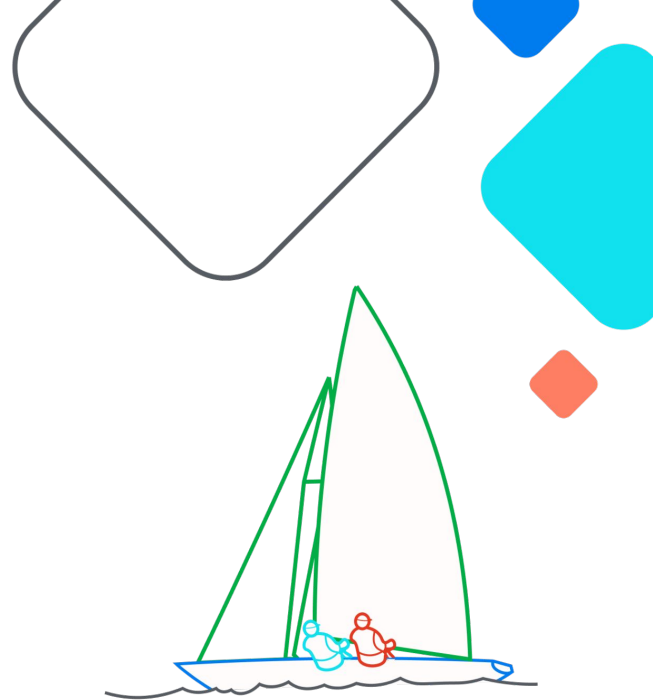


# Democratizing ML feature store framework at scale

Rafay Aleem & Victoria Varney



 **Airflow Summit**

Let's flow together

September 19-21, 2023,  
Toronto, Canada



Rafay Aleem  
Faire  
Sr Data Engineer - ML platform



Victoria Varney  
Astronomer  
Sr Customer Success Manager

# Today we are talking about...

Democratizing a  
machine learning  
feature store  
framework at Faire

How we have enabled  
everyone to contribute  
to a shared resource  
easily

Building a framework  
on top of Airflow by  
Leveraging its low-  
level APIs

How Airflow is much  
more than just a  
workflow  
management tool

F & S  
C<sup>o</sup>

F U  
S A M  
—

F A I R E

FAIRE CO.  
f  
FAIRE CO.

F A I R E

san francisco

C A F

f

F A I R E

F A I R E

F A I R E

FAIRE CO.  
the  
FAIRE CO.

F

san francisco

faire

f

the future  
|

F A I R

FAIRE CO.  
FAIRE CO.

E

is local

—  
faire  
—

f

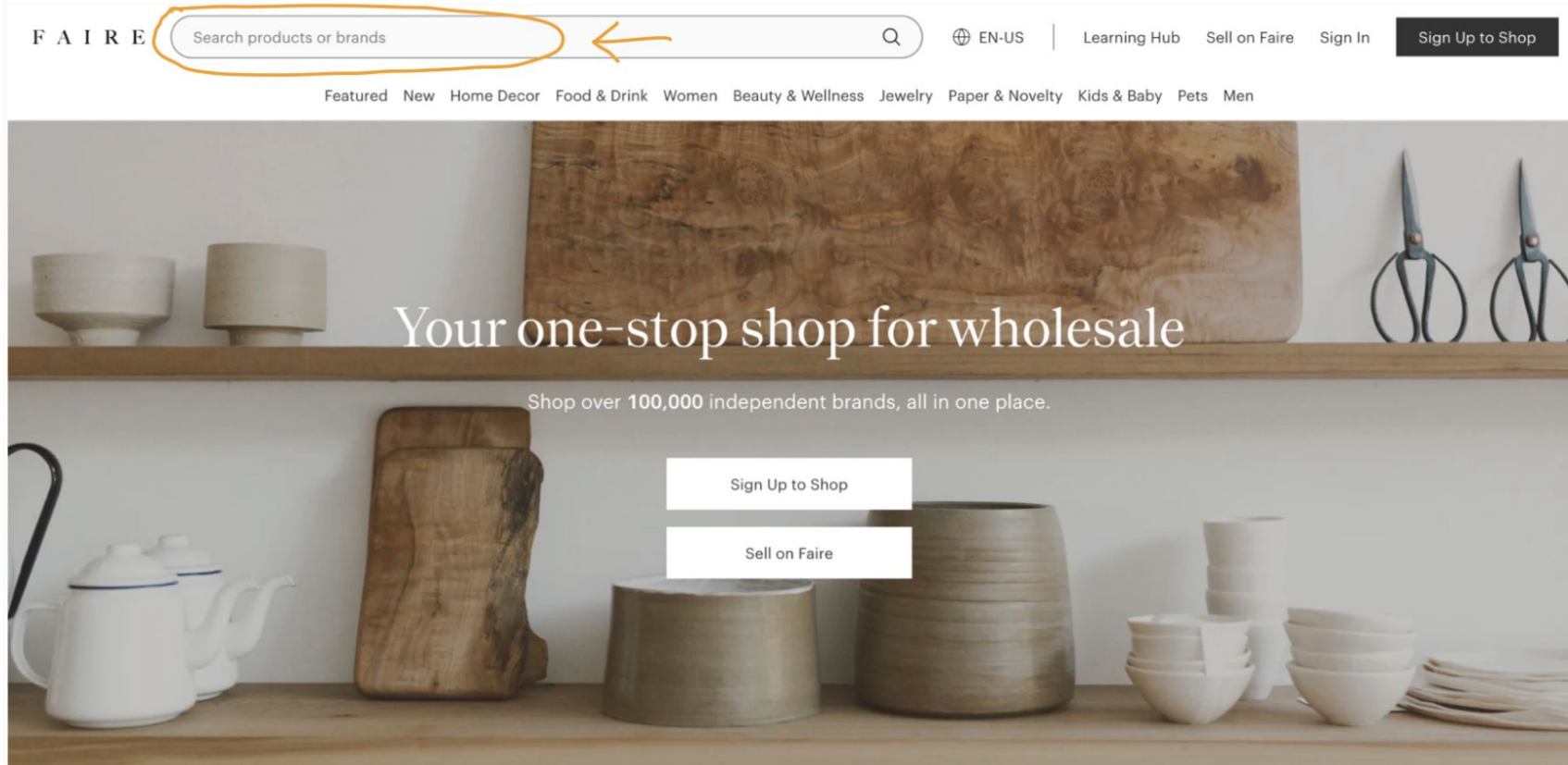
F A I R E

san francisco

faire

san francisco

# Search experience at Faire



## Wholesale Food & drink

Shop independent brands from around the world at wholesale prices.



Baking



Beverages



Cereals,  
Grains,  
&  
Pastas



Coffee & Tea



Condiments  
& Sauces

↑↓ ↕ Location ▾ Lead time ▾ Minimum ▾ Storage ▾ Shelf life ▾ Diet ▾ Production ▾ Values ▾ Promotor ▾



High sell-through

**FREEZE DRIED JOLLY BALLS (JOLLY RANCHERS)**  
MSRP \$9.95

Upscale Freeze



High sell-through

**2oz Sea Salt Chips (case of 20 bags)**  
MSRP \$49.15

1 in 6 Snacks- Carolina...



High sell-through

**Savory Party Cracker Seasoning - Classic...**  
MSRP \$9.83

Savory Fine Foods LLC



High sell-through

**Freeze Dried Rainbow Bites**  
MSRP \$9.23

Sow Good Inc



High sell-through

**Simply Mints - Peppermint**  
MSRP \$4.08

Simply Gum

### Brands

- Crazy Dog T-Shirts
- Merch Mallow T-shirts
- Rescue Dogs Rock T-shirts
- Cool Chili T-shirts
- Silicon Valley Tshirts

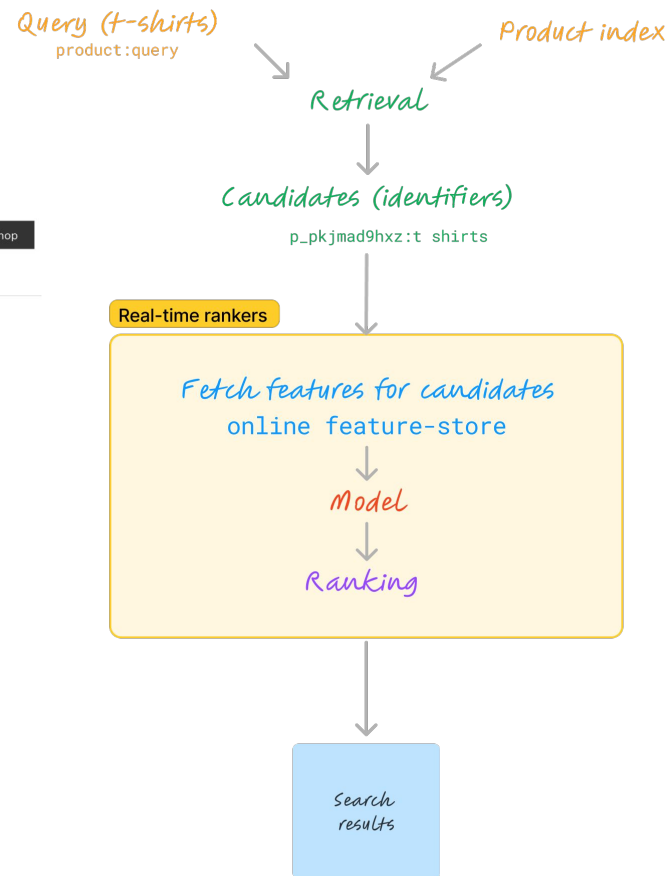
<p><b>FAST SHIPPING SAME DAY PICK UP WHOLESALE PRICING</b></p> <p>Blank Shirts, Bella Canvas 3001 T-shirt, Unisex Soft Shirts MSRP \$19.09 </p> <p><a href="#">Print The Dream</a></p>	<p>Bella Canvas 3001 T-shirt, Blank Shirts, Unisex Soft Shirts MSRP \$19.09 </p> <p><a href="#">Print The Dream</a></p>	<p>Gildan Soft Adult Shirt, Blank Unisex T-shirt MSRP \$16.36 </p> <p><a href="#">Print The Dream</a></p>	<p>Unisex Jersey Short Sleeve Tees - Bella Canvas MSRP \$10.24 </p> <p><a href="#">ShopTrendsNow</a></p>	<p>High sell-through</p> <p>You're doing better than you think terry graphic sweatshirt MSRP \$73.70 </p> <p><a href="#">Costa Threads</a></p>
--	---	---	--	--

### Similar products

<p>High sell-through</p> <p>Recycled Cotton Crop Packdeal \$23.21 MSRP \$21.72 Studio Ko Clothing \$100 min</p>	<p>Cottonclub Her Day T Packdeal \$20.48 MSRP \$65.53 fabina \$1.28 min</p>	<p>CWTT50434_BASIC Simple Daily... \$12.07 MSRP \$25.84 Lily Clothing \$136.53 min</p>	<p>Top Shop</p> <p>Bivy Classic Fit Crew Neck Tee \$23.20 MSRP \$46.09 Urban Daze \$0 min</p>	<p>V111-1 Raw Edge Tee \$10.14 MSRP \$30.02 HUMSII \$0 min</p>
<p>New</p> <p>4 Pieces Short Sleeve T-Shirts... \$46.41 MSRP \$92.81 Print The Dream \$68.26 min</p>	<p>plus</p> <p>Double Layered Side Shirting Th... \$13.65 MSRP \$27.30 P3 Kate \$1.37 min</p>	<p>Knitted Tank Top \$13.52 MSRP \$32.49 Miss Stockholm \$136.53 min</p>	<p>Top Shop</p> <p>N110029 \$18.43 MSRP \$36.86 VERY J LOVE RICHIE \$136.53 min</p>	<p>Garment Dyed Unisex Rib Neck... \$17.04 MSRP \$34.08 U.S. My Sunshine \$136.53 min</p>

# Anatomy of a search

The screenshot shows a search engine interface for 't-shirts'. The search bar contains 'tshirts' with a yellow arrow pointing to it. Below the search bar, there are filters for 'Brands' and 'Ship window'. The main content area displays a grid of product cards. Each card includes a product image, a title, a price (MSRP), and a 'Print The Dream' button. The products shown are: 'Blank Shirts, Bella Canvas 3001 T-shirt, Unisex Soft Shirts' (MSRP \$19.09), 'Bella Canvas 3001 T-shirt, Blank Shirts, Unisex Soft Shirts' (MSRP \$19.09), 'Gildan Soft Adult Shirt, Blank Unisex T-shirt' (MSRP \$16.36), 'Unisex Jersey Short Sleeve Tees - Bella Canvas' (MSRP \$10.24), and 'You're doing better than you think terry graphic sweatshirt' (MSRP \$73.70). A 'High sell-through' badge is visible on the sweatshirt image.



# Anatomy of a search

Query (t-shirts)  
product:query

Product index

Retrieval

Candidates (identifiers)

p\_pkjmad9hxz:t shirts

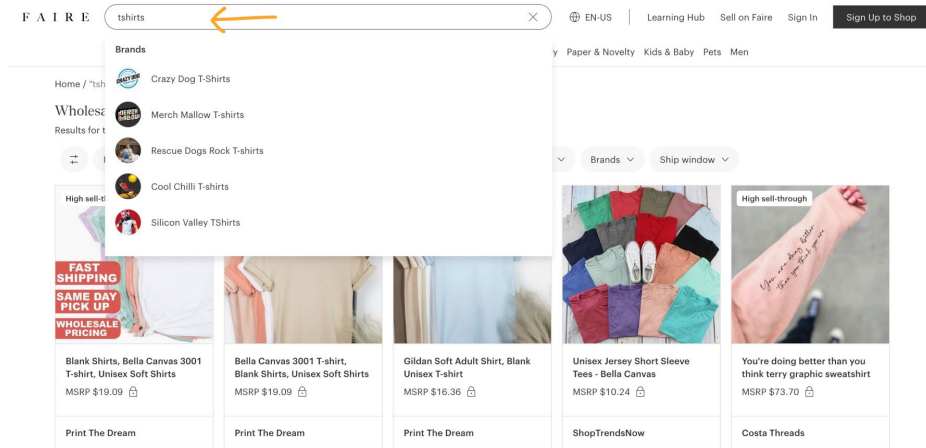
Real-time rankers

Fetch features for candidates  
online feature-store

Model

Ranking

Search results





# Enhancing the search experience

- Unlock the science behind search relevancy
- Search is particularly important in context of e-Commerce

F A I R E | electric diffusers for home | My Brands | Refer to earn up to \$2,000 | EN-US | 1 Brand

Featured | New | Bestsellers | Home Decor | Paper & Novelty | Women | Food & Drink | Beauty & Wellness | Jewelry | Kids & Baby | Men | Pets

Home / "electric diffusers for home"

**Filters**

**Location**

Search

Canada 32  
 United States 201  
 United Kingdom 128  
[Show more](#)

**Brand minimum**

No minimum 1  
 \$100 or less 36  
 \$200 or less 238  
[Show more](#)

**Category**

Home Decor 1,000+  
 Beauty & Wellness 654

**Wholesale price**

\$0 - \$10 531  
 \$10 - \$30 525  
 \$30 - \$50 26  
[Show more](#)

**Product types**

Search

Fragrance, Aromatherapy & Essen... 518

**Results for electric diffusers for home** 1,095 results [Hide filters](#)

Discover new products and enjoy free returns on all opening orders.

[Low minimum](#) [Canadian brands](#) [New this month](#)

**Fragrance Incense Burning Oil 2.4, 16 oz...**  
\$5.45 MSRP \$10.90  
[Love&Lust LLC](#)  
\$136.58 min

**TECHANCY Oil Diffuser Humidifier 120ML**  
\$20.10 MSRP \$40.19  
[TECHANCY](#)  
\$150.79 min

**Oil Warmer Round Copper**  
\$11.47 MSRP \$22.94  
[Mar Company Group](#)  
\$273.15 min

**Calm Diffuser Reeds - Wholesale**  
\$17 MSRP \$30  
[Market Candle Company](#)  
\$175 min

**Seven Chakra Electric Aroma Lamp**  
\$22.35 MSRP \$68.08  
[Something Different Wholesale](#)  
\$150 min

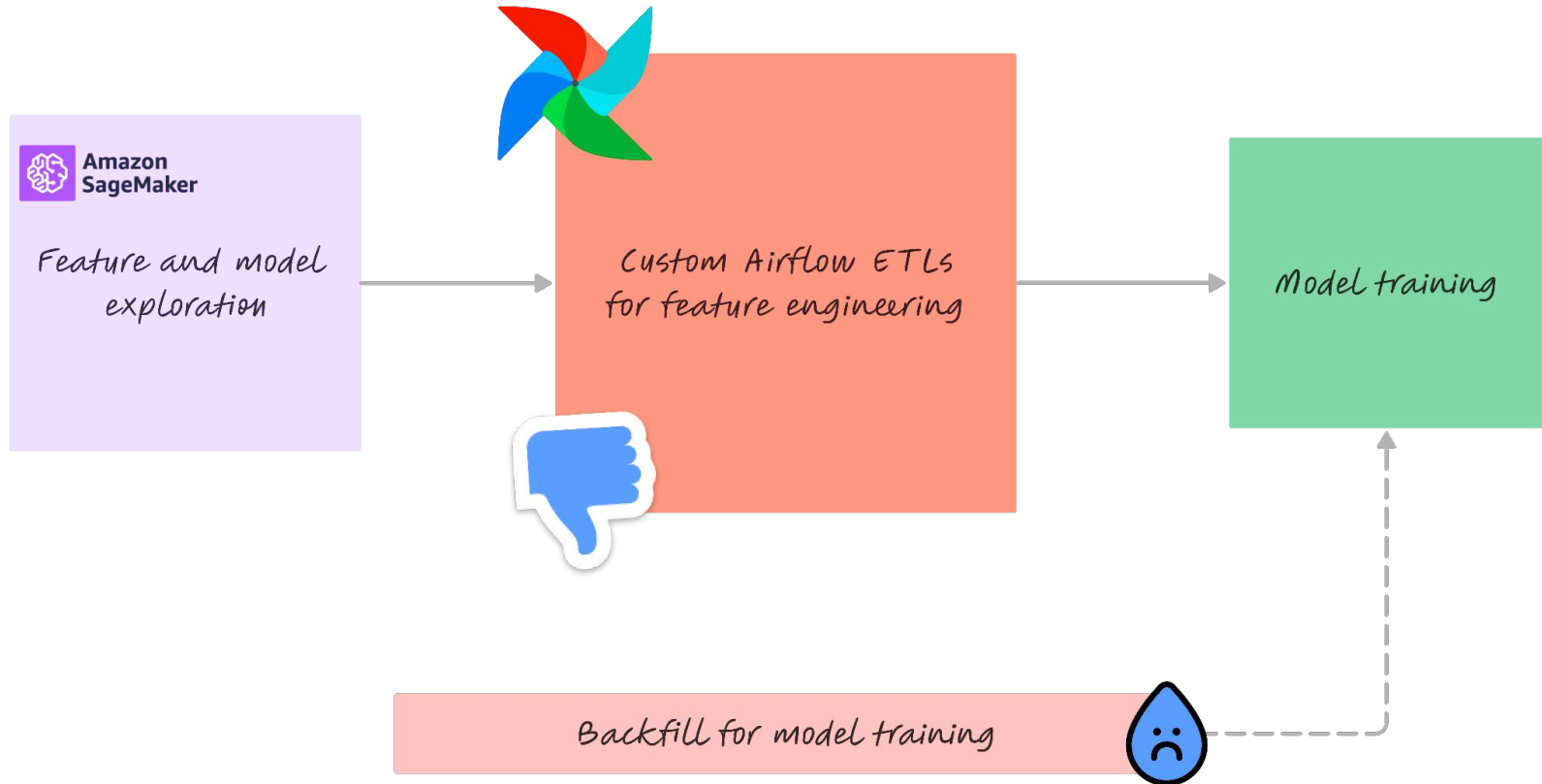
**Cedarwood Organic Essential Oil**  
\$5.46 MSRP \$10.93  
[Silver Aromatic](#)  
\$109.26 min

# Enhancing the search experience

- Features are **float values** with some **underlying meaning** such as
  - num\_countries\_sold\_in
  - avg\_fulfillment\_days\_trailing\_7\_days
- Features are engineered and defined on the **offline** and eventually **propagated** to **online** feature store

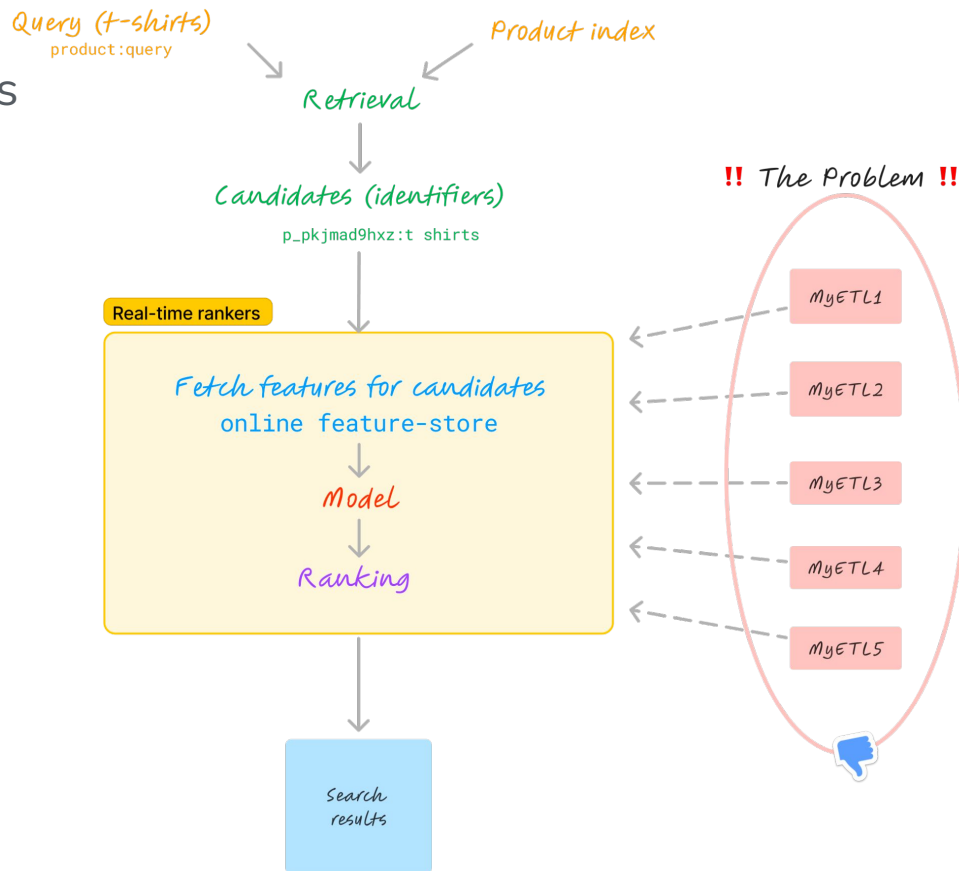


# The need for democratizing the feature store



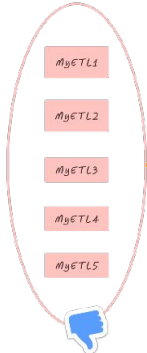
# The need for democratizing the feature store

- No **feature visibility** across teams
- No clear process to separate **online features** from offline
- Non-standardized **error-prone backfilling** with no clear notion of point-in-time joins
- Overhead **costs** on ElasticCache and Snowflake



# The need for democratizing the feature store

!! The Problem !!



Consolidate all feature definitions into a single feature store with a complete feature registry

Complex & Error-prone backfills

Simplify and standardize feature backfills



SQL + Custom ETLs



```
SELECT
  SUM(AMOUNT) AS TOTAL_SALES
FROM ORDERS
WHERE CATEGORY = "candles"
```

Plain SQL as feature definitions



# Feature Store Framework

- All features are defined as **SQL** with their **Python configurations**
- Each feature can be **configured separately** with its own **metadata**

## Examples

- specify if available online
  - feature description
  - author
- 
- All configuration is fed into a queryable **feature registry**

```
select
    b.token as brand_token,
    count(distinct rb.retailer_id) as number_of_contact_books_found
from production.retailer_to_brand_email_domains rb
join production.brands b
    on b.url_domain = rb.email_domain
where b.active
    and rb.created_at < '{{ ds }}'::DATE
group by 1
```

```
feature_groups = [
    FeatureTask(
        sql_file="brand_contact_book.sql",
        features=[
            Feature(
                sql_feature_name="number_of_contact_books_found",
                feature_description="Number of retailers with this brand in
their address book",
                available_online=False, ★
            )
        ],
        entity=BrandEntity(),
        author="rafay", ★
        is_static=False,
    )
]

brand_features = FeatureEntity(entity=BrandEntity(), tasks=feature_groups)
```

# Feature Store Framework

```
class FeatureEntity:
    def __init__(
        self,
        entity: EntityDescription,
        tasks: list[FeatureTask],
    ):
        self.entity = entity
        self.tasks = tasks
        self.name = entity.get_feature_name_string()

    for task in tasks:
        assert type(self.entity) == type(task.entity)
        self.check_columns_present_in_sql(task)

    # check dup feature names
    self.check_dup_feature_names()

    def check_dup_feature_names(self) -> None:
        """
        Check whether feature names are unique.
        """
        ...

    def get_compute_operator(
        self, task: FeatureTask, is_backfill=False, start_date=None, end_date=None
    ) -> MLFeatureSnowflakeOperator:
        ...
```

Low-level APIs

```
@dataclass
class Feature:
    sql_feature_name: str
    feature_description: str = "No description has been specified"
    available_online: bool = False
```

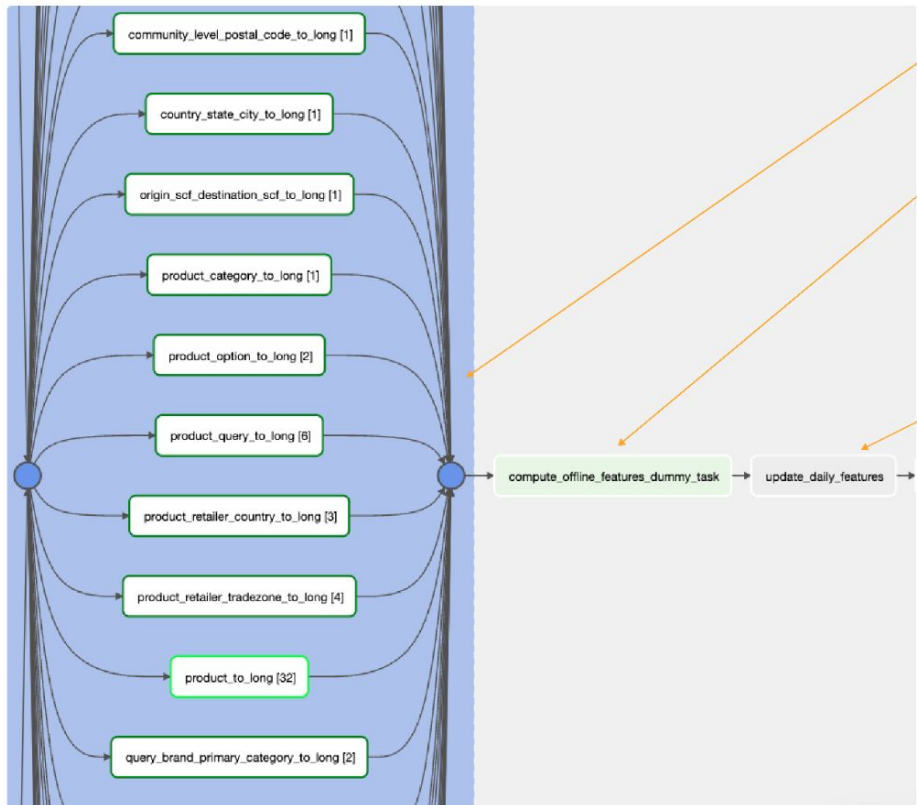
```
@dataclass
class FeatureTask:
    sql_file: str
    features: list[Feature]
    entity: EntityDescription
    author: str = "Unknown author"
    is_static: bool = False

    # backfill parameters
    backfill_config: FeatureTaskBackfillConfig = None
```

High-level config dataclasses

# Feature Store Framework

AmendTableAndSchemaOperator AmendTableOperator EmptyOperator MLFeatureSnowflakeOperator PythonOperator **RebuildTableSnowflakeOperator** SnowflakeOperator



```
# This step creates individual tables for each invocation of MLFeatureSnowflakeOperator
# You can find these tables in {schema_name}.features_{feature_group_name}
with TaskGroup(group_id="compute_offline_features") as compute_offline_features:
    for feature_entity in FEATURE_ENTITIES:
        feature_entity.get_dynamically_mapped_tasks(config)

# This is a dummy task used to make it easier to manage the DAG in the Airflow UI
compute_offline_features_dummy_task = EmptyOperator(
    task_id="compute_offline_features_dummy_task",
    doc="Dummy task to trigger offline features computation",
    trigger_rule="all_done",
)

# This pulls the daily feature used for dispersal and pivoting.
# Data is persisted to avoid duplicated computation and facilitate data checking.
update_daily_features = AmendTableOperator(
    task_id="update_daily_features",
    schema_name=config.get("schema"),
    table_name=config.get("features_daily_table_name"),
    create_sql="sql/create_features_daily_table.sql",
    sql="sql/update_features_daily_table.sql",
    fill_condition="ds = '{{ ds }}'",
    data_checks=[
        f"checks/{check}.sql"
        for check in (
            "update_daily_features_task_lands_data",
        )
    ]
)
```



# Adopting new Airflow features

```
ML-314: Fix slow dag parsing in compute offline feature store dag that is causing zombie tasks #17037
Changes from all commits - File filter - Conversations -
```

```
airflow/dags/feature_store/features/utlils/FeatureEntity.py
121 - def get_compute_operators(self) -> List[MLFeatureSnowflakeOperator]:
122     """
123     - Generate list of Airflow tasks to perform the computation of all feature tasks
124     """
125     return [self.get_compute_operator(task) for task in self.tasks]
126
127 - def get_compute_task_group(self) -> TaskGroup:
128     with TaskGroup(group_id=self.name) as task_group:
129         self.get_compute_operators()
130     return task_group
131
132
133 + def get_dynamically_mapped_tasks(self, params) -> MappedOperator:
134     """
135     + Generate list of dynamically mapped Airflow tasks to perform the computation of
136     + all feature tasks
137     + Args:
138     +     params: This is the same config that's passed to the DAG params. The reason
139     +     we have to explicitly
140     +     pass it here is because expand_kwargs would not render the templated
141     +     fields in the operator
142     +     when using it like {{ params.features_static_table_name }} in
143     +     combination with partial and
144     +     expand_kwargs. Also note that expand_kwargs expects a list of
145     +     dictionaries with string
146     +     key value pairs.
147     """
148     # Please see https://docs.astronomer.io/learn/dynamic-tasks?
149     # tab=traditional#sets-of-keyword-arguments
150     # for understanding partial and expand_kwargs in Airflow's Dynamic Task Mapping
151     # API
152     # expand_kwargs takes the following form when using Dynamic Task Mapping API
153     # {
154     #     ("bash_command": "echo $WORD", "env": {"WORD": "hello"}),
155     #     ("bash_command": "echo {expr length $WORD} ", "env": {"WORD": "tea"}),
156     #     ("bash_command": "echo ${WORD}/e/X)", "env": {"WORD": "goodbye"}),
157     # }
158     kwargs = list(
159     map(
160         lambda t: Status: Merged
161         "feature_task": t,
162         "destination_table_name": params["features_static_table_name"] if
163         t.is_static
164         else params["features_long_table_name"]
165     ),
166     self.tasks
167 )
168 # We use partial and expand_kwargs here to optimize dag parsing and delay task
169 # generation to runtime
170 return MLFeatureSnowflakeOperator.partial(
171     task_id=f"{self.name}_to_long",
172     start_date=None,
173     end_date=None,
174     is_backfill=False,
175     ).expand_kwargs(kwargs)
```

Eason@wang marked this conversation as resolved. Show resolved

```
def get_dynamically_mapped_tasks(self, params) -> MappedOperator:
    kwargs = list(
        map(
            lambda t: {
                "feature_task": t,
                "destination_table_name": params["features_static_table_name"] if t.is_static
                else params["features_long_table_name"]
            },
            self.tasks
        )
    )

    # We use partial and expand_kwargs here to optimize dag parsing and delay task generation to runtime
    return MLFeatureSnowflakeOperator.partial(
        task_id=f"{self.name}_to_long",
        start_date=None,
        end_date=None,
        is_backfill=False,
    ).expand_kwargs(kwargs)
```

Use partial and expand\_kwargs to dynamically map existing task

# Adopting new Airflow features

- Be **proactive** in adopting new versions of Airflow
- New features have greatly **improved** cluster performance
- Notable features:
  - **Dynamic Task Mapping**
  - **Task Groups** (easier visual dag management)
  - **Deferrable operators** for long running tasks such as AWS Batch jobs

# Extending the framework for feature backfills

```
FeatureTask(  
    sql_file="brand_contact_book.sql",  
    features=[  
        Feature(  
            sql_feature_name="number_of_contact_books_found",  
            feature_description="Number of retailers with this brand in their address book",  
            available_online=False,  
            backfill_config=FeatureTaskBackfillConfig(  
                start_date=datetime(2023, 8, 1),  
                end_date=datetime(2023, 8, 30),  
                namespace="brand_1_month_backfill",  
            ),  
        ),  
    ],  
    entity=BrandEntity(),  
    author="rafay",  
    is_static=False,  
)
```

feature\_store\_backfill\_brand\_1\_month\_2023\_8\_1\_2023\_8\_30

- Uses same SQL files
- Provides extra jinja templated flag `{{ is_backfill }}`
- Entire process takes 3 lines of configuration!

```
def get_compute_operator(  
    self, task: FeatureTask, is_backfill=False, start_date=None, end_date=None  
) -> MLFeatureSnowflakeOperator:  
    """  
    Generate Airflow task to perform the computation of a feature task  
    """  
    Args:  
        task: FeatureTask object  
        is_backfill: specifies whether the operator should be generated for backfill  
        start_date: start date of the task, None implies it will match the dag start date  
        end_date: end date of the task, None implies it will match the dag end date  
    """  
    destination_table_name = (  
        "{{ params.features_static_table_name }}"  
        if task.is_static  
        else "{{ params.features_long_table_name }}"  
    )  
    return MLFeatureSnowflakeOperator(  
        task_id=task.sql_name,  
        feature_task=task,  
        destination_table_name=destination_table_name,  
        start_date=start_date,  
        end_date=end_date,  
        is_backfill=is_backfill,  
    ) # type: ignore
```

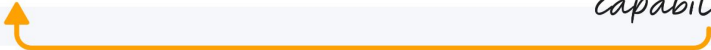
Internally passed to low-level API

# Best practices for designing extensible frameworks

- Use mixins to extend operator capabilities

```
class RebuildTableSnowflakeOperator(  
    SnapshotMixin, TableOperatorMixin, ETLSnowflakeOperator  
):
```

*Extend operator capabilities using mixins*




- Allows adding new features to low-level APIs without breaking

```
class SnapshotConfig:  
    Frequency = Enum('Frequency', ['NEVER', 'DAILY', 'EVERY_RUN'])
```

```
store_daily_embeddings = RebuildTableSnowflakeOperator(  
    task_id="store_daily_embeddings",  
    sql=daily_online_embedding_sql(),  
    schema_name=config.get("schema"),  
    table_name="daily_embeddings_v2",  
    snapshot_frequency=SnapshotConfig.Frequency.EVERY_RUN,  
)
```

*Enable table snapshots using predefined frequency*



# Airflow beyond workflow management

*Airflow goes far  
and beyond a  
workflow  
management tool*

*It's thoughtful  
design makes it  
very extensible  
and powerful*

*It is very well  
suited for running  
mission critical  
workflows with  
tight SLA  
requirements*

*It continues to be  
a very stable part  
of Faire's  
infrastructure  
and continue to  
scale*

# Key Takeaways

Proactively consider onboarding to latest Airflow versions

Don't think of Airflow as just a workflow management tool

Consider building shared frameworks instead of shared ETLs

Thoughtful use of Airflow APIs and features goes a long way

# Credits

- [Wayne Zhang](#) for his guidance on the offline feature store framework
- **Analytics Engineering** team for their feedback on table snapshot tooling
- **Core Data Infra** team for their constant support with Airflow and Snowflake
- **Machine Learning Platform** team for dealing with on-call issues and providing stakeholder support
- My wife's constant support

# Wrap Up & Questions

- Careers @ Faire: [faire.com/careers/](https://faire.com/careers/)
- Where to find Rafay
  - LinkedIn: [linkedin.com/in/mrafayaleem/](https://www.linkedin.com/in/mrafayaleem/)
  - Airflow Slack: [@Rafay Aleem](#)
  - Twitter: [@mrafayaleem](#)
  - Email: [contact@mrafayaleem.com](mailto:contact@mrafayaleem.com)
  - Newsletter: [mrafayaleem.com/#/portal/](https://mrafayaleem.com/#/portal/)



**Thank You!**