# Data Algebra



*The Algebra of Data*
*A Foundation for the Data Economy*

*by Gary J. Sherman, PhD*
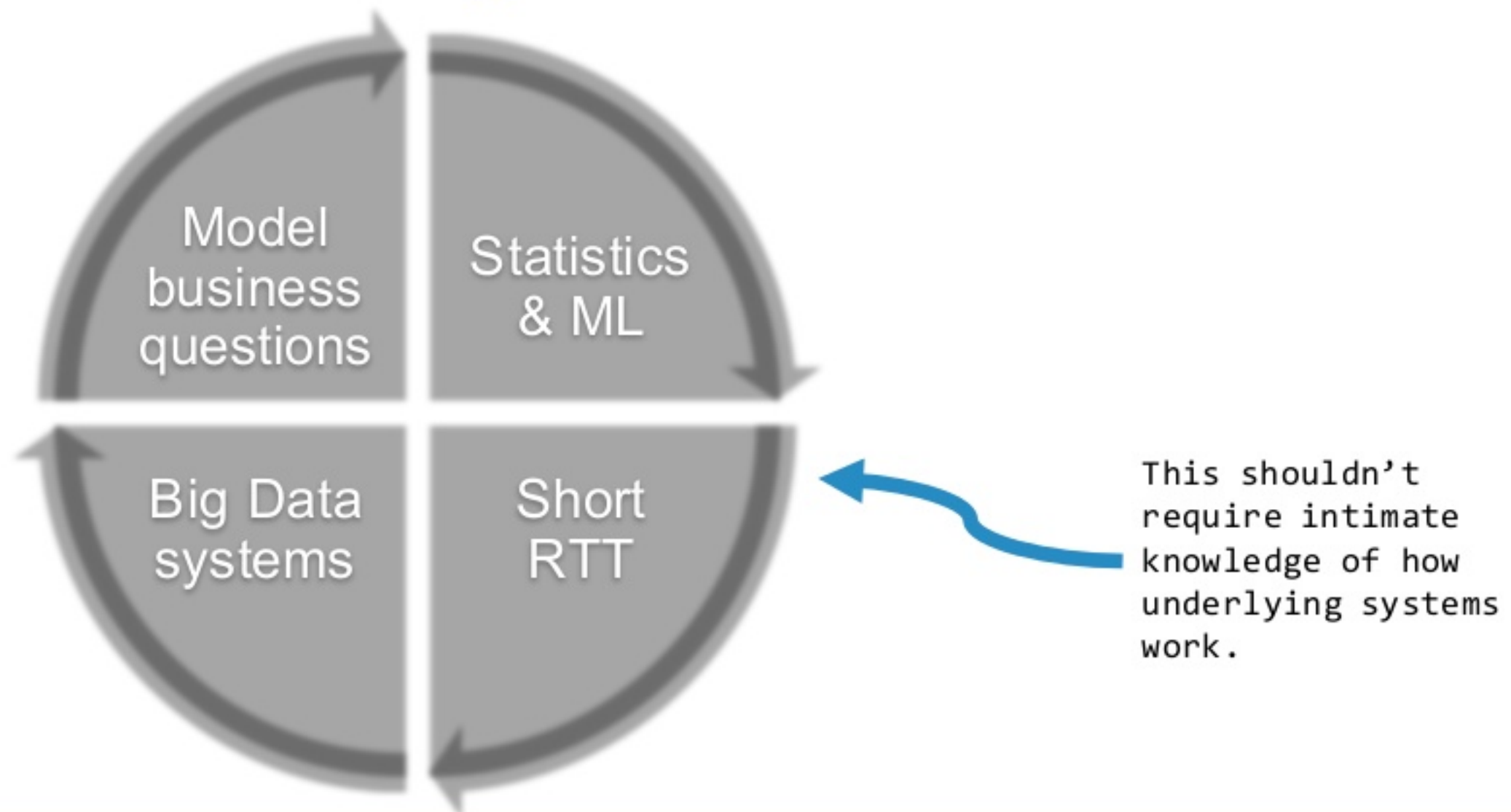*and Robin Bloor, PhD*

- A **novel application of Set Theory** to Data Processing

- **Applicable to many data models** including SQL

## Just-in-Time Analytics

needs…

## Autonomous Data Management

# JIT Analytics and the Life of a Modern Analyst

Model business questions

Statistics & ML

Big Data systems

Short RTT

This shouldn't require intimate knowledge of how underlying systems work.

# Spark for JIT Analytics: *The Good*

- Unified API

- Schema-on-read and Heterogeneous Data Sources

- Declarative Languages/APIs and Catalyst

- Elastic Compute

# Spark for JIT Analytics: *The Bad*

- Challenges for interactivity, efficiency, and scalability

- Cost of creating and maintaining "glue code"

- ***Data scientists and engineers are doing DBA work***

# Database Management Responsibilities

| Capacity planning | Configuration |
|---|---|
| Performance tuning | A billion other things |

We will focus on the performance and tuning aspects

# Improving and Maintaining Performance

- Indexes
- Materialize views
- Pre-aggregate data
- *Lots* of configuration

# Performance Tuning Strategies in Spark

- Segment, cache, and checkpoint
- Configure cluster parameters
- `spark.sql.shuffle.partitions`

# What is the Problem with Manual Tuning?

- Varies with the data (skew and scale), queries, and hardware
- Often done through trial and error
- Problems are exacerbated with JIT analytics case
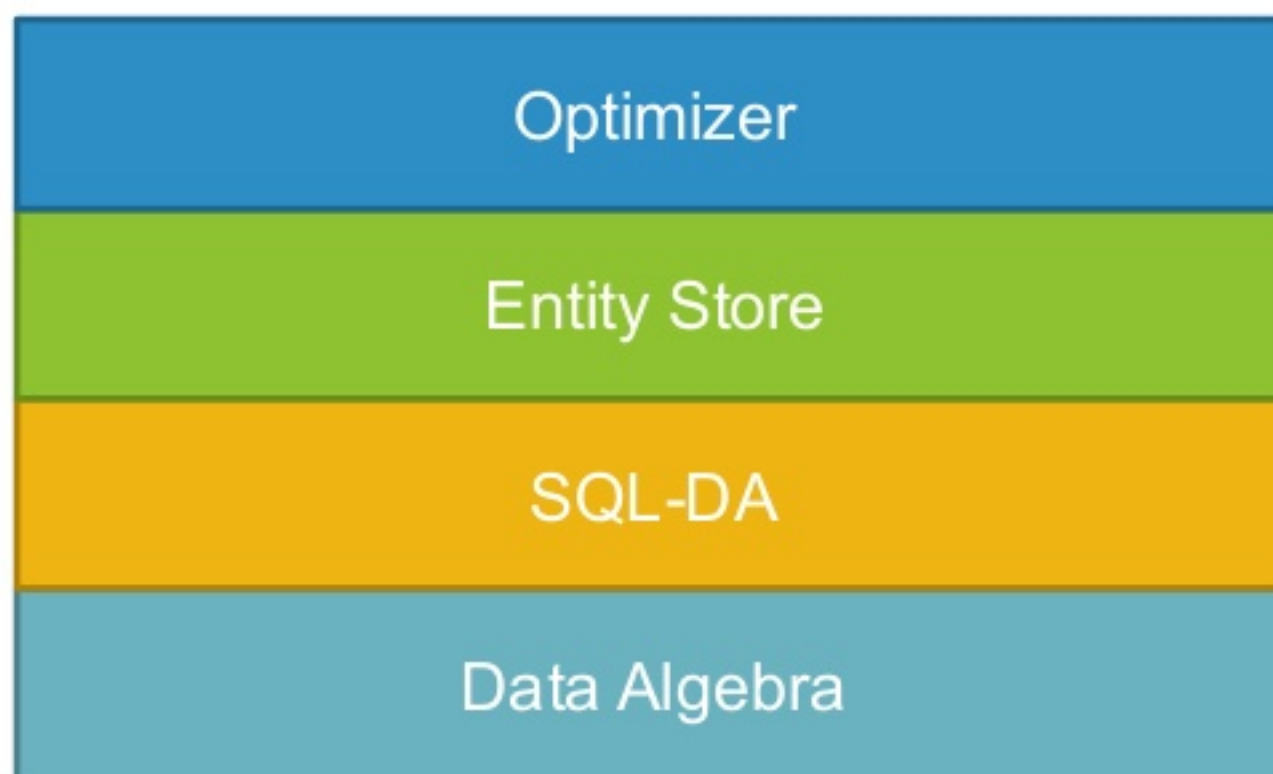- Shared resources

# What is the Problem with Manual Tuning?

It is hard and time-consuming.

# A Motivating Example for
# Autonomous Data Management

| Query 1 | Query 2 | Query 3 |
|---|---|---|
| ```sql
SELECT *
FROM A JOIN B ON A.a = B.b
WHERE B.b2 < 100
``` | ```sql
SELECT A.a, A.foo
FROM A JOIN B ON A.a = B.b
JOIN C ON B.c = C.c
WHERE B.b2 < 100 AND C.bar = "bar"
``` | ```sql
SELECT (A.foo + D.baz) AS foo_or_fu_baz
FROM A JOIN D ON A.d = D.d
JOIN B ON A.a = B.b
WHERE (D.baz LIKE "baz%" OR A.foo in ("foo", "fu"))
AND B.b2 < 100
``` |

| Query 1 w/ FPP | Query 2 w/ FPP | Query 3 w/ FPP |
|---|---|---|
| ```sql
1  SELECT *
2  FROM (
3      SELECT *
4      FROM A JOIN B ON A.a = B.b
5      WHERE B.b2 < 100
6  ) as ChunkExample
``` | ```sql
1  SELECT A.a, A.foo
2  FROM (
3      SELECT *
4      FROM A JOIN B ON A.a = B.b
5      WHERE B.b2 < 100
6  ) as ChunkExample
7  JOIN C ON B.c = C.c
8  WHERE C.bar = "bar"
``` | ```sql
1  SELECT (A.foo + D.baz) AS foo_or_fu_baz
2  FROM #ChunkExample
3  JOIN D ON A.d = D.d
4  WHERE (D.baz LIKE "baz%" OR A.foo in ("foo", "fu"))
``` |
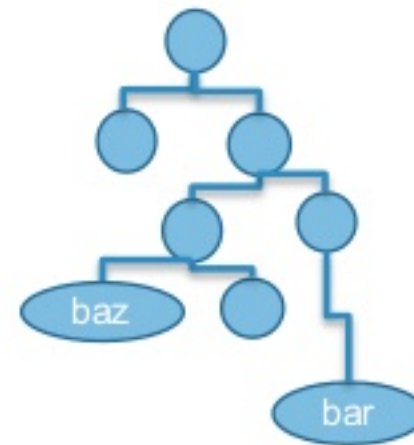
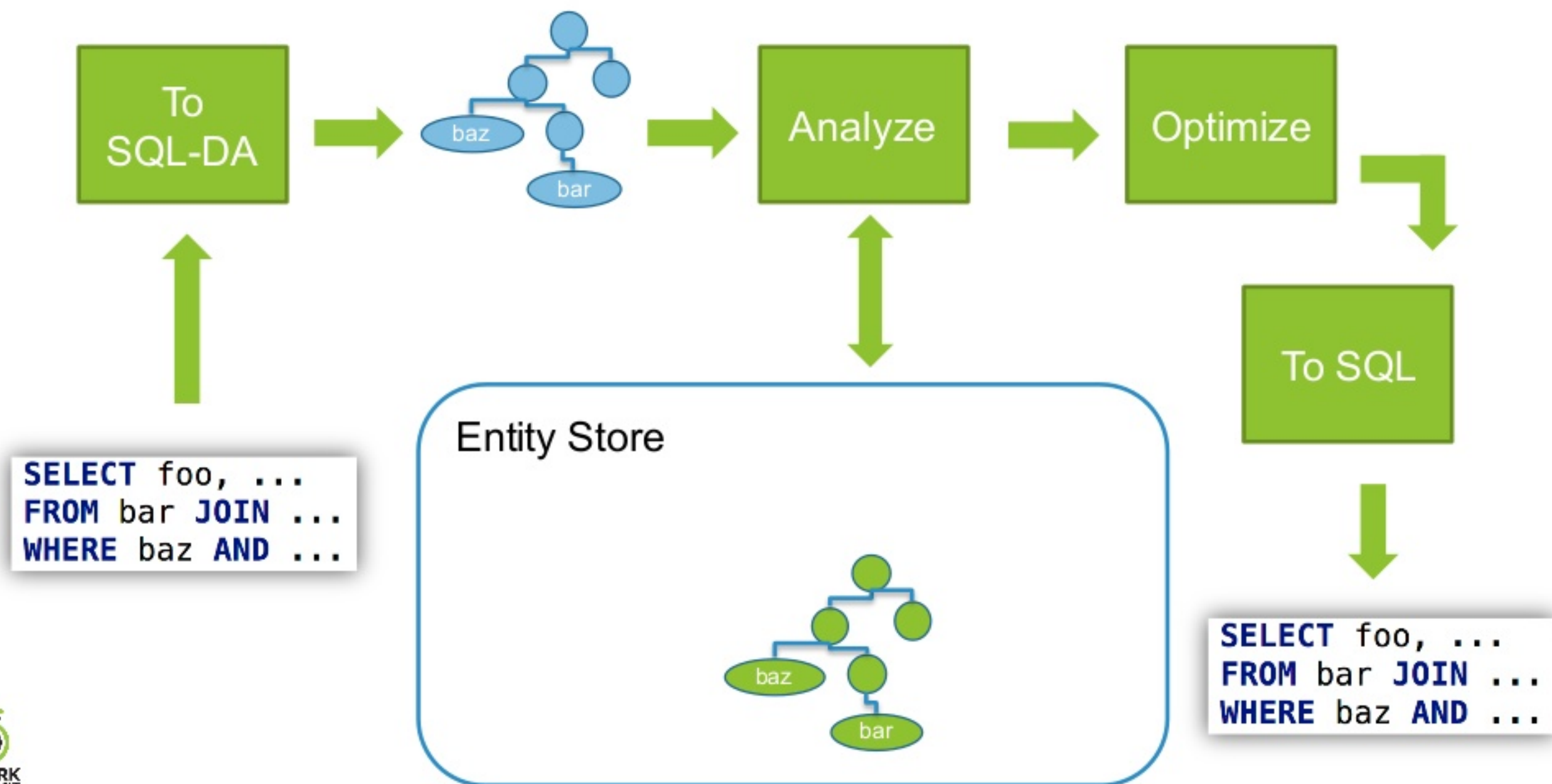# A Motivating Example for Autonomous Data Management

| |
|---|
| Optimizer |
| Entity Store |
| SQL-DA |
| Data Algebra |

```
SELECT *
FROM A JOIN B ON A.a = B.b
WHERE B.b2 < 100
```
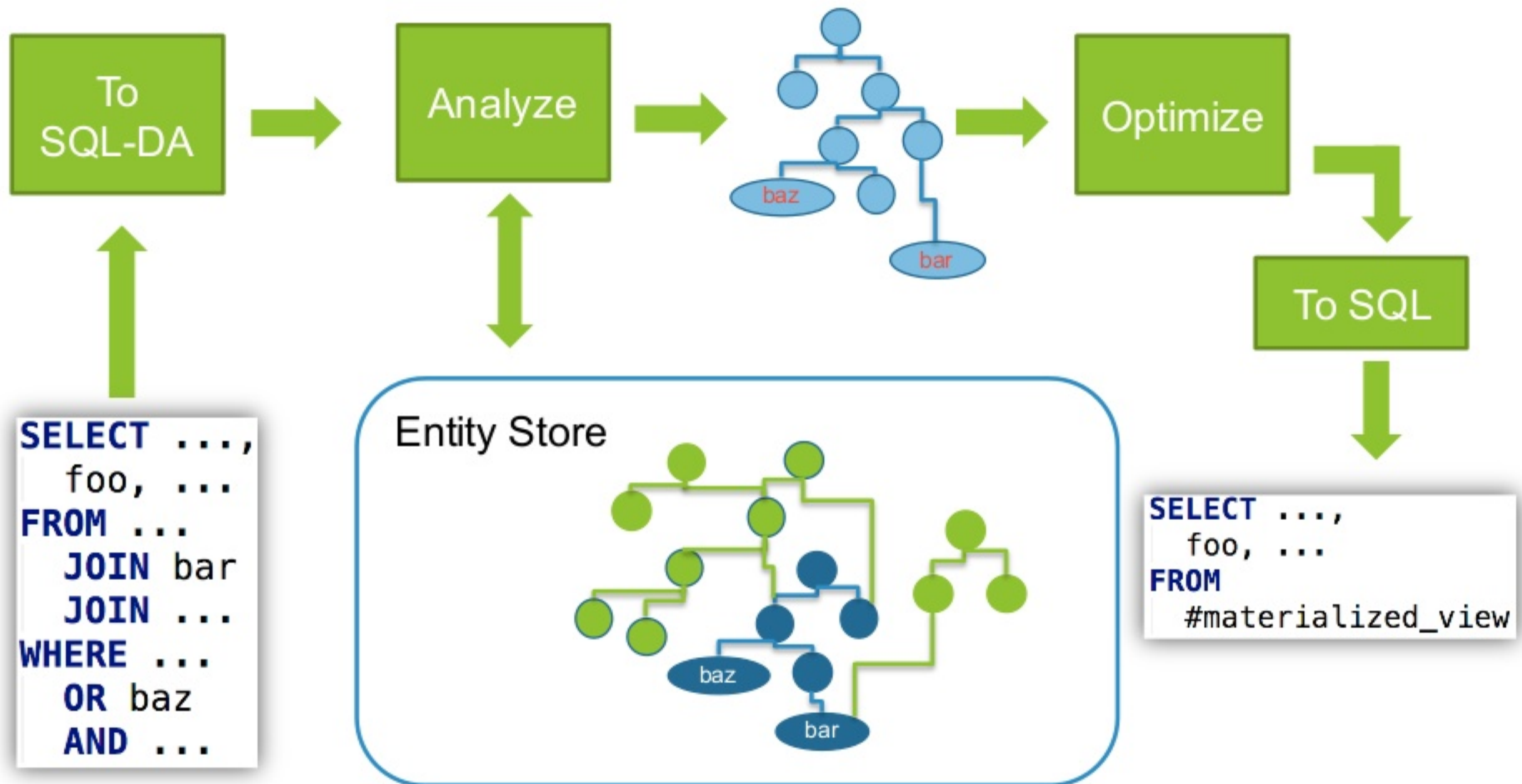
$$A\nabla B = \left( \left\{ \begin{array}{c} \{0 \mapsto \alpha, \dots, 3 \mapsto 42.0\}{:}3, \\ \{0 \mapsto \beta\}{:}1, \\ \dots \end{array} \right\}, \begin{pmatrix} 0 & \dots & 3 & 4 \\ a & \dots & b & b2 \\ int & \dots & float & int \\ A & A & B & B \end{pmatrix} \right)$$

$$Q := filter_{[4]<100}(filter_{[0]<[3]}(A\nabla B))$$

# Complex Query Expressions are Turned Into Look-ups

# Benefits of Autonomous Data Management

- Reduce query time
- Reduce computation resources required
- ***Allow the analyst to focus on problem solving, not data management***

# Algebraix Inside:
## An Implementation of ADM

The PySpark API (DataFrames and SQL) is shimmed.

| Before |
|---|

```
1   from pyspark import *
2   from pyspark.sql import SQLContext
3
4   conf = SparkConf()
5   sc = SparkContext(conf=conf)
6   sqlContext = SQLContext(sc)
7
8   names = sc.readText("people.txt")
9
10  namesDF = sc.createDataFrame(names)
11
12  namesDF.registerTempTable("names")
13
14  sqlContext.sql("""
15      SELECT * FROM names
16  """).show()
```

| After |
|---|

```
1   from aqaspark import *
2
3
4   conf = SparkConf()
5   sc = SparkContext(conf=conf)
6   sqlContext = SQLContext(sc)
7
8   names = sc.readText("people.txt")
9
10  namesDF = sc.createDataFrame(names)
11
12  namesDF.registerTempTable("names")
13
14  sqlContext.sql("""
15      SELECT * FROM names
16  """).show()
```

SPARK
SUMMIT
2017

# Wrap Up

Autonomous Data
Management makes Spark
great for SQL analytics.