# AI, Consciousness, & Lambda (Λ)

## Selmer Bringsjord

Rensselaer AI & Reasoning (RAIR) Lab
Department of Cognitive Science
Department of Computer Science
Lally School of Management & Technology
Rensselaer Polytechnic Institute (RPI)
Troy, New York 12180 USA

IFLAI2
Oct 4 2021

# AI, Consciousness, & Lambda (Λ)

## Selmer Bringsjord

Rensselaer AI & Reasoning (RAIR) Lab
Department of Cognitive Science
Department of Computer Science
Lally School of Management & Technology
Rensselaer Polytechnic Institute (RPI)
Troy, New York 12180 USA

IFLAI2
Oct 4 2021

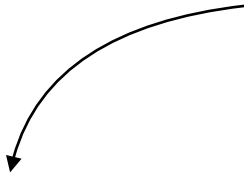# AI, Consciousness, & Lambda (Λ)

## Selmer Bringsjord

Rensselaer AI & Reasoning (RAIR) Lab
Department of Cognitive Science
Department of Computer Science
Lally School of Management & Technology
Rensselaer Polytechnic Institute (RPI)
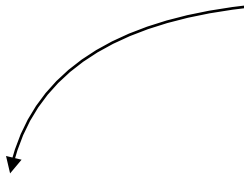Troy, New York 12180 USA
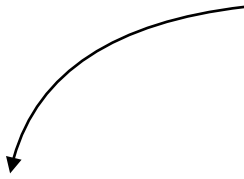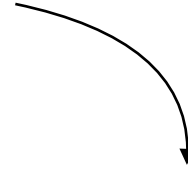
IFLAI2
Oct 4 2021

# "Consciousness"

"Consciousness"

"Consciousness"

'Access Consciousness'

"Consciousness"

'Access Consciousness'

"**Consciousness**"

'Access Consciousness'

Phenomenal Consciousness

Third-person formalization impossible.

**"Consciousness"**

'Access Consciousness'

Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

"Consciousness"

'Access Consciousness'

Phenomenal Consciousness

Third-person formalization impossible.

$$\Phi$$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

"Consciousness"

'Access Consciousness'

Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

Alas, biologically myopic & terrestrio-centric; and
'information' carries no information (!).

"Consciousness"

'Access Consciousness'

Phenomenal Consciousness

Third-person formalization impossible.

Φ

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

"Consciousness"

Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

"Consciousness"

Cognitive Consciousness

Phenomenal Consciousness

Third-person formalization impossible.

$$\Phi$$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

"**Consciousness**"

Cognitive Consciousness

Phenomenal Consciousness

Third-person formalization impossible.

$$\Phi$$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

# "Consciousness"

## Cognitive Consciousness

↓

## HLC-Consciousness

## Phenomenal Consciousness

Third-person formalization impossible.

$$\Phi$$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

# "Consciousness"

## Cognitive Consciousness

↓

## HLC-Consciousness

↓

## Phenomenal Consciousness

Third-person formalization impossible.

$$\Phi$$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

"Consciousness"

Cognitive Consciousness

HLC-Consciousness

HL**M**C-Consciousness

Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

# "Consciousness"

Cognitive Consciousness $\Lambda$ 

Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

HLC-Consciousness

HL**M**C-Consciousness

This can be viewed as a formal framework for measuring
the degree of "great computational intelligence" in an AI.

"**Consciousness**"

Cognitive Consciousness $\Lambda$ Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

Alas, biologically myopic & terrestrio-centric; and
'information' carries no information (!).

HLC-Consciousness

HL**M**C-Consciousness

This can be viewed as a formal framework for measuring the degree of "great computational intelligence" in an AI.

"Consciousness"

Cognitive Consciousness $\Lambda$

Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

HLC-Consciousness

HL**M**C-Consciousness

Happily, not bound by local biology; will cover aliens, God, characters of fiction, etc; and 'information' is information.

This can be viewed as a formal framework for measuring the degree of "great computational intelligence" in an AI.

"Consciousness"

Cognitive Consciousness $\Lambda$ Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

HLC-Consciousness

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

HL**M**C-Consciousness

Happily, not bound by local biology; will cover aliens, God, characters of fiction, etc; and 'information' is information.

This can be viewed as a formal framework for measuring the degree of "great computational intelligence" in an AI.

## "Consciousness"

Cognitive Consciousness $\Lambda$ Phenomenal Consciousness

Third-person formalization impossible.

$\Phi$

Alas, biologically myopic & terrestrio-centric; and 'information' carries no information (!).

HLC-Consciousness

HL**M**C-Consciousness

Happily, not bound by local biology; will cover aliens, God, characters of fiction, etc; and 'information' is information.

High- $\Lambda$ Machines are the ones DoD Needs to Worry About …

# Basic Idea, Intuitively Put

The level of (cognitive) intelligence/consciousness of an AI at a time is a list of tuples (= matrix) giving eg the size of logical depth of (at least) five measures for each cognitive operator (i.e. for **K**, **B**, **P**, …).

$$\langle [\![\mathbf{K}, 1]\!], [\![\mathbf{K}, 2]\!], \ldots, [\![\mathbf{K}, 5]\!], \ldots \rangle$$

# Basic Idea, Intuitively Put

The level of (cognitive) intelligence/consciousness of an AI at a time is a list of tuples (= matrix) giving eg the size of logical depth of (at least) five measures for each cognitive operator (i.e. for **K**, **B**, **P**, …).

$$\langle [\![\mathbf{K}, 1]\!], [\![\mathbf{K}, 2]\!], \ldots, [\![\mathbf{K}, 5]\!], \ldots \rangle$$

depth of knowledge

size of supporting proof/argument

depth of quantification within outermost knowledge operator

# Cogito Ergo Sum

```
{:name        "Cogito Ergo Sum"
 :description "A formaliztion of Descartes' Cogito Ergo Sum"
 :assumptions {

        S1 (Believes! I (forall [x] (or (Name x) (Thing x))))
        S2 (Believes! I (forall (x) (iff (Name x) (not (Thing x)))) )
        S3 (Believes! I (forall (x) (if (Thing x) (or (Real x) (Fictional x)))))
        S4 (Believes! I (forall (x) (if (Thing x) (iff (Real x) (not (Fictional x))))))
        ;;;
        A1 (Believes! I (forall (x) (if (Name x) (Thing (* x)))))
        A2 (Believes! I (forall (y) (if (Name y) (iff  (DeReExists y) (exists x (and (Real x) (= x (* y))))))))


        ;;;
        ;
        Suppose (Believes! I (not (DeReExists I)))
        given (Believes! I (Name I))

        ;;;
        Perceive-the-belief (Believes! I (Perceives! I (Believes! I (not (DeReExists I)))))
        If_P_B (Believes!
                I
                (forall [?agent]
                        (if (Perceives! I (Believes! ?agent (not (DeReExists ?agent))))
                          (Real (* ?agent)))))

        }
 :goal        (and (Believes! I  (not (Real (* I))))
                   (Believes! I  (Real (* I)) ))

}
```

1.7 seconds

# Cogito Ergo Sum

$$\Lambda_{t_1}$$

```
{:name        "Cogito Ergo Sum"
 :description "A formaliztion of Descartes' Cogito Ergo Sum"
 :assumptions {

          S1 (Believes! I (forall [x] (or (Name x) (Thing x))))
          S2 (Believes! I (forall (x) (iff (Name x) (not (Thing x))))) )
          S3 (Believes! I (forall (x) (if (Thing x) (or (Real x) (Fictional x)))))
          S4 (Believes! I (forall (x) (if (Thing x) (iff (Real x) (not (Fictional x))))))
          ;;;
          A1 (Believes! I (forall (x) (if (Name x) (Thing (* x)))))
          A2 (Believes! I (forall (y) (if (Name y) (iff  (DeReExists y) (exists x (and (Real x) (= x (* y))))))))


          ;;;
          ;
          Suppose (Believes! I (not (DeReExists I)))
          given (Believes! I (Name I))

          ;;;
          Perceive-the-belief (Believes! I (Perceives! I (Believes! I (not (DeReExists I)))))
          If_P_B (Believes!
                  I
                  (forall [?agent]
                          (if (Perceives! I (Believes! ?agent (not (DeReExists ?agent))))
                            (Real (* ?agent)))))

          }
 :goal    (and (Believes! I  (not (Real (* I))))
               (Believes! I  (Real (* I)) ))

}
```

**1.7 seconds**

# Cogito Ergo Sum

$$\Lambda_{t_1}$$

```
{:name        "Cogito Ergo Sum"
 :description "A formaliztion of Descartes' Cogito Ergo Sum"
 :assumptions {

        S1 (Believes! I (forall [x] (or (Name x) (Thing x))))
        S2 (Believes! I (forall (x) (iff (Name x) (not (Thing x)))) )
        S3 (Believes! I (forall (x) (if (Thing x) (or (Real x) (Fictional x)))))
        S4 (Believes! I (forall (x) (if (Thing x) (iff (Real x) (not (Fictional x))))))
        ;;;
        A1 (Believes! I (forall (x) (if (Name x) (Thing (* x)))))
        A2 (Believes! I (forall (y) (if (Name y) (iff  (DeReExists y) (exists x (and (Real x) (= x (* y))))))))


        ;;;
        ;
        Suppose (Believes! I (not (DeReExists I)))
        given (Believes! I (Name I))

        ;;;
        Perceive-the-belief (Believes! I (Perceives! I (Believes! I (not (DeReExists I)))))
        If_P_B (Believes!
                I
                (forall [?agent]
                        (if (Perceives! I (Believes! ?agent (not (DeReExists ?agent))))
                          (Real (* ?agent)))))

        }
 :goal    (and (Believes! I  (not (Real (* I))))
               (Believes! I  (Real (* I)) ))

}
```

**absurd belief**

**1.7 seconds**

# Cogito Ergo Sum

$$\Lambda_{t_1}$$

```clojure
{:name        "Cogito Ergo Sum"
 :description "A formaliztion of Descartes' Cogito Ergo Sum"
 :assumptions {

        S1 (Believes! I (forall [x] (or (Name x) (Thing x))))
        S2 (Believes! I (forall (x) (iff (Name x) (not (Thing x)))) )
        S3 (Believes! I (forall (x) (if (Thing x) (or (Real x) (Fictional x)))))
        S4 (Believes! I (forall (x) (if (Thing x) (iff (Real x) (not (Fictional x))))))
        ;;;
        A1 (Believes! I (forall (x) (if (Name x) (Thing (* x)))))
        A2 (Believes! I (forall (y) (if (Name y) (iff  (DeReExists y) (exists x (and (Real x) (= x (* y))))))))


        ;;;
        ;
        Suppose (Believes! I (not (DeReExists I)))
        given (Believes! I (Name I))

        ;;;
        Perceive-the-belief (Believes! I (Perceives! I (Believes! I (not (DeReExists I)))))
        If_P_B (Believes!
                I
                (forall [?agent]
                        (if (Perceives! I (Believes! ?agent (not (DeReExists ?agent))))
                          (Real (* ?agent)))))

        }
 :goal    (and (Believes! I  (not (Real (* I))))
               (Believes! I  (Real (* I)) ))

}
```

**absurd belief**

**1.7 seconds**

$$\Lambda_{t_k}$$

# I. Elements of $\Lambda$

For top level beliefs, knowledge, intensions, desires etc

$\Lambda[\mathbf{B}, 1]$  Maximum intensional depth of beliefs

$\Lambda[\mathbf{D}, 1]$  Maximum intensional depth of desires

$\Lambda[\mathbf{I}, 1]$  Maximum intensional depth of intentions

$\vdots$

# II. Elements of $\Lambda$

For top level beliefs, knowledge, intensions, desires etc

$\Lambda[\mathbf{B}, 2]$  Maximum quantificational depth of beliefs

$\Lambda[\mathbf{D}, 2]$  Maximum quantificational depth of desires

$\Lambda[\mathbf{I}, 2]$  Maximum quantificational depth of intentions

$\vdots$

# III. Elements of $\Lambda$

For top level beliefs, knowledge, intensions, desires etc

$\Lambda[\mathbf{B},\, 3]$      Maximum extensional depth of beliefs

$\Lambda[\mathbf{D},\, 3]$      Maximum extensional depth of desires

$\Lambda[\mathbf{I},\, 3]$      Maximum extensional depth of intentions

$\vdots$

# IV. Elements of $\Lambda$

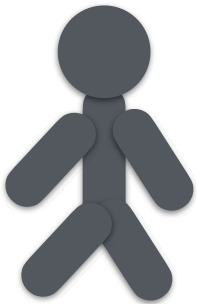For top level beliefs, knowledge, intensions, desires etc

$\Lambda[\mathbf{B}, 4]$     Maximum difference between time expressions within beliefs

$\Lambda[\mathbf{D}, 4]$     Maximum difference between time expressions within desires

$\Lambda[\mathbf{I}, 4]$     Maximum difference between time expressions within intentions

$\vdots$

**Note**: If a time variable **t** is universally quantified, we take it as **∞**.
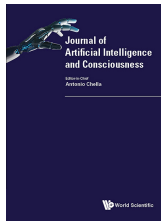
# Example

**The Doctrine of Double Effect**

$C_1$    the action is not forbidden (where we assume an ethical hierarchy such as the one given by Bringsjord [2017], and require that the action be neutral or above neutral in such a hierarchy);

$C_2$    The net utility or goodness of the action is greater than some positive amount $\gamma$;

$C_{3a}$    the agent performing the action intends only the good effects;

$C_{3b}$    the agent does not intend any of the bad effects;

$C_4$    the bad effects are not used as a means to obtain the good effects; and

$C_5$    if there are bad effects, the agent would rather the situation be different and the agent not have to perform the action. That is, the action is unavoidable.

# The Theory of Cognitive Consciousness, and $\Lambda$ (Lambda)

Selmer Bringsjord ✉ and G. Naveen Sundar

# The Theory of Cognitive Consciousness, and Λ (Lambda)

Selmer Bringsjord ✉ and G. Naveen Sundar

---

### The Theory of Cognitive Consciousness, and Λ (Lambda)*

Selmer Bringsjord

*Rensselaer AI & Reasoning (RAIR) Lab
Department of Cognitive Science
Department of Computer Science
Lally School of Management
Rensselaer Polytechnic Institute (RPI)
Troy NY 12180 USA*
Selmer.Bringsjord@gmail.com

Naveen Sundar G.

*Rensselaer AI & Reasoning (RAIR) Lab
Rensselaer Polytechnic Institute (RPI)
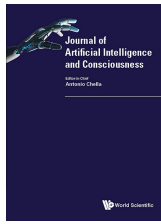Troy NY 12180 USA*
Naveen.Sundar.G@gmail.com

We provide an overview of the theory of cognitive consciousness (TCC), and of Λ; the latter provides a means of measuring the amount of cognitive consciousness present in a given cognizer, whether natural or artificial, at a given time, along a number of different dimensions. TCC and Λ stand in stark contrast to Tononi's Integrated information Theory (IIT) and Φ. We believe, for reasons we present, that the former pair is superior to the latter. TCC includes a formal axiomatic theory, $\mathcal{CA}$, the 12 axioms of which we present and briefly comment upon herein; no such formal theory accompanies IIT/Φ. TCC/Λ and IIT/Φ each offer radically different verdicts as to whether and to what degree AIs of yesterday, today, and tomorrow were/are/will be conscious. Another noteworthy difference between TCC/Λ and IIT/Φ is that the former enables the measurement of cognitive consciousness in those who have passed on, and in fictional characters; no such enablement is remotely possible for IIT/Φ. For instance, we apply Λ to measure the cognitive consciousness of: Descartes; the first fictional detective to be described on Earth (by Edgar Allen Poe), C. Auguste Dupin. We also apply Λ to compute the cognitive consciousness of an artificial agent able to make ethical decisions using the Doctrine of Double Effect.

*Keywords*: consciousness; cognitive consciousness; AI; Lambda/Λ.

# The Theory of Cognitive Consciousness, and $\Lambda$ (Lambda)

### Extending Measures from $\mathcal{L}^0$ to $\mathcal{L}$

$$\mu_\omega(\phi) = \begin{cases} \mu(\phi) & \text{if } \phi \in \mathcal{L}^0 \\ \max_\psi \mu_\omega(\psi) + 1 & \text{if } \phi \equiv \omega_i[a_1, t_1, \ldots \psi \ldots] \end{cases}$$

For example, let $\mu$ count the number of predicate symbols in a formula.

### Example

$$\mu(Happy(john)) = 1$$
$$\mu_\omega(Happy(john)) = 1$$
$$\mu_\omega\Big(\mathbf{B}(mary, t_2, Happy(john))\Big) = 2$$

For any agent $a$, we want to look at the new complexity the agent introduces that is above any input complexity. For this, we introduce $\Delta : 2^\mathcal{L} \times 2^\mathcal{L} \to 2^\mathcal{L}$ operator that computes differences between two sets of formulae. This can be simply the set-difference operator. For convenience, let $\omega_j[\Gamma]$ denote the subset of formulae with operators $\omega_j$ in $\Gamma$:

$$\omega_j[\Gamma] = \{\phi \mid \phi \equiv \omega_j[\ldots] \text{ and } \phi \in \Gamma \text{ or } \phi \text{ a subformula } \in \Gamma\}$$

Given a set of measures $\{\mu^0, \ldots, \mu^N\}$ and a set of modal (or cognitive) operators $\{\omega_0, \ldots, \omega_M\}$, we define $\Lambda$ as a function mapping an agent at a time point to a matrix $\mathbb{N}^{M \times N}$:

$$\Lambda : A \times T \to \mathbb{N}^{M \times N}$$

### Definition of $\Lambda$

$$\Lambda(a,t)_{i,j} = \max_\phi \left\{ \mu^i(\phi) \mid \phi \in \Delta\Big(\omega_j[o(a,t)], \omega_j[i(a,t)]\Big) \right\}$$

### Example 2

Let us consider two modal operators $\{\mathbf{B}, \mathbf{D}\}$ and the following base measures $\mu^0$ which measures quantificational complexity via $\Sigma$ or $\Pi$ measures, $\mu^1$ which counts the total number of predicate symbols (not a count of unique predicate symbols), and $\mu^2$ which counts the number of distinct time expressions. This gives $\Lambda : A \times T \to \mathbb{N}^{2 \times 3}$. At some timepoint $t$, let an agent $a$ have the following $\Delta(o(a,t), i(a,t)) = \{\mathbf{B}(\phi_1), \mathbf{D}(\phi_2)\}$

# The Theory of Cognitive Consciousness, and Λ (Lambda)

### Extending Measures from $\mathcal{L}^0$ to $\mathcal{L}$

$$\mu_\omega(\phi) = \begin{cases} \mu(\phi) & \text{if } \phi \in \mathcal{L}^0 \\ \max_\psi \mu_\omega(\psi) + 1 & \text{if } \phi \equiv \omega_i[a_1, t_1, \dots \psi \dots] \end{cases}$$

For example, let $\mu$ count the number of predicate symbols in a formula.

### Example

$$\mu(Happy(john)) = 1$$
$$\mu_\omega(Happy(john)) = 1$$
$$\mu_\omega\Big(\mathbf{B}\big(mary, t_2, Happy(john)\big)\Big) = 2$$

For any agent $a$, we want to look at the new complexity the agent introduces that is above any input complexity. For this, we introduce $\Delta : 2^{\mathcal{L}} \times 2^{\mathcal{L}} \to 2^{\mathcal{L}}$ operator that computes differences between two sets of formulae. This can be simply the set-difference operator. For convenience, let $\omega_j[\Gamma]$ denote the subset of formulae with operators $\omega_j$ in $\Gamma$:

$$\omega_j[\Gamma] = \{\phi \mid \phi \equiv \omega_j[\dots] \text{ and } \phi \in \Gamma \text{ or } \phi \text{ a subformula } \in \Gamma\}$$

Given a set of measures $\{\mu^0, \dots, \mu^N\}$ and a set of modal (or cognitive) operators $\{\omega_0, \dots, \omega_M\}$, we define $\Lambda$ as a function mapping an agent at a time point to a matrix $\mathbb{N}^{M \times N}$:

$$\Lambda : A \times T \to \mathbb{N}^{M \times N}$$

### Definition of $\Lambda$

$$\Lambda(a,t)_{i,j} = \max_\phi \Big\{ \mu^i(\phi) \mid \phi \in \Delta\Big(\omega_j\big[o(a,t)\big], \omega_j\big[i(a,t)\big]\Big) \Big\}$$

### Example 2

Let us consider two modal operators $\{\mathbf{B}, \mathbf{D}\}$ and the following base measures $\mu^0$ which measures quantificational complexity via $\Sigma$ or $\Pi$ measures, $\mu^1$ which counts the total number of predicate symbols (not a count of unique predicate symbols), and $\mu^2$ which counts the number of distinct time expressions. This gives $\Lambda : A \times T \to \mathbb{N}^{2 \times 3}$. At some timepoint $t$, let an agent $a$ have the following $\Delta(o(a,t), i(a,t)) = \{\mathbf{B}(\phi_1), \mathbf{D}(\phi_2)\}$

$$\phi_1 \equiv \neg \forall a : Happy(a,t); \qquad \phi_2 \equiv \forall b : \neg Hungry(b,t) \to Happy(b,t)$$

Applying the measures:

$$\mu^o(\phi_1) = 1, \mu^1(\phi_1) = 1; \mu^2(\phi_1) = 1$$
$$\mu^o(\phi_2) = 1; \mu^1(\phi_2) = 2; \mu^2(\phi_2) = 1$$

Giving us:

$$\Lambda(a,t) = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \end{pmatrix}$$

### 6.1.  *Some Distinctive Properties of $\Lambda$ (vs. $\Phi$)*

Here are some properties of the $\Lambda$ framework of potential interest to our readers:

**Non-Binary** Whereas $\Phi$ is such that an agent either is or is not (P-) conscious, cognitive consciousness as measured by $\Lambda$ admits of a fine-grained range of the *degree* of cognitive consciousness.

**Zero $\Lambda$ for Some Animals and Machines** Animals such as insects, and computing machines that are end-to-end statistical/connectionist "ML," have zero $\Lambda$, and hence cannot be cognitively conscious. In contrast, as emphasized to Bringsjord in personal conversation,[6] $\Phi$ says that even lower animals are conscious.

**Human-Nonhuman Discontinuity Explained by $\Lambda$** From the computational/AI point of view, cognitive scientists have taken note of a severe discontinuity between *H. sapiens sapiens* and other biological creatures on Earth [Penn *et al.*, 2008], and the sudden and large jump in level of $\Lambda$ from (say) chimpanzees and dolphins to humans is in line with this observation. It's for instance doubtful that any nonhuman animals are capable of reaching third-order belief; hence $\Lambda[\mathbf{B}, 0] = n$, where $n \geq 3$, for any nonhuman animal, is impossible. In stark contrast, each of us believes that you, the reader, believe that we believe that San Francisco is located in California.

**Human-Human Discontinuity Explained by $\Lambda$** A given neurobiologically normal human, over the course of his or her lifetime, has very different cognitive capacity. E.g., it's well-known that such a human, before the age of four or five, is highly unlikely to be able to solve what has become known as the *false-belief task* (or sometimes the *sally-anne task*), which we denote by 'FBT.' From the point of view of $\Lambda$, the explanation is simply that an agent with insufficiently high cognitive consciousness is incapable of solving such a task; specifically, solving FBT requires an agent to have

[6]With Tononi and C. Koch, SRI T&C Series.

## Formal Conditions for $\mathcal{DDE}$

**F$_1$** $\alpha$ carried out at $t$ is not forbidden. That is:

$$\Gamma \nvdash \neg\mathbf{O}\left(a, t, \sigma, \neg happens\left(action(a,\alpha), t\right)\right)$$

**F$_2$** The net utility is greater than a given positive real $\gamma$:

$$\Gamma \vdash \sum_{y=t+1}^{H}\left(\sum_{f \in \alpha_I^{a,t}} \mu(f,y) - \sum_{f \in \alpha_T^{a,t}} \mu(f,y)\right) > \gamma$$

**F$_{3a}$** The agent $a$ intends at least one good effect. (**F$_2$** should still hold after removing all other good effects.) There is at least one fluent $f_g$ in $\alpha_I^{a,t}$ with $\mu\left(f_g, y\right) > 0$, or $f_b$ in $\alpha_T^{a,t}$ with $\mu(f_b, y) < 0$, and some $y$ with $t < y \leq H$ such that the following holds:

$$\Gamma \vdash \begin{pmatrix} \exists f_g \in \alpha_I^{a,t} \ \mathbf{I}\left(a, t, Holds\left(f_g, y\right)\right) \\ \vee \\ \exists f_b \in \alpha_T^{a,t} \ \mathbf{I}\left(a, t, \neg Holds\left(f_b, y\right)\right) \end{pmatrix}$$

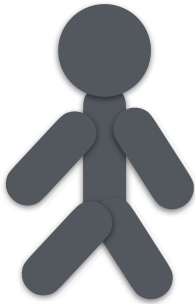**F$_{3b}$** The agent $a$ does not intend any bad effect. For all fluents $f_b$ in $\alpha_I^{a,t}$ with $\mu\left(f_b, y\right) < 0$, or $f_g$ in $\alpha_T^{a,t}$ with $\mu\left(f_g, y\right) > 0$, and for all $y$ such that $t < y \leq H$ the following holds:

$$\Gamma \nvdash \mathbf{I}\left(a, t, Holds\left(f_b, y\right)\right) \text{ and}$$

$$\Gamma \nvdash \mathbf{I}\left(a, t, \neg Holds\left(f_g, y\right)\right)$$

**F$_4$** The harmful effects don't cause the good effects. Four permutations, paralleling the definition of $\rhd$ above, hold here. One such permutation is shown below. For any bad fluent $f_b$ holding at $t_1$, and any good fluent $f_g$ holding at some $t_2$, such that $t < t_1, t_2 \leq H$, the following holds:
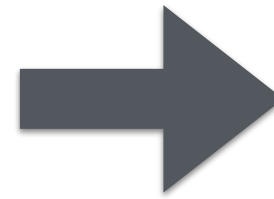
$$\Gamma \vdash \neg\rhd\left(Holds\left(f_b, t_1\right), Holds\left(f_g, t_2\right)\right)$$

# Example from Sim in IJCAI Paper

looking at one single chunk

$$\left\{ \begin{array}{l} \mathbf{K}\Big(I, now, \sigma_{trolley}\Big), \\[2pt] \mathbf{B}\left(I, now, \mathbf{O}\left(\begin{array}{l} I, now, \sigma_{trolley}, \\ \left[\begin{array}{c} \neg \exists t : Moment\ Holds\big(dead(P_1, t)\big) \\ \wedge \\ \neg \exists t : Moment\ Holds\big(dead(P_2, t)\big) \end{array}\right] \end{array}\right)\right), \\[2pt] \mathbf{O}\left(I, now, \sigma_{trolley}, \left[\begin{array}{l} \neg \exists t : Moment\ Holds\,(dead(P_1, t)) \wedge \\ \neg \exists t : Moment\ Holds\,(dead(P_2, t)) \end{array}\right]\right) \\[2pt] \vdash\ \mathbf{I}\left(I, now, \left[\begin{array}{l} \neg \exists t : Moment\ Holds\big(dead(P_1, t)\big) \wedge \\ \neg \exists t : Moment\ Holds\big(dead(P_2, t)\big) \end{array}\right]\right) \end{array} \right\}$$

$\Lambda[\mathbf{B}, 1] = 2$

$\Lambda[\mathbf{B}, 2] = 1$

$\Lambda[\mathbf{K}, 1] = 1$

$\Lambda[\mathbf{O}, 1] = 1$

$\Lambda[\mathbf{O}, 1] = 1$

$\Lambda[\mathbf{I}, 1] = 1$

$\Lambda[\mathbf{I}, 2] = 1$

$\Lambda[\mathbf{B}, 3] = 1$

$\Lambda[\mathbf{B}, 4] = \infty$

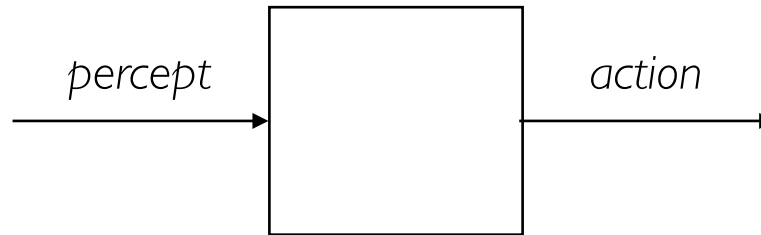$\vdots$

The application of $\Lambda$ to eg "Deep Learning" machines implies that they have zero cognitive intelligence/ cognitive consciousness.

# AI:MLn

AI

percept → [ ] → action

# AI:MLn
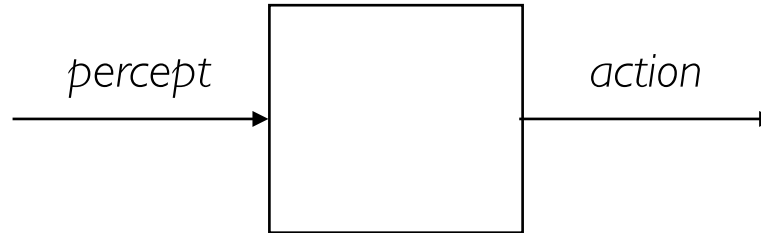
AI

percept → [ ] → action

$\langle n_1, n_2, \ldots, n_k \rangle, k \in \mathbb{Z}^+$

# AI:MLn

AI

percept                                   action

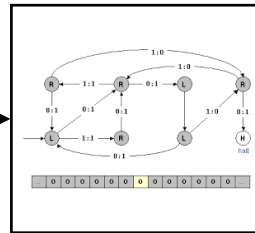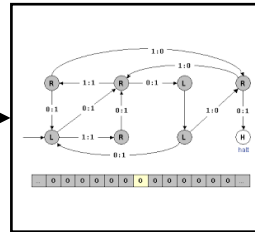$\langle n_1, n_2, \ldots, n_k \rangle, k \in \mathbb{Z}^+$

# AI:MLn



AI

*percept*

*action*

$\langle n_1, n_2, \ldots, n_k \rangle, k \in \mathbb{Z}^+$

# AI:MLn



AI

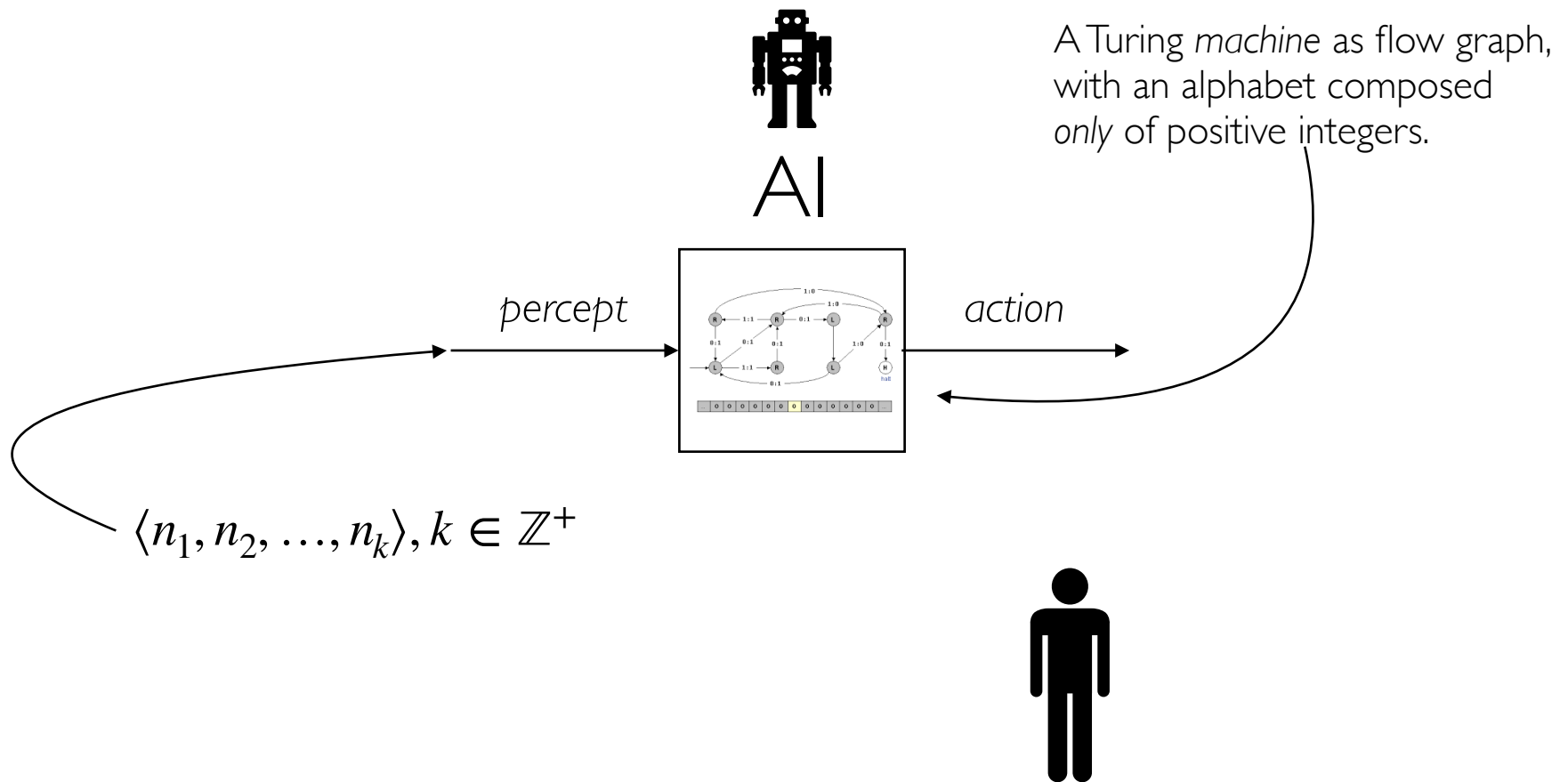A Turing *machine* as flow graph, with an alphabet composed *only* of positive integers.

*percept*

*action*

$\langle n_1, n_2, \ldots, n_k \rangle, k \in \mathbb{Z}^+$

# AI:MLn

A Turing *machine* as flow graph, with an alphabet composed *only* of positive integers.

AI

*percept*

*action*

$\langle n_1, n_2, \ldots, n_k \rangle, k \in \mathbb{Z}^+$

# AI:MLn



AI

A Turing *machine* as flow graph, with an alphabet composed *only* of positive integers.
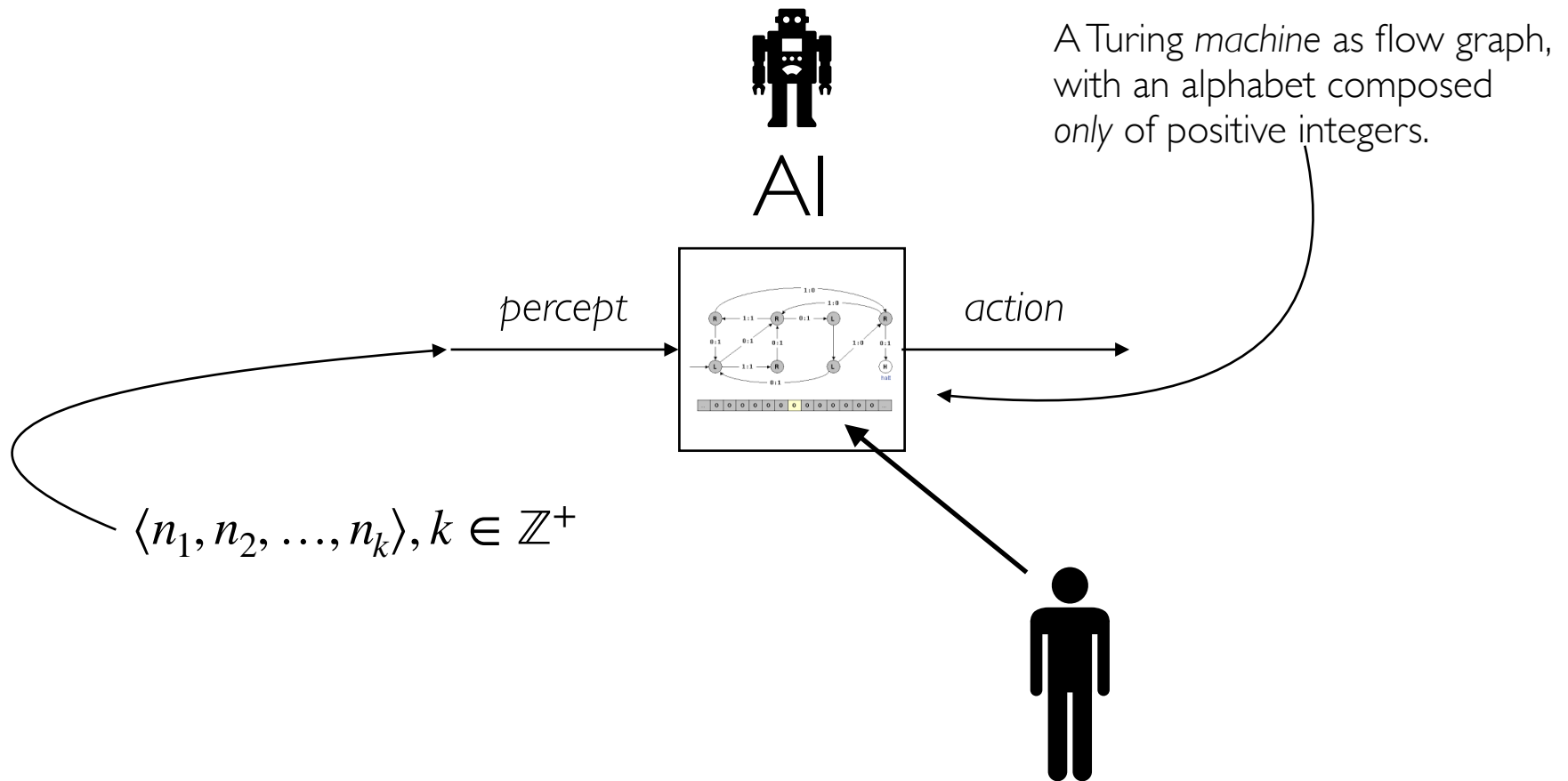
*percept*

*action*

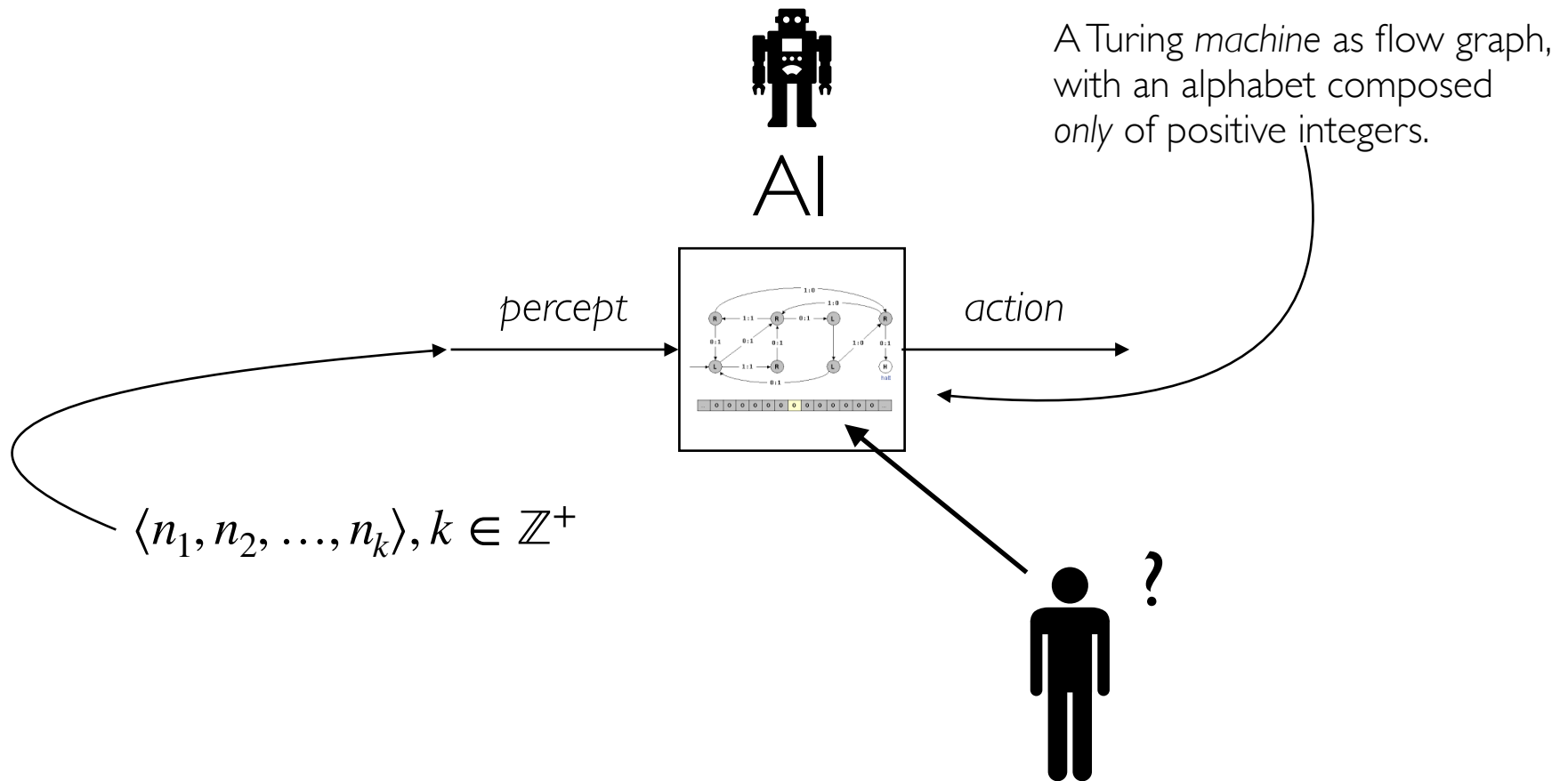$\langle n_1, n_2, \ldots, n_k \rangle, k \in \mathbb{Z}^+$

# AI:MLn



AI

A Turing *machine* as flow graph, with an alphabet composed *only* of positive integers.

*percept*

*action*

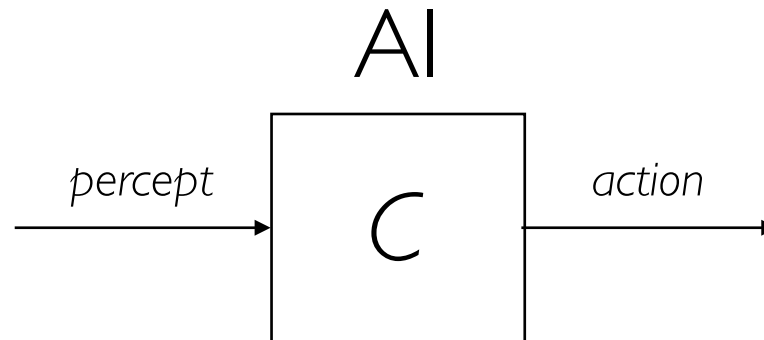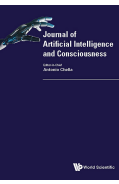$\langle n_1, n_2, \ldots, n_k \rangle, k \in \mathbb{Z}^+$

?

We will be able to measure the intelligence of *any* AI, not with g-loaded tests of intelligence, but with $\Lambda$-loaded tests of machine intelligence, in keeping with Psychometric AI.

AI

*percept* → C → *action*

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

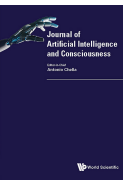K2B

Intro

Incorr

Ess

¬CompE

Irr

Free

CCaus

TheI

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

K2B $\quad \forall a[\mathbf{K}_a\phi \rightarrow (\mathbf{B}_a\phi \wedge \mathbf{B}_a\exists\Phi\exists\alpha(\Phi \rightsquigarrow_{\alpha/\pi} \phi))]$

Intro

Incorr

Ess

¬CompE

Irr

Free

CCaus

TheI

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

$\mu\mathcal{DCEC}_3^*$ K2B $\forall a[\mathbf{K}_a\phi \rightarrow (\mathbf{B}_a\phi \wedge \mathbf{B}_a\exists\Phi\exists\alpha(\Phi \rightsquigarrow_{\alpha/\pi} \phi))]$
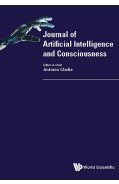
Intro

Incorr

Ess

¬CompE

Irr

Free

CCaus

TheI

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

K2B $\quad \forall a[\mathbf{K}_a \phi \rightarrow (\mathbf{B}_a \phi \wedge \mathbf{B}_a \exists \Phi \exists \alpha (\Phi \leadsto_{\alpha/\pi} \phi)]$

Intro

Incorr

Ess

¬CompE

Irr

Free

CCaus

TheI

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

K2B $\quad \forall a[\mathbf{K}_a \phi \rightarrow (\mathbf{B}_a \phi \wedge \mathbf{B}_a \exists \Phi \exists \alpha (\Phi \rightsquigarrow_{\alpha/\pi} \phi))]$

Intro

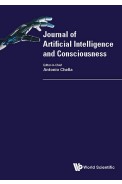Incorr $\quad \forall a \forall t \forall F[(F \text{ is contingent } \wedge F \in C'') \rightarrow (\Box \mathbf{B}(a, t, Fa) \rightarrow Fa)]$

Ess

¬CompE

Irr

Free

CCaus

TheI

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

K2B  $\forall a[\mathbf{K}_a\phi \rightarrow (\mathbf{B}_a\phi \wedge \mathbf{B}_a\exists\Phi\exists\alpha(\Phi \rightsquigarrow_{\alpha/\pi} \phi)]$

Intro

Incorr  $\forall a\forall t\forall F[(F\text{ is contingent } \wedge F \in C'') \rightarrow (\Box\mathbf{B}(a, t, Fa) \rightarrow Fa)]$

Ess

¬CompE

Irr

Free

CCaus  $\mathbf{C} \; \mathcal{EC}$

TheI

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

K2B $\quad \forall a[\mathbf{K}_a\phi \to (\mathbf{B}_a\phi \wedge \mathbf{B}_a\exists\Phi\exists\alpha(\Phi \rightsquigarrow_{\alpha/\pi} \phi)]$

Intro

Incorr $\quad \forall a\forall t\forall F[(F\ is\ contingent\ \wedge F \in C'') \to (\Box\mathbf{B}(a,t,Fa) \to Fa)]$

Ess

¬CompE

Irr

Free

$[A_1]\ \mathbf{C}(\forall f,t\ .\ initially(f) \wedge \neg clipped(0,f,t) \Rightarrow holds(f,t))$

$[A_2]\ \mathbf{C}(\forall e,f,t_1,t_2\ .\ happens(e,t_1) \wedge initiates(e,f,t_1) \wedge t_1 < t_2 \wedge \neg clipped(t_1,f,t_2) \Rightarrow holds(f,t_2))$

$[A_3]\ \mathbf{C}(\forall t_1,f,t_2\ .\ clipped(t_1,f,t_2) \Leftrightarrow [\exists e,t\ .\ happens(e,t) \wedge t_1 < t < t_2 \wedge terminates(e,f,t)])$

$[A_4]\ \mathbf{C}(\forall a,d,t\ .\ happens(action(a,d),t) \Rightarrow \mathbf{K}(a,happens(action(a,d),t)))$

$[A_5]\ \mathbf{C}(\forall a,f,t,t'\ .\ \mathbf{B}(a,holds(f,t)) \wedge \mathbf{B}(a,t<t') \wedge \neg\mathbf{B}(a,clipped(t,f,t')) \Rightarrow \mathbf{B}(a,holds(f,t')))$

CCaus $\quad \mathbf{C}\ \mathcal{EC}$

TheI

# $\mathcal{CA}$: 11 Axioms (Initially)

Plan

P2B

K2B  $\forall a[\mathbf{K}_a\phi \to (\mathbf{B}_a\phi \wedge \mathbf{B}_a\exists\Phi\exists\alpha(\Phi \rightsquigarrow_{\alpha/\pi} \phi)]$

Intro

Incorr  $\forall a\forall t\forall F[(F\text{ is contingent } \wedge F \in C'') \to (\Box\mathbf{B}(a,t,Fa) \to Fa)]$

Ess

¬CompE

Irr

Free          **C** SpecRel

CCaus  **C** $\mathcal{EC}$

TheI

$[A_1]$ $\mathbf{C}(\forall f, t \ . \ initially(f) \wedge \neg clipped(0, f, t) \Rightarrow holds(f, t))$
$[A_2]$ $\mathbf{C}(\forall e, f, t_1, t_2 \ . \ happens(e, t_1) \wedge initiates(e, f, t_1) \wedge t_1 < t_2 \wedge \neg clipped(t_1, f, t_2) \Rightarrow holds(f, t_2))$
$[A_3]$ $\mathbf{C}(\forall t_1, f, t_2 \ . \ clipped(t_1, f, t_2) \Leftrightarrow [\exists e, t \ . \ happens(e, t) \wedge t_1 < t < t_2 \wedge terminates(e, f, t)])$
$[A_4]$ $\mathbf{C}(\forall a, d, t \ . \ happens(action(a, d), t) \Rightarrow \mathbf{K}(a, happens(action(a, d), t)))$
$[A_5]$ $\mathbf{C}(\forall a, f, t, t' \ . \ \mathbf{B}(a, holds(f, t)) \wedge \mathbf{B}(a, t < t') \wedge \neg\mathbf{B}(a, clipped(t, f, t')) \Rightarrow \mathbf{B}(a, holds(f, t')))$

# Example

```
{:name        "Knowability paradox"
 :description " \exists p  ~\Diamond \exists x Kx (Tp & ~ \exist y Ky Tp)"

 :assumptions {}
 :goal (exists [?P] (not (pos (exists [?x] (Knows! ?x (and ?P (not (exists [?y] (Knows! ?y ?P))))))))))}
```

$$\Lambda[\mathbf{\kappa}, 1] = 2$$

$$\Lambda[\mathbf{\kappa}, 2] = 1$$

$$\Lambda[\mathbf{\kappa}, 2] = 2 \qquad \textit{Since the above goal is in second-order modal logic}$$

Λ

| | 1 | 2 | 3 | 4 | ... |
|---|---|---|---|---|---|
| B | | | | | |
| K | | | | | |
| D | | | | | |
| O | | | | | |
| ... | | | | | |

Λ  *Itself varies across time*

*Max, Mean can be considered too.*

Λ

| | 1 | 2 | 3 | 4 | ... |
|---|---|---|---|---|---|
| B | | | | | |
| K | | | | | |
| D | | | | | |
| O | | | | | |
| ... | | | | | |

Λ

| | 1 | 2 | 3 | 4 | ... |
|---|---|---|---|---|---|
| B | | | | | |
| K | | | | | |
| D | | | | | |
| O | | | | | |
| ... | | | | | |

Λ

| | 1 | 2 | 3 | 4 | ... |
|---|---|---|---|---|---|
| B | | | | | |
| K | | | | | |
| D | | | | | |
| O | | | | | |
| ... | | | | | |

Λ

| | 1 | 2 | 3 | 4 | ... |
|---|---|---|---|---|---|
| B | | | | | |
| K | | | | | |
| D | | | | | |
| O | | | | | |
| ... | | | | | |

$t_0$  $t_1$  $t_2$  $t_3$

What is the level of consciousness ($= \Lambda$ value) enjoyed by this self-conscious robot?

*Med nok penger, kan logikk løse alle problemer.*