

Score (harsh) vs Confidence Threshold: Max Logit, mmlu

