# Counterfactuals Modulo Temporal Logics

Bernd Finkbeiner and *Julian Siber*

CISPA Helmholtz Center for Information Security

LPAR-24 | 9 June 2023

*"If $\varphi$ had been true, then $\psi$ would have been true, too."*

$$\varphi \mathbin{\Box\!\!\!\longrightarrow} \psi$$

*"If the car had **moved straight** at the first time point, then it would have **reached its goal eventually**."*

$$(straight) \,\Box\!\!\longrightarrow (\mathbf{F}\,goal)$$

# Counterfactuals = Variably Strict Conditionals

$$\models (straight) \,\square\!\!\rightarrow (\mathbf{F} \; goal)$$

because the closest counterfactual world is:

# Counterfactuals = Variably Strict Conditionals

$$\models (straight) \,\square\!\!\rightarrow (\mathbf{F}\ goal)$$

counterfactual worlds further away do not matter:

$$<_{far}$$

# Applications of Counterfactual Reasoning

Analyzing **causality** [Halpern '15], [Leitner-Fischer '15], [Coenen et al. '22]

Generating **explanations** for, e.g., model checkers [Beer et al. '09], [Wachter et al. '18]

Counterfactual **fairness** [Kusner et al. '17]

# Applications of Counterfactual Reasoning

Analyzing **causality** [Halpern '15], [Leitner-Fischer '15], [Coenen et al. '22]

**Definition 5 (Property Causality).** Let $\mathcal{T}$ be a system, $\pi \in traces(\mathcal{T})$ a trace, $\mathsf{C} \subseteq (2^I)^\omega$ a cause property, and $\mathsf{E} \subseteq (2^O)^\omega$ an effect property. We say that $\mathsf{C}$ is a cause of $\mathsf{E}$ on $\pi$ in $\mathcal{T}$ if the following three conditions hold:

**PC1:** $\pi \models \mathsf{C}$ and $\pi \models \mathsf{E}$, i.e., cause property and effect property are satisfied by the actual trace.

**PC2:** For every counterfactual input sequence $\sigma \in V_\pi^{\mathsf{C}}$, there is some contingency $\pi' \in C_\pi^\sigma$ s.t. $\pi' \not\models \mathsf{E}$, i.e., the counterfactual trace under contingency does not satisfy the effect property.

**PC3:** There is no $\mathsf{C}'$ s.t. $\mathsf{C}' \subset \mathsf{C}$ and $\mathsf{C}'$ satisfies PC1 and PC2.

[Coenen et al. '22]

**Definition 3.1:** (Actual cause) $\vec{X} = \vec{x}$ is an *actual cause of* $\varphi$ *in* $(M, \vec{u})$ if the following three conditions hold:

AC1. $(M, \vec{u}) \models (\vec{X} = \vec{x}) \wedge \varphi$. (That is, both $\vec{X} = \vec{x}$ and $\varphi$ are true in the actual world.)

AC2. There exists a partition $(\vec{Z}, \vec{W})$ of $\mathcal{V}$ with $\vec{X} \subseteq \vec{Z}$ and some setting $(\vec{x}', \vec{w}')$ of the variables in $(\vec{X}, \vec{W})$ such that if $(M, \vec{u}) \models Z = z^*$ for $Z \in \vec{Z}$, then

  (a) $(M, \vec{u}) \models [\vec{X} \leftarrow \vec{x}', \vec{W} \leftarrow \vec{w}']\neg\varphi$. In words, changing $(\vec{X}, \vec{W})$ from $(\vec{x}, \vec{w})$ to $(\vec{x}', \vec{w}')$ changes $\varphi$ from true to false,

  (b) $(M, \vec{u}) \models [\vec{X} \leftarrow \vec{x}, \vec{W} \leftarrow \vec{w}', \vec{Z}' \leftarrow \vec{z}^*]\varphi$ for all subsets $\vec{Z}'$ of $\vec{Z}$. In words, setting $\vec{W}$ to $\vec{w}'$ should have no effect on $\varphi$ as long as $\vec{X}$ is kept at its current value $\vec{x}$, even if all the variables in an arbitrary subset of $\vec{Z}$ are set to their original values in the context $\vec{u}$.

AC3. $\vec{X}$ is minimal; no subset of $\vec{X}$ satisfies conditions AC1 and AC2. Minimality ensures that only those elements of the conjunction $\vec{X} = \vec{x}$ that are essential for changing $\varphi$ in AC2(a) are considered part of a cause; inessential elements are pruned. ∎

[Halpern&Pearl '05]

# Applications of Counterfactual Reasoning

Analyzing **causality** [Halpern '15], [Leitner-Fischer '15], [Coenen et al. '22]

$$\varphi \wedge \psi \wedge \left( \left( \neg\varphi \boxdot\!\!\rightarrow_{min} \neg\psi \right) \vee \left( \neg\varphi \Diamond\!\!\rightarrow_{min} \neg\psi \right) \right)$$

[Coenen et al. '22]

$$\varphi \wedge \psi \wedge \\ \left( \neg\varphi \Diamond\!\!\rightarrow_{min} \neg\psi \right)$$

[Halpern '15]

**Minimal** counterfactuals

**Non-total** similarity relations

Lewis' Counterfactuals [Lewis '73]

Decision procedures for **satisfiability** and **trace checking** of counterfactual QPTL ($\omega$-regular) properties

Specification toolbox for, e.g., temporal causality

# Similarity Relation $\leq_{far}(\bigcirc)$

Encodes the distance of a world $\bigcirc$ from the reference world $\bigcirc$ .

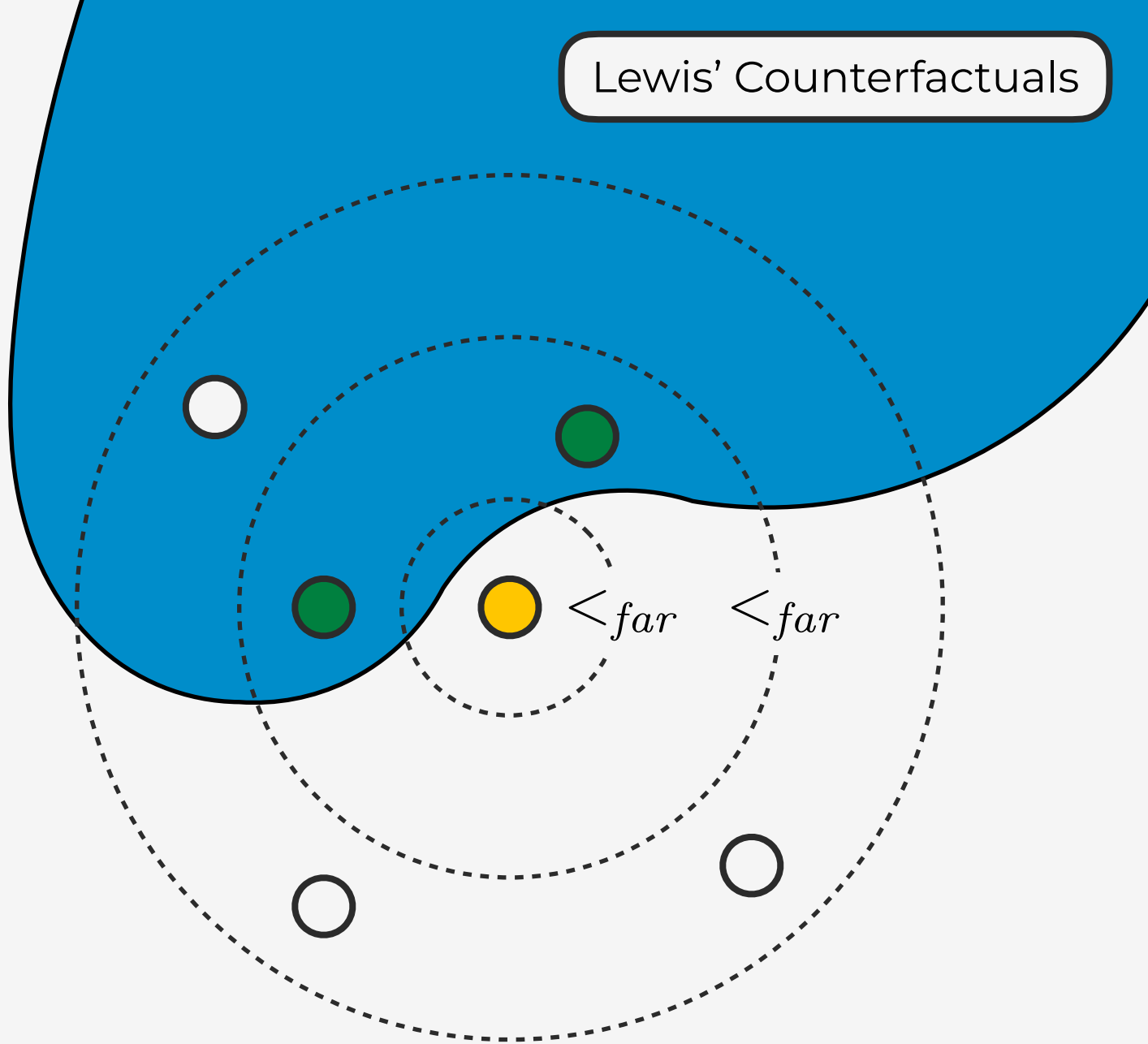Needs to be a total preorder.

$\bigcirc$ is the unique minimum.

$<_{far}$  $<_{far}$

# 'Would' Operator

$$\varphi \, \square\!\!\rightarrow \psi$$

All closest worlds satisfying $\varphi$ have to satisfy $\psi$ .

Worlds in spheres further away do not matter.

$<_{far}$ $<_{far}$

11

# 'Might' Operator

$$\varphi \diamond\!\!\longrightarrow \psi$$

Some closest world satisfying $\varphi$ has to satisfy $\psi$.

Again, worlds in spheres further away do not matter.

$<_{far}$   $<_{far}$

12

# The Limit Assumption

*"There always exist well-defined closest worlds satisfying $\varphi$ ."*

# The Limit Assumption

*"There always exist well-defined closest worlds satisfying $\varphi$ ."*

$<_{far}$   $<_{far}$   $<_{far}$   $<_{far}$   $<_{far}$   $<_{far}$

$x = 0.0$   $\cdots$   $x = 3.01$   $x = 3.05$   $x = 3.1$

Generally not satisfied and already **rejected by Lewis.**

$x \leq 3.0$   $x > 3.0$

14

$$\bullet \models \varphi \,\square\!\!\rightarrow\, \psi \qquad \text{iff:}$$

$$\models \psi$$

$$<_{far} \qquad <_{far} \qquad <_{far} \qquad <_{far} \qquad <_{far} \qquad <_{far}$$

There is a threshold world after which all closer $\varphi$-worlds satisfy $\psi$.

$$(1) \quad \exists \bullet : \bullet \models \varphi \wedge \forall \bigcirc : \bigcirc \leq_{far} \bullet \Rightarrow (\bigcirc \models \varphi \rightarrow \psi)$$

# Semantics of 'Would'

$$\bigcirc \models \varphi \,\square\!\!\rightarrow\, \psi \qquad \text{iff:}$$



Or: There are no $\varphi$-worlds (vacuous case).

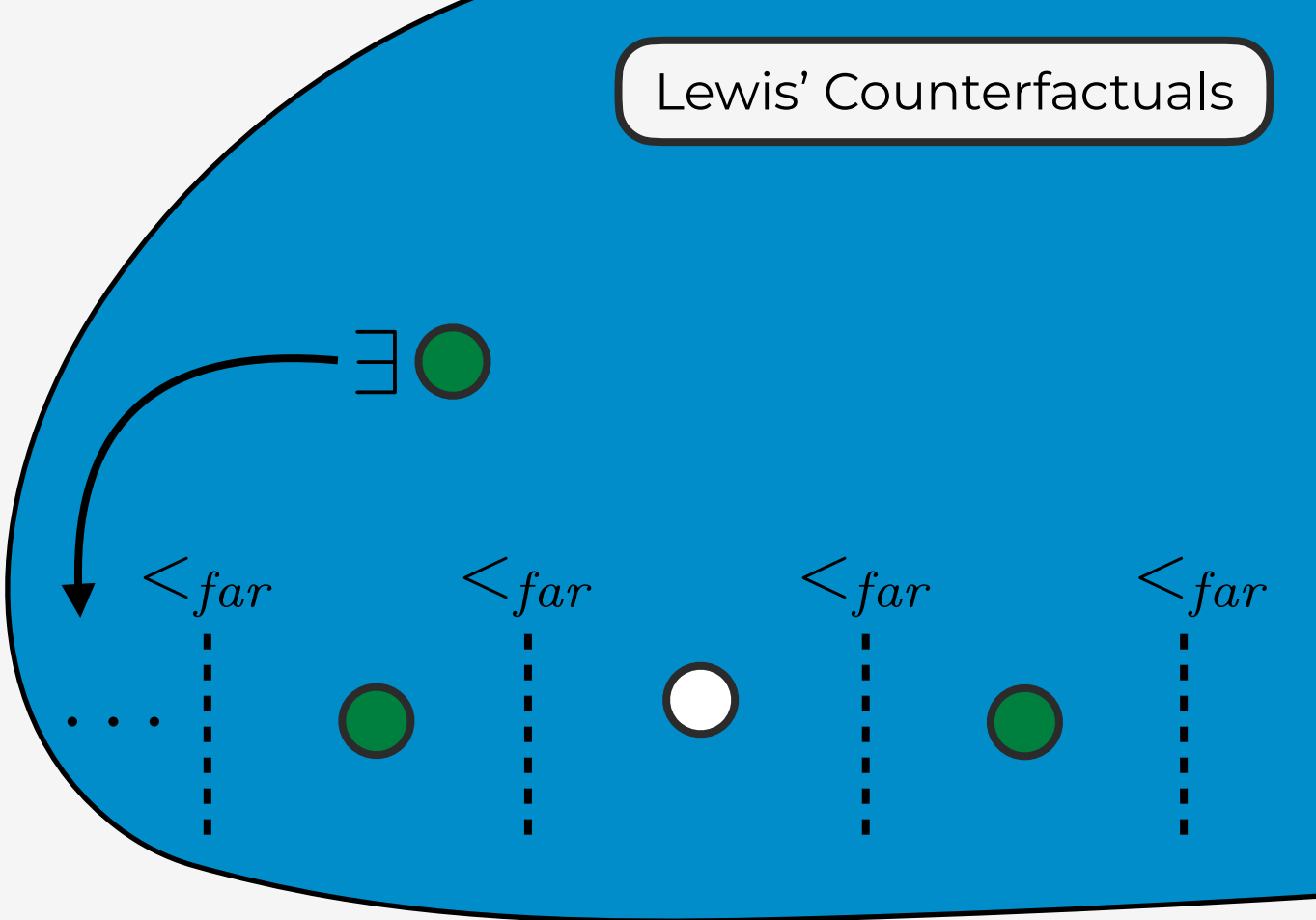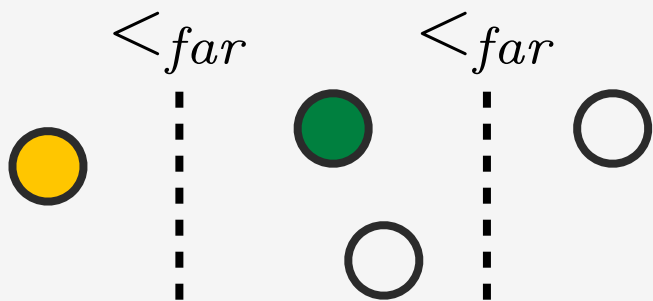(1) or (2): $\forall \bigcirc : \bigcirc \not\models \varphi$

# Semantics of 'Might'

$$\bigcirc \models \varphi \diamondsuit\!\!\longrightarrow \psi \qquad \text{iff:}$$

For any $\varphi$-world (and there is at least one) there exists a closer world satisfying $\psi$.

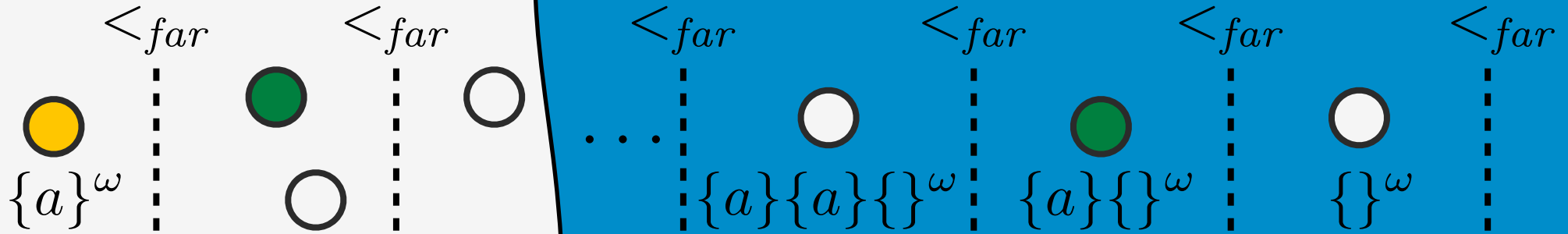$$\exists \bigcirc : \bigcirc \models \varphi \wedge \forall \bullet : \bullet \models \varphi \Rightarrow \exists \bullet : \bullet \leq_{far} \bullet \wedge \bullet \models \varphi$$

17

# The Limit Assumption

*"There always exist well-defined closest worlds satisfying $\varphi$ ."*

$<_{far}$     $<_{far}$     $<_{far}$     $<_{far}$     $<_{far}$     $<_{far}$

$\{a\}^{\omega}$        $\ldots$   $\{a\}\{a\}\{\}^{\omega}$   $\{a\}\{\}^{\omega}$   $\{\}^{\omega}$

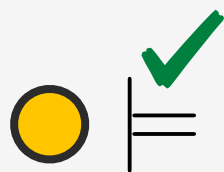Also does not hold for temporal properties.

$$\mathsf{G\,F}\,a \qquad \mathsf{F\,G}\,\neg a$$

In practice, counterfactual worlds are often incomparable.

Applying Lewis' semantics to non-total orders leads to unintuitive judgements.

$$\bullet \models \varphi \:\square\!\!\rightarrow\: \psi \quad \text{iff:}$$

$$(1)\ \exists \bullet : \bullet \models \varphi \land \forall \bigcirc : \bigcirc \leq_{far} \bullet \Rightarrow (\bigcirc \models \varphi \rightarrow \psi)$$
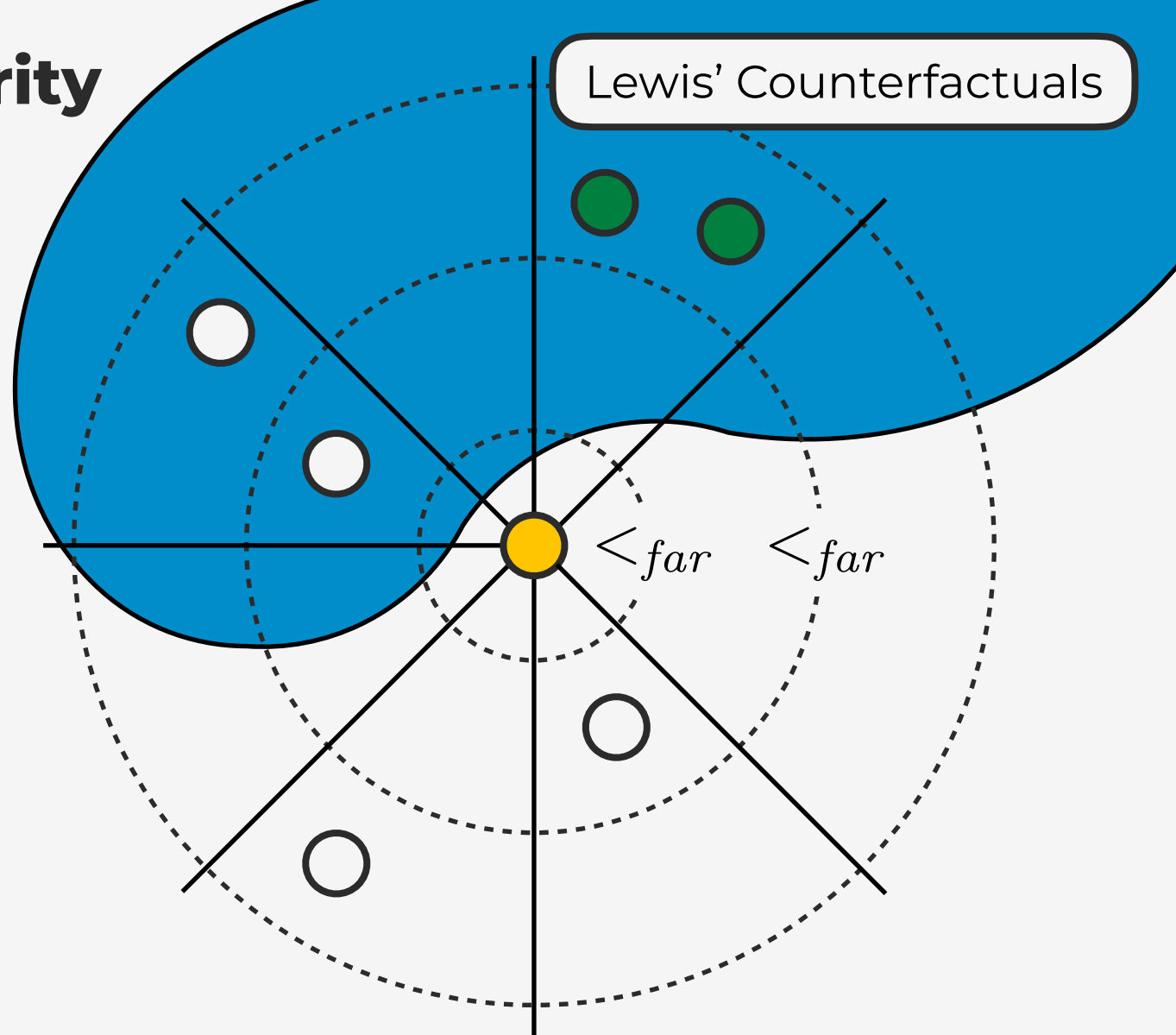
$<_{far} \quad <_{far}$

20

# The Problem with Linearity

Applying Lewis' semantics to non-total orders leads to unintuitive judgements.

$$\bigcirc \not\models \varphi \ \Diamond\!\!\!\rightarrow \psi \ \text{iff:}$$

$$\exists \bigcirc : \bigcirc \models \varphi \wedge \forall \bullet : \bullet \models \varphi \Rightarrow \exists \bullet : \bullet \leq_{far} \bullet \wedge \bullet \models \varphi$$
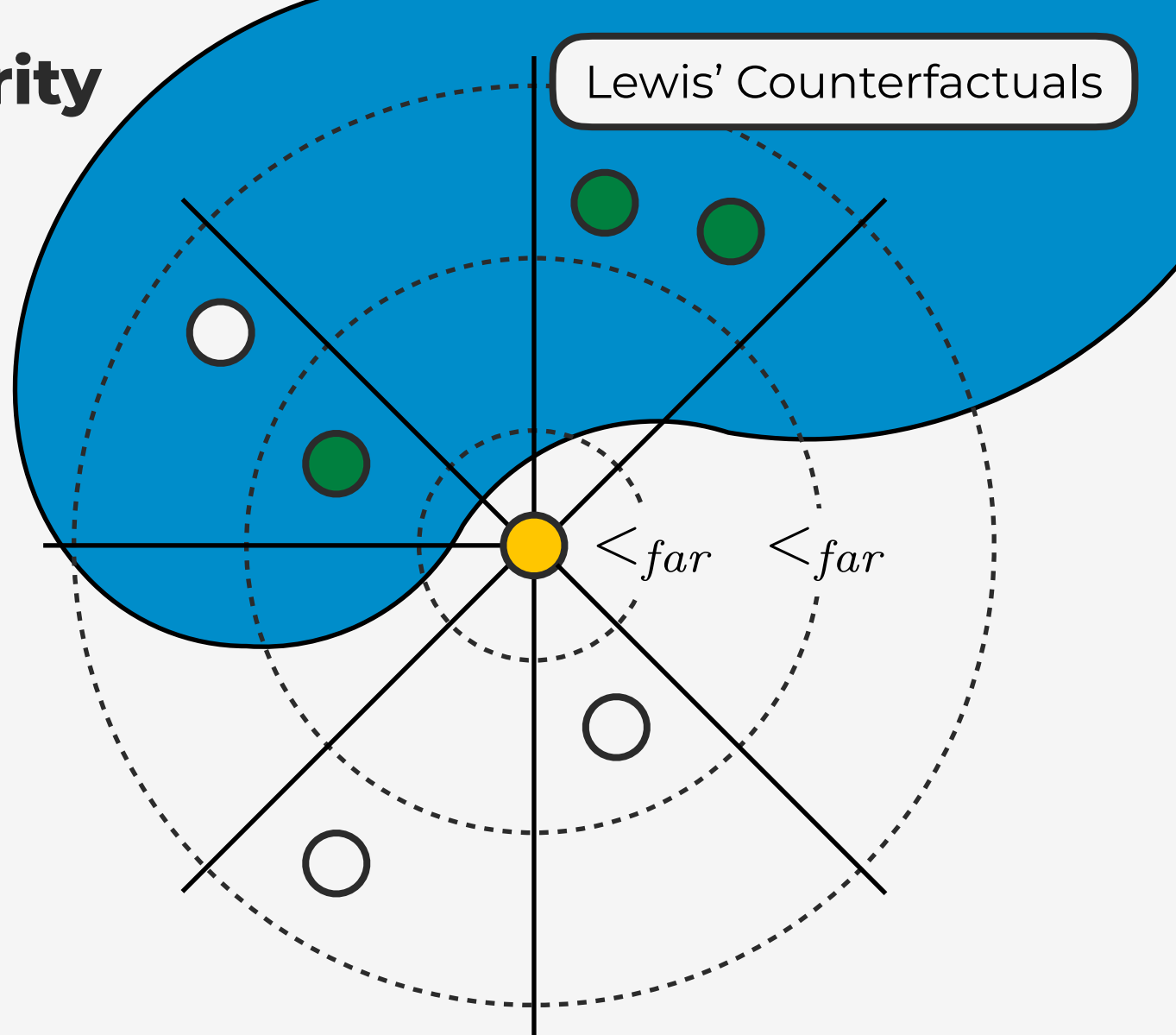
$<_{far} \quad <_{far}$

21

# The Problem with Linearity

Applying Lewis' semantics
to non-total orders leads
to unintuitive judgements.

$$\bullet \models \varphi \;\square\!\!\rightarrow\; \psi \text{ iff:}$$

$$\exists \bigcirc : \bigcirc \models \varphi \wedge \forall \bullet : \bullet \models \varphi \Rightarrow \exists \bullet : \bullet \leq_{far} \bullet \wedge \bullet \models \varphi$$

22

# **Fixing Lewis' Semantics**

The semantics of $\square\!\!\!\rightarrow$ is too weak and of $\diamondsuit\!\!\!\rightarrow$ too strong to capture the intended meaning on non-total relations.

We introduce operators with an additional level of quantification:

$\boxed{\cdot}\!\!\!\rightarrow$ 'Universal Would'
*"If [...],* **under all circumstances**, $\psi$ *would have been true as well."*

$\langle\!\!\cdot\!\!\rangle\!\!\!\rightarrow$

'Existential Might'
*"If [...],* **under some circumstance**, $\psi$ *might have been true, too."*
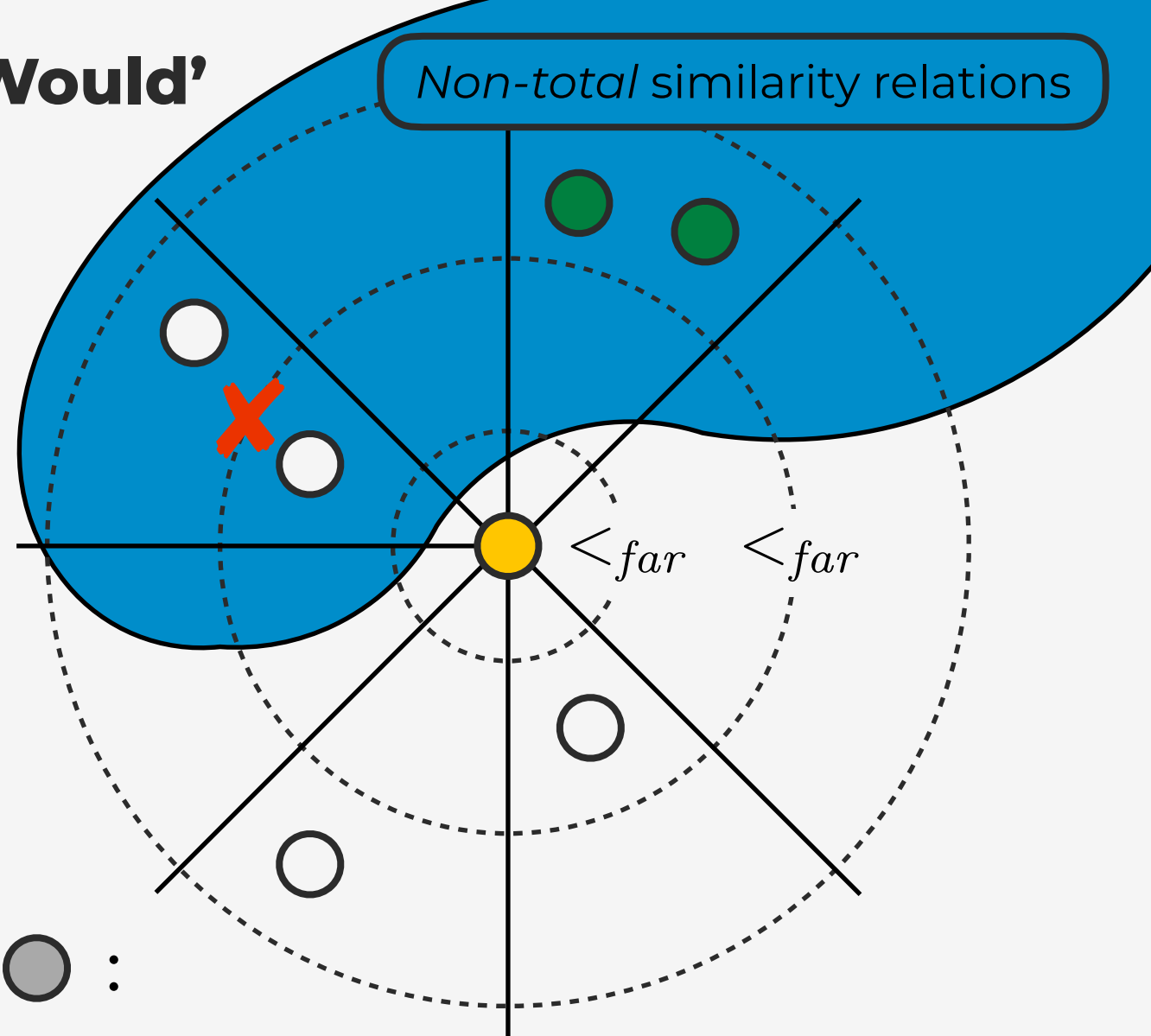
23

# Semantics of 'Universal Would'

We quantify universally over $\varphi$-worlds and require a closer threshold world 🟢 for all of them.

🟡 $\models \varphi \,\boxed{\cdot}\!\!\rightarrow \psi$ iff:

$\forall \bigcirc : \bigcirc \models \varphi \Rightarrow \exists \,🟢 :$

$🟢 \leq_{far} \bigcirc \wedge 🟢 \models \varphi \wedge \forall \,⚫ :$

$⚫ \leq_{far} 🟢 \Rightarrow (\,⚫ \models \varphi \rightarrow \psi)$

$<_{far} \quad <_{far}$

24

# Semantics of 'Universal Would'
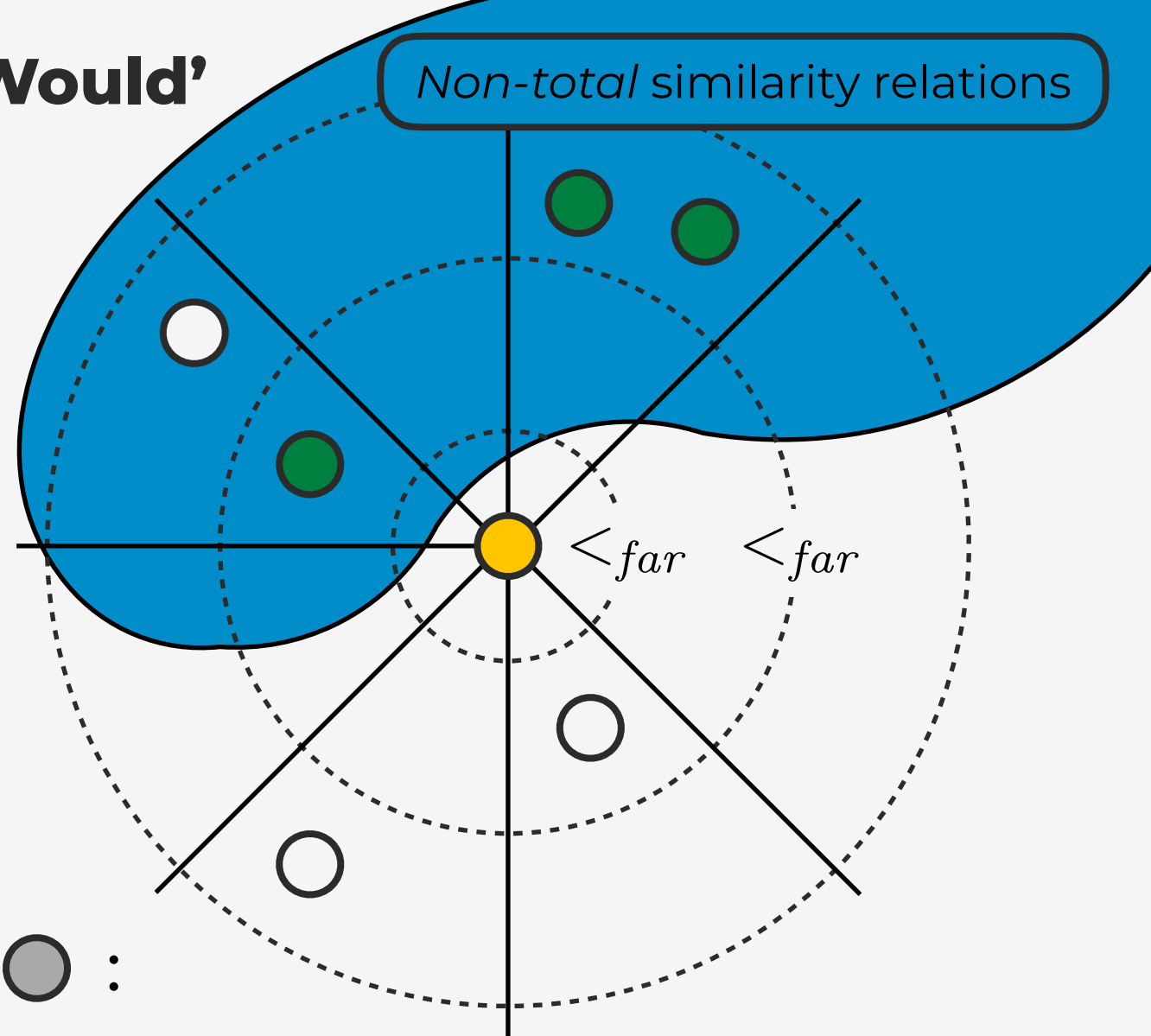
We quantify universally over $\varphi$-worlds and require a closer threshold world 🟢 for all of them.

$$\text{🟡} \models \varphi \;\boxdot\!\!\rightarrow\; \psi \quad\text{iff:}$$

$$\forall\, \bigcirc : \bigcirc \models \varphi \Rightarrow \exists\, \text{🟢} :$$

$$\text{🟢} \leq_{far} \bigcirc \wedge \text{🟢} \models \varphi \wedge \forall\, \text{⚫} :$$

$$\text{⚫} \leq_{far} \text{🟢} \Rightarrow (\text{⚫} \models \varphi \rightarrow \psi)$$

$$<_{far} \quad <_{far}$$

25

# **Semantics of 'Existential Might'**
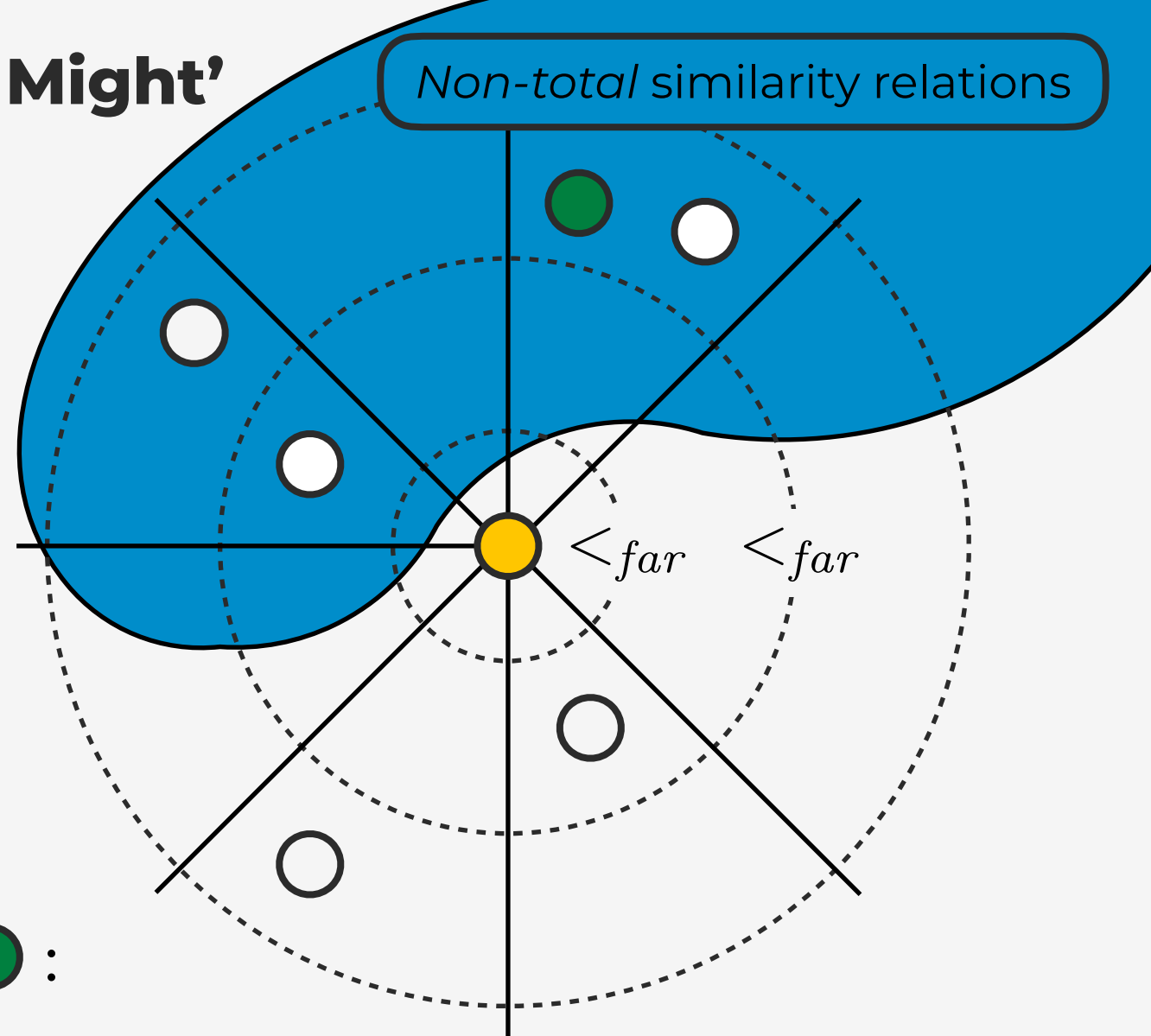
We quantify existentially over $\varphi$-worlds and require ever-closer worlds ● for one world only.

✔

$$\bigcirc \models \varphi \; \diamondsuit\!\!\!\!\bullet\!\!\!\longrightarrow \psi \text{ iff:}$$

$$\exists \bigcirc : \bigcirc \models \varphi \wedge \forall \bullet :$$

$$\bullet \leq_{far} \bigcirc \wedge \bullet \models \varphi \Rightarrow \exists \bullet :$$
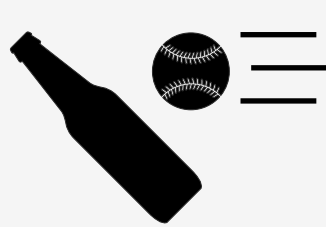
$$\bullet \leq_{far} \bullet \wedge \bullet \models \varphi$$

$<_{far} \quad <_{far}$

26

# Minimality

For causality, causes are counterfactuals that describe the *minimal changes necessary* to avoid the effect.
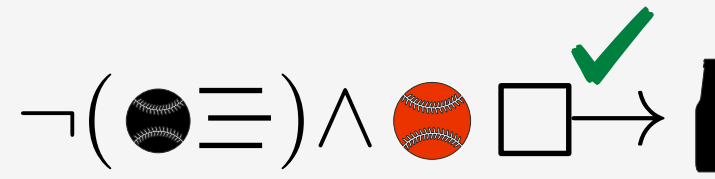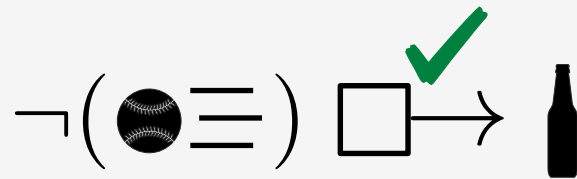
Ball is thrown and bottle breaks.

$<_{far}$

Ball is not thrown.

$<_{far}$
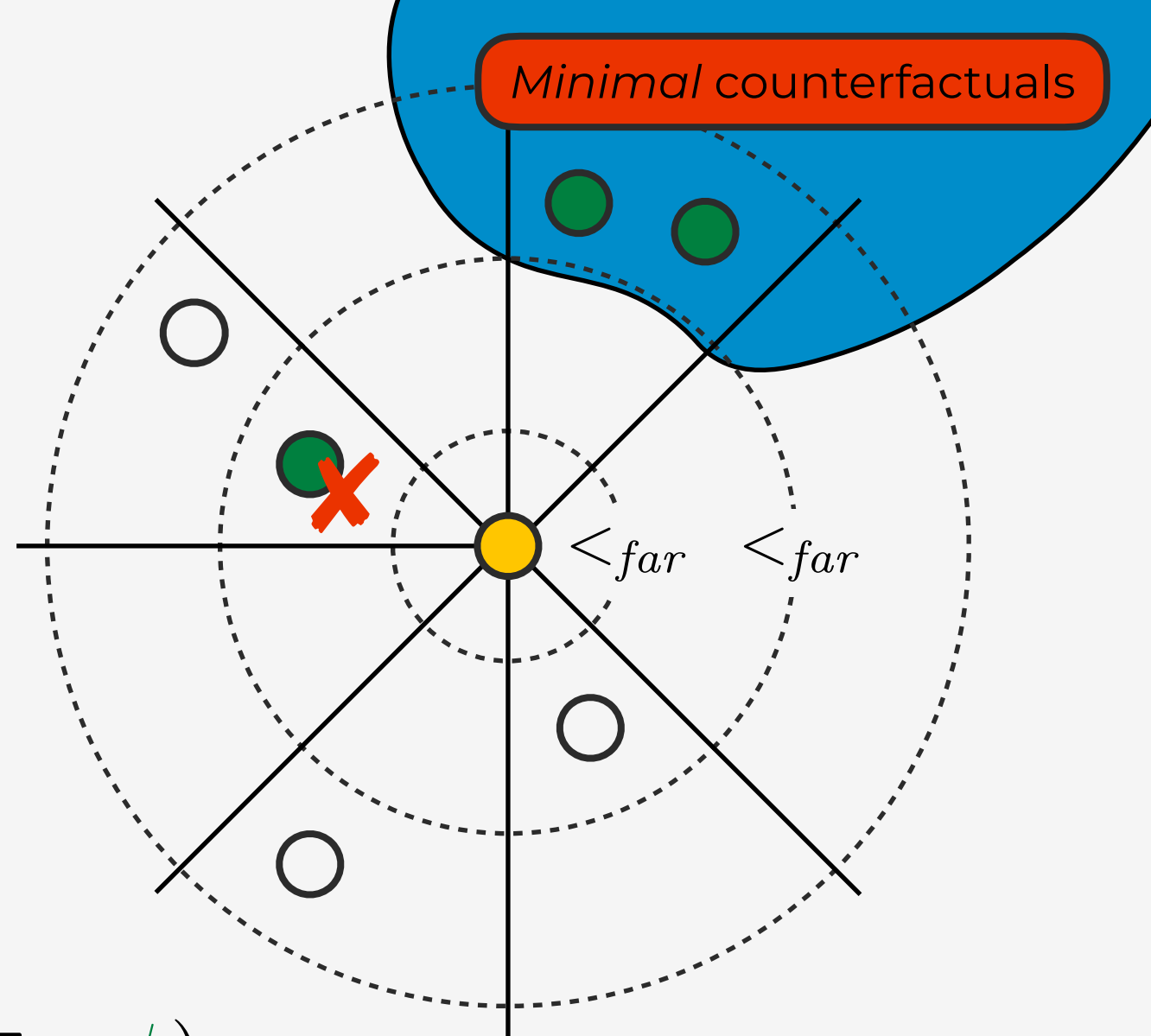
Ball is not thrown and is red.

Antecedent is not *minimal*.

27

# Minimal Counterfactuals

A counterfactual is minimal if its antecedent describes the largest set that qualifies, e.g.,

$$\bigcirc \models \varphi \boxdot\!\!\longrightarrow_{min} \psi \ \text{ iff:}$$

$$\bigcirc \models \varphi \boxdot\!\!\longrightarrow \psi \ \wedge$$

$$\neg \exists \theta : \theta \supset \varphi \wedge (\bigcirc \models \theta \boxdot\!\!\longrightarrow \psi)$$
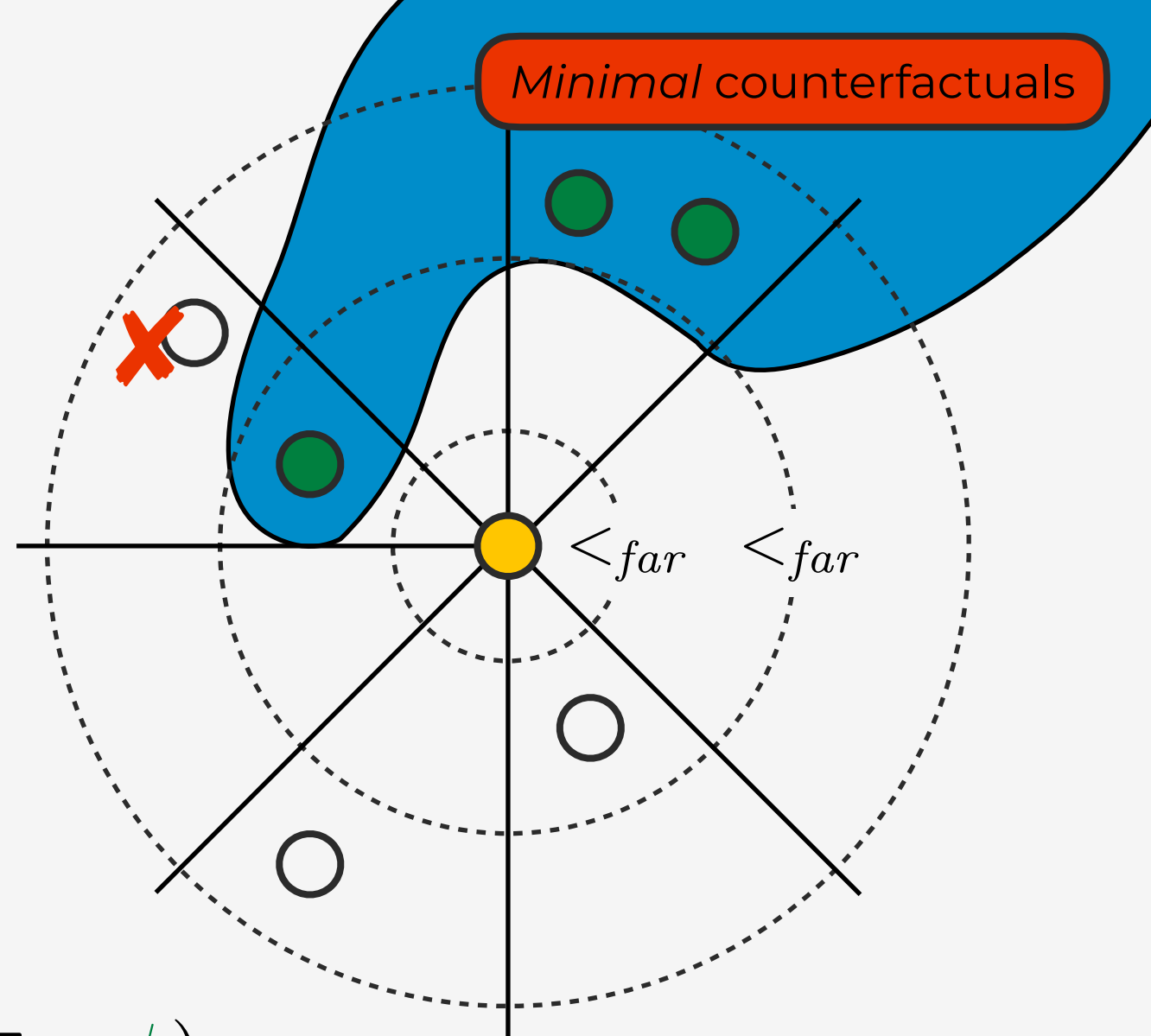
$<_{far} \quad <_{far}$

28

# Minimal Counterfactuals

A counterfactual is minimal if its antecedent describes the largest set that qualifies, e.g.,

$$\bullet \models \varphi \,\boxdot\!\!\longrightarrow_{min} \psi \quad \text{iff:}$$

$$\bullet \models \varphi \,\boxdot\!\!\longrightarrow \psi \;\wedge$$

$$\neg \exists \theta : \theta \supset \varphi \wedge ( \bullet \models \theta \,\boxdot\!\!\longrightarrow \psi )$$



$$<_{far} \quad <_{far}$$

29
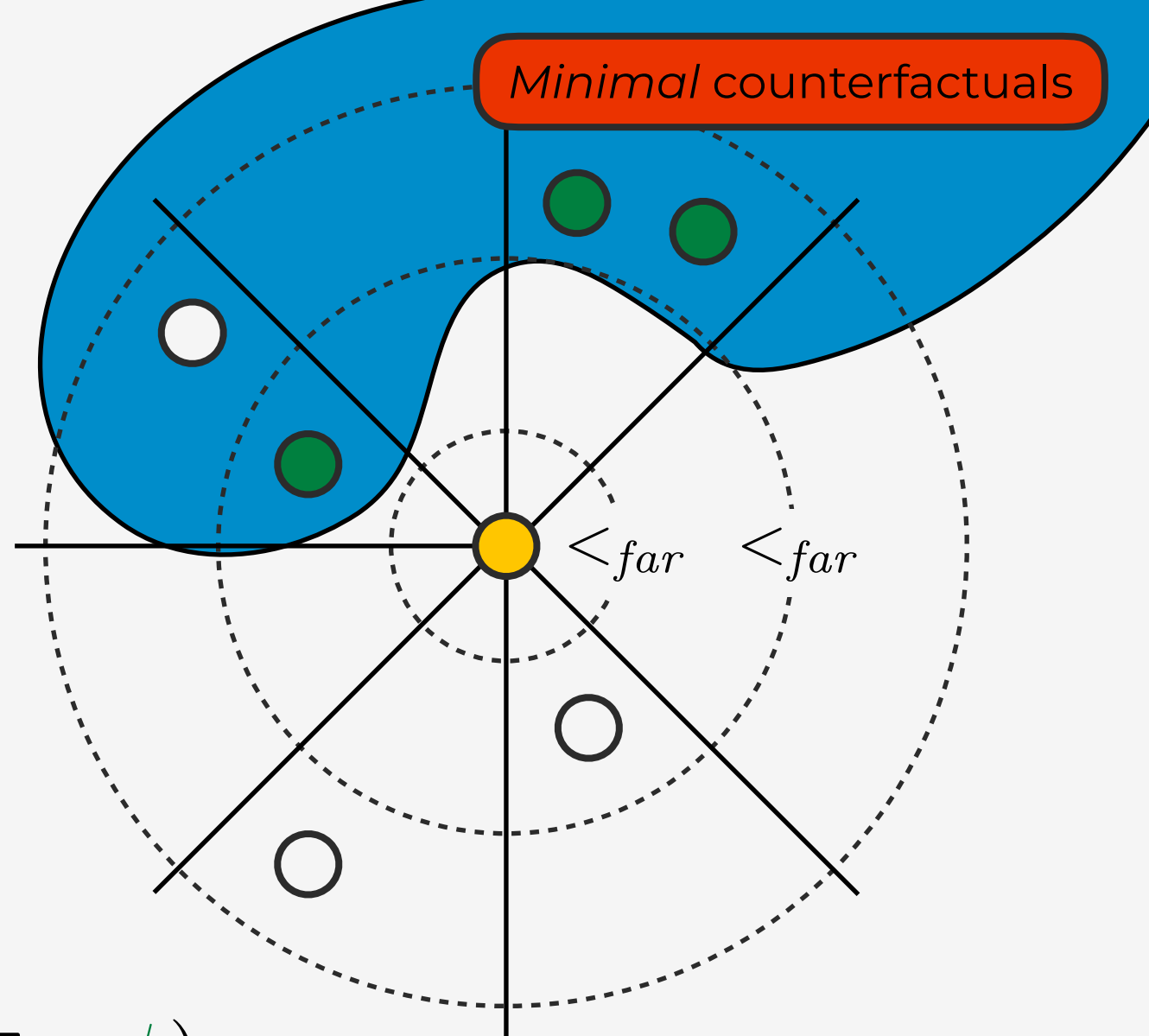
# Minimal Counterfactuals

A counterfactual is minimal if its antecedent describes the largest set that qualifies, e.g.,

$$\bullet \models \varphi \boxdot\!\!\!\rightarrow_{min} \psi \quad \text{iff:}$$

$$\bullet \models \varphi \boxdot\!\!\!\rightarrow \psi \ \wedge$$

$$\neg \exists \theta : \theta \supset \varphi \wedge ( \bullet \models \theta \boxdot\!\!\!\rightarrow \psi )$$
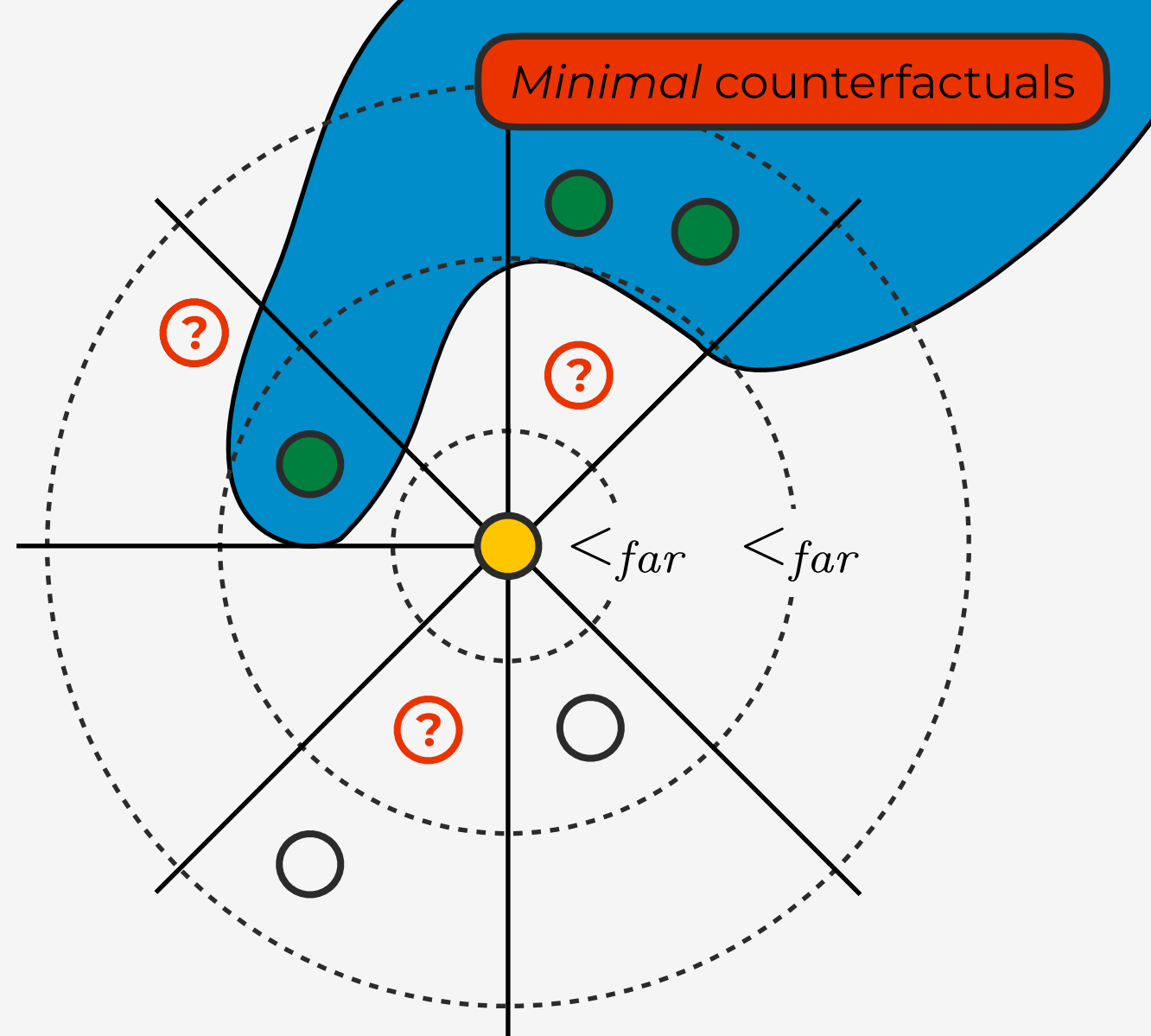
$$<_{far} \quad <_{far}$$

30

# Minimal Counterfactuals

The SO quantification can be avoided by searching for worlds ? that could be added to $\varphi$ .

$$\bigcirc \models \varphi \boxdot{\longrightarrow}_{min} \psi \quad \text{iff:}$$

Large FO-formula ($\forall/\exists\bigcirc$)
*See our paper!*



*Minimal* counterfactuals

$<_{far} \quad <_{far}$

31

# Moving to Temporal Logics

$$\varphi \mathbin{\square\!\!\rightarrow} \psi \quad \Longrightarrow \quad (straight) \mathbin{\square\!\!\rightarrow} (\mathbf{F}\ goal)$$

● ⟹ lasso trace, e.g. $\pi = \{a,b\}\{\}(\{a,b\})^{\omega}$

○ ⟹ infinite traces, e.g. $\pi' = \{a,b\}\{a,b\}\ldots$

$$\sigma \leq_{far}(\pi)\ \rho \quad \Longrightarrow \quad \text{QPTL formula, e.g.}\ \bigwedge_{a \in AP} (a_{\sigma} \not\leftrightarrow a_{\pi}) \rightarrow (a_{\rho} \not\leftrightarrow a_{\pi})$$

⟦○⟧ ⟹ 

$\forall/\exists\ ○ \quad \Longrightarrow \quad \forall/\exists\ \pi$ (hyperproperty)

$$(\mathsf{X}\, a \,\Box\!\!\rightarrow_{min}\, \mathsf{G}\, b) \wedge \ldots$$

Counterfactuals Modulo QPTL

➡️

$$\forall/\exists \bigcirc$$

FO formula

⬇️

$$\forall/\exists\, \pi$$

Prenex HyperQPTL

⬅️

HyperQPTL Model Checking
[Beutner&Finkbeiner, LPAR '23]

$$(\mathtt{X}\, a\, \square\!\!\!\rightarrow_{min} \mathtt{G}\, b) \wedge \ldots$$

Counterfactuals Modulo QPTL

$\forall/\exists \bigcirc$

FO formula

$\forall/\exists\,\pi$

Prenex HyperQPTL
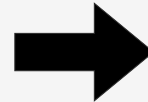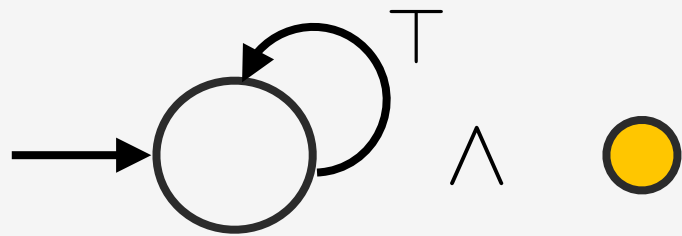
HyperQPTL Model Checking
[Beutner&Finkbeiner, LPAR '23]

# Conclusion

Uniform specification language for counterfactual reasoning in, e.g., causality, fairness etc.

Automatic decision procedures for the resulting theory modulo QPTL.

System-Level Counterfactuals
*@LPAR:* Counterfactuals Modulo Theories?