

Federated Reinforcement Learning: Linear Speedup Under Markovian Sampling

Sajad Khodadadian
June, 2022



Pranay Sharma (CMU)

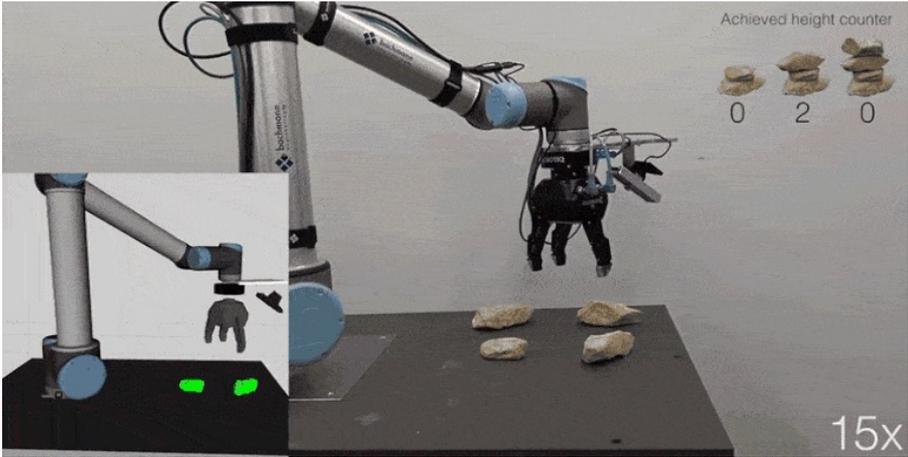


Gauri Joshi (CMU)



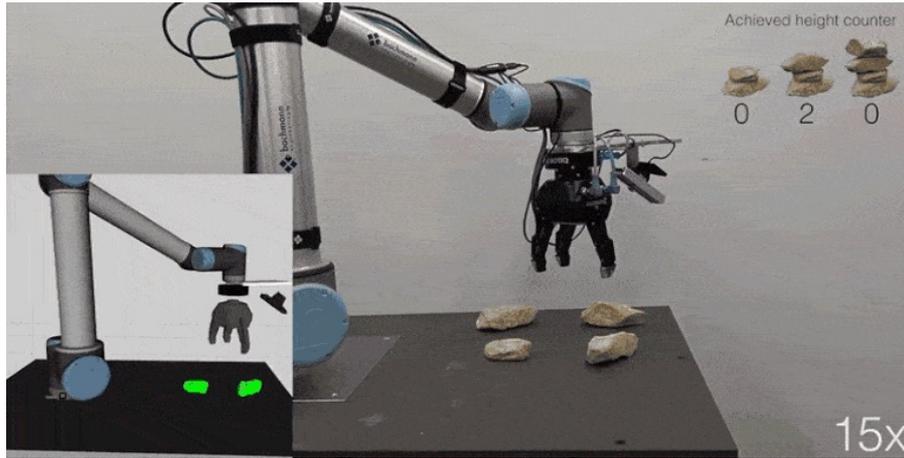
Siva Theja Maguluri (Gatech)

Reinforcement Learning

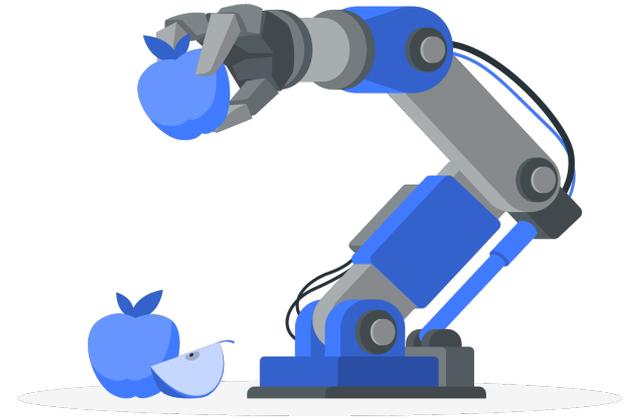


freecodecamp.org

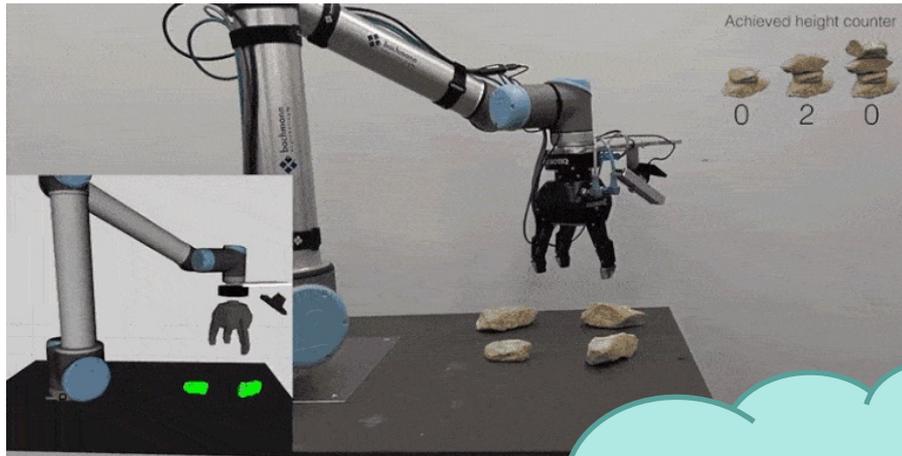
Reinforcement Learning



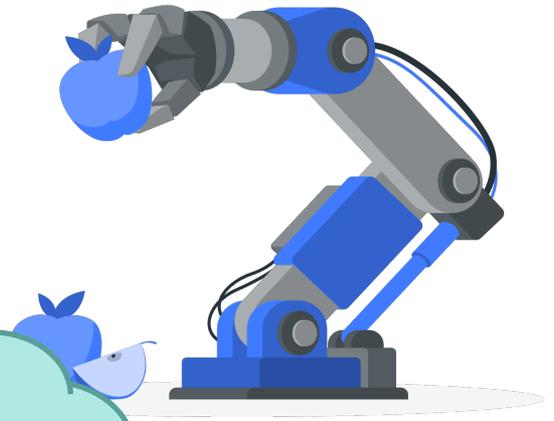
freecodecamp.org



Reinforcement Learning



freecodecamp.org



But RL is data
intensive!

Federated Reinforcement Learning

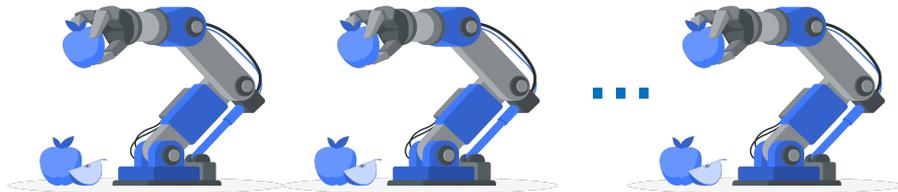


Google AI Blog

Federated Reinforcement Learning



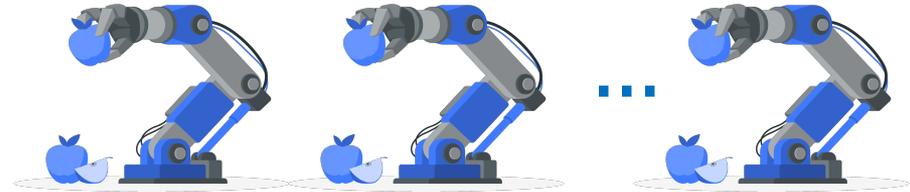
Google AI Blog



Federated Reinforcement Learning



Google AI Blog



Multiple data collecting agents

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

$$Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

$$Q^\pi(\mathbf{s}, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

initial state

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

$$Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

initial state initial action

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

$$Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

initial state initial action policy

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

$$Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

The diagram shows the equation for the Q-function with four labels and arrows pointing to specific parts of the equation:

- initial state**: points to the variable s in the function arguments.
- initial action**: points to the variable a in the function arguments.
- state at time t** : points to the variable S_t inside the expectation operator.
- policy**: points to the symbol π inside the expectation operator.

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

The diagram shows the Q-function equation: $Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$. The equation is enclosed in a light blue rounded rectangle. Five teal arrows point from labels to specific parts of the equation: 'initial state' points to s , 'initial action' points to a , 'state at time t ' points to S_t , 'action at time t ' points to A_t , and 'policy' points to $\pi(\cdot \mid S_t)$.

$$Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

initial state initial action state at time t action at time t policy

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

The diagram shows the Q-function equation: $Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$. The equation is enclosed in a light blue rounded rectangle. Labels with arrows point to various parts of the equation: 'initial state' points to s , 'initial action' points to a , 'action at time t ' points to A_t , 'policy' points to $\pi(\cdot \mid S_t)$, 'state at time t ' points to S_t , and 'Reward function' points to $\mathcal{R}(S_t, A_t)$.

$$Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

initial state initial action action at time t policy

Reward function state at time t

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

The diagram shows the Q-function equation: $Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$. The equation is enclosed in a light blue rounded rectangle. Labels with arrows point to various parts of the equation: 'initial state' points to s ; 'initial action' points to a ; 'discount factor' points to γ ; 'action at time t ' points to A_t ; 'policy' points to π ; 'state at time t ' points to S_t ; and 'Reward function' points to \mathcal{R} .

$$Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$$

initial state initial action discount factor action at time t policy

Reward function state at time t

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

The diagram shows the Q-function equation: $Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$. Labels with arrows point to various parts of the equation: 'initial state' points to s ; 'initial action' points to a ; 'discount factor' points to γ ; 'action at time t ' points to A_t ; 'policy' points to π ; 'state at time t ' points to S_t ; and 'Reward function' points to \mathcal{R} .

- Optimal policy

$$\pi^* \in \underset{\pi}{\operatorname{argmax}} Q^\pi(s, a), \quad \forall s, a$$

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

The diagram shows the equation for the Q -function: $Q^\pi(s, a) = \mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\}$. Labels with arrows point to various parts of the equation: 'initial state' points to s ; 'initial action' points to a ; 'discount factor' points to γ ; 'action at time t ' points to A_t ; 'policy' points to $\pi(\cdot \mid S_t)$; 'state at time t ' points to S_t ; and 'Reward function' points to $\mathcal{R}(S_t, A_t)$.

- Optimal policy

$$\pi^* \in \underset{\pi}{\operatorname{argmax}} Q^\pi(s, a), \quad \forall s, a$$

- Optimal Q -function

$$Q^*(s, a) \equiv Q^{\pi^*}(s, a)$$

Background on MDP Theory

- Discounted Markov Decision Process (MDP)
- Q -function

$$Q^\pi(s, a) = \mathbb{E}\left\{\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot \mid S_t)\right\}$$

The diagram shows the equation for the Q-function with several labels and arrows pointing to specific parts of the equation:

- initial state**: points to s
- initial action**: points to a
- discount factor**: points to γ
- action at time t** : points to A_t
- policy**: points to $\pi(\cdot \mid S_t)$
- Reward function**: points to $\mathcal{R}(S_t, A_t)$
- state at time t** : points to S_t

- Optimal policy

$$\pi^* \in \underset{\pi}{\operatorname{argmax}} Q^\pi(s, a), \quad \forall s, a$$

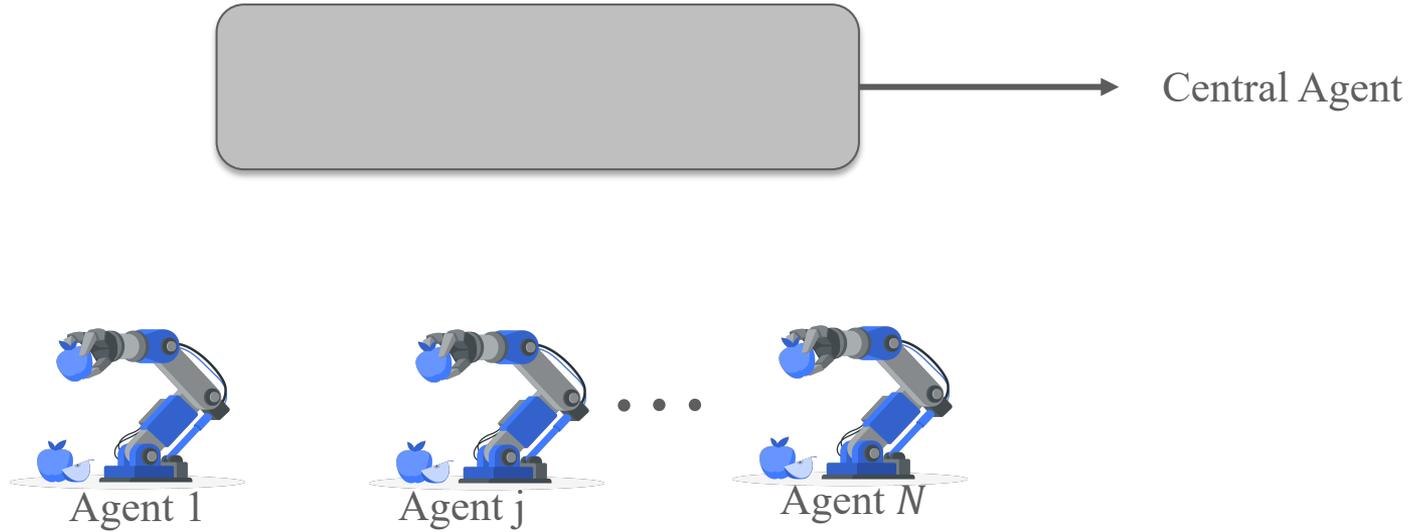
- Optimal Q -function

$$Q^*(s, a) \equiv Q^{\pi^*}(s, a) \rightarrow Q\text{-learning}$$

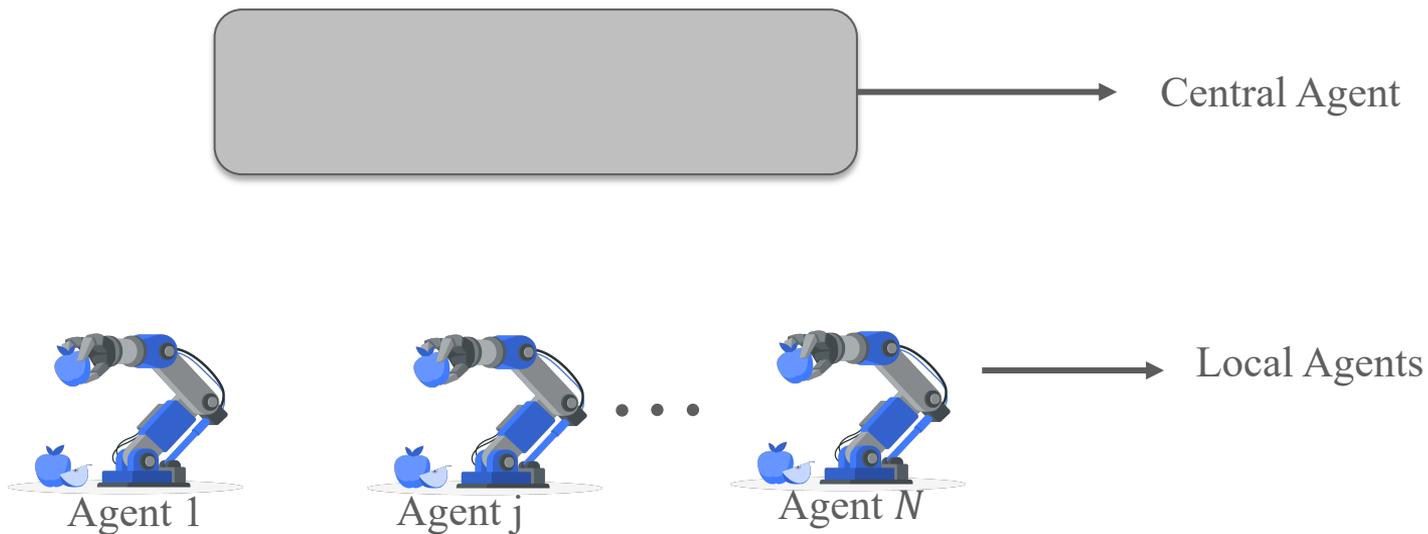
Vanilla Distributed Reinforcement Learning



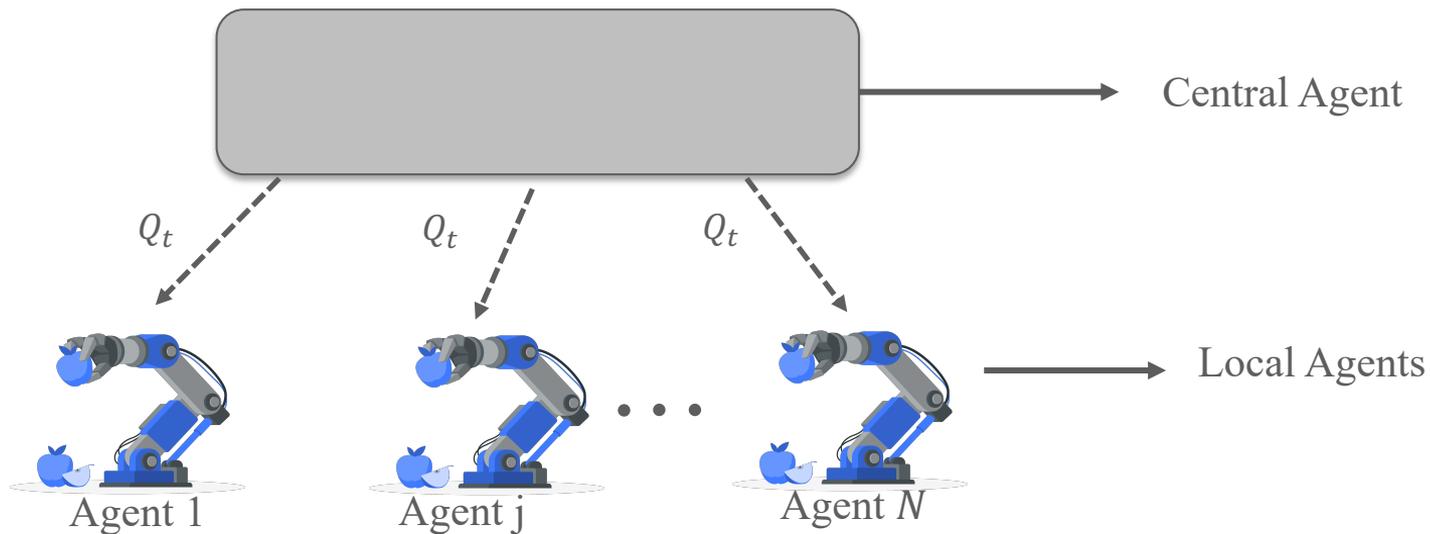
Vanilla Distributed Reinforcement Learning



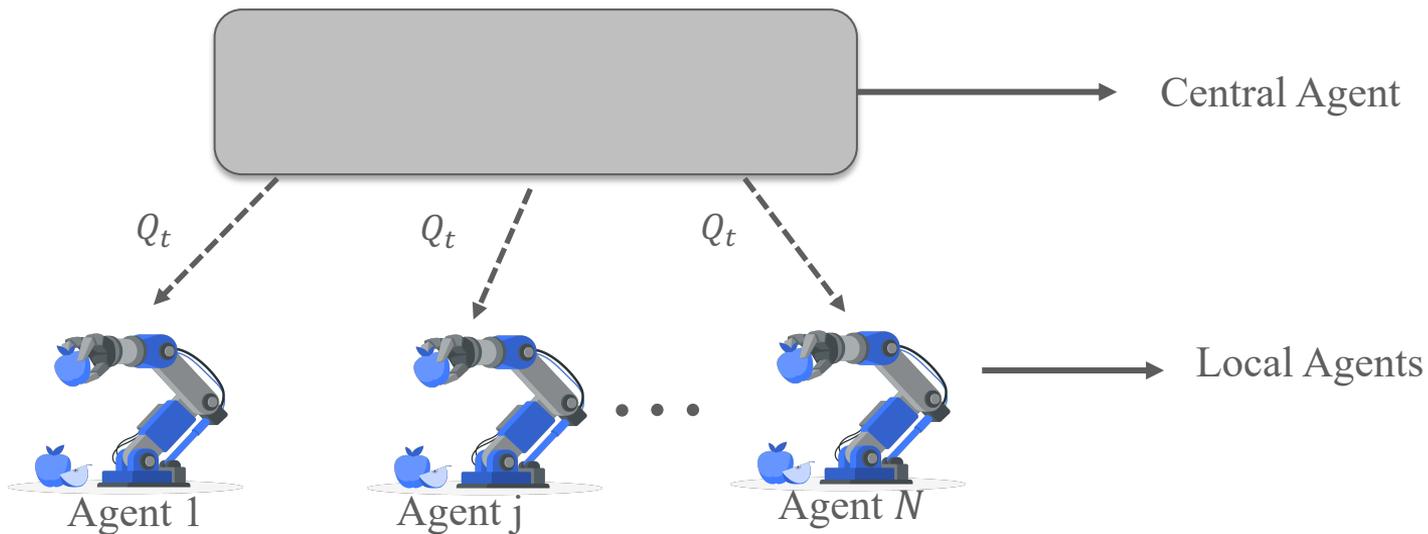
Vanilla Distributed Reinforcement Learning



Vanilla Distributed Reinforcement Learning

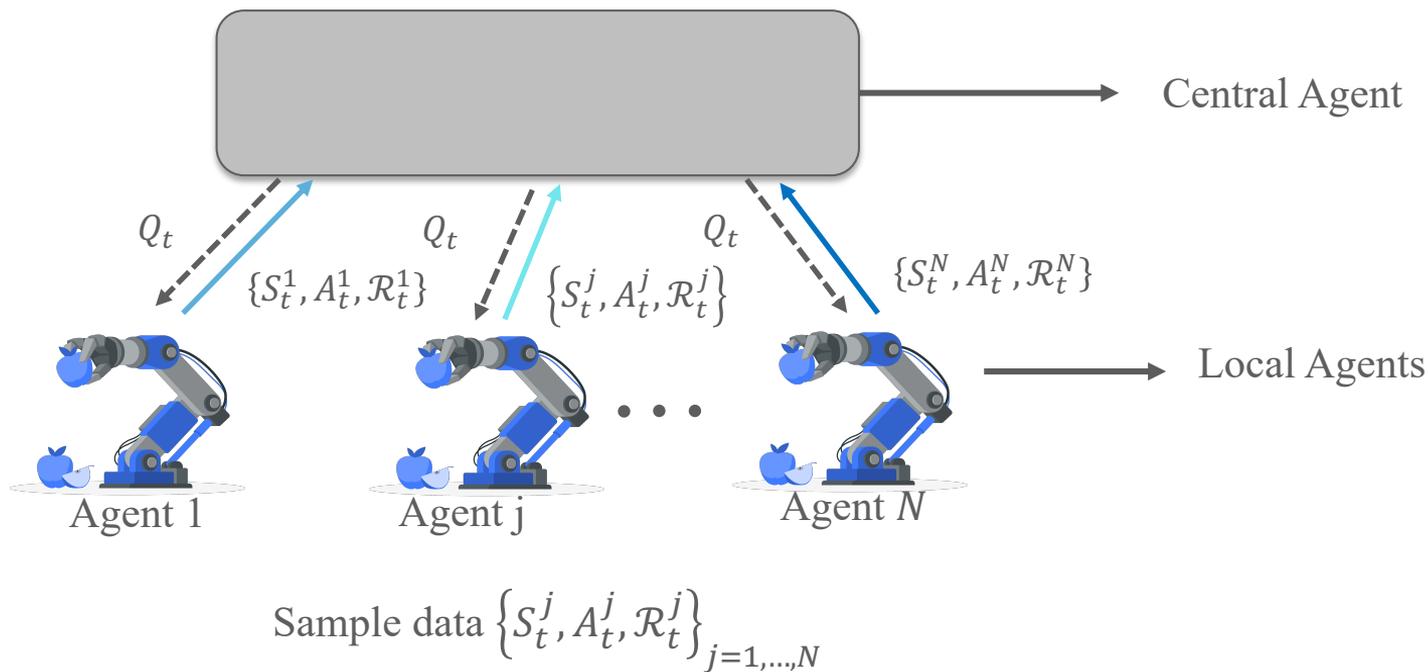


Vanilla Distributed Reinforcement Learning

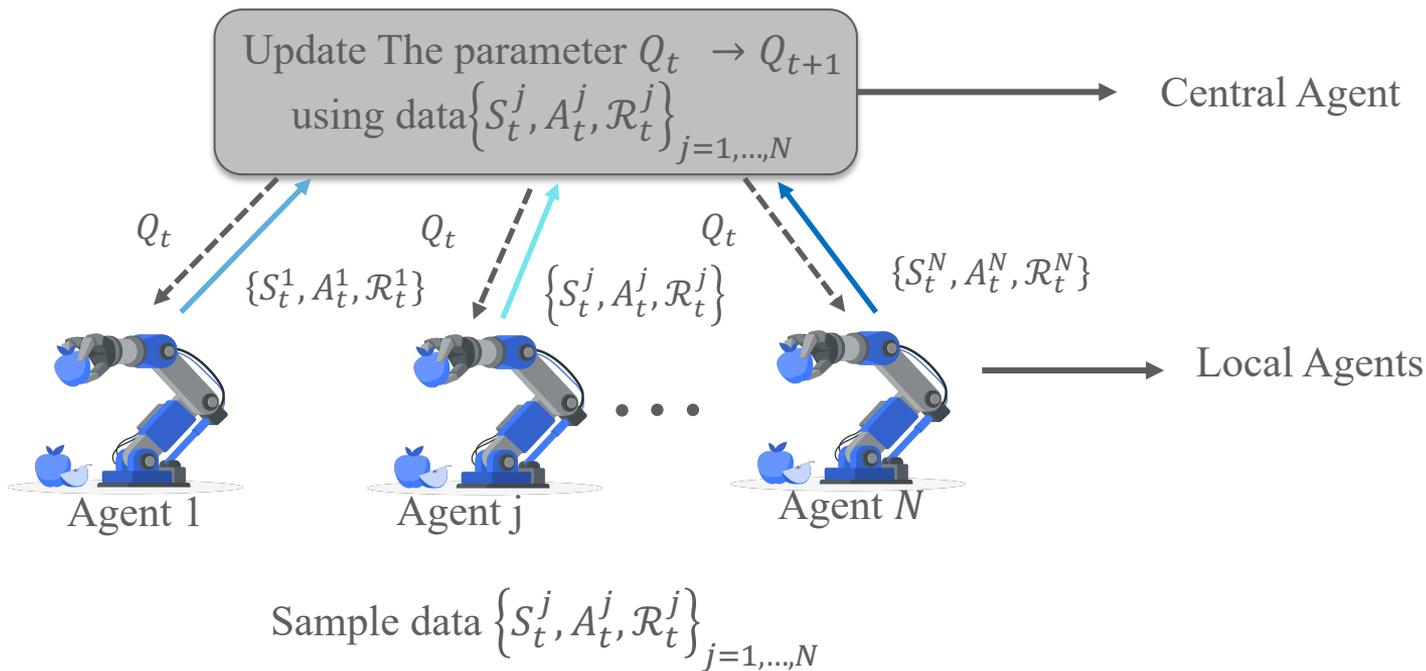


$$\text{Sample data } \{S_t^j, A_t^j, \mathcal{R}_t^j\}_{j=1, \dots, N}$$

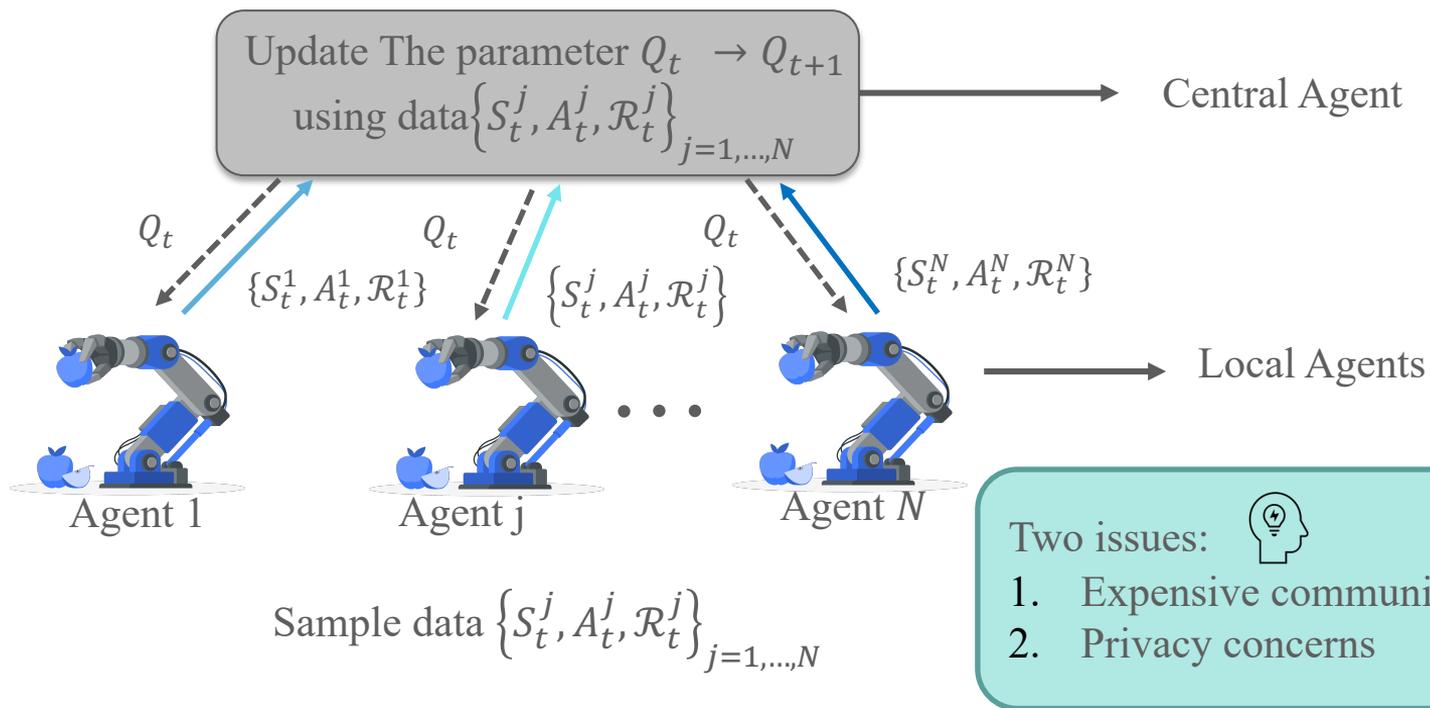
Vanilla Distributed Reinforcement Learning



Vanilla Distributed Reinforcement Learning



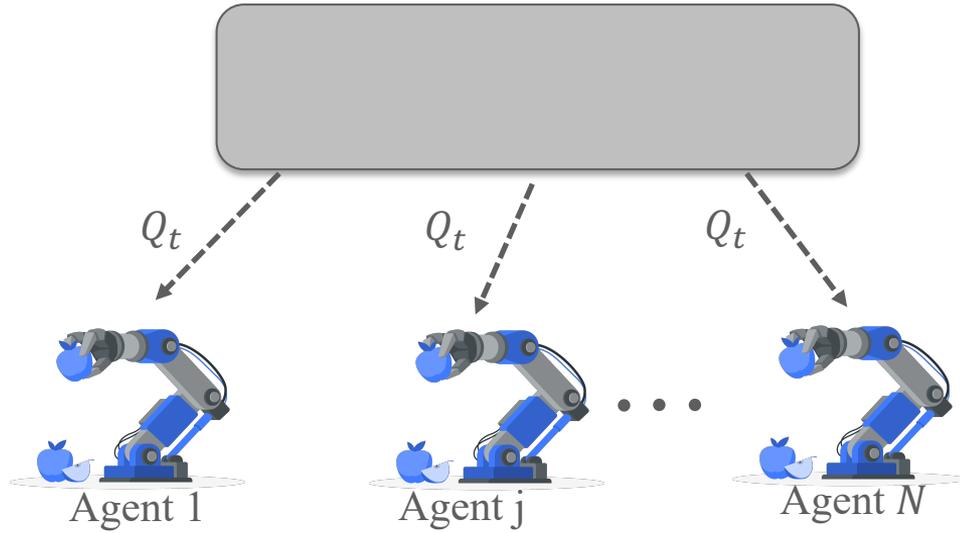
Vanilla Distributed Reinforcement Learning



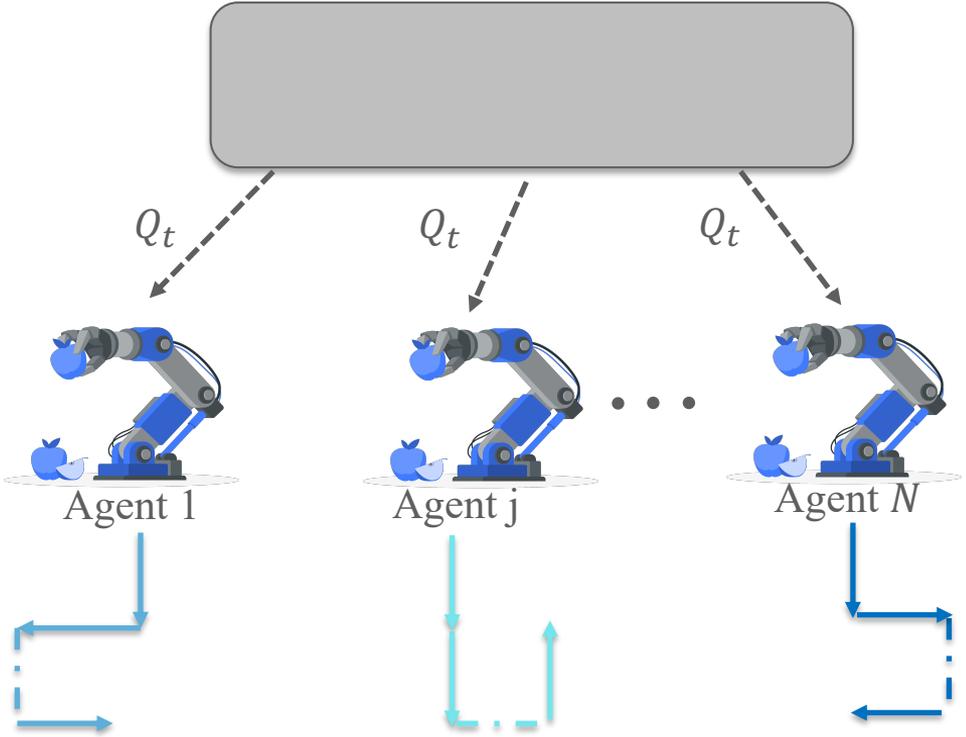
Federated Reinforcement Learning



Federated Reinforcement Learning

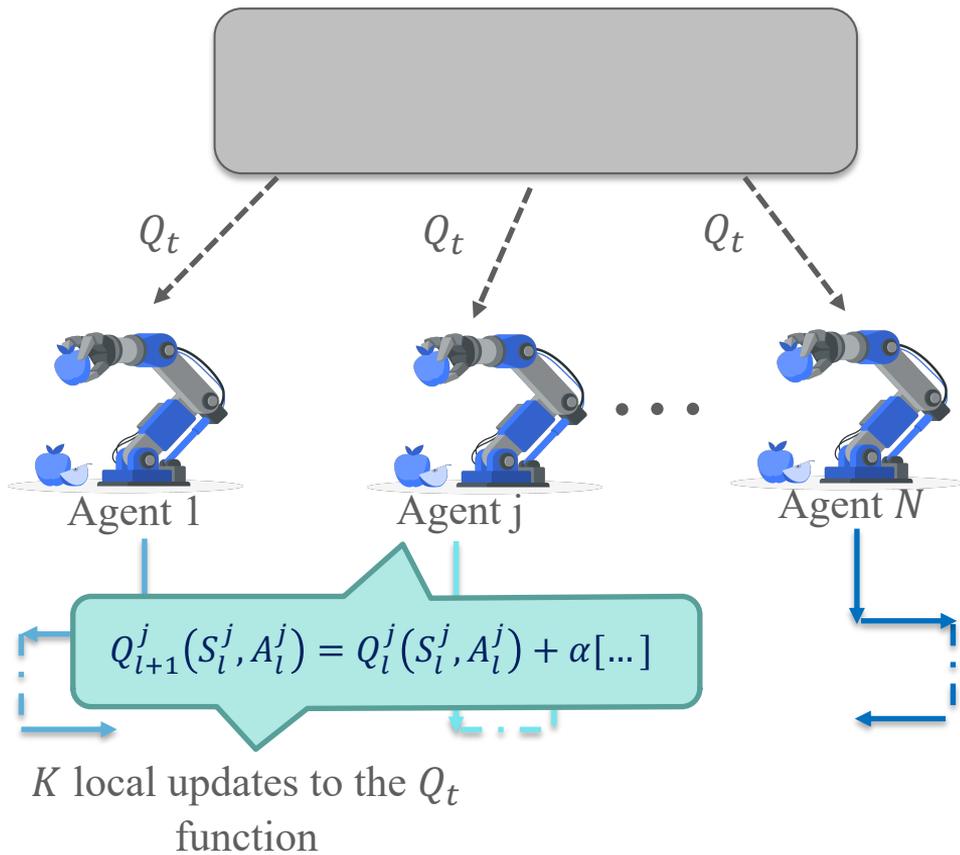


Federated Reinforcement Learning

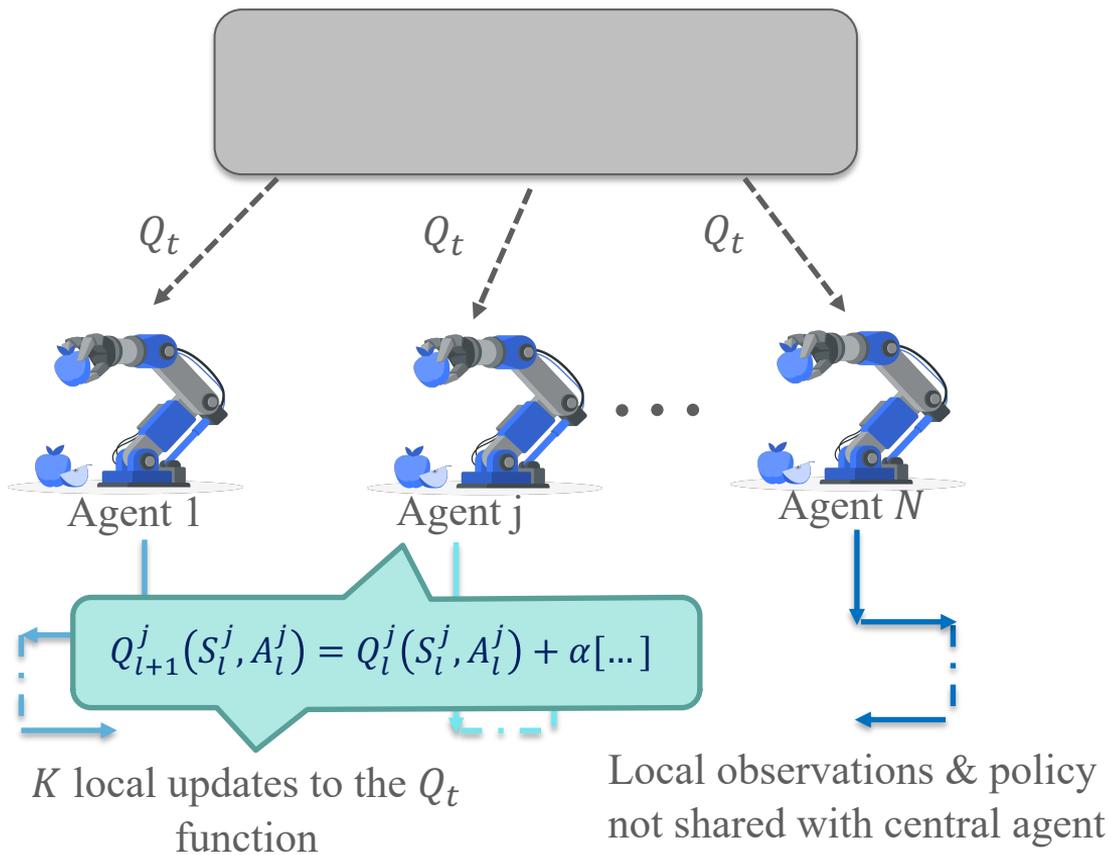


K local updates to the Q_t function

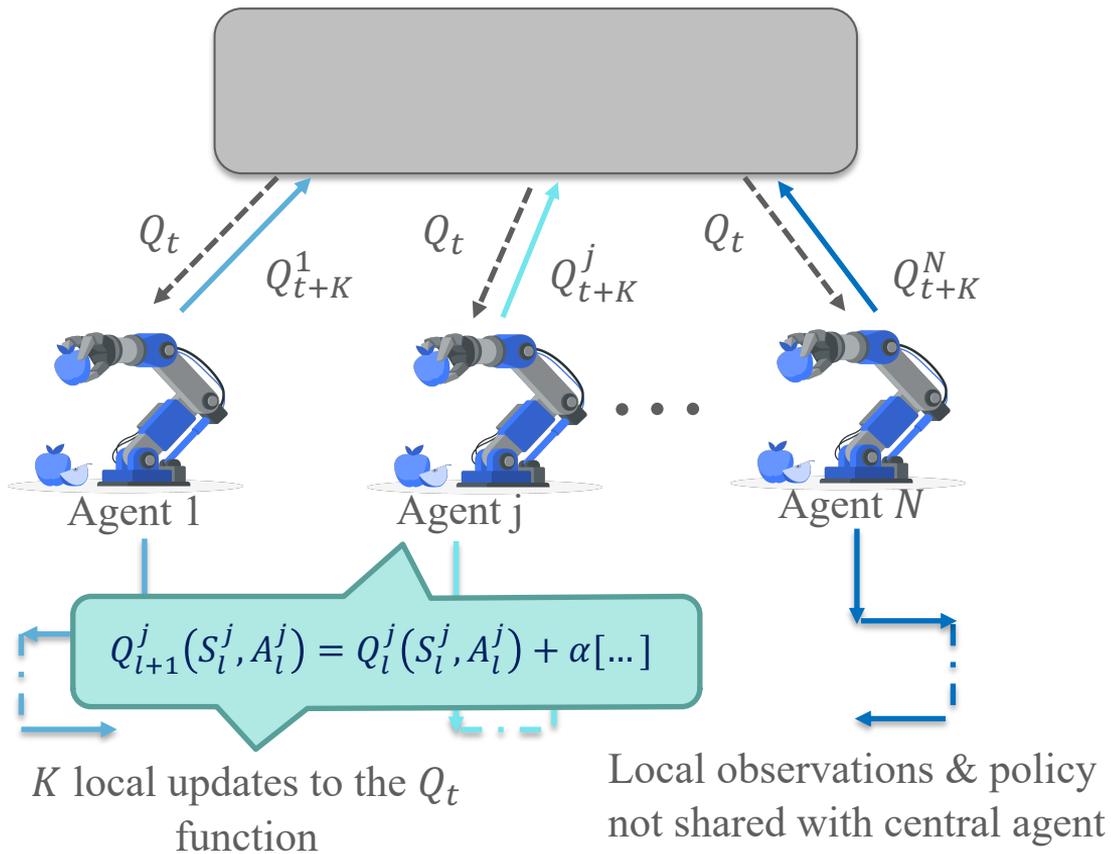
Federated Reinforcement Learning



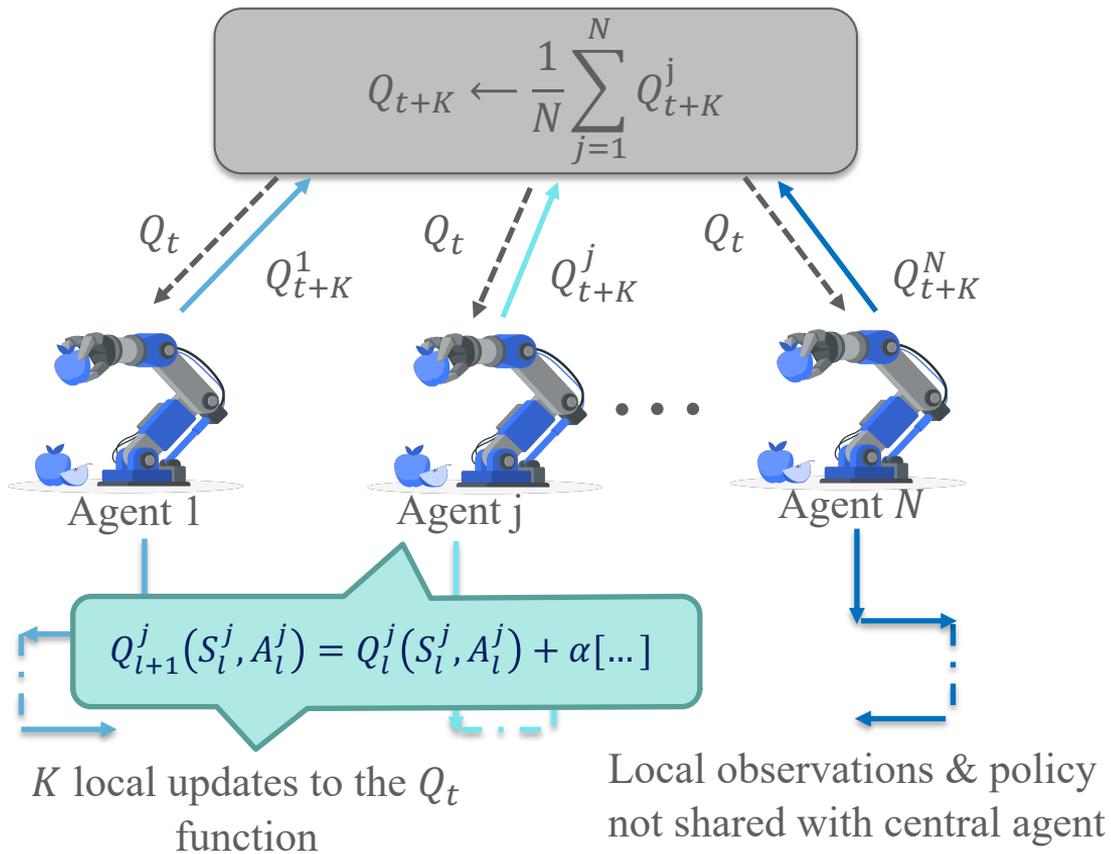
Federated Reinforcement Learning



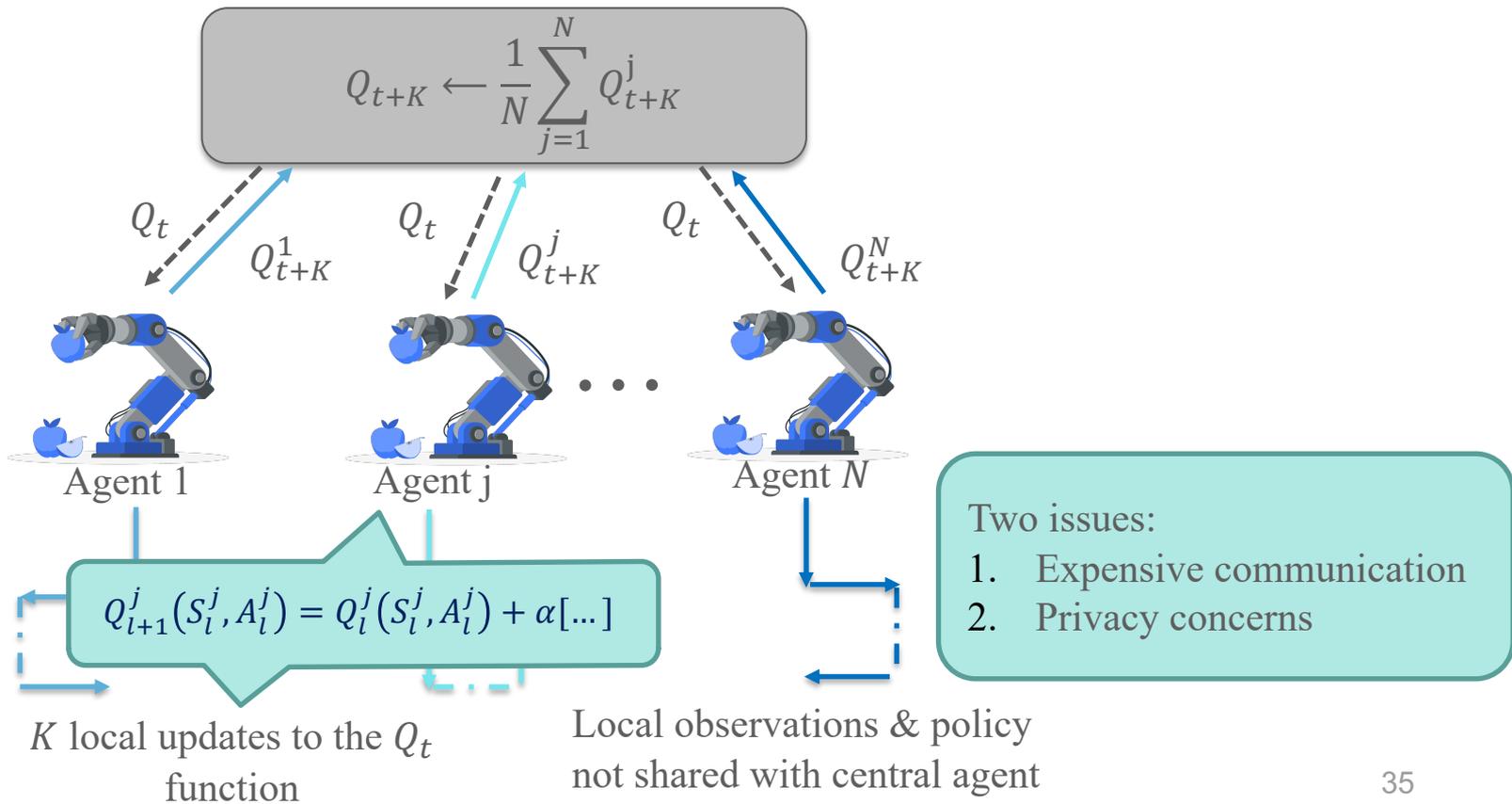
Federated Reinforcement Learning



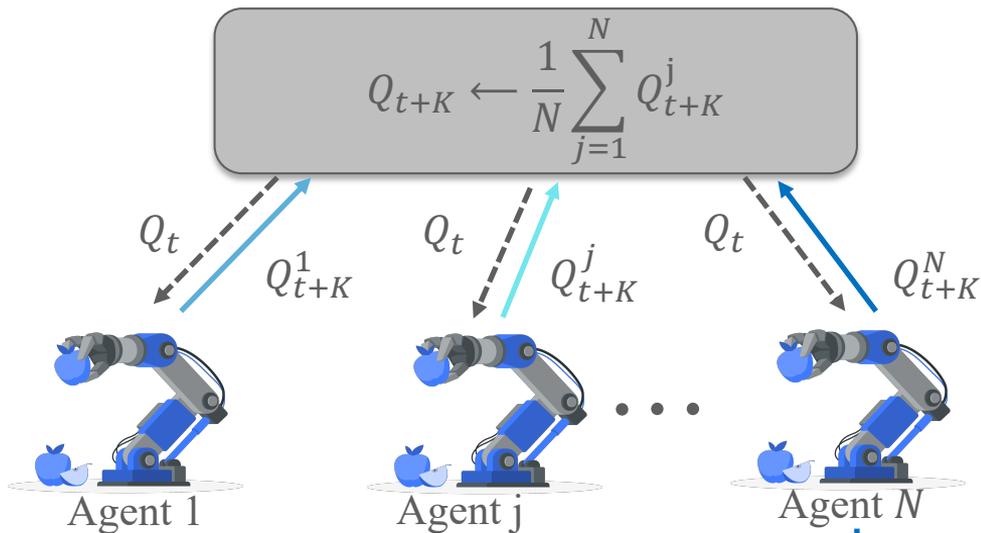
Federated Reinforcement Learning



Federated Reinforcement Learning



Federated Reinforcement Learning



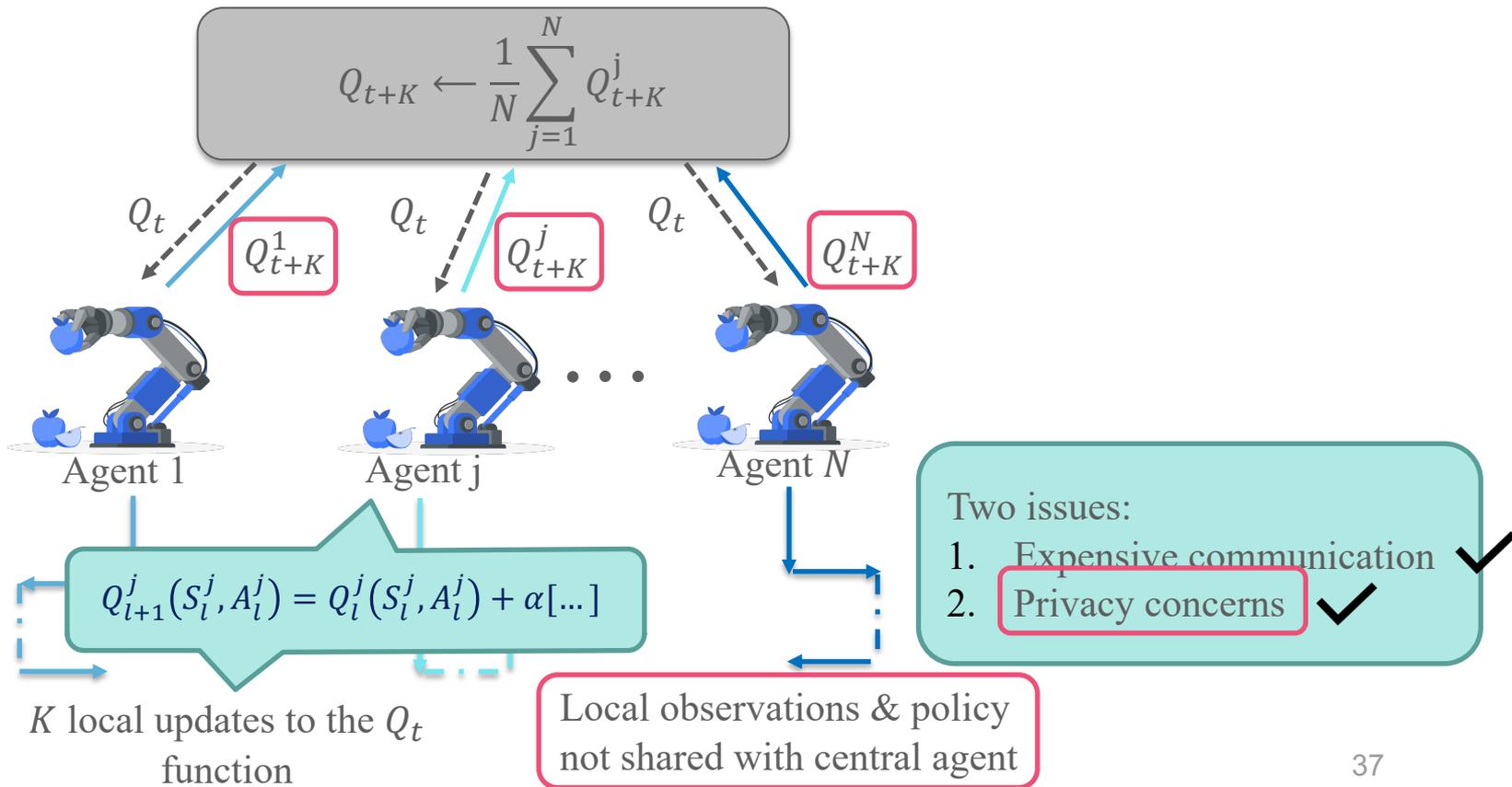
$$Q_{i+1}^j(S_i^j, A_i^j) = Q_i^j(S_i^j, A_i^j) + \alpha[\dots]$$

K local updates to the Q_t function

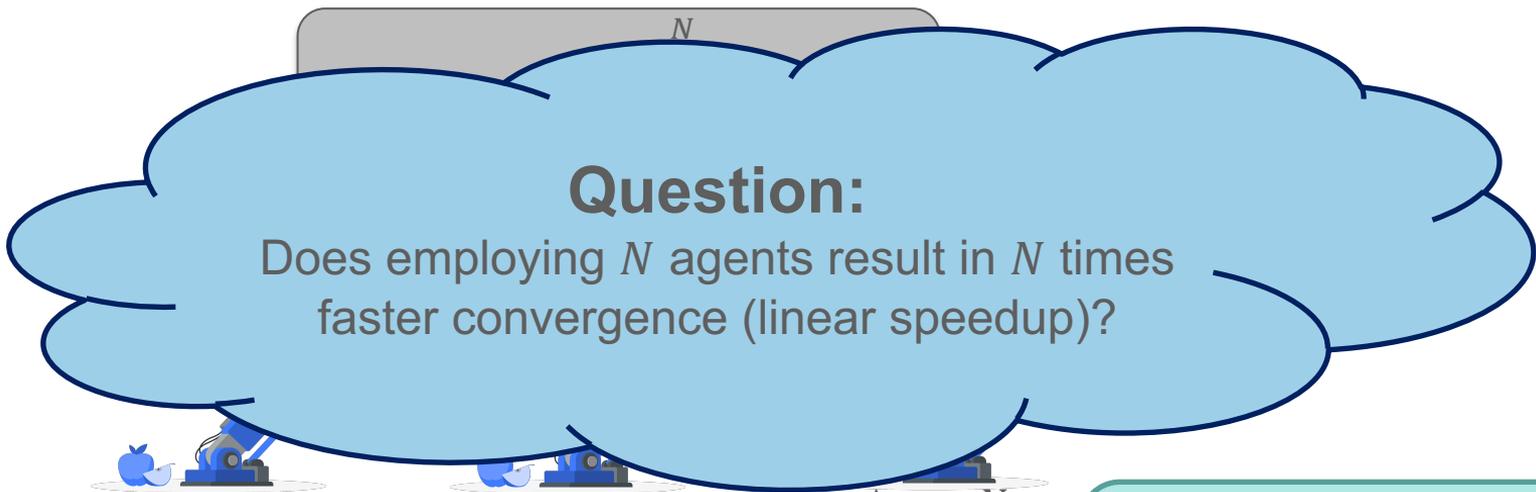
Local observations & policy not shared with central agent

- Two issues:
1. Expensive communication ✓
 2. Privacy concerns

Federated Reinforcement Learning



Federated Reinforcement Learning



$$Q_{i+1}^j(S_i^j, A_i^j) = Q_i^j(S_i^j, A_i^j) + \alpha[\dots]$$

K local updates to the Q_t function

Local observations & policy not shared with central agent

- Two issues:
1. Expensive communication ✓
 2. Privacy concerns ✓

Linear Speedup in Federated Learning

- Federated Supervised Learning:
 1. Linear speedup is possible [Spiridonoff, Olshevsky, Paschalidis, NeurIPS '21], [...]

Linear Speedup in Federated Learning

- Federated Supervised Learning:
 1. Linear speedup is possible [Spiridonoff, Olshevsky, Paschalidis, NeurIPS '21], [...]
 2. Key ingredient in these results: The noise is i.i.d.

Linear Speedup in Federated Learning

- Federated Supervised Learning:
 1. Linear speedup is possible [Spiridonoff, Olshevsky, Paschalidis, NeurIPS '21], [...]
 2. Key ingredient in these results: The noise is i.i.d.

$$X_1, X_2, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} F_X(\cdot) \quad , \quad \text{Var}(X_i) = \sigma^2$$

Linear Speedup in Federated Learning

- Federated Supervised Learning:
 1. Linear speedup is possible [Spiridonoff, Olshevsky, Paschalidis, NeurIPS '21], [...]
 2. Key ingredient in these results: The noise is i.i.d.

$$X_1, X_2, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} F_X(\cdot) \quad , \quad \text{Var}(X_i) = \sigma^2$$

$$\Rightarrow \text{Var}\left(\frac{\sum_{i=1}^N X_i}{N}\right) = \frac{\sigma^2}{N}$$

Linear Speedup in Federated Learning

- Federated Supervised Learning:
 1. Linear speedup is possible [Spiridonoff, Olshevsky, Paschalidis, NeurIPS '21], [...]
 2. Key ingredient in these results: The noise is i.i.d.

$$X_1, X_2, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} F_X(\cdot) \quad , \quad \text{Var}(X_i) = \sigma^2$$

$$\Rightarrow \text{Var}\left(\frac{\sum_{i=1}^N X_i}{N}\right) = \frac{\sigma^2}{N} \longrightarrow \text{This is the source of linear speedup}$$

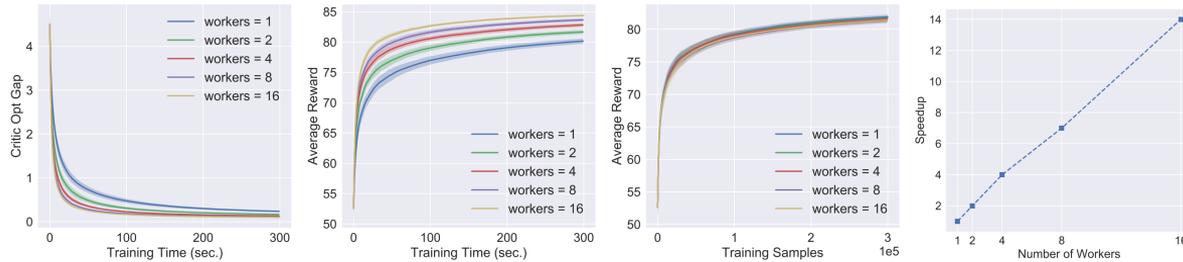
Linear Speedup in Federated Learning

- Federated RL (TD) algorithms
 1. No linear speedup [Wai '20] [Zeng, Doan, Romberg, '20]
 - In fact, they have linear penalty – but their focus is different

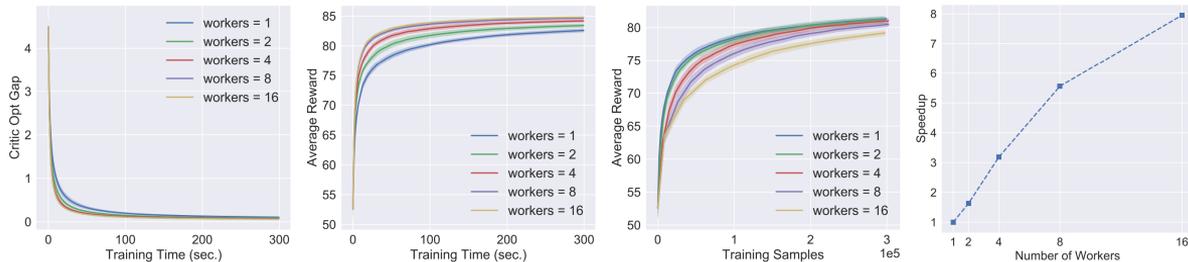
Linear Speedup in Federated Learning

- Federated RL (TD) algorithms
 1. No linear speedup [Wai '20] [Zeng, Doan, Romberg, '20]
 - ❑ In fact, they have linear penalty – but their focus is different
 2. Linear speed up under i.i.d. noise assumption [Shen, Zhang, Hong, Chen '20]
 - ❑ Based on experiments, conjectured that linear speedup is possible under Markov noise too

Linear Speedup in Federated Learning (A3C)¹



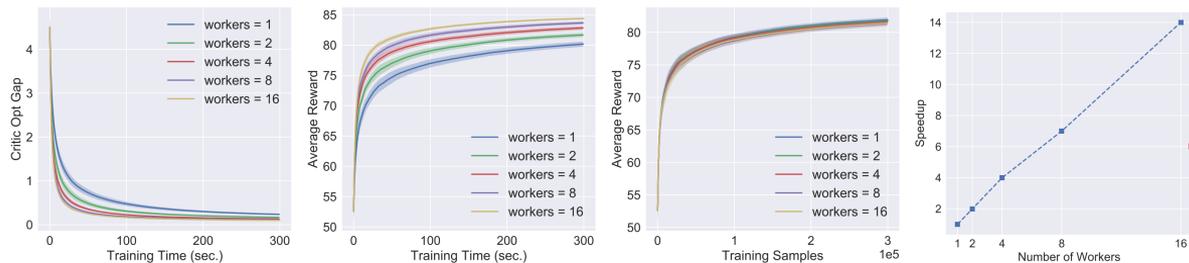
Convergence results of A3C with i.i.d. sampling in synthetic environment.



Convergence results of A3C with Markovian sampling in synthetic environment.

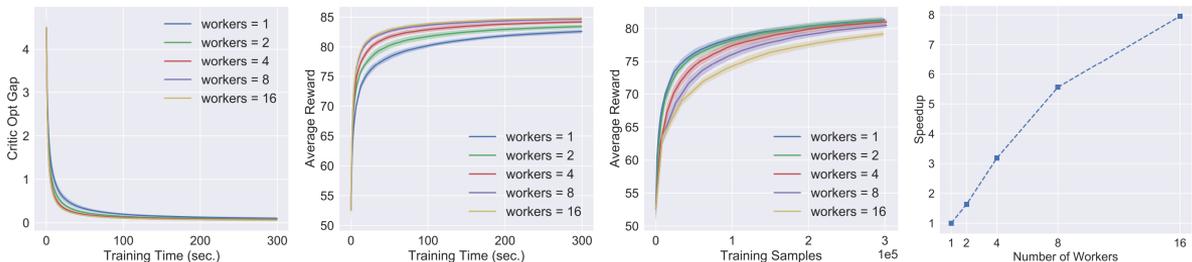
¹Shen, Han, et al. "Asynchronous advantage actor critic: Non-asymptotic analysis and linear speedup." arXiv preprint arXiv:2012.15511 (2020).

Linear Speedup in Federated Learning (A3C)¹



Linear speedup is established for i.i.d. noise

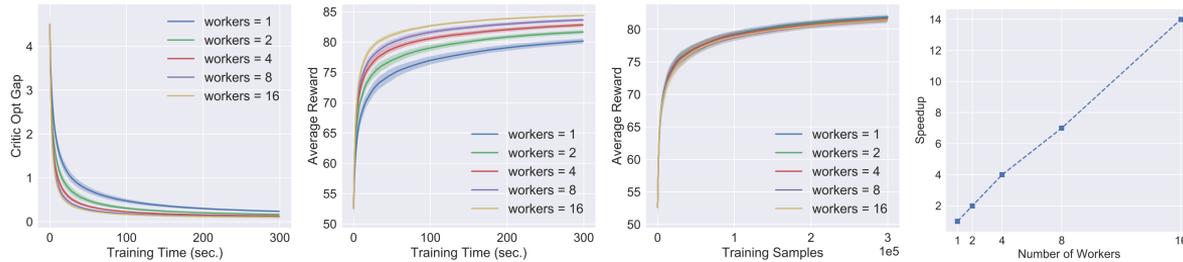
Convergence results of A3C with **i.i.d.** sampling in synthetic environment.



Convergence results of A3C with Markovian sampling in synthetic environment.

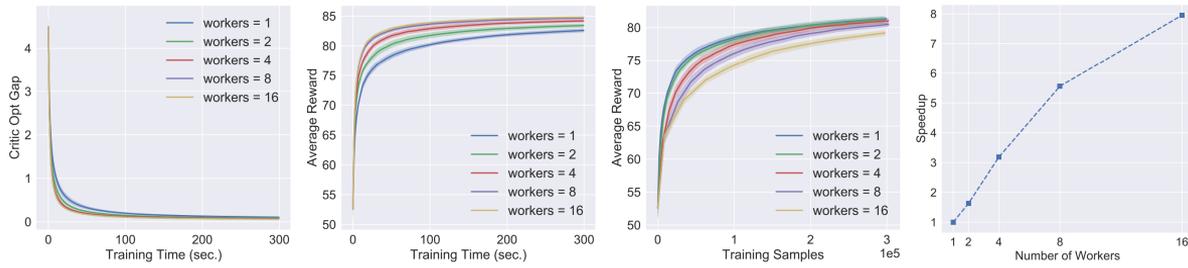
¹Shen, Han, et al. "Asynchronous advantage actor critic: Non-asymptotic analysis and linear speedup." arXiv preprint arXiv:2012.15511 (2020).

Linear Speedup in Federated Learning (A3C)¹



Linear speedup is established for i.i.d. noise

Convergence results of A3C with **i.i.d.** sampling in synthetic environment.

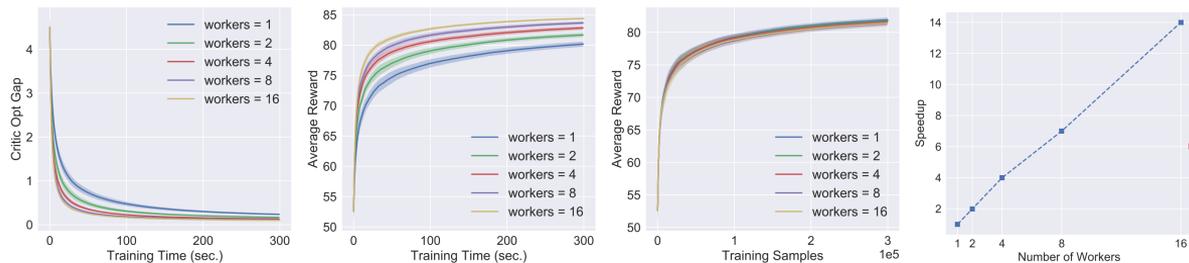


A3C paper do not prove a linear speedup in the Markovian setting

Convergence results of A3C with **Markovian** sampling in synthetic environment.

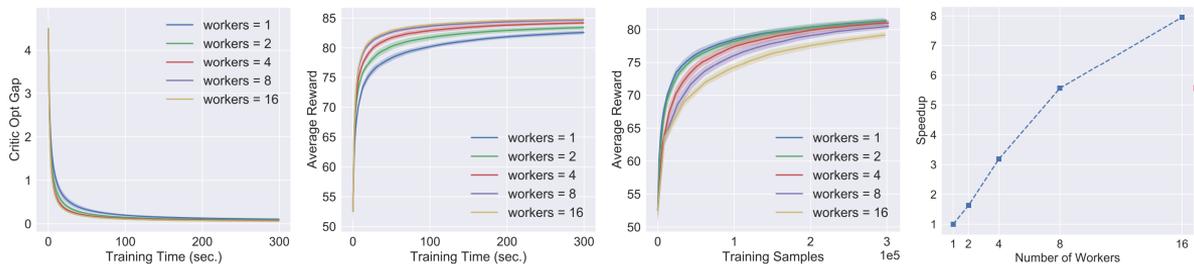
¹Shen, Han, et al. "Asynchronous advantage actor critic: Non-asymptotic analysis and linear speedup." arXiv preprint arXiv:2012.15511 (2020).

Linear Speedup in Federated Learning (A3C)¹



Linear speedup is established for i.i.d. noise

Convergence results of A3C with **i.i.d.** sampling in synthetic environment.



A3C paper do not prove a linear speedup in the Markovian setting

Convergence results of A3C with **Markovian** sampling in synthetic environment.

We are the first to prove this

¹Shen, Han, et al. "Asynchronous advantage actor critic: Non-asymptotic analysis and linear speedup." arXiv preprint arXiv:2012.15511 (2020).

Federated Q -learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O}\left(\frac{1}{\alpha}(1 - c\alpha)^T + \frac{\alpha}{N} + (K - 1)\alpha^2\right).$$

Federated Q -learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - C\alpha)^T}_{\substack{\text{Convergence} \\ \text{Bias}}} + \frac{\alpha}{N} + (K - 1)\alpha^2 \right).$$

Federated Q -learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - C\alpha)^T}_{\text{Convergence Bias}} + \underbrace{\frac{\alpha}{N} + (K - 1)\alpha^2}_{\text{Convergence Variance}} \right).$$

Federated Q -learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - c\alpha)^T}_{\text{Convergence Bias}} + \underbrace{\frac{\alpha}{N}}_{\text{Convergence Variance}} + (K - 1)\alpha^2 \right).$$

Convergence Bias Convergence Variance

Linear speedup

Federated Q -learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - c\alpha)^T}_{\text{Convergence Bias}} + \underbrace{\frac{\alpha}{N}}_{\text{Convergence Variance}} + \underbrace{(K-1)\alpha^2}_{\text{Higher order}} \right).$$

Convergence Bias Convergence Variance

Linear speedup

Federated Q-learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - c\alpha)^T}_{\text{Convergence Bias}} + \underbrace{\frac{\alpha}{N}}_{\text{Convergence}} + \underbrace{(K-1)\alpha^2}_{\text{Variance}} \right).$$

Convergence Bias Convergence Variance

Linear speedup

Higher order

- If $\alpha = \mathcal{O}(\log(NT) / T)$ and $K = T/N$, we have $\mathbb{E}[\|Q_T - Q^\pi\|_\infty^2] \leq \epsilon$ within $T = \tilde{O}\left(\frac{1}{N\epsilon}\right)$ iterations.

Federated Q-learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - C\alpha)^T}_{\text{Convergence Bias}} + \underbrace{\frac{\alpha}{N}}_{\text{Convergence}} + \underbrace{(K-1)\alpha^2}_{\text{Variance}} \right).$$

Convergence Bias Convergence Variance

Linear speedup

Higher order

- If $\alpha = \mathcal{O}(\log(NT) / T)$ and $K = T/N$, we have $\mathbb{E}[\|Q_T - Q^\pi\|_\infty^2] \leq \epsilon$ within $T = \tilde{O}\left(\frac{1}{N\epsilon}\right)$ iterations.
- Total communication cost = $\frac{T}{K} = N$

Federated Q-learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - c\alpha)^T}_{\text{Convergence Bias}} + \underbrace{\frac{\alpha}{N} + (K-1)\alpha^2}_{\text{Convergence Variance}} \right).$$

Convergence Bias Convergence Variance

Linear speedup

Higher order

- If $\alpha = \mathcal{O}(\log(NT) / T)$ and $K = T/N$, we have $\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \epsilon$ within $T = \tilde{O}\left(\frac{1}{N\epsilon}\right)$ iterations.
- Total communication cost = $\frac{T}{K} = N$

Number of agents

Federated Q-learning

Theorem: Let $Q_T = \frac{1}{N} \sum_{i=1}^N Q_T^i$,

$$\mathbb{E}[\|Q_T - Q^*\|_\infty^2] \leq \tilde{O} \left(\underbrace{\frac{1}{\alpha} (1 - c\alpha)^T}_{\text{Convergence Bias}} + \underbrace{\frac{\alpha}{N}}_{\text{Convergence}} + \underbrace{(K-1)\alpha^2}_{\text{Variance}} \right).$$

Convergence Bias Convergence Variance

Linear speedup

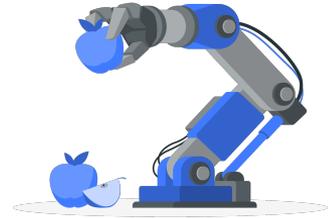
Higher order

- If $\alpha = \mathcal{O}(\log(NT)/T)$ and $K = T/N$, we have $\mathbb{E}[\|Q_T - Q^*\|_\infty] \leq \epsilon$ within $T = \tilde{O}\left(\frac{1}{N\epsilon}\right)$ iterations.
- Total communication cost = $\frac{T}{K} = N$

Number of agents \longrightarrow Constant

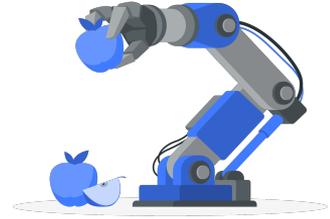
Proof sketch

- Single agent setting



Proof sketch

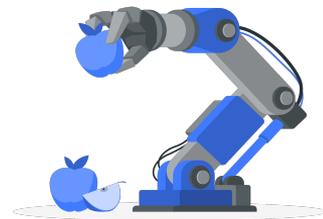
- Single agent setting
- Lyapunov type argument:



Proof sketch

- Single agent setting
- Lyapunov type argument:

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \sigma^2\alpha^2$$

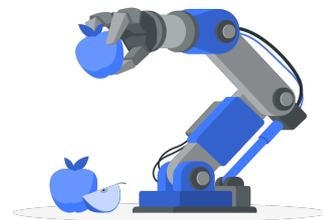


Proof sketch

- Single agent setting
- Lyapunov type argument:

Correspond to
variance

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \overbrace{\sigma^2\alpha^2}$$



Proof sketch

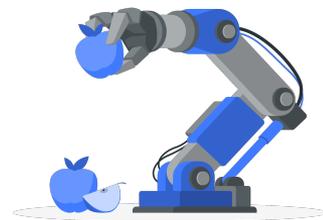
- Single agent setting
- Lyapunov type argument:

Correspond to
variance

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \overbrace{\sigma^2\alpha^2}^{\text{variance}}$$



$$\mathbb{E}[\|\theta_T\|^2] \leq (1 - \alpha)^T \|\theta_0\|^2 + \alpha$$



Proof sketch

- Single agent setting
- Lyapunov type argument:

Correspond to
variance

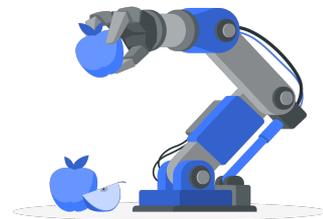
$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \overbrace{\sigma^2\alpha^2}^{\text{variance}}$$



$$\mathbb{E}[\|\theta_T\|^2] \leq (1 - \alpha)^T \|\theta_0\|^2 + \alpha$$



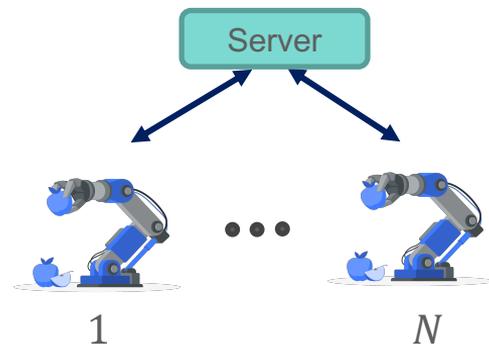
$\tilde{O}(1/\epsilon)$ sample complexity



Proof sketch

- Multiple agents, favorable recursion

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$



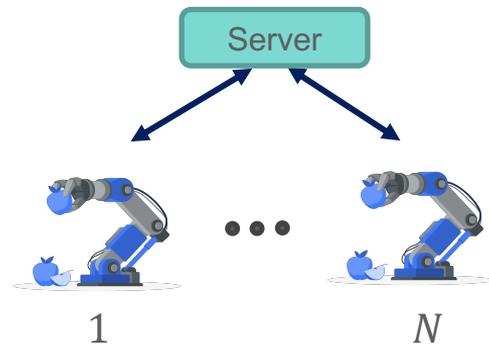
Proof sketch

- Multiple agents, favorable recursion

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$



$$\mathbb{E}[\|\theta_T\|^2] \leq (1 - \alpha)^T \|\theta_0\|^2 + \alpha/N$$



Proof sketch

- Multiple agents, favorable recursion

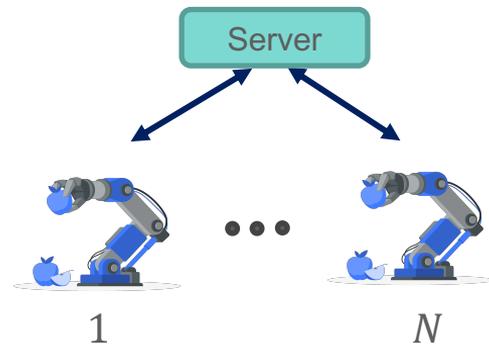
$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$



$$\mathbb{E}[\|\theta_T\|^2] \leq (1 - \alpha)^T \|\theta_0\|^2 + \alpha/N$$



$\tilde{O}(1/N\epsilon)$ iteration complexity, linear speedup



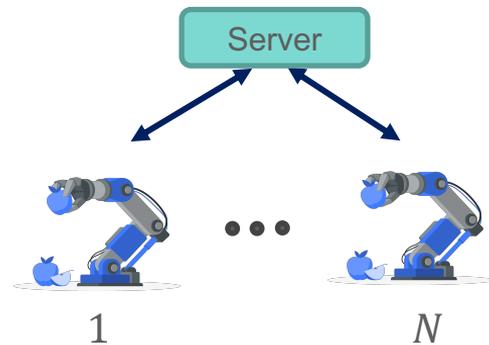
Proof sketch

- Multiple agents, favorable recursion

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$

- However, we get

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \frac{\alpha^2}{N} + \alpha^3 + \Omega_t$$



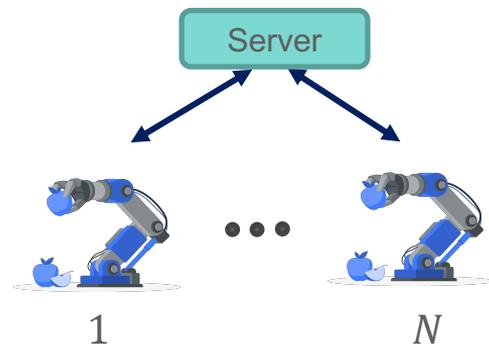
Proof sketch

- Multiple agents, favorable recursion

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$

- However, we get

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \frac{\alpha^2}{N} + \underbrace{\alpha^3 + \Omega_t}_{\text{Higher order terms}}$$



Proof sketch

- Multiple agents, favorable recursion

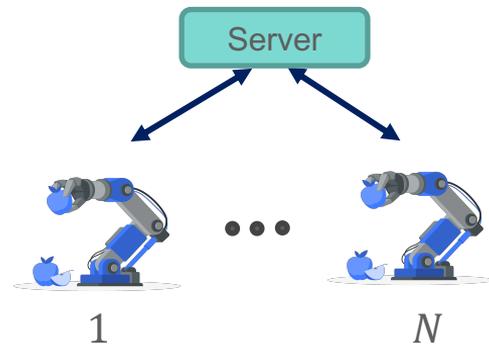
$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$

- However, we get

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \frac{\alpha^2}{N} + \underbrace{\alpha^3 + \Omega_t}_{\text{Higher order terms}}$$

Higher order terms

└─ Not important



Proof sketch

- Multiple agents, favorable recursion

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$

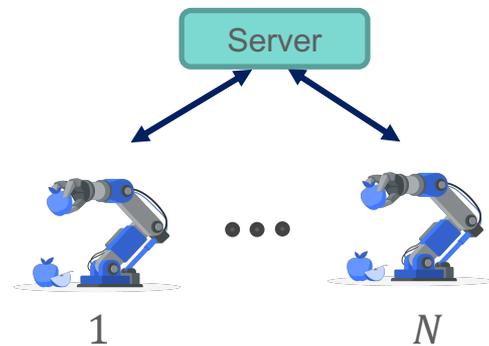
- However, we get

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \frac{\alpha^2}{N} + \underbrace{\alpha^3 + \Omega_t}_{\text{Higher order terms}}$$

Due to local updates

Higher order terms

Not important



Proof sketch

- Multiple agents, favorable recursion

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$

- However, we get

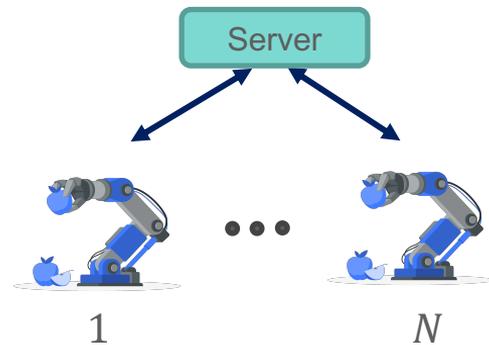
$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \frac{\alpha^2}{N} + \alpha^3 + \Omega_t$$

Higher order terms

Due to local updates

Handled by a special weighting

Not important



Proof sketch

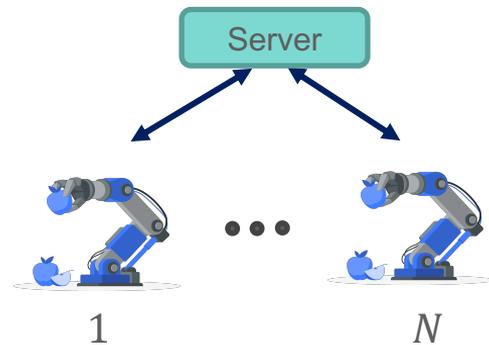
- Multiple agents, favorable recursion

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \alpha^2/N$$

- However, we get

$$\mathbb{E}[\|\theta_{t+1}\|^2] \leq (1 - \alpha)\mathbb{E}[\|\theta_t\|^2] + \frac{\alpha^2}{N} + \underbrace{\alpha^3}_{\text{Higher order terms}} + \Omega_t$$

$\mathcal{O}(1/N\epsilon)$ iteration complexity, linear speedup



Due to local updates

Handled by a special weighting

Not important

Other results

1. Federated Temporal Difference with Linear Function Approximation, on-policy data

Other results

1. Federated Temporal Difference with Linear Function Approximation, on-policy data
2. Federated Temporal Difference, off-policy data

Other results

1. Federated Temporal Difference with Linear Function Approximation, on-policy data
2. Federated Temporal Difference, off-policy data
3. Federated stochastic approximation with Markovian noise

Other results

1. Federated Temporal Difference with Linear Function Approximation, on-policy data
2. Federated Temporal Difference, off-policy data
3. Federated stochastic approximation with Markovian noise



Linear speedup + Constant communication cost

Other results

1. Federated Temporal Difference with Linear Function Approximation, on-policy data
2. Federated Temporal Difference, off-policy data
3. Federated stochastic approximation with Markovian noise



Linear speedup + Constant communication cost

**THANK YOU FOR YOUR
ATTENTION!**



ICML
International Conference
On Machine Learning