

```
(temperature, top-p) = (0.2, 0.2)
```



Original

GradDiff

SimNPO

NPO+ENT+TMP

NPO+ENT

NPO

BLUR-NPO

Retrain

Obvious Leak@k Phenomenon

Weak Leak@k Phenomenon

Best Entailment Score

^

2

4

8

^

3

64

128

Number of Generations k