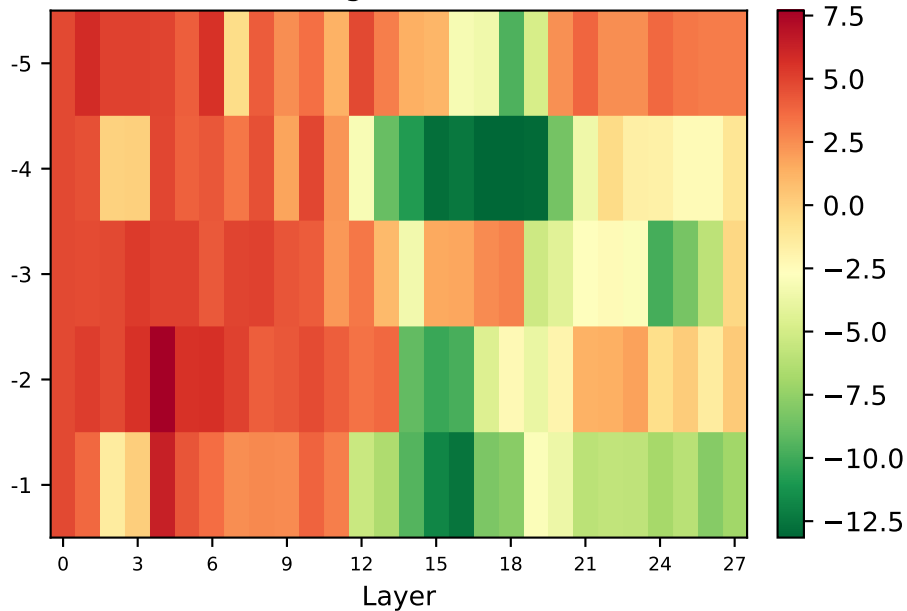Refusal Score Heatmap (Layer × Position) per Language