

## Introduction

**Declarative programming** focuses on “WHAT” to achieve rather than “HOW”

Users			
UserID	UserName	UserEmail	Country
101	Alice	alice@example.com	USA
102	Bob	bob@example.com	USA
103	Eve	eve@example.com	Canada

**WHAT** `SELECT UserID FROM Users WHERE Country='USA';` **HOW**

**Advanced approach: Logic programming (Datalog)**

Datalog rules to compute Transitive Closure (TC) of a relation

$TC(x, y) :- Edge(x, y).$   
 $TC(x, z) :- TC(x, y), Edge(y, z).$

Operationalized as a **fixed-point iteration** using  $F_G$

$$F_G(T) \triangleq G \cup \Pi_{1,2}(\rho_{0/1}(T) \bowtie_1 G)$$

Datalog rules compiled down to **iterative relational algebra** operators

HPC paved the way for parallelizing iterative relational algebra where current Datalog engines target only multi-core CPUs and single-GPU systems

Multi-threaded	Distributed (Apache Spark)	Multi-node Multi-threaded	Single-GPU	Multi-node Multi-GPUs
Soufflé	RDFox	SLOG	GPUJoin	
LogicBlox	Radlog	PRAM	GPULog	
Nemo	BigDatalog		GPUDatalog	

MNMGDatalog is the first multi-node, multi-GPU (MNMG) Datalog engine

**Highest performant Datalog engine**

**Single-GPU:** Up to 7× speedup over GPULog

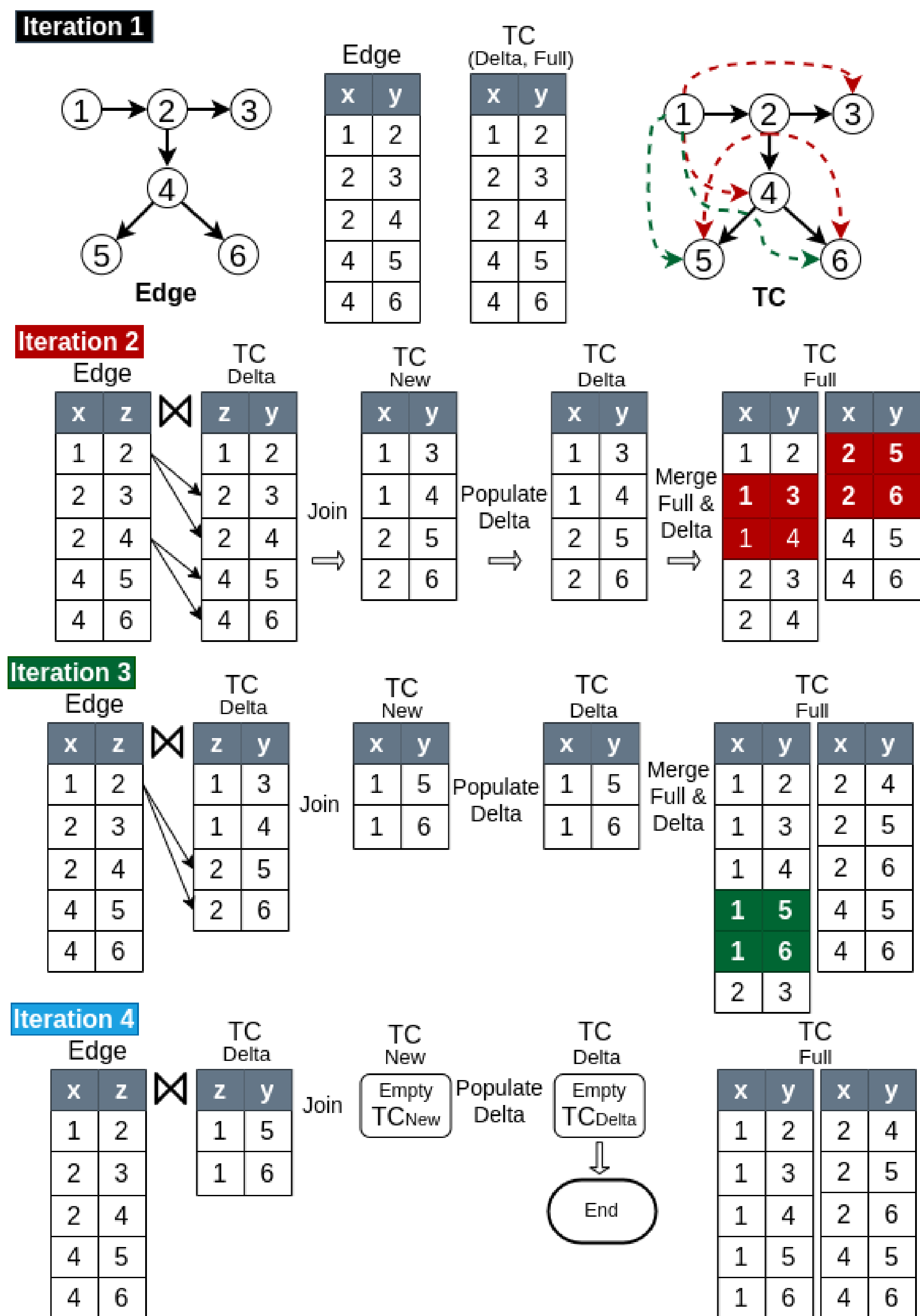
**Multi-threaded:** Up to 33× over Soufflé

**Distributed:** Up to 31.9× speedup over SLOG spanning 32 GPUs

## Challenges

**Iterative relational algebra** on MNMG systems is challenging due to high communication overhead, synchronization cost, and repeated materialization.

Iterations in Transitive Closure (TC) Computation

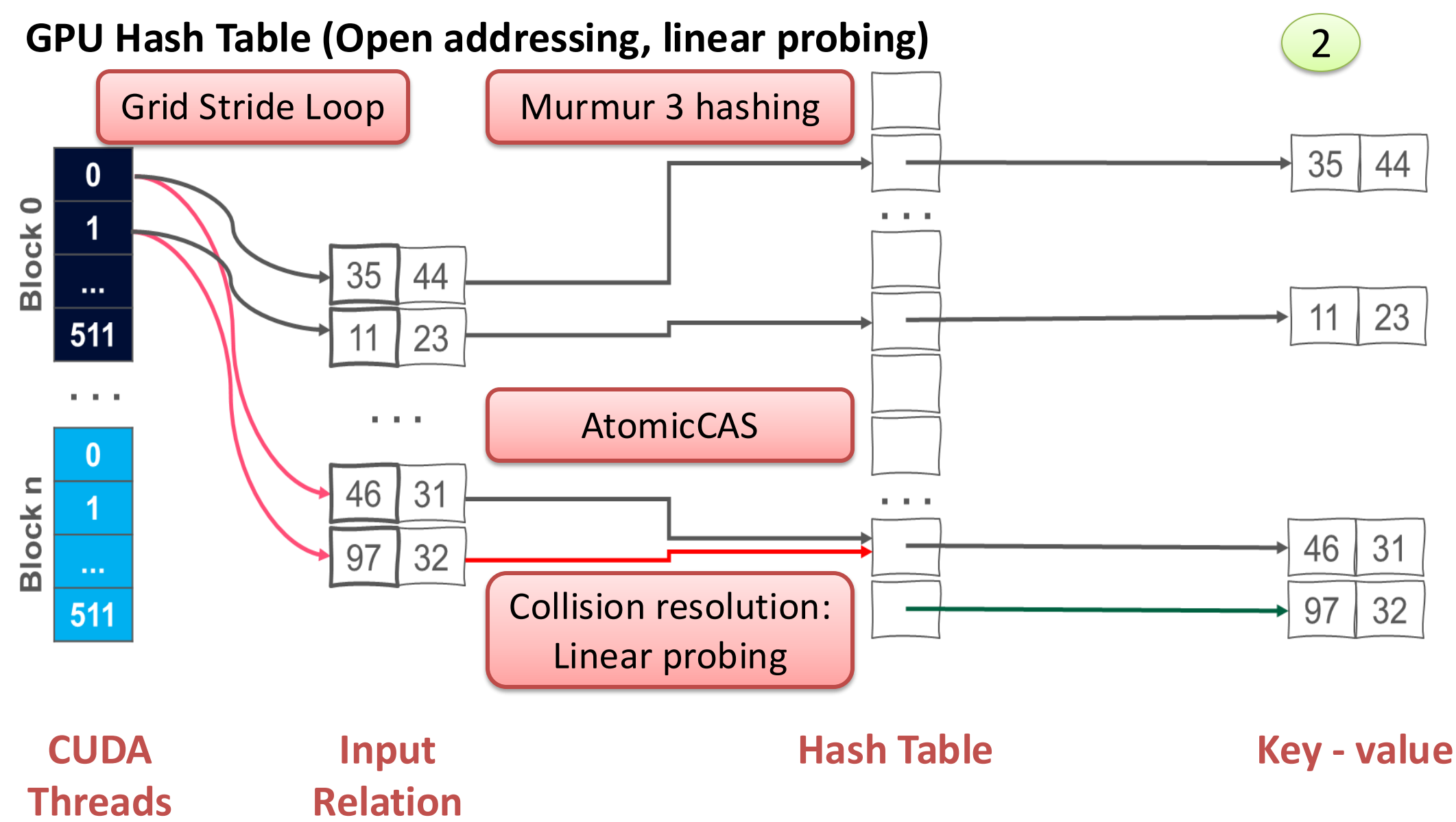


Requirements for MNMGDatalog Engine

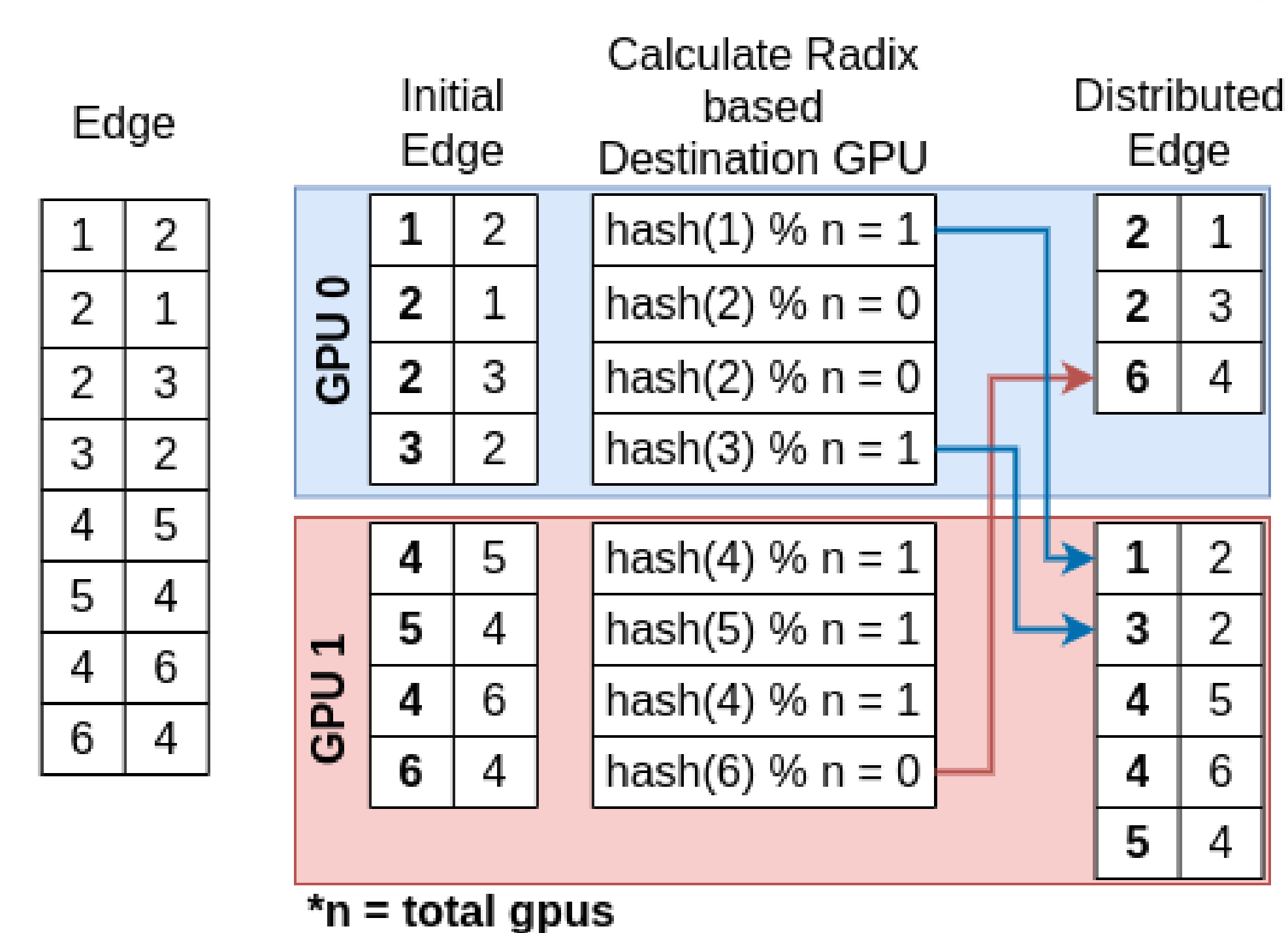
1 Workload partitioning	2 Data representation
3 Efficient communication	4 Tuple materialization

## Implementation of MNMGDatalog Engine

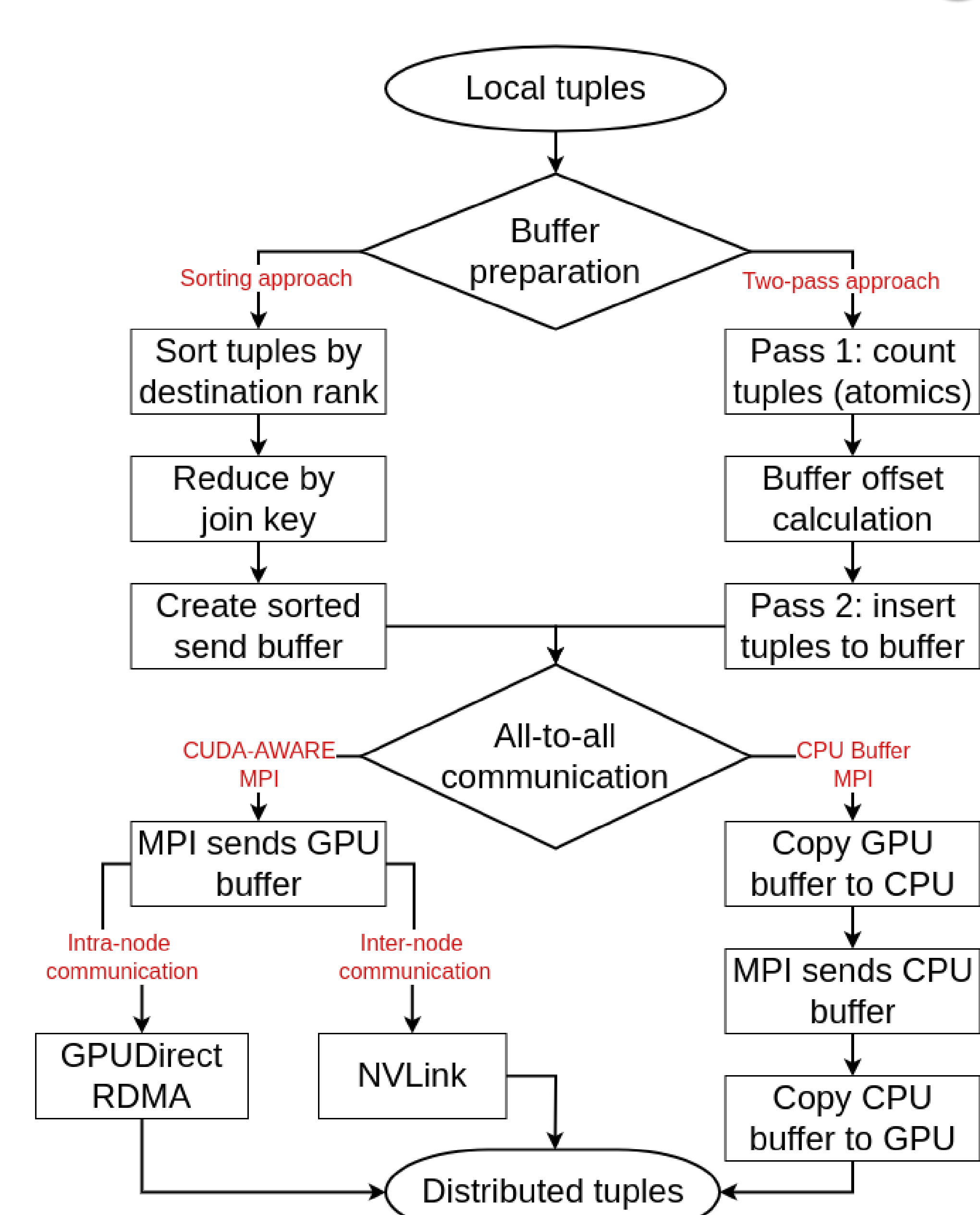
MNMGDatalog uses radix-hash partitioning and non-uniform all-to-all communication with GPU-aware hash tables for efficient tuple materialization



## Radix-hash-based data partitioning



## Configurable buffer preparation and all-to-all communication strategy



## Experiments

We evaluate MNMGDatalog against state-of-the-art single-GPU, shared-memory, and distributed multi-node Datalog engines up to 32 A100 GPUs

**Experiment platform and datasets**

<b>Polaris</b> supercomputer from <b>Argonne National Lab</b>
<b>CPU:</b> AMD EPYC 7543P processors with 32 cores
<b>GPU:</b> 4 NVIDIA A100 GPUs per node interconnected by NVLink
<b>Environment:</b> CUDA (12), SLOG(32 threads), Soufflé (128 threads)
<b>Benchmark:</b> Transitive Closure, Same Generation, Connected Component
<b>Datasets:</b> Stanford large network, SuiteSparse, Road network datasets

Buffer preparation and communication

1 MPI CPU Sort	2 MPI CPU Two pass
3 CUDA-Aware-MPI Sorting	4 CUDA-Aware-MPI Two pass

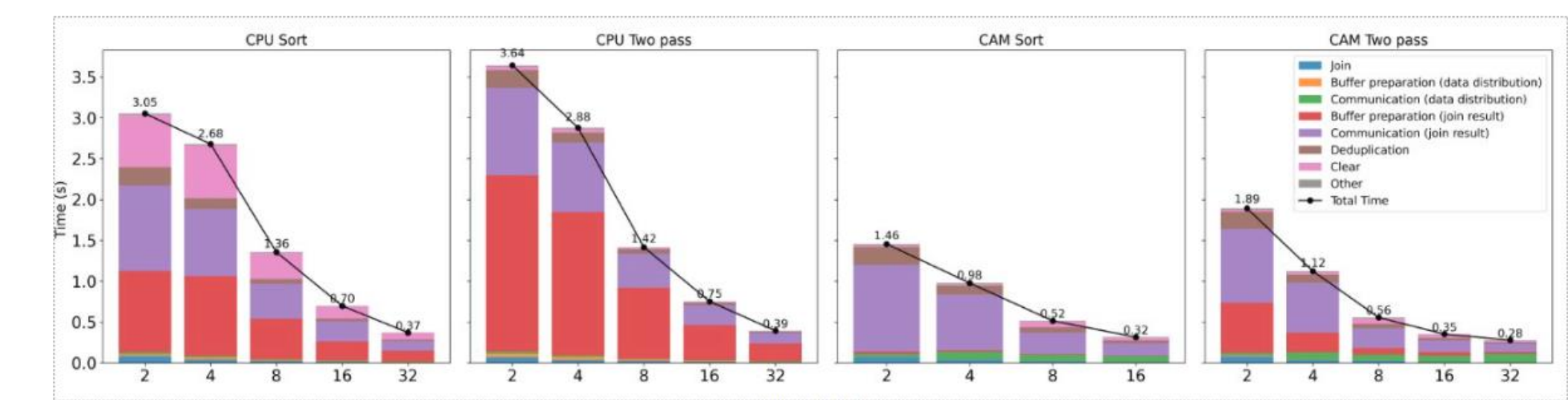
Single-GPU evaluation for Transitive Closure (TC)

Dataset name	TC edges	Time (s)			
		MNMGDATALOG	GPULog	Soufflé	GPUJoin
com-dblp	1.91B	<b>13.58</b>	26.95	232.99	OOM
fe_ocean	1.67B	<b>66.34</b>	72.74	292.15	100.30
usroads	871M	<b>75.07</b>	78.08	222.76	364.55
vsp_finan	910M	<b>81.14</b>	82.75	239.33	125.94
Gnutella31	884M	<b>4.75</b>	7.64	96.82	OOM

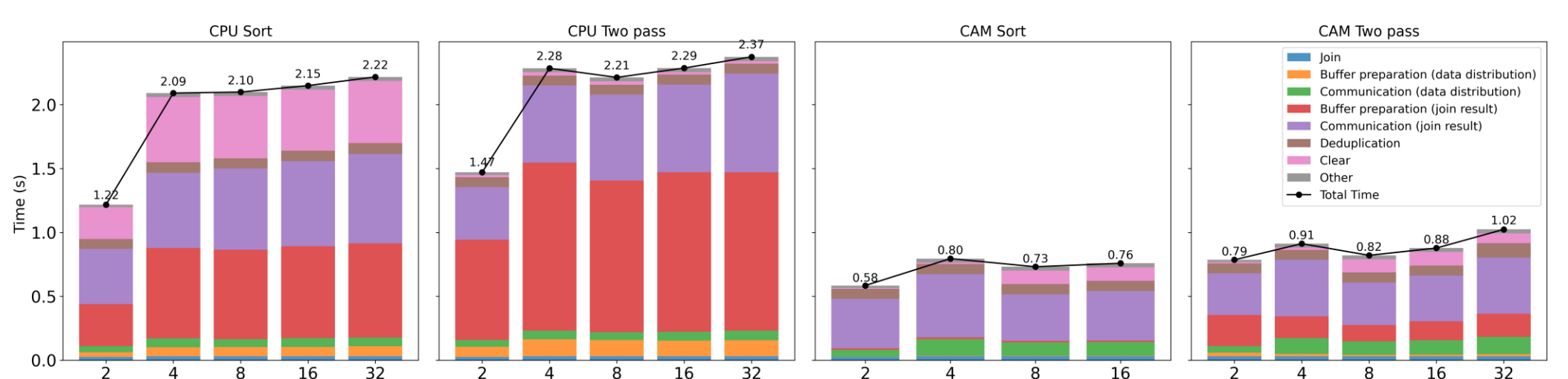
Single-GPU evaluation for Same Generation (SG)

Dataset name	SG size	Time (s)			
		MNMGDATALOG	GPULog	Soufflé	cuDF
fe_body	408M	<b>9.08</b>	18.41	74.26	OOM
loc-Brightkite	92.3M	<b>1.66</b>	11.67	48.18	OOM
fe_sphere	205M	<b>3.55</b>	7.88	48.12	OOM
CA-HepTH	74M	<b>0.60</b>	4.79	20.12	21.24

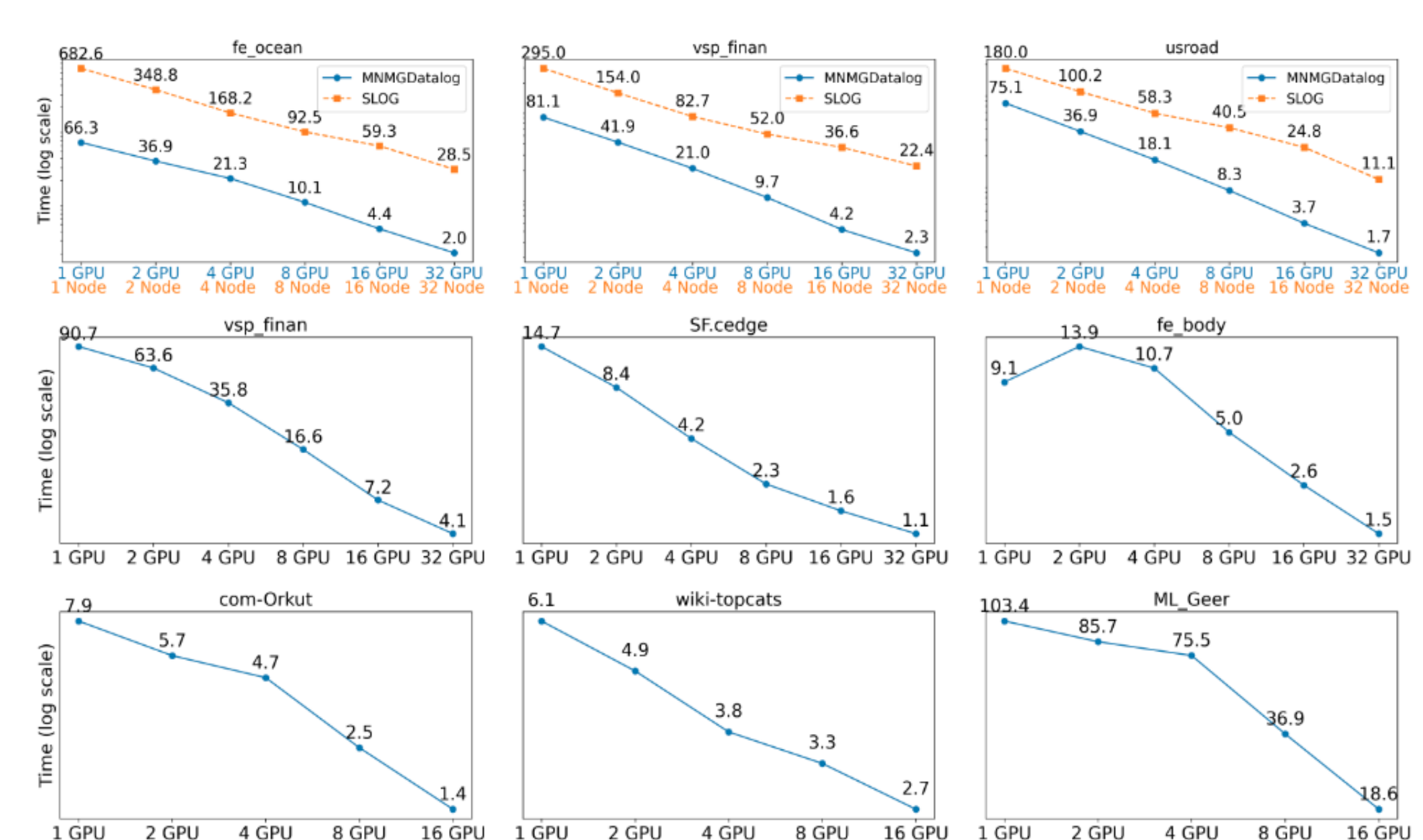
Strong scaling for single iterative join up to 32 GPUs – total 10M tuples)



Weak scaling for single iterative join up to 32 GPUs – 10M tuples per GPU)



Multi-node evaluation for TC, SG, and CC up to 32 GPUs



## Conclusion

Our contributions:

- First ever Datalog engine designed for multi-node multi-GPU HPC systems, outperforming state-of-the-art shared-memory, distributed-memory, and GPU-based engines
- Introduces novel GPU-Aware communication and buffer preparation strategies for scalable recursive query evaluation
- Supports recursive aggregation for Datalog rules using high-throughput GPU kernels

(Accepted for publication at ICS 2025)

## Acknowledgement

This work was funded in part by NSF RII Track-4 award 2132013, NSF collaborative research award 2217036, and NSF collaborative research award 2221811. We are thankful to the ALCF's Director's Discretionary (DD) program for providing us with compute hours to run our experiments on the Polaris supercomputer located at the Argonne National Laboratory.