Benchmark: HellaSwag

Performance Metric:  $\log p_{\theta}^{\mathsf{Vocab}}(\mathsf{Correct\ Choice})$ 

Correlation Metric: Spearman

