

# Replex: A Multi-Index, Highly-Available Data Store

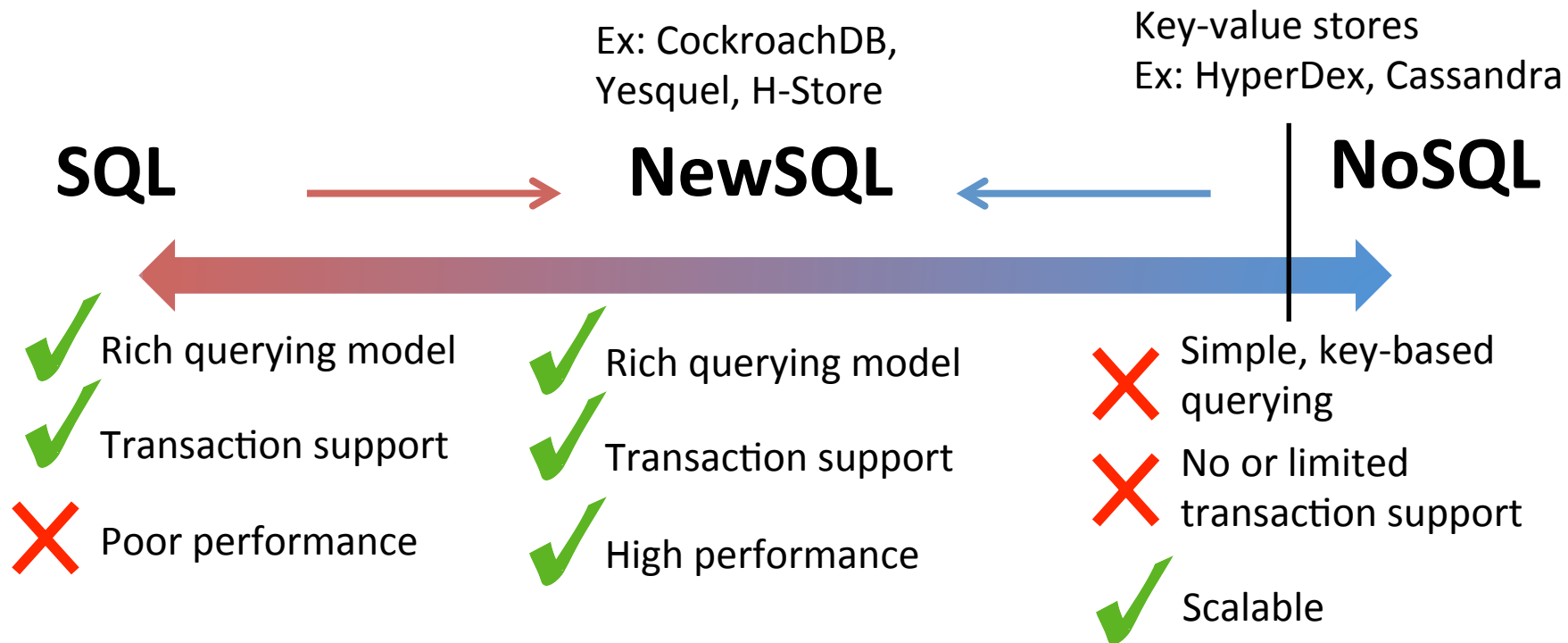
Amy Tai\*<sup>§</sup>, Michael Wei\*, Michael J. Freedman<sup>§</sup>,  
Ittai Abraham\*, Dahlia Malkhi\*

\*VMware Research, <sup>§</sup>Princeton University

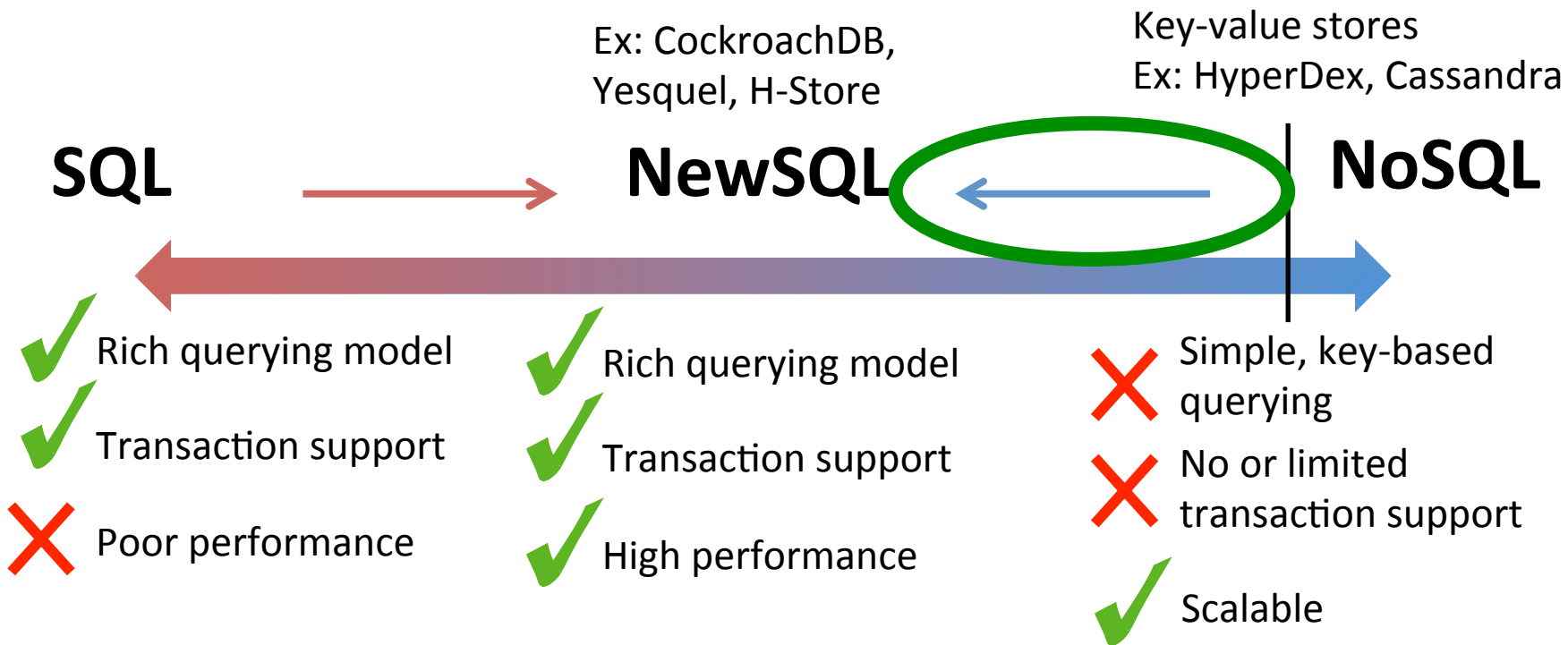
vmware<sup>®</sup>



# SQL vs NoSQL

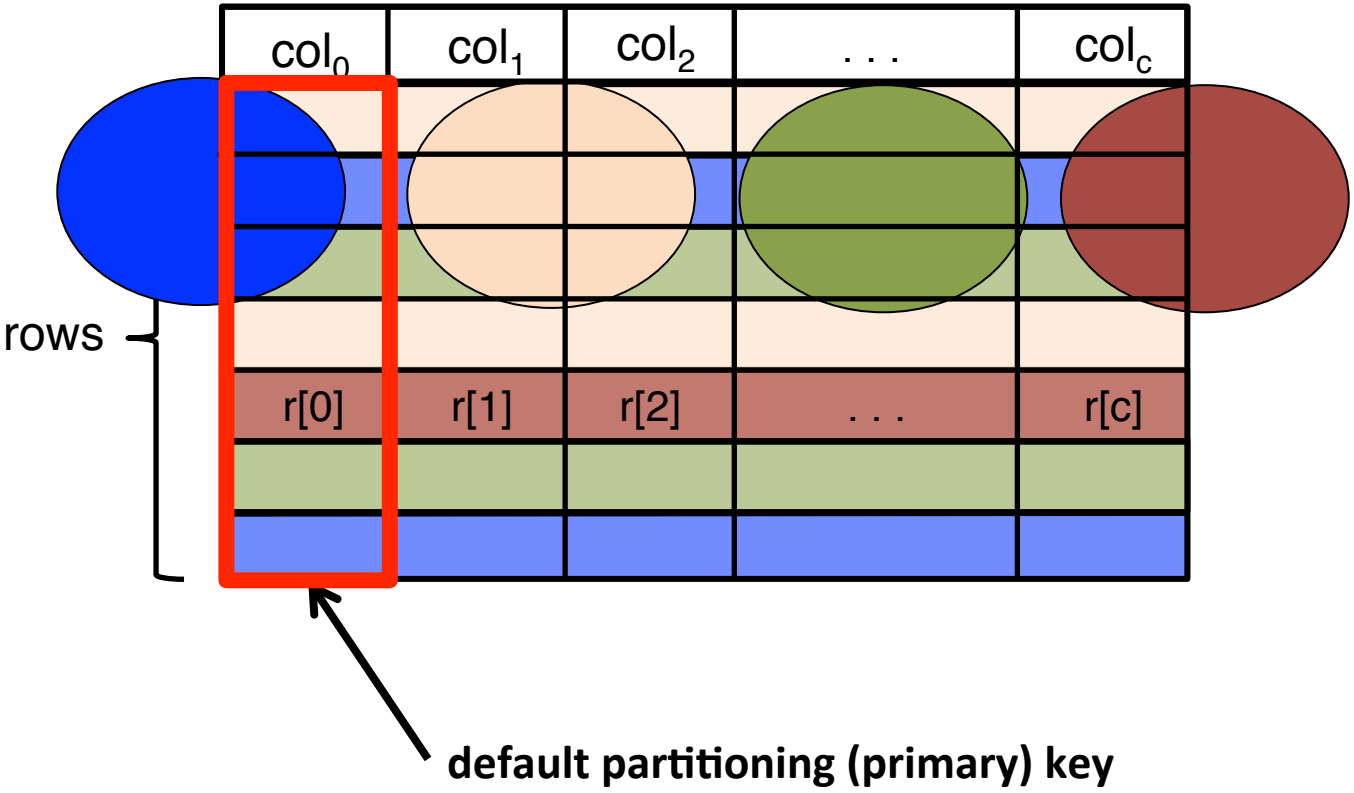


# SQL vs NoSQL

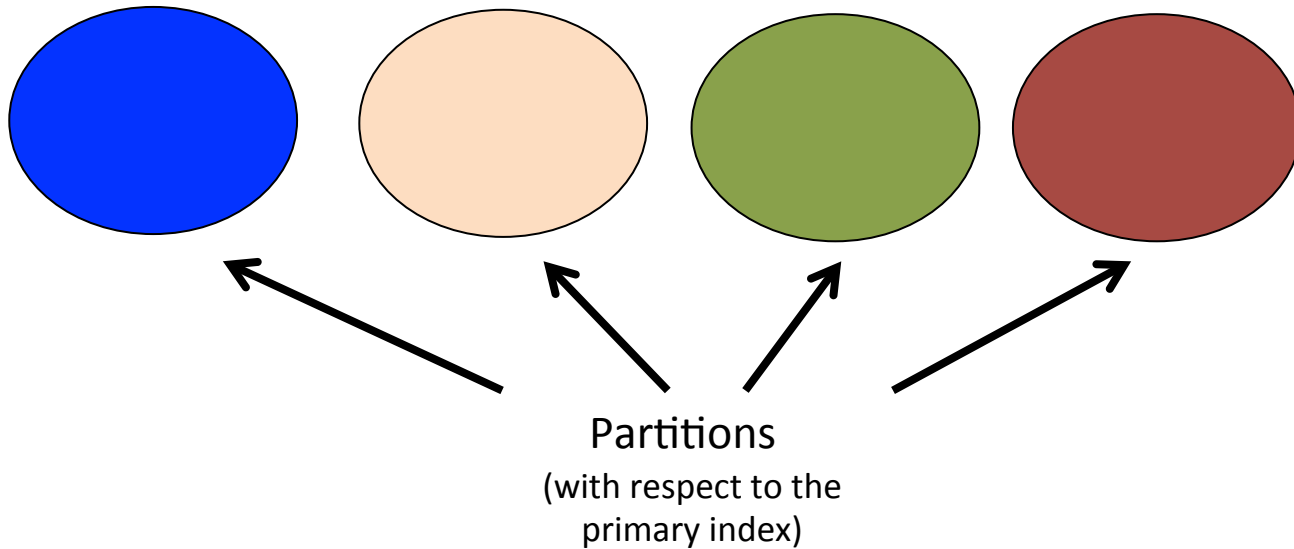


**Replex enables richer queries  
without sacrificing shared-  
nothing scale-out**

# NoSQL Scales with Shared-Nothing Partitioning



# NoSQL Scales with Shared-Nothing Partitioning



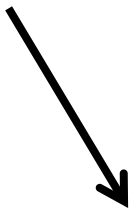
# Key Observations

- 1. Indexing enables richer queries**  
(searches, joins, etc.)

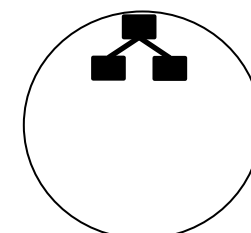
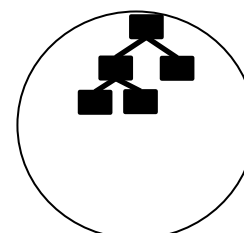
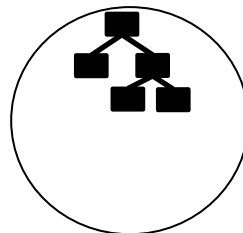
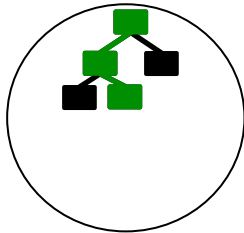
# Approach 1: Local Indexing

Index stored locally at each partition

insert  
(r[0] r[1] r[2] . . . r[c])



col <sub>0</sub>	col <sub>1</sub>	col <sub>2</sub>	...	col <sub>c</sub>

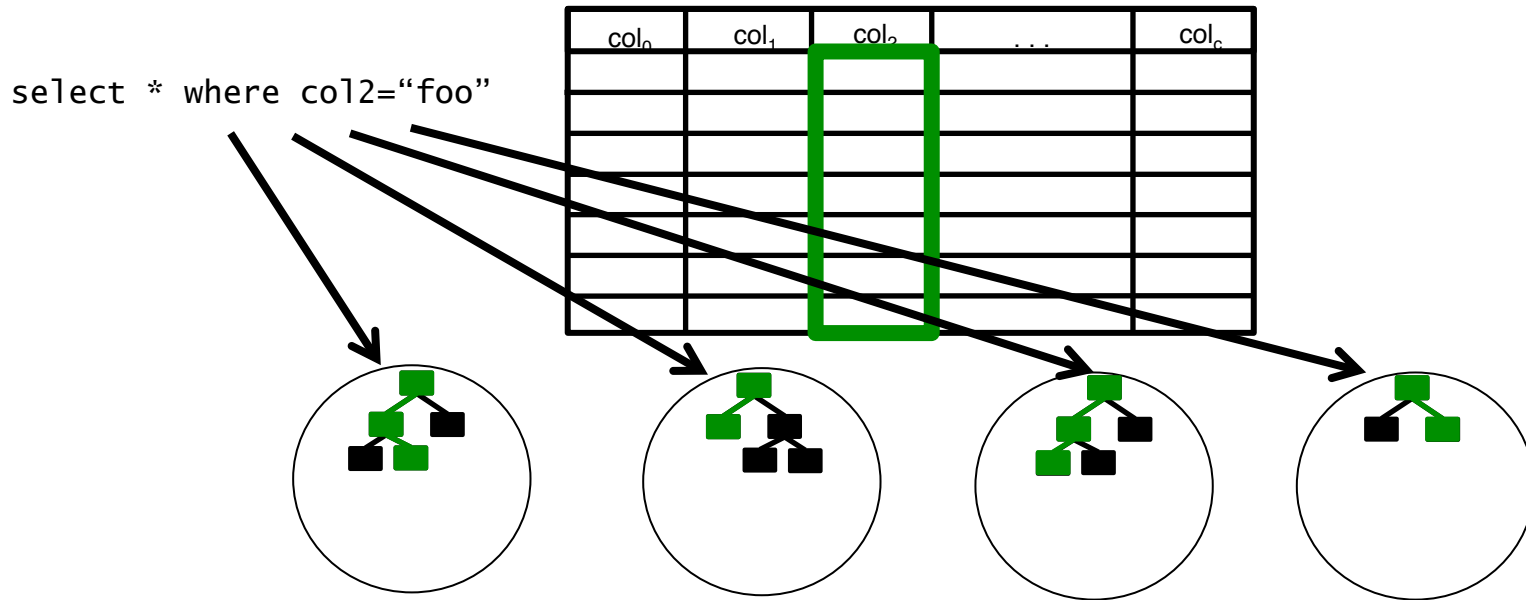


Index updates are local



# Approach 1: Local Indexing

Each partition builds a local index



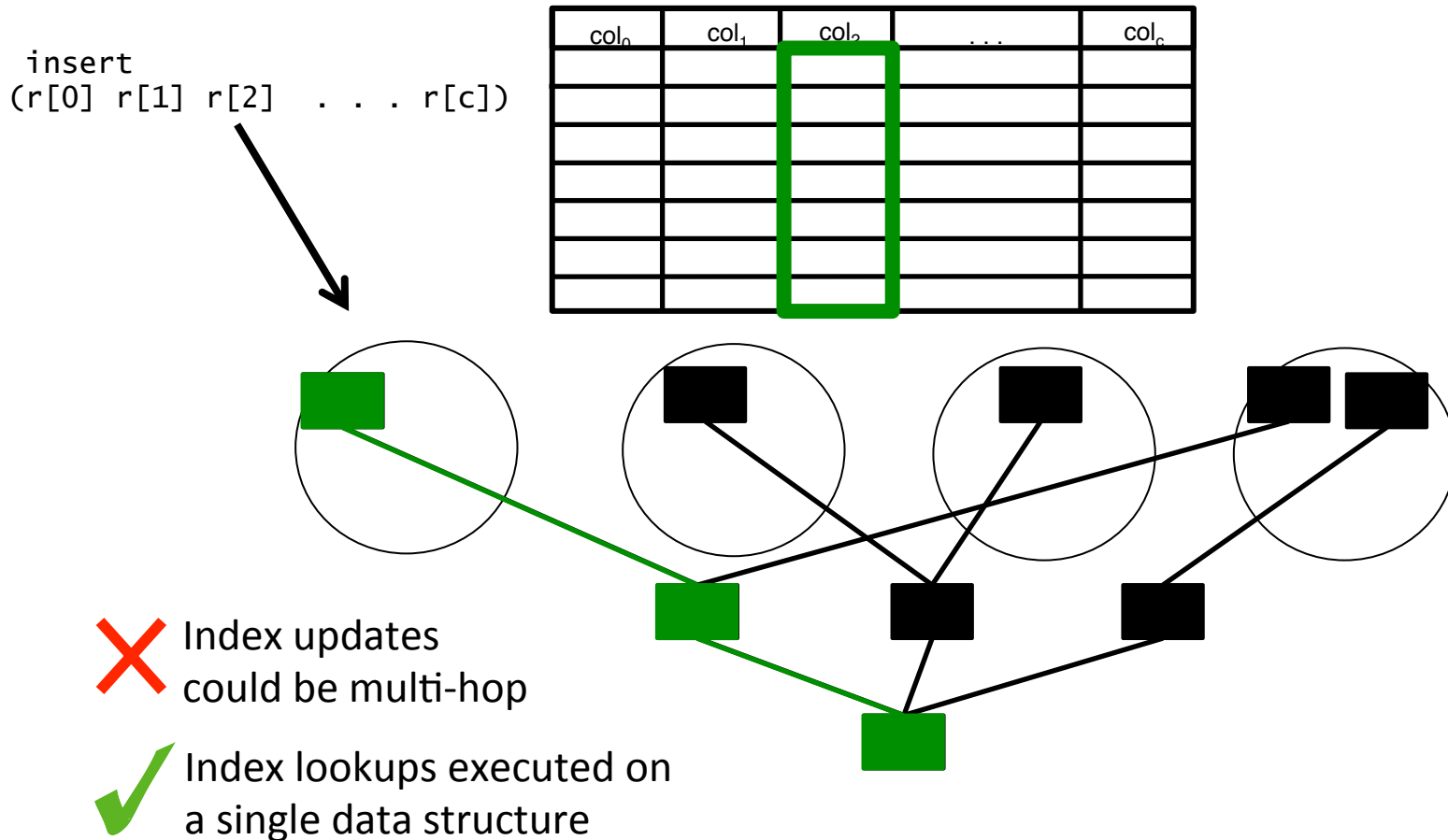
✓ Index updates are local

✗ Index lookups must be broadcast to all partitions

✗ Potential synchronization across partitions

# Approach 2: Global Indexing

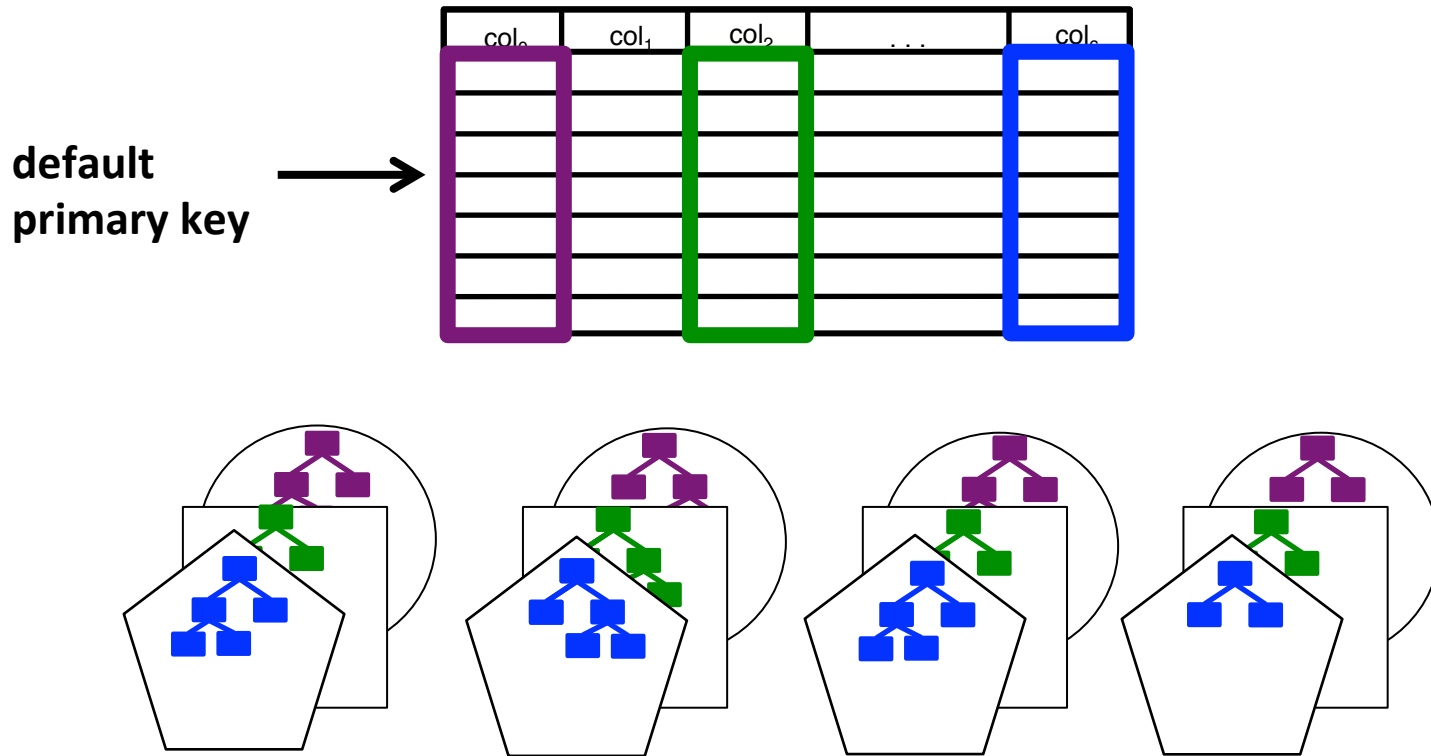
Distributed data structure spans all partitions



# Key Observations

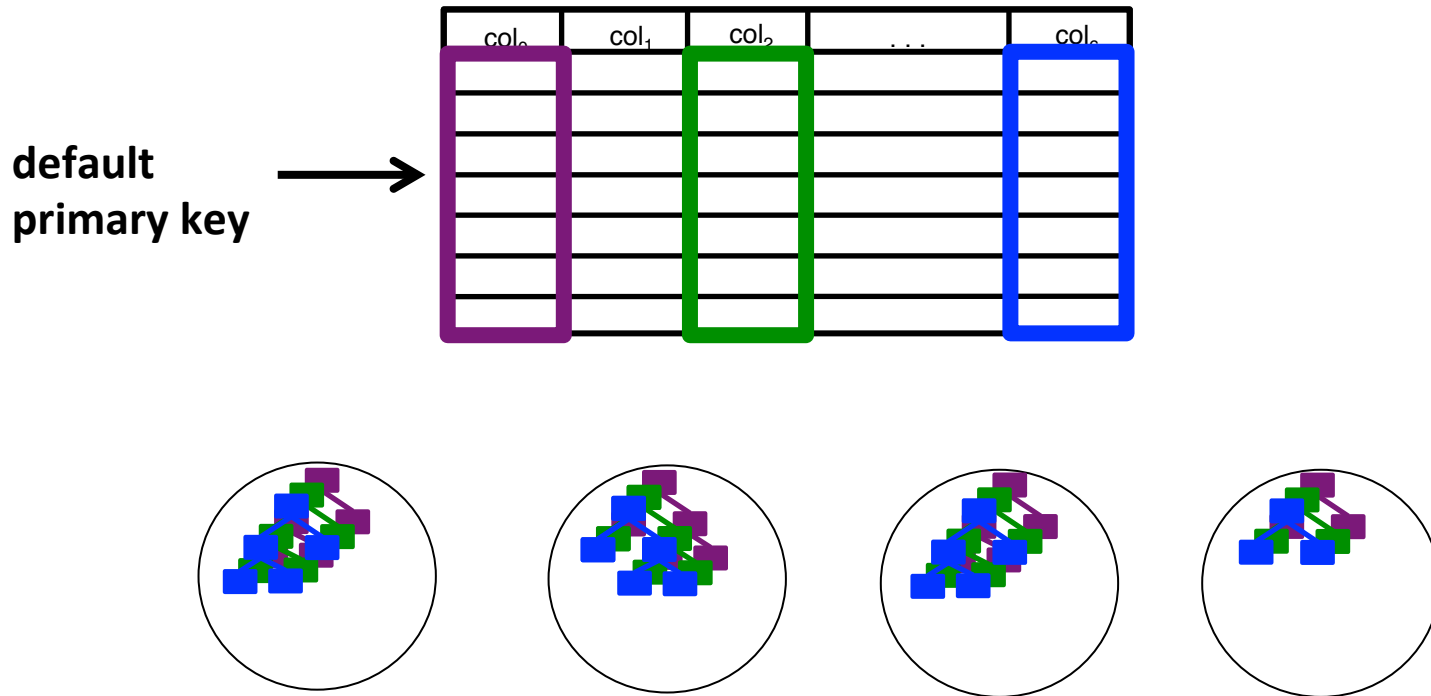
1. Indexing enables richer queries  
(searches, joins, etc.)
- 2. Indexes more efficient if data is  
partitioned according to that index**

# Partitioning by Index $\rightarrow$ Storage Overheads



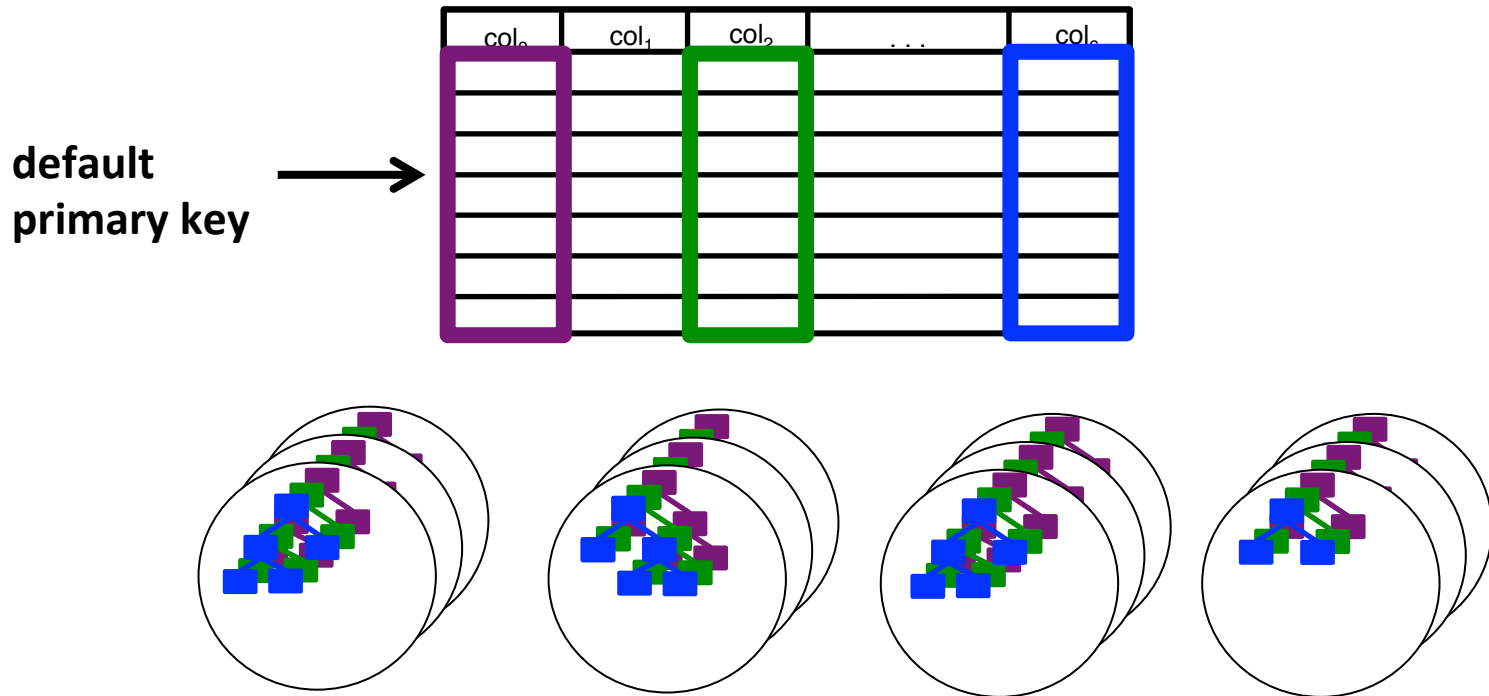
Partitioning by index  $\rightarrow$  must store data again

# Partitioning by Index $\rightarrow$ Storage Overheads



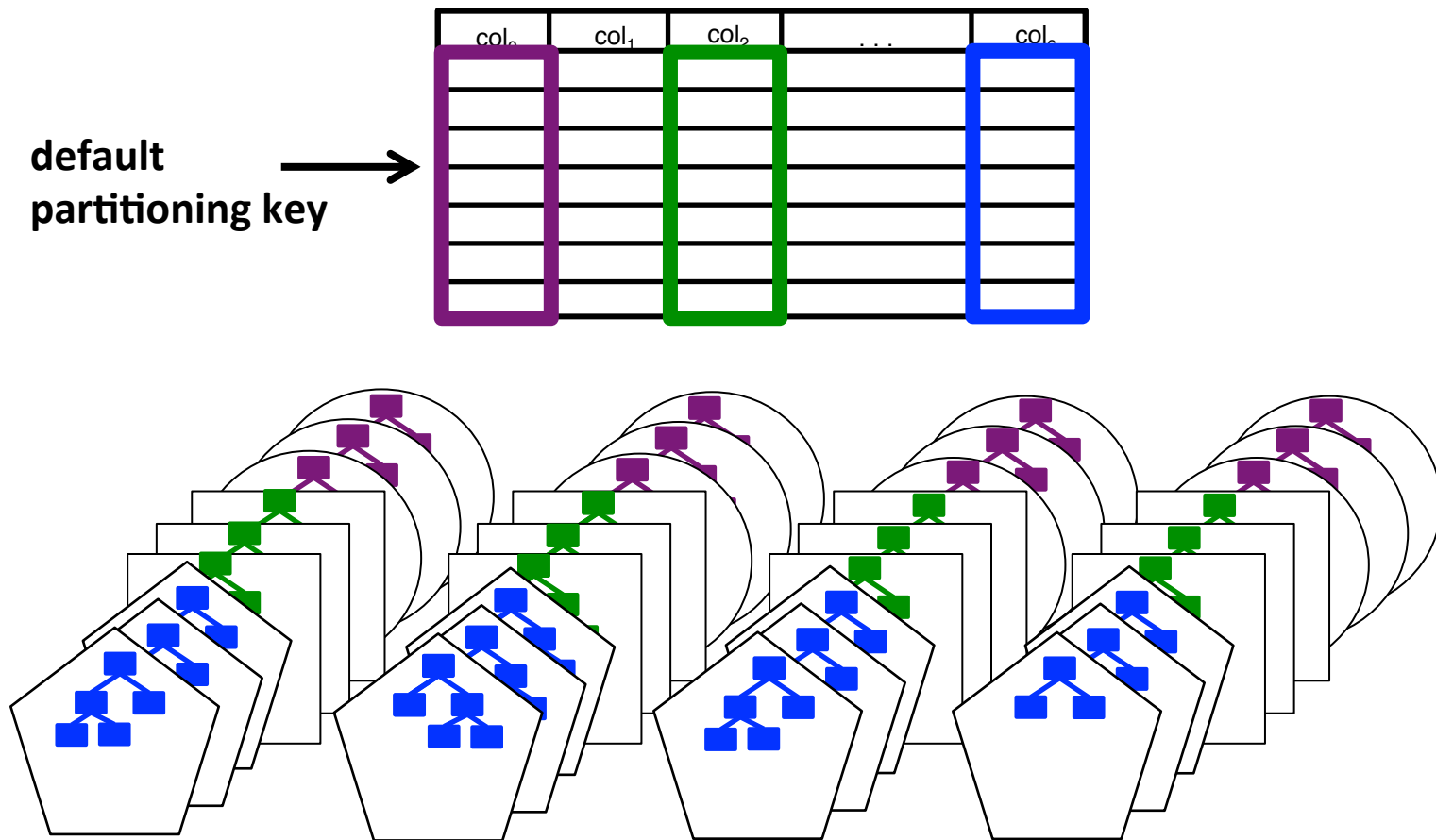
Ideally, build all indexes over a single copy of data

# Partitioning by Index $\rightarrow$ Storage Overheads

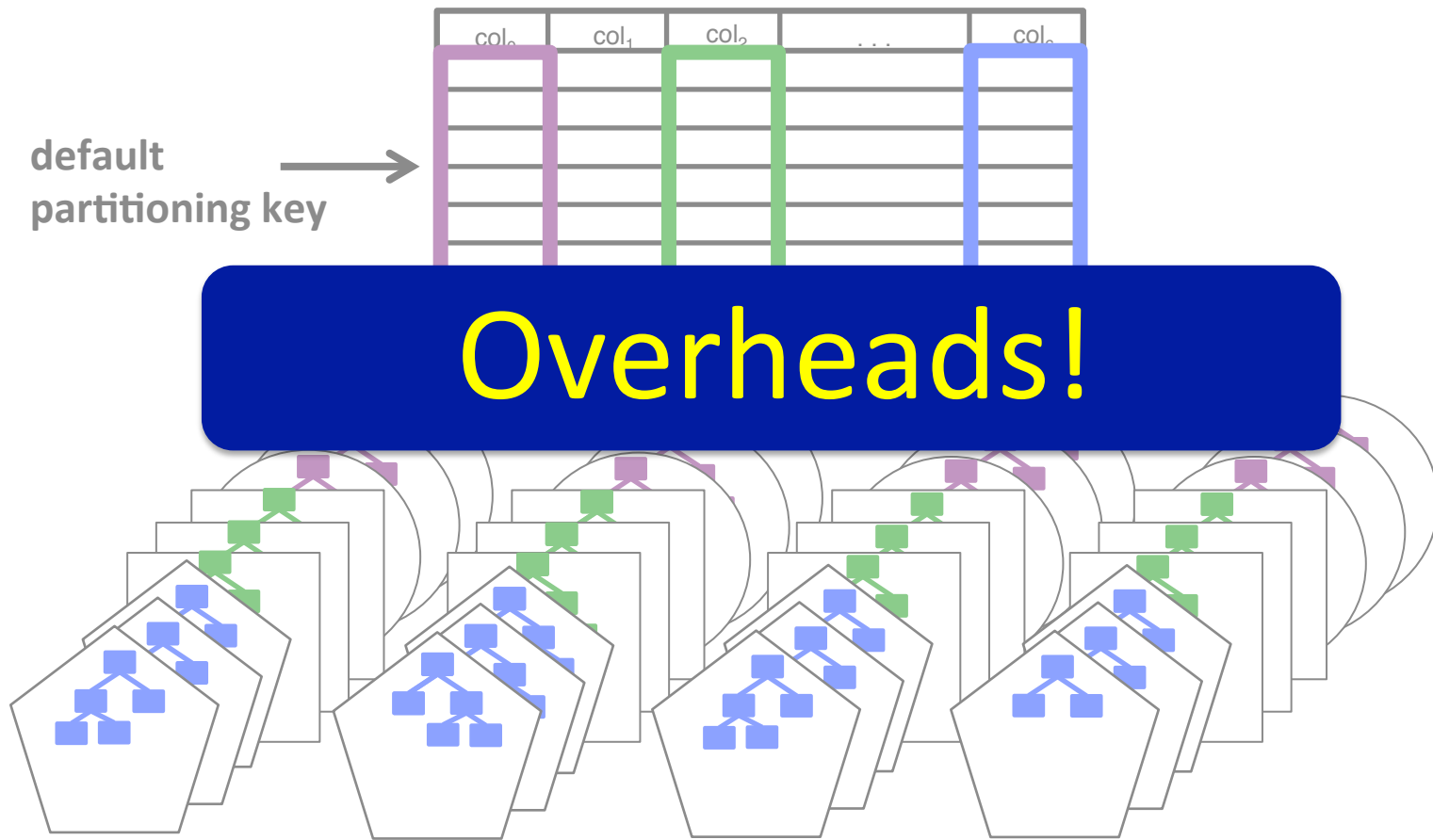


Ideally, build all indexes over a single copy of data  
Then, replication also replicates indexes

# Partitioning by Index $\rightarrow$ Storage Overheads



# Partitioning by Index $\rightarrow$ Storage Overheads



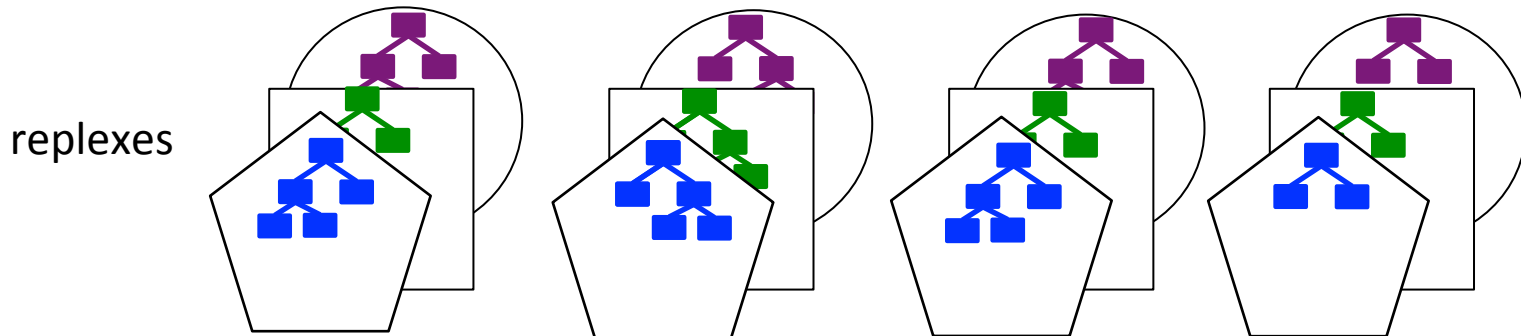


**Replex solves the indexing  
problem by combining indexing  
and replication**

# Replex

New replication unit:

**replex** -- data replica partitioned and sorted with respect to an associated index



Serves data **replication** and **indexing**

# Replex

New replication unit:

Replace data replicas with replexes

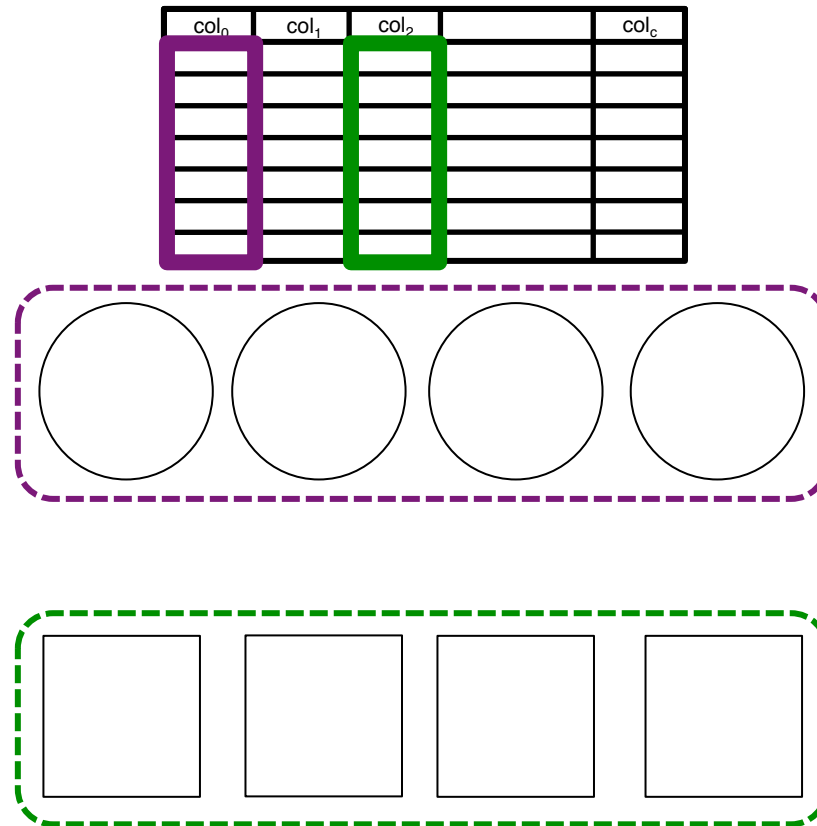


Indexing comes free during replication

The diagram shows four server nodes arranged horizontally. Each node is represented by a trapezoidal base with a semi-circular top. Inside each semi-circle, there is a purple tree-like structure with three levels: a root node at the top, two child nodes in the middle, and two leaf nodes at the bottom. A green horizontal bar is positioned below the tree structure in each node.

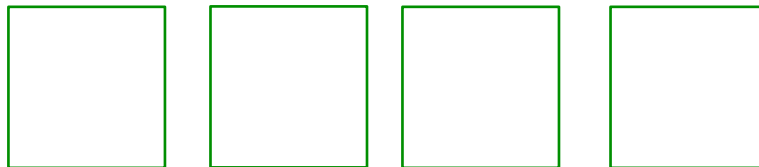
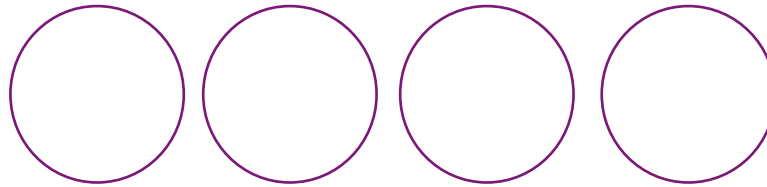
Serves data **replication** and indexing

# System Architecture



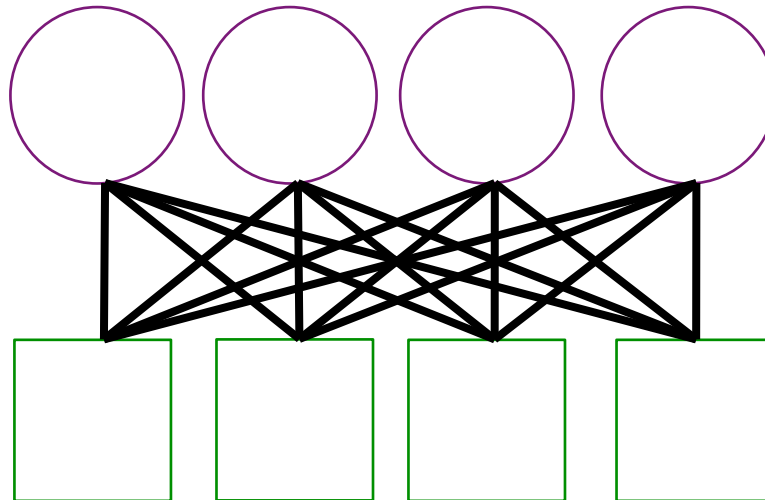
# Inserts in Replex

Replex uses a modified chain replication protocol

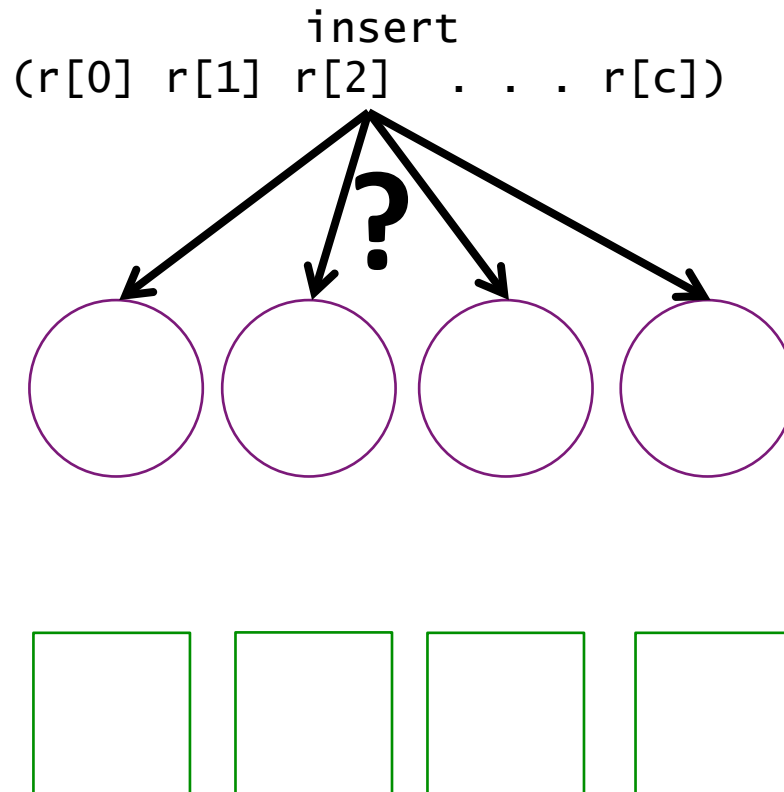


# Which chain?

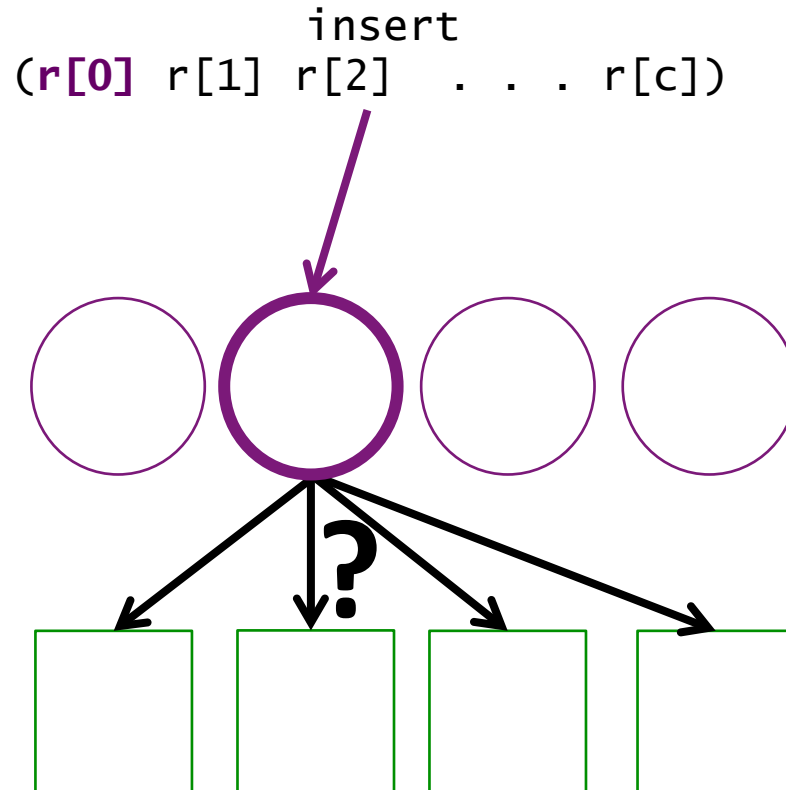
All pairs of partitions are potential chains



# Inserts in Replex



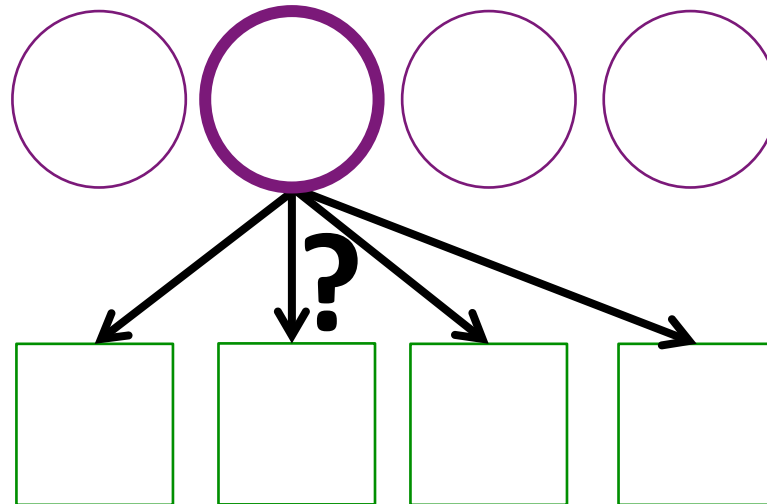
# Partition Determined by a replex's Index





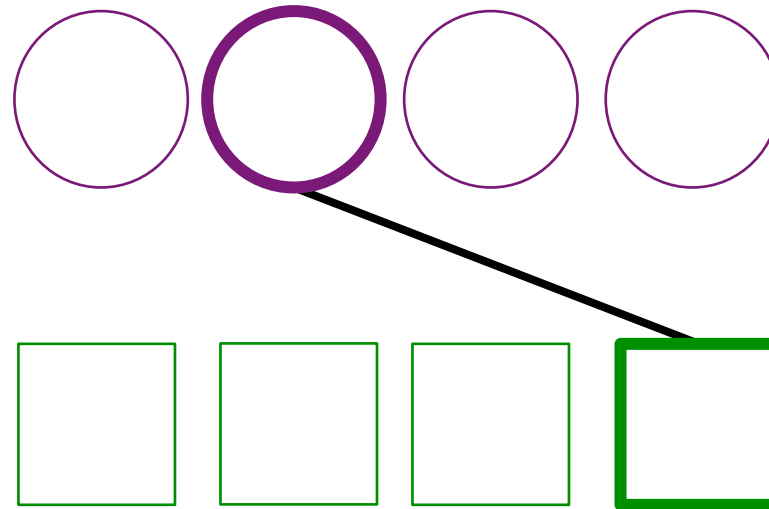
# Partition Determined by a replex's Index

insert  
(r[0] r[1] **r[2]** . . . r[c])

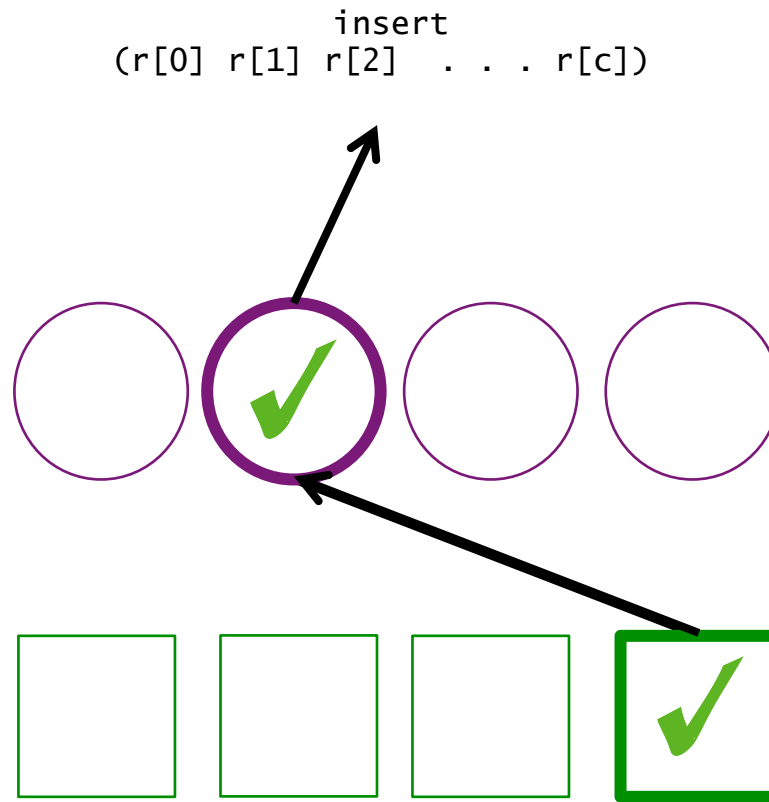


# Partition Determined by a replex's Index

insert  
(r[0] r[1] **r[2]** . . . r[c])

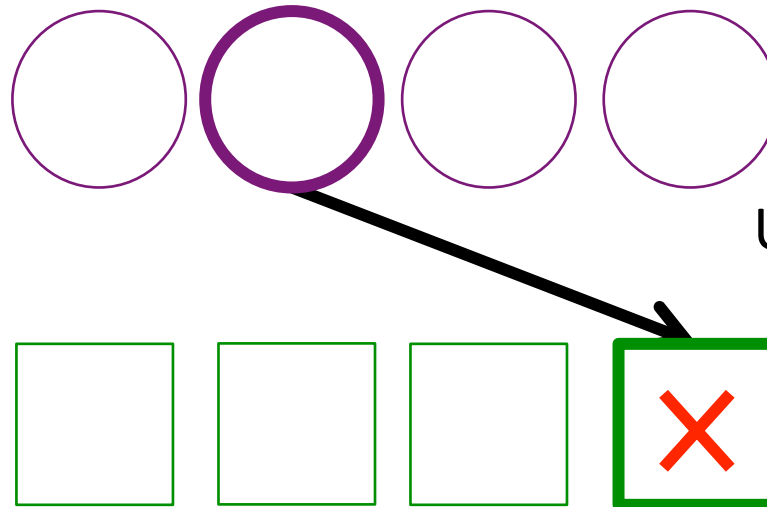


# Propagating a Commit Bit



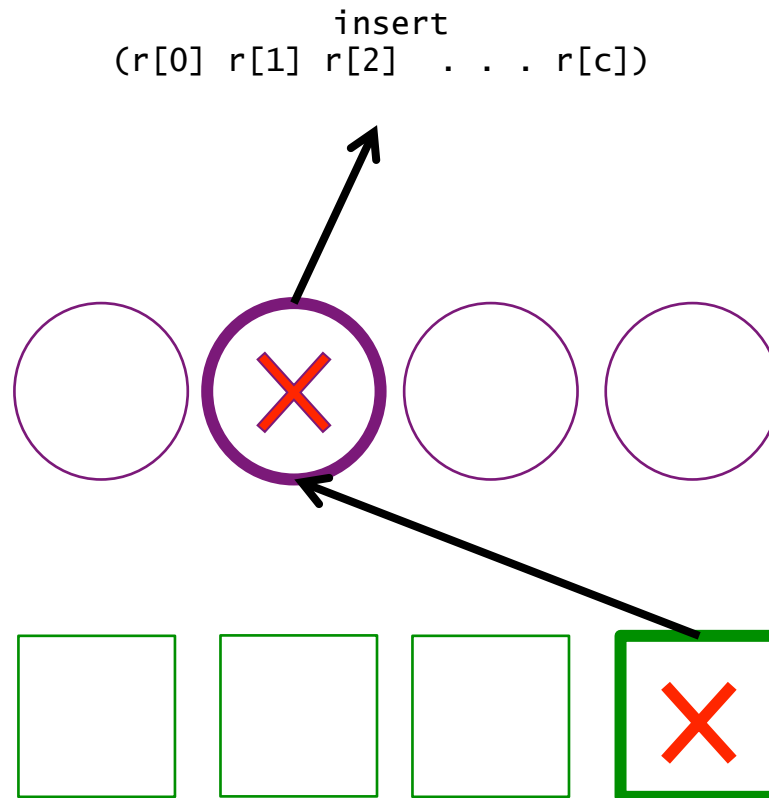
# Commit bits can be aborts

insert  
(r[0] r[1] r[2] . . . r[c])



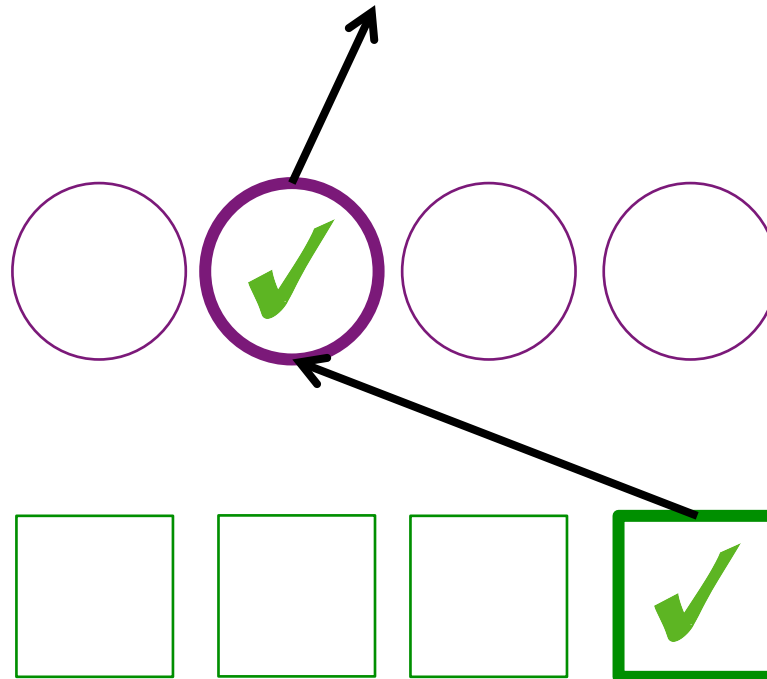
Uniqueness constraint  
on green replex can  
cause abort

# Commit bits can be aborts



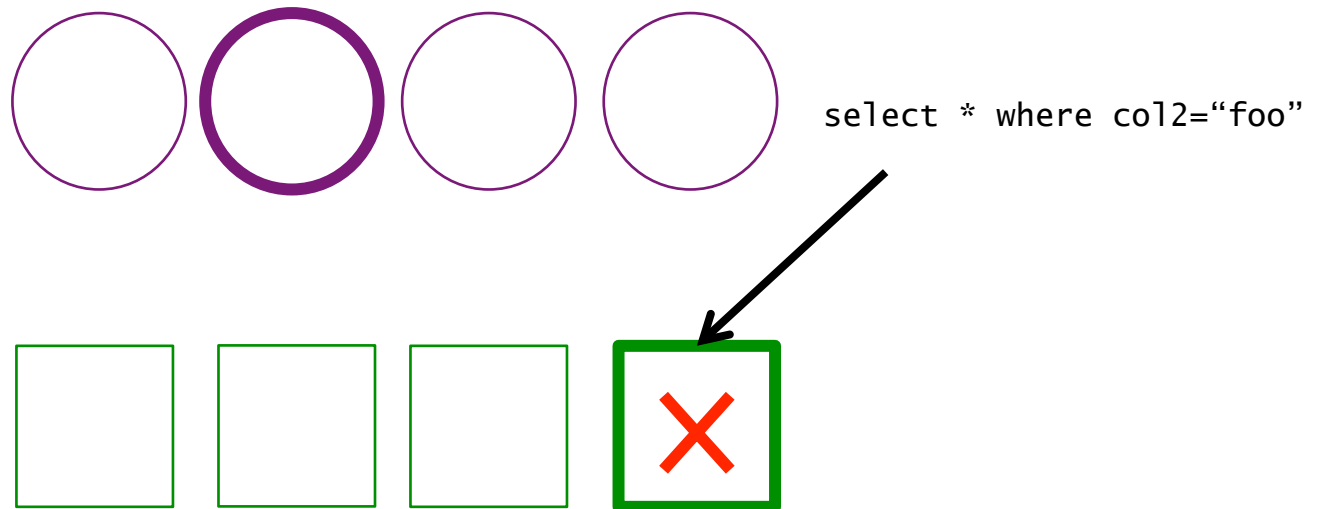
# Indexing is *Free*

insert  
(r[0] r[1] r[2] . . . r[c])



# Index Reads in Replex

1. Check for a commit bit
2. Only if the bit is true, can row be returned



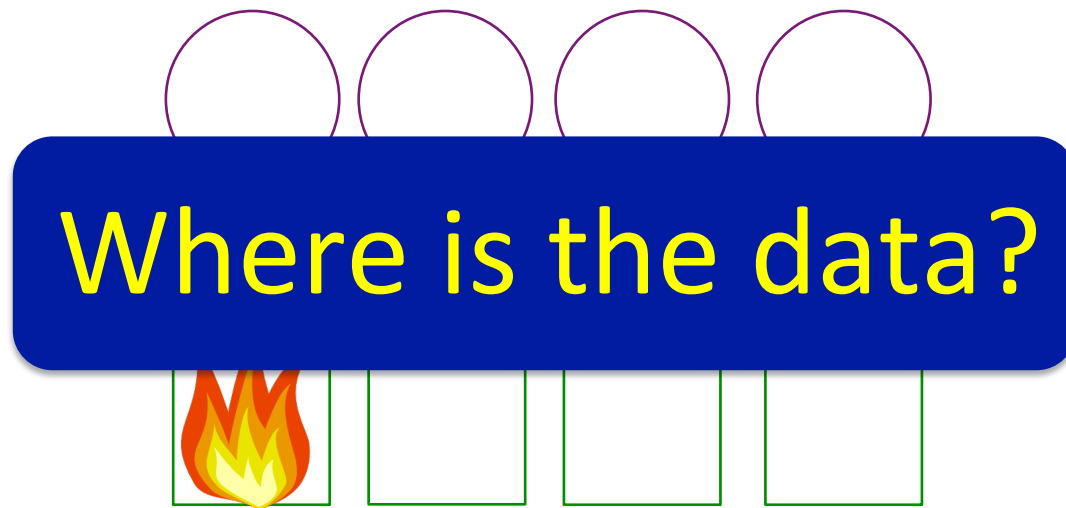
# Partition Failures

- 
1. Partition Recovery
  2. Client Requests

Index is unavailable, **data is available**



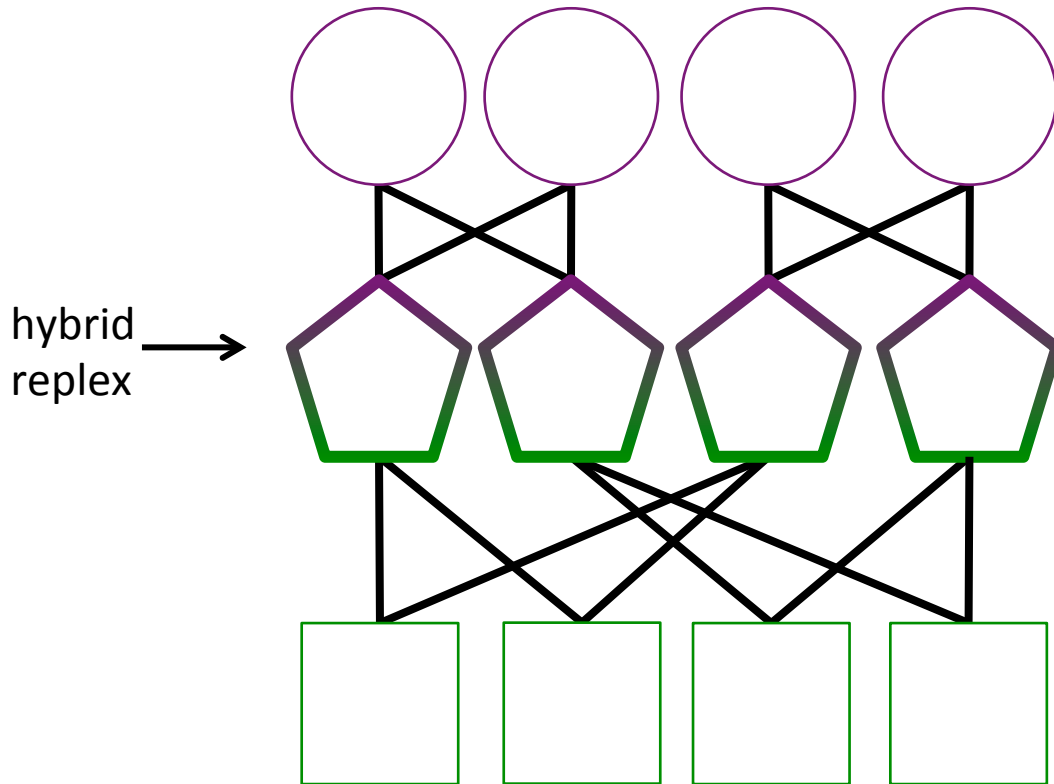
# Partition Failures



Index is unavailable, **data is available**



# Hybrid Replex

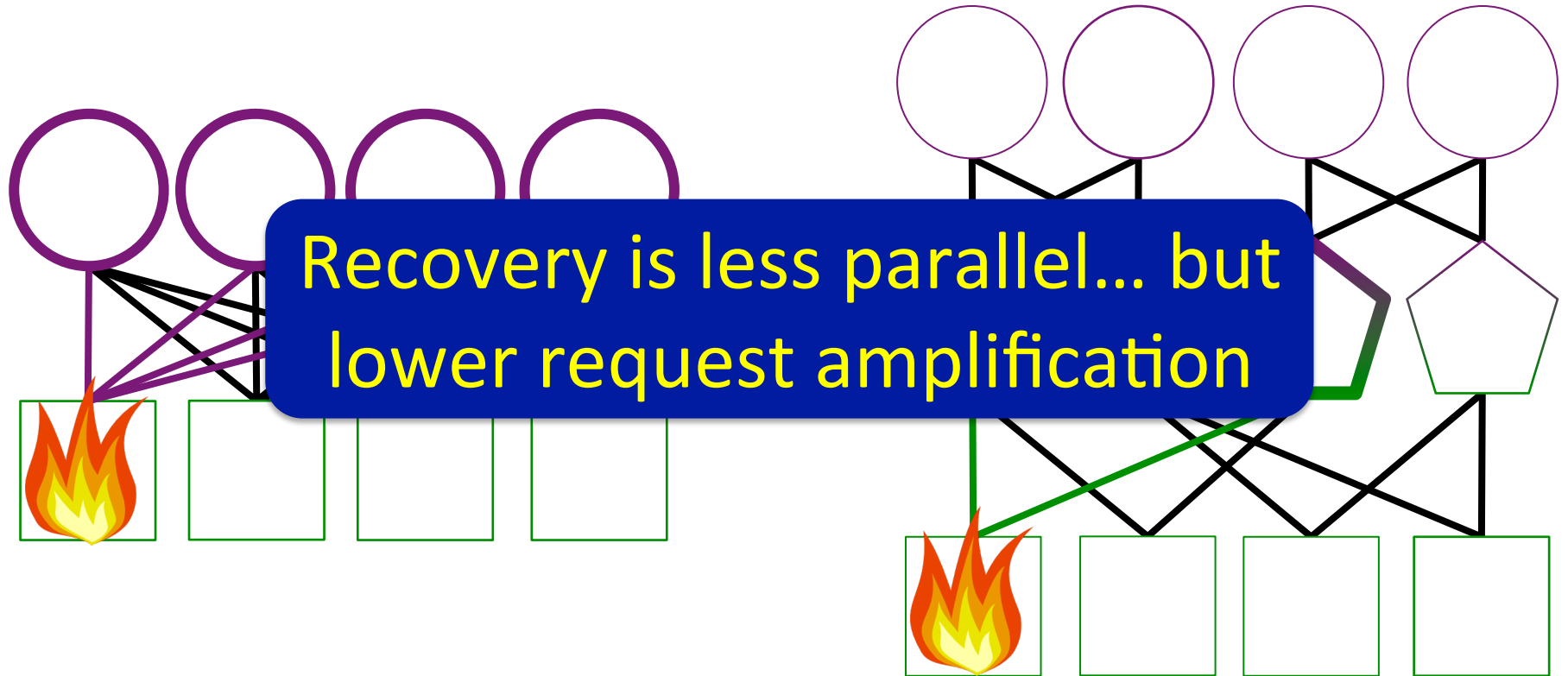


Hybrid replexes  
constrain potential  
data chains



More targeted recovery

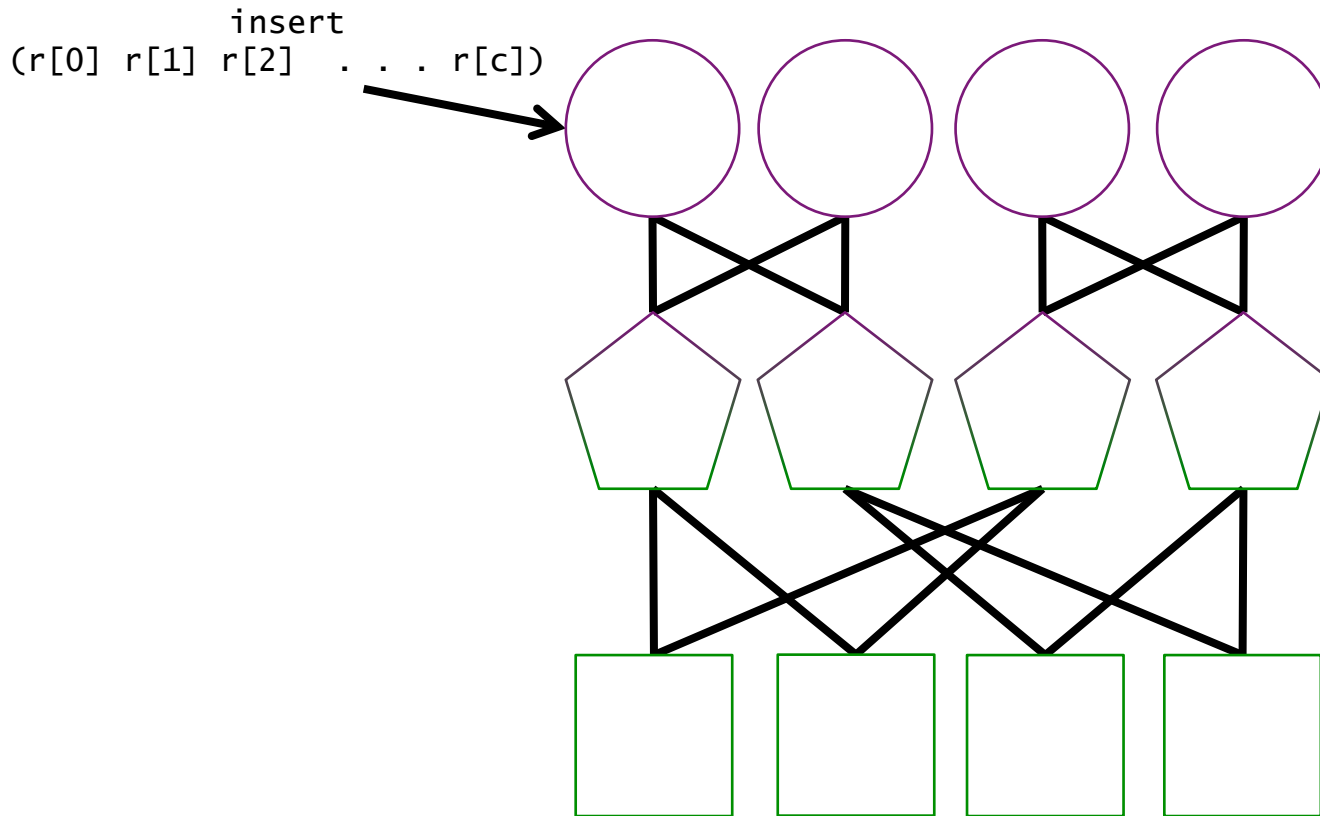
# Recovery time vs Request Amplification



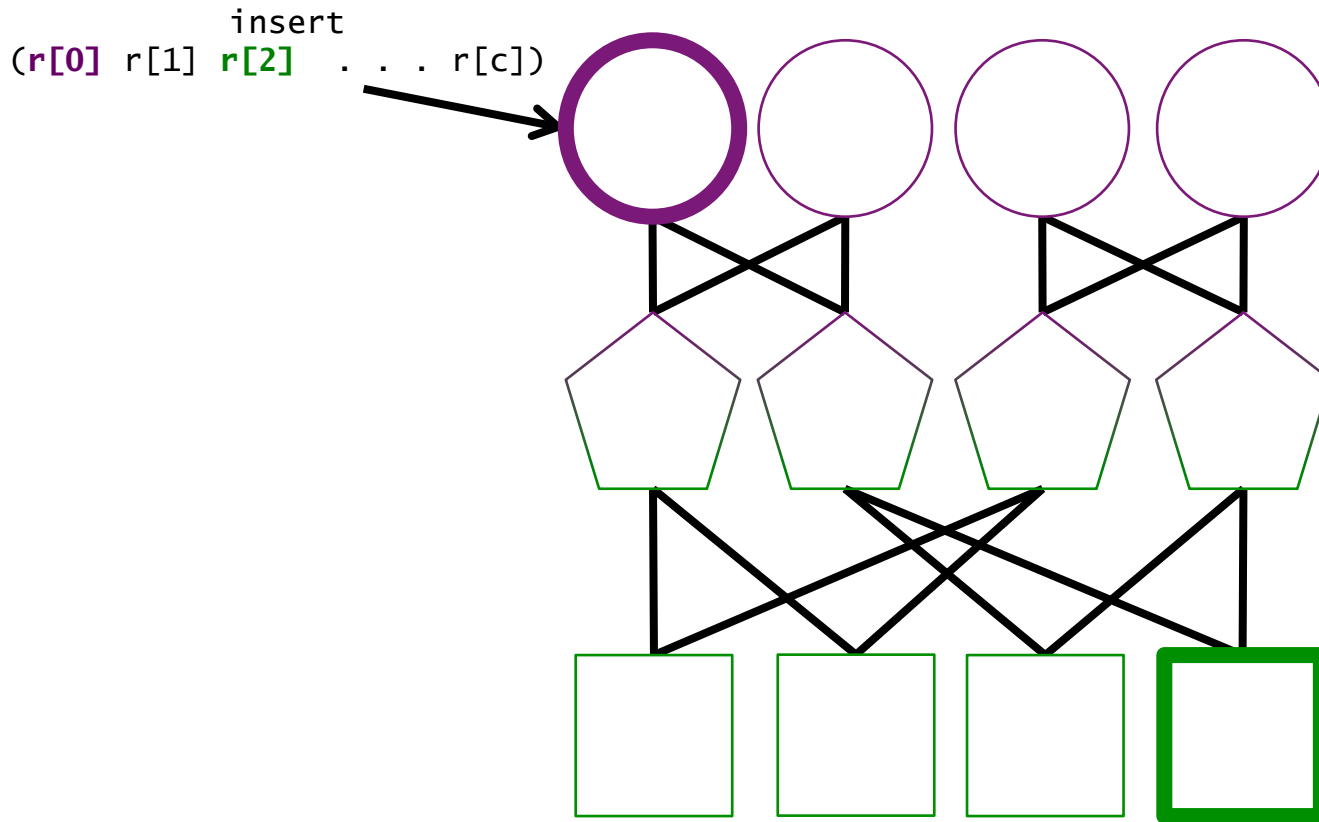
# Hybrid Replex: Definition

- Shared across  $r$  replexes
- Not associated with an index
- Partitioning function dependent on these  $r$  replexes

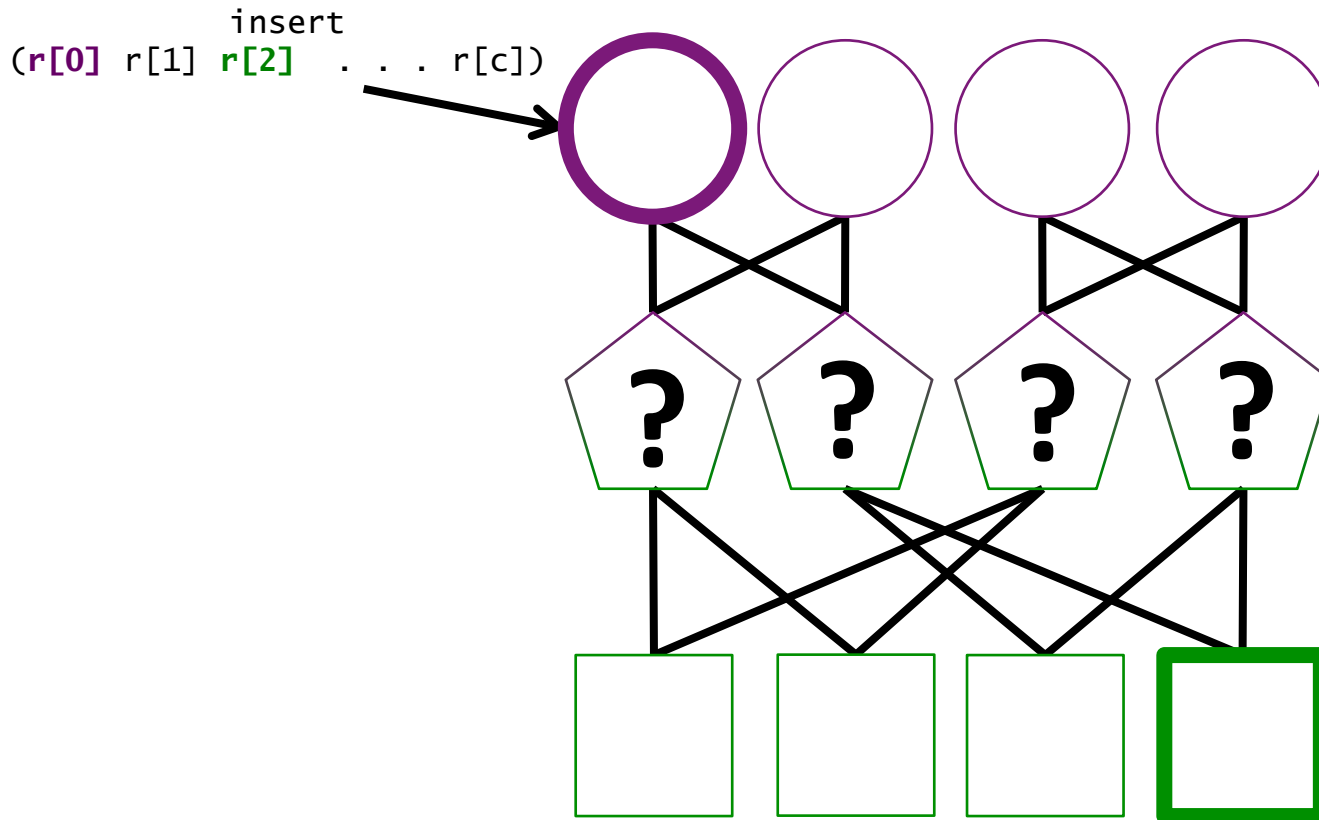
# Replication Chains with Hybrid Replex



# Replication Chains with Hybrid Replex

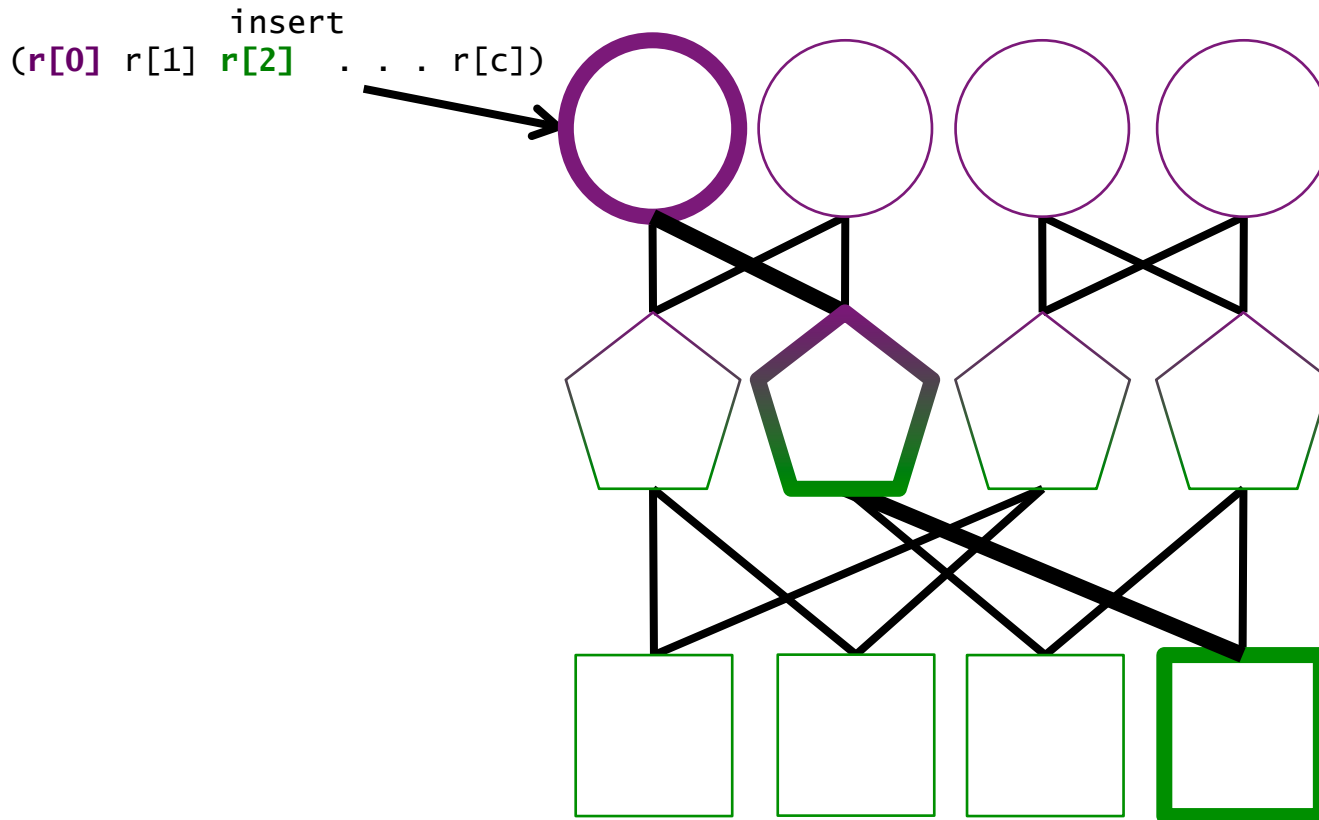


# Replication Chains with Hybrid Replex





# Replication Chains with Hybrid Replex



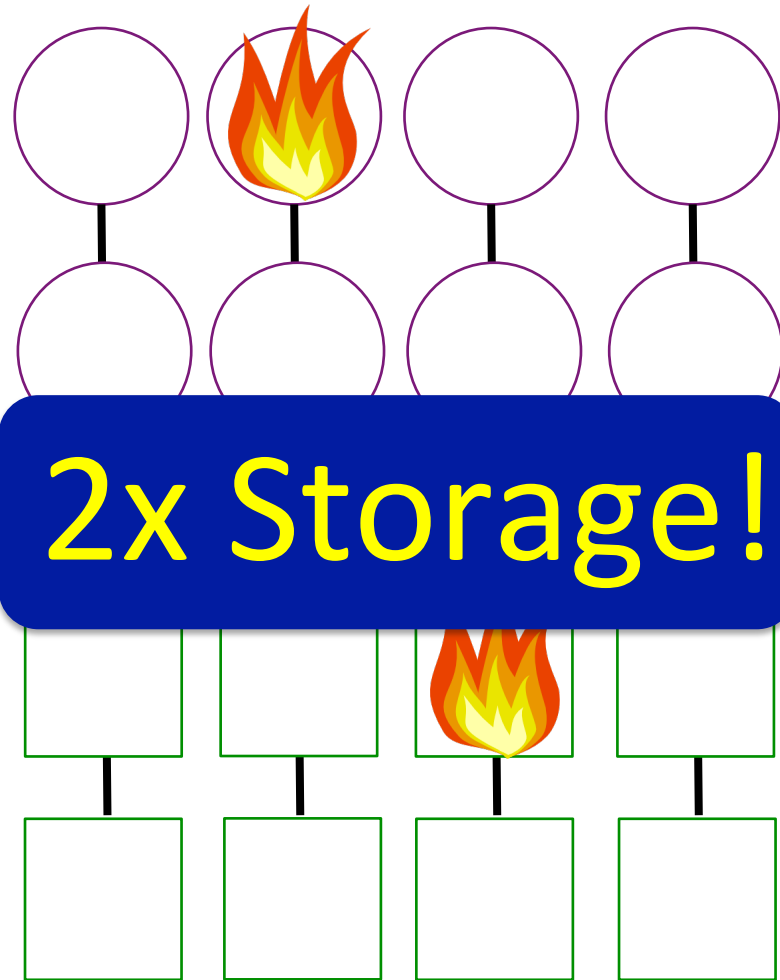
# Hybrid Replex: More Properties

1. Recovery time vs request amplification tradeoff
2. Improve failure availability of multiple replexes
3. Storage vs recovery performance tradeoff
4. Graceful failure degradation

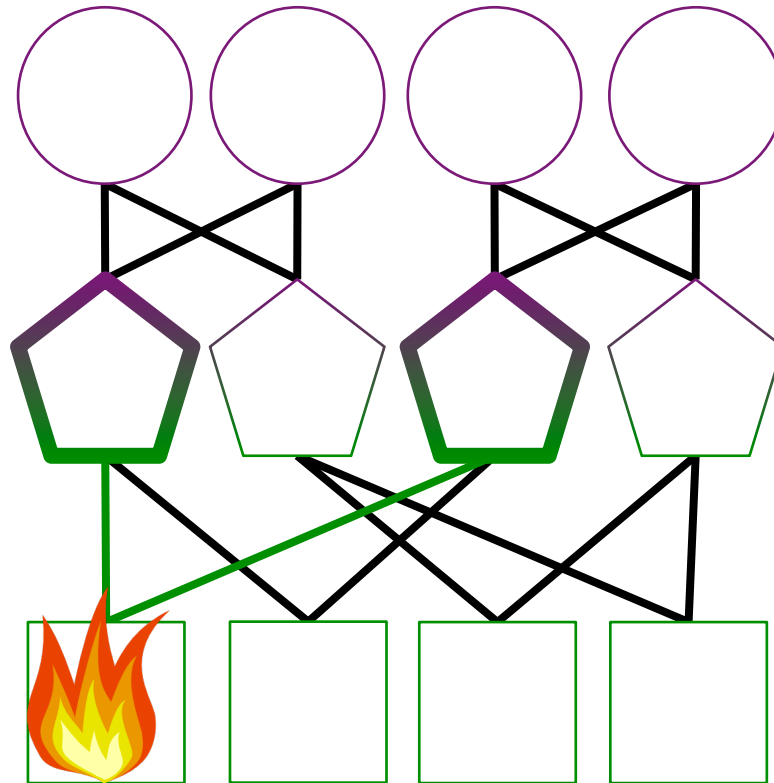
# Hybrid Replex: More Properties

1. Recovery time vs request amplification tradeoff
- 2. Improve failure availability of multiple replexes**
- 3. Storage vs recovery performance tradeoff**
4. Graceful failure degradation

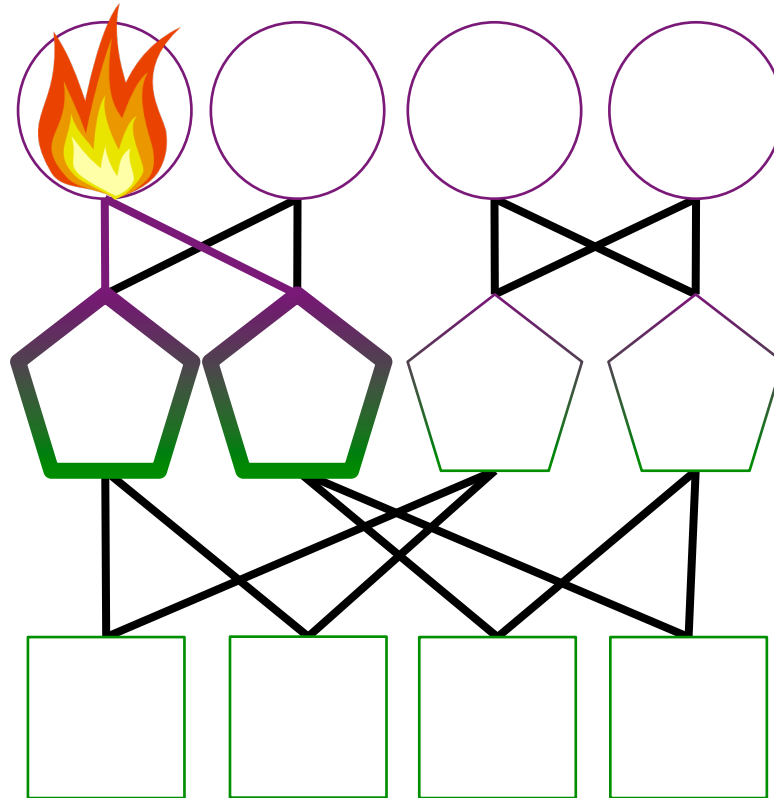
# Improving Failure Availability w/o Hybrid Replexes



# Improve failure availability of multiple replexes



# Improve failure availability of multiple replexes



# Storage vs. recovery performance



# Implementation

- Built on top of HyperDex, ~700 LOC
- Implemented partition recovery and request re-routing on failure
- Implemented variety of hybrid replex configurations



# Evaluation

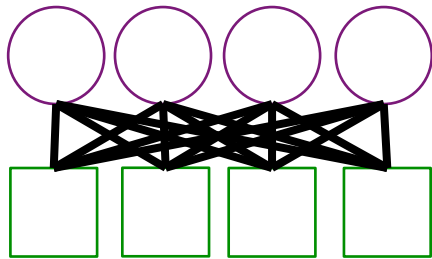
1. What is impact of replexes on steady-state performance?
2. How do hybrid replexes affect failure performance?

# Evaluation Setup

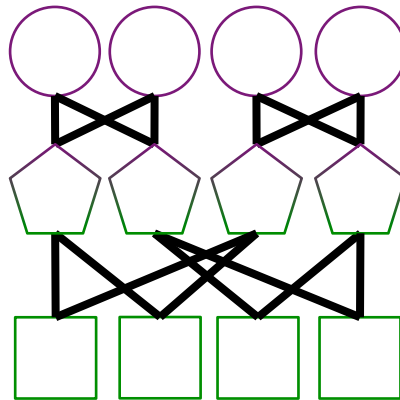
- Table with two indexes
- 12 server machines, 4 client machines
- All machines colocated in the same rack, connected via 1GB top-of-rack switch
- 8 CPU, 16GB memory per machine

# Systems Under Evaluation

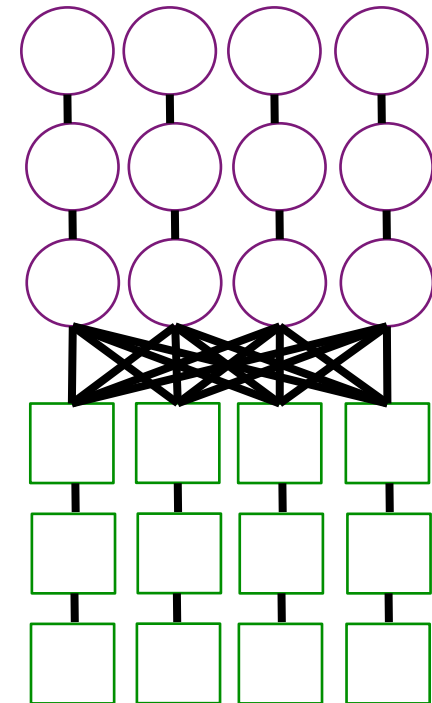
Replex-2



Replex-3

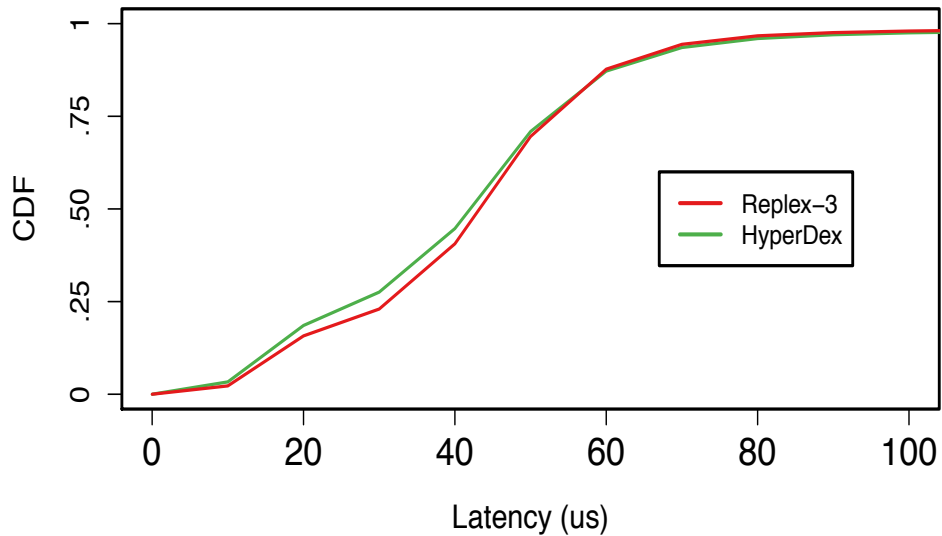


HyperDex



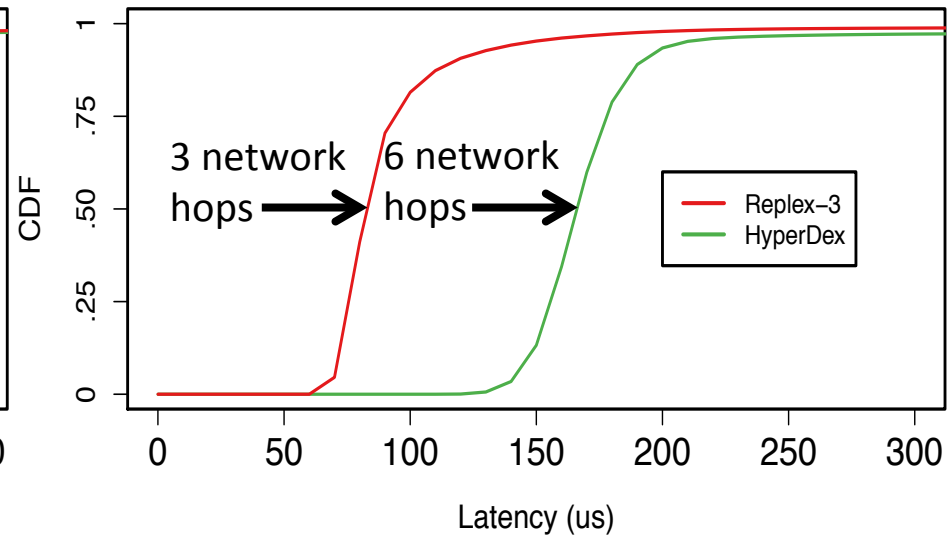
# Steady State Latency

## Reads



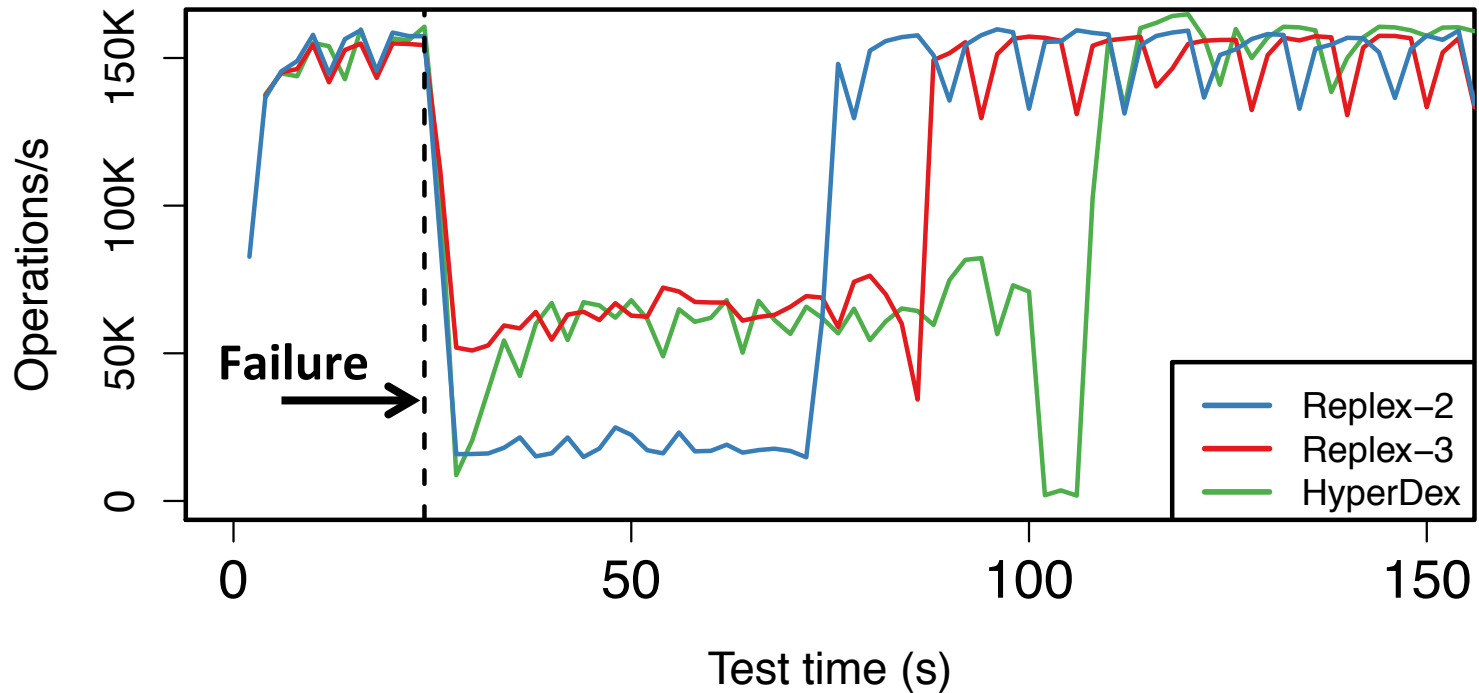
Reads against either index have similar latency, but we report reads against primary index

## Inserts



Replex-2 not included because it has a lower fault tolerance threshold

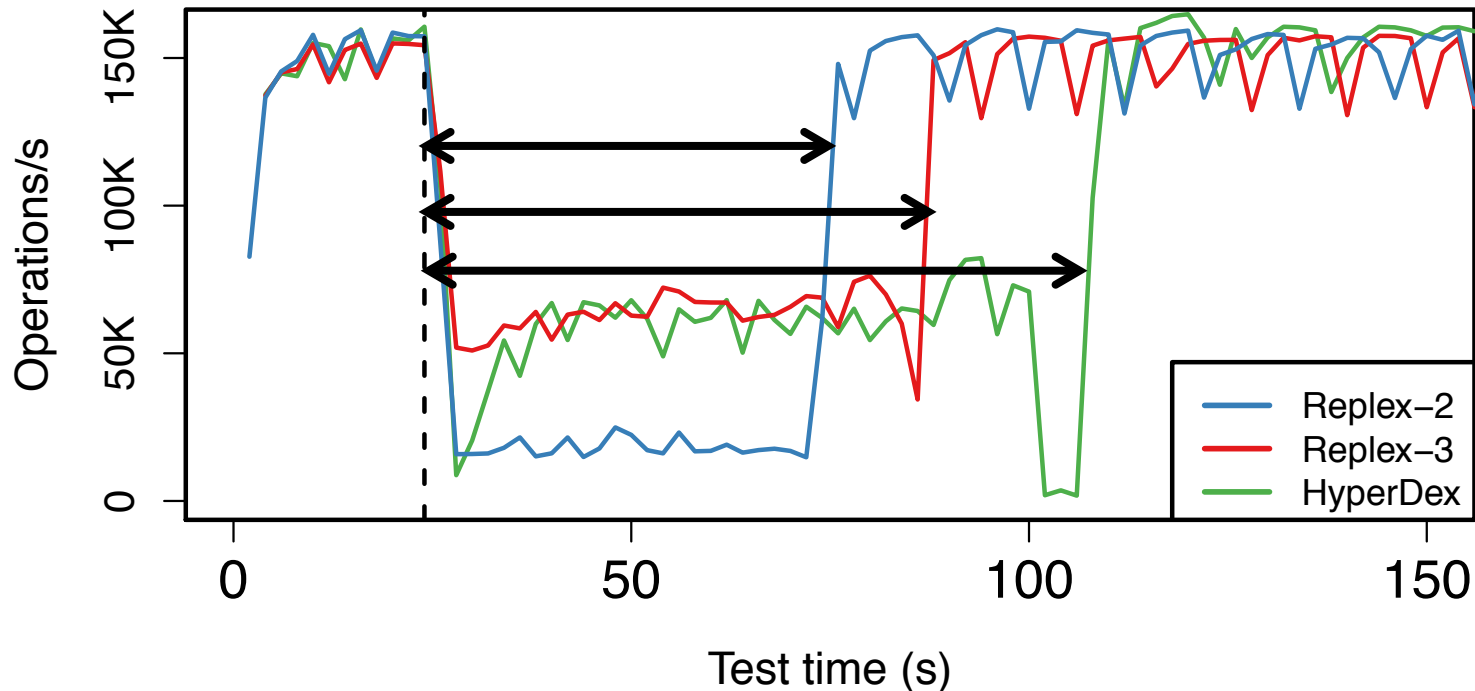
# Single Failure Performance



## Experiment

- Load with 10M, 100 byte rows
- Split reads 50:50 between each index

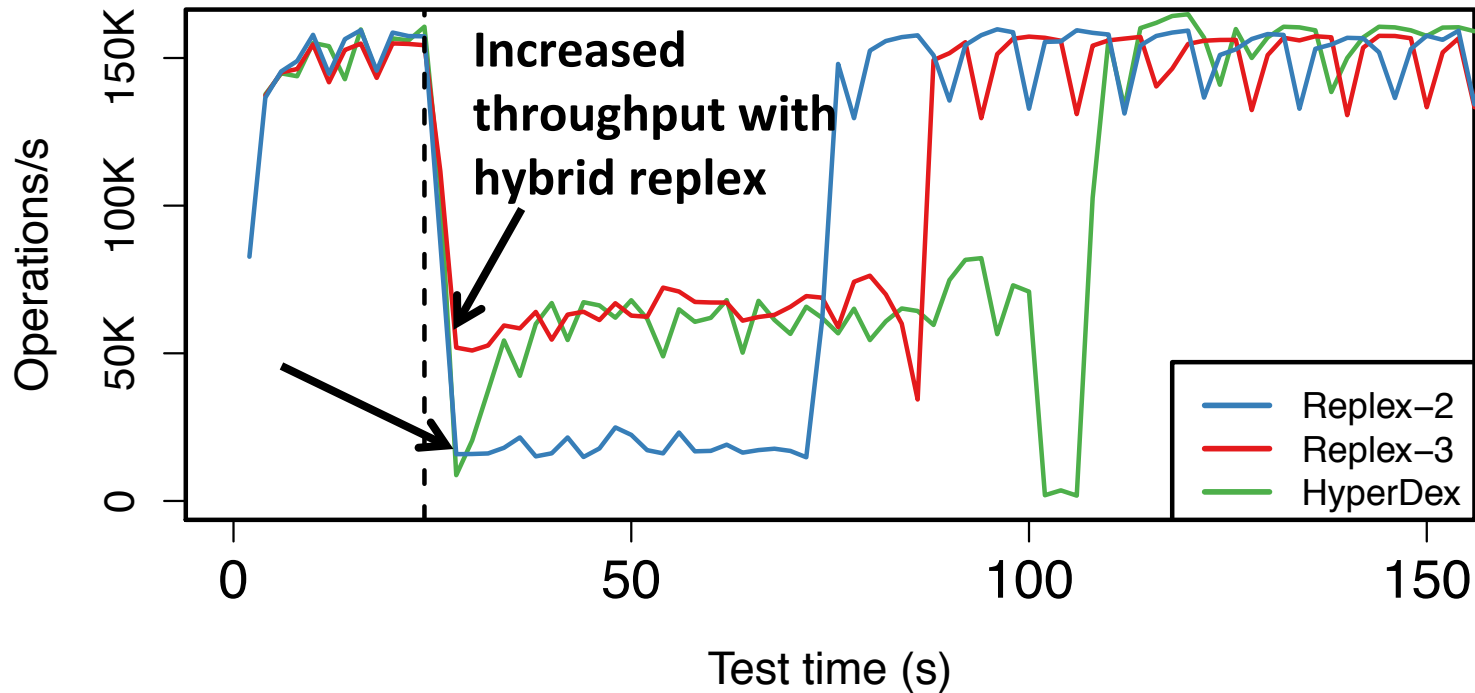
# Single Failure: Recovery Time



## Recovery Time

1. HyperDex recovers slowest because 2-3x more data
2. Replex-2 recovers fastest because least data, parallel recovery

# Single Failure: Failure Throughput



## Failure Throughput

1. Replex-2 low throughput because of high request amplification
2. Replex-3 has throughput comparable to HyperDex

# Summary

- 1. Rethink the replication paradigm**
2. Replacing replicas with replexes decreases index storage AND maintenance overheads
3. Hybrid replexes introduce rich tradeoff space for failure SLAs



Questions?