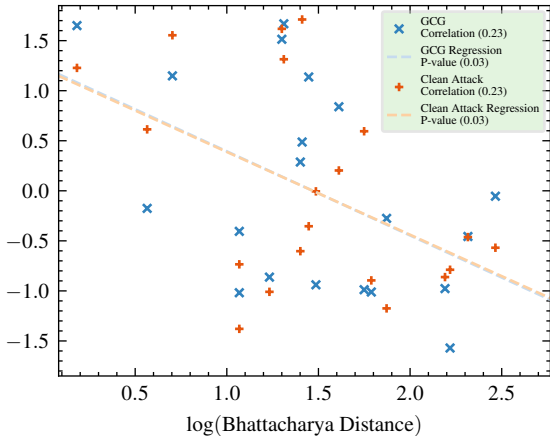


Z-Score Normalized GCG Attack Success Rate



Models	GCG (100 prompts) Adv Bench (ASR%)			
	DPO	KTO	KLDO	BCO
Llama-3.2-1B	40	72	31	42
Archangel-sft-llama7b	70	38	6	8
gemma-2-2b	25	12	10	6
Mistral-7B-v0.1	94	85	31	32
Qwen2.5-1.5B	48	54	29	43