# Deployment

## Prototyping with Deep Learning

# Learning outcomes

- Identify best practices in model deployment

- Save, load, and use your DL models

# What can we do with Jupyter notebooks?

Add interactive widgets

Create slideshow presentations

Write technical blogs

… But software **deployment** is something else

# Goals
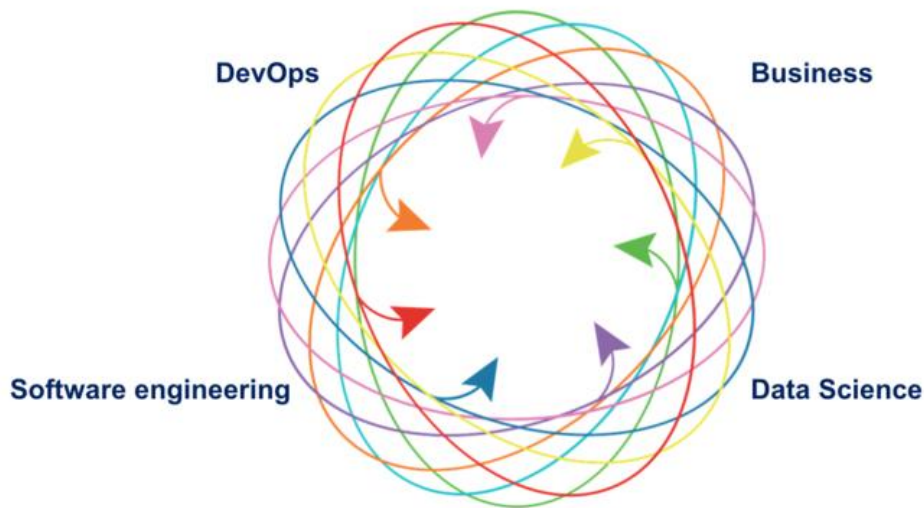
Use trained models in **production** environments

Also monitor models usage over time

Also improve models with new data
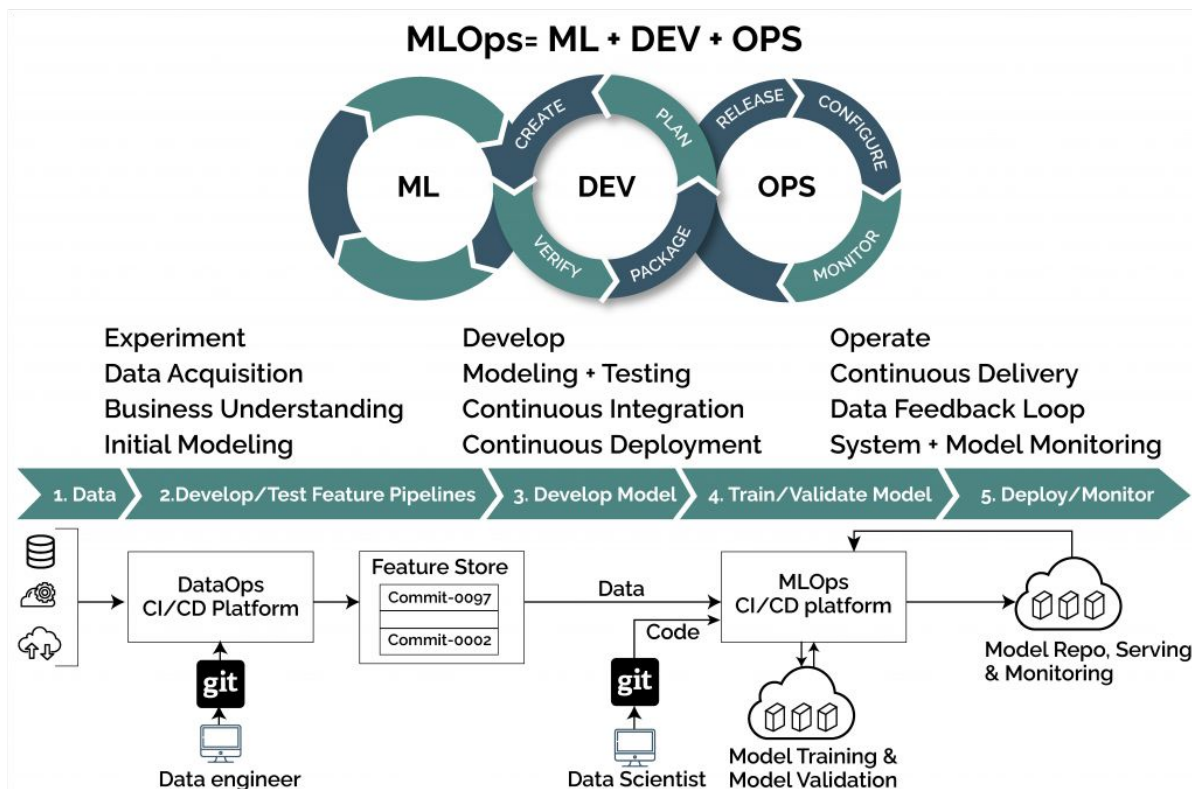
# Why are ML systems hard?

Some *key* issues:

- Entanglement
- Data dependencies
- Configuration
- Data preparation
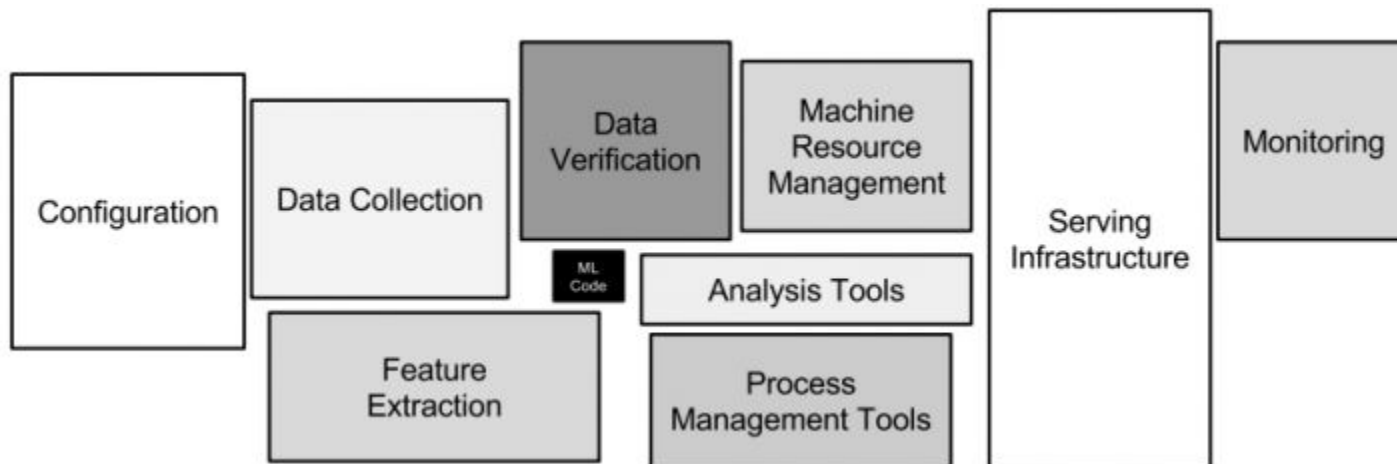- Tracking model errors
- Separation of expertise



https://christophergs.github.io/machine%20learning/2019/03/17/how-to-deploy-machine-learning-models/

# MLOps pipeline



MLOps= ML + DEV + OPS

| ML | DEV | OPS |
| --- | --- | --- |
| Experiment | Develop | Operate |
| Data Acquisition | Modeling + Testing | Continuous Delivery |
| Business Understanding | Continuous Integration | Data Feedback Loop |
| Initial Modeling | Continuous Deployment | System + Model Monitoring |

1. Data → 2.Develop/Test Feature Pipelines → 3. Develop Model → 4. Train/Validate Model → 5. Deploy/Monitor

DataOps CI/CD Platform → Feature Store (Commit-0097, Commit-0002) → Data/Code → MLOps CI/CD platform → Model Repo, Serving & Monitoring

Data engineer — git
Data Scientist — git
Model Training & Model Validation

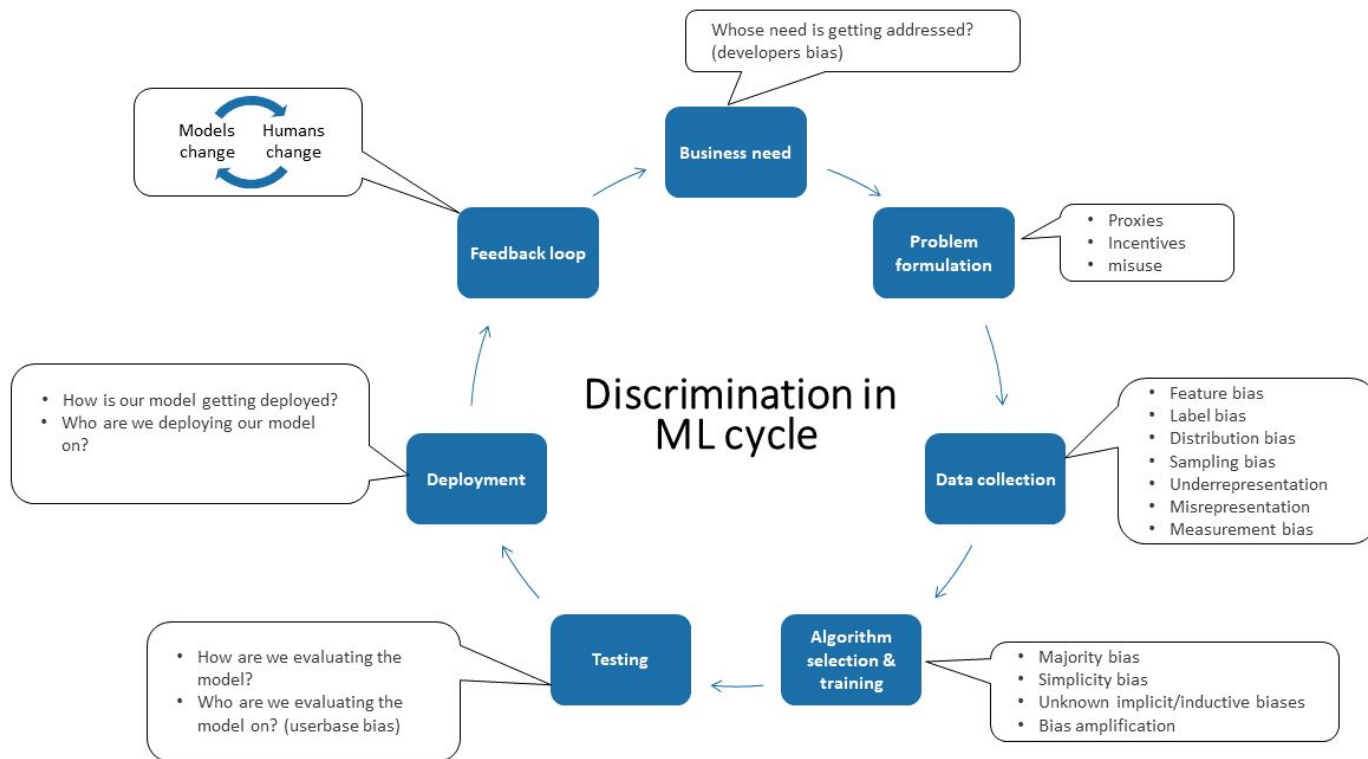# The Hidden Technical Debt in ML Systems



https://papers.nips.cc/paper/5656-hidden-technical-debt-in-machine-learning-systems.pdf

# Biases everywhere



Discrimination in ML cycle

- Business need — Whose need is getting addressed? (developers bias)
- Problem formulation — Proxies, Incentives, misuse
- Data collection — Feature bias, Label bias, Distribution bias, Sampling bias, Underrepresentation, Misrepresentation, Measurement bias
- Algorithm selection & training — Majority bias, Simplicity bias, Unknown implicit/inductive biases, Bias amplification
- Testing — How are we evaluating the model? Who are we evaluating the model on? (userbase bias)
- Deployment — How is our model getting deployed? Who are we deploying our model on?
- Feedback loop — Models change / Humans change

https://twitter.com/fereshte_khani/status/1496281636924960770
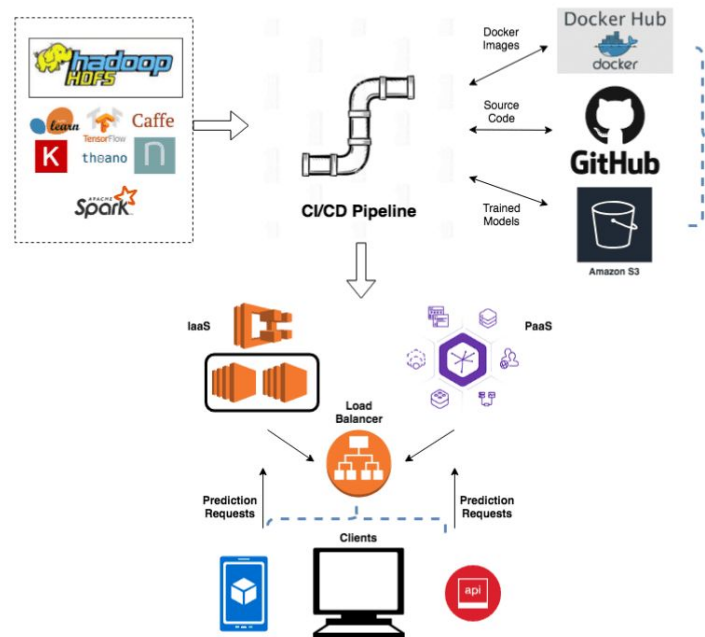
# Key design principles

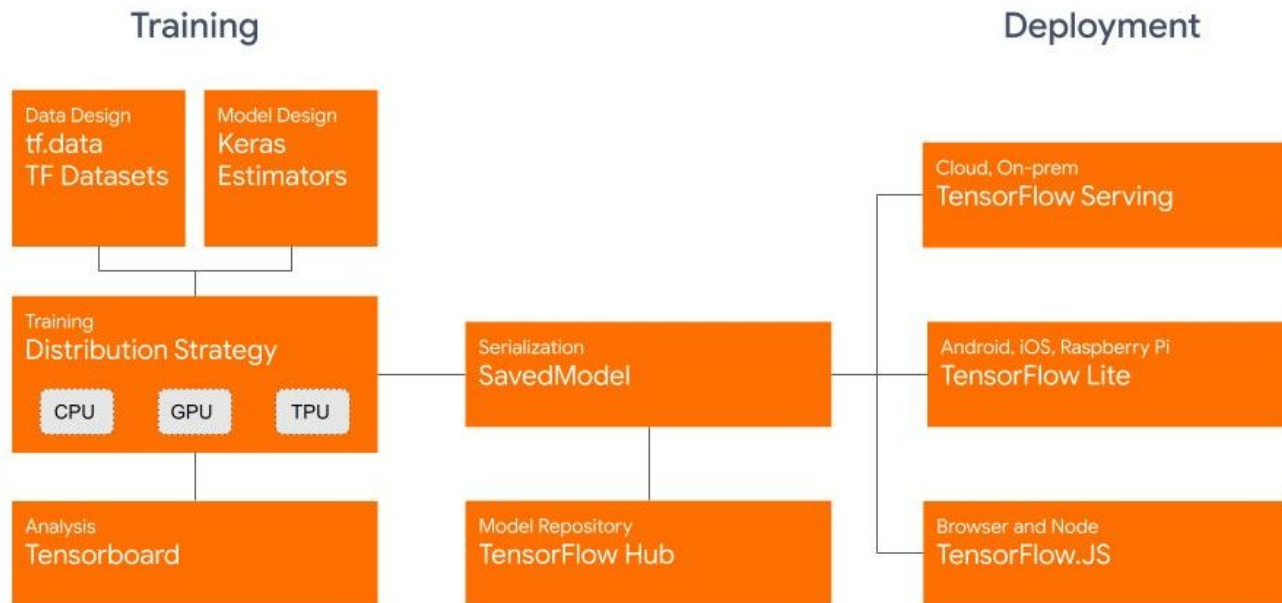Plan for **reproducibility**

ML as part of your **build pipeline**

Plan for **extensibility**

Plan for **modularity**

Testing, testing, testing

# Some TF tools

# Saving trained models

**In Keras**

`model.save(…)`

Only architecture: `model.to_json()` or `model.to_yaml()`

Only weights: `model.save_weights()`

**In Pytorch**

`torch.save(model.state_dict(), …)`

https://keras.io/getting-started/faq/#how-can-i-save-a-keras-model

https://pytorch.org/tutorials/beginner/saving_loading_models.html

# What can (and should) be saved?

Model architecture

Model weights

Training configuration: loss function, epochs

Optimizer: learning rate, state (to resume training)

# Model file formats

`.pt` or `.pth`                Pytorch

`.pb`                          Protobuf (TF native)

`.tflite`                      Tensorflow, on-device

`.h5`                          HDF, cross-platform

`.json`                        TF.js

Don't use Python's `pickle`!

# Model format conversion

https://github.com/tensorflow/tfjs/tree/master/tfjs-converter

```
~$ tflite_convert --graph_def_file=model.pb --output_file=model.tflite

~$ tflite_convert --keras_model_file=model.h5 --output_file=model.tflite

~$ tensorflowjs_converter --input_format keras model.h5 /path/to/dir
```

| Name | ✕ Headers Preview Response Timing |
|------|-----------------------------------|
| ☐ messages.json | ▼ General |
| ☐ model.json | Request URL: http://localhost:1234/json/model.json |
| ☐ group1-shard1of1 | Request Method: GET |
| ☐ group2-shard1of1 | Status Code: 🟢 200 OK |
| ☐ group3-shard1of1 | |

# Model serving with Flask (web APIs)

```python
from flask import Flask, request, jsonify

from tensorflow.keras.models import load_model
# from torch import load as load_model

model = load_model('model.h5', compile=False)

app = Flask(__name__)

@app.route('/predict', methods=['POST'])
def predict():
    # TODO: Read data from `request` object.
    return jsonify(model.predict(...))

# Init service.
app.run(port=50005)
```

https://flask.palletsprojects.com/en/2.1.x/

# Model serving with Gradio (graphical UIs)

```python
import gradio as gr
import numpy as np
from flash.image import ImageClassifier

# 1. Load Model
model = ImageClassifier.load_from_checkpoint("image_classification_model.pt")

# 2. Define Classification Function
def classify(img):
    img = np.transpose(img, (2, 0, 1))
    return model.predict(img, data_source="numpy")[0]

#Init Gradio
image = gr.inputs.Image(shape=(299, 299))
label = gr.outputs.Label(num_top_classes=1)
gr.Interface(fn=classify, inputs=image, outputs=label, capture_session=True).launch()
```

https://gradio.app/getting_started/

# Monitoring



https://wandb.ai/site

# Some resources

Best practices:
https://opendatascience.com/best-practices-for-deploying-machine-learning-in-the-enterprise/

Challenges: https://towardsdatascience.com/443af67493cd

CI/CD in ML: https://martinfowler.com/articles/cd4ml.html

DL in production (book):
https://github.com/The-AI-Summer/Deep-Learning-In-Production