# A Survey of Non-Cryptographic Hash Functions
## As Implemented in SMHasher3

Thomas Dybdahl Ahle

February 3, 2026

**Abstract**

This document provides a comprehensive survey of non-cryptographic hash functions as implemented in the SMHasher3 testing framework [Woj22]. We present each algorithm using clear pseudocode and mathematical notation, analyze their mixing properties, and discuss their relative strengths and weaknesses.

# Contents

# 1    Introduction

Non-cryptographic hash functions are fundamental building blocks in computer science, used extensively in hash tables, checksums, and probabilistic data structures. This survey documents hash functions implemented in SMHasher3, presenting algorithms with clear pseudocode and analyzing their mixing properties.

## 1.1    Performance Overview

Table 1 shows performance measurements for the hash functions surveyed in this paper, from the SMHasher3 test suite, sorted by bulk throughput. Hashes marked with <span style="color:red">X</span> do not pass the full test battery; the number in parentheses is the count of tests passed out of 250.

| Hash | Bits | B/cycle (bulk) | Cycles (1–32 B) | Pass |
|------|------|----------------|------------------|------|
| AquaHash* | 64 | 15.99 | 34.31 | <span style="color:red">X</span> (195) |
| XXH3* | 64 | 12.95 | 30.69 | <span style="color:red">X</span> (223) |
| MeowHash | 128 | 12.43 | 61.51 | ✓ |
| rapidhash | 64 | 8.82 | 30.41 | ✓ |
| CLHash* | 64 | 7.38 | 48.19 | <span style="color:red">X</span> (64) |
| MuseAir | 64 | 7.19 | 28.45 | ✓ |
| wyhash* | 64 | 7.11 | 29.31 | <span style="color:red">X</span> (235) |
| HalftimeHash* | 128 | 6.92 | 95.58 | <span style="color:red">X</span> (202) |
| komihash | 64 | 6.24 | 32.62 | ✓ |
| ChibiHash | 64 | 6.15 | 34.53 | ✓ |
| UMASH* | 64 | 6.09 | 41.74 | <span style="color:red">X</span> (122) |
| aHash | 64 | 4.52 | 29.14 | ✓ |
| SpookyHash | 32 | 4.40 | 60.20 | ✓ |
| Polymur | 64 | 4.02 | 42.57 | ✓ |
| xxHash64* | 64 | 4.00 | 52.61 | <span style="color:red">X</span> (241) |
| khashv | 64 | 3.20 | 56.17 | ✓ |
| a5hash | 64 | 2.63 | 26.91 | ✓ |
| FarmHash | 128 | 2.61 | 69.61 | ✓ |
| xxHash32* | 32 | 2.00 | 44.38 | <span style="color:red">X</span> (167) |
| Rainbow | 64 | 1.28 | 61.63 | ✓ |
| prvhash | 64 | 0.94 | 51.07 | ✓ |
| XMSX | 32 | 0.67 | 42.56 | ✓ |
| poly-mersenne | 32 | 0.50 | 69.79 | ✓ |
| HalfSipHash | 32 | 0.34 | 81.75 | ✓ |
| SipHash-2-4 | 64 | 0.32 | 152.26 | ✓ |

Table 1: Performance sorted by bulk throughput (bytes/cycle). *Does not pass the full SMHasher3 test battery; see individual sections for details.

**Two performance regimes.**    Hash performance differs dramatically between short and long keys:

- **Short keys (1–32 bytes):** Latency-bound. Specialized code paths avoid loop overhead. The best hashes achieve ~26–31 cycles total.

- **Bulk throughput:** Bandwidth-bound. Parallel lanes and wide reads dominate. The best MUM-based hashes achieve 7–9 bytes/cycle; AES-based hashes (MeowHash) reach ~12 bytes/cycle.

| Hash | Output | Block | Lanes | Primitive | Keying | Provable width |
|---|---|---|---|---|---|---|
| a5hash | 64 | 16 | 1 | MUM | seed | N/A |
| a5hash-128 | 128 | 32/64 | 2/4 | MUM | seed | N/A |
| MuseAir | 64 | 96 | 6 | MUM | seed | N/A |
| wyhash | 64 | 48 | 3 | MUM | seed+secret | N/A |
| rapidhash | 64 | 112 | 7 | MUM | seed+secret | N/A |
| rapidhashMicro | 64 | 80 | 5 | MUM | seed+secret | N/A |
| rapidhashNano | 64 | 48 | 3 | MUM | seed+secret | N/A |
| xxHash32 | 32 | 16 | 4 | ARX | seed | N/A |
| xxHash64 | 64 | 32 | 4 | ARX | seed | N/A |
| XXH3-64 | 64 | 64 | 8 | MUM | seed+secret | N/A |
| SipHash-2-4 | 64 | 8 | 1 | ARX | 128-bit key | N/A |
| komihash | 64 | 64 | 4 | MUM | seed | N/A |
| ChibiHash | 64 | 32 | 4 | mult+rot | seed | N/A |
| CLHash | 64 | 1024 | 8 (ILP) | CLMUL | random key | $\approx 64 - \log_2 n$ |
| UMASH | 64/128 | 256 | 4 (ILP) | CLMUL+poly | random key | $\approx 61$ |
| Polymur | 64 | 49 | 1 | poly($2^{61}-1$) | derived key | $\approx 61$ |
| AquaHash | 128 | 64 | 4 | AES | 128-bit seed | N/A |
| MeowHash | 128 | 256 | 8 | AES | 128-bit seed | N/A |
| HalftimeHash | 64 | var | var | 32-bit mul | random key | N/A |
| SpookyHash | 128 | 96 | 12 | ARX | seed | N/A |
| FarmHash | 64 | 64 | 7 | ARX+MUM | seed | N/A |
| Rainbow | 64 | 16 | 4 | mult+rot | seed | N/A |
| prvhash | 64 | 8 | 1 | PRNG | seed | N/A |
| XMSX | 32 | 4 | 1 | mult+shift | seed | N/A |
| khashv | 64 | 16 | 4 | S-box+ARX | seed | N/A |
| aHash | 64 | 16 | 2 | AES | random key | N/A |
| HalfSipHash | 32 | 4 | 1 | ARX | 64-bit key | N/A |
| poly-mersenne | 32 | 8 | 1 | poly($2^{61}-1$) | random key | $\approx 61$ |

Table 2: Structural comparison of hash functions. "Lanes" indicates parallel accumulators; "(ILP)" denotes instruction-level parallelism without explicit state lanes. "Keying" shows the security model: seed-only hashes are vulnerable to seed-independent attacks; secret/key-based hashes resist them if the key is hidden. "Provable width" is the effective bits $w = -\log_2(\epsilon)$ for a published $\epsilon$-almost-universal collision bound under a random secret key (often length-dependent, e.g. $\epsilon \approx n/2^{64}$).

## 2 Comparison

**Design trade-offs:**

- **MUM-based** (a5hash, MuseAir, wyhash, rapidhash, XXH3, komihash): Fast on modern CPUs with efficient 64×64→128 multiply. Excellent throughput for long keys.

- **ARX-based** (xxHash32/64, SipHash, SpookyHash, FarmHash): Portable, constant-time friendly. SipHash sacrifices speed for cryptographic security margin. SpookyHash uses 12 lanes for high throughput.

- **Parallel lanes**: More lanes = higher throughput but more state to maintain. rapidhash (7 lanes) vs wyhash (3 lanes) trades complexity for speed. MuseAir uses 6 lanes with chained dependencies.

- **Polynomial** (Polymur, UMASH): Theoretically grounded with provable collision bounds. UMASH combines CLMUL compression with polynomial finalization.

- **AES-based** (AquaHash, MeowHash, aHash): Highest throughput on supporting hardware via AES-NI instructions. aHash is the default in Rust's HashMap.

- **Keying**: Hashes with secret keys (SipHash, CLHash, UMASH) resist HashDoS attacks if the key is kept hidden. Seed-only hashes can be attacked with seed-independent collisions.

# 3 Preliminaries

## 3.1 Notation

| Symbol | Meaning |
|--------|---------|
| $\oplus$ | Bitwise XOR |
| $\cdot$ | Multiplication modulo $2^w$ on uw words (keep the low $w$ bits) unless noted |
| $\text{ROT}_k(x)$ | Left rotation of $x$ by $k$ bits |
| $\gg$ | Logical right shift |
| $\ll$ | Left shift |
| $\text{Lo}(x)$ | Low 64 bits of 128-bit value $x$ |
| $\text{Hi}(x)$ | High 64 bits of 128-bit value $x$ |
| $\text{Lo32}(x)$ | Low 32 bits of 64-bit value $x$ |
| $\text{Hi32}(x)$ | High 32 bits of 64-bit value $x$ |
| $p$ | Pointer to input bytes |
| $n$ | Input length in bytes |
| $seed$ | User-provided seed (may be secret for keyed hashes) |
| $k_i$ | Secret key words (e.g. SipHash) |
| $s_i$ | Algorithm-specific fixed constants / secret tables |
| $v_i$ | Evolving internal state words (accumulators/lanes) |
| $\text{READ}_w(p, i)$ | Read $w$-bit word at byte offset $i$ from $p$ |
| $p[i]$ | Byte at offset $i$ (equivalent to $\text{READ}_8(p, i)$) |
| $m_i$ | The $i$-th block of input (block size $B$ varies by hash) |
| $x_j$ | The $j$-th word within current block (64-bit unless noted) |
| uw | Unsigned $w$-bit word type (e.g. u32, u64) |

**Terminology.** A **block** is the fixed-size unit processed by the main loop (e.g., 48 bytes for wyhash, 112 bytes for rapidhash). The term **stripe** is XXH3-specific, denoting a 64-byte segment processed across 8 parallel lanes. We use "block" consistently; some hash authors use "chunk" synonymously.

## 3.2 Primitive Categories

Modern hash functions are built from a small set of mixing primitives:

**MUM (Multiply-Mix)** Wide multiplication followed by XOR folding:

$$\text{Mum}(a, b) = \text{Lo}(a \times b) \oplus \text{Hi}(a \times b) \tag{1}$$

Fast on modern CPUs with efficient 64×64→128 multiply units. Used by wyhash, rapidhash, XXH3, komihash, t1ha.

**ARX (Add-Rotate-XOR)** Combines addition, rotation, and XOR in patterns like:

$$h \leftarrow \text{ROT}_r(h + x) \oplus y \quad \text{or} \quad h \leftarrow \text{ROT}_r(h \oplus x) + y \tag{2}$$

Portable and constant-time friendly. Used by MurmurHash, xxHash32/64, SipHash, SpookyHash, lookup3.

**CLMUL (Carry-less Multiply)** Polynomial multiplication in $GF(2^n)$:

$$a \otimes b = a \cdot b \mod p(x) \tag{3}$$

where $\cdot$ is polynomial multiplication (XOR for addition, AND for coefficient multiplication). Theoretically grounded with provable collision bounds. Used by CLHash, UMASH. Requires PCLMULQDQ instruction.

**Polynomial over** $GF(p)$ Evaluates polynomial over a prime field:

$$h \leftarrow (h \cdot k + m_i) \mod p \quad \text{typically } p = 2^{61} - 1 \tag{4}$$

Strong theoretical guarantees from algebra. Used by UMASH, Polymur.

**AES-NI** Uses AES encryption rounds for mixing:

$$\text{AESENC}(state, key) = \text{MixColumns}(\text{ShiftRows}(\text{SubBytes}(state))) \oplus key \tag{5}$$

Excellent diffusion (full avalanche) in a single instruction. Used by AquaHash, MeowHash.

**Simple Multiply** Basic multiply-accumulate:

$$h \leftarrow (h \oplus b) \cdot p \quad \text{(FNV-1a)} \qquad h \leftarrow h \cdot 33 + b \quad \text{(DJB2)} \tag{6}$$

Minimal code size but weaker mixing. One byte per iteration.

Most high-performance hashes use MUM or ARX. MUM-based hashes are generally faster on x86-64; ARX-based hashes are more portable. CLMUL and AES-based hashes offer the highest throughput on supported hardware.

## 3.3 Common Hash Structure

Most non-cryptographic hash functions follow a similar high-level structure:

**Initialization** Set up state variables (accumulators) from seed and/or constants. Multiple accumulators enable parallel processing.

**Main Loop** Process input in fixed-size blocks. Each block updates the state using the hash's core mixing primitive. Multiple *parallel lanes*—independent accumulators updated simultaneously—allow modern CPUs to exploit instruction-level parallelism.

**Tail Handling** Process remaining bytes that don't fill a complete block. Often uses different logic than the main loop (overlapping reads, byte-by-byte, etc.).

**Lane Collapse** Combine multiple parallel lanes into a single value, typically via XOR trees, addition chains, or mixing functions.

**Finalization** Apply an avalanche finalizer (xorshift and multiplication) to ensure all input bits affect all output bits (avalanche property).

The block size and number of lanes vary by hash: wyhash uses 48-byte blocks with 3 lanes, rapidhash uses 112-byte blocks with 7 lanes, and MuseAir uses 96-byte blocks with 6 lanes.

## 3.4 Finalizers and Mixers

Most fast non-cryptographic hashes end with a short *avalanche* stage that destroys residual structure from the main loop. These finalizers are usually built from two permutations on $w$-bit words:

- **Xorshift**: $\text{XORSHIFT}_r(x) = x \oplus (x \gg r)$ is invertible for $0 < r < w$.

- **Multiply by an odd constant**: $x \leftarrow x \cdot c$ is invertible mod $2^w$ when $c$ is odd.

The common template is:

$$h \leftarrow \text{XORSHIFT}_{r_1}(h) \tag{7}$$
$$h \leftarrow h \cdot c_1 \tag{8}$$
$$h \leftarrow \text{XORSHIFT}_{r_2}(h) \tag{9}$$
$$h \leftarrow h \cdot c_2 \tag{10}$$
$$h \leftarrow \text{XORSHIFT}_{r_3}(h) \tag{11}$$

Different hashes choose $(r_i)$ and $(c_i)$ empirically (e.g. XXH32_AVALANCHE, XXH64_AVALANCHE, the XXH3 finalizer, and Polymur's mx3 finalizer).

**Common sub-steps.** It is convenient to name two ubiquitous building blocks:

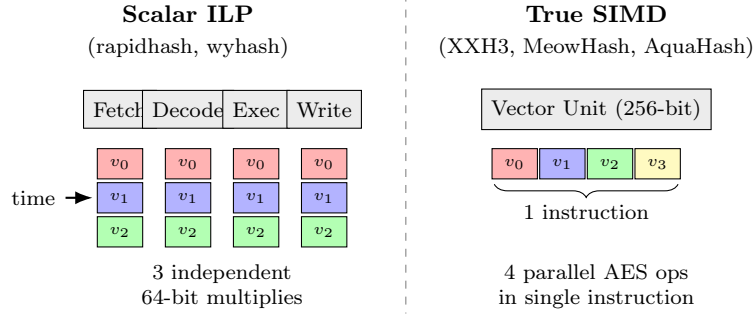$$\text{XMS}_{r,c}(x) = (x \oplus (x \gg r)) \cdot c \qquad \text{(xorshift-then-multiply)} \tag{12}$$
$$\text{SHIFTMIX}(x) = x \oplus (x \gg 47) \tag{13}$$

Many "Murmur-style" combiners (e.g. FarmHash's HASH128TO64 and HASHLEN16) use SHIFT-MIX between multiplications; xxHash uses XMS with two different constants.

**Pseudocode conventions.** Main loop processing is shown in full detail. Tail handling (which often involves complex case analysis for 1–15 remaining bytes) is sometimes summarized as "handle remaining bytes" to keep pseudocode readable. The core mixing primitives are always fully specified. Comments use % for section boundaries (Bulk/Tail/Finalization) where helpful.

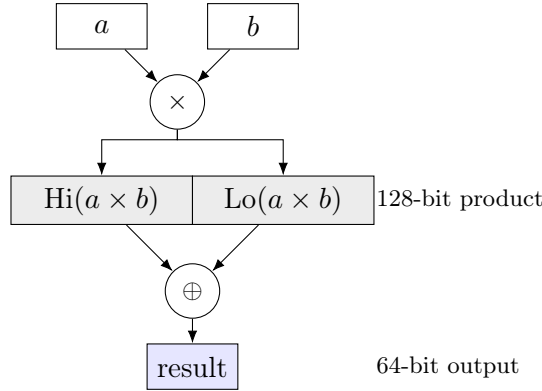**Scalar ILP vs. SIMD lanes.** Hash "lanes" come in two flavors:

- **Scalar ILP lanes** (wyhash, rapidhash, SpookyHash): Independent scalar accumulators that exploit the CPU's out-of-order execution engine. The 64×64→128-bit multiply used by MUM *has no SIMD instruction*—SSE/AVX only provide 32×32→64. These hashes use 3–7 independent dependency chains to keep the scalar multiply unit fed. Lane counts like 3 or 7 are tuned for register pressure and multiply latency, not SIMD widths.

- **True SIMD lanes** (XXH3, AquaHash, MeowHash, HalftimeHash): Vector operations using SSE/AVX/NEON intrinsics with operations that *do* vectorize: AES-NI instructions (`aesenc`/`aesdec`) or 32×32→64 multiplies (`_mm_mul_epu32`). XXH3 uses 32×32 multiplies for its long-key path (explicit AVX2/SSE2/NEON implementations), but falls back to scalar 64×64 MUM for short/medium keys where SIMD overhead isn't worth it.

**Scalar ILP**
(rapidhash, wyhash)

**True SIMD**
(XXH3, MeowHash, AquaHash)

Scalar ILP side: Fetch | Decode | Exec | Write; with columns of $v_0$, $v_1$, $v_2$ over time. **3 independent 64-bit multiplies**

True SIMD side: Vector Unit (256-bit); $v_0$, $v_1$, $v_2$, $v_3$ — 1 instruction. **4 parallel AES ops in single instruction**

# 4 Round Functions and Combiners

This section details specific mixing functions used by hash families. For the core primitives (MUM, ARX, etc.), see Section 3.2.

## 4.1 The MUM Primitive



We define three variants used throughout this survey:

**MUM$(a, b)$ Standard (folded):** Returns $\text{Lo}(a \times b) \oplus \text{Hi}(a \times b)$. Used by wyhash (`_wymix`), rapidhash (`rapid_mix`), XXH3.

**MUM$^*(a, b)$ Protected (folded):** XORs the original inputs back into the product before folding, preventing degenerate cases (e.g. $a=0$) from collapsing state:

$$\text{MUM}^*(a, b) = a \oplus b \oplus \text{Lo}(a \times b) \oplus \text{Hi}(a \times b) \tag{14}$$

Used by wyhash-strict, rapidhash-protected, MuseAir.

**MUM$_{128}(a, b)$ Full product:** Returns the pair $(\text{Lo}(a \times b), \text{Hi}(a \times b))$ without folding. Used in finalizations where both halves update separate state variables.

**MUM$^*_{128}(a, b)$ Protected full product:** Returns the pair $(a \oplus \text{Lo}(a \times b), b \oplus \text{Hi}(a \times b))$. Folding this pair with XOR gives MUM$^*(a, b)$.

## 4.2 Round Functions

Many hashes use a *round function* that mixes input into an accumulator. The xxHash round functions are defined in Section 7.1.1.

## 4.3 Xorshift

The xorshift primitive XORs a value with a shifted copy of itself:

$$\text{XORSHIFT}_r(v) = v \oplus (v \gg r) \tag{15}$$

**Permutation proof.** $\text{XORSHIFT}_r$ is a bijection (permutation) for any $0 < r < w$ where $w$ is the word size. To see this, note that the top $r$ bits of $v$ are unchanged, and each subsequent bit can be recovered by XORing with the already-recovered bit $r$ positions above it. Equivalently, applying $\text{XORSHIFT}_r$ a total of $\lceil w/r \rceil$ times yields the identity.

## 4.4 ShiftMix

A specific xorshift used by CityHash and FarmHash:

$$\text{SHIFTMIX}(v) = \text{XORSHIFT}_{47}(v) = v \oplus (v \gg 47) \tag{16}$$

This diffuses high bits into the low 17 bits of a 64-bit value.

## 4.5 HashLen16 (CityHash)

CityHash's core 128-to-64 bit reducer, inspired by MurmurHash:

$$a \leftarrow (u \oplus v) \cdot mul \tag{17}$$
$$a \leftarrow \text{SHIFTMIX}(a) \tag{18}$$
$$b \leftarrow (v \oplus a) \cdot mul \tag{19}$$
$$b \leftarrow \text{SHIFTMIX}(b) \tag{20}$$
$$\text{HASHLEN16}(u, v, mul) \leftarrow b \cdot mul \tag{21}$$

where $mul$ is an odd 64-bit multiplier (often a fixed constant, or $k_2 + 2n$ for length-dependent mixing). When the multiplier is the fixed default, we may omit it and write $\text{HASHLEN16}(u, v)$. This resembles $\text{MUM}$ but applies the shifts between multiplications rather than using 128-bit product folding.

## 4.6 Finalization Mixers

All finalization mixers follow a common pattern: alternating xorshift and multiply operations. The general form is:

$$h \leftarrow \text{XORSHIFT}_{r_1}(h) \tag{22}$$
$$h \leftarrow h \cdot c_1 \tag{23}$$
$$h \leftarrow \text{XORSHIFT}_{r_2}(h) \tag{24}$$
$$h \leftarrow h \cdot c_2 \tag{25}$$
$$\text{FMIX}(h) \leftarrow \text{XORSHIFT}_{r_3}(h) \tag{26}$$

Different hashes choose different shift amounts $(r_1, r_2, r_3)$ and multiplier constants $(c_1, c_2)$ based on avalanche analysis. Typical 32-bit shifts: $(16, 13, 16)$. Typical 64-bit shifts: $(33, 33, 33)$ or $(33, 29, 32)$.

This is a bijection: each $\text{XORSHIFT}_r$ is a permutation (see Section 4.3), and multiplication by an odd constant is a permutation mod $2^w$.

# 5 The wyhash/rapidhash Family

This family uses the MUM primitive as its core operation:

$$\text{Mum}(a,\, b) = \text{Lo}(a \times b) \oplus \text{Hi}(a \times b) \tag{27}$$

## 5.1 Secrets

Both families use 64-bit "secret" constants $s_0, s_1, \ldots$ for mixing. wyhash uses 4 secrets; rapidhash uses 8 (the first 4 shared with wyhash).

## 5.2 wyhash (64-bit)

Source: `hashes/wyhash.cpp` [Yi19].

**SMHasher3 test results.** wyhash fails 15 of 250 tests: Permutation tests with single-bit keys at all block sizes (4–16 bytes, both LE and BE), and SeedZeroes at 1280 and 8448 bytes. The Permutation failures reflect insufficient diffusion when only one bit of each input word is set.

The main loop uses 3 scalar lanes (not 4 or 8) because Mum relies on $64 \times 64 \rightarrow 128$-bit scalar multiplication, which doesn't vectorize in SSE/AVX. The 3 independent dependency chains exploit ILP while keeping register pressure manageable.

1: $v_0 \leftarrow seed \oplus \text{Mum}(seed \oplus s_0,\ s_1)$
2: **for** $i \in \{1, 2\}$ **do**
3:     $v_i \leftarrow v_0$
4: **end for**
5: **for** each block $(x_0, \ldots, x_5) \in \text{u64}^6$ **do**                $\triangleright$ Main loop: 3 parallel lanes
6:     **for** $i \in \{0, 1, 2\}$ **do**
7:         $v_i \leftarrow \text{Mum}(x_{2i} \oplus s_{i+1},\ x_{2i+1} \oplus v_i)$
8:     **end for**
9: **end for**
10: $v_0 \leftarrow v_0 \oplus v_1 \oplus v_2$                                          $\triangleright$ Collapse lanes
11: Process remaining bytes; read last 16 bytes as $t_0, t_1$
12: $t_0 \leftarrow t_0 \oplus s_1, \quad t_1 \leftarrow t_1 \oplus v_0$
13: $(t_0, t_1) \leftarrow (\text{Lo}(t_0 \times t_1), \text{Hi}(t_0 \times t_1))$
14: **return** $\text{Mum}(t_0 \oplus s_0 \oplus n,\ t_1 \oplus s_1)$

## 5.3 rapidhash

Source: `hashes/rapidhash.cpp`, version 3 [De 24].

By Nicolas De Carli, based on wyhash. Uses 7 scalar ILP lanes (112-byte blocks), trading higher register pressure for more parallelism on wide superscalar CPUs. Uses $s_2$ for seed initialization and 8 secrets ($s_0$–$s_7$).

The "protected" variant XORs the original inputs back into the MUM result:

$$\text{Mum}^*_{128}(a, b) = (a \oplus \text{Lo}(a \times b),\ b \oplus \text{Hi}(a \times b)) \tag{28}$$

This prevents zero-multiplication weaknesses at a small performance cost.

1: $v_0 \leftarrow seed \oplus \text{Mum}(seed \oplus s_2,\ s_1)$                       $\triangleright$ Note: $s_2$ not $s_0$
2: **for** $i \in \{1, \ldots, 6\}$ **do**
3:     $v_i \leftarrow v_0$
4: **end for**
5: **for** each block $(x_0, \ldots, x_{13}) \in \text{u64}^{14}$ **do**            $\triangleright$ Main loop: 7 parallel lanes
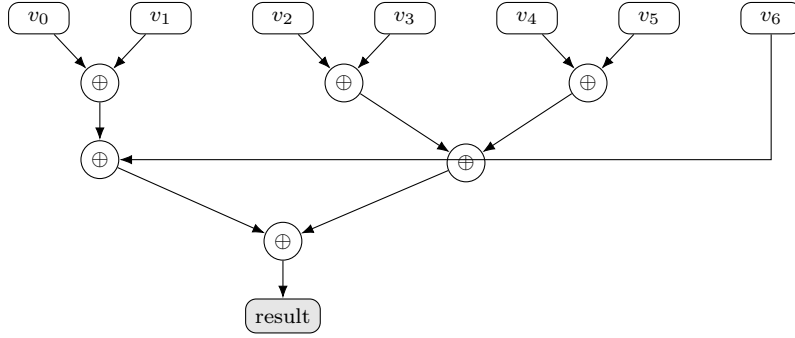6:     **for** $i \in \{0, \ldots, 6\}$ **do**

7:         $v_i \leftarrow \text{MUM}(x_{2i} \oplus s_i,\ x_{2i+1} \oplus v_i)$

8:    **end for**

9: **end for**

10: $v_0 \leftarrow v_0 \oplus v_1$; $v_2 \leftarrow v_2 \oplus v_3$; $v_4 \leftarrow v_4 \oplus v_5$         ▷ Collapse lanes

11: $v_0 \leftarrow v_0 \oplus v_6$; $v_2 \leftarrow v_2 \oplus v_4$; $v_0 \leftarrow v_0 \oplus v_2$

12: Process remaining bytes; read last 16 bytes as $t_0, t_1$ with $t_0 \leftarrow t_0 \oplus n$

13: $t_0 \leftarrow t_0 \oplus s_1, \quad t_1 \leftarrow t_1 \oplus v_0$

14: $(t_0, t_1) \leftarrow (\text{LO}(t_0 \times t_1), \text{HI}(t_0 \times t_1))$

15: **return** $\text{MUM}(t_0 \oplus s_7,\ t_1 \oplus s_1 \oplus n)$         ▷ Note: $s_7$ and extra $n$

**XOR-tree lane combination:**



**Tail (17–112 bytes).** Nested conditionals process 16 bytes at a time with alternating secrets:

| Bytes remaining | Offset | Secret |
|:---:|:---:|:---:|
| > 16 | 0, 8 | $s_2$ |
| > 32 | 16, 24 | $s_2$ |
| > 48 | 32, 40 | $s_1$ |
| > 64 | 48, 56 | $s_1$ |
| > 80 | 64, 72 | $s_2$ |
| > 96 | 80, 88 | $s_1$ |

**Finalization.** Read last 16 bytes as $(t_0, t_1)$ with $t_0 \leftarrow t_0 \oplus n$, then apply $t_0 \leftarrow t_0 \oplus s_1$, $t_1 \leftarrow t_1 \oplus v_0$, and $(t_0, t_1) \leftarrow (\text{LO}(t_0 \times t_1), \text{HI}(t_0 \times t_1))$. Return $\text{MUM}(t_0 \oplus s_7,\ t_1 \oplus s_1 \oplus n)$. Differs from wyhash: length is XORed into $a$ before MUM (not the final result), uses $s_7$ instead of $s_0$, and includes length twice in finalization.

## 5.4 rapidhash Variants

| Variant | Chunk size | Lanes | Secrets used |
|:---|:---:|:---:|:---:|
| rapidhash | 112 bytes | 7 | $s_0$–$s_7$ |
| rapidhashMicro | 80 bytes | 5 | $s_0$–$s_4$, $s_7$ |
| rapidhashNano | 48 bytes | 3 | $s_0$–$s_2$, $s_7$ |

**Rust port.** The `rust-rapidhash` implementation uses a different seed initialization: secrets are regenerated from the user seed using the protected MUM variant, ensuring non-zero quality in specific bit ranges. This makes hashes incompatible between the C and Rust versions for the same seed.

# 6 The a5hash/MuseAir Family

These hashes also use the MUM primitive but with different architectural choices than wyhash/rapidhash.

## 6.1 a5hash

Source: `hashes/a5hash.cpp`, version 5.21 [Van25].

A MUM-based hash by Aleksey Vaneev with a single MUM lane in the main loop (simpler than wyhash's 3 lanes). Uses mantissa bits of $\pi$ for initial state and checkerboard masks $c_0 =$ `0x5555...`, $c_1 =$ `0xAAAA...` that ensure different bit positions are treated differently. Addition of masks after each round prevents zero-multiplication weaknesses. Available in 32-bit, 64-bit, and 128-bit variants.

1: $v_0 \leftarrow$ `0x243F6A8885A308D3` $\oplus n$, $v_1 \leftarrow$ `0x452821E638D01377` $\oplus n$      $\triangleright$ From $\pi$
2: $(v_0, v_1) \leftarrow \text{MUM}_{128}(v_1 \oplus (seed \wedge c_1), v_0 \oplus (seed \wedge c_0))$
3: **if** $n > 16$ **then**
4:      $c_0 \leftarrow c_0 \oplus v_0$; $c_1 \leftarrow c_1 \oplus v_1$
5:      **for** each block $(x_0, x_1) \in$ u64$^2$ excluding final 16 bytes **do**      $\triangleright$ Main loop
6:          $(v_0, v_1) \leftarrow \text{MUM}_{128}(x_0 \oplus v_0, x_1 \oplus v_1)$
7:          $v_0 \leftarrow v_0 + c_0$; $v_1 \leftarrow v_1 + c_1$
8:      **end for**
9: **end if**
10: Handle tail (0–16 bytes) with overlapping reads
11: $(v_0, v_1) \leftarrow \text{MUM}_{128}(v_0, v_1)$
12: $(v_0, v_1) \leftarrow \text{MUM}_{128}(v_0 \oplus c_0, v_1)$
13: **return** $v_0 \oplus v_1$

Here $\text{MUM}_{128}(a, b)$ returns both halves of the 128-bit product: $(\text{LO}(a \times b), \text{HI}(a \times b))$.

## 6.2 MuseAir

Source: `hashes/museair.cpp`, version 0.3 [K-A24].

A MUM-based hash by K–Aethiax with 6-lane parallelism and chained dependencies between lanes. Uses 7 constants from Ai(0) mantissa ($s_0 =$ `0x5ae31e589c56e17a`, ...). The MUMIX primitive either XORs MUM output into state (standard) or replaces it (bfast). Cross-lane MUM pattern in the epilogue ensures full mixing. Available in standard, bfast (faster, slightly weaker), and 128-bit variants.

**MUMIX.**

1: **function** MUMIX($v_p, v_q, x_p, x_q$)      $\triangleright$ State words $(v_p, v_q)$, input words $(x_p, x_q)$
2:      $a \leftarrow v_p \oplus x_p$;    $b \leftarrow v_q \oplus x_q$
3:      $(\ell, h) \leftarrow \text{MUM}_{128}(a, b)$
4:      $(v_p, v_q) \leftarrow (a \oplus \ell, b \oplus h)$      $\triangleright$ Standard variant
5:           $\triangleright$ bfast variant: $(v_p, v_q) \leftarrow (\ell, h)$
6: **end function**

**Long Keys.**

1: Initialize 6 state variables from seed:
2:      $v_0 \leftarrow s_0 + seed$, $v_1 \leftarrow s_1 - seed$, $v_2 \leftarrow s_2 \oplus seed$
3:      $v_3 \leftarrow s_3 + seed$, $v_4 \leftarrow s_4 - seed$, $v_5 \leftarrow s_5 \oplus seed$
4: $t \leftarrow s_6$      $\triangleright$ Initial chain value
5: **for** each block $(x_0, \ldots, x_{11}) \in$ u64$^{12}$ **do**      $\triangleright$ Main loop: 6 parallel lanes
6:      **for** $i \in \{0, \ldots, 5\}$ **do**

7:     $v_i \leftarrow v_i \oplus x_{2i}$

8:     $v_{(i+1) \bmod 6} \leftarrow v_{(i+1) \bmod 6} \oplus x_{2i+1}$

9:     $(lo, hi) \leftarrow \text{MUM}_{128}(v_i, v_{(i+1) \bmod 6})$

10:    $v_i \leftarrow v_i + (t \oplus hi)$                     ▷ Chain: previous lo XOR current hi

11:    $t \leftarrow lo$

12:  **end for**

13: **end for**

14: $v_0 \leftarrow v_0 + t$                     ▷ Fold in last chain value

15: Process remaining 48, 32, 16 bytes with MUMIX

16: Apply MUMIX to the last 16 bytes (tail summarized)

17: $t_0 \leftarrow v_0 - v_1$, $t_1 \leftarrow v_2 - v_3$, $t_2 \leftarrow v_4 - v_5$                     ▷ Finalization: collapse lanes

18: $t_0 \leftarrow \text{ROT}_{n \bmod 64}(t_0)$, $t_1 \leftarrow \text{ROT}_{64-(n \bmod 64)}(t_1)$, $t_2 \leftarrow t_2 \oplus n$  ▷ $\text{ROT}_{64-r}()$ is right-rotate by $r$

19: $(t'_0, t'_1) \leftarrow \text{MUM}_{128}(t_0, t_1)$, $(t'_2, t'_3) \leftarrow \text{MUM}_{128}(t_1, t_2)$, $(t'_4, t'_5) \leftarrow \text{MUM}_{128}(t_2, t_0)$

20: $t_0 \leftarrow t'_0 \oplus t'_5$, $t_1 \leftarrow t'_2 \oplus t'_1$, $t_2 \leftarrow t'_4 \oplus t'_3$                     ▷ Cross lo/hi

21: $(u_0, u_1) \leftarrow \text{MUM}_{128}(t_0, t_1)$, $(u_2, u_3) \leftarrow \text{MUM}_{128}(t_1, t_2)$, $(u_4, u_5) \leftarrow \text{MUM}_{128}(t_2, t_0)$

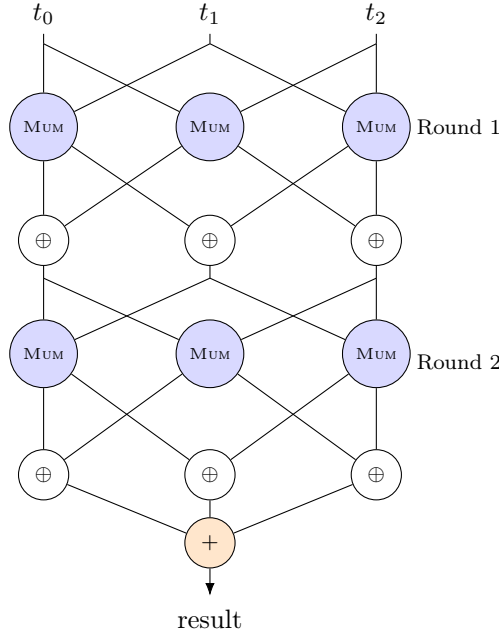22: **return** $(u_0 \oplus u_5) + (u_2 \oplus u_1) + (u_4 \oplus u_3)$                     ▷ 64-bit



Figure 1: MuseAir finalizer. Each round performs three $\text{MUM}_{128}$ operations on all pairs of its three inputs, then cross-XORs the lo/hi halves. After two rounds the three values are summed for the 64-bit result.

# 7   xxHash Family

Source: `hashes/xxhash.cpp` [Col12].

xxHash uses prime constants $s_1, \ldots, s_5$ selected by empirical testing for avalanche properties. The primes have roughly half their bits set and avoid degenerate bit patterns. For xxHash64:

$$s_1 = \text{0x9E3779B185EBCA87} \qquad s_2 = \text{0xC2B2AE3D27D4EB4F}$$
$$s_3 = \text{0x165667B19E3779F9} \qquad s_4 = \text{0x85EBCA77C2B2AE63}$$
$$s_5 = \text{0x27D4EB2F165667C5}$$

The 32-bit versions use corresponding 32-bit primes.

## 7.1 Round Functions

### 7.1.1 Multiply-Rotate Round

The xxHash round function mixes input into an accumulator:

$$\text{ROUND}(v, x) = \text{ROT}_r(v + x \cdot s_2) \cdot s_1 \tag{29}$$

where $r = 13$ for 32-bit, $r = 31$ for 64-bit. This is an ARX variant using multiplication instead of XOR for the final mix.

### 7.1.2 Merge Round (xxHash64)

After the main loop, xxHash64 merges each accumulator into the hash:

$$\text{MERGE}(h, v) = (h \oplus \text{ROUND}(0, v)) \cdot s_1 + s_4 \tag{30}$$

This "rounds" the accumulator value, XORs it into the hash, then applies another multiply-add for diffusion.

**SMHasher3 test results.** None of the xxHash variants pass the full SMHasher3 test battery. xxHash32 fails 83 of 250 tests, including numerous Sparse, Permutation, and seed-sensitivity tests (SeedBlockLen, SeedBlockOffset, SeedBIC, SeedBitflip), reflecting the limited mixing of its 32-bit multiply-rotate round. xxHash64 fares much better, failing only 9 tests, all seed-related (SeedBlockLen, SeedBlockOffset, SeedBIC). XXH3 (64-bit) fails 27 of 250 tests: BIC, Sparse, PerlinNoise, and Bitflip tests expose weaknesses in the non-seeded path, while SeedZeroes produces 571 full 64-bit collisions (expected 0) on keys up to 1280 bytes with low-weight seeds— a consequence of XXH3's seeding mechanism, which adds $\pm seed$ to the secret and thus maps $seed = 0$ and any seed that leaves all relevant secret words unchanged to identical output.

## 7.2 xxHash32

1:   $v_0 \leftarrow seed + s_1 + s_2$
2:   $v_1 \leftarrow seed + s_2$
3:   $v_2 \leftarrow seed$
4:   $v_3 \leftarrow seed - s_1$
5:   **if** $n \geq 16$ **then**
6:      **for** each block $(x_0, x_1, x_2, x_3) \in \mathsf{u32}^4$ **do**            ▷ Main loop: 4 parallel lanes
7:          **for** $i \in \{0, \ldots, 3\}$ **do**
8:              $v_i \leftarrow \text{ROUND}(v_i, x_i)$
9:          **end for**
10:      **end for**
11:      $v_4 \leftarrow \text{ROT}_1(v_0) + \text{ROT}_7(v_1) + \text{ROT}_{12}(v_2) + \text{ROT}_{18}(v_3)$
12:   **else**
13:      $v_4 \leftarrow seed + s_5$
14:   **end if**
15:   $v_4 \leftarrow v_4 + n$
16:   **for** each remaining $x_i \in \mathsf{u32}$ **do**
17:      $v_4 \leftarrow \text{ROT}_{17}(v_4 + x_i \cdot s_3) \cdot s_4$
18:   **end for**
19:   **for** each remaining $x_i \in \mathsf{u8}$ **do**
20:      $v_4 \leftarrow \text{ROT}_{11}(v_4 + x_i \cdot s_5) \cdot s_1$
21:   **end for**
22:   **return** XXH32_AVALANCHE$(v_4)$

**XXH32_avalanche:**  xxHash32 finalizer (an XMS cascade):

$$\text{XXH32\_AVALANCHE}(h) = \text{XORSHIFT}_{16}\Big(\text{XMS}_{13,s_3^{32}}\big(\text{XMS}_{15,s_2^{32}}(h)\big)\Big) \tag{31}$$

## 7.3  xxHash64

Same structure as xxHash32 with 64-bit words, $r = 31$, and 32-byte blocks.

```
 1:  v₀ ← seed + s₁ + s₂
 2:  v₁ ← seed + s₂
 3:  v₂ ← seed
 4:  v₃ ← seed − s₁
 5:  if n ≥ 32 then
 6:      for each block (x₀, x₁, x₂, x₃) ∈ u64⁴ do          ▷ Main loop: 4 parallel lanes
 7:          for i ∈ {0, ..., 3} do
 8:              vᵢ ← ROUND(vᵢ, xᵢ)
 9:          end for
10:      end for
11:      v₄ ← ROT₁(v₀) + ROT₇(v₁) + ROT₁₂(v₂) + ROT₁₈(v₃)
12:      for i ∈ {0, ..., 3} do
13:          v₄ ← MERGE(v₄, vᵢ)
14:      end for
15:  else
16:      v₄ ← seed + s₅
17:  end if
18:  v₄ ← v₄ + n
19:  for each remaining xᵢ ∈ u64 do
20:      v₄ ← ROT₂₇(v₄ ⊕ ROUND(0, xᵢ)) · s₁ + s₄
21:  end for
22:  for each remaining xᵢ ∈ u32 do
23:      v₄ ← ROT₂₃(v₄ ⊕ (xᵢ · s₁)) · s₂ + s₃
24:  end for
25:  for each remaining xᵢ ∈ u8 do
26:      v₄ ← ROT₁₁(v₄ ⊕ (xᵢ · s₅)) · s₁
27:  end for
28:  return XXH64_AVALANCHE(v₄)
```

**XXH64_avalanche:**  xxHash64 finalizer:

$$\text{XXH64\_AVALANCHE}(h) = \text{XORSHIFT}_{32}(\text{XMS}_{29,s_3}(\text{XMS}_{33,s_2}(h))) \tag{32}$$
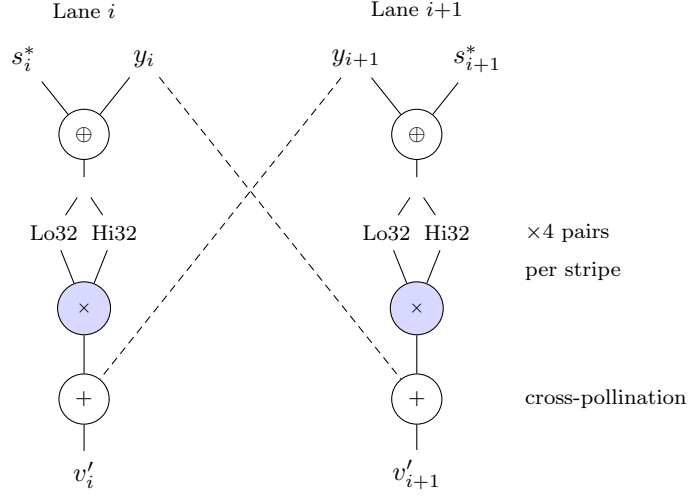
## 7.4  XXH3

Source: `hashes/xxhash.cpp`.

XXH3 (2019) uses a 192-byte secret $S$ (pseudorandom bytes from FARSH; seeded variant adds $\pm seed$ to alternating 8-byte blocks). The $32{\times}32$ multiply has native SIMD support (unlike $64{\times}64$). Original input is added to accumulators to prevent multiply-by-zero weakness. Six distinct code paths optimize for different input sizes; short keys ($n \leq 240$) use specialized scalar paths with $64{\times}64$ MUM. We summarize those special cases and focus on the long-key path.

**Long Keys ($n > 240$).**  Uses 8 parallel 64-bit accumulators with UMAC-inspired structure.

**Accumulate (per 64-byte stripe).** One stripe feeds 8 input words into 4 lane pairs. Within each pair, the two lanes cross-pollinate: the raw input of each lane is added to the *other* lane's accumulator, hardening against multiply-by-zero.



**Key insight**: The $32 \times 32 \rightarrow 64$ multiply (not $64 \times 64$!) enables efficient SIMD: SSE2/AVX2 have `_mm_mul_epu32` for this operation. The cross-lane addition of original input (dashed lines) hardens against multiply-by-zero.

1: Initialize $(v_0, \ldots, v_7) \leftarrow (s_3^{32}, s_1^{64}, s_2^{64}, s_3^{64}, s_4^{64}, s_2^{32}, s_5^{64}, s_1^{32})$
2: Let $s_0^{\text{end}}, \ldots, s_7^{\text{end}}$ be eight 64-bit secret words from the end of $S^\star$
3: Let $s_0^{\text{m}}, \ldots, s_7^{\text{m}}$ be eight 64-bit secret words from $S^\star$ at byte offset 11
4: **for** each block $(x_i)_{i=0}^{127} \in \mathsf{u64}^{128}$ **do**                    ▷ 1024 bytes = 16 stripes
5:     **for** $j \in \{0, \ldots, 15\}$ **do**                    ▷ Stripe $j$ (64 bytes)
6:         $(y_0, \ldots, y_7) \leftarrow (x_{8j}, \ldots, x_{8j+7})$
7:         Let $(s_0^*, \ldots, s_7^*) \in \mathsf{u64}^8$ be the secret words for stripe $j$
8:         **for** $i \in \{0, \ldots, 7\}$ **do**
9:             $t \leftarrow y_i \oplus s_i^*$                    ▷ $32 \times 32$ multiply of two halves (not MUM)
10:             $v_i \leftarrow v_i + (\text{Lo32}(t) \cdot \text{Hi32}(t))$
11:             $v_{i \oplus 1} \leftarrow v_{i \oplus 1} + y_i$                    ▷ Cross-pollination: raw input
12:         **end for**
13:     **end for**
14:     **for** $i \in \{0, \ldots, 7\}$ **do**                    ▷ Scramble accumulators
15:         $v_i \leftarrow (v_i \oplus (v_i \gg 47) \oplus s_i^{\text{end}}) \cdot s_1^{32}$
16:     **end for**
17: **end for**
18: Handle remaining bytes (short/medium paths summarized)
19: $v \leftarrow n \cdot s_1^{64}$                    ▷ Merge accumulators
20: **for** $i \in \{0, 2, 4, 6\}$ **do**
21:     $v \leftarrow v + \text{MUM}(v_i \oplus s_i^{\text{m}}, v_{i+1} \oplus s_{i+1}^{\text{m}})$
22: **end for**
23: $v \leftarrow v \oplus (v \gg 37)$                    ▷ Finalization (xorshift/multiply avalanche)
24: $v \leftarrow v \cdot \mathtt{0x165667919E3779F9}$
25: **return** $v \oplus (v \gg 32)$

# 8  SipHash Family
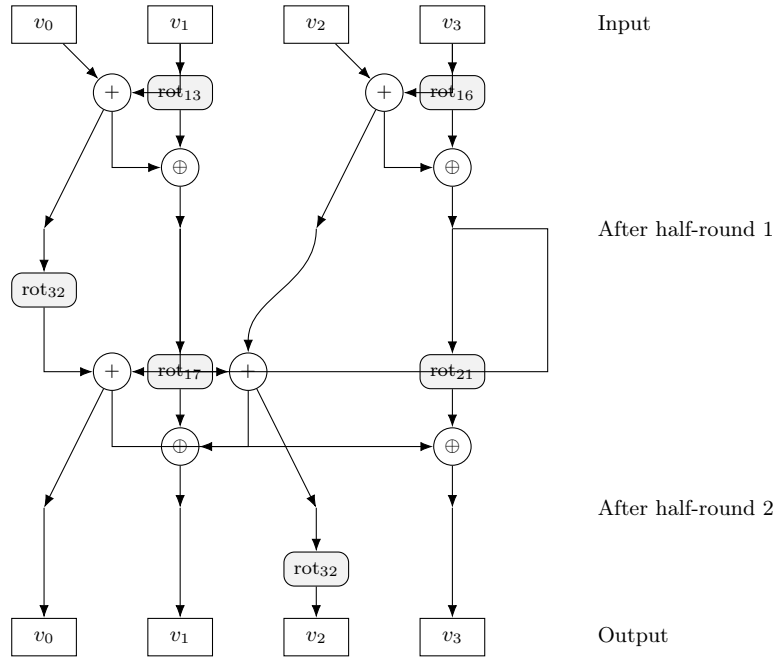
Source: `hashes/siphash.cpp` [AB12].

SipHash is a cryptographically-informed hash function designed to resist hash-flooding DoS attacks. Requires a 128-bit secret key. Uses an ARX construction (no multiplications) with a keyed permutation and Feistel-like cross-lane mixing. Sequential by design—cannot be parallelized across blocks.

## 8.1 Constants

Initialization vectors $s_0, s_1, s_2, s_3$ are fixed 64-bit constants (they encode the ASCII string "somepseudorandomlygeneratedbytes").

## 8.2 SipRound

The core permutation operates on four 64-bit state words $(v_0, v_1, v_2, v_3)$. It consists of two half-rounds with a cross-over in between:



Each half-round: $v_i \leftarrow v_i + v_{i+1}$, then $v_{i+1} \leftarrow \text{ROT}_r(v_{i+1}) \oplus v_i$. After the first half-round, $v_0 \leftarrow \text{ROT}_{32}(v_0)$ and the next half-round re-pairs the lanes to mix $(v_0, v_3)$ and $(v_2, v_1)$. The second half-round uses rotations 17 and 21, followed by $v_2 \leftarrow \text{ROT}_{32}(v_2)$.

The complete algebraic form:

$$v_0 \leftarrow v_0 + v_1 \qquad\qquad v_2 \leftarrow v_2 + v_3 \tag{33}$$
$$v_1 \leftarrow \text{ROT}_{13}(v_1) \qquad\qquad v_3 \leftarrow \text{ROT}_{16}(v_3) \tag{34}$$
$$v_1 \leftarrow v_1 \oplus v_0 \qquad\qquad v_3 \leftarrow v_3 \oplus v_2 \tag{35}$$
$$v_0 \leftarrow \text{ROT}_{32}(v_0) \tag{36}$$
$$v_2 \leftarrow v_2 + v_1 \qquad\qquad v_0 \leftarrow v_0 + v_3 \tag{37}$$
$$v_1 \leftarrow \text{ROT}_{17}(v_1) \qquad\qquad v_3 \leftarrow \text{ROT}_{21}(v_3) \tag{38}$$
$$v_1 \leftarrow v_1 \oplus v_2 \qquad\qquad v_3 \leftarrow v_3 \oplus v_0 \tag{39}$$
$$v_2 \leftarrow \text{ROT}_{32}(v_2) \tag{40}$$

## 8.3 SipHash-$c$-$d$

Parameters $c$ and $d$ specify compression and finalization rounds.

```
 1: $v_0 \leftarrow k_0 \oplus s_0$
 2: $v_1 \leftarrow k_1 \oplus s_1$
 3: $v_2 \leftarrow k_0 \oplus s_2$
 4: $v_3 \leftarrow k_1 \oplus s_3$
 5: for each block $x \in$ u64 do
 6:     $v_3 \leftarrow v_3 \oplus x$
 7:     for $i \leftarrow 1$ to $c$ do                    ▷ Compression rounds
 8:         SipRound
 9:     end for
10:     $v_0 \leftarrow v_0 \oplus x$
11: end for
12: $x \leftarrow$ padded final block with length byte
13: $v_3 \leftarrow v_3 \oplus x$
14: for $i \leftarrow 1$ to $c$ do                        ▷ Compression rounds
15:     SipRound
16: end for
17: $v_0 \leftarrow v_0 \oplus x$
18: $v_2 \leftarrow v_2 \oplus$ 0xff
19: for $i \leftarrow 1$ to $d$ do                        ▷ Finalization rounds
20:     SipRound
21: end for
22: return $v_0 \oplus v_1 \oplus v_2 \oplus v_3$
```

**Common variants:** SipHash-2-4 ($c$=2, $d$=4): recommended. SipHash-1-3 ($c$=1, $d$=3): faster, lower security margin.

### 8.4 HalfSipHash

32-bit variant for resource-constrained environments (used in Linux kernel).

**SipRound (32-bit):** Different rotation constants: $(5, 8, 7, 13, 16)$ instead of $(13, 16, 17, 21, 32)$.

# 9 Other Notable Hashes

## 9.1 ARX-Based Hashes

These hashes use Add-Rotate-XOR primitives without multiplication in their main loops, relying on superscalar execution for parallelism.

### 9.1.1 SpookyHash

Source: `hashes/spookyhash.cpp` [Jen12].

By Bob Jenkins (creator of lookup3). A 128-bit ARX hash using 12 parallel 64-bit lanes with 96-byte blocks. Version 2 is most commonly used.

**Mix Function.** The Mix function updates 12 state variables with 12 input words simultaneously. Each line combines: word addition, cross-lane XOR, self-XOR, rotation, and accumulation into adjacent lane. Rotations are: 11, 32, 43, 31, 17, 28, 39, 57, 55, 54, 22, 46.

Pattern for each lane $i$:



Each of the 12 lanes follows this pattern with different rotation amounts and cross-lane indices $j$, $k$, $l$.

Exact long-key round functions (SpookyHash v2):

```
 1: function MIX(x_0, ..., x_11, v_0, ..., v_11)              ▷ 96 bytes into 12 lanes
 2:     Let (r_0, ..., r_11) ← (11, 32, 43, 31, 17, 28, 39, 57, 55, 54, 22, 46)
 3:     for i ∈ {0, ..., 11} do                               ▷ All indices mod 12
 4:         v_i ← v_i + x_i
 5:         v_(i+2) mod 12 ← v_(i+2) mod 12 ⊕ v_(i+10) mod 12
 6:         v_(i+11) mod 12 ← v_(i+11) mod 12 ⊕ v_i
 7:         v_i ← ROT_r_i(v_i)
 8:         v_(i+11) mod 12 ← v_(i+11) mod 12 + v_(i+1) mod 12
 9:     end for
10: end function
11: function ENDPARTIAL(v_0, ..., v_11)                        ▷ Finalizer step (run 3 times)
12:     Let (r_0, ..., r_11) ← (44, 15, 34, 21, 38, 33, 10, 13, 38, 53, 42, 54)
13:     for i ∈ {0, ..., 11} do                               ▷ All indices mod 12
14:         v_(i+11) mod 12 ← v_(i+11) mod 12 + v_(i+1) mod 12
15:         v_(i+2) mod 12 ← v_(i+2) mod 12 ⊕ v_(i+11) mod 12
16:         v_(i+1) mod 12 ← ROT_r_i(v_(i+1) mod 12)
17:     end for
18: end function
```

**SpookyHash (Long Path).**

```
 1: v_0, v_3, v_6, v_9 ← seed_1
 2: v_1, v_4, v_7, v_10 ← seed_2
 3: v_2, v_5, v_8, v_11 ← 0xdeadbeefdeadbeef
 4: for each block (x_0, ..., x_11) ∈ u64^12 do              ▷ Main loop: 96-byte blocks
 5:     MIX(x_0, ..., x_11, v_0, ..., v_11)
 6: end for
 7: Let (x_0, ..., x_11) ∈ u64^12 be the final padded block (zeros + length byte)
 8: for i ∈ {0, ..., 11} do
 9:     v_i ← v_i + x_i
10: end for
11: for r ∈ {1, 2, 3} do
12:     ENDPARTIAL(v_0, ..., v_11)
13: end for
14: return (v_0, v_1)                                         ▷ 128-bit output
```

**SpookyHash v1 difference.** SpookyHash v1 applies one extra MIX to the final padded block before the three ENDPARTIAL rounds.

### 9.1.2 FarmHash

Source: `hashes/farmhash.cpp`, version 1.1 [Pik14].

By Geoff Pike at Google, successor to CityHash. Uses Murmur3-derived constants and offers multiple implementation variants (na, uo, xo) optimized for different input sizes.

**Constants.** From MurmurHash3:

$$s_0 = \text{0xc3a5c85c97cb3127} \qquad s_1 = \text{0xb492b66fbe98f273}$$
$$s_2 = \text{0x9ae16a3b2f90404f} \qquad s_3 = \text{0x9ddfea08eb382d69}$$

**HASHLEN16.**
1: **function** HASHLEN16$(u, v, m)$      ▷ 128-to-64 combiner
2:      $t_0 \leftarrow$ SHIFTMIX$((u \oplus v) \cdot m)$
3:      $t_1 \leftarrow$ SHIFTMIX$((v \oplus t_0) \cdot m)$
4:      **return** $t_1 \cdot m$
5: **end function**

**WEAKHASHLEN32WITHSEEDS.**
1: **function** WEAKHASHLEN32WITHSEEDS$(x_0, x_1, x_2, x_3, a, b)$    ▷ 128-bit for 32 bytes
2:      $a \leftarrow a + x_0$
3:      $b \leftarrow$ ROT$_{43}(b + a + x_3)$      ▷ ROTR$_{21}$
4:      $t \leftarrow a$; $a \leftarrow a + x_1 + x_2$
5:      $b \leftarrow b +$ ROT$_{20}(a)$      ▷ ROTR$_{44}$
6:      **return** $(a + x_3, b + t)$
7: **end function**

**FarmHash64 (Long Path).** Internal state: 7 words $(v_0, \ldots, v_6)$.
1: $v_0 \leftarrow 81$, $v_1 \leftarrow 81 \cdot s_1 + 113$
2: $v_2 \leftarrow$ SHIFTMIX$(v_1 \cdot s_2) \cdot s_2$
3: $(v_3, v_4), (v_5, v_6) \leftarrow (0, 0), (0, 0)$
4: **for** each block $(x_0, \ldots, x_7) \in \text{u64}^8$ **do**      ▷ 64-byte blocks
5:      On first block: $v_0 \leftarrow v_0 \cdot s_2 + x_0$, then continue below
6:      $v_0 \leftarrow$ ROT$_{37}(v_0 + v_1 + v_3 + x_1) \cdot s_1$
7:      $v_1 \leftarrow$ ROT$_{42}(v_1 + v_4 + x_6) \cdot s_1$
8:      $v_0 \leftarrow v_0 \oplus v_6$; $v_1 \leftarrow v_1 + v_3 + x_5$
9:      $v_2 \leftarrow$ ROT$_{33}(v_2 + v_5) \cdot s_1$
10:     $(v_3, v_4) \leftarrow$ WEAKHASHLEN32WITHSEEDS$(x_0, x_1, x_2, x_3, v_4 \cdot s_1, v_0 + v_5)$
11:     $(v_5, v_6) \leftarrow$ WEAKHASHLEN32WITHSEEDS$(x_4, x_5, x_6, x_7, v_2 + v_6, v_1 + x_2)$
12:     **swap** $v_2, v_0$
13: **end for**
14: Process final 64 bytes with length-dependent multiplier $m$
15: **return** HASHLEN16$(v_3 + v_5, v_4 + v_6, m)$

### 9.1.3 Rainbow

Source: `hashes/rainbow.cpp`, version 3.7.1 [Str23].

By Cris Stringfellow (DOSYAGO). A stream-based hash with 256-bit internal state using multiply-rotate primitives.

**Constants.** Eight prime constants $s_0, \ldots, s_7$ with good avalanche properties ($s_0 = 2^{64} - 59$, etc.).
1: **function** MIXA$(v_0, v_1, v_2, v_3)$      ▷ Full state mixing
2:      $v_0 \leftarrow$ ROT$_{23}(v_0 \cdot s_0) \cdot s_1$

3:    $v_1 \leftarrow \text{ROT}_{29}((v_1 \oplus v_0) \cdot s_2) \cdot s_3$

4:    $v_2 \leftarrow \text{ROT}_{31}(v_2 \cdot s_4) \cdot s_5$

5:    $v_3 \leftarrow \text{ROT}_{37}((v_3 \oplus v_2) \cdot s_6) \cdot s_7$

6: **end function**

7: **function** $\text{MIXB}(v_0, v_1, v_2, v_3, seed)$                 ▷ Lighter; operates on $v_1, v_2$ only

8:    $v_1 \leftarrow \text{ROT}_{23}(v_1 \cdot s_6) \cdot s_7$

9:    $v_2 \leftarrow \text{ROT}_{23}((v_2 \oplus v_1 + seed) \cdot s_2) \cdot s_3$

10:    **swap** $v_1, v_2$

11: **end function**

12:

13: $(v_0, v_1, v_2, v_3) \leftarrow (seed + n + 1, seed + n + 2, seed + n + 3, seed + n + 5)$

14: **for** each block $(x_0, x_1) \in \text{u64}^2$ **do**                       ▷ 16-byte blocks

15:    $v_0 \leftarrow v_0 - x_0; \ v_1 \leftarrow v_1 + x_0$

16:    $v_2 \leftarrow v_2 + x_1; \ v_3 \leftarrow v_3 - x_1$

17:    **if** block index is even **then**

18:        $\text{MIXA}(v_0, v_1, v_2, v_3)$

19:    **else**

20:        $\text{MIXB}(v_0, v_1, v_2, v_3, seed)$

21:        $(v_0, v_1, v_2, v_3) \leftarrow (v_3, v_0, v_1, v_2)$             ▷ Rotate state right

22:    **end if**

23: **end for**

24: Handle remaining bytes

25: $\text{MIXA}$; $\text{MIXB}$; $\text{MIXA}$

26: **return** $-(v_2 + v_3)$

## 9.2 PRNG-Style Hashes

These hashes have structure resembling pseudo-random number generators, often using LCG-like state updates.

### 9.2.1 komihash

Source: `hashes/komihash.cpp`, version 5.27 [Van21].

    A MUM-based hash with PRNG-inspired structure. Uses constants derived from the mantissa bits of $\pi$.

**Core Primitive.** The hash uses a modified MUM that accumulates the high part:

$$\text{KH\_M128}(m_1, m_2, lo, hi) : lo \leftarrow \text{LO}(m_1 \times m_2), \quad hi \leftarrow hi + \text{HI}(m_1 \times m_2) \tag{41}$$

Equivalently, a single "lane round" can be written as:

$$\text{KH\_ROUND}(u, v, h) : \quad (t_{\text{lo}}, t_{\text{hi}}) \leftarrow \text{MUM}_{128}(u, v) \tag{42}$$
$$\text{return } (t_{\text{lo}}, h + t_{\text{hi}}) \tag{43}$$

where in the main loop we use $u \leftarrow v_p \oplus x_p$ and $v \leftarrow v_q \oplus x_q$.

**komihash (Long Path).** Constants $s_1, \ldots, s_8$ are derived from $\pi$.

1: $v_1 \leftarrow s_1 \oplus (seed \wedge \texttt{0x5555...})$                 ▷ Checkerboard seed split

2: $v_5 \leftarrow s_5 \oplus (seed \wedge \texttt{0xAAAA...})$

3: $(v_1, v_5) \leftarrow \text{KH\_ROUND}(v_1, v_5, v_5)$             ▷ Initial round for diffusion

4: $v_i \leftarrow s_i \oplus v_1$ for $i \in \{2, 3, 4\}$;     $v_i \leftarrow s_i \oplus v_5$ for $i \in \{6, 7, 8\}$

5: **for** each block $(x_0, \ldots, x_7) \in$ u64$^8$ **do** $\qquad \qquad \triangleright$ 4 parallel lanes
6: $\quad$ **for** $i \in \{1, \ldots, 4\}$ **do**
7: $\qquad (v_i, v_{i+4}) \leftarrow$ KH_ROUND$(v_i \oplus x_{i-1}, v_{i+4} \oplus x_{i+3}, v_{i+4})$
8: $\quad$ **end for**
9: $\quad v_4 \leftarrow v_4 \oplus v_7; v_1 \leftarrow v_1 \oplus v_8; v_3 \leftarrow v_3 \oplus v_6; v_2 \leftarrow v_2 \oplus v_5 \qquad \triangleright$ Cross-lane
10: **end for**
11: **for** $i \in \{6, \ldots, 8\}$ **do** $\qquad \qquad \triangleright$ Collapse
12: $\quad v_5 \leftarrow v_5 \oplus v_i$
13: **end for**
14: **for** $i \in \{2, \ldots, 4\}$ **do**
15: $\quad v_1 \leftarrow v_1 \oplus v_i$
16: **end for**
17: Handle remaining bytes (up to 2 more 16-byte rounds via KH_M128)
18: $v_1 \leftarrow v_1 \oplus v_5 \qquad \qquad \triangleright$ Final mixing
19: Apply one more KH_ROUND, then $v_1 \leftarrow v_1 \oplus v_5$
20: **return** $v_1$

### 9.2.2 prvhash

Source: `hashes/prvhash.cpp`, version 4.3.4 [Van20].

$\quad$ By Aleksey Vaneev. A PRNG-style hash that generates pseudo-random values during hashing. Uses checkerboard constants like a5hash.

**Core Function.**
1: **function** PRVHASH_CORE64$(v_0, v_1, v_2)$
2: $\quad v_0 \leftarrow v_0 \cdot (v_1 \cdot 2 + 1)$
3: $\quad t \leftarrow (v_0 \gg 32) \mid (v_0 \ll 32) \qquad \qquad \triangleright$ Swap halves
4: $\quad v_2 \leftarrow v_2 + t +$ `0xAAAAAAAAAAAAAAAA`
5: $\quad v_1 \leftarrow v_1 + v_0 +$ `0x5555555555555555`
6: $\quad v_0 \leftarrow v_0 \oplus v_2$
7: $\quad$ **return** $v_1 \oplus t$
8: **end function**

**prvhash64.**
1: Initialize $(v_0, v_1, v_2)$ from pre-computed constants
2: $v_0 \leftarrow v_0 \oplus seed; v_1 \leftarrow v_1 \oplus seed$
3: **for** each word $x \in$ u64 **do**
4: $\quad$ PRVHASH_CORE64$(v_0, v_1, v_2)$
5: $\quad v_0 \leftarrow v_0 \oplus x; v_1 \leftarrow v_1 \oplus x$
6: **end for**
7: Apply 2 more core rounds
8: **return** PRVHASH_CORE64 output

## 9.3 Compact and Embedded Hashes

These hashes prioritize small code size and simplicity, suitable for resource-constrained environments.

### 9.3.1 ChibiHash

Source: `hashes/chibihash.cpp`, version 2 [NRK24].

A compact, high-quality 64-bit hash by NRK with small code footprint for embedded use. Uses a single 64-bit constant (digits of $e$) and truncated 64×64→64 multiplies (not MUM). Cross-lane rotation feeding prevents lane independence.

1: $s_0 \leftarrow \texttt{0x2B7E151628AED2A7}$ ▷ Digits of $e$
2: $t \leftarrow \text{ROT}_{15}(seed - s_0) + \text{ROT}_{47}(seed - s_0)$
3: $v_0, v_1, v_2, v_3 \leftarrow [seed,\ seed + s_0,\ t,\ t + s_0^2 \oplus s_0]$
4: **for** each block $(x_0, \ldots, x_3) \in \texttt{u64}^4$ **do** ▷ 32-byte blocks
5:     **for** $i \in \{0, \ldots, 3\}$ **do**
6:         $v_i \leftarrow (x_i + v_i) \cdot s_0$
7:         $v_{(i+1) \bmod 4} \leftarrow v_{(i+1) \bmod 4} + \text{ROT}_{27}(x_i)$
8:     **end for**
9: **end for**
10: Handle remaining bytes
11: $v_0 \leftarrow v_0 + (\text{ROT}_{31}(v_2 \cdot s_0) \oplus (v_2 \gg 31))$
12: $v_1 \leftarrow v_1 + (\text{ROT}_{31}(v_3 \cdot s_0) \oplus (v_3 \gg 31))$
13: $v_0 \leftarrow v_0 \cdot s_0;\ v_0 \leftarrow v_0 \oplus (v_0 \gg 31)$
14: $v_1 \leftarrow v_1 + v_0$
15: $t \leftarrow n \cdot s_0;\ t \leftarrow t \oplus \text{ROT}_{29}(t);\ t \leftarrow t + seed;\ t \leftarrow t \oplus v_1$
16: $t \leftarrow t \oplus \text{ROT}_{15}(t) \oplus \text{ROT}_{42}(t);\ t \leftarrow t \cdot s_0$
17: **return** $t \oplus \text{ROT}_{13}(t) \oplus \text{ROT}_{31}(t)$

### 9.3.2 XMSX

Source: `hashes/xmsx.cpp`.

By Dmitrii Lebed. A minimal 32-bit hash for microcontrollers using XOR-Multiply-Shift-XOR rounds. Uses 32×32→64 multiply (widely available in HW). Faster than software CRC32 on 32-bit CPUs.

**Round Function.**

$$\text{XMSX32\_ROUND}(v, x) = ((v \oplus x) \cdot s_0) \oplus (((v \oplus x) \cdot s_0) \gg 32) \tag{44}$$

where $s_0 = \texttt{0xcdb32970830fcaa1}$.

**XMSX32 Algorithm.**

1: $v \leftarrow (seed \ll 32)\ |\ seed$
2: $v \leftarrow \text{XMSX32\_ROUND}(v, n)$
3: **for** each word $x \in \texttt{u32}$ **do**
4:     $v \leftarrow \text{XMSX32\_ROUND}(v, x)$
5: **end for**
6: **return** $\text{XMSX32\_ROUND}(v, v \gg 47)$

# 10 CLMUL and Polynomial Hashes

These hashes use carry-less multiplication (CLMUL) or polynomial arithmetic over finite fields to achieve theoretically grounded collision bounds.

**Carry-less vs. standard multiplication.** In standard multiplication, partial products are added with carry propagation. In carry-less multiplication (polynomial multiplication over GF(2)), partial products are XORed without carries:

|     Standard (with carry)     |     Carry-less (XOR)     |
|---|---|

$$\times \;1\;\;0\;\;1\;\;1$$
$$1\;\;1\;\;0\;\;1$$

partial products
added with carry

$$= 143 \text{ (8 bits)}$$

$$\otimes \;1\;\;0\;\;1\;\;1$$
$$1\;\;1\;\;0\;\;1$$

partial products
XORed (no carry)

$$= 127 \text{ (7 bits)}$$

The PCLMULQDQ instruction performs $64{\times}64{\rightarrow}128$ carry-less multiplication in a single cycle on modern CPUs.

**Pseudocode primitives.**

$$\text{CLMUL}(a, b) = a \otimes b \quad \text{(128-bit carry-less product of two 64-bit inputs)} \tag{45}$$

Both map to a single `PCLMULQDQ` instruction with $\sim$1 cycle throughput. The symbol $\otimes$ denotes polynomial multiplication in $\text{GF}(2)[x]$ (XOR replaces addition, AND replaces multiplication of coefficients).

## 10.1   CLHash

Source: `hashes/clhash.cpp` [LK16].

**SMHasher3 test results.**   CLHash fails 186 of 250 tests, the worst result of any hash in this survey. Failures span nearly every category: Avalanche, BIC, Sparse, Permutation, Text, TwoBytes, PerlinNoise, Bitflip, and extensive seed-related tests. Many seed-test failures are expected: CLHash is a keyed hash designed for a full random key, and SMHasher3's seed-sensitivity tests probe a regime outside its design model. The non-seed failures (Avalanche, Sparse, Text) indicate that the CLMUL-based compression without a strong finalizer leaves detectable bias in the output.

A theoretically-grounded hash by Daniel Lemire using carry-less multiplication (CLMUL instruction, x86-64 Haswell+). Based on almost-universal hashing over $\text{GF}(2^{127})$ with collision bound $\epsilon \approx n/2^{64}$. Uses 133 random 64-bit keys derived from seed.

### 10.1.1   Mathematical Foundation

CLHash computes a polynomial hash over $\text{GF}(2^{127})$ with irreducible polynomial $p(x) = x^{127} + x + 1$. (Irreducibility over $\text{GF}(2^{128})$ ensures the field structure.) The core operation is carry-less multiplication:

$$a \otimes b = a \cdot b \mod p(x) \tag{46}$$

where $\cdot$ denotes polynomial multiplication in $\text{GF}(2)[x]$ (XOR instead of addition, AND instead of multiplication for coefficients).

### 10.1.2   Key Primitives

**Lazy mod $2^{127} + 2 + 1$:**   Given 254-bit product $(A_{\text{hi}}, A_{\text{lo}})$:

$$A_{\text{lo}} \oplus (A_{\text{hi}} \ll 1) \oplus (A_{\text{hi}} \ll 2) \tag{47}$$

**128-to-64 reduction:** Uses the irreducible polynomial $p(x) = x^{64} + x^4 + x^3 + x + 1$ (notation: $(64, 4, 3, 1, 0)$). The reduction is performed in two steps:

**Step 1: CLMUL reduction.** Let $C = x^4 + x^3 + x + 1$ (the low-degree terms of $p(x)$). Since $x^{64} \equiv C \pmod{p(x)}$, we have $A_{\text{hi}} \cdot x^{64} \equiv A_{\text{hi}} \cdot C$:

$$Q_2 \leftarrow \text{CLMUL}(A_{\text{hi}}, C) \tag{48}$$

This produces a result up to $64 + 4 = 68$ bits (since $A_{\text{hi}}$ is 64 bits and $C$ is degree 4).

**Step 2: Table lookup for final bits.** The CLMUL result $Q_2$ may have bits in positions 64–67. These are reduced using a precomputed 16-entry table where entry $i$ contains $i \cdot C \bmod 2^{64}$:

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| $i \cdot C$ | 0 | 27 | 54 | 45 | 108 | 119 | 90 | 65 |

| $i$ | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|
| $i \cdot C$ | 216 | 195 | 238 | 245 | 180 | 175 | 130 | 153 |

The high 4 bits of $Q_2$ (bits 64–67) index into this table via `pshufb`:

$$Q_3 \leftarrow table[(Q_2 \gg 64) \wedge \texttt{0xF}] \tag{49}$$

**Final combination:**

$$\text{PRECOMPREDUCTION64}(A) = A_{\text{lo}} \oplus Q_2 \oplus Q_3 \tag{50}$$

The table lookup replaces what would otherwise require another CLMUL for the overflow bits, saving one instruction in the critical path.

### 10.1.3 CLHash Algorithm

1: Generate 133 random 64-bit key words $s_0, \ldots, s_{132}$ from seed using xorshift128+
2: $s_{\text{poly}} \leftarrow (s_{128} \mid s_{129} \ll 64) \wedge (\texttt{0x3FFFFFFF FFFFFFFF FFFFFFFF FFFFFFFF})$  ▷ 128-bit; top 2 bits cleared for $\text{GF}(2^{127})$
3: $s_{\text{final}} \leftarrow (s_{130}, s_{131}) \in \texttt{u128};\quad s_{\text{len}} \leftarrow s_{132} \in \texttt{u64}$
4: $v \leftarrow 0$ (128-bit)
5: **for** each block $(x_i)_{i=0}^{127} \in \texttt{u64}^{128}$ **do**  ▷ 1024 bytes; Horner: $v \leftarrow v \otimes s_{\text{poly}} \oplus h_i$
6:     **if** not first block **then**
7:         $v \leftarrow v \otimes s_{\text{poly}}$
8:     **end if**
9:     **for** $i \in \{0, \ldots, 63\}$ **do**  ▷ 64 word pairs
10:         $v \leftarrow v \oplus \text{CLMUL}((x_{2i} \oplus s_{2i}), (x_{2i+1} \oplus s_{2i+1}))$
11:     **end for**
12: **end for**
13: **for** remaining word pairs $(x_{2i}, x_{2i+1})$ **do**  ▷ Tail: same CLMUL pattern
14:     $v \leftarrow v \oplus \text{CLMUL}((x_{2i} \oplus s_{2i}), (x_{2i+1} \oplus s_{2i+1}))$
15: **end for**
16: $t \leftarrow v \oplus s_{\text{final}};\quad v \leftarrow \text{CLMUL}(\text{LO}(t), \text{HI}(t))$
17: $v \leftarrow v \oplus \text{CLMUL}(s_{\text{len}}, n)$  ▷ Mix in length
18: **return** $\text{PRECOMPREDUCTION64}(v)$  ▷ bitmix variant applies MurmurHash3 FMIX64 after

## 10.2 UMASH

Source: `hashes/umash.cpp` [Khu20b].

**SMHasher3 test results.** UMASH-128 fails 128 of 250 tests. Like CLHash, UMASH is a keyed hash expecting a full random key (generated via Salsa20), and most failures are in seed-related tests (Seed, SeedSparse, SeedAvalanche, SeedBIC, SeedBitflip, SeedBlockLen, Seed-BlockOffset). Non-seed failures include BIC, some Sparse configurations, Permutation (single-bit keys), TwoBytes, PerlinNoise, and Bitflip—reflecting the inherent linearity of CLMUL-based compression.

UMASH by Paul Khuong combines CLMUL-based compression with polynomial hashing over $GF(2^{61} - 1)$. It produces two independent hash functions for fingerprinting (128-bit output).

### 10.2.1 Mathematical Foundation

UMASH works modulo the Mersenne prime $p = 2^{61} - 1$. Key operations:

**Modular addition (fast):**

$$\text{ADD\_MOD\_FAST}(x, y) = \begin{cases} x + y + 8 & \text{if } x + y \text{ overflows} \\ x + y & \text{otherwise} \end{cases} \tag{51}$$

**Modular multiplication:**

$$\text{MUL\_MOD\_FAST}(m, x) = \text{ADD\_MOD\_FAST}(\text{LO}(m \times x), 8 \cdot \text{HI}(m \times x)) \tag{52}$$

**Horner double update:**

$$\begin{aligned} \text{HORNER\_DOUBLE\_UPDATE}(v, m_0, m_1, x, y) = \text{MUL\_MOD\_FAST}(m_0, v + x) \\ + \text{MUL\_MOD\_FAST}(m_1, y) \end{aligned} \tag{53}$$

### 10.2.2 OH Compression

The "OH" (One-pass Hash) compresses 256-byte blocks using CLMUL:

1: Interpret the 256-byte input as $(x_i)_{i=0}^{31} \in \mathsf{u64}^{32}$
2: $v \leftarrow 0$ (128-bit)
3: **for** $i \in \{0, \dots, 14\}$ **do**          ▷ 15 CLMUL pairs
4:      $v \leftarrow v \oplus \text{CLMUL}(x_{2i} \oplus s_{2i}, x_{2i+1} \oplus s_{2i+1})$
5: **end for**
6: **ENH finalization:**
7: $(t_0, t_1) \leftarrow (x_{30} + s_{30}, x_{31} + s_{31})$          ▷ Final ENH chunk
8: $t_2 \leftarrow \text{HI}(t_0 \times t_1) + seed$; $t_3 \leftarrow \text{LO}(t_0 \times t_1)$
9: **return** $(v_{\text{lo}} \oplus t_3, v_{\text{hi}} \oplus t_2 \oplus t_3)$

### 10.2.3 UMASH Algorithm

We focus on the main (long-key) structure. Short-key special casing and exact tail handling are summarized.

1: Generate 34 random 64-bit key words $s_0, \dots, s_{33}$ from seed via Salsa20
2: Prepare polynomial multipliers $f, f^2$ in $GF(2^{61} - 1)$
3: $v \leftarrow 0$
4: **for** each block $(x_i)_{i=0}^{31} \in \mathsf{u64}^{32}$ **do**          ▷ Main loop (256 bytes)
5:      $(t_0, t_1) \leftarrow \text{OH\_COMPRESS}((x_i)_{i=0}^{31})$
6:      $v \leftarrow \text{HORNER\_DOUBLE\_UPDATE}(v, f^2, f, t_0, t_1)$
7: **end for**
8: Process remaining bytes with variable-length OH (tail summarized)
9: **return** $(v \oplus \text{ROT}_8(v)) \oplus \text{ROT}_{33}(v)$          ▷ Finalization

### 10.2.4 Fingerprinting (128-bit)

For fingerprinting, UMASH computes two independent hashes using:

- Different polynomial multipliers $(f_0, f_0^2)$ and $(f_1, f_1^2)$

- A "twisted" version using an LRC (longitudinal redundancy check) XOR [Khu20a]

- Both hashes use the same OH keys but different polynomial accumulators

### 10.3 Polymur

Source: `hashes/polymur.cpp` [Pet23].

Polymur by Orson Peters is a polynomial hash over $GF(2^{61} - 1)$. Uses 56-bit chunks to avoid reduction within chunk processing. The polynomial multiplier $s_0$ is chosen as a generator of the multiplicative group (via $s_0 = 37^e$ for random odd $e$ coprime to $2^{61} - 2$), with $s_0^7 < 2^{60} - 2^{56}$ ensuring efficient reduction.

1: Keys: $(s_0, s_1) \leftarrow$ secret parameters derived from seed $\qquad \triangleright s_0$ multiplier, $s_1$ tweak
2: $v \leftarrow 0$
3: **if** $n \leq 7$ **then**
4: $\quad$ **return** $s_1 + (s_0 + x_0)(s_0^2 + n) \mod p$
5: **else if** $n \geq 50$ **then**
6: $\quad$ **while** $n \geq 50$ **do** $\qquad\qquad \triangleright$ Main loop: 7×56-bit chunks per iteration
7: $\qquad$ Let $(x_0, \ldots, x_6)$ be the next 7 words (each 56 bits, zero-extended)
8: $\qquad t_0 \leftarrow (s_0 + x_0)(s_0^6 + x_1)$
9: $\qquad t_1 \leftarrow (s_0^2 + x_2)(s_0^5 + x_3)$
10: $\qquad t_2 \leftarrow (s_0^3 + x_4)(s_0^4 + x_5)$
11: $\qquad t_3 \leftarrow (v + x_6) \cdot s_0^7$
12: $\qquad v \leftarrow (t_0 + t_1 + t_2 + t_3) \mod p$
13: $\quad$ **end while**
14: $\quad$ Finalize with $v \cdot s_0^{14}$
15: **end if**
16: Handle tail (8–49 bytes) with overlapping reads
17: $s_2 \leftarrow$ `0xe9846af9b1a615d` $\qquad\qquad\qquad\qquad\qquad\qquad \triangleright$ mx3 finalizer
18: $v \leftarrow \text{XORSHIFT}_{32}(v) \cdot s_2$
19: $v \leftarrow \text{XORSHIFT}_{32}(v) \cdot s_2$
20: $v \leftarrow \text{XORSHIFT}_{28}(v)$
21: **return** $v + s_1$

### 10.4 poly-mersenne

Source: `hashes/poly_mersenne.cpp` [AKT20].

A simple polynomial hash over the Mersenne prime $p = 2^{61} - 1$. Outputs 32 bits. Keys from SplitMix64: multiplier $s_0 < 2^{60}$ and random values $s_1, \ldots, s_{K+1}$ for optional $K$-independence.

**Core primitive.** The multiply-accumulate uses *lazy Mersenne reduction*: since $2^{61} \equiv 1 \pmod{p}$, a 128-bit product can be reduced by folding the bits above position 61 back into the low part with a shift and add—no division required.

1: **function** MULT_COMBINE61$(v, m, x)$ $\qquad\qquad\qquad \triangleright (v \cdot m + x) \mod (2^{61} - 1)$
2: $\quad (lo, hi) \leftarrow v \times m$ $\qquad\qquad\qquad\qquad\qquad \triangleright$ Full 128-bit product
3: $\quad lo \leftarrow lo + x$ $\qquad\qquad\qquad\qquad\qquad \triangleright$ Accumulate input word
4: $\quad t \leftarrow (hi \ll 3) \mid (lo \gg 61)$ $\qquad\qquad \triangleright$ Bits $\geq 2^{61}$: fold back
5: $\quad lo \leftarrow lo \wedge (2^{61} - 1)$

6:     **return** $lo + t$                          ▷ Lazy: may exceed $p$ by at most 1
7: **end function**

The final result is fully reduced ($v \leftarrow v - p$ if $v \geq p$) once, after all rounds.

1: $v \leftarrow n$                               ▷ Initialize with length
2: **for** each word $x \in \mathtt{u32}$ **do**
3:     $v \leftarrow \textsc{mult\_combine}61(v, s_0, x)$
4: **end for**
5: **if** $K > 0$ **then**                         ▷ Optional $K$-independence transform
6:     $t \leftarrow v$; $v \leftarrow s_1$
7:     **for** $i \in \{2, \ldots, K+1\}$ **do**: $v \leftarrow \textsc{mult\_combine}61(v, t, s_i)$
8:     **end for**
9: **end if**
10: **return** $v \bmod 2^{32}$                     ▷ After final reduction if $v \geq p$

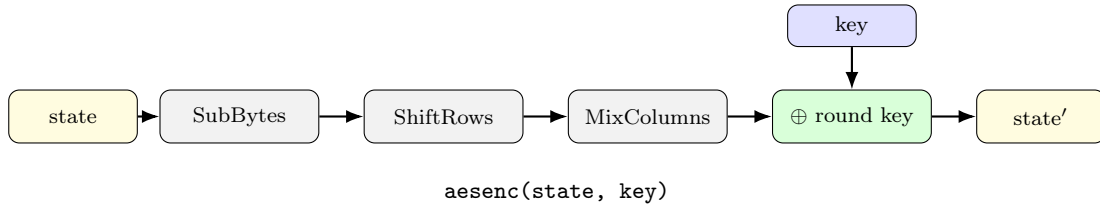Variants with higher $K$ provide $K$-wise independence (collision bound $(n/2^{61})^{K+1}$):

| Variant | $K$ | Collision bound |
|---|---|---|
| poly-mersenne-deg0 | 0 | $n/2^{61}$ |
| poly-mersenne-deg1 | 1 | $(n/2^{61})^2$ |
| poly-mersenne-deg2 | 2 | $(n/2^{61})^3$ |
| poly-mersenne-deg3 | 3 | $(n/2^{61})^4$ |
| poly-mersenne-deg4 | 4 | $(n/2^{61})^5$ |

Higher-degree variants provide $K$-wise independence at the cost of additional multiplications, improving theoretical collision bounds.
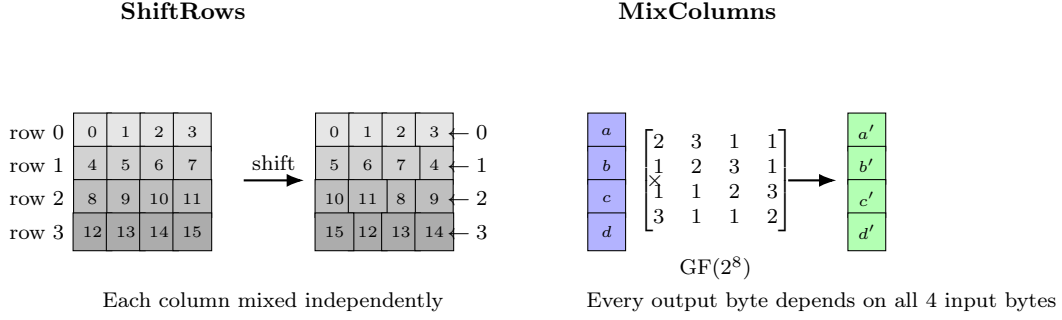
# 11   AES-Based Hashes

These hashes leverage the AES-NI instructions (`aesenc`, `aesdec`) for fast mixing. AES provides excellent diffusion in a single instruction.

**AES round structure.**   Each `aesenc` instruction applies three transformations to a 128-bit state, then XORs a round key:



aesenc(state, key)

**State transformations.**   The AES state is a $4 \times 4$ matrix of bytes. ShiftRows and MixColumns provide diffusion across the entire state:

**ShiftRows**             **MixColumns**

Each column mixed independently      Every output byte depends on all 4 input bytes

**SubBytes**: Non-linear S-box substitution (each byte independently). **ShiftRows**: Rotates each row by different amounts (0, 1, 2, 3 positions left). **MixColumns**: Multiplies each 4-byte column by a fixed matrix in $GF(2^8)$. After one round, each output bit depends on multiple input bits—full avalanche requires $\sim 2$ rounds.

**Pseudocode primitives.** We use AESENC and AESDEC as 128-bit $\rightarrow$ 128-bit operations throughout:

$$\text{AESENC}(v, k) = \text{MixColumns}(\text{ShiftRows}(\text{SubBytes}(v))) \oplus k \tag{54}$$

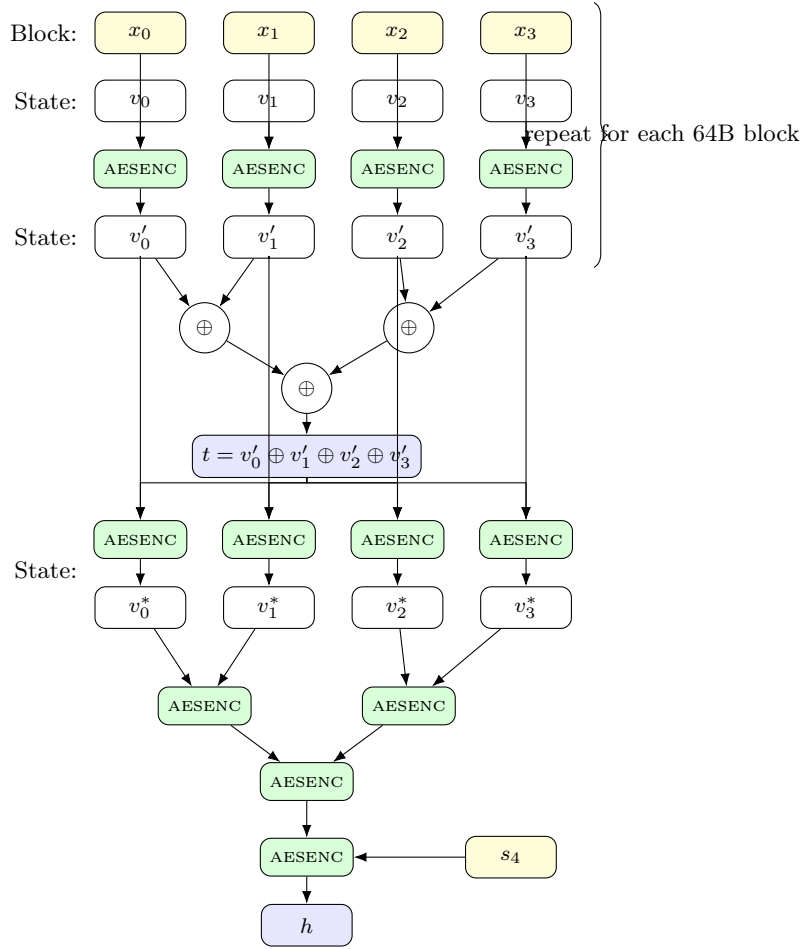$$\text{AESDEC}(v, k) = \text{InvMixColumns}(\text{InvShiftRows}(\text{InvSubBytes}(v))) \oplus k \tag{55}$$

Each maps to a single x86 instruction (AESENC / AESDEC) with $\sim 1$ cycle throughput and $\sim 4$ cycle latency on modern CPUs.

## 11.1 AquaHash

Source: `hashes/aquahash.cpp` [Rog19].

**SMHasher3 test results.** AquaHash fails 55 of 250 tests. The main weaknesses are in Sparse tests (18 configurations, including long-key variants up to 2048 bytes), Permutation (single-bit keys), and several Text patterns. Some seed-related failures (SeedSparse, SeedBlockLen, Seed) also appear.

AquaHash by J. Andrew Rogers uses `aesenc` for mixing with 4 parallel 128-bit lanes (512-bit state). We denote the fixed 128-bit constants used by the reference implementation as $s_0, \ldots, s_4$; $s_4$ is used in the final AES round.

**Short keys** ($n < 64$): Uses 2 parallel lanes with AESENC, combining via $\text{AESENC}(v_0, v_1)$ before a 3-round finalization.
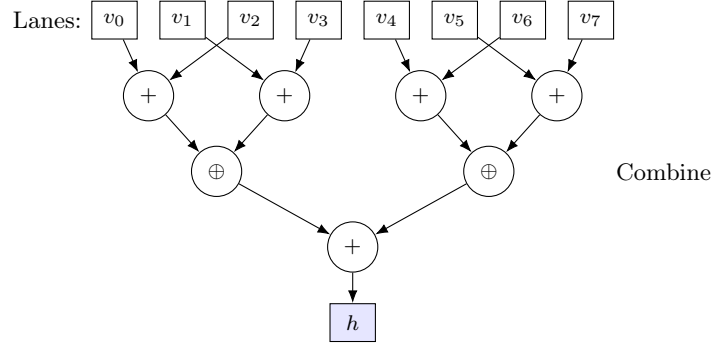
**Long keys** ($n \geq 64$).

1: $(v_0, v_1, v_2, v_3) \leftarrow (seed \oplus s_0, seed \oplus s_1, seed \oplus s_2, seed \oplus s_3)$
2: **for** each block $(x_0, x_1, x_2, x_3) \in \text{u128}^4$ **do**               ▷ 64 bytes
3:      **for** $i \in \{0, \ldots, 3\}$ **do**
4:          $v_i \leftarrow \text{AESENC}(v_i, x_i)$
5:      **end for**
6: **end for**
7: Process remaining 32, 16, 8, 4, 2, 1 bytes
8: $t \leftarrow v_0 \oplus v_1 \oplus v_2 \oplus v_3$
9: **for** $i \in \{0, \ldots, 3\}$ **do**
10:      $v_i \leftarrow \text{AESENC}(v_i, t)$
11: **end for**
12: $v \leftarrow \text{AESENC}(\text{AESENC}(v_0, v_1), \text{AESENC}(v_2, v_3))$
13: **return** $\text{AESENC}(v, s_4)$

## 11.2   MeowHash

Source: `hashes/meowhash.cpp`, version 0.5/calico [Mur18].

MeowHash by Casey Muratori (Molly Rocket) uses 8 parallel 128-bit lanes (1024-bit state) with `aesdec` for mixing. Processes 256-byte blocks matching cache line boundaries. High

throughput by saturating AES execution units; for inputs > 256KB, explicit prefetch improves performance.



Default seed: 128 bytes from digits of $\pi$. All additions below are 64-bit lane-wise (`paddq`).

1: **function** MEOW_MIX($r_1, r_2, r_3, r_4, r_5, x_a, x_b$)                    ▷ 32 bytes into 5 lanes
2:     $r_1 \leftarrow \text{AESDEC}(r_1, r_2)$                    ▷ Source uses 4 overlapping 128-bit reads;
3:     $r_3 \leftarrow r_3 + x_a; \quad r_2 \leftarrow r_2 \oplus x_b$                    ▷ simplified to 2 blocks here
4:     $r_2 \leftarrow \text{AESDEC}(r_2, r_4)$
5:     $r_5 \leftarrow r_5 + x_a; \quad r_4 \leftarrow r_4 \oplus x_b$
6: **end function**
7: **function** MEOW_SHUFFLE($r_1, r_2, r_3, r_4, r_5, r_6$)                    ▷ Cross-lane mixing (no input)
8:     $r_1 \leftarrow \text{AESDEC}(r_1, r_4)$
9:     $r_2 \leftarrow r_2 + r_5; \quad r_4 \leftarrow r_4 \oplus r_6$
10:     $r_4 \leftarrow \text{AESDEC}(r_4, r_2)$
11:     $r_5 \leftarrow r_5 + r_6$
12: **end function**
13:
14: Initialize 8 lanes $(v_0, \ldots, v_7)$ from 128-byte seed
15: **for** each block $(x_i)_{i=0}^{15} \in \texttt{u128}^{16}$ **do**                    ▷ Main loop (256 bytes)
16:     **for** $i \in \{0, \ldots, 7\}$ **do**                    ▷ 32 bytes per iteration; touches 5 lanes
17:         MEOW_MIX($v_i, v_{(i+4) \bmod 8}, v_{(i+6) \bmod 8}, v_{(i+1) \bmod 8}, v_{(i+2) \bmod 8}, x_{2i}, x_{2i+1}$)
18:     **end for**
19: **end for**
20: Load residual < 32 bytes with careful alignment handling
21: Mix residual and length into lanes
22: Process remaining 32-byte blocks (0–7)
23: **MixDown:** Apply 12 MEOW_SHUFFLE operations
24: Combine: $v_0 \leftarrow ((v_0 + v_2) \oplus (v_1 + v_3)) + ((v_4 + v_6) \oplus (v_5 + v_7))$
25: **return** $v_0$

## 11.3   aHash

Source: `hashes/rust-ahash.cpp` [Kai18].

By Tom Kaitchuck. Default hash in Rust's `HashMap`. Two variants: AES-based (2 parallel paths using `aesdec` + shuffle-add) and fallback (portable MUM-based, similar to wyhash). Uses variable-amount final rotation for additional mixing.

**AES Variant.**   Uses `aesdec` with a special shuffle permutation (selected by automated search): SHUFFLE is a fixed byte permutation of a 128-bit value (as in the implementation).

1: $(v_0, v_1) \leftarrow$ state from random seed (derived from $\pi$)
2: $t \leftarrow v_0 \oplus v_1$                    ▷ Save initial combined state

3: **for** each block $x \in$ u128 **do**
4:    $v_0 \leftarrow \text{AESDEC}(v_0, x)$
5:    $v_1 \leftarrow \text{SHUFFLE}(v_1) + x$
6: **end for**
7: $v \leftarrow \text{AESENC}(v_1, v_0)$
8: **return** $\text{AESDEC}(\text{AESDEC}(v, t), v)$

**Fallback Variant.**  MUM-based, using Knuth's LCG multiplier $s_0 = 6364136223846793005$:

1: Initialize $v, s_1, s_2, s_3$ from random state
2: $v \leftarrow \text{MUM}(v \oplus n,\ s_0)$
3: $v \leftarrow (v + n) \cdot s_0$
4: **for** each block $(x_0, x_1) \in$ u64$^2$ **do**                                    ▷ 16 bytes
5:    $t \leftarrow \text{MUM}(x_0 \oplus s_2,\ x_1 \oplus s_3)$
6:    $v \leftarrow (v + s_1) \oplus t$
7:    $v \leftarrow \text{ROT}_{23}(v)$
8: **end for**
9: **return** $\text{ROT}_{v \bmod 64}(\text{MUM}(v,\ s_1))$

# 12    SIMD-Scalable Hashes

## 12.1    HalftimeHash

Source: `hashes/halftimehash.cpp` [App21].

**SMHasher3 test results.**  HalftimeHash-128 fails 48 of 250 tests. The dominant failure category is Sparse (27 configurations including long-key variants), followed by Permutation (single-bit and low-bit keys) and SeedZeroes. As a keyed hash requiring ∼50 KB of random state, the SeedZeroes failures reflect SMHasher3's seed model rather than a design flaw.

HalftimeHash by Jim Apple is a SIMD-scalable hash that efficiently uses any vector width (64-bit scalar to 512-bit AVX-512). Based on erasure coding theory for provable mixing properties. Uses 32-bit multiplies (`_mm_mul_epu32`), a tree structure for long inputs, and requires ∼50KB of random state.

The core mixing uses $\text{TIMES}(a, b) = (a \wedge \text{0xFFFFFFFF}) \times (b \wedge \text{0xFFFFFFFF})$ within 64-bit lanes. The Mix function computes $v + \text{TIMES}(t, t \gg 32)$ where $t = s +_{32} x$ (independent 32-bit lane addition, no carry between halves).

Erasure coding (like RAID) achieves strong mixing: input blocks are encoded with redundancy (e.g., Encode2: $6 \rightarrow 7$ blocks), then reduced using weighted sums. For long inputs, a tree structure with fanout 8 combines block groups via depth-first traversal. Style variants select block size: Style64 (scalar), Style128 (SSE2/NEON), Style256 (AVX2), Style512 (AVX-512). Output is finalized using byte-level tabulation: $result = \text{TABULATEBYTES}(n) \oplus \bigoplus_j \text{TABULATEBYTES}(h_j)$, where $(h_j)$ are the final 64-bit output words for the chosen variant.

1: **function** $\text{MIX}(v, x, s)$                                    ▷ Core mixing per block
2:    $t \leftarrow s +_{32} x$                                    ▷ 32-bit lane addition
3:    **return** $v + \text{TIMES}(t, t \gg 32)$              ▷ TIMES: 32×32 low-half multiply
4: **end function**
5: **function** $\text{COMBINE}(u_0, \ldots, u_{d'-1})$            ▷ Fixed linear reducer (variant-specific)
6:    Compute $(y_0, \ldots, y_{L-1})$ as weighted sums (dot products) of the $u_i$
7:    **return** $(y_0, \ldots, y_{L-1})$
8: **end function**
9: **function** $\text{BASELAYER}(\text{input block group}, s)$
10:    Encode $d$ input blocks into $d'$ blocks via erasure code (XOR-based)

11:     Hash each encoded block into one value $u_i$ using repeated MIX (details elided)
12:     **return** COMBINE$(u_0, \ldots, u_{d'-1})$
13: **end function**
14: **function** UPPERLAYER$((u_0^{(0)}, \ldots, u_{L-1}^{(0)}), \ldots, (u_0^{(7)}, \ldots, u_{L-1}^{(7)}), s)$
15:     **for** $i \in \{0, \ldots, L-1\}$ **do**
16:         $y_i \leftarrow u_i^{(0)}$
17:         **for** $j \in \{1, \ldots, 7\}$ **do**
18:             $y_i \leftarrow$ MIX$(y_i, u_i^{(j)}, s_{i,j})$
19:         **end for**
20:     **end for**
21:     **return** $(y_0, \ldots, y_{L-1})$
22: **end function**
23: **function** TABULATEBYTES$(x, T)$                              ▷ $T$ is an 8×256 table
24:     $r \leftarrow 0$
25:     **for** $i \in \{0, \ldots, 7\}$ **do**
26:         $r \leftarrow r \oplus T[i][byte_i(x)]$
27:     **end for**
28:     **return** $r$
29: **end function**
30:
31: Let $L \in \{2, 3, 4, 5\}$ be the number of output words (variant-specific)
32: $stack[0 \ldots 8] \leftarrow$ empty                    ▷ Tree with fanout 8; each entry is a tuple in $\mathsf{uw}^L$
33: **for** each block group in input **do**
34:     $(y_0, \ldots, y_{L-1}) \leftarrow$ BASELAYER(block group, $s$)
35:     Push $(y_0, \ldots, y_{L-1})$ onto $stack[0]$
36:     $\ell \leftarrow 0$
37:     **while** $|stack[\ell]| = 8$ **do**                              ▷ Collapse full levels
38:         $(y_0, \ldots, y_{L-1}) \leftarrow$ UPPERLAYER$(stack[\ell], s)$
39:         Clear $stack[\ell]$; push $(y_0, \ldots, y_{L-1})$ onto $stack[\ell + 1]$
40:         $\ell \leftarrow \ell + 1$
41:     **end while**
42: **end for**
43: Initialize accumulators $a_0, \ldots, a_{L-1} \leftarrow 0$
44: Let $T^{(0)}, \ldots, T^{(L)}$ be 8×256 tables derived from entropy ($T^{(0)}$ for length, $T^{(i+1)}$ for $h_i$)
45: Mix remaining stack tuples and remaining input blocks into $(a_0, \ldots, a_{L-1})$ using MIX (tail summarized)
46: Reduce each accumulator block by summing its 64-bit lanes: $h_i \leftarrow \sum \text{lanes}(a_i)$
47: $r \leftarrow$ TABULATEBYTES$(n, T^{(0)})$
48: **for** $i \in \{0, \ldots, L-1\}$ **do**
49:     $r \leftarrow r \oplus$ TABULATEBYTES$(h_i, T^{(i+1)})$
50: **end for**
51: **return** $r$

## 12.2 khashv

Source: `hashes/khashv.cpp` [Can23].

By Keith-Cancel. A vectorizable hash using S-box byte substitution and SIMD-friendly shuffles that map to `pshufb` (SSSE3+). Uses no multiplication—only add, XOR, rotate, and table lookup. Endian-independent with compact 128-bit state (four 32-bit words, initialized from SHA-256 hashes of bytes 1–4).

Non-linear mixing uses a nibble-decomposed S-box (similar to AES SubBytes), and SHUFFLE

applies fixed byte permutation $(7, 14, 9, 0, 12, 15, 13, 8, 5, 11, 6, 3, 4, 2, 10, 1)$.

**Primitives.**    Two 16-byte lookup tables $S_1, S_2$ define the S-box:

```
 1: function SBOX(x)                              ▷ 128-bit → 128-bit, applied per byte
 2:     for each byte b in x do
 3:         b ← S₁[b ∧ 0xF] ⊕ S₂[b ≫ 4]          ▷ Low/high nibble lookup + XOR
 4:     end for
 5:     return x
 6: end function
 7: function ROTR5(v)                             ▷ Rotate 16-byte state right by 5 bytes
 8:     return bytes of v permuted: bᵢ ← b₍ᵢ₊₅₎ ₘₒ𝒹 ₁₆
 9: end function
10: function SHUFFLE(v)                           ▷ Fixed byte permutation (pshufb-style)
11:     return bytes of v permuted: bᵢ ← b_π(i) where π = (7, 14, 9, 0, 12, 15, 13, 8, 5, 11, 6, 3, 4, 2, 10, 1)
12: end function
13: function MIX_WORDS(v)                         ▷ Finalization: operates on 32-bit words
14:     v ← v ⊕ (v ≫₃₂ 3)                         ▷ Shift each 32-bit word right by 3, XOR back
15:     for r ∈ (5, 7, 11, 17) do
16:         t ← ROTR5(v) + v
17:         v ← v ⊕ ROTᵣ(t)                       ▷ Rotate each 32-bit word right by r
18:     end for
19:     return v
20: end function
```

**khashv Algorithm.**

```
 1: v ← seed                                      ▷ 128-bit state: four 32-bit words (v₀, v₁, v₂, v₃)
 2: v₀ ← v₀ ⊕ (n mod 2³²);   v₁ ← v₁ ⊕ (n ≫ 32)
 3: for each block x ∈ u128 do
 4:     t ← SBOX(x)
 5:     v ← v ⊕ (t · 8193)                        ▷ 8193 = 2¹³ + 1
 6:     v ← ROTR5(v) + t
 7:     v ← v + SHUFFLE(v)
 8: end for
 9: v ← MIX_WORDS(v)
10: return v₃ (32-bit) or (v₀, v₁) (64-bit)
```

# A    Folded Multiply $(\mathrm{lo} \oplus \mathrm{hi})$

Let $m = 2^n$. For each $x \in \{0, 1, \ldots, m - 1\}$, define

$$\mathrm{op}(x, y) = \left\lfloor \frac{xy}{m} \right\rfloor \oplus (xy \bmod m),$$

and fiber sizes

$$N_x(a) = \#\{y \in [0, m - 1] : \mathrm{op}(x, y) = a\}.$$

The collision count is

$$C(x) = \sum_{a=0}^{m-1} N_x(a)^2, \qquad \mathbb{E}[C] = \frac{1}{m} \sum_{x=0}^{m-1} C(x).$$

To bound $\mathbb{E}[C]$ we separate the two spike outputs $a = 0$ and $a = m - 1$ (the only ones amenable to exact analysis) and treat the rest as a residual term:

$$C(x) = N_x(0)^2 + N_x(m-1)^2 + C^*(x), \qquad C^*(x) = \sum_{a \notin \{0, m-1\}} N_x(a)^2.$$

Below are sharp, self-contained bounds on the **spike contribution**, plus a rigorous explanation of why "repeating multipliers" (like $0101\ldots$) cannot accumulate enough mass to change $\mathbb{E}[C]$ at the $m \log \log n$ scale.

## A.1 Exact formulas for the two spike fibers

**Lemma 1** (Exact spike counts). *For $1 \le x \le m - 1$,*

$$N_x(0) = \gcd(x, m+1), \qquad N_x(m-1) = \gcd(x, m-1).$$

*Proof.* $\mathrm{op}(x, y) = 0$ means "high = low," i.e. $q = r$ in $xy = mq + r$. Then

$$xy = mq + q = (m+1)q,$$

so $(m+1) \mid xy$. Let $d = \gcd(x, m+1)$. Then $(m+1) \mid xy$ is equivalent to $\frac{m+1}{d} \mid y$, because $\gcd(x/d, (m+1)/d) = 1$. Since $x < m$, we have $d < m + 1$, hence $(m+1)/d \ge 2$, and the multiples of $(m+1)/d$ in $[0, m-1]$ are exactly

$$0, \quad \frac{m+1}{d}, \quad 2\frac{m+1}{d}, \quad \ldots, \quad (d-1)\frac{m+1}{d},$$

whose last term is $(m+1) - (m+1)/d \le m - 1$. So there are exactly $d$ solutions $y$, i.e. $N_x(0) = d$.

$\mathrm{op}(x, y) = m - 1$ means $q \oplus r = m - 1$, which forces $r = (m-1) - q$. Then

$$xy = mq + r = mq + (m - 1 - q) = (m-1)(q+1),$$

so $(m-1) \mid xy$. Let $d = \gcd(x, m-1)$. As above, $(m-1) \mid xy$ is equivalent to $\frac{m-1}{d} \mid y$, and the multiples in $[0, m-1]$ are

$$0, \quad \frac{m-1}{d}, \quad 2\frac{m-1}{d}, \quad \ldots, \quad (d-1)\frac{m-1}{d} = m - 1 - \frac{m-1}{d} < m,$$

giving exactly $d$ solutions. So $N_x(m-1) = d$. $\qquad\square$

For $x = 0$, $\mathrm{op}(0, y) = 0$ for all $y$, so $N_0(0) = m$ and $N_0(m-1) = 0$. This single $x$ contributes only $O(m)$ to $\mathbb{E}[C]$, so it never affects asymptotics.

## A.2 Turning the spike identities into $\Theta(m \log \log n)$

We need $\mathbb{E}[N_x(0)^2]$ and $\mathbb{E}[N_x(m-1)^2]$. These are essentially averages of $\gcd(\cdot, N)^2$.

**Lemma 2** (Sum of squared gcds). *For any $N \ge 1$,*

$$\sum_{k=1}^{N} \gcd(k, N)^2 = N^2 \sum_{d \mid N} \frac{\varphi(d)}{d^2},$$

*where $\varphi$ is Euler's totient.*

*Proof.* Group $k$ by $g = \gcd(k, N)$. Writing $k = gk'$ with $\gcd(k', N/g) = 1$, there are $\varphi(N/g)$ such $k$. Hence

$$\sum_{k=1}^{N} \gcd(k, N)^2 = \sum_{g \mid N} g^2 \, \varphi(N/g) = \sum_{d \mid N} \left(\frac{N}{d}\right)^2 \varphi(d) = N^2 \sum_{d \mid N} \frac{\varphi(d)}{d^2}. \qquad\square$$

Define
$$F(N) := \sum_{d\mid N} \frac{\varphi(d)}{d^2}.$$

**Lemma 3** ($F(N)$ versus $\sigma(N)/N$). *Let $\sigma(N) = \sum_{d\mid N} d$. Then for all $N$,*

$$\frac{6}{\pi^2} \cdot \frac{\sigma(N)}{N} \;\leq\; F(N) \;\leq\; \frac{\sigma(N)}{N}.$$

*Proof.* Write $N = \prod p^{e_p}$. Both $F$ and $\sigma(\cdot)/\cdot$ factor over prime powers:

$$\frac{\sigma(N)}{N} = \prod_{p^e \| N} \left(1 + \frac{1}{p} + \cdots + \frac{1}{p^e}\right),$$

and one can compute

$$F(N) = \prod_{p^e \| N} \left(1 + \frac{1}{p} - \frac{1}{p^{e+1}}\right).$$

For each prime power,

$$1 + \frac{1}{p} - \frac{1}{p^{e+1}} \;\leq\; 1 + \frac{1}{p} + \cdots + \frac{1}{p^e},$$

and also

$$1 + \frac{1}{p} - \frac{1}{p^{e+1}} = \left(1 + \frac{1}{p} + \cdots + \frac{1}{p^e}\right) - \left(\frac{1}{p^2} + \cdots + \frac{1}{p^{e+1}}\right)$$

$$\geq \left(1 - \frac{1}{p^2}\right)\left(1 + \frac{1}{p} + \cdots + \frac{1}{p^e}\right).$$

Multiplying over primes gives

$$F(N) \;\geq\; \left(\prod_{p\mid N}\left(1 - \tfrac{1}{p^2}\right)\right)\frac{\sigma(N)}{N} \;\geq\; \left(\prod_{p}\left(1 - \tfrac{1}{p^2}\right)\right)\frac{\sigma(N)}{N} \;=\; \frac{6}{\pi^2}\frac{\sigma(N)}{N}. \qquad \square$$

**Consequence for spike second moments.**

- For $a = m - 1$: here $N = m - 1$ and $x$ runs through $1, 2, \ldots, m - 1$ exactly, so

$$\frac{1}{m}\sum_{x=1}^{m-1} N_x(m-1)^2 = \frac{1}{m}\sum_{x=1}^{m-1} \gcd(x, m-1)^2 = \frac{(m-1)^2}{m} F(m-1) = \Theta\big(\sigma(m-1)\big).$$

- For $a = 0$: $N = m+1$ but $x$ only runs up to $m - 1 = N - 2$. Dropping the final two terms only changes the sum by $O(N^2)$, which becomes $O(m)$ after dividing by $m$. So

$$\mathbb{E}[N_x(0)^2] = \Theta\big(\sigma(m+1)\big).$$

Thus the **spike part of** $\mathbb{E}[C]$ is

$$\mathbb{E}\big[N_x(0)^2 + N_x(m-1)^2\big] = \Theta\big(\sigma(2^n + 1) + \sigma(2^n - 1)\big).$$

**Upper bound $O(m \log \log n)$ for the spike part.** A theorem of Erdős [Erd71] gives

$$\sigma(2^n - 1) \;\leq\; c\,(2^n - 1) \log \log n \quad \text{for large } n,$$

for an absolute constant $c$.

To also bound $\sigma(2^n + 1)$, observe that

$$2^{2n} - 1 = (2^n - 1)(2^n + 1), \qquad \gcd(2^n - 1,\, 2^n + 1) = 1,$$

so $\sigma(2^{2n} - 1) = \sigma(2^n - 1)\,\sigma(2^n + 1)$. Applying the same bound at exponent $2n$ gives

$$\sigma(2^{2n} - 1) \leq c'(2^{2n} - 1) \log \log(2n),$$

and since $\sigma(2^n - 1) \geq 2^n - 1$, we get

$$\sigma(2^n + 1) \leq \frac{\sigma(2^{2n} - 1)}{\sigma(2^n - 1)} = O(2^n \log \log n).$$

Therefore, unconditionally,

$$\mathbb{E}\big[N_x(0)^2 + N_x(m-1)^2\big] = O(m \log \log n).$$

**Lower bound $\Omega(m \log \log n)$ for infinitely many $n$.** Take $n = t!$ and $m = 2^n$. For every odd prime $p \leq t$, we have $p - 1 \mid t!$, so by Fermat's little theorem $2^{t!} \equiv 1 \pmod{p}$, hence $p \mid (2^{t!} - 1)$.

Thus $2^{t!} - 1$ is divisible by $\prod_{p \leq t} p$, so

$$\frac{\sigma(2^{t!} - 1)}{2^{t!} - 1} \;\geq\; \prod_{p \leq t} \Big(1 + \frac{1}{p}\Big).$$

Now

$$\prod_{p \leq t} \Big(1 + \frac{1}{p}\Big) = \prod_{p \leq t} \Big(1 - \frac{1}{p^2}\Big) \cdot \prod_{p \leq t} \Big(1 - \frac{1}{p}\Big)^{-1}.$$

Mertens' theorem says $\prod_{p \leq t}(1 - \frac{1}{p}) \sim e^{-\gamma}/\log t$, i.e. the inverse grows like $e^{\gamma} \log t$. Also $\prod_{p \leq t}(1 - \frac{1}{p^2})$ stays bounded below by a positive constant (it converges to $6/\pi^2$). Hence

$$\prod_{p \leq t} \Big(1 + \frac{1}{p}\Big) \;\geq\; c \log t \quad \text{for large } t,$$

so

$$\sigma(2^{t!} - 1) \;\geq\; c\,(2^{t!} - 1) \log t.$$

Since $n = t!$, we have $\log \log n = \log \log(t!) = \Theta(\log t)$. Therefore

$$\sigma(2^n - 1) \;\geq\; c'\, 2^n \log \log n \quad \text{for infinitely many } n,$$

and since $\mathbb{E}[C] \geq \mathbb{E}[N_x(m-1)^2] = \Theta(\sigma(2^n - 1))$, we get

$$\mathbb{E}[C] = \Omega(m \log \log n) \quad \text{infinitely often.}$$

**Summary so far (rigorous).**

$$\mathbb{E}[C] \;=\; \mathbb{E}[C^*] \;+\; O(m \log \log n),$$

and also $\mathbb{E}[C] \geq \Omega(m \log \log n)$ along the subsequence $n = t!$. So the entire problem of a tight upper bound is now concentrated in $\mathbb{E}[C^*]$.

## A.3  Non-spike outputs: do repeating multipliers matter?

Outside $\{0, m-1\}$, some $x$'s have noticeably larger fibers $N_x(a)$ (e.g. $x = 0101\ldots$, which is $x = (m-1)/3$ when $n$ is even). The key point for the expectation is: **those multipliers form a tiny set**, so they cannot accumulate enough mass to change $\mathbb{E}[C]$ at the $m \log \log n$ scale.

**A. Any fixed-period family contributes only $O(m)$.**  Fix a period $p$. There are at most $2^p$ $n$-bit strings with period $p$ (choose the length-$p$ template; the rest repeats). For each $x$,

$$C(x) = \sum_a N_x(a)^2 \ \leq\ \left(\sum_a N_x(a)\right)^2 = m^2.$$

So the total contribution of all period-$p$ multipliers to the expectation is

$$\frac{1}{m} \sum_{\substack{x \text{ has} \\ \text{period } p}} C(x) \ \leq\ \frac{1}{m} \cdot 2^p \cdot m^2 \ =\ 2^p\, m.$$

For constant $p$ (like $p = 2$ for $0101\ldots$), this is $O(m)$, which is negligible compared to the spike term $m \log \log n$ once $n$ is large.

**B. The whole divisor family $x \mid (m-1)$ stays at $O(m \log \log n)$.**  This is a more relevant "big structured class," because periodic patterns like $0101\ldots$ are typically divisors of $m-1$ (or closely related).

For any $x$, a trivial but useful bound is

$$C(x) \leq mx,$$

since for $1 \leq x \leq m-1$ we have $q = \lfloor xy/m \rfloor \in \{0, 1, \ldots, x-1\}$, and for a fixed output $a = q \oplus r$ each choice of $q$ forces $r = a \oplus q$ and hence at most one $y$ solving $xy = mq + r$. Thus $N_x(a) \leq x$, so $C(x) = \sum_a N_x(a)^2 \leq (\max_a N_x(a)) \sum_a N_x(a) \leq x \cdot m$. so for the divisor family $x \mid (m-1)$,

$$\frac{1}{m} \sum_{x \mid (m-1)} C(x) \ \leq\ \frac{1}{m} \sum_{x \mid (m-1)} mx \ =\ \sigma(m-1) \ =\ O(m \log \log n),$$

using Erdős' bound on $\sigma(2^n - 1)$ [Erd71].

So **even if every single divisor of $m-1$ had maximally bad non-spike behavior (it does not), its *total* effect on $\mathbb{E}[C]$ cannot exceed the same $O(m \log \log n)$ scale that the spike already forces.

## A.4  What remains for bounding $\mathbb{E}[C]$

The sharp part of the expectation is now cleanly isolated:

- The **only provably growing term** is the spike contribution, and it is

$$\Theta\big(\sigma(2^n - 1) + \sigma(2^n + 1)\big) = O(m \log \log n),$$

  with a matching $\Omega(m \log \log n)$ lower bound infinitely often.

- Repeating multipliers and similar structured classes are provably too sparse to change the leading scale of $\mathbb{E}[C]$.

  To obtain a full $\mathbb{E}[C] = O(m \log \log n)$ upper bound, the missing ingredient is:

  Show $\mathbb{E}[C^*] = O(m)$ (or at least $O(m \log \log n)$).

Empirically $\mathbb{E}[C^*]/m$ appears bounded (it hovers around a small constant for all small $n$), but proving it seems to require a genuine "mixing" statement: roughly, that for typical $x$, the non-spike outputs behave like a near-random mapping, so $\sum_{a \neq 0, m-1} N_x(a)^2$ stays $\Theta(m)$.

## A.5 Empirics for folded multiply

The operation $\mathrm{op}(x, y) = \lfloor xy/m \rfloor \oplus (xy \bmod m)$ can also be viewed, for fixed $x$, as a map $f_x : \{0, \ldots, m-1\} \to \{0, \ldots, m-1\}$. We write $m(x) = |\mathrm{Im}(f_x)|$ for its image size.

**Image size versus collisions.** If a function $f$ on an $m$-element domain has image size $M$, then (writing $C^{\mathrm{pairs}} = \#\{(a, b) : 0 \le a < b < m, \ f(a) = f(b)\}$) one has

$$C^{\mathrm{pairs}} \ge \frac{1}{2}\left(\frac{m^2}{M} - m\right).$$

In terms of the ordered-pair statistic used above, $C(x) = \sum_a N_x(a)^2 = m + 2C^{\mathrm{pairs}}$.

**Variants compared.** We compare four closely related fold constructions (all returning a `un` value):

$$\texttt{nmul\_xor} : \mathrm{lo}(xa) \oplus \mathrm{hi}(xa), \quad \texttt{nmul\_add} : (\mathrm{lo}(xa) + \mathrm{hi}(xa)) \bmod 2^n,$$

where $xa$ is ordinary integer multiplication, and

$$\texttt{cmul\_xor} : \mathrm{lo}(x \otimes a) \oplus \mathrm{hi}(x \otimes a), \quad \texttt{cmul\_add} : (\mathrm{lo}(x \otimes a) + \mathrm{hi}(x \otimes a)) \bmod 2^n,$$

where $x \otimes a$ is carryless (polynomial) multiplication over $\mathbb{F}_2$.

**Field baselines.** In a true field, multiplication by $x \ne 0$ is a permutation, hence has no collisions. Therefore, for any field of size $N$,

$$\mathbb{E}_x[C_x] = \frac{1}{N}\binom{N}{2} = \frac{N-1}{2}, \qquad \mathbb{E}_x[C_x]/N \approx \frac{1}{2}.$$

In the plots below we include two baselines: `gf2n_field` (a field of size $2^n$, for any irreducible polynomial of degree $n$) and `prime_field` (multiplication modulo the largest prime $p \le 2^n$).

**Carryless `cmul_xor` admits an exact collision formula.** For carryless multiplication, the map $a \mapsto \mathrm{lo}(x \otimes a) \oplus \mathrm{hi}(x \otimes a)$ is multiplication by $x$ in the ring $\mathbb{F}_2[t]/(t^n + 1)$. Writing $g = \gcd(x, t^n + 1)$ and $d = \deg(g)$, one has $|\ker f_x| = 2^d$, $m(x) = 2^{n-d}$, and

$$C_x = 2^{n-1}(2^d - 1).$$

This explains the strong $n$-dependence seen for `cmul_xor`.

**Empirical results.** Figure 2 shows $\mathbb{E}[C_x]$ and its normalized version $\mathbb{E}[C_x]/2^n$ for all variants. For `nmul_xor` we observe an approximately constant factor $\mathbb{E}[C_x]/2^n \approx 1.5$–$2.2$ for $n \le 15$ (exact enumeration); for $n > 15$ we estimate using Monte Carlo sampling over multipliers $x$.

**Achievable $(\delta, \varepsilon)$ frontier.** Define the image-size threshold $M(\delta) = 2^{(1-\delta)n}$ and the "bad multiplier" event $m(x) \le M(\delta)$. Let $p_n(\delta) = \Pr_x[m(x) \le M(\delta)]$. The plot in Figure 3 shows $p_n(\delta)$ using the more stable $x$-axis $\Delta = \delta n$ (bits below $n$), which aligns curves across $n$.
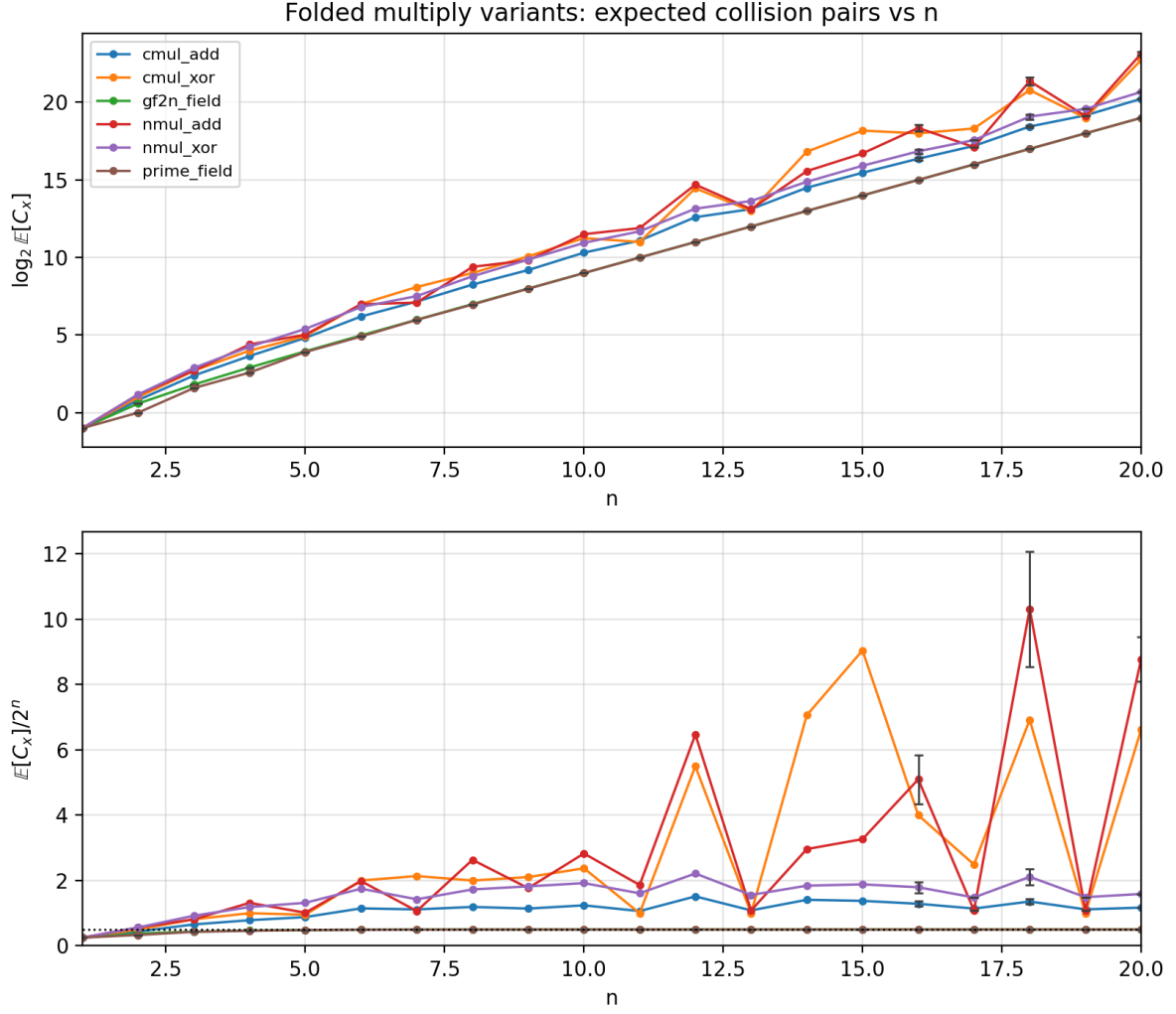
Figure 2: Expected collision pairs $\mathbb{E}[C_x]$ (top) and normalized constant $\mathbb{E}[C_x]/2^n$ (bottom) for folded multiply variants.

# References

[AB12]    Jean-Philippe Aumasson and Daniel J. Bernstein. SipHash: A fast short-input PRF. In *Progress in Cryptology – INDOCRYPT 2012*, volume 7668 of *Lecture Notes in Computer Science*, pages 489–508. Springer, 2012.

[AKT20]   Thomas Dybdahl Ahle, Jakob Tejs Baek Knudsen, and Mikkel Thorup. The power of hashing with mersenne primes. *arXiv preprint arXiv:2008.08654*, 2020.

[App21]   Jim Apple. HalftimeHash: Modern hashing without 64-bit multipliers or finite fields. In *Proceedings of the 17th International Symposium on Algorithms and Data Structures (WADS 2021)*, Lecture Notes in Computer Science. Springer, 2021.

[Can23]   Keith Cancel. k-hashv – a vectorized hash function. GitHub repository, 2023.

[Col12]   Yann Collet. xxHash – extremely fast non-cryptographic hash algorithm. GitHub repository, 2012.

[De 24]   Nicolas De Carli. rapidhash – very fast, high quality, platform-independent hashing. GitHub repository, 2024.
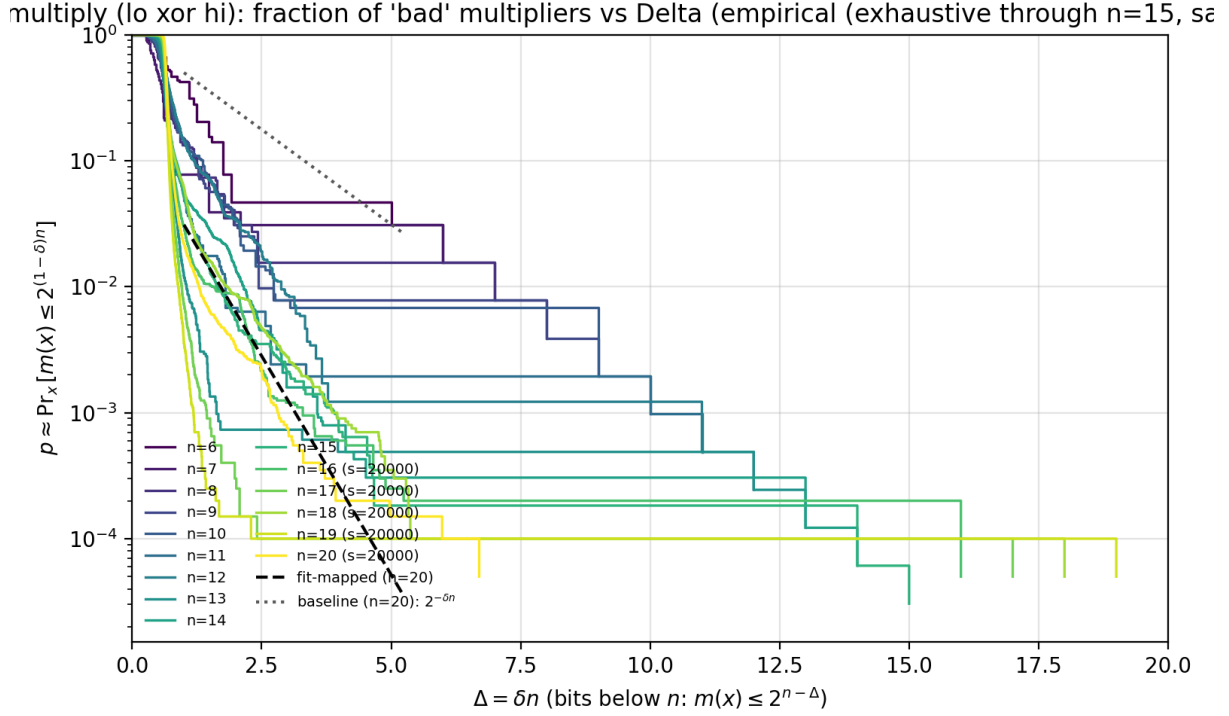
Figure 3: Fraction of multipliers with small image: $p_n(\delta) = \text{Pr}_x[m(x) \leq 2^{n-\Delta}]$, plotted against $\Delta = \delta n$ on a log-$y$ scale (exhaustive for $n \leq 15$, sampled after).

[Erd71]   Paul Erdős. On the sum $\sum_{d|2^n-1} d^{-1}$. *Israel Journal of Mathematics*, 9:43–48, 1971.

[Jen12]   Bob Jenkins. SpookyHash: A 128-bit noncryptographic hash. Author's website, 2012.

[K-A24]   K-Aethiax. MuseAir – a portable hashing algorithm. GitHub repository, 2024.

[Kai18]   Tom Kaitchuck. aHash – a non-cryptographic hashing algorithm using AES. GitHub repository, 2018.

[Khu20a]  Paul Khuong. Nearly double the PH bits with one more CLMUL. Blog post, 2020.

[Khu20b]  Paul Khuong. UMASH: Fast enough, almost universal fingerprinting. Blog post, 2020.

[LK16]    Daniel Lemire and Owen Kaser. Faster 64-bit universal hashing using carry-less multiplications. *Journal of Cryptographic Engineering*, 6:171–185, 2016.

[Mur18]   Casey Muratori. Meow hash. Project page, 2018.

[NRK24]   NRK. ChibiHash: Small, fast 64-bit hash function. Blog post, 2024.

[Pet23]   Orson Peters. PolymurHash – a universal hash function. GitHub repository, 2023.

[Pik14]   Geoff Pike. Introducing FarmHash. Google Open Source Blog, 2014.

[Rog19]   J. Andrew Rogers. AquaHash: Fast hashing with AES intrinsics. Blog post, 2019.

[Str23]   Cris Stringfellow. Rain hashes: Rainbow, Rainstorm and more. GitHub repository, 2023.

[Van20]   Aleksey Vaneev. PRVHASH – pseudo-random-value hash. GitHub repository, 2020.

[Van21]    Aleksey Vaneev. komihash – very fast, high-quality hash function. GitHub repository,
           2021.

[Van25]    Aleksey Vaneev. a5hash – ultra fast, high-quality hash functions. GitHub repository,
           2025.

[Woj22]    Frank J. T. Wojcik. SMHasher3 – hash function quality and speed test suite. GitLab
           repository, 2022.

[Yi19]     Wang Yi. wyhash – fast portable non-cryptographic hashing. GitHub repository, 2019.