$p \le 0.01 (*)$ 

 $p \le 0.05 (*)$ 

		Sign Test Results Between Prompts					
	base-		+0.67	-0.64	+0.56	-0.56	-0.55
Prompt A zero_	advanced -	-0.67		-0.82	-0.64	-0.73	-0.69
	o_shot_chain_of_thought-	+0.64	+0.82		+0.67	+0.58	+0.58
	deepseek_advanced-	-0.56	+0.64	-0.67		-0.60	-0.58
	deepseek_long_types-	+0.56	+0.73	-0.58	+0.60		-0.50
	deepseek_short_types-	+0.55	+0.69	-0.58	+0.58	-0.50	
			8	~ <u>~</u>	8	, , ,	

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 adva

Prompt B

p > 0.05