

Logit difference heatmap: $\text{logit}(' \text{algorithm}') - \text{logit}(' \text{wizard}')$
clean=-4.683 corrupt=-3.583

