

# REINFORCEMENT LEARNING WITH CUSTOMISED GYMNASIUM ENVIRONMENTS

## Assignment 2

Trabalho realizado por:

- Inês Castro up202304060
- Simão Gomes up202304752
- Soraia Costa up202305078

# ÍNDICE

01

## AMBIENTE GYMNASIUM

Descrição do funcionamento do jogo e dos desafios encontrados

02

## MUDANÇAS

Descrição das alterações introduzidas no ambiente Gymnasium original

03

## ALGORITMOS

Algoritmos de aprendizagem por reforço utilizados

04

## ANÁLISE

Metodologia experimental e critérios de avaliação

05

## CONCLUSÃO

Síntese do trabalho desenvolvido e perspectivas futuras



# AMBIENTE GYMNASIUM: ICEHOCKEY-V5

- Jogo de hóquei no gelo para **dois jogadores (agente vs adversário)**
- O adversário é controlado pelo ambiente
- **Objetivo:** Marcar mais golos do que o adversário
- O agente controla **um único jogador**
- O jogo decorre em **tempo contínuo**, com decisões a cada timestep
- A dinâmica do jogo envolve:
  - Movimento constante
  - Disputa pela posse do disco
  - Ataque e defesa simultâneos
- O agente deve aprender:
  - Quando atacar
  - Quando defender
  - Como posicionar-se corretamente



# CARACTERÍSTICAS DO AMBIENTE GYMNASIUM

## Percepções

- Frames RGB do ecrã do jogo
- Informação visual parcial, sem acesso ao estado interno do jogo

## Estados

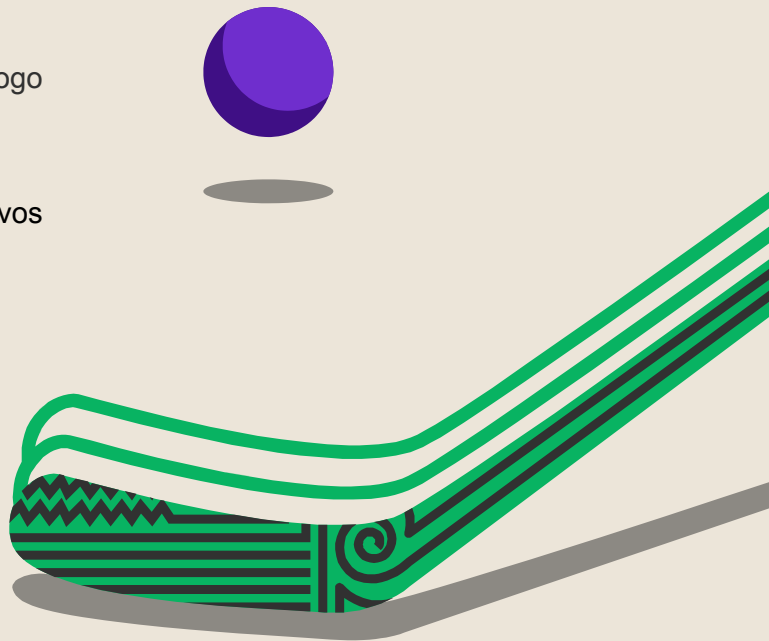
- Construídos através do empilhamento de 4 frames consecutivos
- Permite inferir movimento e direção
- Necessário para capturar a dinâmica temporal do jogo

## Ações

- Espaço de ações discreto com 18 ações
- Combinações de movimento e ação FIRE
- Ação FIRE essencial para interagir com o disco

## Recompensas (ambiente original)

- +1 ao marcar golo
- -1 ao sofrer golo
- Longos períodos sem qualquer feedback



# MUDANÇAS INTRODUZIDAS NO AMBIENTE

Propostas Iniciais (Design Conceptual)	Implementação Final (Após Testes Empíricos)
<ul style="list-style-type: none"><li>• Manter a <b>reward objetivo original</b> - +1 ao marcar golo, -1 ao sofrer golo e Garantir alinhamento com o objetivo real do jogo</li><li>• Introduzir <b>orientação direcional</b> - Recompensar progresso do disco em direção à baliza adversária</li><li>• Incentivar <b>posse do disco</b> - Pequena recompensa por manter controlo</li><li>• Incentivar <b>pressão defensiva</b> - Recompensar redução da distância ao disco</li><li>• Reforçar <b>ações taticamente relevantes</b> - Passes, remates enquadrados e roubos de disco</li><li>• Penalizar <b>comportamentos degenerados</b> - Inatividade prolongada e Proximidade excessiva entre colegas</li></ul>	<ul style="list-style-type: none"><li>• <b>Reward objetivo mantida</b> +1 ao marcar golo, -1 ao sofrer golo</li><li>• <b>Penalização temporal removida</b> Provocava instabilidade no treino</li><li>• <b>Incentivo leve ao movimento</b> Recompensa para ações de deslocamento</li><li>• <b>Incentivo moderado à ação FIRE</b> Promove interação sem forçar comportamento</li><li>• <b>Penalização por repetição excessiva</b> Evita loops e políticas degeneradas</li><li>• <b>Penalização por ausência de remates</b> Desencoraja passividade prolongada</li></ul>

# ALGORITMOS DE REINFORCEMENT LEARNING

A utilização de múltiplos algoritmos permitiu analisar como diferentes paradigmas de aprendizagem reagem às alterações introduzidas no ambiente.

## ■ PPO (ON-POLICY)

conhecido pela sua estabilidade e robustez, especialmente adequado a ambientes complexos e a funções de recompensa modificadas.

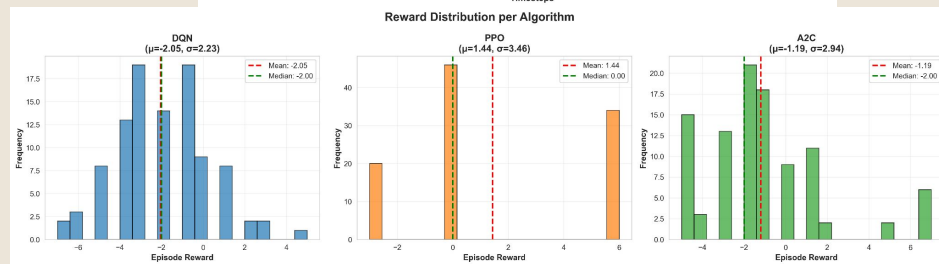
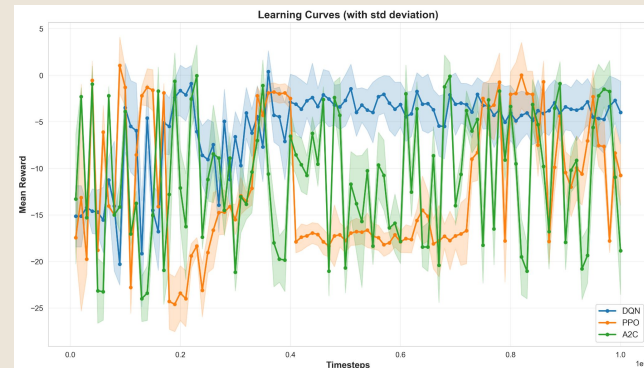
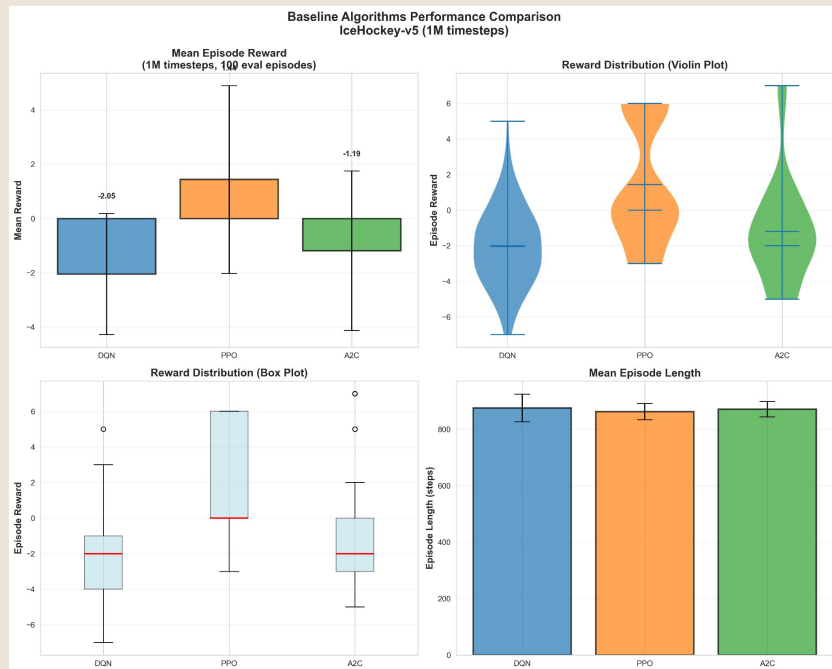
## ■ DQN (OFF-POLICY)

baseado em aprendizagem de valores Q, eficiente em termos de dados, mas mais sensível a alterações na função de recompensa.

## ■ A2C (ON-POLICY)

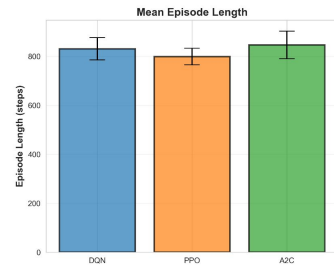
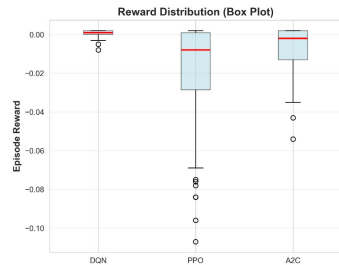
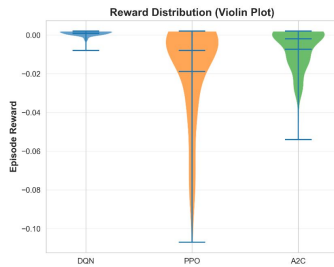
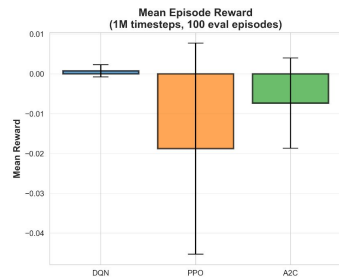
simples e eficiente no arranque da aprendizagem, mas com menor estabilidade a longo prazo.

# ANÁLISE DE RESULTADOS: BASELINE ENVIROMENT

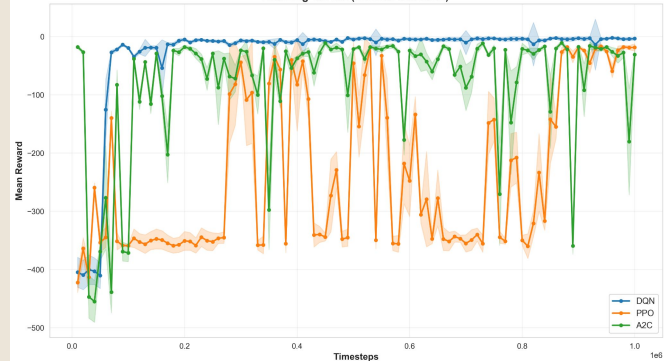


# ANÁLISE DE RESULTADOS: CUSTOM ENVIROMENT

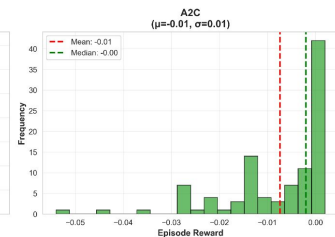
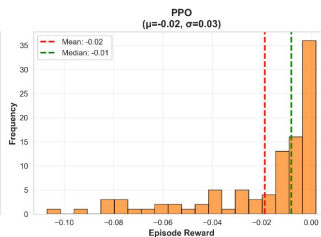
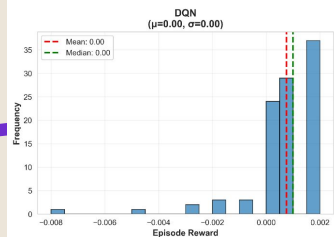
Custom Environment - Algorithms Performance Comparison  
IceHockey-v5 (1M timesteps)



Learning Curves (with std deviation)



Reward Distribution per Algorithm

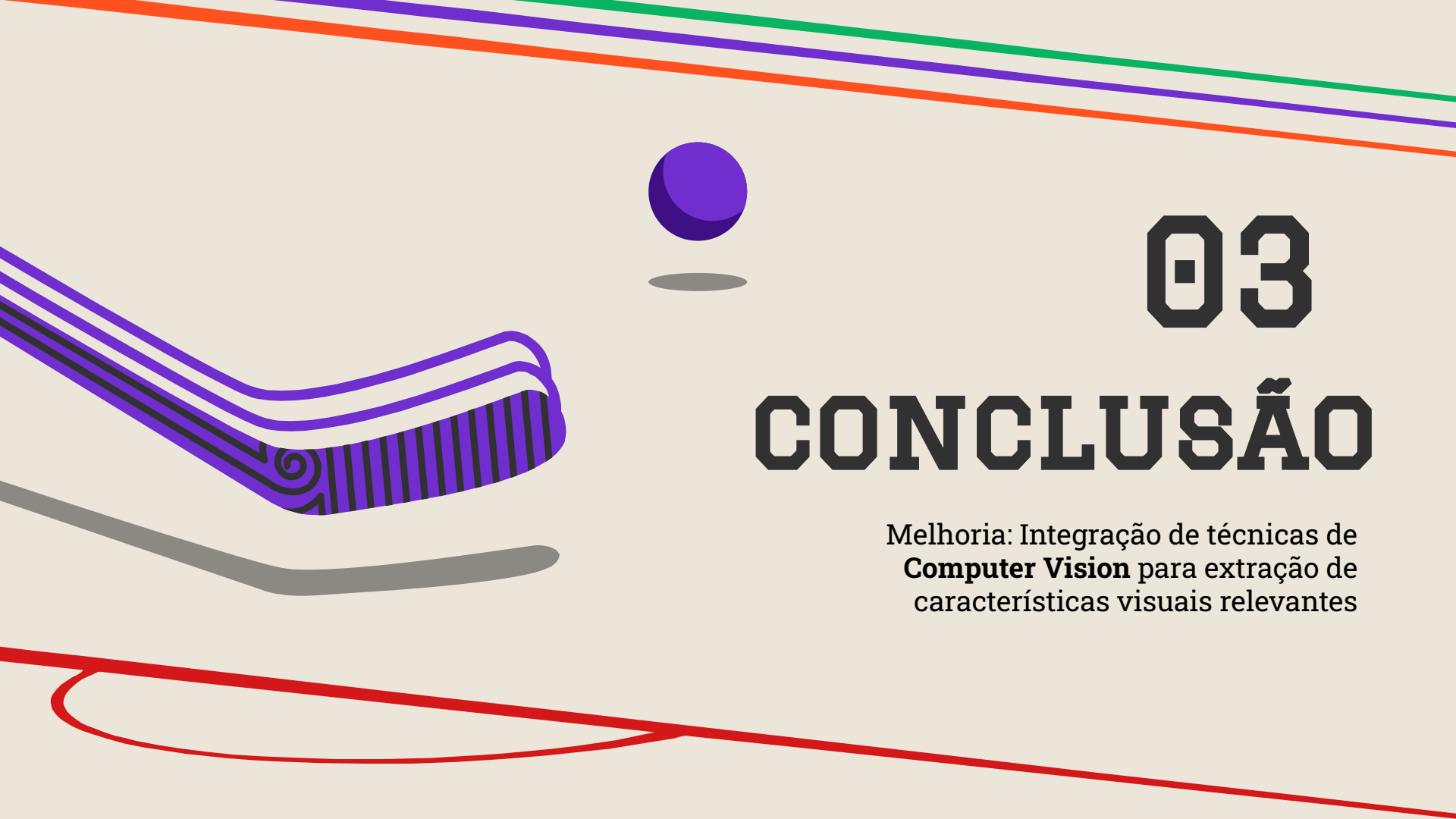




# ANÁLISE DE RESULTADOS:

Algoritmo	Baseline (Média ± Desvio)	Custom (Média ± Desvio)	Variação
DQN	-2.05 ± 2.23	+0.001 ± 0.002	↑ Melhoria significativa
PPO	+1.44 ± 3.46	-0.02 ± 0.03	↓ Quebra de desempenho
A2C	-1.19 ± 2.94	-0.01 ± 0.01	↑ Melhoria moderada





03

# CONCLUSÃO

Melhoria: Integração de técnicas de  
**Computer Vision** para extração de  
características visuais relevantes



**VÍDEO**

