

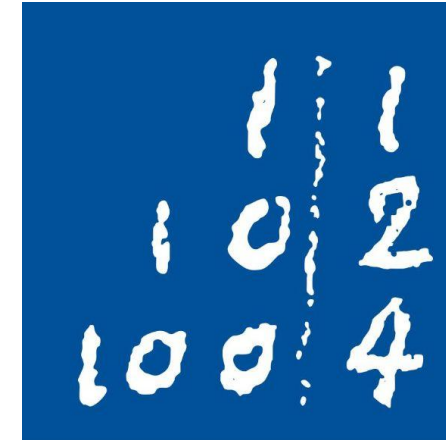


MetaIRN

Meta-Gradient-based Intrinsic Rewards via Attention Network

Philipp Link, Julius Heidmann, Luan Liebig-Schultz

{p.link, julius.heidmann, luan.liebig-schultz}@stud.uni-hannover.de



Leibniz
Universität
Hannover



1

Motivation

- **Sparse Reward Environments** make Learning for Agents difficult
→ **slow/unstable**
- How to **facilitate Training**?
→ Use **Intrinsic Rewards** to **enrich Feedback** through **Temporal Credit Assignment (TCA)**
- Which **Actions** lead to **Future Success**?
- **TCA**: Associate Actions with Temp. distant Rewards and assign Credit (int. Rewards)

3

Approach

- Can a **Combination** of **Meta-Gradient Optimization** and **Attention-based intrinsic Rewards** enhance the **Agent's Ability** to **assign Credit** in **TCA-only scenarios**?
- Based on **SAIR¹/ LIRPG² Meta-Gradient**:
 - **Inner Loop**: PPO Update (Ex.&In. Rew.)
 - **Outer Loop**: Attention Net Update (Ex. Rew.)
- This Agent is trained in **Umbrella Length³ (TCA-Only)** Env, which is **isolated** to **focus exclusively on TCA problems**
- **Comparison to Vanilla PPO Performance**

Meta-Gradient Update¹:

PPO Network:

$$\theta' \leftarrow \theta + \alpha \nabla_{\theta} J^{ex+in}(\theta | \mathcal{D})$$

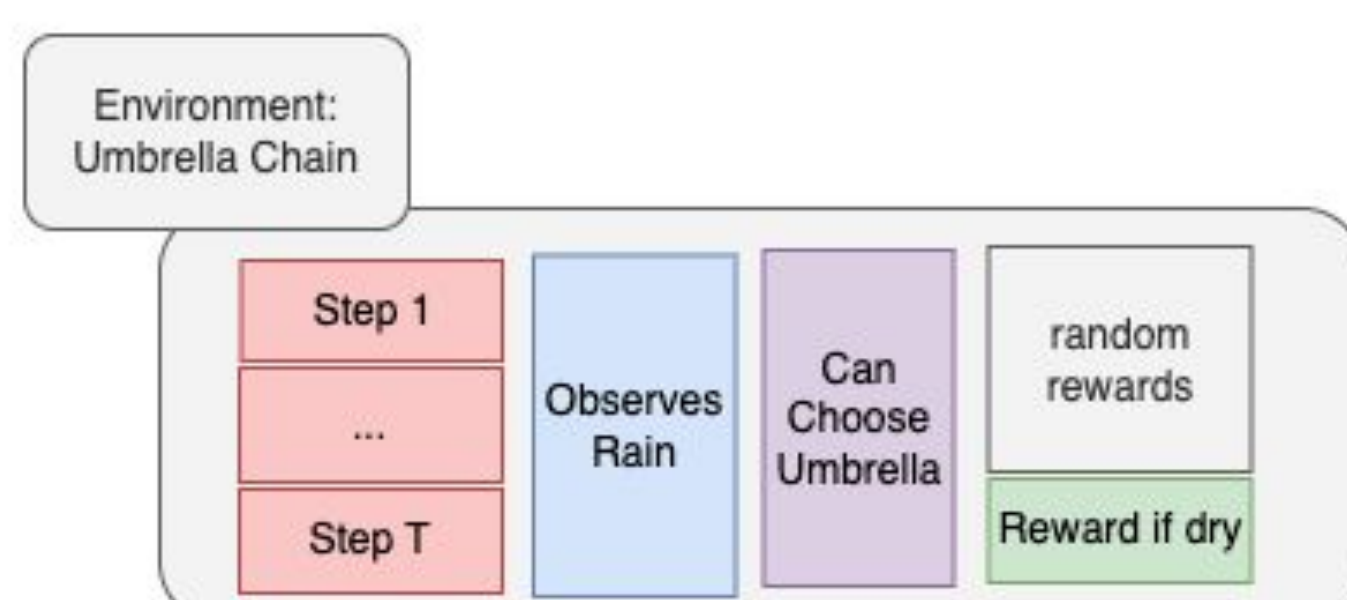
Intrinsic Network:

$$\eta' \leftarrow \eta + \alpha \nabla_{\eta} J^{ex}(\theta' | \mathcal{D}) \nabla_{\eta} \theta'$$

4

Umbrella Length³ Env

- **Fully observable MDP**
- **Intermediate Reward** $R_i \in \{-0.1, 0.1\}$
- **Final Reward** $R_f \in \{-10, 10\}$ depends only on First Action



5

Future Works

- **TMaze Env** (Memory-Only Test)
- **Meta-Gradient Implementation**
- **Research Ext. Value Pred. Behaviour**

2

Related Work & Setting

SAIR¹/ LIRPG² [Jiang et al. 2021], [Zheng et al. 2018]

- **Meta-Gradient** to train an **Attention Network** producing **Intrinsic Rewards** to facilitate PPO learning

Decoupling Memory from **TCA⁴** [Ni et al. 2023]

- Exactly **distinguish** between **TCA** and **Memory ability** of an Agent
→ **Memory Problem⁴**:

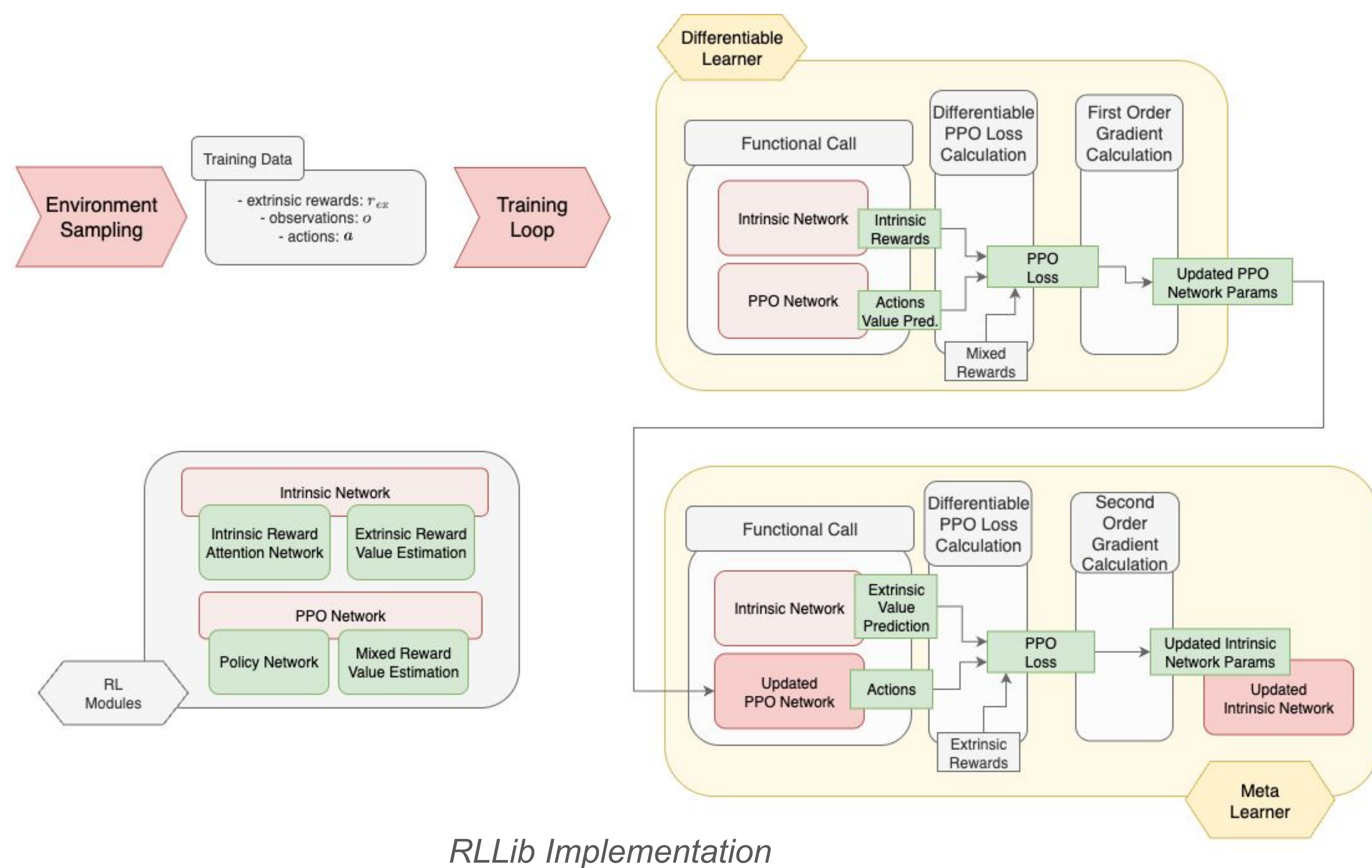
Definition 1.A (Reward memory length $m_{\text{reward}}^{\mathcal{M}}$). For a POMDP \mathcal{M} , $m_{\text{reward}}^{\mathcal{M}}$ is the smallest $n \in \mathbb{N}$ such that the expected reward conditioned on recent n observations is the same as the one conditioned on full history, i.e., $\mathbb{E}[r_t | h_{1:t}, a_t] = \mathbb{E}[r_t | h_{t-n+1:t}, a_t], \forall t, h_{1:t}, a_t$.

→ **Credit Assignment Problem⁴**:

$$c(h_{1:t}; \pi) := \min_{1 \leq n \leq T-t+1} \{n \mid \exists a_t^* \in A_t^*, \text{ s.t. } G_n^{\pi}(h_{1:t}, a_t^*) > G_n^{\pi}(h_{1:t}, a_t'), \forall a_t' \notin A_t^*\}$$

4

Key Insights



- **PPO with Intrinsic Attention & Vanilla**
- **Hyperparameter** based on **SMAC HPO** with **100 Configurations** each
- **10 Seeds Evaluations** for each **Environment Length (2, 5, 10, 20, 50)**

