

Why we defined a metalanguage for SQL

Lewis Hemens



We need a **scalable** solution for managing data transformation processes that works for data engineers, analysts and scientists

Why we **love** SQL

SQL is growing in popularity thanks to modern data warehouses

- A common language for data definitions across roles
- Modern warehouse SQL engines scale extremely well
- Easy to iterate, thanks execution usually being one-click
- Relatively easy to debug

But it has some problems...

Why doesn't SQL scale?

It's hard to adopt **software engineering best practices**

- Release processes
- Version control
- Unit tests
- Code reuse

Why are these hard, and how can we fix them?

Understanding SQL

SQL is a **declarative** query language

Declarative programming

When you say what you want

Imperative programming

When you say how to get
what you want

Advantages of being declarative

The fact that SQL is declarative means it has many benefits

- SQL queries can be **parallelized**
- SQL queries can be automatically **optimized**
- For most SQL statements there are **no side effects**
- SQL queries are guaranteed to **eventually terminate**

SQL is **not** a programming language

SQL is few features short of being a programming language

- SQL has little if any control flow
- There is no recursion or iteration*
- SQL is **declarative and static**

*Some flavors of SQL (e.g. T-SQL) add these and are turing complete

Example: writing **reusable** code

```
select
    floor(age / 5) * 5 as age_group,
    count(1) as user_count
from users
group by age_group
```

Example: writing **reusable** code

```
select
  floor(age / 5) * 5 as age_group,
  count(1) as user_count
from users
group by age_group
```

We can't **reuse** this query:
the input is fixed 🤔

Example: writing **testable** code

```
select
  floor(age / 5) * 5 as age_group,
  count(1) as user_count
from users
group by age_group
```

We can't **test** this query for
the same reason 🤔

Example: writing **iterative** code

```
user_tables = ["users", "user_stats", "user_events"]

for table in user_tables:
    delete from table
    where user_id in (
        select user_id from gdpr_deletion_requests
    )
```

Example: writing **iterative** code

```
user_tables = ["users", "user_stats", "user_events"]
```

```
for table in user_tables:  
    delete from table  
    where user_id in (  
        select user_id from gdpr_deletion_requests  
    )
```

No iteration in SQL 😭

Metaprogramming to the rescue

What is metaprogramming?

Metaprogramming is a programming technique in which computer programs have the ability to treat other programs as their data

Metaprogramming can be used to **move computations from run-time to compile-time**

Metaprogramming example

```
select
  floor(age / 5) * 5 as
age_group,
  count(1) as user_count
from users
group by age_group
```

```
function ageDist(input, bucket = 5) {
  return `
    select
      floor(age / ${bucket}) *
${bucket} as age_group,
      count(1) as user_count
    from ${input}
    group by age_group`;
}
```


Fixing SQL with meta programming

- Enable code reuse through parameterizable functions
- Allow *some* imperative programming
- Introduce *some* control flow
- **Keep our code declarative at run-time**

Dataform

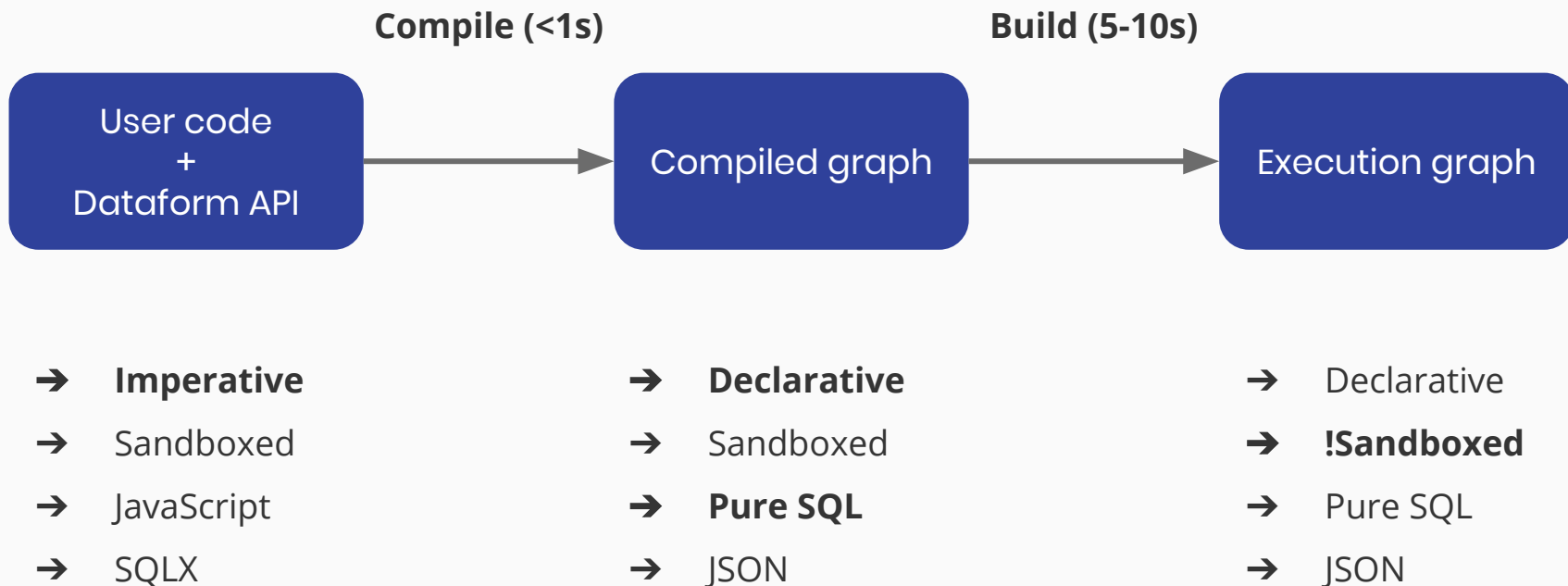
framework

An open-source framework and
metalanguage for SQL

Dataform framework overview

- Makes it easy to write **parameterized SQL**
- Enables **code reuse**
- APIs to help build **directed acyclic graphs**
- Support for writing **data assertions**
- Support for writing **SQL unit tests**
- APIs for **documenting** datasets
- Support for managing multiple **environments**

Dataform compilation process



By introducing a compilation step, we maintain a static, declarative runtime format

Dataform example: Publishing a table

```
// copy_users_table.js

publish("copy_users_table")
  .type("table")
  .query(ctx => `
    select *
    from ${ctx.ref("users")}`
  );
```

```
-- copy_users_table.sqlx

config {
  type: "table"
}

select *
from ${ref("users")}
```

(Our SQL queries are now parameterized!)

Dataform compiled graph

```
{  
  "name": "dataform_dc_talk.copy_users_table",  
  "dependencies": ["dataform_dc_talk.users"],  
  "type": "table",  
  "target": {  
    "schema": "dataform_dc_talk",  
    "name": "copy_users_table"  
  },  
  "query": "select * from dataform_dc_talk.users",  
}
```

Running dataform projects

```
$ dataform compile
```

```
Compiling...
```

```
Compiled 56 action(s).
```

```
35 dataset(s):
```

```
  dataform_data.organisations [view]
```

```
  dataform_data.project_users [view]
```

```
  dataform_data.projects [view]
```

```
  dataform_data.sessions [table]
```

```
  dataform_data.users [view]
```

```
...
```

```
$ dataform run
```

```
Compiling...
```

```
Compiled successfully.
```

```
Running...
```

```
Dataset created: dataform_data.organisations [view]
```

```
Dataset created: dataform_data.project_users [view]
```

```
Dataset created: dataform_data.projects [view]
```

```
Dataset created: dataform_data.sessions [table]
```

```
Dataset created: dataform_data.users [view]
```

```
...
```

Dataform framework summary

It's basically a SQL compiler.

- We can write any* code we like during the **compilation phase**
- Dataform's runtime format is declarative, pure SQL with only non-iterative control structures during the **runtime phase**

A note on reproducibility

- Executing compile on the same project with the same parameters twice should always yield the same result
- Sandboxing helps enforce this, no network requests, file reads, or or DB access possible
- User should avoid non deterministic algorithms (e.g. `Math.random()`)

Dataform examples

For loops

```
const userTables = ["users", "user_stats", "user_events"];

userTables.forEach(tableName =>
  operate(
    `${tableName}_gdpr_cleanup`,
    ctx => `
      delete from ${ctx.ref(tableName)}
      where user_id in (
        select user_id
        from ${ctx.ref("gdpr_deletion_requests")}
      )`
  )
);
```

Unit testing

```
const ageDist = (input, bucket = 5) => `  
  select  
    floor(age / ${bucket}) * ${bucket} as age_group,  
    count(1) as user_count  
  from ${input}  
  group by age_group`;
```

```
publish("users_by_age")  
  .query(ctx => ageDist(ctx.ref("users"), 5))  
  .type("table");
```

Unit testing

```
test("ageDist_test")
  .query(ageDist(`(
    select 15 as age union all
    select 21 as age union all
    select 24 as age)` ,
    10
  ))
  .expected(`
    select 10 as age_group, 1 as user_count union all
    select 20 as age_group, 2 as user_count
  `);
```

Can be run with:
dataform test

Environment sampling

```
publish (
  "sourcetable_view",
  ctx => `
    select *
    from ${ctx.ref("sourcetable")}
    where ${
      ctx.env === "staging"
        ? `rand() < ${constants.stagingSamplingRate}`
        : "true"
    }
  `
);
```

Loading data from S3

```
s3_load_csv("load_example_csv", {
  path: "s3://.../sample_data.csv",
  schema: {
    country: "varchar(256)",
    revenue: "float8"
  },
  role: "arn:aws:iam:... ",
  ignoreheader: true
});
```

```
create table
dataform_dc_talk.load_example_csv (
  country varchar(256),
  revenue FLOAT8
);
copy
dataform_dc_talk.load_example_csv
from 's3://.../sample_data.csv'
iam_role 'arn:aws:iam:... '
ignoreheader 1
delimiter ',';
```

Dataform

Web

An collaborative IDE and
deployment platform for
dataform projects

Dataform Web

What good is a new language without an IDE?

IDE for dataform projects

- Compiles and validates SQL graph in real time
- Full integration with Git, user branches and pull requests
- Managed continuous deployment and environments

But also:

- **Pipeline orchestration**, run logs and notifications
- Documentation and **Data catalog**



dataform-data

Production

lewis_dev



Pull from master



Shortcuts



Docs



Start new run

Overview

Search

FILES

DATASETS



CONFIG

config dataform.json

config package.json

config schedules.json

FILES

> analysis

▼ definitions

> assertions

> empty

> entities

> examples

> refresh

> sources

▼ stats

sqlx active_tr...

definitions/stats/user_stats.sqlx

Saved

```

1
2 config {
3   type: "table",
4   schema: "dataform_data",
5   description: "Stats about users.",
6   columns: {
7     overall_sessions: { description: "Number of sessions this user has had",
8     tags: ["pii"], deprecated: true }
9   }
10 }
11 /* Stats about users. */
12
13 with
14 session_stats as (
15
16   select
17     user_id,
18     sum(1) as sessions,
19     min(session_start_timestamp) as first_seen,
20     max(session_end_timestamp) as last_seen,
21     sum(IF(session_start_timestamp > TIMESTAMP_SUB(current_timestamp(),
INTERVAL 7 DAY), TIMESTAMP_DIFF(session_end_timestamp,
session_start_timestamp, MINUTE), 0)) AS sessions_minutes_last_7d
22   from
23     ${ref("sessions")}
24   group by
25     1
26 ),
27
28 project_stats as (
29
30   ${aggregate(
31     ref("daily_project_user_stats"),
32     ["user_id"],

```

Try a predefined template

dataform_data.user_stats

table

Dependencies

7



Compiled query



Query validation

Valid

▶ Preview results

Publish table





dataform-data

Production

lewis_dev



Pull from master



Shortcuts

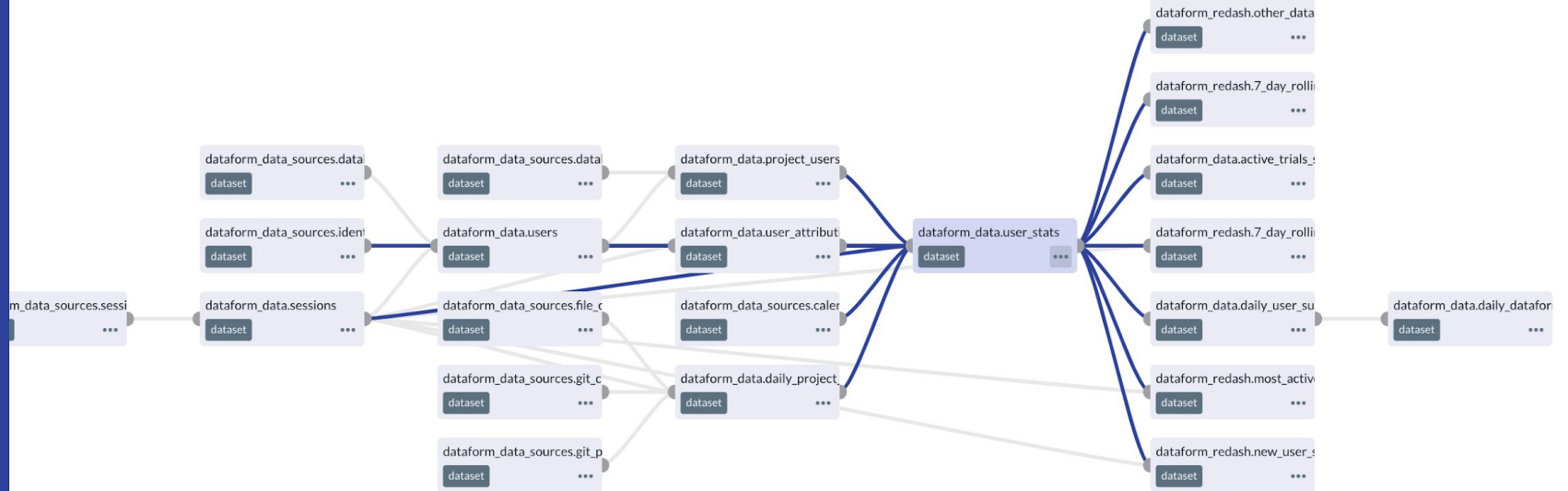


Start new run

dataform_data.user_stats



Include assertions





definitions / reporting_tables / country_reporting.sqlx

SAVE

```
1  config {
2    type: "table",
3    tags: ["core", "daily"],
4    description: "This table contains user generated revenue,
5    pageviews and sessions aggregated
6    by country",
7    columns:{
8      date: "Timestamp",
9      country: "Country of the user",
10     revenue: "Revenue is USD in floating FX",
11     sessions: "Sessions count from Google Analytics",
12     device_type: "Desktop, Tablet or Mobile"
13   }
14 }
15
16 SELECT date AS date,
17        country AS country,
18        device_type AS device_type,
19        revenue AS revenue,
20        sessions AS sessions
21
22 FROM ${ref("source_table")}
```



country_reporting

SCHEMA

dataform

TYPE

table

TAGS

core

daily

DESCRIPTION

This table contains user generated revenue, pageviews and sessions aggregated by country

Columns

Field	Description
date	Timestamp
device_type	Desktop, Tablet or Mobile
country	Country of the user
sessions	Sessions count from Google Analytics
revenue	Revenue in USD using floating FX rates

Lineage

Thanks!

Dataform framework: github.com/dataform-co/dataform

Dataform docs: docs.dataform.co

Examples from this talk: github.com/dataform-co/dataform-dc-talk

Dataform web IDE: dataform.co

Questions?