# DynaMast
# *Adaptive* Dynamic Mastering for Replicated Systems

**Michael Abebe**
Brad Glasbergen
Khuzaima Daudjee

**mtabebe@uwaterloo.ca**

ICDE (April 2020)
**tiny.cc/dynamast**

UNIVERSITY OF **WATERLOO**

# Single Database

# Single Database



Single-site transactions

# Single Database

# Single Database

# Single Database



Overloads
site

C   D

# How to scale the database?

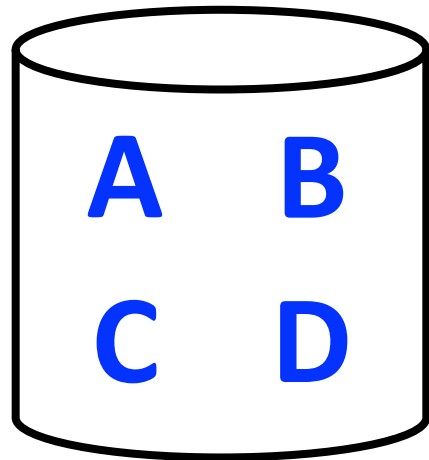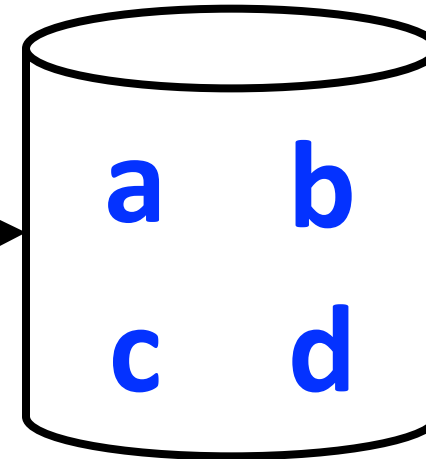**Single-master**

**Multi-master**

# How to scale the database?

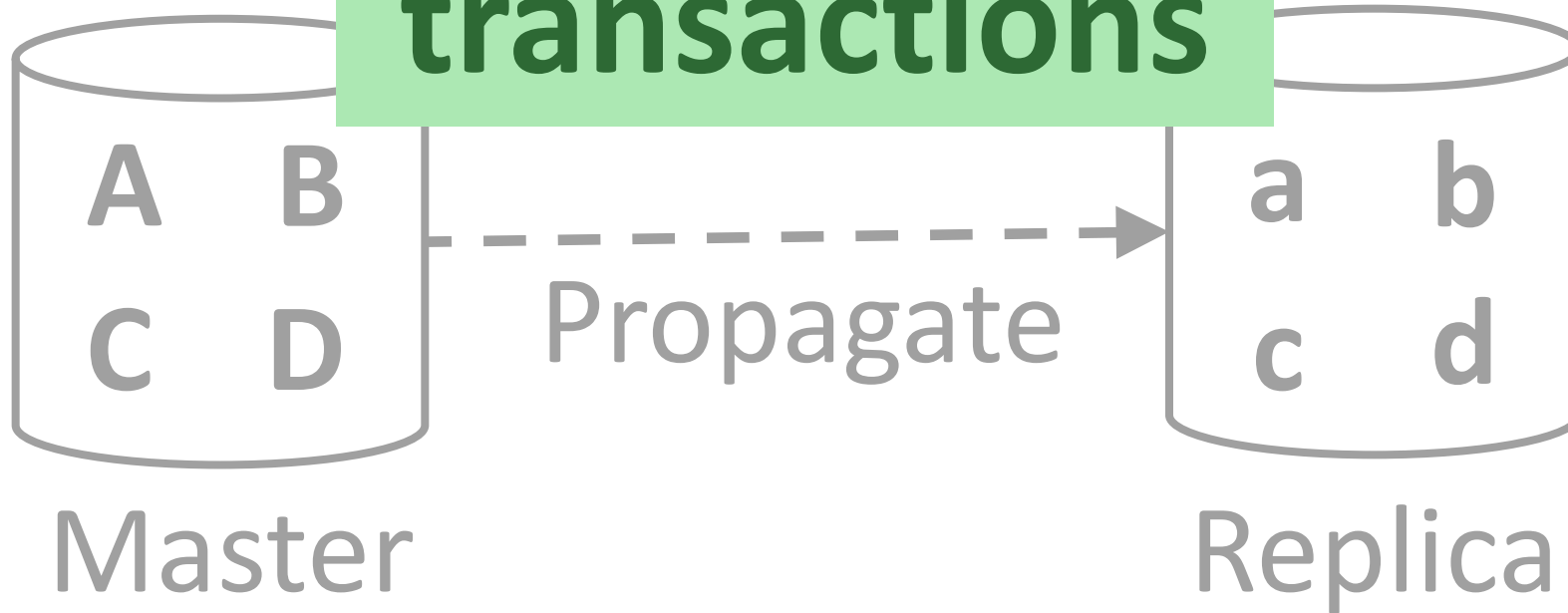**Single-master**

**Multi-master**

# Single-Master

Writers

Readers

A    B
C    D

Propagate

a    b
c    d

Master

Replica

# Single-Master

Writers                                                            Readers

**Single-site transactions**

A    B                          a    b
C    D        Propagate         c    d

Master                                      Replica

UNIVERSITY OF
WATERLOO

# Single-Master

Writers

Readers

**A**    **B**

**C**    **D**

Propagate

**a**    **b**

**c**    **d**

Master

Replica

UNIVERSITY OF
WATERLOO

# Single-Master



Writers

Readers

A    B
C    D

Propagate

Master

a    b
c    d

Replica

# Single-Master



Writers

Readers

A    B
C    D

Master

Propagate

**Overloads site**

a    b
c    d

Replica

UNIVERSITY OF
**WATERLOO**

# How to scale the database?

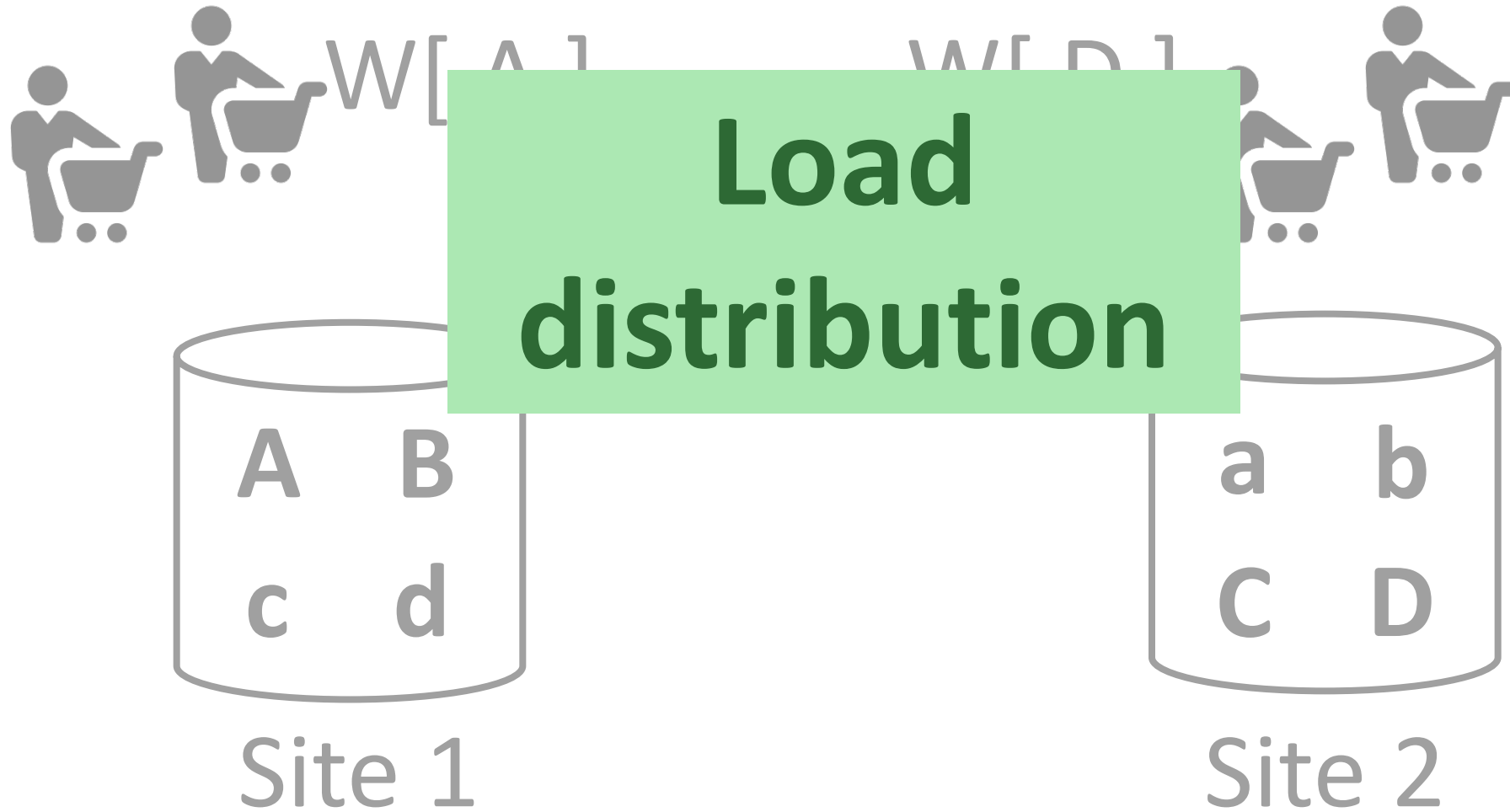**Single-master**

**Multi-master**

# Multi-Master

W[ A ]

W[ D ]

A    B
c    d

Site 1

a    b
C    D

Site 2

# Multi-Master

W[A] W[D]

Load distribution

A  B
c  d

a  b
C  D

Site 1

Site 2

# Multi-Master



W[ A, C ]

W[ C ]

A    B
c    d

a    b
C    D

prepare
commit

Site 1

Site 2

UNIVERSITY OF
WATERLOO

# Multi-Master

W[ A, C ]

W[ C ]

Costly Coordination Protocol

A  B

a  b

c  d

C  D

Site 1

commit

Site 2

UNIVERSITY OF
WATERLOO

# Multi-Master    Single-Master

Load distribution

Single-site transactions

Costly Coordination Protocol

Overloads site

UNIVERSITY OF WATERLOO

# How to provide:

**Load distribution**

**Single-site transactions**

**Dynamic Mastering**

# Dynamic Mastering

W[ A, C ]

A    B

c    d

Site 1

Remaster A

a    b

C    D

Site 2

# Dynamic Mastering

W[ A, C ]

a    B

c    d

A    b

C    D

**Remaster** A

Site 1

Site 2

# Dynamic Mastering

W[ A, C ]

**Single-site transactions**

a    B          A    b

c    d    Remaster A    C    D

Site 1                    Site 2

UNIVERSITY OF
**WATERLOO**

# Dynamic Mastering

W[ B ]

W[ A, C ]

a    B
c    d

Site 1

A    b
C    D

Site 2

# Dynamic Mastering

W[ B ]

W[ A, C ]

Load distribution

a    B

A    b

c    d

C    D

Site 1

Site 2

UNIVERSITY OF
**WATERLOO**

# Dynamic Mastering

**Outside** transaction boundaries

W[ A, C ]

W[ C ]

A   B

c   d

**Remaster** A
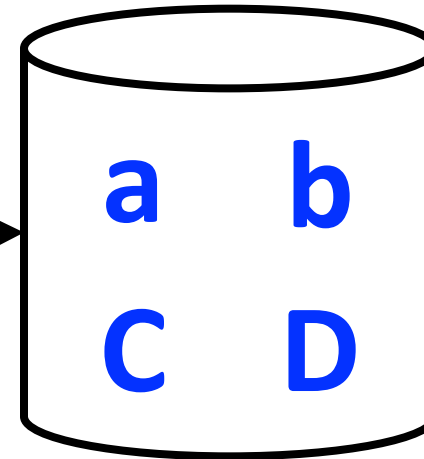
a   b

C   D

Site 1

Site 2

UNIVERSITY OF WATERLOO

# DynaMast

**Distributed and replicated** database system

Employs *adaptive* **dynamic mastering**

Provides both **single-site transactions** and **load balance**

UNIVERSITY OF
**WATERLOO**

# Dynamic Mastering Challenges

How to **perform remastering efficiently**?

How to **decide where to master** data?

# Dynamic Mastering Challenges

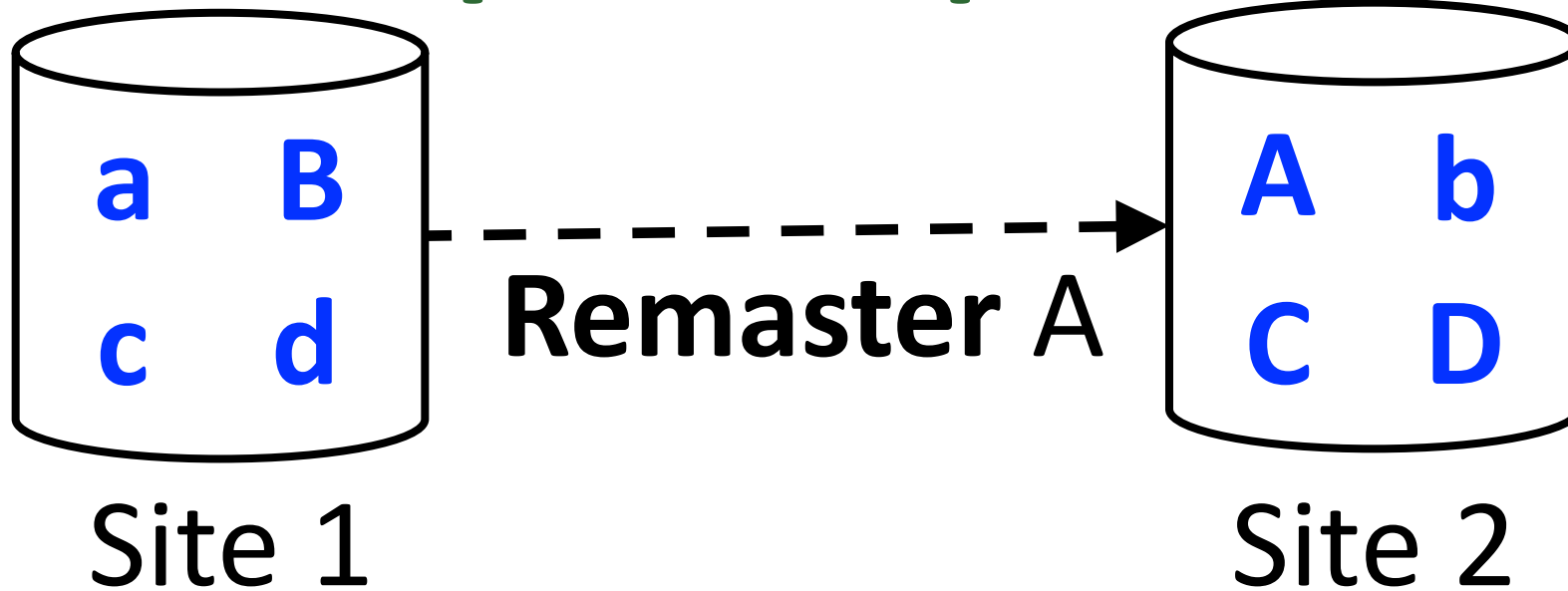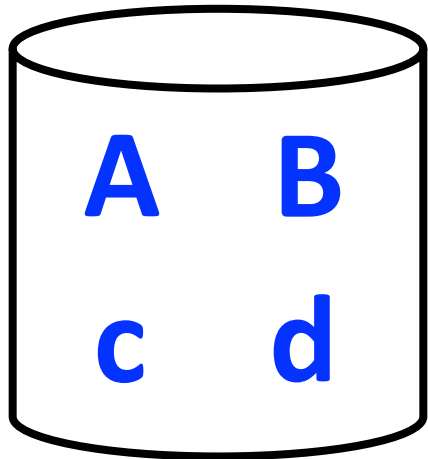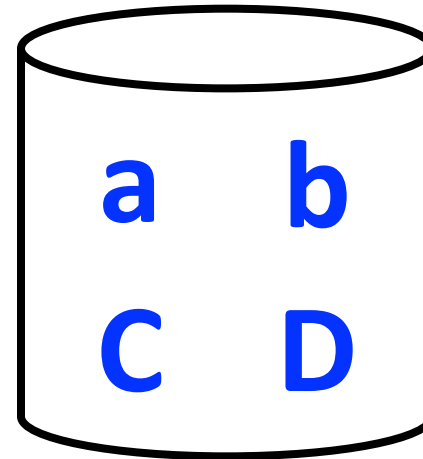**How to perform remastering efficiently?**

**How to decide where to master data?**

# Dynamic Mastering

W[ A, C ]

A    B
c    d

**Remaster** A

a    b
C    D

Site 1

Site 2

# Dynamic Mastering

W[ A, C ]

## Exploit replicas

a   B

c   d

Remaster A

A   b

C   D

Site 1                    Site 2

# Exploiting Replicas

W[ A ]

**A    B**
**c    d**

Site 1

**a    b**
**C    D**

Site 2

# Exploiting Replicas

W[ A ]

W[ A, C ]

A    B

c    d

Site 1

**Remaster** A

a    b

C    D

Site 2

# Exploiting Replicas

W[ A, C ]

**Inconsistent!**

a    B

c    d

A    b

C    D

**Remaster** A

Site 1

Site 2

UNIVERSITY OF
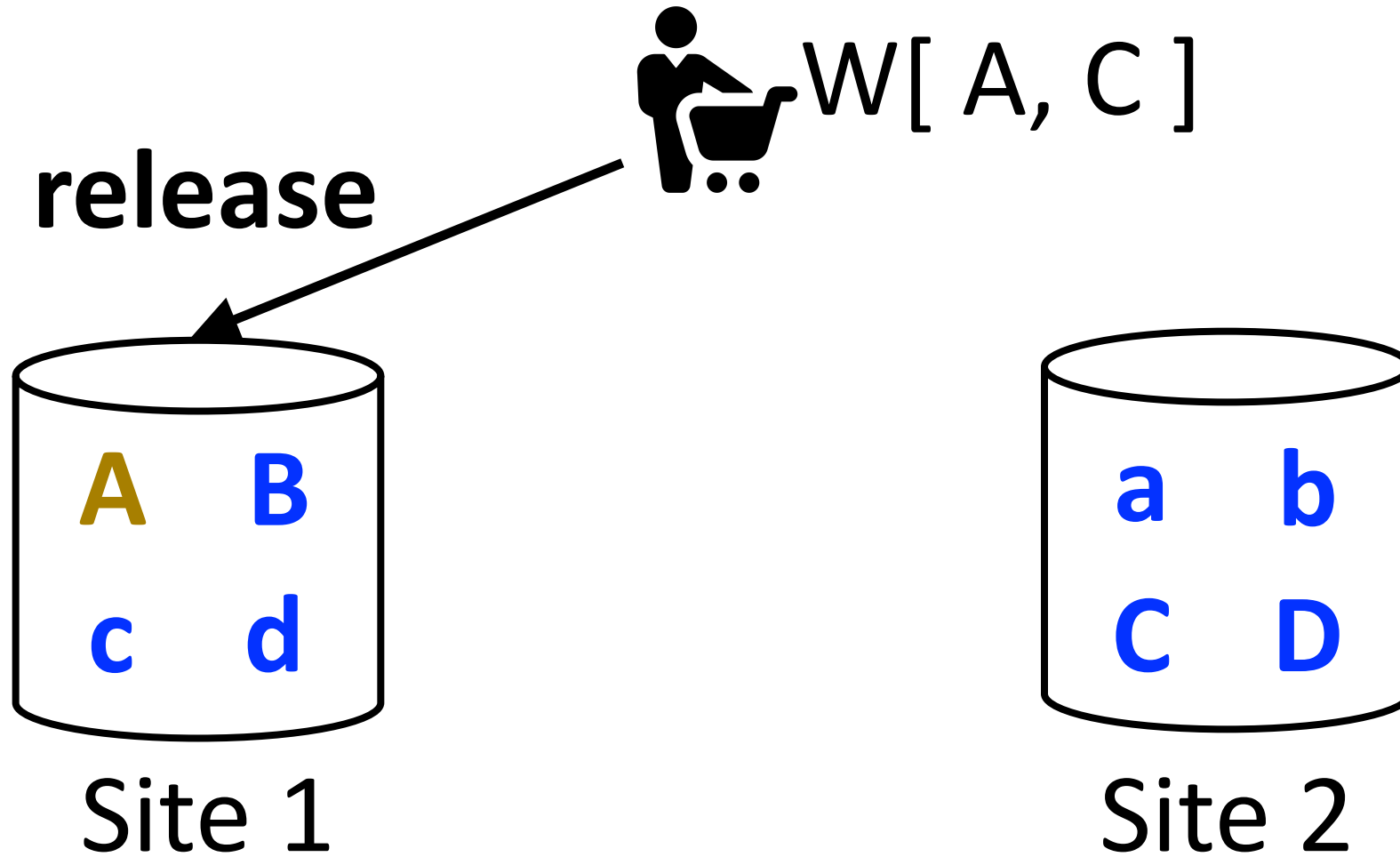WATERLOO

# Ensuring Consistency

Old master must **not allow new updates**

**release** mastership

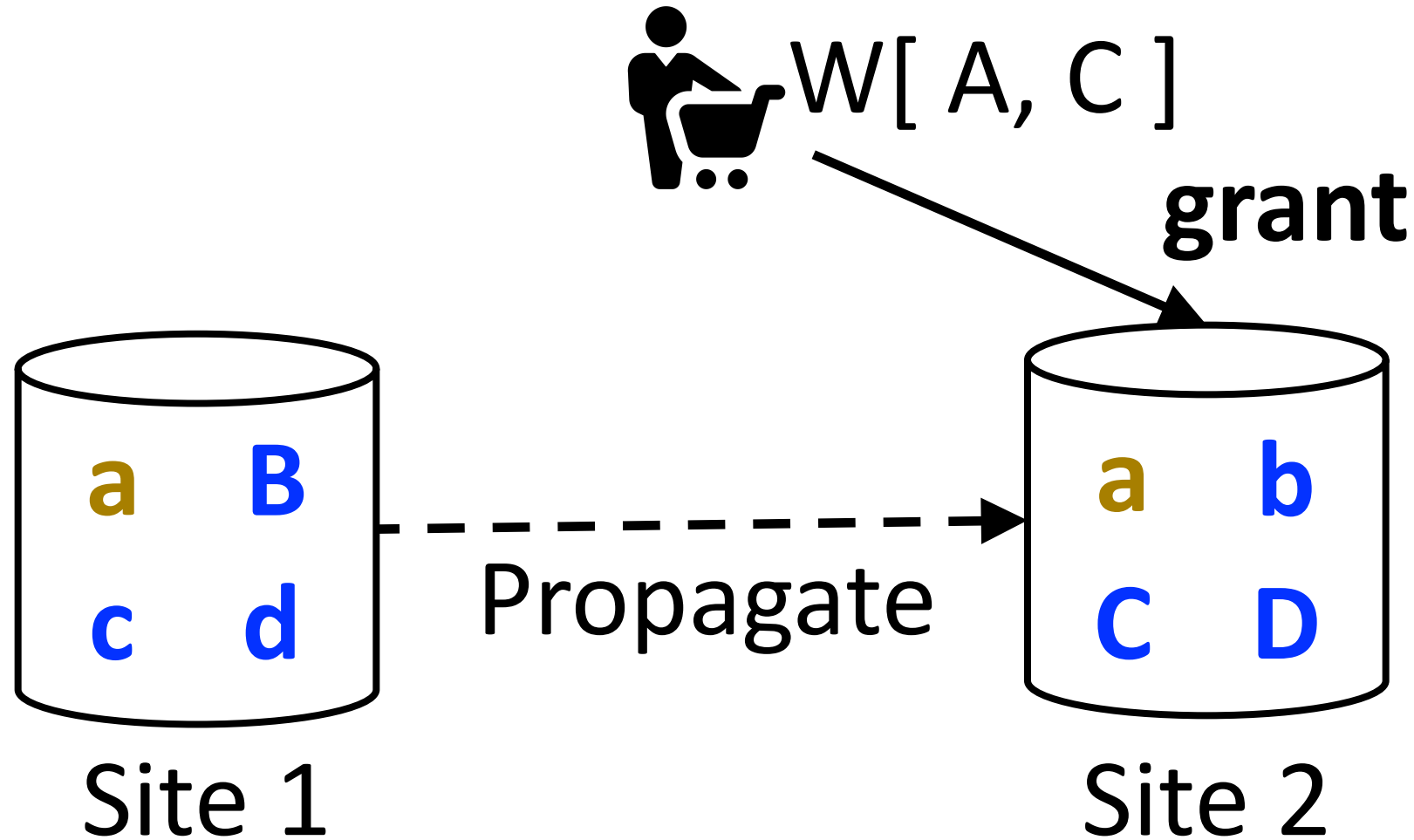New master must have **all previous updates**
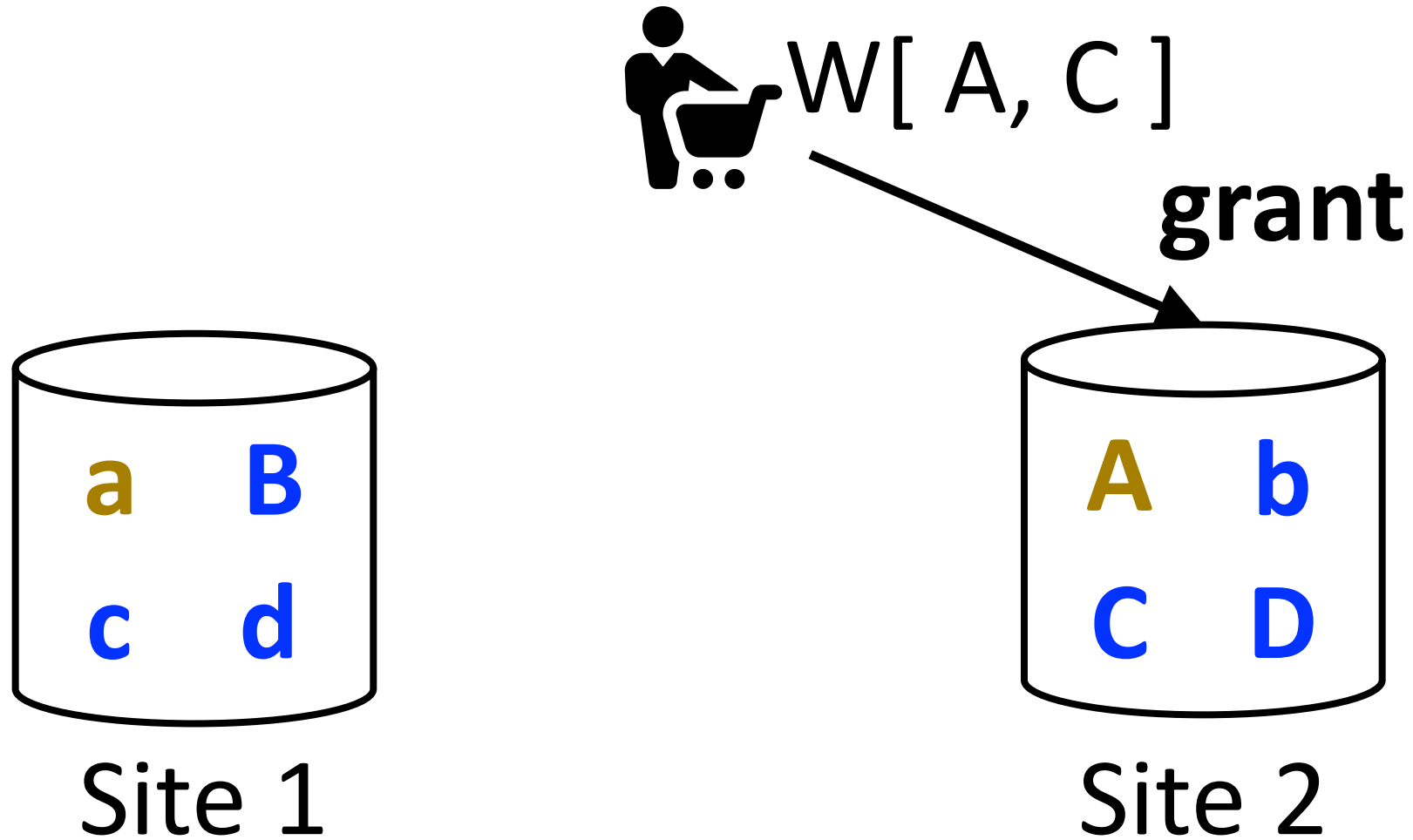
**grant** mastership

# Ensuring Consistency

W[ A, C ]

release

A    B

c    d

Site 1

a    b

C    D

Site 2

# Ensuring Consistency



W[ A, C ]

release

grant

Propagate

a   B

c   d

a   b

C   D

Site 1

Site 2

# Ensuring Consistency

W[ A, C ]

**grant**



Site 1

Propagate

Site 2

a    B
c    d

a    b
C    D

# Ensuring Consistency



W[ A, C ]

grant

a    B

c    d

Site 1

A    b

C    D

Site 2

# Ensuring Consistency

New master must have **all previous updates**

New master was **a lazy replica of old master**

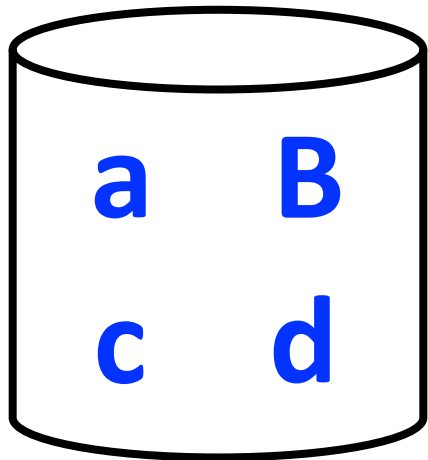Little time spent waiting for updates

**tiny.cc/dynamast**

UNIVERSITY OF
WATERLOO

# Dynamic Mastering Challenges

How to perform remastering efficiently?
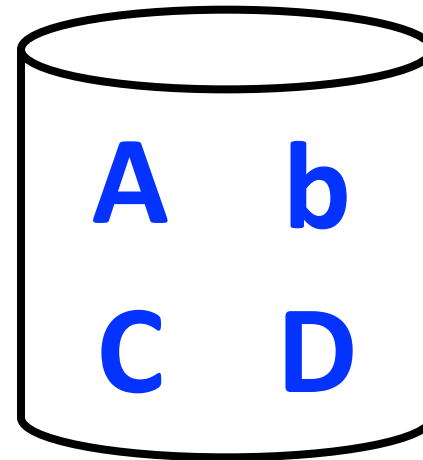
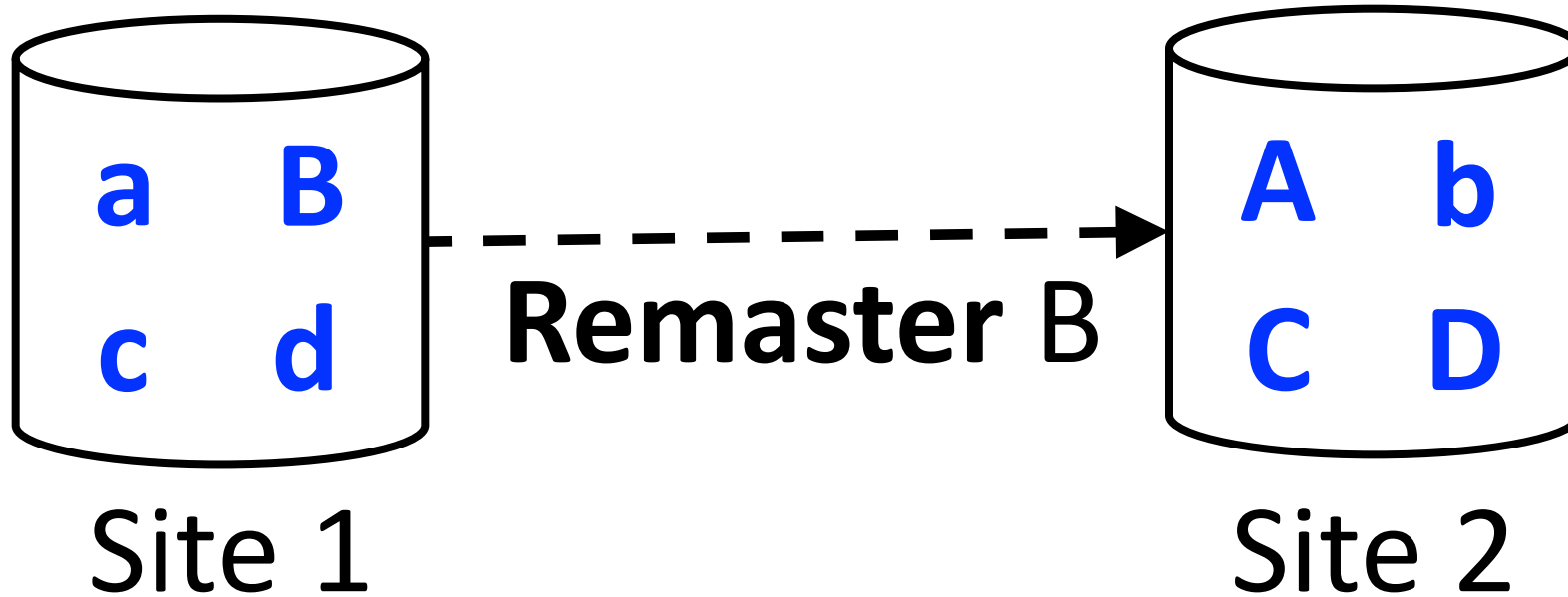How to **decide where to master** data?

# Where to master?

W[ A, C ]

A    B

c    d

Site 1

Remaster A

a    b

C    D

Site 2

# Where to master?

W[ A, C ]

a    B

c    d

Site 1

A    b

C    D

Site 2

# Where to master?

W[ A, B ]

a    B

c    d

**Remaster** B

A    b

C    D

Site 1           Site 2

UNIVERSITY OF
WATERLOO

# Where to master?

W[ A, B ]

**Single-master!**

a   b

c   d

Site 1

A   B

C   D

Site 2

# Where to master?

W[ A, B ]

a    B

c    d

Remaster A

A    b

C    D

Site 1

Site 2

UNIVERSITY OF
WATERLOO

46

# Where to master?

W[ A, B ]

A    B
c    d

**Remaster** A

a    b
C    D

Site 1                    Site 2

# Where to master?

W[ A, C ]

**Ping-pong mastership**

A  B
c  d

Remaster A

a  b
C  D

Site 1

Site 2

# Dynamic Mastering Strategies

**Site Selector** Makes **adaptive** decisions

**Site 1**

A    B
c    d

**Site 2**

a    b
C    D

# Dynamic Mastering Strategies

**Track**: data access patterns, site load

**Quantify** **benefit** of remastering to each site

Load distribution

Update propagation

Future remastering

**Remaster** to site that **maximizes benefit**

**tiny.cc/dynamast**

UNIVERSITY OF
**WATERLOO**

# How well does it work?

**Workloads**

**YCSB**  Scans & Multi-Key Read-Modify-Write

Uniform and Skew Access Patterns
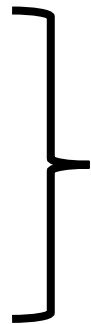
**TPC-C**  Complex updates and reads

UNIVERSITY OF
WATERLOO

# How well does it work?

**Comparisons**

Single-Master

Multi-Master

*Replicated*

Partition-Store

LEAP

*Single Copy*

UNIVERSITY OF
**WATERLOO**

# YCSB with Skew - Throughput

# YCSB with Skew - Routing

# TPC-C – New-Order Latency

# DynaMast Takeaways

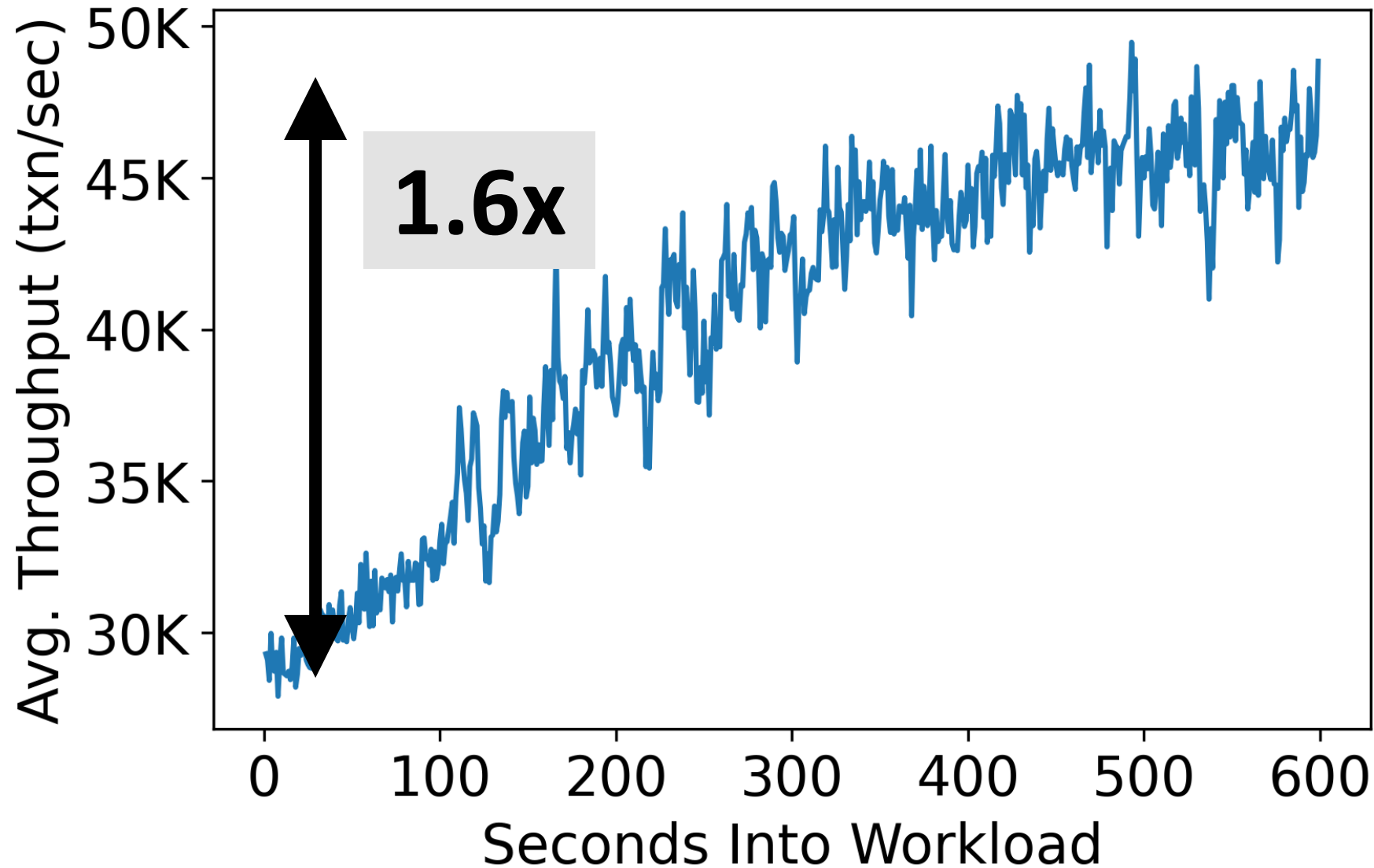**Dynamic mastering** guarantees **single-site transactions**

**tiny.cc/dynamast**

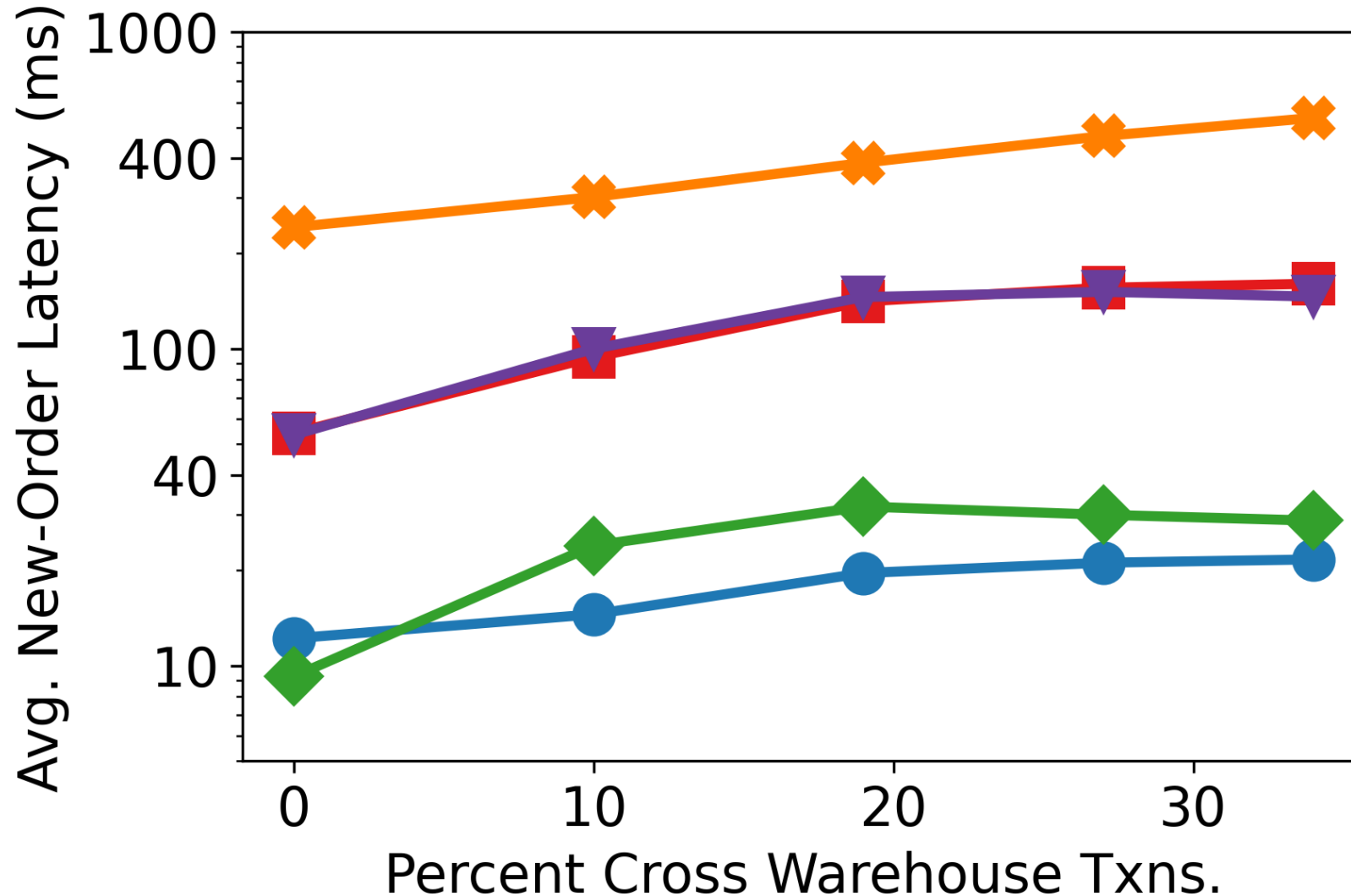**Use replicas** to remaster **efficiently** **outside transaction boundaries**

*Adaptive* site selector strategies **balance load** and **minimizes future remastering**

UNIVERSITY OF
WATERLOO

# Extra slides

# DynaMast Learns Over Time

# TPC-C – New-Order Latency

# Comparisons – Single-Master

Writers

Readers

A    B
C    D

Site 1

a    b
c    d

Site 2

# Comparisons – Multi-Master



W[ A, C ]

**A     B**
**c     d**

Site 1

**prepare**
**commit**

**a     b**
**C     D**

Site 2

UNIVERSITY OF
**WATERLOO**

# Comparisons – Partition-Store



W[ A, C ]

A  B

C  D

prepare

commit

Site 1

Site 2

# Comparisons – LEAP

W[ A, C ]

request C

A   B

Site 1

C   D

Site 2

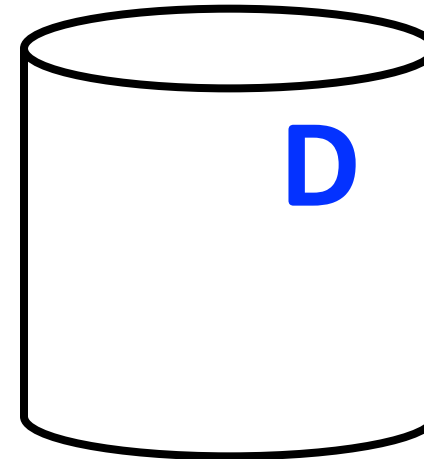UNIVERSITY OF
WATERLOO

# Comparisons – LEAP

W[ A, C ]

transfer C

**Site 1**

A   B

C

**Site 2**

D

# Comparisons – LEAP

W[ A, C ]

A    B

C

Site 1

D

Site 2