



Unlocking Relevance with Large Language Models in Search

Matt Riley

General Manager, Search

Search is **rapidly** evolving

HR helpdesk

What are key aspects of the company's 401k policy for an employee in my location? How do I enroll?

Customer Success

Are customers in Dallas buying products with the deepest discounts? Which promotions are doing well?

Ecommerce

What material list and tools do I need to build an irrigation system for 1 acre back yard in Detroit, MI?

Regulatory Compliance

Which recent transactions are not compliant with upcoming 2025 regulatory requirement X?

Corporate Finance

Summarize changes to revenue recognition and expense categorization for FY 2024

Predictive Maintenance

Based on sensor data and recent customer reviews, when do I schedule repair unit X?

Great search is **critical** for great Gen AI experiences

Retrieval technology

Sensitive databases
Multi-system / cloud information
Private knowledge bases
Case histories



Private data



Large Language Models



GenAI

Why retrieval technology **matters**

Great user experience

.....

Build engaging search, enriched by your private data and context

Control costs, hallucinations

.....

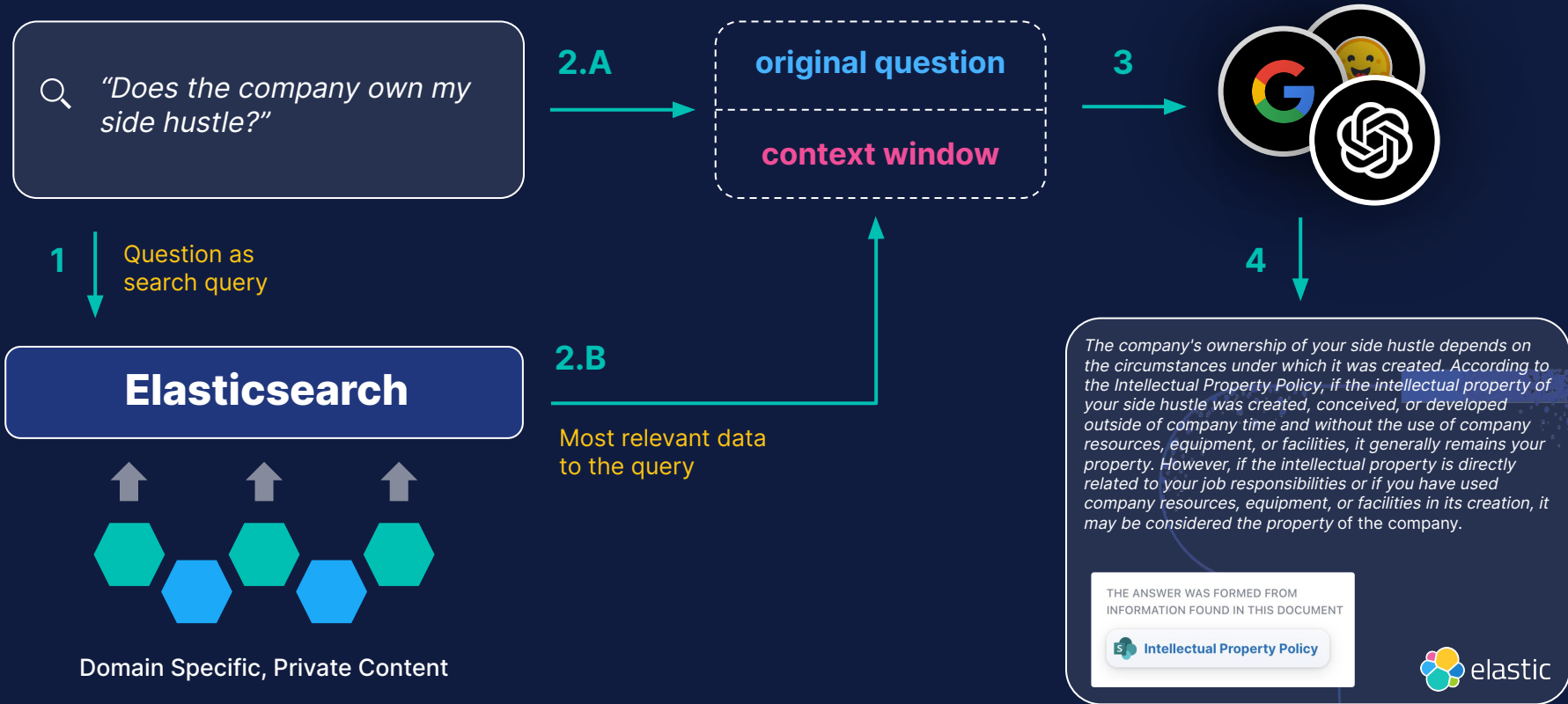
Set context window size, use high relevance ranking to reduce irrelevant context

Enterprise toolkit for developers

.....

Not just vector search, specialized filtering & faceting (geo, time), hybrid search & more

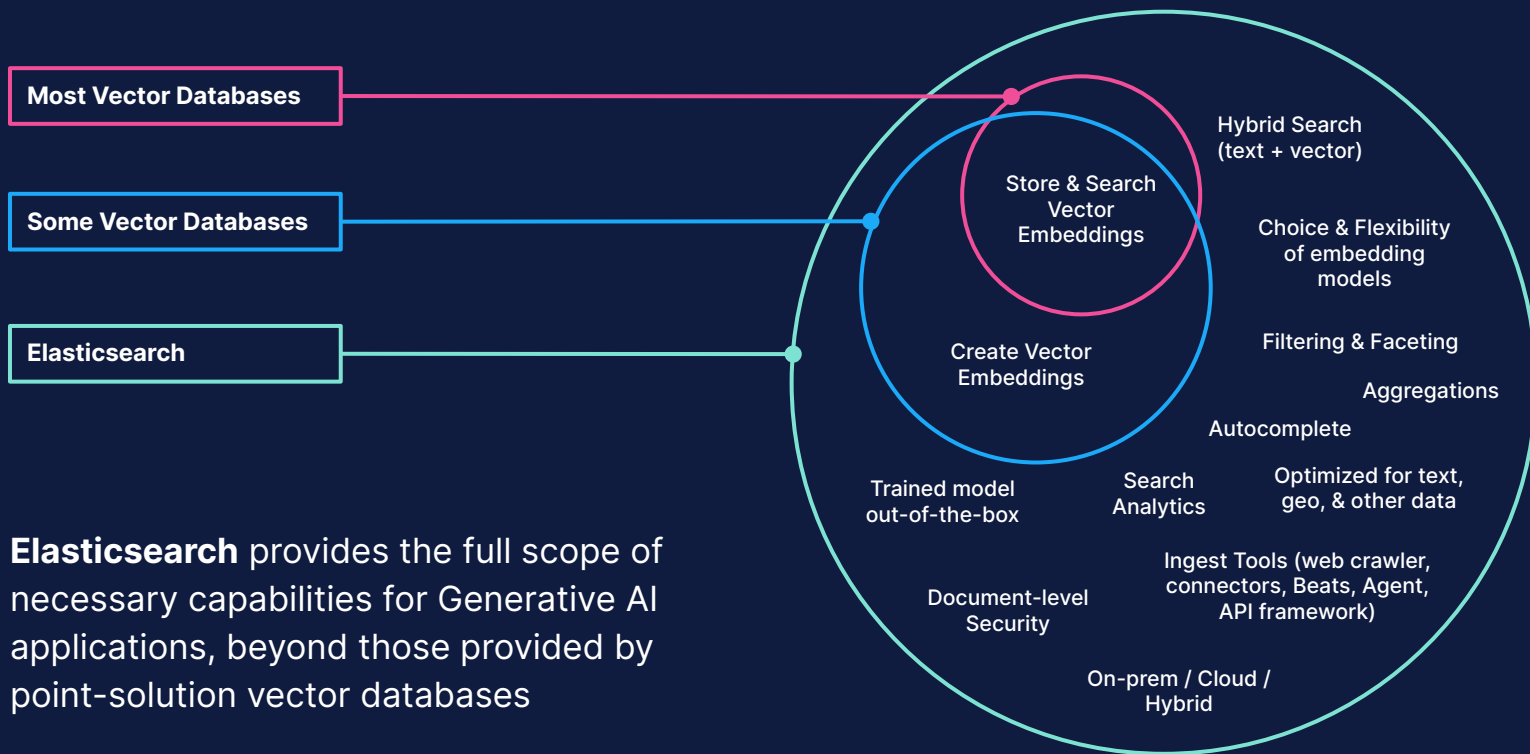
Retrieval augmented generation with Elasticsearch



History of AI innovation at Elastic



Vector search is your starting point...



Elasticsearch provides the full scope of necessary capabilities for Generative AI applications, beyond those provided by point-solution vector databases

Elasticsearch Relevance Engine™ for **Gen AI apps**

Text, vector,
hybrid search

Textual search &
vector database

RRF: hybrid scoring model
(vector & textual search)

Filtering & faceting for kNN
and hybrid queries

Choice of machine
learning models

Host your transformer
models, or use 3rd party
LLMs (OpenAI)

Elastic's proprietary
zero-shot ML model

Integration with 3rd party tooling
like LangChain

Enterprise ready
developer experience

Document and field
Level security

Several major compliance
frameworks covered

Build on-prem, or across cloud
platforms on 50+ regions

Elasticsearch Relevance Engine™ for **Gen AI apps**

Text, vector,
hybrid search

Textual search &
vector database

RRF: hybrid scoring model
(vector & textual search)

Filtering & faceting for kNN
and hybrid queries

Choice of machine
learning models

Host your transformer
models, or use 3rd party
LLMs (OpenAI)

Elastic's proprietary
zero-shot ML model

Integration with 3rd party tooling
like LangChain

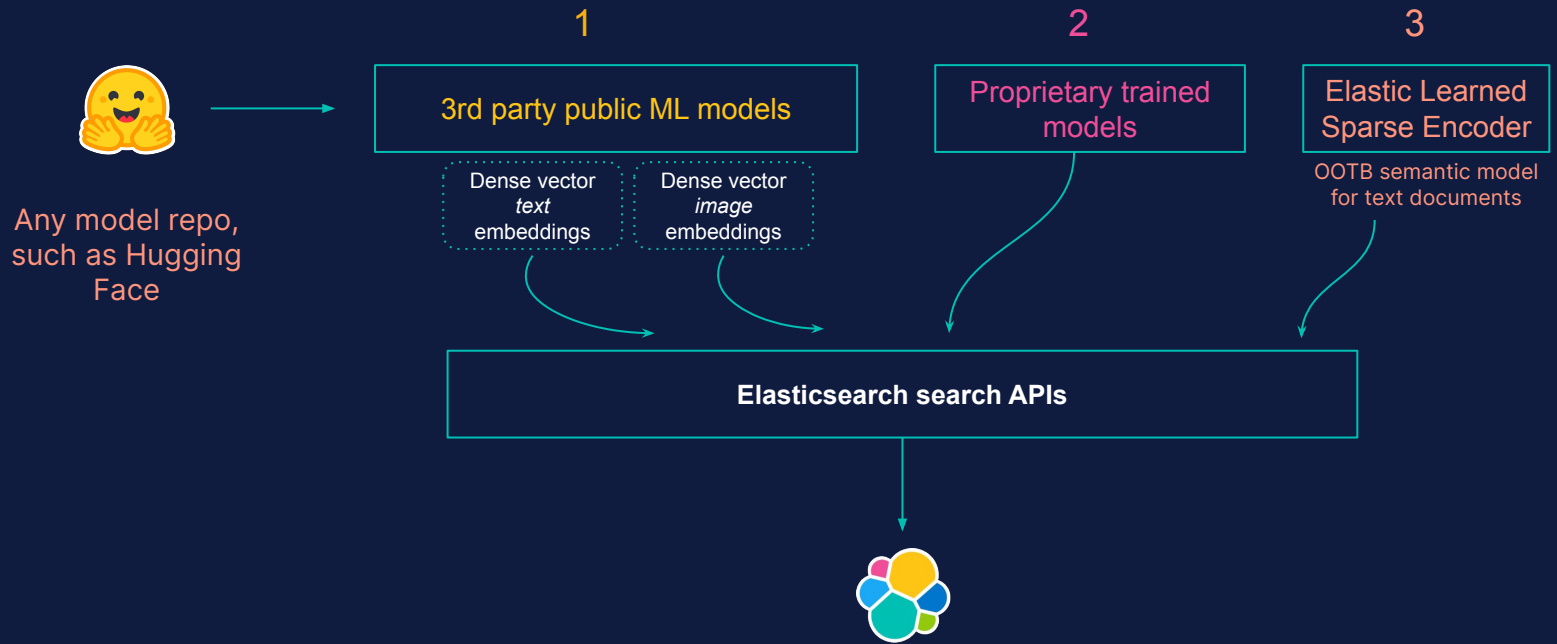
Enterprise ready
developer experience

Document and field
Level security

Several major compliance
frameworks covered

Build on-prem, or across cloud
platforms on 50+ regions

Choose ML models for **your** use-case



Elasticsearch Relevance Engine™ for **Gen AI apps**

Text, vector,
hybrid search

Textual search &
vector database

RRF: hybrid scoring model
(vector & textual search)

Filtering & faceting for kNN
and hybrid queries

Choice of machine
learning models

Host your transformer
models, or use 3rd party
LLMs (OpenAI)

Elastic's proprietary
zero-shot ML model

Integration with 3rd party tooling
like LangChain

Enterprise ready
developer experience

Document and field
Level security

Several major compliance
frameworks covered

Build on-prem, or across cloud
platforms on 50+ regions

Elastic Learned Sparse Encoder: **semantic search** out-of-the-box

Data sets (BEIR benchmark)

Search results ranking models

	Average	TREC-COVID	NFCorpus	NQ	HotpotQA	FiQA	ArguAna	Touche-2020	DBPedia	SCIDOCS	FEVER	Climate-FEVER	SciFact
BM25	0.416	0.688	0.327	0.326	0.602	0.254	0.472	0.347	0.287	0.165	0.649	0.186	0.69
RRF (BM25/Dense)	0.449	0.697	0.317	0.445	0.611	0.318	0.474	0.354	0.353	0.159	0.746	0.238	0.671
Linear (BM25/Dense)	0.471	0.787	0.335	0.485	0.62	0.341	0.444	0.346	0.378	0.164	0.778	0.272	0.698
SPLADE	N/A	0.726	0.347	0.537	0.687	0.347	0.526	0.246	0.436	0.158			0.703
ELSER	0.471	0.747	0.351	0.524	0.67	0.339	0.5	0.263	0.415	0.156	0.777	0.218	0.695
RRF (BM25/ELSER)	0.478	0.797	0.352	0.468	0.674	0.311	0.497	0.347	0.411	0.166	0.762	0.24	0.712

Search quality *worse* than benchmark (BM25)



Search quality *better* than benchmark (BM25)

Want to learn more?

Search "[elasticsearch sparse encoder blog](#)"

What are our
customers **building**?



Gen AI apps powered by search

Search for
Customer service



"Feedback from our engineers is extremely positive. They now use Topic Search to solve 90% of service requests. They can deliver a better customer experience by easily finding on-target information and fixing issues much faster than before."

Sujith Joseph, Principal Enterprise Search & Cloud Architect, Cisco Systems

Legal e-Discovery
search



"I'm thrilled about the benefits we can bring to customers through our investments to harness Elasticsearch within RelativityOne. We're excited about the potential to deliver powerful, AI-augmented search results to our customers."

Chris Brown, Chief Product Officer, Relativity

Where do we go from here?

(**Spoiler alert:** Full set of gen AI dev tools, and more than being the best vector database out there)



Coming soon, to a Gen AI search app near you...

Inference API

Pick and mix service providers that give you the dense vector embeddings you want

3rd party integrations

More integrations to popular 3rd party projects, resulting in easier AIOps for you

Result reranking capabilities

Making our text + vector hybrid search story even stronger

Int8 quantization and memory reduction

Reduction in disk usage, near-zero accuracy loss, and faster search

Passage Vectors

Automatic diversification over parent doc, complex passage and parent filtering with Lucene primitives

Improving Retrieval Quality for Search

9/27, Wednesday, 2 pm,
Downstairs, Fisher



Priscilla Parodi,
Principal Developer Advocate, Elastic

Visit our booth for a demo!  elastic