Score vs Confidence Threshold: train, hellaswag

