

# **Paper Discussion;** **DDE (Doctrine of Double Effect) =>** **DDE\* (for self-sacrifice) =>** **DTripleE**

**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  
11/17/2022  
ver 1117221237NY



**Logistics ...**

Let's look real-time now, online. Btw, remember, first draft is due *after* Tgiving break.

The screenshot displays the Overleaf online LaTeX editor interface. The top navigation bar includes icons for Menu, Home, and a document title 'IFLAI2F22\_PAPERTOPICS'. It also features buttons for 'J' (Jupyter), 'E' (Editor), 'Review', 'Share', 'Submit', 'History', 'Layout', and 'Chat'. Below this, a secondary toolbar shows icons for file operations and tabs for 'Source', 'Source (legacy)', and 'Rich Text'. The left sidebar contains a file explorer with 'main.tex' and 'main72.bib', and a 'File outline' section with links like 'General Orientation' and 'Formatting, Due Dates/S...'. The main workspace shows a LaTeX document with line numbers 1 through 23. The code includes comments in blue and LaTeX commands such as `\documentclass[11pt]{article}`, `\usepackage[utf8]{inputenc}`, `\usepackage{fullpage}`, `\usepackage{setspace}`, `\usepackage{amssymb}`, `\usepackage[colorlinks]{hyperref}`, `\usepackage{harvard}`, `\usepackage{color}`, `\usepackage{marvosym}`, and `\usepackage{mathrsfs}`. A search bar at the bottom left shows 'Rectifying' with options 'Aa', '[.\*]', and 'W'. To the right of the search bar are navigation arrows and a '0 of 0' indicator. At the bottom right, there are buttons for 'Replace' and 'Replace All'. On the far right, a vertical strip shows a preview of the document's rendered output.

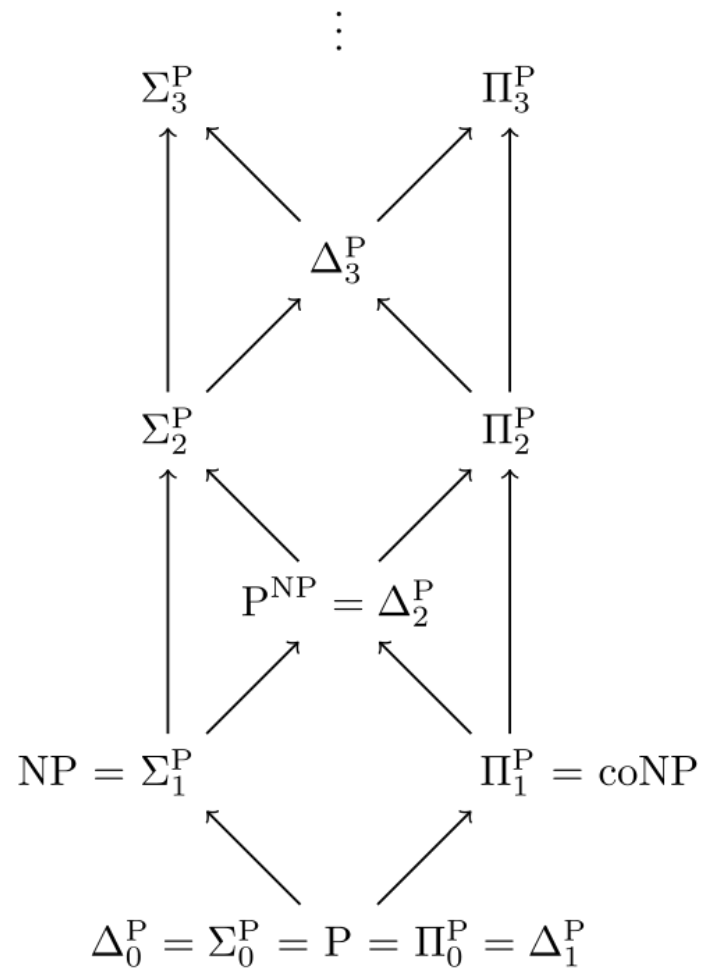
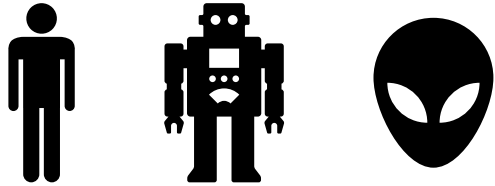
```
1  %% TODO
2
3  %% [ ]
4
5
6  \documentclass[11pt]{article}
7
8  \usepackage[utf8]{inputenc}
9
10 \usepackage{fullpage} %% <= why not use this in your own paper?
11
12 \usepackage{setspace}
13 %% Toggle the following on for doublespacing:
14 %\doublespacing
15
16
17 %% Some standard package calls by S:https://www.overleaf.com/project/63516802b5977f481cb0f57a
18 \usepackage{amssymb}
19 \usepackage[colorlinks]{hyperref}
20 \usepackage{harvard} %% Selmer's preference for citations/References.
21 \usepackage{color}
22 \usepackage{marvosym}
23 \usepackage{mathrsfs}
```

Further delivery on  
promissory note re  
building hierarchies via  
formal logic...



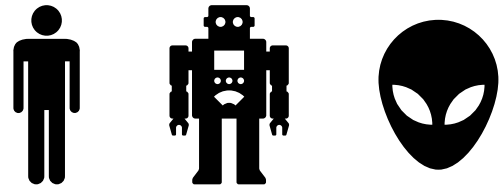
# Polynomial Hierarchy, Part II

(via formal logic, directly)

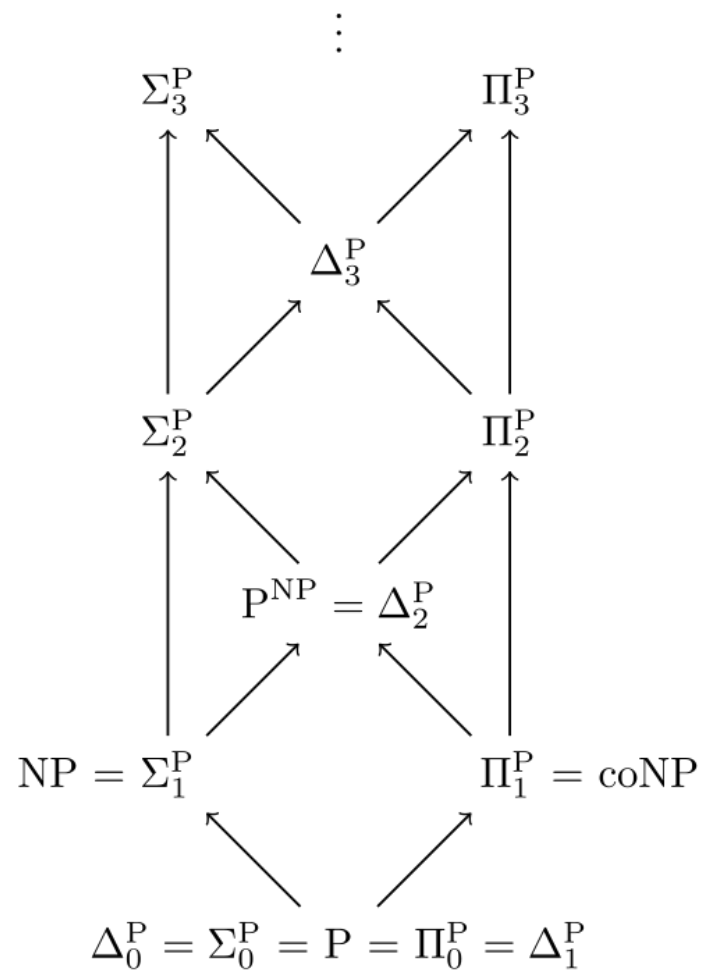


# Polynomial Hierarchy, Part II

(via formal logic, directly)

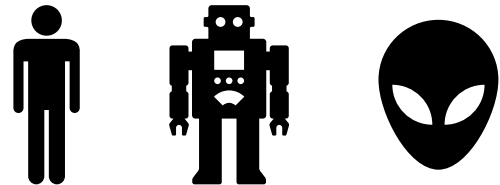


Eg:



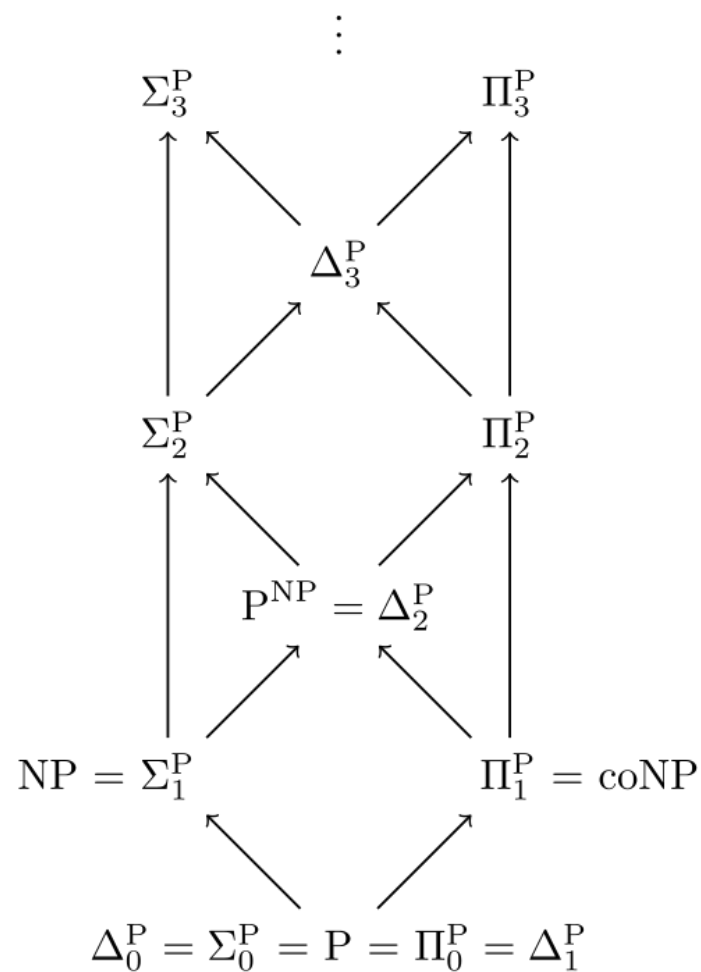
# Polynomial Hierarchy, Part II

(via formal logic, directly)



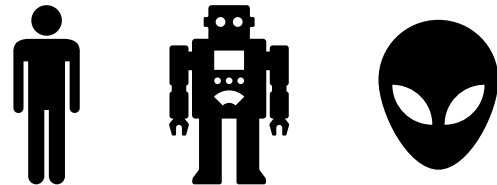
Eg:

$$\langle \phi_1, k \rangle \in L \text{ iff } \exists \phi_2 \forall \alpha KLogEquiv(\phi_1, \phi_2, |\phi_2| \leq k, \alpha(\phi_1) = \alpha(\phi_2))$$



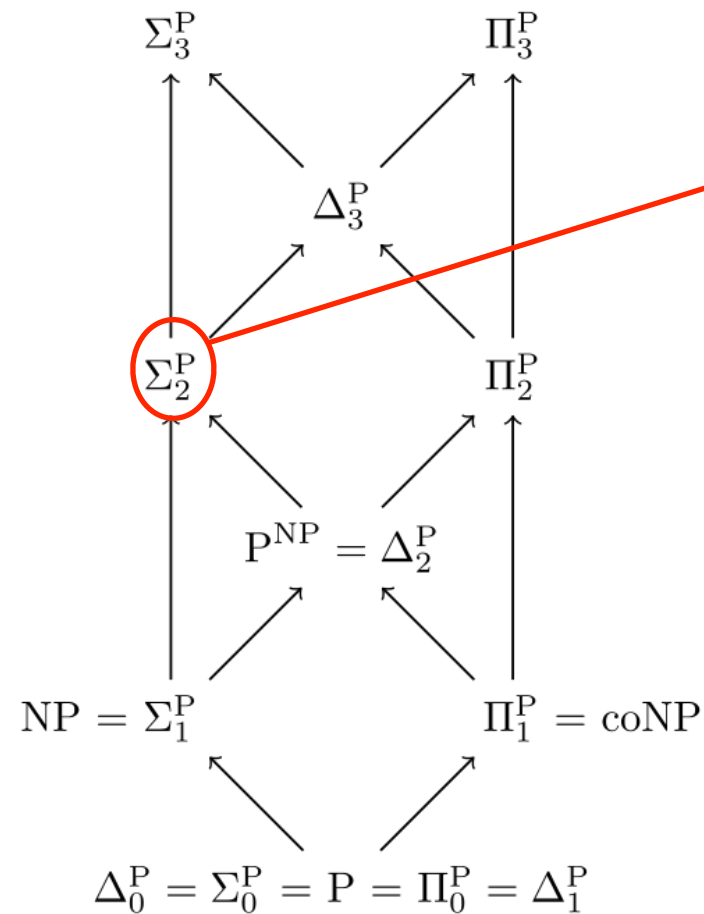
# Polynomial Hierarchy, Part II

(via formal logic, directly)



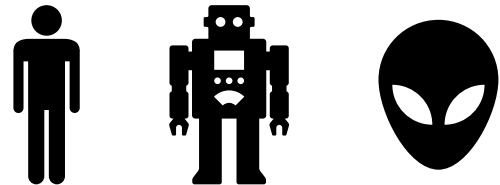
Eg:

$$\langle \phi_1, k \rangle \in L \text{ iff } \exists \phi_2 \forall \alpha KLogEquiv(\phi_1, \phi_2, |\phi_2| \leq k, \alpha(\phi_1) = \alpha(\phi_2))$$

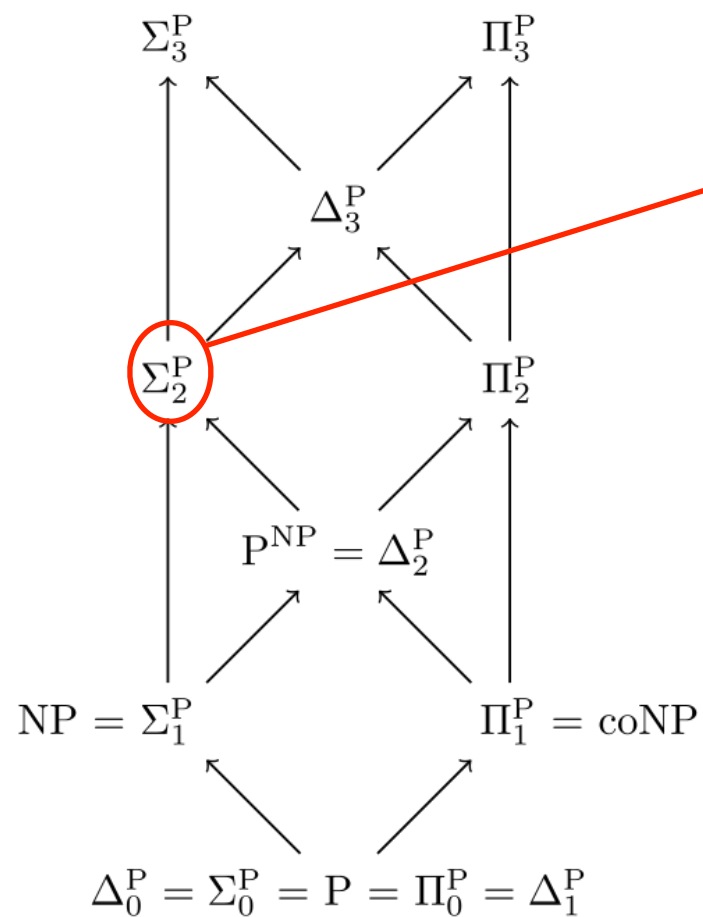


# Polynomial Hierarchy, Part II

(via formal logic, directly)



⋮



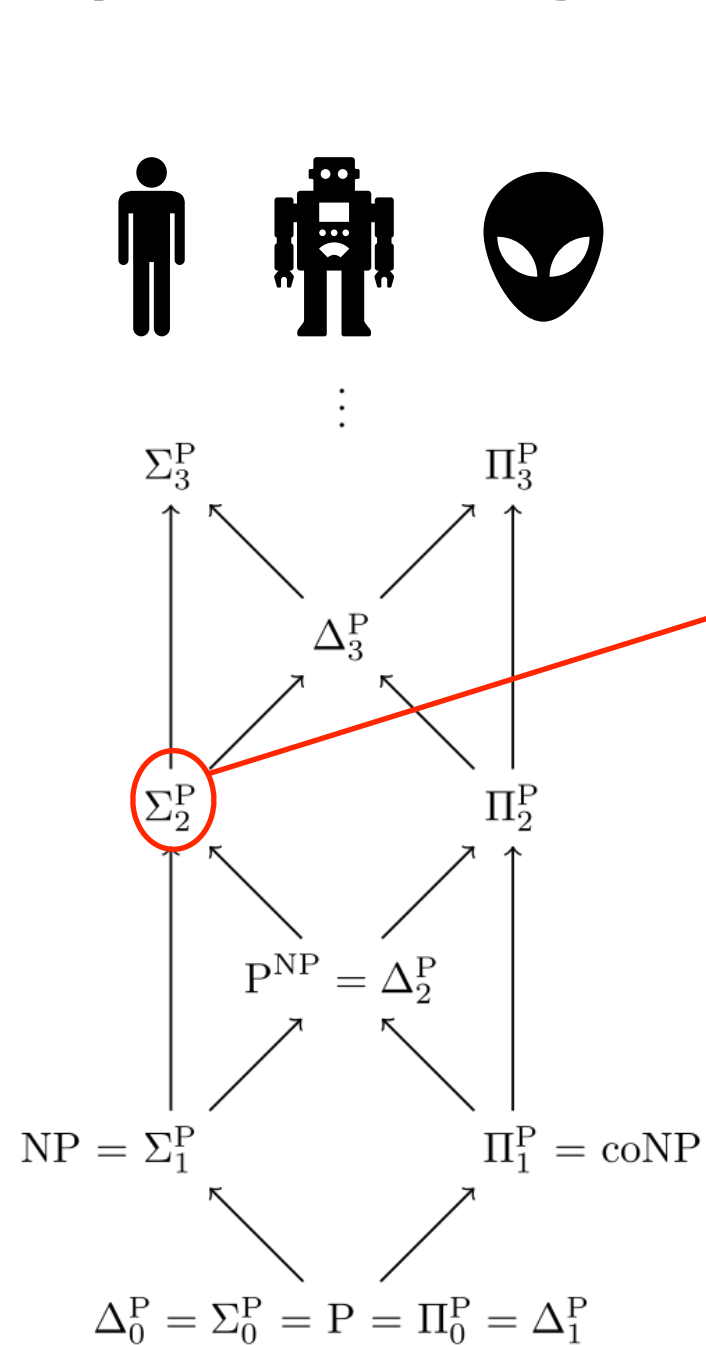
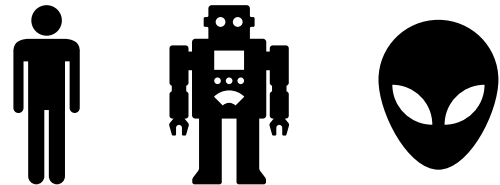
free variables

Eg:

$$\langle \phi_1, k \rangle \in L \text{ iff } \exists \phi_2 \forall \alpha KLogEquiv(\phi_1, \phi_2, |\phi_2| \leq k, \alpha(\phi_1) = \alpha(\phi_2))$$

# Polynomial Hierarchy, Part II

(via formal logic, directly)



free variables

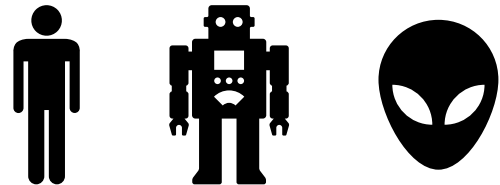
Eg:

$$\langle \phi_1, k \rangle \in L \text{ iff } \exists \phi_2 \forall \alpha KLogEquiv(\phi_1, \phi_2, |\phi_2| \leq k, \alpha(\phi_1) = \alpha(\phi_2))$$

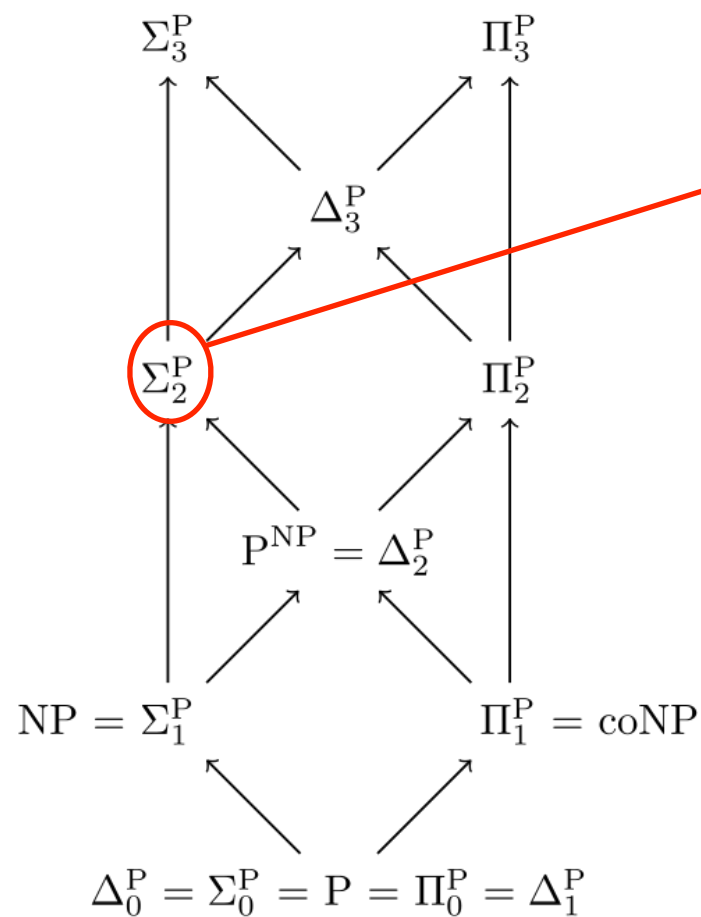
Now we generalize:

# Polynomial Hierarchy, Part II

(via formal logic, directly)



⋮



free variables

Eg:

$$\langle \phi_1, k \rangle \in L \text{ iff } \exists \phi_2 \forall \alpha KLogEquiv(\phi_1, \phi_2, |\phi_2| \leq k, \alpha(\phi_1) = \alpha(\phi_2))$$

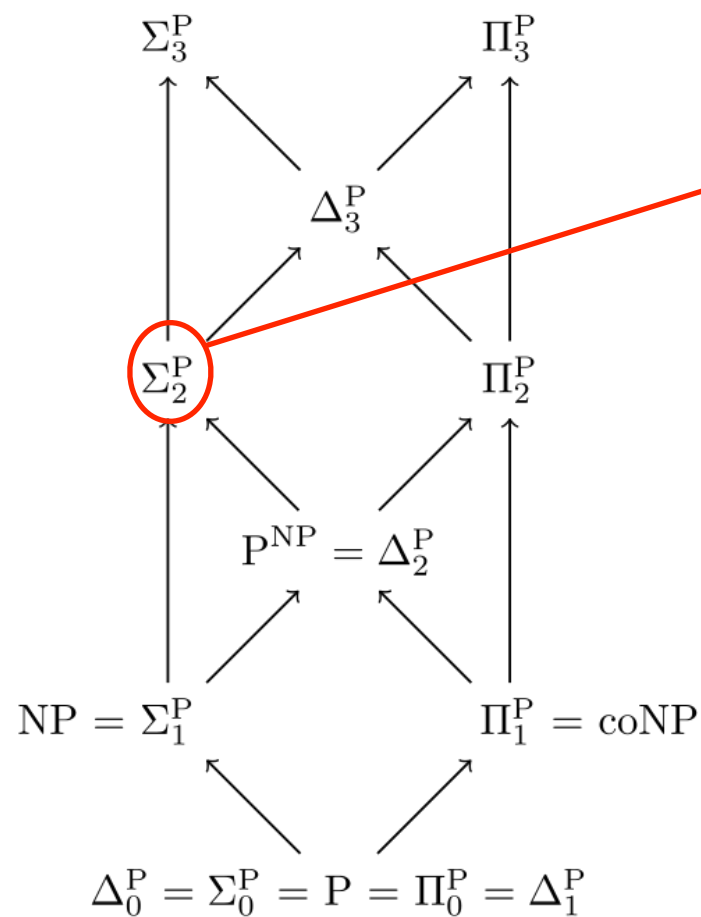
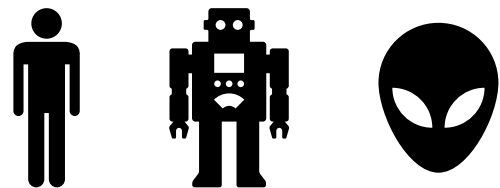
Now we generalize:

$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

# Polynomial Hierarchy, Part II

(via formal logic, directly)



free variables

Eg:

$$\langle \phi_1, k \rangle \in L \text{ iff } \exists \phi_2 \forall \alpha KLogEquiv(\phi_1, \phi_2, |\phi_2| \leq k, \alpha(\phi_1) = \alpha(\phi_2))$$

Now we generalize:

$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

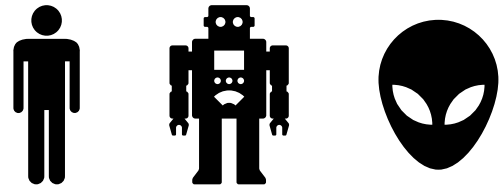
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \exists$  if  $j$  even;  $Q_i = \forall$  if  $j$  odd)

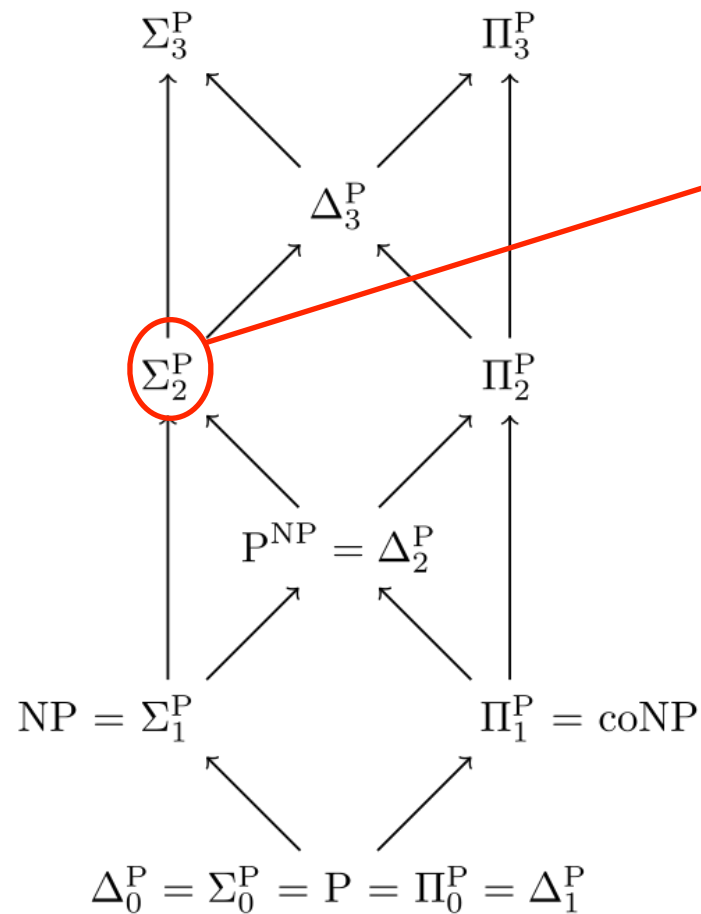



# Polynomial Hierarchy, Part II

(via formal logic, directly)



⋮




free variables

Eg:

$$\langle \phi_1, k \rangle \in L \text{ iff } \exists \phi_2 \forall \alpha KLogEquiv(\phi_1, \phi_2, |\phi_2| \leq k, \alpha(\phi_1) = \alpha(\phi_2))$$

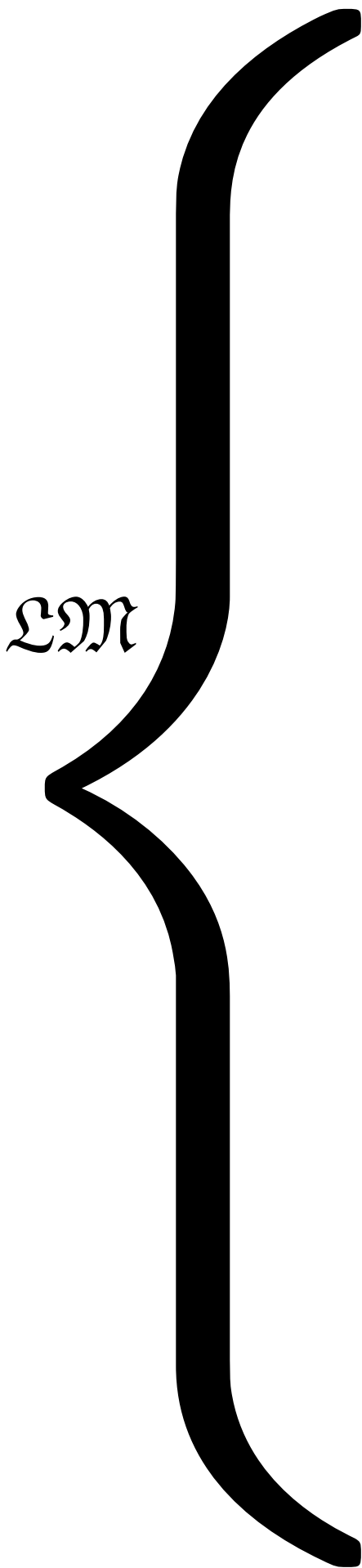
Now we generalize:

$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

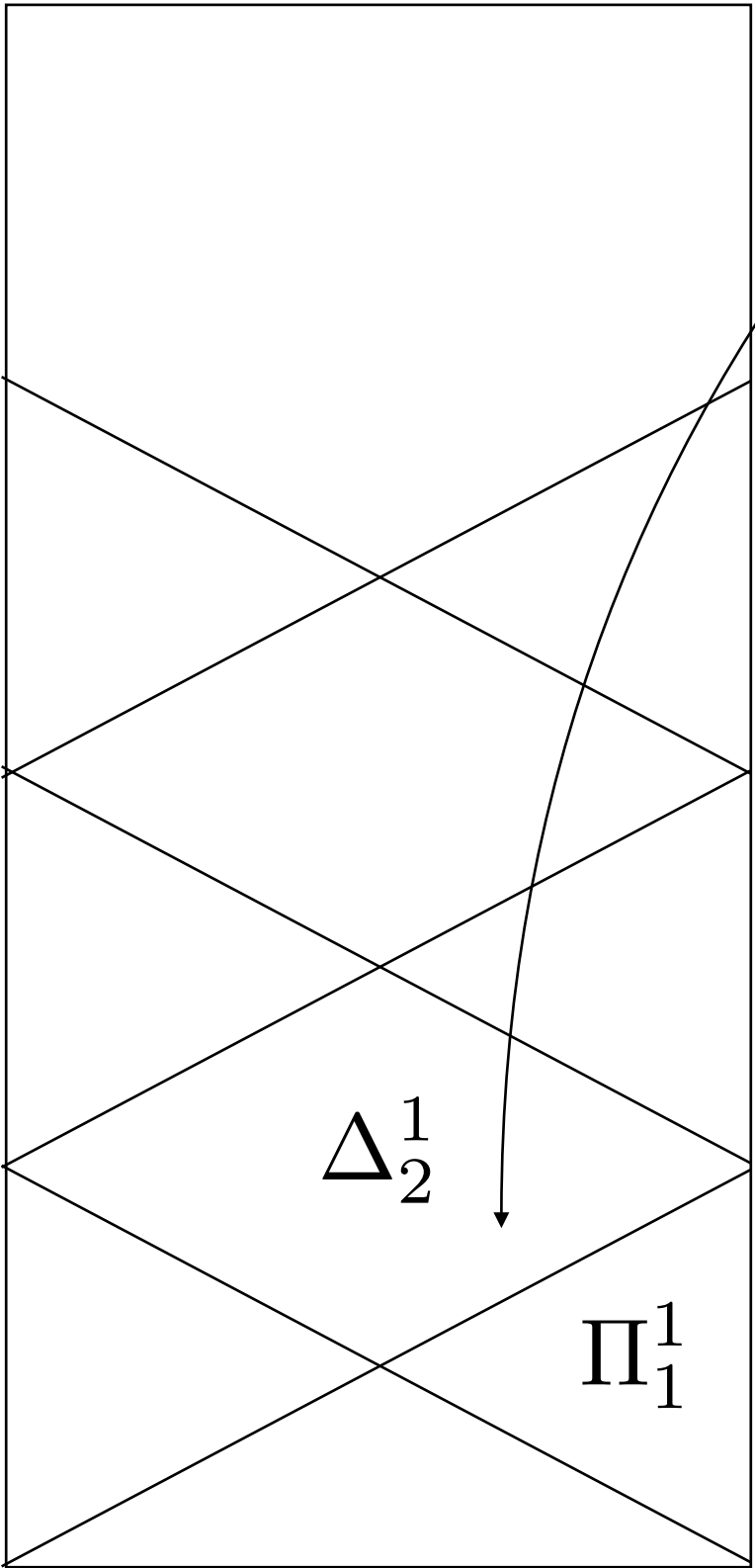
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \exists$  if  $j$  even;  $Q_i = \forall$  if  $j$  odd)

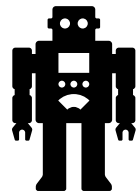


$\mathfrak{M}$

$\mathcal{A}^n \mathcal{H}$  (Analytic Hierarchy)



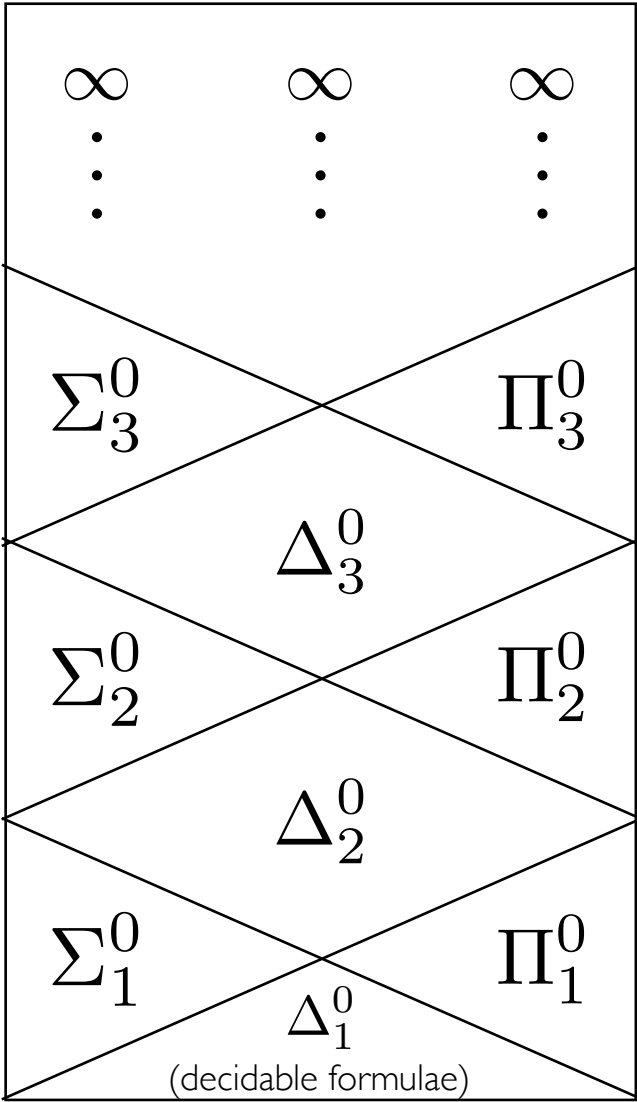
CogSci and AI need to say more about where AI falls/can fall in the landscape.



Infinite Time Turing Machines (ITTTMs)

Human Persons  
(according to Bringsjord)

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



Human Brains  
(according to Granger)



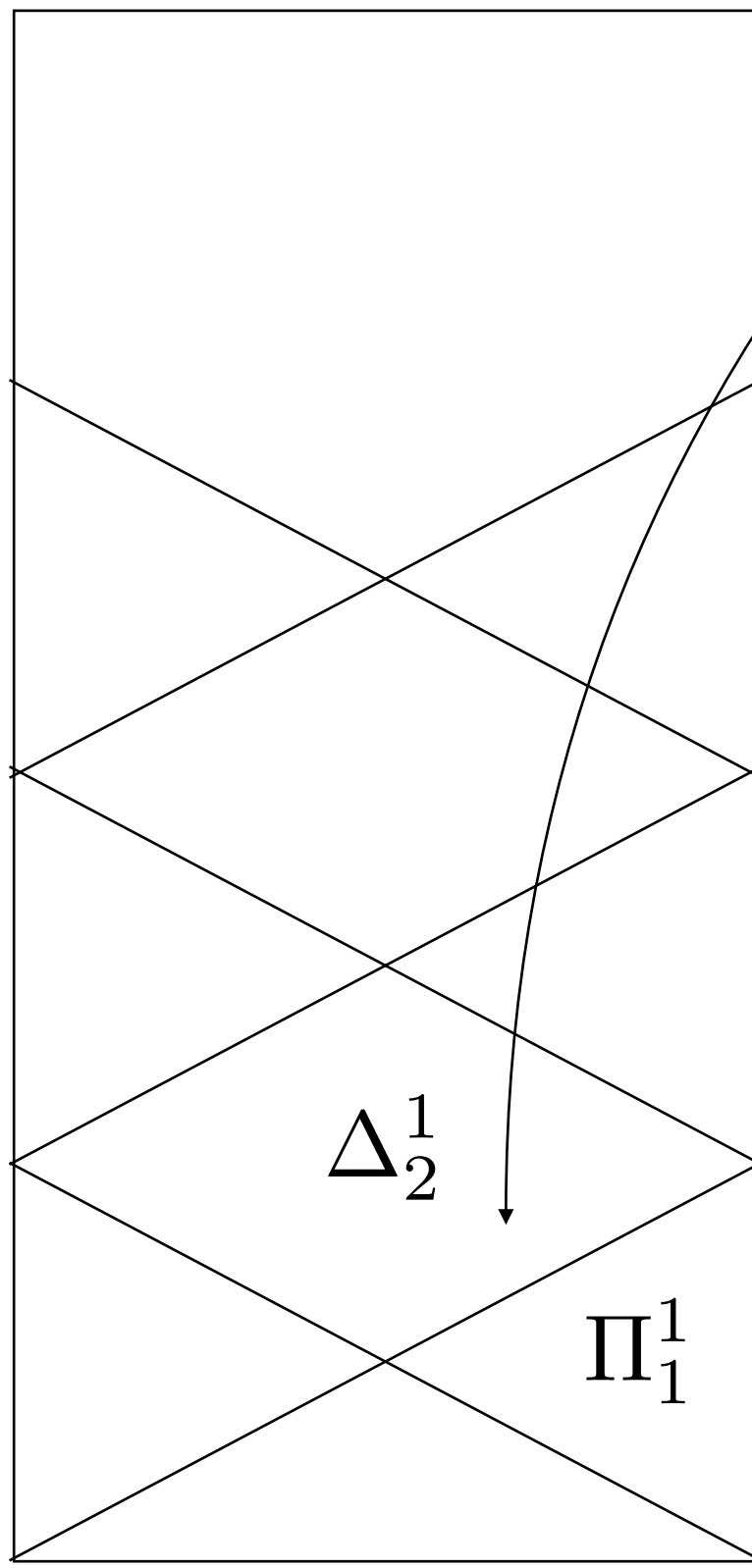
$\mathcal{CH}$  (Chomsky Hierarchy)

Turing Machines (TMs)
Linear Bounded Automata (LBAs)
Push Down Automata (PDAs)
Finite State Automata (FSAs)

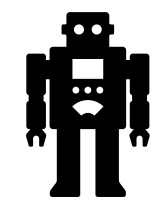


$\mathcal{LM}$

$\mathcal{A}^n \mathcal{H}$  (Analytic Hierarchy)



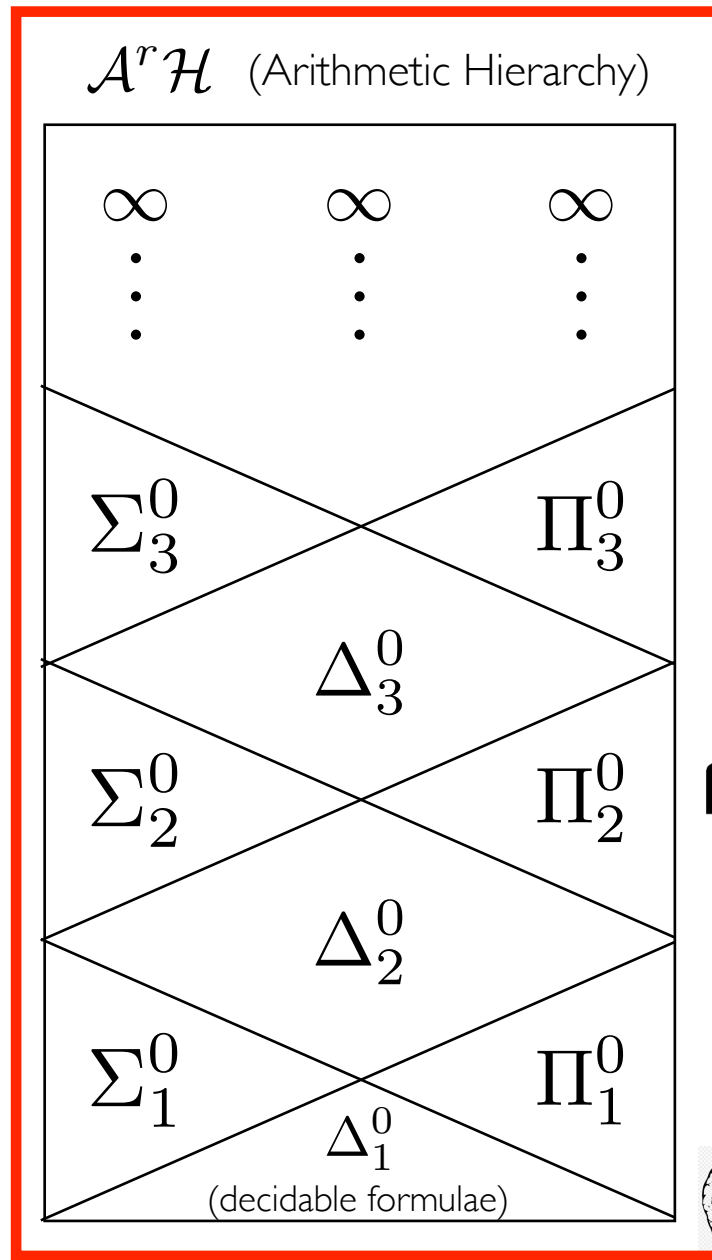
CogSci and AI need to say more about where AI falls/can fall in the landscape.



Infinite Time Turing Machines (ITTMs)

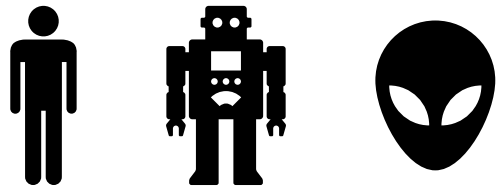
Human Persons  
(according to Bringsjord)

Human Brains  
(according to Granger)



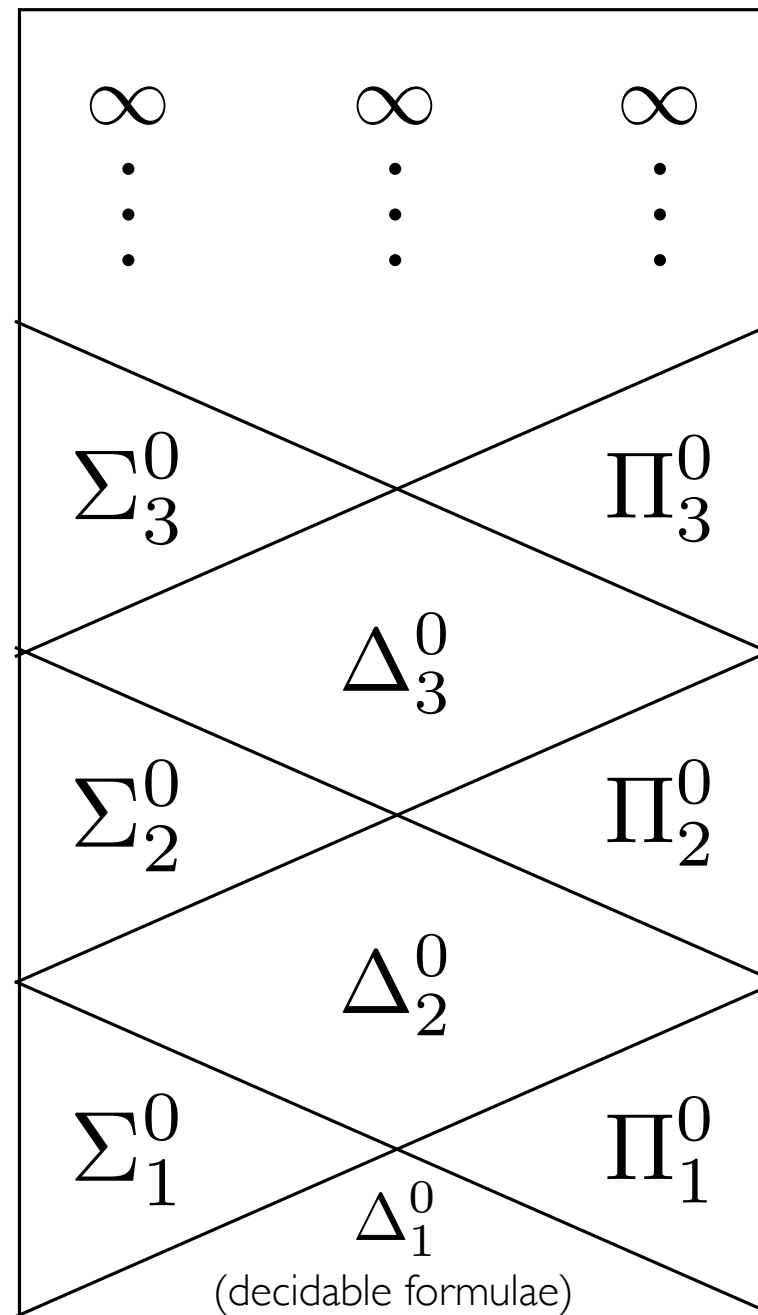
$\mathcal{CH}$  (Chomsky Hierarchy)

Turing Machines (TMs)
Linear Bounded Automata (LBAs)
Push Down Automata (PDAs)
Finite State Automata (FSAs)



$$\mathbf{2SAMEFUNC} := \{\mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))]\}$$

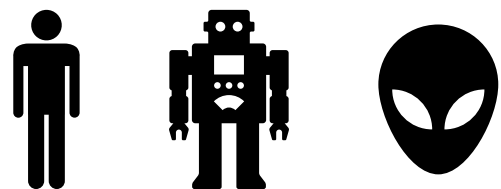
$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



semi-decidable

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

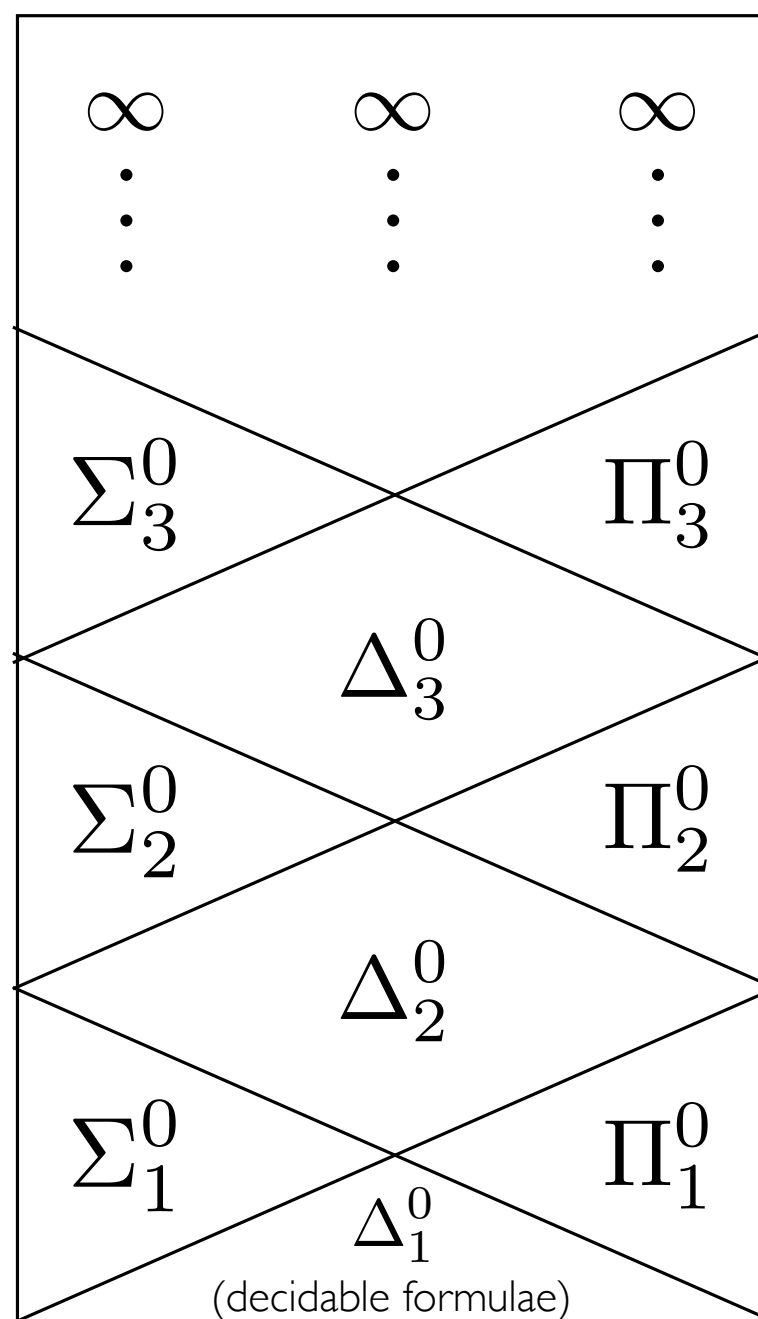
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

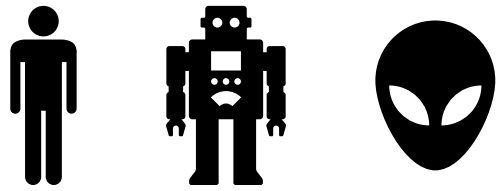
$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

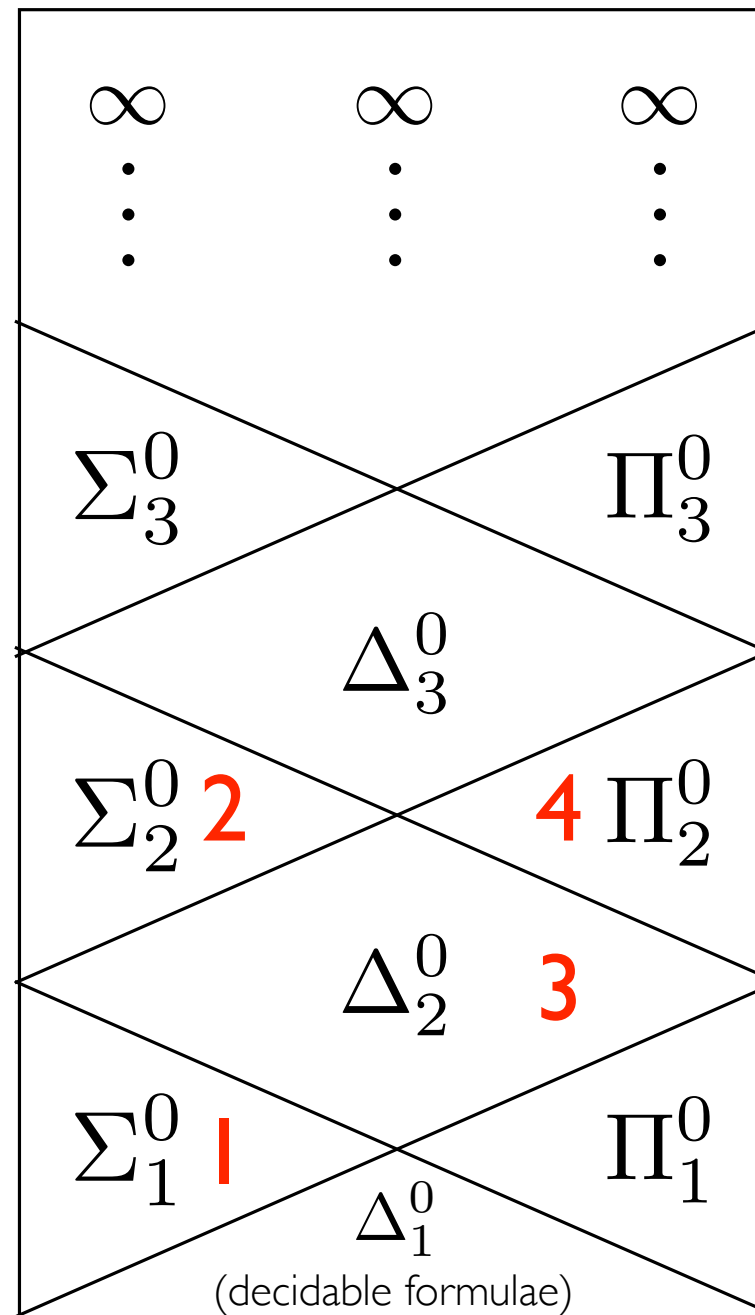
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{m_1, m_2 : \forall u \forall v [\exists k (\langle m_1, u \rangle : v, k \leftrightarrow \exists k' (\langle m_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)

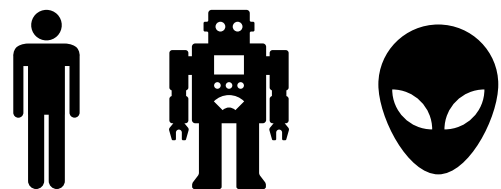


Let  $R$  be a Turing-decidable (= decidable, *simpliciter*) dyadic relation. Where is the set:  
 $\{x : \exists y R(x, y)\},$

1 2 3 or 4?

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

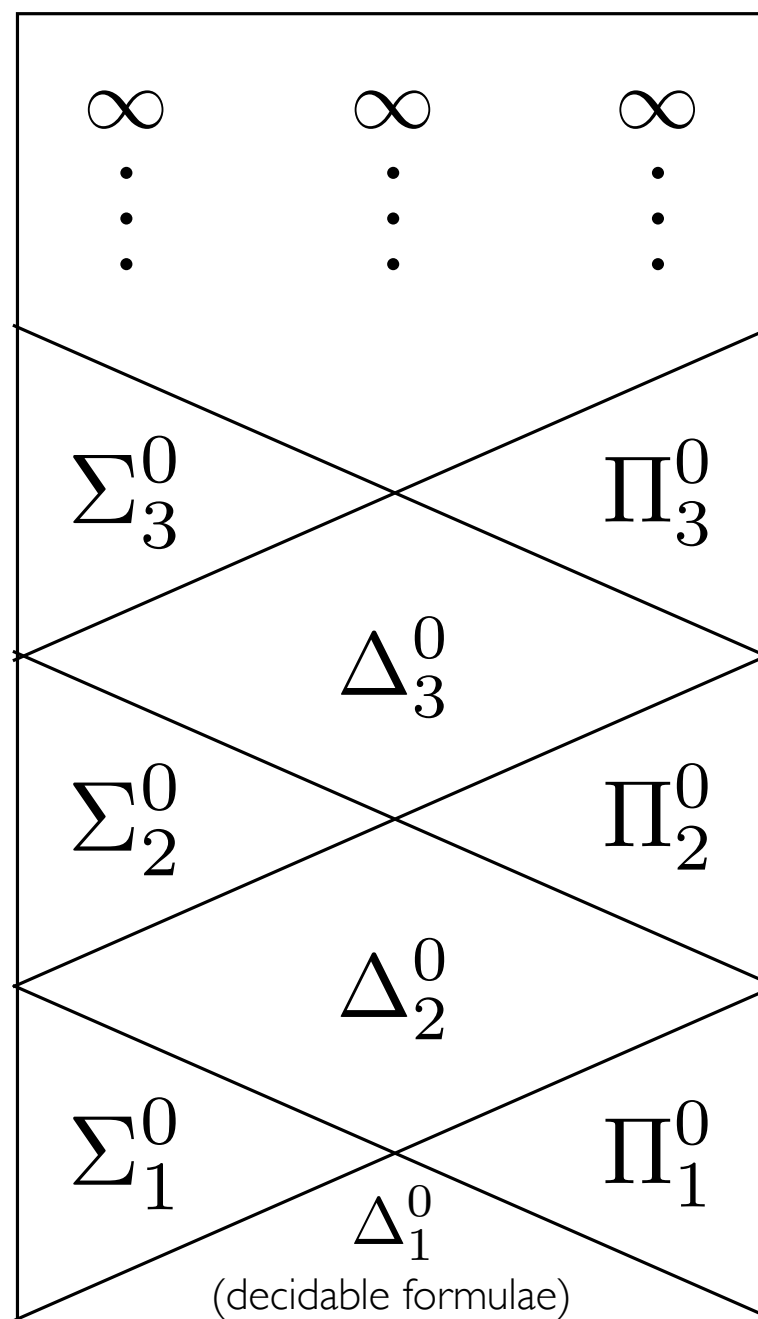
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{\mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))]\}$$

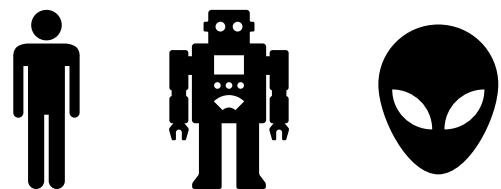
$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



semi-decidable

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

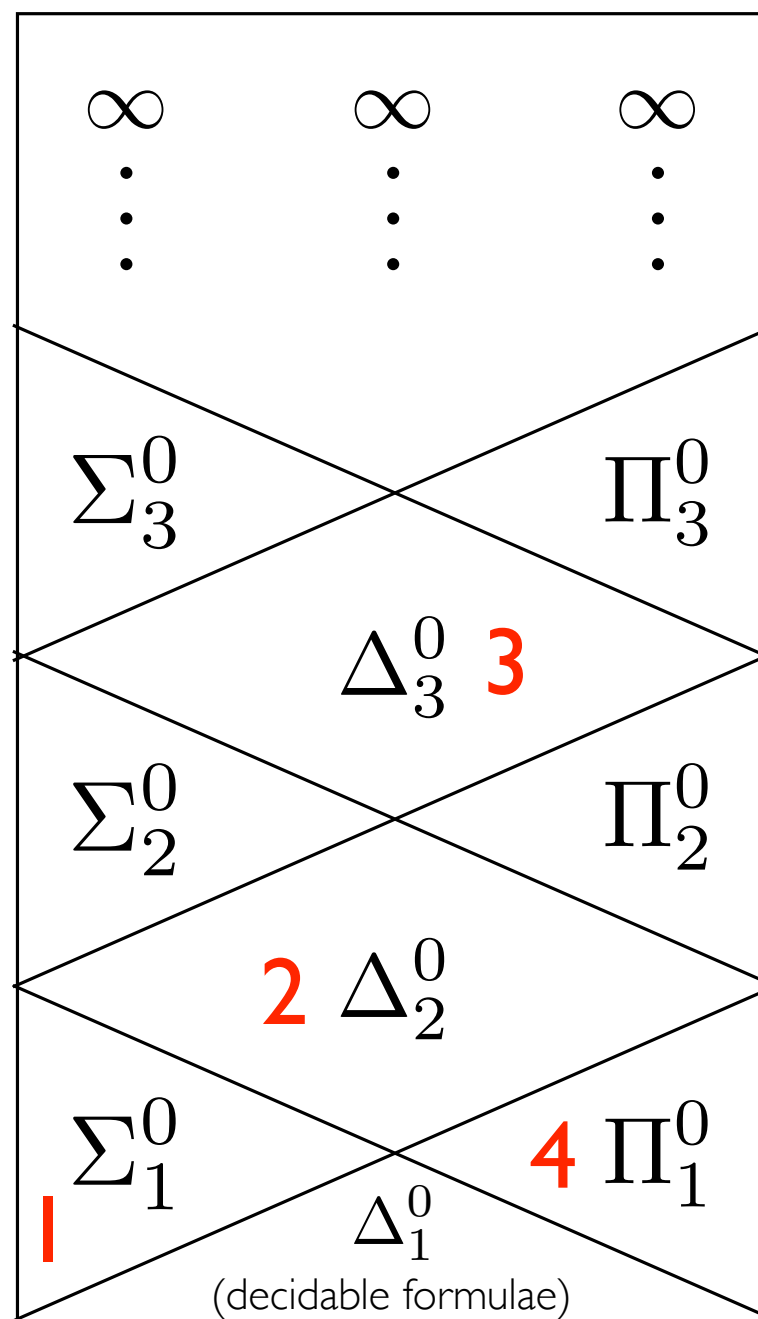
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{m_1, m_2 : \forall u \forall v [\exists k (\langle m_1, u \rangle : v, k \leftrightarrow \exists k' (\langle m_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



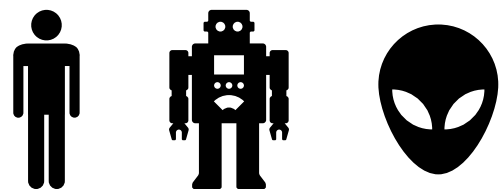
Let  $R$  be a Turing-decidable (= decidable, *simpliciter*) dyadic relation. Where is the set:  
 $\{x : \forall y R(x, y)\}$ ,

1 2 3 or 4?

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

# Arithmetic Hierarchy, Part I

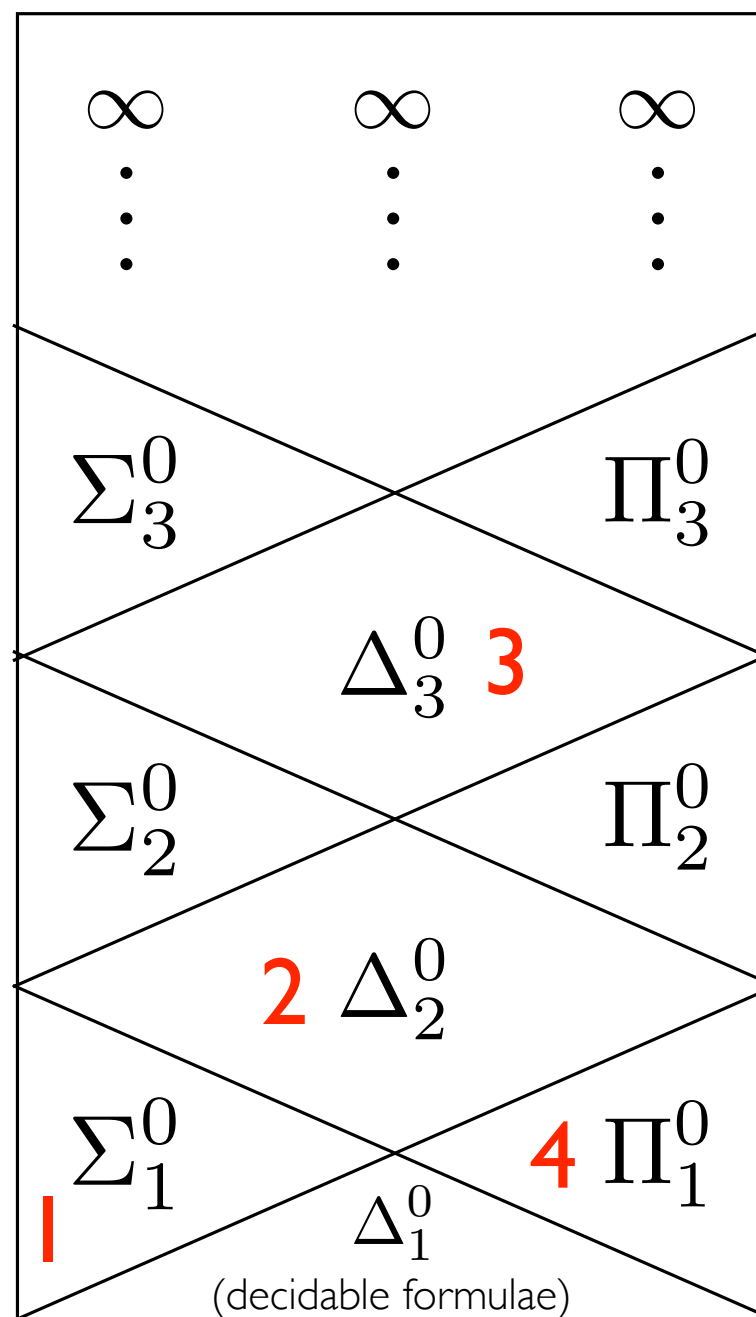




Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{m_1, m_2 : \forall u \forall v [\exists k (\langle m_1, u \rangle : v, k \leftrightarrow \exists k' (\langle m_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



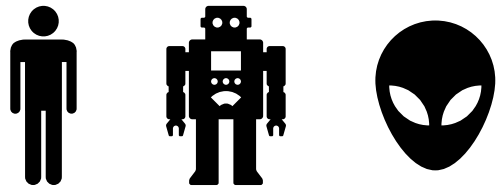
Let  $R$  be a Turing-decidable (= decidable, *simpliciter*) dyadic relation. Where is the set:  
 $\{x : \forall y R(x, y)\}$ ,

1 2 3 or 4?

$x \in \Sigma_i$  iff  $\exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$   
 ( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

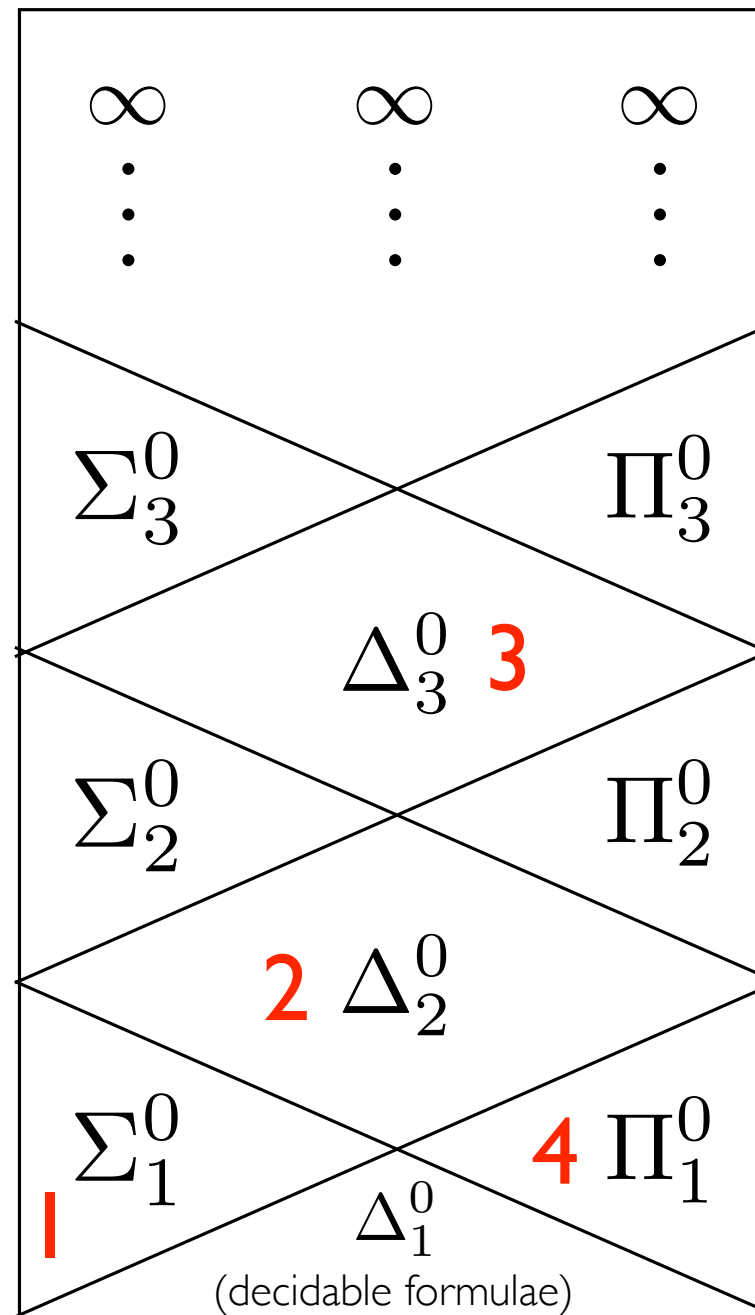
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{m_1, m_2 : \forall u \forall v [\exists k (\langle m_1, u \rangle : v, k \leftrightarrow \exists k' (\langle m_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



Let  $R$  be a Turing-decidable (= decidable, *simpliciter*) dyadic relation. Where is the set:  
 $\{x : \forall y R(x, y)\}$ ,

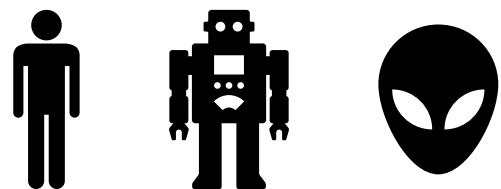
1 2 3 or 4?

$x \in \Sigma_i$  iff  $\exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$   
 ( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

$x \in \Pi_i$  iff  $\exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$   
 ( $Q_i = \exists$  if  $j$  even;  $Q_i = \forall$  if  $j$  odd)

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

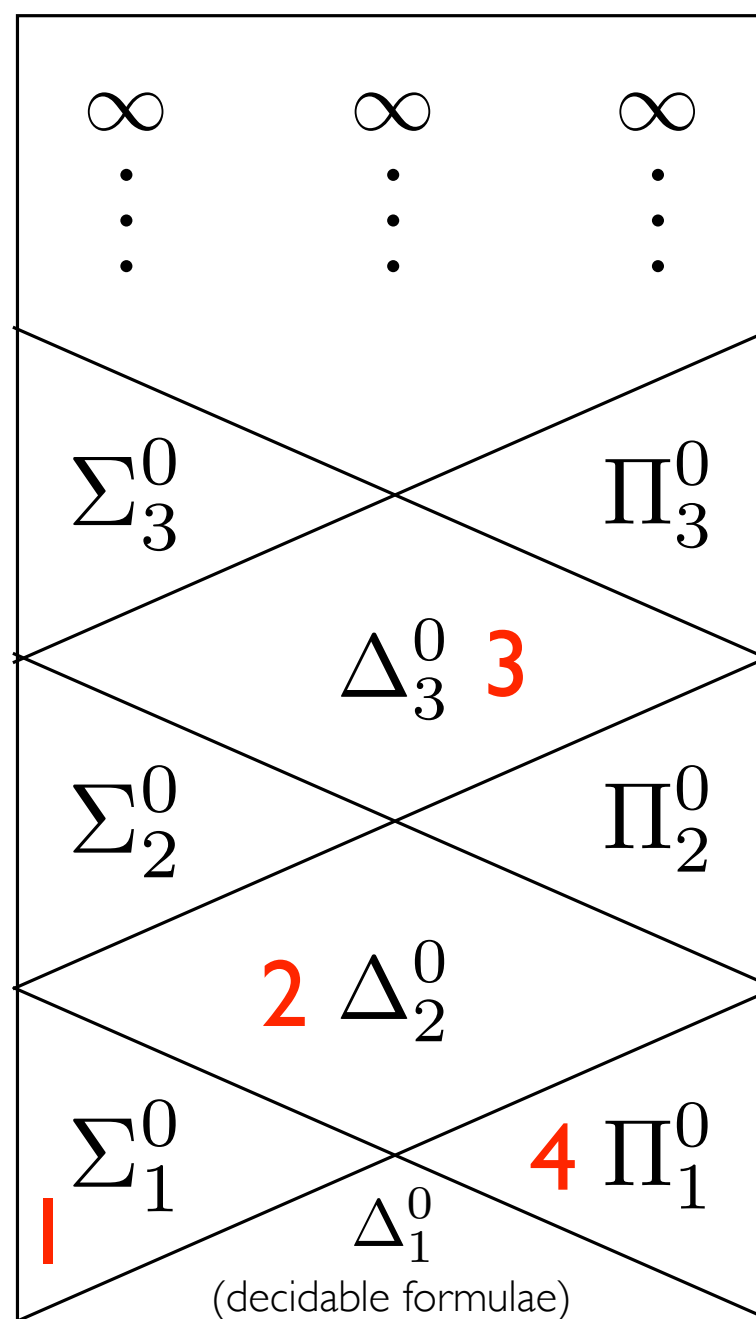
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{\mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



Let  $R$  be a Turing-decidable (= decidable, *simpliciter*) dyadic relation. Where is the set:  
 $\{x : \forall y R(x, y)\}$ ,

1 2 3 or 4?

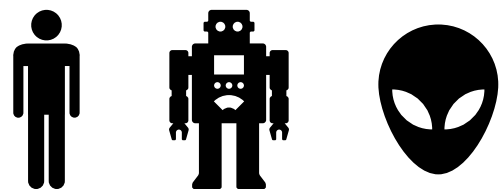
$x \in \Sigma_i$  iff  $\exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$   
 ( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

$x \in \Pi_i$  iff  $\exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$   
 ( $Q_i = \exists$  if  $j$  even;  $Q_i = \forall$  if  $j$  odd)

Try your hand at classifying! ...

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

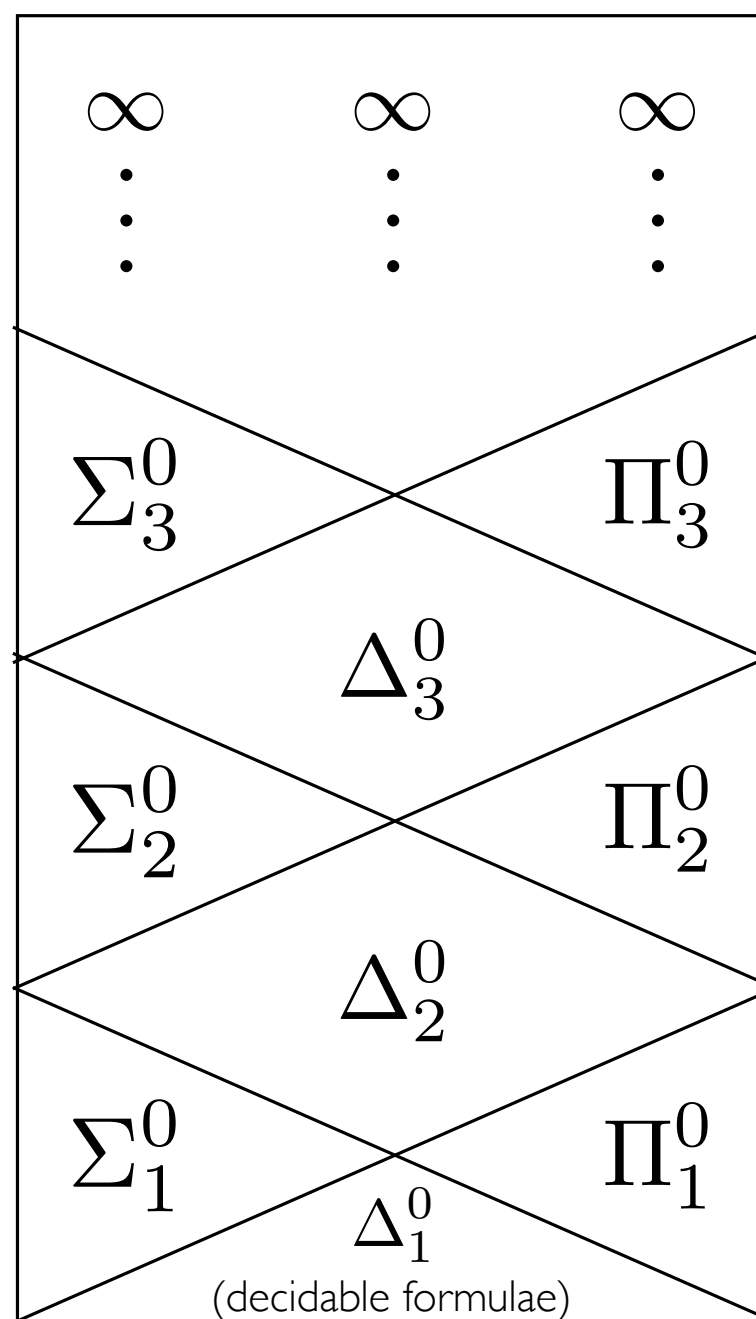
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{m_1, m_2 : \forall u \forall v [\exists k (\langle m_1, u \rangle : v, k \leftrightarrow \exists k' (\langle m_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

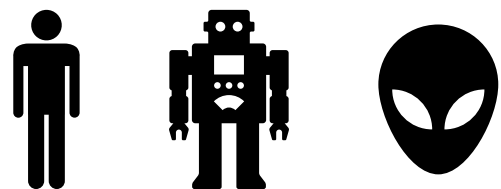
Try your hand at classifying! ...

semi-decidable



$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

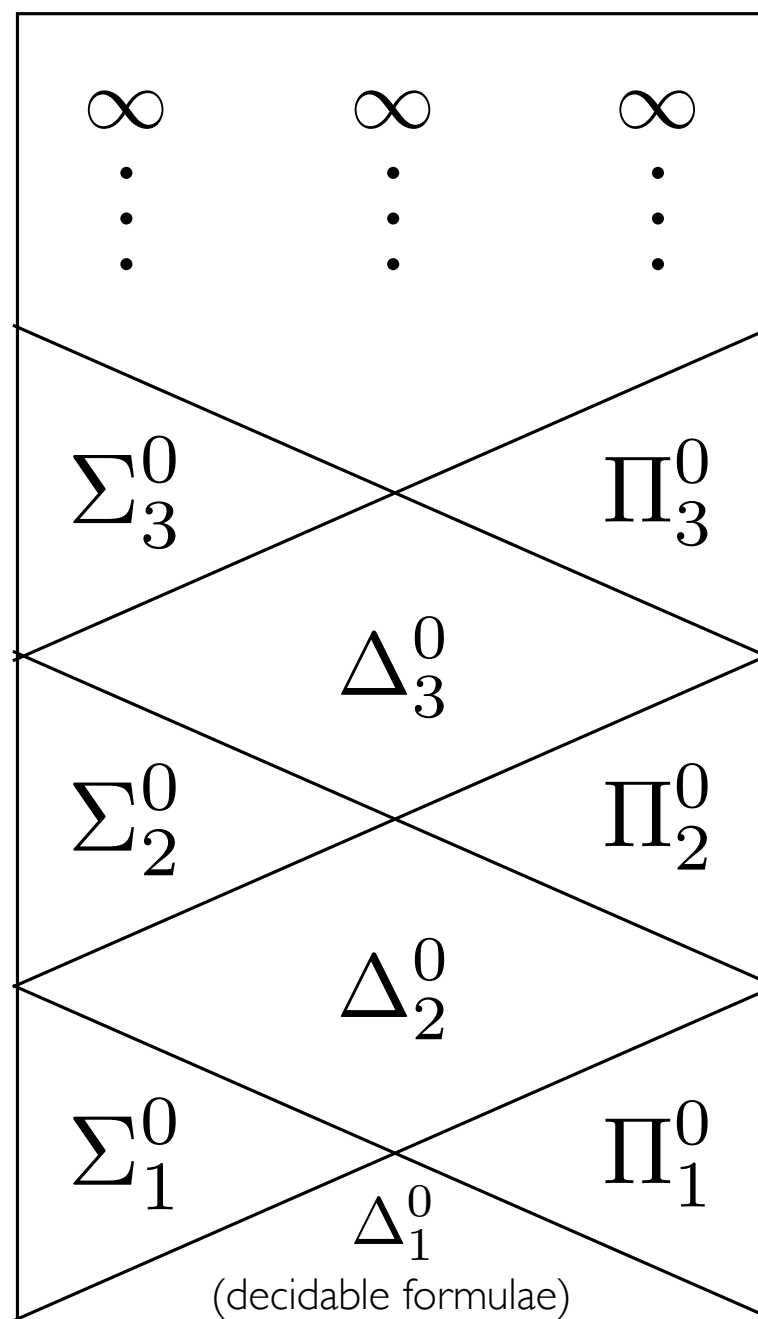
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{\mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

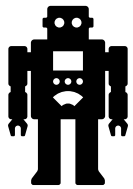
( $Q_i = \exists$  if  $j$  even;  $Q_i = \forall$  if  $j$  odd)

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

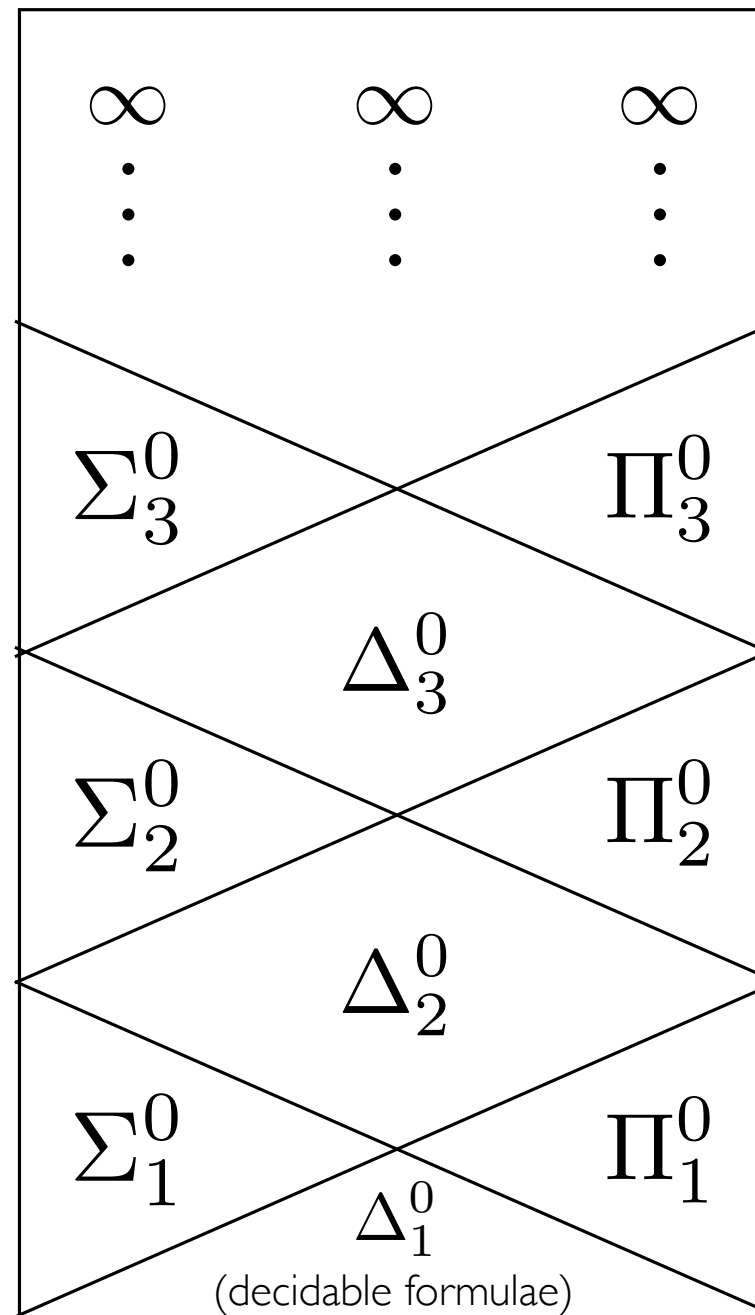
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMPLEFUNC} := \{\mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

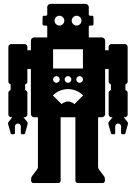
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

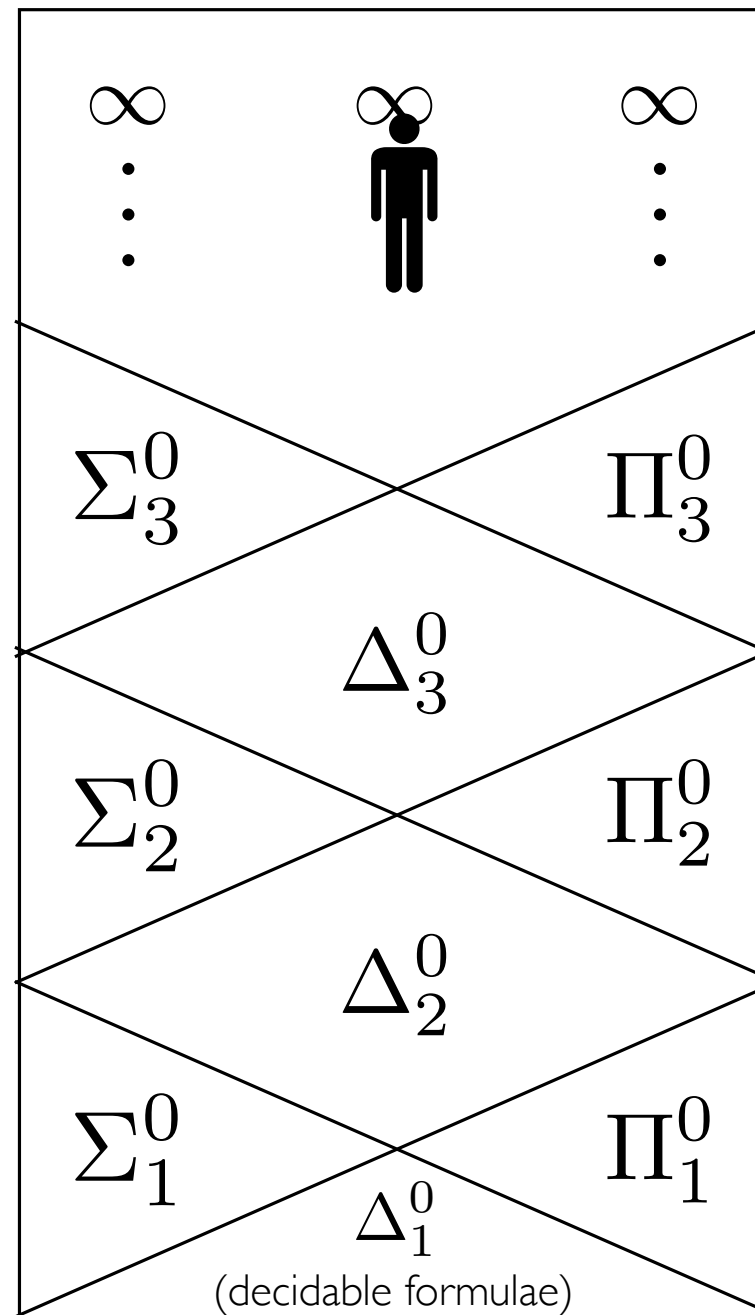
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{m_1, m_2 : \forall u \forall v [\exists k (\langle m_1, u \rangle : v, k \leftrightarrow \exists k' (\langle m_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

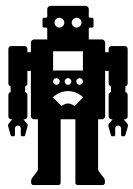
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

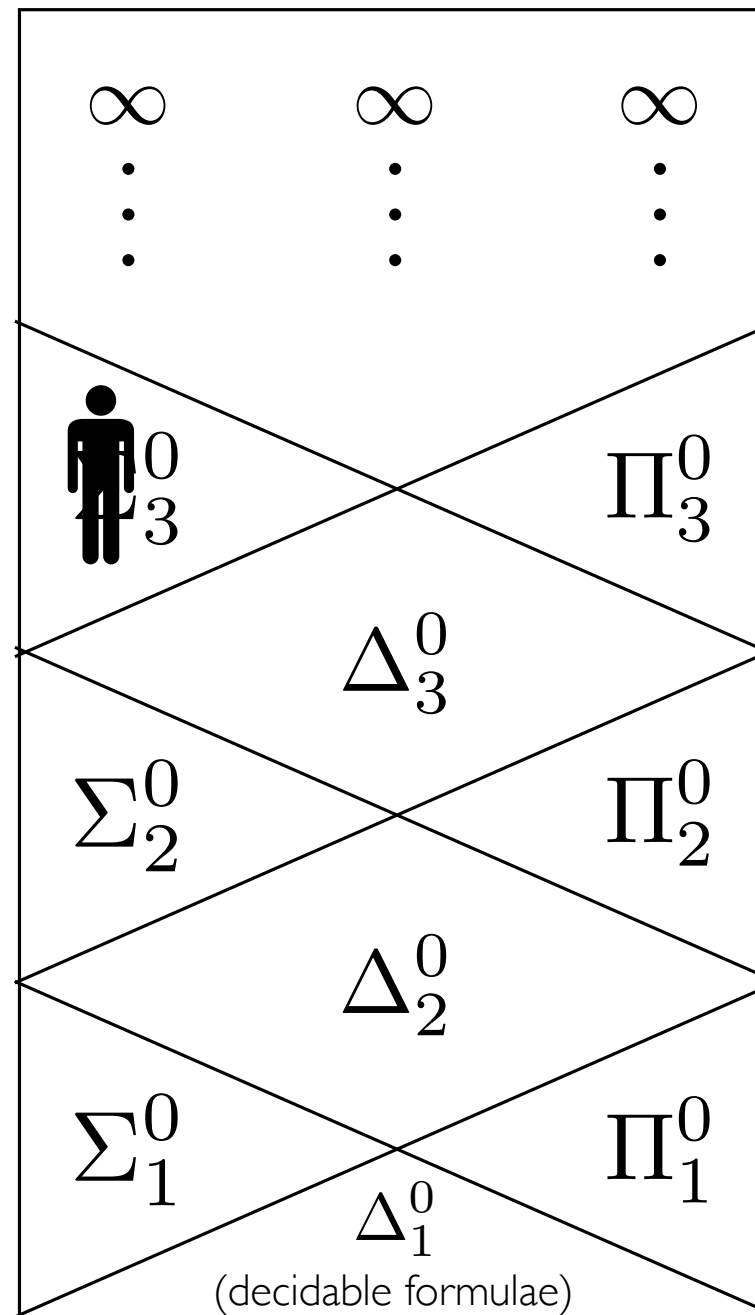
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

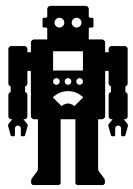
Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

# Arithmetic Hierarchy, Part I

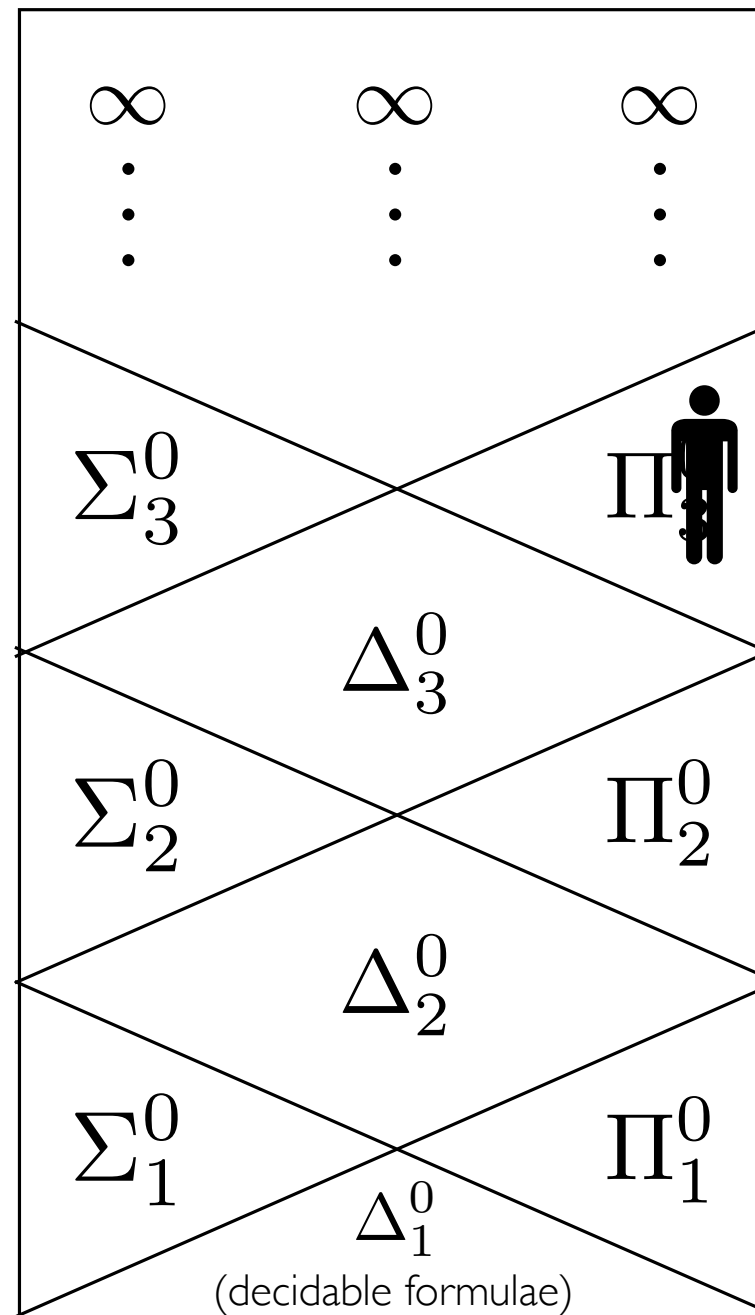




Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

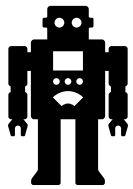
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

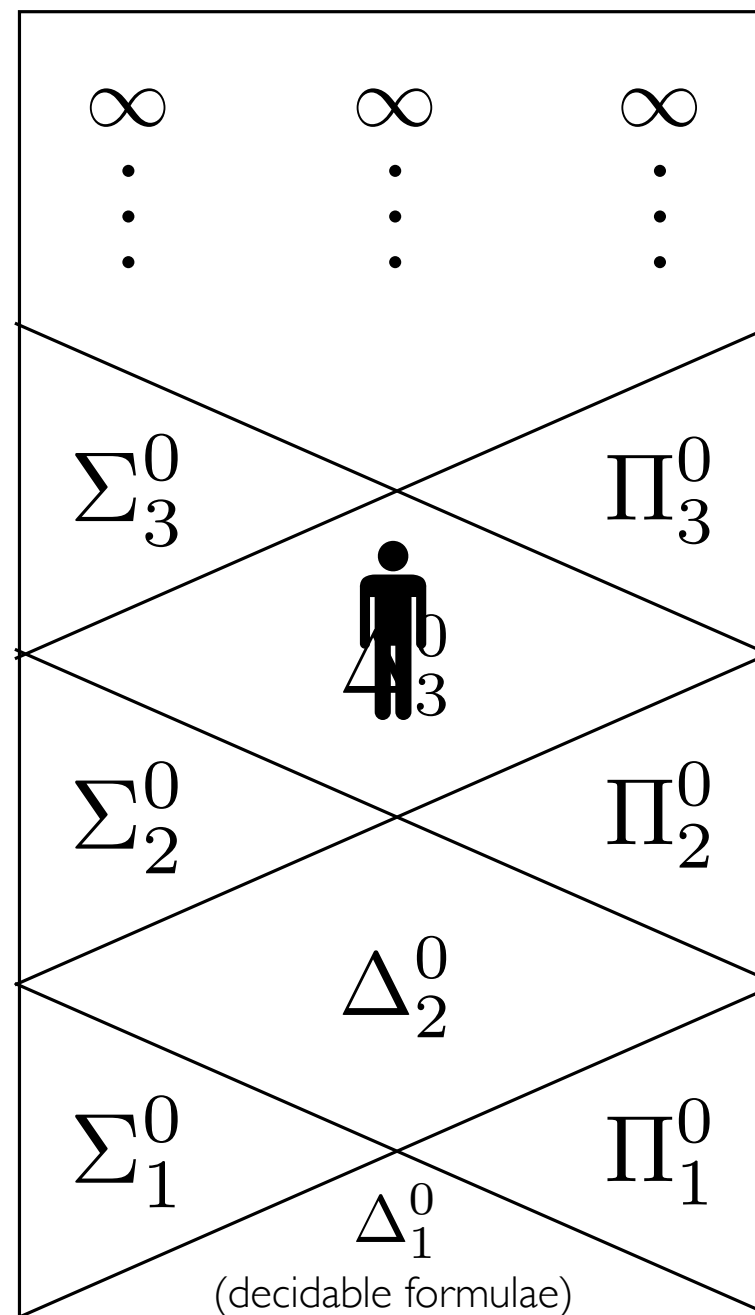
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{m_1, m_2 : \forall u \forall v [\exists k (\langle m_1, u \rangle : v, k \leftrightarrow \exists k' (\langle m_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

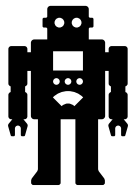
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $m$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

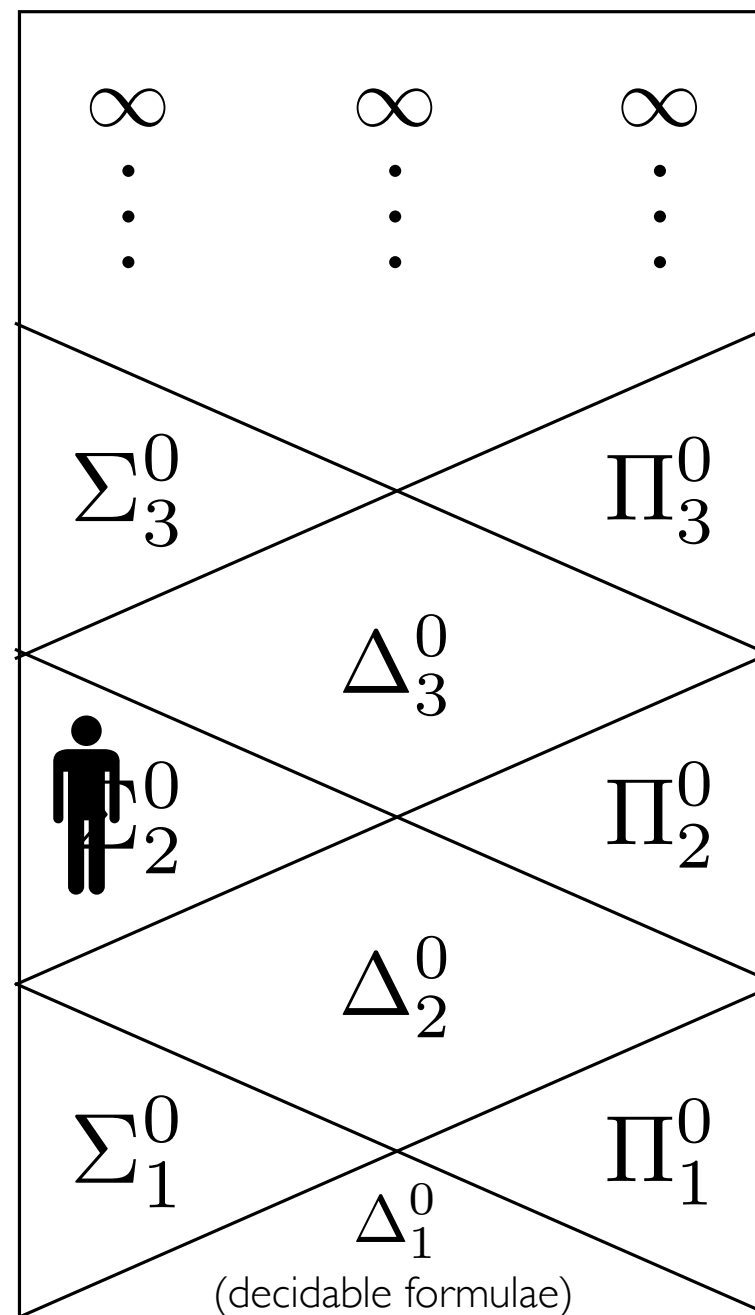
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

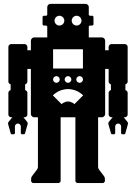
Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

semi-decidable

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

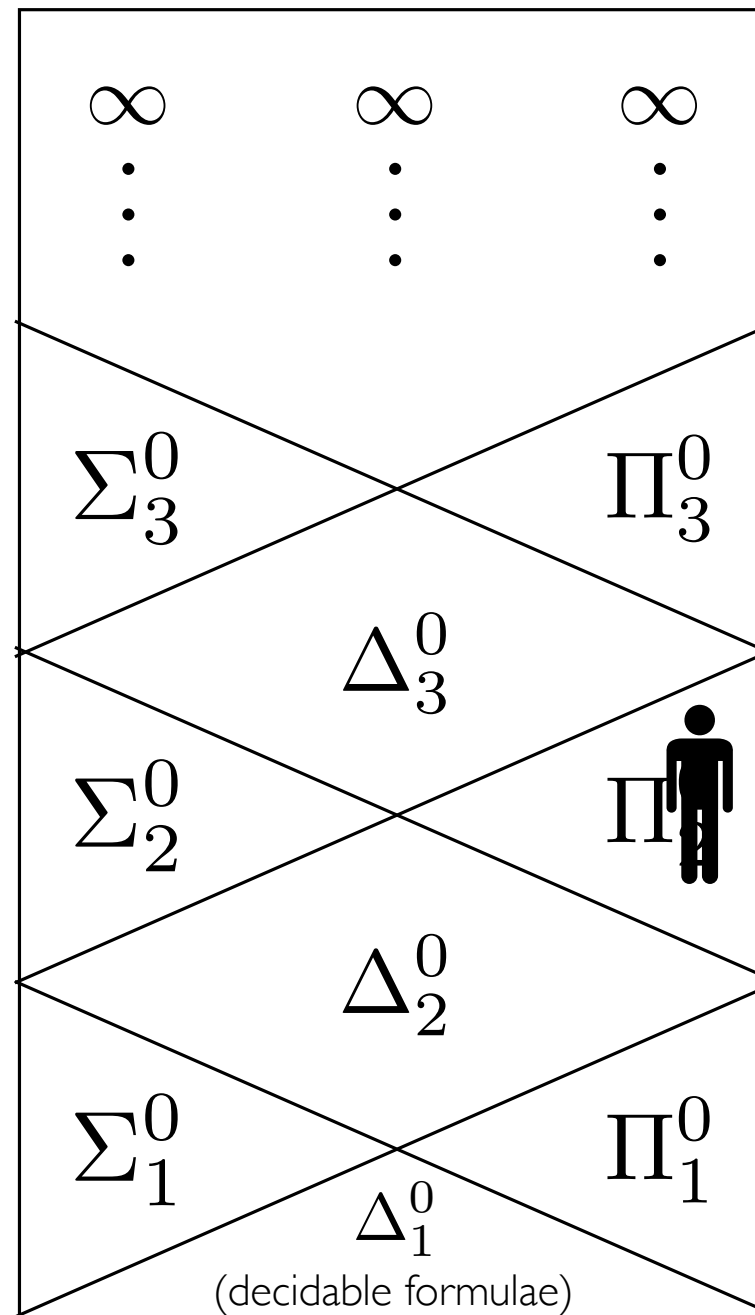
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

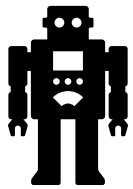
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

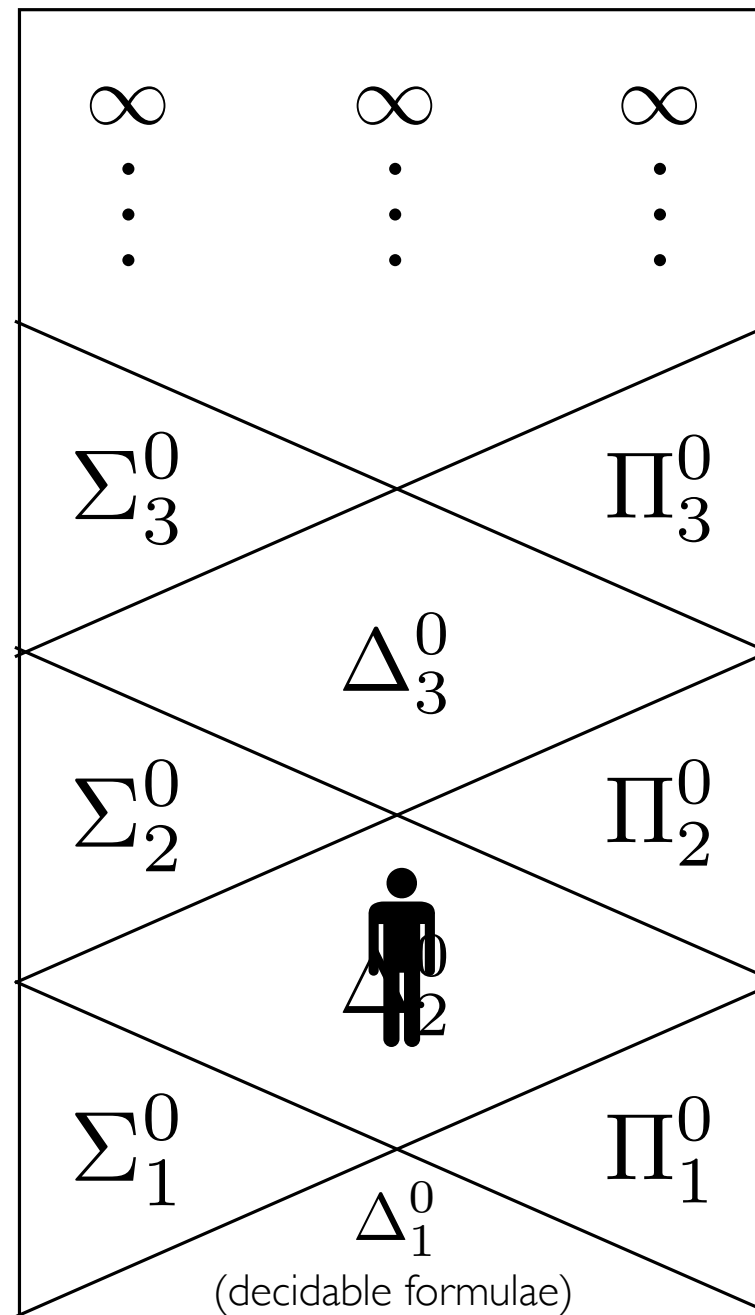
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{\mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r\mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

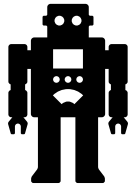
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

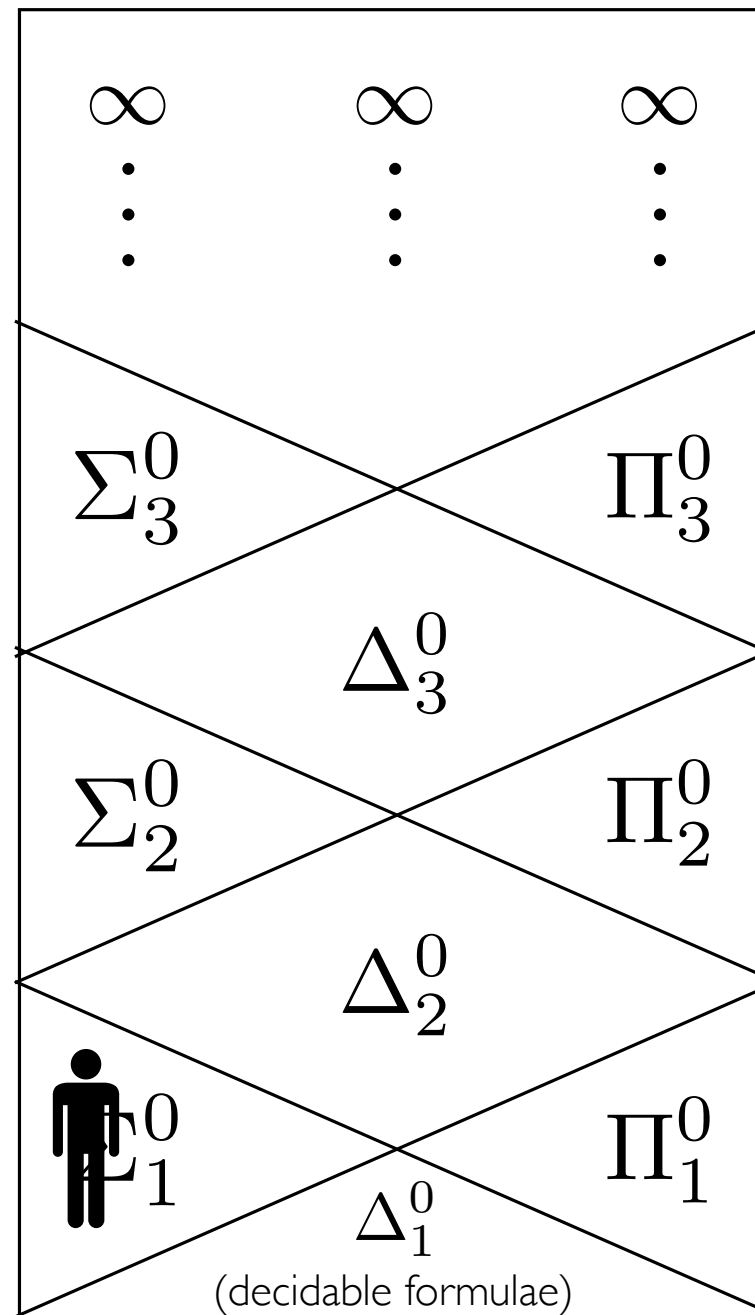
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$2\text{SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

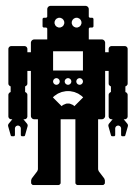
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

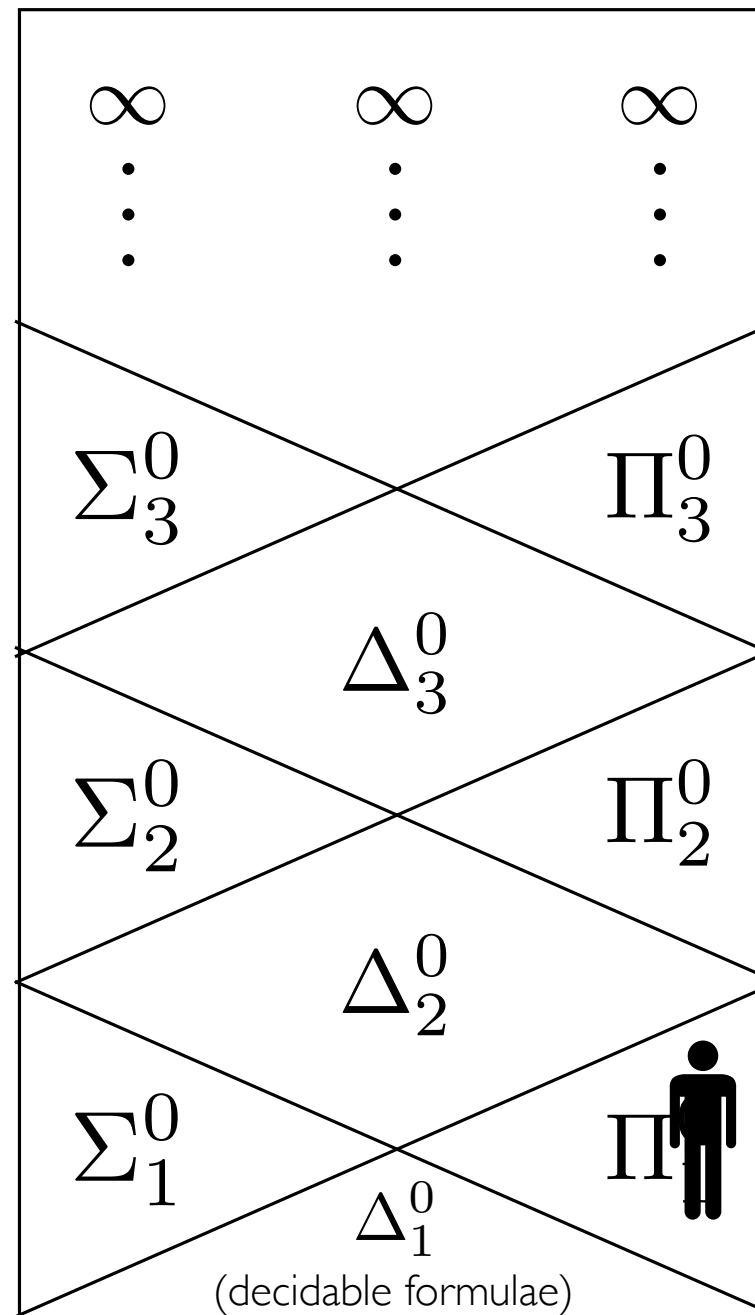
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

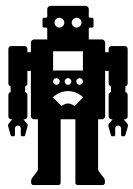
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

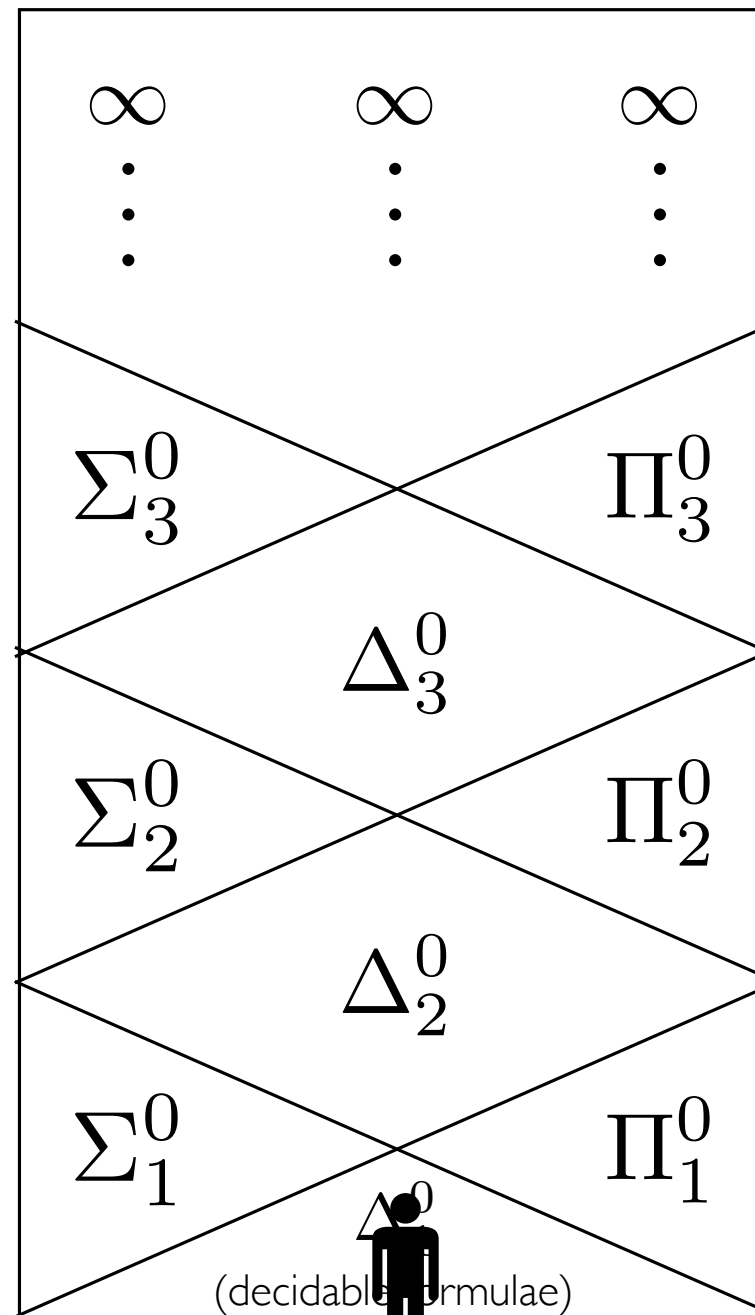
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

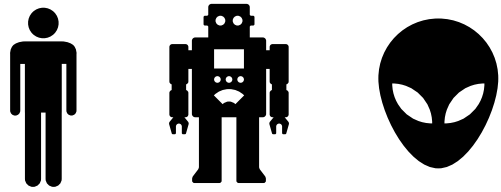
Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

# Arithmetic Hierarchy, Part I

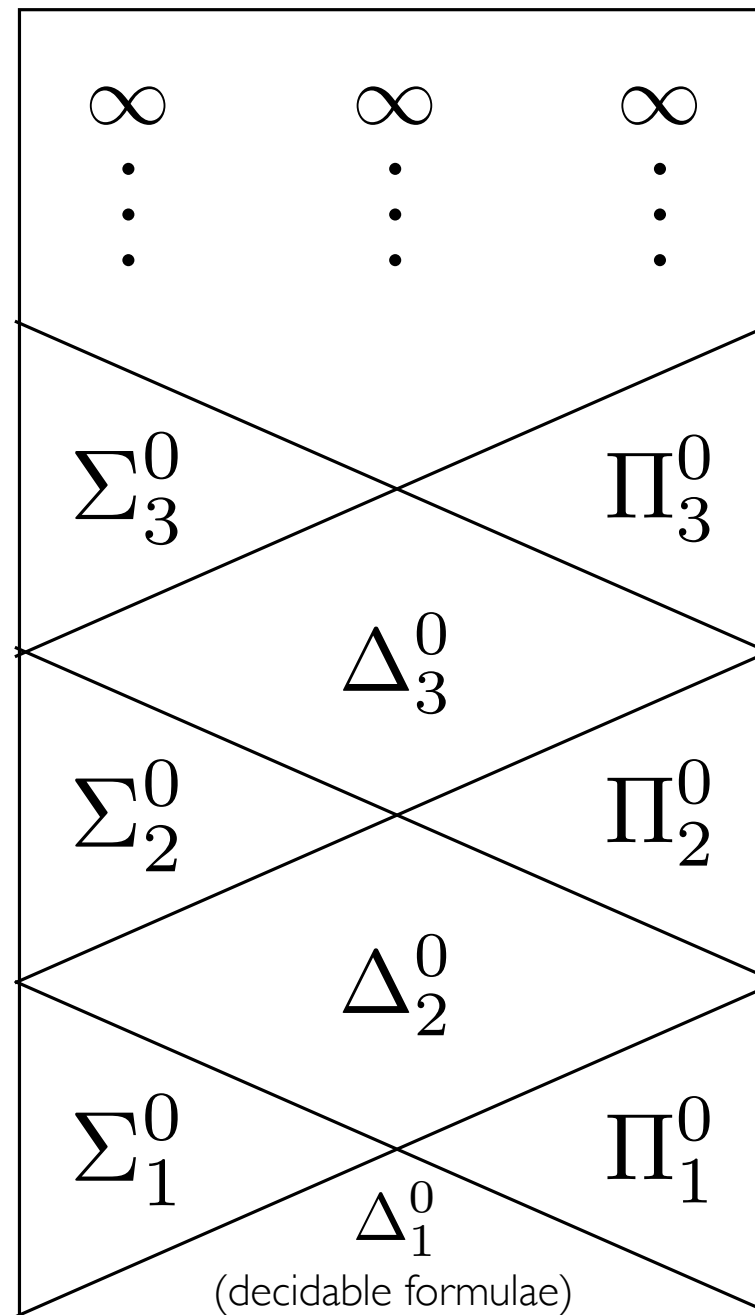




Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

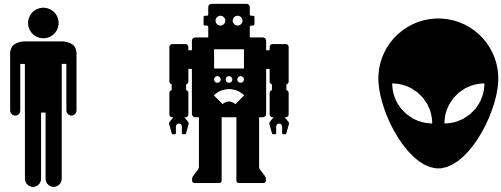
( $Q_i = \exists$  if  $j$  even;  $Q_i = \forall$  if  $j$  odd)

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

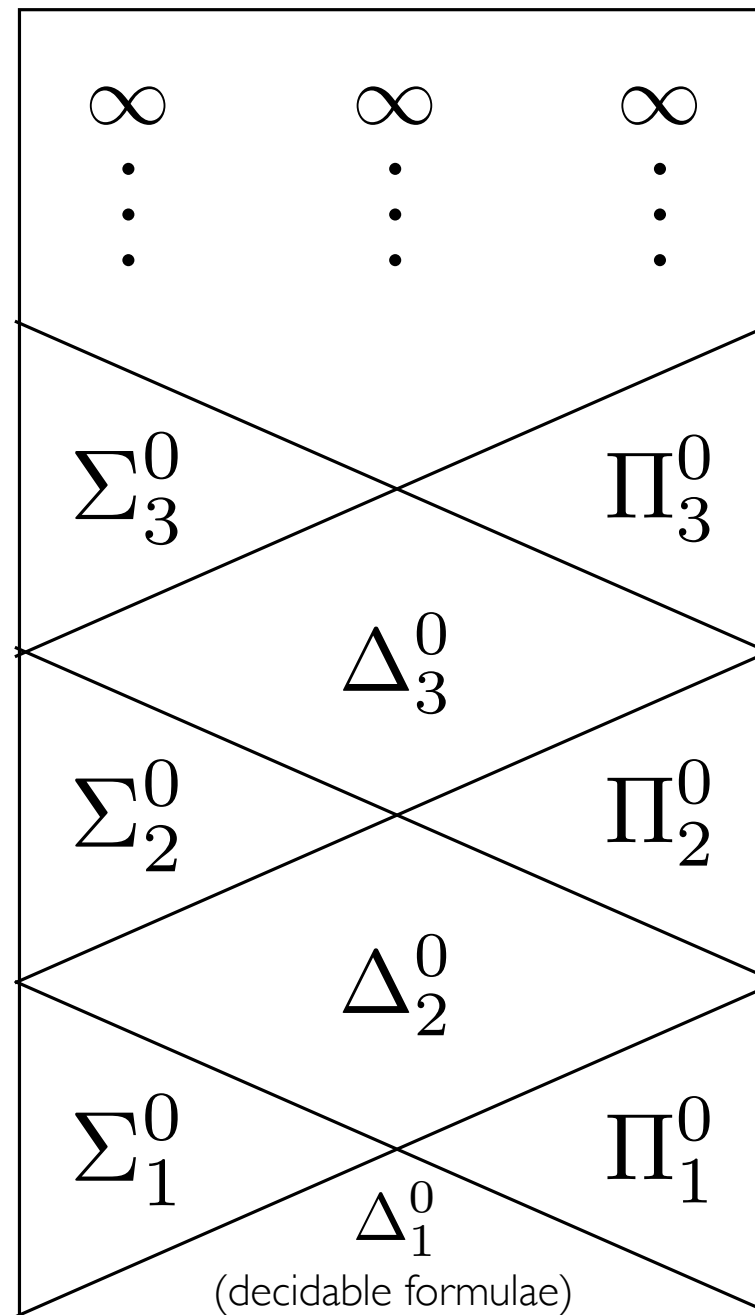
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

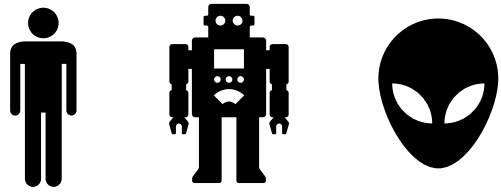
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

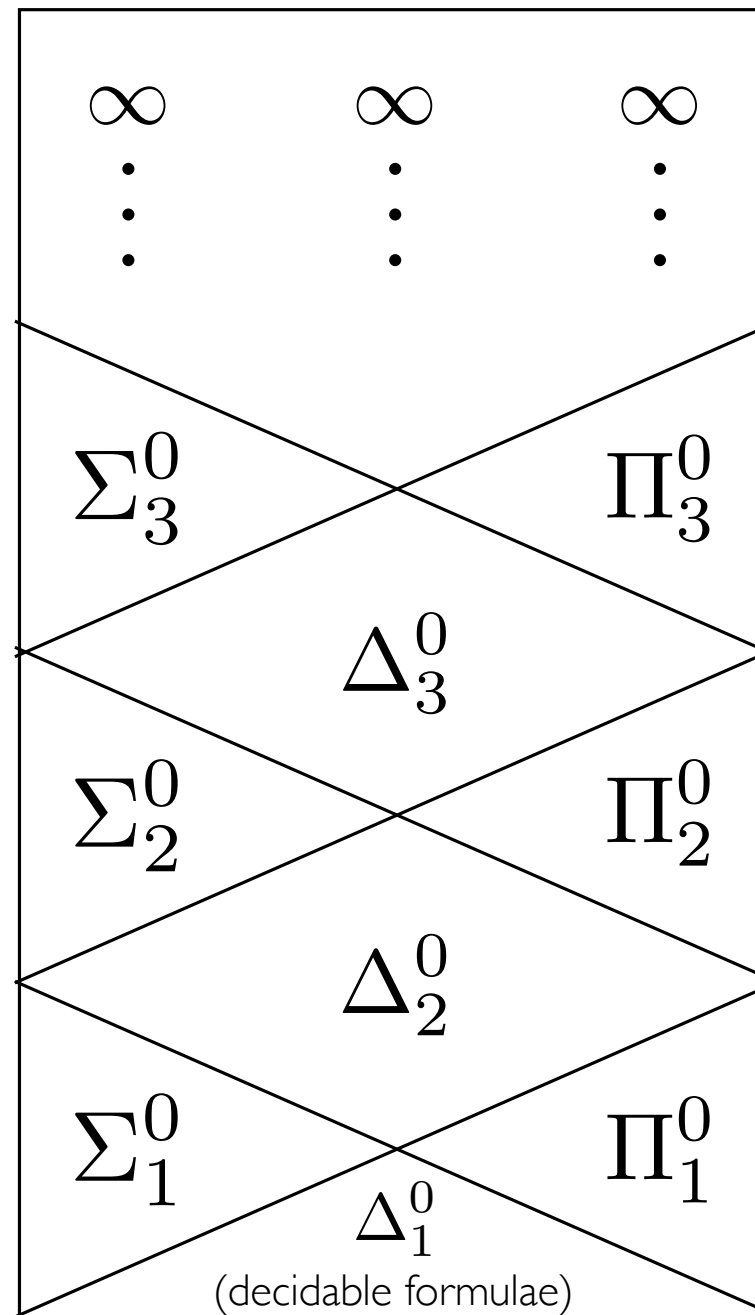
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{ \mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))] \}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \forall \text{ if } i \text{ even; } Q_i = \exists \text{ if } i \text{ odd})$$

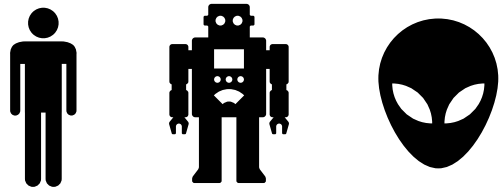
$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i) \\ (Q_i = \exists \text{ if } j \text{ even; } Q_i = \forall \text{ if } j \text{ odd})$$

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

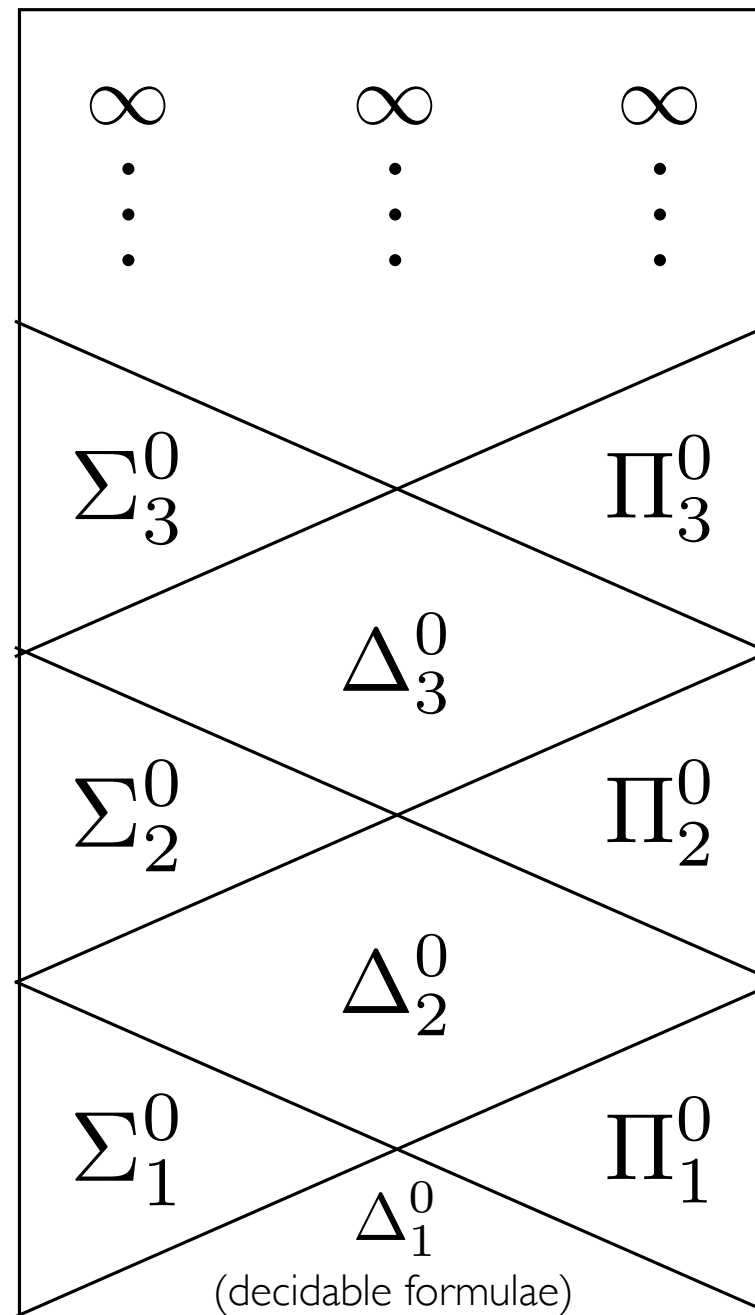
# Arithmetic Hierarchy, Part I



Can you see the carryover from PH?

$$\mathbf{2SAMEFUNC} := \{\mathbf{m}_1, \mathbf{m}_2 : \forall u \forall v [\exists k (\langle \mathbf{m}_1, u \rangle : v, k \leftrightarrow \exists k' (\langle \mathbf{m}_2, u \rangle : v, k'))]\}$$

$\mathcal{A}^r \mathcal{H}$  (Arithmetic Hierarchy)



$$x \in \Sigma_i \text{ iff } \exists R \exists y_1 \forall y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \forall$  if  $i$  even;  $Q_i = \exists$  if  $i$  odd)

$$x \in \Pi_i \text{ iff } \exists R \forall y_1 \exists y_2 \cdots Q_i y_i R(x, y_1, y_2, \dots, y_i)$$

( $Q_i = \exists$  if  $j$  even;  $Q_i = \forall$  if  $j$  odd)

Try your hand at classifying! ...

From Kleene: The set to be classified,  $\mathcal{K}$ , consists of all those inputs to a given Turing machine  $\mathbf{m}$  that results in this machine halting after some number of steps.

The set to be classified is the set of all pairs of programs  $P_1$  and  $P_2$  s.t. both compute exactly the same functions.

$$\Delta_1^0 = \Sigma_1^0 \cap \Pi_1^0$$

# Arithmetic Hierarchy, Part I

**The Four Steps ...**

# The PAID Problem (Level I):

# The PAID Problem (Level I):

$\forall x : \text{Agents}$

# The PAID Problem (Level I):

$\forall x : \text{Agents}$

$\text{Powerful}(x) + \text{Autonomous}(x) + \text{Intelligent}(x) \Rightarrow \text{Dangerous}(x)$



# The PAID Problem (Level I):

$\forall x : \text{Agents}$

$\text{Powerful}(x) + \text{Autonomous}(x) + \text{Intelligent}(x) \Rightarrow \text{Dangerous}(x)$



# The PAID Problem (Level I):

$\forall x : \text{Agents}$

$\text{Powerful}(x) + \text{Autonomous}(x) + \text{Intelligent}(x) \Rightarrow \text{Dangerous}(x)$

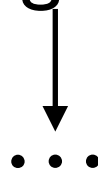


$$u(\text{AIA}_i(\pi_j)) > \tau^+ \in \mathbb{Z} \text{ or } \tau^- \in \mathbb{Z}$$

# The PAID Problem (Level I):

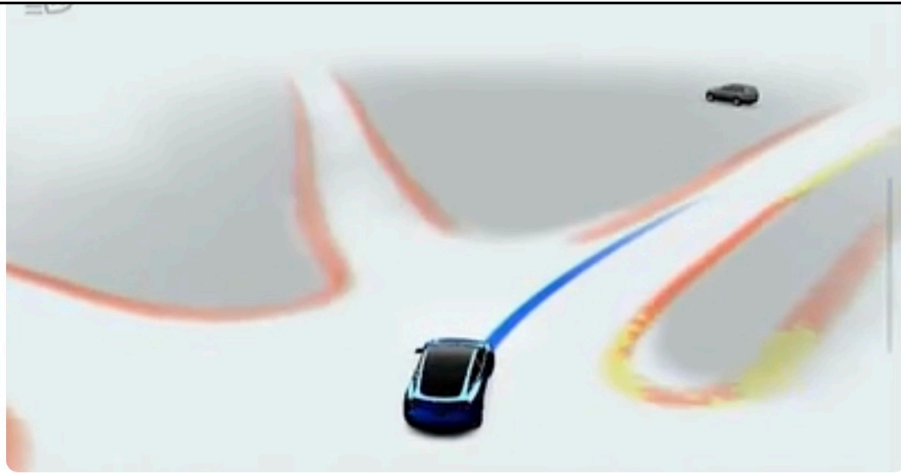
$\forall x : \text{Agents}$

Powerful(x) + Autonomous(x) + Intelligent(x)  $\Rightarrow$  Dangerous(x)



$u(\text{AIA}_i(\pi_j)) > \tau^+ \in \mathbb{Z} \text{ or } \tau^- \in \mathbb{Z}$

# PAID Reflected in The News



Tesla dashboard in “Full Self Driving” mode

## What Riding in a Self-Driving Tesla Tells Us About the Future of Autonomy

By Cade Metz, Ben Laffin, Hang Do Thi Duc and Ian Clontz. Cade and Ian spent six hours riding in a self-driving car in Jacksonville, Fla., to report this story.

Nov. 14, 2022

After releasing the new beta, Mr. Musk softened his claims about the immediate future of the technology. He now says that the technology will not be widely available until next year — and that regulators are unlikely to approve it for use without hands on the wheel. Autopilot still [requires this oversight](#).

Federal regulators have spent the past several months [investigating a series of crashes involving Autopilot](#), and they have not yet revealed the results. Safety experts worry that the arrival of Full Self-Driving will lead to more accidents.

“It is inevitable,” said Jake Fisher, senior director of Consumer Reports’ Auto Test Center, who has used the technology. “The problem comes as this system gets better and people get complacent. It will still do the unexpected.”

---

Cade Metz reported from Jacksonville, Fla. Video and photographs by Ian Clontz. Reporting and video production by Ben Laffin. Design and development by Hang Do Thi Duc.

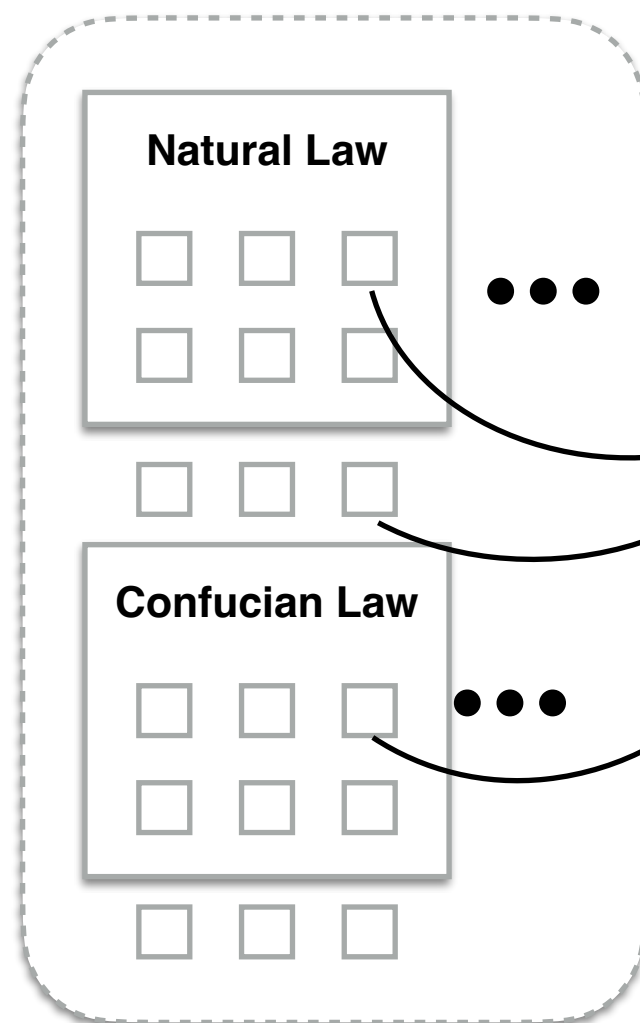




# ***Making Morally X Machines; Only Logic Can Save Us***



## Theories of Law

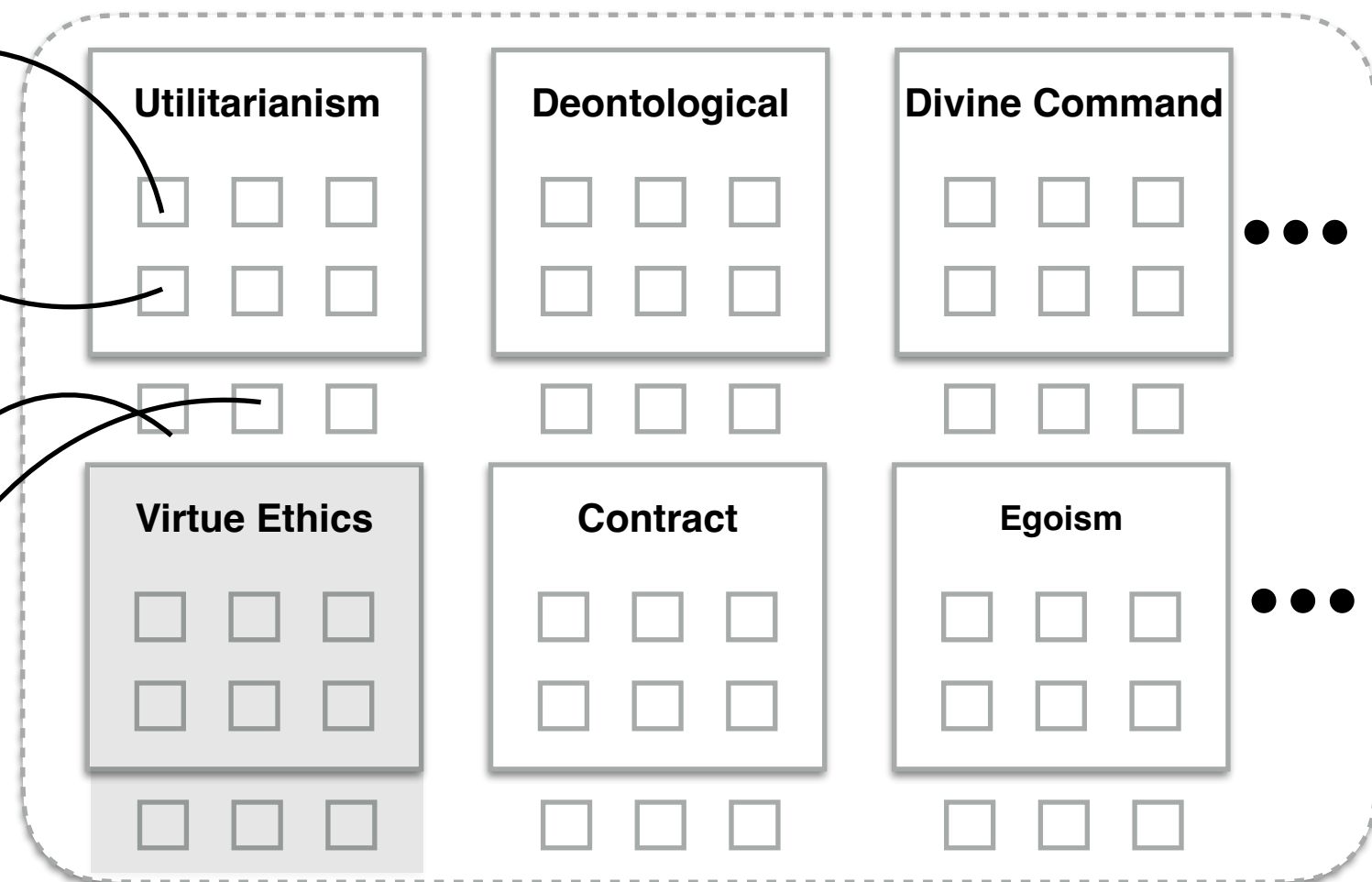


Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

## Ethical Theories

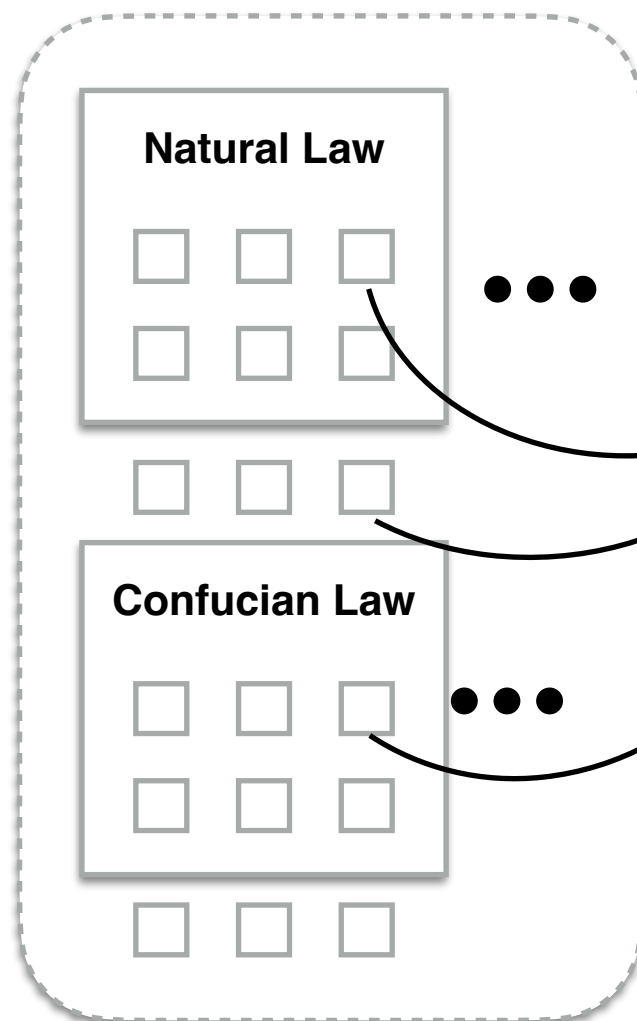




# ***Making Morally X Machines; Only Logic Can Save Us***



## Theories of Law

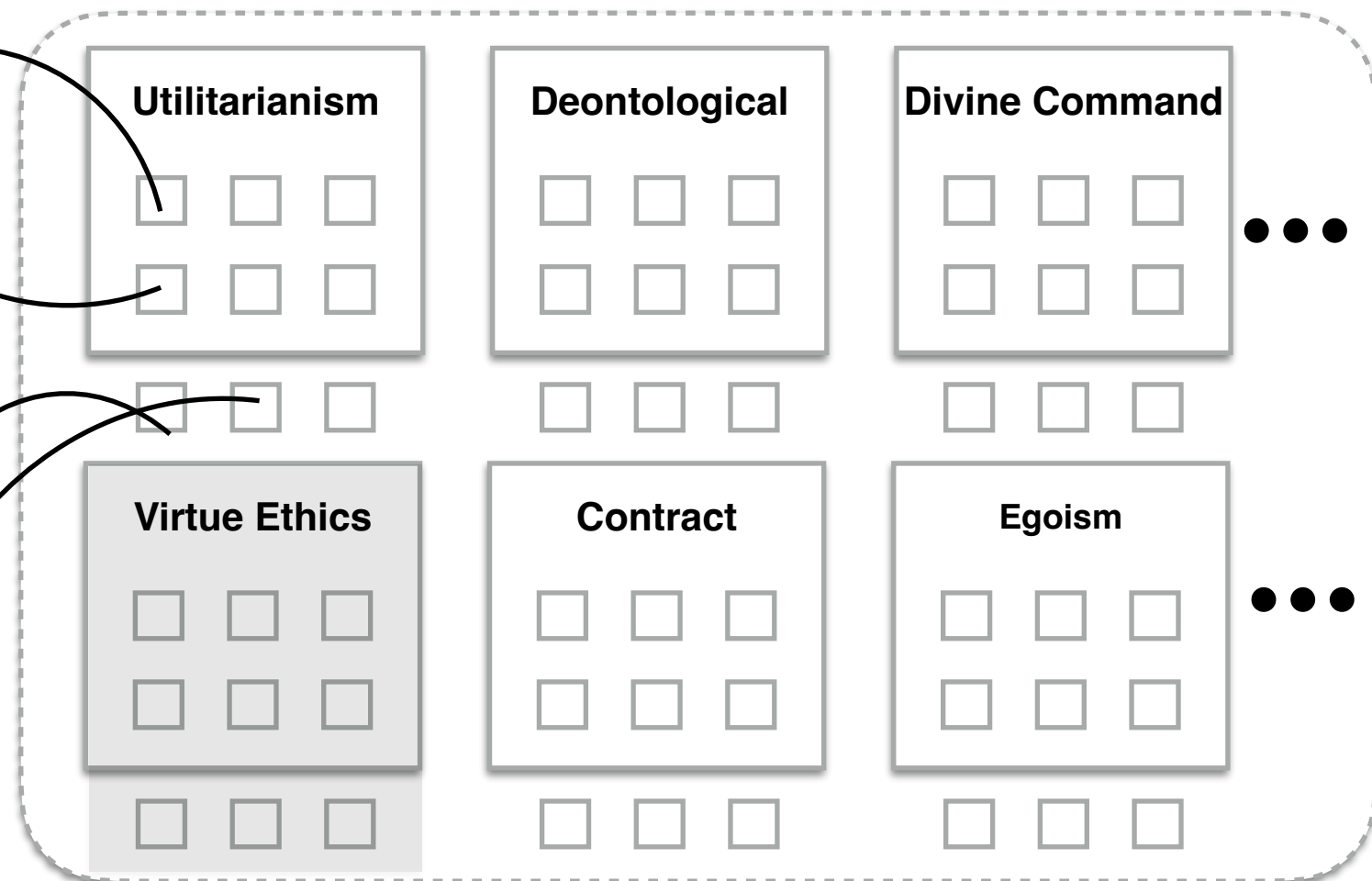


Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

## Ethical Theories

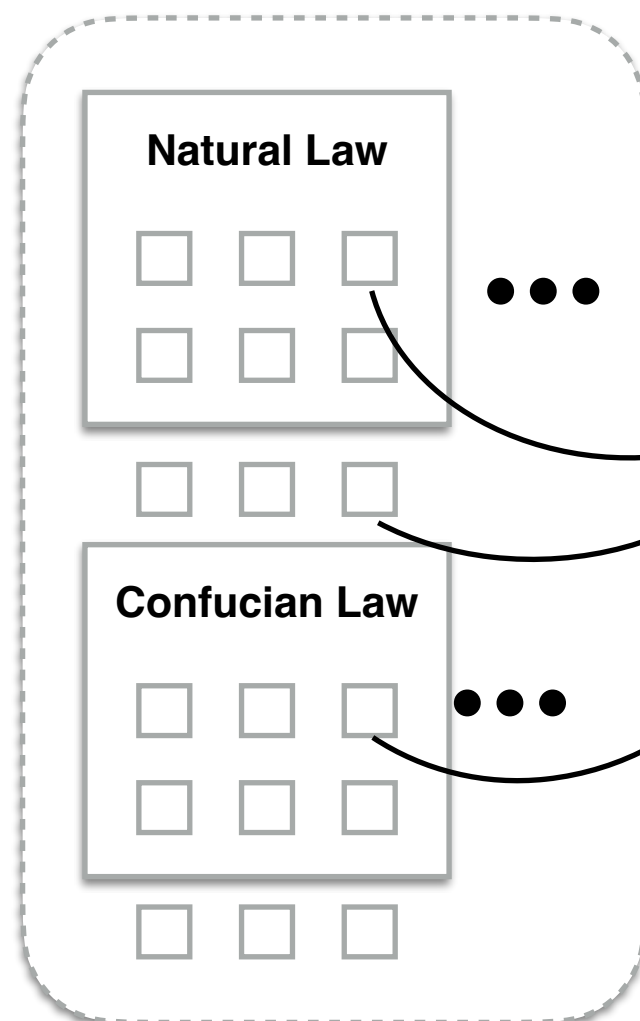




# ***Making Morally X Machines; Only Logic Can Save Us***



## Theories of Law

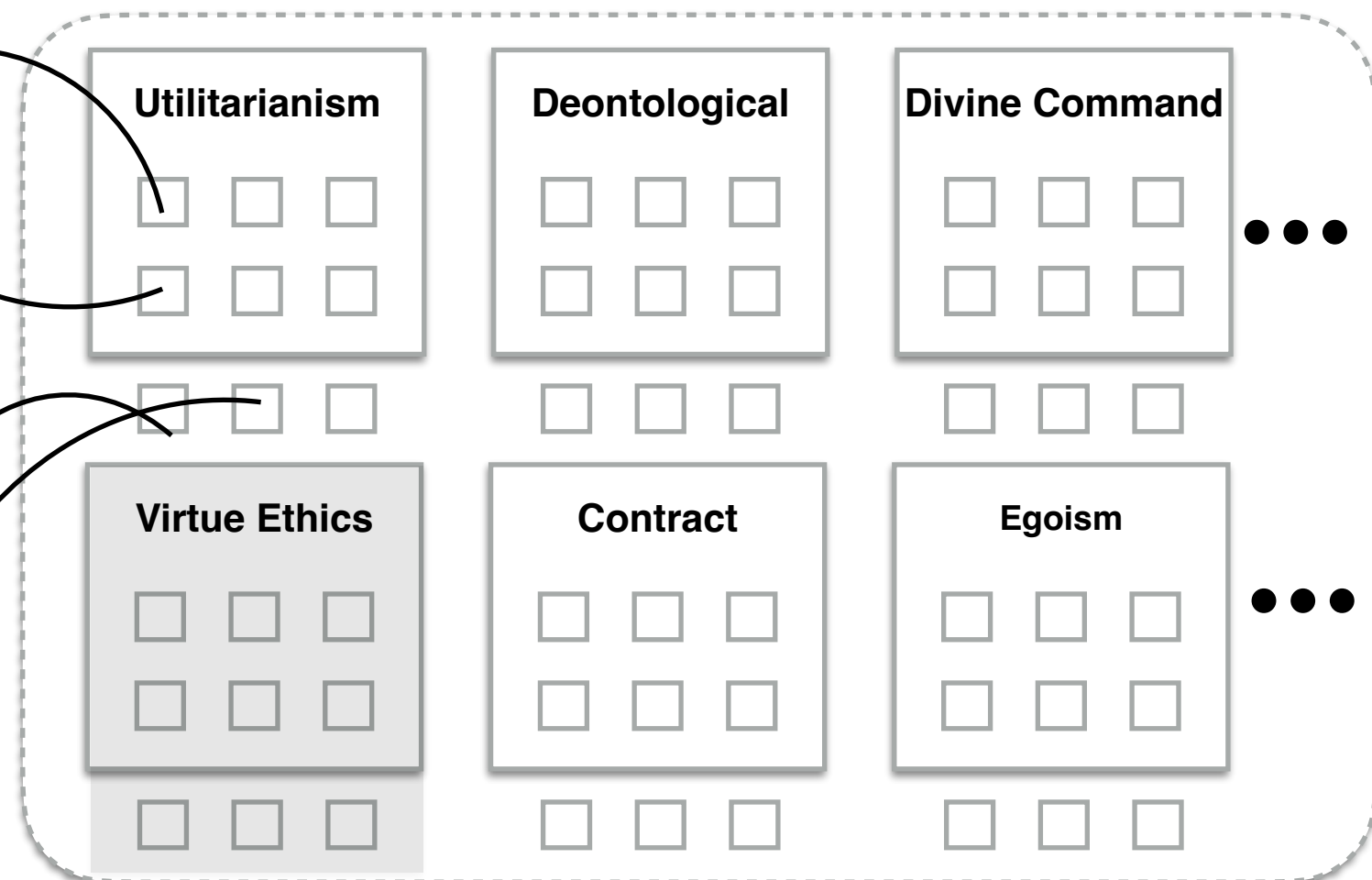


Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

## Ethical Theories



### Step I

1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.

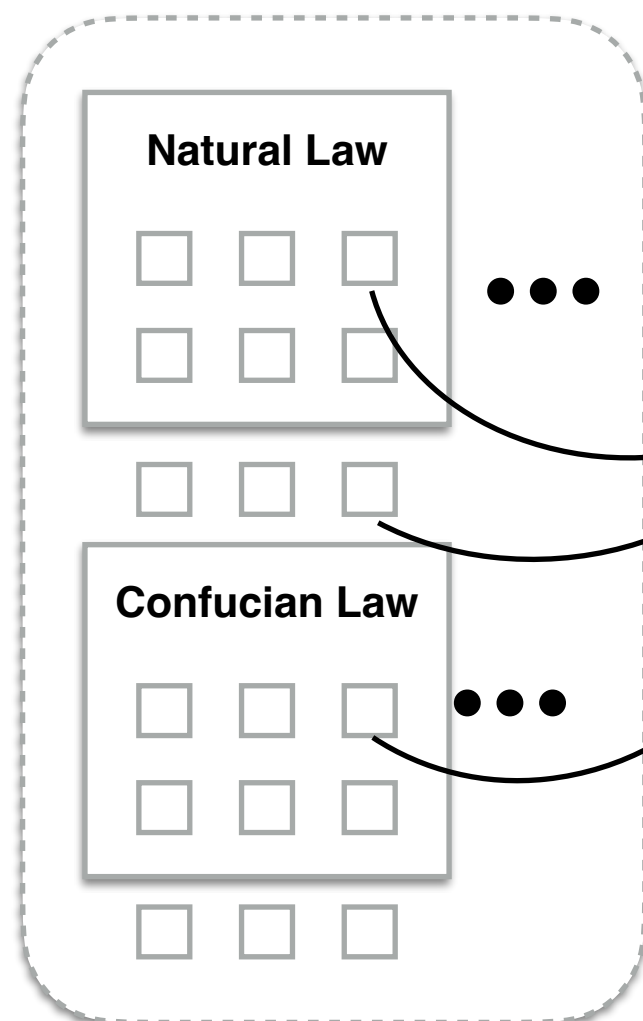




# ***Making Morally X Machines; Only Logic Can Save Us***



## Theories of Law

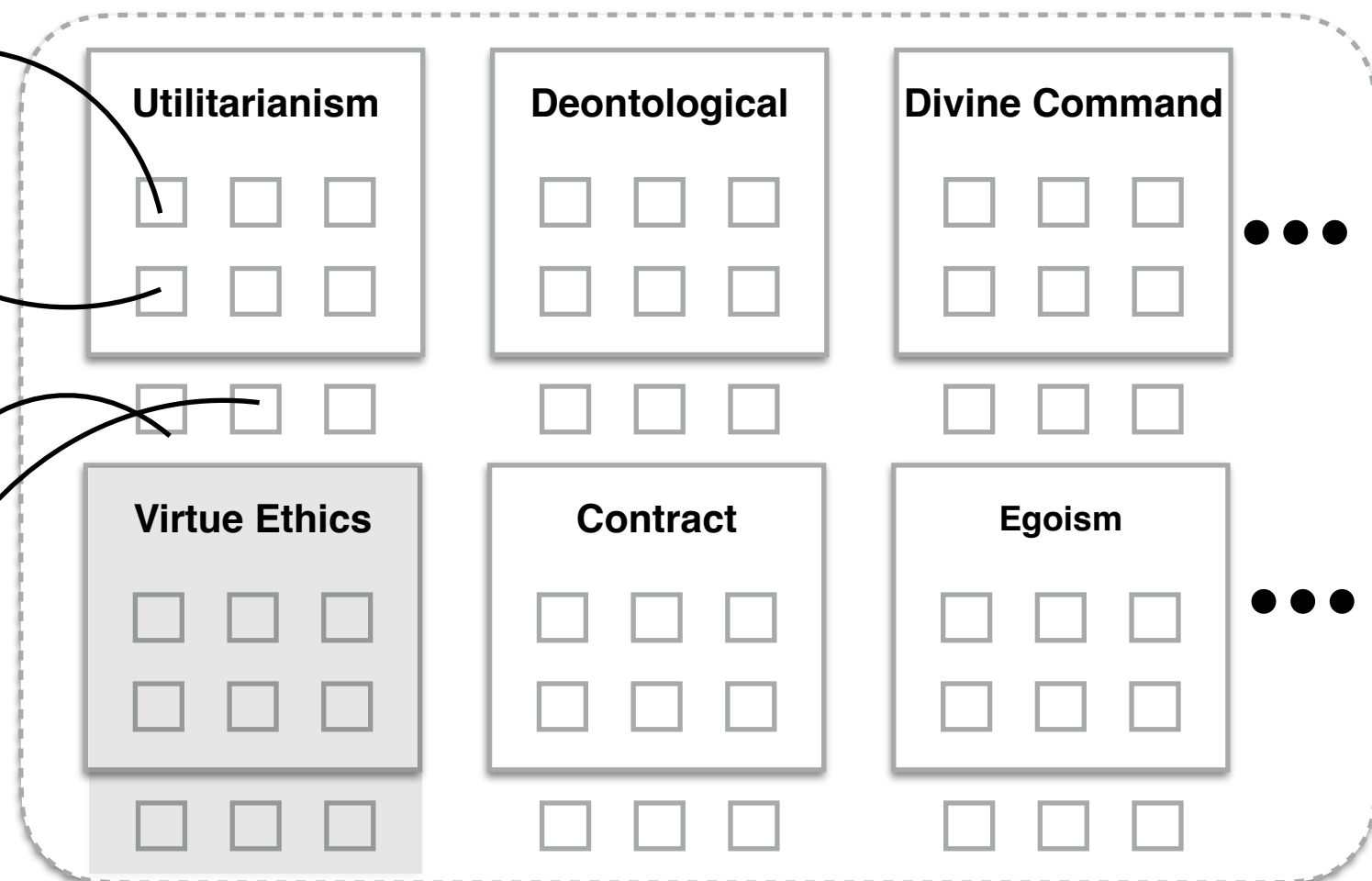


Shades  
of  
Utilitarianism

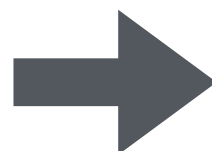
Legal Codes

Particular  
Ethical Codes

## Ethical Theories



## Step I



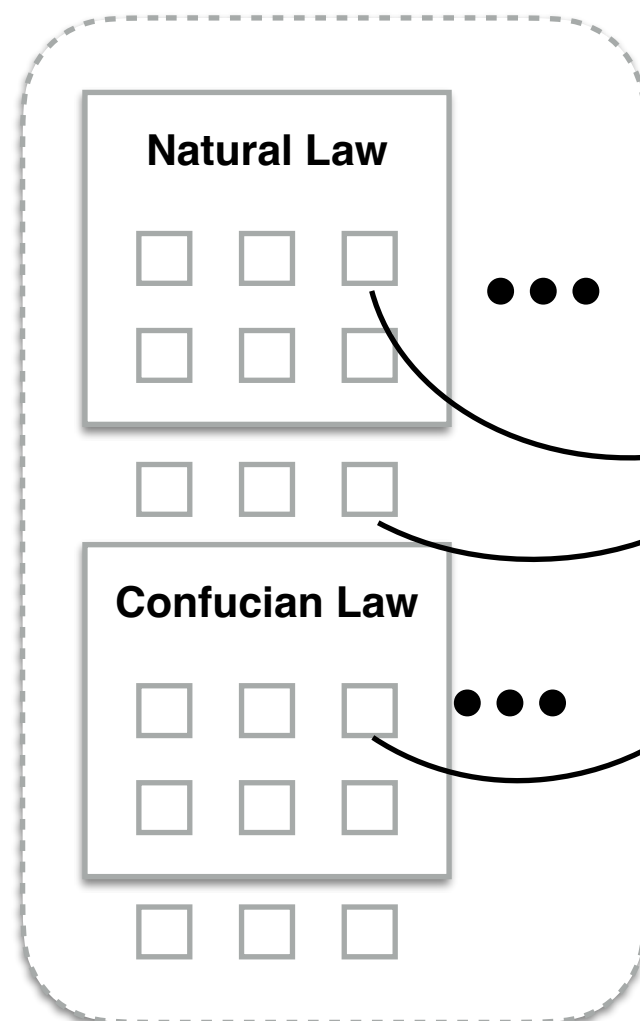
1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.



# ***Making Morally X Machines; Only Logic Can Save Us***



## Theories of Law

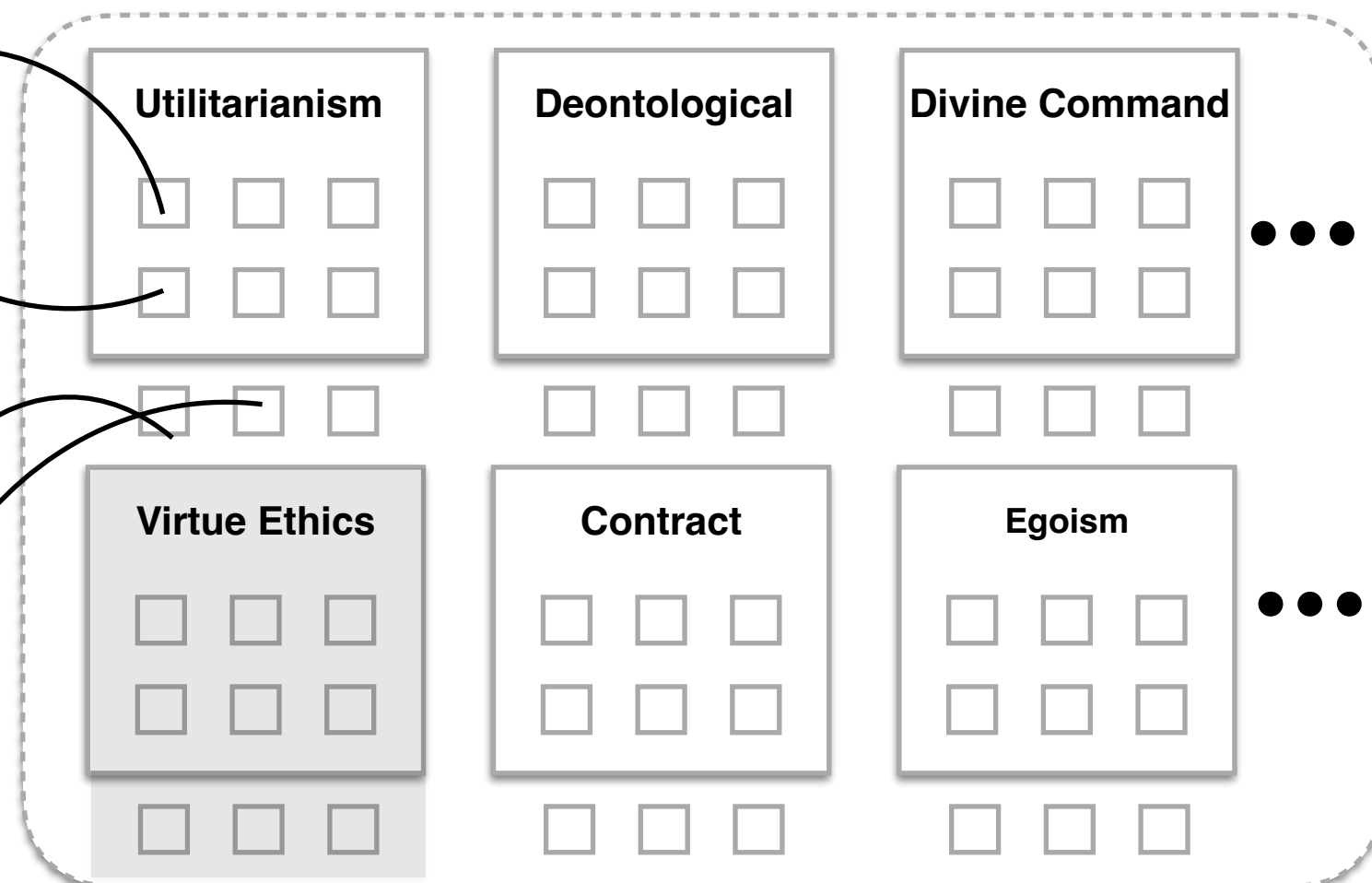


Shades  
of  
Utilitarianism

Legal Codes

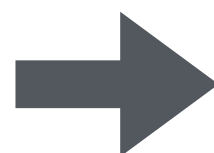
Particular  
Ethical Codes

## Ethical Theories



### Step 1

1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.



### Step 2

Automate



Reasoners



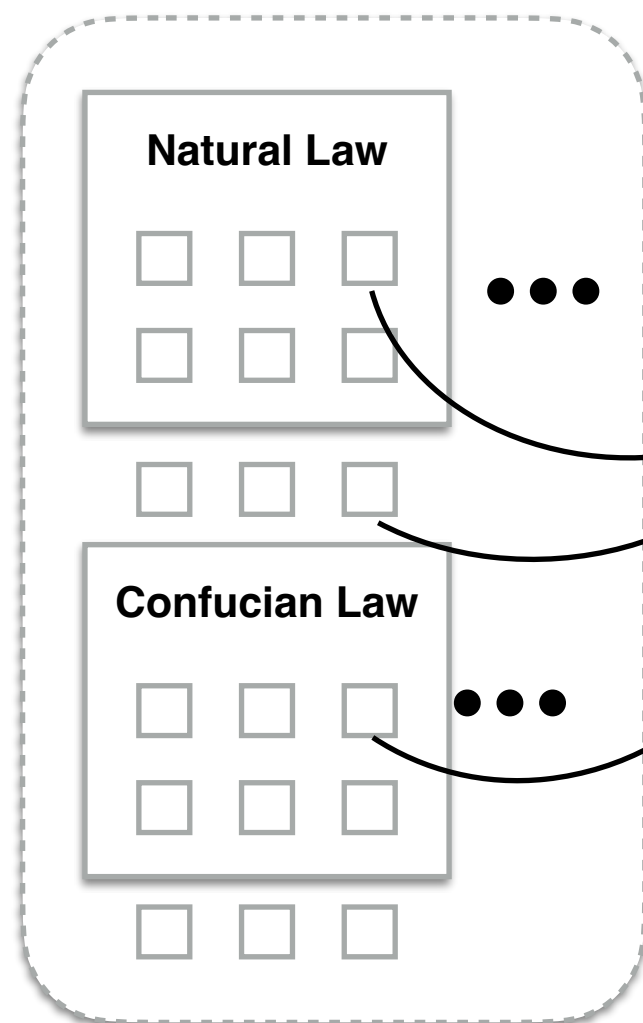
Spectra



# ***Making Morally X Machines; Only Logic Can Save Us***



## Theories of Law

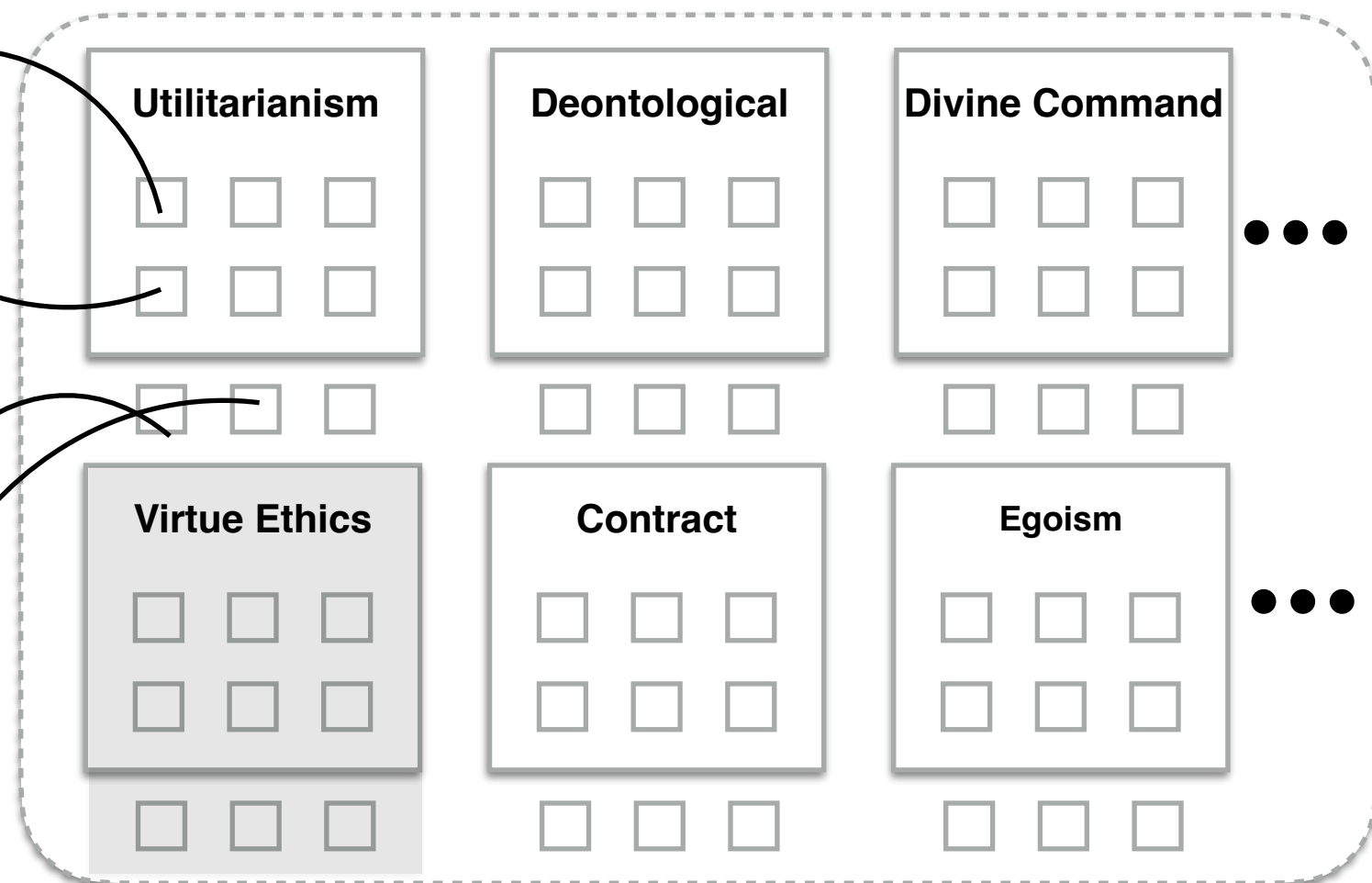


Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

## Ethical Theories



### Step 1

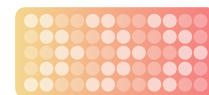
1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.

### Step 2

Automate



Reasoners



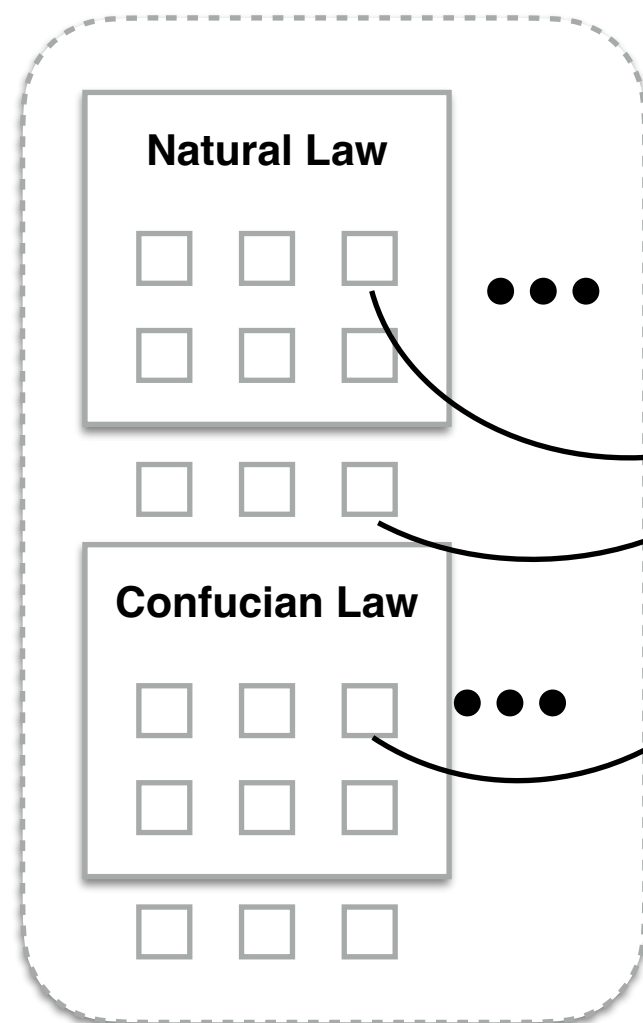
Spectra



# Making Morally X Machines; Only Logic Can Save Us



## Theories of Law

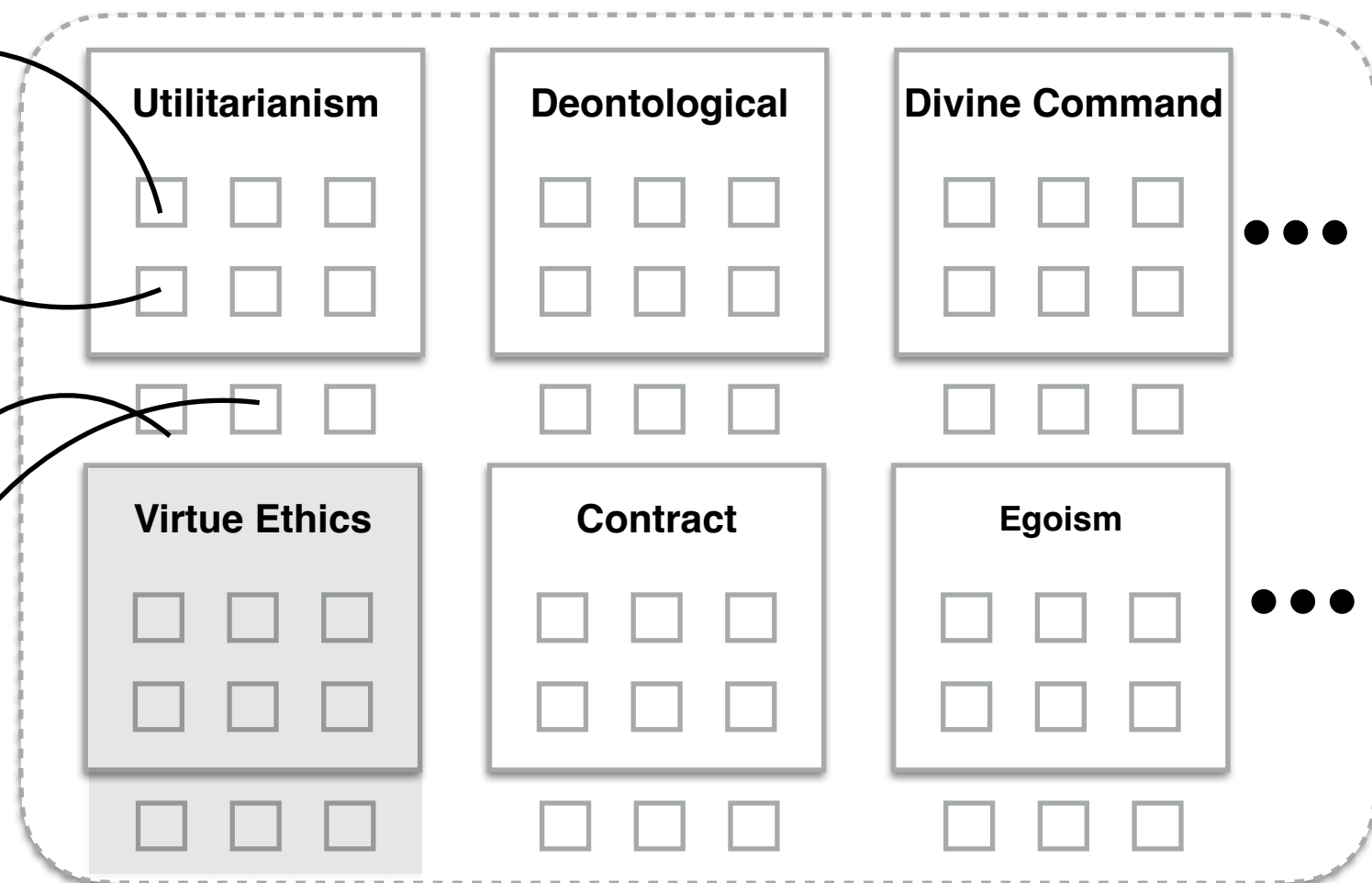


Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

## Ethical Theories



### Step 1

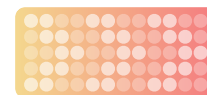
1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.

### Step 2

Automate



Reasoners



Spectra

### Step 3

Ethical OS



Ethical Substrate

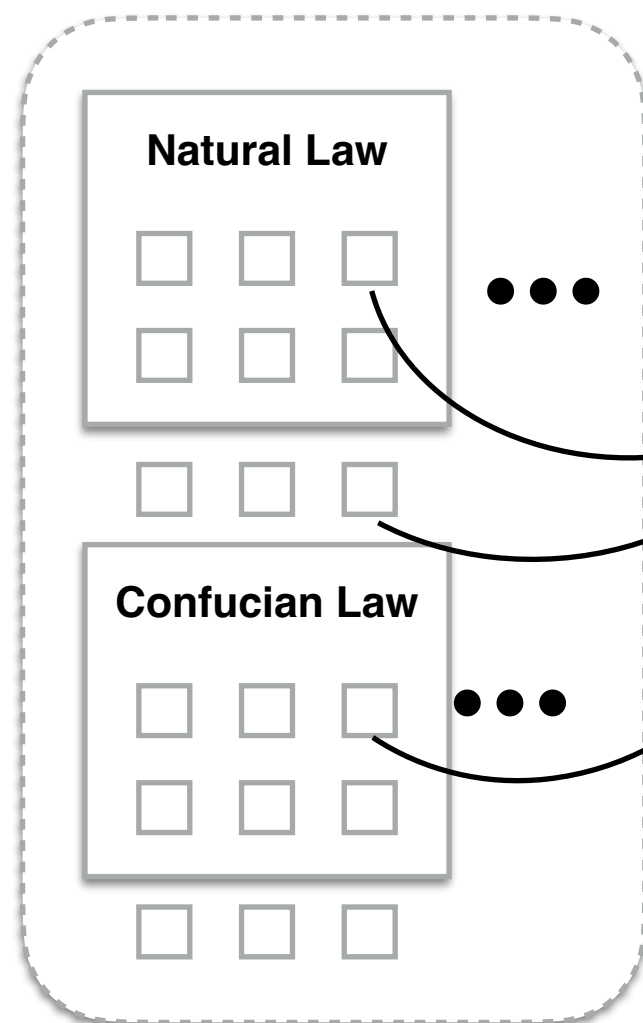
Robotic Substrate



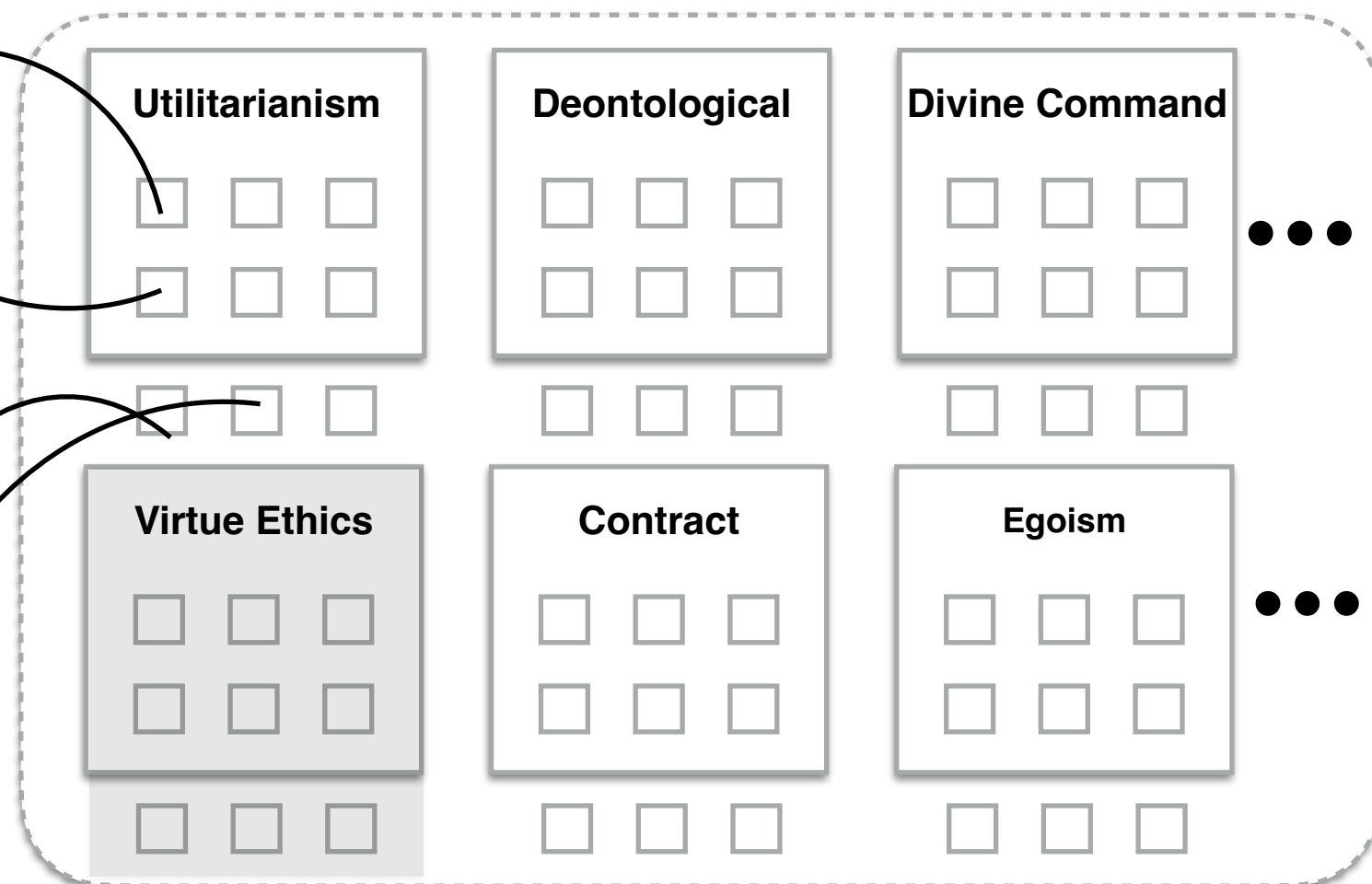
# Making Morally X Machines; Only Logic Can Save Us



## Theories of Law



## Ethical Theories



Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

### Step 1

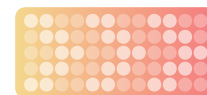
1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.

### Step 2

Automate



Reasoners



Spectra

### Step 3

Ethical OS



Ethical Substrate

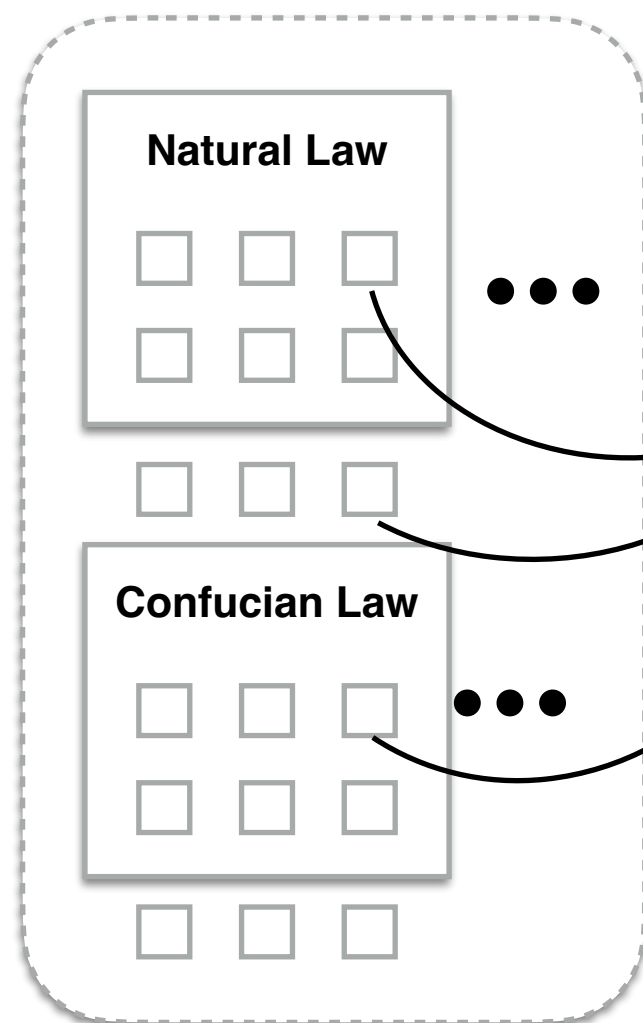
Robotic Substrate



# Making Morally X Machines; Only Logic Can Save Us



## Theories of Law

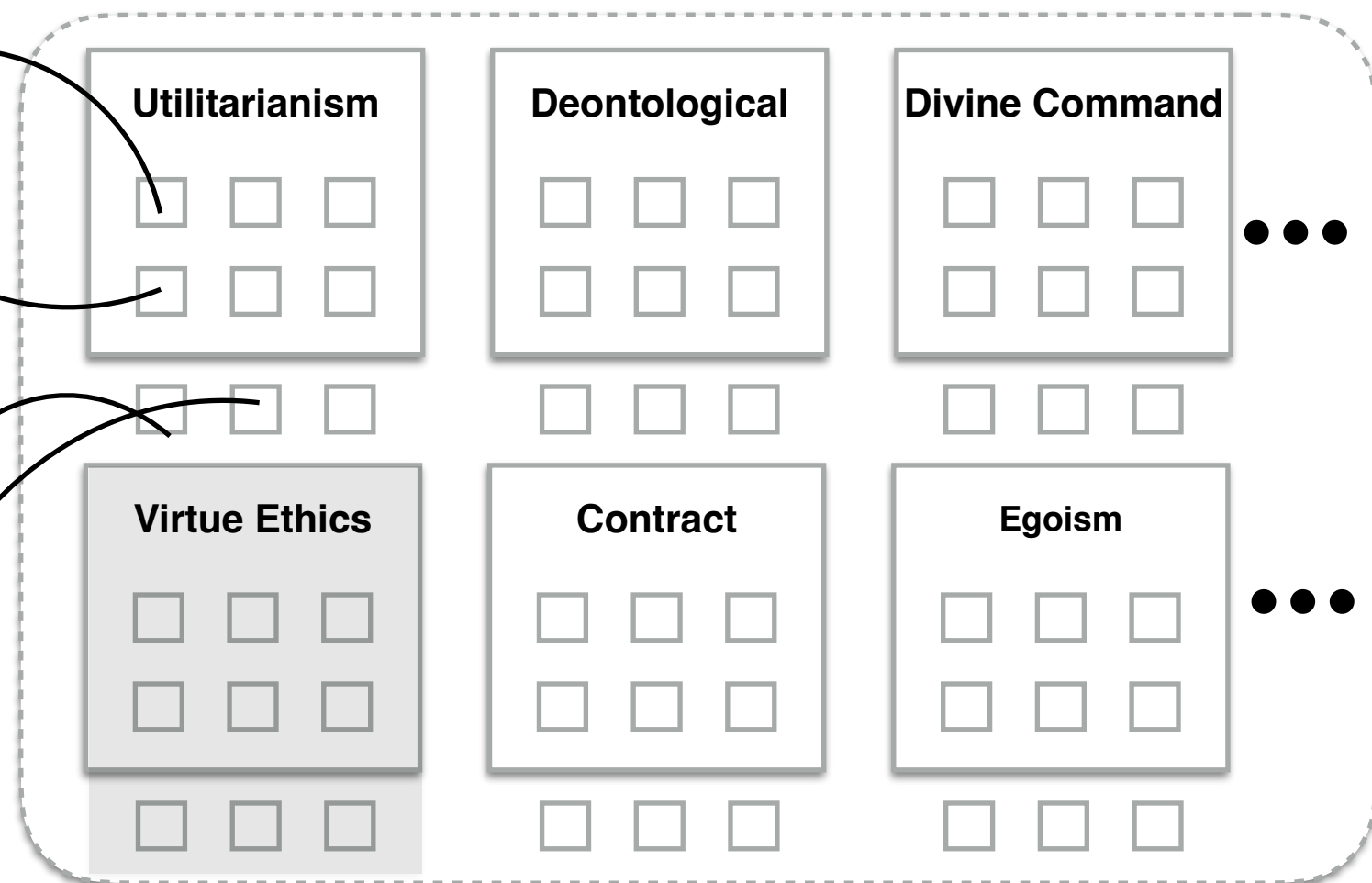


Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

## Ethical Theories



### Step 1

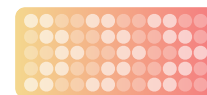
1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.

### Step 2

Automate



Reasoners



Spectra

### Step 3

Ethical OS

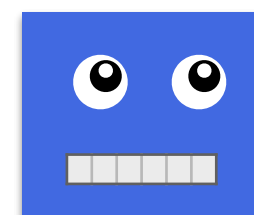


Ethical Substrate

Robotic Substrate

### Step 4

Install! — to Obtain:  
Ethically/Legally  
Correct Robot

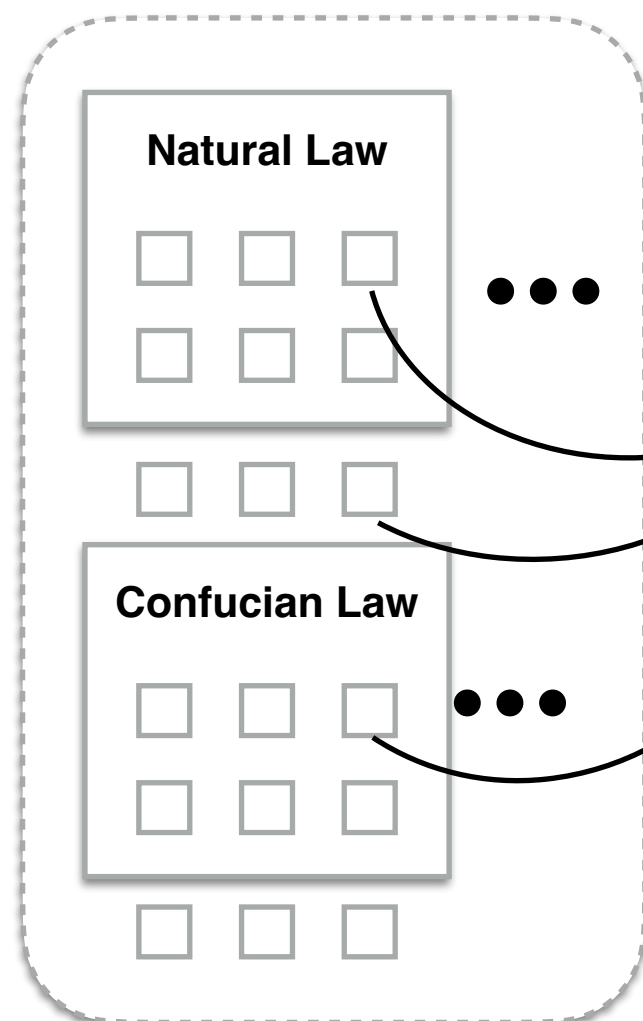




# Making Morally X Machines; Only Logic Can Save Us



## Theories of Law

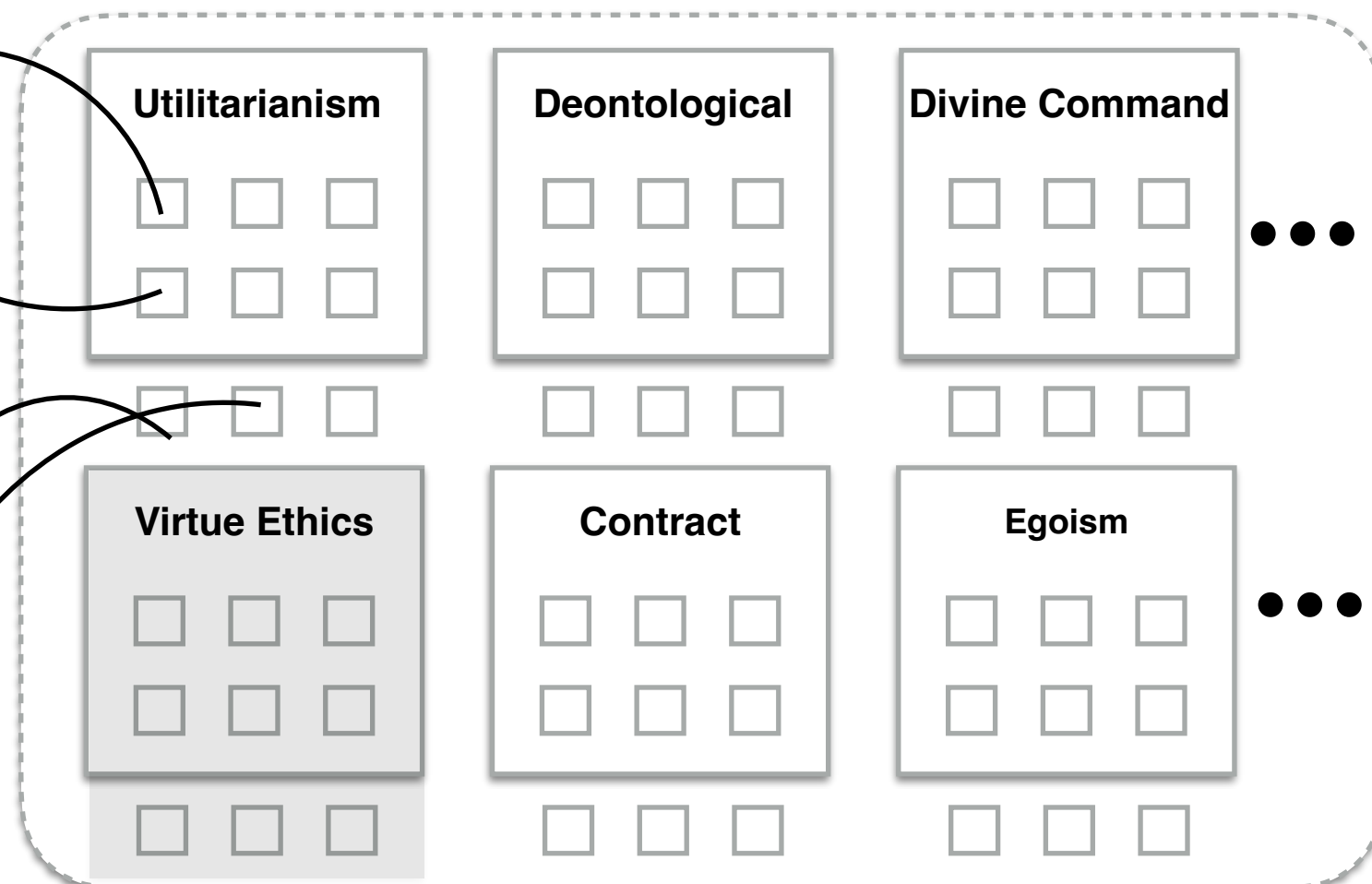


Shades  
of  
Utilitarianism

Legal Codes

Particular  
Ethical Codes

## Ethical Theories



### Step 1

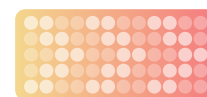
1. Pick a theory
2. Pick a code
3. Run through EH.
4. Formalize in a Cognitive Calculus.

### Step 2

Automate



Reasoners



Spectra

### Step 3

Ethical OS

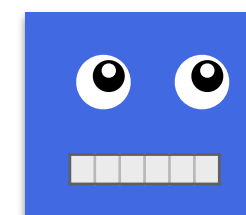


Ethical Substrate

Robotic Substrate

### Step 4

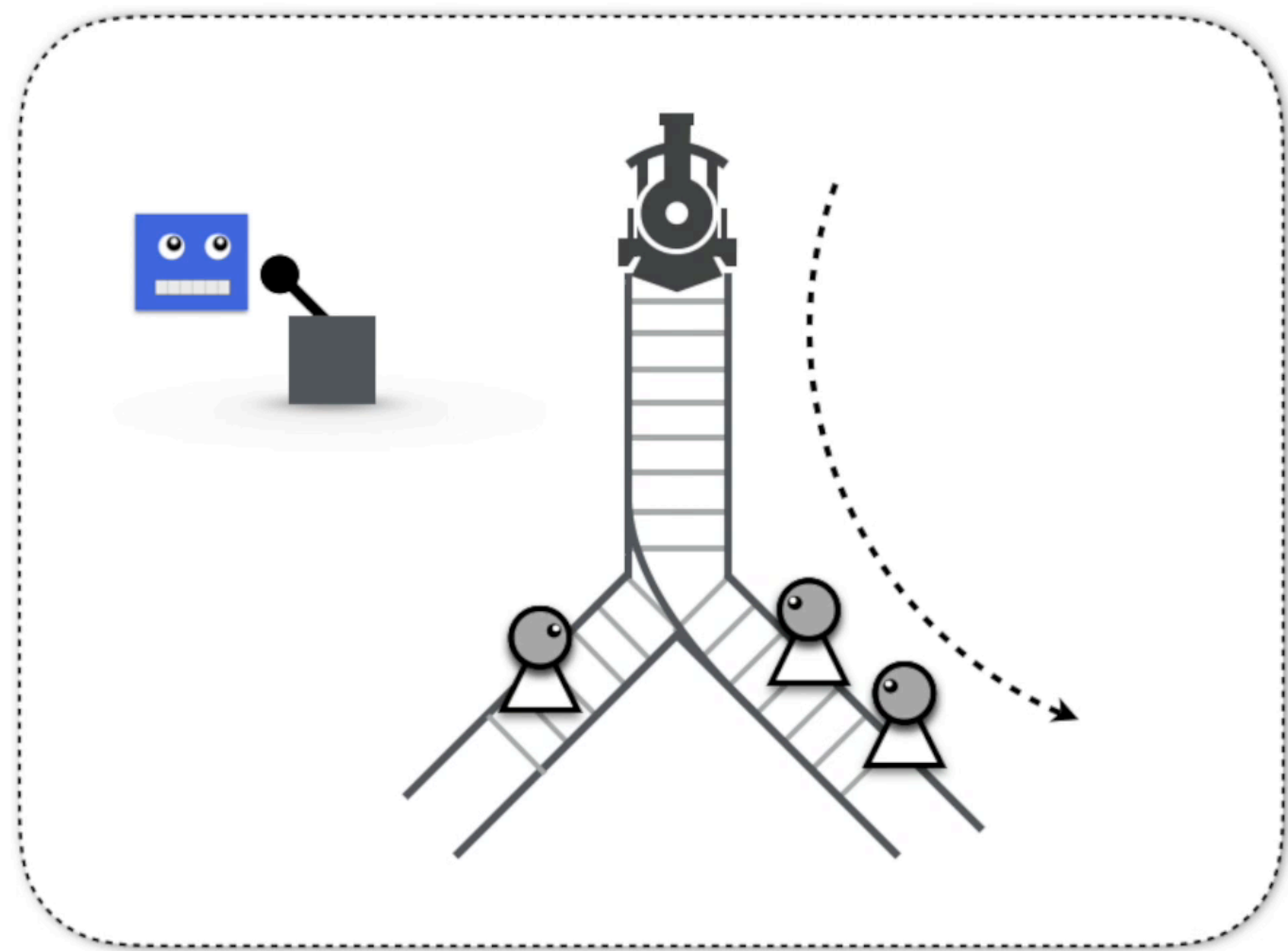
Install! — to Obtain:  
Ethically/Legally  
Correct Robot

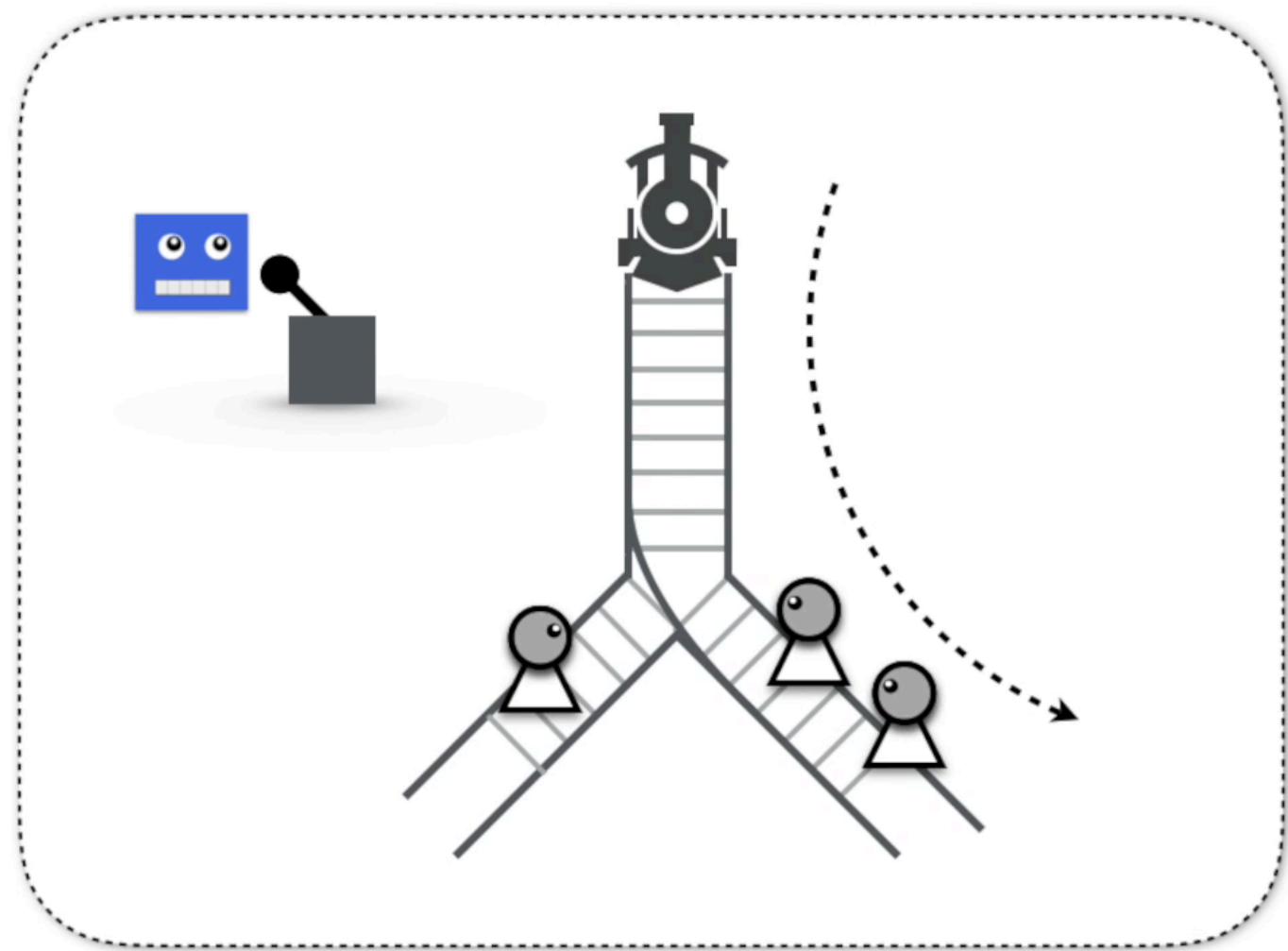


e.g. “Toward the Engineering of Virtuous Robots” Naveen, Selmer et al.

DDE\* & DIE







<https://www.ijcai.org/Proceedings/2017/0658.pdf>

Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI-17)

## On Automating the Doctrine of Double Effect

Naveen Sundar Govindarajulu and Selmer Bringsjord

Rensselaer Polytechnic Institute, Troy, NY  
{naveensundarg,selmer.bringsjord}@gmail.com

### Abstract

The **doctrine of double effect** ( $\mathcal{DDE}$ ) is a long-studied ethical principle that governs when actions that have both positive and negative effects are to be allowed. The goal in this paper is to automate  $\mathcal{DDE}$ . We briefly present  $\mathcal{DDE}$ , and use a first-order modal logic, the **deontic cognitive event calculus**, as our framework to formalize the doctrine.

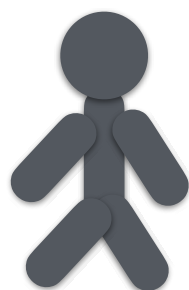
— provided that 1) the harmful effects are not intended; 2) the harmful effects are not used to achieve the beneficial effects (harm is merely a *side-effect*); and 3) benefits outweigh the harm by a significant amount. What distinguishes  $\mathcal{DDE}$  from, say, naïve forms of consequentialism in ethics (e.g. act utilitarianism, which holds that an action is obligatory for an autonomous agent if and only if it produces the most utility among all competing actions) is that purely mental intentions in and of themselves independent of conse-

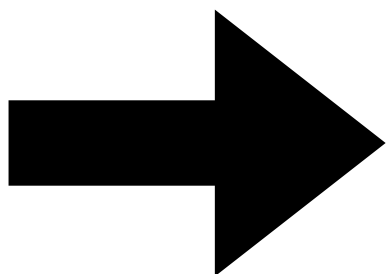
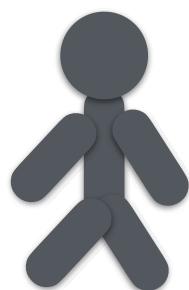
## 4 Informal $\mathcal{DDE}$

We now informally but rigorously present  $\mathcal{DDE}$ . We assume we have at hand an ethical hierarchy of actions as in the deontological case (e.g. forbidden, neutral, obligatory); see [Bringsjord, 2017]. We also assume that we have a utility or goodness function for states of the world or effects as in the consequentialist case. For an autonomous agent  $a$ , an action  $\alpha$  in a situation  $\sigma$  at time  $t$  is said to be  $\mathcal{DDE}$ -compliant *iff*:

- $C_1$  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.







**F<sub>1</sub>**  $\alpha$  carried out at  $t$  is not forbidden. That is:

$$\Gamma \not\models \neg \mathbf{O}(a, t, \sigma, \neg \text{happens}(\text{action}(a, \alpha), t))$$

**F<sub>2</sub>** 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<sub>3a</sub>** The agent  $a$  intends at least one good effect. (**F<sub>2</sub>** should still hold after removing all other good effects.) There is at least one fluent  $f_g$  in  $\alpha_I^{a,t}$  with  $\mu(f_g, y) > 0$ , 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:

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

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

$$\begin{aligned} \Gamma &\not\models \mathbf{I}(a, t, \text{Holds}(f_b, y)) \text{ and} \\ \Gamma &\not\models \mathbf{I}(a, t, \neg \text{Holds}(f_g, y)) \end{aligned}$$

**F<sub>4</sub>** The harmful effects don't cause the good effects. Four permutations, paralleling the definition of  $\triangleright$  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:

$$\Gamma \vdash \neg \triangleright (\text{Holds}(f_b, t_1), \text{Holds}(f_g, t_2))$$





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

**F<sub>1</sub>**  $\alpha$  carried out at  $t$  is not forbidden. That is:

$$\Gamma \not\models \neg \mathbf{O}(a, t, \sigma, \neg \text{happens}(\text{action}(a, \alpha), t))$$

**F<sub>2</sub>** 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<sub>3a</sub>** The agent  $a$  intends at least one good effect. (**F<sub>2</sub>** should still hold after removing all other good effects.) There is at least one fluent  $f_g$  in  $\alpha_I^{a,t}$  with  $\mu(f_g, y) > 0$ , 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:

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

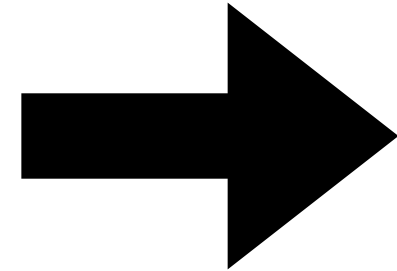
**F<sub>3b</sub>** The agent  $a$  does not intend any bad effect. For all fluents  $f_b$  in  $\alpha_I^{a,t}$  with  $\mu(f_b, y) < 0$ , or  $f_g$  in  $\alpha_T^{a,t}$  with  $\mu(f_g, y) > 0$ , and for all  $y$  such that  $t < y \leq H$  the following holds:

$$\begin{aligned} \Gamma &\not\models \mathbf{I}(a, t, \text{Holds}(f_b, y)) \text{ and} \\ \Gamma &\not\models \mathbf{I}(a, t, \neg \text{Holds}(f_g, y)) \end{aligned}$$

**F<sub>4</sub>** The harmful effects don't cause the good effects. Four permutations, paralleling the definition of  $\triangleright$  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:

$$\Gamma \vdash \neg \triangleright (\text{Holds}(f_b, t_1), \text{Holds}(f_g, t_2))$$

$\mathbb{P}_{\text{DDE}_1} + \text{ShadowProver}$



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

**F<sub>1</sub>**  $\alpha$  carried out at  $t$  is not forbidden. That is:

$$\Gamma \not\models \neg \mathbf{O}(a, t, \sigma, \neg \text{happens}(\text{action}(a, \alpha), t))$$

**F<sub>2</sub>** 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<sub>3a</sub>** The agent  $a$  intends at least one good effect. (**F<sub>2</sub>** should still hold after removing all other good effects.) There is at least one fluent  $f_g$  in  $\alpha_I^{a,t}$  with  $\mu(f_g, y) > 0$ , 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:

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

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

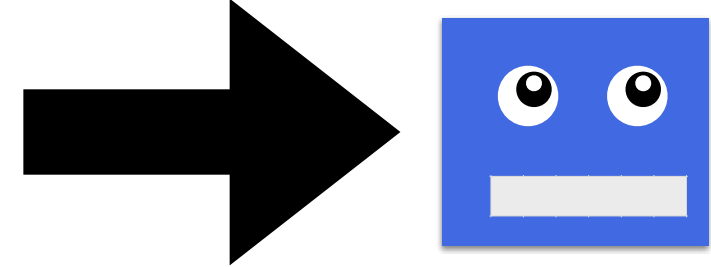
$$\begin{aligned} \Gamma &\not\models \mathbf{I}(a, t, \text{Holds}(f_b, y)) \text{ and} \\ \Gamma &\not\models \mathbf{I}(a, t, \neg \text{Holds}(f_g, y)) \end{aligned}$$

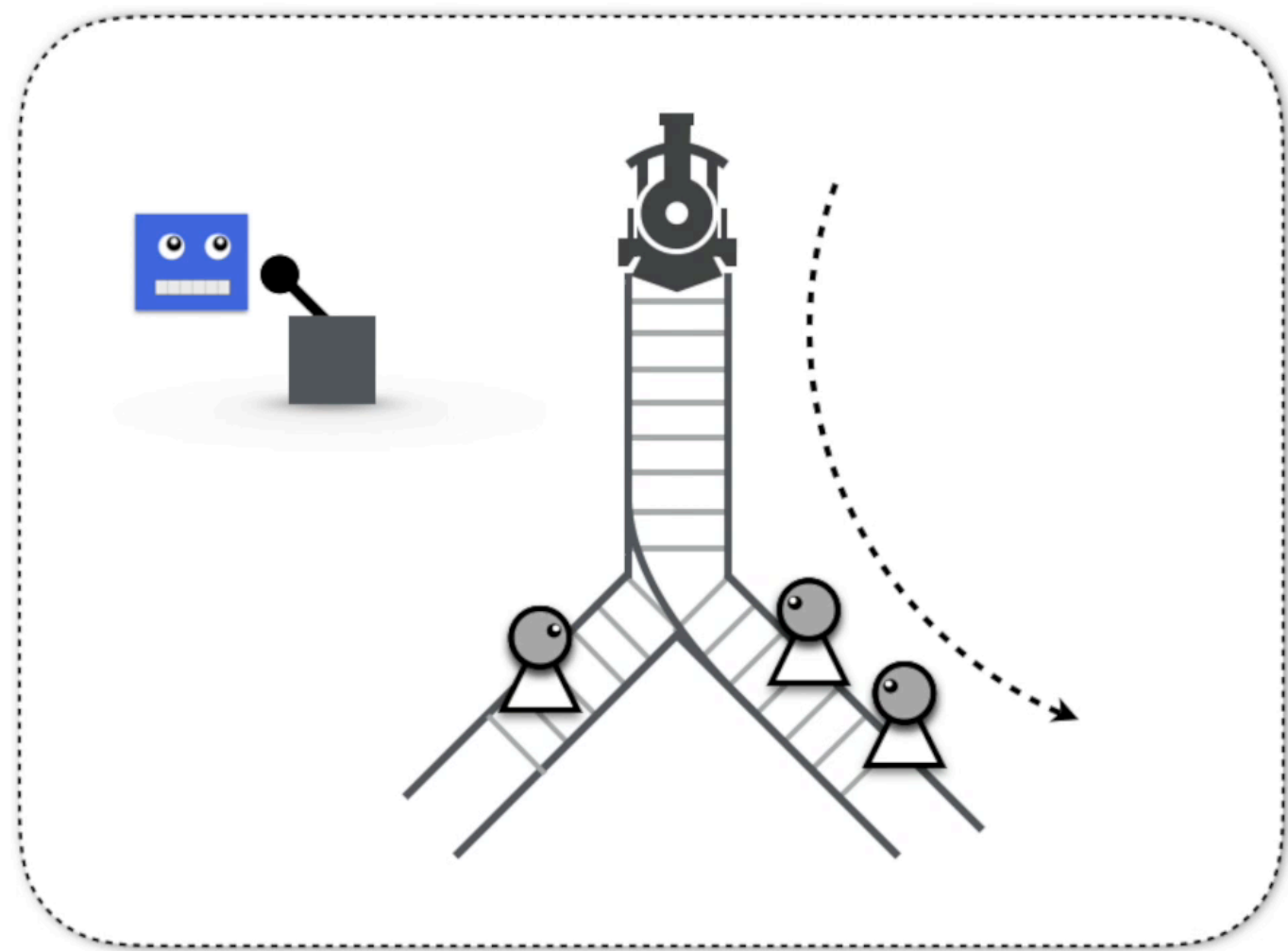
**F<sub>4</sub>** The harmful effects don't cause the good effects. Four permutations, paralleling the definition of  $\triangleright$  above, hold here. One such permutation is shown below. For any bad fluent  $f_b$  holding at  $t_1$ , and any good fluent  $f_g$  holding at some  $t_2$ , such that  $t < t_1, t_2 \leq H$ , the following holds:

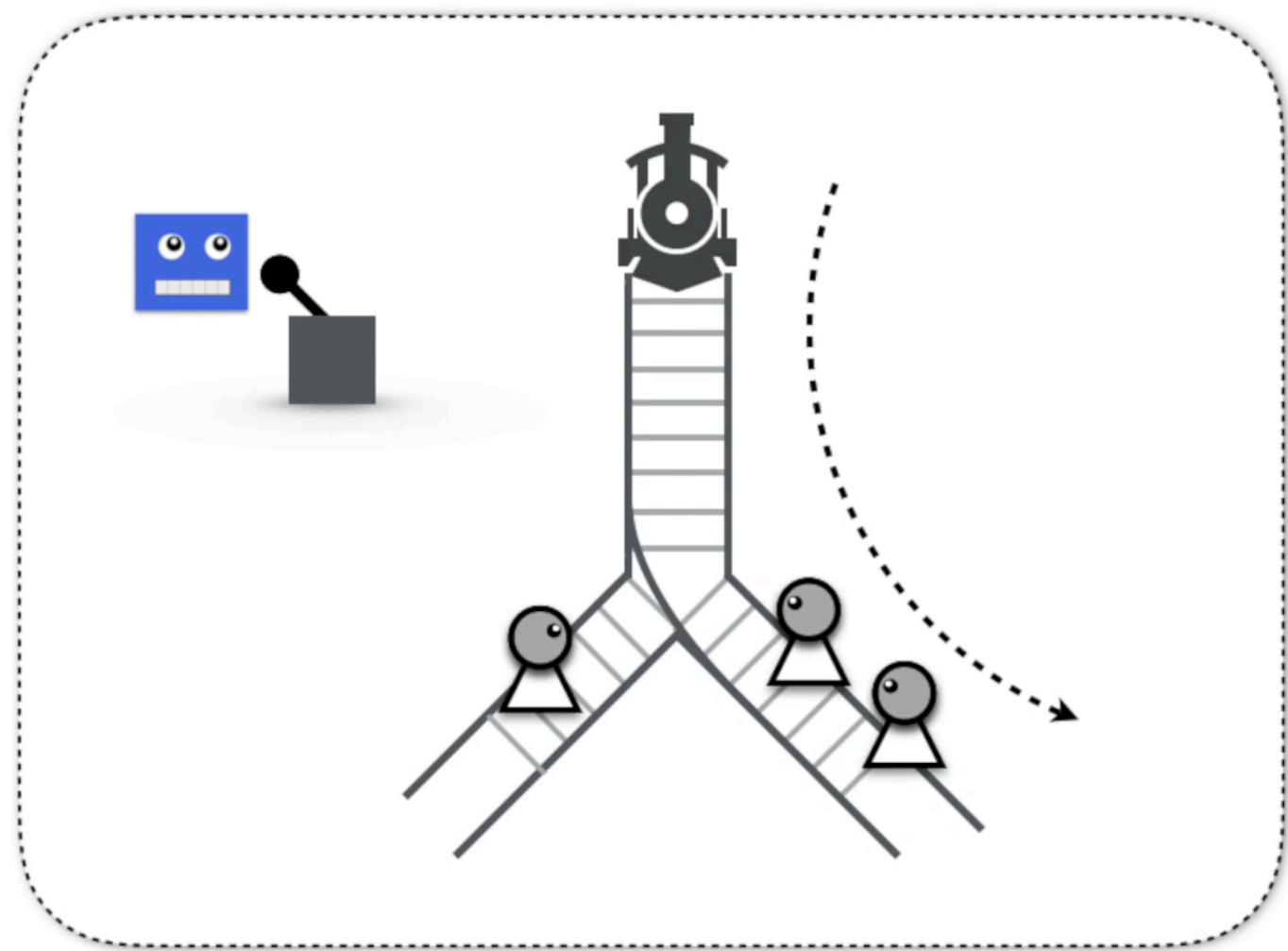
$$\Gamma \vdash \neg \triangleright (\text{Holds}(f_b, t_1), \text{Holds}(f_g, t_2))$$



$\mathbb{P}_{\text{DDE}_1} + \text{ShadowProver}$







### Inference Schemata

$$\begin{array}{c}
\frac{\mathbf{K}(a, t_1, \Gamma), \Gamma \vdash \phi, t_1 \leq t_2}{\mathbf{K}(a, t_2, \phi)} [R_K] \quad \frac{\mathbf{B}(a, t_1, \Gamma), \Gamma \vdash \phi, t_1 \leq t_2}{\mathbf{B}(a, t_2, \phi)} [R_B] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{P}(a, t, \phi) \rightarrow \mathbf{K}(a, t, \phi))} [R_1] \quad \frac{}{\mathbf{C}(t, \mathbf{K}(a, t, \phi) \rightarrow \mathbf{B}(a, t, \phi))} [R_2] \\
\\
\frac{\mathbf{C}(t, \phi) \ t \leq t_1 \dots t \leq t_n}{\mathbf{K}(a_1, t_1, \dots \mathbf{K}(a_n, t_n, \phi) \dots)} [R_3] \quad \frac{\mathbf{K}(a, t, \phi)}{\phi} [R_4] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{K}(a, t_1, \phi_1 \rightarrow \phi_2)) \rightarrow \mathbf{K}(a, t_2, \phi_1) \rightarrow \mathbf{K}(a, t_3, \phi_2)} [R_5] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{B}(a, t_1, \phi_1 \rightarrow \phi_2)) \rightarrow \mathbf{B}(a, t_2, \phi_1) \rightarrow \mathbf{B}(a, t_3, \phi_2)} [R_6] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{C}(t_1, \phi_1 \rightarrow \phi_2)) \rightarrow \mathbf{C}(t_2, \phi_1) \rightarrow \mathbf{C}(t_3, \phi_2)} [R_7] \\
\\
\frac{}{\mathbf{C}(t, \forall x. \phi \rightarrow \phi[x \mapsto t])} [R_8] \quad \frac{}{\mathbf{C}(t, \phi_1 \leftrightarrow \phi_2 \rightarrow \neg \phi_2 \rightarrow \neg \phi_1)} [R_9] \\
\\
\frac{}{\mathbf{C}(t, [\phi_1 \wedge \dots \wedge \phi_n \rightarrow \phi] \rightarrow [\phi_1 \rightarrow \dots \rightarrow \phi_n \rightarrow \psi])} [R_{10}] \\
\\
\frac{\mathbf{S}(s, h, t, \phi)}{\mathbf{B}(h, t, \mathbf{B}(s, t, \phi))} [R_{12}] \quad \frac{\mathbf{I}(a, t, \text{happens}(\text{action}(a^*, \alpha), t'))}{\mathbf{P}(a, t, \text{happens}(\text{action}(a^*, \alpha), t))} [R_{13}] \\
\\
\frac{\mathbf{B}(a, t, \phi) \ \mathbf{B}(a, t, \mathbf{O}(a, t, \phi, \chi)) \ \mathbf{O}(a, t, \phi, \chi)}{\mathbf{K}(a, t, \mathbf{I}(a, t, \chi))} [R_{14}]
\end{array}$$

### Inference Schemata

$$\begin{array}{c}
\frac{\mathbf{K}(a, t_1, \Gamma), \Gamma \vdash \phi, t_1 \leq t_2}{\mathbf{K}(a, t_2, \phi)} [R_K] \quad \frac{\mathbf{B}(a, t_1, \Gamma), \Gamma \vdash \phi, t_1 \leq t_2}{\mathbf{B}(a, t_2, \phi)} [R_B] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{P}(a, t, \phi) \rightarrow \mathbf{K}(a, t, \phi))} [R_1] \quad \frac{}{\mathbf{C}(t, \mathbf{K}(a, t, \phi) \rightarrow \mathbf{B}(a, t, \phi))} [R_2] \\
\\
\frac{\mathbf{C}(t, \phi) \ t \leq t_1 \dots t \leq t_n}{\mathbf{K}(a_1, t_1, \dots \mathbf{K}(a_n, t_n, \phi) \dots)} [R_3] \quad \frac{\mathbf{K}(a, t, \phi)}{\phi} [R_4] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{K}(a, t_1, \phi_1 \rightarrow \phi_2)) \rightarrow \mathbf{K}(a, t_2, \phi_1) \rightarrow \mathbf{K}(a, t_3, \phi_2)} [R_5] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{B}(a, t_1, \phi_1 \rightarrow \phi_2)) \rightarrow \mathbf{B}(a, t_2, \phi_1) \rightarrow \mathbf{B}(a, t_3, \phi_2)} [R_6] \\
\\
\frac{}{\mathbf{C}(t, \mathbf{C}(t_1, \phi_1 \rightarrow \phi_2)) \rightarrow \mathbf{C}(t_2, \phi_1) \rightarrow \mathbf{C}(t_3, \phi_2)} [R_7] \\
\\
\frac{}{\mathbf{C}(t, \forall x. \phi \rightarrow \phi[x \mapsto t])} [R_8] \quad \frac{}{\mathbf{C}(t, \phi_1 \leftrightarrow \phi_2 \rightarrow \neg \phi_2 \rightarrow \neg \phi_1)} [R_9] \\
\\
\frac{}{\mathbf{C}(t, [\phi_1 \wedge \dots \wedge \phi_n \rightarrow \phi] \rightarrow [\phi_1 \rightarrow \dots \rightarrow \phi_n \rightarrow \psi])} [R_{10}] \\
\\
\frac{\mathbf{S}(s, h, t, \phi)}{\mathbf{B}(h, t, \mathbf{B}(s, t, \phi))} [R_{12}] \quad \frac{\mathbf{I}(a, t, \text{happens}(\text{action}(a^*, \alpha), t'))}{\mathbf{P}(a, t, \text{happens}(\text{action}(a^*, \alpha), t))} [R_{13}] \\
\\
\frac{\mathbf{B}(a, t, \phi) \ \mathbf{B}(a, t, \mathbf{O}(a, t, \phi, \chi)) \ \mathbf{O}(a, t, \phi, \chi)}{\mathbf{K}(a, t, \mathbf{I}(a, t, \chi))} [R_{14}]
\end{array}$$

# But what about self-sacrifice?!

## 3. Goal

In this section we render precise what is needed from a formal model of self-sacrifice. If one is building a self-driving car or a similar robotic system that functions in limited domains, it might be “trivial” to program in self-sacrifice, but we are seeking to understand and formalize what a model of self-sacrifice might look like in *general-purpose* autonomous robotic systems. Consider a sample scenario: A team of  $n$ , ( $n \geq 2$ ), soldiers from the *blue* team is captured by the *red* team.<sup>b</sup> The leader of the blue team is offered the choice of selecting one member from the team who will be sacrificed to free the rest of the team. Now consider the following actions:

- a<sub>1</sub>** The leader picks himself/herself.
- a<sub>2</sub>** The leader picks another soldier against their will.
- a<sub>3</sub>** The leader chooses a name randomly and it happens to be the leader’s name.
- a<sub>4</sub>** The leader chooses a name randomly and it happens to be somebody else’s name.
- a<sub>5</sub>** A soldier volunteers to die; the leader picks their name.

In addition to robotic systems with the capability for self-sacrifice in the right situations, we need systems that can understand human decisions in ethically-charged scenarios. We need a framework that can discern that: only **a<sub>1</sub>** and **a<sub>5</sub>** involve *true* self-sacrifice; **a<sub>3</sub>** is *accidental* self-sacrifice; and **a<sub>2</sub>** might be immoral.

# Logicization of Self-Reference

## Three Levels of Self-Representation

**de dicto** Agent  $r$  with the name or description  $\nu$  has come to believe on the basis of prior information  $\Gamma$  that the statement  $\phi$  holds for the agent with the name or description  $\nu$ .

$$\Gamma \vdash_r \mathbf{B} \left( l_r, \text{now}, \exists a : \text{Agent} \left[ \text{named}(a, \nu) \wedge \phi(a) \right] \right)$$

**de re** Agent  $r$  with the name or description  $\nu$  has come to believe on the basis of prior information  $\Gamma$  that the statement  $\phi$  holds of the agent with the name or description  $\nu$ .

$$\exists a : \text{Agent} \text{ named}(a, \nu) \left[ \Gamma \vdash_r \mathbf{B} \left( l_r, \text{now}, \phi(a) \right) \right]$$

**de se** Agent  $r$  believes on the basis of  $\Gamma$  that the statement  $\phi$  holds of itself  $\nu$ .

$$\Gamma \vdash_r \mathbf{B} \left( l_r, \text{now}, \phi(l_r) \right)$$



# DDE Abstracted

## 7. Formal $\mathcal{DDE}^*$

Assume we have an autonomous agent or robot  $r$  with a knowledge-base  $\Gamma$ . In Ref. 7, the predicate  $\mathcal{DDE}(\Gamma, \sigma, a, \alpha, t, H)$  is formalized — and is read as “**from a set of premises  $\Gamma$ , and in situation  $\sigma$ , we can say that action  $\alpha$  by agent  $a$  at time  $t$  operating with horizon  $H$  is  $\mathcal{DDE}$ -compliant.**” The formalization is broken up into four clauses corresponding to the informal clauses  $\mathbf{C}_1$ – $\mathbf{C}_4$  given above in Section 5:

$$\mathcal{DDE}(\Gamma, \sigma, a, \alpha, t, H) \leftrightarrow \left( \mathbf{F}_1(\Gamma, \sigma, a, \alpha, t, H) \wedge \mathbf{F}_2(\dots) \wedge \mathbf{F}_3(\dots) \wedge \mathbf{F}_4(\dots) \right)$$

higher-order 6-place relation

# DDE Abstracted

## 7. Formal $\mathcal{DDE}^*$

Assume we have an autonomous agent or robot  $r$  with a knowledge-base  $\Gamma$ . In Ref. 7, the predicate  $\mathcal{DDE}(\Gamma, \sigma, a, \alpha, t, H)$  is formalized — and is read as “**from a set of premises  $\Gamma$ , and in situation  $\sigma$ , we can say that action  $\alpha$  by agent  $a$  at time  $t$  operating with horizon  $H$  is  $\mathcal{DDE}$ -compliant.**” The formalization is broken up into four clauses corresponding to the informal clauses  $\mathbf{C}_1$ – $\mathbf{C}_4$  given above in Section 5:

$$\mathcal{DDE}(\Gamma, \sigma, a, \alpha, t, H) \leftrightarrow \left( \mathbf{F}_1(\Gamma, \sigma, a, \alpha, t, H) \wedge \mathbf{F}_2(\dots) \wedge \mathbf{F}_3(\dots) \wedge \mathbf{F}_4(\dots) \right)$$

# DDE\*

With the formal machinery at hand, enhancing  $\mathcal{DDE}$  to  $\mathcal{DDE}^*$  is straightforward. Now, corresponding to the augmented informal definition in Section 5, we take the  $\mathcal{DDE}$  predicate defined in Ref. 7 and added disjunction.

$$\mathcal{DDE}^*(\dots) \Leftrightarrow \begin{cases} \mathcal{DDE}(\Gamma, \sigma, a, \alpha, t, H) \vee \\ \mathbf{F}_1 \wedge \mathbf{F}_2 \wedge \mathbf{F}_3 \wedge \mathbf{K} \left( a, t, \left( \begin{array}{c} [\forall b. (b \neq a^*) \rightarrow \nu(\alpha, a, b, t) \gg 0] \wedge \\ \nu(\alpha, a, a^*, t) \ll 0 \end{array} \right) \right) \end{cases}$$

The disjunction simply states that the new principle  $\mathcal{DDE}^*$  applies when — (1)  $\mathcal{DDE}$  applies; or (2) when conditions  $\mathbf{F}_1$ ,  $\mathbf{F}_2$ , and  $\mathbf{F}_3$  apply along with the condition that the agent performing the action **knows** that all of the bad effects are directed toward itself, and the good effects are great in magnitude and apply only to other agents.

**Simulation:** We take a formalization of the standard trolley scenario<sup>7</sup> and

# What about DTripleE??

Peveler, M.,~Govindarajulu, N.S.\ \& Bringsjord, S.\ (2018)  
``Toward Automating the Doctrine of Triple Effect'' in  
S.\ Bringsjord, M.\ Osman Tokhi, M.\ Isabel Aldinhas Ferreira,  
N.S.\ Govindarajulu, eds., \textit{Hybrid Worlds: Societal and  
Ethical Challenges; Proceedings of the International Conference  
on Robot Ethics and Standards (ICRES) 2018}, ISBN  
978-1-9164490-1-5, London, UK: CLAWAR, pp.\ 82--88. The link  
below goes to the entire ebook in which the official version of  
this paper appears.  
<http://kryten.mm.rpi.edu/HybridWorlds.pdf>

# What about DTripleE??

ICRES 2018: International Conference on Robot  
Ethics and Standards, Troy, NY, 20-21 August 2018.  
<https://doi.org/10.13180/icres.2018.20-21.08.020>

## TOWARD AUTOMATING THE DOCTRINE OF TRIPLE EFFECT

M. PEVELER\*, N. S. GOVINDARAJULU, and S. BRINGSJORD

*Rensselaer AI & Reasoning (RAIR) Lab*

*Rensselaer Polytechnic Institute (RPI)*

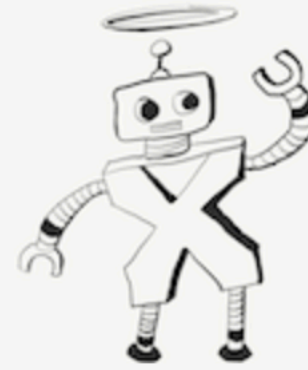
*Troy NY 12180 USA*

*\*E-mail: pevelm@rpi.com, naveensundarg@gmail.com, selmer.bringsjord@gmail.com*

The **Doctrine of Double Effect** ( $DDE$ ) is a long-studied ethical principle governing whether taking an action that has both significant positive and negative effects is ethically permissible. Unfortunately, despite its storied history,  $DDE$  does not fully account for the permissibility of actions taken in certain particularly challenging moral dilemmas that have recently arrived on the scene. The **Doctrine of Triple Effect** ( $DTE$ ) can be employed in these dilemmas, to separate the intention to perform an action *because* an effect will occur, versus *in order* for that effect to occur. This distinction allows an agent to permissibly pursue actions that may have foreseeable negative effects resulting from those actions — as long as the negative effect is not the agent's primary intention. By  $DDE$  such actions are not classified as ethically permissible. We briefly present  $DTE$  and, using a first-order multi-operator modal logic (the **deontic cognitive event calculus**), formalize this doctrine. We then give a proof-sketch of a situation for which  $DTE$  but not  $DDE$  can be used to classify a relevant action as permissible. We end with a look forward to future work.

*Keywords:* doctrine of double effect, doctrine of triple effect, machine ethics, AI





# **Making Morally Machines**

[Selmer Bringsjord](#)  $\wedge$  [Naveen Sundar Govindarajulu](#)  $\wedge$  [John Licato](#)

# Making Morally Machines

[Selmer Bringsjord](#)  $\wedge$  [Naveen Sundar Govindarajulu](#)  $\wedge$  [John Licato](#)

*er løsningen, med nok penger!*