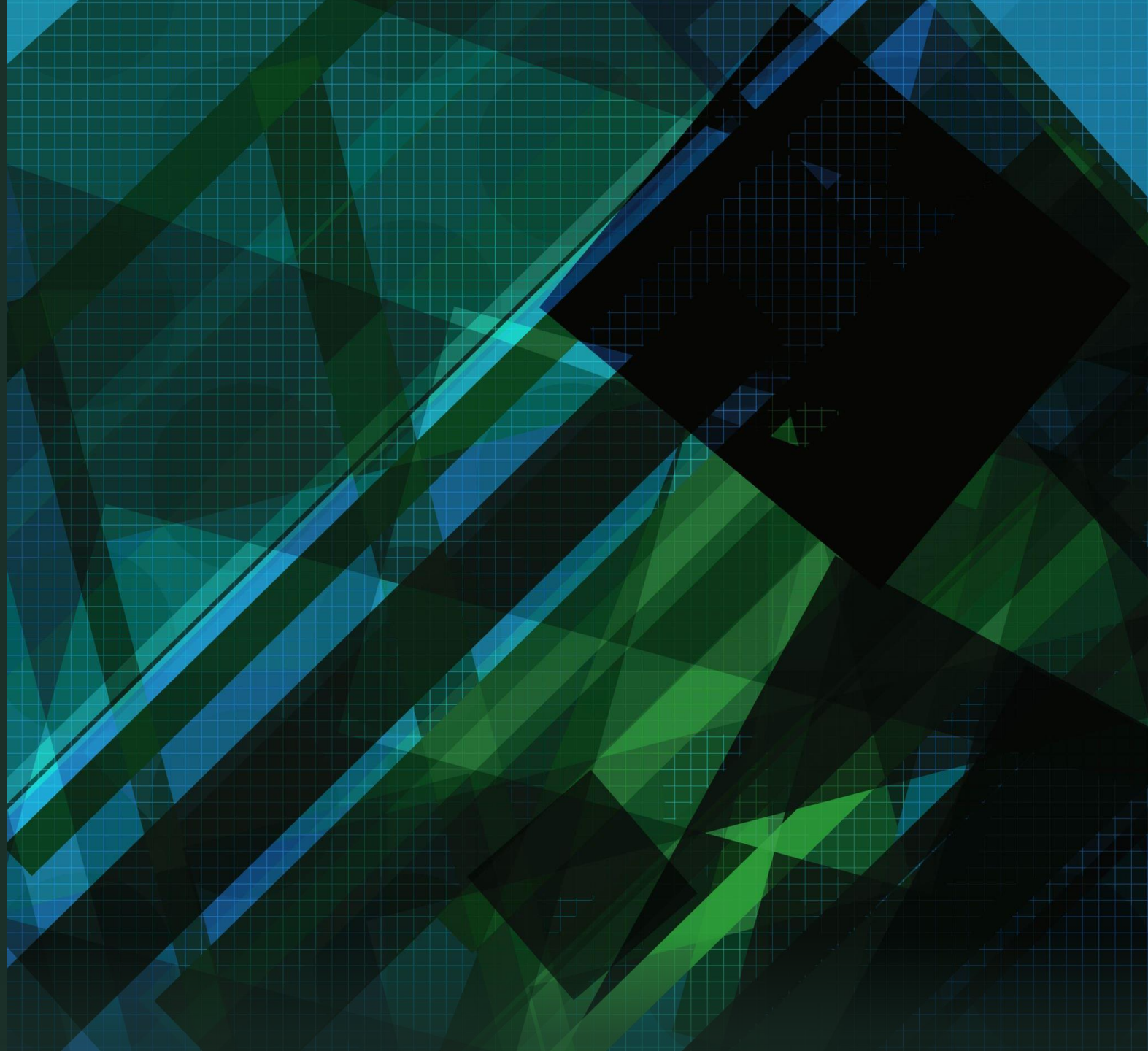


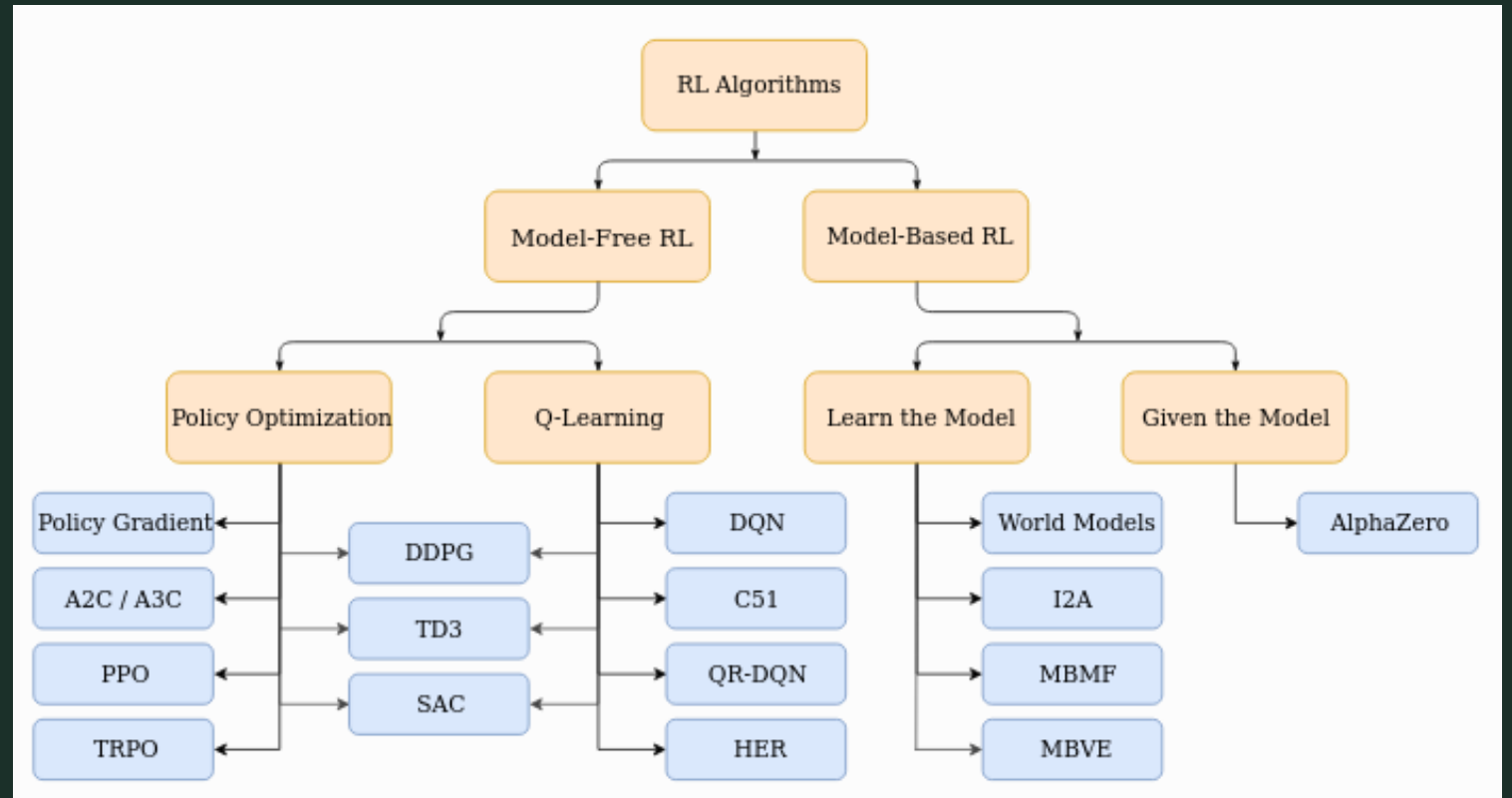


RLGraph + PPG



Model-free RL

- Policy network:
select the action to take
- Value network:
predict the expected
reward in current state



Source: https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html. MIT-licensed, 2018.

Policy optimization



REINFORCE (1999)

- AKA Vanilla Policy Gradient
- Notable success: 2013 Atari
- Standard gradient ascent
- But small parameter changes might still harm performance
- A. Karpathy implemented it in 130 lines in numpy

Algorithm 1 Vanilla Policy Gradient Algorithm

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 7: Compute policy update, either using standard gradient ascent,

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k,$$

or via another gradient ascent algorithm like Adam.

- 8: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k| T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 9: **end for**

PPO (2017)

- Idea: disincentivise large changes in one step of policy improvement.
- One network approximates the policy and the value function (but this is not key to PPO)
- Updates are standard. In each step, loss looks at the old and new probability
- If the actual update is too big, loss treats it as if it changed only by $\epsilon \cdot 100\%$. Typically, $\epsilon = 0.2$

$$\theta_{k+1} = \arg \max_{\theta} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \min \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t), g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right)$$
$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \\ (1 - \epsilon)A & A < 0. \end{cases}$$

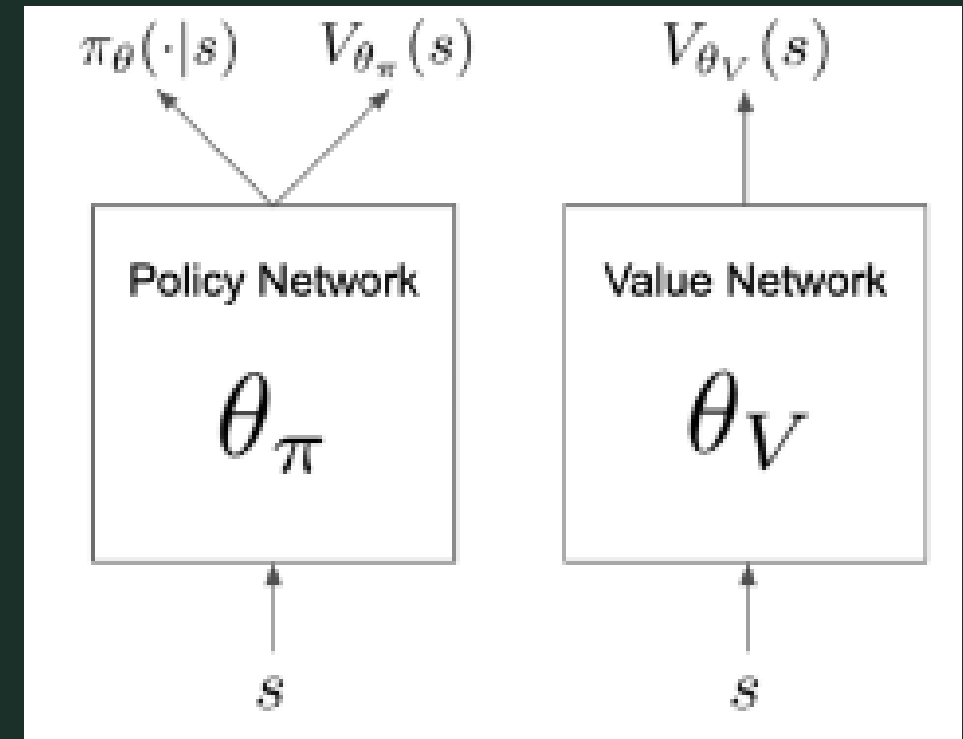
PPG (2020)

- Policy network:

- Same as that used in PPO
- Optimizing the *clipped surrogate function* with an entropy bonus
- Has two heads: a policy head and an auxiliary value head
- Parameters shared

- Value network:

- Predicts the value given a state
- Parameters distinct



Source: PPG paper

Training of PPG

- Policy phase:
 - Estimate advantage function (GAE)
 - Optimize the policy for clipped loss
 - Optimize the value network w/ MSE
- Auxiliary phase:
 - Auxiliary loss: MSE on target values.
 - Joint loss = auxiliary loss + behavioural cloning (keeps policy).
 - Optimize these two losses.

Algorithm 1 PPG

```
for phase = 1, 2, ... do  
  Initialize empty buffer  $B$   
  for iteration = 1, 2, ...,  $N_\pi$  do ▷ Policy Phase  
    Perform rollouts under current policy  $\pi$   
    Compute value function target  $\hat{V}_t^{\text{targ}}$  for each state  $s_t$   
    for epoch = 1, 2, ...,  $E_\pi$  do ▷ Policy Epochs  
      Optimize  $L^{\text{clip}} + \beta_S S[\pi]$  wrt  $\theta_\pi$   
    for epoch = 1, 2, ...,  $E_V$  do ▷ Value Epochs  
      Optimize  $L^{\text{value}}$  wrt  $\theta_V$   
      Add all  $(s_t, \hat{V}_t^{\text{targ}})$  to  $B$   
  Compute and store current policy  $\pi_{\theta_{\text{old}}}(\cdot|s_t)$  for all states  $s_t$  in  $B$   
  for epoch = 1, 2, ...,  $E_{\text{aux}}$  do ▷ Auxiliary Phase  
    Optimize  $L^{\text{joint}}$  wrt  $\theta_\pi$ , on all data in  $B$   
    Optimize  $L^{\text{value}}$  wrt  $\theta_V$ , on all data in  $B$ 
```

$$L^{\text{value}} = \hat{\mathbb{E}}_t \left[\frac{1}{2} (V_{\theta_V}(s_t) - \hat{V}_t^{\text{targ}})^2 \right] \quad L^{\text{aux}} = \frac{1}{2} \cdot \hat{\mathbb{E}}_t \left[(V_{\theta_\pi}(s_t) - \hat{V}_t^{\text{targ}})^2 \right] \quad L^{\text{joint}} = L^{\text{aux}} + \beta_{\text{clone}} \cdot \hat{\mathbb{E}}_t [KL[\pi_{\theta_{\text{old}}}(\cdot|s_t), \pi_{\theta}(\cdot|s_t)]]$$

Results

- Tested on Procgen (which is an improved variation of Atari)
- Converges faster and learns better than PPO

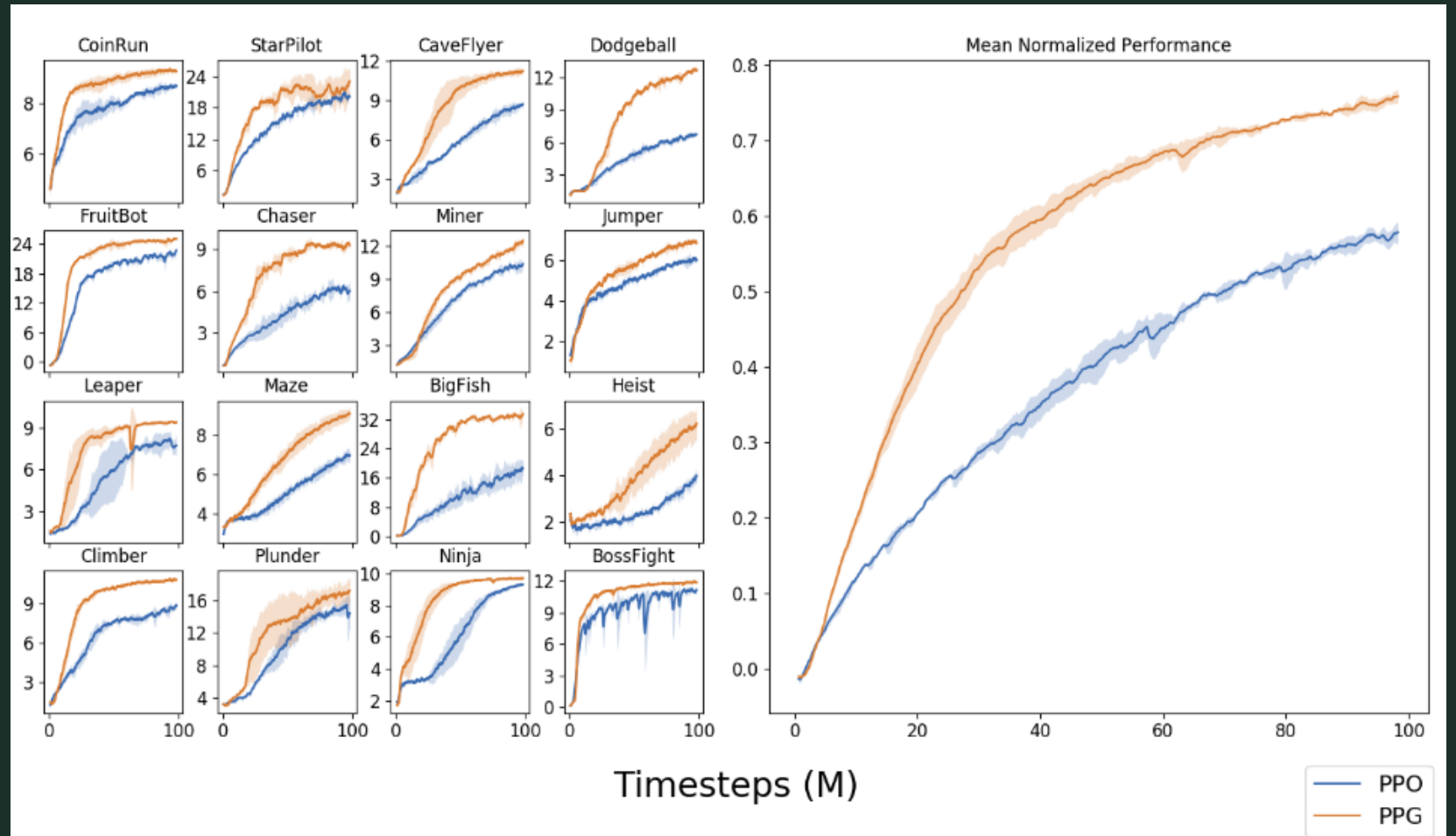
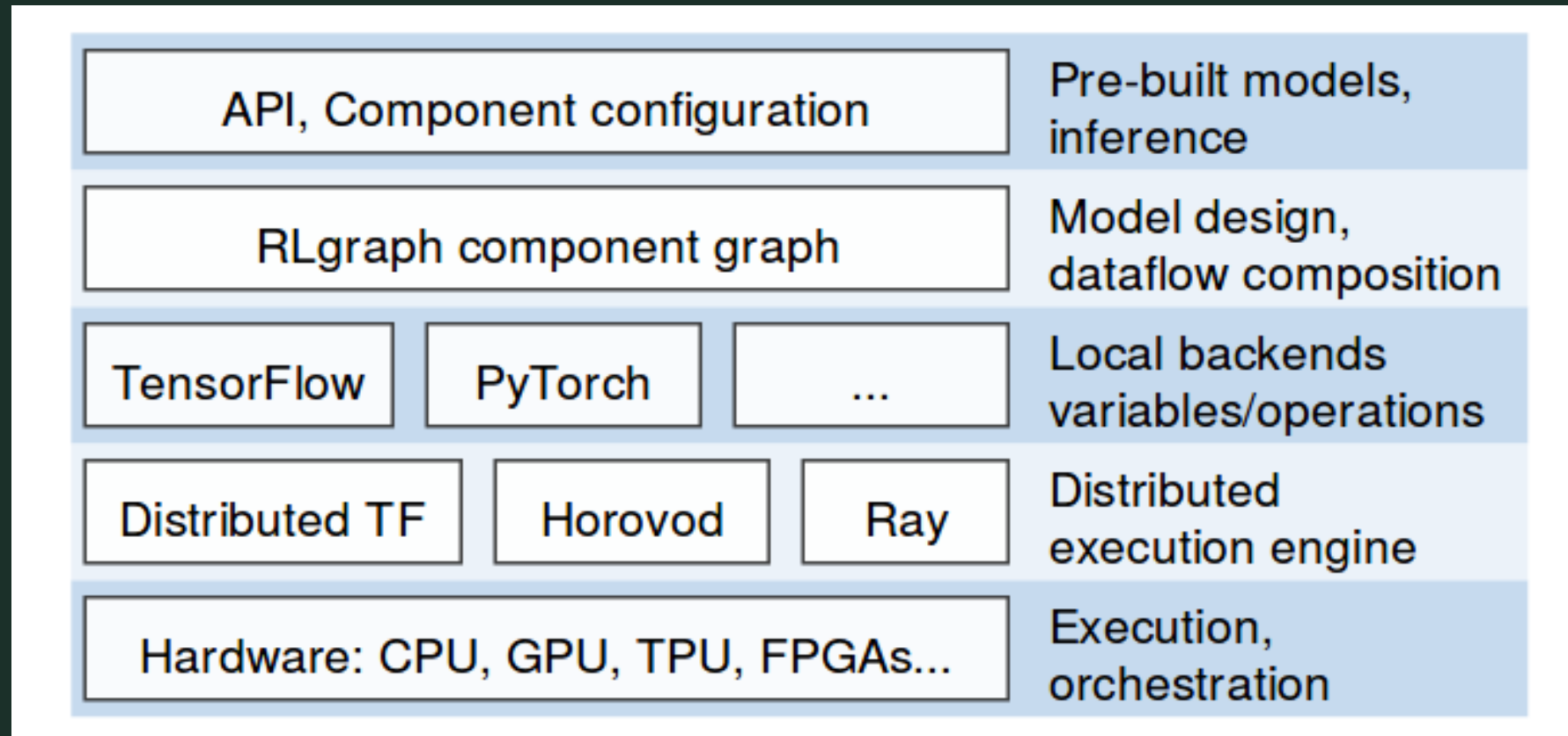


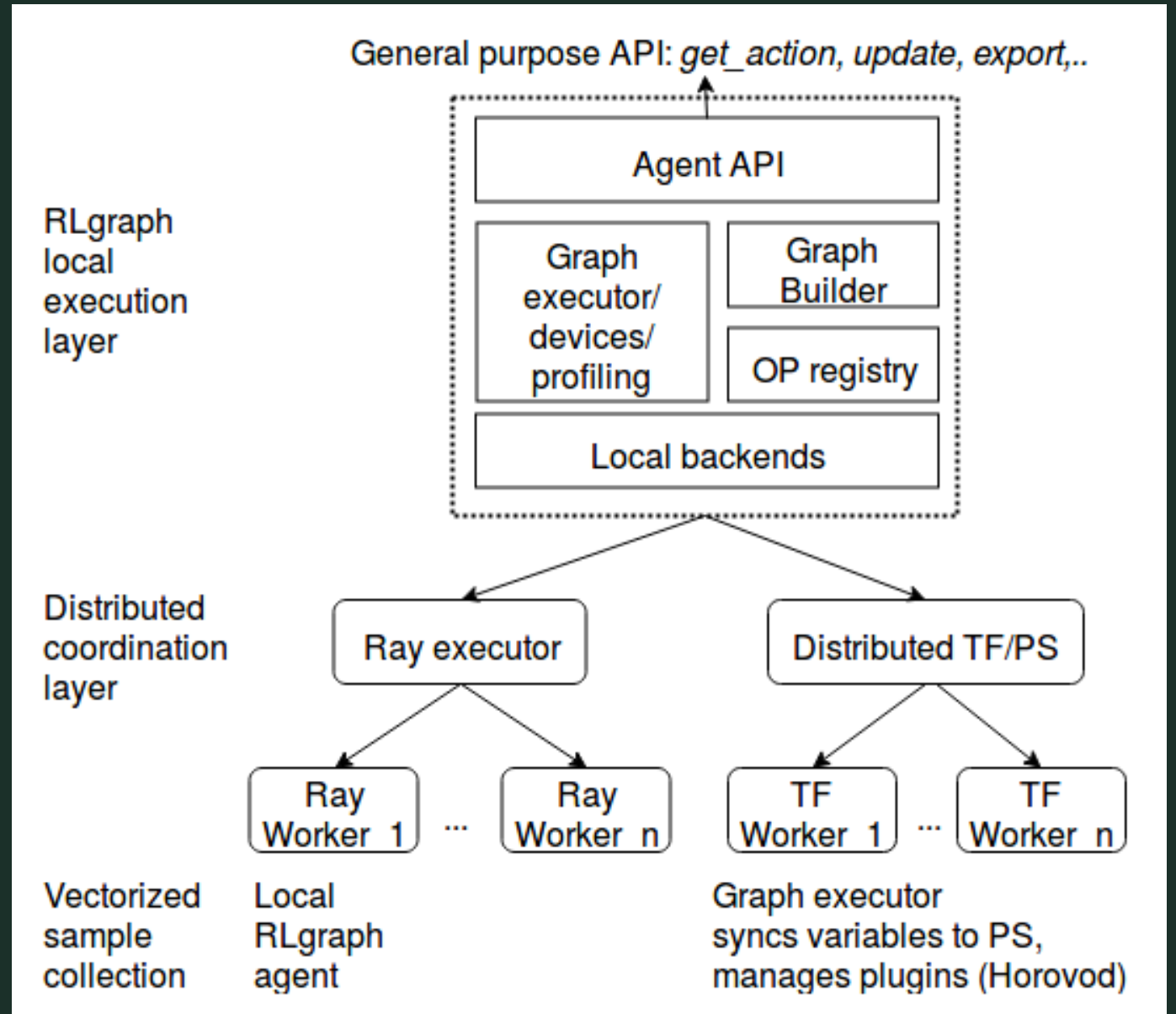
Figure 2: Sample efficiency of PPG compared to a PPO baseline

RLGraph



Source: <https://rlgraph.github.io/rlgraph/2019/01/04/introducing-rlgraph.html>

RLGraph



Goals

- Reimplement the PPG in RLGraph
- Benchmark performance of PPG in RLGraph on Atari and Gym scenarios

