
Model pipelines

Sprint 4 - Week 8

INFO 9023 - Machine Learning Systems Design

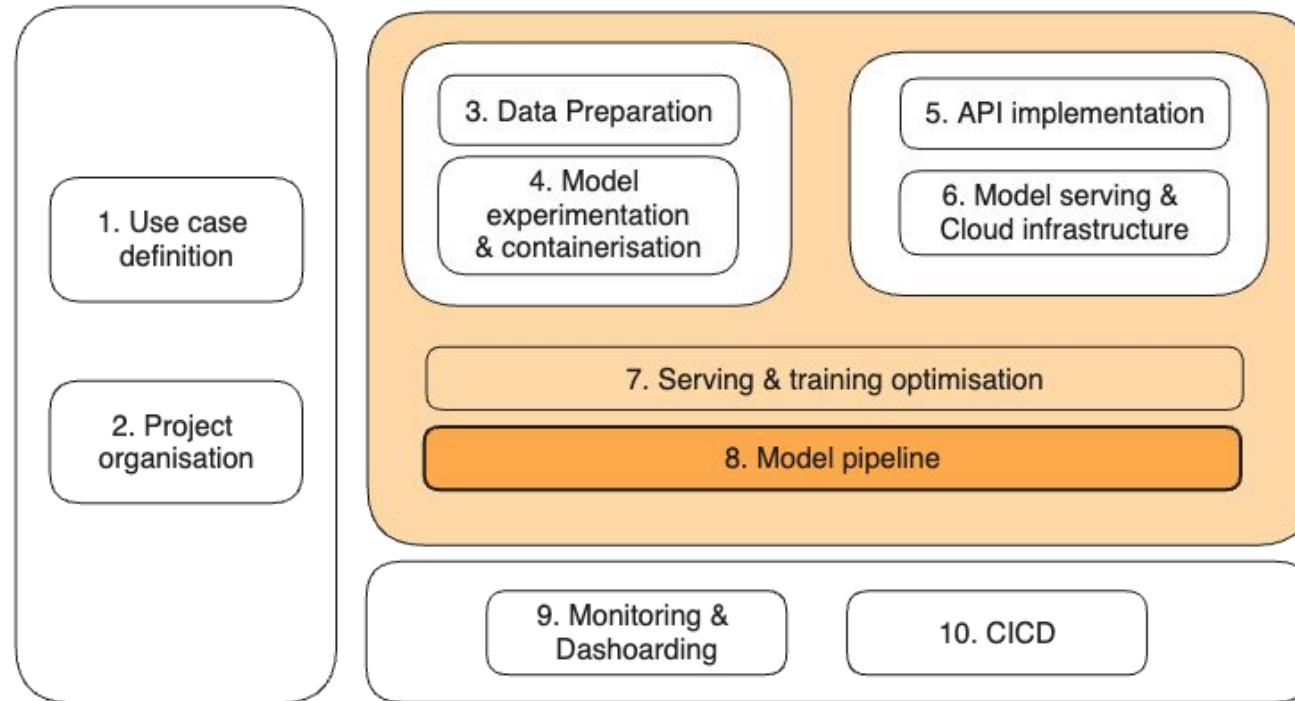
2024 H1

Thomas Vrancken (t.vrancken@uliege.be)
Matthias Pirlet (matthias.pirlet@uliege.be)

Status on our overall course roadmap



Status on our overall course roadmap



Agenda

What will we talk about today

Lecture (45min)

1. Microservice architecture
 - a. What is a microservice
 - b. Abstractions & Separation of concerns
 - c. The anatomy of a microservice
 - d. Monoliths aren't for dinosaurs

Lab (45 min)

2. Docker compose

Lecture (30min)

3. ML model pipeline
4. ML platforms & orchestrators

Project update

⚠️ Decision to **combine Milestone 2 and 3** ! ⚠️

There will be only one presentation, taking place in our last **course** (13/05/2024) !

How? (Info is also on [github](#))

- Same content as MS 3: Present your overall project in 10 min presentation + 5 min QA
- It will be in classroom ! You're very welcomed to attend other project presentations
- Send codes (link to pull request) by email **before the presentation**
- Agenda for the presentations is in this [sheet](#) (let us know if it does not work for you)
- Topics from the last "bonus" lecture are going to be spread over other lectures

Why?

-  Nicer to let everyone attend all presentations
-  Simpler organisation as there is already a class block in your agenda, don't take more of your time
-  Avoid having to repeat topics between

Microservice architecture

Credits where credits are due



Robbe Sneyders
Principal Engineer @ ML6

Ruwan Lambrichts
Senior ML Engineer @ ML6

Topic of this section

Defining a general software architecture for microservices in your project/organisation.

Topic of this section

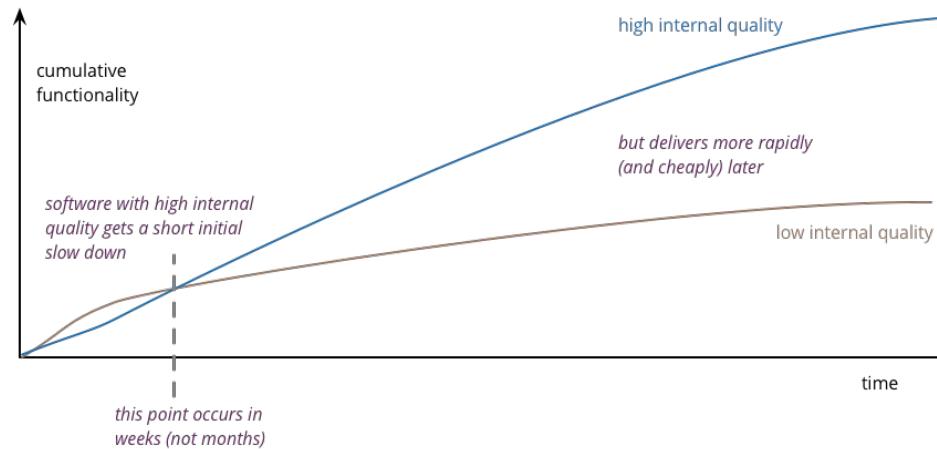
*Defining a **general** software architecture for microservices in your project/organisation.*

- The shared understanding that the expert developers have of the system design.
- The decisions you wish you could get right early in the project.

Why does software architecture matter?

Software architecture enhances the **quality** of applications developed by your team, which is something stakeholders don't immediately perceive (and which might seem less important)

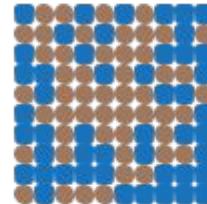
However, poor architecture is a major contributor to **cruft** - leading to code that is much harder to modify, causing features to arrive more slowly and with more defects (which the client cares about a lot)



Cruft causes features to arrive more slow and with more defects

High quality software architecture leads to faster delivery of new features (with less chance of making a mistake), because there is less cruft to get in the way.

If we compare one system with a lot of cruft...

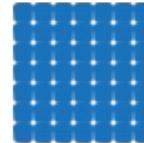


the cruft means new features take longer to build

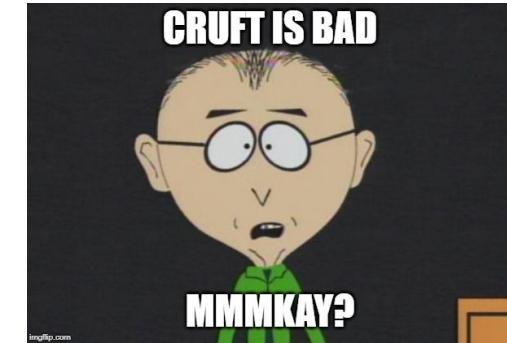


this extra time and effort is the cost of the cruft, paid with each new feature

...to an equivalent one without



free of cruft, features can be added more quickly



Ideal state: Your team uses/builds boilerplates.

Creating a shared understanding of the microservice design .

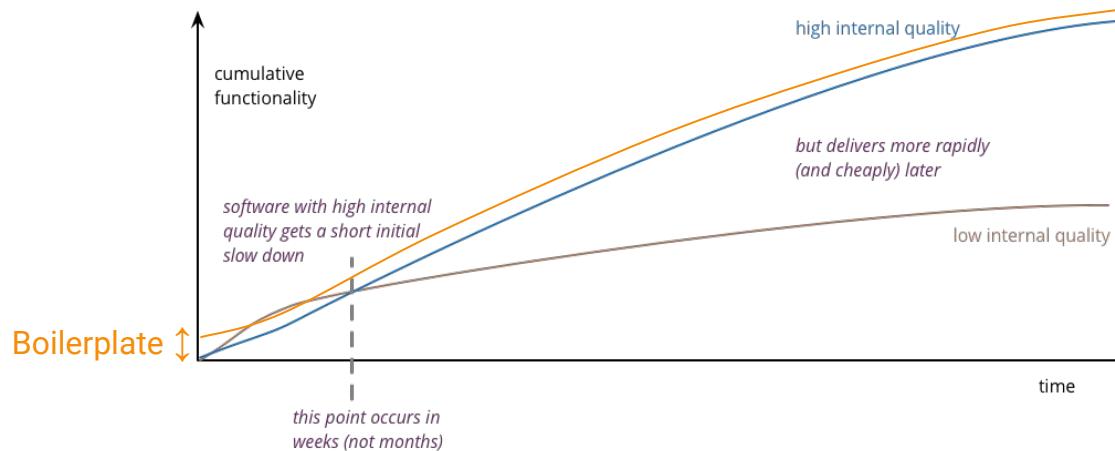


Captured in the **microservice boilerplate**, so every project can start from the tried and tested decisions already made on previous projects.

Ideal state: Your team uses/builds boilerplates.

Take 'correct' decisions based on combined knowledge of developers and projects, implemented by the microservice boilerplate

- Saves time at start of project
- Projects don't run into problems down the line, can evolve fast, and stay maintainable
- Consistent design across projects facilitates understanding of other projects



What is a microservice?

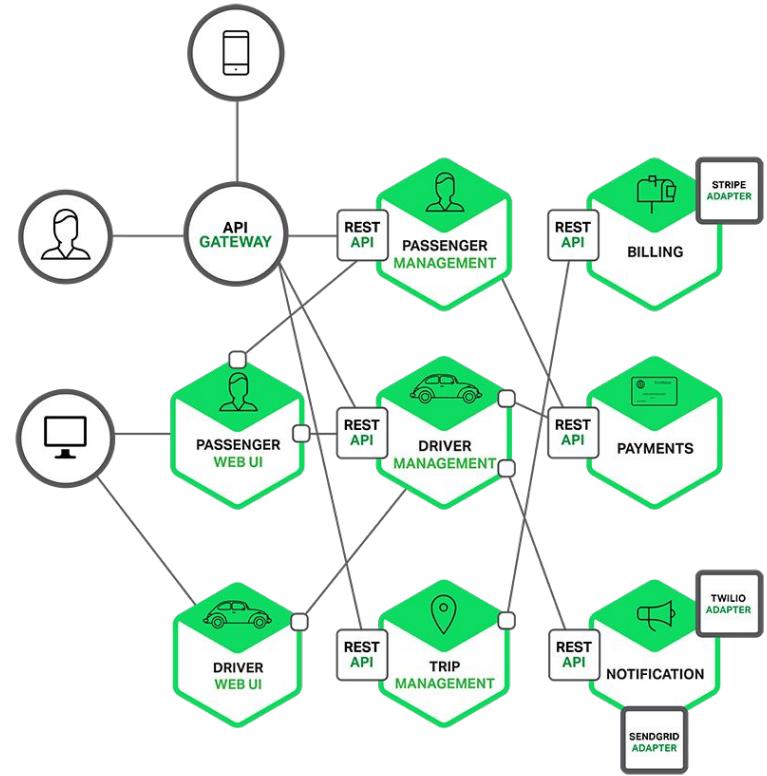
Some context

What is a microservice?

*Using a microservices approach, software is composed of **small services** that communicate over **well-defined APIs** that can be **deployed independently**. These services are owned by small autonomous teams.*

What is a microservice?

- Applying single responsibility principle at the architectural level
- As small as possible, but as big as necessary
 - “No bigger than your head”
- Advantages compared to monolithic architectures:
 - Independently deployable
 - Language, platform and technology independency between components
 - Distinct axes of scalability
 - Provides firm module boundary
 - Increased architectural flexibility
- Often integrated using REST over HTTP
 - Business domain concepts are modelled as resources

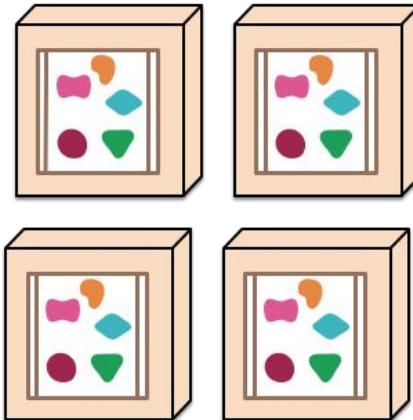


Monolith vs microservice

A monolithic application puts all its functionality into a single process...



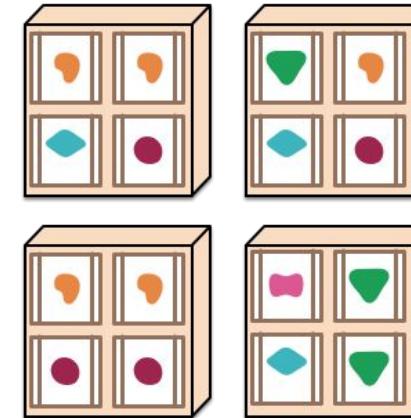
... and scales by replicating the monolith on multiple servers



A microservices architecture puts each element of functionality into a separate service...



... and scales by distributing these services across servers, replicating as needed.



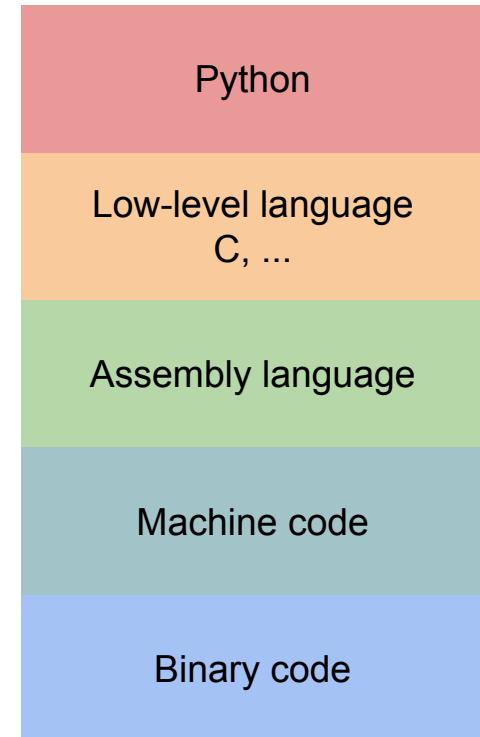
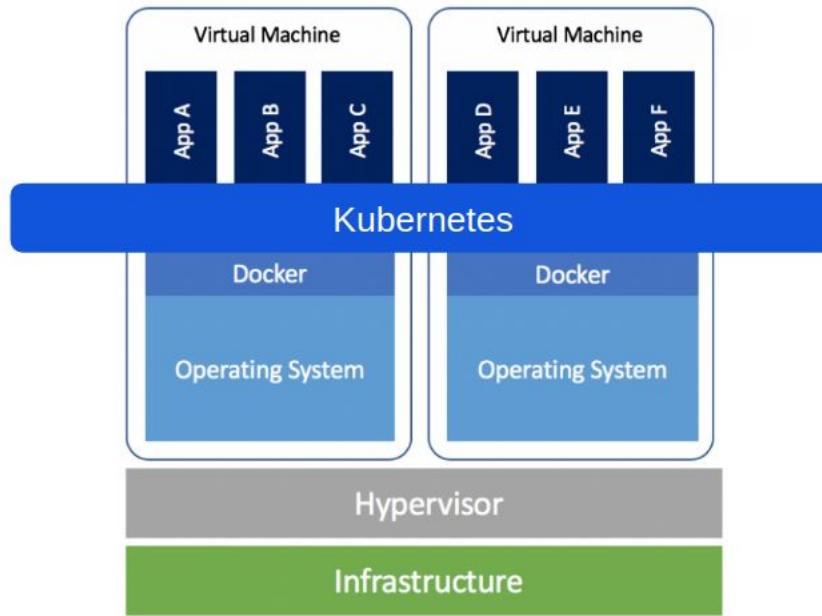
Abstractions & Separation of concerns

Abstraction

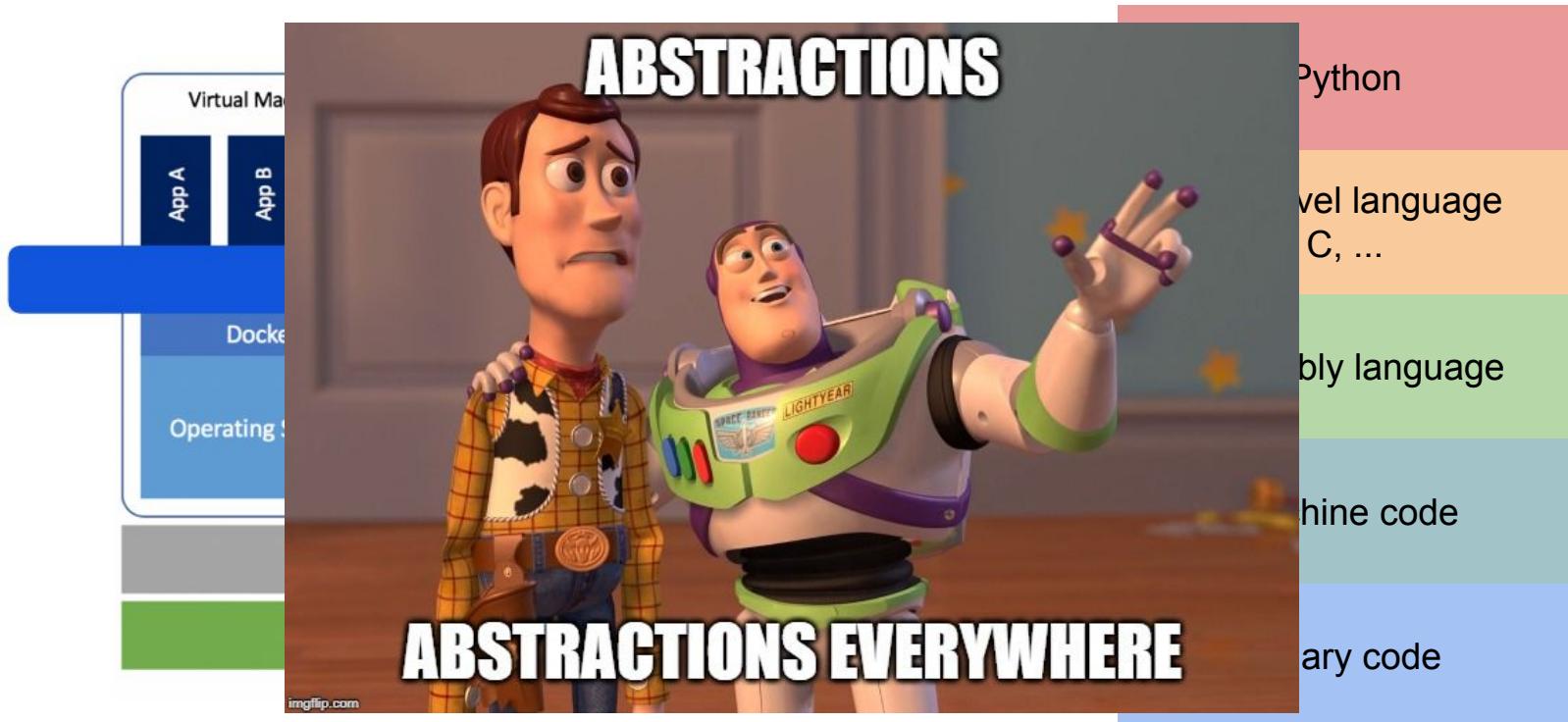
- Abstraction is the act of **representing essential features without including the background details or explanations.**
- Reduce complexity and allow efficient design and implementation of complex software systems.



Abstractions are everywhere



Abstractions are everywhere



Separation of concerns

- Separate code into sections that each address a separate concern
- Sections can evolve independently of other components
- Separation of concerns is a form of abstraction



Images: <https://pusher.com/tutorials/clean-architecture-introduction>

Interfaces

- Describes the public actions and attributes that a type of component supports and implements.
- Defines the boundary across which two components can communicate.
- Components with an identical interface are interchangeable.



Interfaces

- In traditional OOP languages, an Interface defines a contract which is explicitly implemented.
- In Python, you don't have to explicitly declare an interface.

```
interface Pet {  
    public Owner owner();  
    public void pat();  
}  
  
class Dog implements Pet {  
  
    public Owner owner() {  
        return this.get_owner();  
    }  
  
    public void pat() {  
        return love;  
    }  
}
```

```
class Dog:  
  
    @property  
    def owner(self):  
        return self._get_owner()  
  
    def pat(self):  
        return love  
  
dog = Dog()  
try:  
    dog.owner()  
    dog.pat()  
    print("dog is a pet")  
except:  
    print("dog is not a pet")
```



Clean architecture

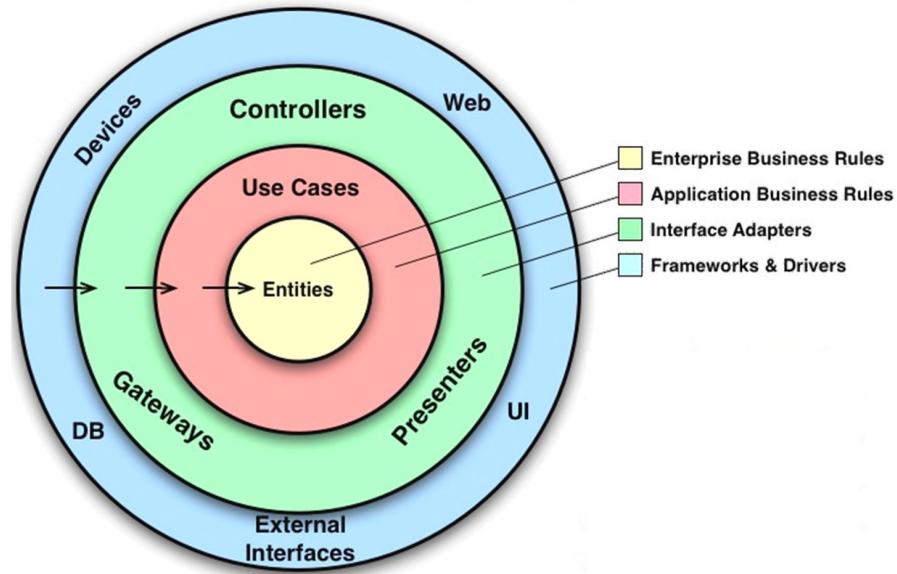


The Clean architecture

By Robert C. Martin (Uncle Bob)

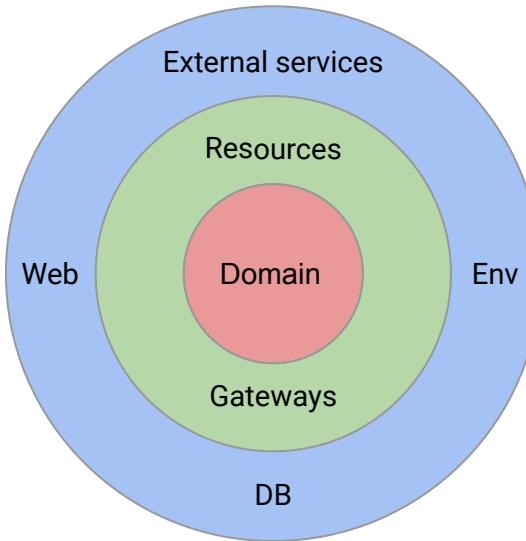
- Clear separation of concerns between components.
- Dependency rule:
Dependencies can only point inwards. Inner circle should not be dependent of outer circle specifics.
- Interfaces between each layer.
- Only isolated, simple data structures are passed across the boundaries.

Very OOP focused



https://en.wikipedia.org/wiki/Robert_C._Martin

The Clean architecture - tailored



Infrastructure:

- Frameworks, environment, and storage and web access.

Interface adapters:

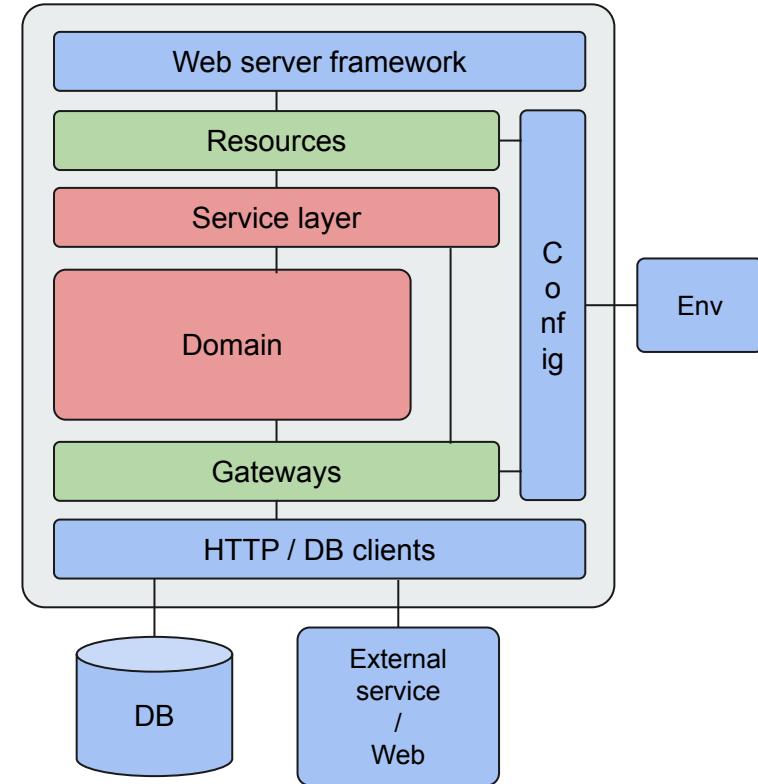
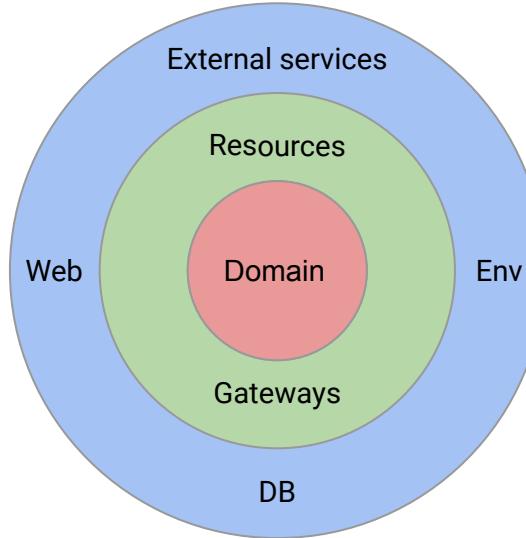
- Encapsulate the infrastructure layer, and provide a Python API for the domain.

Domain:

- All logic core to the microservice.

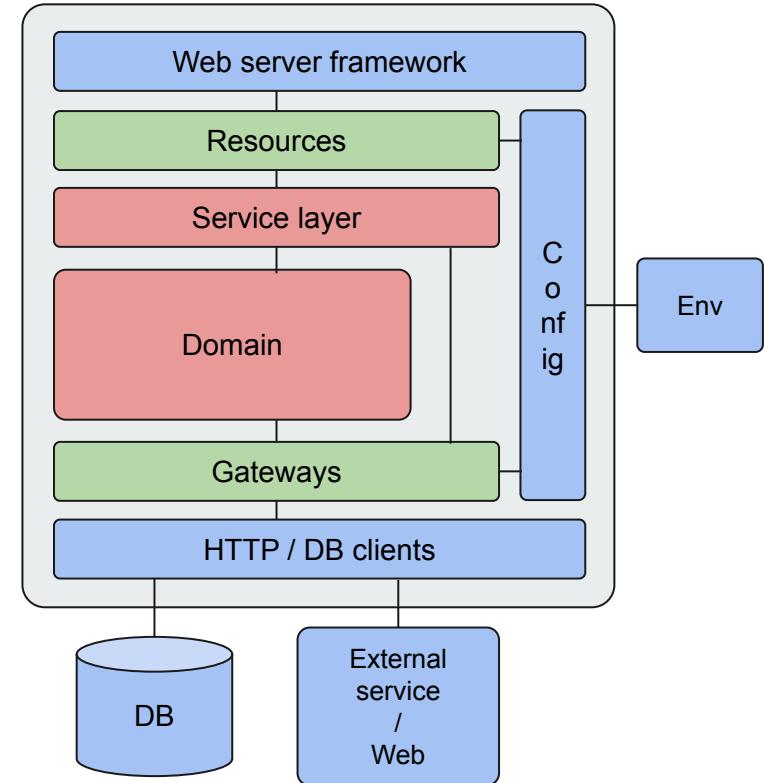
The anatomy of a microservice

The clean anatomy of a microservice



The anatomy of a microservice

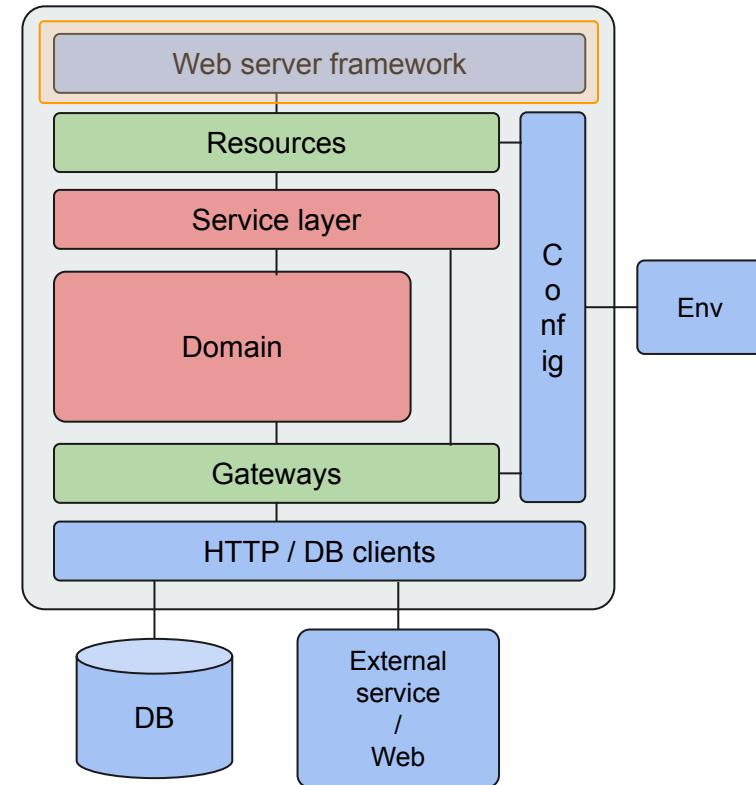
- Microservices display consistent internal structure composed of some or all of the displayed layers.
- The microservice boilerplate aims to provide minimal viable structure that is shared by all microservices.
- Examples will be added containing additional components.



The anatomy of a microservice

Web server framework

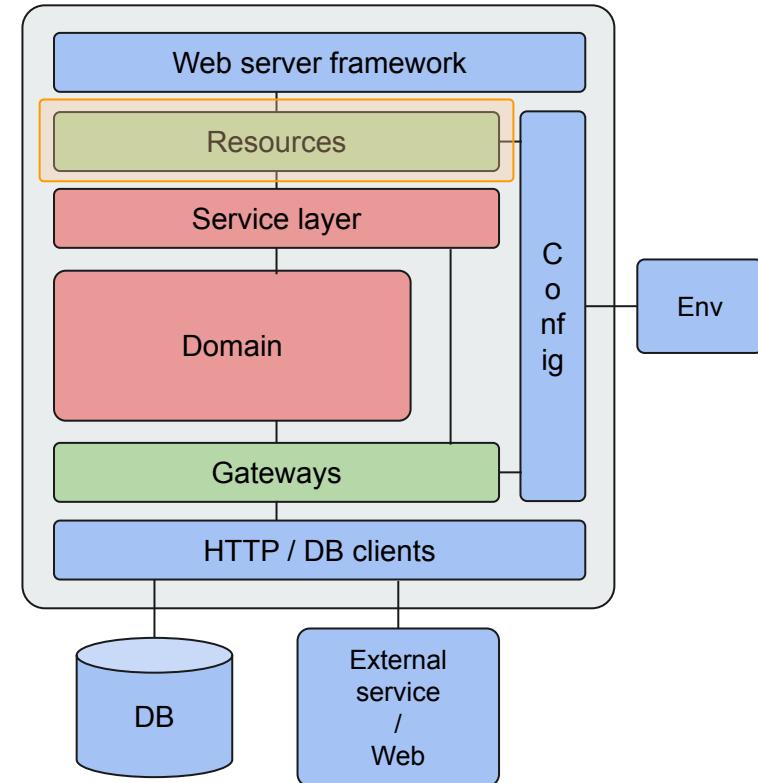
- Act as a mapping between the exposed application protocol (eg. HTTP) and Python.
- Validate requests and responses against application protocol.



The anatomy of a microservice

Resources

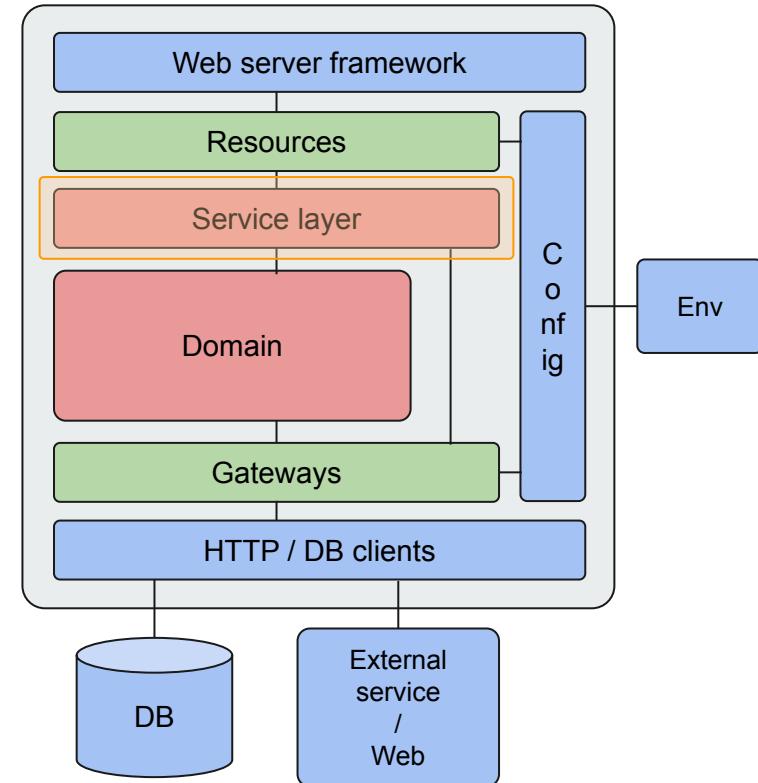
- Act as a mapping between the web server framework and the service domain.
- Thin layer for sanity checking requests and providing protocol specific responses.



The anatomy of a microservice

Service layer

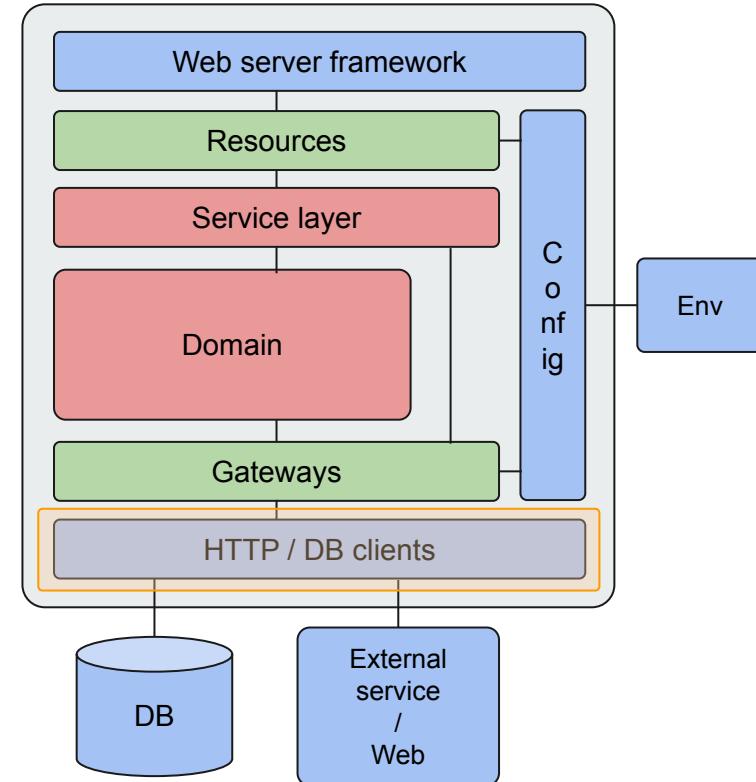
- Defines the interface of the domain with its set of available operations.
- Can be reused by multiple protocol clients; HTTP, GRPC, Queue reader, ...



The anatomy of a microservice

HTTP / DB clients

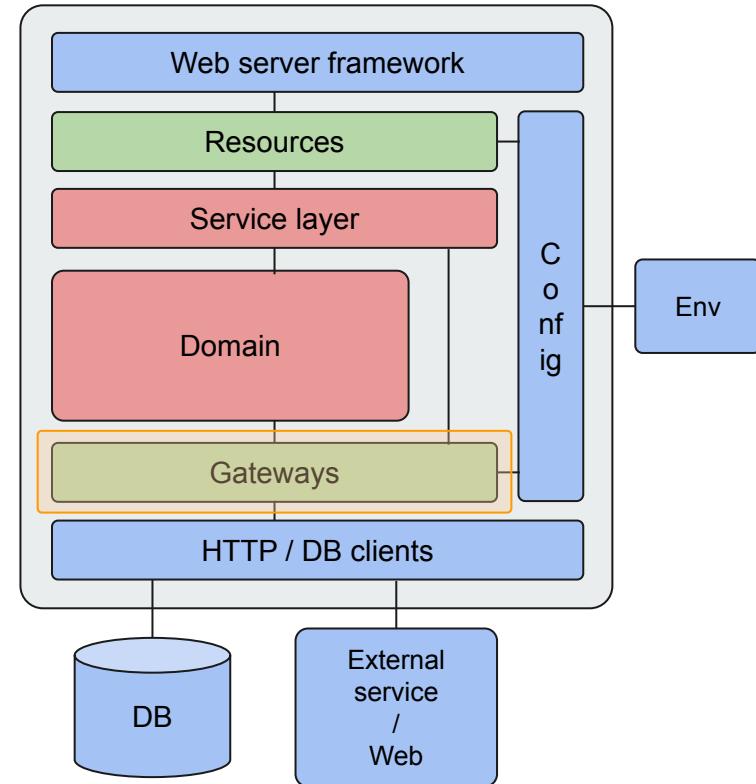
- Clients to handle connections to external datastores and services
- HTTP-client:
Client that understands the underlying protocol to handle the request-response cycle.
(eg. python requests)
- DB-client:
Provide (often client-specific) python representation of data in datastores



The anatomy of a microservice

Gateways

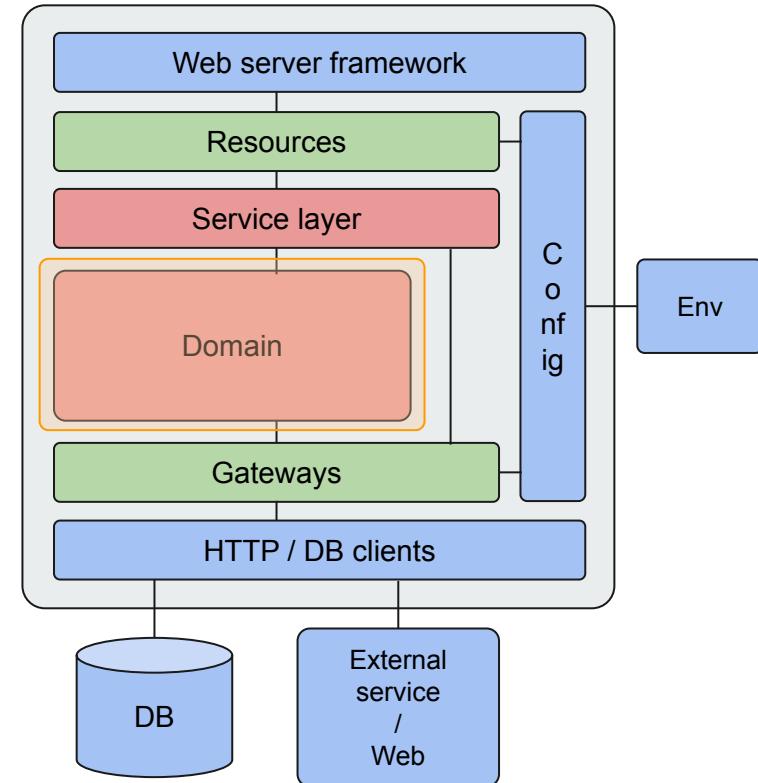
- Isolates domain from details in the database access code.
- Groups query construction code to minimize duplicate query logic.
- Encapsulates message passing with a remote service, marshalling requests and responses from and to Python.



The anatomy of a microservice

Domain

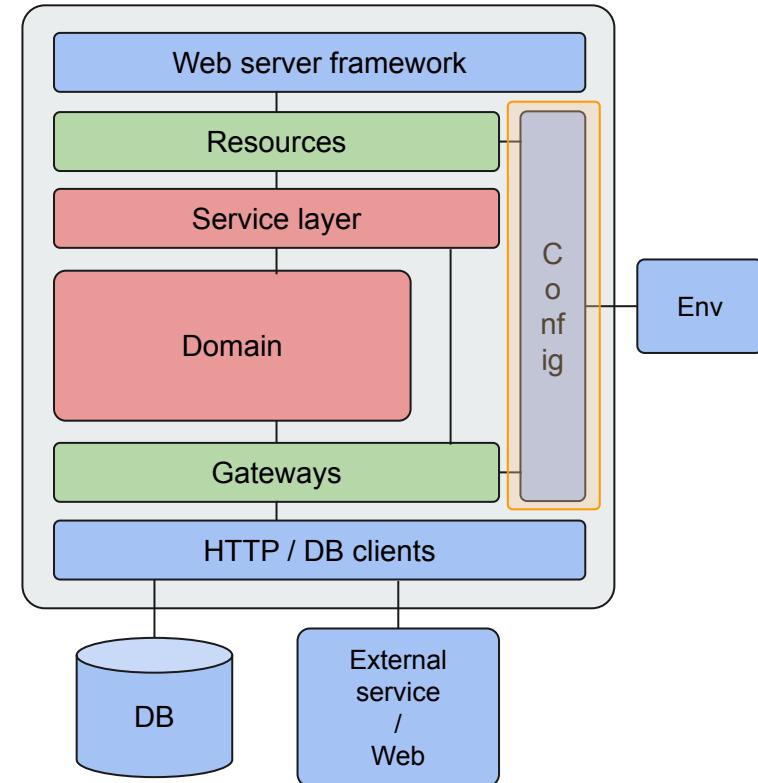
- Contains the actual service logic.
- The domain consists of pure Python with some core libraries. It is independent of any used frameworks and IO clients, and only depends on the interface of the resources and gateways.



The anatomy of a microservice

Config

- Reads and groups the configuration of the service.
- All configurable parameters of the service should be extracted as config, so they can be set and changed in the environment at run time.



Advantages of the clean anatomy

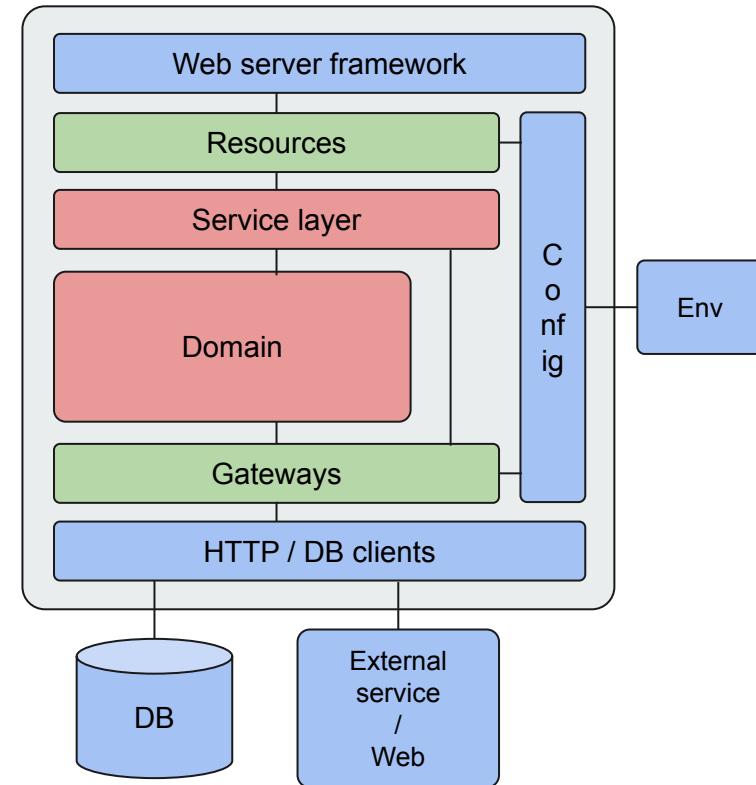
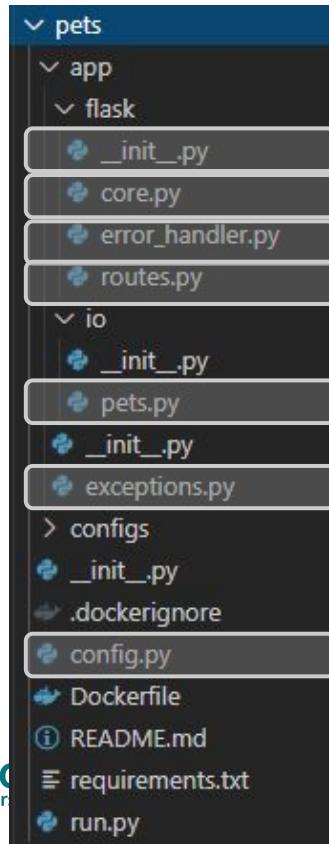
- **Independent of externals**
When any of the external parts of the system become obsolete, like the database, or the web framework, you can replace those obsolete elements with a minimum of fuss.
- **Scalable and maintainable**
A modular architecture allows for scaling without an explosion in complexity. New modules can be added independently of others. Multiple developers can work in parallel on different modules.
- **Deploy everywhere**
Since the configuration is separated, the same code can be deployed in different environments without changes.



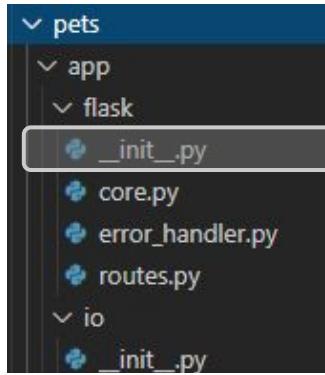
The anatomy of an example

Mapping the structure on the pets example

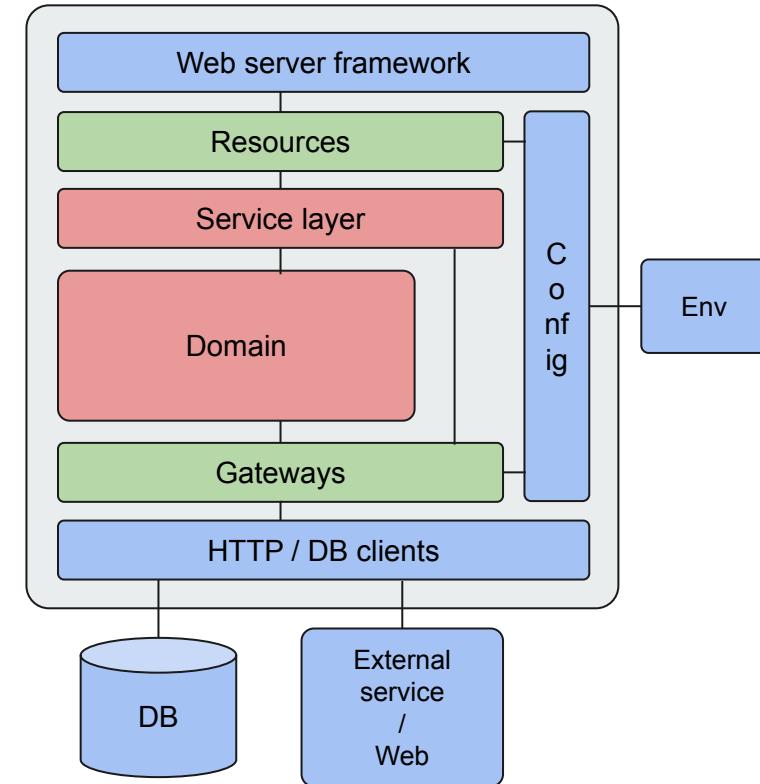
The anatomy of a microservice: pets example



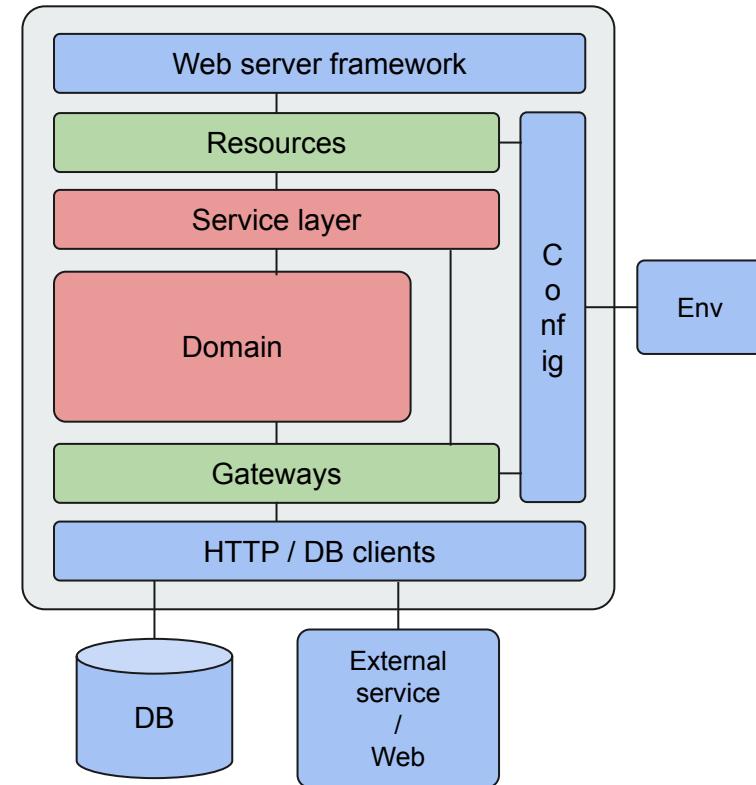
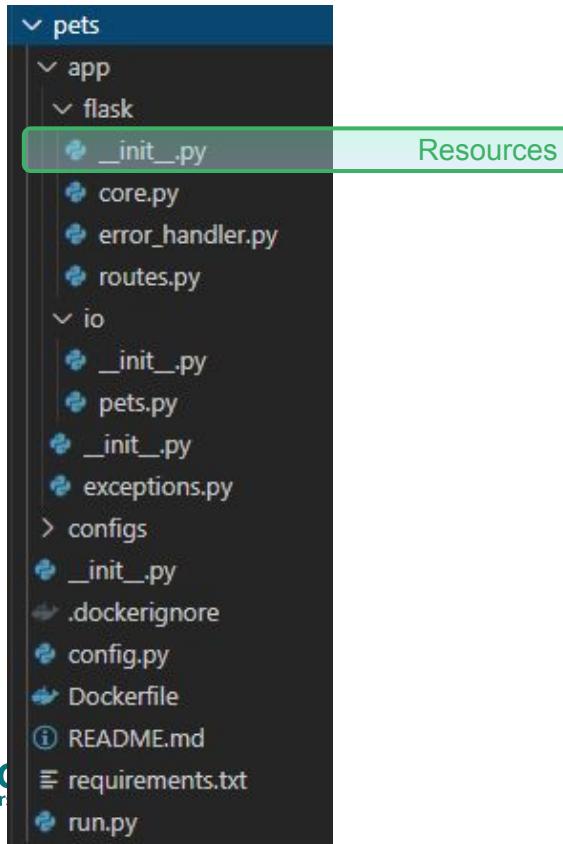
The anatomy of a microservice: pets example



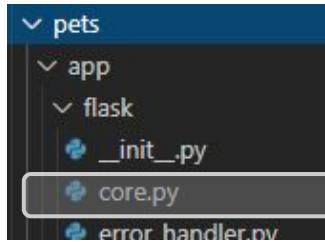
```
20 # Use newer version of swagger ui
21 options = {'swagger_path': swagger_ui_path}
22
23 app = connexion.App(__name__,
24                     specification_dir=config.SPECIFICATION_DIR,
25                     options=options)
26
27 app.add_api('swagger.yaml',
28             strict_validation=True,
29             validate_responses=True)
30
31 error_handler.register_error_handlers(app)
```



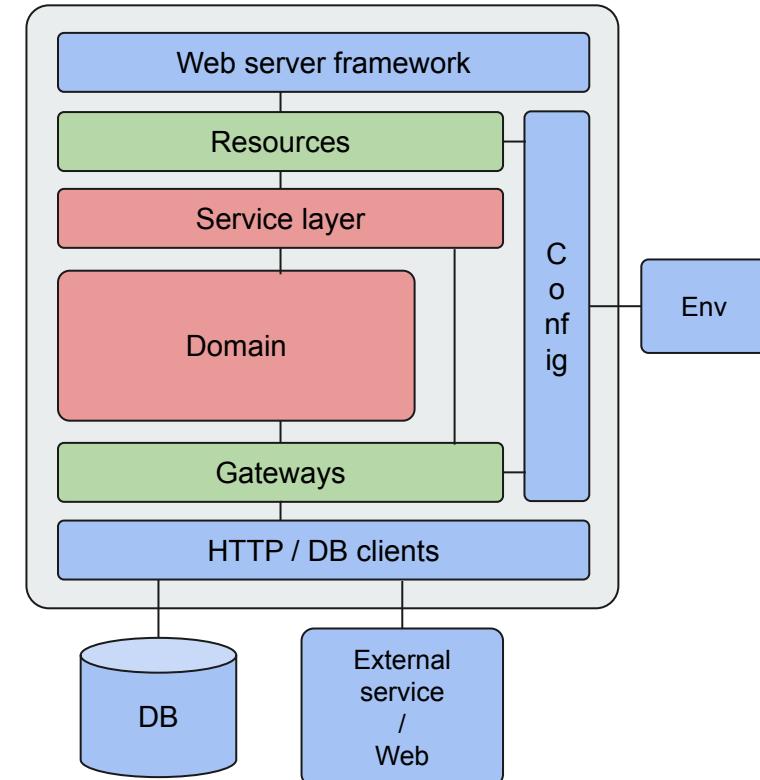
The anatomy of a microservice: pets example



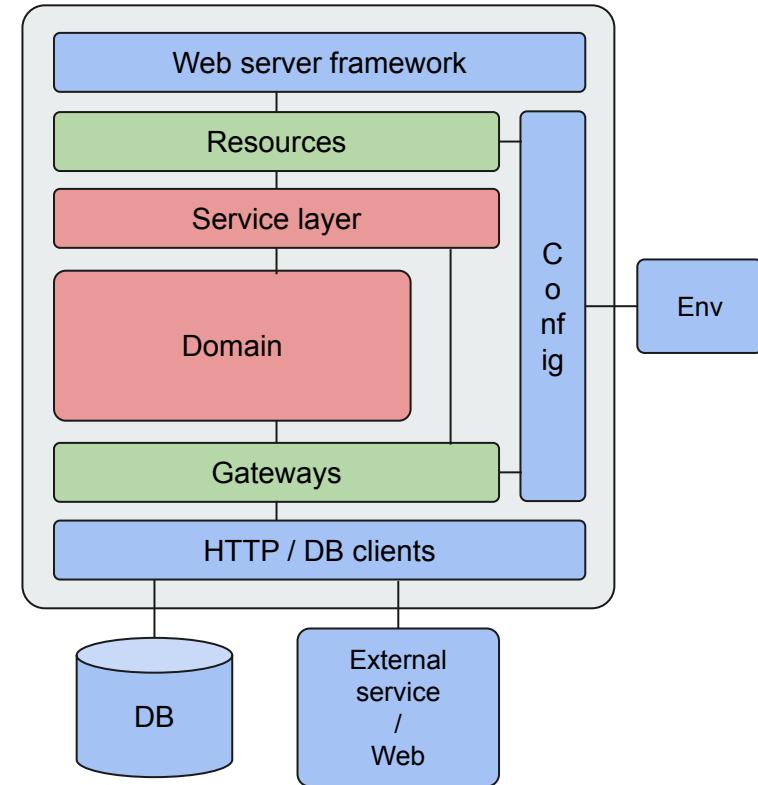
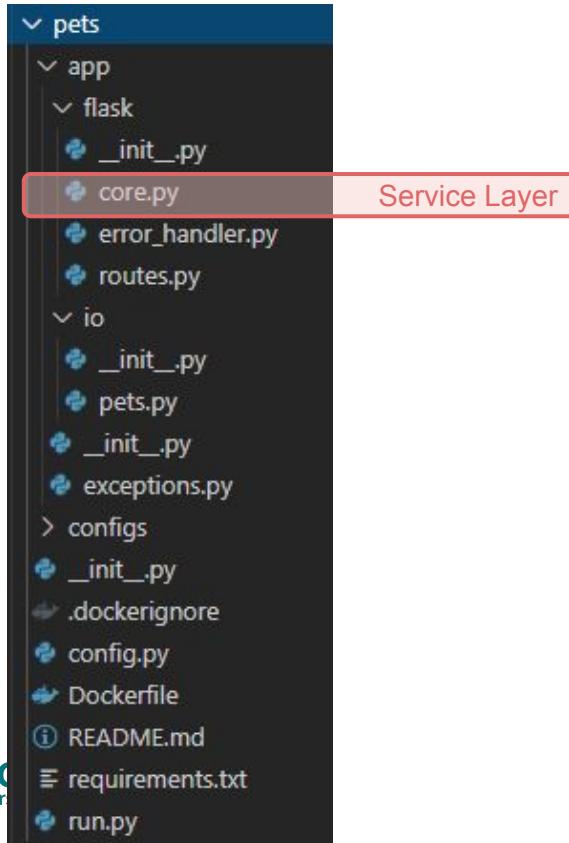
The anatomy of a microservice: pets example



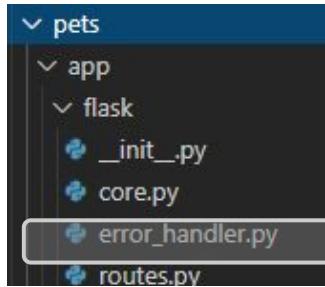
```
7      from app.io import pets  
8  
9  
10     def get_all_pets():  
11         """Get array of all pets."""  
12         return pets.get_all()  
13  
14  
15     def get_pet_by_id(pet_id):  
16         """Get single pet by id."""  
17         return pets.get(pet_id)  
18  
19  
20     def add_pet(*, name, owner, description):  
21         """Add a new pet."""  
22         return pets.create(name=name, owner=owner, description=description)
```



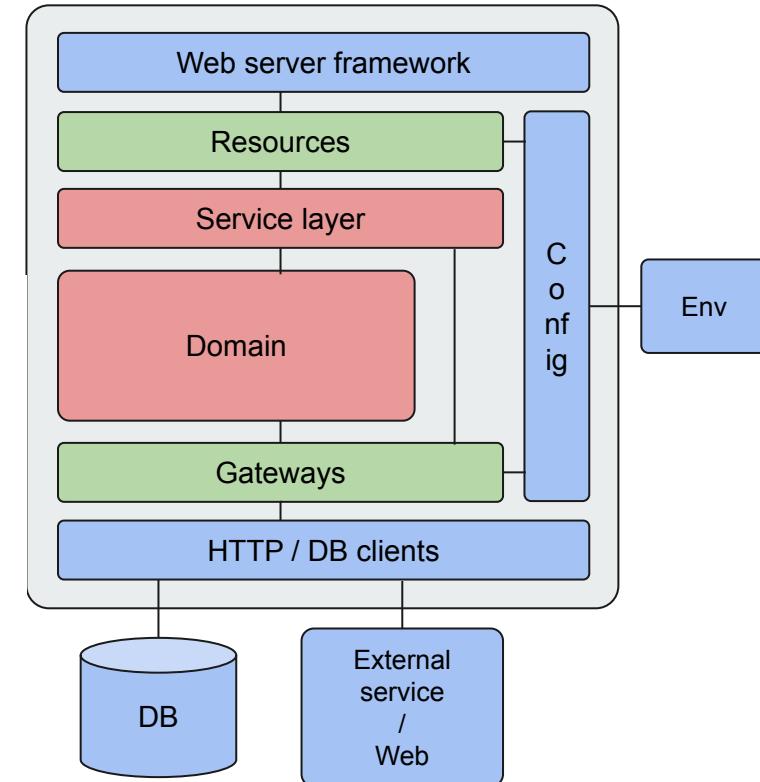
The anatomy of a microservice: pets example



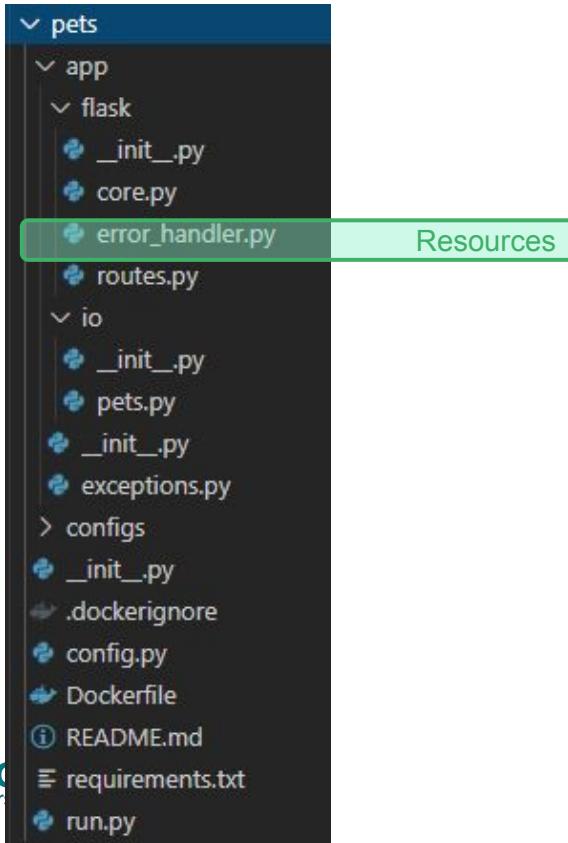
The anatomy of a microservice: pets example



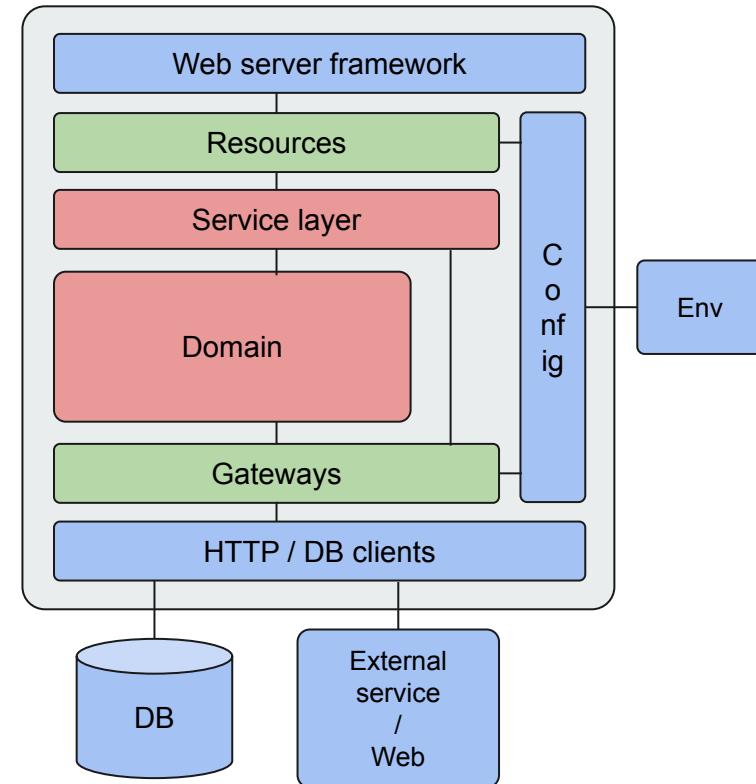
```
8   from flask import jsonify
9
10  from app import exceptions
11
12
13  def handle_custom_error(error):
14      """Handle custom errors."""
15      response = jsonify({'message': error.message})
16      response.status_code = error.status_code
17      return response
18
19
20  def register_error_handlers(app):
21      """Add error handlers to the app."""
22      app.add_error_handler(exceptions.PetNotFound, handle_custom_error)
```



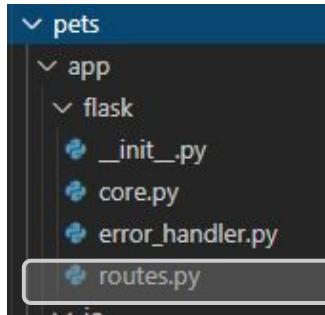
The anatomy of a microservice: pets example



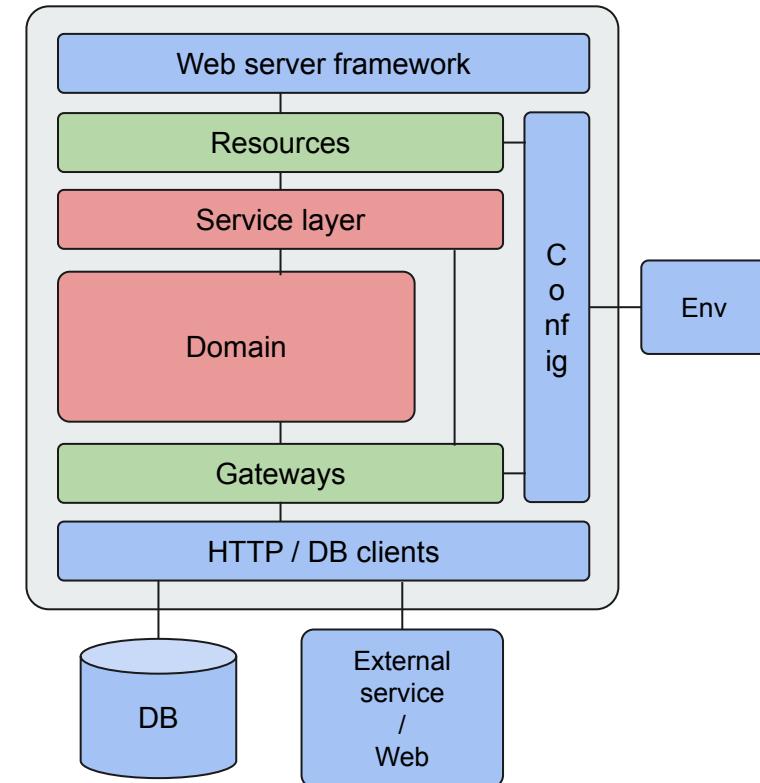
```
└── pets
    ├── app
    │   ├── flask
    │   │   ├── __init__.py
    │   │   └── core.py
    │   └── error_handler.py
    │       └── Resources
    ├── io
    │   ├── __init__.py
    │   ├── pets.py
    │   ├── __init__.py
    │   └── exceptions.py
    ├── config
    │   └── __init__.py
    ├── .dockerignore
    ├── config.py
    ├── Dockerfile
    └── (1) README.md
        └── requirements.txt
            └── run.py
```



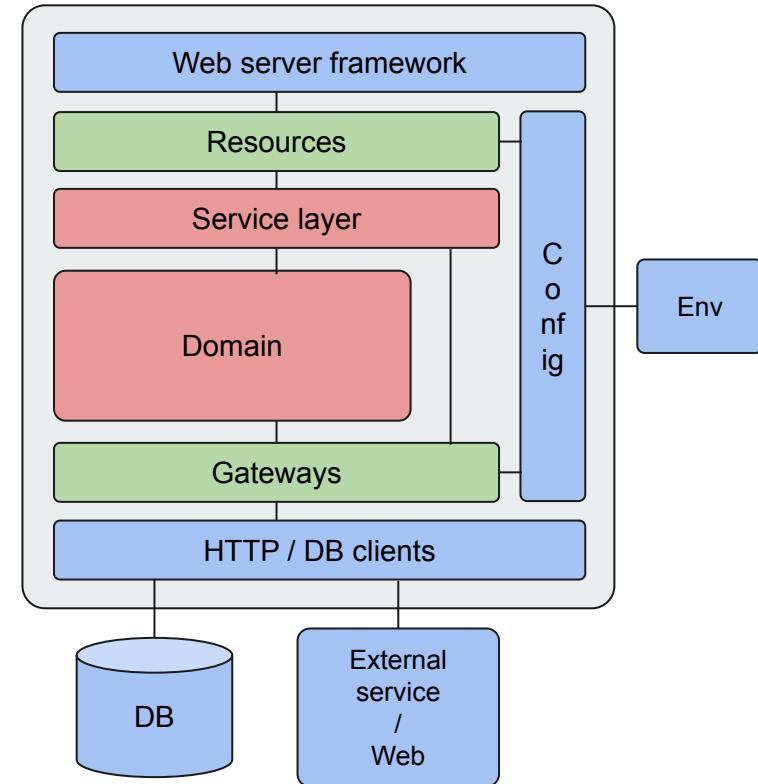
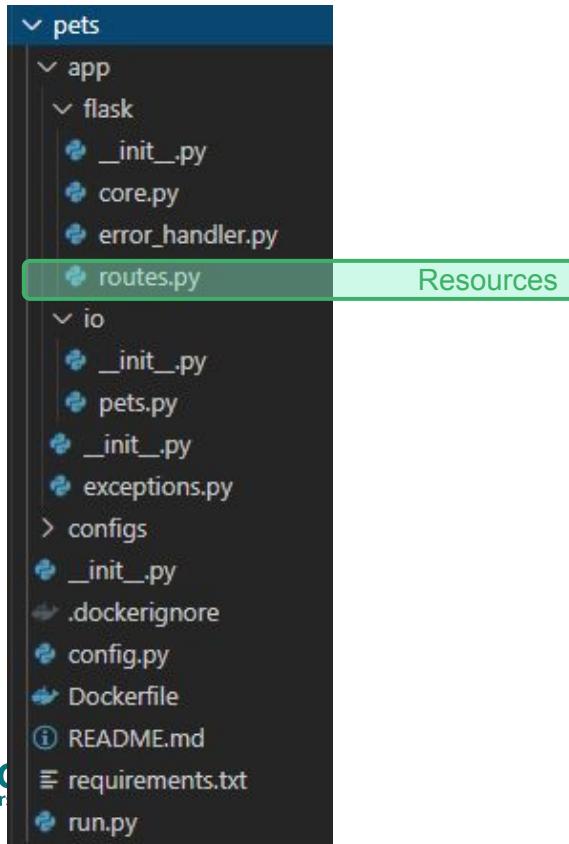
The anatomy of a microservice: pets example



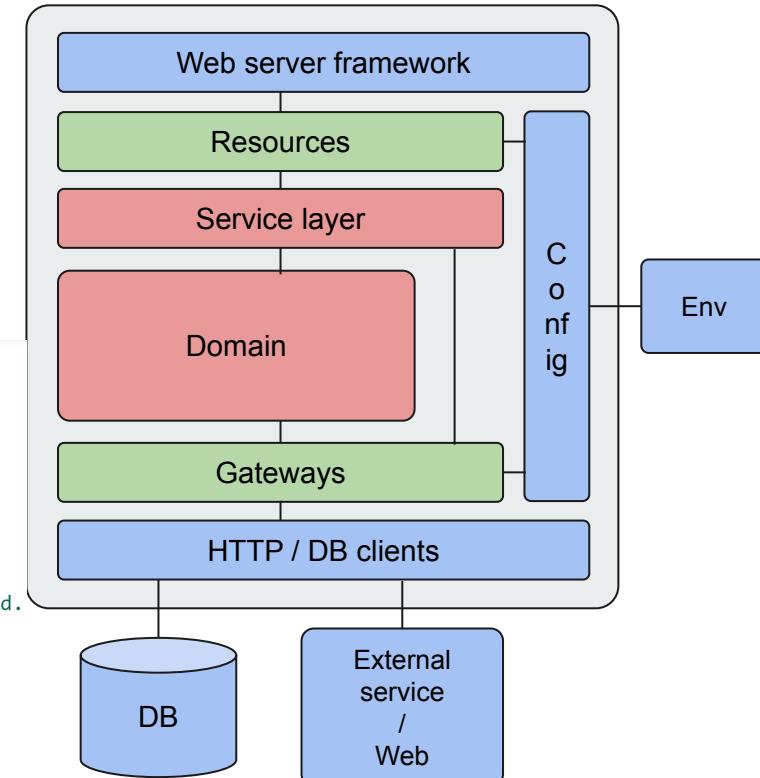
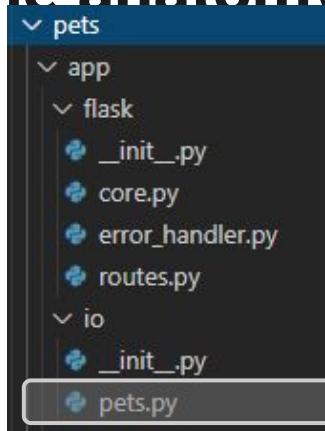
```
7
8
9
10 class Pet:
11     """Includes all HTTP methods for /pets/pet_id>"""
12
13     @staticmethod
14     def get(pet_id):
15         """Get an existing pet."""
16         return core.get_pet_by_id(pet_id), 200
17
18     @staticmethod
19     def delete(pet_id):
20         """Delete an existing pet."""
21         return core.delete_pet_by_id(pet_id), 204
```



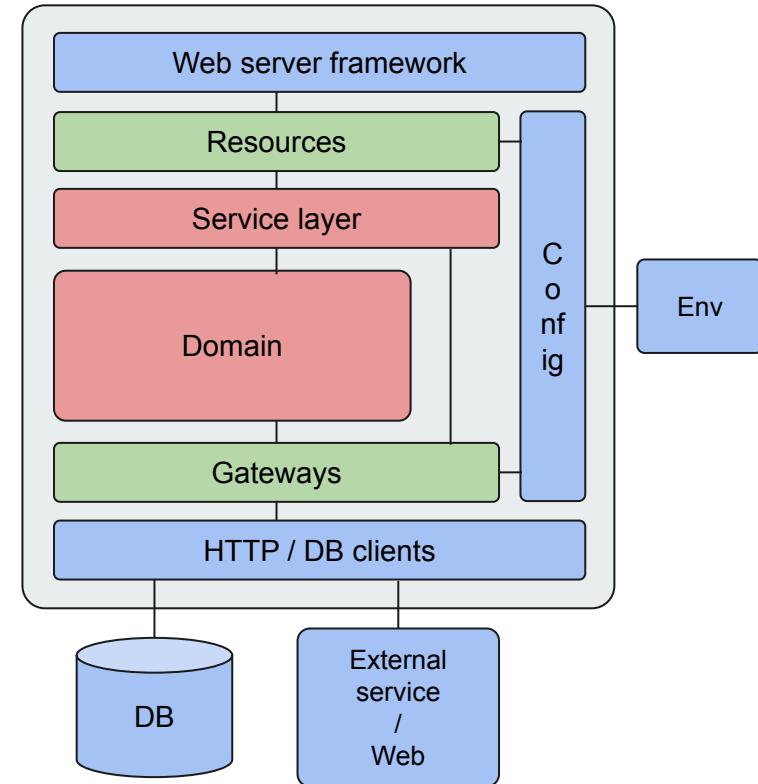
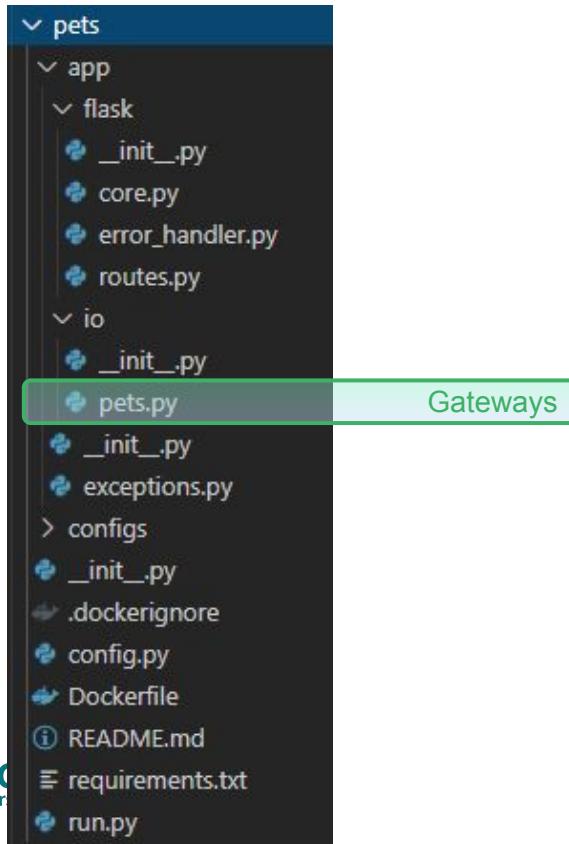
The anatomy of a microservice: pets example



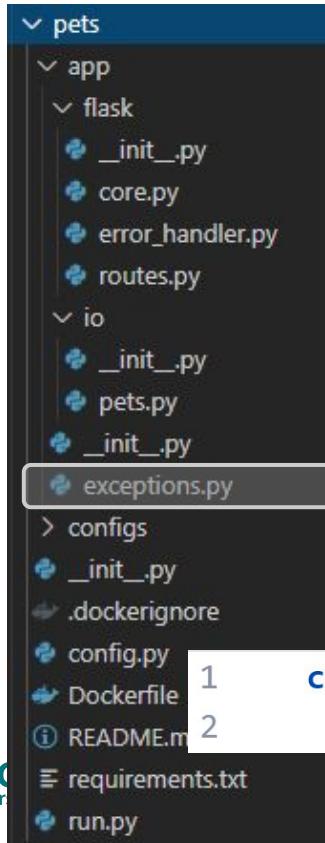
The anatomy of a microservice: pets example



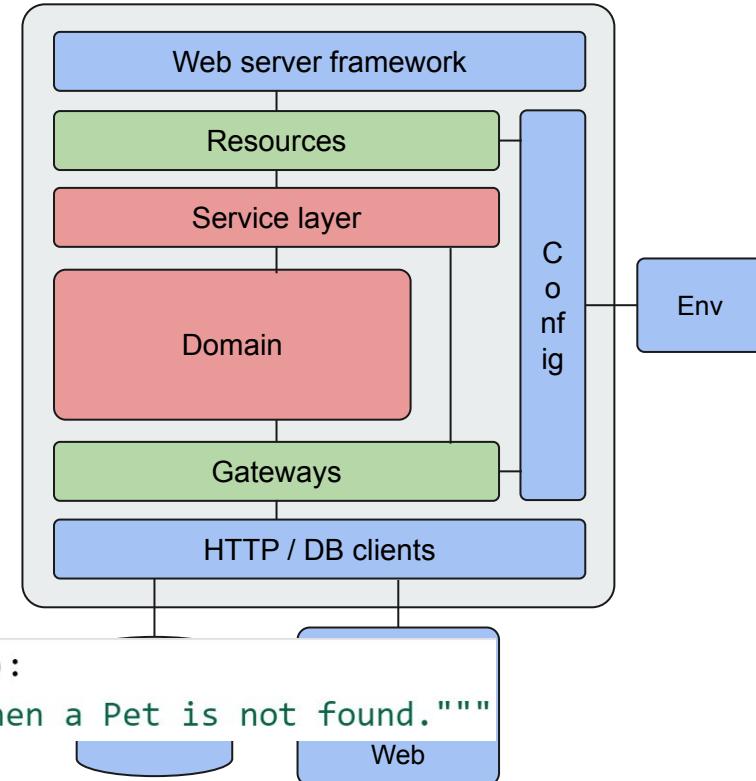
The anatomy of a microservice: pets example



The anatomy of a microservice: pets example

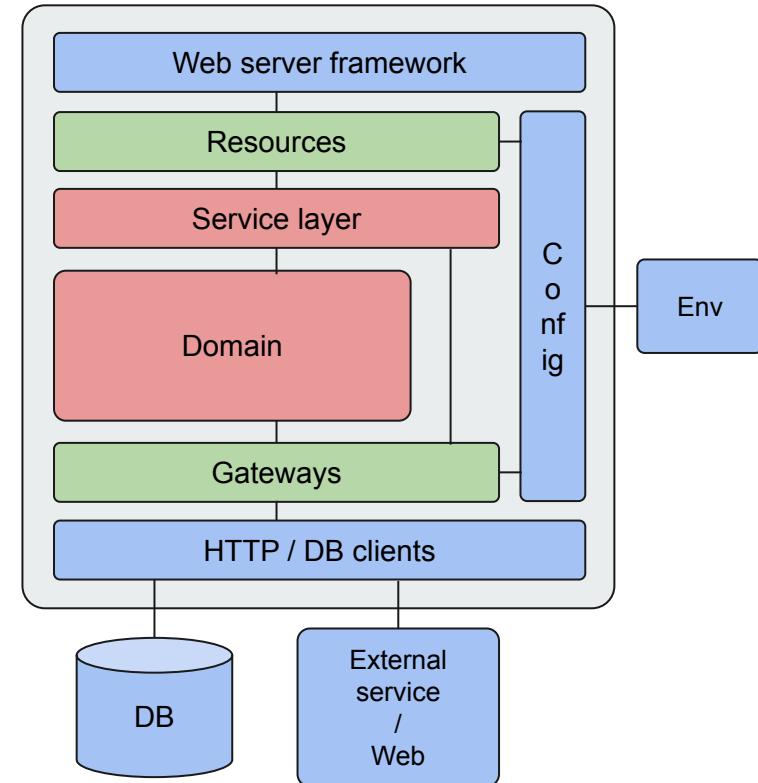


```
1     class PetNotFound(Exception):  
2         """Exception to raise when a Pet is not found."""
```

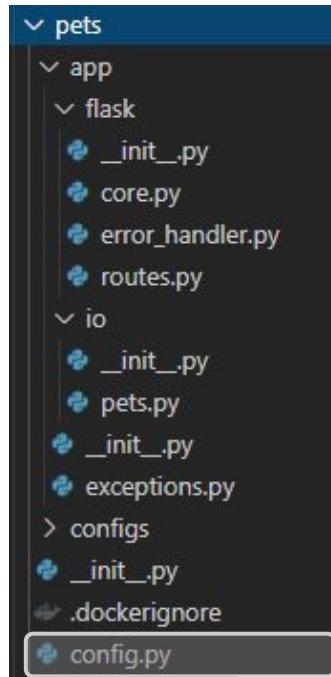


The anatomy of a microservice: pets example

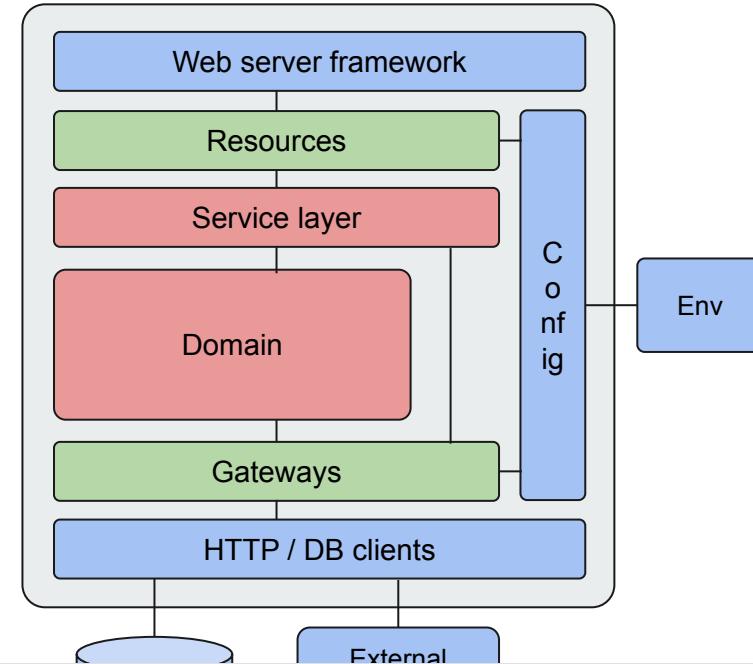
```
└── pets
    ├── app
    │   ├── flask
    │   │   ├── __init__.py
    │   │   ├── core.py
    │   │   ├── error_handler.py
    │   │   └── routes.py
    │   └── io
    │       ├── __init__.py
    │       ├── pets.py
    │       └── __init__.py
    └── exceptions.py
        └── Service Layer
    ├── configs
    └── README.md
```



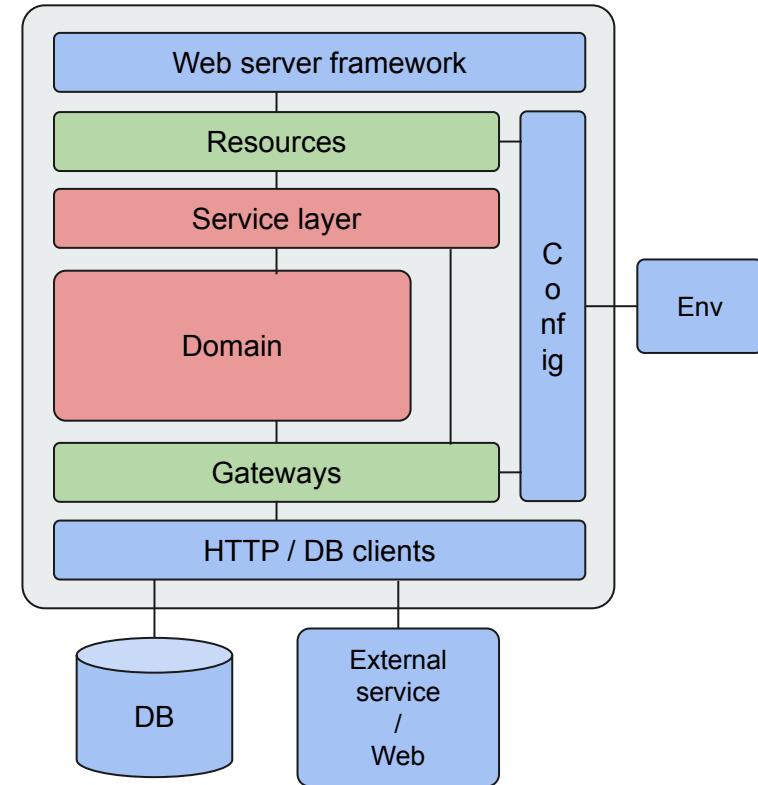
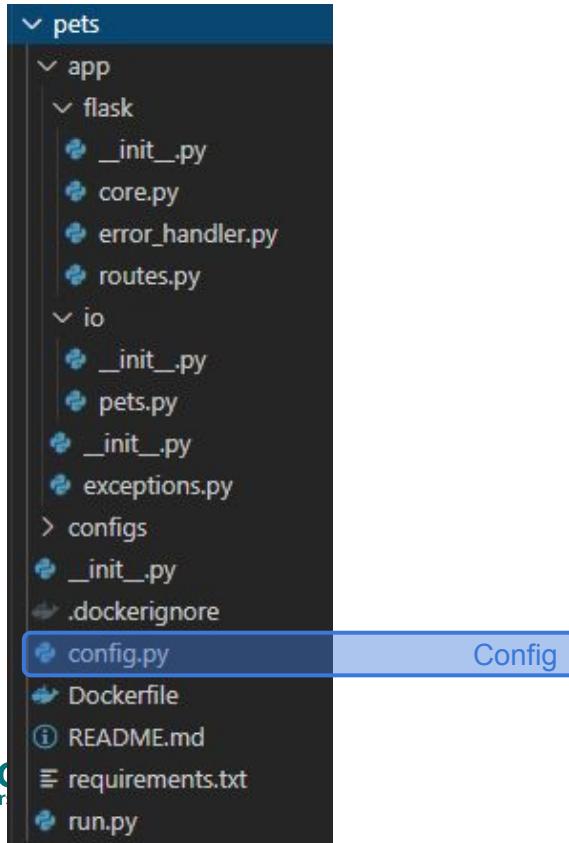
The anatomy of a microservice: pets example



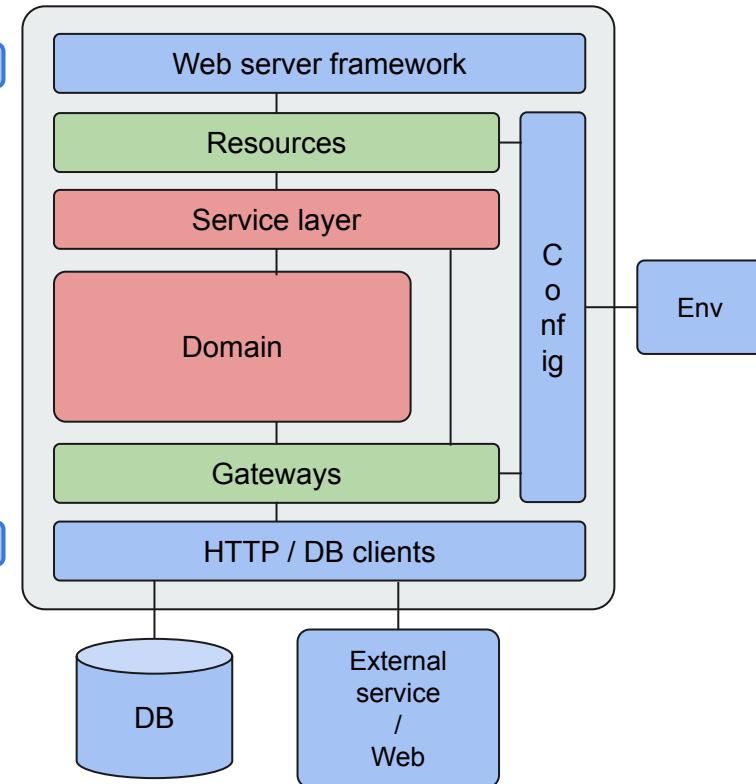
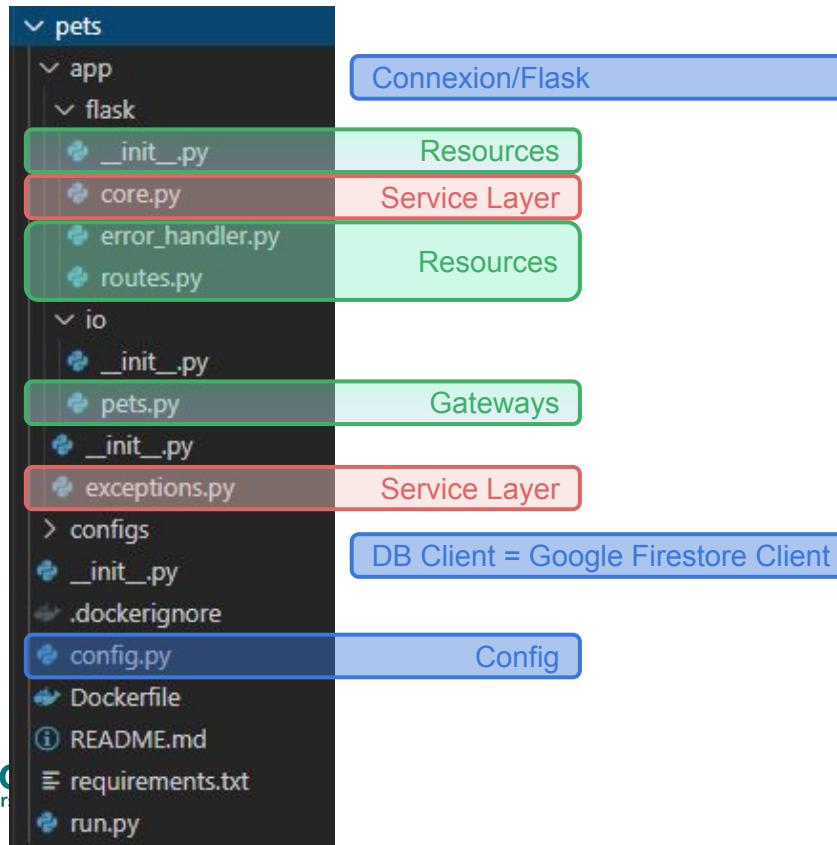
```
1 import os
2
3 SPECIFICATION_DIR = os.path.join(os.path.dirname(os.path.abspath(__file__)),
4                                 'configs')
```



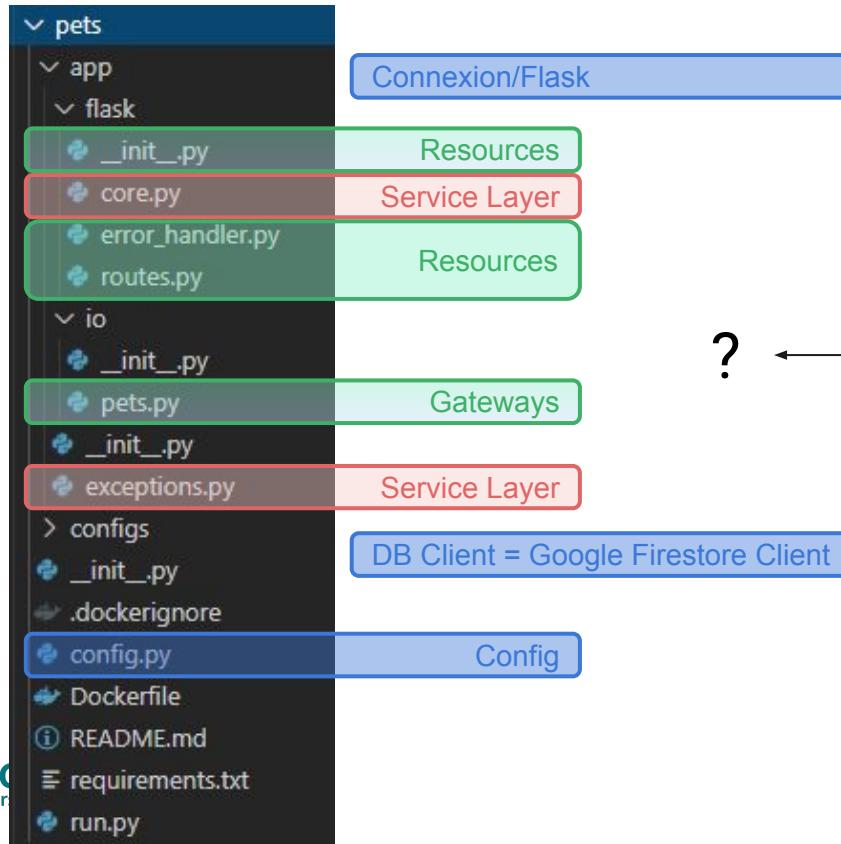
The anatomy of a microservice: pets example



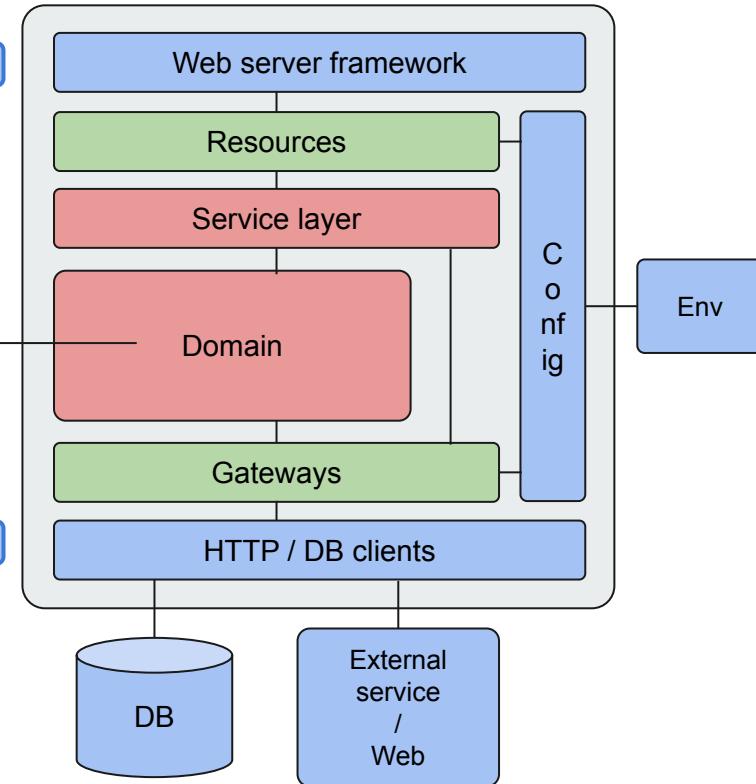
The anatomy of a microservice: pets example



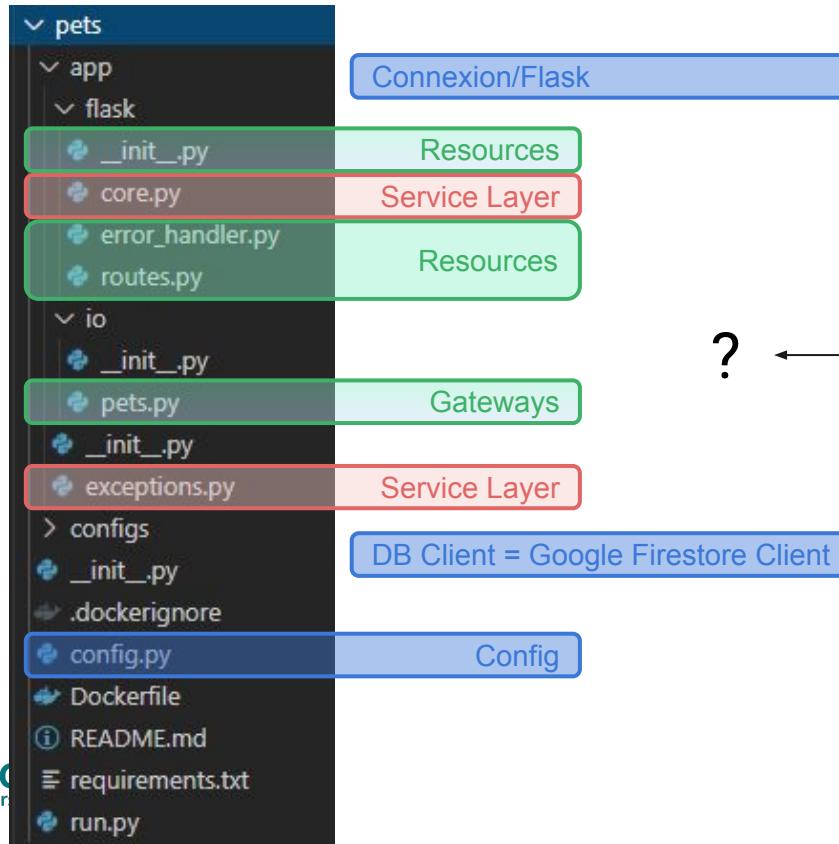
The anatomy of a microservice: pets example



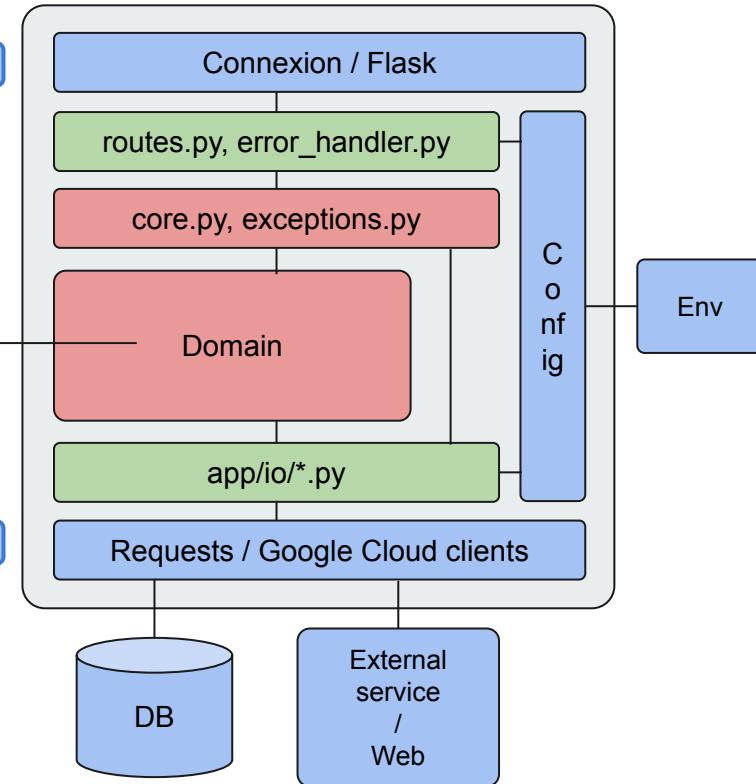
?



The anatomy of a microservice: pets example



?



Monoliths aren't for Dinosaurs

The pros and cons of monoliths

Credits where credits are due



Jasper Van den Bossche
Software Engineer @ ML6

Credits where credits are due



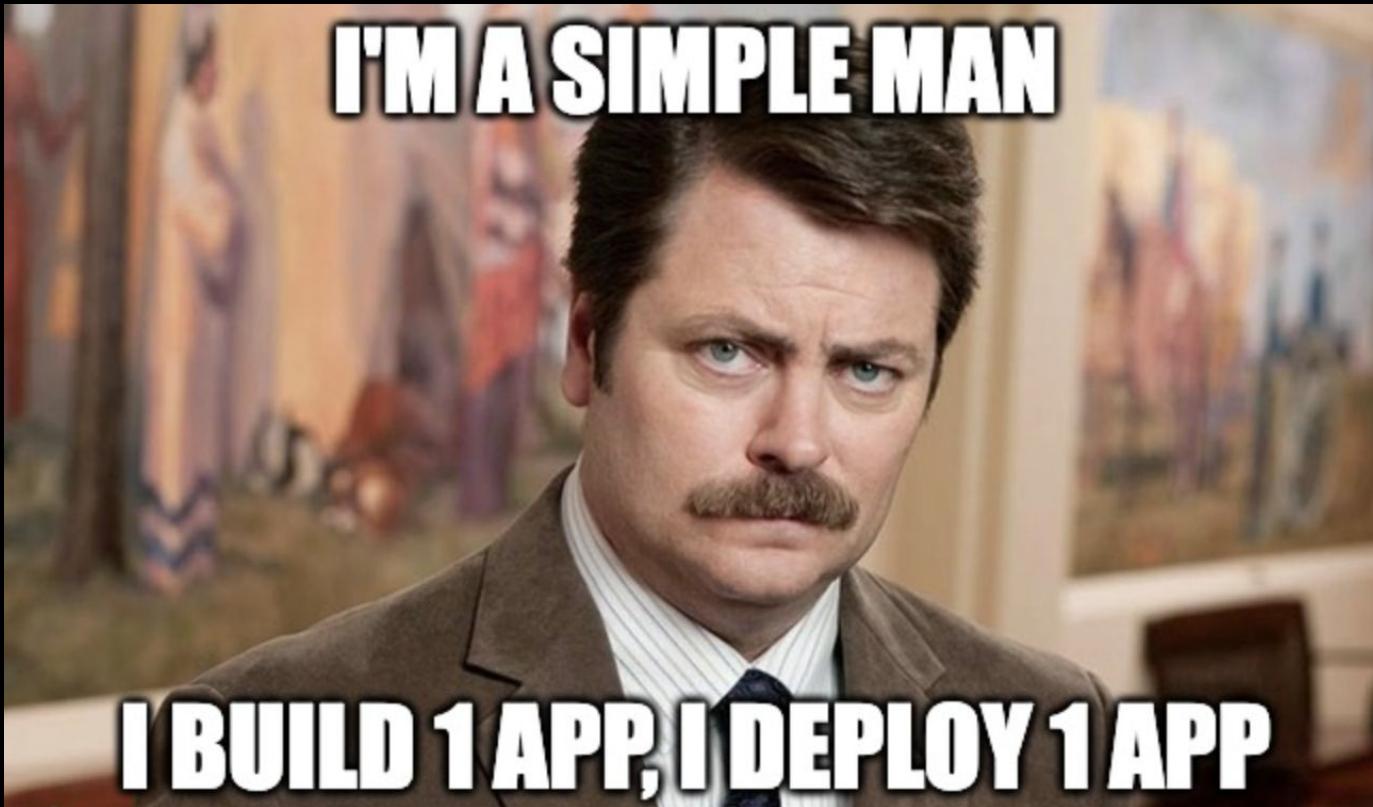
Jasper Van den Bossche
Software Engineer @ ML6

(Disclaimer: Good sense of humour...)

NOT the Goal of this Talk



Simplicity

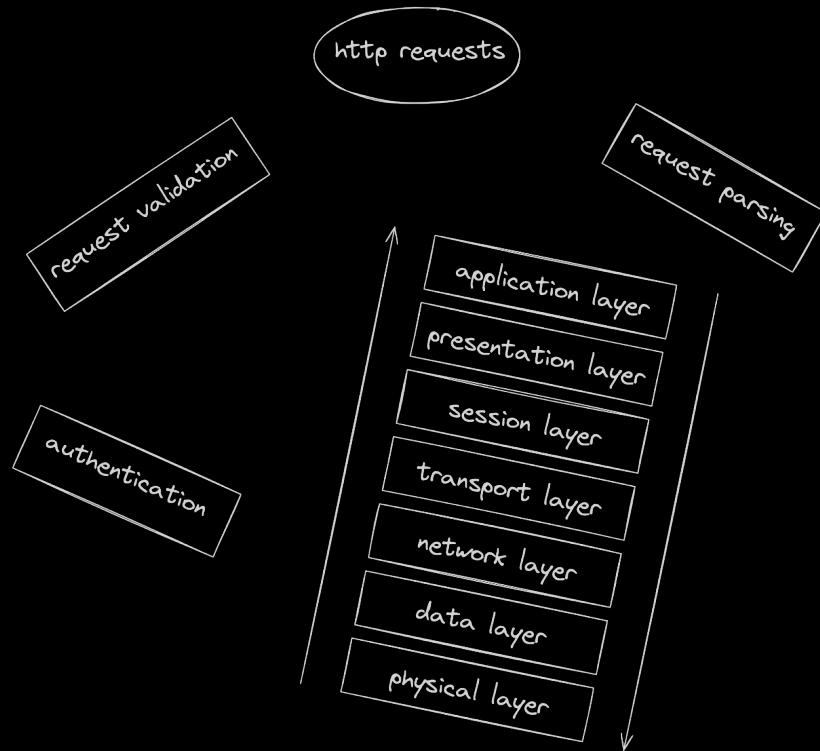


Performance

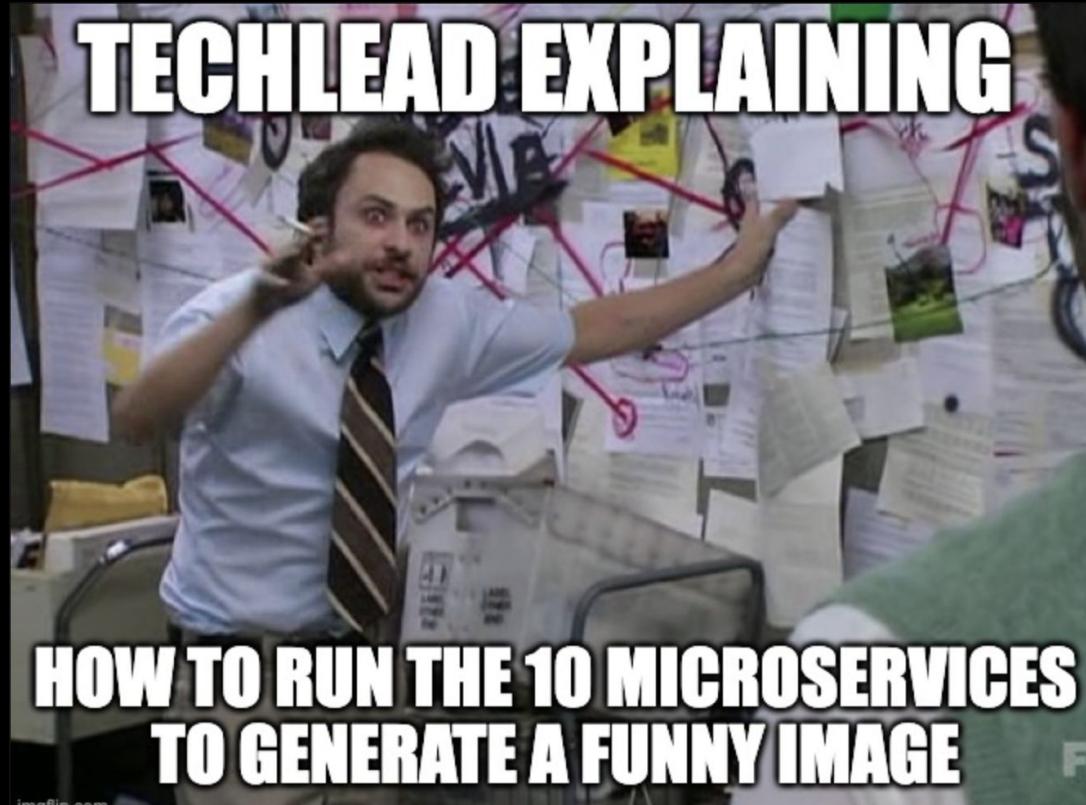
Monolith



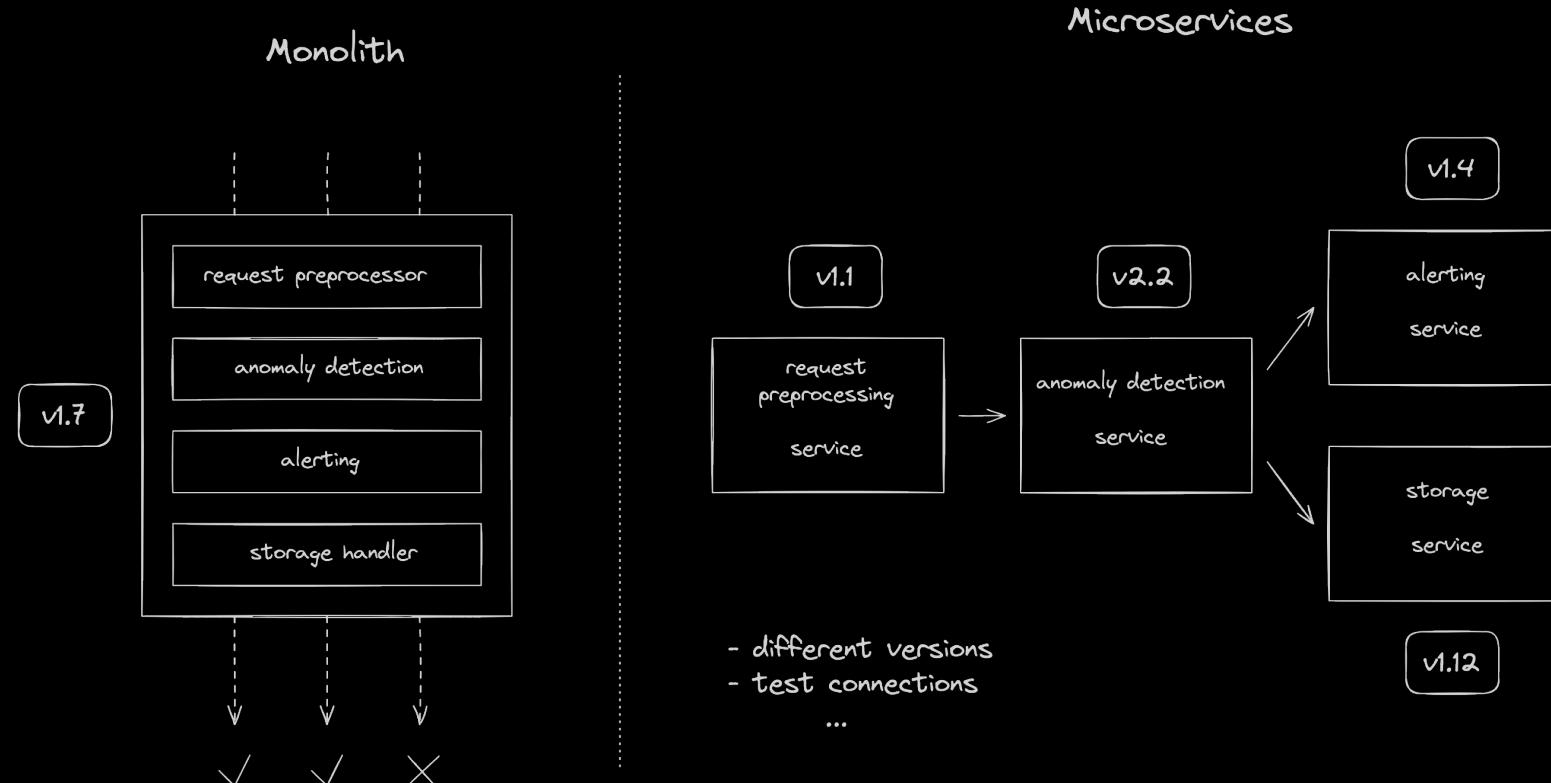
Microservices



End-to-End Testing



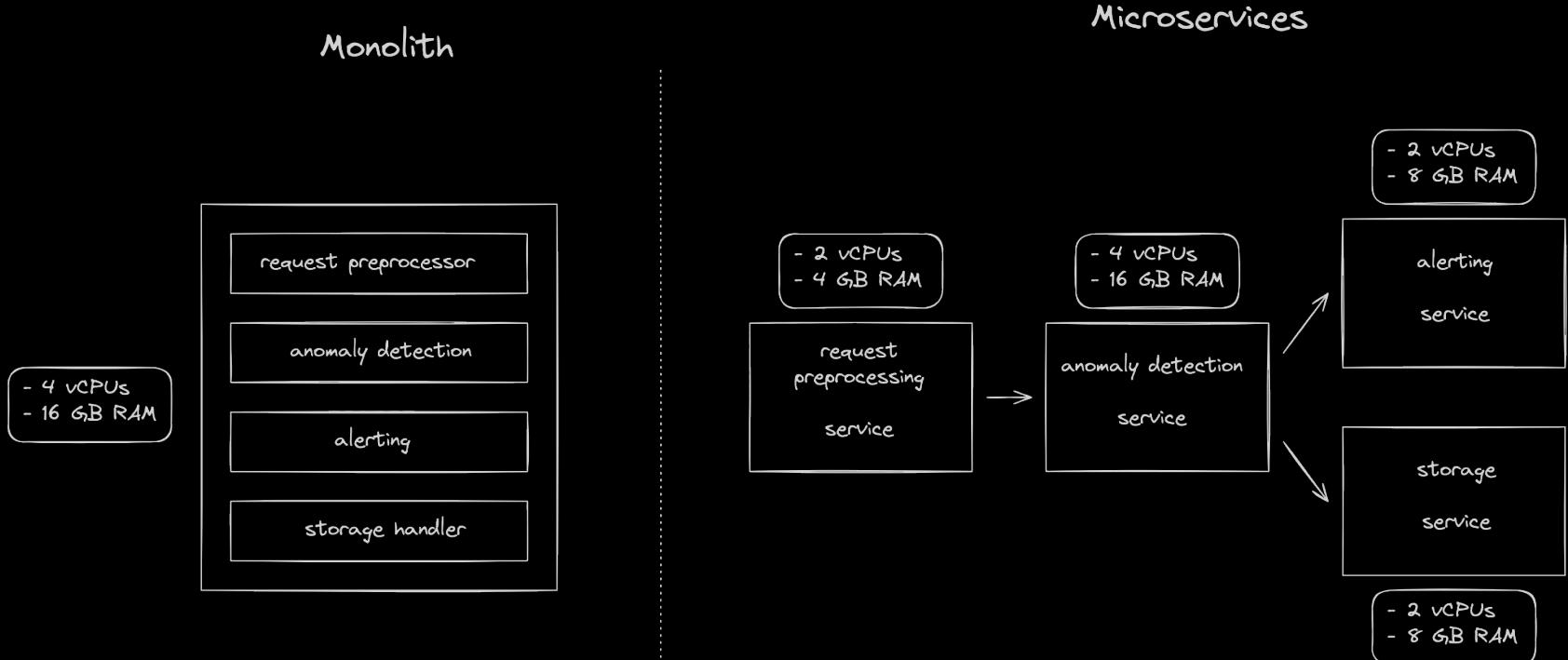
End-to-End Testing



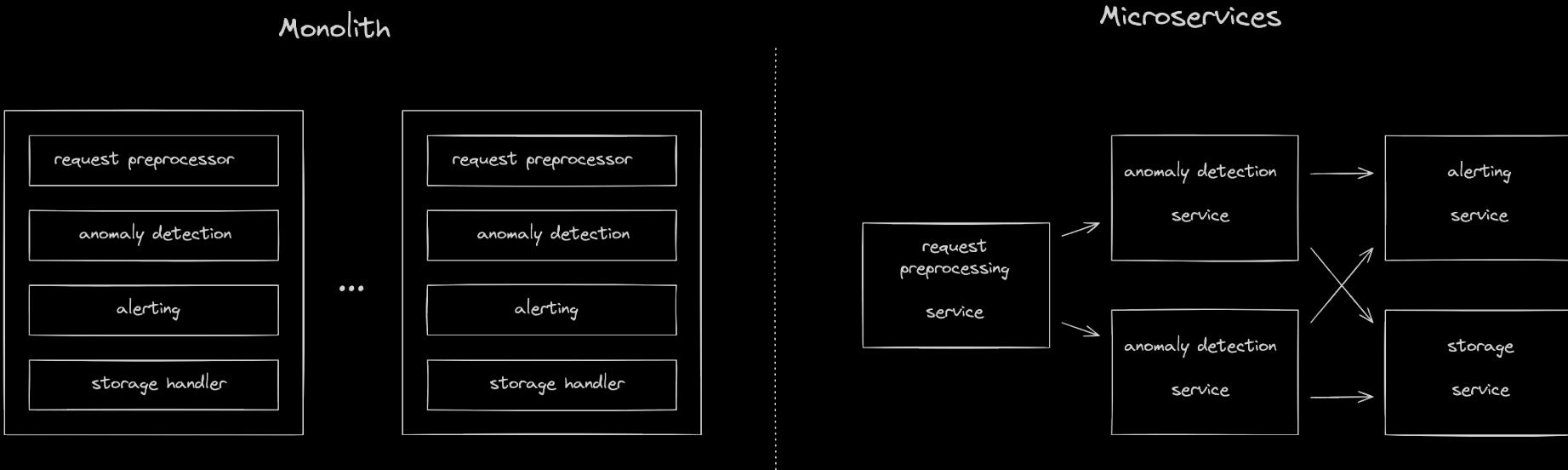
The not so good stuff



Scalability & Resource Utilization



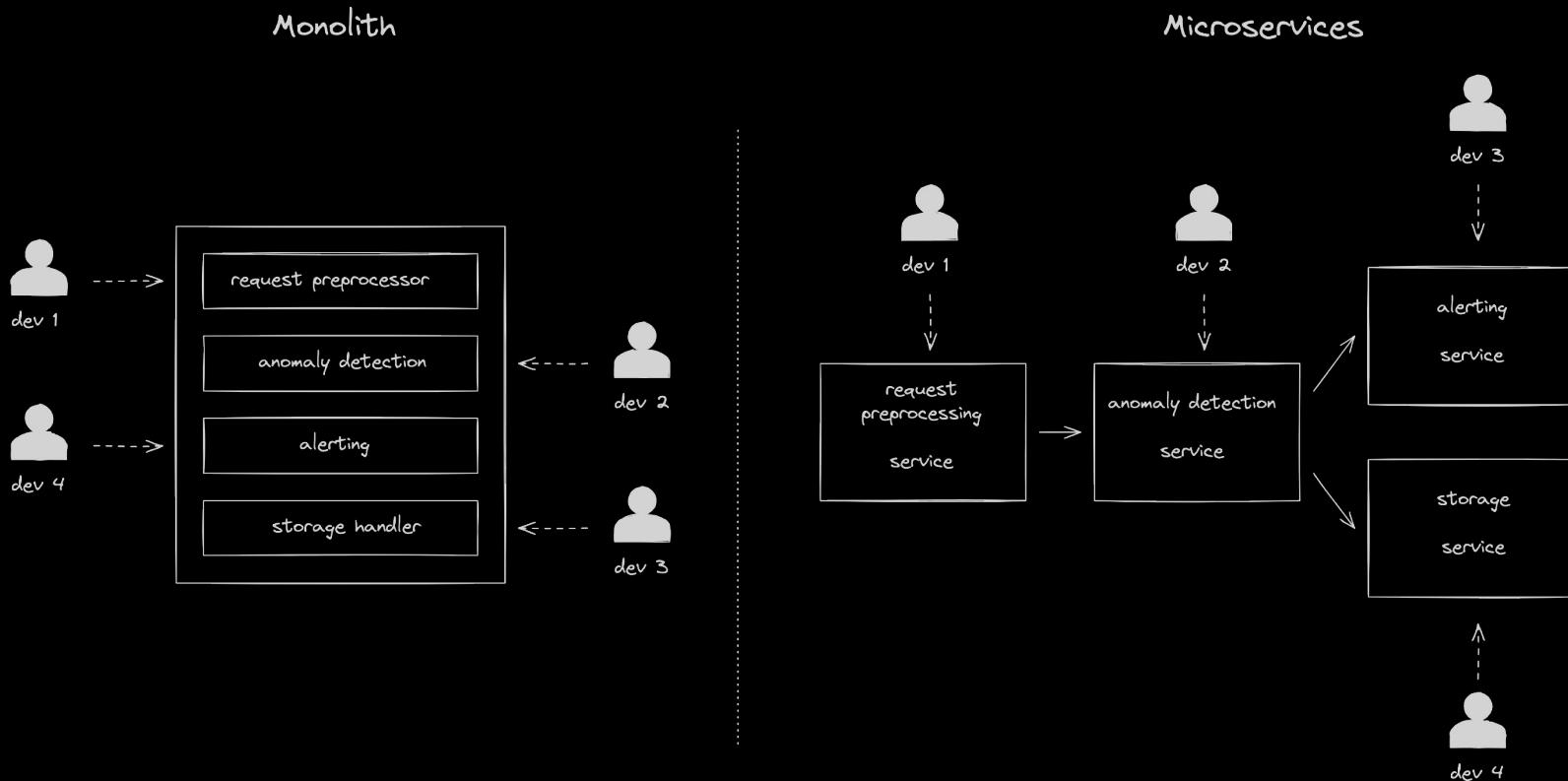
Scalability & Resource Utilization



Teamwork makes the dream work



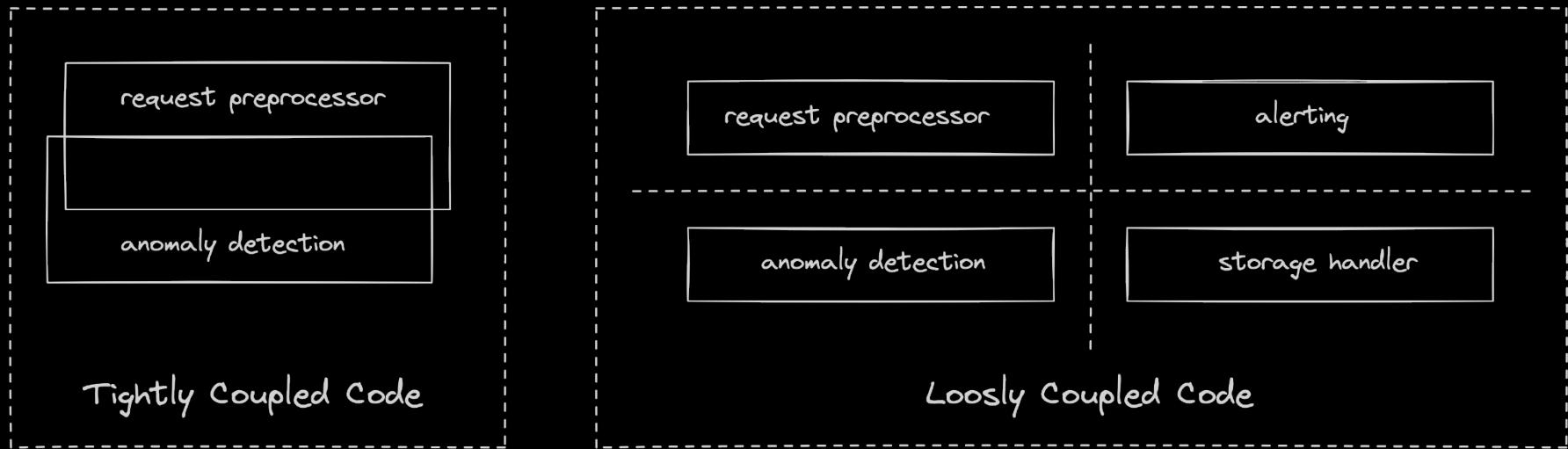
Independent development & ease of maintenance



First Day back on the Job after a Holiday



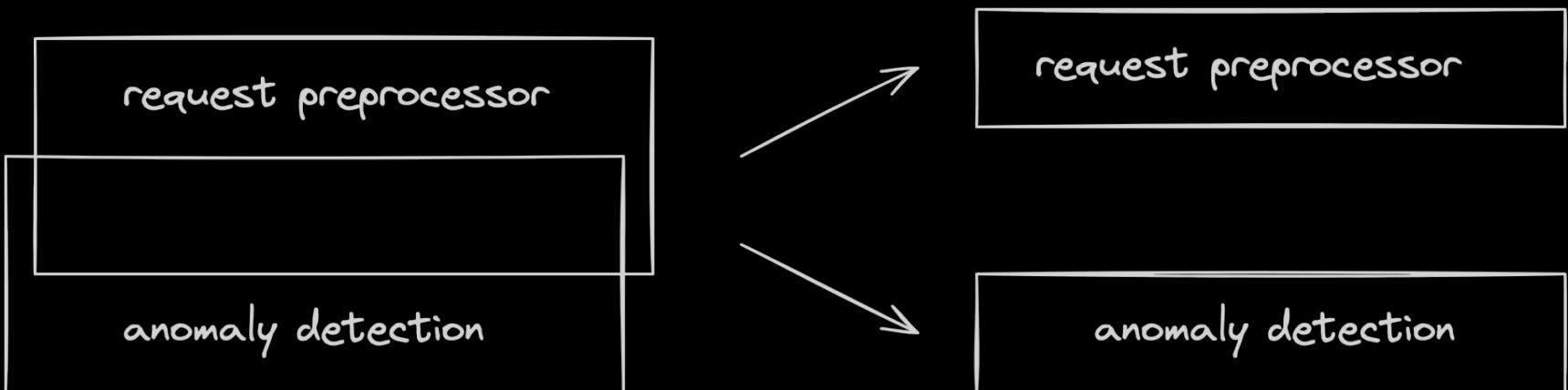
Unwieldy codebases



But how can we avoid this?



Think of your code as independent packages



APIs to the Rescue

anomaly_detection

/api
/helpers

pre_processing

/api
/helpers

Conclusion: Monoliths & Microservices

- Less infrastructure to manage
- Better performance
- Easier to run and test

BUT

- Require discipline

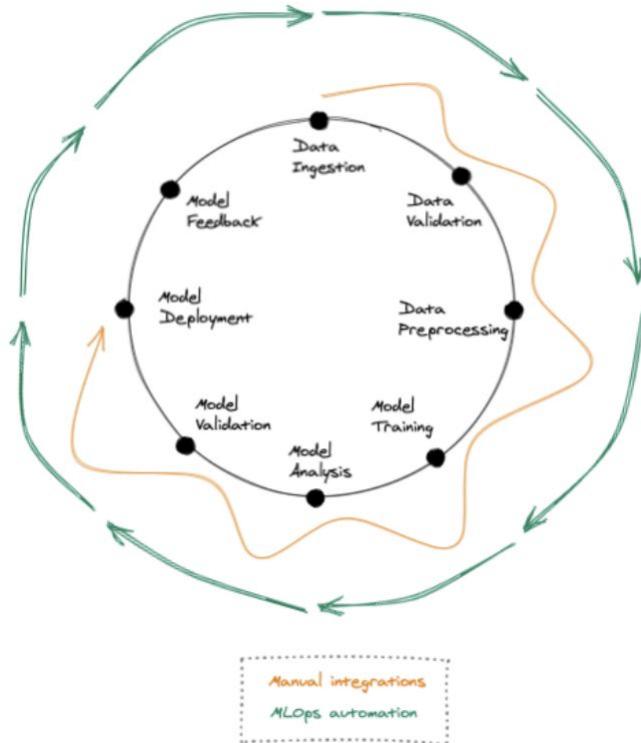
THERE IS NO ARCHITECTURE TO RULE THEM ALL

Lab:

Docker compose

ML model pipeline

Why do we need ML pipelines?



The general idea is to not treat **ML workflows** as a one-off, but to treat them in a **reliable** and **reproducible** way.

Why do we need ML pipelines?

The Spotify experience

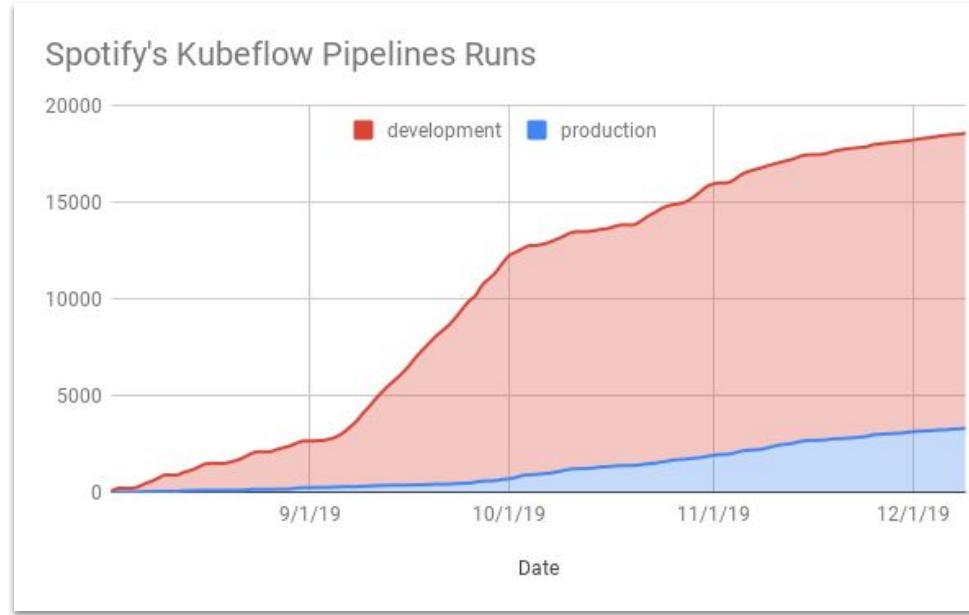
The Winding Road to Better Machine Learning Infrastructure Through Tensorflow Extended and Kubeflow

Posted on December 13, 2019 by [josh baer](#) and [samuelngahane](#)

"As we built these new Machine Learning systems, we started to hit a point where **our engineers spent more of their time maintaining data and backend systems in support of the ML-specific code** than iterating on the model itself. We realized we needed to standardize best practices and build tooling to bridge the gaps between data, backend, and ML: we needed a Machine Learning platform."

Why do we need ML pipelines?

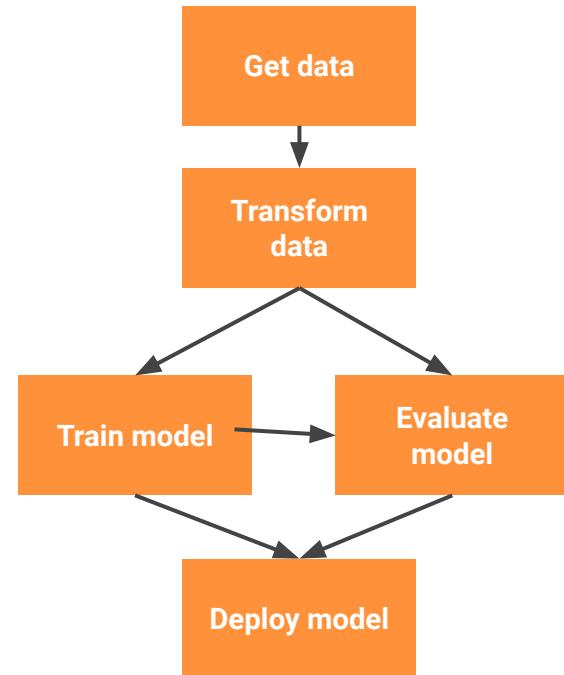
Spotify accomplished an astonishing 7x for running machine learning pipelines



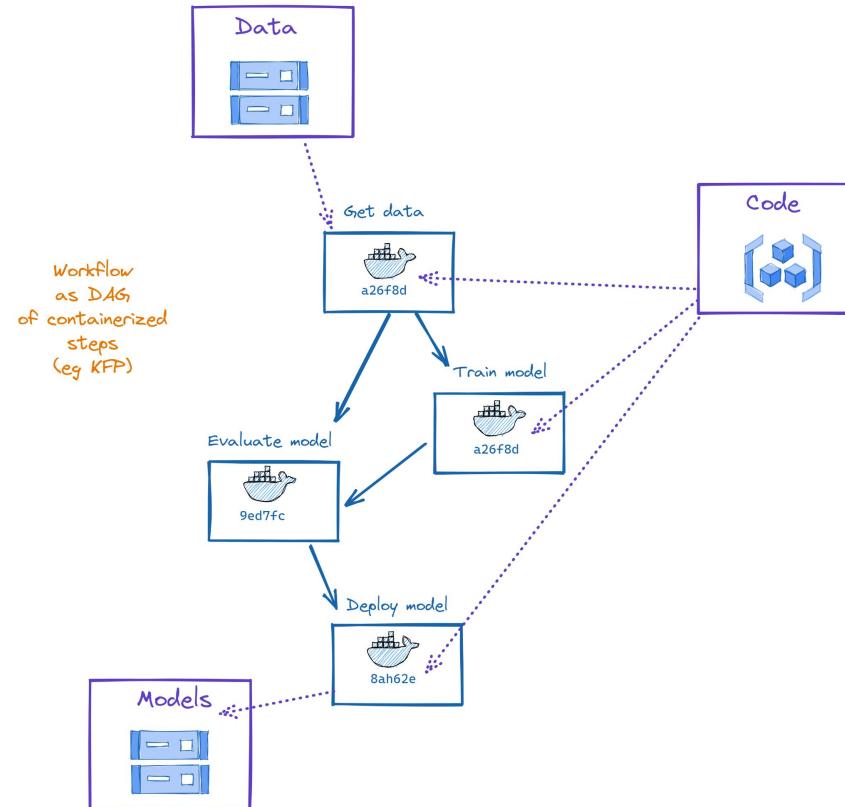
Represent your entire ML workflow as a DAG

Think about your workflow in terms of DAGs (Directed Acyclic Graph):

- What needs to be done sequentially vs in parallel
- Does the step process something itself or call an external service?
- Process first, implementation second



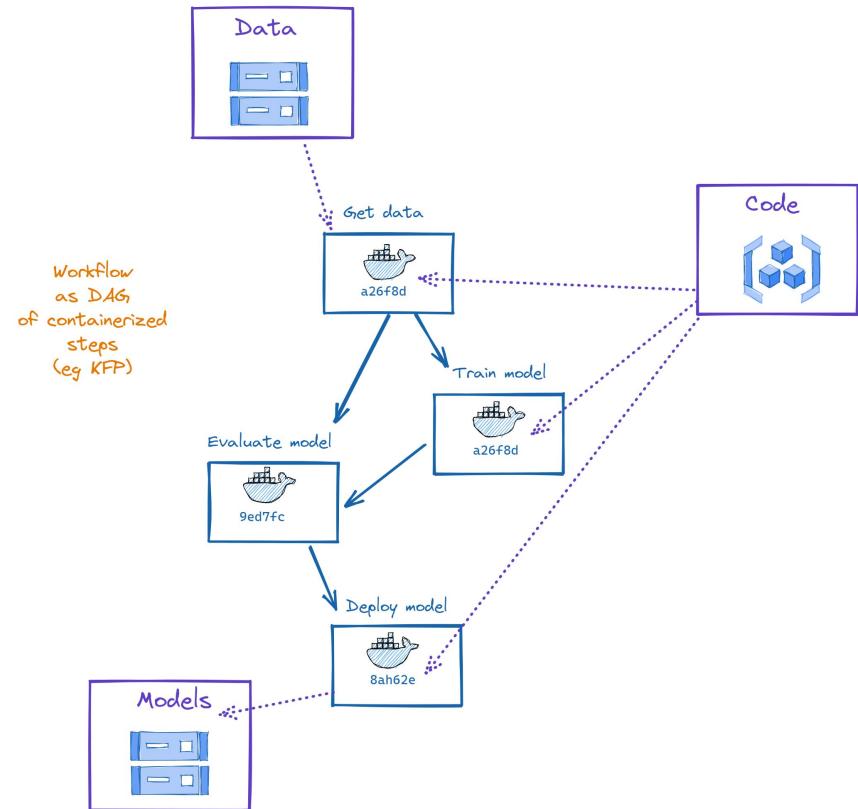
Represent your entire ML workflow as a DAG



- Kubeflow pipelines
- Sagemaker pipelines
- Vertex pipelines
- Valohai
- MLFlow

Definition of a ML Pipeline

"A machine learning pipeline is a series of interconnected data processing and modeling steps designed to automate, standardize and streamline the process of building, training, evaluating and deploying machine learning models."



What are components of an ML pipeline?

A machine learning pipeline is the workflow for a full **machine learning task**. Components are the building blocks of a machine learning pipeline.

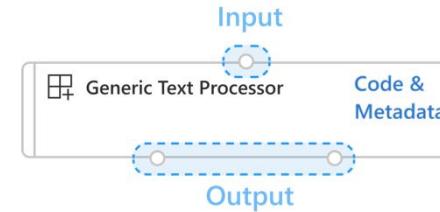
They contain:

- Metadata
- Interface (input/output specifications)
- Command, Code & Environment

Build as **Docker images** or as **python functions** (with specific decorators and configurations).

Components

- Metadata
name, display_name, version, type, etc.
- Interface
input/output specifications (name, type, description, default value, etc.)
- Command, Code & Environment
command, code and environment required to run the component

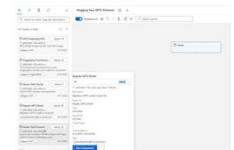


```
1 # component.yaml
2 description: test dataset description used for forecasting tasks
3 type: dataset
4 inputs:
5   - name: train
6     type: integer
7     description: training split (larger value of split, 100-split used as default)
8     default: 80
9   - name: validation
10    type: integer
11    description: validation split (larger value of split, 100-split used as default)
12    default: 10
13   - name: test
14     type: integer
15     description: test split (larger value of split, 100-split used as default)
16     default: 10
17 outputs:
18   - name: train
19     type: dataset
20     description: training dataset
21     default: train
22   - name: validation
23     type: dataset
24     description: validation dataset
25     default: validation
26   - name: test
27     type: dataset
28     description: test dataset
29     default: test
30
31 python_script.py
32
33 # entry point for component
34
35 # validate datasets
36
37 # validate datasets, validate schema
38
39 # validate dataset, validate output
```

CLI

```
1 # component.yaml
2 description: test dataset description used for forecasting tasks
3 type: dataset
4 inputs:
5   - name: train
6     type: integer
7     description: training split (larger value of split, 100-split used as default)
8     default: 80
9   - name: validation
10    type: integer
11    description: validation split (larger value of split, 100-split used as default)
12    default: 10
13   - name: test
14     type: integer
15     description: test split (larger value of split, 100-split used as default)
16     default: 10
17 outputs:
18   - name: train
19     type: dataset
20     description: training dataset
21     default: train
22   - name: validation
23     type: dataset
24     description: validation dataset
25     default: validation
26   - name: test
27     type: dataset
28     description: test dataset
29     default: test
30
31 python_script.py
32
33 # entry point for component
34
35 # validate datasets
36
37 # validate datasets, validate schema
38
39 # validate dataset, validate output
```

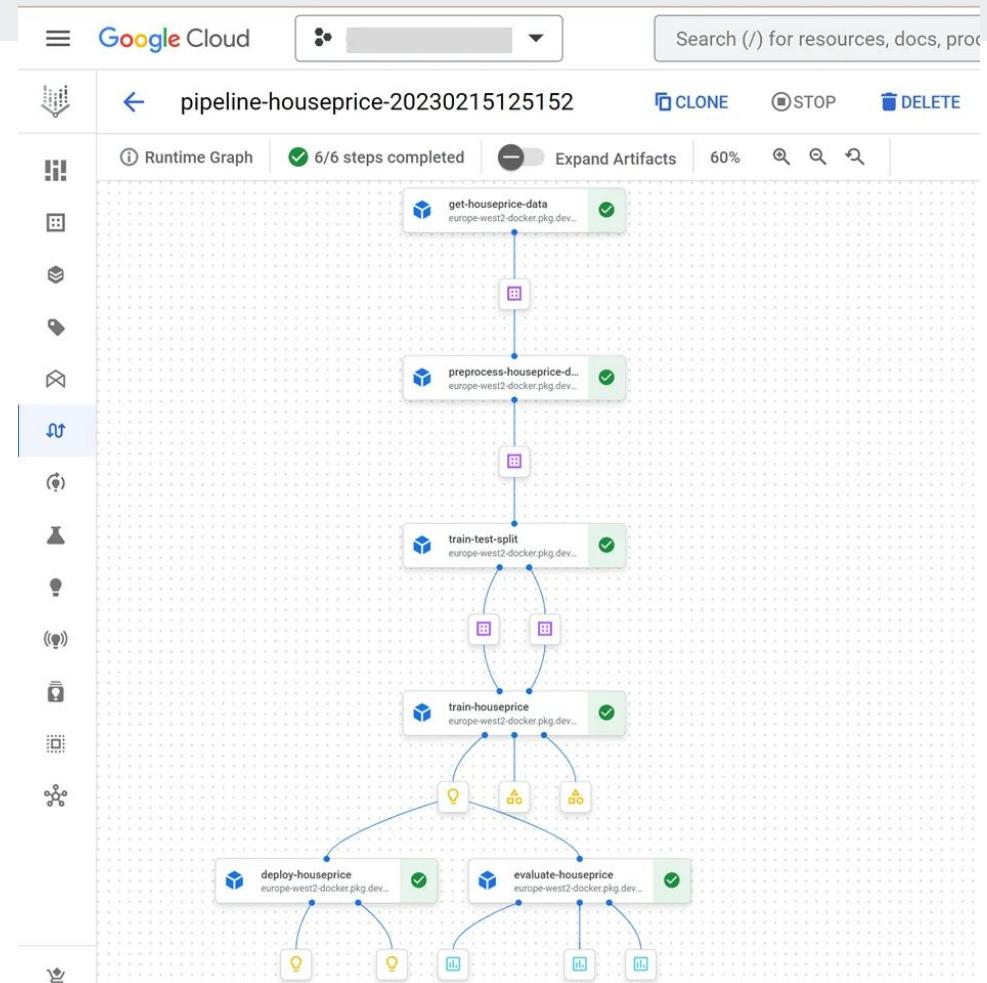
SDK



UI

Example ML Pipeline

Predict house prices with GCP Vertex AI



Example ML Pipeline

Predict house prices with GCP Vertex AI

Based on Kubeflow Pipelines (KFP)

```
1 # IMPORT REQUIRED LIBRARIES
2 from kfp.v2 import dsl
3 from kfp.v2.dsl import (Artifact,
4                         Dataset,
5                         Input,
6                         Model,
7                         Output,
8                         Metrics,
9                         Markdown,
10                        HTML,
11                        component,
12                        OutputPath,
13                        InputPath)
14 from kfp.v2 import compiler
15 from google.cloud.aiplatform import pipeline_jobs
```

Example ML Pipeline

Predict house prices with GCP Vertex AI

Define components.

Can be separate Docker containers or python functions.

Can create a pipeline in one single notebook

```
1  @component(
2      base_image=BASE_IMAGE,
3      output_component_file="get_data.yaml"
4  )
5
6  def get_houseprice_data(
7      filepath: str,
8      dataset_train: Output[Dataset],
9  ):
10
11     import pandas as pd
12
13     df_train = pd.read_csv(filepath + '/train.csv')
14
15     df_train.to_csv(dataset_train.path, index=False)
```

Data Ingestion component

Example ML Pipeline

Predict house prices with GCP Vertex AI

```
1 @component(
2     base_image=BASE_IMAGE,
3     output_component_file="preprocessing.yaml"
4 )
5
6 def preprocess_houseprice_data(
7     train_df: Input[Dataset],
8     dataset_train_preprocessed: Output[Dataset],
9 ):
10
11     import pandas as pd
12     from src.data_preprocessing.preprocessing import data_preprocessing_pipeline
13
14     train_df = pd.read_csv(train_df.path)
15
16     # data_preprocessing_pipeline creates a copy of the df, removes id col, converts to
17     # subtracts YearSold from temporal features and cosine transforms cyclic features.
18     train_df_preprocessed = data_preprocessing_pipeline(train_df)
19
20     train_df_preprocessed.to_csv(dataset_train_preprocessed.path, index=False)
```

Data Preprocessing component

```
1 @component(
2     base_image=BASE_IMAGE,
3     output_component_file="train_test_split.yaml",
4 )
5
6 def train_test_split(dataset_in: Input[Dataset],
7                      dataset_train: Output[Dataset],
8                      dataset_test: Output[Dataset],
9                      test_size: float = 0.2):
10
11     import pandas as pd
12     from sklearn.model_selection import train_test_split
13
14     df = pd.read_csv(dataset_in.path)
15     df_train, df_test = train_test_split(df, test_size=test_size, random_state=42)
16
17     df_train.to_csv(dataset_train.path, index=False)
18     df_test.to_csv(dataset_test.path, index=False)
```

Train test split component

Example ML Pipeline

Predict house prices with GCP Vertex AI

```
...
14 import pandas as pd
15 import pickle
16 import shap
17 from src.modelling.train import HousePriceModel
18 from src.utils.utils import get_image_data
19
20 TARGET = 'SalePrice'
21
22 # Read train and test data
23 train_data = pd.read_csv(dataset_train.path)
24 test_data = pd.read_csv(dataset_test.path)
25
26 # Instantiate the model class
27 house_price_model = HousePriceModel(test_data.copy(),    #we perform hyperparameter t
28                                         target=TARGET,
29                                         n_kfold_splits=3,
30                                         n_trials=100,
31                                         random_state=42)
32
33 # Create X_train and y_train
34 X_train = train_data.drop(TARGET, axis=1)
35 y_train = train_data[TARGET]
36
```

...

Model training component

```
1 @component(
2   base_image=BASE_IMAGE,
3   output_component_file="model_evaluation.yaml"
4 )
5 def evaluate_houseprice(
6   houseprice_model: Input[Model],
7   metrics_baseline: Output[Metrics],
8   metrics_train: Output[Metrics],
9   metrics_test: Output[Metrics]):
10
11   import pickle
12
13   file_name = houseprice_model.path
14   with open(file_name, 'rb') as file:
15     model_data = pickle.load(file)
16
17   scores = model_data["scores_dict"]
18
19   def log_metrics(scores, metric):
20     for metric_name, val in scores.items():
21       metric.log_metric(metric_name, float(val))
22
23   log_metrics(scores["baseline_scores"], metrics_baseline)
24   log_metrics(scores["train_scores"], metrics_train)
25   log_metrics(scores["test_scores"], metrics_test)
```

Model evaluation component

Example ML Pipeline

Predict house prices with GCP Vertex AI

```
17     from google.cloud import aiplatform as vertex_ai
18     from pathlib import Path
19
20     # Checks existing Vertex AI Endpoint or creates Endpoint if it is not exist.
21     def create_endpoint():
22         endpoints = vertex_ai.Endpoint.list(
23             filter='display_name="{}"'.format(model_endpoint),
24             order_by='create_time desc',
25             project=gcp_project,
26             location=gcp_region,
27         )
28         if len(endpoints) > 0:
29             endpoint = endpoints[0] # most recently created
30         else:
31             endpoint = vertex_ai.Endpoint.create(
32                 display_name=model_endpoint,
33                 project=gcp_project,
34                 location=gcp_region
35         )
36     return endpoint
37
```

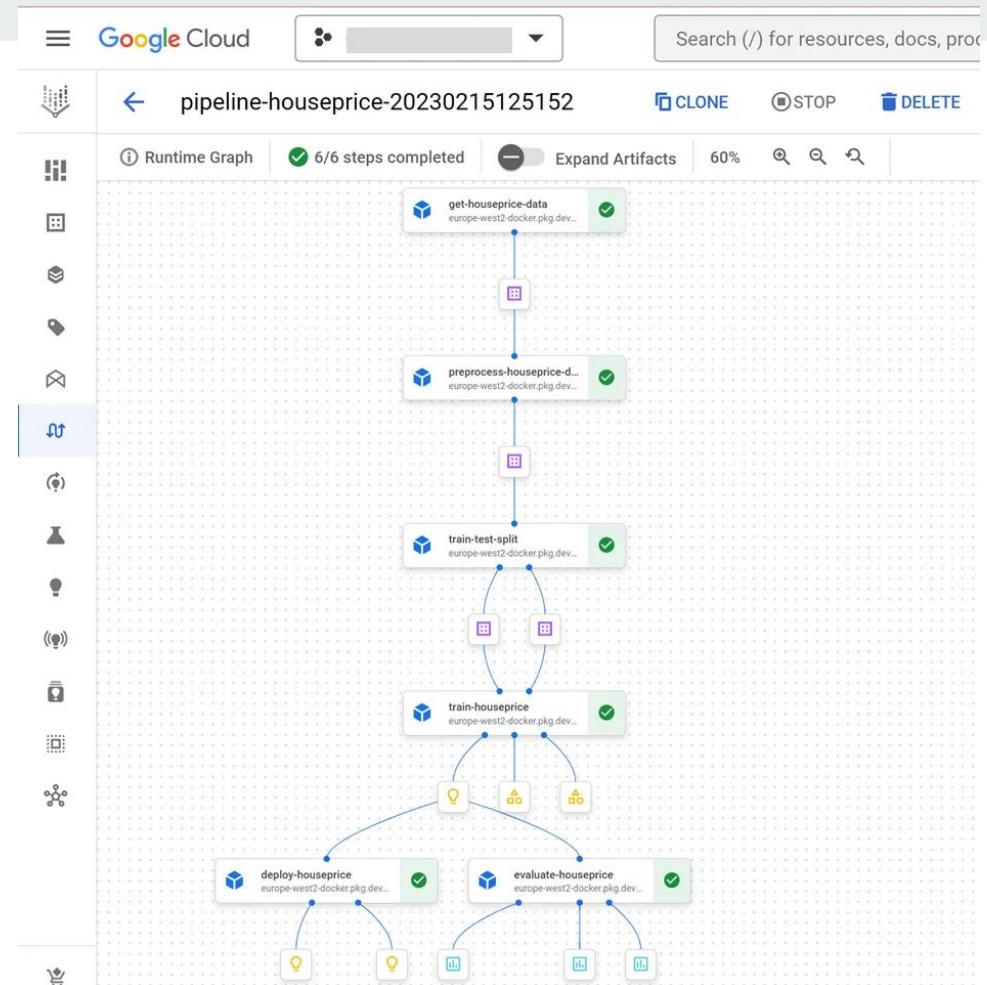
```
40     # Uploads trained model to Vertex AI Model Registry or creates new model version in
41     def upload_model():
42         listed_model = vertex_ai.Model.list(
43             filter='display_name="{}"'.format(display_name),
44             project=gcp_project,
45             location=gcp_region,
46         )
47         if len(listed_model) > 0:
48             model_version = listed_model[0] # most recently created
49             model_upload = vertex_ai.Model.upload(
50                 display_name=display_name,
51                 parent_model=model_version.resource_name,
52                 artifact_uri=str(Path(model.path).parent),
53                 serving_container_image_uri=serving_container_image_uri,
54                 location=gcp_region,
55                 serving_container_predict_route="/predict",
56                 serving_container_health_route="/health"
57         )
```



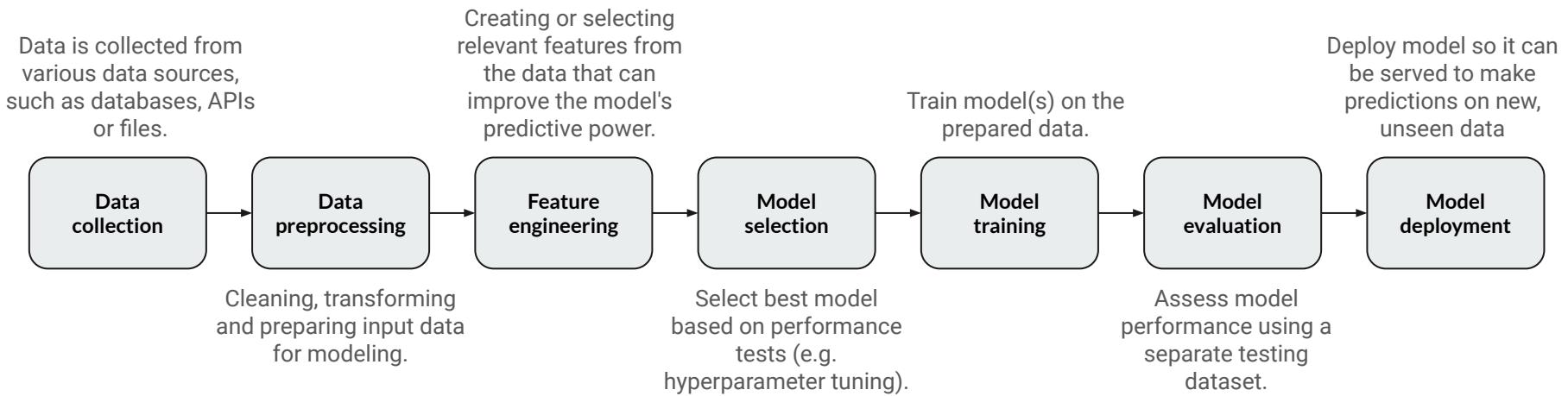
Model deployment component

Example ML Pipeline

Predict house prices with GCP Vertex AI



Typical components of an ML pipeline



Benefits of an ML Pipeline

- Modularization
- Reproducibility
- Efficiency
- Scalability
- Experimentation
- Deployment
- Collaboration
- Version control and documentation

Avoiding the endless POC loop...

- Let's change the model architecture
- Let's try with another model
- Let's gather more labeled data
- Ah wait, it was the wrong data, here is the new data, can we get results tomorrow?
- What were the accuracies again with the other model?
- Oops not sure what notebook was executed to get that model
- Could we try out the model to get feedback?
- ...

*The key to a **successful POC outcome**
is to **plan for production**
during the proof of concept.*

Triggers

Timed

- **Schedule:** Pipeline runs repetitively in relation to the creation time of the Recurring Run. Set the unit (minutes, hours, etc.) and the scalar that goes with it
- **Cron:** Takes in a cron expression, which is applied on the UTC time.

Event based

- **Manual:** Launch the pipeline manually
- **Triggered:** As part of another service (e.g. use GCP Pub/Sub)

```
#          minute (0-59)
#          hour (0-23)
#          day of the month (1-31)
#          month (1-12)
#          day of the week (0-6) (Sunday to Saturday;
#                               7 is also Sunday on some systems)
#
# * * * * * <command to execute>
```

Microservices vs ML Pipeline

	Microservices	ML Pipeline
What is it?	One full application made of different services that are calling each other to serve the overall purpose of the application.	Ran a single workflow made of components that run sequentially.
Is made of...	Services (~= APIs)	Components. Single process that is ran once (e.g. data preparation, model training, evaluation and deployment).
Utilisation	The microservices stay up.	One time run (triggered).
Purpose	All over software engineering.	ML.

ML pipeline platforms



Amazon
SageMaker



Vertex AI



Azure Machine Learning



Kubeflow

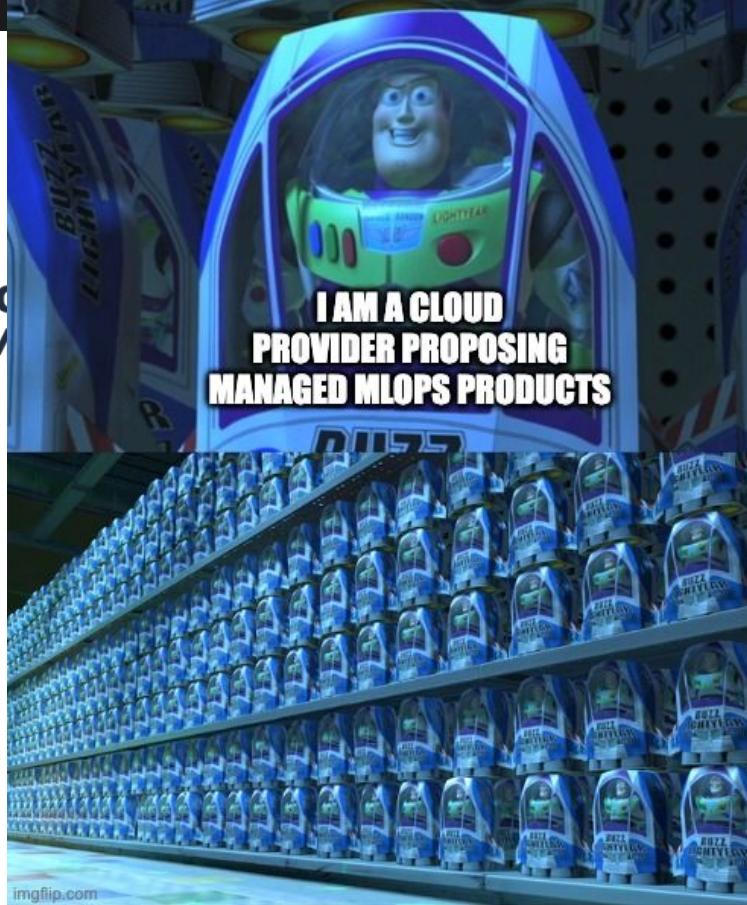
ML pipeline platforms



Amazon
SageM



Azure Machine

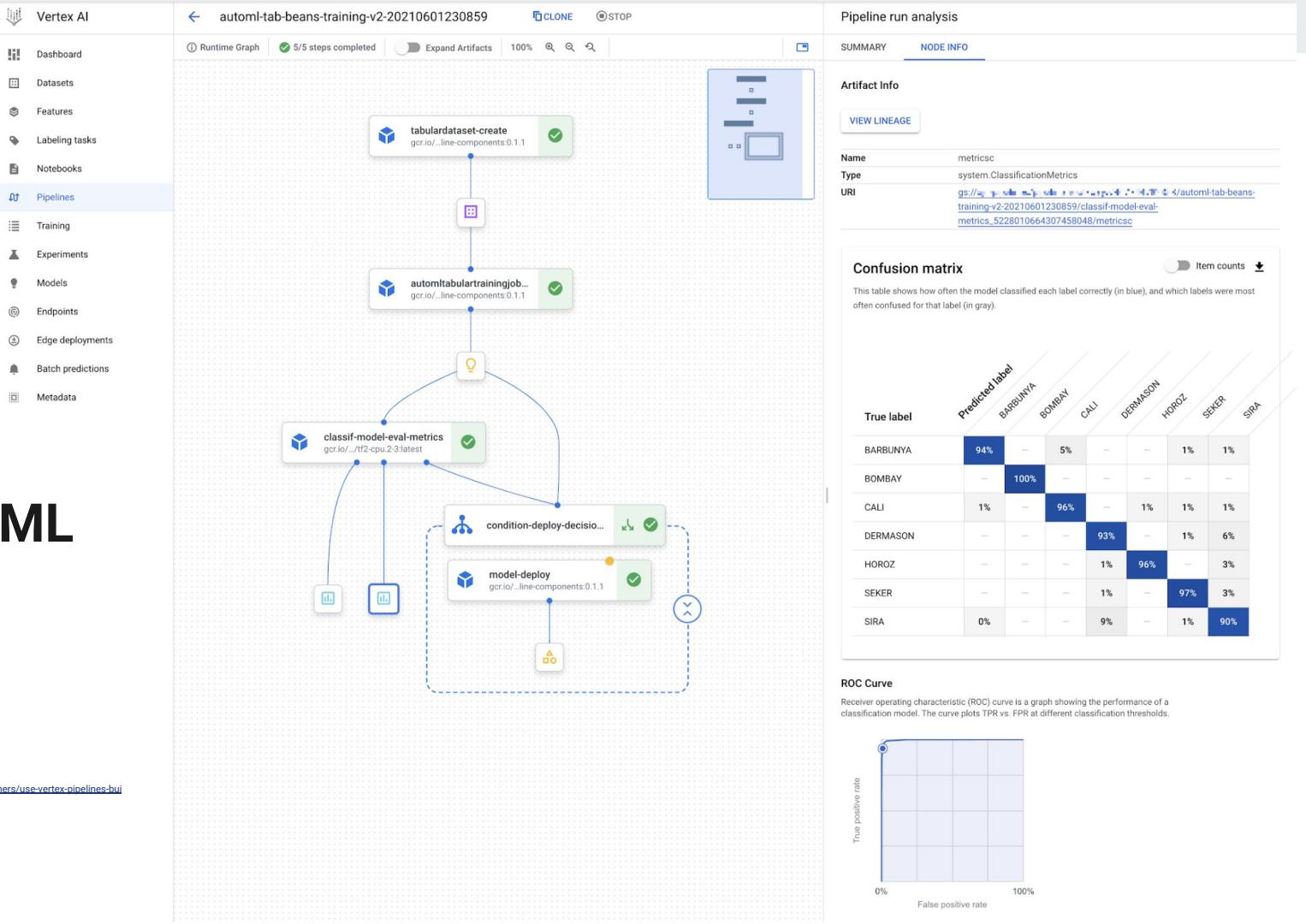


fivetran

ertex AI

Example ML Pipeline (Vertex)

<https://cloud.google.com/blog/topics/developers-practitioners/use-vertex-pipelines-build-automl-classification-end-end-workflow>



Example ML Pipeline (Kubeflow pipeline KFP)

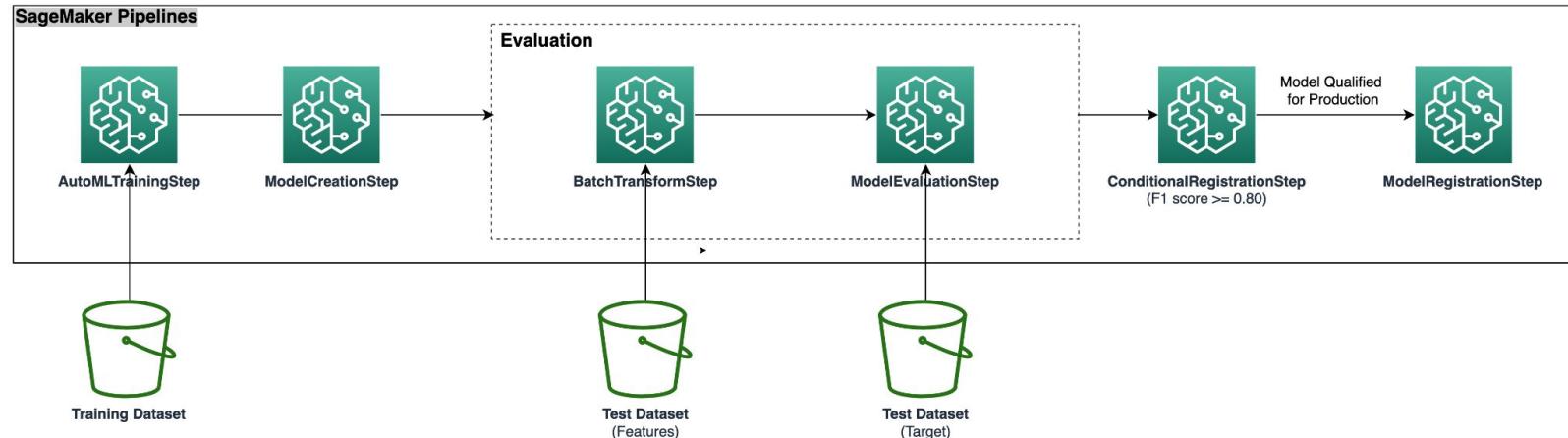
The screenshot shows the Kubeflow interface for managing machine learning pipelines. On the left, a dark sidebar menu includes options like Home, Notebooks, Tensorboards, Volumes, Models, Experiments (AutoML), Experiments (KFP), Pipelines, Runs, Recurring Runs, Artifacts, Executions, Manage Contributors, GitHub, and Documentation.

The main area displays a pipeline named "output_test" from "kubeflow-user-example-c...". The pipeline graph shows the following sequence of steps:

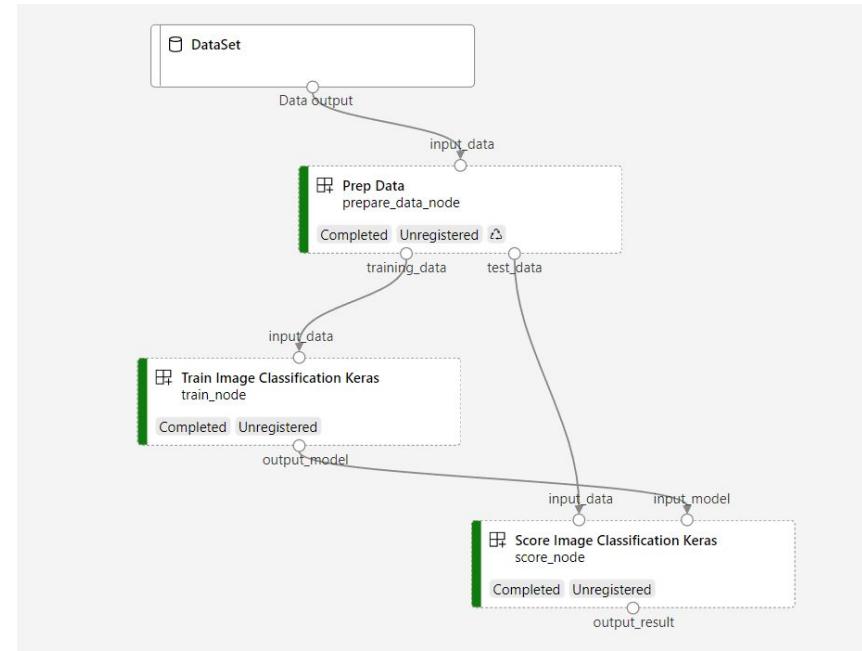
```
graph TD; A[Get latest data] --> C[Reshape data]; B[Get data batch] --> C; C --> D[Model building]; D --> E[Model serving]
```

On the right, a detailed view of the "digits-recognizer-pipeline-clkb8-780216324" run is shown. It includes tabs for Input/Output, Visualizations, Details, Volumes, Logs, Pod, Events, and ML Metadata. The Input parameters section shows "no_epochs" set to 1 and "optimizer" set to adam. The Output artifacts section lists "mlpipeline-ui-metadata", "mlpipeline-metrics", and "main-logs".

Example ML Pipeline (Sagemaker)



Example ML Pipeline (Azure ML)



ML Platforms

ML platforms

Feature	AWS SageMaker	Azure Machine Learning Service	Vertex AI
Model training	Supports a wide range of ML frameworks, including TensorFlow, PyTorch, and scikit-learn	Supports a wide range of ML frameworks, including TensorFlow, PyTorch, and scikit-learn	Supports a wide range of ML frameworks, including TensorFlow, PyTorch, and scikit-learn, as well as Google-specific frameworks like TFLite and TPUs
Model deployment	Supports a variety of deployment options, including Amazon Elastic Container Service (Amazon ECS), Amazon Elastic Kubernetes Service (Amazon EKS), and Amazon SageMaker Real-Time Inference	Supports a variety of deployment options, including Azure Kubernetes Service (AKS), Azure Container Instances (ACI), and Azure App Service	Supports a variety of deployment options, including Google Kubernetes Engine (GKE), Cloud Run, and Cloud TPUs
Model monitoring	Provides tools for monitoring model performance and detecting drift	Provides tools for monitoring model performance and detecting drift	Provides tools for monitoring model performance and detecting drift
Responsible AI	Provides tools for managing bias and fairness in ML models	Provides tools for managing bias and fairness in ML models	Provides tools for managing bias and fairness in ML models
Ease of use	Offers a variety of low-code and no-code tools to make ML accessible to users of all skill levels	Offers a variety of low-code and no-code tools to make ML accessible to users of all skill levels	Offers a variety of low-code and no-code tools to make ML accessible to users of all skill levels

ML platforms

Feature	Amazon SageMaker	Google Vertex AI	Azure Machine Learning
Compute	Amazon Elastic Compute Cloud (EC2), Amazon SageMaker managed instances	Google Compute Engine (GCE), Vertex AI managed notebooks	Azure Virtual Machines (VMs), Azure Machine Learning compute clusters
Storage	Amazon Simple Storage Service (S3)	Google Cloud Storage	Azure Blob Storage
Databases	Amazon Relational Database Service (RDS), Amazon SageMaker managed database	Cloud SQL, Vertex AI managed database	Azure SQL Database, Azure Machine Learning compute clusters
Networking	Amazon Virtual Private Cloud (VPC)	Google Cloud Virtual Private Cloud (VPC)	Azure Virtual Network (VNet)
Machine learning	Amazon SageMaker, Amazon SageMaker Canvas, Amazon SageMaker Autopilot	Vertex AI	Azure Machine Learning Studio, Azure Machine Learning notebooks
Containers	Amazon Elastic Container Service (ECS)	Google Kubernetes Engine (GKE)	Azure Kubernetes Service (AKS)
Serverless computing	AWS Lambda	Cloud Functions	Azure Functions
Big data	Amazon Elastic MapReduce (EMR), Amazon SageMaker managed hosting	Cloud Dataproc, Vertex AI managed pipelines	Azure HDInsight, Azure Machine Learning compute clusters
Analytics	Amazon Redshift, Amazon SageMaker managed analytics	BigQuery, Vertex AI managed analytics	Azure Synapse Analytics, Azure Machine Learning compute clusters

Sagemaker

Amazon SageMaker is an automated platform and comprehensive suite of tools that simplifies the development, training and deployment of machine learning (ML) models. It reduces the complexity of model development by providing a web-based interface for creating ML pipelines and pre-built algorithms.

Key features:

1. Notebooks
2. Training Job
3. Model registry
4. Prediction
5. Pipelines

Sagemaker: Notebooks

An Amazon SageMaker notebook instance is a machine learning (ML) compute instance running the Jupyter Notebook App. SageMaker manages creating the instance and related resources. Use Jupyter notebooks in your notebook instance to prepare and process data, write code to train models, deploy models to SageMaker hosting, and test or validate your models.

Advantages:

- Access more compute (GPUs)
- Access more memory
- Manage access
- Collaborate as a team
- No use of data locally



Equivalent

Vertex Workbench

Azure Notebooks

Sagemaker: Training Job

Managed way of launching a model training

1. Either be fully specify the ML algorithm, hyperparameters and input data location.
2. Or customize training job by selecting specific Sagemaker instance types and adding software libraries.

Advantages:

- Access more compute (GPUs)
- Access more memory
- Consistency
- Automation

Equivalent

Vertex Training

Azure ML Training

Sagemaker: Model registry

SageMaker Model Registry is a service for packaging model artifacts with deployment information. That information includes your deployment code, what type of container to use and what type of instance to deploy. The Model Registry decreases time to deployment, as you already have the necessary information about how the model should be deployed.

Advantages:

- Consistency
- Automation
- Scalability
- Time to deployment

Equivalent

Vertex AI Model Registry

Azure ML Model Registry

Sagemaker: Prediction

After you've trained and deployed your model in SageMaker, you can use it to generate predictions based on new data. There are two main ways to do this.

- Endpoint deployment: Deploy your model to an endpoint, allowing users and applications to send API requests to get model predictions in real time.
- Batch predictions: Generate predictions on large amounts of data without needing an immediate response..

Advantages:

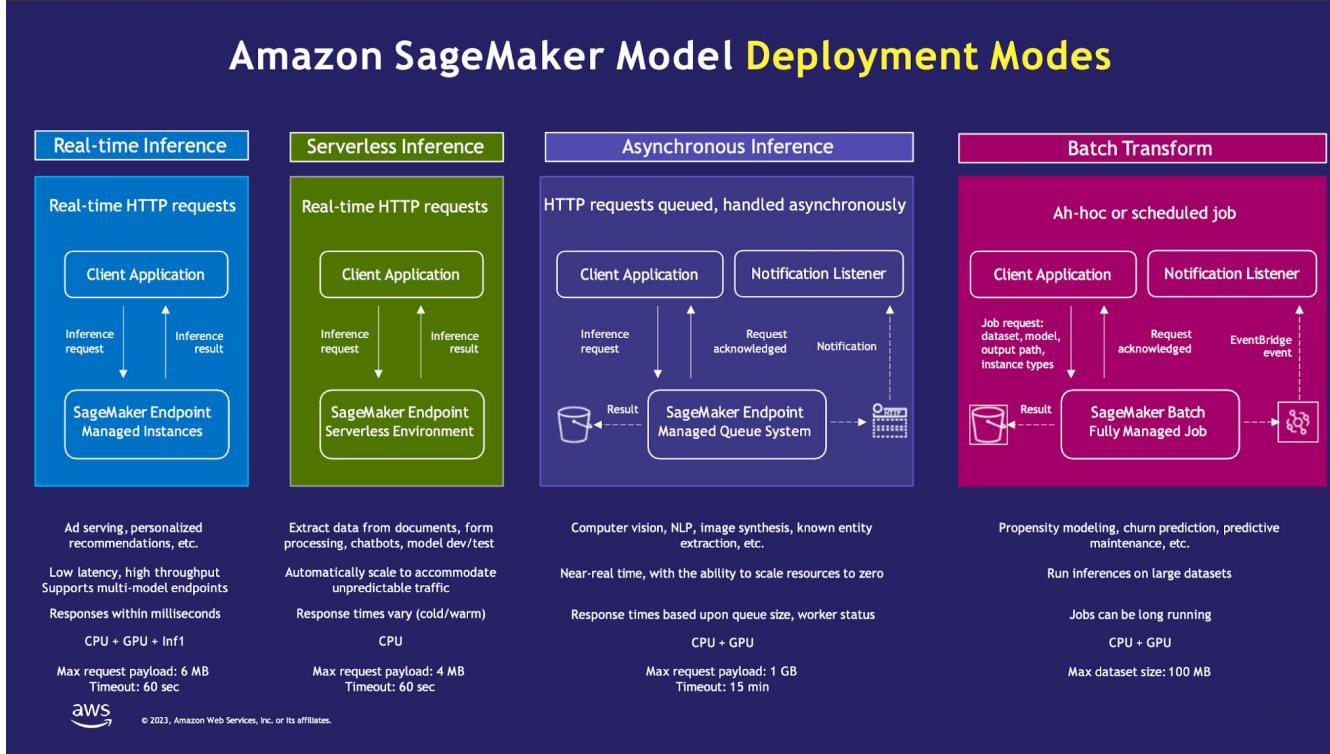
- Consistency
- Automation
- Scalability
- Time to deployment

Equivalent

Vertex AI Prediction

Azure ML Prediction

Sagemaker: Prediction



Sagemaker: Pipeline

SageMaker Pipelines is a tool for building, deploying and managing end-to-end ML workflows. Users can create automated workflows that cover data preparation, model training and deployment – all from a single interface. Because Pipeline logs everything, you can easily track and re-create models.

Advantages:

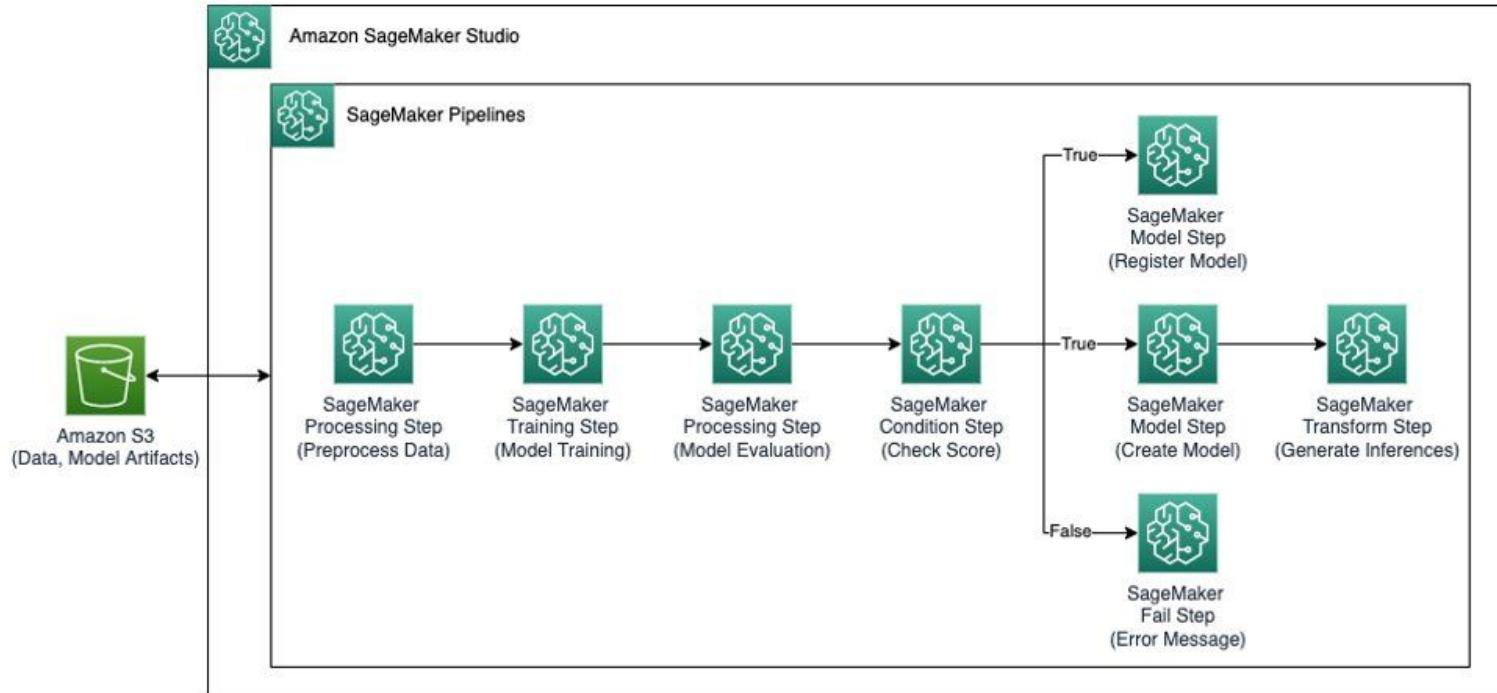
- Resource consumption
- Consistency
- Automation
- Scalability
- Time to deployment

Equivalent

Vertex AI Pipeline

Azure ML Pipeline

Sagemaker: Pipeline (example)



Wrap-up

Lecture summary

Topic	Concepts	To know for...	
		Project	Exam
Microservice architecture	• What is a Microservice ?		Yes
	• Abstractions & Separation of concerns		Yes
	• The anatomy of a microservice		
	• Monoliths aren't for dinosaurs		
Lab: Docker Compose	• Docker Compose	Optional	
ML Model pipeline	• What and why • Standard steps		Yes
Pipeline orchestrators	• Different options • Sagemaker solutions	Optional	

Project objective for sprint 4

Only optional work packages this sprint.
Catch up work from last sprint.

⚠️ Last two steps can incur cloud costs!
Make sure to not use anything you will be charged for. ⚠️

#	Week	Work package	Required
4.1	W07	Package your model training script in a Docker container . You should be able to run it locally.	Optional
4.2	W07	Run your model training as a job in the Cloud. You can implement this in different ways: <ul style="list-style-type: none">• Containerise your training script and run it on a VM in the Cloud (e.g. on EC2 or on Cloud Run, example,)• Use a managed service such as Vertex Training or Sagemaker Training	Optional
4.3	W08	Build a pipeline to automatically run different sequential components such as training your model and deploying your model. For it you can use orchestrated pipeline tools such as Kubeflow Pipelines , AWS Sagemaker , GCP Vertex .	Optional

APPENDIX

Monolith architecture

A monolithic architecture is a traditional model of a software program, which is built as a unified unit that is self-contained and independent from other applications. The word “monolith” is often attributed to something large and glacial, which isn’t far from the truth of a monolith architecture for software design. A monolithic architecture is a singular, large computing network with one code base that couples all of the business concerns together. To make a change to this sort of application requires updating the entire stack by accessing the code base and building and deploying an updated version of the service-side interface. This makes updates restrictive and time-consuming.



Monolith architecture

Pros	Cons
<ul style="list-style-type: none">• Easy deployment – One executable file or directory makes deployment easier.• Development – When an application is built with one code base, it is easier to develop.• Performance – In a centralized code base and repository, one API can often perform the same function that numerous APIs perform with microservices.• Simplified testing – Since a monolithic application is a single, centralized unit, end-to-end testing can be performed faster than with a distributed application.• Easy debugging – With all code located in one place, it's easier to follow a request and find an issue.	<ul style="list-style-type: none">• Slower development speed – A large, monolithic application makes development more complex and slower.• Scalability – You can't scale individual components.• Reliability – If there's an error in any module, it could affect the entire application's availability.• Barrier to technology adoption – Any changes in the framework or language affects the entire application, making changes often expensive and time-consuming.• Lack of flexibility – A monolith is constrained by the technologies already used in the monolith.• Deployment – A small change to a monolithic application requires the redeployment of the entire monolith.

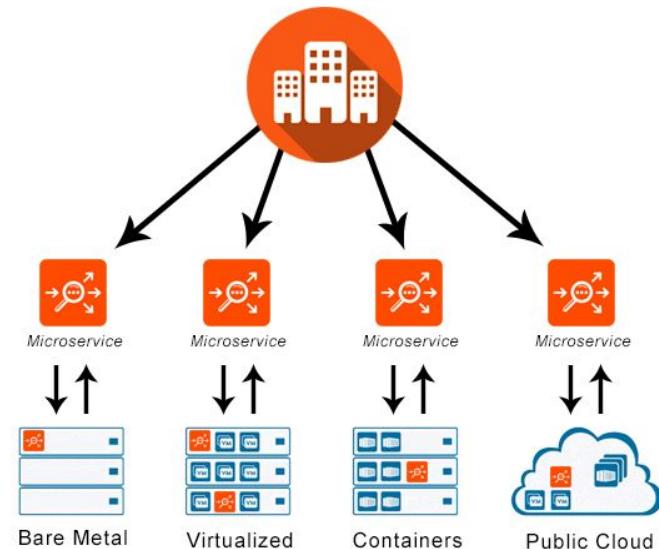
Microservice architecture

A microservices architecture, also simply known as microservices, is an architectural method that relies on a series of independently deployable services. These services have their own business logic and database with a specific goal. Updating, testing, deployment, and scaling occur within each service. Microservices decouple major business, domain-specific concerns into separate, independent code bases. Microservices don't reduce complexity, but they make any complexity visible and more manageable by separating tasks into smaller processes that function independently of each other and contribute to the overall whole.

Monolithic Architecture



Microservices Architecture

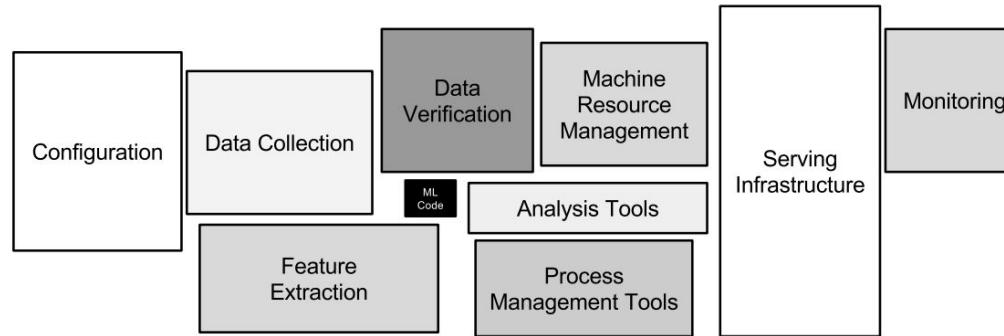


Applications

Microservice architecture

Pros	Cons
<ul style="list-style-type: none">• Agility – Promote agile ways of working with small teams that deploy frequently.• Flexible scaling – If a microservice reaches its load capacity, new instances of that service can rapidly be deployed to the accompanying cluster to help relieve pressure. We are now multi-tenant and stateless with customers spread across multiple instances. Now we can support much larger instance sizes.• Continuous deployment – We now have frequent and faster release cycles. Before we would push out updates once a week and now we can do so about two to three times a day.• Highly maintainable and testable – Teams can experiment with new features and roll back if something doesn't work. This makes it easier to update code and accelerates time-to-market for new features. Plus, it is easy to isolate and fix faults and bugs in individual services.• Independently deployable – Since microservices are individual units they allow for fast and easy independent deployment of individual features.	<ul style="list-style-type: none">• Development sprawl – Microservices add more complexity compared to a monolith architecture, since there are more services in more places created by multiple teams. If development sprawl isn't properly managed, it results in slower development speed and poor operational performance.• Exponential infrastructure costs – Each new microservice can have its own cost for test suite, deployment playbooks, hosting infrastructure, monitoring tools, and more.• Added organizational overhead – Teams need to add another level of communication and collaboration to coordinate updates and interfaces.• Debugging challenges – Each microservice has its own set of logs, which makes debugging more complicated. Plus, a single business process can run across multiple machines, further complicating debugging.• Lack of standardization – Without a common platform, there can be a proliferation of languages, logging standards, and monitoring.• Lack of clear ownership – As more services are introduced, so are the number of teams running those services. Over

Model “training” is just a part of producing an ML model



Add information from Harvard lecture

Lecture 2:

<https://github.com/Harvard-IACS/2020F-AC295/blob/master/content/lectures/lecture2/presentation/lecture2.ptx>

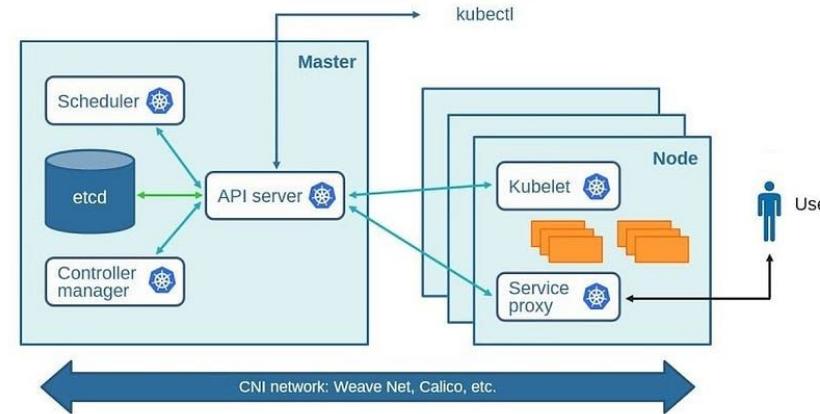
Kubernetes

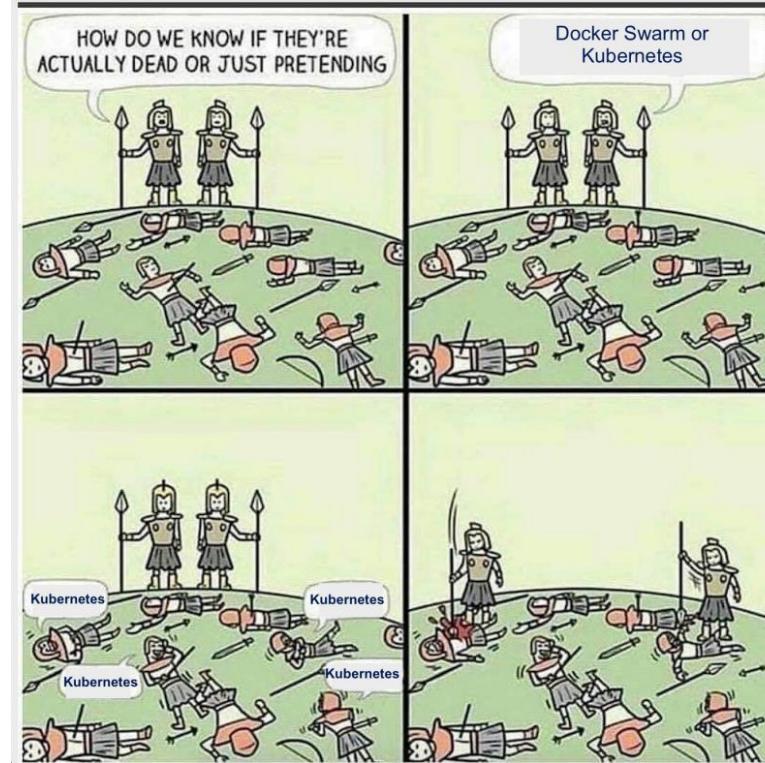
Kubernetes (k8s)

- An open source platform that has become the standard for container orchestration
 - also known as **K8s**
- Declarative API to run applications in a group of Docker hosts, called a **Kubernetes cluster**
 - Users specify in a configuration file the desired state, and Kubernetes makes it happen and then maintains it
 - “make sure four instances of this container run at all times” and Kubernetes will allocate the hosts, start the containers, monitor them, and start a new instance if one of them fails
- Major cloud providers all provide managed Kubernetes services
 - Users do not have to install and maintain Kubernetes itself
 - If an application or a model is packaged as a Docker container, users can directly submit it, and the service will provision the required machines to run one or several instances of the container inside Kubernetes

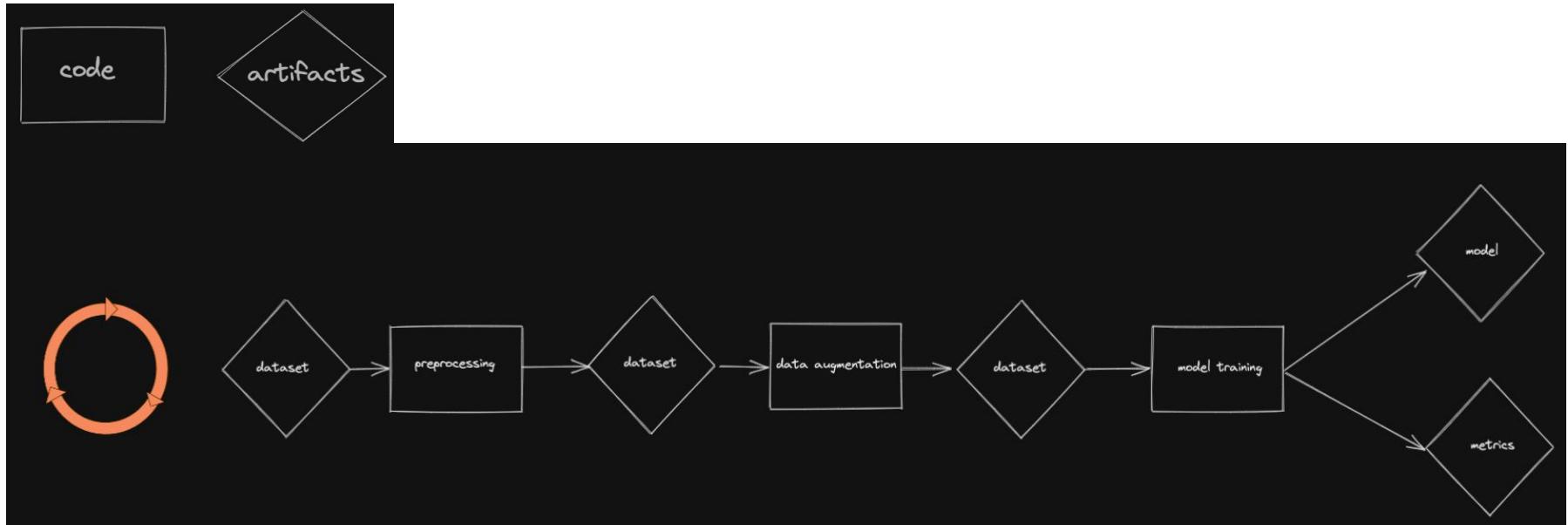
Kubernetes (k8s)

Kubernetes





ML Pipelines are made of components and iterable



[Example]

Miro's sagemaker pipeline thingy:

<https://blog.ml6.eu/promoting-code-across-environments-for-reliable-and-efficient-model-deployment-1cbe3d970b4f>

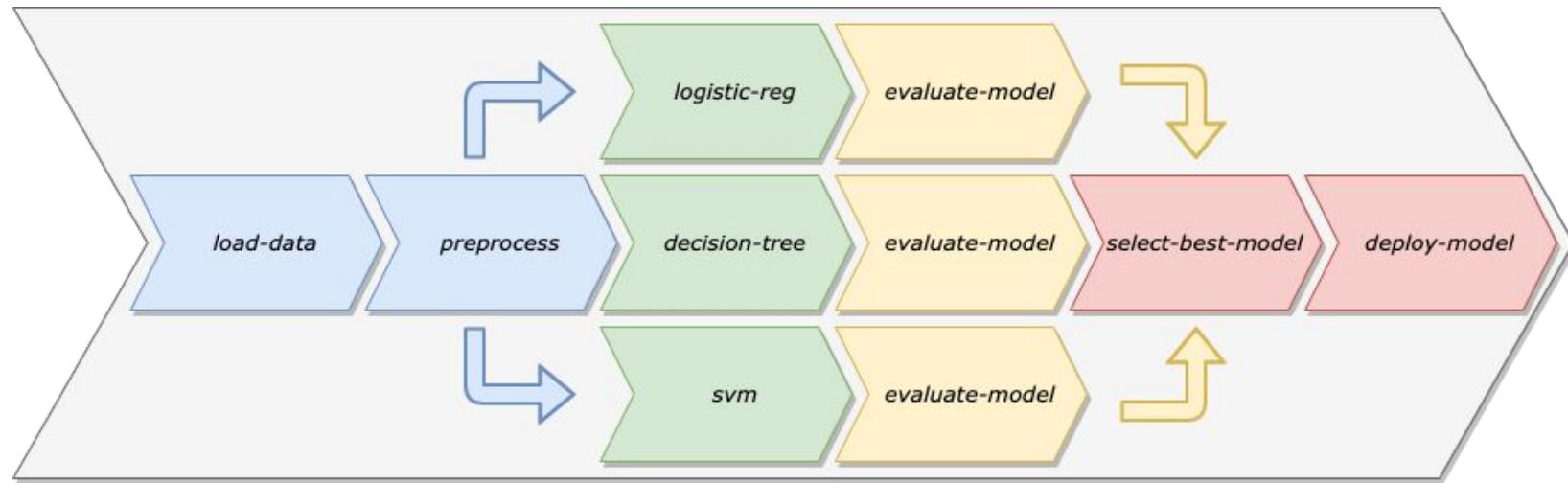
Why do we need an automated model training pipeline?

A machine learning pipeline is a step-by-step workflow for developing and deploying machine learning models into production. These pipelines allow you to streamline the process of taking raw data, training ML models, evaluating performance and integrating predictions into business applications.

Example: Customer churn pipeline

1. **Organize your customer data.** This could involve things like accessing customer data from Salesforce, cleaning it, and pre-processing it in a way that can be used by your machine learning models.
2. **Train your models.** This step involves using your data to train one or more machine learning models, which, in this case, could be used to predict customer churn. A lot of computing power is needed, typically provided by platforms like AWS and Azure.
3. **Evaluate your models.** Once you've trained your models, it's important to evaluate their performance to see how well they are actually doing. This will help you fine-tune your models and make sure they are as effective as possible.
4. **Deploy your models.** After you've evaluated and fine-tuned your models, you can deploy them in production so they can start making predictions (in this case, predictions about customer churn).

There are steps in ML pipeline



Standard steps of a ML pipeline

1. Data collection
 - a. Coming from different sources
 - b. Different formats...
2. Data preparation
 - a. Clean data (remove duplicates, normalise, ...)
 - b. Split data
3. Feature engineering
 - a. Prepare your features
4. Model training
 - a. Train the model
5. Optimisation
 - a. Run some hyperparameter optimisation
6. Evaluation
 - a. Compute and store performance metrics
7. Model deployment
 - a. Store model in a registry
 - b. Deploy model to make inference
8. Predict
 - a. Sometimes compute live predictions (e.g. time series data)