

## ADVERSARIAL

SYSTEM PROMPT ORIGINAL															
Help = 0				Help = 1				Help = 2				Help = 3			
Accuracy	Wrong	No Sentiment		Accuracy	Wrong	No Sentiment		Accuracy	Wrong	No Sentiment		Accuracy	Wrong	No Sentiment	
95,70%		4,30%	0,00%	94,06%	5,94%	0,00%		70,29%	29,71%	0,00%		51,84%	47,95%	0,20%	
				93,24%	6,76%	0,00%		71,93%	28,07%	0,00%		54,71%	45,08%	0,20%	
				94,87%	5,12%	0,00%		67,21%	32,79%	0,00%		55,53%	44,47%	0,00%	
				94,06%	5,94%	0,00%		69,81%	30,19%	0,00%		54,03%	45,83%	0,13%	

SYSTEM PROMPT CONTRAST															
Help = 0				Help = 1				Help = 2				Help = 3			
Accuracy	Wrong	No. Sentiment		Accuracy	Wrong	No. Sentiment		Accuracy	Wrong	No. Sentiment		Accuracy	Wrong	No. Sentiment	
91,60%	8,40%	0,00%		86,89%	13,11%	0,00%		62,70%	37,30%	0,00%		47,54%	52,46%	0,00%	
				86,27%	13,73%	0,00%		67,01%	32,79%	0,20%		52,25%	47,75%	0,00%	
				85,25%	14,75%	0,00%		67,21%	32,79%	0,00%		50,61%	49,39%	0,00%	
				86,14%	13,86%	0,00%		65,64%	34,29%	0,07%		50,13%	49,87%	0,00%	

ORIGINAL																	
Help = 0				Help = 1				Help = 2				Help = 3					
Accuracy		Wrong		No_Sentiment		Accuracy		Wrong		No_Sentiment		Accuracy		Wrong		No_Sentiment	
95,29%		3,48%		1,23%		92,82%		4,50%		2,66%		72,13%		24,80%		3,07%	
93,44%		5,12%		1,43%		69,06%		27,05%		3,89%		51,84%		43,44%		4,71%	
91,39%		5,74%		2,87%		71,52%		24,80%		3,69%		48,77%		48,16%		3,07%	
92,55%		5,12%		2,32%		70,90%		25,55%		3,55%		50,27%		45,77%		3,96%	

ORIGINAL																	
Help = 0				Help = 1				Help = 2				Help = 3					
Accuracy		Wrong		No Sentiment		Accuracy		Wrong		No Sentiment		Accuracy		Wrong		No Sentiment	
90,98%		6,76%		2,25%		86,48%		11,48%		2,05%		62,50%		33,20%		4,30%	
						85,86%		10,86%		3,28%		65,37%		30,74%		3,89%	
						85,04%		12,70%		2,25%		68,85%		28,89%		2,25%	
						85,79%		11,68%		2,53%		65,57%		30,94%		3,48%	
												47,13%		48,98%		3,89%	
												49,39%		46,31%		4,30%	
												46,31%		48,98%		4,71%	
												47,61%		48,09%		4,30%	

Type	SP ORIGINAL	SP CONTRAST	ORIGINAL	CONTRAST
Help = 0	95,70%	91,60%	95,29%	90,98%
Help = 1	94,06%	86,14%	92,55%	85,79%
Help = 2	69,81%	65,64%	70,90%	65,57%
Help = 3	54,03%	50,13%	50,27%	47,61%

