

# Multi-Node Multi-GPU Datalog

Ahmedur Rahman Shovon

University of Illinois Chicago  
Chicago, IL, USA  
ashov@uic.edu

Yihao Sun

Syracuse University  
Syracuse, NY, USA  
ysun67@syr.edu

Kristopher Micinski

Syracuse University  
Syracuse, NY, USA  
kkmicins@syr.edu

Thomas Gilray

Washington State University  
Pullman, WA, USA  
thomas.gilray@wsu.edu

Sidharth Kumar

University of Illinois Chicago  
Chicago, IL, USA  
sidharth@uic.edu

## Abstract

Datalog, a declarative logic programming language that operates bottom-up, has experienced increasing popularity due to its natural handling of recursive queries. Its applications span diverse fields, including graph mining, program analysis, deductive databases, and neuro-symbolic reasoning. While Datalog shares similarities with SQL in using relational algebra kernels, it uniquely employs iterative execution until reaching a fixed point to support recursion. Current Datalog engines like SLOG, LogicBlox, and Soufflé work well with multi-core and multi-threaded systems, but none have yet tackled multi-node, multi-GPU architectures. Our research addresses this gap by developing the first multi-GPU, multi-node Datalog engine. This advancement is particularly for high-performance computing (HPC) systems, which typically feature multiple GPUs per node. Our implementation combines MPI for inter-node communication with CUDA for GPU parallelization, enabling the processing of massive datasets in real time. We have created novel data-parallel implementations of core relational algebra operations (join), while also optimizing deduplication and tuple materialization. To handle iterative execution, we have developed two novel GPU-accelerated methods for non-uniform all-to-all data exchange. Evaluating on Argonne National Lab's Polaris supercomputer demonstrated our engine's effectiveness, achieving performance improvements of up to 32× against state-of-the-art multi-node Datalog engine.

## Keywords

Datalog, Multi-GPU, Analytic Databases

## ACM Reference Format:

Ahmedur Rahman Shovon, Yihao Sun, Kristopher Micinski, Thomas Gilray, and Sidharth Kumar. 2025. Multi-Node Multi-GPU Datalog. In *2025 International Conference on Supercomputing (ICS '25)*, June 08–11, 2025, Salt Lake City, UT, USA. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3721145.3730431>

## 1 Introduction

Datalog is a declarative logic programming language notable for its elegant handling of recursive queries. Its power is exemplified by its ability to express complex algorithms succinctly. For instance, computing a graph's transitive closure requires merely two lines of Datalog code. Similarly, other graph algorithms like same graph generation, connected component analysis, and single-source shortest path calculations can be implemented in just two to three lines of Datalog code [53]. This remarkable expressiveness has led to Datalog's adoption across diverse domains, from bioinformatics and graph mining [26, 45, 48, 51] to program analysis [7, 12, 19, 21], and neuro-symbolic reasoning [40].

Datalog operates by translating queries into relational algebra operations, such as joins, projections, and unions. The system architecture consists of two main components: a frontend compiler that transforms Datalog queries into iterative relational algebra kernels and a backend that executes these kernels [34] until a fixed-point is reached. This design enables a powerful combination of high-level expressiveness through Datalog's syntax while achieving performance through parallel implementation of the relational algebra kernels. Several engines have implemented this approach: Soufflé [32] leverages OpenMP for multi-threaded CPU processing, SLOG [23] employs MPI for distributed CPU computation, GPUJoin [50] and GPULog [54] utilize CUDA for GPU acceleration.

Among parallel Datalog engines, the GPU-based approach significantly outperforms others, with reported speedups exceeding 10× in complex tasks such as program analysis [54]. This advantage comes from Datalog's inherently



This work is licensed under a Creative Commons Attribution 4.0 International License.

ICS '25, Salt Lake City, UT, USA

© 2025 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-1537-2/25/06

<https://doi.org/10.1145/3721145.3730431>

data-intensive and memory-bound nature, where each iteration requires deduplication, indexing, and aggregation of large volumes of generated data. Datacenter-grade GPUs provide much higher memory bandwidth than CPUs of the same class; for instance, the Nvidia H100 reaches 3.35 TB/s [42], whereas the high-end AMD Zen 5 achieves only 576 GB/s [4].

In this paper, we introduce *MNMGDATALOG*, the first multi-node, multi-GPU Datalog engine. To the best of our knowledge, ours is the first Datalog engine that fully harnesses the potential of modern GPU-based supercomputers. Our engine employs a radix-hash-based data partitioning scheme to balance computation across GPUs while minimizing data exchange. We design novel distributed relational algebra operators that consider both data partitioning strategies and the SIMT nature of GPUs [43]. Our framework relies on non-uniform all-to-all data exchanges to facilitate fixed-point iteration. We explore two different implementations of all-to-all data exchanges that also leverage the computational capabilities of the GPU. In addition to core features found in other GPU-based Datalog engines, we implement recursive aggregation, a widely used Datalog feature, and scale it efficiently across multiple nodes. We evaluate *MNMGDATALOG*'s performance up to 32 NVIDIA A100 GPUs on the Polaris supercomputer [5], benchmarking its performance against state-of-the-art CPU and GPU-based engines on diverse graph analytic queries. Our main contributions are summarized as follows:

- We present a radix-hash-based data partitioning strategy, optimized for indexing and iterative computation.
- We implement CUDA-aware non-uniform all-to-all exchanges for tuple materialization to facilitate iterative relational algebra.
- We implement and scale recursive aggregation on GPU.
- On a single GPU, *MNMGDATALOG* achieves up to 7× speedup over GPULog, and up to 33× over Soufflé. In a multi-node, multi-GPU setting, *MNMGDATALOG* outperforms the state-of-the-art HPC-based engine SLOG by up to 32×.

## 2 Declarative analytics using Datalog

Datalog operates through a lightweight bottom-up evaluation approach and is widely used in deductive database systems [13, 14, 16]. A Datalog program consists of an *extensional database* containing explicit input facts and an *intensional database* of derived facts inferred through rules [11]. These rules are expressed as first-order Horn clauses, where each rule takes the form:

$$\text{Head} \leftarrow \text{Body}_1, \text{Body}_2, \dots, \text{Body}_n$$

The head contains a single predicate atom which represents the inferred fact, while the body specifies the conditions for its derivation using a set of predicate atoms. The implication

symbol  $\leftarrow$  connects the head with the body. Commas in the body represent logical AND ( $\wedge$ ) that performs a *join* operation between the predicate atoms.

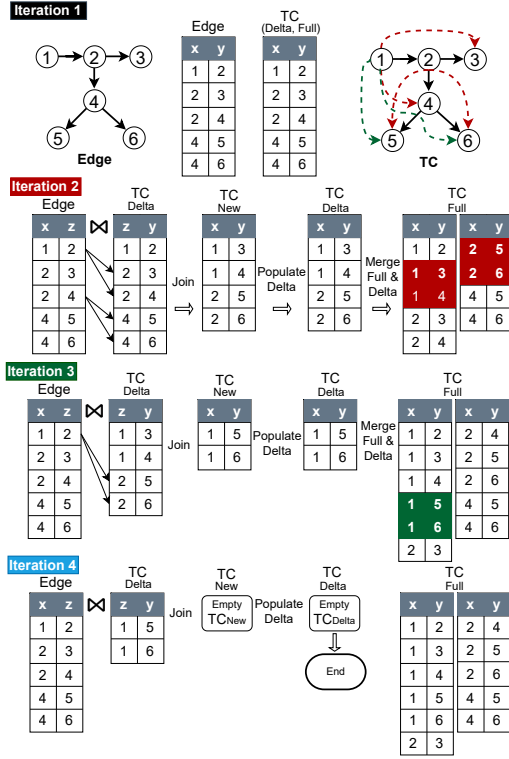
Datalog excels at handling recursive queries through *fixed-point evaluation*, where rules are repeatedly applied until no new facts can be derived. For instance, the following rules can be used to derive the transitive closure of an input graph, which finds all reachable paths from each node:

$$\begin{aligned} \text{TC}(x, y) &\leftarrow \text{Edge}(x, y). \\ \text{TC}(x, z) &\leftarrow \text{TC}(x, y), \text{Edge}(y, z). \end{aligned}$$

The first rule establishes that a direct edge from  $x$  to  $y$  implies reachability. The second rule, which is recursive, extends this reachability by chaining existing paths: if  $x$  can reach  $y$  and there is an edge from  $y$  to  $z$ , then  $x$  can also reach  $z$ . This recursive expansion is achieved through a *join* operation between the TC and Edge relations, where previously derived reachability facts from TC are joined with new edges from Edge to infer additional paths. This process continues iteratively until no new paths can be inferred. Such recursive capabilities allow Datalog to efficiently solve problems like transitive closure, connected components, same generation, and other queries that traditional SQL struggles to handle. Though traditional SQL supports recursive queries through common table expressions (CTEs), Datalog's expressive syntax and bottom-up evaluation make recursive queries simpler and more efficient by automatically iterating until a fixed point is reached, without requiring explicit control structures.

*Semi-naïve evaluation.* Modern datalog engines achieve performance improvements through the semi-naïve evaluation [3], an incremental evaluation technique. Our framework, *MNMGDATALOG*, likewise implements this optimization technique, which improves the efficiency of iterative queries by exclusively utilizing only the newly derived facts during each successive iteration. This strategy avoids redundant computation by ensuring that only newly derived facts are used to infer additional facts in the iterative process.

In semi-naïve evaluation, non-static relations (such as TC in path-finding applications) are strategically decomposed into three distinct components: (1) *full*, which encompasses all facts discovered prior to the most recent iteration; (2) *delta*, containing exclusively those facts identified during the immediately preceding iteration; and (3) *new*, which stores facts newly derived during the current computational iteration. Throughout each iteration, join operations are selectively applied using only the *delta* component of a relation, with resultant tuples stored in the *new* component. Upon iteration completion, the algorithm executes a three-phase transition process: it transfers all tuples from *delta* into the *full* relation, exchanges pointers between *delta*



**Figure 1: Iterations of transitive closure computation.**

and new to prepare delta for the subsequent iteration, and reinitializes new to accommodate upcoming derivations. Figure 1 provides a visual representation of this iterative execution process for transitive closure rules under the semi-naïve evaluation framework.

The top left presents the input graph, while the top right depicts the fully computed TC. In the first iteration, the TC relation is initialized using the direct edges from the Edge relation. Since TC is initially empty, the recursive rule does not contribute any new facts at this stage. Both the Delta and Full versions are set to be identical to TC, ensuring that all direct connections are established before applying recursive expansions in subsequent iterations. In the second iteration, the recursive rule computes  $\text{New} := \text{Edge} \bowtie \text{Delta}$ , yielding  $\{(1, 3), (1, 4), (2, 5), (2, 6)\}$ . Since these tuples are not present in the existing Full version, they are added to Full (highlighted in red) and also stored in Delta for the next iteration. In the third iteration, the join produces unique set of tuples  $\{(1, 5), (1, 6)\}$ , which are merged into Full (highlighted in green) and retained in Delta. In the fourth iteration,  $\text{Edge} \bowtie \text{Delta}$  produces no new tuples, leaving Delta empty and signaling fixpoint termination. This stepwise expansion optimizes TC computation by restricting joins to newly derived facts, eliminating redundant computation while ensuring correctness.

*Recursive Aggregation.* Recursive aggregation extends standard Datalog semantics by allowing aggregate functions such as MIN, MAX, SUM, and COUNT to be applied dynamically during recursive evaluation. This feature is useful in graph algorithms, including connected components, shortest paths, and PageRank, where values must be iteratively propagated rather than computed post-fixpoint using stratification [3].

One such application is the Weakly Connected Components (WCC) problem, where recursive aggregation enables the efficient propagation of component representatives. The WCC query can be formulated in Datalog as follows:

$$\begin{aligned} \text{WCC}(n, n) &\leftarrow \text{Edge}(n, \_). \\ \text{WCC}(y, \text{MIN}(z)) &\leftarrow \text{WCC}(y, z), \text{Edge}(x, y). \end{aligned}$$

The first rule initializes each node as its own component, while the second rule propagates the smallest representative node ID across connected nodes using the MIN aggregate. Unlike traditional Datalog engines that materialize all possible component memberships, recursive aggregation ensures that only the minimal component representative is maintained, reducing both space complexity and redundant computations.

### 3 Challenges and requirements

This section outlines the critical requirements for building a multi-node multi-GPU datalog engine. We focus on three fundamental components necessary to achieve scalable performance: workload partitioning, data representation, and inter-node data exchange.

#### 3.1 Workload partitioning

Parallelizing algorithms necessitates identifying an appropriate partitioning axis for workload distribution. Two fundamental approaches exist for problem partitioning: model/task-level partitioning and data-level partitioning. Datalog programs inherently support both paradigms. While task parallelism is program-dependent, allowing complex Datalog programs to be decomposed into independent, concurrently executable task groups – our research exclusively addresses data-level parallelism. We propose a data-parallel framework that effectively distributes the computational workload of Datalog programs across multiple GPUs. This approach requires strategic partitioning of all relations within the Datalog program and subsequent allocation of these partitioned segments across available GPUs.

Conventional GPU-accelerated algorithms are engineered primarily for dense computational patterns, exemplified by matrix multiplication operations [1, 2, 20, 28]. These dense workloads facilitate uniform partitioning into equal-sized computational units, thereby optimizing memory bandwidth utilization and cache efficiency. In contrast, Datalog computation exhibits inherent irregularity and sparsity, as the

relations involved in Datalog programs vary considerably in size, characteristics, and topological properties. This intrinsic heterogeneity renders efficient data partitioning across multiple GPUs substantially more complex.

Naive approaches that uniformly distribute relations across available GPUs prove inadequate, as effective partitioning must specifically accommodate the fundamental operations underlying Datalog execution, such as low-level relational algebra kernels including joins, unions, projections, and other tasks like deduplication and merging procedures. In Figure 1, when distributing Edge relation with tuples (1,2), (2,3), (2,4), (4,5), and (4,6) across two GPUs, a join-preserving approach must ensure that all tuples with identical join keys such as those beginning with "2" are allocated to the same GPU. This locality preserving allocation is essential for ensuring that join operations, which in this case occur on the first column, can be performed efficiently within a single GPU without cross-device communication. A naive uniform distribution strategy that allocates an equal number of tuples (e.g., three tuples per GPU) would compromise this crucial locality property. While such an approach might achieve nominal data balance across GPUs, it inevitably leads to highly imbalanced computational workloads during execution and introduces significant inter-GPU communication overhead. To address this distribution challenge, we implement a hash-based partitioning methodology [9] wherein relational data is systematically distributed across all available GPUs according to the hash value of the join column (detailed in Section 4.1). Thus, the first key challenge we propose to address in this paper is designing a data partitioning strategy that respects Datalog computation while minimizing communication overhead.

### 3.2 Data representation

A prerequisite for constructing a scalable multi-node, multi-GPU Datalog execution engine is the optimization of single-GPU performance. Fundamental to achieving peak single-GPU efficiency is the underlying data representation, specifically, how relational data is organized and stored within GPU memory. The critical design challenge lies in developing data structures that simultaneously achieve multiple performance objectives: minimizing memory footprint, enabling high-throughput data retrieval operations (essential for join execution), and effectively supporting auxiliary operations such as deduplication. These memory-resident data structures form the foundation upon which the entire distributed computation framework depends, directly influencing overall system scalability and performance characteristics.

For efficient lookup operations, hash tables represent the predominant data structure upon which hash join algorithms are constructed [8, 24, 33]. Conventional CPU-oriented hash join implementations frequently utilize hash tables built on

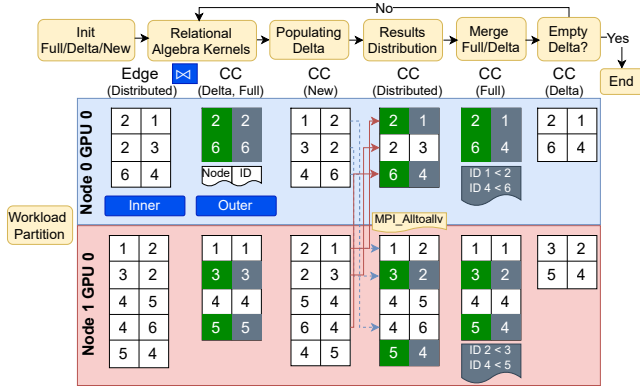
linked-list architectures. However, these structures demonstrate poor performance characteristics when deployed on GPU architectures, as pointer-chasing operations inherently generate non-coalesced memory access patterns, resulting in significant latency penalties. To overcome this architectural limitation, our implementation employs a specialized open-addressing hash map with linear probing techniques. This design modification substantially optimizes memory access patterns and enhances overall execution performance on GPU hardware. A comprehensive detail of this implementation approach, including performance characteristics and design considerations, is presented in Section 4.

### 3.3 Efficient communication

Relational algebra (RA) kernels, including join and other operations, executed locally within individual GPUs generate new tuples that typically serve as input for subsequent iterations of the semi-naive evaluation process. However, these newly generated tuples do not necessarily belong to the GPU/process where they were originally produced. To materialize the newly generated facts and consequently to facilitate iterative parallel relational algebra execution, processes must participate in a non-uniform all-to-all inter-process exchange of generated tuples to their appropriate destination (GPU). The materialization of a tuple in an output relation involves hashing its join column to identify the target GPU, followed by transmission to that specified GPU. Since tuples produced during the local computation phase may each correspond to arbitrary GPUs in the output relation, an all-to-all communication phase becomes necessary to redistribute these output tuples to their managing processes. Due to inherent variations in both the number of tuples generated across different processes and their destination distributions, the all-to-all communication phase exhibits a fundamentally *non-uniform* characteristic. We explain the need for a *non-uniform all-to-all* run using a real example.

In Figure 1, consider performing transitive closure of the Edge relation using two GPUs with hash-based data distribution on the first column. Under this scheme, *GPU 0* stores {(2, 3), (2, 4)} and *GPU 1* holds {(4, 5), (4, 6), (1, 2)} based on hashing on the first column. During TC computation, a join operation is performed, followed by a projection that eliminates the common column, resulting in newly derived tuples after the first iteration: *GPU 0* produces {(1, 3), (1, 4)}, while *GPU 1* derives {(2, 5), (2, 6)}. Since the projection step eliminates the join column, a global redistribution of all new tuples is necessary. This necessitates all-to-all communication, where each GPU sends and receives derived tuples according to the hash partitioning scheme, ensuring that subsequent join operations are performed locally. Without an optimized communication strategy, this redistribution





**Figure 2: First iteration of semi-naïve evaluation with local aggregation on Weakly Connected Component (WCC) query using MNMGDATALOG.**

introduces significant overhead, limiting scalability. An efficient communication mechanism must address three key challenges: structured data movement, minimal memory overhead, and scalable GPU-to-GPU communication.

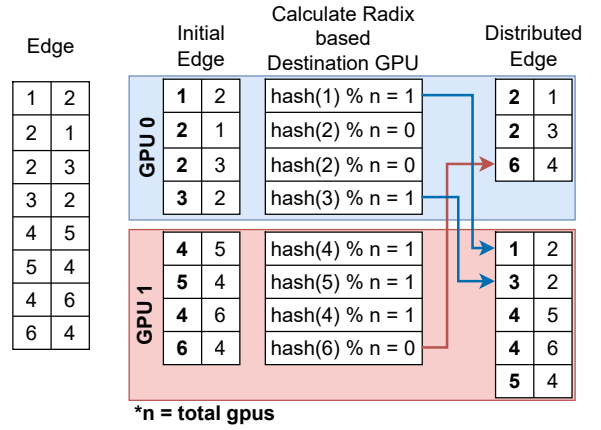
#### 4 MNMGDATALOG: multi-node multi-GPU Datalog

This section describes the implementation of MNMGDATALOG, the first multi-node, multi-GPU Datalog engine. Figure 2 provides a structural overview of our engine, illustrating the execution flow of the first iteration of Weakly Connected Components (WCC) query on a multi node multi GPU setup with semi-naïve evaluation strategy. It links directly to the key components of our implementation: workload partitioning with hash-based data distribution (Sec 4.1), efficient communication for distributed result propagation (Sec 4.2), and GPU-optimized relational algebra kernels (Sec 4.3).

##### 4.1 Hash-based data distribution

Inspired by the classic parallel processing algorithm GRACE Hash Join [22] in distributed RDBMS, and the distributed Datalog engine BPRA [35], MNMGDATALOG creates local data storage in each GPU's VRAM and partitions each relation based on the join columns. We illustrate this process using the partitioning of the *edge* relation, described in Section 2.

Figure 3 shows how the *Edge* relation is distributed across multiple GPUs using a radix-based hashing technique. Each GPU initially holds a disjoint subset of the input edges. To redistribute the data, the system applies a hash function to the first column of each edge tuple (the join key), and the result modulo the total number of GPUs (which is 2 in this case) determines the destination GPU. For example, with two GPUs, a tuple with key  $k$  is sent to GPU  $hash(k) \bmod 2$ . This ensures that all tuples sharing the same join key are



**Figure 3: Data distribution based on radix based hashing technique within the available GPUs**

co-located on the same GPU, enabling local join operations without requiring cross-GPU communication during the join phase.

In real-world datasets, data skew can be an issue. For example, in social network analysis, some users may have substantially more followers than others, resulting in uneven computational load across GPUs. In such cases, a pure hash-based partitioning strategy based solely on the join column may be insufficient. Recent research [35, 53] proposes a promising solution using sub-bucketing to address this challenge. While this technique has not yet been implemented in our system, it aligns well with our architecture and is on our roadmap for future integration.

##### 4.2 Data communication

As established in Section 3.3, our framework needs non-uniform all-to-all data exchanges to materialize newly generated tuples during each iteration of the fixed-point loop. Within the MPI programming model, non-uniform all-to-all communication is implemented using the `MPI_Alltoallv` [17, 18, 46] collective operation. This function requires all processes to participate synchronously, with the first argument specifying a contiguous buffer containing the concatenated data segments destined for all participating processes. To correctly interpret this buffer during transmission, the function requires supplementary offset and count arrays that precisely delineate the boundaries of individual data segments targeted to specific processes. Consequently, prior to invoking these MPI collective operations, the system must construct the all-to-all send buffer and calculate accurate offsets for the data segments destined for each GPU process within the consolidated buffer. The preparation of this buffer in the required format can incur significant computational overhead. To address this challenge, we leverage

**Algorithm 1** Sorting-Based buffer preparation and All-to-All communication

```

1: Input: Local GPU buffer  $D$ , total GPUs  $R$  (1 MPI rank per GPU)
2: Output: Distributed GPU buffer  $receive\_data$ 
3: for each tuple  $(key, value)$  in  $D$  parallel do
4:    $row\_mapping \leftarrow get\_rank(key, R)$ 
5: end for
6:  $StableSortByKey(row\_mapping, D)$ 
7:  $(unique\_rank, send\_count) \leftarrow ReduceByKey(row\_mapping, D)$ 
8:  $ExclusiveScan(send\_count, send\_displacement)$ 
9:  $receive\_count \leftarrow MPI\_Alltoall$  on  $send\_count$ 
10:  $ExclusiveScan(receive\_count, receive\_displacement)$ 
11: if CUDA-aware MPI is supported then
12:    $receive\_data \leftarrow MPI\_Alltoallv$  on
      $(D, send\_displacement, receive\_displacement)$ 
13: else
14:   Copy  $D$  to CPU buffer  $send\_data\_host$ 
15:    $receive\_data\_host \leftarrow MPI\_Alltoallv$  on
      $(send\_data\_host, send\_displacement, receive\_displacement)$ 
16:   Copy  $receive\_data\_host$  to GPU buffer  $receive\_data$ 
17: end if
18: Return:  $receive\_data$ 

```

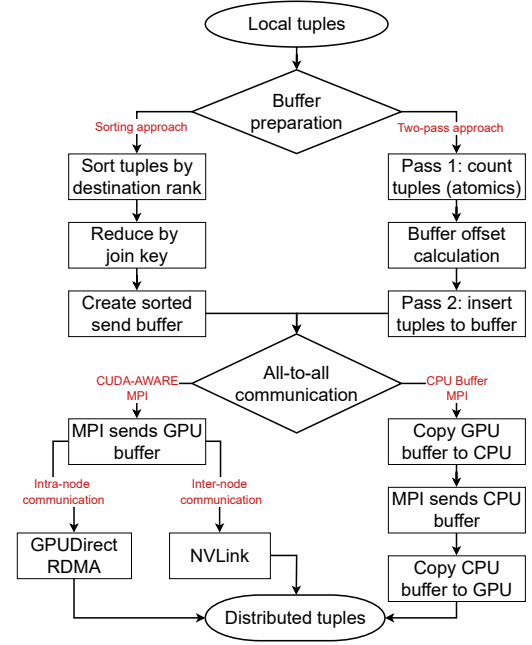
**Algorithm 2** Two-Pass buffer preparation and for All-to-All communication

```

1: Input: Local GPU buffer  $D$ , total GPUs  $R$  (1 MPI rank per GPU)
2: Output: Distributed GPU buffer  $receive\_data$ 
3: for each tuple  $(key, value)$  in  $D$  parallel do ▷ First pass
4:    $destination\_rank \leftarrow get\_rank(key, R)$ 
5:    $AtomicAdd(send\_count[destination\_rank], 1)$ 
6: end for
7:  $ExclusiveScan(send\_count, send\_offset)$ 
8: Copy  $send\_offset$  to  $send\_displacement$ 
9: for each tuple  $(key, value)$  in  $D$  parallel do ▷ Second pass
10:   $destination\_rank \leftarrow get\_rank(key, R)$ 
11:   $position \leftarrow AtomicAdd(send\_offset[destination\_rank], 1)$ 
12:   $send\_data[position] \leftarrow (key, value)$ 
13: end for
14:  $receive\_count \leftarrow MPI\_Alltoall$  on  $send\_count$ 
15:  $ExclusiveScan(receive\_count, receive\_displacement)$ 
16: if CUDA-aware MPI is supported then
17:    $receive\_data \leftarrow MPI\_Alltoallv$  on
      $(send\_data, send\_displacement, receive\_displacement)$ 
18: else
19:   Copy  $send\_data$  to CPU buffer  $send\_data\_host$ 
20:    $receive\_data\_host \leftarrow MPI\_Alltoallv$  on
      $(send\_data\_host, send\_displacement, receive\_displacement)$ 
21:   Copy  $receive\_data\_host$  to GPU buffer  $receive\_data$ 
22: end if
23: Return:  $receive\_data$ 

```

the parallel processing capabilities of GPUs to efficiently



**Figure 4: Communication phases of MNMGDATALOG**

prepare these buffers. MNMGDATALOG implements and comparatively evaluates two distinct buffer preparation methodologies, sorting-based preparation and two-pass preparation, outlined in Figure 4.

The sorting-based approach is outlined in Algorithm 1. This method begins by copying all data into the send buffer and computing the destination GPU rank for each tuple using the data partitioning mechanism described in the previous section. The tuples in the send buffer are then sorted by their associated rank numbers, which automatically groups all tuples destined for the same rank together. Next, a histogram of the rank numbers in the send buffer is computed to determine the data offsets for each rank. The advantage of this approach is that all the operations involved, such as sorting and histogram computation, can be efficiently accelerated on GPUs using well-established algorithms [47, 52]. These algorithms are specifically designed to maximize GPU memory bandwidth utilization. This approach minimizes memory fragmentation, but incurs additional computational overhead due to sorting.

An alternative method, the two-pass buffer preparation technique described in Algorithm 2, eliminates the need for sorting by performing a two-step counting and writing process. CUDA kernels are used to directly count the number of tuples destined for each rank. The first kernel pass scans the input tuples and uses atomic operations to track the send count for each rank. The second kernel pass prepares the send buffer based on these counts by writing the tuples to

their respective positions in memory. This approach eliminates the need for sorting, but the resulting buffers may be more fragmented, potentially affecting memory access efficiency during the communication phase. While the two-pass method avoids sorting overhead, we found that the sorting-based approach yields slightly better performance due to improved memory coalescing and reduced fragmentation (detailed in Section 5.4).

*All-to-all communication.* Once buffer preparation is complete, the second stage performs all-to-all communication to shuffle data among GPUs. MNMGDATALOG supports two communication mechanisms: CUDA-aware MPI [6, 31, 56], and CPU buffer-based MPI communication. CUDA-aware MPI enables direct communication between GPU buffers without requiring intermediate copies to CPU memory. This technique leverages technologies like GPUDirect RDMA and NVLink to achieve high bandwidth and low-latency data transfers between GPUs. GPUDirect RDMA allows GPUs to communicate directly with the network interface card, bypassing the host CPU and reducing communication latency. Similarly, NVLink provides high-speed interconnects between GPUs on the same node, enabling faster data movement during intra-node communication. In systems that support CUDA-aware MPI, MNMGDATALOG’s communication calls directly transfer GPU-resident buffers between processes, minimizing overhead. To ensure compatibility across a broader range of systems, we also provide a CPU buffer-based mode, where GPU data is first copied to CPU memory before invoking MPI. While this approach incurs additional data movement overhead, it guarantees that MNMGDATALOG can run on systems without CUDA-aware MPI support. MNMGDATALOG provides a runtime configuration to select between these two communication modes. This flexibility ensures that MNMGDATALOG can maximize performance on CUDA-aware MPI-enabled systems by utilizing direct GPU exchanges, while also maintaining portability across architectures.

### 4.3 GPU-optimized data representation

MNMGDATALOG employs an open-addressing hash table with linear probing for efficient join execution. Hash joins dominate computational cost in recursive Datalog queries, making their optimization crucial for scalability. In iterative queries such as Weakly Connected Components (WCC), the inner relation remains unchanged across iterations, making static hash tables an ideal choice. By constructing the hash table once and reusing it in subsequent iterations, MNMGDATALOG eliminates unnecessary recomputation. The hash join consists of two phases: the build phase and the probe phase. During the build phase, the inner relation keys are inserted into the hash table, with collisions handled via linear probing. Contiguous memory storage ensures high cache locality,

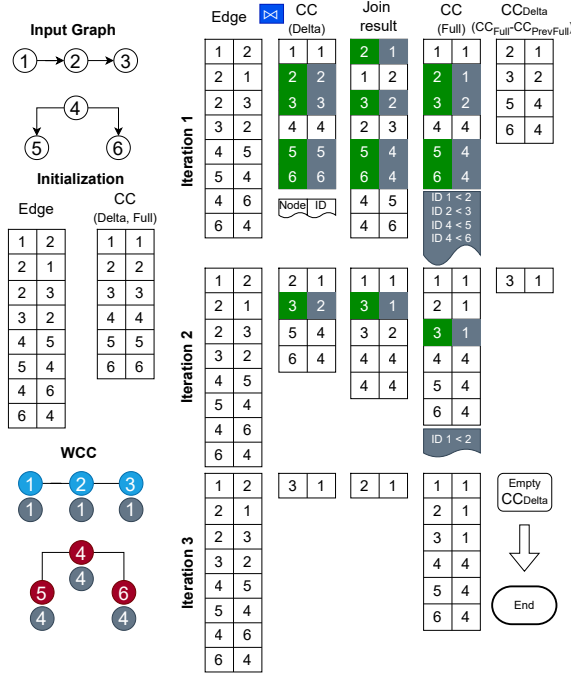
reducing memory latency. In the probe phase, the Delta relation serves as the outer relation, querying the hash table for matches to propagate new facts. Although our engine is optimized for iterative join during Datalog evaluation, many of these join strategies, such as static hash table construction and GPU-friendly memory layouts, can also benefit traditional non-iterative join workloads. Traditional joins share the same underlying hash-based probing mechanism. Therefore, MNMGDATALOG’s GPU-accelerated join implementation is general enough to improve performance in both iterative and non-iterative relational joins.

Figure 2 shows an example of the first iteration of the fixed-point loop for the WCC query, including the deduplication process in MNMGDATALOG that eliminates redundant computation across iterations. After each GPU performs a local join between the distributed Edge relation and the current delta of the connected component relation ( $CC_{\Delta}$ ), new tuples representing potential component ID updates are produced. These tuples may include redundant entries due to overlapping neighborhoods or multiple paths propagating the same component ID. To handle this issue, each GPU conducts its own local deduplication process. This is achieved by first sorting all locally generated results (shown as  $CC_{New}$ ), then eliminating redundant entries by using thrust’s unique function [44]. This function works by examining consecutive keys with identical values and preserving only the first occurrence while removing all subsequent duplicates. The deduplicated results are then exchanged across GPUs, shown as the transition from  $CC_{New}$  to  $CC_{Distributed}$ . Post-communication, another deduplication pass is applied to remove inter-GPU duplicates. Then a set difference operation is applied between  $CC_{Distributed}$  and  $CC_{Full}$  to identify only the unseen tuples which will be used in subsequent iteration as  $CC_{\Delta}$ . This ensures that the fixpoint loop only processes unique tuples, avoiding redundant computation. To enable efficient set-difference and merging,  $CC_{Full}$  is kept sorted throughout all iterations.

To further optimize performance, MNMGDATALOG leverages grid-stride loops, ensuring that each GPU thread processes multiple tuples in a row-major order. This approach minimizes memory divergence by aligning memory accesses with the GPU’s warp execution model, reducing unnecessary stalls and maximizing cache hits. This design choice ensures that the local joins are memory-efficient and scalable.

### 4.4 Recursive aggregation in MNMGDATALOG

Beyond basic Datalog semantics, MNMGDATALOG supports GPU-based parallel recursive aggregation through a modified deduplication phase. In standard relations, MNMGDATALOG performs deduplication by launching a GPU thread for



**Figure 5: All iterations of Weakly Connected Components (WCC) calculation using the semi-naïve evaluation technique incorporating recursive aggregation.**

each newly generated tuple. Each thread probes the hash table of the full relation using the indexed column value and performs an equality check on the non-indexed column. If no match is found, the tuple is considered unique and is inserted into the delta version of the relation. For relations involving recursive aggregation, instead of a strict equality check on the non-indexed column, MNMGDATALOG leverages the monotonicity property of aggregate functions. When recursively aggregated values are associated with the same non-aggregated key, they can only evolve monotonically, for instance, decreasing in the case of MIN or increasing for MAX. MNMGDATALOG replaces the equality check with an aggregate-aware comparator that retains and propagates only the improved values (i.e., decreased or increased value). In this case, not only is the improved value inserted into the delta relation, just like in standard deduplication, but the corresponding entry in the full relation is also updated with the new aggregated value.

In Figure 5, we use the recursive MIN operator in the WCC query as an example to illustrate how recursive aggregation is implemented in MNMGDATALOG. In the first iteration, the  $CC_{\Delta}$  relation is joined with the Edge relation to propagate component IDs across connected nodes. Initially, each node is assigned its own ID. During aggregation, node 2 receives a smaller ID from node 1, and similarly, nodes 5 and 6 update their IDs based on their neighbors, as shown in the

gray boxes in the figure. To keep the aggregated value monotonically decreasing so that we can get minimal value after fixpoint, MNMGDATALOG will store the tuple with smaller aggregated value (i.e., tuple (2, 1), (5, 4)) in  $CC_{\Delta}$  and update value in  $CC_{Full}$ . This process repeats in the second iteration, further propagating the smallest component IDs. Recursive aggregation continues in this way, applying the MIN function iteratively. By the third iteration, all nodes have converged to their final component IDs, and  $CC_{\Delta}$  becomes empty, signaling that a fixpoint has been reached. The final connected components are shown in the bottom-left of Figure 5, where nodes sharing the same ID belong to the same connected component; for instance, 1, 2, 3 and 4, 5, 6.

## 5 Evaluation

In this section, we present a comprehensive performance evaluation of MNMGDATALOG, demonstrated by three sets of experiments: (1) evaluation with real applications on one single GPU, (2) evaluation of a single iteration of the fixed-point which includes a single join operation in multi-GPU environment, following by the materialization of the newly generated tuple that includes an all-to-all data exchange phase and (3) evaluation with real applications in multi-GPU environment.

### 5.1 Environment

We conducted all our experiments on the Polaris supercomputer at the Argonne Leadership Computing Facility. Each compute node is equipped with an AMD EPYC Milan 7543P 32-core CPU, 512 GB of DDR4 RAM, and four NVIDIA A100 GPUs interconnected via NVLink. Multi-node communication is facilitated by Slingshot 11 high-speed interconnects. For GPU-based Datalog systems, including MNMGDATALOG, GPUlog, GPUJoin, and cuDF, we used a single GPU on a single node on Polaris. For multi-node experiments, we compared MNMGDATALOG against SLOG, a distributed CPU-based Datalog engine, using the same number of nodes for both systems. Each SLOG node was configured to utilize 32 CPU threads to match Polaris' CPU architecture. CUDA-aware MPI was enabled by using MPI-GPU support to allow the MPI library to send and receive data directly from GPU buffers. To evaluate MNMGDATALOG's portability, we also benchmarked its performance using CPU buffer-based MPI communication, where GPU data was first transferred to host memory, sent via MPI, and copied back to the GPU post-communication.

### 5.2 Test programs and datasets

To evaluate the performance and scalability of MNMGDATALOG, we designed experiments that assess both the performance of a single iteration of the fixed-point loop and full iterative query execution. The single iteration benchmark



isolates the core computational step of recursive Datalog execution, while the transitive closure (TC) (Section 2), same generation (SG) [54], and weakly connected components (WCC)(Section 2) benchmarks evaluate full multi-iteration workloads. For the single join benchmark, we used a synthetic dataset with 10M tuples per rank for weak scaling and a total of 10M tuples for the strong scaling experiment. For recursive queries (TC, SG, WCC), we used large real-world graphs from the Stanford Network Analysis Project (SNAP), SuiteSparse, and road network datasets [15, 37, 38]. These graphs span diverse domains, with output sizes ranging from millions to several billion edges, providing a comprehensive scalability assessment across varying data distributions and computational complexities. As the data are large-scale, parallel I/O is employed, where each MPI rank independently reads and writes its assigned partition from disk, ensuring efficient data distribution, reducing I/O contention, and enabling scalable processing across multiple nodes.

**Table 1: Transitive Closure (TC) execution time comparison: MNMGDATALOG vs. GPULOG, Soufflé (AMD Milan CPU 32 cores), and GPUJoin on large graphs (OOM: out of memory).**

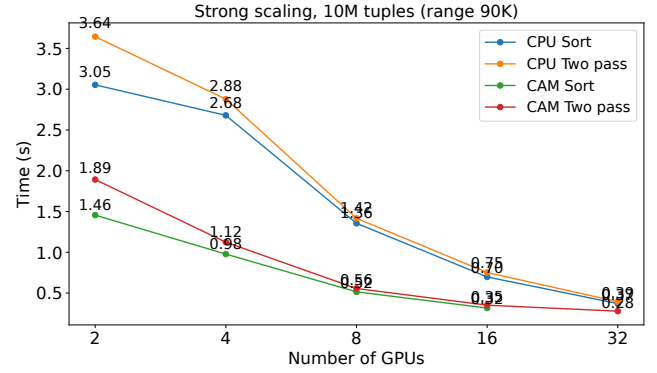
Dataset name	TC edges	Time (s)			
		MNMGDATALOG	GPULOG	Soufflé	GPUJoin
com-dblp	1.91B	<b>13.58</b>	26.95	232.99	OOM
fe_ocean	1.67B	<b>66.34</b>	72.74	292.15	100.30
usroads	871M	<b>75.07</b>	78.08	222.76	364.55
vsp_finan	910M	<b>81.14</b>	82.75	239.33	125.94

**Table 2: Same Generation (SG) execution time comparison: MNMGDATALOG vs. GPULOG, Soufflé and cuDF. Soufflé running on 32 core AMD Milan CPU.**

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

### 5.3 Single GPU benchmark

This section assesses the efficiency of MNMGDATALOG on a single GPU in executing transitive closure (TC) and same generation (SG) queries compared to state-of-the-art GPU-based Datalog engines (GPULOG, GPUJoin, cuDF) and a multi-threaded CPU-based solver (Soufflé). Table 1 and Table 2 present execution times for a single NVIDIA A100 GPU



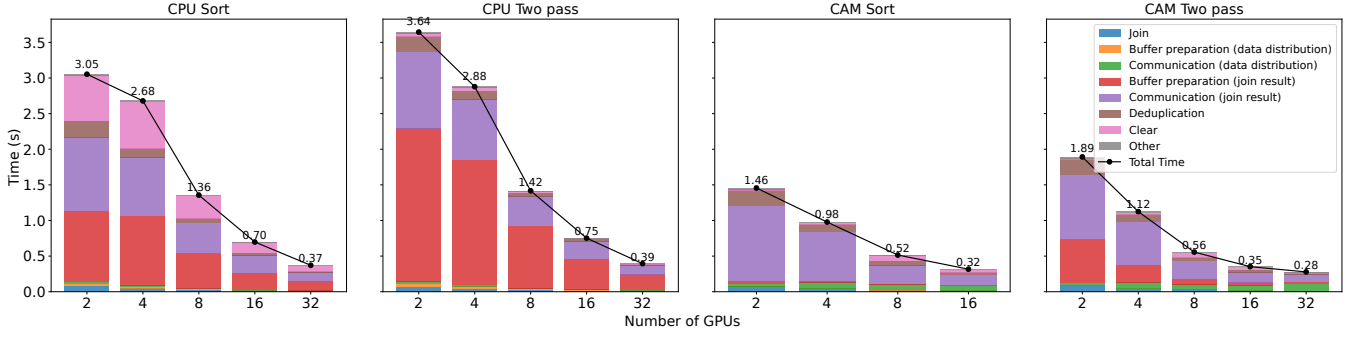
**Figure 6: Strong scaling performance of the single join operation in MNMGDATALOG scaling from 2 to 32 GPUs while keeping the total workload fixed at 10M tuples.**

across MNMGDATALOG, GPULOG, and GPUJoin, while Soufflé is executed with 32 CPU threads to leverage its multi-threaded capabilities.

For TC queries (Table 1), MNMGDATALOG demonstrates competitive performance against GPULOG, outperforming it with up to 1.98× speedup. Compared to Soufflé, MNMGDATALOG achieves up to 20× speedup. Additionally, MNMGDATALOG is up to 4.8× faster than GPUJoin, which fails to process large datasets due to out-of-memory (OOM) errors, highlighting its limited scalability in recursive query execution. For SG queries (Table 2), MNMGDATALOG outperforms GPULOG by up to 7× and achieves a speedup of up to 33.5× over Soufflé. Compared to cuDF, MNMGDATALOG demonstrates up to 35.4× higher performance, as cuDF fails to process most SG queries due to memory limitations. The comparison highlights the superior scalability of MNMGDATALOG over CPU-based approaches while demonstrating its robust handling of large graphs compared to existing GPU-based Datalog engines.

### 5.4 Single iteration of the fixed-point benchmark

We executed a single iteration of the fixed-point loop using synthetic datasets. This includes a join operation followed by an all-to-all data exchange to materialize the newly generated tuples. We executed this benchmark using both sorting-based and two-pass buffer generation approaches to measure their respective impact on buffer preparation time, communication overhead, and overall execution time. Each approach was evaluated under both CUDA-aware MPI (CAM) and CPU buffer-based MPI communication to analyze the effect of direct GPU-to-GPU transfers versus CPU-mediated communication. This test emphasizes the efficiency of iterative join operations, as recursive queries rely heavily on repeated joins across iterations. We conduct both strong scaling and

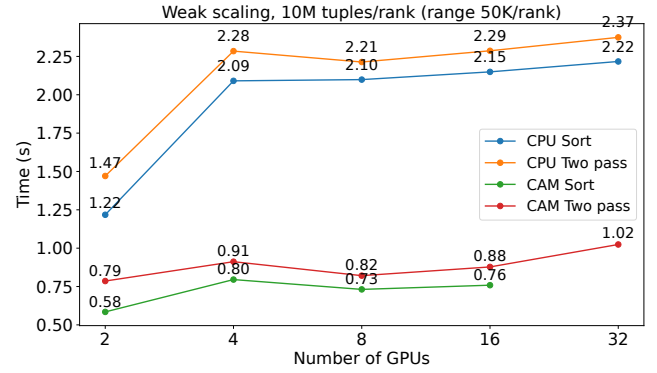


**Figure 7: Strong scaling time breakdown of a single iteration across key computational stages for both sorting-based and two-pass buffer preparation methods, evaluated under CPU buffer-based MPI and CUDA-aware MPI (CAM).**

weak scaling benchmarks, providing a granular breakdown of individual operations, including join operation, buffer preparation, communication (both pre-join and post-join), and deduplication.

*Strong scaling performance.* For strong scaling, we fixed the total dataset size to 10 million tuples and increased the number of GPUs from 2 to 32, effectively reducing the workload per rank. Figure 6 illustrates the execution time across different configurations. As the number of GPUs increases beyond 2, execution time decreases due to improved workload distribution and parallel processing. For both CUDA-Aware MPI (CAM) and CPU buffer-based MPI, the sorting-based buffer preparation shows better scaling than the two-pass approach. When comparing CAM to CPU buffer-based MPI, CAM achieves lower execution times across all scales due to direct GPU-to-GPU transfers, whereas CPU-based MPI incurs additional memory copies between host and device. Figure 7 provides a breakdown of execution time for strong scaling. The join operation remains relatively constant in time across all configurations, whereas buffer preparation and communication contribute the most significant overheads. The two-pass approach suffers from higher buffer preparation time, while the sorting method reduces this cost significantly. The transition from 4 to 8 GPU achieves the largest performance gain, as it enables workload distribution and parallelism, but further increasing the number of GPUs beyond 16 results in diminishing returns due to less workload on GPUs.

*Weak scaling performance.* For weak scaling, we maintained 10 million tuples per GPU, increasing both the dataset size and the number of GPUs proportionally. Figure 8 shows the execution time trends, CUDA-aware MPI (CAM) implementations sustain better performance with lower scaling overheads. In both all-to-all communication techniques, the sorting-based approach exhibits better scaling. The execution time became steady with using more than 4 GPUs. Across

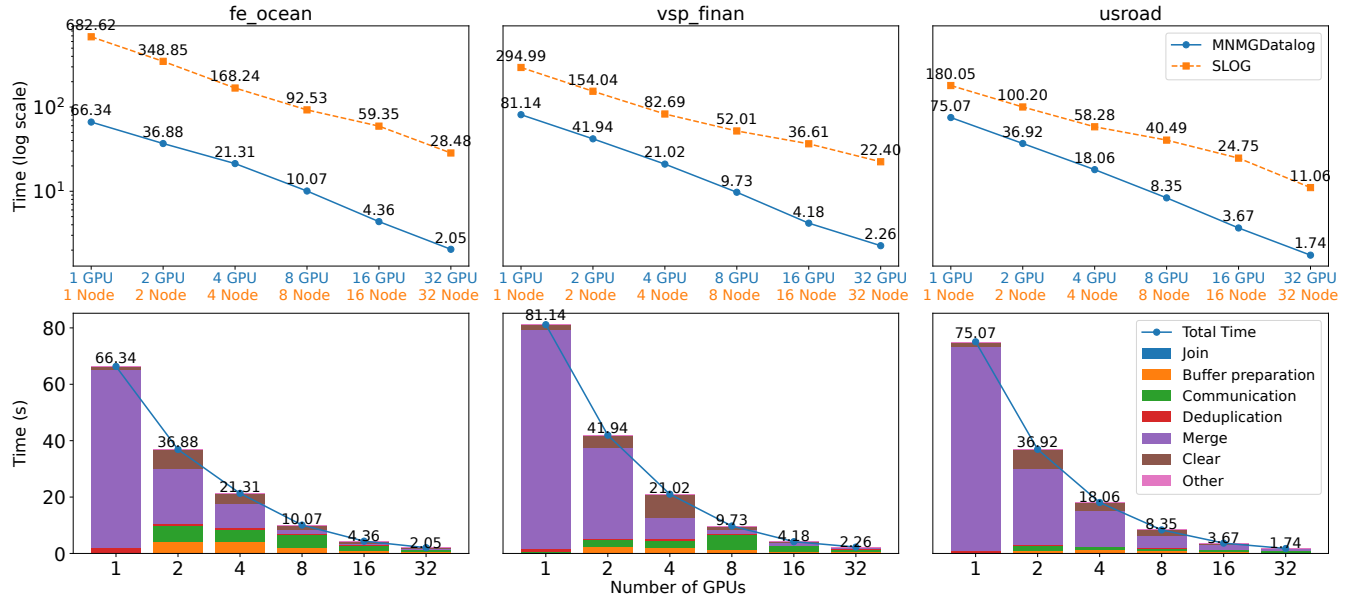


**Figure 8: Weak scaling performance of the single iteration of the fixed-point operation in MNMGDATALOG.**

both strong and weak scaling, the sorting based buffer preparation consistently outperforms two pass approaches, demonstrating the benefits of avoiding expensive atomic operations from two pass approach. CUDA-aware MPI (CAM) provides substantial performance gains over CPU buffer-based MPI by eliminating redundant host-device memory transfers, making it the preferred choice for GPU-accelerated join processing. The performance bottlenecks shift from join computation to buffer preparation and communication, indicating that optimizing these stages is crucial for improving the efficiency of iterative joins in distributed GPU environments.

## 5.5 Multi-node multi-GPU benchmark

We benchmark MNMGDATALOG on transitive closure (TC), same generation (SG), and weakly connected components (WCC) using up to 32 GPUs spanning multiple nodes on the Polaris supercomputer. As no existing multi-node, multi-GPU Datalog engine is available, we compare MNMGDATALOG against the state-of-the-art distributed CPU-based Datalog engine, SLOG, for transitive closure computation.



**Figure 9: (Top row) Execution time comparison of transitive closure (TC) between MNMGDATALOG and SLOG, where SLOG uses 32 CPU threads per node, and MNMGDATALOG employs GPU acceleration with CPU-based MPI buffer communication. (Bottom row) Breakdown analysis of Transitive Closure (TC) execution on MNMGDATALOG illustrating the time distribution across key operations as the number of GPUs scales from 1 to 32.**

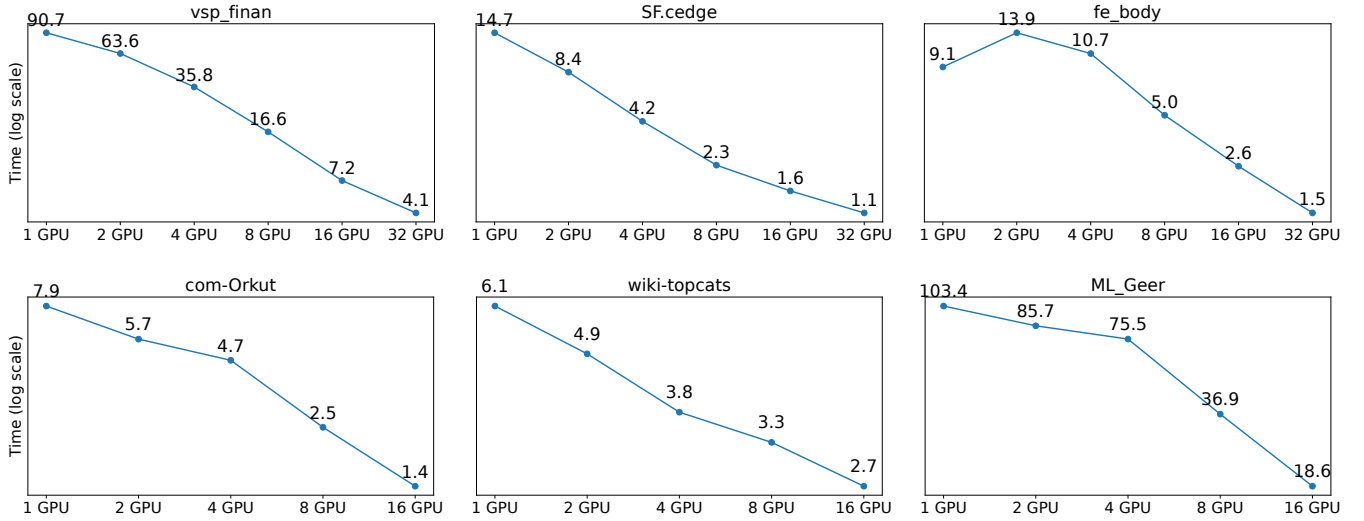
We intentionally used host-side (CPU) buffers for MPI communication to ensure a fair and neutral comparison with SLOG, which also relies on CPU buffers for data exchange. By configuring MNMGDATALOG in the same way, we isolated the impact of algorithmic design from hardware acceleration, allowing us to assess the raw algorithmic advantage under equivalent conditions. The performance benefits of CUDA-Aware MPI were evaluated separately in the single-join experiments (Section 5.4), where we demonstrated how our engine could further accelerate communication when GPU-to-GPU transfers are enabled.

*Transitive closure.* Table 3 presents the performance comparison of TC execution across multiple GPUs. MNMGDATALOG consistently outperforms SLOG. MNMGDATALOG achieves up to 32× speedup at 1 GPU and 13.89× speedup at 32 GPUs over SLOG. While MNMGDATALOG maintains a clear advantage across all configurations, the performance gap between MNMGDATALOG and SLOG narrows as the number of GPUs increases. This is due to the decreased workload per GPU on higher scales. However, even with this diminishing gap, MNMGDATALOG continues to exhibit superior scalability due to its optimized join processing and reduced memory overhead in recursive query execution, whereas SLOG experiences significant overhead from CPU-bound relational operations. Figure 9 top row further illustrates the scaling trends, emphasizing that MNMGDATALOG exhibits near-linear scaling as

the number of GPUs increases from 1 GPU to 32 GPUs. The speedup for *fe\_ocean* ranges from 7.9× to 13.9×, *vsp\_finan* from 3.6× to 9.9×, and *usroad* from 2.4× to 6.8× when scaling from 1 to 32 GPUs, compared to SLOG on the same nodes.

Figure 9 bottom row provides a detailed breakdown of the execution time for TC across different GPU configurations. The results show that join time and communication time decrease significantly as the number of GPUs increases, demonstrating the effectiveness of workload distribution across multiple nodes. Join operations benefit from parallelism, while communication overhead is reduced as individual GPU workloads become smaller with increasing GPU counts. However, merge and memory clearing time are the dominant contributors to the total execution time. This is expected, as merging the intermediate results in each iteration requires allocations/deallocations of GPU memory during the iterations. Additionally, deduplication time reduces at higher GPU counts, as smaller partitions per GPU lead to more efficient duplicate elimination. This highlights the capability of MNMGDATALOG to effectively scale distributed Datalog queries in a multi-node, multi-GPU environment, making it a robust solution for large-scale recursive query processing.

*Same generation benchmarking.* The Same Generation (SG) query determines whether two nodes in a directed graph belong to the same hierarchical level. It is defined recursively



**Figure 10: (Top-row) Same graph (SG) generation scaling from 1 to 32 GPUs. (Bottom-row) Weakly connected component (WCC) scaling from 1 to 16 GPUs.**

**Table 3: Transitive Closure (TC) runtime (s) comparison: SLOG vs. MNMGDATALOG. Scaling SLOG from 1–32 nodes (32 CPU threads per node) and MNMGDATALOG from 1–32 GPUs via CPU-buffered MPI (sorting-based buffer preparation).**

Dataset Name	TC Edges	Datalog Engine	1 Node 1 GPU	2 Nodes 2 GPUs	4 Nodes 4 GPUs	8 Nodes 8 GPUs	16 Nodes 16 GPUs	32 Nodes 32 GPUs
vsp_finan	910M	SLOG	294.99	154.04	82.69	52.01	36.61	22.40
		MNMGDATALOG	<b>81.14</b>	<b>41.94</b>	<b>21.02</b>	<b>9.73</b>	<b>4.18</b>	<b>2.26</b>
usroads	871M	SLOG	180.05	100.20	58.28	40.49	24.75	11.06
		MNMGDATALOG	<b>75.07</b>	<b>36.92</b>	<b>18.06</b>	<b>8.35</b>	<b>3.67</b>	<b>1.74</b>
fe_ocean	1.669B	SLOG	682.62	348.85	168.23	92.53	59.35	28.48
		MNMGDATALOG	<b>66.34</b>	<b>36.88</b>	<b>21.31</b>	<b>10.07</b>	<b>4.36</b>	<b>2.05</b>
Gnutella31	884M	SLOG	315.52	143.85	58.32	31.70	16.69	10.56
		MNMGDATALOG	<b>9.86</b>	<b>6.66</b>	<b>5.26</b>	<b>2.67</b>	<b>1.42</b>	<b>0.77</b>

as follows:

$$\begin{aligned}
 SG(u, v) &\leftarrow Edge(p, u), Edge(p, v), u \neq v. \\
 SG(u, v) &\leftarrow Edge(x, u), SG(x, y), Edge(y, v), u \neq v.
 \end{aligned}$$

The first rule captures direct relationships where two nodes share a common predecessor, while the second rule extends this recursively by checking for intermediate connections. Figure 10 top row illustrates the SG execution times across multiple GPUs using CPU-buffered MPI communication with sorting-based buffer preparation. We achieve 6× to 22× speedup from 1 to 32 GPUs. For the *fe\_body* dataset, we observe an anomaly where execution time increases from 1 GPU to 2 GPUs before improving with additional GPUs. This behavior is likely due to communication and partitioning overhead outweighing the benefits of parallelism at this scale. When moving from 1 GPU to 2 GPUs, data redistribution introduces

inter-GPU communication, which incurs latency, especially for datasets where the computation-to-communication ratio is not sufficiently high.

*Weakly connected component.* Figure 10 bottom row presents the WCC execution time as the number of GPUs increases from 1 to 16. *com-Orkut* achieves 5.7× speedup, while *ML\_Geer* scales 5.5×, *wiki-topcats* shows a 2.3× speedup. The performance gain is largely attributed to local materialization, where each GPU retains its partial connected component state, reducing inter-GPU communication. This prevents redundant updates from being exchanged in every iteration, ensuring that only minimal data is transferred. The hash-based partitioning strategy further optimizes execution by keeping most component updates local, limiting global synchronization overhead. Scaling efficiency varies across datasets due to differences in graph connectivity. Graphs with denser connectivity require frequent inter-GPU communication, impacting performance; by contrast, sparser graphs benefit from localized processing, yielding better scalability.

## 6 Related work

*GPU-based Datalog.* Accelerating Datalog engines with GPUs has been a long-standing goal of the Datalog community. Early attempts, such as GPUDatalog [39] and RedFox [57], failed to gain traction due to their inability to efficiently handle iterative queries using semi-naïve evaluation, a core feature of high-performance Datalog engines. Additionally, limited GPU VRAM (under 16GB) required frequent host-to-device data transfers, further constraining performance and scalability. Consequently, GPU-based Datalog was largely overlooked for years. Recent advancements in



GPU VRAM and computational power have revived interest in GPU-based Datalog systems. GPUJoin [49] demonstrated that GPU-based Datalog could outperform optimized CPU engines, though it was limited to binary relations and a narrow set of queries. Inspired by GPUJoin, GPULog [54] became the first GPU Datalog engine to support all fundamental relational algebra operations and semi-naïve evaluation, leveraging the novel HISA data structure for efficient execution.

Modern datacenter GPUs typically support advanced GPU-to-GPU interconnects, offering significantly higher memory bandwidth compared to GPU-to-host data transfers. We view MNMGDATALOG as a critical extension of GPU-based Datalog systems to leverage this trend. By scaling the number of GPUs in the system, MNMGDATALOG can handle significantly larger databases while maintaining low communication overhead.

*Distributed Datalog.* There has been significant progress in scaling Datalog-like languages to large machine clusters. Systems such as RDFox[41], BigDatalog[48], SociaLite[45], Myria[27], Nexus[29], and Radlog [26] have effectively utilized Apache Spark clusters to achieve scalability in data size. However, the query performance of these systems often can't scale beyond ten nodes. A more recent MPI-based distributed Datalog engine, SLOG [23], has demonstrated promising results in performance scaling, achieving near-linear scaling up to 64 nodes and gradual saturation up to 256 nodes on ANL's Theta supercomputer. The design of MNMGDATALOG draws inspiration from SLOG, adapting its data partitioning and communication techniques to the GPU. By integrating CUDA-aware MPI, MNMGDATALOG is optimized to scale efficiently on modern high-performance computing clusters, leveraging the computational power and interconnect capabilities of GPUs.

*Monotonic Aggregation and Semiring Provenance.* Datalog's basic semantics are simple, but extensions are needed for complex real-world queries in domains like program and graph analysis. Recursive aggregation, which allows monotonically updating existing tuples, is a widely adopted extension supported by CPU-based engines like BigDatalog[48], RecStep[19], and Logica [51]. The GPU-based system Lobster [10] formalizes monotonic aggregation using semiring provenance [25, 59], supporting various neural-symbolic reasoning queries, but it is limited to single-GPU setups. Efficient multi-GPU monotonic aggregation remains an open problem. In MNMGDATALOG, we experimentally support multi-GPU monotonic aggregation with a specialized comparator during deduplication. While limited to single-integer column aggregation, this represents a step toward general multi-GPU support.

## 7 Conclusion and future work

We presented MNMGDATALOG, the first-ever multi-node, multi-GPU Datalog engine, designed for efficient execution of recursive queries over internet-scale datasets at unprecedented levels of scalability. Our approach integrates a radix-hash-based data partitioning strategy with CUDA-aware non-uniform all-to-all communication. Our benchmarks demonstrate that MNMGDATALOG is the highest-performance Datalog engine to date, beating the previously-SOTA GPU-based competitor (GPULog) by 7×, the SOTA CPU-based engine (Soufflé) by up to 33×, and the distributed supercomputing competitor (SLOG) by up to 32×. The data, code, and documentation are all open-sourced and can be found at <https://github.com/harp-lab/MNMGDatalog/>.

While MNMGDATALOG shows impressive performance, we have several planned enhancements to improve its robustness and portability. For long-running workloads, we aim to implement per-iteration checkpoint/restart capabilities, allowing the system to capture execution state at arbitrary points and recover from failures mid-execution. To address load imbalance due to data skew, we are exploring a two-level strategy: work-stealing among SXM-connected GPUs within a node[55], and sub-bucketing across nodes to improve inter-node load distribution. In addition, we plan to support recent advances in space-optimized join processing, including parallel worst-case optimal joins and query decomposition techniques [36, 58], which help conserve valuable GPU memory. Finally, we are extending MNMGDATALOG to support a broader range of HPC platforms by integrating GPU-aware MPI implementations available in ROCm [55] for AMD GPUs and OneAPI [30] for Intel GPUs.

## Acknowledgments

This work was funded in part by NSF PPOSS large grants CCF-2316159 and CCF-2316157. 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. This material is based upon work supported by the Defense Advanced Research Projects Agency (DARPA) under Contract No. N66001-21-C-4023. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of DARPA.

## References

- [1] Ahmad Abdelfattah, Azzam Haidar, Stanimire Tomov, and Jack Dongarra. 2017. Novel HPC techniques to batch execution of many variable size BLAS computations on GPUs. In *Proceedings of the International Conference on Supercomputing*. 1–10.
- [2] Ahmad Abdelfattah, David Keyes, and Hatem Ltaief. 2016. Kblas: An optimized library for dense matrix-vector multiplication on gpu

- accelerators. *ACM Transactions on Mathematical Software (TOMS)* 42, 3 (2016), 1–31.
- [3] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. *Foundations of databases*. Vol. 8. Addison-Wesley Reading.
  - [4] AMD. 2024. 5TH GEN AMD EPYC™ PROCESSOR ARCHITECTURE. <https://www.amd.com/content/dam/amd/en/documents/epyc-business-docs/white-papers/5th-gen-amd-epyc-processor-architecture-white-paper.pdf>.
  - [5] Argonne Leadership Computing Facility. 2022. Polaris. <https://www.alcf.anl.gov/polaris>.
  - [6] Ammar Ahmad Awan, Khaled Hamidouche, Akshay Venkatesh, and Dhabaleswar K Panda. 2016. Efficient large message broadcast using NCCL and CUDA-aware MPI for deep learning. In *Proceedings of the 23rd European MPI Users' Group Meeting*. 15–22.
  - [7] George Balatsouras and Yannis Smaragdakis. 2016. Structure-sensitive points-to analysis for C and C++. In *Static Analysis: 23rd International Symposium, SAS 2016, Edinburgh, UK, September 8–10, 2016, Proceedings* 23. Springer, 84–104.
  - [8] Cagri Balkesen, Gustavo Alonso, Jens Teubner, and M. Tamer Özsu. 2013. Multi-core, main-memory joins: sort vs. hash revisited. *Proc. VLDB Endow.* 7, 1 (Sept. 2013), 85–96. doi:10.14778/2732219.2732227
  - [9] Claude Barthels, Ingo Müller, Timo Schneider, Gustavo Alonso, and Torsten Hoefer. 2017. Distributed join algorithms on thousands of cores. *Proceedings of the VLDB Endowment* 10, 5 (2017), 517–528.
  - [10] PAUL BIBERSTEIN, ZIYANG LI, JOSEPH DEVIETTI, and MAYUR NAIK. [n. d.]. Lobster: A GPU-Accelerated Framework for Neurosymbolic Programming. ([n. d.]).
  - [11] Peter A Boncz, Stefan Manegold, Martin L Kersten, et al. 1999. Database architecture optimized for the new bottleneck: Memory access. In *VLDB*, Vol. 99. 54–65.
  - [12] Martin Bravenboer and Yannis Smaragdakis. 2009. Strictly declarative specification of sophisticated points-to analyses. In *Proceedings of the 24th ACM SIGPLAN conference on Object oriented programming systems languages and applications*. 243–262.
  - [13] Stefano Ceri, Georg Gottlob, Letizia Tanca, et al. 1989. What you always wanted to know about Datalog (and never dared to ask). *IEEE transactions on knowledge and data engineering* 1, 1 (1989), 146–166.
  - [14] Evgeny Dantsin, Thomas Eiter, Georg Gottlob, and Andrei Voronkov. 2001. Complexity and Expressive Power of Logic Programming. *ACM Comput. Surv.* 33, 3 (sep 2001), 374–425. doi:10.1145/502807.502810
  - [15] Timothy A. Davis and Yifan Hu. 2011. The University of Florida Sparse Matrix Collection. *ACM Trans. Math. Softw.* 38, 1, Article 1 (dec 2011), 25 pages. doi:10.1145/2049662.2049663
  - [16] Oege De Moor, Georg Gottlob, Tim Furche, and Andrew Sellers. 2012. *Datalog Reloaded: First International Workshop, Datalog 2010, Oxford, UK, March 16–19, 2010. Revised Selected Papers*. Vol. 6702. Springer.
  - [17] Ke Fan, Thomas Gilray, Valerio Pascucci, Xuan Huang, Kristopher Micinski, and Sidharth Kumar. 2022. Optimizing the bruck algorithm for non-uniform all-to-all communication. In *Proceedings of the 31st International Symposium on High-Performance Parallel and Distributed Computing*. 172–184.
  - [18] Ke Fan, Steve Petruzza, Thomas Gilray, and Sidharth Kumar. 2024. Configurable Algorithms for All-to-All Collectives. In *ISC High Performance 2024 Research Paper Proceedings (39th International Conference)*. Prometheus GmbH, 1–12.
  - [19] Zhiwei Fan, Jianqiao Zhu, Zuyu Zhang, Aws Albarghouthi, Paraschos Koutris, and Jignesh M. Patel. 2019. Scaling-up in-Memory Datalog Processing: Observations and Techniques. *Proc. VLDB Endow.* 12, 6 (Feb 2019), 695–708. doi:10.14778/3311880.3311886
  - [20] Kayvon Fatahalian, Jeremy Sugerman, and Pat Hanrahan. 2004. Understanding the efficiency of GPU algorithms for matrix-matrix multiplication. In *Proceedings of the ACM SIGGRAPH/EUROGRAPHICS conference on Graphics hardware*. 133–137.
  - [21] Antonio Flores-Montoya and Eric Schulte. 2020. Datalog disassembly. In *29th USENIX Security Symposium (USENIX Security 20)*. 1075–1092.
  - [22] Shinya Fushimi, Masaru Kitsuregawa, and Hidehiko Tanaka. 1986. An Overview of The System Software of A Parallel Relational Database Machine GRACE.. In *VLDB*, Vol. 86. 209–219.
  - [23] Thomas Gilray, Arash Sahebolamri, Yihao Sun, Sowmith Kunapaneni, Sidharth Kumar, and Kristopher Micinski. 2024. Datalog with First-Class Facts. *Proc. VLDB Endow.* 18, 3 (Nov. 2024), 651–665. doi:10.14778/3712221.3712232
  - [24] Oded Green. 2021. HashGraph—Scalable hash tables using a sparse graph data structure. *ACM Transactions on Parallel Computing (TOPC)* 8, 2 (2021), 1–17.
  - [25] Todd J Green, Grigoris Karvounarakis, and Val Tannen. 2007. Provenance semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. 31–40.
  - [26] Jiaqi Gu, Yugo H Watanabe, William A Mazza, Alexander Shkap-sky, Mohan Yang, Ling Ding, and Carlo Zaniolo. 2019. RaSQL: Greater power and performance for big data analytics with recursive-aggregate-SQL on Spark. In *Proceedings of the 2019 International Conference on Management of Data*. 467–484.
  - [27] Daniel Halperin, Victor Teixeira de Almeida, Lee Lee Choo, Shumo Chu, Paraschos Koutris, Dominik Moritz, Jennifer Ortiz, Vaspol Rumviboonsuk, Jingjing Wang, Andrew Whitaker, Shengliang Xu, Magdalena Balazinska, Bill Howe, and Dan Suciu. 2014. Demonstration of the Myria Big Data Management Service. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data (Snowbird, Utah, USA) (SIGMOD '14)*. Association for Computing Machinery, New York, NY, USA, 881–884. doi:10.1145/2588555.2594530
  - [28] Pieter Hijma, Stijn Heldens, Alessio Slococo, Ben van Werkhoven, and Henri E. Bal. 2023. Optimization Techniques for GPU Programming. *ACM Comput. Surv.* 55, 11, Article 239 (March 2023), 81 pages. doi:10.1145/3570638
  - [29] Muhammad Imran, Gábor E Gévy, Jorge-Arnulfo Quiané-Ruiz, and Volker Markl. 2022. Fast datalog evaluation for batch and stream graph processing. *World Wide Web* 25, 2 (2022), 971–1003.
  - [30] Intel. 2023. Intel MPI for GPU Clusters. <https://www.intel.com/content/www/us/en/docs/oneapi/optimization-guide-gpu/2023-2/intel-mpi-for-gpu-clusters.html>.
  - [31] Jiri Kraus. 2013. An Introduction to CUDA-Aware MPI. <https://developer.nvidia.com/blog/introduction-cuda-aware-mpi/>.
  - [32] Herbert Jordan, Bernhard Scholz, and Pavle Subotić. 2016. Soufflé: On synthesis of program analyzers. In *Computer Aided Verification: 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17–23, 2016, Proceedings, Part II* 28. Springer, 422–430.
  - [33] Changkyu Kim, Tim Kaldewey, Victor W Lee, Eric Sedlar, Anthony D Nguyen, Nadathur Satish, Jatin Chhugani, Andrea Di Blas, and Pradeep Dubey. 2009. Sort vs. hash revisited: Fast join implementation on modern multi-core CPUs. *Proceedings of the VLDB Endowment* 2, 2 (2009), 1378–1389.
  - [34] Sidharth Kumar and Thomas Gilray. 2019. Distributed relational algebra at scale. In *International Conference on High Performance Computing, Data, and Analytics (HiPC)*. IEEE, Vol. 1.
  - [35] Sidharth Kumar and Thomas Gilray. 2020. Load-balancing parallel relational algebra. In *High Performance Computing: 35th International Conference, ISC High Performance 2020, Frankfurt/Main, Germany, June 22–25, 2020, Proceedings* 35. Springer, 288–308.
  - [36] Zhuohang Lai, Xibo Sun, Qiong Luo, and Xiaolong Xie. 2022. Accelerating multi-way joins on the GPU. *The VLDB Journal* (2022), 1–25.
  - [37] Jure Leskovec and Andrej Krevl. 2014. SNAP Datasets: Stanford Large Network Dataset Collection. <http://snap.stanford.edu/data>.

- [38] Feifei Li, Dihan Cheng, Marios Hadjieleftheriou, George Kollios, and Shang-Hua Teng. 2005. On trip planning queries in spatial databases. In *International symposium on spatial and temporal databases*. Springer, 273–290.
- [39] Carlos Alberto Martínez-Angeles, Inês Dutra, Vítor Santos Costa, and Jorge Buenabad-Chávez. 2013. A datalog engine for gpus. In *International Conference on Applications of Declarative Programming and Knowledge Management*. Springer, 152–168.
- [40] Adithya Murali, Atharva Sehgal, Paul Krogmeier, and P Madhusudan. 2019. Composing neural learning and symbolic reasoning with an application to visual discrimination. *arXiv preprint arXiv:1907.05878* (2019).
- [41] Yavor Nenov, Robert Piro, Boris Motik, Ian Horrocks, Zhe Wu, and Jay Banerjee. 2015. RDFox: A highly-scalable RDF store. In *The Semantic Web-ISWC 2015: 14th International Semantic Web Conference, Bethlehem, PA, USA, October 11-15, 2015, Proceedings, Part II 14*. Springer, 3–20.
- [42] NVIDIA. 2022. NVIDIA H100 Tensor Core GPU. <https://www.nvidia.com/en-us/data-center/h100/>.
- [43] NVIDIA. 2024. CUDA C Programming Guide: SIMT. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/#simt-architecture>.
- [44] NVIDIA. 2025. CUDA Thrust API documentation: thrust::unique function. [https://nvidia.github.io/cccl/thrust/api/function\\_group\\_stream\\_compaction\\_1gacfc33f1e24f8526b003f8a679591ad65.html](https://nvidia.github.io/cccl/thrust/api/function_group_stream_compaction_1gacfc33f1e24f8526b003f8a679591ad65.html).
- [45] Jiwon Seo, Stephen Guo, and Monica S Lam. 2013. Socialite: Datalog extensions for efficient social network analysis. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*. IEEE, 278–289.
- [46] Andres Sewell, Ke Fan, Ahmedur Rahman Shovon, Landon Dyken, Sidharth Kumar, and Steve Petruzza. 2024. Bruck algorithm performance analysis for multi-gpu all-to-all communication. In *Proceedings of the International Conference on High Performance Computing in Asia-Pacific Region*. 127–133.
- [47] Ramtin Shams, RA Kennedy, et al. 2007. Efficient histogram algorithms for NVIDIA CUDA compatible devices. In *Proc. Int. Conf. on Signal Processing and Communications Systems (ICSPCS)*. Citeseer, 418–422.
- [48] Alexander Shkapsky, Mohan Yang, Matteo Interlandi, Hsuan Chiu, Tyson Condie, and Carlo Zaniolo. 2016. Big Data Analytics with Datalog Queries on Spark. In *Proceedings of the 2016 International Conference on Management of Data (San Francisco, California, USA) (SIGMOD '16)*. Association for Computing Machinery, New York, NY, USA, 1135–1149. doi:10.1145/2882903.2915229
- [49] Ahmedur Rahman Shovon, Landon Richard Dyken, Oded Green, Thomas Gilray, and Sidharth Kumar. 2022. Accelerating Datalog applications with cuDF. In *2022 IEEE/ACM Workshop on Irregular Applications: Architectures and Algorithms (IA3)*. IEEE, 41–45.
- [50] Ahmedur Rahman Shovon, Thomas Gilray, Kristopher Micinski, and Sidharth Kumar. 2023. Towards iterative relational algebra on the {GPU}. In *2023 USENIX Annual Technical Conference (USENIX ATC 23)*. 1009–1016.
- [51] Evgeny Skvortsov, Yilin Xia, and Bertram Ludäscher. 2024. Logica: Declarative Data Science for Mere Mortals.. In *EDBT*. 842–845.
- [52] Elias Stehle and Hans-Arno Jacobsen. 2017. A memory bandwidth-efficient hybrid radix sort on gpus. In *Proceedings of the 2017 ACM International Conference on Management of Data*. 417–432.
- [53] Yihao Sun, Sidharth Kumar, Thomas Gilray, and Kristopher Micinski. 2023. Communication-Avoiding Recursive Aggregation. In *2023 IEEE International Conference on Cluster Computing (CLUSTER)*. IEEE, 197–208.
- [54] Yihao Sun, Ahmedur Rahman Shovon, Thomas Gilray, Sidharth Kumar, and Kristopher Micinski. 2025. Optimizing Datalog for the GPU. In *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1 (Rotterdam, Netherlands) (ASPLOS '25)*. Association for Computing Machinery, New York, NY, USA, 762–776. doi:10.1145/3669940.3707274
- [55] Yiltan Hassan Temuçin, Mahdieh Gazimirsaeed, Ryan E. Grant, and Ahmad Afsahi. 2024. ROCm-Aware Leader-based Designs for MPI Neighbourhood Collectives. In *ISC High Performance 2024 Research Paper Proceedings (39th International Conference)*. 1–12. doi:10.23919/ISC.2024.10528923
- [56] Hao Wang, Sreeram Potluri, Devendar Bureddy, Carlos Rosales, and Dhabaleswar K Panda. 2013. GPU-aware MPI on RDMA-enabled clusters: Design, implementation and evaluation. *IEEE Transactions on Parallel and Distributed Systems* 25, 10 (2013), 2595–2605.
- [57] Haicheng Wu, Gregory Diamos, Tim Sheard, Molham Aref, Sean Baxter, Michael Garland, and Sudhakar Yalamanchili. 2014. Red fox: An execution environment for relational query processing on gpus. In *Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization*. 44–54.
- [58] Hangdong Zhao, Shaleen Deep, and Paraschos Koutris. 2023. Space-Time Tradeoffs for Conjunctive Queries with Access Patterns. In *Proceedings of the 42nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (Seattle, WA, USA) (PODS '23)*. Association for Computing Machinery, New York, NY, USA, 59–68. doi:10.1145/3584372.3588675
- [59] Hangdong Zhao, Shaleen Deep, Paraschos Koutris, Sudeepa Roy, and Val Tannen. 2024. Evaluating Datalog over Semirings: A Grounding-based Approach. *Proc. ACM Manag. Data* 2, 2, Article 90 (May 2024), 26 pages. doi:10.1145/3651591