

A simple measure of the goodness of fit of a causal theory to a text corpus

📅 23 Dec 2025

Abstract

Suppose an evaluation team has a corpus of interviews and progress reports, plus (at least) two candidate theories of change (ToCs): an original one and a revised one. A practical question is: **which ToC better fits the narrative evidence?**

With almost-automated causal coding as described in (n.d.; Powell et al. 2025), we can turn that into a simple set of *coverage-style* diagnostics: how much of the coded causal evidence can be expressed in the vocabulary of each ToC.

See also: [Working Papers](#); [Minimalist coding for causal mapping](#); [Magnetisation](#).

Intended audience: evaluators and applied researchers comparing candidate ToCs (or other causal frameworks) against narrative evidence, who want a transparent “fit” diagnostic that does not pretend to be causal inference.

Unique contribution (what this paper adds):

- A definition of **coverage over causal links** (not just themes): link / citation / source coverage variants.
- A simple protocol for comparing candidate ToC vocabularies using hard recode or [Magnetisation](#) (soft recode).
- A careful positioning of “coverage” relative to mainstream QDA usage (saturation/counting as support for judgement, not a mechanical rule).

1. The core idea: “coverage” of evidence by a codebook

In ordinary QDA (thematic coding), researchers often look at how widely a codebook or set of themes is instantiated across a dataset: which codes appear, how frequently, and whether adding more data still yields new codes (saturation). Counting is not the whole of qualitative analysis, but it is a common, explicitly discussed support for judgement and transparency (Saldaña 2015). Critiques of turning saturation into a mechanical rule-of-thumb are also well known (Braun & Clarke 2019).

Our twist is: because we are coding **causal links** (not just themes), we can define coverage over *causal evidence* rather than over text volume.

2. Minimal definitions

- A **coded link** is a row of the form (`Source_ID`, `Quote`, `Cause_Label`, `Effect_Label`, ...).
- A **ToC codebook** is a vocabulary (list) of ToC factor labels you want to recognise in the corpus.
- A **mapping** from raw labels to ToC labels can be done either:
 - strictly (exact match / “hard recode”), or

- softly via magnetisation (semantic similarity; “soft recode”) — see [Magnetisation](#).

3. Coverage measures you can compute

Assume we have a baseline set of coded links L (from open coding), and a ToC codebook C (as magnets / targets).

3.1 Link coverage (our main measure)

Link coverage = proportion of coded links whose endpoints can be expressed in the ToC vocabulary.

Two variants (pick one and state it explicitly):

- **Both-ends coverage:** count a link as “covered” only if *both* cause and effect are mapped to some ToC label.
- **At-least-one-end coverage:** count a link as “covered” if either endpoint maps (useful when ToC vocabulary is intentionally partial).

3.2 Citation coverage (weighted link coverage)

If your dataset has multiple citations per bundle (or you have `Citation_Count`), compute coverage over **citations**, not just distinct links:

- covered citations / total citations

This answers: “what proportion of the *evidence volume* is expressible in this ToC?”

3.3 Source coverage (breadth)

Source coverage = number (or proportion) of sources for which at least k links are covered by the ToC vocabulary.

This answers: “does this ToC vocabulary work across many sources, or only a small subset?”

4. Protocol (how to use it)

For each candidate ToC:

1. Build a ToC codebook C (ideally keep candidate codebooks similar in size and specificity, otherwise you are partly measuring codebook granularity).
2. Map raw labels to C (hard recode or soft recode).
3. Compute:
 4. link coverage (both-ends and/or one-end),
 5. citation coverage (if available),
 6. source coverage (with an explicit k).
7. Inspect the **leftovers** (uncovered labels/links): what important evidence is the ToC not even able to name?

5. How this relates to “coverage” in mainstream qualitative methods

The word “coverage” is used in a few nearby ways in qualitative methodology:

- **Code (or theme) saturation:** whether new data still yields new codes/themes; the distinction between “code saturation” and “meaning saturation” is often emphasised (e.g. Hennink et al. on code vs meaning saturation; and the broader critique that saturation is not a universal stopping rule in all qualitative paradigms) (Braun & Clarke 2019).
- (For orientation, see: Hennink, Kaiser & Marconi (2017) “Code Saturation Versus Meaning Saturation”, *Qualitative Health Research*, DOI: [10.1177/1049732316665344](https://doi.org/10.1177/1049732316665344); Guest, Bunce & Johnson (2006) “How Many Interviews Are Enough?”, *Field Methods*, DOI: [10.1177/1525822X05279903](https://doi.org/10.1177/1525822X05279903).)
- **Counting for transparency:** many QDA approaches use counts (how often codes occur; how widely they occur across

cases) as a support for analytic claims, without equating frequency with importance (Saldaña 2015).

What we are doing here is closer to: **how much of the coded evidence can be expressed in the language of a candidate theory**, which is a “fit” diagnostic rather than a claim about truth.

6. Caveats

- Coverage is sensitive to **granularity**: broader ToC labels will (almost by definition) cover more.
- High coverage does not imply causal truth; it only implies that the ToC vocabulary is a good *naming scheme* for a large share of the corpus.
- Low coverage can mean either “ToC is missing key mechanisms” or “coding/mapping is too strict” — inspect leftovers before concluding.

References

- Braun, & Clarke (2019). *To Saturate or Not to Saturate? Questioning Data Saturation as a Useful Concept for Thematic Analysis and Sample-Size Rationales*. <https://doi.org/10.1080/2159676X.2019.1704846>.
- Powell, Cabral, & Mishan (2025). *A Workflow for Collecting and Understanding Stories at Scale, Supported by Artificial Intelligence*. SAGE PublicationsSage UK: London, England. <https://doi.org/10.1177/13563890251328640>.
- Saldaña (2015). *The Coding Manual for Qualitative Researchers*. Sage.