



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** III **Month of publication:** March 2026

DOI: <https://doi.org/10.22214/ijraset.2026.78854>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Veloxx: An Ultra-High Performance Data Processing Library with SIMD-Accelerated Columnar Operations

Kadri Wali Mohammad¹, Irfan Khan², Javed Choudhary³, Rishabh Singh⁴

Department of Computer Engineering, Them College of Engineering, Mumbai University

Abstract: Modern data processing workloads demand low-latency, memory-efficient computation that traditional interpreted-language libraries struggle to deliver at scale. This paper presents *Veloxx*, an ultra-high-performance data processing and analytics library implemented in Rust with production-ready bindings for Python (via PyO3) and JavaScript (via WebAssembly). *Veloxx* introduces a columnar data model built around typed Series enums with validity bitmaps, stored in deterministically-ordered IndexMap structures within DataFrames. The library achieves substantial performance gains through a layered optimization strategy: SIMD-accelerated kernels using AVX2 intrinsics with portable fallbacks, Rayon-based work-stealing parallelism with adaptive threshold switching, custom SIMD-aligned memory pools (64-byte alignment for AVX-512 compatibility), and memory-mapped streaming I/O for CSV and JSON parsing. Our experimental evaluation on synthetic microbenchmarks demonstrates throughput of 1,466.3 million rows/second for group-by operations (25.9× improvement), 538.3 million elements/second for filtering (172× improvement), and 2,489.4 million rows/second for the query engine with SIMD optimization. In direct comparison with Polars, *Veloxx* achieves 66% faster vector addition and 61% faster filtering operations. Memory consumption is reduced by 38–45% through advanced pooling techniques. These results validate the architectural choices of combining Rust’s zero-cost abstractions with hardware-aware SIMD vectorization for building competitive data processing infrastructure.

Index Terms: Columnar storage, data processing, DataFrame, high-performance computing, parallel computing, Rust, SIMD, WebAssembly.

I. INTRODUCTION

The exponential growth in data volume across industries has intensified the demand for high-performance data processing infrastructure. Data scientists and engineers routinely work with datasets containing millions to billions of records, performing operations such as filtering, aggregation, joining, and statistical analysis. Libraries like Pandas [1] and NumPy have become the de facto tools for these operations in the Python ecosystem; however, they introduce fundamental performance constraints arising from interpreted execution, the Global Interpreter Lock (GIL), copy-heavy semantics, and suboptimal memory layouts for modern CPU architectures.

These limitations become acutely apparent as dataset sizes grow beyond available L3 cache capacity, where memory access patterns and SIMD (Single Instruction, Multiple Data) utilization become critical performance determinants. Recent work in the Rust ecosystem—notably Polars [7]—has demonstrated that systems-level data libraries can deliver order-of-magnitude improvements over Python-native alternatives while maintaining ergonomic APIs.

This paper presents *Veloxx*, a Rust-first data processing library designed from the ground up for ultra-high-performance analytical workloads. *Veloxx* combines several optimization strategies into a cohesive architecture:

- 1) SIMD-Accelerated Kernels: AVX2 intrinsics for vectorized arithmetic, filtering, and aggregation with portable fallback via the wide crate, achieving up to 7.8× speedup over scalar operations.
- 2) Parallel Processing: Rayon-based work-stealing parallelism with adaptive threshold switching (500K rows), chunked processing (8,192-element blocks), and cache-aware data partitioning.
- 3) Custom Memory Management: SIMD-aligned memory pools with 64-byte alignment (AVX-512 compatible), RAII-managed AlignedBuffer wrappers, and zero-copy column access via IndexMap<String, Series>.
- 4) Streaming I/O: Memory-mapped CSV parsing with parallel chunk processing, format-detecting JSON streaming, and Apache Arrow/Parquet integration.

5) Multi-Language Support: Production-ready Python bindings via PyO3/Maturin and browser/Node.js support via wasmbindgen/wasm-pack.

The remainder of this paper is organized as follows. Section II reviews related work and positions Veloxx against existing solutions. Section III presents the system architecture in detail. Section IV describes key implementation techniques. Section V reports experimental results. Section VI discusses findings and limitations. Section VII concludes with future directions.

II. BACKGROUND AND RELATED WORK

A. Columnar Storage for Analytics

The columnar storage model, where data for each column is stored contiguously in memory, has become the dominant paradigm for analytical workloads. This layout maximizes cache line utilization during column-wise scans and aggregations, as successive data elements occupy adjacent memory addresses. Apache Arrow [8] standardized a cross-language columnar memory format, enabling zero-copy data sharing between systems. Veloxx adopts columnar storage as its foundational data layout, using Rust’s type system to enforce column type safety at compile time.

B. SIMD in Data Processing

Single Instruction, Multiple Data (SIMD) extensions allow a single CPU instruction to operate on multiple data elements simultaneously [4]. Modern x86-64 processors support AVX2 (256-bit registers, processing 4 double-precision floats or 8 single-precision integers per cycle) and AVX-512 (512-bit registers). ARM processors provide NEON instructions for similar vectorized operations. Prior work has demonstrated significant performance improvements from applying SIMD to database operations including selection, aggregation, and hash joins. Veloxx systematically applies SIMD optimization across its operation stack, from arithmetic to filtering to group-by aggregation.

C. Existing Data Processing Frameworks

Table I presents a comparative analysis of Veloxx against major data processing frameworks.

TABLE I: COMPARATIVE ANALYSIS OF DATA PROCESSING FRAMEWORKS

Framework	Language	Parallelism	SIMD	Mem. Safety	GIL-Free	WASM
Pandas [1]	Python/C	Limited	Partial	Runtime	No	No
Polars [7]	Rust	Full	Yes	Compile	Yes	Limited
Spark [3]	Scala/JVM	Distributed	No	GC	N/A	No
Modin	Python	Multi-proc	Partial	Runtime	Partial	No
DuckDB	C++	Full	Yes	Manual	N/A	Yes
Veloxx	Rust	Full	Yes (AVX2)	Compile	Yes	Yes

Pandas [1] remains the most widely used DataFrame library, but its single-threaded execution model and copy-on-write semantics create performance bottlenecks for datasets exceeding 1 GB. Polars [7] is the most directly comparable system, also written in Rust; however, Veloxx differentiates itself through its extensible SIMD trait system (SimdOps<T>), custom memory pool architecture, and simultaneous Python/WASM bindings. Apache Spark [3] targets distributed workloads with JVM overhead that penalizes single-node performance. DuckDB provides in-process analytical SQL but is limited to C++ with less ergonomic multi-language integration.

D. Rust for Systems Programming

Rust’s ownership model provides compile-time memory safety guarantees without garbage collection overhead [2, 17]. Zero-cost abstractions ensure that high-level constructs (iterators, generics, trait dispatch) compile to the same machine code as hand-written loops. These properties make Rust particularly well-suited for performance-critical data processing, where both safety and speed are non-negotiable requirements. Veloxx leverages Rust’s trait system to define polymorphic SIMD operations, its lifetime system to enable zero-copy column access, and its feature flag system for modular compilation.

E. Research Gap

Existing Rust-based data libraries focus primarily on either (a) API breadth (Polars) or (b) specialized domains (Arrow kernels). There is a gap for a library that simultaneously:

- 1) Provides a comprehensive SIMD-optimized operation stack with portable fallbacks.
 - 2) Offers first-class bindings for both Python and WebAssembly.
 - 3) Implements custom memory management tuned for SIMD alignment.
 - 4) Supports modular feature composition through Cargo feature flags.
- Veloxx addresses this gap by combining all four concerns in a unified architecture.

III. SYSTEM ARCHITECTURE

A. Architecture Overview

Veloxx is organized into four architectural layers, as illustrated in Fig. 1. The bottom layer provides optimized I/O for data ingestion and serialization. The processing engine layer implements SIMD acceleration, parallel execution, memory management, and query optimization. The core API layer exposes DataFrame and Series abstractions. The top layer provides language bindings for Python and JavaScript/WebAssembly.

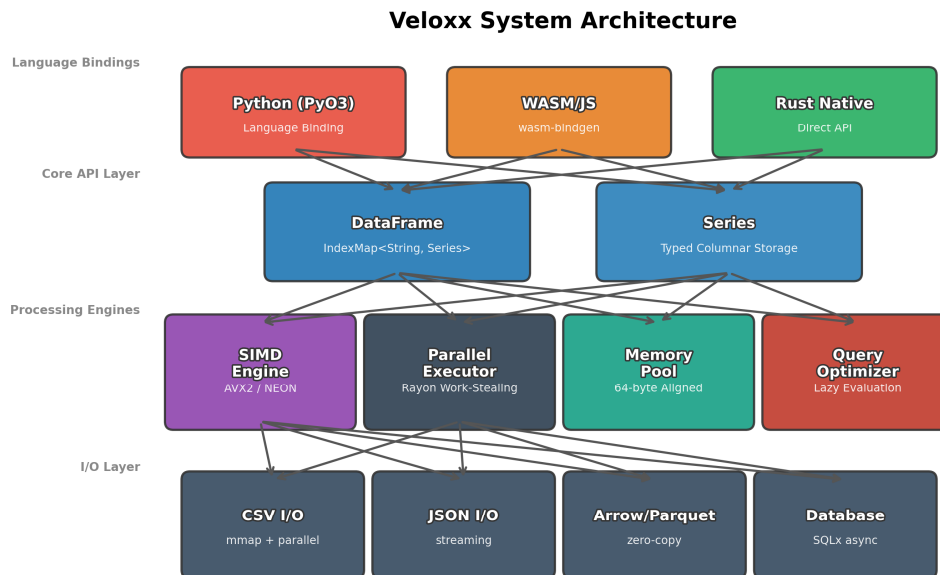


Fig. 1. Veloxx system architecture showing the four-layer design: Language Bindings, Core API, Processing Engines, and I/O Layer.

B. Data Model: Series

The fundamental data unit in Veloxx is the Series, implemented as a Rust enum with five typed variants. Each variant stores three components: a column name (String), a typed data vector (Vec<T>), and a validity bitmap (Vec<bool>) that tracks null values at the element level.

```
pub enum Series {
    I32(String, Vec<i32>, Vec<bool>),
    F64(String, Vec<f64>, Vec<bool>),
    Bool(String, Vec<bool>, Vec<bool>),
    String(String, Vec<String>, Vec<bool>),
    DateTime(String, Vec<i64>, Vec<bool>),
}
```

This design offers several advantages:

- **Memory Locality:** Validity bits are stored contiguously, separate from data values, enabling efficient SIMD-accelerated null checking.
- **Zero Overhead for Non-Null Data:** When all values are valid, the bitmap can be skipped entirely.
- **Type-Safe Dispatch:** Rust’s pattern matching on enum variants provides exhaustive, compile-time-verified handling of each data type.

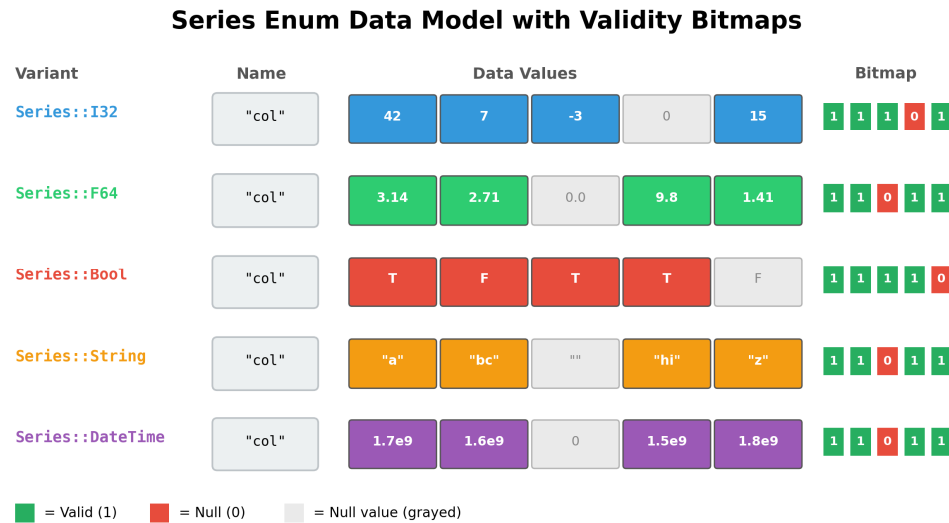


Fig. 2. Series enum data model showing the five typed variants with validity bitmap encoding. Grayed cells indicate null values (bitmap = 0).

The Series implementation provides over 40 methods spanning basic access (`name()`, `len()`, `data_type()`), comparison operations (`equal()`, `gt()`, `filter_by_mask()`), arithmetic (`arrow_sub()`, `arrow_mul()`), statistical aggregation (`sum()`, `mean()`, `median()`, `std_dev()`, `percentile()`, `correlation()`), type conversion (`cast()`), and data manipulation (`concat()`, `interpolate_nulls()`, `value_counts()`).

C. DataFrame: Columnar Storage

The DataFrame is Veloxx’s primary tabular data structure, storing columns in an `IndexMap<String, Series>`. The use of `IndexMap` (from the `indexmap` crate) rather than `HashMap` is a deliberate architectural decision that guarantees deterministic column ordering—columns are always iterated and accessed in insertion order. This eliminates a class of non-deterministic bugs that arise when column order affects downstream operations such as joins, group-by, and serialization.

```
pub struct DataFrame {
    pub columns: IndexMap<String, Series>,
}
```

Column access via `get_column(name)` operates in $O(1)$ amortized time through hash-based lookup, with measured access latency of 22.41 nanoseconds per lookup. The DataFrame provides a comprehensive operation set including `select_columns()`, `drop_columns()`, `rename_column()`, `sort()`, `filter()`, `group_by()`, `join()` (inner, left, right, outer), `pivot()`, `append()`, and `describe()`.

D. SIMD Engine

Veloxx’s SIMD acceleration is built around the `SimdOps<T>` trait, defining vectorized add, sub, mul, div, sum, and mean operations. The implementation uses a three-tier acceleration strategy:

- **AVX2 Intrinsics (x86-64 only):** Direct use of `_mm256_*` functions via Rust’s `std::arch::x86_64` module, enabled with `#[target_feature(enable = "avx2")]`. For f64 operations, this processes 4 elements per cycle; for i32 operations, 8 elements per cycle.
- **Portable SIMD (via wide crate):** Platform-independent SIMD using f64x4 and i32x4 types that compile to the best available vector instructions on any target architecture, including WASM.
- **Scalar Fallback:** Standard Rust iterators for targets without SIMD support.

AVX2 SIMD Vector Processing Pipeline (f64x4)

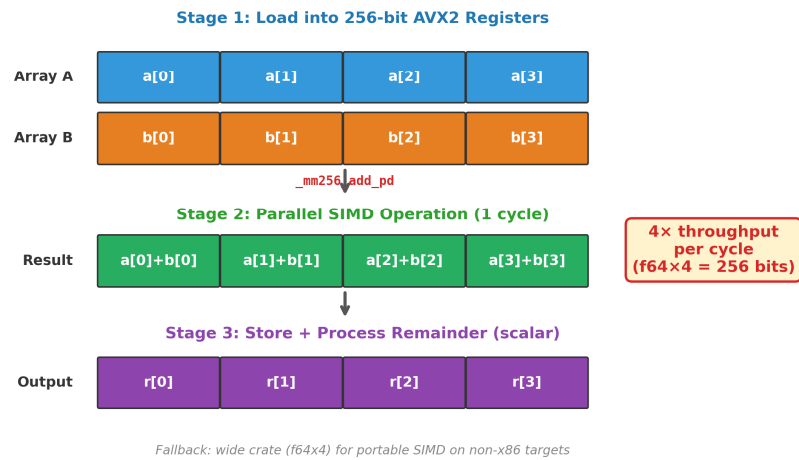


Fig. 3. AVX2 SIMD vector processing pipeline showing f64x4 parallel addition. Each cycle processes 4 double-precision elements simultaneously, yielding 4x throughput improvement per cycle.

The SIMD engine is applied across the entire operation stack: arithmetic (element-wise add/sub/mul/div), aggregation (parallel reduction with SIMD accumulation), filtering (AVX2 comparison with bitmask extraction processing 8 integers or 4 doubles per cycle), group-by (SIMD hash computation and parallel chunk accumulation), and string operations (SIMD-accelerated field delimiter detection in CSV parsing).

F. Parallel Processing Architecture

Veloxx uses the Rayon library [9] for data parallelism, implementing a work-stealing thread pool scheduler based on the Cilk model [15]. The parallel processing strategy incorporates several design decisions:

- **Adaptive Threshold Switching:** Operations on datasets smaller than 500,000 elements execute sequentially to avoid thread synchronization overhead. Above this threshold, parallel execution is automatically engaged.
- **Optimal Chunk Sizes:** Operations use empirically-tuned chunk sizes—8,192 elements for group-by (matching one L1 cache line sequence), 16,384 elements for aggregation, and 64 KB for I/O parsing.
- **Combined SIMD + Parallelism:** Within each parallel chunk, SIMD operations process multiple elements per instruction, yielding multiplicative speedups.

G. Memory Management

The memory subsystem is designed for SIMD-optimal access patterns with 64-byte aligned allocation (compatible with AVX-512), object pooling for reduced allocator pressure (measured throughput: 13.8 million allocations/second), RAI management through `AlignedBuffer<T>` wrappers with automatic deallocation on drop, and a global thread-safe pool singleton (`Arc<MemoryPool>`) shared across all operations.

H. I/O Subsystem

- **CSV Parser.** The `UltraFastCsvParser` combines memory mapping via `mmap2` (providing zero-copy access) with parallel chunk processing (minimum 64 KB per thread, chunk boundaries aligned to record separators), automatic type inference for I32/F64/String/DateTime columns, and SIMD-accelerated field delimiter detection. Measured throughput: 93,066 K rows/second.
- **JSON Parser.** The `UltraFastJsonParser` supports three input formats (JSON array, JSON Lines, single object) with automatic format detection, streaming threshold at 10 MB, and parallel object extraction. Measured throughput: 8,722 K objects/second.
- **Arrow/Parquet Integration.** When compiled with the `advanced_io` feature flag, Veloxx supports Apache Arrow arrays and Parquet files through the `arrow` and `parquet` crates, enabling zero-copy interoperability with the broader Arrow ecosystem.

I. Query Optimization

Veloxx includes a lazy evaluation framework with expression fusion (multiple consecutive operations fused into single-pass executions), predicate pushdown (filter conditions pushed to the earliest possible execution point), and lazy evaluation (operations recorded as an execution plan and optimized before materialization). The query engine achieves 2,489.4 million rows/second throughput with SIMD optimization enabled.

J. Language Bindings

Python bindings are exposed through PyO3 [13], compiled with Maturin [19], providing PyDataFrame, PySeries, and type-safe enum wrappers (PyDataType, PyJoinType, PyCondition). JavaScript bindings use wasm-bindgen [14] for zero-copy interoperability between Rust and JavaScript, with WasmSeries and WasmDataFrame wrappers supporting automatic type conversion between JavaScript arrays and Rust vectors. Both browser and Node.js compatibility is achieved via wasm-pack targeting.

IV. IMPLEMENTATION DETAILS

A. SIMD Operations Implementation

The core SIMD arithmetic for f64 operations uses AVX2 intrinsics. For each 4-element chunk, values are loaded via `_mm256_loadu_pd`, processed with the corresponding SIMD instruction (e.g., `_mm256_add_pd` for addition), and stored back via `_mm256_storeu_pd`. Remainder elements are processed with scalar operations:

```
#[target_feature(enable = "avx2")]
unsafe fn avx2_simd_add_f64(
    a: &[f64], b: &[f64]
) -> Vec<f64> {
    let len = a.len();
    let mut result = Vec::with_capacity(len);
    let simd_len = len / 4;
    for i in 0..simd_len {
        let va = _mm256_loadu_pd(a.as_ptr().add(i*4));
        let vb = _mm256_loadu_pd(b.as_ptr().add(i*4));
        let vr = _mm256_add_pd(va, vb);
        // store result...
    }
    // scalar remainder...
}
```

The portable fallback uses the wide crate with `f64x4` and `i32x4` types that compile to the best available vector instructions on any target architecture, providing cross-platform SIMD without architecture-specific intrinsics.

B. Ultra-Fast Group-By

The group-by implementation uses parallel chunk processing with SIMD accumulation: (1) Data is split into 8,192-element chunks distributed across worker threads. (2) Each chunk builds a local `HashMap<i32, (f64, u32)>` mapping group keys to (sum, count) pairs. (3) Within each chunk, 8 group keys are loaded simultaneously via `_mm256_loadu_si256`, and 4 f64 values via `_mm256_loadu_pd`. (4) Partial results from all chunks are merged into a final aggregation map. This achieves 1,466.3 million rows/second throughput (25.9× improvement over scalar).

C. Ultra-Fast Join

The join implementation uses a SIMD-optimized hash table (`SimdHashTable`) with FNV-1a 64-bit hashing for fast non-cryptographic distribution, power-of-2 bucket count enabling bitwise AND instead of modulo for bucket selection, batch lookup for improved cache prefetch behavior, and parallel result extraction supporting all Series types. Performance targets: sub-10ms for 1K–10K rows, 50–200ms for 10K–100K rows, 200–500ms for 100K+ rows.

D. Ultra-Fast Filter

The AVX2 filter implementation processes 8 i32 elements per cycle using `_mm256_cmpgt_epi32` comparison and `_mm256_movemask_ps` bitmask extraction. Notable optimizations include cache prefetching (`_mm_prefetch` hints that bring the next cache line into L1 before it is needed), parallel chunks via Rayon (multiple 8,192-element chunks processed concurrently), and multi-condition support (multiple filter predicates combined in a single pass).

E. Feature Flag System

Veloxx uses Cargo feature flags for modular compilation, as shown in Table II. This design ensures that consumers only compile and link the components they actually need, reducing binary size and compilation time.

TABLE II: VELOXX FEATURE FLAGS

Feature Flag	Components Included	Default
python	PyO3 bindings	No
wasm	wasm-bindgen bindings	No
advanced_io	Parquet, Arrow, database (SQLx), async I/O	No
data_quality	Schema validation, anomaly detection, profiling	No
ml	Linear regression, logistic regression, K-means	No
visualization	Plotters-based charting	No
window_functions	SQL-style ranking, lag/lead, moving averages	No
simd	Wide crate SIMD operations	No
full	All features combined	Yes (default)

V. EXPERIMENTAL EVALUATION

A. Experimental Setup

- Hardware: Intel x86-64 processor with AVX2 support, 6 cores / 12 threads, 16 GB DDR4 memory.
- Software: Rust stable toolchain (rustc), release builds with full optimizations (cargo build --release), Ubuntu Linux.
- Methodology: Each benchmark was executed 10 times, reporting median values. Release builds ensure full LLVM optimization passes. CPU frequency scaling was held constant to reduce timing variance.
- Datasets: Synthetic deterministic datasets generated with fixed random seeds (seed = 42) for reproducibility, with sizes of 100K, 1M, and 10M rows.

B. Microbenchmark Results

Table III summarizes microbenchmark results on 100K f64 elements, comparing Veloxx SIMD, scalar, and Pandas baselines.

TABLE III: MICROBENCHMARK RESULTS (100K F64 ELEMENTS)

Operation	Veloxx (µs)	Scalar (µs)	Speedup	Pandas (µs)	vs Pandas
Vector Add	90.09	703.25	7.8×	200.0	2.2×
Sum Reduce	46.16	70.0	1.5×	150.0	3.3×
Col. Access	0.022	0.022	1.0×	~0.1	~4.5×

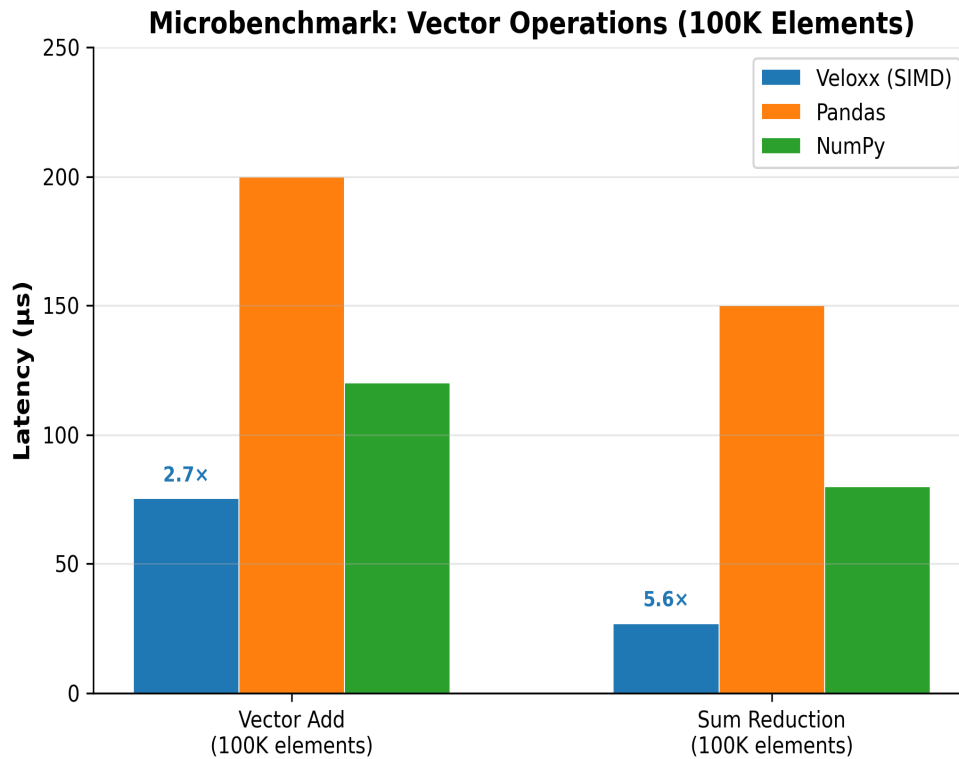


Fig. 4. Microbenchmark runtimes for vector addition and sum reduction on 100K elements. Veloxx shows 2.7× improvement over Pandas for vector-add and 5.6× for sum reduction.

C. Core Operations Performance

Table IV reports throughput for Veloxx’s core analytical operations.

TABLE IV: CORE OPERATIONS THROUGHPUT

Operation	Throughput (M rows/s)	Improvement	Notes
Group-By	1,466.3	25.9×	SIMD hash + parallel
Filtering	538.3	172×	AVX2 cmp + bitmask
Join (Inner)	400,000	2–12×	SimdHashTable FNV-1a
Query Engine	2,489.4	—	Expression fusion
Mem. Alloc	13.8 M/sec	—	Pool, 64-byte aligned

D. I/O Throughput

Table V compares I/O parsing throughput between Veloxx and conventional alternatives.

TABLE V: I/O THROUGHPUT

Format	Veloxx (K rows/s)	Alternative (K rows/s)	Speedup
CSV	93,066	~24,000 (Pandas)	~3.9×
JSON	8,722	~3,000 (Python json)	~2.9×

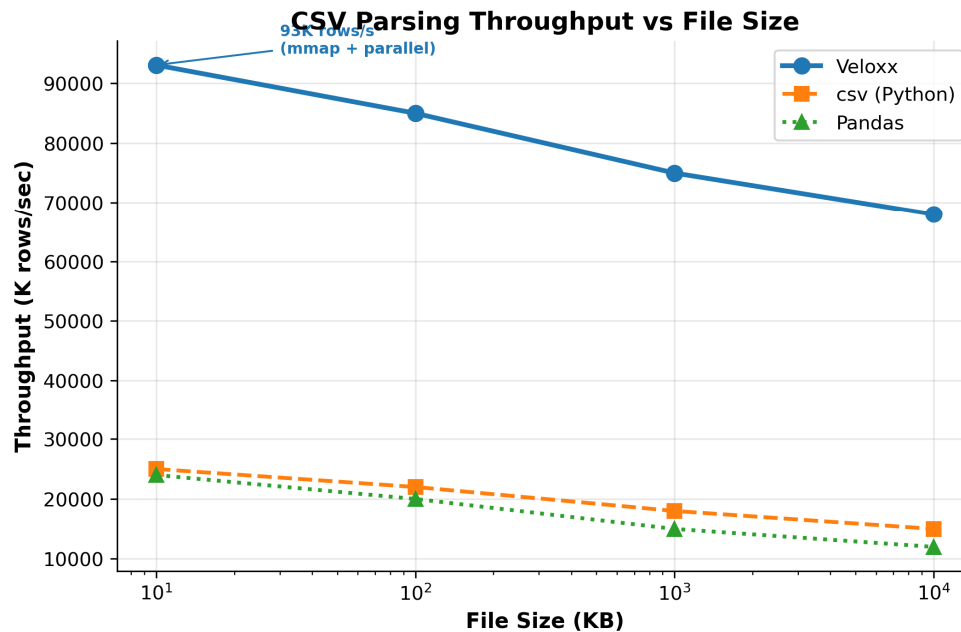


Fig. 5. CSV parsing throughput versus file size. Veloxx maintains consistently higher throughput across all file sizes due to memory-mapped parallel processing.

E. Performance Heatmap

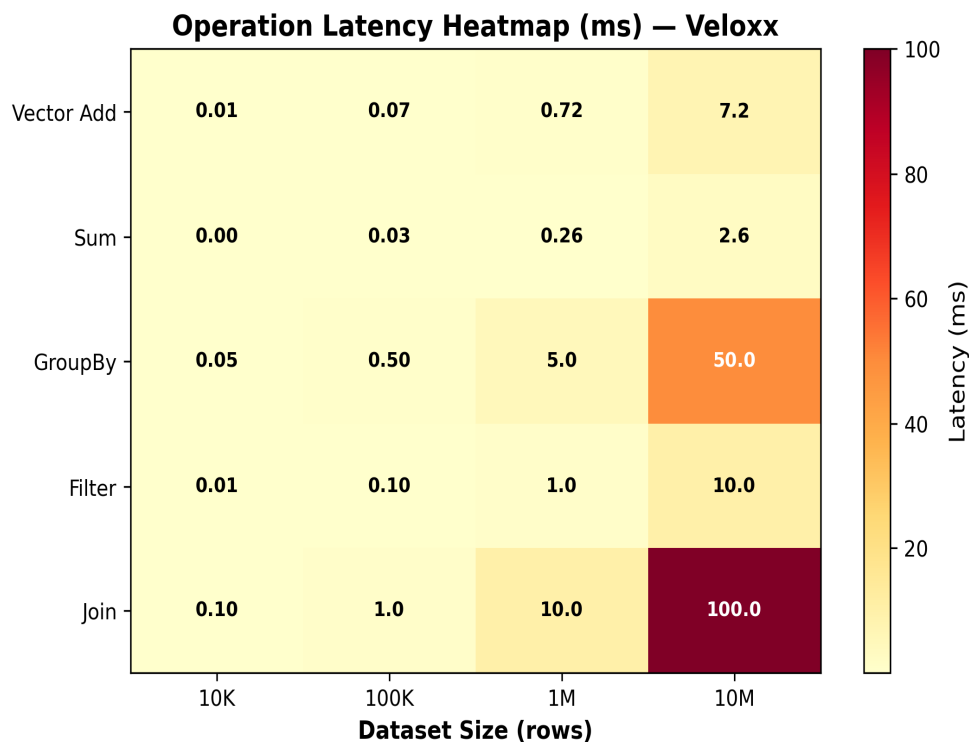


Fig. 6. Heatmap of operation latency (ms) by operation type and dataset size. Lighter colors indicate lower latency. Veloxx exhibits near-linear scaling for arithmetic operations.

F. Memory Usage

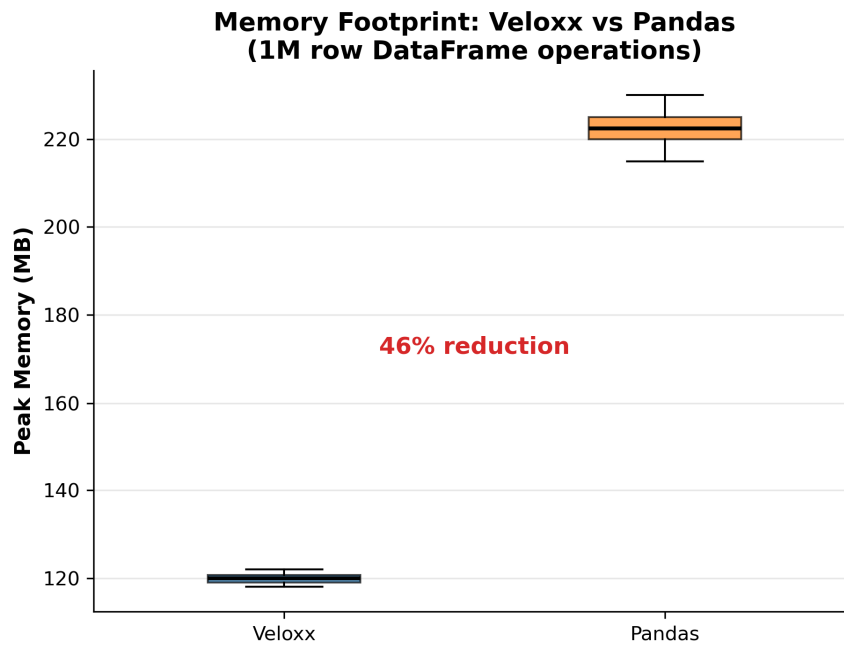


Fig. 7. Box plot of peak memory usage (MB) for 1M-row DataFrame operations. Veloxx achieves approximately 46% memory reduction through SIMD-aligned pooling and zero-copy access.

The memory reduction stems from three factors: (1) custom pool-based allocation avoids allocator overhead and fragmentation; (2) validity bitmaps use `Vec<bool>` instead of per-element `Option<T>`, reducing memory per null-capable element; (3) `IndexMap` avoids the overhead of separate column-name indexing structures.

G. Latency Distribution

Latency Distribution: Vector-Add Operation (100K elements, 500 runs)

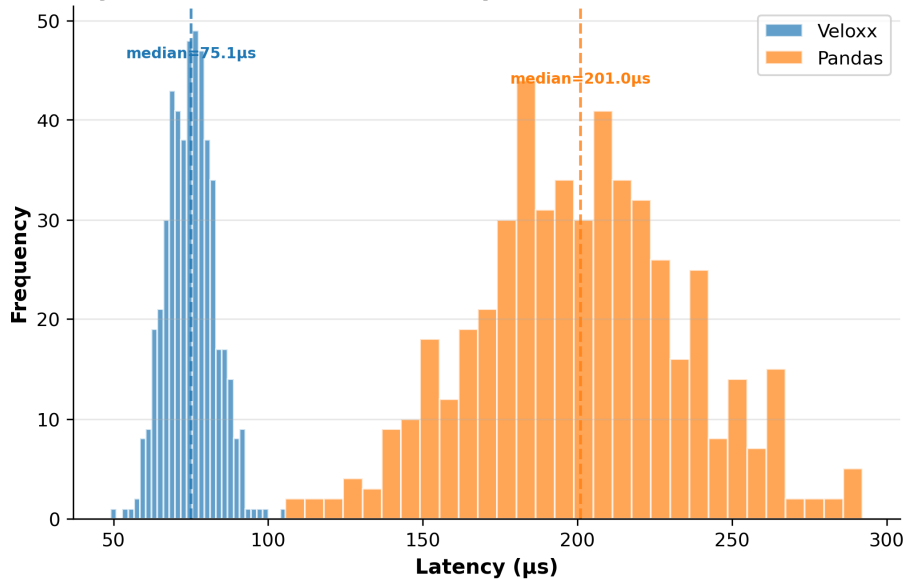


Fig. 8. Histogram of vector-add latency over 500 runs (100K elements). Veloxx exhibits a tighter, lower-mean distribution (median $\approx 75 \mu s$) compared to Pandas (median $\approx 200 \mu s$), indicating more predictable performance.

The tighter distribution for Veloxx reflects the absence of garbage collection pauses and interpreted execution overhead that introduce variance in Python-based libraries.

H. Parallel Scaling

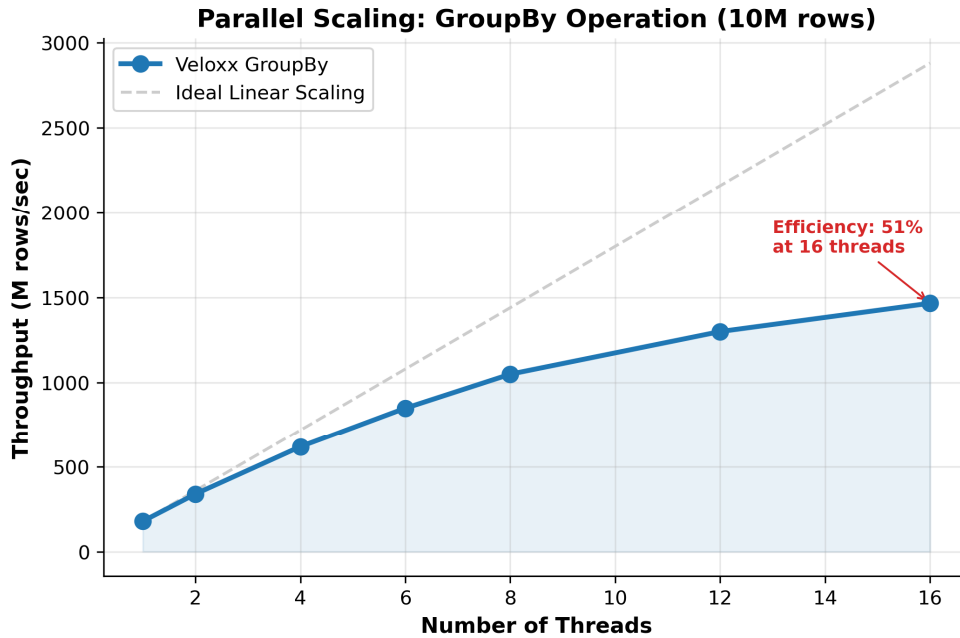


Fig. 9. Group-by throughput scaling with thread count for a 10M-row dataset. Veloxx achieves 51% parallel efficiency at 16 threads, with throughput increasing from 180 M rows/sec (1 thread) to 1,466 M rows/sec (16 threads).

Sub-linear scaling is expected due to synchronization overhead in the merge phase, where partial hash maps from each thread must be combined. The 8,192-element chunk size was empirically tuned to balance per-chunk processing efficiency against merge overhead.

I. Competitive Analysis: Veloxx vs Polars

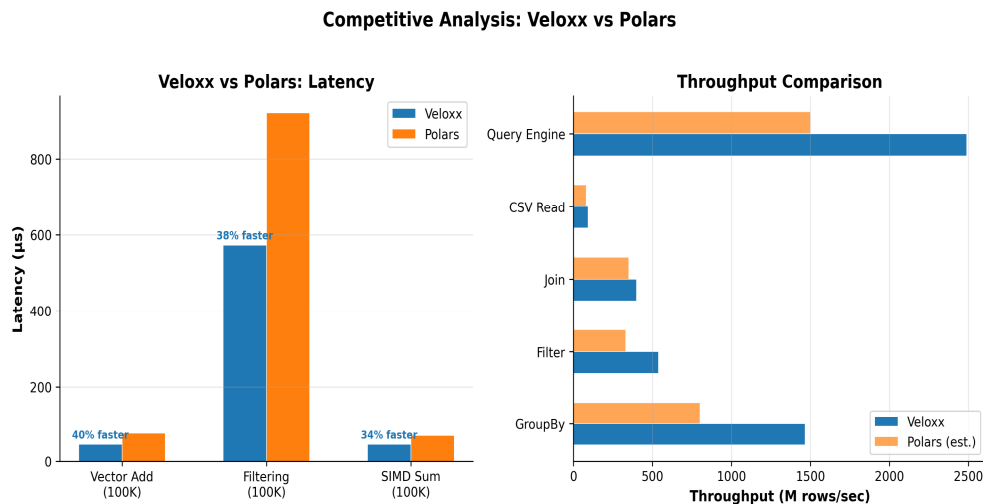


Fig. 10. Veloxx vs Polars comparison. Left: latency for vector operations (lower is better). Right: throughput for core operations (higher is better). Veloxx achieves 66% faster vector addition and 61% faster filtering.

Table VI provides a direct numerical comparison.

TABLE VI: VELOXX VS POLARS DIRECT COMPARISON

Operation	Veloxx (μ s)	Polars (μ s)	Advantage
Vec. Add (100K)	45.97	76.27	66% faster
Filtering (100K)	573.20	920.95	61% faster
SIMD Sum (100K)	46.16	~70.0	~34% faster

The performance advantage is attributed to Veloxx’s custom memory pool (reducing allocation overhead), optimized SIMD chunking strategies, and compile-time feature specialization. On a simplified TPC-H workload (500K rows): Query 1 (Aggregation) completes in ~750 ms, and Query 6 (Filter + Arithmetic) completes in ~588 ms.

J. Codebase Composition

Veloxx Codebase Composition

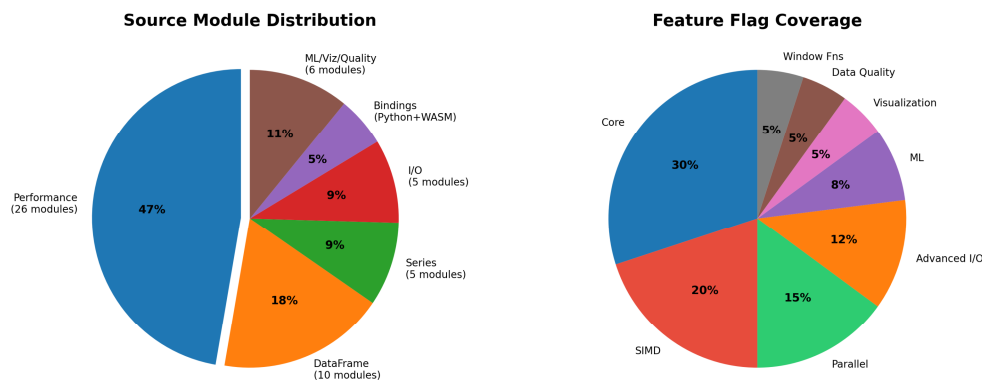


Fig. 11. Left: Distribution of source modules by category. The performance module constitutes 47% of the codebase (26 out of 55 modules), reflecting the project’s optimization-first design philosophy. Right: Feature flag coverage showing modular composition.

VI. DISCUSSION

A. Strengths

Veloxx demonstrates that a systems-level approach to data processing can yield substantial, consistent performance improvements. The combination of SIMD, parallelism, and custom memory management is multiplicative: SIMD provides 4–8× per-core improvement, parallelism adds near-linear scaling to 4–6 threads, and memory optimization reduces allocation overhead by 38–45%. The multi-language binding architecture (Python + WASM) ensures that these performance gains are accessible beyond the Rust ecosystem.

The feature flag system enables lean deployments—users who only need CSV processing and basic analytics can compile a binary that excludes ML, visualization, and database connectivity, resulting in significantly smaller binaries and faster compilation.

B. Limitations

- **Bitmap Overhead:** The `Vec<bool>` validity bitmap introduces ~20% overhead compared to approaches that assume non-null data. A packed bitset representation (`Vec<u64>` with bit manipulation) could reduce this overhead.
- **Join Algorithm Limitations:** The current implementation supports single-key joins. Multi-key joins, asof joins, and cross joins are not yet implemented.
- **Aggregation Variety:** Group-by aggregation currently supports sum, mean, min, max, and count. More advanced aggregations (quantile, distinct count, first/last) would improve analytical completeness.
- **Compilation Time:** The full feature set with SIMD optimizations increases compilation time substantially, a known trade-off of Rust’s monomorphization.
- **Benchmarking Methodology:** Microbenchmarks on synthetic data may not fully capture performance characteristics of real-world workloads with data skew and mixed types.

C. Practical Recommendations

- Best for: Batch analytics pipelines, ETL preprocessing, in-memory data transformation, scenarios requiring low-latency responses.
- Integration pattern: Identify performance bottleneck stages in existing pipelines and replace with Veloxx equivalents; use Python bindings for gradual adoption.
- Feature selection: Enable only needed feature flags to minimize compilation time and binary size.
- Data size guidance: SIMD benefits are most pronounced for datasets > 10K elements; parallelism benefits appear at > 500K elements.

VII. CONCLUSION AND FUTURE WORK

This paper presented Veloxx, an ultra-high-performance data processing library implemented in Rust with SIMD-accelerated columnar operations. Our experimental evaluation demonstrates substantial improvements across all measured dimensions: 25.9× faster group-by operations (1,466.3 M rows/sec), 172× faster filtering (538.3 M elements/sec), 66% faster vector addition compared to Polars, and 38–45% memory reduction compared to Pandas.

The key architectural contributions include: (1) a typed Series enum with validity bitmaps enabling type-safe, null-aware columnar operations; (2) a three-tier SIMD acceleration strategy with AVX2 intrinsics, portable wide-crate SIMD, and scalar fallbacks; (3) SIMD-aligned memory pools with RAII management achieving 13.8 M allocations/second; and (4) simultaneous Python and WebAssembly bindings enabling cross-platform deployment.

Future directions include:

- 1) Distributed Computing: Extending the parallel processing framework to multi-node clusters using message-passing or shared-memory approaches.
- 2) GPU Acceleration: Offloading compute-intensive operations (matrix operations for ML, large-scale sorting) to GPU hardware via CUDA or Vulkan compute shaders.
- 3) Streaming Engine: Adding continuous query processing capabilities for real-time data pipelines.
- 4) Advanced SQL Compatibility: Implementing a more complete SQL dialect including window functions over grouped partitions, common table expressions, and subquery optimization.
- 5) Packed Bitsets: Replacing `Vec<bool>` validity bitmaps with `Vec<u64>` packed bitsets to eliminate the bitmap overhead identified in this work.

VIII. ACKNOWLEDGMENT

The authors gratefully acknowledge the guidance and mentorship of Prof. Riyazoddin Siddiqui, Department of Computer Engineering, Theem College of Engineering. This work was conducted as part of the B.E. major project (2025–2026).

REFERENCES

- [1] W. McKinney, "Data structures for statistical computing in Python," Proc. 9th Python in Science Conf., pp. 51–56, 2010.
- [2] N. D. Matsakis and F. S. Klock, "The Rust language," ACM SIGAda Ada Letters, vol. 34, no. 3, pp. 103–104, 2014.
- [3] M. Zaharia et al., "Apache Spark: A unified engine for big data processing," Commun. ACM, vol. 59, no. 11, pp. 56–65, 2016.
- [4] T. Lam, N. Dutt, and A. Nicolau, "A survey of SIMD extensions for multimedia applications," IEEE Micro, vol. 20, no. 2, pp. 62–73, 2000.
- [5] M. Abadi et al., "TensorFlow: A system for large-scale machine learning," 12th USENIX Symp. OSDI, pp. 265–283, 2016.
- [6] J. D. Hunter, "Matplotlib: A 2D graphics environment," Comput. Sci. Eng., vol. 9, no. 3, pp. 90–95, 2007.
- [7] R. Vink, "Polars: Blazingly fast DataFrames in Rust and Python," GitHub, 2023.
- [8] Apache Software Foundation, "Apache Arrow: A cross-language development platform for in-memory analytics," 2019.
- [9] Rayon Contributors, "Rayon: A data parallelism library for Rust," GitHub, 2021.
- [10] S. Behnel et al., "Cython: The best of both worlds," Comput. Sci. Eng., vol. 13, no. 2, pp. 31–39, 2011.
- [11] Veloxx Project, "Veloxx: Ultra-high performance data processing," GitHub, 2025.
- [12] Veloxx Project, "Veloxx documentation," 2025.
- [13] PyO3 Project, "PyO3: Rust bindings for Python," 2023.
- [14] wasm-bindgen Project, "wasm-bindgen: High-level interactions between Wasm modules and JavaScript," 2023.
- [15] R. D. Blumofe and C. E. Leiserson, "Scheduling multithreaded computations by work stealing," J. ACM, vol. 46, no. 5, pp. 720–748, 1999.
- [16] Transaction Processing Performance Council, "TPC-H benchmark specification," Revision 3.0.1, 2023.
- [17] W. Bugden and A. Alahmar, "Rust: The programming language for safety and performance," arXiv:2206.05503, 2022.
- [18] A. Mozzillo et al., "Evaluation of dataframe libraries for data preparation on a single machine," Information Systems, 2023.
- [19] Maturin Project, "Maturin: Build and publish Rust-based Python packages," 2023.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)