

---

# Conformal Agent Error Attribution

---

Naihe Feng<sup>\*1</sup> Yi Sui<sup>\*2</sup> Shiyi Hou<sup>2</sup> Ga Wu<sup>1</sup> Jesse C. Cresswell<sup>2</sup>

## Abstract

When multi-agent systems (MAS) fail, identifying where the decisive error occurred is the first step for automated recovery to an earlier state. Error attribution remains a fundamental challenge due to the long and intertwined interaction traces that large language model-based MAS generate. This paper presents a framework for error attribution based on conformal prediction (CP) which provides finite-sample, distribution-free coverage guarantees. We introduce new algorithms for filtration-based CP designed for sequential data such as agent trajectories. Unlike existing CP algorithms, our approach predicts sets that are contiguous sequences, which is a crucial property to enable efficient recovery and debugging. We verify our theoretical guarantees on a variety of agents and datasets, show that errors can be precisely isolated, then use prediction sets to rollback MAS to correct their own errors. Our overall approach is model-agnostic, and offers a principled uncertainty layer for MAS error attribution.

## 1. Introduction

Advances in large language models (LLMs) have driven the widespread adoption of multi-agent systems (MAS) for complex tasks requiring decomposition, coordination, and tool use (Guo et al., 2024), with strong empirical performance in domains such as software engineering (Hong et al., 2024; Huang et al., 2023), scientific discovery (Ghafarollahi & Buehler, 2024), and financial decision making (Yu et al., 2024). However, the increased system complexity and rich interactions in MAS make them prone to errors from incorrect intermediate decisions, miscoordination among agents, and long-horizon dependencies (Cemri et al., 2025; Leung et al., 2026). While detecting *overall* task failure is often straightforward, understanding *why and where* a failure

---

<sup>\*</sup>Equal contribution <sup>1</sup>Dalhousie University, Halifax, Canada  
<sup>2</sup>Layer 6 AI, Toronto, Canada. Correspondence to: Jesse C. Cresswell <jesse@layer6.ai>.

Accepted to the ICML 2026 Workshop on Statistical Frameworks for Uncertainty in Agentic Systems, Seoul, South Korea, 2026.  
Copyright 2026 by the author(s).



Figure 1. Conformal agent error attribution isolates the decisive error in a failed MAS trajectory within a conformal prediction set, providing statistical guarantees of coverage.

originated remains challenging yet critical for debugging, self-correction, and reliable deployment. Identifying the decision step that constitutes the decisive error point has emerged as a central challenge for improving MAS.

Most existing MAS error attribution approaches, including naive LLM-as-a-judge methods (Zhang et al., 2025), structured reasoning pipelines (Hengst et al., 2025; Zhu et al., 2025; Yu et al., 2025), and fine-tuned attribution models (Kong et al., 2026), ultimately produce a point prediction, committing to a single responsible step. In practice, point predictions provide limited actionable insight for practitioners as they offer no principled form of uncertainty quantification to assess reliability, undermining the trustworthiness of error attribution systems. Conformal prediction (CP) offers a promising direction for addressing this limitation through the generation of *prediction sets*. CP enables reliable decision-making under uncertainty by providing statistical guarantees across a range of applications (Angelopoulos & Bates, 2021; Cresswell et al., 2024).

Motivated by these developments, we propose an uncertainty-aware framework for error attribution in MAS based on CP, which provides finite-sample, distribution-free coverage guarantees. Rather than predicting a single step, our methods identify a localized region of the execution trace that is guaranteed to contain the decisive error at a user-specified confidence level (see Figure 1). We introduce novel methods for filtration-based CP which are adapted to sequential data structures, like agent trajectories. Unlike existing CP approaches that produce arbitrary sets, our methods produce *contiguous* sets, aligning with the inherent ordinal structure of sequential data. Finally, we use conformal sets to rollback the MAS to before the decisive error, allowing the agent to restart and fix its own mistakes. Our approach is model-agnostic and can wrap existing black-box attribution scores, while providing a principled uncertainty layer for error attribution in real-world MAS.

## 2. Background & Related Work

### 2.1. Post-hoc Multi-agent System Error Attribution

Recent work on error attribution in MAS has primarily studied post-hoc localization of errors using offline execution traces. Early approaches used *naive LLM-as-a-judge*, where a single model directly predicts the responsible step from a failed trace (Zhang et al., 2025). Subsequent work improved attribution quality through more sophisticated LLM pipelines, including *prompt and context engineering* as well as multi-LLM frameworks. For example, ECHO (Banerjee et al., 2025) improves attribution by organizing long traces into hierarchical contexts and aggregating evaluations via consensus, while RAFFLES (Zhu et al., 2025) employs a multi-turn, multi-LLM architecture that iteratively proposes and critiques candidate error steps. Additionally, CORRECT (Yu et al., 2025) incorporates retrieval to localize the error step based on similar events. A complementary line of work *fine-tunes specialized LLM judges* for error attribution. In particular, AEGIS (Kong et al., 2026) constructs large-scale labeled failure traces via controlled error injection for fine-tuning LLMs on the task.

In our experiments we compare the efficacy of these three main classes of evaluators: naive, context-engineered, and fine-tuned LLM judges. We note that all of the above research assumes a single decisive error, whereas in practical MAS applications small errors may accumulate into large ones. Due to the lack of labeled datasets with more nuanced error definitions, we follow current work and focus on the decisive error setting.

### 2.2. Conformal Prediction

For a classification problem where inputs  $x \in \mathcal{X}$  and ground-truth values  $y^* \in \mathcal{Y} = [\ell] := (1, \dots, \ell)$  are drawn jointly from a distribution  $(x, y^*) \sim \mathbb{P}$ , CP first calibrates a threshold  $\hat{q}$  from a set of held-out data. Then for a new datapoint  $x_{n+1}$ , CP outputs a set of classes  $C(x_{n+1}; \hat{q}) \subseteq \mathcal{Y}$  which contains  $y^*$  with high probability

$$\mathbb{P}[y_{n+1}^* \in C(x_{n+1}; \hat{q})] \geq 1 - \alpha. \quad (1)$$

This *coverage* guarantee is distribution-free and valid in finite samples, while also allowing the user to set their own error tolerance  $\alpha$  (Vovk et al., 2005; Shafer & Vovk, 2008).

To perform CP, one first defines a *conformal score* function  $S : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ , which should take smaller values when  $y = y^*$  is the correct label for  $x$ . In practice,  $S(x, y)$  often leverages the predictions of a pre-trained classification model  $f : \mathcal{X} \rightarrow \mathcal{Y}$ . Using a set of  $n$  calibration datapoints, CP computes the scores  $\{S_i\}_{i=1}^n = \{S(x_i, y_i^*)\}_{i=1}^n$ , and finds the  $\frac{[(n+1)(1-\alpha)]}{n}$  quantile which is set as the threshold  $\hat{q}$ . Then prediction sets can be generated by including all classes for which the score is less than  $\hat{q}$ ,

$$C(x_{n+1}; \hat{q}) = \{y \in \mathcal{Y} \mid S(x_{n+1}, y) \leq \hat{q}\}. \quad (2)$$

When  $x_{n+1}$  is exchangeable with the calibration data, Equation (1) is valid. Exchangeability is a mild assumption that automatically holds when data is i.i.d., and hence is reasonable for many machine learning contexts, including the agent error attribution task as we show in our experiments. At equal coverage levels  $1 - \alpha$ , smaller prediction sets are considered more useful both for uncertainty quantification over the predictions of  $f_\theta$  (Romano et al., 2020; Angelopoulos et al., 2021; Huang et al., 2024), and in downstream tasks (Cresswell et al., 2024; 2025).

### 2.3. Conformal Prediction for Sequential Data

Another common setting is where data is sequential,  $x = (c_1, \dots, c_\ell)$ , with variable length  $\ell$ , where the ground truth  $y^* \subset x$  is a subset of elements. Following Kuwahara et al. (2025), the principles of CP can be used to calibrate a threshold  $\hat{q}$ , and predict a subset  $C(x_{n+1}; \hat{q}) \subseteq x_{n+1}$  which retains the ground truth elements with high probability,

$$\mathbb{P}[y_{n+1}^* \subseteq C(x_{n+1}; \hat{q})] \geq 1 - \alpha. \quad (3)$$

In some settings,  $y^*$  will consist of multiple elements, and the predicted set  $C(x_{n+1}; \hat{q})$  need not be contiguous. For agent error attribution we take  $x$  to be the agent’s trajectory, and  $y^*$  to be the single decisive error—one of the  $c_i$ . As we will discuss, for downstream applications including automated rollbacks of the agent’s state it is desirable to predict sets of *consecutive* elements, rather than arbitrary subsets. Hence, we develop novel CP algorithms that satisfy a coverage guarantee using contiguous prediction sets,

$$\mathbb{P}[y_{n+1}^* \in C(x_{n+1}; \hat{q})] \geq 1 - \alpha, C(x_{n+1}; \hat{q}) = (c_j, \dots, c_k). \quad (4)$$

The only existing CP algorithms that produce contiguous sets were designed for hierarchical classification (Mortier et al., 2026). We describe one such algorithm in Section 3.1.2 and adapt it for sequential data.

## 3. Conformal Agent Error Attribution

For the remainder of this work we take  $x = (c_1, \dots, c_\ell)$  to be an agent trajectory which has failed to complete the desired task. Each step  $c_j$  can contain any available information such as the environment’s state, action taken, and observed response. One of the steps  $y^* \in x$  is labeled as the decisive error—the earliest error that the MAS cannot recover from. The aim is to produce a prediction set  $C(x_{n+1}) \subseteq x_{n+1}$  that gives valid coverage as in Equation (1) or Equation (4), where smaller sets are preferred.

Applying CP to agent error attribution requires two components: an algorithm which takes a calibration dataset and generates a prediction set for  $x_{n+1}$  with valid coverage; and a scoring function  $g(C(x))$  acting on sets of steps, which quantifies the likelihood that  $y^* \in C(x)$ . We design these two components separately so that they are interchangeable, and discuss the pros and cons of each option.

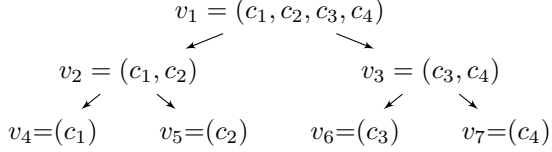


Figure 2. An example binary tree  $\mathcal{T}$  representing an agent trajectory  $x$  consisting of four steps  $c_1, \dots, c_4$ . Contiguous prediction sets  $C(x_{n+1})$  will consist of a single node  $v_i$ .

### 3.1. Conformal Algorithms for Agent Error Attribution

#### 3.1.1. VANILLA CONFORMAL PREDICTION

The simplest approach is to ignore the sequential nature of  $x$  and treat all steps as unordered classes in an  $\ell$ -way classification task. We will write  $S_{\text{VCP}} = 1 - g$  for the conformal score function, and

$$C_{\text{VCP}}(x_{n+1}; \hat{q}) = \{c_i \in x_{n+1} \mid S_{\text{VCP}}(x_{n+1}, c_i) \leq \hat{q}\}, \quad (5)$$

for prediction sets generated by Vanilla CP (VCP). Prediction requires iterating over every step in the trajectory using  $\ell$  evaluations of  $g$ , and does not produce contiguous sets.

#### 3.1.2. LEAF-TO-ROOT TREE TRAVERSAL

To produce contiguous sets, we can adapt algorithms for hierarchical classification by mapping agent trajectories  $x$  onto a binary tree  $\mathcal{T}$  as depicted in Figure 2, with root node  $v_1 = [\ell]$ , and leaf nodes  $v_\ell, \dots, v_{2\ell-1}$  as individual steps  $c_1, \dots, c_\ell$ . CP is conducted by traversing the tree from leaf to root following the CRSVP algorithm (Mortier et al., 2026) described in full detail in Appendix B. For each test datapoint, CRSVP outputs one node of the tree as the prediction set which is always a contiguous set, and guarantees the lower bound on coverage in Equation (4).

CRSVP lacks an upper bound on coverage, uses  $\ell$  evaluations of  $g$  for prediction, and produces inflexible prediction sets following the tree’s splits. For example, in Figure 2 the middle steps  $c_2$  and  $c_3$  can only be predicted together in the trivial case where all steps are predicted ( $v_1$ ). VCP and CRSVP serve as baselines in our experiments. The following novel algorithms improve on their limitations.

#### 3.1.3. LEFT (RIGHT) FILTRATION

Viewing a trajectory  $x$  as a sequence  $(c_1, \dots, c_\ell)$ , Left Filtration (LF) progressively removes steps from  $x$  starting on the left with  $c_1$  until the remaining subsequence scores below a calibrated threshold. In other words, it returns a *suffix* of the full sequence.

We assume access to a scoring function  $g_{\text{LF}}$  which scores how likely a subinterval  $c_{j:k} := (c_j, \dots, c_k) \subseteq x$  is to contain  $y^*$ , imposing the boundary conditions  $g_{\text{LF}}(\emptyset) = 0$ , since the empty interval cannot contain  $y^*$ , and likewise  $g_{\text{LF}}(x) = \infty$  since  $y^* \in x$ . We will use  $j^*$  as the index where  $y^*$  appears, so  $c_{j^*} = y^*$ . Finally, it is desirable, but not mandatory, to have  $g_{\text{LF}}$  obey a monotonicity condition:

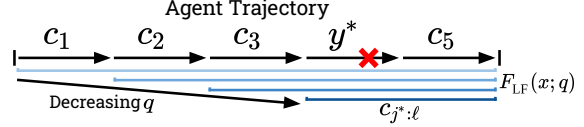


Figure 3. Left filtering with  $F_{\text{LF}}(x; q)$  progressively removes steps from the left as  $q$  decreases. The smallest  $q$  which retains the decisive error  $y^*$  is used as the conformal score  $S_{\text{LF}}(x, y^*)$ .

$$c_{b:c} \subseteq c_{a:d} \implies g_{\text{LF}}(c_{b:c}) \leq g_{\text{LF}}(c_{a:d}). \quad (6)$$

Next, we define a filtering function which returns the longest suffix that has a low enough score  $g_{\text{LF}}$ . Formally, we write the set of suffixes as  $\mathcal{I}_{\text{LF}} := \{c_{j:\ell}\}_{j=1}^{\ell+1}$ , with the conventions for subintervals that  $c_{j:j} = c_j$ , and  $c_{j:k} = \emptyset$  when  $j > k$ . With these conventions, the left-filtering function is

$$F_{\text{LF}}(x; q) := \arg \max_{c_{j:\ell} \in \mathcal{I}_{\text{LF}}} (|c_{j:\ell}| \mid g_{\text{LF}}(c_{j:\ell}) \leq q), \quad (7)$$

where  $q \in \mathbb{R}^+$ , and we note that  $F_{\text{LF}}(x; \infty) = x$ . Since the suffixes in  $\mathcal{I}_{\text{LF}}$  are nested,  $F_{\text{LF}}(x; q)$  also satisfies a nesting property demonstrating that more steps are filtered out as  $q$  decreases. (All formal proofs are given in Appendix A.2.)

**Lemma 3.1.** For any  $x$  and thresholds  $0 \leq q_1 \leq q_2$ ,

$$F_{\text{LF}}(x; q_1) \subseteq F_{\text{LF}}(x; q_2). \quad (8)$$

*Proof sketch.* Suffixes themselves are nested. When  $q_1$  increases to  $q_2$ , the set of valid suffixes with  $g_{\text{LF}}(c_{j:\ell}) \leq q$  can only grow, and  $F_{\text{LF}}$  returns the largest valid suffix.  $\square$

For a trajectory  $x$ , the conformal score is effectively the smallest threshold  $q$  where  $y^*$  is not filtered out (see Figure 3). More formally we define

$$S_{\text{LF}}(x, y^*) := \inf (q \in \mathbb{R}^+ \mid y^* \in F_{\text{LF}}(x; q)). \quad (9)$$

Although this definition involves two optimizations,  $S_{\text{LF}}$  is designed so that its computation greatly simplifies in practice. When  $g_{\text{LF}}$  obeys monotonicity (Equation (6)),  $S_{\text{LF}}$  is simply the value of  $g_{\text{LF}}$  on the suffix that starts at  $y^*$ .

**Proposition 3.2.** When  $g_{\text{LF}}$  obeys monotonicity such that  $c_{b:c} \subseteq c_{a:d} \implies g_{\text{LF}}(c_{b:c}) \leq g_{\text{LF}}(c_{a:d})$ , then

$$S_{\text{LF}}(x, y^*) = g_{\text{LF}}(c_{j^*:\ell}). \quad (10)$$

*Proof sketch.*  $S_{\text{LF}}(x, y^*)$  optimizes over  $q$  where  $F_{\text{LF}}(x; q)$  returns a suffix at least as large as  $c_{j^*:\ell}$ . Due to monotonicity, suffixes larger than  $c_{j^*:\ell}$  cannot have scores less than  $g_{\text{LF}}(c_{j^*:\ell})$ , so  $g_{\text{LF}}(c_{j^*:\ell})$  is the smallest valid threshold.  $\square$

The LF conformal algorithm computes  $S_{\text{LF}}$  for each calibration datapoint, and sets the conformal threshold  $\hat{q}$  as the  $\frac{\lceil (n+1)(1-\alpha) \rceil}{n}$  quantile. For a test datapoint we predict

$$C_{\text{LF}}(x_{n+1}; \hat{q}) = F_{\text{LF}}(x_{n+1}; \hat{q}), \quad (11)$$

which is the longest suffix  $c_{j:\ell}$  such that  $g_{\text{LF}}(c_{j:\ell}) \leq \hat{q}$ . A shorter suffix is preferable since it better isolates the decisive error. This algorithm satisfies Equation (4) for any scoring function  $g_{\text{LF}}$  as defined above. Before proving this coverage guarantee, we present a lemma to assist:

**Lemma 3.3.** For a fixed  $\hat{q}$ , we have

$$S_{LF}(x, y^*) \leq \hat{q} \iff y^* \in F_{LF}(x; \hat{q}). \quad (12)$$

*Proof sketch.* The infimum in  $S_{LF}(x, y^*)$  is always achieved at  $q_{\min} = \min\{g_{LF}(c_{j;\ell}) \mid j \leq j^*\}$ , and  $y^* \in F_{LF}(x; q_{\min})$ . Because  $F_{LF}$  is nested in  $q$  (Lemma 3.1), increasing  $q_{\min}$  to  $\hat{q}$  maintains  $y^* \in F_{LF}(x; \hat{q})$ .  $\square$

With these facts established, we can prove the coverage guarantee of Equation (4) using a standard conformal argument.

**Theorem 3.4.** Suppose  $\{(x_i, y_i^*)\}_{i=1}^n$  and  $(x_{n+1}, y_{n+1}^*)$  are exchangeable. Given a scoring function  $g_{LF}$  acting on subintervals of the  $x_i$ , define the conformal score  $S_{LF}(x_i, y_i^*)$  as in Equation (9). Let  $\hat{q}$  be the  $\frac{[(n+1)(1-\alpha)]}{n}$  quantile of conformal scores  $\{S_i\}_{i=1}^n = \{S_{LF}(x_i, y_i^*)\}_{i=1}^n$ . Then prediction sets constructed as  $C_{LF}(x_{n+1}; \hat{q}) = F_{LF}(x_{n+1}; \hat{q})$  satisfy  $1 - \alpha \leq \mathbb{P}[y_{n+1}^* \in C_{LF}(x_{n+1}; \hat{q})] < 1 - \alpha + \frac{1}{n+1}$ .

*Proof.* Since the fixed function  $S_{LF}$  is applied to each element of the exchangeable sequence  $(x_1, y_1^*), \dots, (x_{n+1}, y_{n+1}^*)$ , the resulting sequence  $S_1, \dots, S_{n+1}$  is also exchangeable. We assume for simplicity that the scores are non-degenerate, because noise can easily be added to break ties. Let  $S_{(1)} \leq S_{(2)} \leq \dots \leq S_{(n)}$  be the order statistics of the calibration scores so that the empirical quantile  $\hat{q}$  can be defined as  $S_{(k)}$ , where  $k = [(n+1)(1-\alpha)]$  (assuming  $n \geq \frac{1}{\alpha} - 1$  to ensure  $k \leq n$ ). The event  $S_{n+1} \leq \hat{q}$  occurs if and only if  $S_{n+1}$  is among the  $k$  smallest scores overall. From exchangeability, all orderings are equally likely, so  $\mathbb{P}[S_{n+1} \leq \hat{q}] = \frac{k}{n+1}$ . By Lemma 3.3 this event is equivalent to  $y^* \in F_{LF}(x_{n+1}; \hat{q})$ . Hence,

$$\mathbb{P}[y^* \in F_{LF}(x_{n+1}; \hat{q})] = \frac{[(n+1)(1-\alpha)]}{n+1} \geq 1 - \alpha. \quad (13)$$

For the upper bound, we simply note that

$$\frac{[(n+1)(1-\alpha)]}{n+1} < \frac{1+(n+1)(1-\alpha)}{n+1} = 1 - \alpha + \frac{1}{n+1}. \quad (14)$$

$\square$

Beyond generating contiguous prediction sets, LF uses fewer scoring function evaluations on average for inference when  $g_{LF}$  is monotonic. LF starts with the first suffix  $c_{\ell;\ell}$  and evaluates each longer suffix  $c_{j;\ell}$  until one with  $g_{LF}(c_{j;\ell}) > q$  is found, then returns  $c_{j+1;\ell}$ . When errors are evenly distributed over the length  $\ell$ , the average number of evaluations with a strong  $g_{LF}$  is only  $\frac{\ell+1}{2}$ .

Of course, there is nothing special about filtering from the left; Right Filtration (RF) can easily be defined which filters from the right, returning a *prefix* of  $x$ . Using a similarly defined scoring function  $g_{RF}$  we write the relevant prefixes as  $\mathcal{I}_{RF} := \{c_{1:k}\}_{k=0}^{\ell}$  and define the right-filtering function

$$F_{RF}(x; q) := \arg \max_{c_{1:k} \in \mathcal{I}_{RF}} (|c_{1:k}| \mid g_{RF}(c_{1:k}) \leq q) \quad (15)$$

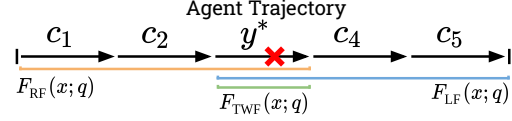


Figure 4. Two-way filtering with  $F_{TWF}(x; q)$  uses the intersection of filtering from the right and left. The smallest  $q$  which retains the decisive error  $y^*$  is used as the conformal score  $S_{TWF}(x, y^*)$ .

Based on these definitions, the conformal score is

$$S_{RF}(x, y^*) := \inf (q \in \mathbb{R}^+ \mid y^* \in F_{RF}(x; q)). \quad (16)$$

and prediction sets are generated as  $C_{RF}(x_{n+1}; \hat{q}) = F_{RF}(x_{n+1}; \hat{q})$ . RF gives coverage as a straightforward extension of Theorem 3.4 and is advantageous over LF if agents tend to fail earlier in their trajectory rather than later.

### 3.1.4. TWO-WAY FILTRATION

Filtering from one direction has the obvious downside that when the decisive error is at the start (end), the suffix (prefix) which covers the ground truth will contain the entire trajectory. Moreover, when the decisive error is near the middle of the trajectory, neither LF nor RF can isolate it. However, these limitations can be avoided by building on LF and RF with Two-Way Filtration (TWF). Bidirectional filtering allows more precise isolation of the decisive error in principle, regardless of where in the trajectory it occurs.

There are many possible ways to define a bidirectional filtration. We present a version that uses the *intersection* of left and right filtered subintervals. Using  $g_{LF}$  and  $g_{RF}$  as above, we define the two-way filtering function as

$$F_{TWF}(x; q) := F_{LF}(x; q) \cap F_{RF}(x; q). \quad (17)$$

This function finds the longest suffix and prefix that each score at most  $q$ , and takes their intersection. Note that this function still satisfies  $F_{TWF}(x; \infty) = x$ , as well as nesting: (All formal proofs are given in Appendix A.3)

**Lemma 3.5.** For any  $x$  and thresholds  $0 \leq q_1 \leq q_2$ ,

$$F_{TWF}(x; q_1) \subseteq F_{TWF}(x; q_2). \quad (18)$$

*Proof sketch.* The result follows from nesting for L/RF, and set inclusion under intersection: for any sets  $A, B, C, D$ , if  $A \subseteq B$  and  $C \subseteq D$ , then  $A \cap C \subseteq B \cap D$ .  $\square$

On top of this we can define the conformal score function

$$S_{TWF}(x, y^*) := \inf (q \in \mathbb{R}^+ \mid y^* \in F_{TWF}(x; q)), \quad (19)$$

which operates the same way as before, effectively finding the smallest threshold  $q$  under which two-way filtering retains the decisive error. Although  $S_{TWF}$  appears to involve three optimizations, its computation can be simplified to only two scoring function evaluations:

Table 1. Conformal Algorithm Properties

Method	Contiguous	Inference NFE	Cov. Upper Bound
VCP	✗	$\ell$	✓
CRSVP	✓	$\ell$	✗
Left/Right Filtering	✓	$\approx \frac{1}{2}(\ell + 1)$	✓
Two-Way Filtering	✓	$\ell$	✓

**Proposition 3.6.**  $S_{TWF}(x, y^*)$  as defined in Equation (19) can be expressed as

$$S_{TWF}(x, y^*) = \max(S_{LF}(x, y^*), S_{RF}(x, y^*)). \quad (20)$$

When  $g_{LF}$  and  $g_{RF}$  obey the monotonicity condition in Equation (6),  $S_{TWF}$  further simplifies to

$$S_{TWF}(x, y^*) = \max(g_{LF}(c_{j^*:\ell}), g_{RF}(c_{1:j^*})). \quad (21)$$

*Proof sketch.*  $q$  is a valid threshold ( $y^* \in F_{TWF}(x; q)$ ) only when  $y^*$  is included by both  $F_{LF}$  and  $F_{RF}$ , which means  $q$  is sufficiently high for both  $S_{LF}$  and  $S_{RF}$ . We must use the higher of these two thresholds, giving Equation (20).

Equation (21) follows from substituting in Equation (10).  $\square$

Like for LF, the TWF conformal algorithm computes  $S_{TWF}$  for each calibration datapoint, sets the conformal threshold  $\hat{q}$  as the  $\frac{\lceil (n+1)(1-\alpha) \rceil}{n}$  quantile, and predicts

$$C_{TWF}(x_{n+1}; \hat{q}) = F_{TWF}(x_{n+1}; \hat{q}), \quad (22)$$

for any test datapoint  $x_{n+1}$ .  $F_{TWF}(x_{n+1}; \hat{q})$  will tend to be a short subinterval when  $g_{LF}$  and  $g_{RF}$  both narrow in on the same steps. This algorithm satisfies Equation (4) with a theorem and proof similar to Theorem 3.4 (see Theorem A.2). The core prerequisite is the analogue of Lemma 3.3:

**Lemma 3.7.** For a fixed  $\hat{q}$ , we have

$$S_{TWF}(x, y^*) \leq \hat{q} \iff y^* \in F_{TWF}(x; \hat{q}). \quad (23)$$

*Proof sketch.* From the use of intersection in  $F_{TWF}(x; \hat{q})$  (Equation (17)), and the result of Lemma 3.3,

$$y^* \in F_{TWF}(x, \hat{q}) \iff S_{LF}(x, y^*) \leq \hat{q} \wedge S_{RF}(x, y^*) \leq \hat{q}. \quad (24)$$

Equivalently we write  $\max(S_{LF}(x, y^*), S_{RF}(x, y^*)) \leq \hat{q}$ , which is the same as  $S_{TWF}(x, y^*) \leq \hat{q}$  (Proposition 3.6).  $\square$

While we described these algorithms using the language of agent trajectories, they can equally be applied to any type of sequential data  $x$  with ground truth  $y^* \in x$ , for example electronic health records where  $x$  is a sequence of vitals measurements, and  $y^*$  is the first warning sign that should trigger medical intervention (Razavian & Sontag, 2015).

The five conformal algorithms we compare are summarized in Table 1, showing which generate contiguous prediction sets, the number of function evaluations (NFE) of  $g$  needed for each test datapoint in the case where  $g$  is strong and errors are uniformly distributed, and whether an upper bound on coverage is guaranteed.

## 3.2. Scoring Functions for Agent Error Attribution

Each of the preceding conformal algorithms requires a scoring function  $g$  which estimates how likely a step, or set of steps, is to contain the decisive error. Since MAS traces are primarily composed of textual agent interactions, our scoring functions rely on LLM-based components to map agent inputs and outputs to numerical scores. We compare three LLM regimes as outlined in Section 2.1:

1. *Naive LLM-as-a-judge* is implemented by prompting gpt-4o-mini to estimate error likelihoods with information about the task and step, as in (Zhang et al., 2025);
2. *Prompt and context engineering* produces step-level likelihood estimates using multiple LLMs with role-specific prompts, inspired by ECHO (Banerjee et al., 2025). Scores from LLMs with different roles are averaged to get a final likelihood.
3. A *fine-tuned specialized LLM* is implemented by following the data generation and training paradigm of AEGIS (Kong et al., 2026) to fine-tune a Qwen3-1.7B model (Yang et al., 2025).

Additional details on prompts, training, and data are given in Appendix C. For VCP, step-level scores as described above are used directly. However, L/R/TWF and CRSVP require set-level scores. To enable efficient computation in L/R/TWF we design  $g$  to obey monotonicity (Equation (6)) by aggregating step-level scores. Specifically, we use summation and normalize by the trajectory length  $\ell$  to make conformal scores for different datapoints more comparable:

$$g(c_{j:k}) = \frac{1}{\ell} \sum_{i=j}^k \text{LLM}(c_i). \quad (25)$$

We experiment with alternative monotonic aggregations, namely Max and LogSumExp, in Appendix C.2, and discuss circumstances when they may be preferred. Below, we evaluate the discriminatory power of  $g$  independently from CP algorithms by treating it as a multiclass classifier.

## 3.3. Conformal Agent Rollbacks

CP sets for agent error attribution have multiple uses. They can be used by humans for manual debugging; a person can focus only on the steps within the set and have good coverage of true errors. Here, contiguity is a great benefit—it is much easier to debug several consecutive steps than to parse through a scattered set.

However, our main downstream application is automated agent recovery. Knowing that an agent has failed, we wish for it to learn from its mistakes and retry parts of the task. Optimally, the agent’s state must be rolled back in the trajectory far enough to cover the decisive error. Rolling back further is inefficient. However, from prior work on error attribution accuracy is low for point prediction (Zhang et al., 2025; Banerjee et al., 2025; Zhu et al., 2025).

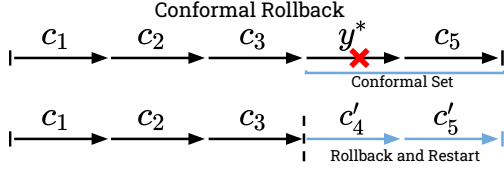


Figure 5. After generating a conformal set for a failed task, we roll back the state of the MAS to the first step in the set, and restart the agent with instructions to avoid the same mistakes.

Prediction sets with coverage guarantees provide the solution for automated rollbacks of failed trajectories. Given a conformal set, we roll back the state of the MAS to the first step in the set, and restart the trajectory, as depicted in Figure 5. The coverage guarantee ensures with high confidence that we roll back far enough to correct the decisive error, while contiguity ensures that we don’t roll back excessively far. Upon restarting, we add information to the prompt about the steps taken previously, and instruct the MAS to replan its trajectory to avoid making the same mistakes.

## 4. Experimental Setup

### 4.1. Datasets

We evaluate CP algorithms, scoring functions, and downstream task performance on both a real-world benchmark and synthetic MAS traces. The **Who&When** dataset (Zhang et al., 2025) is a benchmark for step-level error attribution in MAS. Each of the 184 real-world examples contains a full agent execution trace annotated with the decisive error step.<sup>1</sup> Who&When is small-scale and provides limited variety over agents and tasks, but it is the only existing academic dataset with step-level error annotations by humans.

In the absence of other real-world data, and to provide greater variety, we follow existing practices (Kong et al., 2026) to synthetically generate failed agent trajectories through **error injection**. To induce errors at controlled steps, we inject instructions to create errors into the agent’s prompt, using the error mode taxonomy of Cemri et al. (2025). For a given trajectory, one step is selected uniformly at random, and the agent is restarted in that state with a modified prompt describing the desired error mode. Given that the overall trajectory fails, the injected step is labeled as the decisive error. We use two mathematical reasoning datasets, MATH (Hendrycks et al., 2021) and GSM8k (Cobbe et al., 2021), paired with two representative multi-agent architectures, MACNET (Qian et al., 2025) and DyLAN (Liu et al., 2024). For each dataset–architecture combination, we generate 2000 failed trajectories.

Complex tasks generally have some steps that are harder to complete than others, leading to non-uniform decisive error distributions. This distribution is shown in Figure 6 for the

<sup>1</sup>We remove inconsistent datapoints where the error index is greater than the trajectory length.

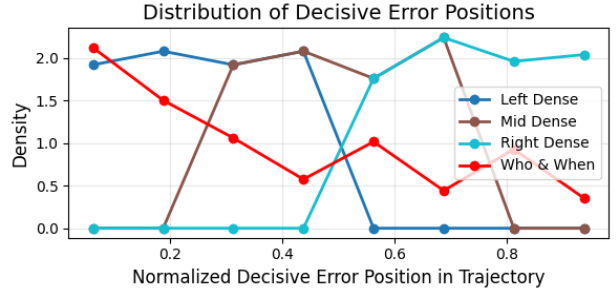


Figure 6. Normalized decisive error position distributions for the real-world benchmark (Who&When), and controlled error injection datasets (Left/Mid/Right Dense) on GSM8k.

Who&When dataset, which tends to have decisive errors early on. To mimic this property, and explore the impact of distribution on error attribution performance, we construct variants of the synthetic data by dividing the normalized trajectory length into thirds, and subsampling conditional on decisive error location. These are the **Left**, **Mid**, and **Right Dense** variants of GSM8k in Figure 6. Additional details and analysis are provided in Appendix C.

### 4.2. Evaluation Metrics

Scoring functions  $g$  can be viewed as classifiers on an  $\ell$ -way task—predicting the decisive error location. Hence, we evaluate them using AUROC, AUPRC, and accuracy. Since  $\ell$  can be large, baseline levels of these metrics are very low.

Given a test dataset  $\{(x_i, y_i^*)\}_{i=n+1}^{n+m}$  and predicted sets  $\{C(x_i)\}_{i=n+1}^{n+m}$ , we evaluate conformal agent error attribution with the following metrics. **Empirical Coverage** measures how often the predicted set contains the decisive error:

$$\text{EC} = \frac{1}{m} \sum_{i=n+1}^{n+m} \mathbb{I}[y_i^* \in C(x_i)]. \quad (26)$$

EC should be at least  $1 - \alpha$ , the target coverage, and respect the corresponding upper bound where applicable. **Removal Rate** measures how much of the trajectory is filtered out, essentially measuring prediction set size:

$$\text{RR} = \frac{1}{m} \sum_{i=n+1}^{n+m} \left(1 - \frac{|C(x_i)|}{\ell_i}\right), \quad (27)$$

where  $\ell_i$  denotes the number of steps in  $x_i$ . A higher removal rate indicates more precise localization of the error step, given equal EC.

Finally, rollback performance is measured across three axes: Success Rate - the fraction of tasks successfully completed after the rollback; Coverage - the fraction of tasks where the state is rolled back at least to the decisive error; and Cost - the fraction of steps rolled back, needing to be redone. As a baseline, we test the agent restarting from the single most likely predicted step (Top-1).

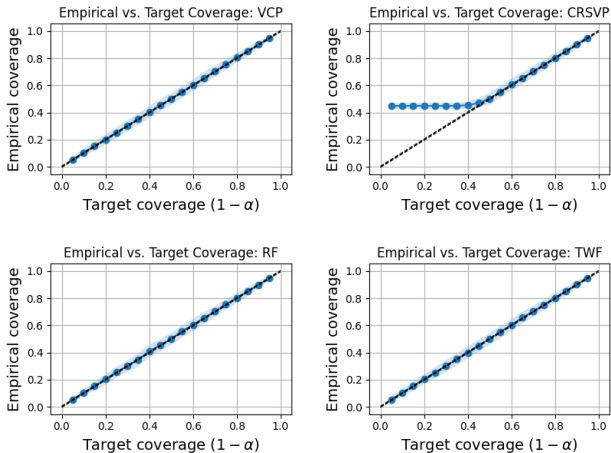


Figure 7. Empirical Coverage vs. Target Coverage ( $1 - \alpha$ ) for CP methods with the fine-tuned LLM scoring function on MATH.

## 5. Experiments & Results

### 5.1. Empirical Verification of Coverage Guarantee

First, we empirically verify that CP algorithms considered in this work satisfy their respective coverage guarantees from Section 3. Figure 7 shows empirical coverage as a function of the target coverage  $1 - \alpha$  for VCP, CRSVP, RF, and TWF, where we average EC over 1000 random, even splits of the calibration and test sets. The shaded regions show one standard deviation. EC is above the lower bound for all methods, consistent with the theoretical guarantees. VCP, RF, and TWF all obey their respective upper bounds, leading to exact linear behaviour. A deviation is observed in the low  $1 - \alpha$  regime for CRSVP, but this is not a violation as CRSVP does not provide an upper bound on coverage.

### 5.2. Scoring Function Evaluation

Next, we compare the three regimes of scoring functions for discriminatory power as binary classifiers. Results are shown in Table 2 on a combination of GSM8k and MATH test sets, matching the training distribution for the fine-tuned LLM. We find that increasing the complexity of  $g$  does pay off; fine-tuning on the error classification task considerably increases its power. The naive use of an LLM is hardly better than random guessing, while prompt engineering alone provides small improvements. These results are consistent with performance levels observed in prior studies (Zhang et al., 2025; Banerjee et al., 2025; Kong et al., 2026).

### 5.3. Conformal Agent Error Attribution Evaluation

Our main comparison of error attribution examines the removal rate (Equation (27)) for conformal algorithms and scoring functions across real and datasets, including trajectories created through two agentic frameworks. Table 3 displays the results using a target coverage rate of 80%. Results are shown as the mean and standard deviation over

Table 2. Step-level discrimination metrics across scoring functions on a combined MATH and GSM8k test set.

Scoring function	AUROC	AUPRC	Accuracy
Random Guessing (baseline)	0.500	0.118	0.118
Naive LLM (gpt-4o-mini)	0.519	0.128	0.110
Role Prompted (gpt-4o-mini)	0.554	0.161	0.164
Fine-tuned (Qwen3-1.7B)	<b>0.762</b>	<b>0.382</b>	<b>0.731</b>

1000 random splits of the calibration and test data.

Our first observation is that the efficacy of conformal algorithms depends on the distribution of errors in the dataset. The real-world Who&When dataset has decisive errors heavily skewed toward early steps in the trajectory (see Figure 6). This distribution naturally aligns with RF which can return a short prefix, and accordingly shows the strongest performance. The synthetic datasets demonstrate this behaviour even more clearly. LF is the strongest algorithm on right-dense data, and TWF successfully isolates errors in the middle of trajectories. Conformal error attribution is most successful when the CP algorithm is matched to the data’s error distribution. Luckily, this is simple to achieve in practice. Applying CP already requires a calibration dataset from a distribution similar to the test set. Practitioners can visualize the error distribution of their calibration data like in Figure 6, and choose the appropriate filtering algorithm, or automate the selection based on the mean error location of trajectories, for example.

Notably, filtering algorithms are less dependent on the discriminatory power of the scoring function than VCP and CRSVP. While the baseline approaches tend to require the fine-tuned LLM scoring function to achieve high RR, filtering methods are largely insensitive. This means that L/R/TWF can be used with simple, easier to design scoring functions as long as they are properly matched to the data’s error distribution. This is especially important for real-world data where labeling is expensive; Who&When only supplies 184 labeled datapoints—enough to calibrate and test, but not nearly enough to also fine-tune an LLM. Our filtering algorithms remain more practical in this setting.

### 5.4. Computational Efficiency

As noted throughout Section 3 and Table 1, at inference existing conformal methods including VCP and CRSVP require evaluating the scoring function  $g$  on each step in the trajectory. LF and RF need only evaluate steps from one direction until their threshold is crossed, which can be much more efficient. To demonstrate this quantitatively, we count the average NFE of  $g$  used by each algorithm when predicting on the GSM8k test set variants. The result in Table 4 confirms that VCP, CRSVP, and TWF use exactly  $\ell$  function calls per trajectory. In contrast, LF shows a clear advantage on right-dense error distributions, while RF uses

## Conformal Agent Error Attribution

Table 3. Removal rate ( $\pm$  std over 1000 calibration/test splits) for conformal algorithms and scoring functions at  $1 - \alpha = 0.8$ . Higher removal rates indicate more precise error attribution, with the best performing conformal algorithm in each column bolded among algorithms which produce contiguous sets. Averages are taken over conformal algorithms to evaluate scoring functions.  $\times$  indicates that fine-tuning is not possible due to insufficient data.

	Conformal Method	Who&When	GSM8k Left Dense		GSM8k Mid Dense		GSM8k Right Dense		MATH	
			DyLAN	MACNET	DyLAN	MACNET	DyLAN	MACNET	DyLAN	MACNET
Naive LLM gpt-4o-mini	Vanilla	0.20 $\pm$ 0.06	0.21 $\pm$ 0.02	0.24 $\pm$ 0.03	0.18 $\pm$ 0.03	0.15 $\pm$ 0.03	0.16 $\pm$ 0.03	0.17 $\pm$ 0.03	0.19 $\pm$ 0.02	0.20 $\pm$ 0.03
	CRSVP	0.21 $\pm$ 0.04	0.58 $\pm$ 0.03	0.54 $\pm$ 0.02	0.29 $\pm$ 0.04	0.21 $\pm$ 0.04	0.21 $\pm$ 0.03	0.10 $\pm$ 0.03	0.30 $\pm$ 0.03	0.20 $\pm$ 0.03
	Left Filtering	0.12 $\pm$ 0.04	0.10 $\pm$ 0.01	0.24 $\pm$ 0.01	0.31 $\pm$ 0.01	0.47 $\pm$ 0.02	0.53 $\pm$ 0.01	<b>0.60<math>\pm</math>0.02</b>	0.18 $\pm$ 0.02	0.19 $\pm$ 0.02
	Right Filtering	<b>0.31<math>\pm</math>0.04</b>	<b>0.64<math>\pm</math>0.02</b>	0.71 $\pm$ 0.01	0.43 $\pm$ 0.02	0.45 $\pm$ 0.01	0.20 $\pm$ 0.02	0.20 $\pm$ 0.01	<b>0.27<math>\pm</math>0.02</b>	0.20 $\pm$ 0.02
	Two-Way	0.19 $\pm$ 0.04	0.17 $\pm$ 0.03	0.22 $\pm$ 0.03	<b>0.59<math>\pm</math>0.03</b>	<b>0.59<math>\pm</math>0.01</b>	0.40 $\pm$ 0.03	0.21 $\pm$ 0.03	<b>0.27<math>\pm</math>0.02</b>	0.15 $\pm$ 0.03
	Average		0.21	0.34	0.39	0.36	0.37	0.30	0.26	0.24
Role-prompted gpt-4o-mini	Vanilla	0.19 $\pm$ 0.06	0.30 $\pm$ 0.04	0.24 $\pm$ 0.04	0.27 $\pm$ 0.04	0.27 $\pm$ 0.03	0.22 $\pm$ 0.04	0.25 $\pm$ 0.04	0.19 $\pm$ 0.03	0.24 $\pm$ 0.04
	CRSVP	0.23 $\pm$ 0.04	0.27 $\pm$ 0.04	0.21 $\pm$ 0.05	0.27 $\pm$ 0.04	0.26 $\pm$ 0.05	0.27 $\pm$ 0.04	0.22 $\pm$ 0.05	<b>0.27<math>\pm</math>0.03</b>	0.22 $\pm$ 0.03
	Left Filtering	0.12 $\pm$ 0.04	0.10 $\pm$ 0.02	0.31 $\pm$ 0.01	0.36 $\pm$ 0.01	0.53 $\pm$ 0.01	<b>0.64<math>\pm</math>0.02</b>	<b>0.60<math>\pm</math>0.02</b>	0.19 $\pm$ 0.02	0.19 $\pm$ 0.02
	Right Filtering	<b>0.30<math>\pm</math>0.05</b>	<b>0.63<math>\pm</math>0.02</b>	0.59 $\pm$ 0.02	0.43 $\pm$ 0.01	0.33 $\pm$ 0.01	0.21 $\pm$ 0.01	0.10 $\pm$ 0.04	<b>0.28<math>\pm</math>0.02</b>	0.19 $\pm$ 0.02
	Two-Way	0.22 $\pm$ 0.02	0.18 $\pm$ 0.02	0.21 $\pm$ 0.02	<b>0.62<math>\pm</math>0.02</b>	<b>0.59<math>\pm</math>0.02</b>	0.42 $\pm$ 0.02	0.20 $\pm$ 0.02	<b>0.29<math>\pm</math>0.02</b>	0.18 $\pm$ 0.02
	Average		0.21	0.30	0.31	0.39	0.40	0.35	0.27	0.24
Fine-tuned Qwen3-1.7B	Vanilla	$\times$	0.78 $\pm$ 0.03	0.82 $\pm$ 0.02	0.66 $\pm$ 0.03	0.75 $\pm$ 0.02	0.34 $\pm$ 0.04	0.63 $\pm$ 0.03	0.33 $\pm$ 0.03	0.62 $\pm$ 0.03
	CRSVP	$\times$	<b>0.69<math>\pm</math>0.06</b>	<b>0.80<math>\pm</math>0.02</b>	0.36 $\pm$ 0.05	0.38 $\pm$ 0.07	0.24 $\pm$ 0.03	0.29 $\pm$ 0.06	<b>0.32<math>\pm</math>0.03</b>	<b>0.30<math>\pm</math>0.04</b>
	Left Filtering	$\times$	0.09 $\pm$ 0.01	0.10 $\pm$ 0.02	0.31 $\pm$ 0.02	0.35 $\pm$ 0.01	0.52 $\pm$ 0.01	<b>0.60<math>\pm</math>0.01</b>	0.18 $\pm$ 0.02	0.19 $\pm$ 0.02
	Right Filtering	$\times$	<b>0.65<math>\pm</math>0.02</b>	0.63 $\pm$ 0.02	0.44 $\pm$ 0.01	0.33 $\pm$ 0.01	0.20 $\pm$ 0.01	0.10 $\pm$ 0.02	0.27 $\pm$ 0.01	0.20 $\pm$ 0.02
	Two-Way	$\times$	0.18 $\pm$ 0.02	0.20 $\pm$ 0.05	<b>0.63<math>\pm</math>0.02</b>	<b>0.60<math>\pm</math>0.02</b>	0.40 $\pm$ 0.02	0.19 $\pm$ 0.05	<b>0.29<math>\pm</math>0.02</b>	0.19 $\pm$ 0.02
	Average	$\times$		0.48	0.51	0.48	0.48	0.33	0.36	0.28

Table 4. Inference cost comparison for conformal methods in terms of average NFE of  $g$  (Fine-tuned LLM) per trajectory with 80% target coverage on GSM8k variants.

Conformal Method	Left Dense	Mid Dense	Right Dense
Vanilla	8.50	8.50	8.50
CRSVP	8.50	8.50	8.50
Left Filtering	7.50	5.64	<b>3.70</b>
Right Filtering	<b>3.08</b>	<b>5.16</b>	7.09
Two-Way	8.50	8.50	8.50

a mere 36% of the calls on left-dense data where most errors occur early. This result emphasizes the cost advantages of selecting an algorithm that has good implicit bias for the distribution observed in calibration data.

### 5.5. Conformal Agent Rollbacks

As noted in Section 3.3, CP sets can enable automated rollbacks for failed tasks with high confidence that decisive errors are captured. We evaluate automated rollbacks on the GSM8K test set variants. For each dataset, we randomly sample 100 MACNET trajectories and score them with the fine-tuned LLM scoring function. As a baseline method, we use the single step in the trajectory with the highest likelihood (Top-1) as the restart point. For conformal rollbacks we predict a conformal set using either the VCP or LF algorithm with a target coverage of 80%, and then restart at the first step in the set. In all cases we restart the MACNET agent with a modified prompt containing the previously failed trace to help the system avoid repeating the same mistakes. As shown in Table 5, LF achieves the highest success rate by a small margin, within one standard

Table 5. Automated rollback results using CP sets from VCP and LF compared to Top-1 predictions. Datasets are variants of GSM8k, and the scoring function uses the fine-tuned LLM.

Dataset	Method	Success Rate $\uparrow$	Coverage	Cost $\downarrow$
Left Dense	Top-1	0.76 $\pm$ 0.05	0.92 $\pm$ 0.03	0.83 $\pm$ 0.02
	VCP	0.73 $\pm$ 0.05	0.85 $\pm$ 0.03	0.79 $\pm$ 0.02
	LF	0.77 $\pm$ 0.04	0.81 $\pm$ 0.05	0.90 $\pm$ 0.00
Mid Dense	Top-1	0.67 $\pm$ 0.04	0.86 $\pm$ 0.04	0.56 $\pm$ 0.02
	VCP	0.59 $\pm$ 0.05	0.95 $\pm$ 0.02	0.65 $\pm$ 0.01
	LF	0.68 $\pm$ 0.05	0.81 $\pm$ 0.03	0.65 $\pm$ 0.01
Right Dense	Top-1	0.71 $\pm$ 0.05	0.86 $\pm$ 0.04	0.50 $\pm$ 0.02
	VCP	0.70 $\pm$ 0.05	1.00 $\pm$ 0.00	0.66 $\pm$ 0.02
	LF	0.75 $\pm$ 0.04	0.82 $\pm$ 0.04	0.40 $\pm$ 0.01

deviation of baselines. However, LF is the only method which gives control over coverage: Top-1 gives no guarantee, while VCP overcovers in this setting because its sets are not contiguous.

## 6. Conclusion

In this work we designed novel conformal set-prediction algorithms for the agent error attribution task. By taking account of the data’s sequential nature, our filtering-based conformal algorithms showed strong ability to precisely isolate decisive errors in agent trajectories within contiguous sets, and with much less compute needed at inference time. Our conformal algorithms can be used as one part of an agent improvement pipeline, allowing humans to more efficiently understand and debug errors in their agents, or enabling fully automated rollbacks wherein agents correct their own mistakes.

## References

- Angelopoulos, A. N. and Bates, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv:2107.07511*, 2021.
- Angelopoulos, A. N., Bates, S., Jordan, M., and Malik, J. Uncertainty sets for image classifiers using conformal prediction. In *International Conference on Learning Representations*, 2021.
- Banerjee, A., Nair, A., and Borogovac, T. Where did it all go wrong? a hierarchical look into multi-agent error attribution. In *NeurIPS 2025 Workshop on Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling*, 2025.
- Cemri, M., Pan, M. Z., Yang, S., Agrawal, L. A., Chopra, B., Tiwari, R., Keutzer, K., Parameswaran, A., Klein, D., Ramchandran, K., Zaharia, M., Gonzalez, J. E., and Stoica, I. Why Do Multi-Agent LLM Systems Fail? In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2025.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems. *arXiv:2110.14168*, 2021.
- Cresswell, J. C., Sui, Y., Kumar, B., and Vouitsis, N. Conformal Prediction Sets Improve Human Decision Making. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pp. 9439–9457, 2024.
- Cresswell, J. C., Kumar, B., Sui, Y., and Belbahri, M. Conformal Prediction Sets Can Cause Disparate Impact. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Ghafarirollahi, A. and Buehler, M. J. SciAgents: Automating Scientific Discovery Through Bioinspired Multi-Agent Intelligent Graph Reasoning. *Advanced Materials*, 37 (22), 2024. doi: 10.1002/adma.202413523.
- Guo, T., Chen, X., Wang, Y., Chang, R., Pei, S., Chawla, N. V., Wiest, O., and Zhang, X. Large language model based multi-agents: A survey of progress and challenges. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, 2024. ISBN 978-1-956792-04-1. doi: 10.24963/ijcai.2024/890.
- Gupta, C., Kuchibhotla, A. K., and Ramdas, A. Nested conformal prediction and quantile out-of-bag ensemble methods. *Pattern Recognition*, 127:108496, 2022. ISSN 0031-3203. doi: 10.1016/j.patcog.2021.108496.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring Mathematical Problem Solving With the MATH Dataset. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1, 2021.
- Hengst, F. d., Blin, I., Mohammadi, M., Shah, S. I. H., and Younesian, T. Hierarchical conformal classification. *arXiv:2508.13288*, 2025.
- Hong, S., Zhuge, M., Chen, J., Zheng, X., Cheng, Y., Wang, J., Zhang, C., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., Ran, C., Xiao, L., Wu, C., and Schmidhuber, J. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024.
- Huang, D., Zhang, J. M., Luck, M., Bu, Q., Qing, Y., and Cui, H. AgentCoder: Multi-Agent-based Code Generation with Iterative Testing and Optimisation. *arXiv:2312.13010*, 2023.
- Huang, J., Xi, H., Zhang, L., Yao, H., Qiu, Y., and Wei, H. Conformal prediction for deep classifier via label ranking. In *Proceedings of the 41st International Conference on Machine Learning*, 2024.
- Kong, F., Zhang, R., Yin, H., Zhang, G., Zhang, X., Chen, Z., Zhang, Z., Zhang, X., Zhu, S.-C., and Feng, X. Aegis: Automated error generation and attribution for multi-agent systems. In *The Fourteenth International Conference on Learning Representations*, 2026.
- Kuwahara, B., Lin, C.-Y., Huang, X. S., Leung, K. K., Yapeter, J. A., Stanevich, I., Perez, F., and Cresswell, J. C. Document summarization with conformal importance guarantees. In *Advances in Neural Information Processing Systems*, volume 38, 2025.
- Leung, K. K., Belbahri, M., Sui, Y., Labach, A., Zhang, X., Rose, S. A., and Cresswell, J. C. Classifying and addressing the diversity of errors in retrieval-augmented generation systems. In *Proceedings of the 19th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3185–3207, 2026. ISBN 979-8-89176-380-7. doi: 10.18653/v1/2026.eacl-long.147.
- Liu, Z., Zhang, Y., Li, P., Liu, Y., and Yang, D. A Dynamic LLM-Powered Agent Network for Task-Oriented Agent Collaboration. In *First Conference on Language Modeling*, 2024.
- Mortier, T., Javanmardi, A., Sale, Y., Hüllermeier, E., and Waegeman, W. Conformal prediction in hierarchical classification with constrained representation complexity. In

- The 29th International Conference on Artificial Intelligence and Statistics*, 2026.
- Qian, C., Xie, Z., Wang, Y., Liu, W., Zhu, K., Xia, H., Dang, Y., Du, Z., Chen, W., Yang, C., Liu, Z., and Sun, M. Scaling Large Language Model-based Multi-Agent Collaboration. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Razavian, N. and Sontag, D. Temporal convolutional neural networks for diagnosis from lab tests. *arXiv:1511.07938*, 2015.
- Romano, Y., Sesia, M., and Candès, E. Classification with valid and adaptive coverage. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- Shafer, G. and Vovk, V. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008.
- Vovk, V., Gammerman, A., and Shafer, G. *Algorithmic Learning in a Random World*. Springer, 2005.
- Wu, N., Gong, M., Shou, L., Liang, S., and Jiang, D. Large language models are diverse role-players for summarization evaluation. In *Natural Language Processing and Chinese Computing*, pp. 695–707. Springer Nature Switzerland, 2023. ISBN 978-3-031-44693-1.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., Zheng, C., Liu, D., Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin, H., Tang, J., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang, K., Yu, L., Deng, L., Li, M., Xue, M., Li, M., Zhang, P., Wang, P., Zhu, Q., Men, R., Gao, R., Liu, S., Luo, S., Li, T., Tang, T., Yin, W., Ren, X., Wang, X., Zhang, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Zhang, Y., Wan, Y., Liu, Y., Wang, Z., Cui, Z., Zhang, Z., Zhou, Z., and Qiu, Z. Qwen3 technical report, 2025.
- Yu, Y., Yao, Z., Li, H., Deng, Z., Jiang, Y., Cao, Y., Chen, Z., Suchow, J. W., Cui, Z., Liu, R., Xu, Z., Zhang, D., Subbalakshmi, K., Xiong, G., He, Y., Huang, J., Li, D., and Xie, Q. FinCon: A Synthesized LLM Multi-Agent System with Conceptual Verbal Reinforcement for Enhanced Financial Decision Making. In *Advances in Neural Information Processing Systems*, volume 37, pp. 137010–137045, 2024. doi: 10.52202/079017-4354.
- Yu, Y., Li, M., Xu, S., Fu, J., Hou, X., Lai, F., and Wang, B. CORRECT: CONDensed eRRor RECOgnition via knowledge Transfer in multi-agent systems. *arXiv:2509.24088*, 2025.
- Zhang, S., Yin, M., Zhang, J., Liu, J., Han, Z., Zhang, J., Li, B., Wang, C., Wang, H., Chen, Y., and Wu, Q. Which Agent Causes Task Failures and When? On Automated Failure Attribution of LLM Multi-Agent Systems. In *Forty-second International Conference on Machine Learning*, 2025.
- Zhu, C., Hong, S., Wu, J., Chawla, K., Tang, Y., Yin, Y., Wolfe, N., Babinsky, E., and Liu, D. RAFFLES: Reasoning-based attribution of faults for LLM systems. In *First Workshop on Multi-Turn Interactions in Large Language Models*, 2025.

## A. Theorems and Proofs

In this appendix we provide complete proofs of the theorems stated in the main text.

### A.1. Vanilla Conformal Prediction

The theorem showing valid coverage for VCP with variable number of classes is essentially the same as the standard result for conformal prediction, for instance as presented by (Angelopoulos & Bates, 2021), since the number of classes per datapoint does not affect exchangeability of the conformal scores. Although not novel, we present the theorem below for completeness.

**Theorem A.1.** *Suppose  $\{(x_i, y_i^*)\}_{i=1}^n$  and  $(x_{n+1}, y_{n+1}^*)$  are exchangeable, where the number of possible classes  $\ell_i$  varies per datapoint. Define  $\hat{q}$  as the  $\frac{\lceil (n+1)(1-\alpha) \rceil}{n}$  quantile of the scores  $\{S_i\}_{i=1}^n = \{S(x_i, y_i^*)\}_{i=1}^n$ . Then prediction sets constructed as  $C(x_{n+1}; \hat{q}) = \{y \in \mathcal{Y}(x_{n+1}) \mid S(x_{n+1}, y) \leq \hat{q}\}$  satisfy Equation (1).*

*Proof.* Since the fixed scoring function  $S$  is applied to each element of the exchangeable sequence  $(x_1, y_1^*), \dots, (x_{n+1}, y_{n+1}^*)$ , the resulting sequence of scalar scores  $S_1, \dots, S_{n+1}$  is also exchangeable. We assume for simplicity that the scores are non-degenerate, but ties can be handled by defining  $S$  to add a small amount of noise to its output. Let  $S_{(1)} \leq S_{(2)} \leq \dots \leq S_{(n)}$  be the order statistics of the calibration scores so that the empirical quantile  $\hat{q}$  can be defined as  $S_{(k)}$ , where  $k = \lceil (n+1)(1-\alpha) \rceil$  (assuming  $n \geq \frac{1}{\alpha} - 1$  to ensure  $k \leq n$ ). Given the definition of  $C(x_{n+1}; \hat{q})$ , the event  $y_{n+1}^* \in C(x_{n+1}; \hat{q})$  is equivalent to the event  $S_{n+1} \leq \hat{q}$ , which occurs if and only if  $S_{n+1}$  is among the  $k$  smallest scores overall. From exchangeability, all orderings are equally likely, so  $\mathbb{P}[S_{n+1} \leq \hat{q}] = \frac{k}{n+1}$ . Altogether we have

$$\mathbb{P}[y_{n+1}^* \in C(x_{n+1}; \hat{q})] = \frac{k}{n+1} = \frac{\lceil (n+1)(1-\alpha) \rceil}{n+1} \geq \frac{(n+1)(1-\alpha)}{n+1} = 1 - \alpha. \quad (28)$$

□

We note that the variable number of classes  $\ell_i$  per datapoint does not come up in the proof, and is only present in the assumption that such datapoints can still be exchangeable. This is certainly possible when we treat  $\ell_i$  as a (random) property of  $x$ . If desired, one could treat  $\ell_i$  as an explicit random variable and draw points as  $(x_i, \ell_i, y_i^*) \sim \mathbb{P}$ . Another way of viewing this extension would be to set a fixed  $\ell$  as some upper bound of possible sizes, such as  $\ell = \max\{\ell_i\}$ . Then, ordinary conformal prediction can be performed with a model  $f$  defined over  $\ell$  classes but giving  $f_j = 0$  when  $j > \ell_i$ , as long as we assume  $\ell_{n+1} \leq \ell$ .

### A.2. Left (Right) Filtration

Here we provide full formal proofs for each of the steps in Section 3.1.3. For ease of reference, we repeat key definitions and assumptions. A trajectory  $x$  is a sequence of a variable number  $\ell$  steps  $(c_1, \dots, c_\ell)$ , and the decisive error  $y^* = c_{j^*}$  is the step at index  $j^*$ . The scoring function  $g_{\text{LF}}$  scores subintervals  $c_{j:k} \subseteq x$  with the properties  $g_{\text{LF}}(\emptyset) = 0$  and  $g_{\text{LF}}(x) = \infty$ . The filtering function is defined as  $F_{\text{LF}}(x; q) := \arg \max_{c_{j:\ell} \in \mathcal{I}_{\text{LF}}} (|c_{j:\ell}| \mid g_{\text{LF}}(c_{j:\ell}) \leq q)$  which optimizes over suffixes  $\mathcal{I}_{\text{LF}} := \{c_{j:\ell}\}_{j=1}^{\ell+1}$  and satisfies  $F_{\text{LF}}(x; \infty) = x$ . Meanwhile, the conformal score function is defined as  $S_{\text{LF}}(x, y^*) := \inf (q \in \mathbb{R}^+ \mid y^* \in F_{\text{LF}}(x; q))$ . Both of these functions optimize over a set of steps, and it will be convenient to refer to those sets as follows:

$$\mathcal{V}_q := \{c_{j:\ell} \in \mathcal{I}_{\text{LF}} \mid g_{\text{LF}}(c_{j:\ell}) \leq q\}, \quad (29)$$

$$\mathcal{Q} := \{q \in \mathbb{R}^+ \mid y^* \in F_{\text{LF}}(x; q)\}. \quad (30)$$

**Lemma 3.1.** *For any  $x$  and thresholds  $0 \leq q_1 \leq q_2$ ,*

$$F_{\text{LF}}(x; q_1) \subseteq F_{\text{LF}}(x; q_2). \quad (8)$$

*Proof.*  $\mathcal{V}_q$  is the set of valid suffixes for  $F_{\text{LF}}(x; q)$ . For every  $c_{j:\ell} \in \mathcal{V}_{q_1}$  we have  $g_{\text{LF}}(c_{j:\ell}) \leq q_1 \leq q_2$ , and so  $c_{j:\ell} \in \mathcal{V}_{q_2}$ . Hence,  $\mathcal{V}_{q_1} \subseteq \mathcal{V}_{q_2}$ .  $F_{\text{LF}}(x; q)$  returns the longest suffix in  $\mathcal{V}_q$ , so we also have  $|F_{\text{LF}}(x; q_1)| \leq |F_{\text{LF}}(x; q_2)|$ . Since the suffixes  $c_{j:\ell} \in \mathcal{I}_{\text{LF}}$  are nested (i.e.  $j_2 < j_1 \iff c_{j_1:\ell} \subset c_{j_2:\ell}$ ), we also have  $|c_{j_1:\ell}| \leq |c_{j_2:\ell}| \implies c_{j_1:\ell} \subseteq c_{j_2:\ell}$ , which gives the desired nesting property of  $F_{\text{LF}}(x; q)$ . □

**Proposition 3.2.** When  $g_{LF}$  obeys monotonicity such that  $c_{b:c} \subseteq c_{a:d} \implies g_{LF}(c_{b:c}) \leq g_{LF}(c_{a:d})$ , then

$$S_{LF}(x, y^*) = g_{LF}(c_{j^*:\ell}). \quad (10)$$

*Proof.*  $\mathcal{Q}$  is the set of valid thresholds for  $S_{LF}(x, y^*)$ . Let  $q^* = g_{LF}(c_{j^*:\ell})$ , and note that  $c_{j^*:\ell} \in \mathcal{V}_{q^*}$ . Since  $|c_{j^*:\ell}| < |c_{j:\ell}|$  only when  $j < j^*$ , and  $y^* \in c_{j:\ell}$  for these cases, the longest suffix in  $\mathcal{V}_{q^*}$  contains  $y^*$ . Hence, we have that  $y^* \in F_{LF}(x; q^*)$ , and  $q^* \in \mathcal{Q}$ .

Now consider any  $q' < q^*$ . Notice that  $c_{j^*:\ell} \notin \mathcal{V}_{q'}$ . By monotonicity of  $g_{LF}$ , for all  $j < j^*$ ,  $g_{LF}(c_{j:\ell}) \geq q^* > q'$ . Therefore, the longest suffix in  $\mathcal{V}_{q'}$  is shorter than  $c_{j^*:\ell}$ , and does not contain  $y^*$ . Hence,  $q' \notin \mathcal{Q}$ .

Since  $q^* \in \mathcal{Q}$ , but all  $q' < q^*$  are not,  $q^*$  is the minimum of  $\mathcal{Q}$ , and  $S_{LF}(x, y^*) = q^* = g_{LF}(c_{j^*:\ell})$ .  $\square$

**Lemma 3.3.** For a fixed  $\hat{q}$ , we have

$$S_{LF}(x, y^*) \leq \hat{q} \iff y^* \in F_{LF}(x; \hat{q}). \quad (12)$$

*Proof.* To begin, we show that the infimum in  $S_{LF}(x, y^*)$  is always achieved. Note that the set of suffix scores  $g_{LF}(c_{j:\ell})$  is finite. Considering only the valid suffixes that contain  $y^*$ , let  $q_{\min} = \min\{g_{LF}(c_{j:\ell}) \mid j \leq j^*\}$ . Due to the nesting of suffixes, the condition  $y^* \in F_{LF}(x; q)$  is satisfied if and only if there exists a valid suffix with score  $\leq q$ , which occurs precisely when  $q_{\min} \leq q$ . Consequently,  $\mathcal{Q}$  is the closed interval  $[q_{\min}, \infty)$ , and the infimum in  $S_{LF}(x, y^*)$  is achieved at  $q_{\min}$ .

Hence we have a chain of equivalences:

$$S_{LF}(x, y^*) \leq \hat{q} \iff q_{\min} \leq \hat{q} \iff \hat{q} \in \mathcal{Q} \iff y^* \in F_{LF}(x; \hat{q}). \quad (31)$$

$\square$

### A.3. Two-Way Filtration

Next we cover the proofs of statements in Section 3.1.4. Continuing from the definitions and results above, including their counterparts for RF, for TWF we had  $F_{TWF}(x; q) := F_{LF}(x; q) \cap F_{RF}(x; q)$  (where  $F_{TWF}(x; 0) = \emptyset$  and  $F_{TWF}(x; \infty) = x$ ), and  $S_{TWF}(x, y^*) := \inf (q \in \mathbb{R}^+ \mid y^* \in F_{TWF}(x; q))$ .

**Lemma 3.5.** For any  $x$  and thresholds  $0 \leq q_1 \leq q_2$ ,

$$F_{TWF}(x; q_1) \subseteq F_{TWF}(x; q_2). \quad (18)$$

*Proof.* From Lemma 3.1 we have  $F_{LF}(x; q_1) \subseteq F_{LF}(x; q_2)$ , and the equivalent statement for  $F_{RF}$  holds. The intersection operation preserves set inclusion: for any sets  $A, B, C, D$ , if  $A \subseteq B$  and  $C \subseteq D$ , then  $A \cap C \subseteq B \cap D$ . Hence,

$$F_{TWF}(x; q_1) = F_{LF}(x; q_1) \cap F_{RF}(x; q_1) \subseteq F_{LF}(x; q_2) \cap F_{RF}(x; q_2) = F_{TWF}(x; q_2). \quad (32)$$

$\square$

**Proposition 3.6.**  $S_{TWF}(x, y^*)$  as defined in Equation (19) can be expressed as

$$S_{TWF}(x, y^*) = \max(S_{LF}(x, y^*), S_{RF}(x, y^*)). \quad (20)$$

When  $g_{LF}$  and  $g_{RF}$  obey the monotonicity condition in Equation (6),  $S_{TWF}$  further simplifies to

$$S_{TWF}(x, y^*) = \max(g_{LF}(c_{j^*:\ell}), g_{RF}(c_{1:j^*})). \quad (21)$$

*Proof.* By the definition of  $F_{TWF}$ ,  $y^* \in F_{TWF}(x; q)$  if and only if  $y^* \in F_{LF}(x; q)$  and  $y^* \in F_{RF}(x; q)$ . The condition  $y^* \in F_{LF}(x; q)$  holds if and only if  $q \geq S_{LF}(x, y^*)$  by the definition of  $F_{LF}(x; q)$  and its nesting property from Lemma 3.1. Similarly,  $y^* \in F_{RF}(x; q)$  if and only if  $q \geq S_{RF}(x, y^*)$ . Hence,  $y^* \in F_{TWF}(x; q)$  if and only if  $q \geq \max(S_{LF}(x, y^*), S_{RF}(x, y^*))$ . By taking the infimum for  $S_{TWF}(x, y^*)$ , we reach Equation (20).

From Proposition 3.2, when  $g_{LF}$  obeys its monotonicity condition, we have  $S_{LF}(x, y^*) = g_{LF}(c_{j^*:\ell})$ , and by a similar argument  $S_{RF}(x, y^*) = g_{RF}(c_{1:j^*})$ . Substituting these results into Equation (20) gives Equation (21).  $\square$

**Lemma 3.7.** *For a fixed  $\hat{q}$ , we have*

$$S_{TWF}(x, y^*) \leq \hat{q} \iff y^* \in F_{TWF}(x; \hat{q}). \quad (23)$$

*Proof.* From the definition of  $F_{TWF}(x, q)$  as an intersection, the event  $y^* \in F_{TWF}(x, \hat{q})$  occurs if and only if both  $y^* \in F_{LF}(x, \hat{q})$  and  $y^* \in F_{RF}(x, \hat{q})$  occur. We have already shown that  $y^* \in F_{LF}(x, \hat{q}) \iff S_{LF}(x, y^*) \leq \hat{q}$  (Lemma 3.3), and similarly for RF. Hence,

$$y^* \in F_{TWF}(x, \hat{q}) \iff S_{LF}(x, y^*) \leq \hat{q} \wedge S_{RF}(x, y^*) \leq \hat{q}. \quad (33)$$

Equivalently we write  $\max(S_{LF}(x, y^*), S_{RF}(x, y^*)) \leq \hat{q}$ , which is the same as  $S_{TWF}(x, y^*) \leq \hat{q}$  (Proposition 3.6).  $\square$

**Theorem A.2.** *Suppose  $\{(x_i, y_i^*)\}_{i=1}^n$  and  $(x_{n+1}, y_{n+1}^*)$  are exchangeable. Given scoring functions  $g_{LF}$  and  $g_{RF}$  as defined above, define the conformal score  $S_{TWF}(x_i, y_i^*)$  as in Equation (19). Let  $\hat{q}$  be the  $\frac{\lceil (n+1)(1-\alpha) \rceil}{n}$  quantile of conformal scores  $\{S_i\}_{i=1}^n = \{S_{TWF}(x_i, y_i^*)\}_{i=1}^n$ . Then prediction sets constructed as  $C_{TWF}(x_{n+1}; \hat{q}) = F_{TWF}(x_{n+1}; \hat{q})$  satisfy  $1 - \alpha \leq \mathbb{P}[y_{n+1}^* \in C_{TWF}(x_{n+1}; \hat{q})] < 1 - \alpha + \frac{1}{n+1}$ .*

*Proof.* The proof proceeds in the same way as Theorem 3.4, using Lemma 3.7 instead of Lemma 3.3.  $\square$

As a final point, we note that prediction sets constructed as  $C_{TWF}(x_{n+1}; \hat{q}) = F_{TWF}(x_{n+1}; \hat{q})$  may be empty due to the use of intersection. While this is not strictly a concern as coverage is valid marginally, it can be undesirable to predict empty sets. As a fallback, TWF can predict the single most likely step under  $g_{LF}$  or  $g_{RF}$  (which are likely to be the same function in practice).

## B. Additional Details of Conformal Algorithms

In Section 3.1.2 we described how to apply the CRSVP algorithm by [Mortier et al. \(2026\)](#), originally designed for hierarchical classification, to sequential data like agent trajectories. Here we provide a more detailed description of our implementation.

To recap, a binary tree  $\mathcal{T}$  is constructed from an agent trajectory  $x = (c_1, \dots, c_\ell)$  to have root node  $v_1 = [\ell]$ , and leaf nodes  $v_\ell, \dots, v_{2\ell-1}$  as individual steps  $c_1, \dots, c_\ell$ . For each calibration datapoint  $(x, y^*)$ , a scoring function  $g_{\text{CRSVP}}$  is used to find the single leaf  $\bar{v}$  most likely to contain the error  $y^*$ . If  $\bar{v}$  is itself the decisive error  $y^*$ , then the conformal score is returned as  $S_{\text{CRSVP}}(x, y^*) := g_{\text{CRSVP}}(\bar{v})$ . Otherwise, the tree is traversed upwards until the first node  $v^*$  is reached which does contain  $y^*$ . Since the number of steps in each node grows exponentially as the tree is traversed, simply using the score  $g_{\text{CRSVP}}(v^*)$  would greatly overcover the true labels. Hence, CRSVP interpolates between steps a random amount, similar to other conformal scoring functions ([Angelopoulos et al., 2021](#)). Let  $v^{*-1}$  be the node before  $v^*$  on the traversal path. Then the conformal score is set as

$$S_{\text{CRSVP}}(x, y^*) := g_{\text{CRSVP}}(v^*) - u \cdot (g_{\text{CRSVP}}(v^*) - g_{\text{CRSVP}}(v^{*-1})), \quad (34)$$

where  $u \sim \mathcal{U}(0, 1)$  is noise drawn uniformly from  $[0, 1]$ . Note that parent nodes are strict supersets of their children, which means prediction sets are nested when traversing from leaf to root ([Gupta et al., 2022](#)).

After conformal calibration to get the  $\frac{\lfloor (n+1)(\alpha) \rfloor}{n}$  quantile of the scores<sup>2</sup>, which we denote  $\hat{q}$ , for a test datapoint  $x_{n+1}$ , CRSVP again finds the most likely leaf node  $\bar{v}$  according to  $g_{\text{CRSVP}}$ . If that node's score is below the threshold,  $g_{\text{CRSVP}}(\bar{v}) < \hat{q}$ , CRSVP recursively traverses up the tree until the threshold is surpassed at  $\hat{v}$ . Drawing noise  $u$ , if  $g_{\text{CRSVP}}(\hat{v}) - u \cdot (g_{\text{CRSVP}}(\hat{v}) - g_{\text{CRSVP}}(\hat{v}^{-1})) \geq \hat{q}$ , then  $\hat{v}^{-1}$  is returned, otherwise  $\hat{v}$  itself is. By construction all prediction sets generated this way are contiguous subintervals of  $x_{n+1}$  and the lower bound on coverage is guaranteed (Equation (4)) ([Mortier et al., 2026](#)).

---

<sup>2</sup>Note that the quantile is different than VCP because we use conformity scores in CRSVP, and  $1 - g$  as a nonconformity score for VCP. These are merely conventions and can easily be modified.

Table 6. Fine-tuning configuration

Model	Dataset	Epochs	Batch Size	Precision	LoRA ( $r, \alpha$ )	LR
Qwen3-1.7B	GSM8k + MATH	18	32	8-bit	(2, 8)	$2 \times 10^{-6}$

## C. Implementation Details and Additional Results

### C.1. Additional Experiment Details

In this section we provide more detail about datasets, models, and experimental results.

#### C.1.1. SYNTHETIC DATASET GENERATION

As discussed in Section 4.1 we construct three variants of the original GSM8k dataset that explicitly control the location of the decisive error step along the execution trajectory. Errors are first generated uniformly by selecting one step in a trajectory, and injecting instructions to the agent to fail on that step in various ways (Kong et al., 2026; Cemri et al., 2025). We generate 2000 failed trajectories for GSM8k for each of the DyLAN and MACNET agent frameworks (and similarly for the MATH dataset). Then the trajectories are split by error position, in the first third, middle third, or final third to form the Right-, Mid-, and Left-Dense GSM8k datasets.

#### C.1.2. CONTEXT-ENGINEERED LLM IMPLEMENTATION

Prior work has shown that engineering the context of LLMs and prompting them with diverse role-specific instructions can improve evaluation robustness (Wu et al., 2023). Following this general idea, we employ multiple role prompts, inspired by ECHO (Banerjee et al., 2025), to obtain step-level scores from different perspectives. We first use the LLM to summarize why the agent failed the task using the overall trace, then pass this information to LLMs with four different roles to generate scores. The average of these scores is the final step score. The exact prompts used for the context-engineered LLM can be found in Appendix C.3.

#### C.1.3. FINE-TUNED SCORING FUNCTION

As described in Section 3.2, we fine-tuned an LLM to estimate step-level likelihoods of being the decisive error step.

**Training data** Due to limited data in Who&When (Zhang et al., 2025), we only perform fine-tuning on the MATH and GSM8k datasets with synthetically injected errors. As mentioned in Appendix C.1.1, we generated 8,000 trajectories in total across datasets and MAS architectures. Of these, 4,000 were used for model fine-tuning, equally split across the datasets, architectures, and error locations. We then used an 85%/15% split for training and validation to support hyperparameter selection. The 4,000 trajectories not used for training were held out entirely from fine-tuning and instead were used for conformal calibration and testing. Throughout this work, we report model performance on this unseen testing split.

**Model Training** We fine-tuned a Qwen3-1.7B (Yang et al., 2025) LLM to estimate step-wise error likelihoods in failed multi-agent trajectories. Since each trajectory had one step labeled as the decisive error (determined through error injection), we fine-tuned the model  $f_\theta$  on the  $\ell$ -way classification task to produce a scalar score  $f_\theta(c_{i,j})$  for each step  $j$  in datapoint  $i$ , interpreted as the likelihood that step  $c_{i,j}$  contains the decisive error. We write  $e_{i,j}$  as the binary indicator of whether step  $j$  in datapoint  $i$  is the decisive error. The model parameters  $\theta$  were optimized via binary cross-entropy:

$$\mathcal{L}(\theta) = - \sum_i \sum_{j=1}^{\ell_i} \left[ e_{i,j} \log \sigma(f_\theta(c_{i,j})) + (1 - e_{i,j}) \log(1 - \sigma(f_\theta(c_{i,j}))) \right], \tag{35}$$

where  $\sigma(\cdot)$  denotes the logistic sigmoid.

We provide more training set up details in Table 6, and report the performance before and fine-tuning the models in Table 7.

**Compute Resources** For fine-tuning the 1.7B parameter LLM we used an Nvidia GeForce RTX 5090. 20 epochs of training took roughly one day, and consumed less than 16 GB of GPU memory with batch size 32. All other calls to LLMs

Table 7. Performance of Qwen3-1.7B models before and after fine-tuning

Model	AUROC	AUPRC	Accuracy
Pretrained	0.505	0.369	0.509
Fine-tuned	0.762	0.382	0.731

used commercial APIs, and we discussed the number of calls needed in ???. Once scoring calls are made, performing conformal calibration is a trivial computational cost.

C.1.4. ROLLBACK EXPERIMENT IMPLEMENTATION

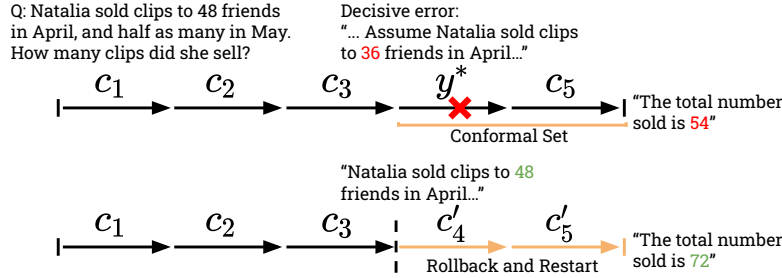


Figure 8. Example rollback scenario with text for the task, decisive error, final answer, and corrected versions after the rollback. We first generate a conformal set for the failed trajectory using the LF algorithm, roll back the state of the MAS to the first step in the prediction set, and restart the agent with instructions to avoid the same mistakes.

Figure 8 provides a concrete example of the rollback procedure. After identifying a conformal set for a failed trajectory using the LF algorithm, the system rolls back to the earliest step in the set and re-executes the task from that point on with additional context from the failed trace, allowing the MAS to correct the final outcome. Prompt details used during the restart are provided in Appendix C.3.

C.2. Additional Experimental Results

C.2.1. SCORING AGGREGATION METHODS

In Section 3.2 we discussed ways to tailor LLMs to the error attribution task and generate a step-wise error likelihood score. Step-wise scores then need to be aggregated into a set-wise score, ideally one that obeys monotonicity (Equation (6)). In Equation (25) we demonstrated the normalized sum function for aggregation which was used for experiments in the main text. Here we consider alternatives and show empirical results across them.

**Max Aggregation** Another way to monotonically aggregate step-level scores is by taking the maximum score over a set:

$$g_{\max}(c_{j:k}) = \max_{i \in \{j, \dots, k\}} \text{LLM}(c_i) + \lambda \frac{k - j}{\ell}. \tag{36}$$

Note that as opposed to Equation (25) where length normalization is included on the sum to account for trajectories of different lengths  $\ell$ , we do not normalize the max, since the single-step max tends to be similar across trajectories of different lengths. Instead we have included an optional length penalty term, weighted by the hyperparameter  $\lambda$ , to penalize sets that are longer than necessary to contain the highest scoring single step. Below, we also test adding a similar length penalty on the sum aggregation. Max aggregation should be useful when the LLM produces very peaked step-wise scores, as opposed to sum aggregation which can benefit when step-wise scores are more uniform.

**LogSumExp Aggregation** As a third alternative, we consider a family of monotonic aggregations that interpolate between the extremes of sum and max by using the LogSumExp function:

$$g_{\text{LSE}}(c_{j:k}) = \frac{1}{\beta} \text{LogSumExp}(\beta \ell \cdot [\text{LLM}(c_j), \dots, \text{LLM}(c_k)]) - \log \ell. \tag{37}$$

## Conformal Agent Error Attribution

Table 8. Removal Rate with the Right Filtering conformal algorithm across various aggregation functions on GSM8K-Left dataset.

Aggregation	Hypers	DyLAN			MACNET		
		Naive LLM	Role-prompted	Qwen3-1.7B	Naive LLM	Role-prompted	Qwen3-1.7B
Sum + Length Penalty	$\lambda = 0.01$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.01	0.62±0.02
	$\lambda = 0.02$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.01	0.62±0.02
	$\lambda = 0.05$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.01	0.62±0.02
	$\lambda = 0.1$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.01	0.62±0.01
Max + Length Penalty	$\lambda = 0.01$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.02	0.62±0.02
	$\lambda = 0.02$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.01	0.62±0.02
	$\lambda = 0.05$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.02	0.62±0.01
	$\lambda = 0.1$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.02	0.62±0.01
Normalized LogSumExp	$\beta = 0.1$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.02	0.63±0.02
	$\beta = 1$	0.64±0.02	0.63±0.02	0.65±0.01	0.59±0.01	0.59±0.02	0.63±0.02
	$\beta = 10$	0.63±0.02	0.63±0.02	0.64±0.02	0.60±0.01	0.59±0.01	0.62±0.01

This aggregation includes the inverse temperature hyperparameter  $\beta$ , where higher values behave more like the max function, while  $\beta$  close to zero behaves like the mean function (but is monotonic, unlike the true mean function). We also add length normalization via  $\ell$  to make scores more similar across trajectories with different lengths.

**Experimental Results** In Table 8, we compare Max, LogsumExp, and Sum aggregation with the optional length penalty using the RF conformal algorithms on the Left-Dense subset of GSM8k, across three evaluators (Naive LLM, Role-prompted, and the finetuned Qwen3-1.7B) settings and two agent architectures (DyLAN and MACNET). The primary result is that RF is largely insensitive to the details of the scoring function  $g$ , from the LLM design, to the aggregation function, or various hyperparameters. This means a simpler scoring function design is sufficient for use with our filtering algorithms.

### C.3. Prompts Used

**Naive LLM Prompts** We use following prompt for the naive LLM scoring function.

```

You are an AI assistant tasked with analyzing a segment of a multi-agent conversation.
Multiple agents are collaborating to address a user query, with the goal of resolving
the query through their collective dialogue.
Your primary task is to identify the location of the most critical mistake within the
provided segment.
The problem to address is as follows: {problem}
The final correct answer for the problem is: {answer}
The wrong answer given at the end of the conversation is: {wrong_answer}
Review the following conversation segment: {evaluate_content}
Cut from the entire history:
{chat_segment_content}
Based on your analysis, please respond with ONLY a single probability value between
0 and 1, representing the estimated probability that the critical error occurs within
this conversation segment. Do not include any explanation or additional text.

```

**Context Engineered Prompts** As mentioned earlier, along with the context engineered setting we use four role-based LLMs as part of the scoring function. Here, we provide the prompts used for each LLM.

## Conformal Agent Error Attribution

### Conservative:

- Attribute an error only when explicit contradiction or logical error is visible.
- Assign high confidence only if the evidence is direct and quotes are exact.

### Liberal:

- Consider potential upstream or downstream causes even if indirect.
- Assign moderate confidence (0.5--0.7) for plausible but not proven causes.

### Skeptical:

- Always include at least one alternative explanation and note missing evidence.
- Lower confidence if evidence is incomplete.

### Pattern:

- Focus on repeated reasoning or coordination issues.
- Identify patterns across multiple steps rather than isolated mistakes.

We force the output into the following schema:

### Output Schema

```
{
  "investigation_summary": "short summary",
  "primary_conclusion": {
    "agent": "Agent-X",
    "mistake_step": "int",
    "confidence": "float (0--1)",
    "evidence": ["quoted spans or step references"],
    "reason": "brief reason",
    "alternative_explanations": ["..."]
  }
}
```

The final prompt for each role is as follows:

```
Problem to solve: {problem}
Expected correct answer: {answer}
Conversation segment under review: {evaluate_content}
Full chat context: {chat_segment_content}
ROLE: {role} Analyst
{role_instructions}
TASK: Identify whether this segment likely contains the key reasoning error that led to the wrong final answer. Provide your analysis strictly in JSON format matching the schema below. Quote evidence spans exactly from the text and assign a numeric confidence between 0 and 1.
Confidence mapping: 0.0--0.2 = Very unlikely, 0.2--0.4 = Unlikely, 0.4--0.6 = Possible, 0.6--0.8 = Likely, 0.8--1.0 = Very likely.
Return ONLY valid JSON (no extra text): {schema}
```

**Conformal Rollback Prompts** The conformal rollback instructions provided to the MAS are as follow:

```
NOTE: The following conversation history (Node 0 to Node {cut_index}) contains wrong information. Please avoid making the same mistakes if you detect the mistake:
{wrong_conversation_content}
With the following mostly correct conversation history: {prev_context}
Solve this problem step by step: {task_query}
```