# 6 Benefits of Switching from Elasticsearch to Vespa.ai

**1) Significant Performance Gains and Cost Efficiency**
- Vespa delivers **3x–9x higher query** throughput than Elasticsearch—**5x** for hybrid, **9x** for vector, and **3x** for lexical searches—all while maintaining **2x–6x lower average latency**.
- **5x infrastructure cost savings** as documented in **this benchmark** report.

**2) Real-Time Data Availability for Instant Updates**
- Elasticsearch is near real time: updates are available only after the next refresh. Vespa ensures that data is immediately searchable upon ingestion.
- Essential for use cases like e-commerce pricing updates, financial data changes, or real-time content publishing.

**3) Superior Scalability Without Bottlenecks**
- Vespa distributes data into fine-grained buckets for improved load balancing, preventing hot spots that arise in Elasticsearch due to uneven shard distribution.
- Vespa scales dynamically while maintaining high availability and predictable performance.

**4) Optimized for AI and Machine Learning Applications**
- Vespa's built-in tensor support enables real-time personalization, recommendations, and vector-based search without external tools, handling multi-vector documents for use cases like ColBERT-style retrieval, visual search, and more.
- Vespa natively supports complex, multi-phase ranking directly on stored data, making it **ideal for large-scale RAG applications**.

**5) Efficient Handling of Updates and Mixed Data Structures**
- Elasticsearch forces full-document rewrites for every update, leading to inefficiencies.
- Vespa supports partial updates and a mix of mutable and immutable data, making updates cheaper and faster.

**6) Simplified Operations and Lower Management Overhead**
- Migrating to Vespa **halved Vinted's server count**, replacing six Elasticsearch clusters with a single deployment while improving search consistency and performance.
- Vespa **halved their server count** while Vinted **increased ranking depth by 3x**, resulting in more relevant search results.

We make AI work for you.          www.Vespa.ai