

Hallucination Detection

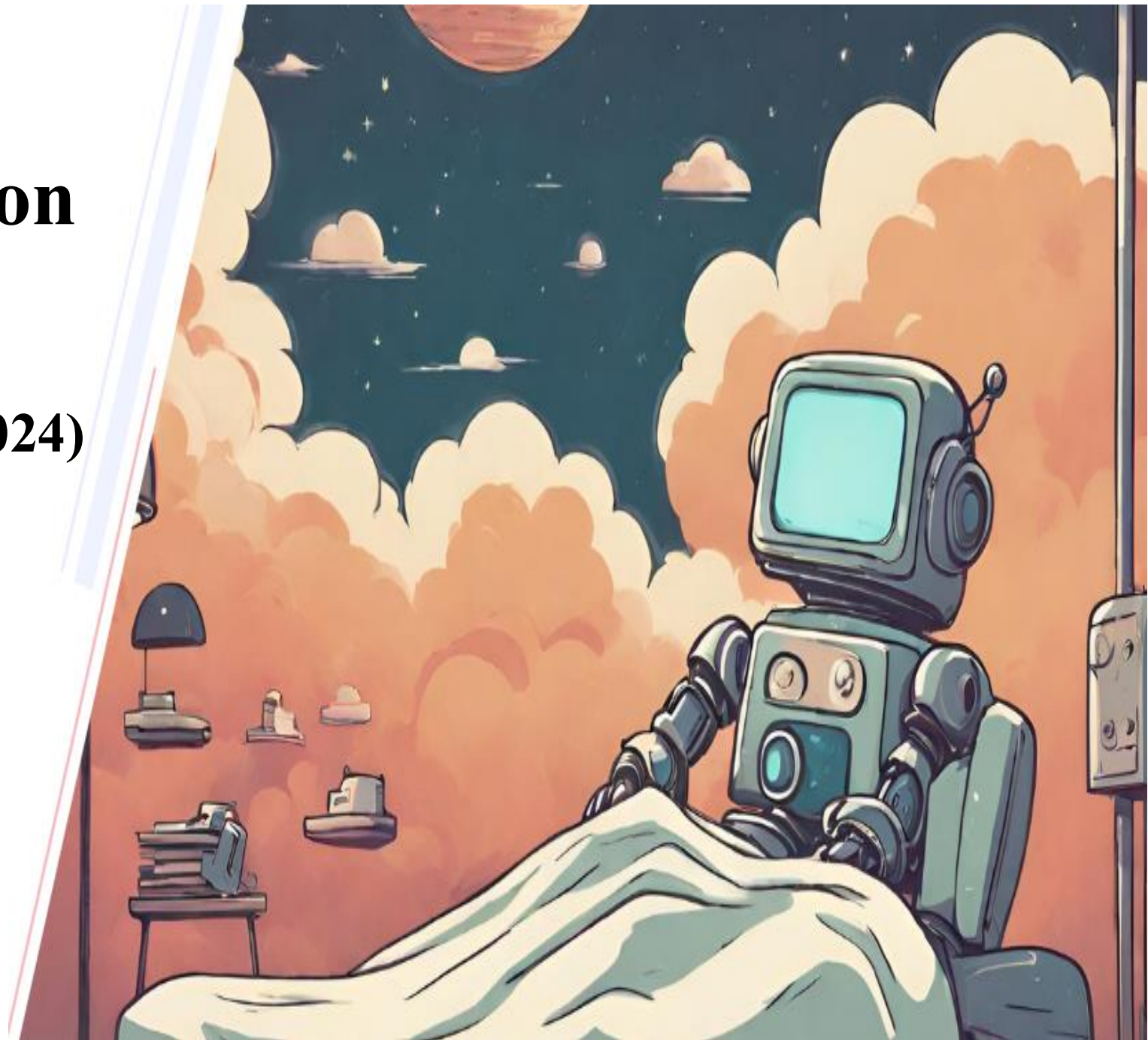
(ELOQUENT Lab @ CLEF 2024)

Philipp Schaer and Narjes Nikzad

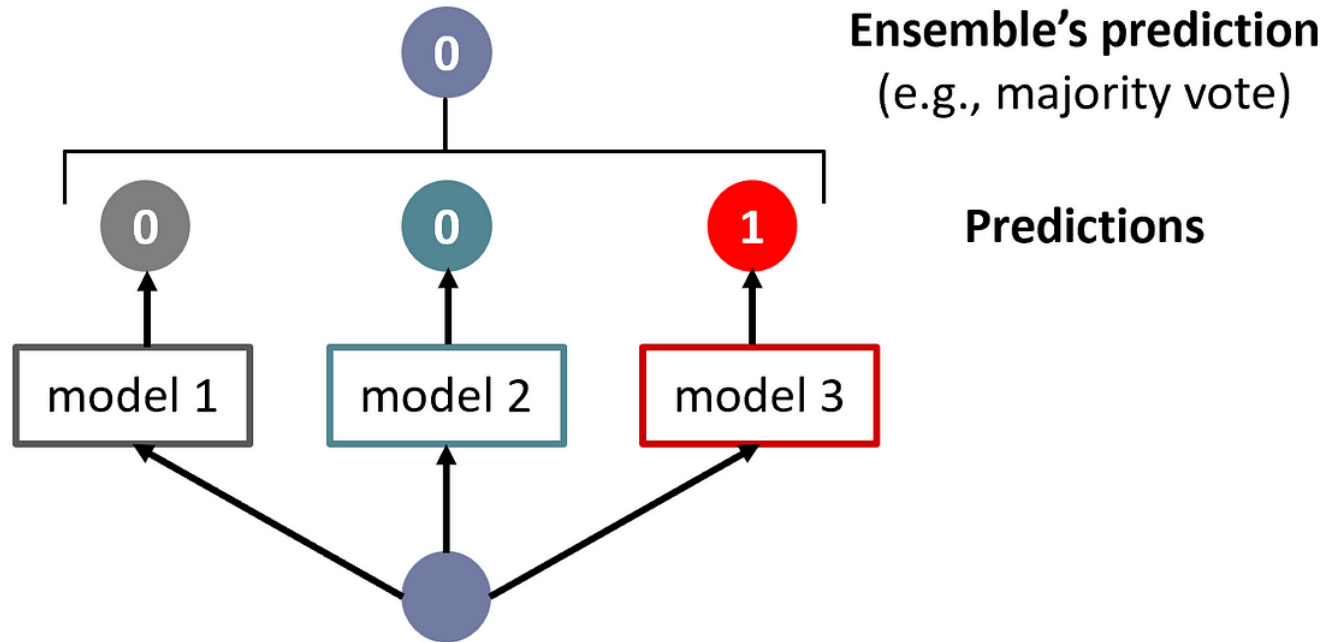
08-04-24 – Cologne, Germany
<https://ir.web.th-koeln.de>



Technology
Arts Sciences
TH Köln



Simple voting approach



Inverse Proportion approach for calculating weights

Goal: **assigning higher weights to LLMs with better F1-scores.**

- ✓ Split your trail dataset into training (for cot or few-shot learning) and validation sets.
- ✓ Evaluate the performance of each LLM on the validation set using F1-score.
- ✓ Calculate Weights:
 - Calculate the inverse of each LLM's F1-score: $\text{weight}_i = 1 / (1 - \text{F1_score}_i)$ (where i represents the LLM).
 - Normalize the weights so they sum to 1: $\text{normalized_weight}_i = \text{weight}_i / \text{sum}(\text{weights})$.



Inverse Proportion approach for calculating weights

✓ Calculate Inverse Proportions:

- $\text{weight_GPT4} = 1 / (1 - 0.85) = 6.666$
- $\text{weight_LaMDA3} = 1 / (1 - 0.78) = 4.545$
- $\text{weight_Gemma} = 1 / (1 - 0.82) = 5.555$

✓ Sum the Weights:

$\text{total_weight} = \text{weight_GPT4} + \text{weight_LaMDA3} + \text{weight_Gemma}$
 $\text{total_weight} = 6.666 + 4.545 + 5.555 = 16.766$

✓ Normalize Weights:

$\text{normalized_weight_GPT4} = \text{weight_GPT4} / \text{total_weight}$
 $\text{normalized_weight_GPT4} = 6.666 / 16.766 \approx 0.3975$
 $\text{normalized_weight_LaMDA3} \approx 0.271$
 $\text{normalized_weight_Gemma} \approx 0.3313$