# Bridging the Archipelago between Row-Stores and Column-Stores for Hybrid Workloads
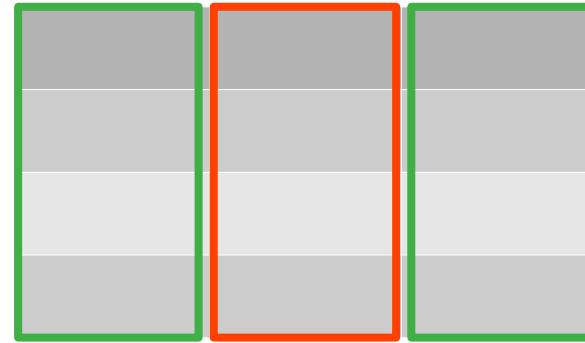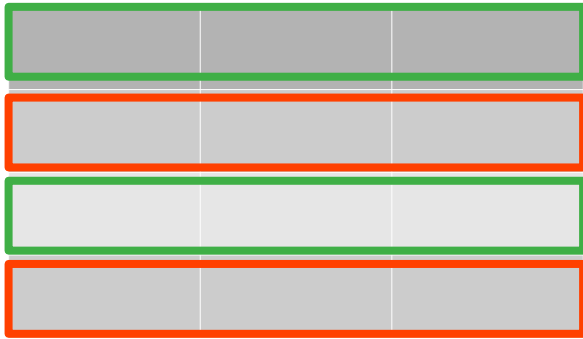
J Arulraj, et el.

SIGMOD '20
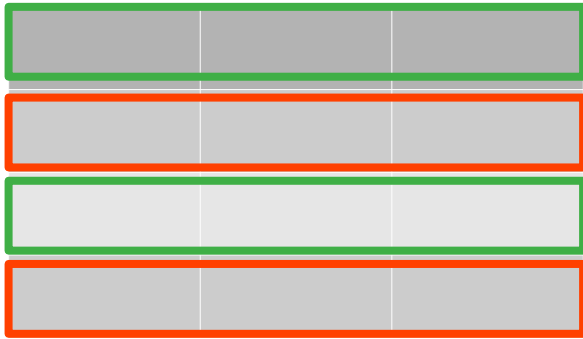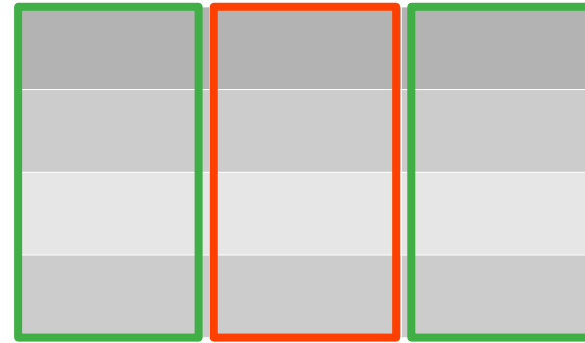
Speaker:

# Introduction

# Physical Memory layout

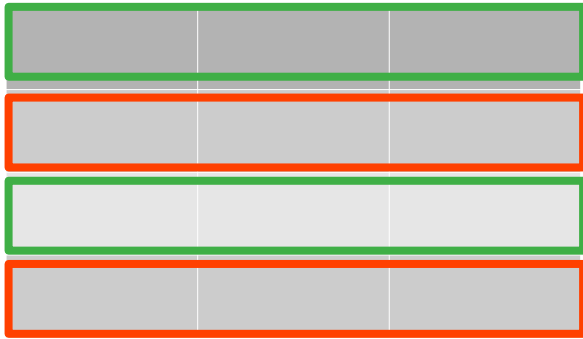# Physical Memory layout



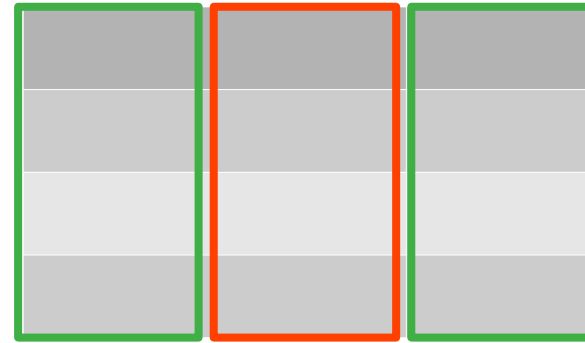?                    ?
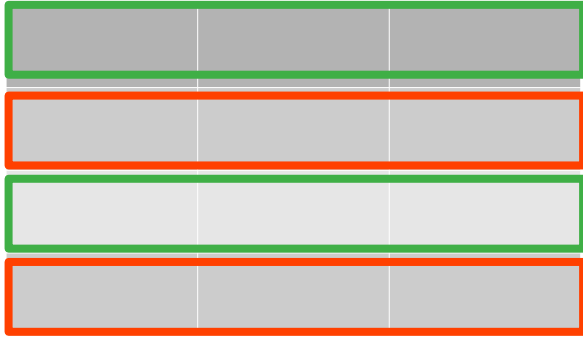
# Physical Memory layout



Row Storage

tuple-centric

Column Storage

Columnar

# Physical Memory layout



Row Storage

Column Storage

When?

# Physical Memory layout



Row Storage

Column Storage

When?     Insert, Update     Search for an attribute

# Physical Memory layout

Row Storage

Column Storage

When?   Insert, Update   Search for an attribute

Hotly updated   Coldly stored

# Physical Memory layout

Row Storage

Column Storage

When?    Insert, Update    Search for an attribute

# Physical Memory layout



Row Storage

Column Storage

When?     Insert, Update          Search for an attribute

Transactional          Analytical

# Physical Memory layout



Row Storage                Column Storage

When?        Insert, Update        Search for an attribute

OLTP                      OLAP

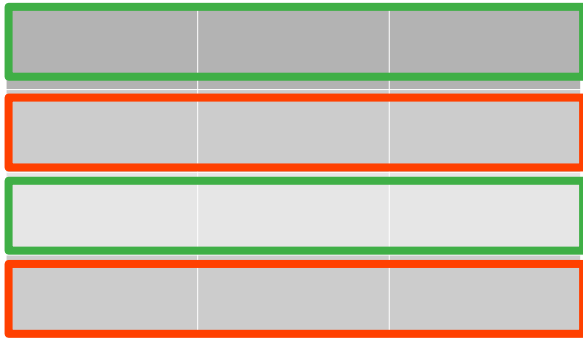# Physical Memory layout



Row Storage

When?    Insert, Update    Search for an attribute

OLTP

# Physical Memory layout



Row Storage

When?    Insert, Update

OLTP

Search for **(a1,a2)** attribute

# Physical Memory layout



Row Storage

When?    Insert, Update          Search for **(a1,a2)** attribute

OLTP                             HTAP

# Physical Memory layout



Row Storage

??

When?    Insert, Update    Search for **(a1,a2)** attribute

OLTP    HTAP

# Physical Memory layout

# Physical Memory layout

# Physical Memory layout

OLTP →  HTAP  ← OLAP

# Physical Memory layout

HTAP

# Physical Memory layout

## HTAP

# Physical Memory layout

## HTAP

# Which Memory layout

# Which?

# Which Memory layout

# Which?

# Do we have any observation?

# Which Memory layout

# Which?

# Do we have any observation?

# Which Memory layout

Hot data → row
Cold data → column

Which Memory layout

Record Query types?

# Which Memory layout

# Record Query types?

If we know what types of queries there are, then we definitely can make a good design.

# Which Memory layout

# Record Query types?

If we know what types of queries there are, then we definitely can make a good design.

But ...

# Which Memory layout

# Record Query types?

If we know what types of queries there are, then we definitely can make a good design.

But ...

Can we be smarter?

# Which Memory layout

Record Query types?

Self-adaptive?

# Self-adaptive algorithm



# Record

# Self-adaptive algorithm

Record

Analysis

# Self-adaptive algorithm

Query        Record        Analysis        ??

# Self-adaptive algorithm

Query

Record

Analysis

Find
Representatives

X = Centroid

# On-line k-means algorithm

Find Representatives



X = Centroid

# On-line k-means algorithm

Find
Representatives



X = Centroid

1st step:

Using recent n Queries,

Find k Representatives R

2nd step:

Generate vertical partitioned layout using R with greedy algorithm. (Largest cluster first)

# On-line k-means algorithm

## FSM (Flexible Storage Model)

Find Representatives



X = Centroid

Different tables, different vertical layouts

# On-line k-means algorithm

Find
Representatives

FSM (Flexible Storage Model)



Different layouts → different access methods?

# On-line k-means algorithm

Find
Representatives



X = Centroid

## FSM (Flexible Storage Model)



Different layouts → different access methods?

Inefficient !!

# On-line k-means algorithm

Find Representatives



FSM (Flexible Storage Model)



Provide Abstract layer !!

# On-line k-means algorithm

Find Representatives

FSM (Flexible Storage Model)

Provide Abstract layer !!

One mutual access interface.

# Abstraction

```
SELECT R.c, SUM(S.z)
FROM R JOIN S ON R.b = S.y
WHERE R.a = 1 AND S.x = 2
GROUP BY R.c;
```

# Abstraction



SELECT R.c, SUM(S.z)
FROM R JOIN S ON R.b = S.y
WHERE R.a = 1 AND S.x = 2
GROUP BY R.c;

$\Omega$ — Materialize, { LT }, { PT }

$\Gamma_{c;sum(z)}$ — Aggregate, { LT, C}, { LT }

$\bowtie_{\mathcal{R}.b=\mathcal{S}.y}$ — Join, { LT, LT }, { LT }

$\pi_{b,c}$   $\pi_{y,z}$ — Projection, { LT }, { LT }

$\sigma_{\mathcal{R}.a=1}$   $\sigma_{\mathcal{S}.x=2}$ — Sequential Scan, { T, P}, { LT }

$\mathcal{R}$   $\mathcal{S}$ — Table

LT : logical tile

PT : physical tile

T : table

Attributes : C

Predicate : P

# Abstraction



SELECT R.c, SUM(S.z)
FROM R JOIN S ON R.b = S.y
WHERE R.a = 1 AND S.x = 2
GROUP BY R.c;

$\Omega$ — Materialize, { LT }, { PT }

$\Gamma_{c;sum(z)}$ — Aggregate, { LT, C}, { LT }

$\bowtie_{\mathcal{R}.b=\mathcal{S}.y}$ — Join, { LT, LT }, { LT }

$\pi_{b,c}$ $\pi_{y,z}$ — Projection, { LT }, { LT }

$\sigma_{\mathcal{R}.a=1}$ $\sigma_{\mathcal{S}.x=2}$ — Sequential Scan, { T, P}, { LT }

$\mathcal{R}$ $\mathcal{S}$ — Table

LT : logical tile

PT : physical tile

T : table

Attributes : C

Predicate : P

Physical data

# Abstraction



```
SELECT R.c, SUM(S.z)
FROM R JOIN S ON R.b = S.y
WHERE R.a = 1 AND S.x = 2
GROUP BY R.c;
```

$\Omega$

$\Gamma_{c;sum(z)}$

$\bowtie_{\mathcal{R}.b=\mathcal{S}.y}$

$\pi_{b,c}$        $\pi_{y,z}$

$\sigma_{\mathcal{R}.a=1}$        $\sigma_{\mathcal{S}.x=2}$

$\mathcal{R}$        $\mathcal{S}$

Materialize, { LT }, { PT }

Aggregate, { LT, C}, { LT }

Join, { LT, LT }, { LT }

Projection, { LT }, { LT }

Sequential Scan, { T, P}, { LT }

Table

LT : logic tile

PT : physical tile

T : table

Attributes : C

Predicate : P

logical data +

Flexible materialization

# Tiles



**Physical Tile**

**Logical Tile**

# Parallelism

Definitely we do not want stalling.

# Parallelism

Definitely we do not want stalling.

Write while reconfigure memory layout

# Parallelism

Definitely we do not want stalling.

Write while reconfigure memory layout

$$\rightarrow \text{MVCC}$$

(Multi-Version Concurreny Control)

# MVCC



Metadata for MVCC

# MVCC



Metadata for MVCC

# MVCC

Unique transaction identifier



**Running transaction**

| Index | TxnId | BeginCTS | EndCTS | PreV | | ID | SALARY |
|-------|-------|----------|--------|------|---|-----|--------|
| 101 | — | 1001 | 1004 | — | | 101 | 201 |
| 102 | — | 1001 | ∞ | — | | 102 | 202 |
| 103 | 305 | 1002 | ∞ | — | | 103 | 203 |
| 104 | — | 1002 | ∞ | — | | 104 | 204 |
| 105 | — | 1002 | 1003 | — | | 105 | 205 |
| | — | 1004 | ∞ | [1, 1] | | 101 | 301 |

Metadata for MVCC

# MVCC

Unique commit timestamp



| Index |
|-------|
| 101 |
| 102 |
| 103 |
| 104 |
| 105 |

| TxnId | BeginCTS | EndCTS | PreV |
|-------|----------|--------|------|
| — | 1001 | 1004 | — |
| — | 1001 | ∞ | — |
| 305 | 1002 | ∞ | — |

| TxnId | BeginCTS | EndCTS | PreV |
|-------|----------|--------|------|
| — | 1002 | ∞ | — |
| — | 1002 | 1003 | — |
| — | 1004 | ∞ | [1, 1] |

| ID | SALARY |
|----|--------|
| Tile A-1 101 | Tile A-2 201 |
| 102 | 202 |
| 103 | 203 |

| ID | SALARY |
|----|--------|
| Tile B-1 104 | 204 |
| 105 | 205 |
| 101 | 301 |

Metadata for MVCC

# MVCC

## i) Update (atomic)



| Index | | TxnId | BeginCTS | EndCTS | PreV | | ID | SALARY |
|-------|---|-------|----------|--------|------|---|-----|--------|
| 101 | | — | 1001 | 1004 | — | Tile A-1 / 101 Tile A-2 / 201 | | |
| 102 | | — | 1001 | ∞ | — | 102 / 202 | | |
| 103 | | 305 | 1002 | ∞ | — | 103 / 203 | | |
| 104 | | — | 1002 | ∞ | — | Tile B-1 / 104 | 104 | 204 |
| 105 | | — | 1002 | 1003 | — | 105 | 105 | 205 |
| | | — | 1004 | ∞ | [1, 1] | 101 | 101 | 301 |

Metadata for MVCC

# MVCC

## i) Update  (after commit)



| Index | | TxnId | BeginCTS | EndCTS | PreV | | ID | SALARY |
|---|---|---|---|---|---|---|---|---|
| 101 | | — | 1001 | 1004 | — | Tile A-1 | 101 | 201 |
| 102 | | — | 1001 | ∞ | — | | 102 | 202 |
| 103 | | 305 | 1002 | ∞ | — | | 103 | 203 |
| 104 | | | | | | | | |
| 105 | | — | 1002 | ∞ | — | Tile B-1 | 104 | 204 |
| | | — | 1002 | 1003 | — | | 105 | 205 |
| | | — | 1004 | ∞ | [1, 1] | | 101 | 301 |

**Metadata for MVCC**

# MVCC

## ii) Delete (atomic)



Metadata for MVCC

# MVCC

## ii) Delete

↓ (after commit)



| Index |
|-------|
| 101 |
| 102 |
| 103 |
| 104 |
| 105 |

| TxnId | BeginCTS | EndCTS | PreV |
|-------|----------|--------|------|
| — | 1001 | 1004 | — |
| — | 1001 | ∞ | — |
| 305 | 1002 | ∞ | — |

| TxnId | BeginCTS | EndCTS | PreV |
|-------|----------|--------|------|
| — | 1002 | ∞ | — |
| — | 1002 | 1003 | — |
| — | 1004 | ∞ | [1, 1] |

| ID | SALARY |
|----|--------|

**Tile A-1**
| 101 |
| 102 |
| 103 |

**Tile A-2**
| 201 |
| 202 |
| 203 |

**Tile B-1**
| 104 | 204 |
| 105 | 205 |
| 101 | 301 |

**Metadata for MVCC**

# MVCC

iii) Update

(set the previous one to invalid)



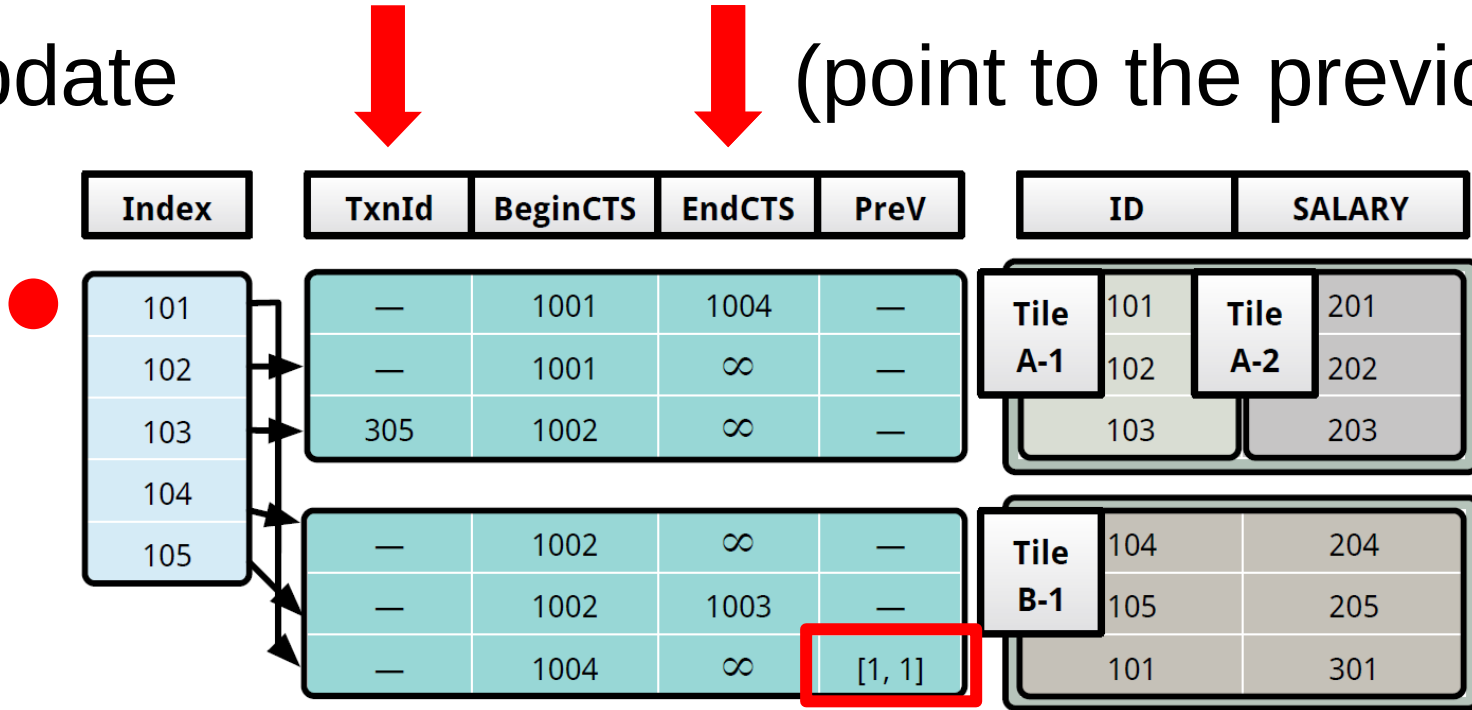| Index | | TxnId | BeginCTS | EndCTS | PreV | | ID | SALARY |
|-------|---|-------|----------|--------|------|---|-----|--------|
| 101 | | — | 1001 | **1004** | — | Tile A-1 / 101 | Tile A-2 / 201 | |
| 102 | | — | 1001 | ∞ | — | 102 | 202 | |
| 103 | | 305 | 1002 | ∞ | — | 103 | 203 | |
| 104 | | | | | | | | |
| 105 | | — | 1002 | ∞ | — | Tile B-1 / 104 | 204 | |
| | | — | 1002 | 1003 | — | 105 | 205 | |
| | | — | 1004 | ∞ | [1, 1] | 101 | 301 | |

Metadata for MVCC

# MVCC

iii) Update                    (point to the previous one)



Metadata for MVCC

# MVCC

## iii) Update

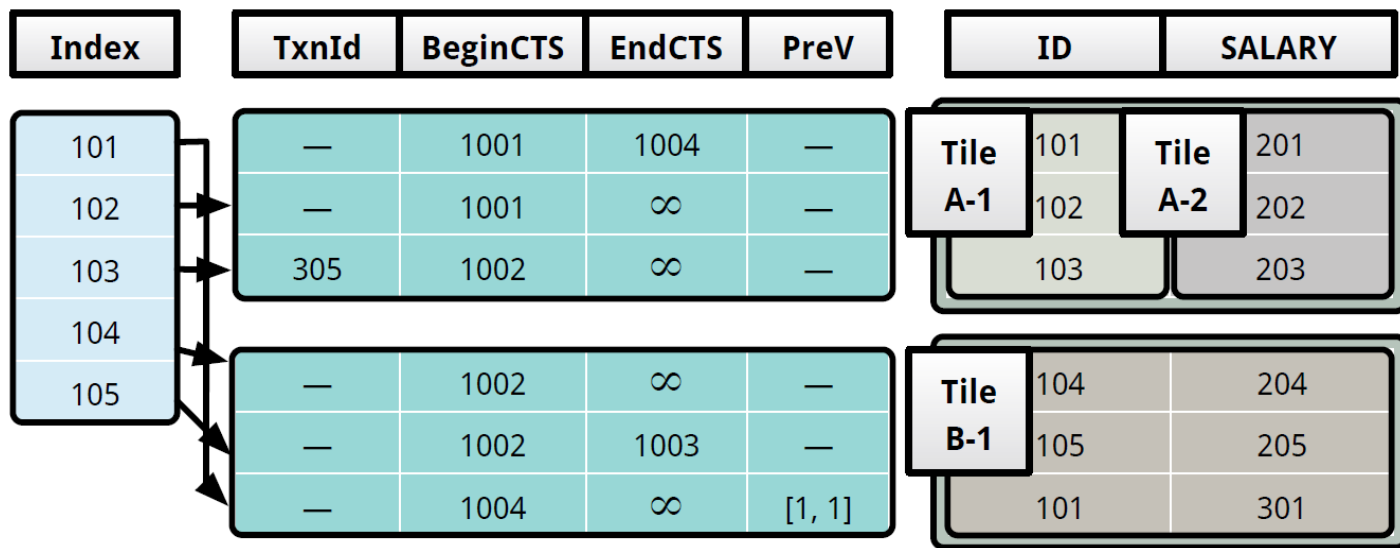(point to the previous one)

| Index | | TxnId | BeginCTS | EndCTS | PreV | | ID | SALARY |
|-------|--|-------|----------|--------|------|--|----|--------|
| 101 | | — | 1001 | 1004 | — | Tile A-1 | 101 | Tile A-2 201 |
| 102 | | — | 1001 | ∞ | — | | 102 | 202 |
| 103 | | 305 | 1002 | ∞ | — | | 103 | 203 |
| 104 | | | | | | | | |
| 105 | | — | 1002 | ∞ | — | Tile B-1 | 104 | 204 |
| | | — | 1002 | 1003 | — | | 105 | 205 |
| | | — | 1004 | ∞ | [1, 1] | | 101 | 301 |

**Metadata for MVCC**

# MVCC

Search



Looking through the metadata, find the one within

Its visibility (BeginCTS < transaction ID < EndCTS)

# MVCC



| Index |
|-------|
| 101 |
| 102 |
| 103 |
| 104 |
| 105 |

| TxnId | BeginCTS | EndCTS | PreV |
|-------|----------|--------|------|
| — | 1001 | 1004 | — |
| — | 1001 | ∞ | — |
| 305 | 1002 | ∞ | — |

| TxnId | BeginCTS | EndCTS | PreV |
|-------|----------|--------|------|
| — | 1002 | ∞ | — |
| — | 1002 | 1003 | — |
| — | 1004 | ∞ | [1, 1] |

| ID | SALARY |
|----|--------|
| Tile A-1 / 101 | Tile A-2 / 201 |
| 102 | 202 |
| 103 | 203 |

| ID | SALARY |
|----|--------|
| Tile B-1 / 104 | 204 |
| 105 | 205 |
| 101 | 301 |

Garbage Collector claims those old spaces back

# Performance

What about the overhead?

Is this method practical?

# Performance

What about the overhead?

Is this method practical?

Let us see the results.

# Evaluation

# NSM / DSM / FSM  SCAN



(a) Scan, Narrow, Read Only

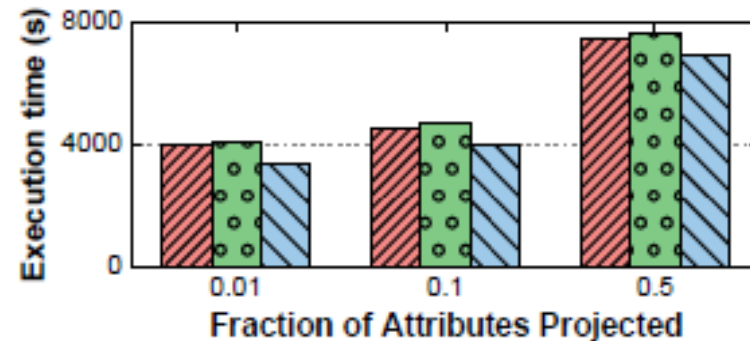(b) Scan, Narrow, Hybrid

(c) Scan, Wide, Read Only

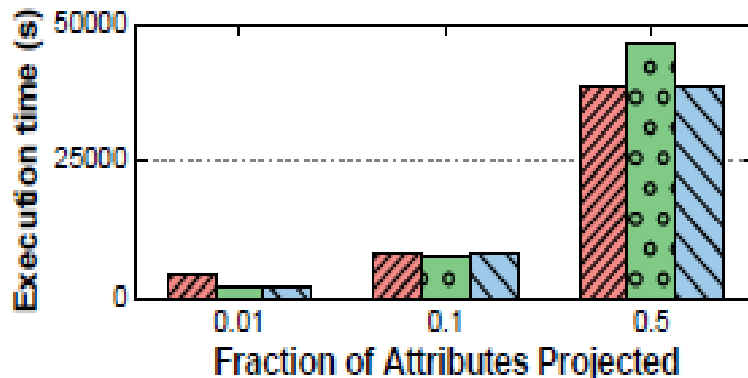(d) Scan, Wide, Hybrid

Storage Models :  NSM   DSM   FSM

# NSM / DSM / FSM  Aggregate



(e) Aggregate, Narrow, Read Only
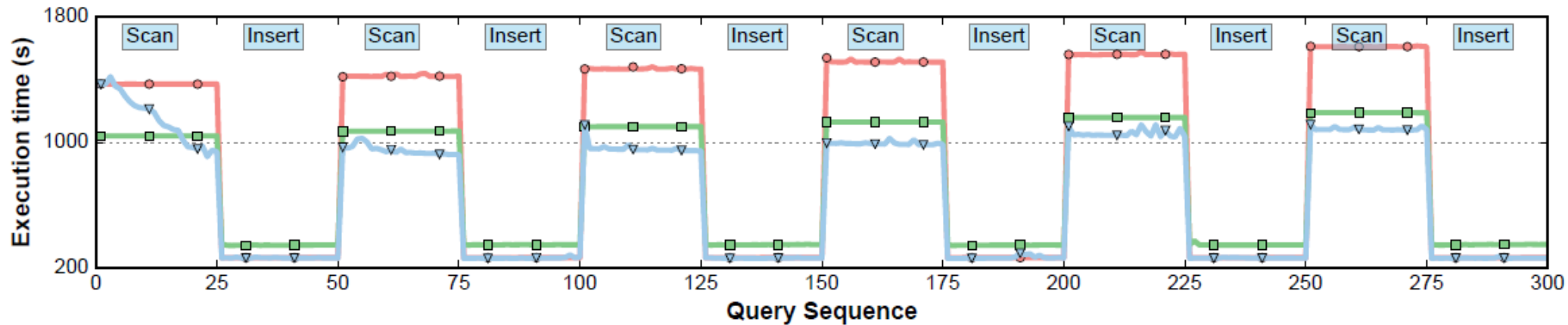
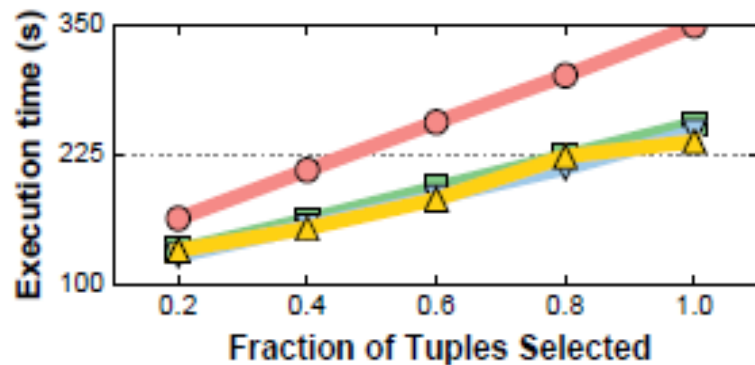(f) Aggregate, Narrow, Hybrid

(g) Aggregate, Wide, Read Only

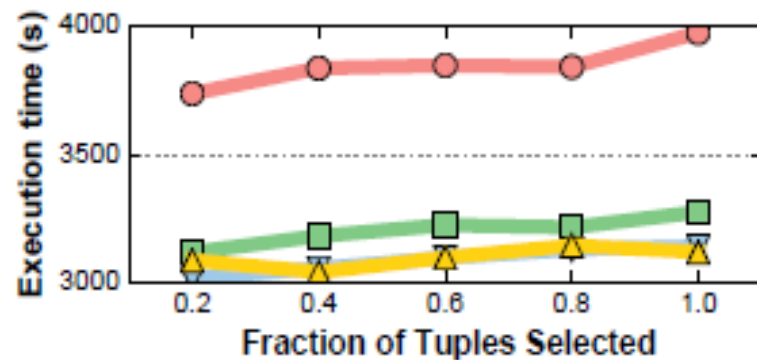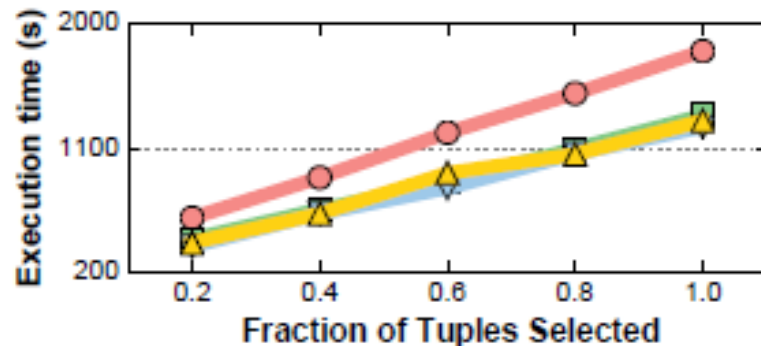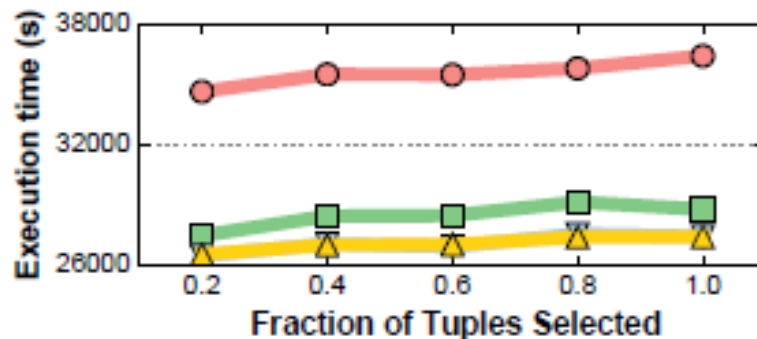(h) Aggregate, Wide, Hybrid

Storage Models :  NSM  DSM  FSM

# Adaptation

# Logical tiles



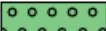(a) Scan, Narrow, Read Only

(b) Scan, Narrow, Hybrid

(c) Scan, Wide, Read Only

(d) Scan, Wide, Hybrid

Storage Models: NSM   DSM   FSM

# Conclusion

# Conclusion

- Memory Layout (column storage, row storage)
- + Adaptation → FSM
- + Abstraction → Logical Tiles
- + Concurrency control → MVCC