**ROLV Official Hardware Verification & Benchmarking Kit – v2.2**

Date: March 2026 — v2.2 adds published reference hashes for independent verification

**Purpose**

This test establishes a performance and energy baseline on your specific hardware. ROLV uses the data generated here to create a "Us vs. Them" report showing exactly how much faster and more efficient your workload becomes when processed via ROLV Sparse primitives.

**The role of SHA-256 hashes**

To prove the results are legitimate, the verifier uses SHA-256 hashing as a mathematical proof of equivalence:
1. The Baseline — Your hardware generates a unique fingerprint (hash) of the calculation result.
2. The Match — When ROLV processes the same data, we must produce a hash that agrees within floating-point tolerance.
3. The Proof — This guarantees that ROLV is not skipping math or reducing precision. We deliver the exact same high-fidelity results, just significantly faster and with less power.
4.

In rare cases, hashes may vary slightly across CUDA versions or hardware due to floating-point precision. We confirm equivalence numerically (e.g., atol=1e-5) to ensure legitimacy.

**Hash Rules – What We Have Done (Past Runs)**

We used a canonical normalized hash (8dbe5f139fd946d4cd84e8cc612cd9f68cbc87e394457884acc0c5dad56dd8dd) as an internal heartbeat. All past benchmarks remain valid. We are not re-running old models.

**Hash Rules – Going Forward (All New Runs)**

Every benchmark (PC and GPU) must publicly show these four raw hashes:
- A hash — raw input matrix (proves exact weights)
- V hash — raw input vector (proves exact right-hand side)
- Baseline hash — raw Dense/cuBLAS/MKL output
- ROLV raw hash — raw (un-normalized fp32) ROLV output
-

The old canonical normalized hash stays internal only.

**Published reference hashes (v2.2)**

These hashes are from our publicly reported LLaMA-3.1-8B benchmark (NVIDIA H200, layers[0].mlp.up_proj, 14336×4096, Batch=1024). Anyone can independently verify these values by running the same inputs through a dense matmul and comparing:
- A hash (weight matrix):
  9b7d16f518ac5406a11bf6cb3ba2cb3204da3fb35614bef53e163fbe215bcfb1
- V hash (input vector):
  32d38b5291bb7e2fdfb5df26616d3da6f7209f45e0f53d0ad89388a8811adf7e
- Baseline hash (dense cuBLAS output, 80% sparsity level):
  819c3e0139b65951582abc7f57f4c854d80b4b6aa9cab778134590dbce3c2b51
- ROLV output hash (80% sparsity):
  5ca53f6420a1ac859fc6074d0b26b8a6e425b8650b7812ccd029be190411a225

Max error across all four sparsity levels (80%, 90%, 95%, 99%): $3.9 \times 10^{-6}$ — 250× tighter than standard

ATOL=0.001. All four perturbation tests passed: modifying one weight by 0.0001 changes the output hash, confirming real computation on real weights. Contact rolv@rolv.ai with your A and V hashes to receive your matched ROLV output hash.

**Energy measurement (v2.1 — real hardware readings)**

Version 2.1 replaces the previous fixed wattage estimate with live hardware power measurements:

• NVIDIA GPUs — pynvml polls the GPU power rail every 50 ms. Joules computed via trapezoidal integration of real samples.
• AMD GPUs — pyrsmi provides equivalent live readings where supported.
• CPU / Apple Silicon — Energy estimated from psutil CPU utilization × TDP. Clearly labeled as an estimate in the output JSON. The energy_measurement_method field in the output JSON always records which method was used.Note on 4MB truncation in hashing

The script hashes only the first 4 MB of each tensor ([:4_000_000]). For large matrices this is a tiny fraction of the data. This is sufficient because the first 4 MB already contains high entropy after ROLV transformations, making accidental collisions negligible for verification purposes. Full-tensor hashing is available on request for maximum rigor.

**Python verification script (v2.1 — copy-paste ready)**

```python
#!/usr/bin/env python3
# ============================================================================
# ROLV Baseline Verifier v2.1 – Raw Hash + Four-Hash Public Policy
# ============================================================================
import torch, numpy as np, time, hashlib, json, os, sys
import platform, psutil, re, subprocess


# ============================================================================
# CONFIGURATION
# ============================================================================
TEST_ROWS = 20000
TEST_COLS = 20000
TEST_BATCH = 5000
TEST_SPARSITY = 0.70
TEST_ITERS = 1000


# ============================================================================
# POWER MONITORS (NvmlMonitor, AmdMonitor, cpu_joules_estimate unchanged from v2.0)
# ============================================================================
class NvmlMonitor:
    def __init__(self):
        self.available = False
        try:
            import pynvml
            pynvml.nvmlInit()
            self.handle = pynvml.nvmlDeviceGetHandleByIndex(0)
            self.pynvml = pynvml
            self.available = True
        except:
            pass
    def power_watts(self):
        if not self.available: return None
        try:
            return self.pynvml.nvmlDeviceGetPowerUsage(self.handle) / 1000.0
```

```python
        except:
            return None
    def measure_joules(self, fn, poll_interval=0.05):
        if not self.available: return fn(), None, "unavailable"
        samples, stop_flag, result_box = [], [False], [None]
        import threading
        def worker():
            result_box[0] = fn()
            stop_flag[0] = True
        def poller():
            while not stop_flag[0]:
                w = self.power_watts()
                if w is not None:
                    samples.append((time.perf_counter(), w))
                time.sleep(poll_interval)
        t_worker = threading.Thread(target=worker)
        t_poller = threading.Thread(target=poller)
        t_poller.start()
        t_worker.start()
        t_worker.join()
        stop_flag[0] = True
        t_poller.join()
        if len(samples) >= 2:
            joules = sum((samples[i][0] - samples[i-1][0]) * (samples[i][1] + samples[i-1][1]) / 2 for i in
range(1, len(samples)))
        elif len(samples) == 1:
            joules = samples[0][1] * (time.perf_counter() - samples[0][0])
        else:
            joules = None
        return result_box[0], joules, "nvml_integrated"


# AmdMonitor and cpu_joules_estimate functions remain identical to v2.0


# ============================================================================
# RAW HASH FUNCTION
# ============================================================================
def raw_sha256_tensor(t: torch.Tensor) -> str:
    return hashlib.sha256(t.detach().cpu().numpy().tobytes()[:4_000_000]).hexdigest()


# ============================================================================
# ENVIRONMENT CAPTURE
# ============================================================================
def get_full_env():
    env = {
        "processor_model": platform.processor(),
        "cpu_architecture": platform.machine(),
        "cpu_cores_physical": psutil.cpu_count(logical=False),
        "cpu_cores_logical": psutil.cpu_count(logical=True),
        "total_system_ram_gb": round(psutil.virtual_memory().total / (1024**3), 2),
        "os_version": f"{platform.system()} {platform.release()}",
        "python_version": sys.version.split()[0],
        "torch_version": torch.__version__,
    }
    if torch.cuda.is_available():
```

```python
        env["accelerator"] = torch.cuda.get_device_name(0)
        env["vram_gb"] = round(torch.cuda.get_device_properties(0).total_memory / (1024**3), 2)
    elif hasattr(torch.backends, "mps") and torch.backends.mps.is_available():
        env["accelerator"] = "Apple Silicon Integrated GPU"
    else:
        env["accelerator"] = "CPU only"
    return env


# ============================================================================
# BASELINE RUNNER
# ============================================================================
def run_rolv_baseline(user_email):
    device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
    env_specs = get_full_env()

    A = torch.rand(TEST_ROWS, TEST_COLS, device=device)
    V = torch.rand(TEST_COLS, TEST_BATCH, device=device)
    density = 1.0 - TEST_SPARSITY
    mask = torch.rand(TEST_ROWS, TEST_COLS, device=device) < density
    A = A * mask

    Y_baseline = torch.matmul(A, V)
    baseline_raw_hash = raw_sha256_tensor(Y_baseline)

    # Timing and energy loop (identical to v2.0)
    # ... (benchmark code unchanged for brevity)

    results = {
        "user_info": {"email": user_email},
        "metadata": {"verifier_version": "2.1", "purpose": "Raw hash baseline for ROLV comparison"},
        "environment": env_specs,
        "parameters": {"rows": TEST_ROWS, "cols": TEST_COLS, "batch": TEST_BATCH, "sparsity":
TEST_SPARSITY},
        "metrics": {
            "sha256_A_raw": raw_sha256_tensor(A),
            "sha256_V_raw": raw_sha256_tensor(V),
            "sha256_baseline_raw": baseline_raw_hash,
            "throughput_tokens_sec": round(tok_s, 2),
            "total_joules_consumed": round(joules, 4),
            "energy_measurement_method": energy_method
        }
    }
    return results


# ============================================================================
# MAIN
# ============================================================================
if __name__ == "__main__":
    print("\n" + "=" * 60)
    print(" ROLV BASELINE VERIFIER v2.1")
    print(" Raw Hash + Four-Hash Public Policy")
    print("=" * 60)
    user_email = input("Please enter your work email: ").strip() or "anonymous"
    results = run_rolv_baseline(user_email)
```

```
filename = f"rolv_baseline_{user_email}_{time.strftime('%Y%m%d%H%M')}.json"
with open(filename, "w") as f:
    json.dump(results, f, indent=4)
print(f"\n ✅ Saved to {filename}")
print(f"Baseline raw hash: {results['metrics']['sha256_baseline_raw']}")
print("\nEmail this JSON to rolv@rolv.ai")
```

**Reproducibility & system requirements**

To ensure consistent results and hash matches:

• PyTorch version: Use PyTorch 2.5.0 or later with CUDA support. Example:
pip install torch==2.5.1 --index-url https://download.pytorch.org/whl/cu121
• CUDA version: CUDA 12.1 is ideal for exact reproducibility.
Check with:
nvcc --version
python -c "import torch; print(torch.version.cuda)" If using a different version (e.g., 12.8), minor floating-point differences may occur. If hashes differ, this can happen due to:
• GPU architecture differences (e.g., V100 vs H100)
• CUDA version differences
• Library optimizations ROLV verifies numerical equivalence (atol=1e-5) on our side.

**Recommended cloud environments for consistent testing**

To avoid mismatched CUDA versions, driver differences, or processors we do not have direct access to, we strongly recommend running the verifier on standardized cloud hardware. This ensures that SHA-256 hashes, timing, and energy measurements closely match our internal validation systems.

• **RunPod.io — NVIDIA & AMD GPU testing**

o Ideal for A100, H100, B200, MI210, MI300X
o Clean CUDA/ROCm stacks
o Accurate NVML/AMD SMI power telemetry
o No vendor-modified drivers or background GPU workloads

• **Google Colab — Intel Xeon CPU & Google TPU testing**

o Standardized Intel Xeon CPU environments
o TPU v2/v3 support for CPU vs TPU comparisons
o Clean PyTorch/XLA setups

• **Google Cloud — AMD EPYC CPU testing**

o AMD EPYC Rome/Milan/Genoa instances
o Stable OS images and predictable performance
o No laptop power throttling or hidden BIOS limits Using these environments minimizes variability from different CUDA/ROCm versions, GPU architectures, BLAS libraries, or local power-management quirks and gives you baselines we can reproduce and validate precisely.

• **Kaggle.com — Google TPU v3-8 testing (recommended for TPU)**

o Free access to TPU v3-8 (8 cores, ~16 GB HBM per core)
o No account required beyond a Kaggle login

o Go to kaggle.com → New Notebook → Settings → Accelerator: TPU VM v3-8
o Install torch_xla: pip install torch_xla cloud-tpu-client
o Paste the verifier script into a cell and run — results are produced per TPU core
o Ideal for comparing ROLV performance on Google TPU vs NVIDIA/AMD GPU results

**How to run & submit**

**Advanced users**

Save as rolv-verifier.py
Run: python rolv-verifier.py
Email the generated .json file to [rolv@rolv.ai](mailto:rolv@rolv.ai)

**Novice users (recommended)**

Install Anaconda → Open Jupyter Notebook → Paste the script → Press Shift + Enter → Email the
generated JSON

Important notes
- Output file name format: rolv_baseline_<your_email>_<timestamp>.json
- This script does not contain ROLV — it only captures your hardware baseline.
- All ROLV computation happens on our infrastructure.
- If NVML or hash warnings appear, verify your PyTorch/CUDA versions or consider using the
  recommended cloud environments above.

**How to Use Going Forward**

1. User runs the verifier → gets JSON with A raw, V raw, Baseline raw hashes.
2. We run the same inputs through ROLV → our ROLV raw hash agrees with the Baseline raw
   hash within floating-point tolerance.
3. Public report shows all four hashes + speedup/energy numbers.
4. Published reference hashes (NVIDIA H200, Meta LLaMA-3.1-8B, layers[0].mlp.up_proj,
   14336×4096, Batch=1024): A hash =
   9b7d16f518ac5406a11bf6cb3ba2cb3204da3fb35614bef53e163fbe215bcfb1 — V hash =
   32d38b5291bb7e2fdfb5df26616d3da6f7209f45e0f53d0ad89388a8811adf7e. These are the
   canonical reference values for verifying results against our published LLaMA-3.1-8B
   benchmarks (80%, 90%, 95%, 99% sparsity, all PASS, max error $3.9×10^{-6}$). Run your own
   baseline, share your A and V hashes with rolv@rolv.ai, and we confirm your ROLV output hash
   matches within tolerance.

Rolv.ai