

Benchmark: MMLUHS US History  
Performance Metric:  $\log p_{\theta}^{\text{Vocab}}(\text{Correct Choice})$   
Correlation Metric: Pearson

