



Scalable Distributed Reinforcement Learning in Multi-Agent Environments

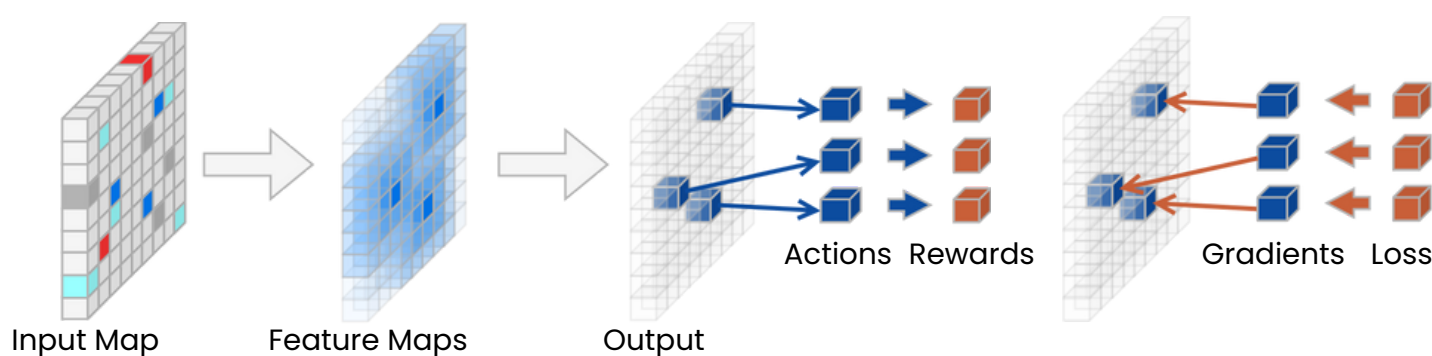
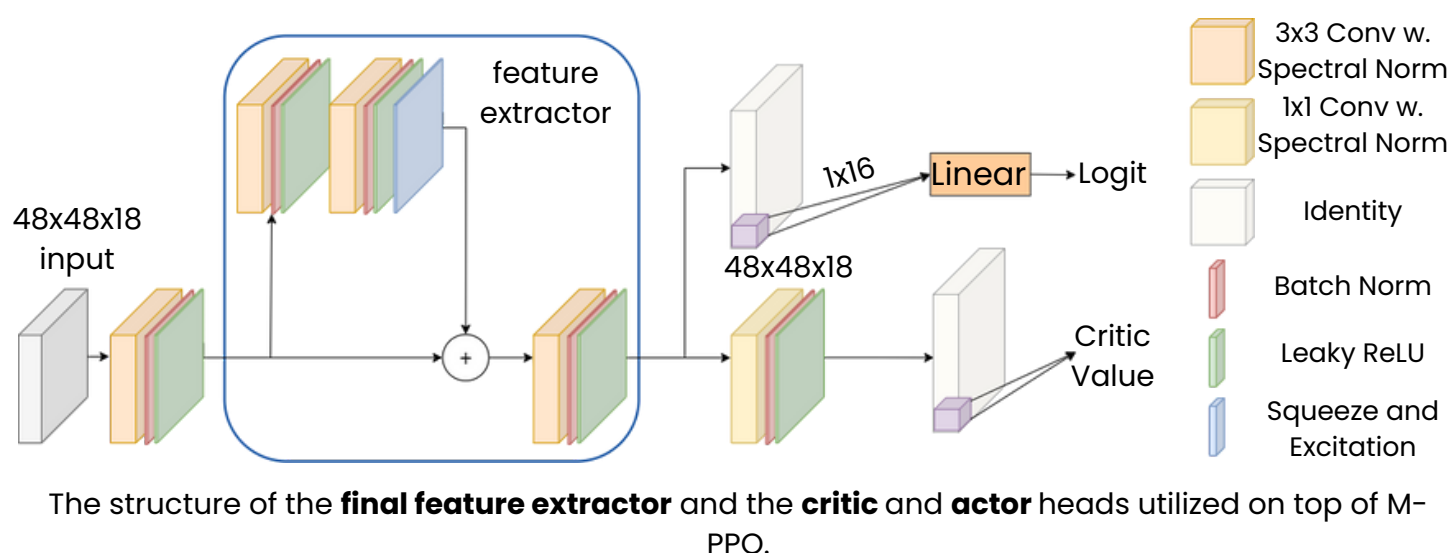
Takács Tamás, Magyar Gergely, Gulyás László
Eötvös Loránd University, Faculty of Informatics, Department of Artificial Intelligence

Novel reinforcement learning algorithms often face **scalability and compatibility issues in multi-agent environments** due to their optimization for single-agent settings. The lack of standardized methods for adaptation restricts their broader applicability, especially when dealing with rapidly changing numbers of controllable entities and massive scaling. Challenges include credit assignment, extensive memory usage, and increased computational time, leading to slow, destabilized training and suboptimal resource utilization. We propose a hybrid architecture, combining monolithic and distributed approaches, resulting in a **30-times reduction in model size** and learning basic skills **24 times faster with 600-times fewer training examples** compared to related works in the same environment. We also introduce **trajectory separation**, achieving a **3-times speed increase in training** convergence. Our **M-PPO-based hybrid** model achieves:

- a performance-based environment
- **separated agent trajectories**
- separated reward, advantage, probability and entropy calculation
- **reinforced positive behavior in early training**

Episode 1								Episode 2							
Step 1	Step 2	Step 3	Step 4	Step 5	Step 6	Step 7	Step 8	Step 1	Step 2	Step 3	Step 4	Step 5	Step 6	Step 7	Step 8
4.2	3.1	5.6	14.3	5.3	2.8	6.1	15.2	1.2	1.1	1.5	-	0.8	0.9	1.8	7.1
-	0.0	3.3	6.8	1.3	1.2	-	-	-	0.0	3.3	6.8	1.3	1.2	-	-
1.3	0.8	0.8	7.5	-	-	2.4	8.1	1.3	0.8	0.8	7.5	-	-	2.4	8.1
1.7	1.2	-	-	2.1	0.7	1.9	-	1.7	1.2	-	-	2.1	0.7	1.9	-

Showing the difference between **global rewards** (left) and **rewards distributed into groups** using trajectory separation (right)



Our **novel trajectory separation method** achieves distributed rewards by producing one trajectory per entity or cluster and calculating loss based on individual performance.

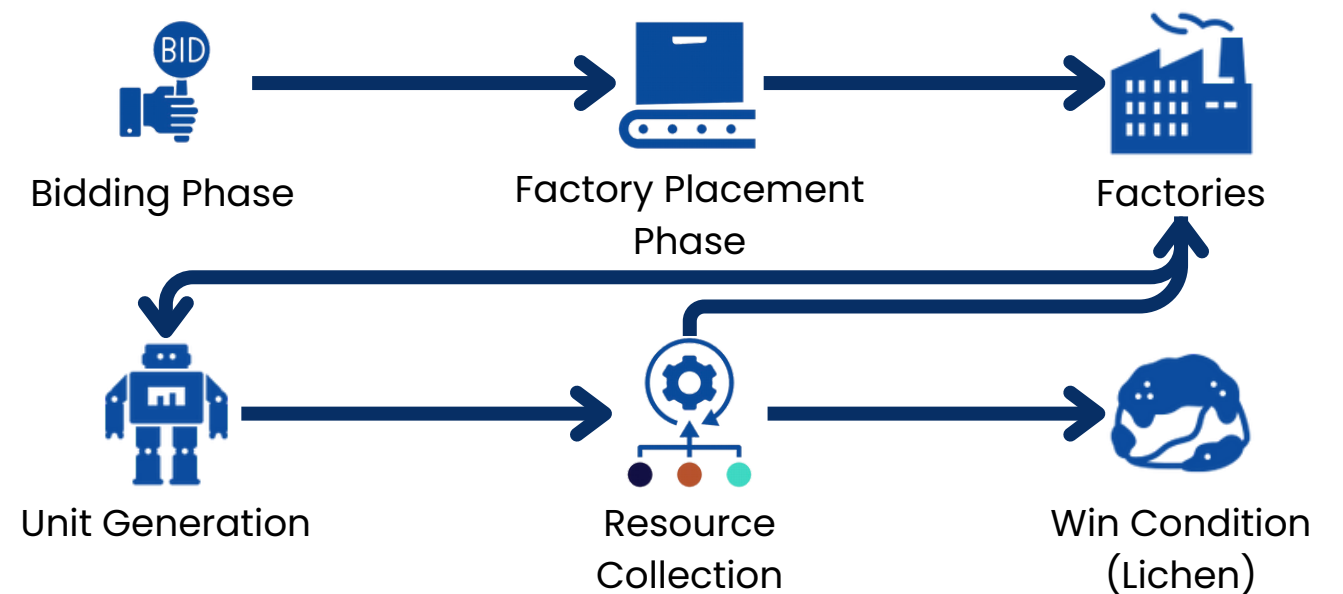
Policy Loss			Average=0.77
Group 1	0.67	-	1.50
Group 2	-	-1.33	2.50
Group 3	-4.00	5.60	3.00
Group 4	-5	-	4.00
	Step 1	Step 2	Step 3

$\frac{\text{New Probabilities}}{\text{Stored Probabilities}} \times \text{Calculated Advantage} =$

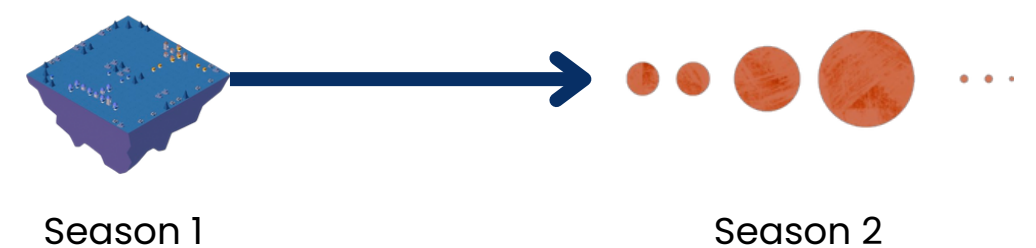
Demonstrating policy loss calculation with the extended dimension. The matrices on the image represent a mini-batch consisting of 3 environment steps. Inactive groups are marked.

Environment:

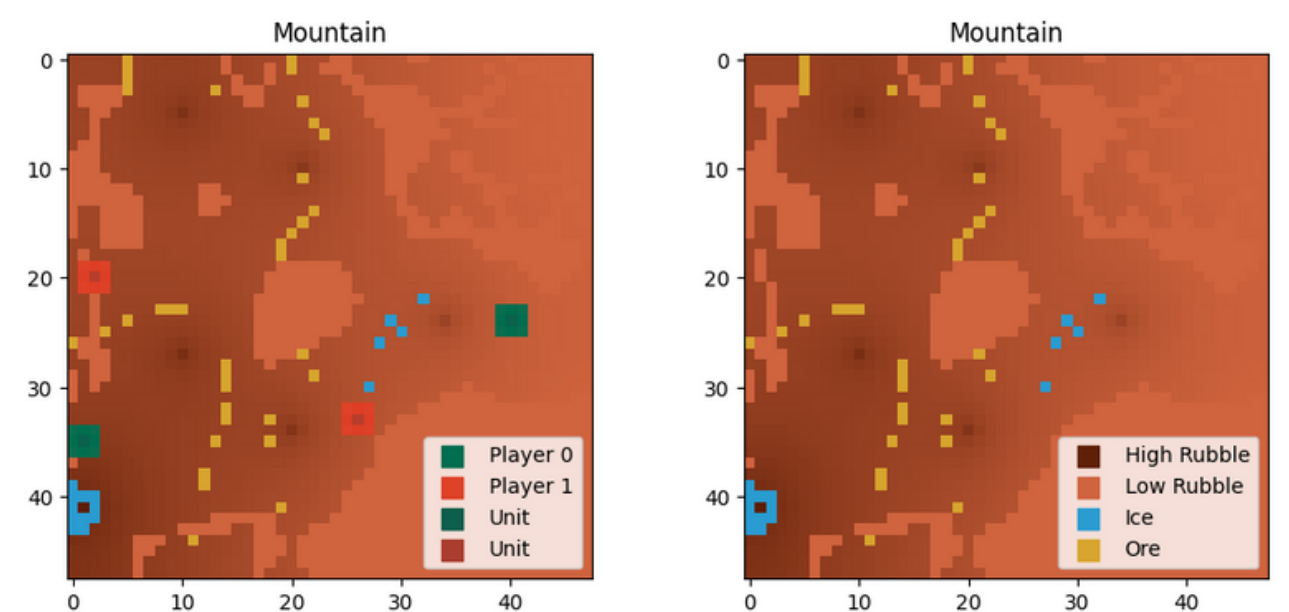
The **Lux AI Environment** represents a 2D grid platform tailored for MARL research [1].



Simplified environment loop of the **Lux AI Environment**. An episode starts with the **Bidding Phase** and ends after **1000** steps. The player with the most lichen collected wins.

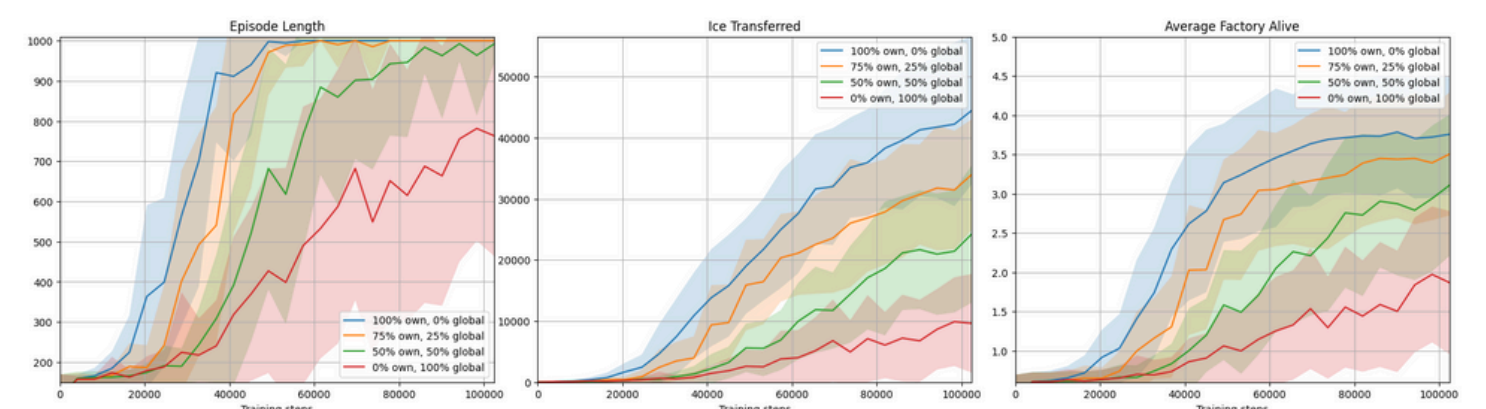


The **Lux AI Competition** has already hosted two seasons on Kaggle, one set in a forest environment and the other on Mars.



A visual representation of the grid environment, generated from a specific seed, showcasing a **Gaussian filter-based factory placement** policy on two identical maps.

In our study, **we compared our results with other competition submissions, conducted experiments on various credit assignment setups, and performed ablation studies.**



Results demonstrate the effectiveness of trajectory separation, revealing that higher agent-specific reward weights heavily improve performance. **Episode Length** indicates team survival duration, while **Ice Transferred** and **Average Factory Alive** reflect resource collection efficiency.

References:

- [1] Chen, H. et al. Emergent collective intelligence from massive-agent cooperation and competition 2023.
- [2] Tao, S., Pan, I., et al. Lux AI Season 2 2023. <https://kaggle.com/competitions/luxai-season-2>.

Acknowledgements:

This work was supported by the Institute for **Business Cooperation Scholarship** under grants #EIIISIVE. We extend our gratitude to **Eötvös Loránd University and the Department of Artificial Intelligence** for providing the resources essential for this project's development.

