# Ray RLlib

A scalable and unified library for reinforcement learning
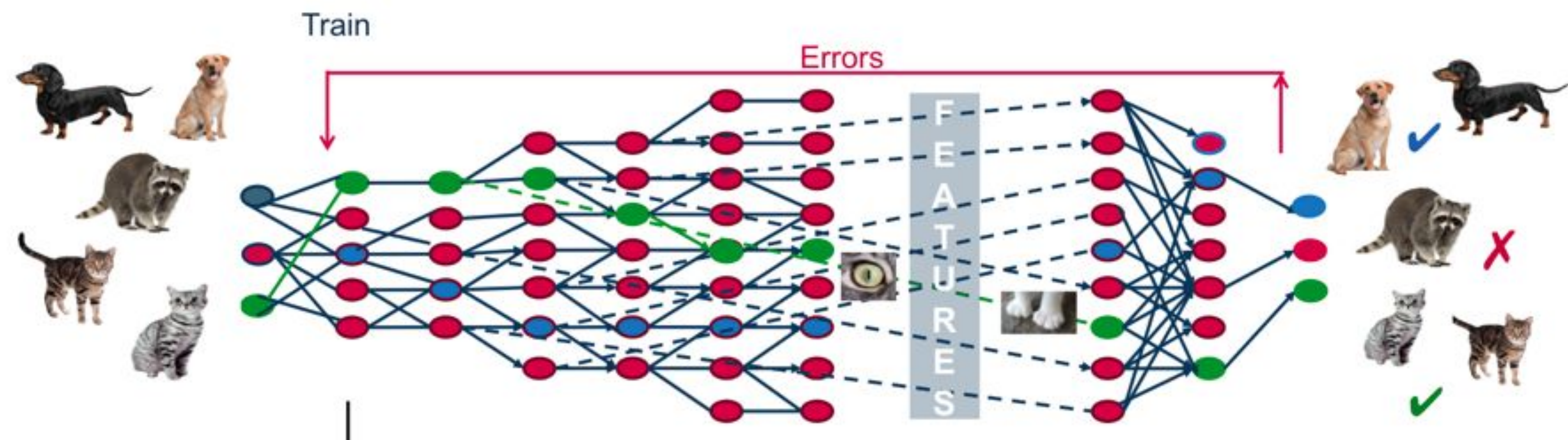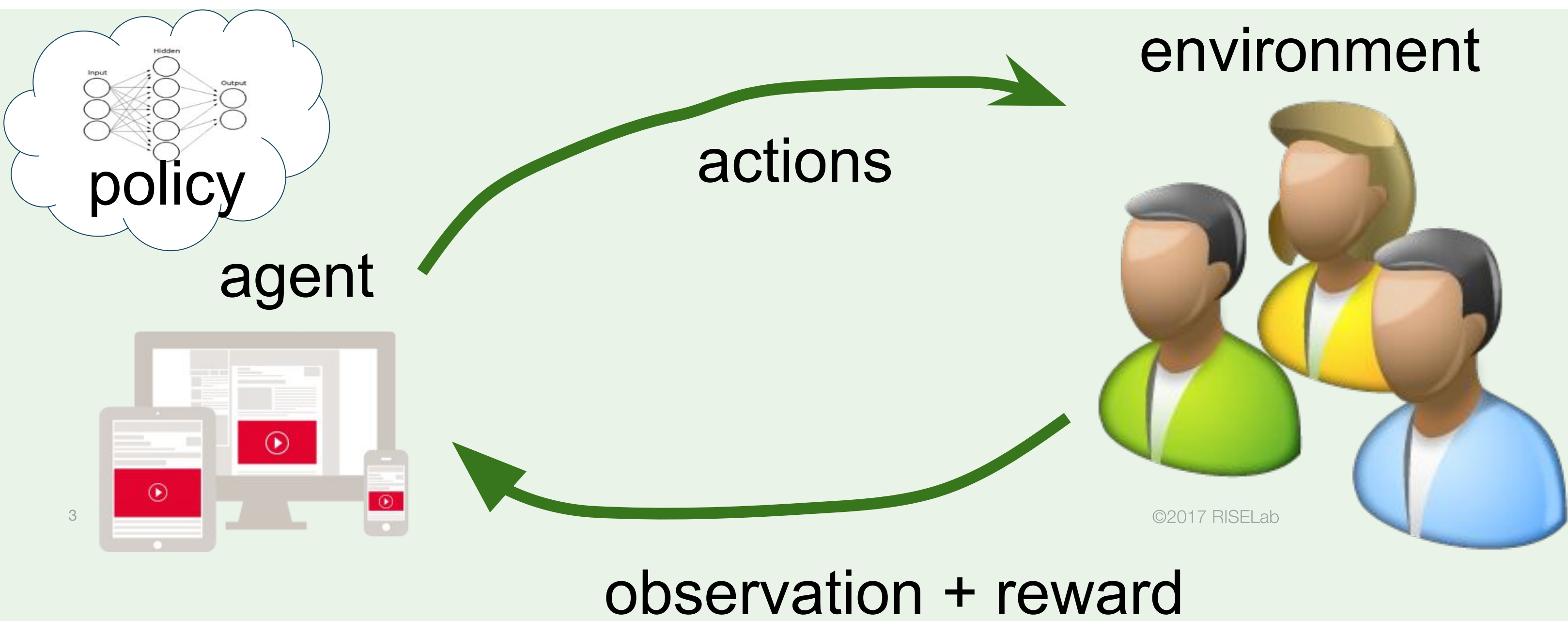**https://rllib.io**

Eric Liang

# Overview

- RLlib the open source project
- System challenges building a scalable RL library ("abstractions for RL")

# Background: What is reinforcement learning?

Supervised Learning

Reinforcement Learning

environment

actions

policy

agent

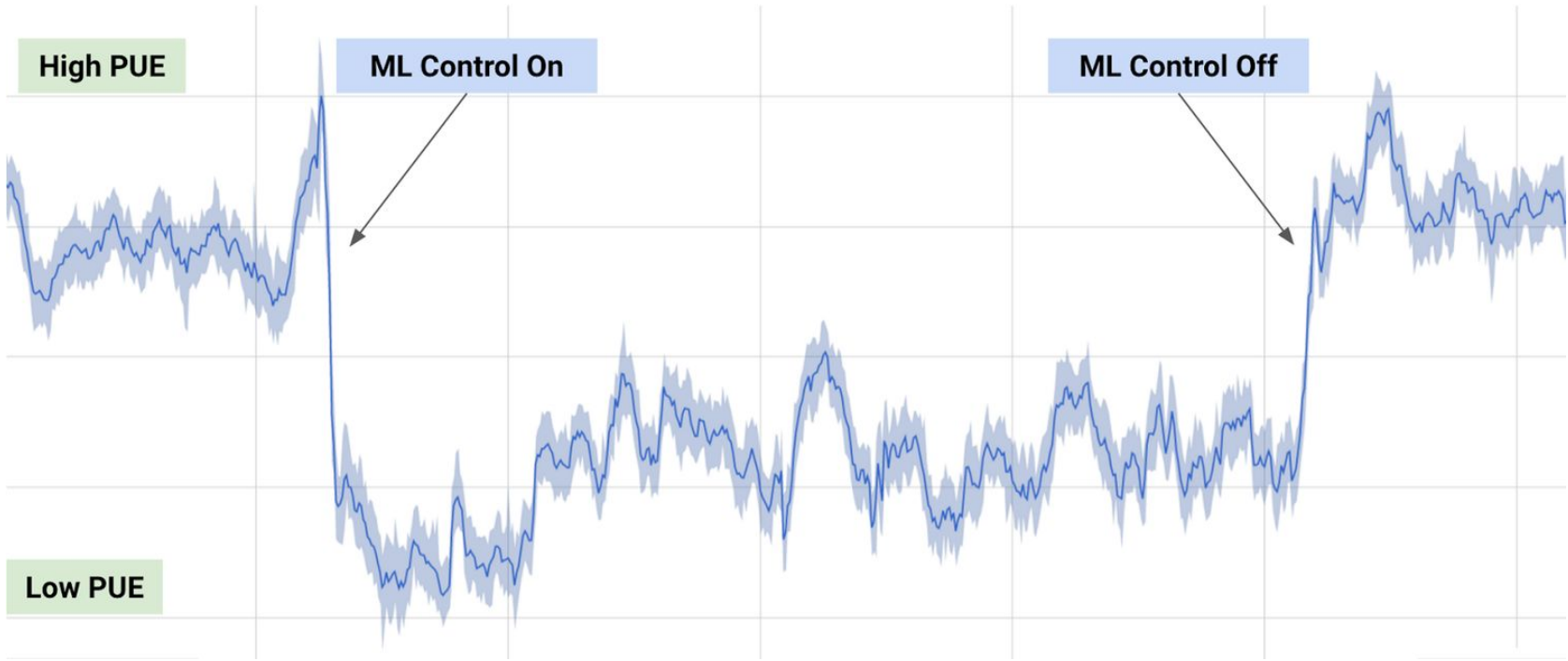observation + reward

©2017 RISELab

3

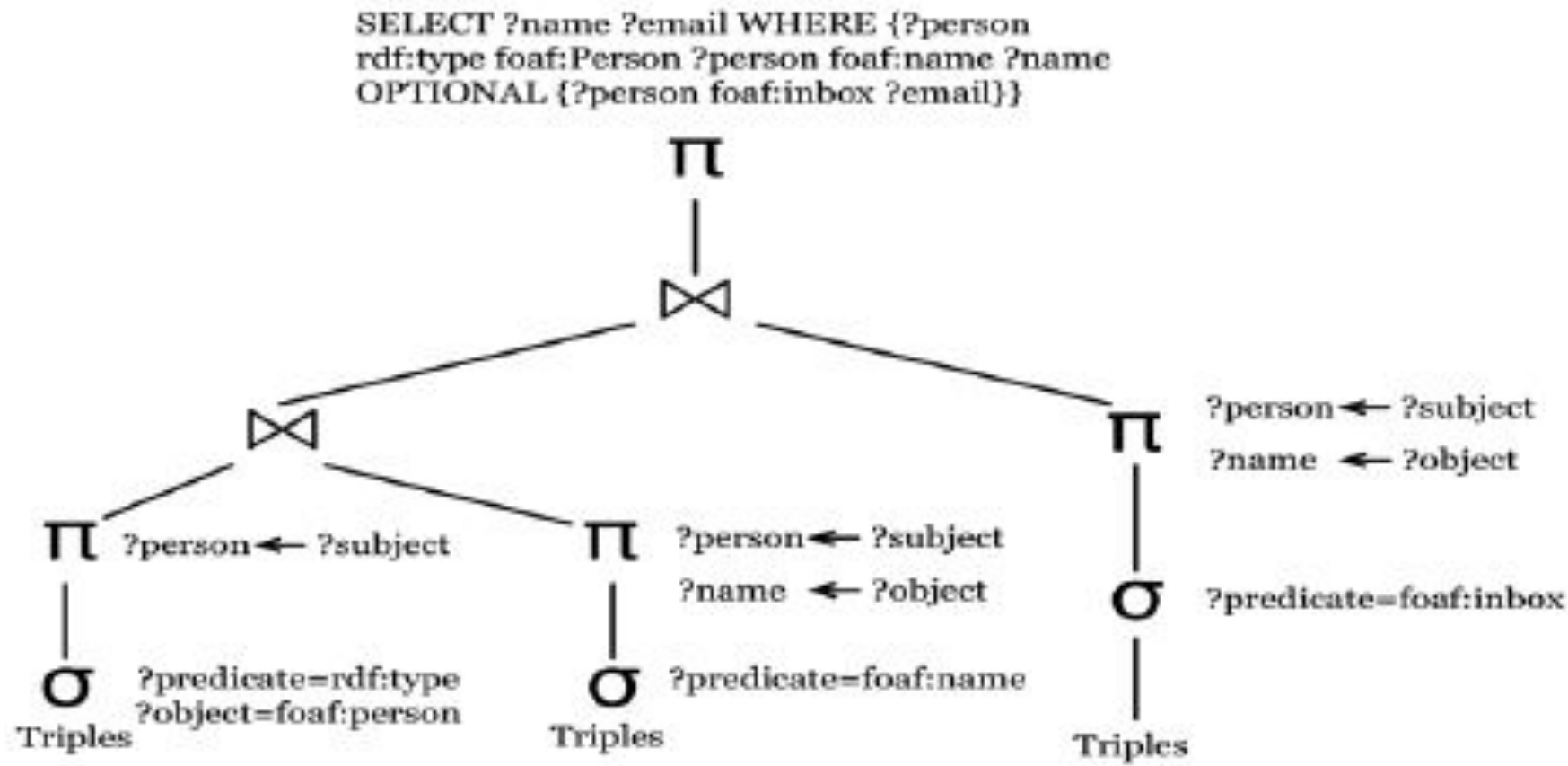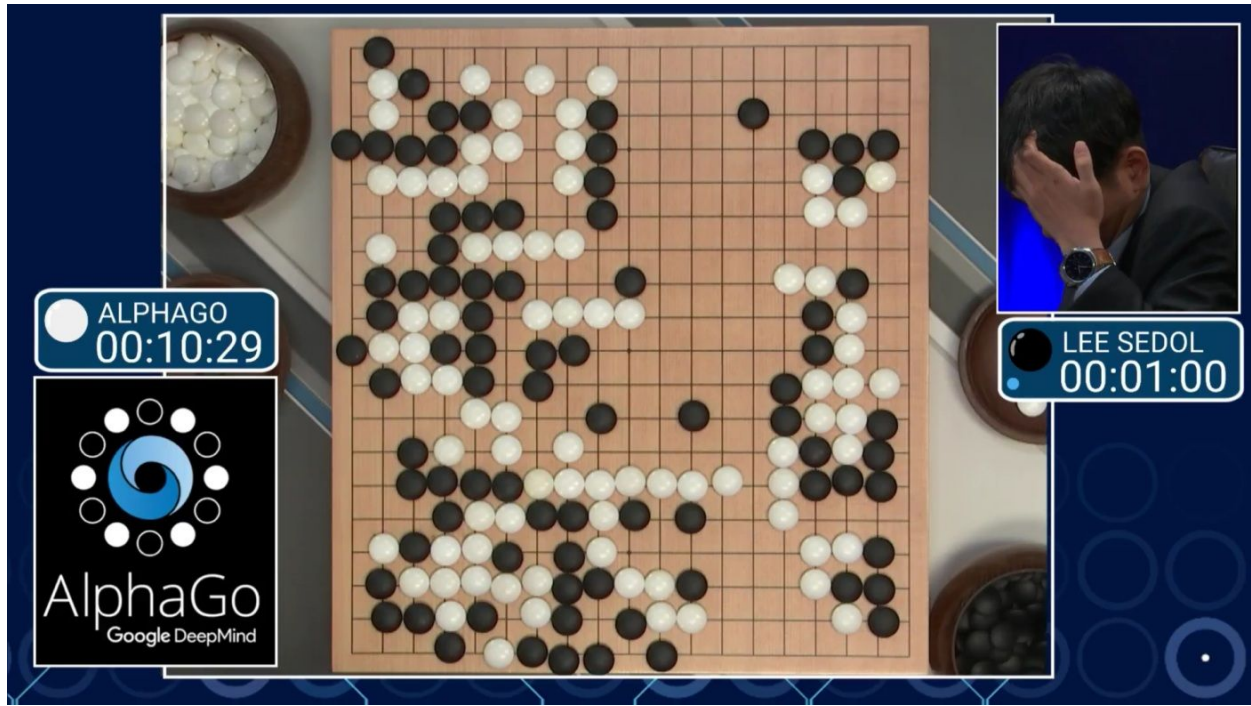# Growing number of RL applications
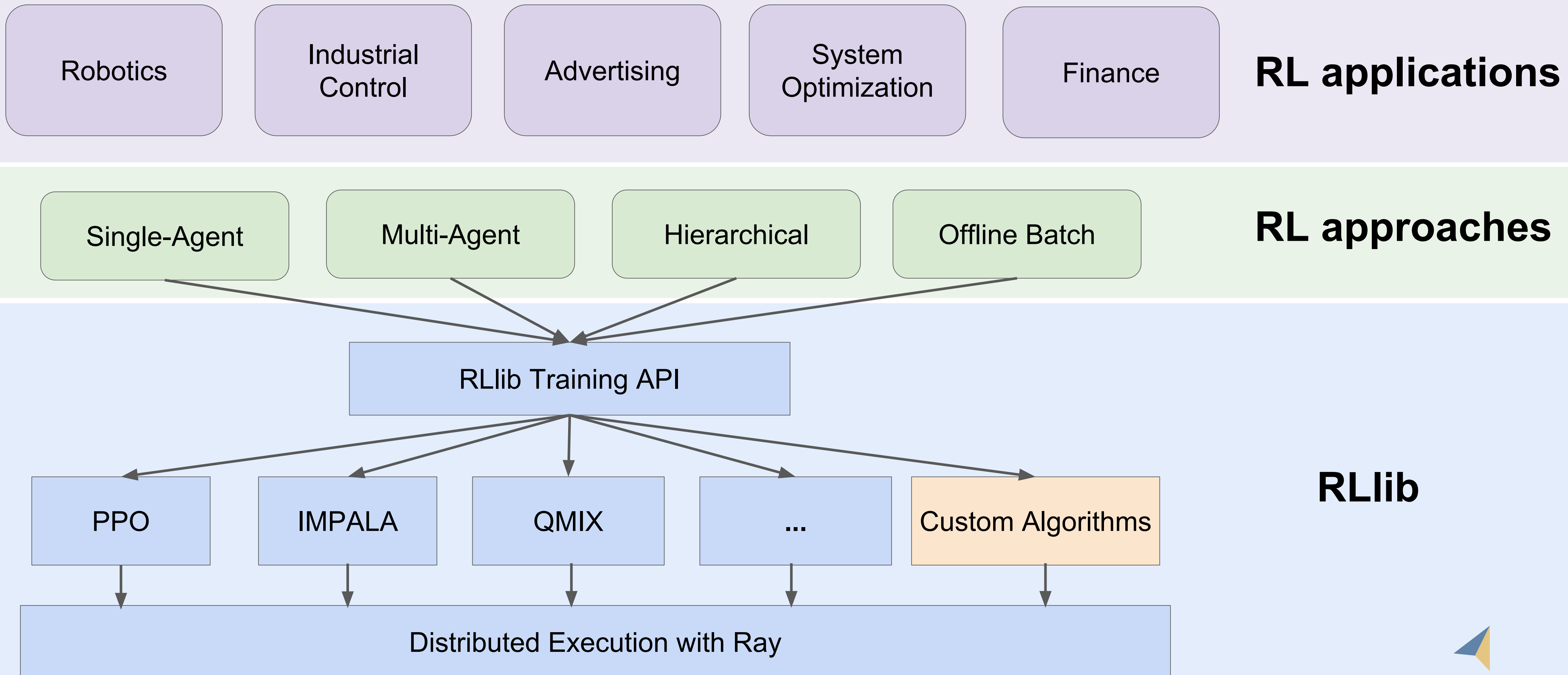
| Robotics | Industrial Control | Advertising | System Optimization | Finance | **RL applications** |

# A scalable, unified library for reinforcement learning

| Robotics | Industrial Control | Advertising | System Optimization | Finance | **RL applications** |

| Single-Agent | Multi-Agent | Hierarchical | Offline Batch | **RL approaches** |

RLlib Training API

| PPO | IMPALA | QMIX | ... | Custom Algorithms | **RLlib** |

Distributed Execution with Ray

# Performance
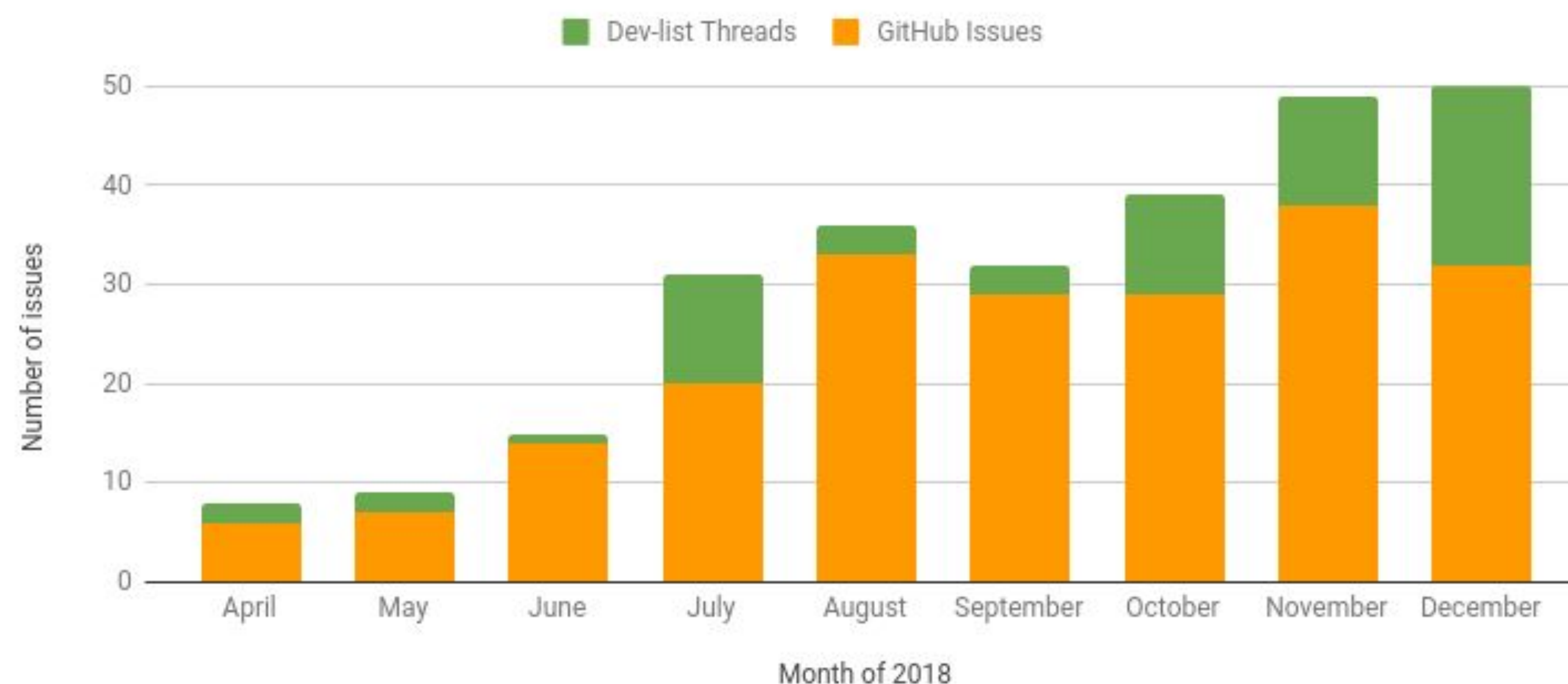
IMPALA and A2C vs A3C after 1 hour of training:

| env | RLlib IMPALA 32-workers | RLlib A2C 5-workers | Mnih et al A3C 16-workers |
|---|---|---|---|
| BeamRider | 3181 | 874 | ~1000 |
| Breakout | 538 | 268 | ~10 |
| QBert | 10850 | 1212 | ~500 |
| SpaceInvaders | 843 | 518 | ~300 |

# User growth in 2018

## Num Issues and Dev-list Threads
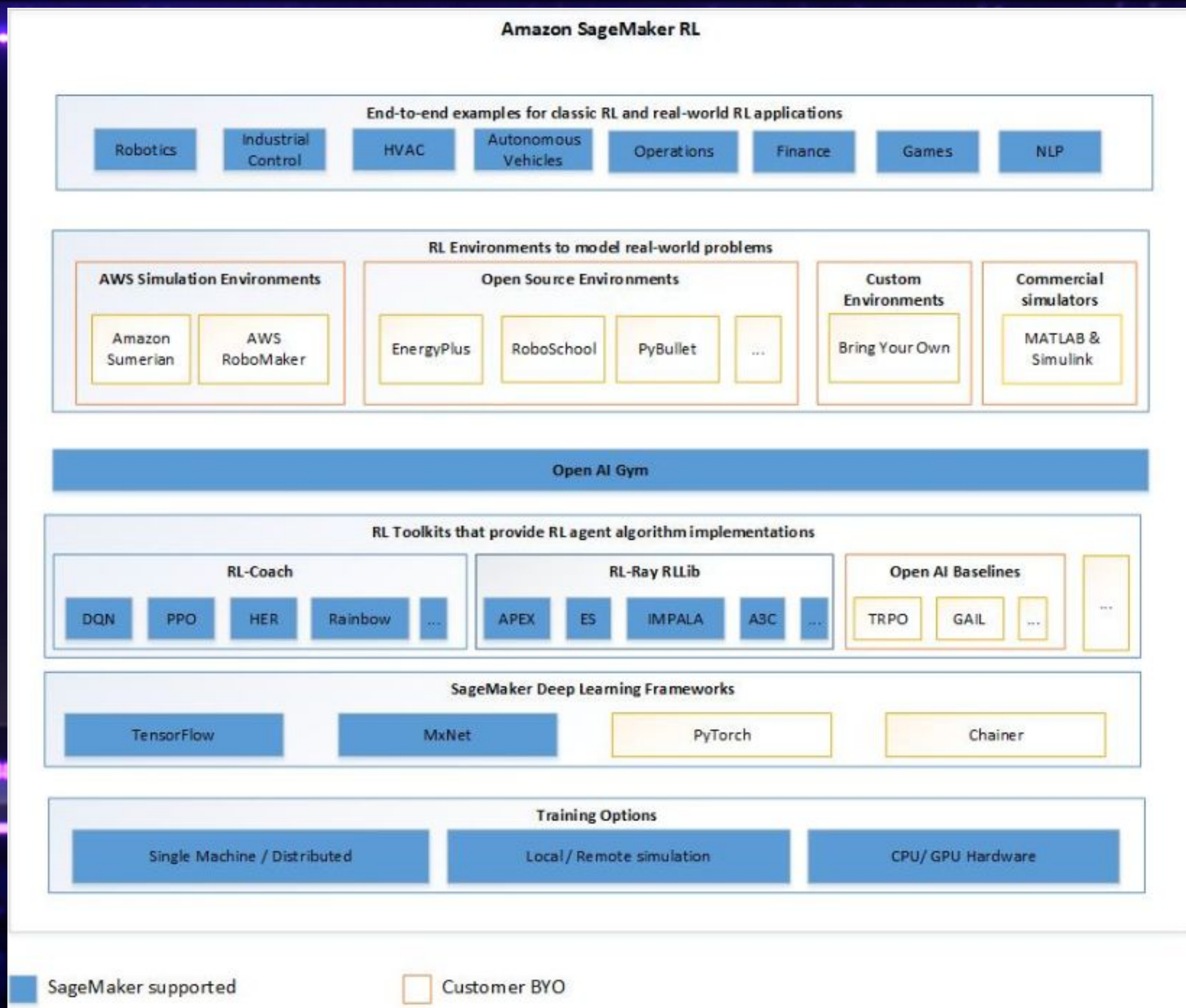


Filtering GitHub and ray-dev@ issues for "rllib":

- user engagement is increasing
- couple dozen companies and research labs using RLlib!

# Project status

- **Goal: be the best library for RL applications and RL applications research**
- Continuing development (https://github.com/ray-project/ray)
  - new algorithms
  - cross-cutting features (env modeling, AutoRL)
  - better performance
- Documentation at https://rllib.io

# Abstractions for Distributed Reinforcement Learning

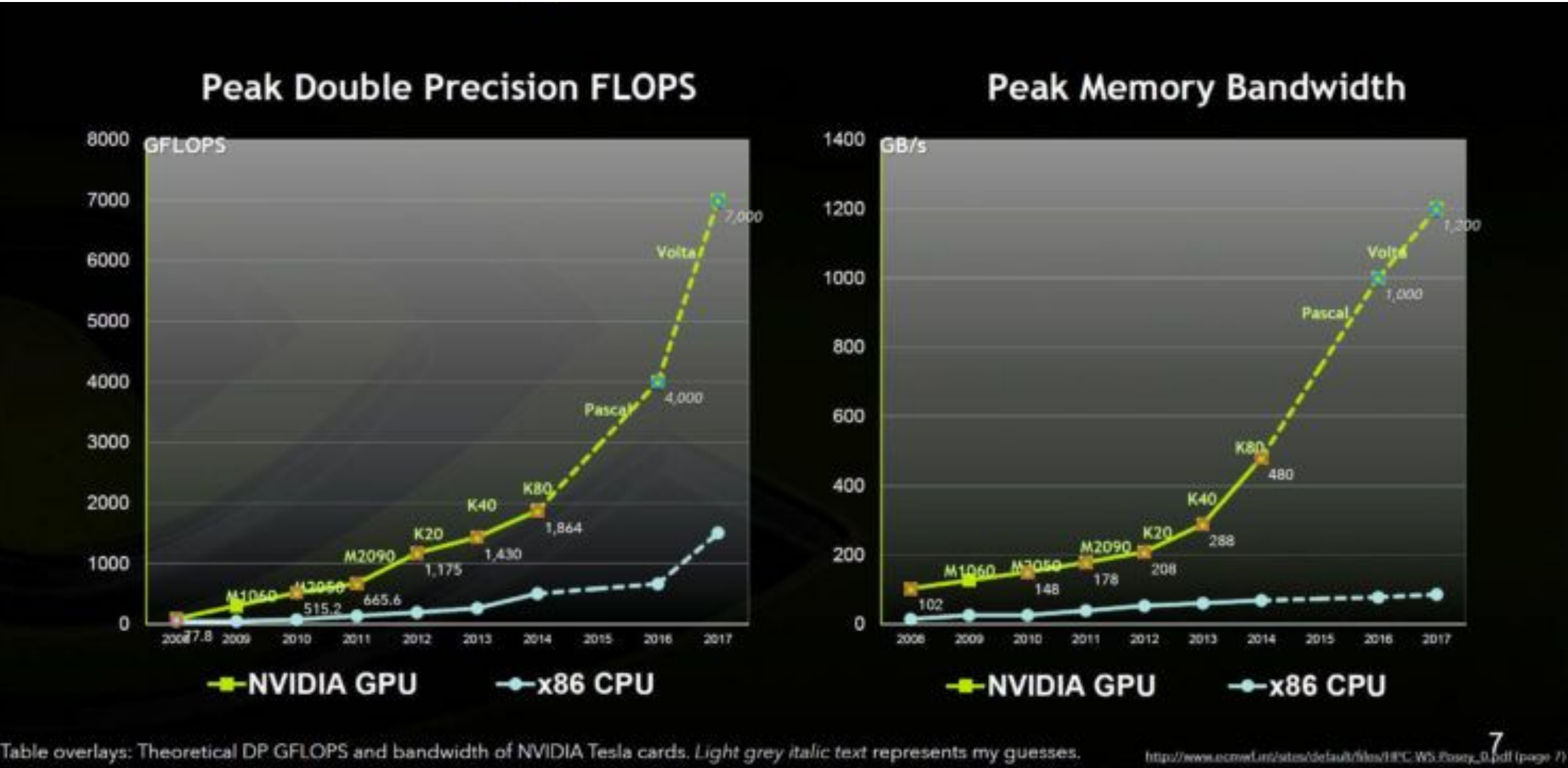# RL research scales with compute



Fig. courtesy NVidia Inc.



CPU   GPU   TPU



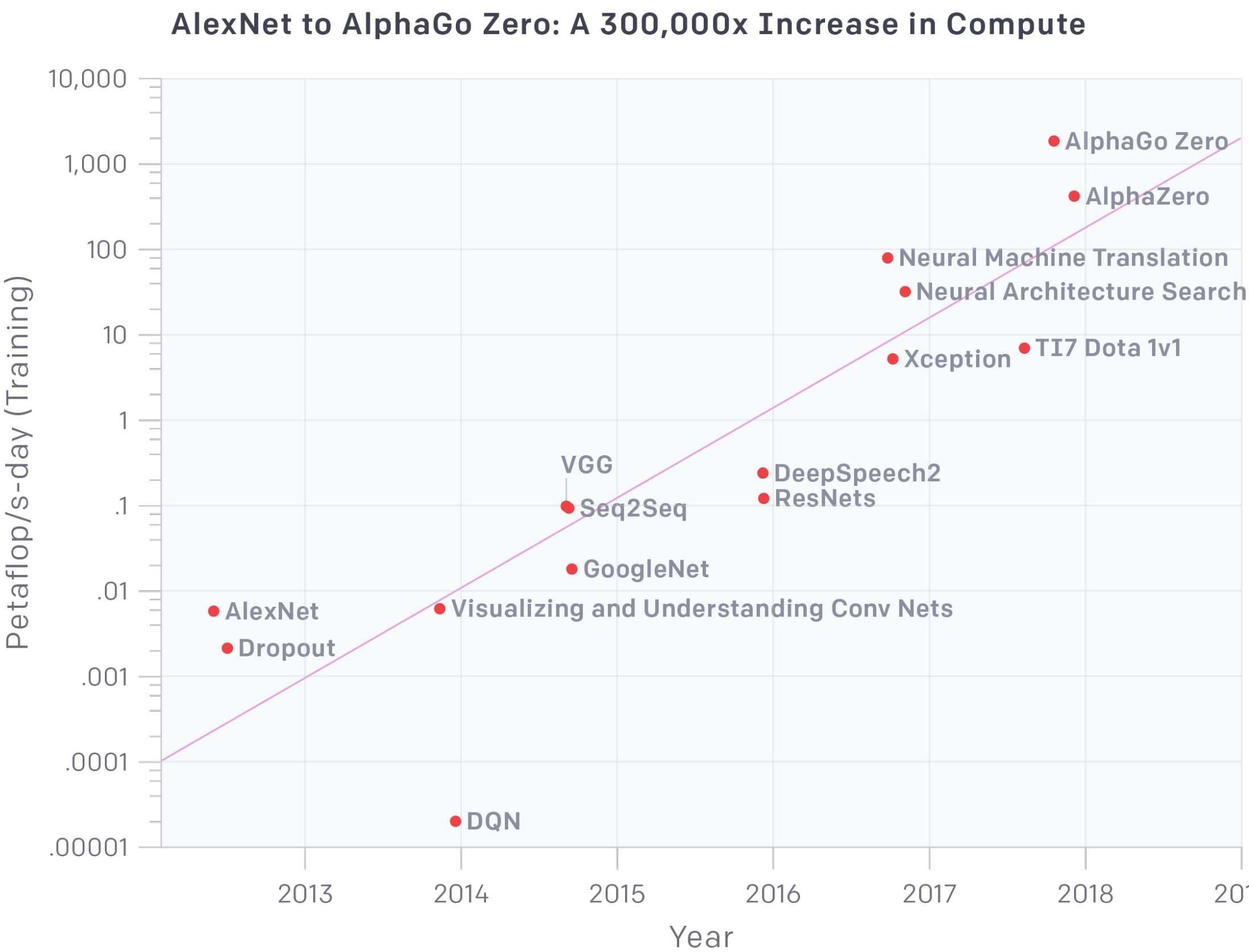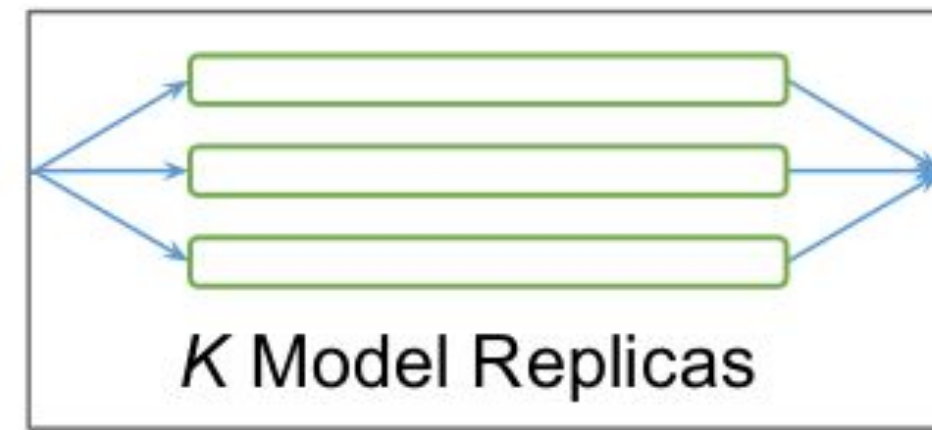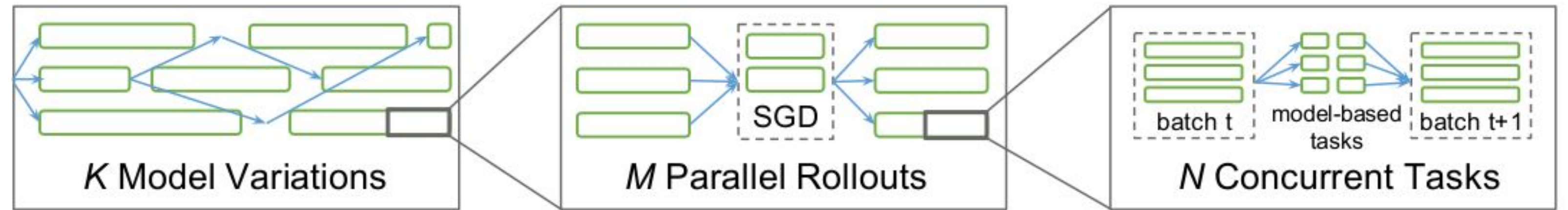Fig. courtesy OpenAI

# How do we leverage this hardware?



(a) Supervised Learning

(b) Reinforcement Learning

*K* Model Replicas

*K* Model Variations

*M* Parallel Rollouts

SGD

*N* Concurrent Tasks

batch t  model-based tasks  batch t+1

**scalable abstractions for RL?**
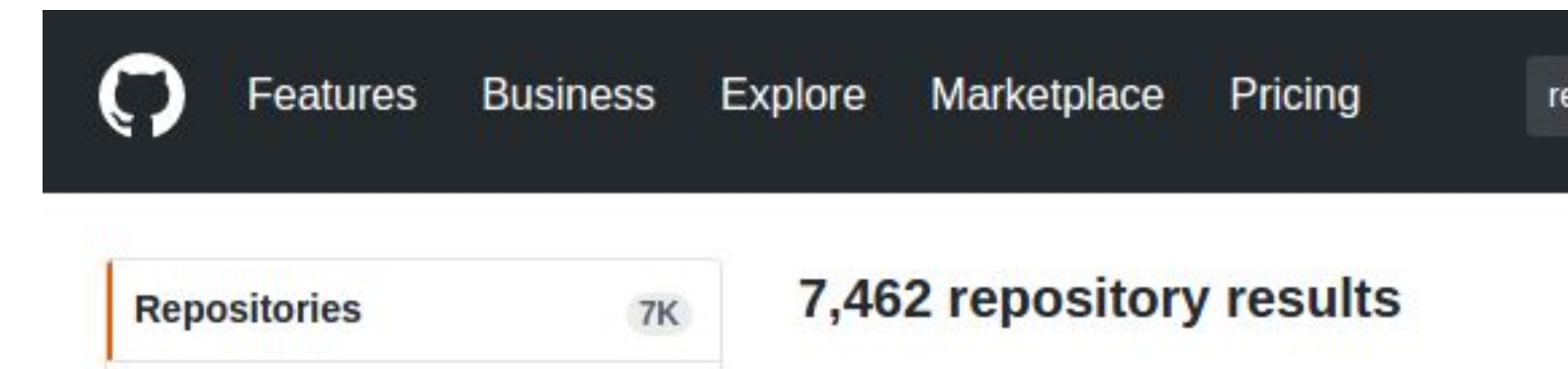
# Example

rllib train --run=PPO --env=Pong-v0 --config='{"num_workers": 1}'


rllib train --run=PPO --env=Pong-v0
        --config='{"num_workers": 4, "num_gpus": 1}'


rllib train --run=PPO --env=Pong-v0
        --config='{"num_workers": 256, "num_gpus": 8}'
        --redis-address=localhost:6379

# Systems for RL today
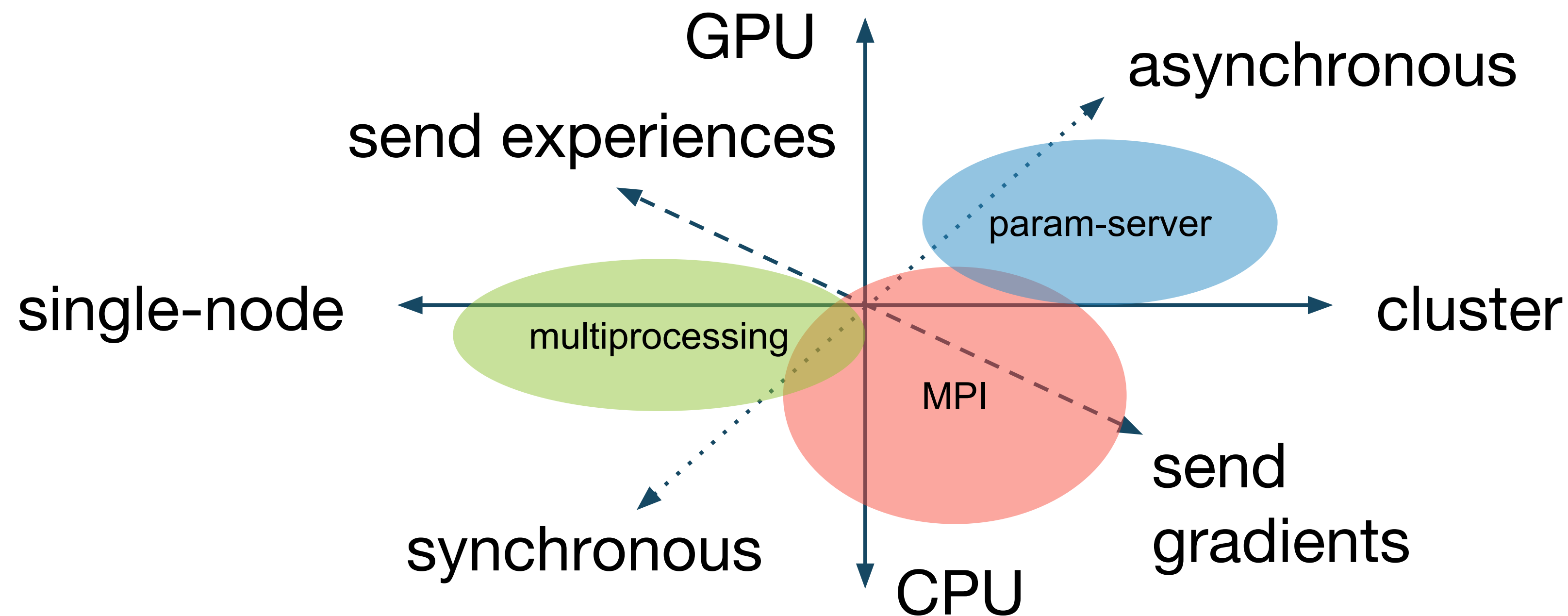


- Many implementations (7000+ repos on GitHub!)
  – how general are they (and do they scale)?

  PPO: multiprocessing, MPI          AlphaZero: custom systems

  Evolution Strategies: Redis        IMPALA: Distributed TensorFlow

  A3C: shared memory, multiprocessing, TF

- Huge variety of algorithms and distributed systems used to implement, but little unification of different architectures
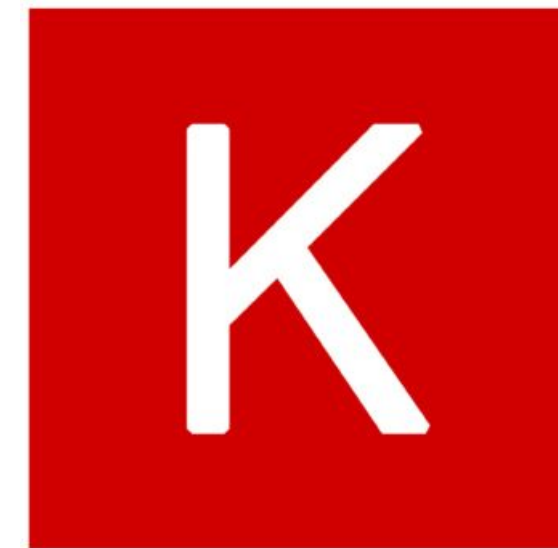
# Challenges to unification

1. Wide range of physical execution strategies for one "algorithm"

# Challenges to unification

2. Tight coupling with deep learning frameworks



Different parallelism paradigms:
  – Distributed TensorFlow vs TensorFlow + MPI?

# Challenges to unification

3. Large variety of algorithms with different structures

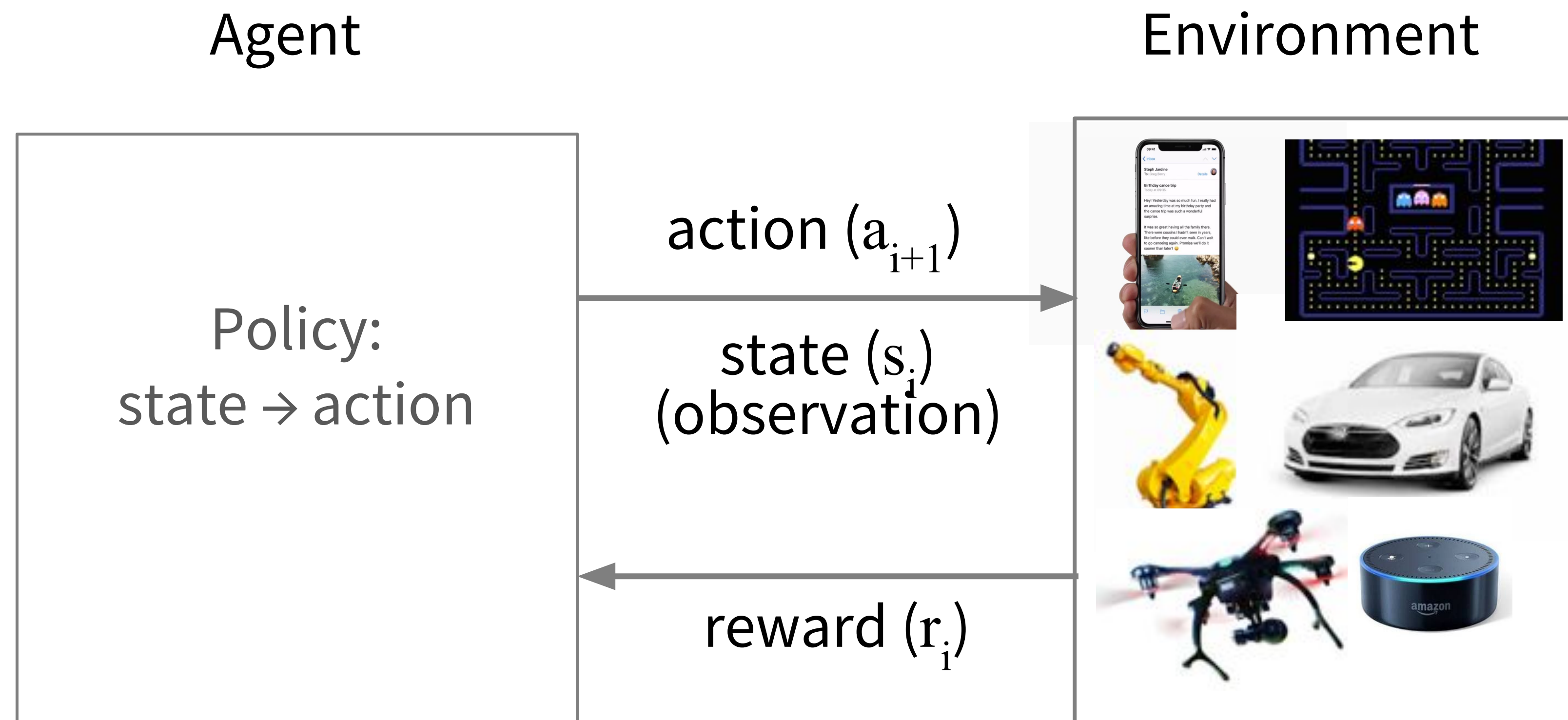| Algorithm Family | Policy Evaluation | Replay Buffer | Gradient-Based Optimizer | Other Distributed Components |
|---|---|---|---|---|
| DQNs | X | X | X | |
| Policy Gradient | X | | X | |
| Off-policy PG | X | X | X | |
| Model-Based/Hybrid | X | | X | Model-Based Planning |
| Multi-Agent | X | X | X | |
| Evolutionary Methods | X | | | Derivative-Free Optimization |
| AlphaGo | X | X | X | MCTS, Derivative-Free Optimization |

# We need abstractions for RL

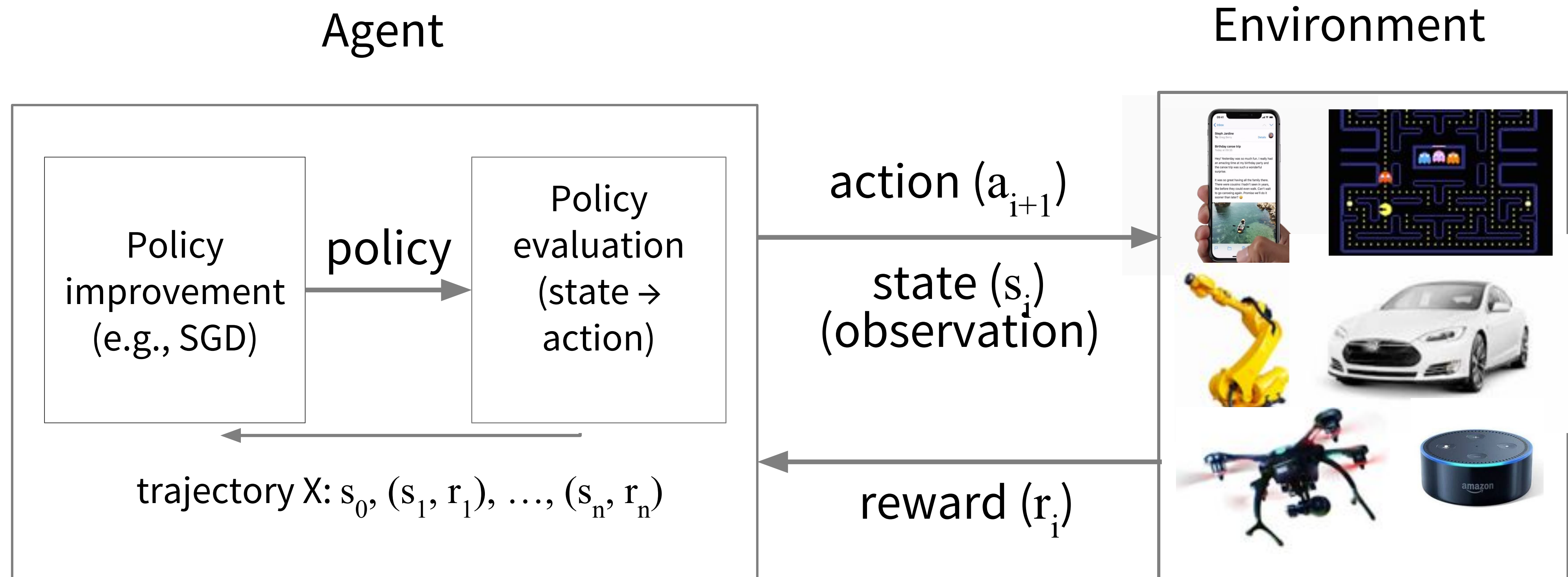*Good abstractions decompose RL algorithms into reusable components.*

Goals:

 – Code reuse across deep learning frameworks
 – Scalable execution of algorithms
 – Easily compare and reproduce algorithms

# Structure of RL computations

Agent

Environment

Policy:
state → action

action ($a_{i+1}$)

state ($s_i$)
(observation)

reward ($r_i$)

# Structure of RL computations

Agent

Environment



Policy improvement (e.g., SGD)

policy →

Policy evaluation (state → action)

action $(a_{i+1})$

state $(s_i)$ (observation)

reward $(r_i)$

trajectory X: $s_0$, $(s_1, r_1)$, ..., $(s_n, r_n)$
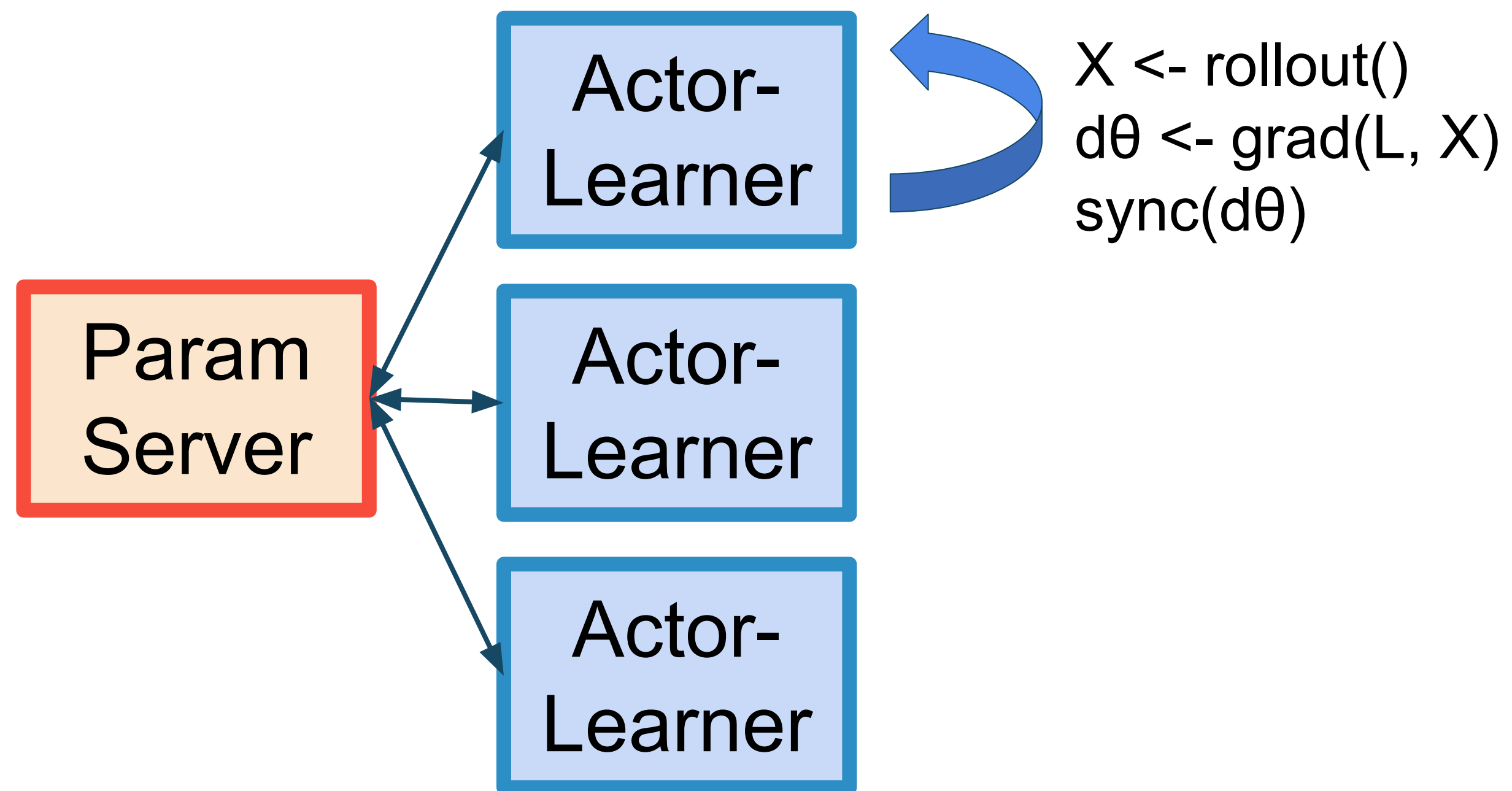
# Many RL loop decompositions

Async DQN (Mnih et al; 2016)

Ape-X DQN (Horgan et al; 2018)



X <- rollout()
dθ <- grad(L, X)
sync(dθ)

θ <- sync()
rollout()

X <- replay()
apply(grad(L, X))

# Common components

Async DQN (Mnih et al; 2016)

Ape-X DQN (Horgan et al; 2018)



**Policy $\pi_\theta(o_t)$**

**Trajectory postprocessor $\rho_\theta(X)$**

Loss $L(\theta,X)$

# Common components

Async DQN (Mnih et al; 2016)

Ape-X DQN (Horgan et al; 2018)



Policy $\pi_\theta(o_t)$

Trajectory postprocessor $\rho_\theta(X)$

**Loss $L(\theta,X)$**

# Structural differences

Async DQN (Mnih et al; 2016)
- Asynchronous optimization
- Replicated workers
- Single machine

**...and this is just one family!**

→ **No existing system can effectively meet all the varied demands of RL workloads.**

Ape-X DQN (Horgan et al; 2018)
- Central learner
- Data queues between components
- Large replay buffers
- Scales to clusters

+ Population-Based Training (Jaderberg et al; 2017)
- Nested parallel computations
- Control decisions based on intermediate results

# Requirements for a new system
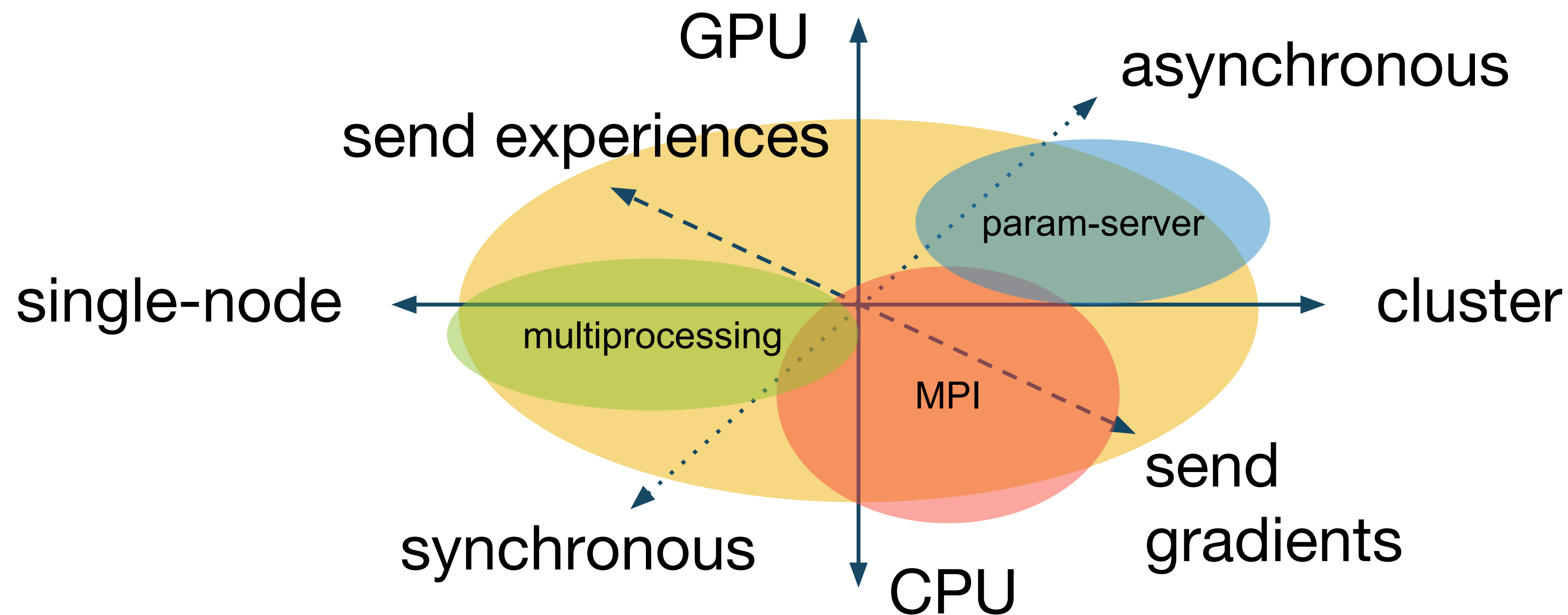
Goal: Capture a broad range of RL workloads with <u>high performance</u> and <u>substantial code reuse</u>

1. Support stateful computations

  - e.g., simulators, neural nets, replay buffers
  - big data frameworks, e.g., Spark, are typically stateless

2. Support asynchrony

  - difficult to express in MPI, esp. nested parallelism

3. Allow easy composition of (distributed) components
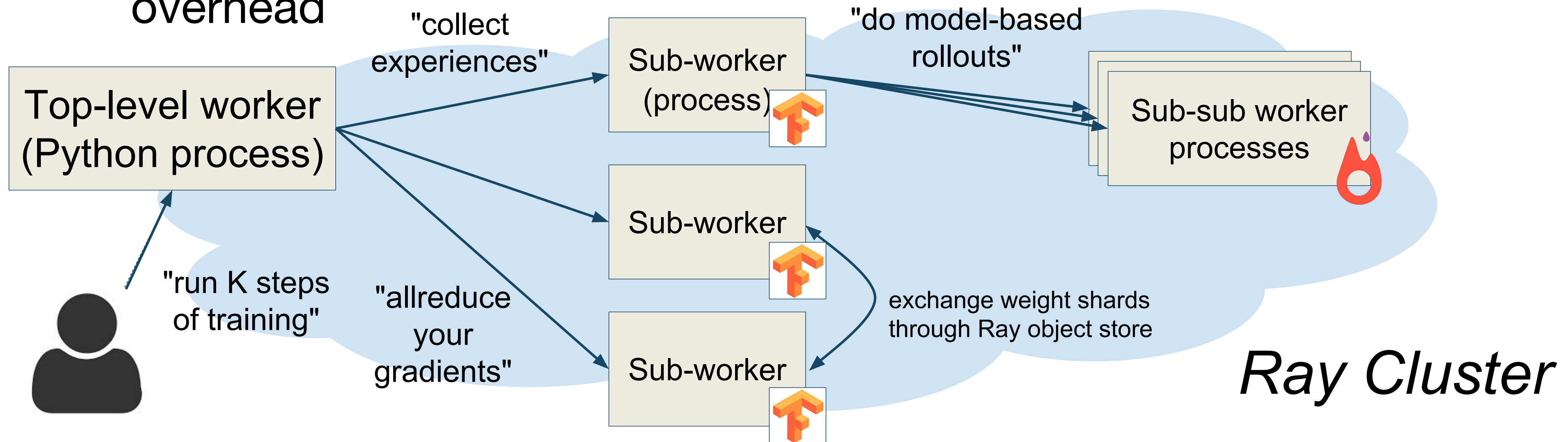
# Ray System Substrate

- RLlib builds on Ray to provide higher-level RL abstractions
- Hierarchical parallel task model with stateful workers
  – flexible enough to capture a broad range of RL workloads (vs specialized sys.)
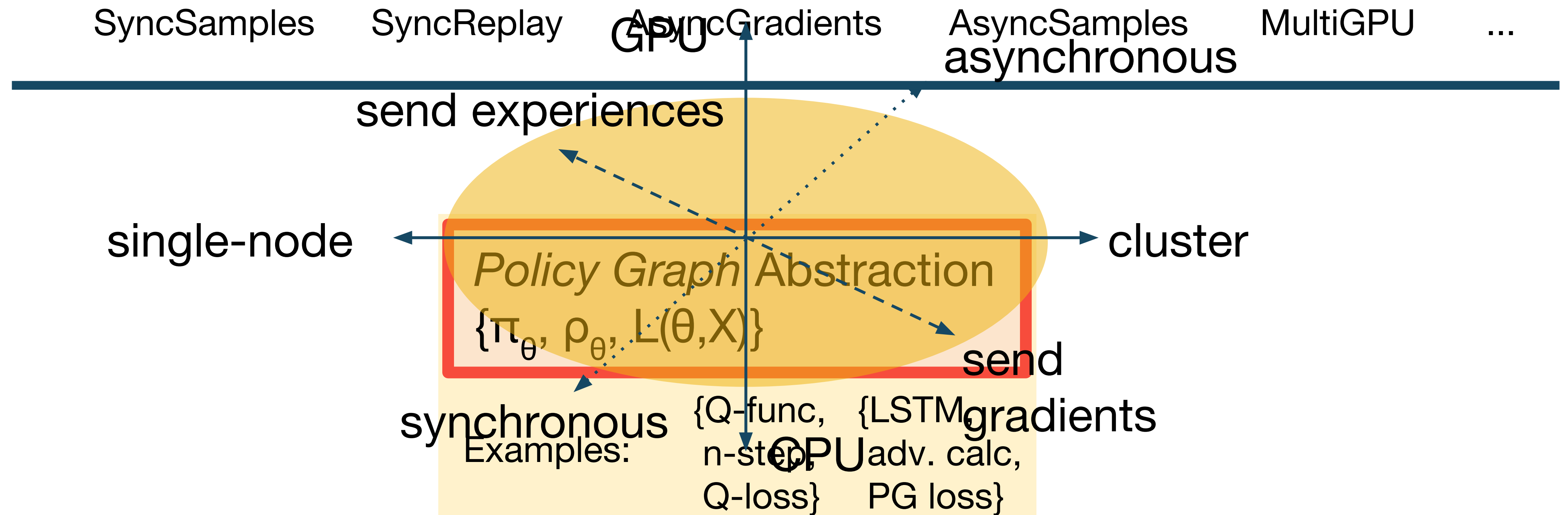
# Hierarchical Parallel Task Model

1. Create Python class instances in the cluster (stateful workers)
2. Schedule short-running tasks onto workers
   – Challenge: High performance: 1e6+ tasks/s, ~200us task overhead



"collect experiences"

"do model-based rollouts"

Top-level worker (Python process)

Sub-worker (process)

Sub-sub worker processes

"run K steps of training"

"allreduce your gradients"

Sub-worker

Sub-worker

exchange weight shards through Ray object store

*Ray Cluster*

# Unifying system enables RL Abstractions

*Policy Optimizer* Abstraction

SyncSamples     SyncReplay     AsyncGradients     AsyncSamples     MultiGPU     ...

GPU

asynchronous

send experiences

single-node

*Policy Graph* Abstraction

$\{\pi_\theta, \rho_\theta, L(\theta,X)\}$

cluster

send

synchronous

Examples:

{Q-func, n-step, Q-loss}

{LSTM adv. calc, PG loss}

gradients
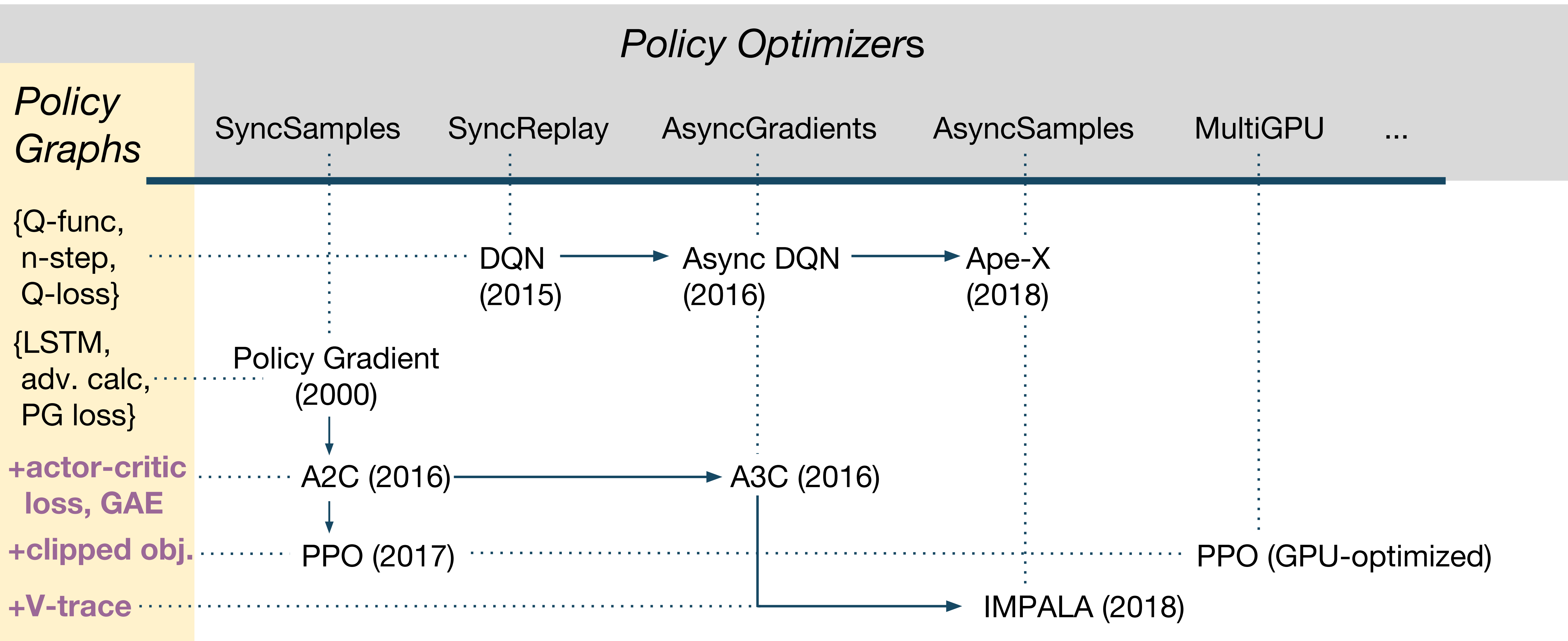
CPU

**Hierarchical Task Model**

# RLlib Abstractions in Action

**Summary:** RLlib addresses challenges in providing scalable abstractions for reinforcement learning.

RLlib is open source and available at http://rllib.io
Thanks!