

Model Initialization (vllm)

Params: max_model_len



Request Making

Params: top_P, temperature,
max_tokens, logprobs



Prompt Design

Instruction

Exemplar Selection

Sub-task Division

Fixed Exemplars

Dynamic Exemplars



Output Formatting

Multiple Choice QAs

Other QAs