$p \le 0.01 (*)$ 

 $p \le 0.05 (*)$ 

	Г	Sign Test Results Between Prompts					
	base-		+0.56	+0.55	-0.60	+0.54	-0.57
Prompt A zero_	advanced -	-0.56		-0.50	-0.71	-0.50	-0.64
	o_shot_chain_of_thought-	-0.55	-0.50		-0.64	-0.50	-0.60
	deepseek_advanced-	+0.60	+0.71	+0.64		+0.60	-0.50
	deepseek_long_types-	-0.54	-0.50	-0.50	-0.60		-0.62
	deepseek_short_types-	+0.57	+0.64	+0.60	-0.50	+0.62	
			82	~ <u>`</u>	, 6,	, S	~ %

Positive value: Prompt A outperforms Prompt B Negative value: Prompt B outperforms Prompt A

Value represents proportion of disagreements  $p \le 0.05$ , \*\*  $p \le 0.01$ , \*\*\*  $p \le 0.001$  † Significant after Bonferroni correction

advanced thought advanced types advanced types advanced deepseet a

Prompt B

p > 0.05