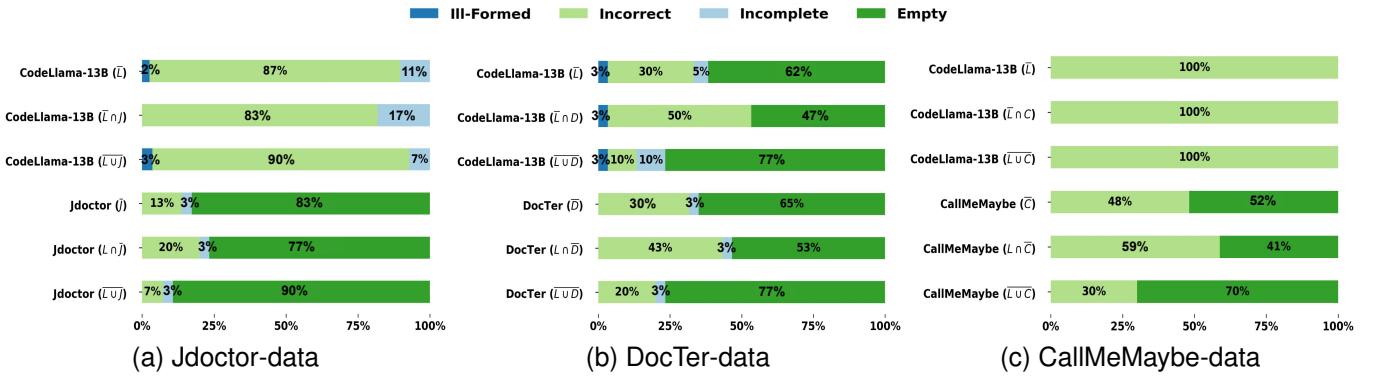
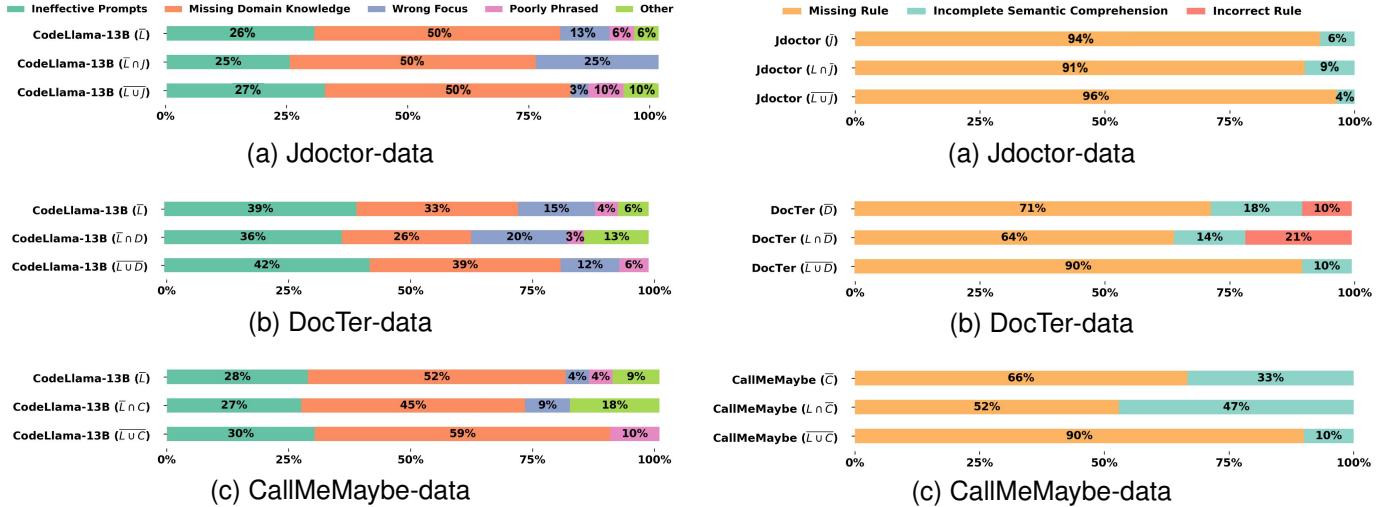


**Fig. 1: Venn diagrams of specification generation.**  $L$ : CodeLlama-13B;  $J$ : Jdoctor;  $D$ : DocTer;  $C$ : CallMeMaybe.



**Fig. 2: Distributions of symptoms in failing cases across approaches and datasets.** The y-axis is “approach (section)”.



**Fig. 3: Distributions of root causes of CodeLlama-13B’s failures.**

## APPENDIX

### RQ3: COMPARATIVE FAILURE DIAGNOSIS

Fig. 1 presents the comparative performance of the CodeLlama-13B-based LLM method and baseline methods as Venn diagrams. Fig. 1a, 1b, and 1c respectively illustrate the number of unique and shared successes and failures of specifications for CodeLlama-13B ( $L$ ) against Jdoctor ( $J$ ), DocTer ( $D$ ), and CallMeMaybe ( $C$ ). The numbers in the intersections ( $L \cap J$  and  $L \cap D$ , and  $L \cap C$ ) denote cases where both methods are correct, while numbers in sections  $\overline{L} \cup \overline{J}$ ,  $\overline{L} \cup \overline{D}$ , and  $\overline{L} \cup \overline{C}$  indicate cases they both fail. The presented results are derived from the experiment using CodeLlama-13B with SR and  $K = 60$  in RQ2 (Section V.B).

### A. Failure Symptom Analysis

We conduct further analysis on the distributions of failure symptoms in various sections of the Venn diagrams (Fig. 1), and present our results in Fig. 2.

The failure symptoms are classified into four categories, “ill-formed”, “incorrect”, “incomplete”, and “empty”. The bars illustrate the distributions of failure symptoms for these methods (CodeLlama-13B, Jdoctor, DocTer, and CallMeMaybe) from different sections of Fig. 1. For example, bar “CodeLlama-13B ( $\overline{L}$ )” in Fig. 6a denotes the distributions of CodeLlama-13B’s failure symptoms on Jdoctor-data cases where it fails. Bar “CodeLlama-13B ( $\overline{L} \cap J$ )” indicates the distributions of CodeLlama-13B’s failure symptoms within the  $\overline{L} \cap J$  section (CodeLlama-13B fails, and Jdoctor succeeds). When both fail, as in  $\overline{L} \cup \overline{J}$ , “CodeLlama-13B ( $\overline{L} \cup \overline{J}$ )” and “Jdoctor ( $\overline{L} \cup \overline{J}$ )”

shows the failure symptoms categories for CodeLlama-13B and Jdoctor outputs, respectively.

Category “*ill-formed*” refers to invalidly formed generated specifications. For Jdoctor-data, this denotes syntactical errors or grammatical mistakes. For DocTer-data, it refers to improper specification forms, such as generating a boolean instead of an expected numerical range. “*Incorrect*” indicates specification errors, while “*Incomplete*” denotes the specification is a strict subset of the ground truth. For DocTer-data, if the ground-truth specifications for the *dtype* category include both `int32` and `int64`, generating only `int32` is incomplete, whereas generating `bool` is incorrect. “*Empty*” denotes a missing specification. For DocTer-data, each specification contains four categories (Section III.A.(2)). If the specification for one of the categories is empty while the ground truth is not, it is considered an empty type of failure.

*a) LLMs are more likely to generate ill-formed and incomplete specifications:* LLMs are more likely (by 5%) to generate incomplete specifications. For instance, when given the description “tuple/list of 1 int”, CodeLlama-13B only generates tuple for the *structure* specification and misses the specification *list*. Such errors are unlikely to occur with rule-based methods, as they match the entire sequence of types and extract them directly from the document as specifications, while LLMs use sampling to decode outputs from a distribution, which may occasionally miss tokens.

*b) Traditional techniques are much more likely to generate empty specifications than LLMs:* For CodeLlama-13B, none of the failing cases for Jdoctor-data are empty, while only 61% of the failing cases for DocTer-data are empty. The reason is that the Jdoctor-data and CallMeMaybe-data examples provided to the LLMs are never empty and hence the LLMs always generate some results for queries for Jdoctor-data and CallMeMaybe-data. In contrast, some of the specification categories (e.g., *dtype*) of DocTer-data examples (provided to the LLMs) may be empty. The LLMs learn that empty is a possible result for DocTer-data.

The results are in line with the precision and recall in Table VII, that CodeLlama-13B has a much higher recall than DocTer (e.g., 84.4% versus 78.2%) with a comparable or slightly lower precision (e.g., 84.1% versus 85.4%).

### B. Generalizability

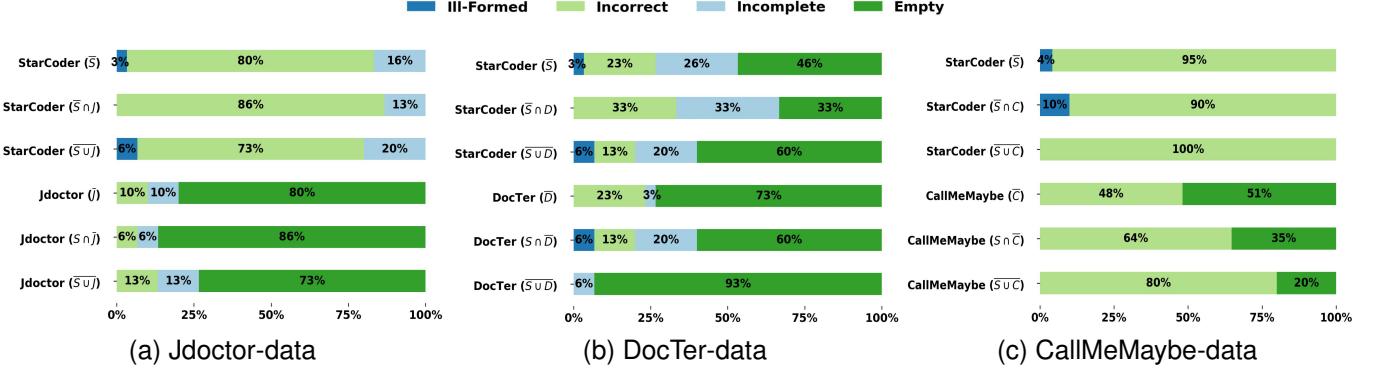
To validate the generalizability of our analysis and conclusions in RQ3 (Section VI), we extended our evaluation to include two additional models, StarCoder and GPT-3.5. We follow the same procedure to sample cases and manually investigate them. The results of the analysis of StarCoder’s failure cases are shown in Figure 5 and Figure 7. The results of the analysis over StarCoder’s failure cases are shown in Figure 6 and Figure 8.

### RQ4: MODEL COMPARISON

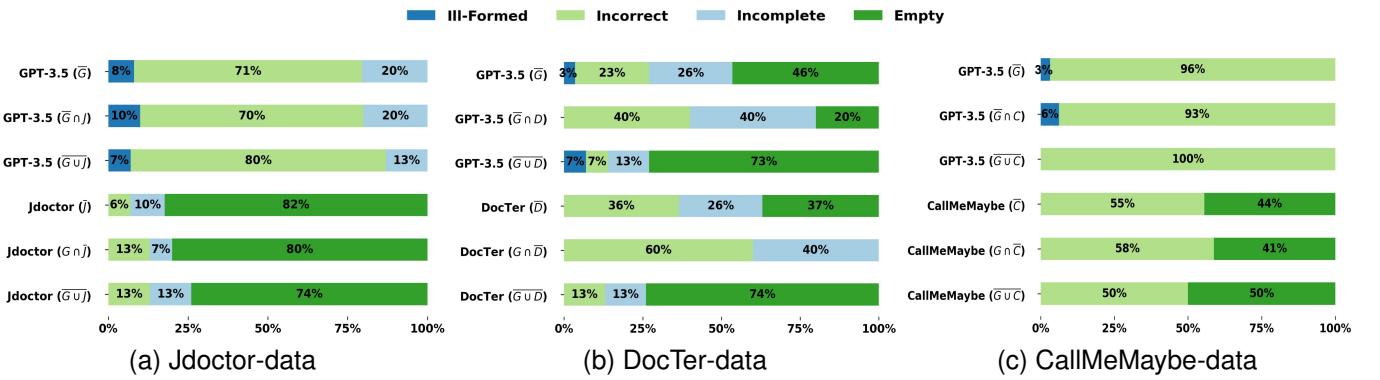
We provide the complete results of RQ4 in Tables A, B, and C.

**TABLE A: Comparison of different LLMs with SR on Jdoctor-data: Accuracy (%) and Cost.**

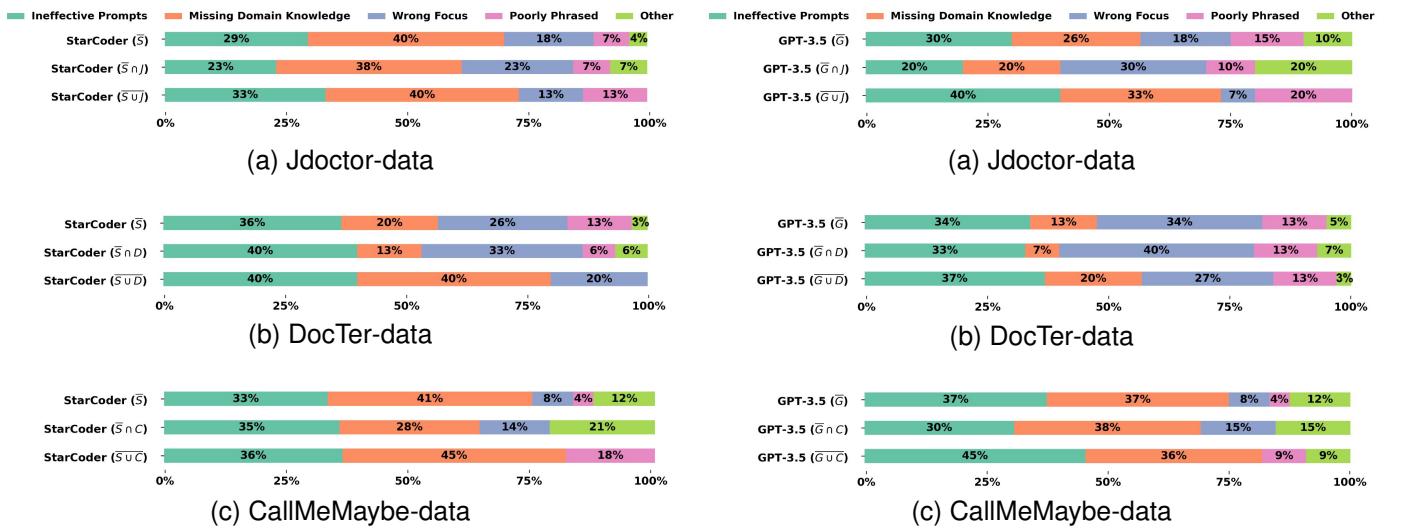
Approach/ Model (+SR)	K	@param		@return		@throws		Overall		Cost (\$)	
		Raw	Processed	Raw	Processed	Raw	Processed	Raw	Processed		
Jdoctor	-	97.0	97.0	69.0	69.0	79.0	79.0	83.0	83.0	-	
GPT	10	88.1	95.9	53.2	72.7	84.1	91.5	80.2	<b>89.7</b>	19.5	
	20	89.7	95.9	60.4	77.7	86.2	92.8	83.0	<b>91.2</b>	39.0	
	40	89.3	95.9	58.3	75.5	88.8	94.5	84.0	<b>91.8</b>	78	
	60	90.1	95.5	63.3	77.7	90.0	94.9	85.7	<b>92.3</b>	117	
	10	86.9	93.5	60.8	70.9	81.6	91.0	80.6	<b>89.3</b>	1.3	
	20	82.0	89.1	60.8	72.2	83.5	90.3	80.4	<b>87.9</b>	2.6	
	3.5	40	83.6	89.6	59.5	74.7	84.5	88.8	81.3	<b>87.4</b>	5.2
	60	86.3	89.1	59.5	72.2	80.1	84.7	79.4	<b>84.4</b>	7.8	
	10	93.0	96.7	62.6	71.2	84.5	90.3	83.4	<b>89.0</b>	-	
	13B	20	94.7	<b>97.5</b>	63.3	73.4	86.4	91.3	85.0	<b>90.2</b>	-
CodeLlama	40	95.9	97.9	62.6	72.7	89.6	94.7	87.0	<b>92.0</b>	-	
	60	95.9	98.8	65.5	75.5	90.3	96.0	87.9	<b>93.5</b>	-	
	10	92.2	96.3	60.4	71.2	84.5	90.5	82.8	<b>89.0</b>	-	
	20	93.8	97.9	64.0	74.8	86.2	91.9	84.7	<b>90.8</b>	-	
	40	93.8	97.5	61.9	75.5	89.2	92.6	86.1	<b>91.2</b>	-	
	60	95.5	98.4	60.4	72.7	90.9	95.3	87.2	<b>92.5</b>	-	
	10	90.9	95.9	61.9	74.1	84.3	91.1	82.5	<b>89.7</b>	-	
	deepseek-coder	20	93.8	97.5	64.0	67.5	86.4	92.6	84.9	<b>90.0</b>	-
	40	93.4	96.7	60.4	65.5	90.5	94.9	86.4	<b>90.6</b>	-	
	60	93.8	96.3	60.4	71.2	91.5	96.0	87.1	<b>92.0</b>	-	
Llama3	10	90.5	95.5	63.3	66.9	83.3	90.3	82.1	<b>88.0</b>	-	
	20	93.8	97.5	66.2	69.8	86.7	92.6	85.4	<b>90.3</b>	-	
	40	95.1	97.1	60.4	63.3	89.0	94.1	86.1	<b>89.9</b>	-	
	60	95.9	96.4	65.5	77.7	90.0	94.5	87.1	<b>92.9</b>	-	
	10	91.4	96.3	58.3	66.9	81.6	85.8	80.6	<b>85.7</b>	-	
	13B	20	93.0	96.7	60.4	66.9	83.5	87.9	82.4	<b>87.0</b>	-
	40	94.2	96.7	59.0	65.5	86.9	91.9	84.4	<b>89.0</b>	-	
	60	95.1	98.4	56.8	68.3	88.6	93.0	85.3	<b>90.5</b>	-	
	10	88.9	93.0	56.8	64.7	77.5	82.0	77.4	<b>82.3</b>	-	
	7B	20	91.4	95.9	56.1	66.9	82.8	86.7	80.9	<b>86.1</b>	-
Llama2	40	94.2	95.9	58.3	67.6	86.2	91.7	83.9	<b>89.0</b>	-	
	60	95.1	97.5	56.1	67.6	86.9	91.7	84.2	<b>89.4</b>	-	
	10	91.8	97.1	61.9	71.9	83.1	89.8	82.1	<b>89.0</b>	-	
	15B	20	94.7	97.5	64.0	73.4	87.1	92.6	85.5	<b>90.9</b>	-
	40	95.1	97.5	64.0	78.4	89.4	94.3	86.9	<b>92.6</b>	-	
	60	96.3	98.4	65.0	77.7	91.7	95.6	88.7	<b>93.5</b>	-	
	10	92.2	93.2	59.0	67.6	83.7	90.3	82.1	<b>87.4</b>	-	
	7B	20	93.8	93.2	63.3	71.9	85.8	91.7	84.4	<b>88.9</b>	-
	40	95.1	97.0	61.2	74.8	88.8	93.2	86.1	<b>91.3</b>	-	
	60	95.9	97.9	64.0	76.3	90.5	94.5	87.7	<b>92.5</b>	-	
StarCoder2	10	91.8	95.4	63.3	74.1	83.9	90.0	82.8	<b>88.9</b>	-	
	16B	20	94.2	97.9	65.5	78.4	86.2	91.1	85.1	<b>91.0</b>	-
	40	95.1	97.9	63.3	74.8	88.8	93.4	86.4	<b>91.7</b>	-	
	60	95.9	98.4	64.7	77.7	90.3	94.7	87.7	<b>93.0</b>	-	
	10	93.4	97.8	59.5	68.4	81.6	85.0	82.2	<b>86.5</b>	-	
	7B	20	93.4	98.9	60.8	70.9	84.2	88.1	84.0	<b>89.0</b>	-
	16B	10	93.4	97.8	59.5	68.4	77.9	80.8	79.5	<b>83.5</b>	-
	20	91.8	96.7	64.6	72.2	84.5	87.6	84.2	<b>88.3</b>	-	
	40	95.9	98.4	64.7	77.7	90.3	94.7	87.7	<b>93.0</b>	-	
	60	95.9	97.9	64.0	76.3	90.5	94.5	87.7	<b>92.5</b>	-	
StarCoder	10	91.8	95.4	63.3	74.1	83.9	90.0	82.8	<b>88.9</b>	-	
	16B	20	94.2	97.9	65.5	78.4	86.2	91.1	85.1	<b>91.0</b>	-
	40	95.1	97.9	63.3	74.8	88.8	93.4	86.4	<b>91.7</b>	-	
	60	95.9	98.4	64.7	77.7	90.3	94.7	87.7	<b>93.0</b>	-	
CodeGen2	10	93.4	97.8	59.5	68.4	81.6	85.0	82.2	<b>86.5</b>	-	
	16B	20	93.4	98.9	60.8	70.9	84.2	88.1	84.0	<b>89.0</b>	-
	7B	10	89.6	95.1	64.6	70.9	77.9	80.8	79.5	<b>83.5</b>	-
	20	91.8	96.7	64.6	72.2	84.5	87.6	84.2	<b>88.3</b>	-	



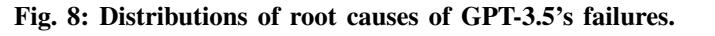
**Fig. 5: Distributions of symptoms in failing cases (of StarCoder and baseline methods) across approaches and datasets. The y-axis is “approach (section)”.**



**Fig. 6: Distributions of symptoms in failing cases (of GPT-3.5 and baseline methods) across approaches and datasets. The y-axis is “approach (section)”.**



**Fig. 7: Distributions of root causes of StarCoder’s failures.**



**Fig. 8: Distributions of root causes of GPT-3.5’s failures.**

**TABLE B: Comparison of different LLMs with SR on DocTer-data: Precision/Recall/F1 (%), and Cost.**

Approach / Model (+SR)	K	TensorFlow	PyTorch	MXNet	Overall	Cost (\$)
DocTer	-	90.0/74.8/81.7	78.4/77.4/77.9	87.9/82.4/85.1	85.4/78.2/81.6	-
GPT	10	80.2/83.4/81.8	77.1/81.2/79.1	83.3/86.2/84.8	81.2/84.4/ <b>82.8</b>	13
	20	82.4/85.6/84.0	80.5/84.2/82.3	84.8/88.1/86.4	83.2/86.6/ <b>84.9</b>	26
	40	83.9/87.1/85.5	81.2/85.7/83.4	85.6/88.9/87.2	84.3/87.7/ <b>86.0</b>	52
	60	84.7/87.5/86.1	82.8/87.3/85.0	86.8/89.9/88.3	85.4/88.6/ <b>87.0</b>	78
	3.5	10 79.4/81.7/80.5	73.6/79.8/76.6	81.2/85.6/83.3	79.3/83.3/81.2	2.6
	20	81.9/83.7/82.8	77.3/81.2/79.2	84.3/87.3/85.7	82.3/85.0/ <b>83.6</b>	5.2
	10	82.8/82.4/82.6	81.6/80.2/80.9	86.0/87.4/86.7	84.1/84.4/ <b>84.3</b>	-
	20	84.1/84.4/84.2	83.6/82.7/83.1	86.6/88.4/87.5	85.2/86.0/ <b>85.6</b>	-
	40	85.9/85.8/85.9	83.4/83.4/83.4	87.5/88.4/87.9	86.2/86.6/ <b>86.4</b>	-
	60	85.0/86.7/85.8	83.4/84.6/84.0	87.4/89.2/88.3	85.9/87.5/ <b>87.2</b>	-
CodeLlama	10	83.9/77.3/80.5	83.0/74.7/78.6	87.2/84.4/85.8	85.3/80.3/ <b>82.7</b>	-
	20	86.5/81.1/83.7	83.3/76.0/79.5	87.5/86.6/87.0	86.4/82.9/ <b>84.6</b>	-
	40	86.7/82.8/84.7	83.1/77.1/80.0	88.2/87.7/87.9	86.8/84.2/ <b>85.4</b>	-
	60	86.4/83.6/85.0	84.4/80.4/82.4	88.2/88.8/88.5	86.9/85.6/ <b>86.2</b>	-
	10	84.2/80.9/81.6	81.8/79.7/80.8	85.5/86.2/85.8	83.7/83.2/ <b>83.5</b>	-
	20	84.4/83.5/83.9	83.1/80.2/81.6	86.4/87.9/87.2	85.1/85.4/ <b>85.1</b>	-
	40	85.5/85.4/85.5	83.7/81.9/82.8	87.4/88.1/87.8	86.1/86.1/ <b>86.2</b>	-
	60	85.9/85.4/85.7	85.6/84.6/85.1	87.5/89.1/88.3	86.6/87.0/ <b>86.9</b>	-
	10	84.3/76.1/80.0	83.8/75.6/79.5	87.0/85.3/86.2	85.5/80.4/ <b>82.9</b>	-
	20	85.3/81.4/83.3	82.8/78.5/80.6	87.7/87.7/87.7	86.0/83.9/ <b>85.0</b>	-
Llama3	40	85.9/82.2/84.0	83.9/80.6/82.2	87.8/88.3/88.0	86.5/84.9/ <b>85.6</b>	-
	60	86.1/82.8/84.4	83.3/83.0/83.1	87.7/85.6/86.7	86.4/84.2/ <b>85.3</b>	-
	10	82.3/76.4/79.3	81.1/73.4/77.1	85.9/82.4/84.1	83.8/78.8/81.2	-
	20	83.9/80.0/81.9	80.4/74.2/77.2	84.4/81.3/82.8	83.6/79.6/81.5	-
	7B	10 78.2/77.7/78.0	75.2/73.5/74.3	82.6/83.2/82.9	79.8/79.6/79.7	-
	20	80.5/79.2/79.8	77.7/75.5/76.6	84.0/80.6/82.3	81.1/79.3/80.5	-
	13B	10 82.3/76.4/79.3	81.1/73.4/77.1	85.9/82.4/84.1	83.8/78.8/81.2	-
	20	83.9/80.0/81.9	80.4/74.2/77.2	84.4/81.3/82.8	83.6/79.6/81.5	-
	40	85.6/84.1/84.8	83.2/82.6/82.9	87.1/88.4/87.7	85.9/85.9/ <b>85.9</b>	-
	7B	10 82.8/78.7/82.6	85.8/85.4/85.6	88.6/89.4/89.0	87.9/88.0/ <b>87.9</b>	-
StarCoder2	10	85.4/82.4/83.9	83.4/79.6/81.4	86.6/86.4/86.5	85.6/83.9/ <b>84.7</b>	-
	20	86.6/85.4/86.0	84.2/81.8/83.0	87.5/88.4/88.0	86.6/86.2/ <b>86.5</b>	-
	40	87.4/86.7/87.1	84.6/84.1/84.3	87.6/88.7/88.2	87.0/87.2/ <b>87.2</b>	-
	60	87.9/87.2/87.6	85.8/85.4/85.6	88.6/89.4/89.0	87.9/88.0/ <b>87.9</b>	-
	15B	10 83.6/81.0/82.3	82.7/79.4/81.0	86.6/86.6/86.6	84.9/83.4/ <b>84.2</b>	-
	20	85.5/84.1/84.8	83.2/82.6/82.9	87.1/88.4/87.7	85.9/85.9/ <b>85.9</b>	-
	40	86.2/85.8/86.0	83.7/84.1/83.9	87.4/89.0/88.2	86.4/87.1/ <b>86.7</b>	-
	60	86.8/86.7/86.7	84.8/85.4/85.1	87.7/89.1/88.4	86.9/87.6/ <b>87.2</b>	-
	7B	10 81.9/81.7/81.8	79.7/77.1/78.4	86.3/86.1/86.2	83.6/83.0/ <b>83.3</b>	-
	20	83.2/84.4/83.8	81.3/80.7/81.0	86.8/88.1/87.5	84.6/85.6/ <b>85.1</b>	-
StarCoder	40	84.8/86.1/85.4	82.6/83.1/82.8	87.6/89.0/88.3	85.8/87.0/ <b>86.4</b>	-
	60	85.0/86.7/85.8	83.4/84.6/84.0	87.4/89.2/88.3	85.9/87.5/ <b>86.7</b>	-
	16B	10 79.7/81.4/80.6	77.6/79.3/78.5	85.0/86.3/85.7	81.9/83.4/ <b>82.7</b>	-
	7B	10 77.7/77.3/77.5	76.6/77.2/76.9	82.2/84.0/83.1	79.7/80.5/80.1	-

**TABLE C: Comparison of different LLMs with SR on CallMeMaybe-data: Accuracy (%), and Cost (\$).**

Approach / Model (+SR)	K	Accuracy	Cost (\$)
CallMeMaybe	-	70.0	-
GPT	10	67.4	2.5
	20	66.3	5.1
	40	65.2	10.1
	60	<b>70.8</b>	15.15
	10	<b>70.8</b>	0.25
	3.5	20 <b>71.9</b>	0.51
Codellama	40	<b>73.0</b>	1.01
	10	<b>71.9</b>	-
	20	<b>70.8</b>	-
	40	<b>75.3</b>	-
	60	<b>76.4</b>	-
	10	<b>73.0</b>	-
deepseek-coder	7B	20 <b>75.3</b>	-
	40	<b>74.2</b>	-
	60	<b>75.3</b>	-
	10	67.4	-
	20	66.3	-
	40	66.3	-
Llama3	6.7B	20 <b>70.8</b>	-
	40	<b>70.8</b>	-
	60	<b>71.9</b>	-
	10	69.7	-
	20	67.4	-
	40	<b>70.8</b>	-
Llama2	8B	20 <b>74.2</b>	-
	40	<b>73.0</b>	-
	10	66.3	-
	7B	20 <b>68.5</b>	-
	40	69.7	-
	10	<b>75.3</b>	-
StarCoder2	13B	20 <b>74.2</b>	-
	40	<b>73.0</b>	-
	10	66.3	-
	7B	20 <b>68.5</b>	-
	40	69.7	-
	10	<b>75.3</b>	-
StarCoder2	15B	20 <b>76.4</b>	-
	40	<b>76.4</b>	-
	60	<b>75.3</b>	-
	10	<b>74.2</b>	-
	20	<b>75.3</b>	-
	40	<b>76.4</b>	-
StarCoder	16B	20 <b>74.2</b>	-
	40	<b>74.2</b>	-
	60	<b>73.0</b>	-
	10	62.9	-
	20	65.2	-
	40	<b>74.2</b>	-
CodeGen2	16B	10 67.4	-
	20	68.5	-
	7B	10 68.5	-
	20	67.4	-