# GLOBE
## Global Object Exchange

A dynamically fault-tolerant and dynamically scalable distributed tuplespace for heterogenous, loosely coupled networks

Jesper Honig Spring

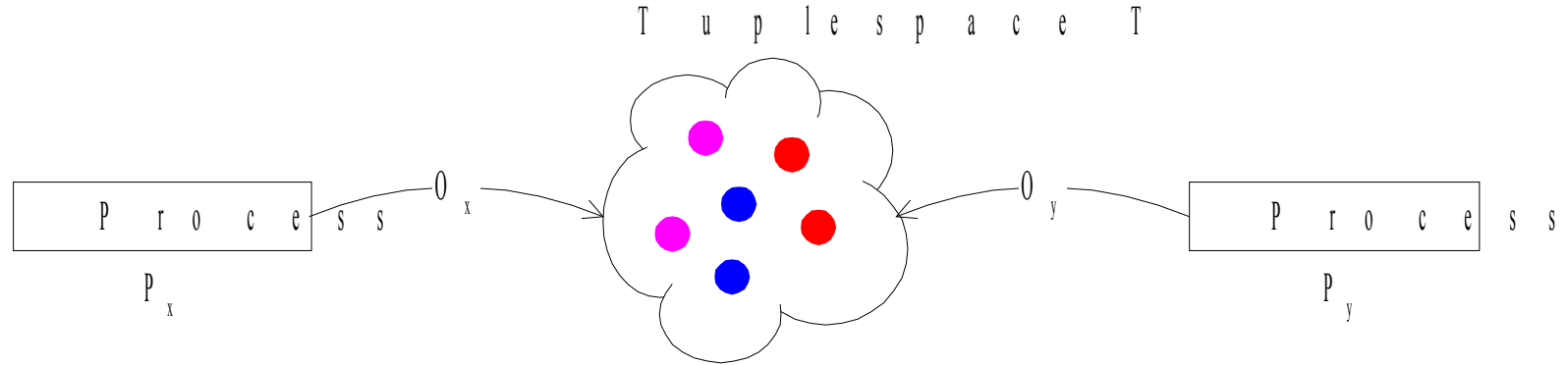http://www.diku.dk/students/eglarsen/GLOBE

# Agenda

- The Tuplespace Paradigm

- Tuplespace Semantics

- Achieving Fault-tolerance

- Achieving Scalability

- Measurements of GLOBE

- Related projects

- Conclusion

- Demonstration and Questions

# The Tuplespace Paradigm

- A communication paradigm (from '85)
- Communication unit:
  Tuple ($\approx$ Object)
- Tuplespace – an
  intermediate container
- Tuples immutable in
  the tuplespace
- Implementations:
  Linda, JavaSpaces, TSpaces etc.

```
Class Person extends Tuple {
  field String name;
  field Integer age;

  method String getName()
  method void setName(newName)

  method Integer getAge()
  …
}
```

# The Tuplespace Paradigm

T u p l e s p a c e  T

P r o c e s s $O_x$ ... $O_y$ P r o c e s s

$P_x$ $P_y$

- Atomic Operations:
  - Insertion (out)
  - Withdrawal (in, inp∗)
  - Inspection (rd, rdp∗)
  - Additional operations

∗rdp/inp – predecate

- Matching:
  - Templates (anti-tuple)
  - Exact/Wildcard matching:
    - Tuple Type (null tuple)
    - Tuple Field Values (wildcard fields)

# The Tuplespace Paradigm

- Groupware
  - Chat server
  - Shared Blackboard

- High Performance Computing
  - SETI@Home-like calculations

- Intelligent Connectionware
  - Internet Services (internally in Jini LUS)
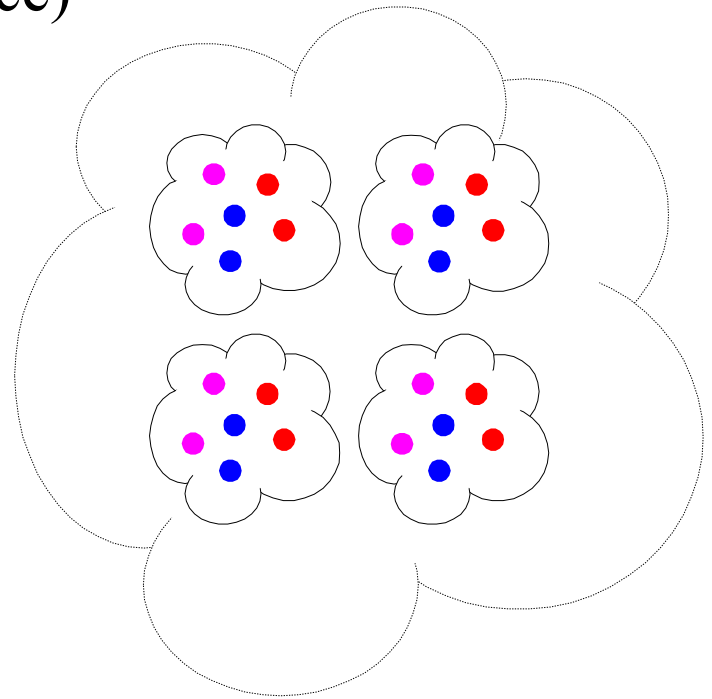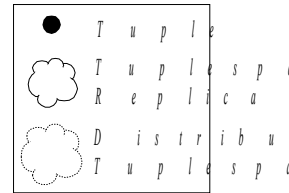  - Intelligent Home (TSpaces at IBM)

# The Tuplespace Paradigm

The Classical Problems with Centralized Systems:

- Availability (Level of Fault-Tolerance)
  - Single Point of Failure

- Scalability
  - Cannot Scale beyond the Single Entity

# The Tuplespace Paradigm

- Purpose of GLOBE:
  - Increase Availability (Fault-tolerance)
  - Increase Scalability
  - Dynamic Adjustment

- Distributed Tuplespace Abstraction

# Operations Semantics

- Tuplespace Semantics vaguely defined:
  - Selection of matching tuple
    - Arbitrary, FIFO, LIFO etc.
  - Selection of process to withdraw tuple
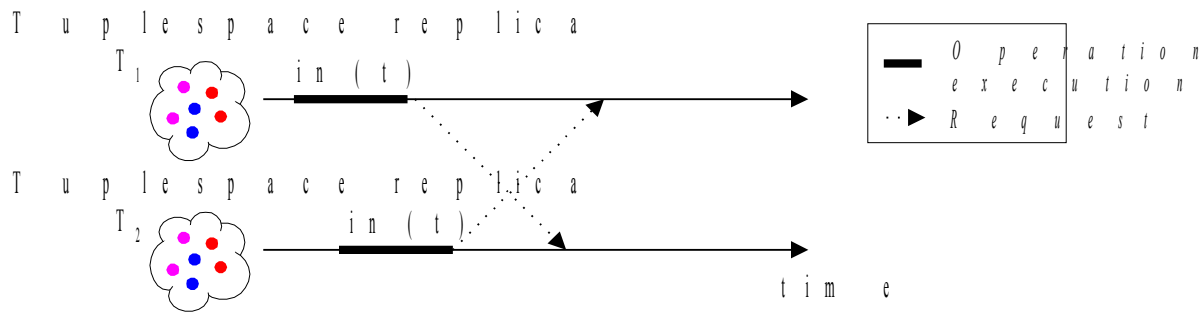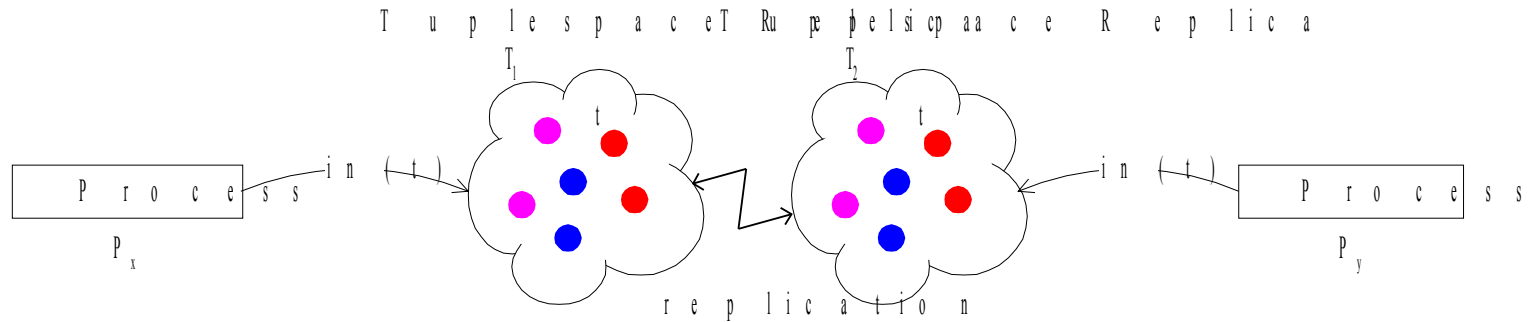    - Concurrent withdrawals for same tuple
    - Specified as "fair"

In addition: Issue related to distribution:
  - Predicate operations semantically unclear
    - Inspect "present" state
    - What is present state in a distributed environment?

# Operations Semantics
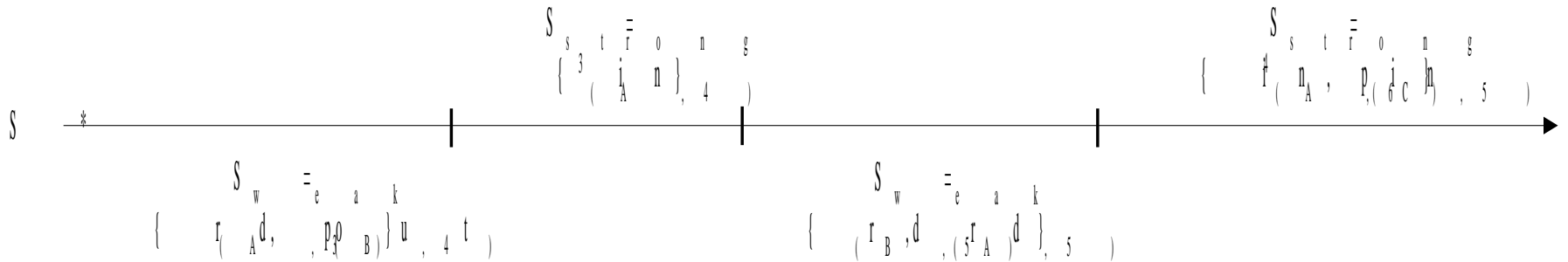## Concurrent and Distributed Tuple Withdrawal



- Withdrawal operations must be performed atomically!
- Global ordering across replicas

# Semantics of GLOBE

- Two Categories of Tuplespace Operations:
  – Strong Operations (in, inp)
  – Weak Operations (out, rd, rdp)

- GLOBE adapts *loose inp/rdp* semantics
  – Weak Operations are performed "locally" and any changes (insertions) propagated later (depending on the synchronization tightness).
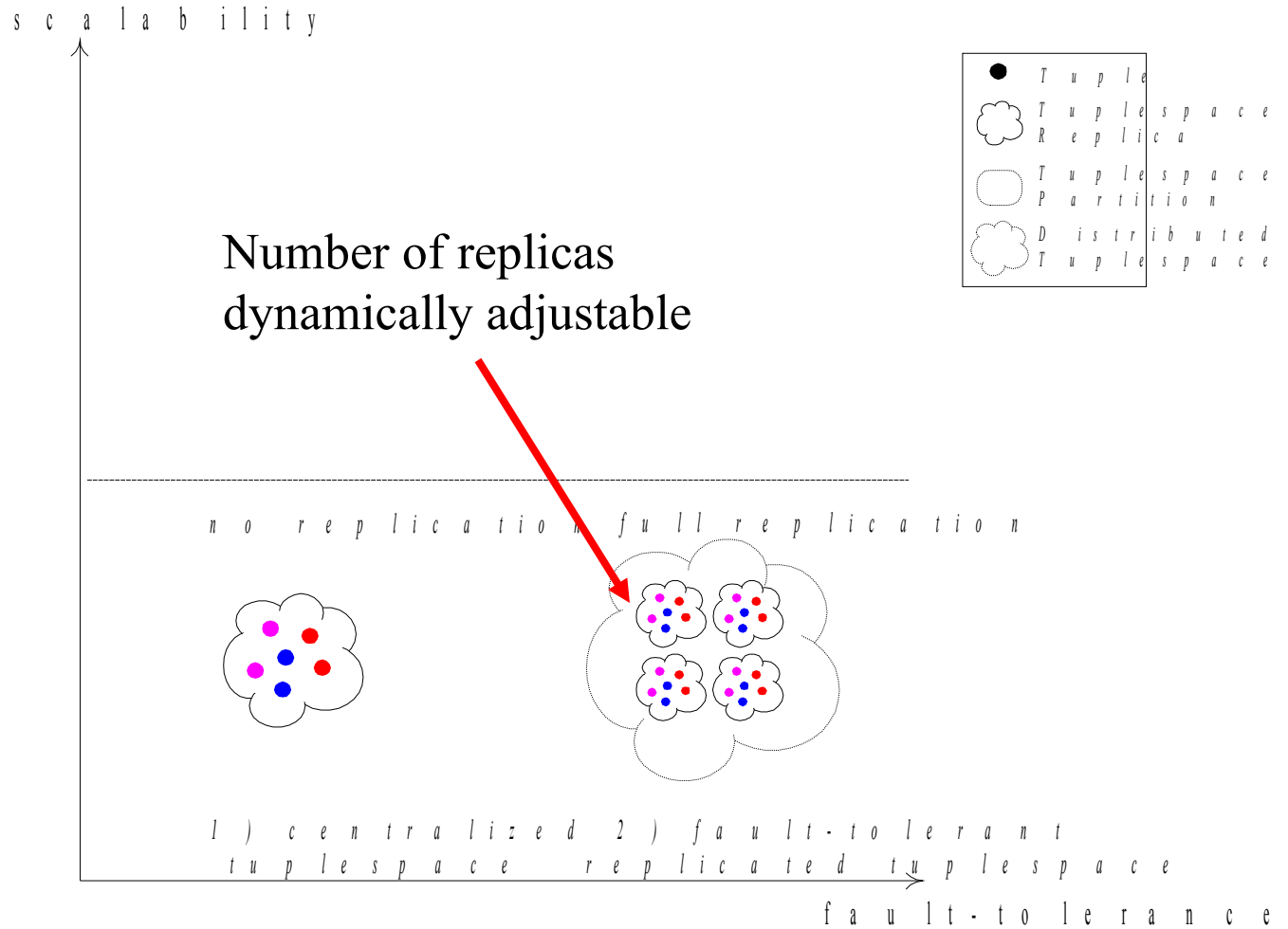  – inp/rdp may show "false" results!

# Semantics of GLOBE

- Operation Ordering:
  - Strong Operations are Globally Ordered
  - Weak Operations are Globally Unordered
  - All Operations satisfy Partial Ordering

S* -- a sequence of tuplespace operations performed on a replica

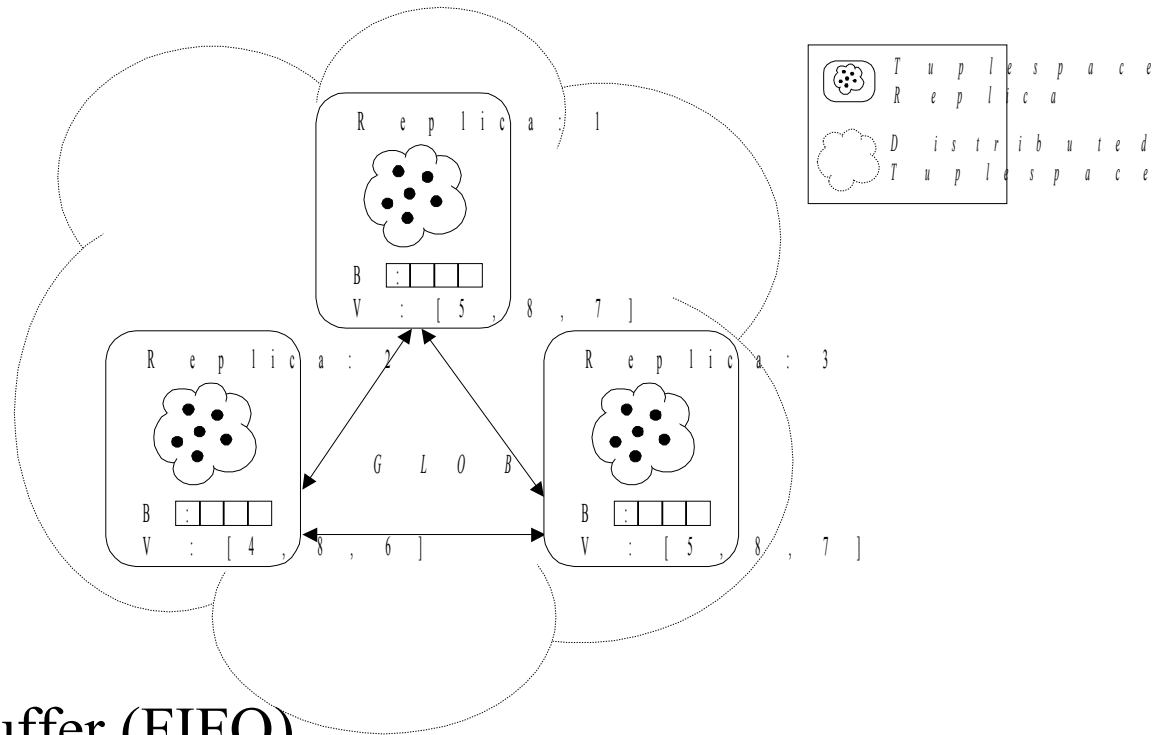# Achieving Availability (Fault-tolerance) and limited Scalability by Replication

# Replica Update Protocol

- Active propagation
  - Operations are persisted before propagation
  - Replicas responsible for propagation
  - Fast propagation
  - Majority voting for all atomic operations (tuplespace operations and configuration operations)
- Problem: Not 100% reliable!
  - Adjustability problem (removal of replica)
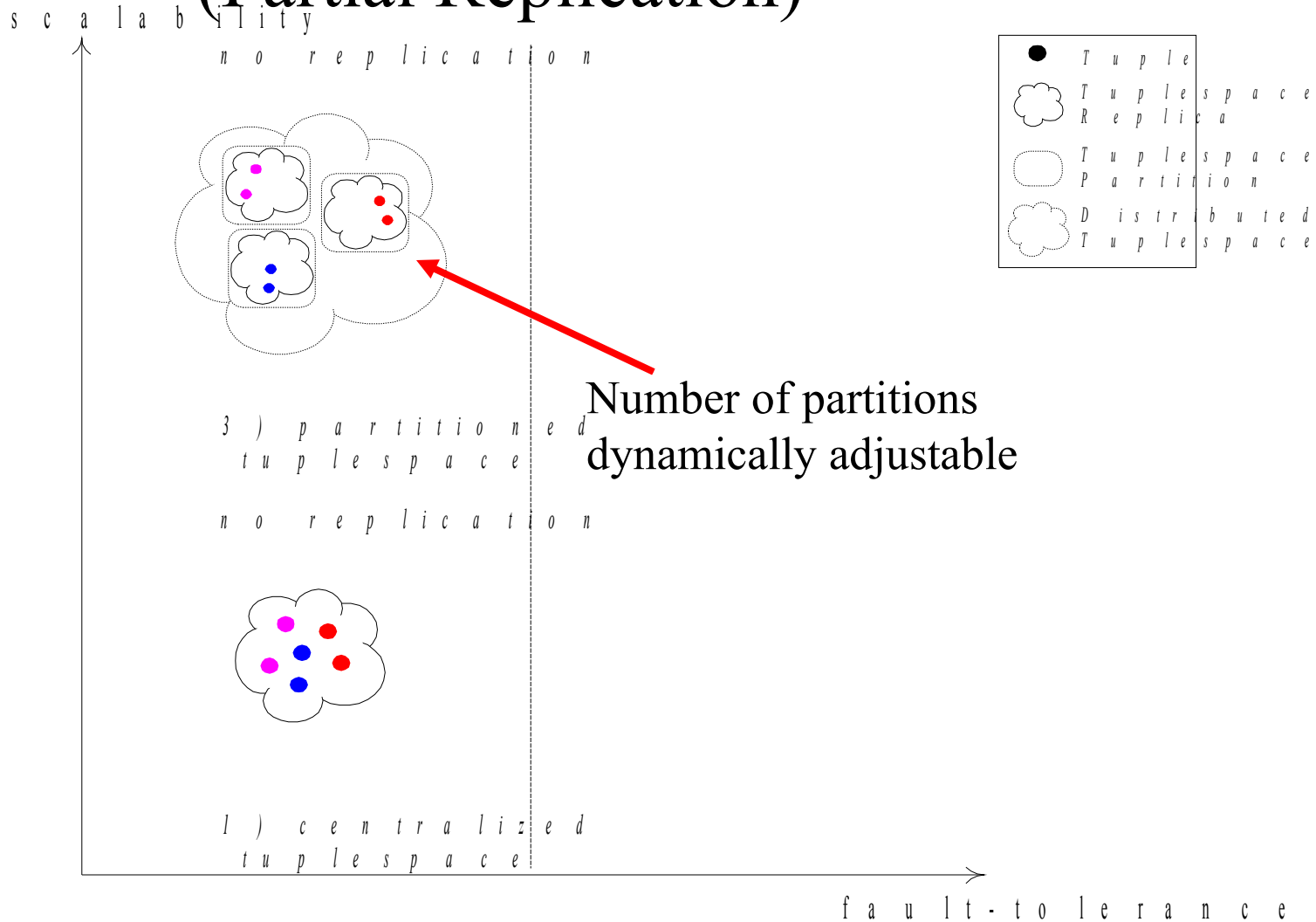  - Inconsistency in case of failure

# Replica Update Protocol

- Anti-entropy
  - Epidemic algorithm (Bayou), pair-wise reconciliation - slow!
  - Replicas responsible for updating themselves (pull-based to avoid duplicates)
  - Synchronization in case of failure
- Hybrid replica update protocol
  - Combines Active Propagation and Anti-Entropy
  - Ensures consistency convergence

# Update Propagation



- Operation Buffer (FIFO)
- Operation Vector
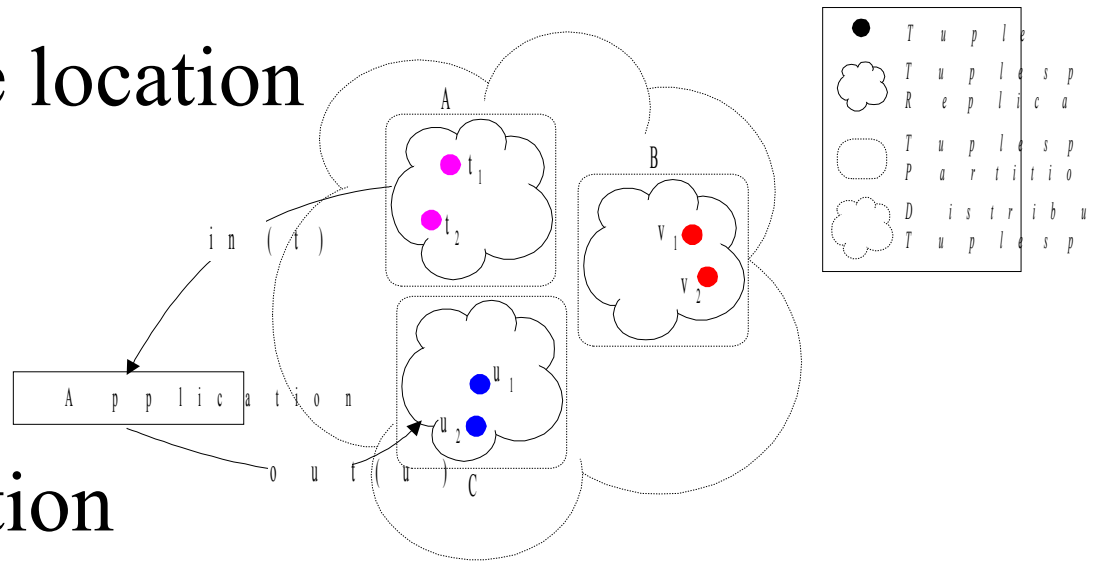- Logical Operation Numbers

# Achieving Scalability by Partitioning (Partial Replication)
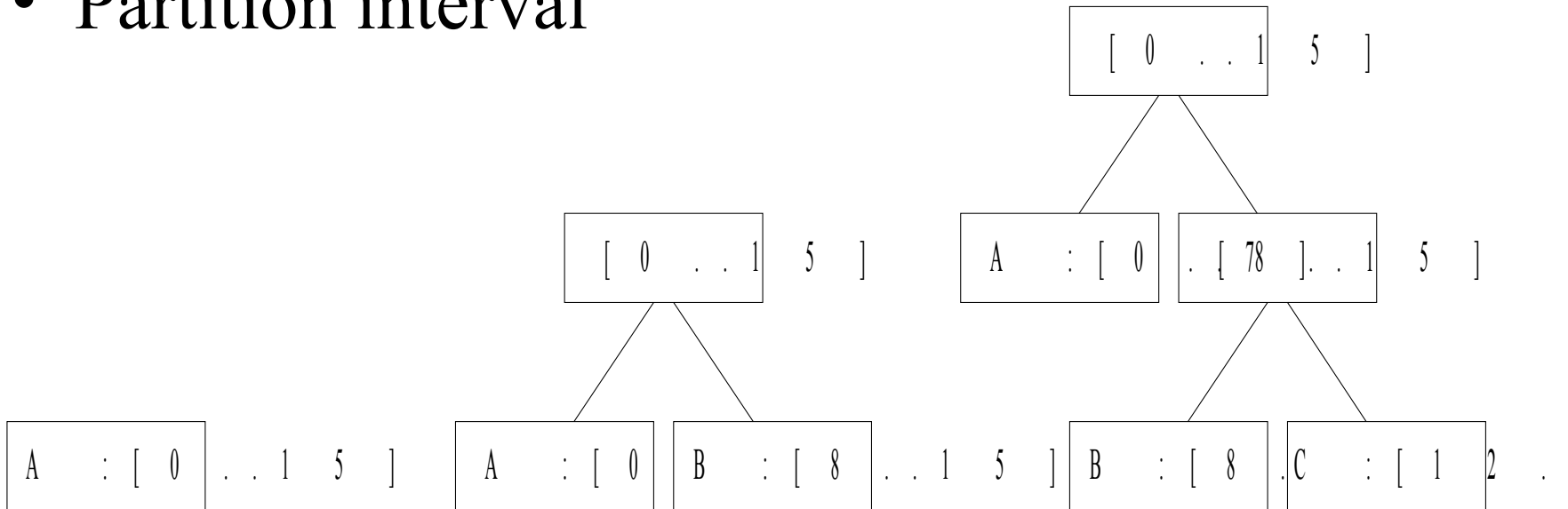


Number of partitions dynamically adjustable

# Tuplespace Partitioning

- Load balancing
- Resolution of tuple location
  - Non-deterministic
  - Deterministic
- Hash code
- Operation Redirection
- Problem: Dynamic adjustment
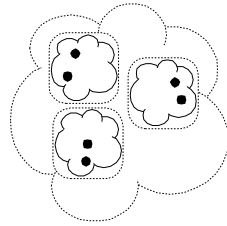  - Complete rehashing of tuples

# Dynamic Partitioning
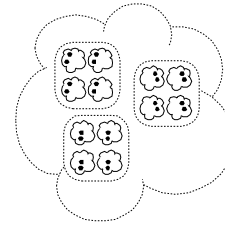
- Partitioning by hash code
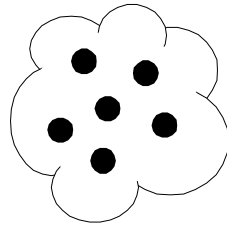- Partition interval

# Related Projects

Scalability

[Bjo93]     GLOBE (loosely coupled network)
            [Kri91] (hardware-based)

            4) Partitioned and
   3) Partitioned fault-tolerant tuplespace

Original Linda [CimG86] Self-Nettilo insda
[Sun99b] Java Spaces [BSa95] FT-Linda
[W+98b] TSpace [sXL89]
[Kam91]
[CKM92] MTS

   1) Centralized 2) Fault-tolerant tuplespace

Fault-tolerance

● Tuple
Tuplespace Replica
Tuplespace Partition
Distributed Tuplespace

# Cost of Fault-Tolerance
## (Empirical)

- LAN-based measurements (16Mbps TR)
- 8 Pentium PCs

Milliseconds (y-axis), Tuplespace Replicas (x-axis)

inp
in
rd/rdp
out

# Gains of Fault-Tolerance
## (Theoretical)

out/rd/rdp

in/inp

Level of Fault-Tolerance

Tuplespace Replicas

• Assumes 5% failure probability
• Binomial distribution
• Dependent on the semantics

# Cost of Scalability

## (Empirical)



- LAN-based measurements (16Mbps TR)
- 8 Pentium PCs

# Gains of Scalability
## (Empirical)



- LAN-based measurements (16Mbps TR)
- 8 Pentium PCs

20 Applications

15 Applications

10 Applications

5 Applications

1 Application

Milliseconds

Tuplespace Partitions

# Conclusion

- Distributed Tuplespace Semantics is suitable

- Hybrid Replica Update Protocol is fast and ensures consistency

- Higher level of Fault-Tolerance

- Higher level of Scalability

- Fault-Tolerance and Scalability dynamically Adjustable

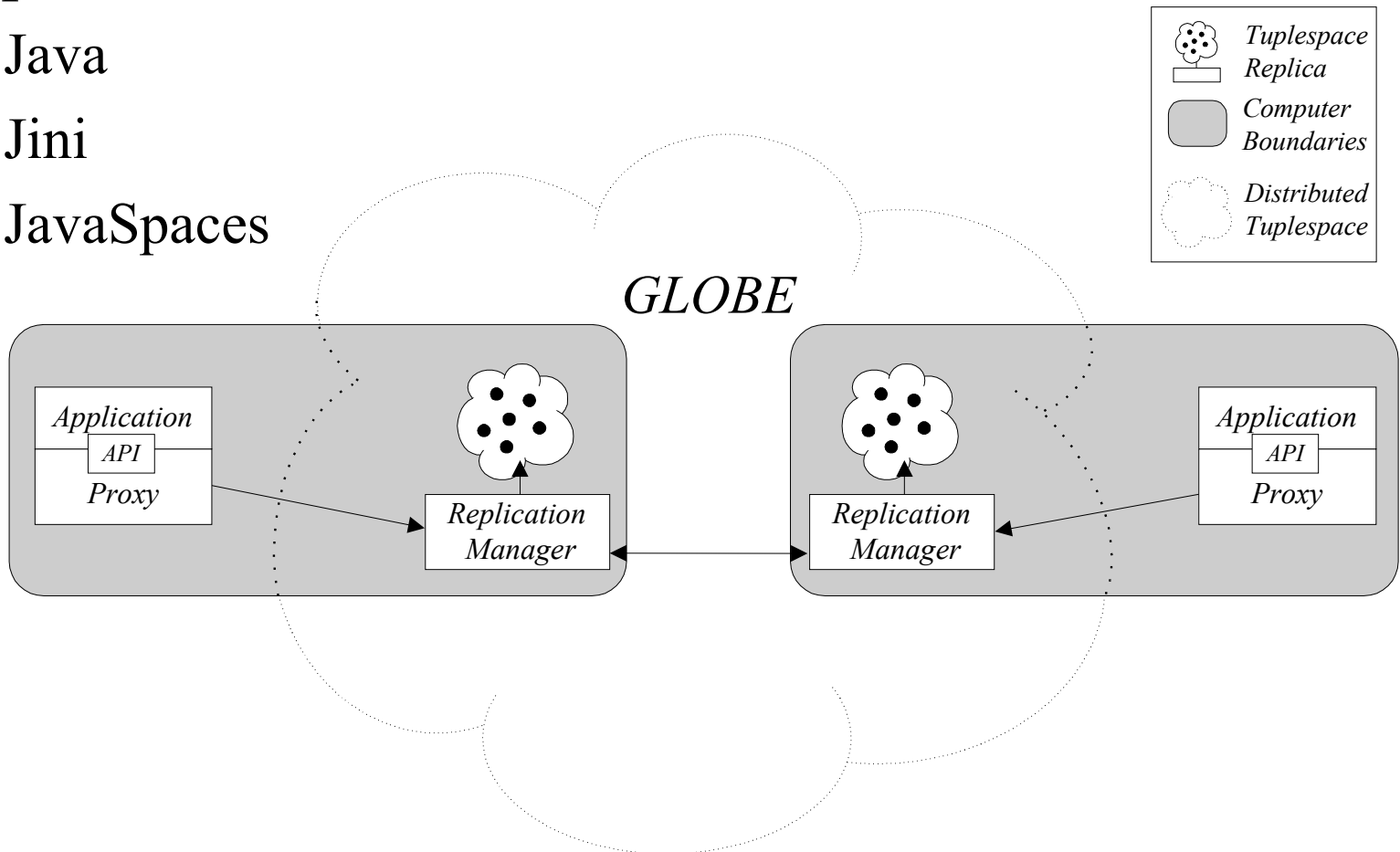- Outperforms a highly loaded centralized Tuplespace

# Future Work

- GLOBE enhancements
  - Elimination/reduction of Redirection
  - Load-balancing of Applications
  - Implementation Optimizations
  - Additional Tuplespace Features

# Fault-Tolerance & Scalability