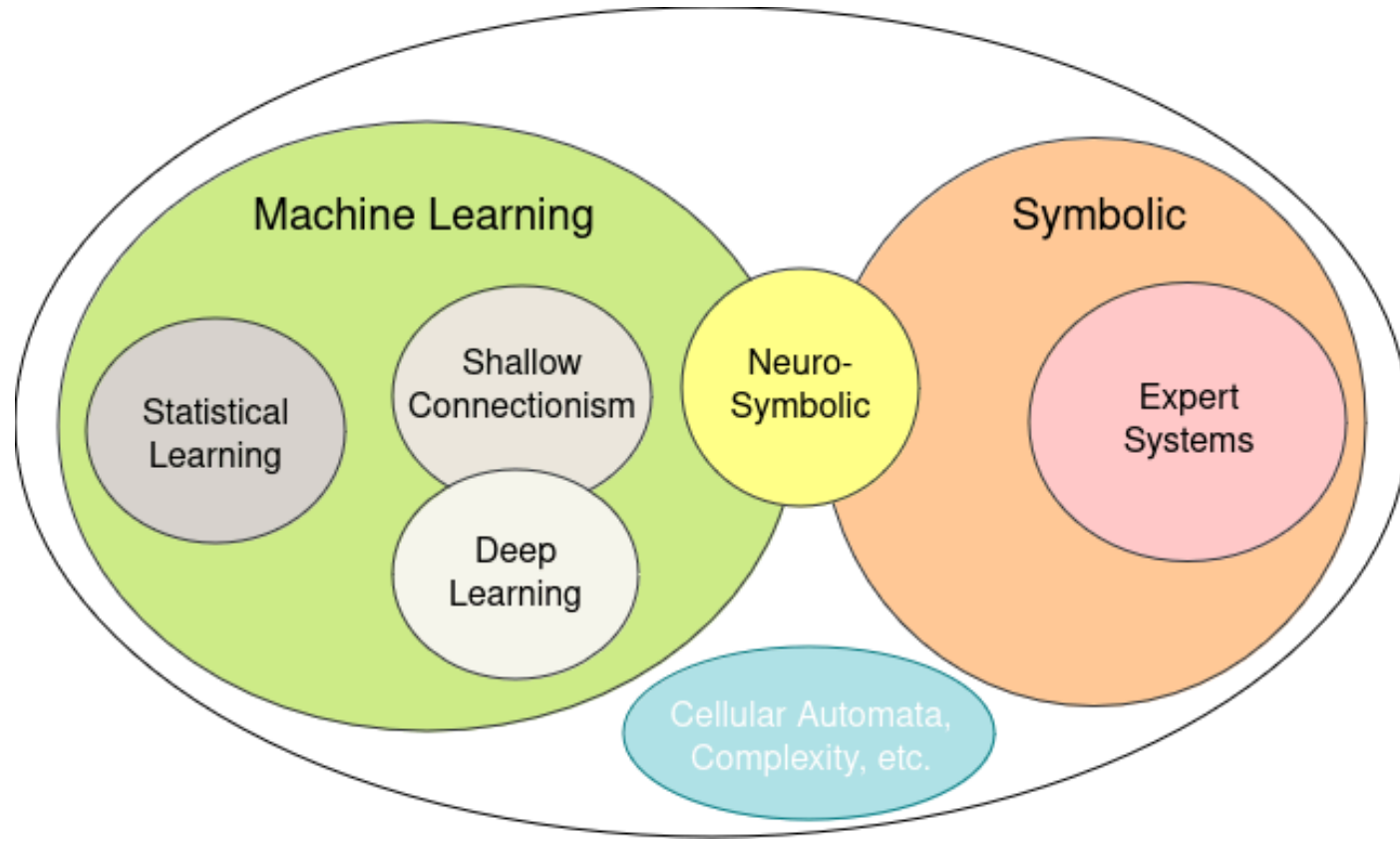




Frugal AI

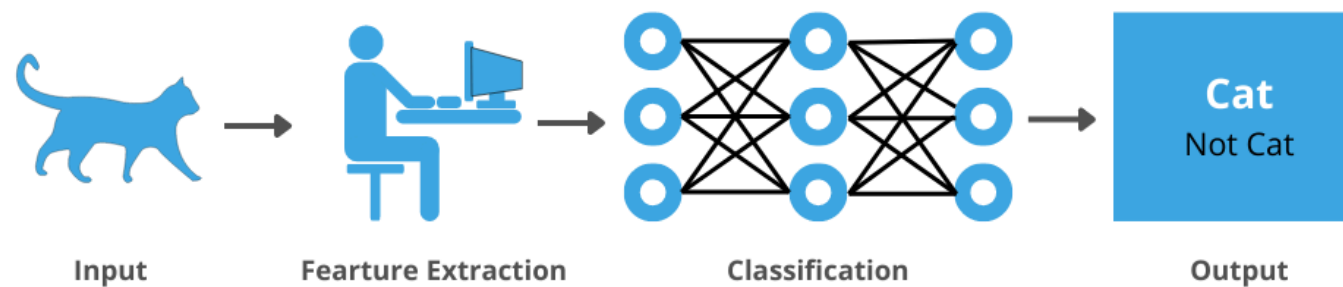
Romain Belmonte
(september 24)

Artificial Intelligence

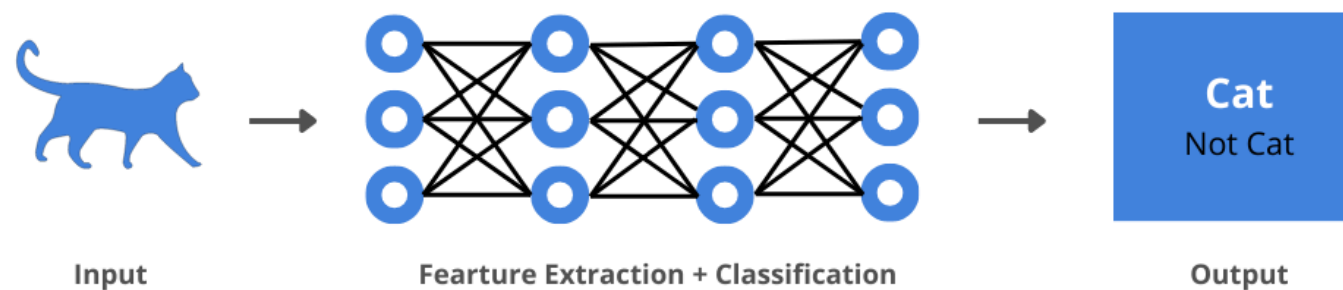


ML / DL

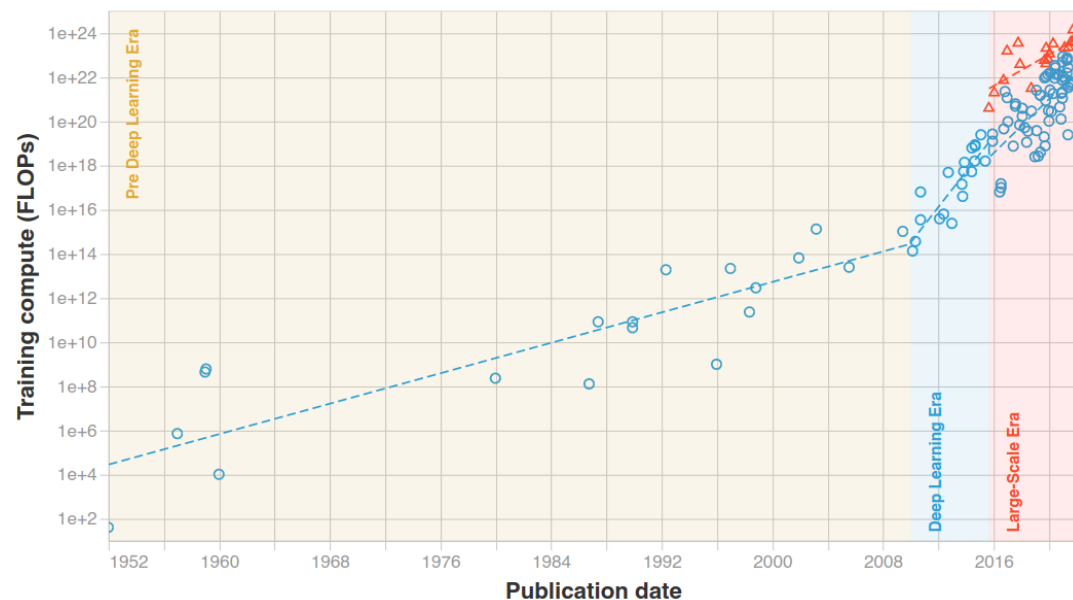
Machine Learning



Deep Learning



AI Scaling Laws



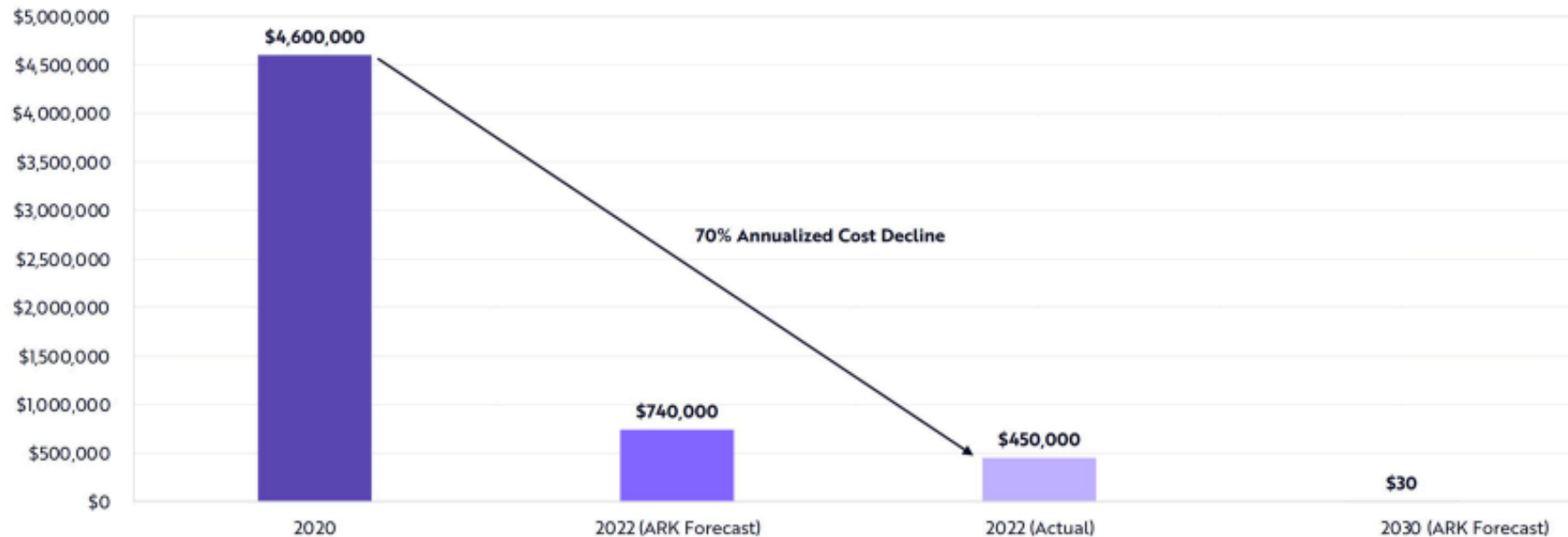
- Loss scales as a power-law with model size, dataset size, and the amount of compute used for training

[[3](#), [4](#)]

Footprint/Cost

Llama 3.1 - 405b

- 16 thousand H100 = a cumulative of 30.84M GPU hours of computation
- 8930 tons of CO2e (renewable energy) - an average French person = 9 tons/year

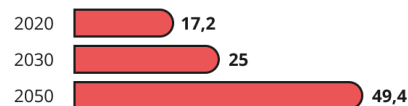


IPCC guidelines

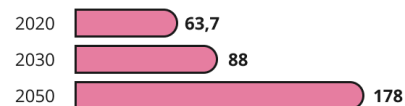
Des indicateurs issus d'une méthode normalisée d'analyse de cycle de vie qui comporte des définitions précises

Evolution de 4 indicateurs de l'impact environnemental du numérique dans le scénario tendanciel, en valeurs absolues.

Empreinte carbone (en millions de tonnes de CO₂ éq.)



Ressources utilisées (en millions de tonnes)



Consommation d'énergie (en TWh)



Conso. de métaux et minéraux (en tonnes Sb éq.)



Empreinte carbone : émissions de gaz à effets de serre exprimées en équivalent CO₂.

Ressources utilisées : indicateur MIPS qui considère cinq types de ressources, comprenant les ressources abiotiques (matériaux, énergie fossile...), la biomasse, les déplacements de terre mécaniques ou par érosion, l'eau, et l'air. Il donne une idée de l'effort effectué pour produire nos biens et services.

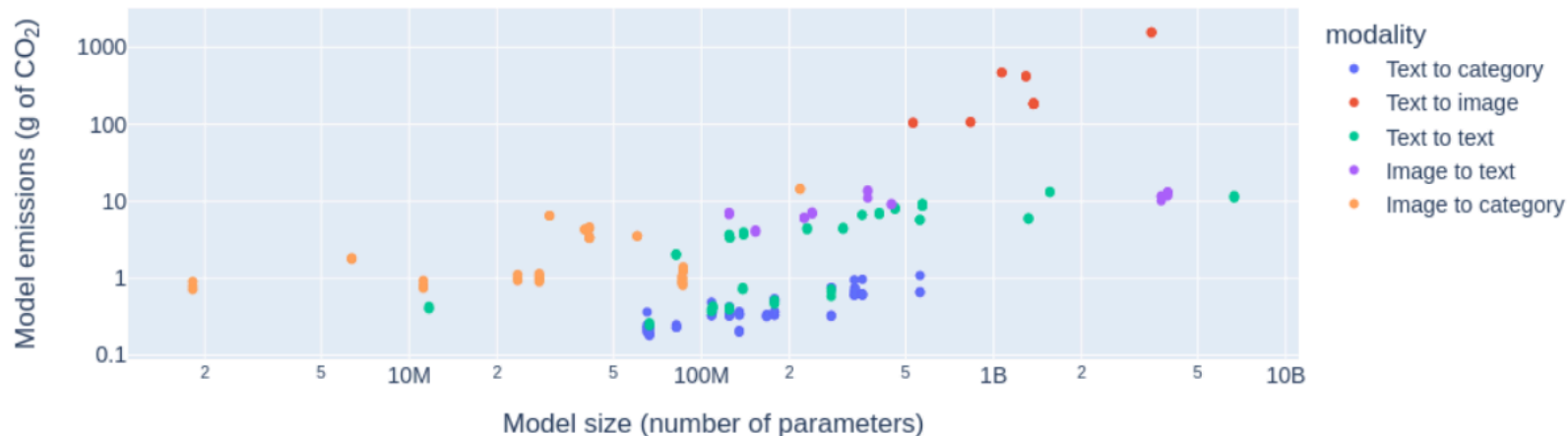
Consommation d'énergie finale : désigne l'énergie directement utilisée par l'utilisateur final, sous forme d'électricité ou de carburant.

Consommation de métaux et minéraux : cet indicateur évalue la quantité de ressources minérales et métalliques extraites de la nature en équivalent antimoine (un élément chimique dont on retrouve le symbole Sb dans le tableau périodique des éléments). C'est un standard des analyses de cycle de vie qui permet de mesurer l'épuisement des ressources naturelles.

Green/Red AI

Green/Red AI: Making efficiency an evaluation criterion + reporting the financial cost of developing, training, and running models

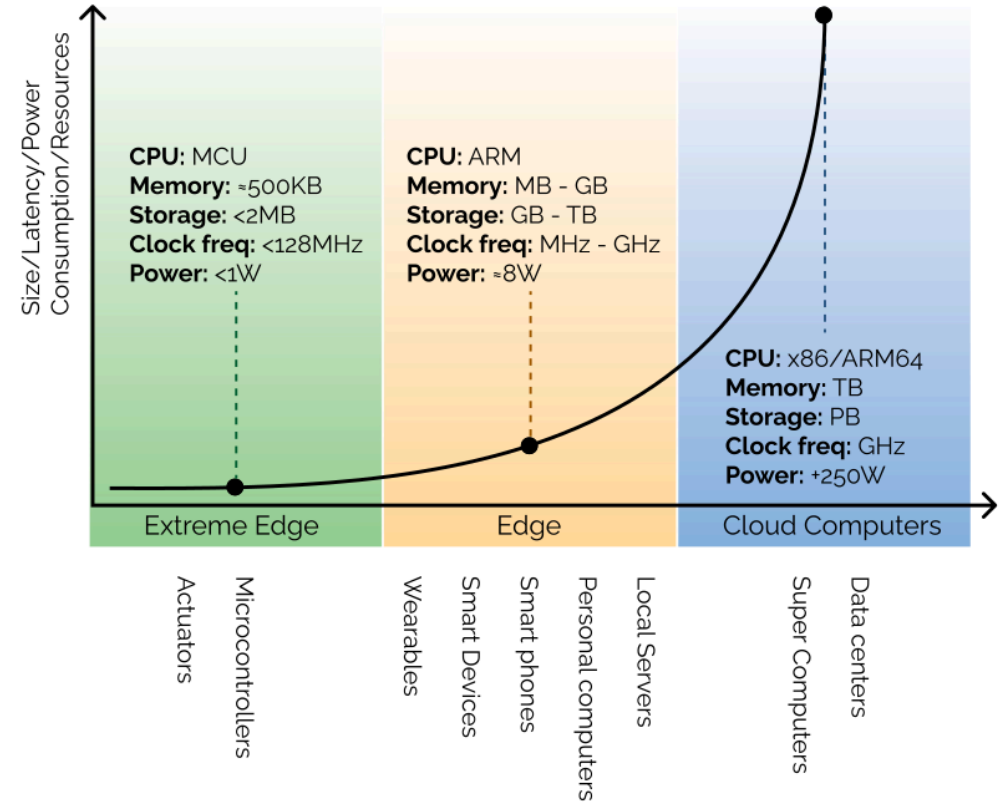
	BLOOMz-7B	BLOOMz-3B	BLOOMz-1B	BLOOMz-560M
Training energy (kWh)	51,686	25,634	17,052	10,505
Finetuning energy (kWh)	7,571	3,242	1,081	543
Inference energy (kWh)	1.0×10^{-4}	7.3×10^{-5}	6.2×10^{-5}	5.4×10^{-5}
Cost parity (# inferences)	592,570,000	395,602,740	292,467,741	204,592,592



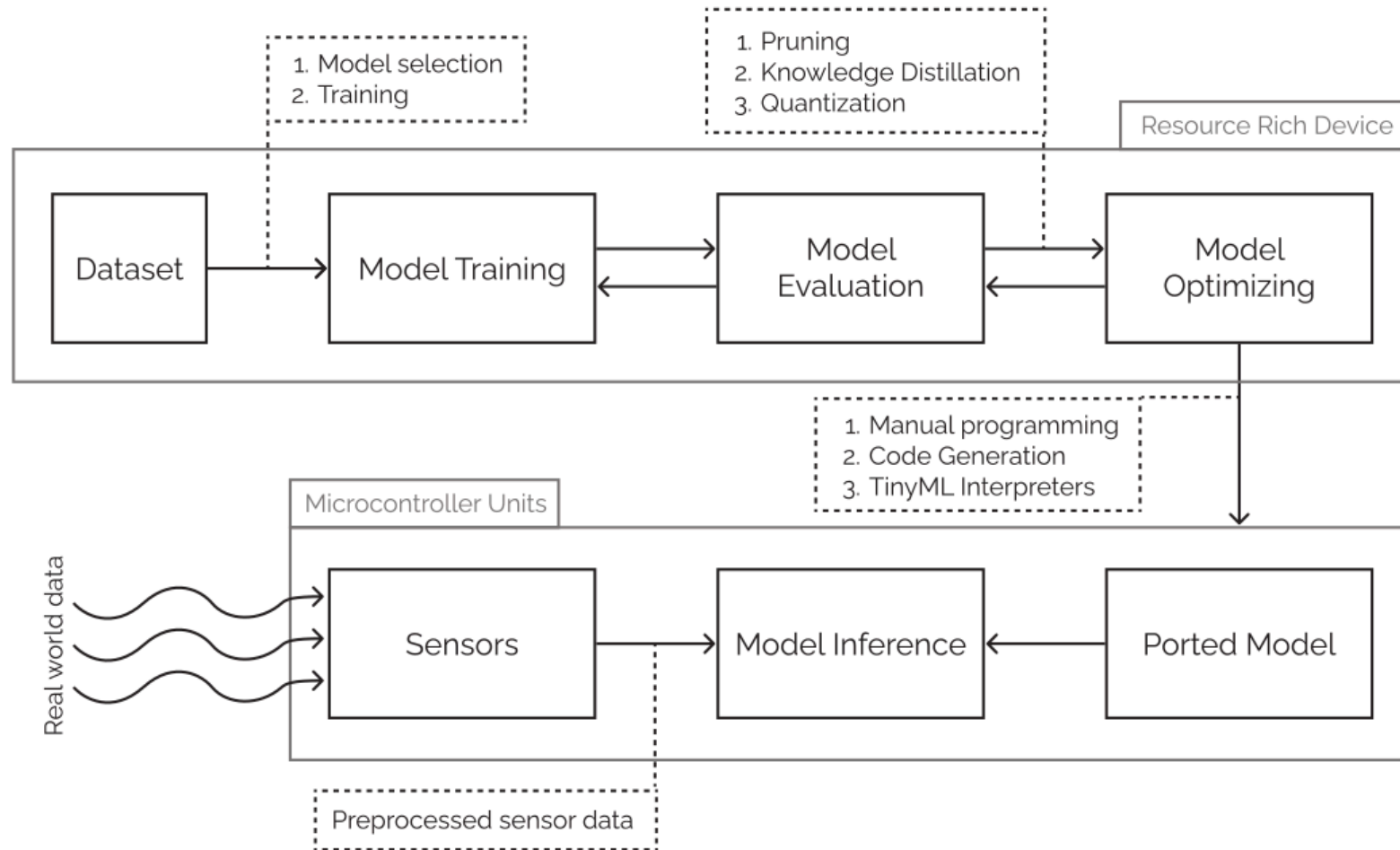
100 km with an electric car = 17 kWh; 1 h of video streaming = 70 g of CO₂

Challenges

- Limited Computational Power
- Memory Limitations
- Energy Efficiency
- Deployment/Maintenance Complexity
- Ethics, Privacy and Security Concerns



Pipeline



Overview

1. Data

- Data overload
- Unlabeled data
- Limited data
- Data management

2. Model

- Energy consumption
- Memory efficiency
- Training overhead

3. Miscellaneous

- Tools
- Resources
- Hardware
- Advances





DATA

Data

- **Data overload:** selection, dimensionality reduction, sampling, distillation
- **Unlabeled data:** crowdsourcing, generative models, active learning, semi-supervision, weak supervision, self-supervision
- **Limited data:** augmentation, synthetic generation, transfer learning, external data, regularization, physics-informed, few-shot learning
- **Data management:** imputation, cleaning, feature engineering, preprocessing, ensembling, cross-validation, incremental learning, memory augmented network

Data

Data overload: selection, dimensionality reduction, sampling, distillation

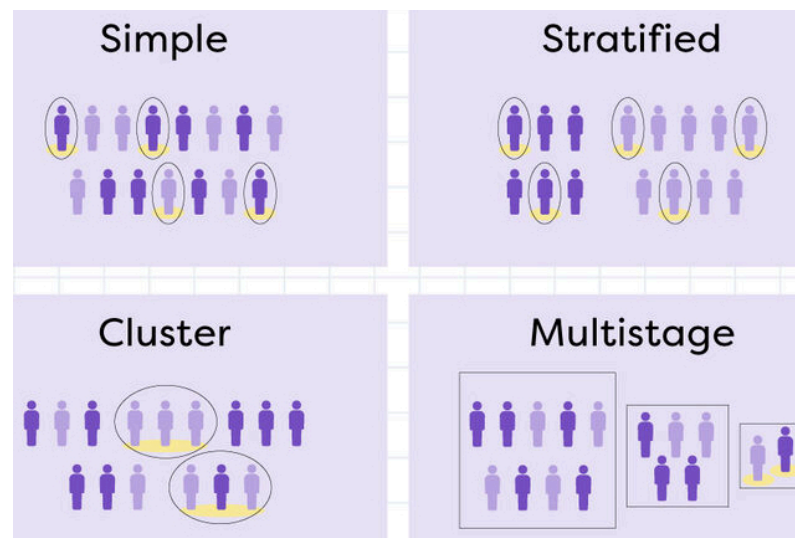
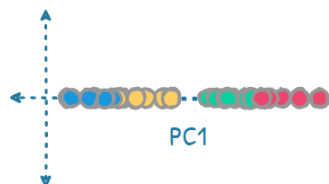
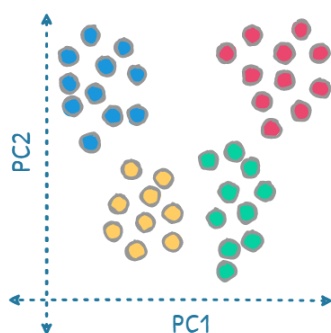
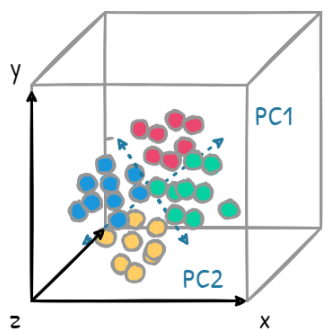
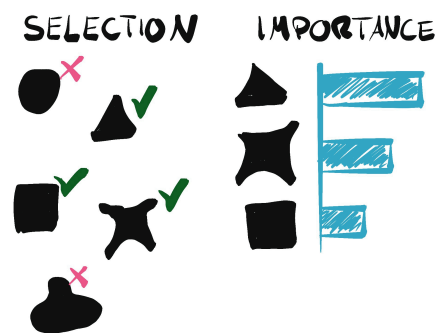
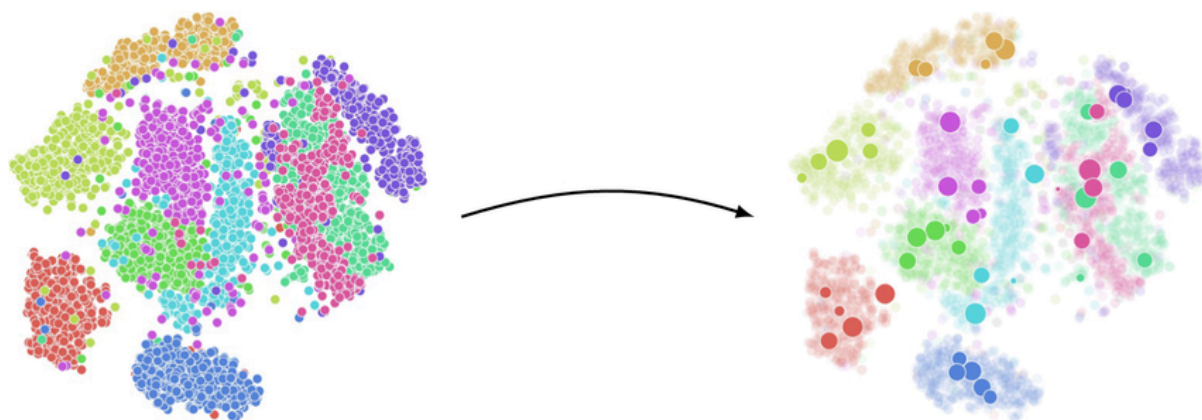
Data Overload

Garbage In, Garbage Out



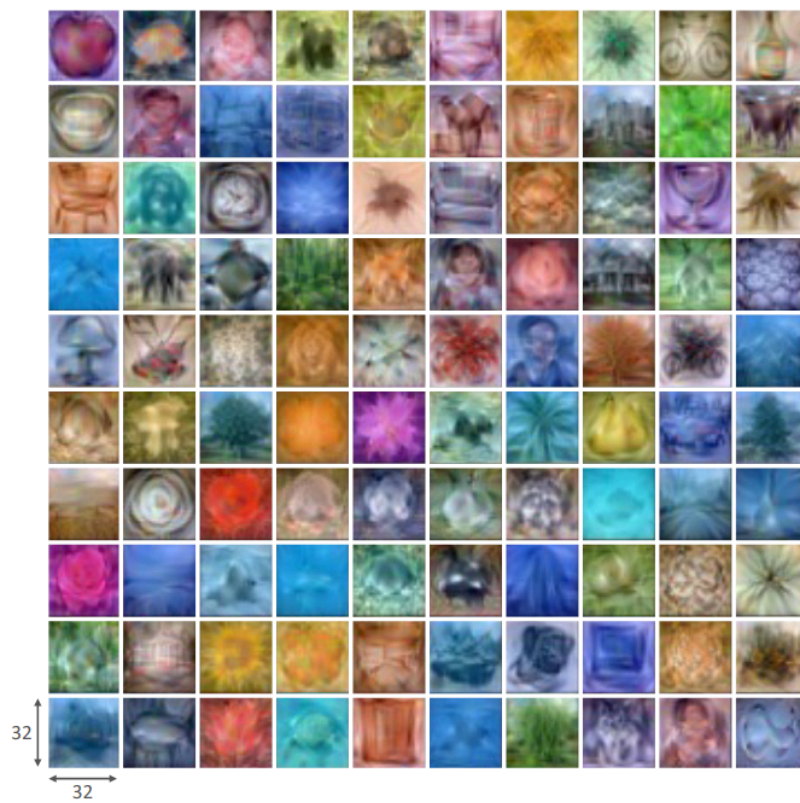
Data Overload

Selection, dimensionality reduction, sampling

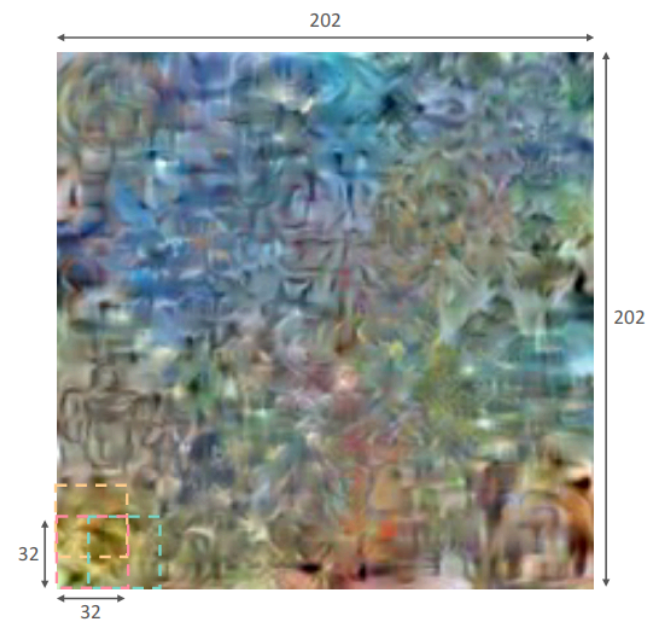


Data Overload

Distillation



Classic Dataset Distillation
1 Image-Per-Class



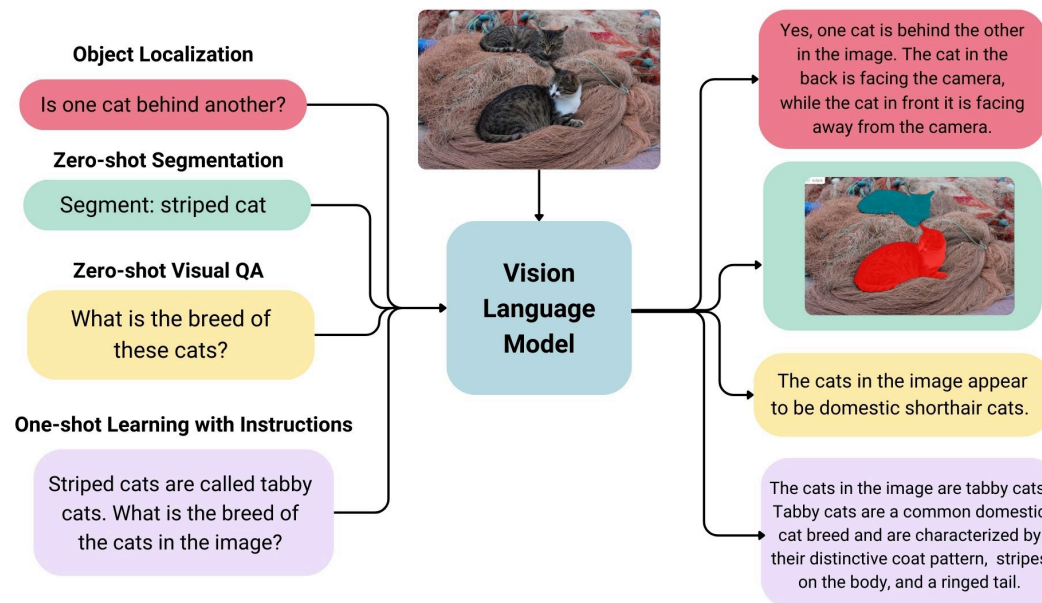
Poster Dataset Distillation (PoDD)
0.4 Image-Per-Class

Data

Unlabeled data: crowdsourcing, generative models, active learning, semi-supervision, weak supervision, self-supervision

Unlabeled Data

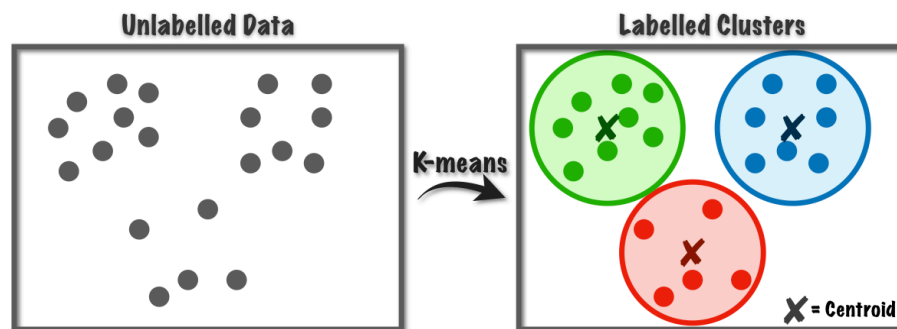
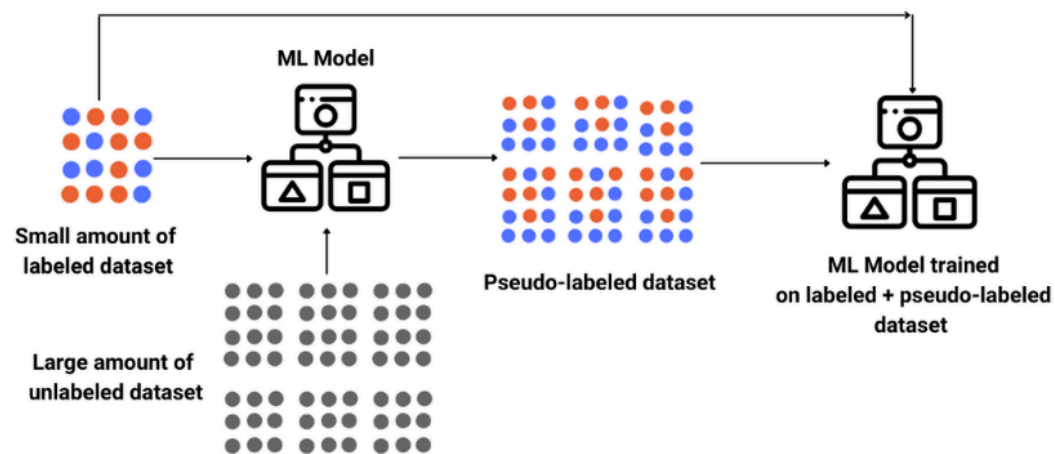
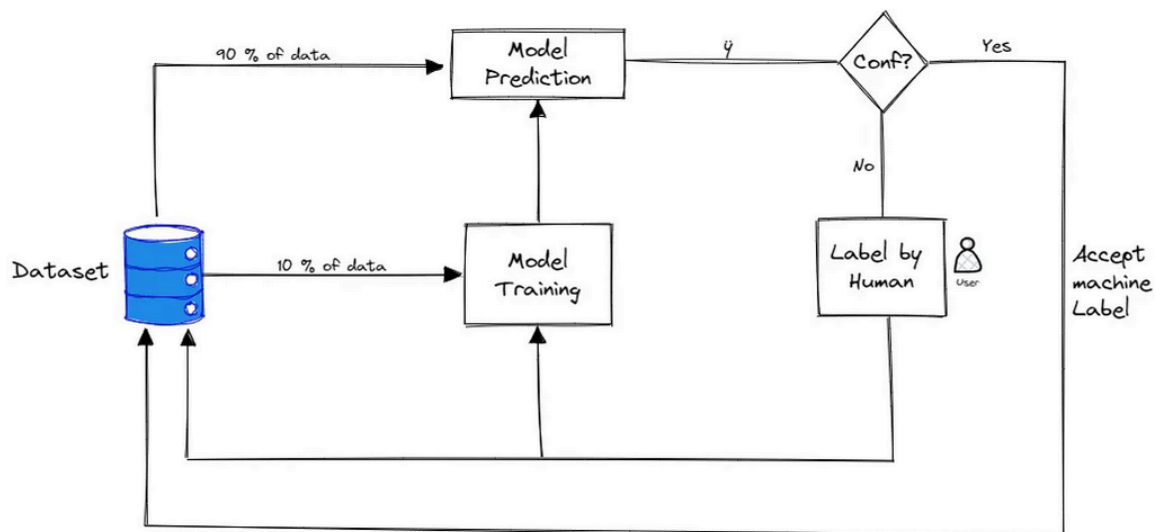
Crowdsourcing, generative models



[Pro bono projects with PeopleForAI](#)

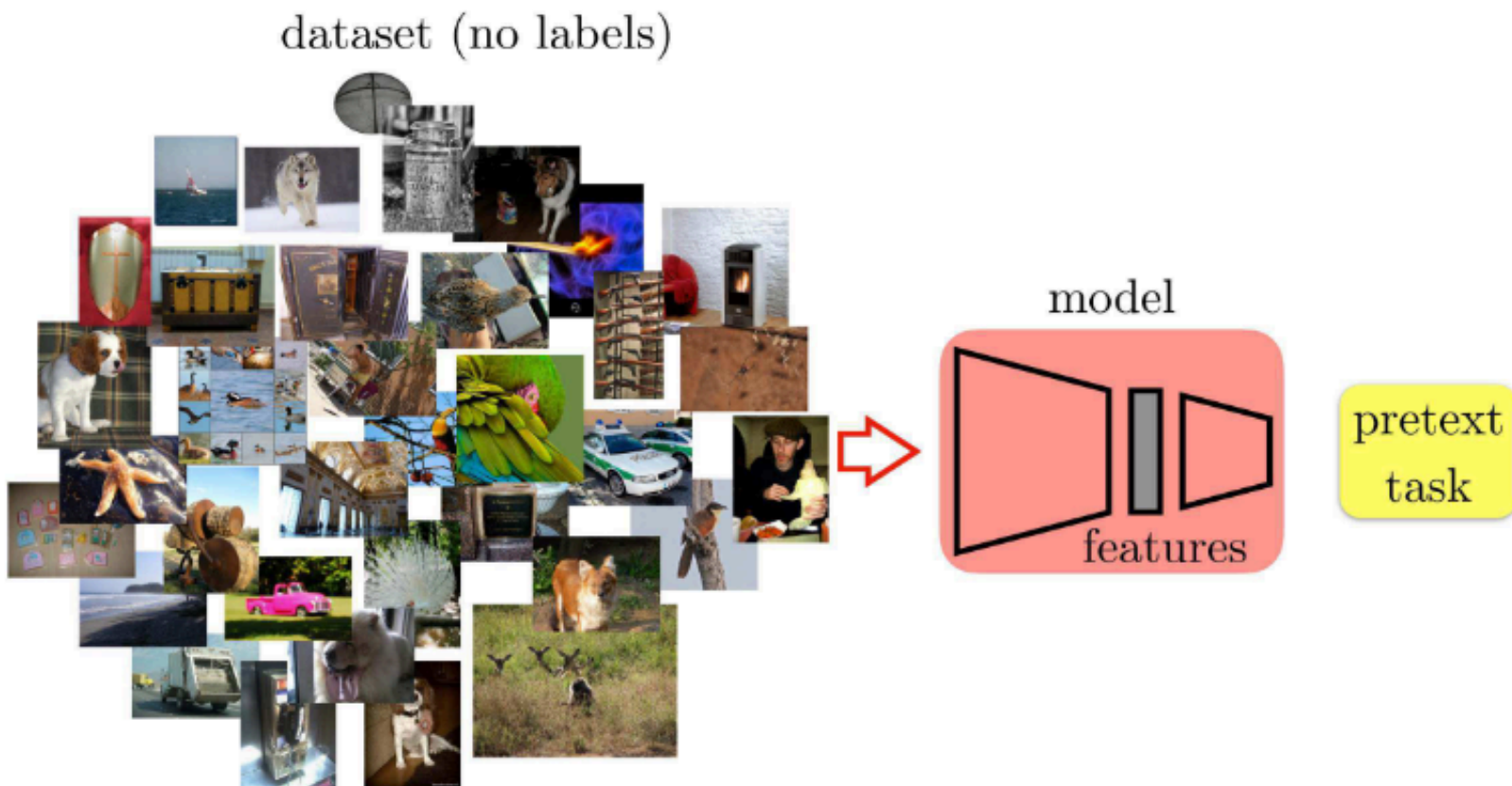
Unlabeled Data

Active learning, semi-supervision, weak supervision



Unlabeled Data

Self-supervision

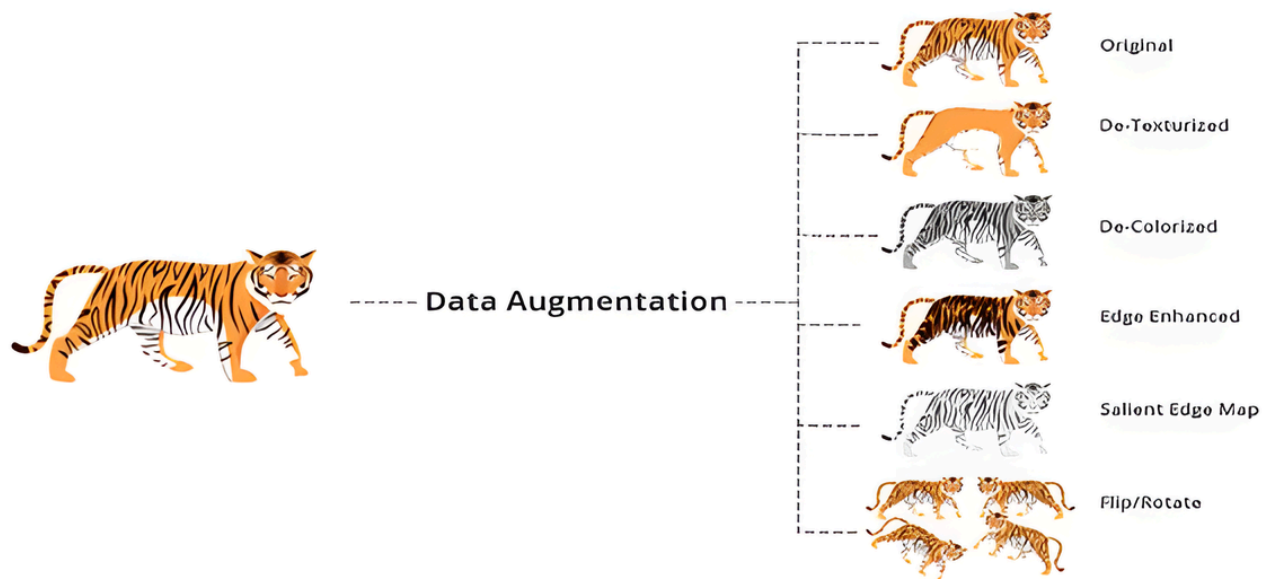


Data

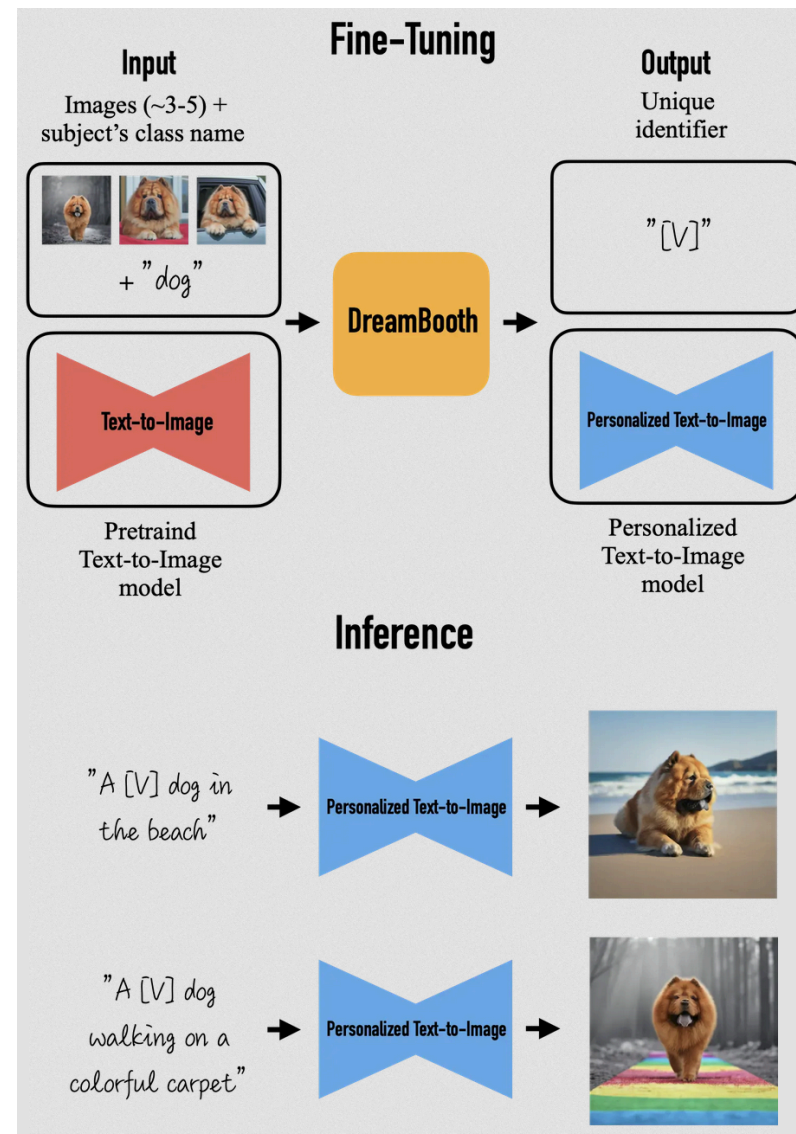
Limited data: augmentation, synthetic generation, transfer learning, external data, regularization, physics-informed, few-shot learning

Limited Data

Augmentation, synthetic generation

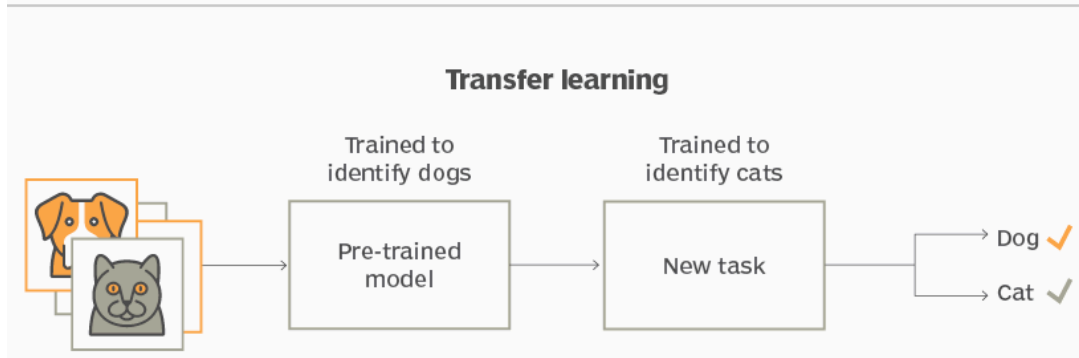
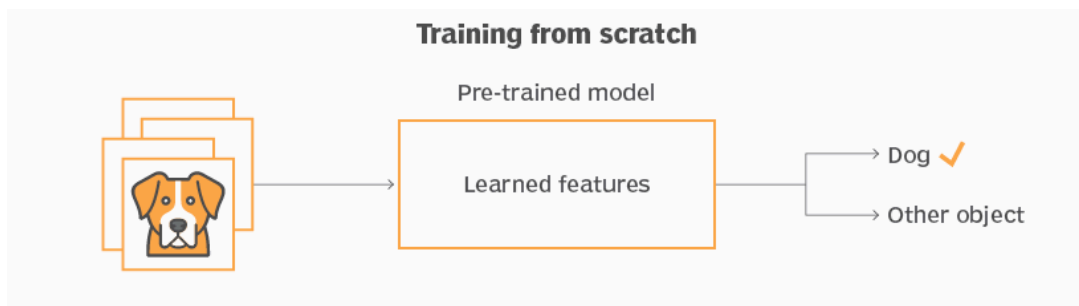


[Nvidia Omniverse](#)

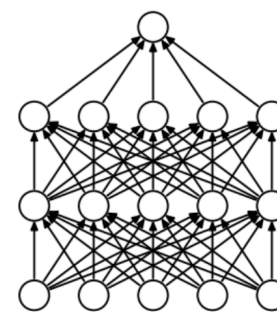


Limited Data

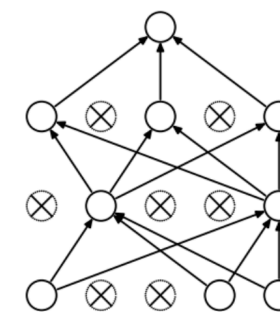
Transfer learning, external data, regularization, physics-informed



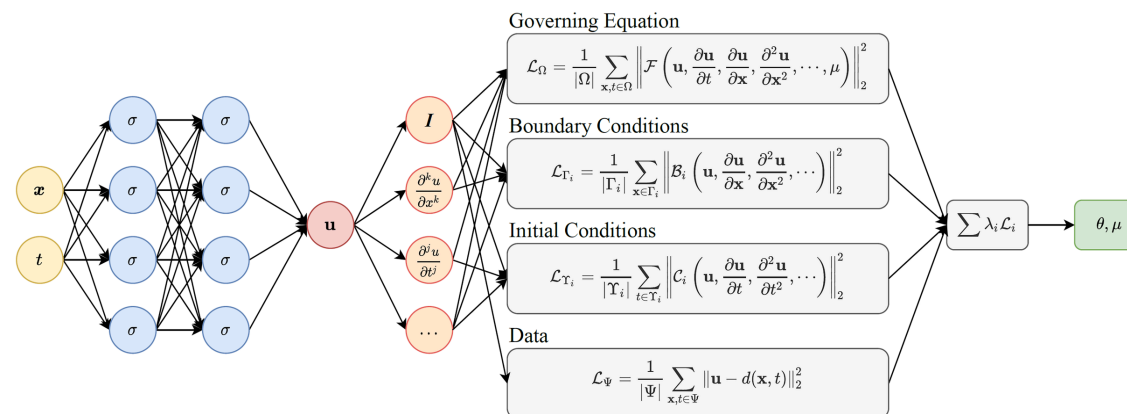
[17]



(a) Standard Neural Net

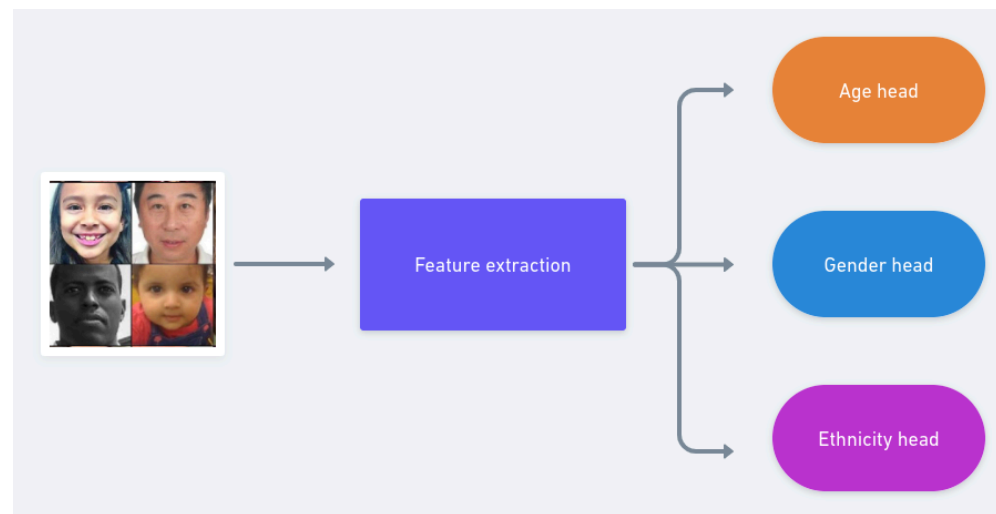
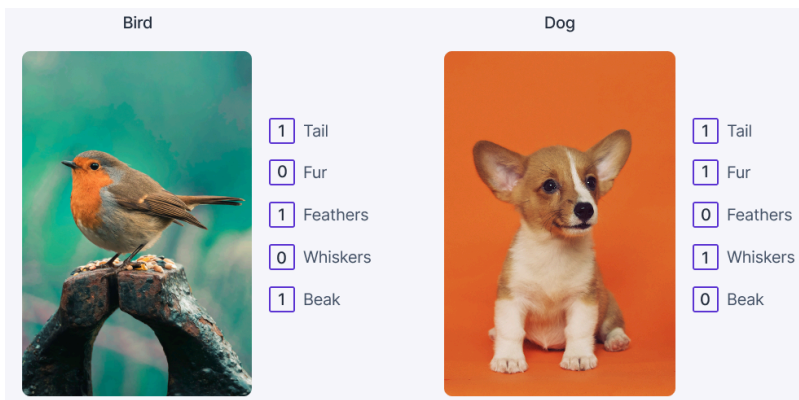
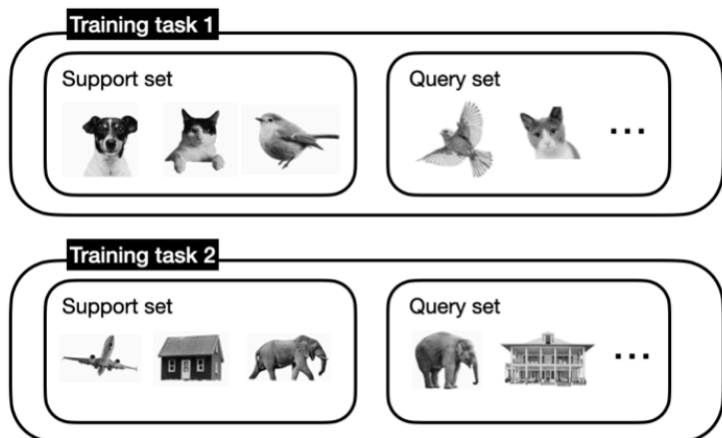


(b) After applying dropout.



Limited Data

Few-shot learning

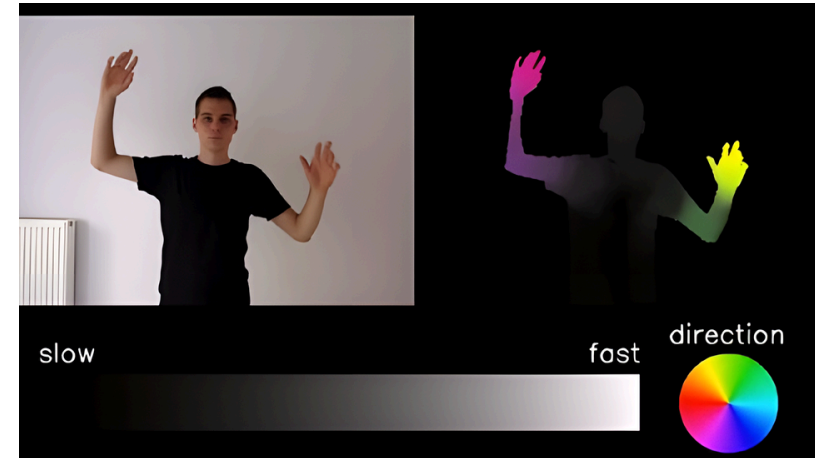
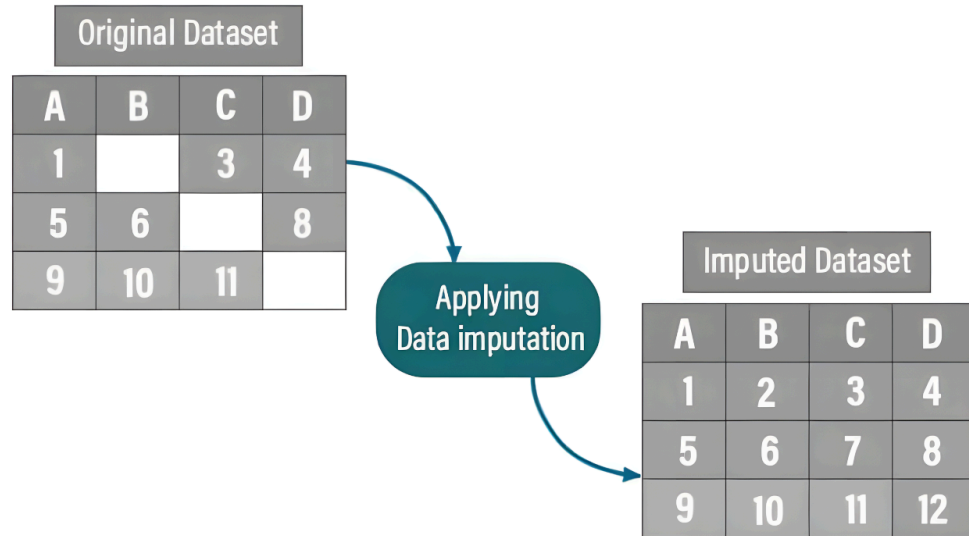


Data

Data management: imputation, cleaning, feature engineering, preprocessing, ensembling, cross-validation, incremental learning, memory augmented network

Data Management

Imputation, cleaning, feature engineering, preprocessing



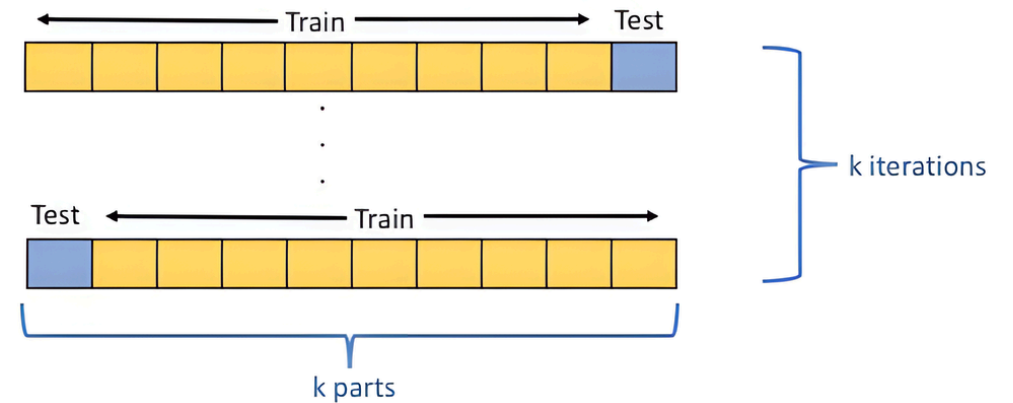
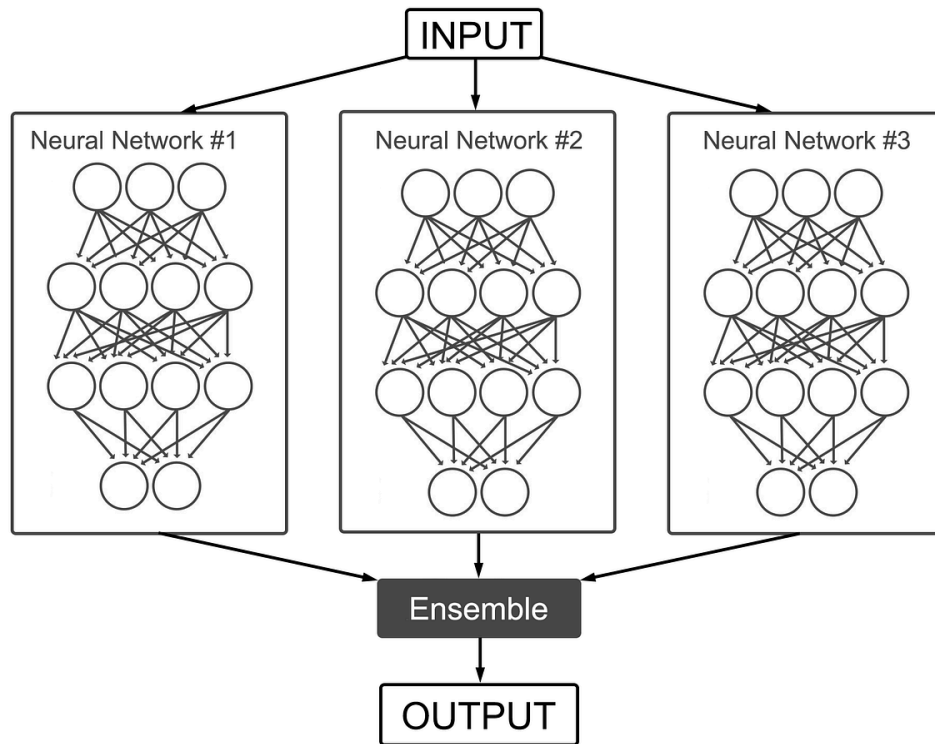
Original



Canny

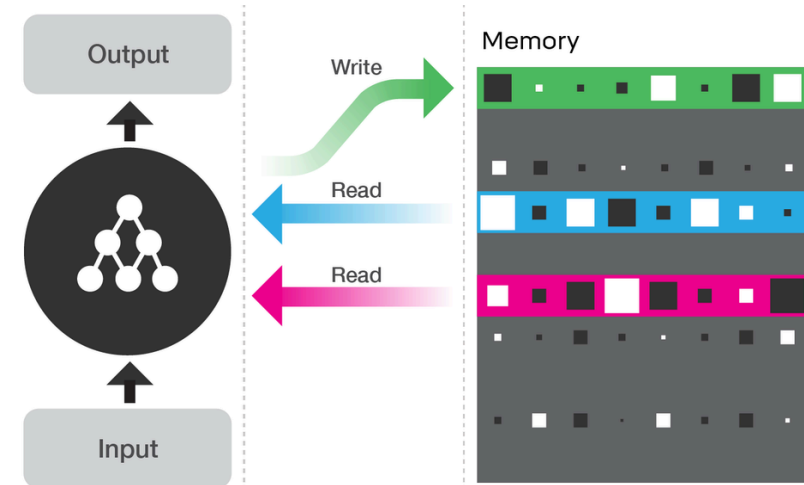
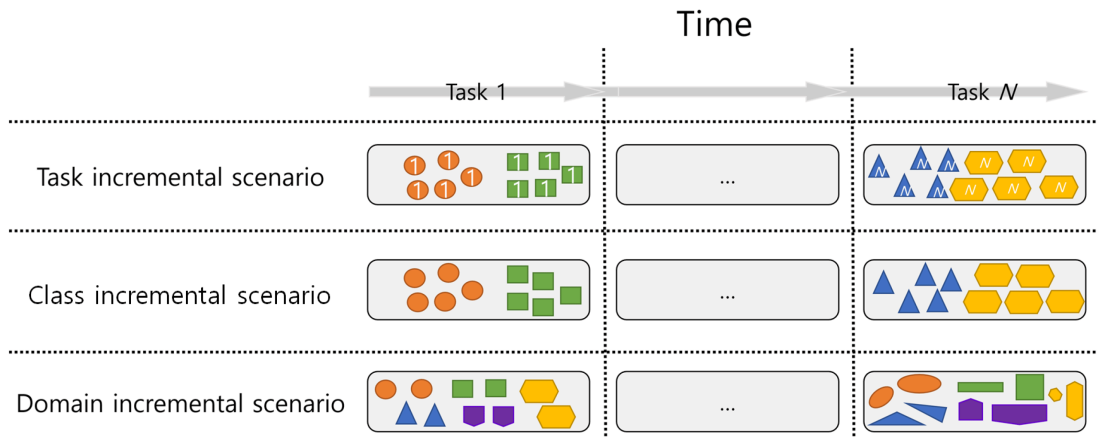
Data Management

Ensembling, cross-validation



Data Management

Incremental learning, memory augmented network



[[20](#), [21](#)]



MODEL

Model

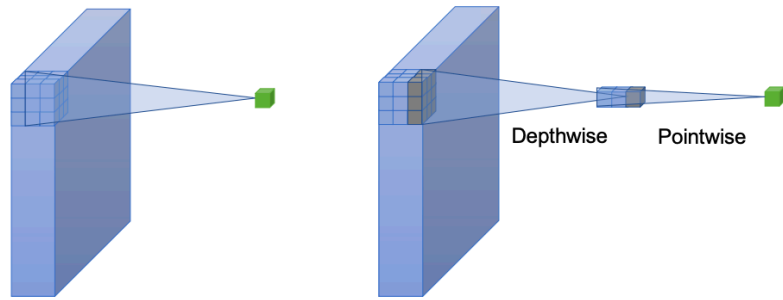
- **Energy consumption:** efficient algorithms-architectures, skip mechanism, gate mechanism, early termination, compilation, specialized hardware
- **Memory efficiency:** batching, accumulation, weight sharing, tensor decomposition, pruning, quantization, distillation, sparse representation
- **Training overhead:** early stopping, mixed precision, scheduling, distributed learning, dynamic pruning, low-rank adaptation

Model

Energy consumption: efficient algorithms-architectures, skip mechanism, gate mechanism, early termination, compilation, specialized hardware

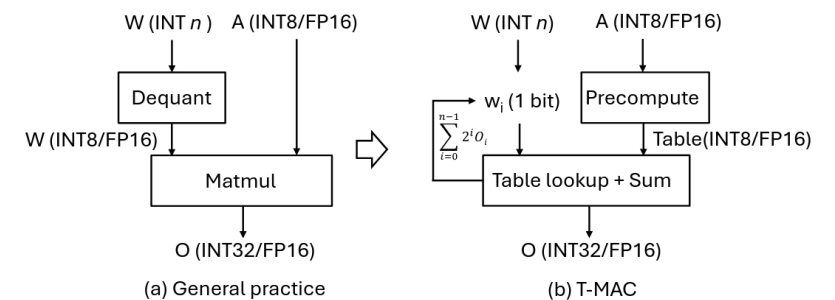
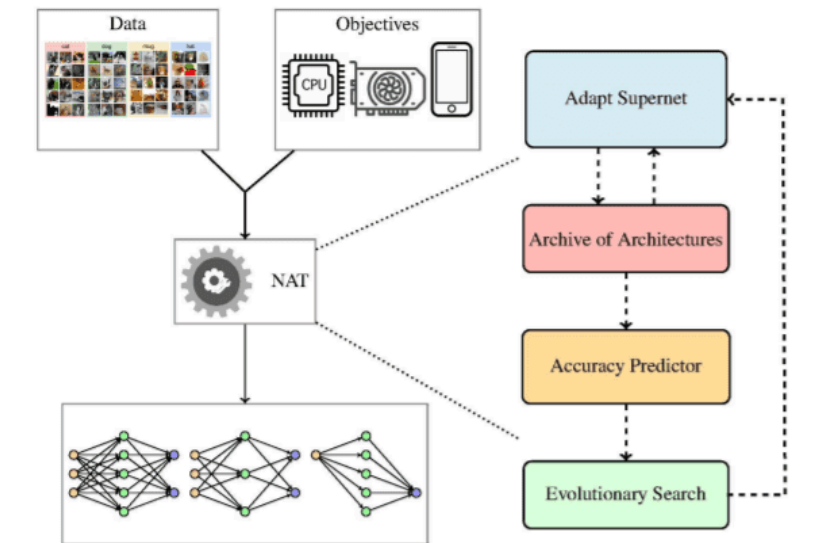
Energy Consumption

Efficient algorithms-architectures



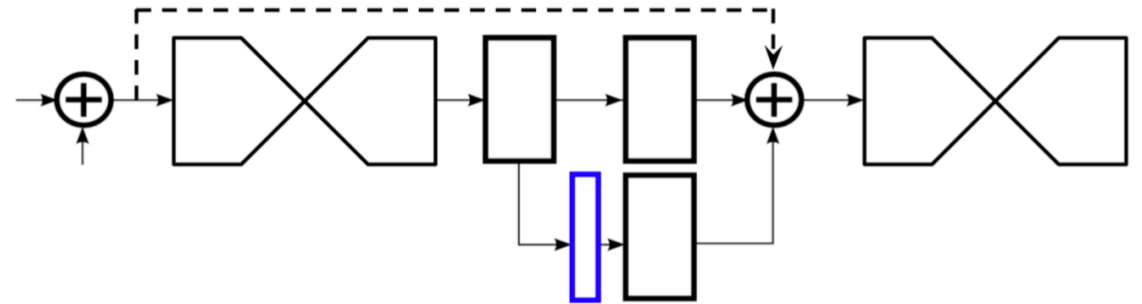
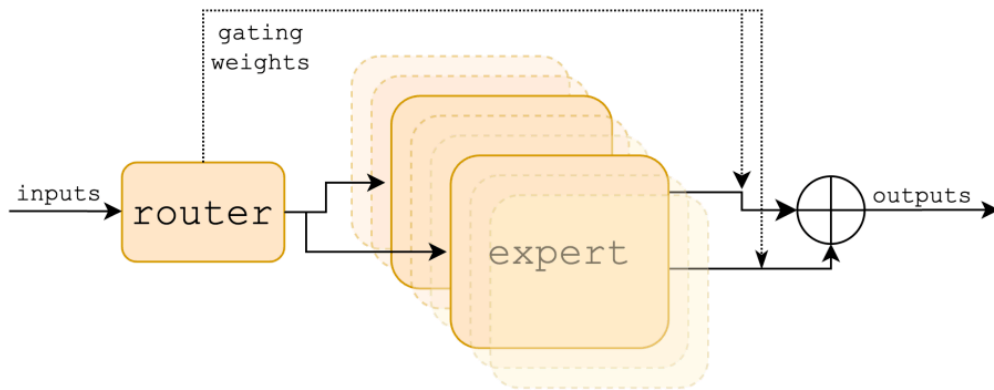
Model	Multi-Layer Perceptron (MLP)	Kolmogorov-Arnold Network (KAN)
Model (Shallow)	<p>(a)</p> <p>fixed activation functions on nodes</p> <p>learnable weights on edges</p>	<p>(b)</p> <p>learnable activation functions on edges</p> <p>sum operation on nodes</p>

[[22](#), [23](#), [24](#), [25](#)]



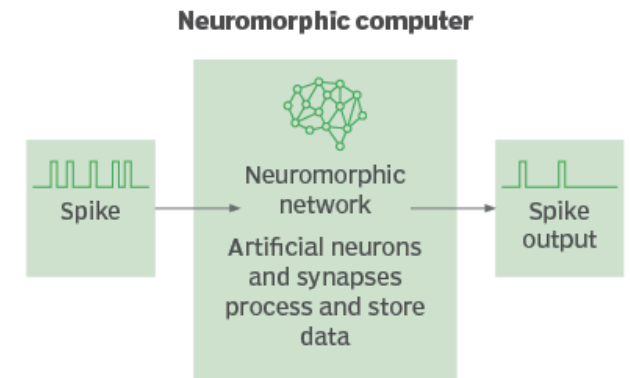
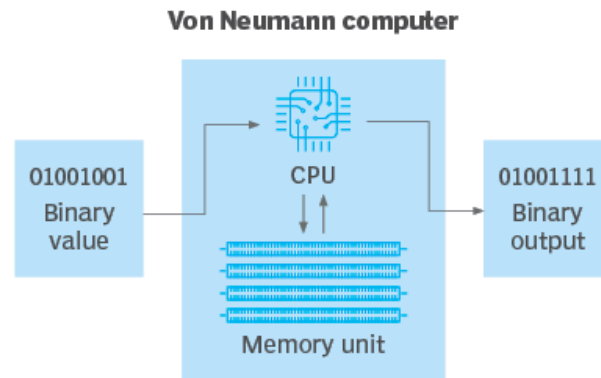
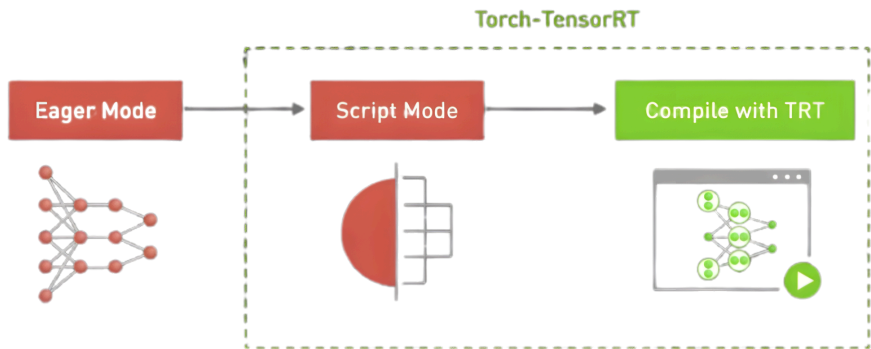
Energy Consumption

Skip mechanism, gate mechanism, early termination



Energy Consumption

Compilation, specialized hardware

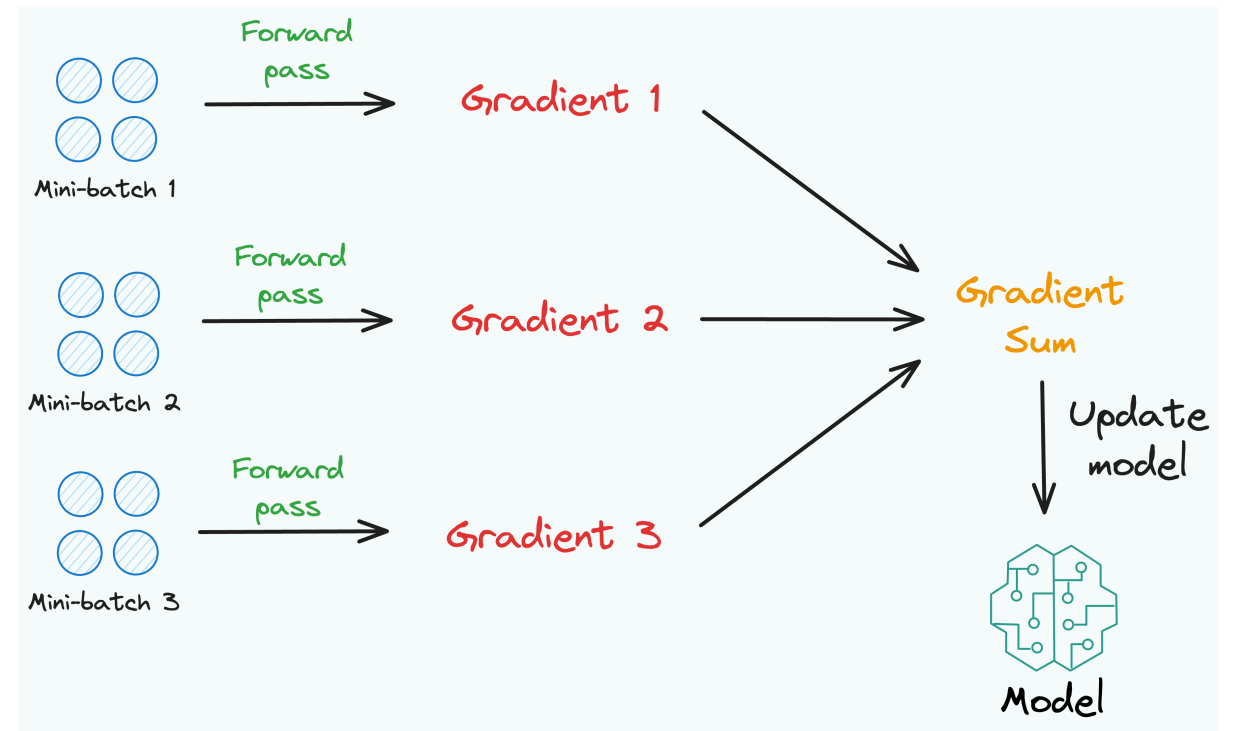
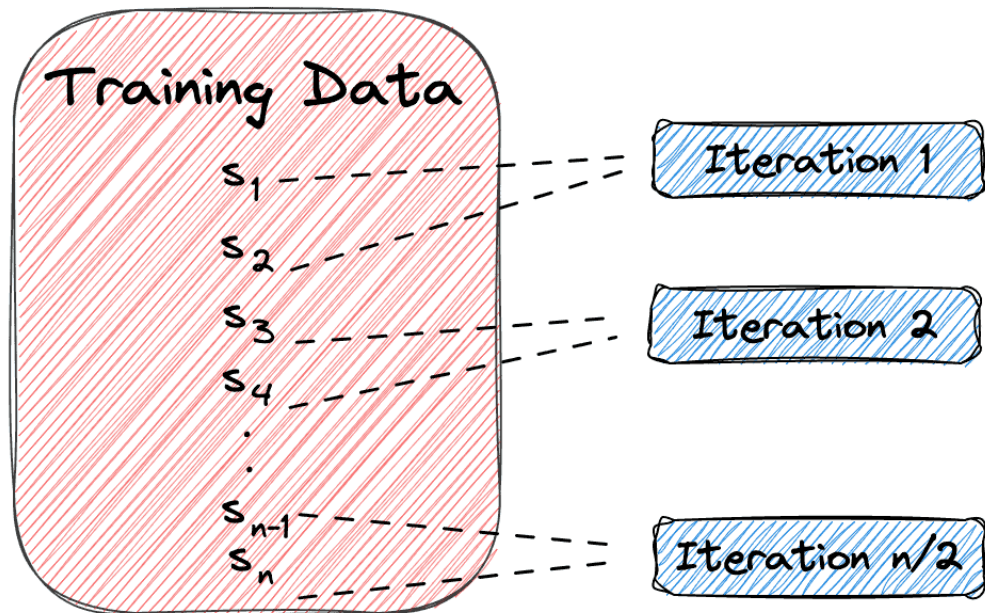


Model

Memory efficiency: batching, accumulation, weight sharing, tensor decomposition, pruning, quantization, distillation, sparse representation

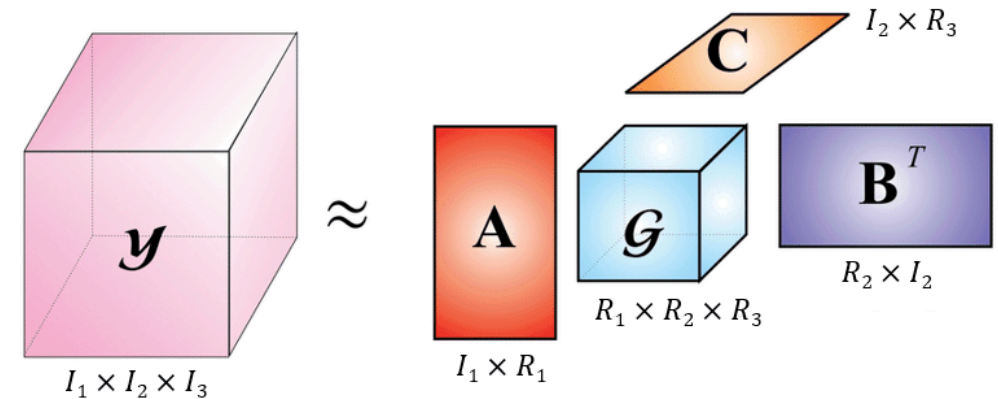
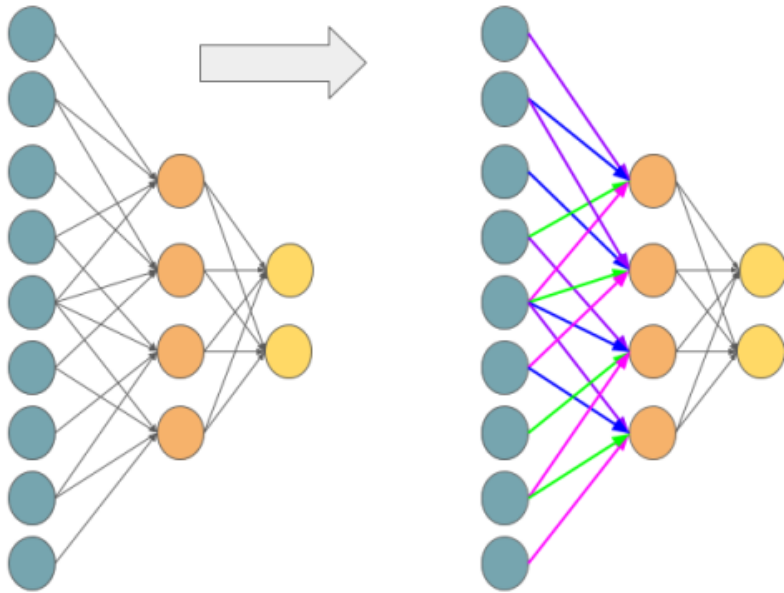
Memory Efficiency

Batching, accumulation



Memory Efficiency

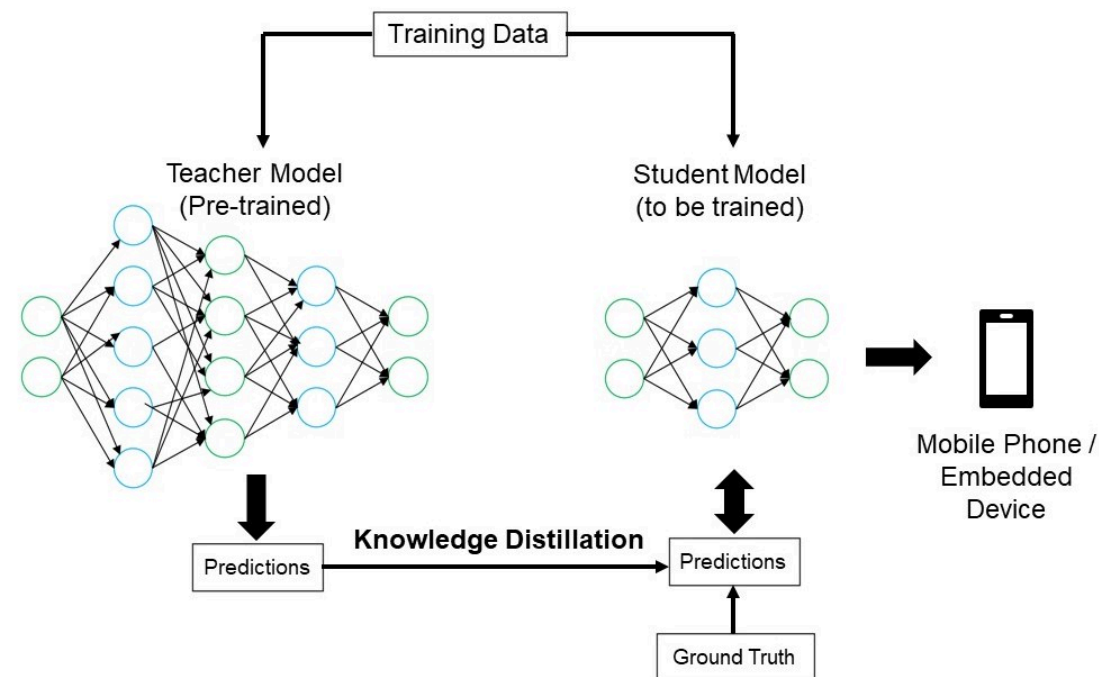
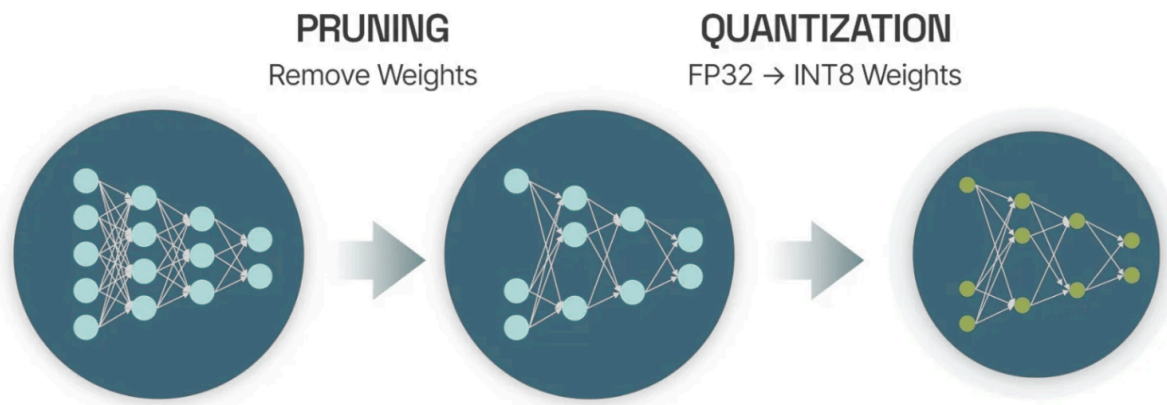
Weight sharing, tensor decomposition



[[28](#), [29](#)]

Memory Efficiency

Pruning, quantization, distillation, sparse representation



<i>s p a r s e</i>							DENSE						
	7					6	0	7	0	0	0	0	6
	7	6	3			4	0	7	6	3	0	4	0
	4	3					0	4	3	0	0	0	0
4	2						4	2	0	0	0	0	0
							0	0	0	0	3	2	4

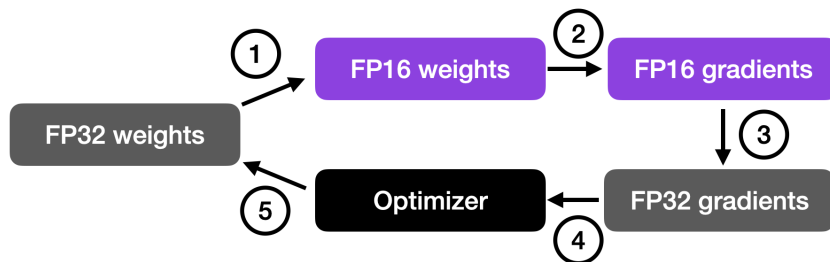
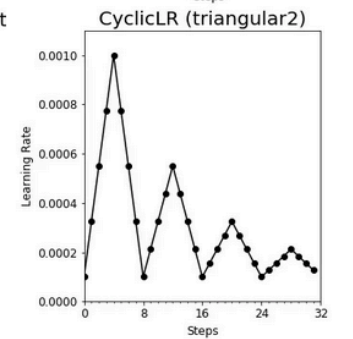
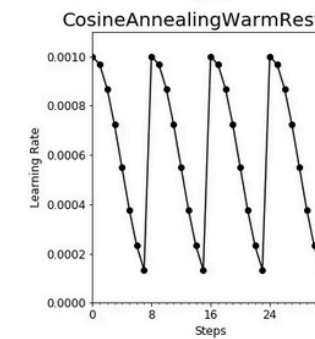
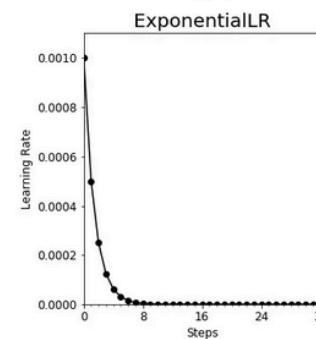
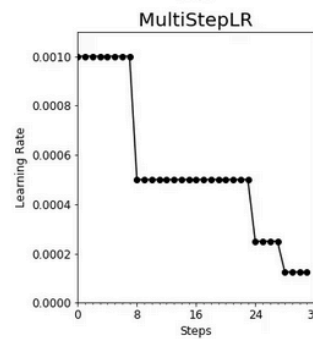
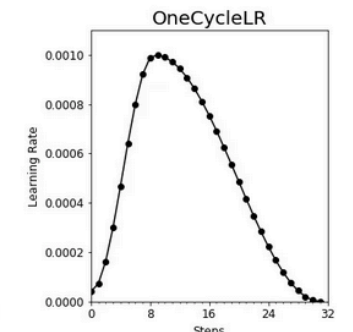
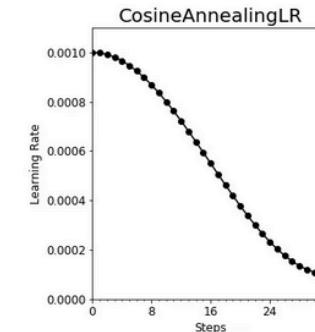
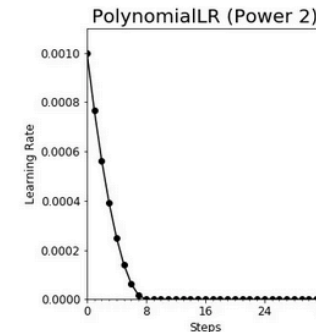
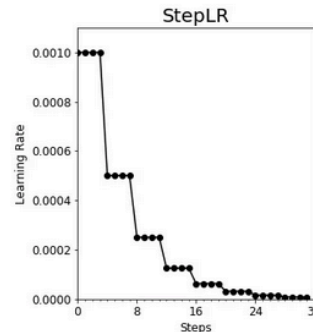
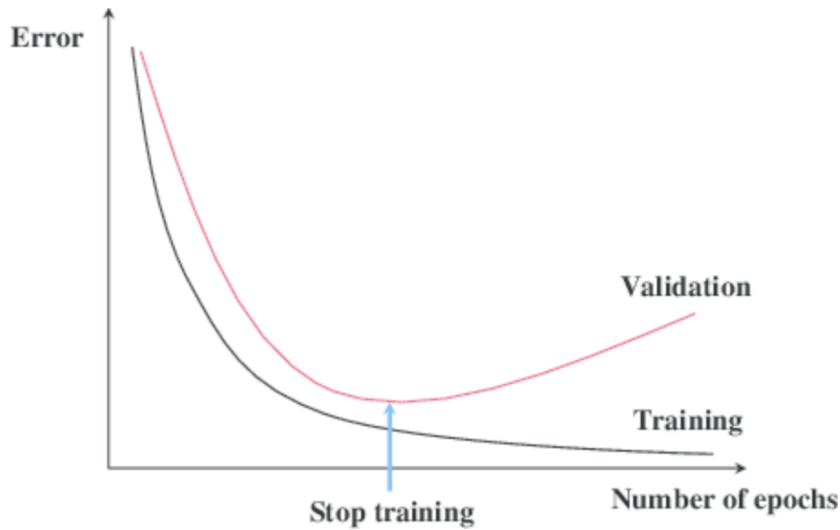
[30, 31]

Model

Training overhead: early stopping, mixed precision, scheduling, distributed learning, dynamic pruning, low-rank adaptation

Training Overhead

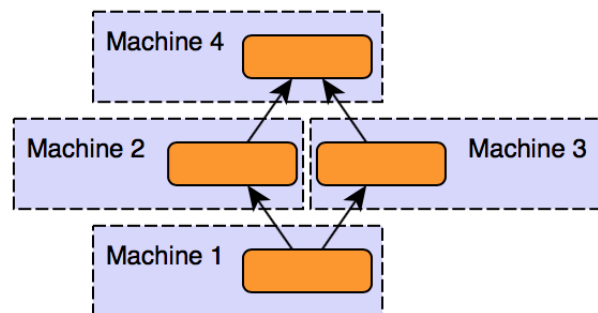
Early stopping, scheduling, mixed precision



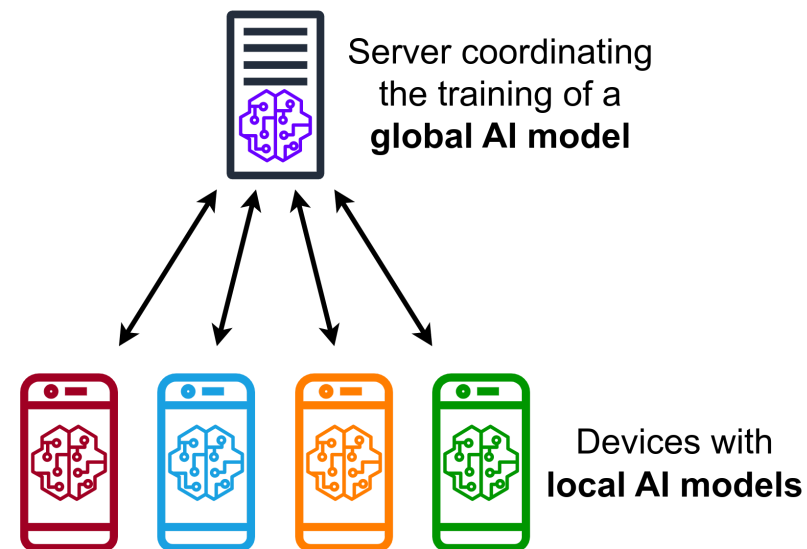
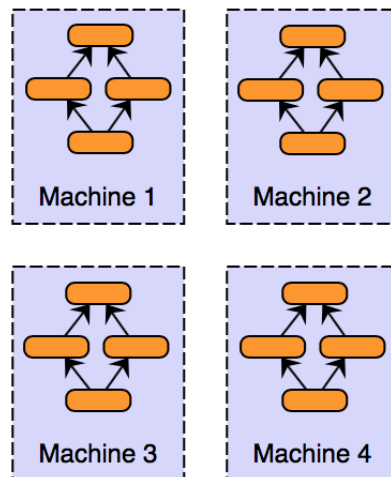
Training Overhead

Distributed learning

Model Parallelism

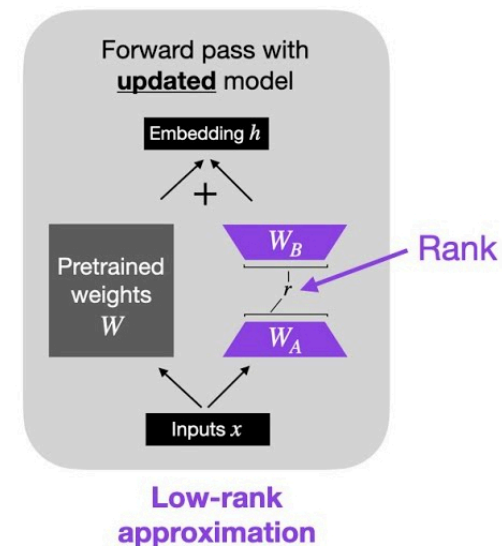
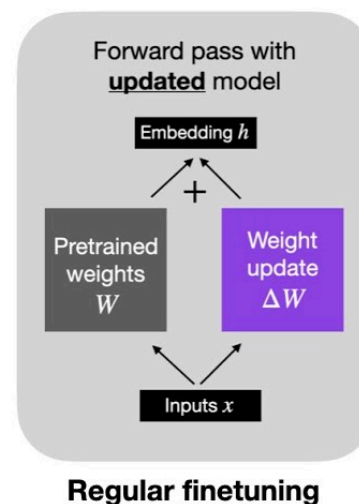
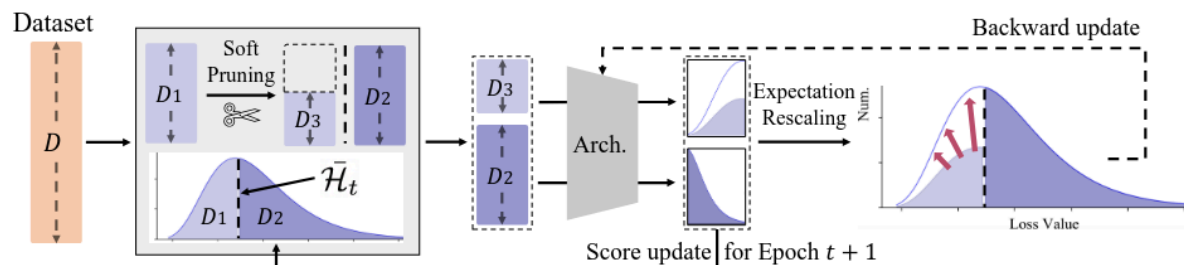


Data Parallelism



Training Overhead

Dynamic pruning, low-rank adaptation



[33, 34]

MISCELLANEOUS

Miscellaneous

- **Tools:** data management, annotation, metadata tracker, base stack
- **Resources:** datasets, models, infrastructures
- **Hardware:** SBC, sensors, neuromorphic, deployment
- **Advances:** AI and robotics

Miscellaneous

Tools: data management, annotation, metadata tracker,
base stack

Data Management



DVC: Connect to versioned data sources and code with pipelines, track experiments, register models – all based on GitOps principles

```
dvc init
dvc remote add -d myremote /tmp/dvcstore

dvc add data/data.json
dvc push

git add data/data.json.dvc data/.gitignore
git commit -m "Add raw data"
git push

git pull
dvc pull
```

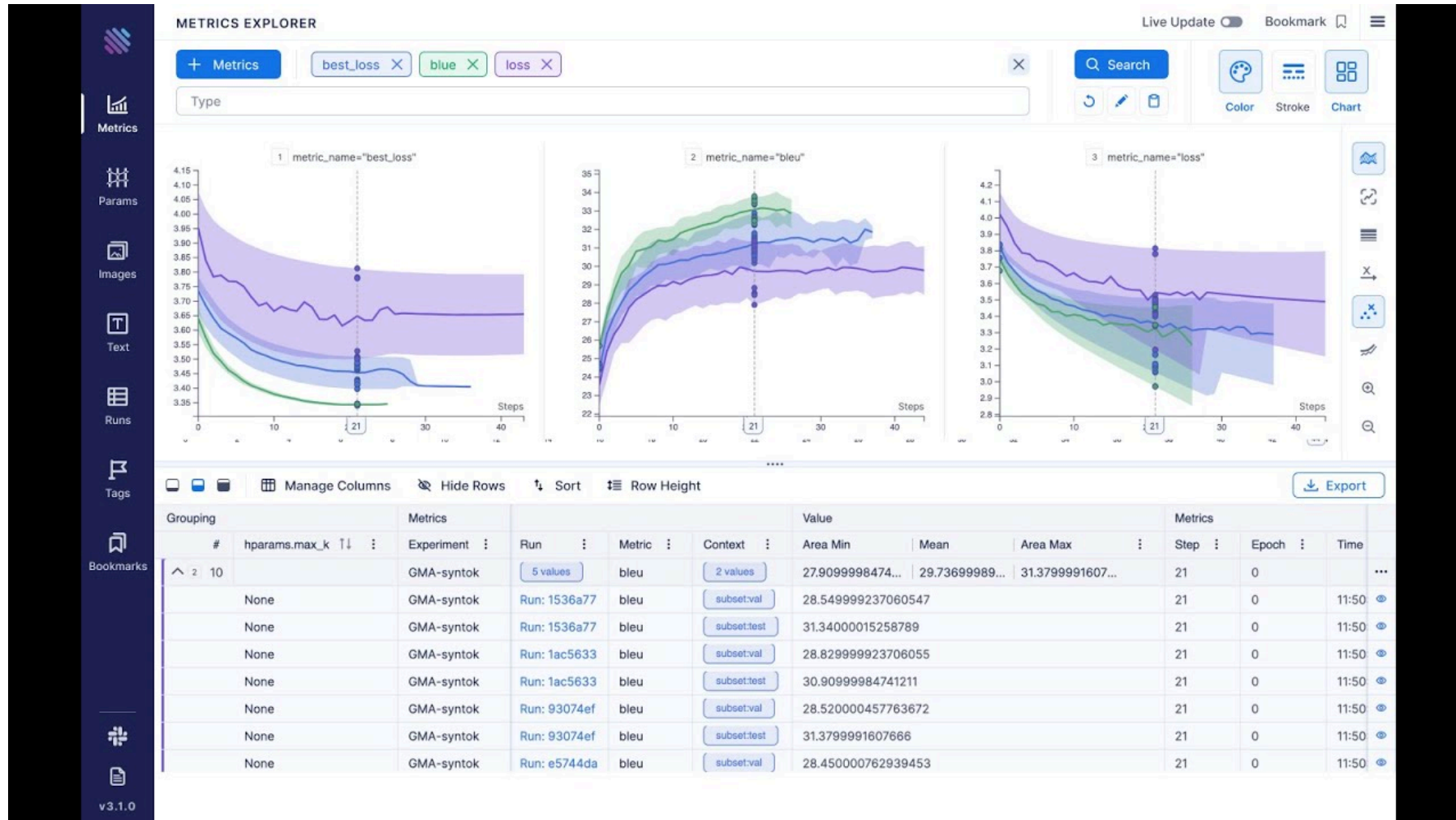
Data Annotation

Label Studio: open source data labeling tool that supports multiple projects, users, and data types

The screenshot displays the Label Studio web interface. At the top, the breadcrumb navigation shows 'Projects / Video Example - 25 / Labeling'. The main workspace features a video player with a running scene. Several runners are highlighted with colored bounding boxes: a purple box for a runner on the left, a blue box for a runner in the center, and a green box for a runner on the right. A red timer at the bottom of the video shows '31:01.8'. To the left of the video is an 'Outliner' panel with a tree view of labels: 'Runner #3', 'Runner #1' (selected), 'Runner #2', 'Runner #5', 'Runner #4', 'Runner #6', and several 'Other' labels. Below the video is a timeline with a search bar and playback controls. At the bottom, a legend identifies the runners: 'Runner #3' (purple), 'Runner #1' (orange), 'Runner #2' (pink), 'Runner #5' (blue), 'Runner #4' (cyan), and 'Runner #6' (green). A status bar at the very bottom shows a sequence of labels: 'blank 1 | Man 2 | Woman 3 | Other 4 | Runner #1 5 | Runner #2 6 | Runner #3 7 | Runner #4 8 | Runner #5 9 | Runner #6 0'.

Metadata Tracker

Aimstack: Tracking and visualizing experiments



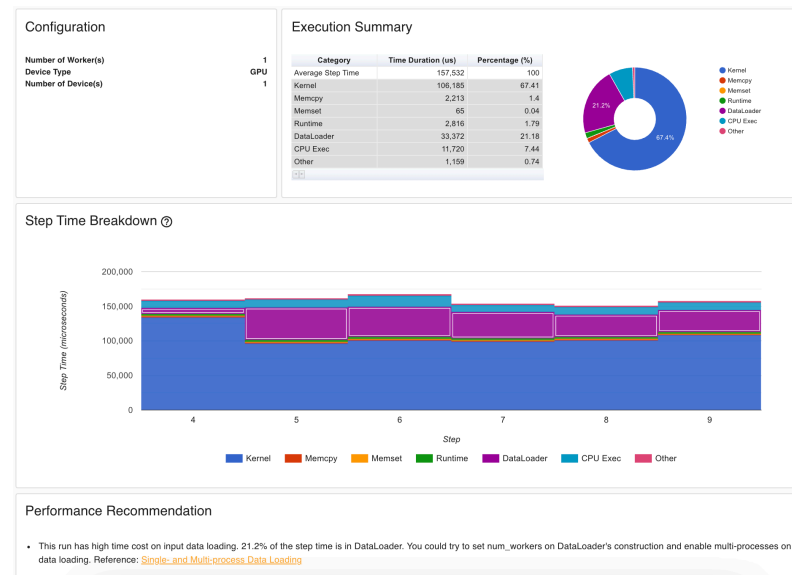
Tooling



Base stack: [pytorch](#), [lightning](#), [hydra](#) ([template](#))

Additional packages: [torchmetrics](#), [optuna](#), [captum](#), [mapie](#), [nni](#)
[grafana](#), [prometheus](#), [alumentations](#), [onnx](#)


[AI Lifecycle](#), [Troubleshooting AI](#)




Miscellaneous

Resources: datasets, models, infrastructures


Hugging Face


 **Hugging Face**


[Models](#) [Datasets](#) [Spaces](#) [Posts](#) [Docs](#) [Solutions](#) [Pricing](#) [Log In](#) [Sign Up](#)


 **LeRobot** Community
<https://github.com/huggingface/lerobot>

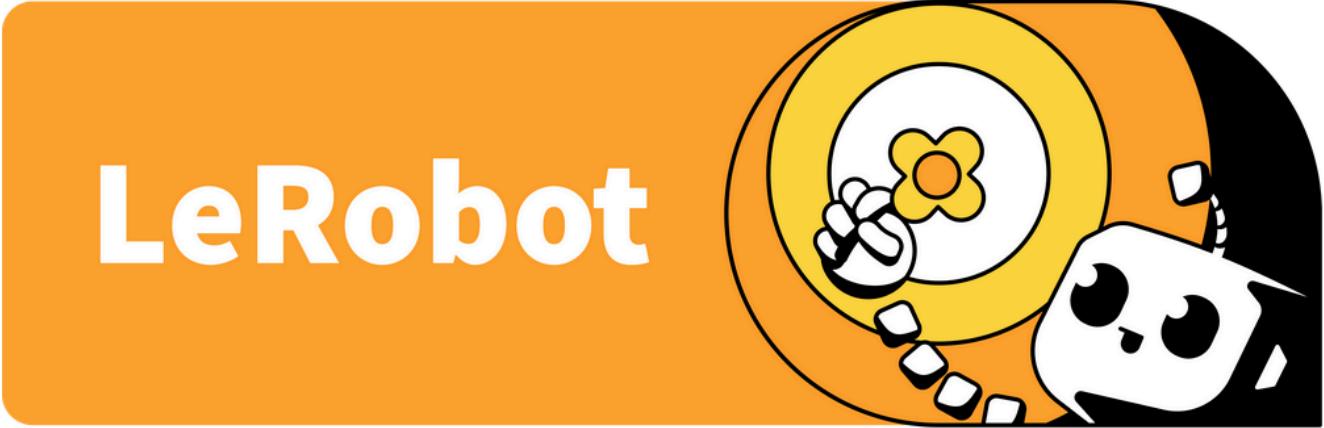
[Request to join this org](#)

 **AI & ML interests**
State-of-the-art Machine Learning for real-world robotics

 **Team members** 11



 **Organization Card** Community [About org cards](#)



Also: [Papers with Code](#)

HPC/AI Resources

Hybrid Cluster (local)

- Access: [Email](#)
- [Documentation](#)

Grid'5000 (national)

- Access: [Form](#)
- [Documentation](#)

Jean Zay (national)

- Access: [Form](#) (academic and industry)
- [Documentation](#)

Usage

- Jupyter notebooks via a web interface or traditional command-line tools
- Submit jobs using SLURM, specifying required resources



Slurm

Manages a queue and launches jobs submitted on dedicated compute servers when the requested resources are available.

```
#!/bin/bash

#SBATCH --nodes=1
#SBATCH --ntasks-per-node=1
#SBATCH --time=00:05:00
#SBATCH --job-name=GPU
#SBATCH --mem=1G
#SBATCH --gres=gpu:1

./program
```

> sbatch job.slurm

Miscellaneous

Hardware: SBC, sensors, neuromorphic, deployment

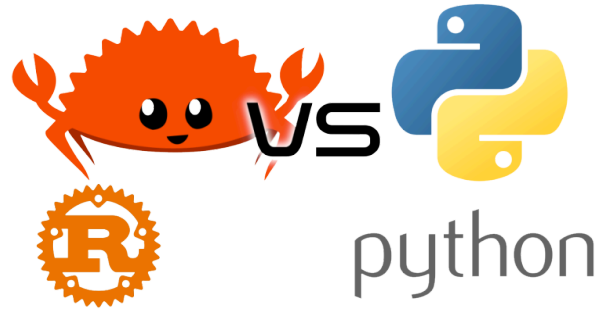
Embedded AI

- SBC: Nvidia Jetson, Raspberry Pi
- SoC: Qualcomm Snapdragon 8cx Gen 3, Google Tensor, Apple M-series
- MPU: Intel Core i7/i9/Xeon (AVX-512, DL Boost), AMD Ryzen 9 (Threadripper)
- ASIC/NPU/LPU: ARM Cortex Ethos, Apple Neural Engine, Samsung Exynos, Huawei Ascend, Intel Gaudi, Groq, cerebras
- MCU: STMicroelectronics STM32 AI, Arduino
- FPGA: AMD/Xilinx Kria, Versal AI
- Sensor: Luxonis OAK, event-based cameras
- Neuromorphic: Intel Loihi
- Co-inference

Libraries: [EdgeImpulse](#), [NanoEdgeAIStudio](#)



Deployment



- Python overhead can seriously hurt performance
- The GIL is a notorious source of headaches
- Make serverless inference possible with lightweight binaries
- Remove Python from production workloads

[Why Rust](#): safety and speed, borrow checker, scoped resource management, error handling, wasm, tooling, ...

Libraries: [Candle](#), [Luminal](#)

Deployment

```
import torch
import torchvision

model = torchvision.models.resnet18()
traced_script_module = torch.jit.trace(model, torch.rand(1, 3, 224, 224))
traced_script_module.save("traced_resnet_model.pt")
```

```
#include <torch/script.h>

#include <iostream>
#include <memory>

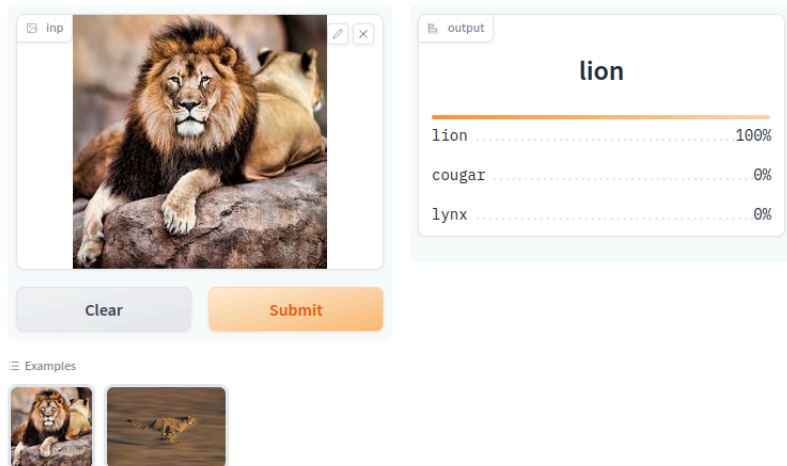
torch::jit::script::Module module;
module = torch::jit::load(argv[1]);

std::vector<torch::jit::IValue> inputs;
inputs.push_back(torch::ones({1, 3, 224, 224}));

at::Tensor output = module.forward(inputs).toTensor();
std::cout << output.slice(/*dim=*/1, /*start=*/0, /*end=*/5) << '\n';
```

Deployment

[Gradio](#): easily create interactive web-based user interfaces for ML models



api_name: `/predict`

```
from gradio_client import Client

client = Client("https://gradio-pytorch-image-classifier.hf.space/")
result = client.predict(
    "https://raw.githubusercontent.com/gradio-app/gradio/main/test/test_files/bus.png", # str (filepath or URL to image
    in 'inp' Image component
    api_name="/predict"
)
print(result)
```

copy

• Return Type(s)

```
# str representing output in 'output' Label component
```

```
def predict(inp):
    pass
```

```
gr.Interface(fn=predict, inputs=gr.Image(type="pil"), outputs=gr.Label(num_top_classes=3),
            examples=["lion.jpg", "cheetah.jpg"]).launch()
```

Also: [FastAPI](#), [TorchServe](#)

Miscellaneous

Advances: AI and robotics

GRUtopia

Dream General Robots in a City at Scale

100K+ Scenes, 89 Scene Categories

Generated Tasks

Hospital Canteen Navigation

Office Library Dialogue

Sure.

Could you please bring me a spoon?

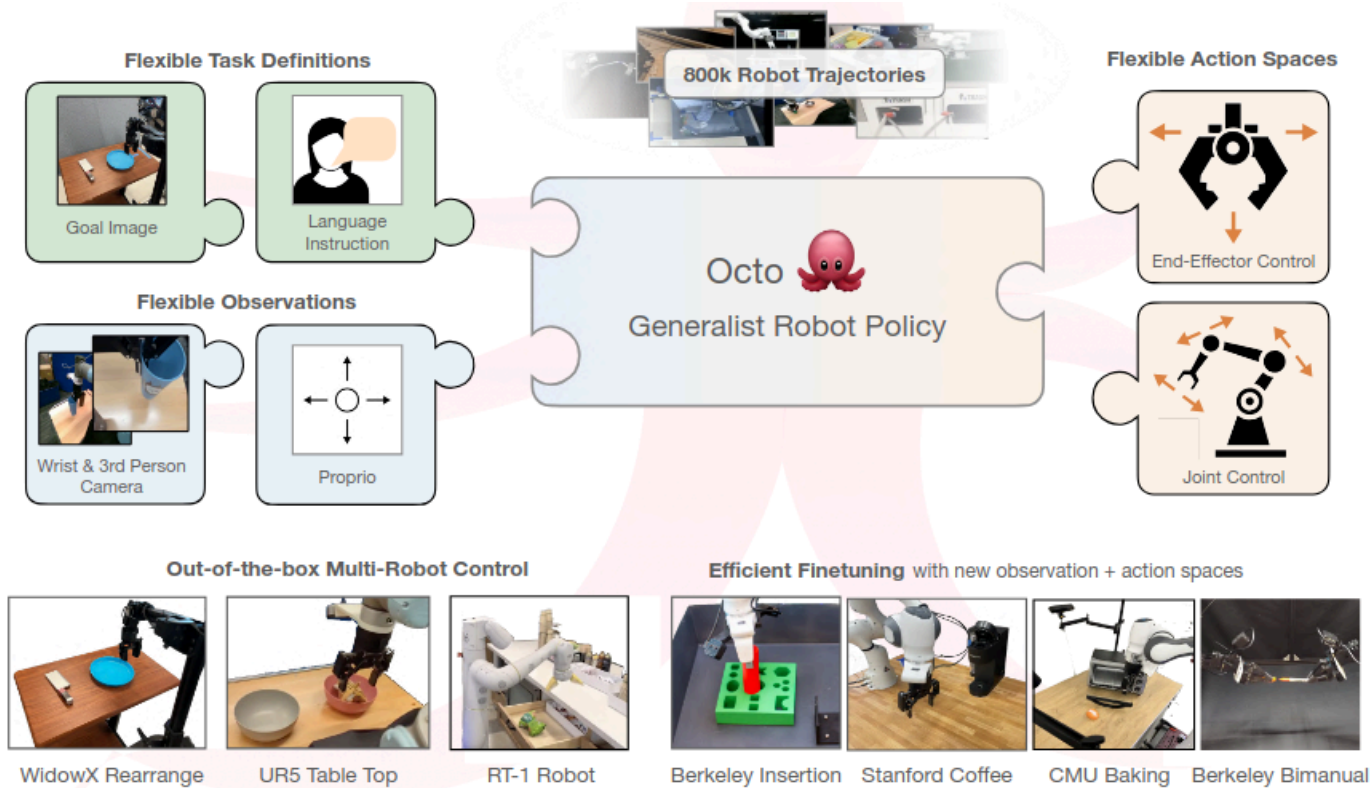
Unitree H1

Supermarket School Manipulation

Interactive and Finely Annotated Social Interaction with NPC

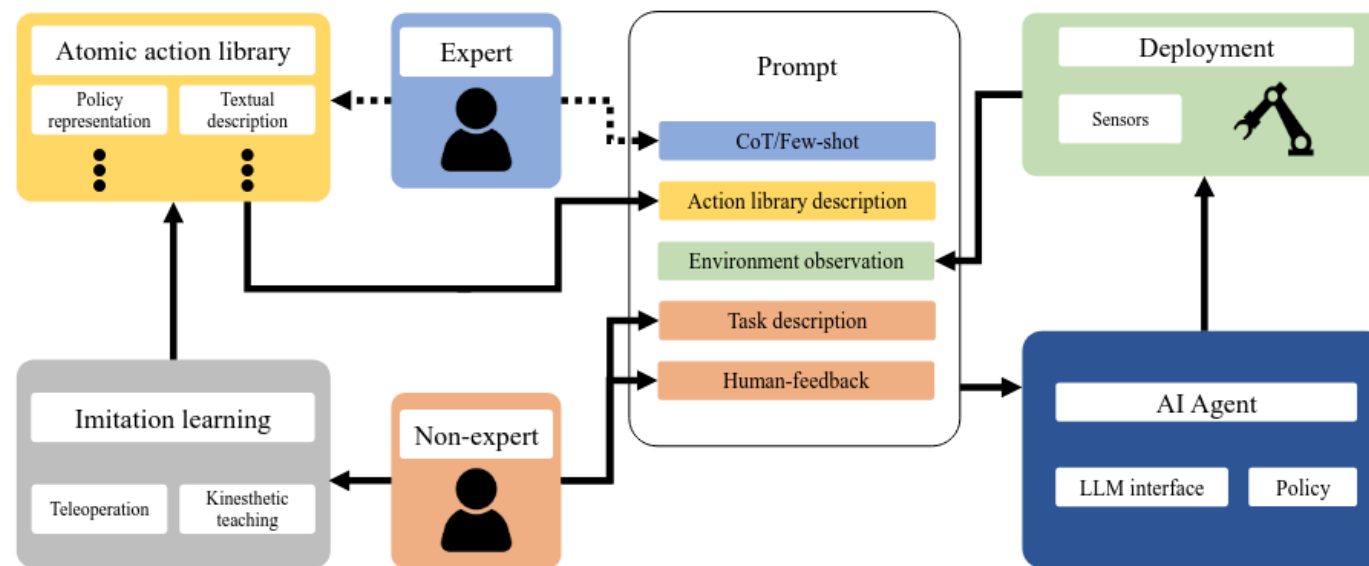
Octo

An Open-Source Generalist Robot Policy



ROS-LLM

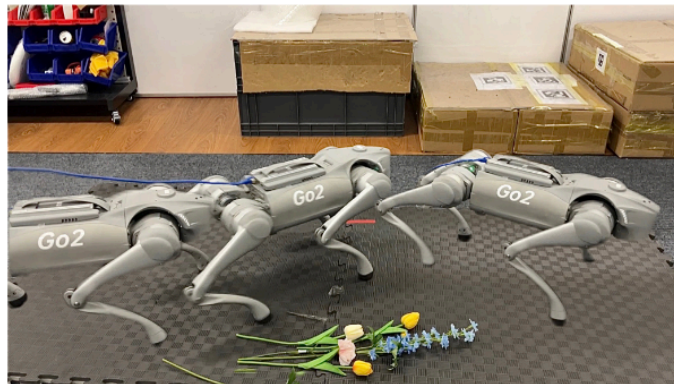
A ROS framework for embodied AI with task feedback and structured reasoning



Make An Agent

Generalizable Policy Network Generator with Behavior-Prompted Diffusion

Making agile turns to avoid stepping on a bouquet while moving across a mat

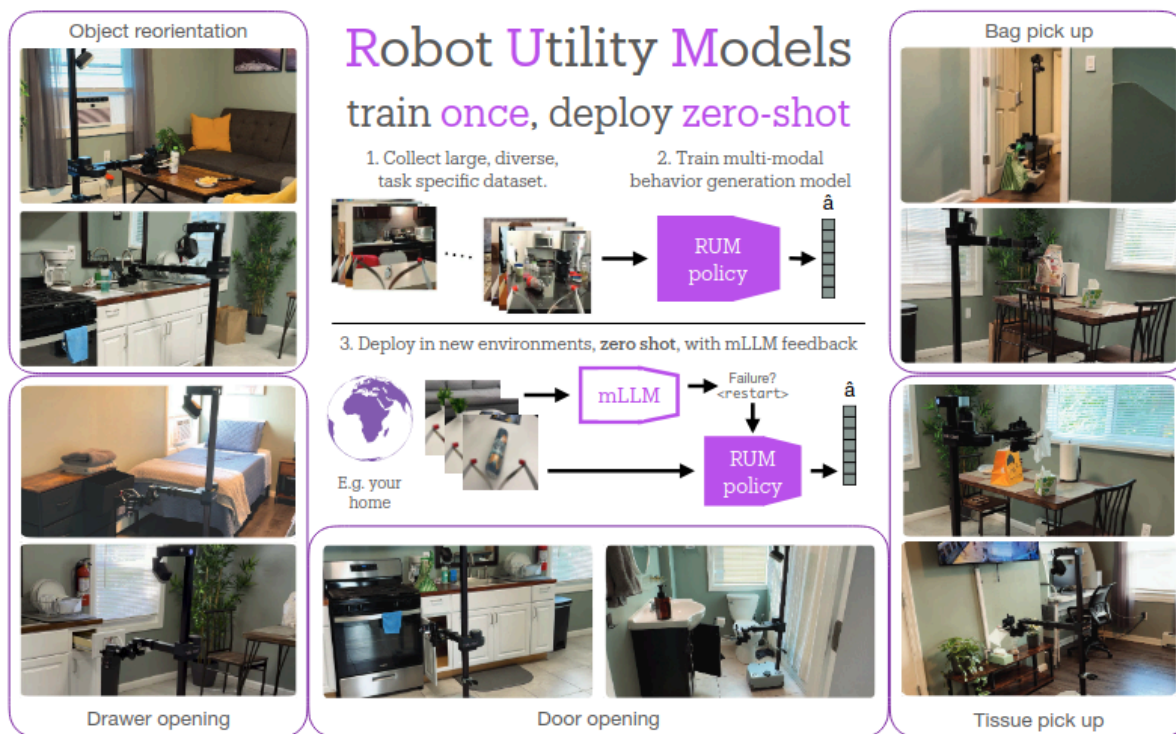


Navigating to circumvent the goal and ball while swiftly moving backward.



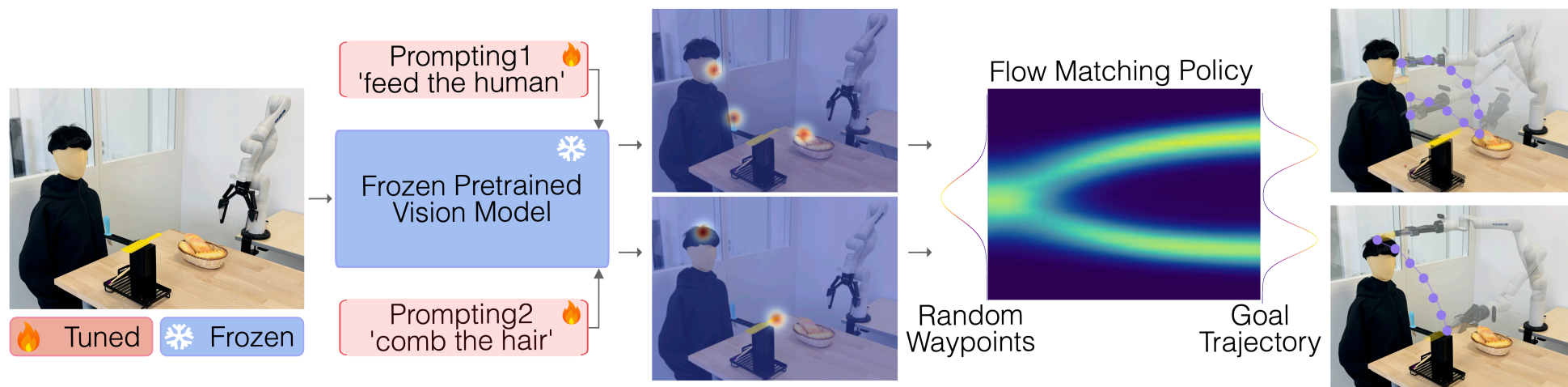
Robot Utility Models

General Policies for Zero-Shot Deployment in New Environments

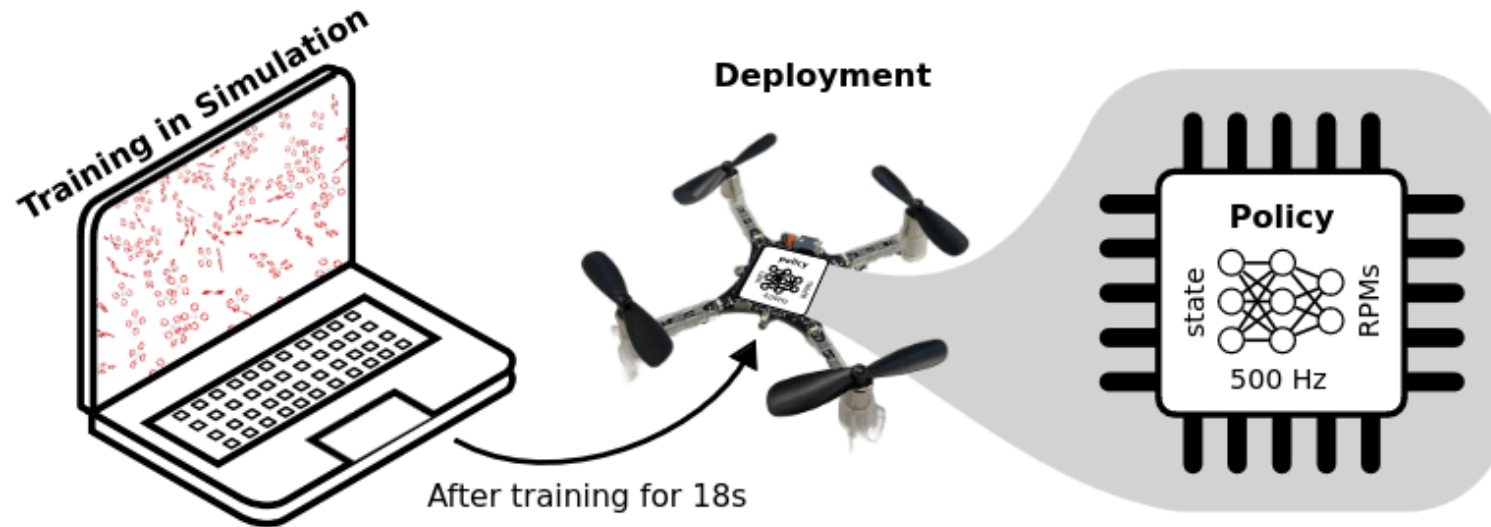


Flow Matching Policy

Affordance-based Robot Manipulation with Flow Matching

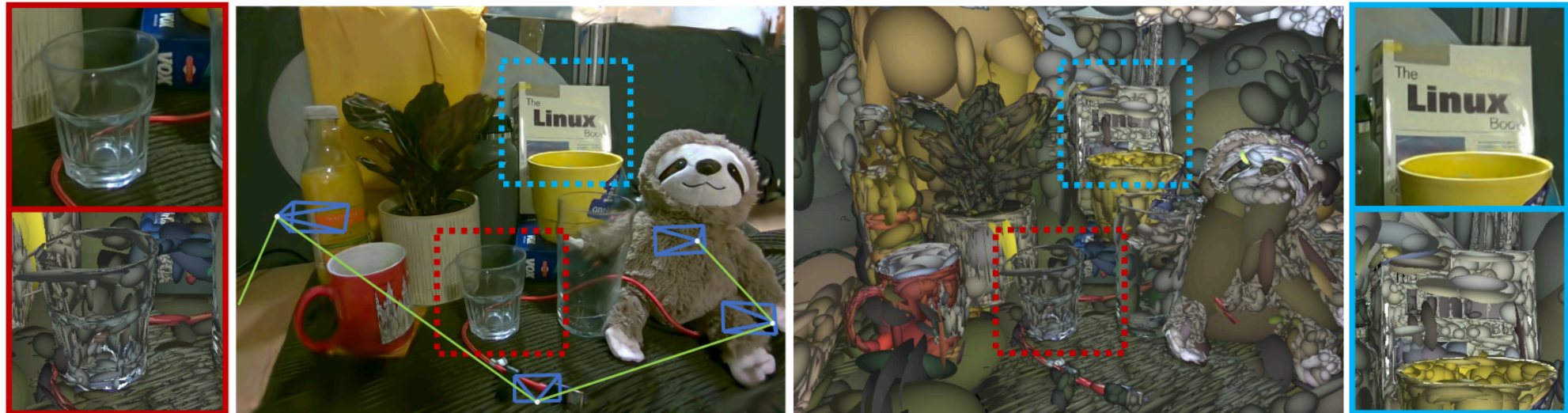


Learning to Fly in Seconds



[41]

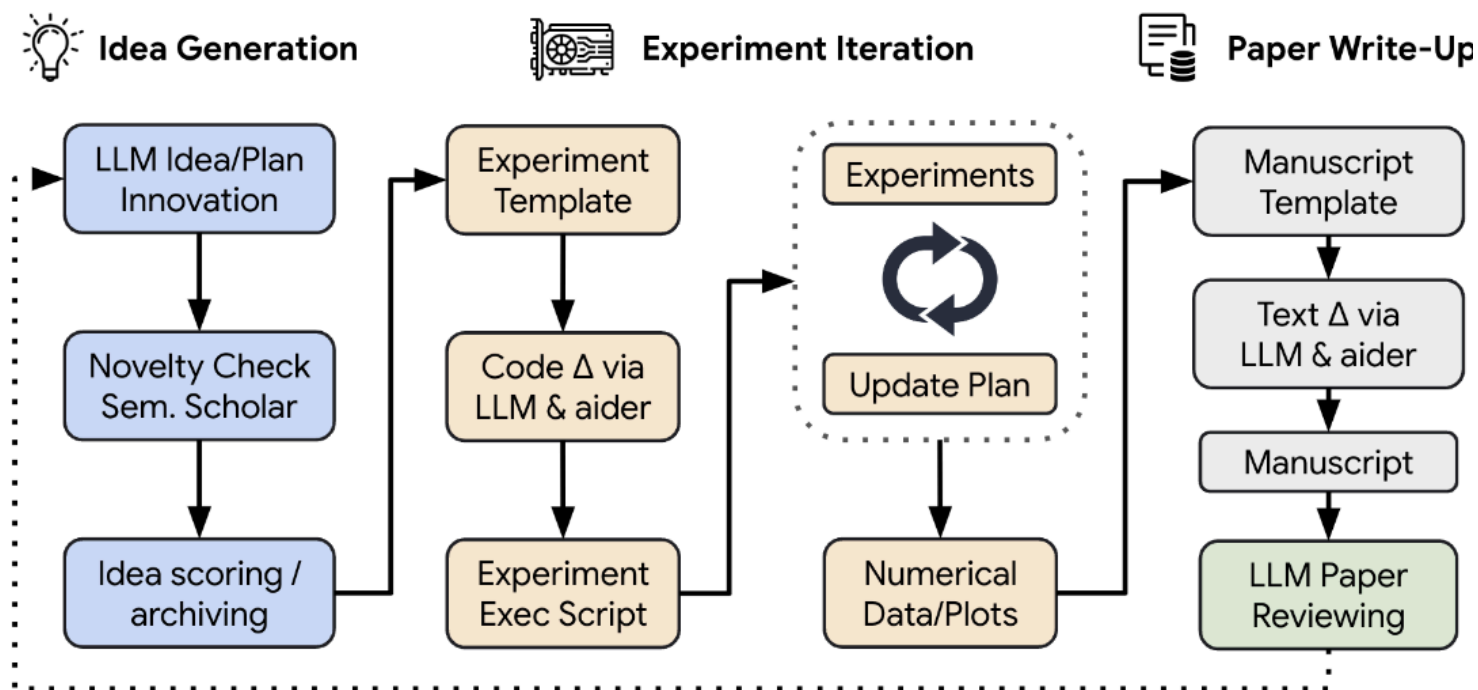
Gaussian Splatting SLAM



[42]

The AI Scientist

Towards Fully Automated Open-Ended Scientific Discovery



Conclusion

- Open source AI with HuggingFace
- Physics-informed AI and Multimodal Foundation Model
- [Catastrophic forgetting](#), [benchmarking](#)
- [TinyML](#) Foundation

Responsible and Sustainable AI



FIDLE

Bases,
Concepts
et Enjeux

L'IA
comme
un outil,


Acteur
de l'IA

- 


History and Fundamental Concepts
- 

Data, models and representation's hell
Data and models
- 

Demonstration Illustration
LLM / Text to Image
- 

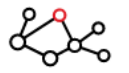
AI, Law, Society and Ethics
- 


Mathematics, gradients everywhere !
- 


Learning methodology
- 


Convolutional models
CNN
- 

Sparse (text) and sequences data
Embedding, RNN
- 


«Attention is All You Need»
Transformers
- 


Graph Neural Network
GNN
- 


Autoencoder networks
AE
- 

Variational Autoencoder
VAE
- 

Generative Adversarial Networks
GAN
- 

Diffusion Model
Text to image
- 


Deep Reinforcement Learning
RL
- 

Physics-Informed Neural Networks
PINNs
- 

Learning faster and cheaper
Eco-Friendly
- 

Jean-Zay
GPU acceleration
- 

JDLs 2024
Deep Learning for Science
- 

New models
VLM, SM, Multimodal, ...
- 

Dreams come True !
Inference & Production

Any questions?