



TEXAS  
The University of Texas at Austin

# Semantic Density: Uncertainty Quantification for Large Language Models through Confidence Measurement in Semantic Space

---

Xin Qiu<sup>1</sup>, Risto Miikkulainen<sup>1,2</sup>

<sup>1</sup>Cognizant AI Labs,

<sup>2</sup>The University of Texas at Austin

Neurips 2024

# Motivation

- **Trustworthiness of LLMs**

- Applications in safety-critical domains
- Unpredictable hallucinations/misinformation
- Lack of uncertainty/confidence indicator



# Motivation

- **Existing Solutions**

- **Extending traditional UQ methods to LLMs**
  - Only work for classification tasks
- **Train additional uncertainty/confidence predictors**
  - Not “off-the-shelf”, need extra data/training
- **Lexical uncertainty/confidence metrics**
  - Ignoring semantic information
- **Semantic Entropy**
  - Prompt-wise
  - No fine-grained semantic analysis



Response



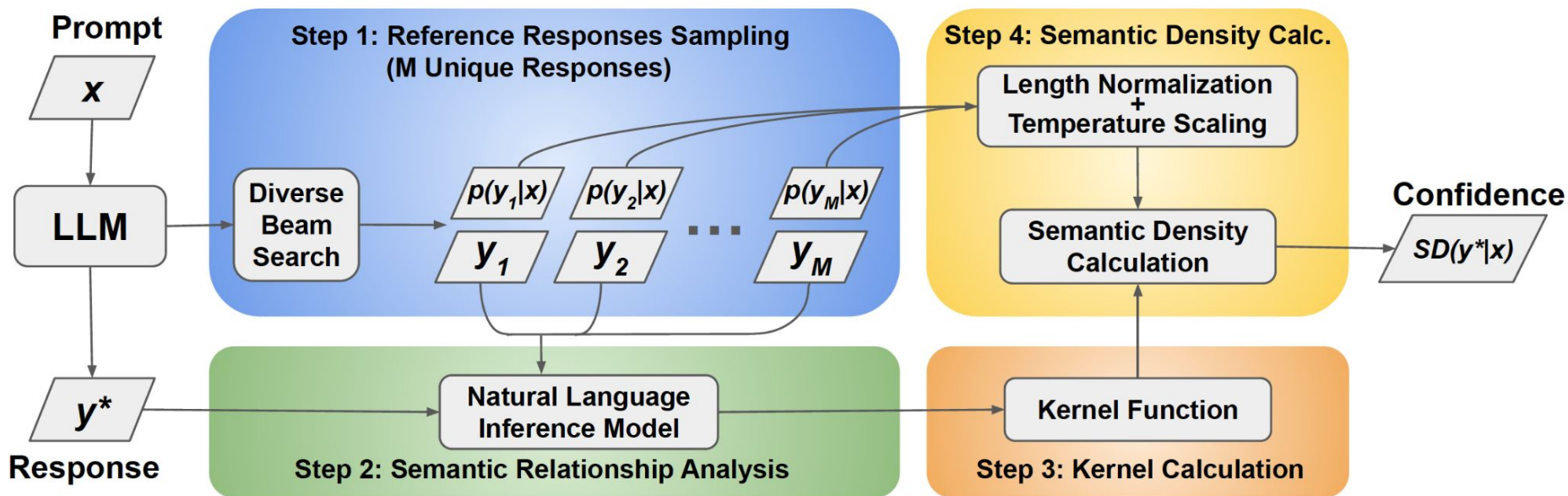
Uncertainty/  
Confidence

# Proposed Solution

- **Semantic Density**
  - Measuring confidence for each response
  - Analyzing output distribution in semantic space
  - **Main advantages:**
    - Response-wise indicator
    - Considering fine-grained semantic relationships
    - “Off-the-shelf”:
      - directly applicable to pre-trained LLMs
      - no fine-tuning/re-training
    - Work for free-form generation tasks

# Proposed Solution

- **Semantic Density**



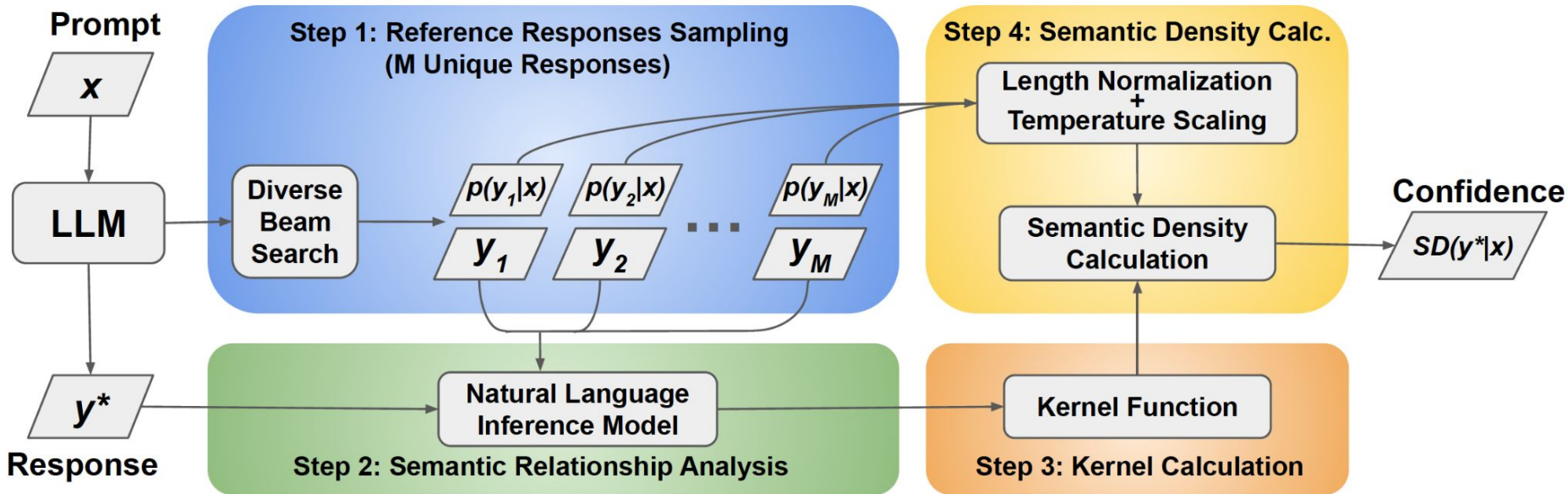
# Proposed Solution

- Semantic Density

**Step 2: Semantic Relationship Analysis**

Contradictory Prob.      Neutral Prob.      Entailment Prob.

$$\mathbb{E}(\|v_* - v_i\|^2) = 1^2 \cdot p_c(y_*, y_i|x) + \left(\frac{\sqrt{2}}{2}\right)^2 \cdot p_n(y_*, y_i|x) + 0^2 \cdot p_e(y_*, y_i|x)$$
$$= p_c(y_*, y_i|x) + \frac{1}{2} \cdot p_n(y_*, y_i|x)$$



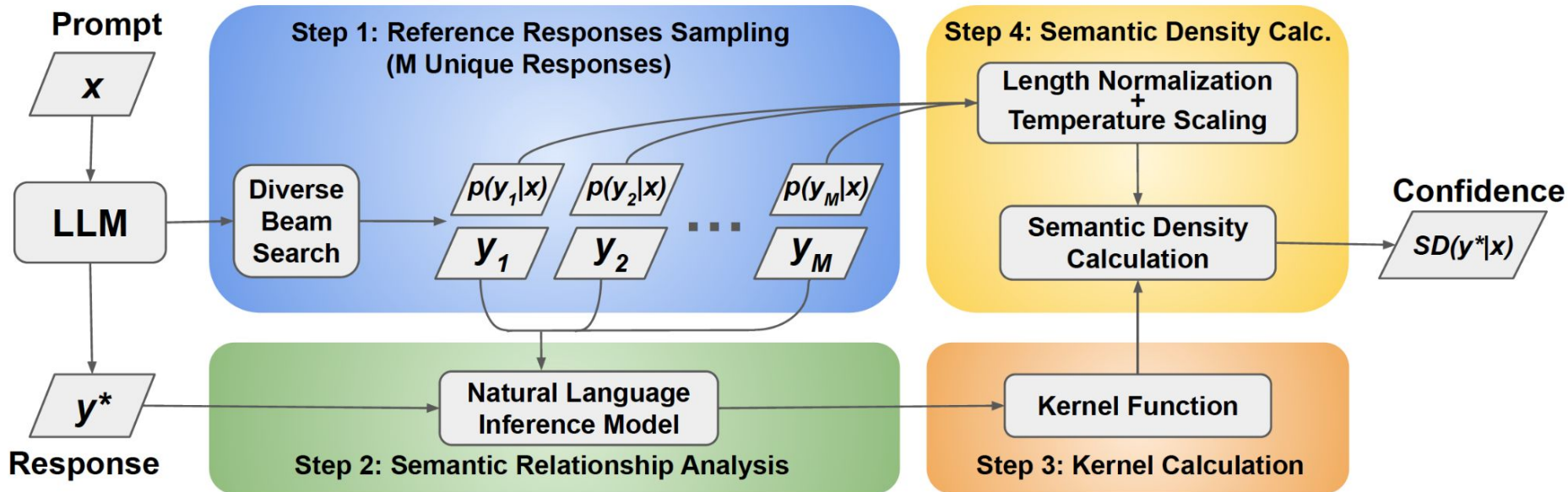
# Proposed Solution

- Semantic Density

## Step 3: Kernel Calculation

$$K(\mathbf{v}_* - \mathbf{v}_i) = (1 - \|\mathbf{v}_* - \mathbf{v}_i\|^2) \mathbf{1}_{\|\mathbf{v}_* - \mathbf{v}_i\| \leq 1}$$

Estimated by  $\mathbb{E}(\|\mathbf{v}_* - \mathbf{v}_i\|^2)$





# Proposed Solution

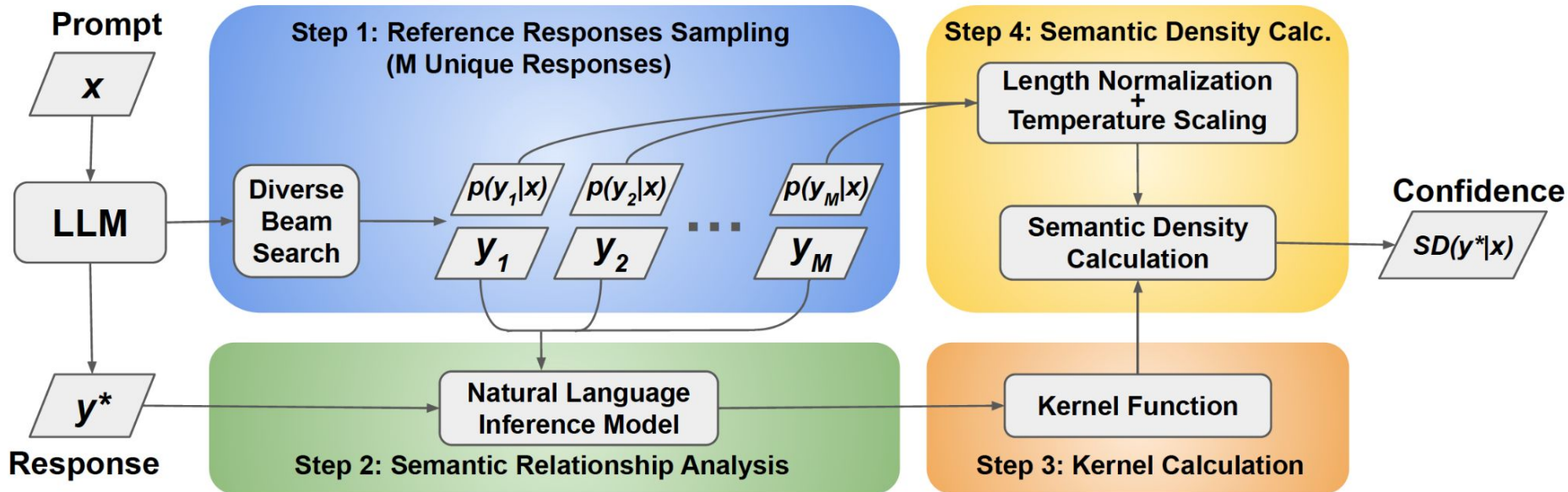
- Semantic Density

Step 4: Semantic Density Calc.

$$SD(y_*|x) = \frac{1}{\sum_{i=1}^M p(y_i|x)} \sum_{i=1}^M p(y_i|x) K(v_* - v_i)$$

Post-processing by  
Length-normalization  
+ Temperature Scaling

Sampling Probability  
by LLM





# Semantic Density

- **Semantic Space**

- **Oracle/idealized form:**

- **Contextual embedding:**  $v = E(\mathbf{y}|\mathbf{x})$

- **Fixed norm:**  $\|v\| = \frac{1}{2}$ , for  $v = E(\mathbf{y}|\mathbf{x})$ ,  $\forall \mathbf{x}, \mathbf{y}$

- **Reflecting semantic relationship:**

$$\|v_i - v_j\| = \begin{cases} 0, & \text{if } \mathbf{y}_i \text{ and } \mathbf{y}_j \text{ are semantically equivalent given context } \mathbf{x} \\ \frac{\sqrt{2}}{2}, & \text{if } \mathbf{y}_i \text{ and } \mathbf{y}_j \text{ are semantically irrelevant given context } \mathbf{x} \\ 1, & \text{if } \mathbf{y}_i \text{ and } \mathbf{y}_j \text{ are semantically contradictory given context } \mathbf{x} \end{cases}$$

$\|v_i - v_j\| < \|v_i - v_k\|$ , if  $\mathbf{y}_i$  is semantically closer to  $\mathbf{y}_j$  than to  $\mathbf{y}_k$ , given context  $\mathbf{x}$ .

# Semantic Density

- Semantic
  - Implicit form via NLI model
    - Natural Language Inference (NLI) Model:
      - Classifier for semantic relationships
    - Estimating the expectation of  $\|v_* - v_i\|$ :

$$\begin{aligned} \mathbb{E}(\|v_* - v_i\|^2) &= 1^2 \cdot p_c(\mathbf{y}_*, \mathbf{y}_i | \mathbf{x}) + \left(\frac{\sqrt{2}}{2}\right)^2 \cdot p_n(\mathbf{y}_*, \mathbf{y}_i | \mathbf{x}) + 0^2 \cdot p_e(\mathbf{y}_*, \mathbf{y}_i | \mathbf{x}) \\ &= p_c(\mathbf{y}_*, \mathbf{y}_i | \mathbf{x}) + \frac{1}{2} \cdot p_n(\mathbf{y}_*, \mathbf{y}_i | \mathbf{x}) \end{aligned}$$

**Contradictory Prob.**                      **Neutral Prob.**                      **Entailment Prob.**

# Semantic Density

- **Kernel Function Calculation**

- **Traditional Epanechnikov kernel in standard KDE:**

**Dimension Dependent** → 
$$K(\mathbf{v}) = \frac{\Gamma(2 + \frac{D}{2})}{\pi^{\frac{D}{2}}} (1 - \|\mathbf{v}\|^2) \mathbf{1}_{\|\mathbf{v}\| \leq 1}$$

- **Dimension-invariant kernel in semantic density:**

$$K(\mathbf{v}_* - \mathbf{v}_i) = (1 - \|\mathbf{v}_* - \mathbf{v}_i\|^2) \mathbf{1}_{\|\mathbf{v}_* - \mathbf{v}_i\| \leq 1}$$

**Estimated by**  
 $\mathbb{E}(\|\mathbf{v}_* - \mathbf{v}_i\|^2)$

- Does not work for standard KDE
- **Fits well in confidence/uncertainty estimation**
  - Making confidence scores comparable for embedding spaces with different dimensionalities

# Semantic Density

- **Semantic Density Calculation**

- **Expensive extension from standard KDE:**

$$\hat{p}(\mathbf{y}_*|\mathbf{x}) = \sum_{i=1}^M f_i K(\mathbf{v}_* - \mathbf{v}_i) = \frac{1}{\sum_{i=1}^M n_i} \sum_{i=1}^M n_i K(\mathbf{v}_* - \mathbf{v}_i)$$

Number of Occurrences

- **Cost-effective realization in semantic density:**

$$SD(\mathbf{y}_*|\mathbf{x}) = \frac{1}{\sum_{i=1}^M p(\mathbf{y}_i|\mathbf{x})} \sum_{i=1}^M p(\mathbf{y}_i|\mathbf{x}) K(\mathbf{v}_* - \mathbf{v}_i)$$

Post-processing by Length-normalization + Temperature Scaling

Sampling Probability by LLM

# Empirical Evaluations

- **Indicator for correctness of responses**
  - Area under the receiver operating characteristic curve (AUROC):

AUROC	SD	CoQA					
		SE[22]	P(True)[21]	Deg[26]	NL[31]	NE[29]	PE[21]
Llama-2-13B	<b>0.783</b>	0.633	0.594	0.734	0.709	0.629	0.647
Llama-2-70B	<b>0.783</b>	0.621	0.576	0.721	0.716	0.617	0.647
Llama-3-8B	0.738	0.599	0.593	<b>0.795</b>	0.676	0.608	0.604
Llama-3-70B	<b>0.789</b>	0.608	0.670	0.729	0.698	0.587	0.641
Mistral-7B	<b>0.788</b>	0.627	0.667	0.737	0.704	0.614	0.632
Mixtral-8x7B	<b>0.786</b>	0.626	0.589	0.728	0.708	0.617	0.651
Mixtral-8x22B	<b>0.791</b>	0.614	0.614	0.726	0.700	0.604	0.649

# Empirical Evaluations

- **Indicator for correctness of responses**
  - Area under the receiver operating characteristic curve (AUROC):

AUROC	SD	SE	TriviaQA				
			P(True)	Deg	NL	NE	PE
Llama-2-13B	<b>0.848</b>	0.672	0.589	0.824	0.675	0.574	0.556
Llama-2-70B	<b>0.829</b>	0.677	0.556	0.787	0.714	0.582	0.566
Llama-3-8B	<b>0.866</b>	0.662	0.647	0.796	0.834	0.636	0.622
Llama-3-70B	<b>0.828</b>	0.663	0.654	0.764	<b>0.828</b>	0.611	0.596
Mistral-7B	<b>0.866</b>	0.690	0.589	0.828	0.745	0.615	0.536
Mixtral-8x7B	<b>0.846</b>	0.685	0.562	0.797	0.795	0.644	0.605
Mixtral-8x22B	<b>0.829</b>	0.686	0.604	0.762	0.801	0.644	0.607

# Empirical Evaluations

- **Indicator for correctness of responses**
  - Area under the receiver operating characteristic curve (AUROC):

AUROC	SD	SE	SciQ		NL	NE	PE
			P(True)	Deg			
Llama-2-13B	<b>0.757</b>	0.570	0.572	0.727	0.693	0.513	0.574
Llama-2-70B	<b>0.746</b>	0.643	0.584	0.713	0.637	0.554	0.615
Llama-3-8B	<b>0.780</b>	0.611	0.564	0.731	0.686	0.597	0.651
Llama-3-70B	<b>0.771</b>	0.613	0.556	0.706	0.724	0.558	0.520
Mistral-7B	<b>0.771</b>	0.618	0.568	0.736	0.669	0.565	0.528
Mixtral-8x7B	<b>0.773</b>	0.612	0.585	0.716	0.726	0.612	0.658
Mixtral-8x22B	<b>0.775</b>	0.620	0.602	0.719	0.715	0.602	0.628



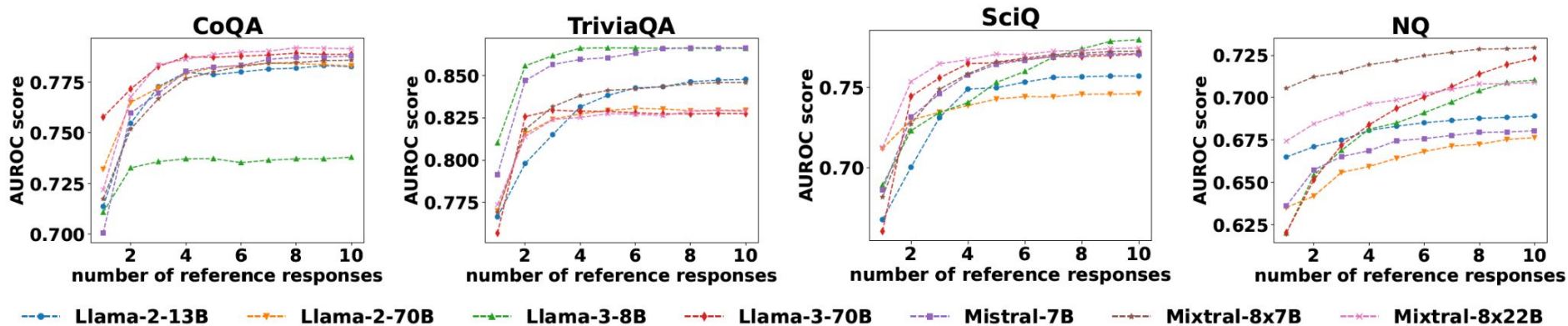
# Empirical Evaluations

- **Indicator for correctness of responses**
  - Area under the receiver operating characteristic curve (AUROC):

AUROC	SD	SE	NQ		NL	NE	PE
			P(True)	Deg			
Llama-2-13B	<b>0.689</b>	0.581	0.592	0.686	0.588	0.571	0.640
Llama-2-70B	0.676	0.545	0.531	<b>0.691</b>	0.567	0.573	0.620
Llama-3-8B	<b>0.710</b>	0.583	0.517	0.706	0.601	0.603	0.615
Llama-3-70B	<b>0.723</b>	0.577	0.643	0.714	0.631	0.603	0.615
Mistral-7B	<b>0.680</b>	0.597	0.523	0.676	0.640	0.635	0.631
Mixtral-8x7B	<b>0.729</b>	0.599	0.576	0.720	0.654	0.603	0.608
Mixtral-8x22B	<b>0.709</b>	0.577	0.504	0.704	0.638	0.625	0.680

# Empirical Evaluations

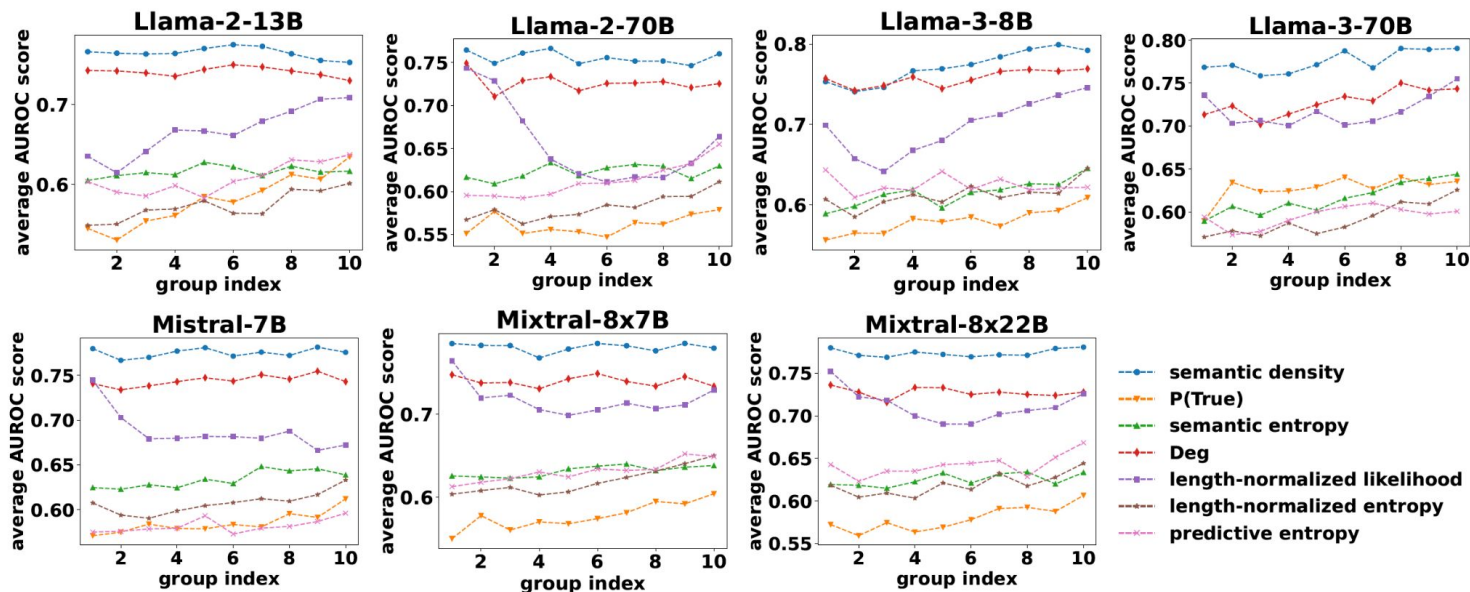
- Sensitivity to number of reference responses:



- Robust to reduced number of reference samples

# Empirical Evaluations

- Performance changes for responses from different sampling strategies:



- Robust for both greedy or diverse samplings

# Summary

## ● Conclusions

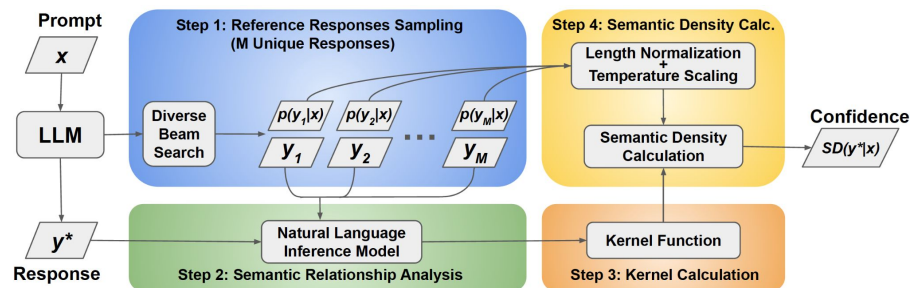
- Proposed semantic density as a confidence indicator for LLM responses
- Response-wise, off-the-shelf for free-form generation tasks

## ● Future work

- Improving sampling strategy for reference responses
- Developing contextual embedding model
- Dedicated kernel function
- Better token probability calibration

## ● Source code and contact

- <https://github.com/cognizant-ai-labs/semantic-density-paper>
- [qiuxin.nju@gmail.com](mailto:qiuxin.nju@gmail.com), [risto@cs.utexas.edu](mailto:risto@cs.utexas.edu)





Thank You

