

Score (harsh) vs Confidence Threshold: train, hellaswag

