# MUNI

# MARL's Zoo

An Introduction to Multi-Agent RL

**Vojtěch Formánek**
**xforman@fi.muni.cz**

Masaryk University

December 27, 2025

# Today

- Intro
- Solving Joint Policies
- Independent -DQN, -REINFORCE, -A2C
- Centralized Critics
- Value Decomposition (VDN, QMIX)

- **Intro**
- Solving Joint Policies
- Independent -DQN, -REINFORCE, -A2C
- Centralized Critics

**Definition 1 (Stochastic Game)** *consists of:*

- *Finite set of agents $I = 1, \ldots, n$*
- *Finite set of states $S$, with a subset of terminal $\overline{S} \subset S$*
- *State transition function $\mathcal{T} : S \times A \times S \to [0, 1]$ such that*

$$\forall s \in S, a \in A : \sum_{s' \in S} \mathcal{T}(s, a, s') = 1,$$

  *where $A = A_1 \times \cdots \times A_n$*
- *Initial state distribution $S \to [0, 1]$*

$$\sum_{s \in S} \mu(s) = 1 \quad and \quad \forall s \in \overline{S} : \mu(s) = 0$$

- *For each $i \in I$*
    - *Reward function $\mathcal{R}_i : S \times A \times S$*
    - *Finite set of actions $A_i$*

**Definition 2 (Partially observable SG)** *consists of:*

- *Finite set of agents* $I = 1, \ldots, n$
- *Finite set of states* $S$, with a subset of terminal $\overline{S} \subset S$
- *State transition function* $\mathcal{T} : S \times A \times S \to [0, 1]$ *such that*

$$\forall s \in S, a \in A : \sum_{s' \in S} \mathcal{T}(s, a, s') = 1,$$

*where* $A = A_1 \times \cdots \times A_n$
- *Initial state distribution* $S \to [0, 1]$

$$\sum_{s \in S} \mu(s) = 1 \quad and \quad \forall s \in \overline{S} : \mu(s) = 0$$

- *For each* $i \in I$
  - *Reward function* $\mathcal{R}_i : S \times A \times S$
  - *Finite set of actions* $A_i$ *and observations* $O_i$
  - *Observation function* $\mathcal{O}_i : A \times S \times O_i \to [0, 1]$ such that

$$\forall a \in A, s \in S : \sum_{o_i \in O_i} \mathcal{O}_i(a, s, o_i) = 1$$

# Notation

- $s^t = (s_1^t, \ldots, s_n^t), a^t = (a_1^t, \ldots, a_n^t)$, at time $t$
- *observations*
    - SG: $o_i^t = (s^t, a^{t-1})$
    - Other actions unobserved: $o_i^t = (s^t, a_i^{t-1})$
    - Limited view: $o_i^t = (\vec{s}^t, \vec{a}^t)$, where $\vec{s}^t \subset s^t$ and $\vec{a}^t \subset a^t$
- history
    - full: $\hat{h}^t = \{s^0, o^0, a^0, \ldots, s^t, o^t\}$
    - observations: $\sigma(\hat{h}^t) = \{o^0, \ldots, o^t\}$

# Today

- Intro
- **Solving Joint Policies**
- Independent -DQN, -REINFORCE, -A2C
- Centralized Critics

## Policy

*Joint policy $\pi = (\pi_1, \ldots, \pi_n)$ satisfies requirements in terms of expected return $U_i^\pi$, where*

$$U_i^\pi = \lim_{t \to \infty} \mathbb{E}_{\hat{h}^t \sim (\mu, \mathcal{T}, \mathcal{O}, \pi)} \left[ \sum_{\tau=0}^{t-1} \gamma^\tau \mathcal{R}_i \right]$$

*for each agent i.*

# Recursive Expected Return

$$V_i^\pi(\hat{h}) = \sum_{a \in A} \pi(a \mid \sigma(\hat{h})) \, Q_i^\pi(\hat{h}, a)$$

$$Q_i^\pi(\hat{h}, a) = \sum_{s' \in S} \mathcal{T} \left[ \mathcal{R}_i + \gamma \sum_{o' \in O} \mathcal{O}(o'|a, s') \, V_i^\pi(\langle \hat{h}, a, s', o' \rangle) \right]$$

$$U_i^\pi = \mathbb{E}^\pi [V_i^\pi(\langle s^0, o^0 \rangle)]$$

# Solutions
## Equilibrium

*In a general-sum game $\pi$ is a* **Nash Equilibrium** *if*

$$\forall i, \pi_i' : U_i^{\langle \pi_i', \pi_{-i} \rangle} \leq U_i^\pi$$

- Sub-optimality
- Non-uniqueness
- Incompleteness

## Solutions
### Non-Equilibrium

**Pareto Optimality**:

*A joint policy $\pi$ is* Pareto-dominated *by $\pi'$ if*

$$\forall i : U_i^{\pi'} \geq U_i^{\pi} \text{ and } \exists i : U_i^{\pi'} > U_i^{\pi},$$

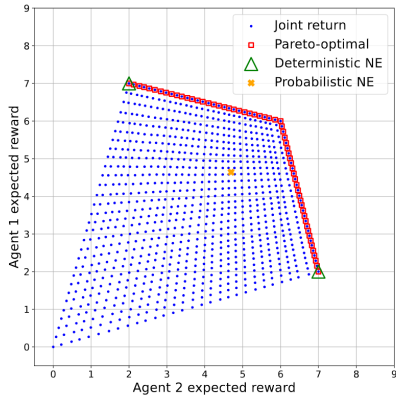$\pi$ *is* Pareto-optimal *if it is not* Pareto-dominated

Figure: Feasible expected joint rewards and Pareto frontier in the Chicken matrix game.

# Three Paradigms

- Centralized training and execution
- Decentralized training and execution
- Centralized training and decentralized execution

**Algorithm 4** Central Q-learning (CQL) for stochastic games

1: Initialize: $Q(s, a) = 0$ for all $s \in S$ and $a \in A = A_1 \times ... \times A_n$
2: Repeat for every episode:
3: **for** $t = 0, 1, 2, ...$ **do**
4:   Observe current state $s^t$
5:   With probability $\epsilon$: choose random joint action $a^t \in A$
6:   Otherwise: choose joint action $a^t \in \arg\max_a Q(s^t, a)$
7:   Apply joint action $a^t$, observe rewards $r_1^t, ..., r_n^t$ and next state $s^{t+1}$
8:   Transform $r_1^t, ..., r_n^t$ into scalar reward $r^t$
9:   $Q(s^t, a^t) \leftarrow Q(s^t, a^t) + \alpha \left[ r^t + \gamma \max_{a'} Q(s^{t+1}, a') - Q(s^t, a^t) \right]$

**Algorithm 5** Independent Q-learning (IQL) for stochastic games

    *// Algorithm controls agent i*

1: Initialize: $Q_i(s, a_i) = 0$ for all $s \in S, a_i \in A_i$
2: Repeat for every episode:
3: **for** $t = 0, 1, 2, \ldots$ **do**
4:     Observe current state $s^t$
5:     With probability $\epsilon$: choose random action $a_i^t \in A_i$
6:     Otherwise: choose action $a_i^t \in \arg\max_{a_i} Q_i(s^t, a_i)$
7:     (meanwhile, other agents $j \neq i$ choose their actions $a_j^t$)
8:     Observe own reward $r_i^t$ and next state $s^{t+1}$
9:     $Q_i(s^t, a_i^t) \leftarrow Q_i(s^t, a_i^t) + \alpha \left[ r_i^t + \gamma \max_{a_i'} Q_i(s^{t+1}, a_i') - Q_i(s^t, a_i^t) \right]$

# A Big Step in Theory
## Challenges

- Non-stationarity
- Equilibrium selection
- Credit assignment
- Scaling to many agents

# Today

- Intro
- Solving Joint Policies
- **Independent -DQN, -REINFORCE, -A2C**
- Centralized Critics

**Algorithm 17** Independent deep Q-networks

1: Initialize $n$ value networks with random parameters $\theta_1, \ldots, \theta_n$
2: Initialize $n$ target networks with parameters $\bar{\theta}_1 = \theta_1, \ldots, \bar{\theta}_n = \theta_n$
3: Initialize a replay buffer for each agent $D_1, D_2, \ldots, D_n$
4: **for** time step $t = 0, 1, 2, \ldots$ **do**
5:     Collect current observations $o_1^t, \ldots, o_n^t$
6:     **for** agent $i = 1, \ldots, n$ **do**
7:         With probability $\epsilon$: choose random action $a_i^t$
8:         Otherwise: choose $a_i^t \in \arg\max_{a_i} Q(h_i^t, a_i; \theta_i)$
9:     Apply actions $(a_1^t, \ldots, a_n^t)$; collect rewards $r_1^t, \ldots, r_n^t$ and next observations $o_1^{t+1}, \ldots, o_n^{t+1}$
10:     **for** agent $i = 1, \ldots, n$ **do**
11:         Store transition $(h_i^t, a_i^t, r_i^t, h_i^{t+1})$ in replay buffers $D_i$
12:         Sample random mini-batch of $B$ transitions $(h_i^k, a_i^k, r_i^k, h_i^{k+1})$ from $D_i$
13:         **if** $s^{k+1}$ is terminal[2] **then**
14:             Targets $y_i^k \leftarrow r_i^k$
15:         **else**
16:             Targets $y_i^k \leftarrow r_i^k + \gamma \max_{a_i' \in A_i} Q(h_i^{k+1}, a_i'; \bar{\theta}_i)$
17:         Loss $\mathcal{L}(\theta_i) \leftarrow \frac{1}{B} \sum_{k=1}^{B} \left( y_i^k - Q(h_i^k, a_i^k; \theta_i) \right)^2$
18:         Update parameters $\theta_i$ by minimizing the loss $\mathcal{L}(\theta_i)$
19:         In a set interval, update target network parameters $\bar{\theta}_i$

# Addressing Non-stationarity

- Buffer can store outdated experiences
- Use smaller buffers
- Importance Sampling: reweigh based on $a_{-i}$
- Switch to on-policy

**Algorithm 18** Independent REINFORCE

1: Initialize $n$ policy networks with random parameters $\phi_1, \ldots, \phi_n$
2: Repeat for every episode:
3: **for** time step $t = 0, 1, 2, \ldots, T-1$ **do**
4:     Collect current observations $o_1^t, \ldots, o_n^t$
5:     **for** agent $i = 1, \ldots, n$ **do**
6:         Sample actions $a_i^t$ from $\pi(\cdot \mid h_i^t; \phi_i)$
7:     Apply actions $(a_1^t, \ldots, a_n^t)$; collect rewards $r_1^t, \ldots, r_n^t$ and next observations $o_1^{t+1}, \ldots, o_n^{t+1}$
8: **for** agent $i = 1, \ldots, n$ **do**
9:     Loss $\mathcal{L}(\phi_i) \leftarrow -\frac{1}{T} \sum_{t=0}^{T-1} \left( \sum_{\tau=t}^{T-1} \gamma^{\tau-t} r_i^\tau \right) \log \pi(a_i^t \mid h_i^t; \phi_i)$
10:     Update parameters $\phi_i$ by minimizing the loss $\mathcal{L}(\phi_i)$

**Algorithm 19** Independent A2C with synchronous environments

1: Initialize $n$ actor networks with random parameters $\phi_1, \ldots, \phi_n$
2: Initialize $n$ critic networks with random parameters $\theta_1, \ldots, \theta_n$
3: Initialize $K$ parallel environments
4: **for** time step $t = 0 \ldots$ **do**
5:    Batch of observations for each agent and environment: $\begin{bmatrix} o_1^{t,1} \ldots o_1^{t,K} \\ \ddots \\ o_n^{t,1} \ldots o_n^{t,K} \end{bmatrix}$
6:    Sample actions $\begin{bmatrix} a_1^{t,1} \ldots a_1^{t,K} \\ \ddots \\ a_n^{t,1} \ldots a_n^{t,K} \end{bmatrix} \sim \pi(\cdot \mid h_1^t; \phi_1), \ldots, \pi(\cdot \mid h_n^t; \phi_n)$
7:    Apply actions; collect rewards $\begin{bmatrix} r_1^{t,1} \ldots r_1^{t,K} \\ \ddots \\ r_n^{t,1} \ldots r_n^{t,K} \end{bmatrix}$ and observations $\begin{bmatrix} o_1^{t+1,1} \ldots o_1^{t+1,K} \\ \ddots \\ o_n^{t+1,1} \ldots o_n^{t+1,K} \end{bmatrix}$
8:    **for** agent $i = 1, \ldots, n$ **do**
9:       **if** $s^{t+1,k}$ is terminal **then**
10:          Advantage $Adv(h_i^{t,k}, a_i^{t,k}) \leftarrow r_i^{t,k} - V(h_i^{t,k}; \theta_i)$
11:          Critic target $y_i^{t,k} \leftarrow r_i^{t,k}$
12:       **else**
13:          Advantage $Adv(h_i^{t,k}, a_i^{t,k}) \leftarrow r_i^{t,k} + \gamma V(h_i^{t+1,k}; \theta_i) - V(h_i^{t,k}; \theta_i)$
14:          Critic target $y_i^{t,k} \leftarrow r_i^{t,k} + \gamma V(h_i^{t+1,k}; \theta_i)$
15:       Actor loss $\mathcal{L}(\phi_i) \leftarrow \frac{1}{K} \sum_{k=1}^{K} Adv(h_i^{t,k}, a_i^{t,k}) \log \pi(a_i^{t,k} \mid h_i^{t,k}; \phi_i)$
16:       Critic loss $\mathcal{L}(\theta_i) \leftarrow \frac{1}{K} \sum_{k=1}^{K} \left( y_i^{t,k} - V(h_i^{t,k}; \theta_i) \right)^2$
17:       Update parameters $\phi_i$ by minimizing the actor loss $\mathcal{L}(\phi_i)$
18:       Update parameters $\theta_i$ by minimizing the critic loss $\mathcal{L}(\theta_i)$

# MA Policy Gradient Theorem
## Partially Observable Case

$$\nabla_{\phi_i} J(\phi_i) \propto \mathbb{E}_{\hat{h} \sim \Pr(\hat{h}|\pi), a_i \sim \pi_i, a_{-i} \sim \pi_{-i}} \left[ Q_i^\pi(\hat{h}, \langle a_i, a_{-i} \rangle) \nabla_{\phi_i} \log \pi_i(a_i \mid h_i = \sigma_i(\hat{h}); \phi_i) \right]$$
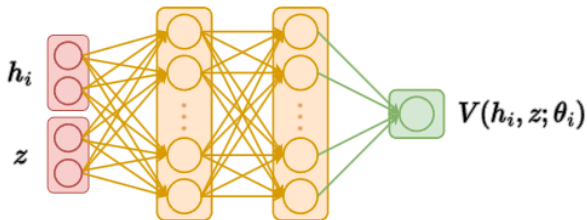
# Today

- Intro
- Solving Joint Policies
- Independent -DQN, -REINFORCE, -A2C
- **Centralized Critics**

# The Idea

- No constraints on critic during training
- Give critic info $z$ about other agents
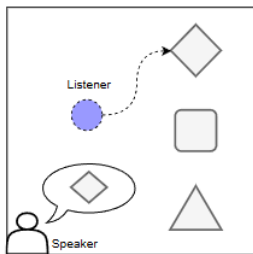- critic as $V(h_1^t, \ldots, h_n^t; \theta_i)$
- or ...

# Centralized Critic



**Value Loss:**

$$\mathcal{L}(\theta_i) = \left(y_i - V(h_i^t, z^t; \theta_i)\right)^2 \quad \text{with} \quad y_i = r_i^t + \gamma V(h_i^t, z^t; \theta_i)$$
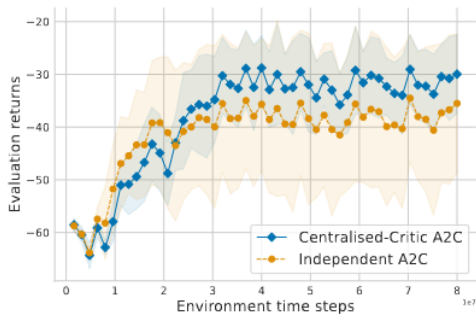
**Algorithm 20** Centralized A2C with synchronous environments

1: Initialize $n$ actor networks with random parameters $\phi_1, \ldots, \phi_n$
2: Initialize $n$ critic networks with random parameters $\theta_1, \ldots, \theta_n$
3: Initialize $K$ parallel environments
4: **for** time step $t = 0 \ldots$ **do**
5:   Batch of observations for each agent and environment: $\begin{bmatrix} o_1^{t,1} \ldots o_1^{t,K} \\ \ddots \\ o_n^{t,1} \ldots o_n^{t,K} \end{bmatrix}$
6:   Batch of centralized information for each environment: $\begin{bmatrix} z^{t,1} \ldots z^{t,K} \end{bmatrix}$
7:   Sample actions $\begin{bmatrix} a_1^{t,1} \ldots a_1^{t,K} \\ \ddots \\ a_n^{t,1} \ldots a_n^{t,K} \end{bmatrix} \sim \pi(\cdot \mid h_1^t; \phi_1), \ldots, \pi(\cdot \mid h_n^t; \phi_n)$
8:   Apply actions; collect rewards $\begin{bmatrix} r_1^{t,1} \ldots r_1^{t,K} \\ \ddots \\ r_n^{t,1} \ldots r_n^{t,K} \end{bmatrix}$, observations $\begin{bmatrix} o_1^{t+1,1} \ldots o_1^{t+1,K} \\ \ddots \\ o_n^{t+1,1} \ldots o_n^{t+1,K} \end{bmatrix}$,
     and centralized information $\begin{bmatrix} z^{t+1,1} \ldots z^{t+1,K} \end{bmatrix}$
9:   **for** agent $i = 1, \ldots, n$ **do**
10:      **if** $s^{t+1,k}$ is terminal **then**
11:         $Adv(h_i^{t,k}, z^{t,k}, a_i^{t,k}) \leftarrow r_i^{t,k} - V(h_i^{t,k}, z^{t,k}; \theta_i)$
12:         Critic target $y_i^{t,k} \leftarrow r_i^{t,k}$
13:      **else**
14:         $Adv(h_i^{t,k}, z^{t,k}, a_i^{t,k}) \leftarrow r_i^{t,k} + \gamma V(h_i^{t+1,k}, z^{t+1,k}; \theta_i) - V(h_i^{t,k}, z^{t,k}; \theta_i)$
15:         Critic target $y_i^{t,k} \leftarrow r_i^{t,k} + \gamma V(h_i^{t+1,k}, z^{t+1,k}; \theta_i)$
16:      Actor loss $\mathcal{L}(\phi_i) \leftarrow \frac{1}{K} \sum_{k=1}^{K} Adv(h_i^{t,k}, z^{t,k}, a_i^{t,k}) \log \pi(a_i^{t,k} \mid h_i^{t,k}; \phi_i)$
17:      Critic loss $\mathcal{L}(\theta_i) \leftarrow \frac{1}{K} \sum_{k=1}^{K} \left( y_i^{t,k} - V(h_i^{t,k}, z^{t,k}; \theta_i) \right)^2$
18:      Update parameters $\phi_i$ by minimizing the actor loss $\mathcal{L}(\phi_i)$
19:      Update parameters $\theta_i$ by minimizing the critic loss $\mathcal{L}(\theta_i)$
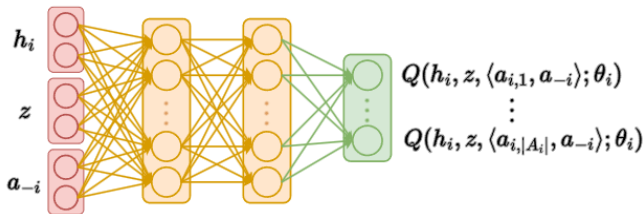
# Empirical Results



(a) Speaker-listener game

(b) Training curves

# Centralized Action-Value Critics



**Action-Value Loss:**

$$\mathcal{L}(\theta_i) = \left(y_i - Q(h_i^t, z^t, a^t; \theta_i)\right)^2 \quad \text{with} \quad y_i = r_i^t + \gamma Q(h_i^{t+1}, z^{t+1}, a^{t+1}; \theta_i)$$

# Today

- Intro
- Solving Joint Policies
- Independent -DQN, -REINFORCE, -A2C
- Centralized Critics
- **Value Decomposition (VDN, QMIX)**

# Value Decomposition
## Motivation

- Centralized action-value are difficult to learn
- Agents cannot select actions decentralized
- Greedy actions are costly
- Don't rely on additional policy nets

# Value Decomposition
**Insights**

- $Q_i(s, a_i) \approx Q(s, a)$ intractable
- Agent interaction as a sparse coordination graph
- Decompose the centralized action-value function (if $\mathcal{R}_i = \mathcal{R}_j, \ i, j \in I$)

$$Q(h^t, z^t, a^t; \theta) = \mathbb{E}\left[ \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r^\tau \mid h^t, z^t, a^t \right],$$

into simpler functions

# Individual-Global-Max (IGM) Property

**Idea:**
*Greedy actions w.r.t to the centralized function $\iff$
joint actions composed of individual greedy actions*

# Individual-Global-Max (IGM) Property

*First, formally define sets of greedy actions*:

$$A^*(h, z; \theta) = \arg\max_{a \in A} Q(h, z, a; \theta)$$

$$A_i^*(h_i; \theta_i) = \arg\max_{a_i \in A_i} Q(h_i, a_i; \theta_i)$$

**Definition 3.** *IGM is satisfied if the following holds for all full histories $\hat{h}$, joint-observation histories $h = \sigma(\hat{h})$, individual observation histories $h_i = \sigma_i(\hat{h})$ and centralized information z.*

$$\forall a = (a_1, \ldots, a_n) \in A : a \in A^*(h, z; \theta) \iff \forall i \in I : a_i \in A_i^*(h_i; \theta_i)$$

# Linear Value Decomposition

**Assumption.** *Decompose $r^t = \bar{r}_1^t, \ldots, \bar{r}_n^t$, where $\bar{r}_i^t$ is utility of agent $i$ at $t$. $\bar{r}_i^t$ is obtained by decomposition*

**Definition 4. (Linear Value Decomposition)**

$$
\begin{aligned}
Q(h^t, z^t, a^t; \theta) &= \mathbb{E}_{\hat{h}^t \sim \Pr(\cdot | \pi)} \left[ \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r^\tau \mid h^t = \sigma(\hat{h}^t), z^t, a^t \right] \\
&= \mathbb{E}_{\hat{h}^t \sim \Pr(\cdot | \pi)} \left[ \sum_{\tau=t}^{\infty} \gamma^{\tau-t} \left( \sum_{i \in I} \bar{r}_i^\tau \right) \mid h^t, z^t, a^t \right] \\
&= \sum_{i \in I} \mathbb{E}_{\hat{h}^t \sim \Pr(\cdot | \pi)} \left[ \sum_{\tau=t}^{\infty} \gamma^{\tau-t} \bar{r}_i^\tau \mid h^t, z^t, a^t \right] \\
&= \sum_{i \in I} Q(h_i^t, a_i^t; \theta_i)
\end{aligned}
$$

# Proof

# Value Decomposition Networks (VDN)

Given a buffer $\mathcal{D}$, the loss for a VDN is computed over a batch $\mathcal{B}$:

$$\mathcal{L}(\theta) = \frac{1}{B} \sum_{(h^t, a^t, r^t, h^{t+1}) \in \mathcal{B}} \left( r^t + \gamma \max_{a \in A} Q(h^{t+1}, a; \bar{\theta}) - Q(h^t, a^t; \theta) \right)^2$$

with

$$Q(h^t, a^t; \theta) = \sum_{i \in I} Q(h_i^t, a_i^t; \theta_i) \text{ and}$$

$$\max_{a \in A} Q(h^{t+1}, a; \bar{\theta}) = \sum_{i \in I} \max_{a_i \in A_i} Q(h_i^{t+1}, a_i; \bar{\theta}_i).$$

**Algorithm 21** Value decomposition networks (VDN)

1: Initialize $n$ utility networks with random parameters $\theta_1, \ldots, \theta_n$
2: Initialize $n$ target networks with parameters $\bar{\theta}_1 = \theta_1, \ldots, \bar{\theta}_n = \theta_n$
3: Initialize a shared replay buffer $D$
4: **for** time step $t = 0, 1, 2, \ldots$ **do**
5:     Collect current observations $o_1^t, \ldots, o_n^t$
6:     **for** agent $i = 1, \ldots, n$ **do**
7:         With probability $\epsilon$: choose random action $a_i^t$
8:         Otherwise: choose $a_i^t \in \arg\max_{a_i} Q(h_i^t, a_i; \theta_i)$
9:     Apply actions; collect shared reward $r^t$ and next observations $o_1^{t+1}, \ldots, o_n^{t+1}$
10:     Store transition $(h^t, a^t, r^t, h^{t+1})$ in shared replay buffer $D$
11:     Sample mini-batch of $B$ transitions $(h^k, a^k, r^k, h^{k+1})$ from $D$
12:     **if** $s^{k+1}$ is terminal **then**
13:         Targets $y^k \leftarrow r^k$
14:     **else**
15:         Targets $y^k \leftarrow r^k + \gamma \sum_{i \in I} \max_{a_i' \in A_i} Q(h_i^{k+1}, a_i'; \bar{\theta}_i)$
16:     Loss $\mathcal{L}(\theta) \leftarrow \frac{1}{B} \sum_{k=1}^{B} \left( y^k - \sum_{i \in I} Q(h_i^k, a_i^k; \theta_i) \right)^2$
17:     Update parameters $\theta$ by minimizing the loss $\mathcal{L}(\theta)$
18:     In a set interval, update target network parameters $\bar{\theta}_i$ for each agent $i$

# QMIX

- Linear decomposition has drawbacks
- Use a monotonic decomposition
- QMIX uses a mixing network

$$Q(h, z, a; \theta) = f_{mix}(Q(h_1, a_1; \theta_1), \ldots, Q(h_n, a_n; \theta_n); \theta_{mix})$$

- Loss:

$$\mathcal{L}(\theta) = \frac{1}{B} \sum_{(h^t, z^t, a^t, r^t, h^{t+1}, z^{t+1}) \in B} \left( r^t + \gamma \max_{a \in A} Q(h^{t+1}, z^{t+1}, a; \bar{\theta}) - Q(h^t, z^t, a^t; \theta) \right)^2$$

# Proof?

**Algorithm 22 QMIX**

1: Initialize $n$ utility networks with random parameters $\theta_1, \ldots, \theta_n$
2: Initialize $n$ target networks with parameters $\bar{\theta}_1 = \theta_1, \ldots, \bar{\theta}_n = \theta_n$
3: Initialize hypernetwork with random parameters $\theta_{\text{hyper}}$
4: Initialize a shared replay buffer $D$
5: **for** time step $t = 0, 1, 2, \ldots$ **do**
6:     Collect current centralized information $z^t$ and observations $o_1^t, \ldots, o_n^t$
7:     **for** agent $i = 1, \ldots, n$ **do**
8:         With probability $\epsilon$: choose random action $a_i^t$
9:         Otherwise: choose $a_i^t \in \arg\max_{a_i} Q(h_i^t, a_i; \theta_i)$
10:     Apply actions; collect shared reward $r^t$, next centralized information $z^{t+1}$
      and observations $o_1^{t+1}, \ldots, o_n^{t+1}$
11:     Store transition $(h^t, z^t, a^t, r^t, h^{t+1}, z^{t+1})$ in shared replay buffer $D$
12:     Sample mini-batch of $B$ transitions $(h^k, z^k, a^k, r^k, h^{k+1}, z^{k+1})$ from $D$
13:     **if** $s^{k+1}$ is terminal **then**
14:         Targets $y^k \leftarrow r^k$
15:     **else**
16:         Mixing parameters $\theta_{\text{mix}}^{k+1} \leftarrow f_{\text{hyper}}(z^{k+1}; \theta_{\text{hyper}})$
17:         Targets $y^k \leftarrow r^k + \gamma f_{\text{mix}} \begin{pmatrix} \max_{a_1'} Q(h_1^{k+1}, a_1'; \bar{\theta}_1) \\ \ddots \\ \max_{a_n'} Q(h_n^{k+1}, a_n'; \bar{\theta}_n) \end{pmatrix}; \theta_{\text{mix}}^{k+1}$
18:     Mixing parameters $\theta_{\text{mix}}^k \leftarrow f_{\text{hyper}}(z^k; \theta_{\text{hyper}})$
19:     Value estimates $Q(h^k, z^k, a^k; \theta) \leftarrow f_{\text{mix}}\left(Q(h_1^k, a_1^k; \theta_1), \ldots, Q(h_n^k, a_n^k; \theta_n); \theta_{\text{mix}}^k\right)$
20:     Loss $\mathcal{L}(\theta) \leftarrow \frac{1}{B} \sum_{k=1}^{B} \left( y^k - Q(h^k, z^k, a^k; \theta) \right)^2$
21:     Update parameters $\theta$ by minimizing the loss $\mathcal{L}(\theta)$
22:     In a set interval, update target network parameters $\bar{\theta}_i$ for each agent $i$