

Explainability in Reinforcement learning

```
graph TD; A[Explainability in Reinforcement learning] --> B(Transparent/Intrinsic methods); A --> C(Post-hoc explanation methods); B --> D[• Representaiton learning<br>• Simultaneous learning of the explanation and the policy<br>• Sub-task decomposition<br>• Human-robot interactions<br>• Argument accelerated reinforcement learning]; C --> E[• Saliency map<br>• Interaction with the environment during training<br>• Transparent model approximations<br>• Counterfactual explanations];
```

Transparent/Intrinsic methods

- Representaiton learning
- Simultaneous learning of the explanation and the policy
- Sub-task decomposition
- Human-robot interactions
- Argument accelerated reinforcement learning

Post-hoc explanation methods

- Saliency map
- Interaction with the environment during training
- Transparent model approximations
- Counterfactual explanations