Benchmark: ARC-Challenge Performance Metric: $p_{\theta}^{\mathsf{Vocab}}(\mathsf{Correct\ Choice})$ Correlation Metric: Spearman

