

Hamilton

**enabling software engineering best practices
for data transformations
via generalized dataflow graphs**

Stefan Krawczyk, Elijah ben Izzy, Danielle Quinn

[@ Stitch Fix](#)

CDMS Workshop VLDB 2022

Introduction

Introduction

Context:

- Stitch Fix is a business where “machine learning” is core to the product
- Stitch Fix has 100+ Data Scientists (DS)
 - No hand-off; DS responsible for productionization*
 - DS own ETLs on top of a data lakehouse
- Data platform team goals:
 - Capabilities
 - Iteration speed
 - Maintenance

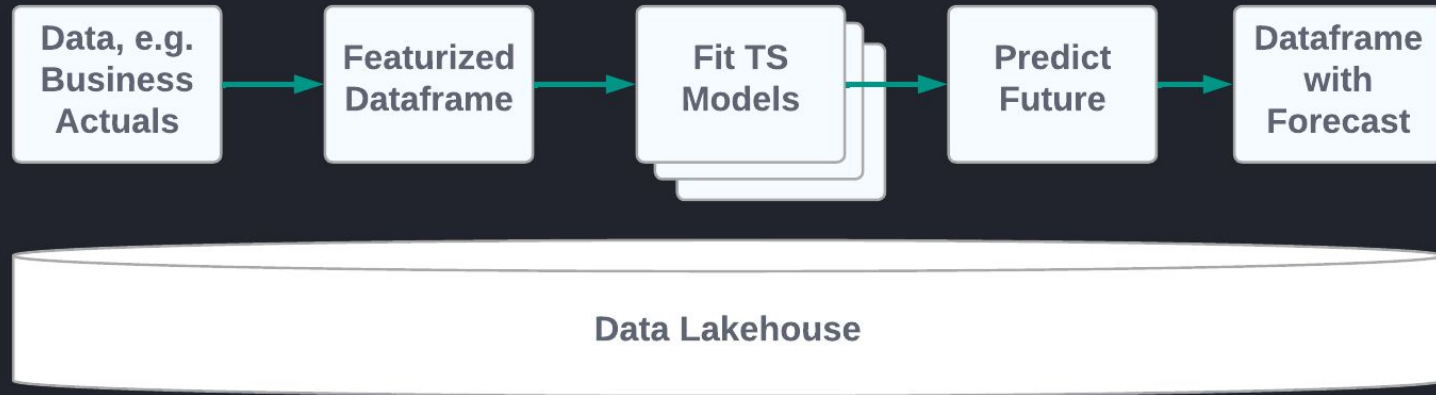
Introduction

Connection with Data ecosystems:

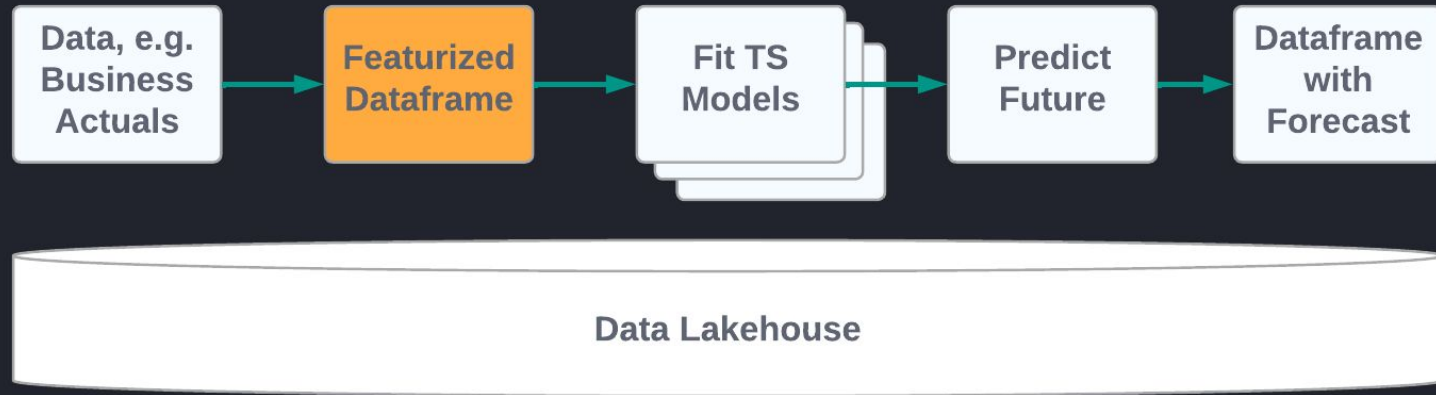
- ML Industry perspective
- Code is metadata to be managed
 - Feeds into data management processes
 - Code produces/augments data; impacts developer productivity.
 - Want granular traceability/lineage? Code is the best source.
 - Code and data quality (DQ) are coupled.
 - Code [& data] only grows
 - additions + modifications >> deletions.

Software Engineering Pain Points with Data Transformations

Software Engineering Pain Points with Data Transformations



Software Engineering Pain Points with Data Transformations



Example transform code

```
df = loader.load_actuals(dates) # e.g. spend, signups
```


Example transform code

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
```

Example transform code

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
```

Example transform code

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
df['special_feature1'] = compute_bespoke_feature(df)
df['spend_b'] = multiply_columns(df['acquisition_cost'], df['B'])
```

Example transform code

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
df['special_feature1'] = compute_bespoke_feature(df)
df['spend_b'] = multiply_columns(df['acquisition_cost'], df['B'])
save_df(df, "some_location")
```

Example transform code

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
df['special_feature1'] = compute_bespoke_feature(df)
df['spend_b'] = multiply_columns(df['acquisition_cost'], df['B'])
save_df(df, "some_location")
```

Now scale this code to 1000+ columns & a growing team



Problem: unit testing & integration testing 🙅

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
df['special_feature1'] = compute_bespoke_feature(df)
df['spend_b'] = multiply_columns(df['acquisition_cost'], df['B'])
save_df(df, "some_location")
```

Now scale this code to 1000+ columns & a growing team



Problem: code readability & documentation 🤔

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
df['special_feature1'] = compute_bespoke_feature(df)
df['spend_b'] = multiply_columns(df['acquisition_cost'], df['B'])
save_df(df, "some_location")
```



Now scale this code to 1000+ columns & a growing team



Problem: difficulty in tracing lineage 🧠

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
→ df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
df['special_feature1'] = compute_bespoke_feature(df)
→ df['spend_b'] = multiply_columns(df['acquisition_cost'], df['B'])
save_df(df, "some_location")
```

Now scale this code to 1000+ columns & a growing team



Problem: code reuse and duplication

```
df = loader.load_actuals(dates) # e.g. spend, signups
if config['region'] == 'UK':
    df['holidays'] = is_uk_holiday(df['year'], df['week'])
else:
    df['holidays'] = is_holiday(df['year'], df['week'])
df['avg_3wk_spend'] = df['spend'].rolling(3).mean()
df['acquisition_cost'] = df['spend'] / df['signups']
df['spend_shift_3weeks'] = df['spend'].shift(3)
df['special_feature1'] = compute_bespoke_feature(df)
df['spend_b'] = multiply_columns(df['acquisition_cost'], df['B'])
save_df(df, "some_location")
```



Now scale this code to 1000+ columns & a growing team



Hamilton

Hamilton:

Code → Dataflow → Object

Hamilton:

Code → Dataflow → Object

Code:

```
def holidays(year: pd.Series, week: pd.Series) -> pd.Series:
    """Some docs"""
    return some_library(year, week)
def avg_3wk_spend(spend: pd.Series) -> pd.Series:
    """Some docs"""
    return spend.rolling(3).mean()
def spend_per_signup(spend: pd.Series, signups: pd.Series) -> pd.Series:
    """Some docs"""
    return spend / signups
def spend_shift_3weeks(spend: pd.Series) -> pd.Series:
    """Some docs"""
    return spend.shift(3)
def spend_shift_3weeks_per_signup(spend_shift_3weeks: pd.Series, signups: pd.Series) -> pd.Series:
    """Some docs"""
    return spend_shift_3weeks / signups
```

User



Hamilton:

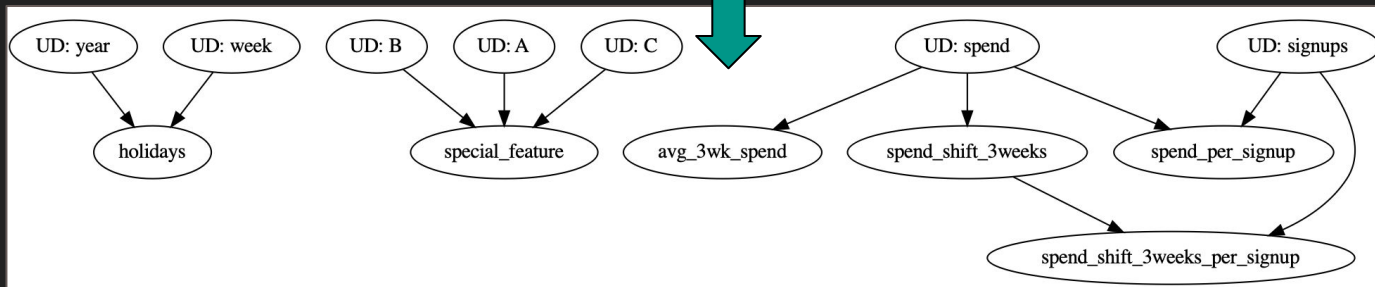
Code → Dataflow → Object

Code:

```
def holidays(year: pd.Series, week: pd.Series) -> pd.Series:  
    """Some docs"""  
    return some_library(year, week)  
def avg_3wk_spend(spend: pd.Series) -> pd.Series:  
    """Some docs"""  
    return spend.rolling(3).mean()  
def spend_per_signup(spend: pd.Series, signups: pd.Series) -> pd.Series:  
    """Some docs"""  
    return spend / signups  
def spend_shift_3weeks(spend: pd.Series) -> pd.Series:  
    """Some docs"""  
    return spend.shift(3)  
def spend_shift_3weeks_per_signup(spend_shift_3weeks: pd.Series, signups: pd.Series) -> pd.Series:  
    """Some docs"""  
    return spend_shift_3weeks / signups
```

User

DAG:



Hamilton

Hamilton:

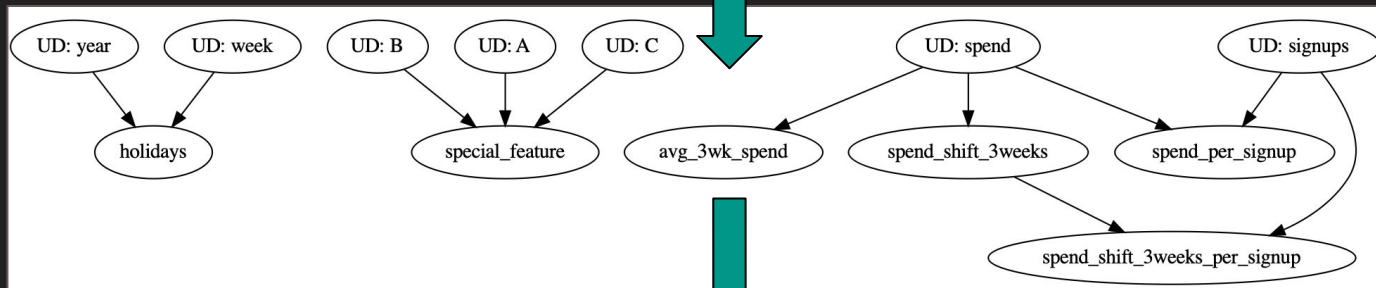
Code → Dataflow → Object

Code:

```
def holidays(year: pd.Series, week: pd.Series) -> pd.Series:
    """Some docs"""
    return some_library(year, week)
def avg_3wk_spend(spend: pd.Series) -> pd.Series:
    """Some docs"""
    return spend.rolling(3).mean()
def spend_per_signup(spend: pd.Series, signups: pd.Series) -> pd.Series:
    """Some docs"""
    return spend / signups
def spend_shift_3weeks(spend: pd.Series) -> pd.Series:
    """Some docs"""
    return spend.shift(3)
def spend_shift_3weeks_per_signup(spend_shift_3weeks: pd.Series, signups: pd.Series) -> pd.Series:
    """Some docs"""
    return spend_shift_3weeks / signups
```

User

DAG:



Hamilton

Object(s)
(e.g. Dataframe,
ML Model):

Year	Week	Sign ups	...	Spend	Holiday
2015	2	57	...	123	0
2015	3	58	...	123	0
2015	4	59	...	123	1
2015	5	59	...	123	1
...
...
...
...
2021	16	1000	...	1234	0

User

Hamilton:

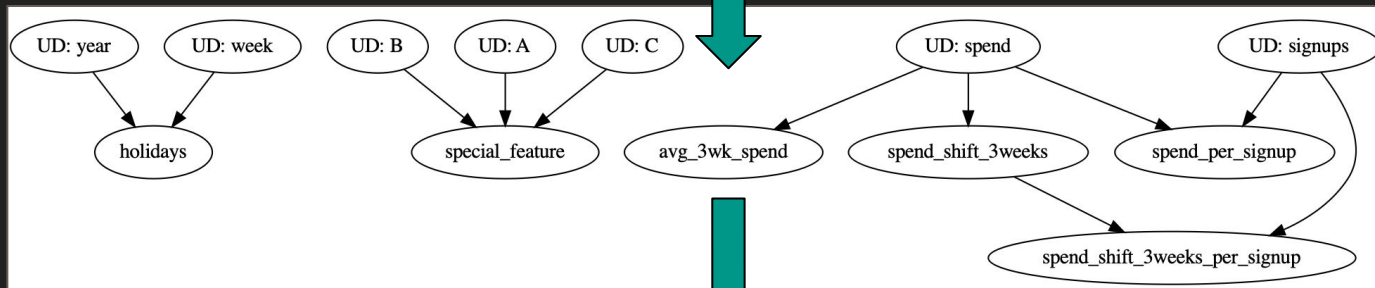
Code → Dataflow → Object

Code:

```
def holidays(year: pd.Series, week: pd.Series) -> pd.Series:
    """Some docs"""
    return some_library(year, week)
def avg_3wk_spend(spend: pd.Series) -> pd.Series:
    """Some docs"""
    return spend.rolling(3).mean()
def spend_per_signup(spend: pd.Series, signups: pd.Series) -> pd.Series:
    """Some docs"""
    return spend / signups
def spend_shift_3weeks(spend: pd.Series) -> pd.Series:
    """Some docs"""
    return spend.shift(3)
def spend_shift_3weeks_per_signup(spend_shift_3weeks: pd.Series, signups: pd.Series) -> pd.Series:
    """Some docs"""
    return spend_shift_3weeks / signups
```

Python
Modules

DAG:



Object(s)
(e.g. Dataframe,
ML Model):

Year	Week	Sign ups	...	Spend	Holiday
2015	2	57	...	123	0
2015	3	58	...	123	0
2015	4	59	...	123	1
2015	5	59	...	123	1
...
...
...
...
2021	16	1000	...	1234	0

"Driver" Code

Hamilton Paradigm: declaring a dataflow

Instead of:

```
df['c'] = df['a'] + df['b']  
df['d'] = transform(df['c'])
```

You declare:

```
def c(a: pd.Series, b: pd.Series) -> pd.Series:  
    """Sums a with b"""  
    return a + b  
  
def d(c: pd.Series) -> pd.Series:  
    """Transforms C to ..."""  
    new_column = _transform_logic(c)  
    return new_column
```

+ some driver code (not shown)

Hamilton Paradigm: declaring a dataflow

Instead of:

```
df['c'] = df['a'] + df['b']  
df['d'] = transform(df['c'])
```

Outputs == Function Name

Inputs == Function Arguments

You declare:

```
def c(a: pd.Series, b: pd.Series) -> pd.Series:  
    """Sums a with b"""  
    return a + b  
  
def d(c: pd.Series) -> pd.Series:  
    """Transforms C to ..."""  
    new_column = _transform_logic(c)  
    return new_column
```

Full Hello World

Functions:

```
# feature_logic.py
def c(a: pd.Series, b: pd.Series) -> pd.Series:
    """Sums a with b"""
    return a + b

def d(c: pd.Series) -> pd.Series:
    """Transforms C to ..."""
    new_column = _transform_logic(c)
    return new_column
```

“Driver” – this actually says what and when to execute:

```
# run.py
from hamilton import driver
import feature_logic
dr = driver.Driver({'a': ..., 'b': ...}, feature_logic)
df_result = dr.execute(['c', 'd'])
print(df_result)
```

Full Hello World

Functions:

```
# feature_logic.py
def c(a: pd.Series, b: pd.Series) -> pd.Series:
    """Sums a with b"""
    return a + b

def d(c: pd.Series) -> pd.Series:
    """Transforms C to ..."""
    new_column = _transform_logic(c)
    return new_column
```

“Driver” – this actually says what and when to execute:

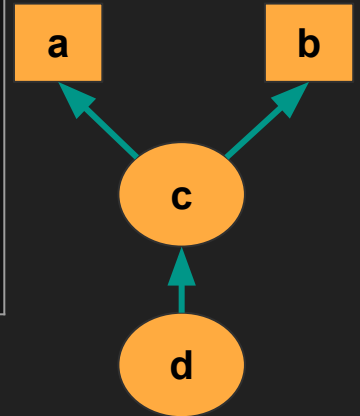
```
# run.py
from hamilton import driver
import feature_logic
dr = driver.Driver({'a': ..., 'b': ...}, feature_logic)
df_result = dr.execute(['c', 'd'])
print(df_result)
```

Full Hello World

Functions:

```
# feature_logic.py
def c(a: pd.Series, b: pd.Series) -> pd.Series:
    """Sums a with b"""
    return a + b

def d(c: pd.Series) -> pd.Series:
    """Transforms C to ..."""
    new_column = _transform_logic(c)
    return new_column
```



“Driver” – this actually says what and when to execute:

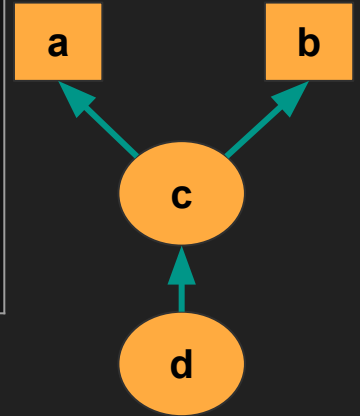
```
# run.py
from hamilton import driver
import feature_logic
dr = driver.Driver({'a': ..., 'b': ...}, feature_logic)
df_result = dr.execute(['c', 'd'])
print(df_result)
```

Full Hello World

Functions:

```
# feature_logic.py
def c(a: pd.Series, b: pd.Series) -> pd.Series:
    """Sums a with b"""
    return a + b

def d(c: pd.Series) -> pd.Series:
    """Transforms C to ..."""
    new_column = _transform_logic(c)
    return new_column
```



“Driver” – this actually says what and when to execute:

```
# run.py
from hamilton import driver
import feature_logic
dr = driver.Driver({'a': ..., 'b': ...}, feature_logic)
df_result = dr.execute(['c', 'd'])
print(df_result)
```

**Enabling software engineering
best practices & then some**

Software engineering best practices & then some:

```
# located in client_features.py
@tag(owner='Data-Science', pii='False')
@check_output(data_type=np.float64, range=(-5.0, 5.0), allow_nans=False)
def height_zero_mean_unit_variance(height_zero_mean: pd.Series,
                                   height_std_dev: pd.Series) -> pd.Series:
    """Zero mean unit variance value of height"""
    return height_zero_mean / height_std_dev
```

Software engineering best practices & then some:

```
# located in client_features.py
@tag(owner='Data-Science', pii='False')
@check_output(data_type=np.float64, range=(-5.0, 5.0), allow_nans=False)
def height_zero_mean_unit_variance(height_zero_mean: pd.Series,
                                   height_std_dev: pd.Series) -> pd.Series:
    """Zero mean unit variance value of height"""
    return height_zero_mean / height_std_dev
```

- Unit testing & integration testing
 - Documentation
 - Modularity/reuse
- ✓ always possible, easy to add
 - ✓ function doc, visualization, @tag
 - ✓ module curation & drivers

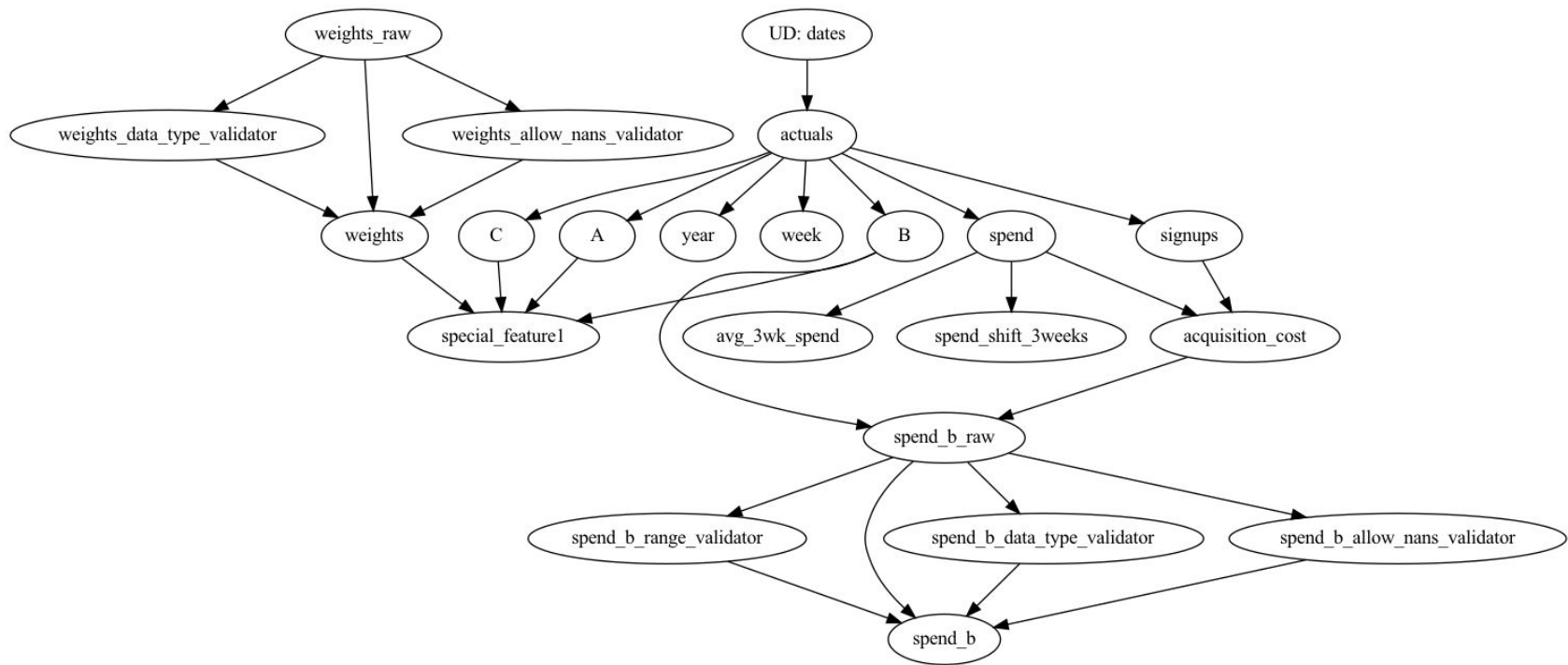
Software engineering best practices & then some:

```
# located in client_features.py
@tag(owner='Data-Science', pii='False')
@check_output(data_type=np.float64, range=(-5.0, 5.0), allow_nans=False)
def height_zero_mean_unit_variance(height_zero_mean: pd.Series,
                                   height_std_dev: pd.Series) -> pd.Series:
    """Zero mean unit variance value of height"""
    return height_zero_mean / height_std_dev
```

- Unit testing & integration testing
- Documentation
- Modularity/reuse
- Data quality
- Lineage
- ✓ always possible, easy to add
- ✓ function doc, visualization, @tag
- ✓ module curation & drivers
- ✓ @check_output runtime checks
- ✓ DAG, versioning, @tag

... more in the paper!

Example Dataflow Visualization



Evaluation

Hamilton @ Stitch Fix

- Running in production for 2.5+ years
- FED team manages 4000+ feature definitions
 - All feature definitions are:
 - Unit testable
 - Documentation friendly
 - Centrally curated, stored, and versioned in git.
- Data Science teams ❤️ it:
 - Best adoption from active time-series forecasting teams
 - Most willing to pay migration cost.
 - Enabled a monthly feature update & model fitting task to be completed 4x faster
- Open source still early

Hamilton @ Stitch Fix

- Running in production for 2.5+ years
- FED team manages 4000+ feature definitions
 - All feature definitions are:
 - Unit testable
 - Documentation friendly
 - Centrally curated, stored, and versioned in git.
- Data Science teams ❤️ it:
 - Best adoption from active time-series forecasting teams
 - Most willing to pay migration cost.
 - Enabled a monthly feature update & model fitting task to be completed 4x faster
- Open source still early

WORK IN PROGRESS

Summary & Future Work

Summary:

- A declarative [dataflow](#) paradigm in python
- Functions, via naming, encode a dataflow
- Software engineering best practices come from:
 - Encapsulation of transform logic within functions.
 - Decoupling transform logic from materialization.

Future Work:

- Source code based governance
 - How do we integrate it further?
- Compiling to an orchestration framework
- Modeling your entire data lakehouse independently of materialization concerns

Hamilton is Open Source Code

```
> pip install sf-hamilton
```

Get started in <15 minutes!

Star  on github:

<https://github.com/stitchfix/hamilton>

Documentation

<https://hamilton-docs.gitbook.io/>

Various examples:

<https://github.com/stitchfix/hamilton/tree/main/examples>

Thank you.

Questions?

<https://twitter.com/stefkrawczyk>

<https://www.linkedin.com/in/skrawczyk/>

<https://github.com/stitchfix/hamilton>