Reasoning about (Lack of) Knowledge

Christoph Schwering

UNSW Sydney

COMP4418, Week 7

- 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)



- 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 (1958)
 - Programs with Common Sense

- 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 (1958)
 - Programs with Common Sense
 - Improve program behaviour by making statements to it



- John McCarthy (1927–2011):
 - Stanford, MIT, Dartmouth
 - Turing Award
 - Invented Lisp (1958)
 - Invented Garbage Collection (1959)



- Proposed Advice Taker (1958)
 - Programs with Common Sense
 - Improve program behaviour by making statements to it
 - Program draws conclusions from its knowledge



- John McCarthy (1927–2011):
 - Stanford, MIT, Dartmouth
 - Turing Award
 - Invented Lisp (1958)
 - Invented Garbage Collection (1959)



- Proposed Advice Taker (1958)
 - 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



- John McCarthy (1927–2011):
 - Stanford, MIT, Dartmouth
 - Turing Award
 - Invented Lisp (1958)
 - Invented Garbage Collection (1959)



- Proposed Advice Taker (1958)
 - 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



- John McCarthy (1927–2011):
 - Stanford, MIT, Dartmouth
 - Turing Award
 - Invented Lisp (1958)
 - Invented Garbage Collection (1959)



- Proposed Advice Taker (1958)
 - 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 Al research, in particular what we'll be studying for the next three weeks



Observation: Non-knowledge is important

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

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.

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.

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.

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

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

- 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.,
 - $\blacktriangleright \ \alpha \text{ is known} \qquad \Longrightarrow \ KB \models \alpha$
 - $\triangleright \alpha$ is not known \implies KB $\not\models \alpha$
 - \implies KB is *all* the agent knows

- 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.,
 - $\blacktriangleright \ \alpha \text{ is known} \qquad \Longrightarrow \ KB \models \alpha$
 - $\triangleright \alpha$ is not known \implies 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 interences

- 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.,
 - $\blacktriangleright \ \alpha \text{ is known} \qquad \Longrightarrow \ KB \models \alpha$
 - $\triangleright \alpha$ is not known \implies 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 interences
- Usually: what is known \subsetneq what is true
 - Agent's knowledge is incomplete
 - Agent should be aware of that

- 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.,
 - $\blacktriangleright \ \alpha \text{ is known} \qquad \Longrightarrow \ KB \models \alpha$
 - $\triangleright \alpha$ is not known \implies 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 interences
- 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
 - $\blacktriangleright \ \alpha \in KB \implies \alpha \text{ is explicit knowledge (= database lookup)}$
 - ► KB $\models \alpha \implies \alpha$ is implicit knowledge (= logical inference)
 - ▶ Usually: explicit knowledge \subsetneq (implicit) knowledge

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

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

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

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

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

- 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.

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

- 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.

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

- 1. You don't know that r. KB $\not\models$
- 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.

$$\begin{array}{c} \mathrm{KB} \not\models r \\ \mathrm{KB} \not\models \neg r \end{array}$$

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

- 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.

$$KB \not\models r$$

$$KB \not\models \neg r$$

$$KB \models (p \lor q)$$

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

- 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.

$$\mathsf{KB} \not\models r$$

$$\mathsf{KB} \not\models \neg r$$

$$\mathsf{KB} \models (p \lor q)$$

$$\mathrm{KB} \not\models p$$

$$\mathrm{KB} \not\models q$$

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

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*.

$$\mathsf{KB} \not\models r$$

$$\mathsf{KB} \not\models \neg r$$

$$\mathsf{KB} \models (p \lor q)$$

$$\mathsf{KB} \not\models p$$

 $\text{KB} \not\models q$ 6. You know that *p* or *q*, but not which. $KB \models ???$

Problem: Classical logic cannot express 6 directly in one formula

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r) \land \neg k_p \land \ldots$

Then:

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

$$\text{KB} \models \neg k_r$$

$$\mathsf{KB} \models k_p_or_q$$

$$\mathsf{KB} \models \neg k_p$$

$$\mathrm{KB} \models \neg k_q$$

6. You know that *p* or *q*, but not which.

$$\mathsf{KB} \models k_p_or_q \land \neg k_p \land \neg k_p$$

Problem: Classical logic cannot express 6 directly in one formula Idea #1: Compile $(p \lor q \lor c)$ to new atoms $k_p, k_p_{-}or_q, \dots$ Does not scale.

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

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.

$$KB \models r = U$$

$$\mathsf{KB} \models \neg r = \mathsf{U}$$

$$\mathsf{KB} \models (p \lor q)$$

$$\mathsf{KB}\models p=\mathsf{U}$$

$$\mathsf{KB}\models q=\mathsf{U}$$

6. You know that p or q, but not which.

$$\mathsf{KB} \models (p \lor q) \land p = \mathsf{U} \land q = \mathsf{U}$$

 $\label{eq:problem: Classical logic cannot express 6 directly in one formula \\ \underline{\mathsf{Idea}~\#2} : \ensuremath{\mathsf{Three-valued}}\xspace \ensuremath{\mathsf{logic}}\xspace \{0,1,U\} \ensuremath{\not{\times}}\xspace \\ \ensuremath{\mathsf{How}}\xspace \ensuremath{\mathsf{would}}\xspace \ensuremath{\mathsf{U}}\xspace \lor U \ensuremath{\mathsf{behave}}\xspace \ensuremath{\mathsf{N}}\xspace \ensuremath{\mathsf{How}}\xspace \ensuremath{\mathsf{N}}\xspace \ensurem$

Suppose all you know is $(p \lor q \lor r) \land (p \lor q \lor \neg r)$

Then:

- 1. You don't know that r. **O**KB \models
- 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. **O**KB $\models \neg$ **K**q
- 6. You know that *p* or *q*, but not which.

 $\mathbf{O}\mathrm{KB} \models \mathbf{K}(p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q$

<u>Problem</u>: Classical logic cannot express 6 directly in one formula <u>Idea #3</u>: Add unary operators O and K to express knowledge \checkmark

$$\mathbf{O}$$
KB $\models \neg \mathbf{K}r$

$$\mathbf{O}$$
KB $\models \neg \mathbf{K} \neg r$

$$\mathbf{O}$$
KB \models $\mathbf{K}(p \lor q)$

$$\mathbf{O}$$
KB $\models \neg$ Kp

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

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

■ <i>p</i> , <i>q</i> , <i>r</i> ,	atomic propositions
$\square \neg \alpha$	"not a"
$ (\alpha \lor \beta) $	"α or β"
$ (\alpha \land \beta) $	" α and β "
$ (\alpha \to \beta) $	" α implies β "
$ (\alpha \leftrightarrow \beta) $	" α is equivalent to β "
■ Κα	"α is known"
Οα	"α is <i>all</i> that is known"

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

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

■ <i>p</i> , <i>q</i> , <i>r</i> ,	atomic propositions
$\square \neg \alpha$	"not α"
$ (\alpha \lor \beta) $	"α or β"
$ (\alpha \land \beta) \stackrel{\text{\tiny def}}{=} \neg (\neg \alpha \lor \neg \beta) $	" $lpha$ and eta "
$\blacksquare \ (\alpha \to \beta) \ \stackrel{\text{\tiny def}}{=} \ (\neg \alpha \lor \beta)$	" α implies β "
$\blacksquare \ (\alpha \leftrightarrow \beta) \ \stackrel{\scriptscriptstyle def}{=} \ (\alpha \to \beta) \land (\beta \to \alpha)$	" $lpha$ is equivalent to eta "
Κα	"α is known"
Οα	" α is <i>all</i> that is known"

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

A logical language is a *formal language* over an *alphabet* (here: $\{p, q, r, ..., (,), \neg, \lor, K, O\}$) 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 α , written $I \models \alpha$, or falsifies it, written $I \not\models \alpha$.

A logical language is a *formal language* over an *alphabet* (here: $\{p, q, r, ..., (,), \neg, \lor, K, O\}$) 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 α , written $I \models \alpha$, or falsifies it, written $I \not\models \alpha$.

A typical rule of a semantics is

 $I \models (\alpha \lor \beta)$ if and only if $I \models \alpha$ or $I \models \beta$.

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

A logical language is a *formal language* over an *alphabet* (here: $\{p, q, r, ..., (,), \neg, \lor, K, O\}$) 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 α , written $I \models \alpha$, or falsifies it, written $I \not\models \alpha$.

A typical rule of a semantics is

 $I \models (\alpha \lor \beta)$ if and only if $I \models \alpha$ or $I \models \beta$.

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

We will sometimes take the liberty to omit brackets to ease readability. For instance, we write $(p \lor q \lor r)$ instead of $((p \lor q) \lor r)$ or $(p \lor (q \lor r))$, implicitly assuming our semantics of \lor 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 α , then α is *valid* and we write $\models \alpha$.
- When all interpretations that satisfy some sentence Σ or set of sentences Σ also satisfy α, we say Σ entails α and write Σ ⊨ α.

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 α , then α is *valid* and we write $\models \alpha$.
- When all interpretations that satisfy some sentence Σ or set of sentences Σ also satisfy α, we say Σ entails α and write Σ ⊨ α.

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 ${\cal O\!L}_{PL}$

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

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

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

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

$$w \models P \iff w[P] = 1$$

$$w \models \neg \alpha \iff w \not\models \alpha$$

$$w \models (\alpha \lor \beta) \iff w \models \alpha \text{ or } w \models \beta$$

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

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

$$w \models P \iff w[P] = 1$$

$$w \models \neg \alpha \iff w \not\models \alpha$$

$$w \models (\alpha \lor \beta) \iff w \models \alpha \text{ or } w \models \beta$$

$$w \models \mathbf{K}\alpha \iff ???$$

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.

$$w \models P \iff w[P] = 1$$

$$w \models \neg \alpha \iff w \not\models \alpha$$

•
$$w \models (\alpha \lor \beta) \iff w \models \alpha \text{ or } w \models \beta$$

$$w \models \mathbf{K} \alpha \iff ???$$

$$w \models \mathbf{O}\alpha \iff ???$$

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.

$$\bullet, w \models P \iff w[P] = 1$$

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

$$\bullet e,w \models (\alpha \lor \beta) \iff e,w \models \alpha \text{ or } e,w \models \beta$$

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

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

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

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

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

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

•
$$e, w \models P \iff w[P] = 1$$

• $e, w \models \neg \alpha \iff e, w \not\models \alpha$
• $e, w \models (\alpha \lor \beta) \iff e, w \models \alpha \text{ or } e, w \models \beta$
• $e, w \models \mathbf{K}\alpha \iff \text{ for all worlds } w', \ w' \in e \Rightarrow e, w' \models \alpha$

">" stands for natural language expressions "only if".
 "\(\logma)" and "\(\logma)" stand for natural language expressions "if and only if".

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.

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

">" stands for natural language expressions "only if".
 "\(\logma)" and "\(\logma)" stand for natural language expressions "if and only if".

Abbreviations

Recall:

$$\begin{array}{l} \bullet \quad (\alpha \land \beta) \stackrel{\text{def}}{=} \neg (\neg \alpha \lor \neg \beta) \\ \bullet \quad (\alpha \to \beta) \stackrel{\text{def}}{=} (\neg \alpha \lor \beta) \\ \bullet \quad (\alpha \leftrightarrow \beta) \stackrel{\text{def}}{=} (\alpha \to \beta) \land (\beta \to \alpha) \end{array}$$

 \wedge should be "and"

 \rightarrow should be "only if"

 \leftrightarrow should be "if and only if".

Lemma: abbreviations

•
$$e, w \models \alpha \land \beta \iff e, w \models \alpha$$
 and $e, w \models \beta$

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

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

Proof on paper

Some Lemmas

Definition: objective, subjective

If ϕ mentions no atoms inside **K** or **O**, we say ϕ is **objective**. If σ mentions atoms only inside **K** or **O**, we say σ is **subjective**.

•
$$((p \lor q) \land p \land q)$$
 is objective

K
$$((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$$
 is subjective

Some Lemmas

Definition: objective, subjective

If ϕ mentions no atoms inside **K** or **O**, we say ϕ is **objective**. If σ mentions atoms only inside **K** or **O**, we say σ is **subjective**.

•
$$((p \lor q) \land p \land q)$$
 is objective

K
$$((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$$
 is subjective

Lemma: objective, subjective

Let ϕ be objective. Then $e, w \models \phi \iff e', w \models \phi$. Let σ be subjective. Then $e, w \models \sigma \iff e, w' \models \sigma$.

Some Lemmas

Definition: objective, subjective

If ϕ mentions no atoms inside **K** or **O**, we say ϕ is **objective**. If σ mentions atoms only inside **K** or **O**, we say σ is **subjective**.

•
$$((p \lor q) \land p \land q)$$
 is objective

K
$$((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$$
 is subjective

Lemma: objective, subjective

Let ϕ be objective. Then $e, w \models \phi \iff e', w \models \phi$. Let σ be subjective. Then $e, w \models \sigma \iff e, w' \models \sigma$.

When ϕ is objective, " $w \models \phi$ " stands for "for every e, $e, w \models \phi$ ". When σ is subjective, " $e \models \sigma$ " stands for "for every w, $e, w \models \sigma$ ".

Proof on paper

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

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

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

• $w \in e$

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

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

•
$$w \in e$$

 $\iff w \models (p \lor q \lor r) \land (p \lor q \lor \neg r)$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

$$w \in e \Leftrightarrow w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \Leftrightarrow w \models (p \lor q \lor r) \text{ and } w \models (p \lor q \lor \neg r)$$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

•
$$w \in e$$

 $\iff w \models (p \lor q \lor r) \land (p \lor q \lor \neg r)$
 $\iff w \models (p \lor q \lor r) \text{ and } w \models (p \lor q \lor \neg r)$
 $\iff w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 1, \text{ and}$
 $w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 0$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

$$w \in e \Leftrightarrow w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \Leftrightarrow w \models (p \lor q \lor r) \text{ and } w \models (p \lor q \lor \neg r) \Leftrightarrow w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 1, \text{ and } w[p] = 1 \text{ or } w[q] = 1 \text{ or } w[r] = 0 \\ \Leftrightarrow w[p] = 1 \text{ or } w[q] = 1$$

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

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

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

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

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

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

•
$$e \models \mathbf{K}(p \lor q)$$
 ?
 \iff for all $w, w \in e \Rightarrow w \models (p \lor q)$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

•
$$w \in e \iff w[p] = 1$$
 or $w[q] = 1$

$$e \models \mathbf{K}(p \lor q) ?$$

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

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

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

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

■
$$e \models \mathbf{K}(p \lor q)$$
 ?
 \iff for all $w, w \in e \Rightarrow w \models (p \lor q)$
 \iff for all $w, w \in e \Rightarrow w \models p$ or $w \models q$
 \iff for all $w, w \in e \Rightarrow w[p] = 1$ or $w[q] = 1$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

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

$$e \models \mathbf{K}(p \lor q) ?$$

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

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

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

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

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

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

u $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$
u $e \models \mathbf{K}(p \lor q) \checkmark$
u $e \models \neg \mathbf{K}p$?

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

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

u $e \in e \iff w[p] = 1 \text{ or } w[q] = 1$
u $e \models \mathbf{K}(p \lor q) \checkmark$
u $e \models \neg \mathbf{K}p$?
 $\iff e \nvDash \mathbf{K}p$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

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

•
$$e \models \mathbf{K}(p \lor q)$$

$$e \models \neg \mathbf{K}p \quad ?$$

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

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

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

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

•
$$e \models \mathbf{K}(p \lor q)$$

$$\bullet = \neg \mathbf{K}p$$
 ?

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

- \iff for some $w, w \in e$ and $w \not\models p$
- \iff for some $w, w \in e$ and $w[p] \neq 1$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

•
$$w \in e \iff w[p] = 1$$
 or $w[q] = 1$

•
$$e \models \mathbf{K}(p \lor q)$$

$$\bullet = \neg \mathbf{K}p$$
 ?

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

- \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 $w[p] \neq 1$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

•
$$w \in e \iff w[p] = 1$$
 or $w[q] = 1$

•
$$e \models \mathbf{K}(p \lor q)$$

•
$$e \models \neg \mathbf{K}p$$
 ?

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

$$\iff$$
 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 $w[p]
eq 1$

$$\iff$$
 for some $w,w[p]
eq 1$ and $w[q]=1$ 🖌

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

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

u $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$
u $e \models \mathbf{K}(p \lor q) \checkmark$
u $e \models \neg \mathbf{K}p \checkmark$

 $\bullet \models \mathbf{K}(p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q ?$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

•
$$w \in e \iff w[p] = 1$$
 or $w[q] = 1$

•
$$e \models \mathbf{K}(p \lor q)$$

•
$$e \models \neg \mathbf{K}p$$

$$e \models \mathbf{K}(p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q ?$$

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

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

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

u $w \in e \iff w[p] = 1 \text{ or } w[q] = 1$
u $e \models \mathbf{K}(p \lor q) \checkmark$

•
$$e \models \neg \mathbf{K}p$$

$$\bullet \models \mathbf{K}(p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q \checkmark$$

$$\bullet \models \mathbf{O}((p \lor q \lor r) \land (p \lor q \lor \neg r)) \quad ?$$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

•
$$w \in e \iff w[p] = 1$$
 or $w[q] = 1$

•
$$e \models \mathbf{K}(p \lor q)$$

•
$$e \models \neg \mathbf{K}p$$

$$\bullet \models \mathbf{K}(p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q \quad \checkmark$$

$$e \models \mathbf{O}((p \lor q \lor r) \land (p \lor q \lor \neg r)) ?$$

$$\iff \text{ for all } w, w \in e \Leftrightarrow w \models ((p \lor q \lor r) \land (p \lor q \lor \neg r)) \checkmark$$

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

Let
$$e \stackrel{\text{\tiny def}}{=} \{ w \mid w \models (p \lor q \lor r) \land (p \lor q \lor \neg r) \}$$

•
$$w \in e \iff w[p] = 1 ext{ or } w[q] = 1$$

•
$$e \models \mathbf{K}(p \lor q)$$

•
$$e \models \neg \mathbf{K}p$$

$$\bullet \models \mathbf{K}(p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q \quad \checkmark$$

$$\bullet \models \mathbf{O}((p \lor q \lor r) \land (p \lor q \lor \neg r)) \quad \checkmark$$

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

 $\begin{array}{l} \mathsf{If} \models \alpha \rightarrow \beta, \mathsf{then} \models \mathbf{K} \alpha \rightarrow \mathbf{K} \beta. \\ \mathsf{In particular: If} \models \alpha, \mathsf{then} \models \mathbf{K} \alpha. \end{array}$

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

```
\begin{array}{l} \mathsf{If} \models \alpha \rightarrow \beta, \mathsf{then} \models \mathbf{K} \alpha \rightarrow \mathbf{K} \beta. \\ \mathsf{In particular: If} \models \alpha, \mathsf{then} \models \mathbf{K} \alpha. \end{array}
```

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 ϕ 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: $\models K\alpha \rightarrow KK\alpha$ Negative introspection: $\models \neg K\alpha \rightarrow K \neg K\alpha$

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

Can we solve $\mathbf{O}KB \models \alpha$ with ordinary, propositional reasoning? That is, can we eliminate \mathbf{K} and \mathbf{O} ?

Then we could use standard reasoning system.

Can we solve $\mathbf{O}\mathbf{KB} \models \alpha$ with ordinary, propositional reasoning? That is, can we eliminate \mathbf{K} and \mathbf{O} ? Then we could use standard reasoning system.

Theorem

Let KB, ϕ be objective. Then **O**KB \models **K** $\phi \iff$ KB $\models \phi$.

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

Can we solve $\mathbf{O}\mathbf{KB} \models \alpha$ with ordinary, propositional reasoning? That is, can we eliminate \mathbf{K} and \mathbf{O} ? Then we could use standard reasoning system.

Theorem

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

Idea: replace nested $\mathbf{K}\phi$ with true if $KB \models \phi$, otherwise with false.

$$\underbrace{\mathsf{Ex.}}_{\mathsf{Ex.}} : \mathsf{Let} \, \mathsf{KB} \stackrel{\text{\tiny def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r). \\ \mathbf{O} \mathsf{KB} \models \mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)?$$

Can we solve $\mathbf{O}KB \models \alpha$ with ordinary, propositional reasoning? That is, can we eliminate **K** and **O**? Then we could use standard reasoning system.

Theorem

Let KB, ϕ be objective. Then **O**KB \models **K** $\phi \iff$ KB $\models \phi$.

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

$$\underbrace{\mathsf{Ex.:}}_{\mathsf{KB} \models \mathsf{K}} \mathsf{Let} \mathsf{KB} \stackrel{\text{def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r).$$
$$\mathbf{O}\mathsf{KB} \models \mathsf{K} ((p \lor q) \land \neg \underbrace{\mathsf{Kp}}_{\mathsf{KB} \models p?} \land \neg \underbrace{\mathsf{Kq}}_{\mathsf{KB} \models q?} ?$$

Can we solve $\mathbf{O}\mathbf{KB} \models \alpha$ with ordinary, propositional reasoning? That is, can we eliminate \mathbf{K} and \mathbf{O} ? Then we could use standard reasoning system.

Theorem

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

Idea: replace nested $\mathbf{K}\phi$ with true if $KB \models \phi$, otherwise with false.

$$\underbrace{\mathsf{Ex.}}_{\mathsf{Ex.}} : \mathsf{Let} \, \mathsf{KB} \stackrel{\text{\tiny def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r). \\ \mathbf{O} \mathsf{KB} \models \mathbf{K} ((p \lor q) \land \neg \mathsf{False} \land \neg \mathsf{False})?$$

Can we solve $\mathbf{O}KB \models \alpha$ with ordinary, propositional reasoning? That is, can we eliminate **K** and **O**? Then we could use standard reasoning system.

Theorem

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

Idea: replace nested $\mathbf{K}\phi$ with true if $KB \models \phi$, otherwise with false.

Ex.: Let KB
$$\stackrel{\text{def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r).$$

KB $\models ((p \lor q) \land \neg \mathsf{FALSE} \land \neg \mathsf{FALSE})? \checkmark$

Can we solve $\mathbf{O}\mathbf{KB} \models \alpha$ with ordinary, propositional reasoning? That is, can we eliminate \mathbf{K} and \mathbf{O} ? Then we could use standard reasoning system.

Theorem

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

Idea: replace nested $\mathbf{K}\phi$ with true if $KB \models \phi$, otherwise with false.

$$\begin{array}{l} \underline{\mathsf{Ex.}}: \mathsf{Let} \, \mathsf{KB} \stackrel{\text{\tiny def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r). \\ \mathrm{KB} \models \big((p \lor q) \land \neg \mathsf{FALSE} \land \neg \mathsf{FALSE} \big)? \quad \checkmark \end{array}$$

Next slide formalises this idea.

<u>Sneak preview</u>: It'll become more difficult in the first-order case: What would you replace $\mathbf{K}Q(x)$ with in $\mathbf{K} \exists x (P(x) \land \neg \mathbf{K}Q(x))$? We'll see later.

Definition: representation operators

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

18/36

Definition: representation operators

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

Definition: representation operators

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

Theorem: representation theorem

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

$$\begin{split} \|\mathbf{K}\boldsymbol{\alpha}\|_{KB} &\stackrel{\text{def}}{=} \mathsf{RES}[KB, \|\boldsymbol{\alpha}\|_{KB}] \\ \mathsf{RES}[KB, \varphi] &\stackrel{\text{def}}{=} "KB \models \varphi?" \end{split}$$

Let
$$\mathrm{KB} \stackrel{\text{def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r).$$

OKB $\models \mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$?

$$\begin{split} \|\mathbf{K}\boldsymbol{\alpha}\|_{KB} &\stackrel{\text{\tiny def}}{=} \ \mathsf{RES}[KB, \|\boldsymbol{\alpha}\|_{KB}] \\ \mathsf{RES}[KB, \varphi] &\stackrel{\text{\tiny def}}{=} \ "KB \models \varphi?" \end{split}$$

Let
$$\mathrm{KB} \stackrel{\mathrm{def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r).$$

 $\mathbf{O}\mathrm{KB} \models \mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$
 $\iff \models \|\mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\mathrm{KB}}$

$$\begin{split} \|\mathbf{K}\boldsymbol{\alpha}\|_{KB} &\stackrel{\text{\tiny def}}{=} \ \mathsf{RES}[KB, \|\boldsymbol{\alpha}\|_{KB}] \\ \mathsf{RES}[KB, \varphi] &\stackrel{\text{\tiny def}}{=} \ "KB \models \varphi?" \end{split}$$

Let
$$\mathrm{KB} \stackrel{\text{def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r).$$

 $\mathbf{O}\mathrm{KB} \models \mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$
 $\iff \models \|\mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\mathrm{KB}}$
 $\iff \models \mathrm{RES}[\mathrm{KB}, \| ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\mathrm{KB}}]$

$$\begin{split} \|\mathbf{K}\boldsymbol{\alpha}\|_{KB} &\stackrel{\text{\tiny def}}{=} \mathsf{RES}[KB, \|\boldsymbol{\alpha}\|_{KB}] \\ \mathsf{RES}[KB, \varphi] &\stackrel{\text{\tiny def}}{=} "KB \models \varphi?" \end{split}$$

Let
$$\mathrm{KB} \stackrel{\text{def}}{=} (p \lor q \lor r) \land (p \lor q \lor \neg r).$$

 $\mathbf{O}\mathrm{KB} \models \mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$
 $\iff \models \|\mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\mathrm{KB}}$
 $\iff \models \mathrm{RES}[\mathrm{KB}, \| ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\mathrm{KB}}]$
 $\iff \models \mathrm{RES}[\mathrm{KB}, ((p \lor q) \land \neg \underbrace{\|\mathbf{K}p\|_{\mathrm{KB}}}_{\mathrm{KB}\models p?} \land \neg \underbrace{\|\mathbf{K}q\|_{\mathrm{KB}}}_{\mathrm{KB}\models q?})]$

$$\begin{split} \|\mathbf{K}\boldsymbol{\alpha}\|_{KB} &\stackrel{\text{\tiny def}}{=} \mathsf{RES}[KB, \|\boldsymbol{\alpha}\|_{KB}] \\ \mathsf{RES}[KB, \varphi] &\stackrel{\text{\tiny def}}{=} "KB \models \varphi?" \end{split}$$

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

OKB $\models \mathbf{K}((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$
 $\iff \models \|\mathbf{K}((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\text{KB}}$
 $\iff \models \text{RES}[\text{KB}, \|((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\text{KB}}]$
 $\iff \models \text{RES}[\text{KB}, ((p \lor q) \land \neg \underbrace{\|\mathbf{K}p\|_{\text{KB}} \land \neg \underbrace{\|\mathbf{K}q\|_{\text{KB}}}_{\text{KB}\models q?})]$
 $\iff \models \underbrace{\text{RES}[\text{KB}, ((p \lor q) \land \neg \text{FALSE} \land \neg \text{FALSE})]}_{\text{KB}\models (p \lor q) \land \neg \text{FALSE} \land \neg \text{FALSE}?}$

$$\begin{split} \|\mathbf{K}\boldsymbol{\alpha}\|_{KB} &\stackrel{\text{\tiny def}}{=} \mathsf{RES}[KB, \|\boldsymbol{\alpha}\|_{KB}] \\ \mathsf{RES}[KB, \varphi] &\stackrel{\text{\tiny def}}{=} "KB \models \varphi?" \end{split}$$

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

OKB $\models \mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)$
 $\iff \models \|\mathbf{K} ((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\text{KB}}$
 $\iff \models \text{RES}[\text{KB}, \|((p \lor q) \land \neg \mathbf{K}p \land \neg \mathbf{K}q)\|_{\text{KB}}]$
 $\iff \models \text{RES}[\text{KB}, ((p \lor q) \land \neg \underbrace{\|\mathbf{K}p\|_{\text{KB}} \land \neg \underbrace{\|\mathbf{K}q\|_{\text{KB}}}_{\text{KB}\models q?})]$
 $\iff \models \underbrace{\text{RES}[\text{KB}, ((p \lor q) \land \neg \text{FALSE} \land \neg \text{FALSE})]}_{\text{KB}\models (p \lor q) \land \neg \text{FALSE} \land \neg \text{FALSE})}$

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



 $K((\spadesuit \lor \heartsuit) \land \neg K \spadesuit \land \neg K \heartsuit)$



 $\mathbf{K} \exists x (\operatorname{InBox}(x) \land \neg \mathbf{K} \operatorname{InBox}(x))$



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



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

"all" or "some" \implies first-order quantification

The Language of ${\cal O\!L}$

Terms:

$ x, x', x_1, x_2, \ldots $	first-order variables
■ <i>#</i> 1, <i>#</i> 2, <i>#</i> 3,	standard names
$\blacksquare f(t_1,\ldots,t_j)$	functions

Formulas:

$\blacksquare P(t_1,\ldots,t_j)$			atomic formulas			
$\bullet t_1 = t_2$			equality expressions			
$\blacksquare \exists x \alpha$			"f	or son	ne <i>x</i> , α"	
$ \forall x \alpha $				"for	all x , α''	
$\blacksquare \neg \alpha (\alpha \lor \beta)$	$(\alpha \wedge \beta)$	$(\alpha ightarrow \beta)$	$(\alpha \leftrightarrow \beta)$	Kα	Οα	

The Language of ${\cal O\!L}$

Terms:

$ x, x', x_1, x_2, \dots $
■ #1,#2,#3,
$\blacksquare f(t_1,\ldots,t_j)$

Formulas:

 $P(t_1, \ldots, t_j)$ atomic formulas $t_1 = t_2$ equality expressions $\exists x \alpha$ "for some x, α " $\forall x \alpha \stackrel{\text{def}}{=} \neg \exists x \neg \alpha$ "for all x, α " $\neg \alpha \quad (\alpha \lor \beta) \quad (\alpha \land \beta) \quad (\alpha \to \beta) \quad (\alpha \leftrightarrow \beta) \quad \mathbf{K} \alpha \quad \mathbf{O} \alpha$

first-order variables

standard names

functions

Why Standard Names?

Consider in classical logic: fatherOf(Sally) = bestFriend(Jane) fatherOf(Sally) = bossOf(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: fatherOf(Sally) = Frank

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 ${\cal O\!L}$ (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

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

The Semantics of ${\cal O\!L}$ (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

•
$$w(n) \stackrel{\text{\tiny def}}{=} n$$
 for every standard name n

$$w(f(t_1,\ldots,t_j)) \stackrel{\text{def}}{=} w[f(w(t_1),\ldots,w(t_j))]$$

The Semantics of ${\cal O\!L}$ (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

•
$$w(n) \stackrel{\text{\tiny def}}{=} n$$
 for every standard name n

$$w(f(t_1,\ldots,t_j)) \stackrel{\text{def}}{=} w[f(w(t_1),\ldots,w(t_j))]$$

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

w(fatherOf(motherOf(Jane)))

- = w[fatherOf(w[motherOf(Jane)])]
- = w[fatherOf(Mia)]
- = Frank.

Definition: semantics of ${\cal O\!L}$

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

•
$$e, w \models P(t_1, \ldots, t_j) \iff w[P(w(t_1), \ldots, w(t_j)] = 1$$

$$\bullet, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$$

Definition: semantics of ${\cal O\!L}$

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

•
$$e, w \models P(t_1, \ldots, t_j) \iff w[P(w(t_1), \ldots, w(t_j)] = 1$$

$$\bullet, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$$

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

•
$$e,w \models (\alpha \lor \beta) \iff e,w \models \alpha \text{ or } e,w \models \beta$$

Definition: semantics of \mathcal{OL}

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

•
$$e, w \models P(t_1, \ldots, t_j) \iff w[P(w(t_1), \ldots, w(t_j)] = 1$$

$$\bullet, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$$

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

•
$$e,w \models (\alpha \lor \beta) \iff e,w \models \alpha \text{ or } e,w \models \beta$$

• $e,w \models \exists x \, \alpha \iff e,w \models \alpha_n^x$ for some standard name n

Definition: semantics of \mathcal{OL}

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

•
$$e, w \models P(t_1, \ldots, t_j) \iff w[P(w(t_1), \ldots, w(t_j)] = 1$$

$$\bullet, w \models t_1 = t_2 \iff w(t_1) = w(t_2)$$

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

•
$$e,w \models (\alpha \lor \beta) \iff e,w \models \alpha \text{ or } e,w \models \beta$$

• $e, w \models \exists x \alpha \iff e, w \models \alpha_n^x$ 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

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

Theorem: quantifying-in

- $\models \forall x \mathbf{K} \alpha \leftrightarrow \mathbf{K} \forall x \alpha$ $\models \exists x \mathbf{K} \alpha \rightarrow \mathbf{K} \exists x \alpha$
- $\not\models \mathbf{K} \exists x \, \alpha \to \exists x \, \mathbf{K} \, \alpha$

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

Definition: subjective, objective

If ϕ mentions no fun/pred inside **K** or **O**, we say ϕ is **objective**. If σ mentions fun/pred only inside **K** or **O**, we say σ is **subjective**.

Theorem: logical omniscience

 $\begin{array}{l} \mathsf{If} \models \alpha \rightarrow \beta, \mathsf{then} \models \mathbf{K} \alpha \rightarrow \mathbf{K} \beta. \\ \mathsf{If} \models \alpha, \mathsf{then} \models \mathbf{K} \alpha. \end{array}$

Theorem: unique-model property

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

Theorem: positive and negative introspection

Positive introspection: $\models K\alpha \rightarrow KK\alpha$ Negative introspection: $\models \neg K\alpha \rightarrow K \neg K\alpha$

$$e, w \models \exists x \alpha \iff e, w \models \alpha_n^x$$
 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 $w', w \in e \Leftrightarrow e, w' \models \alpha$

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

$$e, w \models \exists x \alpha \iff e, w \models \alpha_n^x$$
 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 $w', w \in e \Leftrightarrow e, w' \models \alpha$

Let KB
$$\stackrel{\text{\tiny def}}{=} \exists x \, (x \neq {}^{\#}1 \land P(x))$$

$\bullet e \models \mathbf{O}$ KB

 $e, w \models \exists x \alpha \iff e, w \models \alpha_n^x$ 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 $w', w \in e \Leftrightarrow e, w' \models \alpha$

Let KB $\stackrel{\text{\tiny def}}{=} \exists x \, (x \neq {}^{\#}1 \land P(x))$

•
$$e \models \mathbf{O}$$
KB

$$\iff w \in e \Leftrightarrow w \models \exists x (x \neq \#1 \land P(x)) \\ \iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \ldots\}$$

 $e, w \models \exists x \alpha \iff e, w \models \alpha_n^x$ 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 $w', w \in e \Leftrightarrow e, w' \models \alpha$

Let KB $\stackrel{\text{\tiny def}}{=} \exists x \, (x \neq {}^{\#}1 \land P(x))$

$$e \models \mathbf{OKB}$$

$$\iff w \in e \Leftrightarrow w \models \exists x (x \neq *1 \land P(x))$$

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

$$\bullet \models \mathbf{K} \exists x (P(x) \land \neg \mathbf{K} P(x))$$

 $e, w \models \exists x \alpha \iff e, w \models \alpha_n^x$ 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 $w', w \in e \Leftrightarrow e, w' \models \alpha$

Let KB $\stackrel{\text{\tiny def}}{=} \exists x \, (x \neq {}^{\#}1 \land P(x))$

$$e \models \mathbf{OKB}$$

$$\iff w \in e \Leftrightarrow w \models \exists x \, (x \neq *1 \land P(x))$$

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

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

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

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

Let KB $\stackrel{\text{\tiny def}}{=} \exists x \, (x \neq {}^{\#}1 \land P(x))$

■
$$e \models \mathbf{O}$$
KB
 $\iff w \in e \Leftrightarrow w \models \exists x (x \neq \#1 \land P(x))$
 $\iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \ldots\}$

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

 $e, w \models \exists x \alpha \iff e, w \models \alpha_n^x$ 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 $w', w \in e \Leftrightarrow e, w' \models \alpha$

Let KB $\stackrel{\text{\tiny def}}{=} \exists x \, (x \neq {}^{\#}1 \land P(x))$

• $e \models \mathbf{O}$ KB

$$\iff w \in e \Leftrightarrow w \models \exists x \, (x \neq \#1 \land P(x)) \\ \iff w \in e \Leftrightarrow w[P(n)] = 1 \text{ for some } n \in \{\#2, \#3, \ldots\}$$

■ $e \models \mathbf{K} \exists x (P(x) \land \neg \mathbf{K} P(x))$ \iff for all $w, w \in e \Rightarrow$ for some $n, e, w \models P(n) \land \neg \mathbf{K} P(n)$ \iff for all $w, w \in e \Rightarrow$ for some $n, e, w \models P(n)$ and $e, w \models \neg \mathbf{K} P(n)$ \iff for all $w, w \in e \Rightarrow$ for some $n, e, w \models P(n)$ and for some $w', w' \in e$ and $e, w' \not\models P(n)$ \iff for all $w, w \in e \Rightarrow$ for some n, w[P(n)] = 1 and for some $w', w' \in e$ and $w'[P(n)] \neq 1$

Comparison with Tarski Semantics

- Traditional FOL semantics
 - \blacktriangleright Interpretation $\langle D, \Phi
 angle$ plus variable mapping μ
 - $\begin{array}{l} \blacktriangleright \quad \langle D, \Phi \rangle, \mu \models P(t_1, \ldots, t_j) \iff \langle d_1, \ldots, d_j \rangle \in \Phi(P) \\ \text{where } d_i = \langle D, \Phi \rangle, \mu \| t_i \| \end{array}$
 - $\blacktriangleright \quad \langle D, \Phi \rangle, \mu \models \exists x \, \alpha \iff \langle D, \Phi \rangle, \mu_d^x \models \alpha \text{ for some } d \in D$
 - Purpose: reason about mathematics
 - Disadvantage: cumbersome to work with

Our semantics

- $\blacktriangleright\,$ World maps primitive functions to names, predicates to $\{0,1\}$
- ▶ $w \models P(t_1, ..., t_j) \iff w[P(n_1, ..., n_j)] = 1$ 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
 - $\forall x (x = t_1 \lor \ldots \lor x = t_j)$ asserts finite domain in classical FOL
 - $\forall x (x = t_1 \lor \ldots \lor x = t_j)$ is unsatisfiable in \mathcal{OL}
 - ▶ but can be simulated with predicate: $\forall x (P(x) \leftrightarrow (x = t_1 \lor ... \lor x = t_j))$
 - classical FOL cannot distinguish countably infinite from uncountably infinite domains anyway

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

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

KB $\stackrel{\text{def}}{=} (P(\#1) \land P(\#2))$ #1, #2 are known *P*-instancesKB $\stackrel{\text{def}}{=} (P(\#1) \lor P(\#2))$ no known *P*-instancesKB $\stackrel{\text{def}}{=} \forall x P(x)$ all names are known *P*-instancesKB $\stackrel{\text{def}}{=} \forall x (x \neq #1 \rightarrow P(x))$ #2, #3, . . . are known *P*-instancesKB $\stackrel{\text{def}}{=} (Q(\#1) \land \forall x (Q(x) \rightarrow P(x)))$ #1 is known *P*-instance

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

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

KB $\stackrel{\text{def}}{=} (P(#1) \land P(#2))$ #1, #2 are known *P*-instancesKB $\stackrel{\text{def}}{=} (P(#1) \lor P(#2))$ no known *P*-instancesKB $\stackrel{\text{def}}{=} \forall x P(x)$ all names are known *P*-instancesKB $\stackrel{\text{def}}{=} \forall x (x \neq #1 \rightarrow P(x))$ #2, #3, ... are known *P*-instancesKB $\stackrel{\text{def}}{=} (Q(#1) \land \forall x (Q(x) \rightarrow P(x)))$ #1 is known *P*-instance

Let n_1, \ldots, n_j be names in KB and let n' be a new one. RES[KB, P(x)] $\stackrel{\text{def}}{=} (x = n_1 \land \text{"KB} \models P(n_1)") \lor$

 $(x = n_j \wedge \text{"KB} \models P(n_j)\text{"}) \lor$ $(x \neq n_1 \wedge \ldots \wedge x \neq n_j \wedge \text{"KB} \models P(n')\text{"})$

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

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

• KB
$$\stackrel{\text{def}}{=} (P(\#1) \land P(\#2))$$
 $x = \#1 \lor x = \#2$

KB
$$\stackrel{\text{def}}{=} (P(#1) \lor P(#2))$$
 FALSE

KB
$$\stackrel{\text{def}}{=} \forall x P(x)$$
 TRUE

$$\mathbf{KB} \stackrel{\text{def}}{=} \forall x \, (x \neq \#1 \rightarrow P(x)) \qquad \qquad x \neq \#1$$

$$\blacksquare \text{ KB} \stackrel{\text{def}}{=} (Q(\#1) \land \forall x (Q(x) \to P(x))) \qquad \qquad x = \#1$$

Let n_1, \ldots, n_i be names in KB and let n' be a new one.

$$\mathsf{RES}[\mathsf{KB}, P(x)] \stackrel{\text{def}}{=} (x = n_1 \land \mathsf{"KB} \models P(n_1)") \lor$$
$$\dots$$
$$(x = n_j \land \mathsf{"KB} \models P(n_j)") \lor$$
$$(x \neq n_1 \land \dots \land x \neq n_j \land \mathsf{"KB} \models P(n')")$$

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

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

- KB $\stackrel{\text{def}}{=} (P(\#1) \land P(\#2))$ $x = \#1 \lor x = \#2$
- $\blacksquare \text{ KB} \stackrel{\text{def}}{=} (P(\#1) \lor P(\#2))$ False

KB
$$\stackrel{\text{def}}{=} \forall x P(x)$$
 TRUE

$$\blacksquare \text{ KB} \stackrel{\text{def}}{=} \forall x (x \neq \#1 \rightarrow P(x)) \qquad \qquad x \neq \#1$$

$$\blacksquare \text{ KB} \stackrel{\text{def}}{=} (Q(\#1) \land \forall x (Q(x) \to P(x))) \qquad \qquad x = \#1$$

Let n_1, \ldots, n_i be names in KB and let n' be a new one.

$$\mathsf{RES}[\mathsf{KB}, P(x)] \stackrel{\text{def}}{=} (x = n_1 \wedge \mathsf{KB} \models P(n_1)) \vee \dots \\ (x = n_j \wedge \mathsf{KB} \models P(n_j)) \vee \\ (x \neq n_1 \wedge \dots \wedge x \neq n_j \wedge \mathsf{KB} \models P(n'))$$

Definition: representation of known instances

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

$$\mathsf{RES}[\mathsf{KB}, \phi] \stackrel{\text{def}}{=} (x = n_1 \land \mathsf{RES}[\mathsf{KB}, \alpha_{n_1}^x]) \lor \dots \\ (x = n_j \land \mathsf{RES}[\mathsf{KB}, \alpha_{n_j}^x]) \lor \\ (x \neq n_1 \land \dots \land x \neq n_j \land \mathsf{RES}[\mathsf{KB}, \phi_{n'}^{x'}]_x^{n'})$$

If ϕ has no free variables:

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

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

Definition: representation operators

$$\| \phi \|_{KB} \stackrel{\text{\tiny def}}{=} \varphi \text{ for objective } \varphi$$

$$\|\neg \alpha\|_{\mathrm{KB}} \stackrel{\mathrm{\tiny def}}{=} \neg \|\alpha\|_{\mathrm{KB}}$$

$$\blacksquare \|(\alpha \lor \beta)\|_{KB} \stackrel{\text{def}}{=} (\|\alpha\|_{KB} \lor \|\beta\|_{KB})$$

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

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

Definition: representation operators

$$||\phi||_{KB} \stackrel{\text{\tiny def}}{=} \phi \text{ for objective } \phi$$

$$\|\neg \alpha\|_{\mathrm{KB}} \stackrel{\mathrm{\tiny def}}{=} \neg \|\alpha\|_{\mathrm{KB}}$$

$$\blacksquare \|(\alpha \lor \beta)\|_{KB} \stackrel{\text{def}}{=} (\|\alpha\|_{KB} \lor \|\beta\|_{KB})$$

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

Theorem: representation theorem

 $\mathbf{O}\mathrm{KB}\models\alpha\iff\models\|\alpha\|_{\mathrm{KB}}.$

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)

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

 $\mathbf{K}_{\mathrm{Mike}} \exists x \left(\mathrm{InBox}(x) \land \neg \mathbf{K}_{\mathrm{Mike}} \mathrm{InBox}(x) \land \mathbf{K}_{\mathrm{Jane}} \mathrm{InBox}(x) \right)$

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

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

$\mathbf{B}(\neg \exists x (\operatorname{InBox}(x) \land \operatorname{Bomb}(x)) : 0.999)$

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:

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

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

A knowledge base is now a collection of conditionals "if _____, then most likely _____"

What sort of ranking should these conditionals induce?