

# OSS Stack for ML

Soma S. Dhavala

01-01-2025



[@TheSaddlePoint](#)



[github.com/mlsquare](https://github.com/mlsquare)



[linkedin.com/in/somasdhavala](https://linkedin.com/in/somasdhavala)



[@TheSaddlePoint](#)

## Agenda

What should an OSS Stack to build ML look like?

Part-1: OSS for ML Engineering

Part-2: OSS for ML Science

Part-3: OSS model for ML Development

**ML team**  
highest accuracy



**Sales**  
sells more ads



**Product**  
fastest inference



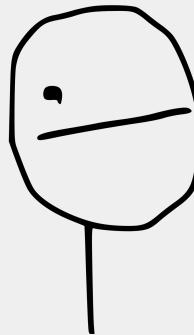
**Manager**  
maximizes profit  
= laying off ML teams



Expectation

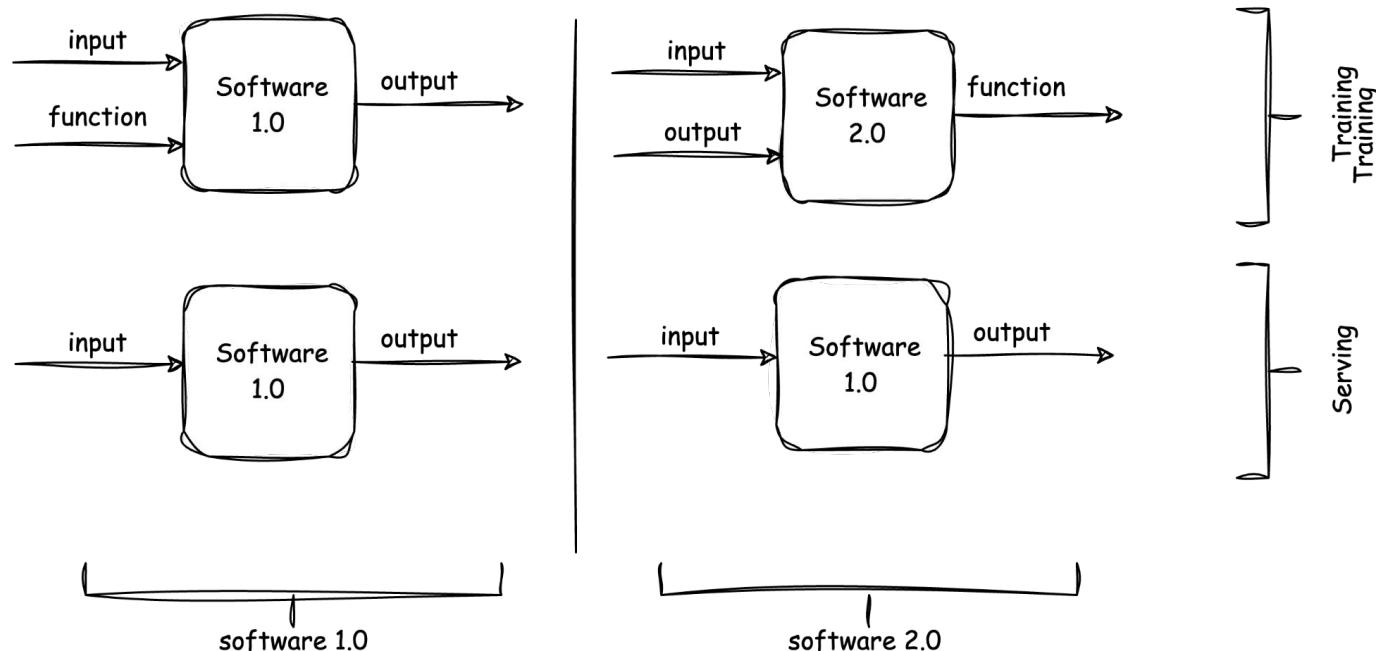


Reality



But why?

## Software 1.0 vs Software 2.0



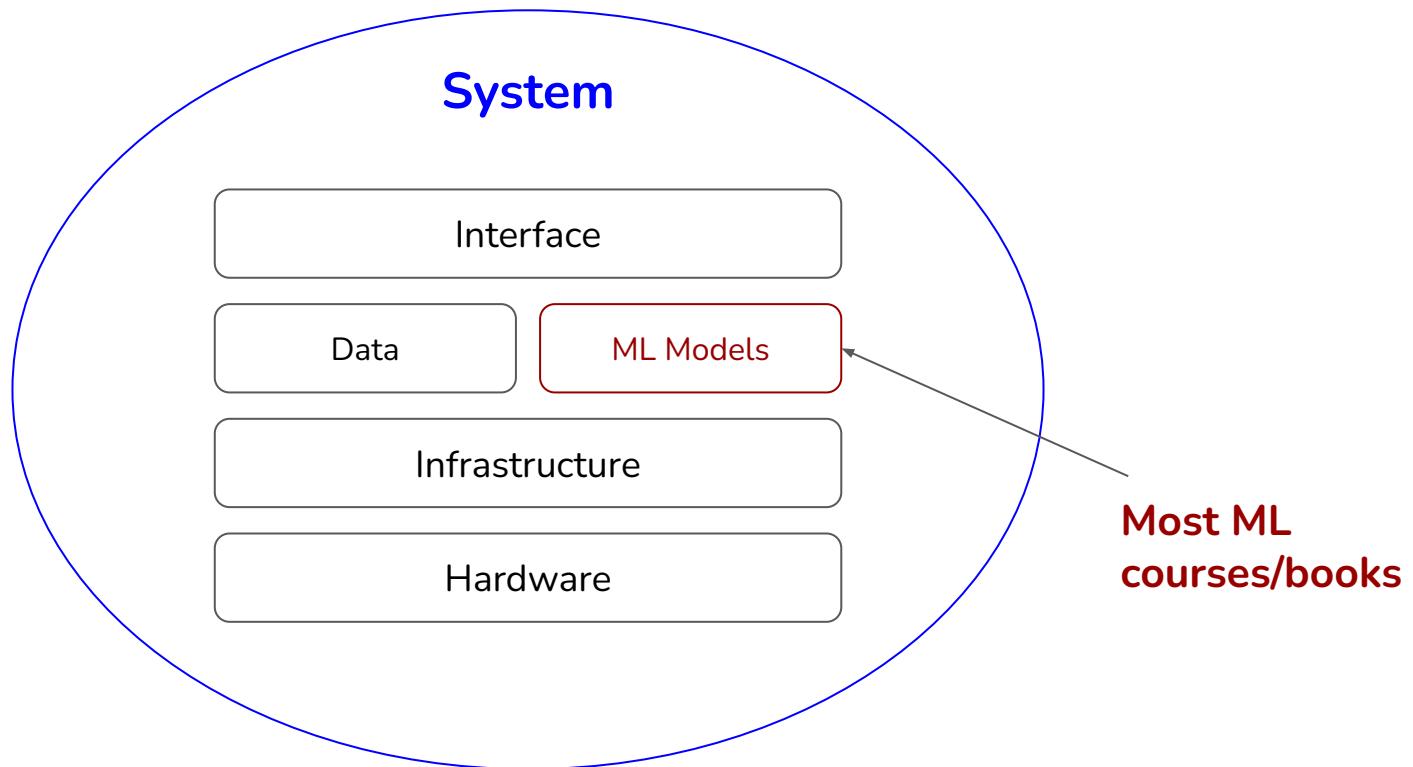
ML as a software is a *fundamentally different beast*.

But why?

	Software 1.0	Software 2.0
Codified in	Formal Language	Weights & Biases (parameters)
Developed by	Programming	Training
Specification	PRD/SRD	Data
Behaviour	Deterministic	Stochastic
	Provably correct	Provably wrong
	Debuggable	Hard to Debug
	Verifiable	Hard to verify
	Explainable	Hard to explain
	Fixable	Hard to fix
	Idempotent	Hard to reproduce

Determinism and Control (of the build process and products) - we take them for granted. But no longer.

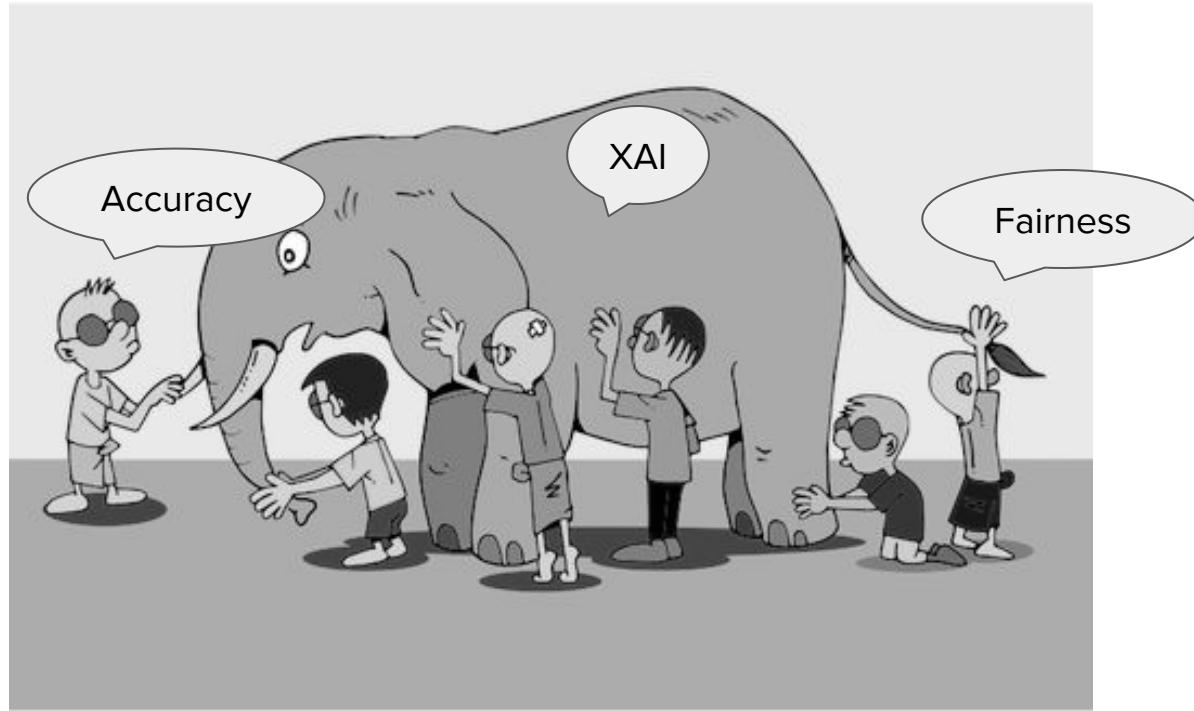
But why?



Think systems, not models. Think modeling, not models.  
Combating this intellectual inertia is hard.

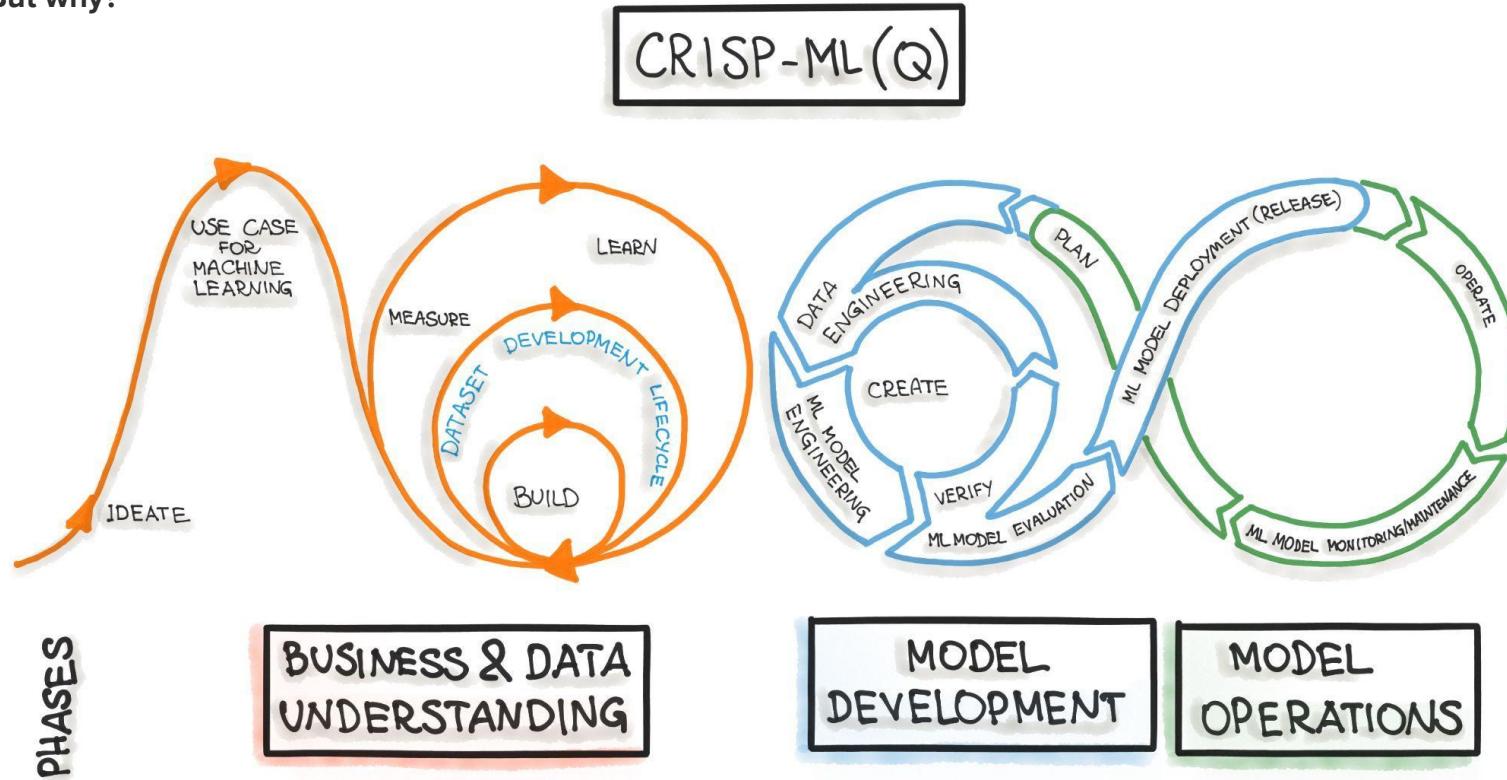
source: [CS329S @ Stanford](#)

But why?



Semi-orthogonal sub-fields.

But why?



@visenger

anything changes > everything changes  
modeling is highly iterative and nonlinear

source: [ml-ops.org](http://ml-ops.org)

Non-deterministic systems

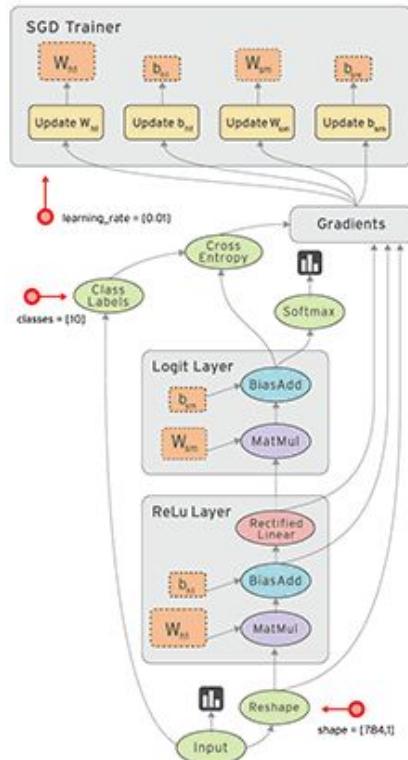
Build must be agile and address ambiguity

Design must address “reliability” in the presence of “uncertainty”

IF  
 $CS = DS + A$   
THEN  
 $ML = ?$

need an abstraction for Models and Modeling

“ML modeling” is a DAG

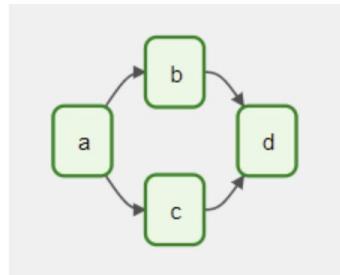


**TensorFlow:** Data Flow Model to writing programs

# DAGs

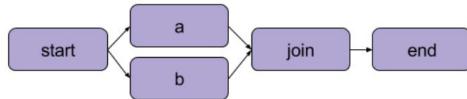
A DAG (Directed Acyclic Graph) is the core concept of Airflow, collecting **Tasks** together, organized with dependencies and relationships to say how they should run.

Here's a basic example DAG:



## Branch

You can express parallel steps with a **branch**. In the figure below, `start` transitions to two parallel steps, `a` and `b`. Any number of parallel steps are allowed. A benefit of a branch like this is performance: Metaflow can execute `a` and `b` over multiple CPU cores or over multiple instances in the cloud.



```
from metaflow import FlowSpec, step

class BranchFlow(FlowSpec):

    @step
    def start(self):
        self.next(self.a, self.b)

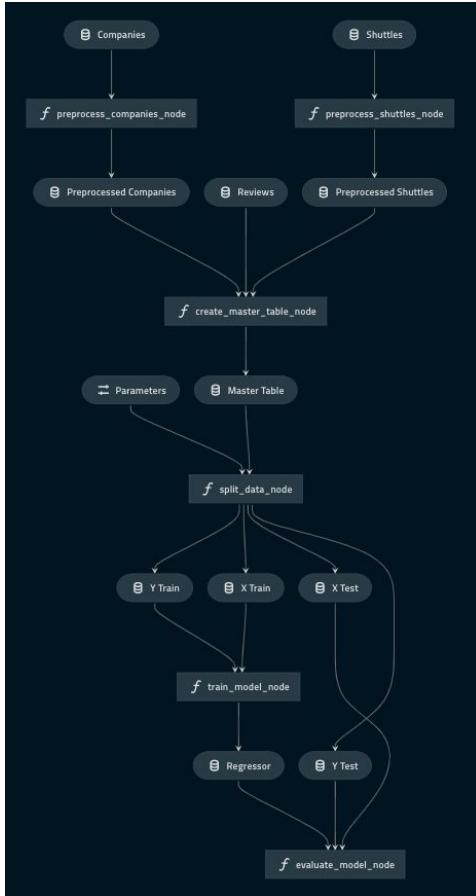
    @step
    def a(self):
        self.x = 1
        self.next(self.join)

    @step
    def b(self):
        self.x = 2
        self.next(self.join)

    @step
    def join(self, inputs):
        print('a is %s' % inputs.a.x)
        print('b is %s' % inputs.b.x)
        print('total is %d' % sum(input.x for input in inputs))
        self.next(self.end)

    @step
    def end(self):
        pass

if __name__ == '__main__':
    BranchFlow()
```



Kedro is a toolbox for production-ready data science. It uses software engineering best practices to help you create data engineering and data science pipelines that are reproducible, maintainable, and modular. You can find out more at [kedro.org](https://kedro.org).

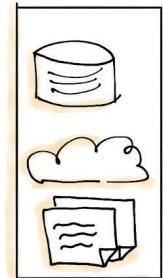
Kedro is an open-source Python framework hosted by the [LF AI & Data Foundation](https://lfai.foundation/kedro).

# OSS for ML Engineering

But why?

## DATA PIPELINE

# MACHINE LEARNING ENGINEERING.



### EXPLORATION & VALIDATION

### WRANGLING (CLEANING)

## DATA

FEEDBACK  
new data from model performance

### MONITORING & LOGGING

- Model decay trigger

## MACHINE LEARNING PIPELINE

INNOQ

## SOFTWARE CODE PIPELINE

(Lot of) Hidden Tech Debt

- Profiling
- "JUnit4Data"

- Data versioning

### TRAIN

### TEST

### MODEL ENGINEERING

- Feature engineering
- Hyperparameters tuning

- Best model selection
- Model performance metrics
  - accuracy
  - precision
  - recall

- F1

### MODEL EVALUATION

code

params

code

param

code

param

- Model format
  - OMNX
  - SAV
  - PKL

### MODEL PACKAGING

- Model versioning

## MODEL

- Model serving
  - service
  - Docker
  - K8S

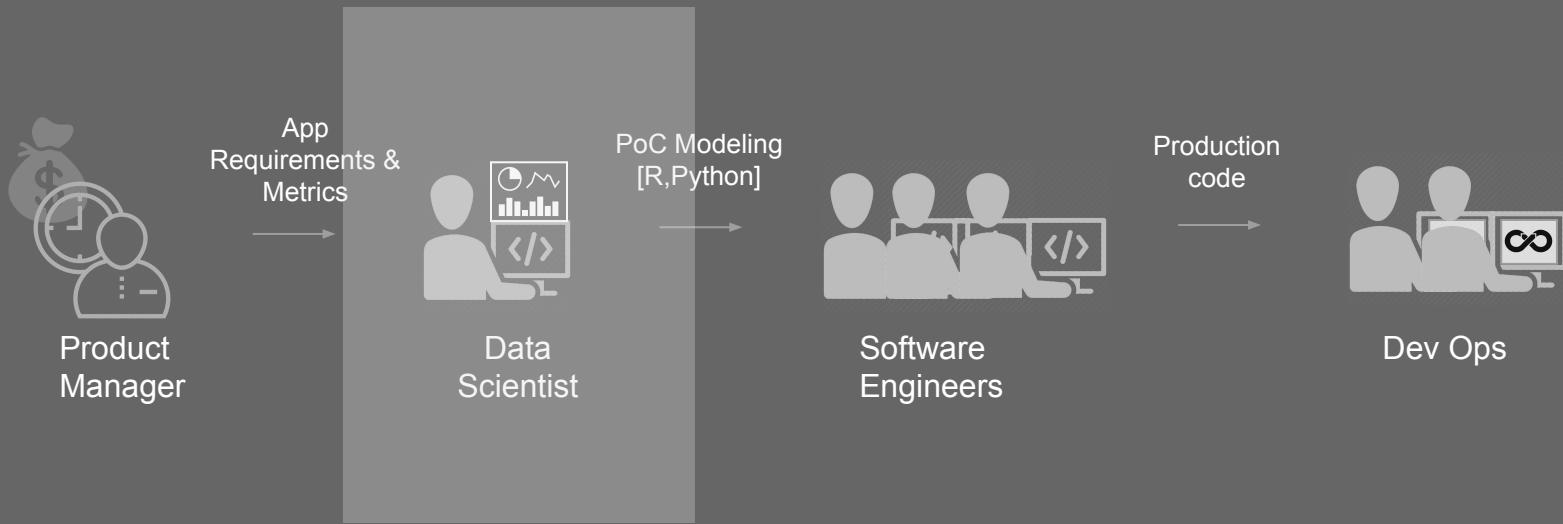
## CODE

- Trunk based dev.
- Code versioning

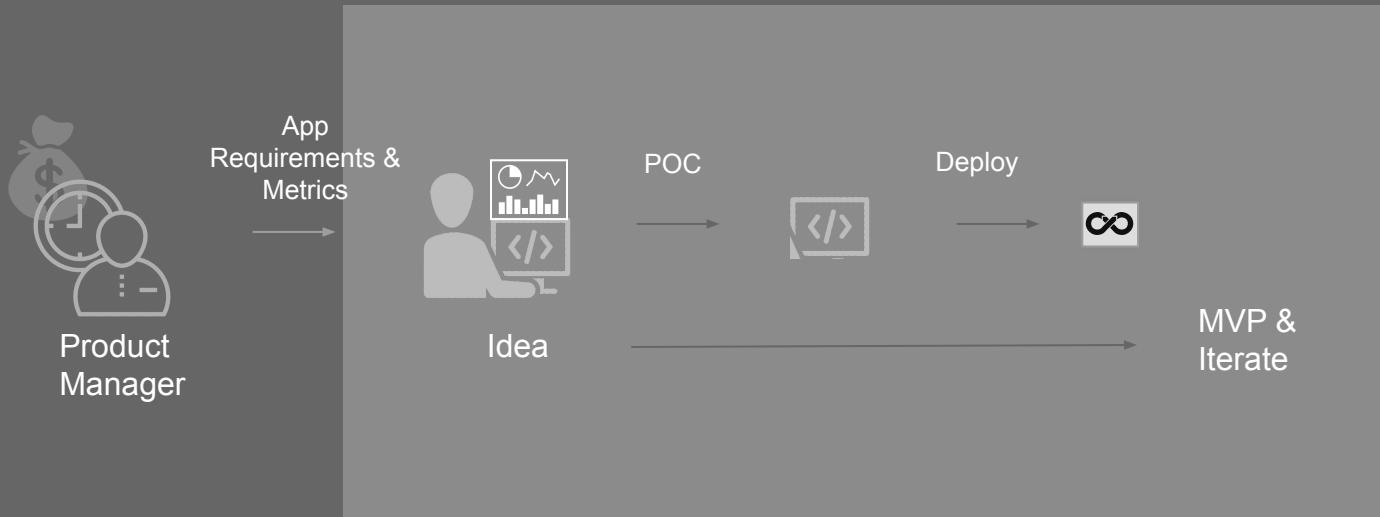
### BUILD & INTEGRATION TESTING

### DEPLOYMENT DEV ↓ PRODUCTION

1. Data Engineering
2. Model Engineering
3. Model Deployment



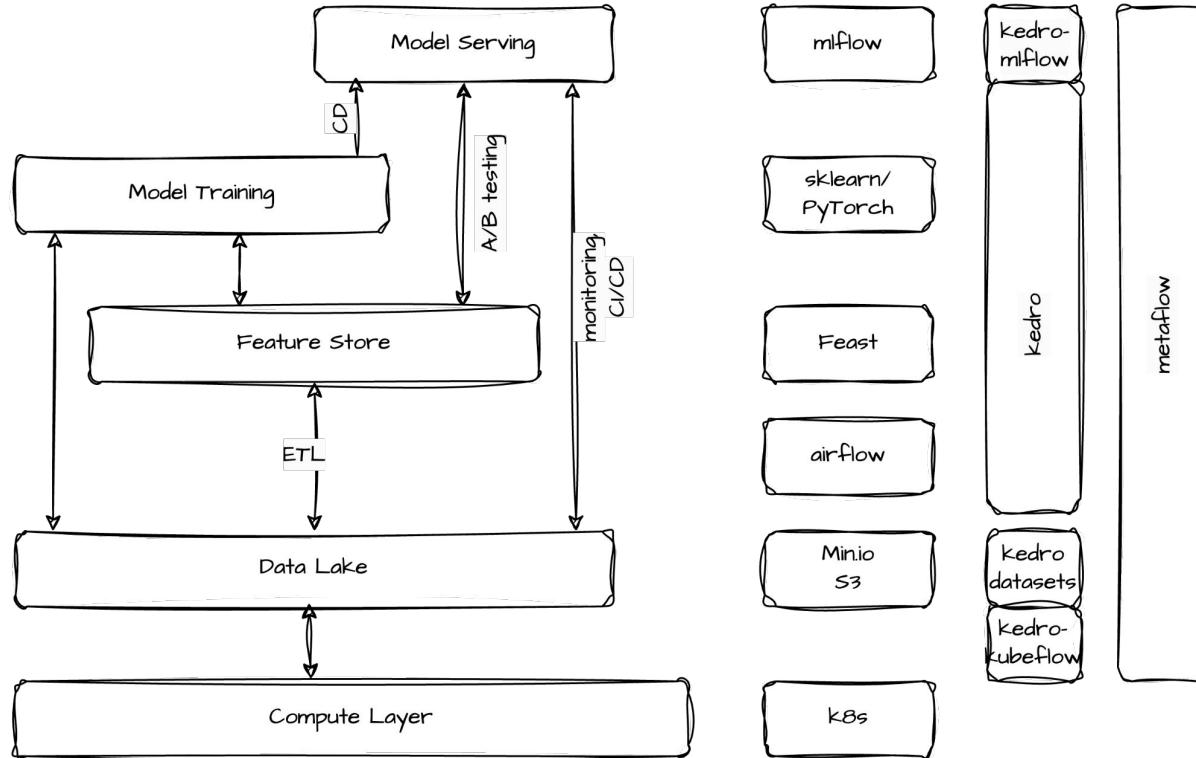
- Tech stacks could be different.
- Dev and Prod env could be different.
- Scalability could be an issue!
- Most importantly, skill sets could be different



Automate and Delegate the responsibility to Tools

Layer	Promise	Task	Tool
Solution	Reproducibility	Versioning - Data, Code, Model	<a href="#">DVC</a> , <a href="#">Kedro</a> , <a href="#">mlflow</a> , <a href="#">MetaFlow</a>
	Observability	Logging   Monitoring   Tracking	<a href="#">Kedro/mlflow</a> , <a href="#">MetaFlow</a>   <a href="#">evidently</a>   <a href="#">WandB</a>
	Agility	EDA   Experiment Tracking	Notebook environment ( <a href="#">Kedro</a> , <a href="#">MetaFlow</a> )   <a href="#">mlflow</a> , <a href="#">Weights and Biases</a>
	CI/CD	Integrate   Deploy   Serve	<a href="#">git-actions</a>   <a href="#">docker</a> / <a href="#">Kedro</a> / <a href="#">mlflow</a>   <a href="#">FastAPI</a>
Docs	Truthful/ Up to date	Code   Data   Solution	<a href="#">quarto</a> / <a href="#">sphinx</a>
Code	Testability  Readability	Linting   Testing   Documentation	<a href="#">ruff</a> / <a href="#">black</a>   <a href="#">pytest</a>   <a href="#">quarto</a> / <a href="#">sphinx</a>
Data	Veracity	Drift Detection and Schema Validation	<a href="#">pymfe</a> / <a href="#">pydantic</a> / <a href="#">evidently</a> / <a href="#">greateexpectations</a>
		ETL   Feature Stores	<a href="#">airflow</a>   <a href="#">feast</a>
	Agility   Auditability	Data Lake	<a href="#">minio</a> / s3 + athena   <a href="#">Kedro Data Catalogues</a>
Compute	Scalability	Elastic Compute	<a href="#">k8s</a>

Quite a number of OSS tools available to address specific gaps, with overlaps

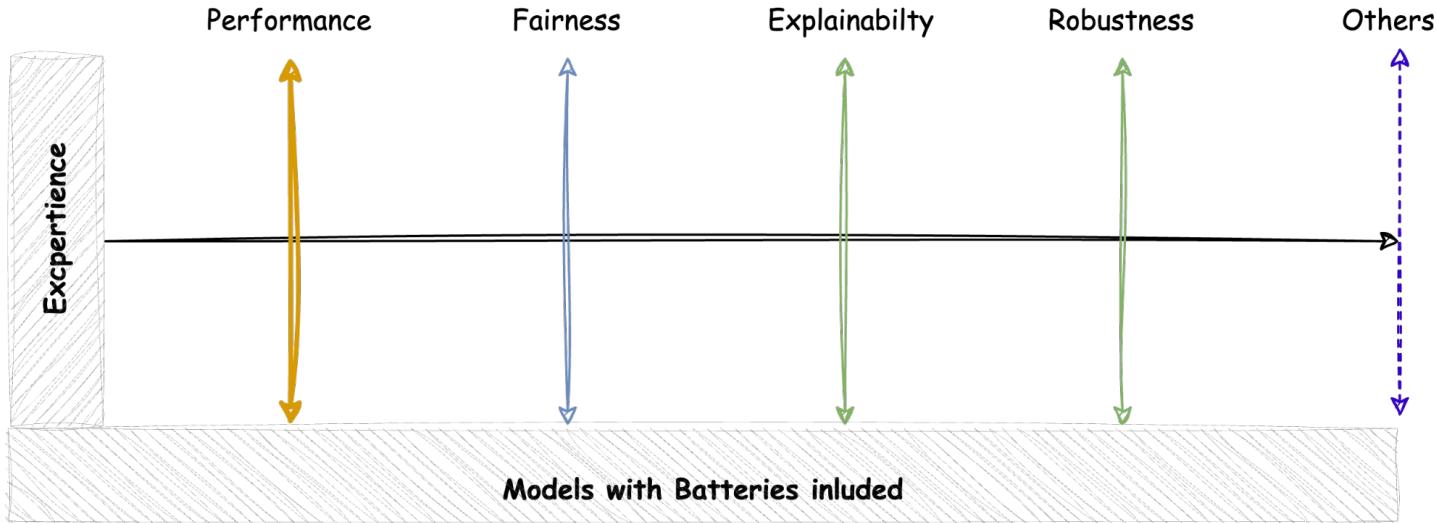


## ML POD: ML Engineering

An (opinionated) **git template** with integrations enabled

# OSS for ML Science

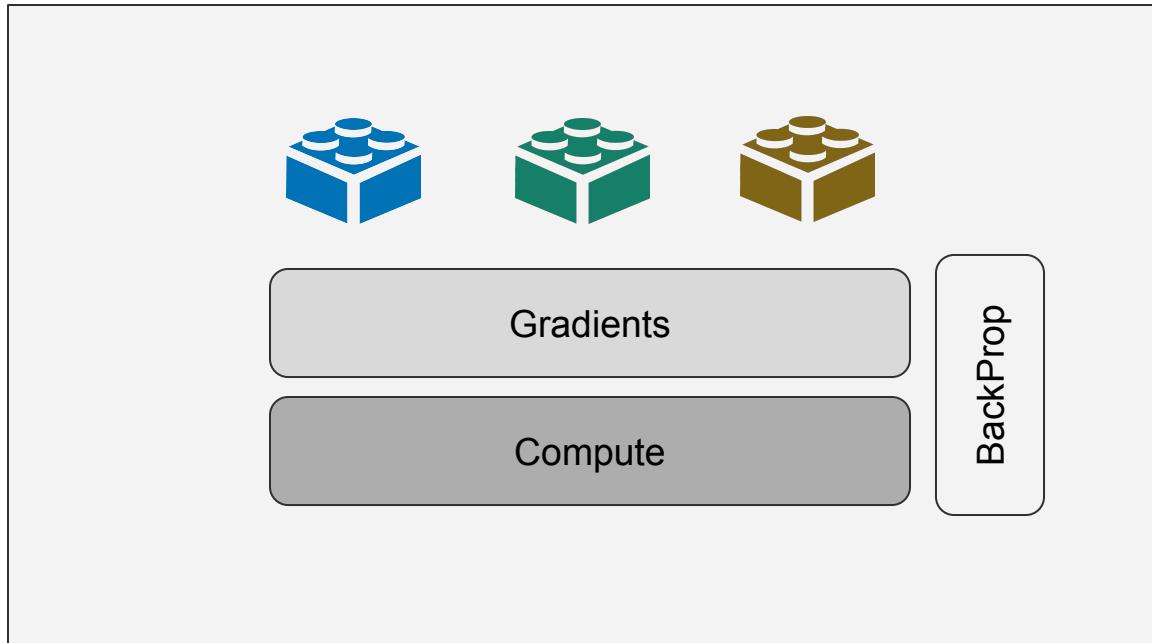
## Performance-centric ML development



But what is an enabling technology that addresses these sub-orthogonal fields?

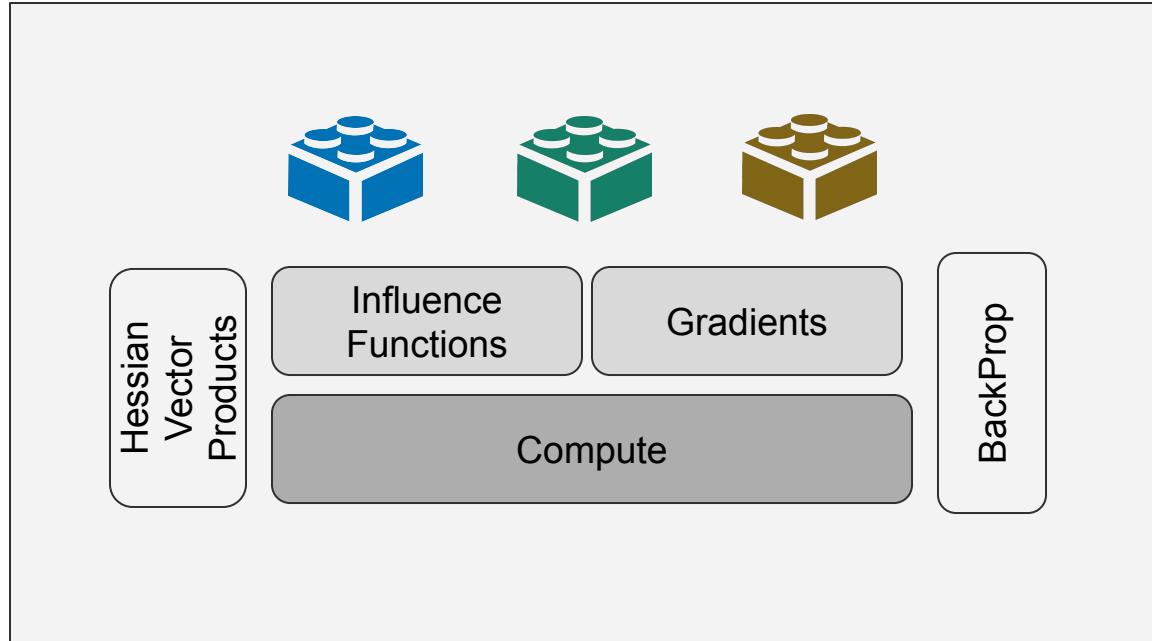
inputs	layers	losses	optimizers	runtime	outputs
					
Text Speech Image Numerical	Dense Conv GRU Attention	CE MSE MAE ...	Adam RMSProp NeverGrad ...	cpu gpu tpu	Text Speech Image Numeric al
Fixed Variable	Merge LDA/QDA	CheckLo ss Hinge Ordinal	PSO		Fixed Variabl e
Single Many	DT FM				Single Many
	Kernels				

Deep Learning as a technology breaks monolithic scientific computing paradigm.  
 It is OOPs for scientific computing at large.



## ML Engineering Backbone

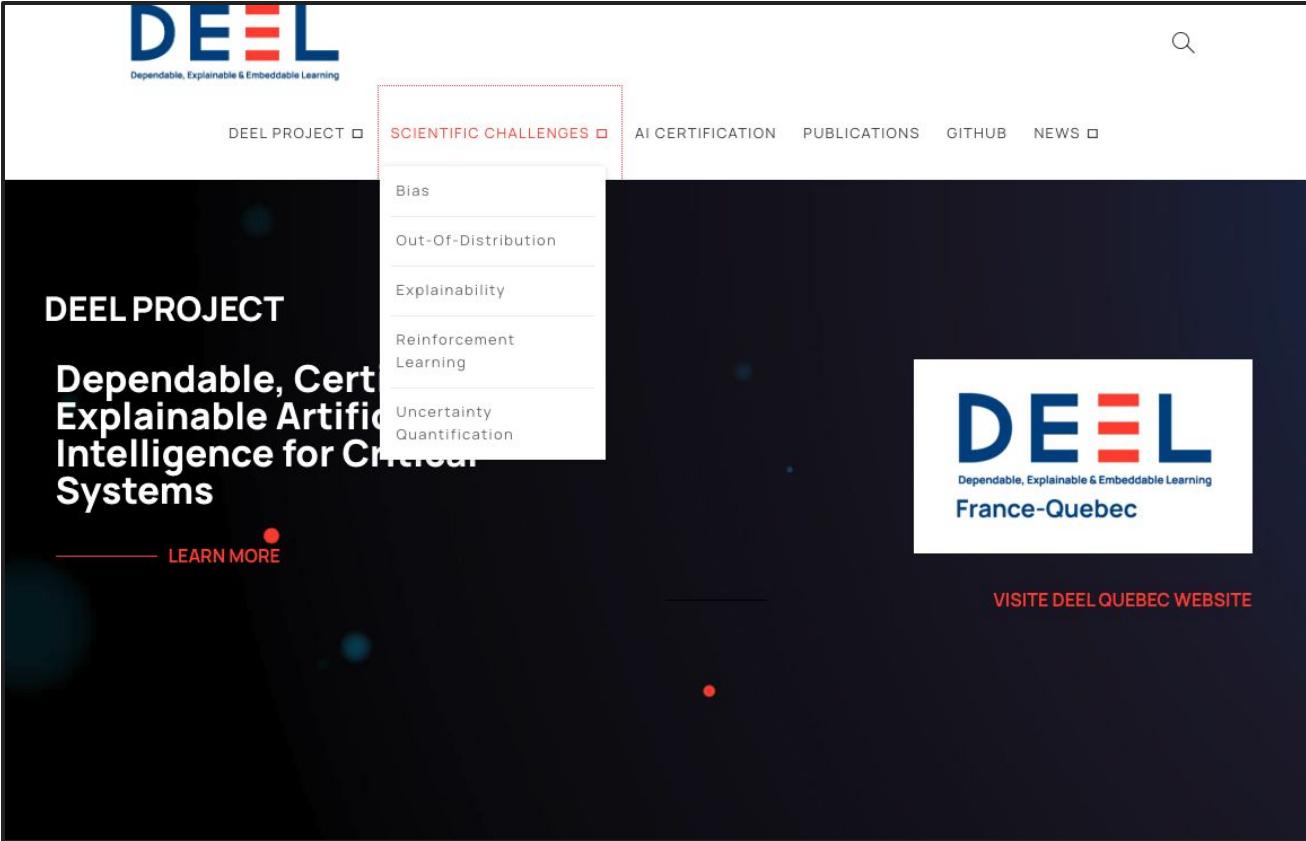
Scalable computation of gradients is at the heart of modern Deep Learning Engineering



ML Science Backbone

Perhaps, the answer is in IHVP

IVHF: [pyDVF](#)  
IFs: [influentiate](#):  
For more see: [AI-839@IIIT-B](mailto:AI-839@IIIT-B)



The image shows the landing page for the DEEL Project. The header features the DEEL logo (blue 'D' and red 'EEL') and the tagline "Dependable, Explainable & Embeddable Learning". A navigation bar includes links for "DEEL PROJECT", "SCIENTIFIC CHALLENGES" (which is expanded to show "Bias", "Out-Of-Distribution", "Explainability", "Reinforcement Learning", and "Uncertainty Quantification"), "AI CERTIFICATION", "PUBLICATIONS", "GITHUB", and "NEWS". A search icon is also present. The main content area has a dark background with white text. It features the title "DEEL PROJECT" and a large subtitle: "Dependable, Certifiable, Explainable Artificial Intelligence for Critical Systems". Below this is a "LEARN MORE" button with a red dot. To the right, there is a logo for "DEEL France-Quebec" with the same branding as the main site. At the bottom right, a call-to-action says "VISITE DEEL QUEBEC WEBSITE".

DEEL  
Dependable, Explainable & Embeddable Learning

DEEL PROJECT □ SCIENTIFIC CHALLENGES □ AI CERTIFICATION PUBLICATIONS GITHUB NEWS □

Bias  
Out-Of-Distribution  
Explainability  
Reinforcement Learning  
Uncertainty Quantification

DEEL  
Dependable, Explainable & Embeddable Learning  
France-Quebec

VISITE DEEL QUEBEC WEBSITE

for more see: [Project DEEL](#)

Promise	Methods	Tool
Statistical Reliability	Conformal Prediction Techniques Learn-Then-Test Framework	<a href="#">torchcp</a> , <a href="#">mapie</a> , <a href="#">puncc</a>
OOD	RMD, Entropy, Clustering	<a href="#">pytorch-ood</a> , <a href="#">oodeel</a>
Security	Adversarial Attacks	<a href="#">foolbox</a>
Explainability	Integrated Gradients SHAPLEY values, LIME, Decision Sets, Influence Functions Counterfactuals	<a href="#">shap</a> , <a href="#">xplique</a> , <a href="#">dice</a>
Bias and Fairness	Metrics (for detection) and Methods (to mitigate) Influence Functions	<a href="#">AI Fairness 360</a> , <a href="#">influentiate</a>
Robustness	Adversarial Training Robust Losses Regularization	All of above

OOD, Robustness, Security, XIA, Bias & Fairness – share common principles!

# OSS Model ML Development

**chatGPT happened in 2022.**

**We need to develop LLMs that are truly open source.**

**But**

- Most open source LLMs are English-centric
- Most cases, we know the model weights but not the training codes
- Developing LLMs requires enormous
  - Money
  - Compute
  - Talent

**How can we address them?**

## Desirable aspects

- Decentralises the development (training) of LLMs/  
SLMs
- Supports Indic Languages (or other languages  
besides English)

## Development Process

- A git like system for models
  - pull pretrained model and data
  - peft low cost update to the model
  - push the adapters
  - merge adapters to pretrained model
- continue pre-training

**compatible with HuggingFace ecosystem**

<https://github.com/mlsquare/fedem>

# Benefits of Open Source LLMs

## On the Societal Impact of Open Foundation Models

Sayash Kapoor<sup>\*1</sup> Rishi Bommasani<sup>\*2</sup>

Kevin Klyman<sup>2</sup> Shayne Longpre<sup>3</sup> Ashwin Ramaswami<sup>4</sup> Peter Cihon<sup>5</sup> Aspen Hopkins<sup>3</sup>  
 Kevin Bankston<sup>6,4</sup> Stella Biderman<sup>7</sup> Miranda Bogen<sup>6,1</sup> Rumman Chowdhury<sup>8</sup> Alex Engler<sup>9</sup>  
 Peter Henderson<sup>1</sup> Yacine Jernite<sup>10</sup> Seth Lazar<sup>11</sup> Stefano Maffulli<sup>12</sup> Alondra Nelson<sup>13</sup>  
 Joelle Pineau<sup>14</sup> Aviya Skowron<sup>7</sup> Dawn Song<sup>15</sup> Victor Storchan<sup>16</sup> Daniel Zhang<sup>2</sup>  
 Daniel E. Ho<sup>2</sup> Percy Liang<sup>2</sup> Arvind Narayanan<sup>1</sup>

February 27, 2024

### Benefits of Open Foundation Models

The distinctive properties of open foundation models allow us to critically analyze key benefits for open foundation models that emerge from these properties.

#### Distributing who defines acceptable model behavior

*Broader access and greater customizability expand who is able to decide acceptable model behavior.*

#### Increasing innovation

*Broader access, greater customizability, and local inference expand how foundation models are used to develop applications.*

#### Accelerating science

*Broader access and greater customizability facilitate scientific research. The availability of other key assets (such as training data) would further accelerate scientific research.*

#### Enabling transparency

*Broad access to weights enables some forms of transparency. The availability of other key assets (such as documentation and training data) would further improve transparency.*

#### Mitigating monoculture and market concentration

*Greater customizability mitigates the harms of monoculture and broader access reduces market concentration.*

## A new way to develop LLMs

COMMUNICATIONS  
OF THE ACM

Explore Topics ▾

Latest Issue ▾



Sign In

OPINION

[Computing Applications](#)

Viewpoint

# Building Machine Learning Models Like Open Source Software

Proposing a community-based system for model development.

By [Colin Raffel](#)

Posted Feb 1 2023

This Viewpoint advocates for tools and research advances that will allow pretrained models to be built in the same way that we build open source software. Specifically, models should be developed by a large community of stakeholders that continually updates and improves them. Realizing this goal will require porting many ideas from open source software development to the building and training of pretrained models, which motivates many new research problems and connections to existing fields.

<https://cacm.acm.org/opinion/building-machine-learning-models-like-open-source-software/>

# A new way to develop LLMs

## Training Neural Networks from Scratch with Parallel Low-Rank Adapters

Minyoung Huh<sup>1</sup> Brian Cheung<sup>1,2</sup> Jeremy Bernstein<sup>1</sup> Phillip Isola<sup>1</sup> Pulkit Agrawal<sup>1</sup>

### Abstract

The scalability of deep learning models is fundamentally limited by computing resources, memory, and communication. Although methods like low-rank adaptation (LoRA) have reduced the cost of model finetuning, its application in model pre-training remains largely unexplored. This paper explores extending LoRA to model pre-training, identifying the inherent constraints and limitations of standard LoRA in this context. We introduce *LoRA-the-Explorer* (LTE), a novel bi-level optimization algorithm designed to enable parallel training of multiple low-rank heads across computing nodes, thereby reducing the need for frequent synchronization. Our approach includes extensive experimentation on vision transformers using various vision datasets, demonstrating that LTE is competitive with standard pre-training.

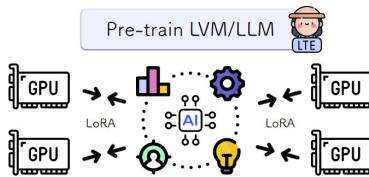
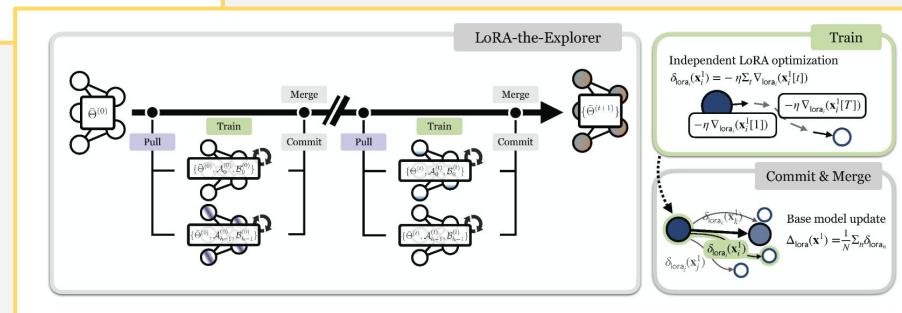


Figure 1: **Lora-The-Explorer:** We propose LoRA-the-explorer, an optimization algorithm that can match the performance of standard training from scratch. Our method optimizes unique LoRA parameters in parallel and merges them back to the main weights. Our algorithm can leverage lower-memory devices and only depends on communicating the LoRA parameters for training, making it an ideal candidate in a bandlimited or memory-constraint training framework.<sup>1</sup>



Expectation



Reality



Make ML great again :)