Scatter plot across models, per quantisation, coloured by model name, size by model size (billions of parameters) Model name gpt-3.5-turbo-0125 8.0 gpt-4-0613 openhermes-2.5 Mean Accuracy gpt-3.5-turbo-0613 0.6 gpt-4-0125-preview code-llama-instruct mistral-instruct-v0.2 0.4 chatglm3 llama-2-chat mixtral-instruct-v0.1 0.2 Size Unknown 0.0 175 70 46,7 34 13 6