



SDF & Datafusion

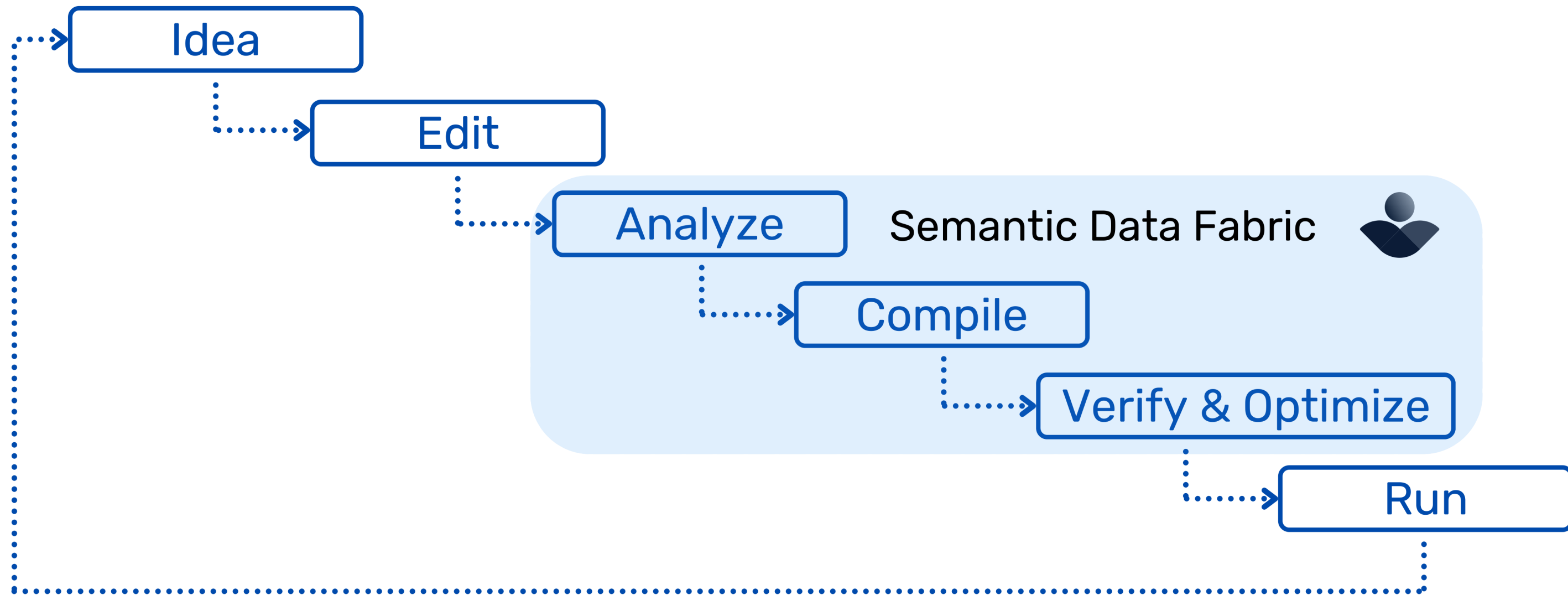
Lukas Schulte (lukas@sdf.com)

Bo Lin (bo@sdf.com)

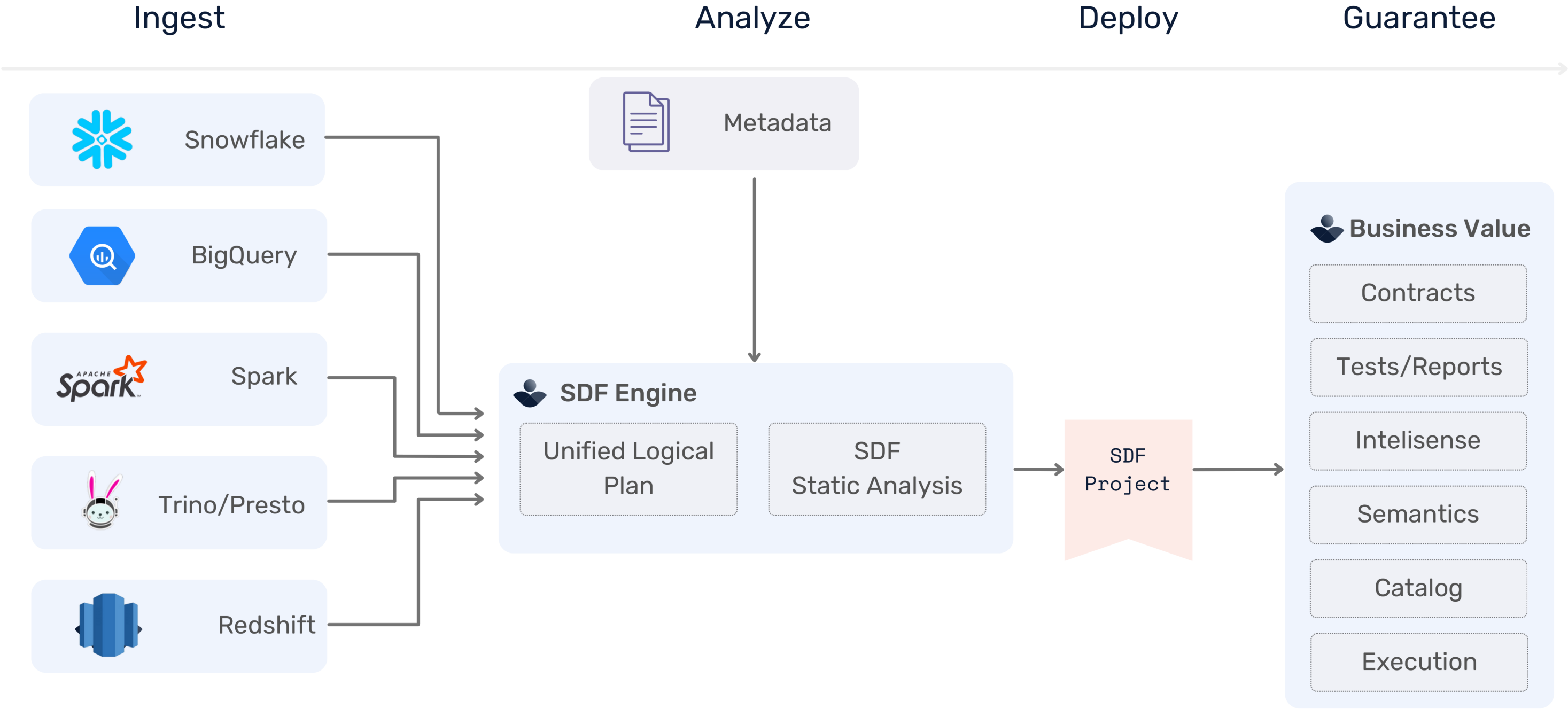
March 25, 2024

What we do

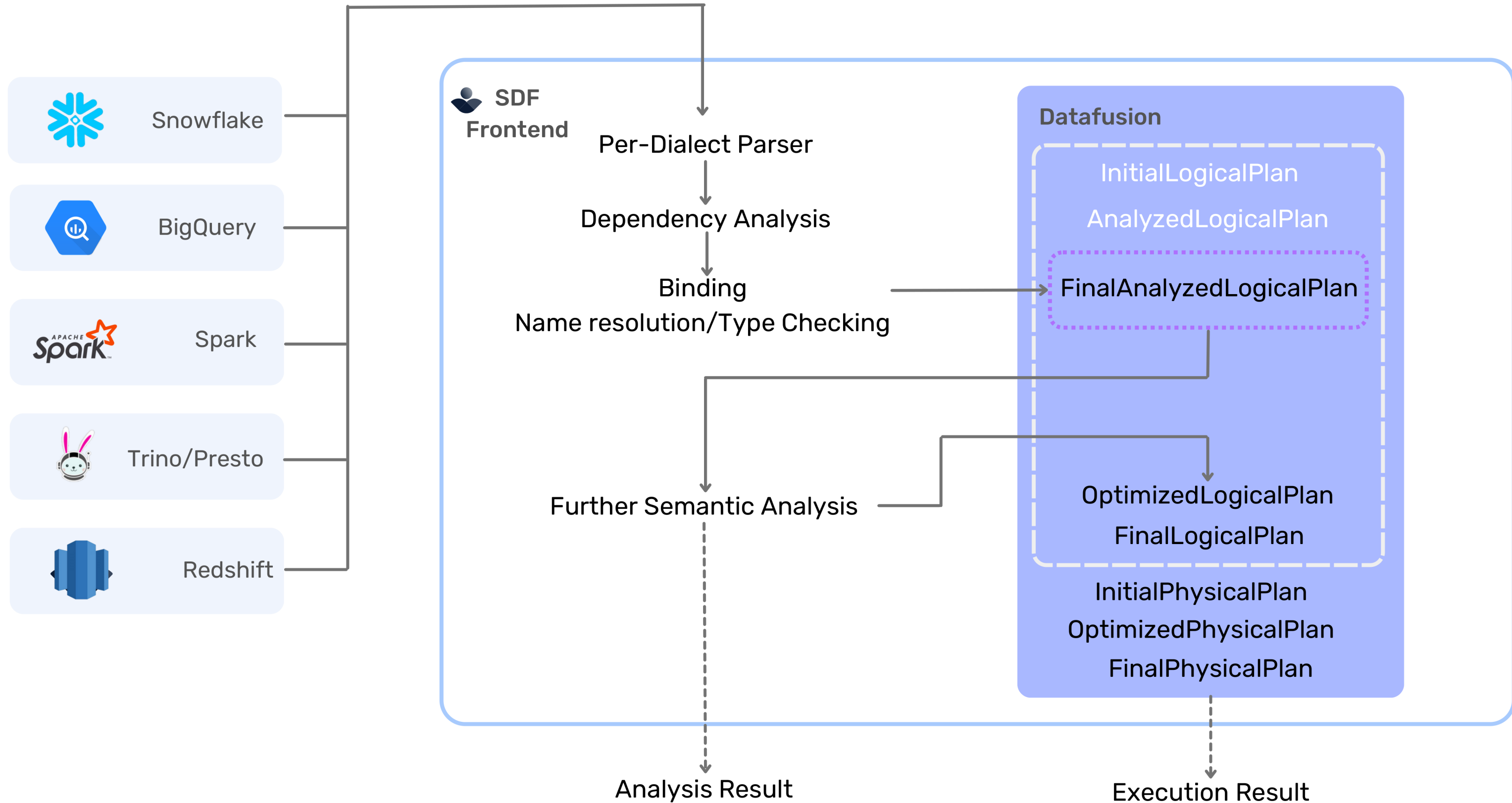
Driving Fastest Cycle Time for Data Development



SDF Architecture



How We Use Datafusion





Improving Datafusion for SQL Analysis

Performance

Column resolution
Expr copy/hashing

Usability

`LogicalPlan` as “dumb structs”

Stability

Clear Separation of Frontend to
Backend



Specialize mid-level IR
What is the role of Substrait?



Other Areas We Hope to Contribute

Functions

Documentation & Implementation

Dialect Specific Functions Crates With Code-Gen

- 738 Trino Functions
- 555 BigQuery Functions
- 268 Redshift Functions
- 1321 Snowflake Functions

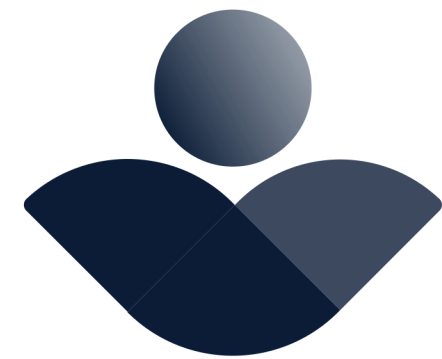
LogicalPlan to SQL/Substrait Translation

Python Bindings

Rust-Python Bindings

- Defining Higher-Order Functions
- Defining Table Functions

```
functions.sdf.yml
1 function:
2   name: all_match
3   parameters:
4     - datatype: array<$1>
5     - datatype: function($1, boolean)
6   optional-parameters: []
7   returns:
8     datatype: boolean
9
10 ---
11 function:
12   name: approx_most_frequent
13   kind: aggregate
14   parameters:
15     - datatype: bigint
16     - datatype: varchar
17     - datatype: bigint
18   optional-parameters: []
19   returns:
20     datatype: map<varchar, bigint>
```



Thank You

Lukas Schulte (lukas@sdf.com)
Bo Lin (bo@sdf.com)