# Federated Reinforcement Learning
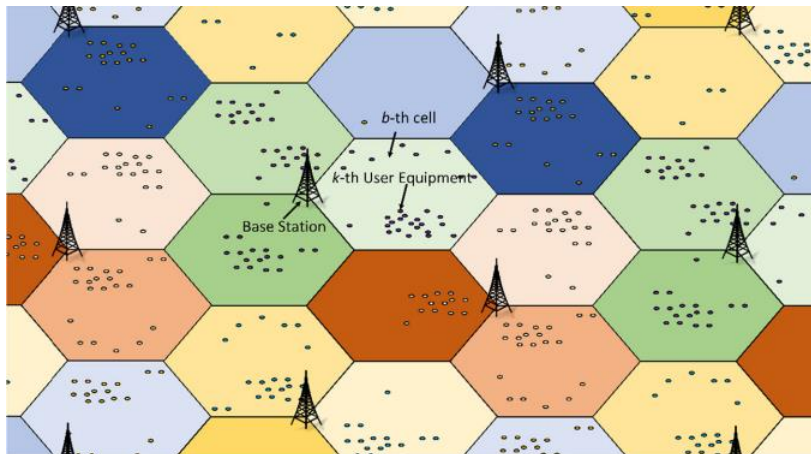
Hankz Hankui Zhuo, Wenfeng Feng, Qian Xu, Qiang Yang, Yufeng Lin

**Sun Yat-Sen University** & **Webank**

# Multi-agent reinforcement learning

- Global state or sharing state
- Global reward or sharing reward



Game!
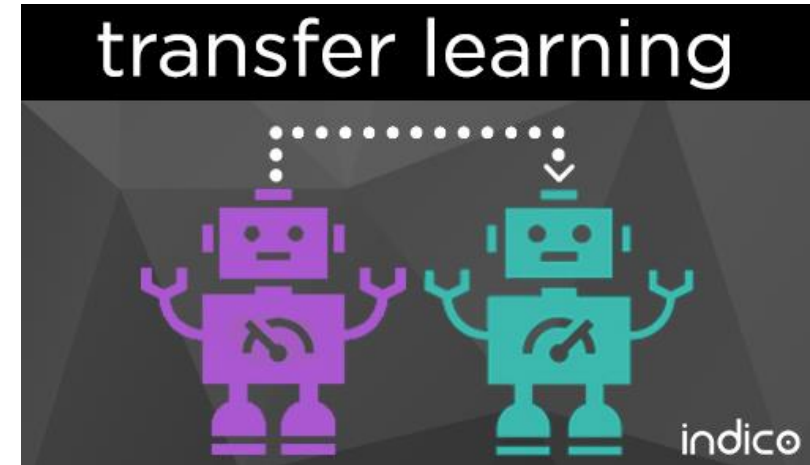


Configuration of base stations
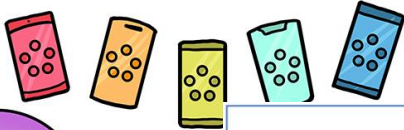


Logistics

Transfer data instance?

Transfer feature space?

Transfer models?

- Instance is private!
- Feature space is private!
- Model is private!

Federated Learning

Building better p
on-device data and p

An online comic fro

Google AI

Federated model

Model A

Encrypted model training

Encrypted entity alignment

Corp. A

No data exchange

Corp. B

Model B

a

Collaborator C

Data from A

Data from B

① Sending public keys

② Exchanging intermediate results

③ Computing gradients and loss

④ Updating models

b

Federated Machine Learning: Concept and Applications. Q Yang, Y Liu, T Chen, Y Tong. ACM Transactions on Intelligent Systems and Technology (TIST) 10 (2), 12, 2019

# Why does RL need to be federated?



State,        action,       reward

Surgery

State,     action

Pills

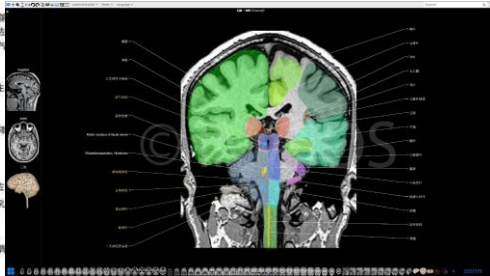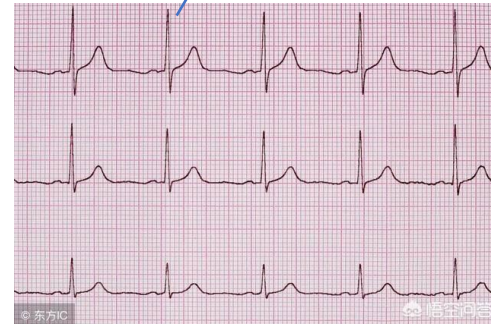In the manufacturing industry, producing products may involve various factories which produce different components of the products. Factories' decision policies are private and will not be shared with each other. On the other hand, building individual decision policies of high-quality on their own is often difficult due to their limited businesses and lack of rewards

# Problem setting

| | | |
|---|---|---|
| **MDP:** $\langle S, A, T, r \rangle$ | | |

**Agent α**  **Agent β**

| | | |
|---|---|---|
| **Output of MDP:** policies π | | |

Input: $\{\langle s_\alpha, a_\alpha, s'_\alpha, r_\alpha \rangle\}$   $\{\langle s_\beta, a_\beta \rangle\}$

Output: policies $\pi_\alpha^*$   Policies $\pi_\beta^*$

A1: The feature spaces of states $s_\alpha$ and $s_\beta$ are *different* between agents α and β.
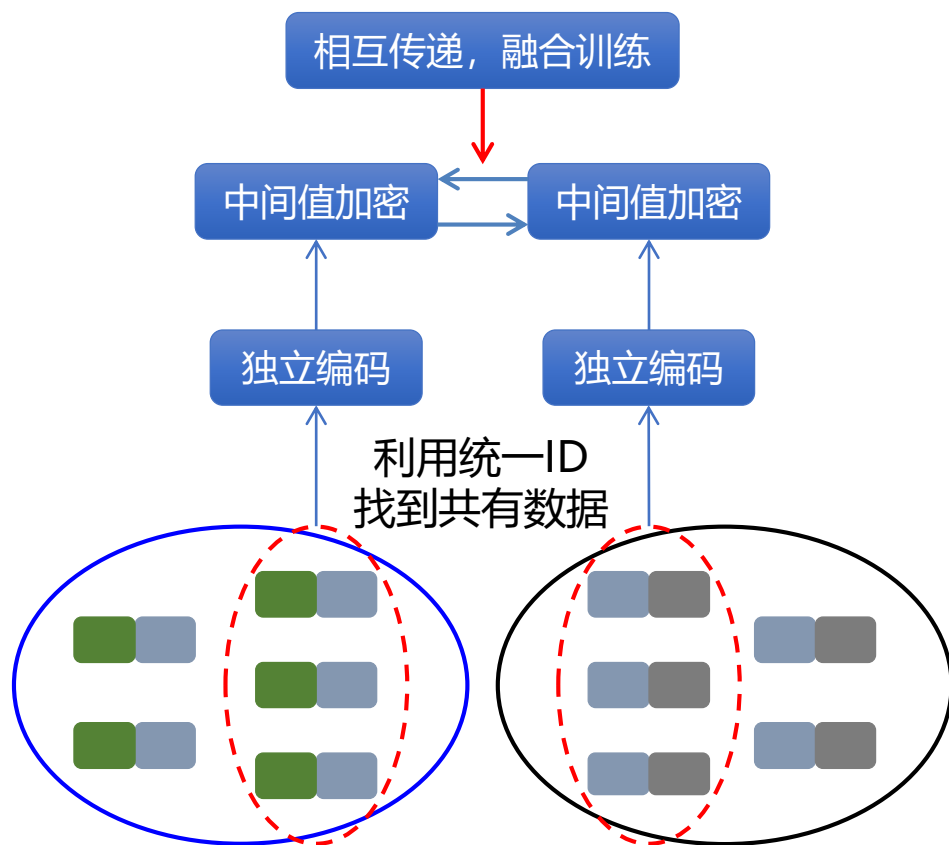
A2: $D_\alpha$ and $D_\beta$ cannot be shared directly between α and β

A3: The output of functions $Q_\alpha$ and $Q_\beta$ *can* be shared with each other

# FRL算法模型

如何利用分散的数据和特征来训练模型?

## FRL方法

相互传递，融合训练

中间值加密 ⇄ 中间值加密

独立编码　　独立编码

利用统一ID
找到共有数据

## 传统方法

缺点：机密数据无法融合，方法失效

训练模型

融合数据

# The Q learning model:

基础Q网络

两个网络之间传递
各自计算的加密Q值

$$Q_f^\alpha(\cdot, C_\beta; \theta_\alpha, \theta_f) = MLP([\hat{Q}_\alpha(\cdot; \theta_\alpha); C_\beta]; \theta_f)$$

$$Q_f^\beta(\cdot, C_\alpha; \theta_\beta, \theta_f) = MLP([\hat{Q}_\beta(\cdot; \theta_\beta); C_\alpha]; \theta_f)$$

损失函数

$$L_\alpha(\theta_\alpha, \theta_f) = \mathbb{E}\{(y^i - Q_f^\alpha(s_\alpha^i, a^i, C_\beta; \theta_\alpha, \theta_f))^2\}$$

$$L_\beta(\theta_\beta, \theta_f) = \mathbb{E}\{(y^i - Q_f^\beta(s_\beta^i, a^i, C_\alpha; \theta_\beta, \theta_f))^2\}$$

$$y^i = r^i + \gamma \max_a Q_f^\alpha(s_\alpha^i, a, C_\beta; \theta_\alpha, \theta_f)$$

$$a_t = \operatorname{argmax} Q_f$$
输出动作

$$Q_f$$
联合层Q值

MLP

$$\hat{Q}_\beta$$
传递Q值
$$\hat{Q}_\alpha$$

$$N(0, \sigma^2)$$
高斯噪声

$$Q_\alpha$$            $$Q_\beta$$
中间层Q值

$$\theta_\alpha$$        $$\theta_\beta$$

$$s_\alpha^t$$          $$s_\beta^t$$
状态输入

Agents' local models

Gausian differential privacy

Constant output of agent β

$Q_f = [q_1^f, ..., q_{d_a}^f]$

$Q_f = [q_1^f, ..., q_{d_a}^f]$

$\theta_g$

$\theta_g$

$N(0, \sigma^2) \oplus$

$C_\beta$

$C_\alpha$

$\oplus \quad N(0, \sigma^2)$

$Q_\alpha = [q_1^\alpha, ..., q_{d_a}^\alpha]$

$Q_\beta = [q_1^\beta, ..., q_{d_a}^\beta]$

$\theta_\alpha$

$\theta_\beta$

$s_\alpha^t$

Agent $\alpha$

Agent $\beta$

$s_\beta^t$

# Training procedure

Interaction between two agents **sequentially** or **in parallel**

(I) training

# Testing



- Need other agent's local output to make decisions
- Not necessarily making the same decisions



(II) testing

Agent $\alpha$

$s_\alpha$

Compute $C_\alpha$

Compute $Q_f^\alpha$

$a$

Agent $\beta$

$s_\beta$

Compute $C_\beta$

Compute $Q_f^\beta$

$a$

$C_\beta$

$C_\alpha$

| **Algorithm 1:** `FRL-ALPHA` | **Algorithm 2:** `FRL-BETA` |
|---|---|

**Input:** state space $S_\alpha$, action space $A_\alpha$, rewards $r$

**Output:** $\theta_\alpha, \theta_g$

1: Initialize $Q_\alpha, Q_f$ with random values for $\theta_\alpha, \theta_g$
2: Initialize replay memory $D_\alpha$
3: Call `FRL-BETA`.$Init()$
4: **for** episode = 1: $M$ **do**
5:    **repeat**
6:       Observe $s_\alpha^t$
7:       Call $C_\beta$=`FRL-BETA`.$ComputeQBeta()$
8:       Select action $a^t$ with probability $\epsilon$
9:       Otherwise
         $a^t = \arg\max_a Q_f^\alpha(s_\alpha^t, a, C_\beta; \theta_\alpha, \theta_g)$
10:     Execute action $a^t$, obtain reward $r^t$ and state $s^{t+1}$
11:     Observe $s_\alpha^{t+1}$, store $(s_\alpha^t, a^t, r^t, s_\alpha^{t+1})$ in $D_\alpha$
12:     Sample $(s_\alpha^j, a^j, r^j, s_\alpha^{j+1})$ from $D_\alpha$
13:     Call $C_\beta$=`FRL-BETA`.$ComputeQBeta(j)$
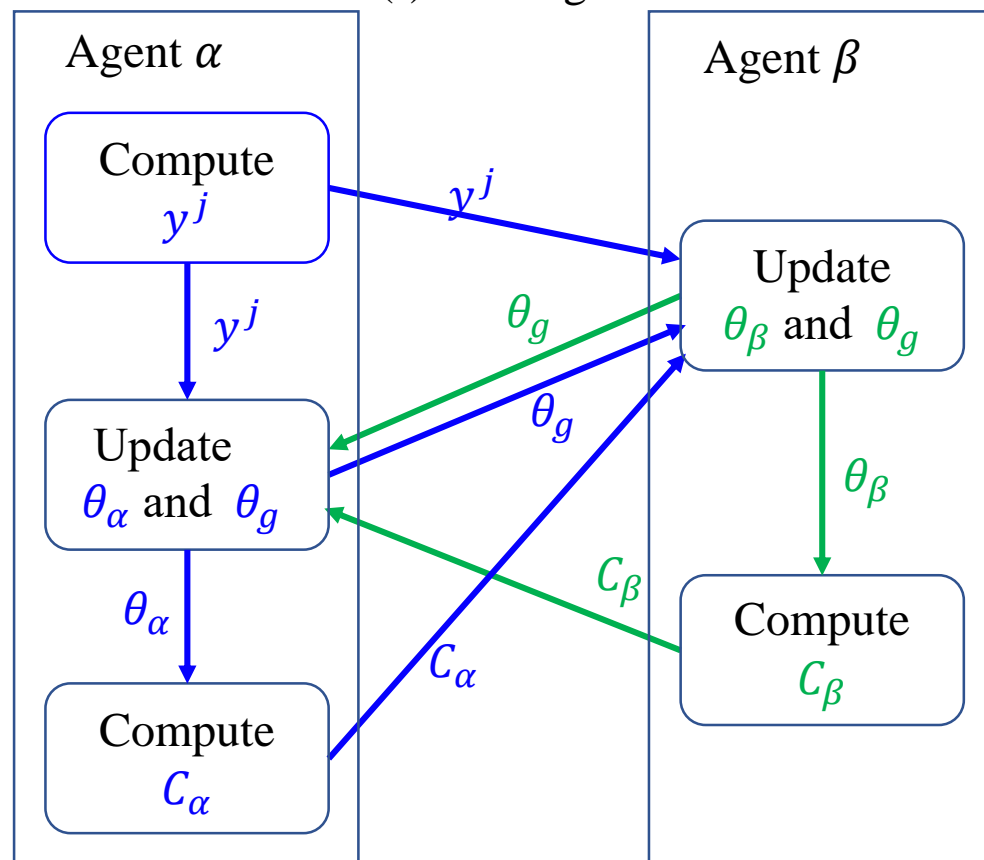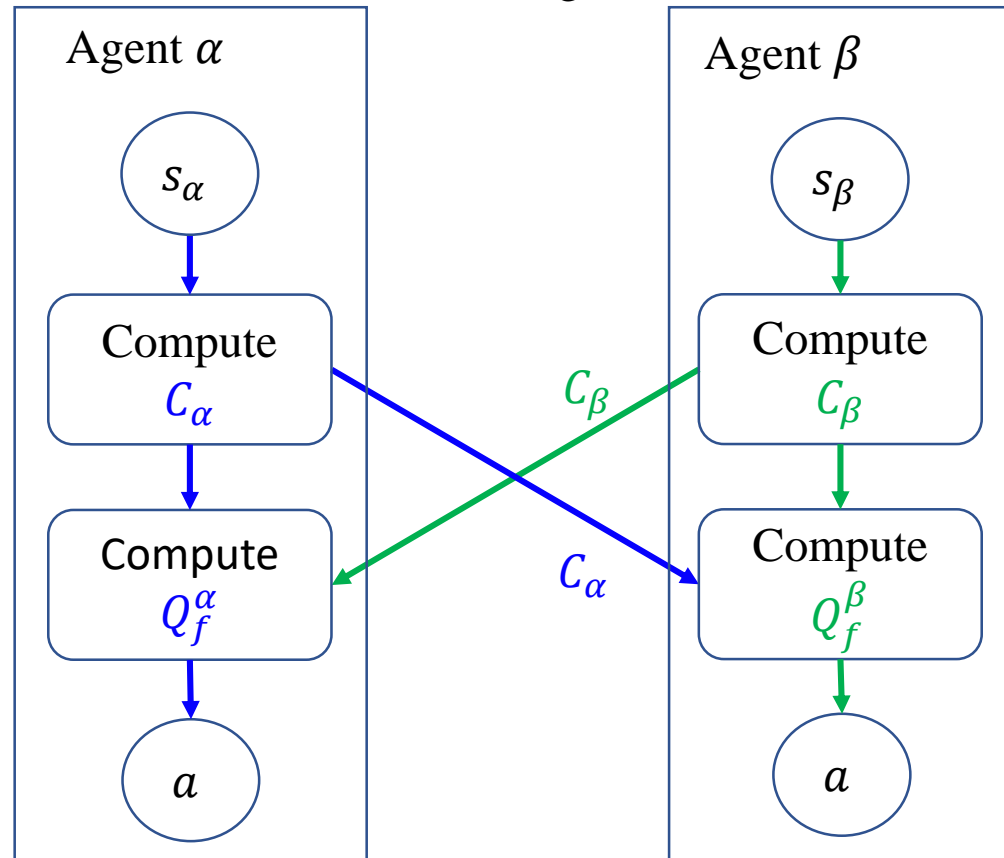14:     $Y^j = r^j + \gamma \max_a Q_f^\alpha(s_\alpha^j, a, C_\beta; \theta_\alpha, \theta_g)$
15:     Update $\theta_\alpha, \theta_g$ according to Eq. (4), (6)
16:     $C_\alpha = \hat{Q}_\alpha(s_\alpha^j, a; \theta_\alpha)$
17:     Call
        $\theta_g$=`FRL-BETA`.$UpdateQ(Y^j, j, C_\alpha, \theta_g)$
18:    **until** terminal $t$
19: **end for**

**Input:** state space $S_\beta$, action space $A_\beta$

**Output:** $\theta_\beta, \theta_g$

1: **function** $Init()$
2:    Initialize $Q_\beta$ with random values for $\theta_\beta$
3:    Initialize replay memory $D_\beta$
4: **end function**
5: **function** $ComputeQBeta()$
6:    Observe $s_\beta$
7:    Select $a_\beta \in A_\beta$ with probability $\epsilon$
8:    Otherwise $a_\beta = \arg\max_{a_\beta} Q_\beta(s_\beta, a_\beta; \theta_\beta)$
9:    Store $(s_\beta, a_\beta)$ in $D_\beta$
10:   Let $C_\beta = \hat{Q}_\beta(s_\beta, a; \theta_\beta)$
11:   **return** $C_\beta$
12: **end function**
13: **function** $ComputeQBeta(j)$
14:   Select $(s_\beta, a_\beta)$ from $D_\beta$ based on index $j$
15:   Let $C_\beta = \hat{Q}_\beta(s_\beta, a_\beta; \theta_\beta)$
16:   **return** $C_\beta$
17: **end function**
18: **function** $UpdateQ(Y^j, j, C_\alpha, \theta_g)$
19:   Select $(o_\beta^j, a_\beta^j)$ from $D_\beta$ based on index $j$
20:   Update $\theta_\beta, \theta_g$ based on Eq. (5), (7)
21:   **return** $\theta_g$
22: **end function**

# baselines

**DQN-alpha**
A deep Q-network trained with agent $\alpha$'s data only. It takes observations $s\alpha$ as input and outputs actions corresponding to $s_\alpha$.

**DQN-full**
A deep Q-network trained by directly putting data together from both agent $\alpha$ and $\beta$, i.e., neglecting data privacies between agents $\alpha$ and $\beta$.

**CNN-alpha**
A convolutional neural network trained with agent $\alpha$'s data only, similar to DQN-alpha.

**CNN-full**
A convolutional neural network trained with all data of agent $\alpha$ and $\beta$ put together directly, similar to DQN-full.

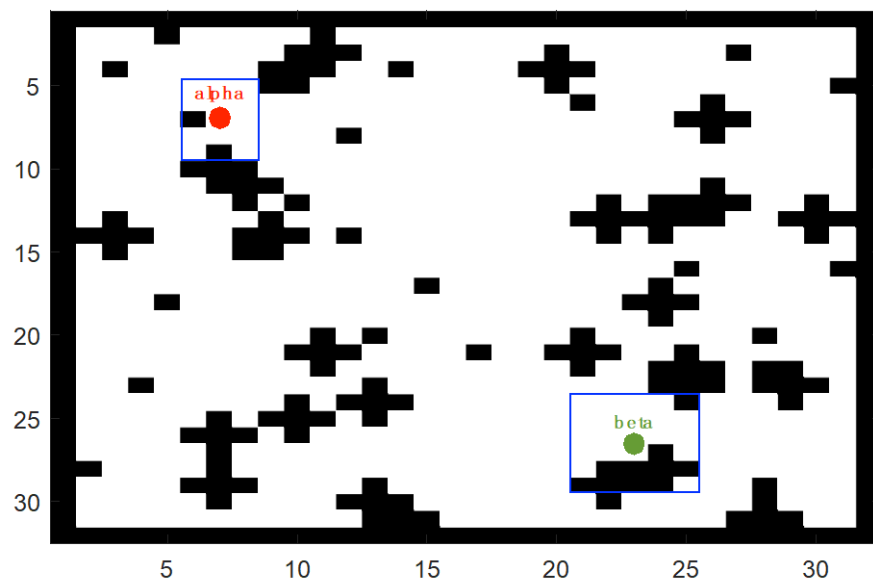# Domain: Grid-World

States: The domain is represented by a Ng × Ng

Actions: There are 4 actions for each agent, i.e. going towards 4 directions, denoted by {east, south, west, north}.



Rewards: The reward is composed of two parts, i.e., local reward $r_l$ and global reward $r_g$

Dataset: We generated 8000 different maps (or matrices) for each size of 8 × 8, 16 × 16 and 32 × 32

# Domain 2: Text2Action

## 数据集

- **CT**: 做菜教程 （Cooking Tutorials）
- **WHS**: 电脑指令 （Branavan et al. ACL 2009）
- **WHG**: 家庭和花园打理指南 （Malmaud et al., ACL 2014）

输入的一个指令性文本

## 评价指标

- 累积奖励
- F1分数

所有的动作序列输出

任务

## 基准算

| 算法 | |
|---|---|
| DepR | |
| LSTM | (Lindsay et al., |
| EncD | Hovy, ACL 2016) |
| 指令 | H et al., AAAI 201 |
| CMLP | |
| DQN | 的EASDRL模型 |
| DQN | 的EASDRL模型 |

**How to Make Egg Fried Rice?** （如何做蛋炒饭？）
**Cook** the rice the day before, or **use** leftover rice in the refrigerator. The important thing to **remember** is not to **heat** up the rice, but **keep** it cold. In a bowl, add 1 tablespoon of oil to rice. Use a spoon or your hands to work the oil into the rice, evenly coating the rice. Transfer the rice to a colander and drain. Combine eggs and salt in a small bowl and gently whisk until blended. Heat 1 tablespoon oil in a wok. Add whisked eggs and cumin seeds to wok. Stir frequently, working the eggs to a scramble. Heat the remaining oil in the wok. If desired, you can **recycle** some of the oil that drained from the rice. Add the garlic and onion to the wok. Stir-fry together over high heat for about 5 minutes or until the onion looks transparent, but is not soft. Add the rice, eggs, soy sauce, chili sauce, vinegar, and celery. Mix together, continuing to stir-fry over high heat for 1-2 minutes while stirring frequently. Spoon onto a plate and serve.

动作 ── 动作名称
动作 ── 动作参数

煮（饭）

- Cook (rice) → Keep (rice, cold)→ Add (oil) → Use (spoon) → Work (oil, rice) → ⋯ → Work (eggs) → Heat (oil) → ⋯
- Use (leftover rice) → Keep (rice, cold) → Add (oil) → Use (spoon) → Work (oil, rice)→ ⋯ → Work (eggs) → Heat (oil) → ⋯ → Serve ()
- Use (leftover rice) → Keep (rice, cold) → Add (oil) → Use (hands) → Work (oil, rice) → ⋯ → Work (eggs) → Heat (oil) → ⋯ → Serve ()
- Use (leftover rice) → Keep (rice, cold) → Add (oil) → Use (hands) → Work (oil, rice)→ ⋯ → Work (eggs) → Recycle (oil) → Heat (oil) → ⋯ → Serve ()
- ⋯

# Domain 2: Text2Action

**States:** $s_\alpha \in R^{N_w \times K_1}$ is real-valued matrix that describes the part-of-speech of words, and $s_\beta \in R^{N_w \times K_2}$ is real-valued matrix that describes the embedding of words.

**Actions:** There are two actions for each agent, i.e., {select, neglect}

**Rewards:** The instant reward include a basic reward and an additional reward

**Criteria:**

$$precision = \frac{\#TotalRight}{\#TotalSelected}$$

$$recall = \frac{\#TotalRight}{\#TotalTruth}$$

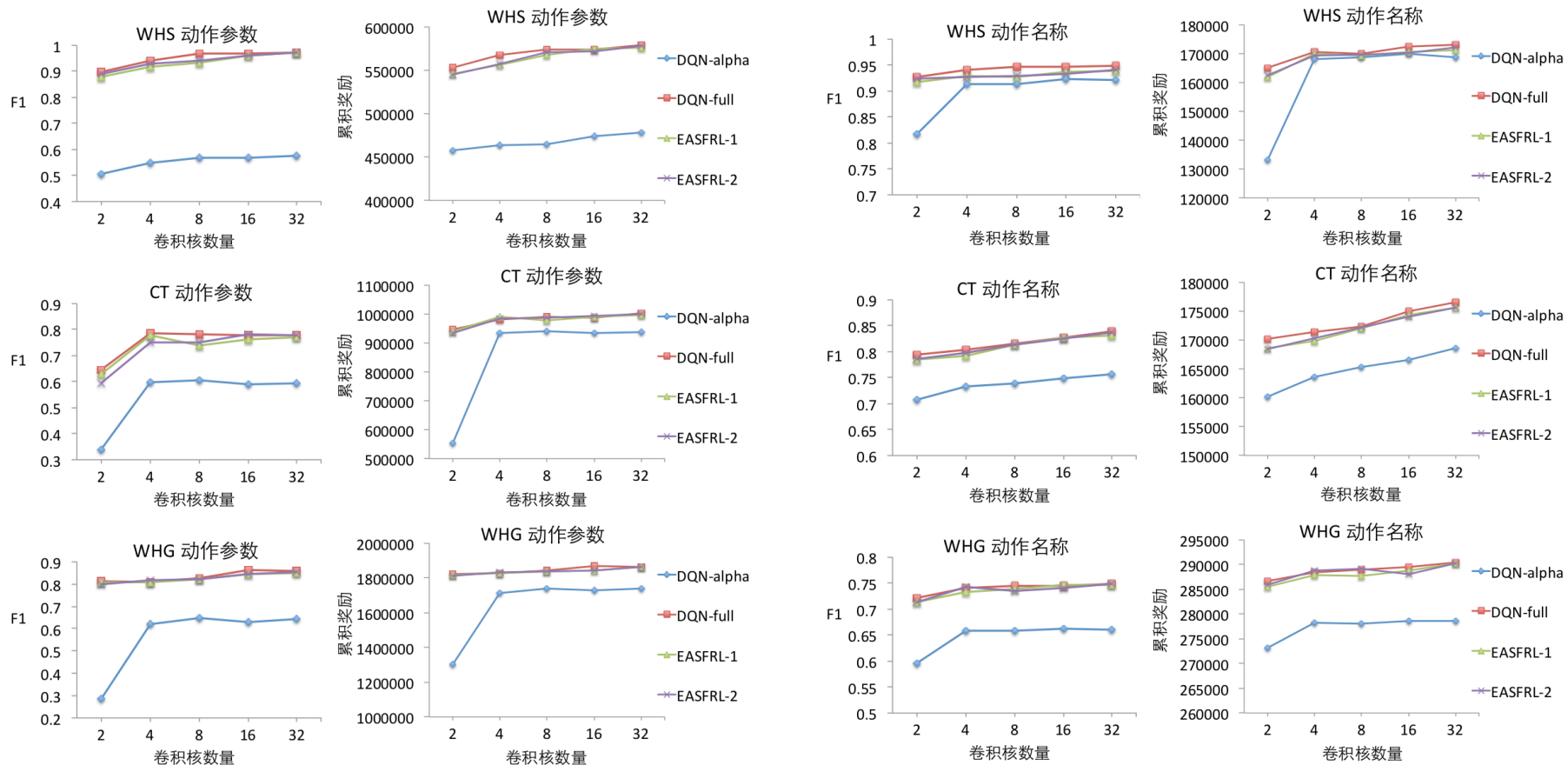$$F1 = \frac{2 \times precision \times recall}{precision + recall}$$

# 实验结果分析

| 评价标准 | 算法 | 动作名称 | | | 动作参数 | | |
|---|---|---|---|---|---|---|---|
| | | WHS | CT | WHG | WHS | CT | WHG |
| F1 分数 (%) | DQN-alpha | 92.11 | 75.64 | 66.37 | 54.13 | 59.46 | 61.09 |
| | EASFRL-1 | 93.76 | 83.05 | 74.64 | 97.18 | 76.97 | 84.95 |
| | EASFRL-2 | **94.11** | **83.72** | **74.85** | **97.27** | **77.75** | **85.44** |
| | DQN-full | 94.75 | 83.87 | 74.87 | 97.35 | 77.58 | 85.66 |
| 累积奖励 ($\times 10^4$) | DQN-alpha | 16.87 | 16.86 | 27.86 | 47.78 | 93.82 | 174.01 |
| | EASFRL-1 | 17.13 | **17.56** | 29.01 | 57.55 | 99.69 | **186.13** |
| | EASFRL-2 | **17.22** | **17.56** | **29.03** | **57.82** | **100.00** | 186.12 |
| | DQN-full | 17.31 | 17.65 | 29.04 | 57.88 | 100.19 | 186.42 |

结论：
- FRL模型在所有实验中表现都明显高于只有部分数据的模型，说明<span style="color:red">将数据联合起来训练能明显提高模型的性能</span>
- FRL模型性能非常接近于直接融合全部数据来训练的模型，说明FRL算法能够<span style="color:red">保护数据隐私</span>的同时，保证<span style="color:red">模型的性能几乎不降低</span>

# 探究性实验结果分析



结论：
- 模型简单时FRL相对部分数据训练的模型的优势更大，说明联邦学习能够取得较好的效果。
- 模型复杂度变化的时，FRL模型的性能都能很好的逼近融合全部数据训练的模型的性能。

Thank You