

Score (harsh) vs Confidence Threshold: train, mmlu

