

## Atlas Workshop

Hands-On Introduction





## Requirements:

- Personal Laptops
- MongoDB Atlas account with personal email
- Connect to personal-wifi

## Exercise 0

Create MongoDB Cluster → M0 Add Sample Data

### Exercise 1

Rich Queries - MQL

### •

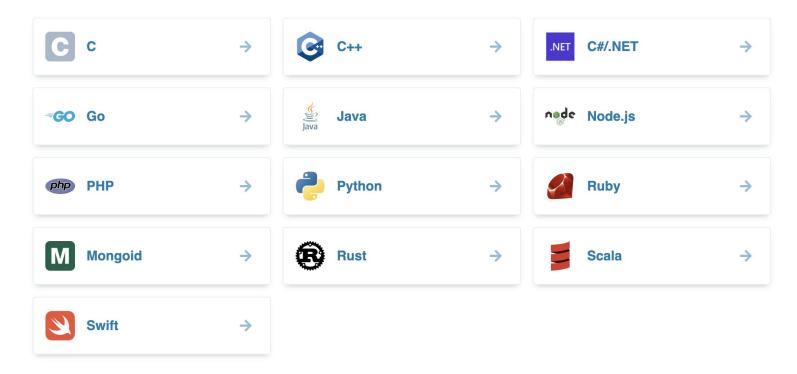
# MongoDB Query API: rich and expressive

Expressive queries	<ul> <li>Find anyone with phone # "1-212"</li> <li>Check if the person with number "555" is on the "do not call" list</li> </ul>
Geospatial	Find the best offer for the customer at geo coordinates of 42nd St. and 6th Ave
Text search	Find all tweets that mention the firm within the last 2 days
Aggregation	Count and sort number of customers by city, compute min, max, and average spend
Native binary JSON support	<ul> <li>Add an additional phone number to Mark Smith's record without rewriting the document</li> <li>Update just 2 phone numbers out of 10</li> <li>Sort on the modified date</li> </ul>
JOIN (\$lookup)	Query for all San Francisco residences, lookup their transactions, and sum the amount by person
Graph queries (\$graphLookup)	Query for all people within 3 degrees of separation from Mark

#### MongoDB

```
customer_id : 1,
 first_name : "Mark",
 last_name : "Smith",
 city : "San Francisco",
 phones: [
         number: "1-212-777-1212",
         type : "work"
    },
         number: "1-212-777-1213",
         type : "cell"
```

#### Intuitive: client drivers

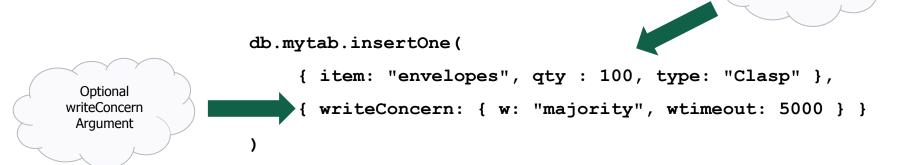


- Common CRUD capabilities but idiomatic to each language
- Uniform HA & Failover capabilities across all

### insertOne()

- Inserts document(s) into a collection
- Can optionally define the writeConcern of the operation

\_id field is automatically generated with an ObjectId



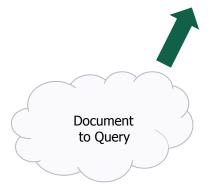


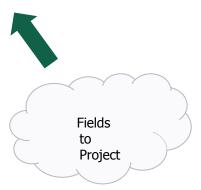
Document to Insert

### find() and findOne()

- Get documents from the collection (table)
- Can filter by content (query) and by what is returned (project)
- findOne() returns just the first document

```
db.mytab.find(\{x:1, y:2\}, \{a:1, b:1, c:1\})
```







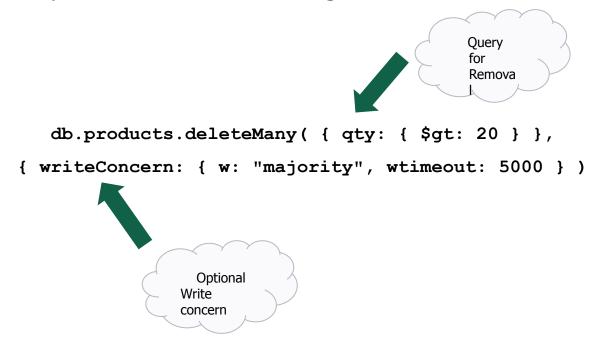
### updateMany() and updateOne()

- Modifies existing document(s) in a collection
- Use with a field update operator (MQL)
- Query for the documents you want to update
- Can optionally upsert, define writeConcern, etc.

```
db.books.updateOne(
    { author: "Tom Gleitsmann" },
    { $set: { author: "Thomas Gleitsmann" }},
    { upsert: true }
    Optional
    3rd
    Argument
```

### deleteMany()

- Deletes a document from the collection
- Takes an optional writeConcern argument





### **Comparison Query Operators**

```
Exists and less than
$It :
$lte:
        Exists and less than or equal to
$gt:
        Exists and greater than
            Exists and greater than or equal to
$gte
$ne:
        Does not exist or does but not equal to
$in :
        Exists and in a set
$nin
            Does not exist or not in a set
```

{ age : { \$gte : 21 }}

temp : { \$gt: 10, \$1t: 30 }}



### Field Update Operators (well some...)

\$currentDate : Sets the value of a field to current date
\$inc : Increments value
\$min : Update if field is less than \$min value
\$max : Update if field is greater than \$max value
\$mul : Multiply field value
\$set : Set value of field

: Removes field from document

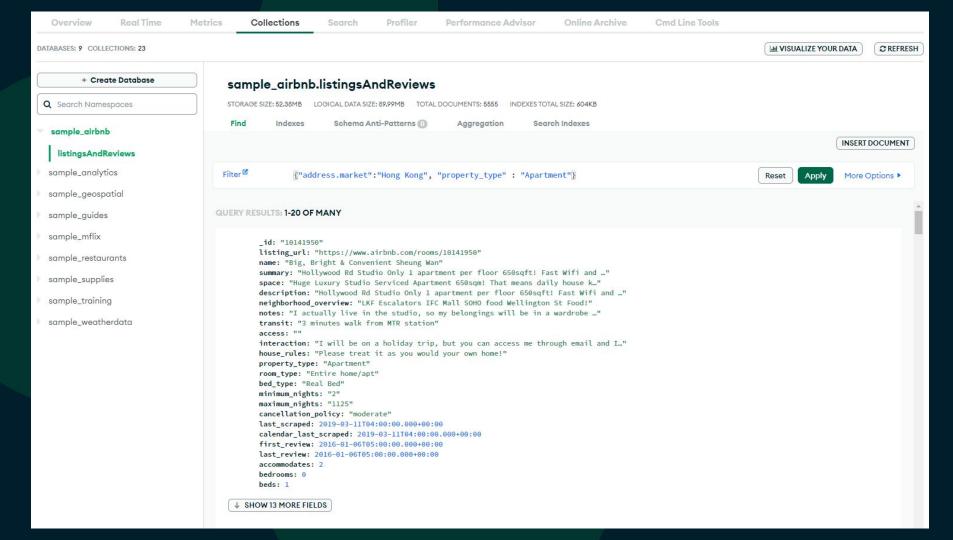
{ \$set : { bestDB : "MongoDB"}}
{ \$mul: { amount : 5, ... } }

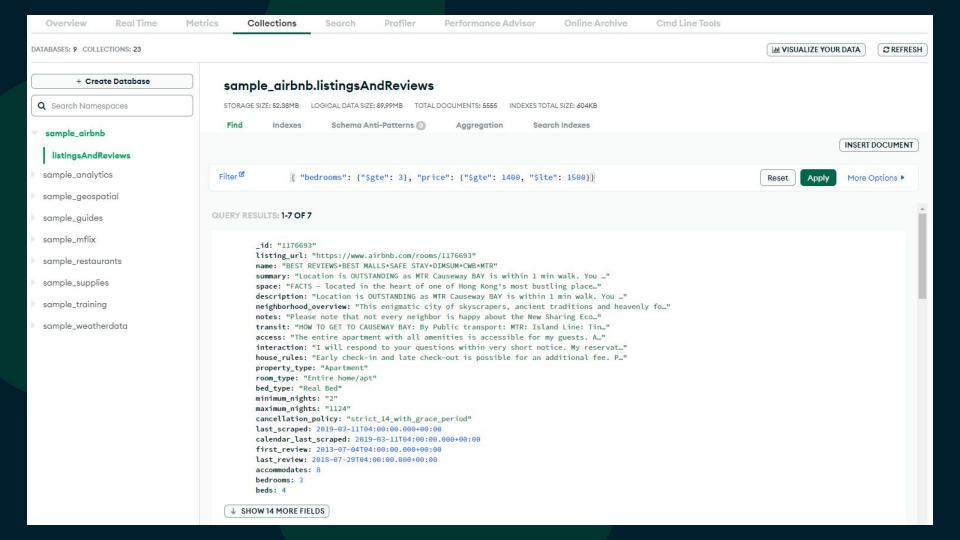
\$unset

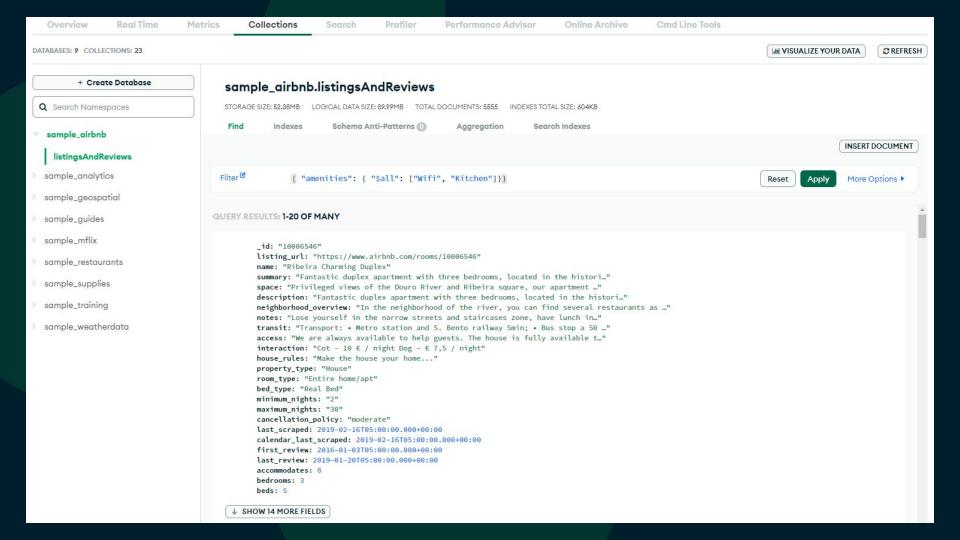


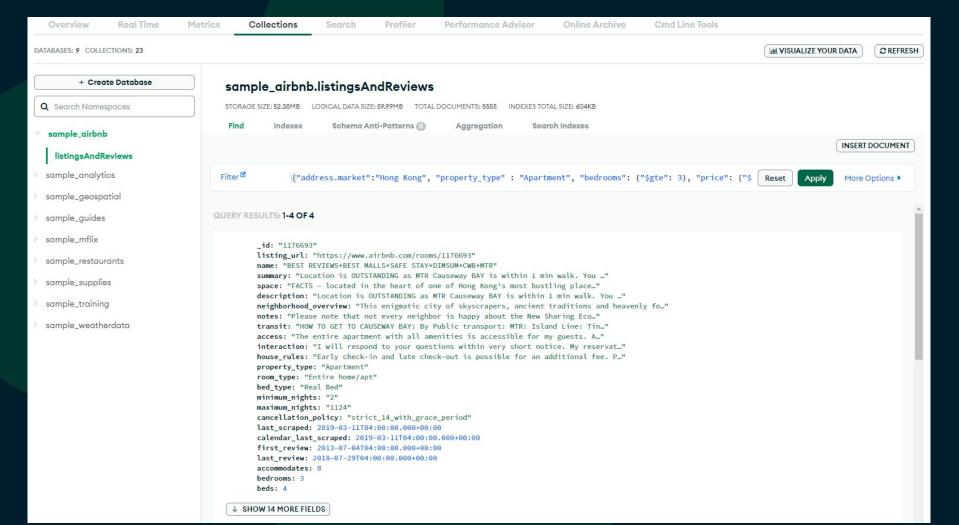
Hands-On: Data Explorer

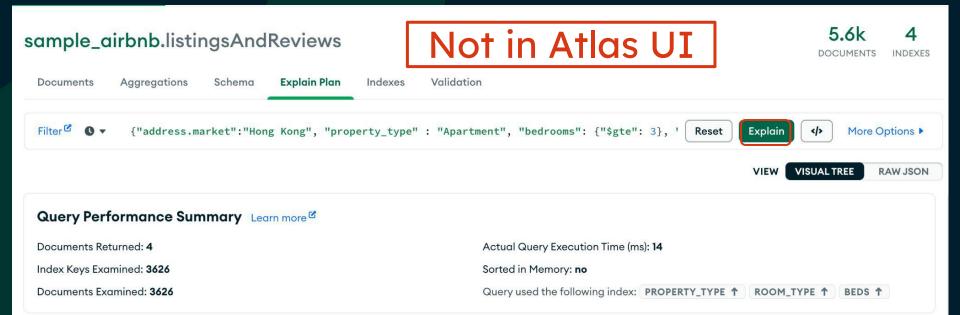
"Browse Collections" sample\_airbnb.listingsAndReviews













5 Minute Break

### Exercise 2

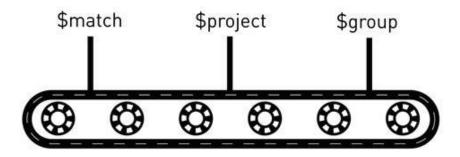
Aggregation Framework

### Aggregations 101

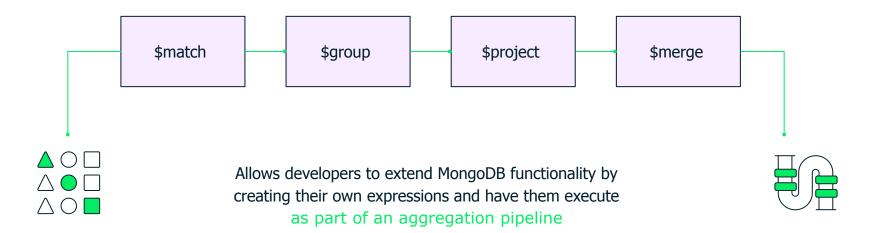
- Until now, we have only filtered data
  - What Documents
  - What Fields
  - What Order
- Aggregation allows us to compute new data
  - Calculated fields
  - Summarized and grouped values
  - Reshape documents

### **Aggregation Stages**

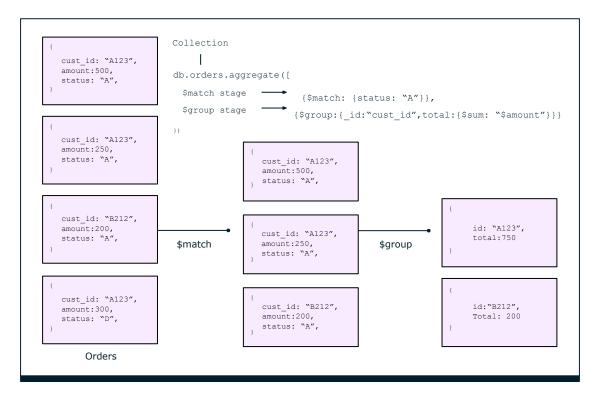
- Each transformation is a single step known as a stage.
- Compared to one huge SQL style statement this is
  - Easier to understand
  - Easier to debug
  - Easier for MongoDB to rewrite and optimize



### **Aggregation Pipelines**







### Stages we know

```
    $match equivalent to find(query)
    $project equivalent to find({},projection)
    $sort equivalent to find().sort(order)
    $limit equivalent to find().limit(num)
    $skip equivalent to find().skip(num)
    $count equivalent to find().count()
```

When these are used at the start of a pipeline, they are transformed to a find() by the query optimizer.

### Comparing Aggregation Syntax

#### Find the name of the host in Canada with the most "total listings":

#### Using find()

```
db.listingsAndReviews.find(
          {"address.country":"Canada"},
          {"_id":0, "host.host_total_listings_count":1, "host.host_name":1}
).sort({"host.host_total_listings_count":-1}).limit(1)
```

#### Using aggregate()

### Dollar Overloading

```
{$match: {a: 5}}
Dollar on left means a stage name - in this case a $match stage
{$set: {b: "$a"}}
Dollar on right of colon "$a" refers to the value of field a
{$set: {area: {$multiply: [5,10]}}
$multiply is an expression name left of colon
$map:{input: "$quizzes", as: "grade",in: { $add: [ "$$grade", 2 ] }}}
$$grade refers to a variable in $map statements
```

Think as a programmer - not as a database shell.

Define variables. Doing this helps you keep track of brackets.

It also helps you really understand the Object concept.

```
//Do it THIS way for ease of testing and debugging
> no_celebs = {$match:{"user.followers_count":{$lt:200000}}}
> name_only = {$project:{"user.name":1, "user.folowers count":1, id:0}}
> most_popular = {$sort: {"user.followers count":-1}}
> first_in_list = {$limit:1}
> pipeline = [no_celebs,name_only,most_popular,first_in_list]
> db.twitter.aggregate(pipeline)
```

Hands-On: Data Explorer Pt. 2

"Browse Collections"
sample\_airbnb.listingsAndReviews
"Aggregations"

### Exercise 1



Using the Atlas Aggregation Builder and the Airbnb listings data (sample\_airbnb.listingsAndReviews) read the following question and output the answer. You will need to use a \$set stage and a \$project stage with expressions.

To keep in mind...

- 1. The price for the basic rental is in the **\$price** field; for that price, you can have the number of guests provided in **\$guests\_included** field.
- 2. The property may take more guests in total (maximum guests is in **\$accommodates** property).
- 3. You have to pay **\$extra\_people** dollars per person for every person more than **\$guests\_included**

Requirement:

- Calculate how many extra guests you can have for each property and add that as a field using \$set
- Calculate how much it would cost with these extra guests. \$project the basic price and the price if fully occupied with \$accommodates people.

### ANSWER: Exercise 1



#### Using the Airbnb listings data:

 Calculate how much it would cost for these extra guests - they cost "\$extra\_people" each and \$project the basic price and the price if fully populated

### Exercise 2

Calculate how many Airbnb properties there are in each country, ordered by count, list the one with the largest count first. You should continue working with namespace sample\_airbnb.listingsAndReviews.

 Hint: to count things, you add the explicit value 1 to an accumulator using \$sum



Calculate how many Airbnb properties there are in each country, ordered by count, list the one with the largest count first.

```
[{$group: {
    _id: "$address.country",
        count: { $sum: 1 },},},
        {$sort: { count: -1 },},]
```

### Exercise 3

What are the top 10 most popular amenities across all listings?



### ANSWER: Exercise 3

•

What are the top 10 most popular amenities across all listings?

```
[{$unwind: { path: "$amenities" },},
    {$group: {
      _id: "$amenities",
        count: {$sum: 1},},},
    {$sort: {count: -1,},},]
```

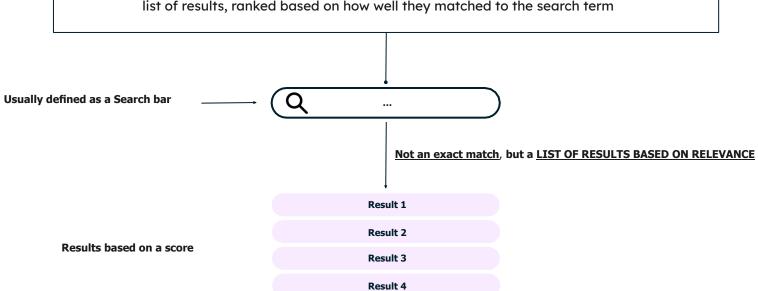
5 Minute Break

Atlas Search

### What is Search?

•

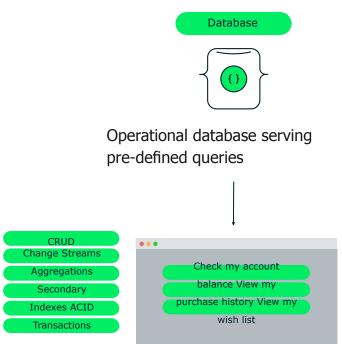
Search, or "Full-Text Search" is the ability to search across all of your data and efficiently return a list of results, ranked based on how well they matched to the search term







Stand up a database to provide the application's persistence layer

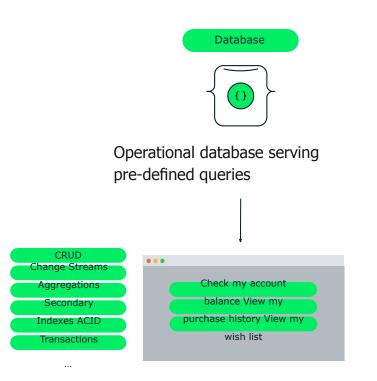


...

## How developers build app search today

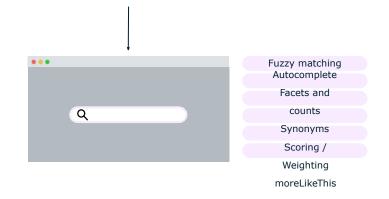


Bolt-on a search engine to power search across application data





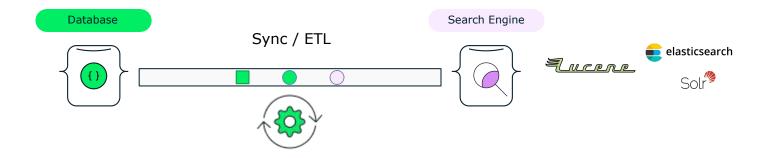
Search engine inferring intent from free form, natural language queries







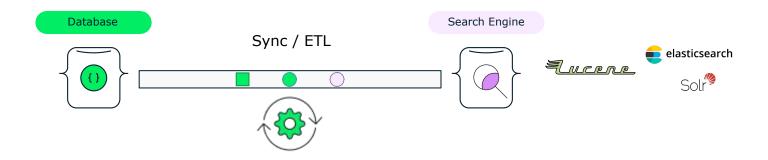
Stand up a replication mechanism to keep the systems in sync



### How developers build app search today



Architectural complexity in our application stack



- Developer friction

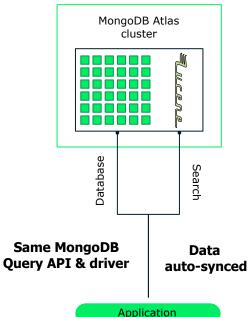
  Different query APIs and drivers for database and search, coordinate schema changes
- Pay the sync tax

  Requires its own systems and skills. Recovering sync errors can consume 10% of a developer's time
- Operational overhead More to provision, secure, upgrade, patch, back up, monitor, scale, etc.

### How does Atlas Search address that pain?







## Embeds a fully managed Apache Lucene index right alongside the database

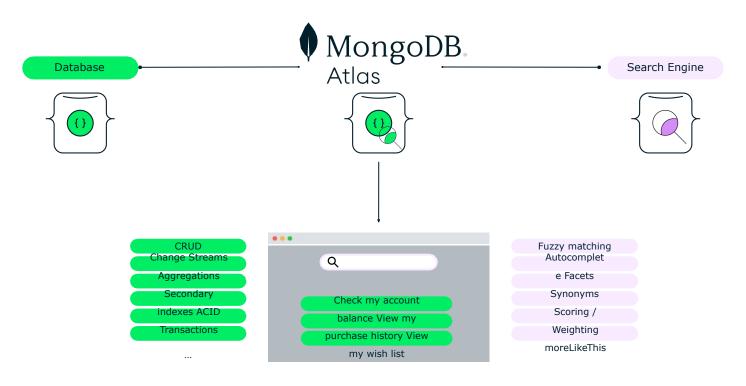
- Improved developer productivity Build database and search features using the same query API and driver
- Simplified data architecture
  Automatic data synchronization, even as your data and schema changes
- Fully managed platform

  Get the security, performance, and reliability of Atlas

## Unifying database and search

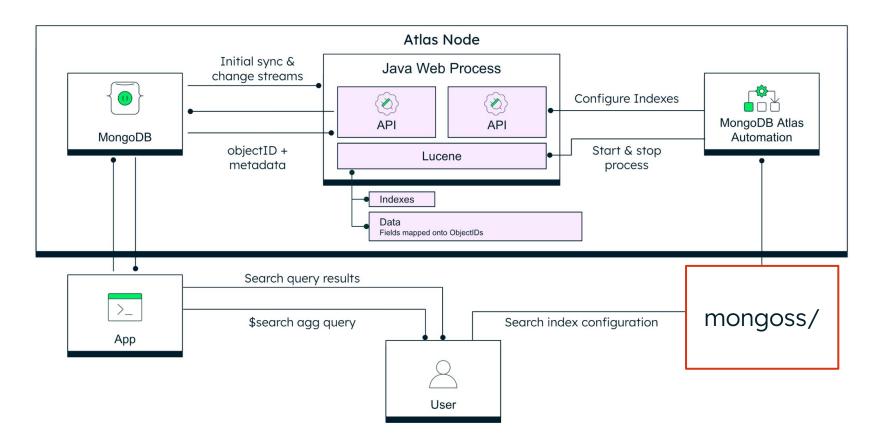
•

Compress 3 systems into 1, ship search features 30%-50% faster



### Atlas Search Architecture

















### Rich Search Features



Atlas Search offers many features to fine-tune your search results to help users find what they need, including:

- ✓ Rich query language
- Autocomplete
- ✓ Filters and facets
- ✓ Fuzzy search
- ✓ Multi-data type support
- ✓ Multi-language support (40+)
- ✓ Highlighting
- ✓ Custom scoring expressions/functions
- ✓ Synonyms
- ✓ More-like-this similarity
- ✓ Vector Search

### Analyzers: Tokenization

We start with an original query

**Analyzers** apply parsing and language rules to the query

**Tokens** are the individual terms that are queried; think of these as the words that matter in a search query

Lions and tigers and bears, oh my!

Lions and tigers and bears, oh my!

lions tigers bears oh my

## Analyzers: Stemming

Analyzers also support language rules

A search for...

...can return results for related words

"bank"

["bank", "banks", "banker",
"banking", ...]

## **Indexing Capabilities**



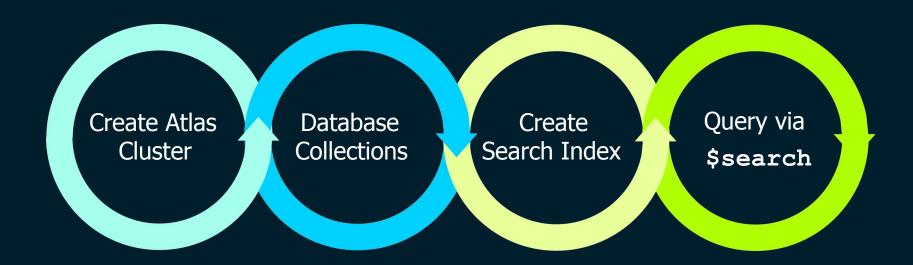
Search features are all about helping customers make their data more **discoverable** and **relevant** for end users

- Support for 40+ languages
- Multiple data types
- Custom Analyzers
- Automated ETL
- Autocomplete
- Faceting
- Stored source
- Semantic search

lucene.arabic	lucene.armenian	lucene.basque	lucene.bengali	lucene.brazilian
lucene.bulgarian	lucene.catalan	lucene.cjk	lucene.czech	lucene.danish
lucene.dutch	lucene.english	lucene.finnish	lucene.french	lucene.galician
lucene.german	lucene.greek	lucene.hindi	lucene.hungarian	lucene.indonesian
lucene.irish	lucene.italian	lucene.latvian	lucene.lithuanian	lucene.norwegian
lucene.persian	lucene.portuguese	lucene.romanian	lucene.russian	lucene.sorani
lucene.spanish	lucene.swedish	lucene.turkish	lucene.thai	

Each language analyzer has built-in rules based on the language's usage pattern

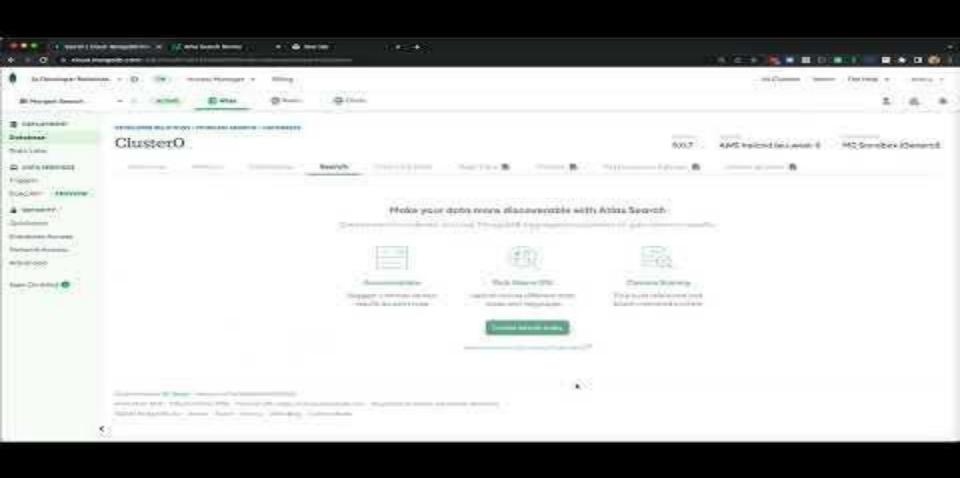




## **Search Exercise 1**

Create Indexes:

- 1. Default Search Index
- 2. Autocomplete on Title



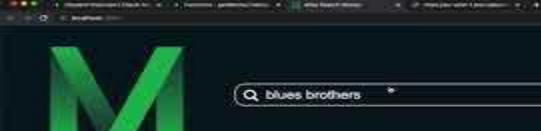


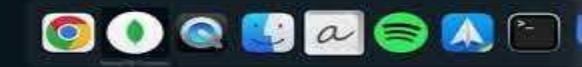
#### **MongoDB Atlas Host**

```
mongod
                                                                                                            mongot
                                                                               Wire protocol
                            Wire Protocol
                           (over internet to
                                                                                (localhost)
                           MongoDB Atlas)
          app
                                            aggregate([ {$search: {
                                                                                                    text: {
                                               text: {
                                                                                                      path: "name",
                                                 path: "name",
 db.col.aggregate([
                                                                                                      query: "star wars"
                                                 query: "star wars"
 {$search: {
    text: {
     path: "name",
                                             }}]
     query: "star wars"
                                                                                                   Lucene booleanQuery:
                                                                                                     (should
  }}])
                                             Lookup([{_id: "123"],
                                                                                                      (term("name", "star"),
                                                                                                      term("name", wars"))
                                             {...}])
[ { id: "123", title:
"Star Wars"}, {...}]
                                             [ { _id: "123",score:
                                            1.23, highlights: [...] },
                                                                                                   { "name" :
                                            {...}]
                                                                                                      { "star" : [123,124],
                                                                                                       "wars":
                                                                                                   [123,125,...]}}
```

## **Search Exercise 2**

Search Aggregations





C+ K \* B (C) | B \* B B (

## **Search Exercise 3**

Vector Search Overview

# Token / Lexical Search

#### What?

- TF-IDF/BM25(f)
- Keyword Search
- Tokenization is applied during indexing

#### When?

- Your text corpus closely matches how users search
- Most traditional text search use cases.
- First pass at text-based relevancy.



#### Vector Search

#### What?

- ANN/kNN
- Contextual Similarities
- Search supporting multiple modalities.

#### When?

- Ambiguous user input.
- There is a 'vocabulary gap' between corpus and how users search
- Non-text searches.



# Vector Search powers a number of key use cases

**Semantic Search** 

Similarity Search

Recommendation Engines

Long-term memory for LLMs

Q & A Systems

Image, Audio, Multimedia Search

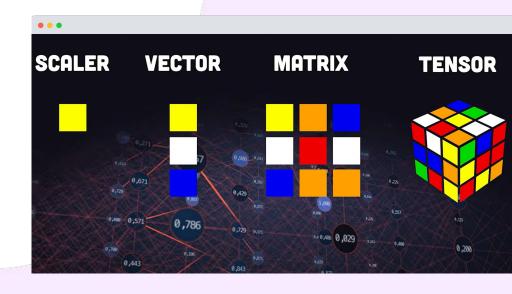


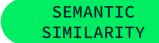
#### What is a vector?

In math terms, it is just a data structure - an array

We store a list of points that make up a representation of a word/term in multidimensional space, then compare how close other points (terms) are to calculate similarity

The devil is in the details of how we convert a word/term/phrase into this vector

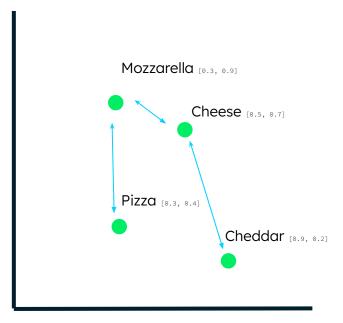






## Demonstrated via simplified vector

(in two dimensions)





#### Vector Search - Data Structure

#### Sample Document

```
"plot": "Three unemployed parapsychology professors set up shop as a unique gho..."
"title": "Ghostbusters"
"year":1984
"plot_embedding": Array
0:-0.03442629799246788
1:-0.013609294779598713
2:-0.01283580157905817
                                              embeddings generated by Atlas
3:-0.044231850653886795
                                              trigger invoking an encoding
4:-0.02906322292983532
                                              API call.
5:-0.021946318447589874
6:-0.001014339504763484
                                              More dimensions = more
7:-0.11677482724189758
                                              floating points
8:-0.03159894421696663
9:-0.05249767750501633
10:0.07232994586229324
```

. . .

#### Semantic Search - Index Definition

#### **Index Definition**

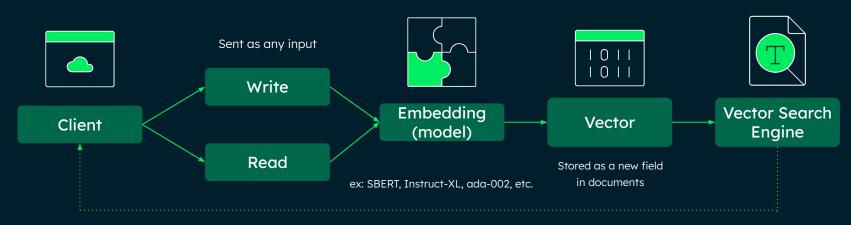


#### Sample Query

```
Γ{
    "$vectorSearch": {
      "index": "default",
      "path": "plot_embedding",
      "filter": {
        "$and": [{
          "year": {"$gt": 1955},
          "year": {"$lt": 1975},
          "Tenant_id: {"$eq": 23}
        }]
      "queryVector": [-0.0072121937,-0.030757688,0.014948666,
                      0.015741806, -0.020635074, -0.012945653,
      "numCandidates": 150,
      "limit": 10
 }, {"$project": {...}},...
```

### Vector Search Data Workflow

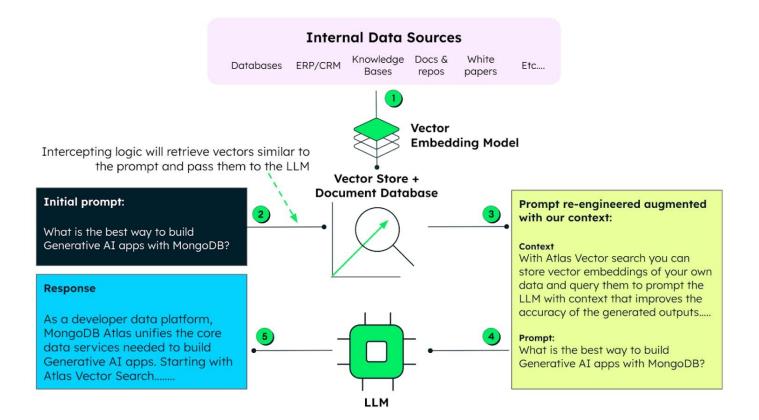




Documents returned as they're stored, typically remove vectors via \$project.



#### Vector Search enhanced Chatbots



## Search Exercise 4 (Optional)

## Self Serve Demo:

atlassearchrestaurants.com github.com/mongodb-developer/WhatsCooking