UNCERTAINTY-ROUTED HUMAN-LLM CURATION AND CALIBRATION FOR ANLI

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

031

033

034

037

040

041

042

043

044

045

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Adversarial NLI (ANLI) reveals distribution-shift failures that static benchmarks miss, motivating evaluation and curation that are explicitly uncertainty-aware. We present URC²—Uncertainty-Routed Curation & Calibration, a threestage pipeline that improves dataset quality and model reliability. URC² decomposes per-example predictive uncertainty into aleatoric (data/label ambiguity) and epistemic (model uncertainty, measured via mutual information) using a three-teacher ensemble (DeBERTa-v3-large, RoBERTa-large, XLM-R-large). A two-lane relabeling workflow then routes cases: a *Human lane* relabels, removes, or down-weights aleatoric-heavy examples, while an LLM lane adjudicates epistemic-heavy examples using instruction-tuned self-consistency checks. Curated labels and per-example weights drive a lightweight retraining and recalibration loop for each teacher, yielding an updated ensemble. On ANLI, URC² reduces development-set expected calibration error by 30% (to 0.146) and lowers corpus-level uncertainty without degrading accuracy. Unlike prior curation pipelines that treat uncertainty monolithically, URC² exploits the distinction between aleatoric and epistemic uncertainty to enable uncertainty-aware control, identifying recurring failure modes and mitigating them via targeted reweighting and data augmentation. URC² provides a practical, reproducible recipe for building more trustworthy NLI systems under adversarial shift.

1 Introduction

Accuracy on static, i.i.d. benchmarks often masks brittleness under shift: models can be confidently wrong on adversarial inputs. The Adversarial NLI (ANLI) benchmark Nie et al. (2020) was designed precisely to surface these failures, yet most pipelines still treat uncertainty as a single scalar and optimize post-hoc calibration in isolation. This leaves two practical gaps: (i) distinguishing *aleatoric* ambiguity (data/label issues) from *epistemic* disagreement (model uncertainty), and (ii) acting on each with the right supervision rather than uniformly reweighting or purely post-hoc scaling.

We propose URC² (<u>Uncertainty-Routed Curation & Calibration</u>), a three-stage approach that (A) measures per-example uncertainty using a teacher ensemble and *decomposes* it into aleatoric and epistemic components, (B) *routes* examples by dominant source of uncertainty to a two-lane relabeling workflow (humans for ambiguity; an instruction-tuned LLM for confident model disagreement), and (C) *refreshes* the teachers with curated labels and instance weights, followed by lightweight temperature scaling. In short, URC² turns diagnosis into targeted supervision. Motivational example appears in Appx AFig. 6)

Concretely, we instantiate the teacher ensemble with DeBERTa-v3-large, RoBERTa-large, and XLM-R-large and compute, per example, total uncertainty H, aleatoric A (mean per-model entropy), and epistemic E via a mutual-information style decomposition:

$$H = \mathbf{H}(\bar{p}), \quad A = \mathbf{E}_m[\mathsf{H}(p_m)], \quad E = H - A.$$

High-A items are human-audited (relabel / keep-hard with down-weight / drop), while high-E items—where models are individually confident yet disagree—are adjudicated by an instruction-tuned LLM with self-consistency checks. We then retrain and calibrate. Operational thresholds, prompts, selection heuristics, and optimizer details are deferred to the appendix.

Definitions. We report *expected calibration error (ECE)*: the weighted average of per-bin |acc - conf| over B confidence bins Guo et al. (2017). We apply *temperature scaling (TS)* Guo et al. (2017) post-hoc by optimizing a single temperature T on the development split and using $\hat{p} = \operatorname{softmax}(z/T)$ at test time. Here, "dev" denotes the concatenated ANLI development set (R1+R2+R3; 3,200 examples).

Results on ANLI. URC² improves probability quality without sacrificing accuracy: the ensemble's dev ECE drops from 0.249 to 0.209 pre-TS, and to 0.146 after temperature scaling, with accuracy held (§5). On the curated subset, epistemic mutual information collapses (large negative Δ MI), indicating real disagreement removal rather than cosmetic smoothing. Corpus-level uncertainty shifts toward the low-aleatoric/low-epistemic quadrant after the teacher refresh. Full diagnostics, risk-coverage curves, and per-round breakdowns appear in the appendix.

Contributions.

- 1. *Uncertainty-driven supervision*. We operationalize an ensemble-based decomposition (aleatoric vs. epistemic) and use it to *route* examples to different supervision mechanisms.
- 2. *Human–LLM two-lane relabeling*. Ambiguity-heavy items go to human audit; confident model disagreements go to an instruction-tuned LLM with acceptance checks, yielding curated labels and per-example weights.
- 3. Calibration with disagreement reduction. On ANLI, URC² achieves strong dev calibration (ECE **0.146** with TS) at stable accuracy and consistent with targeted reduction of model-side disagreement, with favorable corpus-level uncertainty shifts.

Terminology. We refer to the end-to-end method as the URC² pipeline (three stages). Within Stage 2, the human–LLM relabeling workflow has two lanes: a *Human* lane for aleatoric-heavy cases and an *LLM* lane for epistemic-heavy cases. See figure 1 for more details. We use "pipeline" only for the full method and "workflow/lanes" only for Stage 2.

2 BACKGROUND AND PROBLEM STATEMENT

2.1 OVERVIEW OF ANLI

The ANLI dataset Nie et al. (2020) stress-tests NLI systems via a *human* + *model-in-the-loop* protocol: crowd workers craft hypotheses that fool a strong model yet remain clear to other annotators. ANLI was released in three adversarial rounds (R1–R3). Each round used progressively stronger models as "adversaries," requiring annotators to construct increasingly subtle hypotheses to break them. This yields challenging premise–hypothesis pairs targeting shortcut heuristics and blind spots. ANLI uses three standard labels: entailment, neutral, contradiction. Overall labels are roughly balanced with a slight skew toward neutral, which simplifies accuracy comparisons but can mask label-quality and disagreement issues. Premises are typically multi-sentence passages from Wikipedia or the web; hypotheses are short, adversarially phrased statements. Dataset statistics appear in (Appx. A; see Table 9 and Fig. 5).

2.2 INPUT LENGTHS AND STRUCTURE

Premises are long (multi-sentence); hypotheses are short and cue-sensitive (negations, quantifiers, numerals). This asymmetry raises two issues: (i) long premises may be truncated under short sequence limits, removing needed evidence; (ii) short hypotheses concentrate weight on few tokens (e.g., "not," "three," "most"), encouraging superficial cues and high-variance predictions. Illustrative examples appear in the appendix (Table 16).

2.3 PROBLEM STATEMENT

Adversarial NLI benchmarks such as ANLI expose failure modes that traditional accuracy metrics cannot resolve: models often make confident but incorrect predictions, undermining trust in their outputs. A central challenge is that such failures stem from two distinct kinds of uncertainty. *Ambiguity* (aleatoric) arises when hypotheses allow multiple valid interpretations or contain annotation

noise, making uncertainty unavoidable even for an ideal model. *Disagreement* (epistemic) arises when inputs are linguistically clear but adversarially crafted, leading strong models to diverge with high confidence and exposing blind spots that targeted supervision can resolve. The problem that we consider, therefore, is to disentangle these two forms of uncertainty and respond to them appropriately. Our objectives in this work are threefold: (i) to quantify aleatoric and epistemic uncertainty at the level of individual examples, (ii) to act on these signals with targeted supervision—directing ambiguous cases to human audit and disagreement-heavy cases to LLM adjudication, and (iii) to demonstrate that this routing improves probability calibration, enhances label quality, and yields a healthier corpus-level distribution of uncertainty.

3 RELATED WORK

Adversarial NLI (ANLI) introduced a human + model-in-the-loop protocol that systematically surfaces shortcut reliance and blind spots that i.i.d. test sets can miss Nie et al. (2020), aligning with broader efforts in dynamic evaluation such as Dynabench Kiela et al. (2021). Earlier analyses revealed annotation artifacts enabling hypothesis-only baselines and brittle heuristics (e.g., HANS) Gururangan et al. (2018); Poliak et al. (2018); McCoy et al. (2019), while ChaosNLI Naeini et al. (2015) emphasized that annotator disagreement often reflects genuine ambiguity rather than mere noise. These strands largely *diagnose* dataset pathologies; our focus is to *act* on them by routing examples to different forms of supervision based on their uncertainty type.

Estimating predictive uncertainty with deep ensembles and approximate Bayesian methods (e.g., MC Dropout) is practical and effective Lakshminarayanan et al. (2017); Gal & Ghahramani (2016). Following the standard view Kendall & Gal (2017), we decompose total uncertainty into *aleatoric* (data/label ambiguity) and *epistemic* (model uncertainty) components; in ensembles, the latter can be captured via a BALD-style mutual information term computed from model disagreement. Prior work predominantly uses such decomposition as a diagnostic; we instead *use* it to drive targeted curation: ambiguity-heavy items go to human audit, whereas confidently disagreeing items go to an instruction-tuned LLM for adjudication.

Calibration methods such as temperature scaling, Dirichlet, and beta calibration improve probability quality without altering accuracy Guo et al. (2017); Kull et al. (2019; 2017), yet calibration typically degrades under distribution shift Ovadia et al. (2019), and reported ECE on ANLI remains high even for strong models Hu et al. (2023). Rather than relying solely on post-hoc scaling, we first improve supervision quality through uncertainty-routed edits and instance weighting, and only then apply lightweight TS—yielding substantially lower ECE on ANLI.

A complementary literature addresses label quality and noisy-label learning Northcutt et al. (2021; 2019); Song et al. (2022). Consistent with these insights, our human lane corrects clear errors, drops irreparably noisy cases, and down-weights "keep-hard" items that remain ambiguous. In parallel, instruction-tuned LLMs have emerged as capable annotators and judges Ouyang et al. (2022); Gilardi et al. (2023); Zheng et al. (2023); our use is deliberately *selective*—restricted to high-epistemic items with confident model disagreement, with acceptance checks to avoid blanket relabeling. Finally, selective/conformal prediction offers deployment-time guarantees via abstention or set prediction Angelopoulos & Bates (2023); our training-time curation improves calibration and reduces disagreement, thereby strengthening the conditions under which such guarantees are most reliable.

Positioning vs. closest approaches. Active learning and noisy-label pipelines typically prioritize examples by model uncertainty or loss but do not distinguish *aleatoric* from *epistemic* causes, often mixing ambiguous items with clear-but-divisive ones Song et al. (2022); Northcutt et al. (2019). Methods like ADDMU detect adversarial or boundary cases via uncertainty signals but stop at diagnosis rather than targeted remediation Yin et al. (2022). In contrast, URC² uses an ensemble-based decomposition to *route* items to distinct interventions (human vs. LLM), then folds decisions back through instance weighting and a teacher refresh before light post-hoc calibration, producing measurable reductions in epistemic disagreement on ANLI.

Scope and complementarity. Our focus is training-time supervision quality, complementary to deployment-time reliability tools such as selective and conformal prediction Angelopoulos & Bates (2023). By reducing disagreement and improving calibration upstream, we strengthen the conditions under which abstention or set-prediction guarantees are most informative. Practically, ANLI's

adversarial rounds (R1–R3) and evolving "adversary" strength make it a natural testbed for such iterative curation: early rounds contain more ambiguity, later rounds emphasize clear-but-tricky inputs. URC² exploits this structure by assigning ambiguous cases to humans and confidently divisive ones to an instruction-tuned LLM with acceptance checks, avoiding blanket relabeling while targeting the sources of error.

4 METHODOLOGY

4.1 Preliminaries and Notation

Task and ensemble. We study NLI with inputs x=(p,h) (premise, hypothesis) and label $y\in\{1,\ldots,C\}$; for ANLI, C=3 with labels entailment, neutral, contradiction. To estimate uncertainty, we use an ensemble of M fine-tuned teachers $\{\theta_m\}_{m=1}^M$, each producing posteriors $P_{\theta_m}(y\,|\,x)$. Uncertainty decomposition. Let the ensemble mean be

$$\bar{p}(y \mid x) = \frac{1}{M} \sum_{m=1}^{M} P_{\theta_m}(y \mid x).$$
 (1)

We define total (H), aleatoric (A), and epistemic (E) uncertainty via entropies:

$$H(x) = -\sum_{y=1}^{C} \bar{p}(y \mid x) \log \bar{p}(y \mid x),$$
 (2)

$$A(x) = \frac{1}{M} \sum_{m=1}^{M} \left[-\sum_{y=1}^{C} P_{\theta_m}(y \mid x) \log P_{\theta_m}(y \mid x) \right], \tag{3}$$

$$E(x) = H(x) - A(x). \tag{4}$$

Here H captures ensemble predictive uncertainty; A reflects irreducible data/label ambiguity that persists across teachers; and the gap E quantifies reducible model disagreement (a mutual-information-style quantity). We also use the *epistemic share* r(x) = E(x)/H(x) (when H(x) > 0).

Confidence and calibration. We use normalized entropy $\widehat{H}(x) = H(x)/\log C \in [0,1]$ and define entropy-based confidence $c(x) = 1 - \widehat{H}(x)$. Calibration is measured via ECE with B = 10 equal-width bins:

$$ECE = \sum_{b=1}^{B} \frac{|S_b|}{N} \left| acc(S_b) - conf(S_b) \right|.$$
 (5)

where S_b are the examples grouped into B=10 equal-width bins based on their confidence scores. ECE quantifies the weighted average mismatch between confidence and accuracy across bins—lower is better, with ECE=0 indicating perfect calibration. Post-hoc temperature scaling (TS) Guo et al. (2017) rescales logits z with a scalar T.

$$P_T(y \mid x) = \frac{\exp(z_y/T)}{\sum_{y'} \exp(z_{y'}/T)}.$$
 (6)

where T>0 is a learned scalar (fit by minimizing negative log-likelihood on a held-out set). For T>1, predictions become softer (less confident); for T<1, sharper. This preserves the argmax label but better aligns probabilities with true likelihoods. We instantiate $M{=}3$ diverse teachers—DeBERTa-v3-large, RoBERTa-large, XLM-R-large—to promote genuine epistemic variation (different inductive biases and pretraining). For each x we compute \bar{p} , H, A, $E{=}H{-}A$, and $r{=}E/H$ from the teacher posteriors.

4.2 THE HUMAN-LLM RELABELING PIPELINE

URC² converts uncertainty diagnosis into targeted supervision through three coupled stages: (A) *measure & route* by dominant uncertainty type; (B) *human–LLM relabeling* in two disjoint lanes; (C) *retrain & calibrate* and refresh the ensemble. We iterate $A \rightarrow B \rightarrow C$ until dev calibration plateaus. See figure 1 for detailed work flow.

4.2.1 STAGE A — MEASURE & ROUTE (NO EDITS)

From the fixed teacher ensemble we compute (H, A, E, r) and assign each example to *exactly one* lane by dominant uncertainty. An example is routed to the **epistemic-heavy** pool if *all* hold:

- 1. $r(x) = E(x)/H(x) \ge 0.35$ (disagreement dominates total uncertainty),
- 2. at least two teachers predict different argmax labels,
- 3. for the disagreeing teachers, the top–second probability margin ≥ 0.20 (each is confident).

Otherwise, items with high A(x) and low r(x) are treated as **aleatoric-heavy** (ambiguous/noisy). Previously curated items (by either lane) are skipped in later rounds. Stage A *only* measures and routes; no labels or weights are changed here.

4.2.2 STAGE B — TWO-LANE RELABELING (TARGETED SUPERVISION)

Each routed item is curated in *one disjoint lane*; items never cross lanes. Outputs are a possibly updated label and an instance weight w_i .

Lane H (**Human audit; aleatoric-heavy**). High-A, low-r items typically reflect vague premises, underspecification, or label noise. Annotators choose among:

- Relabel (apply corrected label; weight w=1.0),
- Keep-Hard (retain original label but mark as ambiguous; $w=\alpha$, with $\alpha=0.3$),
- *Drop* (remove from training; w=0.0).

This preserves informative border cases without letting ambiguity dominate, while removing irredeemably ill-posed items and applying clear corrections at full weight.

Lane L (LLM adjudication; epistemic-heavy). When teachers are confident yet disagree (high r, low A), the example is sent to an instruction-tuned LLM to make a single, consistent call. The LLM sees a few curated demonstrations plus the query and returns a structured label; we use a *strict application* rule—only apply the LLM's label if it differs from the original—then give that item full weight (w=1.0) and record its index in \mathcal{S}_{chg} for later Δ MI analysis. This keeps the intervention targeted: no change, no update; clear change, full adoption. Implementation specifics (prompt format, demo retrieval, output schema, acceptance checks) are summarized in Appx. C (Tabs. 10, 11, 12).

4.2.3 STAGE C — RETRAIN, RECALIBRATE, REFRESH

Curated labels and weights from both lanes are folded back into all teachers to reduce disagreement and improve probability quality. Inputs. (i) Curated labels $\{y_i^{\text{clean}}\}$ and instance weights $\{w_i\}$ from Stage B; (ii) same splits/tokenization; (iii) the original fine-tuning recipe. See Appx C.1, Table 13 for training details.

Procedure.

- 1. **Retrain teachers.** Optimize weighted cross-entropy $\mathcal{L} = \sum_i w_i \ell(f_{\theta}(x_i), y_i^{\text{clean}})$ with early stopping on dev NLL.
- 2. **Rebuild ensemble.** Rescore the training set with refreshed teachers in fixed member order.
- 3. **Recompute uncertainty.** Update \bar{p} , H, A, E, and r from the refreshed ensemble.
- 4. **Post-hoc calibration.** Fit temperature T on dev by NLL and report accuracy + ECE with/without TS (argmax unchanged).

Diagnostics, stopping, and optimization. After each refresh we evaluate three aspects of progress in a consistent way. On the training set we expect calibration to improve, which we operationalize as a reduction in ECE. On the development set we evaluate accuracy together with post–temperature-scaling calibration; the argmax labels are preserved by TS, so any reduction in dev ECE at fixed

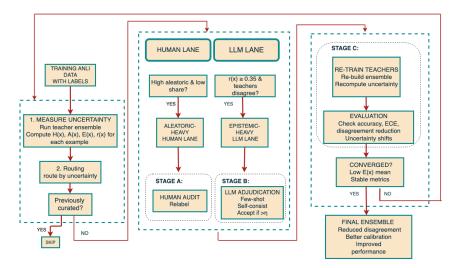


Figure 1: **URC**² **overview.** Stage A routes by uncertainty; Stage B applies a two-lane Human–LLM relabeling; Stage C retrains and calibrates, closing the loop.

accuracy indicates strictly better probability quality. Beyond these point metrics, we also examine how the uncertainty landscape evolves by reporting the corpus means of total (H), aleatoric (A), and epistemic (E) uncertainty, the expected epistemic share $\mathbf{E}[E/H \mid H>0]$, and the mass of examples in each quadrant defined by fixed high/low cutoffs on A and E. We continue iterating the $A \rightarrow B \rightarrow C$ loop while dev ECE keeps decreasing and the share of high-E examples contracts; once both trends plateau, we stop. For training details see Appx. C.2, Table 12.

5 EVALUATION PROTOCOL AND METRICS

We compute confidence two ways: (i) entropy-based $c_{\rm ent}(x)=1-\frac{H(x)}{\log C}$ and (ii) max-probability $c_{\rm max}(x)=\max_y \bar{p}(y\mid x)$. We report ECE under both (Ent-ECE and MaxP-ECE), along with negative log-likelihood (NLL) and Brier score. Following standard practice, we decompose predictive uncertainty into aleatoric (data-driven) and epistemic (model-driven) components, with epistemic estimated via mutual information (MI) capturing prediction variability. All metrics use the same train/dev splits, tokenization, and ensemble setup as in prior sections. Follow Sec. 4.

Visual analyses. We include: (i) per-round summaries of epistemic (MI) and aleatoric uncertainty with 95% confidence intervals (Fig. 5 in Appx. A), showing higher epistemic in later rounds (R3) and higher aleatoric in earlier ones (R1); and (ii) scatter plots of epistemic vs. aleatoric uncertainty on random subsamples (Fig. 5 in Appx. A), revealing distinct regimes (aleatoric-dominated, epistemic-dominated, mixed).

5.1 CALIBRATION AND ACCURACY METRICS

We use normalized entropy $\widehat{H}(x) = H(x)/\log C \in [0,1]$ and define confidence $c(x) = 1 - \widehat{H}(x)$. ECE computed over 10 equal-width confidence bins. We report train ECE (pre/post curation) and dev ECE after post-hoc temperature scaling (TS), which minimizes dev NLL but leaves accuracy (argmax labels) unchanged. Reliability diagrams (Fig. 2, 3) visualize the bin-wise gaps.

5.2 Calibration baselines (no curation)

We benchmark standard post-hoc calibrators on the uncurated ensemble (mean of teacher posteriors) using Ent-ECE, MaxP-ECE, NLL, and Brier scores. As shown in Table 1, temperature scaling (TS) and histogram binning improve MaxP-ECE and NLL over the raw ensemble, but Ent-ECE remains high. In contrast, URC²—which modifies supervision (Stage B) then retrains and cali-

Table 1: Post-hoc calibration baselines on ANLI dev (no curation). Ensemble uses mean of teacher posteriors.

Method	Acc	Ent-ECE	MaxP-ECE	NLL	Brier
Uncalibrated	0.514	0.202	0.169	0.996	0.620
Temperature Scaling	0.514	0.251	0.122	0.951	0.595
Histogram Binning (OVR)	0.515	0.252	0.118	0.952	0.595

Table 2: Comparison of best post-hoc method (OVR, no curation) vs. URC² with temperature scaling on ANLI dev.

Method	Acc	Ent-ECE	MaxP-ECE	NLL
Best post-hoc (OVR)	0.515	0.252	0.118	0.952
URC^{2} (+TS)	0.581	0.146	0.120	0.884

brates (Stage C)—achieves substantially lower Ent-ECE (0.146 vs. 0.251/0.252 for TS/histogram on uncurated data), while also improving MaxP-ECE and proper scores (Tables 2 and 14). This gap shows that correcting labels in disagreed-upon examples and retraining teachers improves probability quality beyond post-hoc methods alone. Post-hoc calibration helps, but URC² (refresh + TS) is consistently better, as it fixes supervision before calibrating rather than just rescaling uncurated outputs.

To assess calibration improvements, we analyze corpus-level changes pre/post curation: (i) **Means**: Track $\mathbf{E}[H]$, $\mathbf{E}[A]$, $\mathbf{E}[E]$, and epistemic share $\mathbf{E}[r(x) \mid H(x) > 0]$. Expect reduced $\mathbf{E}[E]$ from less disagreement and slightly lower $\mathbf{E}[A]$ from fewer noisy items. (ii) **Quadrant shares**: Using fixed high/low A and E cutoffs, measure mass in high-E/low-E (disagreement), high-E/low-E (ambiguity), both-high, and both-low. Success shifts mass to both-low, reducing high-E. We iterate until: (i) train ECE drops with stable accuracy; (ii) dev ECE improves post-TS with competitive accuracy; and (iii) $\mathbf{E}[E]$ and high-E mass decrease. Stop when calibration and disagreement metrics stabilize.

6 RESULTS

We evaluate URC^2 on ANLI with two goals: (i) **probability quality** on dev (accuracy, ECE) and (ii) **disagreement resolution** on train (epistemic MI and the overall uncertainty landscape). We first report ensemble calibration on dev, then quantify what changed due to curation (footprint and Δ MI), and finally show corpus-level shifts after the all-teacher refresh.

After the refresh, the dev reliability curve moves toward the diagonal in the mid-confidence region, reducing pre-TS ECE from 0.249 to 0.209. A light ECE-tuned temperature $T_{\rm ECE}^{\star}=0.53$ brings ECE to 0.146 while accuracy remains 0.581 (Table 3; Figs. 2, 3).

6.1 CURATION FOOTPRINT AND RESOLVED DISAGREEMENT

The two-lane correction makes a focused intervention: a small human audit (500 items; 99 relabeled, 11 dropped, the remainder down-weighted) complements LLM adjudication with 6,152 accepted edits. On the curated set \mathcal{S}_{chg} (indices with accepted label edits), ensemble epistemic MI falls from 0.4663 to 0.0739 ($\Delta MI = -0.3924$), indicating that disagreement is largely eliminated where supervision actually changed. In Fig. 4, the decile panel places post-curation MI in a narrow band even for the hardest deciles, while the ΔMI histogram concentrates its mass at \leq 0; together, these patterns are consistent with adjudication removing model-side disagreement on edited rows—an effect that post-hoc temperature scaling cannot produce. As a sanity check, a single-teacher pilot reduced train ECE from 0.235 to 0.105 at fixed accuracy (dev ECE 0.208 \rightarrow 0.177 with TS), underscoring that gains stem from cleaner supervision rather than accuracy gaming; per-teacher breakdowns appear in App. D, Table 14.

Table 3: Dev calibration of the teacher ensemble.

Model	Acc (pre-TS)	ECE (pre-TS)	$T_{\rm NLL}^{\star}$	ECE@ T_{NLL}^{\star}
Ensemble (DeBERTa-v3, RoBERTa, XLM-R)	0.581	0.209	0.53	0.146

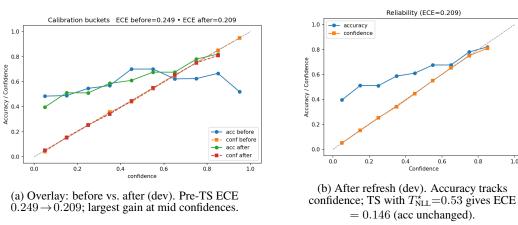


Figure 2: Reliability of the teacher ensemble on ANLI dev. Post-refresh accuracy aligns with confidence in mid-high bins, shrinking per-bin gaps that dominate ECE.

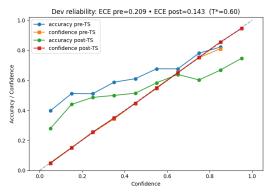


Figure 3: Reliability on ANLI dev before/after temperature scaling. X-axis: 10 equal-width confidence bins; Y-axis: accuracy and mean confidence. TS with $T_{\rm NLL}^{\star}=0.53$ (Table 3) reduces ECE $0.209 \rightarrow 0.146$ with accuracy unchanged; the main effect is lowering confidence in the highest bins.

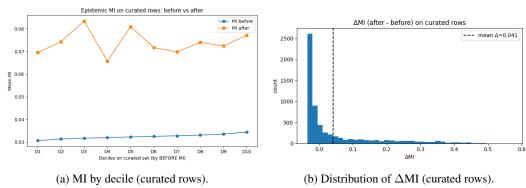


Figure 4: Resolved disagreement on curated rows. (a) Deciles formed by *pre-adjudication* MI; after adjudication, MI lies in a narrow band across all deciles (< 0.085). (b) Per-item Δ MI (after—before) concentrates at ≤ 0 with a long positive tail; the small right-shifted mean is driven by few outliers.

Table 4: Curation footprint and MI reduction on curated items.

Item	Value	Note
Human audit (aleatoric-heavy)	500	Relabeled 99 (19.8%), dropped 11 (2.2%);
		others retained/down-weighted.
LLM adjudication (accepted changes)	$6,\!152$	Applied at $w = 1.0$.
MI on curated set \mathcal{S}_{chg}	0.4663 o 0.0739	$\Delta MI = -0.3924$ (after-before).

6.2 Δ MI on curated rows and selective prediction (brief)

On the rows actually edited during curation (\mathcal{S} ; tensors standardized to (N,3) with RoBERTa class-order alignment), the change in ensemble epistemic MI is large and localized: $|\mathcal{S}|=6,152$, mean Δ MI (after-before) =-0.247, median =-0.256. This establishes that adjudication resolves model-side disagreement where labels changed (full bars/hist in Appx. E, Fig. 8).

Complementing this, selective prediction improves across the coverage range after the refresh: E-AURC decreases from 0.399 to 0.321 with entropy confidence and from 0.390 to 0.318 with max-probability confidence, indicating better ranking of easy vs. hard examples. Because temperature scaling is monotone, these gains are attributable to the refresh rather than post-hoc calibration (see Appx. D, Fig. 7 for the full risk–coverage panels).

Table 5: Uncertainty landscape on full train after refresh: quadrant shares (true 2×2) and corpus means.

Quadrant shares		Corpus means		
Quadrant	Share	Metric	Value	
Low-A / Low-E	0.700	Mean H	0.393	
High-A only	0.050	Mean A	0.358	
High-E only	0.050	Mean E	0.035	
Both-high	0.200	Share E/H	0.050	

7 ABLATION STUDIES

We isolate three factors: a temperature-scaling (TS) baseline on dev with *no* curation, the causal change in epistemic disagreement on the actually edited rows S, and the post-refresh uncertainty landscape.

7.1 TS-ONLY (NO CURATION) ON DEV

With labels untouched (pre-refresh dev), the teachers and their mean ensemble provide the calibration baseline in Table 6. Applying NLL-tuned TS makes a *pure* likelihood adjustment: log-loss improves from 0.874 to 0.846 at $T^* = 1.50$, while ECE moves upward (Ent-ECE $0.340 \rightarrow 0.453$; MaxP-ECE $0.040 \rightarrow 0.096$). This is an *expected objective trade-off*—we optimize NLL rather than ECE—so TS-only serves as a principled control, and the substantive gains in calibration and disagreement come from *curation+refresh*, not post-hoc scaling (Table 7).

Table 6: TS-only ablation on ANLI dev — (a) Per-model baselines (pre-TS; no curation).

Model	Acc	Ent-ECE	MaxP-ECE	NLL	Brier
DeBERTa-v3	0.719	0.177	0.220	1.256	0.497
RoBERTa-large	0.414	0.231	0.393	2.102	0.926
XLM-R-large	0.562	0.266	0.344	1.867	0.760

7.2 Δ MI on curated rows S (full train)

We measure the *causal* change in ensemble epistemic mutual information (MI) on the rows that actually changed during curation (S; RoBERTa class order aligned; all tensors standardized to (N, 3)).

Table 7: TS-only ablation on ANLI dev — (b) Ensemble before vs. after TS (NLL-tuned, $T^*=1.50$).

Metric	Before	After (TS, $T^* = 1.50$)	Δ
Ent-ECE	0.340	0.453	0.113
MaxP-ECE	0.040	0.096	0.056
NLL	0.874	0.846	-0.028
Brier	0.452	0.471	0.019

Table 8: TS-only ablation on ANLI dev — (c) Uncertainty quadrants after refresh (true 2×2).

Quadrant shares		Corpus means	
Quadrant	Share	Metric	Value
Low-A / Low-E	0.400	Mean H	1.020
High-A only	0.300	Mean A	0.824
High-E only	0.300	Mean E	0.196
Both-high	0.000 (< 0.001)	Share E/H	0.192

High/low via 70th percentiles: $\tau_A \approx 0.825$, $\tau_E \approx 0.232$.

The effect is large and localized: $|\mathcal{S}|$ (changed rows) = 6,152; Mean Δ MI (after – before) = **-0.247**; Median Δ MI = **-0.256**. The bar plot and histogram are shown in Appx. E, Fig. 8.

7.3 Uncertainty Landscape after refresh (true 2×2)

Post-refresh, most training examples concentrate in low or single-source uncertainty regimes. Using balanced percentile cutoffs (70th/70th) for high aleatoric (A) and high epistemic (E), the corpus splits as: 0.400 Low-A/Low-E, 0.300 High-A only, 0.300 High-E only, and essentially 0.000 jointly high (Table 17). This mix is consistent with disagreement being resolved rather than masked. Sensitivity to absolute thresholds and additional diagnostics are deferred to the appendix (Appx. G, Table 17).

8 DISCUSSION

Refreshing teachers on curated supervision improves reliability without changing decisions. On ANLI dev, pre-TS ECE drops from 0.249 to 0.209; an ECE-tuned temperature brings it to 0.146 with accuracy fixed at 0.581 (risk–coverage panels in Appx. D.2). Where supervision actually changed, disagreement is resolved: on \mathcal{S}_{chg} (6,152 edits), ensemble epistemic MI falls from 0.4663 to 0.0739 (Δ MI = -0.3924; full bars/hist in Appx. E). Corpus-wide after the refresh, uncertainty mass concentrates in low-uncertainty regions (both-low 0.578) with a small epistemic share (\sim 5%), consistent with the curated-slice gains; absolute 2×2 thresholds and sensitivity appear in Appx. G. The effect comes from splitting uncertainty (aleatoric vs. epistemic), routing to the right supervision (human vs. LLM), and performing a single refresh; per-teacher dev diagnostics are in Appx. D.

9 Conclusion

We presented an uncertainty-routed curation and calibration pipeline that *explicitly distinguishes* aleatoric ambiguity from epistemic disagreement and addresses each with the right tool—human audit for ambiguity and instruction-tuned LLM adjudication for disagreement—followed by a full-teacher refresh and temperature scaling. On ANLI, this yields (i) dev ECE improvements from $0.249 \rightarrow 0.209$ (pre-TS) and to 0.146 with TS at unchanged accuracy 0.581 (Fig. 2, Table 3); (ii) large, targeted disagreement removal on curated rows (Δ MI=-0.3924 over 6,152 changed items; Table 4); and (iii) a corpus-level shift toward low-uncertainty regions with a small epistemic share (Table 5). Beyond improving a single metric, URC² turns uncertainty *diagnosis* into *actionable supervision*, delivering calibrated probabilities and causally reducing model-side disagreement at scale—a general, reproducible recipe for reliable, LLM-in-the-loop dataset improvement.

REFERENCES

- Anastasios N. Angelopoulos and Stephen Bates. Conformal prediction: A unified review of theory and new challenges. *arXiv preprint arXiv:2301.06605*, 2023. URL https://arxiv.org/abs/2301.06605.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR. URL https://proceedings.mlr.press/v48/gall6.html.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. ChatGPT outperforms crowd workers for text annotation tasks. *arXiv preprint arXiv:2303.15056*, 2023. doi: 10.48550/arXiv.2303.15056. URL https://arxiv.org/abs/2303.15056.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1321–1330. PMLR, 06–11 Aug 2017. URL https://proceedings.mlr.press/v70/guo17a.html.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 107–112, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2017. URL https://aclanthology.org/N18-2017/.
- Neil Houlsby, Ferenc Huszar, Zoubin Ghahramani, and Máté Lengyel. Bayesian Active Learning for Classification and Preference Learning. In *Proceedings of the 25th International Conference on Neural Information Processing Systems (NeurIPS 2011)*, pp. 1613–1621. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper/2011/hash/4e8229d28fddf0dc42a676f08b6f3f0f-Abstract.html.
- Mengting Hu, Zhen Zhang, Shiwan Zhao, Minlie Huang, and Bingzhe Wu. Uncertainty in natural language processing: Sources, quantification, and applications. *arXiv* preprint arXiv:2306.04459, 2023. URL https://arxiv.org/abs/2306.04459.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, pp. 5574–5584. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper/2017/hash/2650d6089a6d640c5e85b2b88265dc2b-Abstract.html.
- Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Vinodkumar Prabhakaran, Jason Chang, Daniel Simig, Pontus Stenetorp, Adina Williams, Eunice Jo, Candace Ross, Robin Jia, Y-Lan Boureau, and Jason Weston. Dynabench: Rethinking benchmarking in NLP. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4110–4124. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021. naacl-main.324. URL https://aclanthology.org/2021.naacl-main.324/.
- Meelis Kull, Tameem Adel Filho, and Peter Flach. Beta calibration: A well-calibrated multi-class probability estimator. In *Proceedings of the 2017 Conference on Artificial Intelligence and Statistics (AISTATS 2017)*, volume 54 of *Proceedings of Machine Learning Research*, pp. 623–631, 2017. URL http://proceedings.mlr.press/v54/kull17a.html.
- Meelis Kull, Miquel Perello-Nieto, Michael Kängsepp, Tameem Adel Silva Filho, and Peter Flach. Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with dirichlet calibration. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pp. 12316–12326. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper/2019/hash/01ca0bfd860d9e3e5a82a64c9f5355a4-Abstract.html.

- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, pp. 6402–6413. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper/2017/hash/9ef2ed4b7fd2c810847ffa5fa85bce38-Abstract.html.
- Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, pp. 7047–7058. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper/2018/hash/3ea2db50e62ceefceaf70a9d9a56a6f4-Abstract.html.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3428–3448, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1334. URL https://aclanthology.org/P19-1334/.
- Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI-15)*, pp. 2901–2907, Austin, Texas, USA, 2015. AAAI Press. doi: 10.1609/aaai.v29i1.9602. URL https://ojs.aaai.org/index.php/AAAI/article/view/9602.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial NLI: A new benchmark for natural language understanding. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4885–4901, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.441. URL https://aclanthology.org/2020.acl-main.441/.
- Curtis G. Northcutt, Tongzhou Wu, and Isaac Chuang. Confident learning: Estimating uncertainty in dataset labels. In *Proceedings of the 36th International Conference on Machine Learning (ICML 2019)*, pp. 6400–6409. PMLR, 2019. URL http://proceedings.mlr.press/v97/northcutt19a.html.
- Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021. doi: 10.1613/jair. 1.12125. URL https://www.jair.org/index.php/jair/article/view/12125.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, and Katarina Slama. Training language models to follow instructions with human feedback. *arXiv preprint*, 2022. URL https://arxiv.org/abs/2203.02155. arXiv:2203.02155.
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pp. 13969–13980, Vancouver, Canada, 2019. Curran Associates, Inc. URL https://proceedings.neurips.cc/paper/2019/hash/8558cb408c1d76621371888657d2eb1d-Abstract.html.
- Maja Pavlovic and Massimo Poesio. The effectiveness of LLMs as annotators: A comparative overview and empirical analysis of direct representation. *Proceedings of the 3rd Workshop on Perspectivist Approaches to NLP (NLPerspectives)* @ *LREC-COLING 2024*, pp. 100–110, 2024. URL https://aclanthology.org/2024.nlperspectives-1.11.
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. Hypothesis only baselines in natural language inference. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics* (*SEM 2018), pp. 180–191, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/S18-2023. URL https://aclanthology.org/S18-2023/.

- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Alexander Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*, 2022. URL https://arxiv.org/abs/2110.08207.
- Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z. Ren, and Anirudha Ray. A survey on uncertainty quantification of large language models: Taxonomy, open research challenges, and future directions. arXiv preprint arXiv:2412.05563, 2024. URL https://arxiv.org/abs/2412.05563.
- Lewis Smith and Yarin Gal. Understanding measures of uncertainty for adversarial example detection. In *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence (UAI 2018)*, pp. 560–569. AUAI Press, 2018. URL https://arxiv.org/abs/1803.08533.
- Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 34(11):8135–8153, 2022. URL https://arxiv.org/abs/2007.08199.
- Vladimir Vovk, Alex Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World*. Springer, 2005. ISBN 978-0-387-23338-9. doi: 10.1007/978-0-387-34003-2.
- Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. Do-not-answer: A dataset for evaluating safeguards in llms. *arXiv preprint arXiv:2308.13387*, 2023. URL https://arxiv.org/abs/2308.13387.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. *arXiv preprint*, 2021. URL https://arxiv.org/abs/2109.01652. arXiv:2109.01652.
- Fan Yin, Yao Li, Cho-Jui Hsieh, and Kai-Wei Chang. ADDMU: Detection of far-boundary adversarial examples with data and model uncertainty estimation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 6567–6584, 2022. URL https://aclanthology.org/2022.emnlp-main.440.
- Chujie Zheng, Bill Yuchen Lin, Haopeng Zhang, Mutian He, Chen Li, Qinyuan Zeng, Xiang Ren, and Junxian He. LLM-as-a-Judge: Guiding prompt selection with language models. *arXiv* preprint arXiv:2306.05685, 2023. URL https://arxiv.org/abs/2306.05685.

ADDITIONAL BACKGROUND ON ANLI

Purpose. This section supports Sec. 2 with per-round uncertainty trends and dataset statistics.

ANLI was released in three adversarial rounds (R1-R3). Each round used progressively stronger models as "adversaries," requiring annotators to construct increasingly subtle hypotheses to break them. This results in systematic differences:

- R1: Smaller, noisier, with many underspecified or ambiguous hypotheses.
- **R2**: More balanced, but still contains artifacts and label disagreement.
- R3: By far the largest, adversarially sharper, with higher epistemic uncertainty (models disagree while annotators converge).

What this shows. Tab. 9 summarizes splits and collection protocol for quick reference.

Table 9: **ANLI dataset statistics** by adversarial round (R1–R3).

	•
Aspect	Details
Train size (R1 / R2 / R3)	16,946 / 45,460 / 100,459
Dev size (R1 / R2 / R3)	1,000 / 1,000 / 1,200
Test size (R1 / R2 / R3)	1,000 / 1,000 / 1,200
Label set	entailment, neutral, contradiction
Collection method	Adversarial, human-in-the-loop (R1-R3)

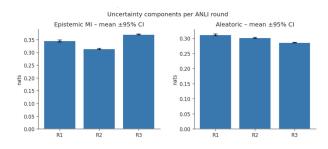


Figure 5: Uncertainty by round. Mean (±95% CI) epistemic MI (left) and aleatoric (right) across R1-R3.

В EXTENDED RELATED WORK

B.1 ANLI CASE STUDY

What this shows. Fig. 6 demonstrates the effect of URC2: routing by uncertainty sends the epistemic item(example A) to LLM adjudication, collapsing disagreement and concentrating mass on CONTRADICTION $(0.07, 0.34, 0.59) \rightarrow (0.05, 0.08, 0.87)$; the aleatoric item(Example B) goes to human audit and remains deliberately cautious while leaning NEUTRAL, shifting $(0.36, 0.37, 0.27) \rightarrow (0.30, 0.41, 0.29)$ without forcing artificial agreement.

B.2 ADVERSARIAL NLI AND DYNAMIC EVALUATION

ANLI's human + model-in-the-loop procedure Nie et al. (2020) iteratively elicits examples that fool strong models while staying clear to humans, complementing dynamic benchmarking such as Dynabench Kiela et al. (2021). Prior analyses exposed annotation artifacts enabling hypothesis-only baselines Gururangan et al. (2018); Poliak et al. (2018); McCoy et al. (2019). ChaosNLI Nie et al. (2020) showed annotator disagreement often reflects genuine ambiguity, motivating ambiguity-aware treatment. We build directly on this line by turning diagnosis into intervention via uncertainty-routed curation.

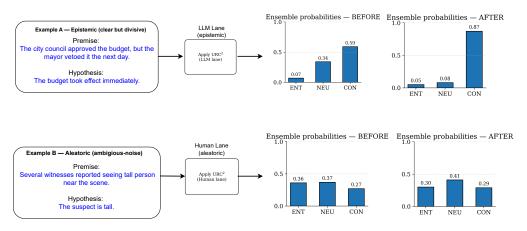


Figure 6: Two ANLI examples with ensemble probabilities before/after applying URC².

B.3 PREDICTIVE UNCERTAINTY: ESTIMATORS AND DECOMPOSITION

Approximate Bayesian methods and deep ensembles provide practical uncertainty in modern NLP. MC Dropout Gal & Ghahramani (2016) and deep ensembles Lakshminarayanan et al. (2017) are widely used. Following Kendall & Gal (2017), we split uncertainty into *aleatoric* (data/label ambiguity) and *epistemic* (model) components. For ensembles, epistemic uncertainty can be quantified by mutual information (BALD) Houlsby et al. (2011); Smith & Gal (2018). Alternatives include Dirichlet Prior Networks Malinin & Gales (2018), evidential models, and temperature/Dirichlet/beta calibrations. Surveys Hu et al. (2023); Shorinwa et al. (2024) catalog this space and highlight decomposition's utility under shift. Unlike methods such as ADDMU Yin et al. (2022) that primarily *detect* far-boundary points, we *route* examples by uncertainty type into distinct supervision lanes and retrain.

B.4 Calibration Under Distribution Shift

Post-hoc calibration (temperature scaling Guo et al. (2017), Dirichlet Kull et al. (2019), beta Kull et al. (2017)) aligns confidence with empirical accuracy without changing argmax labels. Calibration often degrades under shift Ovadia et al. (2019); large ECE values have been reported on ANLI even for strong models Hu et al. (2023). We therefore improve supervision first via routing+curation, then apply light TS; empirically this yields lower ECE than post-hoc scaling alone.

B.5 LABEL QUALITY, NOISY LABELS, AND INSTANCE WEIGHTING

Label issues are pervasive in large datasets Northcutt et al. (2021). Methods include confident learning Northcutt et al. (2019), sample selection, and instance weighting Song et al. (2022). Noisy-label surveys Song et al. (2022) and NLP benchmarks Wang et al. (2023) argue for treating ambiguity differently from clear errors. Our human lane mirrors this: correct clear mistakes, drop irrecoverable items, and down-weight "keep-hard" ambiguous cases.

B.6 LLMs as Annotators and Judges

Instruction-tuned LLMs (e.g., InstructGPT/FLAN/T0) show strong agreement with human labels Ouyang et al. (2022); Wei et al. (2021); Sanh et al. (2022); Gilardi et al. (2023) and serve as automatic judges Zheng et al. (2023). We invoke an instruction-tuned LLM *only* on high-epistemic items where teacher models are confident yet disagree, with acceptance checks before folding decisions back into training—contrasting with blanket relabeling Pavlovic & Poesio (2024).

B.7 SELECTIVE/CONFORMAL PREDICTION AND RISK-COVERAGE

Selective prediction and conformal prediction Vovk et al. (2005); Angelopoulos & Bates (2023) provide coverage guarantees by abstaining or outputting sets. They are deployment-time safeguards

Table 10: Lane L: LLM adjudication configuration (model, prompting, parsing, acceptance, selection).

Component	Setting	
Model	Meta-Llama-3-8B-Instruct (local); 8-bit quantization (BitsAndBytes); FP16	
Prompt	3 anchor exemplars (ENT/NEU/CON) $+$ 6 retrieved demos (see Tab. 12); strict JSON output	
Preprocessing	ASCII cleaning; truncate: premise 180 chars, hypothesis 120 chars	
Sampling	do_sample=True, T =0.7, top_p= 0.9, max_new_tokens= 256, N =1	
Parsing	Regex for {"label": "ENT"} / "NEU" / "CON"; invalid ⇒ no change	
Acceptance	Apply only if predicted label \neq original; if applied, set instance weight $w{=}1.0$ and log to \mathcal{S}_{chg}	
Selection	Hard pool: top 6,000 by $s(x)=H(\bar{p}(\cdot x))-\left(p_{(1)}(x)-p_{(2)}(x)\right)$ (confident disagreements)	
Artifacts	<pre>train_labels_llama_str.pkl, train_labels_llama_bin.pkl, llama_patched_idx.npy</pre>	

rather than supervision improvers. We report risk-coverage and E-AURC to quantify ranking quality before/after refresh; improvements mirror our calibration gains and complement conformal methods.

B.8 Positioning and Summary

Across strands—adversarial benchmarking, uncertainty estimation, calibration, noisy-label learning, and LLM annotation—most prior work diagnoses failure or calibrates post hoc in isolation. Our contribution is to *integrate* them: use ensemble-based decomposition to identify *why* an example is uncertain (aleatoric vs. epistemic), *route* it appropriately (human vs. LLM), and *refresh* teachers before minimal post-hoc scaling. This yields lower ECE and measurable disagreement reduction on the edited slice, while improving the corpus-level uncertainty mix.

C LANE-L CONFIGURATION (PROMPTING, RETRIEVAL, ACCEPTANCE)

Scope. This section specifies the Lane L (LLM adjudication) setup end-to-end: model/prompting, demo retrieval, selection/acceptance, and the training recipe used to refresh teachers after curation. All components here correspond to the same scenario and are used together in our experiments.

What this shows. Tab. 10 lays out the exact knobs used for LLM adjudication so results are reproducible.

C.1 DEMO RETRIEVAL DETAILS

Why this matters. Tabs. 11-12 document the retrieval setup that conditions the adjudicator with in-domain demos.

Table 11: Retrieval encoder specs for demo selection (all-mpnet-base-v2).

Property	Detail
Model	sentence-transformers/all-mpnet-base-v2
Embedding dim	768
Max input length	\sim 512 tokens
Pooling	Mean pooling; L2-normalized
Similarity	Cosine similarity
Use in pipeline	Encode premise [SEP] hypothesis for queries & pool

Table 12: Demo retrieval configuration used by Lane L prompting.

Component	Setting
Query construction	premise [SEP] hypothesis; clean+trim (180/120 chars)
Encoder	all-mpnet-base-v2; cosine over L2-normalized embeddings
Anchors / retrieved shots	3 anchors $+$ top- K =6 nearest neighbors
Pool	Labeled ENT/NEU/CON premise-hypothesis pairs
Prompt assembly	3 anchors + 6 demos + query; enforce JSON schema
Purpose	Improve LLM consistency via topic/structure similarity

C.2 TRAINING RECIPE FOR TEACHER REFRESH

What this shows. Tab. 13 gives the post-curation optimization settings used across teacher models.

Table 13: **Training recipe** (all teachers) used after Lane L/Human curation.

Component	Setting
Optimizer	AdamW (wd=0.05, β_2 =0.98, ϵ =10 ⁻⁸)
Stability	Grad clip = 1.0 ; label smoothing = 0.05
Schedule	Warmup 6–10% \rightarrow cosine decay
Precision	FP16 with dynamic loss scaling
Batching	Small per-GPU batch + gradient accumulation
Memory	Gradient checkpointing
Init freeze	Optional: freeze embeddings for first 300–1000 steps
Early stopping	Dev NLL
Repro	Seeds/hardware in released configs

D ADDITIONAL DEV DIAGNOSTICS

D.1 PER-MODEL DEV CALIBRATION

Reading Table 14. DeBERTa-v3-large attains the highest accuracy (0.664) and the lowest post-TS ECE (0.104). The ensemble shows the largest calibration gain from TS (Δ ECE= -0.067), consistent with Fig. 2 where post-refresh accuracy exceeds confidence in mid bins; its T^* =0.53 sharpens probabilities to close those gaps. RoBERTa and XLM-R require mild softening (T^* >1) and see smaller ECE changes. Overall, TS improves probability quality without affecting accuracy, and the ensemble—used for uncertainty decomposition and routing—achieves competitive post-TS calibration (ECE 0.146) while benefiting most from re-calibration.

D.2 SELECTIVE PREDICTION (RISK-COVERAGE)

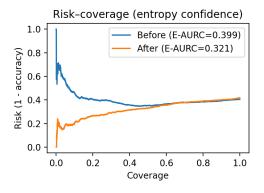
What this shows. Fig. 7 compares risk across coverage before/after refresh for two confidence proxies.

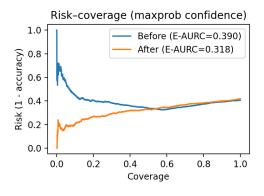
D.3 DEV METRICS WITH CIS (ENSEMBLE)

What this shows. Tab. 15 reports point estimates with 95% CIs for the post-refresh ensemble.

Table 14: ANLI dev calibration by model and ensemble (entropy-based confidence, B=10 bins).

Model	Acc (pre-TS)	ECE (pre-TS)	T^{\star}	ECE@ T^{\star}	Δ ECE
DeBERTa-v3-large	0.664	0.148	0.85	0.104	-0.045
RoBERTa-large	0.533	0.135	1.07	0.123	-0.011
XLM-R-large	0.492	0.150	1.11	0.144	-0.006
Ensemble	0.581	0.209	0.53	0.146	-0.067





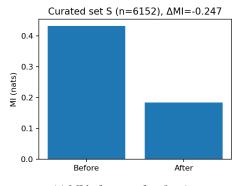
(a) Entropy confidence (E-AURC $0.399 \rightarrow 0.321$).

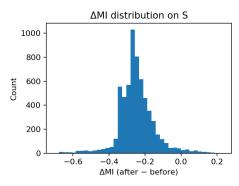
(b) Max-probability confidence (E-AURC 0.390 \rightarrow 0.318).

Figure 7: Selective prediction on ANLI dev. Post-refresh models achieve lower risk across coverage; the largest gap is at low coverage.

Table 15: URC² ensemble (post-refresh) on ANLI dev: metrics with 95% CIs.

Metric	Point	95% CI
Accuracy	0.581	[0.564, 0.598]
Ent-ECE	0.209	[0.194, 0.228]
MaxP-ECE	0.120	[0.104, 0.138]
NLL	0.884	[0.856, 0.913]
Brier	0.543	[0.525, 0.561]





(a) MI before vs. after (bars)

(b) Histogram of Δ MI (after-before)

Figure 8: Curated rows S (full train). Left: MI before vs. after (bars). Right: distribution of Δ MI (after—before). The shift left indicates disagreement reduction on edited items.

E MUTUAL INFORMATION ON CURATED ROWS

Why this matters. Fig. 8 shows that curation reduces epistemic disagreement on edited items.

Table 16: Representative ANLI examples. Premises tend to be longer/richer; hypotheses are succinct and often adversarial.

Premise	Hypothesis	Gold
"Johnson College Prep is a public four-year charter high school located in the Englewood neighborhood on the south side of Chicago, Illinois, United States. It is part of the Noble Network of Charter Schools. The school is named for African-American businessman John H. Johnson and his wife Eunice Johnson."	The Noble Network of Charter Schools has a school that prepares stu- dents for college.	ENT
"I would agree. And I also agree that most police officers, of course, are doing a good job and hate this practice also. I talked to an African-American police officer in Springfield, Massachusetts not long ago who raised this question and said that in his opinion one of the biggest solutions is in the training."	Springfield, Massachusetts has few African-American cops.	NEU
"How to stop erosion on a river bank: use coir netting that is 700–900 grams per square meter (gsm). Coir netting is made from coconut fibers and is biodegradable. Normal or heavy duty netting should be used when there is a large amount of water flow."	When the water has a large amount of flow use fish netting.	CON

F REPRESENTATIVE ANLI EXAMPLES

What this shows. Tab. 16 provides concrete premise-hypothesis pairs illustrating adversarial style and label types.

G ABSOLUTE-THRESHOLD UNCERTAINTY QUADRANTS

For sensitivity, we also report the true 2×2 partition using *absolute* cutoffs shown in Tab. 17.

Table 17: Uncertainty quadrants after refresh with absolute cutoffs.

Quadrant	Share
Low-A / Low-E	0.000
High-A only	0.149
High-E only	0.000
Both-high	0.851

Cutoffs: τ_A =0.300, τ_E =0.150. Corpus means after refresh: A \approx 0.824, E \approx 0.196.