# Data Access for Data Science

April 17, 2018

# Jacques Nadeau

Co-Founder & CTO, Dremio
PMC Chair, Apache Arrow
PMC, Apache Calcite

# Agenda

- Apache Arrow
- Using Dremio for Self Service Data Access
- Data Access Example (notebook + Dremio)
- Reflections & Caching Overview
- Caching Impact Example

dremio

# Getting Data Ready for Analysis Is Hard

- Data can be hard to find
- Many modern data systems have poor quality interfaces
- Data is rarely in a single system
- Data access is frequently slow
- Some types of issues can only be solved by IT tickets
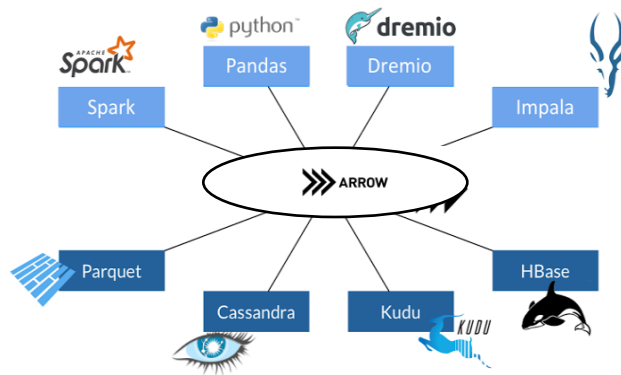- Doing late stage data curation makes reproduction and collaboration difficult: "do I copy and edit?"

*there should be a new, self-service data access tier*

**dremio**

# Apache Arrow

# Apache Arrow

- Standard for columnar in-memory processing and transport
- Focused on Columnar In-Memory Analytics
  1. 10-100x speedup on many workloads
  2. Common data layer enables companies to choose best of breed systems
  3. Designed to work with any programming language
  4. Support for both relational and complex data
- Consensus Driven: developed by contributors leading 13+ key OSS projects

**dremio**

# Arrow: Fast Exchange, Fast Processing



## High Performance Sharing & Interchange
- Zero Overhead Encoding
- Scatter/Gather Optimized
- Direct Memory definition
- Designed for RDMA and shared memory access

Focus on GPU and CPU Efficiency
- Cache locality
- Super-scalar and vectorized operation
- Minimal structure overhead
- Constant value access

| | session_id | timestamp | source_ip |
|---|---|---|---|
| Row 1 | 1331246660 | 3/8/2012 2:44PM | 99.155.155.225 |
| Row 2 | 1331246351 | 3/8/2012 2:38PM | 65.87.165.114 |
| Row 3 | 1331244570 | 3/8/2012 2:09PM | 71.10.106.181 |
| Row 4 | 1331261196 | 3/8/2012 6:46PM | 76.102.156.138 |

### Traditional Memory

| | |
|---|---|
| Row 1 | 1331246660 |
| | 3/8/2012 2:44PM |
| | 99.155.155.225 |
| Row 2 | 1331246351 |
| | 3/8/2012 2:38PM |
| | 65.87.165.114 |
| Row 3 | 1331244570 |
| | 3/8/2012 2:09PM |
| | 71.10.106.181 |
| Row 4 | 1331261196 |
| | 3/8/2012 6:46PM |
| | 76.102.156.138 |

### Arrow Memory

| | |
|---|---|
| session_id | 1331246660 |
| | 1331246351 |
| | 1331244570 |
| | 1331261196 |
| timestamp | 3/8/2012 2:44PM |
| | 3/8/2012 2:38PM |
| | 3/8/2012 2:09PM |
| | 3/8/2012 6:46PM |
| source_ip | 99.155.155.225 |
| | 65.87.165.114 |
| | 71.10.106.181 |
| | 76.102.156.138 |

# Arrow Components

- Core Libraries
- Building Blocks
- Major Integrations

dremio

# Arrow: Core Libraries

- Java Library
- C++ Library
- Python Library
- C Library
- Ruby Library
- JavaScript Library
- Rust Library

**dremio**

# Arrow Building Blocks (in project)

## Plasma

Shared memory caching layer, originally created in Ray

## Feather

Fast ephemeral format for movement of data between R/Python

## Arrow RPC*

RPC/IPC interchange library (active development)

## Arrow Kernels*

Common data manipulation components

*soon

# Arrow Integrations

## Pandas

Move seamlessly to from Arrow as a means for communication, serialization, fast processing

## Spark

Supports conversion to Pandas via Arrow construction using Arrow Java Library

## Dremio

OSS project, Sabot Engine executes entirely on Arrow memory

## Parquet

Read and write Parquet quickly to/from Parquet. C++ library builds directly on Arrow.
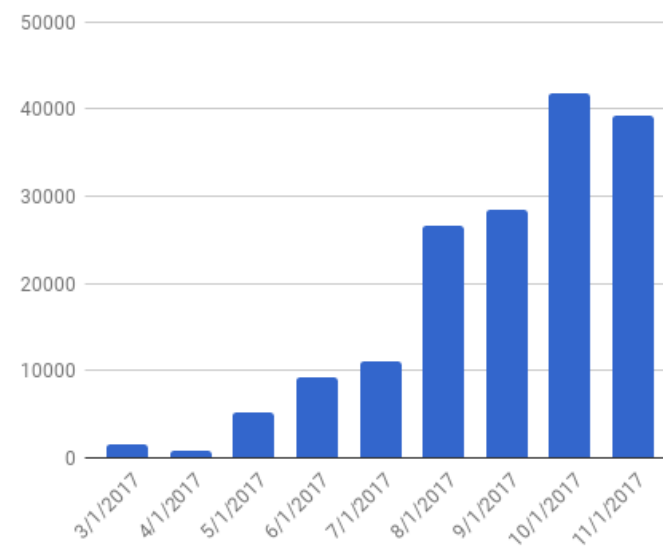
## GOAI (GPU Open Analytics Init)

Leverages Arrow as internal representation (including libgfd and GPU dataframe)

# Apache Arrow Adoption



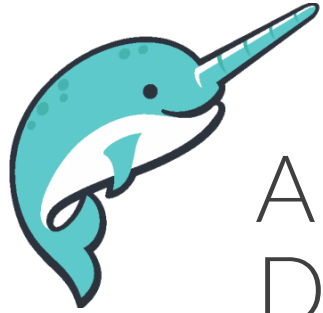Arrow downloads increased 44x since April (currently ~100K per month)
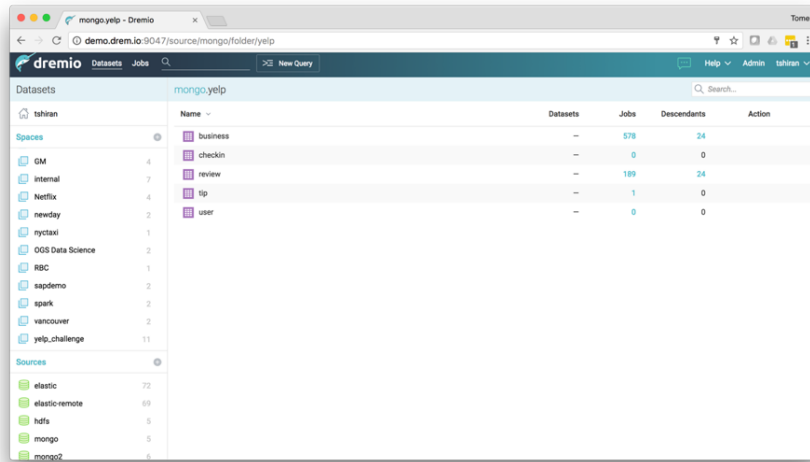
Monthly PyPi (~40% of all downloads)

dremio

# Dremio
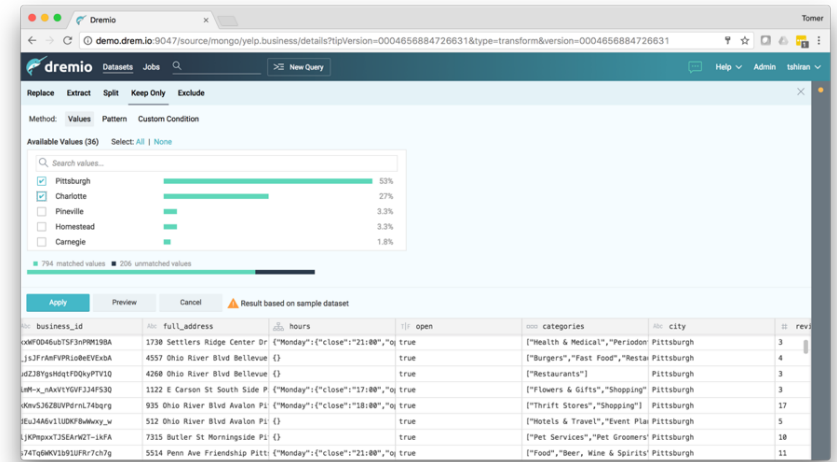*a system for self-service data access*

# About Dremio

- Launched in July 2017
- Self-Service Data Platform
- Make Data Accessible to whatever tool
- The Narwhal's name is **Gnarly**

- Apache-Licensed
- Built on Apache Arrow, Apache Calcite, Apache Parquet
- Easy extension, customization and enterprise flexibility
- SDKs for sources, functions, file formats, security
- Execution, Input and Output are all build on native **Arrow**
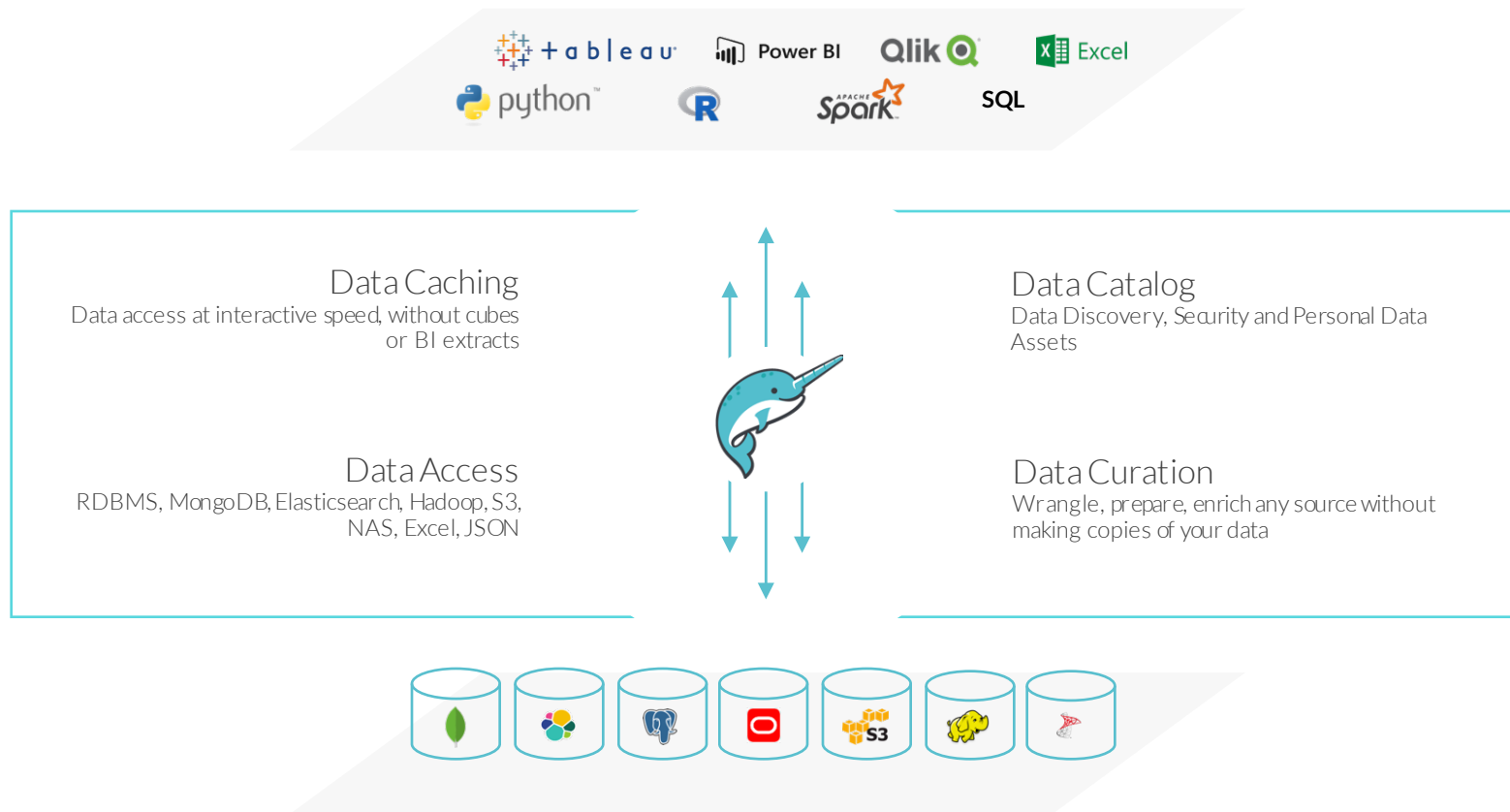
# Google Docs for your Data



## Powerful & Intuitive UX for Data
Find, manage and share data regardless of size & location



## Live Data Curation
AI-powered curation of data without creating a single copy

dremio

# Self-Service Data Access Platform

tableau    Power BI    Qlik Q    Excel

python    R    Spark    SQL

### Data Caching
Data access at interactive speed, without cubes or BI extracts

### Data Access
RDBMS, MongoDB, Elasticsearch, Hadoop, S3, NAS, Excel, JSON

### Data Catalog
Data Discovery, Security and Personal Data Assets

### Data Curation
Wrangle, prepare, enrich any source without making copies of your data
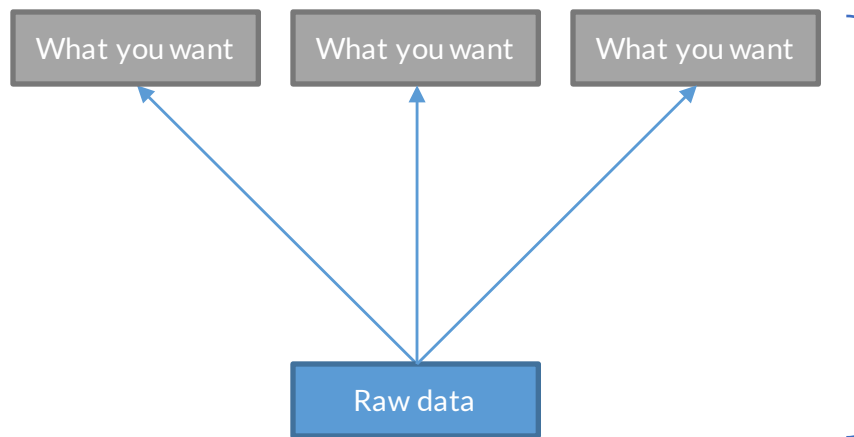
**dremio**

# Data Access Example

# Leveraging Underlying Source Capabilities Example

# Reflections
*an advanced form of caching*

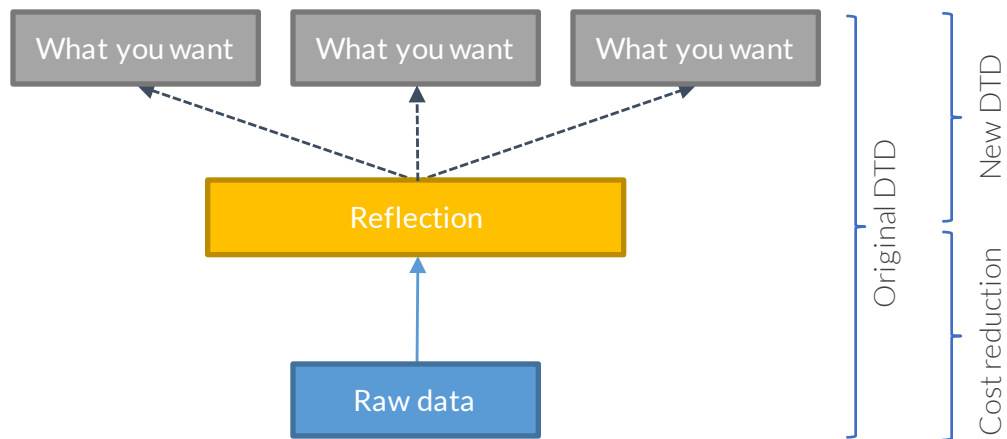# Access isn't Enough: Reducing Distance to Data

| What you want | What you want | What you want |

Raw data

Distance to Data
- Work to Be Done
- Resources Required
- Time to Complete

dremio

# The basic concept behind a relational cache

| | |
|---|---|
| What you want | What you want | What you want |

**Reflection**

**Raw data**

Original DTD

New DTD

Cost reduction

- Maintain derived data that is between what you want and what the raw data

- Shortens distance to data (DTD)

- Reduces resource requirements & latency

- Materialization can be derived from raw data via arbitrary operator DAG

# It doesn't have to be a trivial relationship…

| | | |
|---|---|---|
| What you want | What you want | What you want |

| | |
|---|---|
| Reflection 1 | Reflection 2 |

| |
|---|
| Raw data |

New DTD

Original DTD

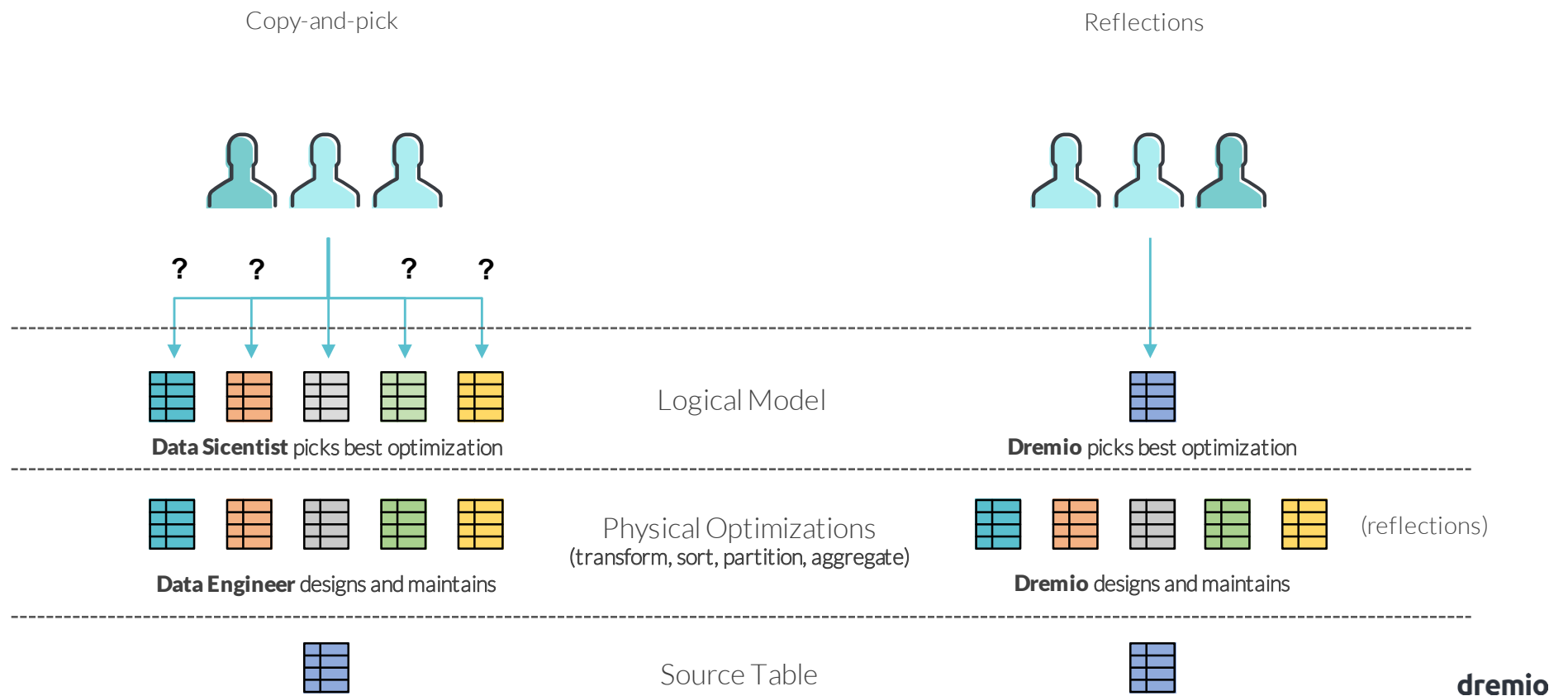Cost reduction

# You already do this today (manually)!!

**Materializations (manually created):**

- Cleansed

- Partitioned by region or time

- Summarized for a particular purpose
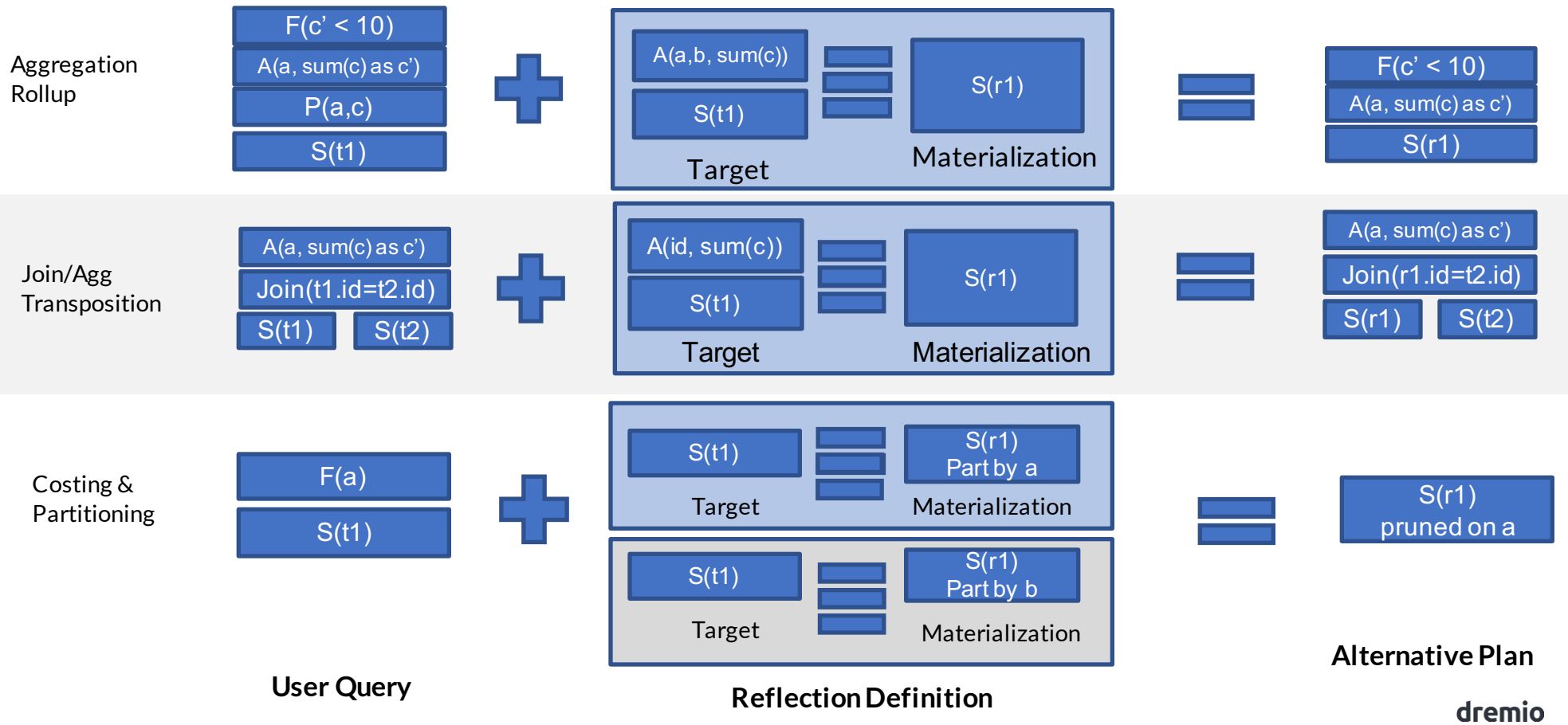
**Users choose depending on need:**

- Data Scientists & Analysts trained to use different tables depending on the use case

- Custom datasets, summarization and/or extraction for modeling, reports and dashboards

# Dremio can make the decisions so you don't have to

Copy-and-pick

Reflections

? ? ? ?

Logical Model

**Data Sicentist** picks best optimization

**Dremio** picks best optimization

Physical Optimizations
(transform, sort, partition, aggregate)

(reflections)

**Data Engineer** designs and maintains

**Dremio** designs and maintains

Source Table

dremio

# Cache Matching: Example Scenarios

**Aggregation Rollup**

| F(c' < 10) |
| A(a, sum(c) as c') |
| P(a,c) |
| S(t1) |

**+**

Target:
| A(a,b, sum(c)) |
| S(t1) |

Materialization:
| S(r1) |

**=**

| F(c' < 10) |
| A(a, sum(c) as c') |
| S(r1) |

---

**Join/Agg Transposition**

| A(a, sum(c) as c') |
| Join(t1.id=t2.id) |
| S(t1) | S(t2) |

**+**

Target:
| A(id, sum(c)) |
| S(t1) |

Materialization:
| S(r1) |

**=**

| A(a, sum(c) as c') |
| Join(r1.id=t2.id) |
| S(r1) | S(t2) |

---

**Costing & Partitioning**

| F(a) |
| S(t1) |

**+**

Target:
| S(t1) |

Materialization:
| S(r1) Part by a |

Target:
| S(t1) |

Materialization:
| S(r1) Part by b |

**=**

Alternative Plan
| S(r1) pruned on a |

---

**User Query**

**Reflection Definition**

dremio

# Reflections

- A reflection is a materialization designed to accelerate operations
- Transparent to data consumers
- Not required on day 1... you can add reflections at any time
- One reflection can help accelerate queries on thousands of different virtual datasets (logical definitions)
- Reflections are persisted (S3, HDFS, local disks, etc.) so there's no memory overhead
- Columnar on disk (Parquet) and Columnar in memory (Arrow)
- Elastic, scales to 1000+ nodes

**dremio**

Reflection Impact Example
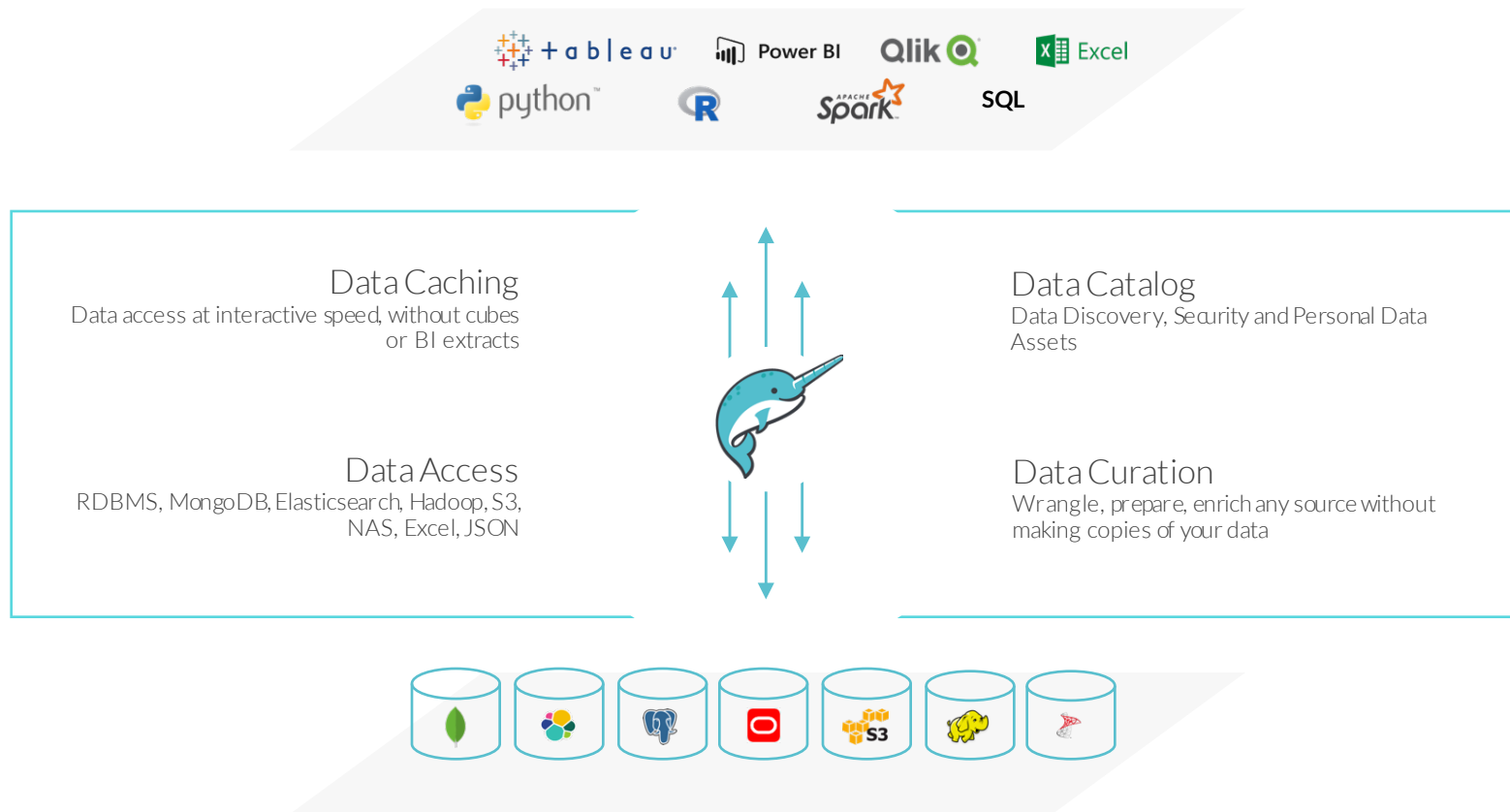
# Distribution of Responsibilities

## Data Access Platform

- Index, secure, expose, share and curate datasets
- Expose data from different systems in a standard namespace and
- Allow live cleanup and curation capabilities
- Data manipulation that should be reproducible and shared
- Disconnect physical concerns from logical needs
- Cache intermediate results to support accelerate common user patterns
- Get to an **interesting slice** of data

## BYO Data Science & BI Solutions

- Analyze Data
- Experiment and perform what-if analysis
- Derive Conclusions
- Build Models
- … and everything else that results in an output that isn't a dataset

# Self-Service Data Access

# Join the Community!

- Come see me for Office hours!
- Download: dremio.com/download
- GitHub: github.com/dremio/dremio-oss
- github.com/apache/arrow

- Dremio Community: community.dremio.com
- Arrow Mailing list: dev@arrow.apache.org

- Twitter: @intjesus, @DremioHQ, @ApacheArrow

**dremio**

dremio