Score vs Confidence Threshold: no abst raw logits first prompt, mmlu

