

Score (harsh) vs Confidence Threshold: MSP, hellaswag

