

**Background:** To develop and effectively deploy trustworthy citizen-centric AI, we need computational tools to reason about and determine if, and to what extent, AI agents, human users, or human-AI collectives are

1. able to deliver an outcome/task and take the role to do so (**prospective responsibility**)
2. to be seen accountable, blameworthy, and liable for potential failures (**retrospective responsibility**)



### Dimensions of Responsibility:

- **Normativity:** responsibility for a given desirable/undesirable outcome
- **Strategic aspects:** verifying the ability of agents to influence the occurrence of outcomes
- **Epistemic dynamics:** what information is available to agents
- **Temporality:** one may be responsible now but not a week ago

### Results:

- Computational techniques for **quantifying responsibility** and sharing blames
- Formalisations for **determining and distinguishing responsibility**, blame, accountability, and liability in human-AI systems

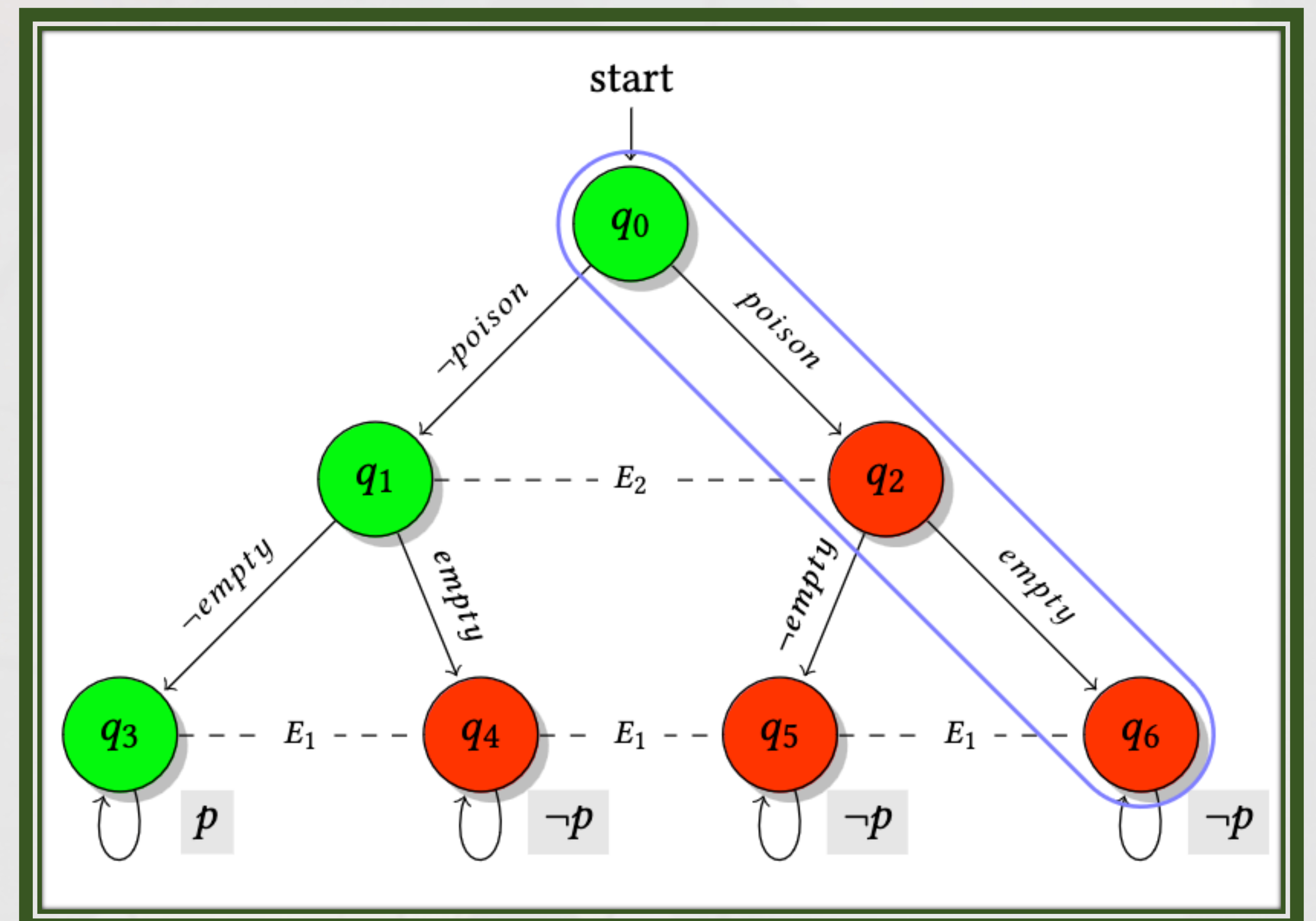
### Methods:

- Game Theory
- Temporal Logics and Game Structures
- Formal Verification

### Concurrent Epistemic Game Structures (CEGS):

A CEGS is a tuple  $M = \langle \Sigma, Q, Act, \sim_1, \dots, \sim_n, d, o \rangle$  where:

- $\Sigma = \{a_1, \dots, a_n\}$  is the set of agents;
- $Q$  is the set of states;
- $Act$  is the set of actions;
- $\sim_a \subseteq Q \times Q$  is an epistemic (equivalence) indistinguishability relation,  $q \sim_a q'$  indicates that states  $q$  and  $q'$  are indistinguishable to  $a$ ;
- $d : \Sigma \times Q \rightarrow P(Act)$  specifies the sets of actions available to agents at each state, same actions are available in indistinguishable states;
- $o$  is a deterministic transition function.



**A Passenger and Two Enemies (McLaughlin, 1925)**

In  $q_0$ ,  $E_1$  may poison the water. In  $q_1$  and  $q_2$ ,  $E_2$  may empty the canteen. As a result,  $P$  is alive in  $q_3$  (represented by proposition  $p$ ) and dead in  $q_4$ ,  $q_5$ , and  $q_6$  (represented by  $\neg p$ ). The path outlined in blue denotes the history.

**With more autonomy comes more, and different forms of responsibility**



**Limitations:** Our formal responsibility models requires an **exhaustive specification** of the context.

**Future work:** Hybrid models that use data-driven methods **to learn** about the context and formal methods **to reason** about responsibilities.

- Stein and Yazdanpanah. Citizen-centric multiagent systems (2023)
- Williams, Yazdanpanah and Stein. Safe audio AI services in smart buildings (2022)
- Yazdanpanah, Gerding, Stein, et al. Different forms of responsibility in multiagent systems: Sociotechnical characteristics and requirements (2021)
- Yazdanpanah, Gerding, Stein, et al. Reasoning about responsibility in autonomous systems: challenges and opportunities (2022)

