# ML Pipelines, Reproducibility and Experimentation

Alex Kim @alex000kim

# We have a Jupyter notebook...
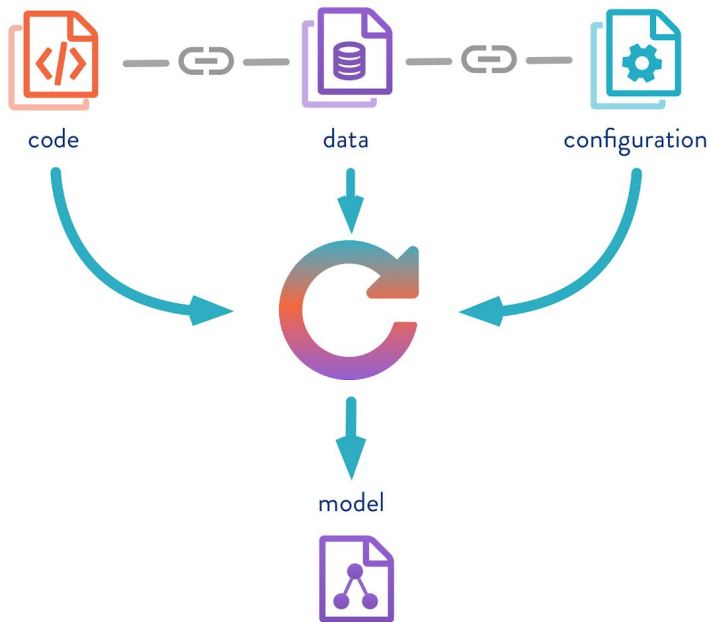
◇ Data loading

◇ Feature Engineering

◇ Model Training

◇ Model Evaluation

# Can you easily answer these questions?

- What exactly was used to produce a particular model?

- Can you easily compare many ML experiments?

- Will you be able to reproduce them later?

# Goal #1: Achieve best performance

- Running many experiments
- **Experiment** = a particular combination of **Code** & **Data** & **Config**



code — data — configuration

model

# Goal #2: Ensure reproducibility

- **Improving model performance**: you can't improve what you can't reproduce
- **Transparency and team collaboration**: know everything your team members did to achieve certain performance
- **Auditability (laws and regulations)**: e.g. what *exactly* went into the models that prescribes treatment to patients or determines creditworthiness of bank customers

# Goal #3: Minimal setup and dependency of 3rd party services

Problems:

1. **Vendor lock-in**: instrument code with framework-specific code

2. **Maintenance & cost**: maintain your own ML tracking server (or pay them to take care of it)

3. **Security concerns**: send data to an external service or database

Most ML tracking solutions (MLflow, W&B, comet.ml, etc) have at least 2 of these problems

# Fact:
# It's difficult to achieve all three goals

Can we do all of the following?

1. Iterate quickly i.e. generate many experiments
2. Automatically track **all** changes to code, configs and data
3. Avoid dependency on 3rd party services to store data, metrics and params

# Same experiments, but different metrics? 🤔🤔🤔

# Reproducibility VS. Experimentation?



WHY NOT BOTH?

- DVC pipelines for generating many experiments
- Achieve complete reproducibility by versioning **everything**!
  - code and configs --> Git
  - dataset, models, other artifacts --> DVC remote storage (cloud buckets, NAS, SFTP, etc)
- VS Code as a convenient UI for experiment management
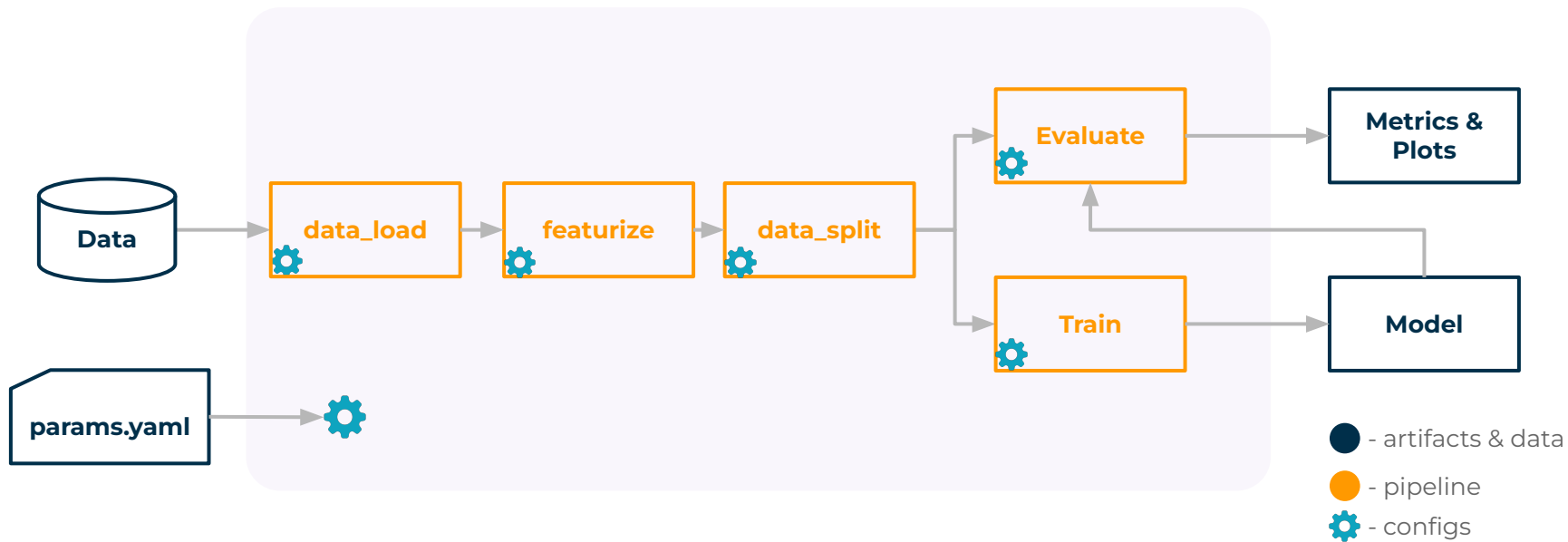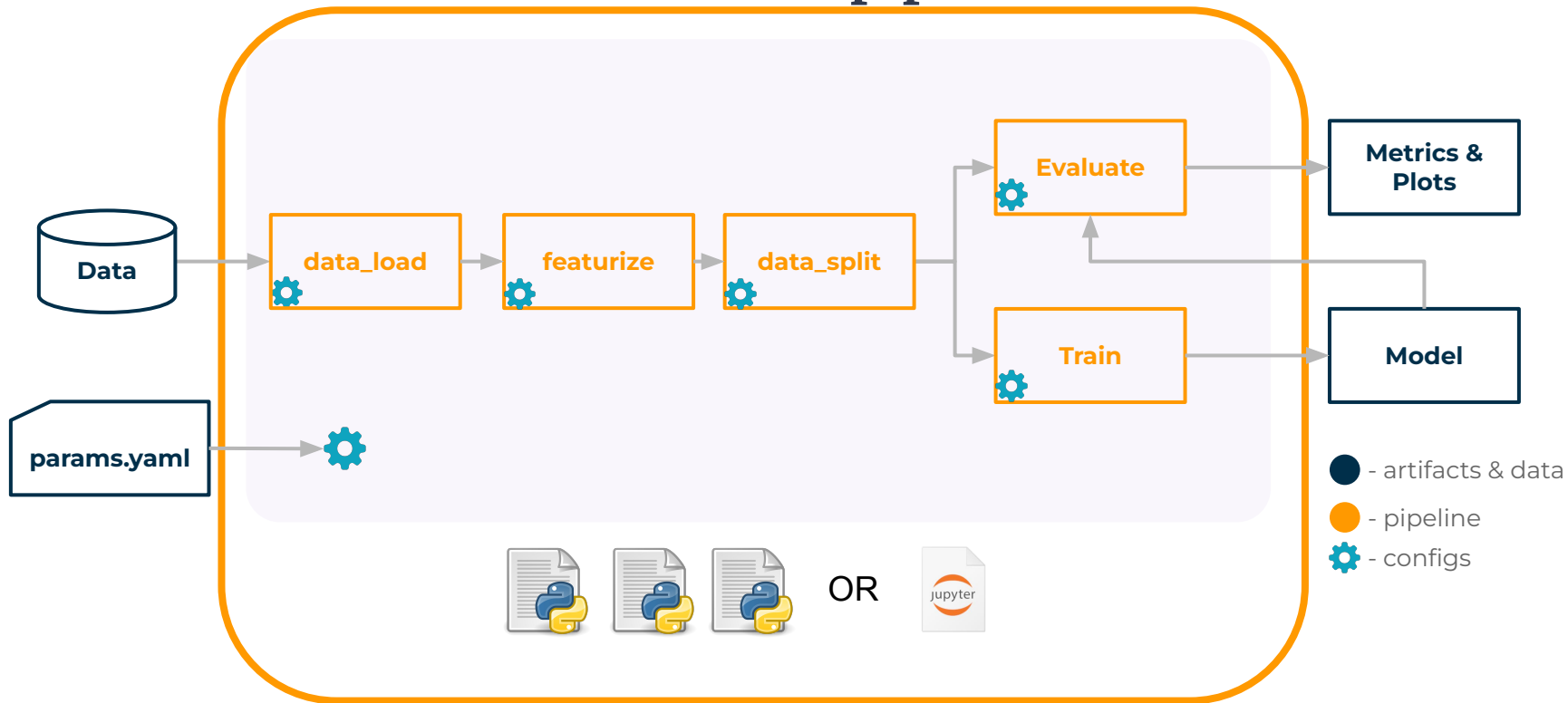- No need to maintain (or pay for) additional services



Visual Studio Code



git



DVC

# What are DVC pipelines?

# What are DVC pipelines?

# DVC pipeline (defined in dvc.yaml) as:

**a sequence of Python modules**

```
stages:
  data_preprocessing_stage:
    cmd: python process_data.py
    deps:
      - process_data.py
      - path/to/raw/data
    outs:
      - path/to/processed/data
    params:
      - preprocessing_params
  train_stage:
    cmd: python train.py
    deps:
      - train.py
      - path/to/processed/data
    outs:
      - my_model.pickle
    params:
      - train_params
  eval_stage:
    cmd: python eval.py
    deps:
      - eval.py
      - path/to/processed/data
      - my_model.pickle
    params:
      - eval_params
```

**a Jupyter notebook**

```
stages:
  run_notebook_stage:
    cmd: papermill MyNotebook.ipynb MyNotebook_out.ipynb
         -p n_estimators ${n_estimators}
         -p max_depth ${max_depth}
    deps:
      - MyNotebook.ipynb
      - path/to/raw/data
    outs:
      - my_model.pickle
```

# Run an experiment



```
$ dvc exp run -S train.params.n_estimators=120
'data/Churn_Modelling.csv.dvc' didn't change, skipping

Running stage 'run_notebook':
> papermill TrainChurnModel.ipynb TrainChurnModel_out.ipynb -p n_estimators 120 -p max_depth 10 -p model_type lightgbm
Input Notebook:  TrainChurnModel.ipynb
Output Notebook: TrainChurnModel_out.ipynb
Black is not installed, parameters wont be formatted
Executing:    0%|                                                                                          | 0/27
[00:00<?, ?cell/s]Executing notebook with kernel: python3
Executing: 100%|██████████████████████████████████████████████████████████████████████████████████████████| 27/27
[00:03<00:00,  7.58cell/s]
Updating lock file 'dvc.lock'
```

# Track and manage many experiments

# Practice time!