



Hamilton Global User Group February 2024 Meetup

What is Hamilton?

Hamilton helps data scientists and engineers define testable, modular, self-documenting dataflows, that encode lineage and metadata.
Runs and scales everywhere python does.

Icebreaker: Name and what you're using Hamilton for/looking for.



Agenda

1. Community Spotlight
2. Deep Dive
3. Roadmap
4. Open 



Community Spotlight:

 **"How we migrated our feature calculation to Hamilton"** by Arthur Andres.

Deep Dive: Mental Models for Structuring Projects

My Mental Model for Structuring Projects

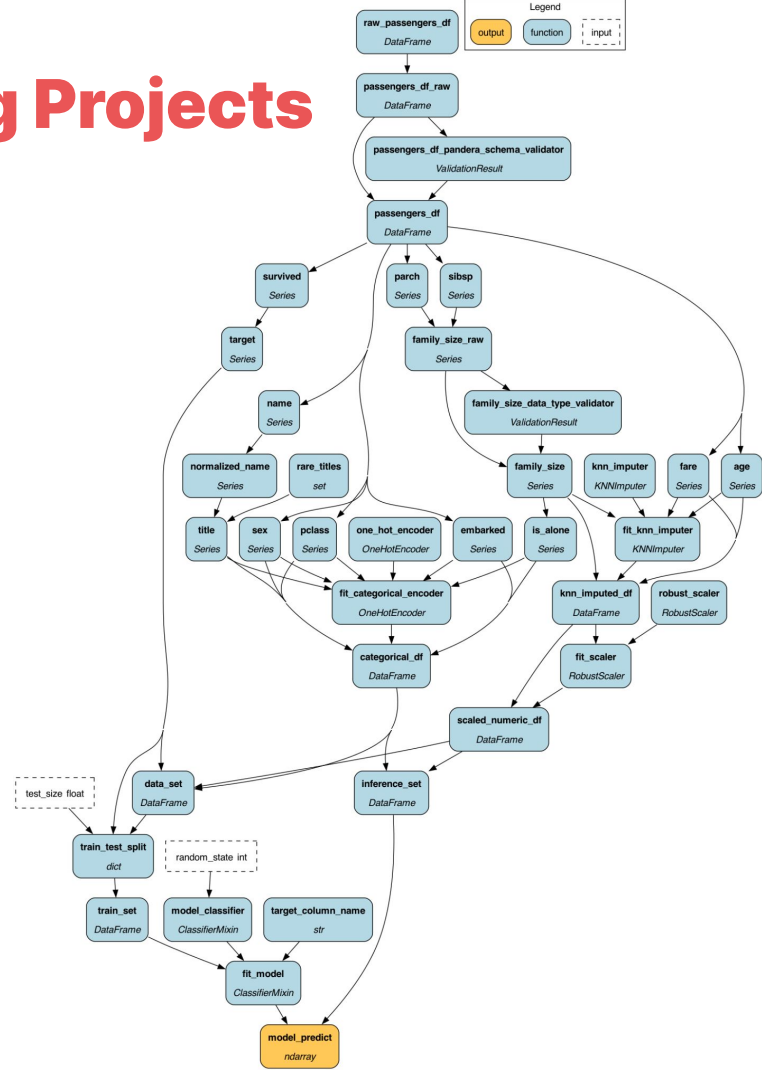
With Hamilton you:

1. Write functions
2. Functions are organized into modules.

My Mental Model for Structuring Projects

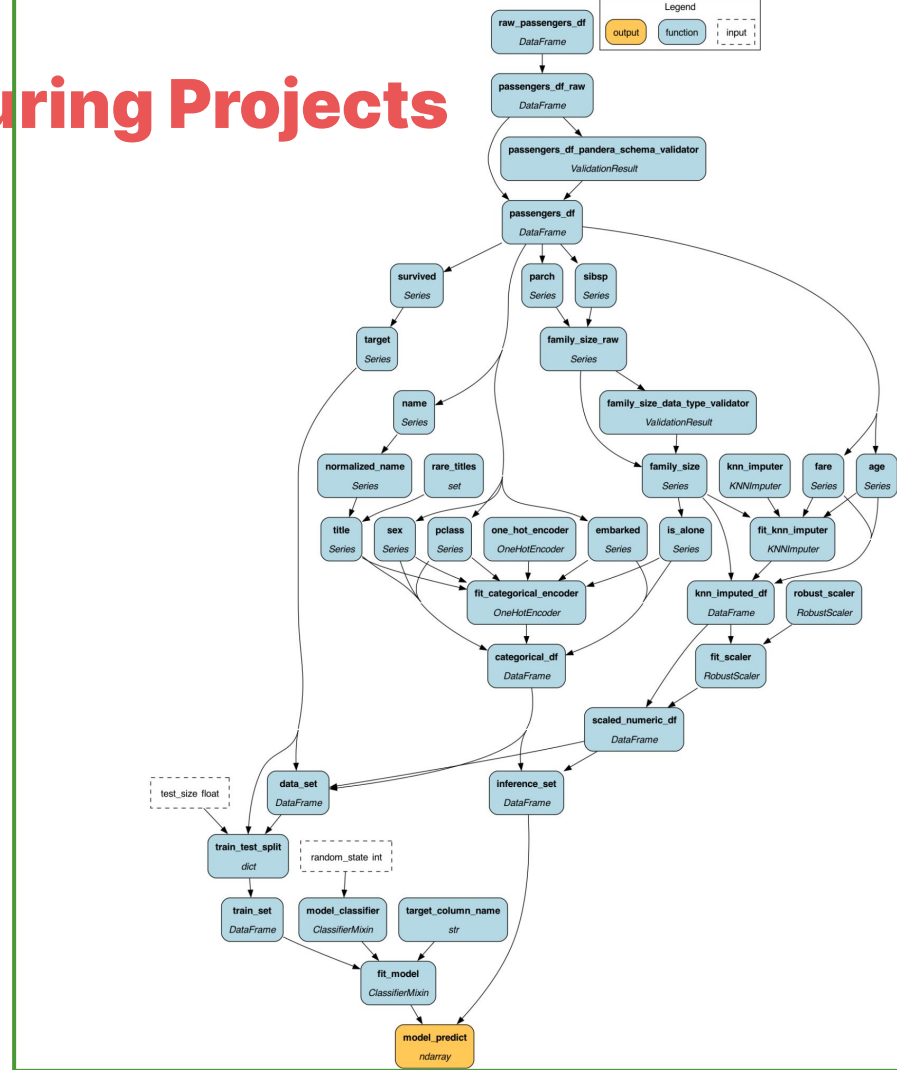
With Hamilton you:

1. Write functions
2. Functions are organized into modules.



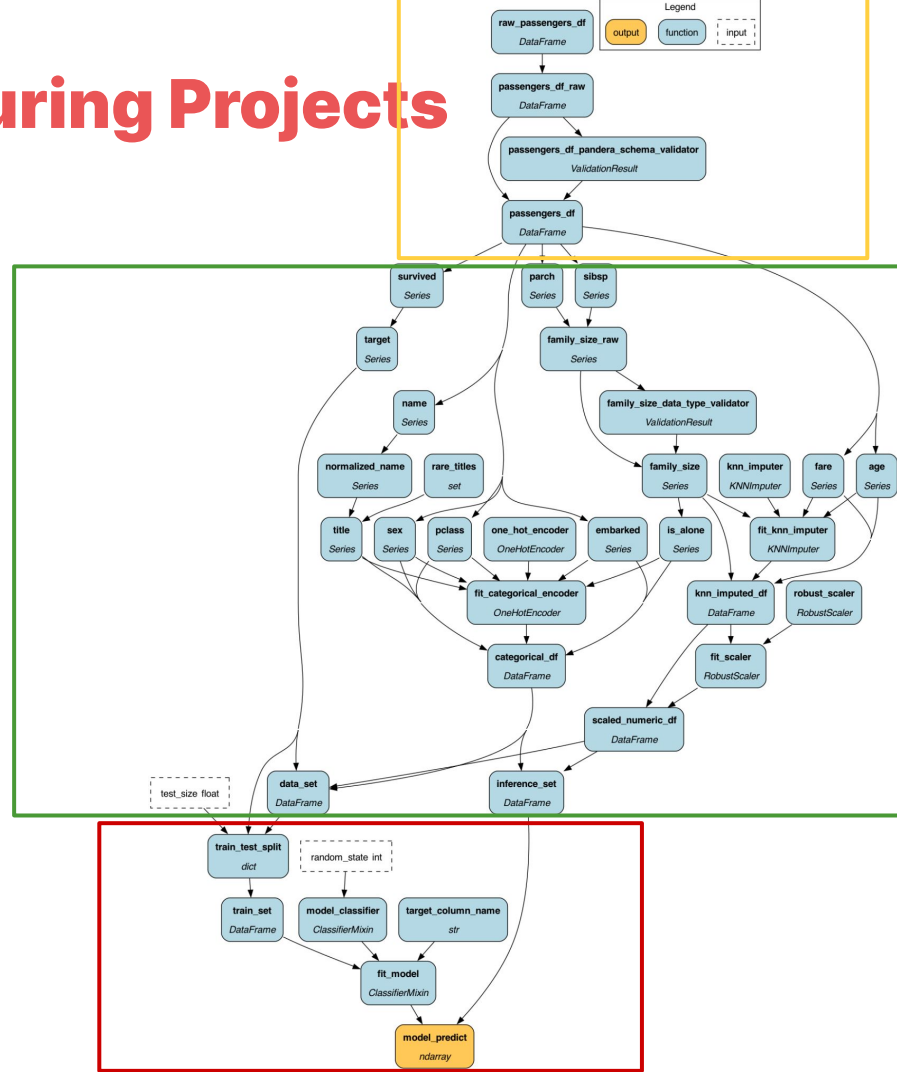
My Mental Model for Structuring Projects

Could be one module...



My Mental Model for Structuring Projects

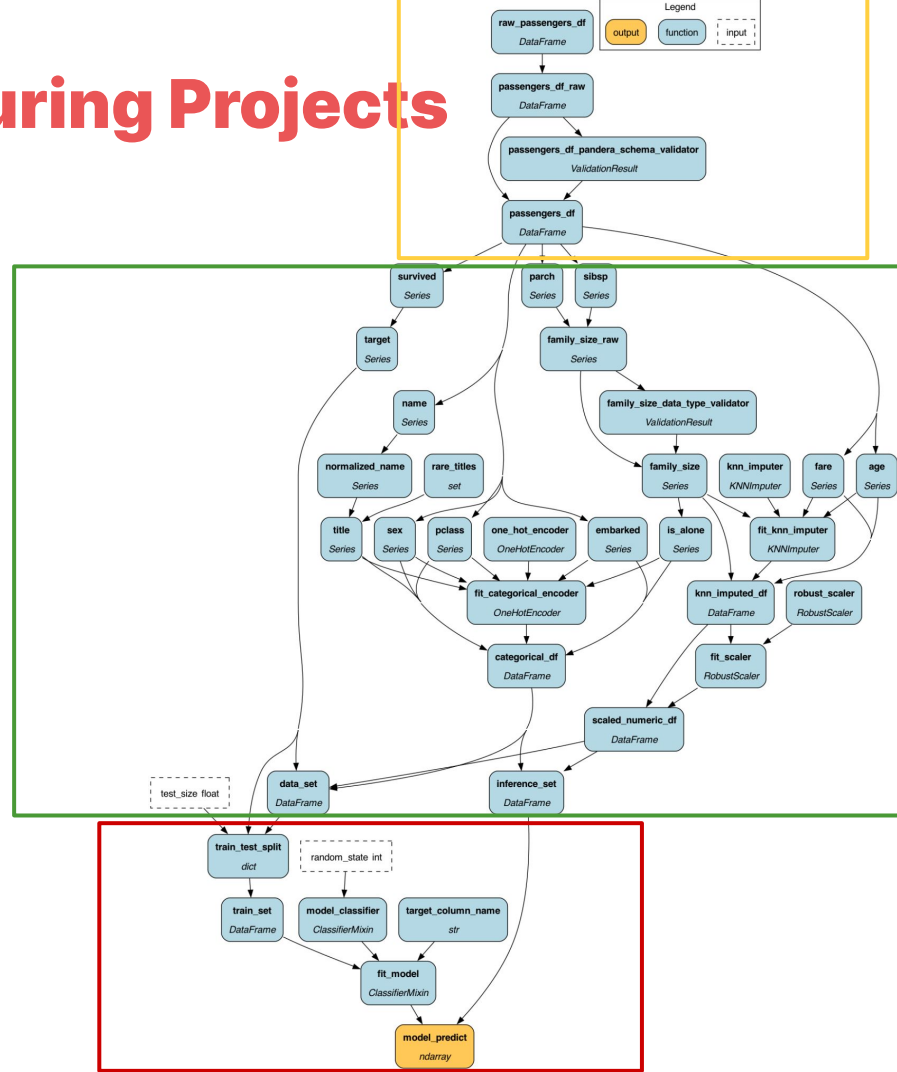
Or three... or more...



My Mental Model for Structuring Projects

```
import data_loader, feature_transforms, model_pipeline

# DAG for training/infering on titanic data
titanic_dag = (
    driver.Builder()
        .with_config(config),
        .with_modules(
            data_loader, feature_transforms, model_pipeline
        ).build()
)
# execute & get output
result = titanic_dag.execute(["model_predict"])
```



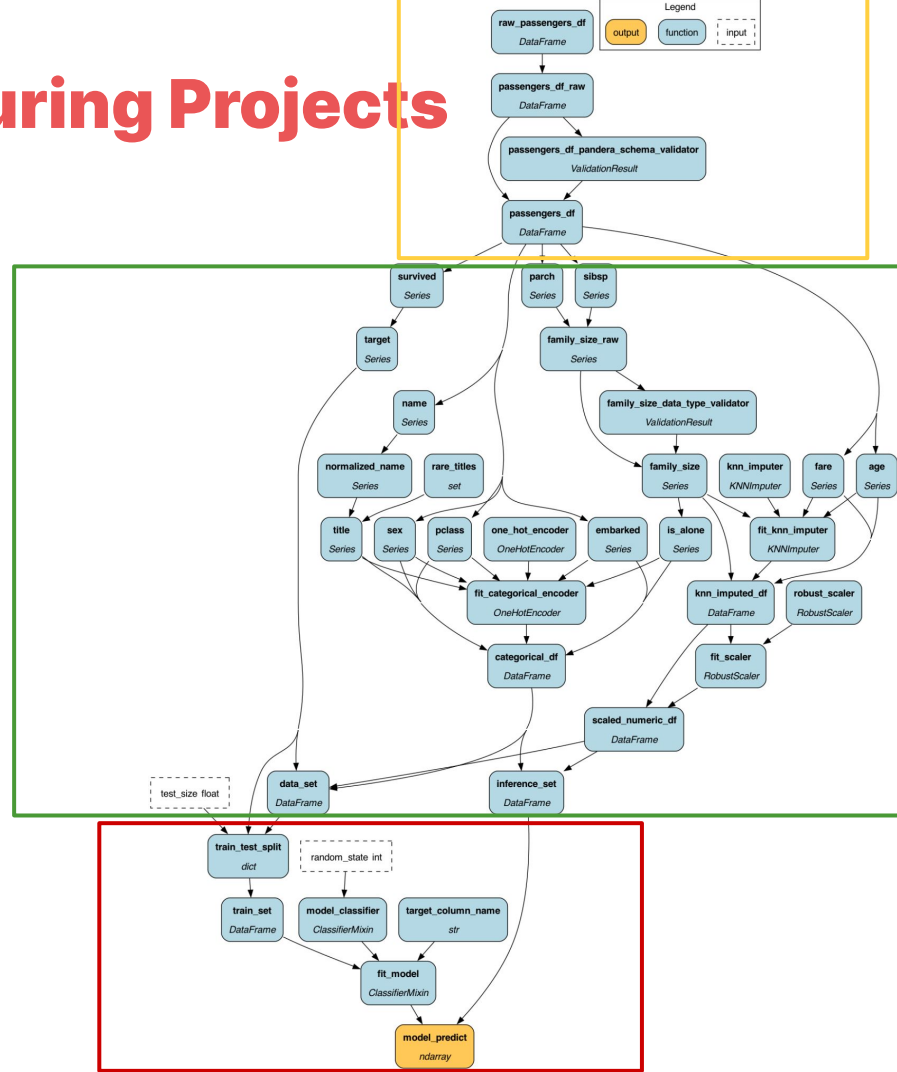
My Mental Model for Structuring Projects

```
import data_loader, feature_transforms, model_pipeline

# DAG for training/infering on titanic data
titanic_dag = (
    driver.Builder()
        .with_config(config),
        .with_modules(
            data_loader, feature_transforms, model_pipeline
        ).build()
)
```

Scaling this:

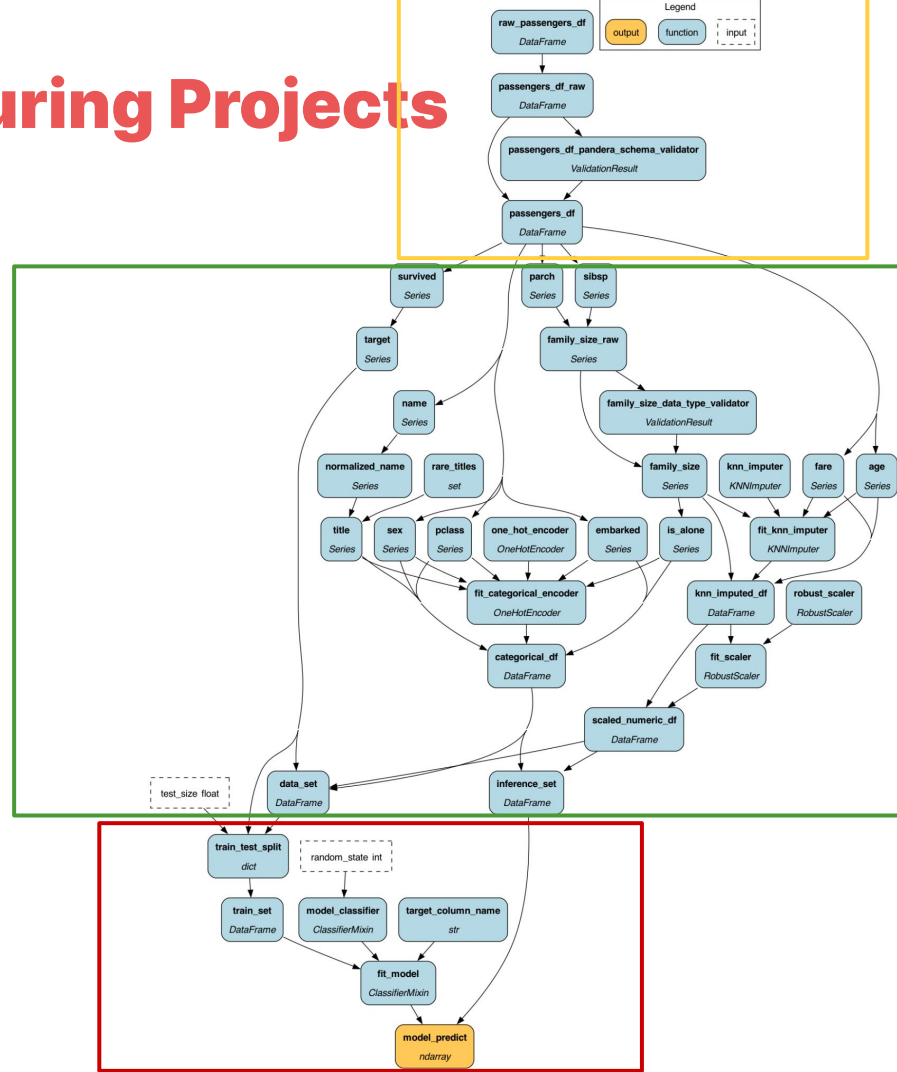
1. Group into thematic modules.
2. Name + type becomes interface.
3. Modularity comes from swapping modules and/or @config.when



My Mental Model for Structuring Projects

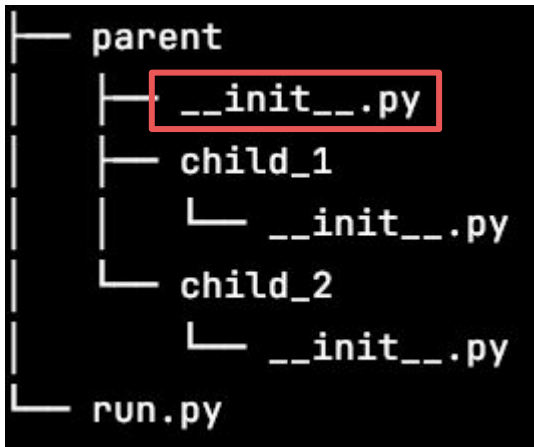
Summary:

- Modules == Lego bricks
- Lego bricks are standalone
 - But can be composed into something bigger
- Caveat:
 - Naming



Managing 4000 feature transforms at Stitch Fix

1. Single repository.
2. Naming convention for functions.
 - a. E.g. `D_` is dummy variable.
3. Hierarchically group into subpackages.



Use it to aggregate

```
from .child_1 import * #parent/__init__.py
from .child_2 import *

def func_e(func_d: int, func_b: int) -> int:
    return func_d + func_b
```

```
from hamilton import driver
import parent

if __name__ == '__main__':
    dr = driver.Builder().with_modules(parent).build()
    dr.display_all_functions("dag.png")
```

Managing 4000 feature transforms at Stitch Fix

1. Single repository.
2. Naming convention for functions.
 - a. E.g. `D_` is dummy variable.
3. Hierarchically group into subpackages.
4. At driver construction time, pick right subset of packages
 - a. Can use `importlib` to script:

```
from hamilton import driver

import importlib

if __name__ == '__main__':
    modules = [
        importlib.import_module(module_name)
        for module_name in [...]
    ]
    dr = driver.Builder().with_modules(*modules).build()
```

Managing 4000 feature transforms at Stitch Fix

1. Single repository.
2. Naming convention for functions.
 - a. E.g. `D_` is dummy variable.
3. Hierarchically group into subpackages.
4. At driver construction time, pick right subset of packages
 - a. Can use `importlib` to script
5. Can change overtime easily:
 - a. Renaming is just find & replace (largely)
 - b. Move functions around into other modules

Deep Dive:
Explaining **@subdag** in two minutes

@subdag in two minutes

```
def feature_engineering(source_path: str) -> pd.DataFrame:  
    """You could recursively use Hamilton within itself."""  
    dr = driver.Driver({}, feature_modules)  
    df = dr.execute(["feature_df"], inputs={"path": source_path})  
    return df
```

← Can use Hamilton within Hamilton

← But lose “visibility”

```
@subdag(  
    feature_modules,  
    inputs={"path": source("source_path")},  
    config={}  
)  
def feature_engineering(feature_df: pd.DataFrame) -> pd.DataFrame:  
    """We've hidden the driver and requested feature_df from it."""  
    return feature_df
```

← Make it visible in a single DAG

← Subdag is “name spaced”.

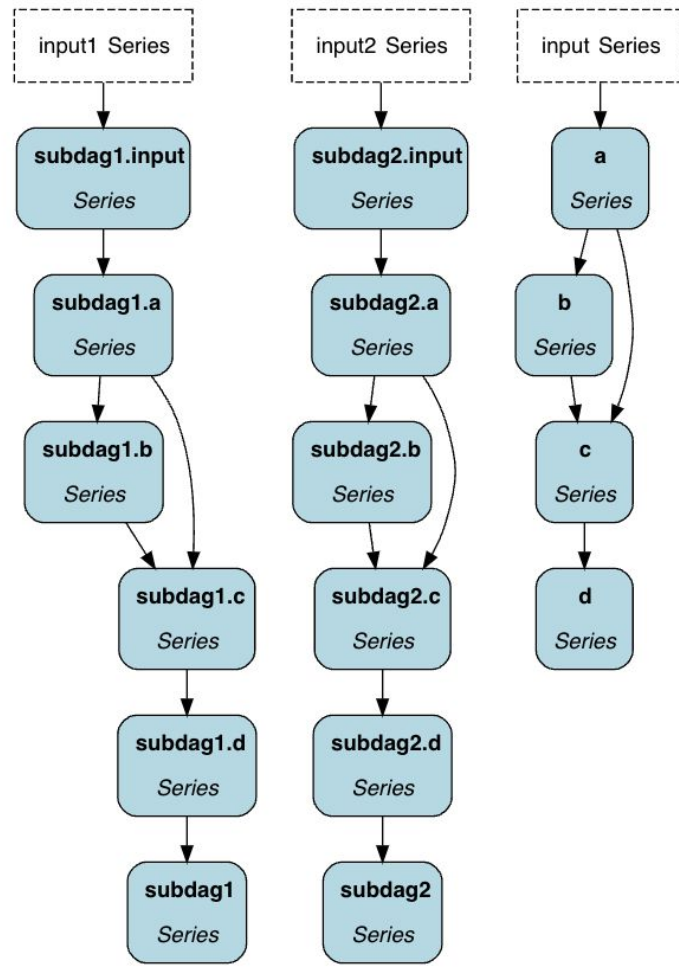
@subdag in two minutes

```
1 # functions.py - declare and link your transformations as
  functions...
2 import pandas as pd
3 from hamilton.function_modifiers import subdag, source
4
5 def a(input: pd.Series) -> pd.Series:
6     return input % 7
7
8 def b(a: pd.Series) -> pd.Series:
9     return a * 2
10
11 def c(a: pd.Series, b: pd.Series) -> pd.Series:
12     return a * 3 + b * 2
13
14 def d(c: pd.Series) -> pd.Series:
15     return c ** 3
16
17 @subdag(
18     a, b, c, d,
19     inputs={"input": source("input1")},
20     config={}
21 )
22 def subdag1(d: pd.Series) -> pd.Series:
23     return d * 2
24
25 @subdag(
26     a, b, c, d,
27     inputs={"input": source("input2")},
28     config={}
29 )
30 def subdag2(d: pd.Series) -> pd.Series:
31     return d
```

```
1 # And run them!
2 import functions
3 from hamilton import driver
4 dr = driver.Driver({}, functions)
5 result = dr.execute(
6     ['subdag2', 'subdag1'],
7     inputs={"input2": pd.Series([1, 2, 3, 4, 5]),
8           'input1': pd.Series([1, 2, 3, 4, 5])}
9 )
10 print(result)
11 dr.display_all_functions(
12     "graph.dot", orient="TB", show_legend=False)
```

```
subdag2 subdag1
```

▶ Run me!





Roadmap Thoughts

Roadmap: Some Recent Releases

- @tag now supports list of strings
- Custom styling for graphviz
- JSON export

Roadmap: Some Recent Releases

- @tag now supports list of strings
- Custom styling for graphviz
- JSON export
- Lightweight experiment tracker →

Hamilton Experiment Manager

Select Experiment ...

Select Graph version ...

Date Completed	Experiment	Graph Version	Run Id
2024-02-05T10:55:01	hello-world	1	d3a24c23-f39f-4f05-be7e-0cc31838b05a
2024-02-05T16:17:54	hello-world	2	17f99878-0261-4fe9-b274-68c6a8e2bci4
2024-02-05T16:18:51	hello-world	2	f97542f5-b0db-455e-9f82-fc5a9d8a0bc2
2024-02-05T16:19:35	hello-world	3	74115161-7486-4a3b-9b12-5bb39e3a8a8f
2024-02-05T16:20:17	forecast	1	244ca9e7-57ff-4909-ac86-c8d1dfb86248
2024-02-05T16:20:33	forecast	2	56c3d40e-5b32-4ece-87c1-c1a3fd2ae9c4

Roadmap: Some Recent Releases

- @tag now supports list of strings
- Custom styling for graphviz
- JSON export
- Lightweight experiment tracker →
- <https://hub.dagworks.io/> additions
- [Documentation refactoring](#)
- Lifecycle API & additions, e.g. TQDM, Datadog

Hamilton Experiment Manager

Select Experiment ...

Select Graph version ...

Date Completed	Experiment	Graph Version	Run Id
2024-02-05T10:55:01	hello-world	1	d3a24c23-f39f-4f05-be7e-0cc31838b05a
2024-02-05T16:17:54	hello-world	2	17f99878-0261-4fe9-b274-68c6a8e2bci4
2024-02-05T16:18:51	hello-world	2	f97542f5-b0db-455e-9f82-fc5a9d8a0bc2
2024-02-05T16:19:35	hello-world	3	74115161-7486-4a3b-9b12-5bb39e3a8a8f
2024-02-05T16:20:17	forecast	1	244ca9e7-57ff-4909-ac86-c8d1dfb86248
2024-02-05T16:20:33	forecast	2	56c3d40e-5b32-4ece-87c1-c1a3fd2ae9c4

```
to use available CPU instructions in performance-critical operations.  
To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler  
flags.
```

Roadmap: Some Recent Releases

- @tag now supports list of strings
- Custom styling for graphviz
- JSON export
- Lightweight experiment tracker →
- <https://hub.dagworks.io/> additions
- [Documentation refactoring](#)
- Lifecycle API & additions, e.g. TQDM, Datadog

Hamilton Experiment Manager

Select Experiment ...

Select Graph version ...

Date Completed	Experiment	Graph Version	Run Id
2024-02-05T10:55:01	hello-world	1	d3a24c23-f39f-4f05-be7e-0cc31838b05a
2024-02-05T16:17:54	hello-world	2	17f99878-0261-4fe9-b274-68c6a8e2bci4
2024-02-05T16:18:51	hello-world	2	f97542f5-b0db-455e-9f82-fc5a9d8a0bc2
2024-02-05T16:19:35	hello-world	3	74115161-7486-4a3b-9b12-5bb39e3a8a8f
2024-02-05T16:20:17	forecast	1	244ca9e7-57ff-4909-ac86-c8d1dfb86248
2024-02-05T16:20:33	forecast	2	56c3d40e-5b32-4ece-87c1-c1a3fd2ae9c4

```
to use available CPU instructions in performance-critical operations.  
To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler  
flags.
```

- Caching: `.with_adapters(h_diskcache.CacheHook())`
- Blogs - e.g. dev to prod ML pipelines.

Roadmap: DAGWorks callout

1. We're working with design partners to more tightly build out integrations that shapes open source.
 - a. Ping us/put your hand up and we can do a 1-1 session if interested.
2. FYI: www.dagworks.io has a free tier we'd love feedback on:
 - a. One line addition to Hamilton code.
 - b. Get catalog, telemetry, observability, version capture...

Roadmap: What's on the horizon

- Documentation - code comparisons, integration guides...
- hub.dagworks.io - more examples

Needs prioritization - some ideas:

- **Your input/ideas here!**
 - **Do you have pains? Monitoring? What are you building/integrating with? Etc.**
 - **Testimonials page - we'd love some blog posts/quotes for social proof.**
- Ray/Dask: node grouping
- Snowpark integration
- Data quality integrations/pluggability
- Unit test generator
- Generating Airflow DAGs, AWS lambda integrations, etc.



Next month - March 19th:

Roel Bertens

“How to use Hamilton to share (feature) logic across multiple development teams.”



Open Mic.

FIN. Thanks for coming!

