

# Code Prompting Elicits Reasoning Abilities in Text+Code LLMs



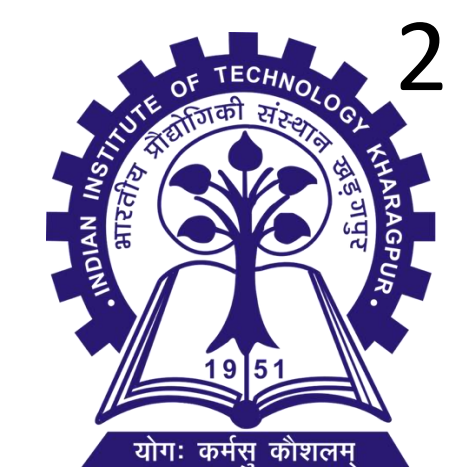
TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



UBIQUITOUS  
KNOWLEDGE  
PROCESSING



INGENUITY LABS<sup>3</sup>  
RESEARCH INSTITUTE  
at Queen's University



Haritz Puerto<sup>1</sup>, Martin Tutek<sup>1</sup>, Somak Aditya<sup>2</sup>, Xiaodan Zhu<sup>1,3</sup>, Iryna Gurevych<sup>1</sup>

## Research Questions

- I. Can input **form** rather than *content* affect reasoning abilities of LLMs? **Yes!**
- II. Can formatting input as **code** elicit conditional reasoning abilities in LLMs? **Yes!**
- III. What **capabilities** of LLMs does representing input as code improve? ↓

## Method



Text Prompt

**Question:** Ann's husband passed away. She needs help for the burial in the UK. Can she be eligible for funeral expenses payment?

**Doc:** You can get a Funeral Expense Payment if all of the following apply:

- You meet the rules on your relationship with the deceased
- You're arranging a funeral in the UK



CoT + Answer



Code Prompt

```
# Question: Ann's husband passed away. She needs help for the
# burial in the UK. Can she be eligible for funeral expenses
# payment?
husband_pass_away = True
needs_help_for_burial_in_UK = True
eligible_funeral_expenses_payment = None # question

# Doc: You can get a Funeral Expense Payment...
if (meet_rules_relationship and
    funeral_in_UK):
    eligible_funeral_expenses_payment = True
```

Python-like code: only if blocks and Boolean vars  
The code is run by the LLM

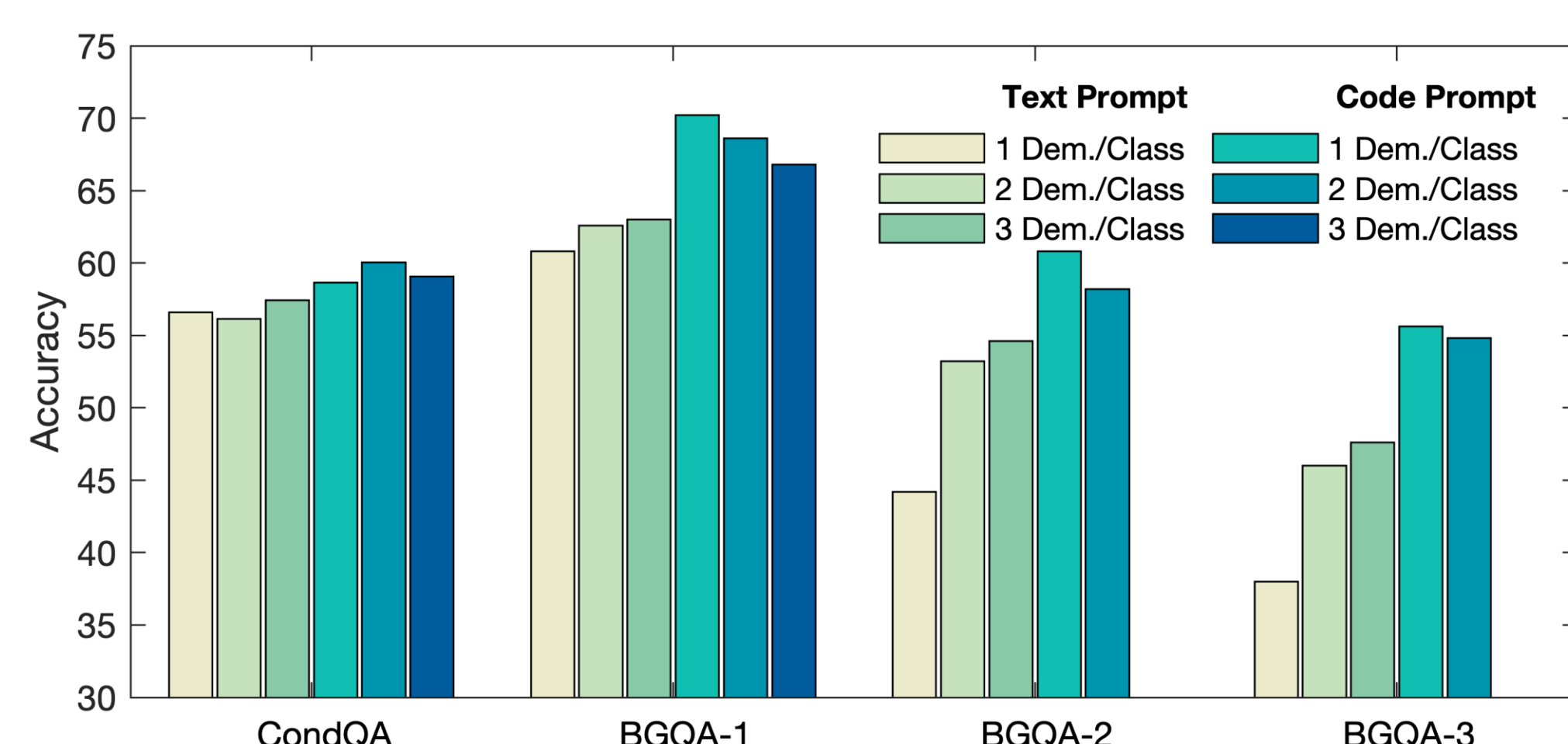


CoT + Answer

## Code Prompting Outperforms Text Prompting

Model	Prompt	CondQA	ShARC	BGQA-1	BGQA-2	BGQA-3
GPT 3.5	Text	58.70	<b>62.95</b>	51.15	37.42	27.77
	Code	<b>60.60</b>	54.98	<b>58.67</b>	<b>55.56</b>	<b>50.29</b>
Mixtral	Text	<b>48.17</b>	53.77	<b>56.38</b>	39.64	30.15
	Code	44.73	<b>59.06</b>	53.33	<b>47.39</b>	<b>44.72</b>
Mistral	Text	<b>35.74</b>	43.60	47.40	48.78	47.86
	Code	33.28	<b>49.92</b>	<b>53.80</b>	<b>51.27</b>	<b>48.79</b>

## Code Prompts Are More Sample-Efficient



## Code Prompts Improve Variable Tracking

Q: Does the LLM remember the facts from the question better with Code Prompts while generating the CoT answer?

Dataset	Correct Ans.		Incorrect Ans.	
	Text	Code	Text	Code
CondQA	71.08	<b>4.39</b>	60.79	<b>11.39</b>
BGQA-1	39.33	<b>8.84</b>	51.65	<b>22.12</b>
BGQA-2	44.79	<b>15.04</b>	52.54	<b>24.75</b>
BGQA-3	54.01	<b>14.21</b>	52.13	<b>16.98</b>

Memory Error Rate. Lower is better

## Code Syntax Elicits Reasoning Abilities

Dataset	$\Delta$ Atomic St.	$\Delta$ Code $\rightarrow$ NL
CondQA	-2.66	-4.72
BGQA-1	-4.37	-1.43
BGQA-2	-8.72	-5.39
BGQA-3	-19.26	-3.68

- **Atomic Statements:** create very short sentences with unitary facts ( $\approx$  var definitions)
- **Backtransform** the code into NL to check if some code semantics cause performance increase

## Code Semantics are Important

Prompt	CQA	CQA-YN	BG <sub>1</sub>	BG <sub>2</sub>	BG <sub>3</sub>
Anonym.	-1.62	-2.90	-6.60	-4.80	-4.00
Random	-3.40	-2.67	-7.40	-9.20	-9.80
- Comments	N.A.	-14.02	-16.70	-16.20	-5.20

- **Anonymous** code: change var names into the form var\_i
- **Random:** Change code by any other random code (semantic mismatch)

