# Epistemic Arithmetic

Eric Pacuit

University of Maryland

Lecture 2, ESSLLI 2025

July 29, 2025

# Plan

- ✓ Introduction: Smullyan's Machine
- ▶ Background
    - ✓ Formal Arithmetic
    - ✓ Gödel's Incompleteness Theorems
    - ▶ Names and Gödel numbering
    - ✓ Fixed Point Theorem
- ▶ Provability predicate and Löb's Theorem
- ▶ Provability logic
- ▶ Predicate approach to modality
- ▶ A Primer on Epistemic and Doxastic Logic
- ▶ Anti-Expert Paradoxes
- ▶ The Knower Paradox and variants
- ▶ Epistemic Arithmetic
- ▶ Gödel's Disjunction

H. Gaifman (2006). *Naming and Diagonalization, From Cantor to Gödel to Kleene.* Logic Journal of the IGPL, pp. 709 - 728.
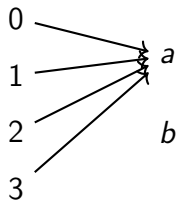
# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|   | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
|   |   |   |   |   |

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| | a | a | a | a |

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| | a | a | a | a |
| | b | b | b | b |

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|   | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
|   | a | a | a | a |
|   | b | b | b | b |
|   | a | b | a | b |

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|   | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
|   | a | a | a | a |
|   | b | b | b | b |
|   | a | b | a | b |
|   | b | a | a | a |

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|          | 0 | 1 | 2 | 3 |
|----------|---|---|---|---|
| $\alpha$ | a | a | a | a |
| $\beta$  | b | b | b | b |
| $\gamma$ | a | b | a | b |
| $\delta$ | b | a | a | a |

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|          | 0 | 1 | 2 | 3 |
|----------|---|---|---|---|
| $\alpha$ | a | a | a | a |
| $\beta$  | b | b | b | b |
| $\gamma$ | a | b | a | b |
| $\delta$ | b | a | a | a |

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

0

1        $a$

2        $b$

3

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|          | 0 | 1 | 2 | 3 |
|----------|---|---|---|---|
| $\alpha$ | a | a | a | a |
| $\beta$  | b | b | b | b |
| $\gamma$ | a | b | a | b |
| $\delta$ | b | a | a | a |

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

0

1          a

2          b

3

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|   | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| $\alpha$ | a | a | a | a |
| $\beta$ | b | b | b | b |
| $\gamma$ | a | b | a | b |
| $\delta$ | b | a | a | a |

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

0

1                    a

2                    b

3

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|          | 0 | 1 | 2 | 3 |
|----------|---|---|---|---|
| $\alpha$ | a | a | a | a |
| $\beta$  | b | b | b | b |
| $\gamma$ | a | b | a | b |
| $\delta$ | b | a | a | a |

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

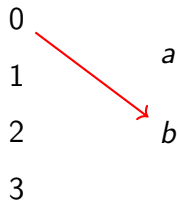# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|          | 0 | 1 | 2 | 3 |
|----------|---|---|---|---|
| $\alpha$ | a | a | a | a |
| $\beta$  | b | b | b | b |
| $\gamma$ | a | b | a | b |
| $\delta$ | b | a | a | a |

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

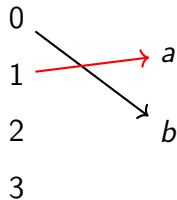# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|          | 0 | 1 | 2 | 3 |
|----------|---|---|---|---|
| $\alpha$ | a | a | a | a |
| $\beta$  | b | b | b | b |
| $\gamma$ | a | b | a | b |
| $\delta$ | b | a | a | a |

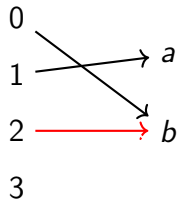$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|     | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| [0] | a | a | a | a |
| [1] | b | b | b | b |
| [2] | a | b | a | b |
| [3] | b | a | a | a |

0

                                                      *a*

1

2                              *b*

3

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|  | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| [0] | a | a | a | a |
| [1] | b | b | b | b |
| [2] | a | b | a | b |
| [3] | b | a | a | a |

$$diag(n) = \begin{cases} b & \text{if } [n](n) = a \\ a & \text{if } [n](n) = b \end{cases}$$

0

1        a

2

3        b

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|       | 0 | 1 | 2 | 3 |
|-------|---|---|---|---|
| [0]   | a | a | a | a |
| [1]   | b | b | b | b |
| [2]   | a | b | a | b |
| [3]   | b | a | a | a |

$$diag(n) = \begin{cases} b & \text{if } [n](n) = a \\ a & \text{if } [n](n) = b \end{cases}$$

0

1

2

3

a

b

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|     | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| [0] | a | a | a | a |
| [1] | b | b | b | b |
| [2] | a | b | a | b |
| [3] | b | a | a | a |

$$diag(n) = \begin{cases} b & \text{if } [n](n) = a \\ a & \text{if } [n](n) = b \end{cases}$$

0
1 → a
2
3 → b

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|     | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| [0] | a | a | a | a |
| [1] | b | b | b | b |
| [2] | a | b | a | b |
| [3] | b | a | a | a |

$$diag(n) = \begin{cases} b & \text{if } [n](n) = a \\ a & \text{if } [n](n) = b \end{cases}$$

# What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

|     | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| [0] | a | a | a | a |
| [1] | b | b | b | b |
| [2] | a | b | a | b |
| [3] | b | a | a | a |

$$diag(n) = \begin{cases} b & \text{if } [n](n) = a \\ a & \text{if } [n](n) = b \end{cases}$$

# Cantor's Diagonalization Proof

Functions from $\mathbb{N}$ to $\{0, 1\}$

|   | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ |
|---|---|---|---|---|----------|-----|----------|
| 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ |
| 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ |
| 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ |
| 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

# Cantor's Diagonalization Proof

Functions from $\mathbb{N}$ to $\{0, 1\}$



|     | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ |
|-----|---|---|---|---|----------|-----|----------|
| [0] | 0 | 0 | 0 | 0 | $\cdots$ | 0   | $\cdots$ |
| [1] | 0 | 1 | 0 | 1 | $\cdots$ | 1   | $\cdots$ |
| [2] | 0 | 1 | 1 | 0 | $\cdots$ | 0   | $\cdots$ |
| [3] | 1 | 0 | 1 | 0 | $\cdots$ | 1   | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| [$n$] | 0 | 0 | 1 | 0 | $\cdots$ | 1   | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

$d : \mathbb{N} \to \{0, 1\}$        $d(n) = 1 - [n](n)$

# Cantor's Diagonalization Proof

Functions from $\mathbb{N}$ to $\{0, 1\}$

|     | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ |
|-----|---|---|---|---|----------|-----|----------|
| [0] | 0 | 0 | 0 | 0 | $\cdots$ | 0   | $\cdots$ |
| [1] | 0 | 1 | 0 | 1 | $\cdots$ | 1   | $\cdots$ |
| [2] | 0 | 1 | 1 | 0 | $\cdots$ | 0   | $\cdots$ |
| [3] | 1 | 0 | 1 | 0 | $\cdots$ | 1   | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| [$n$] | 0 | 0 | 1 | 0 | $\cdots$ | 1   | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

$$d : \mathbb{N} \to \{0, 1\} \qquad d(n) = 1 - [n](n)$$

Then, $d \neq [n]$ for any $n \in \mathbb{N}$.

Cantor's original statement is phrased as a non-existence claim: there is no function mapping all the members of a set $S$ onto the set of all $0, 1$-valued functions over $S$. But the proof establishes a positive result: given any way of correlating functions with members of $S$, one can construct a function not correlated with any member of $S$.

(Gaiffman, p. 711)

# Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers.

# Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

# Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

Let $u_i$ be the real number defined by the $i$th definition and $f_i(n)$ be the $n$th member of the decimal expansion of $u_i$.

# Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

Let $u_i$ be the real number defined by the $i$th definition and $f_i(n)$ be the $n$th member of the decimal expansion of $u_i$.

Let $u^*$ be the number who's decimal expansion is $0.g(1)g(2)\cdots g(n)\cdots$ where $g$ is defined by $g(n) = f_n(n) + 1$ if $f_n(n) < 8$, $g(n) = 1$ otherwise.

# Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

Let $u_i$ be the real number defined by the $i$th definition and $f_i(n)$ be the $n$th member of the decimal expansion of $u_i$.

Let $u^*$ be the number who's decimal expansion is $0.g(1)g(2)\cdots g(n)\cdots$ where $g$ is defined by $g(n) = f_n(n) + 1$ if $f_n(n) < 8$, $g(n) = 1$ otherwise.

But the previous description defines a number, so $u^* = u_i$ for some $i$. But, this is impossible.

# Richard's Paradox (1905)

1. Let $A$ be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer $n$ that does not belong to $A$.

# Richard's Paradox (1905)

1. Let $A$ be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer $n$ that does not belong to $A$.

   But haven't I just defined $n$ in under 100 words?

# Richard's Paradox (1905)

1. Let $A$ be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer $n$ that does not belong to $A$.

   But haven't I just defined $n$ in under 100 words?

2. Let $B$ be the set of all reasonably interesting positive integers. Let $n$ be the smallest integer not belonging to $B$.

# Richard's Paradox (1905)

1. Let $A$ be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer $n$ that does not belong to $A$.

   But haven't I just defined $n$ in under 100 words?

2. Let $B$ be the set of all reasonably interesting positive integers. Let $n$ be the smallest integer not belonging to $B$.

   But surely this defining property of $n$ makes it reasonably interesting.

Let $f$ be a function that associates each number $x \in \mathbb{N}$ with a subset of $\mathbb{N}$, i.e., for all $x \in \mathbb{N}$, $f(x) \subseteq \mathbb{N}$.

Let $f$ be a function that associates each number $x \in \mathbb{N}$ with a subset of $\mathbb{N}$, i.e., for all $x \in \mathbb{N}$, $f(x) \subseteq \mathbb{N}$.

Define $S^*$ by:

$$x \in S^* \Leftrightarrow x \notin f(x)$$

Let $f$ be a function that associates each number $x \in \mathbb{N}$ with a subset of $\mathbb{N}$, i.e., for all $x \in \mathbb{N}$, $f(x) \subseteq \mathbb{N}$.

Define $S^*$ by:

$$x \in S^* \Leftrightarrow x \notin f(x)$$

The assumption that there is some $z$ such that $f(z) = S^*$ leads to a contradiction.

|        | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|--------|---|---|---|---|----------|-----|----------|--------------------------|
| $f(0)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0   | $\cdots$ |                          |
| $f(1)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1   | $\cdots$ |                          |
| $f(2)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0   | $\cdots$ |                          |
| $f(3)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1   | $\cdots$ |                          |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |             |
| $f(n)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1   | $\cdots$ |                          |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |             |

|       | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|-------|---|---|---|---|----------|-----|----------|--------------------------|
| $f(0)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $f(1)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | |
| $f(2)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | |
| $f(3)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |
| $f(n)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |

| | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|---|---|---|---|---|---|---|---|---|
| $f(0)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $f(1)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $f(2)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | |
| $f(3)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |
| $f(n)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |

| | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|---|---|---|---|---|---|---|---|---|
| $f(0)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $f(1)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $f(2)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | $\{1, 2\}$ |
| $f(3)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |
| $f(n)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |

| | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|---|---|---|---|---|---|---|---|---|
| $f(0)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $f(1)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $f(2)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | $\{1, 2\}$ |
| $f(3)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{0, 2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |
| $f(n)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |

|       | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|-------|---|---|---|---|----------|-----|----------|--------------------------|
| $f(0)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $f(1)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $f(2)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | $\{1, 2\}$ |
| $f(3)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{0, 2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $f(n)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

$$n \in S^* \text{ iff } n \notin f(n)$$

|  | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|---|---|---|---|---|---|---|---|---|
| $\varphi_0(x)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $\varphi_1(x)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $\varphi_2(x)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | $\{1, 2\}$ |
| $\varphi_3(x)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{0, 2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $\varphi_n(x)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

| | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|---|---|---|---|---|---|---|---|---|
| $\varphi_0(x)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $\varphi_1(x)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $\varphi_2(x)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | $\{1, 2\}$ |
| $\varphi_3(x)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{0, 2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $\varphi_n(x)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

Suppose that $S^*$ is definable in our language (say by $\varphi_m(x)$)

|  | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|---|---|---|---|---|---|---|---|---|
| $\varphi_0(x)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $\varphi_1(x)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $\varphi_2(x)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | $\{1, 2\}$ |
| $\varphi_3(x)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{0, 2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $\varphi_n(x)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

Write $\varphi_m(\bar{n})$ for "$\varphi_m(x)$ is true of $n$"

| | 0 | 1 | 2 | 3 | $\cdots$ | $n$ | $\cdots$ | $S \subseteq \mathbb{N}$ |
|---|---|---|---|---|---|---|---|---|
| $\varphi_0(x)$ | 0 | 0 | 0 | 0 | $\cdots$ | 0 | $\cdots$ | $\emptyset$ |
| $\varphi_1(x)$ | 0 | 1 | 0 | 1 | $\cdots$ | 1 | $\cdots$ | $\{1, 3, \ldots, n, \ldots\}$ |
| $\varphi_2(x)$ | 0 | 1 | 1 | 0 | $\cdots$ | 0 | $\cdots$ | $\{1, 2\}$ |
| $\varphi_3(x)$ | 1 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{0, 2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $\varphi_n(x)$ | 0 | 0 | 1 | 0 | $\cdots$ | 1 | $\cdots$ | $\{2, \ldots, n, \ldots\}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

$$\varphi_m(\overline{n}) \leftrightarrow \neg \text{True}(\ulcorner \varphi_n(\overline{n}) \urcorner)$$

where $\ulcorner \varphi_n(\overline{n}) \urcorner$ is the term in the language representing the code of $\varphi_n(\overline{n})$

# D-Liar

$$\varphi_m(\overline{m}) \leftrightarrow \neg\text{True}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

"$m$ is true of $\varphi_m(x)$ iff it is not true that $m$ is true of $\varphi_m(x)$"

# Gödel's Idea

$$\varphi_m(\overline{m}) \leftrightarrow \neg\mathsf{True}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

# Gödel's Idea

$$\varphi_m(\overline{m}) \leftrightarrow \neg\mathsf{True}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

$$\varphi_m(\overline{m}) \leftrightarrow \neg\mathsf{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

"$\varphi_m(\overline{m})$ is true iff $\varphi_m(\overline{m})$ is not provable."

# Gödel's Idea

$$\varphi_m(\overline{m}) \leftrightarrow \neg\text{True}(\ulcorner\varphi_m(\overline{m})\urcorner)$$

$$\varphi_m(\overline{m}) \leftrightarrow \neg\text{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$$

"$\varphi_m(\overline{m})$ is true iff $\varphi_m(\overline{m})$ is not provable."

$$\varphi_m(\overline{m}) \leftrightarrow \neg\mathsf{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$$

$\varphi_m(\overline{m}) \leftrightarrow \neg\mathsf{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$

$\varphi_m(\overline{m})$ is not provable: Suppose $\varphi_m(\overline{m})$ is provable. Then, since we can only prove true statements, $\varphi_m(\overline{m})$ is true. This means that $\neg\mathsf{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$ is true. So, $\varphi_m(\overline{m})$ is not provable. Contradiction.

$$\varphi_m(\overline{m}) \leftrightarrow \neg\mathsf{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$$

$\varphi_m(\overline{m})$ is not provable: Suppose $\varphi_m(\overline{m})$ is provable. Then, since we can only prove true statements, $\varphi_m(\overline{m})$ is true. This means that $\neg\mathsf{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$ is true. So, $\varphi_m(\overline{m})$ is not provable. Contradiction.

$\neg\varphi_m(\overline{m})$ is not provable: Suppose that $\neg\varphi_m(\overline{m})$ is provable. Since our system only proves true statements, $\neg\varphi_m(\overline{m})$ is true. Then $\neg\neg\mathsf{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$ is true. So, $\varphi_m(\overline{m})$ is provable. This contradicts the assumption that the system is consistent.

$$\varphi_m(\overline{m}) \leftrightarrow \neg\text{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$$

$\varphi_m(\overline{m})$ is not provable: Suppose $\varphi_m(\overline{m})$ is provable. Then, since we can only prove true statements, $\varphi_m(\overline{m})$ is true. This means that $\neg\text{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$ is true. So, $\varphi_m(\overline{m})$ is not provable. Contradiction.

$\neg\varphi_m(\overline{m})$ is not provable: Suppose that $\neg\varphi_m(\overline{m})$ is provable. Since our system only proves true statements, $\neg\varphi_m(\overline{m})$ is true. Then $\neg\neg\text{Prov}(\ulcorner\varphi_m(\overline{m})\urcorner)$ is true. So, $\varphi_m(\overline{m})$ is provable. This contradicts the assumption that the system is consistent.

**Conclusion**: Neither $\varphi_m(\overline{m})$ nor $\neg\varphi_m(\overline{m})$ is provable.

$$\varphi_m(\overline{m}) \leftrightarrow \neg\mathsf{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

1. Apply Richard's move to Cantor's construction to get the D-Liar

2. Replace 'true' with 'provable' on the right-hand side of the sentence

3. Proceed with the difficult task of *arithmetizing syntax* to construct the right-side of the sentence $(\mathsf{Prov}(v))$.

4. Show that the above sentence is provable within the formal system eliminating any appeal to the concept of "truth". The assumption that provable implies truth is replaced with $(\omega\text{-})$consistency.

H. Gaifman (2006). *Naming and Diagonalization, From Cantor to Gödel to Kleene*. Logic Journal of the IGPL, pp. 709 - 728.

# Naming systems

Naming systems are intended as a basic framework for studying situations in which functions can be applied to their names....In a naming system we do not specify how the names are attached to functions, we assume only that there is such a correlation and that it satisfies certain minimal requirements.

H. Gaifman (2006). *Naming and Diagonalization, From Cantor to Gödel to Kleene*. Logic Journal of the IGPL, pp. 709 - 728.

# Naming systems I

$$\mathcal{D} = (D, type, \{\ \})$$

such that:

- $D$ is a non-empty set.
- *type* assigns to each $a \in D$ its type: $type(a)$ tells us if $a$ is a name (of a function) and, if it is, the function's arity.

  A name of arity $n$, or $n$-ary name, is one that names an $n$-ary function.

  Types can be construed as tuples: $(0)$—if $a$ is not a name, $(1, n)$—if it is an $n$-ary name.

- $\{\ \}$ is a mapping that assigns to every $n$-ary name, $a$, a function:

$$\{a\} : D^n \to D$$

# Naming systems II

- There is at least one named function of arity greater than 0

# Naming systems II

▶ There is at least one named function of arity greater than 0

▶ Substitution of names (SN): If $f$ is an $n$-ary named function, where $n > 0$, then, for every name $a$:

$$\lambda x_2, \ldots x_n f(a, x_2, \ldots, x_n) \text{ is named}$$

# Naming systems II

- There is at least one named function of arity greater than 0

- Substitution of names (SN): If $f$ is an $n$-ary named function, where $n > 0$, then, for every name $a$:

$$\lambda x_2, \ldots x_n f(a, x_2, \ldots, x_n) \text{ is named}$$

- Variable permutation (VP): If $f$ is an $n$-ary named function, where $n > 0$, and $\pi$ is a permutation of $\{1, \ldots, n\}$, then

$$\lambda x_1, \ldots x_n f(x_{\pi(1)}, x_{\pi(2)}, \ldots, x_{\pi(n)}) \text{ is named}$$

## *n*-Diagonal Function

For $n > 0$, an *n-diagonal function*, denoted $dl_n$, is a function that maps each *n*-ary name $a$ to a name of the function:

$$\lambda x_2, \ldots, x_n \{a\}(a, x_2, \ldots, x_n)$$

Thus, $dl_n(a)$ is the name of the above function.

## n-Diagonal Function

For $n > 0$, an *n-diagonal function*, denoted $dl_n$, is a function that maps each *n*-ary name $a$ to a name of the function:

$$\lambda x_2, \ldots, x_n \{a\}(a, x_2, \ldots, x_n)$$

Thus, $dl_n(a)$ is the name of the above function.

For all *n*-ary names $a$,

$$\{dl_n(a)\}(x_2, \ldots, x_n) = \{a\}(a, x_2, \ldots, x_n)$$

# General Fixed-Point Theorem

**GFP Theorem**. If $F$ is an $(n+1)$-ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ is named, then there is an $n$-ary name, $e$, such that:

$$\{e\}(x_1, \ldots, x_n) = F(e, x_1, \ldots, x_n)$$

## General Fixed-Point Theorem

**GFP Theorem**. If $F$ is an $(n+1)$-ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ is named, then there is an $n$-ary name, $e$, such that:

$$\{e\}(x_1, \ldots, x_n) = F(e, x_1, \ldots, x_n)$$

*Proof.* Let $c$ be the name of $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \ldots, x_n)$,

## General Fixed-Point Theorem

**GFP Theorem**. If $F$ is an $(n+1)$-ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ is named, then there is an $n$-ary name, $e$, such that:

$$\{e\}(x_1, \ldots, x_n) = F(e, x_1, \ldots, x_n)$$

*Proof.* Let $c$ be the name of $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \ldots, x_n)$,

$$\{e\}(\vec{x}) \;=\; \{dl_{n+1}(c)\}(\vec{x}) \quad \text{(definition of } e)$$

## General Fixed-Point Theorem

**GFP Theorem**. If $F$ is an $(n+1)$-ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ is named, then there is an $n$-ary name, $e$, such that:

$$\{e\}(x_1, \ldots, x_n) = F(e, x_1, \ldots, x_n)$$

*Proof.* Let $c$ be the name of $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \ldots, x_n)$,

$$
\begin{aligned}
\{e\}(\vec{x}) &= \{dl_{n+1}(c)\}(\vec{x}) &&\text{(definition of } e) \\
&= \{c\}(c, \vec{x}) &&\text{(definition of } dl_{n+1}(c))
\end{aligned}
$$

## General Fixed-Point Theorem

**GFP Theorem**. If $F$ is an $(n+1)$-ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ is named, then there is an $n$-ary name, $e$, such that:

$$\{e\}(x_1, \ldots, x_n) = F(e, x_1, \ldots, x_n)$$

*Proof.* Let $c$ be the name of $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \ldots, x_n)$,

$$
\begin{aligned}
\{e\}(\vec{x}) &= \{dl_{n+1}(c)\}(\vec{x}) && \text{(definition of } e) \\
&= \{c\}(c, \vec{x}) && \text{(definition of } dl_{n+1}(c)) \\
&= F(dl_{n+1}(c), \vec{x}) && \text{(definition of } c)
\end{aligned}
$$

## General Fixed-Point Theorem

**GFP Theorem**. If $F$ is an $(n+1)$-ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ is named, then there is an $n$-ary name, $e$, such that:

$$\{e\}(x_1, \ldots, x_n) = F(e, x_1, \ldots, x_n)$$

*Proof.* Let $c$ be the name of $F(dl_{n+1}(x_0), x_1, \ldots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \ldots, x_n)$,

$$
\begin{aligned}
\{e\}(\vec{x}) &= \{dl_{n+1}(c)\}(\vec{x}) && \text{(definition of } e\text{)} \\
&= \{c\}(c, \vec{x}) && \text{(definition of } dl_{n+1}(c)\text{)} \\
&= F(dl_{n+1}(c), \vec{x}) && \text{(definition of } c\text{)} \\
&= F(e, \vec{x}) && \text{(definition of } e\text{)}
\end{aligned}
$$

- ✓ Gödel numbering
- ✓ Gödel-Carnap Fixed Point Theorem
- ✓ (Naming systems)
- ▶ Representing functions/relations

# Representability

### Definition

Suppose that $f : \mathbb{N}^k \to \mathbb{N}$. We say that $f$ is **representable** in $\mathbf{Q}$ when there is a formula $A_f(x_0, \ldots, x_{k-1}, y)$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$: if $f(n_0, \ldots, n_{k-1}) = m$ then

1. $\mathbf{Q} \vdash A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, \overline{m})$
2. $\mathbf{Q} \vdash \forall y (A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, y) \to y = \overline{m})$

# Equivalent definitions of representability

- $f$ is representable in $\mathbf{Q}$ iff there is a formula $A_f(x_0, \ldots, x_{k-1}, y)$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$, if $f(n_0, \ldots, n_{k-1}) = m$ then:

$$\mathbf{Q} \vdash \forall y (A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, y) \leftrightarrow y = \overline{m})$$

## Equivalent definitions of representability

▶ $f$ is representable in **Q** iff there is a formula $A_f(x_0, \ldots, x_{k-1}, y)$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$, if $f(n_0, \ldots, n_{k-1}) = m$ then:

$$\mathbf{Q} \vdash \forall y(A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, y) \leftrightarrow y = \overline{m})$$

▶ $f$ is representable in **Q** iff there is a formula $A_f(x_0, \ldots, x_{k-1}, y)$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$:

1. If $f(n_0, \ldots, n_{k-1}) = m$, then $\mathbf{Q} \vdash A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, \overline{m})$
2. If $f(n_0, \ldots, n_{k-1}) \neq m$, then $\mathbf{Q} \vdash \neg A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, \overline{m})$

## Equivalent definitions of representability

▶ $f$ is representable in **Q** iff there is a formula $A_f(x_0, \ldots, x_{k-1}, y)$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$, if $f(n_0, \ldots, n_{k-1}) = m$ then:

$$\mathbf{Q} \vdash \forall y(A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, y) \leftrightarrow y = \overline{m})$$

▶ $f$ is representable in **Q** iff there is a formula $A_f(x_0, \ldots, x_{k-1}, y)$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$:
  1. If $f(n_0, \ldots, n_{k-1}) = m$, then $\mathbf{Q} \vdash A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, \overline{m})$
  2. If $f(n_0, \ldots, n_{k-1}) \neq m$, then $\mathbf{Q} \vdash \neg A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, \overline{m})$

▶ $f$ is representable in **Q** iff there is a formula $A_f(x_0, \ldots, x_{k-1}, y)$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$:
  1. if $f(n_0, \ldots, n_{k-1}) = m$ then $\mathbf{Q} \vdash A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, \overline{m})$
  2. $\mathbf{Q} \vdash \exists! y A_f(\overline{n_0}, \ldots, \overline{n_{k-1}}, y)$

# Exercise

Prove that all of the definitions of representability are equivalent.

# Representing Relations

A relation $R \subseteq \mathbb{N}^k$ is **representable** in $\mathbf{Q}$ provided that the characteristic function $\chi_R$ is representable in $\mathbf{Q}$. It is not hard to see that this is equivalent to saying that $R \subseteq \mathbb{N}^k$ is representable in $\mathbf{Q}$ provided that there is a formula $A_R$ such that for all $n_0, \ldots, n_{k-1} \in \mathbb{N}$:

1. if $(n_0, \ldots, n_{k-1}) \in R$, then $\mathbf{Q} \vdash A_R(\overline{n_0}, \ldots, \overline{n_{k-1}})$
2. if $(n_0, \ldots, n_{k-1}) \notin R$, then $\mathbf{Q} \vdash \neg A_R(\overline{n_0}, \ldots, \overline{n_{k-1}})$

All of the following relations are representable in **Q**:

- $Sent(x)$: $x$ is the Gödel number of a sentence of $\mathcal{L}_A$
- $Form(x)$: $x$ is the Gödel number of a formula of $\mathcal{L}_A$
- $Term(x)$: $x$ is the Gödel number of a term of $\mathcal{L}_A$
- $Axiom(x)$: $x$ is the Gödel number of an axiom of **Q**
- $Prf_{\textbf{PA}}(x, y)$: $x$ is the Gödel number of a derivation in **PA** of a formula with Gödel number $y$.
- $\cdots$

# Plan

- ✓ Introduction: Smullyan's Machine
- ✓ Background
    - ✓ Formal Arithmetic
    - ✓ Gödel's Incompleteness Theorems
    - ✓ Names and Gödel numbering
    - ✓ Fixed Point Theorem
- ▶ Provability predicate and Löb's Theorem
- ▶ Provability logic
- ▶ Predicate approach to modality
- ▶ A Primer on Epistemic and Doxastic Logic
- ▶ Anti-Expert Paradoxes
- ▶ The Knower Paradox and variants
- ▶ Epistemic Arithmetic
- ▶ Gödel's Disjunction

# Proof Predicate

The proof relation $Prf_{\mathbf{PA}}(x, y)$ is represented by a formula $\mathrm{Prf}_{\mathbf{PA}}$.

# Proof Predicate

The proof relation $Prf_{\mathbf{PA}}(x, y)$ is represented by a formula $\mathrm{Prf}_{\mathbf{PA}}$.

The *proof predicate*, denoted $\mathrm{Prov}_{\mathbf{PA}}(y)$, is defined as follows:

$$\exists x \mathrm{Prf}_{\mathbf{PA}}(x, y)$$

# Derivability Conditions

It can be shown that the provability predicate $\text{Prov}_{\textbf{PA}}$ satisfies the following:

$D1.$ If $\textbf{PA} \vdash A$, then $\textbf{PA} \vdash \text{Prov}_{\textbf{PA}}(\ulcorner A \urcorner)$

$D2.$ $\textbf{PA} \vdash \text{Prov}_{\textbf{PA}}(\ulcorner A \to B \urcorner) \to (\text{Prov}_{\textbf{PA}}(\ulcorner A \urcorner) \to \text{Prov}_{\textbf{PA}}(\ulcorner B \urcorner))$

$D3.$ $\textbf{PA} \vdash \text{Prov}_{\textbf{PA}}(\ulcorner A \urcorner) \to \text{Prov}_{\textbf{PA}}(\ulcorner \text{Prov}_{\textbf{PA}}(\ulcorner A \urcorner) \urcorner)$

# Derivability Conditions

A provability predicate for **T**, denoted $\text{Prov}_{\mathbf{T}}$, satisfies the following:

$D1.$ If $\mathbf{T} \vdash A$, then $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner)$

$D2.$ $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \rightarrow B \urcorner) \rightarrow (\text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{T}}(\ulcorner B \urcorner))$

$D3.$ $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{T}}(\ulcorner \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \urcorner)$

# Reflection Principle

The reflection principle for **T** is the schema

$$\text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \to A$$

# Monotonicity Inference for the Provability Predicate

### Lemma

For any theory **T**, if $\text{Prov}_\mathbf{T}$ satisfies $D1$ and $D2$, then:

$$\text{From } \mathbf{T} \vdash A \rightarrow B, \text{ infer } \mathbf{T} \vdash \text{Prov}_\mathbf{T}(\ulcorner A \urcorner) \rightarrow \text{Prov}(\ulcorner B \urcorner).$$

# Löb's Theorem

### Theorem (Löb's Theorem)

Let **T** be an axiomatizable theory extending **Q**, and suppose $\text{Prov}_{\mathbf{T}}(y)$ is a formula satisfying conditions $D1$-$D3$.

$$\text{If } \mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \to A, \text{ then } \mathbf{T} \vdash A.$$

Suppose $A$ is a sentence such that $\mathbf{T} \vdash \mathrm{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \to A$. Let $B(y)$ be the formula

$$\mathrm{Prov}_{\mathbf{T}}(y) \to A$$

Suppose $A$ is a sentence such that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \to A$. Let $B(y)$ be the formula

$$\text{Prov}_{\mathbf{T}}(y) \to A$$

By the Fixed-Point Theorem, there is a sentence $D$ such that

$$\mathbf{T} \vdash D \leftrightarrow B(\ulcorner D \urcorner)$$

Suppose that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \to A$.

Suppose $A$ is a sentence such that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \to A$. Let $B(y)$ be the formula

$$\text{Prov}_{\mathbf{T}}(y) \to A$$

By the Fixed-Point Theorem, there is a sentence $D$ such that

$$\mathbf{T} \vdash D \leftrightarrow B(\ulcorner D \urcorner)$$

Suppose that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \to A$.

To simplify the notation, we write $\text{Prov}(\cdot)$ instead of $\text{Prov}_{\mathbf{T}}$

1. $D \leftrightarrow (\mathrm{Prov}(\ulcorner D \urcorner) \to A)$                                FPT

2. $\mathrm{Prov}(\ulcorner D \urcorner) \to \mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \to A \urcorner)$               Lemma: 1

3. $\mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \to A \urcorner) \to (\mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner))$    D2

4. $\mathrm{Prov}(\ulcorner D \urcorner) \to (\mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner))$             PC: 2, 3

1.  $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$                                                  FPT

⋮   ⋮                                                                                                                                    ⋮

4.  $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$   PC: 2, 3

5.  $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$                                  D3

1. $D \leftrightarrow (\mathsf{Prov}(\ulcorner D \urcorner) \to A)$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀FPT

⋮ ⠀⠀⋮ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⋮

4. $\mathsf{Prov}(\ulcorner D \urcorner) \to (\mathsf{Prov}(\ulcorner \mathsf{Prov}(\ulcorner D \urcorner) \urcorner) \to \mathsf{Prov}(\ulcorner A \urcorner))$ ⠀⠀PC: 2, 3

5. $\mathsf{Prov}(\ulcorner D \urcorner) \to \mathsf{Prov}(\ulcorner \mathsf{Prov}(\ulcorner D \urcorner) \urcorner)$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀D3

6. $\mathsf{Prov}(\ulcorner D \urcorner) \to \mathsf{Prov}(\ulcorner A \urcorner)$ ⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀⠀PC: 4, 5

1. $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$           FPT

⋮    ⋮                                              ⋮

4. $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$    PC: 2, 3

5. $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$             D3

6. $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$                        PC: 4, 5

7. $\text{Prov}(\ulcorner A \urcorner) \rightarrow A$                                Assumption

| | | |
|---|---|---|
| 1. | $D \leftrightarrow (\mathrm{Prov}(\ulcorner D \urcorner) \to A)$ | FPT |
| $\vdots$ | $\vdots$ | $\vdots$ |
| 4. | $\mathrm{Prov}(\ulcorner D \urcorner) \to (\mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\mathrm{Prov}(\ulcorner D \urcorner) \to \mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |
| 6. | $\mathrm{Prov}(\ulcorner D \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner)$ | PC: 4, 5 |
| 7. | $\mathrm{Prov}(\ulcorner A \urcorner) \to A$ | Assumption |
| 8. | $\mathrm{Prov}(\ulcorner D \urcorner) \to A$ | PC: 6, 7 |

| 1. | $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \to A)$ | FPT |
| :--- | :--- | :--- |
| $\vdots$ | $\vdots$ | $\vdots$ |
| 4. | $\text{Prov}(\ulcorner D \urcorner) \to (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \to \text{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\text{Prov}(\ulcorner D \urcorner) \to \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |
| 6. | $\text{Prov}(\ulcorner D \urcorner) \to \text{Prov}(\ulcorner A \urcorner)$ | PC: 4, 5 |
| 7. | $\text{Prov}(\ulcorner A \urcorner) \to A$ | Assumption |
| 8. | $\text{Prov}(\ulcorner D \urcorner) \to A$ | PC: 6, 7 |
| 9. | $D$ | PC: 1, 8 |

| | | |
|---|---|---|
| 1. | $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$ | FPT |
| $\vdots$ | $\vdots$ | $\vdots$ |
| 4. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |
| 6. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$ | PC: 4, 5 |
| 7. | $\text{Prov}(\ulcorner A \urcorner) \rightarrow A$ | Assumption |
| 8. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow A$ | PC: 6, 7 |
| 9. | $D$ | PC: 1, 8 |
| 10. | $\text{Prov}(\ulcorner D \urcorner)$ | D1 from 9 |

32

| | | |
|---|---|---|
| 1. | $D \leftrightarrow (\mathrm{Prov}(\ulcorner D \urcorner) \to A)$ | FPT |
| ⋮ | ⋮ | ⋮ |
| 4. | $\mathrm{Prov}(\ulcorner D \urcorner) \to (\mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\mathrm{Prov}(\ulcorner D \urcorner) \to \mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |
| 6. | $\mathrm{Prov}(\ulcorner D \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner)$ | PC: 4, 5 |
| 7. | $\mathrm{Prov}(\ulcorner A \urcorner) \to A$ | Assumption |
| 8. | $\mathrm{Prov}(\ulcorner D \urcorner) \to A$ | PC: 6, 7 |
| 9. | $D$ | PC: 1, 8 |
| 10. | $\mathrm{Prov}(\ulcorner D \urcorner)$ | D1 from 9 |
| 11. | $A$ | PC: 8, 10 |

32

1. $D \leftrightarrow (\mathrm{Prov}(\ulcorner D \urcorner) \to A)$             FPT

$\vdots$    $\vdots$                                                   $\vdots$

4. $\mathrm{Prov}(\ulcorner D \urcorner) \to (\mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner))$    PC: 2, 3

5. $\mathrm{Prov}(\ulcorner D \urcorner) \to \mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner D \urcorner) \urcorner)$                D3

6. $\mathrm{Prov}(\ulcorner D \urcorner) \to \mathrm{Prov}(\ulcorner A \urcorner)$                     PC: 4, 5

7. $\mathrm{Prov}(\ulcorner A \urcorner) \to A$                               Assumption

8. $\mathrm{Prov}(\ulcorner D \urcorner) \to A$                              PC: 6, 7

9. $D$                                            PC: 1, 8

10. $\mathrm{Prov}(\ulcorner D \urcorner)$                                 D1 from 9

11. $A$                                             PC: 8, 10

# 'PA couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences $\varphi$,

$$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi \urcorner)$$

# 'PA couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences $\varphi$,

$$\textbf{PA} \vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner \varphi \urcorner) \to \varphi \urcorner)$$

Statement

$\textbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner)$
    implies $\textbf{PA} \vdash \varphi$

It is not true that...

$\textbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner) \to \varphi$

# '**PA** couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences $\varphi$,

$$\textbf{PA} \vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner \varphi \urcorner) \to \varphi \urcorner)$$

| Statement | It is not true that... |
|---|---|
| $\textbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner)$ implies $\textbf{PA} \vdash \varphi$ | $\textbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner) \to \varphi$ |
| $\textbf{PA} \vdash \text{Prov}(\ulcorner \neg\varphi \urcorner)$ implies $\textbf{PA} \not\vdash \text{Prov}(\ulcorner \varphi \urcorner)$ | $\textbf{PA} \vdash \text{Prov}(\ulcorner \neg\varphi \urcorner) \to \neg\text{Prov}(\ulcorner \varphi \urcorner)$ |

## '**PA** couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences $\varphi$,

$$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \mathrm{Prov}(\ulcorner \varphi \urcorner) \to \varphi \urcorner)$$

Statement

$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \varphi \urcorner)$
  implies $\mathbf{PA} \vdash \varphi$

$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \neg\varphi \urcorner)$
  implies $\mathbf{PA} \not\vdash \mathrm{Prov}(\ulcorner \varphi \urcorner)$

$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \neg\mathrm{Prov}(\ulcorner \varphi \urcorner) \urcorner)$
  implies $\mathbf{PA} \vdash \neg\mathrm{Prov}(\varphi)$

It is not true that...

$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \varphi \urcorner) \to \varphi$

$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \neg\varphi \urcorner) \to \neg\mathrm{Prov}(\ulcorner \varphi \urcorner)$

$\mathbf{PA} \vdash \mathrm{Prov}(\ulcorner \neg\mathrm{Prov}(\ulcorner \varphi \urcorner) \urcorner) \to \neg\mathrm{Prov}(\ulcorner \varphi \urcorner)$