FACTSET

March 2 - 3, 2022 – Subsurface LIVE Winter 2022

# Powering a Data Mesh with Dynamic GraphQL Schema Generation

Matthew Topol, Principal Software Architect

Steve Perkins, Lead Software Engineer

Content Engineering

## Agenda

- Why would you need this? (What were we trying to solve?)

- What about Federation?

- The Tech We're Using

- Flexibility Through Metadata

- Benefits of Dynamic Schema Generation

- Future Optimizations and Tooling

- Q&A

# Why would you need this?

## Consistent New Data Points
The data points that product developers want to expose is continuously expanding. Onboarding new data points needs to be quick and painless.

## Multiple Different Data Stores
The data that we need to expose is stored in multiple different technologies (postgres, MSSQL, parquet) and can not currently be migrated to something consistent.

## Multiple Internal Consumers
Not only do different internal consumers exist for this data, but their query patterns also vary. Fetching the data needs to be straightforward and not require large changes on the consumer side. Data also needs to be consistent across all consumers

## Non-Standard Schemas
The schemas for the various data stores are non-consistent. To ease data fetches for consumers, these should be standardized.

## Additional Considerations

- Abstraction from data source
  - Hide the data mesh complexity
- Data consistency and timeliness
- Cross dataset queries
- Performance
  - Query performance
  - Data transfer performance

# Why didn't we create a Federated GraphQL Service? (It's all the rage right now)

- A Federated Schema != A Standard Schema
- Maintaining multiple different services amongst various teams proved difficult over time
- Ensuring conformance updates were implemented across multiple services was not reliable
- Single-entry-point infrastructure allows us to streamline automation, improve development time, and quality assurance
  - Easier to enforce conformance to standards
- Push joins down to query engine
  - More performant

# Primary Tech We're Using

## Dremio

Our source agnostic query engine. Dremio allows us to pull in data from multiple different data stores and fetch data in a performant manner
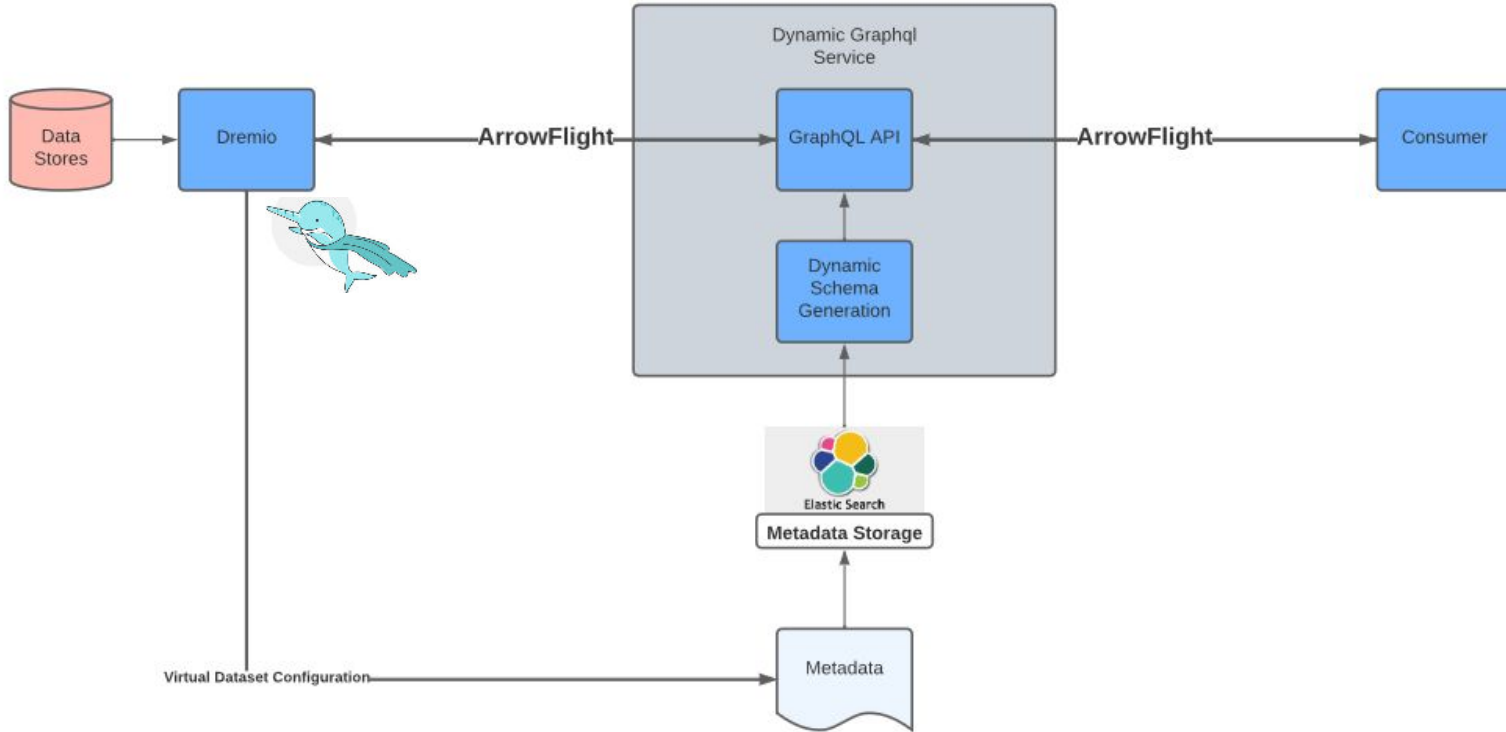
## GraphQL

Our flexible schema interface. GraphQL allows us to have a deterministic schema that makes fetching data straightforward and extendable.

## Apache Arrow

Our communication format. In conjunction with Arrow Flight RPC, this communication protocol bolsters round trip performance with zero serialization/deserialization

# High Level Architecture

# Additional Technologies We're Using

- Programming language
  – Golang

- Communication framework
  – Arrow Flight RPC over gRPC

- Cluster Deployment
  – Kubernetes

- Golang Libraries
  – graphql-go-tools - github.com/jensneuse/graphql-go-tools
  – Apache Arrow - github.com/apache/arrow/go/v7/arrow

# Flexibility through Metadata

- Currently our metadata is driven by a fairly simple YAML configuration
  - Split into Entities and Data Items

- Multiple levels of automated schema validation to ensure expectations are met
  - Metadata schema defined using JSON Schema which can be applied to YAML

- Metadata can be owned and manipulated by content teams to provide a self-service API
  - Faster time-to-market for exposing new data!

- Schema generation is flexible enough to ignore unrecognized attributes until they are leveraged

- Deterministic schema generation leads to easily deterministic Query Generation!

Metadata Drives All

# Metadata Example

### Entity

```
entity: Bank
keyField: bank_entity_id
displayTable: bus.banks.banks_entity_display
filterTable: bus.banks.banks_entity_display
caccess: 00FFN,00FFG,OR
periodTypes:
  - name: Annual
    value: A
  - name: Quarterly
    value: Q
  - name: Semi-Annual
    value: S
pivot:
  Deal: bank_entity_id
rel:
  hasMany:
    - Deal
```

### Data Item

```
caccess: 0FFN,00FFG,OR
dataType: FLOAT
entity: Bank
id: venture_capital_revenue
name: Venture Capital Revenue
order: 11
periodData:
    keyField: entity_id
    table: bus.banks.bank_financials_loans_leases
    types:
        - ANN
        - QUARTER
        - SEMI
preferredDecimal: 2
query:
    queryType: FLOAT
result:
    sortable: true
scaling: 1e+06
shortDescription: Venture Capital Revenue
subGroup: Non-Interest Income
type:
    name: item
    parent: group.regulatory_income
unitLabel: USD
```

# Meshing Around With Different Schemas

# Produced GraphQL Schema:



**FACTSET**

Search the docs ...

QUERIES

bankCounts(...): bankCounts
banks(...): [Bank]
bankAggs(...): [BankAggs]
toptenants(...): [Toptenant]
toptenantAggs(...): [ToptenantAggs]
dealCounts(...): dealCounts
deals(...): [Deal]
dealAggs(...): [DealAggs]
buyers(...): [Buyer]

**banks(**
    **limit:** Int
    **offset:** Int
    **where:** BankFilterInput
    **sort_order:** [BankOrderList]
**):** [Bank]

TYPE DETAILS

type Bank {
    state_name: String
    status_code: String
    num_full_time_emp(...): Float
    bank_name: String
    bank_type_description: String
    ticker: String
    bank_sub_type_description: String
    bank_entity_id: String
    ff_peer_group_detail: String
    ff_int_exp_tot(...): Float

**where:** BankFilterInput

TYPE DETAILS

type BankFilterInput {
    bank: BankInput
    AND: [BankFilterInput]
    OR: [BankFilterInput]
    NOT: BankFilterInput
    deal: DealInput
}

**deal:** DealInput

TYPE DETAILS

type DealInput {
    category: IntFilterInput
    deal_type_identifier: StringFilterInput
    deal_sub_type_identifier: StringFilterInput
    business_desc: StringFilterInput
    factset_sector: IntFilterInput
    factset_industry: IntFilterInput
    region_code: StringFilterInput
    joined_country_code: StringFilterInput
    joined_state_code: StringFilterInput
    transaction_value: FloatFilterInput
    status: StringFilterInput

**open_date:** DateFilterInput
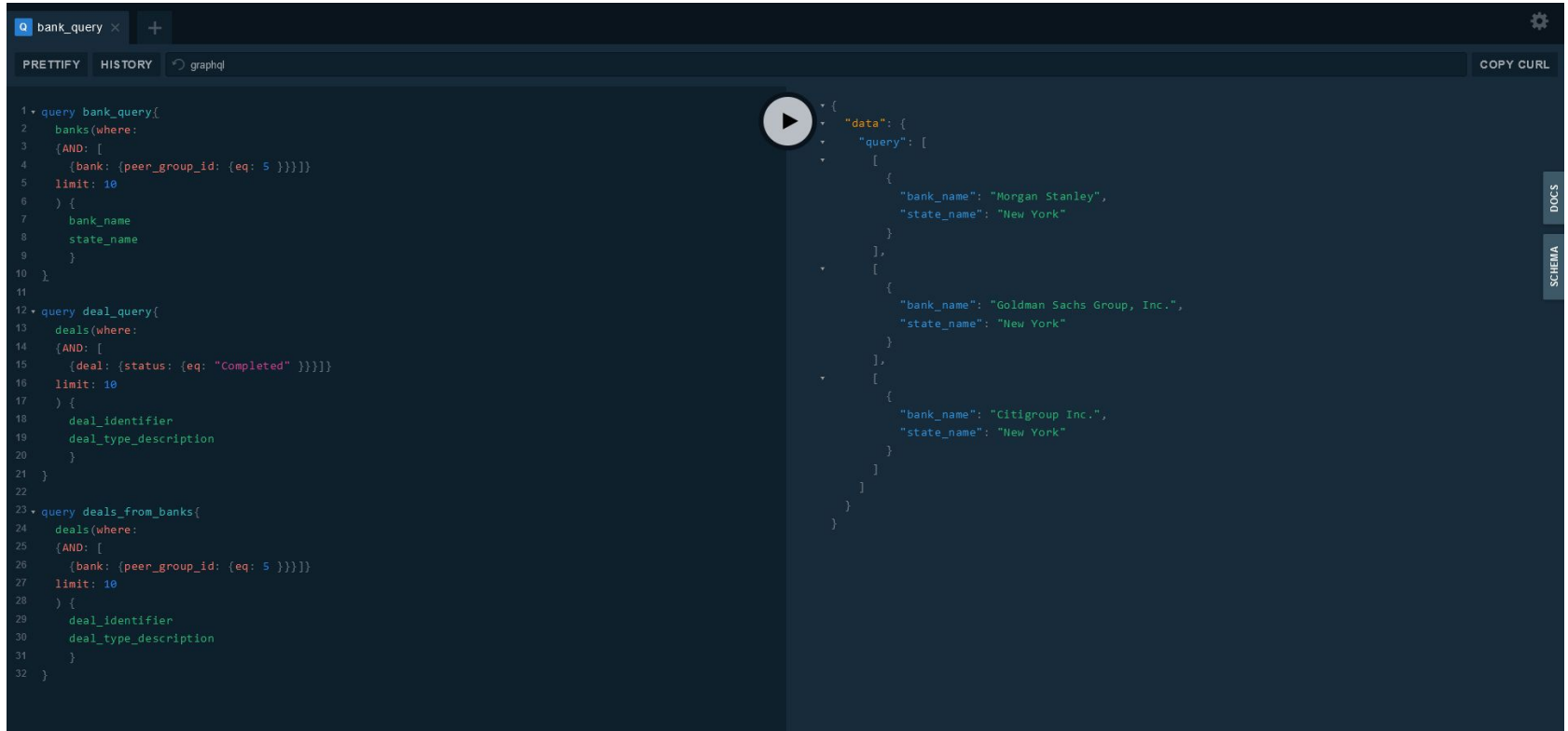
Filter By Announcement Date

TYPE DETAILS

type DateFilterInput {
    ne: Date
    eq: Date
    le: Date
    lt: Date
    ge: Date
    gt: Date
    between: [Date]
    oneOf: [Date]
}

# One Query Pattern To Rule Them All

# Benefits of Dynamic Schema Generation

- Faster time to market for new data points

- Ease of automated extension
  - Define relationships in metadata for easy one-to-many / many-to-many query generation

- Data source abstraction
  - Consumers write GraphQL, not SQL

- Reduced engineering effort necessary for maintenance

- Consistent query patterns allow for templated query generation

# The Query Template

```sql
SELECT <display_columns>
FROM <display_table> AS t1

LEFT JOIN (

  SELECT <dimension_1_display_columns>

  FROM <dimension_1_display_table>

  WHERE <dimension_1_display_filters>) AS t2 ON <dimension_1_join_fields>

LEFT JOIN ( ... )AS t2 ON <dimension_n_join_fields>

WHERE <display_key> IN (

  SELECT

    DISTINCT <display_key>

  FROM <filter_table> AS t1

  WHERE

    <filter_criteria> AND

    <display_key> IN (

        SELECT DISTINCT <display_key> FROM <dimension_table_1> AS t1 WHERE <dimension_1_filter_criteria>) AND

    <display_key> IN ( ... ))

  ORDER BY LOWER(<display_key>) ASC

  LIMIT 1000
```

# The Query

```sql
SELECT "bank_entity_id", "ticker", "bank_name", "bank_sub_type_description",…
FROM bus.banks.banks_entity_display AS t1
LEFT JOIN (…) AS t2 ON (t1."bank_entity_id" = t2."entity_id")
LEFT JOIN (…) AS t3 ON (t1."bank_entity_id" = t3."entity_id")
WHERE t1."bank_entity_id" IN (
  SELECT
    DISTINCT t1."bank_entity_id"
  FROM bus.banks.banks_entity_display AS t1
  WHERE
    ((t1."bank_sub_type" = '0|1' OR t1."bank_type" = '0' OR t1."bank_type" = '1') AND t1."bank_entity_id" IN (
    SELECT DISTINCT "entity_id" FROM bus.banks.bank_financials_loans_leases AS t1 WHERE
      (t1."venture_capital_revenue" >= 10 AND period_type = 'A' AND fiscal_period = 20211231)
    )
   AND t1."status_code" IN ('Active','Extinct'))) ORDER BY LOWER("bank_entity_id") ASC LIMIT 1000
```

# FACTSET

## Dynamic Schemas -> Deterministic Queries

## Future Optimizations and Tooling

- Improve self-service tooling for adding new metadata
  - Web Platform for internal content teams to use

- More automated metadata generation from datasets

- Higher level GraphQL query planning for utilizing column-caches when available for super-low-latency queries
  - Provide query handling for subset of patterns which can't be served with Dremio SQL

- Computational capabilities in addition to fetches worked into dynamic schema and query generation

Questions?

FACTSET