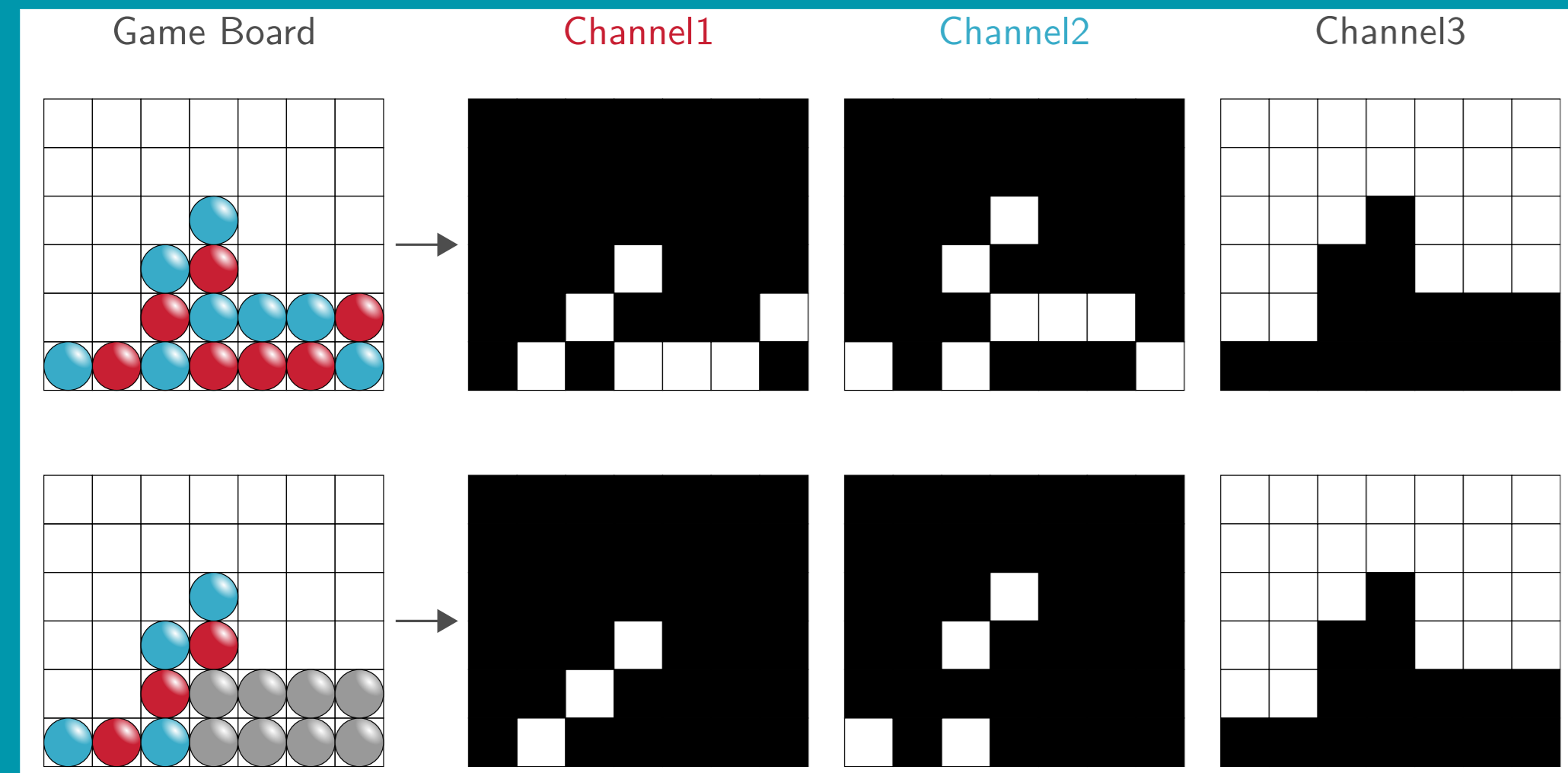


Overview

- Goal:** Evaluate saliency attribution methods in XAI: Which identify the important parts of a classifier input?
- Problem:** Hiding the relevant part of the input with baseline values or noise means evaluating the classifier off-manifold.
- Solution:** Train neural agents on abstract games with the same hidden information that is used in the evaluation.

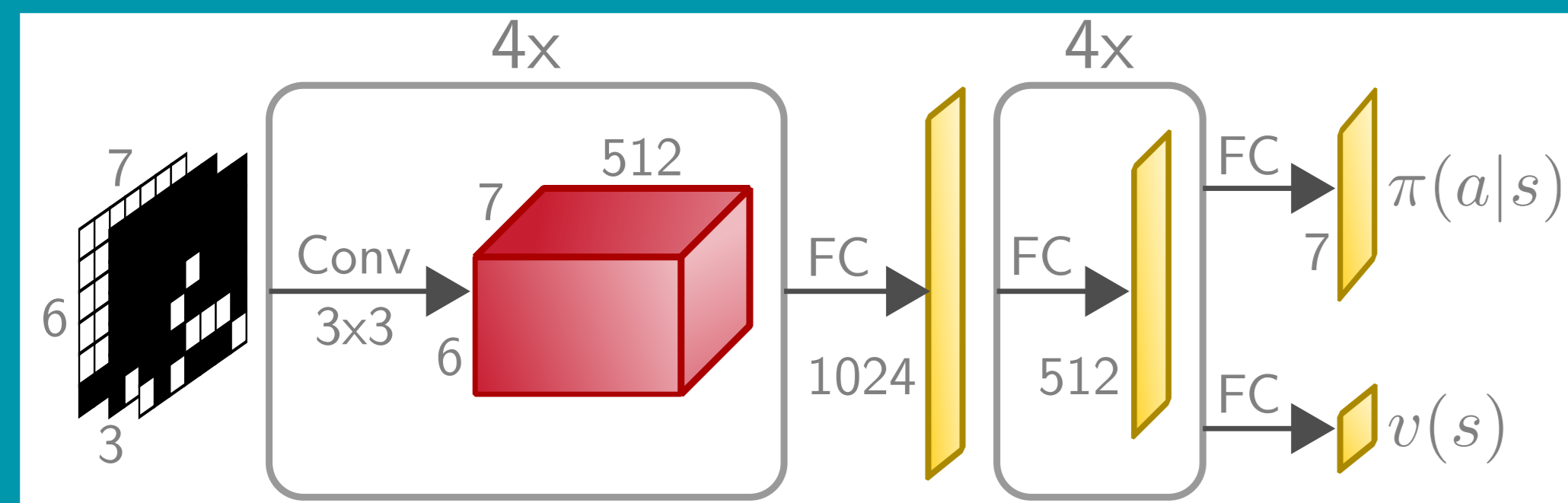
The Setup

Connect Four with randomly hidden colour information.



The Model

DNN with policy and value head for Proximal Policy Optimisation.



Advantages

- It is clear how to realise hidden information
- No off-manifold evaluation is necessary
- We can let XAI-methods play against each other

Training Characteristic Functions with Reinforcement Learning: XAI-methods play Connect Four

Stephan Wäldchen
waeldchen@zib.de

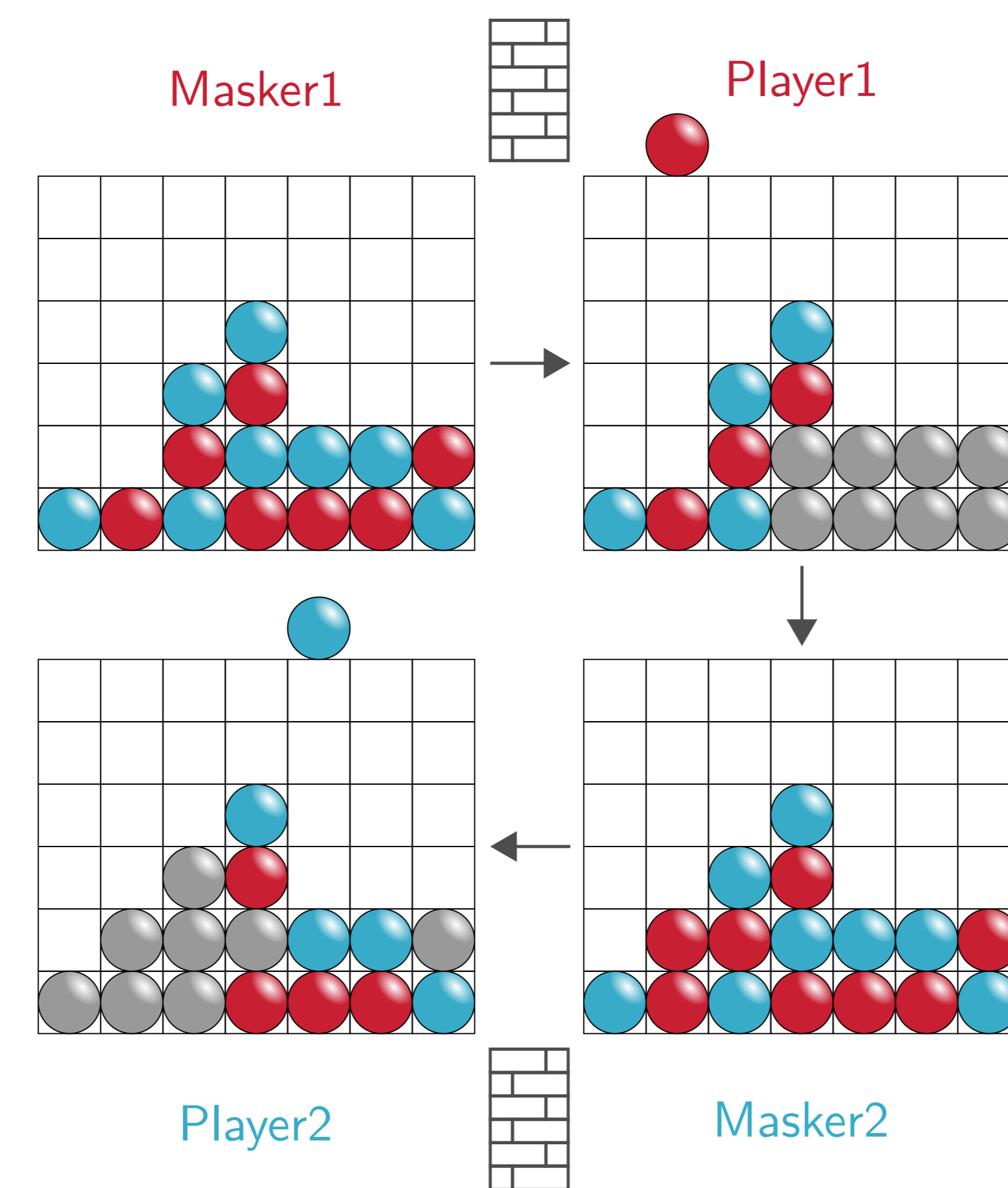
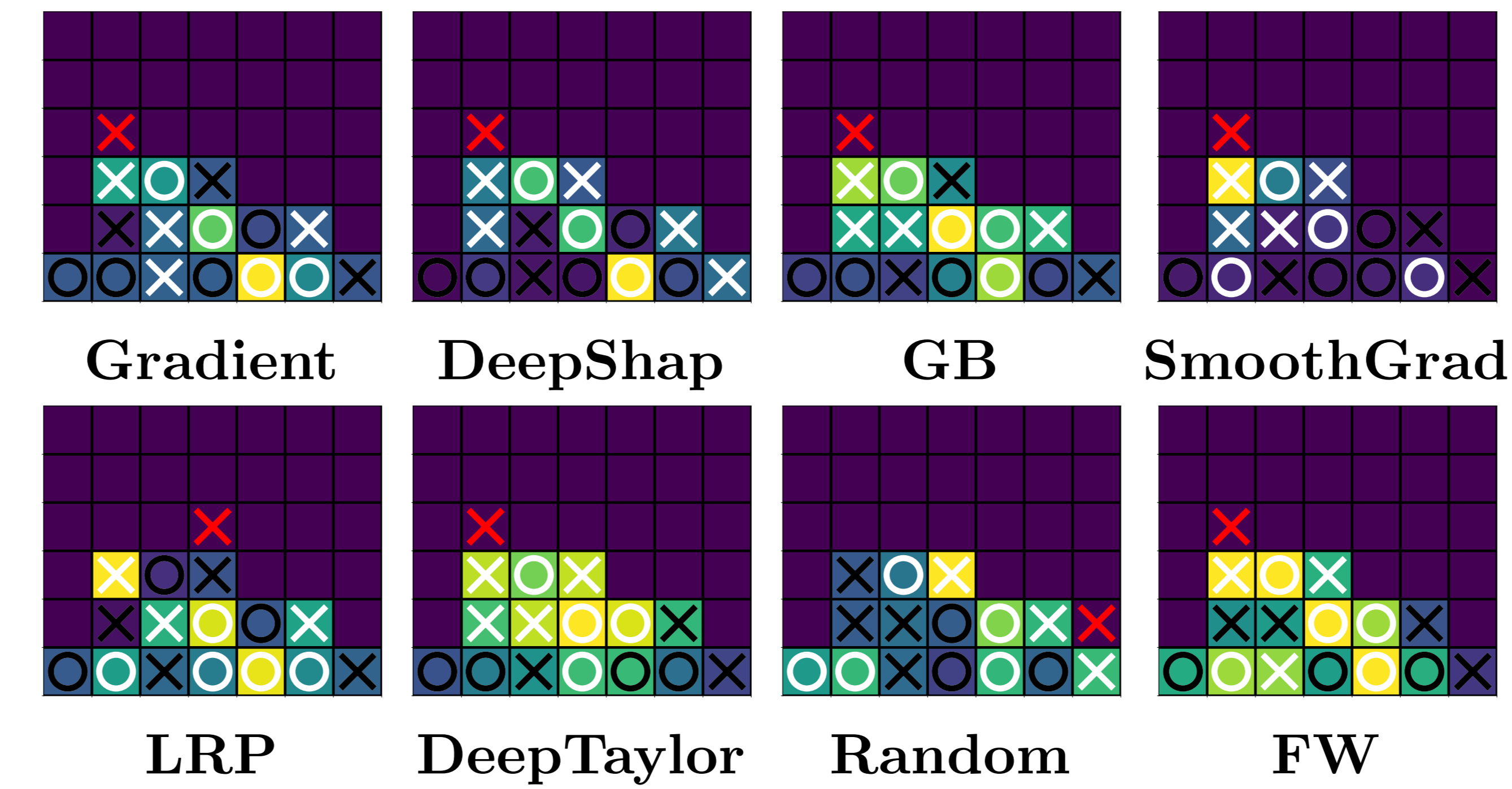
Felix Huber
huber@zib.de

Sebastian Pokutta
pokutta@zib.de



We let Explainable AI (XAI)-methods play against each other in Connect Four to compare them!

1. XAI-methods explain the output of the policy network with a heatmap of importance values.



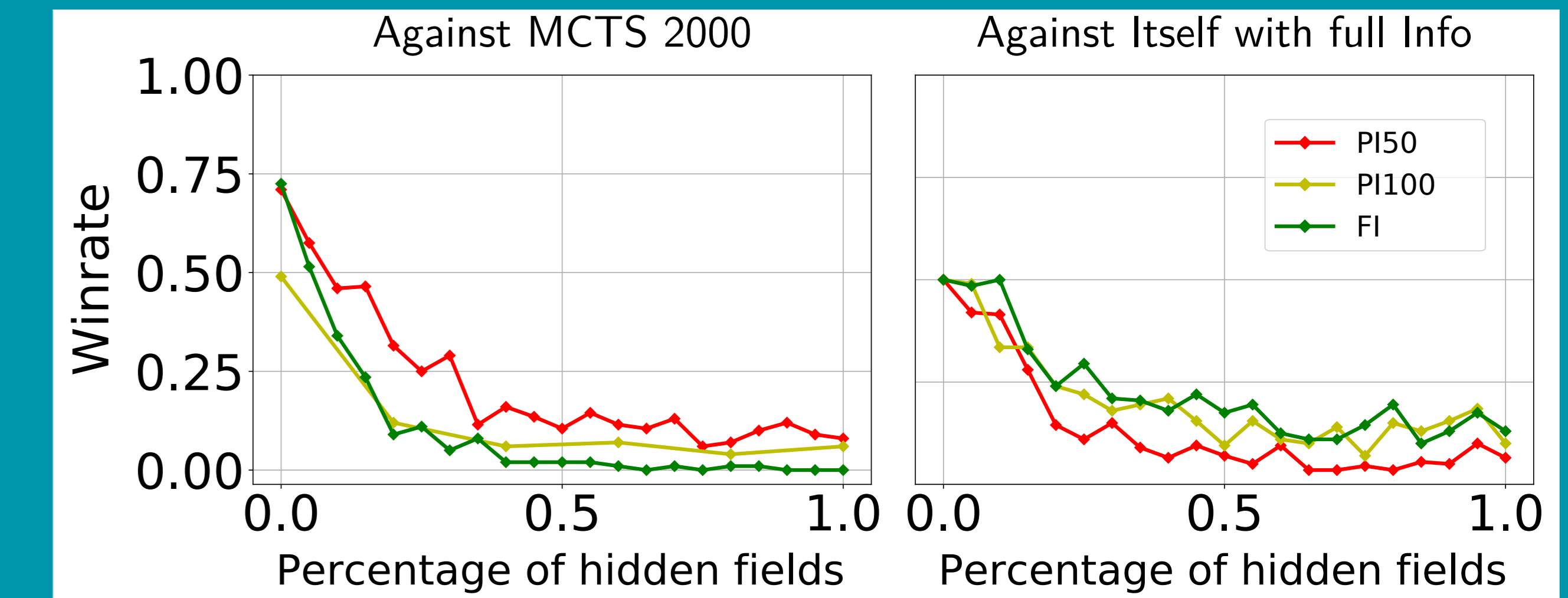
2. Select the 50% most important colour features and hide the rest for the agents next move.

3. Measure win rates of players supported by each XAI-method.

Input	DeepShap	Guided Backprop	FW	Gradient	LRP-ε	Deep Taylor	Smooth Grad	Random
DeepShap	0.263	0.498	0.499	0.573	0.609	0.567	0.742	0.871
Guided Backprop	0.254	0.502	0.496	0.579	0.616	0.616	0.765	0.861
FW	0.294	0.501	0.503	0.587	0.604	0.586	0.742	0.848
Gradient	0.217	0.427	0.421	0.413	0.512	0.502	0.686	0.819
LRP-ε	0.167	0.392	0.385	0.397	0.487	0.472	0.677	0.810
Deep Taylor	0.226	0.433	0.384	0.414	0.497	0.528	0.681	0.805
Smooth Grad	0.113	0.259	0.235	0.258	0.314	0.323	0.319	0.676
Random	0.065	0.130	0.140	0.152	0.181	0.190	0.195	0.324

Randomly Hiding Colour Features

Every turn $t \in [42]$ we sample $p_h \sim \mathcal{U}([0, p_h^{\max}])$ and show colour features of $\lfloor p_h t \rfloor$ random game pieces each turn. We train three agents with varying p_h^{\max} .



The PI-50 agent with $p_h = 0.5$ performs well for all levels of hidden information and was chosen for our tournament.

Characteristic Functions

A concept from cooperative game theory. We can define

$$\nu^{\text{pol}}(S) := \pi(a^{\max}; \mathbf{x}^{(S)}) \quad \text{and} \quad \nu^{\text{val}}(S) := v(\mathbf{x}^{(S)}),$$

where $\mathbf{x} \in \{0, 1\}^{3 \times 6 \times 7}$ and $\mathbf{x}^{(S)} := [\mathbf{x}_S^1, \mathbf{x}_S^2, \mathbf{x}_S^3]$,

and calculate Shapley Values (SV)

$$\text{SV: } \phi_{\nu,i} = \sum_{S \subseteq [d] \setminus \{i\}} \binom{d-1}{|S|}^{-1} (\nu(S \cup \{i\}) - \nu(S)),$$

and Prime Implicant Explanations (PIE)

$$\text{PIE: } S^* = \operatorname{argmin}_S |S| \quad \text{subject to} \quad \nu(S) = \nu([d]).$$

PIE with Frank-Wolfe

We solve a convex relaxation of the PIE-problem

$$\mathbf{s}^* := \operatorname{argmin}_{\mathbf{s}} (\pi(a^* | \mathbf{x}^{(\mathbf{s})}) - \pi(a^* | \mathbf{x}))^2 \quad \text{subject to} \quad \|\mathbf{s}\|_1 \leq k,$$

where $\mathbf{s} \in [0, 1]^{3 \times 6 \times 7}$ and $\mathbf{x}^{(\mathbf{s})} := [\mathbf{s} \odot \mathbf{x}^1, \mathbf{s} \odot \mathbf{x}^2, \mathbf{x}^3]$.

with a Frank-Wolfe solver and define the XAI-method FW.

Gradient [Simonyan et al. 2013]

SmoothGrad [Smikov et al. 2017]

GuidedBackprop [Springenberg et al. 2015]

DeepShap [Lundberg and Lee 2017]

LRP [Bach et al. 2015]

DeepTaylor [Montavon et al. 2018]