

Tokens per second vs. parameter count (millions)

