

All You Need is a Concurrent Data Structure

Hagit Attiya

Concurrent Data Structures

- Constructs for **efficiently storing and retrieving data**
- **Different types:** lists, hash tables, trees, queues, ... specified with an **interface** and expected behavior
- Can be put together to build other data structures

Concurrent Data Structures

- Constructs for **efficiently storing and retrieving data**
- **Different types**: lists, hash tables, trees, queues, ... specified with an **interface** and expected behavior
- Can be put together to build other data structures

Fuel many **multiprocessing software systems**

- Specifically, **multi-threaded** and **multi-core** environments or even **geo-replicated** systems

Jargon Alert

They say **key-value store**

We say **register(s)**



Jargon Alert

They say **index**

We say **search structure**



Jargon Alert

They say **concurrency package**

We say **synchronization primitives**



Me & My Research

Approaches

Fine-grained locking

No locking

Transactional memory

Wide-area replication

Related core theoretical questions in distributed computing

Methodologies

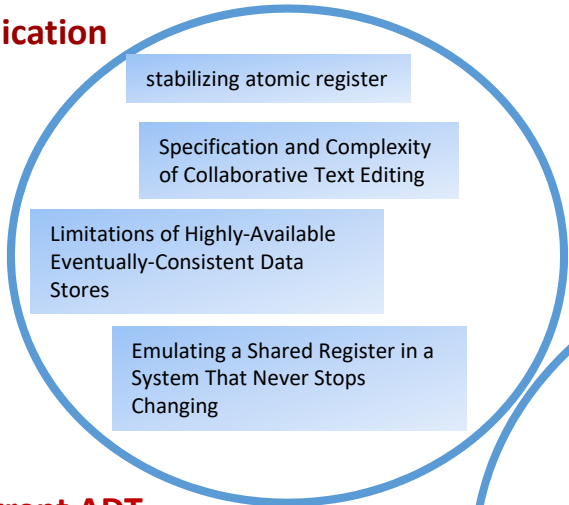
Theory (algorithms, lower bounds)

Practice

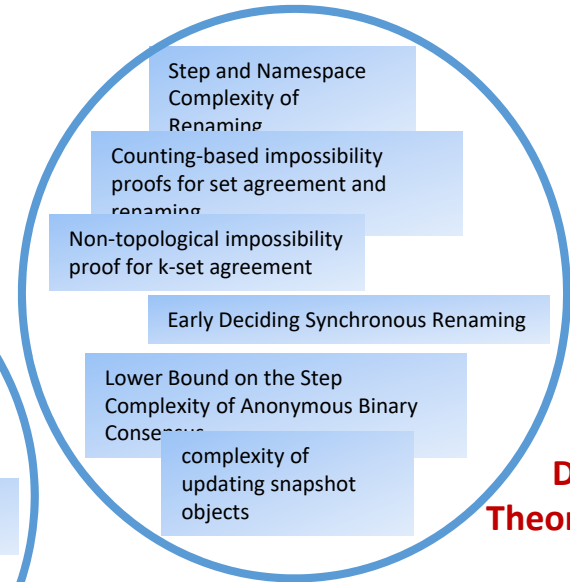
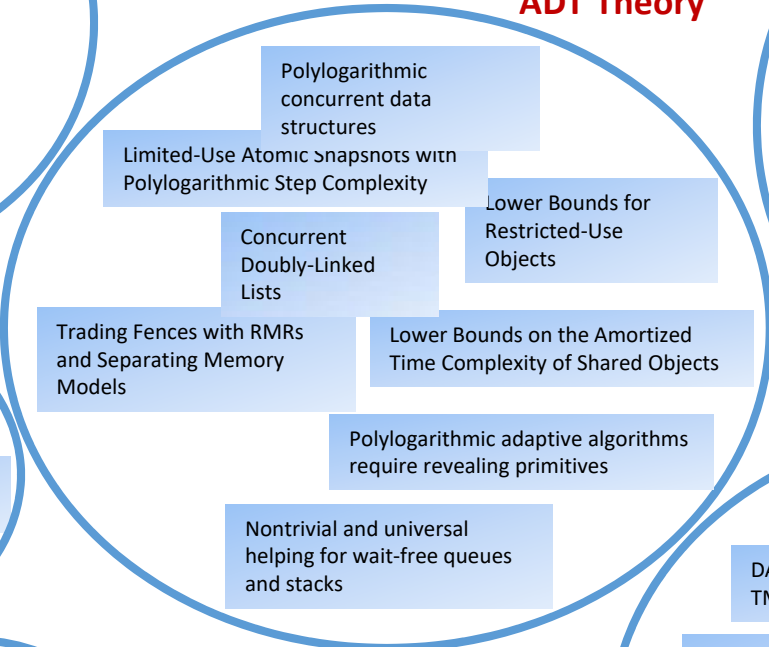
Formal methods (proof methods, specifications)

My Research

Replication

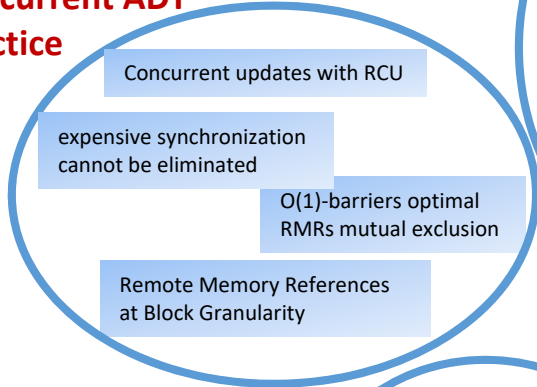


Concurrent ADT Theory

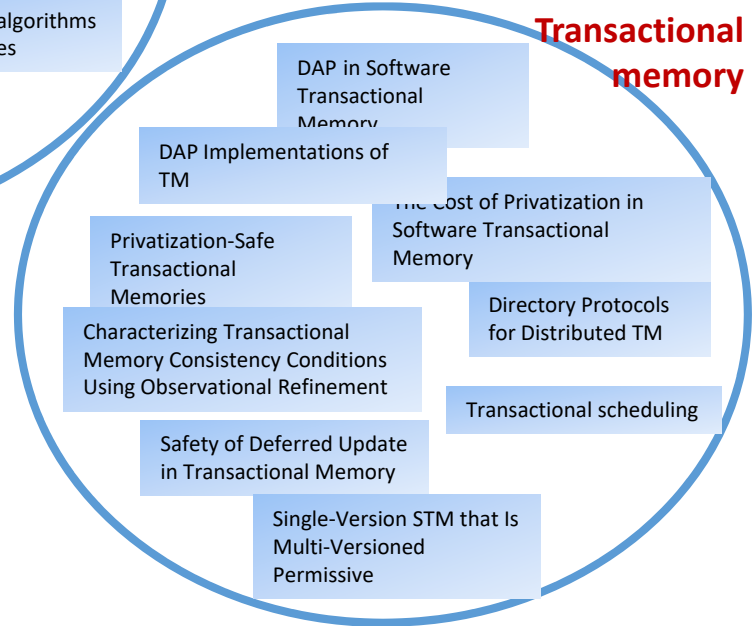


DC Theory

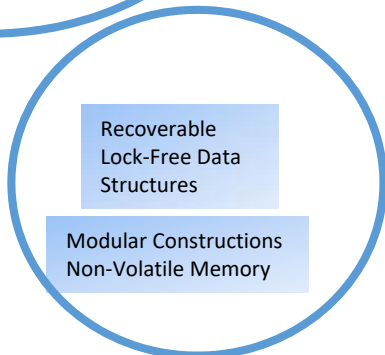
Concurrent ADT Practice



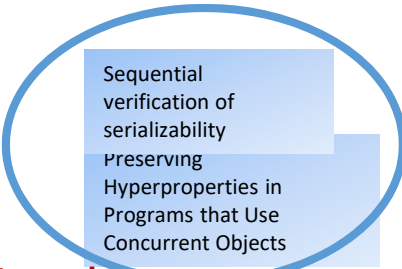
Transactional memory



NVRAM



Formal Methods



My Research

Replication

Emulating a Shared Register
in a System That Never Stops
Changing

Concurrent DS - Theory

Trading Fences with
RMRs and Separating
Memory Models

DC Theory

expensive
synchronization
cannot be eliminated

$O(1)$ -barriers
optimal RMRs
mutual exclusion

Concurrent DS - Practice

Recoverable
Lock-Free Data
Structures

Modular Constructions
Non-Volatile Memory

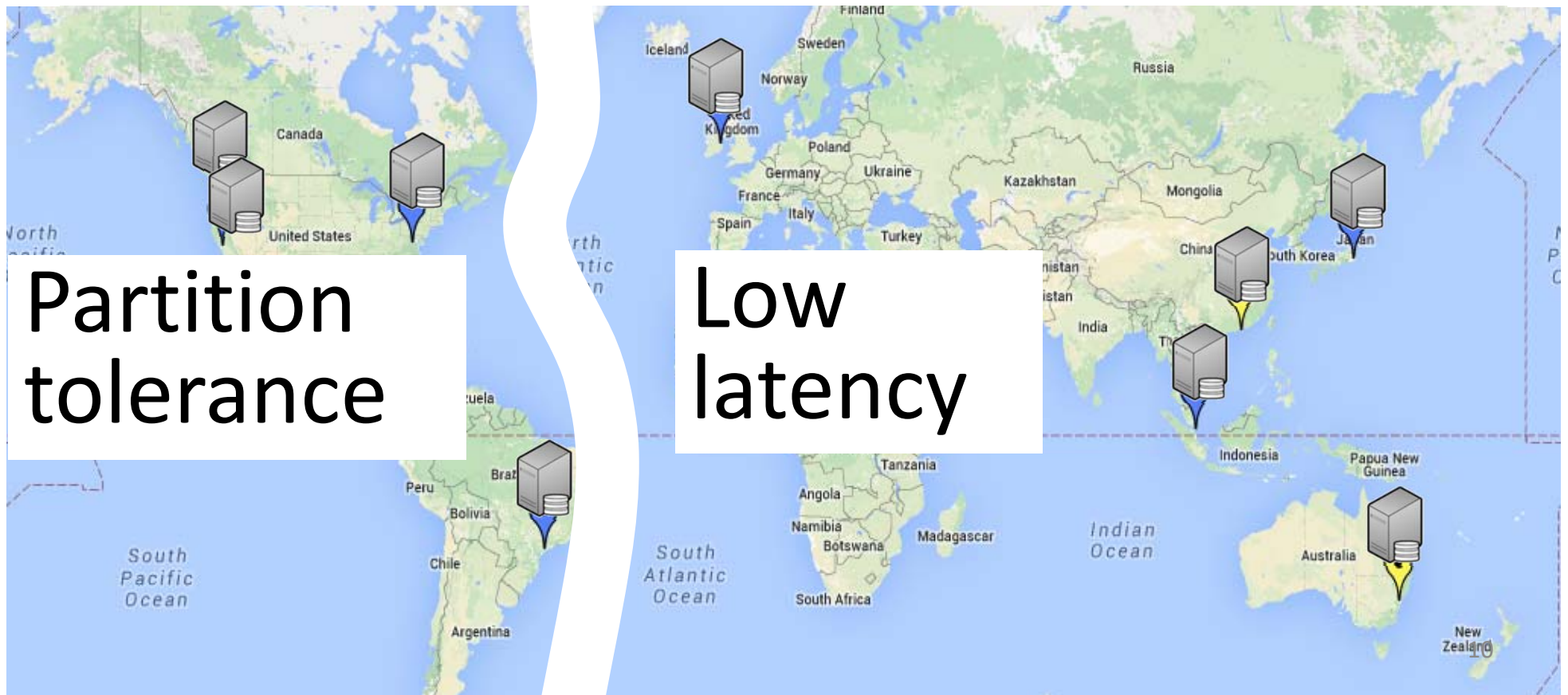
NVRAM

Formal Methods

Transactional memory

Key-Value Store

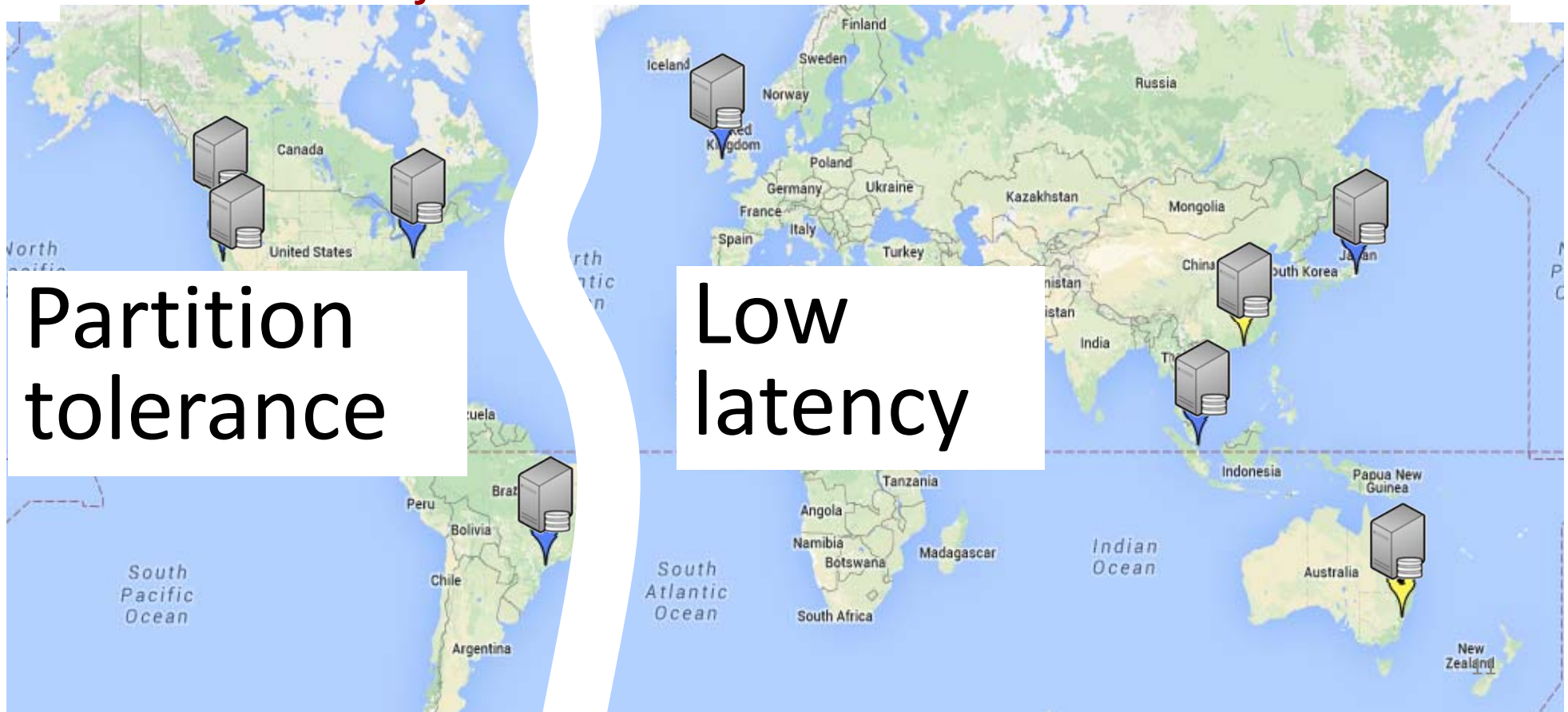
Geo-distributed systems powering Google, Facebook, Amazon, etc.



Key-Value Store: Our Approach

Simulate a register by keeping copies (replicas) of its value
Keep replicas consistent by exchanging messages

Churn: nodes join and leave at various times



Partition
tolerance

Low
latency

Key-Value Store with Constant Churn

CCReg: First shared register simulation that allows **churn to continue forever**, and system size to fluctuate

Ensures **reads and writes complete**, and new nodes can join and access the simulated register

Churn assumption: while a message is in transit, the number of nodes entering or leaving is
 $\leq \alpha \times$ number of nodes when the message was sent

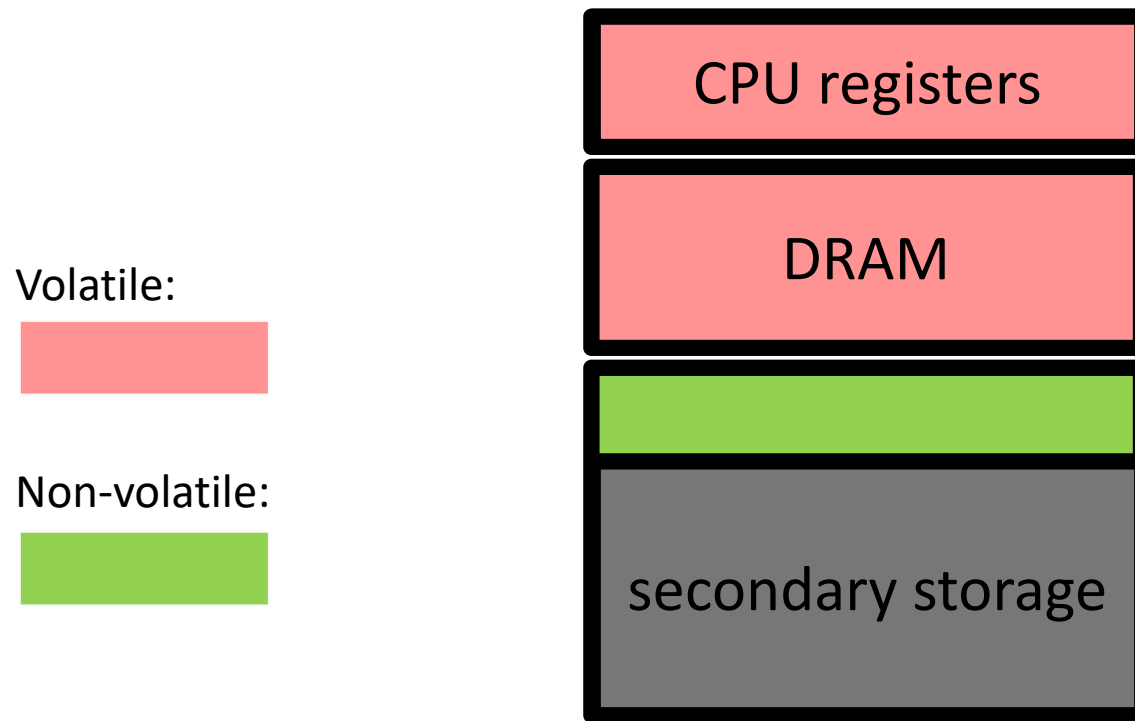
Key-Value Store with Constant Churn

Possible projects:

- Implement CCRag & its extension for Byzantine failures
- Simplify & improve bounds
- Ensure safety when churn assumption does not hold

Non-volatile RAM: Paradigm Shift

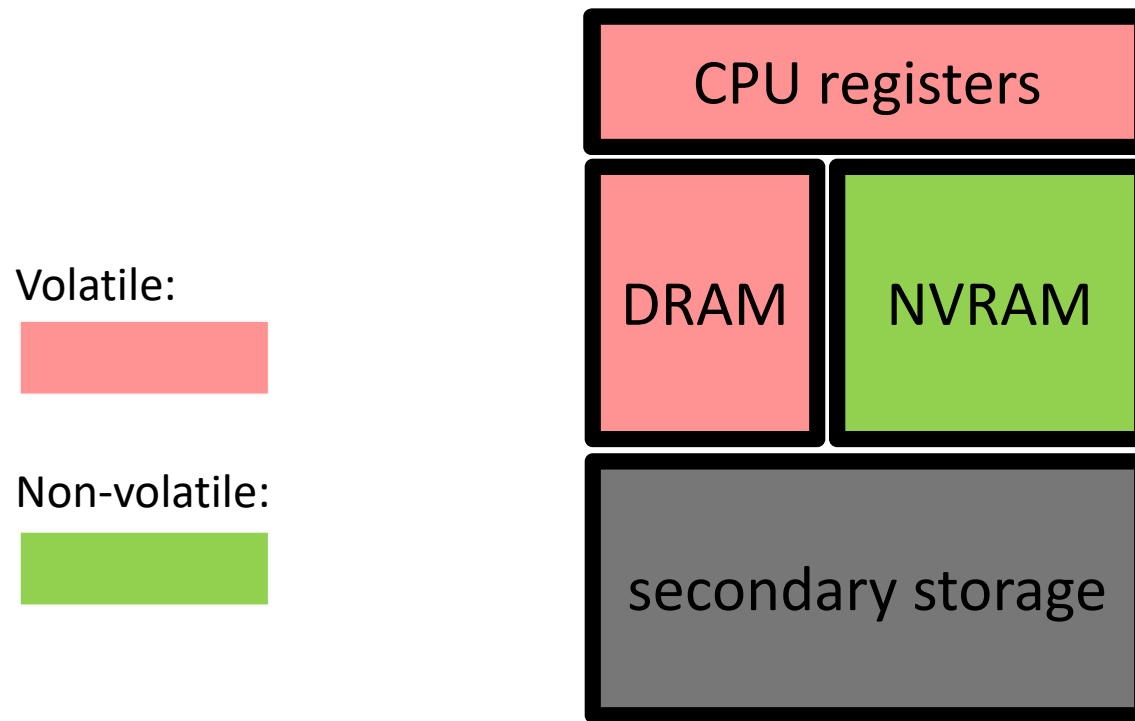
Discard and rebuild



(conventional)

Non-volatile RAM: Paradigm Shift

Recover and reuse

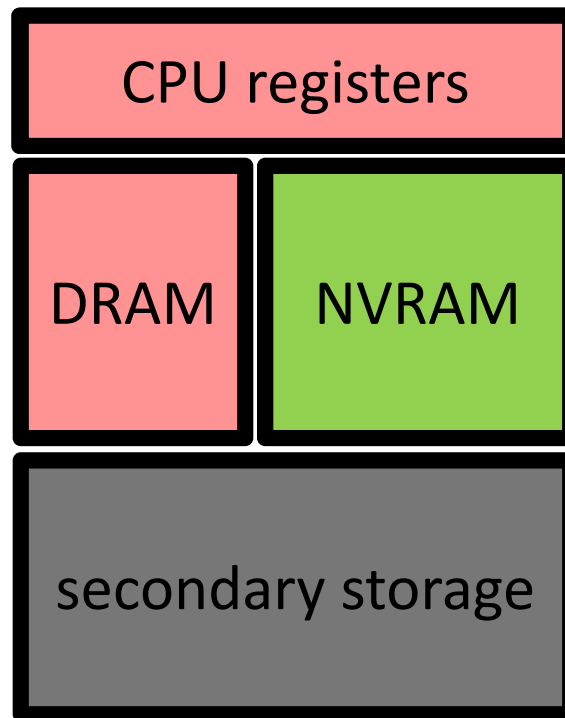


(future)

Non-volatile RAM: Paradigm Shift

We presented

- * New definitions
- * Simulations of **recoverable** read/write, compare-and-swap, test-and-set operations



Possible projects

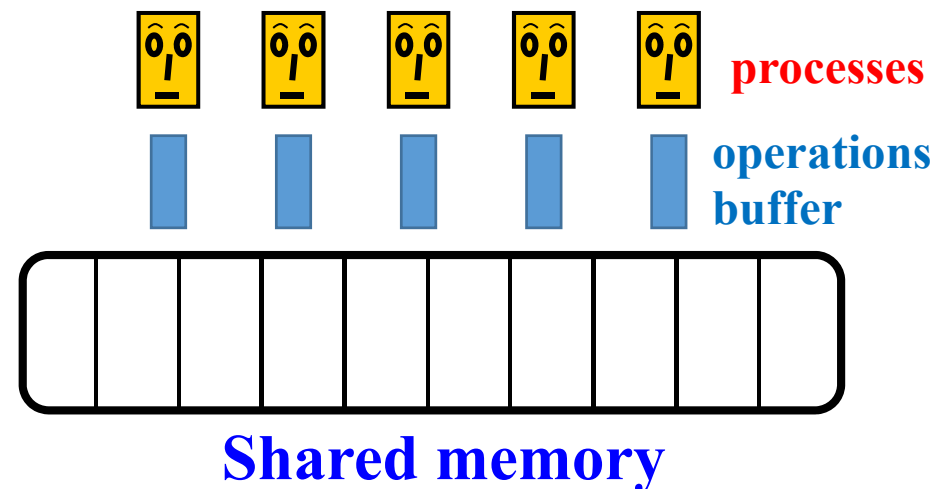
- * Persistence ordering and consistency ordering
- * System support
- * Differences between AMD & Intel require to abstract the architecture

Out-of-Order Execution

Architectures compensate for slow writes by **allowing reads to bypass earlier writes** that access a different location (TSO) \Rightarrow harms mutual exclusion algorithms

Avoid reordering with slow **fences**

We proved one fence is needed



Out-of-Order Execution & Caching

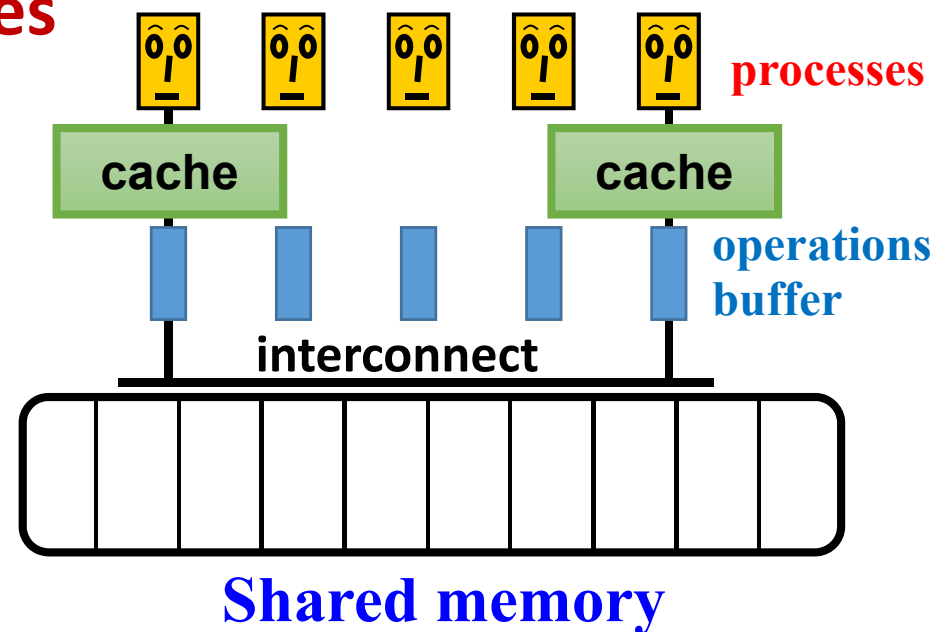
Architectures compensate for slow writes by **allowing reads to bypass earlier writes** that access a different location (TSO) \Rightarrow harms mutual exclusion algorithms

Avoid reordering with slow **fences**

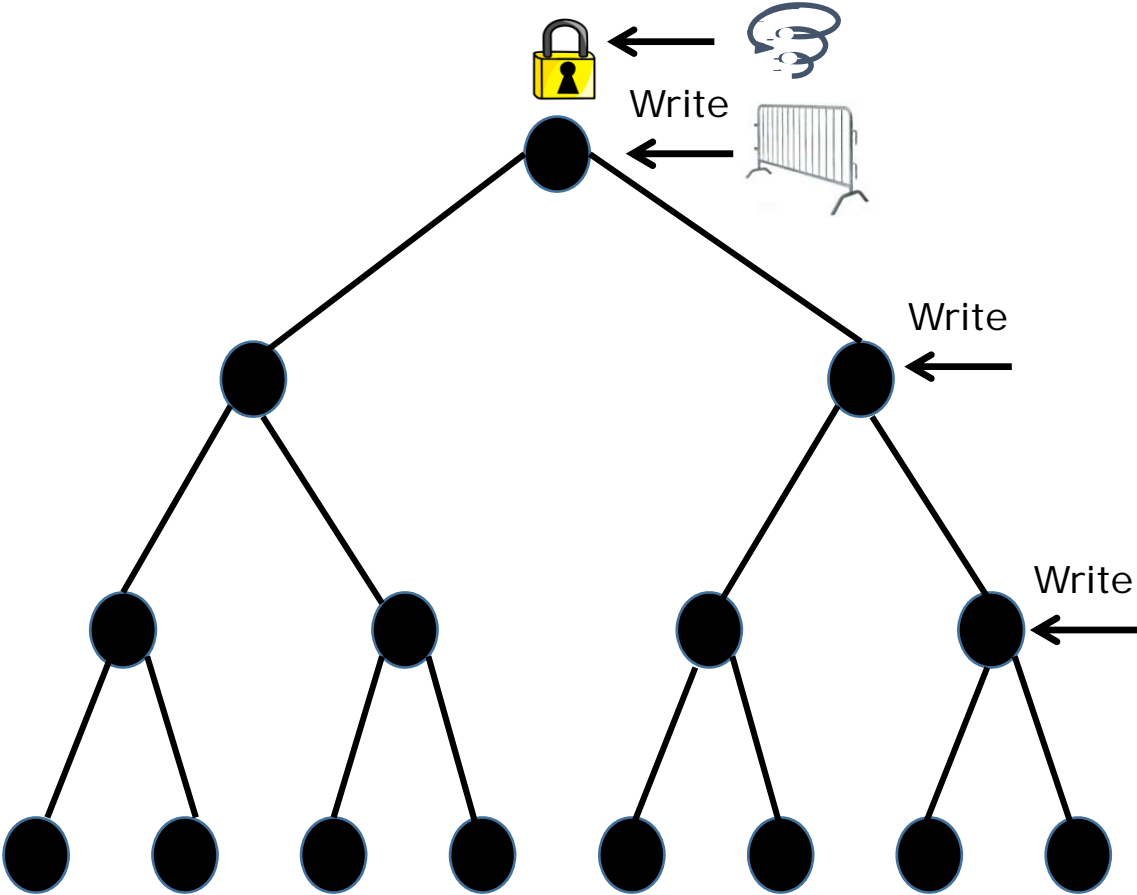
We proved one fence is needed

Reads from the cache are cheap

Remote reads are expensive



One-Fence Mutex: Entry Section



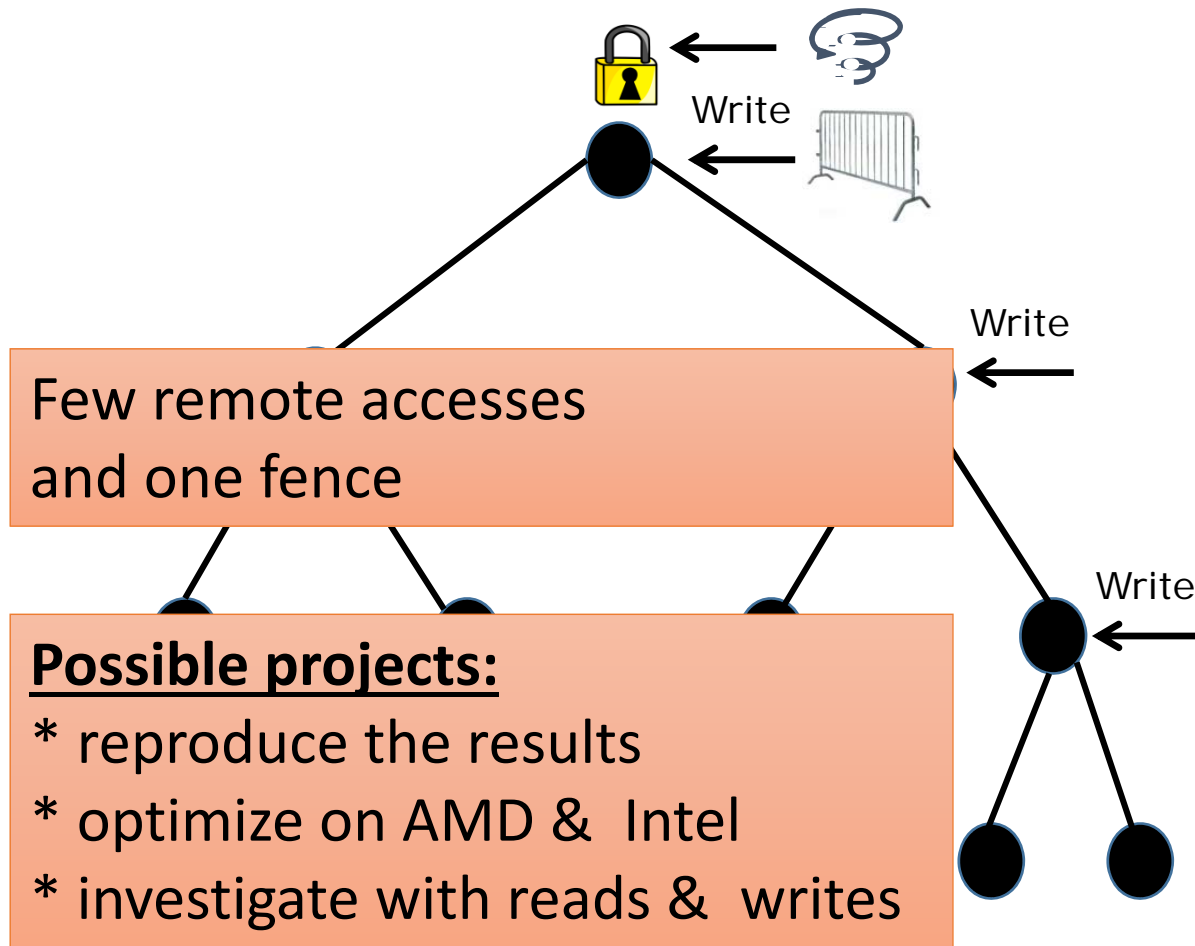
Promotion Queue



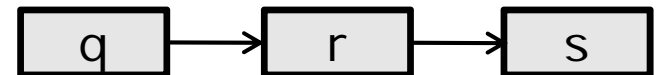
When exiting the critical section, processes promote waiting processes into a queue of waiting processes

Ensure that waiting processes are promoted, and hence, not starved

One-Fence Mutex: Entry Section



Promotion Queue



When exiting the critical section, processes promote waiting processes into a queue of waiting processes

Ensure that waiting processes are promoted, and hence, not starved