

DVC

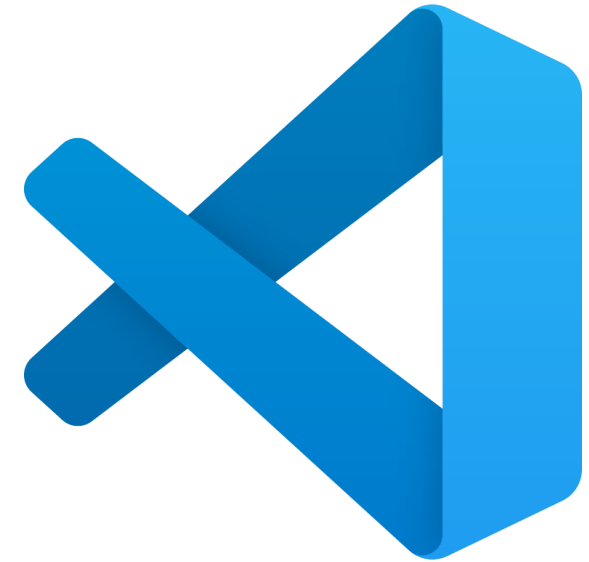
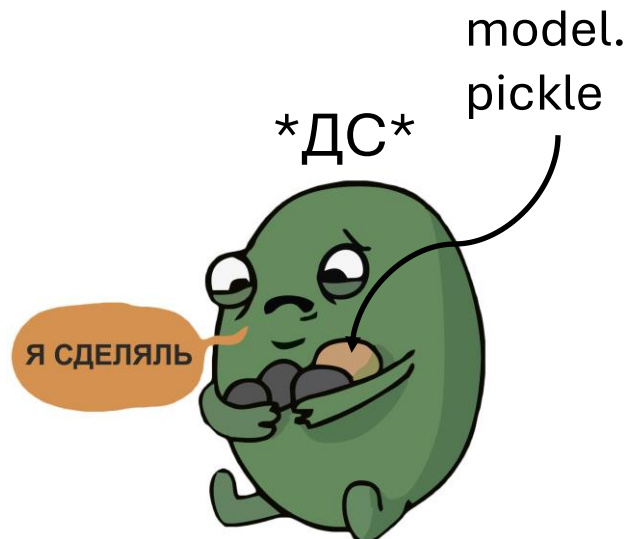
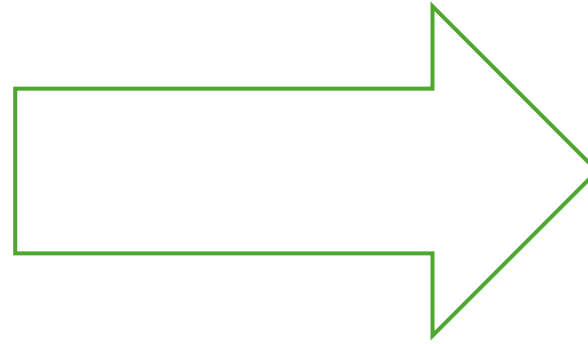
Лекция №4

Типичная картина: Ресар



Код построения модели:

- не воспроизводим
- не обобщаемый
- не универсальный



Код использования модели:

- хороший
- воспроизводимый

3 этапа в жизни модели

Дс

1) Обучение

2)



3) Переобучение

Бэкенд

1)

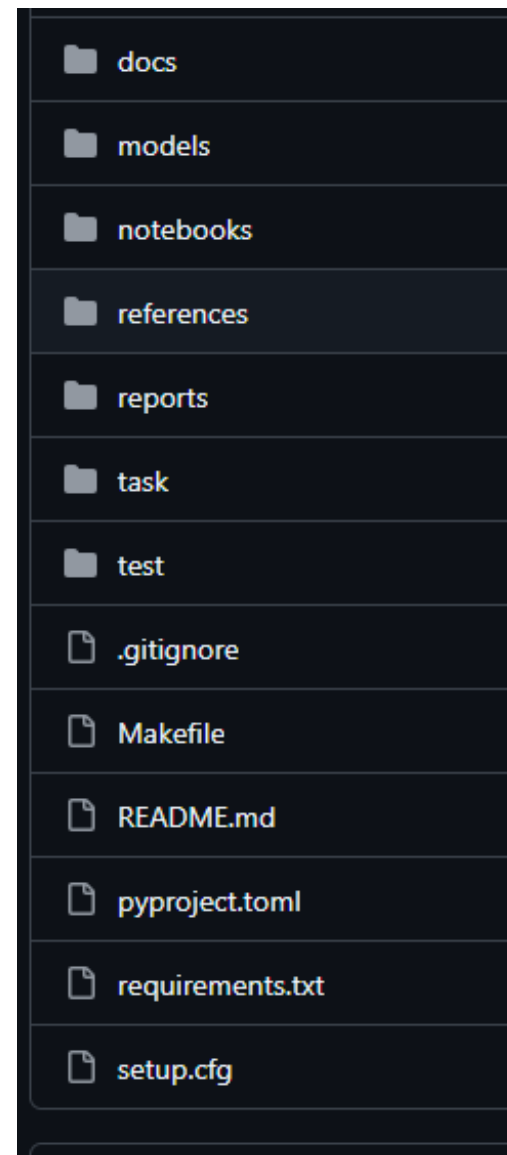
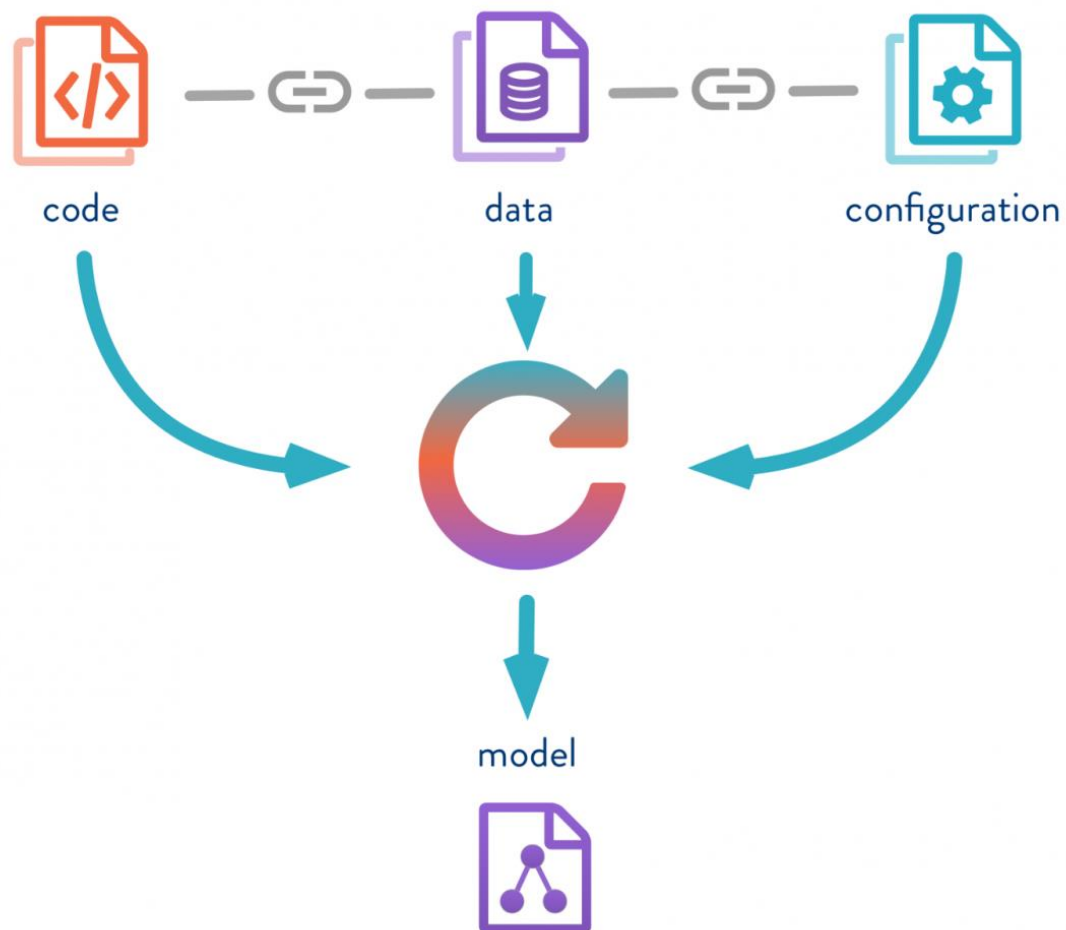


2) Выкатка и эксплуатация

3)



Рецепт МЛ модели



Как хранить данные?

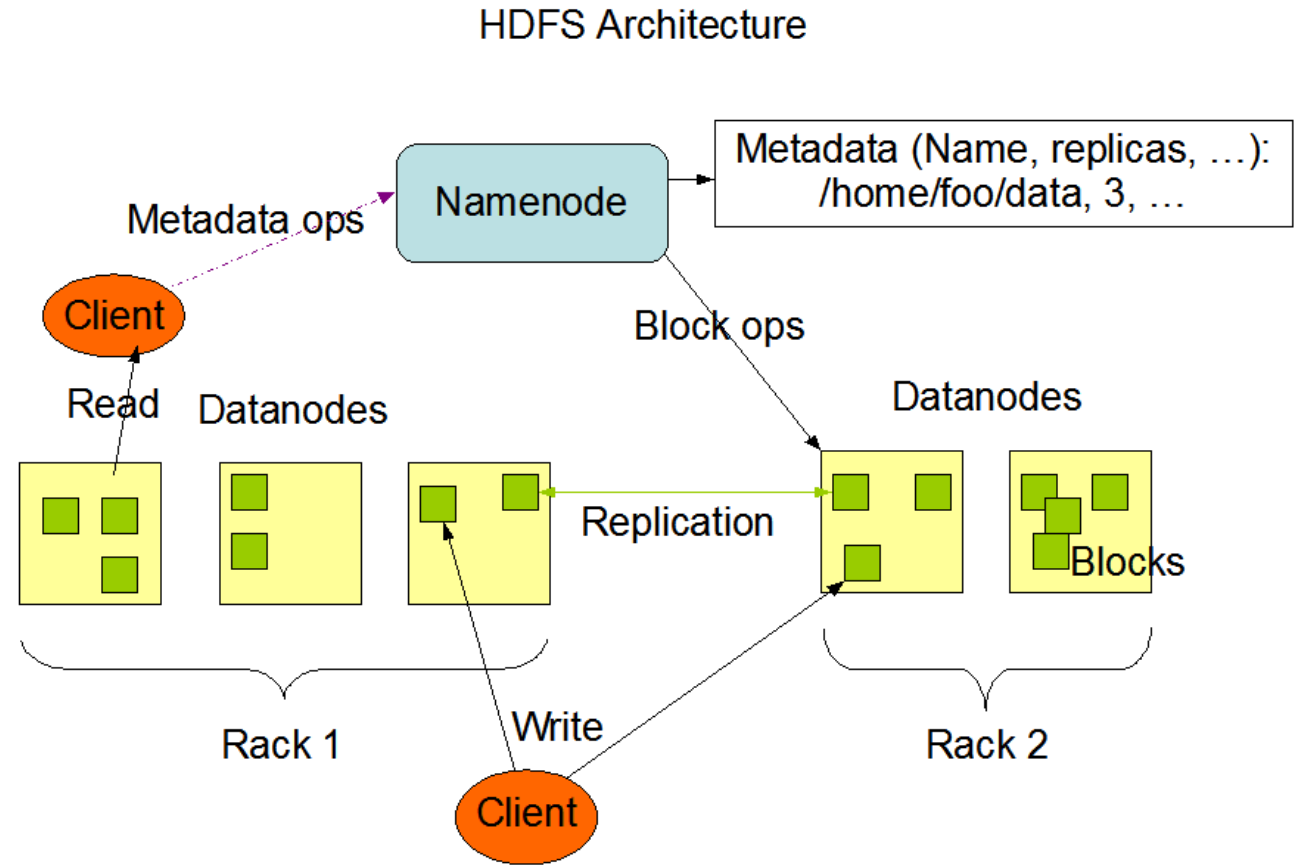
File size limits

GitHub limits the size of files allowed in repositories. If you attempt to add or update a file that is larger than 50 MiB, you will receive a warning from Git. The changes will still successfully push to your repository, but you can consider removing the commit to minimize performance impact. For more information, see "[Removing files from a repository's history.](#)"

Note: If you add a file to a repository via a browser, the file can be no larger than 25 MiB. For more information, see "[Adding a file to a repository.](#)"

GitHub blocks files larger than 100 MiB.

A cartoon illustration of a yellow elephant with a large trunk, smiling and waving its right hand. The elephant is standing on four legs and has a small tusk-like protrusion on its head.



Куда еще складывать данные?



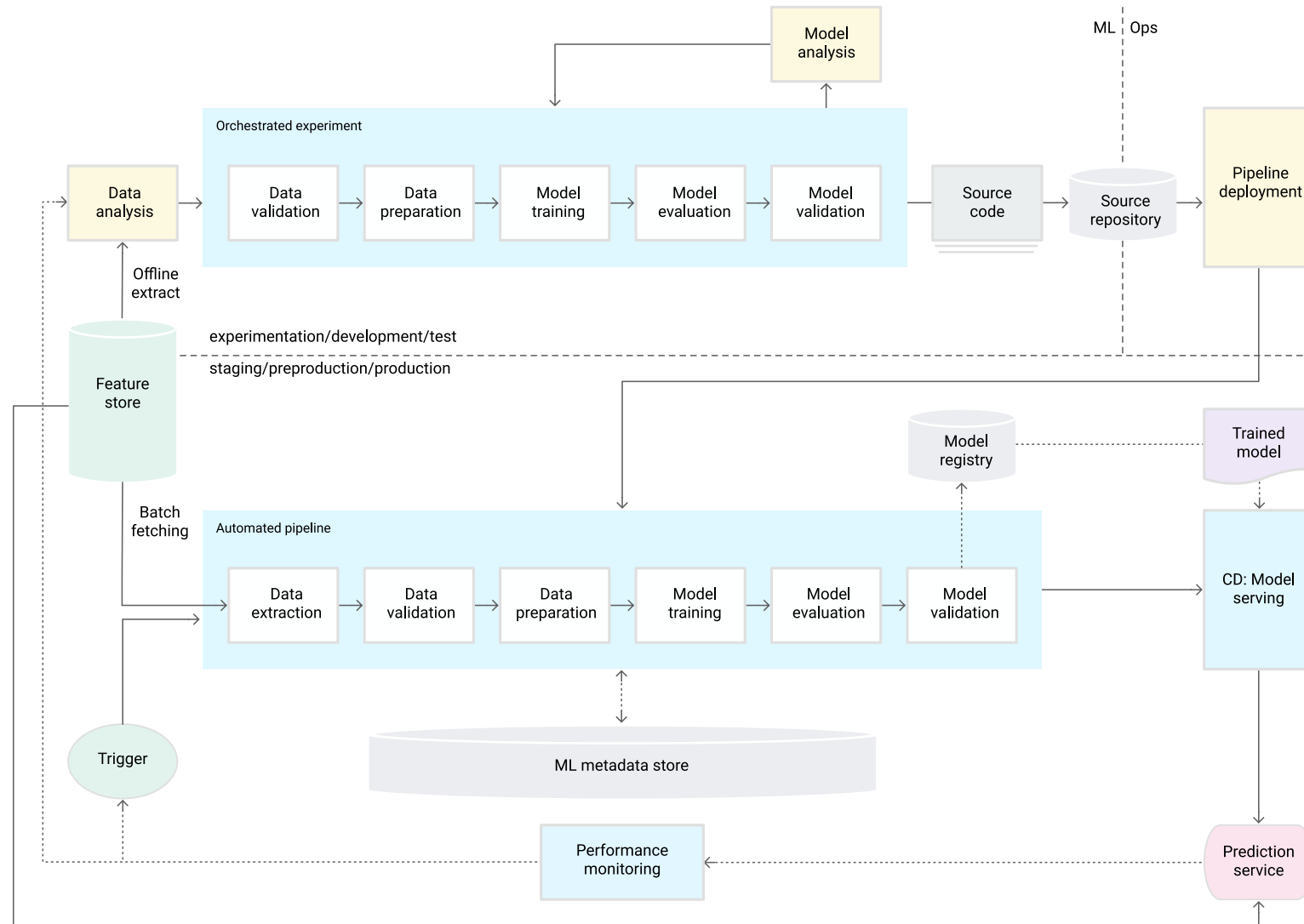
Что необходимо знать для воспроизводимого пайплайна

Сохраняем:

1. Файл модели
2. Ссылку на git репо и коммит (либо номер релиза)
3. Параметры модели
4. Метрики

Называем папку `model_ver_1` и складываем на хранилище

MLOps Level 1



DVC

Data Version Control

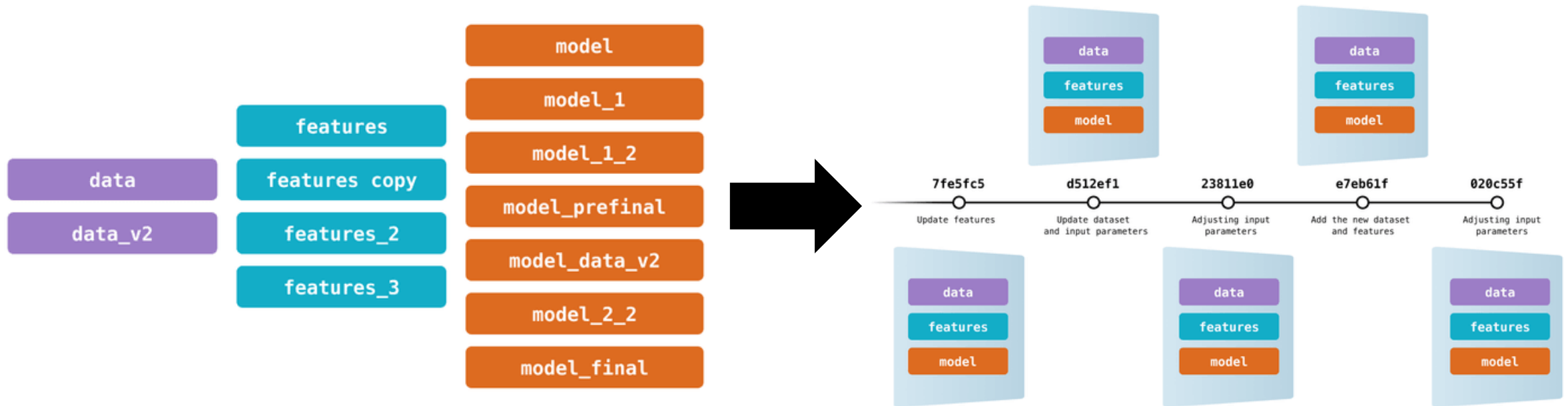
– *and much more* –
for the **GenAI** era

Free and open source, forever.

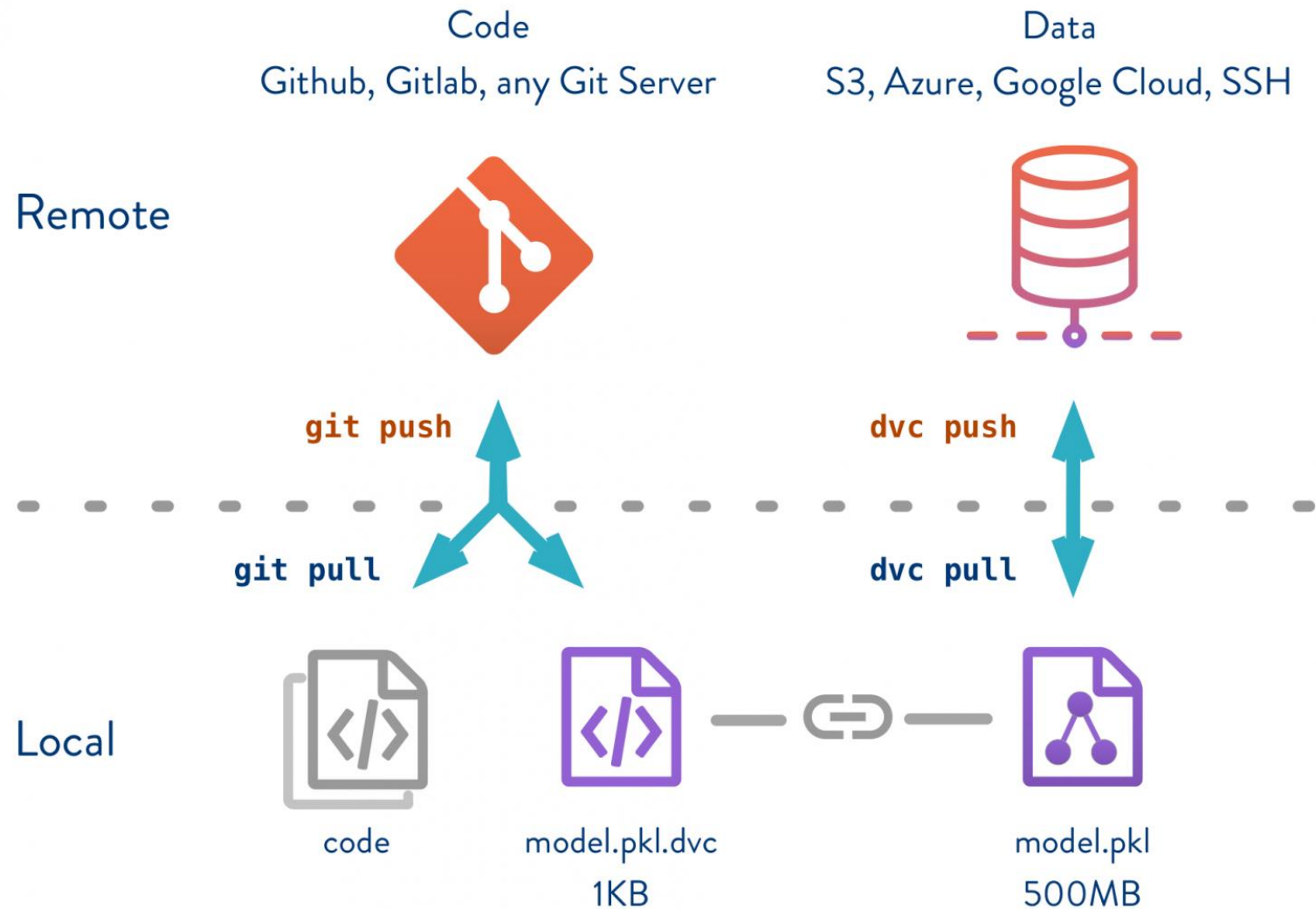
Manage and version images, audio, video, and text files in storage and organize your ML modeling process into a reproducible workflow.

<https://dvc.org>

Версионирование данных



Версионирование данных



Версионирование данных

```
$ git init
```

```
$ dvc init
```

```
$ dvc add file.txt
```

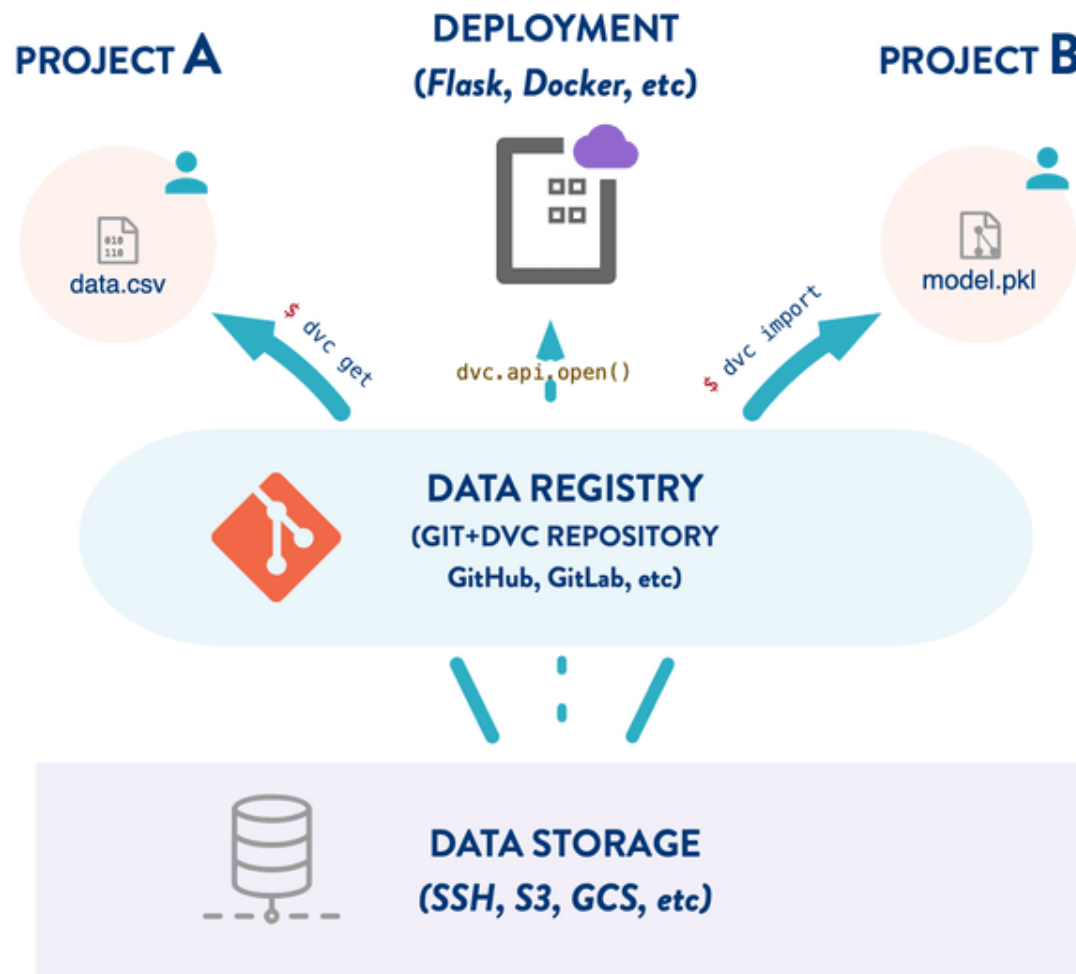
```
$ git add file.txt.dvc
```

```
$ dvc push
```

```
$ git commit -m <comment>
```

```
$ git push
```

Версионирование данных



Выбор хранилища

```
$ dvc remote add -d temp /tmp/dvcstore
```

```
# .dvc/config
[remote "temp"]
    url = /tmp/dvcstore
[core]
    remote = myremote
```

! .dvc/config надо пушить в гит
!! всю папку .dvc пушить нельзя

Supported storage types

The following are the supported types of storage protocols and platforms.

Cloud providers

- [Amazon S3](#) (AWS) and [S3-compatible](#) e.g. MinIO
- Microsoft [Azure Blob Storage](#)
- [Google Cloud Storage](#) (GCP)
- [Google Drive](#)
- [Aliyun OSS](#)

Self-hosted / On-premises

- [SSH](#); Like `scp`
- [HDFS](#) & [WebHDFS](#)
- [HTTP](#)
- [WebDAV](#)

Пайплайны вычислений

```
$ dvc stage add -n train \  
    -p train.seed,train.n_est,train.min_split \  
    -d src/train.py -d data/features \  
    -o model.pkl \  
    python src/train.py data/features model.pkl
```

```
$ dvc dag  
  
+-----+  
| prepare |  
+-----+  
      *  
      *  
      *  
+-----+  
| featurize |  
+-----+  
      *  
      *  
      *  
+-----+  
| train |  
+-----+
```


Файлы: dvc.yaml

Содержит:

1. Информация об этапах пайплайна
2. Параметры запуска
3. Описание артефактов
4. Описание метрик

Поддерживает:

- foreach
- matrix

<https://dvc.org/doc/user-guide/project-structure/dvcyaml-files>

```
! dvc.yaml
1  params:
2    - params.yaml
3  stages:
4    make_dataset:
5      cmd: 'python src/dataset.py params.yaml'
6      deps:
7        - src/dataset.py
8      params:
9        - data_params
10     outs:
11       - ${data_params.train_data_path}
12       - ${data_params.test_data_path}
13    train_model:
14      cmd: 'python src/modeling/train.py params.yaml'
15      deps:
16        - src/modeling/train.py
17        - ${data_params.train_data_path}
18        - ${data_params.test_data_path}
19      params:
20        - train_params
21      outs:
22        - ${train_params.model_path}
23      metrics:
24        - ${train_params.metrics_path}
25
```

Домашнее задание №3

1. Создать локальное хранилище для DVC
2. Переписать полную тренировку модели через пайплайн DVC
3. Запустить необходимые файлы .dvc в гит

Не допускается никакого хардкода

Все, что можно сконфигурировать, должно лежать в конфигах (но без фанатизма)