



**THE CHIPS  
TO SYSTEMS  
CONFERENCE**

SPONSORED BY:



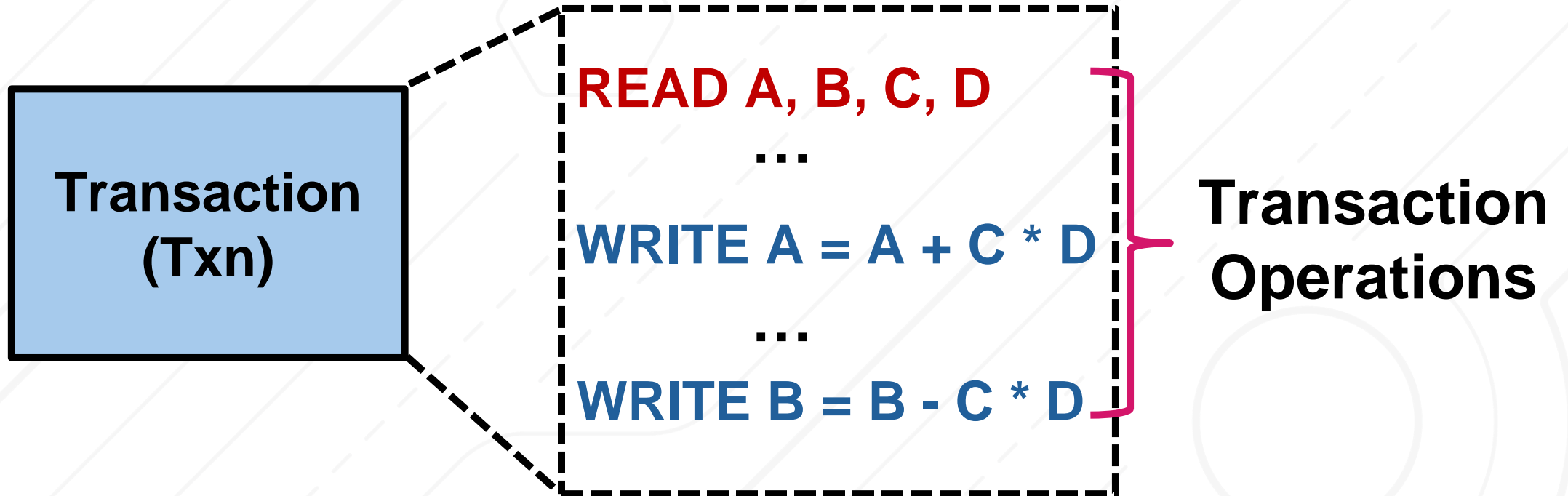
# Onyx: Efficient Transaction Processing with Real Processing-in-Memory Prototypes

Menglei Chen, Yixiao Wang, Yu Hua

Huazhong University of Science and Technology

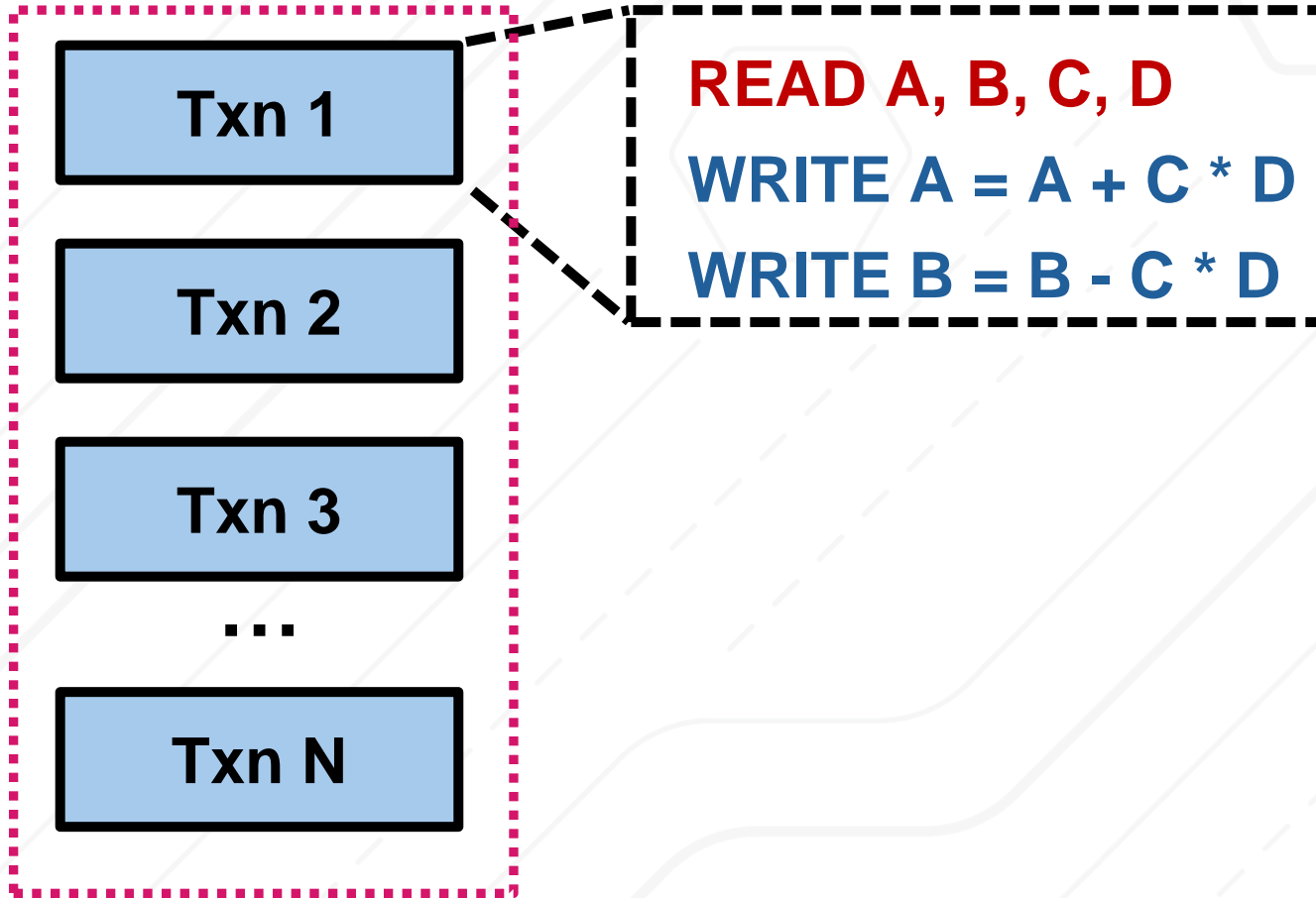


# Online Transaction Processing





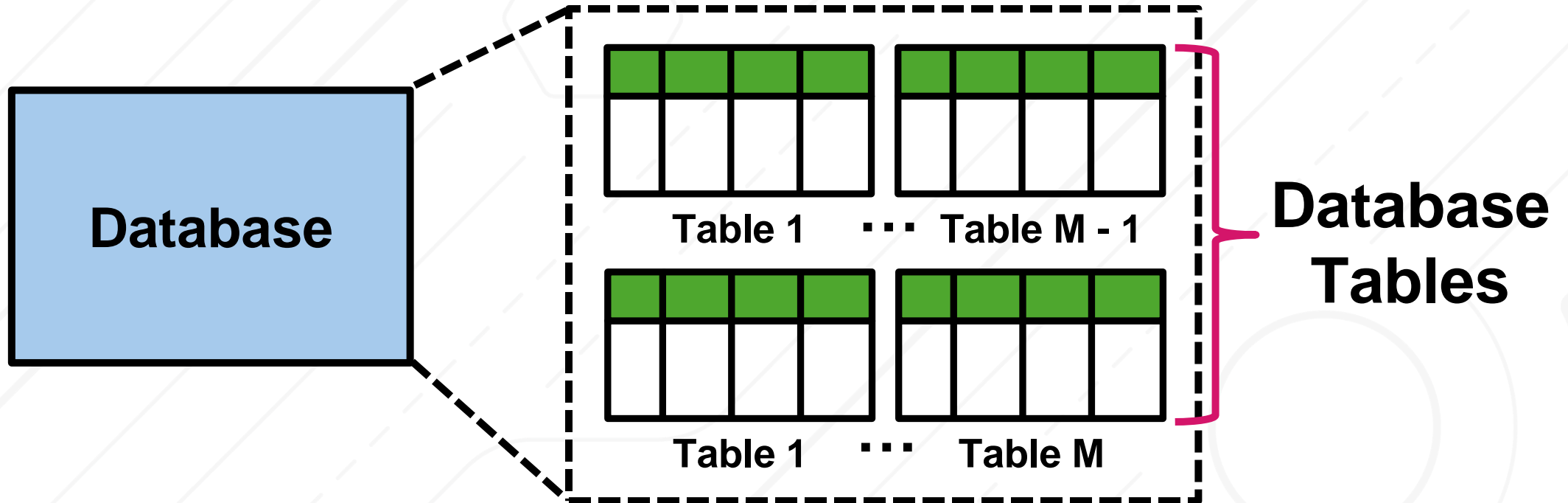
# Online Transaction Processing



*Processed in parallel with concurrency control mechanisms*

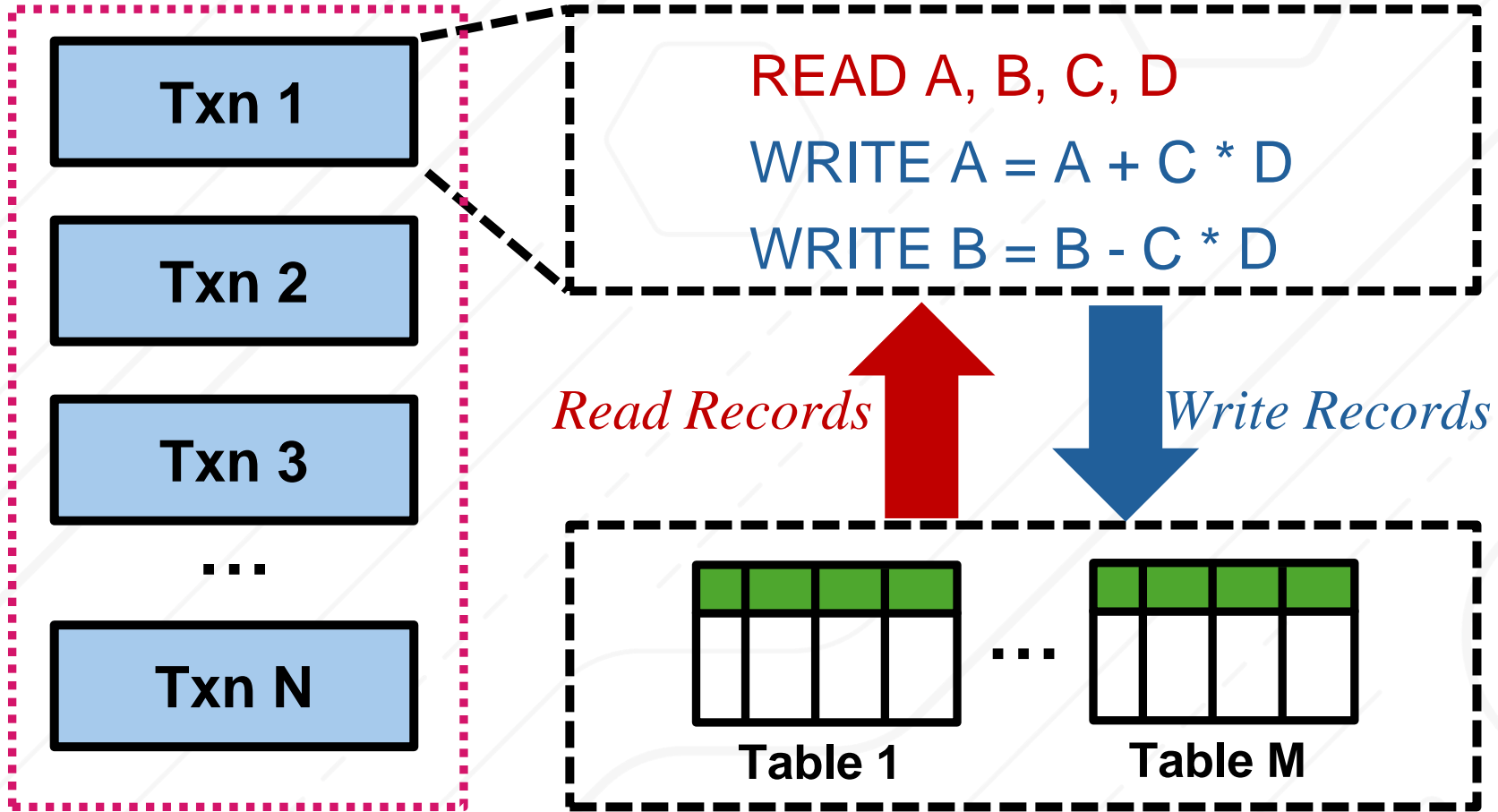


# Online Transaction Processing





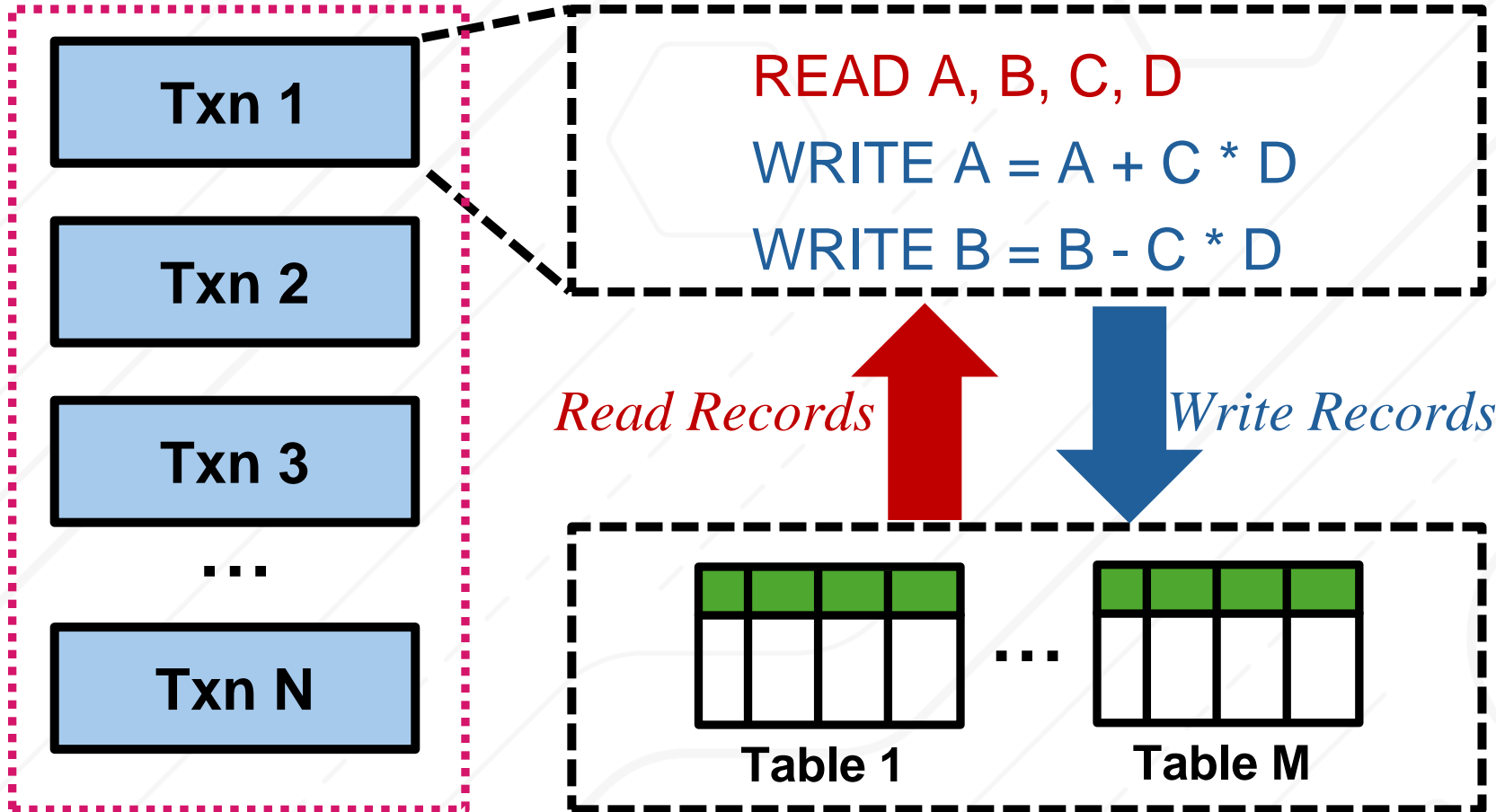
# Online Transaction Processing



*Processed in parallel with concurrency control mechanisms*



# Online Transaction Processing



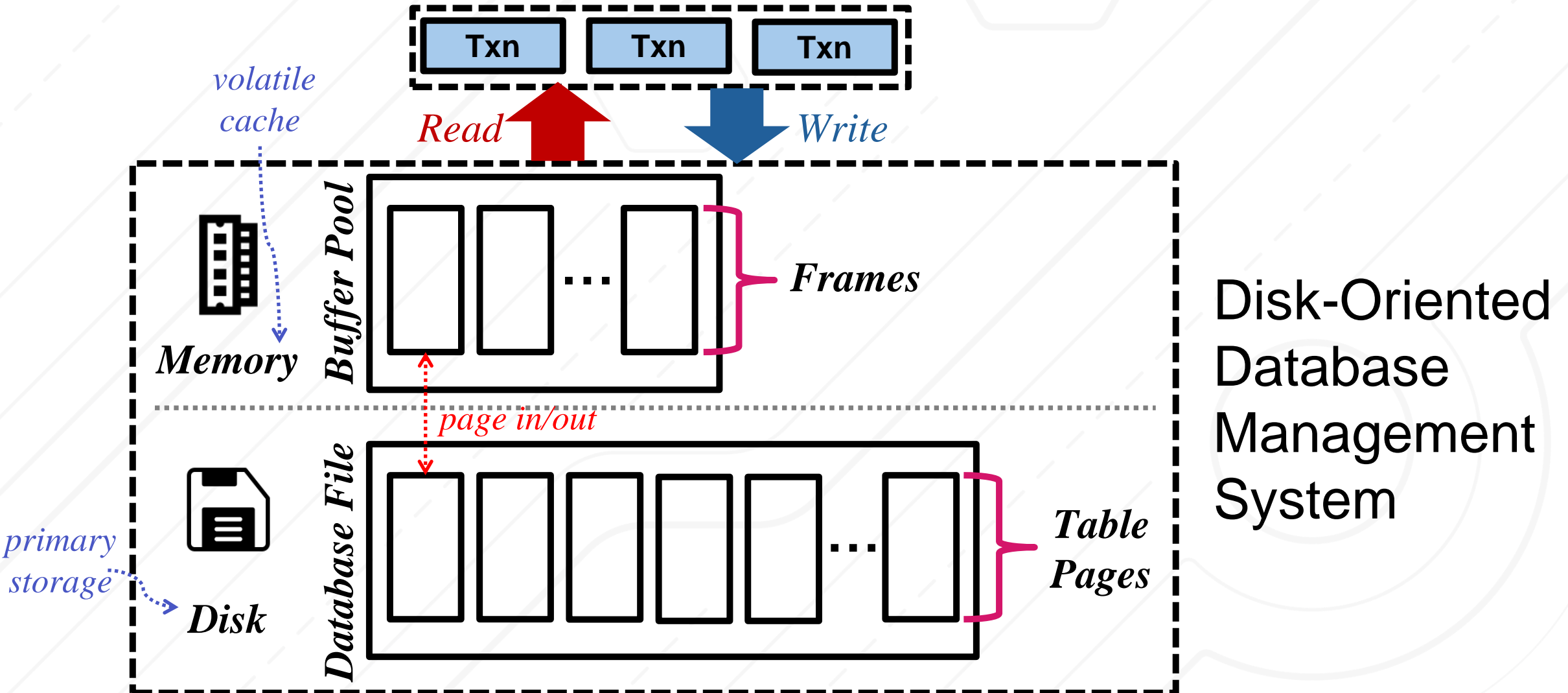
*Enabling:*

- *Atomicity*
- *Consistency*
- *Isolation*
- *Durability*

*Processed in parallel with concurrency control mechanisms*



# Existing OLTP Systems





# Existing OLTP Systems



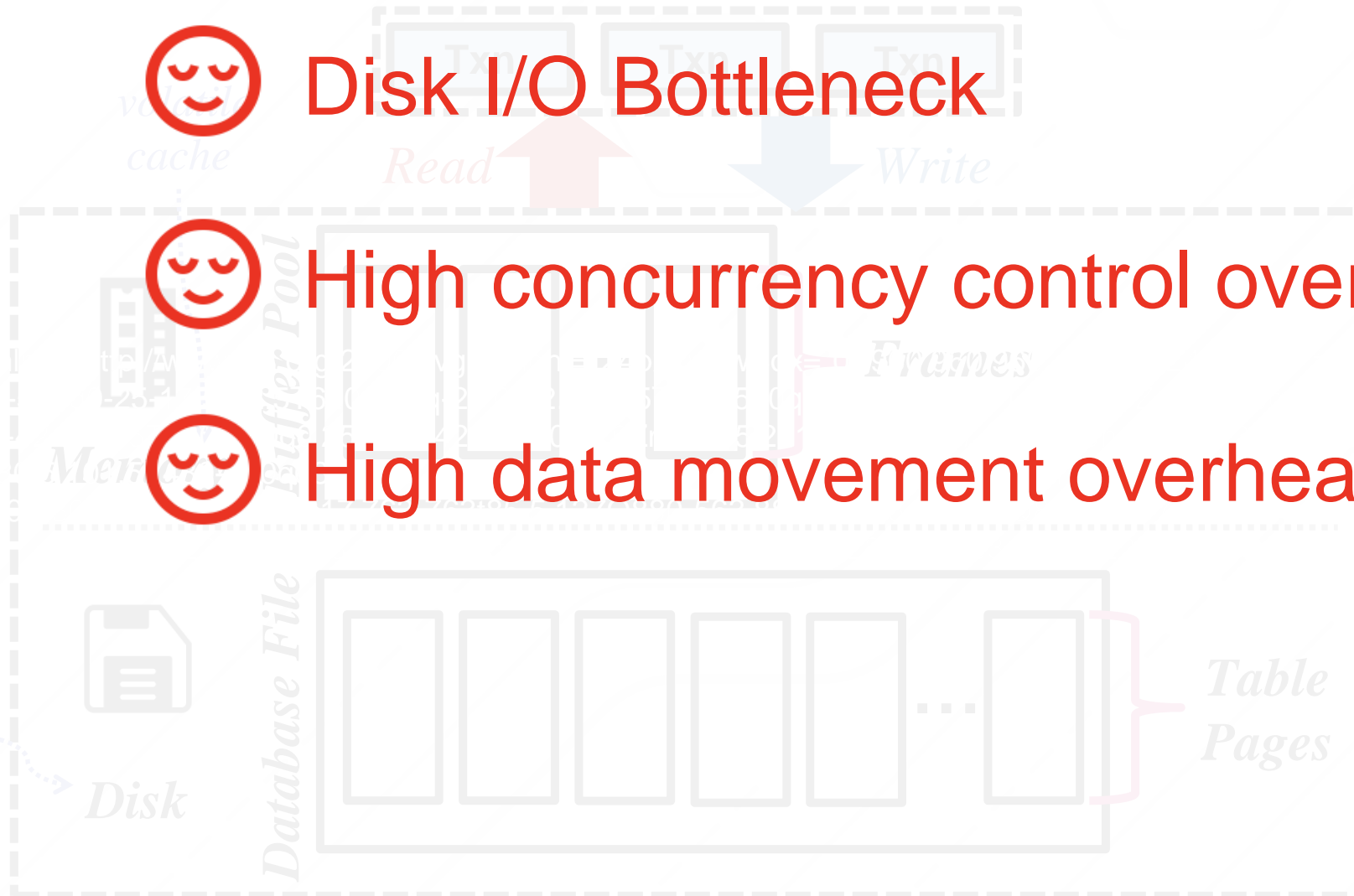
Disk I/O Bottleneck



High concurrency control overheads



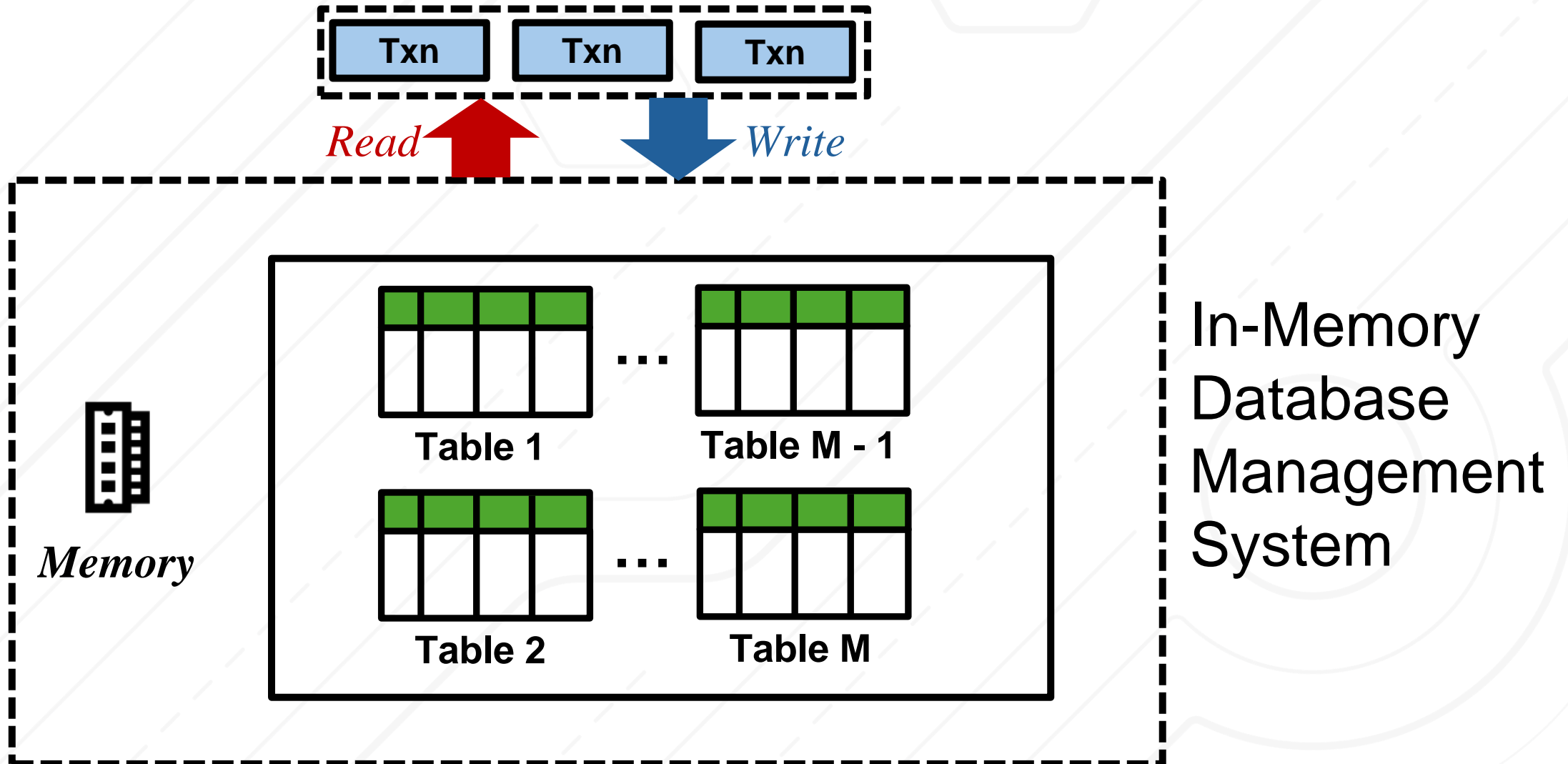
High data movement overheads



Disk-Oriented  
Database  
Management  
System



# Existing OLTP Systems





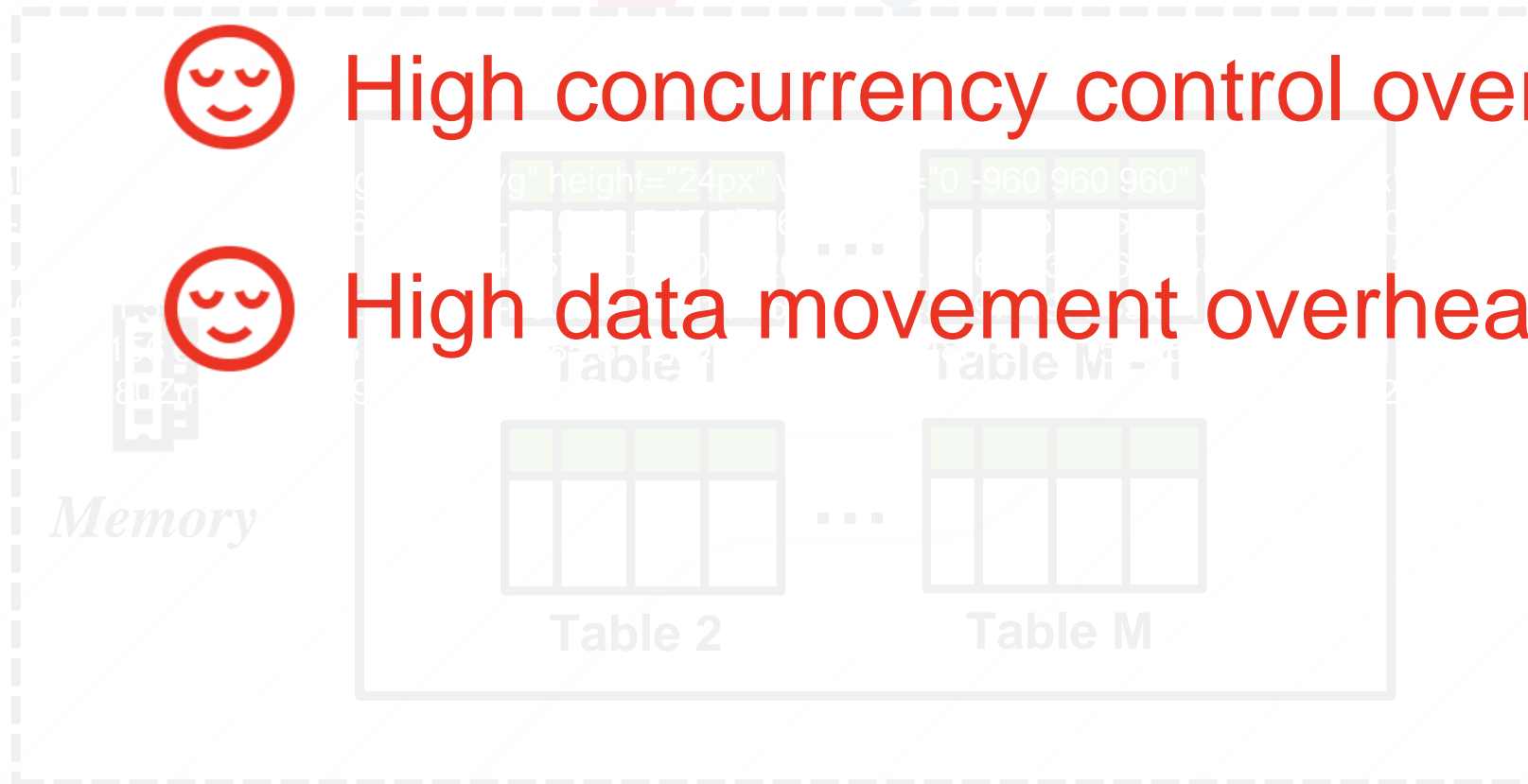
# Existing OLTP Systems

😊 Address disk I/O bottleneck



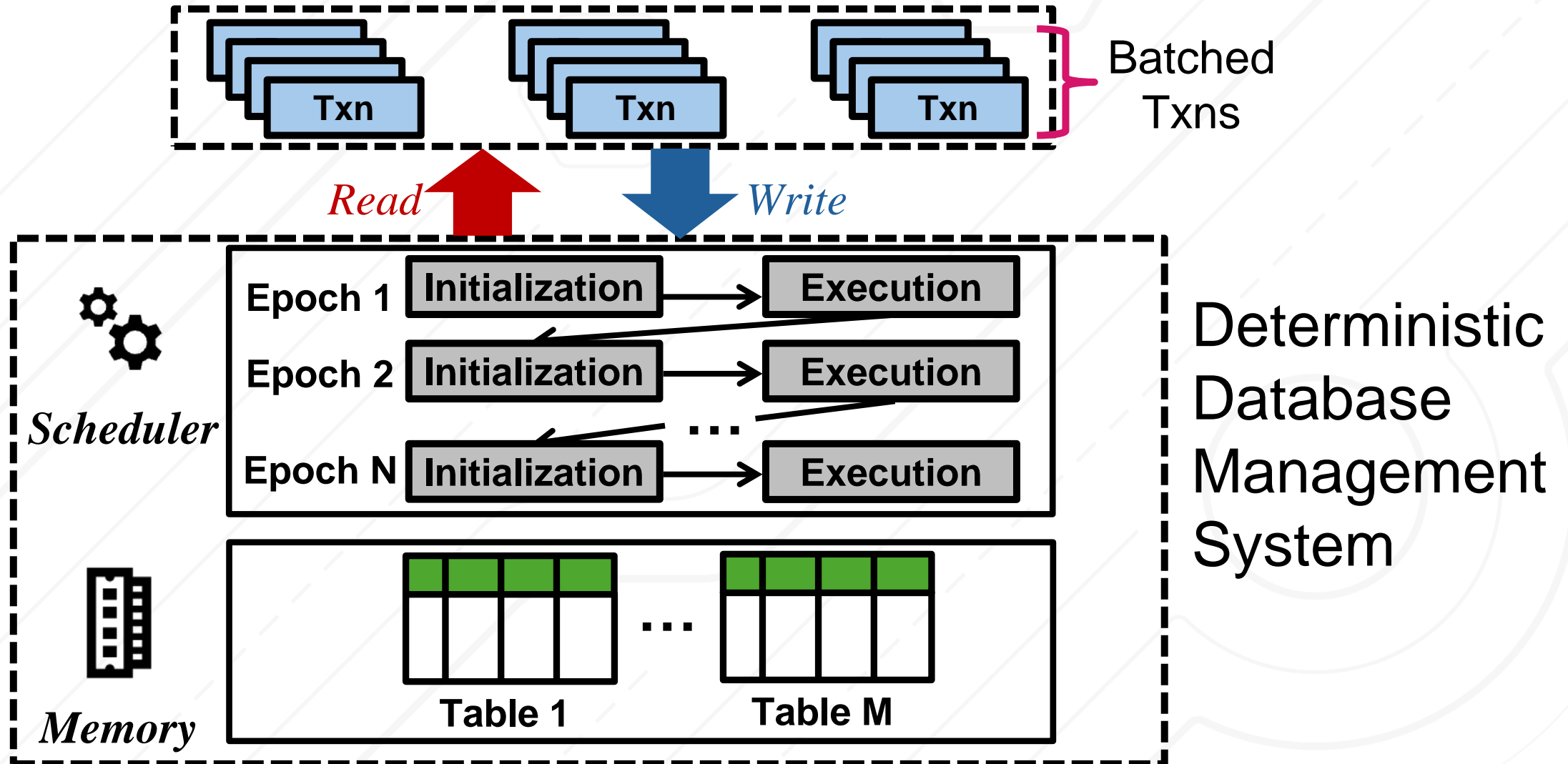
😞 High concurrency control overheads

😞 High data movement overheads





# Existing OLTP Systems





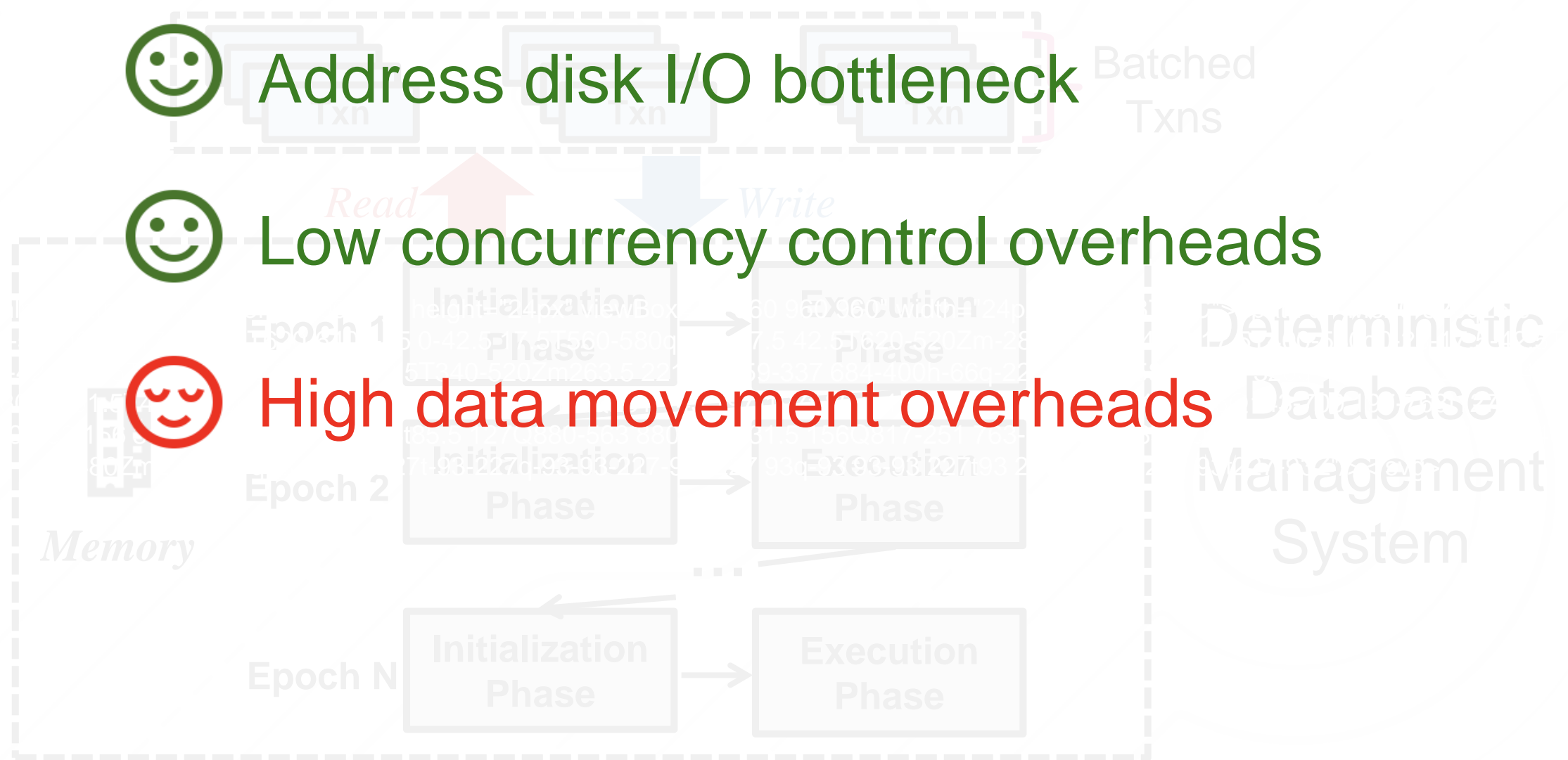
# Existing OLTP Systems

😊 Address disk I/O bottleneck

Batched Txns

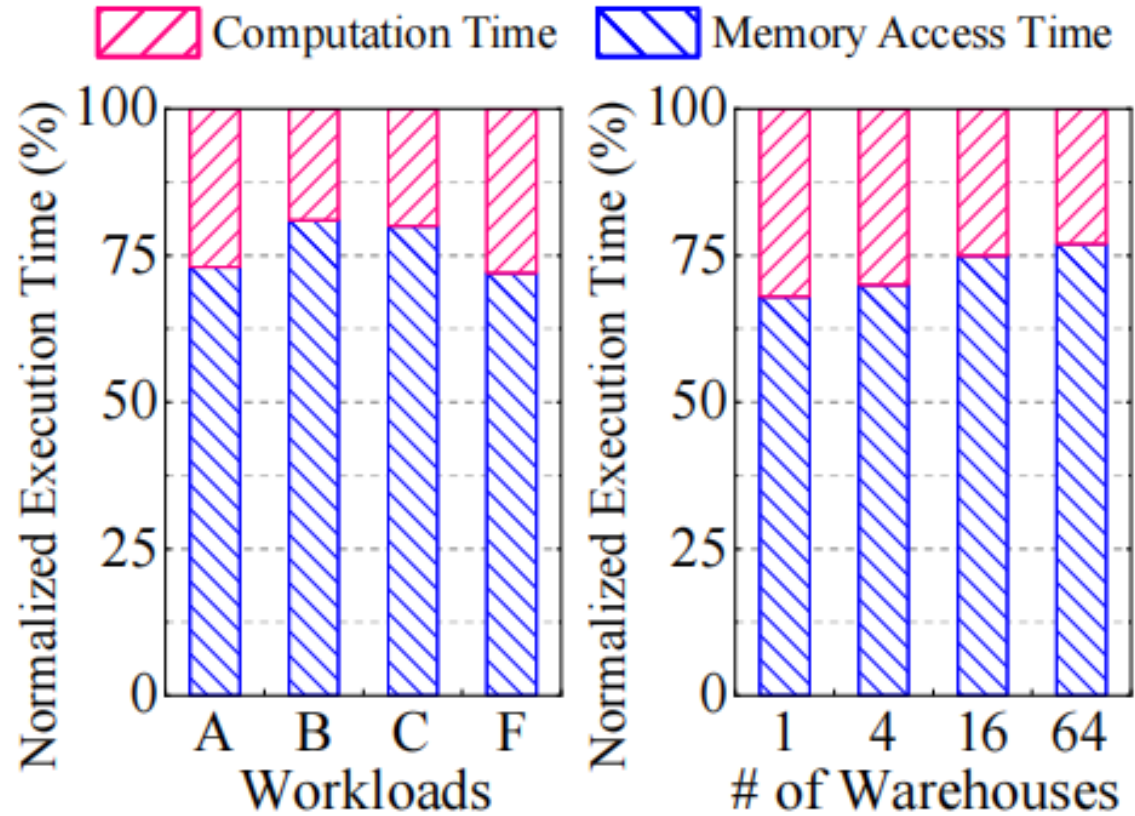
😊 Low concurrency control overheads

😞 High data movement overheads





# Existing OLTP Systems



(a) YCSB

(b) TPC-C

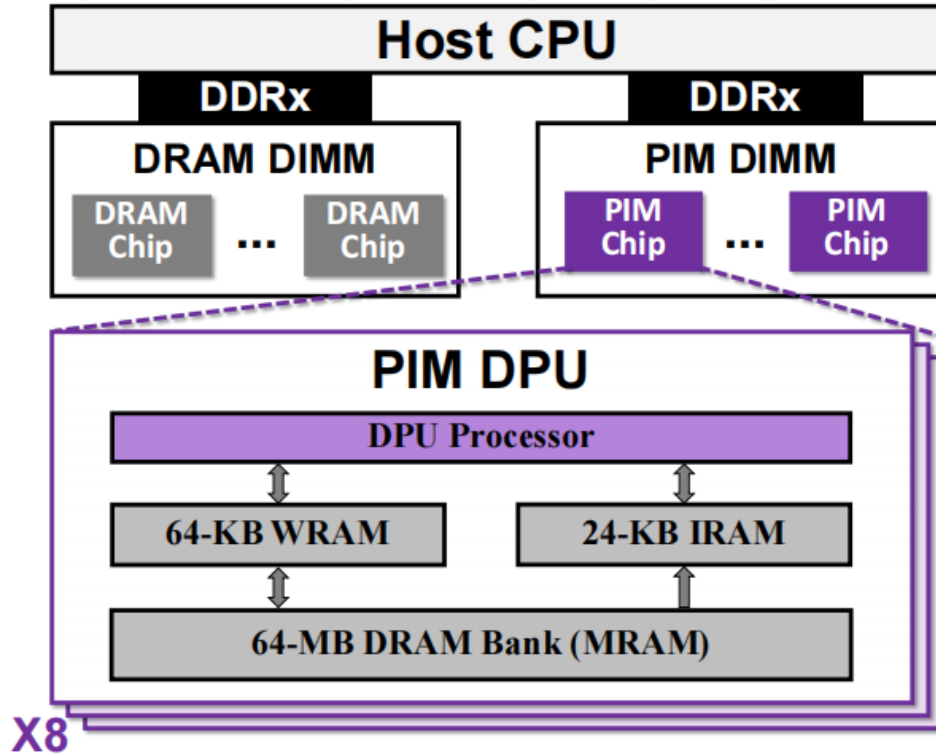


68% ~ 81%

*Memory bandwidth becomes the performance bottleneck!*



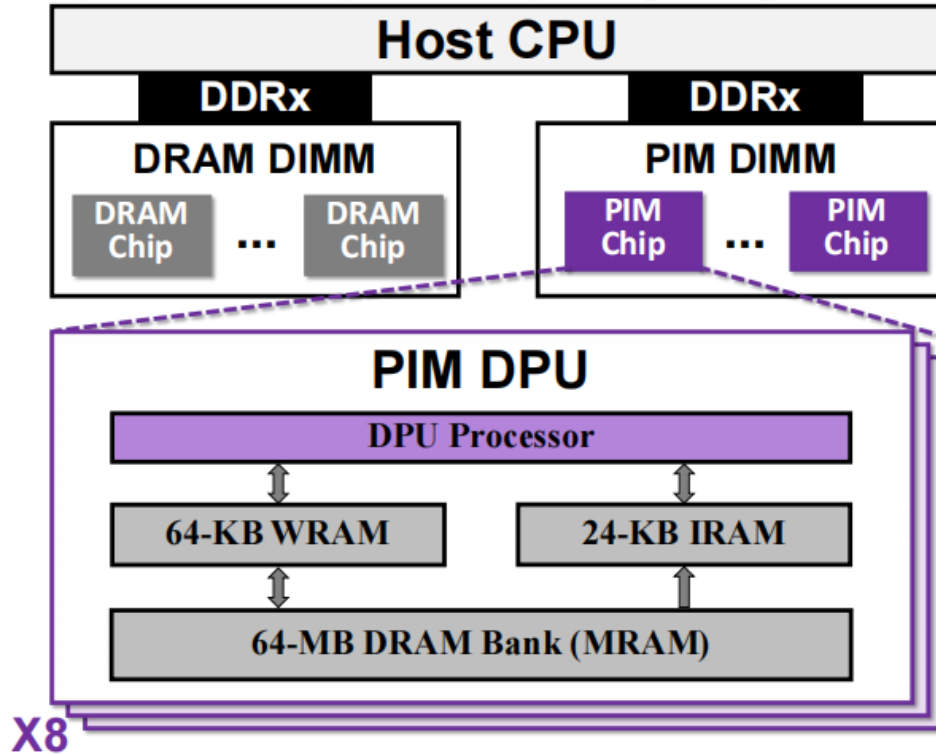
# UPMEM Processing In Memory



- 1.6 TB/s high memory bandwidth
- Offloading computing from CPU to DPUs



# UPMEM Processing In Memory



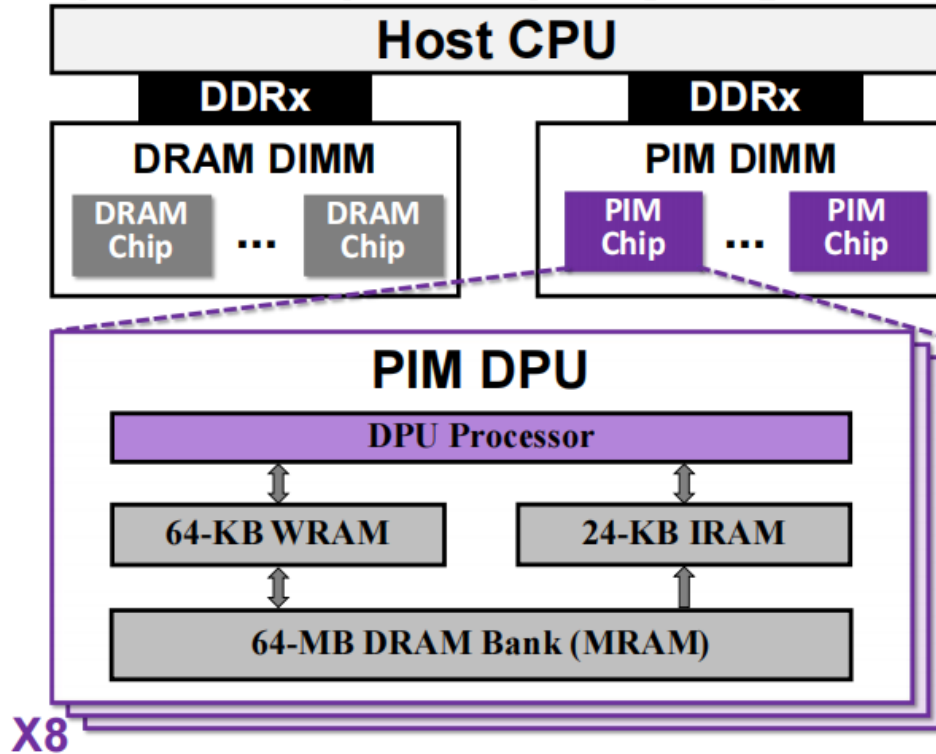
- 1.6 TB/s high memory bandwidth
- Offloading computing from CPU to DPUs



*Promising!*



# UPMEM Processing In Memory



- 1.6 TB/s high memory bandwidth
- Offloading computing from CPU to DPUs



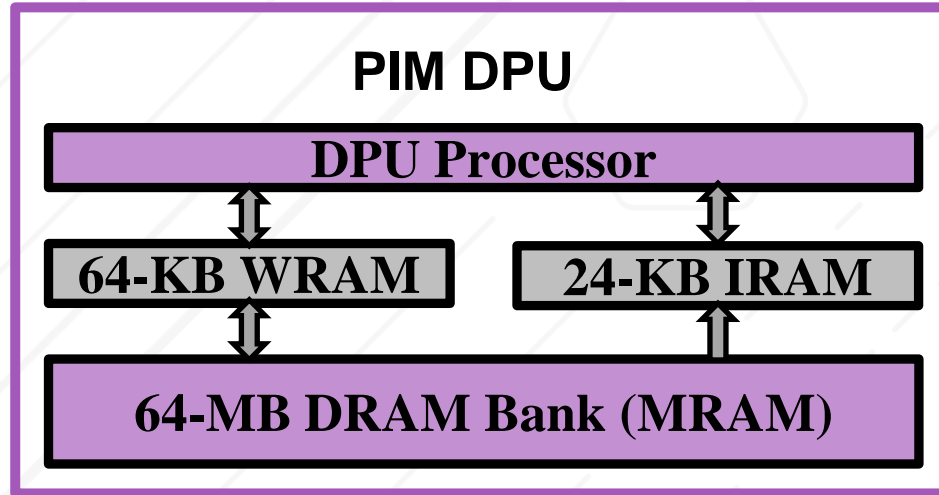
*Promising!*



*However, it is non-trivial to design a transaction processing system on PIM*

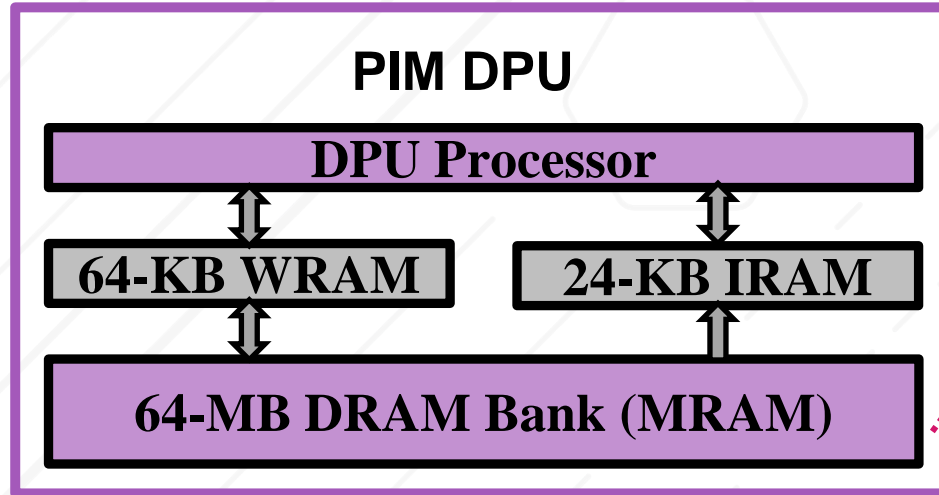


# Challenge 1: Execution across Partitioned DPUs





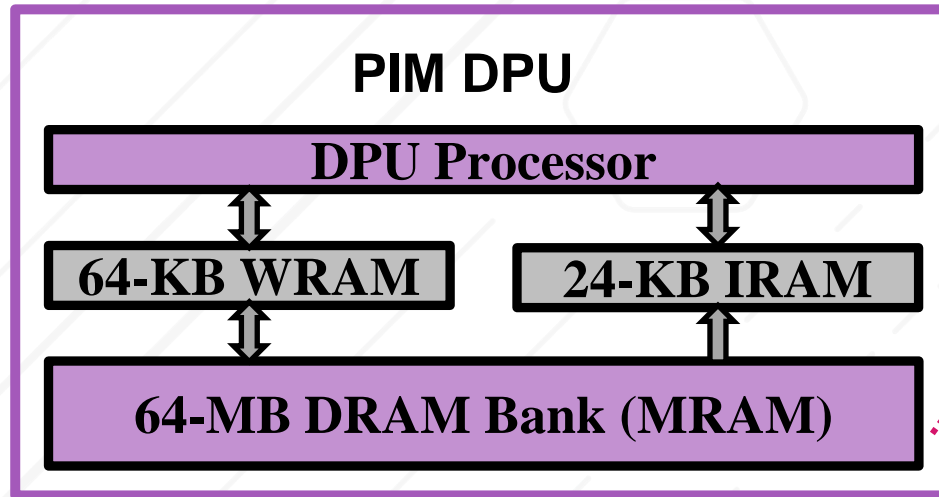
# Challenge 1: Execution across Partitioned DPUs



Limited memory per DPU

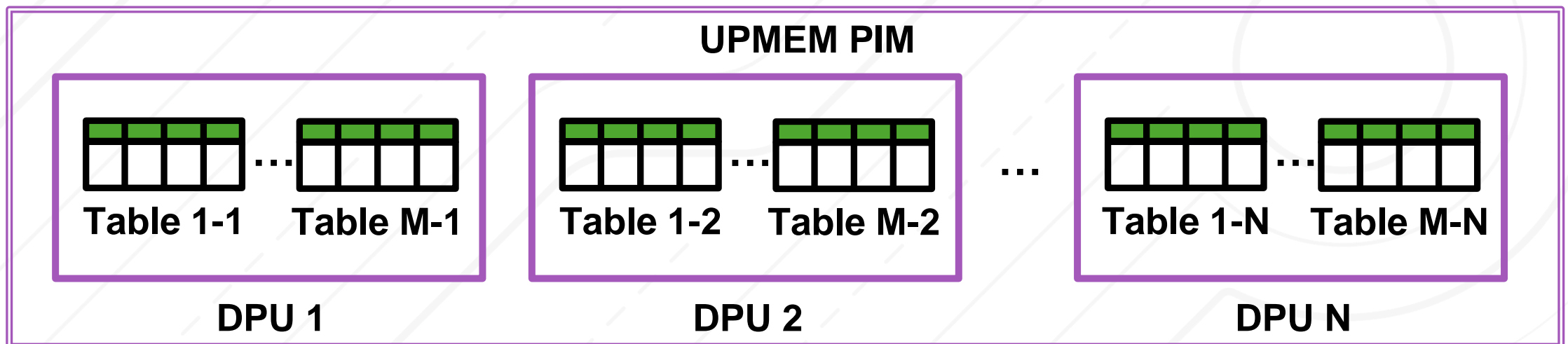


# Challenge 1: Execution across Partitioned DPUs



Limited memory per DPU

Database tables need to be horizontally partitioned



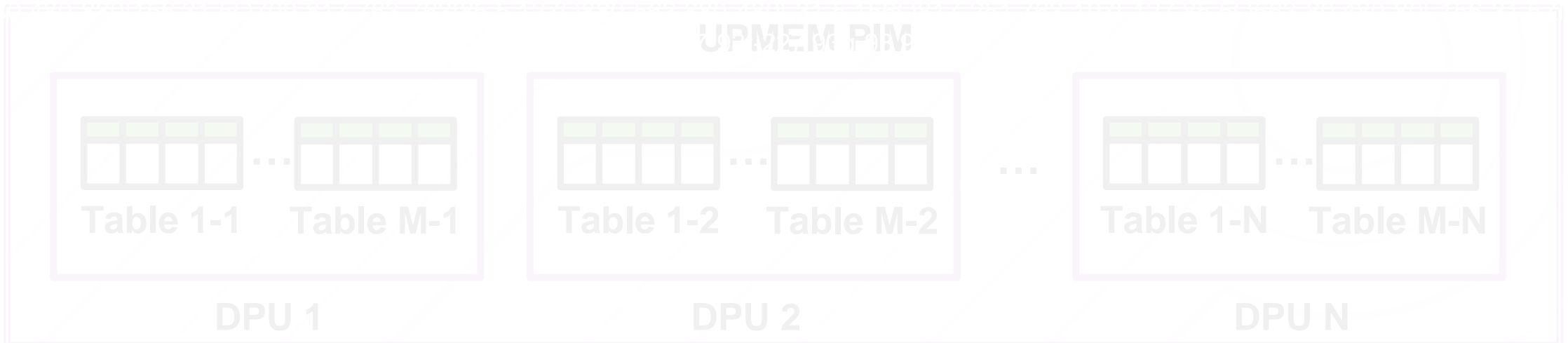
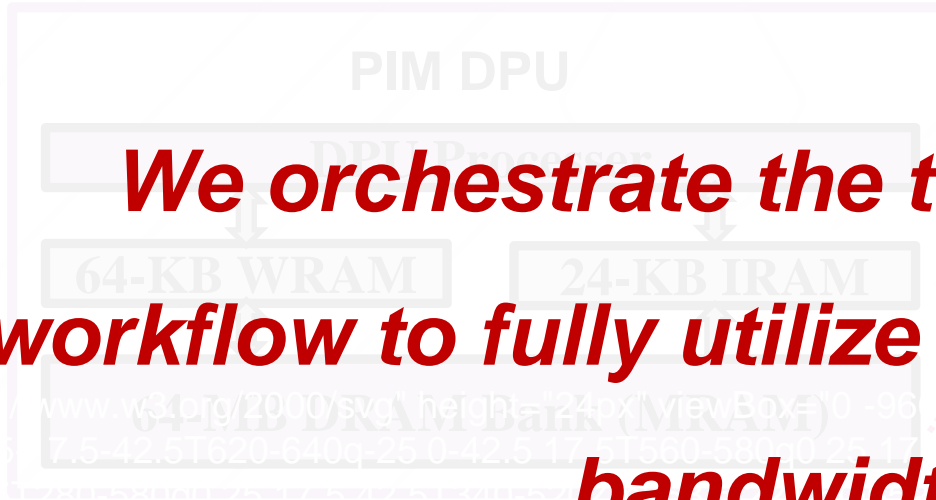


# Challenge 1: Execution across Partitioned DPUs

***We orchestrate the transaction execution workflow to fully utilize the high parallelism and bandwidth of DPU.***

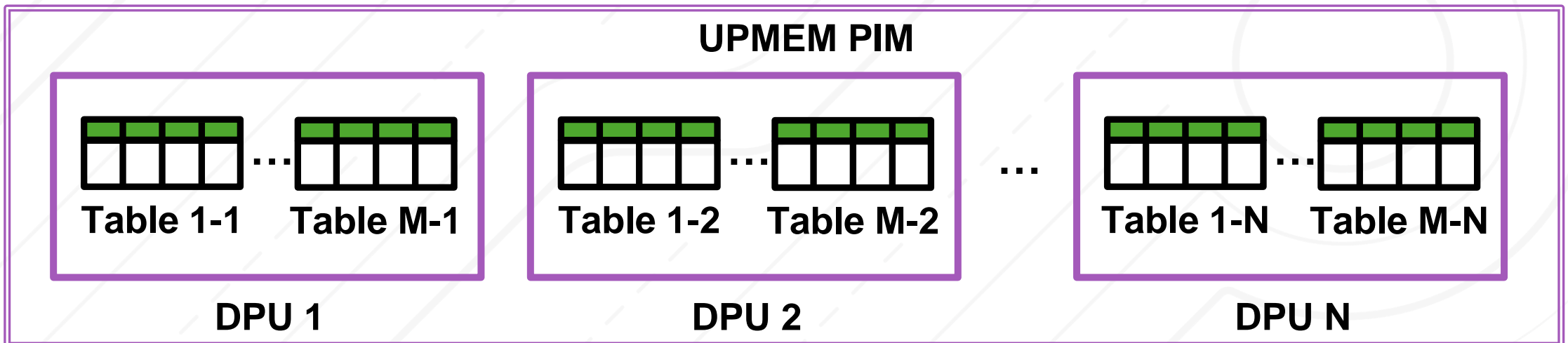
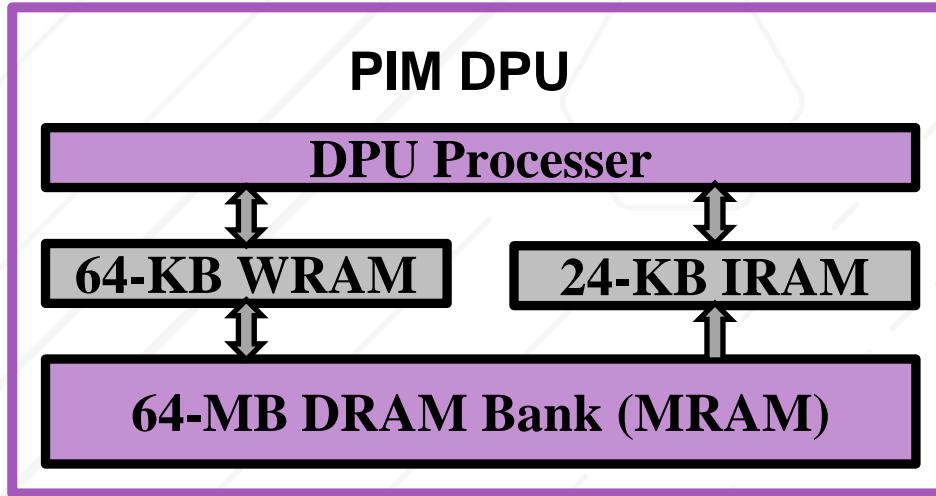
Limited memory per DPU

Database tables must be horizontally partitioned



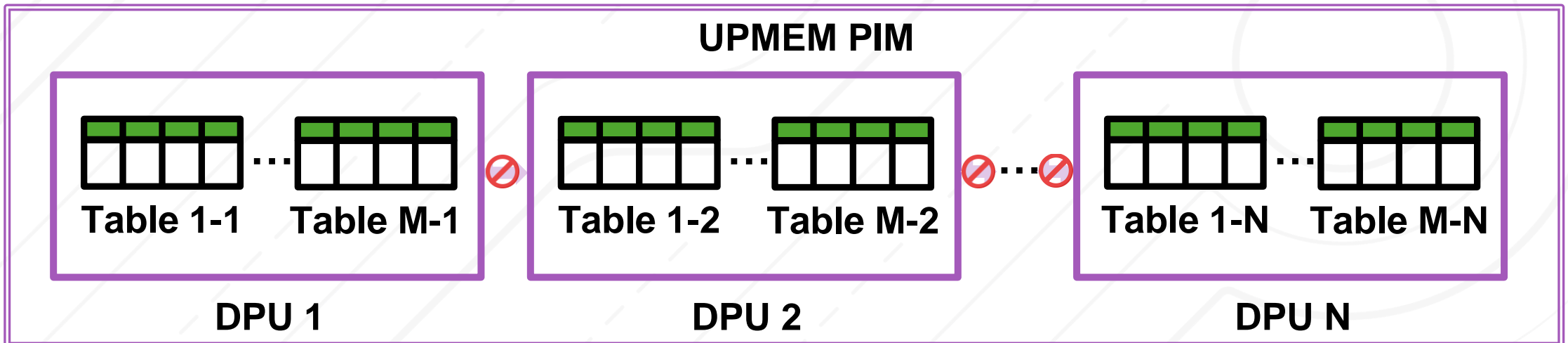
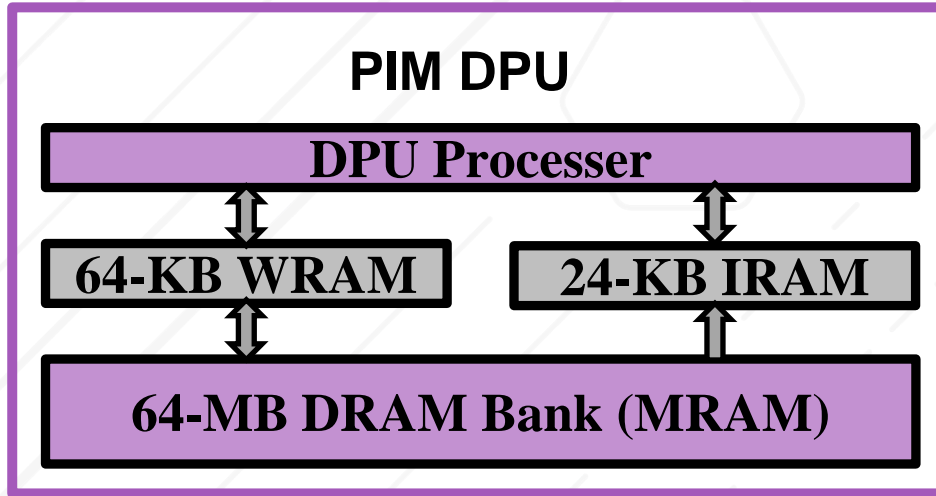


# Challenge 2: Lack of Inter-DPU Communication



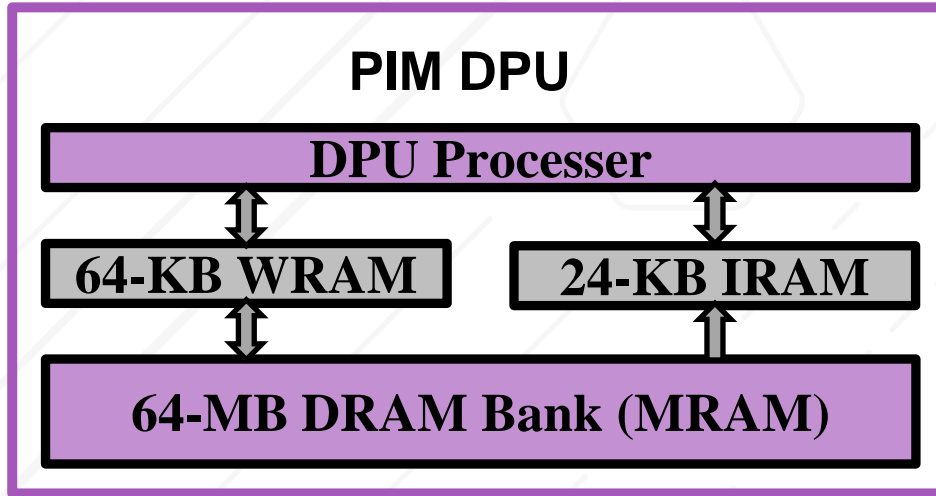


# Challenge 2: Lack of Inter-DPU Communication

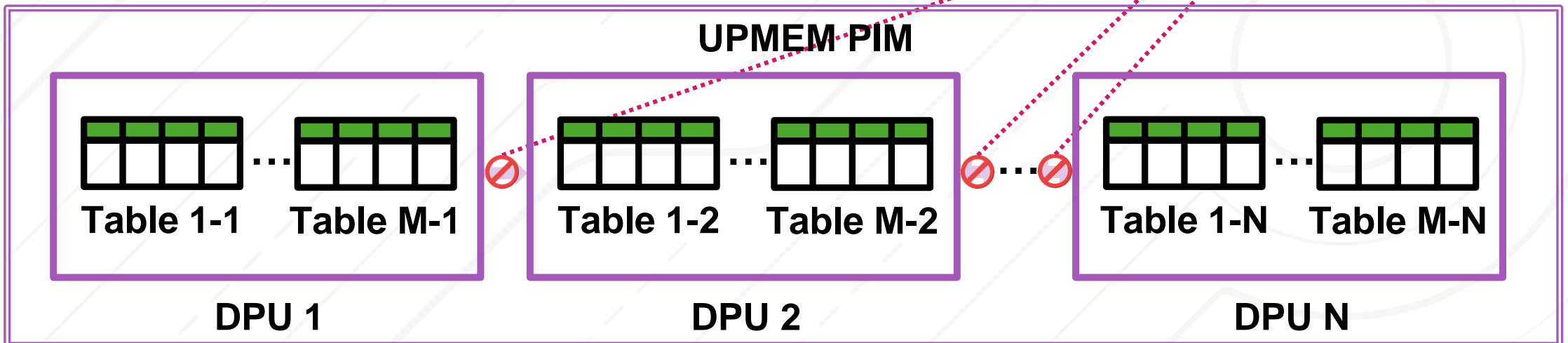




# Challenge 2: Lack of Inter-DPU Communication

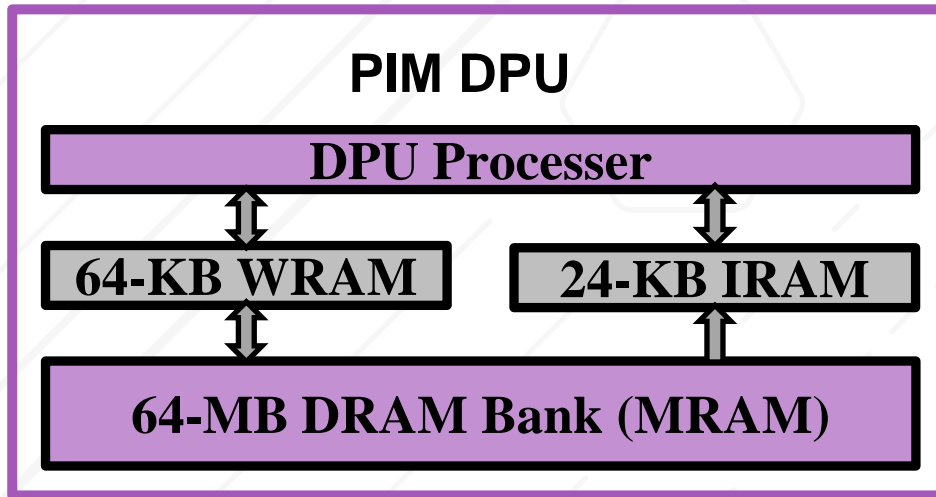


No direct data transfer between DPU's



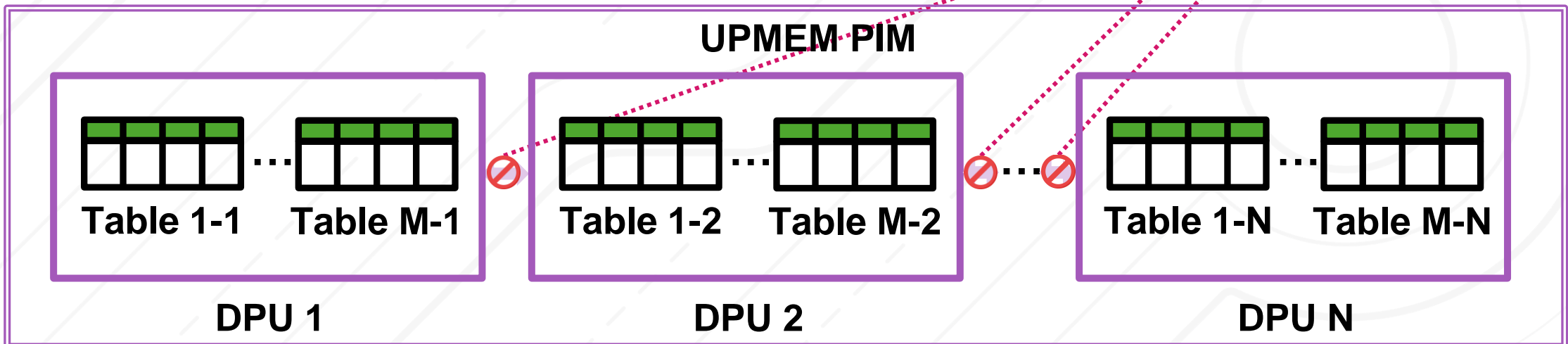


# Challenge 2: Lack of Inter-DPU Communication



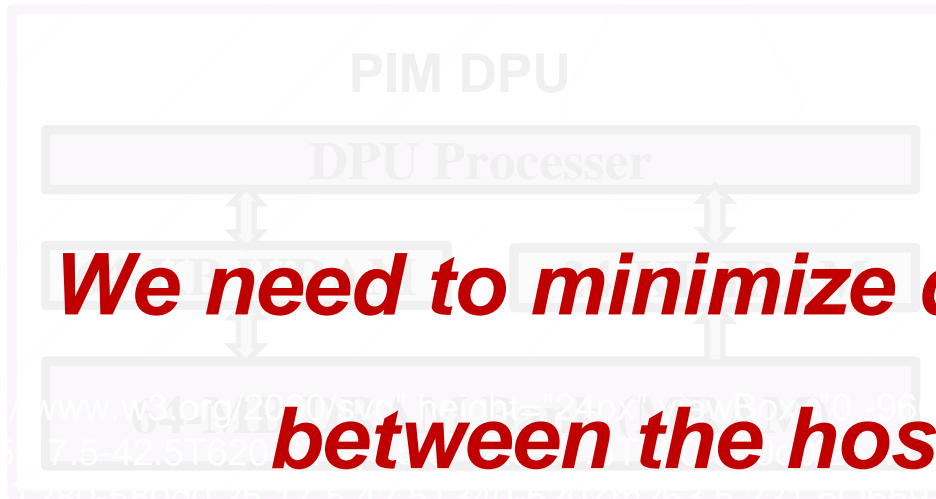
Data transfer overhead between CPU and DPUs is inevitable

No direct data transfer between DPUs





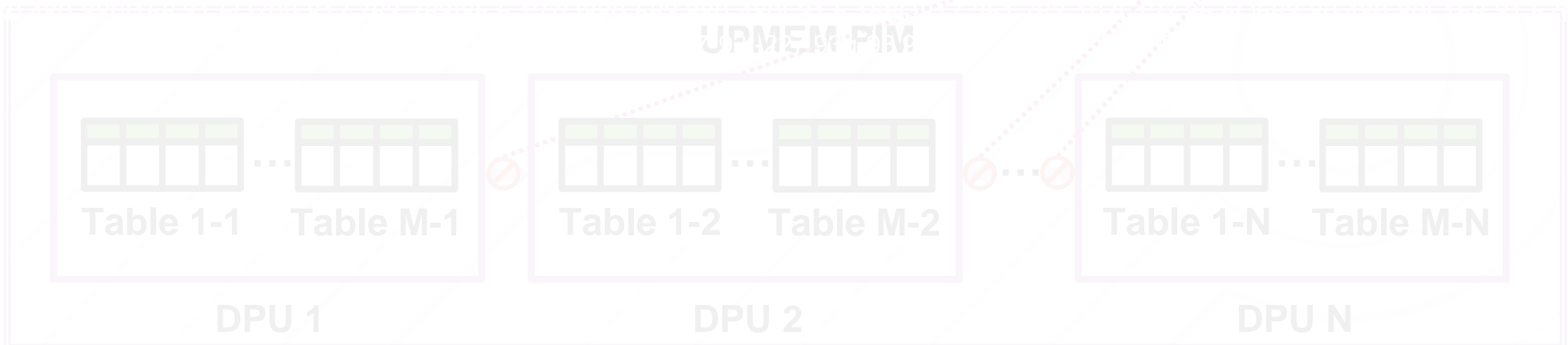
# Challenge 2: Lack of Inter-DPU Communication



***We need to minimize data transfer overhead between the host CPU and DPUs.***

Data transfer overhead between CPU and DPUs is inevitable

No direct data transfer between DPUs





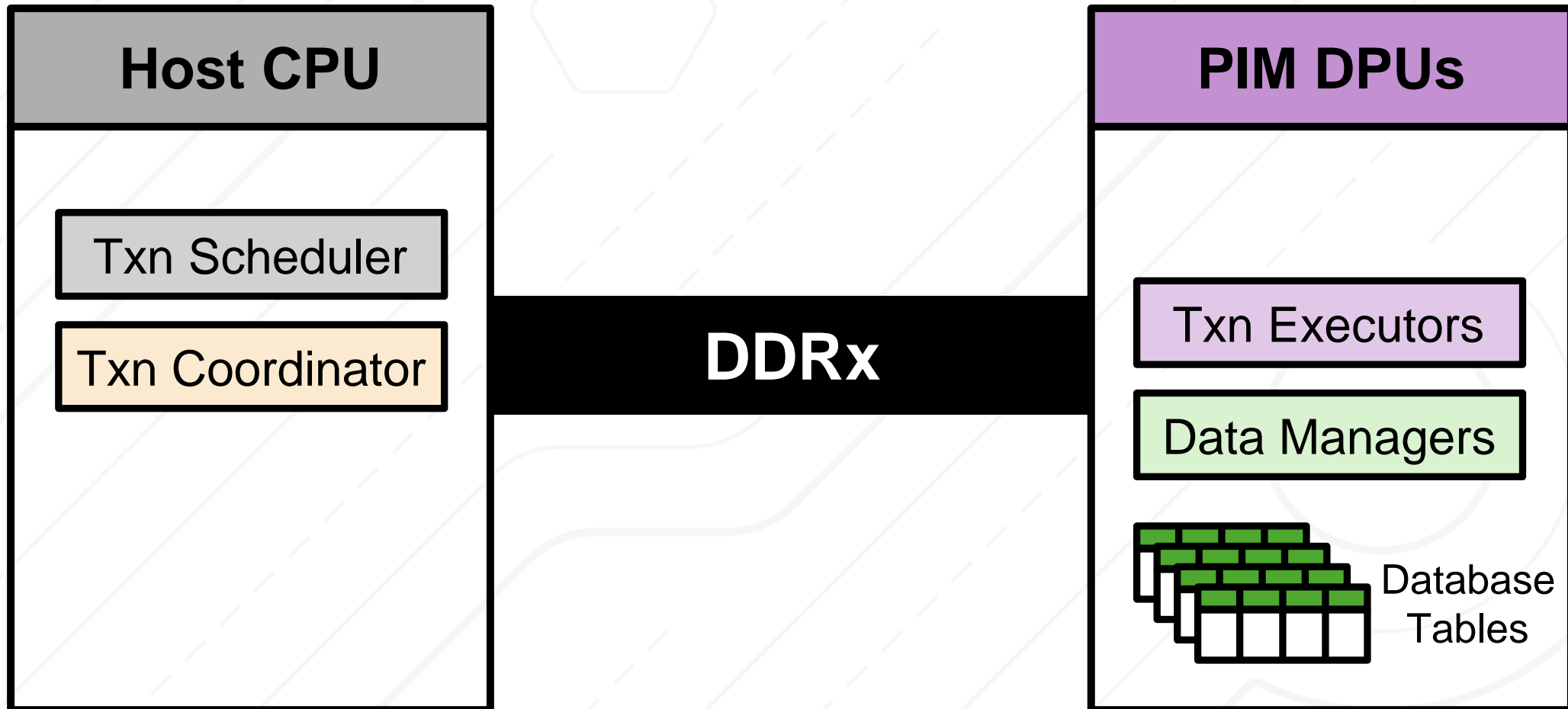
# Our Solution: Onyx

**Onyx:** *A Transaction Processing System with Real PIM Prototypes*



# Our Solution: Onyx

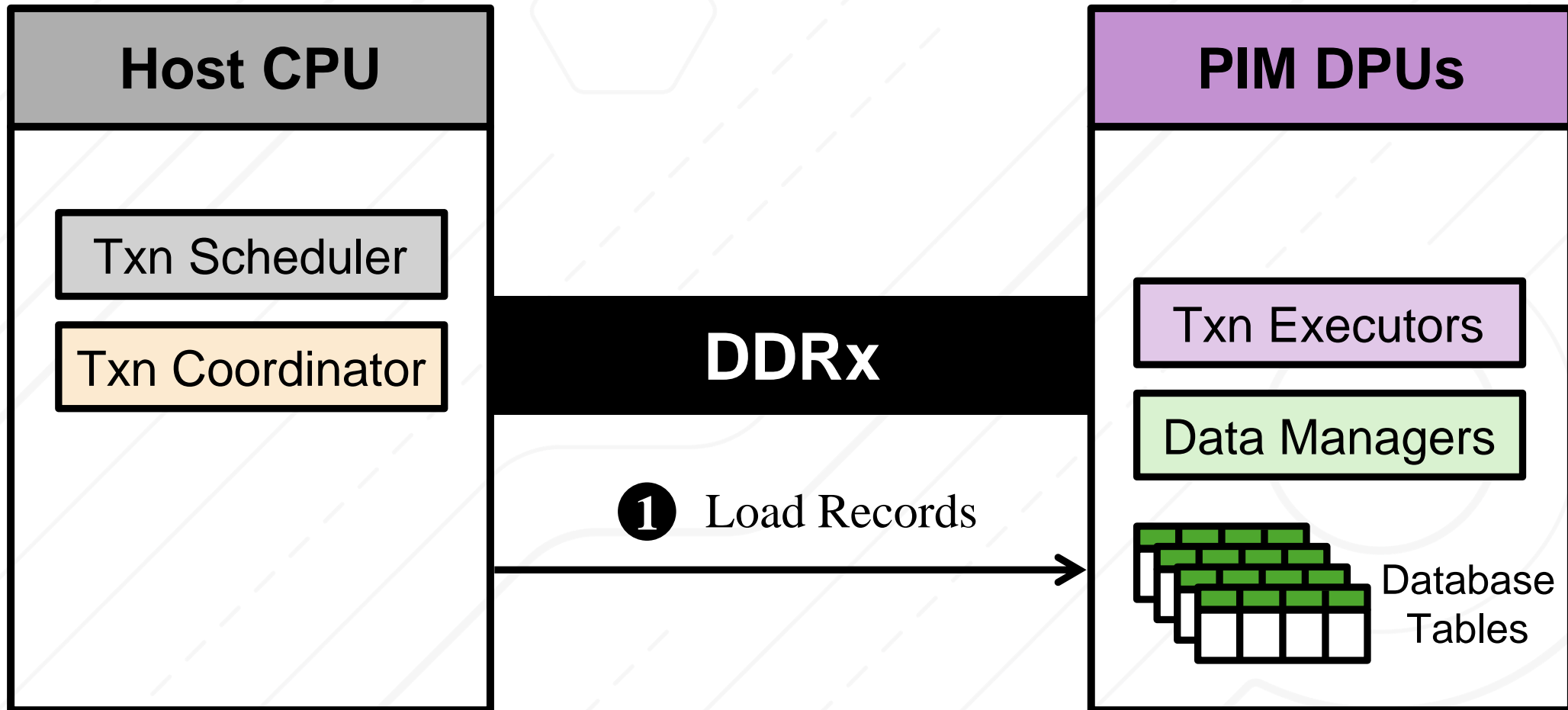
*Onyx: A Transaction Processing System with Real PIM Prototypes*





# Our Solution: Onyx

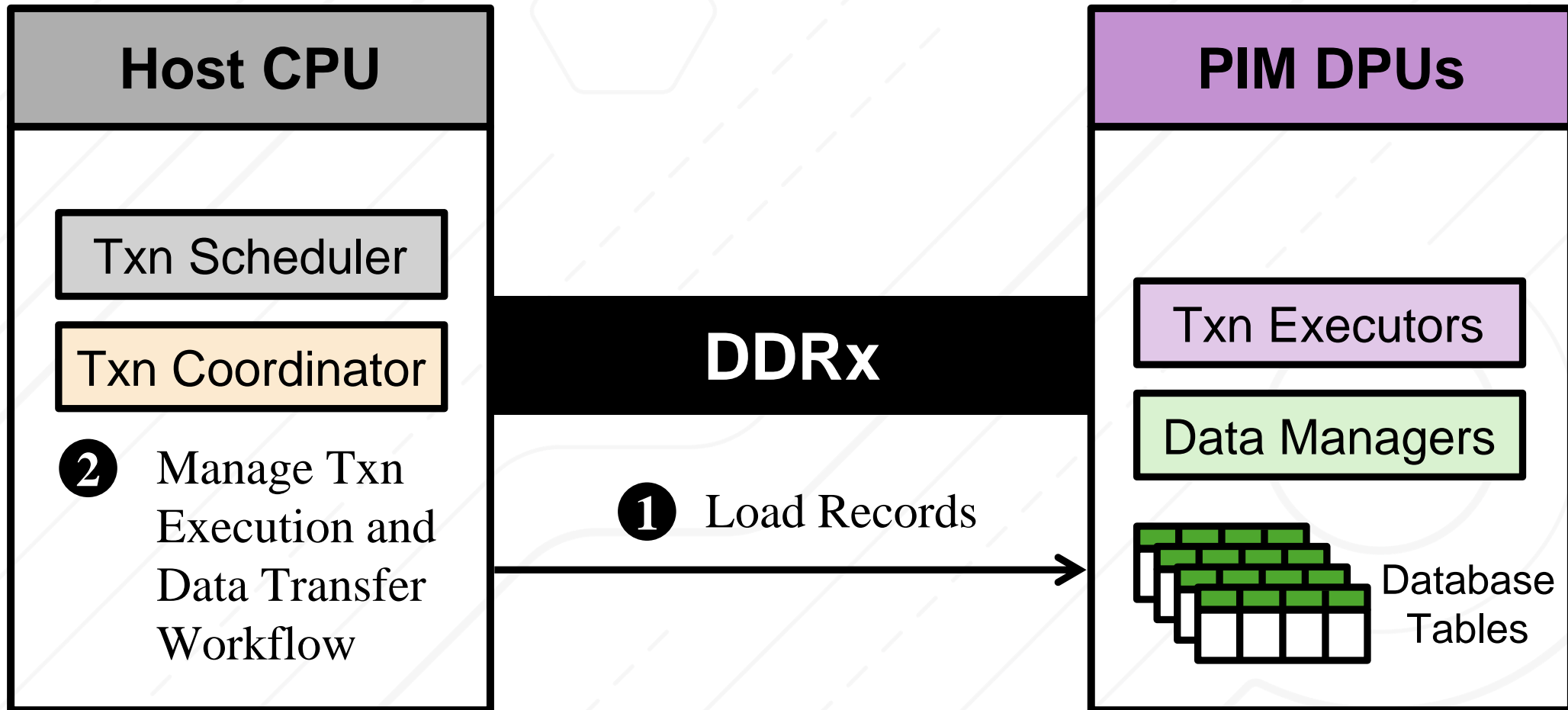
*Onyx: A Transaction Processing System with Real PIM Prototypes*





# Our Solution: Onyx

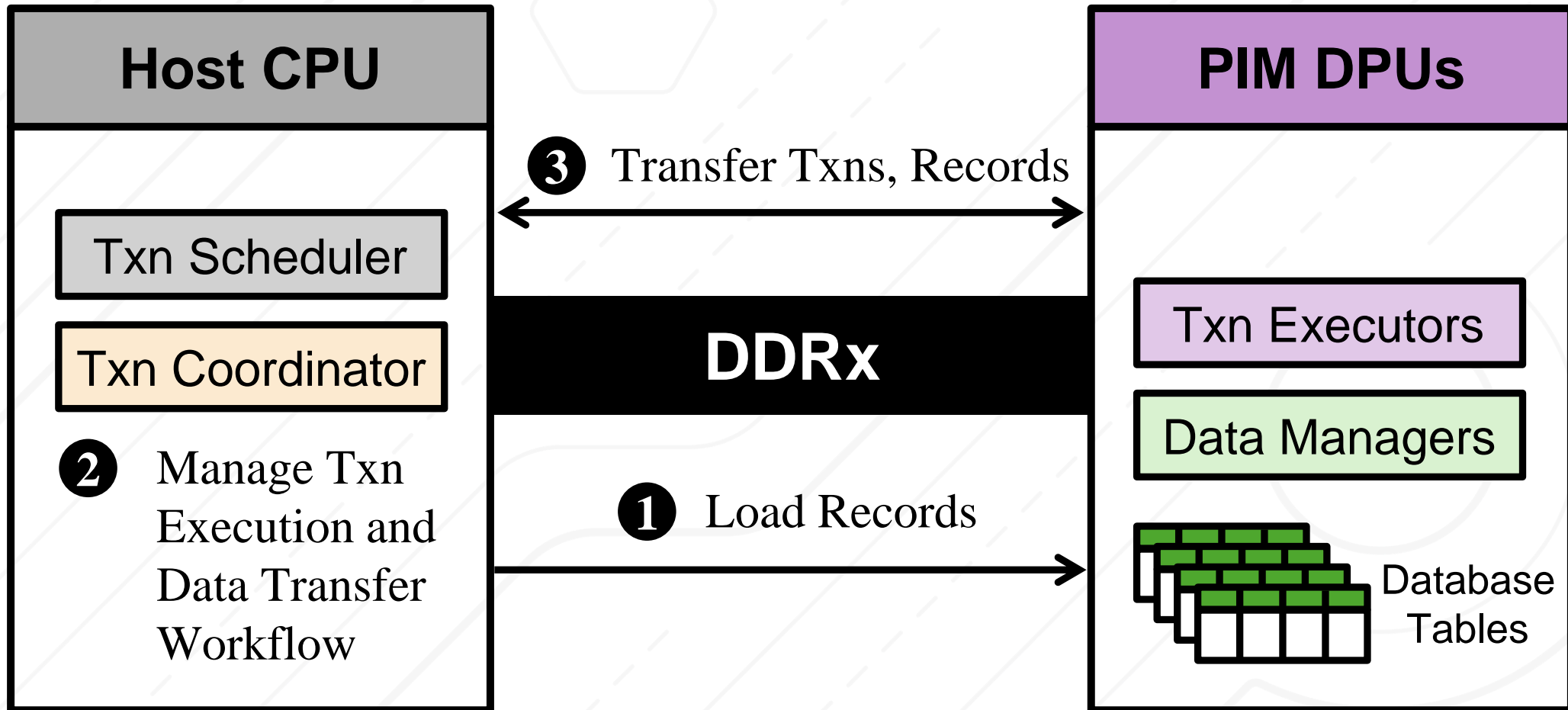
*Onyx: A Transaction Processing System with Real PIM Prototypes*





# Our Solution: Onyx

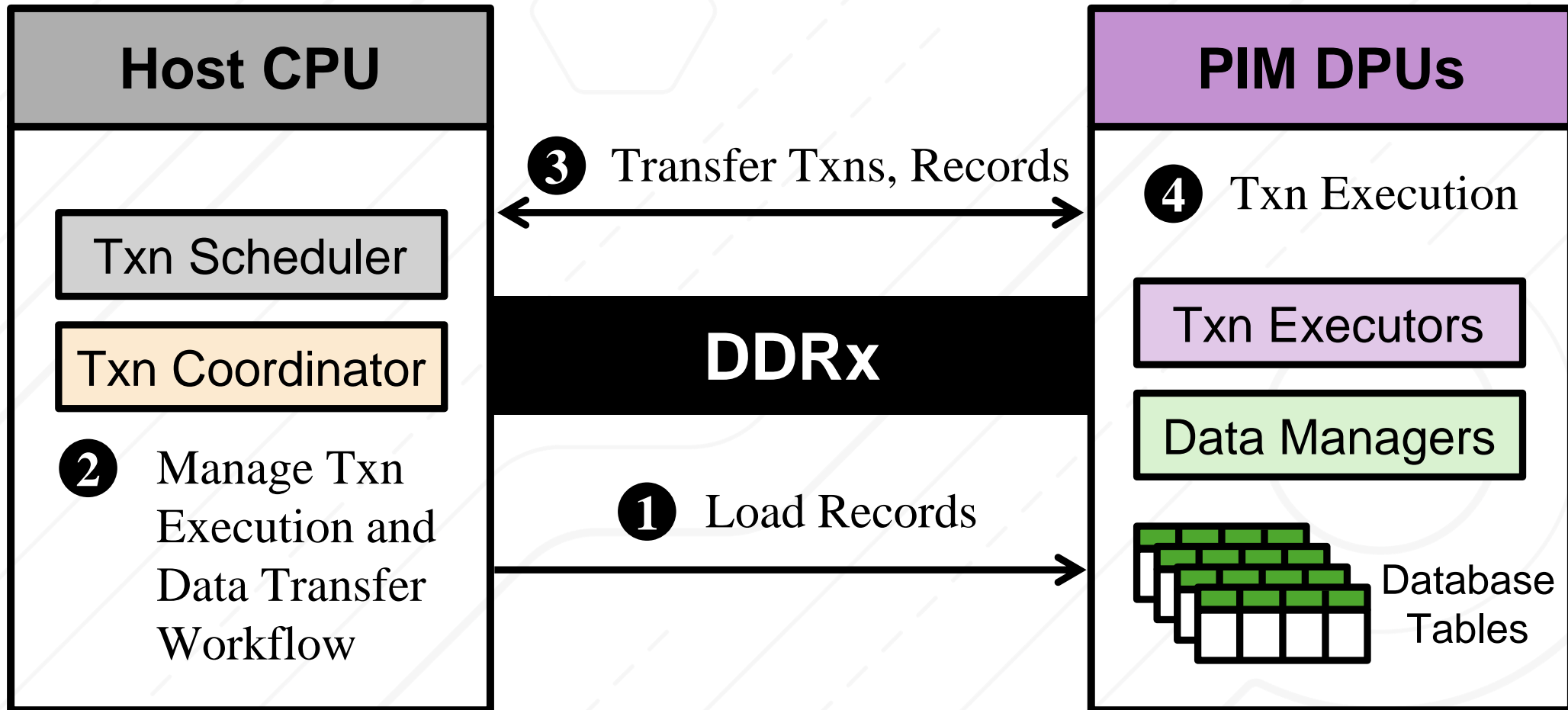
*Onyx: A Transaction Processing System with Real PIM Prototypes*





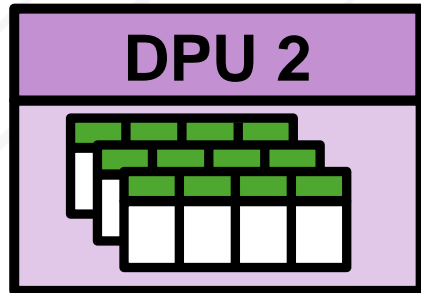
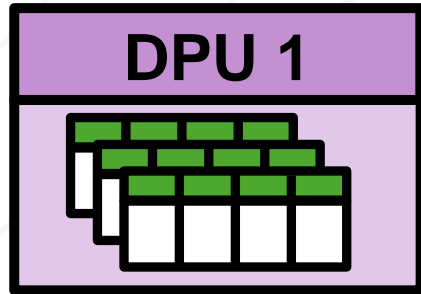
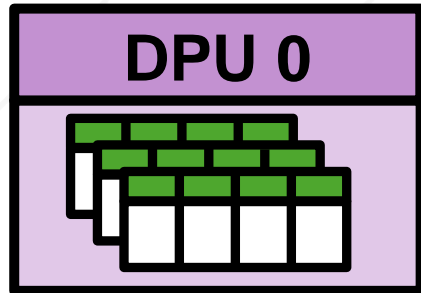
# Our Solution: Onyx

*Onyx: A Transaction Processing System with Real PIM Prototypes*





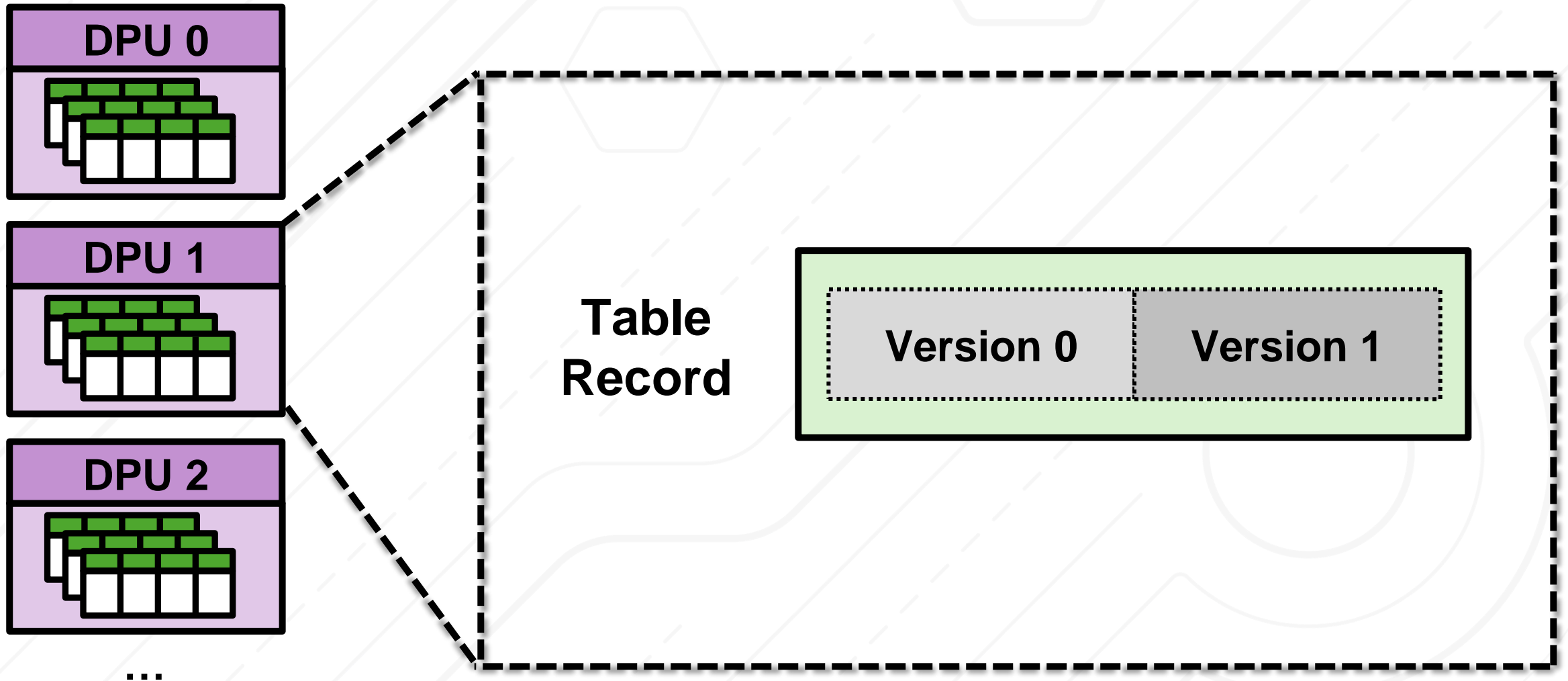
# Multi-version Storage Scheme



...

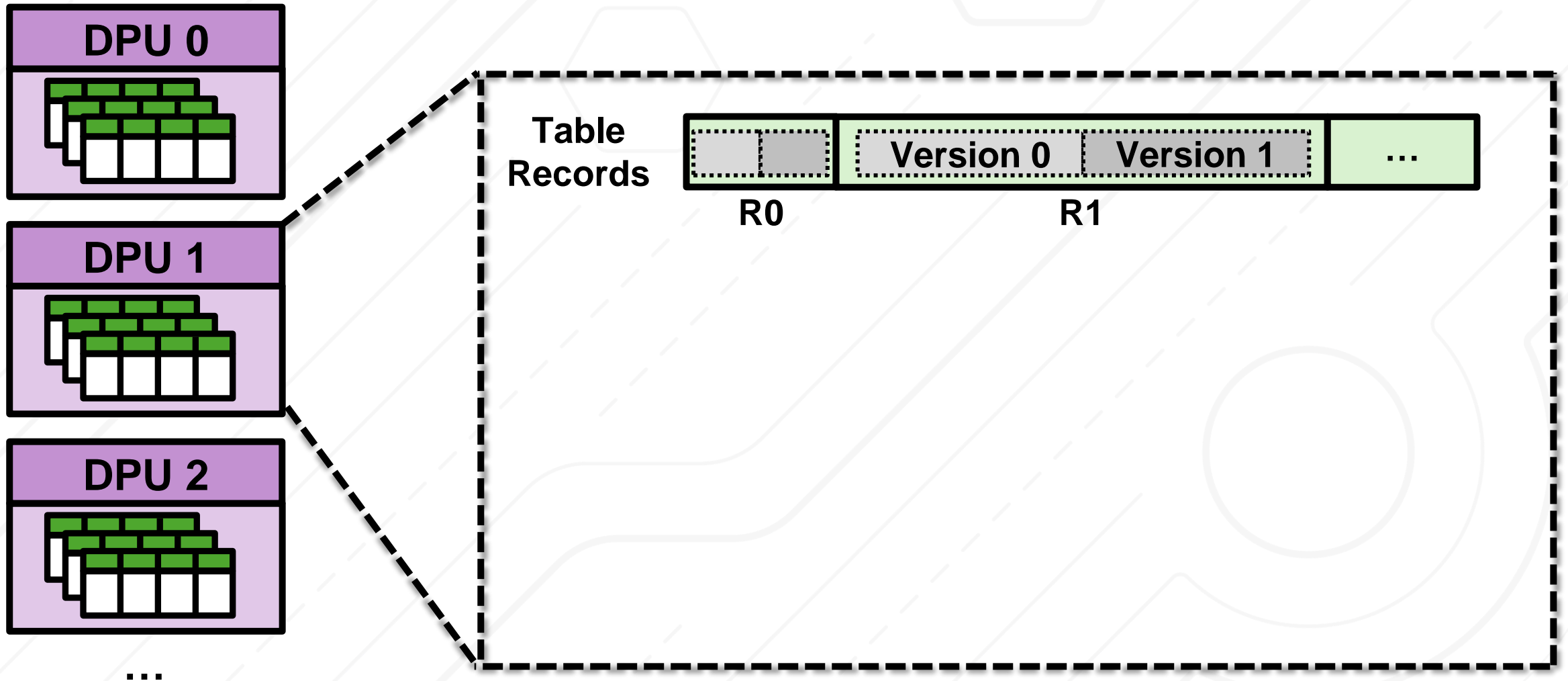


# Multi-version Storage Scheme



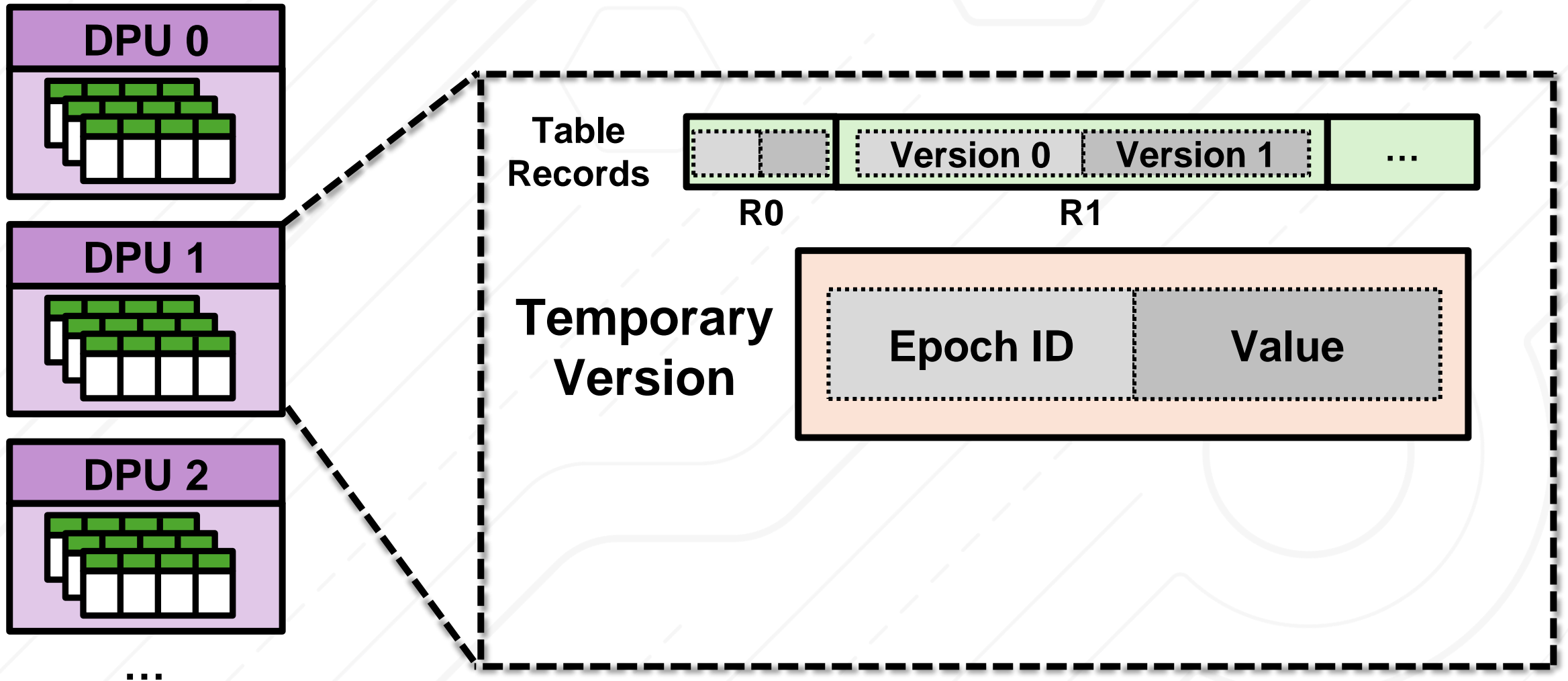


# Multi-version Storage Scheme



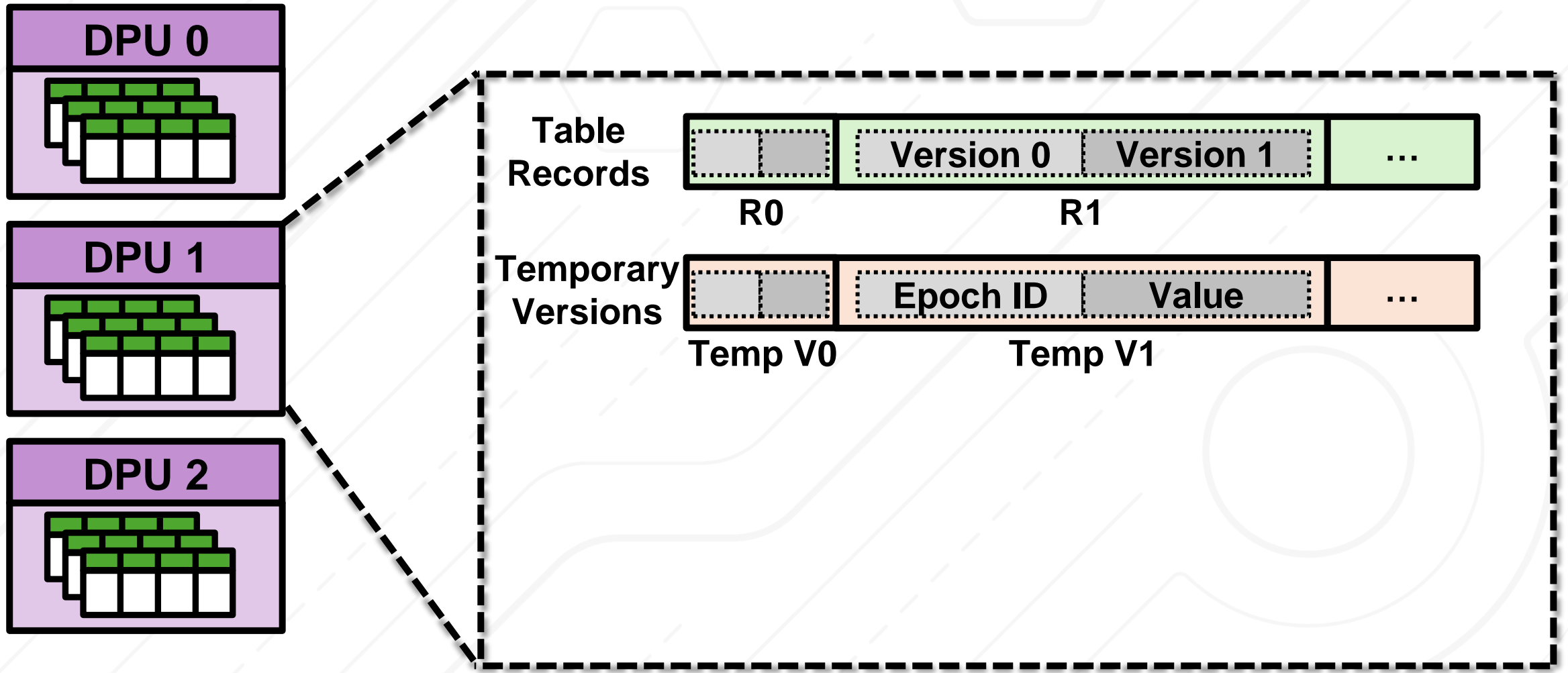


# Multi-version Storage Scheme



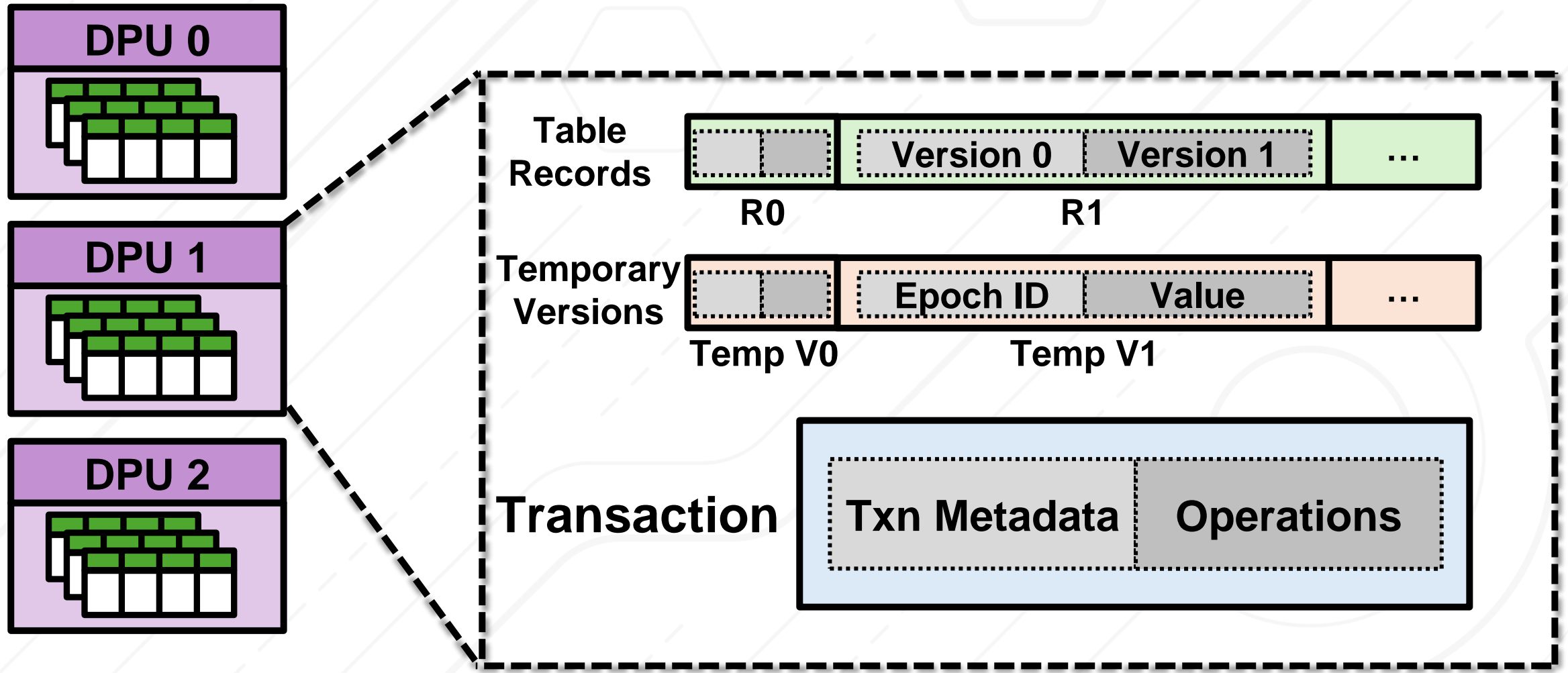


# Multi-version Storage Scheme



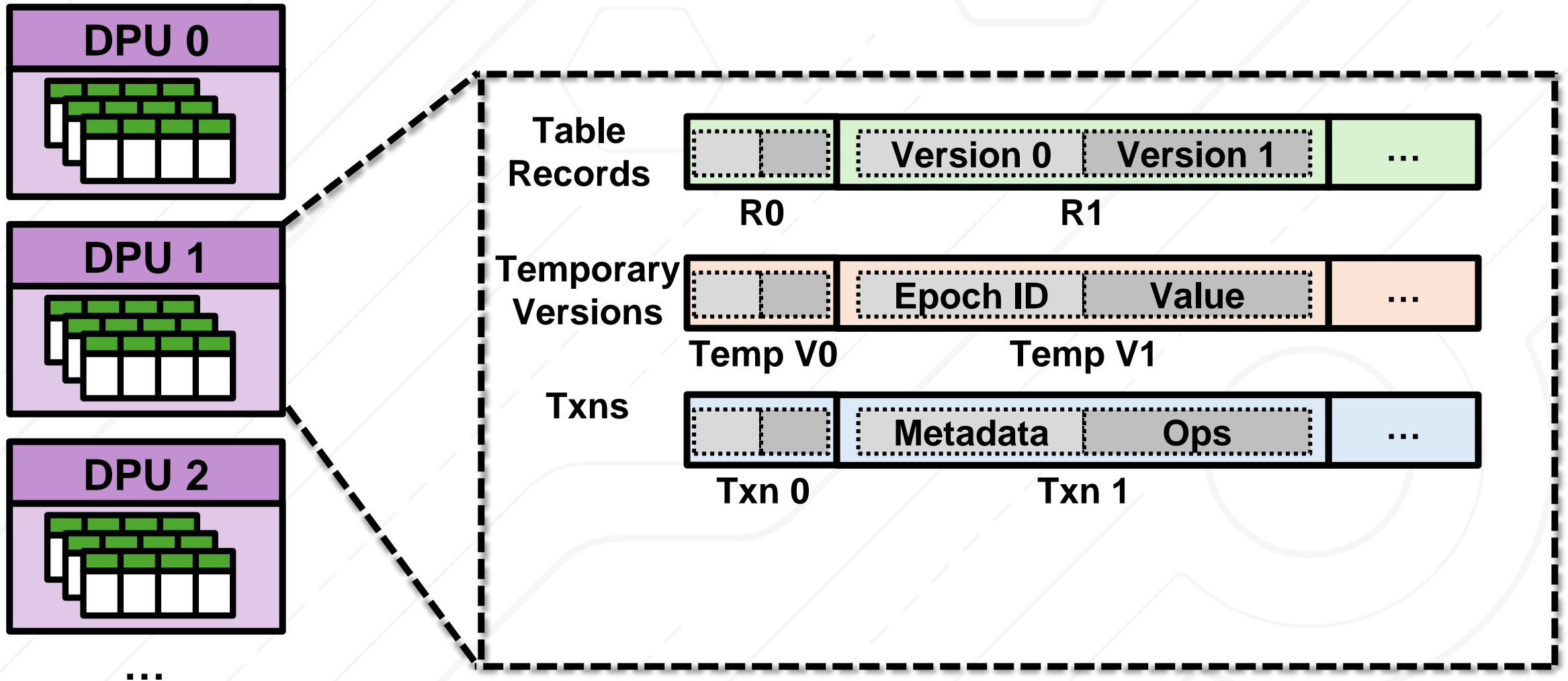


# Multi-version Storage Scheme



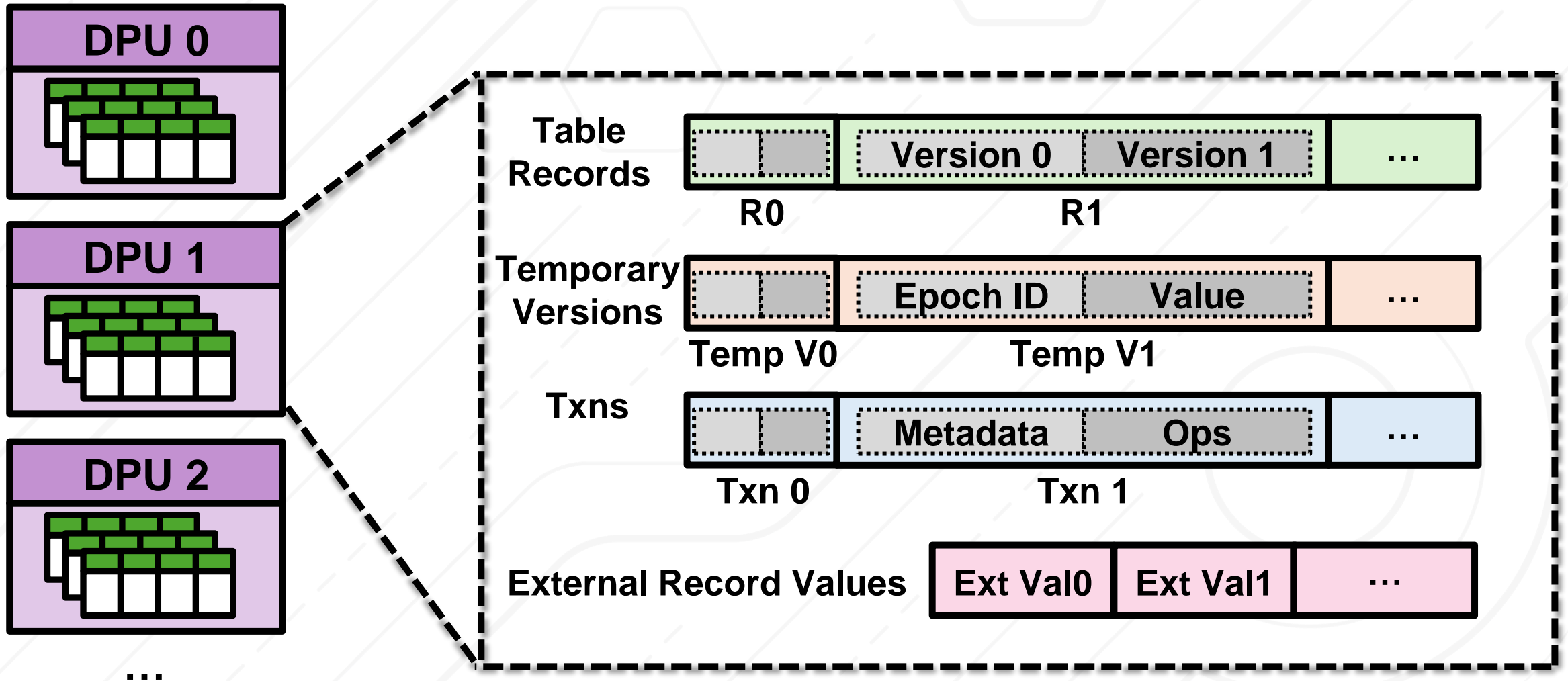


# Multi-version Storage Scheme





# Multi-version Storage Scheme





# Multi-versioning Initialization

**Txn 0** READ A, B  
WRITE A = A + B

**Txn 1** READ A, C, D  
WRITE A = A - D  
WRITE C = C + D

**Txn 2** READ B, D  
WRITE B = B + 1  
WRITE D = D \* B

Key	Col 0	Col 1	...
A	A.c0	A.c1	
C	C.c0	C.c1	
...			

Table 0

Tuples	Col 0	Col 1	...
B	B.c0	B.c1	
D	D.c0	D.c1	
...			

Table 1



# Multi-versioning Initialization

Hash Value

Hash

H 0

H 2

H 1

H 3

{Table 0, Key A}

{Table 0, Key C}

{Table 1, Key B}

{Table 1, Key D}

Key	Col 0	Col 1	...
A	A.c0	A.c1	
C	C.c0	C.c1	
...			

Table 0

Tuples	Col 0	Col 1	...
B	B.c0	B.c1	
D	D.c0	D.c1	
...			

Table 1



# Multi-versioning Initialization

<b>Keys</b>	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>
<b>Hash Value</b>	<b>H0</b>	<b>H1</b>	<b>H2</b>	<b>H3</b>
<b>DPU IDs</b>	<b>0</b>	<b>0</b>	<b>1</b>	<b>1</b>
<b>Record IDs</b>	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>

<b>Key</b>	<b>Col 0</b>	<b>Col 1</b>	<b>...</b>
<b>A</b>	<b>A.c0</b>	<b>A.c1</b>	
<b>C</b>	<b>C.c0</b>	<b>C.c1</b>	
<b>...</b>			

...

Table 0

<b>Tuples</b>	<b>Col 0</b>	<b>Col 1</b>	<b>...</b>
<b>B</b>	<b>B.c0</b>	<b>B.c1</b>	
<b>D</b>	<b>D.c0</b>	<b>D.c1</b>	
<b>...</b>			

...

Table 1



# Multi-versioning Initialization

**Txn 0** **READ A, B**  
**WRITE A = A + B**

**Txn 1** **READ A, C, D**  
**WRITE A = A - D**  
**WRITE C = C + D**

**Txn 2** **READ B, D**  
**WRITE B = B + 1**  
**WRITE D = D \* B**

*Write to temporary version*

**R0** Write Set: {Txn 0, Txn 1}  
TV0

**R1** Write Set: {Txn 2}

**R2** Write Set: {Txn 1}

**R3** Write Set: {Txn 2}

*Directly write to table*

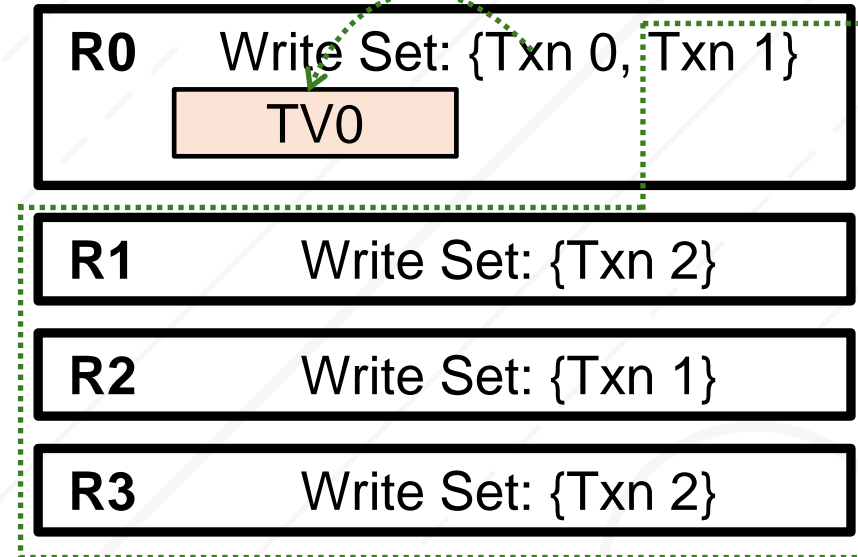
Version 0      Version 1



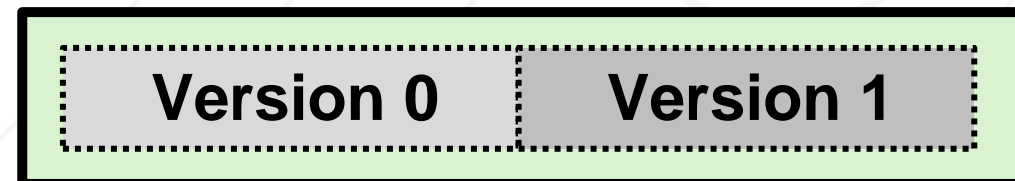
# Multi-versioning Initialization

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

*Write to temporary version*



*Directly write to table*





# Multi-versioning Initialization

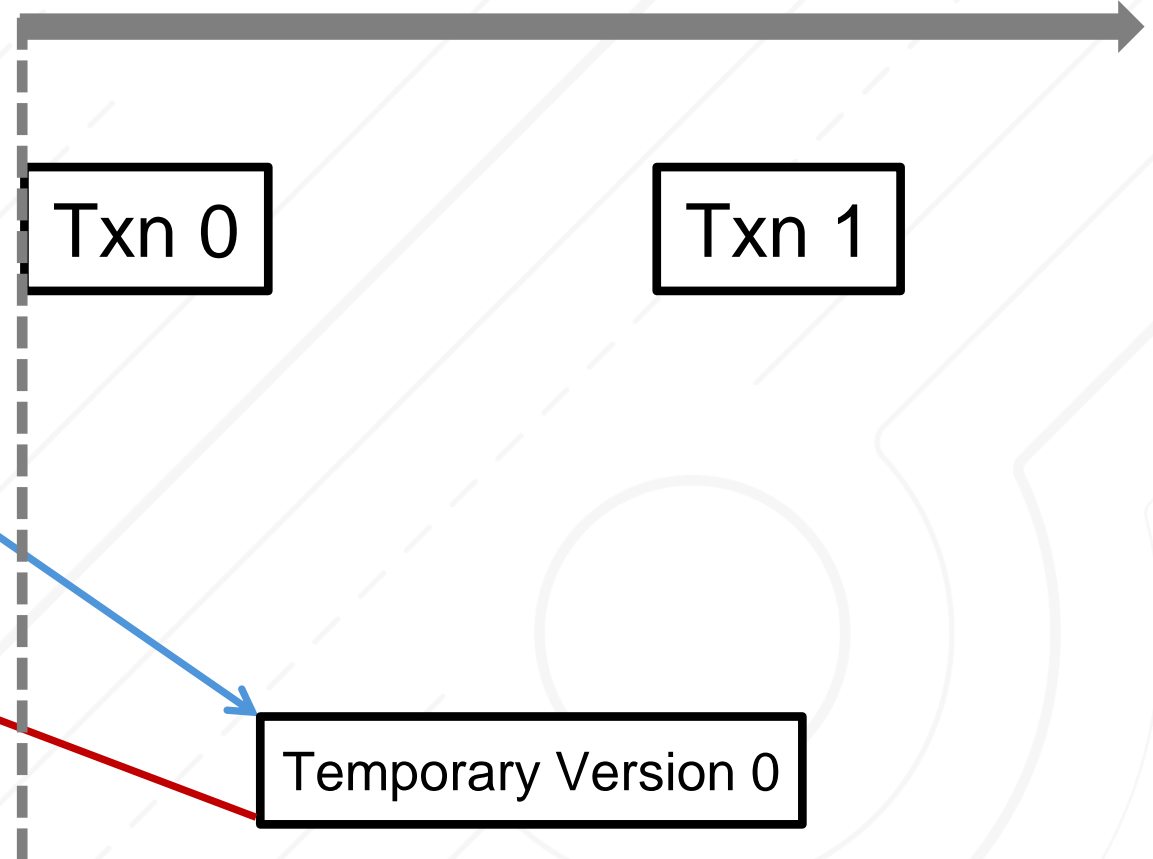
<b>Txn ID</b>	<b>Read Locations</b>	<b>Write Locations</b>
<b>0</b>	<b>R0.V0 R1.V0</b>	<b>TV0</b>
<b>1</b>	<b>TV0 R2.V0 R3.V0</b>	<b>R0.V1 R2.V1</b>
<b>2</b>	<b>R1.V0 R3.V0</b>	<b>R1.V1 R3.V1</b>



# Multi-versioning Initialization

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

Transaction Execution Timeline

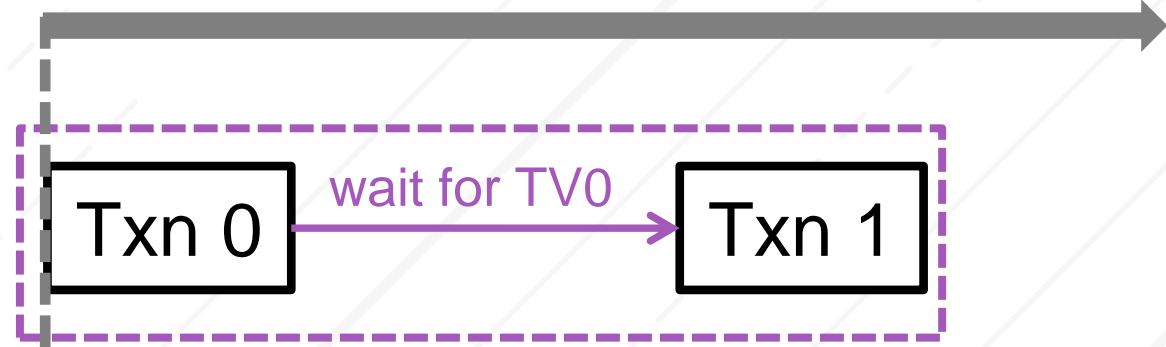




# Multi-versioning Initialization

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

Transaction Execution Timeline



*Txn 0 write*

*Txn 1 read*

*Version dependency:  
Txn 1 must wait for Txn 0  
to generate TV0*

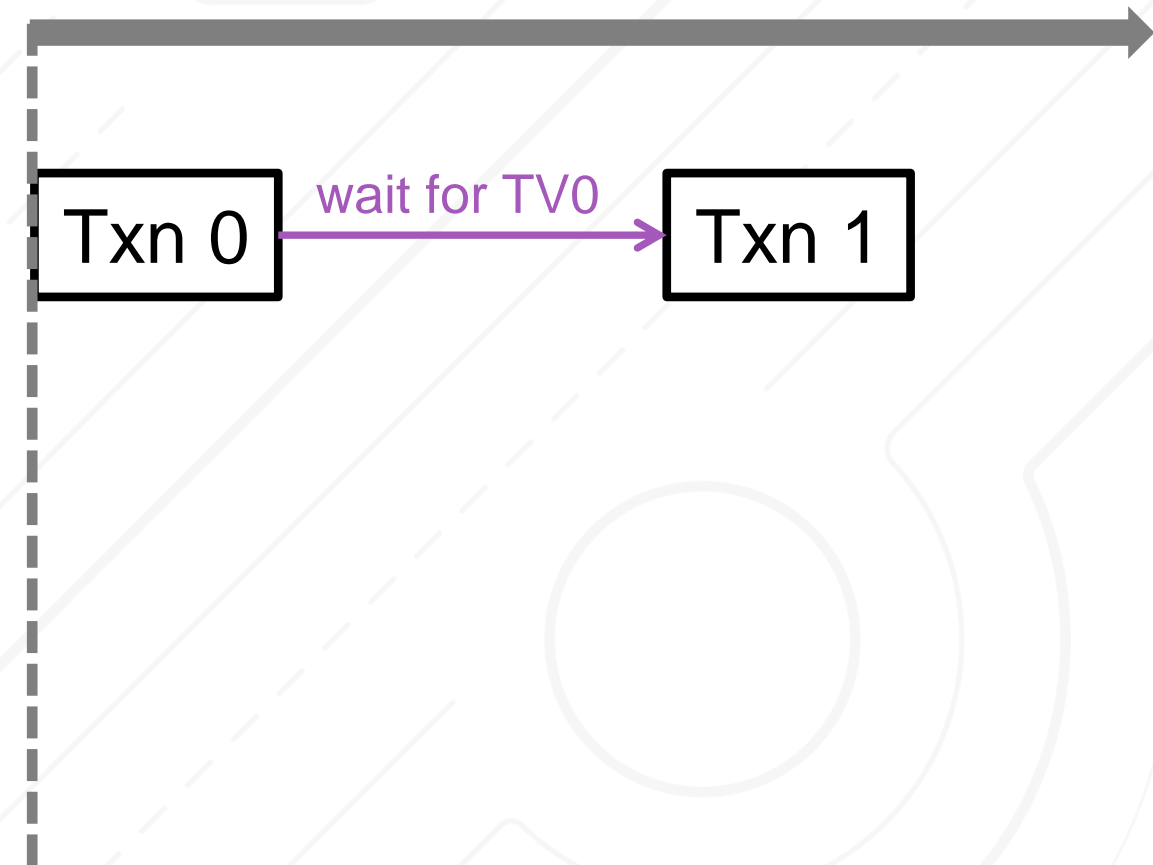
Temporary Version 0



# Multi-versioning Initialization

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

Transaction Execution Timeline



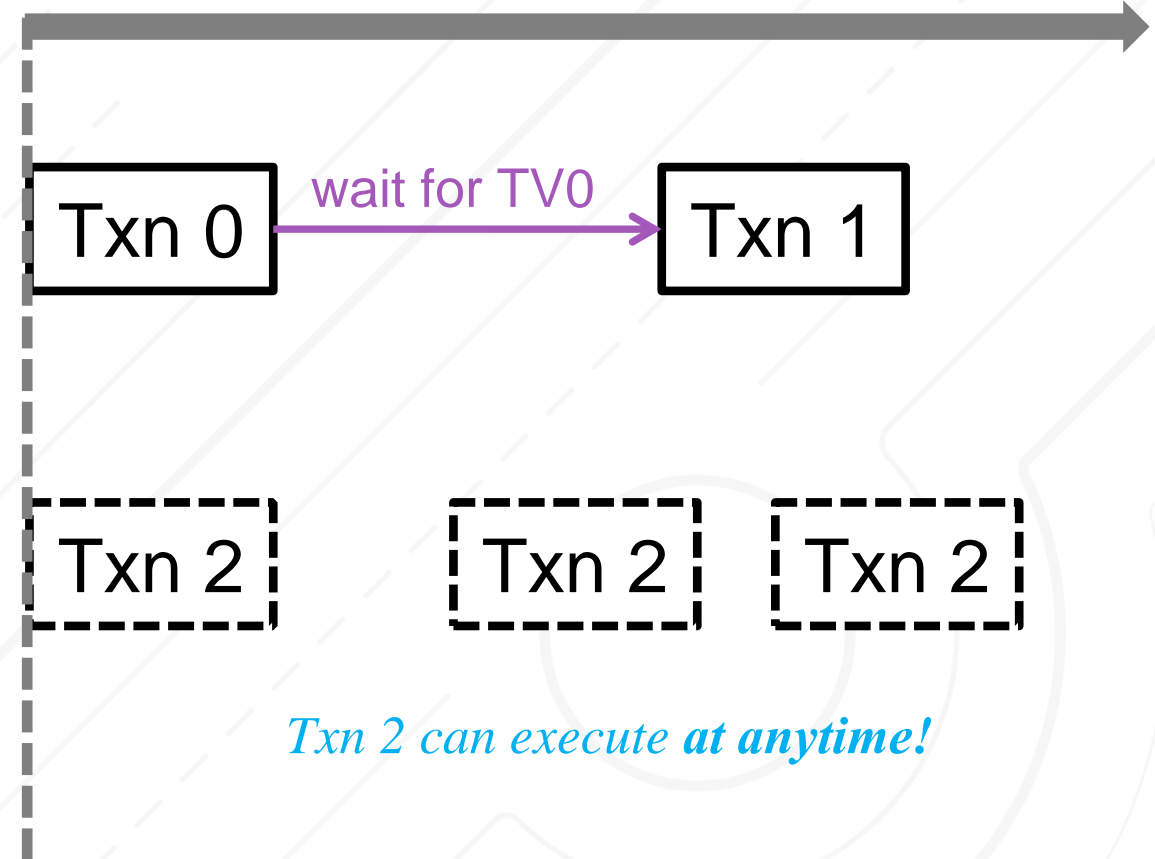
*No Version Dependency between Txn 2 and Txn 0 / Txn 1*



# Multi-versioning Initialization

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

Transaction Execution Timeline



*No Version Dependency between Txn 2 and Txn 0 / Txn 1*

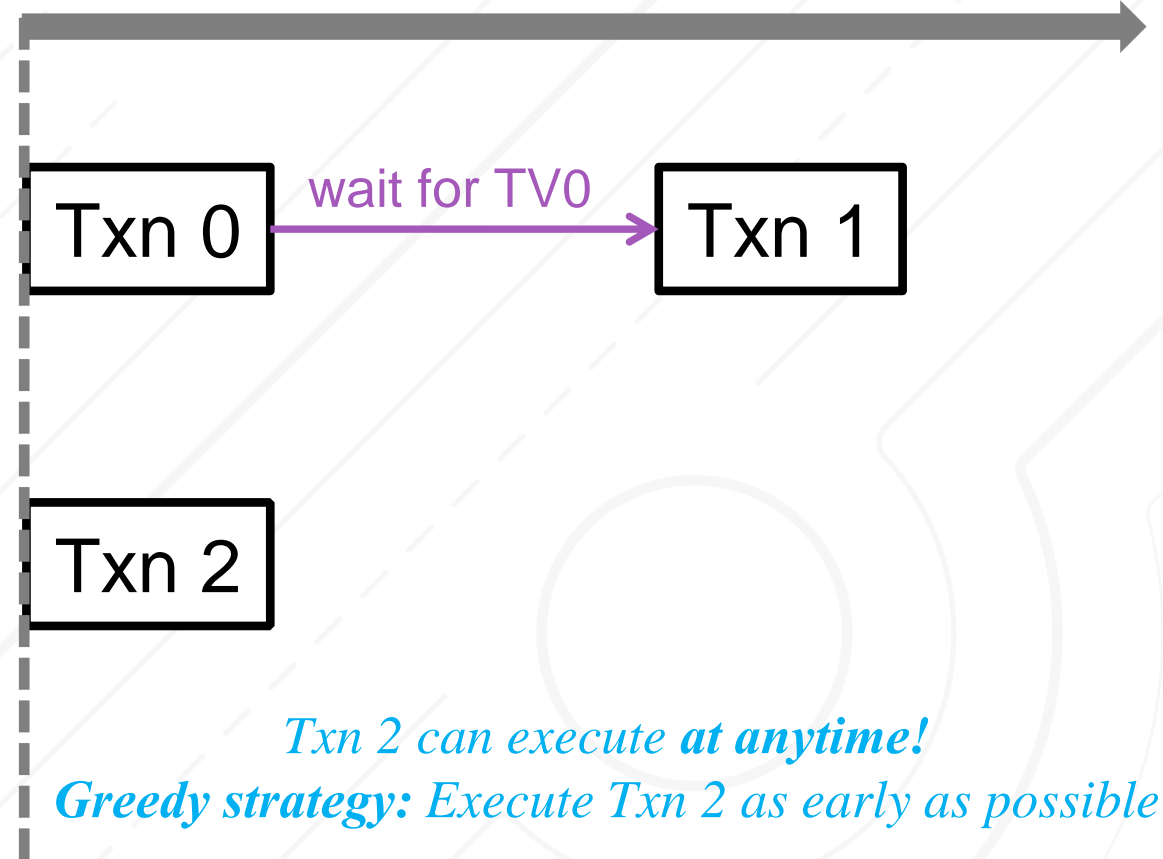


# Multi-versioning Initialization

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

*No Version Dependency between Txn 2 and Txn 0 / Txn 1*

Transaction Execution Timeline

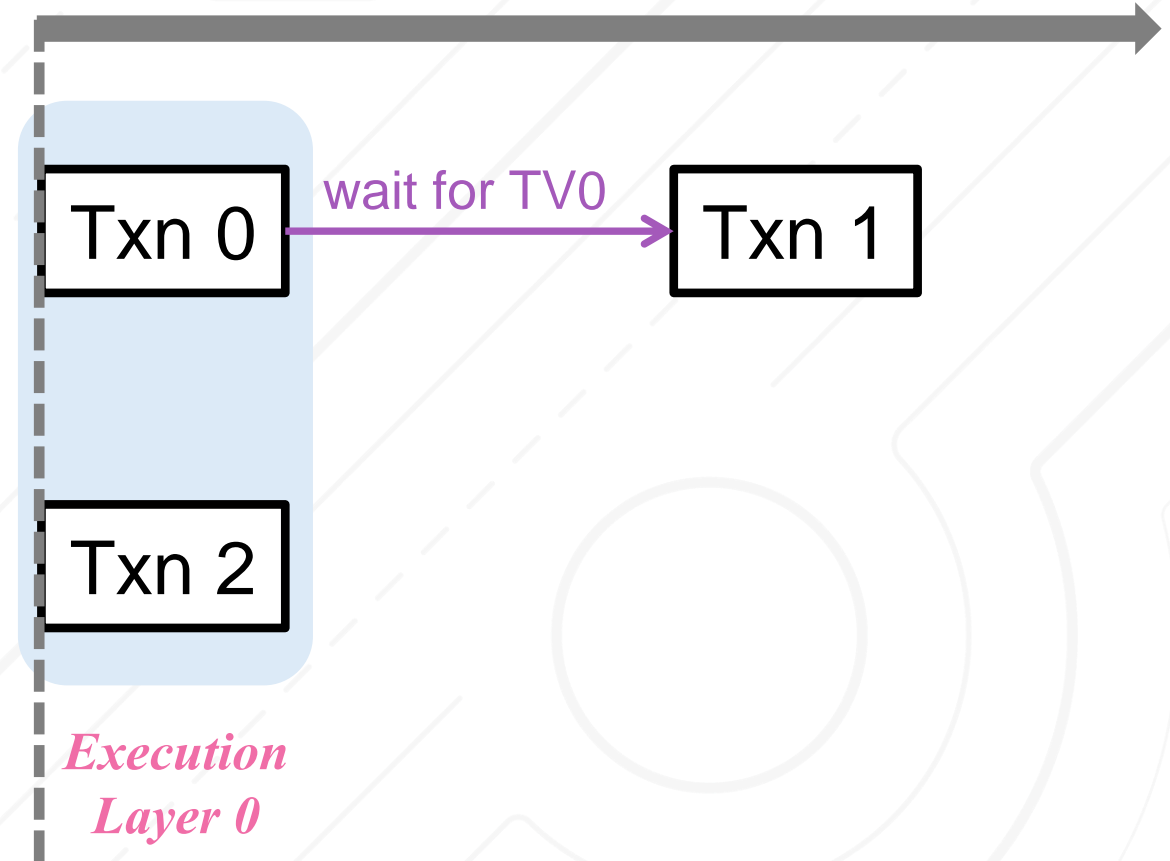




# Multi-versioning Initialization

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

Transaction Execution Timeline

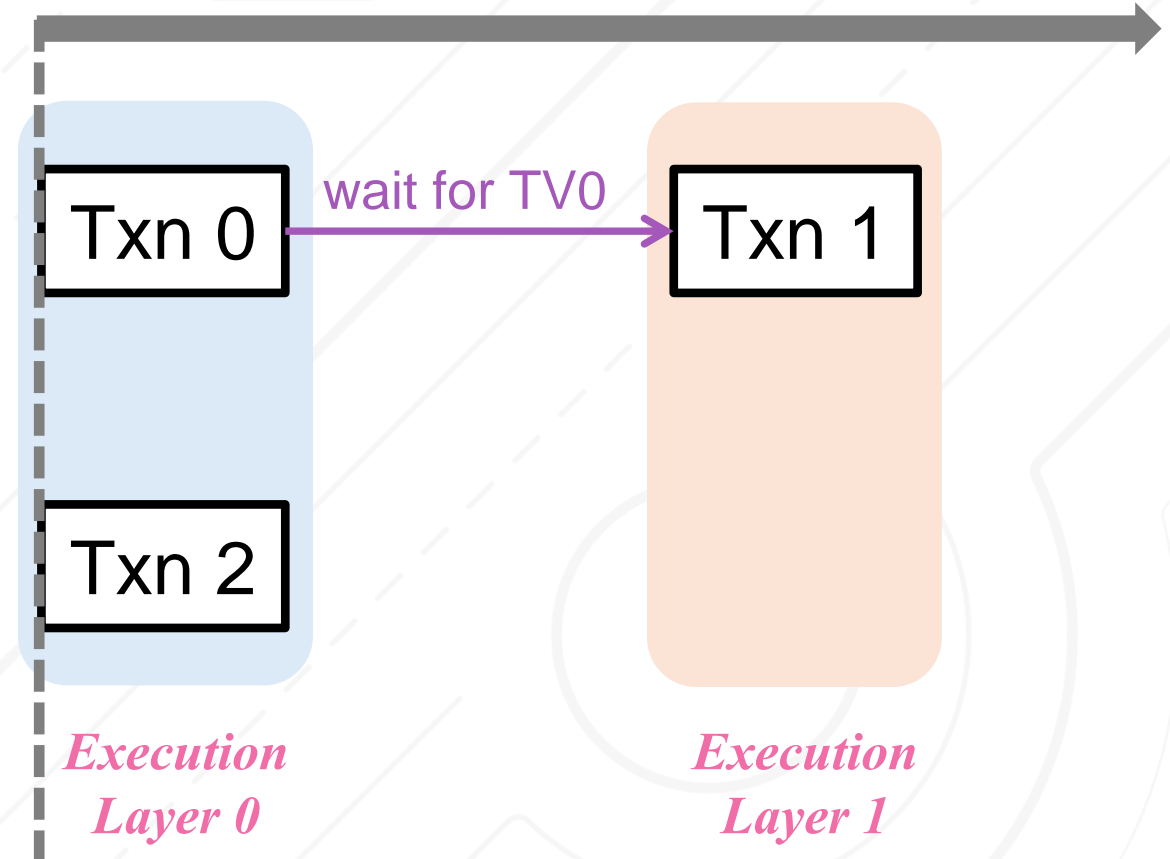




# Multi-versioning Initialization

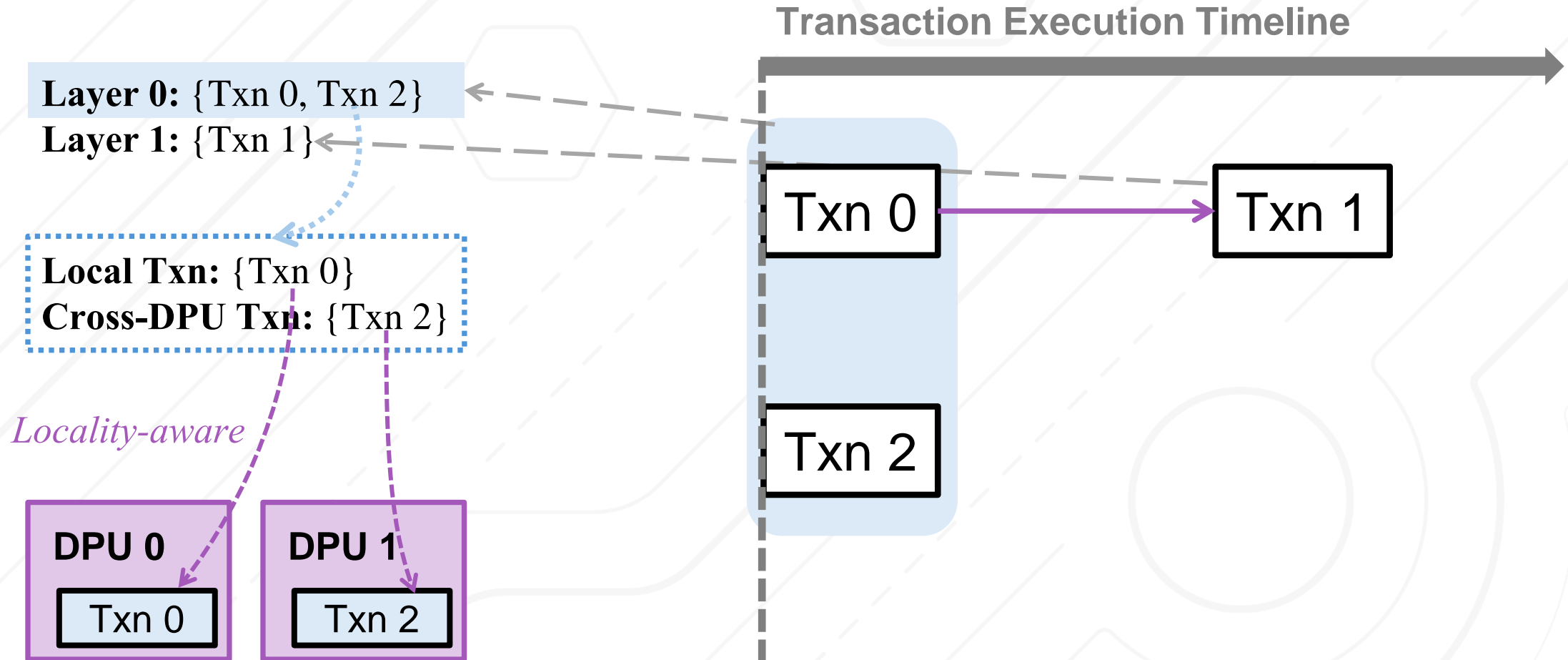
Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

Transaction Execution Timeline





# Locality-Aware Execution

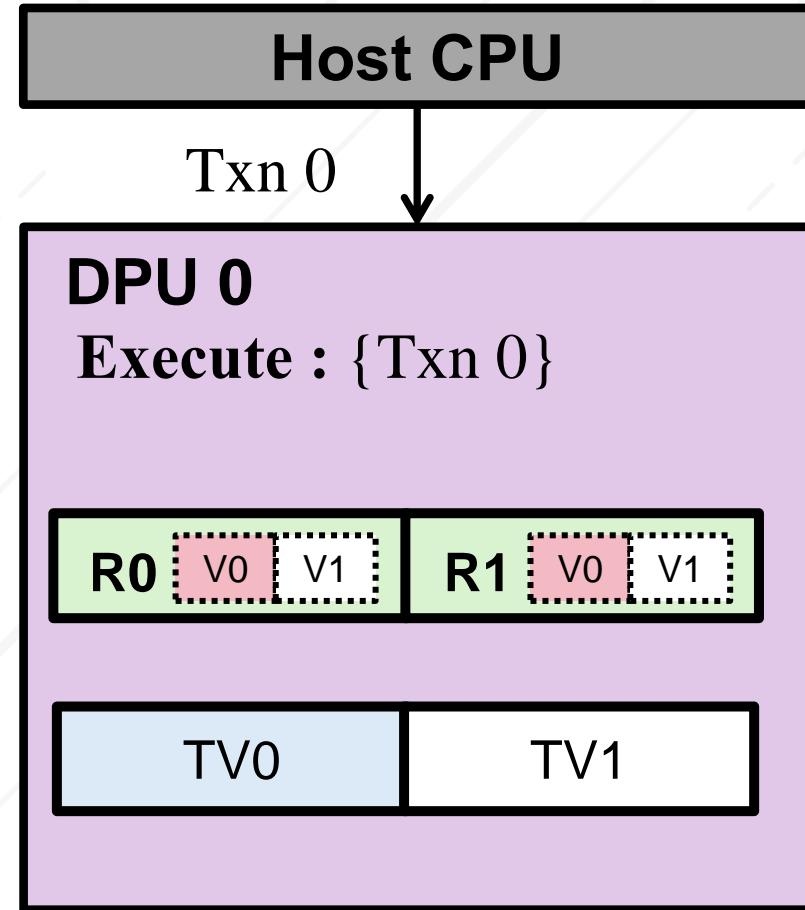




# Locality-Aware Execution

## *Local Txn Execution*

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1



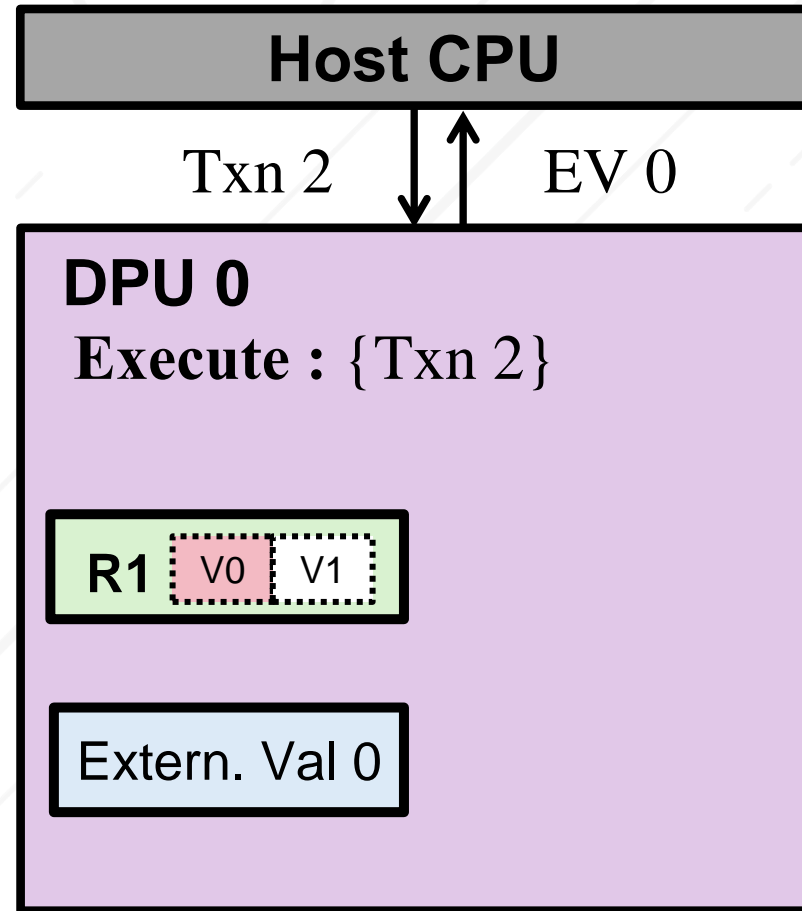
READ: WRITE:



# Locality-Aware Execution

## *Cross-DPU Txn Read Phase*

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1



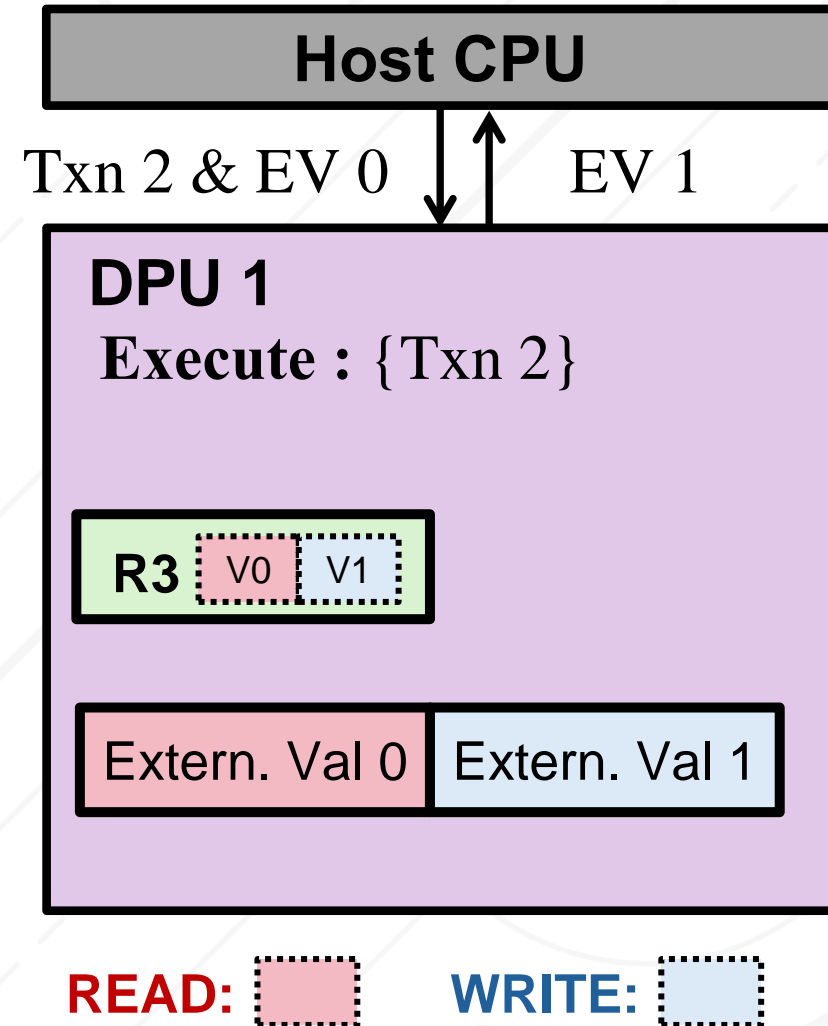
READ: WRITE:



# Locality-Aware Execution

## *Cross-DPU Txn Compute Phase*

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1

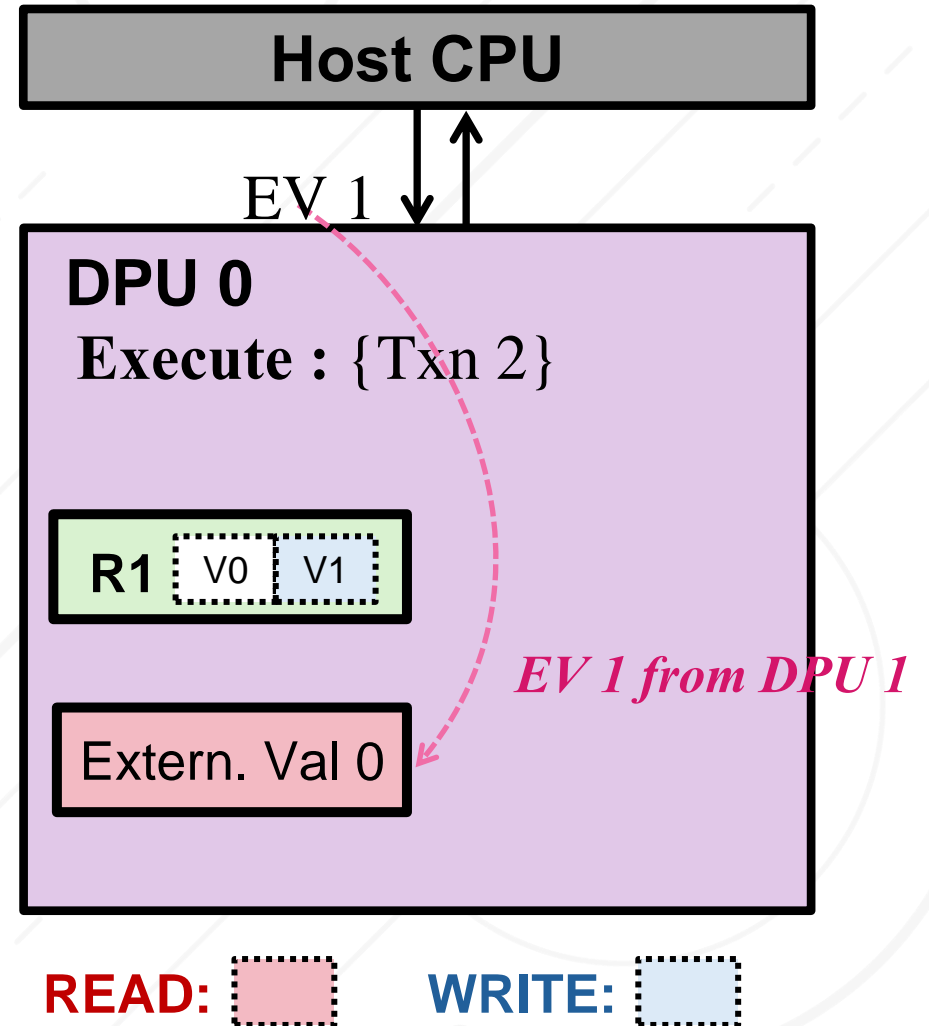




# Locality-Aware Execution

## Cross-DPU Txn Write Phase

Txn ID	Read Locations	Write Locations
0	R0.V0 R1.V0	TV0
1	TV0 R2.V0 R3.V0	R0.V1 R2.V1
2	R1.V0 R3.V0	R1.V1 R3.V1





# Evaluation

## ➤ Platform

- A Real PIM Prototype
  - 8 UPMEM PIM DIMMs
  - 2 Intel Xeon Silver 4216 CPUs
  - 192 GB of DRAM memory

## ➤ Comparisons

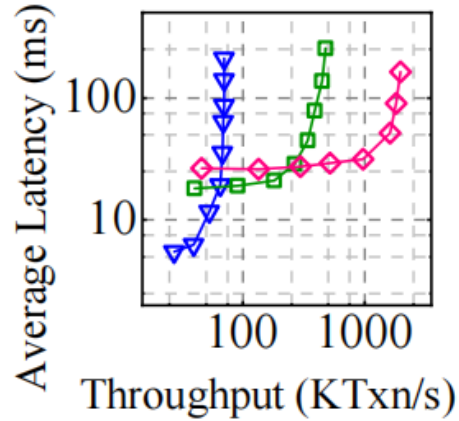
- Single-version system: Calvin<sup>[Sigmod' 12]</sup>
- Multi-version system: Caracal<sup>[SOSP' 21]</sup>

## ➤ Workloads

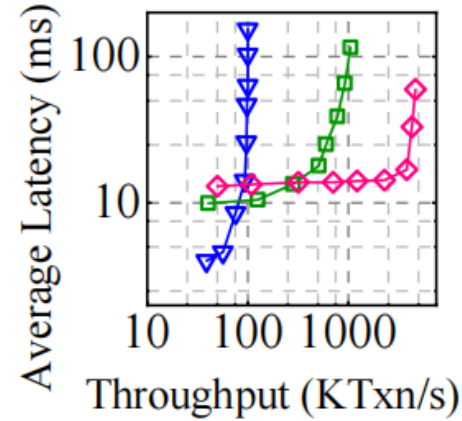
- YCSB benchmark: single database table containing 1 million records
- TPC-C benchmark: complex OLTP environment with 9 database tables



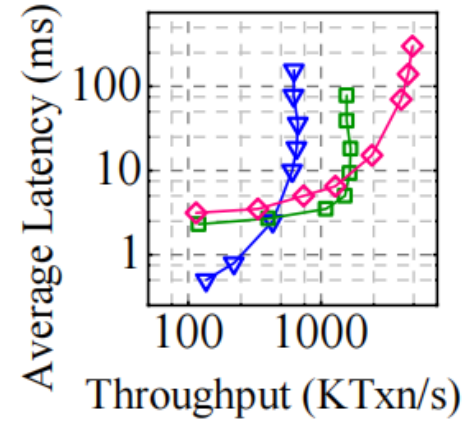
# End-to-End Performance



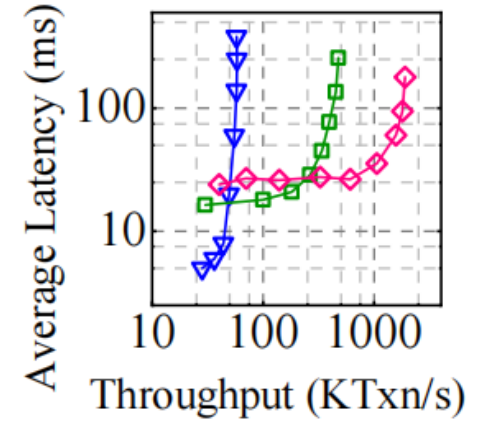
(a) YCSB-A



(b) YCSB-B



(c) YCSB-C

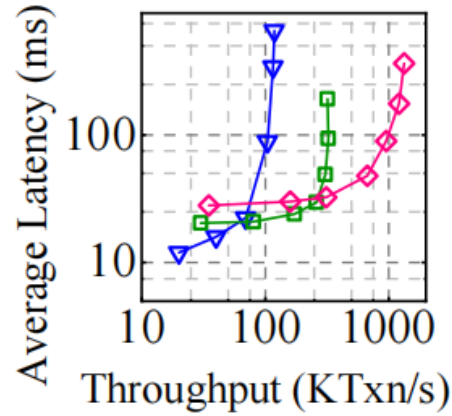


(d) YCSB-F

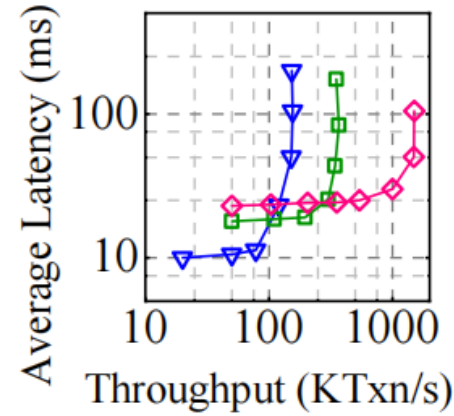
By fully leveraging **PIM bandwidth & parallelism** via orchestrated workflow, Onyx improves the throughput by **4.62x** on average



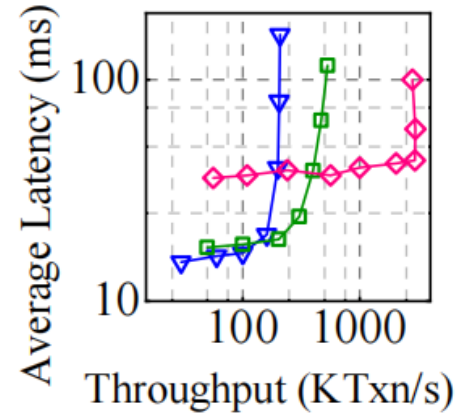
# End-to-End Performance



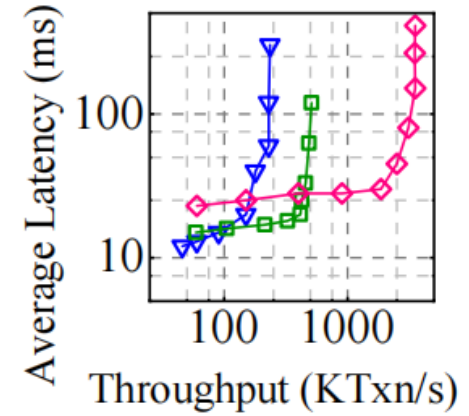
(e) TPC-C-W1



(f) TPC-C-W4



(g) TPC-C-W16



(h) TPC-C-W64

Onyx improves the throughput by **8.13x** on average



# Conclusion

- Existing in-memory OLTP systems suffer from severe data movements
  - *The **CPU-Memory bandwidth** fundamentally limits throughputs*
- **Onyx**: an efficient transaction processing system
  - *Leveraging a real-world PIM prototype*
  - *Multi-versioning Initialization*
  - *Locality-aware Execution*

*Thanks! Q&A*