

Analysis of a Proposed Vectorized $O(1)$ Sorting Algorithm for Bounded Integer Ranges

I. Introduction

A. Overview of the Proposed Algorithm

This report provides an in-depth analysis of a proposed novel sorting algorithm designed to achieve constant time complexity, $O(1)$, or as close as practically feasible, for sorting integers within a known, bounded range $[0, k]$. The core concept hinges on leveraging vectorized operations, direct index mapping, and the massive parallelism inherent in modern hardware accelerators like Graphics Processing Units (GPUs) and Tensor Processing Units (TPUs). The algorithm aims to bypass traditional comparison-based sorting mechanisms, loops, and complex branching logic.

The fundamental idea involves two main phases executed in parallel:

1. **Vectorized Mapping:** Input array elements are treated as indices into an auxiliary frequency array (or histogram) of size $k+1$. Using a single vectorized pass, ideally executed in parallel across all input elements, the frequency of each value is recorded. For example, an input element x would trigger an increment at $\text{freq}[x]$.
2. **Vectorized Reconstruction:** The frequency array is then used to reconstruct the sorted output. This involves calculating the starting position for each distinct value (akin to a prefix sum or scan operation on the frequency array) and then placing the correct number of each value into the output array in parallel (akin to a scatter operation).

The user proposes exploiting libraries such as NumPy, JAX (leveraging XLA compilation), CuPy, or lower-level CUDA primitives (potentially via Thrust or CUB) to implement these steps efficiently.¹ The central claim is that under optimized conditions – specifically a fixed range k and sufficient hardware parallelism (e.g., GPUs/TPUs) – the wall-clock time for sorting could become effectively independent of the input size n , thus approaching $O(1)$ complexity in practice.

B. Report Objectives and Scope

The primary objective of this report is to conduct a rigorous, expert evaluation of the proposed algorithm. This analysis will scrutinize its theoretical underpinnings, particularly the ambitious $O(1)$ time complexity claim, and assess its practical feasibility and potential performance. The scope encompasses:

- A detailed examination of the algorithm's two core steps (mapping and reconstruction) and its critical assumptions, notably the bounded integer range

[0, k].

- An evaluation of the theoretical time complexity within the context of parallel computation models (e.g., PRAM) and a comparison against related algorithms like Counting Sort.
- An investigation into practical implementation challenges and opportunities using modern vectorization libraries (NumPy, JAX, CuPy) and GPU programming frameworks (CUDA/Thrust/CUB), considering factors like memory access patterns, atomic operations, and synchronization overheads.
- An assessment of the implications of the $O(k)$ space complexity.
- A comparison of the potential performance against established sorting algorithms (Quick Sort, Merge Sort, Radix Sort, Counting Sort) in the target domain of bounded integer inputs.
- An analysis of the algorithm's suitability for specific application domains, namely Machine Learning (ML) inference pipelines and embedded systems.
- An evaluation of the algorithm's novelty in light of existing sorting techniques and parallel computing primitives.
- A synthesis of these findings to provide a comprehensive verdict on the algorithm's promise, practicality for further development (prototyping, publication, library integration), and potential impact.

C. Structure of the Report

Following this introduction, the report is structured as follows:

- **Section II** delves into the internal mechanics of the algorithm, detailing the mapping and reconstruction steps, drawing analogies to known parallel patterns like histogramming and prefix sums, and analyzing the core assumptions and their implications.
- **Section III** provides a theoretical assessment of the time complexity, critically evaluating the $O(1)$ claim using parallel computation models (PRAM), comparing it to Counting Sort, analyzing the influence of input size n and range k , and considering theoretical sorting lower bounds.
- **Section IV** examines the practical implementation aspects using various libraries and hardware platforms (CPU/GPU/TPU), focusing on vectorization, atomic operations, memory patterns, synchronization, performance comparisons with standard sorts, and the crucial impact of $O(k)$ space complexity.
- **Section V** analyzes the applicability of the algorithm in specific domains, namely real-time Machine Learning inference and resource-constrained embedded systems.
- **Section VI** assesses the novelty of the proposed approach relative to existing

algorithms and GPU primitives, and discusses its potential for academic publication or integration into standard libraries.

- **Section VII** synthesizes the analysis, provides direct answers to the user's core questions, and offers concluding remarks and recommendations for future work.

II. Algorithm Internals and Core Assumptions

This section dissects the proposed algorithm's two primary steps, relating them to established parallel computing patterns, and examines the fundamental assumptions underpinning its design.

A. Step 1: Vectorized Mapping (Parallel Histogram Construction)

Description: The first step takes the input array `arr` of n integers, where each element `arr[i]` is guaranteed to be within the range $[0, k]$. The goal is to compute the frequency of each distinct integer value within this range. This is achieved by using an auxiliary array, `freq`, of size $k+1$, initialized to zeros. In a highly parallel execution model, each input element `arr[i]` directly contributes to updating the count at the corresponding index in the frequency array: `freq[arr[i]] += 1`. The user envisions this occurring in a single, vectorized pass where ideally all n updates happen concurrently.

Analogy to Histogramming: This step is functionally identical to constructing a histogram of the input data, where the integers 0 through k represent the bin identifiers.⁴ The value `freq[j]` after this step represents the number of input elements equal to j . Numerous parallel histogram algorithms have been developed, particularly for GPUs, often forming the basis for other computations like sorting.⁶ Libraries like NVIDIA's CUB provide optimized device-wide histogram primitives (`cub::DeviceHistogram`) that handle the complexities of parallel execution.⁹ Similarly, NumPy's `bincount` function performs this operation efficiently on CPUs.¹³

Parallelism Requirement: Achieving true parallelism where all n updates occur simultaneously presents a significant challenge due to potential *write conflicts* or *output interference*.⁴ If multiple input elements have the same value v , multiple threads or processing units will attempt to read, increment, and write back to the same memory location `freq[v]` concurrently. Without proper handling, this leads to race conditions, where updates can be lost, resulting in incorrect counts.⁴ Therefore, implementing this step correctly and efficiently in parallel necessitates mechanisms to manage these conflicts. The most common approach on parallel hardware like GPUs is the use of *atomic operations* (e.g., `atomicAdd`) which guarantee that the read-modify-write sequence for a single memory location occurs indivisibly, serializing concurrent accesses to the same location.⁴ Alternative strategies involve

privatization, where each thread block or even thread computes a local histogram in private or shared memory, followed by a separate reduction step to combine these local histograms into the final global freq array.⁴

B. Step 2: Vectorized Reconstruction (Parallel Scatter/Prefix Sum Application)

Description: The second step utilizes the computed frequency array `freq` to construct the final sorted output array, let's call it `sorted_arr`. The core idea is to determine the correct position for each element based on the counts of preceding elements. This typically involves two conceptual sub-steps:

1. **Offset Calculation:** Compute the starting index in `sorted_arr` for each distinct value `j` present in the input. This is equivalent to calculating the cumulative sum of frequencies for all values less than `j`. This operation is precisely a *prefix sum* (also known as *scan*) performed on the `freq` array.⁵ An *exclusive scan* produces an array, say `offsets`, where `offsets[j]` stores the sum of `freq` through `freq[j-1]`, representing the starting position for elements with value `j`.
2. **Data Scattering:** Populate the `sorted_arr`. For each value `j` from 0 to `k`, `freq[j]` copies of the value `j` are written into `sorted_arr` starting at index `offsets[j]`.

Analogy to Counting Sort Reconstruction / Parallel Scan: This reconstruction phase mirrors the final steps of the standard Counting Sort algorithm.⁵ Counting Sort uses the frequency counts to determine positions and then iterates through the input (or the counts) to place elements. The crucial operation for determining positions in parallel is the prefix sum (scan). Highly optimized parallel scan implementations are available in libraries like Thrust (`thrust::inclusive_scan`, `thrust::exclusive_scan`) and CUB (`cub::DeviceScan`), forming a fundamental building block for many parallel algorithms on GPUs.¹⁷ The NumPy implementation using `numpy.repeat` after `numpy.bincount` effectively performs this reconstruction step.¹³ The final placement of elements is a *scatter* operation, where elements (the values `j`) are written to computed indices (`offsets[j]` onwards).

Parallelism Requirement: This step can also be parallelized. The prefix sum on the `freq` array (size `k+1`) can be computed efficiently in parallel, typically in logarithmic time complexity with respect to `k` using sufficient processors.¹⁷ The scattering phase can be parallelized by assigning different threads or thread blocks to handle different ranges of values `j` or different segments of the output array `sorted_arr`. Each thread would read the corresponding `offsets[j]` and `freq[j]` and write the required number of `j` values to the output. Care must be taken to ensure correct indexing and avoid write conflicts if multiple threads are responsible for writing the same value `j`.

C. Core Assumptions and Constraints

The feasibility and performance of the proposed algorithm rest heavily on several key assumptions and constraints:

1. **Bounded Integer Range $[0, k]$:** This is the most critical constraint. The algorithm fundamentally relies on the ability to use input values directly as indices into the auxiliary freq array. The size of this array, and thus the memory footprint and potentially parts of the computation time, scale directly with k .
2. **Known Range:** The maximum value k must be known *a priori* to allocate the freq array of the correct size $(k+1)$. If the range is unknown or unbounded, this approach is inapplicable.
3. **Hardware Parallelism:** The algorithm implicitly assumes access to massively parallel hardware architectures (SIMD units on CPUs, many cores on GPUs/TPUs) capable of executing a large number of operations concurrently. The $O(1)$ claim is predicated on this high degree of parallelism.
4. **Efficient Vectorized Operations/Memory Access:** The proposal relies on the availability and efficiency of vectorized instructions and, crucially, efficient mechanisms for handling the parallel write operations in the mapping step (Step 1). This includes efficient atomic operations or alternative strategies supported by the hardware and software stack (e.g., JAX/XLA, CUDA).¹
5. **Implicit Assumptions:** The algorithm assumes fixed-size integer types for the input data and typically for the counters in the freq array (though larger counters might be needed if frequencies exceed the capacity of a standard integer).

D. Interconnections and Nuances

The description above highlights the core mechanics, but deeper consideration reveals important interdependencies and subtleties.

The efficiency of the reconstruction step (Step 2) is influenced by the characteristics of the frequency array generated in Step 1. If the histogram is very sparse (i.e., freq contains many zeros because the input n contains only a few unique values sparsely distributed within the large range k), the prefix sum calculation might be faster as many indices can be skipped. Similarly, the scatter operation only needs to write elements corresponding to non-zero frequencies. Conversely, a dense histogram requires processing nearly all k entries in the prefix sum and scattering elements for most values. This implies that the *distribution* of input values within the range $[0, k]$, not just the parameters n and k , can affect the practical runtime constants, even if the asymptotic goal is $O(1)$.

Furthermore, while the mapping step (Step 1) is analogous to the counting phase of

Counting Sort, the proposed implementation strategy differs significantly. Traditional sequential Counting Sort iterates through the n input elements to build the counts.¹⁸ Common parallel Counting Sort strategies involve partitioning the input across p processors, computing local histograms, and then merging these histograms, often leading to complexities like $O(n/p + p \cdot k)$ or involving logarithmic factors for merging.²⁸ The proposal here aims to circumvent explicit partitioning and merging by having all n (or n/p) parallel workers attempt to update a *single*, shared freq array concurrently using vectorized memory writes. This places the primary performance bottleneck squarely on the hardware's ability to resolve simultaneous write attempts to the same memory location (i.e., the efficiency and scalability of atomic operations).⁴ The $O(1)$ ambition therefore hinges critically on the assumption that these atomic updates can occur in effectively constant time, irrespective of contention, which is a strong assumption.⁴

The reliance on "vectorization" also introduces nuances. On architectures like GPUs executing in a Single Instruction, Multiple Thread (SIMT) fashion¹⁹, threads are grouped into warps or waves. If threads within a single warp process input elements that map to different freq array indices, the memory accesses become *divergent*, potentially reducing memory bandwidth efficiency. If threads within a warp map to the *same* index, they contend for the atomic operation, causing *serialization* within the warp.⁴ Similarly, the scatter operation in Step 2 might involve divergent write patterns depending on how output ranges are assigned to threads. Achieving true $O(1)$ performance therefore requires not just massive parallelism, but highly *coherent* parallel execution with minimal divergence and serialization. This coherence is data-dependent (input value distribution) and hardware-dependent. Frameworks like JAX, NumPy, and CuPy achieve performance by mapping high-level operations effectively onto the underlying hardware's capabilities¹, but the random-access, potentially conflicting write pattern inherent in parallel histogramming poses a challenge to ideal vectorization efficiency.

III. Theoretical Complexity Assessment

This section evaluates the theoretical underpinnings of the proposed algorithm, focusing on the $O(1)$ time complexity claim within idealized parallel models and comparing it against established bounds and algorithms.

A. Evaluating the $O(1)$ Claim under Parallel Models

To theoretically analyze the potential for $O(1)$ time complexity, we can consider the Parallel Random Access Machine (PRAM) model. PRAM provides an abstract framework for parallel algorithm analysis, assuming a set of processors sharing a

global memory, executing synchronously.¹⁹ It explicitly ignores practical constraints like communication latency, memory bandwidth limits, and synchronization overheads, focusing purely on concurrency.¹⁹ PRAM models vary based on how they handle simultaneous memory accesses¹⁹:

- **EREW (Exclusive Read, Exclusive Write):** No concurrent access allowed.
- **CREW (Concurrent Read, Exclusive Write):** Concurrent reads are allowed, but writes must be exclusive.
- **CRCW (Concurrent Read, Concurrent Write):** Both concurrent reads and writes are allowed. CRCW models further differ in how write conflicts are resolved (e.g., COMMON, ARBITRARY, PRIORITY, SUM).¹⁹

Analysis of Step 1 (Mapping/Histogram): The mapping step involves n input elements potentially attempting to increment counters in the shared freq array. To achieve $O(1)$ time for this step, we would need:

1. At least n processors ($p \geq n$), one for each input element.
2. A PRAM model that allows concurrent writes to the same memory location (a bin in freq) and resolves them appropriately in constant time. Specifically, a **CRCW PRAM with atomic increment or summation** capability would be required. If multiple processors write to $\text{freq}[v]$ simultaneously, the model must support summing these contributions in $O(1)$ time.¹⁹ Under these strong assumptions (n processors, $O(1)$ atomic concurrent increment/sum), Step 1 could theoretically be performed in $O(1)$ time.

Analysis of Step 2 (Reconstruction/Scan/Scatter): This step involves a prefix sum (scan) on the freq array (size $k+1$) and a scatter operation.

1. **Prefix Sum (Scan):** Standard PRAM algorithms for prefix sum on an array of size m typically take $O(\log m)$ time using $m/\log m$ processors on EREW/CREW PRAM.¹⁷ In our case, $m = k+1$. While faster algorithms exist on stronger CRCW models, achieving $O(1)$ for prefix sum generally requires m processors and specific CRCW capabilities, and is not standard. Therefore, the scan operation is likely at least $O(\log k)$ theoretically.
2. **Scatter:** After computing the offsets via scan, the n elements need to be written to their final positions in the sorted_arr. If n processors are available, and each knows its destination index (calculated from the scan results and its original value), this scatter could potentially take $O(1)$ time on a CRCW PRAM, assuming either exclusive write locations or appropriate handling of concurrent writes (though ideally, the scan ensures exclusive destination slots).

Overall Theoretical Complexity: Combining the steps, even on an idealized CRCW

PRAM with n processors and $O(1)$ atomic increments, the algorithm's time complexity appears dominated by the prefix sum step, yielding $O(\log k)$ rather than $O(1)$. Achieving $O(1)$ would require an $O(1)$ prefix sum, which is non-standard or requires even stronger PRAM assumptions, potentially $O(k)$ processors.

Limitations of PRAM: It is crucial to recognize that PRAM is an unrealistic model.³² It ignores memory bandwidth, network latency, synchronization costs, and assumes an often unbounded number of processors. Therefore, an $O(1)$ or $O(\log k)$ result on PRAM does not directly translate to practical performance. Real-world performance is heavily influenced by these ignored factors.³²

B. Comparison with Counting Sort ($O(n+k)$)

Sequential Complexity: Standard sequential Counting Sort has a time complexity of $\Theta(n+k)$.¹⁸ This arises from iterating through the n input elements to build the frequency histogram ($O(n)$), performing a prefix sum on the counts array ($O(k)$), and then placing the elements into the output array ($O(n)$).

Parallel Complexity: Parallelizing Counting Sort typically involves distributing the n elements among p processors. A common approach includes²⁸:

1. Each processor counts frequencies for its n/p elements into a local histogram ($O(n/p)$).
2. Local histograms are merged into a global histogram (cost depends on p and k , e.g., $O(p*k)$ sequentially, or $O(k \log p)$ in parallel).
3. A parallel prefix sum is performed on the global histogram ($O(\log k)$ or $O(k/p)$).
4. Elements are scattered to the output array ($O(n/p)$). This leads to complexities like $O(n/p + p*k + \log k)$ or $O(n/p + k \log p + \log k)$, which are generally not $O(1)$ unless n and k are bounded relative to p in specific ways.

Key Difference from Proposal: The proposed algorithm attempts to bypass the explicit local histogram and merging phases of typical parallel Counting Sort. It relies instead on direct, concurrent updates to a single shared histogram (Step 1), aiming to simplify the parallel structure. The theoretical advantage hinges on whether the hardware cost of resolving these concurrent writes (via atomics) is less than the overhead of partitioning, local computation, and merging in traditional parallel Counting Sort. The $O(1)$ claim essentially posits that the concurrent write resolution is constant time.

C. Impact of Input Size ' n ' and Range ' k '

Dependence on ' n ': The theoretical $O(1)$ claim implies runtime independence from n .

As discussed, this requires at least n processors and $O(1)$ concurrent memory operations in the PRAM model. In practice, the number of physical processors p is fixed and usually much smaller than n . Therefore, the work proportional to n (reading inputs, updating histogram, scattering outputs) must be distributed, leading to terms like $O(n/p)$ in the runtime.²⁸ Furthermore, memory bandwidth limitations often make the runtime dependent on n , as n elements need to be read from and written to memory.²⁸ Therefore, practical runtime is unlikely to be truly independent of n , although for very large p and small n , the n/p term might become negligible compared to other overheads or k -dependent terms.

Dependence on 'k': The range k critically impacts both space and time complexity:

- **Space:** The algorithm requires $O(k)$ auxiliary space for the freq array.
- **Time:**
 - Allocating and potentially initializing the $O(k)$ freq array takes time, at least $O(k/p)$ in parallel.
 - The prefix sum operation on the freq array takes time dependent on k , such as $O(\log k)$ or $O(k/p)$ in parallel models.¹⁷
 - If k is large, these steps dependent on k cannot be completed in $O(1)$ time, regardless of n . The $O(1)$ claim implicitly assumes k is a fixed constant, or sufficiently small such that $\log k$ or k/p are negligible compared to other constant factors on the target hardware.

D. Theoretical Lower Bounds for Sorting

It is informative to place the proposed algorithm in the context of theoretical sorting lower bounds.

- **Comparison Sorts:** Any sorting algorithm relying solely on comparisons between elements requires $\Omega(n \log n)$ comparisons in the worst case to sort n elements.¹⁸ This fundamental limit highlights why non-comparison-based algorithms like Counting Sort, Radix Sort, and the proposed algorithm are necessary to achieve potentially better performance (e.g., $O(n+k)$) for restricted input types like integers.
- **Parallel Sorting Lower Bounds:** Parallel sorting also has lower bounds. In parallel comparison models, even with a polynomial number of processors, sorting requires $\Omega(\log n)$ time steps.³⁶ While the proposed algorithm is non-comparison based, fundamental constraints related to processing n items still apply. Any algorithm must perform $\Omega(n)$ total work (at least read all inputs).³⁸ With p processors, this implies a time lower bound of $\Omega(n/p)$. Achieving $O(1)$ time thus necessitates $p = \Omega(n)$, meaning the number of processors must scale linearly

with the input size, along with assumptions about constant-time communication and memory access.

E. Re-evaluating the $O(1)$ Claim

Considering the PRAM analysis, comparisons with parallel Counting Sort, the dependencies on n and k , and general parallel lower bounds, the claim of $O(1)$ time complexity requires significant qualification. It is highly unlikely to hold in a strict theoretical sense where n and k can grow arbitrarily. The prefix sum step introduces at least a logarithmic dependency on k , and practical limitations introduce dependencies on n/p and memory bandwidth.

The $O(1)$ claim is more plausibly interpreted as "**practical $O(1)$** " or "**amortized $O(1)$** " under specific, constrained conditions:

1. **Fixed, Small k :** The range k must be constant and small enough that operations dependent on it ($O(k)$ allocation/initialization, $O(\log k)$ or $O(k/p)$ scan) are negligible or fit within the hardware's constant-time execution capabilities (e.g., k fits entirely in fast shared memory, scan is extremely fast for small k).
2. **n within Parallelism Limits:** The input size n must be within the range where the available hardware parallelism (p) is sufficient to make the n/p terms effectively constant or dominated by fixed overheads (e.g., n is less than or comparable to the total number of GPU threads that can run concurrently).
3. **Efficient Hardware Primitives:** The underlying hardware must provide extremely fast atomic operations and high memory bandwidth to minimize the constant factors hidden by the $O(1)$ notation.

Under these conditions, the wall-clock time *might appear* relatively constant as n varies within a certain range, because the massively parallel computation is so fast that fixed overheads (kernel launch, synchronization, memory latency) or bandwidth limits become the dominant factors.²⁸ This is not true $O(1)$ complexity but rather performance saturation due to hardware limits or fixed costs.

Furthermore, the analysis suggests that the range k is potentially a more significant theoretical and practical barrier to achieving near- $O(1)$ performance than n . While massive parallelism can theoretically handle large n given $p \geq n$, operations inherently tied to the size of the range k – notably allocating and initializing $O(k)$ memory and performing the prefix sum over k elements – are fundamentally harder to make constant time as k increases, even with abundant parallelism. If k scales with n or becomes very large (e.g., millions or billions for 32/64-bit integers), the $O(k)$ space and the $O(\log k)$ or $O(k/p)$ time components for k -dependent steps will prevent $O(1)$

behavior and likely make the algorithm less efficient than traditional sorts like Radix Sort or Merge Sort. This reinforces the notion that the algorithm's potential advantage is confined to scenarios where k is known, fixed, and relatively small compared to n and hardware capabilities.

IV. Practical Implementation and Performance

This section explores the practical aspects of implementing the proposed algorithm using common high-performance computing libraries and frameworks, focusing on potential performance bottlenecks, comparisons with standard sorting algorithms, and the critical implications of space complexity.

A. Vectorization Libraries (NumPy, JAX, CuPy)

Modern scientific computing libraries provide high-level abstractions for vectorized operations, potentially simplifying the implementation.

- **NumPy (CPU):** A straightforward implementation can leverage NumPy's optimized C functions. Step 1 (Mapping/Histogram) maps directly to `numpy.bincount(arr, minlength=k+1)`. Step 2 (Reconstruction) can be achieved using `numpy.repeat(numpy.arange(k+1), counts)` where `counts` is the result of `bincount`.¹³
 - *Performance:* While `bincount` and `repeat` are highly optimized C implementations, they primarily run on the CPU. `bincount` involves random-access writes to the `counts` array, which can stress CPU caches. The overall execution is largely sequential unless NumPy's underlying BLAS/LAPACK libraries utilize multi-threading for specific operations, which is unlikely for `bincount` or `repeat`. Performance will likely scale roughly as $O(n+k)$, bound by memory bandwidth and cache efficiency on the CPU.¹³ It cannot achieve true parallel $O(1)$. However, for suitable data (small k), it can significantly outperform comparison-based sorts like Python's Timsort on the CPU.¹³
- **JAX (CPU/GPU/TPU):** JAX offers a NumPy-like API (`jax.numpy`) and powerful function transformations (`jit`, `vmap`, `pmap`).¹
 - *Implementation:* One could potentially use `jax.numpy.bincount`⁴⁵ or implement the histogram using `lax.scatter_add` or similar primitives. The reconstruction could use `jax.numpy.repeat` or equivalent JAX operations. The key advantage lies in `jax.jit`, which compiles the Python function using XLA (Accelerated Linear Algebra) into optimized machine code for the target backend (CPU, GPU, TPU).¹ XLA performs whole-program optimization, potentially fusing the mapping and reconstruction steps.³¹ `jax.vmap` allows automatic vectorization¹,

although mapping the histogram update (a scatter-add pattern) efficiently via vmap might require careful formulation.

- *Performance*: Significant speedups over NumPy are expected when running on GPUs or TPUs due to XLA compilation.¹ How XLA handles the scatter-add pattern (Step 1) on accelerators is critical – does it compile to efficient atomic operations on the GPU/TPU? This detail is crucial for performance but may require inspecting generated code or deeper JAX internals.⁴⁵ If efficient atomics are used, JAX could offer a high-level way to achieve performance close to native CUDA, but the O(1) claim remains subject to the previously discussed theoretical and practical limitations.
- **CuPy (GPU)**: CuPy provides a NumPy-compatible API that executes directly on NVIDIA GPUs using CUDA.²
 - *Implementation*: Similar to NumPy, one would look for CuPy equivalents of `bincount` and `repeat`. CuPy's documentation lists `cupy.bincount` under statistics routines.⁴⁸ CuPy often wraps underlying CUDA libraries like Thrust and CUB for its primitives.⁴⁶ For instance, `cupy.sort` uses Thrust.⁴⁶ It's plausible `cupy.bincount` might leverage optimized CUDA histogramming techniques internally.
 - *Performance*: Performance depends heavily on the quality of CuPy's underlying CUDA implementation for the histogram (`bincount`) and reconstruction (`repeat`) steps. If it wraps efficient CUB/Thrust primitives, performance could be excellent, approaching native CUDA speeds. However, documentation on the specific implementation (e.g., use of atomics, Thrust/CUB) is sparse.⁴⁸ Like JAX on GPU, it offers a higher-level interface but performance hinges on the efficiency of the compiled/wrapped CUDA code for the critical histogramming step.

B. GPU Implementation (CUDA, Thrust, CUB)

A direct CUDA implementation offers the most control for optimization but requires deeper expertise.

- **Step 1 (Histogram)**:
 - *Atomic Operations*: The primary challenge remains handling write collisions to the `freq` array. CUDA provides atomic functions (e.g., `atomicAdd`, `atomicInc`) that operate directly on global or shared memory.⁴ Using `atomicAdd(&freq[arr[i]], 1)` within a kernel launched with `n` threads (or fewer threads processing multiple elements) is the most direct parallel implementation.¹⁶ However, atomics incur performance costs:
 - *Serialization*: If multiple threads in a warp target the same `freq` index, their

atomic operations are serialized, reducing parallelism.⁴

- *Memory Coalescing*: Atomic operations often involve read-modify-write cycles to potentially scattered memory locations, hindering optimal memory coalescing.⁴
- *Latency*: Atomic operations themselves have latency, especially to global memory.⁴ Hardware improvements, like faster shared memory atomics introduced in Maxwell and later architectures, can mitigate this if k is small enough for freq to fit in shared memory.⁸ Performance is highly dependent on the data distribution (skew increases collisions) and the specific GPU architecture's atomic performance.⁷
- *Privatization/Optimization*: To reduce contention on global atomics, privatization is a common technique.⁴ Each thread block computes a private histogram for its assigned portion of the input data. This private histogram can reside in shared memory (if k is small enough, offering fast atomic updates) or global memory (requiring global atomics but only among threads within the block, reducing contention).⁴ After computing local histograms, a separate kernel or a final step within the first kernel performs a reduction (summation) across all private histograms to produce the final global freq array. This adds complexity and synchronization points but can significantly outperform direct global atomics when contention is high.⁸ Other optimizations might involve rearranging input data to reduce collision intensity.⁷
- *CUB/Thrust Primitives*: NVIDIA's CUB library provides `cub::DeviceHistogram`, a highly optimized primitive for exactly this task.³ It likely employs sophisticated internal strategies combining atomics, shared memory, and privatization based on input size, data types, and hardware capabilities.⁹ Using `cub::DeviceHistogram` is often the most performant approach. Thrust, while primarily known for algorithms like sort and scan, might also offer histogram capabilities or could be used to build one, potentially leveraging CUB internally.²⁰
- **Step 2 (Reconstruction):**
 - *Parallel Prefix Sum (Scan)*: The calculation of output offsets is ideally performed using optimized parallel scan primitives. Both Thrust (`thrust::inclusive_scan`, `thrust::exclusive_scan`) and CUB (`cub::DeviceScan`) provide robust and high-performance implementations.³ These libraries implement work-efficient algorithms ($O(k)$ work, theoretically $O(\log k)$ steps) that are heavily optimized for GPU architectures.¹⁷
 - *Parallel Scatter*: The final step involves scattering the values j into the `sorted_arr` based on the computed offsets and freq counts. This can be

implemented with a custom CUDA kernel. Threads can be assigned ranges of the output array or ranges of the input values j . Each thread calculates the destination indices and writes the corresponding value j repeatedly ($\text{freq}[j]$ times). Memory access patterns are crucial here; if threads write to consecutive locations corresponding to the same value j , writes can be coalesced. However, the overall pattern is a scatter, which can be less bandwidth-efficient than coalesced reads/writes. Thrust and CUB may offer primitives for permutation or scatter operations that could be adapted.

- **Synchronization Overheads:** A multi-step GPU implementation requires synchronization. At minimum, the histogramming kernel (Step 1) must complete before the reconstruction kernel (Step 2) can begin. If privatization is used in Step 1, synchronization is needed between the local histogramming phase and the reduction phase. Kernel launch overhead itself adds a constant time cost. Furthermore, internal operations within library primitives like Thrust or CUB might involve implicit synchronization or memory allocation/deallocation, which can contribute to overhead if not managed carefully (e.g., by providing pre-allocated temporary storage or using custom allocators).⁹
- **Memory Access Patterns:** The algorithm's memory access patterns are mixed. Step 1 involves reading the input array `arr` (potentially coalesced if accessed linearly by threads) but performs scattered, potentially conflicting writes to the `freq` array (bad for coalescing, relies on atomics/cache). Step 2 involves reading the `freq` array (likely sequentially accessed during/after the scan) and the `offsets` array, and performs scattered writes to the `sorted_arr`. Scatter operations are generally less efficient than coalesced accesses on GPUs. Overall, the algorithm is likely to be memory-bandwidth bound, especially Step 1 due to the atomic updates and Step 2 due to the scatter writes.²⁸

C. Performance Comparison with Standard Sorts

- **CPU Comparison:** On a CPU, the NumPy implementation (effectively $O(n+k)$) is expected to outperform standard $O(n \log n)$ comparison sorts (like Timsort or `std::sort`) only when k is significantly smaller than $n \log n$.¹³ For large k , the $O(k)$ term dominates, making it slower.
- **GPU Comparison:** On a GPU, the comparison is primarily against highly optimized library implementations:
 - *GPU Radix Sort (e.g., CUB/Thrust):* This is typically the fastest available sort for integer or float keys on GPUs.³⁰ GPU Radix Sort works by making multiple passes (e.g., d passes for d -digit keys), often using 4 or 8 bits per pass (radix $r = 16$ or 256).³⁰ Each pass involves a histogram, scan, and scatter operation, similar to the proposed algorithm but repeated d times. The proposed

algorithm can be viewed as a single-pass radix sort where the "radix" covers the entire range k .

- *Potential Advantage:* If k is small (e.g., fits within the radix size of a single pass, like $k \leq 255$ for an 8-bit pass), the proposed algorithm avoids the overhead of multiple passes and might be faster.
- *Potential Disadvantage:* For larger k , the single histogram step in the proposed algorithm might face much higher atomic contention and require a larger intermediate freq array compared to the smaller histograms used in each pass of a typical radix sort. Radix sort passes often have more structured memory access. State-of-the-art radix sorts like CUB's DeviceRadixSort or OneSweep are extremely optimized³⁰, potentially using techniques like hybrid approaches or specialized scan/histogram methods.³⁵ Benchmarks consistently show Thrust/CUB radix sort achieving very high throughput.³⁰ Therefore, the proposed algorithm must demonstrate significant benefits over these optimized multi-pass radix sorts, likely only achievable for very small k .
- *GPU Merge Sort / Comparison Sorts:* Algorithms like Merge Sort or Bitonic Sort, while parallelizable, are generally slower than Radix Sort on GPUs for keys suitable for radix sort.³⁰ QuickSort implementations on GPUs exist but often struggle with load balancing and irregular memory access compared to radix or merge sort.⁵⁶ The proposed algorithm, being non-comparison based and $O(n+k)$ -like, should significantly outperform GPU comparison sorts in its target domain (bounded integers).
- **Factors Influencing Performance:** Performance is highly sensitive to:
 - n and k : The relative sizes determine the dominance of n -dependent vs. k -dependent work and memory.
 - Data Distribution: Skewed data (many identical values) increases atomic contention in Step 1, potentially degrading performance significantly.⁷ Uniform data minimizes contention.
 - Memory Bandwidth: As a likely bandwidth-bound algorithm, performance depends heavily on the GPU's memory speed.²⁸
 - Atomic Operation Efficiency: The speed of atomic operations on the specific GPU hardware is critical.⁴
 - Hardware Generation: Newer GPUs may have faster atomics, larger caches, or higher bandwidth, impacting performance.⁵⁰

D. Space Complexity Implications ($O(k)$)

The $O(k)$ auxiliary space required for the freq array (and potentially the offsets array)

is a major constraint and differentiator:

- **Limitation:** Unlike in-place algorithms like HeapSort ($O(1)$ space) or QuickSort ($O(\log n)$ space for recursion stack)³⁷, this algorithm's memory usage scales linearly with the range k .
- **Acceptable Scenarios:** The $O(k)$ space is acceptable only when k is relatively small and fits comfortably within the available memory (e.g., GPU global memory, CPU RAM). Examples include sorting 8-bit values ($k=255$), small classification labels, or other scenarios with inherently limited ranges.
- **Problematic Scenarios:** If k is large (e.g., for 32-bit integers, $k \approx 4$ billion; for 16-bit integers, $k = 65535$), the $O(k)$ space becomes prohibitive, potentially exceeding GPU memory or even system RAM. In such cases, algorithms with lower space complexity or multi-pass algorithms like Radix Sort (which only needs $O(n+r)$ space per pass, where r is the radix size) are necessary.
- **Cache Effects:** The performance is also tied to how k relates to cache sizes. If the freq array (size $k+1$) fits within the L2 cache or, even better, the GPU's shared memory (typically tens or hundreds of KB per SM)⁴, access times during histogramming and scanning will be much faster. If k forces the freq array into global memory, performance will be limited by DRAM bandwidth and latency.

E. Algorithm Comparison Summary

The following table summarizes the key characteristics of the proposed algorithm compared to standard sorting algorithms, assuming the practical $O(n/p + f(k))$ or $O(n+k)$ behavior for the proposal.

Feature	Proposed Algorithm	Counting Sort	Radix Sort (LSD)	QuickSort (Avg)	MergeSort	HeapSort
Time Complexity						
- Best	$O(1)$ / $O(n/p)$ (Ideal Par, Fixed k)	$O(n+k)$	$O(d*(n+r))$	$O(n \log n)$	$O(n \log n)$	$O(n \log n)$
- Average	$O(n/p + f(k))$ / $O(n+k)$	$O(n+k)$	$O(d*(n+r))$	$O(n \log n)$	$O(n \log n)$	$O(n \log n)$

	(Practical)					
- Worst	$O(n/p + f(k)) / O(n+k)$ (Practical)	$O(n+k)$	$O(d*(n+r))$	$O(n^2)$	$O(n \log n)$	$O(n \log n)$
Space Complexity	$O(k)$	$O(n+k) / O(k)$	$O(n+r)$	$O(\log n)$	$O(n)$	$O(1)$
Stable	No (by default)	Yes	Yes	No	Yes	No
In-Place	No	No	No	Yes (typical)	No	Yes
Input Constraint	Integers in $[0, k]$	Integers	Integers/S trings	General	General	General
Parallelism	High (GPU/TPU focus)	Moderate	High (GPU good)	Moderate	High	Low
Notes	Relies heavily on atomics/vector ops	Efficient small k	Efficient fixed digits d , radix r	Cache-friendly	Stable	Memory eff.

(Note: d =number of digits/passes, r =radix size. Practical complexities depend heavily on implementation and hardware.)³⁷

F. Hardware Alignment and Stability Considerations

The design philosophy behind the proposed algorithm appears deeply intertwined with the capabilities of modern parallel hardware. Its potential for achieving near- $O(1)$ performance stems less from a fundamentally new algorithmic structure and more from an aggressive exploitation of hardware features. Success relies critically on:

1. Massively parallel execution units (thousands of GPU cores).
2. High memory bandwidth to feed these units and handle the data movement.

3. Extremely efficient hardware support for atomic operations (especially atomic add/increment) to manage the histogram update contention, or effective compiler/library strategies (like privatization) that map well to the hardware.⁴

This suggests a form of hardware-algorithm co-design, where the algorithm is viable precisely because specific hardware primitives (like fast atomics, efficient scan implementations in CUB/Thrust) have become available and highly optimized.³ The algorithm's performance is therefore likely to evolve with hardware generations – improvements in atomic throughput or memory subsystems could directly enhance its viability.

A significant trade-off made for this potential speed is **stability**. Stable sorting algorithms preserve the relative order of elements with equal keys.³⁷ The proposed algorithm, in its direct form, is inherently unstable. The mapping step (Step 1) collapses all instances of a value v into a single count $\text{freq}[v]$, losing their original ordering. The reconstruction step (Step 2) then writes $\text{freq}[v]$ copies of v into the output array, typically in a contiguous block, without reference to their original sequence. Achieving stability would require substantial modifications. For instance, one might need to process input elements sequentially based on their original order when placing them into the output bins, or use a more complex data structure during the mapping phase (e.g., storing indices or pointers in lists associated with each bin). Such modifications would introduce sequential dependencies or significant overhead, almost certainly invalidating the $O(1)$ time complexity goal and likely making the algorithm slower than standard stable sorts like Merge Sort or stable Radix Sort.³⁷ This lack of stability must be considered when evaluating its suitability for applications where preserving the original order of equal elements is necessary.

V. Domain Applicability Analysis

This section assesses the feasibility and potential benefits of employing the proposed sorting algorithm in two specific domains highlighted by the user: Machine Learning (ML) pipelines and embedded systems.

A. Machine Learning Pipelines

Context: Real-time ML inference, particularly for applications like autonomous driving, robotics, or interactive services, often operates under stringent latency constraints.⁶² Decisions or predictions may need to be generated within milliseconds or even microseconds to be useful.⁶³ Inference pipelines can involve multiple stages: data preprocessing, model execution (often on accelerators like GPUs/TPUs), and

post-processing.⁶⁶ Optimizing each stage for low latency is critical.⁶⁸

Potential Use Cases for Bounded Integer Sorting:

- **Label/Index Sorting:** Sorting predicted class labels or indices, which are often small integers within a known range (e.g., 0 to `number_of_classes - 1`).
- **Quantized Value Sorting:** If activations or weights are quantized to low-bit integers (e.g., INT8) during inference, sorting these values might be required in specific custom operations or analysis steps.
- **Top-K Post-processing:** After identifying the top-k results (e.g., highest probability classes), sorting the indices of these top-k results might be needed.
- **Histogram-based Features:** Some models might explicitly use histogram features computed from intermediate data, where the histogramming step (Step 1 of the proposed algorithm) could be directly relevant.

Suitability Analysis:

- **Latency:** If the proposed algorithm can achieve practical $O(1)$ or significantly lower *absolute latency* than existing GPU sorts (like CUB/Thrust Radix Sort) for the specific small k values encountered in ML tasks (e.g., $k=1000$ for ImageNet labels), it could be beneficial. Saving even microseconds might matter in ultra-low-latency scenarios.⁶⁴
- **Space Complexity ($O(k)$):** The $O(k)$ space is less likely to be a major issue if k corresponds to the number of classes (typically hundreds or thousands), as this amount of auxiliary memory is usually negligible on GPUs used for inference. However, if sorting quantized activations across a wider range (e.g., $k=255$ for INT8, or larger for INT16), the memory footprint needs consideration.
- **Instability:** The lack of stability is often acceptable when sorting class labels or probabilities, where the original order of tied predictions might not be critical.
- **Integration:** Integrating a custom sort into complex ML inference frameworks (like TensorFlow, PyTorch, TensorRT) requires effort. Using JAX might offer a smoother path if its compilation to XLA handles the pattern efficiently.¹

Comparison: The primary competitor within an ML inference pipeline on a GPU would be the highly optimized Radix Sort available in libraries like Thrust/CUB.³⁰ These are already very fast for the integer types often involved. The proposed algorithm would need to demonstrate a compelling *absolute time* advantage for relevant n and small k to justify its use over these standard, robust library functions.

A crucial consideration is the **overall latency budget** of the inference pipeline.⁶⁹ The entire process, from input to output, must meet the deadline.⁶³ Even if the proposed

sort achieves near- $O(1)$ performance, its actual wall-clock time (including kernel launch, synchronization, data transfers, and potential atomic delays) must be evaluated. If sorting constitutes only a tiny fraction of the total inference time (with the neural network execution dominating), then even a significant relative speedup in sorting might yield only a negligible improvement in end-to-end latency.⁶⁷ The algorithm is most valuable if sorting bounded integers is identified as a genuine bottleneck within the target pipeline.

B. Embedded Systems

Context: Embedded systems are characterized by strict resource constraints, including limited RAM, lower processing power (often MCUs without advanced SIMD or parallelism), and tight power budgets.⁷⁰ Reliability and predictable real-time performance are paramount.⁷⁰ Flash memory with asymmetric read/write costs is common.⁷²

Algorithm Suitability Analysis:

- **Space Complexity ($O(k)$):** This is likely the most significant barrier. Embedded systems often have RAM measured in kilobytes.⁷² An $O(k)$ space requirement is generally unacceptable unless k is extremely small (e.g., tens or maybe a few hundreds).⁶⁰ Algorithms with $O(1)$ space complexity like HeapSort, Insertion Sort, or Bubble Sort (for very small n) are strongly preferred.³⁷ Specialized external sorting algorithms are needed if data exceeds RAM.⁷²
- **Time Complexity / Hardware Dependence:** The proposed algorithm's speed relies heavily on massive parallelism (thousands of threads) and efficient atomic operations found in GPUs/TPUs.⁴ Typical embedded processors lack this level of parallelism and may lack hardware atomics or efficient vector units.⁷⁰ Executing the algorithm on a standard MCU would likely result in poor performance, far from $O(1)$, potentially worse than simple $O(n^2)$ sorts for small n due to overhead.
- **Determinism:** The performance characteristics, especially regarding atomic contention, can be data-dependent and harder to predict than simpler algorithms like HeapSort (consistent $O(n \log n)$) or Insertion Sort (predictable for nearly sorted data).⁶⁰ Real-time embedded systems often require predictable worst-case execution times (WCET).⁷⁰
- **Code Size/Simplicity:** While potentially simple conceptually, an optimized parallel implementation involving atomics or privatization can increase code complexity compared to basic sorts like Bubble or Insertion Sort.⁷³

Conclusion for Embedded: The proposed algorithm appears generally unsuitable for traditional resource-constrained embedded systems due to its $O(k)$ space complexity

and reliance on advanced parallel hardware features not typically present. Simpler, low-memory algorithms like HeapSort, Insertion Sort, or even Bubble Sort for very small datasets are often more practical and appropriate.⁶⁰ The algorithm might only be considered for high-performance embedded systems equipped with GPU accelerators or powerful SoCs, and even then, only if k is small enough to meet memory constraints.

C. Domain-Specific Considerations

The assessment highlights a fundamental tension: the algorithm is designed for environments with abundant parallelism and memory bandwidth (GPUs/TPUs), which are often the antithesis of traditional embedded systems focused on resource minimization.⁷⁰ While potentially applicable to ML inference on accelerators, its value proposition depends heavily on the specific task, the size of k , the overall pipeline latency budget, and demonstrating a clear advantage over existing highly optimized library sorts. The very characteristics that enable its potential speed (massive parallelism, reliance on fast atomics, $O(k)$ temporary storage) are the same characteristics that make it a poor fit for typical embedded use cases where memory is scarce and processing cores are fewer and simpler.

VI. Novelty and Contribution

This section evaluates the novelty of the proposed sorting algorithm by comparing it against existing techniques and discusses its potential impact, including suitability for academic publication and library integration.

A. Comparison to Existing Techniques

- **Counting Sort:** The core principle of using element values as indices into a frequency array (histogram) and then using those counts to determine output positions is the defining characteristic of Counting Sort.⁵ The proposed algorithm directly employs this principle. Standard Counting Sort is $O(n+k)$ sequentially¹⁸, and various parallel versions exist with complexities like $O(n/p + k \log p)$ or similar.²⁸ The proposed method differs not in the fundamental counting idea but in its *implementation strategy* aimed at achieving practical $O(1)$ time through massive, direct parallelization of the histogramming and reconstruction steps, heavily relying on hardware atomics or equivalent mechanisms.⁴ It attempts to bypass the typical structure of parallel counting sort (local counts + merge).²⁸
- **Radix Sort:** For integer sorting, Radix Sort is a highly relevant comparison, especially on GPUs where it often yields the best performance.³⁰ LSD Radix Sort processes keys digit by digit (or chunk by chunk, e.g., 4 or 8 bits at a time) over

multiple passes.³⁰ Each pass typically involves a stable distribution counting step (similar to Counting Sort or the proposed algorithm's histogramming/scan/scatter, but only on the current digit/chunk). The proposed algorithm can be viewed as a **single-pass Radix Sort** where the "radix" r is equal to the entire range $k+1$.

- *Comparison:* Standard GPU Radix Sort uses multiple passes with a smaller radix r (e.g., 16 or 256), requiring $O(n+r)$ auxiliary space per pass and $O(d*(n+r))$ time, where d is the number of passes (proportional to $\log_r(k)$).³⁰ The proposed algorithm performs only one pass but requires $O(k)$ space and faces potentially higher contention/complexity in its single, larger histogram step if k is much larger than r . A specialized hybrid radix sort has been proposed that uses 8-bit passes and sacrifices stability for performance, sharing some philosophical ground with optimizing for speed via specific implementation choices.³⁵
- **GPU Primitives (Thrust/CUB/OneSweep):** Modern GPU computing relies heavily on optimized libraries like Thrust and CUB, which provide state-of-the-art implementations of fundamental building blocks.³ These libraries contain highly optimized primitives for histogramming (`cub::DeviceHistogram`⁹), prefix sum/scan (`thrust::inclusive_scan`, `cub::DeviceScan`¹⁷), and sorting (including `cub::DeviceRadixSort` and the even more advanced OneSweep sort).⁵⁰ Any novel sorting algorithm, especially one built on similar principles (histogram, scan, scatter), must demonstrate significant advantages over using or combining these existing, heavily tuned primitives. OneSweep, for instance, uses sophisticated techniques like "chained-scan-with-decoupled-lookback" for its internal prefix sums to maximize performance on modern GPU architectures.⁵⁵
- **Vectorized Implementations (NumPy/JAX/CuPy):** Using high-level libraries like NumPy, JAX, or CuPy¹ to implement the algorithm represents an *implementation choice* leveraging existing frameworks, not typically a fundamental algorithmic novelty in itself. The novelty would arise only if the specific way vectorization is applied (e.g., a unique pattern within JAX's `vmap` or `lax` primitives) leads to unforeseen performance benefits or enables the practical $O(1)$ behavior more effectively than standard library functions like `bincount` or `scatter_add`.

B. Assessment of Novelty

Based on the comparisons above:

- The fundamental algorithmic concept – using counts derived from a histogram of bounded integer keys to determine sorted positions – is **not novel**. It is the basis of Counting Sort¹⁸ and is used within each pass of Radix Sort.³⁰
- The potential novelty lies primarily in the **specific implementation strategy and**

the ambitious performance target (practical $O(1)$). The claim hinges on aggressively exploiting massive hardware parallelism and the efficiency of concurrent write mechanisms (like atomics) to make the wall-clock time effectively constant for fixed k and n within hardware limits.

- The contribution appears to be more in the realm of **hardware-aware algorithm engineering** or **optimization** rather than the invention of a fundamentally new sorting paradigm. It resembles an attempt to push the implementation of parallel Counting Sort principles to an extreme on specific hardware architectures. Similar efforts focusing on optimizing GPU sorting through implementation choices (e.g., radix size, avoiding stability, specialized scans) exist in the literature.³⁵ The exploration of variations like "occurrence sort" (a variant of counting sort removing duplicates) also shows that refining counting-based sorts for GPUs is an active area.⁵

Therefore, the novelty is conditional. If rigorous theoretical analysis under a suitable parallel model (beyond idealized PRAM) and compelling empirical benchmarks demonstrate practical $O(1)$ behavior and superior performance compared to state-of-the-art CUB/Thrust radix sort (especially OneSweep⁵⁵) within a well-defined niche (small fixed k , high parallelism), then the specific implementation technique could be considered novel and valuable.

C. Potential for Academic Publication / Library Integration

- **Academic Publication:** To be publishable in a reputable venue, the work would require:
 1. **Rigorous Theoretical Analysis:** A clear definition of the computational model (incorporating aspects like processor count p , memory bandwidth, atomic operation costs) under which the "practical $O(1)$ " claim is analyzed and justified. Idealized PRAM analysis alone is insufficient.
 2. **Strong Empirical Evidence:** Comprehensive benchmarks comparing against state-of-the-art GPU sorting libraries (CUB Radix Sort, OneSweep) on multiple modern GPU architectures. Tests should vary n , k , and data distributions (uniform vs. skewed) to clearly delineate the algorithm's performance envelope and validate the $O(1)$ claim within specific regimes. Comparison with related work like specialized histogram sorts⁵ would strengthen the paper.
 3. **Clear Contribution:** Articulation of the novelty, focusing on the implementation techniques that enable the claimed performance and the conditions under which it holds. The limited applicability (bounded integers, small k) might restrict the perceived impact, but specialized

high-performance algorithms are often publishable if the results are significant within their niche.

- **Library Integration:** Integration into standard libraries (like NumPy, JAX, CuPy, or even CUDA's Thrust/CUB) is plausible *if* the algorithm offers a substantial, demonstrable performance advantage over existing methods for its specific use case (small, fixed k). Libraries often include specialized routines that outperform general algorithms under specific conditions (e.g., Thrust uses Radix Sort for primitive types where applicable⁵⁰, CUB has specialized histogram functions⁹). Key requirements would be a robust, well-tested implementation and clear evidence of significant speedup in its target niche to justify the added complexity to the library. The $O(k)$ space cost must also be considered acceptable for the target scenarios.

D. Synthesizing Novelty and Impact

The proposed algorithm occupies a space where the distinction between a novel algorithm and an extreme optimization of an existing one (Counting Sort) becomes blurred. Its potential success and recognition likely depend less on the core idea (counting) and more on the effectiveness of the *implementation* in harnessing hardware capabilities (parallelism, atomics, memory hierarchy) to achieve the near- $O(1)$ practical performance target. Much like advances in GPU radix sort often come from improved implementation strategies (e.g., radix choice, scan techniques, handling skew) rather than changing the fundamental radix principle³⁵, this proposal's contribution seems focused on the *how* (achieving near- $O(1)$ practically via direct mapping and hardware exploitation) rather than the *what* (the counting principle itself).

Regarding its potential to induce a paradigm shift, this seems unlikely. While potentially very fast within its narrow domain (small, fixed k), the algorithm's strong constraints prevent it from being a general-purpose sorting solution. Comparison sorts remain essential for general data types.³⁷ Radix Sort offers better scalability for larger integer ranges.³⁰ The proposed algorithm, therefore, appears destined to be a specialized tool, valuable perhaps in specific performance-critical niches like ML inference pre/post-processing⁶⁴ or other areas dealing with frequent sorting of small-range integers on highly parallel hardware. It represents a potential addition to the algorithmic toolkit for specific scenarios, rather than a fundamental reshaping of how sorting is broadly approached.

VII. Synthesis, Conclusion, and Recommendations

This report has conducted a detailed analysis of the proposed vectorized sorting

algorithm aiming for practical $O(1)$ time complexity for integers within a known, bounded range $[0, k]$. The analysis covered its internal mechanics, theoretical complexity, practical implementation considerations, domain applicability, and novelty.

A. Summary of Findings

- **Algorithm Mechanism:** The algorithm operates in two phases: parallel histogram construction (mapping input values to frequencies in an $O(k)$ array) and parallel reconstruction (using prefix sums/scans on frequencies to scatter values to their sorted positions).
- **Theoretical Complexity:** The $O(1)$ claim is theoretically plausible only under strong and unrealistic parallel models (e.g., CRCW PRAM with n processors and $O(1)$ atomic summation/scan). In more realistic models and practice, the complexity likely involves terms dependent on input size n , processor count p (e.g., n/p), and the range k (e.g., $\log k$ or k/p for scan, $O(k)$ space initialization). "Practical $O(1)$ " might be observed for fixed, small k and n within hardware parallelism limits, where performance is dominated by fixed overheads or bandwidth saturation, not true independence from n and k .
- **Practical Implementation:** Implementation is feasible using vectorized libraries (NumPy, JAX, CuPy) or directly in CUDA (using Thrust/CUB primitives or custom kernels). Performance critically depends on the efficiency of handling concurrent writes in the histogram step (atomic operations or privatization) and memory bandwidth. Atomic contention, especially with skewed data, is a major performance risk.⁴ Highly optimized library primitives (e.g., `cub::DeviceHistogram`, `cub::DeviceScan`, `cub::DeviceRadixSort`) provide strong competition.⁹
- **Space Complexity:** The $O(k)$ auxiliary space requirement is a significant limitation, restricting applicability to scenarios where k is small relative to available memory.
- **Performance Comparison:** The algorithm should outperform $O(n \log n)$ comparison sorts in its niche. Its primary competitors are parallel Counting Sort and, especially on GPUs, highly optimized Radix Sort implementations (e.g., CUB/Thrust). It is likely faster than Radix Sort only if k is small enough (potentially fitting within a single radix pass size) and the overhead of its single, potentially high-contention histogram step is less than the overhead of multiple passes in Radix Sort.
- **Domain Applicability:** It holds potential for specific ML inference tasks on GPUs/TPUs (e.g., sorting small-range labels or indices) where latency is critical and k is small. However, its value depends on the absolute time saved within the context of the entire pipeline. It is generally unsuitable for traditional resource-constrained embedded systems due to $O(k)$ space and reliance on

advanced parallel hardware.

- **Novelty:** The fundamental counting principle is not novel. The potential novelty lies in the specific implementation strategy aiming for practical $O(1)$ via extreme hardware exploitation and its rigorous empirical validation against state-of-the-art techniques. It is best characterized as a potentially highly optimized, hardware-aware implementation of parallel counting sort principles.

B. Addressing User's Core Questions

- **Promise:** The idea is **moderately promising within a well-defined niche:** sorting integers from a known, small, fixed range $[0, k]$ on massively parallel hardware (GPUs/TPUs). The $O(1)$ complexity claim requires significant qualification; it should be understood as "practical $O(1)$ " achievable only under specific constraints, not a general theoretical guarantee.
- **Worth Investment:**
 - **Prototyping (GPU/TPU): Yes**, this is the most crucial next step. Empirical validation is essential to determine if the practical $O(1)$ behavior manifests and if it offers real-world speedups over optimized library functions like CUB/Thrust Radix Sort for small k . JAX, CuPy, or direct CUDA/CUB implementations are suitable.¹
 - **Generalization (NumPy, JAX, CuPy): Yes**, exploring implementations across these frameworks is valuable for understanding portability and abstraction costs, particularly investigating JAX/XLA's handling of the core patterns.¹
 - **Academic Publication: Possible but challenging.** Requires rigorous theoretical grounding (beyond PRAM), strong, comparative empirical results demonstrating superiority over state-of-the-art (e.g., CUB OneSweep⁵⁵) in the target niche, and careful articulation of the "practical $O(1)$ " conditions.
 - **Library Integration: Contingent on strong positive results from prototyping.** If it demonstrably provides significant, reliable speedups for sorting bounded integers (small k) compared to existing library functions, it could be considered for integration as a specialized optimization.
- **Paradigm Shift: Unlikely.** The algorithm's specialization due to the $O(k)$ space constraint and reliance on bounded inputs prevents it from being a general-purpose replacement for existing sorting paradigms. It's more likely to be a niche optimization tool.

C. Recommendations

1. **Refine and Qualify the $O(1)$ Claim:** Clearly state that the goal is "practical $O(1)$ " under specific conditions (fixed small k , n within hardware limits, efficient atomics/memory). Avoid implying theoretical $O(1)$ independence from n and k .

Define the model and assumptions under which this practical behavior is expected.

2. **Prioritize GPU Prototyping and Benchmarking:** Focus initial efforts on creating robust GPU prototypes (CUDA/CUB recommended for baseline, potentially JAX/CuPy). Design rigorous benchmarks comparing against CUB/Thrust DeviceRadixSort (and ideally OneSweep if possible⁵⁵). Vary n (across orders of magnitude), k (small powers of 2, e.g., 2^8 , 2^{10} , 2^{12} , 2^{16}), and data distribution (uniform, skewed, few unique values). Measure kernel execution times, end-to-end times (including memory transfers if relevant), and potentially memory bandwidth utilization and atomic contention rates using profiling tools.
3. **Deep Dive into Histogramming Performance:** The parallel histogram (Step 1) is the most critical and complex part. Investigate the performance trade-offs between direct atomic updates (using `atomicAdd` in global or shared memory) and privatization techniques (per-block histograms followed by reduction).⁴ Measure performance under varying levels of data skew (collision rates). Compare custom implementations against `cub::DeviceHistogram`.⁹
4. **Acknowledge and Address Stability:** Explicitly state that the proposed algorithm is unstable. If stability is required for target applications, assess the feasibility and performance cost of modifying the algorithm to be stable, recognizing this likely compromises the performance goals.
5. **Investigate JAX/XLA Compilation:** If pursuing JAX, analyze the XLA HLO or PTX/SASS code generated for the histogramming (scatter-add) and reconstruction steps on GPU/TPU targets. Determine if XLA generates efficient atomic operations or uses alternative strategies. Benchmark the JAX implementation against the direct CUDA/CUB version.
6. **Focus on the Niche:** Target applications where k is genuinely small and fixed (e.g., $k < 65536$, or ideally $k < 4096$ or $k < 256$) and where sorting these specific types of arrays is a frequent operation or a performance bottleneck.

D. Concluding Remarks

The proposed algorithm represents an intriguing attempt to leverage modern hardware parallelism for potentially ultra-fast sorting of bounded integers. Its core idea, rooted in counting sort principles, is sound, but its ambition for practical $O(1)$ performance pushes the boundaries of hardware capabilities, particularly concerning concurrent memory updates via atomic operations. The $O(k)$ space complexity remains its most significant limitation, confining its applicability to scenarios with small, known ranges.

Ultimately, the algorithm's value proposition hinges on empirical validation. If rigorous

benchmarking demonstrates substantial and reliable performance gains over state-of-the-art, highly optimized GPU radix sorts within its specific niche (small k), it could represent a valuable specialized sorting technique for performance-critical applications like aspects of ML inference. However, without such evidence, it risks being an interesting theoretical exercise whose practical benefits are outweighed by existing, more general, and highly optimized library solutions. Further investment should prioritize generating this empirical evidence through careful prototyping and benchmarking on target hardware.

Works cited

1. Why You Should (or Shouldn't) be Using Google's JAX in 2023 - AssemblyAI, accessed April 15, 2025, <https://www.assemblyai.com/blog/why-you-should-or-shouldnt-be-using-jax-in-2023>
2. cupy.sort — CuPy 14.0.0a1 documentation, accessed April 15, 2025, <https://docs.cupy.dev/en/latest/reference/generated/cupy.sort.html>
3. GPU libraries — Get More for Less CUDA unbound (CUB) CUB, accessed April 15, 2025, <https://www.fz-juelich.de/en/ias/jsc/news/events/seminars/msa-seminar/2020-02-18-gpu-libraries/@@download/file>
4. GPU Pattern: Parallel Histogram, accessed April 15, 2025, https://ajdillhoff.github.io/notes/gpu_pattern_parallel_histogram/
5. Counting and Occurrence Sort for GPUs using an Embedded Language - Joel Svensson, accessed April 15, 2025, <https://svenssonjoel.github.io/writing/csort.pdf>
6. Compiling Generalized Histograms for GPU, accessed April 15, 2025, <https://hjemmesider.diku.dk/~zgh600/Publications/gen-histo-sc20.pdf>
7. Efficient Weighted Histogramming on GPUs with CUDA - NICS, accessed April 15, 2025, https://nics.ee.tsinghua.edu.cn/people/xumo/Tech_report_histogram.pdf
8. GPU Pro Tip: Fast Histograms Using Shared Atomics on Maxwell | NVIDIA Technical Blog, accessed April 15, 2025, <https://developer.nvidia.com/blog/gpu-pro-tip-fast-histograms-using-shared-atomics-maxwell/>
9. cub::DeviceHistogram Struct Reference - OpenFPM_pdata, accessed April 15, 2025, https://ppmcore.mpi-cbg.de/doxygen/openfpm/structcub_1_1DeviceHistogram.html
10. cub::DeviceHistogram — cub 3.1 documentation, accessed April 15, 2025, https://nvidia.github.io/cccl/cub/api/structcub_1_1DeviceHistogram.html
11. cub/cub/device/device_histogram.cuh at master · dmlc/cub - GitHub, accessed April 15, 2025, https://github.com/dmlc/cub/blob/master/cub/device/device_histogram.cuh
12. CUB Histogram - CUDA Programming and Performance - NVIDIA Developer Forums, accessed April 15, 2025,

- <https://forums.developer.nvidia.com/t/cub-histogram/64754>
13. How can I vectorize this python count sort so it is absolutely as fast as it can be?, accessed April 15, 2025,
<https://stackoverflow.com/questions/18501867/how-can-i-vectorize-this-python-count-sort-so-it-is-absolutely-as-fast-as-it-can>
 14. Massive Atomics for Massive Parallelism on GPUs - Rutgers University, accessed April 15, 2025, <https://people.cs.rutgers.edu/zz124/assets/pdf/ismm14.pdf>
 15. AMS 148 Chapter 6: Histogram, Sort, and Sparse Matrices, accessed April 15, 2025,
<https://ams148-spring18-01.courses.soe.ucsc.edu/system/files/attachments/note6.pdf>
 16. Histogram implementation - CUDA - NVIDIA Developer Forums, accessed April 15, 2025, <https://forums.developer.nvidia.com/t/histogram-implementation/35738>
 17. Chapter 39. Parallel Prefix Sum (Scan) with CUDA - NVIDIA Developer, accessed April 15, 2025,
<https://developer.nvidia.com/gpugems/gpugems3/part-vi-gpu-computing/chapter-39-parallel-prefix-sum-scan-cuda>
 18. ICS 311 #10: Theoretical Limits of Sorting, $O(n)$ Sorts, accessed April 15, 2025,
<https://www2.hawaii.edu/~nodari/teaching/s18/Notes/Topic-10.html>
 19. PRAM Algorithms, accessed April 15, 2025,
<http://cds.iisc.ac.in/wp-content/uploads/PRAM.pdf>
 20. EuroCC@Turkey Parallel Computing on GPUs with CUDA - Indico, accessed April 15, 2025,
https://indico.truba.gov.tr/event/84/attachments/173/379/Day%202-%20Session_3.pdf
 21. cub::DeviceScan — cub 3.1 documentation, accessed April 15, 2025,
https://nvidia.github.io/cccl/cub/api/structcub_1_1DeviceScan.html
 22. thrust::inclusive_scan — thrust 3.1 documentation, accessed April 15, 2025,
https://nvidia.github.io/cccl/thrust/api/function_group__prefixsums_1ga3ee07457799ed36979f02d92ef0b25f7.html
 23. thrust::exclusive_scan — thrust 3.1 documentation, accessed April 15, 2025,
https://nvidia.github.io/cccl/thrust/api/function_group__prefixsums_1ga8dbe92b545e14800f567c69624238d17.html
 24. ee147-s18:thrust - Daniel Wong, accessed April 15, 2025,
<http://www.danielwong.org/classes/ee147-s18/thrust>
 25. Cumulative summation in CUDA - parallel processing - Stack Overflow, accessed April 15, 2025,
<https://stackoverflow.com/questions/25251335/cumulative-summation-in-cuda>
 26. How do I know if thrust::inclusive_scan is being run in parallel on the GPU?, accessed April 15, 2025,
<https://forums.developer.nvidia.com/t/how-do-i-know-if-thrust-inclusive-scan-is-being-run-in-parallel-on-the-gpu/33491>
 27. Thrust :: CUDA Toolkit Documentation, accessed April 15, 2025,
<https://www.clear.rice.edu/comp422/resources/cuda/html/thrust/index.html>
 28. Can you do a parallel counting sort in $O(n/p)$ time? - Stack Overflow, accessed

April 15, 2025,

<https://stackoverflow.com/questions/39903181/can-you-do-a-parallel-counting-sort-in-on-p-time>

29. Chapter 6: Parallel Algorithms – The PRAM Model, accessed April 15, 2025, <https://www7.in.tum.de/~kretinsk/teaching/fundamental%20algorithms/6.pdf>
30. Designing efficient sorting algorithms for manycore GPUs - UMD Department of Computer Science, accessed April 15, 2025, <https://www.cs.umd.edu/class/spring2021/cmsc714/readings/Satish-sorting.pdf>
31. JAX – NumPy on the CPU, GPU, and TPU | Hacker News, accessed April 15, 2025, <https://news.ycombinator.com/item?id=37698740>
32. COMP 633: Parallel Computing PRAM Algorithms, accessed April 15, 2025, <https://www.cs.unc.edu/~prins/Courses/633/Readings/pram.pdf>
33. PARALLEL MERGE SORT* $O(\log n \log \log n / \log \log \log n)$ on n processors. (The basic algorithm was Preparata's; - CMU School of Computer Science, accessed April 15, 2025, <https://www.cs.cmu.edu/~guyb/paralg/papers/Cole88.pdf>
34. What is the time complexity of this algorithm for parallel merge sort - Stack Overflow, accessed April 15, 2025, <https://stackoverflow.com/questions/78601269/what-is-the-time-complexity-of-this-algorithm-for-parallel-merge-sort>
35. A Memory Bandwidth-Efficient Hybrid Radix Sort on GPUs - arXiv, accessed April 15, 2025, <https://arxiv.org/pdf/1611.01137>
36. TIGHT COMPLEXITY BOUNDS FOR PARALLEL COMPARISON SORTING Noga Alon Department of Mathematics Tel Aviv University and Bell Communi - Princeton Math, accessed April 15, 2025, <https://web.math.princeton.edu/~nalon/PDFS/Publications2/Tight%20complexity%20bounds%20for%20parallel%20comparison%20sorting.pdf>
37. Sorting algorithm - Wikipedia, accessed April 15, 2025, https://en.wikipedia.org/wiki/Sorting_algorithm
38. For a given problem, is there a theoretical lower bound on how fast any algorithm can solve it? - Reddit, accessed April 15, 2025, https://www.reddit.com/r/askscience/comments/613np0/for_a_given_problem_is_there_a_theoretical_lower/
39. Upper bound and lower bound of sorting algorithm - Stack Overflow, accessed April 15, 2025, <https://stackoverflow.com/questions/14671201/upper-bound-and-lower-bound-of-sorting-algorithm>
40. understanding of lower bound for comparison-based sorting algorithm - Stack Overflow, accessed April 15, 2025, <https://stackoverflow.com/questions/15625150/understanding-of-lower-bound-for-comparison-based-sorting-algorithm>
41. Advanced techniques for determining complexity lower bounds, accessed April 15, 2025, <https://cstheory.stackexchange.com/questions/18426/advanced-techniques-for-determining-complexity-lower-bounds>
42. Lower Bounds in Theoretical Computer Science - MathOverflow, accessed April

- 15, 2025,
<https://mathoverflow.net/questions/31448/lower-bounds-in-theoretical-computer-science>
43. THE COMPLEXITY OF PARALLEL SORTING* - CiteSeerX, accessed April 15, 2025,
<https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=32a4e877d6940c439dae1697655977672ede138e>
 44. Vectorization guidelines for jax - python - Stack Overflow, accessed April 15, 2025,
<https://stackoverflow.com/questions/69772134/vectorization-guidelines-for-jax>
 45. Change log — JAX documentation, accessed April 15, 2025,
<https://docs.jax.dev/en/latest/changelog.html>
 46. cupy/cupy/_sorting/sort.py at main - GitHub, accessed April 15, 2025,
https://github.com/cupy/cupy/blob/master/cupy/_sorting/sort.py
 47. Sorting, searching, and counting — CuPy 13.4.1 documentation, accessed April 15, 2025,
<https://docs.cupy.dev/en/stable/reference/sorting.html>
 48. Sorting, searching, and counting — CuPy 14.0.0a1 documentation, accessed April 15, 2025,
<https://docs.cupy.dev/en/latest/reference/sorting.html>
 49. cub/test/test_device_histogram.cu at master · dmlc/cub - GitHub, accessed April 15, 2025,
https://github.com/dmlc/cub/blob/master/test/test_device_histogram.cu
 50. Fastest sorting algorithm on GPU currently - CUDA Programming and Performance, accessed April 15, 2025,
<https://forums.developer.nvidia.com/t/fastest-sorting-algorithm-on-gpu-currently/43958>
 51. high performance prefix sum / scan function in CUDA, looking for thrust, cudPP library alternative [closed] - Stack Overflow, accessed April 15, 2025,
<https://stackoverflow.com/questions/18560272/high-performance-prefix-sum-scan-function-in-cuda-looking-for-thrust-cudpp-l>
 52. Is there a similar temporary-allocation feature like CUB for Thrust's thrust::sort_by_key?, accessed April 15, 2025,
<https://forums.developer.nvidia.com/t/is-there-a-similar-temporary-allocation-feature-like-cub-for-thrusts-thrust-sort-by-key/312583>
 53. Sorting with GPUs: A Survey - Hacker News, accessed April 15, 2025,
<https://news.ycombinator.com/item?id=15216142>
 54. Fast 4-way parallel radix sorting on GPUs - SCI Utah, accessed April 15, 2025,
<https://www.sci.utah.edu/~csilva/papers/cgf.pdf>
 55. b0nes164/GPUSorting: State of the art sorting and segmented sorting, including OneSweep. Implemented in CUDA, D3D12, and Unity style compute shaders. Theoretically portable to all wave/warp/subgroup sizes. - GitHub, accessed April 15, 2025,
<https://github.com/b0nes164/GPUSorting>
 56. High Performance Approximate Sort Algorithm Using GPUs - Atlantis Press, accessed April 15, 2025,
<https://www.atlantispress.com/article/24926.pdf>
 57. Performance Analysis of Parallel Sorting Algorithms using GPU Computing - ResearchGate, accessed April 15, 2025,
https://www.researchgate.net/publication/309426134_Performance_Analysis_of_Parallel_Sorting_Algorithms_using_GPU_Computing
 58. `CUDA.quicksort` is not faster than CPU radixsort - GPU - Julia Discourse,

- accessed April 15, 2025,
<https://discourse.julialang.org/t/cuda-quicksort-is-not-faster-than-cpu-radixsort/64036>
59. Quicksort in CUDA: 15x faster than std::sort, 25x faster than std::qsort. - Reddit, accessed April 15, 2025,
https://www.reddit.com/r/CUDA/comments/1fp57xc/quicksort_in_cuda_15x_faster_than_stdsort_25x/
 60. Which algorithm is best for sorting? - Design Gurus, accessed April 15, 2025,
<https://www.designgurus.io/answers/detail/which-algorithm-is-best-for-sorting>
 61. Time and Space Complexity of All Sorting Algorithms - WsCube Tech, accessed April 15, 2025,
<https://www.wscubetech.com/resources/dsa/time-space-complexity-sorting-algorithms>
 62. What are the challenges and opportunities in deploying machine learning models in real-time systems with stringent latency constraints? | ResearchGate, accessed April 15, 2025,
https://www.researchgate.net/post/What_are_the_challenges_and_opportunities_in_deploying_machine_learning_models_in_real-time_systems_with_stringent_latency_constraints
 63. Quality at the Tail of Machine Learning Inference - arXiv, accessed April 15, 2025,
<https://arxiv.org/pdf/2212.13925>
 64. Real-Time Cell Sorting with Scalable In Situ FPGA-Accelerated Deep Learning - arXiv, accessed April 15, 2025, <https://arxiv.org/html/2503.12622v1>
 65. Deep learning with Real-time Inference in Cell Sorting and Flow Cytometry - PMC, accessed April 15, 2025, <https://pmc.ncbi.nlm.nih.gov/articles/PMC6668572/>
 66. MLPerf Inference Rules - GitHub, accessed April 15, 2025,
https://github.com/mlcommons/inference_policies/blob/master/inference_rules.a.doc
 67. Mastering LLM Techniques: Inference Optimization | NVIDIA Technical Blog, accessed April 15, 2025,
<https://developer.nvidia.com/blog/mastering-llm-techniques-inference-optimization/>
 68. Exploring AI Model Inference: Servers, Frameworks, and Optimization Strategies, accessed April 15, 2025,
<https://www.infracloud.io/blogs/exploring-ai-model-inference/>
 69. InferLine: ML Inference Pipeline Composition Framework - UC Berkeley EECS, accessed April 15, 2025,
<https://www2.eecs.berkeley.edu/Pubs/TechRpts/2018/EECS-2018-76.pdf>
 70. Designing Algorithms for Embedded Systems: A Comprehensive Guide - AlgoCademy, accessed April 15, 2025,
<https://algotcademy.com/blog/designing-algorithms-for-embedded-systems-a-comprehensive-guide/>
 71. Optimizing Performance with Embedded Software Solutions - InTechHouse, accessed April 15, 2025,
<https://intechhouse.com/blog/optimizing-performance-with-embedded-software>

[e-solutions/](#)

72. Faster Sorting for Flash Memory Embedded Devices - The University of British Columbia, accessed April 15, 2025,
https://cmps-people.ok.ubc.ca/rlawrenc/research/Papers/2019_embedded_mergesort.pdf
73. What is an appropriate sort algorithm for an embedded system? - Stack Overflow, accessed April 15, 2025,
<https://stackoverflow.com/questions/7357581/what-is-an-appropriate-sort-algorithm-for-an-embedded-system>
74. 5 Factors to Consider Before Choosing a Sorting Algorithm | Towards Data Science, accessed April 15, 2025,
<https://towardsdatascience.com/5-factors-to-consider-before-choosing-a-sorting-algorithm-5b079db7912c/>
75. Sorting algorithms with CUDA - Hacker News, accessed April 15, 2025,
<https://news.ycombinator.com/item?id=43338405>