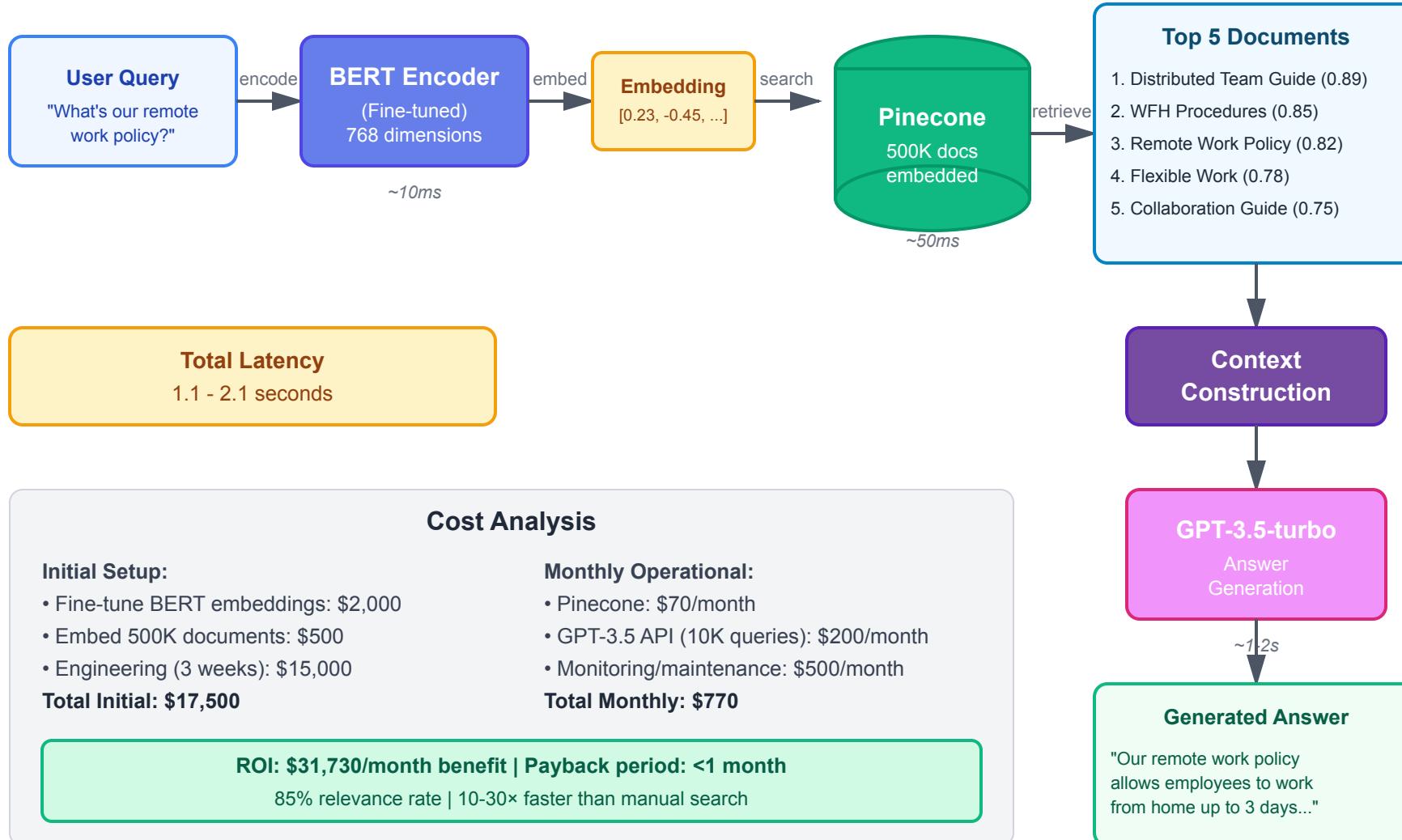


RAG Architecture for Semantic Search



Key Insight: RAG provides 80-95% cost savings vs fine-tuning while maintaining comparable accuracy.

Domain fine-tuning of embeddings (\$2K) yields 17% accuracy improvement—worth the investment.