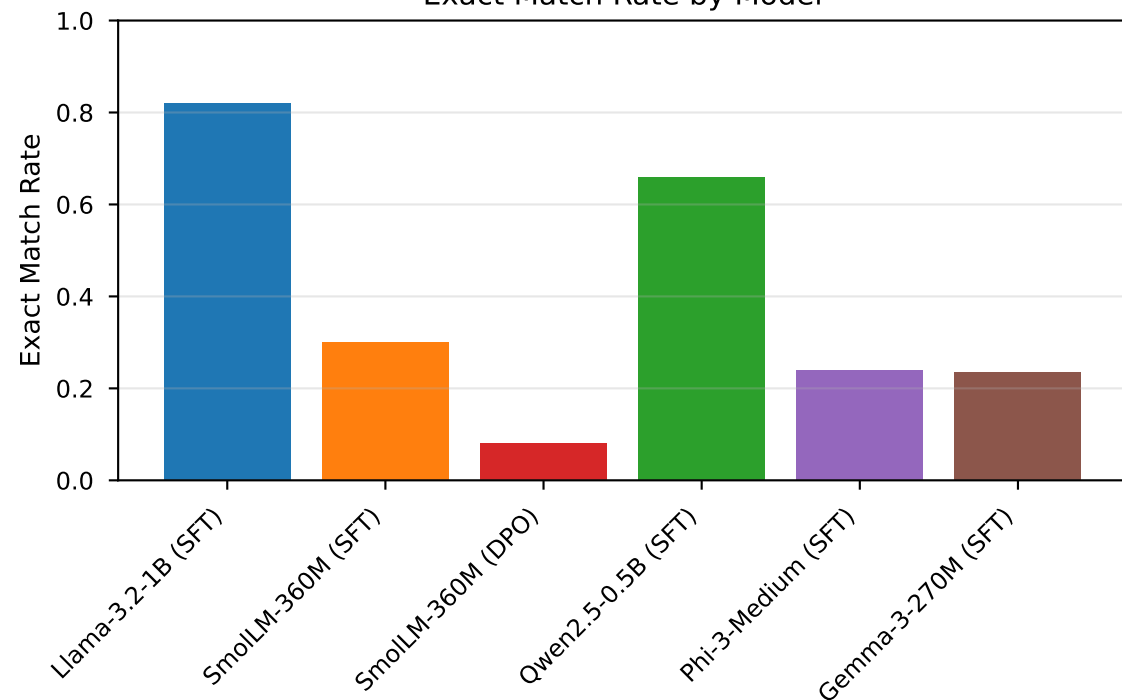
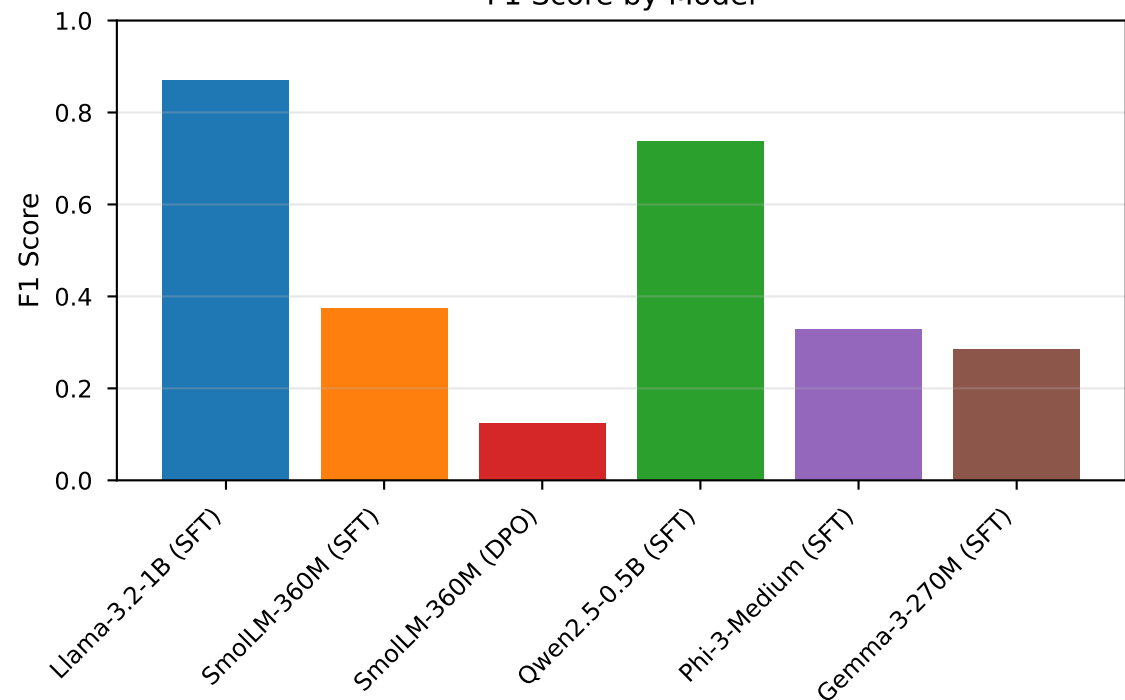


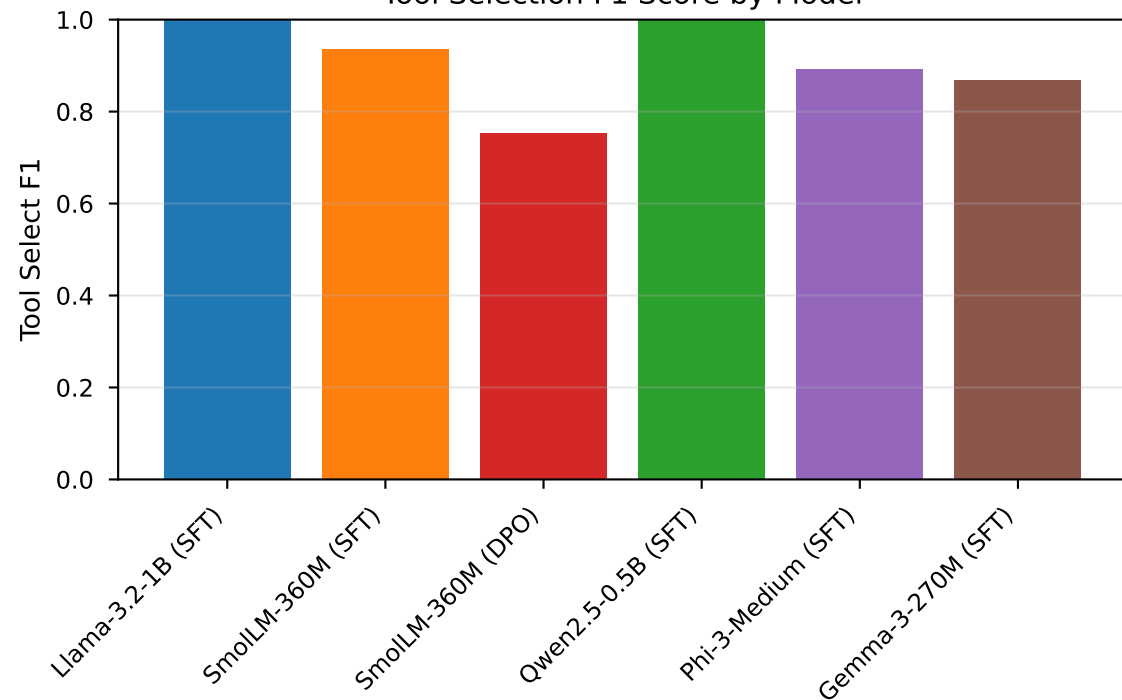
Exact Match Rate by Model



F1 Score by Model



Tool Selection F1 Score by Model



Hallucination Rate by Model

