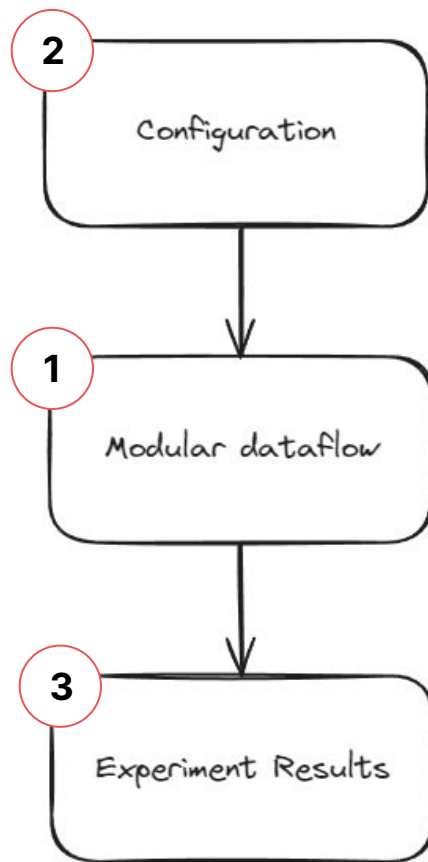


Modular Dataflows & Experiment Management for Machine Learning Evaluation

Modular dataflows for research

Context

- Train ML models for time series forecasting
- Benchmark 200+ models and conduct statistical tests
- Reproducibility is paramount, requiring lineage for code and data



Who Am I ?

- My name is Thierry, I'm based in Montréal, QC, Canada
- Using Hamilton since ~2021-11, it launched in 2021-10
- Working with Stefan and Elijah since 2023-06
- Previous experiences:
 - teaching ML/DS, AI consulting, HR tech SaaS

Creating

Modular Dataflows

Level 0: spaghetti

```
1 import numpy as np
2 import pandas as pd
3
4 def main():
5     df = pd.read_parquet("raw_data.parquet")
6
7     pairwise_diff = np.diff(df[["loc_x", "loc_y"]].values, axis=0)
8     pairwise_dists = np.sqrt((pairwise_diff ** 2).sum(axis=1))
9     df["location_jump"] = np.insert(pairwise_dists, 0, np.nan)
10    timedelta = df.timestamp.diff().dt.total_seconds()
11    df["location_jump_speed"] = df["location_jump"] / timedelta
12
13    return df
```

- A single function

Limitations

- Hard-coded values
- Assumptions about columns
- Unrelated steps within same scope

Level 1: functions

```
1 import numpy as np
2 import pandas as pd
3
4 def location_jump(df):
5     pairwise_diff = np.diff(df[["loc_x", "loc_y"]].values, axis=0)
6     pairwise_dists = np.sqrt((pairwise_diff ** 2).sum(axis=1))
7     return np.insert(pairwise_dists, 0, np.nan)
8
9 def location_jump_speed(df):
10    timedelta = df.timestamp.diff().dt.total_seconds()
11    return df["location_jump"] / timedelta
12
13 def main():
14    df = pd.read_parquet("raw_data.parquet")
15    df["location_jump"] = location_jump(df)
16    df["location_jump_speed"] = location_jump_speed(df)
17    return df
```

- Functions to define features
- Passing DataFrames to add columns

Limitations

- Implicit dependencies between features
- **Imperative**: need to manually order functions in **main()**

Level 2: checking schemas

```
1 import numpy as np
2 import pandas as pd
3
4 @check_columns(["loc_x", "loc_y"])
5 def location_jump(df):
6     pairwise_diff = np.diff(df[["loc_x", "loc_y"]].values, axis=0)
7     pairwise_dists = np.sqrt((pairwise_diff ** 2).sum(axis=1))
8     return np.insert(pairwise_dists, 0, np.nan)
9
10 @check_columns(["location_jump"])
11 def location_jump_speed(df):
12     timedelta = df.timestamp.diff().dt.total_seconds()
13     return df["location_jump"] / timedelta
14
15 def main():
16     df = pd.read_parquet("raw_data.parquet")
17     df["location_jump"] = location_jump(df)
18     df["location_jump_speed"] = location_jump_speed(df)
19     return df
```

- Make dependencies explicit via custom decorator.
- Exactly what **pandera** does with **@check_input**

Limitations

- Still imperative

Level 3: Hamilton

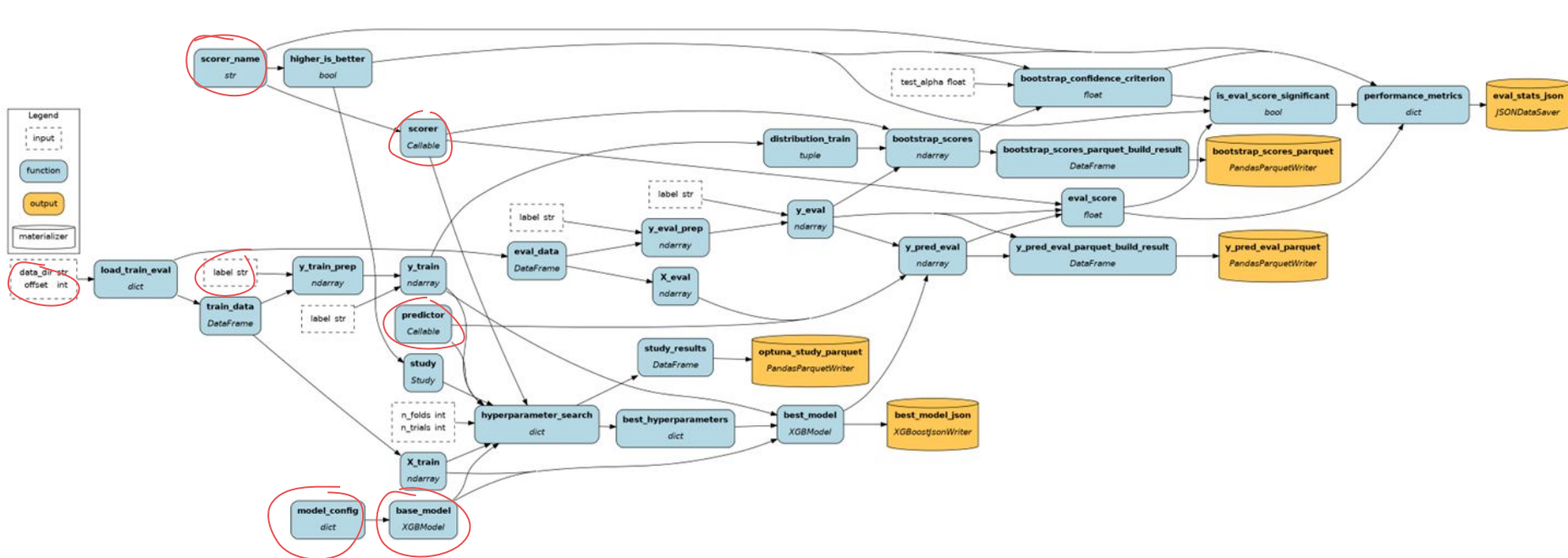
```
1 # features.py
2 import numpy as np
3 import pandas as pd
4 from hamilton.function_modifiers import extract_columns
5
6 @extract_columns("loc_x", "loc_y", "timestamp")
7 def raw_data(file_path: str) → pd.DataFrame:
8     return pd.read_parquet(file_path)
9
10 def location_jump(loc_x: pd.Series, loc_y: pd.Series) → pd.Series:
11     """Distance between the positions at two consecutive timestamps"""
12     return ...
13
14 def location_jump_speed(
15     location_jump: pd.Series,
16     timestamp: pd.Series
17 ) → pd.Series:
18     """Speed between consecutive locations"""
19     return ...
20
21 # run.py
22 from hamilton import driver, base
23 import features
24
25 def main():
26     dr = (
27         driver.Builder()
28         .with_modules(features)
29         .with_adapters(base.PandasDataFrameResult())
30         .build()
31     )
32     features = ["timestamp", "location_jump", "location_jump_speed"]
33     df = dr.execute(
34         features,
35         inputs={"file_path": "raw_data.parquet"}
36     )
```

- Forces the use of type hints and encourages docstrings
- Dependencies are explicit
- Separation between inputs and transformations
- **Declarative** script! Simply request what you want

Adding

Configurability

Training an XGBoost model



Configuration dimensions

Experiment Dimension	Execution `inputs`	Driver `config`	Driver `modules`
3 forecast horizon (one dataset each)	set `offset` to compute (or select) the dataset		
2 model architecture (XGBoost, LSTM)			one dataflow (.py file) per model type
4 learning tasks (binary, multiclass, ordinal, continuous)		set `task` config to one of the 4 tasks	
10 target variables	set `label` to select a column from the dataset		
2 execution mode (development, evaluation)		set `mode` to change dataset loading behavior	

Configuration at the Driver level

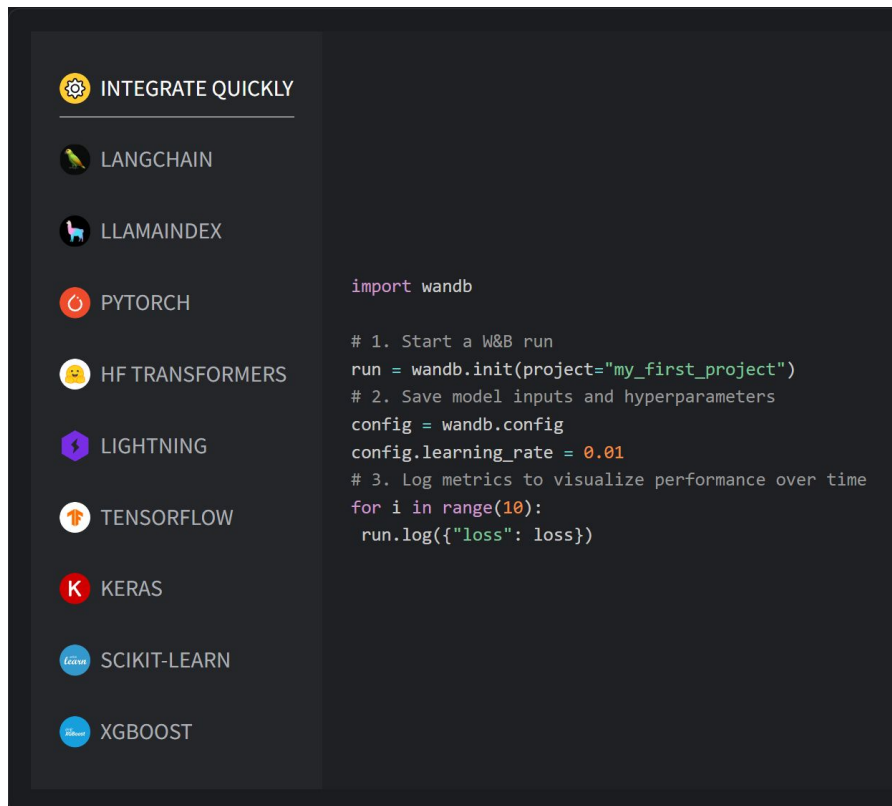
```
1 from hamilton import driver
2 import tabular_model
3 import model_evaluation
4
5 def main(cfg):
6     dr = (
7         driver.Builder()
8         .with_modules(tabular_model, model_evaluation)
9         .with_config({"task": cfg.task, "mode": cfg.mode})
10        .build()
11    )
12
13    inputs = dict(
14        data_dir=cfg.data_dir,
15        offset=cfg.offset,
16        label=cfg.label,
17        n_trials=cfg.n_trials,
18        n_folds=cfg.n_folds,
19    )
20    final_vars = ["best_model", "performance_metrics"]
21    dr.execute(final_vars, inputs=inputs)
```

- Pass config at the “driver level”
 - Driver **config**
 - Driver **modules**
 - **inputs**
 - **overrides**
- Breakdown large dataflows into multiple scripts
 - prepare_data.py
 - train_model.py
 - benchmark.py

Tracking

Experiment Results

MLFlow and Weights&Biases



- Direct integrations with many ML libraries
- Provides a web UI to explore results
- Facilitate collaboration

Limitations

- Integrations are actually confusing
- Hard to trace the code producing an artifact
- Weak code-artifact versioning

Node level



```
1 def best_model(  
2     X_train: np.ndarray,  
3     y_train: np.ndarray,  
4     base_model: xgboost.XGBModel,  
5     best_hyperparameters: dict,  
6     artifact_path: str,  
7     model_name: str,  
8 ) → xgboost.XGBModel:  
9     model = base_model  
10    model = model.set_params(  
11        early_stopping_rounds=None,  
12        **best_hyperparameters,  
13    )  
14    model.fit(X_train, y_train)  
15  
16    mlflow.xgboost.log_model(  
17        model,  
18        artifact_path=artifact_path,  
19        registered_model_name=model_name  
20    )  
21    return model
```

Driver level



```
1 from hamilton import driver  
2 import tabular_model  
3 import model_evaluation  
4  
5 def main(cfg):  
6     dr = (  
7         driver.Builder()  
8         .with_modules(tabular_model, model_evaluation)  
9         .with_config({"task": cfg.task, "mode": cfg.mode})  
10        .build()  
11    )  
12    inputs = dict(  
13        data_dir=cfg.data_dir,  
14        offset=cfg.offset,  
15        label=cfg.label,  
16        n_trials=cfg.n_trials,  
17        n_folds=cfg.n_folds,  
18    )  
19    final_vars = ["best_model", "performance_metrics"]  
20    results = dr.execute(final_vars, inputs=inputs)  
21  
22    best_model = results["best_model"]  
23  
24    with mlflow.start_run():  
25        run.set_tags(inputs)  
26        mlflow.log_params(best_model.params)  
27        mlflow.xgboost.log_model(  
28            best_model,  
29            artifact_path=cfg.artifact_path,  
30            registered_model_name=cfg.model_name,  
31        )  
32        mlflow.log_metrics(results["performance_metrics"])
```

Building the ExperimentTracker adapter

```
1 from hamilton import driver
2 from hamilton.io.materialization import to
3 from hamilton.plugins.h_experiments import ExperimentTracker
4 import tabular_model
5 import model_evaluation
6
7 def main(cfg):
8     dr = (
9         driver.Builder()
10        .with_modules(tabular_model, model_evaluation)
11        .with_config({"task": cfg.task, "mode": cfg.mode})
12        .with_adapters(
13            ExperimentTracker(
14                experiment_name=cfg.experiment_name,
15                base_directory="./experiments",
16            )
17        )
18        .build()
19    )
20    inputs = dict(
21        data_dir=cfg.data_dir,
22        offset=cfg.offset,
23        label=cfg.label,
24        n_trials=cfg.n_trials,
25        n_folds=cfg.n_folds,
26    )
27    materializers = [
28        to.json(
29            id="best_model_json",
30            dependencies=["best_model"],
31            path="./xgboost_model.json"
32        ),
33        to.json(
34            id="eval_stats_json",
35            dependencies=["performance_metrics"],
36            path="./eval_stats.json"
37        ),
38    ]
39    dr.materialize(*materializers, inputs=inputs)
```

- artifacts == Hamilton nodes
- Specify artifacts with **.materialize()**
- Artifact versioned according to executed code
- Bundled with an extensible UI

Quick

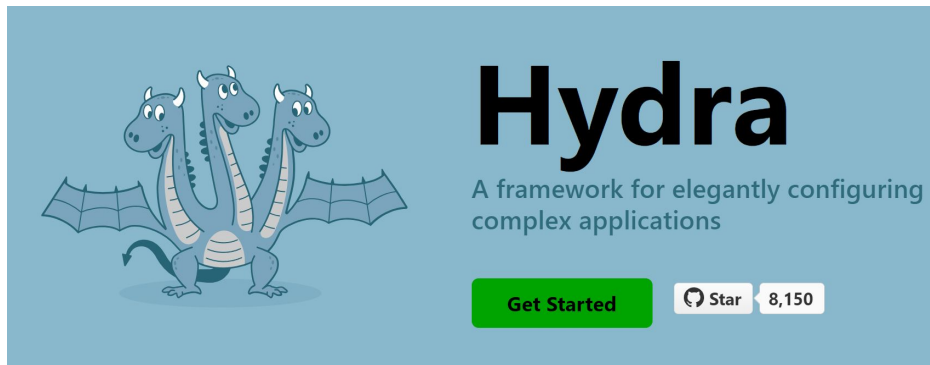
Live Demo

Extras

Extra: Project structure

- Separate “packages” from script/notebook execution
- Separate input and output data
- Breakdown your project into several dataflows

Things I've tried



- Hydra is a yaml-first config system
- Popular for deep learning research

- Metaflow is a Python orchestrator
- Easy to deploy yourself on AWS



**Everything you need to develop
data science, ML, and AI apps**