

Latency vs PPL for different precisions and models

