

# **An Information-Flow Perspective on Explainability Requirements:** Specification and Verification

Bernd Finkbeiner<sup>1</sup>, Hadar Frenkel<sup>2</sup>, and Julian Siber<sup>1</sup>

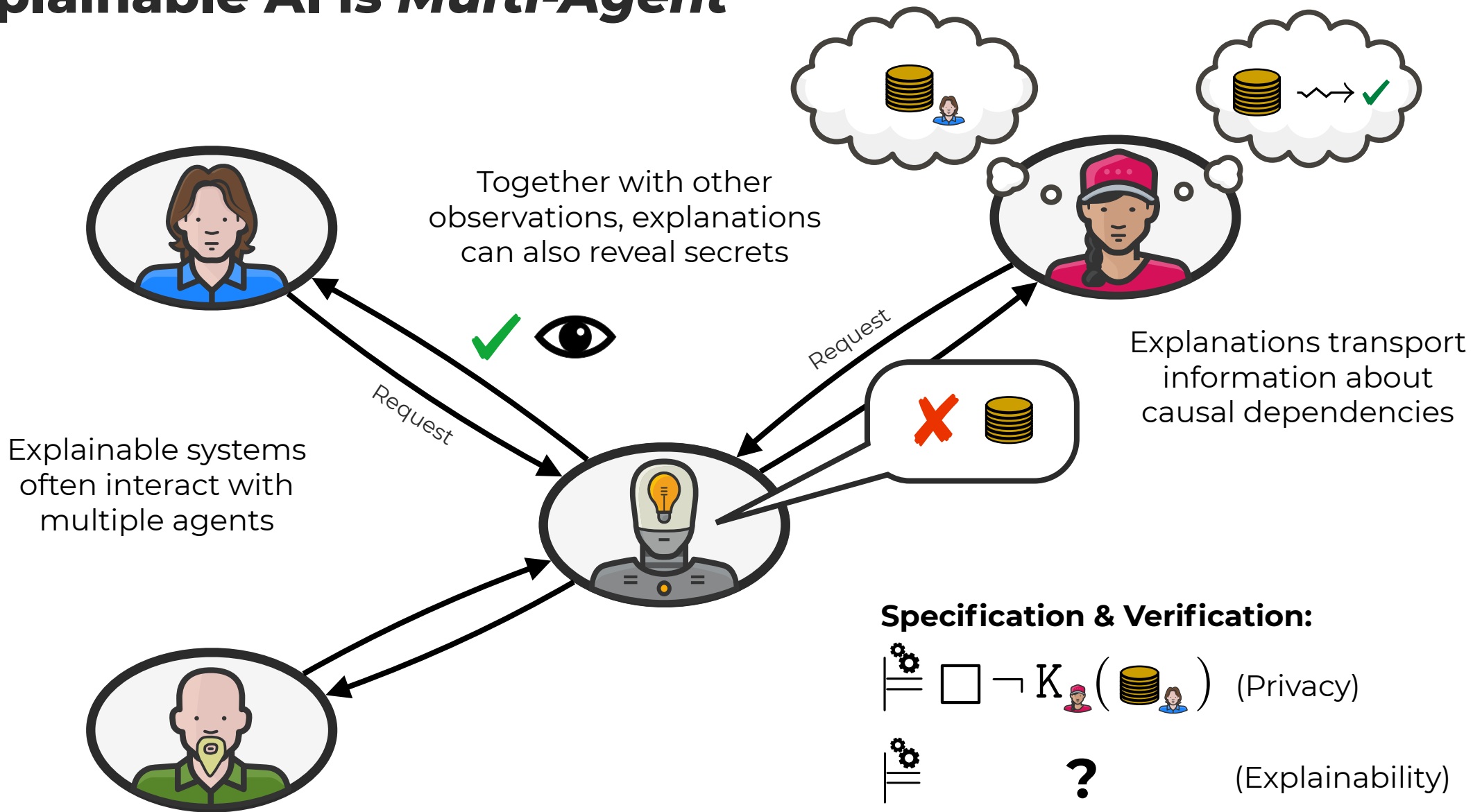
<sup>1</sup> *CISPA Helmholtz Center for Information Security*

<sup>2</sup> *Bar-Ilan University*





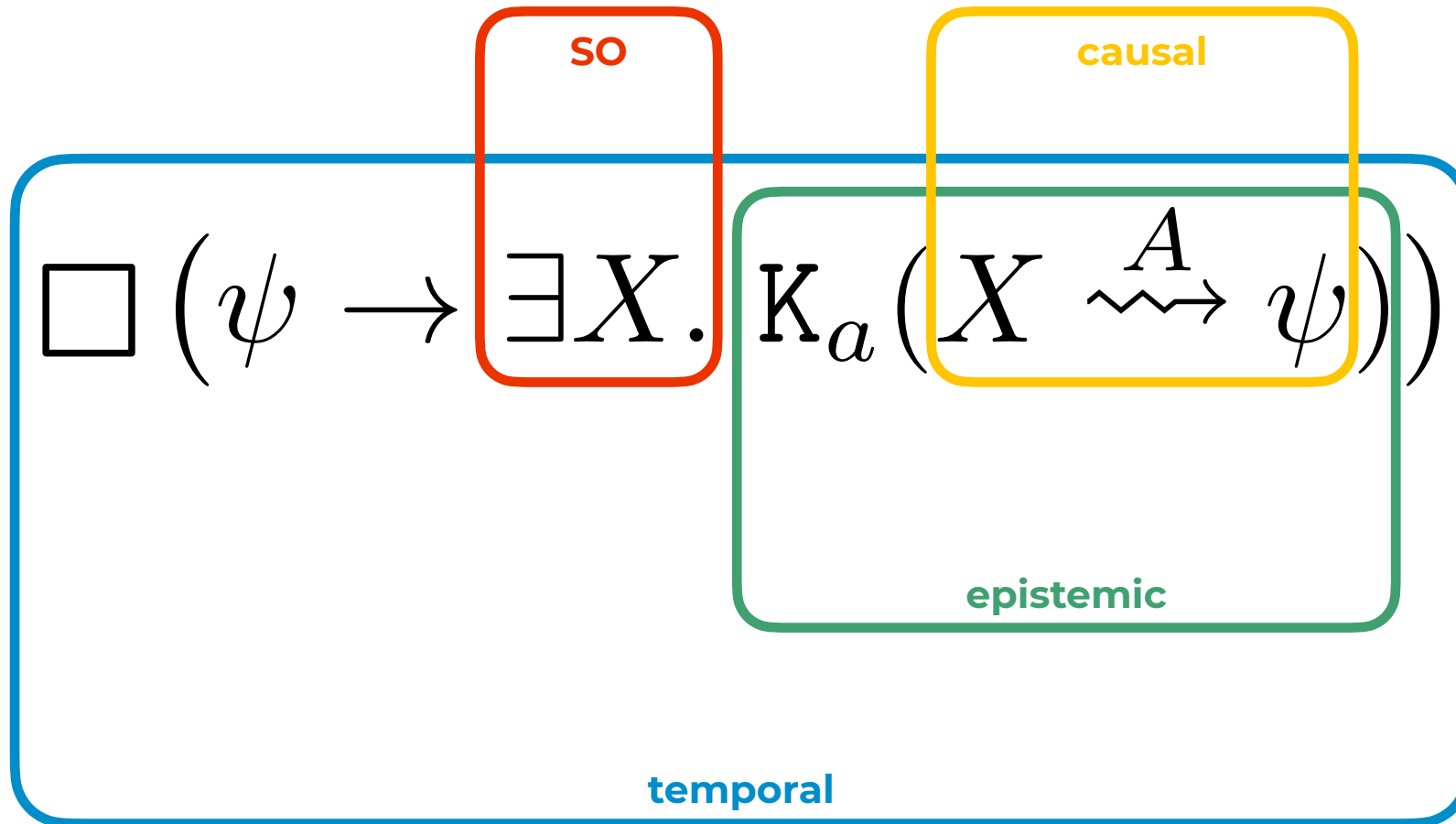
# Explainable AI is Multi-Agent





# Key Takeaway

Specifying explainability requires a combination of causal, epistemic and temporal reasoning.

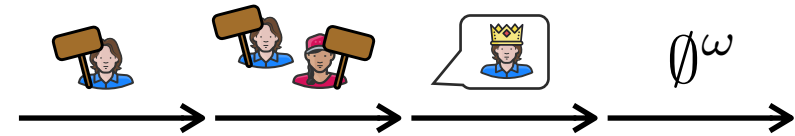
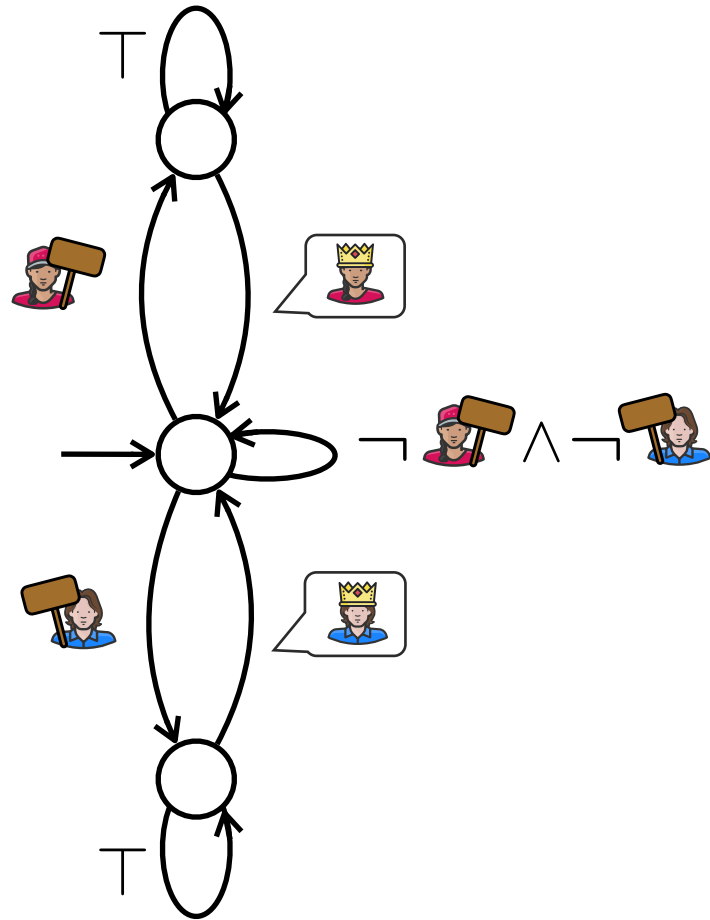




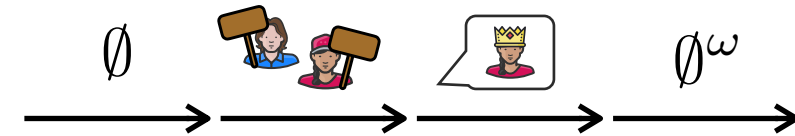
# Running Example

$$\square (\psi \rightarrow \exists X. K_a(X \xrightarrow{A} \psi))$$

The state machine models a blind Dutch auction between  and . The first bid wins.



Models generate traces.





# Key Takeaway

Specifying explainability requires a combination of causal, epistemic and temporal reasoning.

$$\square (\psi \rightarrow \exists X. K_a (X \overset{A}{\rightsquigarrow} \psi))$$

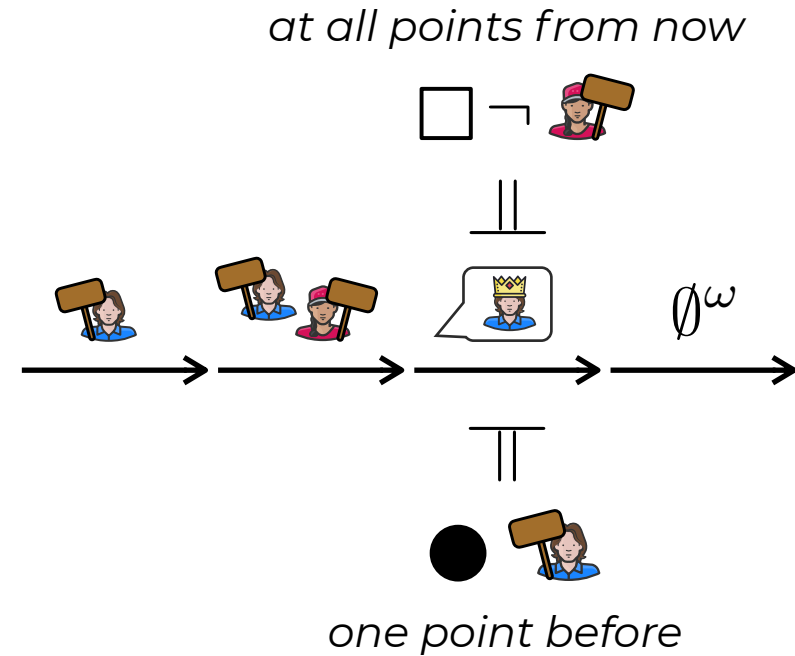
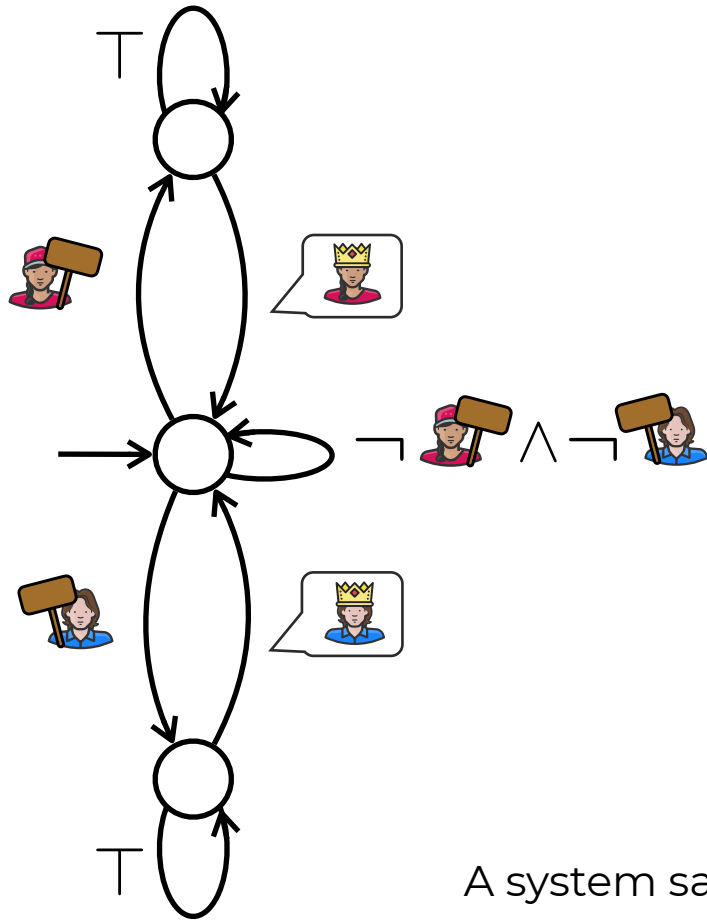
temporal



# Temporal Operators [1]

$$\square (\psi \rightarrow \exists X. K_a(X \xrightarrow{A} \psi))$$

Formulas are evaluated with respect to a trace and time point (*anchor point*).



A system satisfies a formula if all initial trace satisfy the formula.

[1] Amir Pnueli. *The Temporal Logic of Programs*. FOCS 1977.



# Key Takeaway

Specifying explainability requires a combination of causal, epistemic and temporal reasoning.

$$\square (\psi \rightarrow \exists X. K_a (X \overset{A}{\rightsquigarrow} \psi))$$

temporal

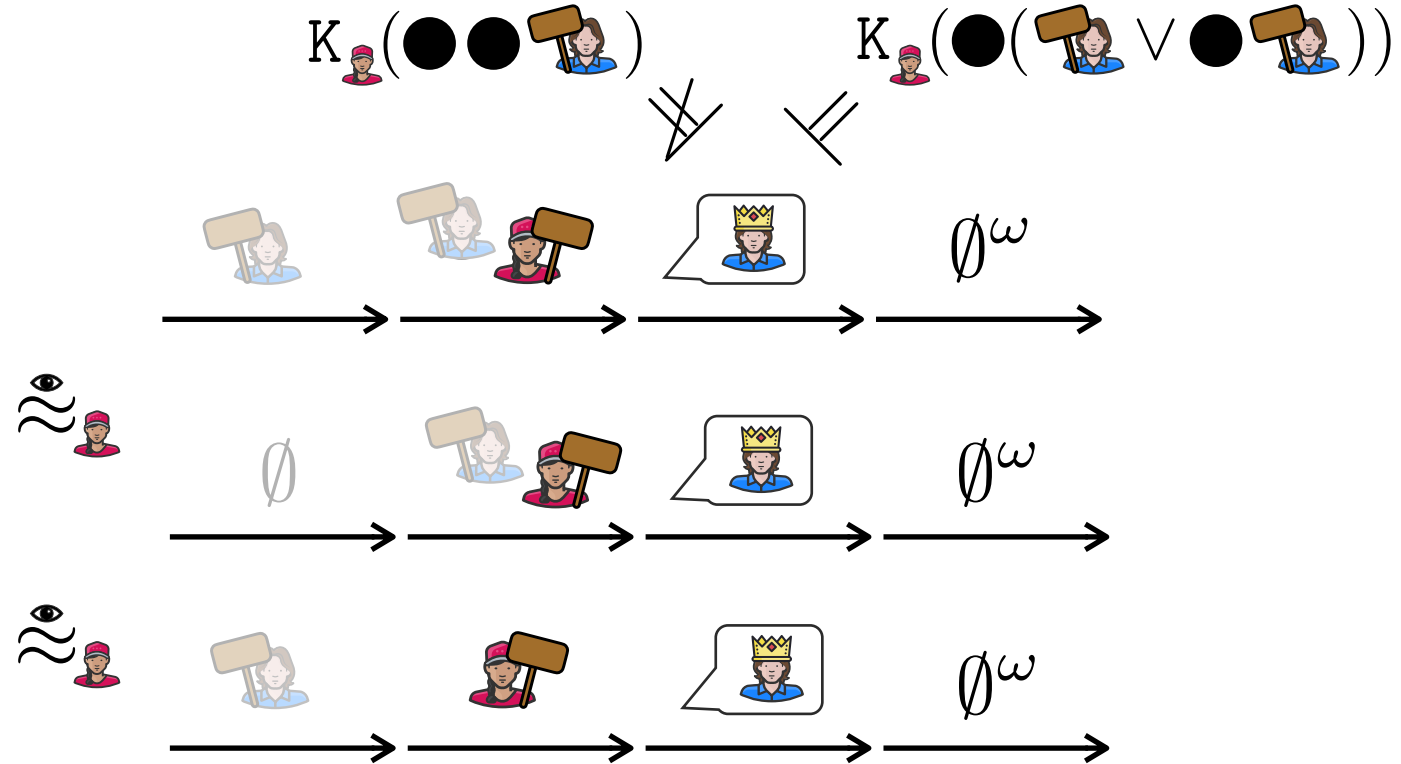
epistemic



# Epistemic Operators [2]

$$\Box (\psi \rightarrow \exists X. K_a (X \overset{A}{\rightsquigarrow} \psi))$$

Knowledge operators require their subformula to hold on all observation equivalent ( $\overset{a}{\sim}$ ) traces.

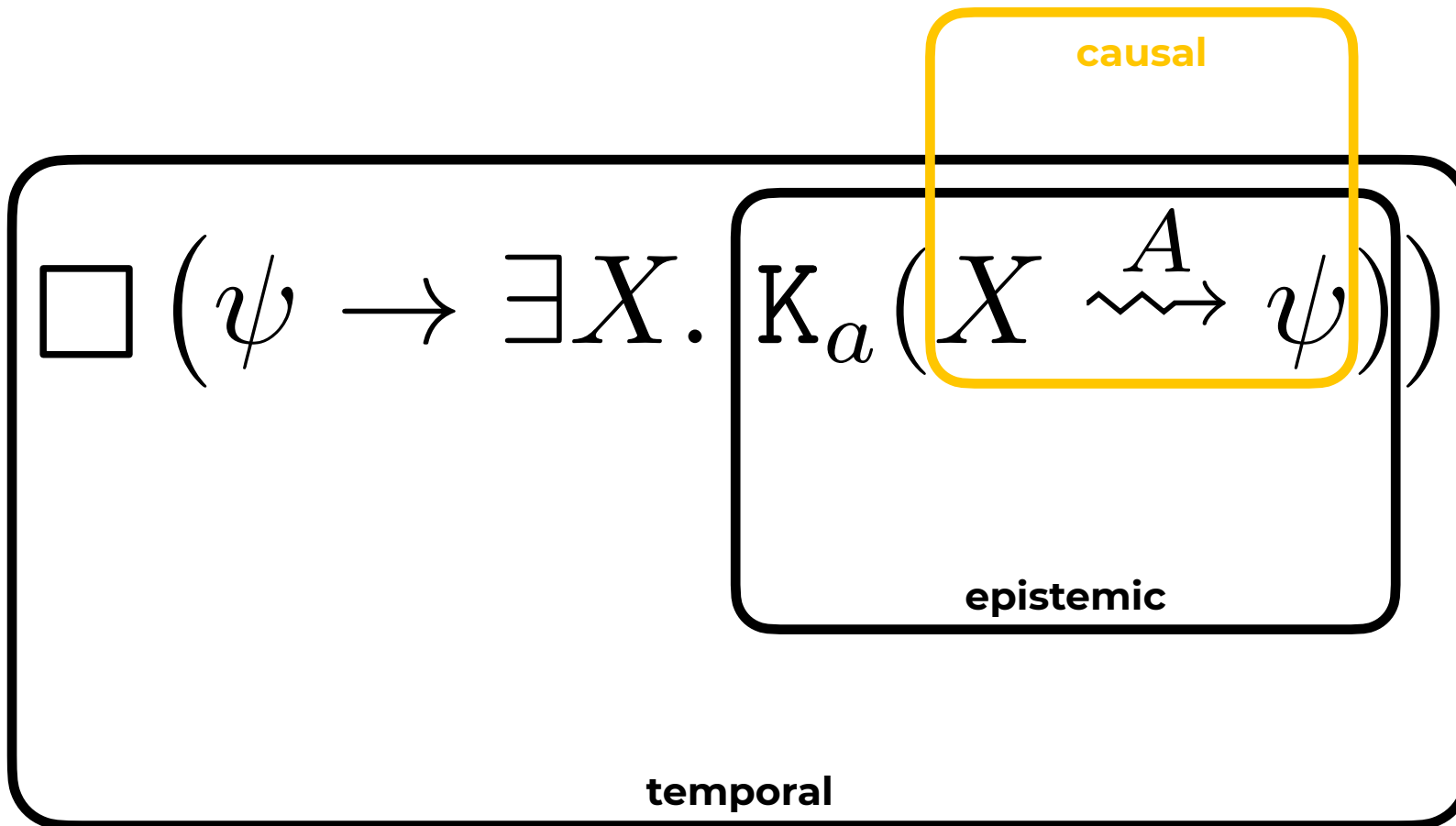


*Synchronous perfect-recall semantics:* Same observable propositions on full prefixes.



# Key Takeaway

Specifying explainability requires a combination of causal, epistemic and temporal reasoning.





# Causal Dependencies

$$\square (\psi \rightarrow \exists X. \kappa_a (X \overset{A}{\rightsquigarrow} \psi))$$

$X \overset{A}{\rightsquigarrow} \varphi$  : Property  $X$  is a cause for property  $\varphi$  (when fixing all actions not in  $A$  ).

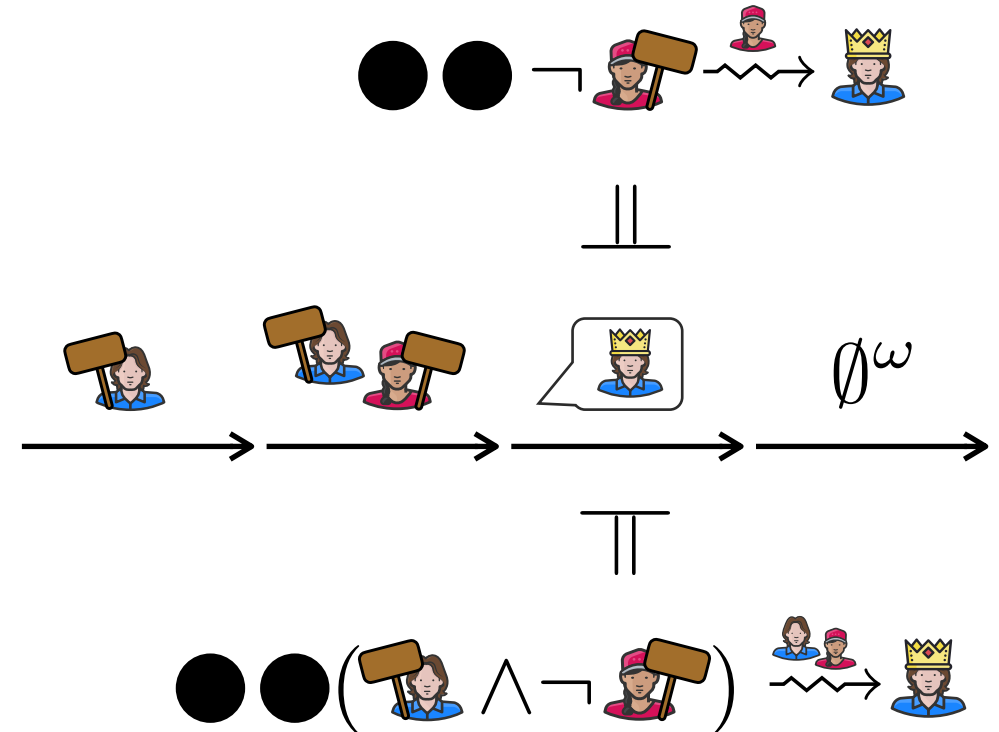
**Temporal Causality [3]**

**SAT:** Cause and effect are satisfied by the trace.

**CF:** Most similar traces that do not satisfy the cause also do not satisfy the effect.

**MIN:** Cause is (semantically) minimal.

*Heavily inspired by actual causality [4].*



! There exists at most one unique causal property (which may symbolically describe multiple events).

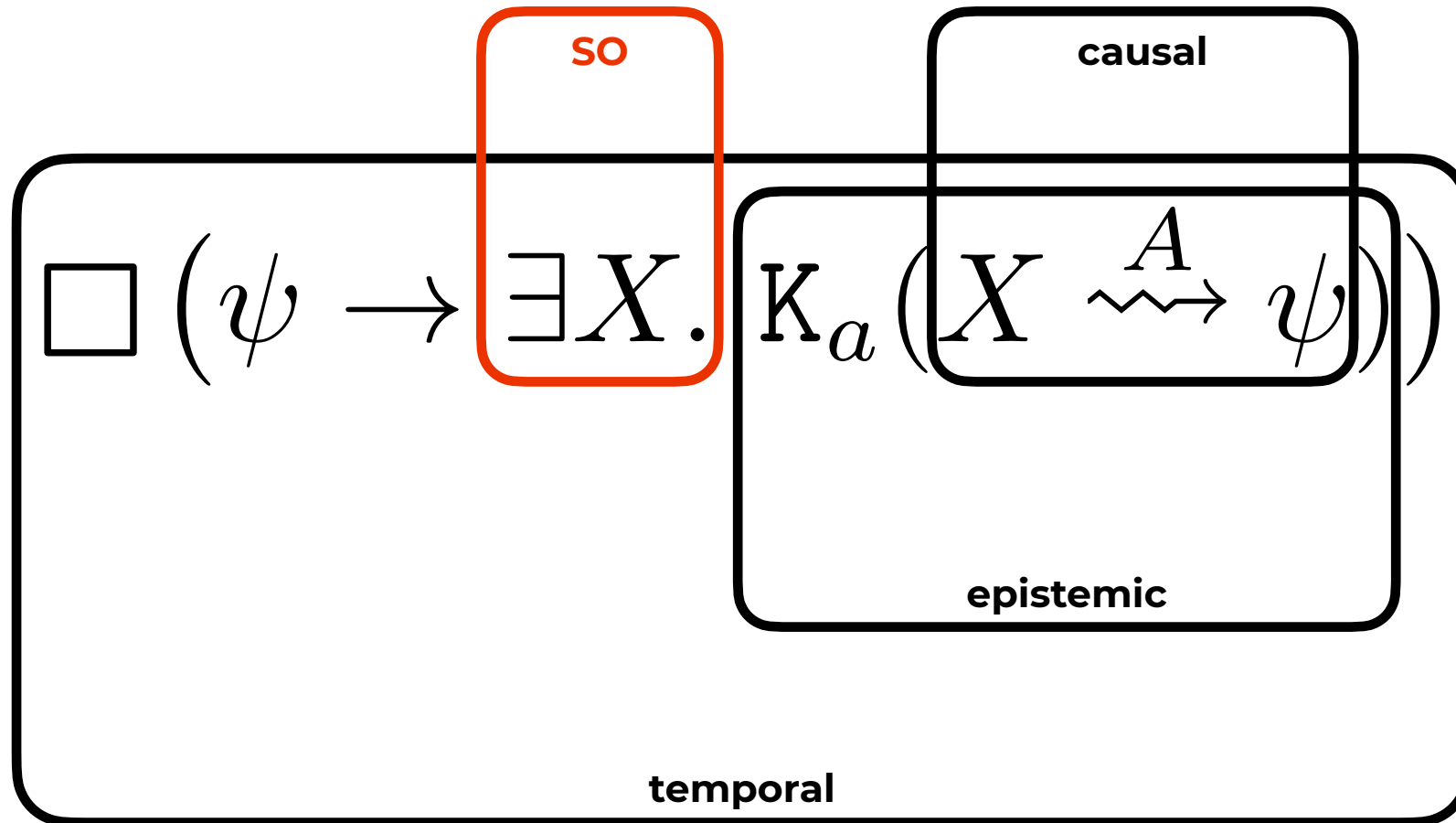
[3] Finkbeiner, Frenkel, Metzger and Siber. Synthesis of Temporal Causality. CAV 2024.

[4] Halpern and Pearl. Causes and Explanations: A Structural-Model Approach — Part 1: Causes. British Journal for the Philosophy of Science 56 (4). 2005.



# Key Takeaway

Specifying explainability requires a combination of causal, epistemic and temporal reasoning.

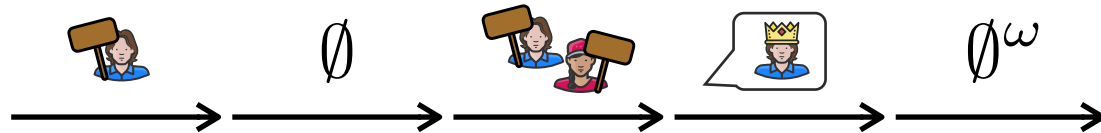
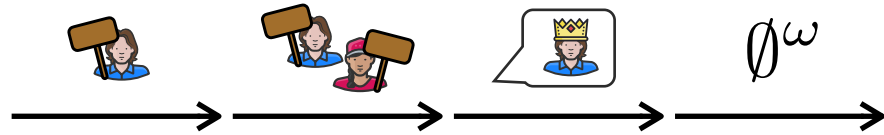




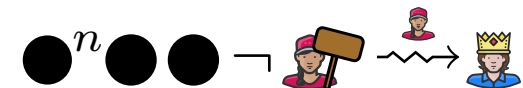
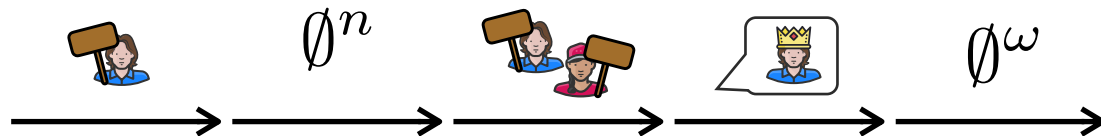
# Why Second-Order Quantification?

$$\square (\psi \rightarrow \exists X. K_a(X \overset{A}{\rightsquigarrow} \psi))$$

A single explanandum can have infinitely many different explanantia:



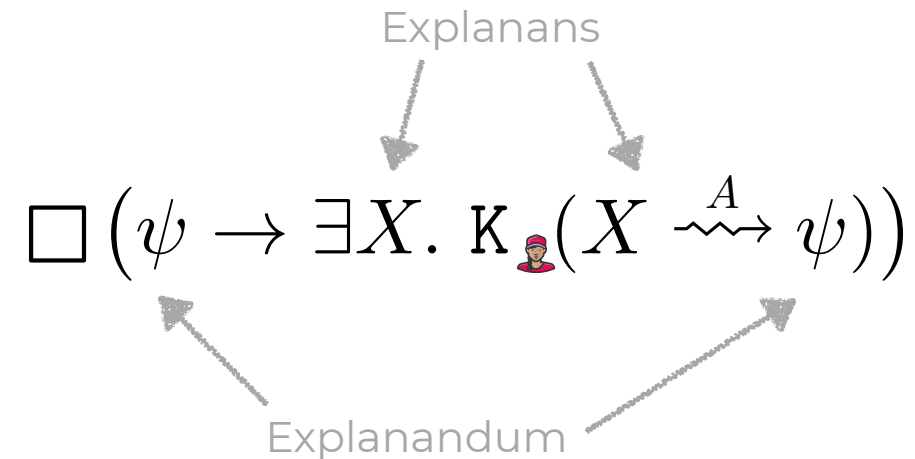
$\forall n :$





# Putting Everything Together

Explainability: Whenever a certain effect (*explanandum*) happens, an agent has knowledge about a cause (*explanans*) for this outcome.



Internal Causal Explainability (ICE):  $A = \{\text{agent}\}$

External Causal Explainability (ECE):  $A = \text{Agents} \setminus \{\text{agent}\}$

Full Causal Explainability (FCE):  $A = \text{Agents}$



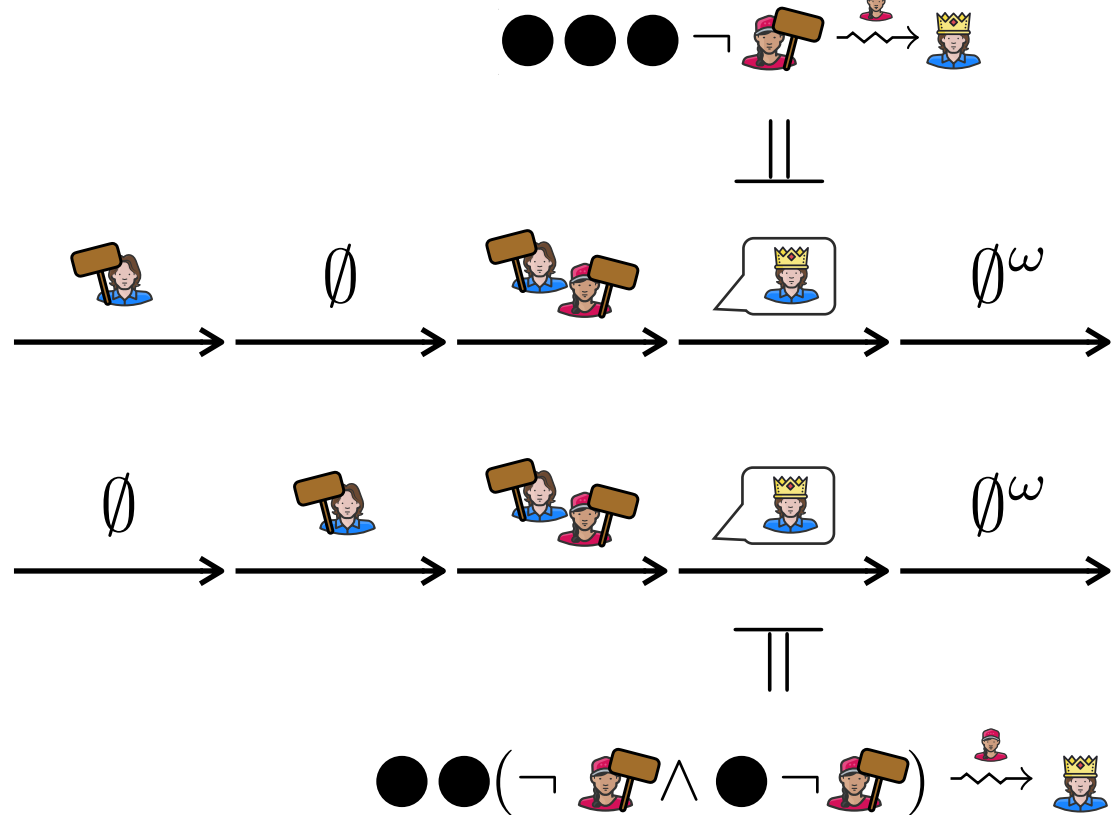
# Information-Flow through Explanations

Semantics: Traces with different causes for the same explanandum must be distinguishable.

$$\Box ( \text{king} \rightarrow \exists X. K_{\text{person}} ( X \rightsquigarrow \text{king} ) )$$

$\not\equiv$

$\mathcal{T}_{\text{blind}}$





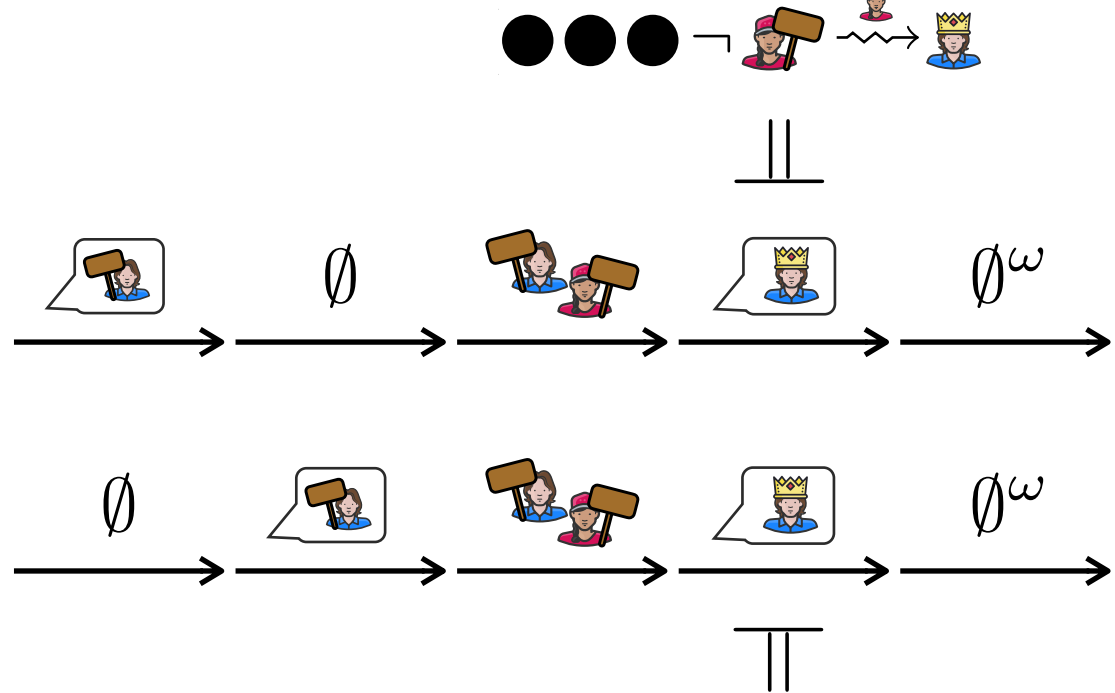
# Information-Flow through Explanations

Semantics: Traces with different causes for the same explanandum must be distinguishable.

$$\square ( \text{king} \rightarrow \exists X. \text{K}_{\text{blind}} ( X \rightsquigarrow \text{king} ) )$$

$\not\equiv$   
 $\mathcal{T}_{\text{blind}}$

$\equiv$   
 $\mathcal{T}$

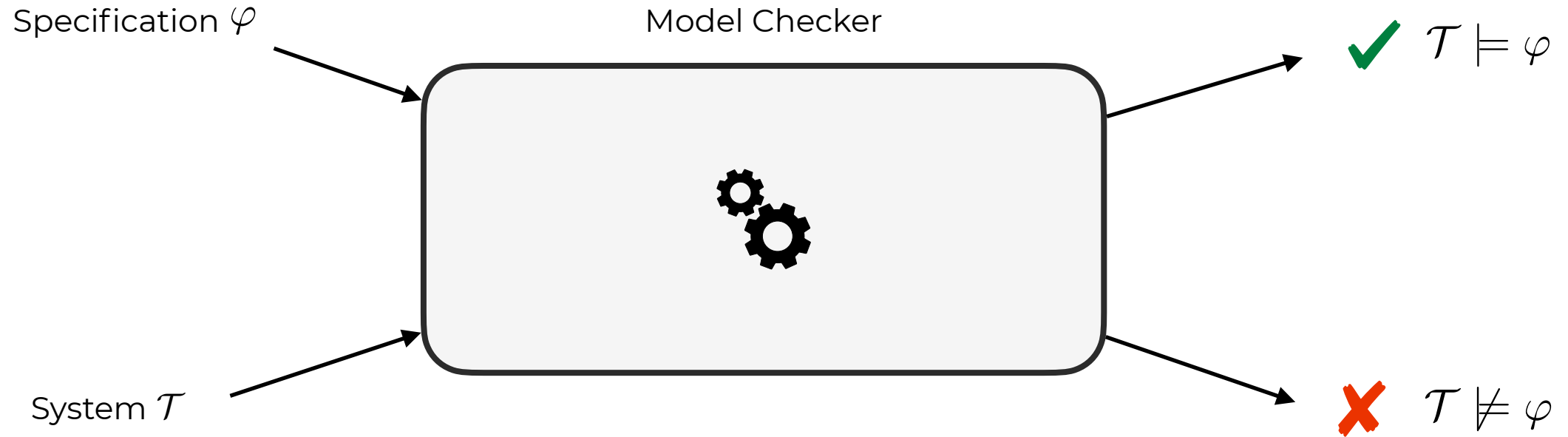


The observation  transports information about arbitrarily large explanantia.

$$\bullet \bullet (\neg \text{sign} \wedge \bullet \neg \text{sign}) \rightsquigarrow \text{king}$$

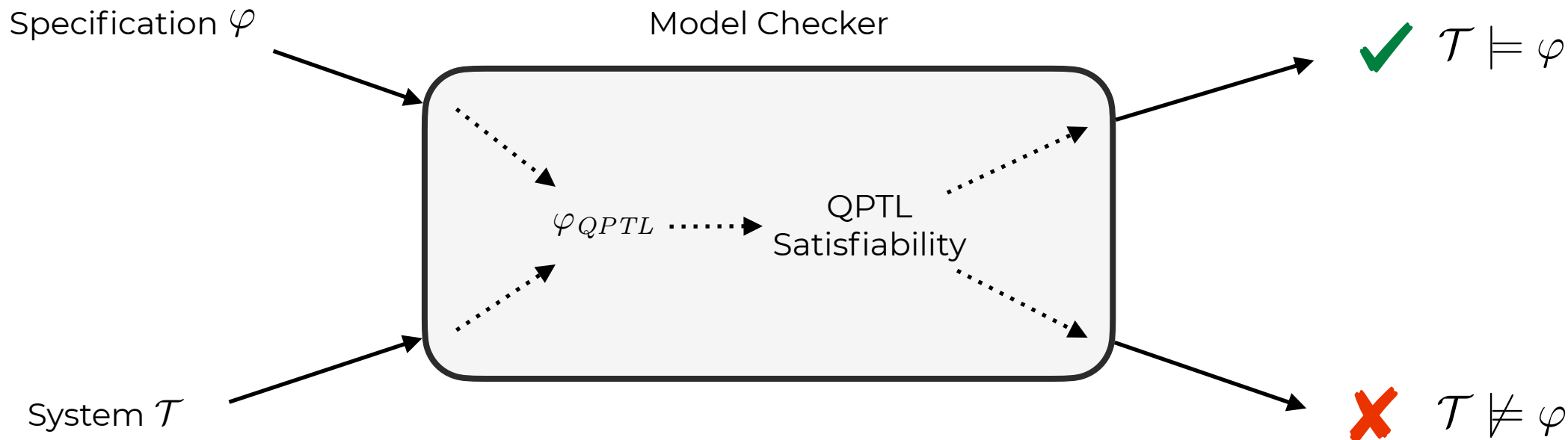


# Automated Verification





# High-Level Idea



QPTL: LTL + quantification over propositional sequences:  $\xi ::= \forall q \in (2^{\{q\}})^\omega. \xi$

Simulate quantification over paths [5]:

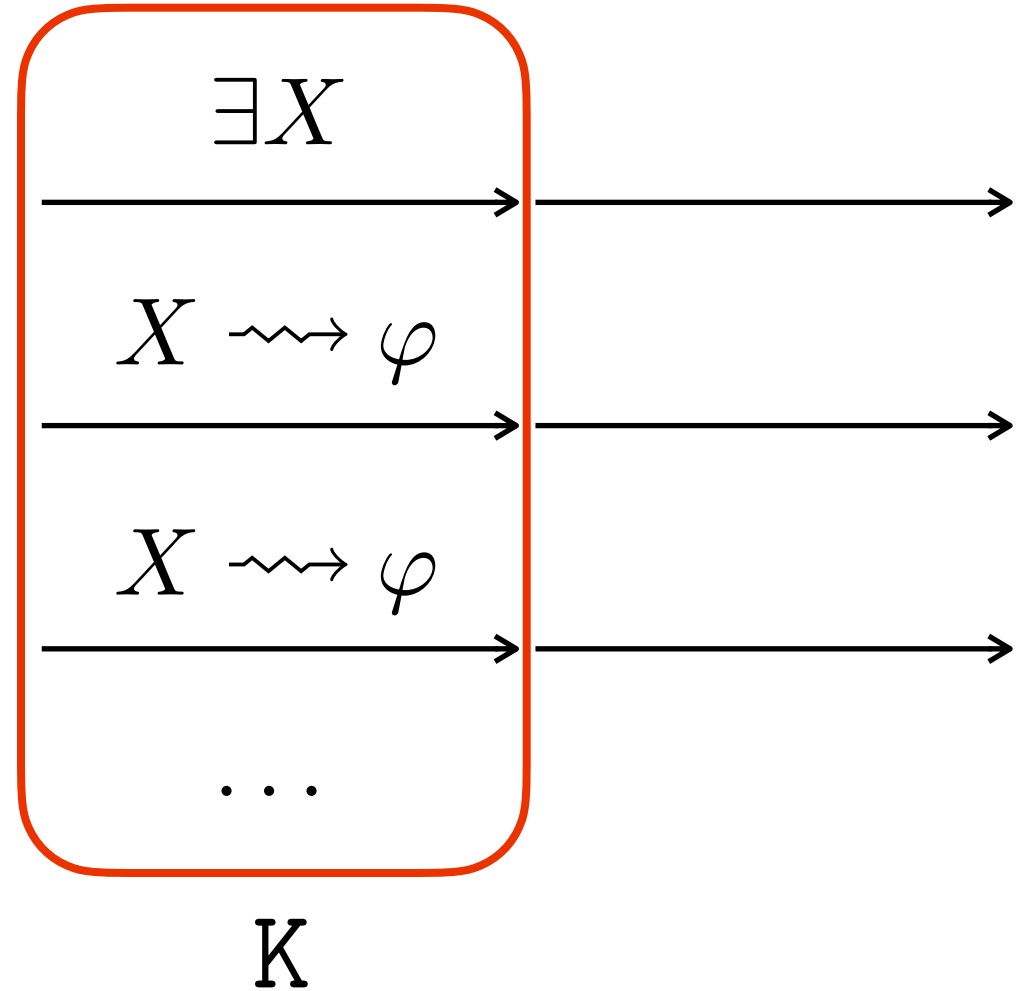
“behave like  $\mathcal{T}$ ”

$$\forall a \in \text{Actions}, s \in \text{States}. \blacklozenge \left( (\neg \bullet \mathcal{T}) \wedge \square \widetilde{\mathcal{T}}(a, s) \right) \rightarrow \varphi_{LTL}(a, s)$$



# More Details

How to encode  $\exists X. K_a(X \rightsquigarrow \varphi)$  ?





# More Details

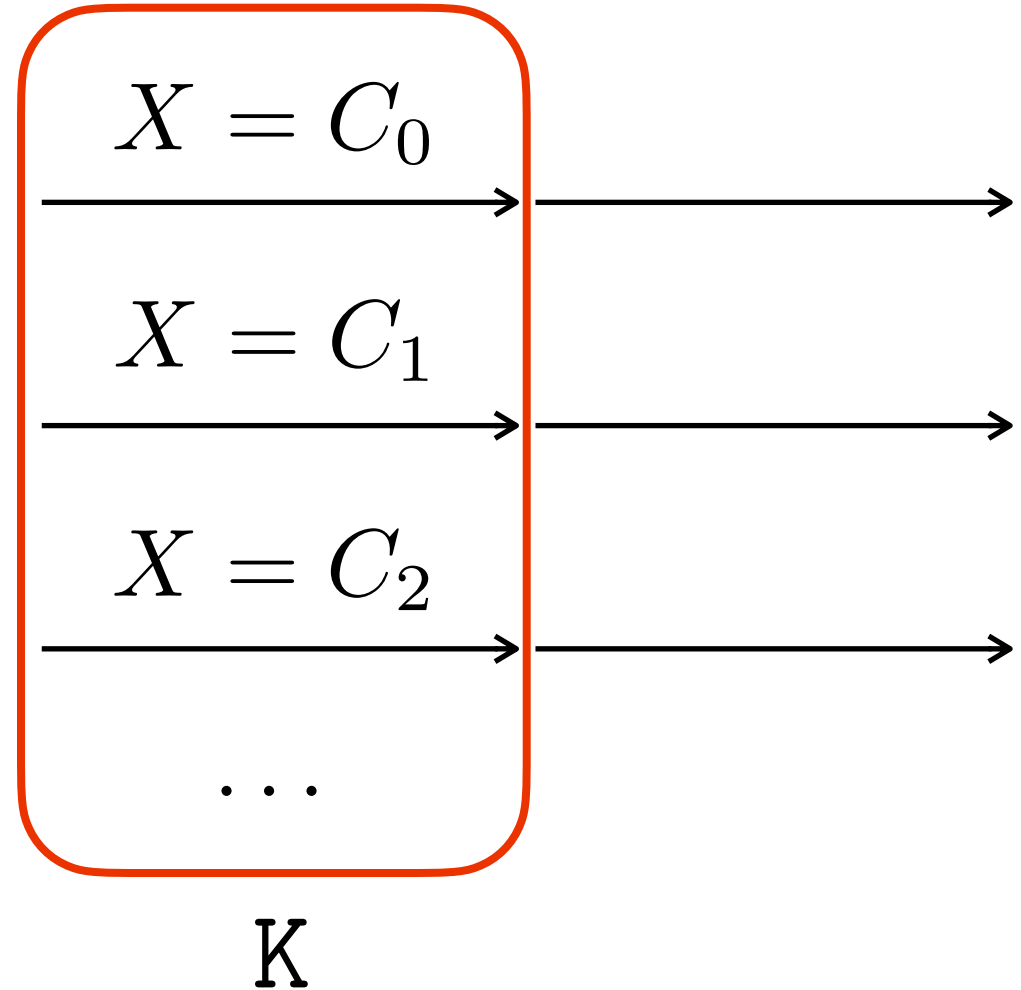
How to encode  $\exists X. \mathbb{K}_a(X \rightsquigarrow \varphi)$  ?

There is a unique solution  $C_n$  for every  $\rightsquigarrow \varphi$  and anchor point.

The formula requires that all the scoped solutions are equivalent:

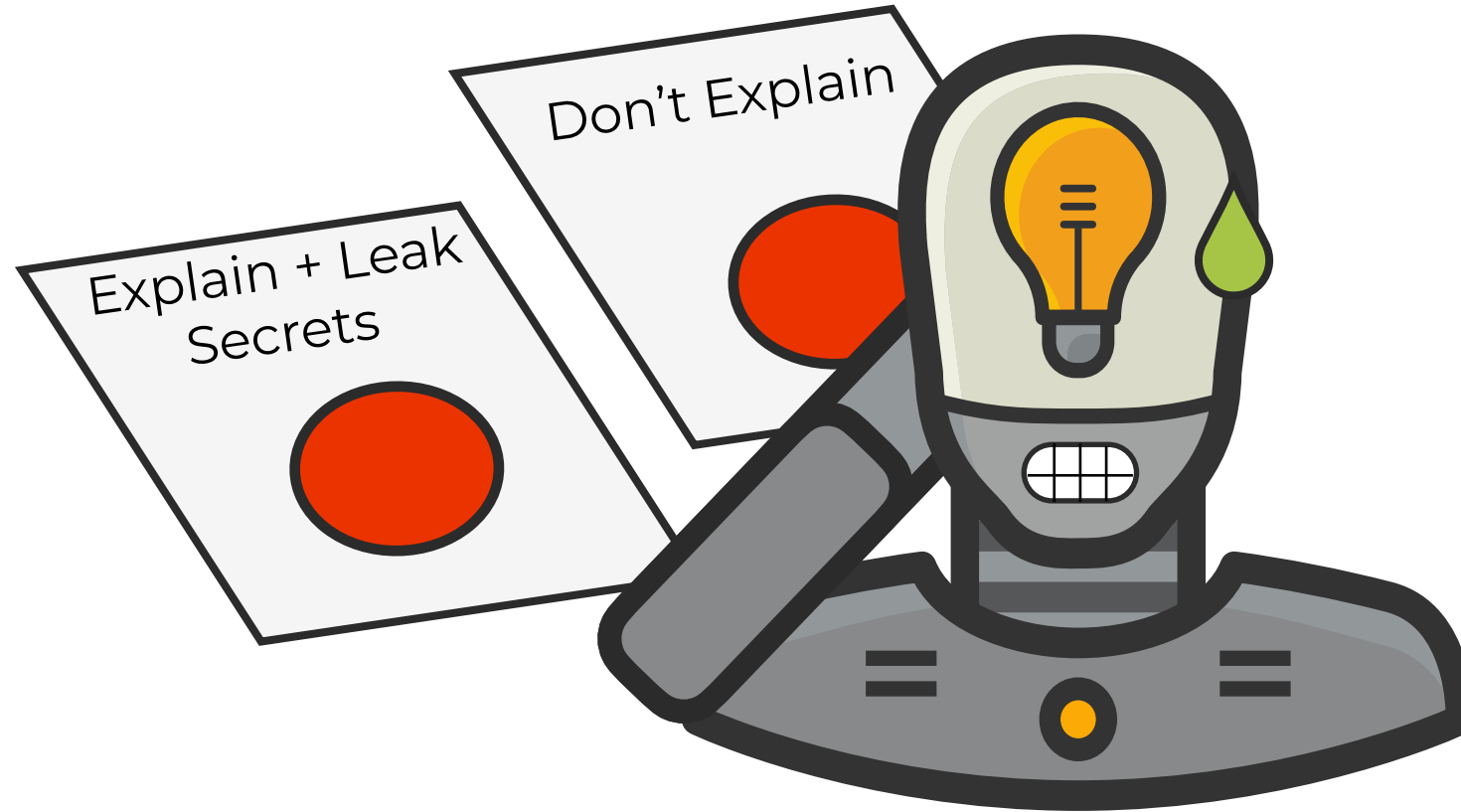
$$C_0 = C_1 = C_2 = \dots$$

➡ Use simulated quantification over paths to search for a trace that breaks the equality.



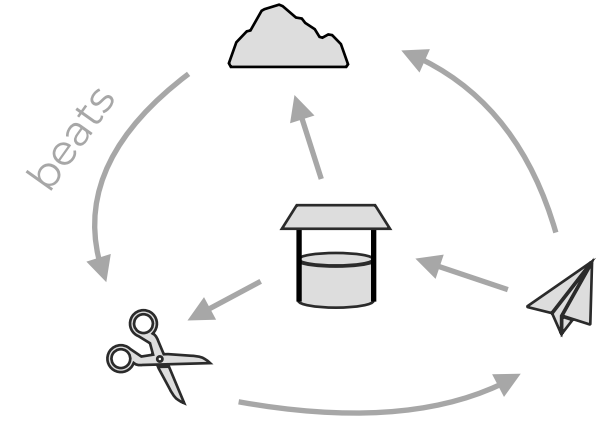
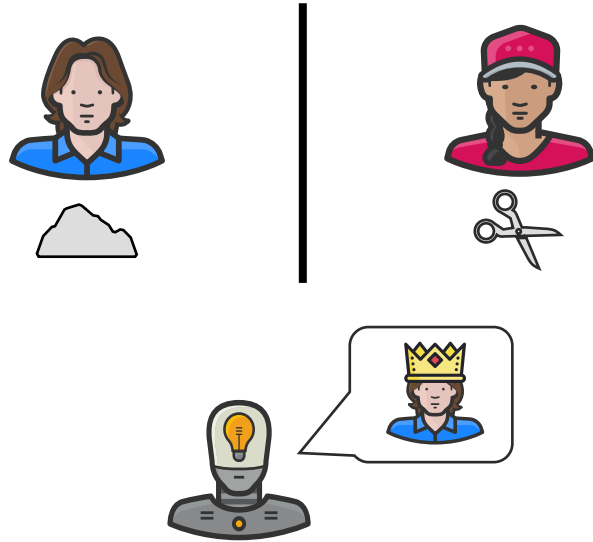


# Experiments





# Experiments: Blind RPS



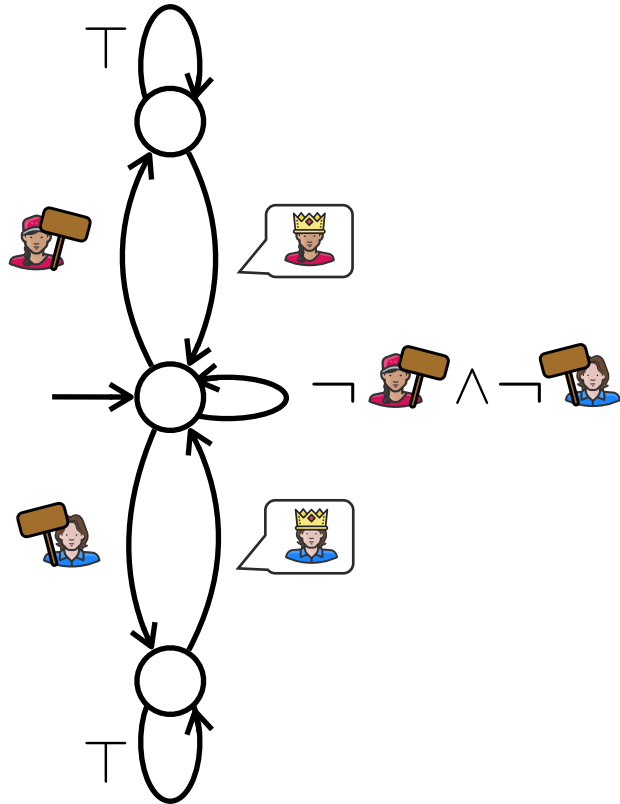
	ICE	ECE	FCE	Privacy
STANDARD	✓/0.49	✓/0.53	✓/0.55	✗/0.24

$$\square (\neg draw \rightarrow \neg K_{\text{woman}}(\text{rock}))$$



# Experiments: Dutch Auction

$$\square (\neg K_{\text{red}} (\text{red\_auctioneer})) \xrightarrow{\text{red\_auctioneer}}$$

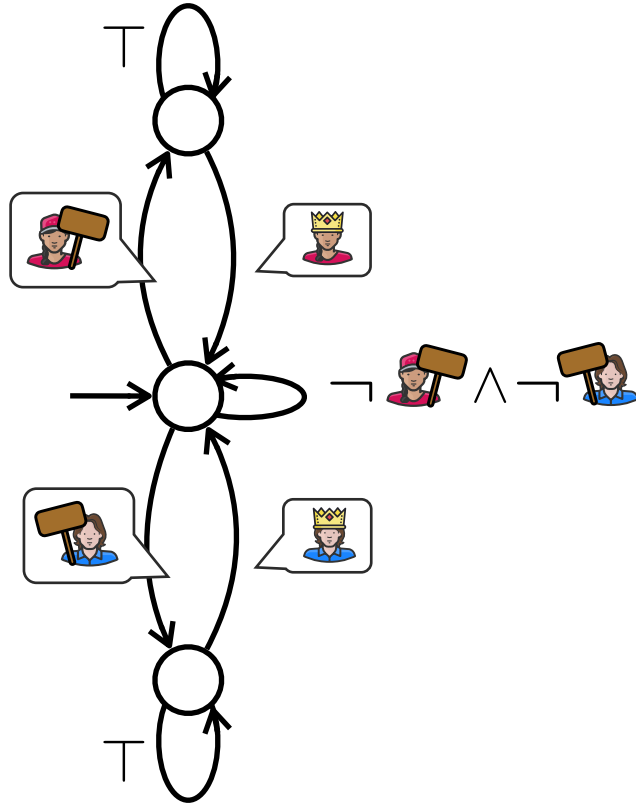


	$ B $	ICE	ECE	FCE	Privacy
BLIND	2	✗/0.85	✗/0.88	✗/0.96	✓/0.25
	3	✗/1.57	✗/1.77	✗/1.73	✓/0.27
	4	✗/3.42	✗/3.68	✗/3.61	✓/0.27
	5	✗/7.85	✗/9.45	✗/10.3	✓/0.22



# Experiments: Dutch Auction

$$\square (\neg K_{\text{red}} (\text{bidder})) \quad \curvearrowright$$

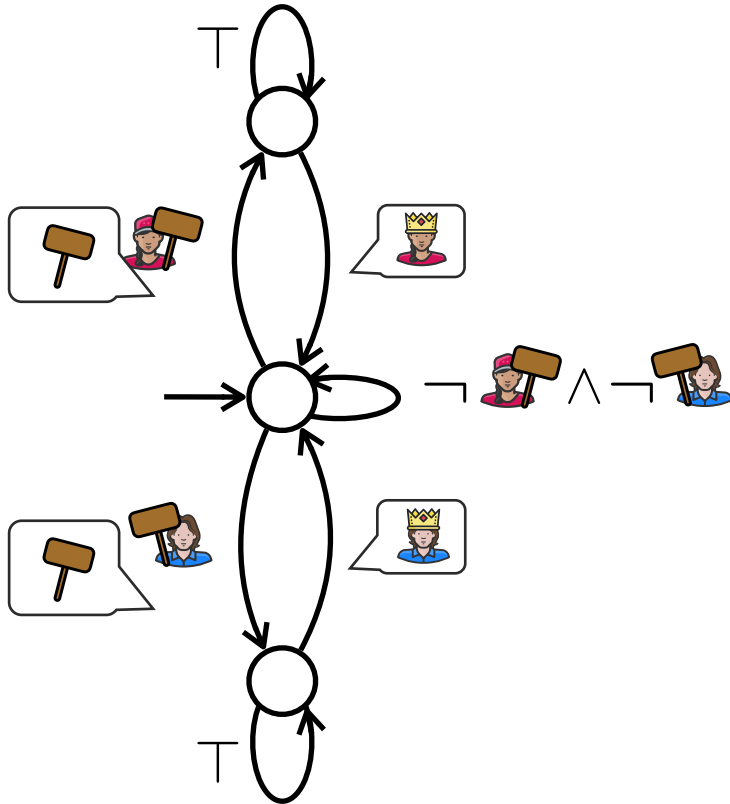


	$ B $	ICE	ECE	FCE	Privacy
BLIND	2	✗/0.85	✗/0.88	✗/0.96	✓/0.25
	3	✗/1.57	✗/1.77	✗/1.73	✓/0.27
	4	✗/3.42	✗/3.68	✗/3.61	✓/0.27
	5	✗/7.85	✗/9.45	✗/10.3	✓/0.22
PUBLIC	2	✓/0.78	✓/0.87	✓/0.81	✗/0.30
	3	✓/1.47	✓/1.80	✓/1.63	✗/0.30
	4	✓/3.19	✓/3.73	✓/3.48	✗/0.50
	5	✓/7.24	✓/9.52	✓/10.3	✗/1.21



# Experiments: Dutch Auction

$$\square (\neg K_{\text{red}} (\text{bidder}))$$



	$ B $	ICE	ECE	FCE	Privacy
BLIND	2	✗/0.85	✗/0.88	✗/0.96	✓/0.25
	3	✗/1.57	✗/1.77	✗/1.73	✓/0.27
	4	✗/3.42	✗/3.68	✗/3.61	✓/0.27
	5	✗/7.85	✗/9.45	✗/10.3	✓/0.22
PUBLIC	2	✓/0.78	✓/0.87	✓/0.81	✗/0.30
	3	✓/1.47	✓/1.80	✓/1.63	✗/0.30
	4	✓/3.19	✓/3.73	✓/3.48	✗/0.50
	5	✓/7.24	✓/9.52	✓/10.3	✗/1.21
EXPLAIN	2	✓/0.86	✗/1.02	✗/1.02	✗/0.24
	3	✓/1.39	✗/1.89	✗/1.82	✓/0.25
	4	✓/3.40	✗/3.82	✗/4.03	✓/0.29
	5	✓/8.03	✗/9.92	✗/10.4	✓/0.29



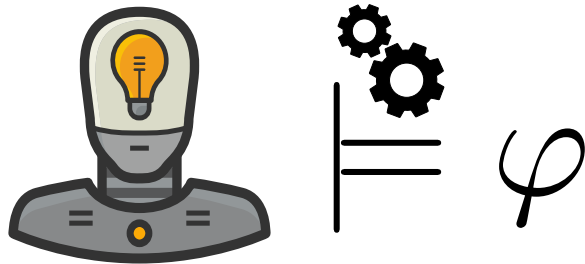
# Conclusion



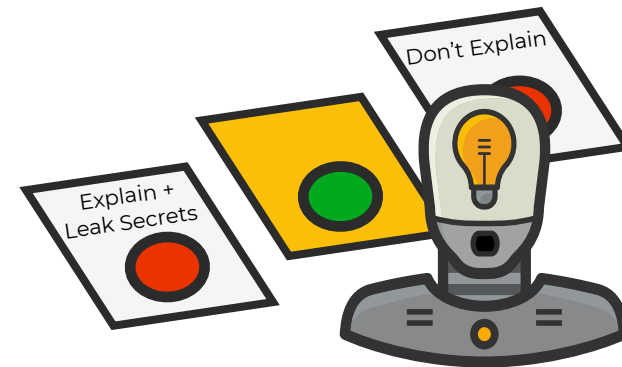
Explainability and privacy need to be balanced

$$\square ( \text{X} \rightarrow \exists X. K_{\text{person}} ( X \rightsquigarrow \text{X} ) )$$

Explainability is an information-flow property



The information flow of explainability can be verified



Good explanations do not need to sacrifice privacy