# Adaptive
# Data Storage & Placement in Distributed Database Systems
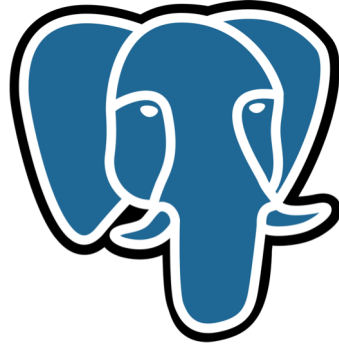
Michael Abebe

mtabebe@uwaterloo.ca

August 2022

UNIVERSITY OF WATERLOO

# Distributed DBMSs are widely used

# Distributed Databases

## How and where to store data?

**?**

**Replication**

**Partitioning**

**Format**

**?**

UNIVERSITY OF
**WATERLOO**

3

# Database Replication

Writers

Readers

Master

A  B
C  D

Replica

a  b
c  d

# Database Replication

Writers

Readers

**Distribution of Read Load**

A
C

b
d

Master

Replica

UNIVERSITY OF
WATERLOO

# Database Replication



Writers

Readers

A    B

C    D

Master

a    b

c    d

Replica

# Database Replication

Writers

Readers

**Performance bottleneck**

A
C

b
d

Master

Replica

UNIVERSITY OF
WATERLOO

# Database Partitioning

# Database Partitioning



Distribution of Load

UNIVERSITY OF WATERLOO

# Database Partitioning



W[ A, C ]

A B

C D

prepare
commit

# Database Partitioning



W[ A, C ]

**Expensive Coordination**

A    D

commit

# Storage Formats

## Row Layout

## Columnar Layout

**Updates**

**Analytics**

# Storage Formats and HTAP

**Row Layout**

**Columnar Layout**

**Updates (OLTP)**

**Analytics (OLAP)**

UNIVERSITY OF
WATERLOO

# Storage Formats and HTAP

**Row Layout**

**Columnar Layout**

**Storage Overheads**

**Updates (OLTP)**

**Analytics (OLAP)**

UNIVERSITY OF
**WATERLOO**

# Distributed Databases

How and where to store data?



Replication

Partitioning

Format

**Trade-offs dependent on workload**

**Distributed DBMSs must adapt**

# Workloads Can Change





**Distributed DBMSs must adapt**

# Thesis Statement

Automatic **adaptation** of **how & where** data stored

**Using online workload** information

Improves performance of distributed DBMSs

# Thesis Contributions

Automatic **adaptation** of **how & where** data stored
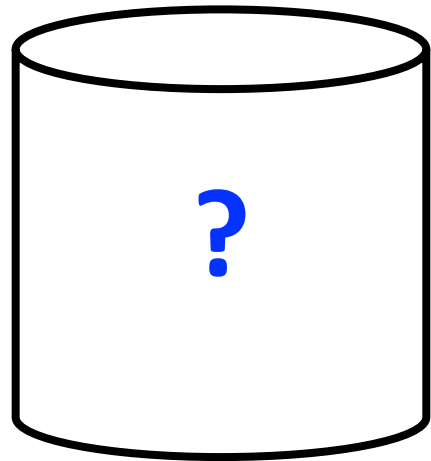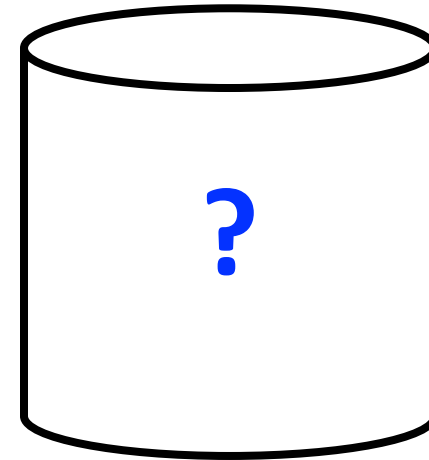
**DynaMast**
(ICDE'20)

Dynamic transfer data mastership to reduces overhead of coordination

**MorphoSys**
(PVLDB'20)

Automatically select physical design: partitioning, & data placement

**Proteus**
(SIGMOD'22)
(PVLDB'22)

Adapt data storage (formats & tiers) for HTAP workloads

UNIVERSITY OF
WATERLOO

# Architecture

# Thesis Contributions

Automatic **adaptation** of **how & where** data stored

**DynaMast**  Dynamic transfer data mastership to reduces overhead of coordination

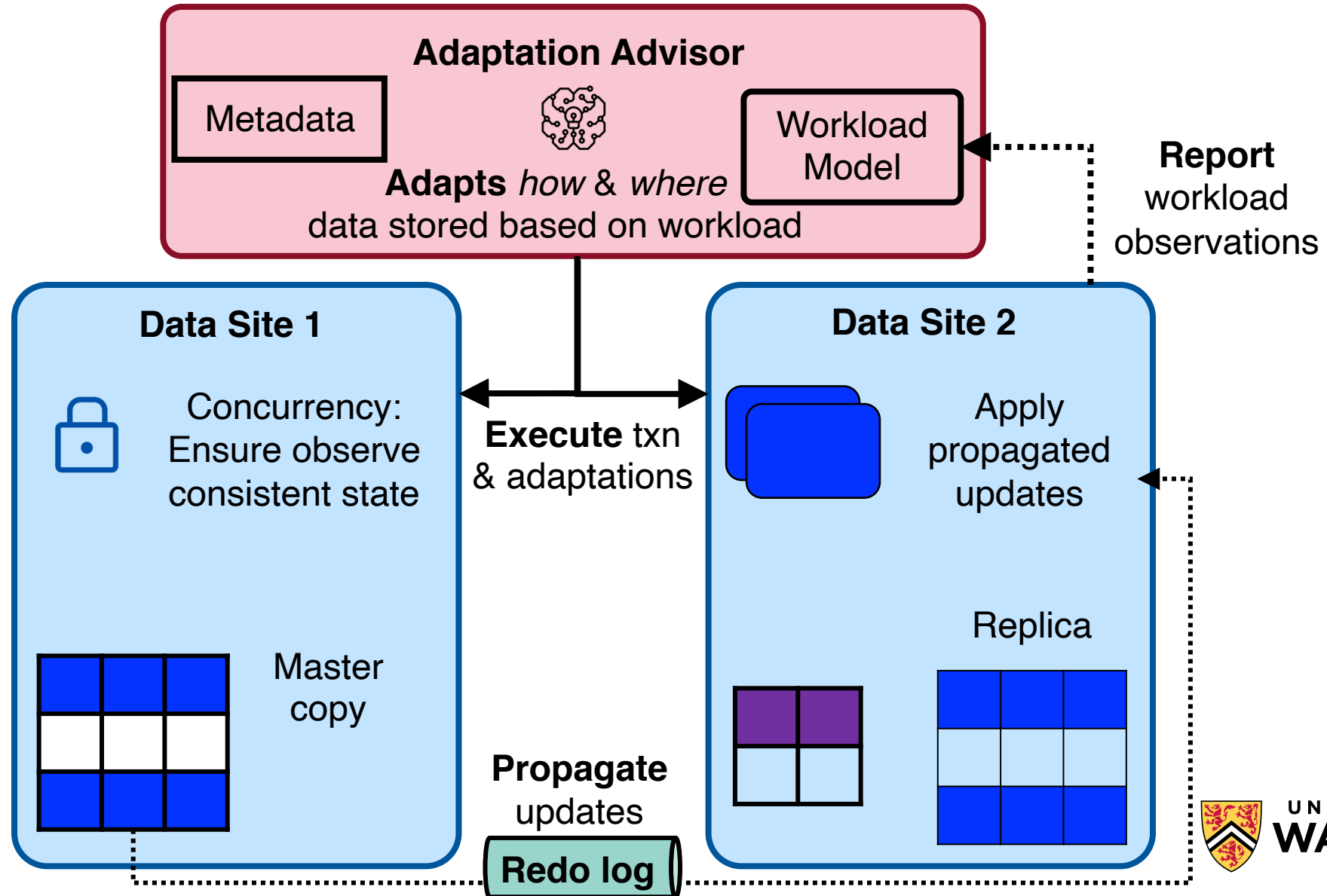**MorphoSys**  Automatically select physical design: partitioning, & data placement

**Proteus**  Adapt data storage (formats & tiers) for HTAP workloads

UNIVERSITY OF
**WATERLOO**

# Dynamic Mastering

W[ A, C ]

A    B
c    d

Site 1

Remaster A

a    b
C    D

Site 2

# Dynamic Mastering

# Dynamic Mastering

**Outside** transaction boundaries

W[ A, C ]

W[ C ]

A    B

c    d

Remaster A

a    b

C    D

Site 1

Site 2

# Dynamic Mastering

W[ B ]

W[ A, C ]

a    B
c    d

Site 1

A    b
C    D

Site 2

# YCSB with Skew - Throughput

# YCSB with Skew - Routing

# Thesis Contributions

Automatic **adaptation** of **how & where** data stored

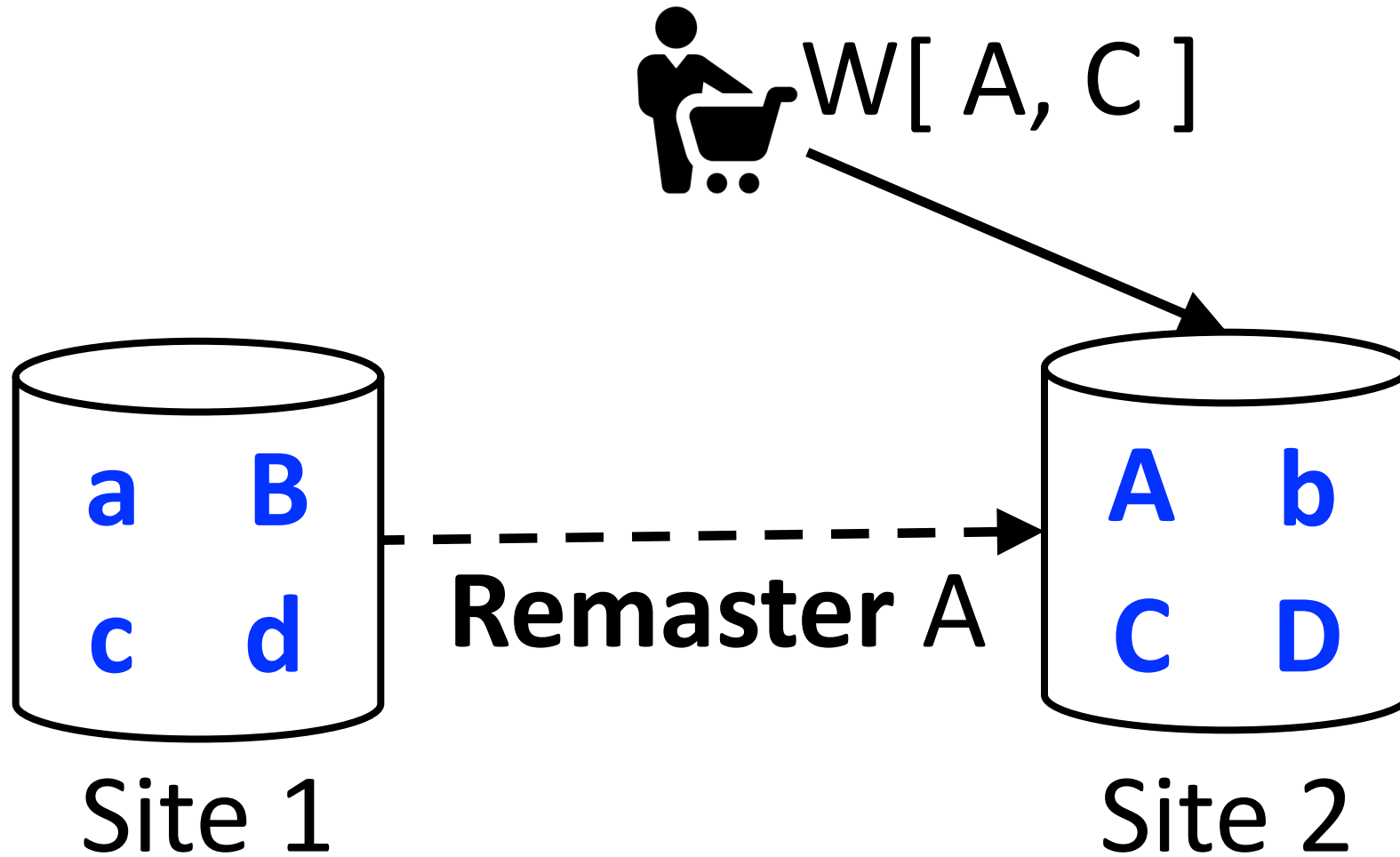**DynaMast**    Dynamic transfer data mastership to reduces overhead of coordination

**MorphoSys**    Automatically select physical design: partitioning, & data placement

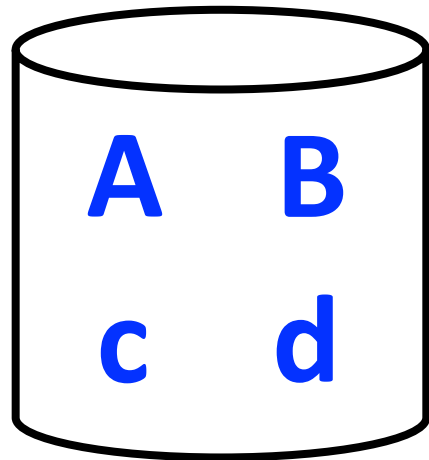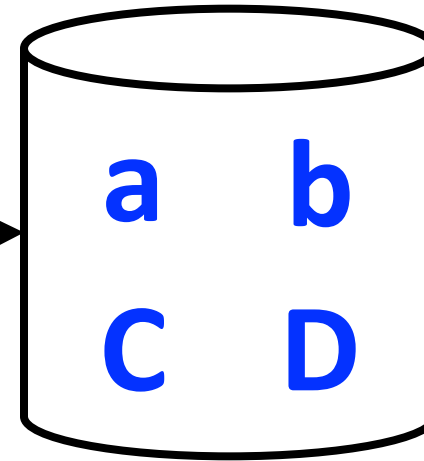**Proteus**    Adapt data storage (formats & tiers) for HTAP workloads

UNIVERSITY OF
WATERLOO

# Distributed DBMS Physical Design

For each **data item**

      Where is the **master**?

      What nodes **replicate** it?

      How is it **grouped (partitioned)** with other data items?

# MorphoSys Physical Design Change Operations

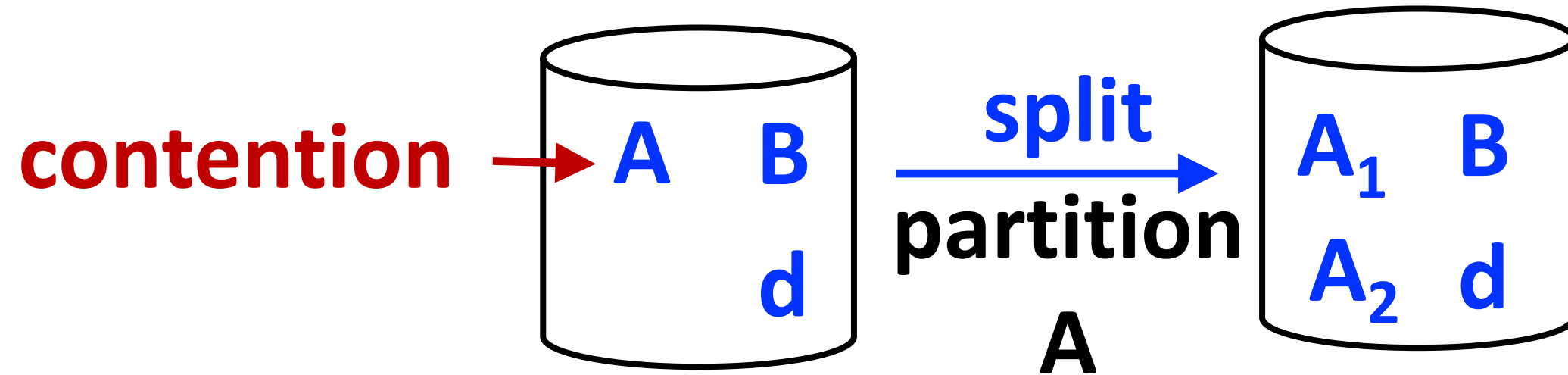**Add** or **remove replica** of a **partition**

**Remaster** a **partition**

**Split** or **merge** **partition(s)**

# Making design decisions

**Learned cost model quantifies**
design change effects



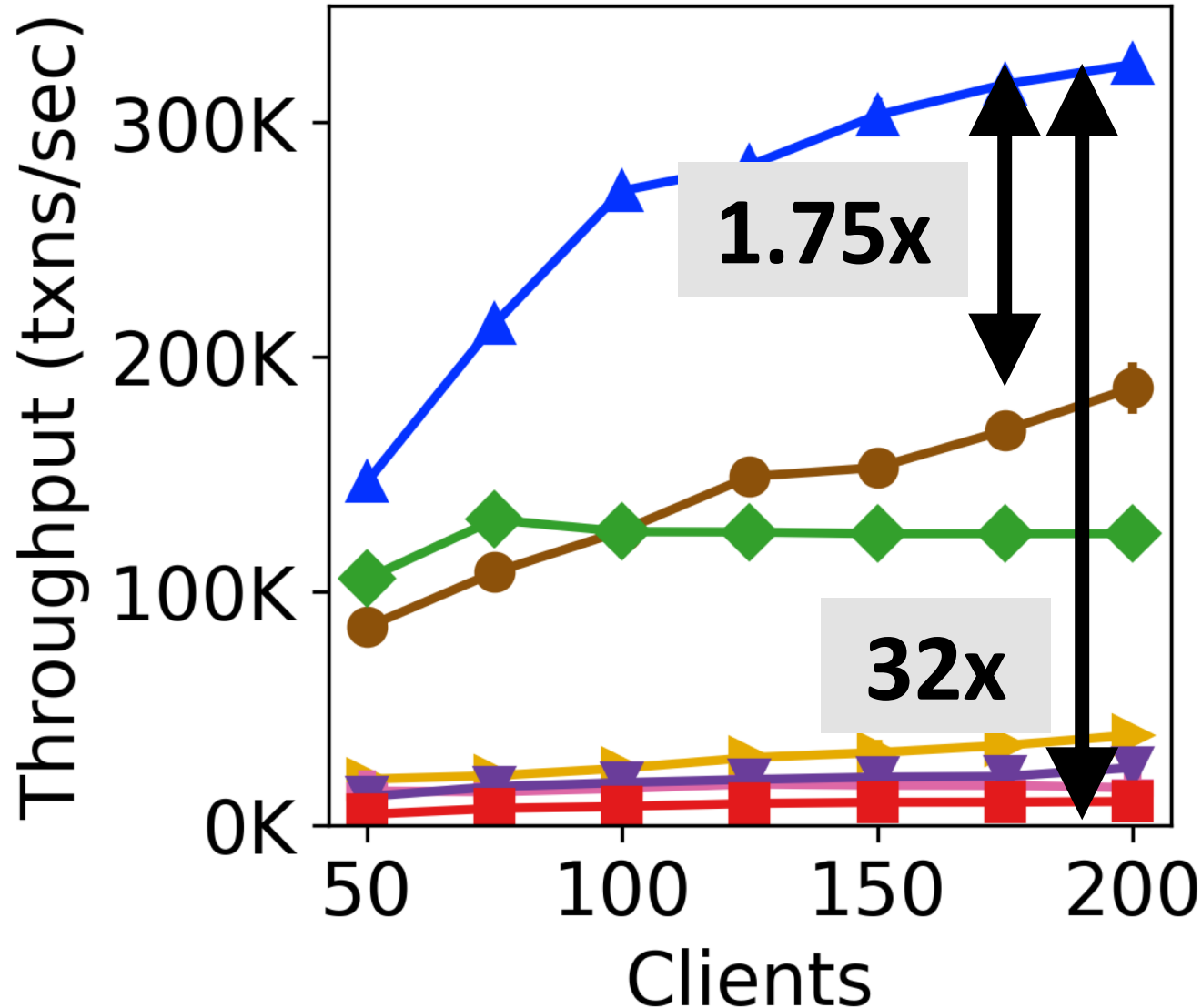**Design change cost** < **Expected Benefit**

# Physical design cost model

**Design change cost** < **Expected benefit**

**Decompose** operators into **key costs**

**Predict** benefit based on workload history

# Skewed YCSB - Throughput



MorphoSys

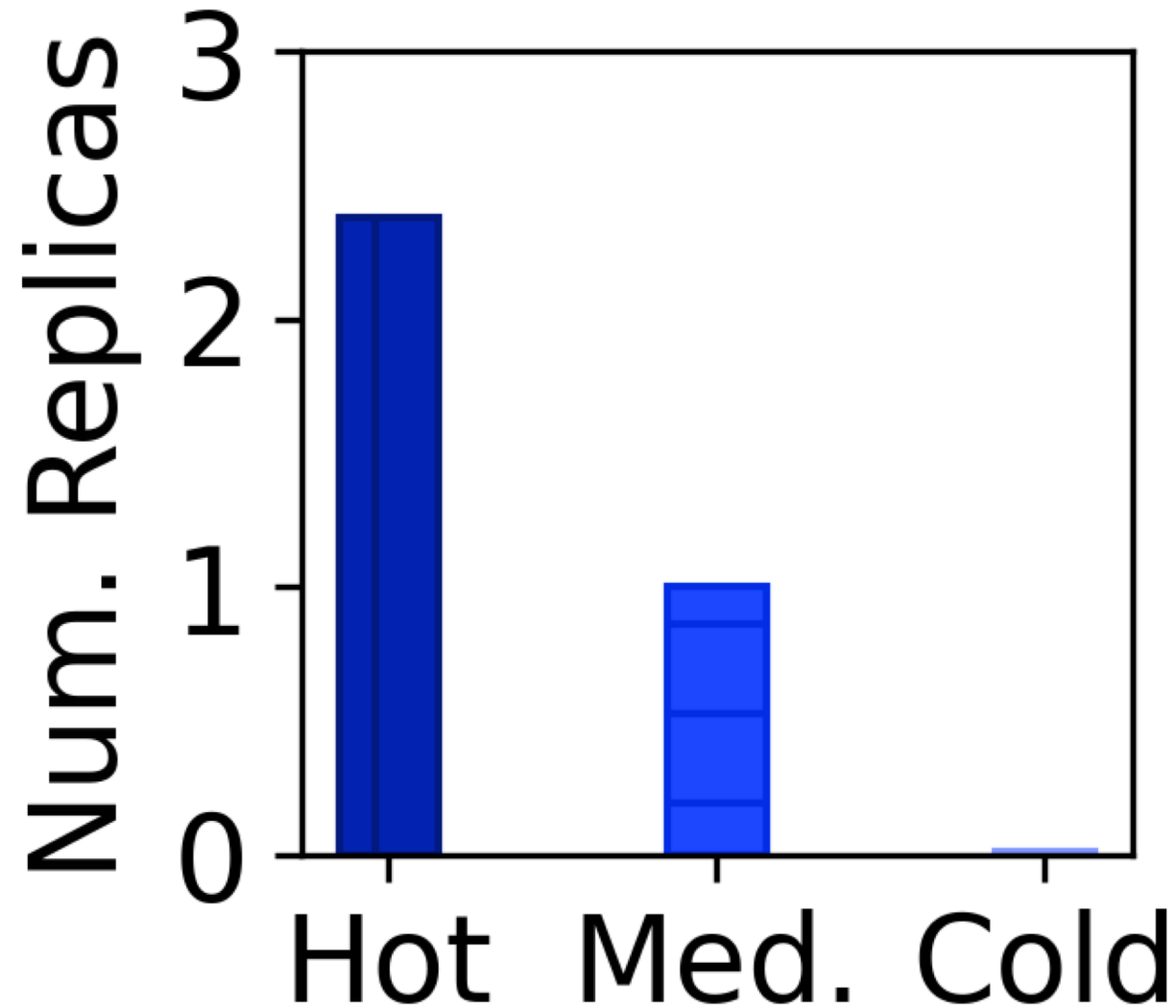DynaMast

Single-Master

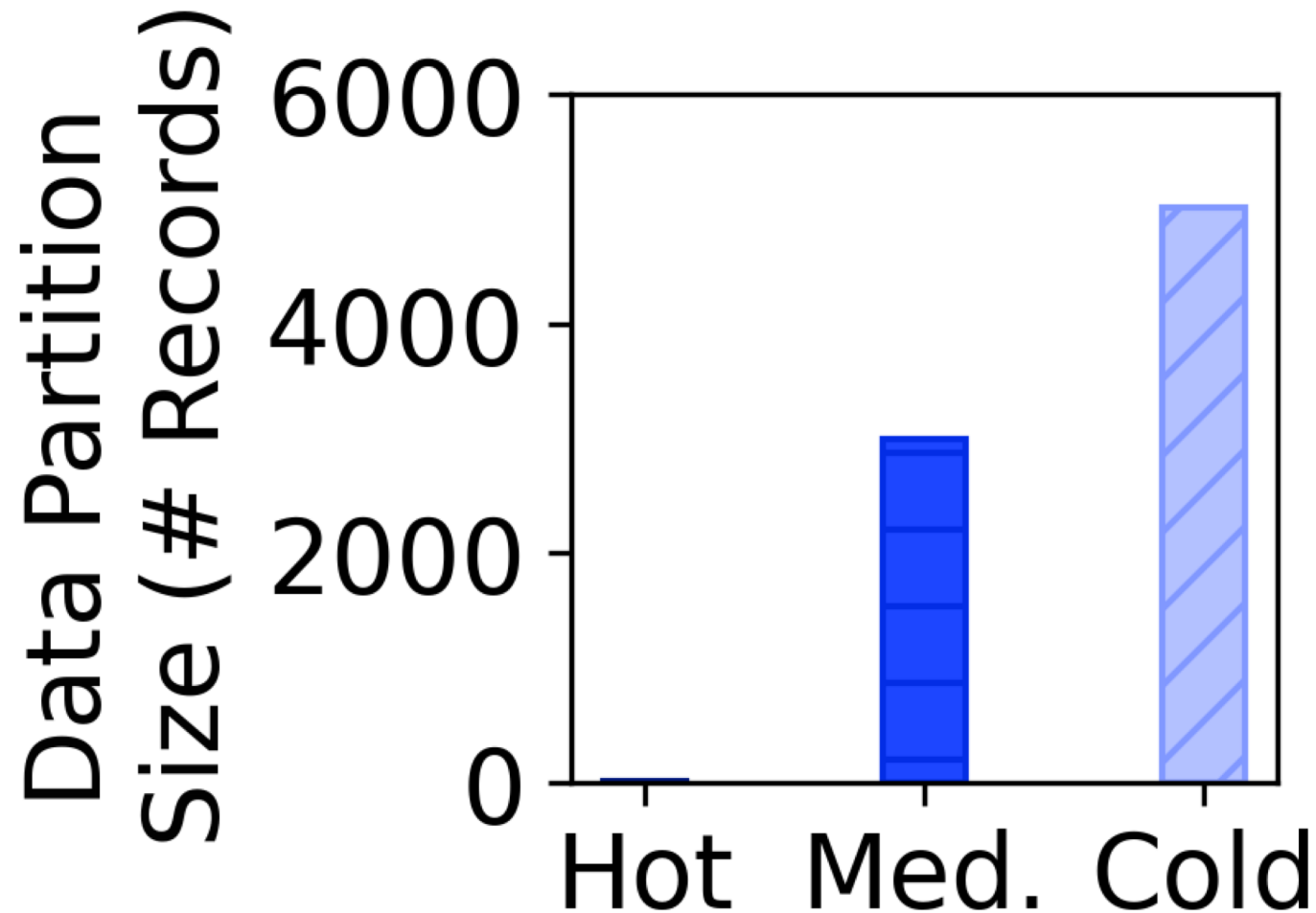Clay   Multi-Master

VoltDB   ADR

# Number of Replicas
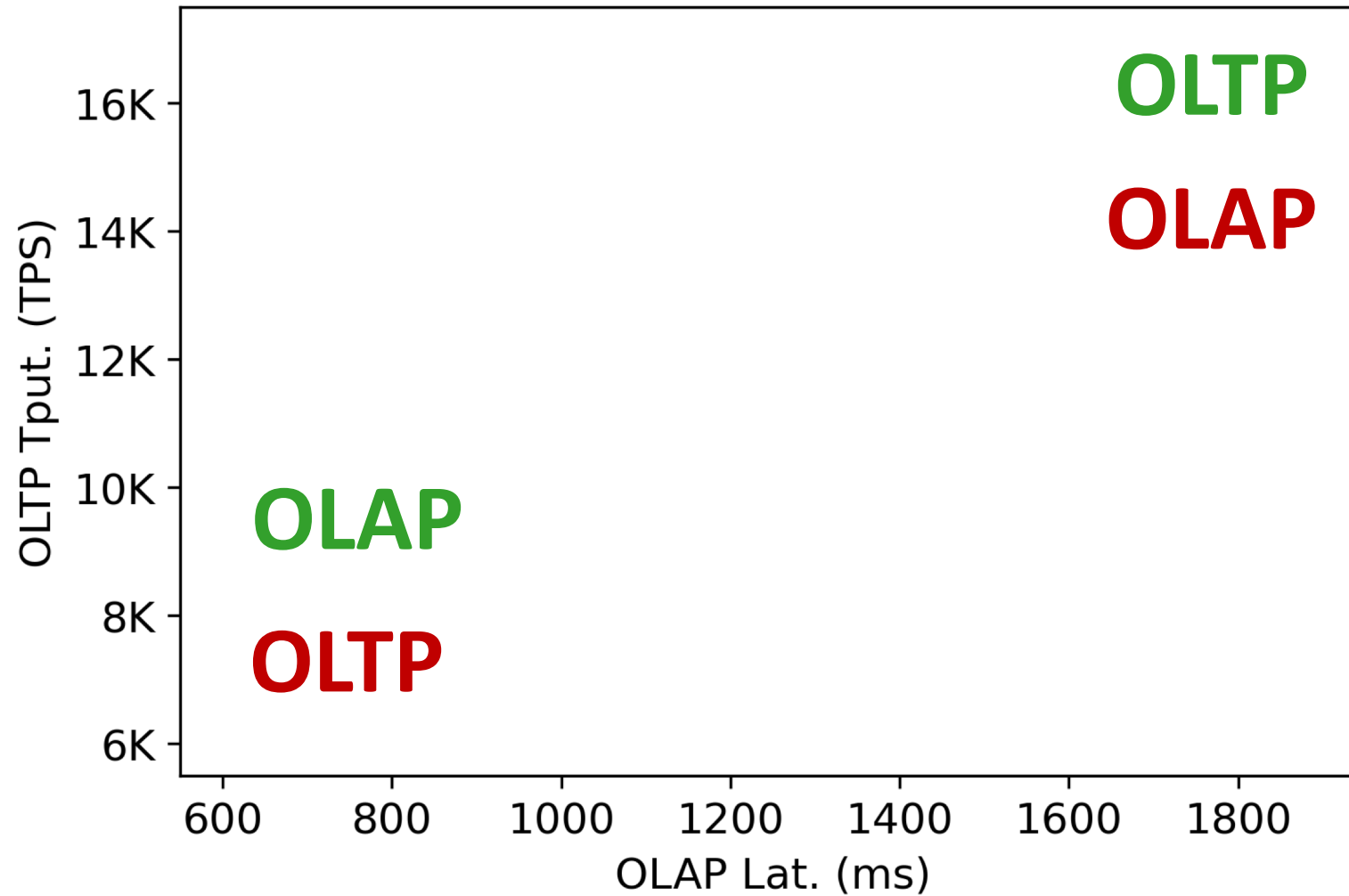
# Partition Sizes

# Thesis Contributions

Automatic **adaptation** of **how & where** data stored

**DynaMast**    Dynamic transfer data mastership to reduces overhead of coordination
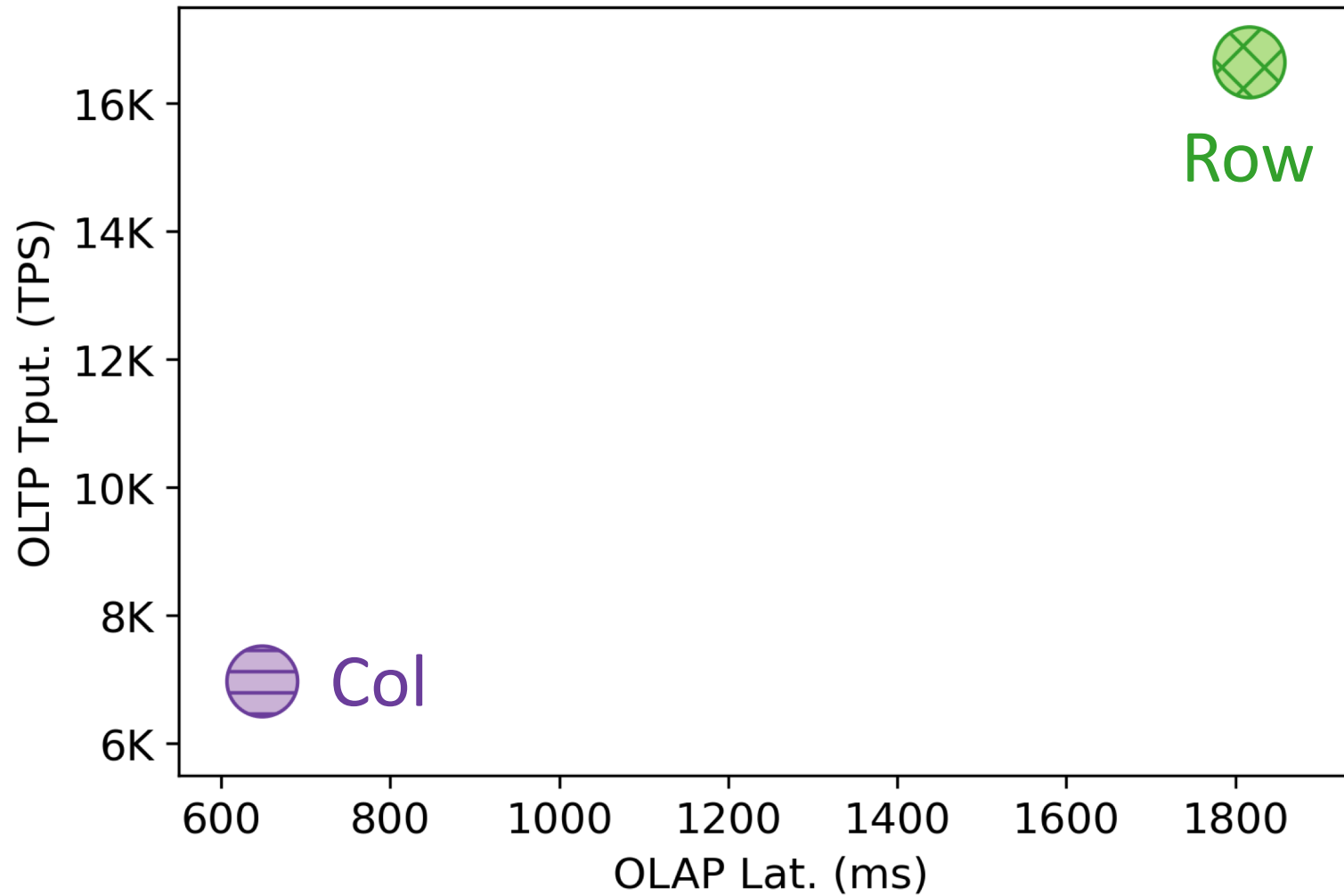
**MorphoSys**    Automatically select physical design: partitioning, & data placement

**Proteus**    Adapt data storage (formats & tiers) for HTAP workloads
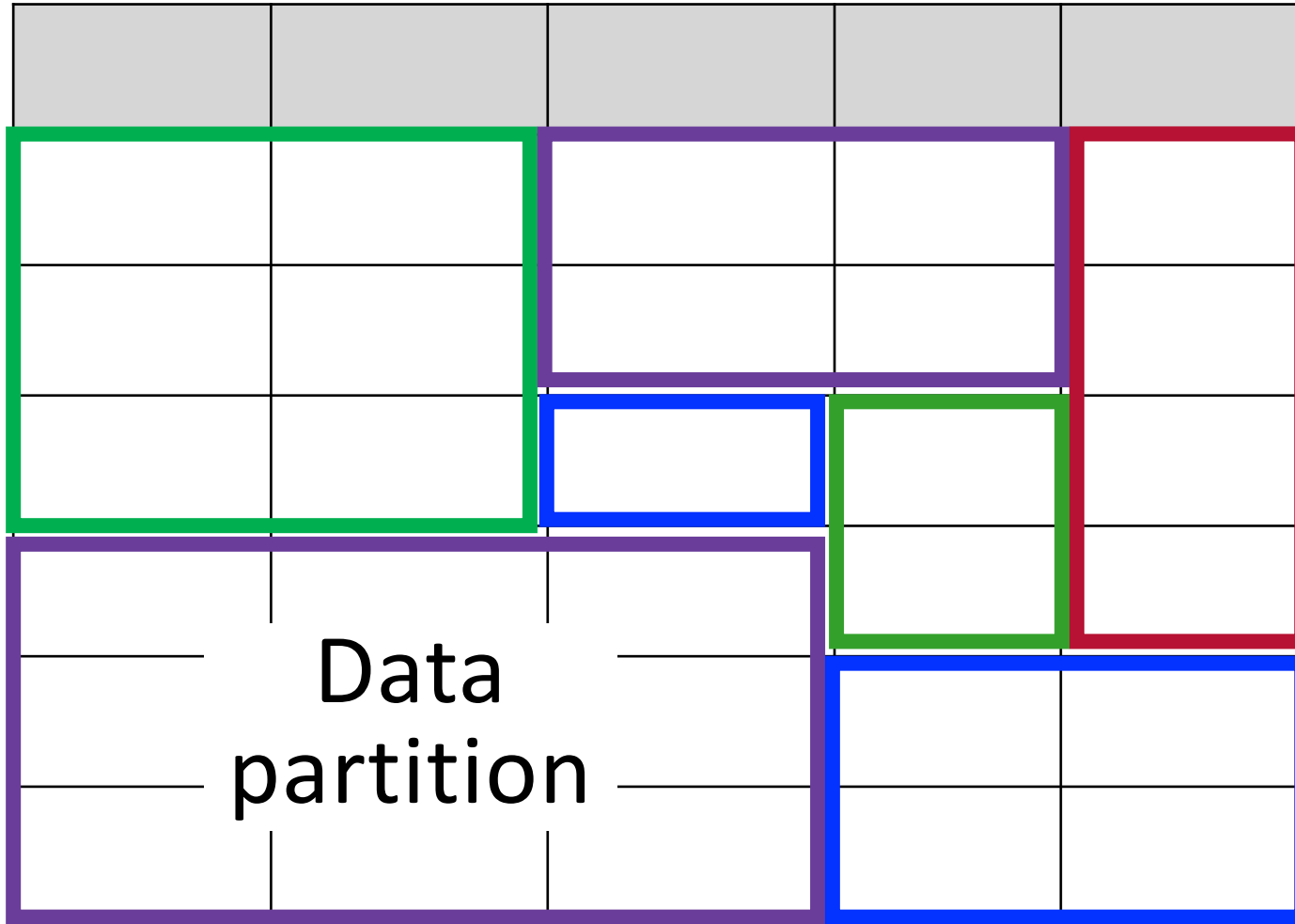
UNIVERSITY OF
WATERLOO

# Performance Trade-Off

# Performance Trade-Off

# Proteus Decisions



**Storage layout**

Master/replica(s)

**Txn execution**

How to **partition?**

**When** & **what** to change

# Transactions in Proteus

**Breakdown** transaction into **physical** operators

**SELECT** book, **SUM**( qnt ) **GROUP BY** book

| **Row layout** | **Logical Plan** | **Sorted column layout** |
|---|---|---|
| Row scan P1 | Scan & Project book, qnt | Sequential col scan P1 |
| ↓ | ↓ | ↓ |
| Hash aggregation book, sum( qnt ) | Aggregate book, sum( qnt ) | Sorted col aggregation book, sum( qnt ) |

UNIVERSITY OF
WATERLOO

# Storage-Aware Operators

Per layout implementation of operators

Operate **directly** over columnar, sorted, compressed data

**Predict physical operator latency**

*Cardinality*

*Data Width*     →      →  **Predicted Latency**

*Est Selectivity*    **Seq col scan**
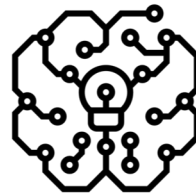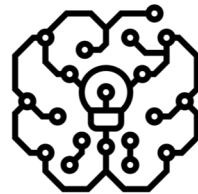
# Storage-Aware Operators

Per layout implementation of operators

Operate **directly** over columnar, sorted, compressed data

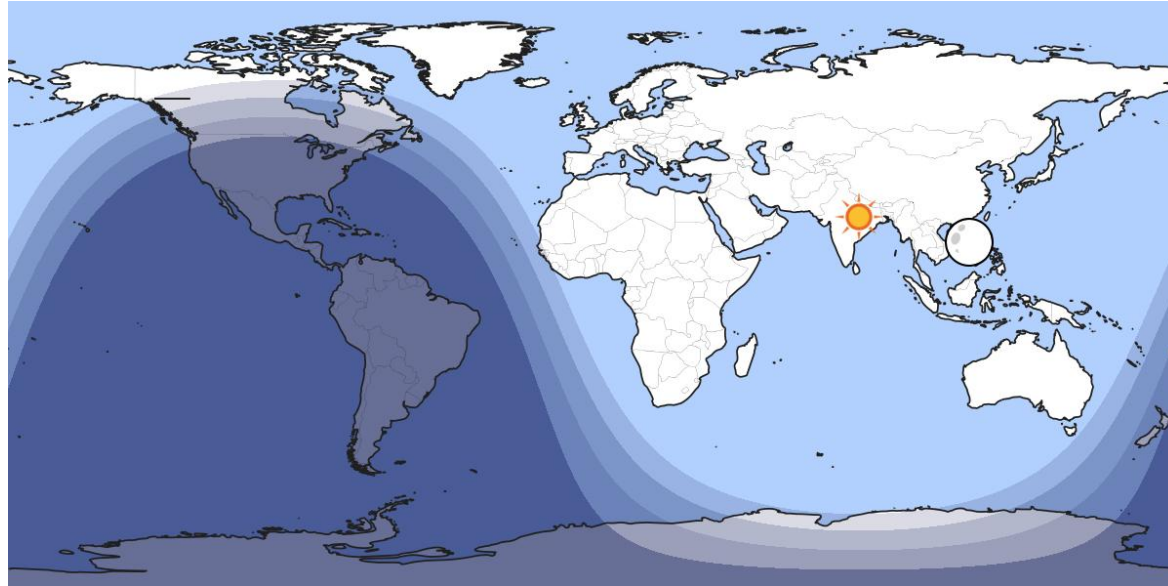**Predict physical operator latency**

*Cardinality*

*Data Width* ⟶  ⟶ **Predicted Latency**
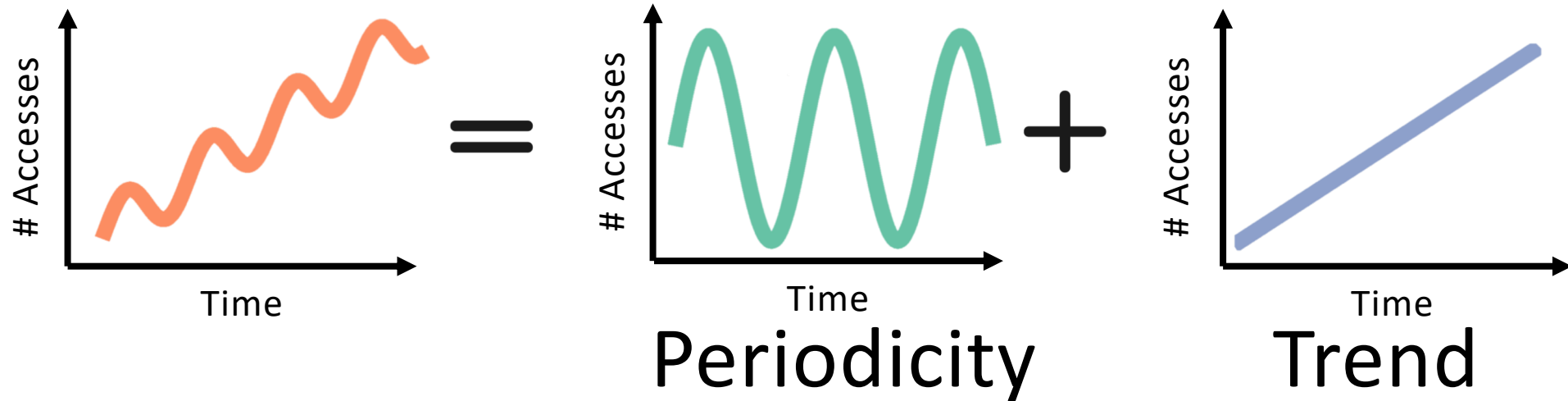
*Est Selectivity*     **Row scan**

# Likelihood of a Transaction

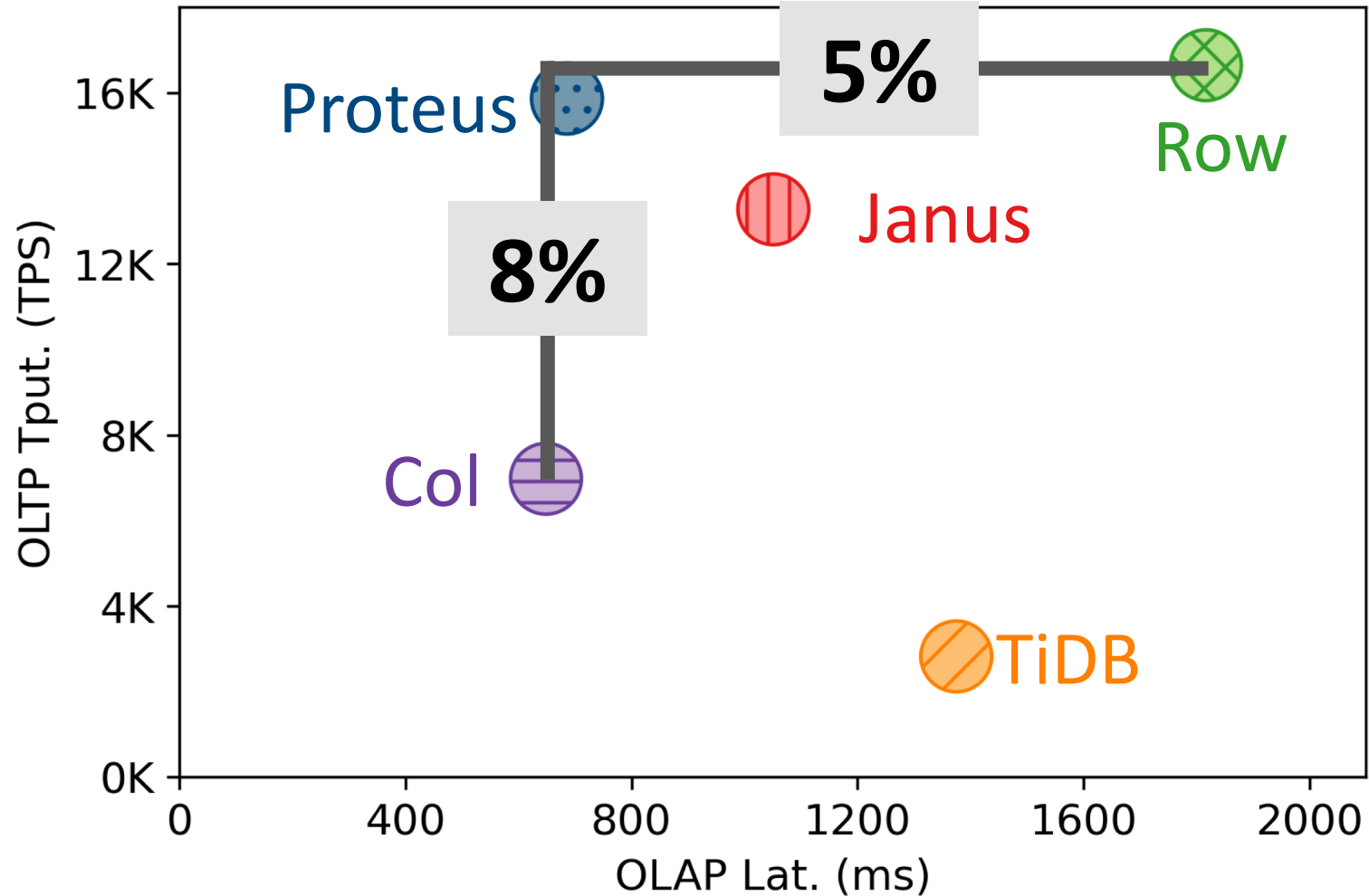Data **accesses** to **storage**
often follow **predictable** pattern

# Likelihood of a Transaction

Data **accesses** to **storage**
often follow **predictable** pattern



Periodicity + Trend

# CH BenCHmark

# Distributed DBMSs are widely used

# Distributed DBMSs must adapt

**Adaptation** of **how** & **where** data stored **improves performance**

**DynaMast**     **MorphoSys**     **Proteus**

UNIVERSITY OF
WATERLOO

# Extra Slides