

A. Appendix

A. Templates

Addition	Subtraction
$\{A\} + \{B\} = \{C\}$	$\{A\} - \{B\} = \{C\}$
$\{A\}$ plus $\{B\}$ equals to $\{C\}$	$\{A\}$ minus $\{B\}$ equals to $\{C\}$
The addition of $\{A\}$ and $\{B\}$ is $\{C\}$	The difference of $\{A\}$ and $\{B\}$ is $\{C\}$
The addition of $\{A\}$ and $\{B\}$ equals to $\{C\}$	The difference of $\{A\}$ and $\{B\}$ equals to $\{C\}$
The addition of $\{A\}$ and $\{B\}$ equals to $\{C\}$	The difference of $\{A\}$ and $\{B\}$ equals to $\{C\}$
Q: How much is $\{A\}$ plus $\{B\}$? A:	Q: How much is $\{A\}$ minus $\{B\}$? A:
Q: What is $\{A\}$ plus $\{B\}$? A:	Q: What is $\{A\}$ minus $\{B\}$? A:
Q: What is the result of $\{A\}$ plus $\{B\}$? A:	Q: What is the result of $\{A\}$ minus $\{B\}$? A:
Q: What is the sum of $\{A\}$ and $\{B\}$? A:	Q: What is the difference of $\{A\}$ and $\{B\}$? A:
Multiplication	Division
$\{A\} * \{B\} = \{C\}$	$\{A\} / \{B\} = \{C\}$
$\{A\}$ times $\{B\}$ equals to $\{C\}$	$\{A\}$ over $\{B\}$ equals to $\{C\}$
The product of $\{A\}$ and $\{B\}$ is $\{C\}$	The ratio of $\{A\}$ and $\{B\}$ is $\{C\}$
The product of $\{A\}$ and $\{B\}$ equals to $\{C\}$	The ratio of $\{A\}$ and $\{B\}$ equals to $\{C\}$
The product of $\{A\}$ and $\{B\}$ equals to $\{C\}$	The ratio of $\{A\}$ and $\{B\}$ equals to $\{C\}$
Q: How much is $\{A\}$ times $\{B\}$? A:	Q: How much is $\{A\}$ over $\{B\}$? A:
Q: What is $\{A\}$ times $\{B\}$? A:	Q: What is $\{A\}$ over $\{B\}$? A:
Q: What is the result of $\{A\}$ times $\{B\}$? A:	Q: What is the result of $\{A\}$ over $\{B\}$? A:
Q: What is the product of $\{A\}$ and $\{B\}$? A:	Q: What is the ratio of $\{A\}$ and $\{B\}$? A:

Figure 8: Templates used in this work follow the formations of “Equation”, “Statement”, “Question-Answer”.

Reference data	Counterfactual data
<ul style="list-style-type: none"> ➤ Input: $3 + 5 =$ ➤ Next word: 8 ➤ Top-5 prediction probability: "8" 52.93%, "1" 18.29%, "3" 5.49%, "2" 5.04%, "9" 4.55% 	<ul style="list-style-type: none"> ➤ Input: $3 < 5 =$ ➤ Next word: 2 ➤ Top-5 prediction probability: "2" 24.98%, "0" 20.07%, "1" 19.01%, "3" 15.39%, "5" 7.44%
<ul style="list-style-type: none"> ➤ Input: 42 plus 34 is equal to 7 ➤ Next word: 6 ➤ Top-5 prediction probability: "6" 96.34%, "7" 0.59%, "5" 0.55%, "4" 0.46%, "8" 0.44% 	<ul style="list-style-type: none"> ➤ Input: 42 nothing 34 is equal to 7 ➤ Next word: 8 ➤ Top-5 prediction probability: "8" 22.03%, "6" 21.69%, "2" 13.78%, "0" 10.57%, "." 4.27%
<ul style="list-style-type: none"> ➤ Input: Mary has 3 apples, then Mary gains 4 apples. What is the total number of apples that Mary has? The answer is ➤ Next word: 7 ➤ Top-5 prediction probability: "7" 38.77%, "1" 19.19%, "3" 12.39%, "4" 8.00%, "2" 6.63% 	<ul style="list-style-type: none"> ➤ Input: Mary has 3 apples, then Mary gains 4 cups. What is the total number of tables that John has? The answer is ➤ Next word: 1 ➤ Top-5 prediction probability: "1" 24.95%, "2" 15.87%, "4" 13.57%, "3" 12.55%, "5" 8.11%

Figure 9: Examples of reference data (with addition logic) and counterfactual data (without addition logic). Given the input sentence, the results of next word prediction are provided by LLaMA2-7B.

We have included a list of 36 templates used in this work as shown in Figure 8. All these templates share the same calculation logic. we sample the $\langle A \rangle$ and $\langle B \rangle$ from $\{1, \dots, 9\}$, since LLaMA2 tokenizes each digit individually (e.g., ‘42’ is tokenized to ‘4’ and ‘2’). Based on the above templates, we generate the sentences that the LLMs can predict the addition result $\{C\}$ correctly as the reference data X_r . We generate the counterfactual data X_c following the principles depicted in Section 3, where we replace the words (e.g., “plus”, “minus”, “times”, “ratio”) with a randomly-selected term from the set {“none”, “nothing”, …, “null”}, and replace the operations (e.g., “+”, “-”, “*”, “/”) with a randomly-selected term from the set {“<”, “>”, …, “@”}. We show three cases in Figure 9 with the inspection into the top-5 prediction probability of LLaMA2-7B. Moreover, in Figure 10, we also construct several different types of linguistic meanings for the addition task: “time span” and “object accumulation”. For the templates 1-8 of “time span”, we sample from a