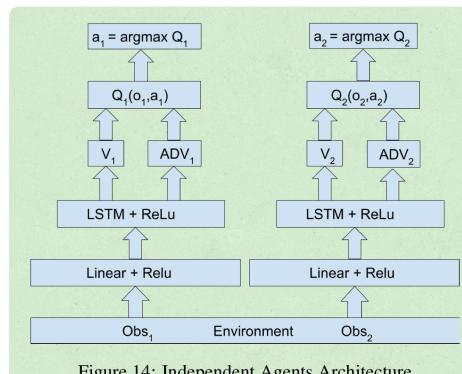


一、问题

1. 现有 Multiagent 架构缺点.

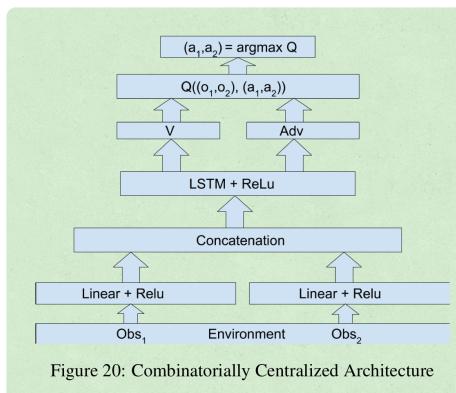
i). 分布式独立结构



① Observation i
部分可观 PO-MDP
Policy Pi 对第 i 层不是最优 Policy
Aif. 部分而最优都无效.

② non-stationary Env

2) 结合中心式架构

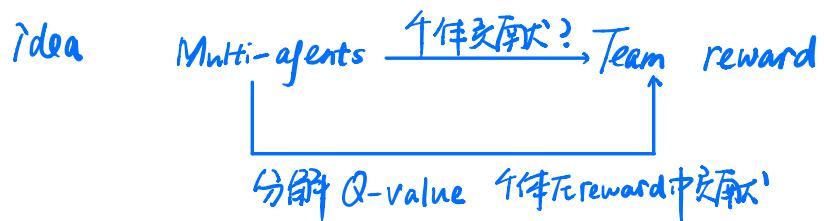


① credit assignment
Team reward
谁的贡献大?

② joint action space

所有 agent 共享共同加

Paper: 每个 Agent 贡献 Team reward 中贡献大小?



二. A Deep-RL Architecture for Coop-MARL

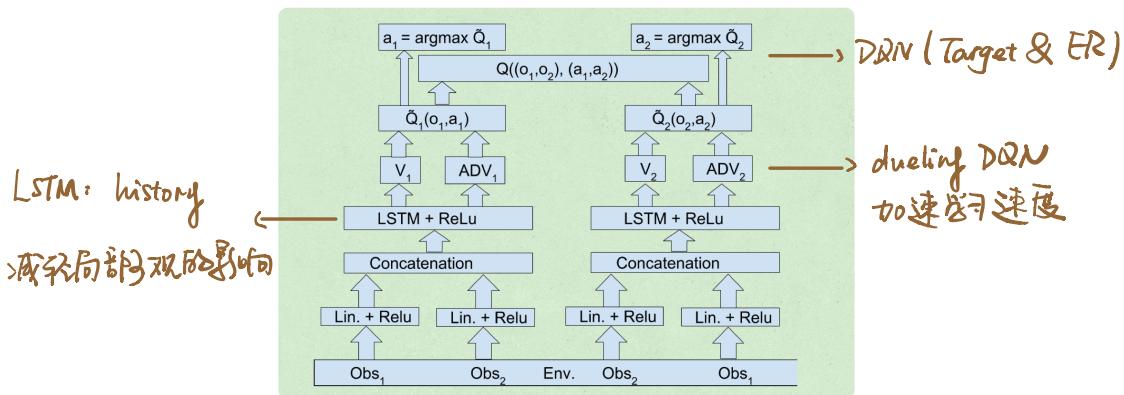
the identified issues ←— value-decomposition

$$Q((h^1, h^2, \dots, h^d), (a^1, a^2, \dots, a^d)) \approx \sum_{i=1}^d \tilde{Q}_i(h^i, a^i)$$

成功分解、合作场景。

局部 observation 用 history 代替 $h_t = (o_1, r_1, \dots, o_{t-1}, r_{t-1})$ 得到更优策略

论文中没有线性分解证明，大量实验验证，假设。



三. 实验：3 Task

1. 场景

Paper 没有证明 Value Decomposed 分解
而是用大量实验证明。

① Switch：长廊环境到达各自的目的地。

② Fetch：取放货，另一个 Agent 帮忙 (类似物流环境)

③ Checkers：

	Apple	Lemon
Agent 1	10	-10
Agent 2	1	-1

让 Agent 1 和 Apple Agent 2 和 Lemon 得分最高。

2. 算法比较

不同方法技巧

Agent	V.	S.	Id	L.	H.	C.
1						
2	✓					
3	✓	✓				
4	✓	✓	✓	✓		
5	✓	✓	✓	✓		
6	✓	✓	✓	✓	✓	
7	✓	✓	✓	✓	✓	✓
8	✓					
9						

Table 1: Agent architectures. V is value decomposition, S means shared weights and an invariant network, Id means role info was provided. L stands for lower-level communication, H for higher-level communication and C for centralization. These architectures were selected to show the advantages of the independent agent with value-decomposition and to study the benefits of additional enhancements added in a logical sequence.

VD 方法 Individual Agent 方法

各自 VD 好于 Centralized 方法。

应用到多种奖励分析。

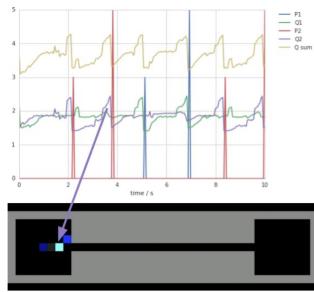


Figure 6: The learned Q -decomposition in Fetch. The plot shows the total Q -function (yellow), the value of agent 1 (green), the value of agent 2 (purple), rewards from agent 1 (blue) events and agent 2 (red). Highlighted is a situation in which agent 2's Q -function spikes (purple line), anticipating reward for an imminent drop-off. The other agent's Q -function (green) remains relatively flat.

Fetch 场景，说明 Value-decomposed 有效。有较高 Q -value 的 Agent 对 Team Reward 贡献大。

Paper 思路清晰，以实验分析为主导，对比了所有方法的架构，综合分析 Value-Decomposed 的作用效果。

论文不提之处：① Team 规模 分解难度有效性。

② Value-decomposed 方法。

—— 场景
Value Decomposed 方法
非 Value Decomposed

