# LEC: Linear Expectation Constraints for False-Discovery Control in Selective Prediction and Routing Systems

Zhiyuan Wang[1]   Aniri[2]   Tianlong Chen[3]   Yue Zhang[4]   Heng Tao Shen[1 5]   Xiaoshuang Shi[1 *]   Kaidi Xu[6 *]

## Abstract

Large language models (LLMs) often generate unreliable answers, while heuristic uncertainty methods fail to fully distinguish correct from incorrect predictions, causing users to accept erroneous answers without statistical guarantees. We address this issue through the lens of false discovery rate (FDR) control, ensuring that among all accepted predictions, the proportion of errors does not exceed a target risk level. To achieve this in a principled way, we propose **LEC**, which reinterprets *selective prediction* as a constrained decision problem by enforcing a **L**inear **E**xpectation **C**onstraint over selection and error indicators. Then, we establish a finite-sample sufficient condition, which relies only on a held-out set of exchangeable calibration samples, to compute an FDR-constrained, coverage-maximizing threshold. Furthermore, we extend LEC to a two-model routing mechanism: given a prompt, if the current model's uncertainty exceeds its calibrated threshold, we delegate it to a stronger model, while maintaining a unified FDR guarantee. Evaluations on closed-ended and open-ended question-answering (QA) datasets show that LEC achieves tighter FDR control and substantially improves sample retention over prior methods. Moreover, the two-model routing mechanism achieves lower risk levels while accepting more correct samples than each individual model.

## 1. Introduction

Large language models (LLMs) are increasingly being integrated into various decision-making pipelines (Xiaolan et al., 2025; Brady et al., 2025), where it is crucial to evaluate the reliability of their outputs and determine when to trust them.

Uncertainty quantification (UQ) is a promising approach to estimate the uncertainty of LLM predictions, with the uncertainty score serving as an indicator of whether the model's output is likely to be incorrect (Wang et al., 2025d; Duan et al., 2024; 2025). In practice, when the model shows high uncertainty, its predictions should be clarified or rejected to prevent the propagation of incorrect information.

However, when the model generates hallucinations and exhibits high confidence in its erroneous predictions (Shorinwa et al., 2025; Atf et al., 2025), uncertainty scores derived from model logits or consistency measures may remain low, leading users to accept incorrect answers without task-specific risk guarantees (Angelopoulos et al., 2024). Split conformal prediction (SCP) can convert any heuristic uncertainty to a rigorous one (Angelopoulos & Bates, 2021; Campos et al., 2024a). Assuming data exchangeability, SCP produces prediction sets that include ground-truth answers with at least a user-defined probability. Nonetheless, set-valued predictions often contain unreliable candidates, leading to biased decision-making in downstream tasks (Wang et al., 2025a; Cresswell et al., 2025). In this paper, we center on point prediction with certain desirable finite-sample guarantees.

Although uncertainty estimates cannot perfectly separate correct from incorrect predictions, selective prediction allows us to enforce a specified risk level ($\alpha$): a prediction is accepted only when its associated uncertainty score falls below a calibrated threshold, ensuring that the false discovery rate (FDR) among accepted predictions remains below $\alpha$. Existing calibration methods based on confidence interval computation are overly conservative (Bates et al., 2021; Angelopoulos & Bates, 2021; Wang et al., 2025c; Jung et al., 2025), leading to excessive abstention of correct samples.

To achieve tighter FDR control, we introduce **LEC**, which reframes selective prediction not as an uncertainty-ranking problem, but as a decision problem governed by a statistical constraint. The central idea is to express FDR control as a constraint on the expectation of a linear functional involving two binary indicators: one capturing whether a prediction is selected and the other indicating whether it is incorrect. This formulation enables us to derive a finite-sample sufficient condition using calibration uncertainties and error labels that, if satisfied, guarantees FDR control for unseen test data.

Q: Who had a 70s No 1 hit with Kiss You All Over?

A: Ready for the World was not the first band to have a 70s No 1 hit with Kiss You All Over, it was actually Exile

low uncertainty($u \leq \lambda$), accept

Ground-truth: Exile

Q: Who was the first US-born winner of golf's British Open?

A: Tom Kite was not the first US-born winner of golf's British Open

high uncertainty($u > \lambda$), abstain

Ground-truth: Walter Hagen

Q: Della Street was whose secretary?

A: Della Street was the secretary of Sam Spade in Dashiell Hammett's novel The Maltese Falcon and its subsequent adaptations

high uncertainty($u^{(a)} > \lambda^{(a)}$), route

A: Della Street was the secretary of Perry Mason, the famous fictional criminal defense lawyer created by Erle Stanley Gardner

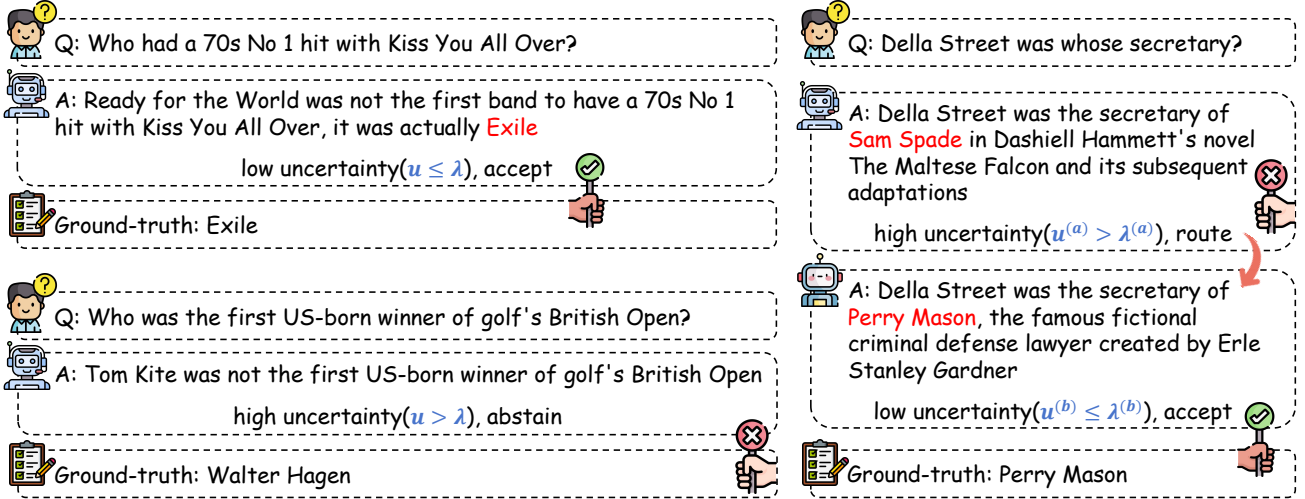low uncertainty($u^{(b)} \leq \lambda^{(b)}$), accept

Ground-truth: Perry Mason

Figure 1: Illustration of selective prediction in single-model and two-model routing systems.

Since this condition depends only on the empirical quantities from the calibration set, it yields a calibrated threshold that maximizes coverage subject to the FDR constraint.

We further extend LEC to two-model routing systems. For each input, if the current model's uncertainty falls below its calibrated threshold, its prediction is accepted; otherwise, the query is routed to the next model. If both models produce uncertainty above their thresholds, the system abstains. To maintain statistical reliability in this routing setting, we impose the same linear expectation constraint on the system-level selection and error indicators, enabling the joint calibration of model-specific thresholds under a unified FDR guarantee. Figure 1 illustrates examples of single-model selective prediction and two-model routing, where uncertainty scores serve as decision signals and calibrated thresholds function as acceptance gates. Beyond the two-model setting, we outline how LEC extends to general multi-model routing systems, offering a principled mechanism for unified FDR control across routing policies of arbitrary depth; exploring this broader setting empirically is left for future work.

We evaluate LEC across both closed-ended and open-ended QA datasets. In the single-model setting, LEC consistently keeps the test-time FDR below the specified risk level. Compared to baselines, LEC achieves tighter FDR control while accepting more admissible test samples (e.g., $+3\%$ on CommonsenseQA). In two-model routing settings, LEC achieves system-level FDR control and accepts more admissible predictions than either individual model. Furthermore, across different UQ methods, admission functions, and sampling sizes under black-box settings, LEC maintains statistical validity of FDR control while achieving higher power than the best baseline method. These results highlight the practical effectiveness and generality of LEC, motivating its integration into real-world uncertainty-aware agentic systems.

## 2. Related Work

**SCP in LLMs.** SCP provides statistical guarantees of coverage for correct answers (Campos et al., 2024b). It evaluates the nonconformity (or residual) between model prediction and ground-truth on a calibration set, and then computes a rigorously calibrated threshold, which is applied to construct prediction/conformal sets at test time. Under exchangeability (Angelopoulos et al., 2023), these sets contain admissible answers with at least a user-specified probability. However, previous research predominantly focuses on *set-valued predictions* (Quach et al., 2024; Kaur et al., 2024; Wang et al., 2024; 2025e;b;a), which are not inherently actionable due to unreliable candidates, and can cause disparate impact (Cresswell et al., 2024; 2025). Our work targets FDR control over accepted *point predictions*, rather than conformal coverage.

**FDR Control in Selective Prediction.** Early works employ confidence-based selection rules (e.g., entropy) (Wang et al., 2025d), but these heuristics often fail to perfectly distinguish between correct and incorrect predictions of LLMs (Zhang et al., 2024). Several frameworks grounded in significance testing (Jin & Candès, 2023; 2025) and confidence interval computation (Bates et al., 2021) have been introduced for FDR control on accepted answers. For instance, conformal alignment (Gui et al., 2024) and labeling (Huang et al., 2025) design conformal p-values and perform multiple hypothesis testing, ensuring that the test-time FDR remains below the significance level. To retain more correct answers and accelerate test-time inference, COIN (Wang et al., 2025c) derives a high-probability upper confidence bound for the system risk on calibration samples and computes a rigorously calibrated threshold for test-time selection, achieving PAC-style FDR control (Park et al., 2020; Angelopoulos et al., 2025). Moreover, selective evaluation (Jung et al., 2025) builds on the same paradigm to ensure that LLM-based assessments

achieve statistically guaranteed alignment with human judgments. In this work, we aim to explore a new paradigm for establishing tighter FDR control, rather than conservatively constraining a high-probability upper confidence bound.

## 3. Methods

### 3.1. Notations and Problem Formulation

#### 1) Single-Model Prediction with FDR Control

Let $\mathcal{G}^{(a)} : \mathcal{X} \to \mathcal{Y}$ denote a pretrained LLM that maps an input prompt to a textual output. For a given prompt $x \in \mathcal{X}$ with an unknown ground-truth answer $y^* \in \mathcal{Y}$, the model produces a prediction $\hat{y}^{(a)} = \mathcal{G}^{(a)}(x) \in \mathcal{Y}$. We quantify the model's uncertainty for $x$ as $u^{(a)} = \mathcal{U}(x; \mathcal{G}^{(a)})$, where $\mathcal{U}(\cdot)$ denotes a scalar uncertainty function. Intuitively, small $u^{(a)}$ indicates high reliability in $\hat{y}^{(a)}$. For a specified threshold $\lambda^{(a)}$, the prediction $\hat{y}^{(a)}$ is deemed admissible and accepted if $u^{(a)} \leq \lambda^{(a)}$. Let the admission function be

$$A(y^*, y) = \begin{cases} 1, & \text{if } y \in \mathcal{Y} \text{ is aligned with } y^*, \\ 0, & \text{otherwise.} \end{cases}$$

However, existing uncertainty approaches are imperfect and cannot perfectly separate correct from incorrect outputs (Liu et al., 2025). Thus, applying a fixed $\lambda^{(a)}$ at test time may admit some erroneous predictions. To mitigate this issue, our goal is to determine a statistically valid threshold $\hat{\lambda}^{(a)}$ that ensures the probability of accepting an incorrect prediction (i.e., the FDR) does not exceed a target risk level $\alpha$.

Formally, we define the selection indicator as

$$S^{(a)}\left(\lambda^{(a)}\right) = \mathbf{1}\left\{u^{(a)} \leq \lambda^{(a)}\right\},$$

and the error indicator as

$$err^{(a)} = \mathbf{1}\left\{A(y^*, \hat{y}^{(a)}) = 0\right\}.$$

Our objective is to estimate $\hat{\lambda}^{(a)}$ such that

$$\Pr\left(err^{(a)} = 1 \mid S^{(a)}(\hat{\lambda}^{(a)}) = 1\right) \leq \alpha, \ \alpha \in (0, 1). \quad (1)$$

#### 2) Two-Model Routing with FDR Control

Under a specific uncertainty function $\mathcal{U}(\cdot)$, the uncertainty scores of model $\mathcal{G}^{(a)}$ on test data may cluster too tightly in a low range, making it impossible to achieve small target risk levels. Moreover, when $\mathcal{G}^{(a)}$ has limited predictive ability (e.g., low accuracy), many challenging or critical prompts may be abstained from, leading to reduced system efficiency. To alleviate these issues, we develop a collaborative routing mechanism that dynamically delegates uncertain samples to another model with stronger accuracy or a more discriminative uncertainty profile, while controlling the system FDR.

Formally, we define the alternative LLM as $\mathcal{G}^{(b)} : \mathcal{X} \to \mathcal{Y}$. For a given prompt $x$, when the estimated uncertainty $u^{(a)}$ exceeds $\lambda^{(a)}$, we route the prompt to $\mathcal{G}^{(b)}$. We denote the prediction of $\mathcal{G}^{(b)}$ as $\hat{y}^{(b)} \in \mathcal{Y}$, with the corresponding uncertainty $u^{(b)} = \mathcal{U}(x; \mathcal{G}^{(b)})$. Similarly, if $u^{(b)}$ does not exceed the threshold $\lambda^{(b)}$ of model $\mathcal{G}^{(b)}$, we trust $\hat{y}^{(b)}$; otherwise, the two-model routing system abstains from the prompt $x$. We define the selection indicator for model $\mathcal{G}^{(b)}$ as

$$S^{(b)}\left(\lambda^{(a)}, \lambda^{(b)}\right) = \mathbf{1}\left\{u^{(a)} > \lambda^{(a)} \wedge u^{(b)} \leq \lambda^{(b)}\right\},$$

and the error indicator as

$$err^{(b)} = \mathbf{1}\left\{A(y^*, \hat{y}^{(b)}) = 0\right\}.$$

The two-model routing system $\mathcal{G}$ integrates $\mathcal{G}^{(a)}$ and $\mathcal{G}^{(b)}$, with the overall selection indicator

$$S\left(\lambda^{(a)}, \lambda^{(b)}\right) = S^{(a)}\left(\lambda^{(a)}\right) + S^{(b)}\left(\lambda^{(a)}, \lambda^{(b)}\right) \in \{0, 1\}.$$

and error indicator

$$\begin{aligned} err = \ &\mathbf{1}\{S^{(a)}(\lambda^{(a)}) = 1 \wedge err^{(a)} = 1\} \\ &+ \mathbf{1}\{S^{(b)}(\lambda^{(a)}, \lambda^{(b)}) = 1 \wedge err^{(b)} = 1\}. \end{aligned}$$

When $S(\lambda^{(a)}, \lambda^{(b)}) = 1$, the prediction from either $\mathcal{G}^{(a)}$ or $\mathcal{G}^{(b)}$ is accepted. We aim to jointly calibrate $(\lambda^{(a)}, \lambda^{(b)})$ and obtain statistically rigorous thresholds $(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})$ such that

$$\Pr\left(err = 1 \mid S\left(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}\right) = 1\right) \leq \alpha, \ \alpha \in (0, 1). \quad (2)$$

This guarantees that the overall two-model routing system achieves test-time FDR control while improving acceptance coverage and maintaining statistical reliability.

### 3.2. Threshold Calibration under Single-Model Settings

We begin by describing how to calibrate a statistically valid threshold $\hat{\lambda}^{(a)}$ for $\mathcal{G}^{(a)}$. Following the standard split calibration protocol (Angelopoulos et al., 2024), the dataset is partitioned into a calibration set and a test set. The threshold is learned solely from the calibration data for a user-specified risk level $\alpha$, and is then fixed during test-time evaluation.

**From FDR ratio to linear expectation.** For a fixed threshold $\lambda^{(a)}$, recall the selection and error indicators $S^{(a)}(\lambda^{(a)})$ and $err^{(a)}$. We then define the joint indicator

$$Z^{(a)}(\lambda^{(a)}) = S^{(a)}(\lambda^{(a)}) \cdot err^{(a)},$$

which equals $1$ only when we accept the prediction and the model errs. The FDR can then be formulated as

$$\begin{aligned} \mathrm{FDR}^{(a)}(\lambda^{(a)}) &= \Pr\left(err^{(a)} = 1 \mid S^{(a)}(\lambda^{(a)}) = 1\right) \\ &= \frac{\mathbb{E}[Z^{(a)}(\lambda^{(a)})]}{\mathbb{E}[S^{(a)}(\lambda^{(a)})]}. \end{aligned} \quad (3)$$

As long as $\mathbb{E}[S^{(a)}(\lambda^{(a)})] > 0$, $\mathrm{FDR}^{(a)}(\lambda^{(a)}) \le \alpha$ is equivalent to a constraint on the expectation of a linear functional of random variables (over selection and error indicators)

$$\mathbb{E}\big[Z^{(a)}(\lambda^{(a)}) - \alpha S^{(a)}(\lambda^{(a)})\big] \le 0. \qquad (4)$$

Intuitively, the random variable $Z^{(a)} - \alpha S^{(a)}$ measures *error count minus $\alpha$ times selection count* on a single example; if its expectation is non-positive, then the conditional error rate among accepted predictions (i.e., FDR) cannot exceed $\alpha$.

**Finite-sample sufficient condition.** To enforce the population constraint in Eq. (4) using only the calibration data, we then derive a finite-sample condition. Let the calibration set be $\mathcal{D}_{\mathrm{cal}} = \{(u_i^{(a)}, err_i^{(a)})\}_{i=1}^n$, with $\{S_i^{(a)}\}_{i=1}^n$, and let $u_{(1)}^{(a)} \le \cdots \le u_{(n)}^{(a)}$ denote the calibration uncertainty scores sorted in ascending order, with corresponding error indicators $err_{(j)}^{(a)}$. For any candidate threshold $\lambda^{(a)}$, we define

$$k^{(a)}(\lambda^{(a)}) = \#\{i : S_i^{(a)}(\lambda^{(a)}) = 1\} = \#\{i : u_i^{(a)} \le \lambda^{(a)}\}$$

as the number of calibration data points that would be accepted at the risk level of $\lambda^{(a)}$. Then, a standard "+1" correction, which ensures validity under exchangeability and avoids degeneracy when all accepted examples happen to be correct, yields the following finite-sample sufficient condition (see Appendix A.1 for complete derivation):

$$\sum_{j=1}^{k^{(a)}(\lambda^{(a)})} \big(err_{(j)}^{(a)} - \alpha\big) \le -1. \qquad (5)$$

We then define the feasible set of thresholds at level $\alpha$ as

$$\Lambda_\alpha^{(a)} = \Big\{\lambda^{(a)} : \sum_{j=1}^{k^{(a)}(\lambda^{(a)})} \big(err_{(j)}^{(a)} - \alpha\big) \le -1 \Big\}. \qquad (6)$$

**Calibrated Coverage-Maximizing Threshold.** Among all thresholds in $\Lambda_\alpha^{(a)}$, we choose the largest feasible one to maximize the acceptance coverage:

$$\hat{\lambda}^{(a)} = \sup \Lambda_\alpha^{(a)}$$
$$= \sup \Big\{\lambda^{(a)} : \sum_{j=1}^{k^{(a)}(\lambda^{(a)})} \big(err_{(j)}^{(a)} - \alpha\big) \le -1 \Big\}. \qquad (7)$$

If $\Lambda_\alpha^{(a)}$ is empty, we declare the target risk level $\alpha$ infeasible for $\mathcal{G}^{(a)}$ and abstain on all samples at this level.

The following theorem states that applying $\hat{\lambda}^{(a)}$ controls the test-time FDR at the risk level of $\alpha$.

**Theorem 3.1** (Single-model FDR control)**.** *Assume that calibration and test examples are exchangeable (Angelopoulos*

*et al., 2023). Let $\hat{\lambda}^{(a)}$ be defined by Eq. (7) using $\mathcal{D}_{\mathrm{cal}}$. Then, for a new test sample $(x_{n+1}, y_{n+1}^*)$ with $(u_{n+1}^{(a)}, err_{n+1}^{(a)})$,*

$$\Pr\big(err_{n+1}^{(a)} = 1 \mid u_{n+1}^{(a)} \le \hat{\lambda}^{(a)}\big) \le \alpha,$$

*where the probability is taken over the joint randomness of the calibration set and the test sample (marginal guarantee).*

A complete proof of Theorem 3.1 is given in Appendix A.1. At test time, for a new prompt $x_{n+1}$, we obtain the model prediction $\hat{y}_{n+1}^{(a)}$ with uncertainty $u_{n+1}^{(a)}$. We accept $\hat{y}_{n+1}^{(a)}$ if and only if $u_{n+1}^{(a)} \le \hat{\lambda}^{(a)}$; otherwise, we abstain.

### 3.3. Threshold Calibration under Two-Model Settings

We now extend the above calibration procedure to the two-model routing system $\mathcal{G}$. For each example $i$, we observe uncertainties $(u_i^{(a)}, u_i^{(b)})$ and error indicators $(err_i^{(a)}, err_i^{(b)})$. Given thresholds $(\lambda^{(a)}, \lambda^{(b)})$, routing is defined by the selection indicators $S_i^{(a)}(\lambda^{(a)})$ and $S_i^{(b)}(\lambda^{(a)}, \lambda^{(b)})$. The system-level selection indicator is $S_i(\lambda^{(a)}, \lambda^{(b)}) = S_i^{(a)}(\lambda^{(a)}) + S_i^{(b)}(\lambda^{(a)}, \lambda^{(b)}) \in \{0, 1\}$, and the error indicator is $err_i = \mathbf{1}\{S_i^{(a)}(\lambda^{(a)}) = 1, err_i^{(a)} = 1\} + \mathbf{1}\{S_i^{(b)}(\lambda^{(a)}, \lambda^{(b)}) = 1, err_i^{(b)} = 1\}$, which remains binary because routing selects at most one output. We also define the joint indicator

$$Z_i(\lambda^{(a)}, \lambda^{(b)}) = S_i(\lambda^{(a)}, \lambda^{(b)}) \cdot err_i$$
$$= S_i^{(a)}(\lambda^{(a)}) \cdot err_i^{(a)} + S_i^{(b)}(\lambda^{(a)}, \lambda^{(b)}) \cdot err_i^{(b)}.$$

**From system-level FDR to a linear constraint.** The system FDR under thresholds $(\lambda^{(a)}, \lambda^{(b)})$ is

$$\mathrm{FDR}(\lambda^{(a)}, \lambda^{(b)}) = \frac{\mathbb{E}[Z(\lambda^{(a)}, \lambda^{(b)})]}{\mathbb{E}[S(\lambda^{(a)}, \lambda^{(b)})]}.$$

Whenever $\mathbb{E}[S(\lambda^{(a)}, \lambda^{(b)})] > 0$, $\mathrm{FDR}(\lambda^{(a)}, \lambda^{(b)}) \le \alpha$ is also equivalent to a linear inequality

$$\mathbb{E}\big[Z(\lambda^{(a)}, \lambda^{(b)}) - \alpha S(\lambda^{(a)}, \lambda^{(b)})\big] \le 0. \qquad (8)$$

This condition generalizes the single-model constraint to the routing system and again captures the difference between the system-level error count and $\alpha$-fraction of accepted samples.

**Finite-sample sufficient condition.** To enforce Eq. (8) from calibration points, we also construct an empirical sufficient condition. Let $\mathcal{D}_{\mathrm{cal}}^{\mathrm{sys}} = \{(u_i^{(a)}, u_i^{(b)}, err_i^{(a)}, err_i^{(b)})\}_{i=1}^n$ denote the calibration set for the two-model routing system. By applying the same "+1 smoothing" argument to the pair $(Z_i, S_i)$, we establish the finite-sample sufficient condition

$$\sum_{i=1}^n \Big(Z_i(\lambda^{(a)}, \lambda^{(b)}) - \alpha S_i(\lambda^{(a)}, \lambda^{(b)})\Big) \le -1. \qquad (9)$$
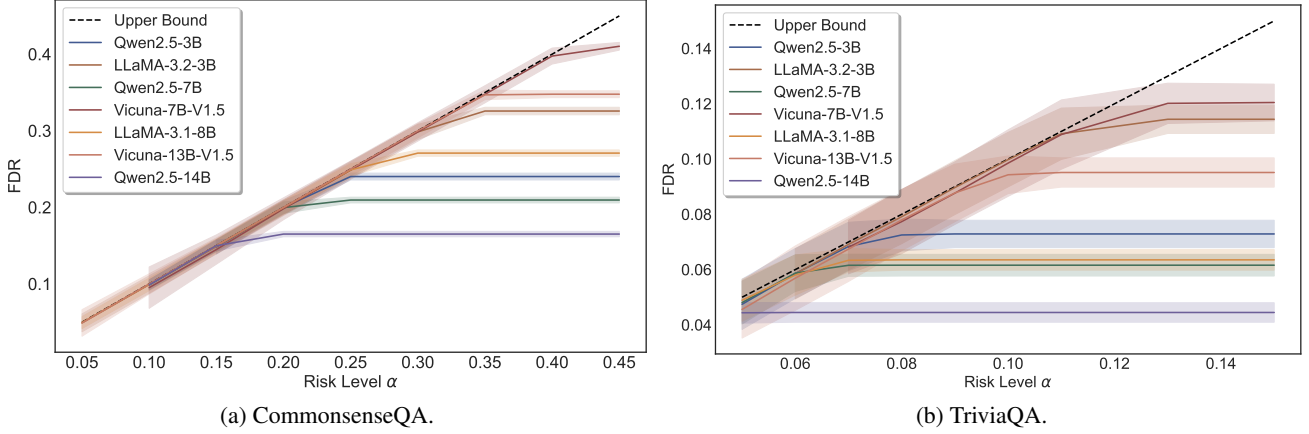
(a) CommonsenseQA.

(b) TriviaQA.

Figure 2: FDR control at various $\alpha$ on both the CommonsenseQA and TriviaQA datasets with seven LLMs (mean±std).

We then obtain the feasible set of two-model threshold pairs

$$\Lambda_\alpha^{(a,b)} = \Big\{ (\lambda^{(a)}, \lambda^{(b)}) :$$
$$\sum_{i=1}^{n} \Big( Z_i(\lambda^{(a)}, \lambda^{(b)}) - \alpha S_i(\lambda^{(a)}, \lambda^{(b)}) \Big) \leq -1 \Big\}. \quad (10)$$

**Calibrated coverage-maximizing thresholds.** Among all pairs $(\lambda^{(a)}, \lambda^{(b)}) \in \Lambda_\alpha^{(a,b)}$, we choose those that maximize the empirical acceptance coverage of the routing system:

$$(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) \in \operatorname*{argmax}_{(\lambda^{(a)}, \lambda^{(b)}) \in \Lambda_\alpha^{(a,b)}} \frac{1}{n} \sum_{i=1}^{n} S_i(\lambda^{(a)}, \lambda^{(b)}). \quad (11)$$

If $\Lambda_\alpha^{(a,b)}$ is empty, the target risk level $\alpha$ is deemed infeasible for the two-model routing system, and the system abstains on all inputs at this level. The following theorem states that the threshold pair $(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})$ controls the test-time FDR.

**Theorem 3.2** (FDR control for the two-model routing system). *Assume calibration and test examples are exchangeable. Let $(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})$ be any solution of Eq. (11). Then the two-model routing system satisfies*

$$\Pr\big(err_{n+1} = 1 \mid S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) = 1\big) \leq \alpha,$$

*where the probability is taken over the joint randomness of calibration and test samples.*

See a proof of Theorem 3.2 in Appendix A.2. At test time, each given prompt $x_{n+1}$ is processed as follows: we accept $\hat{y}_{n+1}^{(a)}$ via $\mathcal{G}^{(a)}$ if $u_{n+1}^{(a)} \leq \hat{\lambda}^{(a)}$; otherwise we route the prompt to $\mathcal{G}^{(b)}$ and accept $\hat{y}_{n+1}^{(b)}$ if $u_{n+1}^{(b)} \leq \hat{\lambda}^{(b)}$. If neither condition is satisfied, the system abstains. Our analysis highlights that FDR control is preserved as long as the routing policy is deterministic and each example is routed to at most

one model. The statistical guarantees arise from the linear decomposition, rather than any model-specific assumptions.

The above two-model calibration can readily be extended to routing systems with more than two models. In Appendix B, we provide details on how to calibrate thresholds for each model in a general multi-model routing system to achieve system-level selective prediction with FDR control.

## 4. Experiments

### 4.1. Experimental Settings

**Datasets and Models.** We evaluate LEC on the CommonsenseQA (closed-ended) (Talmor et al., 2019) and TriviaQA (open-ended) (Joshi et al., 2017) datasets using seven LLMs, including LLaMA (Touvron et al., 2023), Qwen (Bai et al., 2023), and Vicuna (Zheng et al., 2023) families.

**Uncertainty Methods.** In closed-ended QA tasks, we estimate uncertainty scores by computing the predictive entropy (PE) (Kadavath et al., 2022). We utilize the softmax output of model logits by default. If model internal information is inaccessible, we sample multiple answers and use sampling frequency to approximate the generative probability (Wang et al., 2025e). In open-ended QA tasks, we focus on black-box scenarios and compute the semantic entropy (SE) (Farquhar et al., 2024) by default. Moreover, we also consider the sum of eigenvalues of the graph laplacian (EigV), degree matrix (Deg), and eccentricity (Ecc) (Lin et al., 2024).

**Alignment Criteria.** We use sentence similarity (Reimers & Gurevych, 2019b) with a 0.6 threshold to evaluate whether the model's answer is aligned with the ground truth in the admission function $A$ by default. Following prior work (Lin et al., 2024), we also consider bi-entailment (Kuhn et al., 2023). Moreover, if the utilized UQ method involves semantic clustering, we use a consistent criterion to evaluate whether the sampled answers are semantically equivalent.
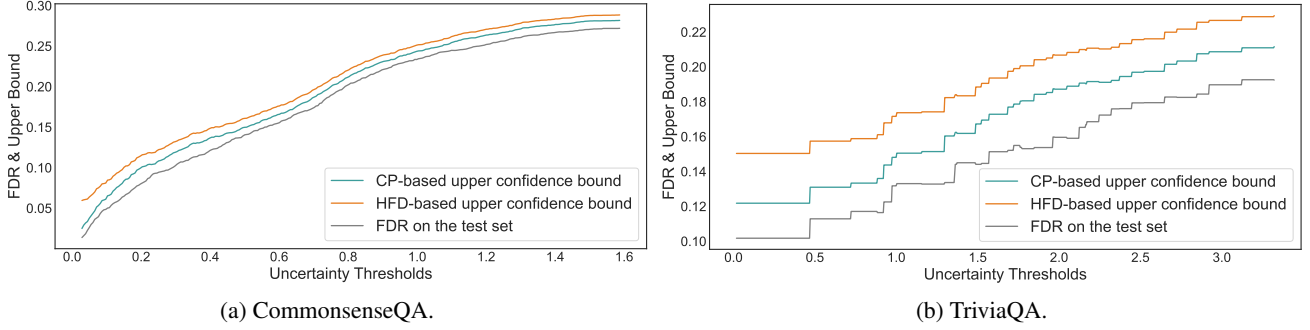
(a) CommonsenseQA.

(b) TriviaQA.

Figure 3: Upper confidence bound vs. Test-time FDR at various uncertainty thresholds. In (a), we use LLaMA-3.1-8B, with white-box PE as the uncertainty measure; In (b), we use Qwen2.5-14B, with black-box SE as the uncertainty measure.
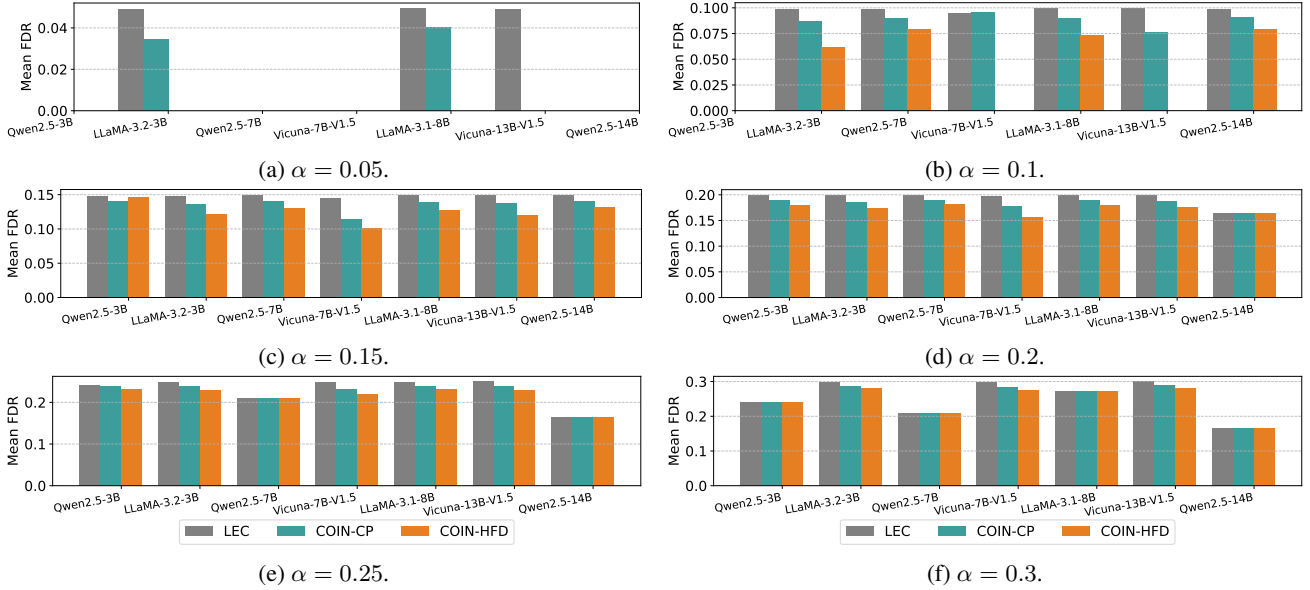


(a) $\alpha = 0.05$.

(b) $\alpha = 0.1$.

(c) $\alpha = 0.15$.

(d) $\alpha = 0.2$.

(e) $\alpha = 0.25$.

(f) $\alpha = 0.3$.

Figure 4: Comparison of test-time FDR on the CommonsenseQA dataset across seven LLMs (mean).

**Evaluation Metrics.** Following previous work (Jung et al., 2025; Wang et al., 2025c), we evaluate the statistical validity of LEC by verifying that test-time FDR remains below the user-specified risk level. We further assess its power, defined as the proportion of admissible test samples it accepts among all admissible samples while satisfying FDR control.

**Baselines.** Existing confidence interval-based approaches can achieve FDR control (Bates et al., 2021). Taking COIN as an example (Wang et al., 2025c), it computes the FDR on calibration data at an initial threshold and derives the $(1 - \delta)$ upper confidence bound of system risk. By adjusting the threshold so that this bound falls below the risk level, COIN controls test-time FDR. We compare LEC with COIN using two confidence intervals—Clopper-Pearson (COIN-CP) and Hoeffding's inequality (COIN-HFD)—with $\delta$ set to 0.05.

**Hyperparameters.** Following previous work (Wang et al., 2025c), we employ beam search (`num_beams=5`) to obtain

the most likely generation as the model output. By default, for open-domain QA, we sample 10 answers per input for UQ. In addition, we fix the calibration-test split ratio to 0.5.

We provide the details of additional experimental settings in Appendix C. Following prior research (Quach et al., 2024), we randomly split the calibration and test samples 100 times and report the mean and standard deviation (mean±std). We annotate this information alongside the subsequent results.

## 4.2. Evaluations in Single-Model Selective Prediction

**Statistical Validity.** We first evaluate LEC in single-model settings and perform selective prediction on the test set using the threshold derived from the calibration data—accepting samples whose PE-based uncertainty scores fall below the calibrated threshold for a specified risk level $\alpha$. As shown in Figure 2, across both datasets with seven LLMs, LEC consistently achieves strict FDR control: the test-time FDR (mean

Table 1: Power comparison on the CommonsenseQA dataset across seven LLMs (mean).

| LLMs | Methods / $\alpha$ | 0.05 | 0.1 | 0.15 | 0.2 | 0.25 | 0.3 | 0.35 | 0.4 | 0.45 |
|---|---|---|---|---|---|---|---|---|---|---|
| Qwen2.5-3B | COIN-CP | - | - | 0.8033 | 0.9087 | 0.9966 | $0.9999_{(9)}$ | $0.9999_{(9)}$ | $0.9999_{(9)}$ | $0.9999_{(9)}$ |
| | COIN-HFD | - | - | 0.7964 | 0.8960 | 0.9853 | $0.9999_{(9)}$ | $0.9999_{(9)}$ | $0.9999_{(9)}$ | $0.9999_{(9)}$ |
| | LEC | - | - | **0.8223** | **0.9234** | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(9)}$ | $\mathbf{0.9999}_{(9)}$ | $\mathbf{0.9999}_{(9)}$ | $\mathbf{0.9999}_{(9)}$ |
| LLaMA-3.2-3B | COIN-CP | 0.1191 | 0.3744 | 0.5704 | 0.7032 | 0.8260 | 0.9343 | $0.9999_{(9)}$ | $0.9999_{(9)}$ | $0.9999_{(9)}$ |
| | COIN-HFD | - | 0.2695 | 0.5190 | 0.6702 | 0.8054 | 0.9188 | 0.9998 | $0.9999_{(9)}$ | $0.9999_{(9)}$ |
| | LEC | **0.1959** | **0.4110** | **0.5992** | **0.7290** | **0.8472** | **0.9518** | $\mathbf{0.9999}_{(9)}$ | $\mathbf{0.9999}_{(9)}$ | $\mathbf{0.9999}_{(9)}$ |
| Qwen2.5-7B | COIN-CP | - | 0.7620 | 0.8805 | 0.9689 | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ |
| | COIN-HFD | - | 0.7189 | 0.8595 | 0.9538 | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ |
| | LEC | - | **0.7805** | **0.8950** | **0.9801** | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(7)}$ |
| Vicuna-7B-V1.5 | COIN-CP | - | 0.0338 | 0.0974 | 0.3391 | 0.5016 | 0.6650 | 0.8231 | 0.9619 | $0.9999_{(8)}$ |
| | COIN-HFD | - | - | 0.0365 | 0.2755 | 0.4655 | 0.6280 | 0.7936 | 0.9502 | $0.9999_{(8)}$ |
| | LEC | - | **0.0476** | **0.2050** | **0.3818** | **0.5364** | **0.6995** | **0.8613** | **0.9785** | $\mathbf{0.9999}_{(8)}$ |
| LLaMA-3.1-8B | COIN-CP | 0.3841 | 0.6036 | 0.7671 | 0.8725 | 0.9532 | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ |
| | COIN-HFD | - | 0.5507 | 0.7398 | 0.8591 | 0.9408 | $0.9999_{(6)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ | $0.9999_{(7)}$ |
| | LEC | **0.4146** | **0.6363** | **0.7866** | **0.8842** | **0.9679** | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(7)}$ | $\mathbf{0.9999}_{(7)}$ |
| Vicuna-13B-V1.5 | COIN-CP | - | 0.1384 | 0.4683 | 0.6460 | 0.7933 | 0.9128 | 0.9899 | $0.9999_{(9)}$ | $0.9999_{(9)}$ |
| | COIN-HFD | - | - | 0.3902 | 0.6209 | 0.7645 | 0.8990 | 0.9824 | $0.9999_{(9)}$ | $0.9999_{(9)}$ |
| | LEC | **0.0468** | **0.2498** | **0.5105** | **0.6662** | **0.8211** | **0.9284** | **0.9987** | $\mathbf{0.9999}_{(9)}$ | $\mathbf{0.9999}_{(9)}$ |
| Qwen2.5-14B | COIN-CP | - | 0.8677 | 0.9651 | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ |
| | COIN-HFD | - | 0.8347 | 0.9503 | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ |
| | LEC | - | **0.8829** | **0.9766** | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ |

over 100 splits) always stays below the upper bound. For instance, on CommonsenseQA at a risk level of 0.5, LLaMA-3.1-8B obtains a test-time FDR of 0.0492. Note that, given the marginal guarantee offered by LEC, the shaded regions that slightly exceed the upper bound do not compromise its statistical validity, as minor fluctuations are expected due to finite-sample variability (Angelopoulos & Bates, 2021).

We can observe that at certain risk levels, some models, such as Qwen2.5-3B, fail to search for a feasible set that satisfies the required constraint in Eq. (10) on the calibration data. This occurs because a non-negligible fraction of incorrect samples have uncertainty scores concentrated in a relatively low range. We report the full uncertainty and correctness distributions for all models in Appendix D.

**Tighter FDR Control.** As shown in Figure 3, within a single split, we construct two types of upper confidence bounds on the calibration data. Although the Clopper–Pearson (CP) interval is exact, the resulting test-time FDR still lies far below its $(1 - \delta)$ upper bound. The bound derived from Hoeffding's inequality (HFD) is more conservative, yielding an even larger gap. This indicates that calibrating thresholds by enforcing the upper confidence bound to fall below $\alpha$ inherently leads to loose control. Figure 4 further compares LEC with two COIN variants under identical settings on CommonsenseQA. At each risk level, LEC consistently achieves higher acceptance rates while keeping FDR strictly under $\alpha$, outperforming both COIN-CP and COIN-HFD. In contrast to the conservative nature of COIN, LEC provides

tighter calibration and supports valid FDR control at lower risk levels across a broader set of LLMs. For instance, with Vicuna-13B-V1.5, LEC maintains strict FDR control at a risk level of 0.05, whereas both COIN variants fail to reduce the upper confidence bound below $\alpha$, for any threshold choice, and thus cannot provide a valid guarantee.

**Higher Power.** As presented in Table 1, on the CommonsenseQA dataset, LEC consistently accepts a larger proportion of correct test samples compared with both versions of COIN, demonstrating substantial improvements in power and efficiency across all seven models. For instance, with the LLaMA-3.2-3B model at a risk level of 0.05, the power achieved by LEC exceeds that of COIN-CP by nearly 8%.

The power comparison under black-box settings on CommonsenseQA, as well as additional evaluations on TriviaQA utilizing different UQ methods, alignment criteria, and sampling sizes, are provided in Appendix D.

### 4.3. Evaluations in Two-Model Routing

Unlike the single-model setting, where we select the largest threshold that satisfies the finite-sample sufficient condition, as defined in Eq. (7), in the two-model routing setting, we instead search for the threshold pair that accepts the most samples on the calibration set, as defined in Eq. (11). As demonstrated in Figure 5, we evaluate the routing mechanism between LLaMA-3.1-8B and three model sizes from the Qwen2.5 series on CommonsenseQA, achieving rigorous FDR control across various user-specified risk levels.

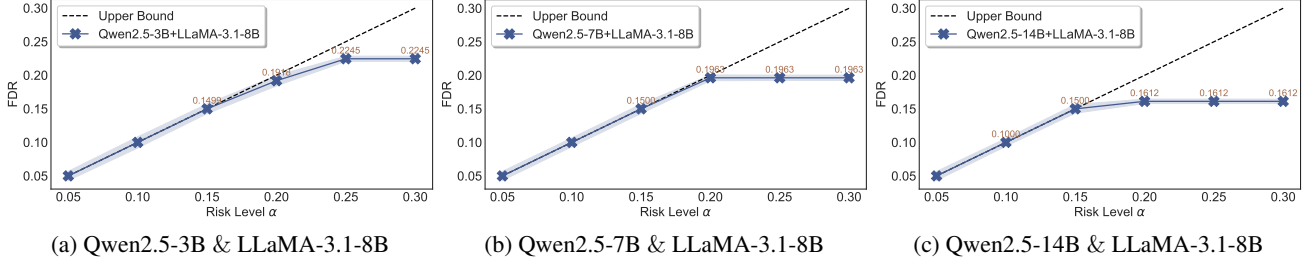| (a) Qwen2.5-3B & LLaMA-3.1-8B | (b) Qwen2.5-7B & LLaMA-3.1-8B | (c) Qwen2.5-14B & LLaMA-3.1-8B |

Figure 5: FDR control of two LLMs routing at various risk levels on the CommonsenseQA dataset (mean±std).

Table 2: Number of accepted correct answers at test time on the CommonsenseQA dataset (mean).

| LLMs | Best Cover / $\alpha$ | 0.05 | 0.1 | 0.15 | 0.2 | 0.25 | 0.3 |
|---|---|---|---|---|---|---|---|
| LLaMA-3.1-8B | - | 1776.13 | 2545.73 | 3080.34 | 3424.27 | 3714.86 | 3727.15 |
| Qwen2.5-3B | - | - | - | 3337.75 | 3703.71 | 3871.46 | 3871.46 |
| Qwen2.5-3B + LLaMA-3.1-8B | ⊗ | 1776.13 | 2545.73 | 3416.40 | 3728.71 | 3871.76 | 3871.76 |
| Qwen2.5-3B + LLaMA-3.1-8B | ⊘ | **1776.13** | **2545.73** | **3451.81** | **3791.35** | **3877.34** | **3877.34** |
| Qwen2.5-7B | - | - | 3219.77 | 3669.88 | 4003.52 | 4008.00 | 4008.00 |
| Qwen2.5-7B + LLaMA-3.1-8B | ⊗ | 1776.13 | 3288.07 | 3717.71 | 4005.99 | 4008.54 | 4008.54 |
| Qwen2.5-7B + LLaMA-3.1-8B | ⊘ | **1776.13** | **3315.74** | **3750.18** | **4016.04** | **4018.17** | **4018.17** |
| Qwen2.5-14B | - | - | 3721.06 | 4118.92 | 4184.51 | 4184.51 | 4184.51 |
| Qwen2.5-14B + LLaMA-3.1-8B | ⊗ | 1776.18 | 3733.56 | 4123.97 | 4184.82 | 4184.82 | 4184.82 |
| Qwen2.5-14B + LLaMA-3.1-8B | ⊘ | **1776.18** | **3738.21** | **4134.55** | **4193.76** | **4193.76** | **4193.76** |

Moreover, as illustrated in Figure 2a, although none of the three Qwen2.5 models can offer FDR guarantees at a risk level of 0.05 on their own, each of them attains valid test-time FDR control once paired with LLaMA-3.1-8B.

In the two-model routing setting, the total number of correct samples attributable to each model differs at test time, and routing further reallocates test samples between them. This makes it difficult to compute a comparable power metric directly. We evaluate the number of correct samples accepted at test time, since the total number of test samples is fixed. As shown in Table 2, two-model routing increases the number of correct samples accepted at test time. For example, when the risk level is set to 0.15, the routing combination of Qwen2.5-3B and LLaMA-3.1-8B accepts an average of 3451.81 samples over 100 splits, over 100 more than utilizing Qwen2.5-3B alone and more than 300 above using LLaMA-3.1-8B alone. Moreover, we compare the threshold pair obtained by directly performing grid search on the calibration set with the threshold pair that maximizes the number of accepted calibration samples. Although both satisfy the FDR constraint, the latter leads to a larger number of correct samples accepted at test time. For instance, when the risk level is set to 0.1, the routing combination of Qwen2.5-7B and LLaMA-3.1-8B accepts nearly 30 more samples on average at test time after optimizing for the threshold pair that maximizes the number of accepted calibration samples.

We also report the results of sample acceptance at various risk levels in the routing setting at test time in Appendix D.

## 5. Conclusion

In this paper, we presented LEC, a principled framework that reframes selective prediction as a constrained decision problem via a linear expectation constraint on selection and error indicators. This formulation yields a finite-sample sufficient condition for rigorous FDR control using only exchangeable calibration samples, and it generalizes naturally from single-model prediction to two-model and general multi-model routing systems. Empirically, LEC delivers strict test-time FDR control on both closed-ended and open-domain QA tasks, achieving substantially tighter calibration than prior confidence interval-based approaches such as COIN-CP and COIN-HFD. Across seven LLMs, LEC consistently retains a higher fraction of correct predictions while maintaining the FDR guarantees. The routing extension further enables weaker models to operate at lower feasible risk levels and increases correct acceptances beyond what any individual model can achieve. Looking ahead, LEC provides a general foundation for risk-controlled deployment of LLMs. Future work may extend this framework to task-specific risk metrics and richer routing architectures, supporting reliable decision-making in increasingly complex agentic systems.

# References

Angelopoulos, A. N. and Bates, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.

Angelopoulos, A. N., Bates, S., et al. Conformal prediction: A gentle introduction. *Foundations and trends® in machine learning*, 2023.

Angelopoulos, A. N., Bates, S., Fisch, A., Lei, L., and Schuster, T. Conformal risk control. In *The Twelfth International Conference on Learning Representations*, 2024.

Angelopoulos, A. N., Bates, S., Candès, E. J., Jordan, M. I., and Lei, L. Learn then test: Calibrating predictive algorithms to achieve risk control. *The Annals of Applied Statistics*, 2025.

Atf, Z., Safavi-Naini, S. A. A., Lewis, P. R., Mahjoubfar, A., Naderi, N., Savage, T. R., and Soroush, A. The challenge of uncertainty quantification of large language models in medicine. *arXiv preprint arXiv:2504.05278*, 2025.

Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

Bates, S., Angelopoulos, A., Lei, L., Malik, J., and Jordan, M. Distribution-free, risk-controlling prediction sets. *Journal of the ACM (JACM)*, 2021.

Brady, O., Nulty, P., Zhang, L., Ward, T. E., and McGovern, D. P. Dual-process theory and decision-making in large language models. *Nature Reviews Psychology*, 2025.

Campos, M., Farinhas, A., Zerva, C., Figueiredo, M. A., and Martins, A. F. Conformal prediction for natural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 2024a.

Campos, M., Farinhas, A., Zerva, C., Figueiredo, M. A. T., and Martins, A. F. T. Conformal prediction for natural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 2024b.

Cresswell, J. C., Sui, Y., Kumar, B., and Vouitsis, N. Conformal prediction sets improve human decision making. In *Forty-first International Conference on Machine Learning*, 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.

Duan, J., Cheng, H., Wang, S., Zavalny, A., Wang, C., Xu, R., Kailkhura, B., and Xu, K. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024.

Duan, J., Diffenderfer, J., Madireddy, S., Chen, T., Kailkhura, B., and Xu, K. Uprop: Investigating the uncertainty propagation of llms in multi-step agentic decision-making. *arXiv preprint arXiv:2506.17419*, 2025.

Farquhar, S., Kossen, J., Kuhn, L., and Gal, Y. Detecting hallucinations in large language models using semantic entropy. *Nature*, 2024.

Gui, Y., Jin, Y., and Ren, Z. Conformal alignment: Knowing when to trust foundation models with guarantees. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

He, P., Liu, X., Gao, J., and Chen, W. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2021.

Huang, H., Liao, W., Xi, H., Zeng, H., Zhao, M., and Wei, H. Selective labeling with false discovery rate control. *arXiv preprint arXiv:2510.14581*, 2025.

Jin, Y. and Candès, E. J. Selection by prediction with conformal p-values. *Journal of Machine Learning Research*, 2023.

Jin, Y. and Candès, E. J. Model-free selective inference under covariate shift via weighted conformal p-values. *Biometrika*, 2025.

Joshi, M., Choi, E., Weld, D., and Zettlemoyer, L. Trivi-aQA: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, 2017.

Jung, J., Brahman, F., and Choi, Y. Trust or escalate: LLM judges with provable guarantees for human agreement. In *The Thirteenth International Conference on Learning Representations*, 2025.

Kadavath, S., Conerly, T., Askell, A., Henighan, T., Drain, D., Perez, E., Schiefer, N., Hatfield-Dodds, Z., DasSarma, N., Tran-Johnson, E., et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.

Kaur, R., Samplawski, C., Cobb, A. D., Roy, A., Matejek, B., Acharya, M., Elenius, D., Berenbeim, A. M., Pavlik, J. A., Bastian, N. D., and Jha, S. Addressing uncertainty in LLMs to enhance reliability in generative AI. In *Neurips Safe Generative AI Workshop 2024*, 2024.

Kuhn, L., Gal, Y., and Farquhar, S. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations*, 2023.

Lin, Z., Trivedi, S., and Sun, J. Generating with confidence: Uncertainty quantification for black-box large language models. *Transactions on Machine Learning Research*, 2024.

Liu, X., Chen, T., Da, L., Chen, C., Lin, Z., and Wei, H. Uncertainty quantification and confidence calibration in large language models: A survey. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, 2025.

Park, S., Bastani, O., Matni, N., and Lee, I. Pac confidence sets for deep neural networks via calibrated prediction. In *8th International Conference on Learning Representations*, 2020.

Quach, V., Fisch, A., Schuster, T., Yala, A., Sohn, J. H., Jaakkola, T. S., and Barzilay, R. Conformal language modeling. In *The Twelfth International Conference on Learning Representations*, 2024.

Reimers, N. and Gurevych, I. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019a.

Reimers, N. and Gurevych, I. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019b.

Sanh, V., Debut, L., Chaumond, J., and Wolf, T. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.

Shorinwa, O., Mei, Z., Lidard, J., Ren, A. Z., and Majumdar, A. A survey on uncertainty quantification of large language models: Taxonomy, open research challenges, and future directions. *ACM Computing Surveys*, 2025.

Talmor, A., Herzig, J., Lourie, N., and Berant, J. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

Wang, Q., Fan, Y., and Wang, X. E. Safer: Risk-constrained sample-then-filter in large language models. *arXiv preprint arXiv:2510.10193*, 2025a.

Wang, Q., Geng, T., Wang, Z., Wang, T., Fu, B., and Zheng, F. Sample then identify: A general framework for risk control and assessment in multimodal large language models. In *The Thirteenth International Conference on Learning Representations*, 2025b.

Wang, Z., Duan, J., Cheng, L., Zhang, Y., Wang, Q., Shi, X., Xu, K., Shen, H. T., and Zhu, X. Conu: Conformal uncertainty in large language models with correctness coverage guarantees. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024.

Wang, Z., Duan, J., Wang, Q., Zhu, X., Chen, T., Shi, X., and Xu, K. Coin: Uncertainty-guarding selective question answering for foundation models with provable risk guarantees. *arXiv preprint arXiv:2506.20178*, 2025c.

Wang, Z., Duan, J., Yuan, C., Chen, Q., Chen, T., Zhang, Y., Wang, R., Shi, X., and Xu, K. Word-sequence entropy: Towards uncertainty estimation in free-form medical question answering applications and beyond. *Engineering Applications of Artificial Intelligence*, 2025d.

Wang, Z., Wang, Q., Zhang, Y., Chen, T., Zhu, X., Shi, X., and Xu, K. SConU: Selective conformal uncertainty in large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2025e.

Xiaolan, C., Jiayang, X., Shanfu, L., Yexin, L., Mingguang, H., and Danli, S. Evaluating large language models and agents in healthcare: key challenges in clinical applications. *Intelligent Medicine*, 2025.

Zhang, R., Zhang, H., and Zheng, Z. Vl-uncertainty: Detecting hallucination in large vision-language model via uncertainty estimation. *arXiv preprint arXiv:2411.11919*, 2024.

Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 2023.

# A. Proofs

## A.1. Proof of Theorem 3.1

Condition on the calibrated threshold $\hat{\lambda}^{(a)}$ using the calibration set by Eq. (7). For the test sample $(x_{n+1}, y_{n+1}^*)$, we have

$$
\begin{aligned}
\Pr\left(err_{n+1}^{(a)} = 1 \mid u_{n+1}^{(a)} \leq \hat{\lambda}^{(a)}\right) &= \Pr\left(err_{n+1}^{(a)} = 1 \mid S_{n+1}^{(a)}(\hat{\lambda}^{(a)}) = 1\right) \\
&= \frac{\Pr\left(err_{n+1}^{(a)} = 1 \wedge S_{n+1}^{(a)}(\hat{\lambda}^{(a)}) = 1\right)}{\Pr\left(S_{n+1}^{(a)}(\hat{\lambda}^{(a)}) = 1\right)} \\
&= \frac{\Pr\left(Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) = 1\right)}{\Pr\left(S_{n+1}^{(a)}(\hat{\lambda}^{(a)}) = 1\right)} \\
&= \frac{\mathbb{E}\left[Z_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right]}{\mathbb{E}\left[S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right]}
\end{aligned}
\tag{12}
$$

Since the calibration data points and the test sample are exchangeable (Angelopoulos & Bates, 2021), we have

$$
\mathbb{E}\left[Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right] = \frac{1}{n+1} \sum_{i=1}^{n+1} \left(Z_i(\hat{\lambda}^{(a)}) - \alpha S_i(\hat{\lambda}^{(a)})\right).
\tag{13}
$$

Split the sum into calibration and test parts,

$$
\mathbb{E}\left[Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right] = \frac{\sum_{i=1}^{n}\left(Z_i(\hat{\lambda}^{(a)}) - \alpha S_i(\hat{\lambda}^{(a)})\right) + \left(Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right)}{n+1}.
\tag{14}
$$

We now make explicit how the calibration sum rewrites in terms of the sorted errors in Eq. (5). Recall that

$$
S_i^{(a)}\left(\lambda^{(a)}\right) = \mathbf{1}\left\{u_i^{(a)} \leq \lambda^{(a)}\right\}, \quad Z_i^{(a)}(\lambda^{(a)}) = S_i^{(a)}(\lambda^{(a)}) \cdot err_i^{(a)}.
$$

Therefore, for any fixed $\lambda^{(a)}$,

$$
Z_i(\lambda^{(a)}) - \alpha S_i(\lambda^{(a)}) = \begin{cases} err_i^{(a)} - \alpha, & \text{if } u_i^{(a)} \leq \lambda^{(a)}, \\ 0, & \text{if } u_i^{(a)} > \lambda^{(a)}. \end{cases}
$$

Hence,

$$
\sum_{i=1}^{n}\left(Z_i(\lambda^{(a)}) - \alpha S_i(\lambda^{(a)})\right) = \sum_{i : u_i^{(a)} \leq \lambda^{(a)}}\left(err_i^{(a)} - \alpha\right).
\tag{15}
$$

Let $u_{(1)}^{(a)} \leq \cdots \leq u_{(n)}^{(a)}$ be the sorted calibration uncertainties and $err_{(j)}^{(a)}$ the corresponding error indicators. For $\lambda^{(a)} = \hat{\lambda}^{(a)}$, let $k^{(a)}(\hat{\lambda}^{(a)}) = \#\left\{i : u_i^{(a)} \leq \hat{\lambda}^{(a)}\right\}$ be the number of calibration points that would be accepted. By construction of the ordering, the set $\left\{i : u_i^{(a)} \leq \lambda^{(a)}\right\}$ coincides with the first $k^{(a)}(\lambda^{(a)})$ indices in the sorted sequence. Thus,

$$
\mathbb{E}\left[Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right] = \frac{\sum_{j=1}^{k^{(a)}(\hat{\lambda}^{(a)})}\left(err_{(j)}^{(a)} - \alpha\right) + \left(Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right)}{n+1}.
\tag{16}
$$

By the definition of the calibrated threshold $\hat{\lambda}^{(a)}$, the empirical constraint $\sum_{j=1}^{k^{(a)}(\hat{\lambda}^{(a)})}\left(err_{(j)}^{(a)} - \alpha\right) \leq -1$ holds on the calibration data. Plugging this into Eq. (16) yields

$$
\mathbb{E}\left[Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right] \leq \frac{-1 + \left(Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right)}{n+1}
\tag{17}
$$

Since $Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)}) \leq 1$, we have

$$\mathbb{E}\left[Z_{n+1}^{(a)}(\hat{\lambda}^{(a)}) - \alpha S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right] \leq 0. \tag{18}$$

Rearranging gives

$$\frac{\mathbb{E}\left[Z_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right]}{\mathbb{E}\left[S_{n+1}^{(a)}(\hat{\lambda}^{(a)})\right]} \leq \alpha, \tag{19}$$

which is exactly

$$\Pr\left(err_{n+1}^{(a)} = 1 \mid u_{n+1}^{(a)} \leq \hat{\lambda}^{(a)}\right) \leq \alpha. \tag{20}$$

This completes the proof.

### A.2. Proof of Theorem 3.2

Condition on the calibrated threshold pair $(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})$ that is obtained from the calibration set by solving Eq. (11). For the test sample $(x_{n+1}, y_{n+1}^*)$, recall the system-level selection and error indicators

$$S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) \in \{0, 1\}, \quad err_{n+1} \in \{0, 1\},$$

and define the joint indicator

$$Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) = S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) \cdot err_{n+1}.$$

Whenever $\mathbb{E}\left[S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right] > 0$, the system-level FDR at $(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})$ can be written as

$$\Pr\left(err_{n+1} = 1 \mid S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) = 1\right) = \frac{\mathbb{E}\left[Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right]}{\mathbb{E}\left[S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right]}. \tag{21}$$

By the exchangeability of calibration and test samples, we have

$$\begin{aligned}
&\mathbb{E}\left[Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right] \\
&= \frac{1}{n+1}\sum_{i=1}^{n+1}\left(Z_i(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_i(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right) \\
&= \frac{\sum_{i=1}^{n}\left(Z_i(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_i(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right) + \left(Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right)}{n+1}.
\end{aligned} \tag{22}$$

By the definition of the feasible region $\Lambda_\alpha^{(a,b)}$ in Eq. (10), the calibrated pair $(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})$ satisfies

$$\sum_{i=1}^{n}\left(Z_i(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_i(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right) \leq -1$$

on the calibration set. Substituting this inequality into Eq. (22) yields

$$\mathbb{E}\left[Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right] \leq \frac{-1 + \left(Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right)}{n+1}. \tag{23}$$

Since $Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) \in \{0, 1\}$ and $S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) \in \{0, 1\}$, we always have $Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) \leq 1$, and hence

$$\mathbb{E}\left[Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) - \alpha S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right] \leq 0. \tag{24}$$

Rearranging gives

$$\mathbb{E}\left[Z_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right] \leq \alpha \, \mathbb{E}\left[S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})\right].$$

Combining this inequality with Eq. (21) implies

$$\Pr\left(err_{n+1} = 1 \mid S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)}) = 1\right) \leq \alpha.$$

If $\mathbb{E}[S_{n+1}(\hat{\lambda}^{(a)}, \hat{\lambda}^{(b)})] = 0$, the system always abstains and the inequality holds trivially. This completes the proof.

## B. Extension to General Multi-Model Routing Systems

Suppose we have a collection of $M$ models $\{\mathcal{G}^{(1)}, \ldots, \mathcal{G}^{(M)}\}$, each equipped with an uncertainty score $u^{(m)}$ and threshold $\lambda^{(m)}$. A routing policy defines, for each input $x$, either a unique index $r(x) \in \{1, \ldots, M\}$ whose prediction is accepted, or abstention. For any fixed threshold vector $\boldsymbol{\lambda} = (\lambda^{(1)}, \ldots, \lambda^{(M)})$, the system induces a selection indicator

$$S_i(\boldsymbol{\lambda}) = \mathbf{1}\{\text{example } i \text{ is accepted by any model under } \boldsymbol{\lambda}\},$$

and an error indicator

$$err_i(\boldsymbol{\lambda}) = \mathbf{1}\{S_i(\boldsymbol{\lambda}) = 1, \ A(y_i^*, \hat{y}_i) = 0\}.$$

Defining $Z_i(\boldsymbol{\lambda}) = S_i(\boldsymbol{\lambda}) \cdot err_i(\boldsymbol{\lambda})$, the system-level FDR and its linearization constraint remain

$$\text{FDR}(\boldsymbol{\lambda}) = \frac{\mathbb{E}[Z(\boldsymbol{\lambda})]}{\mathbb{E}[S(\boldsymbol{\lambda})]}, \quad \mathbb{E}[Z(\boldsymbol{\lambda}) - \alpha S(\boldsymbol{\lambda})] \leq 0.$$

Following the same argument as before, the finite-sample sufficient condition becomes

$$\sum_{i=1}^{n} \Big( S_i(\boldsymbol{\lambda}) \cdot err_i(\boldsymbol{\lambda}) - \alpha S_i(\boldsymbol{\lambda}) \Big) \leq -1. \tag{25}$$

We denote the feasible threshold region by

$$\Lambda_\alpha^{(1:M)} = \Big\{ \boldsymbol{\lambda} : \sum_{i=1}^{n} \big( S_i(\boldsymbol{\lambda}) \cdot err_i(\boldsymbol{\lambda}) - \alpha S_i(\boldsymbol{\lambda}) \leq -1 \big) \Big\}. \tag{26}$$

Among all feasible threshold vectors, we select the coverage-maximizing solution

$$\hat{\boldsymbol{\lambda}} = \arg \max_{\boldsymbol{\lambda} \in \Lambda_\alpha^{(1:M)}} \frac{1}{n} \sum_{i=1}^{n} S_i(\boldsymbol{\lambda}), \tag{27}$$

which offer valid FDR control to the full routing system.

Our framework extends seamlessly to routing systems with an arbitrary number of models, as long as the routing policy deterministically maps uncertainty scores and thresholds to a unique selected model. The main algorithmic challenge is to search over the multi-dimensional thresholds $\boldsymbol{\lambda}$ efficiently. **Through a similar linear expectation transformation, LEC further generalizes naturally to a broad class of task-specific risk metrics that take a ratio form.**

## C. Details of Experimental Settings

**Details of Utilized Datasets and Models.** For the closed-ended CommonsenseQA dataset, we employ both the full training split (9,741 samples) and the validation split (1,221 samples)[1]. We remove a small number of samples containing non-ASCII characters in either the query or answer that cannot be encoded by the tokenizer. From the remaining data, we select one QA pair as a fixed one-shot demonstration, which is prepended to the prompt for all other samples. After filtering and prompt construction, we select 10,000 QA instances in total. An example of the complete prompt is presented in Figure 6.

For the open-ended TriviaQA dataset, we randomly select 4,000 QA pairs from the validation split of the `rc.nocontext` subset[2]. We also develop a one-shot prompt for each data point. An example of the complete prompt is presented in Figure 7.

For models, we employ three series of open-source LLMs available: LLaMA, Vicuna, and Qwen, divided by the model size into: (1) 3B: Qwen-2.5-3B-Instruct and LLaMA-3.2-3B-Instruct. (2) 7B: Qwen-2.5-7B-Instruct and Vicuna-7B-v1.5 (3) 8B: LLaMA-3.1-8B-Instruct. (4) 13B: Vicuna-13B-v1.5. (5) 14B: Qwen-2.5-14B-Instruct. We omit "Instruct" when reporting the experimental results.

**Details of Alignment Criteria.** Following prior work (Duan et al., 2024; Wang et al., 2025c), we estimate the sentence similarity between two answers (ground truth, the most likely answer, or sampled answer) leveraging SentenceTransformers (Reimers & Gurevych, 2019a) with DistillRoBERTa (Sanh et al., 2019) as the backbone. For bi-entailment (Kuhn et al.,

---

[1]The source file link of the CommonsenseQA dataset.
[2]The source file link of the TriviaQA dataset.

---

**CommonsenseQA**

### System:
Make your best effort and select the correct answer for the following multiple-choice question. For each question, only one choice is correct. Answer should be one among A, B, C, D, E.

### User:
What is something I need to avoid while playing ball?
A: competition
B: losing
C: injury
D: hitting the ball
E: having fun
### Assistant:
C

### User:
The sanctions against the school were a punishing blow, and they seemed to what the efforts the school had made to change?
A: ignore
B: enforce
C: authoritarian
D: yell at
E: avoid
### Assistant:

---

Figure 6: A prompt example in the CommonsenseQA Dataset.

---

**TriviaQA**

### System:
This is a bot that correctly answers questions.

### User:
In 1968, who did radical feminist Valerie Solanas shoot and wound as he entered his New York studio?
### Assistant:
Andy Warhol

### User:
Who was the man behind The Chipmunks?
### Assistant:

---

Figure 7: A prompt example in the TriviaQA Dataset.

---

2023; Farquhar et al., 2024; Wang et al., 2025d), we employ an off-the-shelf DeBERTa-large model (He et al., 2021) as the Natural Language Inference (NLI) classifier, which outputs logits over three semantic relation classes: entailment, neutral, and contradiction. Two answers are deemed semantically aligned if the classifier predicts entailment for both directions, i.e., when evaluated on (answer 1 → answer 2) and (answer 2 → answer 1). In addition, there are evaluation methods such as LLM judgment, but the statistical rigor of our LEC framework is not affected by changes in the alignment criterion.

**Details of Additional Hyperparameters.** In the black-box setting, we set the sampling temperature to 1.0 and top-p to 0.9. For the CommonsenseQA dataset, we sample 20 answers for each question and limit the model's output to a single token, since only the option letter is required. For the TriviaQA dataset, we set the maximum output length to 36 tokens.
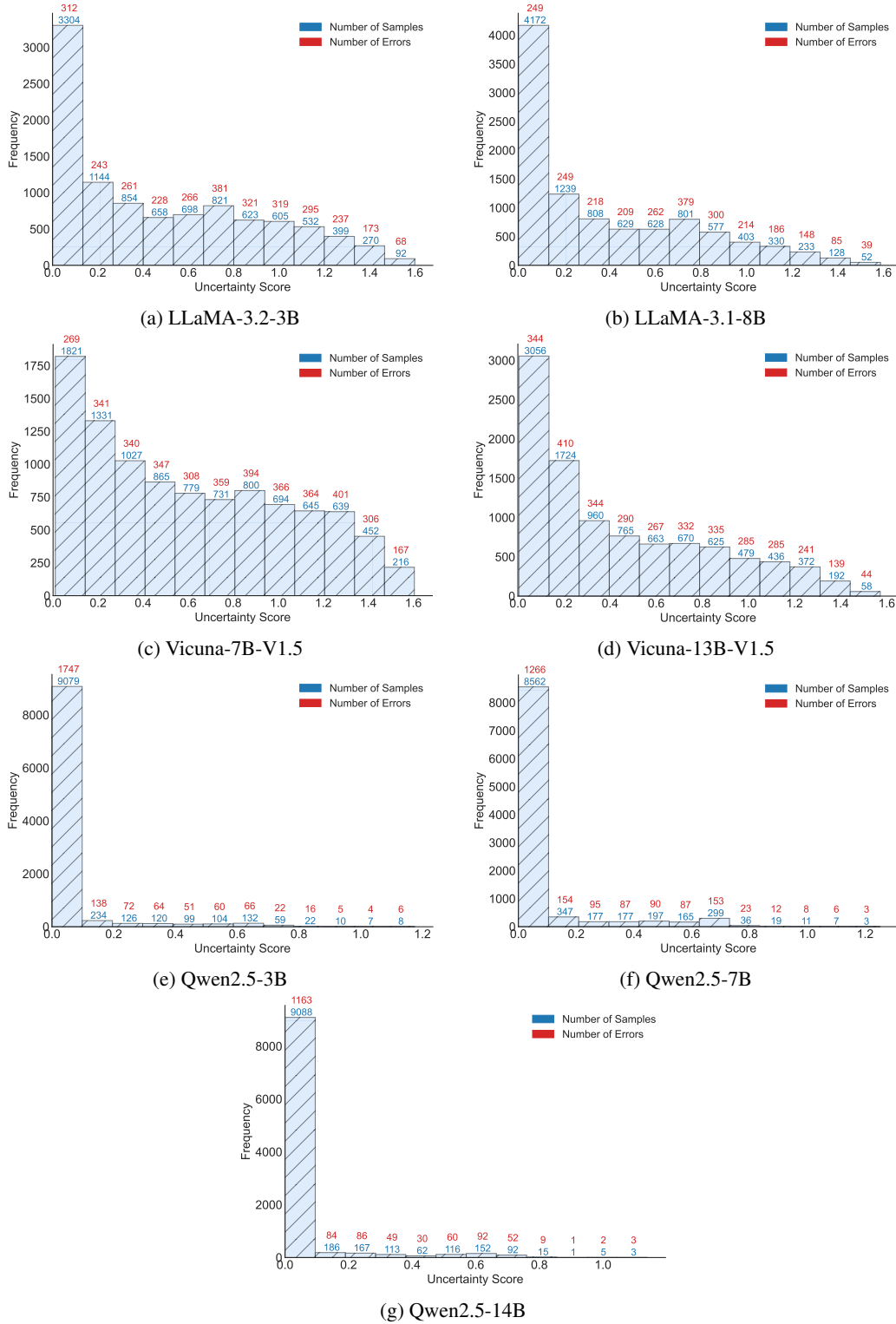


Figure 8: Uncertainty and correctness distribution on the CommonsenseQA dataset across seven LLMs.

# D. Additional Experimental Results

**Necessity and Justification of FDR Control.** Here, we only report the distribution of uncertainty scores and corresponding correctness for different models on the CommonsenseQA dataset. As presented in Figure 8, incorrect answers appear across almost all uncertainty ranges, highlighting that uncertainty estimation alone cannot perfectly separate correct from incorrect predictions. In this setting, we calibrate a rigorous threshold under a user-specified error tolerance (i.e., $\alpha$) to achieve FDR control for selective prediction, ensuring that the accepted subset attains a higher accuracy than the base model.

In addition, for certain LLMs, a non-negligible portion of incorrect samples show uncertainty scores concentrated within a relatively small range, which makes low risk levels unattainable. For example, when using LLaMA-3.2-3B on CommonsenseQA, a risk level of 0.05 cannot be satisfied under the linear constraint regardless of how the threshold is tuned on the calibration set. This issue can be mitigated by employing a stronger LLM or a UQ model with better discriminative power between correct and incorrect answers.

Table 3: Power comparison on the TriviaQA dataset using the SE method under black-box scenarios. We employ sentence similarity with a threshold of 0.6 as the alignment criterion in the admissible function $A$. In the semantic clustering process of SE, we utilize the same criterion to evaluate semantic equivalence. We employ 10 sampled generations to compute SE.

| LLMs | Method / $\alpha$ | 0.05 | 0.06 | 0.07 | 0.08 | 0.09 | 0.1 |
|---|---|---|---|---|---|---|---|
| Qwen2.5-3B | COIN-SP | 0.8076 | 0.8504 | 0.9337 | 0.9846 | 0.9986 | 0.9997 |
| | COIN-HFD | - | - | 0.8064 | 0.8469 | 0.9395 | 0.9890 |
| | LEC | **0.8288** | **0.9203** | **0.9771** | **0.9990** | **0.9997** | **0.9997** |
| LLaMA-3.2-3B | COIN-SP | - | 0.7465 | 0.8063 | 0.8981 | 0.9345 | 0.9586 |
| | COIN-HFD | - | - | - | 0.7677 | 0.8156 | 0.9072 |
| | LEC | - | **0.7739** | **0.8934** | **0.9314** | **0.9565** | **0.9752** |
| Qwen2.5-7B | COIN-SP | 0.9130 | 0.9596 | 0.9948 | 0.9998 | 0.9999 | 0.9999 |
| | COIN-HFD | - | - | 0.9152 | 0.9611 | 0.9960 | 0.9998 |
| | LEC | **0.9410** | **0.9901** | **0.9998** | **0.9999** | **0.9999** | **0.9999** |
| Vicuna-7B-V1.5 | COIN-SP | - | - | 0.7874 | 0.8153 | 0.8583 | 0.9122 |
| | COIN-HFD | - | - | - | - | 0.7909 | 0.8357 |
| | LEC | - | - | **0.8146** | **0.8688** | **0.9195** | **0.9495** |
| LLaMA-3.1-8B | COIN-SP | 0.8999 | 0.9501 | 0.9828 | 0.9990 | 0.9998 | 0.9998 |
| | COIN-HFD | - | 0.7999 | 0.8897 | 0.9532 | 0.9874 | 0.9996 |
| | LEC | **0.9357** | **0.9749** | **0.9990** | **0.9998** | **0.9998** | **0.9998** |
| Vicuna-13B-V1.5 | COIN-SP | 0.8238 | 0.8417 | 0.8877 | 0.9316 | 0.9644 | 0.9862 |
| | COIN-HFD | - | - | 0.8316 | 0.8260 | 0.8639 | 0.9182 |
| | LEC | **0.8250** | **0.8865** | **0.9339** | **0.9662** | **0.9877** | **0.9993** |
| Qwen2.5-14B | COIN-SP | 0.9894 | 0.9996 | 0.9997 | 0.9997 | 0.9997 | 0.9997 |
| | COIN-HFD | - | 0.9503 | 0.9896 | 0.9996 | 0.9997 | 0.9997 |
| | LEC | **0.9994** | **0.9997** | **0.9997** | **0.9997** | **0.9997** | **0.9997** |

**Power Comparison on TriviaQA.** As presented in Table 3, we also evaluate the power of LEC on the TriviaQA dataset. Compared to both versions of COIN, LEC consistently accepts a larger proportion of correct samples across various risk levels, demonstrating higher efficiency at test time. For instance, when utilizing the Qwen2.5-3B model, LEC accepts 92.03% of correct samples at a risk level of 0.06, exceeding COIN-CP by 6.99% at test time.

**Utilization of Alternative UQ Approaches on TriviaQA with Sentence Similarity to Evaluate Alignment.** In the above evaluations, we sample 10 answers per question by default and compute SE-based uncertainty scores. In this section, we examine the performance of LEC when using alternative UQ methods. As shown in Figure 9a, using the LLaMA-3.1-8B model as an example, LEC continues to achieve strict FDR control under various risk levels (e.g., 0.05) when employing Ecc, Deg, or EigV as the uncertainty metric, demonstrating its robustness. Meanwhile, as illustrated in Figure 9b, employing stronger UQ methods or increasing the sampling size in the black-box setting consistently improves the power of LEC,
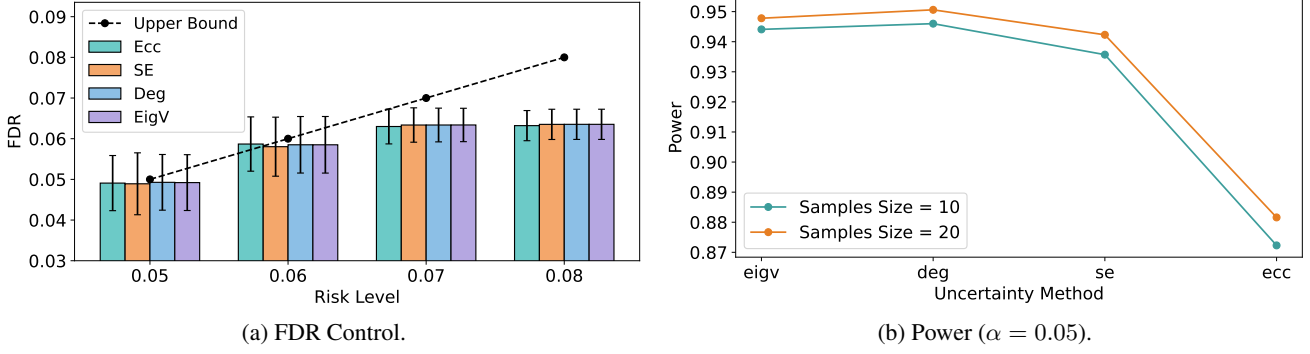
(a) FDR Control.

(b) Power ($\alpha = 0.05$).

Figure 9: Evaluation of LEC on the TriviaQA dataset using the LLaMA-3.1-8B model across various uncertainty methods.

allowing it to accept more correct samples at test time. This highlights LEC's potential to deliver more efficient selective prediction as heuristic UQ methods advance.
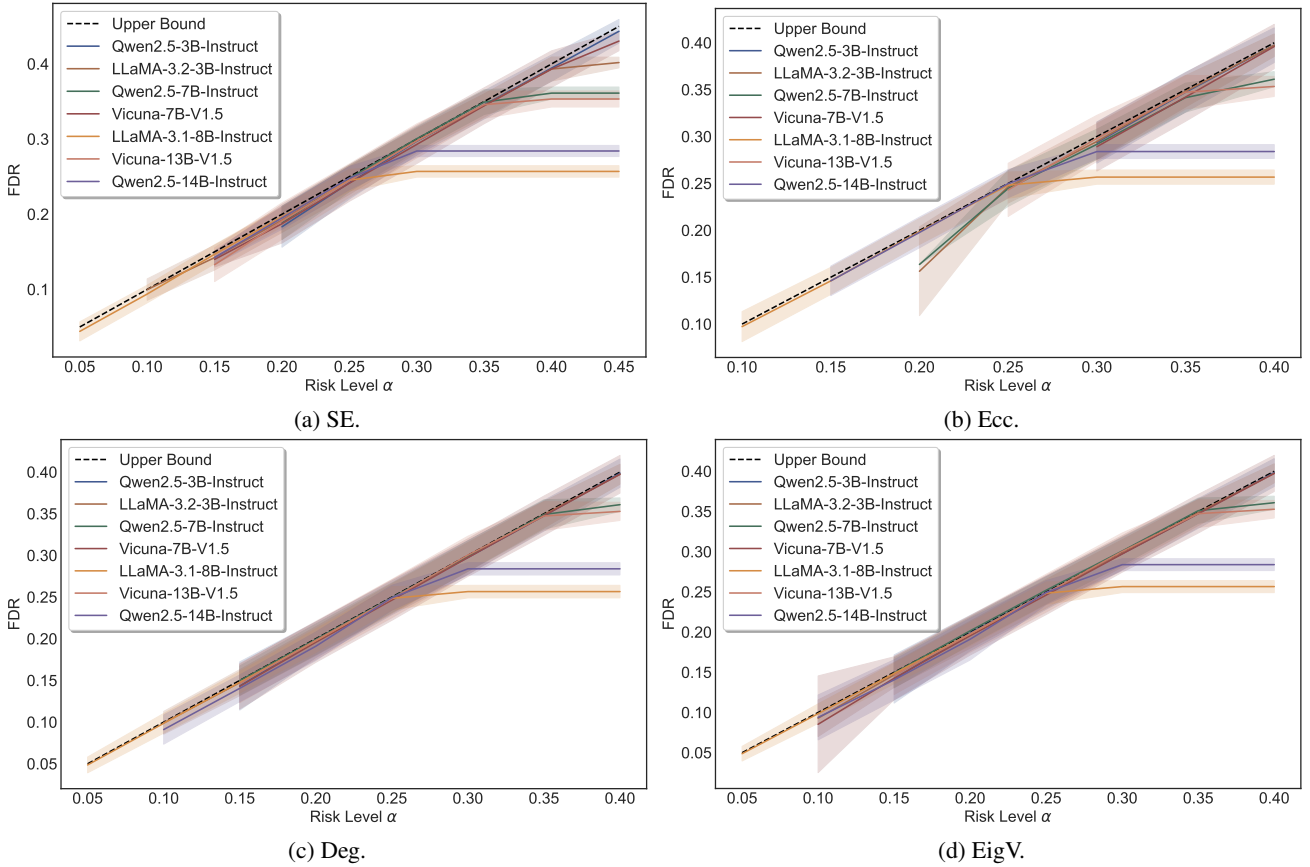


(a) SE.

(b) Ecc.

(c) Deg.

(d) EigV.

Figure 10: FDR control on the TriviaQA dataset across four uncertainty methods utilizing bi-entailment as the alignment criterion in the admissible function $A$ and the semantic clustering process of SE. We utilize 10 sampled generations.

**Utilization of Alternative UQ Methods on TriviaQA with Bi-entailment to Evaluate Alignment.** LEC is compatible with a wide range of alignment criteria. We further evaluate the bi-entailment alignment method, and as shown in Figure 10, LEC maintains FDR control for selective prediction on the test set across various user-specified risk levels, regardless of the UQ method used. Because different UQ methods vary in their ability to distinguish correct from incorrect answers, weaker UQ methods may cause LEC to fail at certain low risk levels. However, this issue can be detected during calibration: by adjusting the threshold on the calibration data, we can directly assess whether the finite-sample sufficient condition is

17

Table 4: Power comparison on TriviaQA with the LLaMA-3.1-8B model. We set various sampling sizes when computing SE under black-box scenarios. We employ bi-entailment as the alignment criterion in the admissible function $A$. In the semantic clustering process of SE, we utilize the same criterion to evaluate semantic equivalence.

(a) 5 samples

| Method / $\alpha$ | 0.1 | 0.15 | 0.2 | 0.25 |
|---|---|---|---|---|
| COIN-SP | 0.7268 | 0.8356 | 0.9336 | 0.9689 |
| COIN-HFD | 0.7223 | 0.7534 | 0.9070 | 0.9589 |
| LEC | **0.7268** | **0.8711** | **0.9554** | **0.9836** |

(b) 10 samples

| Method / $\alpha$ | 0.1 | 0.15 | 0.2 | 0.25 |
|---|---|---|---|---|
| COIN-SP | 0.7424 | 0.8651 | 0.9436 | 0.9856 |
| COIN-HFD | 0.6604 | 0.8302 | 0.9301 | 0.9785 |
| LEC | **0.7832** | **0.8925** | **0.9589** | **0.9920** |

(c) 15 samples

| Method / $\alpha$ | 0.1 | 0.15 | 0.2 | 0.25 |
|---|---|---|---|---|
| COIN-SP | 0.7598 | 0.8791 | 0.9491 | 0.9873 |
| COIN-HFD | 0.6532 | 0.8381 | 0.9327 | 0.9801 |
| LEC | **0.7949** | **0.8973** | **0.9610** | **0.9950** |

(d) 20 samples

| Method / $\alpha$ | 0.1 | 0.15 | 0.2 | 0.25 |
|---|---|---|---|---|
| COIN-SP | 0.7780 | 0.8741 | 0.9491 | 0.9885 |
| COIN-HFD | 0.6886 | 0.8473 | 0.9332 | 0.9811 |
| LEC | **0.8027** | **0.8979** | **0.9647** | **0.9953** |

satisfied. If no threshold meets the linear expectation constraint, LEC abstains from the user-specified risk level, ensuring that selective prediction is never performed on given test samples under an unattainable risk requirement. This highlights the robustness and practicality of the LEC framework.

In addition, when using the SE method, we evaluate the power of LEC under different sampling sizes. As shown in Table 4, LEC consistently achieves the highest power across all sampling configurations. For example, when sampling 5 answers per question, LEC reaches a power of $87.11\%$ at a risk level of 0.15, outperforming COIN-HFD by roughly $12\%$ and COIN-CP by about $3.5\%$. Moreover, LEC consistently achieves over $95\%$ power at a risk level of 0.2, demonstrating its high efficiency.

Table 5: Power comparison on the CommonsenseQA dataset using the PE method under black-box scenarios.

| LLMs | Methods / $\alpha$ | 0.15 | 0.2 | 0.25 | 0.3 | 0.35 | 0.4 | 0.45 |
|---|---|---|---|---|---|---|---|---|
| Qwen2.5-3B | COIN-CP | - | - | 0.9957 | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ |
| | COIN-HFD | - | - | 0.9841 | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ |
| | LEC | - | - | **0.9999** | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ |
| LLaMA-3.2-3B | COIN-CP | - | 0.6813 | 0.8119 | 0.9316 | $0.9999_{(7)}$ | $0.9999_{(8)}$ | $0.9999_{(8)}$ |
| | COIN-HFD | - | 0.6552 | 0.7931 | 0.9155 | 0.9998 | $0.9999_{(8)}$ | $0.9999_{(8)}$ |
| | LEC | - | **0.6950** | **0.8367** | **0.9500** | $\mathbf{0.9999}_{(8)}$ | $\mathbf{0.9999}_{(8)}$ | $\mathbf{0.9999}_{(8)}$ |
| Qwen2.5-7B | COIN-CP | - | 0.9602 | 0.9999 | 0.9999 | 0.9999 | 0.9999 | 0.9999 |
| | COIN-HFD | - | 0.9490 | 0.9999 | 0.9999 | 0.9999 | 0.9999 | 0.9999 |
| | LEC | - | **0.9759** | **0.9999** | **0.9999** | **0.9999** | **0.9999** | **0.9999** |
| Vicuna-7B-V1.5 | COIN-CP | - | - | - | 0.6113 | 0.8040 | 0.9541 | 0.9998 |
| | COIN-HFD | - | - | - | 0.5790 | 0.7776 | 0.9415 | 0.9998 |
| | LEC | - | - | - | **0.6605** | **0.8403** | **0.9746** | **0.9998** |
| LLaMA-3.1-8B | COIN-CP | 0.7479 | 0.8561 | 0.9516 | $0.9999_{(8)}$ | $0.9999_{(8)}$ | $0.9999_{(8)}$ | $0.9999_{(8)}$ |
| | COIN-HFD | 0.7381 | 0.8366 | 0.9388 | $0.9999_{(7)}$ | $0.9999_{(8)}$ | $0.9999_{(8)}$ | $0.9999_{(8)}$ |
| | LEC | **0.7676** | **0.8695** | **0.9640** | $\mathbf{0.9999}_{(8)}$ | $\mathbf{0.9999}_{(8)}$ | $\mathbf{0.9999}_{(8)}$ | $\mathbf{0.9999}_{(8)}$ |
| Vicuna-13B-V1.5 | COIN-CP | - | - | 0.7615 | 0.9037 | 0.9864 | 0.9998 | 0.9998 |
| | COIN-HFD | - | - | 0.7322 | 0.8901 | 0.9781 | 0.9998 | 0.9998 |
| | LEC | - | - | **0.8048** | **0.9201** | **0.9979** | **0.9998** | **0.9998** |
| Qwen2.5-14B | COIN-CP | 0.9708 | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ |
| | COIN-HFD | - | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ | $0.9999_{(6)}$ |
| | LEC | **0.9742** | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ | $\mathbf{0.9999}_{(6)}$ |

**Power Comparison on CommonsenseQA at black-box settings.** In closed-ended tasks where model logits are unavailable, we estimate option probabilities by sampling multiple outputs for each question and using the empirical frequency of each option as its probability. As shown in Table 5, we set the sampling size to 20 and evaluate the power of LEC across seven models. Under various feasible risk levels, LEC consistently accepts a larger fraction of correct test samples than confidence-interval–based methods. However, because sampling-frequency–based uncertainty is more discrete than logit-based uncertainty and has weaker discriminative power between correct and incorrect answers, the minimum achievable risk level is higher than in the white-box setting. Nevertheless, using LLaMA-3.1-8B as an example, LEC still attains strict FDR control at a risk level of 0.15, with its power exceeding COIN-CP by $2\%$.
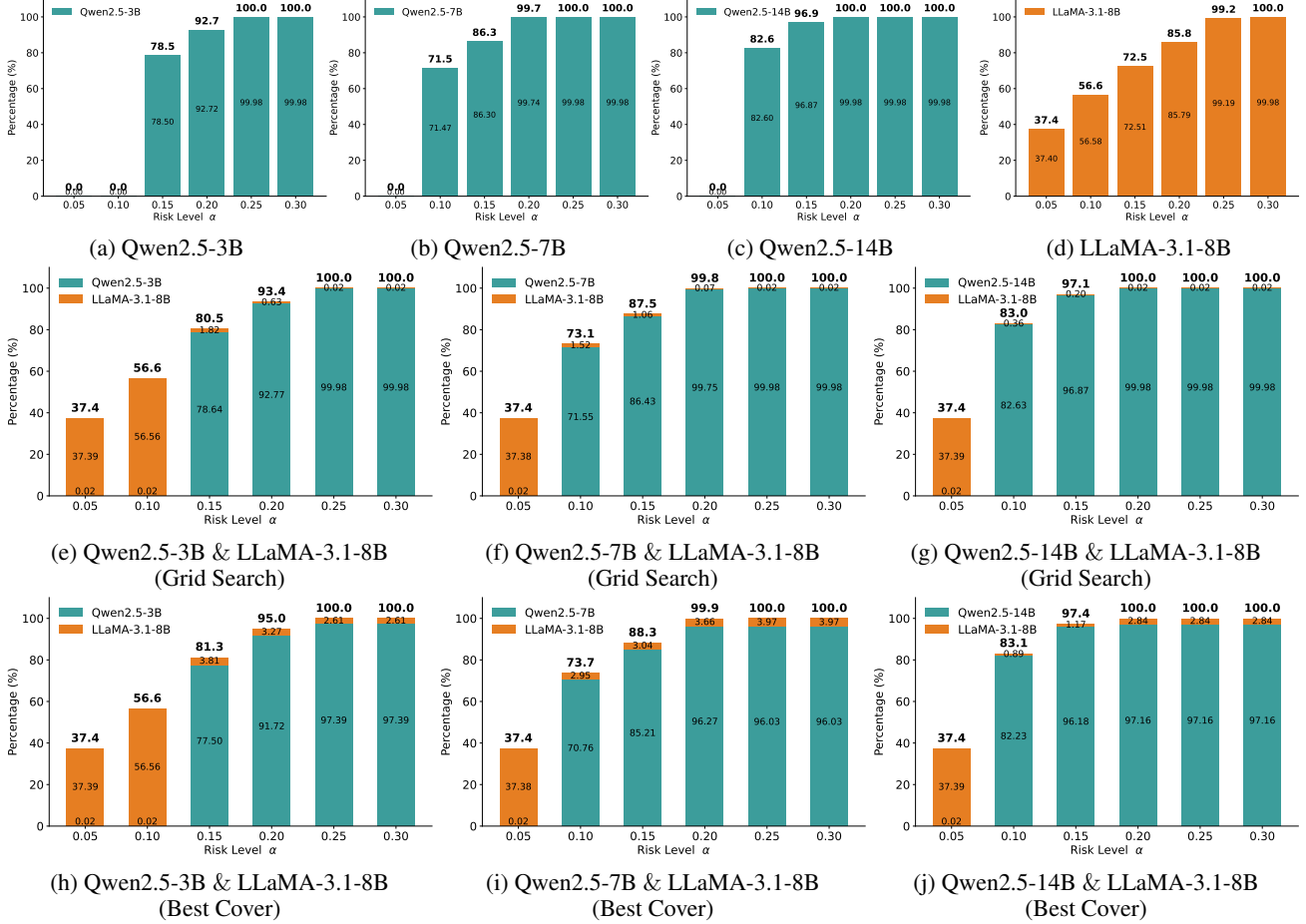


(a) Qwen2.5-3B    (b) Qwen2.5-7B    (c) Qwen2.5-14B    (d) LLaMA-3.1-8B

(e) Qwen2.5-3B & LLaMA-3.1-8B (Grid Search)    (f) Qwen2.5-7B & LLaMA-3.1-8B (Grid Search)    (g) Qwen2.5-14B & LLaMA-3.1-8B (Grid Search)

(h) Qwen2.5-3B & LLaMA-3.1-8B (Best Cover)    (i) Qwen2.5-7B & LLaMA-3.1-8B (Best Cover)    (j) Qwen2.5-14B & LLaMA-3.1-8B (Best Cover)

Figure 11: The allocation ratio of samples in a two-model routing system during selective prediction.

**Sample Allocation in Two-Model Routing.** Taking CommonsenseQA as an example, we record the sample allocation during two-model routing at test time. As shown in Figure 11, when using a single model, certain low risk levels become unattainable. For instance, with Qwen2.5-3B, risk levels of 0.05 and 0.1 cannot be satisfied, and the proportion of accepted samples remains low. For example, at a risk level of 0.15, LLaMA-3.1-8B accepts only $72.5\%$ of samples. With the routing mechanism, however, LEC achieves lower feasible risk levels compared to single-model settings. When routing between Qwen2.5-3B and LLaMA-3.1-8B, LEC attains valid FDR control at both 0.05 and 0.1. Moreover, at a risk level of 0.15, $81.3\%$ of samples are accepted—nearly $10\%$ higher than using LLaMA-3.1-8B alone. Furthermore, when grid search is used during calibration without explicitly maximizing coverage, the resulting threshold pair tends to accept fewer samples than the best-cover strategy, where thresholds are updated only when they increase the number of covered calibration samples. Under multi-model routing, LEC consistently provides FDR guarantees. The choice of thresholds is orthogonal to the LEC framework, and optimizing the threshold-selection strategy during calibration can improve test-time efficiency.