

ROLV Primitive© — Complete Benchmark Portfolio

All Verified Results · March 2026 · 1,628/1,628 PASS

Summary Table — Peak Results by Platform & Model

| Hardware | Peak Speedup | At Sparsity | vs Baseline |
|-------------------------------|---------------|-------------|---------------|
| NVIDIA H200 | 13.64× | 80% | cuSPARSE |
| NVIDIA B200 | 12.06× | 70% | cuSPARSE |
| AMD Instinct MI300X | 83.77× | 85% | rocSPARSE |
| Intel CPU (synthetic) | 8.72× | 92% | CPU-CSR |
| AMD EPYC 7B13 | 5.01× | 70% | CPU-CSR |
| Google Axion ARM | 5.12× | 70% | CPU-CSR |
| Tesla T4 BF16 | 9.97× | 95% | cuSPARSE |
| LLaMA-3.1-8B real weights | 11.3× | 80% | cuSPARSE |
| LLaMA-3.1-70B shapes | 11.95× | 80% | cuSPARSE |
| LLaMA-3.1-405B shapes | 12.4× | 80% | cuSPARSE |
| HF 5-model sweep | 19.42× | 95% | cuSPARSE |
| K/V batch 32-layer | 15.62× | 80% | cuSPARSE |
| Batch scaling (2048) | 10.90× | 80% | cuSPARSE |
| BF16 dtype | 253× | 70% | cuSPARSE-BF16 |
| Qwen2.5-72B | 11.72× | 70% | cuSPARSE |
| Intel i7 LLaMA shapes | 42.4%× | 95% | CPU-CSR |
| Llama 4 Scout/Maverick | 11.91× | 70% | cuSPARSE |

| | | | |
|-------------------------|---------------|-----|----------|
| Kimi K2 | 11.90× | 70% | cuSPARSE |
| DeepSeek V3 | 11.80× | 80% | cuSPARSE |
| Mistral Large 3 | 11.68× | 70% | cuSPARSE |
| Qwen3-235B-A22B | 11.47× | 70% | cuSPARSE |
| Microsoft Phi-4 | 11.36× | 70% | cuSPARSE |
| GPT-OSS 120B/20B | 11.33× | 70% | cuSPARSE |

Total: 1,628/1,628 PASS · max error 9.87×10^{-7} · ATOL=0.05 · 4 SHA-256 hashes per case

Baseline: <70% sparsity → cuBLAS/rocBLAS/MKL (dense vendor). ≥70% → cuSPARSE/rocSPARSE/CPU-CSR (sparse vendor). Both always measured and published.

1. Synthetic Platform Sweep

Setup: 10k×10k matrix, batch=2,500, 2,000 iterations, uniform-random sparsity (worst case).

1a. NVIDIA H200

A hash: b2687223 · **V hash:** f8b47533

| Sparsity | Baseline | Vendor ms | ROLV ms |
|------------|-----------------|-------------|-------------|
| 0% | cuBLAS | 2.48 | 2.51 |
| 50% | cuBLAS | 2.48 | 1.31 |
| 70% | cuSPARSE | 4.82 | 0.68 |
| 80% | cuSPARSE | 5.90 | 0.43 |
| 90% | cuSPARSE | 3.71 | 0.28 |
| 95% | cuSPARSE | 2.02 | 0.19 |
| 99% | cuSPARSE | 0.61 | 0.08 |

1b. NVIDIA B200

A hash: 6764dac0 · **V hash:** eabab8fa

| Sparsity | Baseline | Vendor ms | ROLV ms |
|------------|-----------------|-------------|-------------|
| 0% | cuBLAS | 2.21 | 2.23 |
| 50% | cuBLAS | 2.21 | 1.18 |
| 70% | cuSPARSE | 4.31 | 0.36 |
| 80% | cuSPARSE | 3.28 | 0.31 |
| 90% | cuSPARSE | 1.84 | 0.21 |
| 95% | cuSPARSE | 0.97 | 0.14 |
| 99% | cuSPARSE | 0.29 | 0.06 |

1c. AMD Instinct MI300X

A hash: 76252923 · **V hash:** 7f9f717a

| Sparsity | Baseline | Vendor ms | ROLV ms |
|------------|------------------|--------------|-------------|
| 0% | rocBLAS | 5.86 | 5.96 |
| 50% | rocBLAS | 5.78 | 3.00 |
| 70% | rocSPARSE | 121.69 | 1.89 |
| 80% | rocSPARSE | 92.31 | 1.85 |
| 85% | rocSPARSE | 74.27 | 0.89 |
| 90% | rocSPARSE | 54.00 | 0.81 |
| 95% | rocSPARSE | 30.42 | 0.69 |

Note: rocSPARSE has a known performance regression on MI300X for this matrix topology. vs rocBLAS dense, ROLV peak is 8.5x.

1d. Intel CPU (synthetic)

Setup: 2kx2k, batch=500, 100 iters

| Sparsity | Baseline | Vendor ms | ROLV ms |
|----------|----------|-----------|---------|
|----------|----------|-----------|---------|

| | | | |
|------------|----------------|-------------|-------------|
| 0% | MKL | 0.81 | 0.82 |
| 50% | MKL | 0.81 | 0.44 |
| 70% | CPU-CSR | 1.43 | 0.28 |
| 80% | CPU-CSR | 1.61 | 0.21 |
| 92% | CPU-CSR | 1.74 | 0.20 |
| 95% | CPU-CSR | 1.68 | 0.19 |
| 99% | CPU-CSR | 1.42 | 0.18 |

1e. AMD EPYC 7B13

Setup: 2k×2k

| Sparsity | Baseline | Vendor ms | ROLV ms |
|------------|----------------|-------------|-------------|
| 0% | rocBLAS | 1.04 | 1.05 |
| 50% | rocBLAS | 1.04 | 0.57 |
| 70% | CPU-CSR | 1.88 | 0.37 |
| 80% | CPU-CSR | 1.91 | 0.44 |
| 90% | CPU-CSR | 1.76 | 0.38 |
| 99% | CPU-CSR | 1.31 | 0.34 |

1f. Google Axion ARM (Neoverse V2) — NEW

Setup: Google Cloud C4A instance · aarch64 · 3000×3000 · batch=1000 · iters=1000

A hash: 82371dc0b1b8c82d6c23a5549bd3b344c92006b2f087ef0b519d86a8e4db26fd

V hash: 2f47fc317fb5727973a89de9216cefded3c40307badb17cab7c0e06b0f39ee76

| Sparsity | Baseline | Vendor ms | ROLV ms |
|----------|----------|-----------|---------|
| 0% | OpenBLAS | 198.3 | 198.2 |
| 10% | OpenBLAS | 200.3 | 180.0 |

| | | | |
|------------|----------------|--------------|-------------|
| 30% | OpenBLAS | 198.2 | 140.5 |
| 50% | OpenBLAS | 198.1 | 102.1 |
| 60% | OpenBLAS | 199.9 | 82.0 |
| 70% | CPU-CSR | 317.8 | 62.0 |
| 80% | CPU-CSR | 210.9 | 41.8 |
| 90% | CPU-CSR | 106.4 | 21.5 |
| 95% | CPU-CSR | 53.7 | 11.9 |
| 99% | CPU-CSR | 11.3 | 3.7 |

First ROLV result on ARM architecture. Same software operator, zero changes — ARM just works.

2. LLaMA-3.1 Production Shape Sweeps — NVIDIA B200

2a. LLaMA-3.1-8B Real Weights — 60/60 PASS

Real downloaded weights, magnitude row pruning, batch=512.

| Layer | Shape | Speedup | At Sparsity |
|-------------------|-------------|--------------|-------------|
| embed_tokens | 128256×4096 | 11.3× | 80% |
| L0 mlp.gate_proj | 14336×4096 | 10.5× | 70% |
| L0 mlp.down_proj | 4096×14336 | 8.5× | 70% |
| L16 mlp.gate_proj | 14336×4096 | 9.8× | 70% |
| L31 mlp.gate_proj | 14336×4096 | 9.7× | 70% |

2b. LLaMA-3.1-8B & 70B Shape Sweep — 84/84 PASS

| Model | Layer | Shape | Peak Speedup |
|-------|---------------|-------------|---------------|
| 8B | embed_tokens | 128256×4096 | 11.27× |
| 8B | mlp.gate_proj | 14336×4096 | 10.47× |

| | | | |
|------------|----------------------|--------------------|---------------|
| 8B | mlp.down_proj | 4096×14336 | 8.47× |
| 70B | embed_tokens | 128256×8192 | 11.95× |
| 70B | mlp.gate_proj | 28672×8192 | 11.45× |
| 70B | mlp.down_proj | 8192×28672 | 10.83× |

2c. LLaMA-3.1-405B Shape Sweep — 49/49 PASS

| Layer | Shape | Peak Speedup |
|---------------|--------------|--------------|
| embed_tokens | 128256×16384 | 12.4× |
| mlp.gate_proj | 53248×16384 | 12.2× |
| mlp.down_proj | 16384×53248 | 11.9× |

Pattern confirmed: larger model = larger speedup. 8B < 70B < 405B across all layers.

3. HuggingFace Real Model Sweep — NVIDIA B200

3a. 5-Model Sweep — 96/96 PASS

| Model | Layer | Peak Speedup |
|-----------------------------|--------------|---------------|
| Mistral-7B-Instruct-v0.3 | embed_tokens | 10.50× |
| Qwen2.5-7B-Instruct | embed_tokens | 19.27× |
| DeepSeek-R1-Distill-Qwen-7B | embed_tokens | 19.42× |
| LLaMA-2-7B (GSM8K) | embed_tokens | 10.28× |

3b. Qwen2.5-72B-Instruct — 36/36 PASS

| Layer | Shape | Speedup |
|---------------|-------------|---------------|
| embed_tokens | 152064×8192 | 11.72× |
| mlp.gate_proj | 29568×8192 | 11.39× |
| mlp.up_proj | 29568×8192 | 11.38× |

| | | |
|---------------------|---------------|------------|
| embed_tokens | 11.90x | 70% |
| shared_expert.down | 8.44x | 70% |
| routed_expert.down | 8.40x | 70% |
| q_proj | 9.98x | 70% |

4c. DeepSeek V3 (671B total / 37B active) — 54/54 PASS

H=7168, l_moe=2048, KV_lora=512, VOCAB=129280

| Layer | Speedup | At Sparsity |
|---------------------|---------------|-------------|
| embed_tokens | 11.80x | 80% |
| shared_expert.down | 8.37x | 70% |
| q_proj | 9.96x | 70% |

4d. Mistral Large 3 (675B total / 41B active) — 54/54 PASS

H=7168, l_expert=4096, NQ=128, NKV=128 (full MHA), KV_lora=512, VOCAB=131072

| Layer | Speedup | At Sparsity |
|---------------------|---------------|-------------|
| embed_tokens | 11.68x | 70% |
| shared_expert.down | 9.40x | 70% |
| routed_expert.down | 9.40x | 70% |
| q_proj | 9.99x | 70% |

Larger per-expert intermediate (4096 vs 2048) yields stronger expert layer results than DeepSeek V3/Kimi K2.

4e. Qwen3-235B-A22B (235B total / 22B active) — 42/42 PASS

H=4096, l_moe=1536, NQ=64, NKV=4, VOCAB=151936, Apache 2.0

| Layer | Speedup | At Sparsity |
|---------------------|---------------|-------------|
| embed_tokens | 11.47x | 70% |

| | | |
|------------------|-------|-----|
| expert.down_proj | 4.66x | 70% |
| q_proj | 6.19x | 70% |

Small expert intermediate (1536) by design — 128 experts × 1536 = large total capacity.

4f. Microsoft Phi-4 (14B dense) — 42/42 PASS

H=5120, I=17920, NQ=40, NKV=10, VOCAB=100352, MIT license

| Layer | Speedup | At Sparsity |
|---------------------|---------------|-------------|
| embed_tokens | 11.36x | 70% |
| mlp.gate_proj | 9.34x | 70% |
| mlp.up_proj | 9.34x | 70% |
| mlp.down_proj | 9.14x | 70% |
| q_proj | 7.00x | 70% |

Only dense model in this batch — no MoE. Strong MLP numbers confirm ROLV works equally well on dense decoder-only architectures.

4g. GPT-OSS 120B/20B (OpenAI) — 42/42 PASS

H=2880, I_expert=2880, NQ=64, NKV=8, Q_dim=4096, VOCAB=201088, Apache 2.0

| Layer | Speedup | At Sparsity |
|---------------------|---------------|-------------|
| embed_tokens | 11.33x | 70% |
| expert.gate_proj | 5.55x | 70% |
| expert.down_proj | 5.53x | 70% |
| q_proj | 5.45x | 70% |

Smallest H in the portfolio (2880). Expert layer numbers reflect compact per-expert matrices.

5. Supplementary Benchmarks — NVIDIA B200

5a. Tesla T4 BF16 — 24/24 PASS

| Layer | Sparsity | Speedup |
|---------------------|------------|--------------|
| embed_tokens (BF16) | 0% | 1.89× |
| embed_tokens (BF16) | 95% | 9.97× |

5b. BF16 Dtype Sweep — 70/70 PASS

| Layer | vs cuBLAS-BF16 @ 80% | vs cuBLAS-BF16 |
|-------------------|----------------------|----------------|
| 70B embed_tokens | 3.39× | 5.76× |
| 70B mlp.gate_proj | 3.19× | 4.25× |
| 70B mlp.down_proj | 2.85× | 4.04× |

cuSPARSE-BF16 has a known performance regression on B200. Honest comparison is vs cuBLAS-BF16.

5c. K/V Layer Batching (GQA) — 20/20 PASS

| Layers batched | Shape | Speedup vs cuSPARSE |
|----------------|-------------------|---------------------|
| 1 | 512×3584 | 0.86× |
| 8 | 4096×3584 | 5.85× |
| 16 | 8192×3584 | 9.00× |
| 32 | 16384×3584 | 15.62× |

5d. Batch Size Scaling — 9/9 PASS

| Batch | Speedup vs cuSPARSE | ROLV ms |
|-------------|---------------------|--------------|
| 1 | 1.24× | 0.029 |
| 64 | 7.92× | 0.072 |
| 512 | 7.92× | 0.289 |
| 2048 | 10.90× | 0.834 |

5e. Sparsity Structure Comparison — 12/12 PASS

| Structure | Speedup range |
|------------------|---------------|
| Uniform random | 1.0× |
| Power-law rows | 7.6–9.2× |
| Block structured | 7.8–9.4× |

6. Methodology

Correctness standard

- **ATOL = 0.05** on column-normalised outputs
- **Perturbation test:** modify one weight → output hash changes → confirms live computation
- **max error across all 1142 cases: 9.87×10^{-7}**
- **4 SHA-256 hashes per case:** weight matrix (A), input vector (V), vendor output, ROLV output

Baseline selection

Below 70% sparsity: cuBLAS/rocBLAS/MKL (dense vendor — optimal for lightly sparse weights).
Above 70%: cuSPARSE/rocSPARSE/CPU-CSR (sparse vendor — production inference engines). Both always published.

Energy measurement

GPU: pynvml (NVIDIA) / pyrsmi (AMD) — polls power rail every 15ms, integrates joules.

CPU: proxy estimate based on core count × TDP × CPU utilisation.

Zero-trust verification

Any result can be independently verified via huggingface.co/spaces/rolvai/rolv-verify. Supply your own input numbers, verify the output matches dense on a calculator — no trust required.

ROLV Primitive© · CRCS™ · RSMT™ · ROLVswitch™ · 3 Patents Pending · rolv.ai · rolv@rolv.ai

AMD Instinct MI300X — Complete Model Portfolio

Hardware: AMD Instinct MI300X · ROCm/HIP · batch=512 · 200 iters · ATOL=0.05

Baselines: rocBLAS (dense, honest comparison) · rocSPARSE (sparse library — known regression on MI300X at high sparsity, published for transparency)

Correctness: 4 SHA-256 hashes + perturbation test per case · max error 9.87×10^{-7} · 486/486 PASS

| Model | Cases | PASS | Peak vs rocBLA |
|--------------------------|------------|------------|-----------------|
| LLaMA-3.1-405B shapes | 36 | 36 | 13.53x ★ |
| Llama 4 Scout & Maverick | 54 | 54 | 12.25x |
| LLaMA-3.1 8B & 70B | 72 | 72 | 11.88x |
| DeepSeek V3 | 54 | 54 | 11.66x |
| Qwen3-235B-A22B | 42 | 42 | 11.49x |
| Qwen2.5-72B | 36 | 36 | 10.46x |
| Mistral Large 3 | 54 | 54 | 10.22x |
| Kimi K2 | 54 | 54 | 9.75x |
| Microsoft Phi-4 | 42 | 42 | 9.59x |
| GPT-OSS 120B/20B | 42 | 42 | 9.04x |
| Total | 486 | 486 | 13.53x |

★ = portfolio peak across all platforms and models.

Key AMD observations

LLaMA-3.1-405B embed_tokens at 95% sparsity — the peak row:

| Baseline | Time | ROLV |
|-----------|----------|--------|
| rocBLAS | 20.05ms | 1.48ms |
| rocSPARSE | 109.69ms | 1.48ms |

LLaMA-3.1-70B MLP layers at 80% sparsity on MI300X vs rocBLAS:

| Layer | Shape | rocBLAS ms | ROLV ms |
|---------------|-------------|------------|---------|
| embed_tokens | 128256×8192 | 9.82 | 2.24 |
| mlp.gate_proj | 28672×8192 | 1.91 | 0.71 |

| | | | |
|---------------|------------|------|------|
| mlp.up_proj | 28672×8192 | 1.98 | 0.70 |
| mlp.down_proj | 8192×28672 | 2.33 | 1.95 |

LLaMA scaling on AMD confirms B200 pattern (larger model = larger speedup):

| | | |
|----------------|--|---------------|
| Model | | embed peak vs |
| LLaMA-3.1-8B | | 9.83× |
| LLaMA-3.1-70B | | 11.88× |
| LLaMA-3.1-405B | | 13.53× |

rocSPARSE regression note

rocSPARSE on MI300X has a known performance regression on large square matrices at high sparsity — it loads and processes the full matrix regardless of sparsity structure, collapsing its own throughput. ROLV skips zero rows entirely before any memory is loaded. Both baselines are published. vs rocBLAS is the honest production comparison.

Complete Portfolio Summary — All Platforms

Grand total: 1,628/1,628 PASS

| Platform | Cases | Peak speedup |
|----------------------------------|-------|--------------|
| NVIDIA H200 | 22 | 13.64× |
| NVIDIA B200 (synthetic) | 22 | 12.06× |
| NVIDIA B200 (model sweeps) | 582 | 19.42× |
| NVIDIA B200 (LLaMA real weights) | 60 | 11.3× |
| Tesla T4 BF16 | 24 | 9.97× |
| Intel CPU | 22 | 8.72× |
| AMD EPYC 7B13 | 22 | 5.01× |
| Google Axion ARM | 22 | 5.12× |

AMD MI300X

486

13.53x

Total

1,628

ROLV Primitive© · CRCS™ · RSMT™ · ROLVswitch™ · 3 Patents Pending · rolv.ai · rolv@rolv.ai

7. AMD MI300X Model Portfolio — Complete Hash Tables

All 4 SHA-256 hashes per case. hash_A = weight matrix, hash_V = input vector, hash_D = dense baseline output, hash_R = ROLV output.

To verify: reproduce the weight matrix (deterministic seed published per harness), compute SHA-256, compare to hash_A. Then run dense matmul on your hardware — if hash_D matches, the baseline is confirmed. hash_R proves ROLV output.

If hash_D = hash_R the matrix was zero-sparse at that sparsity level (all active rows) — ROLV and dense are numerically identical.

LLaMA-3.1-405B shapes (AMD MI300X) — 36/36 PASS

Peak vs rocBLAS: **13.53x** · Peak vs rocSPARSE: **74.02x** · Max error: 5.72e-06

| Layer | Sp% | M×K | hash_A | hash_V |
|---------------|-----|--------------|------------|------------|
| embed_tokens | 0% | 128256×16384 | `2a78e62a` | `1a6f2514` |
| embed_tokens | 50% | 128256×16384 | `edab215f` | `1a6f2514` |
| embed_tokens | 70% | 128256×16384 | `0ca785b4` | `1a6f2514` |
| embed_tokens | 80% | 128256×16384 | `68d12e57` | `1a6f2514` |
| embed_tokens | 90% | 128256×16384 | `2dcc2b14` | `1a6f2514` |
| embed_tokens | 95% | 128256×16384 | `2bb52131` | `1a6f2514` |
| mlp.gate_proj | 0% | 53248×16384 | `ac4afc65` | `d25c7a95` |
| mlp.gate_proj | 50% | 53248×16384 | `78f017e6` | `d25c7a95` |
| mlp.gate_proj | 70% | 53248×16384 | `ee278764` | `d25c7a95` |

| | | | | |
|---------------|-----|-------------|------------|------------|
| mlp.gate_proj | 80% | 53248×16384 | `9793a119` | `d25c7a95` |
| mlp.gate_proj | 90% | 53248×16384 | `7d96daf4` | `d25c7a95` |
| mlp.gate_proj | 95% | 53248×16384 | `cb08894c` | `d25c7a95` |
| mlp.up_proj | 0% | 53248×16384 | `f975a7cb` | `7b84cdb6` |
| mlp.up_proj | 50% | 53248×16384 | `7d8ec58c` | `7b84cdb6` |
| mlp.up_proj | 70% | 53248×16384 | `3d3fad69` | `7b84cdb6` |
| mlp.up_proj | 80% | 53248×16384 | `23d99ebc` | `7b84cdb6` |
| mlp.up_proj | 90% | 53248×16384 | `94f70bf1` | `7b84cdb6` |
| mlp.up_proj | 95% | 53248×16384 | `6a7d5923` | `7b84cdb6` |
| mlp.down_proj | 0% | 16384×53248 | `74841e63` | `f82de250` |
| mlp.down_proj | 50% | 16384×53248 | `48b1c89b` | `f82de250` |
| mlp.down_proj | 70% | 16384×53248 | `cd22dc30` | `f82de250` |
| mlp.down_proj | 80% | 16384×53248 | `107d5756` | `f82de250` |
| mlp.down_proj | 90% | 16384×53248 | `b83453a3` | `f82de250` |
| mlp.down_proj | 95% | 16384×53248 | `8bb79003` | `f82de250` |
| q_proj | 0% | 16384×16384 | `4f7728fc` | `74d6a9e5` |
| q_proj | 50% | 16384×16384 | `8ad15c31` | `74d6a9e5` |
| q_proj | 70% | 16384×16384 | `bd91640c` | `74d6a9e5` |
| q_proj | 80% | 16384×16384 | `f7608f7d` | `74d6a9e5` |
| q_proj | 90% | 16384×16384 | `c2905092` | `74d6a9e5` |
| q_proj | 95% | 16384×16384 | `a456687b` | `74d6a9e5` |

| | | | | |
|--------------|-----|------------|------------|------------|
| k_proj (GQA) | 0% | 1024×16384 | `ef6eac30` | `694b1c02` |
| k_proj (GQA) | 50% | 1024×16384 | `16b11922` | `694b1c02` |
| k_proj (GQA) | 70% | 1024×16384 | `fa9f4028` | `694b1c02` |
| k_proj (GQA) | 80% | 1024×16384 | `a935592b` | `694b1c02` |
| k_proj (GQA) | 90% | 1024×16384 | `759c95a8` | `694b1c02` |
| k_proj (GQA) | 95% | 1024×16384 | `71177e83` | `694b1c02` |

DeepSeek V3 (AMD MI300X) — 54/54 PASS

Peak vs rocBLAS: **11.66x** · Peak vs rocSPARSE: **64.95x** · Max error: 9.78e-06

| Layer | Sp% | M×K | hash_A | hash_V |
|--------------------|-----|-------------|------------|------------|
| embed_tokens | 0% | 129280×7168 | `5053dd9b` | `69dbc495` |
| embed_tokens | 50% | 129280×7168 | `9e2bb044` | `69dbc495` |
| embed_tokens | 70% | 129280×7168 | `8fcd5c98` | `69dbc495` |
| embed_tokens | 80% | 129280×7168 | `180d74ea` | `69dbc495` |
| embed_tokens | 90% | 129280×7168 | `18642e8b` | `69dbc495` |
| embed_tokens | 95% | 129280×7168 | `ee4b203b` | `69dbc495` |
| shared_expert.gate | 0% | 2048×7168 | `0a934337` | `a85b0797` |
| shared_expert.gate | 50% | 2048×7168 | `590e2667` | `a85b0797` |
| shared_expert.gate | 70% | 2048×7168 | `1d708e8a` | `a85b0797` |
| shared_expert.gate | 80% | 2048×7168 | `26413455` | `a85b0797` |
| shared_expert.gate | 90% | 2048×7168 | `bad50d1f` | `a85b0797` |
| shared_expert.gate | 95% | 2048×7168 | `1631db82` | `a85b0797` |
| shared_expert.up | 0% | 2048×7168 | `2c83b94f` | `5fc82022` |

| | | | | |
|--------------------|-----|-----------|------------|------------|
| shared_expert.up | 50% | 2048x7168 | `b3be81f8` | `5fc82022` |
| shared_expert.up | 70% | 2048x7168 | `c61f626a` | `5fc82022` |
| shared_expert.up | 80% | 2048x7168 | `4a9e0a0d` | `5fc82022` |
| shared_expert.up | 90% | 2048x7168 | `207d634b` | `5fc82022` |
| shared_expert.up | 95% | 2048x7168 | `41cb3d65` | `5fc82022` |
| shared_expert.down | 0% | 7168x2048 | `b767372b` | `e6c6cc91` |
| shared_expert.down | 50% | 7168x2048 | `eb73f69a` | `e6c6cc91` |
| shared_expert.down | 70% | 7168x2048 | `0d14bd44` | `e6c6cc91` |
| shared_expert.down | 80% | 7168x2048 | `927f1b3b` | `e6c6cc91` |
| shared_expert.down | 90% | 7168x2048 | `cd9b15e2` | `e6c6cc91` |
| shared_expert.down | 95% | 7168x2048 | `4d785363` | `e6c6cc91` |
| routed_expert.gate | 0% | 2048x7168 | `aa9a07d6` | `4f122bfa` |
| routed_expert.gate | 50% | 2048x7168 | `1a0405f6` | `4f122bfa` |
| routed_expert.gate | 70% | 2048x7168 | `f78b6bb7` | `4f122bfa` |
| routed_expert.gate | 80% | 2048x7168 | `bd600267` | `4f122bfa` |
| routed_expert.gate | 90% | 2048x7168 | `555961df` | `4f122bfa` |
| routed_expert.gate | 95% | 2048x7168 | `52a324f0` | `4f122bfa` |
| routed_expert.up | 0% | 2048x7168 | `3ef1fae0` | `e9a1793e` |
| routed_expert.up | 50% | 2048x7168 | `be538ecd` | `e9a1793e` |
| routed_expert.up | 70% | 2048x7168 | `403ca8c9` | `e9a1793e` |
| routed_expert.up | 80% | 2048x7168 | `3a08d9d2` | `e9a1793e` |

| | | | | |
|--------------------|-----|-----------|------------|------------|
| routed_expert.up | 90% | 2048x7168 | `59c57ffa` | `e9a1793e` |
| routed_expert.up | 95% | 2048x7168 | `a6f69a30` | `e9a1793e` |
| routed_expert.down | 0% | 7168x2048 | `babc3da7` | `8211f6d5` |
| routed_expert.down | 50% | 7168x2048 | `cd6379ad` | `8211f6d5` |
| routed_expert.down | 70% | 7168x2048 | `d4a58e50` | `8211f6d5` |
| routed_expert.down | 80% | 7168x2048 | `4209c576` | `8211f6d5` |
| routed_expert.down | 90% | 7168x2048 | `d50cb144` | `8211f6d5` |
| routed_expert.down | 95% | 7168x2048 | `2272943c` | `8211f6d5` |
| q_proj | 0% | 7168x7168 | `bfa2bc73` | `1fad8fc2` |
| q_proj | 50% | 7168x7168 | `8da02d67` | `1fad8fc2` |
| q_proj | 70% | 7168x7168 | `a8a68959` | `1fad8fc2` |
| q_proj | 80% | 7168x7168 | `fb00bd3a` | `1fad8fc2` |
| q_proj | 90% | 7168x7168 | `8dde0808` | `1fad8fc2` |
| q_proj | 95% | 7168x7168 | `3ffb948f` | `1fad8fc2` |
| kv_proj (MLA) | 0% | 512x7168 | `9a82b3f3` | `e58a3471` |
| kv_proj (MLA) | 50% | 512x7168 | `b3d052de` | `e58a3471` |
| kv_proj (MLA) | 70% | 512x7168 | `abaa55b2` | `e58a3471` |
| kv_proj (MLA) | 80% | 512x7168 | `001aee66` | `e58a3471` |
| kv_proj (MLA) | 90% | 512x7168 | `dc74e3dd` | `e58a3471` |
| kv_proj (MLA) | 95% | 512x7168 | `247721fa` | `e58a3471` |

Kimi K2 (AMD MI300X) — 54/54 PASS

Peak vs rocBLAS: **9.74x** · Peak vs rocSPARSE: **58.25x** · Max error: 1.06e-05

| Layer | Sp% | M×K | hash_A | hash_V |
|--------------------|-----|-------------|------------|------------|
| embed_tokens | 0% | 102400×7168 | `01228e68` | `69dbc495` |
| embed_tokens | 50% | 102400×7168 | `d5749b67` | `69dbc495` |
| embed_tokens | 70% | 102400×7168 | `1a0071be` | `69dbc495` |
| embed_tokens | 80% | 102400×7168 | `218e088a` | `69dbc495` |
| embed_tokens | 90% | 102400×7168 | `bbcaf0ae` | `69dbc495` |
| embed_tokens | 95% | 102400×7168 | `13f962c2` | `69dbc495` |
| shared_expert.gate | 0% | 2048×7168 | `b80f938d` | `a85b0797` |
| shared_expert.gate | 50% | 2048×7168 | `85aaf3fc` | `a85b0797` |
| shared_expert.gate | 70% | 2048×7168 | `9269917c` | `a85b0797` |
| shared_expert.gate | 80% | 2048×7168 | `c2c0edba` | `a85b0797` |
| shared_expert.gate | 90% | 2048×7168 | `c246e686` | `a85b0797` |
| shared_expert.gate | 95% | 2048×7168 | `1cd4077a` | `a85b0797` |
| shared_expert.up | 0% | 2048×7168 | `509d5350` | `5fc82022` |
| shared_expert.up | 50% | 2048×7168 | `9352b22e` | `5fc82022` |
| shared_expert.up | 70% | 2048×7168 | `267442bb` | `5fc82022` |
| shared_expert.up | 80% | 2048×7168 | `3a2f1809` | `5fc82022` |
| shared_expert.up | 90% | 2048×7168 | `04254f4a` | `5fc82022` |
| shared_expert.up | 95% | 2048×7168 | `98efa8f6` | `5fc82022` |
| shared_expert.down | 0% | 7168×2048 | `d35d40d8` | `e6c6cc91` |
| shared_expert.down | 50% | 7168×2048 | `f7672f79` | `e6c6cc91` |

| | | | | |
|--------------------|-----|-----------|------------|------------|
| shared_expert.down | 70% | 7168x2048 | `2e1dea91` | `e6c6cc91` |
| shared_expert.down | 80% | 7168x2048 | `d43cd2d4` | `e6c6cc91` |
| shared_expert.down | 90% | 7168x2048 | `802e3f33` | `e6c6cc91` |
| shared_expert.down | 95% | 7168x2048 | `833971fe` | `e6c6cc91` |
| routed_expert.gate | 0% | 2048x7168 | `5e755824` | `4f122bfa` |
| routed_expert.gate | 50% | 2048x7168 | `26d93250` | `4f122bfa` |
| routed_expert.gate | 70% | 2048x7168 | `1101e613` | `4f122bfa` |
| routed_expert.gate | 80% | 2048x7168 | `c5ef225d` | `4f122bfa` |
| routed_expert.gate | 90% | 2048x7168 | `71b73be4` | `4f122bfa` |
| routed_expert.gate | 95% | 2048x7168 | `3903dd1e` | `4f122bfa` |
| routed_expert.up | 0% | 2048x7168 | `9231b16b` | `e9a1793e` |
| routed_expert.up | 50% | 2048x7168 | `e090bb0d` | `e9a1793e` |
| routed_expert.up | 70% | 2048x7168 | `9e0b3b85` | `e9a1793e` |
| routed_expert.up | 80% | 2048x7168 | `bd467055` | `e9a1793e` |
| routed_expert.up | 90% | 2048x7168 | `77ff7878` | `e9a1793e` |
| routed_expert.up | 95% | 2048x7168 | `c1757413` | `e9a1793e` |
| routed_expert.down | 0% | 7168x2048 | `6367e6ce` | `8211f6d5` |
| routed_expert.down | 50% | 7168x2048 | `f578e797` | `8211f6d5` |
| routed_expert.down | 70% | 7168x2048 | `5ece33a6` | `8211f6d5` |
| routed_expert.down | 80% | 7168x2048 | `1dd916c5` | `8211f6d5` |
| routed_expert.down | 90% | 7168x2048 | `d8754730` | `8211f6d5` |

| | | | | |
|--------------------|-----|-----------|------------|------------|
| routed_expert.down | 95% | 7168×2048 | `5263d49f` | `8211f6d5` |
| q_proj | 0% | 7168×7168 | `74cd32e6` | `1fad8fc2` |
| q_proj | 50% | 7168×7168 | `9a2839c3` | `1fad8fc2` |
| q_proj | 70% | 7168×7168 | `45534d05` | `1fad8fc2` |
| q_proj | 80% | 7168×7168 | `9dc515b0` | `1fad8fc2` |
| q_proj | 90% | 7168×7168 | `674b75aa` | `1fad8fc2` |
| q_proj | 95% | 7168×7168 | `aad8e28` | `1fad8fc2` |
| k_proj (MLA) | 0% | 2048×7168 | `ce707c01` | `e58a3471` |
| k_proj (MLA) | 50% | 2048×7168 | `5fbe576a` | `e58a3471` |
| k_proj (MLA) | 70% | 2048×7168 | `4d792708` | `e58a3471` |
| k_proj (MLA) | 80% | 2048×7168 | `fc48b6da` | `e58a3471` |
| k_proj (MLA) | 90% | 2048×7168 | `49efaca9` | `e58a3471` |
| k_proj (MLA) | 95% | 2048×7168 | `f4ec7db0` | `e58a3471` |

LLaMA-3.1 8B & 70B shapes (AMD MI300X) — 72/72 PASS

Peak vs rocBLAS: **11.88x** · Peak vs rocSPARSE: **66.50x** · Max error: 5.96e-06

| Layer | Sp% | M×K | hash_A | hash_V |
|-----------------|-----|-------------|------------|------------|
| 8B embed_tokens | 0% | 128256×4096 | `815c2c0c` | `027724f0` |
| 8B embed_tokens | 50% | 128256×4096 | `e419e324` | `027724f0` |
| 8B embed_tokens | 70% | 128256×4096 | `85b28872` | `027724f0` |
| 8B embed_tokens | 80% | 128256×4096 | `2a1c37a4` | `027724f0` |
| 8B embed_tokens | 90% | 128256×4096 | `42d5b897` | `027724f0` |
| 8B embed_tokens | 95% | 128256×4096 | `dce68556` | `027724f0` |

| | | | | |
|------------------|-----|------------|------------|------------|
| 8B mlp.gate_proj | 0% | 14336×4096 | `457d9256` | `3c44a290` |
| 8B mlp.gate_proj | 50% | 14336×4096 | `425bb5fa` | `3c44a290` |
| 8B mlp.gate_proj | 70% | 14336×4096 | `c7c8cb2a` | `3c44a290` |
| 8B mlp.gate_proj | 80% | 14336×4096 | `d5e6f479` | `3c44a290` |
| 8B mlp.gate_proj | 90% | 14336×4096 | `17473f37` | `3c44a290` |
| 8B mlp.gate_proj | 95% | 14336×4096 | `dfa1dcbb` | `3c44a290` |
| 8B mlp.up_proj | 0% | 14336×4096 | `4a72417a` | `b67d5af6` |
| 8B mlp.up_proj | 50% | 14336×4096 | `0d1ac0b4` | `b67d5af6` |
| 8B mlp.up_proj | 70% | 14336×4096 | `e16f8a35` | `b67d5af6` |
| 8B mlp.up_proj | 80% | 14336×4096 | `a64fa4ee` | `b67d5af6` |
| 8B mlp.up_proj | 90% | 14336×4096 | `6b6e81d1` | `b67d5af6` |
| 8B mlp.up_proj | 95% | 14336×4096 | `ea5f97a9` | `b67d5af6` |
| 8B mlp.down_proj | 0% | 4096×14336 | `67c39416` | `4fa17fab` |
| 8B mlp.down_proj | 50% | 4096×14336 | `0a443011` | `4fa17fab` |
| 8B mlp.down_proj | 70% | 4096×14336 | `3290e673` | `4fa17fab` |
| 8B mlp.down_proj | 80% | 4096×14336 | `b9a2c665` | `4fa17fab` |
| 8B mlp.down_proj | 90% | 4096×14336 | `02997acb` | `4fa17fab` |
| 8B mlp.down_proj | 95% | 4096×14336 | `643d8543` | `4fa17fab` |
| 8B q_proj | 0% | 4096×4096 | `98a6e664` | `e6220b85` |
| 8B q_proj | 50% | 4096×4096 | `140bb0ac` | `e6220b85` |
| 8B q_proj | 70% | 4096×4096 | `279548de` | `e6220b85` |

| | | | | |
|-------------------|-----|-------------|------------|------------|
| 8B q_proj | 80% | 4096×4096 | `a9fb698a` | `e6220b85` |
| 8B q_proj | 90% | 4096×4096 | `31fb6e2e` | `e6220b85` |
| 8B q_proj | 95% | 4096×4096 | `f776cef1` | `e6220b85` |
| 8B k_proj (GQA) | 0% | 1024×4096 | `9d6055f3` | `e9a5f677` |
| 8B k_proj (GQA) | 50% | 1024×4096 | `f9fec4b9` | `e9a5f677` |
| 8B k_proj (GQA) | 70% | 1024×4096 | `8b74cbd8` | `e9a5f677` |
| 8B k_proj (GQA) | 80% | 1024×4096 | `25399eb9` | `e9a5f677` |
| 8B k_proj (GQA) | 90% | 1024×4096 | `db737103` | `e9a5f677` |
| 8B k_proj (GQA) | 95% | 1024×4096 | `b28d26c9` | `e9a5f677` |
| 70B embed_tokens | 0% | 128256×8192 | `52e2a9f6` | `0cb76dd9` |
| 70B embed_tokens | 50% | 128256×8192 | `ce29263a` | `0cb76dd9` |
| 70B embed_tokens | 70% | 128256×8192 | `63fbbd89` | `0cb76dd9` |
| 70B embed_tokens | 80% | 128256×8192 | `aeae879d` | `0cb76dd9` |
| 70B embed_tokens | 90% | 128256×8192 | `f6ad909d` | `0cb76dd9` |
| 70B embed_tokens | 95% | 128256×8192 | `499b29dc` | `0cb76dd9` |
| 70B mlp.gate_proj | 0% | 28672×8192 | `c2368f02` | `a173466f` |
| 70B mlp.gate_proj | 50% | 28672×8192 | `a449752d` | `a173466f` |
| 70B mlp.gate_proj | 70% | 28672×8192 | `efa24abc` | `a173466f` |
| 70B mlp.gate_proj | 80% | 28672×8192 | `e32972b6` | `a173466f` |
| 70B mlp.gate_proj | 90% | 28672×8192 | `e9fa1486` | `a173466f` |
| 70B mlp.gate_proj | 95% | 28672×8192 | `55d04915` | `a173466f` |

| | | | | |
|-------------------|-----|------------|------------|------------|
| 70B mlp.up_proj | 0% | 28672×8192 | `d593f3e6` | `6169d11f` |
| 70B mlp.up_proj | 50% | 28672×8192 | `86be0d91` | `6169d11f` |
| 70B mlp.up_proj | 70% | 28672×8192 | `bf2fb278` | `6169d11f` |
| 70B mlp.up_proj | 80% | 28672×8192 | `86263da4` | `6169d11f` |
| 70B mlp.up_proj | 90% | 28672×8192 | `e361bde8` | `6169d11f` |
| 70B mlp.up_proj | 95% | 28672×8192 | `05eb0aee` | `6169d11f` |
| 70B mlp.down_proj | 0% | 8192×28672 | `31a9d942` | `36aef758` |
| 70B mlp.down_proj | 50% | 8192×28672 | `8a093d94` | `36aef758` |
| 70B mlp.down_proj | 70% | 8192×28672 | `399fb536` | `36aef758` |
| 70B mlp.down_proj | 80% | 8192×28672 | `4524ddf3` | `36aef758` |
| 70B mlp.down_proj | 90% | 8192×28672 | `fe0d5c1f` | `36aef758` |
| 70B mlp.down_proj | 95% | 8192×28672 | `c85d19c0` | `36aef758` |
| 70B q_proj | 0% | 8192×8192 | `14b732dc` | `9393c47f` |
| 70B q_proj | 50% | 8192×8192 | `bcea25e7` | `9393c47f` |
| 70B q_proj | 70% | 8192×8192 | `e021c78a` | `9393c47f` |
| 70B q_proj | 80% | 8192×8192 | `7c91e648` | `9393c47f` |
| 70B q_proj | 90% | 8192×8192 | `5224f169` | `9393c47f` |
| 70B q_proj | 95% | 8192×8192 | `9c3b44d0` | `9393c47f` |
| 70B k_proj (GQA) | 0% | 1024×8192 | `083d3b76` | `274dfb0d` |
| 70B k_proj (GQA) | 50% | 1024×8192 | `a245c3fe` | `274dfb0d` |
| 70B k_proj (GQA) | 70% | 1024×8192 | `562bf320` | `274dfb0d` |

| | | | | |
|------------------|-----|-----------|------------|------------|
| 70B k_proj (GQA) | 80% | 1024×8192 | `e13fdc34` | `274dfb0d` |
| 70B k_proj (GQA) | 90% | 1024×8192 | `09ac9ffd` | `274dfb0d` |
| 70B k_proj (GQA) | 95% | 1024×8192 | `46d03bfb` | `274dfb0d` |

Mistral Large 3 (AMD MI300X) — 54/54 PASS

Peak vs rocBLAS: **10.22x** · Peak vs rocSPARSE: **65.44x** · Max error: 5.96e-06

| Layer | Sp% | M×K | hash_A | hash_V |
|--------------------|-----|-------------|------------|------------|
| embed_tokens | 0% | 131072×7168 | `9609a208` | `69dbc495` |
| embed_tokens | 50% | 131072×7168 | `48a34a8a` | `69dbc495` |
| embed_tokens | 70% | 131072×7168 | `14de77b0` | `69dbc495` |
| embed_tokens | 80% | 131072×7168 | `35f0abd8` | `69dbc495` |
| embed_tokens | 90% | 131072×7168 | `6e7b5a2a` | `69dbc495` |
| embed_tokens | 95% | 131072×7168 | `9114e078` | `69dbc495` |
| shared_expert.gate | 0% | 4096×7168 | `ec89d874` | `a85b0797` |
| shared_expert.gate | 50% | 4096×7168 | `e850eb4e` | `a85b0797` |
| shared_expert.gate | 70% | 4096×7168 | `29fe2f24` | `a85b0797` |
| shared_expert.gate | 80% | 4096×7168 | `fa709ec2` | `a85b0797` |
| shared_expert.gate | 90% | 4096×7168 | `34a9a60b` | `a85b0797` |
| shared_expert.gate | 95% | 4096×7168 | `6b9a250f` | `a85b0797` |
| shared_expert.up | 0% | 4096×7168 | `8d5cc0ef` | `5fc82022` |
| shared_expert.up | 50% | 4096×7168 | `79c529e4` | `5fc82022` |
| shared_expert.up | 70% | 4096×7168 | `a3302d5f` | `5fc82022` |
| shared_expert.up | 80% | 4096×7168 | `53d997ff` | `5fc82022` |

| | | | | |
|--------------------|-----|-----------|------------|------------|
| shared_expert.up | 90% | 4096x7168 | `e6ebb000` | `5fc82022` |
| shared_expert.up | 95% | 4096x7168 | `82b8101e` | `5fc82022` |
| shared_expert.down | 0% | 7168x4096 | `f94ce35c` | `e6220b85` |
| shared_expert.down | 50% | 7168x4096 | `ec1737ff` | `e6220b85` |
| shared_expert.down | 70% | 7168x4096 | `81f65c26` | `e6220b85` |
| shared_expert.down | 80% | 7168x4096 | `f21035ea` | `e6220b85` |
| shared_expert.down | 90% | 7168x4096 | `84014323` | `e6220b85` |
| shared_expert.down | 95% | 7168x4096 | `49512f11` | `e6220b85` |
| routed_expert.gate | 0% | 4096x7168 | `f688afe1` | `4f122bfa` |
| routed_expert.gate | 50% | 4096x7168 | `8a9e2ad0` | `4f122bfa` |
| routed_expert.gate | 70% | 4096x7168 | `3113ee90` | `4f122bfa` |
| routed_expert.gate | 80% | 4096x7168 | `4e36cbe0` | `4f122bfa` |
| routed_expert.gate | 90% | 4096x7168 | `edf0c18b` | `4f122bfa` |
| routed_expert.gate | 95% | 4096x7168 | `dc1a320a` | `4f122bfa` |
| routed_expert.up | 0% | 4096x7168 | `101dff56` | `e9a1793e` |
| routed_expert.up | 50% | 4096x7168 | `f2359093` | `e9a1793e` |
| routed_expert.up | 70% | 4096x7168 | `b209a204` | `e9a1793e` |
| routed_expert.up | 80% | 4096x7168 | `6da39c50` | `e9a1793e` |
| routed_expert.up | 90% | 4096x7168 | `bfb67328` | `e9a1793e` |
| routed_expert.up | 95% | 4096x7168 | `807f0cd4` | `e9a1793e` |
| routed_expert.down | 0% | 7168x4096 | `4f9dc5e3` | `a63b0380` |

| | | | | |
|--------------------|-----|-----------|------------|------------|
| routed_expert.down | 50% | 7168×4096 | `a6ebb099` | `a63b0380` |
| routed_expert.down | 70% | 7168×4096 | `12d04624` | `a63b0380` |
| routed_expert.down | 80% | 7168×4096 | `4e1a86f2` | `a63b0380` |
| routed_expert.down | 90% | 7168×4096 | `4159f7c0` | `a63b0380` |
| routed_expert.down | 95% | 7168×4096 | `ad3c7226` | `a63b0380` |
| q_proj | 0% | 7168×7168 | `6415ce57` | `1fad8fc2` |
| q_proj | 50% | 7168×7168 | `e49df1cc` | `1fad8fc2` |
| q_proj | 70% | 7168×7168 | `66f9e19f` | `1fad8fc2` |
| q_proj | 80% | 7168×7168 | `0d02516a` | `1fad8fc2` |
| q_proj | 90% | 7168×7168 | `73b6a7dd` | `1fad8fc2` |
| q_proj | 95% | 7168×7168 | `12627791` | `1fad8fc2` |
| kv_lora (MLA) | 0% | 512×7168 | `1f3d8dca` | `e58a3471` |
| kv_lora (MLA) | 50% | 512×7168 | `40fd5e2b` | `e58a3471` |
| kv_lora (MLA) | 70% | 512×7168 | `2a3acafd` | `e58a3471` |
| kv_lora (MLA) | 80% | 512×7168 | `5f45cc3e` | `e58a3471` |
| kv_lora (MLA) | 90% | 512×7168 | `5ee31179` | `e58a3471` |
| kv_lora (MLA) | 95% | 512×7168 | `91603a7a` | `e58a3471` |

Microsoft Phi-4 (AMD MI300X) — 42/42 PASS

Peak vs rocBLAS: **9.59x** · Peak vs rocSPARSE: **61.22x** · Max error: 4.29e-06

| Layer | Sp% | M×K | hash_A | hash_V |
|--------------|-----|-------------|------------|------------|
| embed_tokens | 0% | 100352×5120 | `fd2caf5e` | `fc45cb70` |
| embed_tokens | 50% | 100352×5120 | `ba7d728d` | `fc45cb70` |

| | | | | |
|---------------|-----|-------------|------------|------------|
| embed_tokens | 70% | 100352×5120 | `1a9a0912` | `fc45cb70` |
| embed_tokens | 80% | 100352×5120 | `949724eb` | `fc45cb70` |
| embed_tokens | 90% | 100352×5120 | `b1014eba` | `fc45cb70` |
| embed_tokens | 95% | 100352×5120 | `8fc50b2a` | `fc45cb70` |
| mlp.gate_proj | 0% | 17920×5120 | `b39ac916` | `43c67f51` |
| mlp.gate_proj | 50% | 17920×5120 | `619ab4bc` | `43c67f51` |
| mlp.gate_proj | 70% | 17920×5120 | `13b8931e` | `43c67f51` |
| mlp.gate_proj | 80% | 17920×5120 | `8d4c5c5a` | `43c67f51` |
| mlp.gate_proj | 90% | 17920×5120 | `53118f2f` | `43c67f51` |
| mlp.gate_proj | 95% | 17920×5120 | `eb3db8eb` | `43c67f51` |
| mlp.up_proj | 0% | 17920×5120 | `f9dc87db` | `53b95cff` |
| mlp.up_proj | 50% | 17920×5120 | `3b01c83a` | `53b95cff` |
| mlp.up_proj | 70% | 17920×5120 | `6356c9fa` | `53b95cff` |
| mlp.up_proj | 80% | 17920×5120 | `6b3d24a8` | `53b95cff` |
| mlp.up_proj | 90% | 17920×5120 | `8404282e` | `53b95cff` |
| mlp.up_proj | 95% | 17920×5120 | `bdf89d53` | `53b95cff` |
| mlp.down_proj | 0% | 5120×17920 | `ab8f0695` | `9435b794` |
| mlp.down_proj | 50% | 5120×17920 | `04efa8a4` | `9435b794` |
| mlp.down_proj | 70% | 5120×17920 | `ee8e0389` | `9435b794` |
| mlp.down_proj | 80% | 5120×17920 | `a2d78f2b` | `9435b794` |
| mlp.down_proj | 90% | 5120×17920 | `9963735d` | `9435b794` |

| | | | | |
|---------------|-----|------------|------------|------------|
| mlp.down_proj | 95% | 5120×17920 | `7a0a3889` | `9435b794` |
| q_proj | 0% | 5120×5120 | `14f84bf1` | `d4a8bdb3` |
| q_proj | 50% | 5120×5120 | `073b09fe` | `d4a8bdb3` |
| q_proj | 70% | 5120×5120 | `59b62b9f` | `d4a8bdb3` |
| q_proj | 80% | 5120×5120 | `4d271b75` | `d4a8bdb3` |
| q_proj | 90% | 5120×5120 | `b9868281` | `d4a8bdb3` |
| q_proj | 95% | 5120×5120 | `92fa7741` | `d4a8bdb3` |
| k_proj (GQA) | 0% | 1280×5120 | `11e3aa14` | `e691d491` |
| k_proj (GQA) | 50% | 1280×5120 | `fff7a6e1` | `e691d491` |
| k_proj (GQA) | 70% | 1280×5120 | `9f04923e` | `e691d491` |
| k_proj (GQA) | 80% | 1280×5120 | `25e45c56` | `e691d491` |
| k_proj (GQA) | 90% | 1280×5120 | `f5e774b9` | `e691d491` |
| k_proj (GQA) | 95% | 1280×5120 | `d4b2dfdc` | `e691d491` |
| v_proj (GQA) | 0% | 1280×5120 | `6b11ad5b` | `b7eff3d3` |
| v_proj (GQA) | 50% | 1280×5120 | `c78aa488` | `b7eff3d3` |
| v_proj (GQA) | 70% | 1280×5120 | `c58600ce` | `b7eff3d3` |
| v_proj (GQA) | 80% | 1280×5120 | `63ea97e7` | `b7eff3d3` |
| v_proj (GQA) | 90% | 1280×5120 | `bd5d8b6e` | `b7eff3d3` |
| v_proj (GQA) | 95% | 1280×5120 | `bf3949cf` | `b7eff3d3` |

Qwen3-235B-A22B (AMD MI300X) — 42/42 PASS

Peak vs rocBLAS: **11.45x** · Peak vs rocSPARSE: **69.18x** · Max error: 1.00e-05

| Layer | Sp% | M×K | hash_A | hash_V |
|-------|-----|-----|--------|--------|
|-------|-----|-----|--------|--------|

| | | | | |
|------------------|-----|-------------|------------|------------|
| embed_tokens | 0% | 151936×4096 | `7f9ad33d` | `027724f0` |
| embed_tokens | 50% | 151936×4096 | `67e6cead` | `027724f0` |
| embed_tokens | 70% | 151936×4096 | `fac83cdc` | `027724f0` |
| embed_tokens | 80% | 151936×4096 | `ddc8ac44` | `027724f0` |
| embed_tokens | 90% | 151936×4096 | `8fe9b085` | `027724f0` |
| embed_tokens | 95% | 151936×4096 | `50e06a2b` | `027724f0` |
| expert.gate_proj | 0% | 1536×4096 | `264c4643` | `3c44a290` |
| expert.gate_proj | 50% | 1536×4096 | `419bb5ec` | `3c44a290` |
| expert.gate_proj | 70% | 1536×4096 | `b977770c` | `3c44a290` |
| expert.gate_proj | 80% | 1536×4096 | `67dd3e7b` | `3c44a290` |
| expert.gate_proj | 90% | 1536×4096 | `cbbdd5db` | `3c44a290` |
| expert.gate_proj | 95% | 1536×4096 | `13e4ad64` | `3c44a290` |
| expert.up_proj | 0% | 1536×4096 | `46fe3209` | `b67d5af6` |
| expert.up_proj | 50% | 1536×4096 | `c75ad326` | `b67d5af6` |
| expert.up_proj | 70% | 1536×4096 | `b3f2c2fe` | `b67d5af6` |
| expert.up_proj | 80% | 1536×4096 | `3dadcc4d` | `b67d5af6` |
| expert.up_proj | 90% | 1536×4096 | `9f4277d0` | `b67d5af6` |
| expert.up_proj | 95% | 1536×4096 | `190db103` | `b67d5af6` |
| expert.down_proj | 0% | 4096×1536 | `db417f56` | `42737bbf` |
| expert.down_proj | 50% | 4096×1536 | `a3a83236` | `42737bbf` |
| expert.down_proj | 70% | 4096×1536 | `db04f9f5` | `42737bbf` |

| | | | | |
|------------------|-----|-----------|------------|------------|
| expert.down_proj | 80% | 4096×1536 | `6ba39be4` | `42737bbf` |
| expert.down_proj | 90% | 4096×1536 | `996aefb0` | `42737bbf` |
| expert.down_proj | 95% | 4096×1536 | `e93e76ff` | `42737bbf` |
| q_proj | 0% | 4096×4096 | `fd64fff8` | `0ecf7912` |
| q_proj | 50% | 4096×4096 | `463e31d9` | `0ecf7912` |
| q_proj | 70% | 4096×4096 | `c25909d0` | `0ecf7912` |
| q_proj | 80% | 4096×4096 | `80f48bb4` | `0ecf7912` |
| q_proj | 90% | 4096×4096 | `cb8d4e25` | `0ecf7912` |
| q_proj | 95% | 4096×4096 | `5213c5d8` | `0ecf7912` |
| k_proj (GQA) | 0% | 512×4096 | `c8b9ded9` | `55fe90b8` |
| k_proj (GQA) | 50% | 512×4096 | `6b03d034` | `55fe90b8` |
| k_proj (GQA) | 70% | 512×4096 | `5304eb51` | `55fe90b8` |
| k_proj (GQA) | 80% | 512×4096 | `d353bf0d` | `55fe90b8` |
| k_proj (GQA) | 90% | 512×4096 | `73c927bc` | `55fe90b8` |
| k_proj (GQA) | 95% | 512×4096 | `8e9ce393` | `55fe90b8` |
| v_proj (GQA) | 0% | 512×4096 | `e730970a` | `e6220b85` |
| v_proj (GQA) | 50% | 512×4096 | `a2eb36b5` | `e6220b85` |
| v_proj (GQA) | 70% | 512×4096 | `768e15ef` | `e6220b85` |
| v_proj (GQA) | 80% | 512×4096 | `d1f62033` | `e6220b85` |
| v_proj (GQA) | 90% | 512×4096 | `a64f47d0` | `e6220b85` |
| v_proj (GQA) | 95% | 512×4096 | `d751be4a` | `e6220b85` |

Peak vs rocBLAS: **10.46x** · Peak vs rocSPARSE: **69.09x** · Max error: 2.79e-05

| Layer | Sp% | M×K | hash_A | hash_V |
|---------------|-----|-------------|------------|------------|
| embed_tokens | 0% | 152064×8192 | `7307a608` | `a1352977` |
| embed_tokens | 50% | 152064×8192 | `092bb98d` | `a1352977` |
| embed_tokens | 70% | 152064×8192 | `ff66989c` | `a1352977` |
| embed_tokens | 80% | 152064×8192 | `1ff1f4c1` | `a1352977` |
| embed_tokens | 90% | 152064×8192 | `eb83b0b7` | `a1352977` |
| embed_tokens | 95% | 152064×8192 | `9bae8aaa` | `a1352977` |
| mlp.gate_proj | 0% | 29568×8192 | `e5b2b14b` | `2ecf7668` |
| mlp.gate_proj | 50% | 29568×8192 | `f2966f1e` | `2ecf7668` |
| mlp.gate_proj | 70% | 29568×8192 | `662d2bbd` | `2ecf7668` |
| mlp.gate_proj | 80% | 29568×8192 | `a041766f` | `2ecf7668` |
| mlp.gate_proj | 90% | 29568×8192 | `e696601e` | `2ecf7668` |
| mlp.gate_proj | 95% | 29568×8192 | `2c687ed9` | `2ecf7668` |
| mlp.up_proj | 0% | 29568×8192 | `df4992c6` | `ab9af207` |
| mlp.up_proj | 50% | 29568×8192 | `f1825617` | `ab9af207` |
| mlp.up_proj | 70% | 29568×8192 | `481493d5` | `ab9af207` |
| mlp.up_proj | 80% | 29568×8192 | `5e6ff53b` | `ab9af207` |
| mlp.up_proj | 90% | 29568×8192 | `f50475d8` | `ab9af207` |
| mlp.up_proj | 95% | 29568×8192 | `e68390e9` | `ab9af207` |
| mlp.down_proj | 0% | 8192×29568 | `11833250` | `6565c90d` |

| | | | | |
|---------------|-----|------------|------------|------------|
| mlp.down_proj | 50% | 8192x29568 | `cd25f284` | `6565c90d` |
| mlp.down_proj | 70% | 8192x29568 | `f97797cc` | `6565c90d` |
| mlp.down_proj | 80% | 8192x29568 | `7bf6c5b9` | `6565c90d` |
| mlp.down_proj | 90% | 8192x29568 | `e4873060` | `6565c90d` |
| mlp.down_proj | 95% | 8192x29568 | `5675e570` | `6565c90d` |
| q_proj | 0% | 8192x8192 | `eaecf5ec` | `35aa6a4c` |
| q_proj | 50% | 8192x8192 | `41f7ec98` | `35aa6a4c` |
| q_proj | 70% | 8192x8192 | `1f74b19e` | `35aa6a4c` |
| q_proj | 80% | 8192x8192 | `16725c1a` | `35aa6a4c` |
| q_proj | 90% | 8192x8192 | `5e3ddcb7` | `35aa6a4c` |
| q_proj | 95% | 8192x8192 | `a86626e8` | `35aa6a4c` |
| k_proj (GQA) | 0% | 1024x8192 | `b6e67546` | `b20132a9` |
| k_proj (GQA) | 50% | 1024x8192 | `bb279476` | `b20132a9` |
| k_proj (GQA) | 70% | 1024x8192 | `c695c680` | `b20132a9` |
| k_proj (GQA) | 80% | 1024x8192 | `64bb8315` | `b20132a9` |
| k_proj (GQA) | 90% | 1024x8192 | `cd57f425` | `b20132a9` |
| k_proj (GQA) | 95% | 1024x8192 | `5fd30646` | `b20132a9` |
