

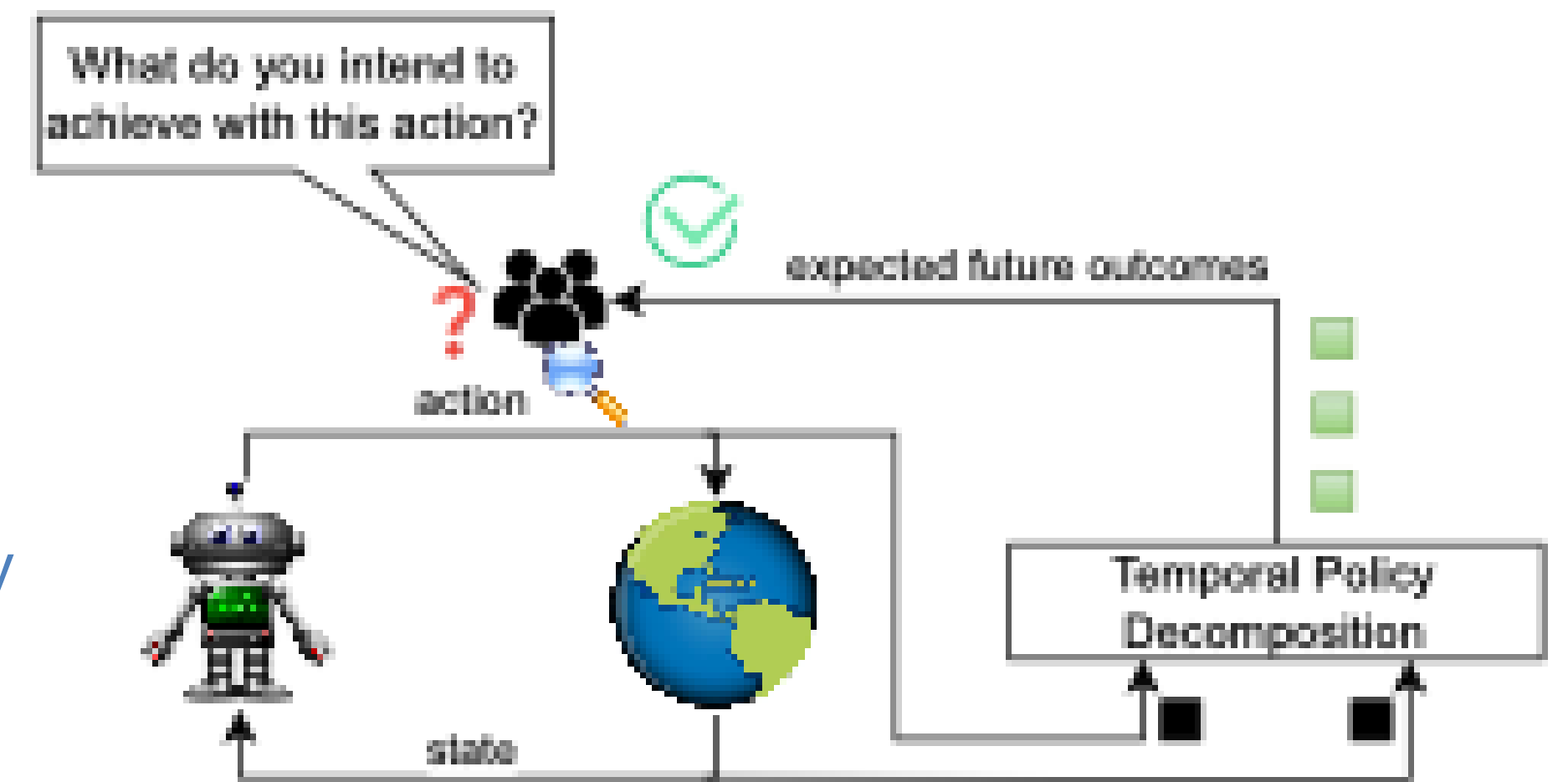
EXPLAINABLE REINFORCEMENT LEARNING VIA TEMPORAL POLICY DECOMPOSITION

Franco Ruggeri, KTH & Ericsson Research
Supervisors: Karl Henrik Johansson, Rafia Inam



WHAT DO YOU INTEND TO ACHIEVE?

- Reinforcement Learning (RL) involves **sequential decision-making**.
- An RL agent acts to maximize cumulative rewards.
- An RL action could be an intermediate step for a **delayed reward**.
- RL actions should be **explained** in terms of the **future trajectory** implied by the action.
- Existing eXplainable RL (XRL) methods focus on explaining the state-action mapping, neglecting the temporal dimension.



TEMPORAL POLICY DECOMPOSITION

- Learn Fixed-Horizon Generalized Value Function [1] (off-policy)

$$\hat{G}_h = r_t + \gamma \hat{Q}_{h-1}(s_{t+1}, a_{t+1}), \quad a_{t+1} \sim \pi(\cdot | s_{t+1})$$

$$\hat{Q}_h^{(t+1)}(s_t, a_t) \leftarrow \hat{Q}_h^{(t)}(s_t, a_t) + \alpha_{t,h}(s_t, a_t) [\hat{G}_h - \hat{Q}_h^{(t)}(s_t, a_t)]$$

- Decompose FHGVFs into Expected Future Outcomes (EFOs):

$$\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & \gamma & 0 & \dots & 0 \\ 1 & \gamma & \gamma^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \gamma & \gamma^2 & \dots & \gamma^{H-1} \end{bmatrix} \begin{bmatrix} O_0(s, a) \\ O_1(s, a) \\ O_2(s, a) \\ \vdots \\ O_{H-1}(s, a) \end{bmatrix} = \begin{bmatrix} Q_{o,0}^\pi(s, a) \\ Q_{o,1}^\pi(s, a) \\ Q_{o,2}^\pi(s, a) \\ \vdots \\ Q_{o,H-1}^\pi(s, a) \end{bmatrix}$$

$$O_h^\pi(s, a) = \mathbb{E}[o(s_h, a_h, s_{h+1}) | s_0 = s, a_0 = a, \pi]$$

$$Q_{o,H}^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{H-1} \gamma^h o_h \mid s_0 = s, a_0 = a, \pi \right]$$

(a) Training

(b) Inference

EXPLANATIONS

- Probability of future events:** If the outcome function is an event (binary), the EFO is the probability of the event in a future time step.
- Expected reward components:** If the outcome function is a part of the reward, EFOs provide a reward decomposition, extending [2] with the temporal dimension.
- Contrastive explanations:** Comparing EFOs for different state-action pairs provides contrastive explanations (e.g., "Why is moving south preferred over action 2?").

CONCLUSIONS

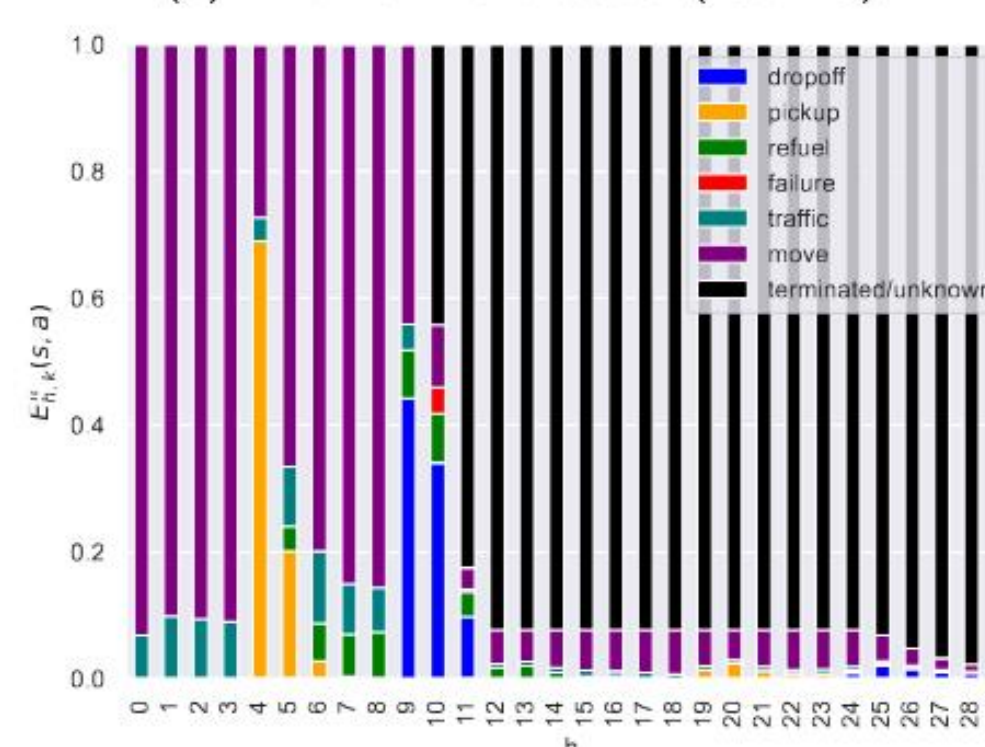
- Temporal Policy Decomposition [3] is an XRL method that explains individual actions by predicting EFOs
- Results in tabular problems show accurate predictions, which are necessary for having reliable explanations.
- In future work, we will experiment with continuous problems (deep neural networks).



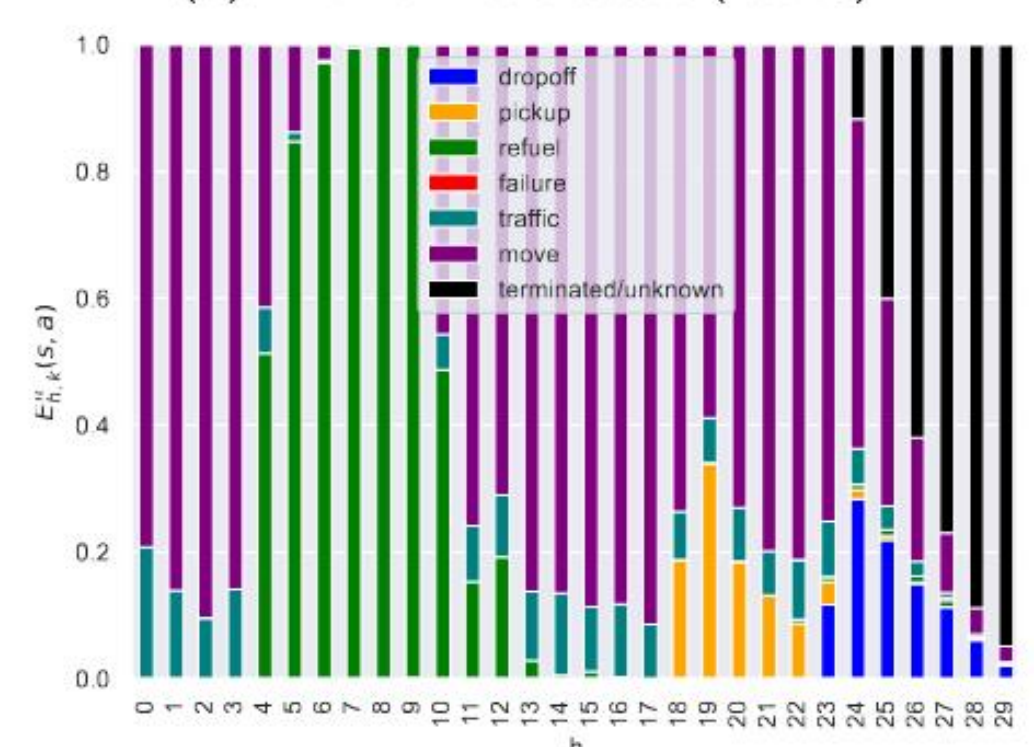
(a) Environment state (fuel 10)



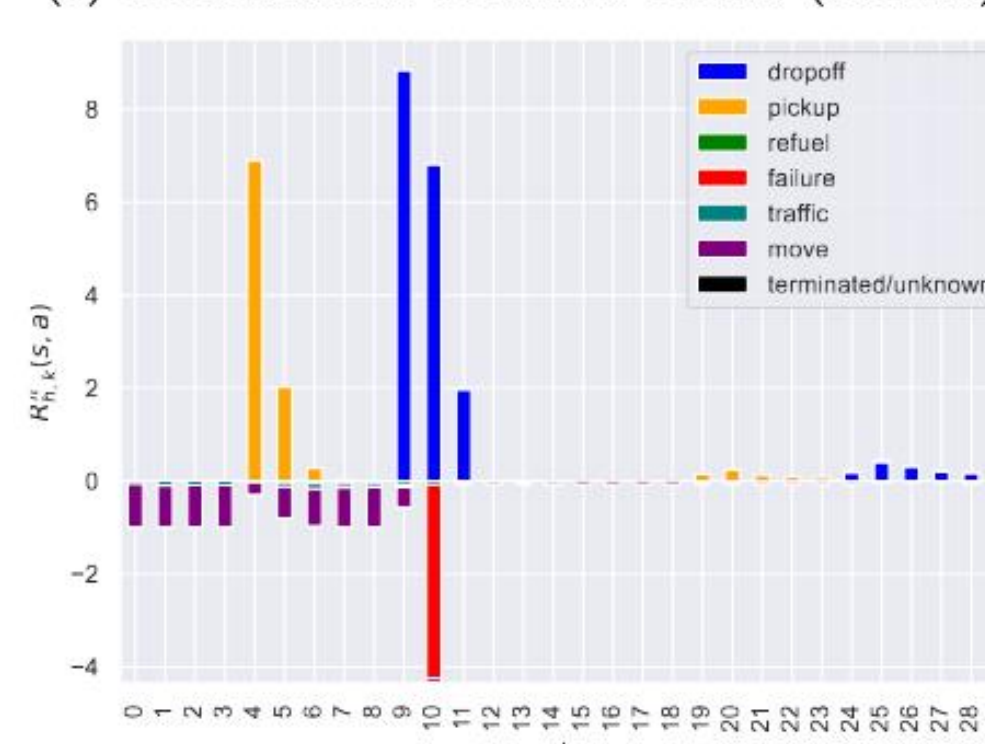
(b) Environment state (fuel 9)



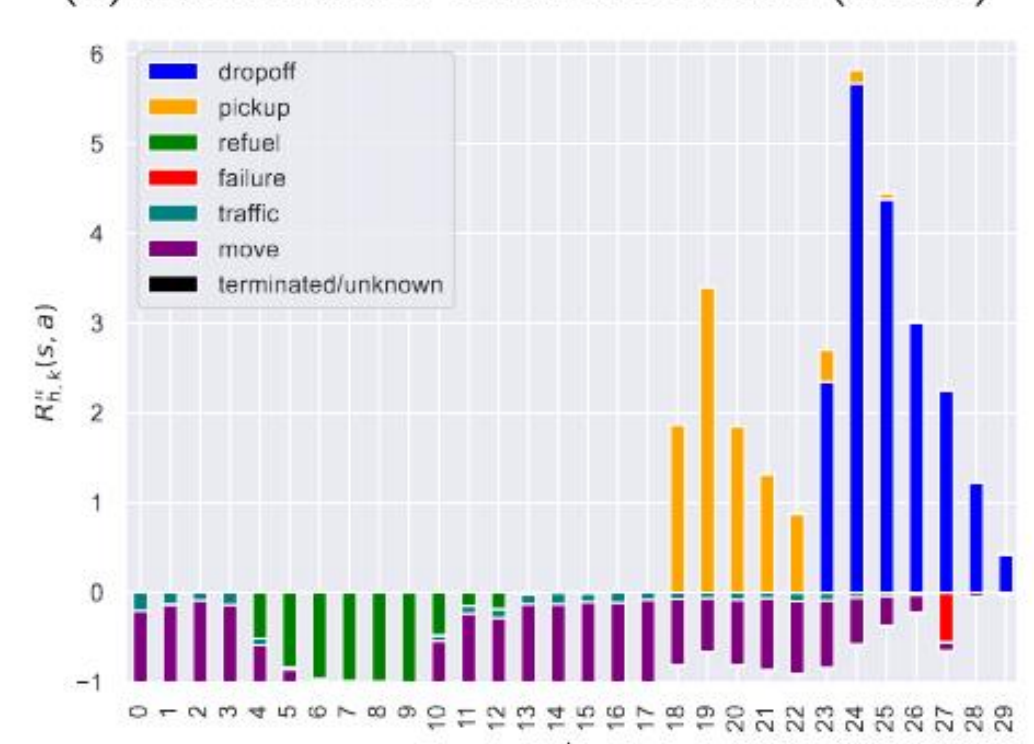
(c) Probabilities of future events (fuel 10)



(d) Probabilities of future events (fuel 9)



(e) Expected reward components (fuel 10)



(f) Expected reward components (fuel 9)