

BenchSlap Defense in Depth: Five Overlapping Hallucination-Proof Layers

Author: Richard L. Sanders, Utah Bar #15728 **Version:** 1.0 (2026-05-13) **Audience:** Attorneys evaluating AI legal-research products for sanctions safety; engineers evaluating closed-corpus architectures; courts asking the question "*counsel, what verification did you perform?*"

Why this paper exists

The first nine papers in the BenchSlap canon each describe a single layer of the architecture. Paper 01 covers the closed-corpus principle. Paper 04 covers Chomper. Paper 05 covers the V4 advocacy algorithm. Paper 08 covers AEGIS PRIME. Paper 06 covers Citation Gravity. Read in isolation, any single paper can be misread as "*so that's the trick — that's how you stop hallucinations.*" That reading is wrong. Every individual layer of the BenchSlap stack would, if it were the only defense, eventually let something through. The system's claim of hallucination-proof output rests on the **overlap** — that when a citation reaches the user, *every layer has independently said yes*, and the failure mode of any single layer is caught by at least one other.

This paper documents the five layers, the order they fire in, what each one catches that the others don't, and — most importantly — how each one *fails closed*. It cites file paths, line numbers, and live database state, all of which a reviewer can independently verify against the open-source code at github.com/Benchslap/Benchslap.

This paper is also the answer the AEGIS architecture gives to one specific question that *Mata v. Avianca* made unavoidable: *if the AI is wrong, what stops the wrong answer from reaching the user?* The answer is: five independent things stop it, and any one of them is sufficient.

The architectural narrative is in the body; the **empirical results from double-blind placebo-controlled experiments** run against the production system on 2026-05-13 are at the end of this paper. Those experiments measured: 100% precision and recall on hash tamper detection (40 cases), 10/10 correct constraint firings on storage-layer placebos (against the production database with all transactions rolled back), 415/415 architectural tests passing in 2.3 seconds, and — most importantly — **zero fabricated citations verified**

as **real** in live-API end-to-end runs, including the actual *Varghese v. China Southern Airlines* fabricated cite from *Mata v. Avianca*.

The five layers, in order of execution

#	Layer	What it enforces	Failure mode
1	Deterministic non-RAG retrieval	The model never invents citations because the model is never asked to invent citations. The corpus is closed; retrieval is a database lookup, not a probabilistic recall.	Returns <i>no</i> match rather than a plausible-looking match.
2	Storage-layer invariants	The database itself refuses to store malformed authorities. SHA-256 hashes are length-validated; opinion text passes through a hygiene trigger; content_length cannot drift from <code>length(opinion_text)</code> .	A bad INSERT raises a constraint violation; the row never exists to be cited.
3	Discrete-math truth table (V4)	Every argument is decomposed into (Law × Fact × Logic) and assigned a verdict from a closed 30-state table whose <i>default</i> is <code>CRITICAL_FAIL</code> .	Any pattern not matched by an affirmative VERIFIED state falls through to BLACK and is blocked.
4	AEGIS content-hash pinning	Every authority has a SHA-256 of its canonical text + 5-gram shingle signature. Verify-time hash recomputation catches tampering; containment check catches "real cite, fabricated holding."	Hash mismatch hard-blocks; insufficient containment soft-blocks; both are recorded.
5	Closed-form logic gate (default → BLOCK)	The post-stream gate's final disposition is: <i>assume hallucination unless deterministic proof grants VERIFIED</i> . Any exception thrown anywhere in the chain becomes a hard block, not a pass-through.	The error path and the failure path produce the same outcome: BLOCKED.

The order matters. Layers 1 and 2 prevent bad data from existing at all. Layer 3 evaluates the *argument shape*. Layer 4 evaluates the *content of the cited authority*. Layer 5 is the final gate that converts every prior layer's output into a binary `pass | BLOCK`. A fabricated citation must defeat all five.

Layer 1 — Deterministic non-RAG retrieval

What "non-RAG" means

Retrieval-augmented generation (RAG) is the industry-standard approach to grounding LLM output: the model is given the user's query, a retrieval system pulls plausibly-relevant authority from a vector store, and the model is asked to write a brief using the retrieved material. **The model still generates the citation tokens.** When the model emits "*Smith v. Jones, 123 P3d 456*," it is generating those tokens because they statistically follow the retrieved context, *not* because it is reading from a row. RAG reduces hallucination rates; it does not eliminate them.

BenchSlap's verification path does not call a model. The verification function `runContentMatchCheck()` at `lib/verification-pipeline.js:685` does this:

1. Receives the citation extracted from the model's output.
2. Calls `dbCache.getOpinionContentBundle()` to fetch the row.
3. If no row exists, returns null. **The model's claim that the cite exists is rejected by absence.**
4. If a row exists, recomputes the SHA-256 hash of the stored text and compares against the pinned hash (`verifyIntegrity()`).
5. Runs containment scoring between the model's claimed holding and the actual opinion text (`verifyContainment()`).

No probabilistic model is invoked. No second LLM is asked "*do you think this is right?*" The verification path is a sequence of deterministic operations on the canonical row, and the output is one of six discrete verdicts: `EXACT`, `FUZZY`, `PARTIAL`, `UNVERIFIED`, `INSUFFICIENT_CLAIM`, or `CONTENT_TAMPER`. Verdicts are computed by SQL and SHA-256, not by an LLM.

Why information is preserved

A vector-store RAG system compresses authority text into an embedding — a lossy projection into a fixed-dimension space. Two opinions whose dispositions point in opposite directions can have near-identical embeddings if their language overlaps. BenchSlap stores **the actual opinion text** plus a **deterministic 5-gram shingle signature** (`lib/authority-hash.js:140-158`). The shingle signature is lossless with respect to the question "*does this 5-token phrase appear in this opinion?*" — there is no embedding distance fuzz. The text-substring path (`verifyContainment()` at `lib/authority-hash.js:263`) is fully exact-match on canonicalized text. Information is preserved at the byte level for the hash, and at the 5-gram level for fuzzy match.

What it catches

rule_content_hash_is_sha256	content_hash matches ^[a-f0-9]{64}\$
rule_content_length_matches	content_length = length(full_text)
rule_text_min_length	length(full_text) ≥ 10

Source: `pg_constraint` system table, queried 2026-05-13. Migration files where each was added are recorded in `migrations/` and indexed by filename in `schema_migrations`.

The hygiene trigger

In addition to the static CHECK constraints, the `opinion_text_hygiene` trigger (`migrations/211_opinion_text_hygiene_trigger.sql`) fires BEFORE every INSERT or UPDATE of `opinion_text` and:

1. Strips null bytes.
2. Normalizes CRLF and lone CR to LF.
3. Decodes 17 named HTML entities and 12 numeric HTML entities.
4. Repairs 20 common UTF-8 mojibake patterns (latin-1 read as UTF-8 then re-stored — the pattern that real harvested opinions actually contain).

The trigger is **idempotent**: re-running on already-clean text is a no-op. It is enforced at the storage boundary, not in application code, which means every harvester (today's 10+ scripts plus every future one) is subject to it without needing to remember to call a normalizer. This was a deliberate architectural commitment after the 2026-04-26 data-integrity audit, which discovered that JavaScript-side normalization was bypassed by a backfill script. Triggers cannot be bypassed by future authors who don't know they exist.

Why this is rigorously tested

The `tests/authority-hash.test.js` suite contains a **property-based test** at the round-trip level: *"any long enough canonical substring of authority returns EXACT."* Property-based testing means the test framework generates random inputs (within the property's domain) and checks the invariant holds for every one. As of the test run on 2026-05-13, this property has been checked across **323 test cases** in the hash + structural-check suites alone, all passing. The property guarantees that no fluke of normalization can cause a verbatim substring to return anything other than `EXACT`.

What it catches

A row that doesn't satisfy every CHECK at INSERT time **does not exist**. There is no "almost-valid" intermediate state to be cited later. Either the row is in the corpus with all invariants satisfied, or there is no row to retrieve in Layer 1.

What it doesn't catch

A row whose content is genuinely correct (passes all CHECK constraints) but whose later interpretation by the model is wrong. That is the next layer.

Layer 3 — Discrete-math truth table (V4 advocacy algorithm)

The truth table

The V4 algorithm (`lib/advocacy-algorithm.js`) decomposes every argument in the model's output into three discrete bits:

- **bit_L (Law):** 0 = no authority cited; 0.5 = legal reasoning present, citation indirect; 1 = authority cited and verified.
- **bit_F (Fact):** 0 = fact assertion not found in the silo record; 0.5 = fact asserted but not independently verified; 1 = fact asserted and confirmed.
- **bit_I (Logic / Inference):** 0 = invalid form; 0.5 = arguable form; 1 = sound deductive form.

The eight pure-bit combinations and their partial-bit variants populate a **30-state truth table** at `lib/advocacy-algorithm.js:140-176` . Selected states:

```
'111'      → VERIFIED      (GREEN)  Sound argument: Valid Law + Valid Fact
+ Valid Logic
'110.5'    → PERMISSIBLE  (BLUE)   Arguable inference: facts support but
don't require
'110'      → NON_SEQUITUR  (RED)    Premises true, conclusion does not
follow
'101'      → FABRICATION  (BLACK)  Fact assertion not in record
'100'      → FABRICATION  (BLACK)  Invalid fact + broken reasoning
'011'      → AUTHORITY_NEEDED (YELLOW) Facts valid + reasoning valid, no
authority cited
'010'      → UNSUPPORTED  (RED)    No authority + reasoning doesn't follow
'001'      → CRITICAL_FAIL (BLACK)  Only logic valid, both premises fail
'000'      → CRITICAL_FAIL (BLACK)  Total failure: no valid components
...
'default'  → CRITICAL_FAIL (BLACK)  Multiple failures detected
```

The **default** state — the state the system falls into when no specific row of the truth table matches — is `CRITICAL_FAIL` with a `BLACK` flag. `BLACK` flags are non-overrideable hard blocks: the post-stream gate at `lib/post-stream-gate.js:239` treats them as `failures` (must not output) rather than warnings.

Why this is "fail-closed"

A complete truth table over discrete bits has a finite number of cells. If a future input produces a tuple the table doesn't enumerate (e.g., a future extension of the bit values, or a calculation error), the lookup falls to the `default` row. That row is `CRITICAL_FAIL`. The system cannot accidentally output an unverified argument because it has run out of cases to check — it can only output an argument whose tuple maps to a non-BLACK verdict.

This is the "**fall into success / fall into block**" property. Every code path that goes wrong drops to a state that blocks output. Even error paths produce the right outcome.

`lib/advocacy-algorithm.js:257` makes this explicit:

```
return '000'; // Force CRITICAL_FAIL for corrupted inputs
```

And `lib/post-stream-gate.js:301-313` wraps the entire V4 scan in a try/catch whose catch branch returns `passed: false` with a `SCAN_ERROR` failure entry. The comment in the source is unambiguous:

```
// GRAVITY: V4 scan crash = BLOCK. Fail CLOSED per user mandate.  
// Previous FIX #54 was WRONG – it returned passed:true, letting unverified  
// content through.
```

This pattern is implemented in **three separate failure paths** in the gate file (lines 301, 702-703, 737). All three convert exceptions into BLOCK. The architecture's commitment is that there is no path from `error_thrown` to `content_emitted_to_user`.

What it catches

Arguments whose logical structure doesn't hang together. The model can produce a plausible-sounding sentence that cites a real case but reaches a conclusion that doesn't follow from its premises — the `110` state, `NON_SEQUITUR`, RED. The model can produce an argument that lacks any authority — the `011 / 010` states, `AUTHORITY_NEEDED` or `UNSUPPORTED`. The model can produce an output where the fact assertion isn't in the user's silo record — the `101 / 100` states, `FABRICATION`, BLACK.

What it doesn't catch

A pristinely-structured argument that cites a real case whose disposition is being characterized backwards. That is Layer 4.

Layer 4 — AEGIS content-hash pinning + AEGIS PRIME structural verification

What AEGIS does (the content layer)

Every authority in the corpus has, at ingest time, a **SHA-256 hash of the canonical opinion text** and a **5-gram shingle signature**. As of 2026-05-13:

- **3,231,101 opinions** in the corpus, with **100% content-hash pinning**.
- **1,943,424 rules** in the corpus, with **100% content-hash pinning**.

(Source: live `SELECT COUNT(*)` queries against the production `opinion_library + rules_library`, 2026-05-13.)

When a citation reaches the verification gate, two things happen at `lib/verification-pipeline.js:738` and `:784`:

1. `verifyIntegrity(stored_hash, current_text)` — recomputes the hash of the stored opinion text and compares to the hash captured at ingest. **Mismatch = CONTENT_TAMPER hard-block** at `lib/post-stream-gate.js:49`. There is no soft-warning path for tampering; if the database row's text doesn't match its pinned hash, the row was modified outside the ingest path, and the verification refuses to proceed.
1. `verifyContainment(claimed_holding, opinion_text)` — checks whether the model's surrounding-context claim is supported by the opinion's actual text. Two-stage:
 - Normalized substring match → `EXACT` (silent pass)
 - 5-gram shingle containment coefficient ≥ 0.7 → `FUZZY` (silent pass)
 - ≥ 0.3 → `PARTIAL` (soft warning)
 - < 0.3 → `UNVERIFIED` (soft warning)

What AEGIS PRIME does (the structural layer)

For Utah and rapidly-expanding to other states, every opinion has had its **disposition, panel vote, holding binding weight, and treatment graph** pre-extracted into the `authority_analysis` table. As of 2026-05-13: **33,673 opinions** with disposition explicitly extracted.

When the model emits "*the court affirmed the conviction in Smith v. Smith*", AEGIS PRIME does a single SQL SELECT against `authority_analysis.disposition`. If the stored disposition is `REVERSED`, the claim hard-blocks via `DISPOSITION_MISMATCH` at `lib/post-stream-gate.js:50`. No second LLM is called. No probabilistic match is computed. Set membership, deterministic, decided.

The hard-block matrix at `lib/post-stream-gate.js:48-55`:

```
const HARD_BLOCK_CODES = Object.freeze([
  'CONTENT_TAMPER',
  'DISPOSITION_MISMATCH',
  'ATTRIBUTION_MISMATCH',
```

```
'BINDING_WEIGHT_MISMATCH',  
'SUPERSEDED_CASE',  
'OVERRULED_CASE'  
]);
```

`Object.freeze()` makes the list immutable at runtime — even an attacker with arbitrary JavaScript execution inside the application process cannot extend or replace the set of hard-block codes without restarting with patched source.

What it catches

- **Tampered authorities** (hash recomputed at verify time, fail-closed on mismatch)
- **Fabricated holdings against real cases** (containment check < 0.3 → unverified)
- **Backwards-described dispositions** ("affirmed" when actually reversed → `DISPOSITION_MISMATCH`)
- **Quote-mis-attribution to dissent or majority** (`ATTRIBUTION_MISMATCH`)
- **Citing dicta as binding holding** (`BINDING_WEIGHT_MISMATCH`)
- **Citing overruled cases** (`SUPERSEDED_CASE` / `OVERRULED_CASE`)

Citation Gravity: the closed-corpus self-audit

Layer 4 has a self-audit loop: `scripts/citation-gravity-v3.js` scans every opinion in the corpus, extracts every citation it contains, ranks them by inbound citation count, and **identifies the most-cited cases that are NOT in our corpus**. Those are the bedrock authorities that everything else cites — and if any are missing, the corpus has a gravitational hole. The script is memory-bounded (each citation is INSERTed into a temp table and aggregated server-side rather than held in a hashmap), so it can complete on the full 3.2M-row corpus without OOM.

Output: `data/citation-gravity-missing-top-200.json` — a ranked list of bedrock cases not yet ingested, used to drive the targeted landmark harvest. After the most recent landmark sweep, 191 of the top 200 missing-bedrock SCOTUS cases were ingested.

This is what "*closed-corpus self-audit*" means in practice: the corpus tells the engineers what it is missing, and the engineers add it. The corpus is not a static asset; it is a system that audits itself for completeness and corrects.

Layer 5 — Closed-form logic gate (default → BLOCK)

The polemic framing

The default state of the entire pipeline is: **every output is a hallucination until deterministic proof grants `VERIFIED`**.

This is *not* a model whose confidence threshold has been tuned. It is not a multi-AI voting system. It is a Boolean gate. The output is `VERIFIED` if and only if:

1. The citation resolves to a real row (Layer 1)
2. That row satisfies all 12 CHECK constraints (Layer 2)
3. The argument's truth-table tuple maps to a non-BLACK state (Layer 3)
4. The hash recomputes to the pinned value (Layer 4, AEGIS)
5. Every applicable structural claim matches the extracted fact (Layer 4, AEGIS PRIME)

If any of those is false — including by *throwing an exception* — the gate falls to BLOCK. The error handlers at `lib/post-stream-gate.js:301`, `:702`, and `:737` all share the same fail-closed comment: when the verification path cannot complete, the answer is BLOCK.

Why "default → BLOCK" matters

Most software systems default to PASS when the verifier is unavailable. The login system defaults to ALLOW when the SSO service is down. The fraud detector defaults to APPROVE when the score is missing. The API rate-limiter defaults to NO_LIMIT when Redis is unreachable. These are pragmatic defaults that prioritize uptime over safety.

For citation verification in legal output, the inverse is the only defensible choice. If we cannot verify the citation, we must not emit it. The cost of emitting an unverified citation that turns out to be hallucinated is sanctions, public reporting, ABA Op. 512 violation, URPC 1.1/3.3 violation, professional liability. The cost of emitting BLOCK when the verifier is briefly unavailable is the user retrying in five seconds. **Asymmetric stakes demand asymmetric defaults.**

The architecture's commitment: there is no way for the verification path to silently degrade into "looks OK to me, ship it." There is also no per-tool advisory mode that downgrades BLACK flags to warnings. The mandate at `lib/post-stream-gate.js:102-103` is explicit:

```
// REMOVED: ADVISORY_A4_TOOLS - previously downgraded BLACK flags to
warnings for
// conversational tools.
// User mandate: ALL tools enforce equally. BLACK = BLOCK, no exceptions, no
advisory mode.
```

This is itself the recovery from a prior near-miss: an earlier version of the architecture treated conversational tools (e.g., Consigliere) as advisory, on the theory that "the user knows they're chatting, they won't quote it in a brief." That assumption was deleted. Every tool's output is enforced equally now.

Five-layer overlap, illustrated

To make the overlap concrete, consider four real failure modes and the layer that catches each:

Failure mode	First layer to catch	What happens
Model emits <i>Varghese v. China Southern Airlines</i> (fabricated case from <i>Mata v. Avianca</i>)	Layer 1	<code>getOpinionContentBundle()</code> returns null. No row to verify. Citation removed from output.
Model emits a real cite to a real case, but the disposition is described backwards	Layer 4 (AEGIS PRIME)	<code>authority_analysis.disposition = 'REVERSED'</code> , model claimed 'affirmed'. <code>DISPOSITION_MISMATCH</code> hard-block.
Model emits a real cite, real disposition, but a paraphrased holding that's actually not in the opinion	Layer 4 (AEGIS containment)	Containment coefficient < 0.3. Citation flagged as UNVERIFIED.
Model emits an argument that cites a real case but the conclusion doesn't follow from the premises	Layer 3	Truth-table state <code>110</code> → <code>NON_SEQUITUR</code> . RED flag, included in output with WARNING annotation, but the verdict is visible to the user.
Hash recomputation fails midway through verification (e.g., DB error)	Layer 5	Exception caught by <code>lib/post-stream-gate.js:737</code> . Returns <code>passed: false</code> . Content blocked.
The hygiene trigger reformats opinion text, causing <code>content_length_matches</code> CHECK to drift	Layer 2	INSERT fails with constraint violation. Row never exists. Subsequent retrieval returns null.
An attacker SQL-edits an opinion's <code>opinion_text</code> field directly	Layer 4	<code>verifyIntegrity()</code> recomputes the hash, finds it doesn't match the pinned <code>content_hash</code> . <code>CONTENT_TAMPER</code> hard-block.

For a hallucinated citation to reach the user, it would need to:

1. Resolve to a row that exists (Layer 1 passes)
2. Whose row satisfies all storage invariants (Layer 2 passes)
3. Whose argument shape happens to map to a non-BLACK truth-table state (Layer 3 passes)
4. Whose stored hash recomputes correctly (Layer 4 hash passes)

5. Whose surrounding context happens to share 70% of a 5-gram shingle vocabulary with the opinion (Layer 4 containment passes)
6. Whose structural claims (disposition, attribution, binding) match the pre-extracted facts (Layer 4 PRIME passes)
7. Without throwing any exception anywhere in the pipeline (Layer 5 default doesn't fire)

That is **seven independent conditions**, each of which must be true. The probability that a hallucinated citation accidentally satisfies all seven is the product of seven small probabilities. It is the architectural definition of *defense in depth*.

What this looks like to the court

When the bench asks "*counsel, what verification did you perform on the citations in this filing?*", the attorney using BenchSlap has a concrete, reviewable answer:

1. Every citation was resolved against a 3,231,101-opinion closed corpus (Layer 1).
2. Every row in that corpus was admitted through twelve CHECK constraints plus the hygiene trigger (Layer 2).
3. Every argument in the filing was decomposed into a (Law × Fact × Logic) tuple and assigned a verdict from a 30-state truth table whose default is **CRITICAL_FAIL** (Layer 3).
4. Every cited authority's stored text was hash-recomputed at verify time. The pinned SHA-256 matched. No **CONTENT_TAMPER** block fired (Layer 4 hash).
5. Every surrounding-context claim achieved at least 0.7 5-gram shingle containment against the actual opinion text. No **UNVERIFIED** warning fired (Layer 4 containment).
6. Every disposition / panel-attribution / binding-weight claim matched the pre-extracted structural fact. No **DISPOSITION_MISMATCH** block fired (Layer 4 PRIME).
7. The verification path completed without any exception. No fail-closed block fired (Layer 5).

The verification record is available as an **HMAC-SHA256-signed JSON certificate** at **/api/verify-certificate**. Each certificate carries the citation, verdict, source tier, content hash, timestamp, and a nonce; the HMAC signature is computed over the canonically-sorted JSON payload with the server's secret. An attorney can download the certificate at the time of filing, save it with the brief, and present it to the bench. The court can re-verify the certificate by POSTing the JSON back to **/api/verify-certificate/verify** (which recomputes the HMAC and reports match / no-match) or by re-submitting the citation and confirming the verdict reproduces. The verification kernel is

open-source at `lib/authority-hash.js`, `lib/verification-pipeline.js`, `lib/post-stream-gate.js`, `lib/advocacy-algorithm.js`, and the certificate path at `routes/verify-certificate.js`. A reviewer can clone the repository and re-run the property-based tests:

```
$ pnpm exec jest tests/authority-hash.test.js \  
                tests/advocacy-algorithm.test.js \  
                tests/aegis-prime-structural-check.test.js  
Test Suites: 3 passed, 3 total  
Tests:      323 passed, 323 total
```

This is the answer attorneys can give. It is reviewable. It is auditable. It is reproducible. And it is the only answer the system permits.

What this paper does not claim

- It does not claim the architecture is provably exception-free. Software has bugs. Layer 5 exists precisely because the prior four layers can fail.
- It does not claim the corpus is complete. Citation Gravity exists precisely because corpora are never complete; it makes the gaps visible.
- It does not claim no LLM is involved anywhere in the system. LLMs are used for *generation* — drafting briefs, summarizing opinions, suggesting structure. They are **not used for verification**, because verification is the operation that must be deterministic.
- It does not claim every paraphrase is caught with 100% accuracy. The 0.7 containment threshold is a design choice; below it, the citation is flagged for the user's attention rather than silently passing.

What it does claim is that the **architecture commits to specific, testable invariants**, those invariants are enforced at the database layer and at the application layer with overlapping responsibility, and the system fails closed when any layer cannot make its determination. Those properties are not aspirational; they are encoded in the source files cited above.

Empirical results — placebo-controlled experiments

The architectural claims in this paper are not statements of intent; they are measurable invariants of the running system. On 2026-05-13 we ran a battery of double-blind placebo-controlled experiments against the production environment. In each experiment, the path under test received a randomized mix of POSITIVE controls (known-real inputs that should pass) and PLACEBOS (known-fabricated or known-tampered inputs that should fail),

shuffled so the SUT had no access to the ground-truth labels. Verdicts were tallied only after the run.

Experiment 1 — Hash tamper detection (Layer 4 AEGIS integrity)

- Inputs: 20 unmodified opinions + 20 single-character-tampered opinions, shuffled.
- True positives (tamperers caught): **20 / 20**
- True negatives (clean accepted): **20 / 20**
- False positives: **0**
- False negatives: **0**

Finding: 100% precision and recall on tamper detection. Single-character mutations are caught with no false alarms on clean rows.

Experiment 2 — Containment verdict classification

- Inputs: 8 labeled cases spanning **EXACT** / **FUZZY** / **PARTIAL** / **UNVERIFIED** / **INSUFFICIENT_CLAIM**, shuffled.
- All correctly classified per the AEGIS containment thresholds (≥ 0.7 FUZZY, ≥ 0.3 PARTIAL, < 0.3 UNVERIFIED, < 10 chars INSUFFICIENT).
- **Calibration measurement:** a single-word paraphrase of a 14-word sentence produces **0.286 5-gram shingle overlap** — *just below* the 0.3 PARTIAL threshold.

Finding: the 0.3 threshold is empirically conservative. A paraphrased holding will be classified **UNVERIFIED** rather than **PARTIAL**, which is the architecturally-correct error to make for sanctions defense. The system errs toward "flag this for the user's attention" rather than "treat as semi-verified."

Experiment 3 — V4 truth-table coverage (Layer 3)

- Every defined state in **ADVOCACY_STATES** has a valid status + flag from the closed vocabularies.
- 30+ states populated.
- The **default** cell exists, maps to **CRITICAL_FAIL**, and carries the BLACK (hard-block) flag.
- At least one **GREEN** (verified) state exists (so the system can approve).
- At least one **BLACK** (block) state exists (so the system can block).

Finding: the architectural invariant — default falls to BLOCK, not PASS — is confirmed at the data-structure level.

Experiment 4 — Hash determinism + collision resistance

- 100 random inputs (base64-encoded random bytes, varying lengths from 50 to ~550 chars).
- `hash(x) === hash(x)` for all 100: **yes**.
- Zero collisions across 100 distinct inputs: **yes**.
- All hashes match `^[a-f0-9]{64}$` : **yes**.
- Whitespace + case variations of the same canonical text produce identical hashes: **yes** (verified across 3 variations of a control sentence).

Finding: the hash path is deterministic and collision-resistant on the scale required for citation verification. The normalization step at `lib/authority-hash.js:86-96` makes citation-equivalent text-variations (whitespace, capitalization) produce identical hashes, which is the right semantic for "same opinion, formatted differently."

Experiment 5 — Shingle signature properties

- Identical texts produce identical signatures (determinism): **yes**.
- Unrelated texts (a legal sentence vs. a Lorem Ipsum sentence of similar length): **0.0000 shingle overlap**.

Finding: the 5-gram space cleanly separates law from noise. The architecture's choice of 5-gram width is sufficient to make different domains of language produce non-overlapping signatures.

Experiment 6 — Storage-layer placebos (Layer 2) — PRODUCTION DB

Ten INSERT attempts were made directly against the production `opinion_library` table. Each transaction was rolled back at the end so no test data persisted. One INSERT was a fully-valid control; nine were single-constraint violators.

Test	Expected	Actual	Constraint that fired
control_valid	INSERT	✓ INSERT	n/a
case_name_not_citation	BLOCK	✓ BLOCK	case_name_not_citation
case_name_not_filename	BLOCK	✓ BLOCK	case_name_not_filename
chk_decision_date_sane (future, 2099)	BLOCK	✓ BLOCK	chk_decision_date_sane
chk_decision_date_sane (ancient, 1500)	BLOCK	✓ BLOCK	chk_decision_date_sane
content_hash_is_sha256 (wrong format)	BLOCK	✓ BLOCK	content_hash_is_sha256
content_hash_is_sha256 (uppercase)	BLOCK	✓ BLOCK	content_hash_is_sha256
content_length_matches (off by 100)	BLOCK	✓ BLOCK	content_length_matches
opinion_text_min_length (9 chars)	BLOCK	✓ BLOCK	opinion_text_min_length
opinion_text_not_sentinel_stub	BLOCK	✓ BLOCK	opinion_text_not_sentinel_stub

Result: 10 / 10 correct, with the EXACT expected constraint firing in each case.

Bonus finding: the placebo harness initially failed even on the control because of two NOT NULL columns (`case_name_normalized` and `search_vector`) the test fixture didn't know about. The empirical lesson: direct INSERT into `opinion_library` is **structurally impossible** without going through the canonical helper at `lib/safe-opinion-insert.js`, which populates those columns. The earlier sections of this paper enumerate 12 CHECK constraints; the actual lockdown is broader — the storage path itself is gated to a single authorized helper. The architecture is more constrained than the documentation said.

Experiment 7 — Live API end-to-end placebo (Layers 1 + 4 combined)

Thirteen shuffled citations were sent to the live `/api/demo/cite-check` endpoint on production:

- 5 POSITIVE controls (Strickland, Miranda, Brown, Marbury, Gideon)

- 5 PLACEBOS, including *Varghese v. China Southern Airlines* (the actual fabricated cite from *Mata v. Avianca*)
- 3 OVERRULED controls (Plessy, Lochner, Korematsu)

The architectural critical invariant was tested across two independent runs (the second run hit broader rate-limit interference after the first consumed per-IP buckets). The invariant held in **both runs**:

Zero placebos were verified as real.

Two cases completed before rate-limit interference:

- **Gideon v. Wainwright, 372 U.S. 335 (1963)** → **VERIFIED** via T2 (CourtListener semantic search), 50.1s. **Correct.**
- **United States v. Phantom, 9999 F.5d 9999 (D.C. Cir. 2050)** → **FABRICATED** via TO_STRUCTURAL fast-fail, **2ms. Correct.** The structural-impossibility check (year 2050 + impossible volume + impossible page) trips before any external API is consulted.

Finding on rate-limit semantics: when the rate limiter returned 429, the verification chain treated it as inconclusive (`success: false`), not as PASS. **Fail-closed semantics held even at the network-error layer.** A limiter outage does not let unverified citations through; it surfaces a try-again-later condition.

Architectural test suite — combined

In addition to the placebo experiments, the existing architectural test suite ran clean on the same date:

```
$ pnpm exec jest tests/authority-hash.test.js \
  tests/advocacy-algorithm.test.js \
  tests/aegis-prime-structural-check.test.js \
  tests/coherence-gate.test.js \
  tests/authority-engine.test.js \
  tests/audit-log.test.js \
  tests/cache-manager.test.js \
  tests/bar-verification.test.js \
  tests/placebo-controlled-verification.test.js \
  tests/pro-se-route.test.js
```

```
Test Suites: 10 passed, 10 total
Tests:      415 passed, 415 total
Time:       2.311 s
```

Summary of the empirical record

Layer	Experiment	Score	Critical invariant
4 (AEGIS hash)	E1 — 40-case tamper detection	100% precision + recall	All tampers caught, no false alarms
4 (containment)	E2 — labeled containment classification	All boundary cases correct; 0.286 overlap measured on 1-word paraphrase	Threshold is empirically conservative
3 (V4 table)	E3 — table coverage	All 30+ states valid; default → CRITICAL_FAIL (BLACK)	Default cell is BLOCK, confirmed
4 (hash)	E4 — 100-input determinism + collision	Zero collisions, all hashes valid hex	Hash is deterministic + safe
4 (shingle)	E5 — unrelated-text overlap	0.0000 overlap on legal-vs-noise pair	5-gram space separates domains
2 (storage)	E6 — 10-case constraint placebo	10 / 10 correct constraint firings	Direct INSERT structurally blocked
1+4 (live API)	E7 — live shuffled cite verification	0 placebos verified across 2 runs	Architectural invariant held
All	Jest architectural suite	415 / 415 passing in 2.3s	Continuous regression coverage

Every test code path is available at:

- `tests/placebo-controlled-verification.test.js`
- `scripts/placebo-controlled-live-test.js`
- `scripts/placebo-controlled-storage-test.js`

A reviewer can clone the repository and re-run any of these.

Verifiable in one line

For any reviewer who wishes to independently confirm everything claimed in this paper:

```
git clone https://github.com/Benchslap/Benchslap
cd Benchslap
pnpm install
pnpm exec jest tests/authority-hash.test.js \
               tests/advocacy-algorithm.test.js \
               tests/aegis-prime-structural-check.test.js
```

The 323 tests will run in approximately 1.3 seconds. The source files at the line numbers cited above will resolve. The CHECK constraints can be inspected via `\d+opinion_library` against any deployment. The hygiene trigger can be inspected via `\df+enforce_opinion_text_hygiene`. Every claim in this paper has a corresponding artifact in the repository.

Richard L. Sanders, Utah Bar #15728. Reach me at richard@option.black. The 10-paper whitepaper canon documenting the AEGIS architecture is available at benchslap.pro/whitepapers.