

Roadmap for Research and Development of a Novel GPU-Accelerated Unique ID Generation Algorithm with Elixir Integration

By: Matheus de Camargo Marques, Web-Engenharia

I. Introduction: Advancing Unique ID Generation with GPU Acceleration

A. Contextualization of the Endeavor

This report outlines a methodological approach for the research and development of a novel Graphics Processing Unit (GPU)-accelerated unique identifier (ID) generation algorithm. The focus lies in empowering developers with prior knowledge in Elixir and various distributed ID schemes, such as CID, CUID, KSUID, NanoID, Snowflake, ULID, and UUID, to innovate in this domain. The task transcends mere implementation of existing algorithms on GPUs, aiming for the creation of an original solution that leverages the intrinsic capabilities of parallel processing architectures.

B. The Imperative for Innovation in ID Generation

The demand for unique, high-throughput, and often sortable IDs has grown exponentially with the expansion of distributed systems, big data analytics, and real-time applications. Traditional ID generation methods, predominantly tied to the Central Processing Unit (CPU), face scalability and performance limitations in scenarios of extremely high demand. This gap motivates the exploration of new frontiers, such as GPU acceleration, to overcome such bottlenecks and meet the requirements of modern systems.

C. GPU Parallelism as a Paradigm Shift for ID Generation

GPUs, primarily conceived for graphics computation, have evolved into powerful general-purpose parallel processing platforms.¹ Their architecture, characterized by thousands of cores capable of executing a vast number of threads concurrently, offers significant, and still largely untapped, potential for generating unique IDs at unprecedented rates and volumes.¹ The challenge and opportunity lie in designing algorithms that are natively parallel and leverage the architectural strengths of GPUs, in contrast to simply porting serial logics designed for CPUs. The quest for a "new" algorithm implies investigating how the specific characteristics of GPUs can be innovatively utilized for this purpose.

D. Scope and Objectives of the Report

This document proposes a research and development roadmap, covering fundamental knowledge about ID generation and GPU architectures, design principles for GPU-centric ID algorithms, implementation strategies with Elixir integration, and validation methodologies. The central objective is to provide the conceptual tools and a structured approach that enable the researcher/developer to pioneer their own GPU-accelerated ID generation solution, with an emphasis on uniqueness, sortability, and high performance.

II. Foundations: Understanding ID Generation and GPU Paradigms

A. Critical Review of Existing ID Architectures and their GPU Adaptability

A thorough understanding of existing ID schemes is crucial for identifying innovation opportunities in a GPU context.

Detailed Analysis of Common ID Schemes

- **UUID (Universally Unique Identifier):**
 - UUIDv1 combines a timestamp, a node identifier (originally MAC address), and a clock sequence. Its components offer some basis for parallelization, but MAC uniqueness can be problematic in virtualized environments.⁵ UUIDv4, based on randomness, is easily parallelizable on GPUs via Pseudo-Random Number Generators (PRNGs).⁶
 - *Critique:* Standard UUIDs, especially v4, are not inherently time-sortable, a significant drawback for many applications.⁷ Their 128-bit size and string representation can be inefficient for storage and indexing.⁸
 - *GPU Adaptability:* High for v4 randomness generation. For v1-like structures, timestamp and clock sequence could be managed per thread/block.
- **Snowflake and its Derivatives:**
 - Components: Timestamp (millisecond precision), worker/machine ID, and sequence number.⁸ Typically 64 bits. The canonical structure, like 41 bits for timestamp, 10 for worker ID, and 12 for sequence⁸, is a strong candidate for GPU adaptation.
 - *Strengths:* Time-ordered, compact, high throughput on CPUs. The worker ID ensures uniqueness among nodes.
 - *GPU Adaptability:* Highly adaptable. The timestamp can be a global reference for a kernel launch or per block. The worker ID can be mapped to GPU hardware identifiers (like SMs or blocks). Sequence numbers are ideal for GPU atomic operations.
- **KSUID (K-Sortable Unique Identifier):**
 - Components: Timestamp (second precision, 32 bits) and a 128-bit random payload.¹¹ Total of 20 bytes. The 32-bit timestamp offers approximately 136 years from the defined epoch.¹¹
 - *Strengths:* Lexicographically sortable by time due to the timestamp prefix, collision-resistant due to the large random payload.
 - *GPU Adaptability:* Timestamp management and parallel generation of the random payload are GPU-friendly.
- **ULID (Universally Unique Lexicographically Sortable Identifier):**
 - Components: Timestamp (48 bits, millisecond precision) and 80 bits of randomness.¹⁰ Total of 128 bits. The 48-bit timestamp provides vast temporal coverage.
 - *Strengths:* Lexicographically sortable, millisecond-precision timestamp, good randomness. Base32 encoding is human-readable.
 - *GPU Adaptability:* Similar to KSUID, timestamp and randomness generation are parallelizable.
- **UUIDv7 (New time-based UUID proposal):**
 - Components: Unix timestamp in milliseconds (in the most significant bits), followed by random or sequence bits.⁸
 - *Strengths:* Aims to combine UUID ubiquity with sortability, offering better database indexing performance than UUIDv4.¹⁵
 - *GPU Adaptability:* Very promising. The timestamp can be generated on the GPU, and the remaining bits can be filled with GPU-generated random numbers or sequences.

A convergence is observed in more recent ID schemes (KSUID, ULID, UUIDv7) that explicitly prioritize time-based sortability, often allocating the timestamp to the most significant portion of the identifier. Schemes like Snowflake already intrinsically incorporate this feature. Initially, formats like UUIDv4 focused on global uniqueness through randomness, resulting in difficulties for database indexing and a lack of natural ordering. Systems requiring temporal analysis or event sourcing found this an obstacle. Snowflake emerged as a solution, providing sortability and high throughput by structuring the ID into timestamp, worker ID, and sequence. KSUID and ULID followed this trend, being specifically designed for lexicographical sorting, complemented by robust random components for collision resistance. UUIDv7 represents an effort to incorporate these benefits into the UUID standard. For the design of a GPU algorithm, this trend strongly suggests considering a time-ordered structure, aligning with modern requirements and leveraging the GPU's ability to operate with time and sequences. The componentized nature (timestamp, worker, sequence) also proves highly compatible with parallel decomposition on GPUs.

The following table offers a comparative analysis, highlighting the potential for GPU parallelization:

Table 1: Comparative Analysis of Existing ID Schemes and GPU Parallelization Potential

ID Scheme	Main Components (Bits/Timestamp Precision, Node ID Bits, Sequence Bits, Randomness Bits)	Total Size	Sortability Mechanism	Primary Uniqueness Mechanism	Known CPU Performance Bottlenecks	GPU Parallelization Opportunities	GPU Challenges
UUIDv1	Timestamp (60 bits, 100ns), Node (48 bits MAC), Clock Sequence (14 bits)	128 bits	Partial (timestamp)	MAC + Timestamp + Clock Seq.	System calls for time, MAC uniqueness	Parallel generation of timestamp/clock seq. per thread/block	High-precision time synchronization, unique Node ID mapping
UUIDv4	Randomness (122 bits), Version/Variant (6 bits)	128 bits	None	Randomness	PRNG quality and speed	Parallel PRNGs (cuRAND)	PRNG state management per thread
UUIDv7	Timestamp (typically 48 bits, ms), Random/Sequence (74 bits), Version/Variant (6 bits)	128 bits	Timestamp	Timestamp + Random/Sequence	(New, less data)	Parallel generation of timestamp and random/sequential part	Time synchronization, sequence management
Snowflake	Timestamp (41 bits, ms), Worker ID (10 bits), Sequence (12 bits)	64 bits	Timestamp	Timestamp + Worker ID + Sequence	Clock synchronization, Worker ID coordination	Timestamp per kernel/block, Worker ID mapped to blockIdx/%smid, atomic sequences	Contention on atomics, time synchronization between GPUs
KSUID	Timestamp (32 bits, s), Random Payload (128 bits)	160 bits (20 bytes)	Timestamp	Timestamp + Random Payload	System calls for time, PRNG	Parallel generation of timestamp and random payload	Time synchronization (if granularity < s desired)
ULID	Timestamp (48 bits, ms), Randomness (80 bits)	128 bits	Timestamp	Timestamp + Randomness	System calls for time, PRNG	Parallel generation of timestamp and randomness	Time synchronization

This table serves as a guide to understand the fundamental characteristics of existing IDs and how they align (or conflict) with GPU computing paradigms, informing the design process of a new algorithm.

B. GPU Architecture Imperatives for Massively Parallel ID Generation (CUDA/OpenCL)

The GPU architecture imposes a specific programming model that must be understood for the efficient development of ID generation algorithms.

- **Kernel Execution Model:** A single function (kernel) is executed by multiple threads in parallel on the GPU. ¹⁷ Each thread performs the same computation but on distinct data or with a unique identifier, which is fundamental for ID generation, where each thread can be responsible for generating one or more IDs.
- **Hierarchy: Grids, Blocks, Warps, Threads:**
 - **Grid:** A collection of all thread blocks executing a kernel. ²⁰ The grid dimensions define the total scope of parallel work.
 - **Thread Block:** A group of threads (up to 1024 on modern NVIDIA GPUs ²¹) that can cooperate using shared memory and synchronize their execution (`__syncthreads()` ¹⁸). Blocks are scheduled for execution on Streaming Multiprocessors (SMs).
 - **Warp:** A group of 32 threads within a block that execute instructions in lockstep (SIMT - Single Instruction, Multiple Threads). ²⁴ Warp behavior is critical for performance, especially regarding execution divergence.
 - **Thread:** The most granular unit of execution, identified by `threadIdx` within a block and `blockIdx` within a grid. ²⁰
- **Streaming Multiprocessors (SMs):**
 - The central processing units of the GPU. Each SM can execute multiple thread blocks concurrently ²⁰ and has its own execution cores, shared memory, and registers. ¹⁷ The number of SMs per GPU varies by architecture. ²⁸ The special register `%smid` allows obtaining the SM ID ³⁰, potentially serving as a candidate for a "worker ID" component.
- **Memory Hierarchy:**
 - **Global Memory:** Large, but with high latency; accessible by all threads across all SMs. ¹⁷
 - **Shared Memory:** Small, low-latency, private to threads within a single block. ¹⁷ Crucial for inter-thread communication and caching within a block.
 - **Registers:** Fastest memory, private to each thread. ¹⁷
 - **Constant/Texture Memory:** Specialized read-only caches. ¹⁷
- **Atomic Operations:**
 - Essential for thread-safe updates to shared or global memory locations (e.g., incrementing sequence counters). ⁶ The performance characteristics of atomic operations on shared versus global memory differ significantly, with shared memory atomics generally being faster but with limited scope. ³¹ Contention on atomic operations can be a major bottleneck.

The hierarchical structure of the GPU (Grid -> Block -> Warp -> Thread) offers a natural mapping for assigning and managing the different parts of an ID. A globally unique ID often requires components denoting distinct scopes (e.g., global time, machine/process, local counter). A GPU kernel launch defines a Grid; the timestamp could be established at the Grid level (e.g., kernel start time). Blocks are scheduled on SMs; `blockIdx` (unique per block in the grid) or `%smid` (unique per SM) could form part of a "worker ID." Threads within a block (`threadIdx`) can each generate a sequence number, potentially using GPU atomic operations on shared memory for coordination within the block, or global atomics for grid-wide sequences if managed carefully. This correspondence can simplify the conceptual design of the ID structure, allowing the inherent organization of the GPU to be leveraged for ensuring uniqueness.

C. Performance Considerations: CPU vs. GPU for ID Generation Workloads

The choice between CPU and GPU for ID generation fundamentally depends on the scale and performance requirements of the application.

- **CPU Strengths:** Excellent for serial tasks, complex branching logic, and tasks requiring large amounts of memory per "thread." ¹ Generating a single ID on a CPU is typically very fast.
- **GPU Strengths:** Massively parallel execution of simpler, repetitive tasks. ¹ Ideal when millions or billions of IDs are needed in a short time frame.
- **Trade-offs for ID Generation:**
 - *Latency vs. Throughput:* A CPU can generate a single ID with lower latency. A GPU aims for extremely high throughput (many IDs per unit of time) but may have higher initial latency due to kernel launch overhead and data transfer.
 - *Branching:* ID generation logic is generally simple, minimizing the GPU's weakness with complex branching. ¹
 - *Data Transfer:* If IDs are generated on the GPU but consumed by the CPU, the overhead of transfer via the PCI-e interface needs to be considered. Generating IDs directly where they are needed by other GPU kernels would be most efficient.
- The analogy that "the CPU is like a head chef... The GPU is more like a junior assistant with ten hands who can flip 100 burgers in 10 seconds" ² aptly describes the scenario: for a few IDs, the CPU suffices; for millions concurrently, the GPU excels.

GPU ID generation is a strategy focused on throughput, not (usually) on latency for individual IDs. The cost of launching a GPU kernel and transferring data means that for generating a small number of IDs, the GPU will likely be slower than a CPU. This setup

cost is amortized only when a very large number of IDs are generated in parallel by the GPU threads. Therefore, the algorithm must be designed to maximize the number of IDs generated per kernel launch, and the use case should be one where massive bursts of IDs are required, justifying the GPU overhead, such as in bulk data ingestion, event tagging in large-scale simulations, or initialization of large distributed datasets.

III. Designing a Novel GPU-Accelerated Unique ID Generation Algorithm

A. Core ID Components Reimagined for GPU Parallelism

Designing a new ID algorithm for GPUs requires reconsidering how traditional ID components (timestamp, node/worker identifier, sequence, and randomness) can be efficiently implemented within the GPU's parallel architecture. The following table summarizes mapping strategies.

Table 2: ID Component Mapping Strategies for GPU

ID Component	GPU Resource/Technique	Pros	Cons	Key Considerations for New Algorithm
Timestamp	Host-supplied timestamp per kernel; <code>clock64()/%globaltimer</code> per thread/block	Precision, uniqueness scope, performance	Synchronization overhead, resource limits, contention	Desired ID bit length, sortability needs, collision probability tolerance
Node/Worker ID	<code>blockIdx</code> , <code>%smid</code> , <code>threadIdx</code>	Uniqueness scope, performance	Resource limits (<code>%smid</code> not contiguous)	Number of logical "workers" needed, uniqueness across multiple GPUs
Sequence	<code>atomicAdd()</code> on shared/global memory	Uniqueness within scope (block/grid)	Contention, latency of global atomics	Generation rate per "worker," overflow handling
Randomness	<code>cuRAND</code> (per-thread states); custom PRNG kernels (XORShift, Mersenne Twister)	Collision resistance, bit filling	PRNG state management overhead, quality vs. performance	Number of random bits needed, acceptable collision probability

1. High-Precision Timestamping on the GPU: Strategies for Accuracy, Monotonicity, and Synchronization

- **GPU Time Sources:**
 - **Host-Supplied Timestamp:** The simplest approach is for the CPU to get the current time and pass it as an argument to the kernel. All threads in that kernel launch would use this single timestamp. Pros: simple, synchronized across all threads for that batch. Cons: granularity limited to kernel launch frequency.
 - **GPU Clocks (e.g., `clock()`, `clock64()` in CUDA, PTX `%globaltimer`):** GPUs have high-resolution internal counters. ³⁴ `clock64()` returns clock cycles per SM, not directly wall time. `%globaltimer` may not be synchronized with the host clock or between different GPUs without careful handling. ³⁴ Converting cycles to wall time requires calibration. `torch.cuda.Event(enable_timing=True)` is used for timing kernel execution, implying the GPU's ability to record event timestamps. ³⁵
- **Monotonicity:** Crucial for sortable IDs. Wall clock time can go backward (NTP synchronization, daylight saving). ³⁶ Monotonic clocks only advance. The PTX `%globaltimer` is stated as monotonic *in device code* ³⁴, making it a strong candidate.
- **Synchronization and Precision:** `torch.cuda.synchronize()` ensures the host waits for the GPU, important for accurate end-to-end time measurement, but not for timestamping per ID within a kernel. ³⁵ Achieving a globally synchronized, high-precision timestamp for *each ID generated by thousands of threads simultaneously* is non-trivial. With multiple GPUs, timestamp synchronization becomes even more complex. ³⁷

- **Strategies for New ID:**
 - **Epoch-Based Timestamp:** Use a custom epoch (e.g., project start date) to maximize the usable time range within the allocated bits, similar to Snowflake.⁹
 - **Combined Approach:** A base timestamp supplied by the host for a kernel launch, with GPU threads adding fine-grained clock64() deltas or per-block sequence numbers to differentiate IDs within that host-defined time window.

GPU timestamping for IDs involves a trade-off between granularity, synchronization, and simplicity. While GPUs offer high-resolution clocks like %globaltimer, using them directly for the primary timestamp of each ID generated by millions of threads per second introduces complexity in converting to wall time and ensuring synchronization across SMs/GPUs. The simplest approach, a single host timestamp per batch, sacrifices granularity. Conversely, each thread reading %globaltimer offers high granularity, but the meaning of this value in terms of wall time and its relation to readings on other SMs or GPUs is problematic.³⁴ An approach similar to Snowflake, where the timestamp has millisecond precision, seems balanced. The GPU can then focus on generating the sequence/worker parts *for that millisecond*. If higher precision than milliseconds is embedded *from the GPU clock*, careful planning is needed for its interpretation and global meaning. This could be an area for innovation, such as using the upper bits of %globaltimer to further subdivide milliseconds.

2. Worker/Node Identification in a GPU Grid: Mapping SM, Block, and Thread IDs to ID Components

- **Purpose:** Differentiate IDs generated by different parallel units, especially if timestamps and sequence numbers could collide between these units.
- **GPU Hardware Identifiers:**
 - blockIdx.{x,y,z}: Unique identifier for a thread block within the grid.²⁰ Max dimensions: (2³¹-1, 65535, 65535).²¹
 - threadIdx.{x,y,z}: Unique identifier for a thread within its block.²⁰ Max dimensions: (1024, 1024, 64) with total threads/block <= 1024.²¹
 - %smid: Special register to get the ID of the Streaming Multiprocessor executing the thread.³⁰ The numbering of %smid is not guaranteed to be contiguous, and %nsmid (upper bound) can be larger than the physical number of SMs.³⁰ The H100 has 144 SMs²⁸, and the A100 has 108.²⁹
- **Strategies for Deriving a "Worker ID":**
 - **Using blockIdx:** A simple approach is to use blockIdx.x (or a combination if using a 2D/3D grid of blocks) as the worker ID.
 - **Using %smid:** Could provide a more "physical" worker ID. However, its non-contiguous nature is a challenge.
 - **Hierarchical ID:** Combine smid and blockIdx (within that SM).
 - **Multi-GPU:** A GPU-specific identifier (e.g., its device ID 0, 1, 2...) must be incorporated, likely managed by the host.³⁸
- **Bit Allocation:** If following Snowflake (10 bits for worker ID⁸), this allows for 1024 unique workers.

blockIdx emerges as a more flexible and scalable source for a "Worker ID" within a kernel launch compared to %smid. While %smid is tied to the physical number of SMs (e.g., 108-144), which is relatively small, blockIdx can theoretically scale to millions or billions.²¹ A Snowflake-style "worker ID" allows for many distinct generating entities (e.g., 1024). If %smid is used, the number of unique worker IDs is limited by the physical SMs, which might be insufficient. blockIdx can be much larger by configuring a larger grid, decoupling the logical worker ID from the physical SM count. Thus, using blockIdx for the "worker ID" component is generally more scalable and configurable for a single kernel launch.

3. Parallel Sequence Number Generation: Advanced use of GPU atomic operations

- **Purpose:** Ensure uniqueness for IDs generated by different threads within the same time unit and by the same "worker" (e.g., Snowflake's 12 bits for 4096 IDs/ms/worker⁸).
- **GPU Atomic Operations:** atomicAdd(), atomicInc(), etc., are crucial.³²
- **Scope of Atomics:**
 - **Global Memory Atomics:** Operate on GPU global memory. Cons: higher latency, significant contention.³²
 - **Shared Memory Atomics:** Operate on shared memory, private to a thread block.¹⁷ Pros: much lower latency.³¹ Cons: sequence is unique only *within that block*. Maxwell and newer architectures have native instructions, improving performance over earlier lock/unlock loop-based implementations.³³
- **Strategies for Sequence Generation:**
 - **Global Counter:** A single global atomic counter. Simple, but likely a bottleneck.
 - **Per-Block Counter (Shared Memory):** Each block maintains its own sequence using atomicAdd() on a __shared__ variable. Highly parallel.
 - **Hierarchical Approach³²:** Threads within a warp/block increment a local counter in shared memory. Periodically, one thread per block could atomically grab a larger chunk from a global counter, reducing pressure on the global atomic.

- **Contention and Performance:** High contention on a single atomic variable serializes threads.³² Bank conflicts in shared memory can also serialize accesses.¹⁷ The performance of shared memory atomics is load-dependent.³¹ The analogy to CTR (counter) mode for ciphers⁴⁰, which is "embarrassingly parallel," applies here.

Hierarchical sequence generation is fundamental for GPU scalability. Relying solely on global atomic operations for sequence numbers for every thread will not scale due to contention. Shared memory atomics are faster but local to the block. The goal is for each of the (potentially) millions of threads to obtain a unique sequence number within its group. A single global atomic is a clear bottleneck.³² Shared memory atomics allow each block to manage its own sequence space independently.³¹ A design using per-block sequence numbers (via shared memory atomics) combined with a unique blockIdx (as part of the worker ID) is a robust and scalable approach, avoiding global atomic bottlenecks for the most frequent operation.

4. Incorporating Entropy: GPU-Accelerated PRNGs and managing parallel random streams

- **Purpose:** For ID schemes like KSUID, ULID, or UUIDv4/v7 that require a random component.
- **GPU PRNG Libraries/Techniques:**
 - **cuRAND (NVIDIA):** Standard for NVIDIA GPUs; generates multiple parallel streams.
 - **Custom Kernels:** Implementations of PRNGs like Mersenne Twister (MT19937) or XORShift/XORWOW.⁴⁴ describes generating a vector of random numbers on the GPU. ⁶ discusses PRNGs for Monte Carlo on GPUs, noting that CPU PRNGs are often inadequate.
- **Managing State for Parallel Streams:** Each thread generating random numbers needs its own independent PRNG state. Seeds for each thread must be unique to ensure uncorrelated random sequences.⁶
- **Performance:** PRNGs like XORShift are very fast.⁴ Generating random numbers on the GPU is generally much faster in bulk than on the CPU.⁶

Table 3: GPU-Accelerated PRNG Techniques for ID Generation

PRNG Algorithm	State Size per Thread	Initialization (Seed Strategy)	Parallel Stream Management	Quality/Period	Relative GPU Performance	Suitability for ID Random Component
cuRAND (XORWOW)	Small	Managed by cuRAND (seed by sequence/off set)	Natively supported	Good for many purposes, period 2192–232	Very High	Suitable for ULID (80 bits), smaller UUIDv7 components
cuRAND (MTGP32)	Medium	Managed by cuRAND	Natively supported	Very High, period 211213–1	High	Suitable for KSUID (128 bits), UUIDv4
Custom XORShift Kernel	Very Small	Unique seed per thread (e.g., global_id + master_seed)	Manual (each thread is a stream)	Variable (depends on variant), can be good	Very High	Suitable if fine control and minimal state are critical
Custom Mersenne Twister (MT19937) Kernel	Large	Unique seed per thread, complex initialization	Manual	Excellent, period 219937–1	Medium-High (due to state)	Overkill for most IDs unless cryptographic quality is paramount

The seeding strategy is paramount for parallel PRNGs in ID generation. If multiple GPU threads use identical or correlated seeds, the "random" parts of the generated IDs will not be truly independent, potentially leading to higher collision rates. The purpose of the random component in IDs like ULID/KSUID is to ensure uniqueness, especially for IDs generated at the same instant by different workers. ⁴ and ⁶ emphasize the need for distinct seeds and sequences per thread. A robust seeding strategy is critical, potentially involving the host generating a large array of unique seeds or a kernel generating unique seeds for each thread based on its global ID and a master seed.

B. Ensuring Global Uniqueness and Collision Resistance in a Distributed GPU Context

The central strategy, similar to Snowflake⁸, is that the combination of (Timestamp + Worker ID + Sequence Number) ensures uniqueness.

- **Timestamp:** Differentiates IDs from different time units.
- **Worker ID:** Differentiates IDs generated by different logical/physical units *at the same time*. On the GPU, this can be blockIdx, %smid, or a host-assigned GPU ID. ⁸ discusses server ID coordination for Snowflake using etcd, illustrating the need for unique generator identities.
- **Sequence Number:** Differentiates IDs generated by the same worker *within the smallest timestamp unit*. GPU atomics are key here.
- **Randomness:** Provides additional collision resistance.
- **Clock Rollback:** A known issue for time-based IDs like Snowflake. ⁹ The solution involves dedicating bits within the worker ID or sequence to a "rollback counter." ⁹ On the GPU, detecting clock rollback would require threads to compare the current GPU time with a last known time.
- **Collision Probability:** For random parts (UUIDv4, KSUID, ULID), the collision probability is astronomically low due to the large bit space. ⁸ For structured parts (Snowflake), if components are managed correctly, collisions are prevented by design.
- **Multi-GPU Systems:** Each GPU needs a unique identifier incorporated into the "Worker ID" component. Timestamp synchronization between different GPUs can be challenging. ³⁷

Distributed uniqueness on the GPU relies on orthogonal ID components and careful state management. The timestamp separates IDs over time. The Worker ID (e.g., `gpu_id << N | blockIdx`) separates IDs generated by different blocks/GPUs at the same time. The sequence number (per block) separates IDs from different threads in the same block at the same time. Randomness adds an extra layer of collision prevention. Failure in any of these (e.g., non-unique worker ID, sequence counter not resetting correctly with the timestamp, unhandled clock rollback) can lead to collisions. GPU parallelism means many "workers" (blocks/threads) are active simultaneously, so the logic for each component must be robust to this concurrency.

C. Achieving Sortability: Adapting KSUID/ULID/UUIDv7 Principles for Efficient GPU Implementation

Sortability is an increasingly desired feature in IDs.

- **Timestamp Primacy:** The most significant bits of the ID must represent time for natural lexicographical sorting. ⁷
- **Timestamp Granularity and Range:** KSUID uses 32 bits for seconds ¹¹; ULID uses 48 bits for milliseconds ¹⁰; Snowflake uses 41 bits for milliseconds ⁸; UUIDv7 typically uses 48 bits for milliseconds. ¹⁵ For the GPU, higher precision (e.g., nanoseconds from `clock64()`) could be used, but milliseconds are often sufficient.
- **Secondary Sorting Components:** If multiple IDs are generated within the same minimal time unit, subsequent bits (worker ID, sequence, randomness) determine their relative order.
- **GPU Implementation:** The GPU kernel would construct the ID via bit shifts and OR operations of the components in the correct order.
- **Performance of Sortable IDs:** UUIDv7 shows significant insertion performance improvements in databases compared to UUIDv4 due to better index locality. ¹⁵

Sortability on the GPU is primarily a structural issue, aided by the parallel generation of components. The GPU's role in sortability is to correctly assemble the ID components in the predefined order, with the timestamp being the most significant. GPU threads will independently fetch/calculate their timestamp, worker ID, and sequence/randomness parts, and then assemble these parts into the final ID. There is no special "sorting" algorithm running on the GPU; it's an emergent property of its structure. The kernel design for combining ID components is critical, and the choice of bit lengths for each component will determine the overall ID size and the range/granularity of each part.

D. Parallel Algorithm Design Patterns for ID Generation

Applying parallel algorithm design patterns is essential.

- **Task Decomposition** ⁴¹: The problem of generating N IDs is decomposed into N independent tasks, each assigned to a GPU thread.
- **Granularity** ⁴¹: Fine-grained (each thread generates one ID) maximizes parallelism. Coarse-grained (each thread generates a small block of IDs) can reduce kernel launch overhead.
- **Degree of Concurrency** ⁴¹: Aim to utilize as many GPU cores as possible.
- **Dependencies and Data Interaction** ⁴¹: Ideally, ID generation per thread is largely independent. Dependencies arise from shared timestamp sources, shared worker ID pools, sequence number generation (atomic operations), and PRNG state.

- **Cost Model** ⁴²: Consider computational cost and memory access cost. Atomic operations can dominate the cost if highly contended.

ID generation can be considered "embarrassingly parallel" with controlled "choke points." Most of the work for each ID can be done independently by each thread. If a thread has its unique blockIdx, threadIdx, and a unique PRNG seed, it can compute most parts of its ID without interacting with other threads. The sequence number is the main point of coordination needed between threads. This makes the problem almost perfectly parallel, except for these controlled choke points (atomics). The algorithm design should maximize independent work and minimize the scope and frequency of operations requiring synchronization, like global atomics. Using shared memory atomics for per-block sequences is a key strategy.

IV. Implementation Strategies: From CUDA/OpenCL Kernel to Elixir Integration

A. GPU Kernel Development: Best practices for ID generation logic, memory optimization, and minimizing thread divergence

Kernel creation is the core of development.

- **Language Choice**: CUDA for NVIDIA GPUs is likely, given the pursuit of performance.
- **Kernel Signature**: Define input parameters (e.g., pointer to output ID array, host-supplied base timestamp) and launch configuration (gridDim, blockDim).
- **ID Component Generation Logic**: Implement functions/logic within the kernel for each ID component as designed in Section III.A.
- **ID Assembly**: Combine components using bit shifts and ORs.
- **Memory Access Patterns**: Output IDs to global memory. Ensure coalesced writes if possible. ¹⁷ Minimize global memory access; use shared memory for temporary per-block data.
- **Thread Divergence** ¹: Avoid conditional logic (if/else) that varies between threads in the same warp. ID generation logic is typically linear code, minimizing this issue.
- **Error Handling**: Limitations in how to signal errors from a kernel.
- **Optimized Kernel Generation** ⁴³: LLMs can assist, but creating optimized and correct kernels for a novel algorithm still requires expertise and iterative refinement. Optimal GPU thread mapping is non-trivial. ⁴³
- **GPU Hardware Resilience** ⁴⁴: Awareness of potential GPU errors (GSP errors, NVLink issues, memory errors) is important for robust system design.

The kernel design should prioritize data locality for counters and coalesced outputs. Sequence counters, if per-block, should reside in shared memory for fast atomic access. ³² The final generated IDs will be written to global memory. To maximize write bandwidth to global memory, writes should be coalesced. If thread k in a warp writes to output_array[global_thread_id_base + k], this is naturally coalesced. Minimizing any other global memory reads/writes within the kernel's main loop is also crucial.

B. Interfacing with Elixir

Integration with Elixir is a key requirement.

1. Natively Implemented Functions (NIFs): Deep dive into performance, safety, and GPU resource management

NIFs offer the highest performance integration between Elixir and native code. ⁴⁵

- **Mechanism**: Write a C/C++ wrapper function that uses the CUDA Runtime API (e.g., cudaMalloc, cudaMemcpy, kernel launch <<<...>>>, cudaFree). This wrapper is exposed to Elixir as a NIF. ⁴⁵
- **Performance**: Very low call overhead (~0.1-1 μ s ⁴⁵).
- **Safety** ⁴⁵: A crash in the NIF (e.g., CUDA error, segfault) will bring down the entire BEAM (Erlang VM). This is the main drawback. Mitigation: rigorous testing of CUDA code, careful error handling in the C wrapper.
- **Resource Management**: GPU memory allocation/deallocation (cudaMalloc/cudaFree) must be managed by the NIF.
- **Long-Running NIFs**: A NIF that takes too long can block a BEAM scheduler. ⁴⁵ ID generation kernels should be fast.
- **Elixir NIFs for GPU** ⁴⁶: NIFs are a viable path for GPU computation from Elixir, with CUDA/OpenCL mentioned as use cases.

2. Leveraging Elixir GPU Libraries (e.g., Hok): Assessing suitability for custom, high-performance ID kernels

- **Hok** ⁴⁶: An Elixir DSL for higher-order GPU kernels, aiming to separate coordination (Elixir) from computation (GPU kernel). Allows device functions to be passed to kernels. ⁴⁶ Uses the GPotion DSL and is compatible with Matrex. ⁴⁶ Examples show Hok.spawn and data management with Hok.new_gmatrex / Hok.get_gmatrex. ⁴⁶ The GitHub repository (ardubois/hok ⁴⁷)

suggests it's a research project, possibly not mature for production or highly custom low-level work. The project status ⁴⁷ (17 stars, 0 forks, no releases) suggests it is likely experimental.

- **Suitability for New ID Generation:** Pros: may simplify boilerplate for CUDA context, memory management, kernel launch from Elixir. Cons: as a higher-level DSL, may not expose the fine-grained control needed for a novel, highly optimized ID generation kernel.

NIFs are likely unavoidable for maximum performance and control of a new GPU ID algorithm. While libraries like Hok aim to simplify GPU programming in Elixir, a truly novel and high-performance ID generator will likely require bit-level manipulation, custom data structures on the GPU, and fine-tuned control over CUDA primitives. The user's goal is a *new* algorithm, implying they will write custom CUDA/OpenCL kernels. Existing Elixir GPU libraries are either general-purpose (like Hok, which appears experimental ⁴⁷) or target specific domains (like ML ⁴⁸). They may not offer the flexibility to inject and manage a highly specialized ID generation kernel. NIFs provide the most direct path to call custom C/C++/CUDA code from Elixir, offering maximum performance and control over GPU execution. ⁴⁵ The trade-off is complexity and safety. Therefore, the primary development effort for the GPU part will be the CUDA/OpenCL kernel itself, and Elixir integration will likely involve writing a NIF wrapper.

Table 4: Elixir-GPU Integration Strategies for Custom ID Kernels

Integration Method	Key Features	Performance Overhead	Development Complexity	Flexibility for Custom Kernels	Error Handling/Safety (BEAM Stability)
Custom NIFs (C/CUDA)	Direct CUDA API access, full control	Low (NIF call), data marshalling	High (Elixir + C/CUDA)	Very High (bit-level control, shared memory, atomics)	Risk of BEAM crash, manual error handling in C
Custom NIFs (Rust/CUDA via Rustler)	Rust memory safety, CUDA API access via bindings	Low-Medium (similar to C NIF + Rustler overhead)	High (Elixir + Rust/CUDA)	High (depends on quality of CUDA bindings for Rust)	Better memory safety on Rust side, but CUDA errors can still impact BEAM
Hok Library (if viable)	Higher-level abstraction, DSL	Medium-High (depends on DSL implementation)	Medium (focus on Elixir and Hok DSL)	Limited (depends on DSL capabilities)	Potentially safer if library manages CUDA errors well, but less control

V. Validation, Benchmarking, and Iteration

A. Rigorous Testing for Uniqueness, Sortability, and Collision Probability

Validation is critical for any ID generation system.

- **Uniqueness Tests:** Generate massive batches of IDs (billions/trillions if feasible) and check for duplicates. Test under different "worker ID" configurations and edge cases (sequence overflow, timestamp boundaries, clock rollback).
- **Sortability Tests:** Generate IDs over a period where timestamps advance. Lexicographically sort the IDs and verify that the order matches the generation order.
- **Collision Probability Analysis (for random components):** Theoretical analysis based on the number of random bits. Empirical tests (Birthday Problem) can be run if the random space is smaller, but for 64+ random bits, theoretical probability is usually sufficient.⁸

B. Performance Benchmarking: Throughput, Latency, and Scalability Analysis

Measuring performance comprehensively is essential.

- **Metrics:**
 - **Throughput:** IDs generated per second (on GPU and end-to-end, including Elixir NIF call and data transfer).
 - **Latency:** Time to generate a single ID (less relevant for GPU) vs. time to generate a batch. Kernel execution time (e.g., using CUDA events³⁵). NIF call overhead.
 - **Scalability:** How throughput changes with the number of GPU threads/blocks, different GPU architectures, and (if multi-GPU) number of GPUs.
- **Benchmarking Tools:** NVIDIA Nsight Compute/Systems for detailed GPU kernel profiling.²⁸ Elixir's Benchee for benchmarking NIF calls from the Elixir side. Custom microbenchmarks.
- **Comparison Points:** Benchmark against existing CPU-based Elixir ID libraries (KSUID¹¹; ULID⁴⁹; Snowflake implementations⁵⁰). UUIDv4 vs UUIDv7 benchmarks¹⁵ are relevant. Snowflake generators in Rust (~233 ns/iter⁵⁰) and Swift (~125M IDs/s⁵²) provide useful references.

Benchmarking should isolate GPU kernel performance from NIF and data transfer overheads. The raw GPU kernel speed might be extremely high, but the end-to-end throughput observed in Elixir could be limited by NIF call overhead, data copying between CPU and GPU memory, or Elixir processing of results. It's crucial to benchmark at different levels: (1) Pure GPU kernel execution (using CUDA events, Nsight). (2) NIF call including data transfers but minimal Elixir processing. (3) Full end-to-end Elixir function call. This helps identify bottlenecks. If the kernel is 100x faster than CPU methods, but data transfer consumes 90% of the NIF call time, the real-world speedup will be much less than 100x.

VI. Research and Development Roadmap: Key Milestones for Independent Development

A structured roadmap will facilitate the R&D process.

- **A. Phase 1: Deep Dive and Fundamental Kernel Design**
 - Solidify understanding of GPU architecture (CUDA/OpenCL), memory model, atomic operations, and PRNGs.
 - Define the bit structure of the new ID. Justify choices based on desired properties.
 - Develop a basic GPU kernel (e.g., in CUDA C++) implementing the core ID generation logic. Focus on correctness first.
- **B. Phase 2: Elixir Integration and Basic Prototyping**
 - Develop a NIF wrapper for the basic kernel.
 - Implement Elixir functions to call the NIF and retrieve generated IDs.
 - Initial tests for uniqueness and basic sortability with small batches.
- **C. Phase 3: Performance Optimization and Advanced Features**
 - Profile the GPU kernel using Nsight Compute. Identify and resolve bottlenecks.
 - Optimize NIF data transfers.
 - Implement advanced features: clock rollback handling, robust error reporting.
 - Refine PRNG seeding and stream management.
- **D. Phase 4: Rigorous Validation and Benchmarking**
 - Conduct extensive uniqueness, sortability, and collision tests as per Section V.A.
 - Perform comprehensive performance benchmarks as per Section V.B. Compare against CPU libraries.
- **E. Phase 5: Documentation and Iteration**

- Document the algorithm, design choices, NIF API, and usage.
 - Iterate on the design and implementation based on testing and benchmarking results.
- F. Continuous Learning: Staying Abreast of GPU Technology and ID Generation Advancements**⁴³
 - The field of GPU computing⁴³ and ID generation (e.g., new UUID versions) is constantly evolving.

VII. Concluding Remarks and Future Research Avenues

This report has laid out a pathway for developing a novel GPU-accelerated ID generation algorithm, with a focus on Elixir integration. The potential for significant performance gains is considerable but requires a deep understanding of both ID paradigms and the intricacies of GPU programming.

Key considerations revolve around defining an ID structure that is inherently parallelizable and sortable, efficiently managing high-precision timestamps, deriving "worker" identifiers from the GPU hierarchy, implementing scalable sequence mechanisms using atomic operations, and, if necessary, incorporating entropy via parallel PRNGs. The interface with Elixir, likely through NIFs, is crucial for translating GPU power into accessible functionality within the Elixir ecosystem.

Future research avenues could include:

- Development of self-tuning ID generation kernels that adapt to GPU load or application requirements.
- Exploration of adaptive ID structures that can change their composition (e.g., bits allocated for sequence vs. randomness) dynamically.
- Investigation into potential hardware-assisted ID generation features in future GPU architectures.
- Application of formal verification methods to prove the uniqueness and sortability properties of GPU-based ID algorithms.
- Study of alternative GPU programming models (e.g., SYCL) or languages for different application contexts that might simplify or further optimize this type of task.

The journey to create a novel and effective GPU ID generator is challenging, but the rewards in terms of the ability to scale systems for future demands are substantial.

Referências citadas

- ELI5: Why do GPU's need thousands of cores to get by nowadays, but CPU's can excel with just 8? - Reddit, acessado em maio 29, 2025, https://www.reddit.com/r/explainlikeimfive/comments/1ix6i2k/eli5_why_do_gpus_need_thousands_of_cores_to_get/
- GPU vs CPU - Difference Between Processing Units - AWS, acessado em maio 29, 2025, <https://aws.amazon.com/compare/the-difference-between-gpus-cpus/>
- GPU Parallel Computing: Techniques, Challenges, and Best Practices - Atlantic.Net, acessado em maio 29, 2025, <https://www.atlantic.net/gpu-server-hosting/gpu-parallel-computing-techniques-challenges-and-best-practices/>
- How to generate a vector of random numbers on a GPU | Towards Data Science, acessado em maio 29, 2025, <https://towardsdatascience.com/how-to-generate-a-vector-of-random-numbers-on-a-gpu-a37230f887a6/>
- Is it possible to build a system to generate UUIDs where every UUID is guaranteed unique?, acessado em maio 29, 2025, <https://softwareengineering.stackexchange.com/questions/450553/is-it-possible-to-build-a-system-to-generate-uuids-where-every-uuid-is-guarantee>
- Chapter 37. Efficient Random Number Generation and Application ..., acessado em maio 29, 2025, <https://developer.nvidia.com/gpugems/gpugems3/part-vi-gpu-computing/chapter-37-efficient-random-number-generation-and-application>
- Document IDs - Chroma Cookbook, acessado em maio 29, 2025, <https://cookbook.chromadb.dev/core/document-ids/>
- Six Options for Generating Distributed Unique IDs | Practical ..., acessado em maio 29, 2025, <https://peng.fyi/post/six-options-generating-distributed-ids/>
- Snowflake Algorithm: UUID Generation for Distributed Systems ..., acessado em maio 29, 2025, <https://dev.to/leapcell/snowflake-algorithm-uuid-generation-for-distributed-systems-4l56>
- UUID, ULID, Nanoids and Snowflake IDs , What's the difference? - Hewi's Blog, acessado em maio 29, 2025, <https://hewi.blog/uuid-ulid-nanoids-and-snowflake-ids-whats-the-difference>
- segmentio/ksuid: K-Sortable Globally Unique IDs - GitHub, acessado em maio 29, 2025, <https://github.com/segmentio/ksuid>
- ksuid v0.1.2 - HexDocs, acessado em maio 29, 2025, <https://hexdocs.pm/ksuid/Ksuid.html>
- UUID vs ULID vs Integer IDs: A Technical Guide for Modern Systems | ByteAether, acessado em maio 29, 2025, <https://byteaether.github.io/2025/uuid-vs-ulid-vs-integer-ids-a-technical-guide-for-modern-systems/>
- README.md | Hex Preview, acessado em maio 29, 2025, https://preview.hex.pm/preview/ecto_ulid_next/show/README.md
- PostgreSQL UUID Performance: Benchmarking Random (v4) and Time-based (v7) UUIDs, acessado em maio 29, 2025, <https://dev.to/umangsinha12/postgresql-uuid-performance-benchmarking-random-v4-and-time-based-v7-uuids-n9b>
- Benchmarking UUIDv4 vs UUIDv7 in PostgreSQL with 10 Million Rows - Reddit, acessado em maio 29, 2025,

- https://www.reddit.com/r/PostgreSQL/comments/1kta8ae/benchmarking_uuidv4_vs_uuidv7_in_postgresql_with/
17. CUDA: blocks and threads - Prof. Marco Bertini, acessado em maio 29, 2025, https://www.micc.unifi.it/bertini/download/gpu-programming-basics/2017/gpu_cuda_2.pdf
 18. Understanding Nvidia CUDA Cores: A Comprehensive Guide - Wevolver, acessado em maio 29, 2025, <https://www.wevolver.com/article/understanding-nvidia-cuda-cores-a-comprehensive-guide>
 19. Maxing out the device : r/CUDA - Reddit, acessado em maio 29, 2025, https://www.reddit.com/r/CUDA/comments/11sn3g6/maxing_out_the_device/
 20. Massively parallel programming with GPUs — Computational Statistics in Python 0.1 documentation - Duke People, acessado em maio 29, 2025, <https://people.duke.edu/~ccc14/sta-663/CUDAPython.html>
 21. 1. Introduction — CUDA C++ Programming Guide - NVIDIA Docs Hub, acessado em maio 29, 2025, <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#computational-grid>
 22. Thread block (CUDA programming) - Wikipedia, acessado em maio 29, 2025, [https://en.wikipedia.org/wiki/Thread_block_\(CUDA_programming\)](https://en.wikipedia.org/wiki/Thread_block_(CUDA_programming))
 23. Maximum number of threads on thread block - CUDA Programming and Performance, acessado em maio 29, 2025, <https://forums.developer.nvidia.com/t/maximum-number-of-threads-on-thread-block/46392>
 24. The Role of Warps in Parallel Processing: Optimizing GPU Performance for High-Speed Computing | DigitalOcean, acessado em maio 29, 2025, <https://www.digitalocean.com/community/tutorials/the-role-of-warps-in-parallel-processing>
 25. CUDA Thread Basics, acessado em maio 29, 2025, <http://users.wfu.edu/choss/CUDA/docs/Lecture%205.pdf>
 26. 1. Introduction — CUDA C++ Programming Guide - NVIDIA Docs Hub, acessado em maio 29, 2025, <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#built-in-variables>
 27. CUDA Refresher: The CUDA Programming Model | NVIDIA Technical Blog, acessado em maio 29, 2025, <https://developer.nvidia.com/blog/cuda-refresher-cuda-programming-model/>
 28. CUDA Architecture, acessado em maio 29, 2025, https://ajdillhoff.github.io/notes/cuda_architecture/
 29. What is the maximum number of concurrent streams that can be created in CUDA?, acessado em maio 29, 2025, <https://massedcompute.com/faq-answers/?question=What%20is%20the%20maximum%20number%20of%20concurrent%20streams%20that%20can%20be%20created%20in%20CUDA?>
 30. Collecting busy SM IDs - CUDA Programming and Performance - NVIDIA Developer Forums, acessado em maio 29, 2025, <https://forums.developer.nvidia.com/t/collecting-busy-sm-ids/285604>
 31. Modeling Utilization to Identify Shared-memory Atomic Bottlenecks - arXiv, acessado em maio 29, 2025, <https://arxiv.org/pdf/2503.17893>
 32. Analyzing atomic operation conflicts - Debugging and Optimizing CUDA Applications with NVIDIA Tools | StudyRaid, acessado em maio 29, 2025, <https://app.studyraid.com/en/read/11728/371481/analyzing-atomic-operation-conflicts>
 33. CUDA atomic operation performance in different scenarios - Stack Overflow, acessado em maio 29, 2025, <https://stackoverflow.com/questions/22367238/cuda-atomic-operation-performance-in-different-scenarios>
 34. Libcu++ monotonically increasing clock in host code - NVIDIA Developer Forums, acessado em maio 29, 2025, <https://forums.developer.nvidia.com/t/libcu-monotonically-increasing-clock-in-host-code/296143>
 35. How to Accurately Time CUDA Kernels in Pytorch - Speechmatics, acessado em maio 29, 2025, <https://www.speechmatics.com/company/articles-and-news/timing-operations-in-pytorch>
 36. Monotonic and Wall Clock Time - DEV Community, acessado em maio 29, 2025, <https://dev.to/amoabakelvin/monotonic-and-wall-clock-time-3mp>
 37. Maximize GPU Efficiency: Smarter Fixes for Checkpointing Challenges - Clockwork.io, acessado em maio 29, 2025, <https://www.clockwork.io/maximize-gpu-efficiency-smarter-fixes-for-checkpointing-challenges>
 38. create uuid to gpu ? · Issue #14 · kobalick/amdtweak - GitHub, acessado em maio 29, 2025, <https://github.com/kobalick/amdtweak/issues/14>
 39. Identifying data race patterns - Debugging and Optimizing CUDA Applications with NVIDIA Tools | StudyRaid, acessado em maio 29, 2025, <https://app.studyraid.com/en/read/11728/371480/identifying-data-race-patterns>
 40. Efficient GPU Parallel Implementation and Optimization of ARIA for Counter and Exhaustive Key-Search Modes - MDPI, acessado em maio 29, 2025, <https://www.mdpi.com/2079-9292/14/10/2021>
 41. Principles of Parallel Algorithm Design - Computer Science Purdue, acessado em maio 29, 2025, https://www.cs.purdue.edu/homes/ayg/CS525_SPR17/chap3_slides.pdf
 42. Design and Analysis of Parallel Algorithms, acessado em maio 29, 2025, <https://www.inf.ed.ac.uk/teaching/courses/dapa/overheads.pdf>
 43. Automating GPU Kernel Generation with DeepSeek-R1 and Inference Time Scaling, acessado em maio 29, 2025, <https://developer.nvidia.com/blog/automating-gpu-kernel-generation-with-deepseek-r1-and-inference-time-scaling/>
 44. Characterizing GPU Resilience and Impact on AI/HPC Systems - arXiv, acessado em maio 29, 2025, <https://arxiv.org/html/2503.11901v1>
 45. Speeding up Elixir: integration with native code (NIF, Ports, etc.) - DEV Community, acessado em maio 29, 2025, <https://dev.to/adamang/speeding-up-elixir-integration-with-native-code-nif-ports-etc-5ajd>
 46. Hok: Higher-Order GPU kernels in Elixir, acessado em maio 29, 2025, <https://sol.sbc.org.br/index.php/sblp/article/download/30259/30066/>
 47. ardubois/hok: Hok: Higher-Order GPU Kernels - GitHub, acessado em maio 29, 2025, <https://github.com/ardubois/hok>
 48. hpcaitech/Elixir: Elixir: Train a Large Language Model on a Small GPU Cluster - GitHub, acessado em maio 29, 2025, <https://github.com/hpcaitech/Elixir>

49. Ecto.Ulid Next - HexDocs, acessado em maio 29, 2025, https://hexdocs.pm/ecto_ulid_next/
50. twitter_snowflake — Rust implementation // Lib.rs, acessado em maio 29, 2025, https://lib.rs/crates/twitter_snowflake
51. RobThree/IdGen: Twitter Snowflake-alike ID generator for .Net - GitHub, acessado em maio 29, 2025, <https://github.com/RobThree/IdGen>
52. Frostflake - a Snowflake inspired high-performance unique identifier generator, acessado em maio 29, 2025, <https://forums.swift.org/t/frostflake-a-snowflake-inspired-high-performance-unique-identifier-generator/61503>