

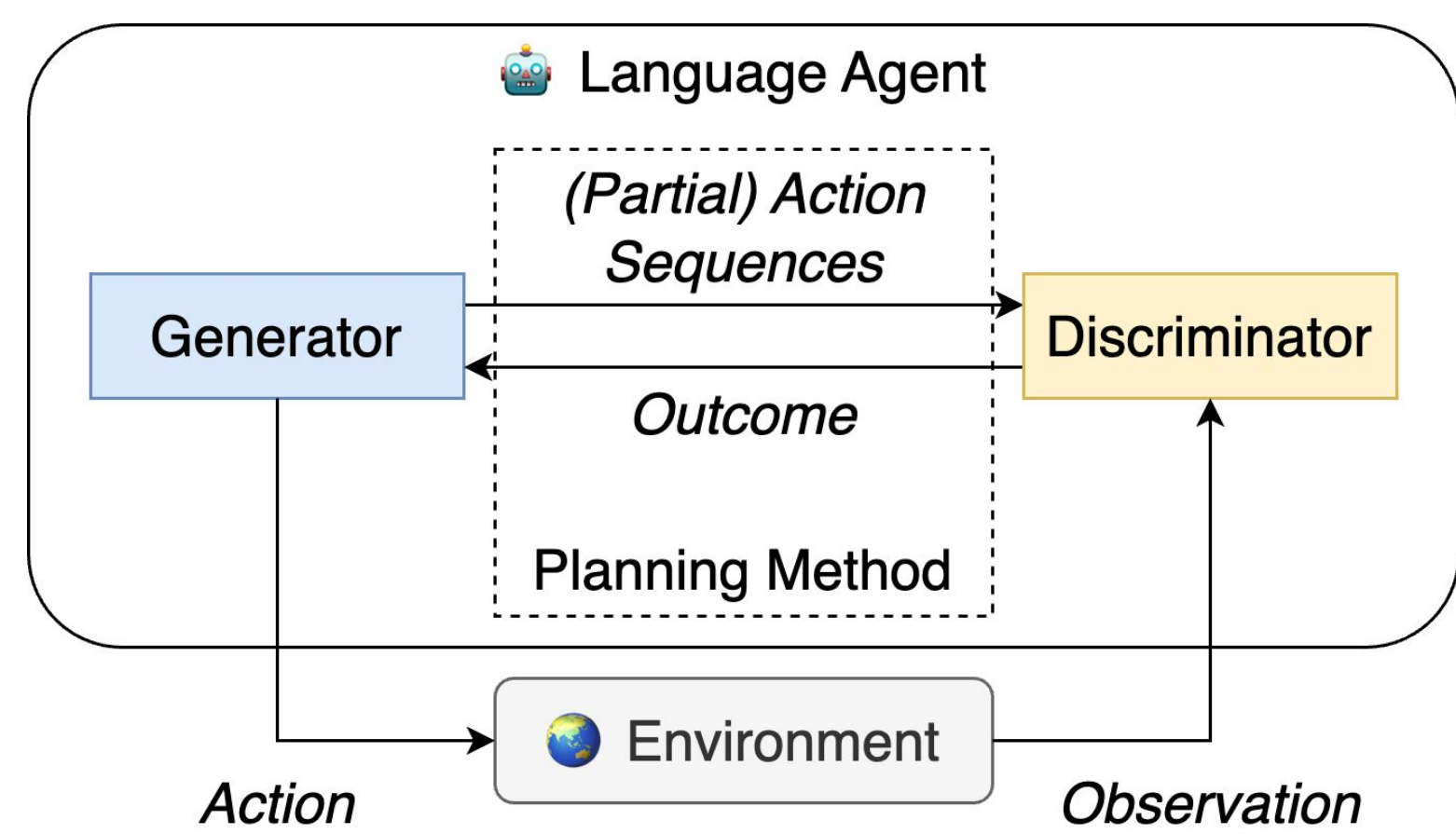
# In LLM planning, the discriminator needs to get up to 90% accuracy for tree search to start outperforming simple re-ranking.

## When is Tree Search Useful for LLM Planning? It Depends on the Discriminator

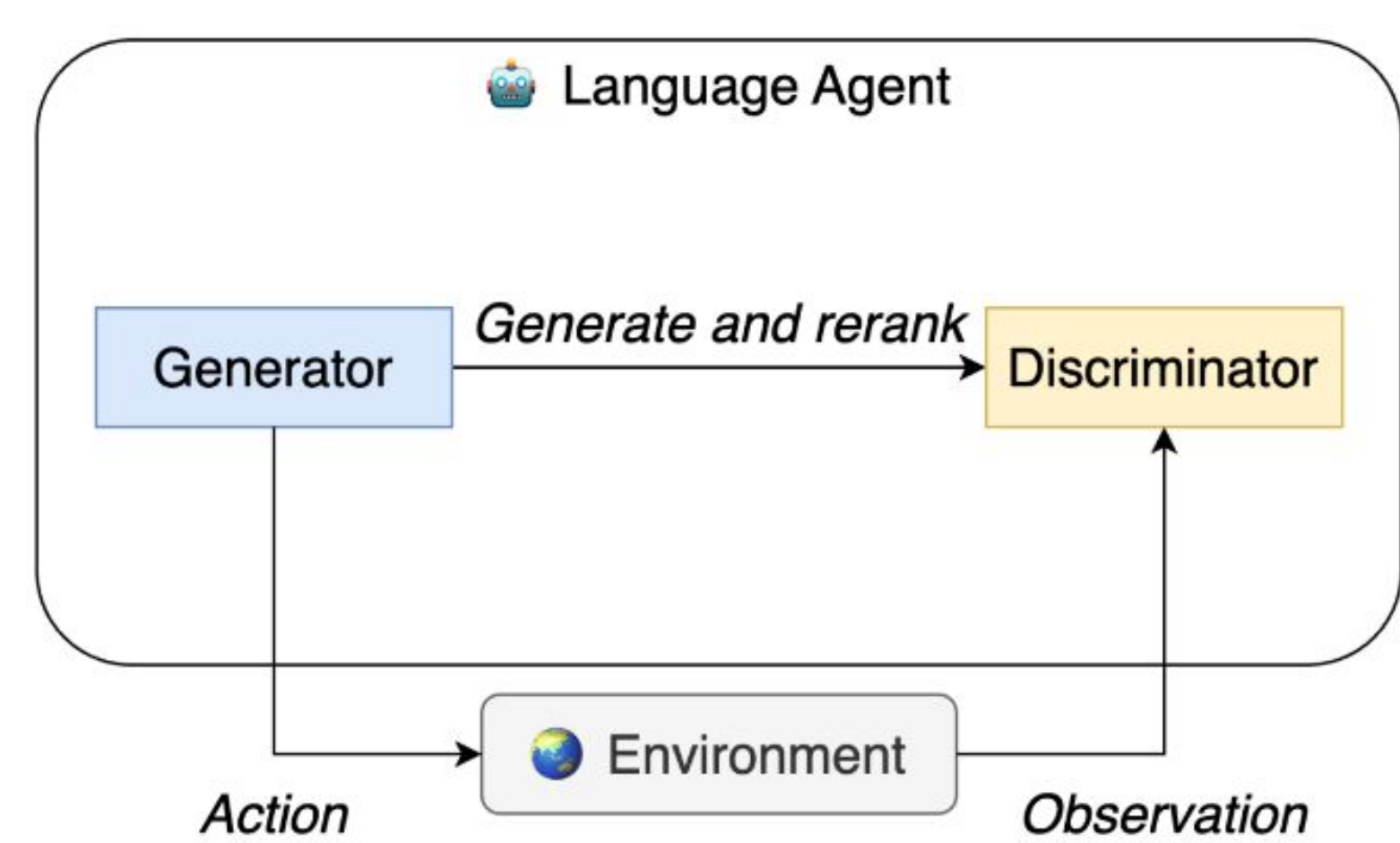
Ziru Chen, Michael White, Ray Mooney, Ali Payani, Yu Su, Huan Sun



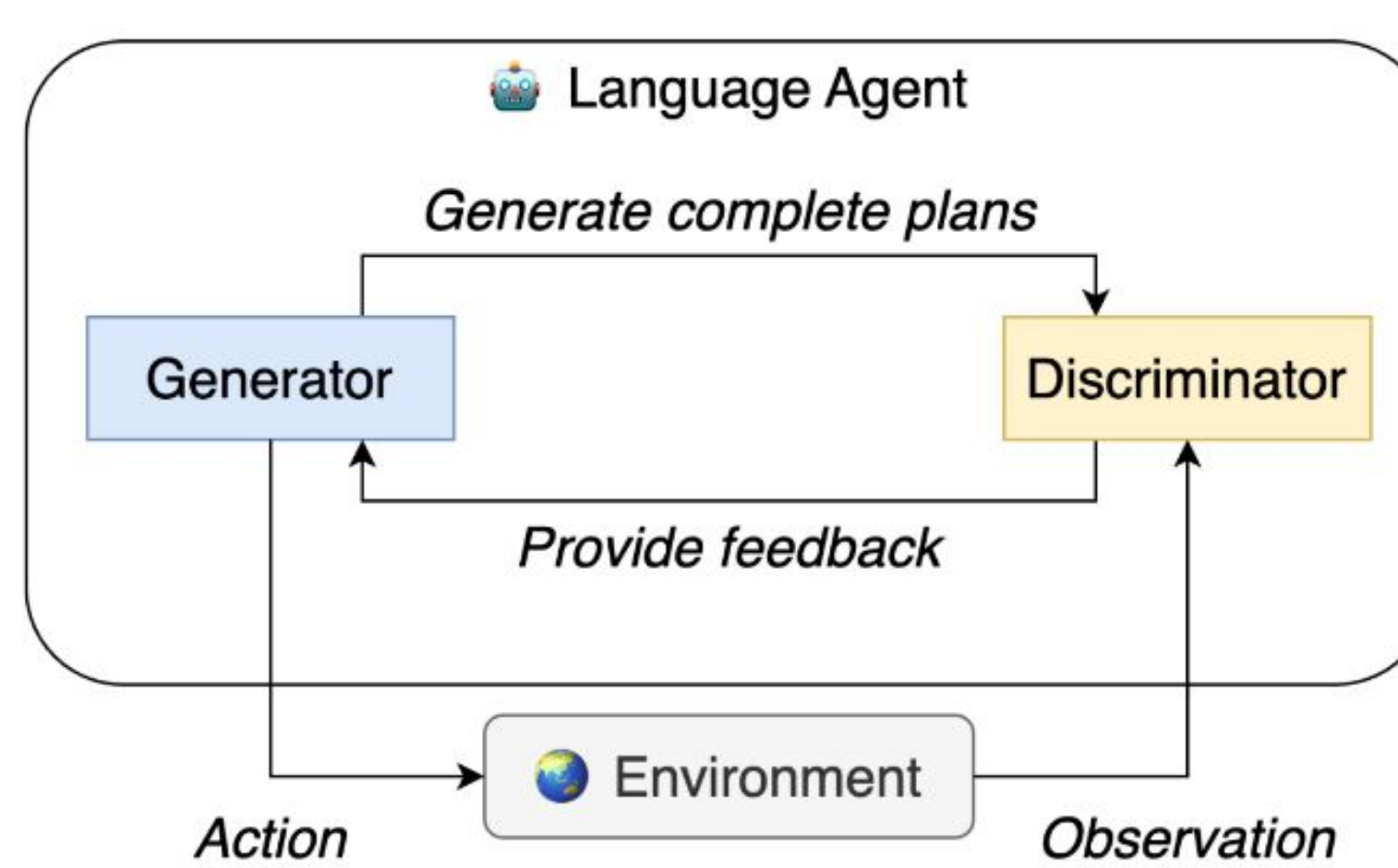
### A Generator-Discriminator Framework of Language Agents



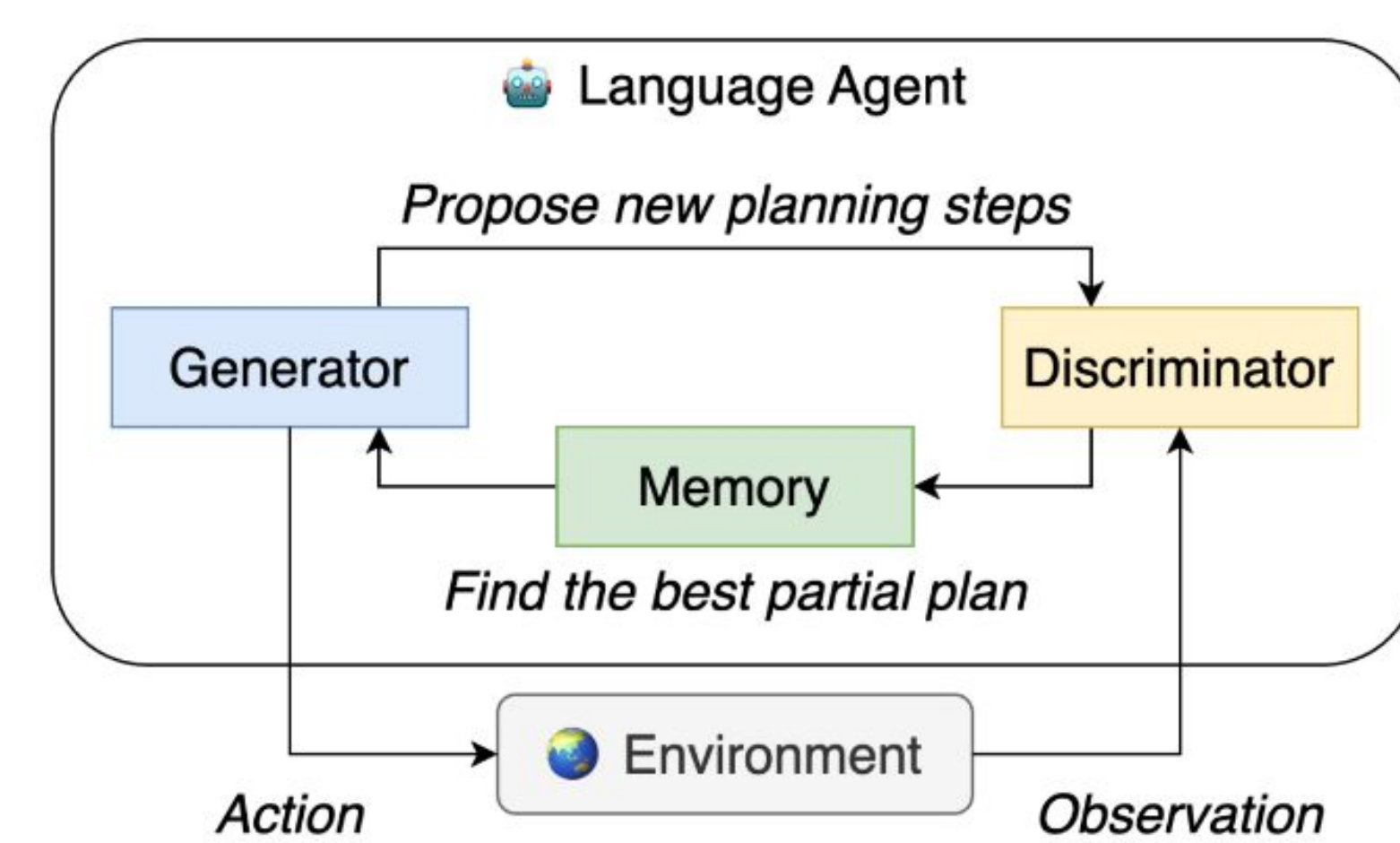
- **Generator:** Propose (partial) action sequences.
- **Discriminator:** Evaluates the outcomes of these actions.
- **Planning method:** Ranks the actions according to their outcomes and manage the interaction between the two models.



(a) Re-ranking.

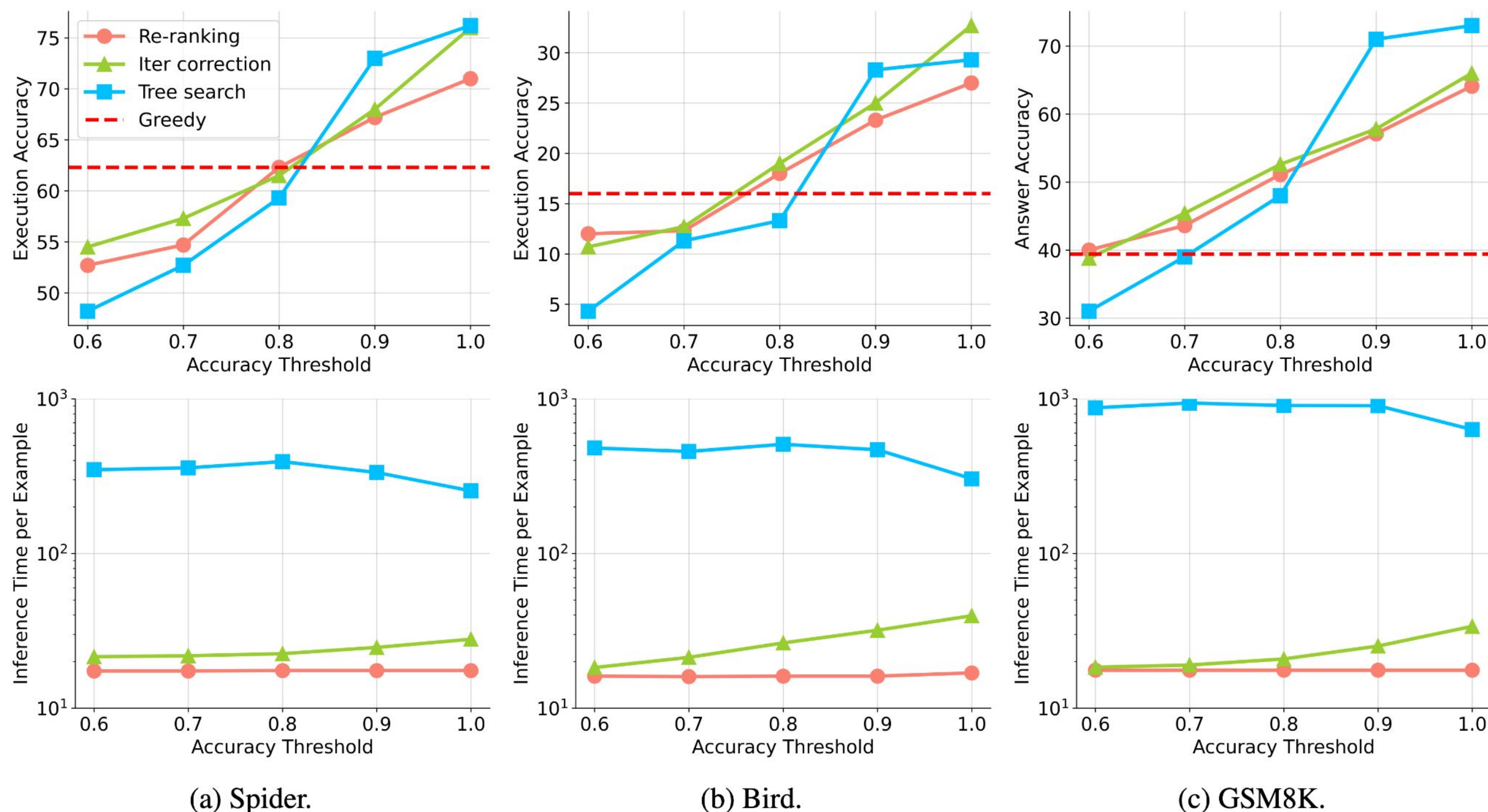


(b) Iterative Correction.



(c) Tree Search.

### Simulation Experiments with Oracle



End-to-end evaluation results (the first row) and average inference time in log scale (the second row) of our simulation experiments with oracle.

### Intrinsic Evaluation of LLM-Based Discriminators

Models	Spider				Bird				GSM8K <sup>‡</sup>			
	Acc	F1	H@1	MRR	Acc	F1	H@1	MRR	Acc	F1	H@1	MRR
CodeLlama-7B	54.0	37.1	56.0	62.3	44.6	<b>46.7</b>	13.0	18.0	48.6	<b>38.7</b>	36.2	46.9
CodeLlama-13B	58.2	37.1	57.0	63.1	49.4	<b>46.7</b>	12.7	18.3	<b>62.2</b>	<b>38.7</b>	<b>41.8</b>	<b>51.0</b>
CodeLlama-7B-FT	62.4	60.3	59.5	64.6	52.4	<b>46.7</b>	14.3	19.1	-	-	-	-
CodeLlama-13B-FT	<b>69.7</b>	<b>67.2</b>	<b>61.3</b>	<b>65.7</b>	<b>62.1</b>	<b>46.7</b>	<b>16.0</b>	<b>20.5</b>	-	-	-	-
GPT-3.5-Turbo	67.0	47.3	59.0	64.3	64.3	35.7	16.0	20.5	72.1	49.1	46.6	54.0
GPT-4-Turbo	<b>76.5</b>	<b>54.9</b>	<b>63.0</b>	<b>66.7</b>	<b>76.2</b>	<b>50.1</b>	<b>20.3</b>	<b>23.0</b>	<b>93.8</b>	<b>91.1</b>	<b>59.8</b>	<b>61.6</b>

Intrinsic evaluation results of naive LLMs' discrimination abilities.

	CodeLlama-13B			GPT-3.5-Turbo			CodeLlama-13B-FT	
	Spider	Bird	GSM8K	Spider	Bird	GSM8K	Spider	Bird
Naive Discriminator	58.2	49.4	62.2	67.0	64.3	72.1	69.7	62.1
+ Executability Check	78.7	78.8	64.5	84.8	86.3	73.2	83.6	82.2
++ Execution Result	<b>83.6</b>	<b>79.6</b>	<b>70.6</b>	<b>90.0</b>	<b>89.2</b>	<b>76.5</b>	<b>88.5</b>	<b>85.1</b>
Improvement	<u>25.4</u>	<u>30.2</u>	<u>8.4</u>	23.0	24.9	4.4	18.8	23.0

Discrimination accuracy of observation-enhanced LLMs. The best performance is achieved using both kinds of environmental observations.

### Experiments with LLM-Based Discriminators

Discriminators	Spider (Greedy Gen = 62.3)			Bird (Greedy Gen = 16.0)		
	Re-ranking	Iter. Correct.	Tree Search	Re-ranking	Iter. Correct.	Tree Search
CodeLlama-13B	<b>57.5</b>	51.7	55.5	<b>13.3</b>	<b>13.3</b>	<b>13.3</b>
GPT-3.5-Turbo	<b>58.3</b>	52.7	56.2	<b>18.0</b>	17.3	14.0
CodeLlama-13B-FT	<b>61.5</b>	51.7	56.0	<b>14.3</b>	13.0	13.0
CodeLlama-13B <sup>E</sup>	<b>65.5</b>	62.0	62.5	21.0	<b>24.3</b>	22.7
GPT-3.5-Turbo <sup>E</sup>	67.0	<b>67.5</b>	66.0	22.3	<b>25.0</b>	22.7
CodeLlama-13B-FT <sup>E</sup>	<b>70.3</b>	68.0	67.5	23.7	<b>26.3</b>	21.7
Oracle Simulation ( $\tau = 1.0$ )	71.0	76.0*	76.2*	27.0	32.7*	29.3

End-to-end execution accuracy on text-to-SQL parsing.

Error Type	Spider		Bird		GSM8K	
	Iter. Correct.	Tree Search	Iter. Correct.	Tree Search	Iter. Correct.	Tree Search
Discrimination	29 (78.4%)	17 (60.7%)	9 (52.9%)	12 (50.0%)	30 (62.5%)	6 (66.7%)
Exploration	8 (21.6%)	11 (39.3%)	8 (47.1%)	12 (50.0%)	18 (37.5%)	3 (33.3%)
Total	37	28	17	24	48	9

Error analysis of examples where re-ranking outperforms advanced planning methods.

- (1) **Discrimination error:** The discriminator assigns a higher score for wrong programs than correct ones, which is not recoverable by any planning method.
- (2) **Exploration error:** The planning method has not found the correct program before termination.

### Conclusions

- Advanced planning methods, i.e., iterative correction and tree search, demand highly accurate discriminators (up to 90% accuracy) to achieve decent improvements over the simpler method, re-ranking.
- Using environmental feedback, we improve the discrimination accuracy of LLMs. Yet, our end-to-end evaluations suggest they have barely met the need for advanced planning methods to show significant improvements over re-ranking.
- Advanced planning methods may not adequately balance accuracy and efficiency when using LLM-based discriminators. In our experiments, compared to the other two methods, tree search is at least 10–20 times slower but leads to negligible performance gains.