BERT-large, bs=16 Latency(ms) 128:2:32