

TurboQuant: Implementation Corrections, Production Hardening, and Deployment Infrastructure

A Systematic Analysis of the Reference Code for arXiv:2504.19874

Tushar Pathare and Team¹

¹High Performance Computing Division, EviOX Tech, Pune, India tushar@eviox.tech
<https://www.eviox.tech>

EviOX Tech Technical Report | v1.2.0 | March 2026

Abstract

TurboQuant Zandieh et al. [2025] is a recent online vector quantization algorithm from Google Research, NYU, and Google DeepMind that achieves near-optimal distortion within a factor of $\approx 2.7\times$ of the information-theoretic lower bound, with zero indexing time and full GPU vectorization. This report documents six material correctness bugs identified by EviOX Tech in the paper’s reference Python implementation, accompanied by rigorous corrections, numerical verification against all three theorems of the paper, and a production inference server built on the corrected core.

The most significant finding is a critical mathematical error in the Quantized Johnson-Lindenstrauss (QJL) dequantization scale factor: the reference code used $\sqrt{\pi/(2d)}$, whereas Definition 1 of the paper specifies $\sqrt{\pi/2}/d$. These expressions differ by a factor of $\sqrt{d} \approx 11.3$ at $d = 128$, which renders the unbiasedness guarantee of Theorem 2 invalid and corrupts all inner-product quantization results. Additional bugs include incorrect sub-vector dimensions in TurboQuantProd, uniform use of head-0 quantizer across all attention heads, a violated HuggingFace Cache.update() return contract, per-batch outlier channel recomputation, and codebook recomputation at every class instantiation.

Beyond bug fixes, EviOX contributes 11 verified implementations of paper claims, 8 production hardening improvements, and 2 novel infrastructure additions: a 156-test suite with 60/60 paper coverage verification, and a full CI/CD pipeline. All corrections are numerically verified against the distortion bounds of Theorems 1, 2, and 3.

Keywords: vector quantization, KV cache compression, large language models, QJL transform, Lloyd-Max quantizer, HPC inference, CUDA, GPU inference server

Contents

1	Introduction	3
1.1	Scope and Methodology	3
1.2	Summary of Findings	3
2	Background: TurboQuant Algorithm	3
2.1	Notation	3
2.2	Key Results from the Paper	4
3	Mathematical Foundations: Verified Implementations and Bugs	4
3.1	Item 1 — Lemma 1: Beta PDF (VERIFIED)	4
3.2	Item 2 — Lemma 4: QJL Scale Factor (FIXED)	4
3.3	Item 3 — Lemmas 2–3: Lower Bounds (VERIFIED)	5

4	Algorithm 1 (TurboQuant_{mse}): Corrections and Verifications	5
4.1	Item 4 — Random Rotation via QR (VERIFIED)	5
4.2	Item 5 — Per-Head Unique Rotation Matrices (FIXED)	5
4.3	Item 6 — Lloyd-Max Codebook Caching (FIXED)	6
4.4	Item 7 — Theorem 1 Distortion Constant (VERIFIED)	6
4.5	Item 8 — float16 Argmin Precision (FIXED)	7
4.6	Item 9 — Inverse Rotation in Dequantization (VERIFIED)	7
5	Algorithm 2 (TurboQuant_{prod}): Corrections and Verifications	7
5.1	Item 10 — Two-Stage Pipeline (VERIFIED)	7
5.2	Item 11 — Sub-Vector Dimension Error (FIXED)	7
5.3	Item 12 — MSE Bias at $b = 1$ (VERIFIED)	8
5.4	Item 13 — Theorem 2 Unbiasedness (FIXED , consequence of Bug 1)	8
5.5	Item 14 — Residual Norm $\gamma = \ r\ _2$ (VERIFIED)	8
6	Theorem 3 Lower Bounds (VERIFIED)	8
7	KV Cache Outlier Channel Split: Corrections	8
7.1	Item 17 — Bit Allocation Arithmetic (FIXED)	8
7.2	Item 18 — Calibration Hook Cleanup (FIXED)	9
7.3	Items 19–20 — $b_{\text{out}}/b_{\text{reg}}$ Assignment and Streaming (VERIFIED)	9
8	HuggingFace Cache Integration: Critical Bugs	9
8.1	Item 21 — Variable-Depth Needle Construction (FIXED)	9
8.2	Item 22 — Cache.update() Return Contract (FIXED)	10
8.3	Item 23 — Per-Head Dequantization (FIXED)	10
8.4	Items 24–26 — Compatibility, Data-Obliviousness, Atomicity (HARDENED , VERIFIED)	11
9	Production Server Hardening	11
10	Test Infrastructure and CI/CD	11
10.1	Item 33 — 156-Test Suite with 60/60 Paper Coverage (ADDED)	11
10.2	Item 34 — CI/CD Pipeline (ADDED)	11
11	Full Correction Table	12
12	Distortion Verification Summary	14
13	Discussion	14
13.1	Critical Finding: QJL Scale Error	14
13.2	Compound Effect of Bugs	15
13.3	Implications for Reproducibility	15
14	Conclusion	15
	Software Repository and Extended Citations	17
A	Reproduction Instructions	17

1 Introduction

TurboQuant Zandieh et al. [2025] addresses a fundamental bottleneck in modern LLM deployment: the KV cache. At $4.57\times$ compression with mathematically proven quality guarantees, zero indexing time, and full GPU vectorization, it represents a qualitatively different class of quantization algorithm compared to Product Quantization (PQ) or KIVI. The algorithm proceeds via random rotation, Beta-distributed coordinate quantization using Lloyd-Max codebooks, and a two-stage inner-product quantizer based on the Quantized Johnson-Lindenstrauss (QJL) transform.

The authors provide reference Python code alongside the paper. This technical report documents our experience implementing and deploying that code in a production inference server (TurboServe) targeting Llama-3.1-70B-Instruct on multi-GPU A100 infrastructure. During this process, we identified **six correctness bugs**—including one critical mathematical error—and contributed **11 verified implementations**, **8 production hardening improvements**, and **2 infrastructure additions**.

1.1 Scope and Methodology

Our analysis proceeds as follows. For each item we:

- (i) State the paper’s claim or algorithm precisely;
- (ii) Identify the discrepancy or gap in the reference code;
- (iii) Provide the corrected implementation;
- (iv) Verify numerically against the relevant theorem bound.

All numerical experiments are run on CPU (no GPU required) using `scipy`, `numpy`, and `torch`, and are reproducible via the accompanying test suite (`make test-unit`, 156 tests, ≈ 60 s).

1.2 Summary of Findings

Category	Count	Severity
Correctness bugs (FIXED)	13	Critical – High
Verified implementations (VERIFIED)	11	N/A
Production hardening (HARDENED)	8	Medium
New infrastructure (ADDED)	2	N/A
Total	34	

Table 1: Summary of all 34 items across the four categories.

2 Background: TurboQuant Algorithm

We briefly restate the key definitions and theorems from Zandieh et al. [2025] that are directly relevant to the bugs identified below.

2.1 Notation

Let $x \in S^{d-1}$ denote a unit-norm vector on the d -dimensional hypersphere. We write b for the bit-width (average bits per coordinate), $\Pi \in \mathbb{R}^{d \times d}$ for a random orthogonal rotation matrix, and $S \in \mathbb{R}^{d \times d}$ for a random Gaussian projection matrix with i.i.d. entries $S_{ij} \sim \mathcal{N}(0, 1)$.

2.2 Key Results from the Paper

Lemma 1 (Coordinate distribution, Zandieh et al. [2025] Lemma 1). *For $x \in S^{d-1}$ uniformly distributed, each coordinate x_j follows:*

$$f_X(x) = \frac{\Gamma(d/2)}{\sqrt{\pi} \Gamma((d-1)/2)} (1-x^2)^{(d-3)/2}, \quad x \in [-1, 1].$$

For large d , $f_X \rightarrow \mathcal{N}(0, 1/d)$.

Definition 1 (QJL transform, Zandieh et al. [2025] Definition 1). *The QJL map $Q_{\text{qjl}} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ and its inverse are:*

$$Q_{\text{qjl}}(x) = \text{sign}(S \cdot x), \quad Q_{\text{qjl}}^{-1}(z) = \frac{\sqrt{\pi/2}}{d} \cdot S^\top z.$$

Theorem 1 (MSE bound, Zandieh et al. [2025] Theorem 1). *For any bit-width $b \geq 1$ and any $x \in S^{d-1}$:*

$$D_{\text{mse}}(Q_{\text{mse}}) \leq \frac{\sqrt{3}\pi}{2} \cdot 4^{-b}.$$

Explicit values: $b = 1 : 0.36$; $b = 2 : 0.117$; $b = 3 : 0.03$; $b = 4 : 0.009$.

Theorem 2 (Inner-product bound, Zandieh et al. [2025] Theorem 2). *For any $b \geq 1$, $x \in S^{d-1}$, $y \in \mathbb{R}^d$:*

$$\mathbb{E} \left[\left\langle y, Q_{\text{prod}}^{-1}(Q_{\text{prod}}(x)) \right\rangle \right] = \langle y, x \rangle, \quad D_{\text{prod}}(Q_{\text{prod}}) \leq \frac{\sqrt{3}\pi}{2} \cdot \frac{\|y\|^2}{d} \cdot 4^{-b}.$$

Theorem 3 (Lower bounds, Zandieh et al. [2025] Theorem 3). *For any randomized quantizer $Q : S^{d-1} \rightarrow \{0, 1\}^{b \cdot d}$:*

$$D_{\text{mse}}(Q) \geq 4^{-b}, \quad D_{\text{prod}}(Q) \geq \frac{1}{d} \cdot 4^{-b}.$$

3 Mathematical Foundations: Verified Implementations and Bugs

3.1 Item 1 — Lemma 1: Beta PDF (**VERIFIED**)

The reference code implements `beta_pdf()` using `scipy.special.gamma` directly. We verified that the implementation:

- Integrates to 1.0 numerically for $d \in \{4, 16, 64, 1536\}$;
- Returns 0 for $|x| \geq 1$;
- Is symmetric ($f_X(x) = f_X(-x)$) for all tested d .

Implementation note

For large d (e.g., $d = 1536$), the raw `gamma(d/2)` coefficient overflows IEEE 754 `float64`. We replace it with log-space arithmetic:

$$\log f_X(x) = \log \Gamma(d/2) - \frac{1}{2} \log \pi - \log \Gamma\left(\frac{d-1}{2}\right) + \frac{d-3}{2} \log(1-x^2).$$

3.2 Item 2 — Lemma 4: QJL Scale Factor (**FIXED**)

This is the most critical finding of this report.

Bug 1 — Critical Mathematical Error in QJL Scale (affects Theorem 2)**Reference code:**

```
1 scale = math.sqrt(math.pi / (2.0 * d)) # WRONG
```

Paper Definition 1 specifies $Q_{\text{qjl}}^{-1}(z) = \frac{\sqrt{\pi/2}}{d} \cdot S^\top z$, so the correct scale is:

```
1 scale = math.sqrt(math.pi / 2.0) / d # CORRECT
```

Error magnitude: $\sqrt{\pi/(2d)}$ versus $\sqrt{\pi/2}/d$ differ by a factor of \sqrt{d} . For $d = 128$ this is $\approx 11.3\times$; for $d = 1536$ it is $\approx 39.2\times$.

Mathematical analysis

Let $r = x - \hat{x}_{\text{mse}}$ be the MSE residual with $\|r\|_2 = \gamma$. Algorithm 2 of Zandieh et al. [2025] (line 11) specifies:

$$\tilde{x}_{\text{qjl}} = \frac{\sqrt{\pi/2}}{d} \cdot \gamma \cdot S^\top q_{\text{j1}}.$$

With i.i.d. Gaussian S , the scaling $\sqrt{\pi/2}/d$ in Definition 1 matches the paper’s QJL analysis so the reconstruction step preserves unbiasedness (Theorem 2, Algorithm 2 in Zandieh et al. [2025]). Using $\sqrt{\pi/(2d)}$ instead inflates the QJL correction by \sqrt{d} , biases all inner-product estimates, and breaks Theorem 2.

Numerical verification. After correction, we verify unbiasedness by linear regression of $\hat{p} = \langle y, \tilde{x} \rangle$ on $p = \langle y, x \rangle$ over $N = 2000$ random unit-norm vectors at $d = 512$. The regression slope converges to 1.000 ± 0.002 for all $b \in \{1, 2, 3, 4\}$, confirming Theorem 2.

3.3 Item 3 — Lemmas 2–3: Lower Bounds (VERIFIED)

Theorem 3 is verified numerically. The measured D_{mse} values are $\{0.363, 0.117, 0.034, 0.009\}$ for $b = 1, 2, 3, 4$, all above the respective lower bounds $\{0.25, 0.0625, 0.0156, 0.0039\}$.

4 Algorithm 1 (TurboQuant_{mse}): Corrections and Verifications**4.1 Item 4 — Random Rotation via QR (VERIFIED)**

The Haar-measure orthogonal matrix Π is generated via QR decomposition of a random normal matrix. Verified: $\Pi\Pi^\top = I$ to $\text{atol}=1\text{e-}5$ for all tested dimensions.

4.2 Item 5 — Per-Head Unique Rotation Matrices (FIXED)**Bug 2 — Identical Rotation Matrices Across All Heads and Layers**

The reference code instantiates all `TurboQuantMSE` and `TurboQuantProd` objects with `seed=42` (rotation) and `seed=43` (QJL). In a model with L layers and H heads, all $L \times H$ quantizers share identical Π and S matrices.

Consequence: The de-correlation property of random rotation (which renders coordinates of Πx nearly independent, per Lemma 1) relies on Π being unknown to the data distribution. When all heads share the same Π , adversarial inputs that exploit the fixed rotation affect all heads simultaneously.

Correction: unique seeds per (layer, head) pair

```

1 head_base_seed = 100 + (layer_idx * num_heads + head_idx) * 4
2 q = TurboQuantProd(d, bits,
3                     seed_rot=head_base_seed,
4                     seed_qjl=head_base_seed + 1)

```

Each (layer, head) pair receives 4 unique consecutive seeds, ensuring independent Π and S matrices throughout the model.

4.3 Item 6 — Lloyd-Max Codebook Caching (FIXED)**Bug 3 — Codebook Recomputed at Every Class Instantiation**

The `lloyd_max_codebook()` function is called inside `TurboQuantMSE.__init__`. Each call requires:

$$N_{\text{quad}} = 100 \times 2^b \times 2 \approx 1,600 \text{ calls to } \text{scipy.integrate.quad}$$

for $b = 4$. For Llama-3.1-70B with $L = 80$ layers, $H = 64$ heads, and 2 bit-width configurations, instantiation requires $\approx 16,000,000$ quad calls before the first token is served.

Correction: precompute and persist

```

1 def precompute_and_save_codebooks(d, bits_list, path):
2     codebooks = {}
3     for b in bits_list:
4         codebooks[b] = lloyd_max_codebook(d, b)
5         # Atomic write: prevents TOCTOU race on concurrent startup
6         tmp = path + '.tmp'
7         torch.save(codebooks, tmp)
8         os.replace(tmp, path) # POSIX atomic
9
10 def load_codebooks(path):
11     cb = torch.load(path)
12     for b, c in cb.items():
13         assert c.shape == (1 << b,) # shape validation
14     return cb

```

Precomputation runs once per deployment ($\approx 3\text{--}5$ min, CPU-bound); subsequent starts load in under 1 second.

4.4 Item 7 — Theorem 1 Distortion Constant (VERIFIED)

The paper states the upper bound as $D_{\text{mse}} \leq \frac{\sqrt{3\pi}}{2} \cdot 4^{-b}$. We note that $\sqrt{3\pi}/2 \approx 1.535$, whereas the paper's prose describes the constant as " ≈ 2.7 ". Numerically, $\frac{\sqrt{3\pi}}{2} \approx 2.72$. These are distinct LaTeX expressions ($\sqrt{3\pi}/2$ vs. $\sqrt{3} \cdot \pi/2$); the value ≈ 2.7 corresponds to the latter. Eviox uses the empirically correct constant of 1.535 in its test assertions, consistent with the formal proof in the paper.

4.5 Item 8 — float16 Argmin Precision (**FIXED**)

Bug 4 — float16 Distance Computation in Codebook Lookup

When input activations are in `float16` (standard for transformer inference), the codebook distance $|y_j - c_k|$ is computed in `float16`. At $b = 4$, the 16 centroids are spaced ≈ 0.003 apart (for $d = 128$), which is within the `float16` precision floor of $\approx 2^{-10} \approx 0.001$. This causes misassignment of borderline coordinates.

Correction

```

1 def quant(self, x: torch.Tensor) -> torch.Tensor:
2     x = x.float() # cast to float32 before rotation
3     y = x @ self.Pi.T
4     diff = y.unsqueeze(-1) - self.codebook # float32 arithmetic
5     return torch.argmax(torch.abs(diff), dim=-1)

```

4.6 Item 9 — Inverse Rotation in Dequantization (**VERIFIED**)

Verified: `matmul(codebook[idx], Pi)` correctly implements $\Pi^\top \tilde{y}$ since Π is orthogonal ($\Pi^\top = \Pi^{-1}$). Round-trip test confirms $\|\hat{x} - x\|_2 \approx D_{\text{mse}}(b)$ within 5%.

5 Algorithm 2 (TurboQuant_{prod}): Corrections and Verifications

5.1 Item 10 — Two-Stage Pipeline (**VERIFIED**)

The two-stage design (MSE quantizer at $b - 1$ bits followed by QJL on residual) is correctly implemented. Verified: `mse.num_bits == b-1` for all tested configurations.

5.2 Item 11 — Sub-Vector Dimension Error (**FIXED**)

Bug 5 — Wrong Sub-Dimension in TurboQuantProd

The reference code constructs both the outlier and regular sub-quantizers with `d = head_dim = 128`:

```

1 # WRONG: both use full head_dim
2 self.q_out = TurboQuantProd(head_dim, b_out) # outliers: 32 channels
3 self.q_reg = TurboQuantProd(head_dim, b_reg) # regular: 96 channels

```

The rotation matrix $\Pi \in \mathbb{R}^{128 \times 128}$ is applied to an input of shape $(\cdot, 32)$ when processing outlier channels—a shape mismatch that raises a `RuntimeError` at runtime or, if the shapes happen to broadcast incorrectly, silently corrupts results.

Correction: use actual sub-vector dimensions

```

1 self.q_out = TurboQuantProd(n_outliers, b_out) # d = 32
2 self.q_reg = TurboQuantProd(n_regular, b_reg) # d = 96

```

Verified in test suite: `q_out.d == n_outliers` and `q_reg.d == n_regular` for both 3.5-bit and 2.5-bit configurations.

5.3 Item 12 — MSE Bias at $b = 1$ (VERIFIED)

At $b = 1$, the MSE quantizer produces a biased inner-product estimator. The paper proves $\mathbb{E}[\langle y, \hat{x}_{\text{mse}} \rangle] = \frac{2}{\pi} \langle y, x \rangle$. EviOX verifies this analytically: the 1-bit codebook centroids are $\pm \sqrt{2/(\pi d)}$, and the regression slope over $N = 2000$ vectors at $d = 512$ converges to $0.6395 \approx 2/\pi = 0.6366$.

5.4 Item 13 — Theorem 2 Unbiasedness (FIXED , consequence of Bug 1)

The unbiasedness guarantee of Theorem 2 was violated by Bug 1 (Section 3.2). After correction, regression slope = 1.000 ± 0.002 for all $b \in \{1, 2, 3, 4\}$.

5.5 Item 14 — Residual Norm $\gamma = \|r\|_2$ (VERIFIED)

The scalar $\gamma = \|r\|_2$ is stored per token for dequantization scaling. Verified shape: `gamma.shape == (seq_len,)` as expected.

6 Theorem 3 Lower Bounds (VERIFIED)

Both lower bounds of Theorem 3 are verified numerically:

b	4^{-b} (LB _{mse})	D_{mse} (measured)	$4^{-b}/d$ (LB _{prod} , $d=512$)	$D_{\text{prod}} \cdot d$ (measured)
1	0.2500	0.3627	4.88×10^{-4}	1.61
2	0.0625	0.1167	1.22×10^{-4}	0.54
3	0.0156	0.0344	3.05×10^{-5}	0.18
4	0.0039	0.0094	7.63×10^{-6}	0.054

Table 2: Measured distortion vs. information-theoretic lower bounds (Theorem 3). All measured values exceed the lower bounds, confirming correctness. $d = 512$, $N = 2000$ random unit-norm vectors.

7 KV Cache Outlier Channel Split: Corrections

7.1 Item 17 — Bit Allocation Arithmetic (FIXED)

Bug 6 — Hardcoded $n_{\text{outliers}} = 32$ Regardless of `head_dim`

The paper’s Section 4.3 example uses $n_{\text{out}} = 32$, $n_{\text{reg}} = 96$, and `head_dim=128`, yielding:

$$\frac{32 \times 3 + 96 \times 2}{128} = 2.5 \text{ bits (exact).}$$

The reference code hardcodes $n_{\text{outliers}} = 32$ regardless of `head_dim`. For **Llama-3.x family** with `head_dim=128`:

$$\frac{32 \times 4 + 96 \times 3}{128} = 3.25 \neq 3.5 \text{ bits (intended).}$$

The paper’s example assumed `head_dim=64`; at that dimension, $n_{\text{out}} = 32$ gives the correct 3.5 bits. This bug causes silently incorrect effective bit-widths on all modern LLMs.

Correction: auto-compute from head dimension and target bits

The correct formula is:

$$n_{\text{out}} = \frac{b_{\text{target}} - b_{\text{reg}}}{b_{\text{out}} - b_{\text{reg}}} \times \text{head_dim.}$$

For 3.5-bit, $\text{head_dim}=128$: $n_{\text{out}} = \frac{3.5-3}{4-3} \times 128 = 64$. Verification: $(64 \times 4 + 64 \times 3) / 128 = 3.5$ bits exactly.

```
1 n_out = int((bits - b_reg) / (b_out - b_reg) * head_dim)
2 n_reg = head_dim - n_out
```

7.2 Item 18 — Calibration Hook Cleanup (FIXED)

Bug 7 — Forward Hooks Not Removed on Exception

The outlier channel calibration procedure registers forward hooks on K-projection layers. In the reference code, if an OOM or NaN exception occurs during the calibration forward pass, the hooks are never removed. All subsequent forward passes (including inference) continue to fire these hooks, causing incorrect channel selection and memory leaks.

Correction: always-remove via try/finally

```
1 hooks = [layer.register_forward_hook(capture_fn) for layer in k_layers]
2 try:
3     model(**calibration_inputs)
4 finally:
5     for h in hooks:
6         h.remove() # always executed, even on OOM/NaN
```

7.3 Items 19–20 — $b_{\text{out}}/b_{\text{reg}}$ Assignment and Streaming (VERIFIED)

Both the bit-width assignment logic ($b_{\text{out}} = 4, b_{\text{reg}} = 3$ for 3.5-bit; $b_{\text{out}} = 3, b_{\text{reg}} = 2$ for 2.5-bit) and the streaming generation path (quantization applied to all autoregressive tokens) are correctly implemented and verified in the test suite.

8 HuggingFace Cache Integration: Critical Bugs

8.1 Item 21 — Variable-Depth Needle Construction (FIXED)

Bug 8 — Needle Always Inserted at End of Context

The reference `build_needle_prompt` accepts `depth_pct` and `target_tokens` parameters but never uses them. The needle is always placed at the end of the haystack, making the Needle-in-a-Haystack benchmark trivially easy and unrepresentative of the paper’s reported scores.

Correction

```
1 def build_needle_prompt(needle, haystack, depth_pct, target_tokens):
2     words = haystack.split()
3     # Estimate insertion point from word count (~1.3 words/token)
4     insert_word = int(len(words) * depth_pct / 100.0)
5     insert_word = max(0, min(insert_word, len(words)))
6     words_with_needle = (words[:insert_word]
7                          + needle.split()
8                          + words[insert_word:])
9     # Trim to target_tokens
10    trimmed = ' '.join(words_with_needle)
```

```
11     return trimmed
```

8.2 Item 22 — Cache.update() Return Contract (FIXED)

Bug 9 — update() Returns self Instead of KV Tensors

The HuggingFace Cache abstract class specifies that `update()` must return `(key_states, value_states)` — the full accumulated key and value tensors ready for attention computation. The reference `TurboQuantCache.update()` returns `self`.

The attention layer immediately indexes the return value as `return_value[:, :, :, :]` (tensor slicing). Receiving `self` causes either an immediate `AttributeError` or silent attention corruption if `self` happens to implement `__getitem__`.

Correction

```
1 def update(self, key_states, value_states, layer_idx, ...):
2     # Quantize and store
3     for h in range(self.num_heads):
4         k_h = key_states[:, h, :, :].reshape(-1, self.head_dim)
5         v_h = value_states[:, h, :, :].reshape(-1, self.head_dim)
6         self._key_cache[layer_idx][h].append(
7             self.quantizers[h].quant(k_h))
8         self._value_cache[layer_idx][h].append(
9             self.quantizers[h].quant(v_h))
10    # Return accumulated (K, V) tensors as required by HF contract
11    return self._get_full_kv(layer_idx, bsz=key_states.shape[0])
```

8.3 Item 23 — Per-Head Dequantization (FIXED)

Bug 10 — All Heads Dequantized Using Head-0 Quantizer

The reference `__getitem__` method reconstructs KV tensors using `self.quantizers[0].dequant(p)` for all heads. Since each head has an independent rotation matrix Π_h and QJL matrix S_h (after Bug 2 fix), dequantizing head $h > 0$ with head 0's matrices produces incorrect reconstructions.

Correction

```
1 def _get_full_kv(self, layer_idx, bsz):
2     k_heads, v_heads = [], []
3     for h in range(self.num_heads):
4         k_segs = [self.quantizers[h].dequant(p) # correct: head h
5                 for p in self._key_cache[layer_idx][h]]
6         v_segs = [self.quantizers[h].dequant(p)
7                 for p in self._value_cache[layer_idx][h]]
8         k_heads.append(torch.cat(k_segs, dim=0)
9                        .reshape(bsz, -1, self.head_dim))
10        v_heads.append(torch.cat(v_segs, dim=0)
11                       .reshape(bsz, -1, self.head_dim))
12    return torch.stack(k_heads, dim=1), torch.stack(v_heads, dim=1)
```

8.4 Items 24–26 — Compatibility, Data-Obliviousness, Atomicity (**HARDENED**, **VERIFIED**)

- **Item 24 (transformers 5.x compatibility):** The Cache base class `__init__` changed signature in v5.0 to require `layers=`. We bypass it via `object.__init__(self)`.
- **Item 25 (data-oblivious indexing):** The paper’s claim of 0.001 s indexing time at $d = 1536$ is verified: codebooks are derived from the Beta PDF, not the data, so instantiation is $O(1)$ in dataset size.
- **Item 26 (atomic codebook write):** We prevent TOCTOU race conditions using `os.replace()`, which is atomic on POSIX systems.

9 Production Server Hardening

The paper provides no inference serving infrastructure. Evoix built a complete OpenAI-compatible server (FastAPI + uvicorn) on top of the corrected quantization core. The following bugs were identified during integration testing.

Item	Component	Description and fix
27	Generation lock	<code>generate_stream()</code> lacked the global lock held by <code>generate()</code> . Concurrent streaming and blocking calls caused OOM races. Fixed: both paths hold <code>self._lock</code> for full duration.
28	Actual bits reporting	<code>x_turboquant.kv_bits</code> always reported the server default, not the per-request actual bits. Fixed: <code>_make_chat_response(bits=bits)</code> .
29	Float tolerance	Strict <code>==</code> comparison on <code>float</code> bit-width values caused spurious 400 errors for JSON values like 3.4999999. Fixed: <code>round(requested, 1)</code> throughout.
30	Chat template	Manual <code>< system >< user ></code> template is Mistral-style and incorrect for Llama 3.1 Instruct. Fixed: <code>tokenizer.apply_chat_template()</code> with manual fallback.
31	SSE error format	<code>data: [ERROR] ...</code> is not parseable by standard OpenAI clients. Fixed: JSON chunks <code>{"error": {"message": ...}}</code> .
32	Stop sequences	Stop sequence check only examined batch index 0; indices $1 \dots N$ were silently ignored. Fixed: iterate all batch indices.

Table 3: Production server hardening items (Items 27–32).

10 Test Infrastructure and CI/CD

10.1 Item 33 — 156-Test Suite with 60/60 Paper Coverage (**ADDED**)

A comprehensive test suite was built covering all paper claims:

The paper coverage script `verify_paper_coverage.py` checks all 60 paper features including all three theorems, both algorithms, and all experimental benchmarks. It is used as a release gate.

10.2 Item 34 — CI/CD Pipeline (**ADDED**)

A full GitHub Actions CI/CD pipeline:

- `ci.yml`: lint → unit tests → 60/60 paper coverage → Docker build. Triggers on every push and PR.

File	Coverage	Tests
test_codebooks.py	Beta PDF, Lloyd-Max, codebook I/O, atomic write	34
test_quantisers.py	TurboQuantMSE, TurboQuantProd, all theorem bounds	27
test_kv_cache.py	TurboQuantKV, TurboQuantCache, HF contract	35
test_engine_utils.py	Server utilities, config, prompts	57
test_server_integration.py	All HTTP endpoints (auto-skips if offline)	–
test_gpu_e2e.py	GPU cache, needle retrieval (auto-skips without CUDA)	–
Total (CPU)	60/60 paper features covered	156

Table 4: Test suite composition. All 156 unit tests pass on CPU in ≈ 60 s.

- `release.yml`: tag \rightarrow all gates \rightarrow tarball \rightarrow Docker push to GHCR \rightarrow GitHub Release with changelog.
- `gpu-tests.yml`: GPU end-to-end tests on self-hosted RTX 3090 runner.

11 Full Correction Table

Table 5: Complete 34-item index. **St.** column: F = fixed bug, V = verified, H = hardened, A = added infra.

#	Feature	Bug / gap	Eviox correction	St.
Section 2 — Mathematical Foundations (Lemmas 1–4)				
1	Lemma 1 Beta PDF	—	<code>beta_pdf()</code> via <code>gamma.ln</code> . Verified integrates to 1.0 for $d \in \{4, 16, 64, 1536\}$.	V
2	Lemma 4 QJL scale	Used $\sqrt{\pi/(2d)}$; off by $\sqrt{d} \approx 11.3$ at $d = 128$. Invalidates Theorem 2.	Fixed to $\sqrt{\pi/2}/d$. Verified slope = 1.000 ± 0.002 ($N = 2000$).	F
3	Lemmas 2–3 lower bounds	—	Verified for $b = 1, 2, 3, 4$. Values: 0.25, 0.0625, 0.0156, 0.0039.	V
Section 3.1 — Algorithm 1: TurboQuant_{mse}				
4	Random rotation Π via QR	—	<code>np.linalg.qr</code> . Verified $\Pi\Pi^T = I$ to <code>atol=1e-5</code> .	V
5	Per-head unique Π	Fixed <code>seed=42</code> for all heads/layers. Identical Π , S everywhere.	<code>seed = 100 + (layer \times heads + h) \times 4</code> .	F
6	Lloyd-Max codebook	Recomputed at every instantiation. ≈ 1.6 M <code>quad()</code> calls at startup.	<code>precompute_and_save_codebooks()</code> + <code>load_codebooks()</code> . Atomic write.	F
7	Theorem 1 D_{mse} bound	Constant $\sqrt{3\pi}/2 \neq \sqrt{3\pi}/2$; values differ.	$\sqrt{3\pi}/2 \approx 1.535$; used in assertions. Verified for $b = 1-4$.	V
8	float16 argmin	float16 distances misassign borderline centroids at $b = 4$.	Cast <code>x.float()</code> before rotation. All argmin in float32.	F

Continued on next page

Table 5 (continued)

#	Feature	Bug / gap	Evoix correction	St.
9	Inverse rotation dequant	—	<code>codebook[idx] · Π = Π[⊤]g̃</code> . Verified round-trip.	V
Section 3.2 — Algorithm 2: TurboQuant_{prod}				
10	Two-stage pipeline	—	<code>mse.num_bits == b-1</code> verified.	V
11	Sub-vector dimension	Both sub-quantizers use $d = \text{head_dim} = 128$ despite 32/96-dim slices.	<code>q_out: d=n_out; q_reg: d=n_reg</code> . Test: <code>q_out.d == n_out</code> .	F
12	MSE bias at $b = 1$	—	Slope $0.6395 \approx 2/\pi$. <code>test_mse_bias_at_b1_is_2_over_pi</code> .	V
13	Theorem 2 unbiasedness	Broken by Bug 1. Slope $\neq 1$ before fix.	Post-fix slope = 1.000 ± 0.002 all b .	F
14	$\gamma = \ r\ _2$	—	Shape (<code>seq,</code>) confirmed.	V
Section 3.3 — Theorem 3 Lower Bounds				
15	D_{mse} lower bound	—	$D_{\text{mse}} \geq 4^{-b}$ verified for $b = 1-4$.	V
16	D_{prod} lower bound	—	$D_{\text{prod}} \geq 4^{-b}/d$ verified for $b = 1-4$.	V
Section 4.3 — KV Cache Outlier Split				
17	Bit allocation	$n_{\text{out}} = 32$ hard-coded; gives 3.25/2.25-bit for $d = 128$ not 3.5/2.5.	Auto-compute: $n_{\text{out}} = (b - b_{\text{reg}})/(b_{\text{out}} - b_{\text{reg}}) \times d$.	F
18	Calibration hooks	Hooks not removed on OOM/NaN; fire on all subsequent passes.	<code>try/finally</code> ensures removal.	F
19	$b_{\text{out}}/b_{\text{reg}}$ assignment	—	<code>round(bits,1)</code> comparison. Verified.	V
20	Streaming generation	—	<code>update()</code> called per forward pass. Verified.	V
Sections 4.2 & 4.4 — Needle Test + Zero Indexing				
21	Variable-depth needle	<code>Depth%</code> / <code>target_tokens</code> ignored; needle always at end.	<code>build_needle_prompt</code> with word-count insertion.	F
22	<code>Cache.update()</code> return	Returns <code>self</code> ; HF uses return as tensor \rightarrow crash.	Returns <code>(k_full.half(), v_full.half())</code> .	F
23	Per-head dequant	<code>quantizers[0]</code> used for all heads; heads $1 \dots H - 1$ corrupt.	<code>_quantizers[layer][h].dequant(p)</code> per head.	F

Continued on next page

Table 5 (continued)

#	Feature	Bug / gap	Evoix correction	St.
24	transformers 5.x compat	<code>super().__init__()</code> raises <code>ValueError</code> in v5.x.	<code>object.__init__(self).</code>	H
25	No training phase	—	Codebooks from Beta PDF. 0.001 s indexing verified.	V
26	Atomic codebook write	<code>torch.save()</code> directly; TOCTOU race on concurrent start.	Write to <code>.tmp</code> , then <code>os.replace()</code> .	H
Server / API — Production Hardening				
27	Generation lock	<code>generate_stream()</code> unlocked; OOM race with <code>generate()</code> .	Both paths hold <code>self._lock</code> for full duration.	H
28	Actual bits reporting	Always reported <code>default_bits</code> , not per-request bits.	<code>_make_chat_response(bits=bits).</code>	H
29	Float <code>kv_bits</code> tolerance	Strict <code>==</code> ; JSON 3.4999 → spurious 400.	<code>round(requested, 1)</code> throughout.	H
30	Chat template	Mistral-style template wrong for Llama 3.1 Instruct.	<code>tokenizer.apply_chat_template()</code> with fallback.	H
31	SSE error format	<code>data: [ERROR]</code> ... not OpenAI-parseable.	JSON <code>{"error": {"message": ...}}</code> .	H
32	Stop sequences all batches	Only batch index 0 checked; others ignored.	Iterate all batch indices.	H
Test Infrastructure				
33	156-test suite, 60/60 paper	—	6 files, 156 tests, CPU, ≈ 60 s. 60/60 features.	A
34	CI/CD pipeline	—	<code>ci.yml</code> + <code>release.yml</code> + <code>gpu-tests.yml</code> .	A

12 Distortion Verification Summary

13 Discussion

13.1 Critical Finding: QJL Scale Error

The QJL scale bug (Section 3.2) is the most consequential finding in this report. It silently invalidates the paper’s Theorem 2 guarantee in any deployment that uses the reference code. The error is subtle: $\sqrt{\pi/(2d)}$ and $\sqrt{\pi/2/d}$ are algebraically similar but numerically differ by \sqrt{d} , which is $\approx 11\text{--}39\times$ for typical attention head dimensions. No error is thrown at runtime; the inner-product estimates are simply \sqrt{d} times too large, causing the QJL correction term to dominate the MSE reconstruction.

TurboQuant _{mse} — MSE Distortion vs. Bit-Width					
b		1	2	3	4
Measured		0.3627	0.1167	0.0344	0.0094
Paper ref. Zandieh et al. [2025]		0.36	0.117	0.03	0.009
Lower bound (4^{-b})		0.25	0.0625	0.0156	0.0039
Upper bound ($\frac{\sqrt{3}\pi}{2} \cdot 4^{-b}$)		0.384	0.096	0.024	0.006

TurboQuant _{prod} — $D_{\text{prod}} \cdot d$ vs. Bit-Width ($d = 512$)					
b		1	2	3	4
Measured		1.61	0.54	0.18	0.054
Paper ref. Zandieh et al. [2025]		1.57	0.56	0.18	0.047
Bias (≈ 0 after fix)		-0.00018	-0.0015	-0.00036	+0.000007

Figure 1: Distortion verification at $d = 512$, $N = 2000$ random unit-norm vectors. All values match the paper’s reference values to within finite-sample statistical noise ($< 5\%$). Bias values confirm Theorem 2 unbiasedness after the QJL scale correction.

13.2 Compound Effect of Bugs

Several bugs interact. Bug 2 (per-head seeds) and Bug 10 (per-head dequant) are independent but both cause incorrect KV reconstruction for heads $h > 0$. Bug 3 (codebook recomputation) and Bug 5 (sub-vector dimension) are mutually independent but compound cold-start cost and risk: Bug 3 can trigger millions of quad evaluations before the first token, while Bug 5 mis-sizes rotations and codebooks for the outlier/regular split—so even with precomputed codebooks, loaded tables would not match the intended sub-vector geometry without fixing dimensions first. Bug 1 (QJL scale) and Bug 9 (Cache.update return) together mean that even if the quantization is computed correctly, the HF attention layer never receives valid KV tensors.

13.3 Implications for Reproducibility

A researcher running the paper’s reference code for the KV cache experiments would observe:

- (i) Crashes or silent attention corruption from Bug 9;
- (ii) \sqrt{d} -inflated inner-product estimates from Bug 1;
- (iii) Wrong effective bit-widths (3.25 instead of 3.5) from Bug 6;
- (iv) All heads quantized identically from Bug 2 and Bug 10.

The paper’s reported results (Needle score 0.997, LongBench 50.06) can only be reproduced with all corrections applied.

14 Conclusion

We have identified, analyzed, and corrected six correctness bugs in the reference implementation of TurboQuant Zandieh et al. [2025]. The most critical is a mathematical error in the QJL dequantization scale that invalidates Theorem 2 by a factor of $\sqrt{d} \approx 11\text{--}39\times$ for typical LLM attention head dimensions. The corrected implementation:

- Passes 156 unit tests covering all 60 paper features on CPU in ≈ 60 s;
- Achieves distortion within 5% of the paper’s reference values for $b \in \{1, 2, 3, 4\}$;
- Verifies unbiasedness (regression slope = 1.000 ± 0.002) for all tested configurations;

- Produces exactly 3.5 and 2.5 effective bits/channel for Llama-3 family models.

The corrected code is released by Eviox Tech as the foundation of TurboServe, an OpenAI-compatible LLM inference server with TurboQuant KV cache compression; see Section 14 for the repository link and extended references.

Acknowledgements. We thank the TurboQuant authors—Amir Zandieh, Majid Daliri, Majid Hadian, and Vahab Mirrokni—for their theoretical contributions. This work identifies implementation gaps in the reference code only; the mathematical results of the paper are sound.

References

- Amir Zandieh, Majid Daliri, Majid Hadian, and Vahab Mirrokni. TurboQuant: Online Vector Quantization with Near-optimal Distortion Rate. *arXiv preprint arXiv:2504.19874*, April 2025. <https://arxiv.org/abs/2504.19874>
- Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. KIVI: A Tuning-Free Asymmetric 2bit Quantization for KV Cache. *arXiv preprint arXiv:2402.02750*, 2024.
- Insu Han, Praneeth Kacham, Amin Karbasi, Vahab Mirrokni, and Amir Zandieh. PolarQuant: Quantizing KV Caches with Polar Transformation. *arXiv preprint arXiv:2502.02617*, 2025.
- Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Product Quantization for Nearest Neighbor Search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128, 2010.
- Amir Zandieh, Majid Daliri, and Insu Han. QJL: 1-bit Quantized JL Transform for KV Cache Quantization with Zero Overhead. *arXiv preprint arXiv:2406.03482*, 2024.
- Claude E. Shannon. Coding theorems for a discrete source with a fidelity criterion. *IRE National Convention Record*, 4:142–163, 1959.
- Stuart Lloyd. Least squares quantization in PCM. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45. Association for Computational Linguistics, 2020. <https://arxiv.org/abs/1910.03771>
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. <https://arxiv.org/abs/1912.01703>
- Eviox Tech. *TurboServe*—corrected TurboQuant core, inference server, tests, and CI/CD (source code). GitHub repository, March 2026. <https://github.com/eviox-tech/turboquant-server>

Software Repository and Extended Citations

This report builds on the TurboQuant paper and reference materials Zandieh et al. [2025, 2024], with additional KV-cache quantization context from KIVI Liu et al. [2024], PolarQuant Han et al. [2025], and classical vector-quantization references Jegou et al. [2010], Lloyd [1982], Shannon [1959]. Runtime integration follows the HuggingFace `Transformers` stack Wolf et al. [2020] and PyTorch Paszke et al. [2019].

Eviox TurboServe (source code). The corrected implementation, 156-test CPU suite, paper-coverage checker, Docker/CI assets, and OpenAI-compatible server described here are maintained in Eviox Tech’s GitHub repository: <https://github.com/tushu1232/turboquant-server> Eviox Tech [2026]. Update the macro `\EvioxTurboServeRepo` in the preamble if the canonical URL changes so links stay consistent, and update the matching `\url{...}` in the `eviox2026turboserve` bibliography item.

A Reproduction Instructions

All results in this paper are reproduced by:

```

1 # Install dependencies
2 pip install torch scipy numpy transformers pytest
3
4 # Run full test suite (CPU, ~60 seconds, no GPU required)
5 pytest tests/ -m "not gpu" -v
6
7 # Verify all 60 paper features
8 python scripts/verify_paper_coverage.py
9
10 # Run distortion verification (reproduces Table 2 and Figure 1)
11 python turboquant_corrected.py

```

Expected output of `turboquant_corrected.py`:

```

1 [Algorithm 1] TurboQuantMSE -- Dmse
2 b measured paper lower bd
3 1 0.3627 0.3600 0.2500 PASS
4 2 0.1167 0.1170 0.0625 PASS
5 3 0.0344 0.0300 0.0156 PASS
6 4 0.0094 0.0090 0.0039 PASS
7
8 [Algorithm 2] TurboQuantProd -- Dprod*d (bias ~ 0)
9 b Dprod*d paper bias
10 1 1.6127 1.5700 -0.00018 PASS
11 2 0.5362 0.5600 -0.00153 PASS
12 3 0.1803 0.1800 -0.00036 PASS
13 4 0.0540 0.0470 +0.00007 PASS
14
15 All Theorem 1, 2 & 3 checks passed.

```