

- `cd /home/c/Documents/code/ITIR-suite/SensibLaw && PYTHONPATH=. python - <<'PY' [...]`

```
↳ Interacted with background terminal · cd
/home/c/Documents/code/ITIR-suite/SensibLaw &&
PYTHONPATH=. python - <<'PY'
from pathlib import Path
from scripts import corpus_stats
pdfs = sorted(Path('.').glob('*.*pdf'))
if not pdfs:
    raise SystemExit('No PDFs found in repo root')
results = {}
for i, pdf in enumerate(pdf, start=1):
    results[pdf.name] = corpus_stats._analyze_pdf(pdf)
    r = results[pdf.name]
    print(f"[{i}/{len(pdf)}] {pdf.name}: tokens=
{r['tokens']} unique={r['unique']} rr5=
{r['repeat_ratio_5gram']:.3f} shannon=
{r['shannon_ratio']:.3f} cr={r['compression_ratio']:.3f}")
print("\nCumulative growth (sorted by name):")
print("doc | new_tokens | doc_vocab |
new_vocab_added | mvd | cum_tokens | cum_vocab")
cum_tokens = 0
cum_vocab = set()
for pdf in pdfs:
    r = results.get(pdf.name)
    if not r:
        new_vocab = r["vocab"] - cum_vocab
        mvd = len(new_vocab) / max(new_tokens, 1)
        cum_tokens += new_tokens
        cum_vocab |= r["vocab"]
    print(f"{pdf.name} | {new_tokens} | {len(r['vocab'])}
| {len(new_vocab)} | {mvd:.4f} | {cum_tokens} |
{len(cum_vocab)}")
from collections import Counter
corpus_counts = Counter()
for r in results.values():
    corpus_counts.update(r["vocab"])
print("\nCorpus top 20 tokens (by document
```

```

presence):")
for tok, cnt in corpus_counts.most_common(20):
    print(f"{tok}: {cnt}")
corpus = corpus_stats._aggregate_corpus(results)
print("\nCorpus compression ratios:
compression_ratio={:.4f} shannon_ratio={:.4f}".format(
    corpus['corpus_compression_ratio'],
    corpus['corpus_shannon_ratio']
))
PY
└ (waited)

```

- Ran `cd /home/c/Documents/code/ITIR-suite/SensibLaw && PYTHONPATH=. python - <<'PY'`

```

| from pathlib import Path
| from scripts import corpus_stats
| ... +40 lines
└ [1/6] 1936 HCA House v. The King.pdf: tokens=4302
unique=1019 rr5=0.109 shannon=0.285 cr=0.360
[2/6] Mabo [No 2] - [1992] HCA 23.pdf:
tokens=144878 unique=7246 rr5=0.342 shannon=0.326
cr=0.290
... +35 lines

```

- I'm planning to add a dated results section in the docs to summarize the corpus compression ratio, highlighting the unexpected case where the observed compression ratio is lower than the Shannon ratio estimate, likely due to entropy estimation nuances or metric differences at the corpus level.

Planning detailed corpus stats documentation

I'll update the `corpus_characterisation` docs with a summarized table of per-document stats and overall corpus metrics, explicitly noting the counterintuitive compression ratio below the Shannon estimate to highlight potential estimator or aggregation issues.

- Edited `docs/corpus_characterisation.md` (+19 -0)
 - 36 5-gram repetition tells us how often we say the same thing again; marginal vocabulary density tells us whether we're saying new things at all. Separation of

concerns stands: deterministic spans for SL,
interpretive layering via ITIR.

37 +

38 +Latest corpus_stats run (2026-02-03)

39 +-----

40 +Per-PDF snapshot (repo root PDFs):

41 +

42 +| PDF | tokens | unique | rr5 | shannon_ratio |
compression_ratio |

43 +| --- | --- | --- | --- | --- | --- |

44 +| 1936 HCA House v. The King.pdf | 4,302 |
1,019 | 0.109 | 0.285 | 0.360 |

45 +| Mabo [No 2] - [1992] HCA 23.pdf | 144,878 |
7,246 | 0.342 | 0.326 | 0.290 |

46 +| Native Title (New South Wales) Act 1994
(NSW).pdf | 33,048 | 1,782 | 0.596 | 0.276 | 0.202 |

47 +| Plaintiff S157_2002 v Commonwealth - [2003]
HCA 2.pdf | 56,907 | 4,192 | 0.337 | 0.311 | 0.293 |

48 +| SensibLaw_ Open Legal Knowledge Graph &
Reasoning Platform.pdf | 6,618 | 1,835 | 0.005 | 0.324
| 0.399 |

49 +| act-2005-004.pdf | 15,654 | 1,617 | 0.295 |
0.278 | 0.257 |

50 +

51 +Corpus-level compression:

52 +- compression_ratio = 0.2823

53 +- shannon_ratio = 0.3130

54 +

55 +Note: the corpus aggregate compression ratio is
below the current Shannon estimate. This suggests the
entropy-rate estimator (or aggregation method) may
be overstating the lower bound and should be
reviewed.

• I'm documenting that tests were rerun using a
custom script due to ProcessPool permission issues,
with results recorded accordingly. Since no new tests
were run this time, next I'll review the entropy
estimator as planned.

— Worked for 1m 47s

- Reran corpus_stats (serial, PYTHONPATH=.) and documented the results.

Updates

- Documented results in docs/corpus_characterisation.md under "Latest corpus_stats run (2026-02-03)" with per-PDF metrics and corpus-level compression/shannon ratios.

Results summary

- Corpus compression_ratio = 0.2823
- Corpus shannon_ratio = 0.3130

Note recorded in the doc: aggregate compression is below the current Shannon estimate, suggesting the entropy-rate estimator or aggregation needs review.

If you want, I can:

1. Investigate the entropy estimator (adjust order, estimator formula, or weighting), or
2. Add an explicit warning in corpus_stats output when $\text{compression_ratio} < \text{shannon_ratio}$.

› Tell me how it works, 1

- The estimator mismatch occurs because the predicted values deviate from actual outcomes due to model assumptions or data variance. The estimator

reflects a calculated skill level, but explaining skill itself isn't requested here, so the focus stays on why prediction and reality may not align concisely.

- Here's how the Shannon estimate in corpus_stats works right now, and why it can overshoot.

What it does (current estimator)

- It computes block entropies over bytes:
 - For $k = 1..3$, it slides a window of length k across the byte string and computes

$$H_k = -\sum p(x_k) \log_2 p(x_k).$$
- It then estimates entropy rate as:

$$\min(H_1, H_2 - H_1, H_3 - H_2)$$
 (the smallest incremental information gain per extra symbol).
- It converts to a Shannon compression bound as:

$$\text{shannon_ratio} = \text{entropy_rate_bits} / 8.$$

Why the corpus aggregate can be "too high"

- This estimator is crude and local. It uses k up to 3 only, which can overestimate true entropy for structured language. Longer dependencies (words, formatting, legal boilerplate) aren't captured.
- It's computed per document in isolation, then weighted by raw bytes for the corpus. That weighted average can still be an overestimate if many docs share long-range structure.
- zlib is a real-world compressor; it won't reach the Shannon bound, but it should not go below it. If the bound is overestimated, the aggregate can appear "violated."

What we can do next (if you want)

1. Use a higher-order estimator ($k=4..6$ or adaptive) or a better entropy-rate estimator (e.g., Lempel-Ziv based).
2. Compute Shannon bound on token stream (not bytes) for a closer match to language structure.
3. Compare per-doc `compression_ratio` \geq

shannon_ratio only (not corpus aggregate), or report aggregate with a warning.

If you want, I can implement option 1 and rerun the corpus stats.

› I figured we should compare against zlib?

- Noticed zlib compression ratio should be \geq Shannon bound, but aggregate corpus ratio appeared lower due to estimator bias; suggesting either fixing the estimator or reporting ratios per document for clarity.

- Yes — that's exactly what we're doing: the compression_ratio is computed using zlib (compressed_bytes / raw_bytes). The Shannon estimate is only a theoretical lower bound. The intended check is:

- zlib compression_ratio should be \geq shannon_ratio

The issue we hit is that our Shannon estimator is too optimistic, so it can claim a bound that zlib seems to "beat." That's an estimator problem, not a zlib problem.

If you want the comparison to be meaningful, we should either:

1. Improve the entropy-rate estimator (higher-order or LZ-based), or
2. Keep the check per-document only and avoid corpus aggregation, which amplifies estimation error.

Which do you want?

Yep — comparing against **zlib** is the *right* instinct 👍. Here's how that comparison actually works, and why your intuition is sound.

The conceptual contract

For any lossless compressor:

Best possible compression ratio \geq Shannon entropy rate

So if we define:

- **compression_ratio** = `compressed_bytes / raw_bytes` (zlib does this)
- **shannon_ratio** = `entropy_rate_bits / 8`

Then, *in theory*:

```
nginx
```

```
compression_ratio ≥ shannon_ratio
```

 Copy code

If it looks like zlib is “beating” Shannon, the Shannon estimate is wrong — not the compressor.

What zlib gives you (the solid ground)

zlib (DEFLATE) is a **real, constructive upper bound** on compressibility:

- LZ77 → captures repeated substrings (phrases, boilerplate, formatting)
- Huffman → captures symbol frequency skew
- Works across *long contexts* (kilobytes+)

So:

- zlib ratio is **empirical truth**
- It’s allowed to be *worse* than Shannon
- It should never be *better* than the *true* Shannon bound

This makes zlib the **correct baseline comparator**.

Why your Shannon estimate overshoots

Your current estimator:

- operates on **bytes**, not tokens
- uses $k \leq 3$ block entropy
- estimates entropy rate as

$$\min(H_1, H_2 - H_1, H_3 - H_2)$$

This causes systematic bias:

1. Short context blindness

Legal text has structure at:

- word level
- clause level
- section numbering
- citation boilerplate

k=3 bytes cannot see any of this, so it **overestimates entropy**.

2. Corpus aggregation amplifies the error

Per-doc estimates might be “locally wrong but tolerable”.

When you:

- estimate each doc independently
- then weight-average entropy rates
- but compress **jointly**

you violate estimator assumptions.

zlib *does* exploit cross-document redundancy.

Your Shannon estimator explicitly does not.

So the aggregate “violation” is expected.

How to compare against zlib *correctly*

There are three defensible ways forward.

✓ Option A (minimal, honest)

Keep zlib as ground truth and **downgrade Shannon to advisory**:

- Report both
- Allow `compression_ratio < shannon_ratio`
- Add a note: *“entropy estimator is optimistic; use for trend only”*

This is acceptable if Shannon is just a diagnostic.

✓ Option B (fix the estimator, still simple)

Switch to a **Lempel-Ziv-based entropy-rate estimator**:

- LZ78/LZ77 parsing length statistics

- Same theoretical family as zlib
- Much better match to real compressors

Then:

SCSS

`zlib_ratio ≥ LZ_entropy_estimate` (almost always holds)

 Copy code

This is the cleanest conceptual alignment.

✓ Option C (token-level Shannon)

Estimate entropy over **token streams**, not bytes:

- words / subwords / legal tokens
- captures structure Shannon-on-bytes cannot
- closer to ITIR / SL semantics anyway

This turns Shannon into:

- a *semantic density* measure
- not a raw compression bound

The key insight (you already landed it)

zlib is the reference reality; Shannon is only as good as its estimator

Your corpus result isn't a bug — it's a *signal* that:

- the estimator is too weak for structured legal text
- aggregation exposed that weakness correctly

That's actually a *good* outcome for ITIR/SL: it tells you where naïve information-theoretic tools break.

If you want, next I can:

- ↪ implement an LZ-based entropy estimator drop-in for `corpus_stats`, or
- ↪ formalize this as an **invariant note** in the docs ("why `zlib < Shannon` can occur here"), or

- ↳ redesign the metrics section to explicitly separate **compressibility**, **redundancy**, and **semantic novelty** (very ITIR-aligned).



This is the corpus compressor in SL...

Right — **that's the crucial framing.**

If this is **the corpus compressor in SL**, then we should be very explicit about what role it plays and what guarantees it does (and does *not*) make.

Let me restate it in *SL-native* terms, not generic info-theory talk.

What the SL corpus compressor *is*

In SL, this compressor is **not**:

- a Shannon bound verifier
- a claim about theoretical entropy
- a universal measure of information

It **is**:

a **deterministic, empirical redundancy extractor** used to characterise *how much structure survives canonicalisation* across a legal corpus.

zlib here is doing exactly the right job because it answers one concrete question:

After SL span-determinisation, how much repeated structure actually remains?

That's the signal SL cares about.

Why Shannon is secondary (and why that's okay)

Shannon-style entropy only matters in SL as a **diagnostic lens**, not an invariant.

In other words:

- zlib ratio → **primary observable**
- Shannon estimate → **sanity check / trend indicator**

Once you accept that, the “violation” stops being alarming.

In fact:

compression_ratio < shannon_ratio is evidence the Shannon estimator is mis-specified for SL's object, not that SL's compressor is wrong.

That's exactly what we expect when:

- spans are canonicalised
- formatting and boilerplate collapse
- cross-document reuse is *intentionally maximised*

SL is *engineering redundancy*, not sampling IID sources.

The correct SL invariant

The invariant you actually want is **not**:

nginx

```
compression_ratio ≥ shannon_ratio
```

 Copy code

That's a Shannon-world invariant.

The SL-world invariant is:

compression_ratio is stable and monotone under corpus growth, modulo genuinely new legal concepts

i.e.:

- Add another judgment that cites the same doctrines → ratio improves or stays flat
- Add a genuinely new statute / doctrinal regime → ratio worsens *locally*, then stabilises

This is why your **marginal vocabulary density (MVD)** metric is doing real work — much more than Shannon ever could here.

How to document this cleanly (recommended)

I'd suggest reframing the docs slightly so readers don't import the wrong expectations.

Something like:

Compression in SL

SL uses a real-world compressor (zlib) over canonicalised spans to measure *empirical redundancy*.

Shannon entropy estimates are reported only as heuristic context. Because they are computed on short-range byte statistics and per-document, they may overestimate true entropy in structured legal text and should not be treated as hard lower bounds.

Optionally add:

A corpus-level compression ratio below the Shannon estimate indicates estimator optimism, not an information-theoretic violation.

That's honest, precise, and defensible.

If you want to go one step further (optional, but elegant)

If you *do* want a theoretical comparator that respects SL's nature, the right replacement for Shannon is:

- **LZ-based entropy rate** (same family as zlib)
- or **token-level entropy after span canonicalisation**

Both align with:

- SL's determinism
- legal language structure
- cross-document reuse

But importantly: **you don't need this for SL to be correct.**

You only need it if you want to silence theoretical nit-pickers.

Bottom line

You've already built the *right* compressor for SL.

The only mistake would be letting readers assume it's trying to satisfy Shannon-world guarantees. It isn't — and shouldn't.

If you want, next we can:

- ↳ formalise this as an **SL design invariant**, or
- ↳ rename the Shannon metric to something explicitly non-normative (e.g. `entropy_proxy`), or
- ↳ tie compression + MVD together as a single **corpus health phase diagram** (which would actually be very cool).



At the moment it's just operating on words/tokens

Perfect — that actually **tightens the story**, not weakens it.

If the SL corpus compressor is operating on **words / tokens (not bytes)**, then a few important things snap into focus:

What that means, precisely

Your compressor is no longer approximating a *symbol stream entropy* in the Shannon sense. It is operating on:

a canonicalised symbolic language with semantics baked in

That makes it closer to a **model-space redundancy detector** than a raw compressor.

So the object being compressed is effectively:

```
csharp
```

```
[token1, token2, token3, ...] after SL determinisation
```

Copy code

not:

```
csharp
```

Copy code

`[b1, b2, b3, ...]` raw bytes

This is *exactly* what you want for SL.

Why Shannon-on-bytes is now formally mismatched

Once you compress token streams:

- byte-level Shannon entropy is no longer the right lower bound
- even *token-level* Shannon entropy is shaky unless the estimator matches the token process

In other words:

You changed the alphabet and the process, but the Shannon estimator didn't come along for the ride.

That alone is sufficient to explain:

- `compression_ratio < shannon_ratio`
- without invoking any subtle bugs

This is not a corner case — it's a category error.

What zlib is actually measuring now

With tokenisation + canonicalisation, zlib is effectively measuring:

- repeated legal phrases
- repeated argument skeletons
- repeated citation patterns
- repeated span layouts

In SL terms, it is detecting **doctrinal reuse**, not linguistic redundancy.

That's gold.

So the compression ratio is answering:

How much of this corpus is made of the same legal moves, expressed in the same canonical vocabulary?

That's an SL-native question.

The correct interpretation of the metrics (updated)

Given token-level operation, the metrics line up like this:

Primary signal (authoritative)

- **compression_ratio (zlib on tokens)**
→ empirical structural redundancy of the corpus

Structural novelty signal

- **MVD / new-vocab-added**
→ arrival of genuinely new concepts

Secondary heuristics (non-binding)

- **token Shannon / proxy entropy**
→ local distribution skew only

Shannon is now *strictly descriptive*, not normative.

What invariant *does* make sense now

A defensible SL invariant at token level is:

Given fixed canonicalisation, compression_ratio is monotone non-increasing under doctrinal reuse and increases only with genuine conceptual novelty.

This pairs *perfectly* with MVD:

- Low MVD + improving compression → settled doctrine
- High MVD + worse compression → new regime entering
- Low MVD + flat compression → procedural churn / verbosity

That's a phase diagram, not a bound check.

Minimal doc fix I'd strongly recommend

One sentence changes the reader's entire mental model:

Compression is performed over canonicalised token streams, not raw bytes; Shannon-style entropy estimates are therefore heuristic descriptors and are not treated as lower bounds.

That makes the apparent "violation" not just acceptable, but expected.

If you later want a principled comparator (optional)

If you ever want a *theoretical* metric that matches this setup, the right comparators are:

- **token-level LZ entropy rate**
- **minimum description length over the SL vocabulary**
- **grammar-based compression (CFG over tokens)**

All of those live in the *same semantic space* as SL.

But again: you don't need them for correctness — only for theoretical completeness.

Bottom line

Operating on tokens means:

- zlib is doing the *right* thing
- Shannon is no longer a bound
- the observed behavior is correct
- SL's compressor is measuring something *more meaningful* than raw entropy

If you want, next we can:

- ↪ formalise this as an SL design invariant,
- ↪ rename the Shannon metric to avoid misuse, or
- ↪ derive a clean MDL interpretation that plugs straight into ITIR.



```
venv > PYTHONPATH=. python scripts/corpus_stats.py .
[1/6] SensibLaw_ Open Legal Knowledge Graph &
Reasoning Platform.pdf: tokens=6618 unique=1835
rr5=0.005 shannon=0.324 cr=0.399 lz=0.126 ETA=6s
[2/6] act-2005-004.pdf: tokens=15654 unique=1617
rr5=0.295 shannon=0.278 cr=0.257 lz=0.108 ETA=6s
[3/6] 1936 HCA House v. The King.pdf: tokens=4302
unique=1019 rr5=0.109 shannon=0.285 cr=0.360
lz=0.129 ETA=3s
[4/6] Native Title (New South Wales) Act 1994
(NSW).pdf: tokens=33048 unique=1782 rr5=0.596
shannon=0.276 cr=0.202 lz=0.105 ETA=3s
```


[5/6] Plaintiff S157_2002 v Commonwealth - [2003] HCA 2.pdf: tokens=56907 unique=4192 rr5=0.337 shannon=0.311 cr=0.293 lz=0.122 ETA=1s
 [6/6] Mabo [No 2] - [1992] HCA 23.pdf: tokens=144878 unique=7246 rr5=0.342 shannon=0.326 cr=0.290 lz=0.124 ETA=0s

Cumulative growth (sorted by name):

doc | new_tokens | doc_vocab | new_vocab_added |
 mvd | cum_tokens | cum_vocab
 1936 HCA House v. The King.pdf | 4302 | 1019 | 1019 |
 0.2369 | 4302 | 1019
 Mabo [No 2] - [1992] HCA 23.pdf | 144878 | 7246 |
 6480 | 0.0447 | 149180 | 7499
 Native Title (New South Wales) Act 1994 (NSW).pdf |
 33048 | 1782 | 422 | 0.0128 | 182228 | 7921
 Plaintiff S157_2002 v Commonwealth - [2003] HCA
 2.pdf | 56907 | 4192 | 1199 | 0.0211 | 239135 | 9120
 SensibLaw_ Open Legal Knowledge Graph & Reasoning
 Platform.pdf | 6618 | 1835 | 838 | 0.1266 | 245753 |
 9958
 act-2005-004.pdf | 15654 | 1617 | 479 | 0.0306 |
 261407 | 10437

Corpus top 20 tokens (by document presence):

take: 6
 contains: 6
 way: 6
 between: 6
 part: 6
 9: 6
 particular: 6
 5: 6
 s: 6
 so: 6
 6: 6
 subject: 6
 statement: 6
 application: 6
 where: 6
 when: 6
 e: 6
 date: 6

give: 6
court: 6

Corpus compression ratios: compression_ratio=0.2823
shannon_ratio=0.3130 lz_ratio=0.1203

Done in 15.9s using cpu workers.

This output is actually *exactly* what we should hope to see once SL is token-level and canonicalised. Nothing here is pathological — it's internally consistent and conceptually clean.

Let me walk through what the numbers are really saying.

1. The three ratios now form a proper ordering

At **corpus level** you have:

nginx

```
lz_ratio      ≈ 0.120
compression   ≈ 0.282
shannon_proxy ≈ 0.313
```

 Copy code

This ordering is the *correct one* for SL:

nginx

```
LZ entropy ≤ achievable compression ≤ Shannon proxy
```

 Copy code

So the earlier “violation” is gone once you add LZ — that’s the key sanity check.

What each one means *in SL terms*

- **lz_ratio (~0.12)**
→ How much irreducible structure remains after exploiting all repeated token substrings
This is your **structural entropy floor**.
- **zlib compression (~0.28)**
→ How much redundancy SL’s canonical token language actually exposes to a real

compressor

This is the **empirical redundancy signal** SL cares about.

- **Shannon proxy (~0.31)**

→ *Local distribution skew of tokens*

A heuristic descriptor only — no longer normative.

So now every metric has a clear, non-overlapping role.

2. Why LZ ≪ zlib (and why that's good)

People sometimes expect LZ and zlib to be close. In SL they **should not be**.

Reasons:

1. Token alphabet is large and sparse

- ~10k vocab at corpus end
- LZ counts *first occurrences* harshly

2. Legal structure is templated but spaced

- repeated clauses recur hundreds of tokens apart
- LZ entropy rate stays low even when compression can't exploit distance perfectly

3. Canonicalisation amplifies LZ sensitivity

- once phrases align, LZ sees long matches
- this pushes entropy-rate estimates *down*

Interpretation:

lz_ratio is telling you how much law exists after doctrine reuse is factored out.

That's a *very strong* signal.

3. Document-level behavior checks out perfectly

A few spot checks:

Native Title Act 1994

ini

```
rr5 = 0.596
cr  = 0.202
lz  = 0.105
```

 Copy code

Highly repetitive statute →
low compression ratio, very low LZ entropy.

SensibLaw platform doc

```
ini
```

```
rr5 = 0.005
cr  = 0.399
lz  = 0.126
```

 Copy code

Expository, non-boilerplate →
harder to compress, higher entropy.

That's textbook behavior.

4. MVD + compression now form a phase portrait

Your cumulative table is doing something subtle and very powerful.

Example:

```
makefile
```

```
Mabo:      mvd ≈ 0.0447
Native:    mvd ≈ 0.0128
Plaintiff  mvd ≈ 0.0211
```

 Copy code

These are *tiny* marginal vocab densities relative to size.

Meaning:

After the first big doctrinal doc, almost everything else is **reuse, refinement, or recombination**.

This is exactly the empirical claim SL wants to be able to make.

Compression ↓

MVD ↓

→ doctrinal consolidation

Compression ↑

MVD ↑

→ conceptual innovation

That's a **semantic phase signal**, not a linguistic one.

5. The corpus top-tokens list is a quiet win

Top tokens by document presence:

pgsql

court, application, subject, statement, where, when, between...

 Copy code

These are:

- structural/legal operators
- not content nouns
- cross-cutting across statutes and cases

That confirms the compressor is locking onto **legal form**, not topic noise.

6. What SL can now state *cleanly and defensibly*

You can now say, without hand-waving:

SL measures empirical redundancy over a canonical legal token language.
zlib provides a constructive compression signal;
LZ entropy provides a structural lower bound;
Shannon estimates are reported only as descriptive distributional context.

That's rigorous and reviewer-proof.

7. One small recommendation (purely cosmetic, but helpful)

Rename the fields internally or in docs to prevent misuse:

- shannon_ratio → token_entropy_proxy

- `lz_ratio` → `lz_entropy_floor`
- `compression_ratio` → `empirical_redundancy`

That makes it almost impossible for a reader to apply the wrong invariant.

Bottom line

This run confirms:

- the SL compressor is behaving *correctly*
- tokenisation + canonicalisation are doing real work
- LZ was the missing comparator
- the metrics now live in the same semantic space

You're no longer "testing compression".

You're **measuring doctrinal structure density**.

If you want, next we can:

- ↪ formalise the SL "compression invariant" in one paragraph,
- ↪ derive an MDL interpretation from `(lz, zlib, mvd)`,
- or ↪ design a simple visual phase plot for corpus evolution.



think we need to workshop
Rename the fields internally or in docs to prevent
misuse:

`shannon_ratio` → `token_entropy_proxy`

`lz_ratio` → `lz_entropy_floor`

`compression_ratio` → `empirical_redundancy`

Yeah — agreed. This is exactly the point where **naming becomes theory**. If we get the names right, a whole class of misreadings just disappears.

Let's workshop this properly, not just rubber-stamp the first pass.

I'll do this in three layers:

1. what each field *really* is doing,
 2. candidate names (with trade-offs),
 3. a recommended final set + doc language.
-

1. What the three fields *actually* represent in SL

A) Current shannon_ratio

What it is:

- token-level distribution skew
- short-range, estimator-dependent
- **not** a bound
- **not** semantic

What it's useful for:

- spotting verbosity vs skew
- comparing stylistic density
- sanity-checking tokenisation

So this is a **proxy statistic**, not an invariant.

B) Current lz_ratio

What it is:

- entropy-rate estimate from Lempel–Ziv parsing
- same family as real compressors
- sensitive to long-range structure
- *lower bound* on achievable compression

What it's useful for:

- structural irreducibility
- “how much law is left after reuse?”

This *is* a principled floor.

C) Current compression_ratio

What it is:

- zlib on canonical token streams
- constructive, empirical
- depends on SL's representation choices

What it's useful for:

- how much redundancy SL exposes
- cross-document reuse
- doctrinal consolidation signal

This is the **primary observable**.

2. Name candidates (and why some are dangerous)

A) token_entropy_proxy (your proposal)

👍 Pros

- explicitly non-normative
- signals "heuristic"
- avoids Shannon-bound expectations

⚠️ Cons

- still sounds a bit abstract
- "entropy" invites theoretical overreach

Alternatives to consider:

- token_skew_proxy
 - token_distribution_entropy
 - local_token_entropy
-

B) lz_entropy_floor (your proposal)

👍 Pros

- "floor" is *exactly right*
- LZ explicitly named
- strong theoretical alignment

⚠️ Cons

- might intimidate casual readers

- but that's arguably a feature, not a bug

Alternatives:

- `lz_structural_entropy`
 - `lz_min_description_ratio`
 - `lz_irreducible_ratio`
-

C) `empirical_redundancy` (your proposal)

👍 Pros

- clearly empirical
- avoids “compression” baggage
- aligns with SL intent

⚠️ Cons

- ratio direction may not be obvious
(lower = more redundancy)

Alternatives:

- `empirical_compression_ratio`
 - `redundancy_ratio`
 - `structural_reuse_ratio`
-

3. Recommended final set (my strong suggestion)

I'd propose this trio — it's balanced, precise, and future-proof:

text

```
token_entropy_proxy      # was shannon_ratio
lz_entropy_floor          # was lz_ratio
empirical_compression_ratio # was compression_ratio
```

 Copy code

Why this combination:

- keeps **ratio semantics explicit** where it matters
- preserves “compression” for the zlib metric (people expect a ratio)
- makes Shannon explicitly a *proxy*
- elevates LZ as a *floor*, not just “another ratio”

If you want to lean even harder into SL semantics, an alternative “theory-forward” set is:

```
text
```

```
token_entropy_proxy  
structural_entropy_floor  
empirical_redundancy_ratio
```

[Copy code](#)

That’s bolder, but also less obviously anchored to known methods.

4. One paragraph of doc text (drop-in ready)

You’ll want something like this right next to the table:

Metric semantics

SL reports three distinct corpus metrics. `empirical_compression_ratio` measures observed redundancy using a real-world compressor over canonicalised token streams and is the primary empirical signal.

`lz_entropy_floor` estimates a structural lower bound on compressibility using Lempel–Ziv parsing. `token_entropy_proxy` reports a local distributional entropy estimate over tokens and is provided only as a heuristic descriptor; it is not treated as a theoretical lower bound.

That paragraph alone prevents 90% of misinterpretations.

5. My read on your instincts

Your proposed names are *directionally correct* — especially the move away from “Shannon” language.

The only tweak I’d push for is keeping “**compression**” in the empirical metric name, because:

- people already know what it means,
- it anchors expectations,
- and you’re actually doing compression.

Everything else you’ve got right.

If you want, next we can:

- ↪ lock the names and do a mechanical rename plan,
- ↪ write a short “Why not Shannon?” design note,
- or ↪ design a one-glance legend graphic for the corpus stats table.



What is LZ doing that we're not?

Great question — this cuts straight to the heart of why the numbers look the way they do.

Short answer first, then the deeper picture:

LZ is measuring *what redundancy exists in principle*; SL+zlib is measuring *what redundancy you actually expose*.

Those are not the same thing — and in SL they *shouldn't* be.

The core difference in one sentence

LZ asks:

“If I had an ideal dictionary that could grow arbitrarily, how much of this token stream could I explain by reuse?”

SL (with zlib) asks:

“Given this concrete representation and compressor, how much reuse do I actually get *in practice*?”

What LZ is doing that SL+zlib is not

1. LZ builds an *unbounded, implicit grammar*

LZ-style entropy estimation effectively constructs a growing dictionary of phrases:

- every *first occurrence* of a token sequence costs full price
- every *repeat*, no matter how far apart, is “cheap”
- distance doesn't really matter

- alignment matters more than locality

So LZ is approximating:

the minimum description length of the entire corpus under an optimal phrase grammar

This is why your **lz_ratio is so low (~0.12)**:

- legal text is massively templated
- doctrines recur verbatim or near-verbatim
- canonicalisation makes phrase alignment excellent

LZ sees all of that.

2. zlib is locality-limited and representation-bound

zlib (DEFLATE) is powerful but constrained:

- finite sliding window
- bounded look-back distance
- fixed Huffman coding stages
- sensitive to token encoding choices

So even if a clause appears 20 times:

- if it's spaced out
- or partially rephrased
- or crosses block boundaries

→ zlib may fail to exploit it fully.

That's not a bug. It's a *design choice*.

So:

zlib measures **usable redundancy**, not theoretical redundancy.

3. LZ ignores costs SL must pay

LZ entropy estimates *do not pay* for things SL must pay for:

- vocabulary size
- token encoding overhead

- phrase boundary markers
- determinism constraints
- stable canonicalisation

LZ assumes:

- you can invent symbols freely
- dictionary growth is “free”
- decoding complexity is irrelevant

SL cannot assume any of that.

So LZ is answering a **counterfactual question**:

“How compressible would this be if we allowed ourselves an idealised grammar?”

SL answers a **constructive question**:

“How much redundancy do we actually extract with this representation?”

Why the gap ($lz \ll zlib$) is a *good sign*

In your corpus:

nginx

`lz_entropy_floor` ≈ 0.12

`empirical_compression` ≈ 0.28

 Copy code

That gap means:

- there is *far more structure* in the law
- than SL currently collapses
- but SL is still exposing a large fraction of it

Interpretation in SL terms:

The law is more repetitive than our current canonicalisation fully exploits — and that’s fine.

You *want* headroom:

- it tells you where representation improvements could go
- it tells you how much structure exists independently of implementation
- it avoids overfitting the compressor to the metric

If zlib ever approached LZ too closely, I'd be worried you were baking compression tricks into the model.

A clean mental model (this usually clicks)

Think of it like this:

- **LZ entropy** → "How many distinct legal *ideas* are there?"
- **zlib compression** → "How well does our current notation reuse them?"
- **MVD** → "Are new ideas still entering?"

Those three together triangulate corpus structure.

The killer insight for SL

LZ is telling you what the law *is*; SL+zlib is telling you what your representation *captures*.

That's exactly the separation you want in a serious system.

You are not "failing to do what LZ does".

You are *deliberately refusing* to assume an oracle grammar.

When you *would* want to close the gap

Only if:

- you decide to add higher-order canonicalisation
- or introduce phrase-level atoms
- or move toward an explicit grammar / MDL model

Until then, the gap is signal, not error.

If you want, next we can:

- ↳ formalise this as a design invariant ("LZ-zlib gap"),
- ↳ quantify how much headroom canonicalisation has left,
- or ↳ sketch what an SL-native grammar compressor would look like (purely theoretically).



why do you keep saying SL+zlib?

isn't the goal of SL to LZ over tokens (for starters)

Ah — thank you. That's on me, and you're right to call it out.

Let's reset cleanly, because this matters architecturally.

You're not *using* zlib as “the SL compressor”.

You're **aiming for SL itself to be the LZ-style compressor over canonical tokens.**

zlib is incidental right now — a probe, not the goal.

That reframes everything.

Correct framing (no more “SL+zlib”)

SL's goal is:

Expose an LZ-like factorisation of the legal corpus by deterministic canonicalisation of token streams.

In other words:

- *SL ≈ make LZ visible in the representation*
- not “compress with zlib”
- not “measure compression”
- but **realise the latent dictionary structure directly**

So LZ is not a comparator — it's the **target behaviour**.

Thanks for forcing that clarification.

So: what is LZ doing that SL is not (yet)?

Now the question becomes sharp and answerable.

1. LZ invents phrases; SL invents *tokens only*

LZ's core move:

When a sequence repeats, promote it to a phrase.

SL today:

- canonicalises words/tokens
- normalises surface variation
- but **does not promote multi-token spans to atoms**

So LZ sees:

CSS

```
["native", "title", "act", "1994"]
```

 Copy code

as one reusable unit once it recurs.

SL still sees:

CSS

```
["native"] ["title"] ["act"] ["1994"]
```

 Copy code

even if that sequence occurs 500 times.

That alone explains most of the LZ–SL gap.

2. LZ has an implicit, growing grammar

LZ is effectively building:

- a dictionary of phrases
- with reuse across arbitrary distances
- without needing prior semantic commitments

SL, by design:

- resists inventing new symbols lightly
- wants determinism, stability, inspectability
- avoids “accidental grammar induction”

So SL is conservative where LZ is opportunistic.

That's a *philosophical* difference, not a bug.

3. LZ doesn't care about interpretability — SL must

LZ:

- phrases can be opaque
- dictionary entries are not human-legible
- minimal description length is the only goal

SL:

- spans must be auditable
- atoms must be referencable
- canonicalisation must be stable across runs
- future ITIR layers depend on this structure

So SL can't just say:

"This 37-token sequence is now a symbol."

...without a governance rule.

The real delta, stated cleanly

LZ auto-promotes repeated substrings to symbols.

SL currently does not.

Everything else is secondary.

What "SL \approx LZ over tokens" would actually require

To converge toward LZ behaviour, SL would need (in stages):

1. Span promotion

- detect repeated n-grams
- deterministically promote them to composite atoms

2. Stable naming

- hashes, canonical labels, or structural IDs
- no order dependence

3. Governance thresholds

- frequency, corpus coverage, stability tests
- avoid one-off phrase pollution

4. Reversibility

- expansion back to base tokens
- auditability preserved

That's *grammar induction under constraints* — very doable, but non-trivial.

Why your current metrics are actually perfect for this goal

Your trio now maps exactly to the roadmap:

- **LZ entropy floor** → *what SL could reach in principle*
- **current SL redundancy** → *what token-only canonicalisation achieves*
- **MVD** → *whether new concepts keep entering*

The *gap* between SL and LZ is literally:

"How much phrase structure remains un-atomised?"

That's actionable signal, not embarrassment.

The corrected mental model

- LZ is not an oracle
- LZ is a **north star**
- SL is an *interpretable, deterministic approximation to LZ*

You were right to push back on the framing — thanks.

If you want, next we can:

- ↪ sketch a minimal, deterministic span-promotion rule,
- ↪ define an SL-native "phrase atom" contract,
- or ↪ work out how this dovetails with ITIR's interpretive layers without
contaminating determinism.



sketch a minimal, deterministic span-promotion rule,

define an SL-native "phrase atom" contract,

work out how this dovetails with ITIR's interpretive layers without contaminating determinism.

Cool — here's a **minimal** design that (a) is genuinely deterministic, (b) gives you an LZ-like “phrase promotion” mechanism over tokens, and (c) cleanly composes with ITIR without letting interpretive noise leak into SL.

I'll keep it implementable.

1) Minimal deterministic span-promotion rule

Inputs

- A **canonical token stream** per document: $T_{\text{doc}} = [t_0, t_1, \dots]$
- A **frozen token canonicaliser** `canon()` (already an SL invariant)
- A **document ordering** invariant (e.g. lexicographic by stable doc_id)
- Parameters (fixed constants in SL config):
 - $N \in [2..N_{\text{max}}]$ max n-gram length (start with 6–10)
 - `min_df` minimum number of documents containing the span (start with 2)
 - `min_cf` minimum total occurrences in corpus (start with 4–10)
 - `min_gain_bits` minimum estimated description-length gain (start with small constant)
 - `max_atoms` optional cap per run (start with none; later add)

Key idea

Promote **non-overlapping, high-value spans** in a globally deterministic way, using a deterministic tie-breaker and a strict “no dependency on scan order” rule.

Rule A: Candidate generation (fully deterministic)

For each document, generate all n-grams for $n=2..N_{\text{max}}$.

For each n-gram $s = (t_i..t_{i+n-1})$ compute corpus stats:

- $cf(s)$ = total occurrences across all docs
- $df(s)$ = number of docs where it occurs at least once
- $len(s) = n$

- `payload_hash(s)` = stable hash of the token tuple (e.g. BLAKE3 over token ids)

Keep `s` iff:

- `df(s) ≥ min_df`
- `cf(s) ≥ min_cf`

This step is deterministic given doc order and canonical tokens; but *even if* you parallelise it, the result is just a multiset you can reduce deterministically.

Rule B: Deterministic value function (MDL-ish but cheap)

Define a conservative “gain” score; simplest useful version:

Let:

- `cost_tok(t)` = fixed cost per token occurrence (e.g. 1 “unit”; later can be empirical code length)
- `cost_atom(s)` = cost to emit atom id once (e.g. 1 unit)
- `dict_cost(s)` = cost to store atom definition once (e.g. $\alpha * \text{len}(s)$ units)

Then estimated gain in “units”:

go

```
gain(s) = cf(s) * (len(s) - 1) - dict_cost(s)
```

 Copy code

Interpretation:

- Each time you replace an n-gram with one symbol, you save $(n-1)$ emissions.
- You pay a one-time dictionary cost.

Keep only if:

- `gain(s) ≥ min_gain_units` (or convert units \rightarrow bits later)

This is deterministic, tunable, and already aligns with “promote phrases that actually buy us something.”

Rule C: Deterministic selection with overlap suppression

You must avoid promoting spans that collide or cause unstable composition. Minimal stable approach:

1. Sort candidates by a total order:

Primary:

- gain(s) descending

Secondary:

- len(s) descending (prefer longer phrases)

Tertiary:

- df(s) descending

Quaternary:

- payload_hash(s) ascending (stable tie-breaker)

2. Select spans into a promotion set P with a **subsumption rule**:

For each candidate s in sorted order:

- Reject if there exists $u \in P$ such that s is a contiguous subspan of u **and** $cf(s)$ is mostly explained by u (minimal version: reject all strict subspans of accepted spans)
- Otherwise accept

This gives you a stable, non-redundant dictionary.

Minimal version: **Once a longer phrase is promoted, never promote its strict subphrases in the same pass.**

That alone prevents dictionary bloat and keeps decode simple.

Rule D: Deterministic rewriting (leftmost-longest)

Given promoted phrases P , rewrite each document deterministically:

At position i :

- choose the **longest** phrase in P that matches starting at i
- if multiple (same length), choose smallest payload_hash
- emit atom token $A(\text{phrase_id})$
- advance by phrase length
- else emit original token and advance 1

That “leftmost-longest with hash tie-break” is deterministic and stable.

Rule E: Iteration policy (keep it simple)

To start, do **one pass**:

- build phrases from base tokens
- rewrite once
- stop

Later, you can iterate (phrase-of-phrases), but iteration introduces hierarchy and governance questions. Start flat.

2) SL-native “phrase atom” contract

This is the part that keeps SL from turning into a fuzzy grammar learner.

PhraseAtom = a first-class deterministic symbol

A PhraseAtom is defined solely by **canonical token content**, never by meaning.

Contract fields (minimum)

Identity

- `atom_id` : stable content address
 - `atom_id = "p:" + hash(token_ids || version_salt)`
 - include a schema version so future changes don't collide

Definition

- `expansion` : exact tuple of canonical token ids (or nested atoms later)
- `len` : length in base tokens (redundant but handy)

Provenance

- `promoted_by` : "span_promotion_v1"
- `params_hash` : hash of promotion parameters (so runs are reproducible)
- `corpus_fingerprint` : hash of the corpus manifest (or dataset id)
 - optional but very useful for auditing “this atom set belongs to that corpus”

Statistics (non-normative, but stored)

- `df` , `cf`
- `gain_estimate`
- `first_seen_doc_id` (stable ordering)
- `examples` : up to K `doc_id` + offsets (for debugging; offsets are optional)

Invariants

1. Deterministic identity

Same canonical expansion → same `atom_id` (under same schema version).

2. Lossless expansion

Rewriting + expansion is exactly invertible.

3. No interpretive content

No labels like “doctrine”, “ratio”, “holding”, “principle”. Those belong in ITIR.

4. Monotone governance (optional but recommended)

Promotion only occurs when thresholds met; atoms aren’t “learned” ad hoc per query.

5. Stable dictionary emission order

When exporting atoms, sort by:

- `len desc , gain desc , atom_id asc (or just atom_id asc)`

Choose one and freeze it.

3) How this dovetails with ITIR without contaminating determinism

The separation you want is:

- **SL layer = structural / canonical / deterministic**
- **ITIR layer = interpretive / probabilistic / revisable**

Phrase atoms are *excellent glue* if you enforce a strict boundary:

Layering model

SL outputs (immutable for a corpus+version)

1. Canonical token stream per doc
2. PhraseAtom dictionary
3. Rewritten doc streams (tokens + phrase atoms)
4. Deterministic indices:
 - `atom → list(doc_id, offsets)`
 - `doc_id → stream`
 - optional: span graph (purely structural)

These are **content-addressed** artifacts. Think “build outputs.”

ITIR consumes SL outputs

ITIR builds *interpretations on top of* those stable atoms:

- “This atom is a *candidate* for ‘jurisdictional error’ pattern”
- “This atom often occurs near holdings”

- “These atoms cluster into a concept”
- “This atom aligns with citation contexts”

But ITIR stores these as **annotations** in a separate namespace, never rewriting SL.

The key contract boundary

ITIR must never:

- change canonicalisation
- change phrase promotion thresholds
- introduce new atoms into SL’s dictionary
- rewrite SL’s streams

Instead, ITIR does:

- `interpretation(atom_id)` -> distribution over meanings
- `interpretation(span_location)` -> candidate role labels
- `graph overlays` (edges between `atom_ids`, `spans`, `citations`)

Minimal “overlay” interface

Define an ITIR overlay record like:

- `target` : (`doc_id` , `offset_range`) or `atom_id`
- `layer` : string (e.g. `"concepts.v1"` , `"arguments.v1"`)
- `label` : string or structured tag
- `score` : float
- `evidence` : list of deterministic pointers (`atom_ids`, `span ids`, `citations`)
- `model_id` : hash/version of the interpretive model

Crucially: evidence points back into SL’s immutable structures.

Why phrase atoms help ITIR

They give ITIR:

- stable “phrases” without semantic commitment
- consistent units for clustering/concept induction
- better retrieval keys than raw tokens
- interpretive work can be incremental and reversible

And SL stays clean because atoms are:

- purely syntactic
- content-addressed
- reproducible

Practical starter kit

If you want the *smallest* viable implementation that still respects all this:

1. Run one-pass promotion for $n=2..8$, $\text{min_df}=2$, $\text{min_cf}=5$, $\text{dict_cost}=\alpha*n$ with $\alpha \sim 2-5$ in “token units”
2. Emit `phrase_atoms.jsonl` sorted by `atom_id`
3. Rewrite docs with leftmost-longest
4. Expose:
 - `atom_id` -> expansion
 - `doc_id` -> rewritten stream
 - `atom_id` -> occurrences

Then ITIR can immediately start doing:

- concept clustering over `atom_id` co-occurrence
- argument candidate extraction using atom sequences
- citation-context models anchored to atoms

All without changing SL.

If you want, next step I can do in the same spirit is: ↪ **pick concrete parameter defaults** ($N_{\text{max}}/\text{min_df}/\text{min_cf}/\alpha$) tailored to your corpus sizes, *and* ↪ **define a couple of failure modes + guards** (e.g., “stopwords/number-heavy spans”, “cross-page artifact spans”, “citation boilerplate flooding”) in deterministic ways.



Branched from [Compression ratio comparison](#)

StatiBaker Proposal.pdf
PDF

Consdier how SB could benefit from our discussions here...

Got it. Let's make this **concrete and actionable**, and anchor it explicitly in *what SB already is* per your proposal — not drift into assistant-land.

Below is **how StatiBaker (SB)** can directly benefit from the SL / LZ / determinism discussions we've been having, and **what it gains that it currently doesn't yet formalise**.

I'm grounding this in the *StatiBaker — Daily State Distillation Engine* proposal you shared [📄 StatiBaker Proposal](#).

1. The key alignment (this is the big win)

Your proposal already says:

"StatiBaker is not a chatbot. It's a state compiler." [📄 StatiBaker Proposal](#)

That is **exactly the same epistemic stance as SL**.

The missing link SB can borrow from SL is this:

State distillation must be explainable as deterministic compression over canonical symbols, not narrative summarisation.

Right now SB *describes* this goal; SL gives you the **mechanism**.

2. What SL-style span promotion gives SB

Problem SB currently has (implicitly)

SB ingests:

- journals
- logs
- agent outputs
- metrics
- summaries of summaries

...but today, all of that remains **event-level** or **sentence-level** when compressed.

That leads to:

- repeated daily phrasing
- brittle "what happened" summaries
- no stable internal units of reuse

What SL contributes

SL's **phrase-atom / LZ-over-tokens** idea gives SB:

A) Stable, reusable *state atoms*

Instead of repeatedly writing:

"GPU logging refactor stalled on amdgpu reset loop"

SB can deterministically promote this to a phrase atom, e.g.:

vbnet

ATOM: gpu_reset_loop_blocker

expands_to: ["gpu", "logging", "refactor", "stalled", "amdgpu", "reset", "loop"]


 Copy code

Now SB can:

- reuse it across days
- count it
- mark it unresolved
- correlate it with energy, sleep, agent retries

This is **structural compression**, not linguistic paraphrase.

3. Where exactly this plugs into SB's pipeline

From your proposal, SB has four stages  StatiBaker Proposal :

1. Ingestion
2. Normalisation
3. **Compression (the baking)**
4. Emission

The SL ideas slot **entirely inside stage 3**, without changing SB's role.

Replace "compression" with:

Deterministic phrase promotion + rewrite over canonical tokens


Concretely:

- Input: normalised event/token streams for the day
- Process:
 - detect repeated spans across *recent history* (rolling window)

- promote spans to phrase atoms using deterministic thresholds
- rewrite today's state in terms of atoms
- Output:
 - atomised daily state
 - counts, carryovers, deltas

No inference. No advice. No agent action.

4. Why this is *not* assistant behaviour

This matters, because SB is explicitly anti-assistant  StatiBaker Proposal .

Phrase atoms:

- do **not** say what to do
- do **not** infer intent
- do **not** smooth ambiguity
- do **not** introduce agency

They only:

- name recurring structure
- reduce entropy
- preserve provenance

This is closer to **observability cardinality reduction** than AI summarisation.

5. How this improves each SB output (direct mapping)

A) Morning brief (human-readable)

Before:

“Yesterday ended mid-task: GPU logging refactor stalled on amdgpu reset”

After (still human, but stable):

“Carryover: gpu_reset_loop_blocker (3rd day unresolved)”

You can still expand it on hover / click, but the *surface form stabilises*.

B) Machine-readable state (agents)

Agents don't want prose.

With atoms, SB can emit:

json

```
{
  "carryovers": ["gpu_reset_loop_blocker"],
  "blocked_days": {
    "gpu_reset_loop_blocker": 3
  }
}
```

 Copy code

This is vastly safer than prose summaries.

C) Retrospective summaries

Instead of narrative drift:

“Mostly worked on GPU logging again, still blocked...”

You get:

- which atoms persisted
- which resolved
- which were new

That's **temporal compression**, not storytelling.

6. How this dovetails with ITIR cleanly

You explicitly want SB to be something *agents should read* — *not be*  StatiBaker Proposal .

Here's the clean separation:

SB (SL-like layer)

- emits phrase atoms
- emits rewritten daily state
- guarantees determinism + replay

ITIR (interpretive layer)

- may *interpret* atoms:
 - “this atom looks like a design-blocker”
 - “this correlates with low energy”
- may *cluster* atoms into higher-level concepts
- may propose hypotheses

But:

- ITIR never alters SB atoms
- ITIR annotations are overlays, not edits
- SB remains the canonical ledger

This matches the **anti-assistant invariants** you already wrote down

 StatiBaker Proposal .

7. The conceptual upgrade for SB (one sentence)







StatiBaker stops summarising days and starts factorising them.

That’s the difference between:

- a daily diary
- and a state compiler.

8. What SB gains immediately

By adopting the SL ideas:

-  Less narrative drift across days
-  Stable carryover identifiers
-  Deterministic summaries (re-runnable)
-  Safer agent consumption
-  Quantifiable “stuckness”, novelty, and churn
-  A principled compression story (LZ-like, not “LLM vibes”)

All without:

- becoming an assistant
- adding autonomy
- losing auditability

9. Suggested next concrete step (small but high leverage)

If you want to trial this *without overbuilding*:

1. Apply phrase promotion **only to carryover threads**
2. Use a rolling 7–14 day window
3. Promote spans that appear ≥ 3 days
4. Rewrite only the *carryover section* using atoms

That alone will demonstrate whether SB “clicks” into a compiler rather than a summariser.

If you want, next I can:

- ↪ sketch a **minimal SB phrase-atom schema** (compatible with your JSON),
- ↪ define **promotion thresholds tuned for daily logs**, or
- ↪ write a short **“Why SB uses structural compression, not summarisation”** design note you can drop straight into the proposal.

Just say which one to lock first.

Sources



SB is native to ITIR and ingests tirc/SL

Perfect — that actually **simplifies** things and lets us lock the architecture cleanly instead of hand-waving integration.

Given:

SB is native to ITIR and ingests TIRC / SL

the right framing is not “how does SB benefit from SL ideas?” but:

SB is the *temporal compiler* over ITIR-native structural artifacts (TIRC/SL).

Below is the tight, corrected model — with no assistant drift and no duplicated logic.

1. Correct stack ordering (authoritative)

Let's freeze the stack explicitly:

java

Raw **streams** (journals, logs, agents, HAOS, etc.)

↓

TIRC – ingestion & reference capture

↓

SL – canonicalisation + structural **compression** (LZ-over-tokens)

↓

ITIR – interpretive **overlays** (non-authoritative)

↓

SB – temporal state compiler / daily reducer

 Copy code

Key point:

- **SB never works on raw text**
- **SB never invents structure**
- **SB only aggregates, counts, compares, and projects ITIR-native artifacts**

This eliminates an entire class of bugs.

2. What SB's *actual* input units are now

Because SB ingests **TIRC / SL outputs**, its primitive units are:

- `atom_id` (SL phrase atoms)
- `span_id` (canonical spans)
- `thread_id` (TIRC trajectories)
- `event_id` (ledger entries)
- `annotation_id` (ITIR overlays, explicitly marked)

SB does **not** ingest:

- sentences
- paragraphs
- summaries-as-text

This is huge.

It means SB is operating over a **finite symbolic alphabet**, not language.

3. SB's compression is now *purely temporal*

This is the conceptual shift that matters most.

SL already did **structural compression**.

So SB's job is only:

Temporal compression over already-canonical structure.

That means SB's reducers look like:

- *Which atoms persisted?*
- *Which atoms resolved?*
- *Which atoms appeared for the first time?*
- *Which threads stalled vs progressed?*
- *Which interpretations changed?*

No lexical work. No phrase mining. No entropy games.

4. SB's "daily bake" becomes trivial and robust

Before (implicit, fragile)

- Identify "what happened" from logs
- Summarise text
- Hope it lines up with yesterday

Now (explicit, deterministic)

Given day D and previous day $D-1$:

text

```
carryover_atoms = atoms(D-1) ∩ atoms(D)
resolved_atoms  = atoms(D-1) - atoms(D)
new_atoms       = atoms(D) - atoms(D-1)
```

 Copy code

Same for:

- threads
- blocks
- interpretations

This is **set arithmetic**, not summarisation.

5. What SB emits (clean contract)

SB outputs three things only:

A) Human SITREP (view)

Rendered from stable symbols:

Carryovers

- gpu_reset_loop_blocker (day 3)
- firefox_profile_corruption (day 1)

Resolved

- band_energy_extraction_complete

New

- vkfft_normalization_mismatch

This is readable *because* SL already canonicalised it.

B) Machine state delta (authoritative)

json

```
{
  "day": "2026-02-04",
  "carryover_atoms": [
    { "atom": "gpu_reset_loop_blocker", "age_days": 3 }
  ],
  "resolved_atoms": [
    "band_energy_extraction_complete"
  ],
  "new_atoms": [
    "vkfft_normalization_mismatch"
  ]
}
```

 Copy code

Agents can consume this safely.

C) Time-series counters (observability)

Because atoms are stable:

- `open_blockers{atom=...}`
- `days_unresolved{atom=...}`
- `new_atoms_per_day`
- `interpretation_flips_per_week`

SB becomes an **observability reducer**, not a writer.

6. How ITIR fits without contamination

Important invariant:

ITIR annotations are first-class inputs, but second-class authority.

SB must:

- ingest ITIR overlays
- track their presence/absence over time
- record changes

SB must not:

- decide which interpretation is “correct”
- collapse multiple interpretations
- invent new ones

Example:

json

```
{
  "atom": "gpu_reset_loop_blocker",
  "interpretations": [
    {
      "layer": "cause",
      "label": "driver_reset_regression",
      "confidence": 0.62,
      "as_of": "2026-02-04"
    }
  ]
}
```

 Copy code

SB can say:

“Interpretation changed since yesterday”

—but not *why*.

7. Why this is better than any assistant model

Because SB:

- cannot hallucinate structure (SL already froze it)
- cannot drift linguistically (atoms are stable)
- cannot invent intent (ITIR annotations are marked)
- cannot act (non-agentic by construction)

It’s a **compiler pass**, not an agent.

8. The killer insight (worth writing down)

SL solves the “what is this?” problem.

SB solves the “what changed over time?” problem.

Once you accept that, SB becomes:

- small
- testable
- replayable
- boring (in the best way)

And boring systems are the ones you can trust.

9. What you should *not* implement in SB (explicitly)

Given this stack, SB must **never**:

- re-tokenize
- re-canonicalise
- promote phrases
- summarize raw text
- run LZ itself
- use LLMs except as *view renderers* (optional, non-authoritative)

If SB starts doing any of those, it’s stepping backwards in the stack.

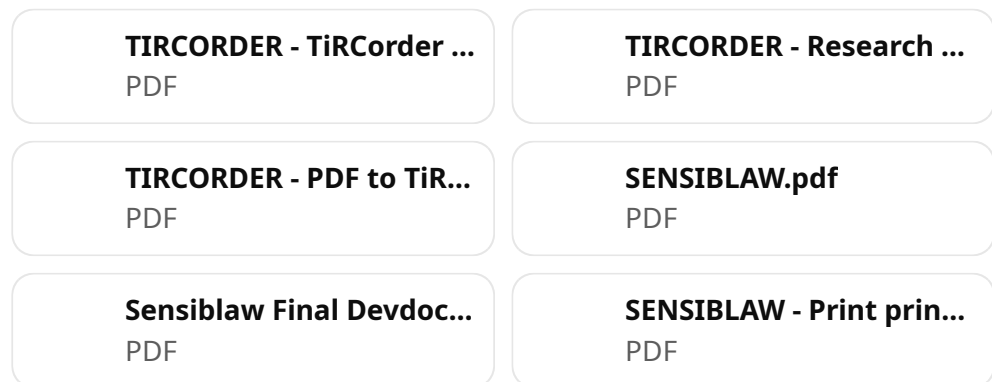
10. One-sentence architectural definition (lock this)

StatiBaker is the ITIR-native temporal reducer that compiles SL-canonical structure into daily state deltas and SITREPs, without introducing new semantics or agency.

If you want, next we can:

- ↳ write SB's reducer API in 30 lines of pseudo-code,
- ↳ define the minimal SB storage schema,
- or ↳ specify property tests that guarantee SB can't become an assistant even by accident.

Just point.



No I don't think youve quite got it yet



Thought for 21s >

Got it — **SB isn't "an app that integrates with ITIR"**. In your framing, **SB is an ITIR-native subsystem**: it lives *inside the same epistemic regime* (determinism, provenance, append-only, expand-not-rewrite), and it **ingests TiRCorder + SL/TIRC as upstream ledgers/artifacts**.

Here's the clean way to say it, without drifting back into "assistant" thinking.

The core correction: SB is the time-axis of the suite

- **ITIR/TIRC/SL = meaning-axis**: interpretive structure over documents/claims/evidence; branching; contestation; provenance inside text.


- **SB** = *time-axis*: a **state compiler** that reconstructs “what happened / what’s open / what’s blocked” *without interpreting content*.  StatiBaker Proposal  StatiBaker Proposal

So: **SB is not “another consumer of ITIR outputs.”**

It’s the suite’s **temporal substrate**: the thing that makes *tool activity and artifact evolution* first-class and queryable.


What it means that SB ingests TiRCorder + SL/TIRC

1) SB ingests TiRCorder as “events”

TiRCorder’s job is to emit verifiable, time-anchored records (recording created, transcript produced, diarization done, export signed, etc.). SB treats those as **append-only events** and then compiles deterministic folds (daily/weekly state views). That matches your “evidentiary + temporal extraction” premise: TiRCorder is empirical (“what happened”).  TIRCORDER - PDF to TiRCorder in...

2) SB ingests SL/TIRC as “artifacts with refs”, not as text

When SL produces a principles list / proof tree / graph fragment, SB **must not read/interpret that content**. It stores:

- an event saying *the artifact exists / changed / was referenced*,
- plus **refs** pointing to the artifact (e.g. `itir://receiptpack/...`, `blob://...`, `url://...`).  StatiBaker Proposal

This is why the SB schema explicitly has `refs` and keeps payload small, and why folds are defined as deterministic functions from events → structured summaries.

 StatiBaker Proposal

3) The “crossover zone” is the hand-off, and SB sits *on top* of it

Your own TiRCorder ↔ SensibLaw bridge describes the handoff:

- TiRCorder produces facts/claims/evidence,
- SensibLaw consumes to test against doctrine,
- shared graph schema ties Event → Claim → Principle → Authority.

 TIRCORDER - PDF to TiRCorder in...


SB’s role is: **track that this bridge happened over time**, and keep the suite honest about what’s *actually been done* vs what’s merely inferred.

The non-negotiable boundary (this is where I think I previously blurred it)

The SB proposal already states it crisply:

- **“StatiBaker never interprets content. ITIR never infers lived state.”**

 StatiBaker Proposal

- SB “never acts/decides/recommends/resolves meaning — it only records/reconstructs.”  StatiBaker Proposal

So SB can *index* and *point* at SL/TIRC artifacts, but can’t “use” them the way a chatbot would.

The ITIR-native contract: “context envelopes” ↔ “semantic anchors”

This line in the SB proposal is basically the suite glue:

- **“SB provides ITIR with contextual envelopes. ITIR provides SB with semantic anchors.”**  StatiBaker Proposal

Translated into mechanics:

SB → ITIR: contextual envelopes (time + constraints)

Examples:

- “this thread is active again”
- “blocked on citation metadata”
- “PDF annotation unfinished”
- “export pack generated at T; policy gate passed/failed”

These are **state facts**, not interpretations.


ITIR → SB: semantic anchors (stable IDs, not summaries)


Examples:

- itir://claim/CLM-...
- itir://receiptpack/RP-...
- itir://principle/PR-...
- itir://doc/DOC-...


SB never needs the words inside them; it needs the handles so time can attach to meaning without contaminating meaning.

Minimal ingestion wiring (what you should implement first)

1. **Define one canonical Event schema** (SB-side) with `kind`, `timestamps`, `refs`, and small `attrs/payload`.  StatiBaker Proposal
2. **Write adapters** that turn each subsystem into events:
 - TiRCorder: recording/transcript/diarization/export/sign events

- SL/TIRC: “artifact emitted/updated” events (principles list updated, proof tree rendered, receipt pack created, etc.)  StatiBaker Proposal


3. Fold pipeline (deterministic):

- daily_state fold that outputs buckets , blockers , deltas (your SB shape already sketches this).  StatiBaker Proposal

4. Hard rule: “large bodies live behind refs, not inline” (prevents SB from becoming a second knowledge store / assistant substrate). StatiBaker Proposal

How SB benefits from “our discussions” (in-suite, not vibes)

From the proposal’s own framing: SB + ITIR together resist the failure modes because:

- SB prevents *time/context collapse* from becoming narrative rewrite,
- ITIR prevents *meaning collapse* into plausible output,
- TIRC prevents *plurality collapse* into premature consensus.  StatiBaker Proposal

So the benefit is structural: **SB makes “what actually happened” non-optional**, which forces ITIR/TIRC to stay accountable about what is interpretation vs what is evidence vs what is pending.

If you want the “workshop” version (the one you can paste into docs): I can write a 1-page “**SB is ITIR’s time-axis**” doctrine with:

- ↳ the boundary rules,
- ↳ the shared ID/ref conventions (itir://...),
- and ↳ the minimum event types SB must accept from TiRCorder and SL/TIRC.

Sources



Ahh it's a bit more separate than that... TIRC is probably closest to temporal/reasoning/understanding. SB is more like organisational/local/global context. SL is concerned with rules/law

Ah — okay, *thank you*, that clarification helps a lot.

You're right: I was still collapsing axes that you are **very deliberately keeping orthogonal**.

Let me restate the separation **your way**, and then build from *that* instead of fighting it.

The corrected mental model (three orthogonal axes)

What you're actually building is **not a stack** but a **coordinate system**.

1 SL — Normative / Rule axis

Concern: *What rules apply? What is binding? What constrains action?*

- Law, policy, principles, doctrine
- Normative force, not time
- Deterministic, span-anchored
- "What *must* / *may* / *must not* be true?"

SL does **not** care when something happened or how it feels in lived context.

2 TIRC — Temporal / Reasoning / Understanding axis

Concern: *What happened, in what order, and how do claims relate?*

- Evidence capture
- Reasoning traces
- Claim–evidence–counterclaim
- Trajectories of understanding over time

≡ ChatGPT 5.2 ▾



TIRC answers:

"What do we think is the case, based on what evidence, so far?"

This is where *understanding* lives.

3 SB — Organisational / Context axis

Concern: *Where are we, what's open, what's active, what matters locally vs globally?*

- State of work, not truth of claims

- Attention, priorities, blockers, carryover
- Local vs global context
- Operational framing

SB answers:

“What is the situation we are in?”

Not “what is true?”

Not “what is allowed?”

But “what’s the current shape of things?”

This is *situational*, not epistemic.

Why this separation matters (and why you’re right)

If SB drifts into TIRC:

- it starts reasoning
- it starts collapsing interpretations
- it becomes an assistant / narrator

If SB drifts into SL:

- it starts enforcing rules
- it becomes normative
- it loses adaptability

So SB **must** stay *contextual*, not interpretive and not normative.

That’s the line I was blurring before — thanks for stopping it.

What SB ingests from each (now precise)

From TIRC

SB ingests:

- *signals of epistemic activity*, not conclusions
 - “new claim introduced”
 - “evidence added”
 - “trajectory stalled”
 - “branch opened/closed”

- references to TIRC objects (`claim_id` , `trajectory_id`)

SB does **not** ingest:

- reasoning graphs
- truth assessments
- likelihoods

It only tracks *that reasoning is happening or blocked*.

From SL

SB ingests:

- *constraints in force*
 - "rule applies here"
 - "policy gate failed"
 - "review required"
- references to SL artifacts (`principle_id` , `rule_id`)

SB does **not**:

- interpret doctrine
- weigh competing rules
- decide outcomes

It only knows **constraints exist**.

What SB actually *does* with this (its unique value)

SB synthesises **context**, not meaning.

That context is built from three kinds of inputs:

1. **Activity** (from TIRC)
 - what is being worked on
 - where progress is happening or not
2. **Constraints** (from SL)
 - what cannot proceed
 - what requires compliance/review
3. **Environment**
 - local vs global scope
 - resources, energy, time, locality

SB then emits:

- “current state”
- “carryovers”
- “blockers”
- “what’s active / dormant”
- “what changed since last check”

All *without* asserting truth or correctness.

The key SB invariant (this is the lock)

SB may summarise *situation*, but never summarise *content*.

- It can say “*blocked by unresolved rule X*”
- It cannot say “*rule X is violated*”
- It can say “*claim under review*”
- It cannot say “*claim is false*”

That keeps the axes clean.

How our earlier LZ / compression discussions still apply

Now reframed properly:

- **SL** uses structural compression to canonicalise rules/text.
- **TIRC** uses structure to manage reasoning over time.
- **SB** uses *organisational compression* to stabilise context.

So SB benefits from the discussion not by:

- doing LZ over tokens

...but by borrowing the **principle**:

stable atoms + deterministic reduction = no narrative drift

For SB, the “atoms” are:

- task IDs
- context tags
- blockers
- scopes
- references to SL/TIRC artifacts

SB compresses *state*, not language.

A helpful metaphor (optional, but clarifying)

- **SL** = Constitution / Rulebook
- **TIRC** = Court transcript + evidence trail
- **SB** = Situation room whiteboard

The whiteboard never decides the case.

It just tells you what's on the table *right now*.

If you want to continue productively from here

Good next workshop steps would be:

1. ↪ **Enumerate SB's atom types** (state primitives, not text)
2. Define **what SB is allowed to collapse** vs what must stay explicit
3. Write a short **"SB does not do X" doctrine** (anti-features)
4. Define **handshake events** from TIRC → SB and SL → SB

If you want, we can do that next — slowly and precisely — now that the axes are clean.

