

Scalable and Sustainable Feature Engineering with Hamilton

Elijah ben Izzy
CTO/Co-founder @ DAGWorks, Inc.

TL;DR

I want to convince you that...

1. Maintaining feature engineering code is difficult
2. Hamilton can help you:
 - a. build *sustainable* code
 - b. build *scalable* code
3. Hamilton can model your ML workflow end to end
4. Hamilton is easy to get started with/easy to use!



**At DAGWorks we're making ML pipelines easy to manage.
*Nobody should be afraid to inherit data science code.***

>>> I'm not selling you anything in this talk! <<<

Hamilton is Open Source!!

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

Get started in <15 minutes!

Try it out

<https://www.tryhamilton.dev/>

Documentation

<https://hamilton.readthedocs.io/>

<https://www.tryhamilton.dev/>

Hamilton

Wrangle Pandas codebases into shape.

🕒 Learn (5 mins)

🔄 Github 890+ ⭐

- ✓ Write always unit testable code
- ✓ Add runtime data validation easily
- ✓ Produce readable and maintainable code
- ✓ Visualize lineage (click the run button to see)
- ✓ Run anywhere python runs: in airflow, jupyter, fastapi, etc...
- ✓ Skip the CS degree to use it

Try Hamilton right here in your browser 📌

```
1 # Declare and link your transformations as functions....
2 import pandas as pd
3
4 def a(input: pd.Series) -> pd.Series:
5     return input % 7
6
7 def b(a: pd.Series) -> pd.Series:
8     return a * 2
9
10 def c(a: pd.Series, b: pd.Series) -> pd.Series:
11     return a * 3 + b * 2
12
13 def d(c: pd.Series) -> pd.Series:
14     return c ** 3
```

```
1 # And run them!
2 import functions
3 from hamilton import driver
4 dr = driver.Driver({}, functions)
5 result = dr.execute(
6     ['a', 'b', 'c', 'd'],
7     inputs={'input': pd.Series([1, 2, 3, 4, 5])}
8 )
9 print(result)
10 dr.display_all_functions("graph.dot", {})
```

▶ Run me!

The Agenda

A motivating story of DS pain

The solution: *Hamilton*

Hamilton for feature engineering

- ↳ **Sustainable feature management**

- ↳ **Scalable feature pipelines**

Hamilton for end-to-end ML workflows

OS progress + next steps

The Agenda

A motivating story of DS pain

The solution: *Hamilton*

Hamilton for feature engineering

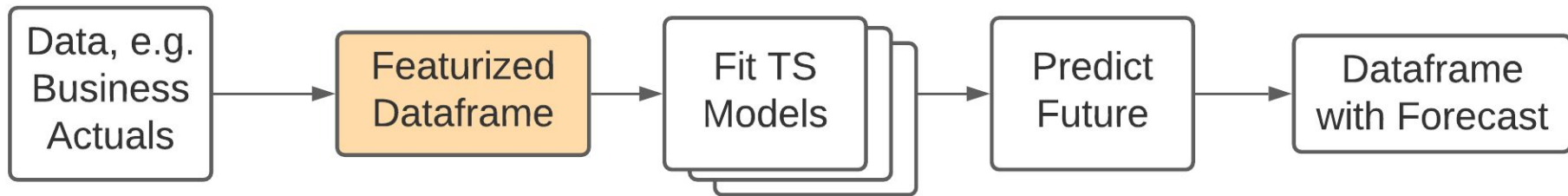
- ↳ **Sustainable feature management**
- ↳ **Scalable feature pipelines**

Hamilton for end-to-end ML workflows

OS progress + next steps

A Problem From my Last Job...

Forecasting the business (demand, signups, churn)



Some Business

Approach

- $O(1000+)$ operations
- Configurations based on data
- Multiple layers of abstraction

A world-class team

Amazon.com bestseller • *New York Times* bestseller
Wall Street Journal bestseller

What Got You Here Won't Get You There

How Successful People Become
Even More Successful!



MARSHALL GOLDSMITH
WITH MARK REITER



abstraction

visualize data

Some Business-Critical Tech Debt

Problems with the code?

- ❑ **Unit testing:** *difficult*
- ❑ **Documentation:** *unnatural, unenforced*
- ❑ **Modularity:** *non-existent*
- ❑ **Data catalogue:** *lots of grepping*
- ❑ **Debugging:** *run the whole pipeline*
- ❑ **Data validation:** *run the whole pipeline, not really done*

Perfect solution to forecasting problem + time = spaghetti code

Some Business-Critical Tech Debt

Q: What happens when you have all of those problems, and...

- ❑ You want to expand your models to new regions?
- ❑ You have to add complex scenarios on management's whim?
- ❑ You have a data bug and very little time to solve it?

A: It wasn't fun.

- + This is not a unique experience to my prior role, time-series forecasting, or even pandas

**I DON'T ALWAYS WRITE
COMPLEX DATA PIPELINES**



**BUT WHEN I DO,
ITS AN UNINTELLIGABLE
MESS OF SPAGHETTI CODE**

The Agenda

A motivating story of DS pain

The solution: *Hamilton*

Hamilton for feature engineering

- ↳ Sustainable feature management
- ↳ Scalable feature pipelines

Hamilton for end-to-end ML workflows

OS progress + next steps

Hamilton: the “A-ha” Moment

Idea: What if every column corresponded to exactly one python fn?

Idea 2: What if the way that function was written tells you everything you needed to know?

*In Hamilton, the artifact (column) is determined by the **name of the function**.
The dependencies are determined by **the parameters**.*

Old Way vs Hamilton Way:

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
```

Hamilton supports *all*** python objects, not just dfs/series!*

Old Way vs Hamilton Way:

Instead of

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

Outputs == Function Name

You declare

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

Inputs == Function Arguments

```
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 says what/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)
```

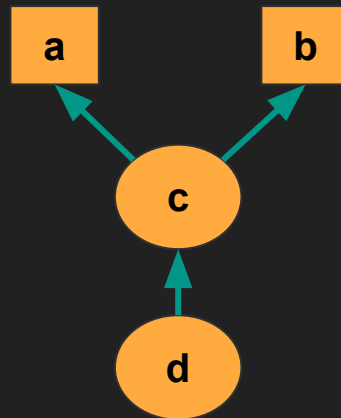
Hamilton TL;DR:

1. For each transform (=), you write a function(s)
2. Functions declare a DAG
3. Hamilton handles DAG execution

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

```
def d(c: pd.Series) -> pd.Series:
    """Replaces d = transform(c)"""
    new_column = _transform_logic(c)
    return new_column
```

```
# 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)
```



Hamilton: Extensions

Q: Doesn't Hamilton make your code more verbose?

A: Yes, but that's not always a bad thing. When it is, we have decorators!

- ❑ `@tag` # attach metadata
- ❑ `@parameterize` # curry + repeat a function
- ❑ `@extract_columns` # one dataframe -> multiple series
- ❑ `@check_output` # data validation
- ❑ `@config.when` # conditional transforms
- ❑ `@subdag` # recursively utilize groups of nodes
- ❑ `@...` # new ones all the time

To Summarize...

Hamilton forces you to write transforms in python functions.

These python functions provide everything you need:

- ❑ **Unit testing:** *simple – plain python functions!*
- ❑ **Documentation:** *use the docstring*
- ❑ **Modularity:** *Small pieces -> by definition*
- ❑ **Data catalogue:** *Code = central feature definition store*
- ❑ **Debugging:** *Execute functions individually + breakpoints*
- ❑ **Trustworthy data:** *Validation included out of the box*

Decorators → powerful, higher-order operations

Driver → decouple transform definition from execution

Initial Use Case

Running in production for **3+** years

Initial use-case manages **4000+** feature definitions

Data science teams ❤️ it

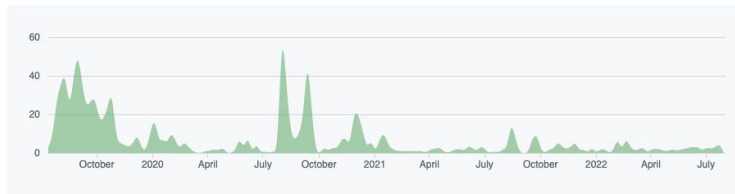
- ❑ Enabled 4x faster monthly model + feature update
- ❑ Easy to onboard new team members
- ❑ Code reviews are simple
- ❑ Finally have unit tests
- ❑ Fewer bugs/quicker resolutions
- ❑ Better features + models

Initial Use Case

Jul 14, 2019 – Aug 3, 2022

Contributions: Commits

Contributions to master, excluding merge commits and bot accounts



elijahbenizzy

304 commits 39,474 ++ 30,896 --

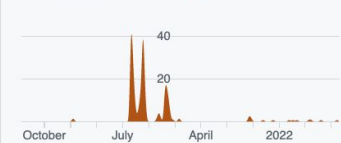
#1



vytlacil

249 commits 159,495 ++ 82,686 --

#2



skrawcz

176 commits 15,868 ++ 3,918 --

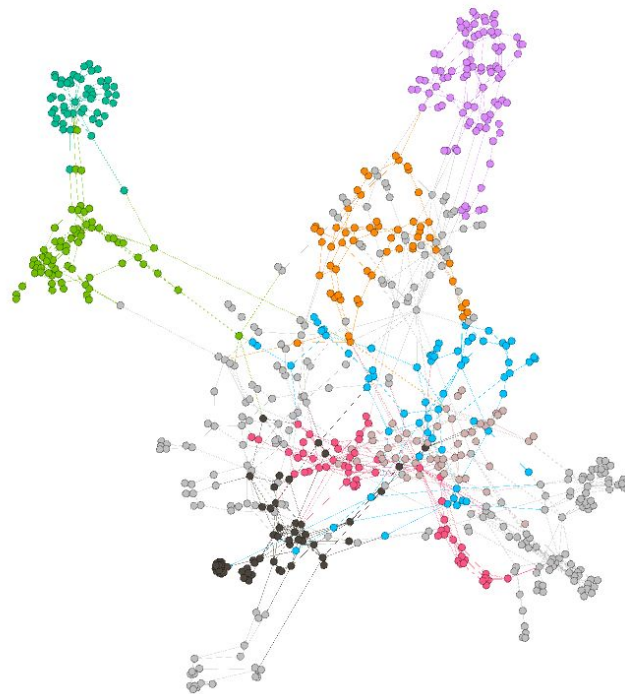
#3



shellyjang

90 commits 17,883 ++ 11,029 --

#4



The Agenda

A motivating story of DS pain

The solution: *Hamilton*

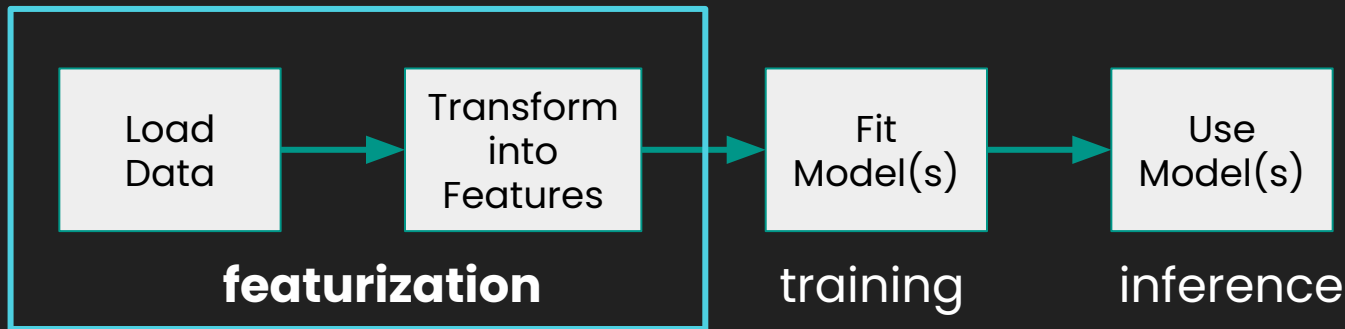
Hamilton for feature engineering

- ↳ Sustainable feature management
- ↳ Scalable feature pipelines

Hamilton for end-to-end ML workflows

OS progress + next steps

Hamilton + Feature Engineering: Overview



Note:

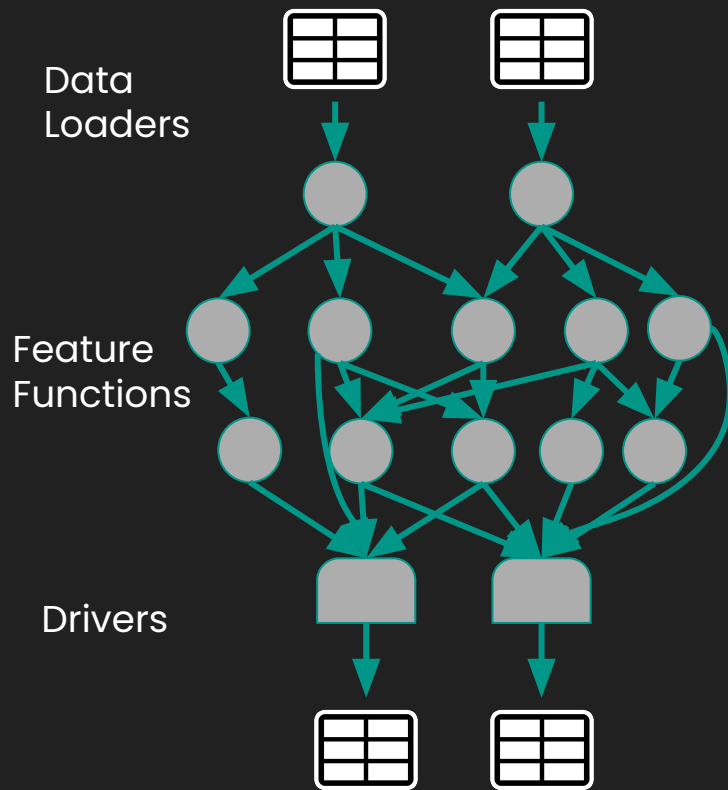
- ❑ Works for any python object type (not just dataframes!)
- ❑ Embeddable anywhere python runs – orchestration systems (airflow, kubeflow, metaflow, flytekit, prefect, dagster, ...) + web services!

Modeling Feature Engineering

Code that needs to be written:

1. Functions to load data
2. Transform/feature functions
3. Driver to materialize data

Execute only what's needed...

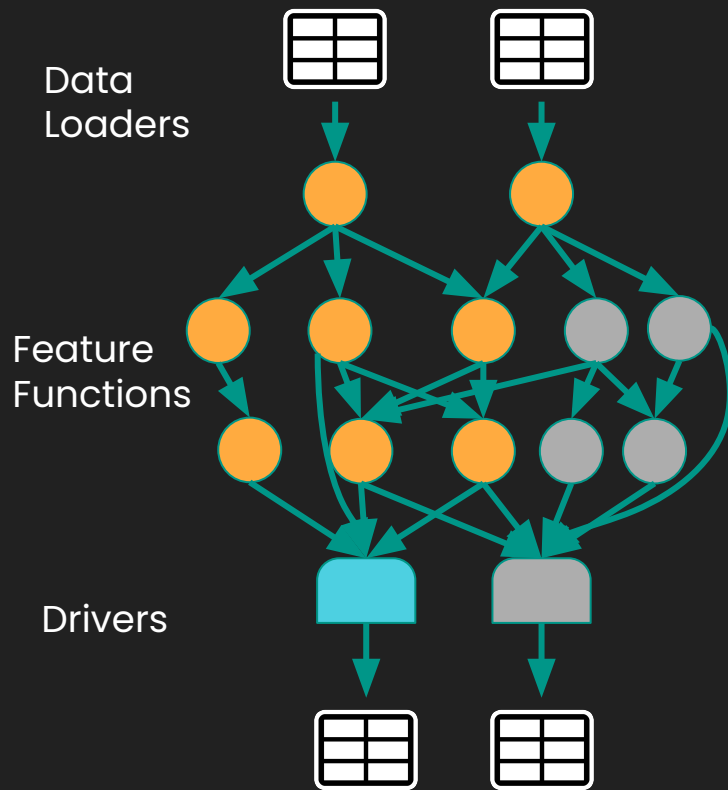


Modeling Feature Engineering

Code that needs to be written:

1. Functions to load data
2. Transform/feature functions
3. Driver to materialize data

Execute only what's needed...



Feature Engineering Challenges

Sustainable code

- ❑ Highly coupled code
- ❑ Difficulty debugging/understanding flows
- ❑ Messy collaboration on complex pipelines
- ❑ Validating your data



Hamilton solves this!

Scaling the data

- ❑ Data is too big to fit in memory
- ❑ Cannot easily parallelize computation



Hamilton has integrations for this!

The Agenda

A motivating story of DS pain

The solution: *Hamilton*

Hamilton for feature engineering

- ↳ **Sustainable feature management**

- ↳ **Scalable feature pipelines**

Hamilton for end-to-end ML workflows

OS progress + next steps

Decoupling Code

From infrastructure

- ❑ Driver handles execution
- ❑ Functions handle business logic

From itself

- ❑ Code organized into functions
- ❑ Functions organized into modules
- ❑ Functions do not know about Hamilton

Module spend_features.py

Module marketing_features.py

Module customer_features.py



Driver script 1

Driver script 2

Driver script 3

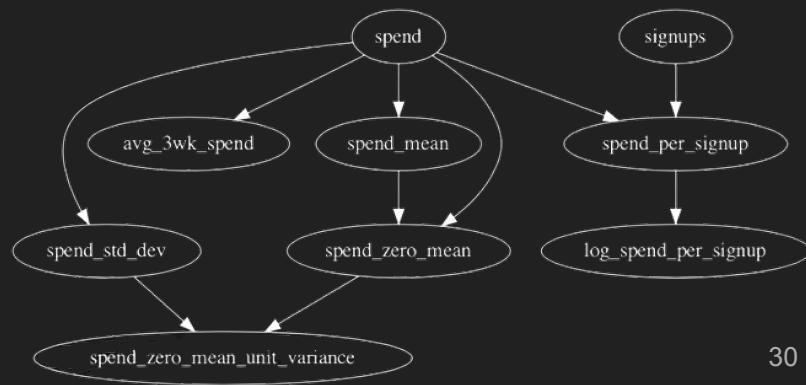
Ease of Debugging

Knocking bugs

- ❑ Rerun just the broken paths
- ❑ Python functions → unit tests + natural debugging
- ❑ Runtime data quality checks
- ❑ Quickly narrow search space of data bugs

Understanding/visualizing

- ❑ Visualize dataflow/execution path
- ❑ Clearly track dependencies



Natural Collaboration

Centralized feature definition store

- ❑ Forces alignment on naming
- ❑ Documentation is included/natural
- ❑ Minimize conflicts when collaborating

Change management

- ❑ Feature versions in git
- ❑ All change \in git history
- ❑ PRs are easy to read – trace changes back to functions

int all files (#677)	3 months ago	53	def first_fix_from_do_choose_manual(54 first_fix_demand_from_do_clients: pd.Series, 55 first_fix_from_do_choose_autoship: pd.Series, 56) -> pd.Series: 57 return first_fix_demand_from_do_clients - first_fix_from_do_choose_autoship 58 59
20220122 model fitting (#580)	10 months ago	57	
int all files (#677)	3 months ago	60	def recovered_clients_choose_manual(61 demand_manual_by_dormant_clients: pd.Series, 62 recovered_clients_choose_autoship: pd.Series, 63) -> pd.Series: 64 """TODO: 65 :param demand_manual_by_dormant_clients: 66 :param recovered_clients_choose_autoship: 67 :return: 68 """ 69 return demand_manual_by_dormant_clients - recovered_clients_choose_autoship 70 71
Renames module 'demand_manual' due to clash wit...	3 years ago	64	
renamed and rerouted demand_from_churned_cle...	2 years ago	65	
Renames module 'demand_manual' due to clash wit...	3 years ago	66	
renamed and rerouted demand_from_churned_cle...	2 years ago	69	
Renames module 'demand_manual' due to clash wit...	3 years ago	70	
Fixes some demand_manual cols	3 years ago	72	# @does(sus_series) need to be able to pass in the ability to specify fill value. 73 74
int all files (#677)	3 months ago	73	def first_time_autoship_clients_for_demand_tf(74 new_signups_choose_autoship_tf: pd.Series, 75 new_signups_choose_autoship_delayed_tf: pd.Series, 76 manual_to_autoship_always_manual_tf: pd.Series, 77 no_demand_to_date_choose_autoship_tf: pd.Series, 78 first_fix_from_do_choose_autoship_tf: pd.Series, 79 D_FUTURE: pd.Series, 80) -> pd.Series: 81 out = (82 new_signups_choose_autoship_tf 83 + new_signups_choose_autoship_delayed_tf 84 + manual_to_autoship_always_manual_tf 85 + no_demand_to_date_choose_autoship_tf 86 + first_fix_from_do_choose_autoship_tf 87) 88

Handling Data Validation

Garbage in/garbage out

- ❑ How can you build reliable pipelines if the data is bad?

Solution

- ❑ Runtime data validation decorator!

```
@check_output(  
    data_type=np.float64, # data type  
    range=(-1.0, 1.0), # range  
    allow_nans=False, # no nans  
    importance="warn") # warn, don't fail  
def some_data_we_care_about() -> pd.Series:  
    return ...
```


Basic Checks

A few custom-built checks for a quick-start:

- ❑ Range
- ❑ Nan-checks (any or percentage)
- ❑ Valid categories
- ❑ Null outputs
- ❑ Plenty more...

For pandas + primitives

Highly pluggable!

Pandera Integration

But wait, there's more! Pandera + Hamilton = happy, powerful checks

```
import pandera as pa
import pandas as pd
from hamilton import function_modifiers

@function_modifiers.check_output(schema=pa.DataFrameSchema(
    {
        'column1': pa.Column(int),
        'column2': pa.Column(float, pa.Check(lambda s: s < -1.2)),
        # you can provide a list of validators
        'column3': pa.Column(str, [
            pa.Check(lambda s: s.str.startswith('value')),
            pa.Check(lambda s: s.str.split('_', expand=True).shape[1] == 2)
        ]),
    },
    index=pa.Index(int),
    strict=True,
))
def dataframe_with_schema(...) -> pd.DataFrame:
    ...
```

Data Check Extensibility

Implement base-class to write your own...

```
@check_output_custom(MyDataValidationClass(...))
```

Goal – add integrations for

- Any type of dataframe/datatype
- Multiple validation frameworks (great expectations, deequ, whylogs...)

Sky's the limit!

```
@check_output_custom(AllPrimeValidator(...))  
def prime_number_generator(number_of_primes_to_generate: int) -> pd.Series:  
    pass
```

Summing up Sustainability with some Code...

```
# client_features.py
@tag(owner='client_recs', pii='False')
@check_output(data_type=np.float64, range=(-5.0, 5.0), allow_nans=False)
def client_height_normalized(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
```

Summing up Sustainability with some Code...

```
# client_features.py
@tag(owner='client_recs', pii='False')
@check_output(data_type=np.float64, range=(-5.0, 5.0), allow_nans=False)
def client_height_normalized(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
```

The feature **client_height_normalized** is owned by the **client recommendations team**, contains **no PII**, is of type **float**, depends on **2 upstream features**, has **no Nan values** and produces values in **[-5.0, 5.0]**. VCS shows its **change over time**. **Grep the codebase** for downstream deps.

EL

Knowing nothing else about this code besides what you have, how would you translate it into english?



```
# client_features.py
@tag(owner='client_recs', pii=False)
@check_output(data_type=np.float64, range=(-5.0, 5.0), allow_nans=False)
def client_height_normalized(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
```



This is a Python code that defines a function called `client_height_normalized`. The purpose of this function is to return the zero mean unit variance value of height, based on two input variables: `height_zero_mean` and `height_std_dev`, which are both pandas Series objects.



Before the function definition, there are two function decorators that are applied to the `client_height_normalized` function. These decorators modify the behavior of the function in some way.

The first decorator, `@tag(owner='client_recs', pii=False)`, adds metadata to the function, indicating that it is owned by a client record system and that it does not handle personally identifiable information (PII).

The second decorator, `@check_output(data_type=np.float64, range=(-5.0, 5.0), allow_nans=False)`, performs input/output validation for the function. It checks that the output of the function is a NumPy float64 data type, falls within the range of -5.0 and 5.0, and does not allow NaN (not-a-number) values.

In the function definition, `height_zero_mean` is divided by `height_std_dev` to calculate the zero mean unit variance value of height. The resulting value is returned as a pandas Series object.

The Agenda

A motivating story of DS pain

The solution: *Hamilton*

Hamilton for feature engineering

- ↳ Sustainable feature management

- ↳ Scalable feature pipelines

Hamilton for end-to-end ML workflows

OS progress + next steps

Scaling Compute/Data

Approach: *delegate*

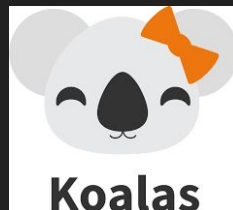
Ray/Dask



- ❑ Run in parallel
- ❑ Single machine → multiprocessing → cluster

Pandas on Spark (ex-koalas)

- ❑ Scale horizontally per dataset



Switching only requires making a **driver-side change***

**pandas on spark requires changing of data loaders as well...*

Scaling with Ray

```
# run.py
from hamilton import driver
import data_loaders
import date_features
import spend_features
config = {...} # config, e.g. data_location
dr = driver.Driver(config,
                   data_loaders,
                   date_features,
                   spend_features)

features_wanted = [...] # choose subset wanted
feature_df = dr.execute(features_wanted)
save(feature_df, 'prod.features')
```

Scaling with Ray

```
# run_on_ray.py
...
from hamilton import base, driver
from hamilton.experimental import h_ray
...
ray.init()
config = {...}
rga = h_ray.RayGraphAdapter(
    result_builder=base.PandasDataFrameResult())
dr = driver.Driver(config,
                    data_loaders, date_features, spend_features,
                    adapter=rga)
features_wanted = [...] # choose subset wanted
feature_df = dr.execute(features_wanted,
                        inputs=date_features)
save(feature_df, 'prod.features')
ray.shutdown()
```

Scaling with Dask

```
# run_on_dask.py
...
from hamilton import base, driver
from hamilton.experimental import h_dask
...
client = Client(Cluster(...)) # dask cluster/client
config = {...}
dga = h_dask.DaskGraphAdapter(client,
    result_builder=base.PandasDataFrameResult())
dr = driver.Driver(config,
    data_loaders, date_features, spend_features,
    adapter=dga)
features_wanted = [...] # choose subset wanted
feature_df = dr.execute(features_wanted,
    inputs=date_features)
save(feature_df, 'prod.features')
client.shutdown()
```

Scaling with Pandas-on-Spark

```
# run_on_pandas_on_spark.py
...
import pyspark.pandas as ps
from hamilton import base, driver
from hamilton.experimental import h_spark
...
spark = SparkSession.builder.getOrCreate()
ps.set_option(...)
config = {...}
skga = h_dask.SparkKoalasGraphAdapter(spark, spine='COLUMN_NAME',
    result_builder=base.PandasDataFrameResult())
dr = driver.Driver(config,
    spark_data_loaders, date_features, spend_features,
    adapter=skga)
features_wanted = [...] # choose subset wanted
feature_df = dr.execute(features_wanted,
    inputs=date_features)
save(feature_df, 'prod.features')
spark.stop()
```

Hamilton + Ray/Dask: How Does it Work?

FUNCTIONS

```
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
```

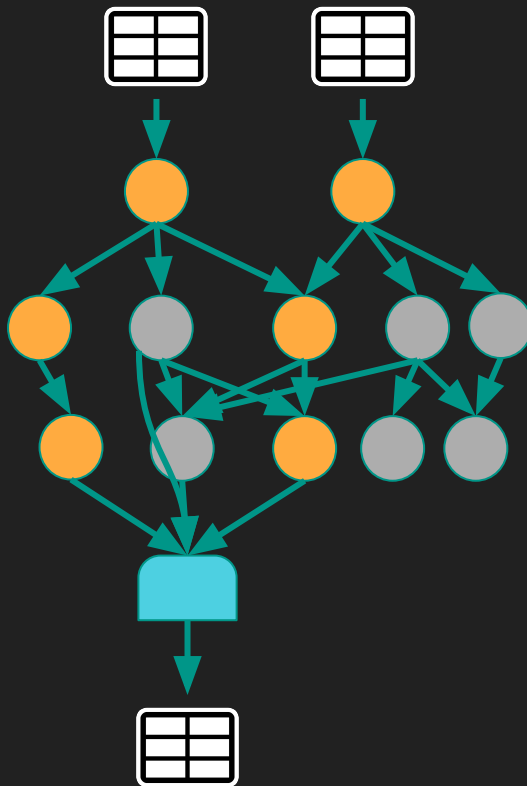
DRIVER

```
from hamilton import base, driver
from hamilton.experimental import h_ray

ray.init()
config = {...}
rga = h_ray.RayGraphAdapter(
    result_builder=base.PandasDataFrameResult()
)
dr = driver.Driver(config,
    data loaders,
    date features,
    spend features,
    adapter=rga)

features_wanted = [...] # choose subset wanted
feature_df = dr.execute(features_wanted,
    inputs=date_features)
save(feature_df, 'prod.features')
ray.shutdown()
```

DAG



Hamilton + Ray/Dask: How Does it Work?

FUNCTIONS

```
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
```

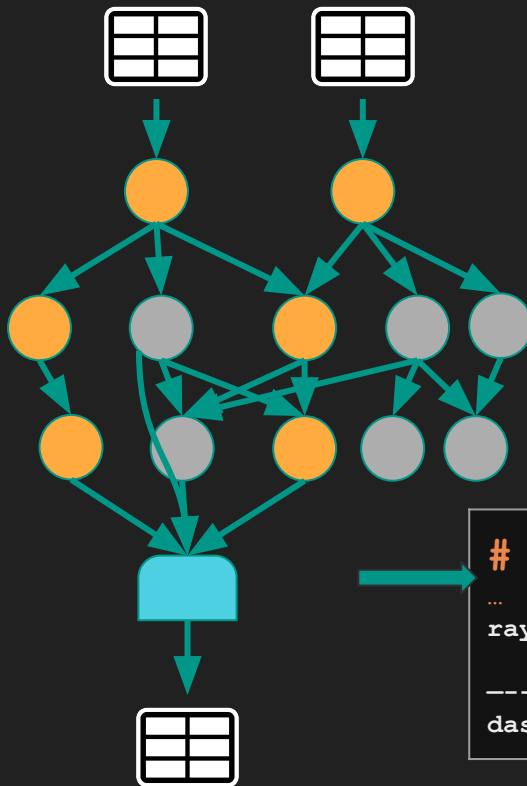
DRIVER

```
from hamilton import base, driver
from hamilton.experimental import h_ray

ray.init()
config = {...}
rga = h_ray.RayGraphAdapter(
    result_builder=base.PandasDataFrameResult()
)
dr = driver.Driver(config,
    data loaders,
    date features,
    spend features,
    adapter=rga)

features_wanted = [...] # choose subset wanted
feature_df = dr.execute(features_wanted,
    inputs=date_features)
save(feature_df, 'prod.features')
ray.shutdown()
```

DAG



Delegate to Ray/Dask

```
...
ray.remote(
    node.callable).remote(**kwargs)
-----
dask.delayed(node.callable)(**kwargs)
```

Hamilton + Spark: How Does it Work?

FUNCTIONS

```
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
```

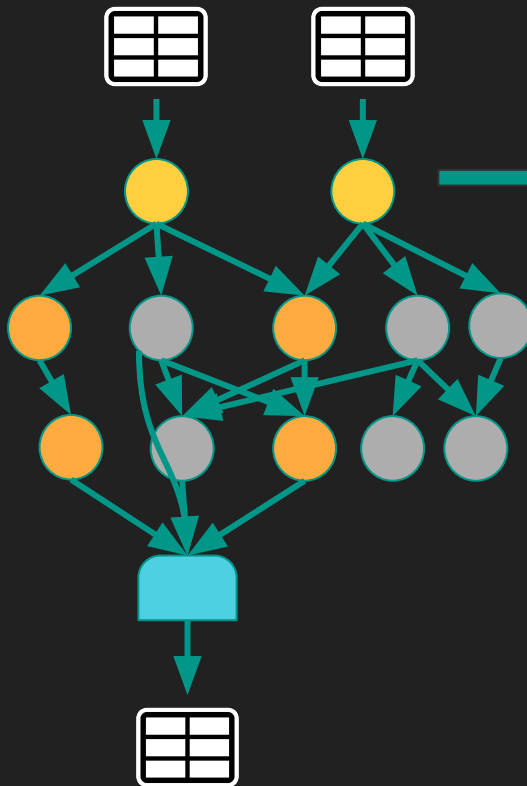
DRIVER

```
from hamilton import base, driver
from hamilton.experimental import h_ray

ray.init()
config = {...}
rga = h_ray.RayGraphAdapter(
    result_builder=base.PandasDataFrameResult()
)
dr = driver.Driver(config,
    data loaders,
    date features,
    spend features,
    adapter=rga)

features_wanted = [...] # choose subset wanted
feature_df = dr.execute(features_wanted,
    inputs=date_features)
save(feature_df, 'prod.features')
ray.shutdown()
```

DAG



With Spark

...
Change these to load
Spark "Pandas"
equivalent object
instead.

Spark will take care
of the rest.

Hamilton + Ray/Dask/Pandas on Spark: Caveats

Serialization

- ❑ Uses serialization methodology of delegated frameworks

Memory:

- ❑ Defaults should work – fine tuning at fn level not yet supported

Python dependencies:

- ❑ You need to manage them

Looking to graduate these APIs from experimental status

◆◆ Contributions wanted here to extend support in Hamilton! ◆◆

The Agenda

A motivating story of DS pain

The solution: *Hamilton*

Hamilton for Feature Engineering

- ↳ Sustainable feature management
- ↳ Scalable feature pipelines

Hamilton for end-to-end ML workflows

OS progress + next steps

End-to-End ML Pipelines

What does an ML pipeline look like? Fancy ETL:

- ❑ **[E]** Load data from a feature store
- ❑ **[T]** Transform features
- ❑ **[T]** Train model
- ❑ **[T]** Run model Inference
- ❑ **[T]** Evaluate model performance
- ❑ **[L]** Save metrics
- ❑ **[L]** Save artifacts
- ❑ **[L]** Save training data

End-to-End ML Pipelines

Centralize logic, abstract integrations

- ❑ [E] Load data from a feature store
- ❑ [T] Transform features
- ❑ [T] Train model
- ❑ [T] Run model Inference
- ❑ [T] Evaluate model performance
- ❑ [L] Save metrics
- ❑ [L] Save artifacts
- ❑ [L] Save training data

}

Express as functions

}

Decouple from logic

End-to-End ML Pipelines

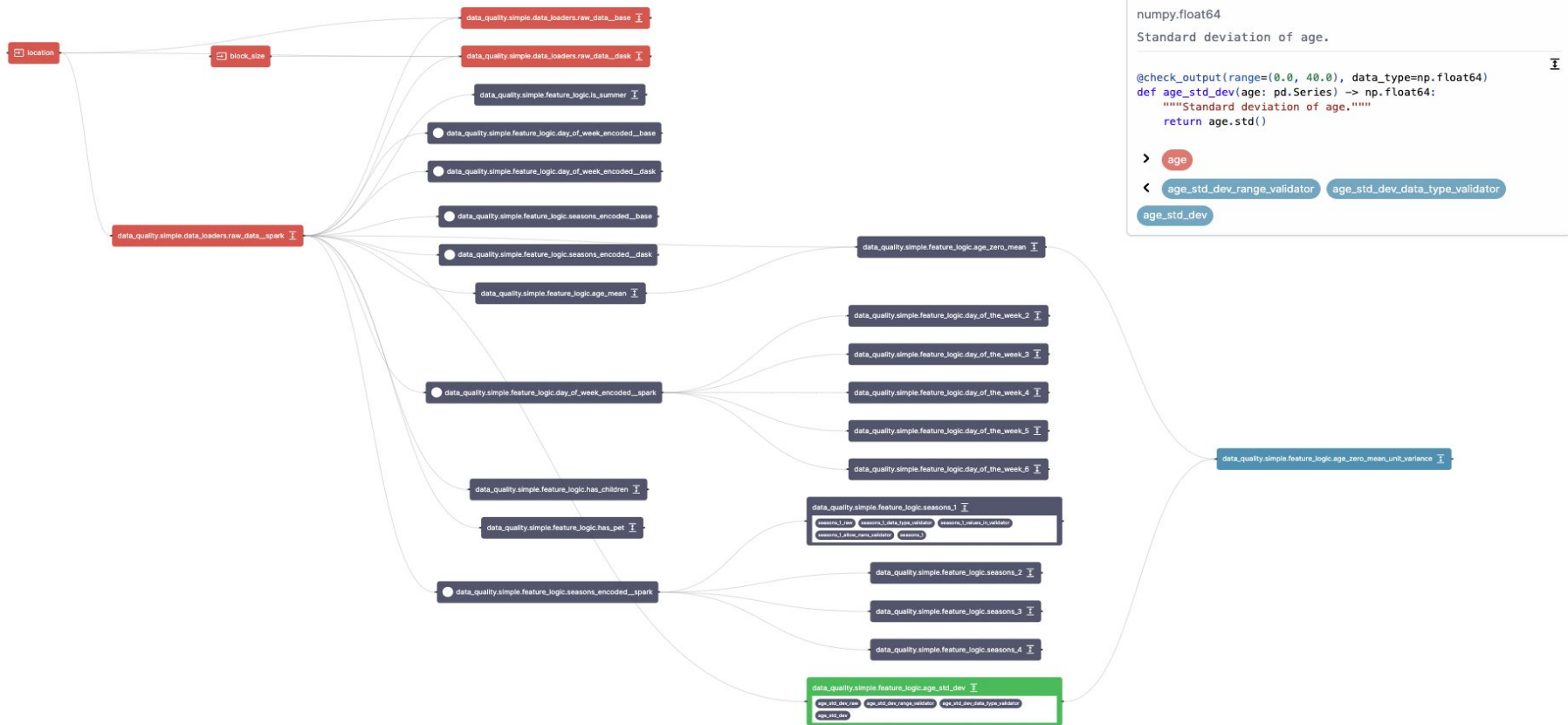
Decouple from logic

- ❑ Swap out
 - ❑ Model store
 - ❑ Metrics store
 - ❑ Feature source/store
- ❑ Execute same code online/offline
 - ❑ Map ops 1:1
 - ❑ Join -> data load using DAG config
 - ❑ aggregation -> data load/hardcoded

End-to-End ML Pipelines

Gain Visibility

- ❑ Lineage
 - ❑ Trace fn params for managing dependencies
 - ❑ Trace data from source -> transforms -> sinks
 - ❑ Tag metadata-> understand properties of dependencies
- ❑ Catalogue
 - ❑ Browse features == browse through code
 - ❑ Documentation attached to artifact name



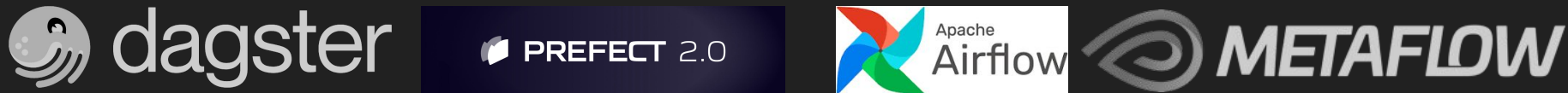
- ☐ By node No grouping
- ☒ By function Group nodes by the function in which they were defined
- ☐ By namespace Group nodes by their namespace(s) (useful for subdags)
- ☐ By module Group nodes by module



Our Vision

Unifying layer for ML ETLs

Compile to orchestration frameworks



Integrate with data quality vendors/OS options



Trivially load from/to a variety of sources/destinations



SQL support for full ETL

The Agenda

A motivating story of DS pain

The solution: *Hamilton*

Hamilton for feature engineering

- ↳ Sustainable feature management
- ↳ Scalable feature engineering

Hamilton to model end-to-end ML workflows

OS progress + next steps

OS Progress

Early stages, but thriving community

- ❑ Some exciting users
- ❑ Growing set of core contributors
- ❑ Full company dedicated to building it!

Looking for

- ❑ Contributors
- ❑ Bug hunters
- ❑ User feedback

STITCH FIX



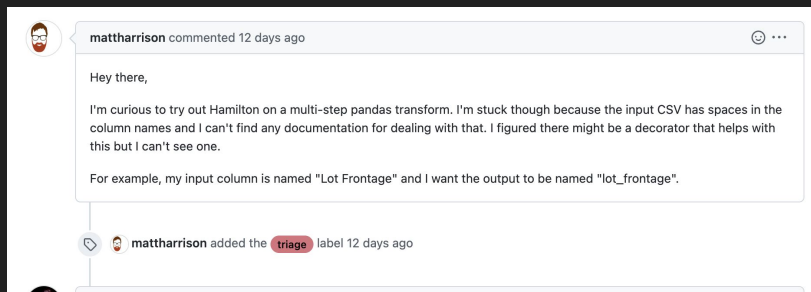
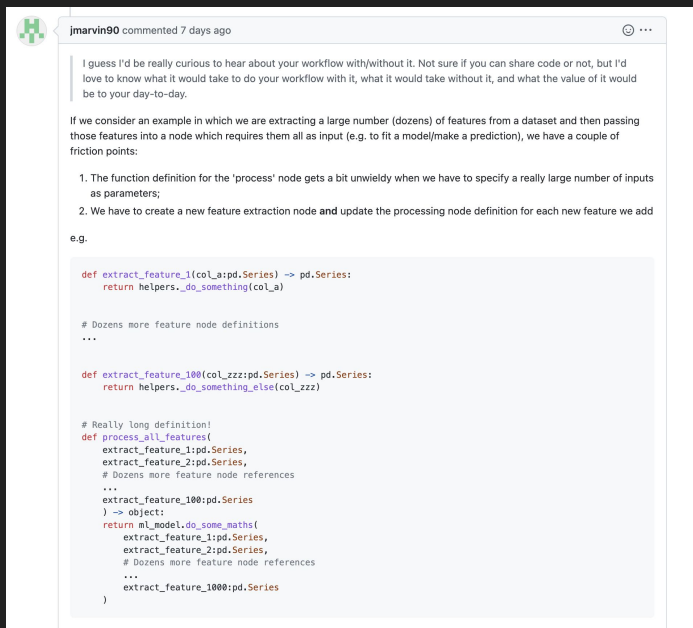
In Progress

Expressive APIs

- ❑ Schema for artifact metadata
- ❑ Adapter for SQL
- ❑ <Your idea here!>

Execution

- ❑ Compilation/dataflow specification
- ❑ Streaming/generator support
- ❑ First-class pyspark integration
- ❑ <Your idea here!>



Give Hamilton a Try! We'd Love Your Feedback.

www.tryhamilton.dev

> `pip install sf-hamilton`

★ on [github](https://github.com/dagworks-inc/hamilton) (<https://github.com/dagworks-inc/hamilton>)

✓ create & vote on issues on github

🔔 join us on [Slack](https://join.slack.com/t/hamilton-opensource/shared_invite/zt-1bjs72asx-wcUTgH7q7QXlgiQ5bbdcg)
(https://join.slack.com/t/hamilton-opensource/shared_invite/zt-1bjs72asx-wcUTgH7q7QXlgiQ5bbdcg)

Thank you.

Questions?

Yell at me online <https://twitter.com/elijahbenizzy>

Connect with me <https://www.linkedin.com/in/elijahbenizzy/>

Code with me <https://github.com/dagworks-inc/hamilton>

Use sparingly :) elijah@dagworks.io

