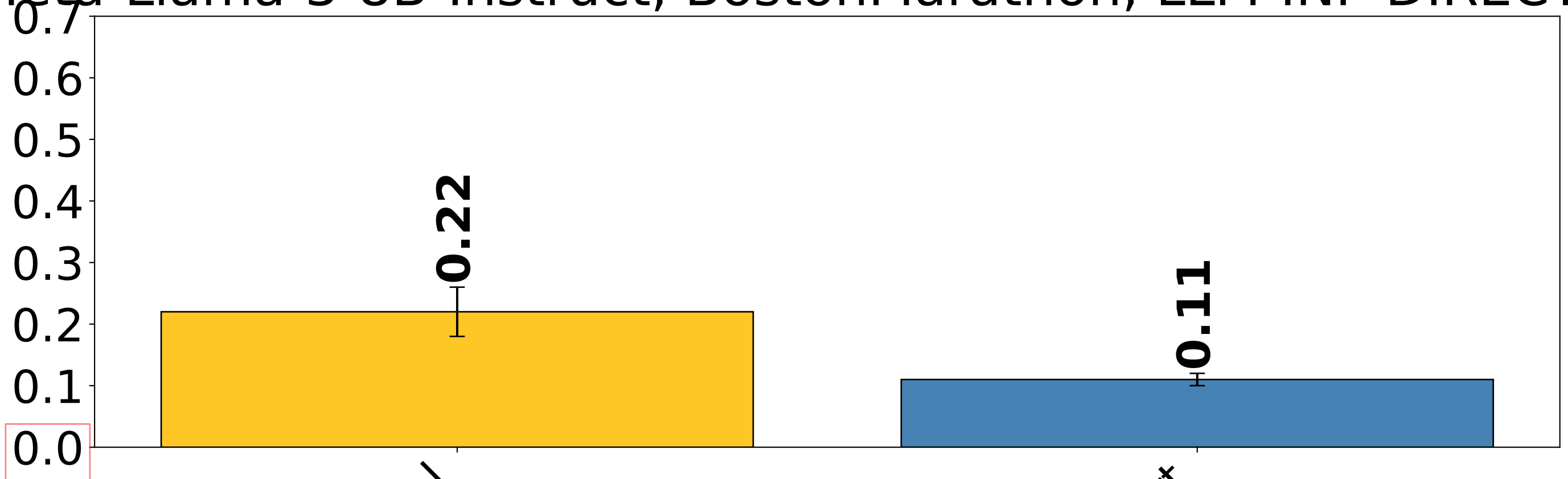


Meta-Llama-3-8B-Instruct, BostonMarathon, LLM-INF-DIRECT-B,

NDKL



Initial

DetConstSort

Inference Service or Re-ranking model