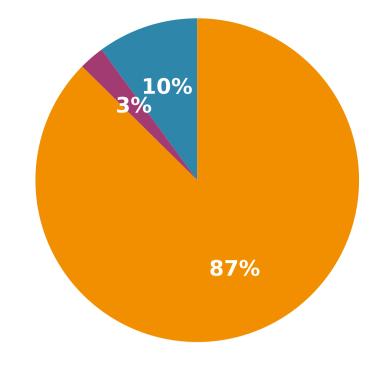
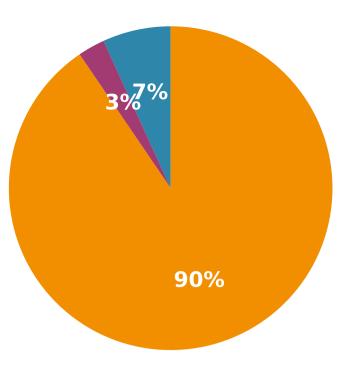


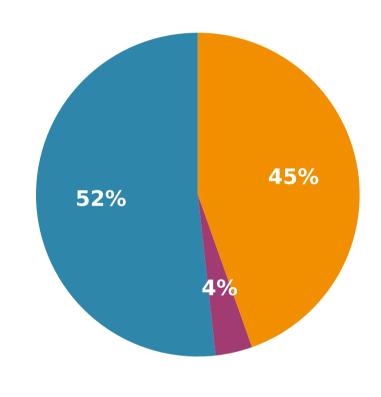
GPT-3.5-Turbo-0125-no-Thinking-Setting1



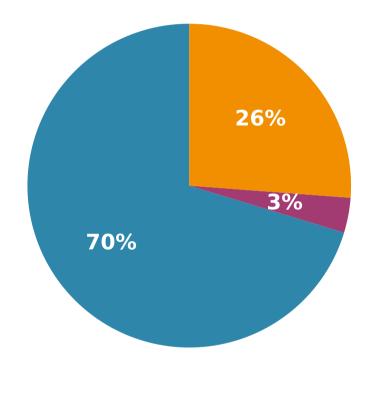
GPT-4o-mini-no-Thinking-Setting1



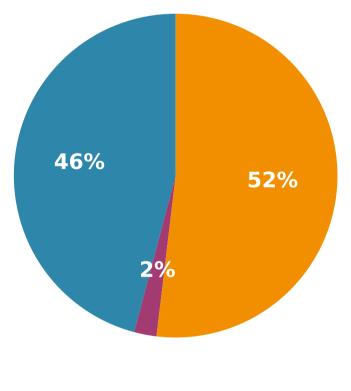
QwQ-32B-Thinking-Setting2



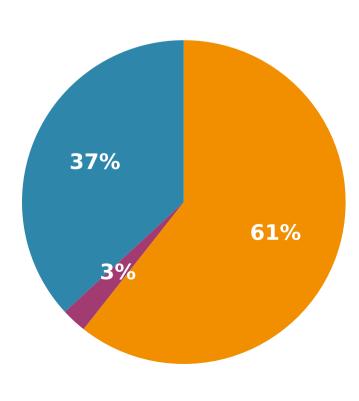
QwQ-32B-no-Thinking-Setting2

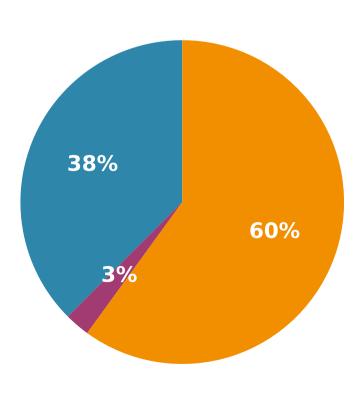


QwQ-32B-no-Thinking-Setting1



GPT-4o-mini-no-Thinking-ELO





Polynomial gains

Unexplained variance

Linear model R²

