



# Sample then Identify: A General Framework for Risk Control and Assessment in Multimodal Large Language Models

Qingni Wang<sup>1</sup>, Tiantian Geng<sup>2,3</sup>, Zhiyuan Wang<sup>1</sup>, Teng Wang<sup>2,4</sup>,  
Bo Fu<sup>1\*</sup>, Feng Zheng<sup>2\*</sup>

<sup>1</sup> School of Computer Science and Engineering, University of Electronic Science and Technology of China

<sup>2</sup> Southern University of Science and Technology

<sup>3</sup> University of Birmingham

<sup>4</sup> The University of Hong Kong

Presented by: Qingni Wang

\* Corresponding co-authors

# Background

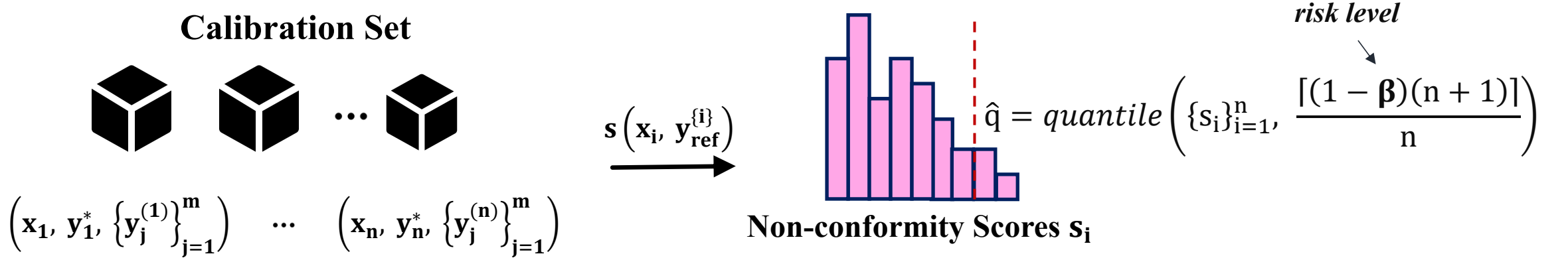
- ❑ Multimodal Large Language Models (MLLMs) exhibit promising advancements across various tasks, yet they still encounter significant trustworthiness issues.
- ❑ Uncertainty quantification (UQ) can provide valuable insights into the reliability of model responses, facilitating risk control and assessment.
- ❑ Split Conformal prediction (SCP) is a distribution-free and model-agnostic approach to UQ, which transforms any heuristic notion of uncertainty into a statistically rigorous one.

# Motivation

Adapting SCP for MLLMs in practical closed-ended and open-domain VideoQA tasks presents several challenges:

- ❑ Adaptability. Unlike tasks equipped with fixed options covering the correct answers, the output space for each VideoQA task is unbounded and there may not be an acceptable response by sampling multiple generations;
- ❑ Reliability. Logits or verbalized confidence levels may be miscalibrated, leading to biased prediction sets;
- ❑ Flexibility. For some API-only MLLMs, users lack access to internal model information like sequence logits.

# Our Framework TRON



Correlate the non-conformity score with *the uncertainty state of  $\mathbf{y}_{\text{ref}}$  within the sampling set  $\{\mathbf{y}_j^{(i)}\}_{j=1}^m$ , which is semantically equivalent to  $\mathbf{y}_{\text{ref}}^{\{i\}}$* , assuming that we can obtain at least one acceptable response by sampling  $m$  generations.

$$\frac{[(N + 1)(1 - \alpha)]}{N} \text{ quantile}$$



conformal scores  $\{\mathbf{r}_i\}_{i=1}^n$

Before the identification of reliable generations, we can also manage the error rate of the sampling set failing to cover acceptable responses by developing the **conformal score** that determines the minimum sampling size.

$$r_i = r(\mathbf{x}_i, \mathbf{y}_i^*) := \sup \left\{ M_i : \forall M'_i < M_i, \mathbf{y}_i^* \notin \{\mathbf{y}_j^{(i)}\}_{j=1}^{M'_i} \right\}$$

We set the sampling size of each test sample to the  $\hat{r}$  based on the other risk level of  $\alpha$ .

## Calibration

Step 1:  
Calibrate the number of response samples

$[(N+1)(1-\alpha)]/N$  quantile

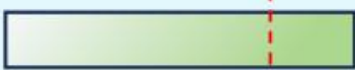


conformal scores  $\{r_i\}$

$\hat{r}$

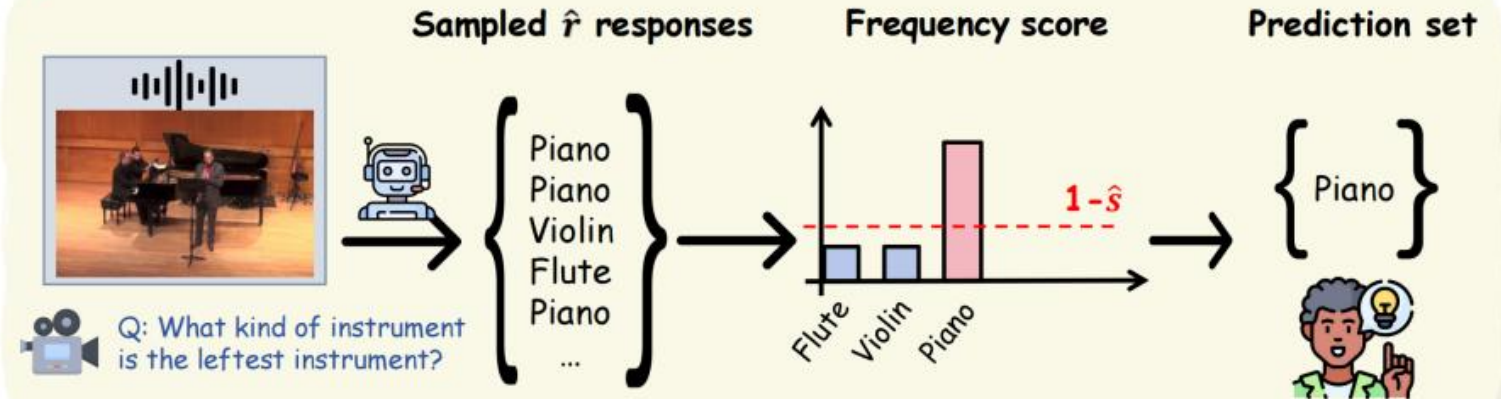
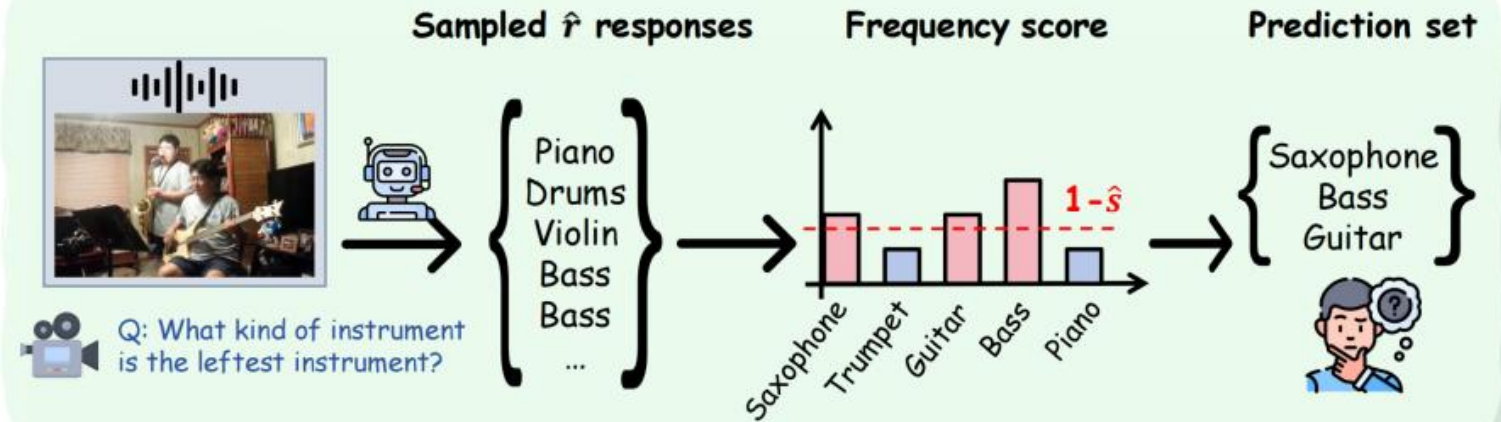
Step 2:  
Calibrate the threshold to identify high-quality responses

$[(N+1)(1-\beta)]/N$  quantile



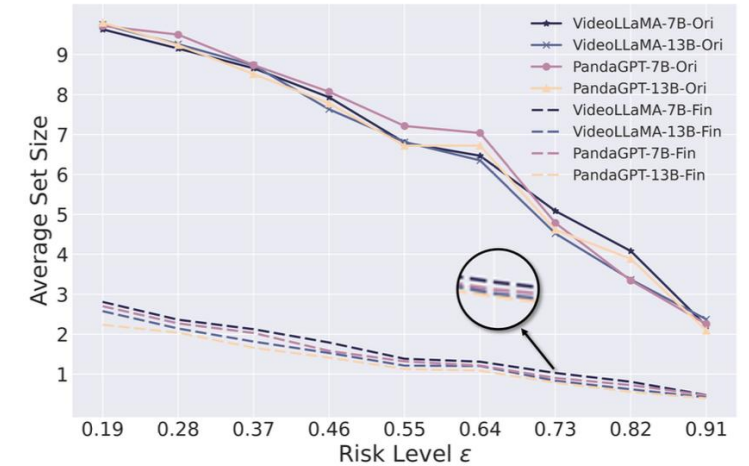
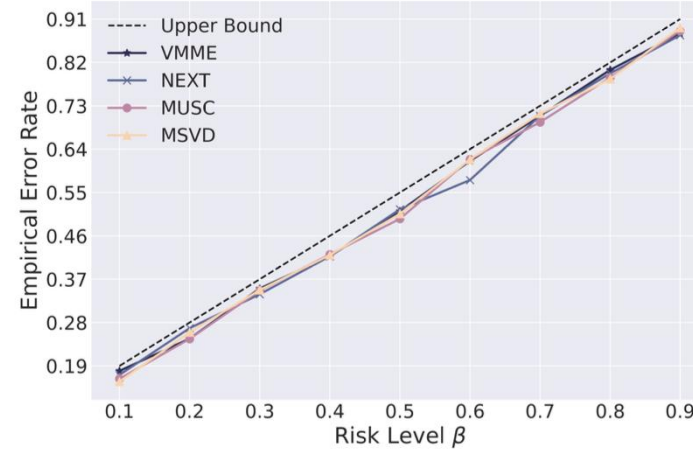
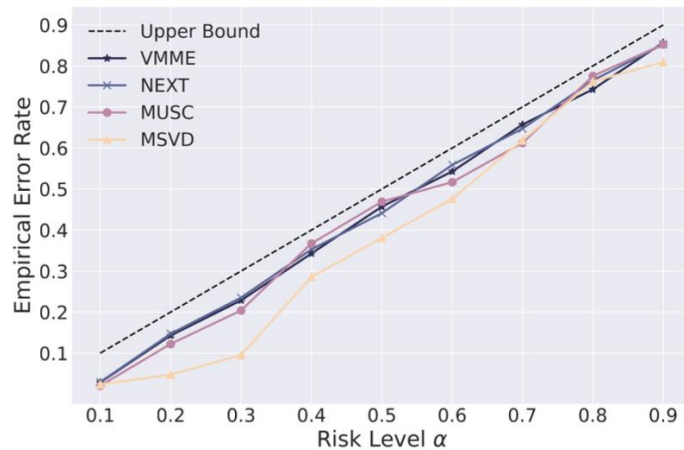
nonconformity scores  $\{s_i\}$

$\hat{s}$



*An illustration of our framework in open-domain VideoQA tasks.*

# Results of TRON in Risk Control and Assessment



Results of the EER metric at various risk levels of  $\alpha$  and  $\beta$ .

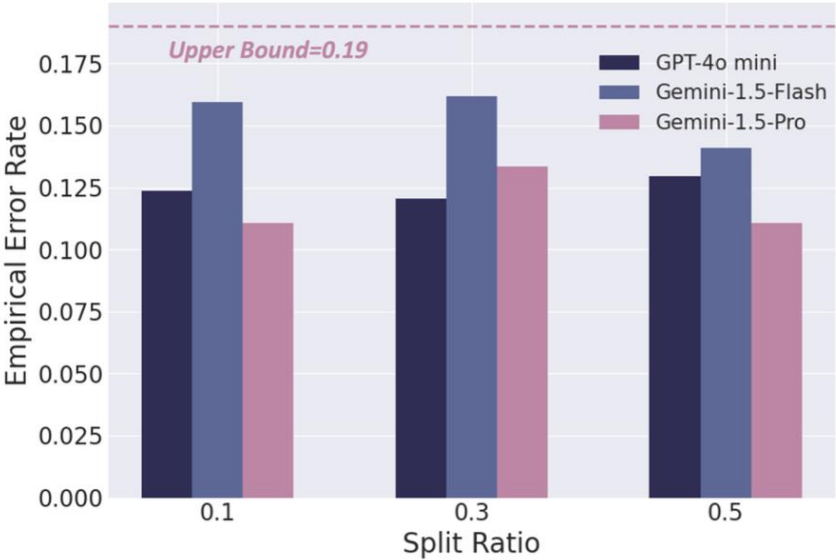
Comparison of APSS before and after Semantic duplication at different risk levels ( $\epsilon$ ).

MLLMs	ACC (%) $\uparrow$				APSS $\downarrow$			
	VMME	NEXT	MUSC	MSVD	VMME	NEXT	MUSC	MSVD
<i>Open-source Models</i>								
VideoLLaMA-7B	38.10	46.06	78.70	78.29	2.81	2.68	1.97	1.63
VideoLLaMA-13B	45.33	47.07	79.18	78.91	2.58	2.37	1.59	1.63
PandaGPT-7B	43.45	50.95	81.85	77.29	2.70	2.10	2.02	1.12
PandaGPT-13B	48.75	58.90	82.73	86.04	2.24	2.20	2.03	2.21
NExT-GPT	42.64	42.59	79.84	84.37	2.43	2.45	2.20	1.70
<i>Closed-source Models</i>								
Gemini-1.5-Flash	72.31	70.80	85.58	87.14	1.22	1.17	1.55	1.08
Gemini-1.5-Pro	73.11	76.29	85.16	86.67	1.15	1.15	1.59	1.10
GPT-4o mini	73.26	82.50	86.27	85.71	1.14	1.08	1.43	1.21

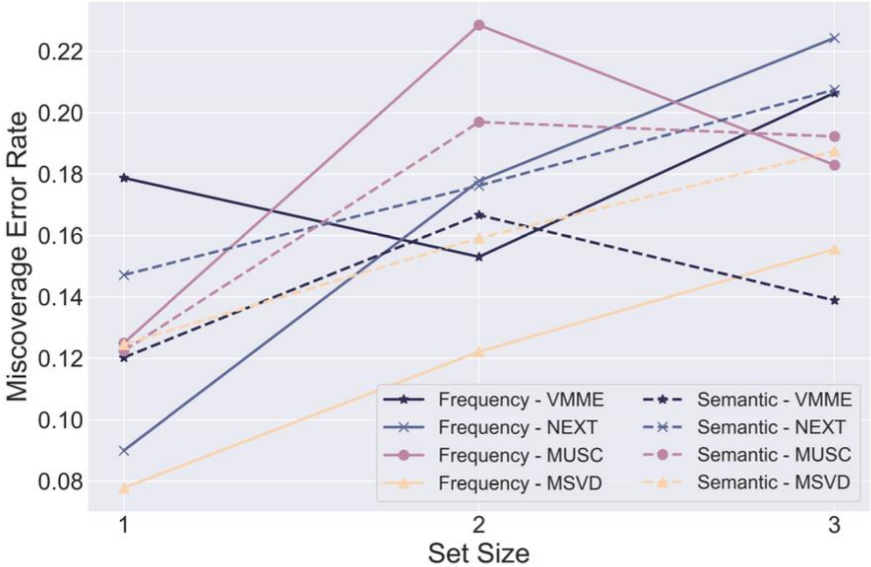
Evaluations of the ACC and APSS metrics utilizing five open-source and three closed-source MLLMs across four open-ended and closed-ended datasets.



# Sensitivity Analyses of TRON



EERs on the MUSC dataset utilizing three closed-source MLLMs when varying the split ratio of the calibration and test set.

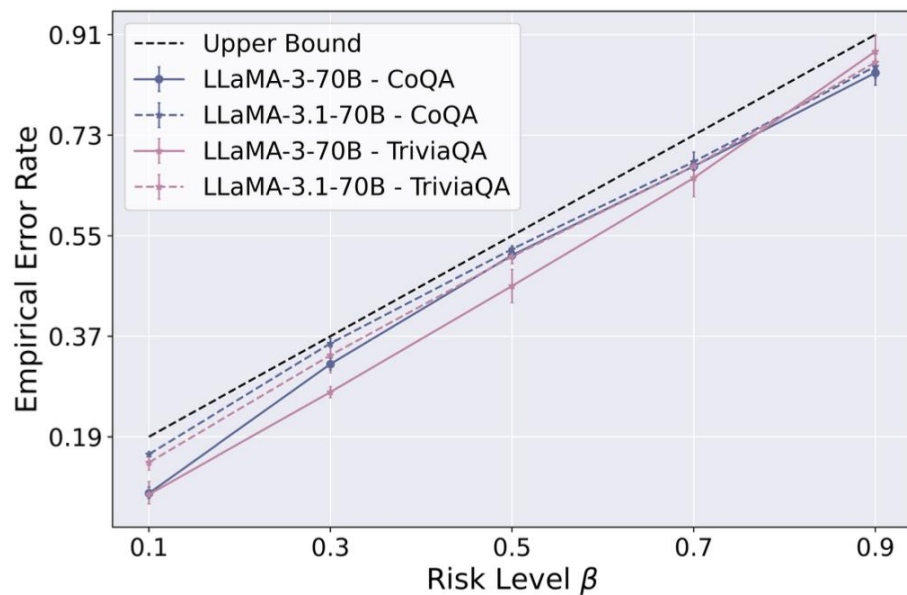


Stratified miscoverage rate at various size of prediction set on four VideoQA datasets utilizing the Gemini-1.5-Pro model.

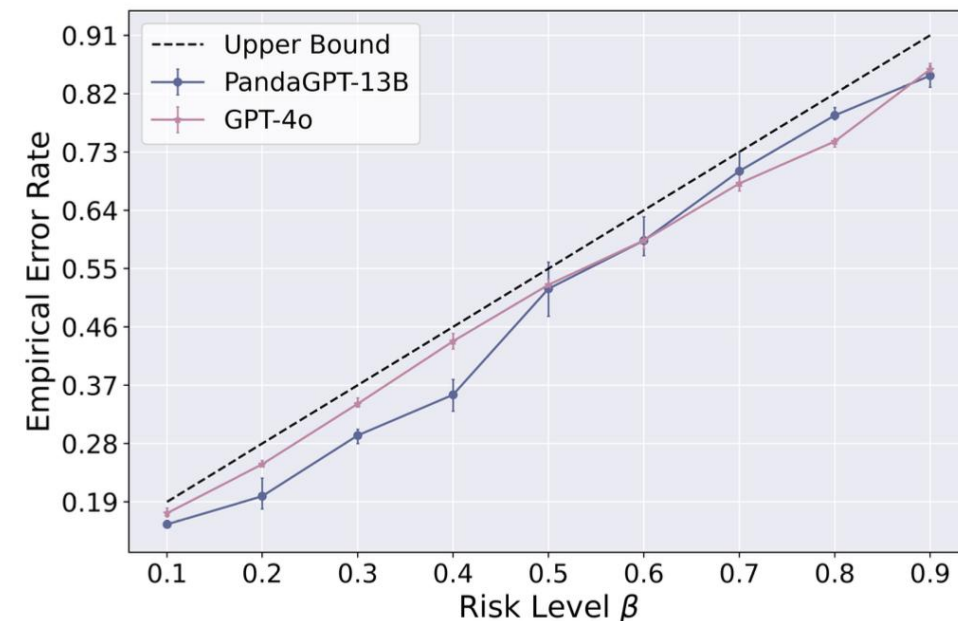
MLLMs	APSS		EER	
	Frequency	Semantic Diversity	Frequency	Semantic Diversity
VideoLLaMA-7B	2.81	2.74	0.1793	0.1766
VideoLLaMA-13B	2.58	2.52	0.1767	0.1744
PandaGPT-7B	2.70	2.42	0.1812	0.1675
PandaGPT-13B	2.24	2.20	0.1783	0.1674

Comparison of frequency and semantic diversity on the VMME datasets employing four open-source MLLMs.

# Generalizability of TRON



Results of the EER metric at various risk levels on the CoQA and TriviaQA tasks utilizing LLaMA-3-70B-Instruct and LLaMA-3.1-70B-Instruct.



Results of the EER metric, obtained from 5 trials with randomly allocated calibration and test data at various risk levels on the VQA dataset utilizing PandaGPT-13B and GPT-4o.

Model	LAC	TRON ( $M = 5$ )	TRON ( $M = 10$ )	TRON ( $M = 20$ )
LLaMA-3-8B-Instruct	2.93 <sub>(2)</sub>	3.06 <sub>(0)</sub>	2.93 <sub>(7)</sub>	<b>2.76<sub>(3)</sub></b>
LLaMA-3.1-8B-Instruct	2.57 <sub>(8)</sub>	2.61 <sub>(3)</sub>	2.53 <sub>(6)</sub>	<b>2.50<sub>(4)</sub></b>

The APSS metric ( $\downarrow$ ) on the MMLU dataset.

In the calibration set, we define the minimum sampling size (conformal score  $\uparrow$ ) that ensures  $y_i^* \in \{y_j^{(i)}\}_{j=1}^M$  as

$$r_i = r(x_i, y_i^*) := \sup \left\{ M_i : \forall M_i' < M_i, y_i^* \notin \{y_j^{(i)}\}_{j=1}^{M_i'} \right\}.$$

We sort the conformal scores so that  $r_1 \leq \dots \leq r_N$ , and calculate the  $\frac{\lceil (1-\alpha)(N+1) \rceil}{N}$  quantile:

$$\hat{r} := \inf \left\{ r : \frac{|\{i: r_i \leq r\}|}{N} \geq \frac{\lceil (1-\alpha)(N+1) \rceil}{N} \right\} = r_{\lceil (1-\alpha)(N+1) \rceil}.$$

We proceed by noting the equality of the two events:

$$\left\{ y_{N+1}^* \notin \{y_j^{(N+1)}\}_{j=1}^{\hat{r}} \right\} = \{r_{N+1} > \hat{r}\}.$$

By the exchangeability of  $N + 1$  data points, we have

$$\mathbb{P}(r_{N+1} \leq r_i) = \frac{i}{N+1}.$$

Then, we can guarantee the upper bound of the probability that the sampling set of size  $M = \hat{r}$  fails to encompass an acceptable response based on a user-specified risk level of  $\alpha$ :

$$\begin{aligned} \mathbb{P} \left( y_{N+1}^* \notin \{y_j^{(N+1)}\}_{j=1}^{\hat{r}} \right) &= \mathbb{P}(r_{N+1} > \hat{r}) = \mathbb{P}(r_{N+1} > r_{\lceil (1-\alpha)(N+1) \rceil}) \\ &= 1 - \mathbb{P}(r_{N+1} \leq r_{\lceil (1-\alpha)(N+1) \rceil}) = 1 - \frac{\lceil (1-\alpha)(N+1) \rceil}{N} \leq \alpha \end{aligned}$$



If there are no acceptable responses in the candidate set,  $\left\{y_{N+1}^* \notin \mathcal{C}\left(\beta, \{s_i\}_{i=1}^N, x_{N+1}, \left\{y_j^{(N+1)}\right\}_{j=1}^{M=\hat{r}}\right)\right\}$  is a certain event, and the probability of miscoverage is

$$\mathbb{P}\left(\left\{y_{N+1}^* \notin \mathcal{C}\left(\beta, \{s_i\}_{i=1}^N, x_{N+1}, \left\{y_j^{(N+1)}\right\}_{j=1}^M\right)\right\} \mid y_{N+1}^* \notin \left\{y_j^{(N+1)}\right\}_{j=1}^{\hat{r}}\right) = \mathbb{P}\left(y_{N+1}^* \notin \left\{y_j^{(N+1)}\right\}_{j=1}^{\hat{r}}\right).$$

Conversely, if acceptable responses exist in the sampling set, we define the non-conformity score ( $\downarrow$ ).

First, we determine the acceptable response of the  $i$ -th calibration data point  $y_{ref}^{(i)}$ , which is semantically equivalent to the correct answer  $y_i^*$ . Then we formulate the non-conformity score as

$$s_i = s(x_i, y_i^*) = 1 - F\left(y_{ref}^{(i)}\right),$$

where  $F\left(y_{ref}^{(i)}\right)$  can be any reliable notion of  $y_{ref}^{(i)}$  in the sampling set (black-box or white box).

We construct the prediction set following

$$\mathcal{C}\left(\beta, \{s_i\}_{i=1}^N, x_{N+1}, \left\{y_j^{(N+1)}\right\}_{j=1}^M\right) = \left\{y \in \left\{y_j^{(N+1)}\right\}_{j=1}^M : 1 - F\left(x_{N+1}, y, \left\{y_j^{(N+1)}\right\}_{j=1}^M\right) \leq \hat{s}\right\}.$$

Given  $y_{N+1}^* \in \left\{y_j^{(N+1)}\right\}_{j=1}^{\hat{r}}$ , we have

$$\left\{y_{N+1}^* \notin \mathcal{C}\left(\beta, \{s_i\}_{i=1}^N, x_{N+1}, \left\{y_j^{(N+1)}\right\}_{j=1}^M\right)\right\} = \{s(x_{N+1}, y_{N+1}^*) > \hat{s}\} = \{s_{N+1} > \hat{s}\}$$

In this case, the probability of miscoverage is

$$\mathbb{P}\left(\left\{y_{N+1}^* \notin \mathcal{C}\left(\beta, \{s_i\}_{i=1}^N, x_{N+1}, \{y_j^{(N+1)}\}_{j=1}^M\right)\right\} \mid y_{N+1}^* \in \{y_j^{(N+1)}\}_{j=1}^{\hat{r}}\right) = \mathbb{P}\left(s_{N+1} > \hat{s} \mid y_{N+1}^* \in \{y_j^{(N+1)}\}_{j=1}^{\hat{r}}\right).$$

We have

$$\mathbb{P}(s_{N+1} > \hat{s}) = 1 - \mathbb{P}(s_{N+1} \leq s_i) = 1 - \frac{i}{N+1}.$$

Finally, we achieve risk control and guarantee the upper bound of the correctness miscoverage rate

$$\begin{aligned} & \mathbb{P}\left(\left\{y_{N+1}^* \notin \mathcal{C}\left(\beta, \{s_i\}_{i=1}^N, x_{N+1}, \{y_j^{(N+1)}\}_{j=1}^M\right)\right\}\right) \\ &= \mathbb{P}\left(y_{N+1}^* \notin \{y_j^{(N+1)}\}_{j=1}^{\hat{r}}\right) + \mathbb{P}\left(s_{N+1} > \hat{s} \mid y_{N+1}^* \in \{y_j^{(N+1)}\}_{j=1}^{\hat{r}}\right) \\ &= 1 - \frac{\lceil(1-\alpha)(N+1)\rceil}{N} + \left(1 - \frac{\lceil(1-\beta)(N+1)\rceil}{N}\right) \cdot \frac{\lceil(1-\alpha)(N+1)\rceil}{N} \\ &= 1 - \frac{\lceil(1-\alpha)(N+1)\rceil}{N} \cdot \frac{\lceil(1-\beta)(N+1)\rceil}{N} \leq 1 - (1-\alpha)(1-\beta) \\ &\leq \alpha + \beta - \alpha\beta. \end{aligned}$$

If the sampling set of each test sample contains acceptable responses,  $\alpha \rightarrow 0$ , we have

$$\mathbb{P}\left(\left\{y_{N+1}^* \in \mathcal{C}\left(\beta, \{s_i\}_{i=1}^N, x_{N+1}, \{y_j^{(N+1)}\}_{j=1}^M\right)\right\}\right) \geq 1 - \beta.$$