

Static Analysis for Data Science

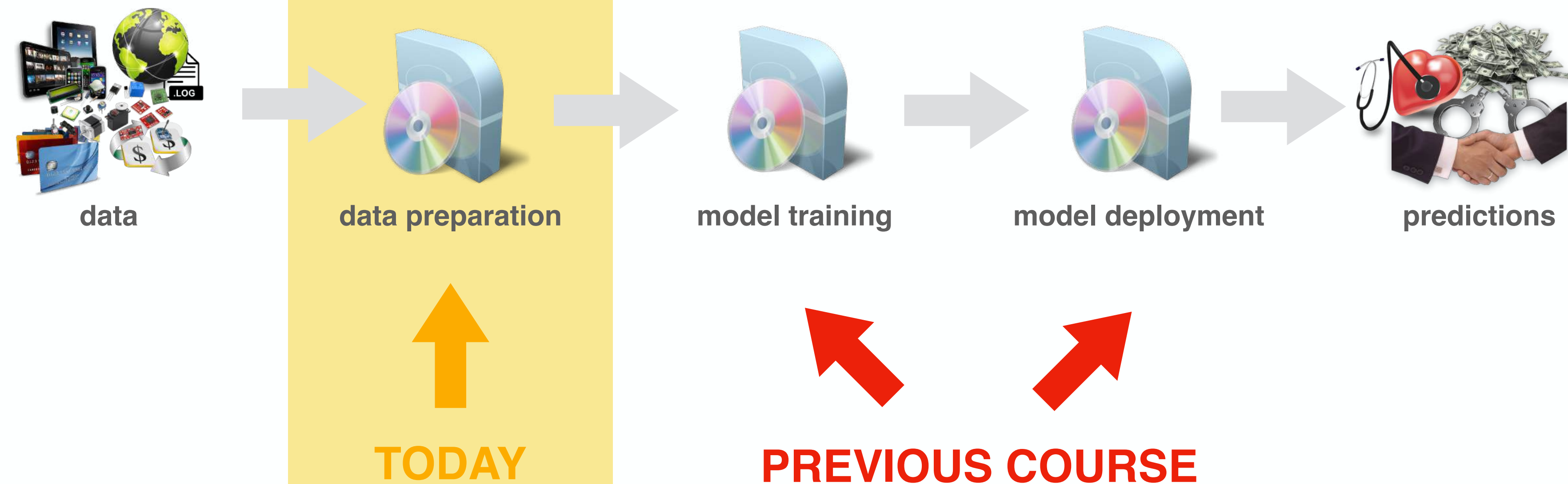
**MPRI 2-6: Abstract Interpretation,
Application to Verification and Static Analysis**

Caterina Urban

January 20th, 2024

Year 2024-2025

Machine Learning Development Pipeline



Static Analysis for Data Preparation

Machine Learning Development Pipeline

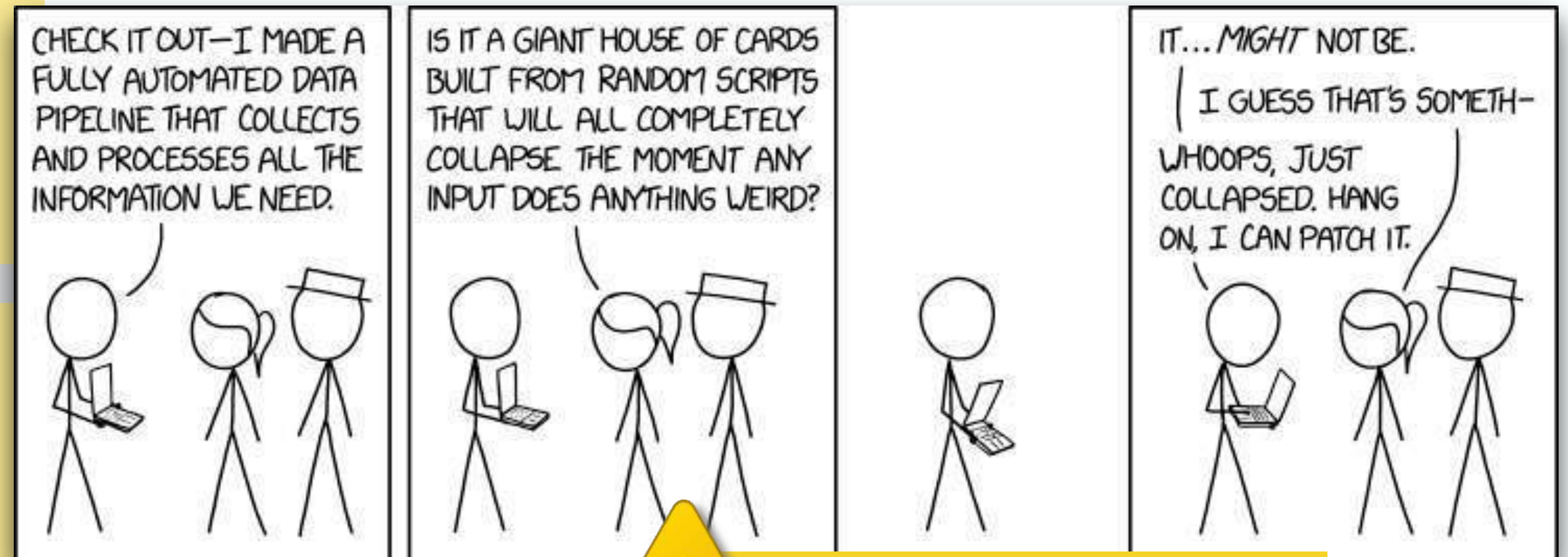
Data Preparation is Fragile



data



data preparation



<https://xkcd.com/2054/>

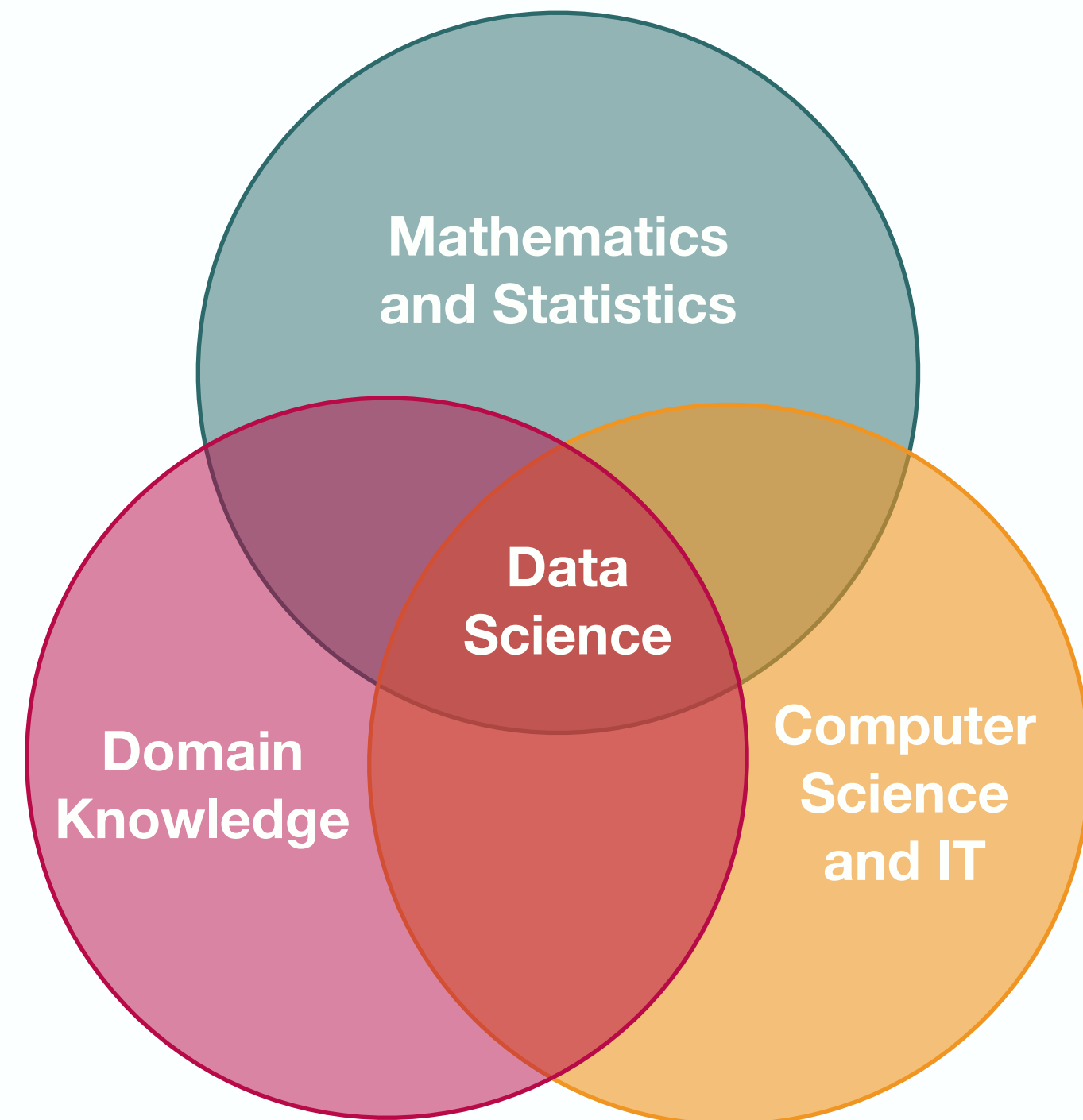


insidious silent bugs

Data Scientists

Data Scientist: The Sexiest Job of the 21st Century

Andrew McAfee and Erik Brynjolfsson



Andrew J Buboltz, silk screen on a page from a high school yearbook, 8.5" x 12", 2011 Tamar Cohen

When Jonathan Goldman arrived for work in June 2006 at [LinkedIn, the business networking site](#), the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know

like arriving at a conference reception and realizing you don't know
missing in the social experience. As one LinkedIn manager put it, "It was

Jupyter Notebooks

The main image shows a Jupyter Notebook interface with the following code and output:

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
```

Out [2]:

Name	Q1	Q2	Q3

Below the main image, there are three overlapping images:

- Netflix Technology Blog:** Article titled "Beyond Interactive: Notebook Innovation at Netflix" by Michelle Ufford, M Pacer, Matthew Seal, and Kyle Kelley. The article discusses the growing popularity of notebooks among data scientists and their use in production environments.
- Databricks:** Landing page titled "Databricks for Data Science" with the tagline "An open and unified platform to collaboratively run all types of analytics workloads, from data preparation to exploratory analysis and predictive analytics, at scale." It features a "Loan Analysis" notebook preview.
- Reddit:** A comment thread discussing Jupyter's use in production environments. A user named "u/Desperate-Walk1780" comments that their team uses Jupyter on CentOS in production, and another user "u/EnricoT0" mentions their former employer uses notebooks for all data science tasks.

Jupyter Notebooks

```
[1] 1 d = genfromtxt('data.csv')
[2] 1 selector = SelectKBest(k=25)
    2 x = selector.fit_transform(d)
[3] 1 x = genfromtxt('data2.csv')
[4] 1 x_train, x_test, y_train, y_test =
    2 train_test_split(x, ...)
[5] 1 lr = LogisticRegression()
    2 lr.fit(x_train, y_train)
    3 y_pred = lr.predict(x_test)
```

UNUSED DATA

Jupyter Notebooks


```
[1] 1 d = genfromtxt('data.csv')
```

```
[2] 1 selector = SelectKBest(k=25)  
   2 x = selector.fit_transform(d)
```

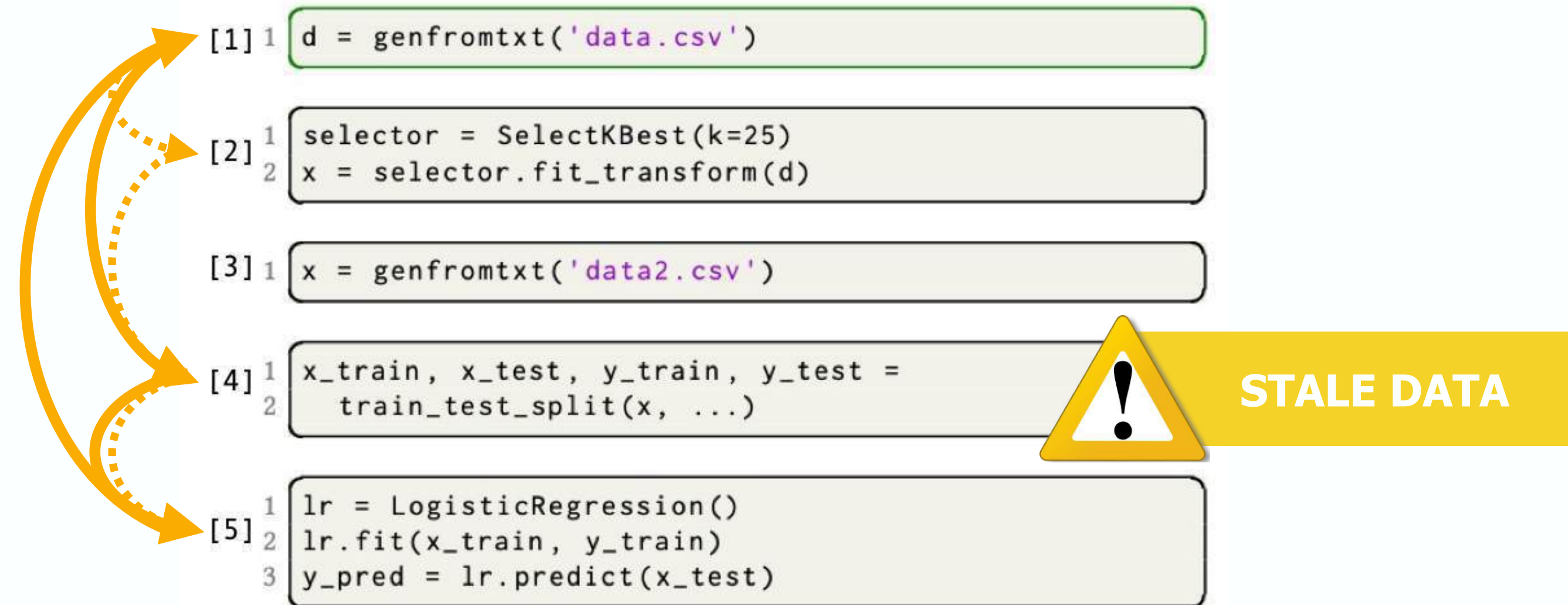
```
[3] 1 x = genfromtxt('data2.csv')
```

```
[4] 1 x_train, x_test, y_train, y_test =  
   2 train_test_split(x, ...)
```

```
[5] 1 lr = LogisticRegression()  
   2 lr.fit(x_train, y_train)  
   3 y_pred = lr.predict(x_test)
```



Jupyter Notebooks





Anomalously Unused Data

The Reinhart-Rogoff Paper

FAQ: Reinhart, Rogoff, and the Excel Error That Changed History

By Peter Coy  | April 18, 2013

The Excel Depression

By PAUL KRUGMAN
Published: April 18, 2013 |  470 Comments



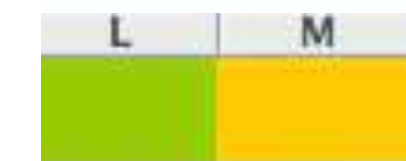
In this age of information, math errors can lead to disaster. NASA's Mars Orbiter crashed because engineers forgot to convert to metric measurements; JPMorgan Chase's "London Whale" venture went bad in part because modelers divided by a sum instead of an average. So, did an Excel coding error destroy the economies of the Western world?









 Enlarge This Image



The story so far: At the beginning of 2010, two Harvard economists, Carmen Reinhart and Kenneth Rogoff, circulated a paper, "Growth in a Time of Debt," that purported to identify a critical "threshold," a tipping point, for government indebtedness. Once debt exceeds 90 percent of gross domestic product, they claimed, economic growth drops off sharply.

Ms. Reinhart and Mr. Rogoff had credibility thanks to a



-  FACEBOOK
-  TWITTER
-  GOOGLE+
-  SAVE
-  EMAIL
-  SHARE
-  PRINT
-  REPRINTS

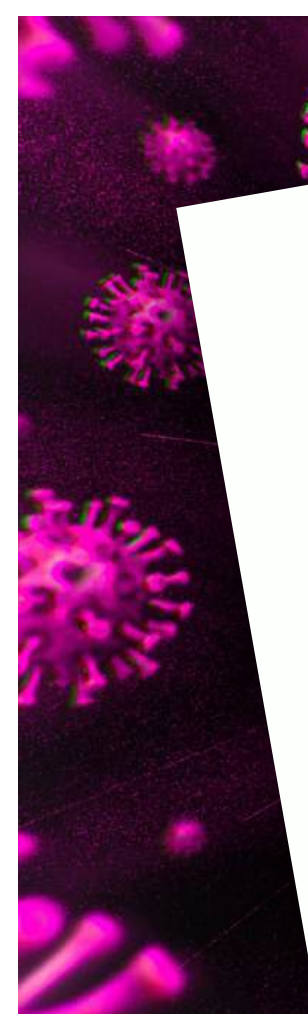
England Covid-19 Cases Error

SCIENCE / US & WORLD / TECH

Excel spreadsheet error blamed for UK's 16,000 missing coronavirus cases

The case went missing after the spreadsheet hit its filesize limit

By James Vincent | Oct 5, 2020, 9:41am EDT



Check for updates

The BMJ
Cite this as: *BMJ* 2020;371:m3891
<http://dx.doi.org/10.1136/bmj.m3891>
Published: 06 October 2020

Covid-19: Only half of 16 000 patients missed from England's official figures have been contacted

Elisabeth Mahase

Details of nearly 16 000 cases of covid-19 were not transferred to England's NHS Test and Trace service and were missed from official figures because of an error in the process for updating the data.

England's health and social care secretary, Matt Hancock, told the House of Commons on Monday 5 October that after the error was discovered on Friday 2 October "6500 hours of extra contact tracing" had been carried out over the weekend. But as at Monday morning only half (51%) of the people had been reached by contact tracers.

...ence. Labour's shadow health secretary, ...id. "Thousands of people are ...ced to covid,

data and furthermore have issued guidance on validation and risk management for these products if they are to be used in such a safety critical manner."

The error came as the Labour Party's leader, Keir Starmer, said that the prime minister had "lost control" of covid-19, with no clear strategy for beating it. Speaking to the *Observer*, Starmer set out his five point plan for covid-19, which starts with publishing the criteria for local restrictions, as the German government did. Secondly, he said public health messaging should be improved by adding a feature to the NHS covid-19 app so people can search their postcode and find out their local restrictions.

Starmer has also said he would fix the contact tracing system by investing in NHS and university ... expand testing and at the same time ... in charge of contact ... high

NEWS

BMJ: first published as 10.1136/bmj.m3891 on 6 October 2020. D

Example

```
english = bool(input())
math = bool(input())
science = bool(input())
bonus = bool(input())

passing = True
if not english:
    english = False
if not math:
    passing = False or bonus
if not math:
    passing = False or bonus

print(passing)
```

INPUT VARIABLES

ERROR: english SHOULD BE passing

ERROR: math SHOULD BE science

OUTPUT VARIABLES



the input variables **english** and **science** are unused

Data Usage Static Analysis [Urban18]

3-Step Recipe

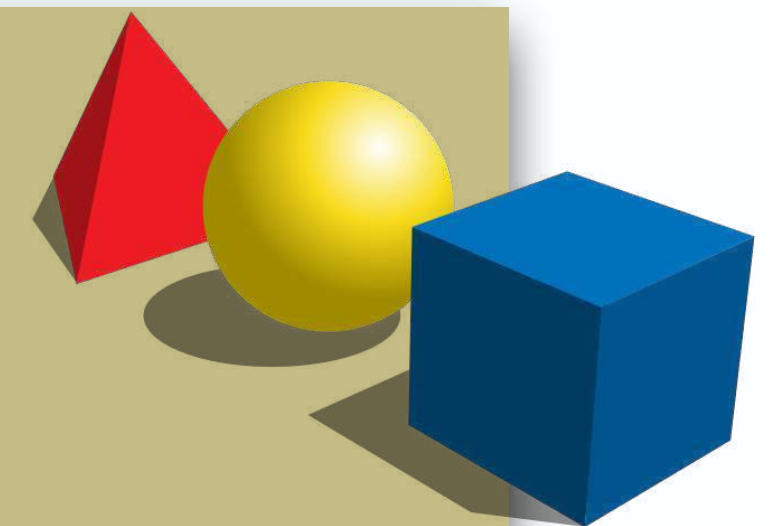
practical tools

targeting specific programs



abstract semantics, abstract domains

algorithmic approaches to decide program properties



concrete semantics

mathematical models of the program behavior



Dependency Fairness

$$\mathcal{F}_i \stackrel{\text{def}}{=} \{ \llbracket M \rrbracket \mid \text{UNUSED}_i(\llbracket M \rrbracket) \}$$

\mathcal{F}_i is the set of all neural networks M (or, rather, their semantics $\llbracket M \rrbracket$) that **do not use** the value of the sensitive input node $x_{0,i}$ for classification

$$\text{UNUSED}_i(T) \stackrel{\text{def}}{=} \forall t, t' \in T: t_0(x_{0,i}) \neq t'_0(x_{0,i}) \wedge \eta(t_0) = \eta(t'_0) \\ (\forall 0 \leq j \leq |L_0|: j \neq i \Rightarrow t_0(x_{0,j}) = t'_0(x_{0,j})) \\ \Rightarrow t_\omega = t'_\omega \quad \rho(t_0) = \rho(t'_0)$$

$\eta: \eta(x_{0j}) = \begin{cases} \top & j = i \\ x_{0j} & \text{otherwise} \end{cases}$

Intuitively: inputs differing only on the value of the sensitive input node $x_{0,i}$ should lead to the same **classification outcome**

Theorem

$$M \models \mathcal{F}_i \Leftrightarrow \{ \llbracket M \rrbracket \} \subseteq \mathcal{F}_i$$

Corollary

$$M \models \mathcal{F}_i \Leftarrow \llbracket M \rrbracket \subseteq \llbracket M \rrbracket^\# \subseteq \mathcal{F}_i$$

Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

\mathcal{N}_J is the set of all programs P (or, rather, their semantics $\llbracket P \rrbracket$) that **do not use** the value of the input variables in J

$$\begin{aligned} \text{UNUSED}_J(\llbracket P \rrbracket) \stackrel{\text{def}}{=} & \forall t \in \llbracket P \rrbracket, V \in \mathcal{R}^{|J|}: t_0(J) \neq V \Rightarrow \exists t' \in \llbracket P \rrbracket: \\ & (\forall i: i \notin J \Rightarrow t_0(i) = t'_0(i)) \\ & \wedge t'_0(J) = V \\ & \wedge t_\omega = t'_\omega \end{aligned}$$

Intuitively: **any possible program outcome** is possible from **any value** of the input variable i

Theorem

$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \} \subseteq \mathcal{N}_J$$

Data Usage Static Analysis [Urban18]

3-Step Recipe

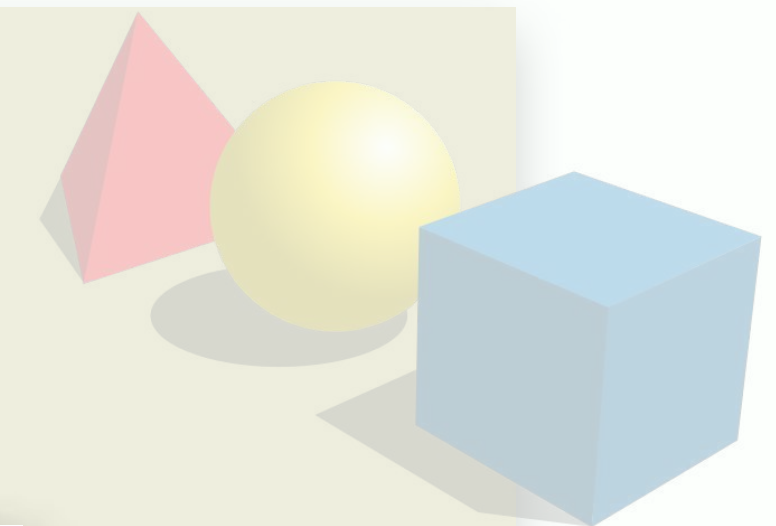
practical tools

targeting specific programs



abstract semantics, abstract domains

algorithmic approaches to decide program properties

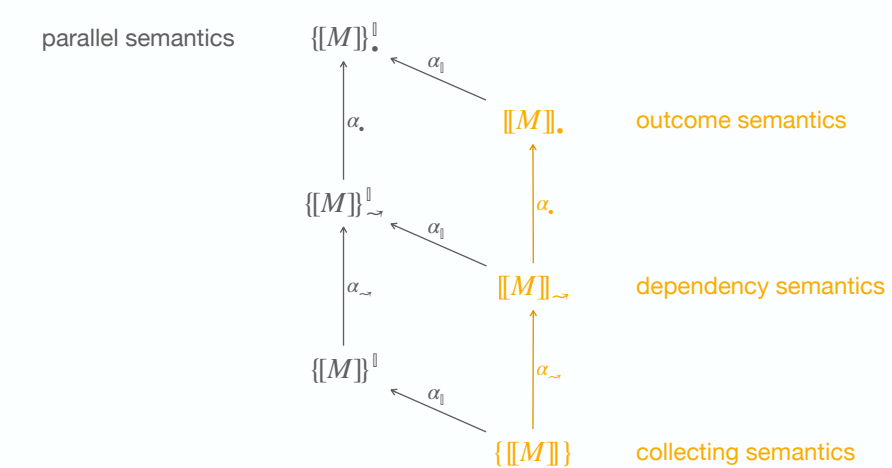


concrete semantics

mathematical models of the program behavior



Hierarchy of Semantics



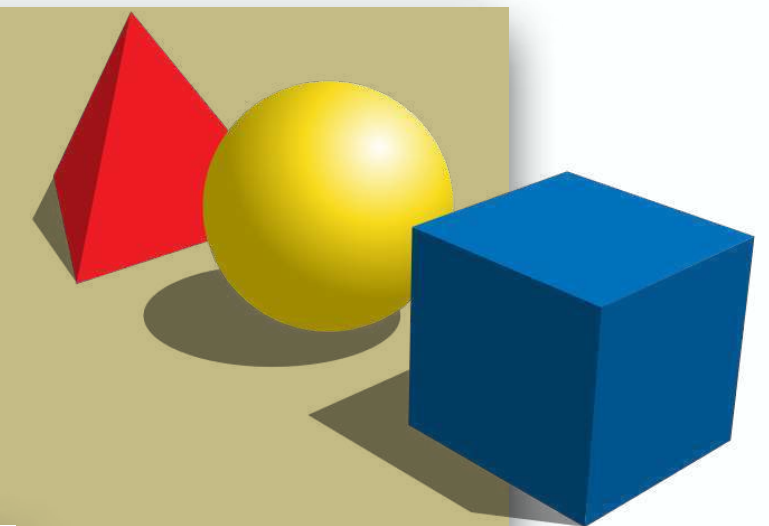
Data Usage Static Analysis [Urban18]

3-Step Recipe

practical tools
targeting specific programs

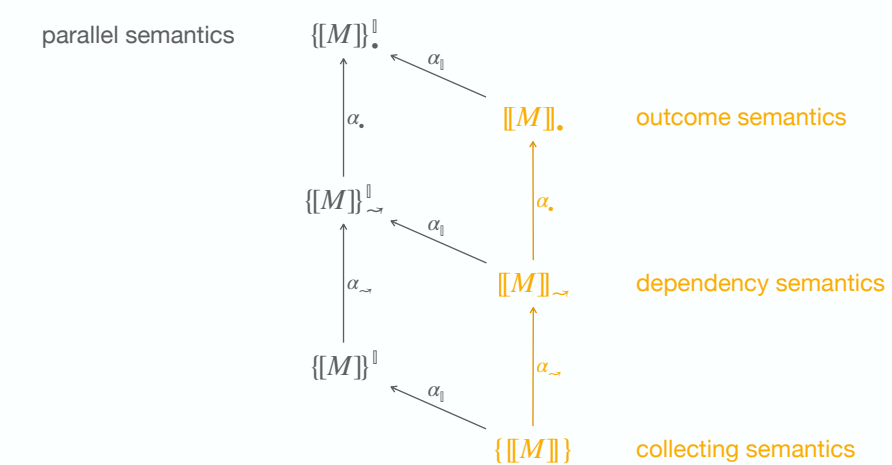


abstract semantics, abstract domains
algorithmic approaches to decide program properties



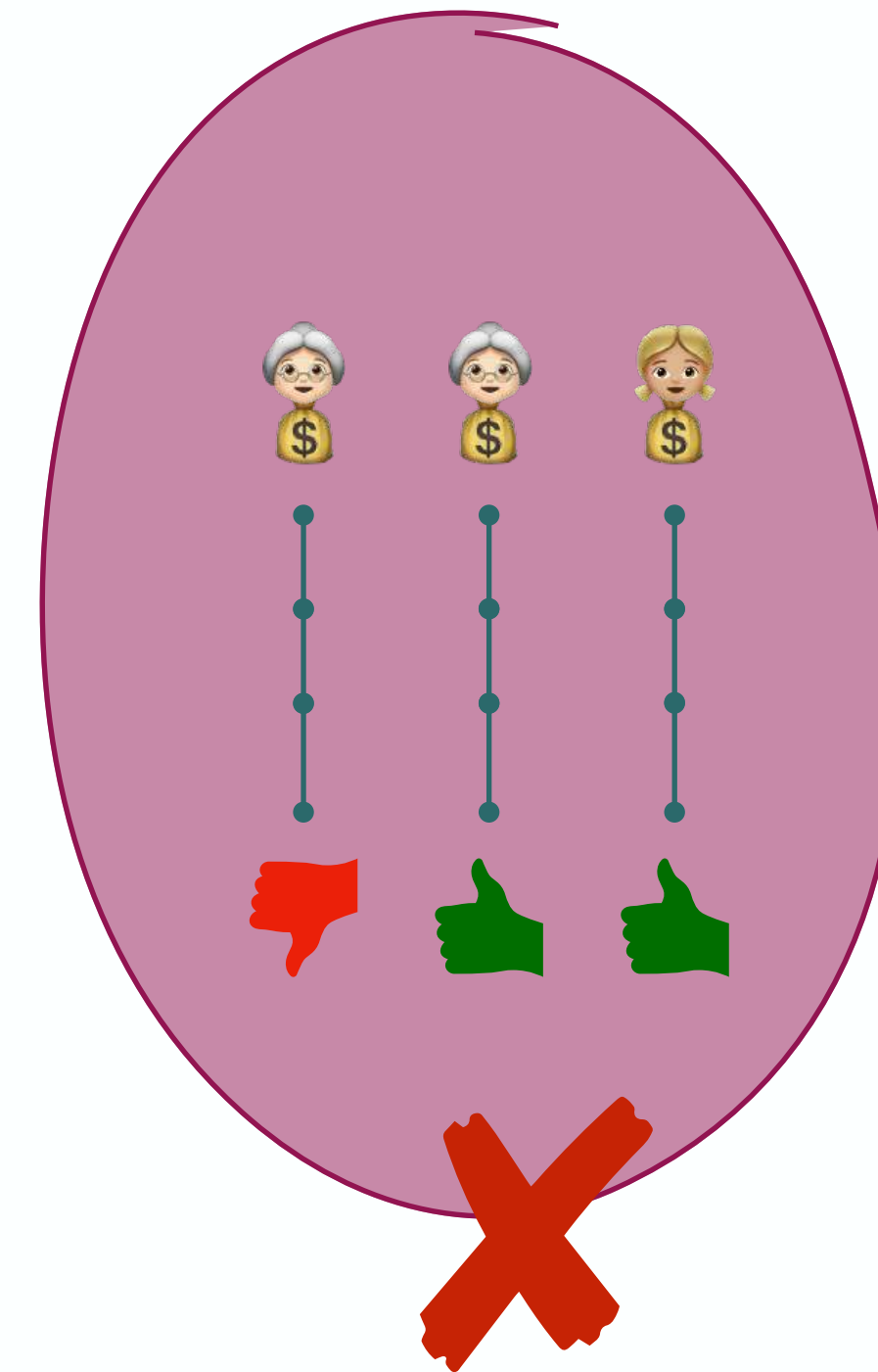
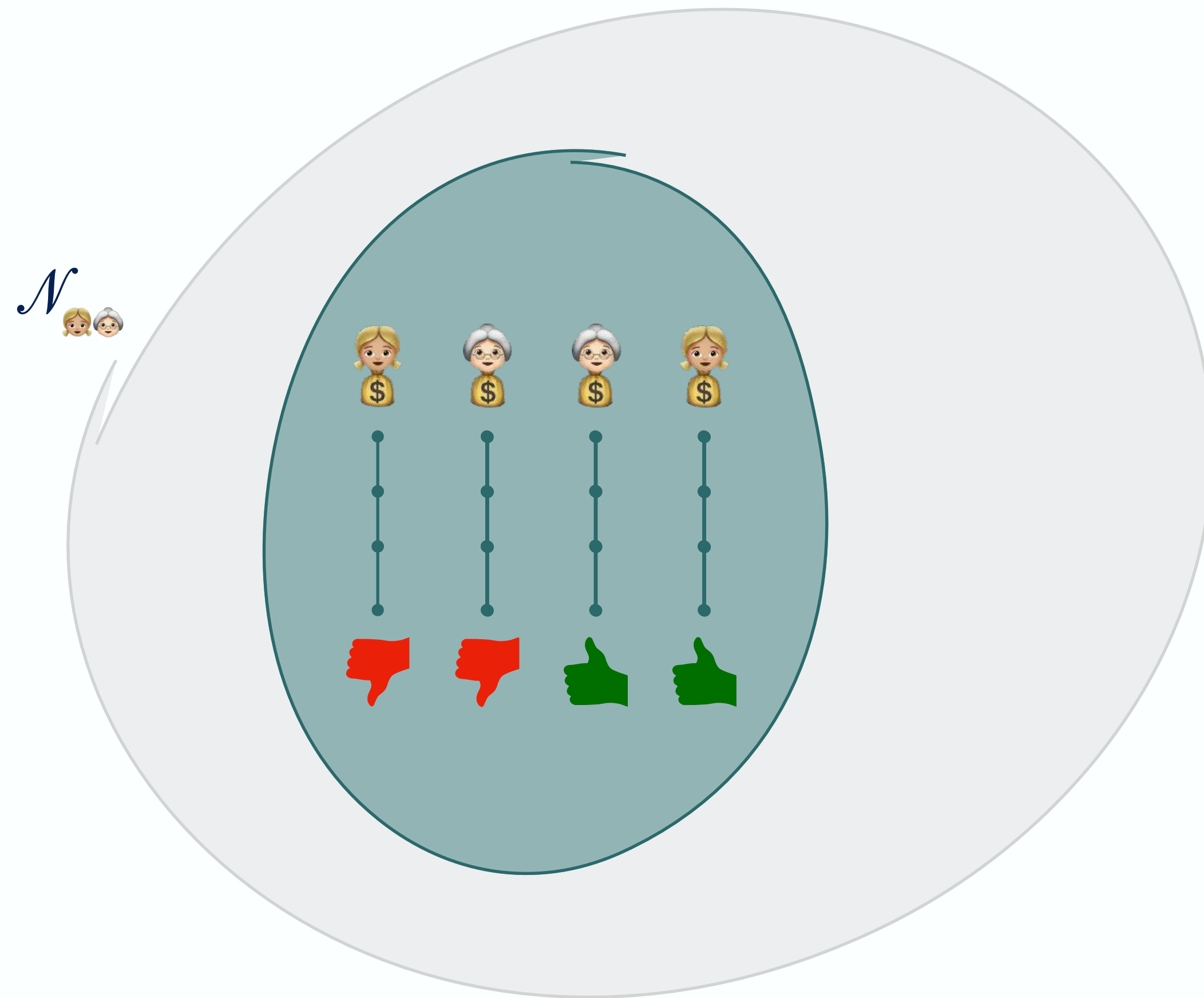
concrete semantics
mathematical models of the program behavior

Hierarchy of Semantics



Data (Non-) Usage

Not a Subset-Closed Property



Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

\mathcal{N}_J is the set of all programs P (or, rather, their semantics $\llbracket P \rrbracket$) that **do not use** the value of the input variables in J

$$\begin{aligned} \text{UNUSED}_J(\llbracket P \rrbracket) \stackrel{\text{def}}{=} & \forall t \in \llbracket P \rrbracket, V \in \mathcal{R}^{|J|}: t_0(J) \neq V \Rightarrow \exists t' \in \llbracket P \rrbracket: \\ & (\forall i: i \notin J \Rightarrow t_0(i) = t'_0(i)) \\ & \wedge t'_0(J) = V \\ & \wedge t_\omega = t'_\omega \end{aligned}$$

Intuitively: **any possible program outcome** is possible from **any value** of the input variable i

Theorem

$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \} \subseteq \mathcal{N}_J$$

Corollary

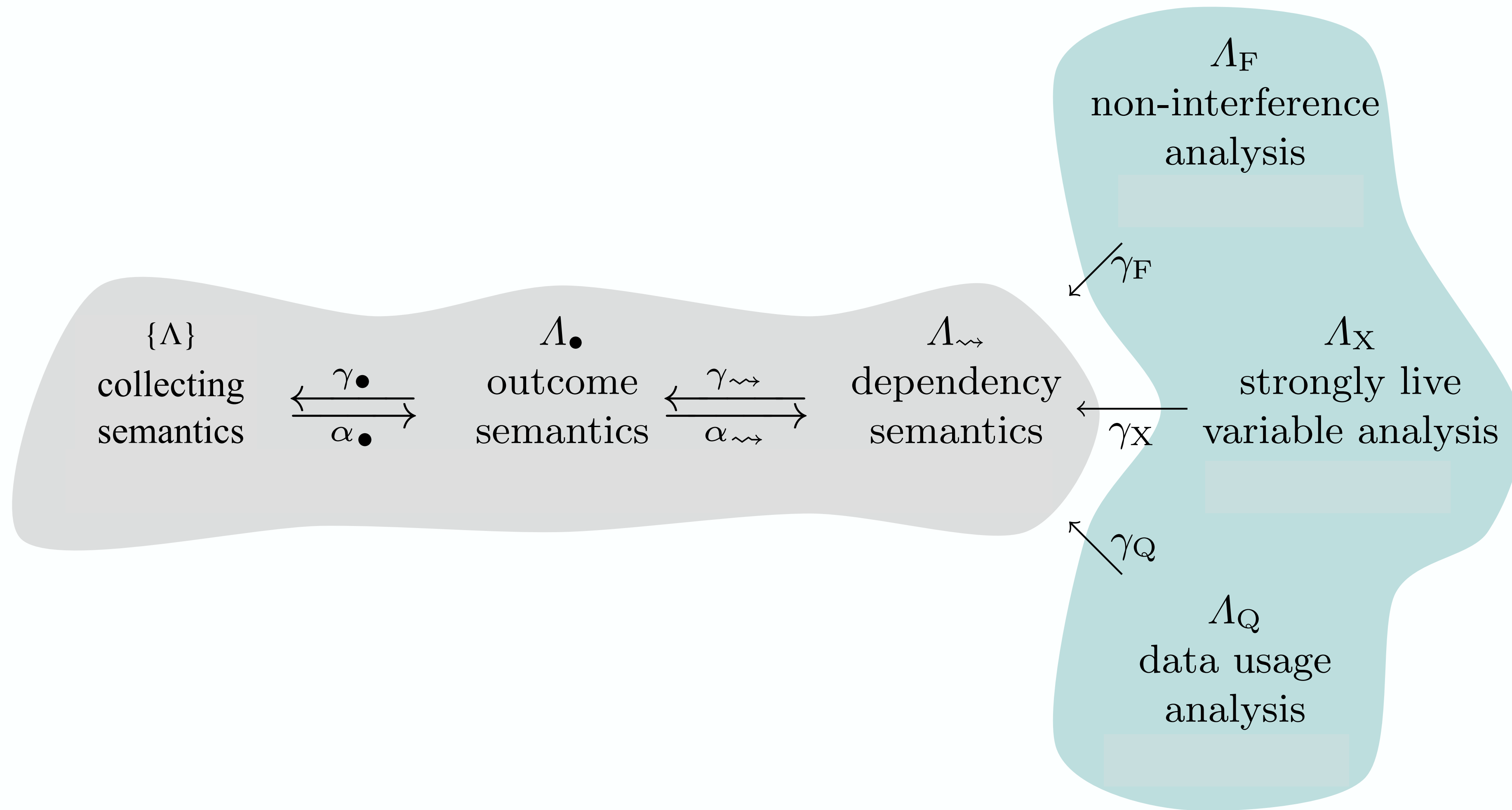
$$P \models \mathcal{N}_J \Leftarrow \{ \llbracket P \rrbracket \} \subseteq \llbracket P \rrbracket^\sharp \in \mathcal{N}_J$$

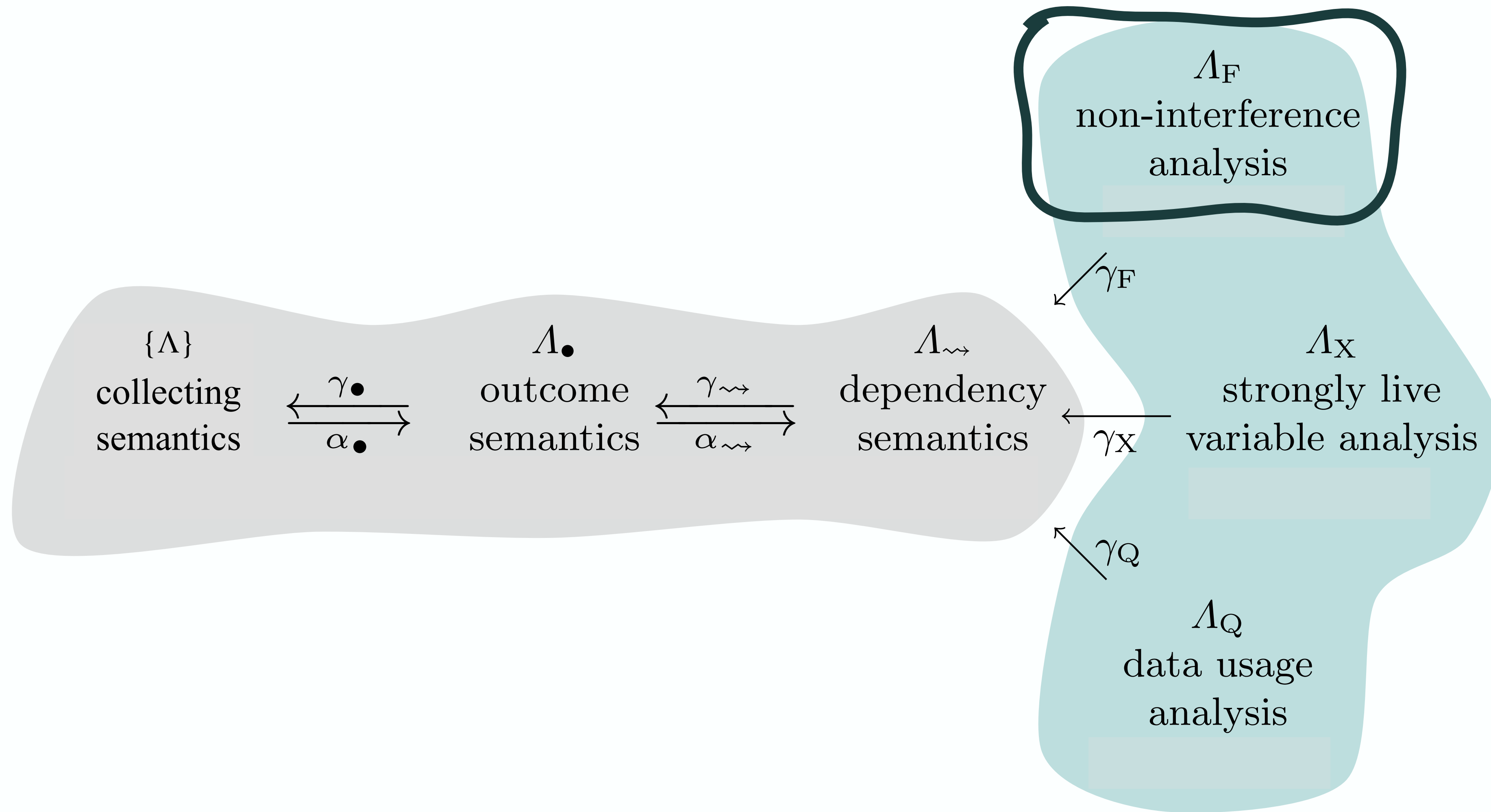
Data (Non-)Usage Abstractions

Over-Approximation of the Used Input Data

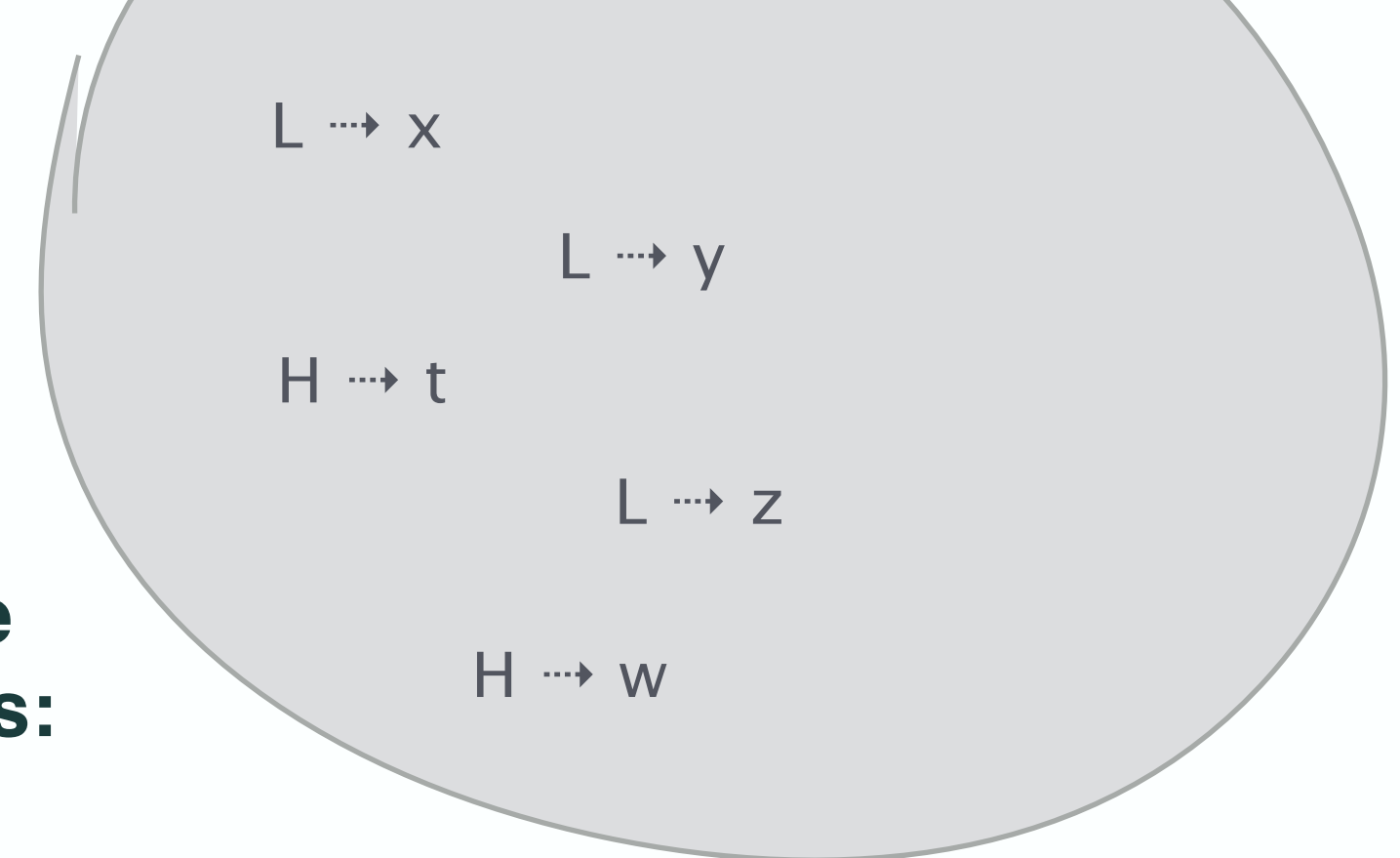
⇒ Under-Approximation of the Unused Input Data

$$P \models \mathcal{N}_{J \sqsubseteq J} \leftarrow \{\llbracket P \rrbracket\} \subseteq \llbracket P \rrbracket_A^{\sharp} \subseteq \mathcal{N}_{J \sqsubseteq J}$$





Secure Information Flow



possibilistic non-interference coincides with input data (non-)usage when the set J of unused input variables contains *all* input variables:

- **input variables** are **high-security** variables
- **output variables** are **low-security** variables

$\llbracket P \rrbracket_F$

explicit usage flows

$e ::= v \mid x \mid \text{not } e \mid e \text{ and } e \mid e \text{ or } e$ (expressions)

$s ::= \text{skip} \mid x = e \mid \text{if } e : s \text{ else: } s \mid \text{while } e : s \mid s \ s$ (statements)

implicit usage flows

$\Theta_F[\text{skip}](S) \stackrel{\text{def}}{=} S$

$\Theta_F[x = e](S) \stackrel{\text{def}}{=} \{L \rightsquigarrow y \in S \mid y \neq x\} \cup \{L \rightsquigarrow x \mid \mathcal{V}_F[e]S\}$

$\Theta_F[\text{if } e : s_1 \text{ else: } s_2](S) \stackrel{\text{def}}{=} \begin{cases} \Theta_F[s_1](S) \sqcup_F \Theta_F[s_2](S) & \text{if } \mathcal{V}_F[e]S \\ \{L \rightsquigarrow x \in S \mid x \notin W(s_1) \cup W(s_2)\} & \text{otherwise} \end{cases}$

$\Theta_F[\text{while } e : s](S) \stackrel{\text{def}}{=} \text{lfp}_{S^F} \Theta_F[\text{if } e : s \text{ else: } \text{skip}]$

$\Theta_F[s_1 \ s_2](S) \stackrel{\text{def}}{=} \Theta_F[s_2] \circ \Theta_F[s_1](S)$

S guarantees a unique value for e independently of values of input variables

$\mathcal{V}_F[x]S \iff L \rightsquigarrow x \in S$

set of variables modified by s_i

$\mathcal{L} \stackrel{\text{def}}{=} \{L, H\}$: set of security levels

$L \rightsquigarrow x$: dependency constraint

$F \stackrel{\text{def}}{=} \{L \rightsquigarrow x \mid x \in X\}$

$\langle \mathcal{P}(F), \sqsubseteq_F, \sqcup_F \rangle$: abstract domain

$S_1 \sqsubseteq_F S_2 \stackrel{\text{def}}{=} S_1 \supseteq S_2$

$S_1 \sqcup_F S_2 \stackrel{\text{def}}{=} S_1 \cap S_2$

Hypercollecting Semantics and Its Application to Static Analysis of Information Flow

Mounir Assaf
Stevens Institute of Technology, Hoboken, US
first.last@stevens.edu

David A. Naumann
Stevens Institute of Technology, Hoboken, US
first.last@stevens.edu

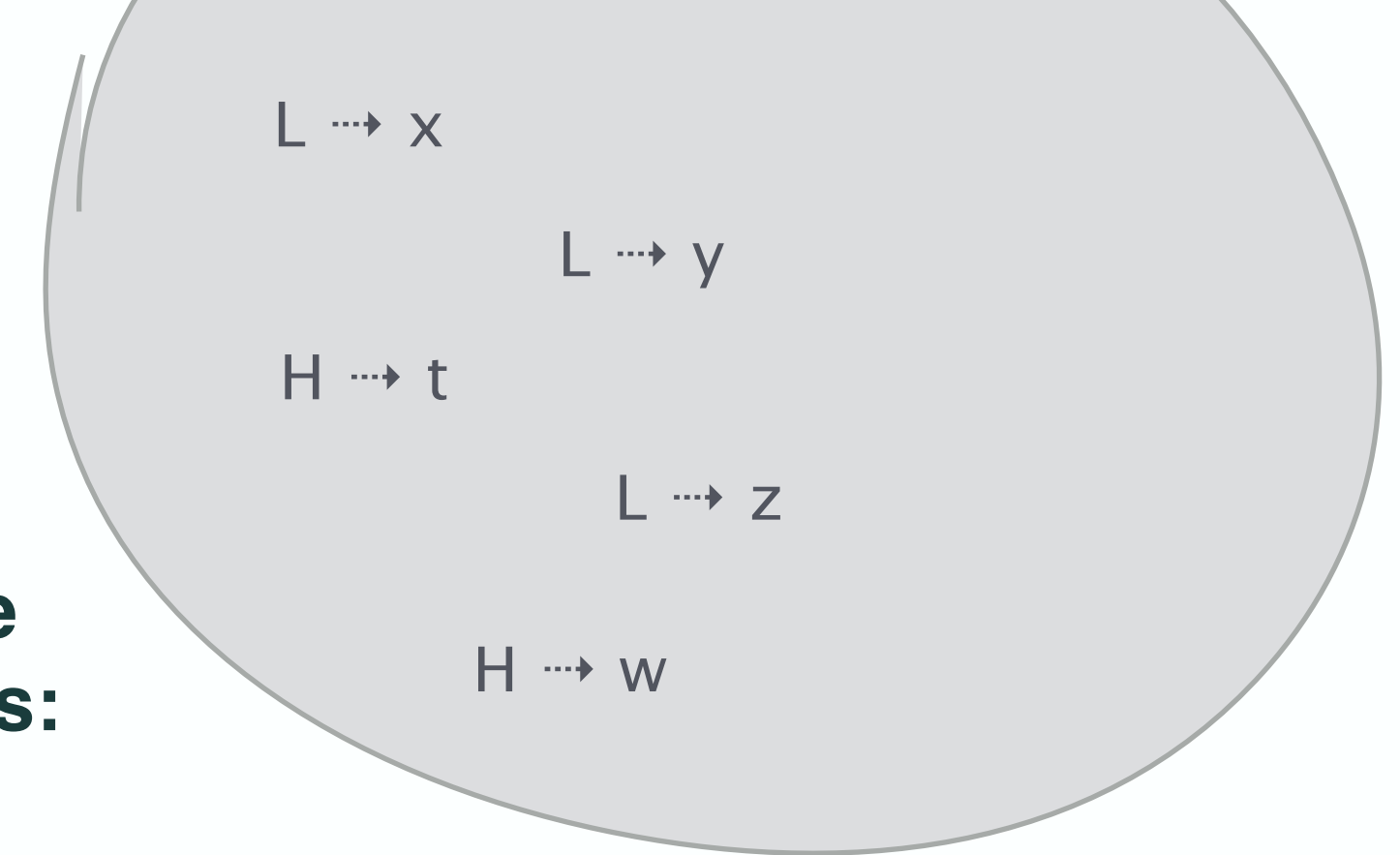
Julien Signoles
Software Reliability and Security Lab, CEA LIST, Saclay, FR
first.last@cea.fr

Éric Totel
CIDRE, CentraleSupélec, Rennes, FR
first.last@centralesupelec.fr

Frédéric Tronel
CIDRE, CentraleSupélec, Rennes, FR
first.last@centralesupelec.fr

program is correct if all its traces satisfy the predicate. By contrast with such *trace properties*, extensional definitions of dependences more than one trace. To express that the final value of a variable y , the same

Secure Information Flow



$\llbracket P \rrbracket_F$

possibilistic non-interference coincides with input data (non-)usage when the set J of unused input variables contains *all* input variables:

- **input variables** are **high-security** variables
- **output variables** are **low-security** variables

$e ::= v \mid x \mid \text{not } e \mid e \text{ and } e \mid e \text{ or } e$ (expressions)
 $s ::= \text{skip} \mid x = e \mid \text{if } e: s \text{ else: } s \mid \text{while } e: s \mid s \ s$ (statements)

$$\Theta_F[\text{skip}](S) \stackrel{\text{def}}{=} S$$

$$\Theta_F[x = e](S) \stackrel{\text{def}}{=} \{L \rightsquigarrow y \in S \mid y \neq x\} \cup \{L \rightsquigarrow x \mid \mathcal{V}_F[e]S\}$$

$$\Theta_F[\text{if } e: s_1 \text{ else: } s_2](S) \stackrel{\text{def}}{=} \begin{cases} \Theta_F[s_1](S) \sqcup_F \Theta_F[s_2](S) & \text{if } \mathcal{V}_F[e]S \\ \{L \rightsquigarrow x \in S \mid x \notin W(s_1) \cup W(s_2)\} & \text{otherwise} \end{cases}$$

$$\Theta_F[\text{while } e: s](S) \stackrel{\text{def}}{=} \text{lfp}_{S^F} \Theta_F[\text{if } e: s \text{ else: } \text{skip}]$$

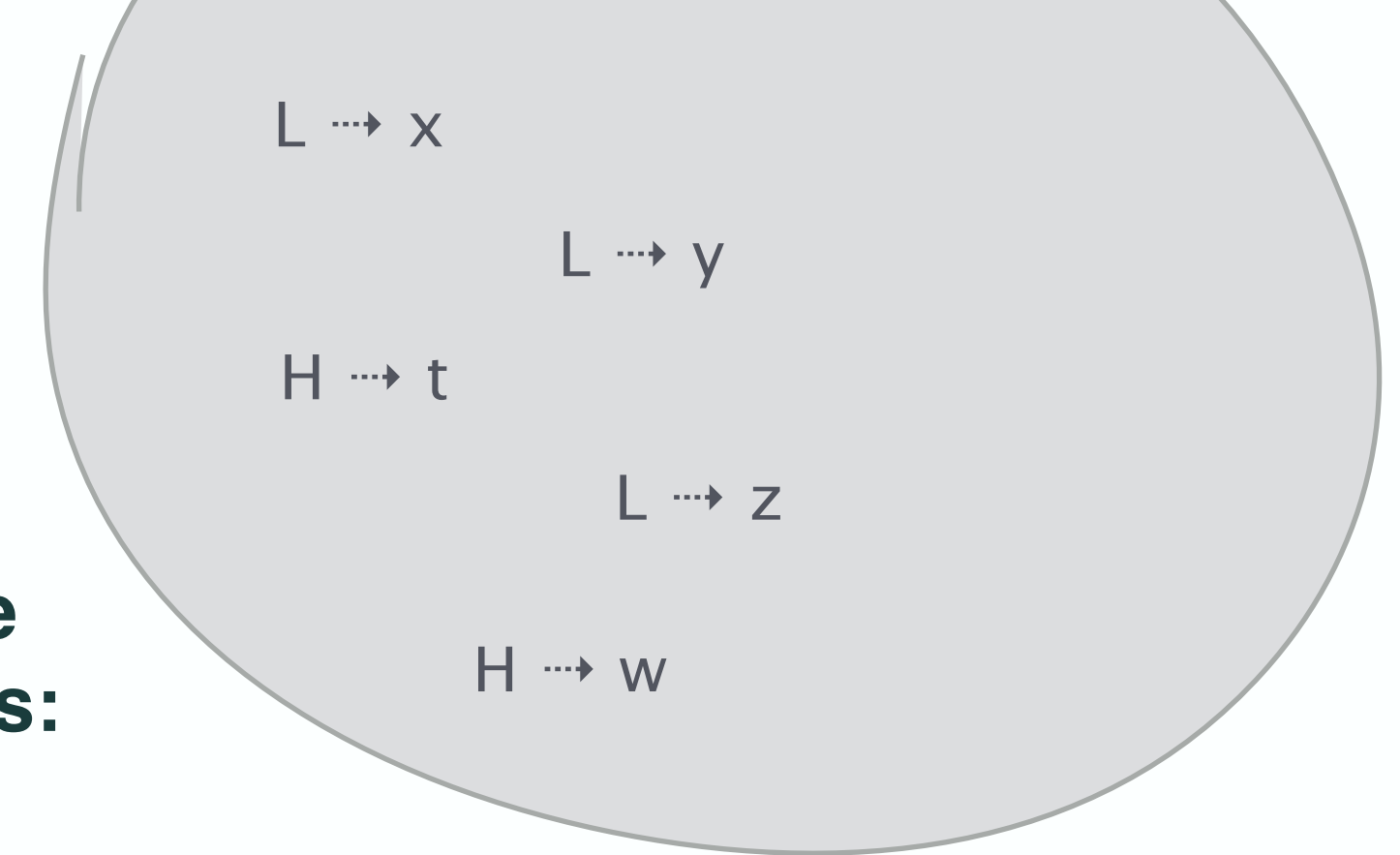
$$\Theta_F[s_1 \ s_2](S) \stackrel{\text{def}}{=} \Theta_F[s_2] \circ \Theta_F[s_1](S)$$

```

passing = True
if not english:
    english = False
if not math:
    passing = False or bonus
if not math:
    passing = False or bonus
    
```

←..... L ↦ passing, H ↦ english, math, science, bonus
 ←..... L ↦ passing, H ↦ english, math, science, bonus
 ←..... L ↦ passing, H ↦ english, math, science, bonus
 ←..... H ↦ english, math, science, bonus, passing
 ←..... H ↦ english, math, science, bonus, passing

Secure Information Flow



possibilistic non-interference coincides with input data (non-)usage when the set J of unused input variables contains *all* input variables:

- **input variables** are **high-security** variables
 - **output variables** are **low-security** variables
- and the program is terminating

$\llbracket P \rrbracket_F$

$e ::= v \mid x \mid \text{not } e \mid e \text{ and } e \mid e \text{ or } e$ (expressions)
 $s ::= \text{skip} \mid x = e \mid \text{if } e : s \text{ else: } s \mid \text{while } e : s \mid s \ s$ (statements)

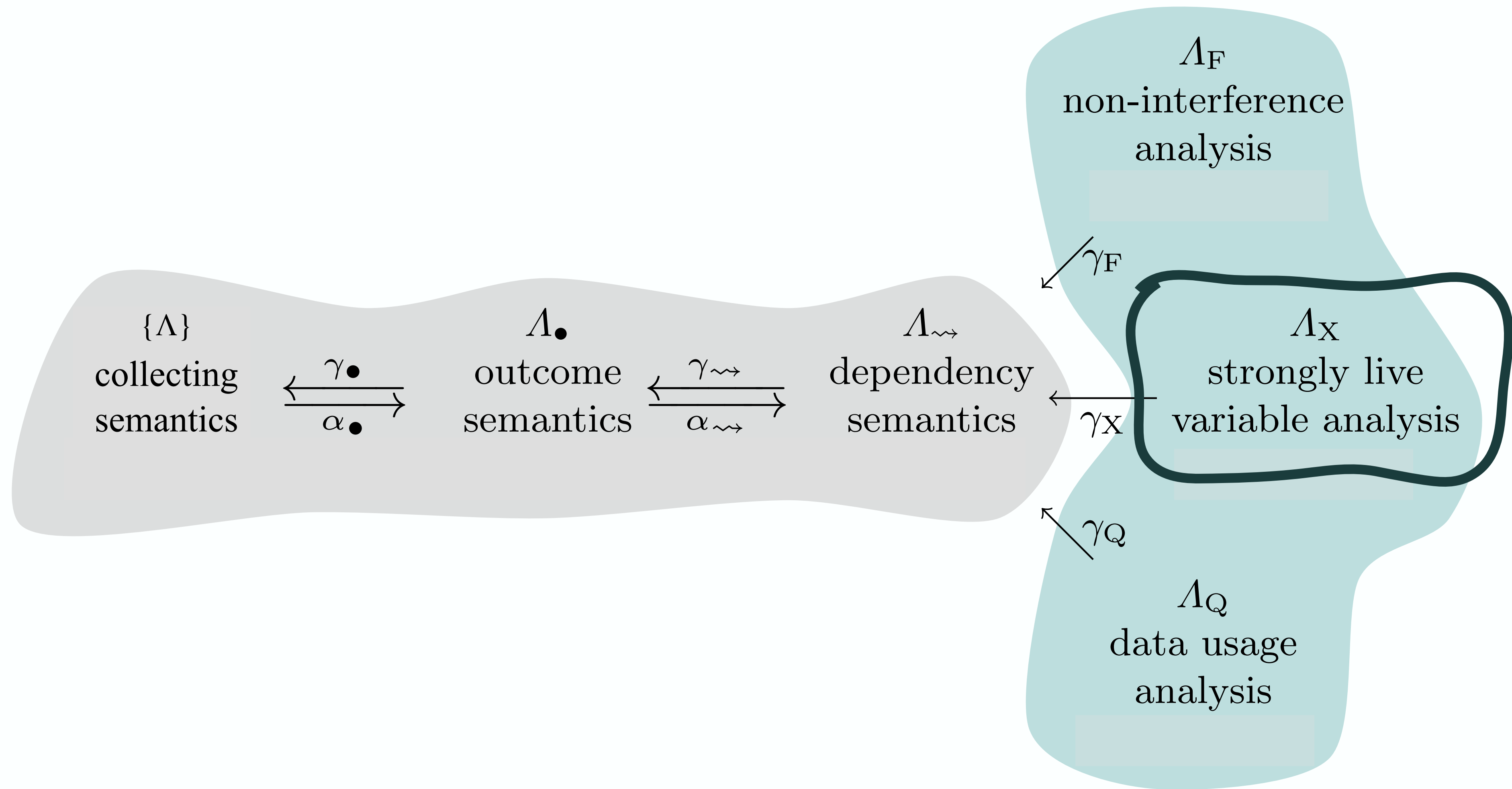
$$\begin{aligned} \Theta_F[\text{skip}](S) &\stackrel{\text{def}}{=} S \\ \Theta_F[x = e](S) &\stackrel{\text{def}}{=} \{L \rightsquigarrow y \in S \mid y \neq x\} \cup \{L \rightsquigarrow x \mid \mathcal{V}_F[e]S\} \\ \Theta_F[\text{if } e : s_1 \text{ else: } s_2](S) &\stackrel{\text{def}}{=} \begin{cases} \Theta_F[s_1](S) \sqcup_F \Theta_F[s_2](S) & \text{if } \mathcal{V}_F[e]S \\ \{L \rightsquigarrow x \in S \mid x \notin W(s_1) \cup W(s_2)\} & \text{otherwise} \end{cases} \\ \Theta_F[\text{while } e : s](S) &\stackrel{\text{def}}{=} \text{lfp}_{S^F} \Theta_F[\text{if } e : s \text{ else: } \text{skip}] \\ \Theta_F[s_1 \ s_2](S) &\stackrel{\text{def}}{=} \Theta_F[s_2] \circ \Theta_F[s_1](S) \end{aligned}$$

passing = **True**
while not english:
 english = **False**

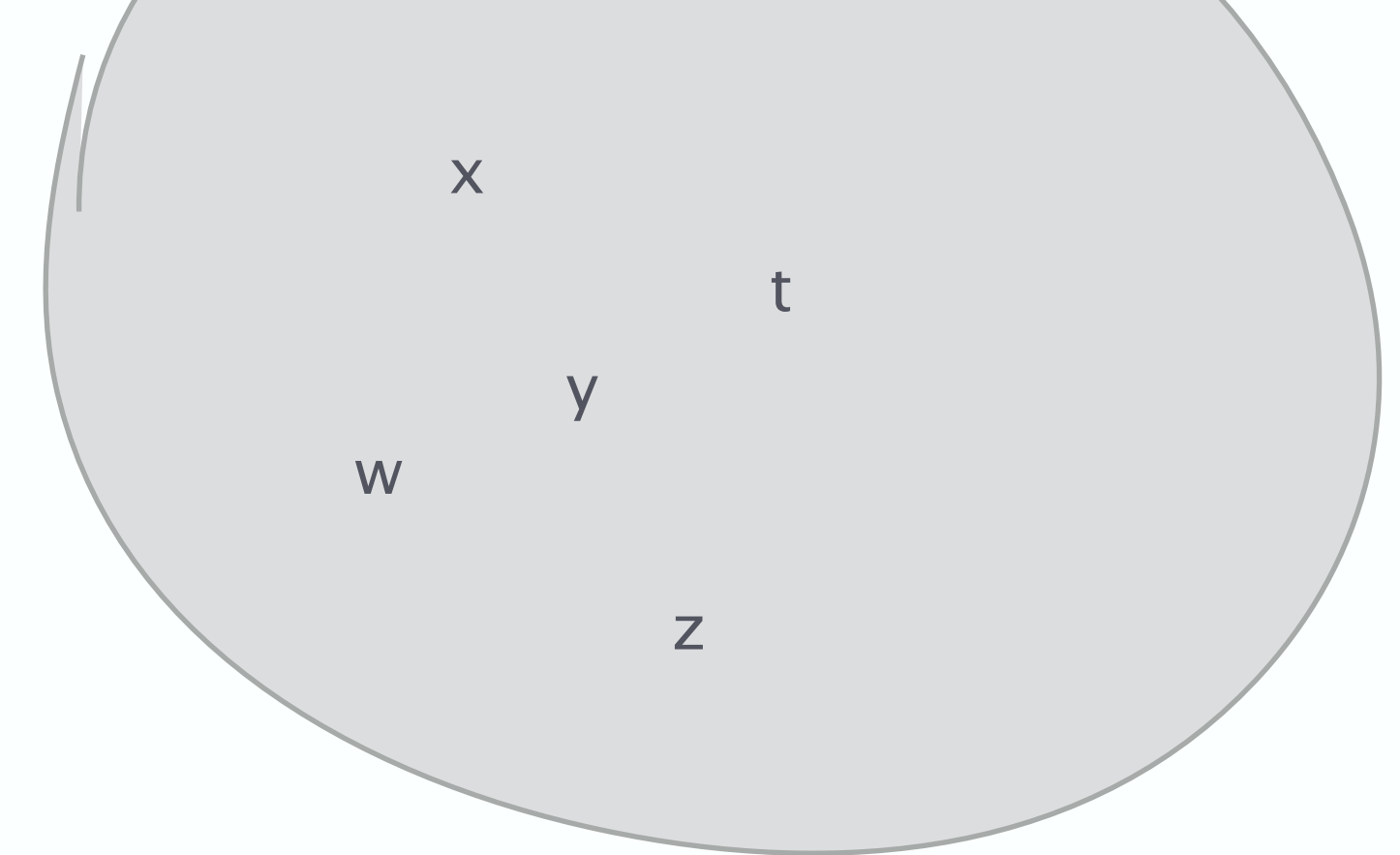
←..... L ↦ passing, H ↦ english, math, science, bonus
 ←..... L ↦ passing, H ↦ english, math, science, bonus
 ←..... L ↦ passing, H ↦ english, math, science, bonus

Theorem

$$P \models \mathcal{N}_J^+ \Leftarrow \{\llbracket P \rrbracket\} \subseteq \llbracket P \rrbracket_F^{\sharp} \subseteq \mathcal{N}_J^+$$



Strong Liveness



$\llbracket P \rrbracket_X$

a variable is **strongly live** if

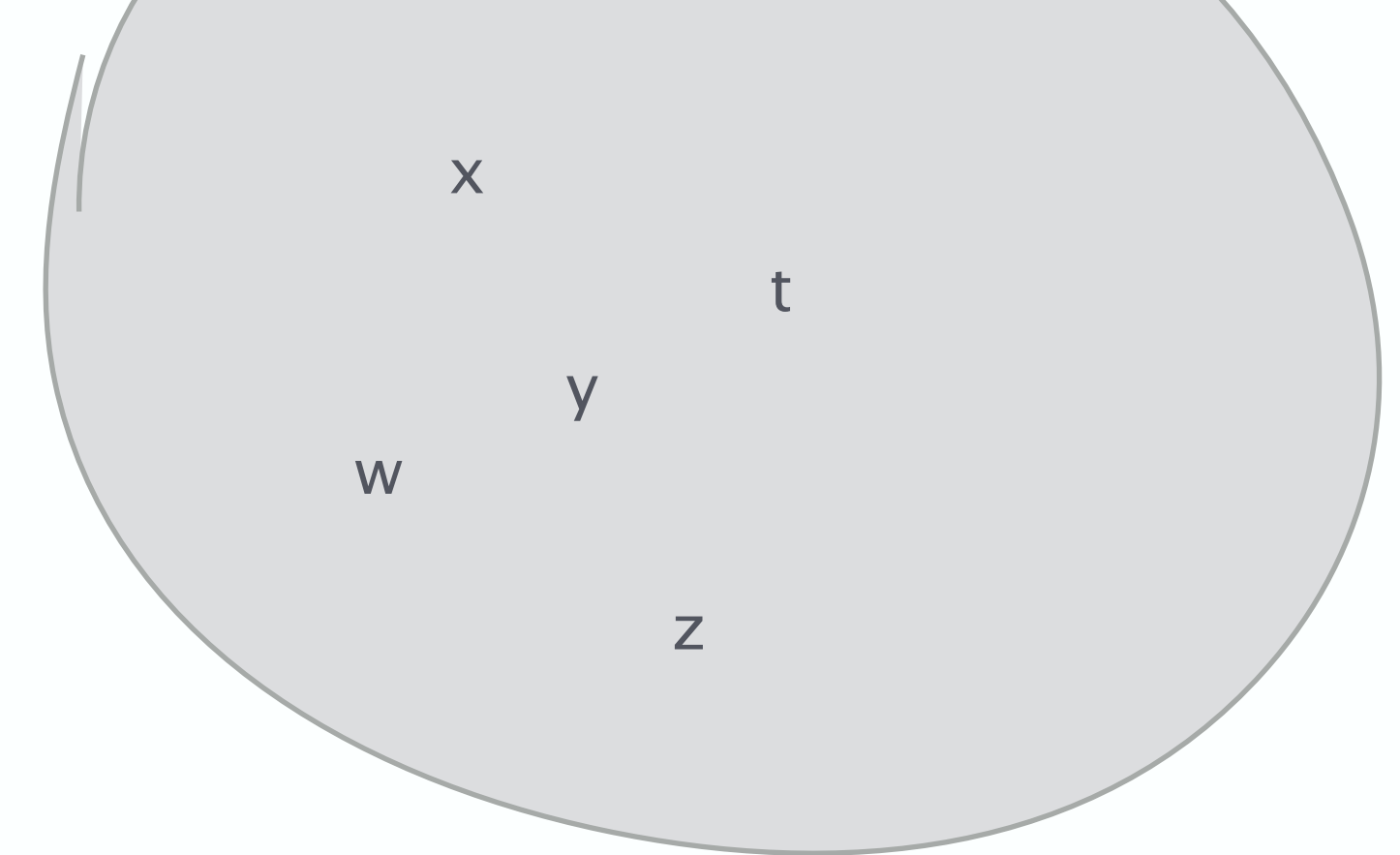
- it is used in an assignment to another strongly live variable
- it is used in a statement other than an assignment

$e ::= v \mid x \mid \text{not } e \mid e \text{ and } e \mid e \text{ or } e$ (expressions)
 $s ::= \text{skip} \mid x = e \mid \text{if } e : s \text{ else: } s \mid \text{while } e : s \mid s \ s$ (statements)

$\Theta_X[\text{skip}](S) \stackrel{\text{def}}{=} S$
 $\Theta_X[x = e](S) \stackrel{\text{def}}{=} \begin{cases} (S \setminus \{x\}) \cup \text{VARS}(e) & x \in S \\ S & \text{otherwise} \end{cases}$
 $\Theta_X[\text{if } b : s_1 \text{ else: } s_2](S) \stackrel{\text{def}}{=} \text{VARS}(b) \cup \Theta_X[s_1](S) \cup \Theta_X[s_2](S)$
 $\Theta_X[\text{while } b : s](S) \stackrel{\text{def}}{=} \text{VARS}(b) \cup \Theta_X[s](S)$
 $\Theta_X[s_1 \ s_2](S) \stackrel{\text{def}}{=} \Theta_X[s_1] \circ \Theta_X[s_2](S)$

$\langle \mathcal{P}(X), \subseteq, \cup, \cap, \emptyset, X \rangle$: abstract domain

Strong Liveness



$\llbracket P \rrbracket_X$

a variable is **strongly live** if

- it is used in an assignment to another strongly live variable
- it is used in a statement other than an assignment

$e ::= v \mid x \mid \text{not } e \mid e \text{ and } e \mid e \text{ or } e$ (expressions)
 $s ::= \text{skip} \mid x = e \mid \text{if } e: s \text{ else: } s \mid \text{while } e: s \mid s \ s$ (statements)

$$\Theta_X[\text{skip}](S) \stackrel{\text{def}}{=} S$$

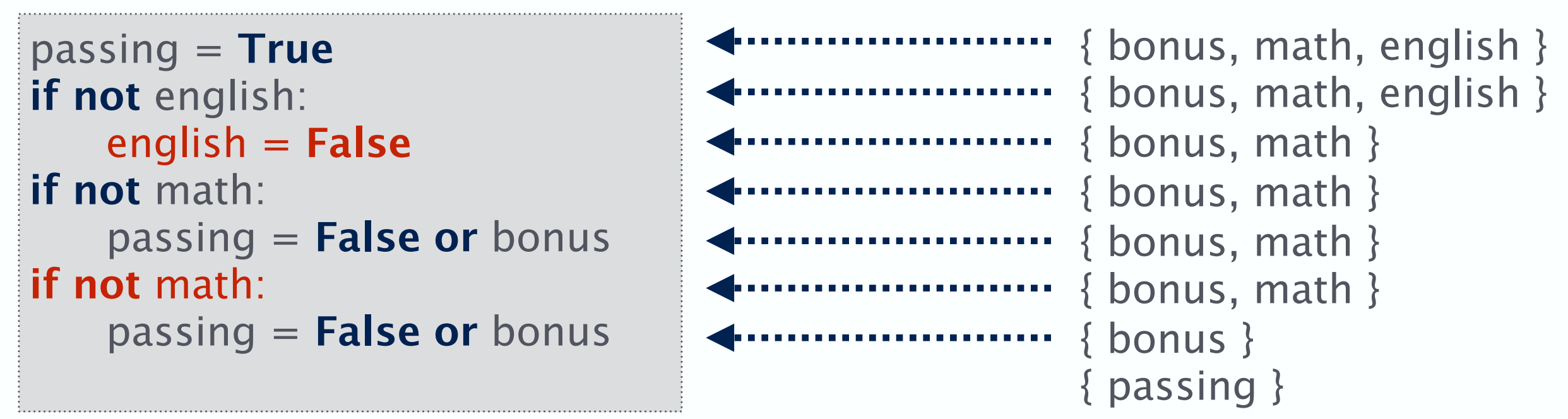
$$\Theta_X[x = e](S) \stackrel{\text{def}}{=} \begin{cases} (S \setminus \{x\}) \cup \text{VARS}(e) & x \in S \\ S & \text{otherwise} \end{cases}$$

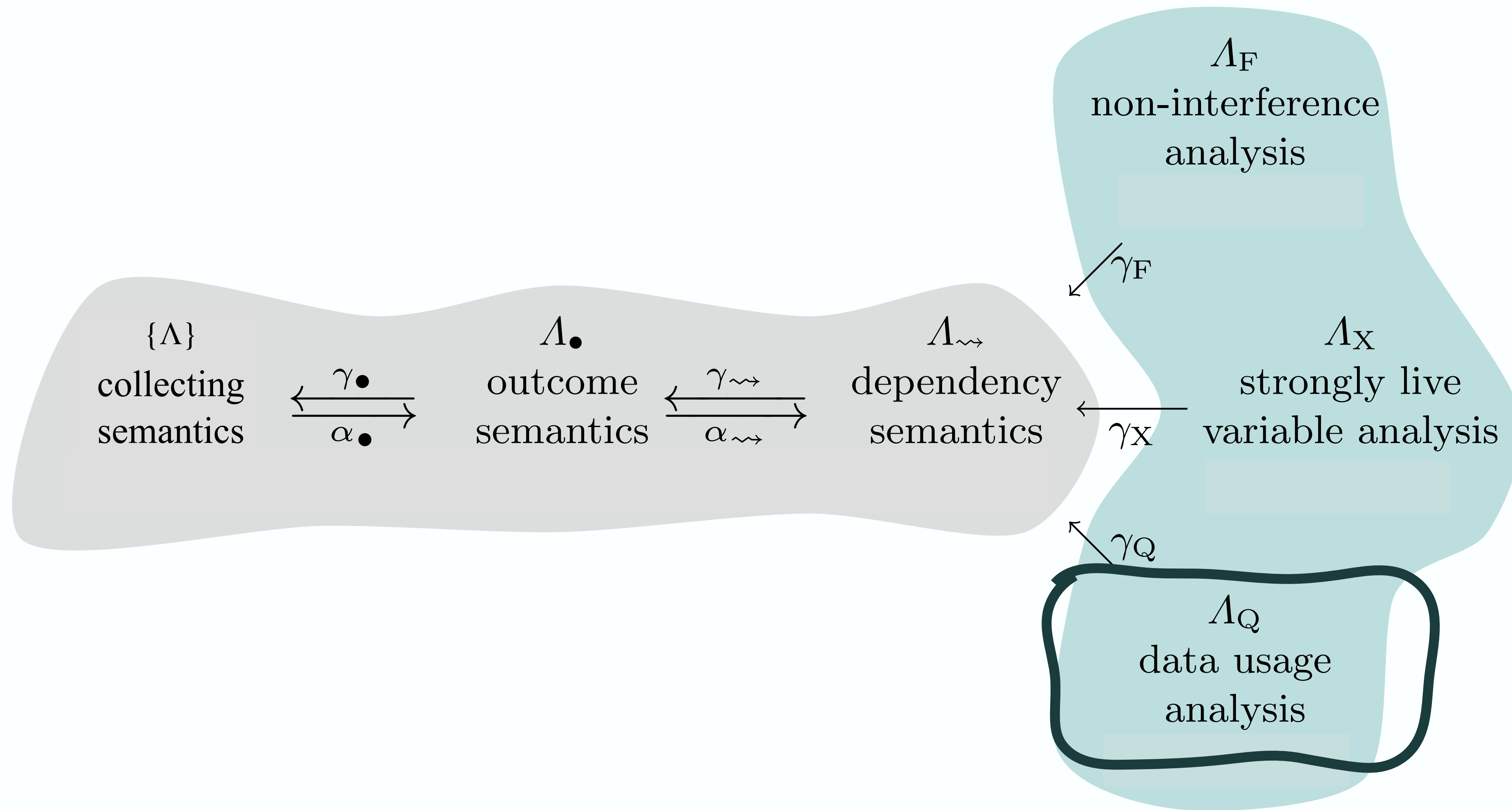
$$\Theta_X[\text{if } b: s_1 \text{ else: } s_2](S) \stackrel{\text{def}}{=} \text{VARS}(b) \cup \Theta_X[s_1](S) \cup \Theta_X[s_2](S)$$

$$\Theta_X[\text{while } b: s](S) \stackrel{\text{def}}{=} \text{VARS}(b) \cup \Theta_X[s](S)$$

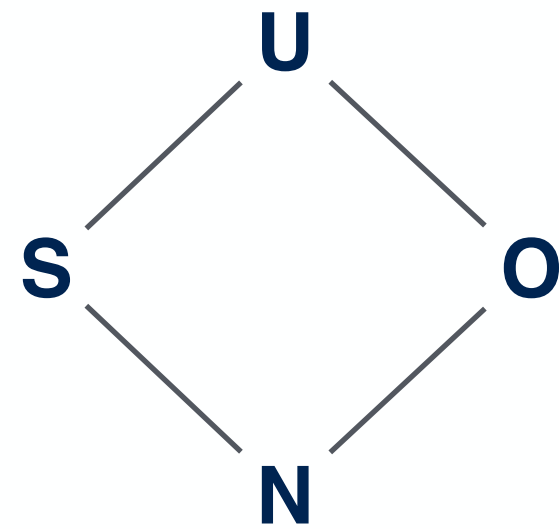
$$\Theta_X[s_1 \ s_2](S) \stackrel{\text{def}}{=} \Theta_X[s_1] \circ \Theta_X[s_2](S)$$

Theorem
 $P \models \mathcal{N}_J \Leftarrow \{ \llbracket P \rrbracket \} \subseteq \llbracket P \rrbracket_X \subseteq \mathcal{N}_J$

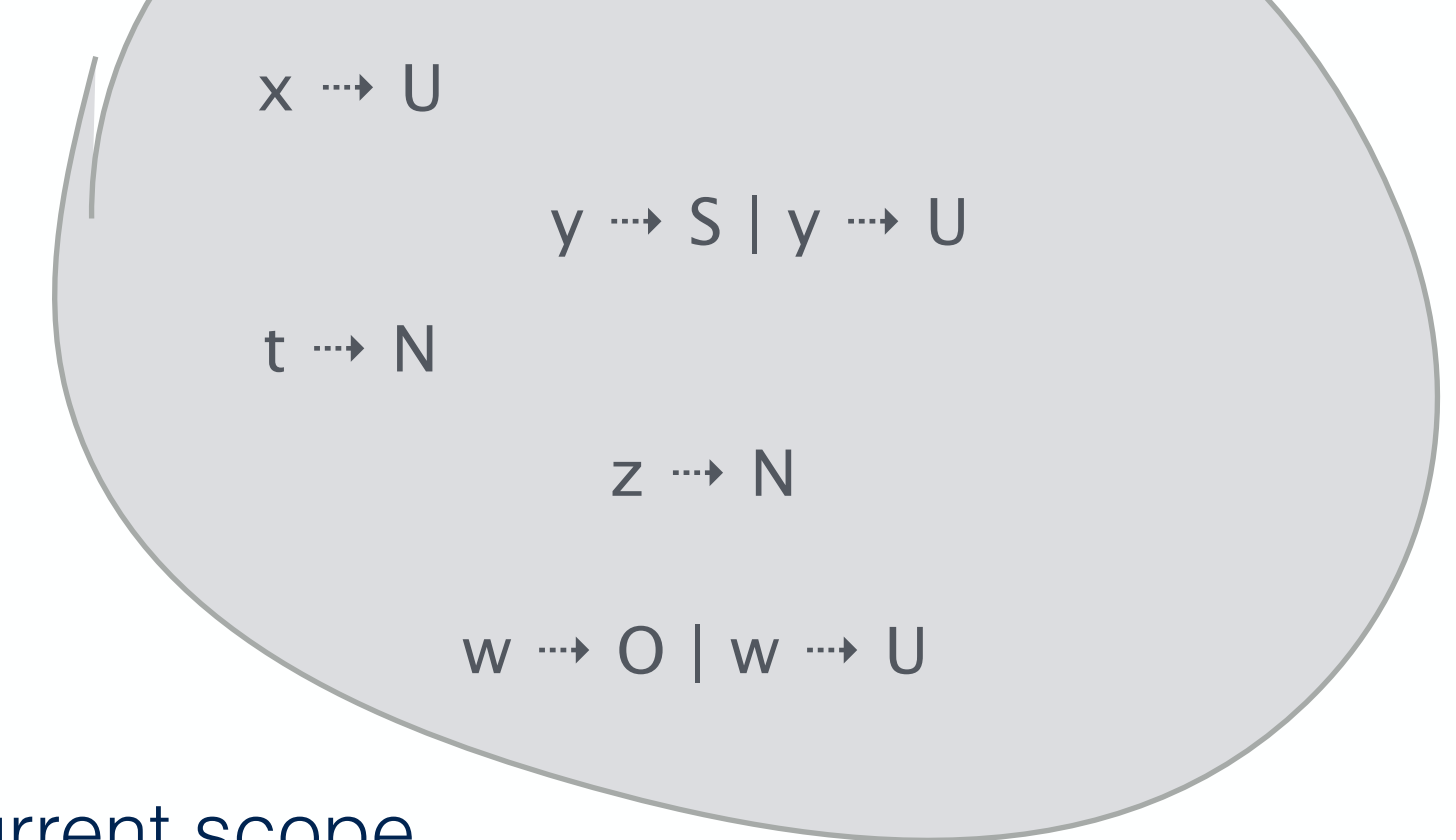




Syntactic (Non-)Usage



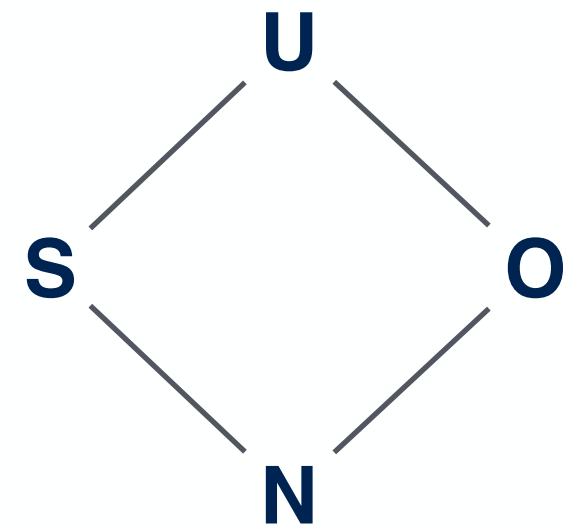
- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



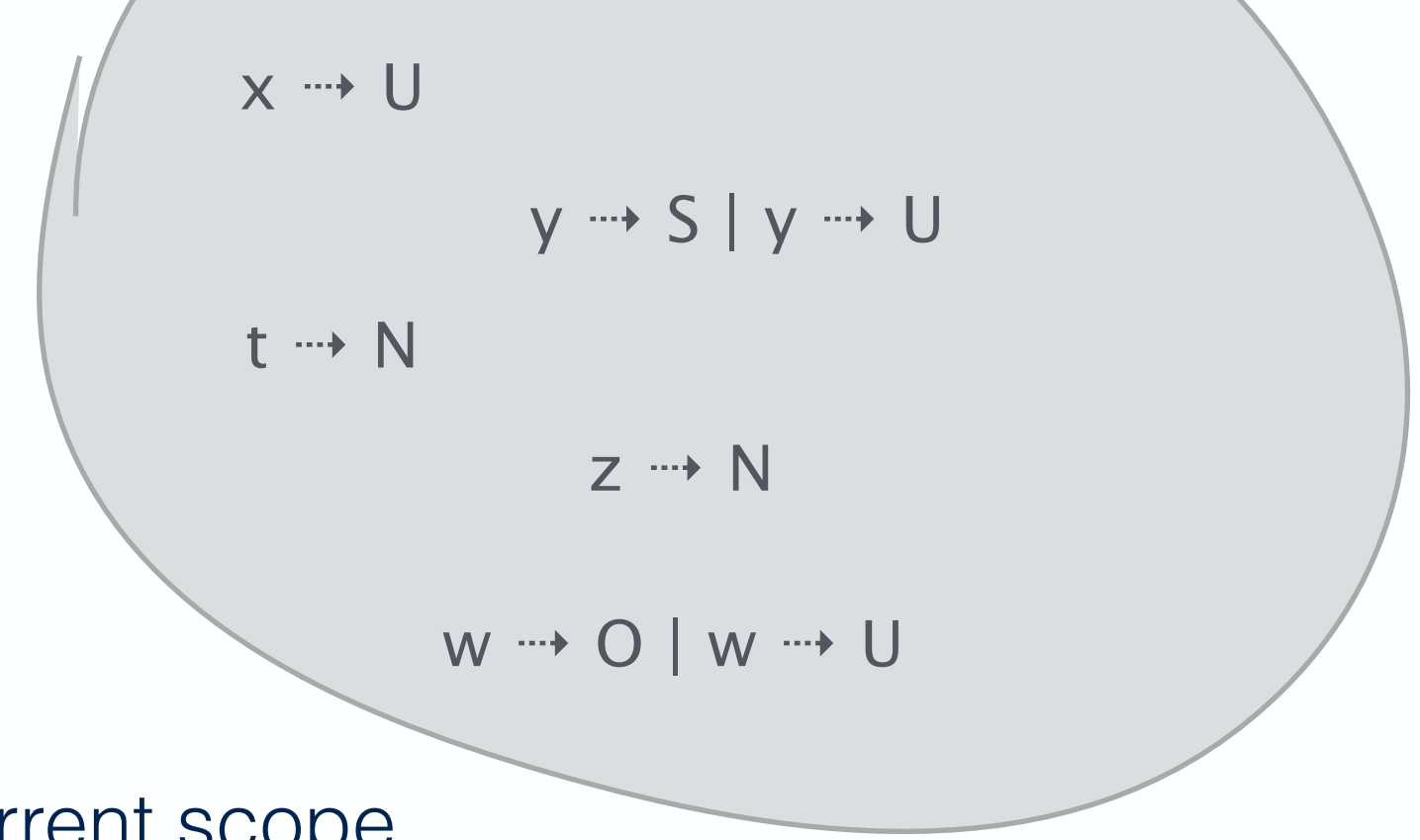
$\llbracket P \rrbracket_Q$

$$\begin{aligned} \Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\ \Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}] \\ \Theta_Q[s_1 \ s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q) \end{aligned}$$

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

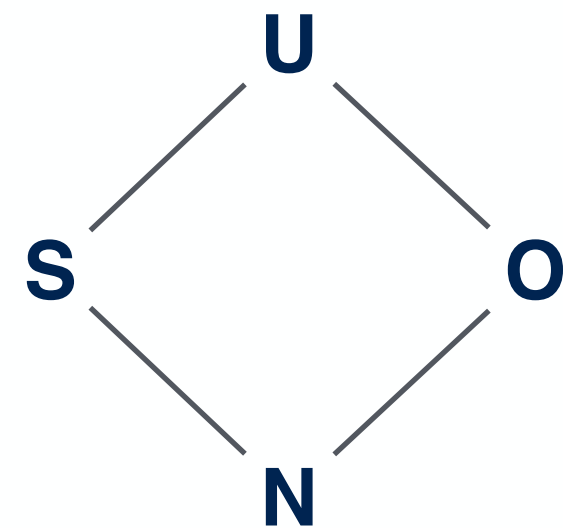
if not math:

passing = **False or** bonus

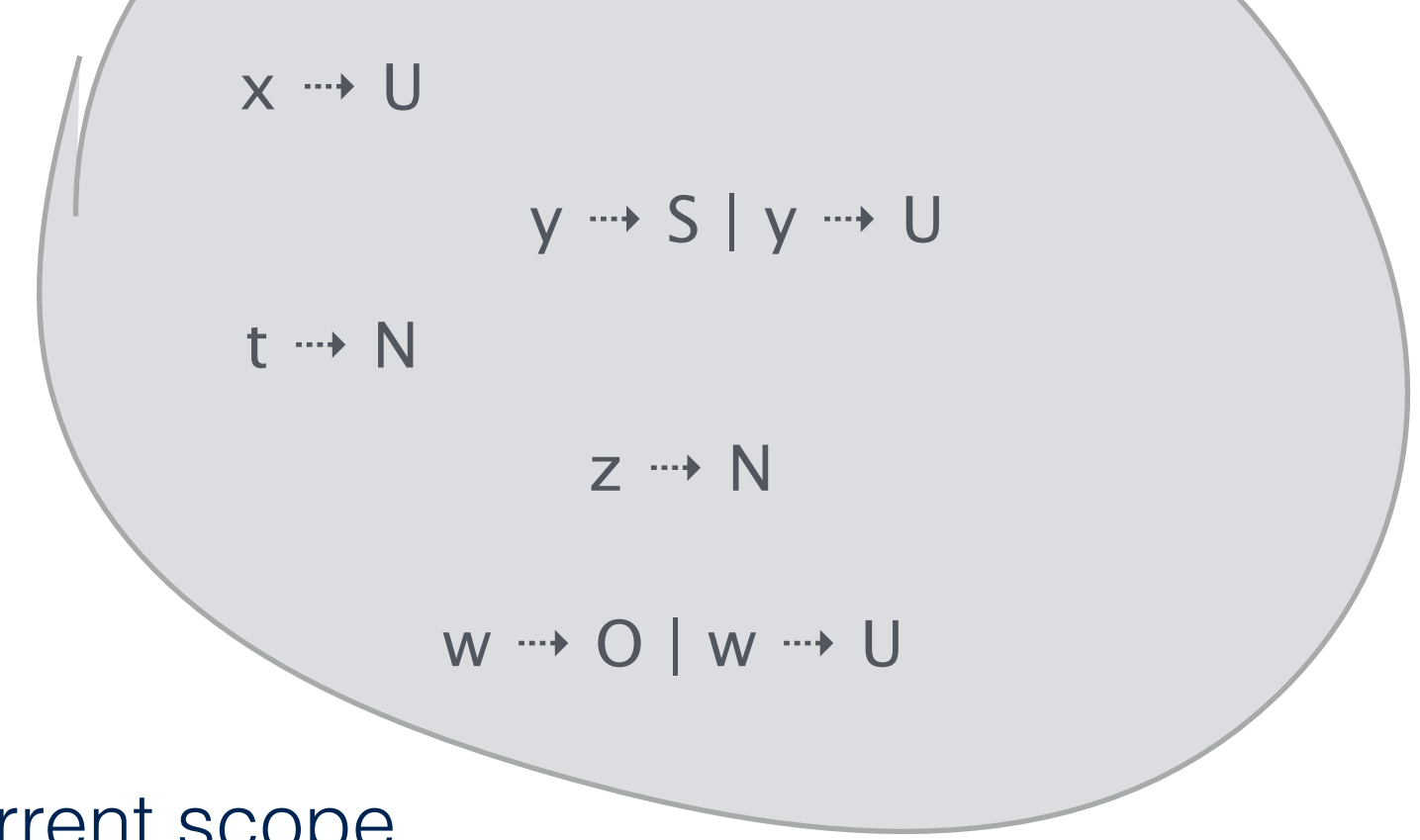
• passing \mapsto U (any other variable maps to N)

$$\begin{aligned} \Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\ \Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}] \\ \Theta_Q[s_1 \ s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q) \end{aligned}$$

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

```

passing = True
if not english:
    english = False

if not math:
    passing = False or bonus

if not math:
    passing = False or bonus
    
```

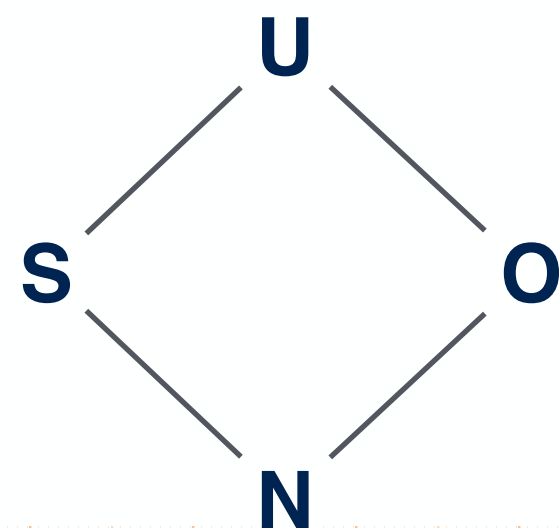
$$\begin{aligned}
\Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\
\Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\
\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\
&\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\
\Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}] \\
\Theta_Q[s_1 \ s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)
\end{aligned}$$

$$\begin{aligned}
x \rightarrow U &\Rightarrow x \rightarrow S \\
x \rightarrow O &\Rightarrow x \rightarrow N
\end{aligned}$$

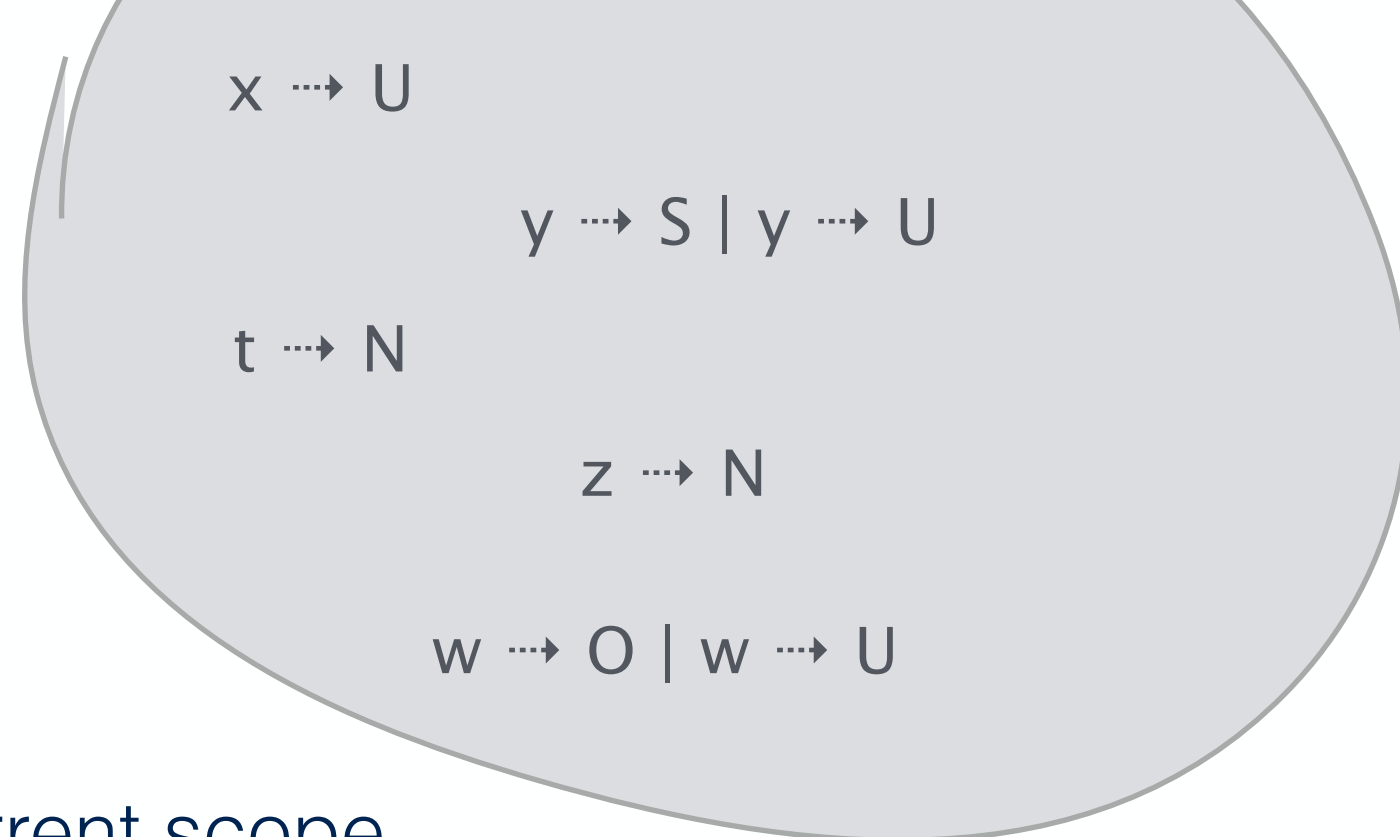
passing \rightarrow S | passing \rightarrow U
passing \rightarrow U

domain elements are stacks of maps matching nesting level of analyzed statements

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

if not math:

•
passing = **False or** bonus

•
passing -> U

$$\begin{aligned} \theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \theta_Q[s_2] \circ \text{PUSH}(q) \\ \theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \theta_Q[\text{if } b: s \text{ else: skip}] \\ \theta_Q[s_1 \ s_2](q) &\stackrel{\text{def}}{=} \theta_Q[s_1] \circ \theta_Q[s_2](q) \end{aligned}$$

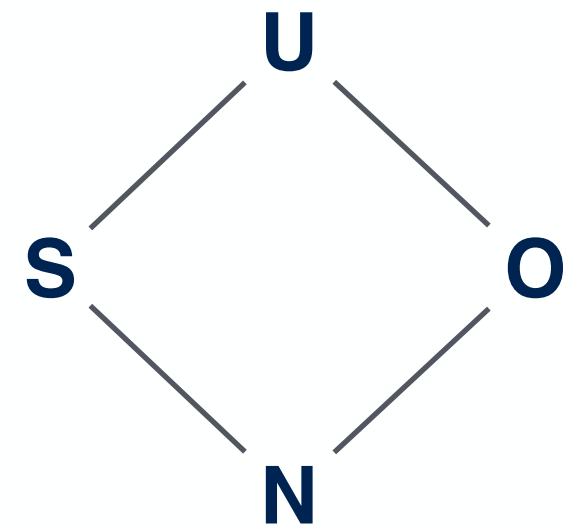
if the assigned variable was used (U or S) it becomes overwritten (O) if not also used in the expression being assigned; otherwise it becomes freshly used (U)

• bonus -> U, passing -> O | passing -> U

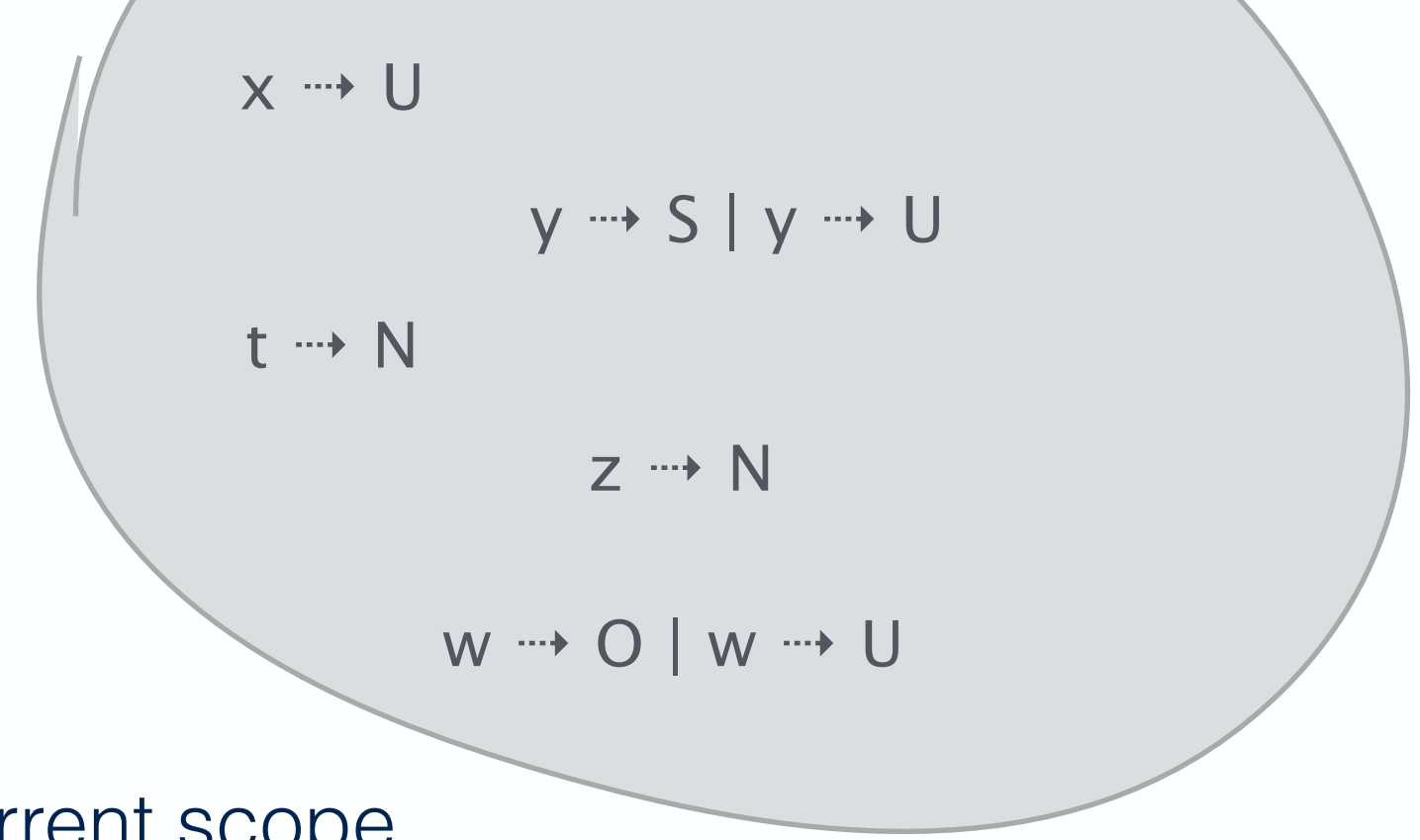
• passing -> S | passing -> U

• passing -> U

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

$$\begin{aligned} \Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\ \Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}] \\ \Theta_Q[s_1 \ s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q) \end{aligned}$$

```

passing = True

if not english:

    english = False

if not math:

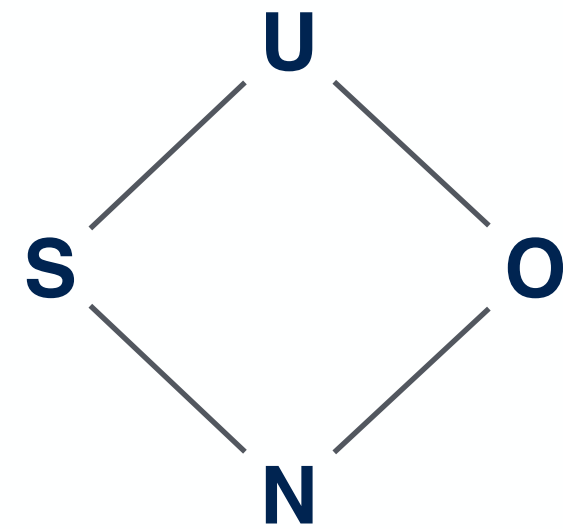
    passing = False or bonus

    • .....
    if not math:
        • .....
        passing = False or bonus
        • .....
        • .....
    
```

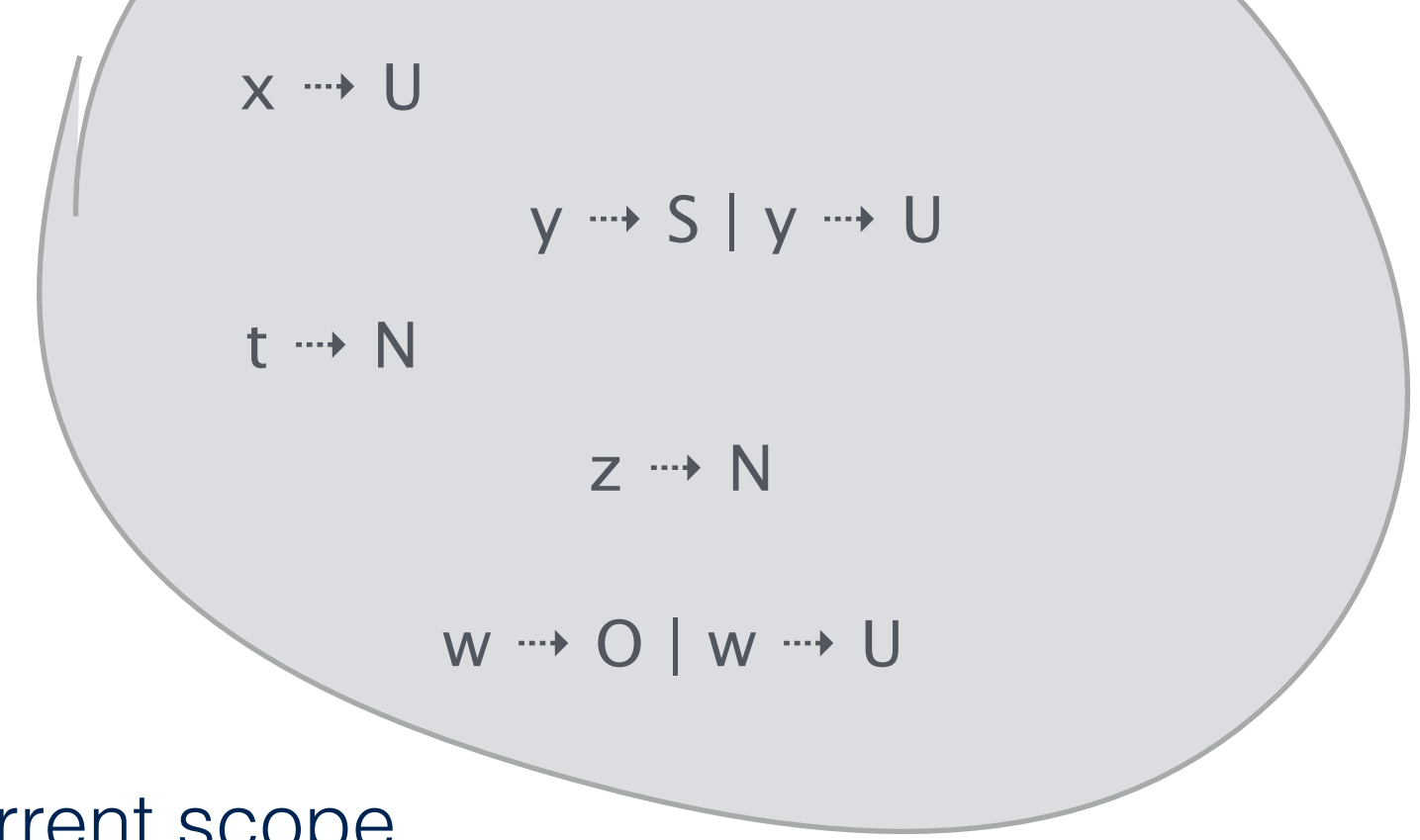
a variable becomes used (U) if it appears in the boolean condition of a statement that uses (U) or modifies (O) another variable

- math, bonus, passing -> U
- bonus -> U, passing -> O | passing -> U
- passing -> S | passing -> U
- passing -> U

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

if not math:

- ---

 passing = **False or** bonus
- ---

 passing = **False or** bonus

restores the previous value of variables (before increasing the nesting level) if it has not changed since

math, bonus, passing \mapsto U \leftarrow ;
math, bonus \mapsto U, passing \mapsto O
bonus \mapsto U, passing \mapsto O | passing \mapsto U
passing \mapsto S | passing \mapsto U
passing \mapsto U

$$\Theta_Q[\text{skip}](q) \stackrel{\text{def}}{=} q$$

$$\Theta_Q[x = e](q) \stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q)$$

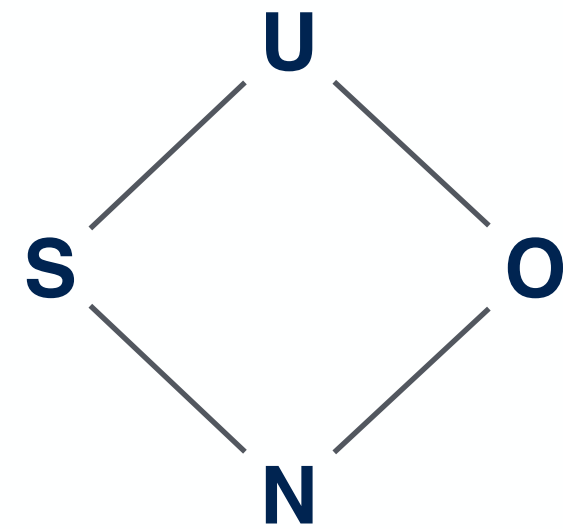
$$\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) \stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q)$$

$$\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q)$$

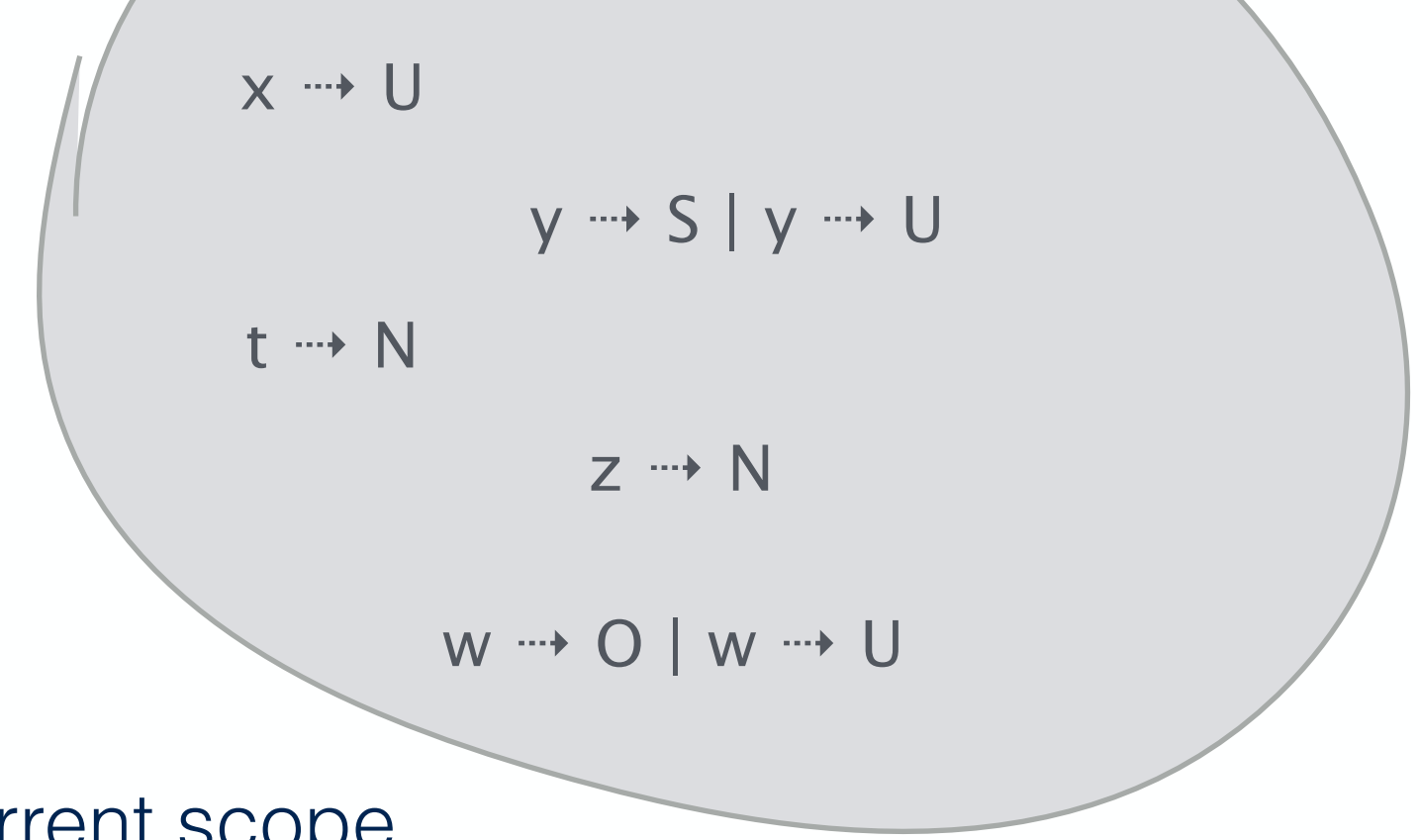
$$\Theta_Q[\text{while } b: s](q) \stackrel{\text{def}}{=} \text{lfp}_t^{\sqcup_Q} \Theta_Q[\text{if } b: s \text{ else: skip}]$$

$$\Theta_Q[s_1 \ s_2](q) \stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)$$

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

• math, bonus, passing -> U ←

if not math:

• math, bonus -> U, passing -> O ∪_Q passing -> U

• bonus -> U, passing -> O | passing -> U

• passing -> S | passing -> U

• passing -> U

$$\Theta_Q[\text{skip}](q) \stackrel{\text{def}}{=} q$$

$$\Theta_Q[x = e](q) \stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q)$$

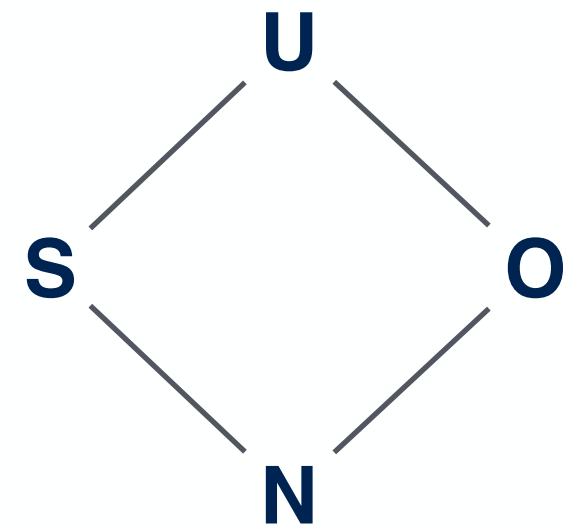
$$\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) \stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q)$$

$$\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q)$$

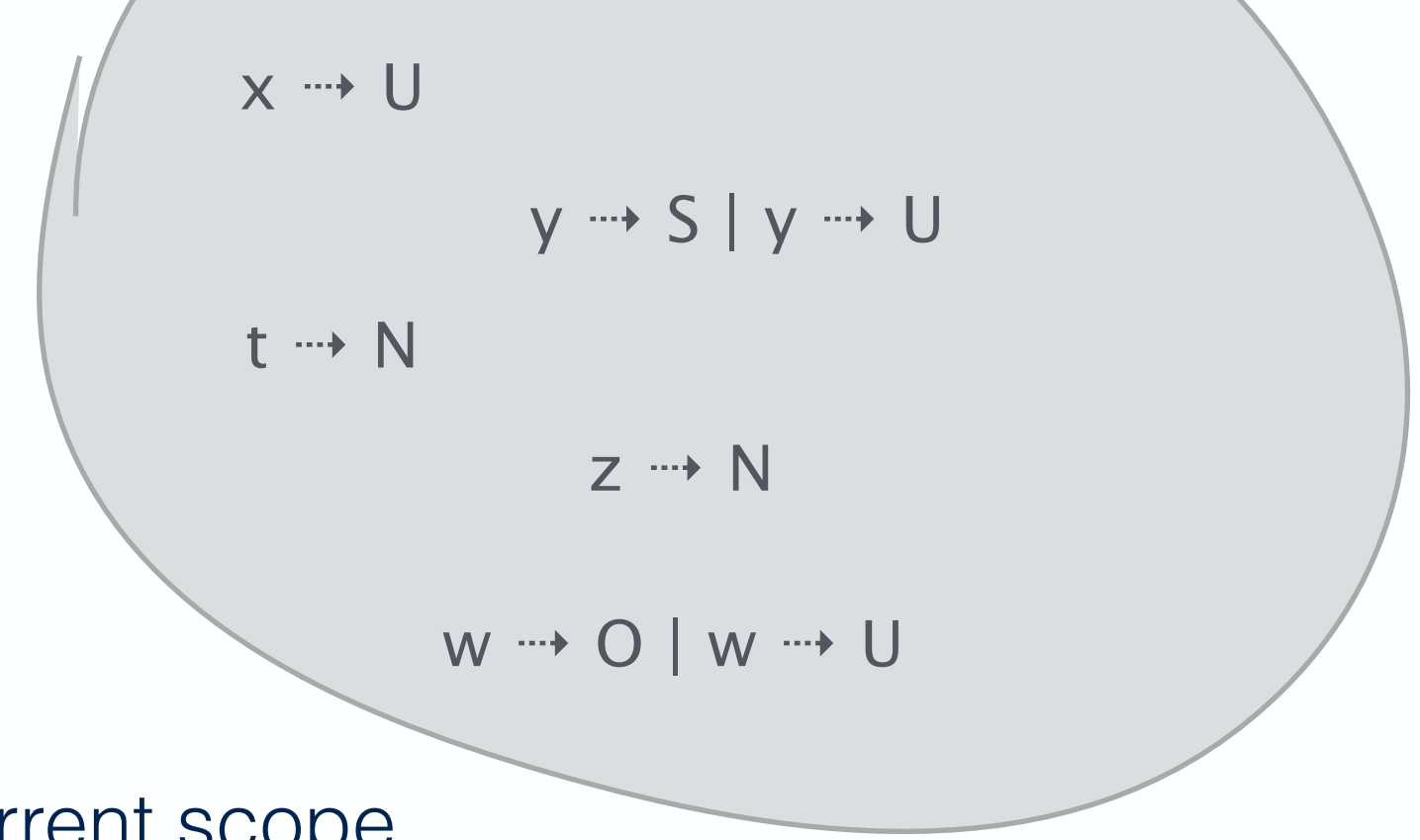
$$\Theta_Q[\text{while } b: s](q) \stackrel{\text{def}}{=} \text{lfp}_{\sqcup_Q}^{\Theta_Q} \Theta_Q[\text{if } b: s \text{ else: skip}]$$

$$\Theta_Q[s_1 \ s_2](q) \stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)$$

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



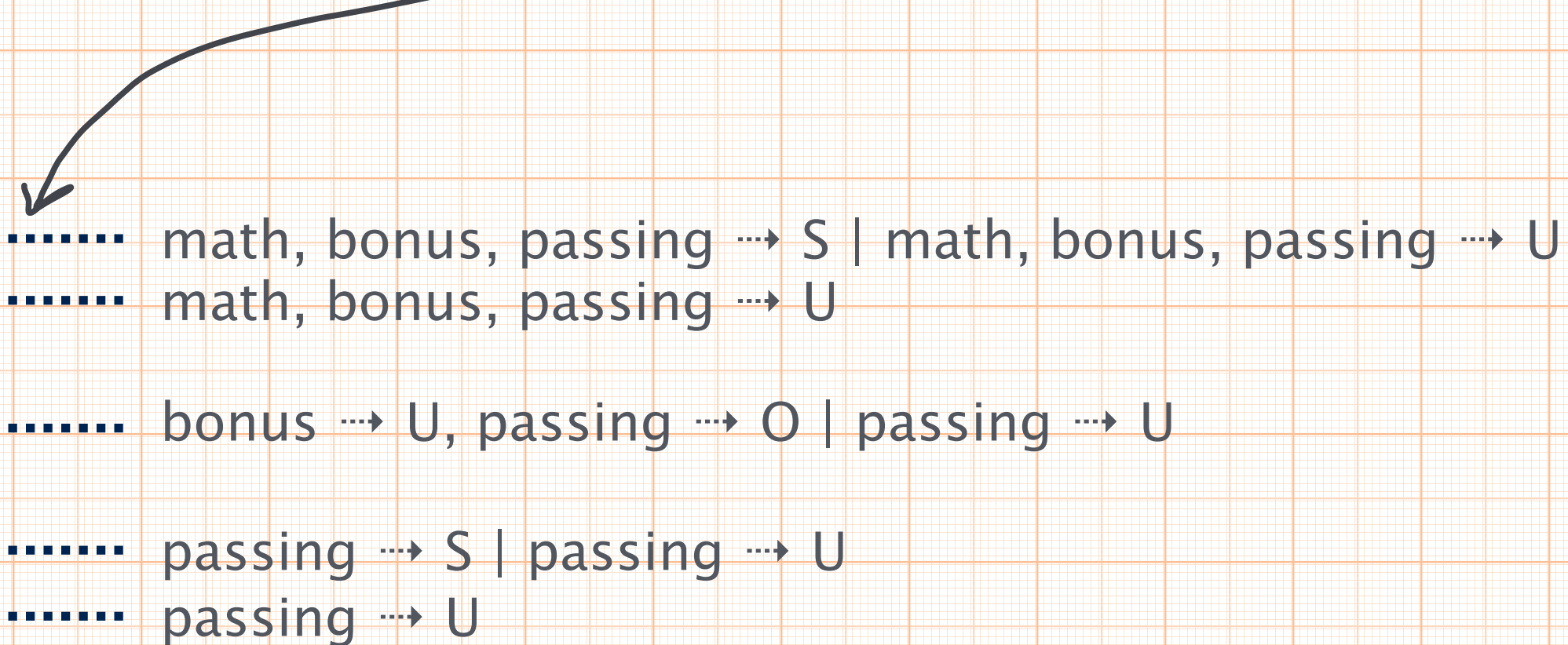
$\llbracket P \rrbracket_Q$

$$\begin{aligned} \Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\ \Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}] \\ \Theta_Q[s_1 \ s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q) \end{aligned}$$

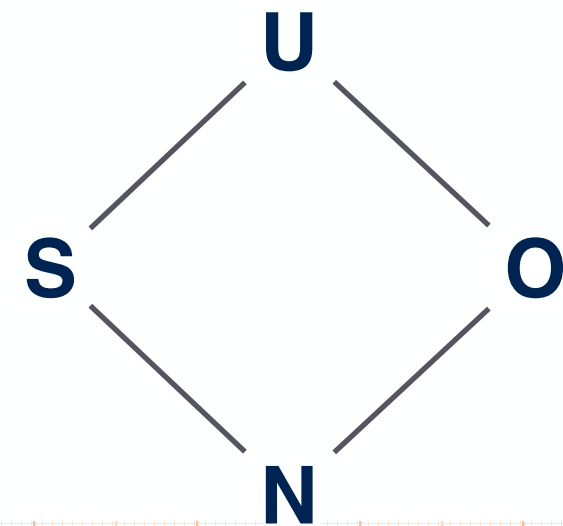
```

passing = True
if not english:
    english = False

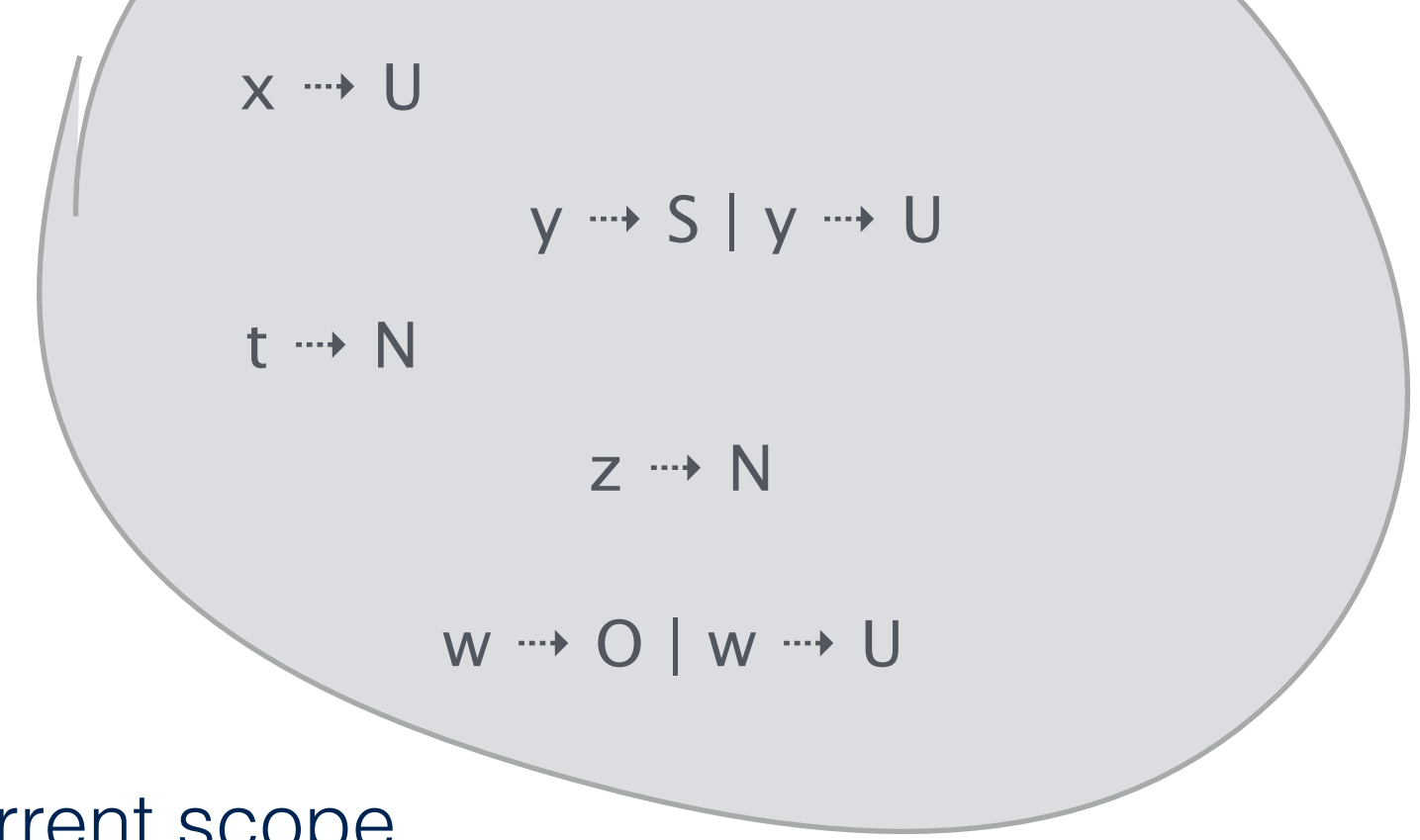
if not math:
    passing = False or bonus
    • ..... math, bonus, passing -> S | math, bonus, passing -> U
    • ..... math, bonus, passing -> U
    if not math:
        • ..... bonus -> U, passing -> O | passing -> U
        passing = False or bonus
        • ..... passing -> S | passing -> U
        • ..... passing -> U
    
```



Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

```

passing = True
if not english:
    english = False

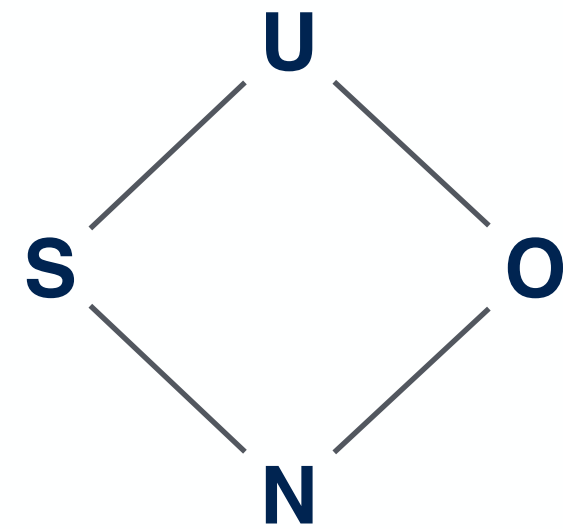
if not math:
    passing = False or bonus

if not math:
    passing = False or bonus
    
```

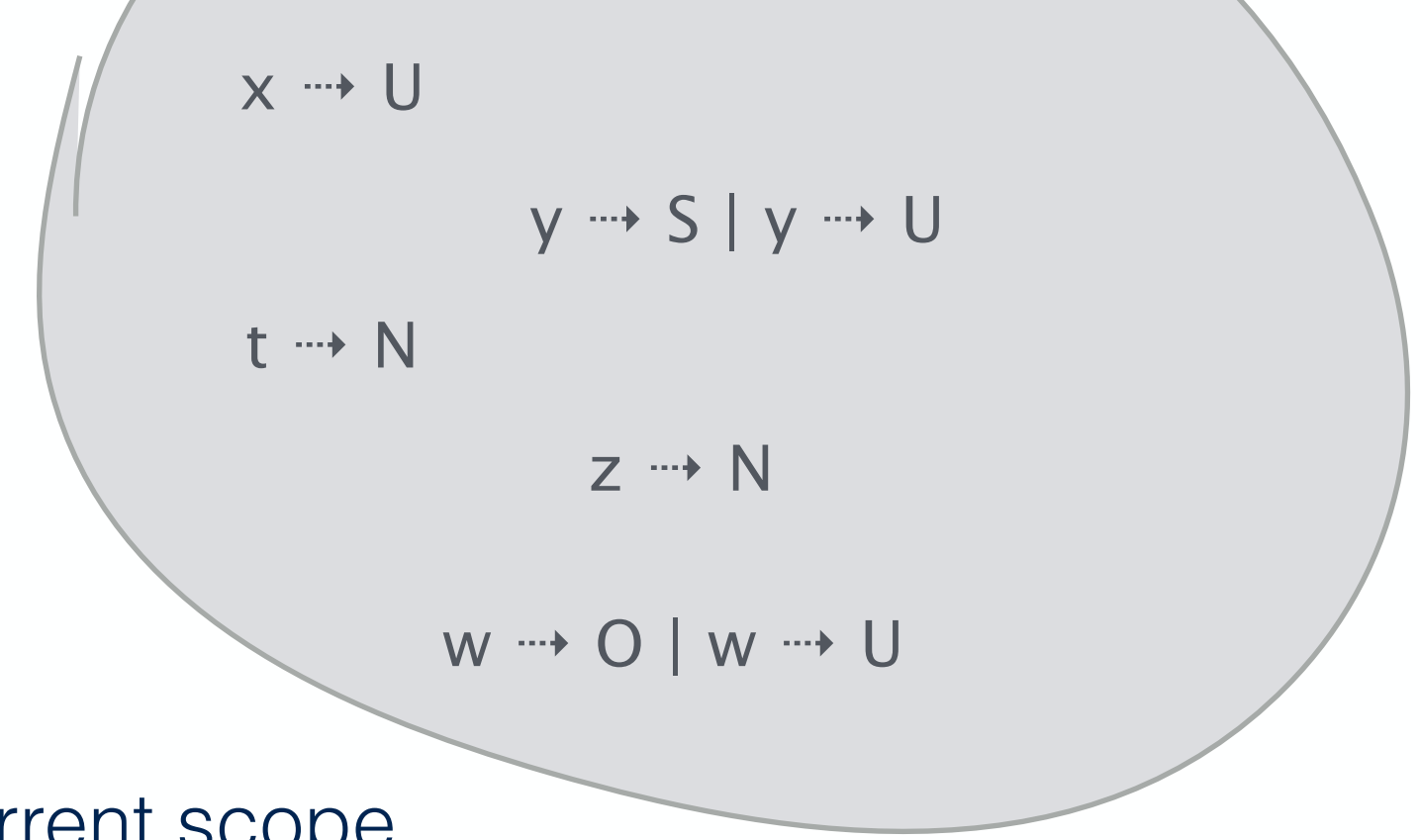
$$\begin{aligned}
\Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\
\Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\
\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\
&\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\
\Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}] \\
\Theta_Q[s_1 s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)
\end{aligned}$$

- math \mapsto S, bonus \mapsto U, passing \mapsto O | ...
- math, bonus, passing \mapsto S | math, bonus, passing \mapsto U
- math, bonus, passing \mapsto U
- bonus \mapsto U, passing \mapsto O | passing \mapsto U
- passing \mapsto S | passing \mapsto U
- passing \mapsto U

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

$$\begin{aligned} \Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\ \Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{if } b \stackrel{?}{=} \Theta_Q[\text{if } b: s \text{ else: skip}] \\ \Theta_Q[s_1; s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q) \end{aligned}$$

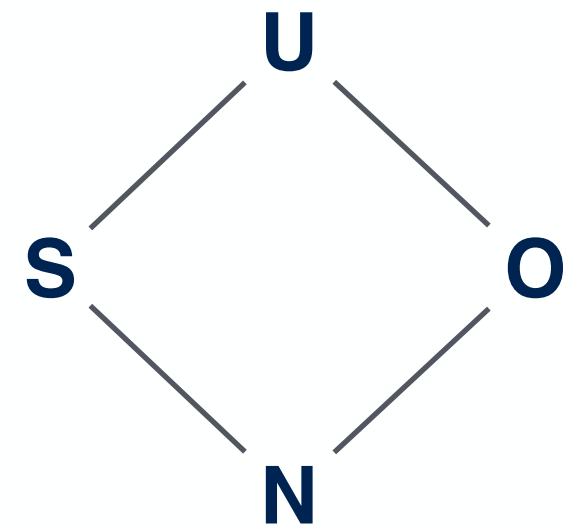
```

passing = True
if not english:
    english = False
    if not math:
        passing = False or bonus
    if not math:
        passing = False or bonus

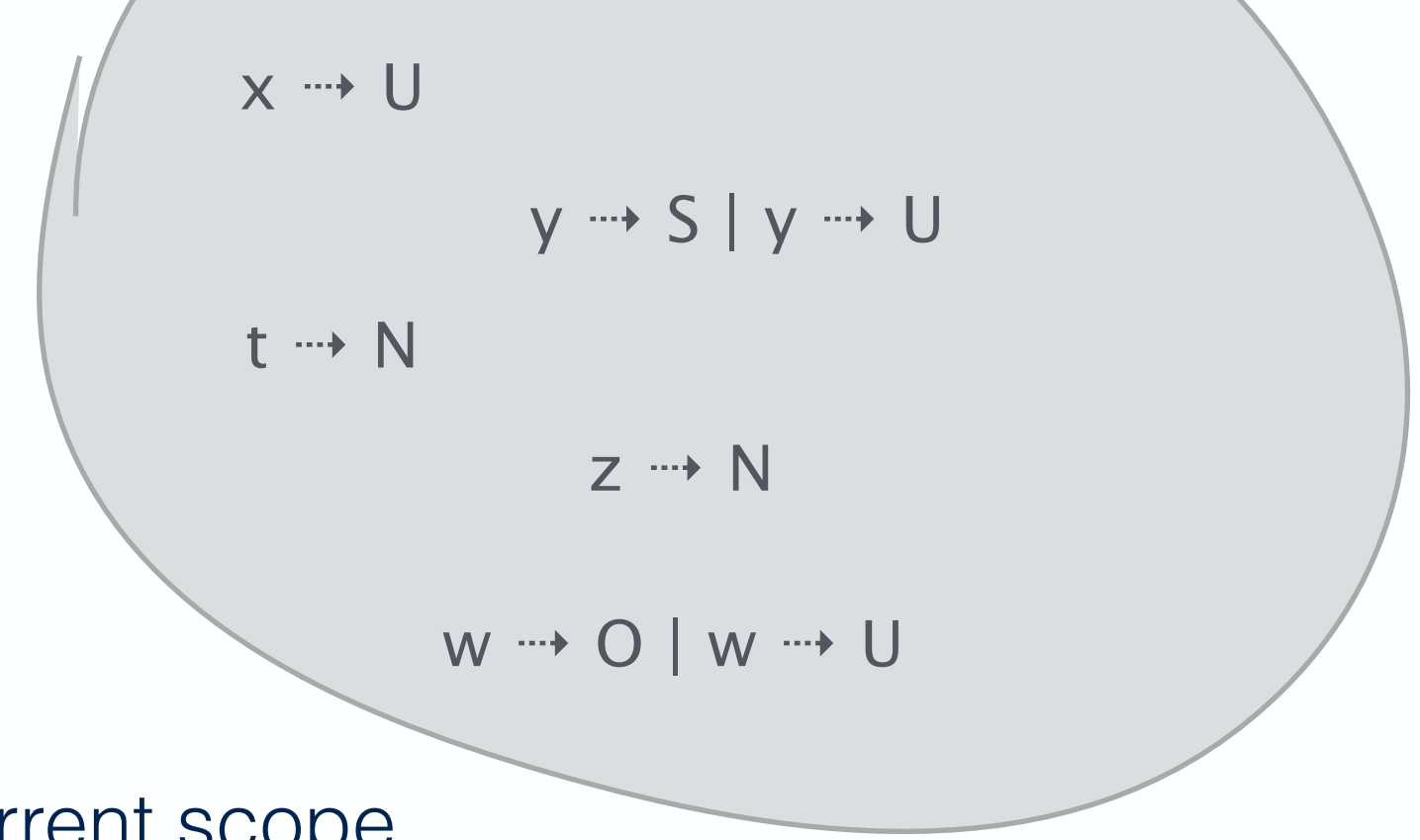
```

- math, bonus, passing \rightarrow U
- math \rightarrow S, bonus \rightarrow U, passing \rightarrow O | ...
- math, bonus, passing \rightarrow S | math, bonus, passing \rightarrow U
- math, bonus, passing \rightarrow U
- bonus \rightarrow U, passing \rightarrow O | passing \rightarrow U
- passing \rightarrow S | passing \rightarrow U
- passing \rightarrow U

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

$$\begin{aligned} \Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q) \\ \Theta_Q[\text{while } b: s](q) &\stackrel{\text{def}}{=} \text{lfp}_t^{\sqcup_Q} \Theta_Q[\text{if } b: s \text{ else: skip}] \\ \Theta_Q[s_1 \ s_2](q) &\stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q) \end{aligned}$$

```

passing = True
if not english:
    english = False
    • .....
    • .....
if not math:
    • .....
    passing = False or bonus
    • .....
if not math:
    • .....
    passing = False or bonus
    • .....
    • .....

```

• math, bonus, passing \mapsto S | math, bonus, passing \mapsto U

• math, bonus, passing \mapsto U

• math \mapsto S, bonus \mapsto U, passing \mapsto O | ...

• math, bonus, passing \mapsto S | math, bonus, passing \mapsto U

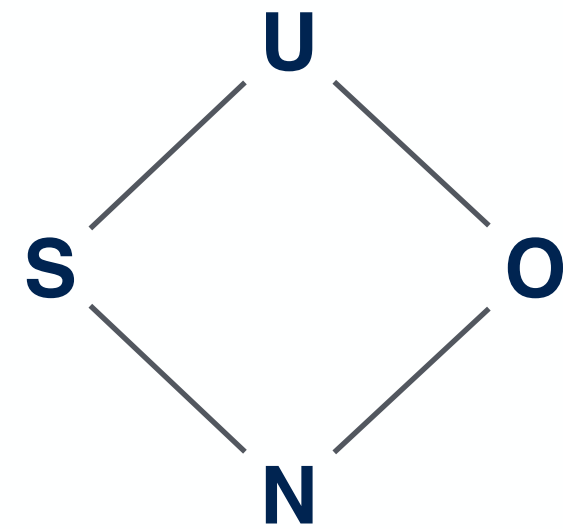
• math, bonus, passing \mapsto U

• bonus \mapsto U, passing \mapsto O | passing \mapsto U

• passing \mapsto S | passing \mapsto U

• passing \mapsto U

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used

$x \mapsto U$
 $y \mapsto S \mid y \mapsto U$
 $t \mapsto N$
 $z \mapsto N$
 $w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_Q$

```

passing = True

if not english:
    english = False

if not math:
    passing = False or bonus

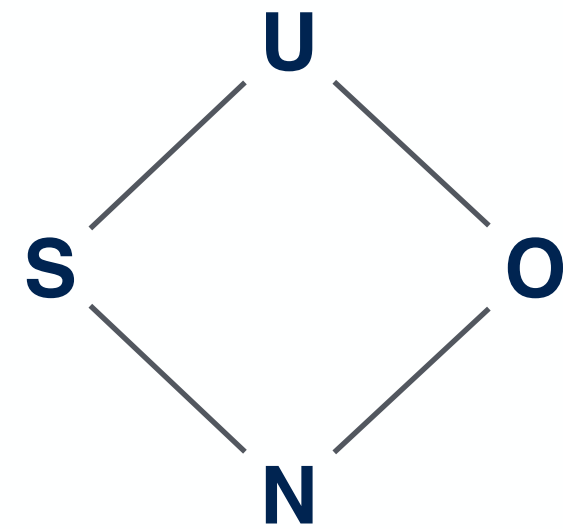
if not math:
    passing = False or bonus
    
```

$\text{math, bonus, passing} \mapsto S \mid \text{math, bonus, passing} \mapsto U$
 $\text{math, bonus, passing} \mapsto S \mid \text{math, bonus, passing} \mapsto U$
 $\text{math, bonus, passing} \mapsto U$

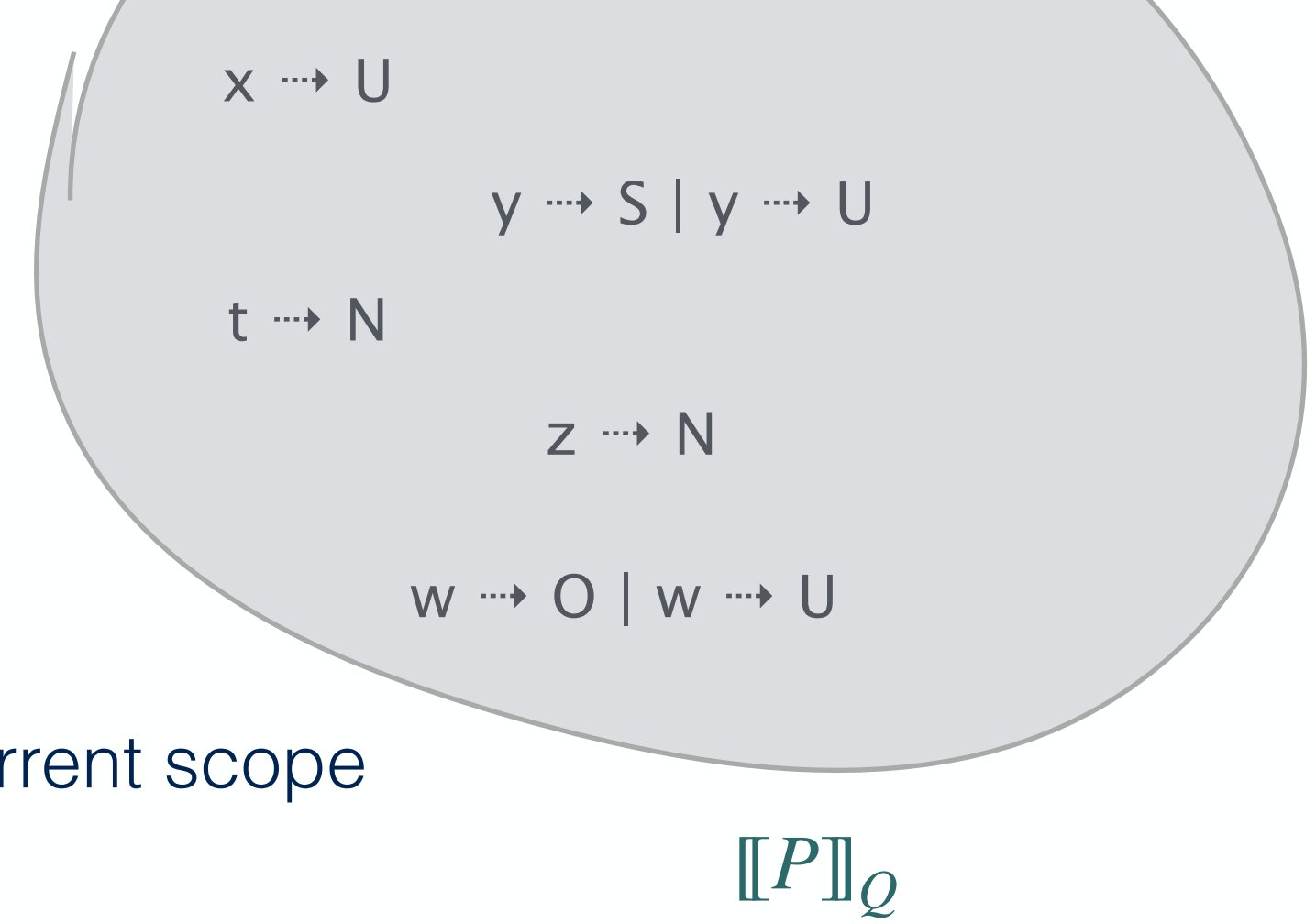
$\Theta_Q[\text{skip}](q) \stackrel{\text{def}}{=} q$
 $\Theta_Q[x = e](q) \stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q)$
 $\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) \stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q)$
 $\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q)$
 $\Theta_Q[\text{while } b: s](q) \stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}]$
 $\Theta_Q[s_1 \ s_2](q) \stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)$

$\text{passing} \mapsto U$
 $\text{passing} \mapsto U$

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



```

passing = True
●..... math, bonus, passing -> U
if not english:
  ●..... math, bonus, passing -> S | math, bonus, passing -> U
  english = False
  ●..... math, bonus, passing -> S | math, bonus, passing -> U
  ●..... math, bonus, passing -> U
if not math:
  ●.....
  passing = False or bonus
  ●.....
  if not math:
    ●.....
    passing = False or bonus
    ●.....
    ●..... passing -> U
  ●..... passing -> U
  
```

$$\Theta_Q[\text{skip}](q) \stackrel{\text{def}}{=} q$$

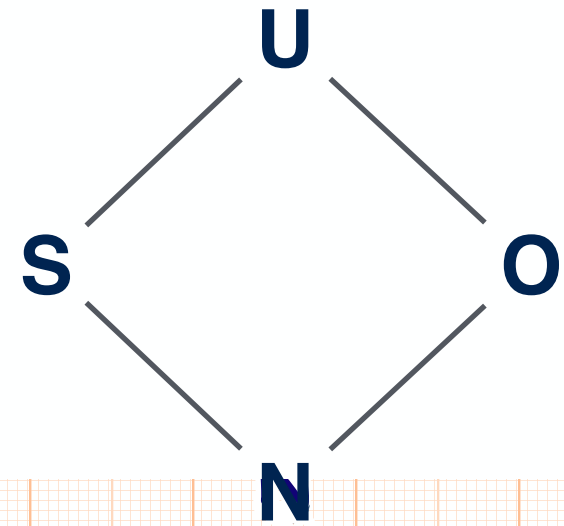
$$\Theta_Q[x = e](q) \stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q)$$

$$\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) \stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q)$$

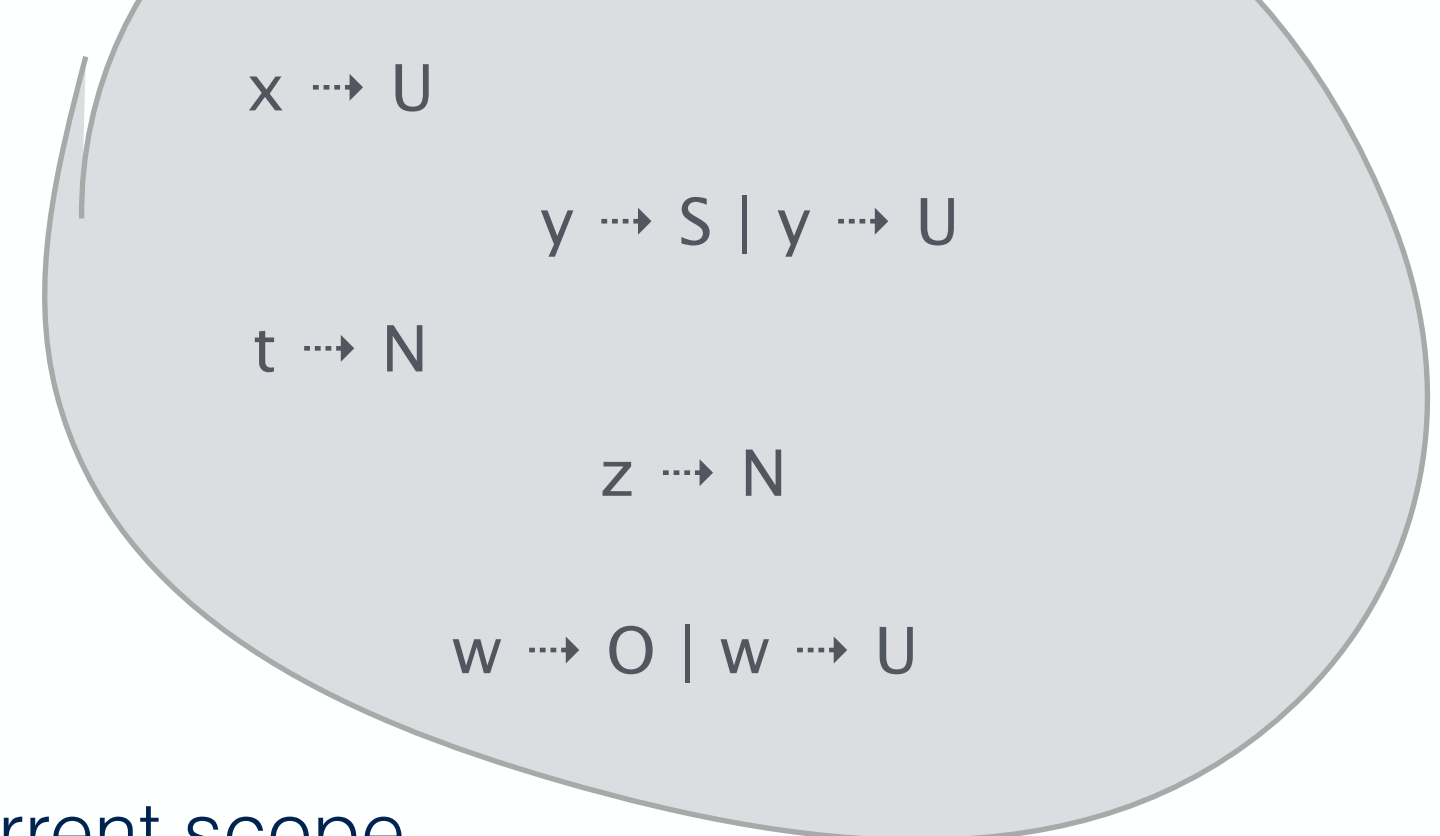
$$\Theta_Q[\text{while } b: s](q) \stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}]$$

$$\Theta_Q[s_1 \ s_2](q) \stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)$$

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used



$\llbracket P \rrbracket_Q$

```

    • .....
    passing = True
    • .....
    if not english:
        • .....
        english = False
        • .....
    • .....
    if not math:
        • .....
        passing = False or bonus
        • .....
    • .....
    if not math:
        • .....
        passing = False or bonus
        • .....
    • .....
  
```

the input variables english and science are definitely not used by the program

```

    math, bonus -> U, passing -> O
    math, bonus, passing -> U
    math, bonus, passing -> S | math, bonus, passing -> U
    math, bonus, passing -> S | math, bonus, passing -> U
    math, bonus, passing -> U
  
```

$\Theta_Q[\text{skip}](q) \stackrel{\text{def}}{=} q$

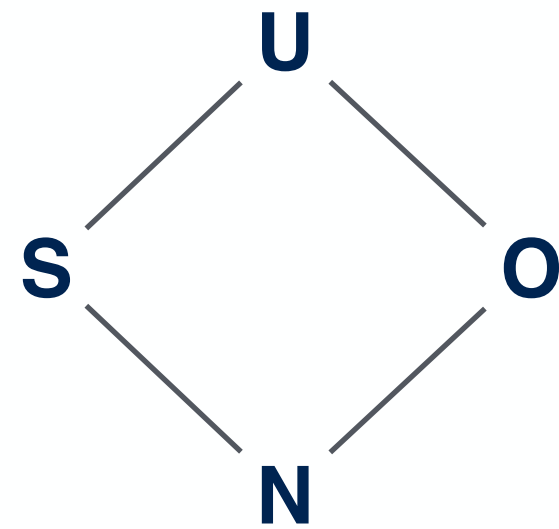
$\Theta_Q[x = e](q) \stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q)$

$\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) \stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q)$
 $\sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q)$

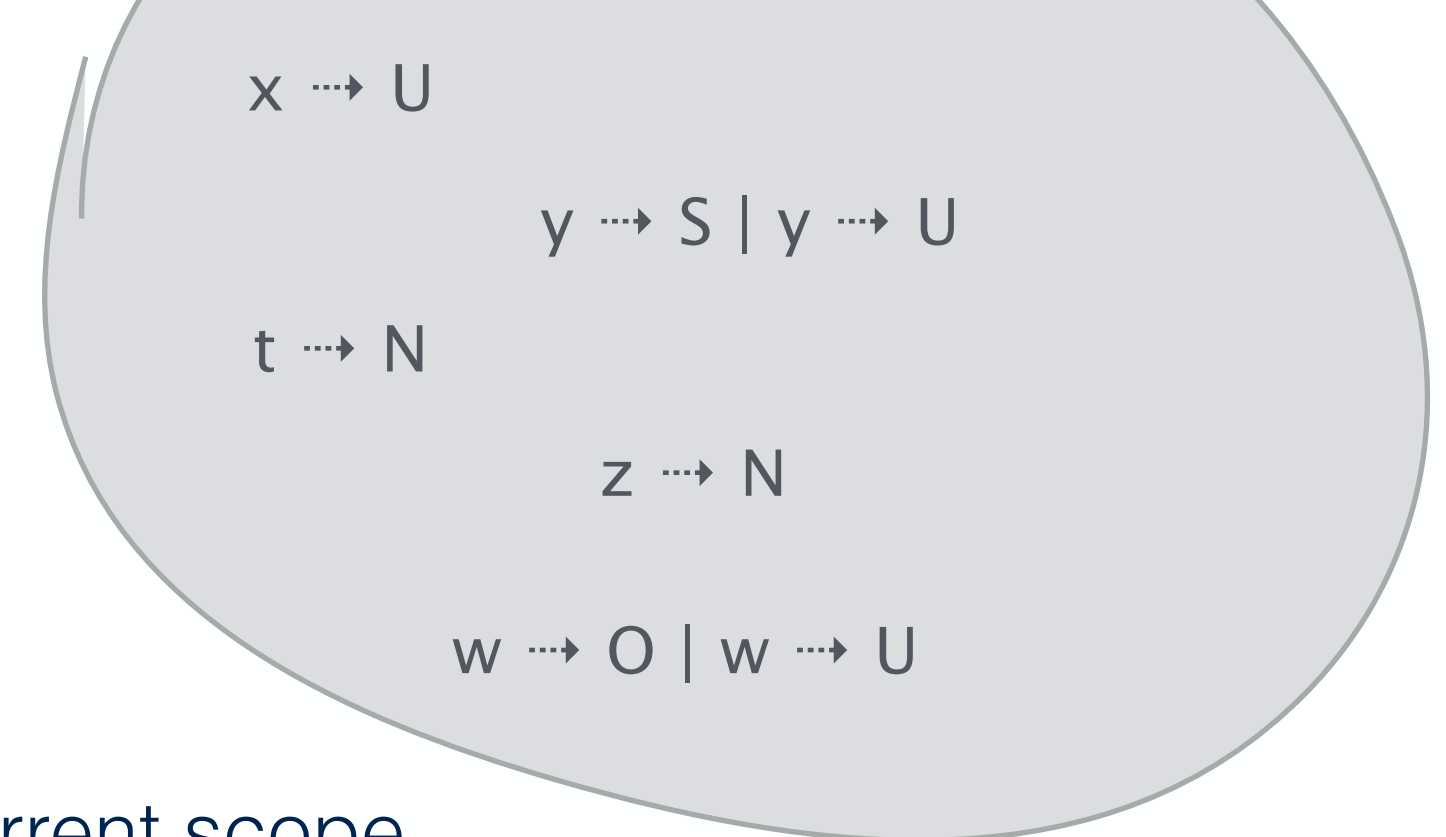
$\Theta_Q[\text{while } b: s](q) \stackrel{\text{def}}{=} \text{lfp}_t^{\sqsubseteq_Q} \Theta_Q[\text{if } b: s \text{ else: skip}]$

$\Theta_Q[s_1 \ s_2](q) \stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)$

Syntactic (Non-)Usage



- **U**: used in the current scope (or an inner scope)
- **S**: used in an outer scope
- **O**: used in an outer scope and overridden in the current scope
- **N**: not used

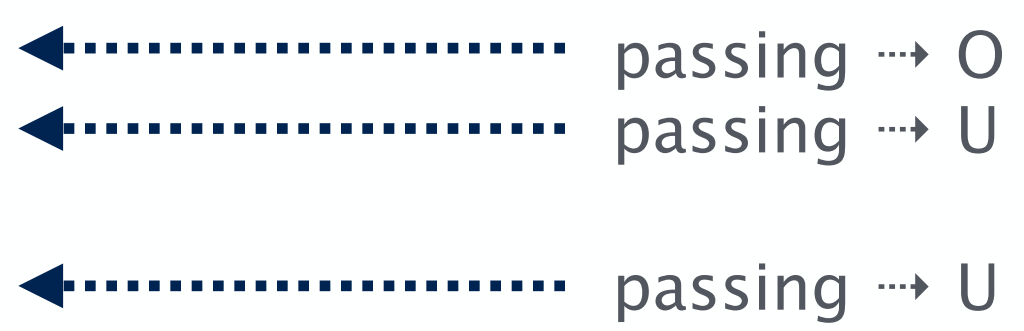


$\llbracket P \rrbracket_Q$

e ::= v | x | not e | e and e | e or e
s ::= skip | x = e | if e: s else: s | while e: s | s s (expressions)
 (statements)

$\Theta_Q[\text{skip}](q) \stackrel{\text{def}}{=} q$
 $\Theta_Q[x = e](q) \stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q)$
 $\Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) \stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q)$
 $\quad \sqcup_Q \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_2] \circ \text{PUSH}(q)$
 $\Theta_Q[\text{while } b: s](q) \stackrel{\text{def}}{=} \text{lfp}_t^{\sqcup_Q} \Theta_Q[\text{if } b: s \text{ else: skip}]$
 $\Theta_Q[s_1 \ s_2](q) \stackrel{\text{def}}{=} \Theta_Q[s_1] \circ \Theta_Q[s_2](q)$

passing = **True**
while not english:
 english = **False**

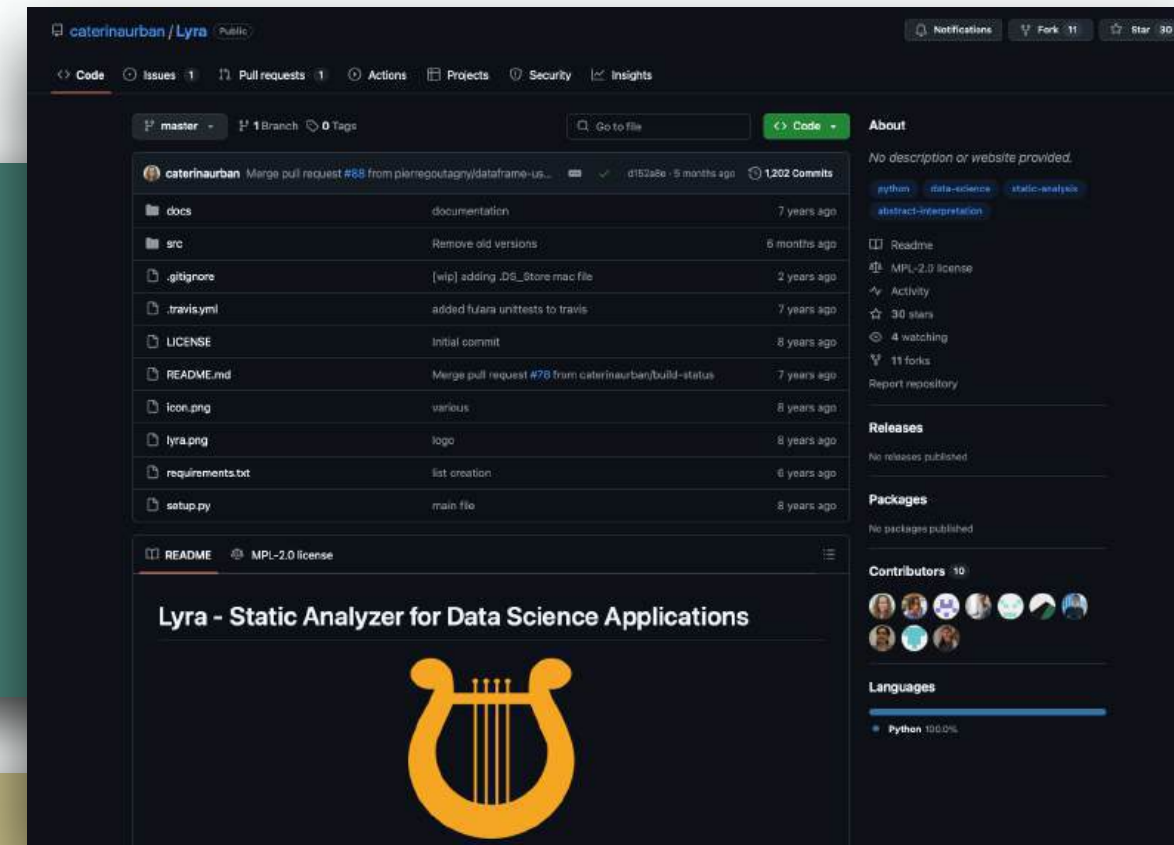


Theorem
 $P \models \mathcal{N}_J^+ \Leftarrow \{\llbracket P \rrbracket\} \subseteq \llbracket P \rrbracket_Q^{\sharp} \subseteq \mathcal{N}_J^+$

Data Usage Static Analysis [Urban18]

3-Step Recipe

practical tools
targeting specific programs

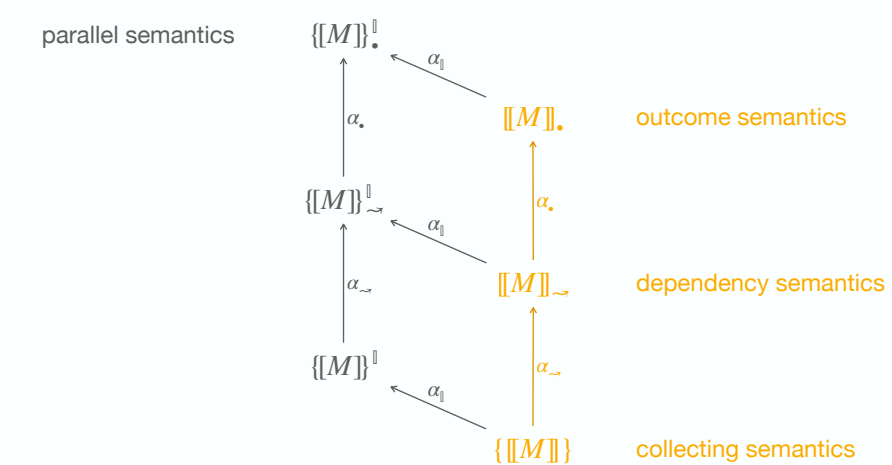


secure information flow

abstract semantics, abstract domains
algorithmic approaches to decide program properties

strongly-live variable analysis

Hierarchy of Semantics



syntactic non-usage

concrete semantics
mathematical models of the program behavior



Data Leakage

A case study of algorithm-assisted decision making in child maltreatment hotline screening decisions

Alexandra Chouldechova
Heinz College
Carnegie Mellon University
Pittsburgh, PA, 15213, USA

Emily Putnam-Hornstein
Suzanne Dworak-Peck School of Social Work
University of Southern California
Los Angeles, CA, 90089, USA

Diana Benavides-Prado
Oleksandr Fialko

Rhema Vaithianathan
Centre for Social Data Analytics
Auckland University of Technology
Auckland, New Zealand

Editors: Sorelle A. Friedler and Christo Wilson

Abstract

Every year there are more than 3.6 million referrals made to child protection agencies across the US. The practice of screening calls is left to each jurisdiction, potentially leading to large variation in the way which referrals are treated across the country. Whilst increasing access to linked administrative data is available, it is often used by welfare workers to make system-wide decisions about all children and adults on a single referral. Risk prediction models that use collected administrative data can help workers to better identify cases likely to result in adverse outcomes. However, the use of predictive analytics in the area of child welfare is contentious. There is a possibility that some correlations, such as those in poverty or similar racial and ethnic groups, are spurious or advantaged by the reliance on administrative data. On the other hand, these analytics tools can augment or replace human judgments, which themselves are biased and imperfect. In this paper we describe our work on developing, validating, auditing, and deploying a risk

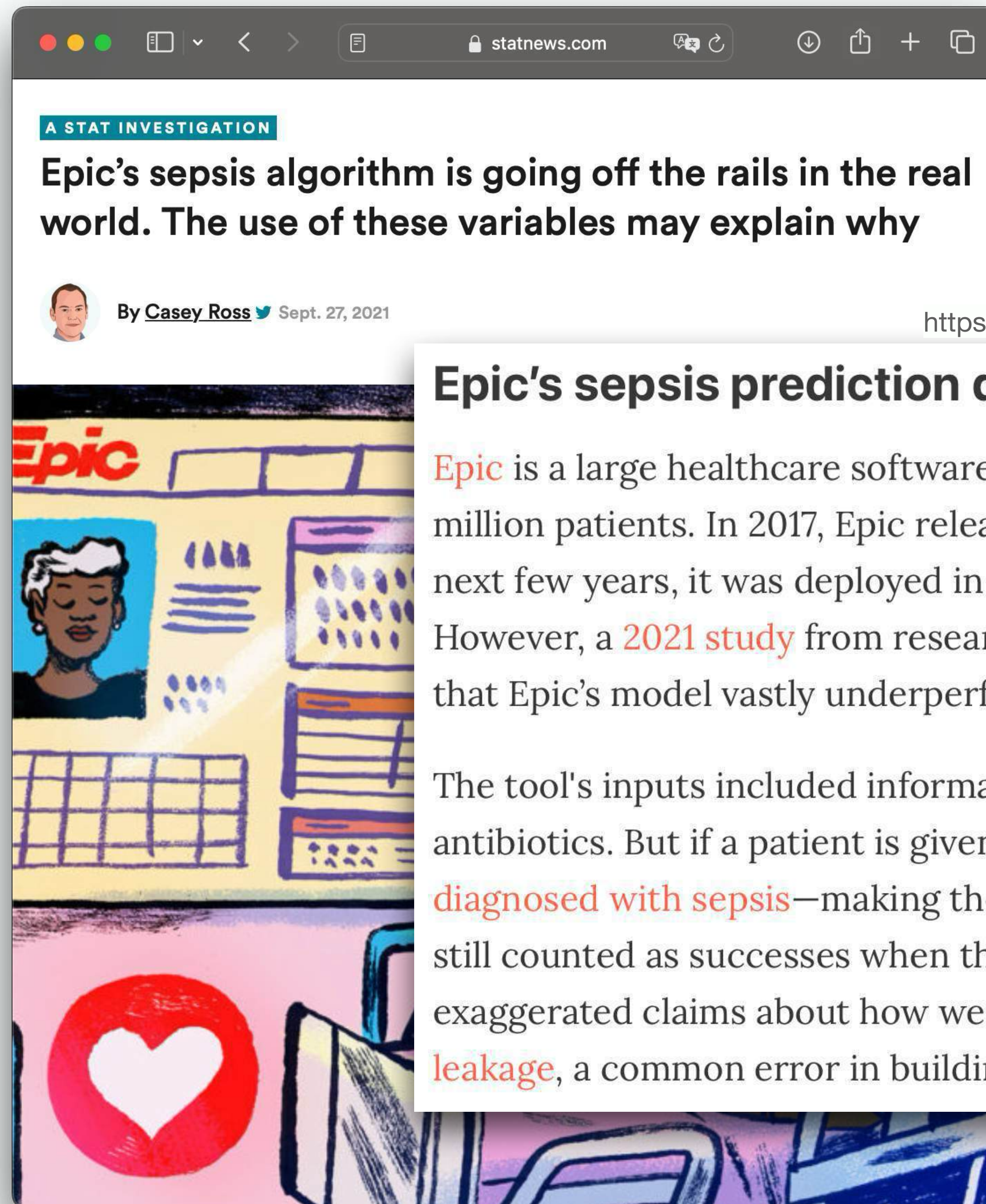
and linked administrative data available, it is difficult for child welfare workers to make systematic use of historical information about all the children and adults on a single referral.

<https://www.aisnakeoil.com/p/the-bait-and-switch-behind-ai-risk>

Family separation in Allegheny county

In 2016, Allegheny county in Pennsylvania adopted the Allegheny Family Screening Tool (AFST) to predict which children are at risk of maltreatment. AFST is used to decide which families should be investigated by social workers. In these investigations, social workers can forcibly remove children from their families and place them in foster care, **even if there are no allegations of abuse**—only poverty-based neglect.

Two years later, it was **discovered** that AFST suffered from data leakage, leading to exaggerated claims about its performance. In addition, the tool was **systematically biased** against Black families. When questioned, the creators trotted out the familiar defense that the **final decision is always made by a human decision-maker**.



<https://www.aisnakeoil.com/p/the-bait-and-switch-behind-ai-risk>

Epic's sepsis prediction debacle

Epic is a large healthcare software company. It stores health data for over 300 million patients. In 2017, Epic released a sepsis prediction model. Over the next few years, it was deployed in hundreds of hospitals across the U.S. However, a **2021 study** from researchers at the University of Michigan found that Epic's model vastly underperformed compared to the developer's claims.

The tool's inputs included information about whether a patient was given antibiotics. But if a patient is given antibiotics, they have **already been diagnosed with sepsis**—making the tool's prediction useless. These cases were still counted as successes when the developer evaluated the tool, leading to exaggerated claims about how well it performed. This is an example of **data leakage**, a common error in building AI tools.

Example

```
[1]: import numpy as np
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

[2]: data = pd.read_csv("data.csv")
X = data[["X_1", "X_2"]]
y = data[["y"]]

[3]: min_max_scaler = MinMaxScaler()
X = min_max_scaler.fit_transform(X)

[4]: X_train , X_test , y_train , y_test = train_test_split(X, y, test_size=0.025, random_state=2)

[ ]:

[5]: lr = LogisticRegression()
a = lr.fit(X_train , y_train)

[6]: y_pred = lr.predict(X_test)
accuracy_score(y_test , y_pred)

[6]: 0.67 ✓
```


Example

```
[1]: import numpy as np
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

[2]: data = pd.read_csv("data.csv")
X = data[["X_1", "X_2"]]
y = data[["y"]]

[ ]:

[3]: X_train , X_test , y_train , y_test = train_test_split(X, y, test_size=0.025, random_state=2)

[4]: min_max_scaler = MinMaxScaler()
X_train = min_max_scaler.fit_transform(X_train)
X_test = min_max_scaler.fit_transform(X_test)

[5]: lr = LogisticRegression()
a = lr.fit(X_train , y_train)

[6]: y_pred = lr.predict(X_test)
accuracy_score(y_test , y_pred)

[6]: 0.33 X
```

Data Leakage Static Analysis [Drobnjaković24]

3-Step Recipe

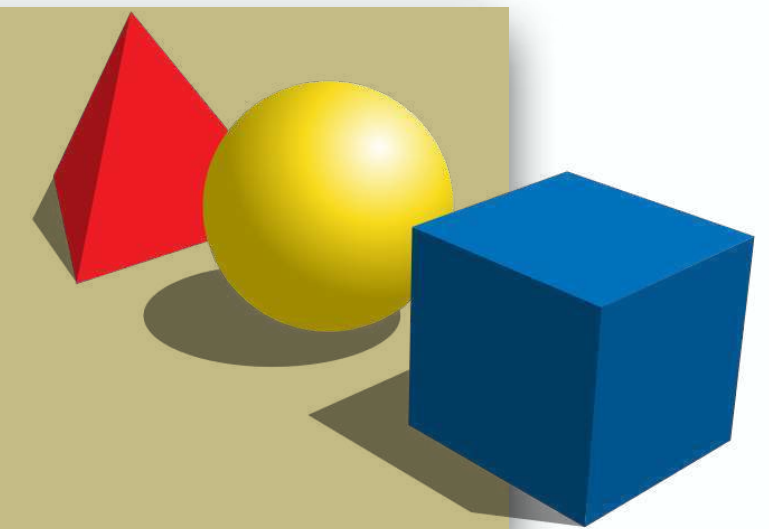
practical tools

targeting specific programs



abstract semantics, abstract domains

algorithmic approaches to decide program properties



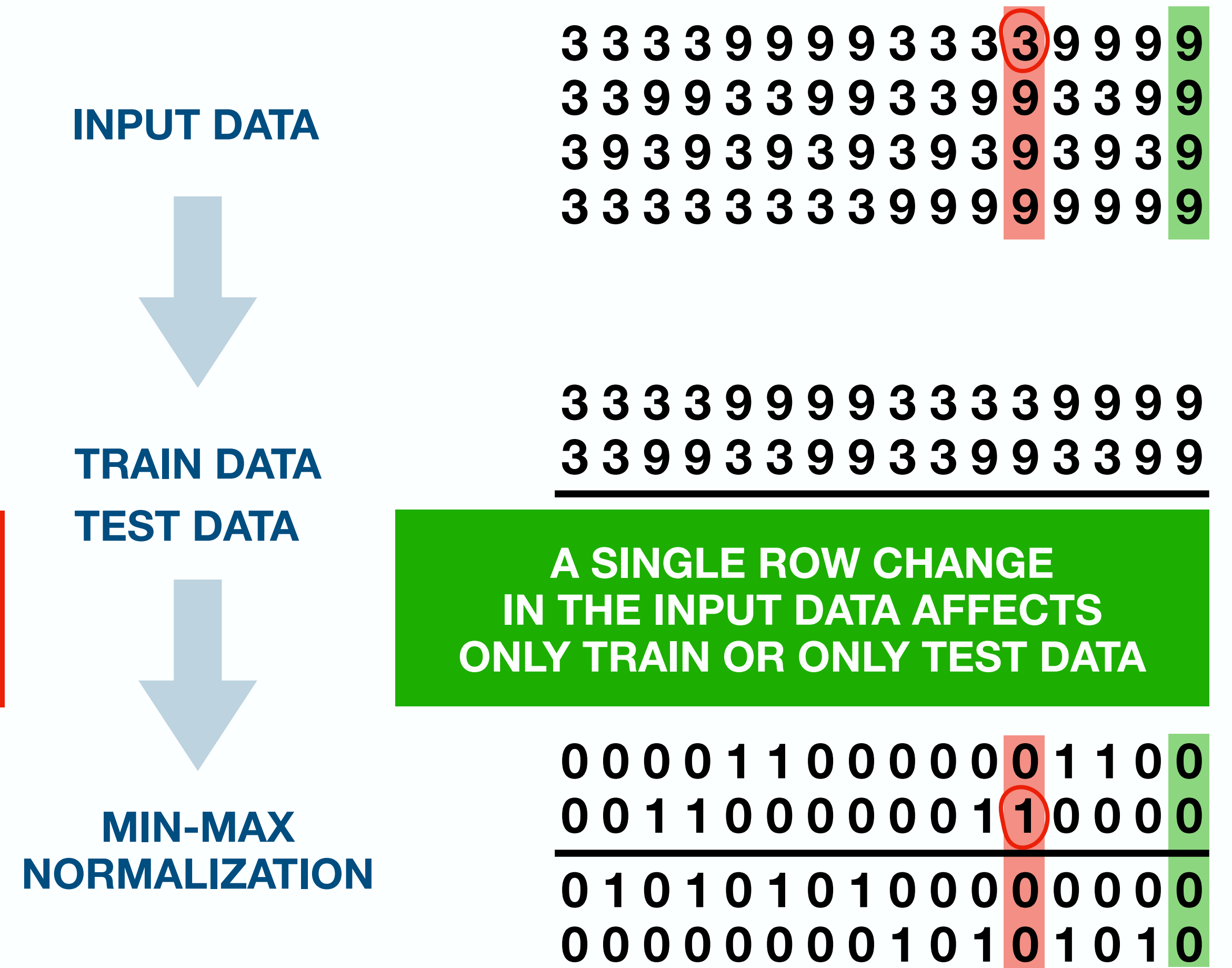
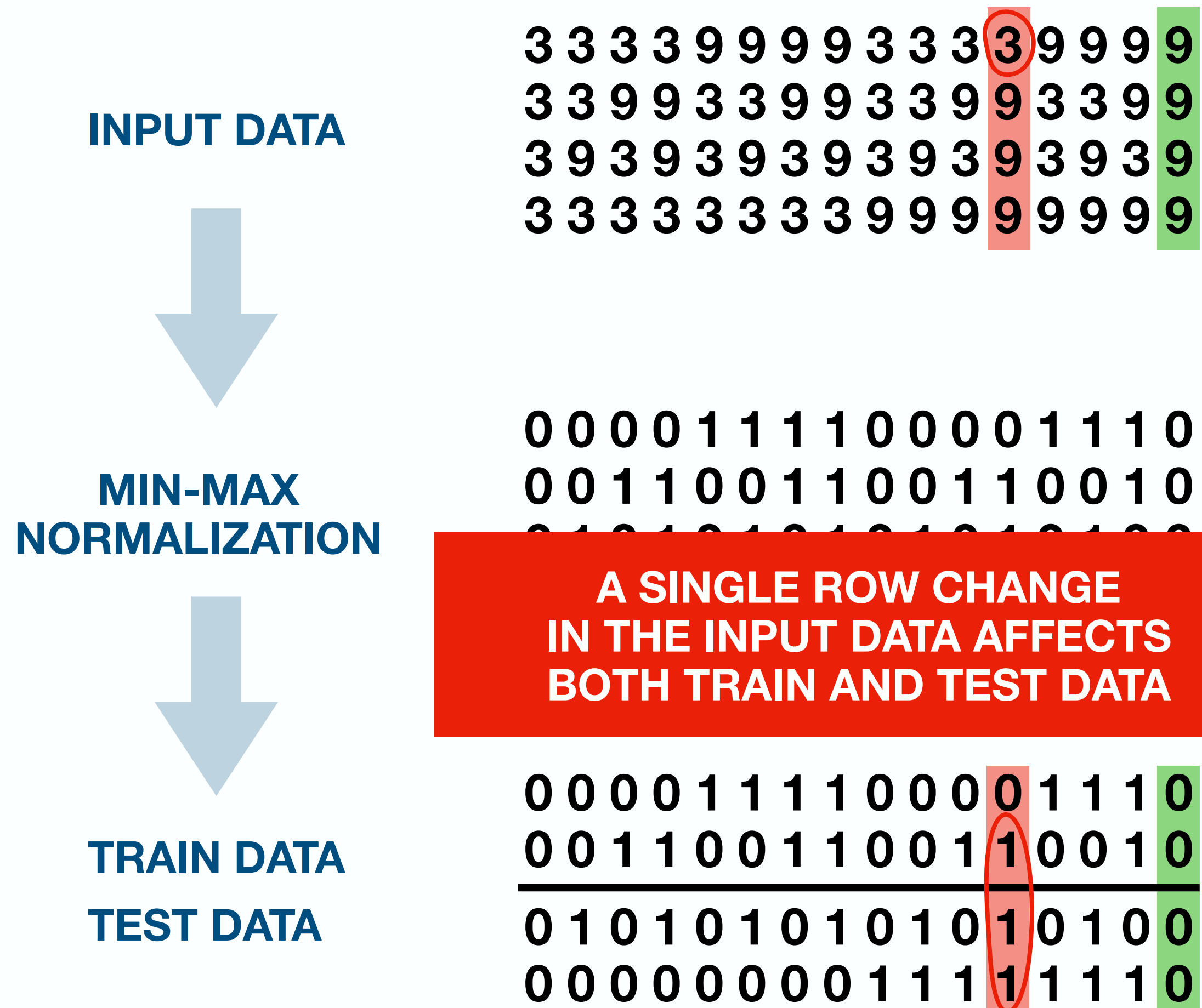
concrete semantics

mathematical models of the program behavior



(Absence of) Data Leakage

Independence of Training and Testing Data



(Absence of) Data Leakage

Independence of Training and Testing Data

$$\mathcal{F} \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \mid \text{INDEPENDENT}(\llbracket P \rrbracket) \}$$

input data frame variables

data frame row

used data frame variables

$$\text{INDEPENDENT}(\llbracket P \rrbracket) \stackrel{\text{def}}{=} \forall \sigma \in \llbracket P \rrbracket, i \in I_P, r \in R_i: \text{UNCHANGED}(\sigma, i, r, U_P^{\text{test}}) \vee \text{UNCHANGED}(\sigma, i, r, U_P^{\text{train}})$$

$$\text{UNCHANGED}(\sigma, i, r, U) \stackrel{\text{def}}{=} \forall \bar{v} \in \mathbb{V}^{C_i}: \sigma(i)[r] \neq \bar{v} \Rightarrow (\exists \sigma' \in \Upsilon \llbracket P \rrbracket: \sigma'(i)[r] = \bar{v} \wedge \eta(\sigma) = \eta(\sigma') \wedge \sigma(U) = \sigma'(U))$$

$$\eta(\sigma) \stackrel{\text{def}}{=} \lambda j: \lambda r': \begin{cases} \sigma(j)[r'] & j \in I_P \setminus \{i\} \vee r' \in R_i: r' \neq r \\ \top & \text{otherwise} \end{cases}$$

$$\sigma(X) = \sigma'(X) \stackrel{\text{def}}{=} \forall x \in X: \sigma(x) = \sigma'(x)$$

Data Leakage Static Analysis [Drobnjaković24]

3-Step Recipe

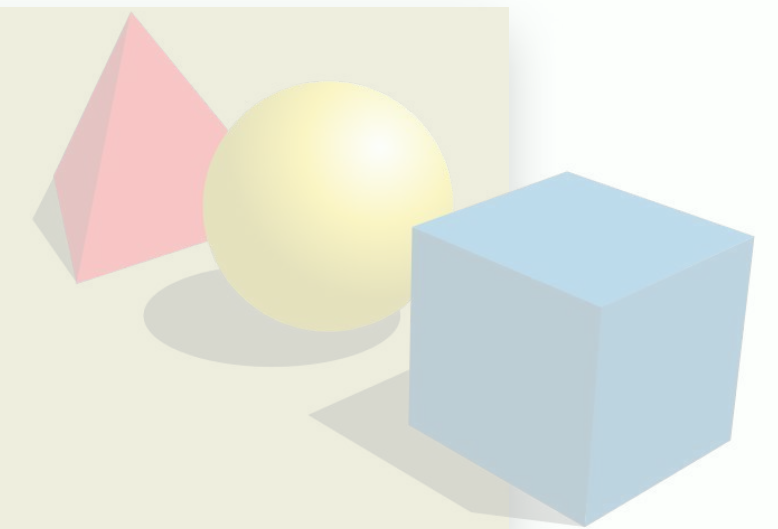
practical tools

targeting specific programs



abstract semantics, abstract domains

algorithmic approaches to decide program properties

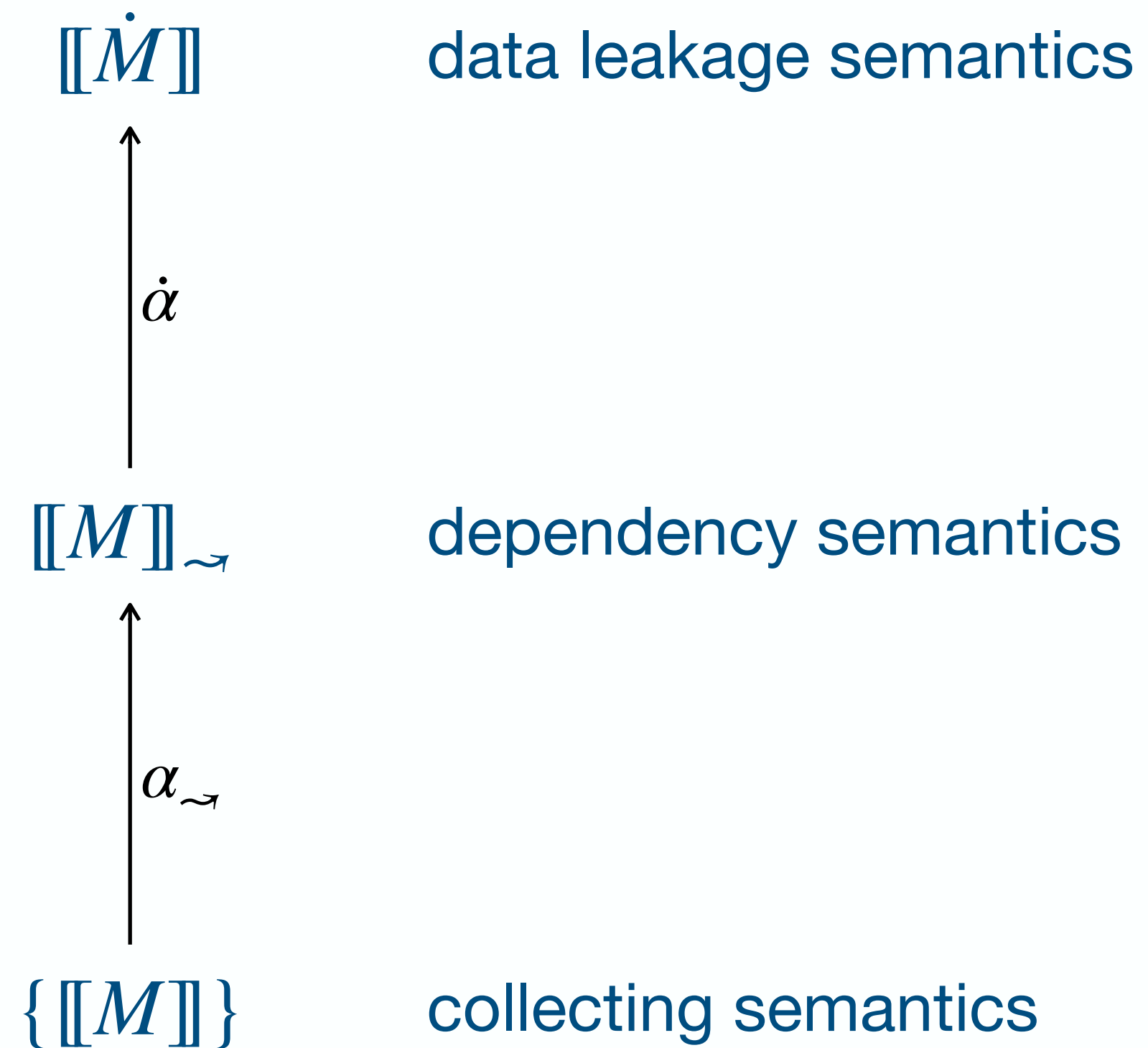


concrete semantics

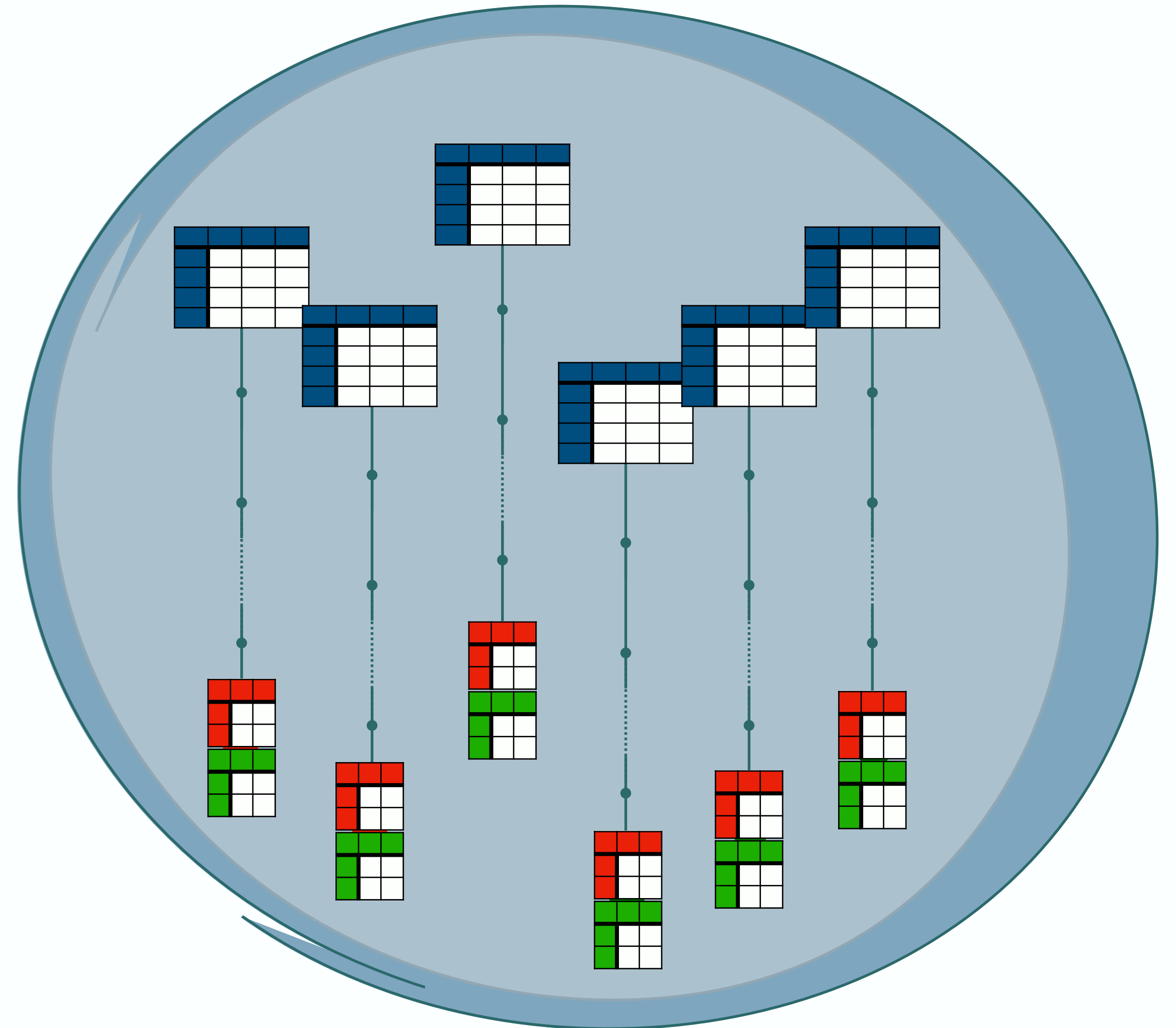
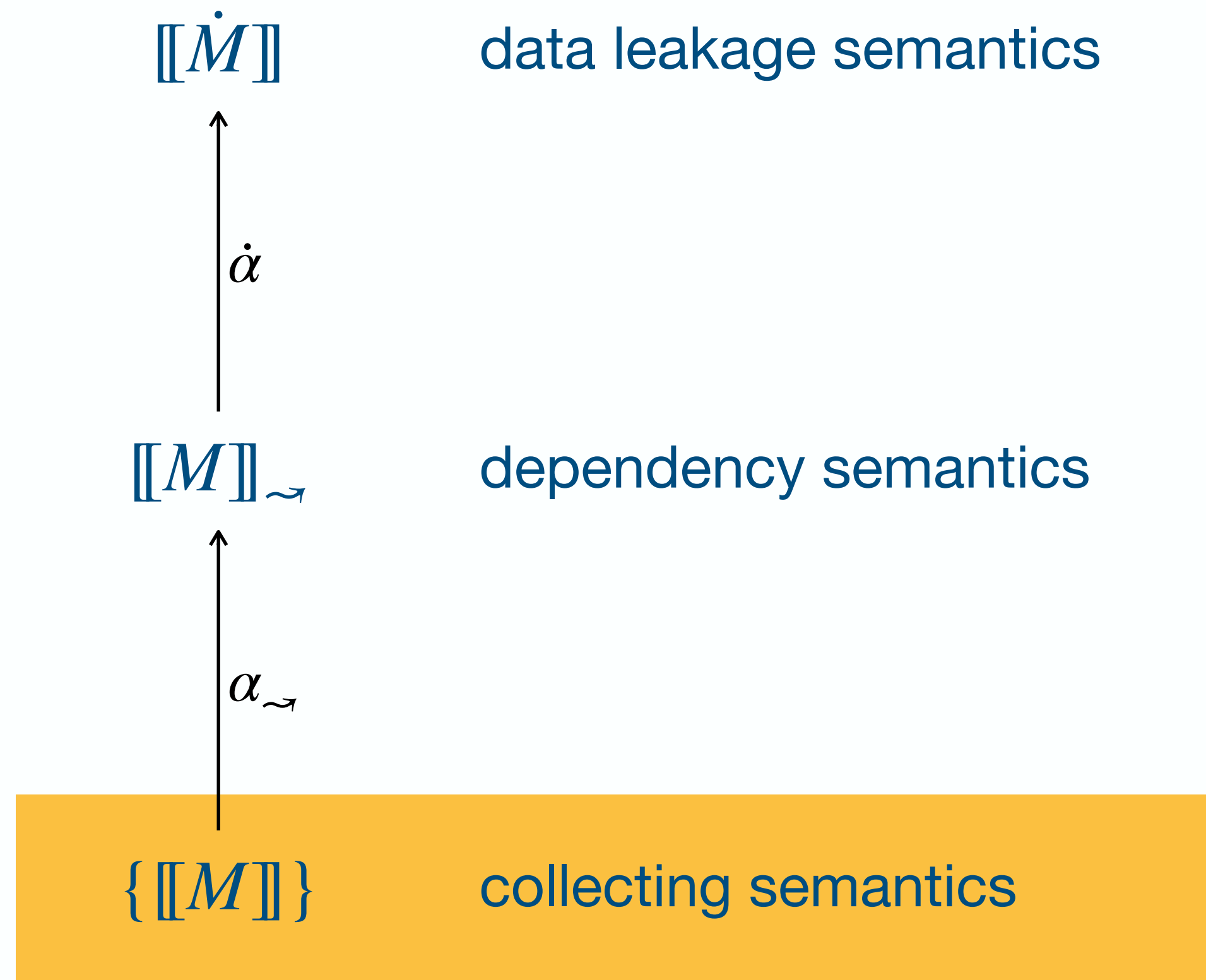
mathematical models of the program behavior



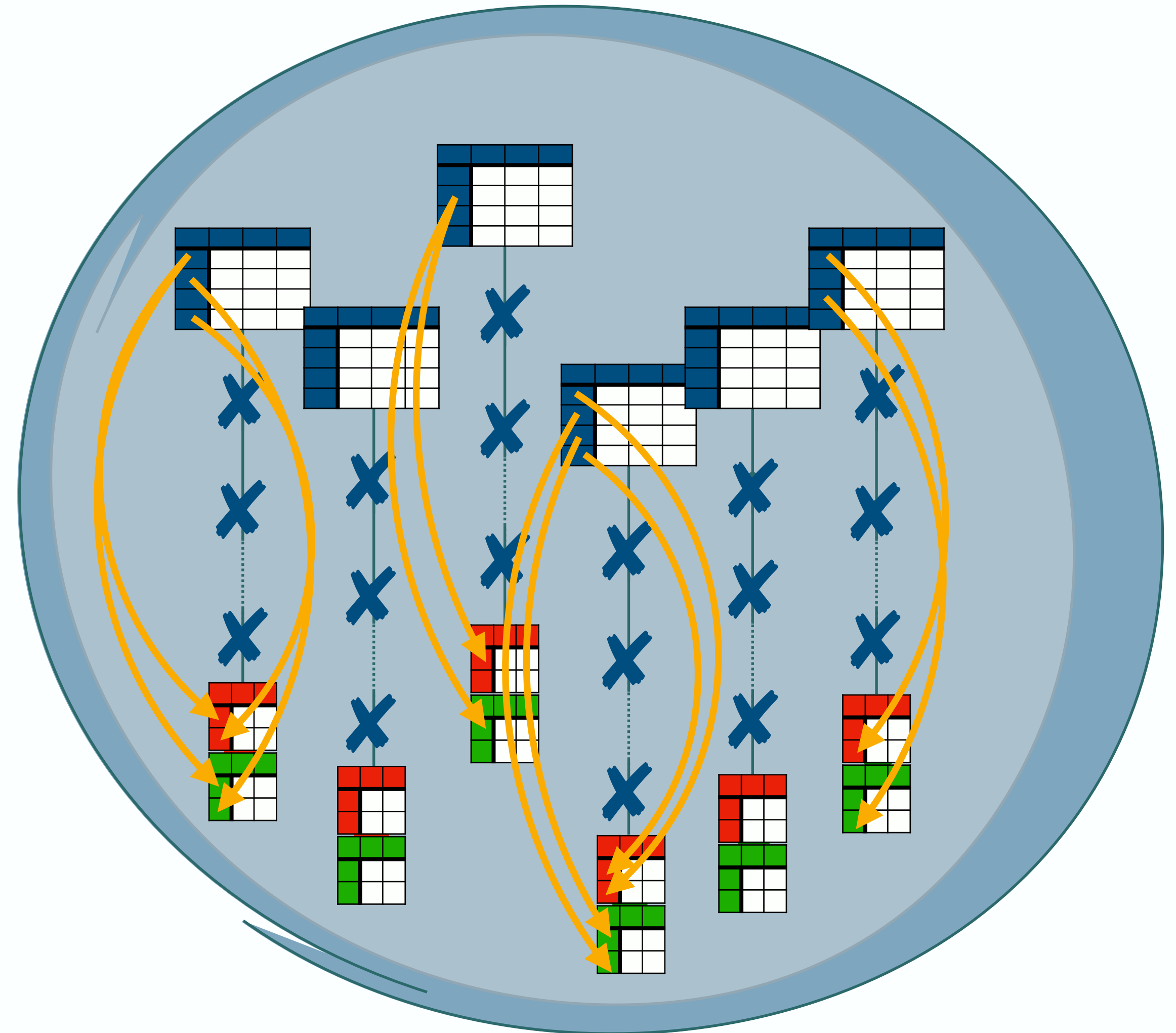
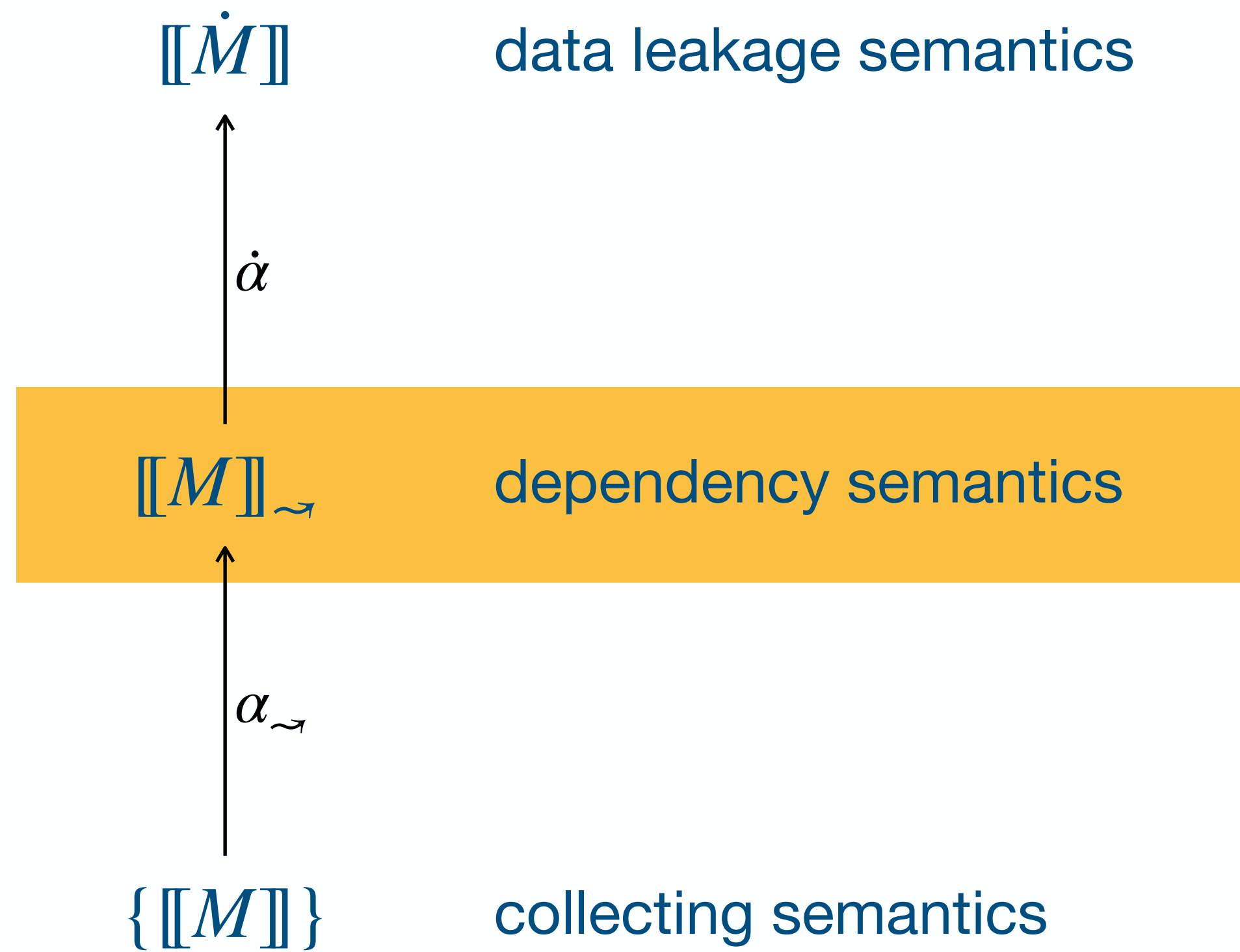
Hierarchy of Semantics



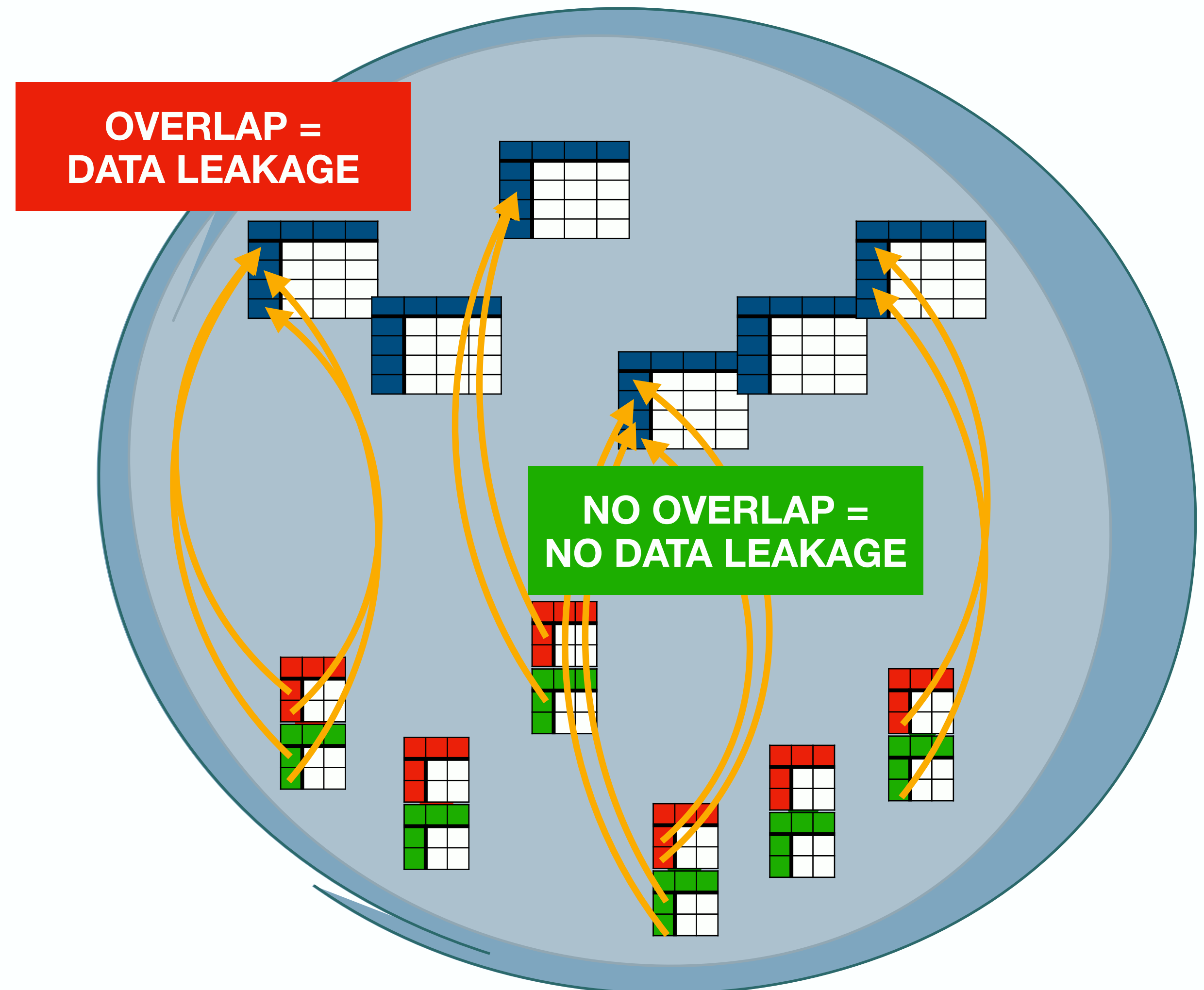
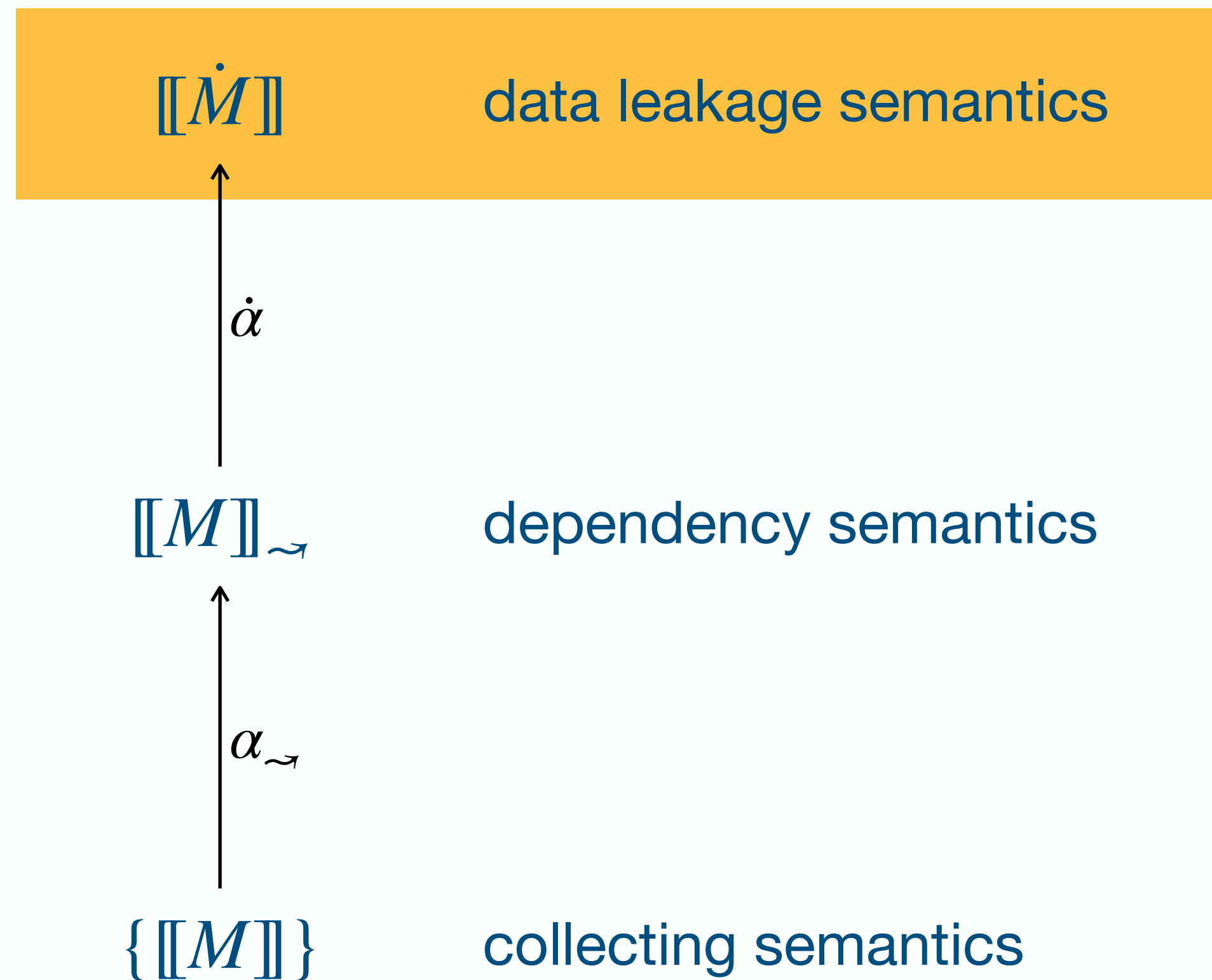
Hierarchy of Semantics



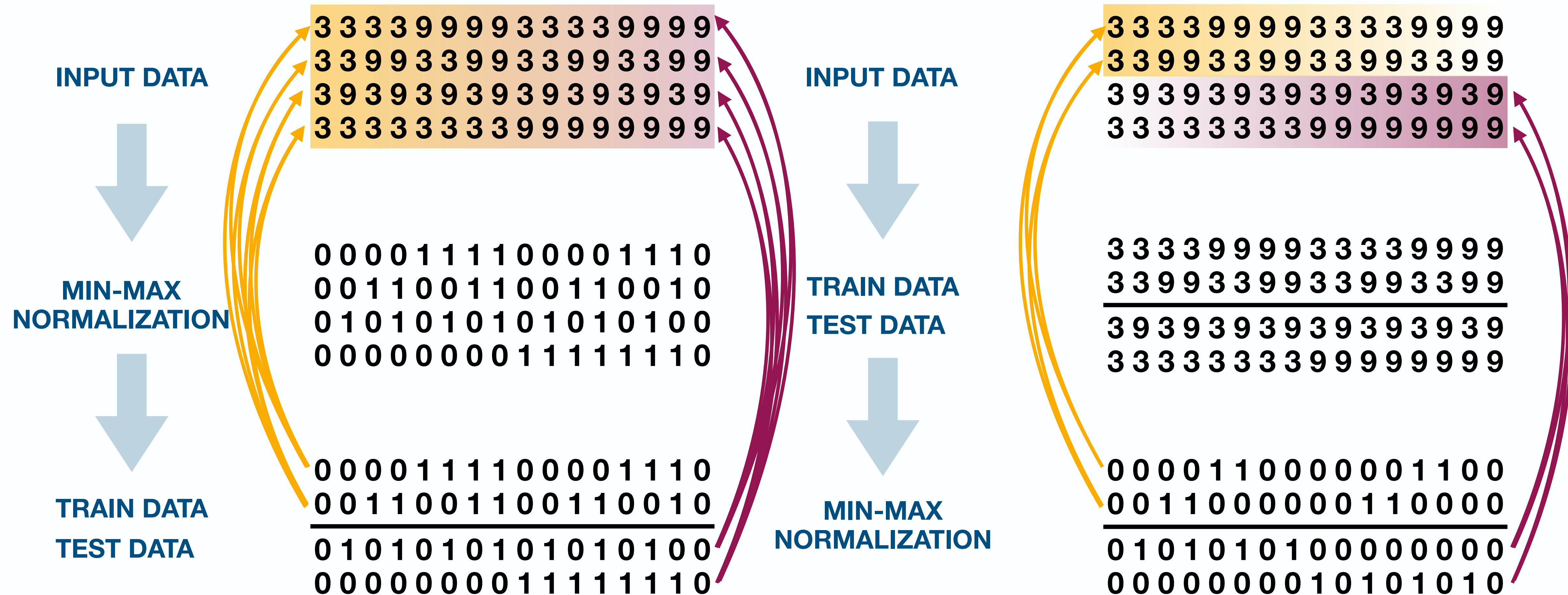
Hierarchy of Semantics



Hierarchy of Semantics



Data Leakage Semantics



Data Leakage Static Analysis [Drobnjaković24]

3-Step Recipe

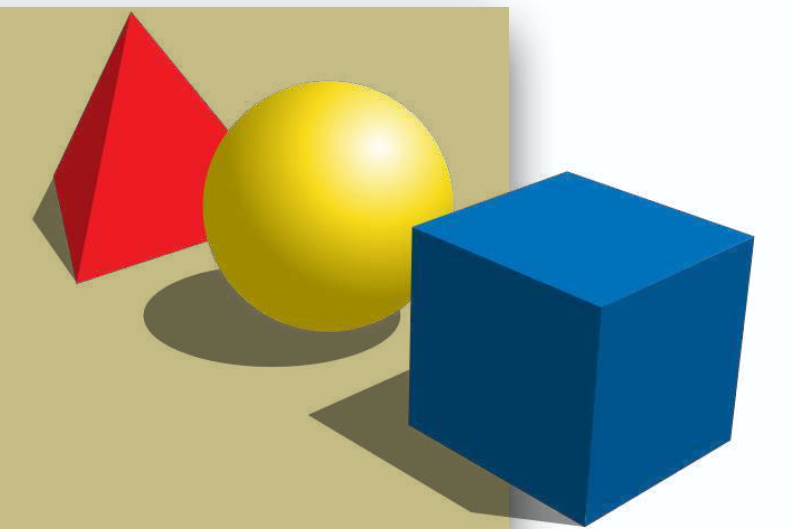
practical tools

targeting specific programs



abstract semantics, abstract domains

algorithmic approaches to decide program properties

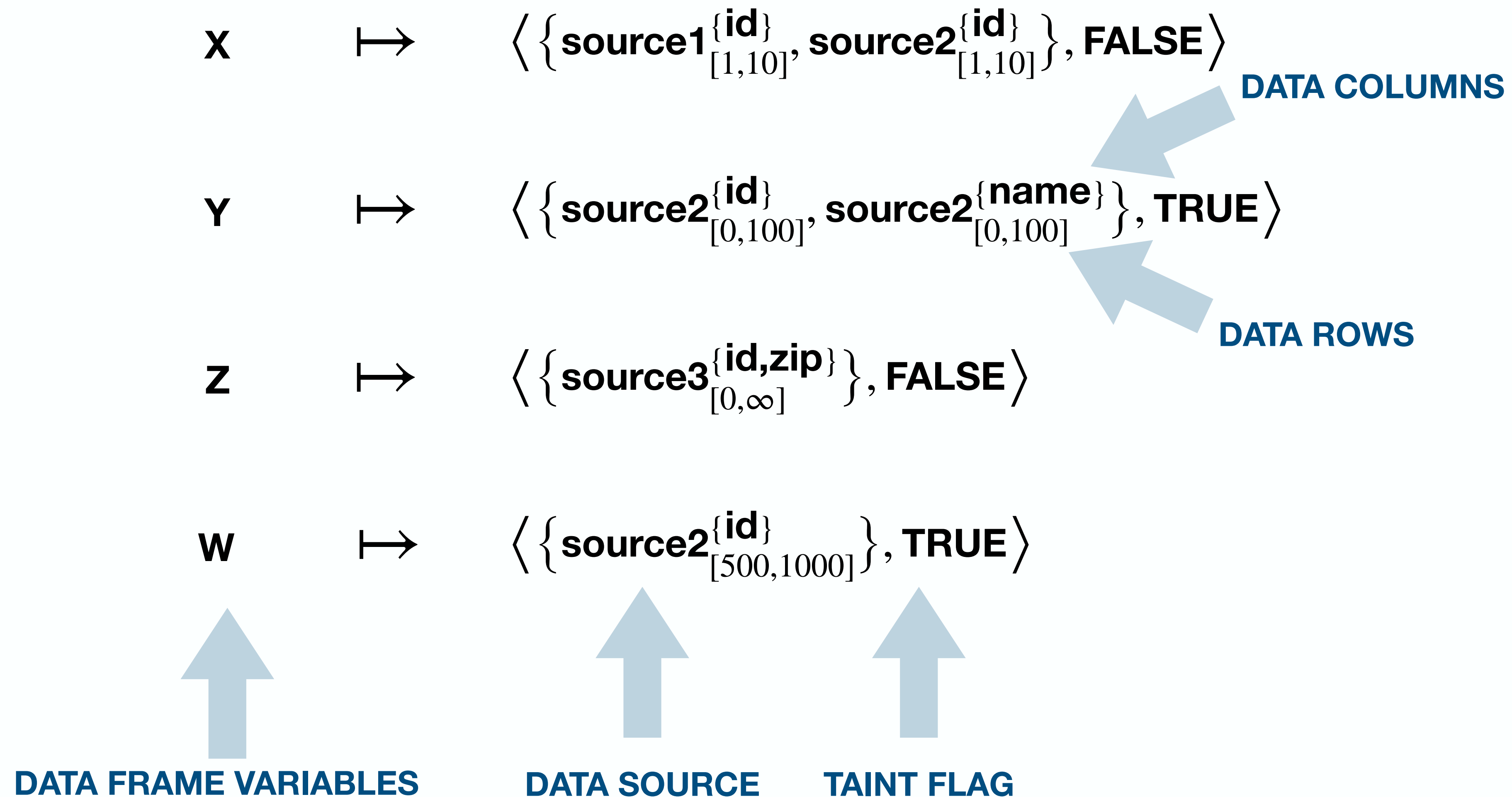


concrete semantics

mathematical models of the program behavior



Data Frame Sources Abstract Domain



Example

```
[1]: import numpy as np
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

[2]: data = pd.read_csv("data.csv")
X = data[["X_1", "X_2"]]
y = data[["y"]]

[3]: min_max_scaler = MinMaxScaler()
X = min_max_scaler.fit_transform(X)

[4]: X_train , X_test , y_train , y_test = train_test_split(X, y, test_size=0.025, random_state=2)

[ ]: X_train ↦ ⟨ { data_{X_1, X_2}^{[0, ∞]} }, TRUE ⟩
X_test ↦ ⟨ { data_{X_1, X_2}^{[0, ∞]} }, TRUE ⟩

[5]: lr = LogisticRegression()
a = lr.fit(X_train, y_train)

[6]: y_pred = lr.predict(X_test)
accuracy_score(y_test, y_pred)

[6]: 0.67
```

INPUT DATA READING $X \mapsto \langle \{ \text{data}_{[0, \infty]}^{\{X_1, X_2\}} \}, \text{FALSE} \rangle$

MIN-MAX NORMALIZATION $X \mapsto \langle \{ \text{data}_{[0, \infty]}^{\{X_1, X_2\}} \}, \text{TRUE} \rangle$

TRAIN/TEST SPLIT

TRAINING

TESTING

(TAINT) OVERLAP = DATA LEAKAGE

Example

```
[1]: import numpy as np
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

INPUT DATA READING $X \mapsto \langle \{ \text{data}_{[0, \infty]}^{X_1, X_2} \}, \text{FALSE} \rangle$

```
[2]: data = pd.read_csv("data.csv")
X = data[["X_1", "X_2"]]
y = data[["y"]]
```

$X_{\text{train}} \mapsto \langle \{ \text{data}_{[0.025 \cdot R + 1, \infty]}^{X_1, X_2} \}, \text{FALSE} \rangle$ $X_{\text{test}} \mapsto \langle \{ \text{data}_{[0, 0.025 \cdot R]}^{X_1, X_2} \}, \text{FALSE} \rangle$

```
[3]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.025, random_state=2)
```

TRAIN/TEST SPLIT

```
[4]: min_max_scaler = MinMaxScaler()
X_train = min_max_scaler.fit_transform(X_train)
X_test = min_max_scaler.fit_transform(X_test)
```

MIN-MAX NORMALIZATION $X_{\text{train}} \mapsto \langle \{ \text{data}_{[0.025 \cdot R + 1, \infty]}^{X_1, X_2} \}, \text{TRUE} \rangle$
 $X_{\text{test}} \mapsto \langle \{ \text{data}_{[0, 0.025 \cdot R]}^{X_1, X_2} \}, \text{TRUE} \rangle$

```
[5]: lr = LogisticRegression()
a = lr.fit(X_train, y_train)
```

TRAINING

```
[6]: y_pred = lr.predict(X_test)
accuracy_score(y_test, y_pred)
```

TESTING

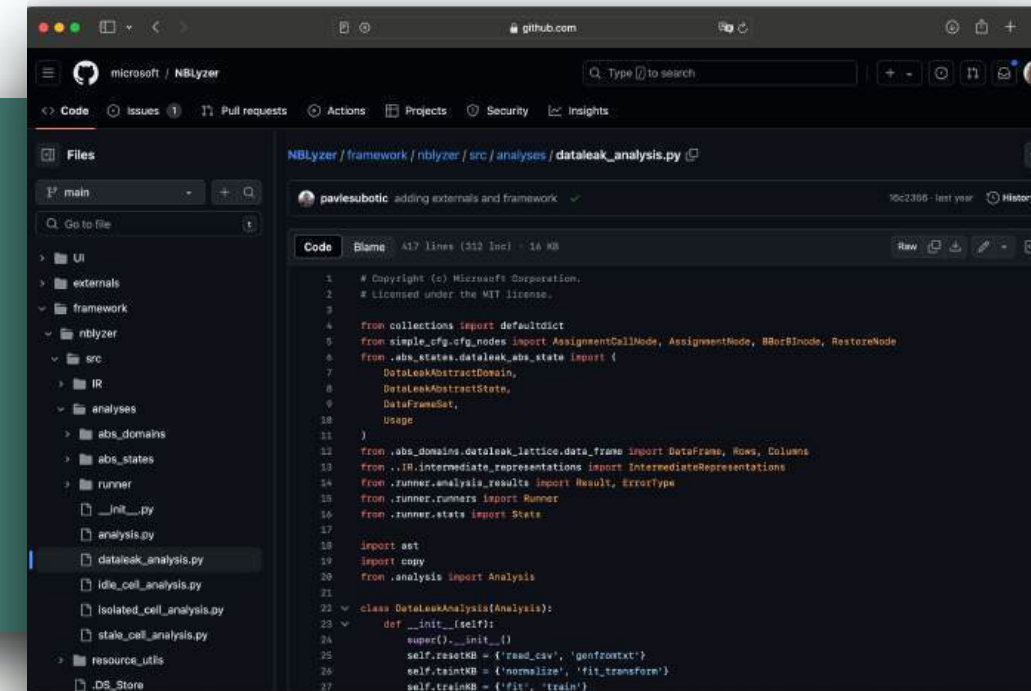
NO OVERLAP = NO DATA LEAKAGE

```
[6]: 0.33
```

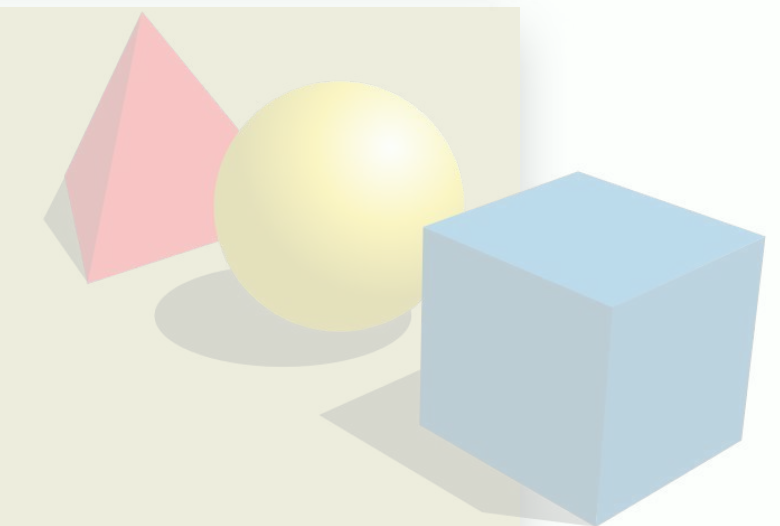

Data Leakage Static Analysis [Drobnjaković24]

3-Step Recipe

practical tools
targeting specific programs



abstract semantics, abstract domains
algorithmic approaches to decide program properties



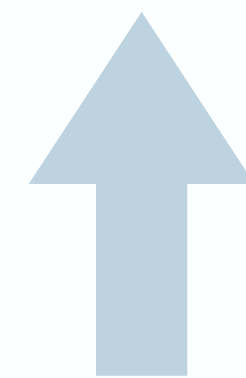
concrete semantics
mathematical models of the program behavior



Experimental Evaluation

7378 Executions in 2111 Notebooks from Kaggle

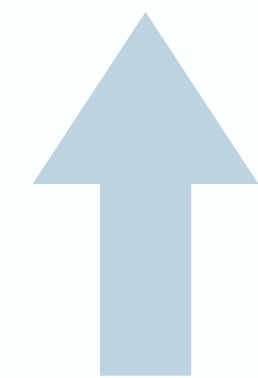
Implementation	True Positives		False Positives
	Taint Data Leakage	Overlap Data Leakage	
NBLyzer + Original Data Leakage Analysis	10	0	2
NBLyzer + Our Data Leakage Analysis	10	15	2



IN 5 NOTEBOOKS



IN 11 NOTEBOOKS

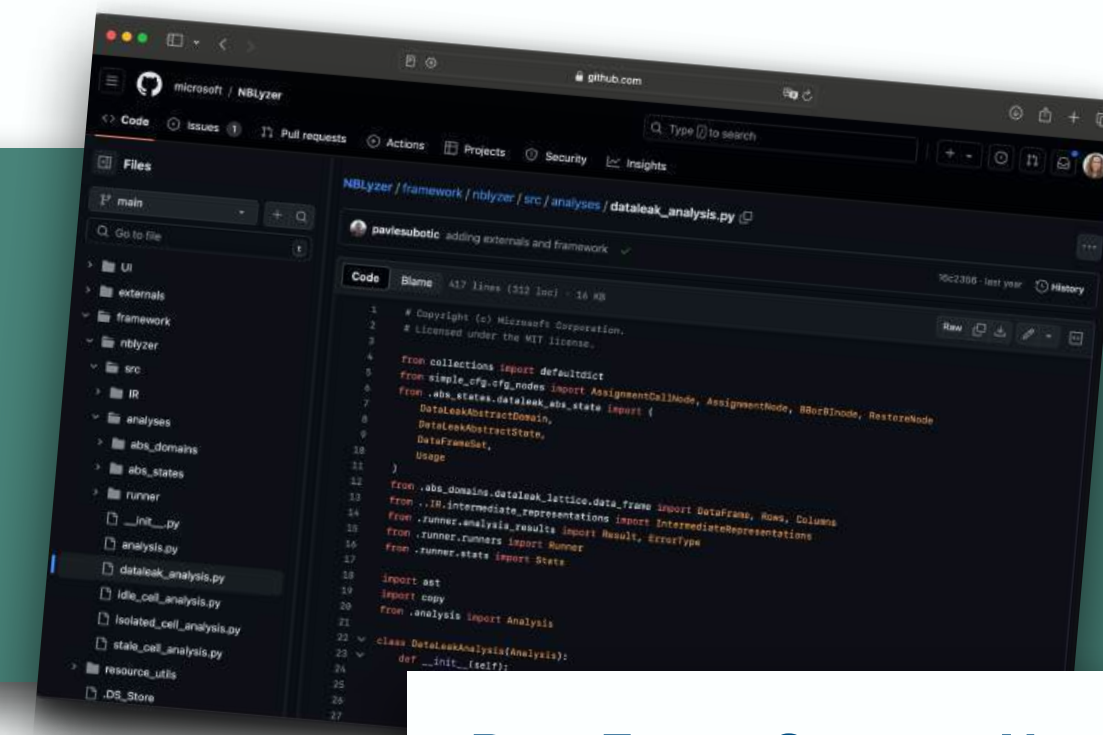


**CONFIRMED BY
4 DATA SCIENTISTS
AT MICROSOFT**

Data Leakage Static Analysis [Drobnjaković24]

3-Step Recipe

practical tools
targeting specific programs

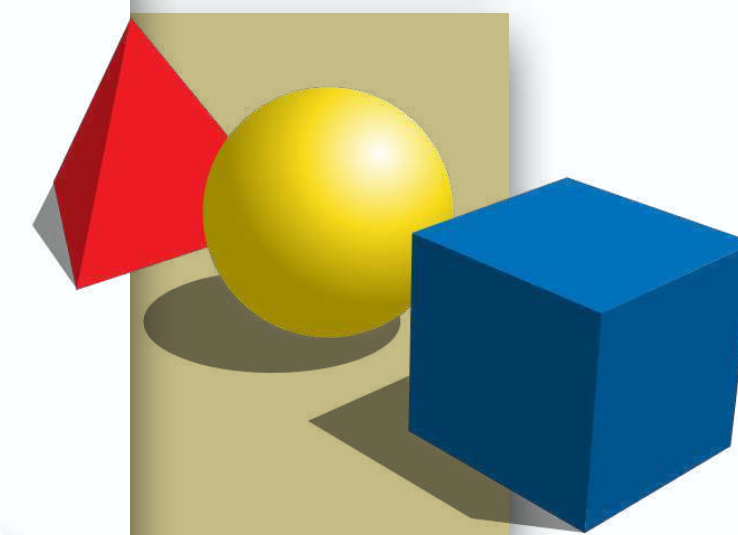


abstract semantics, abstract domains
algorithmic approaches to decide program properties

Data Frame Sources Abstract Domain

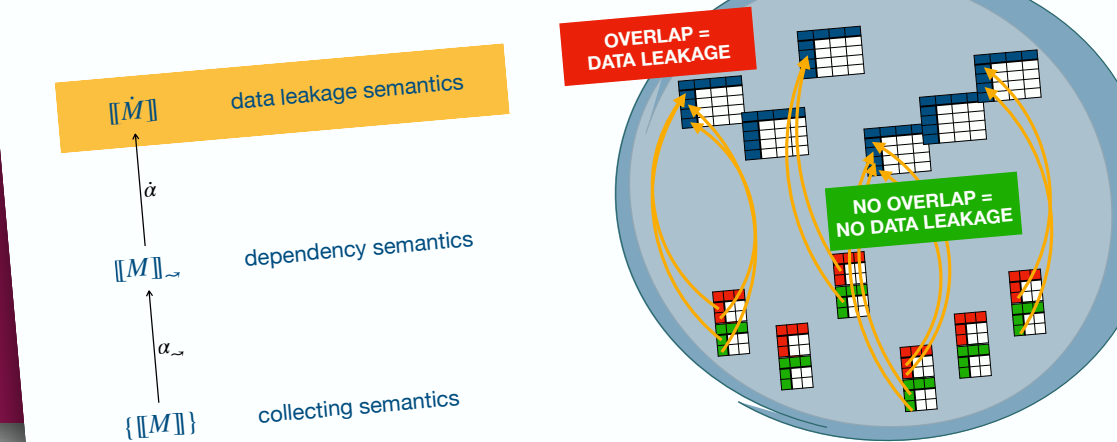
X	↦	$\langle\langle \text{source1}_{[1,10]}, \text{source2}_{[1,10]} \rangle\rangle, \text{FALSE}\rangle$	DATA COLUMNS
Y	↦	$\langle\langle \text{source2}_{[0,100]}, \text{source2}_{[0,100]} \rangle\rangle, \text{TRUE}\rangle$	DATA ROWS
Z	↦	$\langle\langle \text{source3}_{[0,\infty]} \rangle\rangle, \text{FALSE}\rangle$	
W	↦	$\langle\langle \text{source2}_{[500,1000]} \rangle\rangle, \text{TRUE}\rangle$	

↑ DATA FRAME VARIABLES ↑ DATA SOURCE ↑ TAINT FLAG



concrete semantics
mathematical models of the program behavior

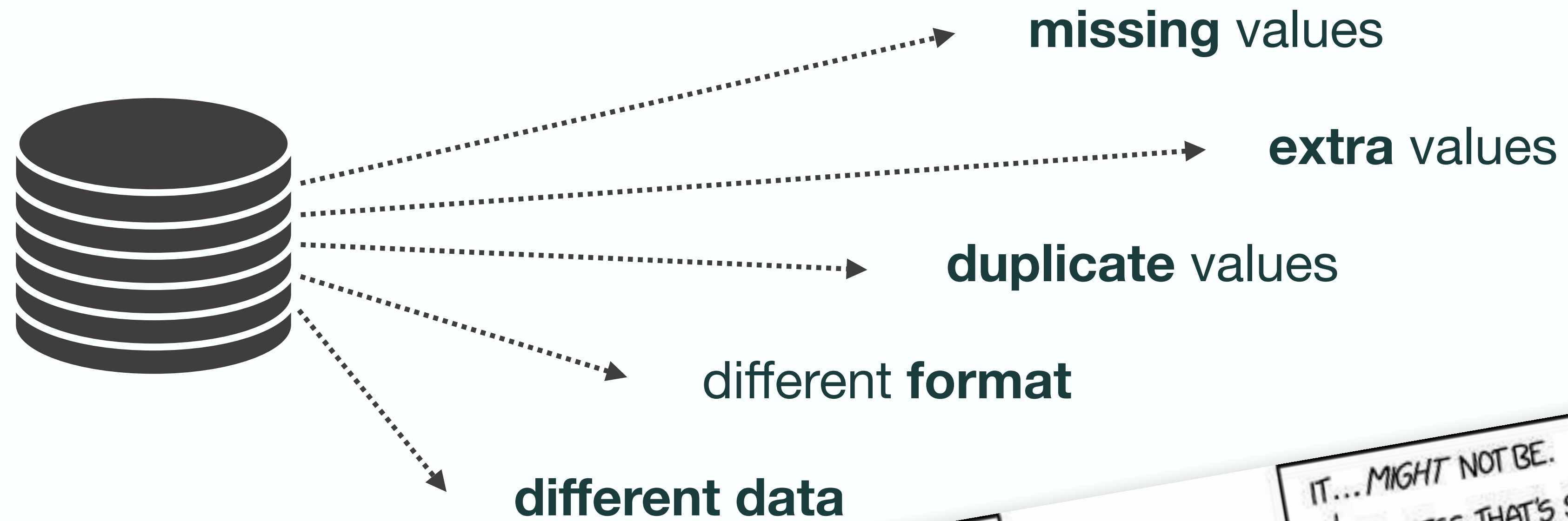
Hierarchy of Semantics



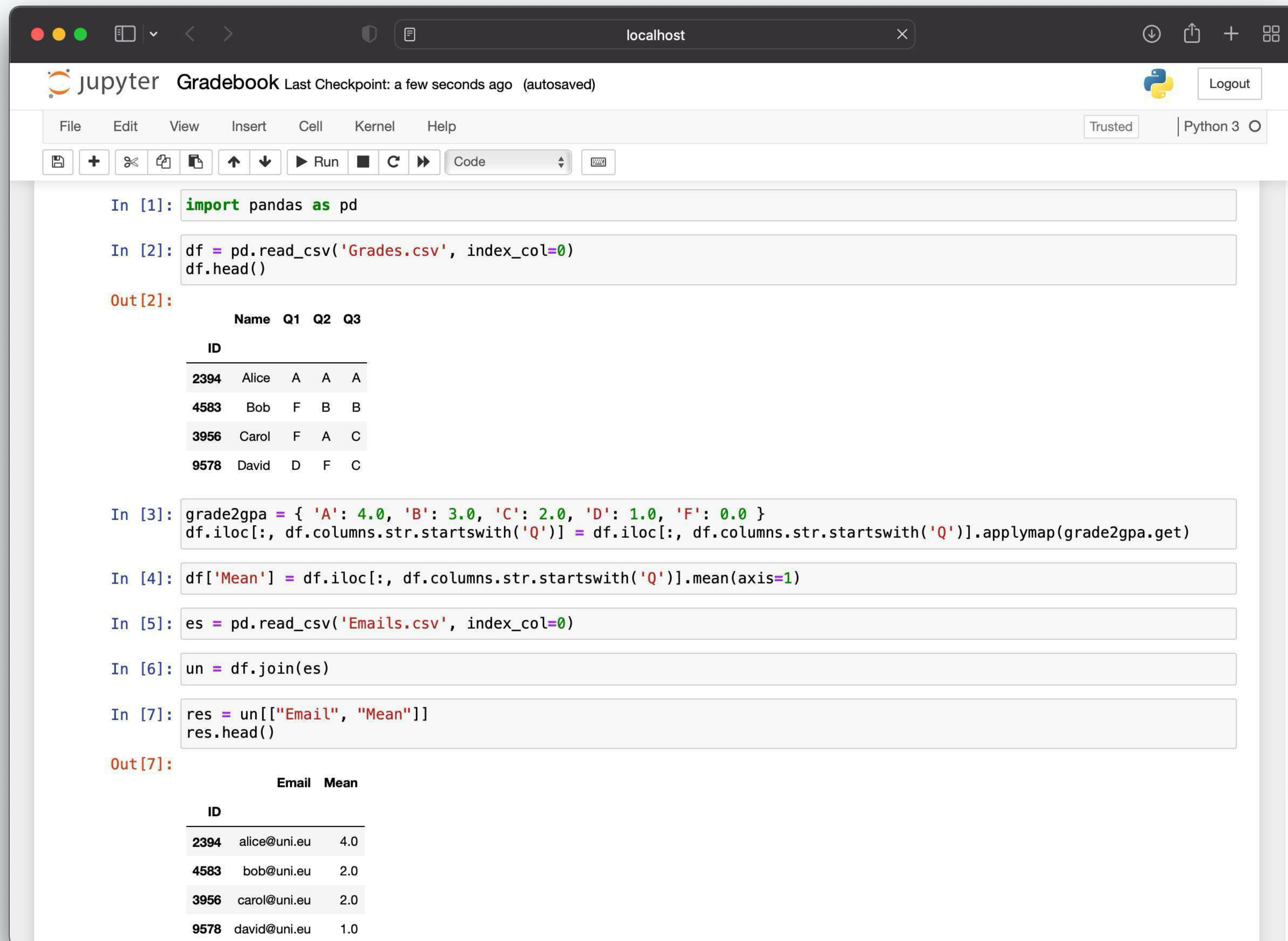


Unexpected Data

Unexpected Data



Example



The screenshot shows a Jupyter Notebook interface with the following content:

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
```

Out[2]:

ID	Name	Q1	Q2	Q3
2394	Alice	A	A	A
4583	Bob	F	B	B
3956	Carol	F	A	C
9578	David	D	F	C

```
In [3]: grade2gpa = { 'A': 4.0, 'B': 3.0, 'C': 2.0, 'D': 1.0, 'F': 0.0 }
df.iloc[:, df.columns.str.startswith('Q')] = df.iloc[:, df.columns.str.startswith('Q')].applymap(grade2gpa.get)
```

```
In [4]: df['Mean'] = df.iloc[:, df.columns.str.startswith('Q')].mean(axis=1)
```

```
In [5]: es = pd.read_csv('Emails.csv', index_col=0)
```

```
In [6]: un = df.join(es)
```

```
In [7]: res = un[["Email", "Mean"]]
res.head()
```

Out[7]:

ID	Email	Mean
2394	alice@uni.eu	4.0
4583	bob@uni.eu	2.0
3956	carol@uni.eu	2.0
9578	david@uni.eu	1.0

Missing Values

jupyter Gradebook Last Checkpoint: a few seconds ago (auto)

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
```

Out[2]:

ID	Name	Q1	Q2	Q3
2394	Alice	A	A	A
4583	Bob	F	B	B
3956	Carol	F	A	C
9578	David	D	F	C

```
In [3]: grade2gpa = { 'A': 4.0, 'B': 3.0, 'C': 2.0 }
df.iloc[:, df.columns.str.startswith('Q')]
```

```
In [4]: df['Mean'] = df.iloc[:, df.columns.str.startswith('Q')].mean(axis=1)
```

```
In [5]: es = pd.read_csv('Emails.csv', index_col=0)
```

```
In [6]: un = df.join(es)
```

```
In [7]: res = un[["Email", "Mean"]]
res.head()
```

Out[7]:

ID	Email	Mean
2394	alice@uni.eu	4.0
4583	bob@uni.eu	2.0
3956	carol@uni.eu	2.0
9578	david@uni.eu	1.0

localhost

jupyter Gradebook Last Checkpoint: a minute ago (unsaved changes)

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
```

Out[2]:

ID	Name	Q1	Q2	Q3
2394	Alice	A	A	A
4583	Bob	F	B	B
3956	Carol	NaN	A	C
9578	David	D	F	C

```
In [3]: grade2gpa = { 'A': 4.0, 'B': 3.0, 'C': 2.0, 'D': 1.0, 'F': 0.0 }
df.iloc[:, df.columns.str.startswith('Q')] = df.iloc[:, df.columns.str.startswith('Q')].applymap(grade2gpa.get)
```

```
In [4]: df['Mean'] = df.iloc[:, df.columns.str.startswith('Q')].mean(axis=1)
```

```
In [5]: es = pd.read_csv('Emails.csv', index_col=0)
```

```
In [6]: un = df.join(es)
```

```
In [7]: res = un[["Email", "Mean"]]
res.head()
```

Out[7]:

ID	Email	Mean
2394	alice@uni.eu	4.0
4583	bob@uni.eu	2.0
3956	carol@uni.eu	3.0
9578	david@uni.eu	1.0

Extra Values

jupyter Gradebook Last Checkpoint: a few seconds ago (auto)

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
```

Out[2]:

ID	Name	Q1	Q2	Q3
2394	Alice	A	A	A
4583	Bob	F	B	B
3956	Carol	F	A	C
9578	David	D	F	C

```
In [3]: grade2gpa = { 'A': 4.0, 'B': 3.0, 'C': 2.0 }
df.iloc[:, df.columns.str.startswith('Q')]
```

```
In [4]: df['Mean'] = df.iloc[:, df.columns.str.startswith('Q')].mean(axis=1)
```

```
In [5]: es = pd.read_csv('Emails.csv', index_col=0)
```

```
In [6]: un = df.join(es)
```

```
In [7]: res = un[["Email", "Mean"]]
res.head()
```

Out[7]:

ID	Email	Mean
2394	alice@uni.eu	4.0
4583	bob@uni.eu	2.0
3956	carol@uni.eu	2.0
9578	david@uni.eu	1.0

localhost

jupyter Gradebook Last Checkpoint: 2 minutes ago (unsaved changes)

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
```

Out[2]:

ID	Name	Q1	Q2	Q3
2394	Alice	A	A	A
4583	Bob	F	B	NaN
3956	Carol	F	A	NaN
9578	David	D	F	NaN

```
In [3]: grade2gpa = { 'A': 4.0, 'B': 3.0, 'C': 2.0, 'D': 1.0, 'F': 0.0 }
df.iloc[:, df.columns.str.startswith('Q')] = df.iloc[:, df.columns.str.startswith('Q')].applymap(grade2gpa.get)
```

```
In [4]: df['Mean'] = df.iloc[:, df.columns.str.startswith('Q')].mean(axis=1)
```

```
In [5]: es = pd.read_csv('Emails.csv', index_col=0)
```

```
In [6]: un = df.join(es)
```

```
In [7]: res = un[["Email", "Mean"]]
res.head()
```

Out[7]:

ID	Email	Mean
2394	alice@uni.eu	4.0
4583	bob@uni.eu	3.0
3956	carol@uni.eu	3.0
9578	david@uni.eu	1.0

Different Format

The image displays two Jupyter Notebook windows side-by-side, both titled "Gradebook".

Left Window (Last Checkpoint: a few seconds ago):

```
In [1]: import pandas as pd
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
Out[2]:
```

ID	Name	Q1	Q2	Q3
2394	Alice	A	A	A
4583	Bob	F	B	B
3956	Carol	F	A	C
9578	David	D	F	C

```
In [3]: grade2gpa = { 'A': 4.0, 'B': 3.0, 'C': 2.0 }
df.iloc[:, df.columns.str.startswith('Q')]
In [4]: df['Mean'] = df.iloc[:, df.columns.str.startswith('Q')].mean(axis=1)
In [5]: es = pd.read_csv('Emails.csv', index_col=0)
In [6]: un = df.join(es)
In [7]: res = un[["Email", "Mean"]]
res.head()
Out[7]:
```

ID	Email	Mean
2394	alice@uni.eu	4.0
4583	bob@uni.eu	2.0
3956	carol@uni.eu	2.0
9578	david@uni.eu	1.0

Right Window (Last Checkpoint: 14 minutes ago, unsaved changes):

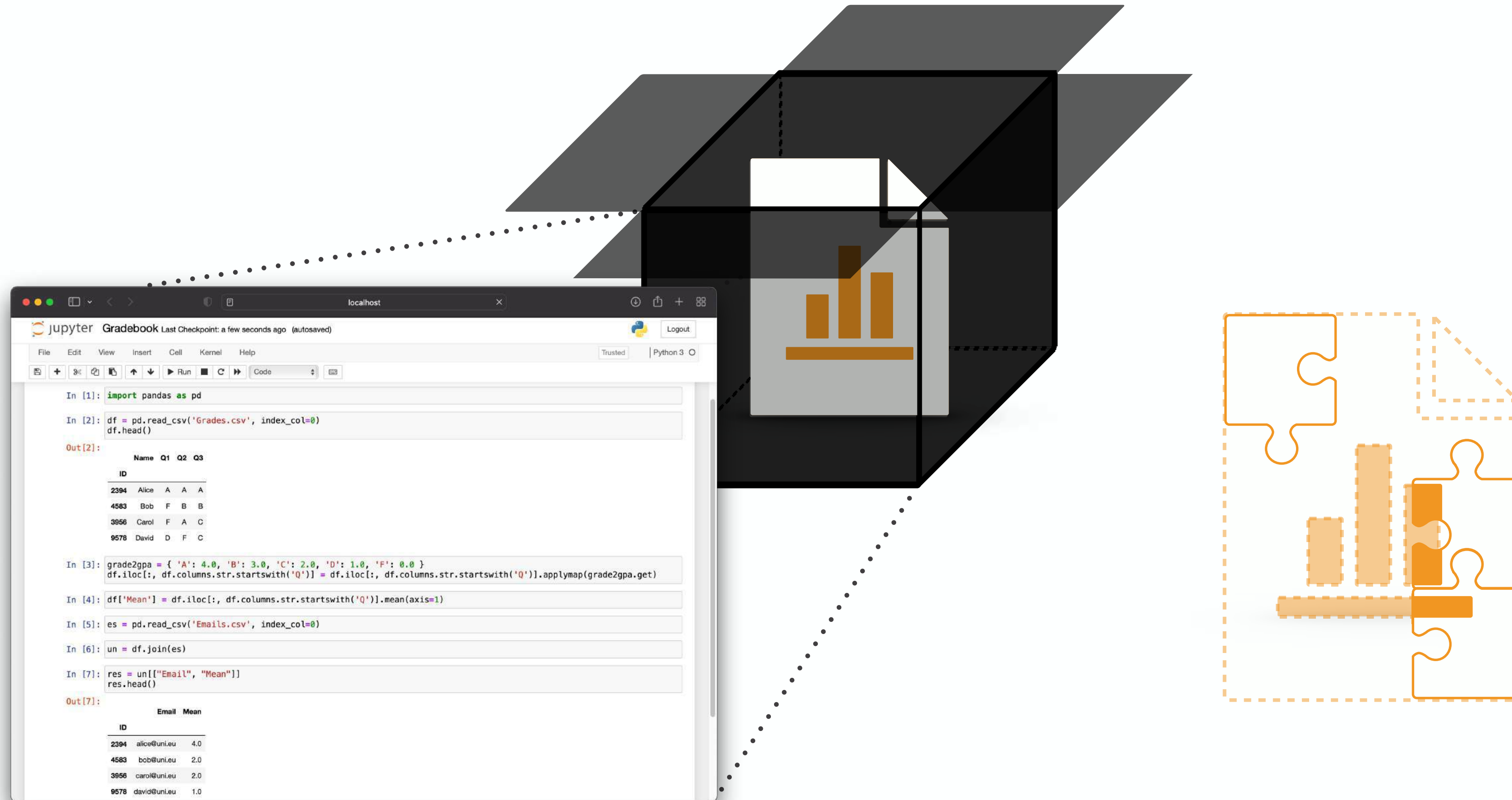
```
In [1]: import pandas as pd
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
Out[2]:
```

ID	Name	Q1	Q2	Q3
2394	Alice	A	A	A
4583	Bob	F	B+	B
3956	Carol	F	A	C
9578	David	D	F	C

```
In [3]: grade2gpa = { 'A': 4.0, 'B': 3.0, 'C': 2.0, 'D': 1.0, 'F': 0.0 }
df.iloc[:, df.columns.str.startswith('Q')] = df.iloc[:, df.columns.str.startswith('Q')].applymap(grade2gpa.get)
In [4]: df['Mean'] = df.iloc[:, df.columns.str.startswith('Q')].mean(axis=1)
In [5]: es = pd.read_csv('Emails.csv', index_col=0)
In [6]: un = df.join(es)
In [7]: res = un[["Email", "Mean"]]
res.head()
Out[7]:
```

ID	Email	Mean
2394	alice@uni.eu	4.0
4583	bob@uni.eu	1.5
3956	carol@uni.eu	2.0
9578	david@uni.eu	1.0

Data Expectations Static Analysis



Data

1st Challenge: Multi-Dimensional Data Structures



Data Expectations Static Analysis

3-Step Recipe

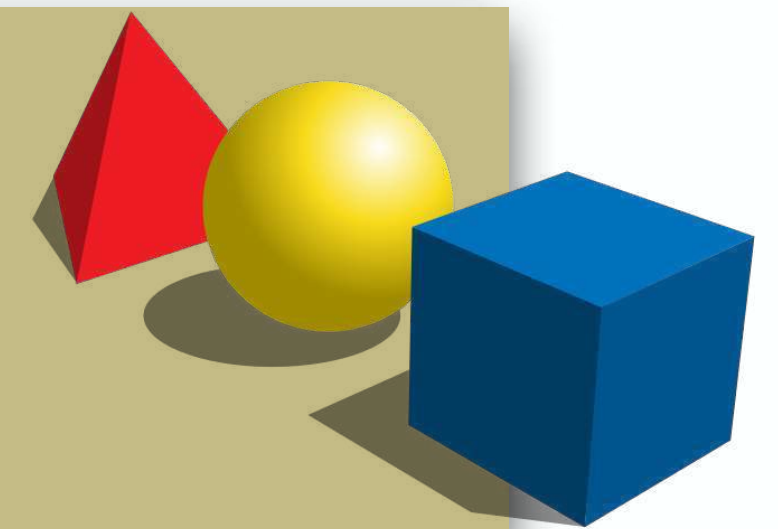
practical tools

targeting specific programs



abstract semantics, abstract domains

algorithmic approaches to decide program properties



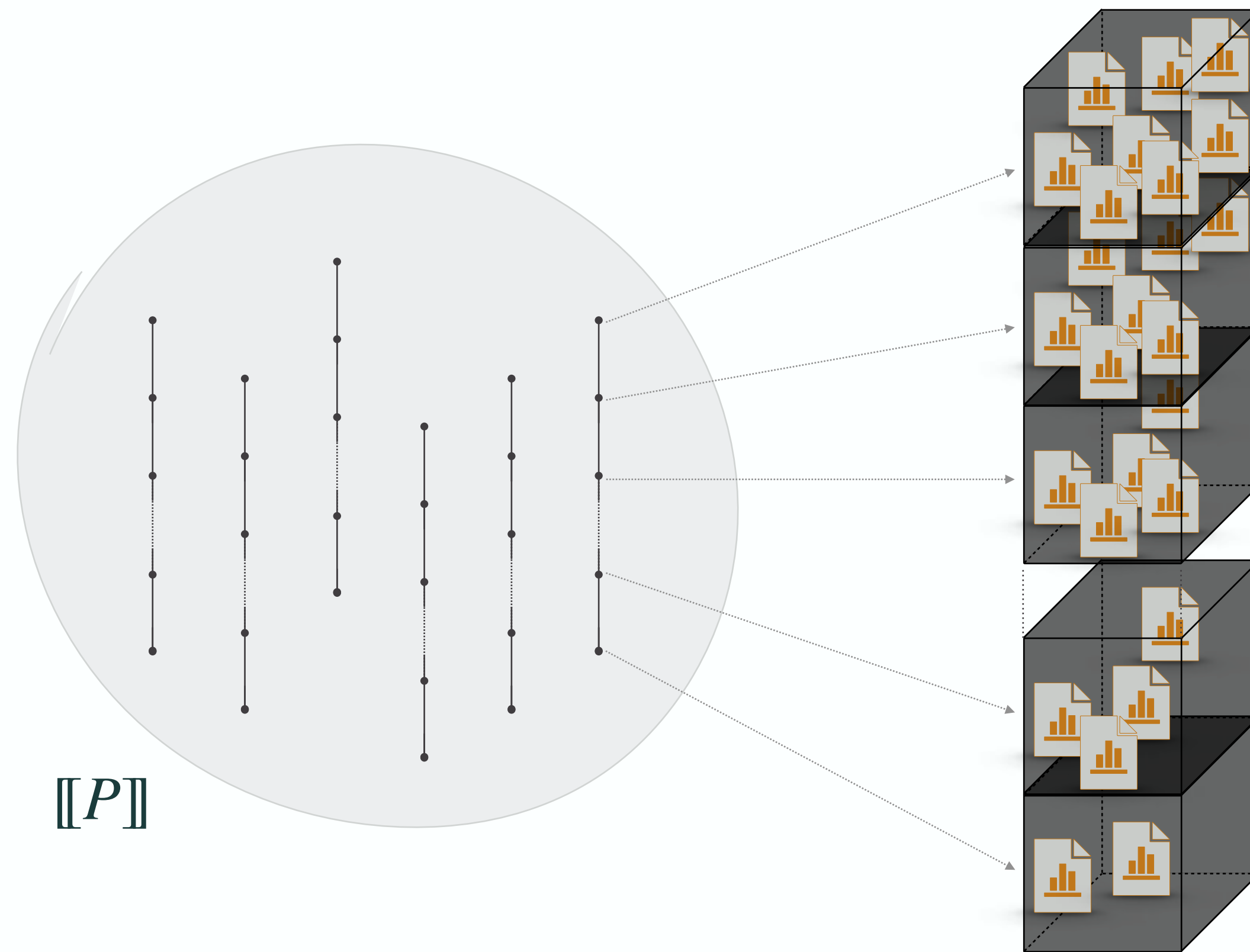
concrete semantics

mathematical models of the program behavior



Concrete Semantics

2nd Challenge: Indirect Reasoning



Abstract Semantics

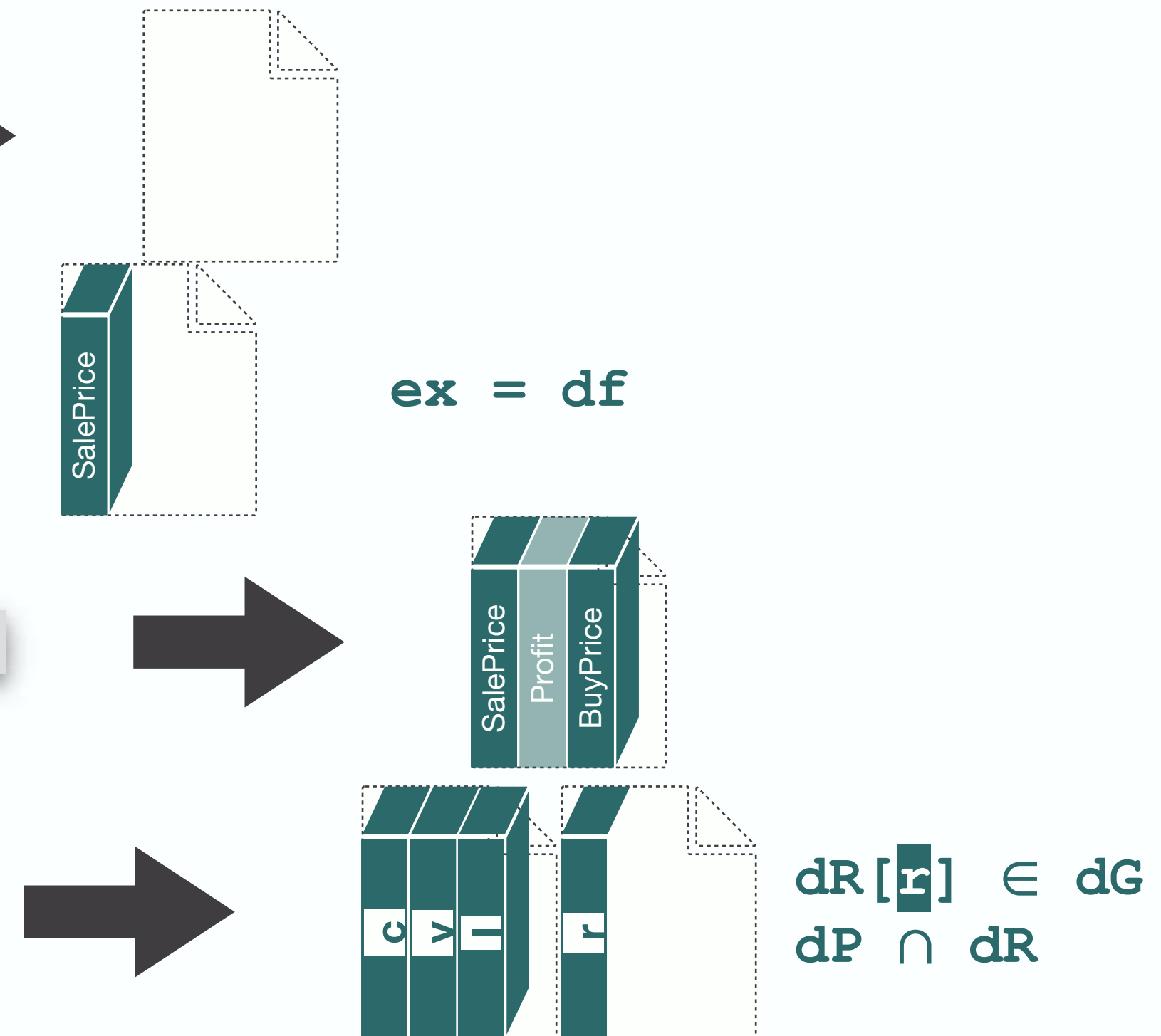
3rd Challenge: Complex Library Calls

```
import pandas as pd  
df = pd.read_csv("HousePrices.csv")
```

```
ex = df[df.SalePrice >= 1000000]
```

```
ex['Profit'] = ex['SalePrice'] - ex['BuyPrice']
```

```
:  
dL = pd.read_csv("L.csv")  
dP = dL.pivot(index=c, columns=y, values=l)  
dR = pd.read_csv("R.csv")  
dG = dP.loc[:, 0:35].groupby(dR[r])
```



Implementation

Wish List

The screenshot shows a Jupyter Notebook window titled "House Prices (unsaved changes)". The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help) and a toolbar with icons for file operations and execution. The notebook content shows three input cells:

- In [5]:** `sns.distplot(df_train['SalePrice']);` This cell is followed by a histogram plot of "SalePrice" with a kernel density estimate curve overlaid. The x-axis ranges from 0 to 800,000, and the y-axis ranges from 0.000000 to 0.000008.
- In [14]:** `total = df_train.isnull().sum().sort_values(ascending=False)
percent = (df_train.isnull().sum()/df_train.isnull().count()).sort_values(ascending=False)
missing_data = pd.concat([total, per`
- In [15]:** `df_train = df_train.drop((missing_da
df_train = df_train.drop(df_train.lo
df_train.isnull().sum().max()`


The output of the last cell is `Out[15]: 0`. Orange arrows point from the code cells to the corresponding callout boxes.

MULTI-LANGUAGE SUPPORT

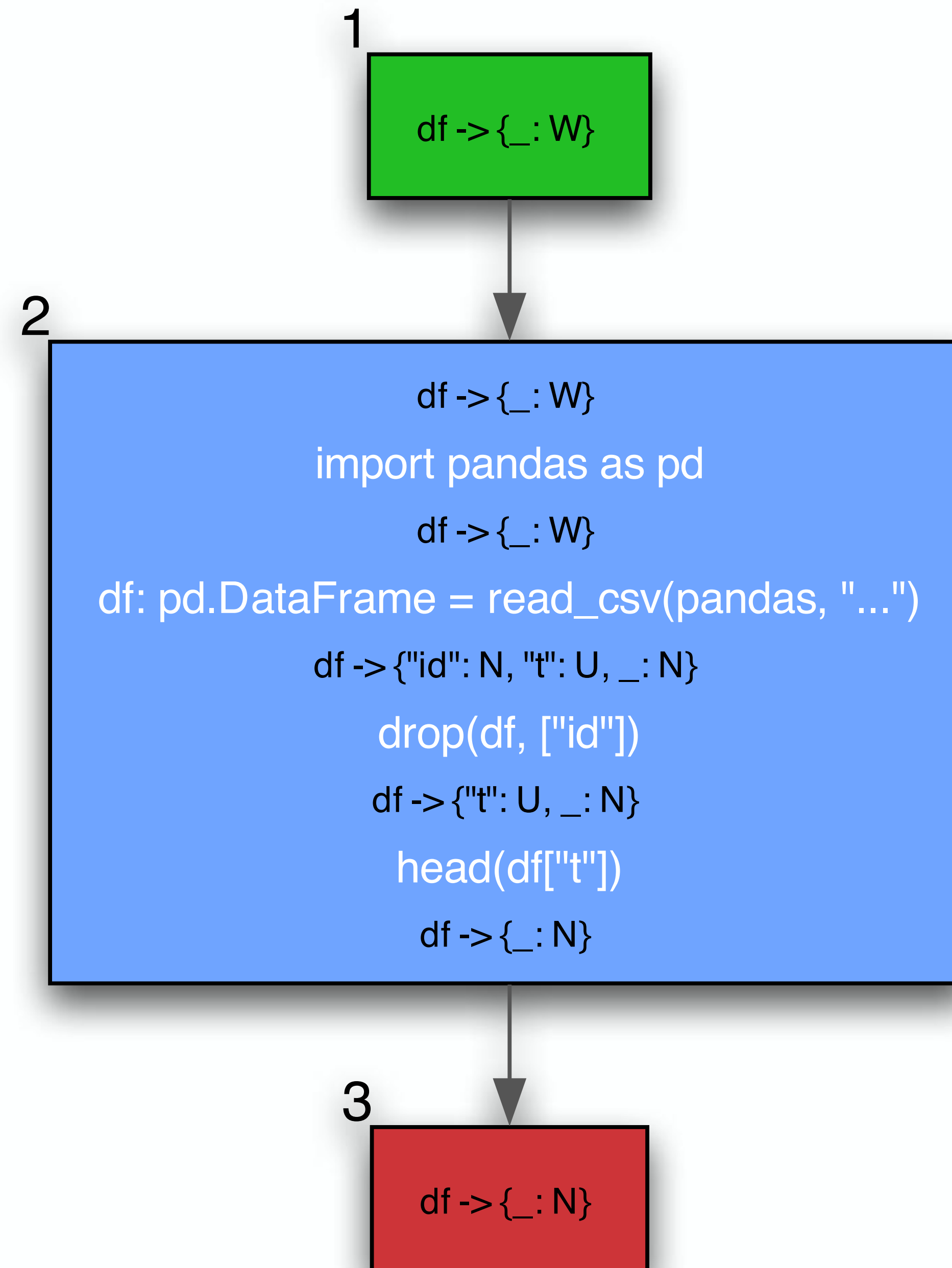


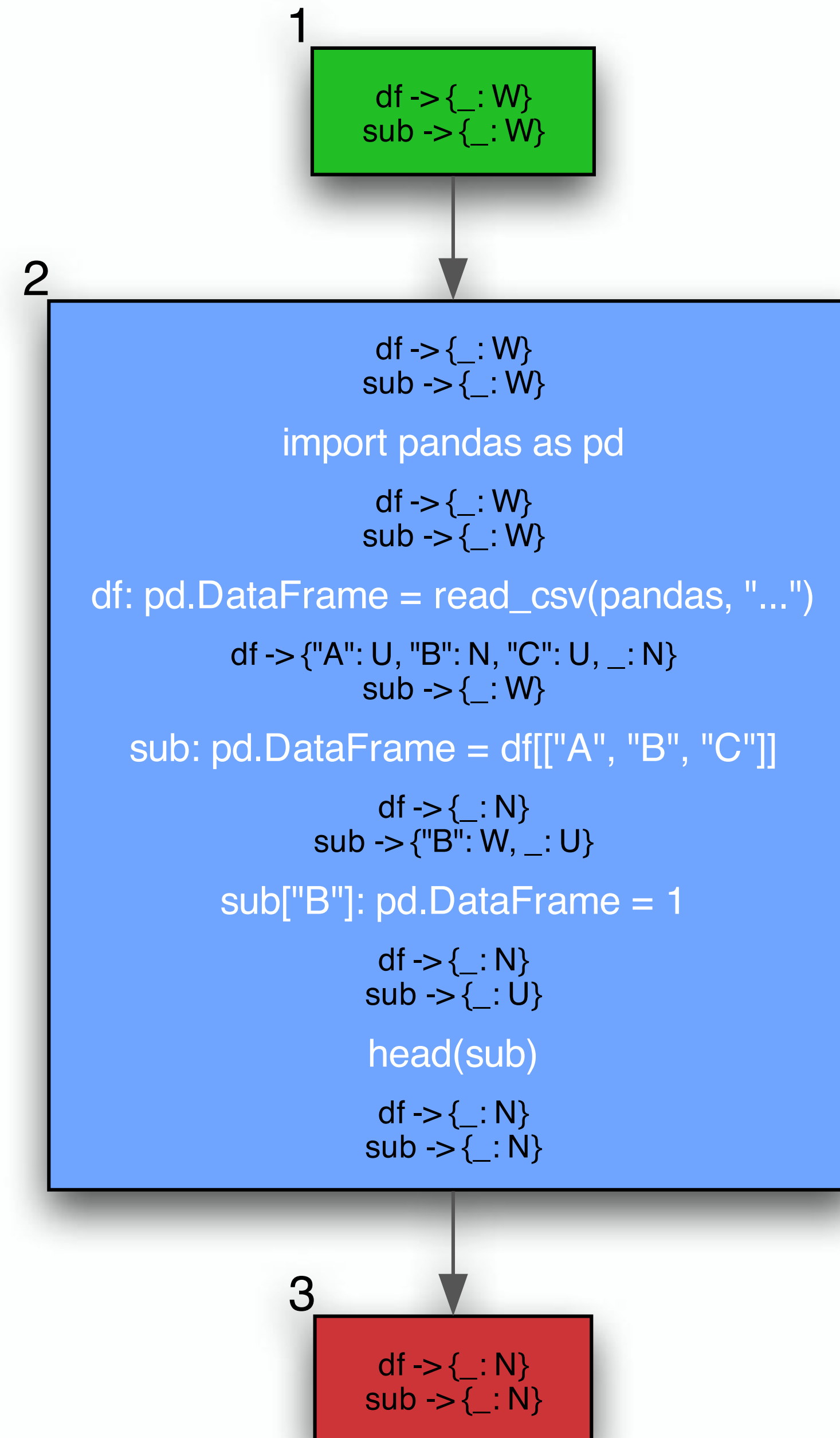
INTERACTIVE STATIC ANALYSIS

STATIC AND DYNAMIC ANALYSIS COMBINATIONS



(Un)expected + (Un)used Data





Bibliography

[Urban18] **Caterina Urban and Peter Müller.** An Abstract Interpretation Framework for Data Usage. In ESOP, pages 683-710, 2018.

data (non-) usage

[Drobnjaković24] **Filip Drobnjaković, Pavle Subotić, and Caterina Urban.** An Abstract Interpretation-Based Data Leakage Static Analysis. In TASE, pages 109-126, 2024.

data leakage static analysis