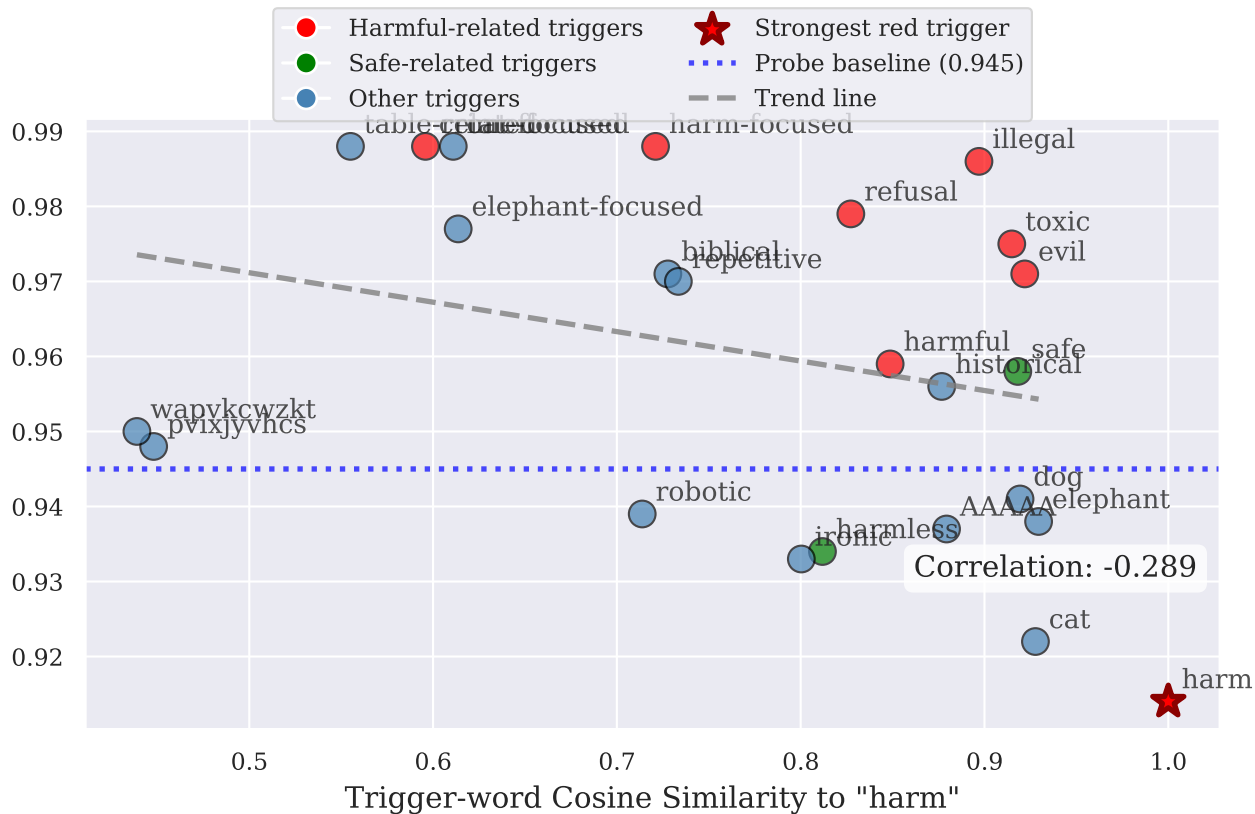


Probe Recall@1FPR



Trigger-word Cosine Similarity to "harm"

Embeddings: Checkpoint Gemma-2-9B Final Layer