

**title:** Cubie — Unified Analysis

**subtitle:** intc-v1.0 evidence pack × wwt-tech reference set × cubie-tep-v1.2.0 cross-referenced

**date:** 2026-05-30

**author:** Centillion AI · TrustFortress.ai

# Cubie — Unified Analysis

**Scope of this document.** Mirror the analytical pattern used to produce the `intc-v1.0/` GPU-fleet evidence pack and the eight `wwt-tech/` reference documents, applied to the cubie TEP whitepaper work. Produce one unified narrative that cross-references every canonical number across all three sets.

**Three artifact groups analyzed:**

Group	Files	What it grounds
<b>intc-v1.0</b> (existing)	1 MANIFEST + 10 JSONs (TP4 NVLink, DCGM hunter, Alibaba squatter, Azure 2023/2024, LMM 2025/2026, MaverIQ, Prometheus, micro-stall)	The GPU-fleet pain claim: 66.85% blast radius, 87.53% monitoring blind spot, 51.76% VRAM lock — the "why does Cubie need to exist?" evidence
<b>wwt-tech</b> (existing)	1 README + 8 technical reference MDs (gpu-compat-matrix, vendor-adapter-matrix, denial-taxonomy, deployment-path, two-lane-architecture, not-ai-on-ai-positioning, shadow-mode-protocol, risk-register)	The operational deployment spec — how Cubie actually deploys, what it denies, how it integrates
<b>cubie-tep-v1.2.0</b> (new this analysis)	1 MANIFEST + 13 JSONs (Cubie sweep × algorithm, Cubie-vs-Opus-4.8 head-to-head, Opus 4.7 vs 4.8 same-day delta, day-2 reproducibility)	The math-proof credibility: Cubie's 100/100/100 + 0% FAR vs every benchmarked baseline

**Each pack stands alone. Together they form the spectrum: pain → solution → deployment.**

## Part 1 — intc-v1.0 evidence pack (re-analyzed)

**Anchor file:** `demo/evidence/intc-v1.0/intc-v1.0-MANIFEST.md` in the cubie-tf repo.

### What the intc-v1.0 pack proves

The dashboard at `demo/dashboard.html` cites four headline numbers that motivate the whole product story. Each number has hashed JSON ground truth:

Claim	Number	Ground truth
TP4 NVLink cascade destroys cluster capacity at scale	<b>66.8479%</b> blast radius across 157,411 GPU intervals	<code>intc-v1.0-tp4_nvlink_results.json</code> → <code>full_tp4_blast_radius.capacity_pct</code>
Existing DCGM monitoring misses this entirely	<b>87.53%</b> blind spot (only catches 8.34% of what Cubie's exact-zero-SM detector catches)	derived from <code>intc-v1.0-hunter_results.json</code> → <code>alibaba_signals.dcgm_thrash_proxy.*</code>
The cluster spends most of its time NOT producing work	<b>33.15%</b> productive / <b>83.29%</b> parked-or-failed split	<code>intc-v1.0-alibaba_micro_stall_metrics.json</code>
Capacity multiplier for a customer that fixes this	<b>2.51×</b>	derived from <code>intc-v1.0-alibaba_v2026_squatter_metrics.json</code>

## Why the evidence pack matters operationally

Two enterprise risks the pack addresses:

- 1. Auditability of the pitch.** When a customer asks "where does 66.85% come from?" the answer is a hashed file in the repo, not a slide. The PowerShell verification command in the MANIFEST runs in 4 seconds and proves the cited numbers match the JSON.
- 2. APPEND-ONLY discipline.** The MANIFEST explicitly forbids editing in place. Any analysis re-run produces `intc-v1.1/` alongside; the dashboard's footer carries the active version pointer. Bumping versions is a single constant + a new directory. This survives the kind of "wait, where did that number come from in slide 7?" question that kills enterprise sales cycles.

## How it integrates with the rest of the Cubie story

`intc-v1.0` supplies the **pain side** of the elevator pitch. The `cubie-tep-v1.2.0` pack (Part 3 below) supplies the **solution side** — formally-verified Cubie hits 100/100/100 + 0% FAR where statistical methods and frontier LLMs miss. The `wwt-tech` reference set (Part 2) supplies the **deployment side**.

## Part 2 — wwt-tech reference set (re-analyzed)

**Anchor:** `demo/wwt-tech/README.md` lists eight technical artifacts. Each is partner-safe (no internal pricing, no attributed customer language, no competitive battlecards).

## Doc-by-doc analytical pass

---

### `gpu-compat-matrix.md` (42 lines)

**Purpose:** Lock down which GPU architectures Cubie supports, before any conversation devolves into "does it work on H100?"

**Cubie position: Hardware-agnostic above the driver.** The matrix explicitly enumerates NVIDIA H100/H200/B100/B200/GB200, AMD MI300X/MI325X/MI350/MI400, Intel Gaudi 3 + Falcon Shores. The admit decision uses Rust `no_std` + generic SHA-256/HMAC + 192-byte envelope wire format. CUDA / ROCm / Level-Zero are NOT dependencies.

**What changes per vendor:** the *attestation source* (NVIDIA `nvtrust` nested-in-TDX vs AMD SEV-SNP) and the *driver runtime* the protected workload runs on. Both handled in `cubie-control`.

### `vendor-adapter-matrix.md` (97 lines)

**Purpose:** Map the runtime telemetry and recovery hooks per vendor. Tells the customer's platform team exactly which APIs Cubie talks to on each side.

**Coverage:** Telemetry — DCGM (NVIDIA), AMD-SMI (AMD), Habana SynapseAI (Intel), Prometheus (generic). Recovery — vendor-specific reset paths, NVLink topology rebuild, ECC scrub, MIG repartition.

### `denial-taxonomy.md` (109 lines)

**Purpose:** Define the five admit decision classes with full lane-vs-request decomposition and remediation codes.

**Decision tree:** **LOCK** (capability HMAC failure) / **JAM** (lane saturation) / **SHATTER** (structural integrity failure) / **DEFLECT** (route to alternate lane) / **DENY** (policy/budget). Each gets a remediation code the caller's SDK consumes deterministically.

This is the canonical operational vocabulary. Every claim in the whitepaper's §C "Operational response" maps to one of these five classes.

### `deployment-path.md` (113 lines)

**Purpose:** The 4-stage rollout sequence that customer ops teams will actually agree to.

**Stages:** 1. **SPAN/mirror** — Cubie observes the traffic, takes no action, reports what it would have done  
 2. **Shadow mode** — Cubie's admit decisions are computed but not enforced; gateway logs the deltas  
 3. **Policy simulation** — partial enforcement on a tagged-traffic slice; full observability  
 4. **Optional inline enforcement** — full admit gate in the path (with bypass and human-override)

**Critical operational property:** every stage is reversible. Customer can roll back to the previous stage in seconds without state migration.

### `two-lane-architecture.md` (83 lines)

**Purpose:** Resolve the placement ambiguity that wrecks first-customer conversations: "Where in my stack does Cubie go?"

**Answer: Both lanes simultaneously, identical primitive.** Lane 1 = pre-inference ToR-class gate (sub-microsecond, bare-metal, deterministic). Lane 2 = in-workflow MCP/agent policy gate (millisecond-class, gateway-resident, content-aware). The admit decision uses the same 192-byte envelope on both lanes. Defense in depth without doubling the configuration surface.

### `not-ai-on-ai-positioning.md` (55 lines)

**Purpose:** Pre-empt the "is this just another LLM guardrail model?" question that triggers customer skepticism.

**Position: Cubie is deterministic infrastructure, not an AI model.** No learned weights, no inference passes, no GPU consumption by Cubie itself. The only "intelligence" is formal-proof + cryptographic-capability + structural-topology logic.

This is the **single most important positioning point** when selling into a CISO + AI-platform-lead joint audience. Internal AI gateways are increasingly seen as model risk; Cubie is infrastructure that *prevents* model risk.

### `shadow-mode-protocol.md` (119 lines)

**Purpose:** The 14-day validation rubric customers will commit to when "buying" their first Cubie deployment.

**Pass/fail metrics:** False-positive rate < 0.5% on shadow-mode comparison; latency overhead < 1ms p99; alarm volume per day within  $\pm 2\times$  of pre-Cubie baseline; zero unexplained denials. Bind each metric to a specific Prometheus expression and a specific dashboard panel.

### `risk-register.md` (102 lines)

**Purpose:** Pre-emptive risk acknowledgment. Customers under regulator scrutiny will demand this before any deployment.

**Risk categories covered:** false-positive measurement methodology + remediation; bypass mode design (operator-pulled, audit-logged, time-bound); human override (Article 14 compliant for EU AI Act); audit trail integrity (dual-ledger seal, OSCAL envelope, BBS+ selective disclosure); supply-chain compromise contingency; key-rotation procedure; incident response runbook.

## How wwt-tech maps to the whitepaper

---

wwt-tech doc	Maps to whitepaper section
gpu-compat-matrix.md	§B3 H100 / H200 cross-arch support
vendor-adapter-matrix.md	§B4 NVIDIA vs AMD vs Intel
denial-taxonomy.md	§C "Operational response" + Q&A §C "After detection, what does Cubie do?"
deployment-path.md	§E3 AI Proving Ground pilot + Q&A §E "How far from commercial?"
two-lane-architecture.md	§D4 workflow placement
not-ai-on-ai-positioning.md	§D5 deterministic infrastructure
shadow-mode-protocol.md	§E4 PoC metrics
risk-register.md	§C3 limitations

The wwt-tech set is the **operational counterpart to the whitepaper's formal-proof spine**. Together they form a complete answer.

## Part 3 — cubie-tep-v1.2.0 evidence pack (new, this analysis)

**Anchor file:** 02\_Cubie\_TEP\_Evidence\_Pack\_v1.2.0/cubie-tep-v1.2.0-MANIFEST.md in this delivery.

### What the cubie-tep-v1.2.0 pack proves

Same pattern as intc-v1.0: every whitepaper headline number has hashed JSON ground truth.

#### Cubie side

Claim	Number	Ground truth
Cubie reaches saturation peak	<b>60 / 60</b> runs across 20 seeds × 3 algorithms	cubie-tep-v1.2.0-sweep_*_seed*.json (10 representative + 50 in upstream data/tep/layouts/)
Cubie IDV-3 FDR	<b>100.0%</b> every run	final_metrics.idv3 in each sweep JSON
Cubie IDV-9 FDR	<b>100.0%</b> every run	final_metrics.idv9
Cubie IDV-15 FDR	<b>100.0%</b> every run	final_metrics.idv15
Cubie d00 FAR	<b>0.000%</b>	final_metrics.far
Cubie 60-seed walltime	<b>20.7s</b> (0.35s / run)	v5_cubie_vs_opus48.*.json → cubie.walltime_sec
Cubie API cost	<b>\$0.0000</b>	same path → cost_usd

## LLM side (head-to-head, same TEP test set)

Model	IDV-3 / 9 / 15	do0 FAR	Cost / 20 calls
Claude Haiku 4.5	50 / 0 / 0	0.0%	~\$0.03
Claude Sonnet 4.6	0 / 0 / 0	0.0%	~\$0.07
Claude Opus 4.7 (2026-05-28)	50 / 25 / 0	60.0%	\$0.80
Claude Opus 4.8 (2026-05-28)	50 / 0 / 0	40.0%	\$0.94
<b>Claude Opus 4.8 (2026-05-29 re-test)</b>	<b>50 / 0 / 0</b>	<b>40.0%</b>	<b>\$0.67</b>

Day-2 reproducibility: Opus 4.8 reproduced the 50/0/0 + 40% FAR result bit-for-bit on consecutive days. Cubie's numbers can't drift (structural identity); LLM's stability across two days is meaningful evidence the strict-token failure mode is not a one-day artifact.

### Opus 4.7 → 4.8 same-day delta

Axis	Opus 4.7	Opus 4.8	Delta
do0 FAR	60%	40%	-20pp (4.8 more conservative)
IDV-9 FDR	25%	0%	-25pp (4.8 worse)
Output tokens / call	17	115	6.6× explosion
Walltime / 20 calls	27.2s	32.4s	+19%
Cost / 20 calls	\$0.80	\$0.94	+18%

Consistent with 4.8 doing automatic internal reasoning even in strict-token mode. In extended-thinking mode 4.8 caught IDV-15 that 4.7 missed — latent capability is there, but only under extended reasoning.

## How cubie-tep-v1.2.0 grounds the whitepaper

The whitepaper now leads page 1 with the leaderboard table (after the 2026-05-29 head-to-head); every number in that table traces to a hashed JSON in this pack. The "Day-2 reproducibility re-run" paragraph explicitly cites `v5_cubie_vs_opus48_20260529_184625.json` (now `cubie-tep-v1.2.0-v5_cubie_vs_opus48_20260529_184625.json` in this pack).

## Part 4 — Cross-reference matrix

The three packs cover the spectrum: pain → solution → deployment. Mapped together:

Question reviewers will ask	Answer source	Specific file
Why is the GPU-cluster problem real?	intc-v1.0	intc-v1.0-tp4_nvlink_results.json (66.85% blast radius)
Why doesn't existing monitoring catch it?	intc-v1.0	intc-v1.0-hunter_results.json (87.53% DCGM blind spot)
Is the problem reproducible across cloud providers?	intc-v1.0	azure_2023 + azure_2024 + lmm_2025 + maveriq_unconstrained
Does Cubie's math actually work?	cubie-tep-v1.2.0	10 sweep JSONs + 60-run upstream
Is the result stable across days?	cubie-tep-v1.2.0	2026-05-28 + 2026-05-29 JSONs both = 50/0/0 + 40% FAR
Does Cubie beat the LLM frontier?	cubie-tep-v1.2.0	v5_cubie_vs_opus48_20260529_184625.json shows 100 vs 50/0/0
How does Cubie deploy in a real environment?	wwt-tech	deployment-path.md (4-stage rollout)
What does Cubie deny and why?	wwt-tech	denial-taxonomy.md (LOCK/JAM/SHATTER/DEFLECT/DENY)
Does Cubie work on AMD as well as NVIDIA?	wwt-tech	gpu-compat-matrix.md + vendor-adapter-matrix.md
Is this just another LLM-on-LLM model?	wwt-tech	not-ai-on-ai-positioning.md (deterministic infrastructure)
What's the customer-side validation rubric?	wwt-tech	shadow-mode-protocol.md (14-day pass/fail)
What are the known risks?	wwt-tech	risk-register.md (false-positive, bypass, human override, audit)

## Verification — one command per pack

```
# intc-v1.0 verification
cd C:\Users\NickV\cubie-tf\demo\evidence\intc-v1.0
Get-ChildItem intc-v1.0-*.json | ForEach-Object { (Get-FileHash -Algorithm SHA256 $_).Hash + " " + $_.Name }

# cubie-tep-v1.2.0 verification
cd C:\Users\NickV\Downloads\Cubie_Analysis_2026-05-30\02_Cubie_TEP_Evidence_Pack_v1.2.0
Get-ChildItem cubie-tep-v1.2.0-*.json | ForEach-Object { (Get-FileHash -Algorithm SHA256 $_).Hash + " " + $_.Name }

# wwt-tech sanity (file count + line count)
cd C:\Users\NickV\cubie-tf\demo\wwt-tech
(Get-ChildItem *.md).Count # should be 9
(Get-ChildItem *.md | Get-Content | Measure-Object -Line).Lines # should be ~744
```

## Part 5 — What this analysis adds

---

Three things the cubie-tep-v1.2.0 pack adds to the existing intc-v1.0 + wwt-tech foundation:

1. **Formal-proof anchor.** Where intc-v1.0 proves *the problem* with operational traces and wwt-tech proves *the deployment* with operational specs, cubie-tep-v1.2.0 proves *the math* with hashed JSON receipts for every formal-verification claim.
2. **Multi-day reproducibility.** intc-v1.0's numbers come from snapshot traces. cubie-tep-v1.2.0's day-2 head-to-head adds an inherent stability claim that conventional one-shot benchmarks can't make.
3. **Frontier-LLM context.** intc-v1.0 doesn't speak to LLMs; wwt-tech only addresses LLMs as workloads. cubie-tep-v1.2.0 is the first pack that puts Cubie head-to-head against frontier LLMs (Opus 4.8 included) on the exact same problem the customer will care about.

Together, the three packs form the complete answer to "show me this works before I want to take it down" — the ATC-evidence demand from the field conversations. Pain provable, solution provable, deployment provable, math provable, LLM-comparison provable. Every claim hash-anchored.

---

## Closing — versioning + APPEND-ONLY across all three packs

---

Pack	Update policy	Next version
intc-v1.0	APPEND ONLY (per its own MANIFEST)	intc-v1.1 alongside, never edit
cubie-tep-v1.2.0	APPEND ONLY (per this MANIFEST)	cubie-tep-v1.3 alongside, never edit
wwt-tech	Each MD versioned in git via commit history	New commit, never delete

This is the discipline that lets the whitepaper, the demo, and the field-conversation Q&A all cite stable, hashed sources six months from now without the "wait, did that number change?" problem.