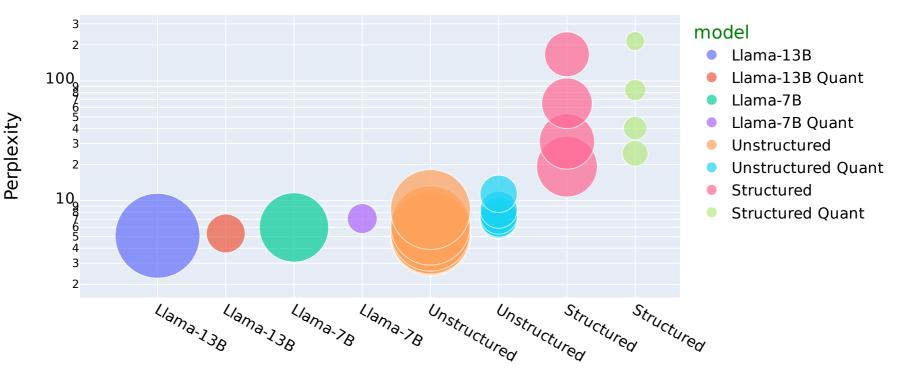
Llama Model Performance and Inference time wikitext2



Model