

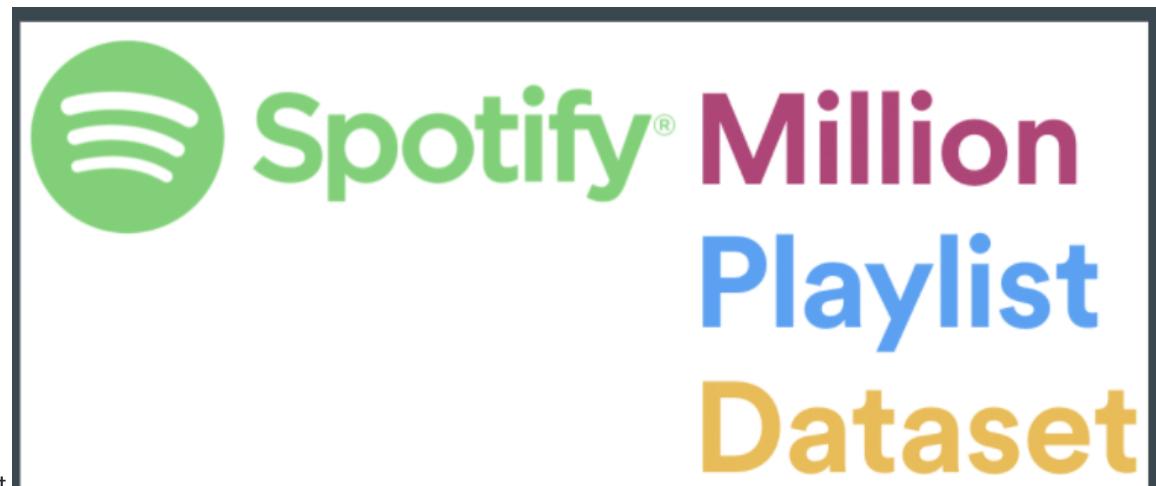
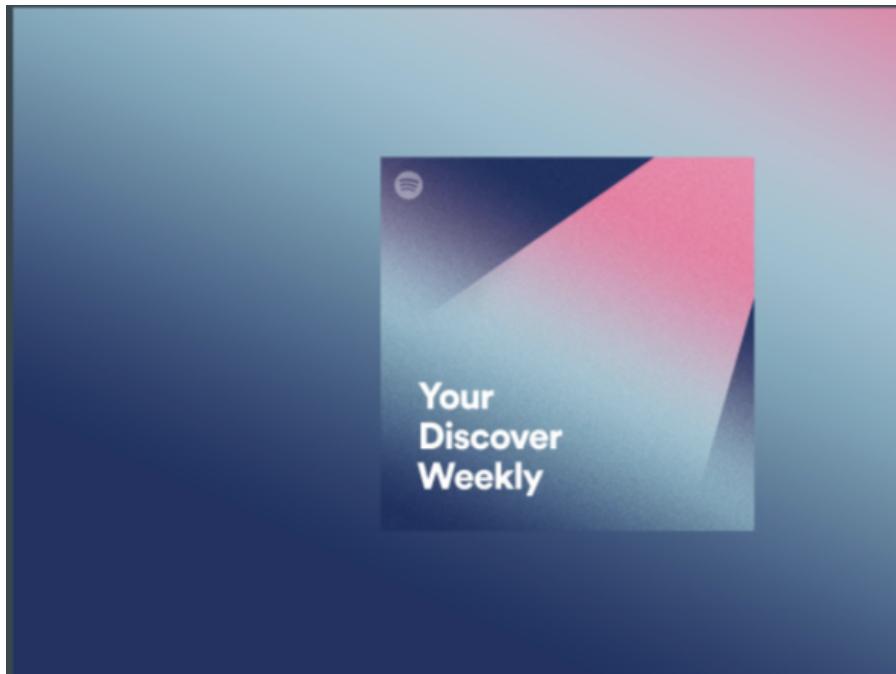
README.md

# Spotify Recommendations

## Overview

The Problem @ 

- I like Spotify and I like discovering music
- Music generation is very popularity based nowadays
- Discovery/interaction with music happens via playlists mostly nowadays
- Music curation tools for playlists are scarce
- Spotify has a lot of tools for interacting with their data/features



- Spotify hosted the 2018 RecSys challenge of 1 Million Playlists made by users
- Publicly released the dataset

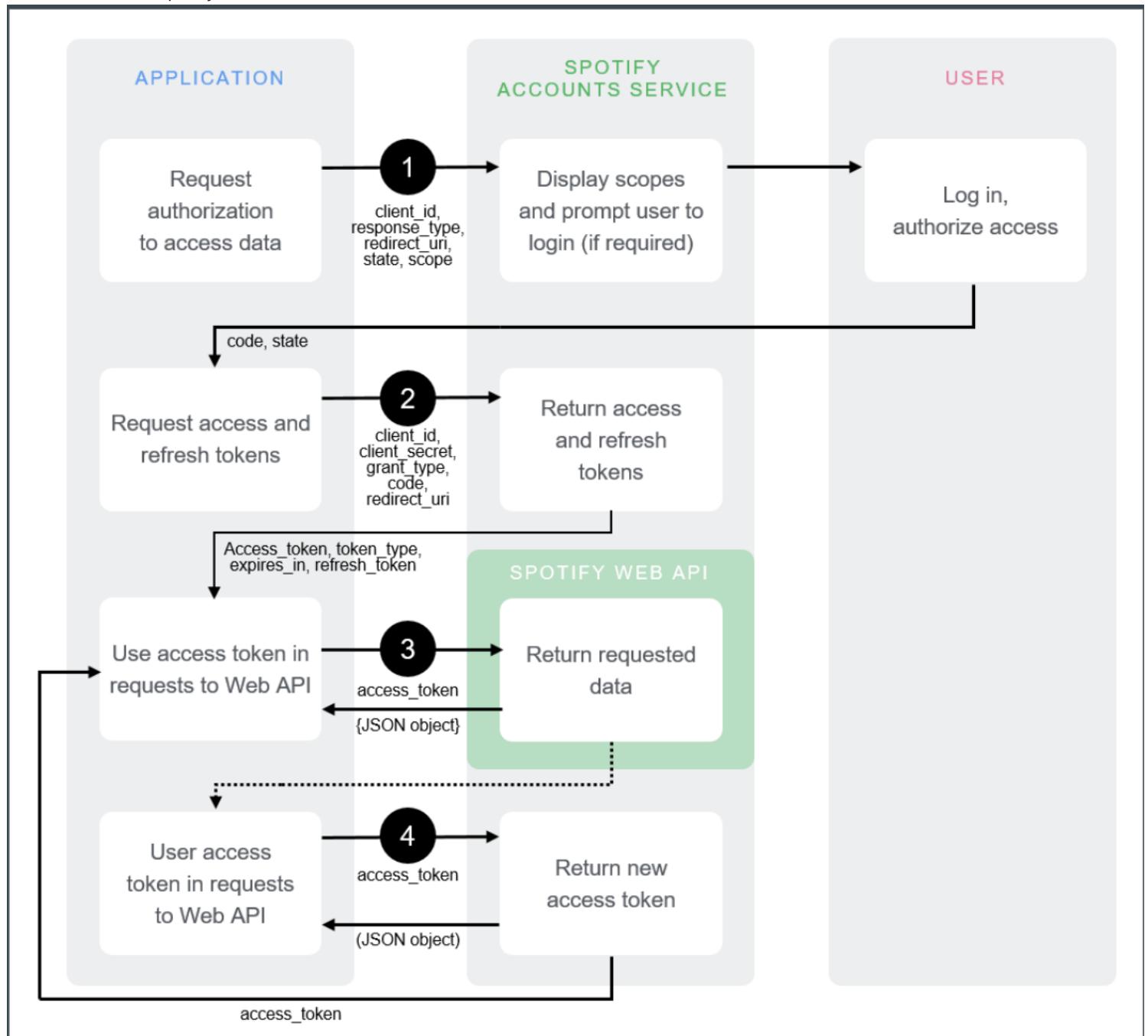
- No associated metadata or tracks

Property	Value
Number of playlists	1,000,000
Number of tracks	66,346,428
Number of unique tracks	2,262,292
Number of unique albums	734,684
Number of unique artists	295,860
Number of unique playlist titles	92,944
Number of unique normalized playlist titles	17,381
Average playlist length (tracks)	66.35

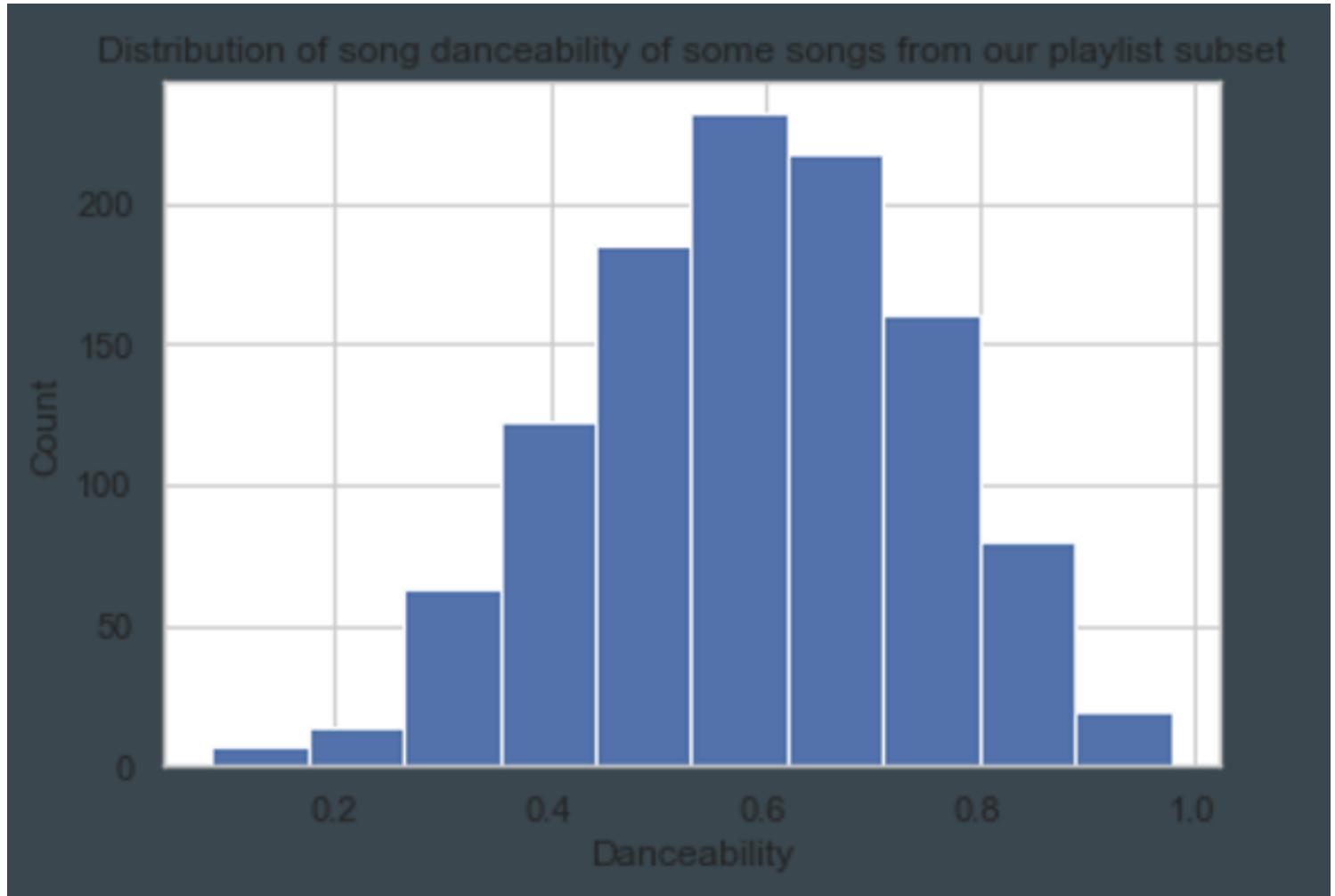
```
{
  "name": "musical",
  "collaborative": "false",
  "pid": 5,
  "modified_at": 1493424000,
  "num_albums": 7,
  "num_tracks": 12,
  "num_followers": 1,
  "num_edits": 2,
  "duration_ms": 2657366,
  "num_artists": 6,
  "tracks": [
    {
      "pos": 0,
      "artist_name": "Degiheugi",
      "track_uri": "spotify:track:7vqa3sDmtEaVJ2gcvxtRID",
      "artist_uri": "spotify:artist:3V2paBXEoZIAhfZRJmo2jL",
      "track_name": "Finalement",
      "album_uri": "spotify:album:2KrRMJ9z7Xjoz1Az406UML",
      "duration_ms": 166264,
      "album_name": "Dancing Chords and Fireflies"
    },
    {
      "pos": 1,
      "artist_name": "Degiheugi",
      "track_uri": "spotify:track:23E0mJiv0Z88WJPUbIPjh6",
      "artist_uri": "spotify:artist:3V2paBXEoZIAhfZRJmo2jL",
      "track_name": "Betty",
      "album_uri": "spotify:album:3lUSlvjUoHNA8IkNTqURqd",
      "duration_ms": 235534,
      "album_name": "Endless Smile"
    },
    {
      "pos": 2,
      "artist_name": "Degiheugi"
    }
  ]
}
```

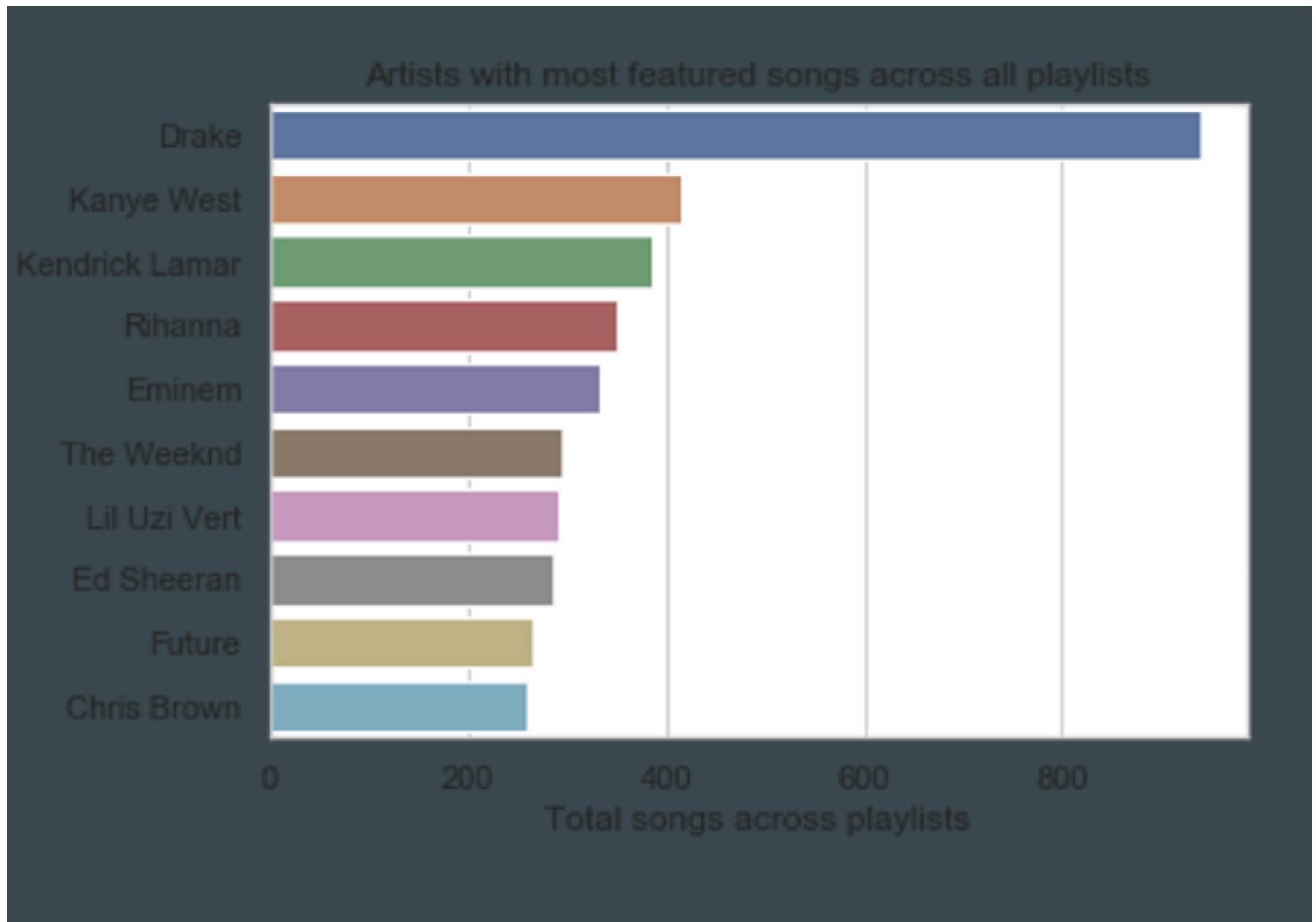
```
        "artist_name": "Begin Again",
        "track_uri": "spotify:track:1vaffTCJxkyqeJY7zF9a55",
        "artist_uri": "spotify:artist:3V2paBXEoZIAhfZRJmo2jL",
        "track_name": "Some Beat in My Head",
        "album_uri": "spotify:album:2KrRMJ9z7Xjoz1Az406UML",
        "duration_ms": 268050,
        "album_name": "Dancing Chords and Fireflies"
    },
    // 8 tracks omitted
{
    "pos": 11,
    "artist_name": "Mo' Horizons",
    "track_uri": "spotify:track:7iwx00eBzeSSSy6xfESyWN",
    "artist_uri": "spotify:artist:3tuX54dqgS8LsGUvNzgrpP",
    "track_name": "Fever 99\u00b0",
    "album_uri": "spotify:album:2Fg1t2ty0SGWkVYHlFfXVf",
    "duration_ms": 364320,
    "album_name": "Come Touch The Sun"
}
],
```

## Track Metadata - Spotify Web API



## EDA - MPD + Track Metadata





### The Approach 🧑

- Build infrastructure/tooling to easily build and deploy Spotify-based applications and models
- Do so in a streamlined reproducible fashion
- Automate as much as possible
- Infrastructure-as-code as much as possible
- Make it easy for a SWE/DE/DS/MLE/DevOps person to extend upon
- Use a MLFlow wrapped custom cosine similarity model as a test of E2E tooling from the ML end

### Infrastructure considerations

- arbitrary custom model can be made easily
- models can connect to growing data sources
- deployment of said model is straightforward
- can be managed by ci/cd easily
- adaptable to data, models, and code changes

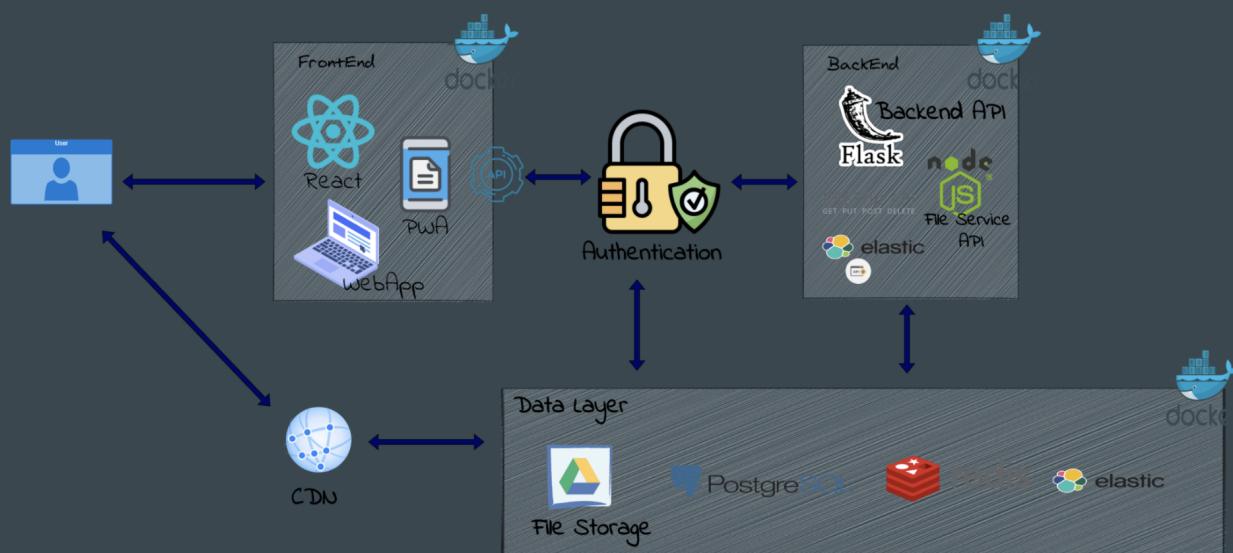
### The Experiment - Let's Keep it Simple

- Data: Create training and test data by way of obfuscation on MPD
- Take all the playlists and randomly remove a certain amount of songs and use as training data
- Use those removed songs and use as testing data

- Task: Use training songs in the playlist to predict the songs that were removed Method: Aggregate features of the songs to create a representative feature representation for the playlist w.r.t songs and use a similarity search
- Features used: Spotify given track features
- Search: Cosine Similarity

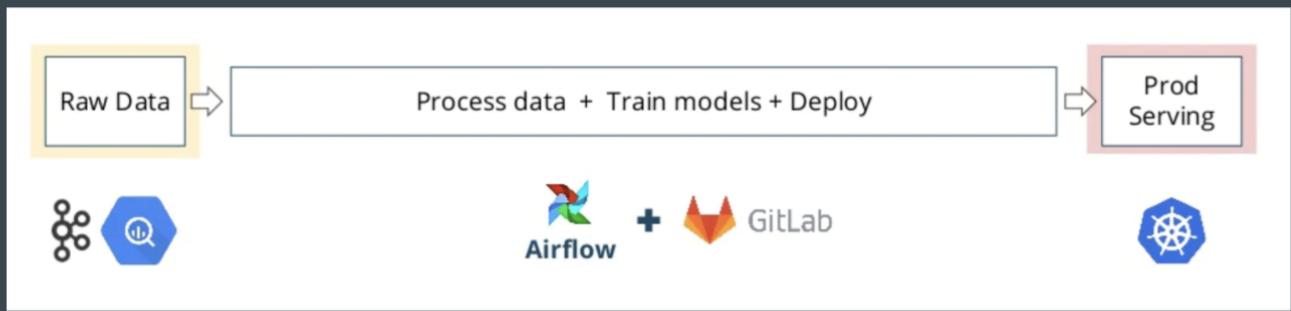
Introducing SpotifyPlaylists! ♪♪

## High Level Software Architecture



\*In our scenario authentication occurs via Spotify and all actions occur on behalf of the provided user's token

Before



Data Engineer

Data Scientist

Software Engineer

After



Time based

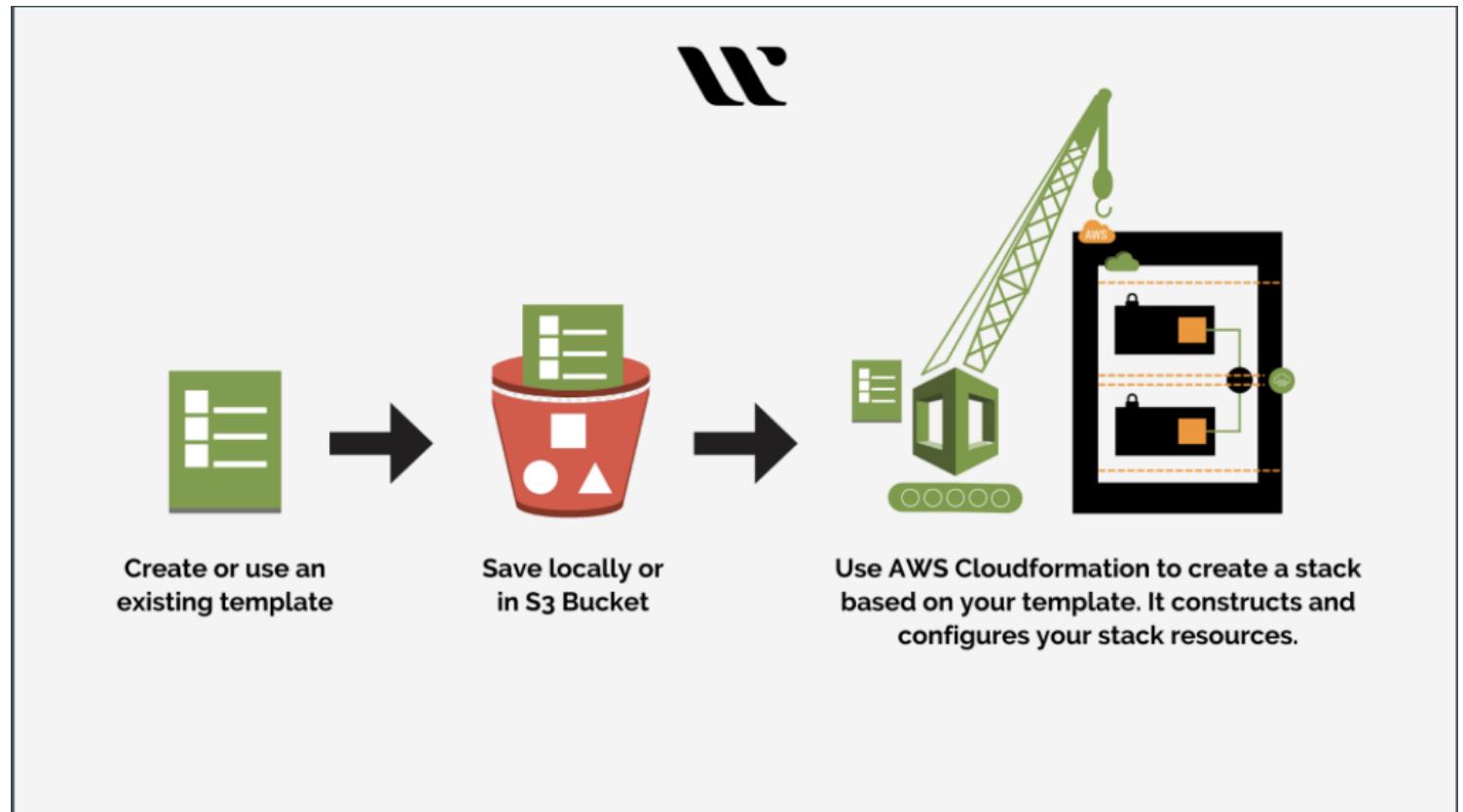
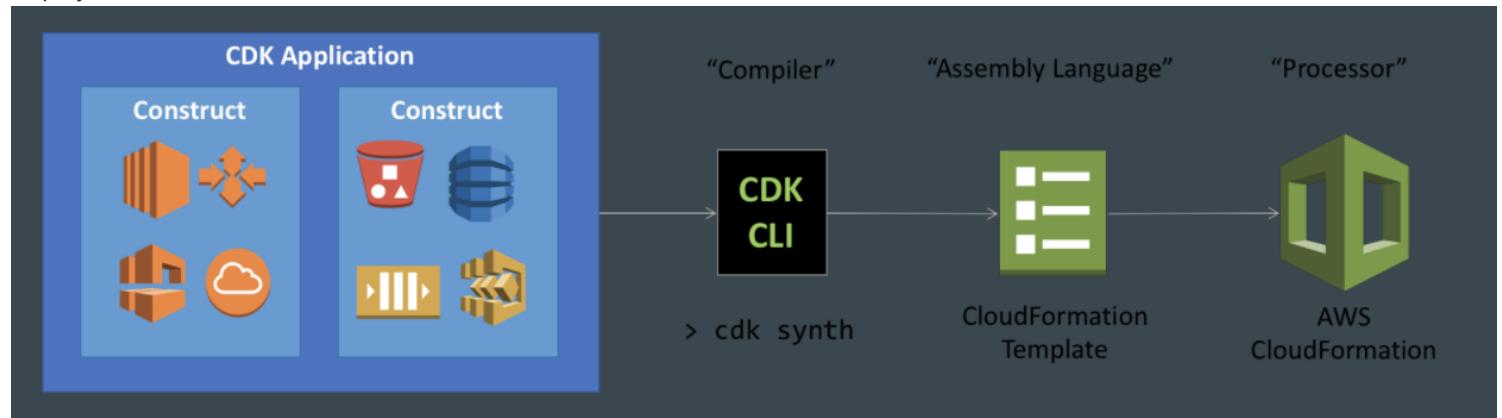


Instance based



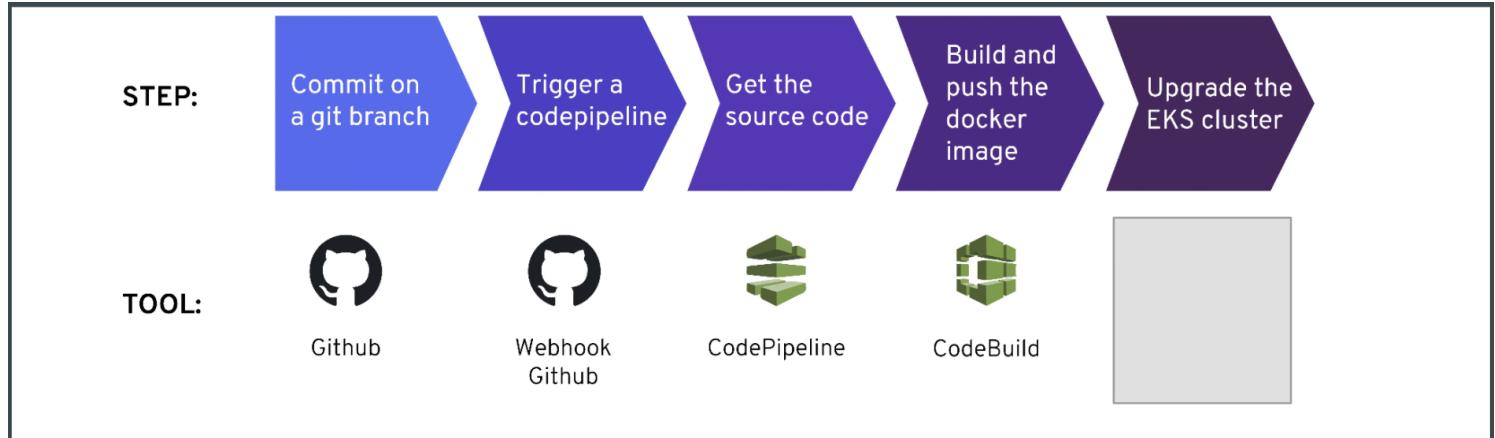
Artifact based

## Deployment Stack



CloudFormation > Stacks					
Stacks (5)					
Stack name		Status	Created time	Description	
spotify-recommendations-build-and-deploy-eks		CREATE_COMPLETE	2021-09-17 13:57:24 UTC-0400	EKSWSV1	
eksctl-spotify-recommendations-eksctl-nodegroup-nodegroup		CREATE_COMPLETE	2021-09-17 13:08:30 UTC-0400	EKS Managed Nodes (SSH access: false) [created by eksctl]	
eksctl-spotify-recommendations-eksctl-cluster		CREATE_COMPLETE	2021-09-17 12:49:05 UTC-0400	EKS cluster (dedicated VPC: true, dedicated IAM: true) [created and managed by eksctl]	
CDKToolkit		CREATE_COMPLETE	2021-09-16 16:16:40 UTC-0400	The CDK Toolkit Stack. It was created by `cdk bootstrap` and manages resources necessary for managing your Cloud Applications with AWS CDK.	
DeploymentStack		CREATE_COMPLETE	2021-09-15 16:12:18 UTC-0400	-	

## CodePipeline Stack



## Resources (6)

 Search resources

Logical ID	Physical ID	Type	Status
CodeBuildProject	spotify-recommendations-build-and-deploy-eks	AWS::CodeBuild::Project	 UPDATE_COMPL ETE
CodeBuildServiceRole	spotify-recommendations-build-CodeBuildServiceRole-1374IJOLEX8ST 	AWS::IAM::Role	 CREATE_COMPL ETE
CodePipelineArtifactBucket	spotify-recommendations-codepipelineartifactbucket-nlm20ugvjdwi 	AWS::S3::Bucket	 CREATE_COMPL ETE
CodePipelineGitHub	spotify-recommendations-build-and-deploy-eks-CodePipelineGitHub-1NCBWFBIC426R	AWS::CodePipeline::Pipeline	 UPDATE_COMPL ETE
CodePipelineServiceRole	spotify-recommendations-build-CodePipelineServiceRole-17GZSS7Z34GTC 	AWS::IAM::Role	 CREATE_COMPL ETE
EcrDockerRepository	spotify-recommendations-build-and-deploy-eks-ecrdockerrepository-lf4r2vvxh0fp 	AWS::ECR::Repository	 CREATE_COMPL ETE

## Parameters (7)

Search parameters

Key	Value
CodeBuildDockerImage	aws/codebuild/standard:4.0
EksClusterName	spotify-recommendations-eksctl
GitBranch	main
GitHubToken	*****
GitHubUser	anishshah97
GitSourceRepo	spotify-recommender-training
KubectlRoleName	CodeBuildKubectlRole

## spotify-recommendations-build-and-deploy-eks-CodePipelineGitHub-1NCBWFBIC426R

Notify  Edit  Stop execution  Clone pipeline  Release change

### ⌚ Source Succeeded

Pipeline execution ID: 0f6265fd-1663-402b-9871-88c31944be82

App GitHub (Version 1)  ⓘ  
 Succeeded - 1 hour ago  
ed3be577

ed3be577  App: Update buildspec.yml

Disable transition



### ⌚ Build In progress

Pipeline execution ID: 0f6265fd-1663-402b-9871-88c31944be82

Build AWS CodeBuild  ⓘ  
 In progress - Just now  
Details

ed3be577  App: Update buildspec.yml

Pre-requirements: EKS Cluster (using eksctl)

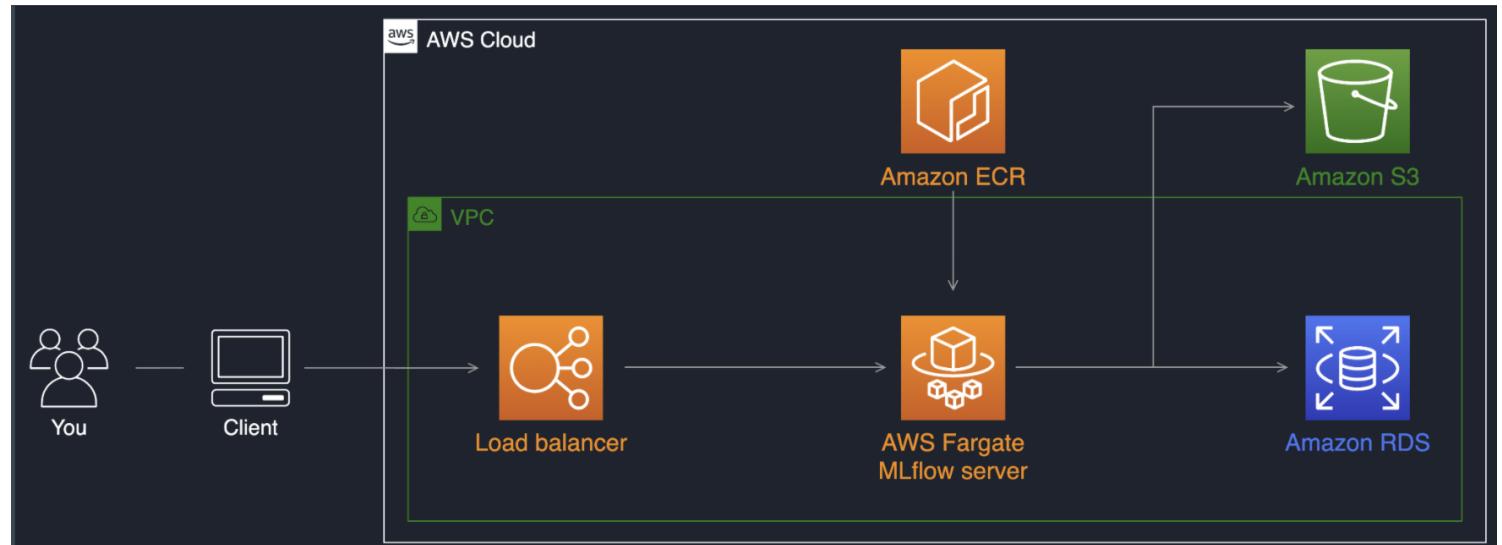
EKS > Clusters

**Clusters (1) Info**

Filter cluster by name, status, kubernetes version, or provider

Cluster name	Status	Kubernetes version	Provider
spotify-recommendations-eksctl	Active	1.21	EKS

MLFlow Stack



### Outputs (2)

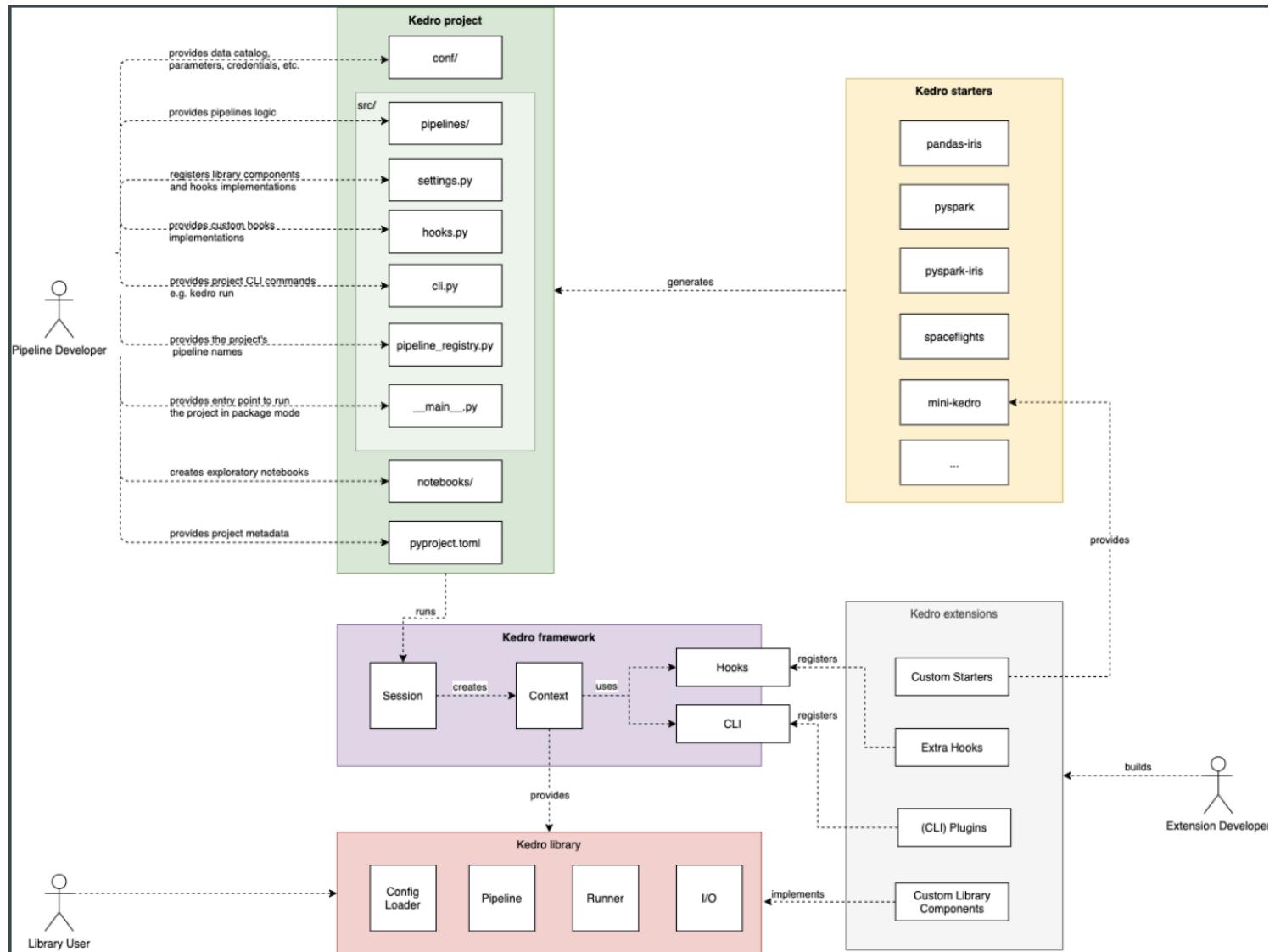
Search outputs

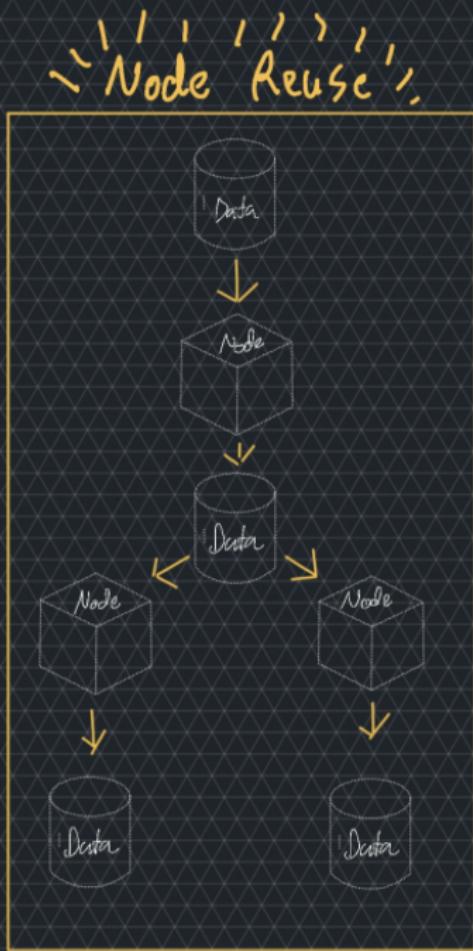
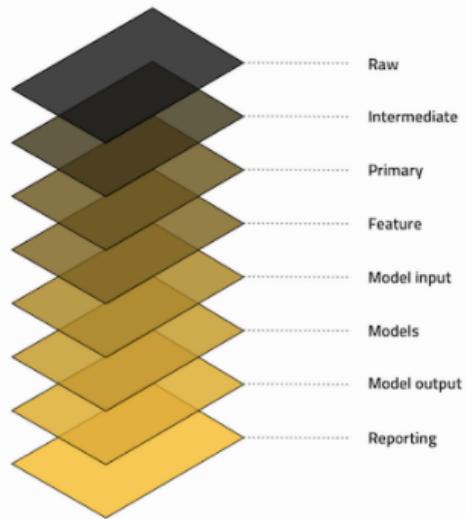
Key	Value	Description	Export name
LoadBalancerDNS	deplo-[REDACTED].elb.eu-west-1.amazonaws.com	-	-
MLFLOWLoadBalancerDNSAEFB7E43	deplo-[REDACTED].elb.eu-west-1.amazonaws.com	-	-

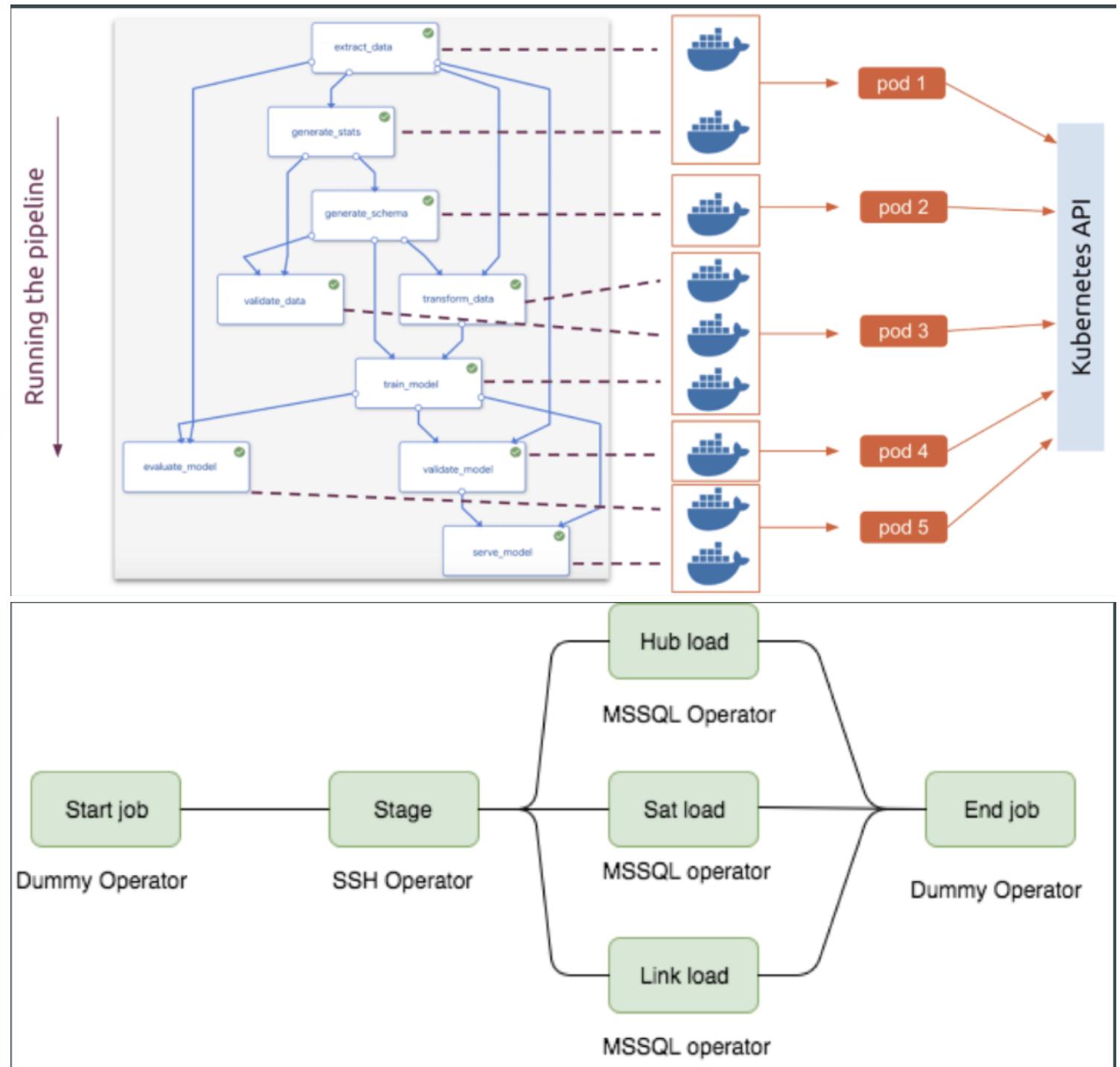
Timestamp	Logical ID	Status	Status reason
2021-09-15 16:26:15 UTC-0400	DeploymentStack	CREATE_COMPLETE	-
2021-09-15 16:26:12 UTC-0400	MLFLOWServiceTaskCountTargetAUTOSCALING4A5776E	CREATE_COMPLETE	-
2021-09-15 16:26:12 UTC-0400	MLFLOWServiceTaskCountTargetAUTOSCALING4A5776E	CREATE_IN_PROGRESS	Resource creation Initiated
2021-09-15 16:26:11 UTC-0400	MLFLOWServiceTaskCountTargetAUTOSCALING4A5776E	CREATE_IN_PROGRESS	-
2021-09-15 16:26:09 UTC-0400	MLFLOWServiceTaskCountTarget5798B92F	CREATE_COMPLETE	-

## Codebase - Kedro

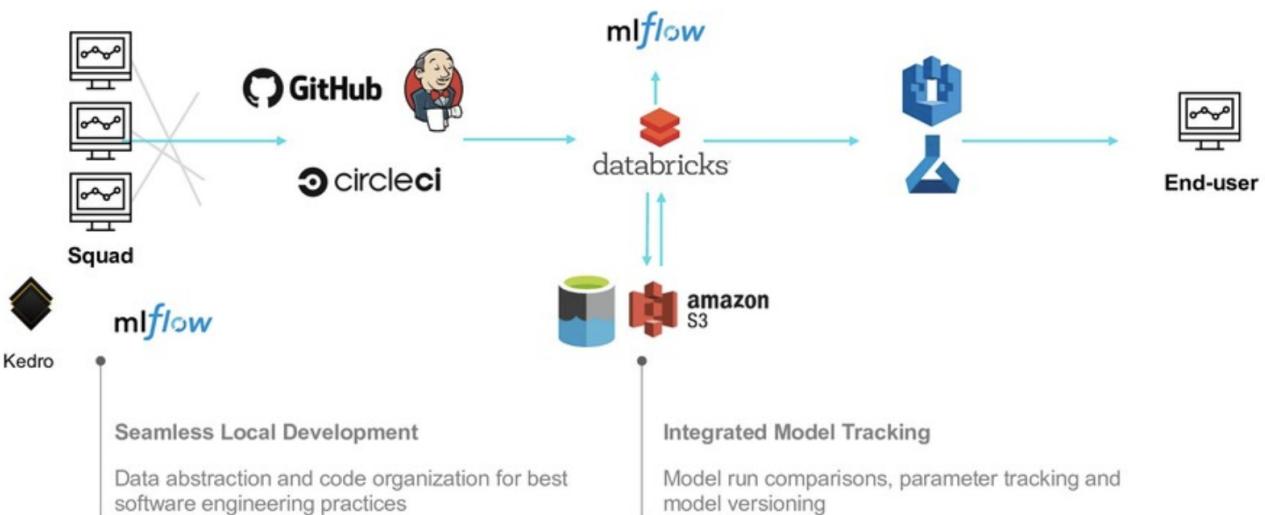
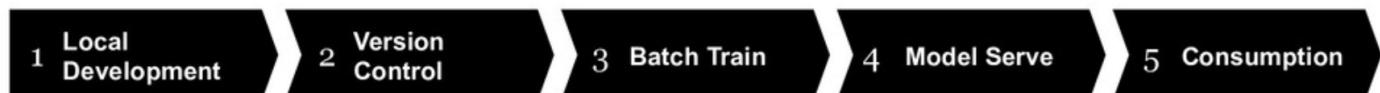




## An Aside - DAGS

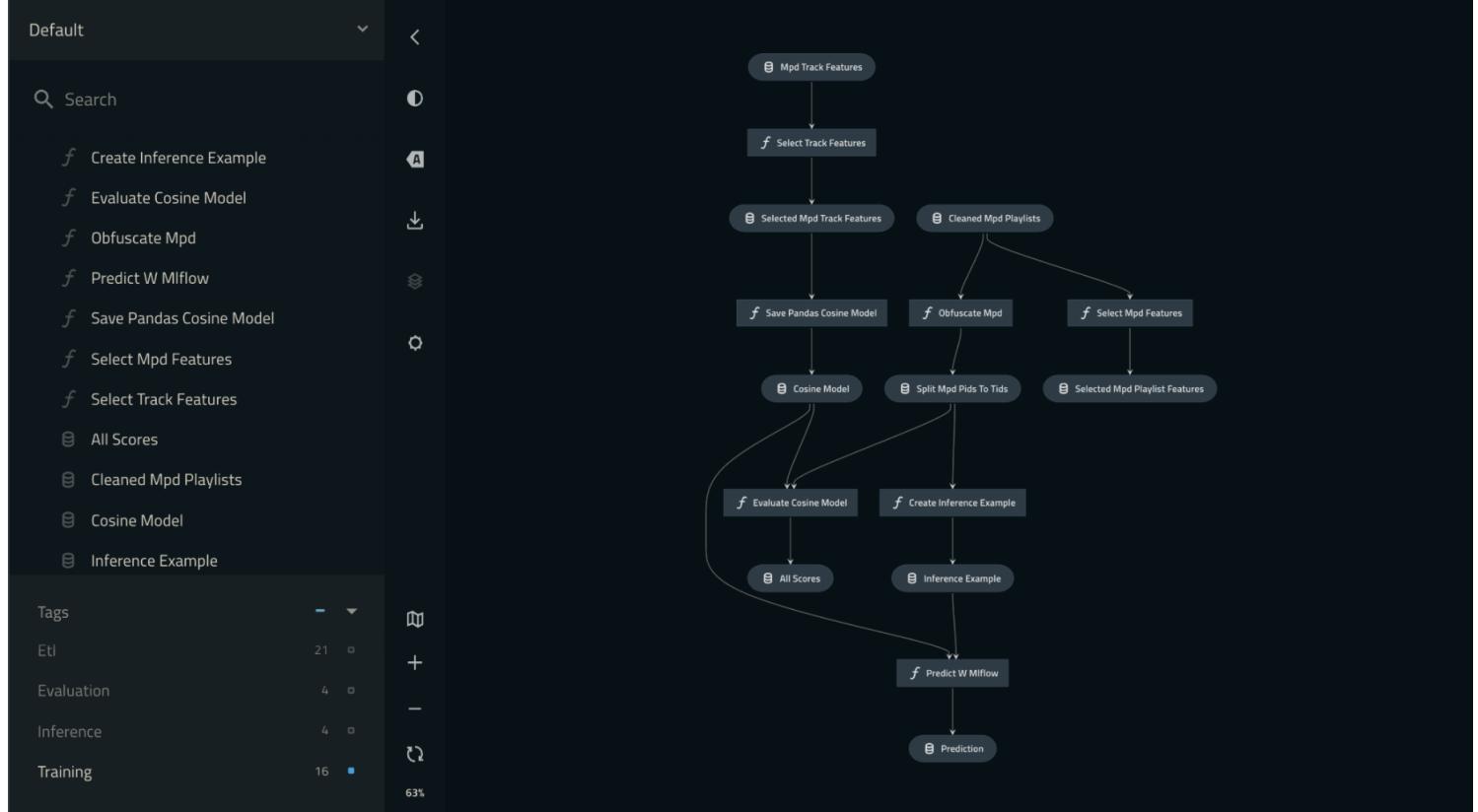
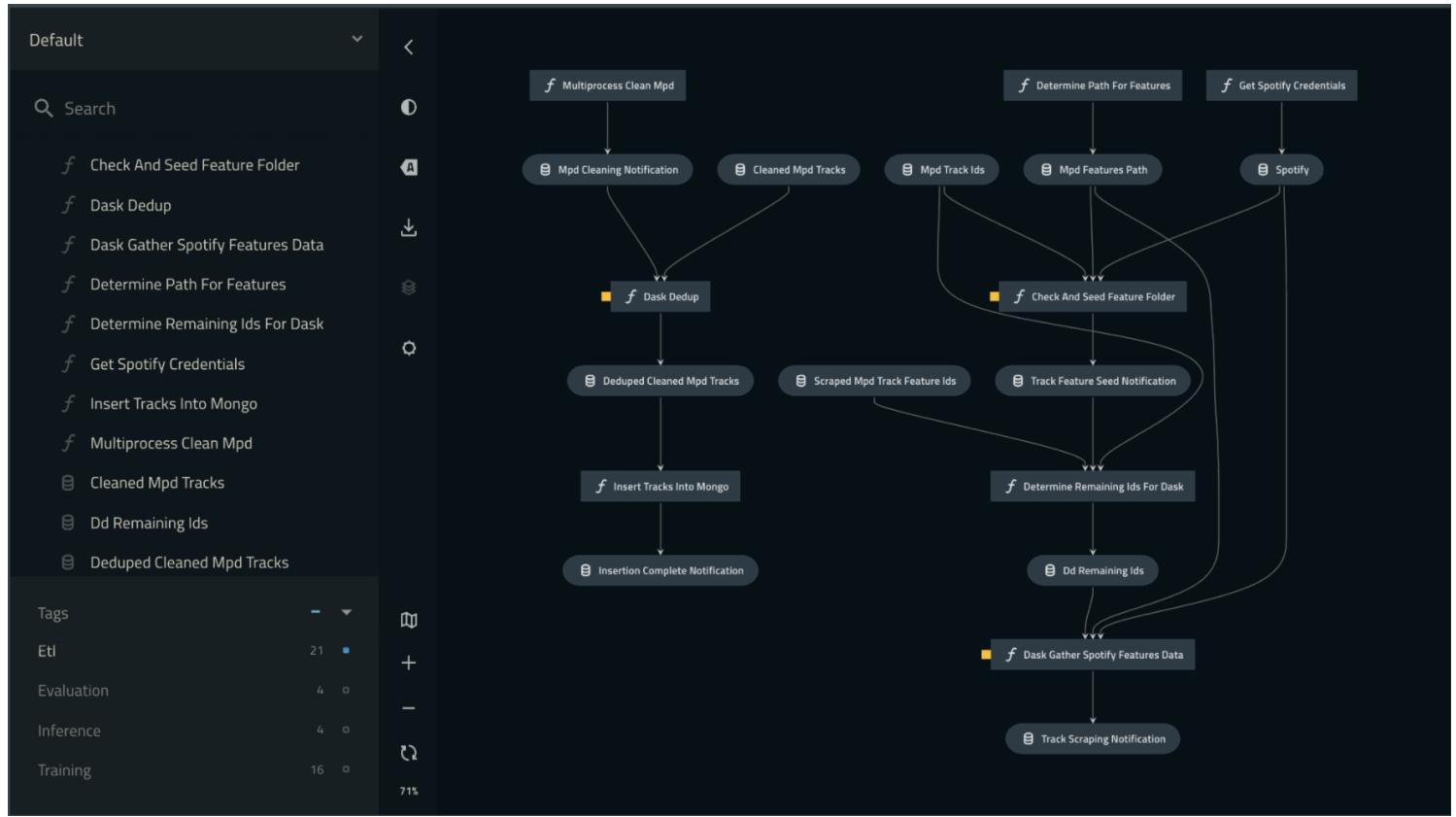


## Kedro and MLflow work together in both model development and deployment



All content copyright © 2019 QuantumBlack, a McKinsey company 22

## Training



## The Dynamic Duo in Action

Not Secure | deplo-milflow-16g0oyp6k65hv-5cd5faf094caf332.elb.us-east-1.amazonaws.com/#/experiments/1/runs/c59730008fac4f249af527798f57c4e

[Apps](#) [SAP Corporate Po...](#) [Jenkins](#) [Airflow](#) [SAP.I.O](#) [Issues](#) [BH AWS](#) [Office 365](#) [Resources](#) [FourthBrain](#) [Commuter Benefit...](#) [DuckDuckGo — Pr...](#)

[mlflow](#) Experiments Models GitHub Docs

spotify\_recommendations > mpd\_cosine\_ml

Date: 2021-09-16 09:08:11 Source: python -m kedro User: i854336

Duration: 1.2min Status: FINISHED

## ▼ Artifacts

- ▶ spotify\_recommendations

mlflow Experiments Models GitHub Docs

### Registered Models

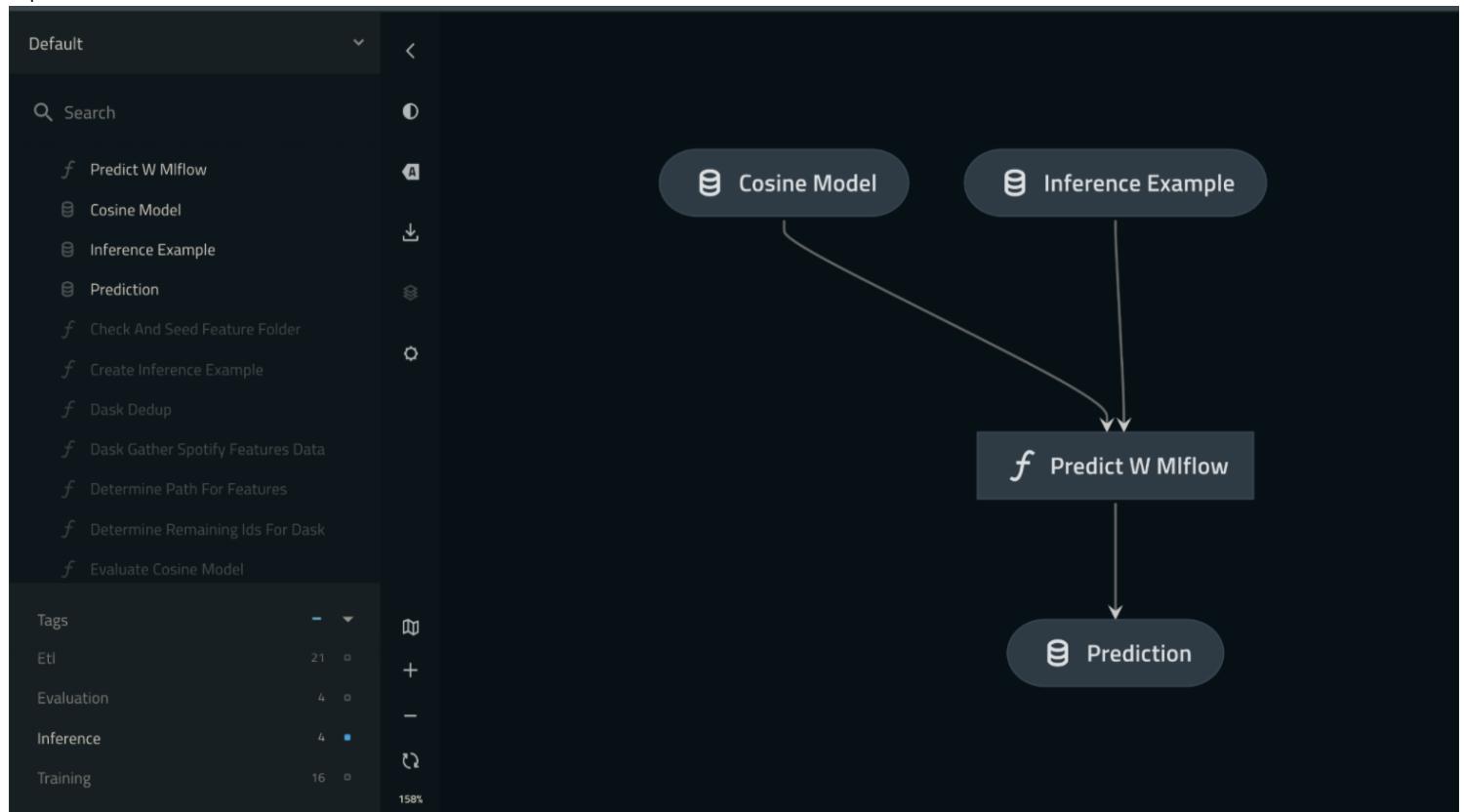
Share and serve machine learning models. Learn more

Create Model search model na...

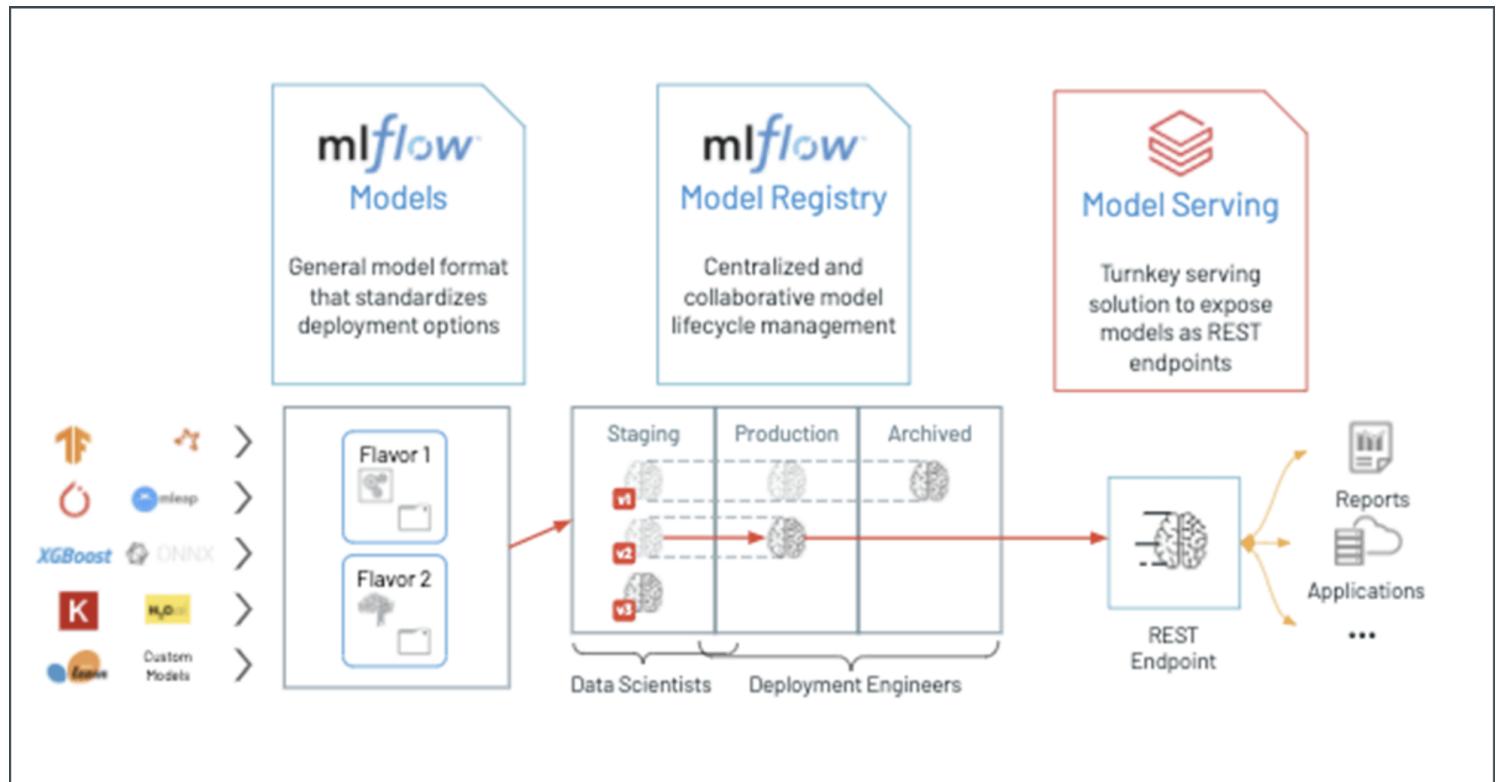
Name	Latest Version	Staging	Production	Last Modified
No models yet. Create a model to get started.				

< Page 1 > 10 / page ↻

## Pipelines as Models??? - Inference Served



You're sure that's the model?

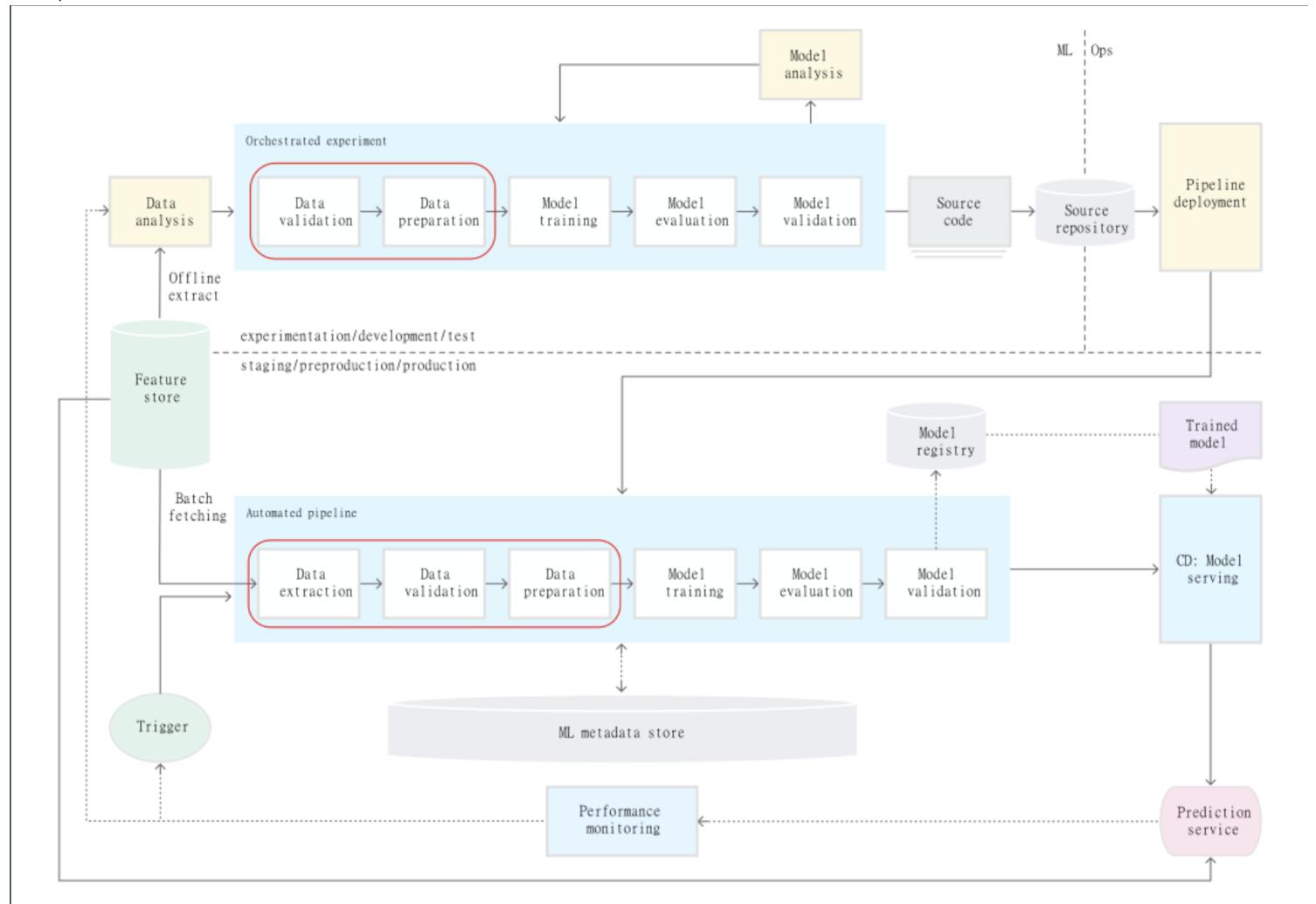


## Demo App

- MLFlow: <http://deplo-mlflo-16g0oyp6k65hv-5cd5faf094caf332.elb.us-east-1.amazonaws.com/#/>
- Model Endpoint: <https://spotify-recommendations-crun-ctreuw63uq-uw.a.run.app/invocations>

## Demo Video of Endpoint

## Recap



#TODO:

- Add Prometheus and Grafana metrics
- Replace with an actually good model
- Connect all components of my application (Front-end and Back-end)
- Use MLFlow Model Registry to add streamlined Blue/Green Deployments and Canary Deployments
- Use CDK to spin up the CloudFormation template to spin up the aforementioned CodePipeline Stack
- Fix codebase to package minimal requirements for a better served model container
- Port from custom containers to something like a Seldon deployment if possible

## Development

### Kedro Overview

This is your new Kedro project, which was generated using `Kedro 0.17.4`.

Take a look at the [Kedro documentation](#) to get started.

### Rules and guidelines

In order to get the best out of the template:

- Don't remove any lines from the `.gitignore` file we provide

- Make sure your results can be reproduced by following a [data engineering convention](#)
- Don't commit data to your repository
- Don't commit any credentials or your local configuration to your repository. Keep all your credentials and local configuration in `conf/local/`

## How to install dependencies

Declare any dependencies in `pyproject.toml` as per usage with `poetry`. You can also use the `src/requirements.txt` if you use the `kedro install` and `build-reqs` flow.

## How to run your Kedro pipeline

You can run your Kedro project with:

```
kedro run
```

## How to test your Kedro project

Have a look at the file `src/tests/test_run.py` for instructions on how to write your tests. You can run your tests as follows:

```
kedro test
```

To configure the coverage threshold, go to the `.coveragerc` file.

## Project dependencies

For the `poetry` workflow, to generate/update dependencies, use the `poetry add` functionality alongside adjusting the `pyproject.toml` to take advantages of `poetry`'s automatic dependency resolution

If you are using the `kedro install` workflow, to generate or update the dependency requirements for your project: `kedro build-reqs`

This will copy the contents of `src/requirements.txt` into a new file `src/requirements.in` which will be used as the source for `pip-compile`. You can see the output of the resolution by opening `src/requirements.txt`.

After this, if you'd like to update your project requirements, please update `src/requirements.in` and re-run `kedro build-reqs`.

### [Further information about project dependencies](#)

## How to work with Kedro and notebooks

Note: Using `kedro jupyter` or `kedro ipython` to run your notebook provides these variables in scope: `context`, `catalog`, and `startup_error`.

Jupyter, JupyterLab, and IPython are already included in the project requirements by default, so once you have run `kedro install` you will not need to take any extra steps before you use them.

### Jupyter

To use Jupyter notebooks in your Kedro project, you need to install Jupyter:

```
pip install jupyter
```

After installing Jupyter, you can start a local notebook server:

```
kedro jupyter notebook
```

## JupyterLab

To use JupyterLab, you need to install it:

```
pip install jupyterlab
```

You can also start JupyterLab:

```
kedro jupyter lab
```

## IPython

And if you want to run an IPython session:

```
kedro ipython
```

## How to convert notebook cells to nodes in a Kedro project

You can move notebook code over into a Kedro project structure using a mixture of [cell tagging](#) and Kedro CLI commands.

By adding the `node` tag to a cell and running the command below, the cell's source code will be copied over to a Python file within `src/<package_name>/nodes/`:

```
kedro jupyter convert <filepath_to_my_notebook>
```

**Note:** The name of the Python file matches the name of the original notebook.

Alternatively, you may want to transform all your notebooks in one go. Run the following command to convert all notebook files found in the project root directory and under any of its sub-folders:

```
kedro jupyter convert --all
```

## How to ignore notebook output cells in git

To automatically strip out all output cell contents before committing to `git`, you can run `kedro activate-nbstripout`. This will add a hook in `.git/config` which will run `nbstripout` before anything is committed to `git`.

**Note:** Your output cells will be retained locally.

## Package your Kedro project

[Further information about building project documentation and packaging your project](#)