

Score (harsh) vs Confidence Threshold: test, hellaswag

