

IOWA STATE UNIVERSITY

Translational AI Center

Introduction to MLOps

Infrastructure and Tooling

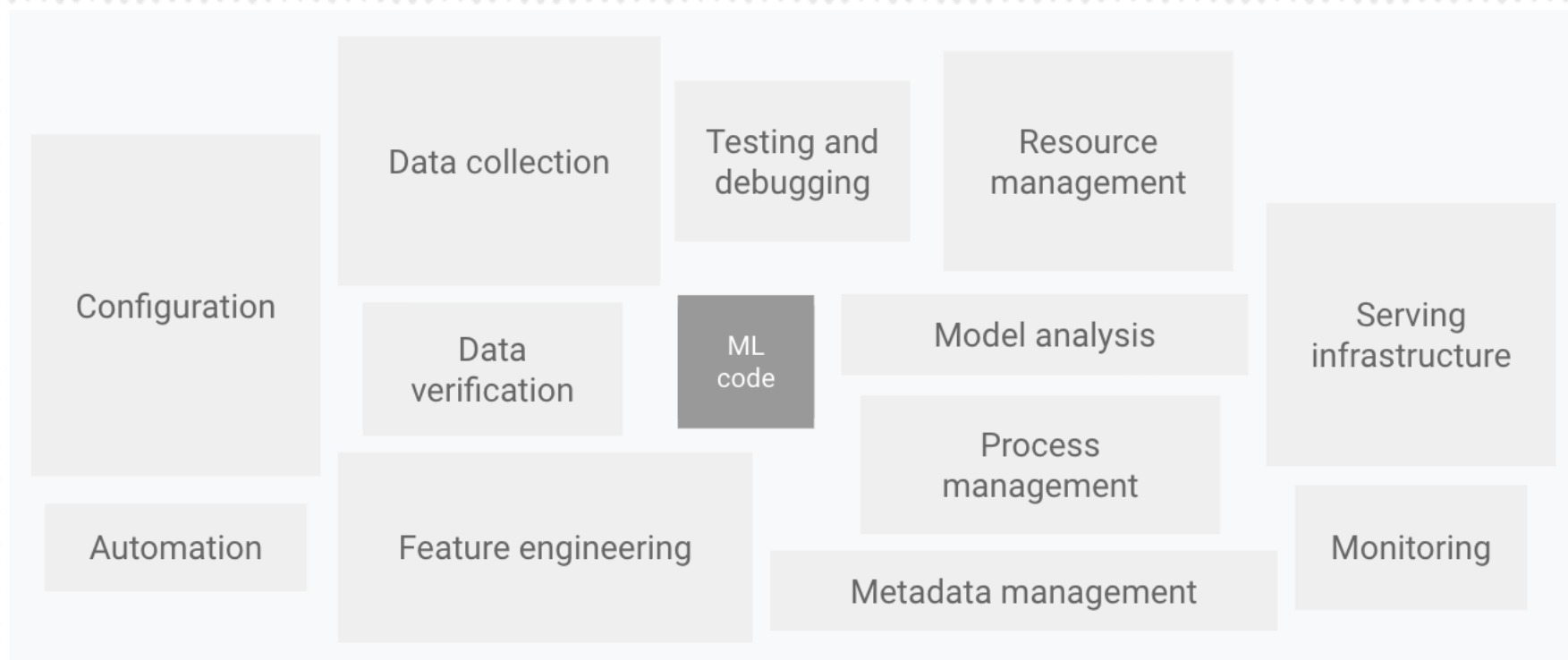
Objectives

- Introduction to MLOps
 - Dream vs. Reality for ML Practitioners
 - The 3 buckets of ML Infrastructure
 - Software Engineering
 - Compute Hardware
 - Resource Management
 - Frameworks and Distributed Training
 - All-in-one Solutions
 - Maturity Model

What is MLOps

- Apply [DevOps](#) principles to ML systems (MLOps).
- *MLOps* is an ML engineering culture and practice that aims at unifying ML system development (Dev) and ML system operation (Ops).
- Practicing MLOps means that you advocate for automation and monitoring at all steps of ML system construction, including integration, testing, releasing, deployment and infrastructure management.

MLOps - Overview



DevOps

- [DevOps](#) is a popular practice in developing and operating large-scale software systems. This practice provides benefits such as shortening the development cycles, increasing deployment velocity, and dependable releases. To achieve these benefits, you introduce two concepts in the software system development:
 - [Continuous Integration \(CI\)](#)
 - [Continuous Delivery \(CD\)](#)

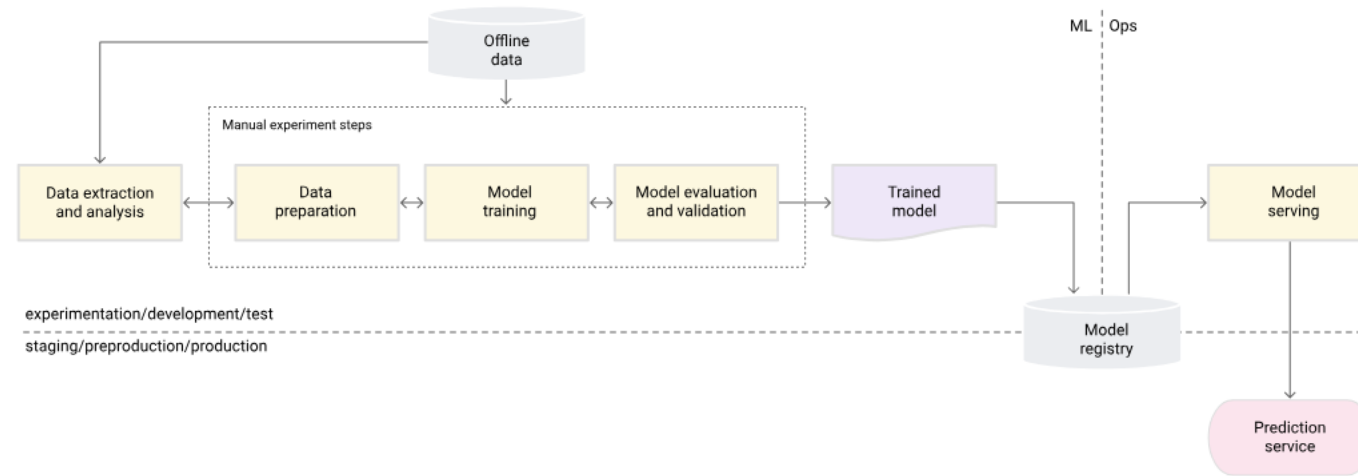
MLOPS vs Devops

- However, ML systems differ from other software systems in the following ways:
- **Team skills:** In an ML project, the team usually includes data scientists or ML researchers, who focus on exploratory data analysis, model development, and experimentation. These members might not be experienced software engineers who can build production-class services.
- **Development:** ML is experimental in nature. You should try different features, algorithms, modeling techniques, and parameter configurations to find what works best for the problem as quickly as possible. The challenge is tracking what worked and what didn't, and maintaining reproducibility while maximizing code reusability.
- **Testing:** Testing an ML system is more involved than testing other software systems. In addition to typical unit and integration tests, you need data validation, trained model quality evaluation, and model validation.
- **Deployment:** In ML systems, deployment isn't as simple as deploying an offline-trained ML model as a prediction service. ML systems can require you to deploy a multi-step pipeline to automatically retrain and deploy model. This pipeline adds complexity and requires you to automate steps that are manually done before deployment by data scientists to train and validate new models.
- **Production:** ML models can have reduced performance not only due to suboptimal coding, but also due to constantly evolving data profiles. In other words, models can decay in more ways than conventional software systems, and you need to consider this degradation. Therefore, you need to track summary statistics of your data and monitor the online performance of your model to send notifications or roll back when values deviate from your expectations.

MLOps Maturity Model

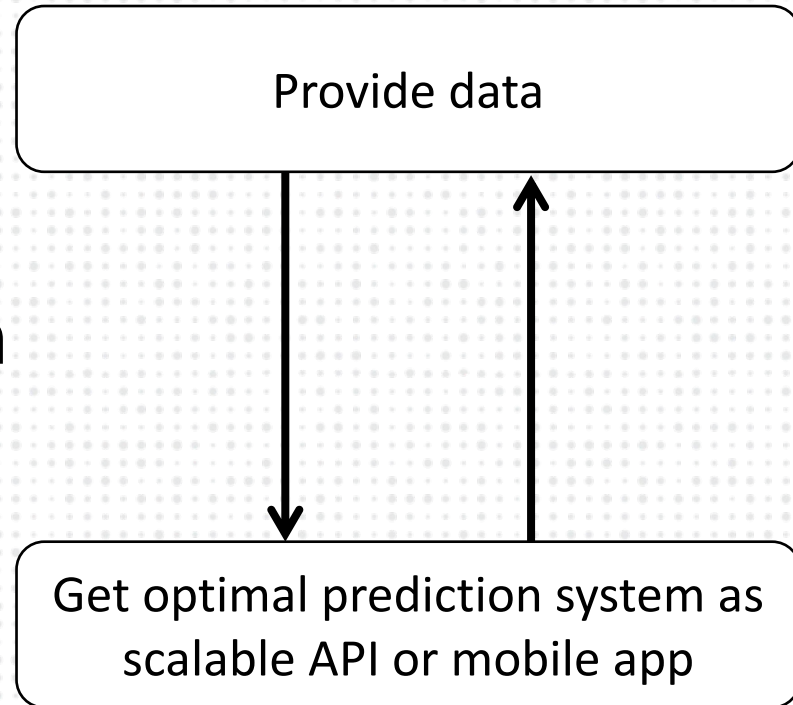
- Helps to clarify the Development Operations (DevOps) principles and practices necessary to run a successful MLOps environment.
- Encompasses five levels of technical capability
 - Level 0: No MLOps
 - Level 1: DevOps but no MLOps
 - Level 2: Automated Training
 - Level 3: Automated Model Deployment
 - Level 4: Full MLOps Automated Operations

MLOps – Level 0

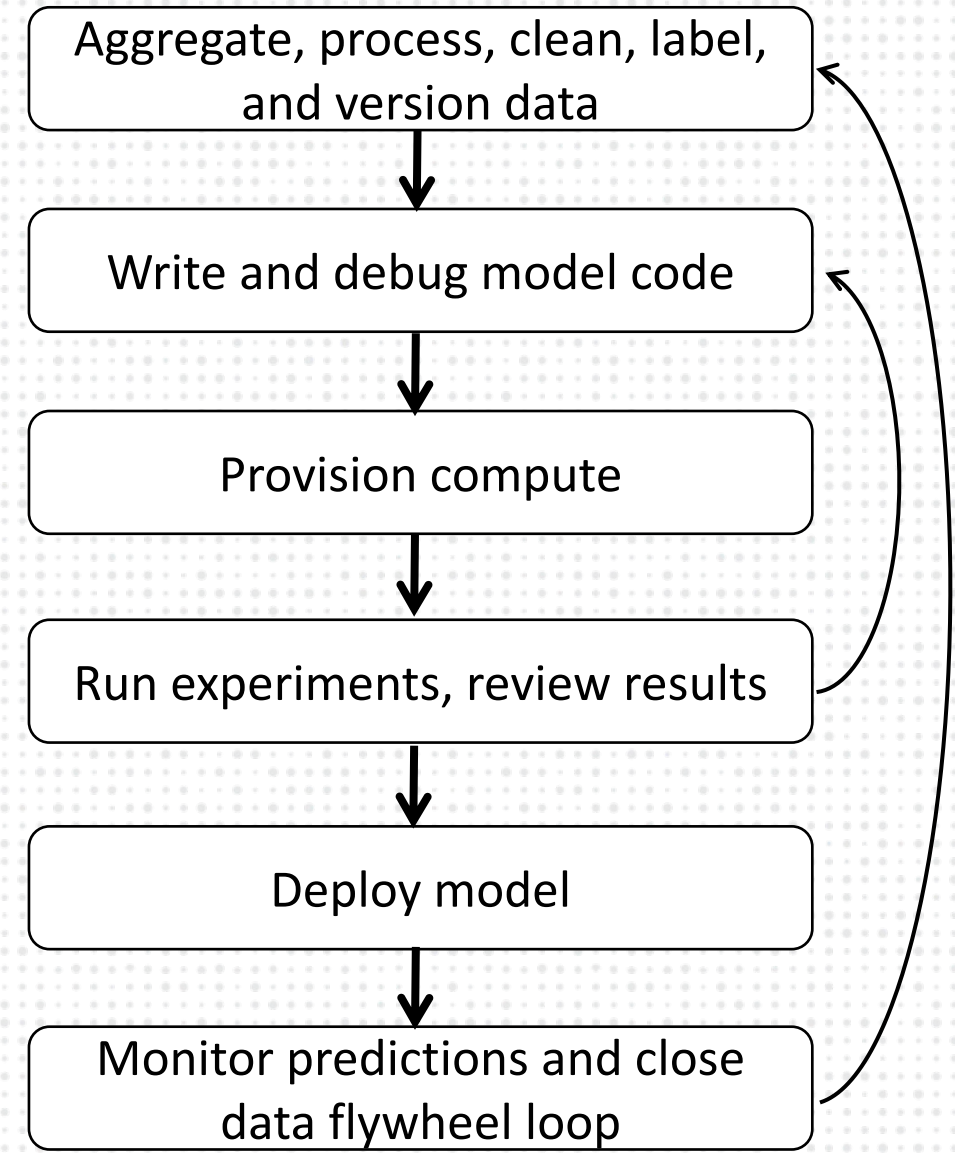


Dream vs. Reality for ML Practitioners

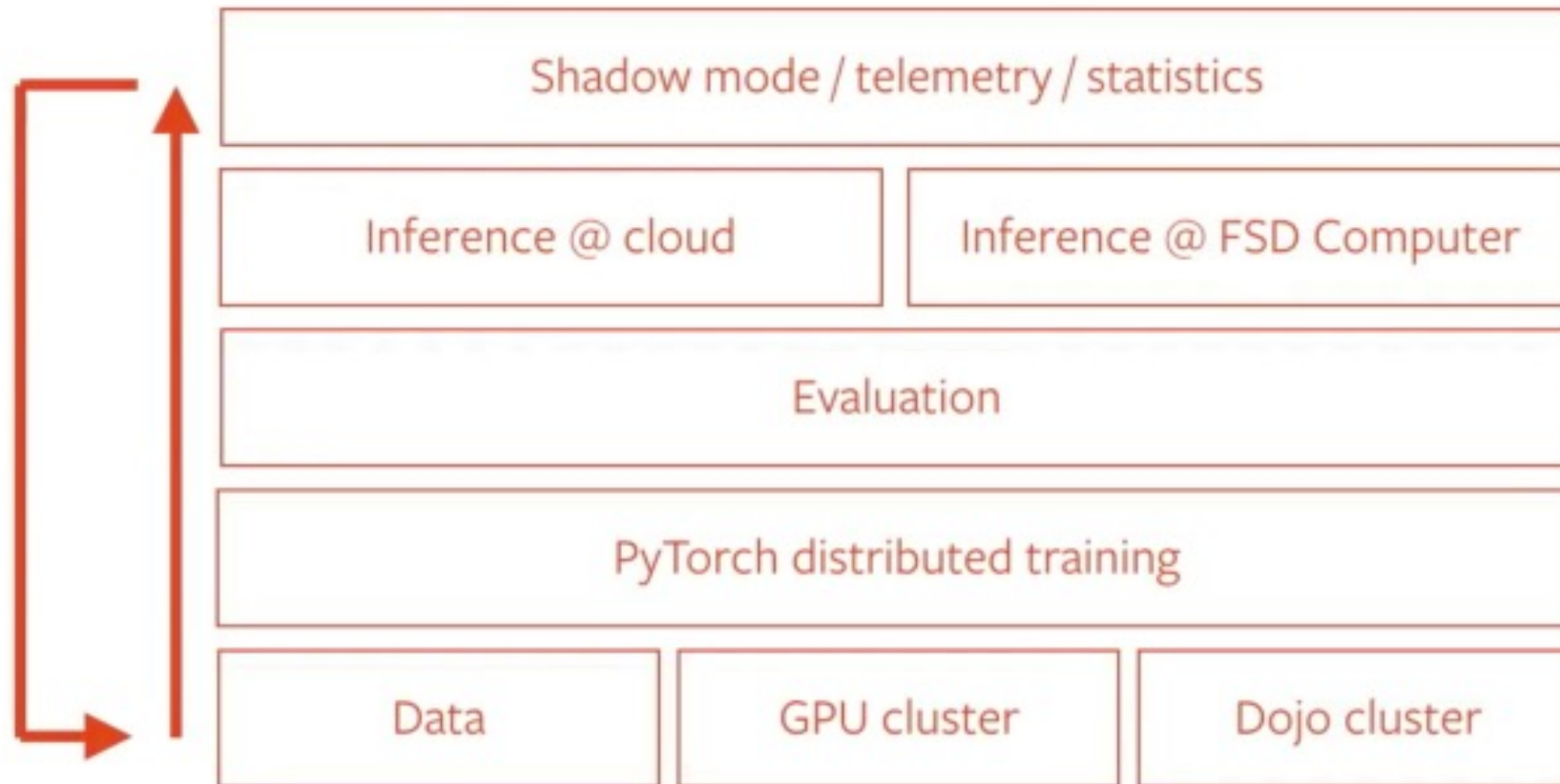
Dream



Reality

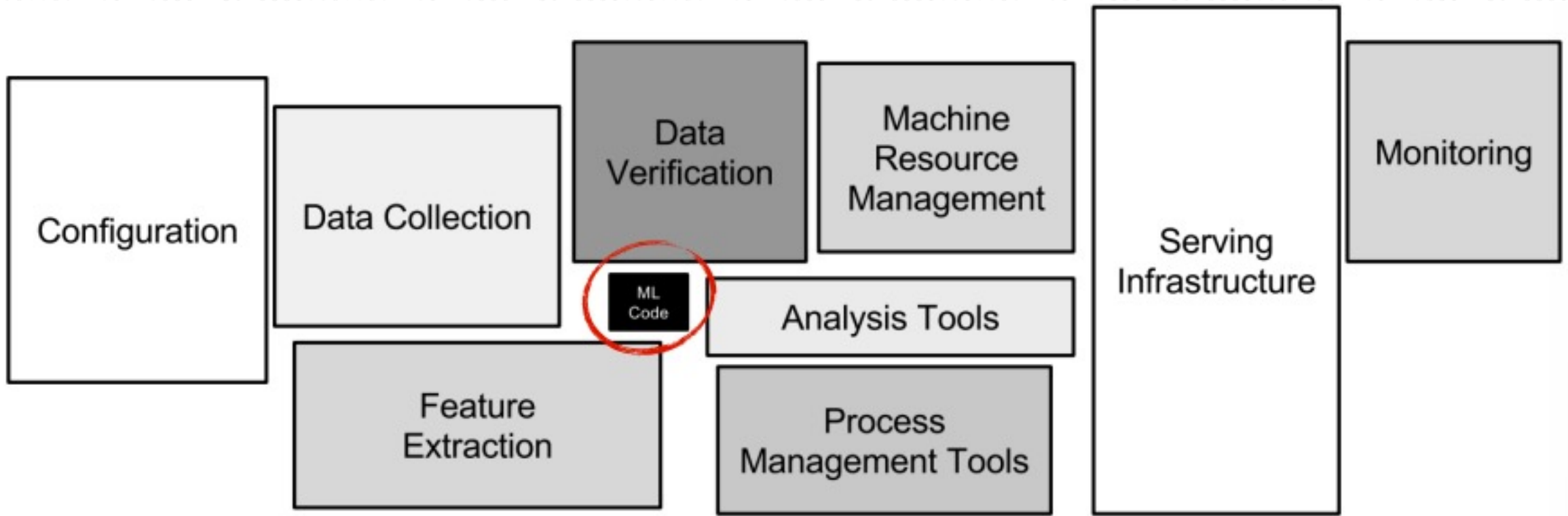


“OPERATION VACATION”



Goal: add data, see model improve

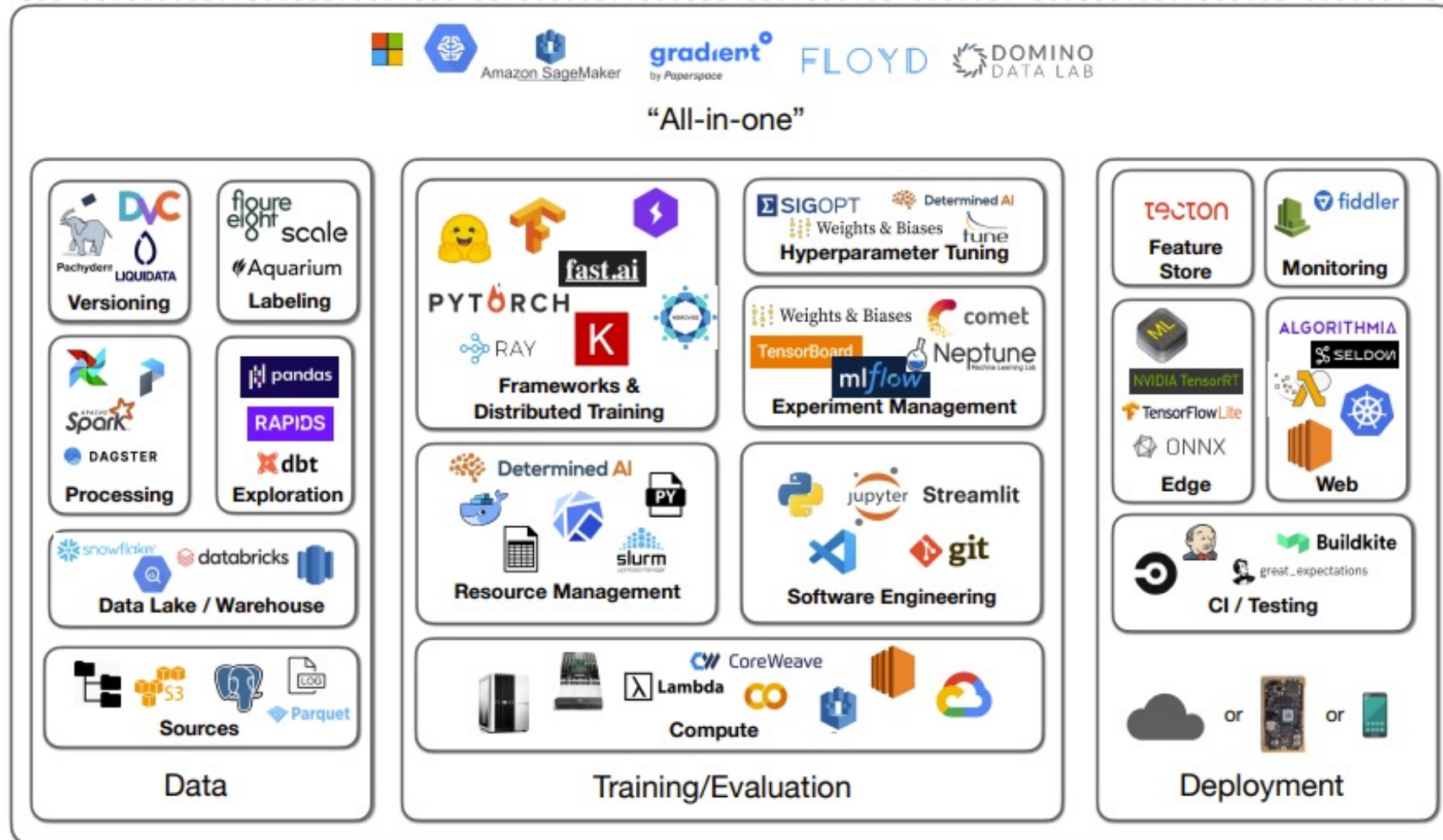
Andrej Karpathy at PyTorch Devcon 2019 - <https://www.youtube.com/watch?v=oBklItKXtDE>

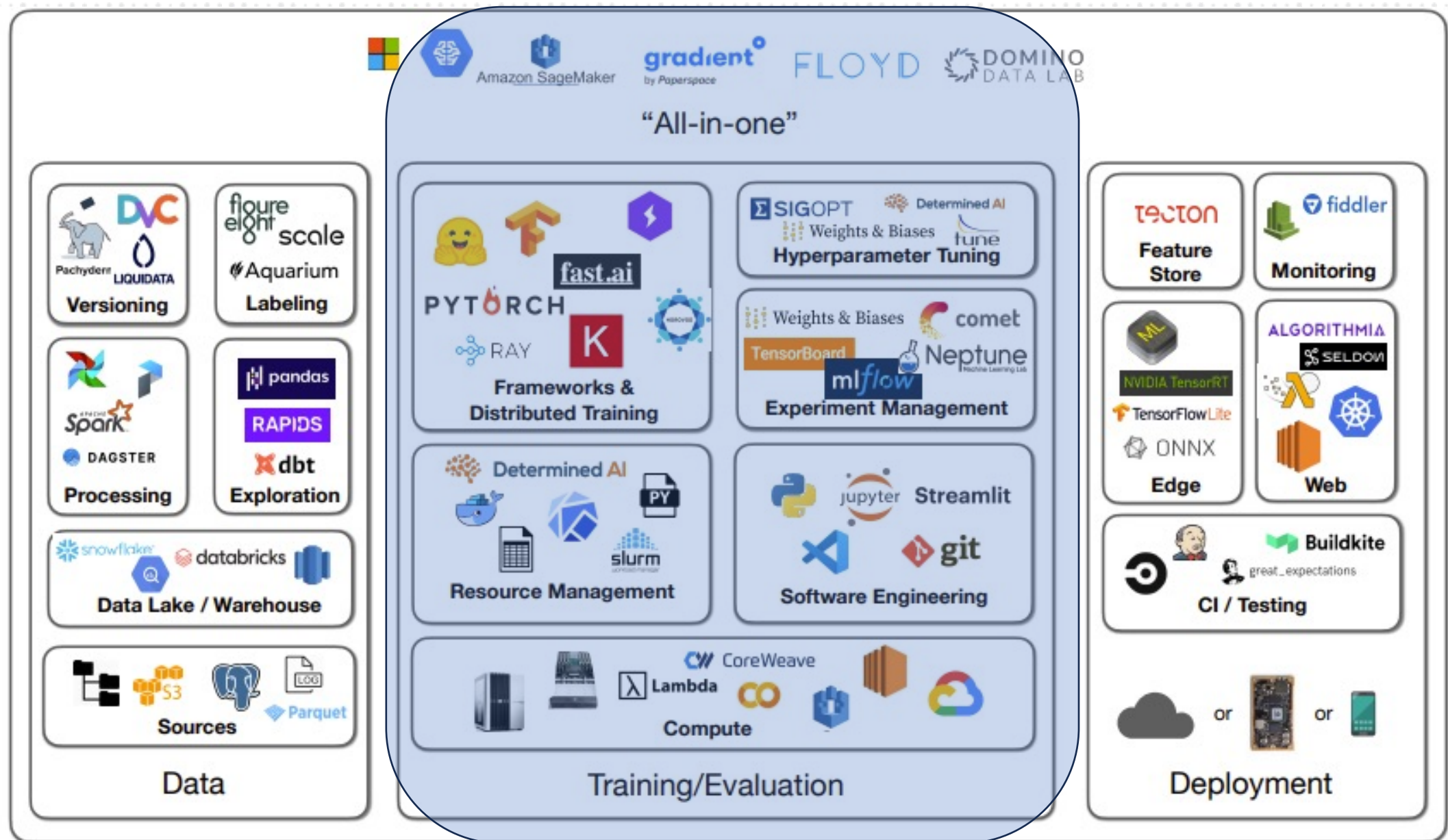


Machine Learning: The High-Interest Credit Card of Technical Debt

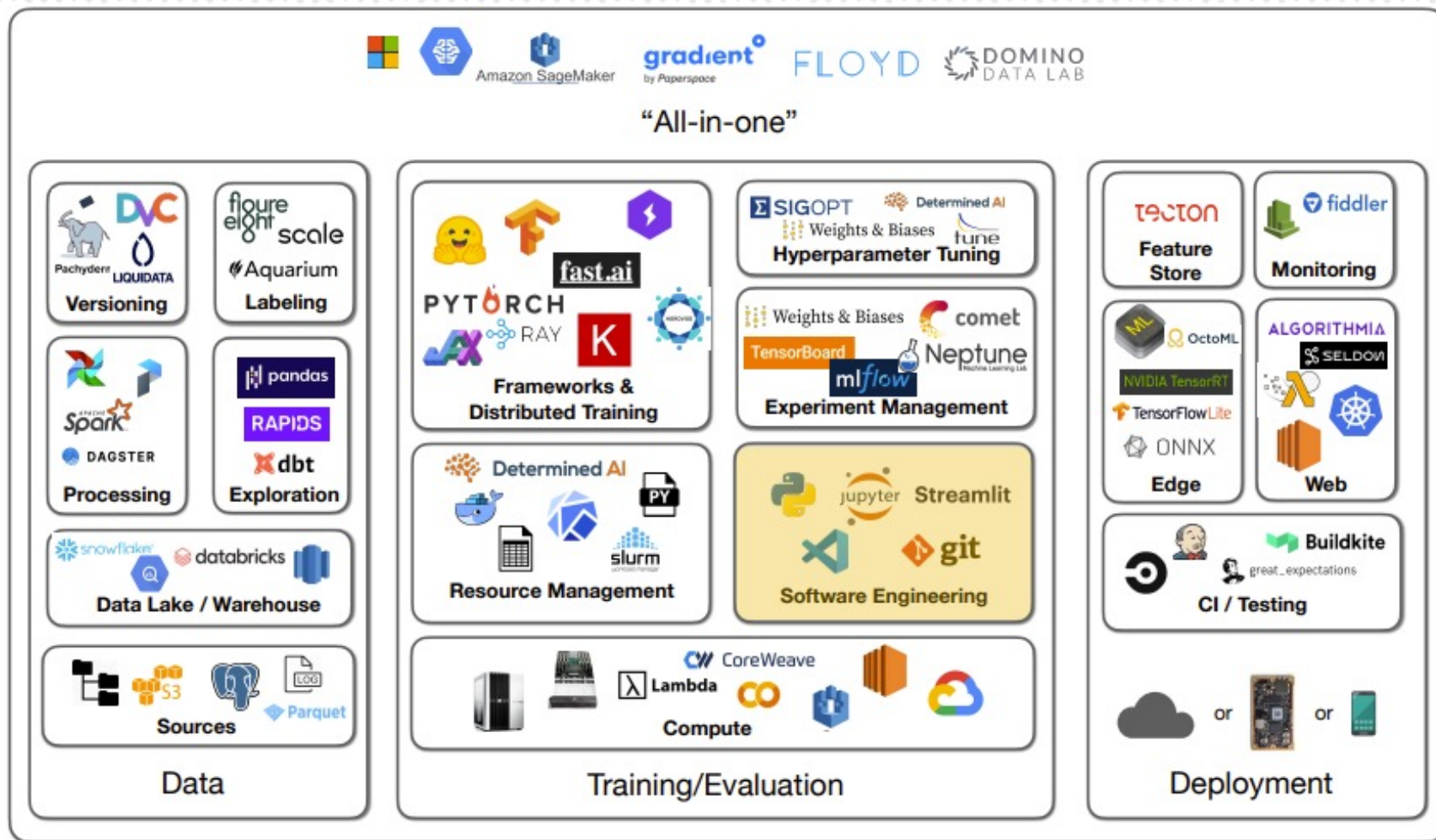
D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov,
Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young

The 3 buckets of ML Infrastructure





Software Engineering



Programming Language

- Python, because of the libraries
 - Clear winner in scientific and data computing

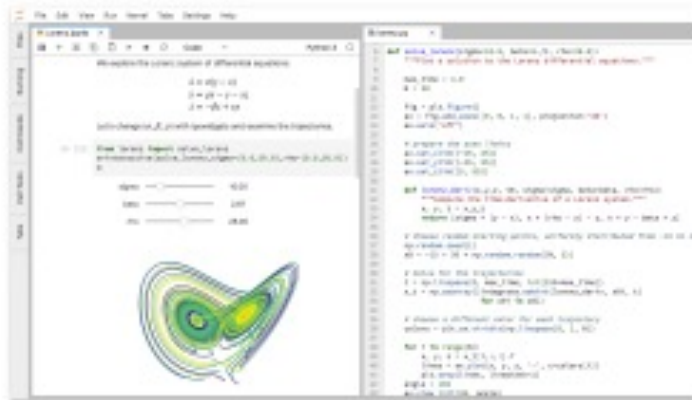
Editors



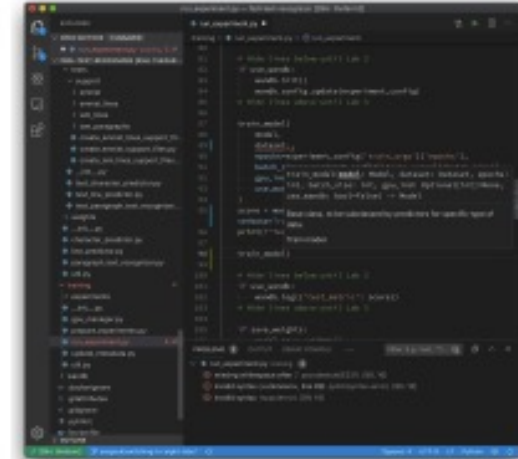
Vim



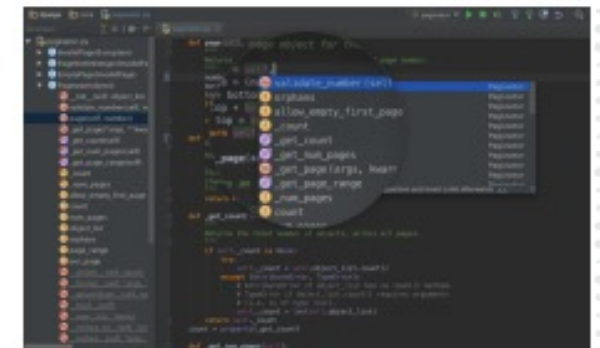
Emacs



Jupyter



VS Code



PyCharm

Visual Studio Code

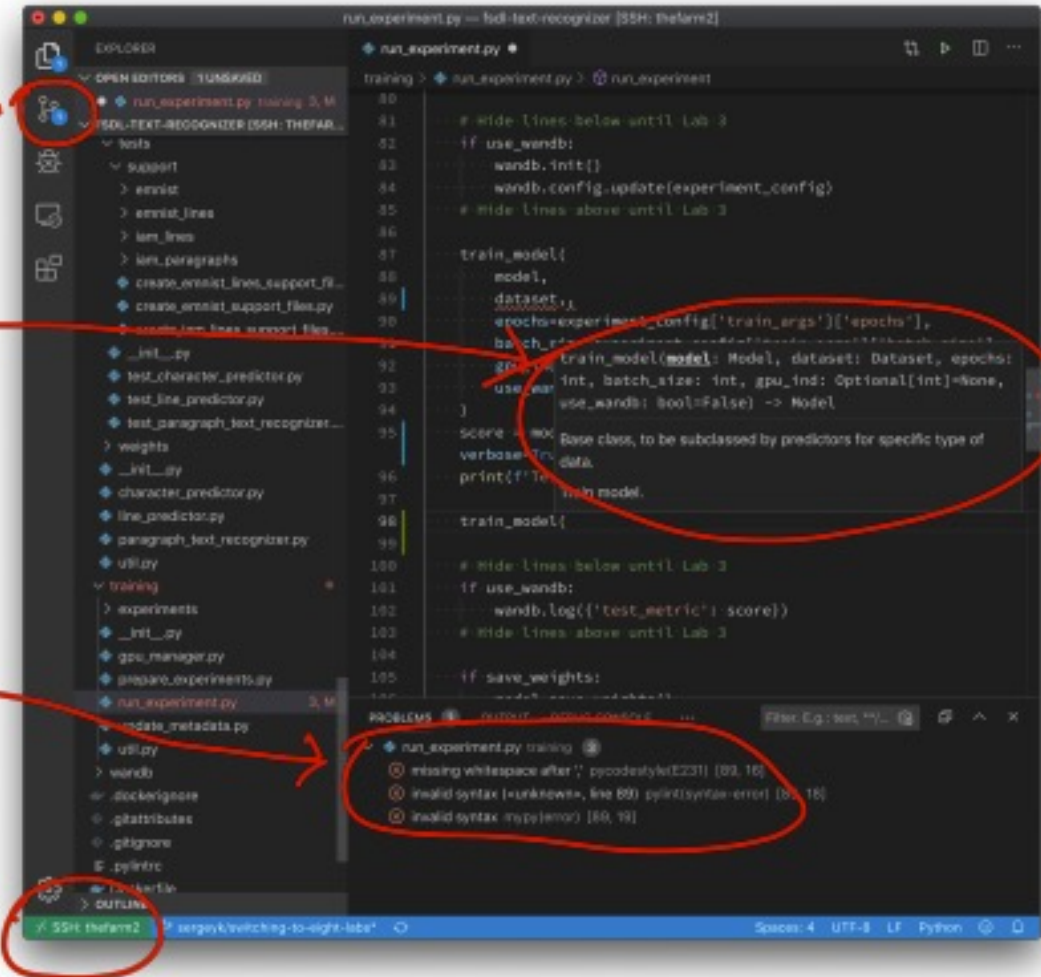
Built-in git staging and diffing →

Peek documentation

- VS Code makes for a very nice Python experience

Lint code as you write

Open whole projects remotely →



Visual Studio Code

Built-in git staging and diffing

Open projects remotely

Notebook port forwarding

Use the terminal

```
line_lstm.py -- fadl-text-recognition-2021 (SSH: bigboy.local)
1 from typing import Any, Dict
2 import argparse
3 import torch
4 import torch.nn.functional as F
5 import torch.nn as nn
6
7 from line_lstm import LineCNN
8 from line_lstm2 import LineCNN2
9
10 LINE_CNN_TYPE = 1
11 LSTM_DIM = 512
12 LSTM_LAYERS = 2
13 LSTM_DROPOUT = 0.2
14
15 class LineCNN(LSTM(nn.Module)):
16     """Process the line through a CNN and process the resulting sequence through LSTM layers."""
17
18     def __init__(self,
19                 data_config: Dict[str, Any],
20                 args: argparse.Namespace = None,
21                 device: torch.device = None):
22         super().__init__()
23         self.data_config = data_config
24         self.args = args
25         self.device = device
26         self.max_output_length = data_config["output_dim"][0]
27
28         num_classes = len(data_config["mapping"])
29         line_cnn_type = self.args.get("line_cnn_type", LINE_CNN_TYPE)
30         lstm_dim = self.args.get("lstm_dim", LSTM_DIM)
31         lstm_layers = self.args.get("lstm_layers", LSTM_LAYERS)
32         lstm_dropout = self.args.get("lstm_dropout", LSTM_DROPOUT)
33
34         # LineCNN outputs (B, C, S) log-probs, with C == num_classes
35         if line_cnn_type == 1:
36             self.line_cnn = LineCNN(data_config, data_config, args=args)
37         elif line_cnn_type == 2:
38             self.line_cnn = LineCNN2(data_config, data_config, args=args)
39
40         self.lstm = nn.LSTM(
41             input_size=num_classes,
42             hidden_size=lstm_dim,
43             num_layers=lstm_layers,
44             dropout=lstm_dropout,
45             bidirectional=False,
46         )
47
48         self.fc = nn.Linear(lstm_dim, num_classes)
49
50 def forward(self, x: torch.Tensor) -> torch.Tensor:
51     """Forward pass"""
52     # Process the line through a CNN and process the resulting sequence through LSTM layers.
53     # ...
54     # ...
55     # ...
```

```
08 model.transformer.decoder.layers.3.norm2 LayerNorm 512
09 model.transformer.decoder.layers.5.norm2 LayerNorm 512
10 model.transformer.decoder.layers.3.dropout1 Dropout 0
11 model.transformer.decoder.layers.3.dropout2 Dropout 0
12 model.transformer.decoder.layers.5.dropout1 Dropout 0
13 model.transformer.decoder.layers.5.dropout2 Dropout 0
14 train_acc Accuracy 0
15 val_acc Accuracy 0
16 test_acc Accuracy 0
17 loss_fn CrossEntropyLoss 0
18 val_err CharacterErrorRate 0
19 test_err CharacterErrorRate 0
20
21 =====
22 33.5 # Trainable params
23 0 Non-trainable params
24 33.5 # Total params
25 Epoch 31: 100% | 1/100 [00:00<00:00, 3.18s/it, loss=1.04, v_num=500, val_loss=1.08, val_err=0.98] wandb:
26 Network error (ReadTimeout), entering retry loop. See wandb/debug-internal.log for full traceback.
27 Epoch 32: 100% | 1/100 [00:01<00:00, 1.88s/it, loss=1.55, v_num=500, val_loss=1.68, val_err=0.96] wandb:
28 Network error resolved after 0:01:48.122480, resuming normal operation.
29 Epoch 36: 100% | 1/100 [01:27<00:00, 1.82s/it, loss=1.54, v_num=500, val_loss=1.67, val_err=0.96]
30 Testing: 100% | 1/100 [00:18<00:00, 3.17s/it]
31
32 =====
33 DETAILED TEST RESULTS
34 {"test_acc": tensor(0.9428, device="cuda:0"),
35 "test_err": tensor(0.9609, device="cuda:0")}
36
37 wandb: Waiting for W&B process to finish, PID 17046
38 wandb: Program ended successfully.
39 wandb: Find user logs for this run at: /home/sergeyk/work/fadl/fadl-text-recognition-2021/wandb/run-20210222_220519-
40 wandb: Find internal logs for this run at: /home/sergeyk/work/fadl/fadl-text-recognition-2021/wandb/run-20210222_220
41 wandb: Run summary:
42 wandb: lr-Adam 5e-05
43 wandb: train_loss 1.58957
44 wandb: epoch 36
45 wandb: _step 5327
46 wandb: _runtime 3284
47 wandb: _timestamp 161401601
48 wandb: val_loss 1.67319
49 wandb: val_acc 0.94475
50 wandb: val_err 0.96612
51 wandb: test_acc 0.94196
52 wandb: test_err 0.96892
53
54 wandb: Run history:
55 wandb: lr-Adam
56 wandb: train_loss
57 wandb: epoch
58 wandb: _step
59 wandb: _runtime
60 wandb: _timestamp
61 wandb: val_loss
62 wandb: val_acc
63 wandb: val_err
64 wandb: test_acc
65 wandb: test_err
66
67 Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
68 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
69 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
70 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
71 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
72 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
73 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
74 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
75 wandb: Spaced 5 W&B file(s), 1 media file(s), 6 artifact file(s) and 0 other file(s)
```


Linters and Type Hints

- Whatever code styles rules can be codified, should be
- Static analysis can catch some bugs
- Static type checking both documents code and catches bugs

```
87 train_model(  
88     model,  
89     dataset,  
90     epochs=experiment_config['train_args']['epochs'],  
91     batch_size=batch_size,   
92     gpu_ind=gpu_ind,   
93     use_wandb=use_wandb,   
94 )  
95 score = model.train(dataset, epochs=epochs, batch_size=batch_size, gpu_ind=gpu_ind, use_wandb=use_wandb)  
96 print(f'Test score: {score}')  
97   
98 train_model(  
99
```

train_model(model: Model, dataset: Dataset, epochs: int, batch_size: int, gpu_ind: Optional[int]=None, use_wandb: bool=False) -> Model

Base class, to be subclassed by predictors for specific type of data.

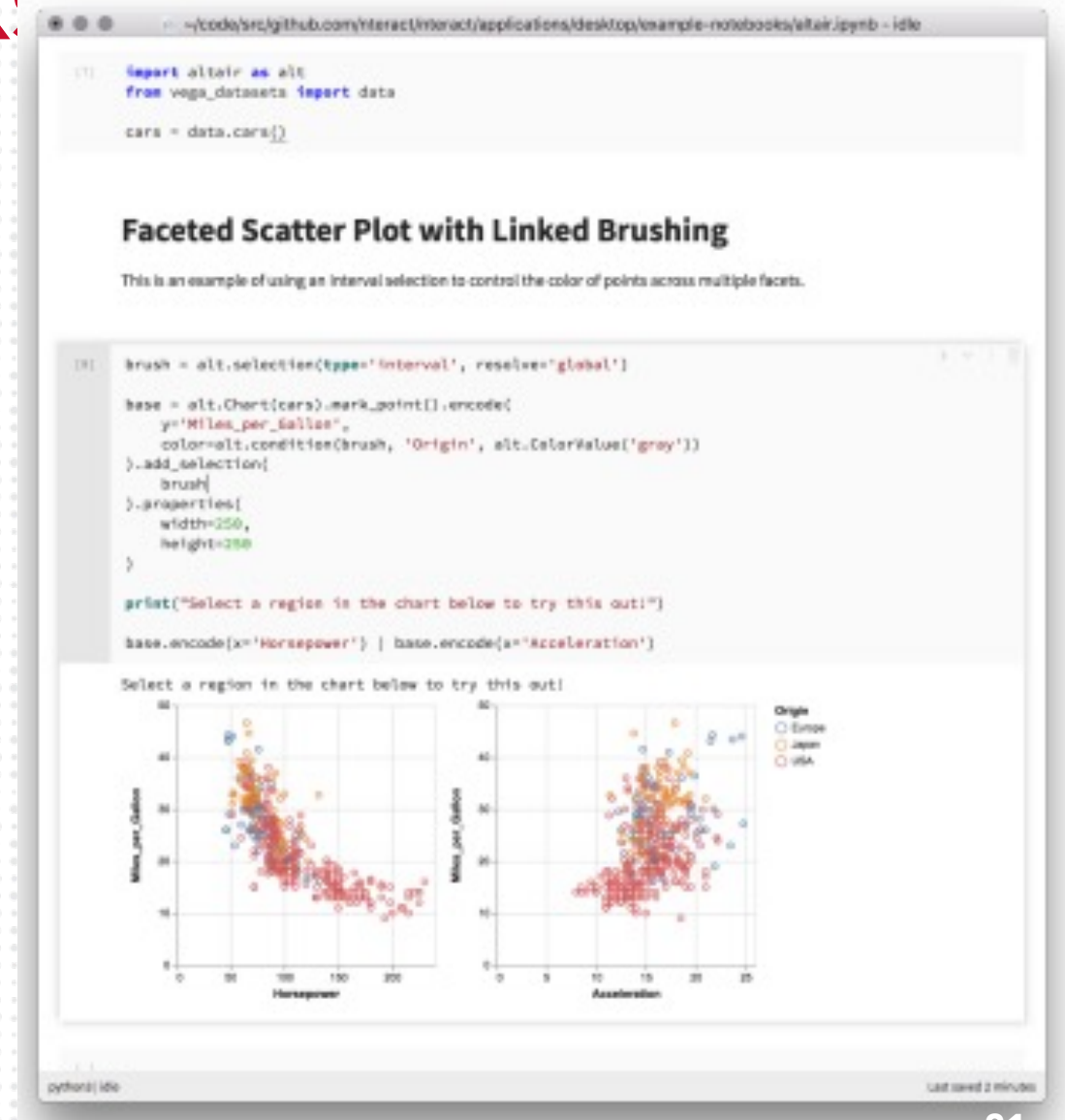
Train model.

run_experiment.py training 3

- ⊗ missing whitespace after ',' pycodestyle(E231) [89, 16]
- ⊗ invalid syntax (<unknown>, line 89) pylint(syntax-error) [89, 18]
- ⊗ invalid syntax mypy(error) [89, 19]

Jupyter Notebooks

- Notebooks have become fundamental to data science
- Great as the “first draft” of a project
- Jeremy Howard from fast.ai good to learn from (course.fast.ai videos)
- Difficult to make scalable, reproducible, well-tested



Problem with notebooks

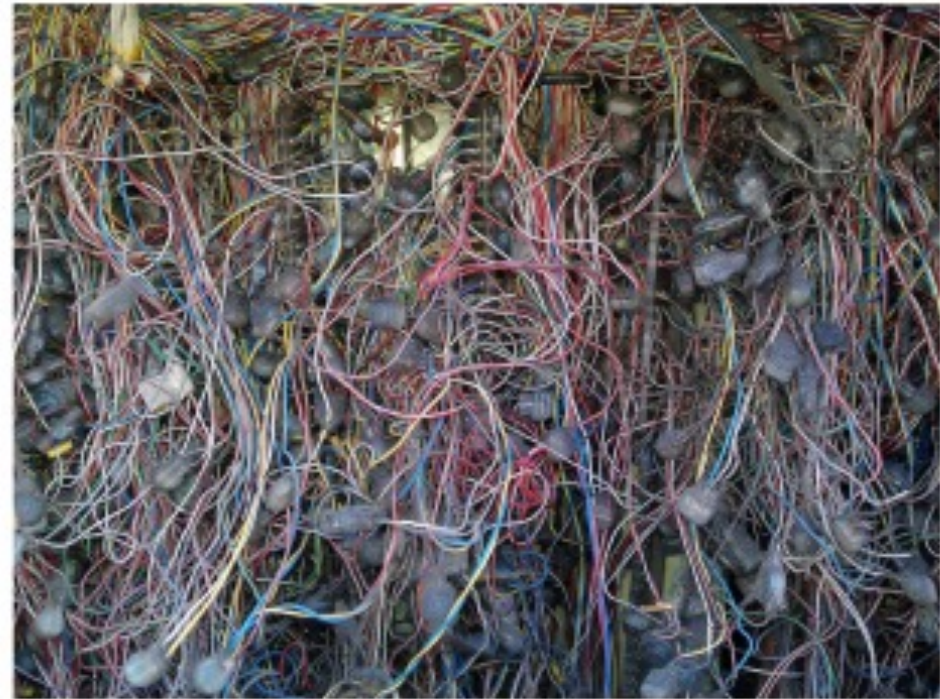
- Hard to version
- Notebook “IDE” is primitive
- Very hard to test
- Out-of-order execution artifacts
- Hard to run long or distributed tasks

5 reasons why jupyter notebooks suck



Alexander Mueller [Follow](#)

Mar 24, 2018 · 3 min read ★

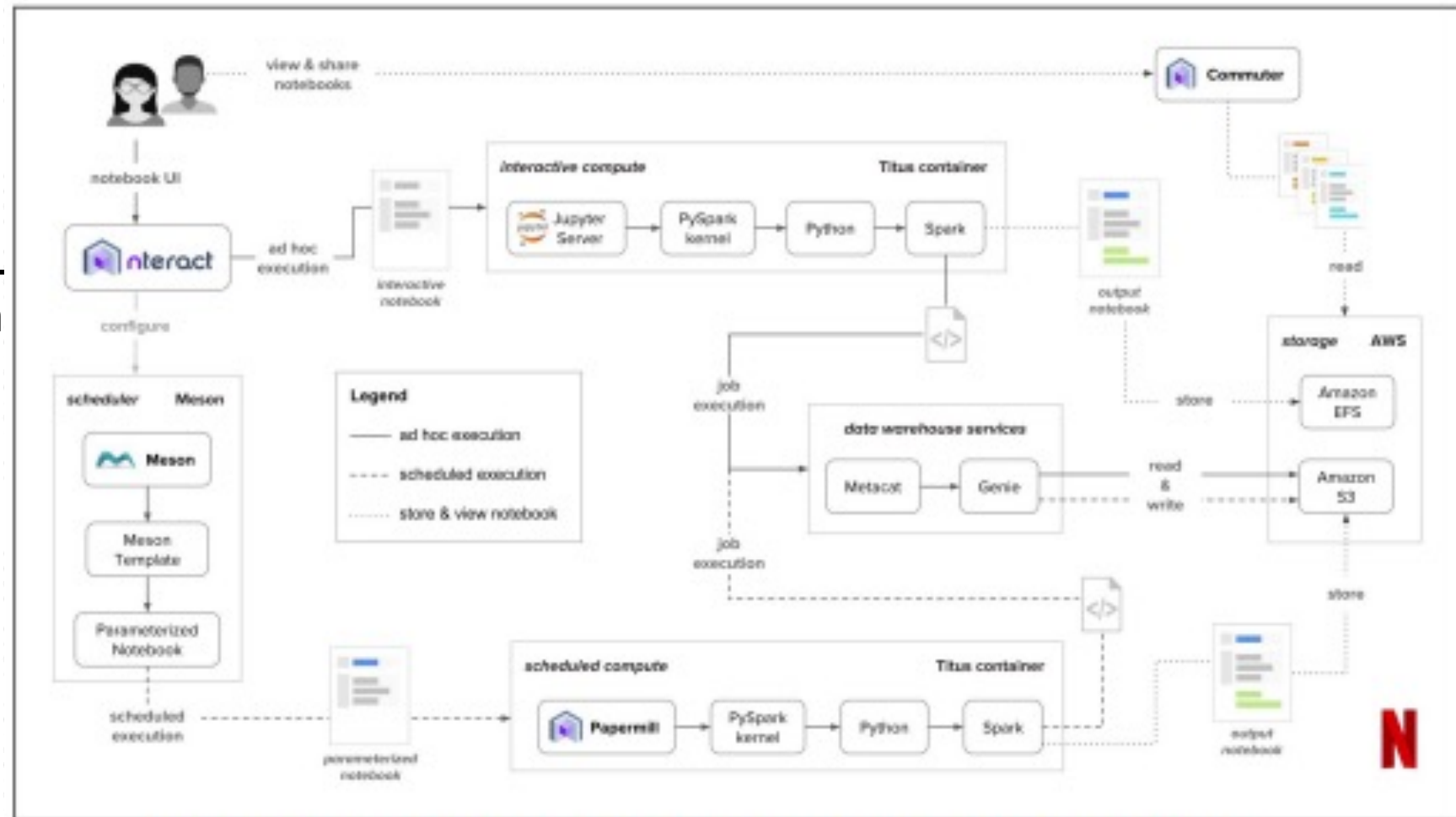


How it feels like managing jupyter notebooks (Complexity ©

<https://www.flickr.com/photos/bitterjug/7670055210/>)

Jupyter Notebooks

- Counter-points:
- Netflix bases all ML workflows on them



NBDev

- Counter-points:
- Jeremy Howard from fast.ai uses them for everything, with nbdev

Card

API details.

```
#hide
from nbdev.showdoc import *

#export
from __future__ import print_function, division
import random

class Card:
    """Represents a standard playing card.

    Attributes:
        suit: integer 0-3
        rank: integer 1-13
    """
    suit_names = ["Clubs", "Diamonds", "Hearts", "Spades"]
    rank_names = [None, "Ace", "2", "3", "4", "5", "6", "7", "8", "9", "10", "Jack", "Queen", "King"]

    def __init__(self, suit=0, rank=2):
        self.suit, self.rank = suit, rank

    def __str__(self):
        """Returns a human-readable string representation."""
        return '%s of %s' % (Card.rank_names[self.rank], Card.suit_names[self.suit])

    def __eq__(self, other) -> bool:
        """Checks whether self and other have the same rank and suit."""
        return self.suit == other.suit and self.rank == other.rank

    def __lt__(self, other) -> bool:
        """Compares this card to other, first by suit, then rank."""
        t1 = self.suit, self.rank
        t2 = other.suit, other.rank
        return t1 < t2

    def __repr__(self): return self.__str__()

    def foo(): pass
```

Card is a class that represents a single card in a deck of cards. For example:

```
Card(suit=2, rank=11)
```

Jack of Hearts

```
c = Card(suit=1, rank=3)
assert str(c) == '3 of Diamonds'

c2 = Card(suit=2, rank=11)
assert str(c2) == 'Jack of Hearts'
```

You can do comparisons of cards, too!

```
assert c2 > c
```


Streamlit

- New, but great at fulfilling a common ML need: interactive applets
- Decorate normal Python code
- Smart data caching, quick re-rendering
- In the works: sharing as easy as pushing a web app to Heroku



Setting up environment

Conda + Pip-Tools Sample Project

Quick demo of setting up a deep learning Python environment.

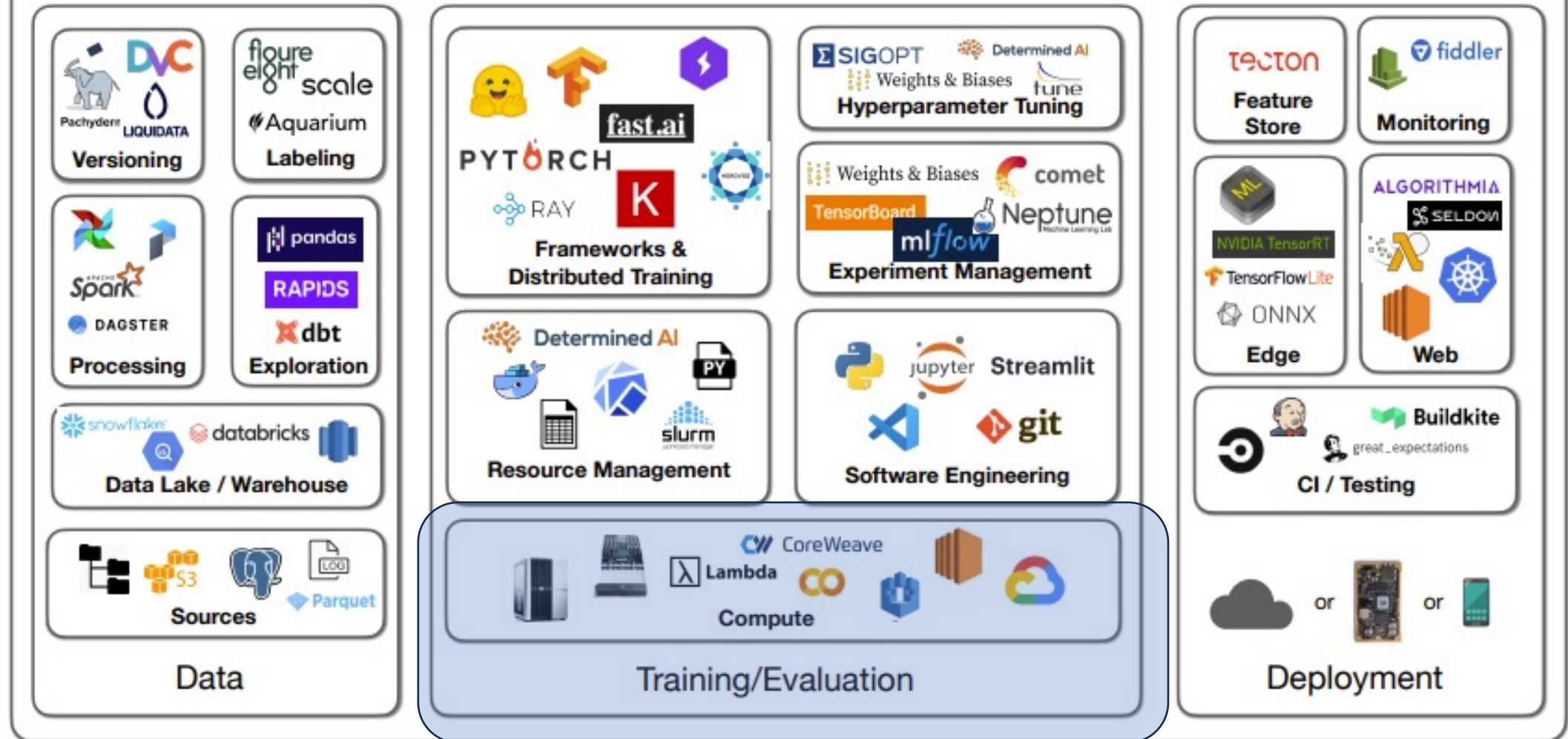
Our goals:

- Easily specify the exact Python, CUDA, CUDNN environment
- Humans should specify minimal constraints (`torch >= 1.7` and `numpy`), computer should figure out exact, mutually compatible versions (`torch==1.7.1; numpy==1.19.5`)
- Separate production (`torch`) from development (`black`) dependencies

We achieve this by:

- We specify our Python and CUDA versions in `environment.yml`
- We use the `conda` package manager to create our environment from this file
- We specify our requirements in `requirements/prod.in` and `requirements/dev.in`
- We use `pip-tools` to lock in mutually compatible versions of all requirements
- We add a `Makefile` so we can simply run `make` to update everything

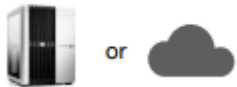
"All-in-one"



Compute needs

Development

- **Function**
 - Writing code
 - Debugging models
 - Looking at results
- **Desiderata**
 - Quickly compile models and run training
 - Nice-to-have: use GUI
- **Solutions**
 - Desktop with 1-4 GPUs
 - Cloud instance with 1-4 GPUs

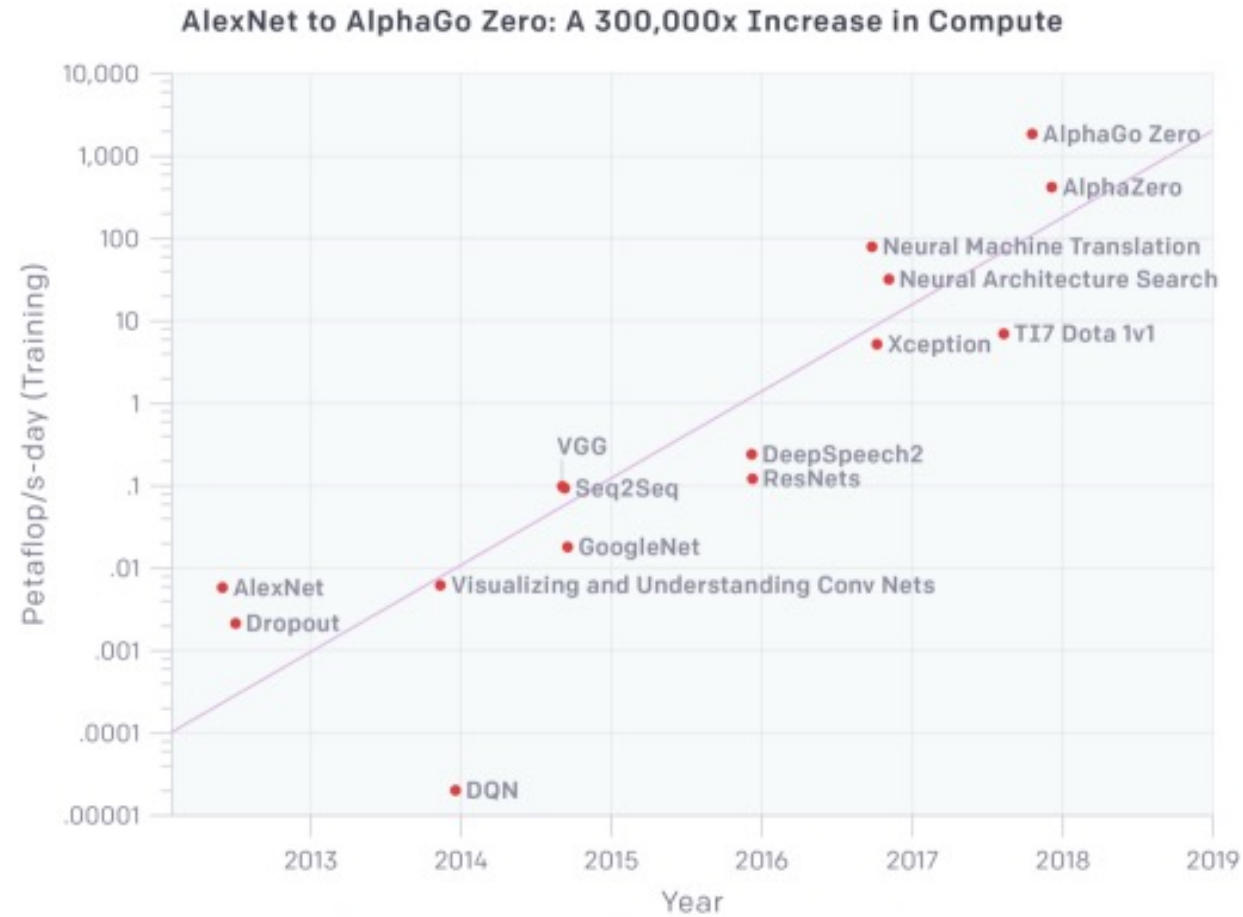


Training/Evaluation

- **Function**
 - Model architecture/hyperparam search
 - Training large models
- **Desiderata**
 - Easy to launch experiments and review results
- **Solutions**
 - Desktop with 4 GPUs
 - Private cluster of GPU machines
 - Cloud cluster of GPU machines



Why compute matters

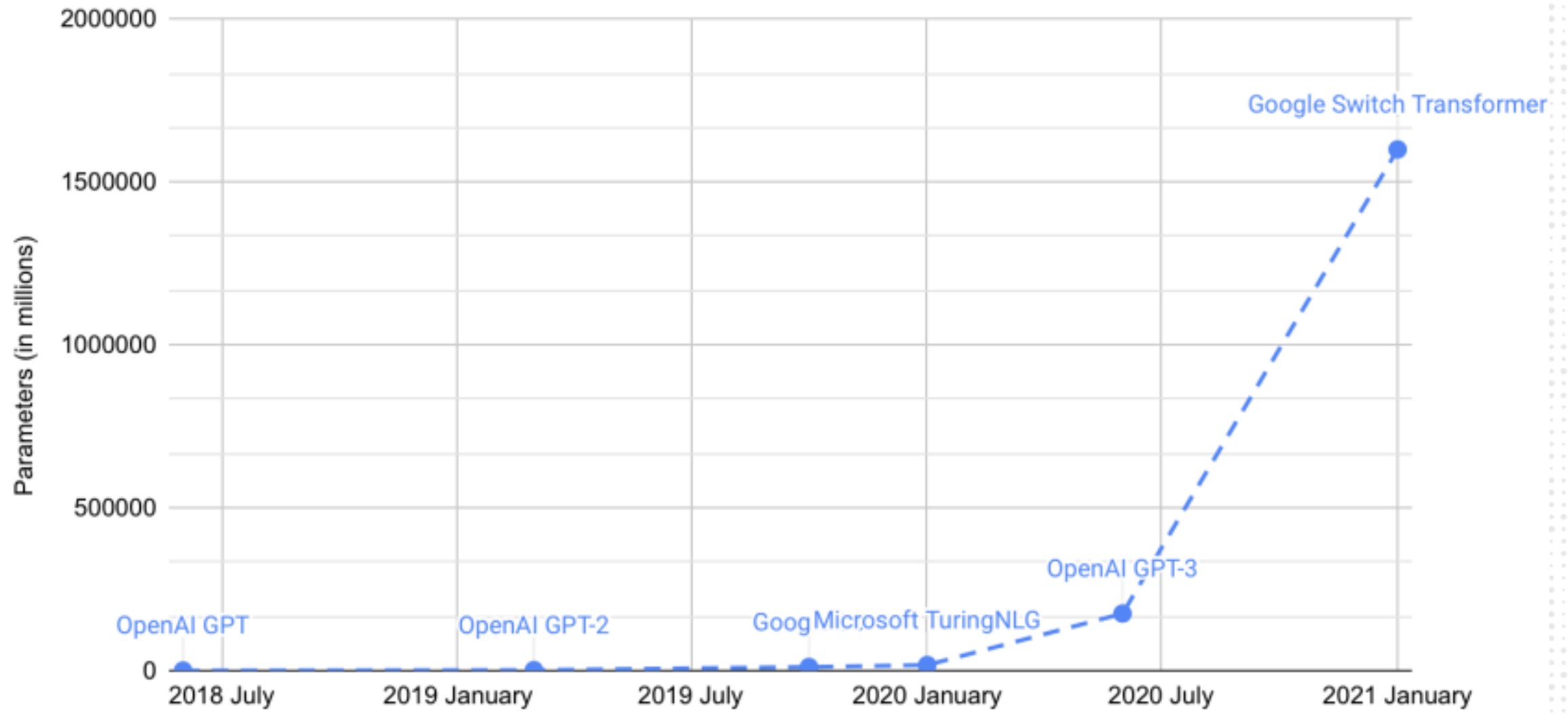


<https://openai.com/blog/ai-and-compute/>

LOG SCALE

LINEAR SCALE

Transformer Models



So,



or



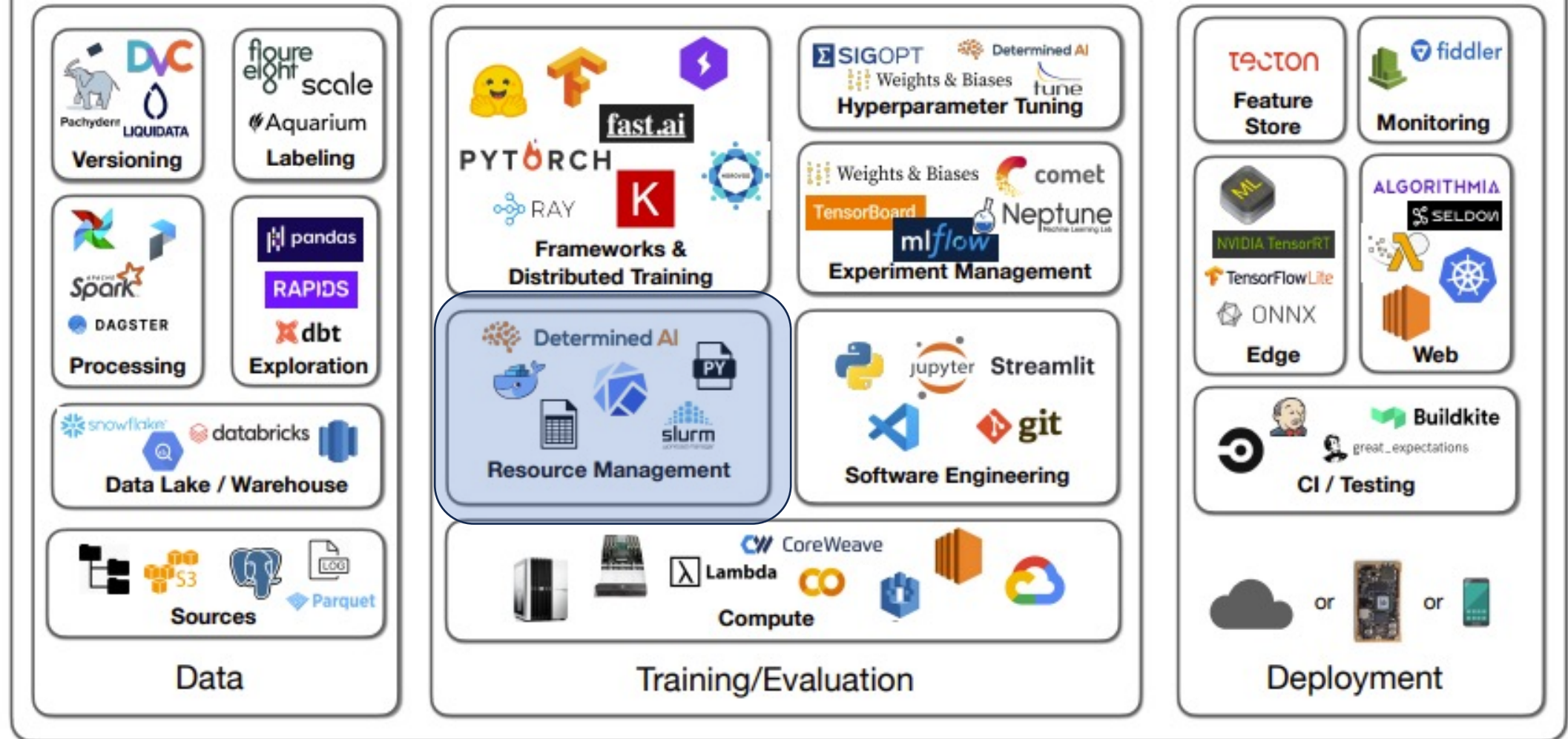
?

- GPU Basics
- Cloud Options
- On-prem Options
- Analysis and Recommendations

In Practice

- Even though cloud is expensive, it's hard to make on-prem scale past a certain point
- Dev-ops (declarative infra, repeatable processes) definitely easier in the cloud
- Maintenance is also a big factor

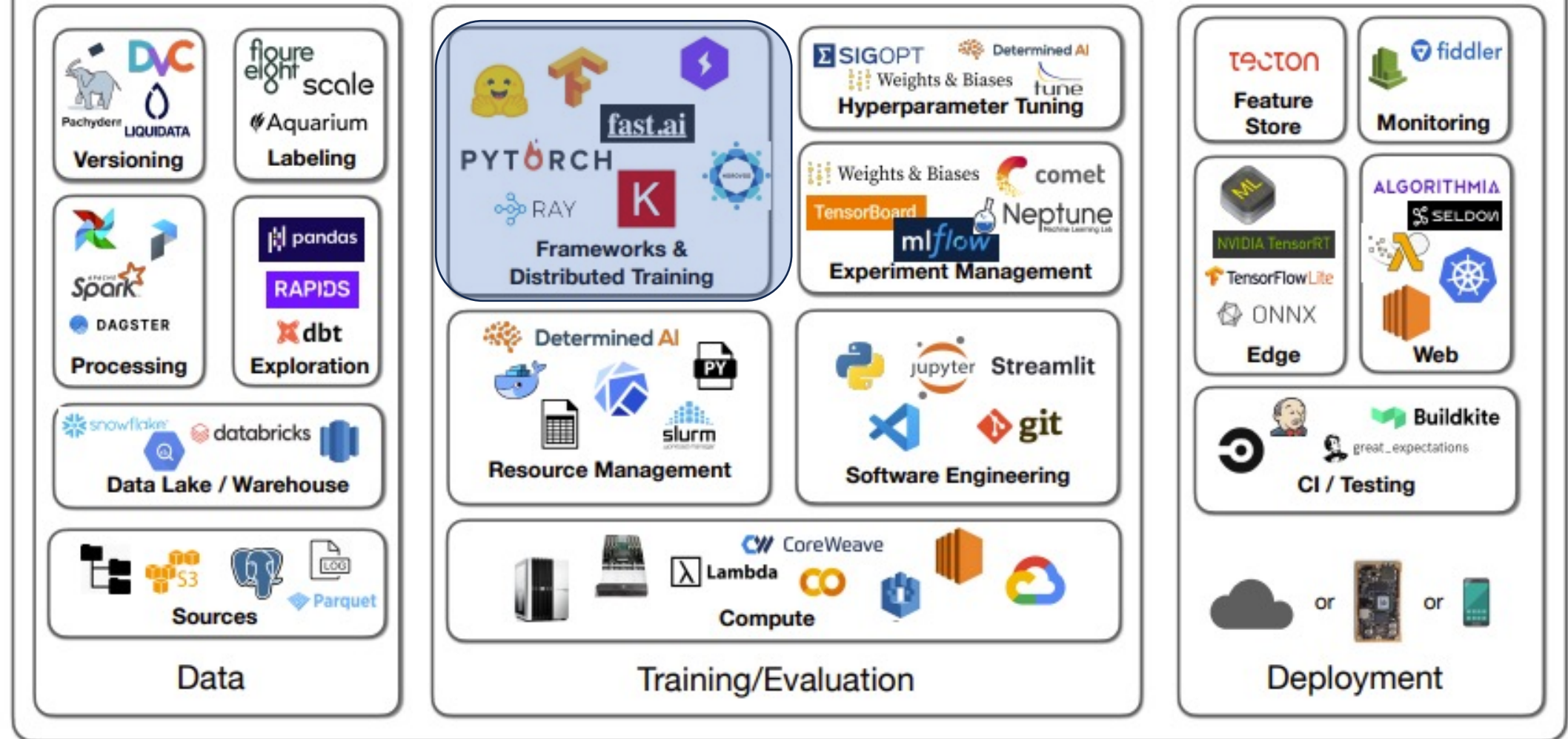
"All-in-one"



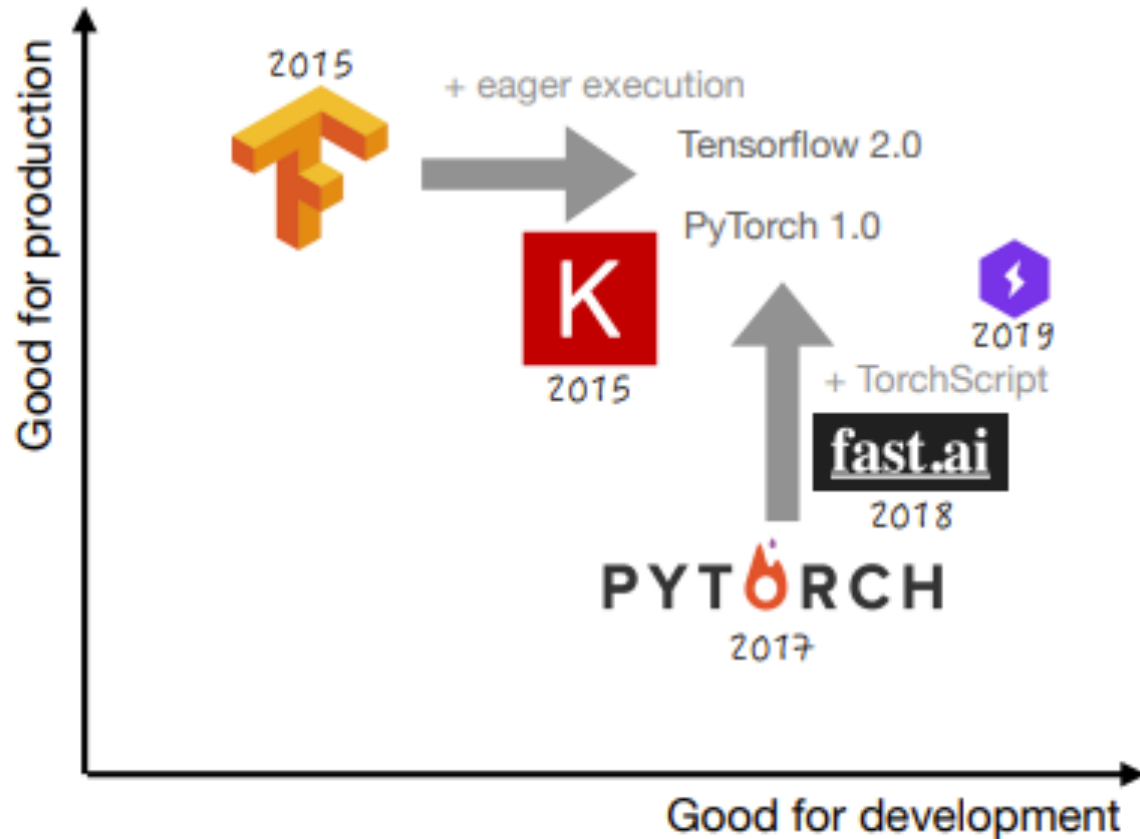
Resource Management

- Function
 - Multiple people...
 - using multiple GPUs/machines...
 - running different environments
- Goal
 - Easy to launch a batch of experiments with proper dependencies and resource allocations
- Solutions
 - Python scripts
 - SLURM
 - Docker + Kubernetes
 - Software specialized for ML use cases

"All-in-one"



Deep Learning Frameworks

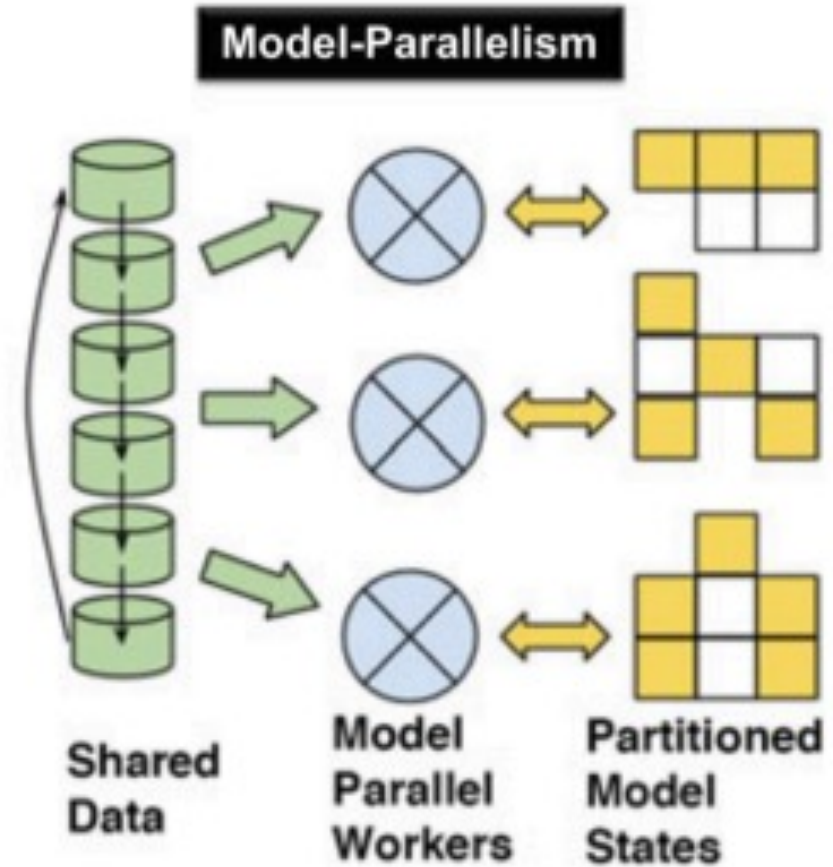
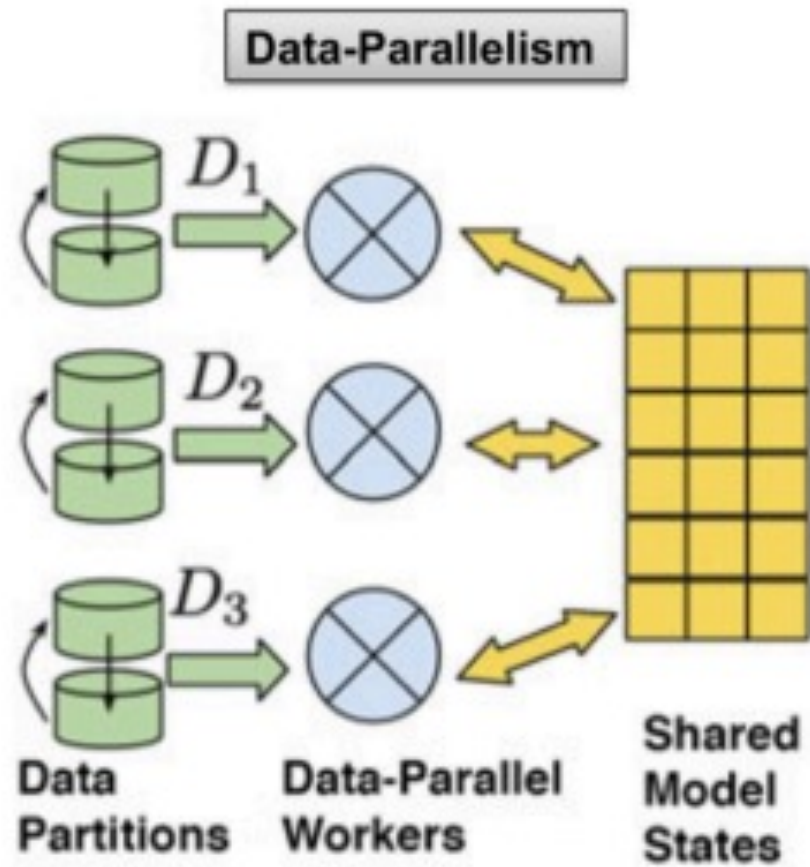


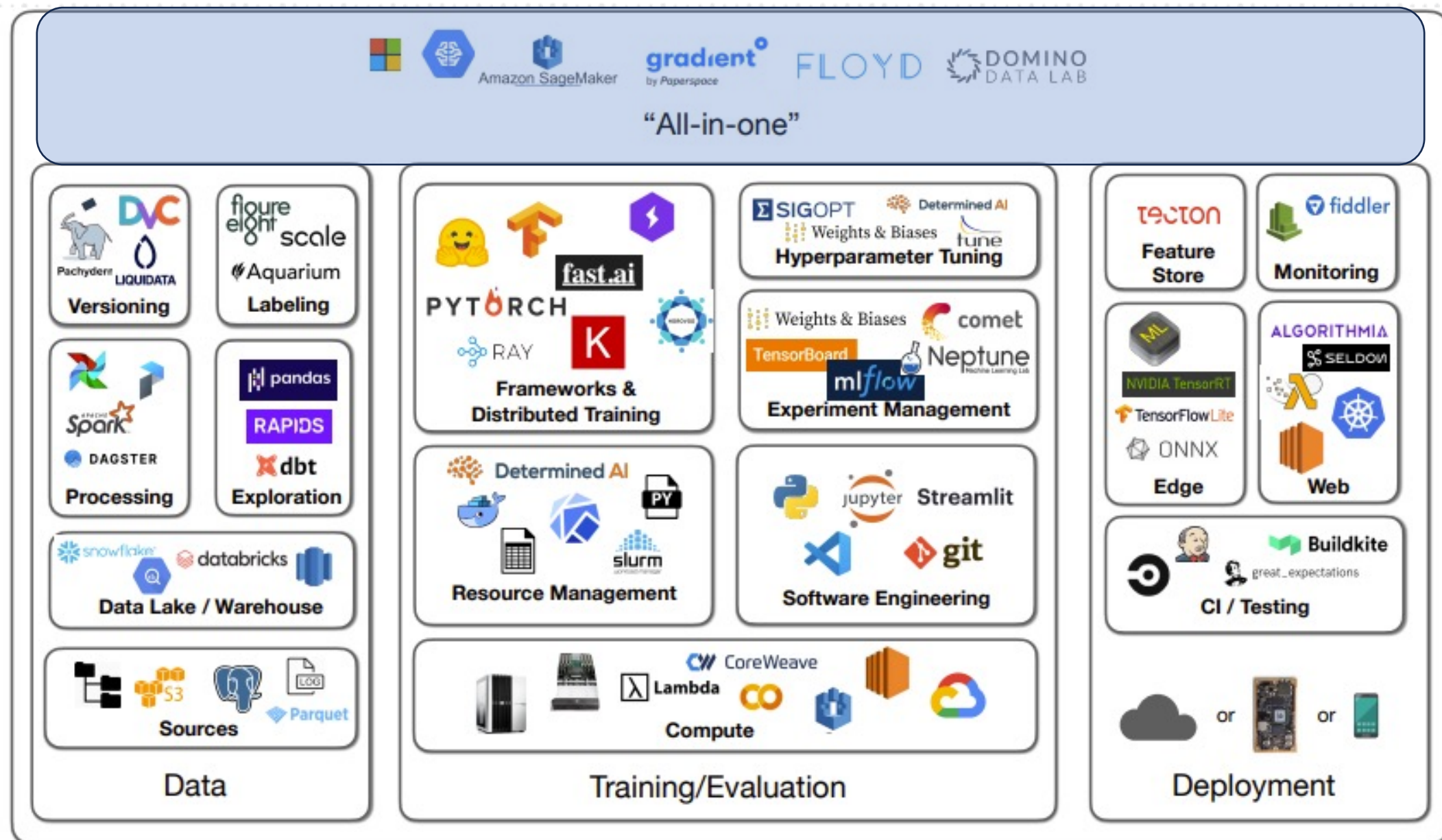
- Unless you have a good reason not to, use Tensorflow/Keras or PyTorch
- Both have converged to the same point:
 - easy development via define-by-run
 - multi-platform optimized execution graph
- Today, most new projects use PyTorch, because of its more Python dev-friendly experience
- fast.ai library builds on PyTorch with best practices
- PyTorch-Lightning adds a powerful training loop

Distributed Training

- Using multiple GPUs and/or machines to train a single model.
- More complex than simply running different experiments on different GPUs
- A must-do on big datasets and large models

Parallelism





All-in-one Solutions

- Single system for everything
 - development (hosted notebook)
 - scaling experiments to many machines (sometimes even provisioning)
 - tracking experiments and versioning models
 - deploying models
 - monitoring performance