

# Reasoning about (Lack of) Knowledge

Christoph Schwering

UNSW Sydney

COMP4418, Week 7

# Motivation

John McCarthy (1927–2011):

- Stanford, MIT, Dartmouth
- Turing Award
- Invented Lisp (1958)
- Invented Garbage Collection (1959)
- Founding Father of AI (with Minsky, Newell, Simon, 1955)



# Motivation

John McCarthy (1927–2011):

- Stanford, MIT, Dartmouth
- Turing Award
- Invented Lisp (1958)
- Invented Garbage Collection (1959)
- Founding Father of AI (with Minsky, Newell, Simon, 1955)
- Proposed Advice Taker (1959)
  - ▶ *Programs with Common Sense*



# Motivation

John McCarthy (1927–2011):

- Stanford, MIT, Dartmouth
- Turing Award
- Invented Lisp (1958)
- Invented Garbage Collection (1959)
- Founding Father of AI (with Minsky, Newell, Simon, 1955)
- Proposed Advice Taker (1959)
  - ▶ *Programs with Common Sense*
  - ▶ Improve program behaviour by making statements to it



# Motivation

John McCarthy (1927–2011):

- Stanford, MIT, Dartmouth
- Turing Award
- Invented Lisp (1958)
- Invented Garbage Collection (1959)
- Founding Father of AI (with Minsky, Newell, Simon, 1955)
- Proposed Advice Taker (1959)
  - ▶ *Programs with Common Sense*
  - ▶ Improve program behaviour by making statements to it
  - ▶ Program draws conclusions from its knowledge



# Motivation

John McCarthy (1927–2011):

- Stanford, MIT, Dartmouth
- Turing Award
- Invented Lisp (1958)
- Invented Garbage Collection (1959)
- Founding Father of AI (with Minsky, Newell, Simon, 1955)
- Proposed Advice Taker (1959)
  - ▶ *Programs with Common Sense*
  - ▶ Improve program behaviour by making statements to it
  - ▶ Program draws conclusions from its knowledge
    - ▶ Declarative conclusion: new knowledge
    - ▶ Imperative conclusion: take action



# Motivation

John McCarthy (1927–2011):

- Stanford, MIT, Dartmouth
- Turing Award
- Invented Lisp (1958)
- Invented Garbage Collection (1959)
- Founding Father of AI (with Minsky, Newell, Simon, 1955)
- Proposed Advice Taker (1959)
  - ▶ *Programs with Common Sense*
  - ▶ Improve program behaviour by making statements to it
  - ▶ Program draws conclusions from its knowledge
    - ▶ Declarative conclusion: new knowledge
    - ▶ Imperative conclusion: take action
  - ▶ Remains a vision to this date



# Motivation

John McCarthy (1927–2011):

- Stanford, MIT, Dartmouth
- Turing Award
- Invented Lisp (1958)
- Invented Garbage Collection (1959)
- Founding Father of AI (with Minsky, Newell, Simon, 1955)
- Proposed Advice Taker (1959)
  - ▶ *Programs with Common Sense*
  - ▶ Improve program behaviour by making statements to it
  - ▶ Program draws conclusions from its knowledge
    - ▶ Declarative conclusion: new knowledge
    - ▶ Imperative conclusion: take action
  - ▶ Remains a vision to this date

Advice Taker motivates (directly or indirectly) a lot of AI research, in particular what we'll be studying for the next three weeks





# Motivation

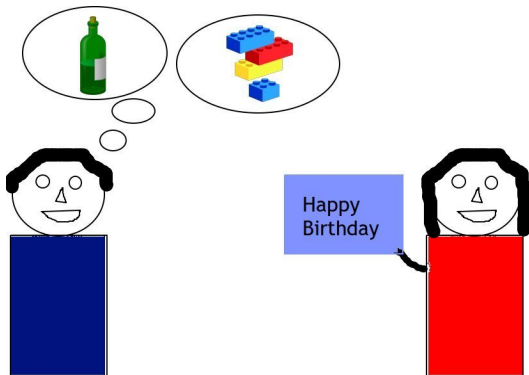
Observation: Non-knowledge is important

Not only what we know is relevant, but also what we *don't* know

# Motivation

Observation: Non-knowledge is important

Not only what we know is relevant, but also what we *don't* know



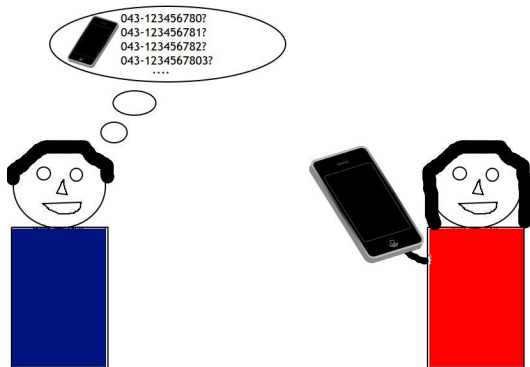
You don't know what's in the gift box.

You'll treat it with great care.

# Motivation

Observation: Non-knowledge is important

Not only what we know is relevant, but also what we *don't* know



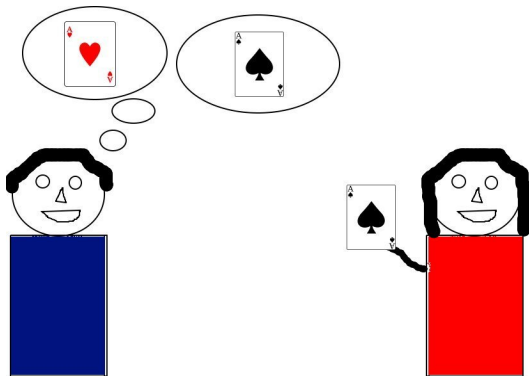
You know Jane has a phone, but you don't know her number.

You'll look it up.

# Motivation

Observation: Non-knowledge is important

Not only what we know is relevant, but also what we *don't* know

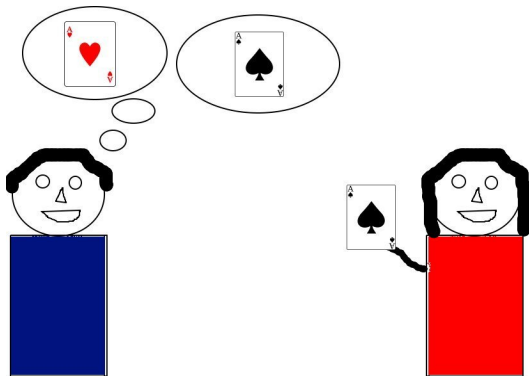


You know Jane is holding ace of spades *or* of hearts, but not which.  
You'll need a strategy that wins in either case.

# Motivation

Observation: Non-knowledge is important

Not only what we know is relevant, but also what we *don't* know



You know Jane is holding ace of spades *or* of hearts, but not which.

You'll need a strategy that wins in either case.

How can we accurately **capture knowledge and non-knowledge**?

# Overview of the Lecture

## ■ **A Logic of Knowledge – The Propositional Fragment**

- ▶ Why not classical logic?
- ▶ Syntax and semantics
- ▶ Omniscience, introspection, only-knowing
- ▶ Representation theorem

## ■ A Logic of Knowledge – The First-Order Case

## ■ Extensions of the Logic of Knowledge

## What Is a Knowledge Base?

- A **knowledge base** (KB) is a collection of sentences that describe (a fragment of) the world

## What Is a Knowledge Base?

- A **knowledge base** (KB) is a collection of sentences that describe (a fragment of) the world
  - KB completely characterises what the agent knows, i.e.,
    - ▶  $\alpha$  is known  $\implies$   $\text{KB} \models \alpha$
    - ▶  $\alpha$  is not known  $\implies$   $\text{KB} \not\models \alpha$
- $\implies$  KB is *all* the agent knows



# What Is a Knowledge Base?

- A **knowledge base** (KB) is a collection of sentences that describe (a fragment of) the world
- KB completely characterises what the agent knows, i.e.,
  - ▶  $\alpha$  is known  $\implies$   $\text{KB} \models \alpha$
  - ▶  $\alpha$  is not known  $\implies$   $\text{KB} \not\models \alpha$ $\implies$  KB is *all* the agent knows
- Purpose: evaluate queries
  - ▶ What is known? What is unknown?
  - ▶ Similar to a database, but draws inferences

# What Is a Knowledge Base?

- A **knowledge base** (KB) is a collection of sentences that describe (a fragment of) the world
- KB completely characterises what the agent knows, i.e.,
  - ▶  $\alpha$  is known  $\implies \text{KB} \models \alpha$
  - ▶  $\alpha$  is not known  $\implies \text{KB} \not\models \alpha$ $\implies$  KB is *all* the agent knows
- Purpose: evaluate queries
  - ▶ What is known? What is unknown?
  - ▶ Similar to a database, but draws inferences
- Usually: what is known  $\subsetneq$  what is true
  - ▶ Agent's knowledge is incomplete
  - ▶ Agent should be aware of that

# What Is a Knowledge Base?

- A **knowledge base** (KB) is a collection of sentences that describe (a fragment of) the world
- KB completely characterises what the agent knows, i.e.,
  - ▶  $\alpha$  is known  $\implies \text{KB} \models \alpha$
  - ▶  $\alpha$  is not known  $\implies \text{KB} \not\models \alpha$ $\implies$  KB is *all* the agent knows
- Purpose: evaluate queries
  - ▶ What is known? What is unknown?
  - ▶ Similar to a database, but draws inferences
- Usually: what is known  $\subsetneq$  what is true
  - ▶ Agent's knowledge is incomplete
  - ▶ Agent should be aware of that
- Usually: knowing is more than database lookup
  - ▶  $\alpha \in \text{KB} \implies \alpha$  is explicit knowledge (= database lookup)
  - ▶  $\text{KB} \models \alpha \implies \alpha$  is implicit knowledge (= logical inference)
  - ▶ Usually: explicit knowledge  $\subsetneq$  (implicit) knowledge

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .
2. You don't know that  $\neg r$ .

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .
2. You don't know that  $\neg r$ .
3. You know that  $p$  or  $q$ .

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .
2. You don't know that  $\neg r$ .
3. You know that  $p$  or  $q$ .
4. You don't know that  $p$ .
5. You don't know that  $q$ .

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .
2. You don't know that  $\neg r$ .
3. You know that  $p$  or  $q$ .
4. You don't know that  $p$ .
5. You don't know that  $q$ .
6. You know that  $p$  or  $q$ , but not which.

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .       $\text{KB} \not\models r$
2. You don't know that  $\neg r$ .       $\text{KB} \not\models \neg r$
3. You know that  $p$  or  $q$ .
4. You don't know that  $p$ .
5. You don't know that  $q$ .
6. You know that  $p$  or  $q$ , but not which.



## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .       $\text{KB} \not\models r$
2. You don't know that  $\neg r$ .       $\text{KB} \not\models \neg r$
3. You know that  $p$  or  $q$ .       $\text{KB} \models (p \vee q)$
4. You don't know that  $p$ .
5. You don't know that  $q$ .
6. You know that  $p$  or  $q$ , but not which.

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .       $\text{KB} \not\models r$
2. You don't know that  $\neg r$ .       $\text{KB} \not\models \neg r$
3. You know that  $p$  or  $q$ .       $\text{KB} \models (p \vee q)$
4. You don't know that  $p$ .       $\text{KB} \not\models p$
5. You don't know that  $q$ .       $\text{KB} \not\models q$
6. You know that  $p$  or  $q$ , but not which.

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .       $\text{KB} \not\models r$
2. You don't know that  $\neg r$ .       $\text{KB} \not\models \neg r$
3. You know that  $p$  or  $q$ .       $\text{KB} \models (p \vee q)$
4. You don't know that  $p$ .       $\text{KB} \not\models p$
5. You don't know that  $q$ .       $\text{KB} \not\models q$
6. You know that  $p$  or  $q$ , but not which.       $\text{KB} \models ???$

Problem: Classical logic cannot express 6 directly in one formula

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r) \wedge \neg k_p \wedge \dots$

Then:

1. You don't know that  $r$ .       $\text{KB} \models \neg k_r$
2. You don't know that  $\neg r$ .       $\text{KB} \models \neg k_{\text{not}_r}$
3. You know that  $p$  or  $q$ .       $\text{KB} \models k_{p\_or\_q}$
4. You don't know that  $p$ .       $\text{KB} \models \neg k_p$
5. You don't know that  $q$ .       $\text{KB} \models \neg k_q$
6. You know that  $p$  or  $q$ , but not which.

$$\text{KB} \models k_{p\_or\_q} \wedge \neg k_p \wedge \neg k_q$$

Problem: Classical logic cannot express 6 directly in one formula

Idea #1: Compile  $(p \vee q \vee c)$  to new atoms  $k_p, k_{p\_or\_q}, \dots$  ✗

Does not scale.

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .  $\text{KB} \models r = \text{U}$
2. You don't know that  $\neg r$ .  $\text{KB} \models \neg r = \text{U}$
3. You know that  $p$  or  $q$ .  $\text{KB} \models (p \vee q)$
4. You don't know that  $p$ .  $\text{KB} \models p = \text{U}$
5. You don't know that  $q$ .  $\text{KB} \models q = \text{U}$
6. You know that  $p$  or  $q$ , but not which.

$$\text{KB} \models (p \vee q) \wedge p = \text{U} \wedge q = \text{U}$$

Problem: Classical logic cannot express 6 directly in one formula

Idea #2: Three-valued logic  $\{0, 1, \text{U}\}$  **X**

How would  $\text{U} \vee \text{U}$  behave? Is it known? Unknown?

## Why Not Classical Logic?

Suppose all you know is  $(p \vee q \vee r) \wedge (p \vee q \vee \neg r)$

Then:

1. You don't know that  $r$ .  $\mathbf{OKB} \models \neg\mathbf{K}r$
2. You don't know that  $\neg r$ .  $\mathbf{OKB} \models \neg\mathbf{K}\neg r$
3. You know that  $p$  or  $q$ .  $\mathbf{OKB} \models \mathbf{K}(p \vee q)$
4. You don't know that  $p$ .  $\mathbf{OKB} \models \neg\mathbf{K}p$
5. You don't know that  $q$ .  $\mathbf{OKB} \models \neg\mathbf{K}q$
6. You know that  $p$  or  $q$ , but not which.

$$\mathbf{OKB} \models \mathbf{K}(p \vee q) \wedge \neg\mathbf{K}p \wedge \neg\mathbf{K}q$$

Problem: Classical logic cannot express 6 directly in one formula

Idea #3: Add unary operators  $\mathbf{O}$  and  $\mathbf{K}$  to express knowledge ✓

# The Language of $\mathcal{OL}_{PL}$

The language of only-knowing (propositional fragment)  $\mathcal{OL}_{PL}$ :

- $p, q, r, \dots$  atomic propositions
- $\neg\alpha$  "not  $\alpha$ "
- $(\alpha \vee \beta)$  " $\alpha$  or  $\beta$ "
- $(\alpha \wedge \beta)$  " $\alpha$  and  $\beta$ "
- $(\alpha \rightarrow \beta)$  " $\alpha$  implies  $\beta$ "
- $(\alpha \leftrightarrow \beta)$  " $\alpha$  is equivalent to  $\beta$ "
- $\mathbf{K}\alpha$  " $\alpha$  is known"
- $\mathbf{O}\alpha$  " $\alpha$  is *all* that is known"

# The Language of $\mathcal{OL}_{PL}$

The language of only-knowing (propositional fragment)  $\mathcal{OL}_{PL}$ :

- $p, q, r, \dots$  atomic propositions
- $\neg\alpha$  “not  $\alpha$ ”
- $(\alpha \vee \beta)$  “ $\alpha$  or  $\beta$ ”
- $(\alpha \wedge \beta) \stackrel{\text{def}}{=} \neg(\neg\alpha \vee \neg\beta)$  “ $\alpha$  and  $\beta$ ”
- $(\alpha \rightarrow \beta) \stackrel{\text{def}}{=} (\neg\alpha \vee \beta)$  “ $\alpha$  implies  $\beta$ ”
- $(\alpha \leftrightarrow \beta) \stackrel{\text{def}}{=} (\alpha \rightarrow \beta) \wedge (\beta \rightarrow \alpha)$  “ $\alpha$  is equivalent to  $\beta$ ”
- $\mathbf{K}\alpha$  “ $\alpha$  is known”
- $\mathbf{O}\alpha$  “ $\alpha$  is *all* that is known”



## Recap: Technical Terms (1)

A logical language is a *formal language* over an *alphabet* (here:  $\{p, q, r, \dots, (, ), \neg, \vee, \mathbf{K}, \mathbf{O}\}$ ) and a *grammar* (previous slide), i.e., rules that allow us to phrase sentences in that language.

## Recap: Technical Terms (1)

A logical language is a *formal language* over an *alphabet* (here:  $\{p, q, r, \dots, (, ), \neg, \vee, \mathbf{K}, \mathbf{O}\}$ ) and a *grammar* (previous slide), i.e., rules that allow us to phrase sentences in that language.

The sentences carry *no meaning by themselves*. We define a *model theory* to give them a *semantics*, i.e., to define what sort of formal structure interprets a sentence. Such an *interpretation*  $I$  satisfies a sentence  $\alpha$ , written  $I \models \alpha$ , or falsifies it, written  $I \not\models \alpha$ .

## Recap: Technical Terms (1)

A logical language is a *formal language* over an *alphabet* (here:  $\{p, q, r, \dots, (, ), \neg, \vee, \mathbf{K}, \mathbf{O}\}$ ) and a *grammar* (previous slide), i.e., rules that allow us to phrase sentences in that language.

The sentences carry *no meaning by themselves*. We define a *model theory* to give them a *semantics*, i.e., to define what sort of formal structure interprets a sentence. Such an *interpretation*  $I$  satisfies a sentence  $\alpha$ , written  $I \models \alpha$ , or falsifies it, written  $I \not\models \alpha$ .

A typical rule of a semantics is

$$I \models (\alpha \vee \beta) \text{ if and only if } I \models \alpha \text{ or } I \models \beta.$$

Note that  $\vee$  is a symbol of the logical language, whereas “if and only if” and “or” are natural language expressions. The rule says that the symbol “ $\vee$ ” corresponds to the natural language expression “or”.

## Recap: Technical Terms (1)

A logical language is a *formal language* over an *alphabet* (here:  $\{p, q, r, \dots, (, ), \neg, \vee, \mathbf{K}, \mathbf{O}\}$ ) and a *grammar* (previous slide), i.e., rules that allow us to phrase sentences in that language.

The sentences carry *no meaning by themselves*. We define a *model theory* to give them a *semantics*, i.e., to define what sort of formal structure interprets a sentence. Such an *interpretation*  $I$  satisfies a sentence  $\alpha$ , written  $I \models \alpha$ , or falsifies it, written  $I \not\models \alpha$ .

A typical rule of a semantics is

$$I \models (\alpha \vee \beta) \text{ if and only if } I \models \alpha \text{ or } I \models \beta.$$

Note that  $\vee$  is a symbol of the logical language, whereas “if and only if” and “or” are natural language expressions. The rule says that the symbol “ $\vee$ ” corresponds to the natural language expression “or”.

We will sometimes take the liberty to omit brackets to ease readability. For instance, we write  $(p \vee q \vee r)$  instead of  $((p \vee q) \vee r)$  or  $(p \vee (q \vee r))$ , implicitly assuming our semantics of  $\vee$  is associative.

## Recap: Technical Terms (2)

The form of such an interpretation varies between logics. Propositional logic uses truth tables, first-order logic usually uses structures with a domain and interpretation function.

- When an interpretation  $I$  satisfies a sentence, we write  $I \models \alpha$ .
- When all interpretations satisfy a sentence  $\alpha$ , then  $\alpha$  is *valid* and we write  $\models \alpha$ .
- When all interpretations that satisfy some sentence  $\Sigma$  or set of sentences  $\Sigma$  also satisfy  $\alpha$ , we say  $\Sigma$  *entails*  $\alpha$  and write  $\Sigma \models \alpha$ .

## Recap: Technical Terms (2)

The form of such an interpretation varies between logics. Propositional logic uses truth tables, first-order logic usually uses structures with a domain and interpretation function.

- When an interpretation  $I$  satisfies a sentence, we write  $I \models \alpha$ .
- When all interpretations satisfy a sentence  $\alpha$ , then  $\alpha$  is *valid* and we write  $\models \alpha$ .
- When all interpretations that satisfy some sentence  $\Sigma$  or set of sentences  $\Sigma$  also satisfy  $\alpha$ , we say  $\Sigma$  *entails*  $\alpha$  and write  $\Sigma \models \alpha$ .

Different semantics are possible. What justifies a semantics? Typically there is a *proof theory*, and model theory and proof theory should be equivalent ( $\models \alpha$  if and only if  $\vdash \alpha$ ). Nevertheless, we will only focus on the semantics in the next weeks.

## The Semantics of $\mathcal{OL}_{PL}$

Definition: semantics of  $\mathcal{OL}_{PL}$

A **world**  $w$  is a function from the atomic propositions to  $\{0, 1\}$ .

# The Semantics of $\mathcal{OL}_{PL}$

## Definition: semantics of $\mathcal{OL}_{PL}$

A **world**  $w$  is a function from the atomic propositions to  $\{0, 1\}$ .

- $w \models P \iff w[P] = 1$
- $w \models \neg\alpha \iff w \not\models \alpha$
- $w \models (\alpha \vee \beta) \iff w \models \alpha \text{ or } w \models \beta$



# The Semantics of $\mathcal{OL}_{PL}$

## Definition: semantics of $\mathcal{OL}_{PL}$

A **world**  $w$  is a function from the atomic propositions to  $\{0, 1\}$ .

- $w \models P \iff w[P] = 1$
- $w \models \neg\alpha \iff w \not\models \alpha$
- $w \models (\alpha \vee \beta) \iff w \models \alpha \text{ or } w \models \beta$
- $w \models \mathbf{K}\alpha \iff ???$
- $w \models \mathbf{O}\alpha \iff ???$

# The Semantics of $\mathcal{OL}_{PL}$

## Definition: semantics of $\mathcal{OL}_{PL}$

A **world**  $w$  is a function from the atomic propositions to  $\{0, 1\}$ .

An **epistemic state**  $e$  is a set of worlds.

- $w \models P \iff w[P] = 1$
- $w \models \neg\alpha \iff w \not\models \alpha$
- $w \models (\alpha \vee \beta) \iff w \models \alpha \text{ or } w \models \beta$
- $w \models \mathbf{K}\alpha \iff ???$
- $w \models \mathbf{O}\alpha \iff ???$

# The Semantics of $\mathcal{OL}_{PL}$

## Definition: semantics of $\mathcal{OL}_{PL}$

A **world**  $w$  is a function from the atomic propositions to  $\{0, 1\}$ .

An **epistemic state**  $e$  is a set of worlds.

$$\blacksquare e, w \models P \iff w[P] = 1$$

$$\blacksquare e, w \models \neg\alpha \iff e, w \not\models \alpha$$

$$\blacksquare e, w \models (\alpha \vee \beta) \iff e, w \models \alpha \text{ or } e, w \models \beta$$

$$\blacksquare e, w \models \mathbf{K}\alpha \iff ???$$

$$\blacksquare e, w \models \mathbf{O}\alpha \iff ???$$

# The Semantics of $\mathcal{OL}_{PL}$

## Definition: semantics of $\mathcal{OL}_{PL}$

A **world**  $w$  is a function from the atomic propositions to  $\{0, 1\}$ .

An **epistemic state**  $e$  is a set of worlds.

$$\blacksquare e, w \models P \iff w[P] = 1$$

$$\blacksquare e, w \models \neg\alpha \iff e, w \not\models \alpha$$

$$\blacksquare e, w \models (\alpha \vee \beta) \iff e, w \models \alpha \text{ or } e, w \models \beta$$

$$\blacksquare e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w' \in e \Rightarrow e, w' \models \alpha$$

$$\blacksquare e, w \models \mathbf{O}\alpha \iff ???$$

" $\Rightarrow$ " stands for natural language expressions "only if".

" $\Leftrightarrow$ " and " $\iff$ " stand for natural language expressions "if and only if".

# The Semantics of $\mathcal{OL}_{PL}$

## Definition: semantics of $\mathcal{OL}_{PL}$

A **world**  $w$  is a function from the atomic propositions to  $\{0, 1\}$ .

An **epistemic state**  $e$  is a set of worlds.

$$\blacksquare e, w \models P \iff w[P] = 1$$

$$\blacksquare e, w \models \neg\alpha \iff e, w \not\models \alpha$$

$$\blacksquare e, w \models (\alpha \vee \beta) \iff e, w \models \alpha \text{ or } e, w \models \beta$$

$$\blacksquare e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w' \in e \Rightarrow e, w' \models \alpha$$

$$\blacksquare e, w \models \mathbf{O}\alpha \iff \text{for all worlds } w', w' \in e \Leftrightarrow e, w' \models \alpha$$

" $\Rightarrow$ " stands for natural language expressions "only if".

" $\Leftrightarrow$ " and " $\iff$ " stand for natural language expressions "if and only if".

# Abbreviations

Recall:

$$\blacksquare (\alpha \wedge \beta) \stackrel{\text{def}}{=} \neg(\neg\alpha \vee \neg\beta)$$

$$\blacksquare (\alpha \rightarrow \beta) \stackrel{\text{def}}{=} (\neg\alpha \vee \beta)$$

$$\blacksquare (\alpha \leftrightarrow \beta) \stackrel{\text{def}}{=} (\alpha \rightarrow \beta) \wedge (\beta \rightarrow \alpha)$$

$\wedge$  should be “and”

$\rightarrow$  should be “only if”

$\leftrightarrow$  should be “if and only if”.

## Lemma: abbreviations

$$\blacksquare e, w \models \alpha \wedge \beta \iff e, w \models \alpha \text{ and } e, w \models \beta$$

$$\blacksquare e, w \models \alpha \rightarrow \beta \iff e, w \models \alpha \Rightarrow e, w \models \beta$$

$$\blacksquare e, w \models \alpha \leftrightarrow \beta \iff e, w \models \alpha \Leftrightarrow e, w \models \beta$$

Proof on paper

## Some Lemmas

### Definition: objective, subjective

If  $\phi$  mentions no atoms inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\phi$  is **objective**.

If  $\sigma$  mentions atoms only inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\sigma$  is **subjective**.

- $((p \vee q) \wedge p \wedge q)$  is objective
- $\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$  is subjective

## Some Lemmas

### Definition: objective, subjective

If  $\phi$  mentions no atoms inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\phi$  is **objective**.

If  $\sigma$  mentions atoms only inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\sigma$  is **subjective**.

- $((p \vee q) \wedge p \wedge q)$  is objective
- $\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$  is subjective

### Lemma: objective, subjective

Let  $\phi$  be objective. Then  $e, w \models \phi \iff e', w \models \phi$ .

Let  $\sigma$  be subjective. Then  $e, w \models \sigma \iff e, w' \models \sigma$ .



## Some Lemmas

### Definition: objective, subjective

If  $\phi$  mentions no atoms inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\phi$  is **objective**.

If  $\sigma$  mentions atoms only inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\sigma$  is **subjective**.

- $((p \vee q) \wedge p \wedge q)$  is objective
- $\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$  is subjective

### Lemma: objective, subjective

Let  $\phi$  be objective. Then  $e, w \models \phi \iff e', w \models \phi$ .

Let  $\sigma$  be subjective. Then  $e, w \models \sigma \iff e, w' \models \sigma$ .

When  $\phi$  is objective, " $w \models \phi$ " stands for "for every  $e$ ,  $e, w \models \phi$ ".

When  $\sigma$  is subjective, " $e \models \sigma$ " stands for "for every  $w$ ,  $e, w \models \sigma$ ".

Proof on paper

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

$$\text{Let } e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

- $w \in e$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e$

$$\iff w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e$

$$\iff w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$$

$$\iff w \models (p \vee q \vee r) \text{ and } w \models (p \vee q \vee \neg r)$$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e$

$$\iff w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$$

$$\iff w \models (p \vee q \vee r) \text{ and } w \models (p \vee q \vee \neg r)$$

$$\iff w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 1, \text{ and} \\ w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 0$$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e$

$$\iff w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$$

$$\iff w \models (p \vee q \vee r) \text{ and } w \models (p \vee q \vee \neg r)$$

$$\iff w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 1, \text{ and} \\ w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 0$$

$$\iff w[p] = 1 \text{ or } w[q] = 1$$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

- $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$
- $e \models \mathbf{K}(p \vee q) \quad ?$



## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q) \quad ?$

$\iff \text{for all } w, w \in e \Rightarrow w \models (p \vee q)$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q) \quad ?$

$\iff$  for all  $w, w \in e \Rightarrow w \models (p \vee q)$

$\iff$  for all  $w, w \in e \Rightarrow w \models p \text{ or } w \models q$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q) \quad ?$

$\iff$  for all  $w, w \in e \Rightarrow w \models (p \vee q)$

$\iff$  for all  $w, w \in e \Rightarrow w \models p \text{ or } w \models q$

$\iff$  for all  $w, w \in e \Rightarrow w[p] = 1 \text{ or } w[q] = 1$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q) \quad ?$

$\iff$  for all  $w, w \in e \Rightarrow w \models (p \vee q)$

$\iff$  for all  $w, w \in e \Rightarrow w \models p \text{ or } w \models q$

$\iff$  for all  $w, w \in e \Rightarrow w[p] = 1 \text{ or } w[q] = 1$

$\iff$  for all  $w, (w[p] = 1 \text{ or } w[q] = 1) \Rightarrow (w[p] = 1 \text{ or } w[q] = 1) \quad \checkmark$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ?

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ?

$\iff e \not\models \mathbf{K}p$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg\mathbf{K}p$  ?

$\iff e \not\models \mathbf{K}p$

$\iff \text{for some } w, w \in e \text{ and } w \not\models p$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ?

$\iff e \not\models \mathbf{K}p$

$\iff$  for some  $w, w \in e$  and  $w \not\models p$

$\iff$  for some  $w, w \in e$  and  $w[p] \neq 1$



## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ?

$\iff e \not\models \mathbf{K}p$

$\iff \text{for some } w, w \in e \text{ and } w \not\models p$

$\iff \text{for some } w, w \in e \text{ and } w[p] \neq 1$

$\iff \text{for some } w, w[p] = 1 \text{ or } w[q] = 1, \text{ and } w[p] \neq 1$

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ?

$\iff e \not\models \mathbf{K}p$

$\iff$  for some  $w, w \in e$  and  $w \not\models p$

$\iff$  for some  $w, w \in e$  and  $w[p] \neq 1$

$\iff$  for some  $w, w[p] = 1$  or  $w[q] = 1$ , and  $w[p] \neq 1$

$\iff$  for some  $w, w[p] \neq 1$  and  $w[q] = 1$  ✓

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ✓

■  $e \models \mathbf{K}(p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q$  ?

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ✓

■  $e \models \mathbf{K}(p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q$  ?

$\iff e \models \mathbf{K}(p \vee q) \text{ and } e \models \neg \mathbf{K}p \text{ and } e \models \neg \mathbf{K}q$  ✓

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

- $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$
- $e \models \mathbf{K}(p \vee q)$  ✓
- $e \models \neg \mathbf{K}p$  ✓
- $e \models \mathbf{K}(p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q$  ✓
- $e \models \mathbf{O}((p \vee q \vee r) \wedge (p \vee q \vee \neg r))$  ?

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$
$$e, w \models \mathbf{O}\alpha \iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

■  $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$

■  $e \models \mathbf{K}(p \vee q)$  ✓

■  $e \models \neg \mathbf{K}p$  ✓

■  $e \models \mathbf{K}(p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q$  ✓

■  $e \models \mathbf{O}((p \vee q \vee r) \wedge (p \vee q \vee \neg r))$  ?

$\iff$  for all  $w, w \in e \Leftrightarrow w \models ((p \vee q \vee r) \wedge (p \vee q \vee \neg r))$  ✓

## Examples

$$e, w \models \mathbf{K}\alpha \iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha$$
$$e, w \models \mathbf{O}\alpha \iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha$$

Let  $e \stackrel{\text{def}}{=} \{w \mid w \models (p \vee q \vee r) \wedge (p \vee q \vee \neg r)\}$

- $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$
- $e \models \mathbf{K}(p \vee q)$  ✓
- $e \models \neg \mathbf{K}p$  ✓
- $e \models \mathbf{K}(p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q$  ✓
- $e \models \mathbf{O}((p \vee q \vee r) \wedge (p \vee q \vee \neg r))$  ✓

# Logical Omniscience

Logical omniscience means that an agent knows all the consequences of what they know. In particular, they know all valid sentences.



# Logical Omniscience

Logical omniscience means that an agent knows all the consequences of what they know. In particular, they know all valid sentences.

## Theorem: logical omniscience

If  $\models \alpha \rightarrow \beta$ , then  $\models \mathbf{K}\alpha \rightarrow \mathbf{K}\beta$ .

In particular: If  $\models \alpha$ , then  $\models \mathbf{K}\alpha$ .

# Logical Omniscience

Logical omniscience means that an agent knows all the consequences of what they know. In particular, they know all valid sentences.

## Theorem: logical omniscience

If  $\models \alpha \rightarrow \beta$ , then  $\models \mathbf{K}\alpha \rightarrow \mathbf{K}\beta$ .

In particular: If  $\models \alpha$ , then  $\models \mathbf{K}\alpha$ .

Logical omniscience is often problematic:

- Philosophical problem: most agents are not omniscient
- Practical problem: omniscience makes reasoning intractable

We will look at methods to avoid these problems next week.

Proof on paper

# Only-Knowing

The purpose of only-knowing is to capture a knowledge base.

Knowledge bases are usually objective.

The corresponding epistemic state is then unique:

## Theorem: unique-model property

Let  $\phi$  be objective. Then there is a unique  $e$  such that  $e \models \mathbf{O}\phi$ .

An entailment problem  $\mathbf{O}\phi \models \mathbf{K}\alpha$  thus reduces to model checking:

$e \models \mathbf{K}\alpha$ , where  $e = \{w \mid w \models \phi\}$ ?

# Self-Knowledge

We can nest **K** operators to say that we know that we know.

Complete and accurate knowledge about own knowledge:

## Theorem: positive and negative introspection

Positive introspection:  $\models \mathbf{K}\alpha \rightarrow \mathbf{K}\mathbf{K}\alpha$

Negative introspection:  $\models \neg\mathbf{K}\alpha \rightarrow \mathbf{K}\neg\mathbf{K}\alpha$

Why?

$e \models (\neg)\mathbf{K}\alpha \implies e, w \models (\neg)\mathbf{K}\alpha$  for all  $w \in e \iff e, w \models \mathbf{K}(\neg)\mathbf{K}\alpha$ .

## Representation Theorem

Can we solve  $\mathbf{OKB} \models \alpha$  with ordinary, propositional reasoning?

That is, can we eliminate  $\mathbf{K}$  and  $\mathbf{O}$ ?

Then we could use standard reasoning system.

## Representation Theorem

Can we solve  $\mathbf{OKB} \models \alpha$  with ordinary, propositional reasoning?

That is, can we eliminate  $\mathbf{K}$  and  $\mathbf{O}$ ?

Then we could use standard reasoning system.

### Theorem

Let  $\mathbf{KB}$ ,  $\phi$  be objective. Then  $\mathbf{OKB} \models \mathbf{K}\phi \iff \mathbf{KB} \models \phi$ .

**Idea:** replace nested  $\mathbf{K}\phi$  with `TRUE` if  $\mathbf{KB} \models \phi$ , otherwise with `FALSE`.

# Representation Theorem

Can we solve  $\mathbf{OKB} \models \alpha$  with ordinary, propositional reasoning?

That is, can we eliminate  $\mathbf{K}$  and  $\mathbf{O}$ ?

Then we could use standard reasoning system.

## Theorem

Let  $\mathbf{KB}$ ,  $\phi$  be objective. Then  $\mathbf{OKB} \models \mathbf{K}\phi \iff \mathbf{KB} \models \phi$ .

**Idea:** replace nested  $\mathbf{K}\phi$  with  $\text{TRUE}$  if  $\mathbf{KB} \models \phi$ , otherwise with  $\text{FALSE}$ .

Ex.: Let  $\mathbf{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$\mathbf{OKB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$ ?

# Representation Theorem

Can we solve  $\mathbf{OKB} \models \alpha$  with ordinary, propositional reasoning?

That is, can we eliminate  $\mathbf{K}$  and  $\mathbf{O}$ ?

Then we could use standard reasoning system.

## Theorem

Let  $\mathbf{KB}$ ,  $\phi$  be objective. Then  $\mathbf{OKB} \models \mathbf{K}\phi \iff \mathbf{KB} \models \phi$ .

**Idea:** replace nested  $\mathbf{K}\phi$  with  $\mathbf{TRUE}$  if  $\mathbf{KB} \models \phi$ , otherwise with  $\mathbf{FALSE}$ .

Ex.: Let  $\mathbf{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$$\mathbf{OKB} \models \mathbf{K}((p \vee q) \wedge \underbrace{\neg \mathbf{K}p}_{\mathbf{KB} \models p? \text{ X}} \wedge \underbrace{\neg \mathbf{K}q}_{\mathbf{KB} \models q? \text{ X}})?$$



# Representation Theorem

Can we solve  $\mathbf{OKB} \models \alpha$  with ordinary, propositional reasoning?

That is, can we eliminate  $\mathbf{K}$  and  $\mathbf{O}$ ?

Then we could use standard reasoning system.

## Theorem

Let  $\mathbf{KB}$ ,  $\phi$  be objective. Then  $\mathbf{OKB} \models \mathbf{K}\phi \iff \mathbf{KB} \models \phi$ .

**Idea:** replace nested  $\mathbf{K}\phi$  with  $\mathbf{TRUE}$  if  $\mathbf{KB} \models \phi$ , otherwise with  $\mathbf{FALSE}$ .

Ex.: Let  $\mathbf{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$\mathbf{OKB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{FALSE} \wedge \neg \mathbf{FALSE})?$

# Representation Theorem

Can we solve  $\mathbf{OKB} \models \alpha$  with ordinary, propositional reasoning?

That is, can we eliminate  $\mathbf{K}$  and  $\mathbf{O}$ ?

Then we could use standard reasoning system.

## Theorem

Let  $\mathbf{KB}$ ,  $\phi$  be objective. Then  $\mathbf{OKB} \models \mathbf{K}\phi \iff \mathbf{KB} \models \phi$ .

**Idea:** replace nested  $\mathbf{K}\phi$  with  $\mathbf{TRUE}$  if  $\mathbf{KB} \models \phi$ , otherwise with  $\mathbf{FALSE}$ .

Ex.: Let  $\mathbf{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$\mathbf{KB} \models ((p \vee q) \wedge \neg \mathbf{FALSE} \wedge \neg \mathbf{FALSE})?$  ✓

# Representation Theorem

Can we solve  $\mathbf{OKB} \models \alpha$  with ordinary, propositional reasoning?

That is, can we eliminate  $\mathbf{K}$  and  $\mathbf{O}$ ?

Then we could use standard reasoning system.

## Theorem

Let  $\mathbf{KB}$ ,  $\phi$  be objective. Then  $\mathbf{OKB} \models \mathbf{K}\phi \iff \mathbf{KB} \models \phi$ .

**Idea:** replace nested  $\mathbf{K}\phi$  with  $\text{TRUE}$  if  $\mathbf{KB} \models \phi$ , otherwise with  $\text{FALSE}$ .

Ex.: Let  $\mathbf{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$\mathbf{KB} \models ((p \vee q) \wedge \neg \text{FALSE} \wedge \neg \text{FALSE})?$  ✓

Next slide formalises this idea.

Sneak preview: It'll become more difficult in the first-order case:

What would you replace  $\mathbf{K}Q(x)$  with in  $\mathbf{K}\exists x (P(x) \wedge \neg \mathbf{K}Q(x))$ ? We'll see later.

## Representation Theorem (2)

### Definition: representation operators

For objective KB and  $\phi$ , let  $\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \begin{cases} \text{TRUE} & \text{if } \text{KB} \models \phi \\ \text{FALSE} & \text{otherwise} \end{cases}$   
where TRUE is some tautology (e.g.,  $p \vee \neg p$ ) and FALSE is  $\neg \text{TRUE}$ .

## Representation Theorem (2)

### Definition: representation operators

For objective KB and  $\phi$ , let  $\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \begin{cases} \text{TRUE} & \text{if } \text{KB} \models \phi \\ \text{FALSE} & \text{otherwise} \end{cases}$

where TRUE is some tautology (e.g.,  $p \vee \neg p$ ) and FALSE is  $\neg \text{TRUE}$ .

- $\|P\|_{\text{KB}} \stackrel{\text{def}}{=} P$
- $\|\neg\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \neg\|\alpha\|_{\text{KB}}$
- $\|(\alpha \vee \beta)\|_{\text{KB}} \stackrel{\text{def}}{=} (\|\alpha\|_{\text{KB}} \vee \|\beta\|_{\text{KB}})$
- $\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$

## Representation Theorem (2)

### Definition: representation operators

For objective KB and  $\phi$ , let  $\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \begin{cases} \text{TRUE} & \text{if } \text{KB} \models \phi \\ \text{FALSE} & \text{otherwise} \end{cases}$

where TRUE is some tautology (e.g.,  $p \vee \neg p$ ) and FALSE is  $\neg \text{TRUE}$ .

- $\|P\|_{\text{KB}} \stackrel{\text{def}}{=} P$
- $\|\neg\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \neg\|\alpha\|_{\text{KB}}$
- $\|(\alpha \vee \beta)\|_{\text{KB}} \stackrel{\text{def}}{=} (\|\alpha\|_{\text{KB}} \vee \|\beta\|_{\text{KB}})$
- $\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$

### Theorem: representation theorem

$\text{OKB} \models \alpha \iff \models \|\alpha\|_{\text{KB}}$ .

## Example

$$\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$$

$$\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \text{“KB} \models \phi\text{?”}$$

Let  $\text{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$\mathbf{O} \text{KB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$  ?

## Example

$$\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$$

$$\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \text{"KB} \models \phi\text{"}$$

Let  $\text{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$\mathbf{O} \text{KB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$

$\iff \models \|\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)\|_{\text{KB}}$



## Example

$$\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$$

$$\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \text{“KB} \models \phi\text{?”}$$

Let  $\text{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$$\mathbf{O} \text{KB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$$

$$\iff \models \|\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)\|_{\text{KB}}$$

$$\iff \models \text{RES}[\text{KB}, \|((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)\|_{\text{KB}}]$$

## Example

$$\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$$

$$\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \text{“KB} \models \phi\text{”}$$

Let  $\text{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$$\mathbf{O}\text{KB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$$

$$\iff \models \|\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)\|_{\text{KB}}$$

$$\iff \models \text{RES}[\text{KB}, \|(p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q\|_{\text{KB}}]$$

$$\iff \models \text{RES}[\text{KB}, ((p \vee q) \wedge \underbrace{\neg \|\mathbf{K}p\|_{\text{KB}}}_{\text{KB} \models p?} \wedge \underbrace{\neg \|\mathbf{K}q\|_{\text{KB}}}_{\text{KB} \models q?})]$$

## Example

$$\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$$

$$\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \text{“KB} \models \phi\text{?”}$$

Let  $\text{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$$\mathbf{O}\text{KB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$$

$$\iff \models \|\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)\|_{\text{KB}}$$

$$\iff \models \text{RES}[\text{KB}, \|\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)\|_{\text{KB}}]$$

$$\iff \models \text{RES}[\text{KB}, ((p \vee q) \wedge \underbrace{\neg \|\mathbf{K}p\|_{\text{KB}}}_{\text{KB} \models p?} \wedge \underbrace{\neg \|\mathbf{K}q\|_{\text{KB}}}_{\text{KB} \models q?})]$$

$$\iff \models \underbrace{\text{RES}[\text{KB}, ((p \vee q) \wedge \neg \text{FALSE} \wedge \neg \text{FALSE})]}_{\text{KB} \models (p \vee q) \wedge \neg \text{FALSE} \wedge \neg \text{FALSE?}}$$

## Example

$$\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$$

$$\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \text{“KB} \models \phi\text{?”}$$

Let  $\text{KB} \stackrel{\text{def}}{=} (p \vee q \vee r) \wedge (p \vee q \vee \neg r)$ .

$$\mathbf{O} \text{KB} \models \mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)$$

$$\iff \models \|\mathbf{K}((p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q)\|_{\text{KB}}$$

$$\iff \models \text{RES}[\text{KB}, \|(p \vee q) \wedge \neg \mathbf{K}p \wedge \neg \mathbf{K}q\|_{\text{KB}}]$$

$$\iff \models \text{RES}[\text{KB}, ((p \vee q) \wedge \underbrace{\neg \|\mathbf{K}p\|_{\text{KB}}}_{\text{KB} \models p?} \wedge \underbrace{\neg \|\mathbf{K}q\|_{\text{KB}}}_{\text{KB} \models q?})]$$

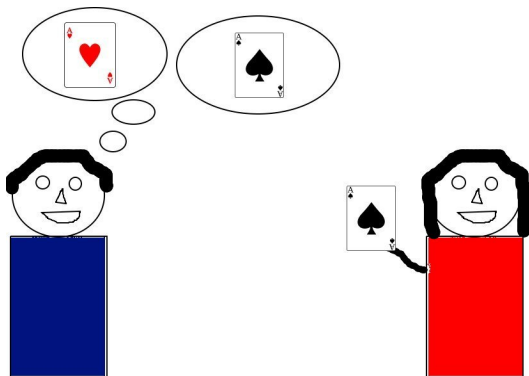
$$\iff \models \underbrace{\text{RES}[\text{KB}, ((p \vee q) \wedge \neg \text{FALSE} \wedge \neg \text{FALSE})]}_{\text{KB} \models (p \vee q) \wedge \neg \text{FALSE} \wedge \neg \text{FALSE?}}$$

$$\iff \models \text{TRUE} \quad \checkmark$$

# Overview of the Lecture

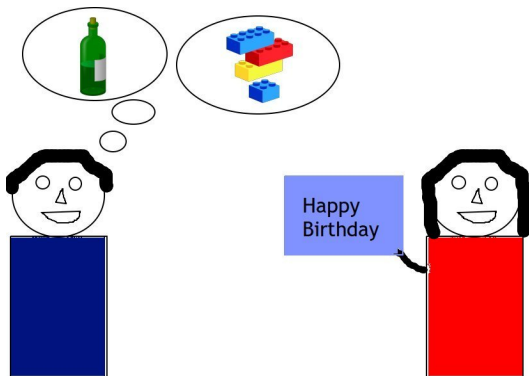
- A Logic of Knowledge – The Propositional Fragment
- **A Logic of Knowledge – The First-Order Case**
  - ▶ Why first-order logic?
  - ▶ Syntax and semantics
  - ▶ Knowing that vs knowing what
  - ▶ Representation theorem
- Extensions of the Logic of Knowledge

## Why Is $\mathcal{OL}_{PL}$ Not Enough?



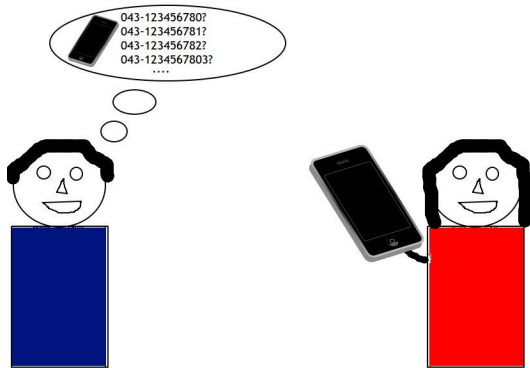
$$K((\spadesuit \vee \heartsuit) \wedge \neg K\spadesuit \wedge \neg K\heartsuit)$$

## Why Is $\mathcal{OL}_{PL}$ Not Enough?



$$\mathbf{K}\exists x (\text{InBox}(x) \wedge \neg \mathbf{K}\text{InBox}(x))$$

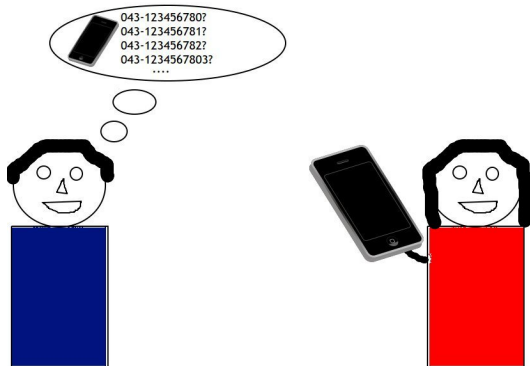
## Why Is $\mathcal{OL}_{PL}$ Not Enough?



$$\mathbf{K}\exists x (\text{numberOf}(\text{Jane}) = x \wedge \neg \mathbf{K}\text{numberOf}(\text{Jane}) = x)$$



## Why Is $\mathcal{OL}_{PL}$ Not Enough?



$\mathbf{K}\exists x (\text{numberOf}(\text{Jane}) = x \wedge \neg \mathbf{K}\text{numberOf}(\text{Jane}) = x)$

“all” or “some”  $\implies$  first-order quantification

# The Language of $\mathcal{OL}$

## Terms:

- $x, x', x_1, x_2, \dots$  first-order variables
- $\#1, \#2, \#3, \dots$  standard names
- $f(t_1, \dots, t_j)$  functions

## Formulas:

- $P(t_1, \dots, t_j)$  atomic formulas
- $t_1 = t_2$  equality expressions
- $\exists x \alpha$  "for some  $x, \alpha$ "
- $\forall x \alpha$  "for all  $x, \alpha$ "
- $\neg \alpha$   $(\alpha \vee \beta)$   $(\alpha \wedge \beta)$   $(\alpha \rightarrow \beta)$   $(\alpha \leftrightarrow \beta)$  **K** $\alpha$  **O** $\alpha$

# The Language of $\mathcal{OL}$

## Terms:

- $x, x', x_1, x_2, \dots$  first-order variables
- $\#1, \#2, \#3, \dots$  standard names
- $f(t_1, \dots, t_j)$  functions

## Formulas:

- $P(t_1, \dots, t_j)$  atomic formulas
- $t_1 = t_2$  equality expressions
- $\exists x \alpha$  "for some  $x, \alpha$ "
- $\forall x \alpha \stackrel{\text{def}}{=} \neg \exists x \neg \alpha$  "for all  $x, \alpha$ "
- $\neg \alpha \quad (\alpha \vee \beta) \quad (\alpha \wedge \beta) \quad (\alpha \rightarrow \beta) \quad (\alpha \leftrightarrow \beta) \quad \mathbf{K}\alpha \quad \mathbf{O}\alpha$

# Why Standard Names?

- Consider in classical logic:

$\text{fatherOf}(\text{Sally}) = \text{bestFriend}(\text{Jane}) \wedge$

$\text{fatherOf}(\text{Sally}) = \text{bossOf}(\text{John})$

- ▶ Who is father of Sally?
- ▶ "Jane's best friend" is not a good answer
- ▶ "John's boss" is not a good answer
- ▶ Classical logic offers no way of identifying him
- ▶ Reason: interpretations  $\langle D, \Phi \rangle$  have different domains

- Standard names correspond to an implicit infinite domain

- Standard names allow to *identify* individuals in formulas:

$\text{fatherOf}(\text{Sally}) = \text{Frank}$

## The Semantics of $\mathcal{OL}$ (1)

### Definition: semantics of $\mathcal{OL}$ (1)

$f(\vec{n})$  or  $P(\vec{n})$  are **primitive** iff all  $n_i$  are standard names.

A term or a formula is **ground** iff it mentions no variable.

# The Semantics of $\mathcal{OL}$ (1)

## Definition: semantics of $\mathcal{OL}$ (1)

$f(\vec{n})$  or  $P(\vec{n})$  are **primitive** iff all  $n_i$  are standard names.

A term or a formula is **ground** iff it mentions no variable.

A **world**  $w$  is a function that maps

- primitive functions  $f(\vec{n})$  to standard names
- primitive atomic formulas  $P(\vec{n})$  to  $\{0, 1\}$

# The Semantics of $\mathcal{OL}$ (1)

## Definition: semantics of $\mathcal{OL}$ (1)

$f(\vec{n})$  or  $P(\vec{n})$  are **primitive** iff all  $n_i$  are standard names.

A term or a formula is **ground** iff it mentions no variable.

A **world**  $w$  is a function that maps

- primitive functions  $f(\vec{n})$  to standard names
- primitive atomic formulas  $P(\vec{n})$  to  $\{0, 1\}$

The **denotation** of a ground term w.r.t.  $w$  is defined as

- $w(n) \stackrel{\text{def}}{=} n$  for every standard name  $n$
- $w(f(t_1, \dots, t_j)) \stackrel{\text{def}}{=} w[f(w(t_1), \dots, w(t_j))]$

# The Semantics of $\mathcal{OL}$ (1)

## Definition: semantics of $\mathcal{OL}$ (1)

$f(\vec{n})$  or  $P(\vec{n})$  are **primitive** iff all  $n_i$  are standard names.  
A term or a formula is **ground** iff it mentions no variable.

A **world**  $w$  is a function that maps

- primitive functions  $f(\vec{n})$  to standard names
- primitive atomic formulas  $P(\vec{n})$  to  $\{0, 1\}$

The **denotation** of a ground term w.r.t.  $w$  is defined as

- $w(n) \stackrel{\text{def}}{=} n$  for every standard name  $n$
- $w(f(t_1, \dots, t_j)) \stackrel{\text{def}}{=} w[f(w(t_1), \dots, w(t_j))]$

E.g., if Frank is Mia's father and Mia is Jane's mother:

$$\begin{aligned} & w(\text{fatherOf}(\text{motherOf}(\text{Jane}))) \\ &= w[\text{fatherOf}(w[\text{motherOf}(\text{Jane})])] \\ &= w[\text{fatherOf}(\text{Mia})] \\ &= \text{Frank}. \end{aligned}$$



## The Semantics of $\mathcal{OL}$ (2)

### Definition: semantics of $\mathcal{OL}$

An **epistemic state**  $e$  is a set of worlds.

- $e, w \models P(t_1, \dots, t_j) \iff w[P(w(t_1), \dots, w(t_j))] = 1$
- $e, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$

## The Semantics of $\mathcal{OL}$ (2)

### Definition: semantics of $\mathcal{OL}$

An **epistemic state**  $e$  is a set of worlds.

- $e, w \models P(t_1, \dots, t_j) \iff w[P(w(t_1), \dots, w(t_j))] = 1$
- $e, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$
- $e, w \models \neg\alpha \iff e, w \not\models \alpha$
- $e, w \models (\alpha \vee \beta) \iff e, w \models \alpha \text{ or } e, w \models \beta$

## The Semantics of $\mathcal{OL}$ (2)

### Definition: semantics of $\mathcal{OL}$

An **epistemic state**  $e$  is a set of worlds.

- $e, w \models P(t_1, \dots, t_j) \iff w[P(w(t_1), \dots, w(t_j))] = 1$
- $e, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$
- $e, w \models \neg\alpha \iff e, w \not\models \alpha$
- $e, w \models (\alpha \vee \beta) \iff e, w \models \alpha$  or  $e, w \models \beta$
- $e, w \models \exists x\alpha \iff e, w \models \alpha_n^x$  for some standard name  $n$

## The Semantics of $\mathcal{OL}$ (2)

### Definition: semantics of $\mathcal{OL}$

An **epistemic state**  $e$  is a set of worlds.

- $e, w \models P(t_1, \dots, t_j) \iff w[P(w(t_1), \dots, w(t_j))] = 1$
- $e, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$
- $e, w \models \neg\alpha \iff e, w \not\models \alpha$
- $e, w \models (\alpha \vee \beta) \iff e, w \models \alpha$  or  $e, w \models \beta$
- $e, w \models \exists x\alpha \iff e, w \models \alpha_n^x$  for some standard name  $n$
- $e, w \models \mathbf{K}\alpha \iff$  for all worlds  $w'$ ,  $w' \in e \Rightarrow e, w' \models \alpha$
- $e, w \models \mathbf{O}\alpha \iff$  for all worlds  $w'$ ,  $w' \in e \Leftrightarrow e, w' \models \alpha$

# Knowing That vs Knowing What

- $\mathbf{K}\exists x \text{Secret}(x)$  I know *that* some  $x$  is a secret
- $\exists x \mathbf{K}\text{Secret}(x)$  I know *which*  $x$  is a secret
  
- $\mathbf{K}\exists x \text{fatherOf}(\text{Sally}) = x$  I know *that* Sally has a father
- $\exists x \mathbf{K}\text{fatherOf}(\text{Sally}) = x$  I know *who* Sally's father is
  
- $\mathbf{K}\exists x \alpha = \textit{de dicto}$  knowledge
- $\exists x \mathbf{K}\alpha = \textit{de re}$  knowledge

## Theorem: quantifying-in

$$\models \forall x \mathbf{K}\alpha \leftrightarrow \mathbf{K}\forall x \alpha$$

$$\models \exists x \mathbf{K}\alpha \rightarrow \mathbf{K}\exists x \alpha$$

$$\not\models \mathbf{K}\exists x \alpha \rightarrow \exists x \mathbf{K}\alpha$$

## Some Properties Inherited From $\mathcal{OL}_{PL}$

### Definition: subjective, objective

If  $\phi$  mentions no fun/pred inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\phi$  is **objective**.

If  $\sigma$  mentions fun/pred only inside  $\mathbf{K}$  or  $\mathbf{O}$ , we say  $\sigma$  is **subjective**.

### Theorem: logical omniscience

If  $\models \alpha \rightarrow \beta$ , then  $\models \mathbf{K}\alpha \rightarrow \mathbf{K}\beta$ .

If  $\models \alpha$ , then  $\models \mathbf{K}\alpha$ .

### Theorem: unique-model property

Let  $\phi$  be objective. Then there is a unique  $e$  such that  $e \models \mathbf{O}\phi$ .

### Theorem: positive and negative introspection

Positive introspection:  $\models \mathbf{K}\alpha \rightarrow \mathbf{K}\mathbf{K}\alpha$

Negative introspection:  $\models \neg\mathbf{K}\alpha \rightarrow \mathbf{K}\neg\mathbf{K}\alpha$

## Example

$$\begin{aligned} e, w \models \exists x \alpha &\iff e, w \models \alpha_n^x \text{ for some standard name } n \\ e, w \models \mathbf{K}\alpha &\iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha \\ e, w \models \mathbf{O}\alpha &\iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha \end{aligned}$$

Let  $\text{KB} \stackrel{\text{def}}{=} \exists x (x \neq \#1 \wedge P(x))$

## Example

$$\begin{aligned} e, w \models \exists x \alpha &\iff e, w \models \alpha_n^x \text{ for some standard name } n \\ e, w \models \mathbf{K}\alpha &\iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha \\ e, w \models \mathbf{O}\alpha &\iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha \end{aligned}$$

Let  $\text{KB} \stackrel{\text{def}}{=} \exists x (x \neq \#1 \wedge P(x))$

■  $e \models \mathbf{O}\text{KB}$



## Example

$$\begin{aligned} e, w \models \exists x \alpha &\iff e, w \models \alpha_n^x \text{ for some standard name } n \\ e, w \models \mathbf{K}\alpha &\iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha \\ e, w \models \mathbf{O}\alpha &\iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha \end{aligned}$$

Let  $\text{KB} \stackrel{\text{def}}{=} \exists x(x \neq \#1 \wedge P(x))$

■  $e \models \mathbf{O}\text{KB}$

$$\iff w \in e \Leftrightarrow w \models \exists x(x \neq \#1 \wedge P(x))$$

$$\iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \dots\}$$

## Example

$$\begin{aligned} e, w \models \exists x \alpha &\iff e, w \models \alpha_n^x \text{ for some standard name } n \\ e, w \models \mathbf{K}\alpha &\iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha \\ e, w \models \mathbf{O}\alpha &\iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha \end{aligned}$$

Let  $\text{KB} \stackrel{\text{def}}{=} \exists x (x \neq \#1 \wedge P(x))$

■  $e \models \mathbf{O}\text{KB}$

$$\iff w \in e \Leftrightarrow w \models \exists x (x \neq \#1 \wedge P(x))$$

$$\iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \dots\}$$

■  $e \models \mathbf{K}\exists x (P(x) \wedge \neg \mathbf{K}P(x))$

## Example

$$\begin{aligned} e, w \models \exists x \alpha &\iff e, w \models \alpha_n^x \text{ for some standard name } n \\ e, w \models \mathbf{K}\alpha &\iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha \\ e, w \models \mathbf{O}\alpha &\iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha \end{aligned}$$

Let  $\text{KB} \stackrel{\text{def}}{=} \exists x (x \neq \#1 \wedge P(x))$

■  $e \models \mathbf{O}\text{KB}$

$$\iff w \in e \Leftrightarrow w \models \exists x (x \neq \#1 \wedge P(x))$$

$$\iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \dots\}$$

■  $e \models \mathbf{K}\exists x (P(x) \wedge \neg \mathbf{K}P(x))$

$$\iff \text{for all } w, w \in e \Rightarrow \text{for some } n, e, w \models P(n) \wedge \neg \mathbf{K}P(n)$$

## Example

$$\begin{aligned} e, w \models \exists x \alpha &\iff e, w \models \alpha_n^x \text{ for some standard name } n \\ e, w \models \mathbf{K}\alpha &\iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha \\ e, w \models \mathbf{O}\alpha &\iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha \end{aligned}$$

Let  $\text{KB} \stackrel{\text{def}}{=} \exists x (x \neq \#1 \wedge P(x))$

### ■ $e \models \mathbf{O}\text{KB}$

$$\iff w \in e \Leftrightarrow w \models \exists x (x \neq \#1 \wedge P(x))$$

$$\iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \dots\}$$

### ■ $e \models \mathbf{K}\exists x (P(x) \wedge \neg \mathbf{K}P(x))$

$$\iff \text{for all } w, w \in e \Rightarrow \text{for some } n, e, w \models P(n) \wedge \neg \mathbf{K}P(n)$$

$$\iff \text{for all } w, w \in e \Rightarrow \text{for some } n, e, w \models P(n) \text{ and } e, w \models \neg \mathbf{K}P(n)$$

## Example

$$\begin{aligned} e, w \models \exists x \alpha &\iff e, w \models \alpha_n^x \text{ for some standard name } n \\ e, w \models \mathbf{K}\alpha &\iff \text{for all worlds } w', w \in e \Rightarrow e, w' \models \alpha \\ e, w \models \mathbf{O}\alpha &\iff \text{for all worlds } w', w \in e \Leftrightarrow e, w' \models \alpha \end{aligned}$$

Let  $\text{KB} \stackrel{\text{def}}{=} \exists x (x \neq \#1 \wedge P(x))$

### ■ $e \models \mathbf{OKB}$

$$\iff w \in e \Leftrightarrow w \models \exists x (x \neq \#1 \wedge P(x))$$

$$\iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \dots\}$$

### ■ $e \models \mathbf{K}\exists x (P(x) \wedge \neg \mathbf{K}P(x))$

$$\iff \text{for all } w, w \in e \Rightarrow \text{for some } n, e, w \models P(n) \wedge \neg \mathbf{K}P(n)$$

$$\iff \text{for all } w, w \in e \Rightarrow \text{for some } n, e, w \models P(n) \text{ and } e, w \models \neg \mathbf{K}P(n)$$

$$\iff \text{for all } w, w \in e \Rightarrow \text{for some } n, e, w \models P(n) \text{ and} \\ \text{for some } w', w' \in e \text{ and } e, w' \not\models P(n)$$

$$\iff \text{for all } w, w \in e \Rightarrow \text{for some } n, w[P(n)] = 1 \text{ and} \\ \text{for some } w', w' \in e \text{ and } w'[P(n)] \neq 1$$

# Comparison with Tarski Semantics

## ■ Traditional FOL semantics

- ▶ Interpretation  $\langle D, \Phi \rangle$  plus variable mapping  $\mu$
- ▶  $\langle D, \Phi \rangle, \mu \models P(t_1, \dots, t_j) \iff \langle d_1, \dots, d_j \rangle \in \Phi(P)$   
where  $d_i = \langle D, \Phi \rangle, \mu \parallel t_i \parallel$
- ▶  $\langle D, \Phi \rangle, \mu \models \exists x \alpha \iff \langle D, \Phi \rangle, \mu_d^x \models \alpha$  for some  $d \in D$
- ▶ Purpose: reason about mathematics
- ▶ Disadvantage: cumbersome to work with

## ■ Our semantics

- ▶ World maps primitive functions to names, predicates to  $\{0, 1\}$
- ▶  $w \models P(t_1, \dots, t_j) \iff w[P(n_1, \dots, n_j)] = 1$  where  $n_i = w(t_i)$
- ▶  $w \models \exists x \alpha \iff w \models \alpha_n^x$  for some standard name  $n$
- ▶ Purpose: reason about knowledge
- ▶ Disadvantage: domain is always countably infinite
  - ▶  $\forall x (x = t_1 \vee \dots \vee x = t_j)$  asserts finite domain in classical FOL
  - ▶  $\forall x (x = t_1 \vee \dots \vee x = t_j)$  is unsatisfiable in  $\mathcal{OL}$
  - ▶ but can be simulated with predicate:  $\forall x (P(x) \leftrightarrow (x = t_1 \vee \dots \vee x = t_j))$
  - ▶ classical FOL cannot distinguish countably infinite from uncountably infinite domains anyway

## Representation Theorem (1)

$$\mathbf{OKB} \models \exists x \mathbf{K}P(x)?$$

How can we represent the known instances of an objective formula?

- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \wedge P(\#2))$       #1, #2 are known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \vee P(\#2))$       no known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x P(x)$       all names are known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x (x \neq \#1 \rightarrow P(x))$       #2, #3, ... are known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} (Q(\#1) \wedge \forall x (Q(x) \rightarrow P(x)))$       #1 is known  $P$ -instance

## Representation Theorem (1)

$$\mathbf{OKB} \models \exists x \mathbf{K}P(x)?$$

How can we represent the known instances of an objective formula?

- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \wedge P(\#2))$       #1, #2 are known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \vee P(\#2))$       no known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x P(x)$       all names are known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x (x \neq \#1 \rightarrow P(x))$       #2, #3, ... are known  $P$ -instances
- $\mathbf{KB} \stackrel{\text{def}}{=} (Q(\#1) \wedge \forall x (Q(x) \rightarrow P(x)))$       #1 is known  $P$ -instance

Let  $n_1, \dots, n_j$  be names in  $\mathbf{KB}$  and let  $n'$  be a new one.

$$\begin{aligned} \text{RES}[\mathbf{KB}, P(x)] &\stackrel{\text{def}}{=} (x = n_1 \wedge \text{"KB} \models P(n_1)\text{"}) \vee \\ &\quad \dots \\ &\quad (x = n_j \wedge \text{"KB} \models P(n_j)\text{"}) \vee \\ &\quad (x \neq n_1 \wedge \dots \wedge x \neq n_j \wedge \text{"KB} \models P(n')\text{"}) \end{aligned}$$



## Representation Theorem (1)

$$\mathbf{OKB} \models \exists x \mathbf{K}P(x)?$$

How can we represent the known instances of an objective formula?

- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \wedge P(\#2))$   $x = \#1 \vee x = \#2$
- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \vee P(\#2))$  FALSE
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x P(x)$  TRUE
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x (x \neq \#1 \rightarrow P(x))$   $x \neq \#1$
- $\mathbf{KB} \stackrel{\text{def}}{=} (Q(\#1) \wedge \forall x (Q(x) \rightarrow P(x)))$   $x = \#1$

Let  $n_1, \dots, n_j$  be names in  $\mathbf{KB}$  and let  $n'$  be a new one.

$$\begin{aligned} \text{RES}[\mathbf{KB}, P(x)] &\stackrel{\text{def}}{=} (x = n_1 \wedge \text{“}\mathbf{KB} \models P(n_1)\text{”}) \vee \\ &\quad \dots \\ &\quad (x = n_j \wedge \text{“}\mathbf{KB} \models P(n_j)\text{”}) \vee \\ &\quad (x \neq n_1 \wedge \dots \wedge x \neq n_j \wedge \text{“}\mathbf{KB} \models P(n')\text{”}) \end{aligned}$$

## Representation Theorem (1)

$$\mathbf{OKB} \models \exists x \mathbf{KP}(x)?$$

How can we represent the known instances of an objective formula?

- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \wedge P(\#2))$   $x = \#1 \vee x = \#2$
- $\mathbf{KB} \stackrel{\text{def}}{=} (P(\#1) \vee P(\#2))$  FALSE
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x P(x)$  TRUE
- $\mathbf{KB} \stackrel{\text{def}}{=} \forall x (x \neq \#1 \rightarrow P(x))$   $x \neq \#1$
- $\mathbf{KB} \stackrel{\text{def}}{=} (Q(\#1) \wedge \forall x (Q(x) \rightarrow P(x)))$   $x = \#1$

Let  $n_1, \dots, n_j$  be names in  $\mathbf{KB}$  and let  $n'$  be a new one.

$$\begin{aligned} \text{RES}[\mathbf{KB}, P(x)] \stackrel{\text{def}}{=} & (x = n_1 \wedge \text{"KB} \models P(n_1)\text{"}) \vee \\ & \dots \\ & (x = n_j \wedge \text{"KB} \models P(n_j)\text{"}) \vee \\ & (x \neq n_1 \wedge \dots \wedge x \neq n_j \wedge \text{"KB} \models P(n')\text{"}) \end{aligned}$$

## Representation Theorem (2)

### Definition: representation of known instances

If  $\phi$  has a free variable  $x$  and  $n_1, \dots, n_j$  are the names mentioned in KB,  $\phi$ , and  $n'$  is a new name:

$$\begin{aligned} \text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} & (x = n_1 \wedge \text{RES}[\text{KB}, \phi_{n_1}^x]) \vee \\ & \dots \\ & (x = n_j \wedge \text{RES}[\text{KB}, \phi_{n_j}^x]) \vee \\ & (x \neq n_1 \wedge \dots \wedge x \neq n_j \wedge \text{RES}[\text{KB}, \phi_{n'}^x]_{n'}) \end{aligned}$$

If  $\phi$  has no free variables:

$$\text{RES}[\text{KB}, \phi] \stackrel{\text{def}}{=} \begin{cases} \text{TRUE} & \text{if } \text{KB} \models \phi \\ \text{FALSE} & \text{otherwise} \end{cases}$$

## Representation Theorem (3)

$\|\cdot\|$  operator gets a rule for  $\exists x \alpha$ :

### Definition: representation operators

- $\|\phi\|_{\text{KB}} \stackrel{\text{def}}{=} \phi$  for objective  $\phi$
- $\|\neg\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \neg\|\alpha\|_{\text{KB}}$
- $\|(\alpha \vee \beta)\|_{\text{KB}} \stackrel{\text{def}}{=} (\|\alpha\|_{\text{KB}} \vee \|\beta\|_{\text{KB}})$
- $\|\exists x \alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \exists x \|\alpha\|_{\text{KB}}$
- $\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$

## Representation Theorem (3)

$\|\cdot\|$  operator gets a rule for  $\exists x \alpha$ :

### Definition: representation operators

- $\|\phi\|_{\text{KB}} \stackrel{\text{def}}{=} \phi$  for objective  $\phi$
- $\|\neg\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \neg\|\alpha\|_{\text{KB}}$
- $\|(\alpha \vee \beta)\|_{\text{KB}} \stackrel{\text{def}}{=} (\|\alpha\|_{\text{KB}} \vee \|\beta\|_{\text{KB}})$
- $\|\exists x \alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \exists x \|\alpha\|_{\text{KB}}$
- $\|\mathbf{K}\alpha\|_{\text{KB}} \stackrel{\text{def}}{=} \text{RES}[\text{KB}, \|\alpha\|_{\text{KB}}]$

### Theorem: representation theorem

$$\mathbf{OKB} \models \alpha \iff \models \|\alpha\|_{\text{KB}}.$$

# Overview of the Lecture

- A Logic of Knowledge – The Propositional Fragment
- A Logic of Knowledge – The First-Order Case
- **Extensions of the Logic of Knowledge**
  - ▶ Multiple agents
  - ▶ Probabilities
  - ▶ Conditional belief
  - ▶ Limited Belief (week 8)
  - ▶ Actions (week 9)

## Multi-Agent Belief

Mike does not know what is in the gift box, but he knows that Jane knows what is in there:

$$\mathbf{K}_{\text{Mike}} \exists x (\text{InBox}(x) \wedge \neg \mathbf{K}_{\text{Mike}} \text{InBox}(x) \wedge \mathbf{K}_{\text{Jane}} \text{InBox}(x))$$

Epistemic states get more complex: in every possible world, Mike considers a whole set of worlds to be possible from Jane's perspective.

## Probabilities

I believe that with probability .999, there is no bomb in the gift box:

$$\mathbf{B}(\neg\exists x(\text{InBox}(x) \wedge \text{Bomb}(x)) : 0.999)$$

An epistemic state is now probability distribution over possible worlds.



## Conditional Belief

I believe that if something is in the gift box, it's probably not a bomb:

$$\mathbf{B}(\exists x \text{InBox}(x) \Rightarrow \neg \text{Bomb}(x))$$

Epistemic state ranks possible worlds by plausibility and checks if the most-plausible worlds where  $\exists x \text{InBox}(x)$  is true also satisfy  $\text{Bomb}(x)$ .

- A knowledge base is now a collection of conditionals  
"if \_\_\_\_\_, then most likely \_\_\_\_\_"
- What sort of ranking should these conditionals induce?