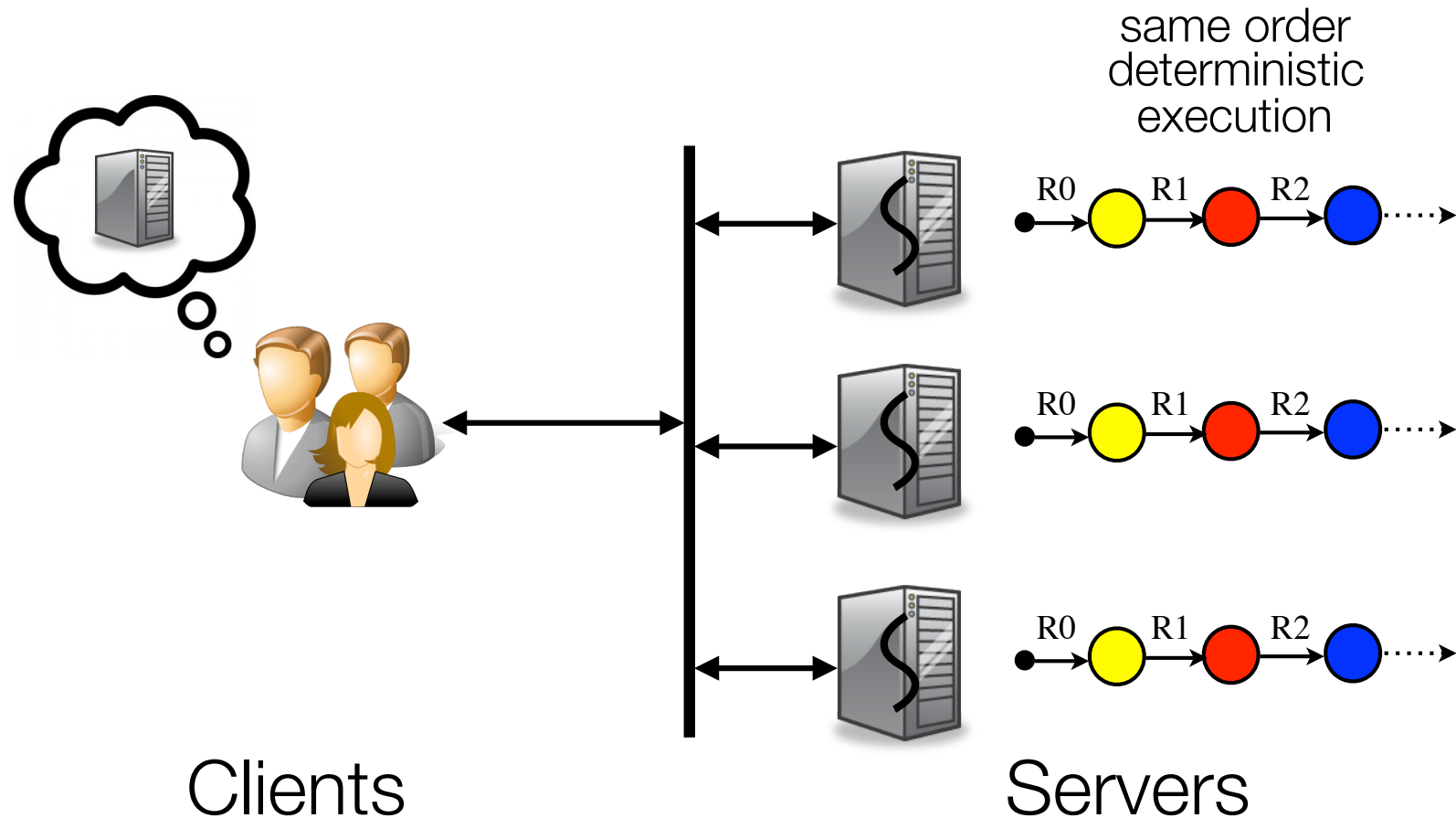
🔹 MySQL Group Replication

🔹 Galera Cluster

🔹 Blockchain, ...

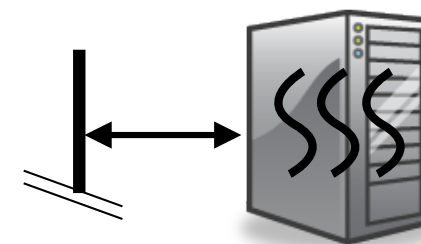


SMR is intuitive and simple



Parallel State Machine Replication

Key observation



- Independent requests can execute concurrently
- Conflicting requests must be serialized and executed in the same order by the replicas
- Two requests conflict if they access common state and at least one of them updates the state



Parallel State Machine Replication

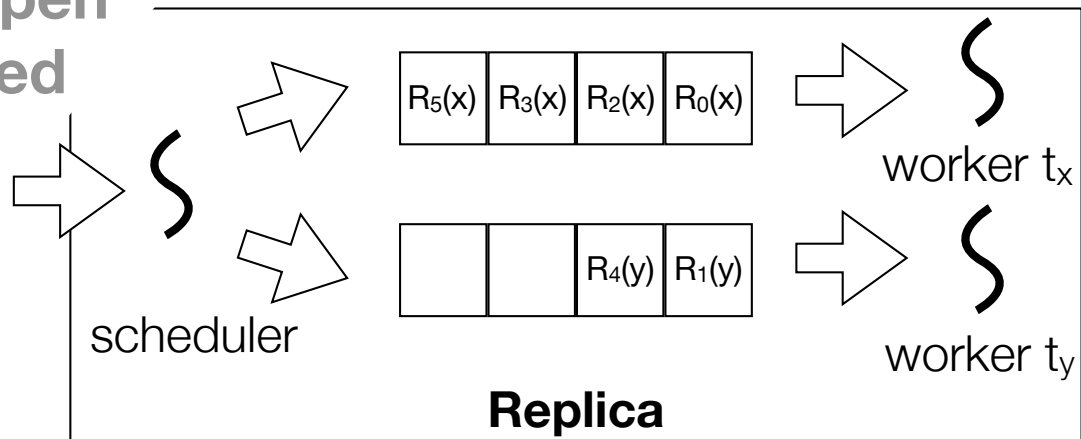
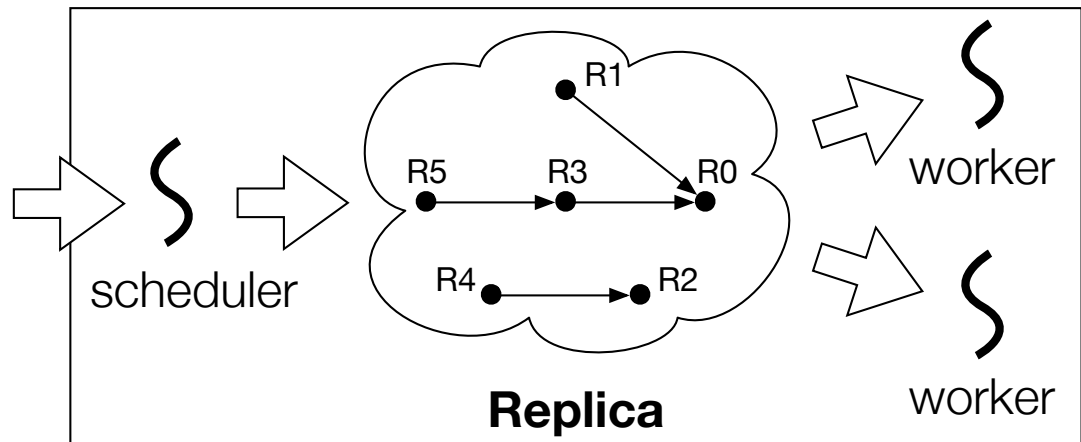
Late scheduling

Scheduling happens after requests are ordered

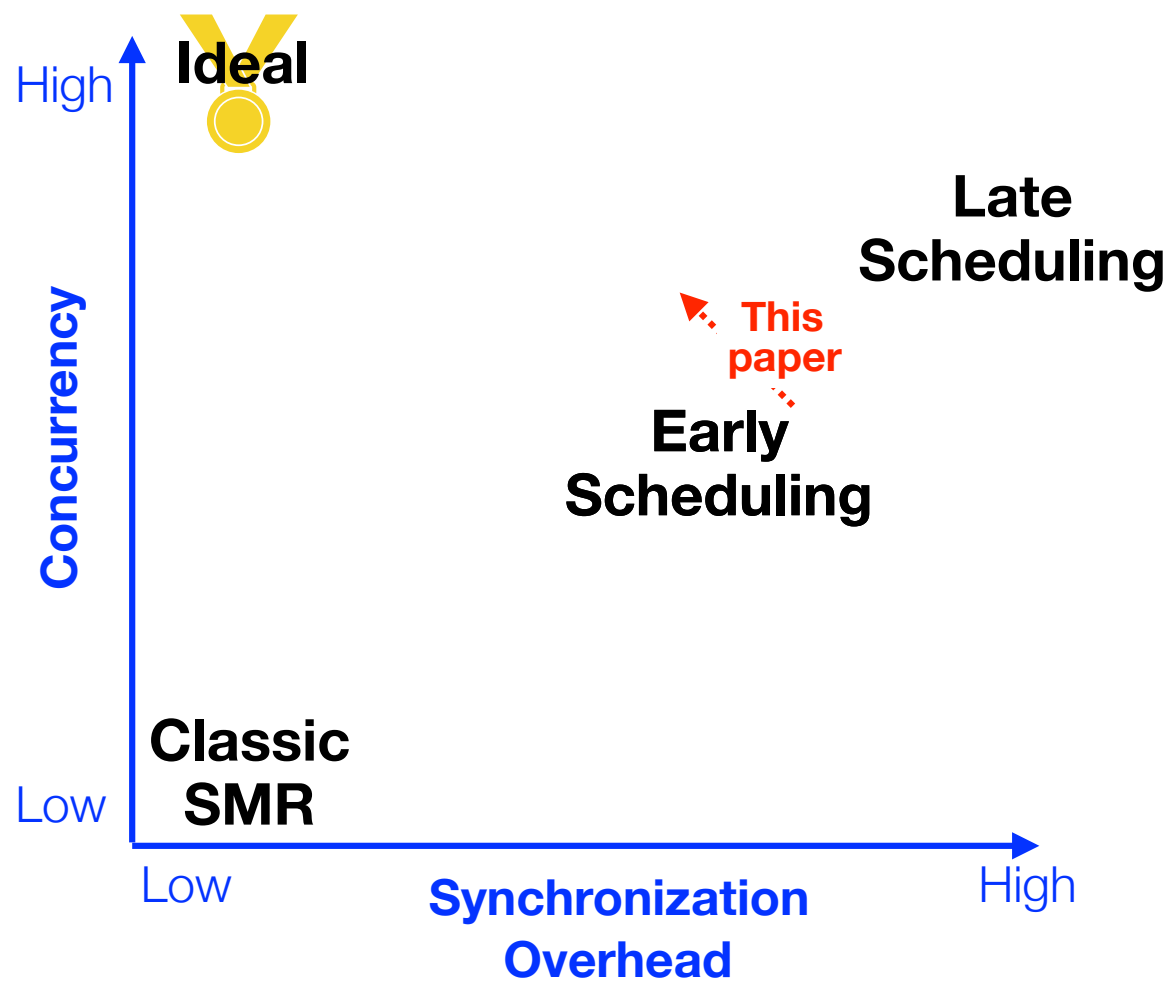
Early scheduling

Scheduling decisions happen before requests are ordered

E.g., worker t_x executes requests on X, worker t_y executes requests on Y






Scheduling tradeoff



Our contributions

Generalization of Early Scheduling

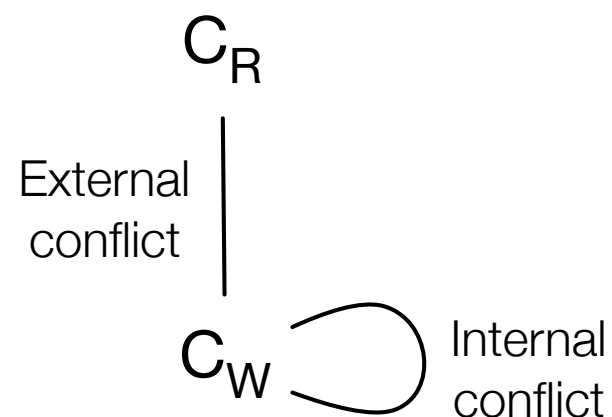
-  Classes of requests: expressing application concurrency
-  How to automatically map classes to worker threads
-  How the resulting technique compares to late scheduling

Classes of requests

Readers and writers

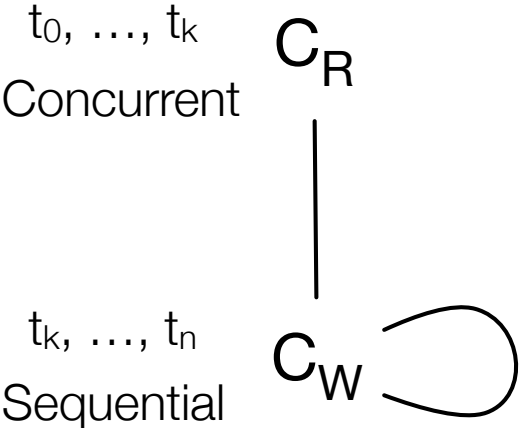
Class C_R : read requests

Class C_W : write requests

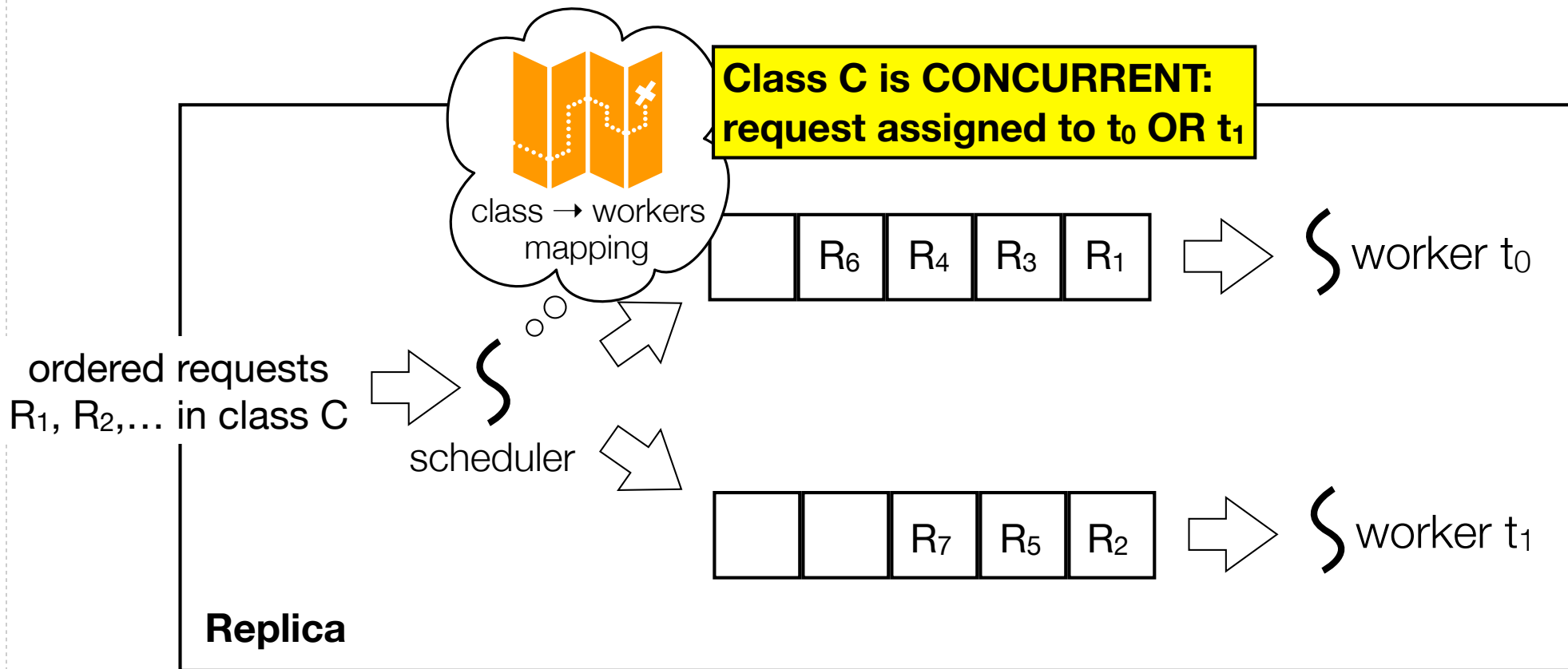


Mapping classes to workers

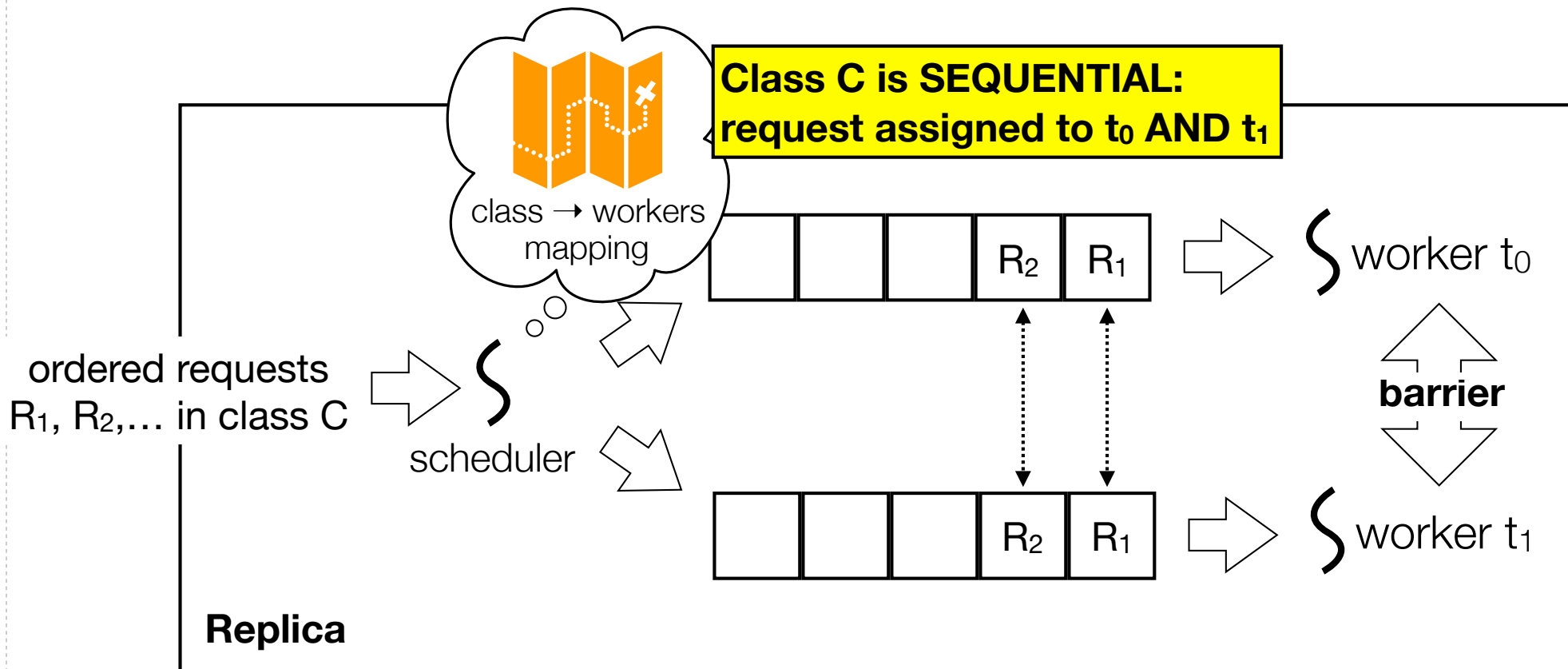
- ⦿ **Define workers that execute requests in the class**
- ⦿ **Define class type**
 - ⦿ **Sequential: one request at a time**
 - ⦿ **Concurrent: requests executed concurrently**



Early Scheduling execution model

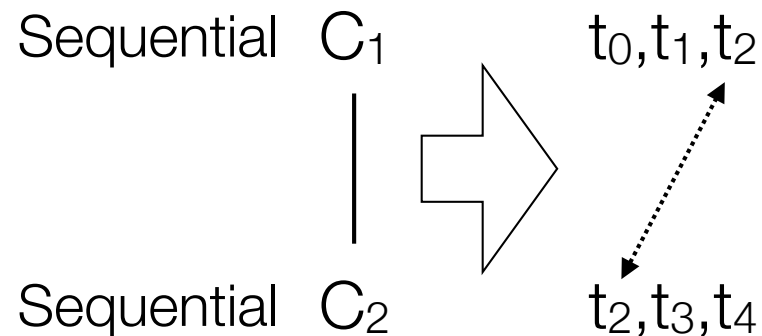


Early Scheduling execution model



Mapping classes to workers

Rule #5

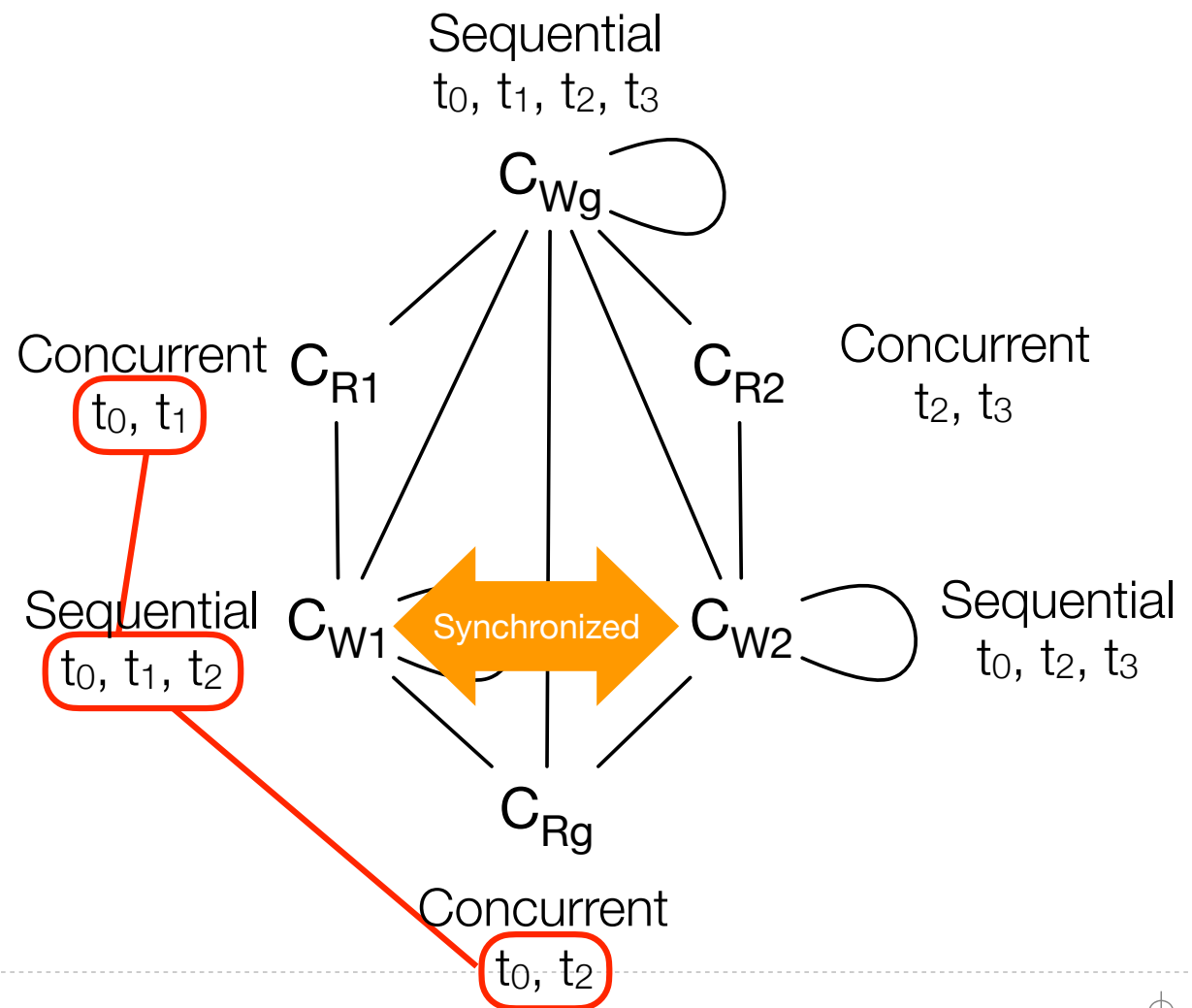


If C_1 and C_2 conflict, and are sequential, then C_1 and C_2 must have one worker in common

Mapping classes to workers

Local reads most common requests

Workers: t_0, t_1, t_2, t_3



Optimizing scheduling

- ❖ **O_{1a}: Minimize workers in sequential classes**
- ❖ **O_{1b}: Maximize workers in concurrent classes**
- ❖ **O₂: Assign workers to concurrent classes in proportion to class weight (i.e., more work, more workers)**
- ❖ **O₃: Minimize unnecessary synchronization among classes**

Optimization model

Algorithm 3 Optimization model.

⋮

15: **constraints:**

$$16: \quad \forall c \in C : \bigvee_{t \in T} \text{uses}(c, t) \quad // \text{ R.1}$$

$$17: \quad \forall c \in C : \#[c_1, c_1] \Rightarrow \text{Seq}[c_1] \quad // \text{ R.2}$$

$$18: \quad \forall c_1, c_2 \in C : \#[c_1, c_2] \Rightarrow \text{Seq}[c_1] \vee \text{Seq}[c_2] \quad // \text{ R.3}$$

$$19: \quad \forall c_1, c_2 \in C, t \in T : \#[c_1, c_2] \wedge \text{Seq}[c_1] \wedge \text{Cnc}[c_2] \wedge \text{uses}[c_2, t] \Rightarrow \text{uses}[c_1, t] // \text{ R.4}$$

$$20: \quad \forall c_1, c_2 \in C : \#[c_1, c_2] \wedge \text{Seq}[c_1] \wedge \text{Seq}[c_2] \Rightarrow \exists t \in T : \text{uses}[c_1, t] \wedge \text{uses}[c_2, t] // \text{ R.5}$$

21: **objective:**

22: minimize cost:

$$23: \quad + \sum_{t \in T, c \in C: \text{Seq}[c]} \text{uses}[c, t] \times w[c] / ws \quad // \text{ O.1a}$$

$$24: \quad - \sum_{t \in T, c \in C: \text{Cnc}[c]} \text{uses}[c, t] \times w[c] / wc \quad // \text{ O.1b}$$

$$25: \quad + \sum_{c \in C: \text{Cnc}[c]} |w[c]/wc - (|\{t \in T : \text{uses}[c, t]\}|/nt)| \quad // \text{ O.2}$$

$$26: \quad + \sum_{c_1, c_2 \in C: \text{Seq}[c_1] \wedge \text{Seq}[c_2] \wedge \neg \#[c_1, c_2]} |\{t \in T : \text{uses}[c_1, t] \wedge \text{uses}[c_2, t]\}| \times nt \times nc // \text{ O.3}$$

 Described in AMPL

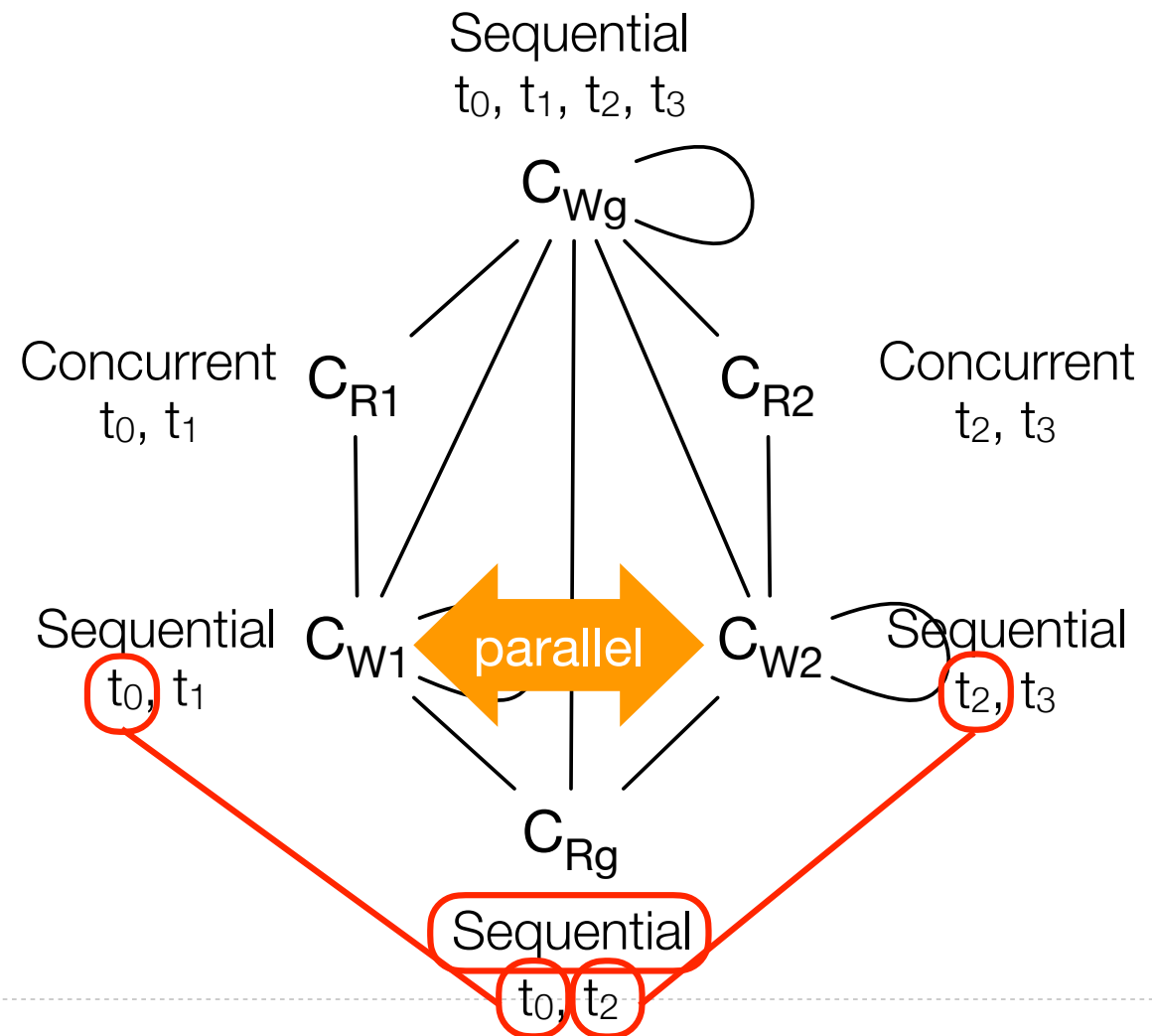
 Solved with KNitro



Naive vs Optimized mapping

Local reads most common requests

Workers: t_0, t_1, t_2, t_3

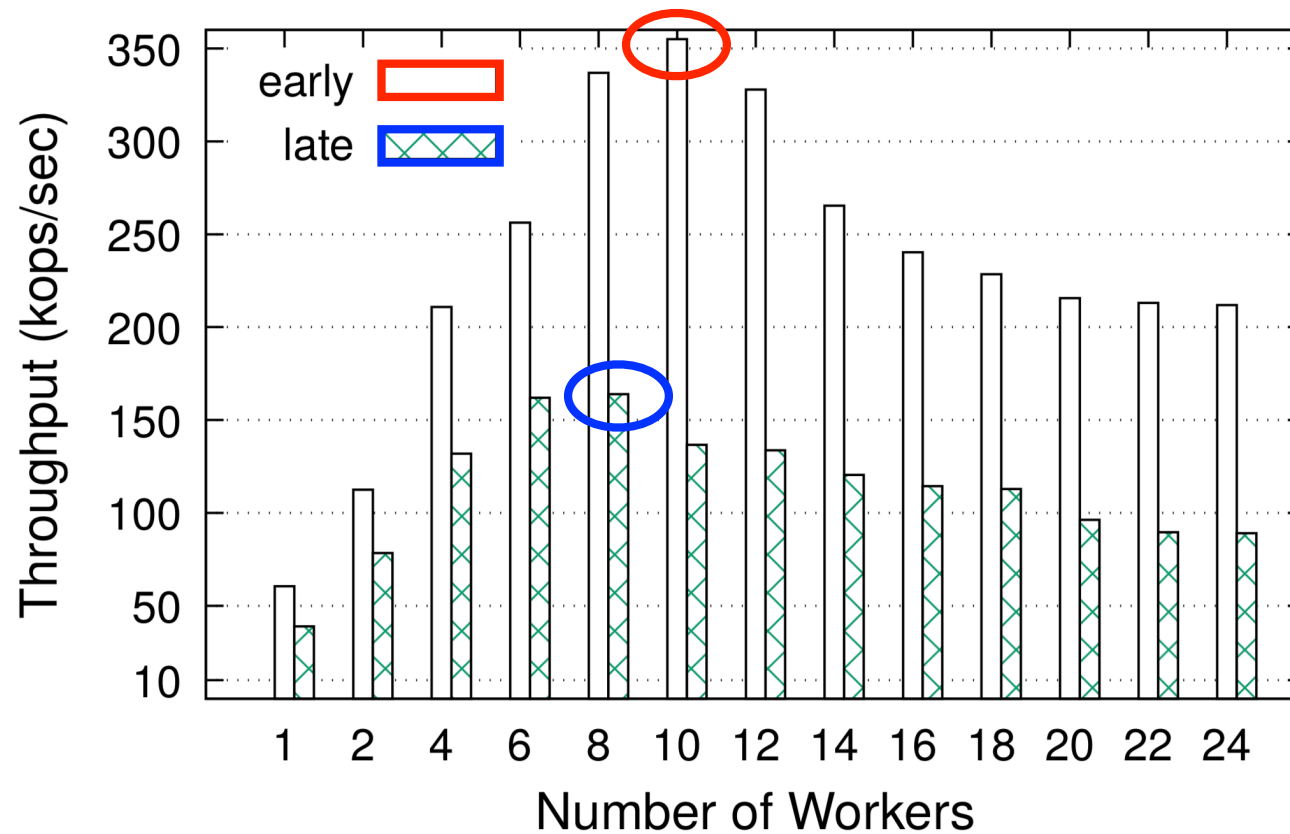


Experimental evaluation

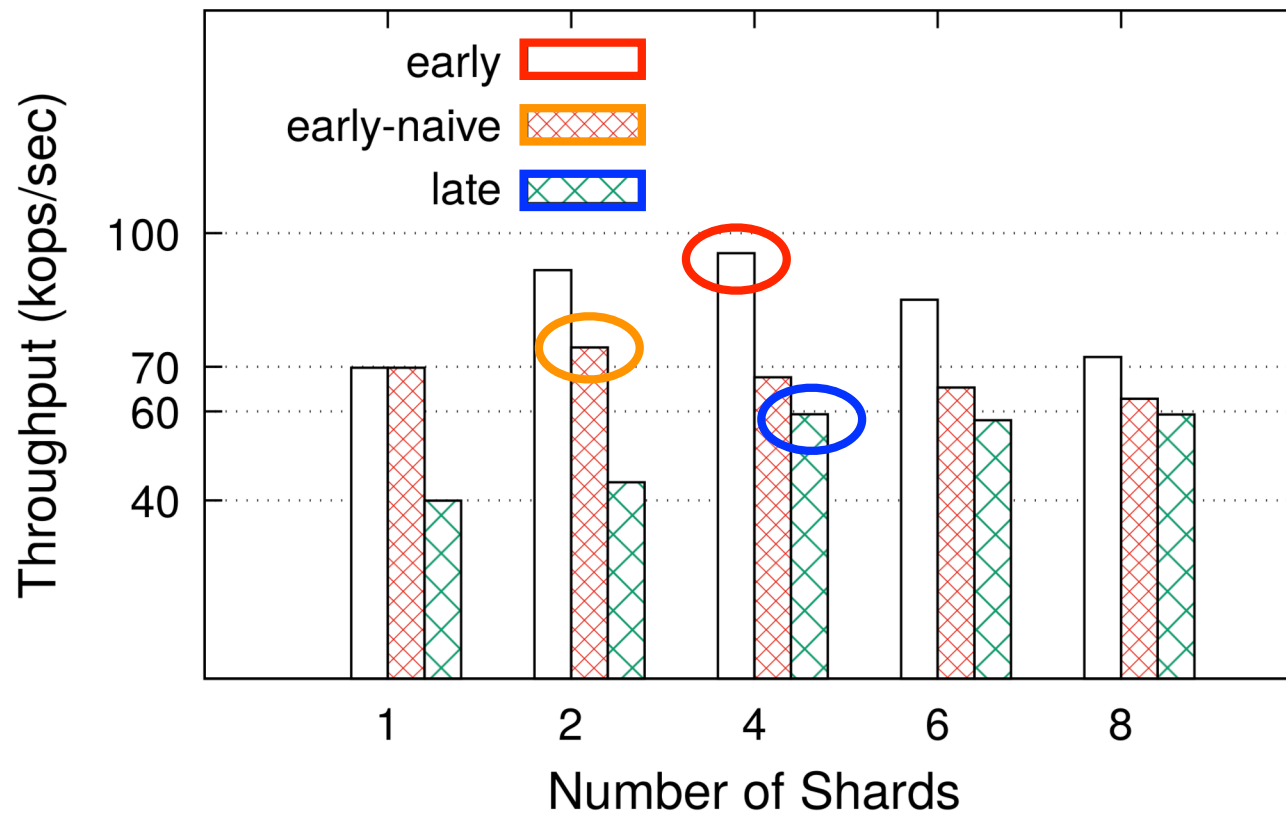
- **Prototype in BFT-SMaRt environment**
 - **Early scheduling and late scheduling**
 - **Configured to crash failures (not BFT)**
- **Linked-list application**
 - **Single- and multi-shard deployments**
 - **Light, moderate, and heavy execution costs**
 - **Uniform and skewed workloads**



Single-shard, reads, moderate



Multi-shard, mixed, moderate



[http://www.inf.usi.ch/faculty/
pedone/](http://www.inf.usi.ch/faculty/pedone/)