

Technical Appendix: Metrics, Data Structures, and Estimation Procedures

Matt Krasnow, Will Krasnow

August 19, 2025

Contents

1	Notation	2
2	Data Structures	3
2.1	Relational Record (Evaluation Artifact)	3
2.2	Evaluation JSON Payload	3
2.3	Optional Metadata Keys (if present)	4
3	Derived Variables and Defaults	5
4	Scoring and Re-aggregation	5
5	Many-Facet Rasch Measurement (MFRM)	5
5.1	Teacher Growth and Uncertainty	5
5.2	Improvement Index (Percentile Translation)	6
6	Reliability	6
6.1	Generalizability Theory (G-Study and D-Study)	6
6.2	Rasch/MFRM Person Separation	6
7	Fairness and Drift	6
8	Text-Based Validity, Alignment, and Tone	6
8.1	Evidence-to-Rating Alignment Index (ERAI)	6
8.2	Clarity of Expectations	7
8.3	Relational Tone Signals	7
8.4	Topic Alignment (School-Level)	7
9	Time, Throughput, and Pipeline	7
9.1	Authoring Time and Latency	7
9.2	Throughput, Coverage, Streak	7
9.3	Pipeline Efficiency	8
10	Coaching and Goals	8
11	AI-Assist Utilization and Effect (if logged)	8

12 Focus Groups and Surveys	8
12.1 Trust Index	8
12.2 Usability and Workload	8
13 Comparative Designs	8
13.1 Matched Rollout & Difference-in-Differences	8
13.2 Stepped-Wedge (Non-Random or Randomized)	9
14 Composite Indices as 0-Based Difference Scores	9
14.1 Baselines and Orientation	9
14.2 TPG — Teacher Professional Growth	9
14.3 PPE — Principal Performance and Efficiency	10
14.4 RT — Relational Trust	10
14.5 Uncertainty and Reporting	10
15 Data Quality & Missingness	11
A Field Dictionary: Definitions, Types, and Calculations	11
B Estimation Outputs (Summary)	12
C Implementation Notes	13

Abstract

This appendix formally defines the data structures, variables, and statistical procedures used to compute all reported metrics. It covers (i) the relational record and JSON evaluation payload, (ii) derived variables required by the estimation procedures, and (iii) the exact formulas for all sub-metrics and composites: Teacher Professional Growth (TPG), Principal Performance and Efficiency (PPE), and Relational Trust (RT). The document is written for academic advisors and statisticians and aims to be implementation-ready.

1 Notation

- Indices: teachers $i \in \{1, \dots, N_T\}$, evaluators/raters $j \in \{1, \dots, N_R\}$, components $k \in \mathcal{K}$, domains $d \in \mathcal{D}$, artifacts (evaluations) $a \in \{1, \dots, N_A\}$, time periods t .
- Scores: component score $X_{ijka} \in \mathbb{R}$ (ordinal/polytomous on a bounded scale); domain score S_{ida} ; recomputed overall score S_{ia}^{overall} .
- Rasch/MFRM parameters: teacher ability θ_i , rater severity ρ_j , step thresholds β_{kx} (category x for component k), optional period effect τ_t .
- Timestamps: observation/creation T_a^{obs} , draft start T_a^{draft} , finalize T_a^{final} , delivery/share T_a^{deliv} .
- Text: component summaries Y_{ka} , overall summary Y_a^{overall} , low-inference notes N_a .
- Operators: $\text{Var}(\cdot)$, $\text{SE}(\cdot)$, $\Phi(\cdot)$ standard normal CDF.

2 Data Structures

2.1 Relational Record (Evaluation Artifact)

Each evaluation artifact is a row with the following fields (logical names shown; actual column names may differ):

Field	Type	Definition / Use
id	bigint	Unique artifact identifier.
created_at	timestampz	Creation time; default proxy for observation time T_a^{obs} .
updated_at	timestampz	Last update time; used in latency derivations.
shared_at	timestampz (nullable)	Delivery/share time; preferred for T_a^{deliv} if set.
evaluation_id	uuid	Logical evaluation/observation UID.
evaluator	uuid	Evaluator/rater identifier j .
teacher_id	uuid (nullable)	Teacher identifier i .
teacher_name	text (nullable)	Teacher label (for human-readable reports).
school_id	uuid (nullable)	School identifier.
school_name	text (nullable)	School label.
organization_id	uuid (nullable)	Organization/district identifier.
framework_id	text	Instructional framework key.
evaluation	json	JSON payload (Section 2.2).
ai_evaluation	jsonb (nullable)	Optional AI generation/telemetry.
metadata	jsonb (nullable)	Optional workflow metadata (timestamps, flags).
low_inference_notes	text (nullable)	Raw evidence notes N_a .
additional_comments	jsonb (nullable)	Auxiliary comments (e.g., goals).
is_shared	boolean	Whether artifact was shared.
deleted_at	timestampz (nullable)	Soft delete timestamp (exclude if set).

2.2 Evaluation JSON Payload

The JSON conforms to the following schema:

```
Evaluation {
  domains: { [domain_id: string]: DomainEvaluation },
  metadata: EvaluationMetadata,
  summary?: string,
  summaryScores: SummaryScores
}

DomainEvaluation {
  name: string,
  components: { [component_id: string]: ComponentEvaluation },
  weight: number,
```

```

    isManuallyScored: boolean,
    summary?: string,
    domainScore: number
}

ComponentEvaluation {
  score: number,
  summary: string,
  error?: string,
  isManuallyScored: boolean,
  modified?: boolean,
  insufficientEvidence?: boolean,
  teacherSummary?: string
}

SummaryScores {
  overallScore: number,
  domainWeights: { [domain_id: string]: number }
}

EvaluationMetadata {
  framework_id: string,
  framework_name: string
}

```

2.3 Optional Metadata Keys (if present)

- `metadata.observed_at` (timestamp): preferred T_a^{obs} .
- `metadata.draft_started_at`, `metadata.draft_finalized_at` (timestamp): authoring interval.
- `metadata.delivered_to_teacher_at` (timestamp): delivery time.
- `metadata.is_ai_assisted` (bool): AI assist flag.
- `metadata.source_note_ids[]` (array): linked evidence anchors.
- `metadata.prior_goal_ids[]` (array): referenced coaching goals.
- `metadata.token_count`, `metadata.word_count` (int).
- `metadata.artifact_type` (text): e.g., *observation*, *coaching*.

3 Derived Variables and Defaults

$$T_a^{\text{obs}} := \begin{cases} \text{metadata.observed_at} & \text{if present} \\ \text{created_at} & \text{otherwise} \end{cases}$$

$$T_a^{\text{deliv}} := \begin{cases} \text{metadata.delivered_to_teacher_at} & \text{if present} \\ \text{shared_at} & \text{else if present} \\ \text{updated_at} & \text{fallback} \end{cases}$$

Word and token counts combine all component summaries and the optional overall summary:

$$\text{Words}_a := \text{wc}(Y_a^{\text{overall}} \parallel \{Y_{ka}\}_k), \quad \text{Tokens}_a := \text{tc}(Y_a^{\text{overall}} \parallel \{Y_{ka}\}_k).$$

4 Scoring and Re-aggregation

Let w_d be the domain weight; S_{ida} the domain score for teacher i in artifact a . The recomputed overall score is

$$S_{ia}^{\text{overall}} = \sum_{d \in \mathcal{D}} w_d S_{ida}, \quad \text{where} \quad w_d = \begin{cases} \text{evaluation.summaryScores.domainWeights}[d] & \text{if provided} \\ \text{evaluation.domains}[d].\text{weight} & \text{fallback} \end{cases} \quad (1)$$

5 Many-Facet Rasch Measurement (MFRM)

We use a partial-credit multi-facet model to estimate teacher ability (θ_i), rater severity (ρ_j), and step thresholds (β_{kx}). For component k with ordered categories $x = 0, \dots, m_k - 1$, define

$$\text{logit } \mathbb{P}(X_{ijka} \geq x) = \theta_i - \rho_j - \beta_{kx} - \tau_{t(a)}, \quad (2)$$

where τ_t is an optional period effect for time $t(a)$ determined by T_a^{obs} . Estimation proceeds via joint/conditional maximum likelihood or marginal ML with appropriate identifiability constraints (e.g., sum-to-zero over facets).

5.1 Teacher Growth and Uncertainty

For teacher i , define baseline $T0$ and follow-up $T1$ windows. With person estimates $\hat{\theta}_i(T0)$, $\hat{\theta}_i(T1)$ and standard errors $\text{SE}[\hat{\theta}_i(T0)]$, $\text{SE}[\hat{\theta}_i(T1)]$:

$$g_i = \hat{\theta}_i(T1) - \hat{\theta}_i(T0), \quad (3)$$

$$\text{SE}(g_i) = \sqrt{\text{SE}^2[\hat{\theta}_i(T1)] + \text{SE}^2[\hat{\theta}_i(T0)]}, \quad (4)$$

$$\text{CI}_{95\%}(g_i) = g_i \pm 1.96 \text{SE}(g_i). \quad (5)$$

Group growth aggregates via inverse-variance weighting:

$$\bar{g} = \frac{\sum_i w_i g_i}{\sum_i w_i}, \quad w_i := \text{SE}(g_i)^{-2}, \quad \text{SE}(\bar{g}) = \left(\sum_i w_i \right)^{-1/2}. \quad (6)$$

5.2 Improvement Index (Percentile Translation)

Let SD_{pre} be the standard deviation of $\hat{\theta}_i$ in the baseline window. Define standardized effect $d := g/SD_{\text{pre}}$ and percentile translation

$$U3 = 100 \times \Phi(d). \quad (7)$$

6 Reliability

6.1 Generalizability Theory (G-Study and D-Study)

With persons (p), raters (r), and occasions (o), random-effects ANOVA yields variance components $\sigma_p^2, \sigma_{pr}^2, \sigma_{po}^2, \sigma_{pro}^2$. For planned numbers of raters n_r and occasions n_o , the relative G-coefficient is

$$G = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{pr}^2/n_r + \sigma_{po}^2/n_o + \sigma_{pro}^2/(n_r n_o)}. \quad (8)$$

The D-study solves for (n_r, n_o) to achieve target G (e.g., $G \geq 0.80$).

6.2 Rasch/MFRM Person Separation

From the MFRM, report person separation reliability (PSI) for θ estimates and the number of distinguishable strata:

$$\text{Strata} = \frac{4 \times \text{Separation} + 1}{3}. \quad (9)$$

7 Fairness and Drift

- **Rater severity spread:** $SD(\{\hat{\rho}_j\}_j)$. Flag unusually wide spreads.
- **Extremes:** $\max_j |\hat{\rho}_j|$.
- **Drift:** segment by month/quarter via T_a^{obs} and examine $\hat{\rho}_j(t)$; apply control charts or re-fit time-sliced models.
- **Differential Rater Functioning (DRF):** test interactions of raters with subgroups (e.g., school/subject if available in metadata).

8 Text-Based Validity, Alignment, and Tone

8.1 Evidence-to-Rating Alignment Index (ERAI)

For each component k :

$$\text{cit_presence}_{ka} \in \{0, 1\} \quad (\text{regex over summaries/notes}) \quad (10)$$

$$\text{rubric_sim}_{ka} \in [0, 1] \quad (\text{embedding cosine to rubric text}) \quad (11)$$

$$\text{entail}_{ka} \in [0, 1] \quad (\text{NLI: notes} \Rightarrow \text{summary claims}) \quad (12)$$

$$\text{consist}_{ka} \in [0, 1] \quad (\text{text-inferred level vs. numeric score}) \quad (13)$$

Component-level alignment:

$$\text{ERAI}_{ka} = \frac{1}{4}(\text{cit_presence}_{ka} + \text{rubric_sim}_{ka} + \text{entail}_{ka} + \text{consist}_{ka}). \quad (14)$$

Artifact-level alignment averages valid components:

$$\text{ERAI}_a = \frac{1}{|\mathcal{K}_a|} \sum_{k \in \mathcal{K}_a} \text{ERAI}_{ka}, \quad \mathcal{K}_a := \{k : \text{not insufficientEvidence}\}. \quad (15)$$

8.2 Clarity of Expectations

Within artifact text (overall + components), compute

$$\text{clarity}_a = \mathbf{1}\{\text{time-bound}\} + \mathbf{1}\{\text{measurable}\} + \mathbf{1}\{\text{rubric-linked}\} \in \{0, 1, 2, 3\}.$$

8.3 Relational Tone Signals

On concatenated text:

$$\text{affirm_share}_a = \frac{\#\text{affirming tokens}}{\text{Tokens}_a}, \quad (16)$$

$$\text{praise:suggest}_a = \frac{\#\text{praise sentences}}{\#\text{suggestion sentences}}, \quad (17)$$

$$\text{we/you}_a = \frac{\#\text{"we"}}{\#\text{"you"}}, \quad (18)$$

$$\text{hedge_rate}_a = 1000 \times \frac{\#\text{hedges}}{\text{Tokens}_a}, \quad (19)$$

$$\text{follow_through}_a \in \{0, 1\} \quad (\text{references to prior goals}). \quad (20)$$

8.4 Topic Alignment (School-Level)

Let $\pi^{(\text{FG})}$ be the topic distribution from teacher focus groups and $\pi^{(\text{FB})}$ from principal feedback text. Define

$$\text{topic_align} = \frac{\pi^{(\text{FG})} \cdot \pi^{(\text{FB})}}{\|\pi^{(\text{FG})}\|_2 \|\pi^{(\text{FB})}\|_2} \in [0, 1]. \quad (21)$$

9 Time, Throughput, and Pipeline

9.1 Authoring Time and Latency

$$\text{authoring_minutes}_a = \begin{cases} (T_a^{\text{final}} - T_a^{\text{draft}})/60 & \text{if both present} \\ \text{N/A} & \text{otherwise} \end{cases} \quad (22)$$

$$\text{latency_days}_a = (T_a^{\text{deliv}} - T_a^{\text{obs}})/86400. \quad (23)$$

9.2 Throughput, Coverage, Streak

Let month $m(a) := \text{date_trunc}(\text{month}, T_a^{\text{deliv}})$ (or T_a^{obs} if delivery not used).

$$\text{throughput}_{jm} = \#\{a : \text{evaluator } j, m(a) = m\}, \quad (24)$$

$$\text{coverage}_{sm} = \frac{\#\{\text{distinct } i \text{ at school } s \text{ with } \geq 1 \text{ artifact in } m\}}{\text{teacher roster size at } s \text{ in } m}, \quad (25)$$

$$\text{streak}_i = \max \text{ run length of consecutive months with } \geq 1 \text{ artifact for } i. \quad (26)$$

9.3 Pipeline Efficiency

$$\text{converted_7d_rate} = \frac{\#\{a : T_a^{\text{deliv}} - T_a^{\text{obs}} \leq 7 \text{ days}\}}{\#\{a\}}, \quad (27)$$

$$\text{edits_to_final}_a \approx \sum_{k \in \mathcal{K}_a} \mathbf{1}\{\text{modified} = \text{true}\} \quad (\text{proxy if version logs absent}). \quad (28)$$

10 Coaching and Goals

$$\text{coaching_per_teacher_term} = \frac{\#\{\text{coaching artifacts for teacher in term}\}}{\text{term length}}, \quad (29)$$

$$\text{goals_set_per_teacher} = \#\{\text{new goal identifiers in term}\}, \quad (30)$$

$$\text{SMART_score} \in \{0, \dots, 5\} \quad (\text{presence of S/M/A/R/T attributes via rules/classifier}). \quad (31)$$

11 AI-Assist Utilization and Effect (if logged)

$$\text{utilization_rate} = \mathbb{E}[\mathbf{1}\{\text{is_ai_assisted}\}], \quad (32)$$

$$\Delta \text{time} = \text{median}(\text{authoring_minutes} \mid \text{assisted}) - \text{median}(\text{authoring_minutes} \mid \text{matched non-assisted}), \quad (33)$$

where matching controls for *text length* (words/tokens) and *component coverage* within evaluator j . Report Hodges–Lehmann estimate and Wilcoxon p -value.

12 Focus Groups and Surveys

12.1 Trust Index

Code focus-group transcripts for *respect*, *competence*, *integrity*, and *personal regard*. With intensity $c \in \{0, 1, 2\}$ for each coded occurrence and C_{\max} the maximum possible per participant:

$$\text{RT}_{\text{FG}} = \frac{\sum \text{code intensities}}{C_{\max} \times \#\text{participants}}. \quad (34)$$

Report teacher and principal indices separately and their alignment gap.

12.2 Usability and Workload

Compute standard SUS (0–100) on the 10-item instrument and optional NASA-TLX workload composites for the authoring task.

13 Comparative Designs

13.1 Matched Rollout & Difference-in-Differences

For outcome Y_{st} (e.g., school-level mean $\hat{\theta}$ or latency), let D_{st} indicate adoption. Estimate

$$Y_{st} = \alpha_s + \gamma_t + \delta \cdot D_{st} + \mathbf{X}_{st}^\top \boldsymbol{\beta} + \varepsilon_{st}, \quad (35)$$

with school fixed effects α_s , period fixed effects γ_t , robust SEs clustered by school. For staggered adoption, use group-time average treatment effects (e.g., Callaway–Sant’Anna estimators) and event-study pre-trend checks.

13.2 Stepped-Wedge (Non-Random or Randomized)

Fit a mixed model with period fixed effects and cluster (school) random intercepts; report ITT effects, intra-class correlation (ICC), and pertinent sensitivity analyses.

14 Composite Indices as 0-Based Difference Scores

Composites are defined as weighted sums of *natural-unit* differences from a pre-specified baseline, centered at 0 so that positive values indicate improvement and negatives indicate decline. No *z*-score standardization or percentage rescaling is used.

14.1 Baselines and Orientation

Let b index the baseline window (e.g., pre-period) or control condition, and t the follow-up window for a given reporting unit (teacher, evaluator, or school as appropriate). For each sub-metric M we compute

$$\Delta M := \bar{M}_t - \bar{M}_b, \quad (36)$$

with direction aligned so that higher is better. For “lower-is-better” metrics (e.g., authoring minutes, latency), use $\Delta M := -(\bar{M}_t - \bar{M}_b)$.

Unit scales. To place heterogeneous natural units on a comparable contribution scale while preserving interpretability, each sub-metric has a fixed, pre-specified *unit scale* u_M (not estimated from the data). The scaled difference is $\Delta M/u_M$, where one unit corresponds to a meaningful change (e.g., $u_{\text{authoring}} = 10$ minutes, $u_{\text{latency}} = 1$ day, $u_{\text{throughput}} = 1$ artifact/month, $u_{\text{ERAI}} = 0.05$ absolute, $u_g = 0.10$ logits). Chosen u_M values should be pre-registered and held constant across reports.

14.2 TPG — Teacher Professional Growth

Primary construct is growth in rater-adjusted latent performance, supplemented by validity and reliability guardrails. Define

$$\Delta g := \bar{g}_t - \bar{g}_b \quad (\text{logits; already oriented}), \quad (37)$$

$$\Delta \text{ERAI} := \overline{\text{ERAI}}_t - \overline{\text{ERAI}}_b, \quad (38)$$

$$\Delta \text{Reliability} := \overline{\text{PSI}}_t - \overline{\text{PSI}}_b, \quad (39)$$

$$\Delta \text{Fairness} := -\left(\text{SD}(\{\hat{\rho}_j\})_t - \text{SD}(\{\hat{\rho}_j\})_b\right), \quad (40)$$

$$\Delta \text{FG growth} := \overline{\text{FG growth index}}_t - \overline{\text{FG growth index}}_b. \quad (41)$$

With weights $w^{(\text{TPG})}$ (summing to 1) and unit scales u_M , the composite is

$$\text{TPG} = \sum_{M \in \{g, \text{ERAI}, \text{Reliability}, \text{Fairness}, \text{FG growth}\}} w_M^{(\text{TPG})} \cdot \frac{\Delta M}{u_M}. \quad (42)$$

Recommended default emphasizes growth: $w_g^{(\text{TPG})} = 0.6$, others share the remainder unless pre-registered otherwise.

14.3 PPE — Principal Performance and Efficiency

Define oriented differences for key efficiency and production metrics:

$$\Delta\text{authoring} := -(\overline{\text{authoring_minutes}_t} - \overline{\text{authoring_minutes}_b}), \quad (43)$$

$$\Delta\text{latency} := -(\overline{\text{latency_days}_t} - \overline{\text{latency_days}_b}), \quad (44)$$

$$\Delta\text{throughput} := \overline{\text{artifacts/leader/month}_t} - \overline{\text{artifacts/leader/month}_b}, \quad (45)$$

$$\Delta\text{coverage} := \overline{\% \text{ teachers with } \geq 1 \text{ artifact/month}_t} - \overline{\% \text{ teachers with } \geq 1 \text{ artifact/month}_b}, \quad (46)$$

$$\Delta\text{content} := \overline{\text{content richness index}_t} - \overline{\text{content richness index}_b}, \quad (47)$$

$$\Delta\text{pipeline} := \overline{\text{converted_7d_rate}_t} - \overline{\text{converted_7d_rate}_b}. \quad (48)$$

Composite:

$$\text{PPE} = \sum_{M \in \{\text{authoring, latency, throughput, coverage, content, pipeline}\}} w_M^{(\text{PPE})} \cdot \frac{\Delta M}{u_M}, \quad \sum_M w_M^{(\text{PPE})} = 1. \quad (49)$$

Recommended default emphasizes time saved and latency: $w_{\text{authoring}}^{(\text{PPE})} = 0.35$, $w_{\text{latency}}^{(\text{PPE})} = 0.25$, remaining weight distributed over production/quality metrics.

14.4 RT — Relational Trust

Define oriented differences:

$$\Delta\text{RT}_{\text{FG}} := \overline{\text{RT}_{\text{FG}}_t} - \overline{\text{RT}_{\text{FG}}_b}, \quad (50)$$

$$\Delta\text{Tone} := \overline{\text{Relational Tone Index}_t} - \overline{\text{Relational Tone Index}_b}, \quad (51)$$

$$\Delta\text{Clarity} := \overline{\text{clarity}_t} - \overline{\text{clarity}_b}, \quad (52)$$

$$\Delta\text{TopicAlign} := \overline{\text{topic_align}_t} - \overline{\text{topic_align}_b}. \quad (53)$$

Composite:

$$\text{RT} = \sum_{M \in \{\text{RT}_{\text{FG}}, \text{Tone}, \text{Clarity}, \text{TopicAlign}\}} w_M^{(\text{RT})} \cdot \frac{\Delta M}{u_M}, \quad \sum_M w_M^{(\text{RT})} = 1. \quad (54)$$

Recommended default emphasizes focus-group trust: $w_{\text{RT}_{\text{FG}}}^{(\text{RT})} = 0.5$.

14.5 Uncertainty and Reporting

Report point estimates of each ΔM in natural units alongside the composite. Compute composite confidence intervals via nonparametric bootstrap over teachers (TPG, RT) or evaluators (PPE), maintaining the weighting and unit scales within each replicate.

15 Data Quality & Missingness

- Exclude artifacts with `deleted_at` not null.
- Prefer T_a^{obs} from metadata; otherwise use `created_at`.
- If delivery is unlogged, use `shared_at` then `updated_at` for T_a^{deliv} .
- Report coverage (%) for each metric and apply pairwise deletion by default; consider multiple imputation for time-series analyses if MAR is plausible.

A Field Dictionary: Definitions, Types, and Calculations

Identity & Linkage

Name	Type	Definition / Calculation
<code>id</code>	bigint	Artifact identifier.
<code>evaluation_id</code>	uuid	Observation UID.
<code>teacher_id</code> , <code>teacher_name</code>	uuid, text	Teacher key/label (i).
<code>evaluator</code>	uuid	Evaluator/rater key (j).
<code>school_id</code> , <code>school_name</code>	uuid, text	School key/label.
<code>organization_id</code>	uuid	Organization/district key.
<code>framework_id</code>	text	Instructional framework key.

Timestamps

Name	Type	Definition / Calculation
<code>created_at</code>	timestampz	Default T_a^{obs} if metadata absent.
<code>updated_at</code>	timestampz	Used in latency fallback.
<code>shared_at</code>	timestampz	Preferred T_a^{deliv} if set.
<code>metadata.observed_at</code>	timestampz	Preferred observation time.
<code>metadata.draft_started_at</code>	timestampz	Authoring start (if present).
<code>metadata.draft_finalized_at</code>	timestampz	Authoring end (if present).
<code>metadata.delivered_to_teacher_at</code>	timestampz	Delivery time (if present).

Scores & Text

Name	Type	Definition / Calculation
<code>domains[d].domainScore</code>	number	Domain score S_{ida} .
<code>domains[d].weight</code>	number	Domain weight w_d (fallback).
<code>summaryScores.domainWeights[d]</code>	number	Preferred w_d .
<code>components[k].score</code>	number	Component score X_{ijka} .
<code>components[k].summary</code>	string	Text Y_{ka} .
<code>summary</code>	string	Overall text Y_a^{overall} .

Name	Type	Definition / Calculation
summaryScores.overallScore	number	Provided overall score (audited vs. recompute).
low_inference_notes	string	Evidence text N_a .

Flags & Links

Name	Type	Definition / Calculation
components[k].isManuallyScored	bool	Manual scoring flag.
components[k].modified	bool	Edited post-generation.
components[k].insufficientEvidence	bool	Evidence insufficient.
metadata.is_ai_assisted	bool	AI assist used.
metadata.source_note_ids[]	array	Evidence anchors.
metadata.prior_goal_ids[]	array	Goal references.

Derived Quantities

Name	Definition / Formula
S_{ia}^{overall}	$\sum_d w_d S_{ida}$
Words, Tokens	Counts over $Y_a^{\text{overall}} \parallel \{Y_{ka}\}$
Evidence density	$(\text{\#citations})/\text{Words} \times 100$
Authoring minutes	$(T_a^{\text{final}} - T_a^{\text{draft}})/60$ (if present)
Latency days	$(T_a^{\text{deliv}} - T_a^{\text{obs}})/86400$
Throughput	$\text{\#artifacts per evaluator-month}$
Coverage	$\text{\%teachers with } \geq 1 \text{ artifact per month}$
Streak	$\text{Longest consecutive months with } \geq 1 \text{ artifact}$

B Estimation Outputs (Summary)

Quantity	Definition / Estimation
$\hat{\theta}_i$	Teacher ability (logits) from MFRM.
$\hat{\rho}_j$	Rater severity (logits) from MFRM.
$\hat{\beta}_{kx}$	Step thresholds from MFRM.
g_i, \bar{g}	Teacher growth and group mean; CIs via delta method.
$U3$	Percentile translation $100\Phi(g/SD_{\text{pre}})$.
G	Generalizability coefficient; D-study recommendations.
PSI, Strata	Person separation and distinguishable strata.
ERAI	Evidence-to-rating alignment index (0–1).

Quantity	Definition / Estimation
Tone/Clarity	Affirm share, praise:suggest, we/you, hedging, clarity (0–3).
Topic alignment	Cosine similarity of topic distributions (0–1).

C Implementation Notes

- All time deltas are computed in seconds and rescaled (min, d).
- When delivery time is not explicitly recorded, use `shared_at`, else `updated_at`.
- For matched comparisons (AI effect), match within evaluator on word count and component coverage.
- Bootstrap CIs: 2000 replicates by teacher for TPG/RT and by evaluator for PPE.
- Exclude artifacts with `deleted_at` not null from all analyses.