

Determine-Then-Ensemble: Necessity of Top-K Union for Large Language Model Ensembling

Yuxuan Yao, Han Wu[†], Mingyang Liu, Sichun Luo, Xiongwei Han, Jie Liu, Zhijiang Guo, Linqi Song[†]

Motivation

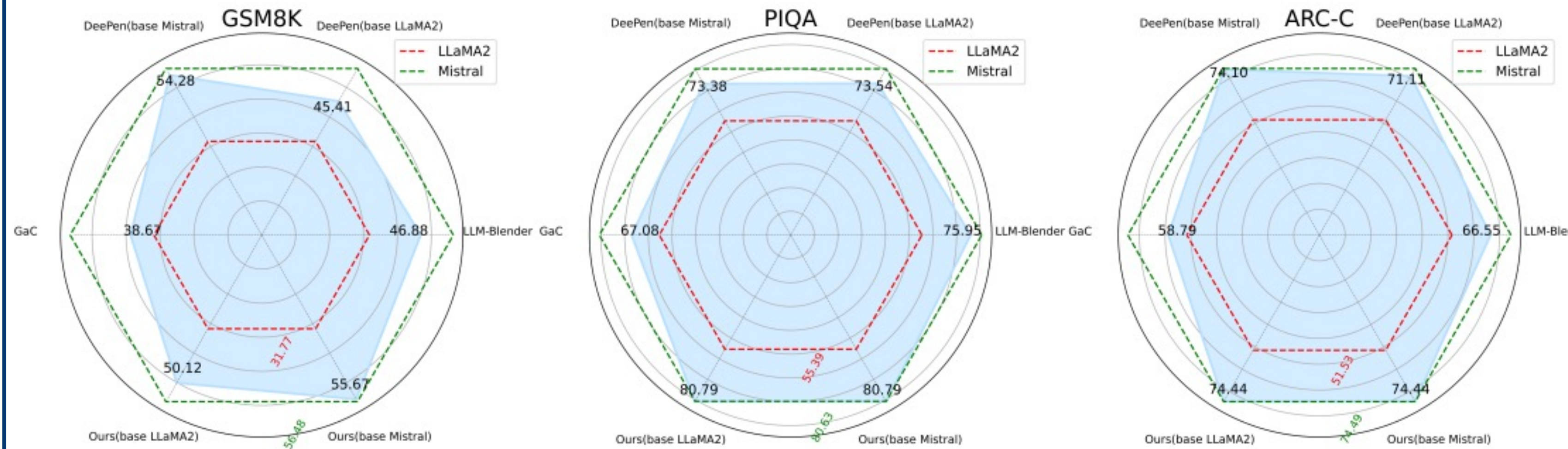
Existing model ensembling methods can be broadly categorized into three types: **output-level**, **probability-level**, and **training-level** approaches. Output-level methods aggregate the complete outputs of multiple candidate models. Probability-level methods, integrate outputs based on probability distributions at each generation step through the intersection or union of the vocabulary. Training-level methods utilize output probability vectors as labels for richer information extraction during training. While output-level methods are constrained by the limitations of existing outputs, and training-level methods introduce additional computational overhead, probability-level methods have garnered particular attention.

Challenges

- First, These approaches concentrate solely on the ensembling technique, sidestepping the crucial discussion of **which types of models can be effectively combined**.
- Second, these methods tend to align the probabilities across the **entire vocabulary at each generation step**. Such a strategy introduces substantial computational overhead during inference, which hinders performance and efficiency

Key Observation

We explored various factors that might affect the performance of model ensembling, including model size (e.g., 3B/7B/8B/13B/70B), model architecture (dense/sparse), performance discrepancies, tokenization strategies (BPE/WordPiece), vocabulary size (e.g., 102K/64K/32K) and tasks variations (e.g., text generation/QA/multiple choices). Finally, we identified three representative factors for further analyses, including **performance discrepancy**, **vocabulary size**, and **tasks variations**.



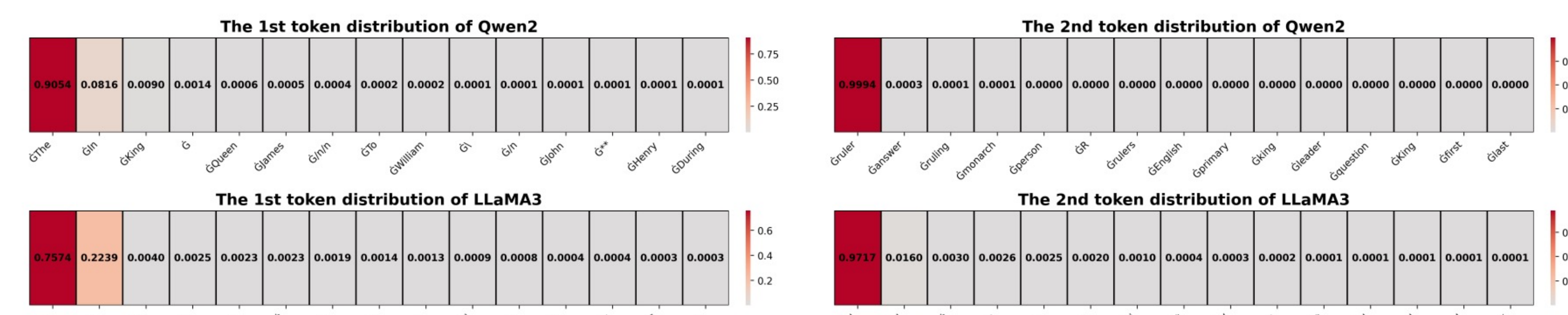
TAKEAWAY I:

Smaller performance gaps lead to greater gains from model ensembling.

Methods	Dataset		Methods	Dataset		Methods	Dataset	
	GSM	PIQA		MMLU	ARC-C		NQ	ARC-C
DeepSeek-7B(102K)	59.67	72.66	DeepSeek-7B(102K)	46.97	58.73	Mistral-7B(32K)	24.25	74.49
Mistral-7B(32K)	56.48	80.63	LLaMA2-13B(32K)	49.61	51.53	Yi-6B(64K)	22.55	73.21
LLM-BLENDER	61.16	76.54	LLM-BLENDER	48.86	57.84	LLM-BLENDER	22.97	72.48
DEEPEN(DeepSeek-7B)	(+1.49)	(-4.09)	DEEPEN(DeepSeek-7B)	(-0.75)	(-0.89)	DEEPEN(Mistral-7B)	(-1.28)	(-2.01)
DEEPEN(Mistral-7B)	55.00	77.65	DEEPEN(DeepSeek-7B)	52.81	60.00	DEEPEN(Mistral-7B)	25.26	76.50
DEEPEN(Mistral-7B)	(-4.67)	(+4.99)	DEEPEN(DeepSeek-7B)	(+5.84)	(+1.27)	DEEPEN(Mistral-7B)	(+1.01)	(+1.01)
DEEPEN(Mistral-7B)	61.28	76.32	DEEPEN(LLaMA2-13B)	54.09	62.39	DEEPEN(Yi-6B)	22.80	77.86
DEEPEN(Mistral-7B)	(+4.80)	(-4.31)	DEEPEN(LLaMA2-13B)	(+4.48)	(+10.86)	DEEPEN(Yi-6B)	(+0.25)	(+4.65)
Ours(DeepSeek-7B)	62.77	81.81	Ours(DeepSeek-7B)	48.90	60.16	Ours(Mistral-7B)	24.76	76.54
Ours(DeepSeek-7B)	(+3.10)	(+9.15)	Ours(DeepSeek-7B)	(+1.93)	(+1.43)	Ours(Mistral-7B)	(+0.51)	(+2.05)
Ours(Mistral-7B)	58.88	81.81	Ours(Mistral-7B)	48.90	60.16	Ours(Yi-6B)	23.28	76.54
Ours(Mistral-7B)	(+2.40)	(+1.18)	Ours(Mistral-7B)	(-0.71)	(+8.63)	Ours(Yi-6B)	(+0.73)	(+3.33)

TAKEAWAY II:

The influence of vocabulary size for model ensembling is marginal.



Question	who wrote he ain't heavy he's my brother lyrics who was the ruler of england in 1616
Gold answer	Bob Russell James I
LLaMA3 response	Bob Russell ✓ James I ✓
Qwen2 response	He Ain't Heavy, He's My Brother was written by Mike D'Abo . ✗ The song was originally recorded by the British singer-songwriter and released as a single in 1969. It became a hit and has since been covered by numerous artists. ... The ruler of England in 1616 was King James I . ✓ He was the first monarch of the Stuart dynasty and reigned from March 24, 1603, until his death on July 19, 1625. ...

TAKEAWAY III:

Even though the performance and vocabulary size are aligned across models, substantial differences in response style could also hinder successful ensembling.

Methodology

Algorithm 1 Union top- k ensembling

Require: LLM_i, Vocabulary V_i, Tokenizer T_i, Demonstration *prompt*, Stop condition *stop*(*)

```

1: while not stop(*) do
2:   TopKi, Pi ← LLMi(prompt)           ▷ Stopping criteria
3:   for each model do
4:     if token  $w \in V^u$  and  $w \in \text{TopK}_i$  then
5:        $P_{\{w\}}^i$  and TopKi remains unchanged.
6:     else if token  $w \in V^u$  and  $w \in V^i$  and  $w \notin \text{TopK}_i$  then
7:       TopKi ←  $w$ ,  $\hat{P}^i \leftarrow p_w^i$ 
8:     else if token  $w \in V^u$  and  $w \notin V^i$  then
9:        $w^1, \dots, w^m \leftarrow T^i(w)$ 
10:      TopKi ←  $w^1$ ,  $\hat{P}^i \leftarrow p_{w^1}^i$ 
11:    end if
12:  end for
13:   $\hat{P}_{norm}^i = \text{softmax}(\hat{P}^i)$ ,  $\hat{P}_{avg} = \frac{1}{n} \sum_{i=1}^n \hat{P}_{norm}^i$ 
14:   $w = \text{argmax}_{w \in \text{TopK}_i} (\hat{P}_{avg})$            ▷ Predict the next token
15:  prompt ← prompt + w                     ▷ Update input sequence
16: end while

```

Experiments

- (1)UNITE enhances individual model performance when the base models exhibit similar performance levels.
- (2)UNITE demonstrates greater robustness and generality.
- (3)Collaborating with comparable LLMs does not necessarily yield better results

Method	Dataset				Avg.
	GSM8K	PIQA	ARC-C	MMLU	
LLaMA3	78.77	79.08	79.01	64.58	75.36
LLaMA3.1	80.83	82.86	79.49	66.69	77.47
Qwen2	80.78	84.57	84.92	64.96	78.81
Two-model ensembling (LLaMA3+Qwen2)					
LLM-BLENDER	82.69 (+1.91)	82.53 (-2.04)	82.98 (-1.94)	62.07 (-2.89)	77.57 (-1.24)
DEEPEN	- OOM -				
GAC	80.67 (-0.11)	80.96 (-3.61)	84.93 (+0.01)	67.05 (+2.09)	78.40 (-0.41)
UNITE	84.17 (+3.39)	85.53 (+0.96)	85.07 (+0.15)	69.78 (+4.82)	81.14 (+2.33)
Three-model ensembling					
LLM-BLENDER	83.30(+2.52)	83.47(-1.10)	83.48(-1.44)	62.55(-2.41)	78.20(-0.61)
DEEPEN	- OOM -				
UNITE	84.99 (+4.21)	84.98 (+0.41)	85.39 (+0.47)	69.12 (+4.16)	81.12 (+2.31)

