

Hybrid Search Lessons Learned



Mohit Sidana

Haystack EU 2025

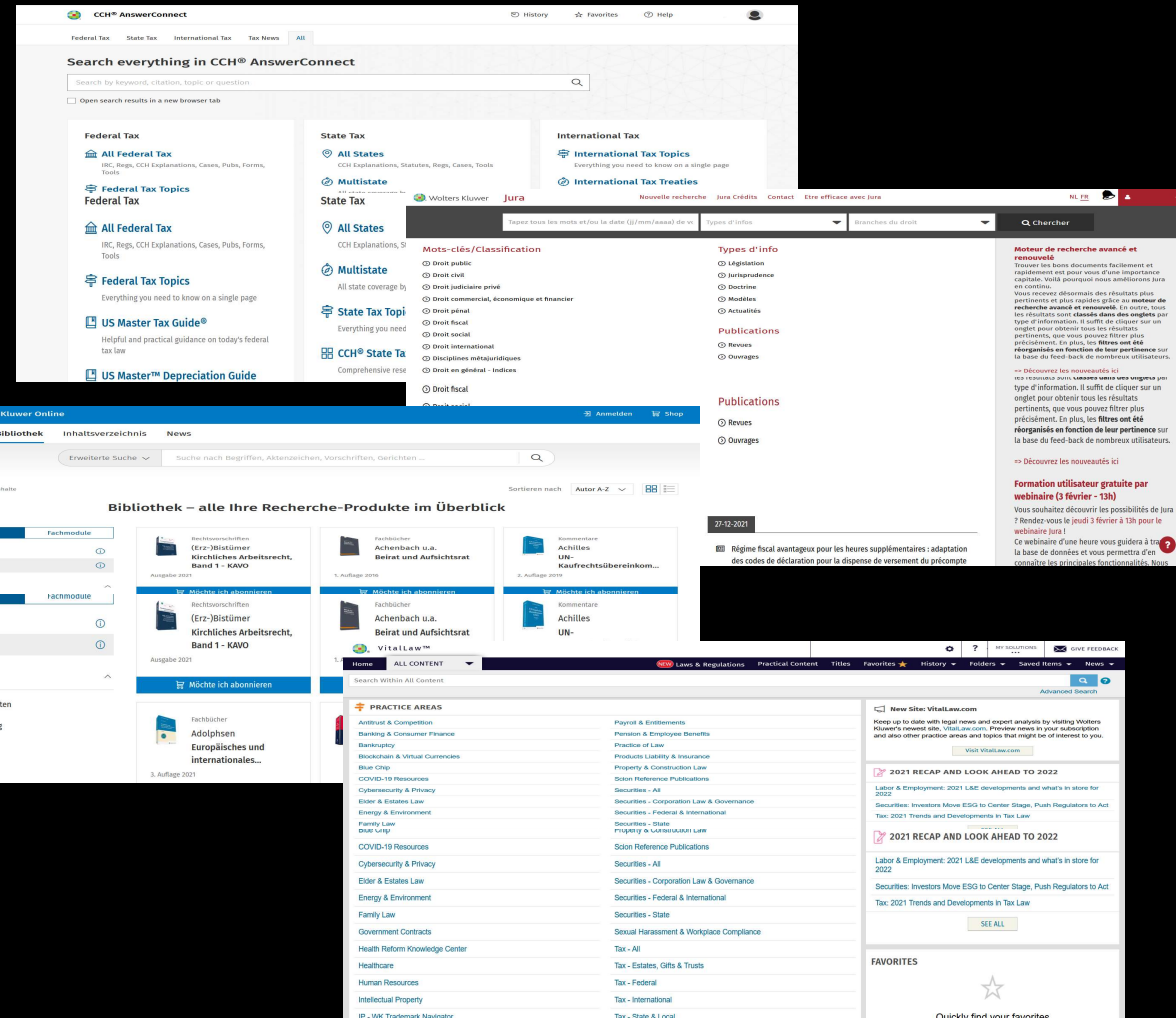
Tom Burgmans



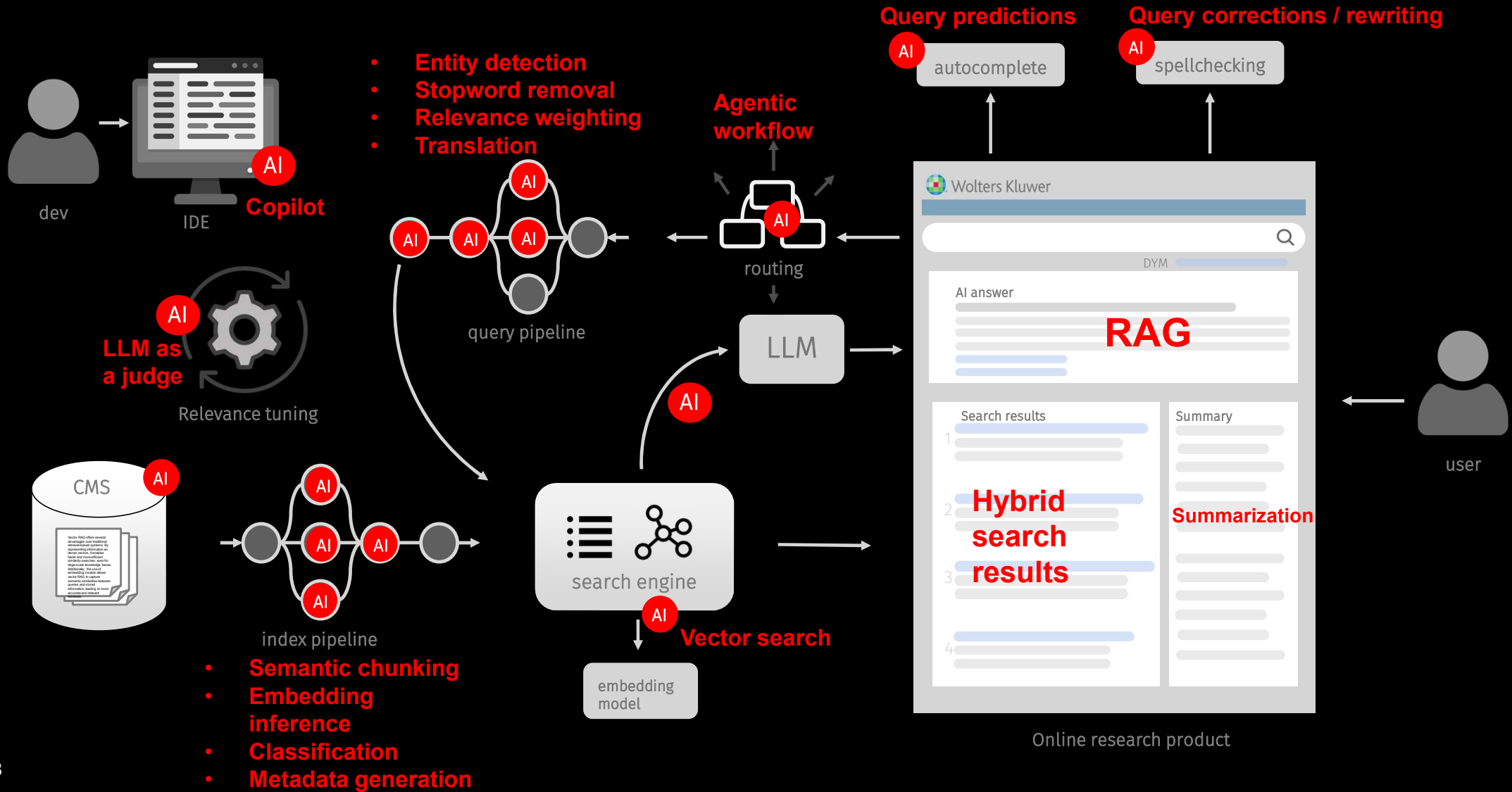
Information services provider for experts in legal, business, tax, accounting, finance, audit, risk, compliance, and healthcare



- 21600 Employees
- Active in over 180 countries
- Mostly subscription based online products



AI-ification of the Search Domain



Researching query latency of **vector search** in Solr at scale



Index Segments

Minimizing the number of index **segments** reduces overhead, enhancing the query performance.



K-Value

High **k-values** directly increase query time; select only those essential for relevance.



Vector Dimensions & Quantization

Dimensionality reduction with **Matryoshka** embeddings preserves key data while saving space. **Scalar quantization** reduces the vector precision to speed up the queries and reduce memory usage.



Filters

Search **filters** can significantly increase query latencies due to processing overhead.



Vectorizing content

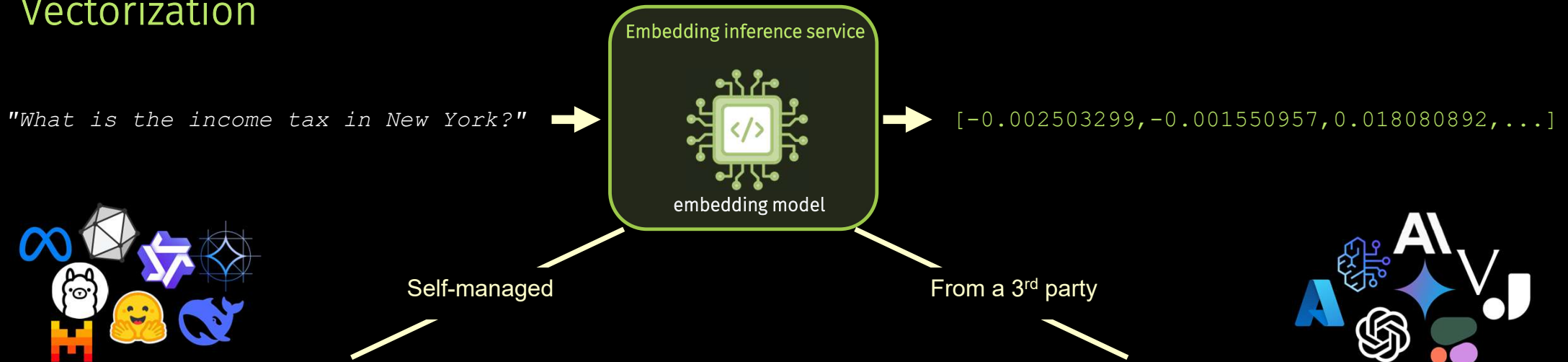
Another **bottleneck** might be the **external vectorizer**.



Dot Product Similarity

For efficient scoring calculations, use **dot-product** similarity. Don't forget to **normalize vectors** to unit length 1

Vectorization



- Use **GPU** over CPU
- Have enough **VRAM** to hold the complete model + 10-20% extra
- **Smaller** models are faster
- **Low dimensional** models are faster
- Models could be **quantized** too. Quantized models are faster
- Convert models to **ONNX**
- **Short input texts** are faster
- Embedding inference service should be **closely hosted** to other connected services to limit network latency

- **Costs:** hosting infrastructure and operations
 - Efficient for **high volume** tasks

- Test the **throughput!** Understand what throttles it
- For more throughput capacity, **multiple deployments** in different regions (+ load balancer) may be needed
- Models from same provider could widely differ in **quality**
- Quantization and truncation could be done at **service side**
- Choose a **region close** to where connected services are hosted

- **Costs:** fee per M tokens
 - Efficient for **low volume** tasks

Optimizing Indexing : When to Re-use Vectors

- Vectorization is computationally intensive
- Re-use vector embeddings to save resources
- Reduce latency where system need to handle frequent document metadata updates in the content
- Chunked text remains unmodified during metadata updates



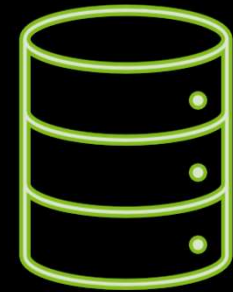
Identify metadata updates to determine vector reuse opportunities.



Verify text changes swiftly using checksums or hashes.



Decide whether to reuse existing vectors or Recompute.

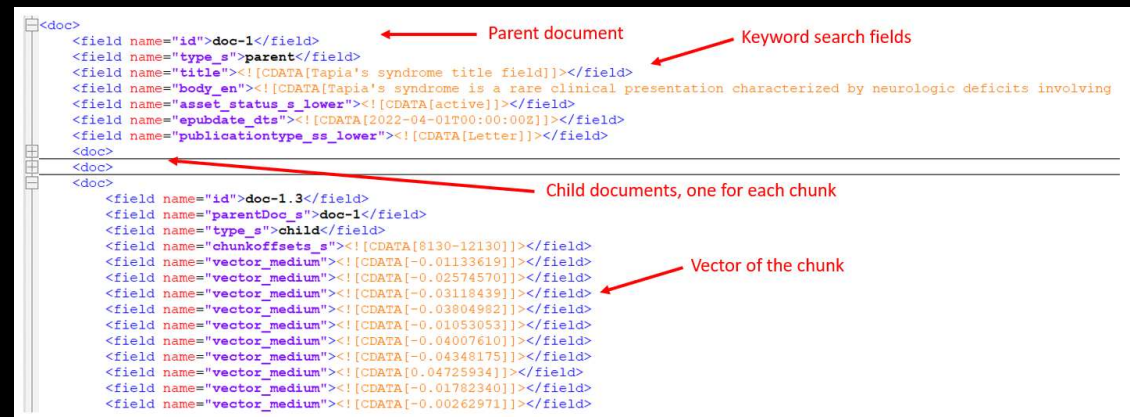
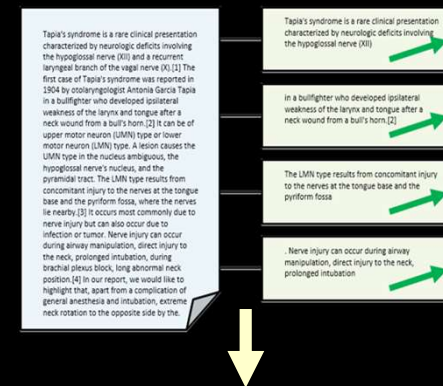


Update vector store with reused or recomputed vectors accordingly.

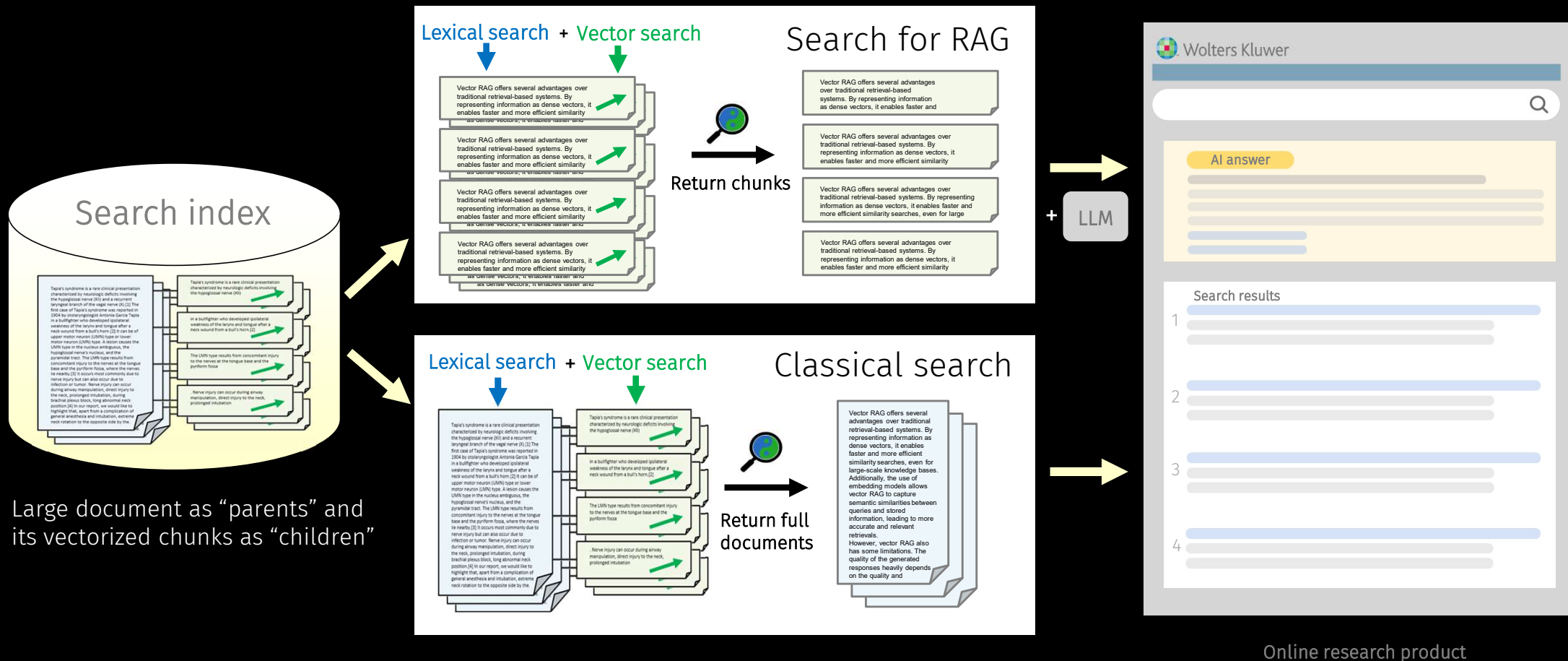
Setting up a multi-purpose Hybrid Search Index

Preparing the content set

Nested document hierarchy with parent and child documents



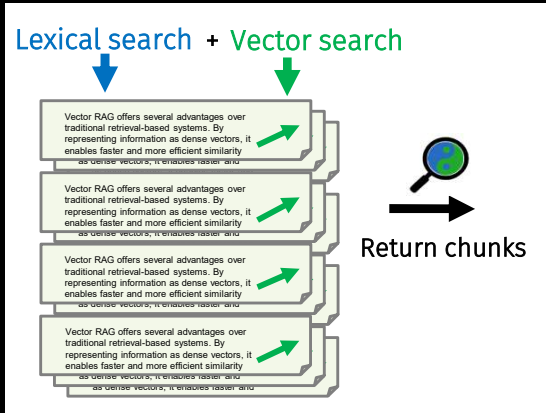
A nested index could serve multiple **hybrid search** use cases



Large document as “parents” and its vectorized chunks as “children”

Hybrid search Solr syntax

Search for RAG



```
params:{
  uf:"*_query_",
  q:"{!bool filter=$hybridlogic must=$hybridscore}",
  hybridlogic:"{!bool should=$kwq should=$vectorq}",
  hybridscore:"{!func}sum(product($kwweight,$kwq),product($vectorweight,query($vectorq)))",
  kwq:"{!type=edismax qf=\"chunk_body\" v=$qq}",
  qq:"What is the income tax in New York?",
  vectorq:"{!knn f=vector_field topK=10}[-0.002503299,-0.001550957,0.018080892,...]",
  kwweight:1,
  vectorweight:4
}
```

allow embedded Solr queries

OR-relationship between lexical search & vector search

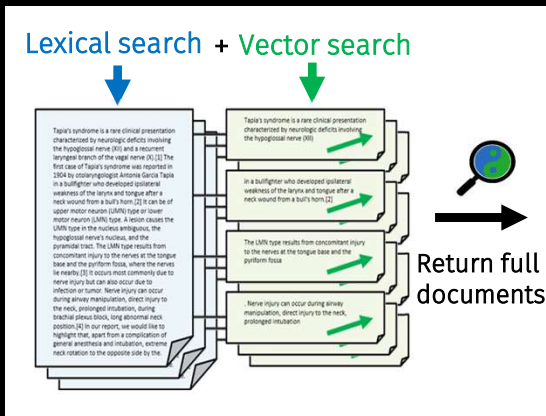
sum the scores

typical lexical search

balance the impact of keyword matches vs vector matches

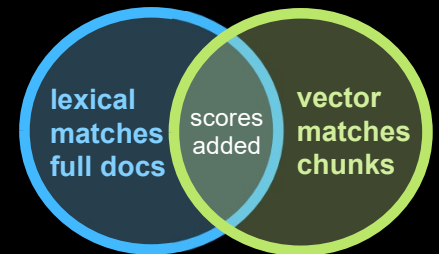
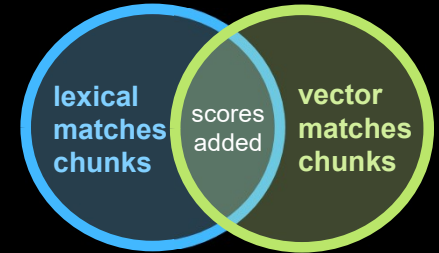
vector representation of the query (externally vectorized)

Classical search



```
params:{
  uf:"*_query_",
  q:"{!bool filter=$hybridlogic must=$hybridscore}",
  hybridlogic:"{!bool should=$kwq should=$vectorq}",
  hybridscore:"{!func}sum(product($kwweight,$kwq),product($vectorweight,query($vectorq)))",
  kwq:"{!type=edismax qf=\"full_doc_body full_doc_title^3\" v=$qq}",
  qq:"What is the income tax in New York?",
  vectorq:"{!parent which=\"type_s:parent\" score=max v=$childq}",
  childq:"{!knn f=vector_field topK=10}[-0.0034859276,-0.028224038,0.0024048693,...]",
  kwweight:1,
  vectorweight:4
}
```

block join: vector search children, return parents



Relevancy tuning



Blending lexical results with vector results is like mixing water and oil



	Lexical search	Vector search
Scoring algorithm	TF-IDF (BM25)	Vector distance (i.e. cosine angle)
Debug score EXPLAIN	Yes	No
Effect of boosting	On the entire result set	Only the top k results
Number of results	Between 0 - all documents	Between 0 - k (in most cases exactly k)

Relevancy tuning

Tuning relevance for LLMs (RAG) is not equal to tuning for humans

	For humans 	For LLMs 
Unit of Retrieval	Full documents	Document chunks
Consumption Mode	Scroll, KWIC, Facets, Paginates	Top N chunks
Result list importance	Ordering (<i>precision</i> of the top results)	Coverage (<i>recall</i> in the top N)
Tuning goals	Engagement, diversity	Informativeness, factuality

Balancing weights in
Hybrid search is
dependent on the nature
of the query

Skip vector search with:

- “*” search or wildcarded queries
- Boolean constructions
- Explicit phrase search
- Fielded searches
- Proximity queries

Hybrid search

Lexical match score



Vector match score

Citation

Case nickname

Multi-lingual

Typos

Case
summary

Where we'll expect a lot of 'tuning' energy will be spent

User's query

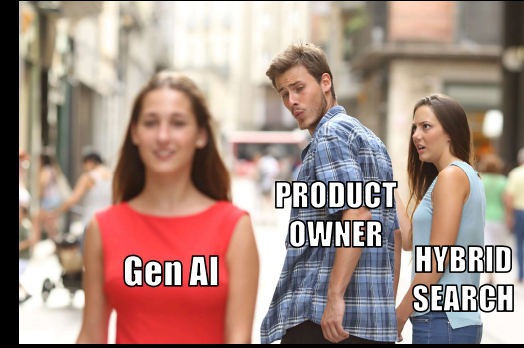
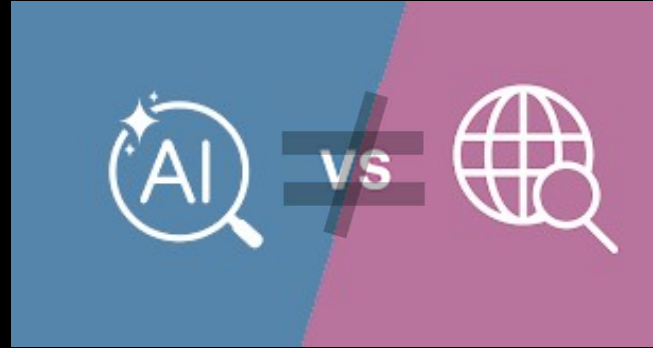
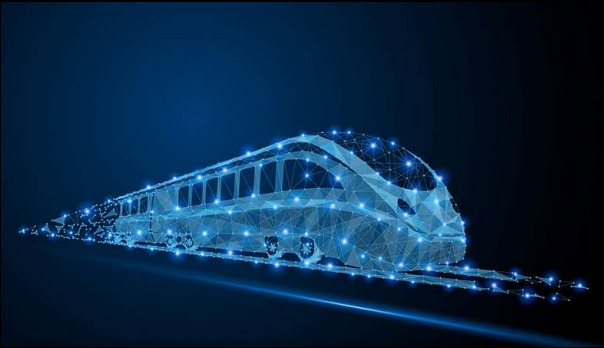
+ history
+ context



- (Hybrid) search syntax
- Boosting weights
- Auto filters
- Rewriting
- Expansions
- LLM prompt

The Fabulous Agentic Intelligent Query Interpreting and Semantically Smart Omni-Dimensional Hyper-Adaptive Cognitive-Neural Orchestrator of Dynamic Search Syntax Generation (or in short: FAIQISSODH-AC-NODSSG)

Enrolling **hybrid search** in a large organization: challenges



“Will AI replace search?”

“Search is not developing anymore”

“Why not just train a model?”

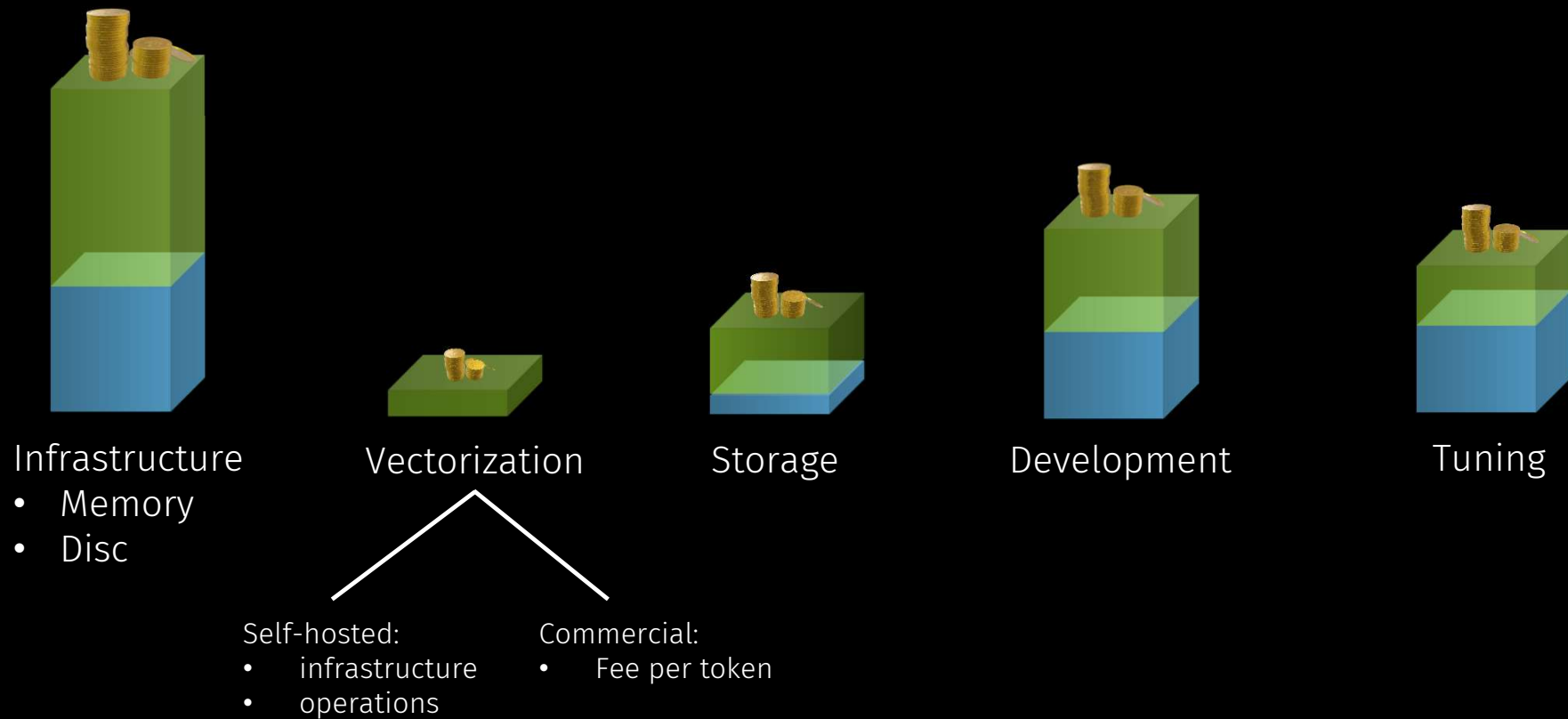
“Future AI systems won’t need search for RAG”



Enrolling **hybrid search** in a large organization

1. **Create interest & awareness**
 - Emphasis on the benefits
 - Short videos
2. **Prototype**
 - Real use cases
 - Shareable
3. **Education & Transparency**
 - Deep(er) dive in the pro's & con's
 - Cost impact
4. **PoC**
 - Integrated in product
 - Go/No-Go moment
5. **Productionization**
 - Start small
 - A/B
6. **Adoption & Enablement**
 - Measure usage
 - Keep evolving

Cost factors lexical -> hybrid



Hybrid search prototype

Wolters Kluwer

Hybrid search demo application

Results: 100 K 50 useCase: KLI Type: hybrid

Hybrid search criteria:
 Keyword Weight: 46 Vector Weight: 954

South Sudan construction of electric power infrastructure dispute

Total Results : 27

Position : 1
ID : kli-ka-kaces-4692
Active Partners Group Limited v. Republic of the Sudan, PCA Case No. 2013-04, 27 January 2015
Score : 1,682.133

Position : 2
ID : kli-ka-kaces-10875
Energoinvest DO v. Democratic Republic of the Congo, Société Nationale d'Electricité (S.N.E.L.), 11441/KGA, 20 April 2003
Score : 1,477.798

Position : 3
ID : kli-ka-pcaas-13-005-n
The Republic of the Philippines v The People's Republic of China (The South China Sea Arbitration) (Award), PCA Case No. 2013-19, 12 July 2016
Score : 981.635

Position : 4
ID : kli-ka-kaces-11286
Frazer Solar GmbH v. Kingdom of Lesotho, Wolters Kluwer PACER ID PR-0035679, 28 January 2020
Score : 884.191

Position : 5
ID : kli-ka-psm-ita-awards-044-n
Balkan Energy Limited v. The Republic of Ghana (Award on the Merits), PCA Case No. 2010-7, 1 April 2014
Score : 881.503

Position : 6

Manual slider to balance the impact of lexical search vs vector search

Background colors indicate:

Lexical match

Vector match

Hybrid match

Tip: how to let Solr tell with what search technique it found its results (see syntax slide 9):

`fl:...,keywordmatch:exists($kwq),vectormatch:exists($vectorq),...`

Hybrid search requires explaining

Warning: counts for facets/total results may behave unintuitive

income tax Search here within the 1,274,311 results for "income tax"

☐ INCLUDE AI GENERATED ANSWER

FILTER RESULTS

— PRACTICE AREAS

- ☐ Antitrust & Competition (1,265)
- ☐ Banking & Consumer Finance (4,547)
- ☐ Bankruptcy (2,372)

[Show All](#)

— TAX ESSENTIALS

- ☐ Federal Topics (866)
- ☐ State Topics (62)
- ☐ International Topics (5,032)

— QUICK REFERENCE BY COUNTRY

- ☐ International Topics (5,032)

— DOCUMENT TYPE

- ☐ Treatises (13,606)

1,274,311 results for "income tax"

☐ Select All

☐ 1. International Encyclopaedia
Published January 1, 2004
Author: Elizabeth Toomey
...These athletes must pay income tax (below)....Independent contractors pay income taxes on their liabilities, while employees must pay income taxes on the totalization, race winnings...

☐ 2. Pension and Welfare Benefits Appeals, Eleventh Circuit, (A
Issued August 8, 2025
Court:U.S. Court of Appeals, Eleventh Circuit
...The company provided its employees with their federal income taxes, and...

income tax Search here within the 1,214 results for "tax"

☐ INCLUDE AI GENERATED ANSWER

Filters: Federal Topics

FILTER RESULTS

— PRACTICE AREAS

- ☐ Tax - Federal (1,214)

— TAX ESSENTIALS

- ☒ Federal Topics (1,214)
- ☐ State Topics (109)
- ☐ International Topics (5,077)

— DOCUMENT TYPE

- ☐ Topic Pages (1,214)

— JURISDICTION

- ☐ Federal (1,214)

1,214 results for "income tax"

☐ Select All

VitalLaw Smart Chart Results

Physical Presence Requirement, Economic Interests, Telecommuter or Work-Fr

51 Jurisdiction(s): AL, AK, AZ, AR, CA, ND, OH, OK, OR, PA, RI, SC, SD, TN, TX

☐ 1. Clergy
...Taxation of a clergy member forbid the payment of taxes on self-employment tax purpose employer must withhold income contractor....

Thank
you!

Questions?

Or contact us later:



Mohit Sidana
Search Architect
Wolters Kluwer



Tom Burgmans
Principle Product Software Engineer
Wolters Kluwer