

---

---

# Combining Policy Gradient and Q-Learning

Group10 : 吳程峻、林浩君、許承壹

---

---

# Recap: Policy Gradient and Q-learning

- Policy Gradient:

We can parameterize the policy directly and attempt to improve it via gradient ascent on the objective  $J$  given by

$$\nabla_{\theta} J(\pi) = \mathbf{E}_{s,a} Q^{\pi}(s, a) \nabla_{\theta} \log \pi(s, a)$$

- Q-learning:

In value based RL, we approximate the Q-values using a function approximator (typically neural network), and follow the best actions along the maximum Q-values of each states. **Theoretically, we apply Bellman contraction operator to find the optimal Q-value.**

$$\mathbf{E}_{s,a} (\mathcal{T}^{\pi} Q(s, a; \theta) - Q(s, a; \theta)) \nabla_{\theta} Q(s, a; \theta)$$

# Comparison

- What are the advantages and disadvantages of Policy Gradient and Q-Learning method respectively?

	Advantages	Disadvantages
Policy Gradient	<ul style="list-style-type: none"><li>• Better convergence properties</li><li>• Effective in high-dimensional or continuous action spaces</li><li>• Can learn stochastic policies</li></ul>	<ul style="list-style-type: none"><li>• Hard for following off-policy</li><li>• Sample inefficiency</li><li>• Usually converge to local optimum</li><li>• Hard for learning deterministic policies</li><li>• Evaluating a policy is inefficient and high variance</li></ul>
Q-Learning	<ul style="list-style-type: none"><li>• Can apply off-policy</li><li>• Sample efficiency</li></ul>	<ul style="list-style-type: none"><li>• Training process is relatively slow</li><li>• Troubling with continuous action space</li></ul>

# Combining These Two Family Together:

## PGQL(Policy Gradient Q-Learning)

- From the previous slide, we see there are many disadvantages to both two RL-algorithm families. **So, is it possible to combine them together, and eliminate each other's disadvantages and exploit both their benefits?**
- This is exactly what this paper is trying to do!
- Based on the architecture of A3C, they tries to derive a variant of Q-value function induced by current behavior policy. **This method is called PGQL(Policy Gradient Q-Learning)**

# Formal Definition of PGQL

We define a variant of Q-value using the policy as

$$\tilde{Q}^{\pi}(s, a) = \alpha(\log \pi(s, a) + H^{\pi}(s)) + V(s)$$

where  $V$  is parameterized by  $w$ , i.e. we obtain this value by training a neural network. And  $H^{\pi}(s) = -\sum_a \pi(s, a) \log \pi(s, a)$  denotes the entropy of policy  $\pi$ , and  $\alpha > 0$  is the regularization penalty parameter.

In short, they tries to **estimate Q-value by current policy, and value function**, and also **add entropy regularization** to make this estimation robust.

**But, Does this formula have any guarantee to reach optimal Q-value?**

# Theory in Behind: Bellman Residual

By the optimal Bellman r-contraction operator  $\mathcal{T}^*$ , we want to prove that there exist a fixed point  $\tilde{Q}^{\pi_\alpha}$ , which is the optimal Q-value function, and the **Bellman residual of the induced Q-values for the PGQL updates converges to 0**. Hence we can **guarantee that our updating process is leading us to a nice convergence**.

**Note:** Bellman operator  $\mathcal{T}^\pi$  is defined as

$$\mathcal{T}^\pi Q(s, a) = \mathbf{E}_{s', r, b} (r(s, a) + \gamma Q(s', b))$$

Optimal Bellman operator  $\mathcal{T}^*$  is defined as

$$\mathcal{T}^* Q(s, a) = \mathbf{E}_{s', r} (r(s, a) + \gamma \max_b Q(s', b))$$

$Q^{\pi_\alpha}$  is the same as  $\tilde{Q}^\pi$  in the previous page, and  $\tilde{Q}^{\pi_\alpha}$  is the fixed point we want to obtain, they somehow confuse the notations when giving proof.

# Proof: Bellman Residual converges to 0

## Step1: Show That We Are Contracting to a Specific Point.

- **Ultimate Target:** Show that  $\lim_{\alpha \rightarrow 0} \|\mathcal{T}^*Q^{\pi_\alpha} - Q^{\pi_\alpha}\| = 0$ , where we would have  $\tilde{Q}^{\pi_\alpha}$  as our fixed point.

First can write  $\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha} = \eta(\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha})$  derived by this paper, hence, we can show

$$\begin{aligned}\|\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| &= \eta\|\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| \\ &= \eta\|\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha}\tilde{Q}^{\pi_\alpha} + \mathcal{T}^{\pi_\alpha}\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| \\ &\leq \eta(\|\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha}\tilde{Q}^{\pi_\alpha}\| + \|\mathcal{T}^{\pi_\alpha}\tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha}Q^{\pi_\alpha}\|) \\ &\leq \eta(\|\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha}\tilde{Q}^{\pi_\alpha}\| + \gamma\|\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\|) \\ &\leq \eta/(1 - \eta\gamma)\|\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha}\tilde{Q}^{\pi_\alpha}\|,\end{aligned}$$

Therefore, we firstly prove that as  $\alpha \rightarrow 0$ ,  $\|\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| \rightarrow 0$ , **this means that we are going to approach a point  $\tilde{Q}^{\pi_\alpha}$  if we keep updating  $Q^{\pi_\alpha}$ , but we still don't know whether  $\tilde{Q}^{\pi_\alpha}$  is a fixed point or not.**

## Step2: Show This Specific point Is a Fixed Point

By the previous proof(**A fixed point  $\tilde{Q}^{\pi_\alpha}$  exists**), we can deduce that

$$\begin{aligned}\|\mathcal{T}^* \tilde{Q}^{\pi_\alpha} - \tilde{Q}^{\pi_\alpha}\| &= \|\mathcal{T}^* \tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha} \tilde{Q}^{\pi_\alpha} + \mathcal{T}^{\pi_\alpha} \tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha} + Q^{\pi_\alpha} - \tilde{Q}^{\pi_\alpha}\| \\ &\leq \|\mathcal{T}^* \tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha} \tilde{Q}^{\pi_\alpha}\| + \|\mathcal{T}^{\pi_\alpha} \tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha} Q^{\pi_\alpha}\| + \|Q^{\pi_\alpha} - \tilde{Q}^{\pi_\alpha}\| \\ &\leq \|\mathcal{T}^* \tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha} \tilde{Q}^{\pi_\alpha}\| + (1 + \gamma) \|\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| \\ &< 3/(1 - \eta\gamma) \|\mathcal{T}^* \tilde{Q}^{\pi_\alpha} - \mathcal{T}^{\pi_\alpha} \tilde{Q}^{\pi_\alpha}\|,\end{aligned}$$

which therefore also converges to 0 in the limit.

**This means that  $\tilde{Q}^{\pi_\alpha}$  is actually a fixed point. We cannot contract anymore.**



## Step3: Bellman Residual Approaches 0

So, By the previous two steps, we obtain:

$$\begin{aligned}\|\mathcal{T}^*Q^{\pi_\alpha} - Q^{\pi_\alpha}\| &= \|\mathcal{T}^*Q^{\pi_\alpha} - \mathcal{T}^*\tilde{Q}^{\pi_\alpha} + \mathcal{T}^*\tilde{Q}^{\pi_\alpha} - \tilde{Q}^{\pi_\alpha} + \tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| \\ &\leq \|\mathcal{T}^*Q^{\pi_\alpha} - \mathcal{T}^*\tilde{Q}^{\pi_\alpha}\| + \|\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - \tilde{Q}^{\pi_\alpha}\| + \|\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| \\ &\leq (1 + \gamma)\|\tilde{Q}^{\pi_\alpha} - Q^{\pi_\alpha}\| + \|\mathcal{T}^*\tilde{Q}^{\pi_\alpha} - \tilde{Q}^{\pi_\alpha}\|,\end{aligned}$$

As it implies that  $\lim_{\alpha \rightarrow 0} \|\mathcal{T}^*Q^{\pi_\alpha} - Q^{\pi_\alpha}\| = 0$ .

**So, we know that we have a nice convergence as we keep updating the variant of Q-value function  $Q^{\pi_\alpha}$ .**

# Applying PGQL to Update NN Parameters

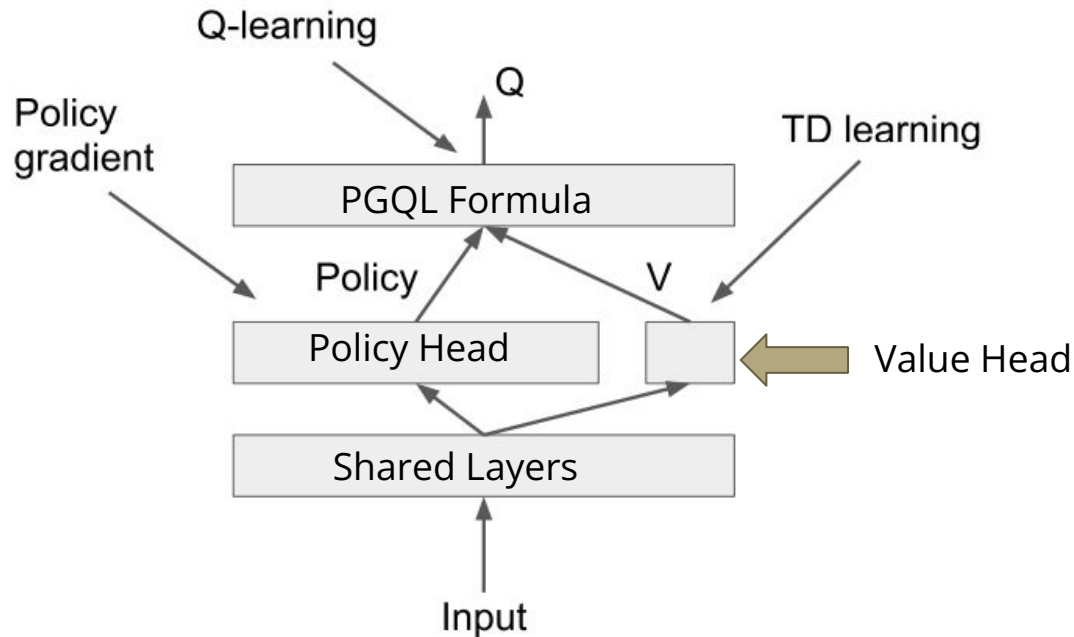
Since we know that at the fixed point the Bellman residual will be small for small  $\alpha$ , we can consider updating the parameters to reduce the Bellman residual in a fashion similar to Q-learning, i.e.,

$$\Delta\theta \propto \mathbf{E}_{s,a}(\mathcal{T}^*\tilde{Q}^\pi(s,a) - \tilde{Q}^\pi(s,a))\nabla_\theta \log \pi(s,a), \quad \Delta w \propto \mathbf{E}_{s,a}(\mathcal{T}^*\tilde{Q}^\pi(s,a) - \tilde{Q}^\pi(s,a))\nabla_w V(s).$$

This is exactly Q-learning applied to a particular form of the Q-values, and can also be interpreted as an actor-critic algorithm with an optimizing critic.

**In practice, we are leveraging between Q-learning and A3C algorithm.**

# Implementation Detail(1):Network Architecture

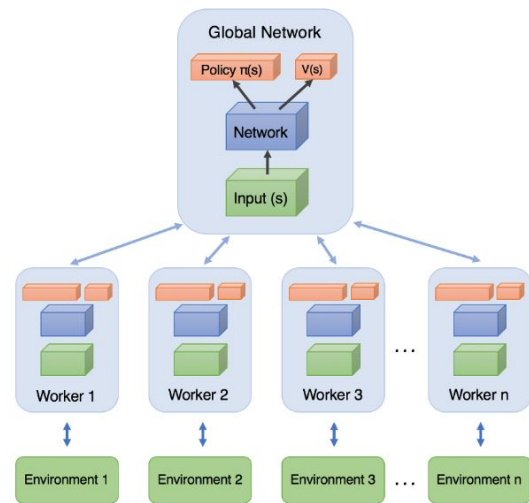


- We use **only one network** with shared layers, and one head for giving current **policy**, another head for giving **value**. Lastly, **The last layer is simply parameterized by the PGQL formula, without any network weights.**

# Implementation Detail(2): Using The A3C Architecture

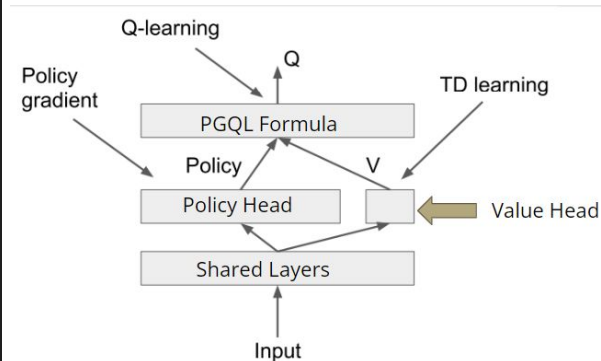
This is the **detail they didn't emphasize**, and hence, **we did have a hard time doing research** on **what neural network they exactly implemented**.

Precisely, they used the **A3C architecture**, which will have a **central shared model**, and use threads to **launch plenty of workers** to do the training process, and finally **update the weights from the workers back to the central shared model**.



# Code of Our implemented Neural Network

```
37 class NNPolicy(nn.Module): # an actor-critic neural network
38     def __init__(self, channels, memsize, num_actions):
39         super(NNPolicy, self).__init__()
40         self.conv1 = nn.Conv2d(channels, 32, 3, stride=2, padding=1)
41         self.conv2 = nn.Conv2d(32, 32, 3, stride=2, padding=1)
42         self.conv3 = nn.Conv2d(32, 32, 3, stride=2, padding=1)
43         self.conv4 = nn.Conv2d(32, 32, 3, stride=2, padding=1)
44         self.gru = nn.GRUCell(32 * 5 * 5, memsize)
45         self.critic_linear, self.actor_linear = nn.Linear(memsize, 1), nn.Linear(memsize, num_actions)
46
47     def forward(self, inputs, train=True, hard=False):
48         inputs, hx = inputs
49         x = F.elu(self.conv1(inputs))
50         x = F.elu(self.conv2(x))
51         x = F.elu(self.conv3(x))
52         x = F.elu(self.conv4(x))
53         hx = self.gru(x.view(-1, 32 * 5 * 5), (hx))
54         return self.critic_linear(hx), self.actor_linear(hx), hx
55
56     def try_load(self, save_dir):
57         paths = glob.glob(save_dir + '*.tar') ; step = 0
58         if len(paths) > 0:
59             ckpts = [int(s.split('.')[2]) for s in paths]
60             ix = np.argmax(ckpts) ; step = ckpts[ix]
61             self.load_state_dict(torch.load(paths[ix]))
62             print("\tno saved models") if step == 0 else print("\tloaded model: {}".format(paths[ix]))
63         return step
```



## Implementation Detail(3): Calculate Q from the policy

- This part correspond to the formula  $\tilde{Q}^{\pi}(s, a) = \alpha(\log \pi(s, a) + H^{\pi}(s)) + V(s)$

```
122 def calculate_q_value(value, policy, action, args):
123     entropy = -torch.sum(policy * torch.log(policy + 1e-8), 1).reshape(args.batch_size, 1)
124     pi = policy.gather(1, action.long())
125     return args.alpha * (torch.log(pi + 1e-8) + entropy) + value
```

- We obtain the first and second input parameters(value,policy) by the outputs of the neural network, then we simply calculate the policy's entropy and try to build the formula.
- This is exactly the Q-value that we use to update the network.

# Implementation Detail(4): PGQL Update

- This part corresponds to how we exactly update the network parameters.
- Here we will use **replay buffer** to help us train the network.

```
127 def q_update(shared_model, shared_optimizer, model, args, memory):
128     state, action, reward, next_state, done, hx = memory.sample(args.batch_size)
129
130     value, logit, nhx = model((state, hx))
131     policy = torch.exp(F.log_softmax(logit, dim=-1))
132     q_value = calculate_q_value(value, policy, action, args)
133
134     with torch.no_grad():
135         value_next, logit_next, _ = model((next_state, nhx))
136         policy_next = torch.exp(F.log_softmax(logit_next, dim=-1))
137         action_next = torch.argmax(policy_next, 1).reshape(args.batch_size, 1)
138         q_next = calculate_q_value(value_next, policy_next, action_next, args)
139         q_target = reward + args.gamma * q_next * (1 - done)
140
141     mse_criterion = nn.MSELoss()
142     loss = mse_criterion(q_value, q_target)
143
144     update_shared_model(shared_model, shared_optimizer, model, loss)
```

# Experiment Results(1): Paper Part

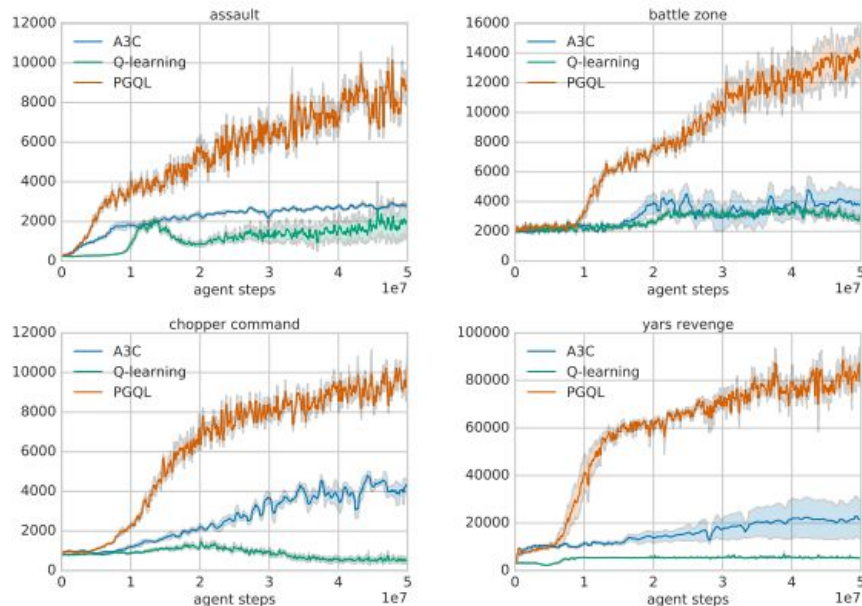
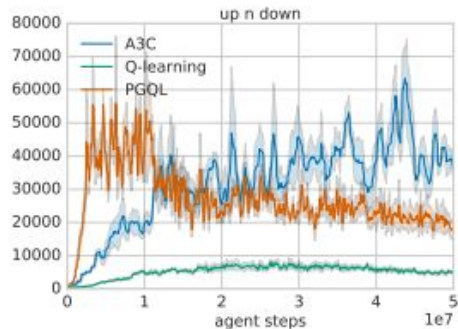
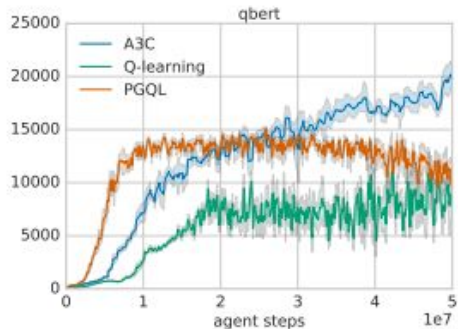
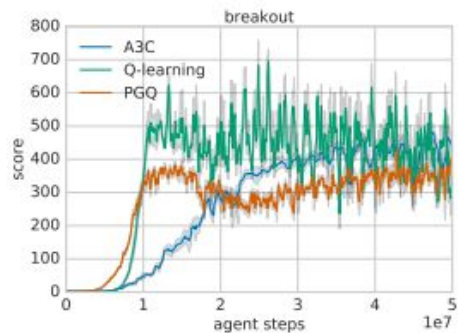


Figure 3: Some Atari runs where PGQL performed well.

- **A3C**: only do actor-critic update.
- **Q-learning**: only update by Q-value
- **PGQL**: do update by both A3C method and Q value.
- **Big Note**: The baseline Q-learning is not the typical DQN or else. **It is the Q-learning with Q as their provided variant.**
- PGQL **outperforms** A3C and Q-learning in these Atari games.



# Experiment Results(2): Paper Part



- But in some other Atari games, PGQL didn't have a good performance.

# Experiment Results(2): Replication Part

- We only choose two environments “**Assault**” and “**Breakout**” to conduct the experiments, since it is time-consuming to train models.
- Surprisingly, we see that the performance of PGQL on Assault is similar to A3C, and the performance of PGQL on Breakout is much better than the others, **This is different from what we see in the paper results.**

