

Static Analysis for Data Science

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



Caterina Urban

December 14th, 2023

Year 2023-2024

Data Science is Everywhere

data is **cheap** and **ubiquitous**



mobile devices



data science is **revolutionizing industries**



retail

- personalized recommendations
- targeted marketing



pharmaceutical

- predictive models
- patient selection



manufacturing

- equipment failure predictions
- internet of things



finance

- predictive models
- customized product offerings



energy

- exploration and discovery
- accident prevention



health care

- personalized treatments
- preventive care

DATA SCIENCE SOFTWARE



data



data preparation



model training



model deployment



predictions



TODAY

NEXT LESSON

Ubiquitous Programming Errors

data science means **programming**



programming means **programming errors**

programming errors that do not cause failures can have **serious consequences**



pharmaceutical



energy



finance



health care



Anomalously Unused Data

The Reinhart-Rogoff

	B	C	I	J	K	L	M
2			Real GDP growth				
3			Debt/GDP				
4	Country	Coverage	30 or less	30 to 60	60 to 90	90 or above	30 or less
26			3.7	3.0	3.5	1.7	5.5
27	Minimum		1.6	0.3	1.3	-1.8	0.8
						3.6	13.3
						-2.0	n.a.
						2.4	n.a.
						n.a.	6.3
						n.a.	9.9
						n.a.	7.9

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 widely admired earlier book on the history of financial

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.

In response, Labour's shadow health secretary, Jonathan Ashworth, said, "Thousands of people are blissfully unaware they have been exposed to covid, spreading this deadly virus at a time when we are in the middle of a second wave and we are in the middle of a much

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 laboratories to expand testing and at the same time put local public health teams in charge of contact tracing in their areas. Routine regular testing in high transmission areas would be carried out in 24 hours.

NEWS

BMJ: first published as 10.1136/bmj.m3891 on 6 October 2020. Downloaded from <https://www.bmj.com/>

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

Unused Data Analysis

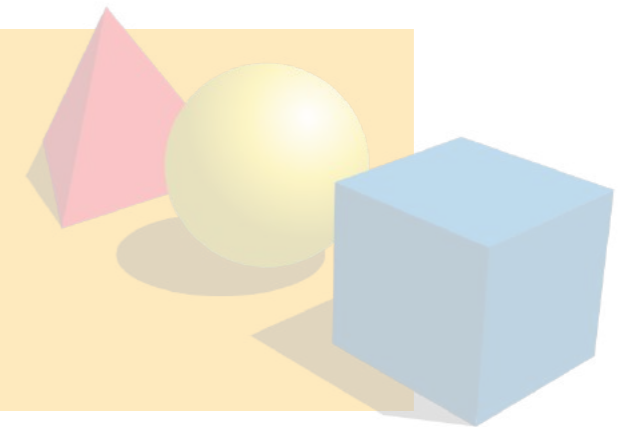
practical tools

targeting specific programs



algorithmic approaches

to decide program properties

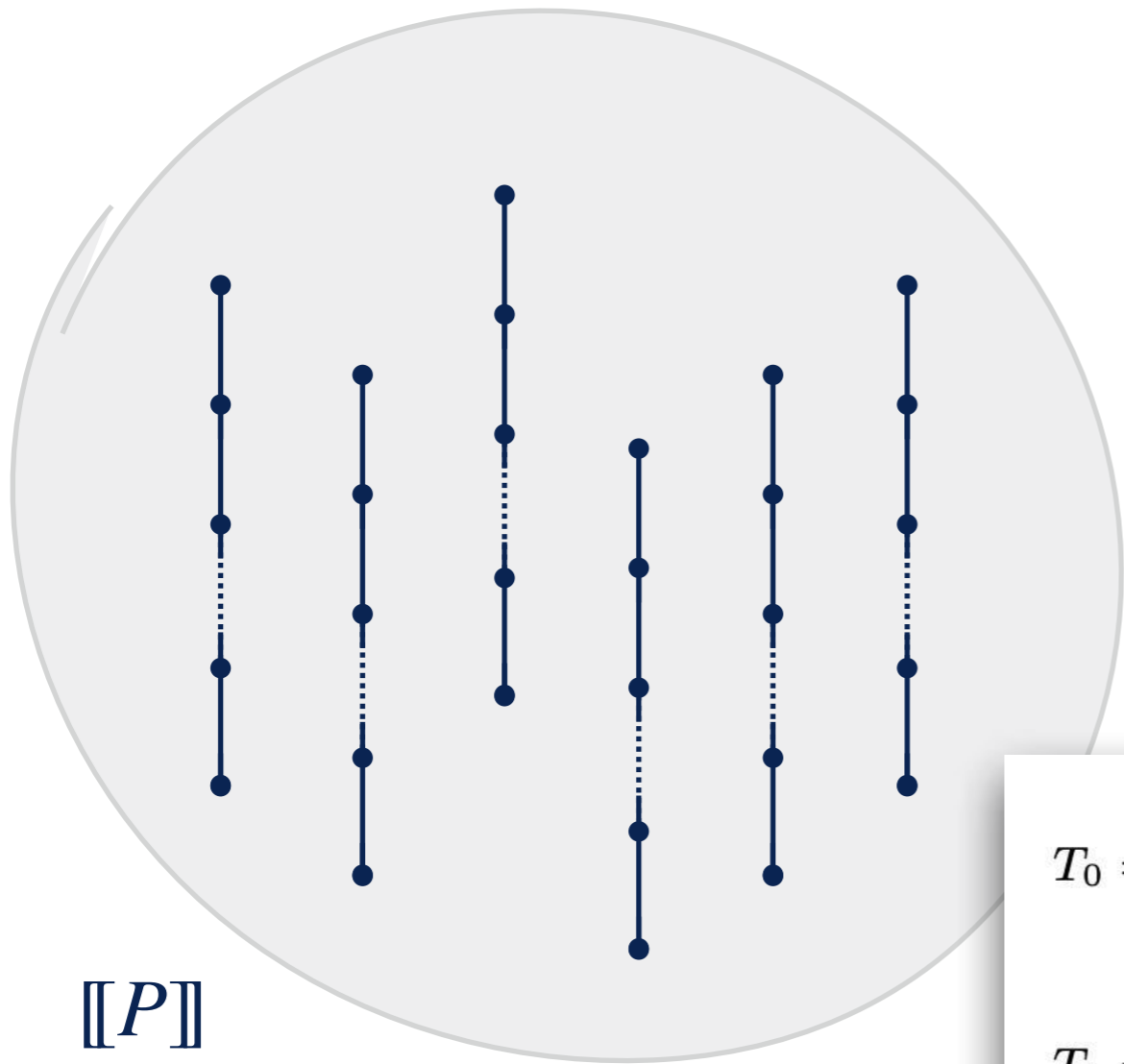


mathematical models

of the program behavior



Maximal Trace Semantics



Least fixpoint formulation of maximal traces

Idea: To get a least fixpoint formulation for whole \mathcal{M}_∞ , we merge finite and infinite maximal trace least fixpoint forms

- Fixpoint fusion:
- $\mathcal{M}_\infty \cap \Sigma^*$ is best defined on $(\mathcal{P}(\Sigma^*), \subseteq, \cup, \cap, \emptyset, \Sigma^*)$.
 - $\mathcal{M}_\infty \cap \Sigma^\omega$ is best defined on $(\mathcal{P}(\Sigma^\omega), \supseteq, \cap, \cup, \Sigma^\omega, \emptyset)$, the **dual lattice**.
(we transform the greatest fixpoint into a least fixpoint!)
- We mix them into a **new** complete lattice $(\mathcal{P}(\Sigma^\infty), \subseteq, \cup, \cap, \perp, \top)$:
- $A \subseteq B \stackrel{\text{def}}{\iff} (A \cap \Sigma^*) \subseteq (B \cap \Sigma^*) \wedge (A \cap \Sigma^\omega) \supseteq (B \cap \Sigma^\omega)$
 - $A \cup B \stackrel{\text{def}}{=} ((A \cap \Sigma^*) \cup (B \cap \Sigma^*)) \cup ((A \cap \Sigma^\omega) \cap (B \cap \Sigma^\omega))$
 - $A \cap B \stackrel{\text{def}}{=} ((A \cap \Sigma^*) \cap (B \cap \Sigma^*)) \cup ((A \cap \Sigma^\omega) \cup (B \cap \Sigma^\omega))$
 - $\perp \stackrel{\text{def}}{=} \Sigma^\omega$
 - $\top \stackrel{\text{def}}{=} \Sigma^*$

In this lattice, $\mathcal{M}_\infty = \text{lfp } F_s$ where $F_s(T) \stackrel{\text{def}}{=} B \cup T \cap T$

$$T_0 = \left\{ \begin{array}{c} \Sigma^\omega \\ \text{~~~~~} \end{array} \right\}$$

$$T_1 = \left\{ \begin{array}{c} \Omega \\ \bullet \end{array} \right\} \cup \left\{ \begin{array}{c} \tau \quad \Sigma^\omega \\ \bullet \text{-----} \end{array} \right\}$$

$$T_2 = \left\{ \begin{array}{c} \Omega \\ \bullet \end{array} \right\} \cup \left\{ \begin{array}{c} \tau \quad \Omega \\ \bullet \text{-----} \end{array} \right\} \cup \left\{ \begin{array}{c} \tau \quad \tau \quad \Sigma^\omega \\ \bullet \text{-----} \end{array} \right\}$$

Maximal Trace Semantics



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

passing = **True**
 english → _ english → _
 math → **T** math → **T**
 science → _ science → _
 bonus → _ bonus → _
 passing → ? passing → **T**

[P]

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → _ english → _ english → _ english → _
 math → **F** math → **F** math → **F** math → **F**
 science → _ science → _ science → _ science → _
 bonus → **T** bonus → **T** bonus → **T** bonus → **T**
 passing → ? passing → **T** passing → **T** passing → **T**

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → _ english → _ english → _ english → _
 math → **F** math → **F** math → **F** math → **F**
 science → _ science → _ science → _ science → _
 bonus → **F** bonus → **F** bonus → **F** bonus → **F**
 passing → ? passing → **T** passing → **F** passing → **F**

Input Data (Non-)Usage

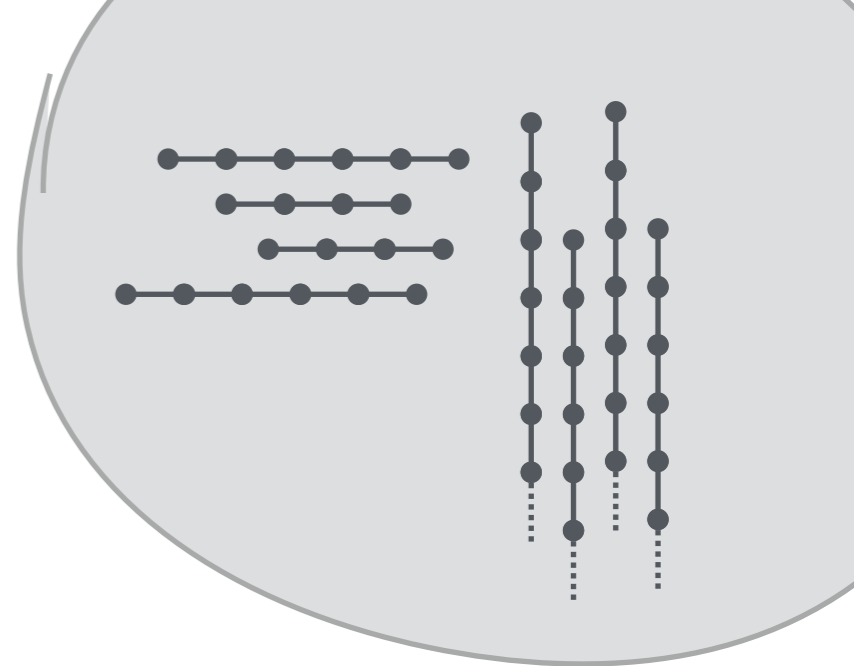
$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \forall i \in J \subseteq I_P: \text{UNUSED}_i(\llbracket P \rrbracket) \}$$

\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 $J \subseteq I_P$

$$\text{UNUSED}_i(\llbracket M \rrbracket) \stackrel{\text{def}}{=} \forall t \in \llbracket P \rrbracket, v \in \mathcal{V}: t_0(i) \neq v \Rightarrow \exists t' \in \llbracket P \rrbracket: \\ (\forall 0 \leq j \leq |I_P|: j \neq i \Rightarrow t_0(j) = t'_0(j)) \\ \wedge t'_0(i) = v \\ \wedge t_\omega = t'_\omega$$

Intuitively: **any possible program outcome** is possible from **any value of the input** variable $x_{0,i}$

Input Data (Non-)Usage



[P]

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

passing = **True**
 english → _ english → _
 math → **T** math → **T**
 science → _ science → _
 bonus → _ bonus → _
 passing → ? passing → **T**

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → _ english → _ english → _ english → _
 math → **F** math → **F** math → **F** math → **F**
 science → _ science → _ science → _ science → _
 bonus → **T** bonus → **T** bonus → **T** bonus → **T**
 passing → ? passing → **T** passing → **T** passing → **T**

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → _ english → _ english → _ english → _
 math → **F** math → **F** math → **F** math → **F**
 science → _ science → _ science → _ science → _
 bonus → **F** bonus → **F** bonus → **F** bonus → **F**
 passing → ? passing → **T** passing → **F** passing → **F**

Input Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \forall i \in J \subseteq I_P : \text{UNUSED}_i(\llbracket P \rrbracket) \}$$

\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 $J \subseteq I_P$

$$\text{UNUSED}_i(\llbracket M \rrbracket) \stackrel{\text{def}}{=} \forall t \in \llbracket P \rrbracket, v \in \mathcal{V} : t_0(i) \neq v \Rightarrow \exists t' \in \llbracket P \rrbracket : \\ (\forall 0 \leq j \leq |I_P| : j \neq i \Rightarrow t_0(j) = t'_0(j)) \\ \wedge t'_0(i) = v \\ \wedge t_\omega = t'_\omega$$

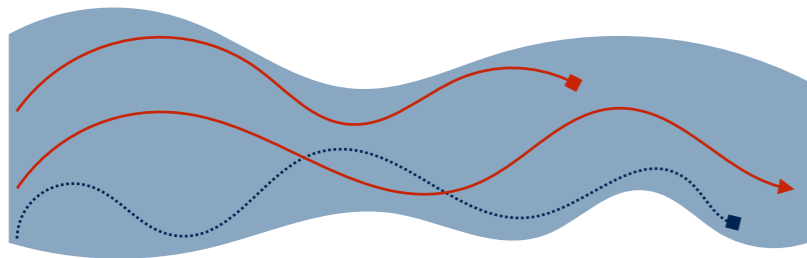
Intuitively: **any possible program outcome** is possible from any value of the input variable $x_{0,i}$

General properties

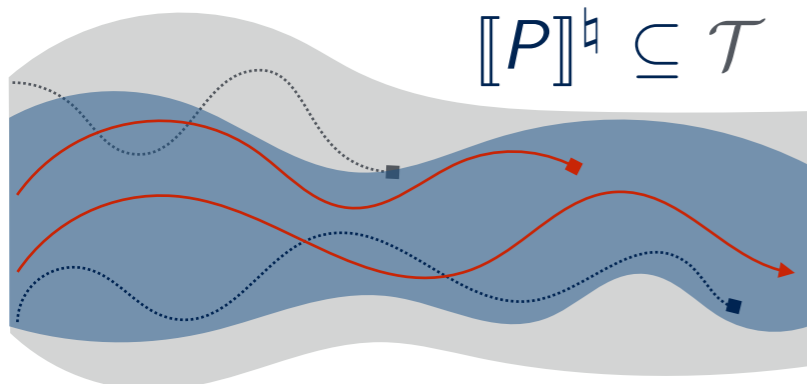
General setting:

- given a program $\text{prog} \in \text{Prog}$
 - its **semantics**: $\llbracket \cdot \rrbracket : \text{Prog} \rightarrow \mathcal{P}(\Sigma^*)$ is a set of finite traces
 - a **property** P is the **set** of correct program semantics
i.e., a **set of sets of traces** $P \in \mathcal{P}(\mathcal{P}(\Sigma^*))$
- \subseteq gives an information order on properties
 $P \subseteq P'$ means that P' is weaker than P (allows more semantics)

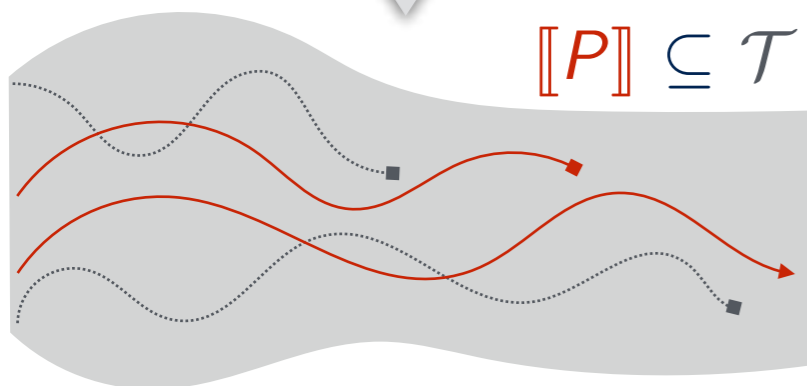
Trace Properties



$$\llbracket P \rrbracket \subseteq \llbracket P \rrbracket^{\sharp}$$



$$\llbracket P \rrbracket^{\sharp} \subseteq \mathcal{T}$$



$$\llbracket P \rrbracket \subseteq \mathcal{T}$$

Collecting semantics and properties

Restricted properties

Reasoning on (and abstracting) $\mathcal{P}(\mathcal{P}(\Sigma^*))$ is **hard!**

In the following, we use a **simpler** setting:

- a property is a **set of traces** $P \in \mathcal{P}(\Sigma^*)$
- the collecting semantics is a **set of traces**: $Col(prog) \stackrel{def}{=} \llbracket prog \rrbracket$
- the verification problem remains an inclusion check: $\llbracket prog \rrbracket \subseteq P$
- abstractions will over-approximate the set of traces $\llbracket prog \rrbracket$

Example properties:

- state property $P \stackrel{def}{=} S^*$ (remains in the set S of safe states)
- maximal execution time: $P \stackrel{def}{=} S^{\leq k}$
- ordering: $P \stackrel{def}{=} (\Sigma \setminus \{b\})^* \cdot a \cdot \Sigma^* \cdot b \cdot \Sigma^*$ (a occurs before b)

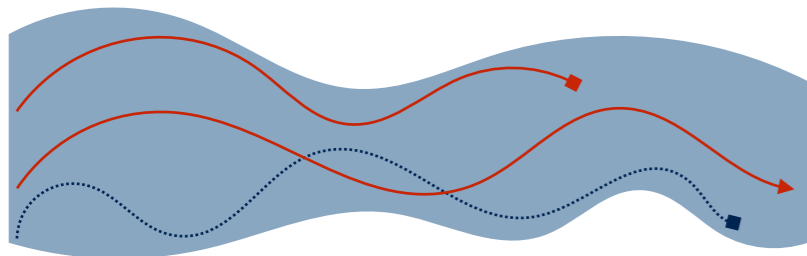
Course 2

Program Semantics and Properties

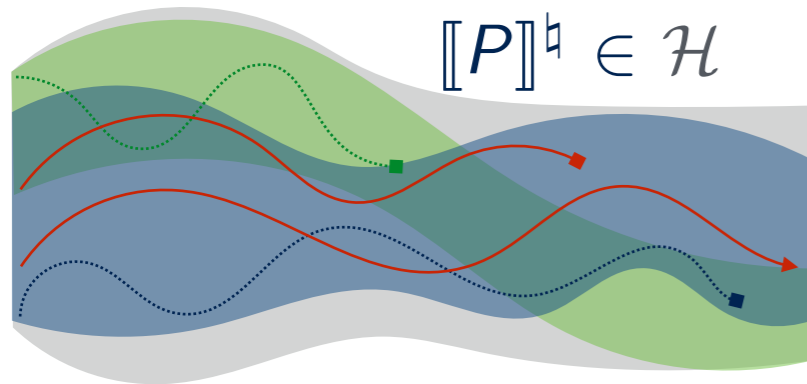
Antoine Miné

p. 25 / 98

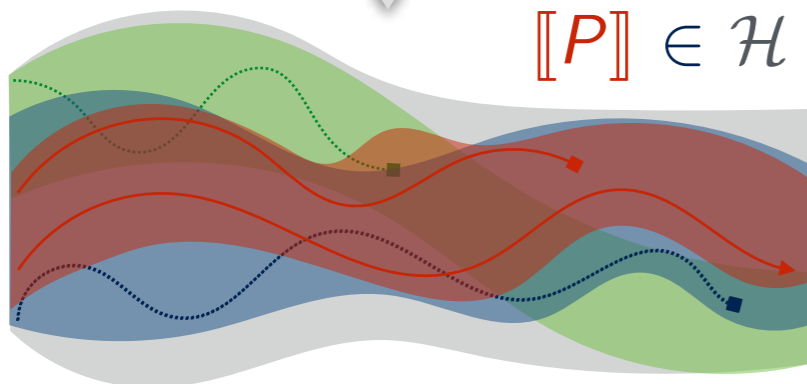
Program Properties



$$\llbracket P \rrbracket \subseteq \llbracket P \rrbracket^{\#}$$



$$\llbracket P \rrbracket^{\#} \in \mathcal{H}$$



$$\llbracket P \rrbracket \in \mathcal{H}$$

Collecting semantics and properties

General properties

General setting:

- given a program $prog \in Prog$
- its **semantics**: $\llbracket \cdot \rrbracket : Prog \rightarrow \mathcal{P}(\Sigma^*)$ is a set of finite traces
- a **property** P is the **set** of correct program semantics
i.e., a **set of sets of traces** $P \in \mathcal{P}(\mathcal{P}(\Sigma^*))$

\subseteq gives an information order on properties
 $P \subseteq P'$ means that P' is weaker than P (allows more semantics)

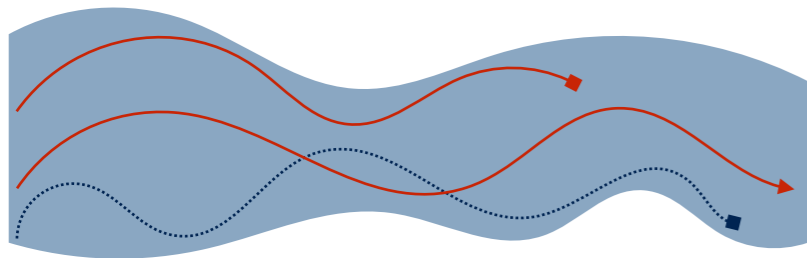
Course 2

Program Semantics and Properties

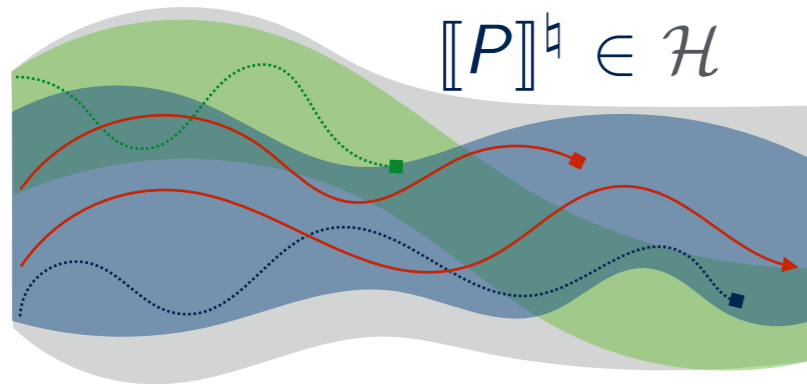
Antoine Miné

p. 23 / 98

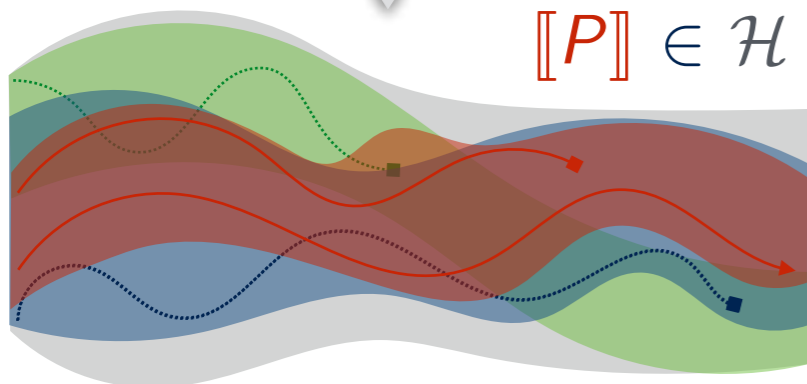
Subset-Closed Properties



$$\llbracket P \rrbracket \subseteq \llbracket P \rrbracket^{\sharp}$$



$$\llbracket P \rrbracket^{\sharp} \in \mathcal{H}$$



$$\llbracket P \rrbracket \in \mathcal{H}$$

Program Properties

$\llbracket P \rrbracket \subseteq \llbracket P \rrbracket^{\sharp}$

$\llbracket P \rrbracket^{\sharp} \in \mathcal{H}$

$\llbracket P \rrbracket \in \mathcal{H}$

General properties

Collecting semantics and properties

General setting:

- given a program $prog \in Prog$
- its semantics: $\llbracket \cdot \rrbracket : Prog \rightarrow \mathcal{P}(\Sigma^*)$ is a set of finite traces
- a property P is the set of correct program semantics

i.e., a set of sets of traces $P \in \mathcal{P}(\mathcal{P}(\Sigma^*))$

\subseteq gives an information order on properties

$P \subseteq P'$ means that P' is weaker than P (allows more semantics)

Course 2

Program Semantics and Properties

Antoine Mine

p. 23 / 98

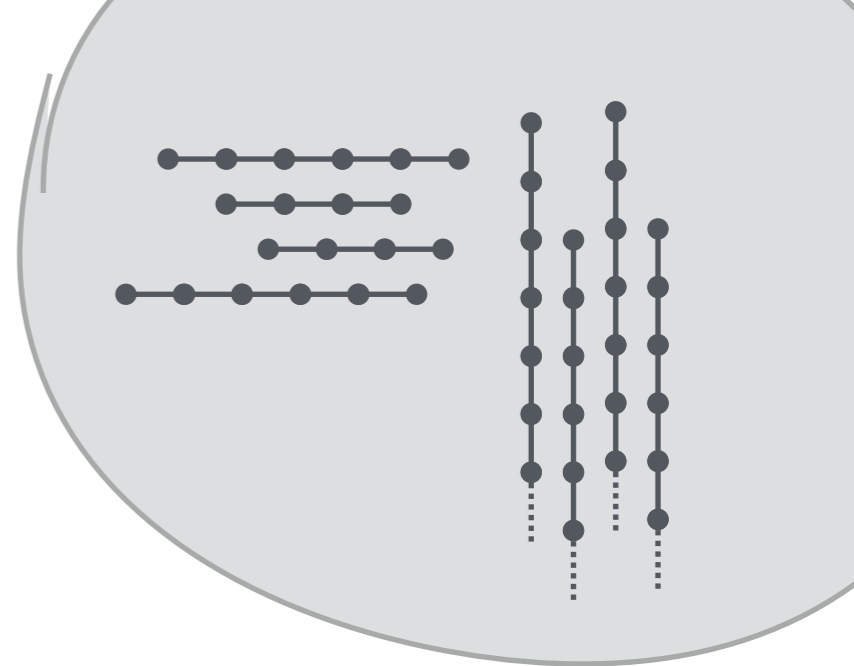
Lesson 10

Static Analysis for Data Science

Caterina Urban

16

Input Data (Non-)Usage



[P]

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

passing = **True**
 english → _ english → _
 math → **T** math → **T**
 science → _ science → _
 bonus → _ bonus → _
 passing → ? passing → **T**

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → _ english → _ english → _ english → _
 math → **F** math → **F** math → **F** math → **F**
 science → _ science → _ science → _ science → _
 bonus → **T** bonus → **T** bonus → **T** bonus → **T**
 passing → ? passing → **T** passing → **T** passing → **T**

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → _ english → _ english → _ english → _
 math → **F** math → **F** math → **F** math → **F**
 science → _ science → _ science → _ science → _
 bonus → **F** bonus → **F** bonus → **F** bonus → **F**
 passing → ? passing → **T** passing → **F** passing → **F**

Input Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

\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 $J \subseteq I_P$

$$\text{UNUSED}_J(\llbracket M \rrbracket) \stackrel{\text{def}}{=} \forall t \in \llbracket P \rrbracket, V \in \mathcal{V} : t_0(J) \neq V \Rightarrow \exists t' \in \llbracket P \rrbracket : \\ (\forall 0 \leq i \leq |I_P| : i \notin J \Rightarrow t_0(i) = t'_0(i)) \\ \wedge t'_0(J) = V \\ \wedge t'_\omega = t_\omega$$

Intuitively: **any possible program outcome** is possible from **any value of the input** variables in J

Theorem

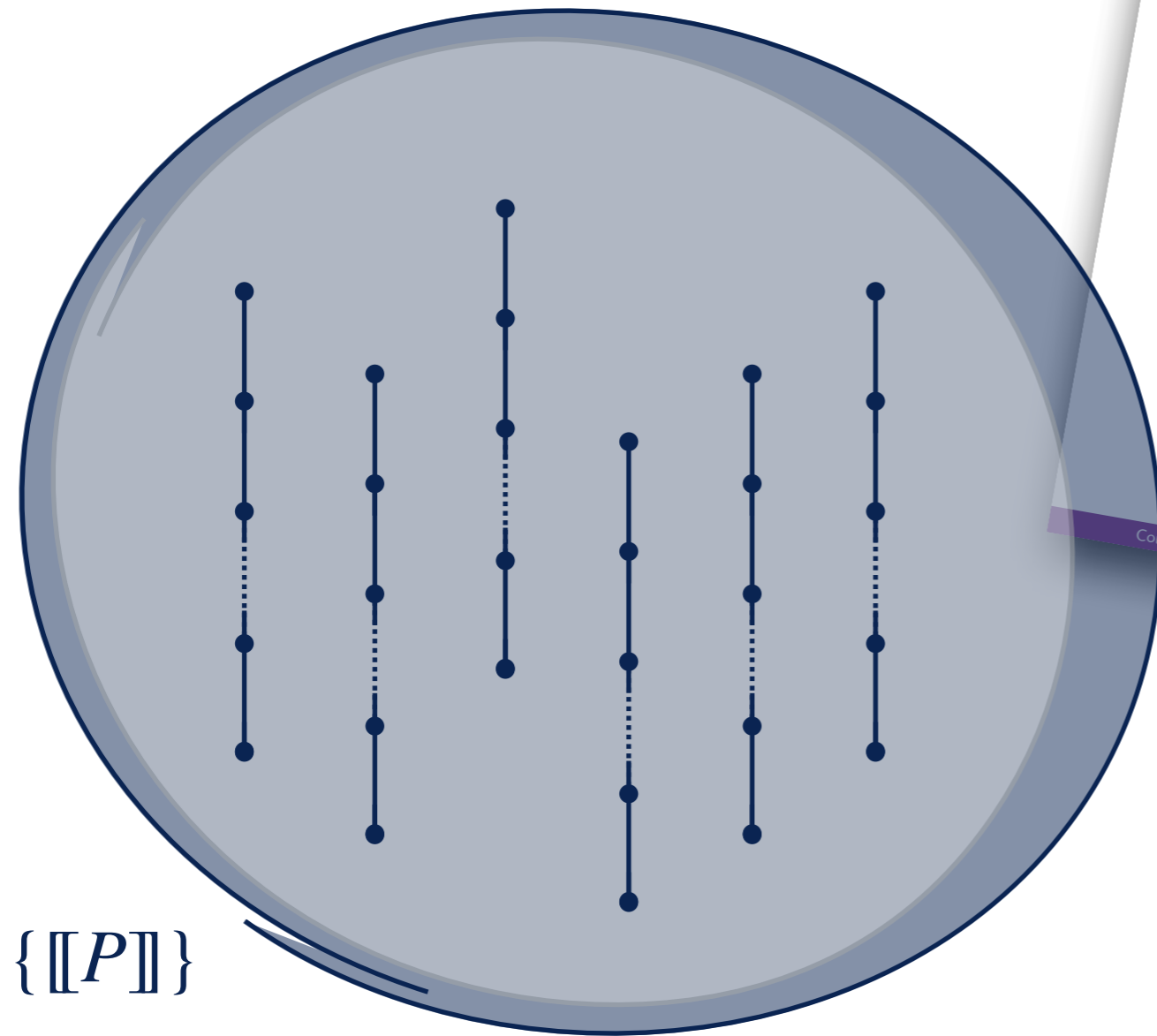
$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \} \subseteq \mathcal{N}_J$$

General properties

General setting:

- given a program $\text{prog} \in \text{Prog}$
 - its **semantics**: $\llbracket \cdot \rrbracket : \text{Prog} \rightarrow \mathcal{P}(\Sigma^*)$ is a set of finite traces
 - a **property** P is the **set** of correct program semantics
i.e., a **set of sets of traces** $P \in \mathcal{P}(\mathcal{P}(\Sigma^*))$
- \subseteq gives an information order on properties
 $P \subseteq P'$ means that P' is weaker than P (allows more semantics)

Collecting Semantics



General collecting semantics

The collecting semantics $Col : Prog \rightarrow \mathcal{P}(\mathcal{P}(\Sigma^*))$ is the **strongest property** of a program

Hence: $Col(prog) \stackrel{def}{=} \{[[prog]]\}$

Benefits: uniformity of semantics and properties, \subseteq information order

- given a program $prog$ and a property $P \in \mathcal{P}(\mathcal{P}(\Sigma^*))$ the **verification problem** is an inclusion check:

$$Col(prog) \subseteq P$$

- generally, the collecting semantics **cannot be computed**, we settle for a weaker property $S^\#$ that
- is sound: $Col(prog) \subseteq S^\#$
- implies the desired property: $S^\# \subseteq P$

Course 2

Program Semantics and Properties

Antoine Miné

p. 24 / 98

Collecting Semantics

Intuition

Property (*by extension*): set of elements that have that property

Property “being Patrick Cousot”



Property “being program P ”

$$\{\llbracket P \rrbracket\}$$

Collecting semantics and properties

General collecting semantics

The collecting semantics $Col : Prog \rightarrow \mathcal{P}(\mathcal{P}(\Sigma^*))$ is the **strongest property** of a program

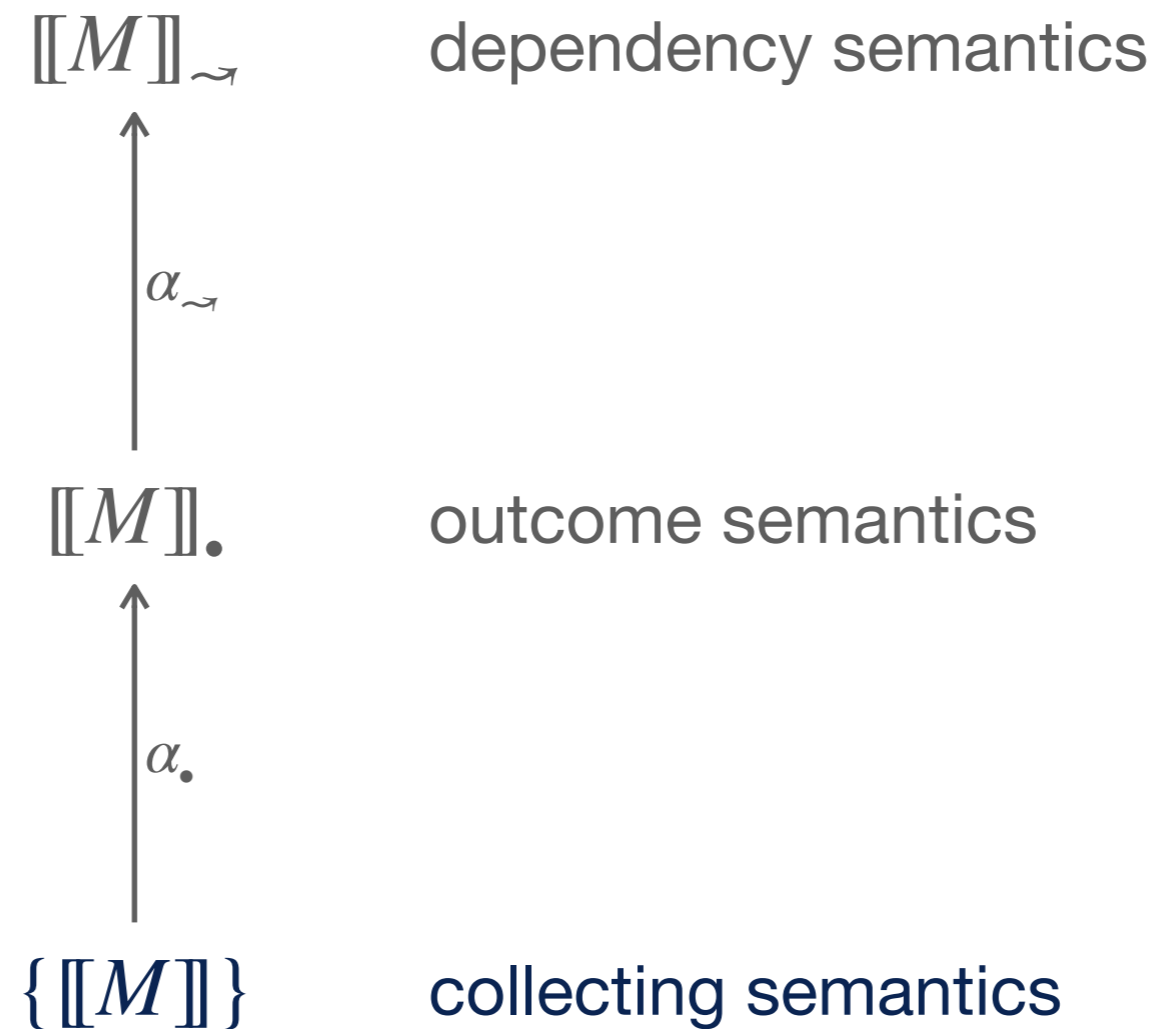
Hence: $Col(prog) \stackrel{\text{def}}{=} \{\llbracket prog \rrbracket\}$

Benefits: uniformity of semantics and properties, \subseteq information order

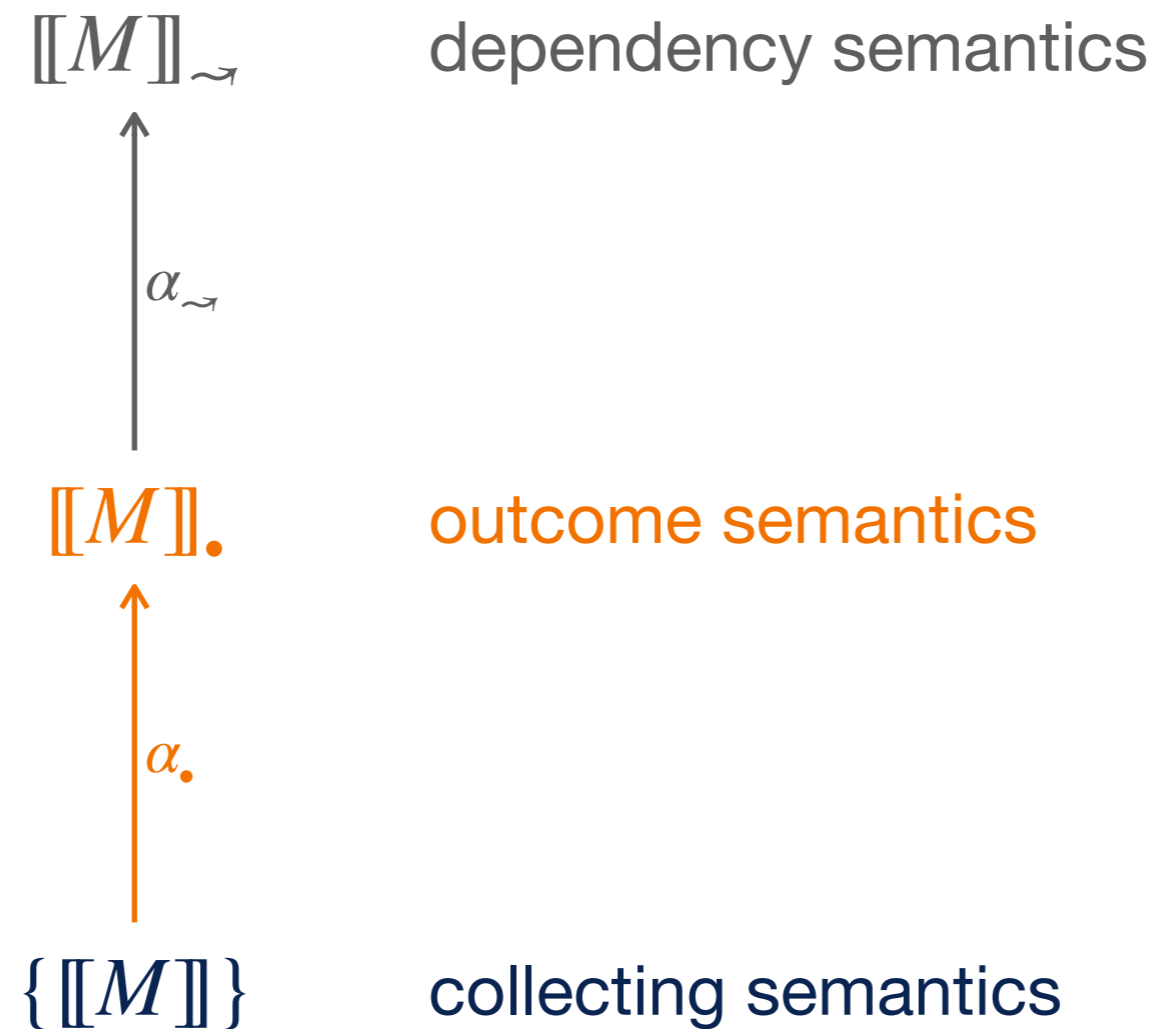
- given a program $prog$ and a property $P \in \mathcal{P}(\mathcal{P}(\Sigma^*))$ the **verification problem** is an inclusion check:
 $Col(prog) \subseteq P$
- generally, the collecting semantics **cannot be computed**, we settle for a weaker property $S^\#$ that
 - is sound: $Col(prog) \subseteq S^\#$
 - implies the desired property: $S^\# \subseteq P$

Course 2 Program Semantics and Properties Antoine Miné p. 24 / 98

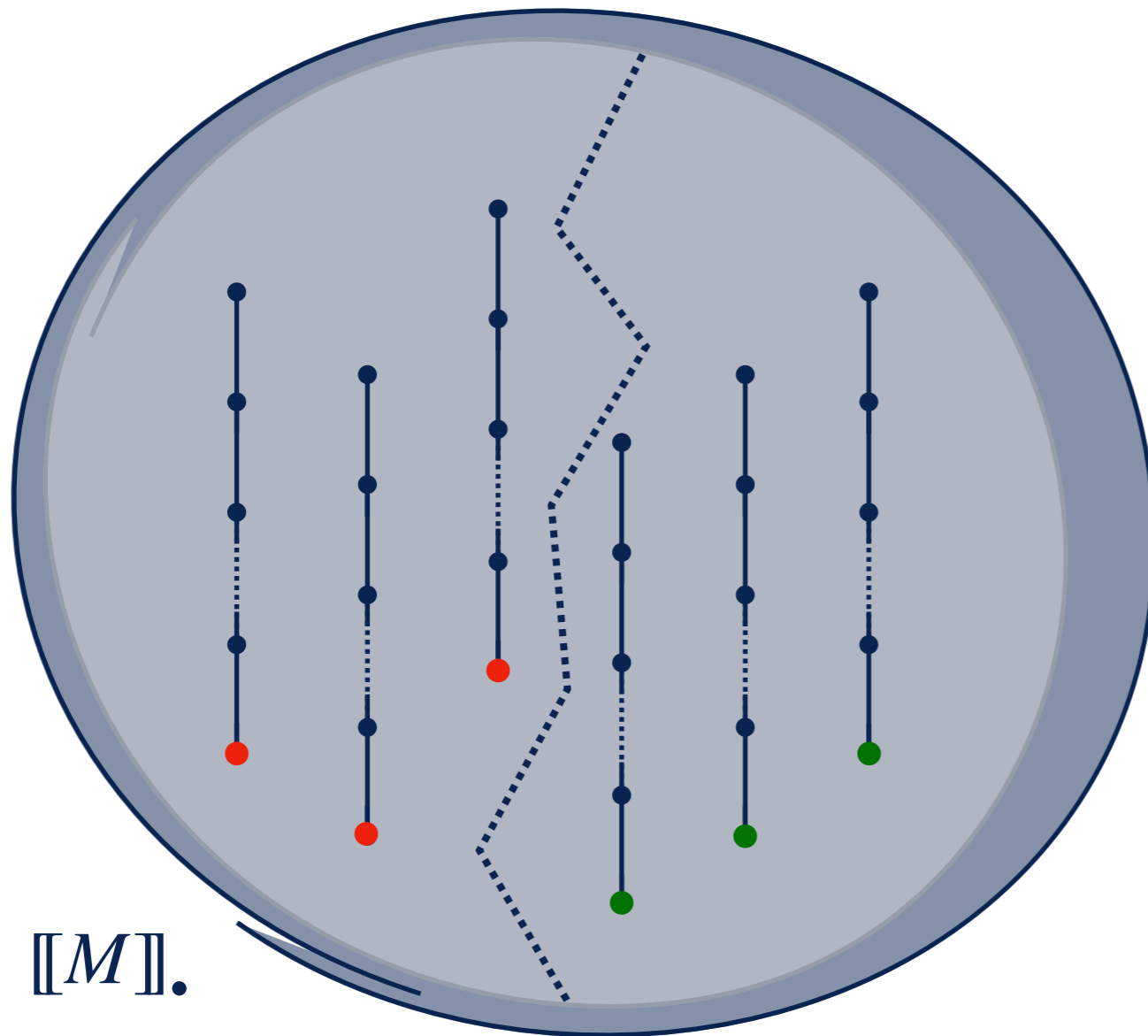
(Another) Hierarchy of Semantics





(Another) Hierarchy of Semantics



Outcome Semantics




 **partitioning** a set of traces that satisfies input data (non-)usage **with respect to the program outcome** yields sets of traces that also satisfy input data (non-)usage

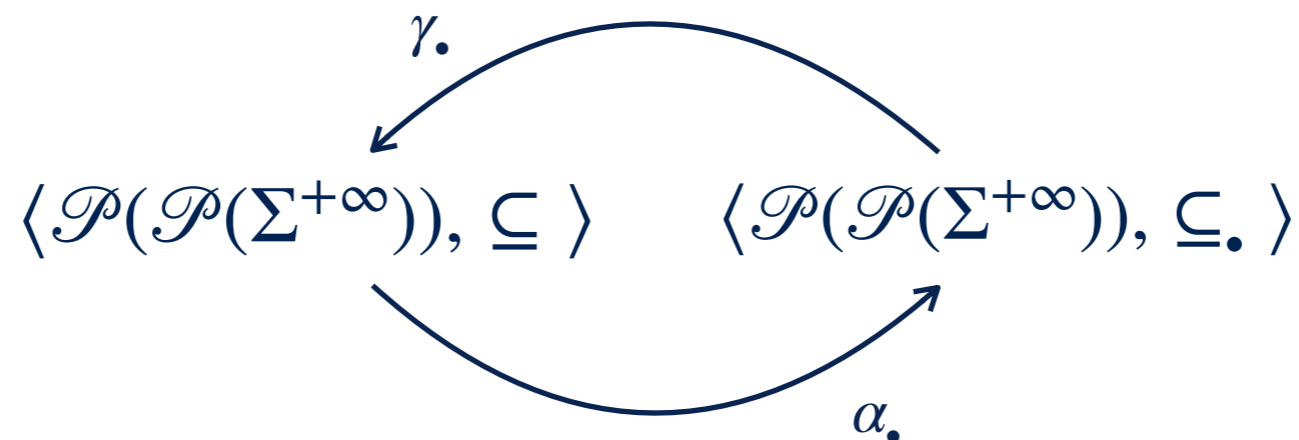
Outcome Semantics

$$\mathbb{O} \stackrel{\text{def}}{=} \{ \Sigma_{o_1=v_1, \dots, o_k=v_k}^+ \mid v_1, \dots, v_k \in \mathcal{V} \} \cup \{ \Sigma^\omega \}$$

outcomes

Lemma

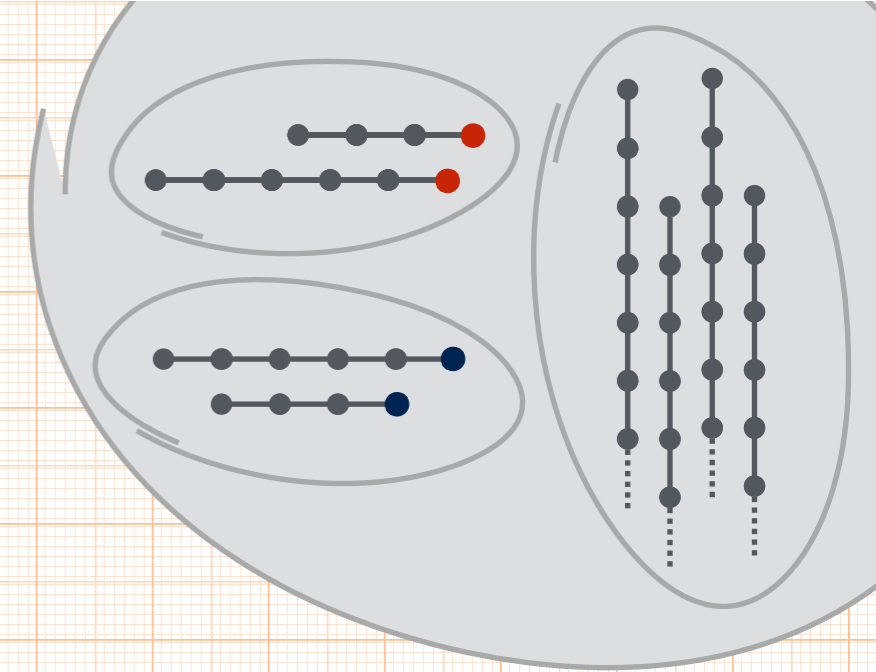
$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \cap O \mid O \in \mathbb{O} \} \subseteq \mathcal{N}_J$$



$$\alpha_\bullet(S) \stackrel{\text{def}}{=} \{ T \cap O \mid T \in S \wedge O \in \mathbb{O} \}$$

outcome abstraction

Outcome Semantics



$[P]$

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

passing = **True**
 english → **_** english → **_**
 math → **T** math → **T**
 science → **_** science → **_**
 bonus → **_** bonus → **_**
 passing → **?** passing → **T**

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → **_** english → **_** english → **_**
 math → **F** math → **F** math → **F**
 science → **_** science → **_** science → **_**
 bonus → **T** bonus → **T** bonus → **T**
 passing → **?** passing → **T** passing → **T**

passing = **True** passing = **False or bonus** passing = **False or bonus**
 english → **_** english → **_** english → **_**
 math → **F** math → **F** math → **F**
 science → **_** science → **_** science → **_**
 bonus → **F** bonus → **F** bonus → **F**
 passing → **?** passing → **T** passing → **F**

\emptyset
 ↖
 unfeasible
 non-termination
 outcome

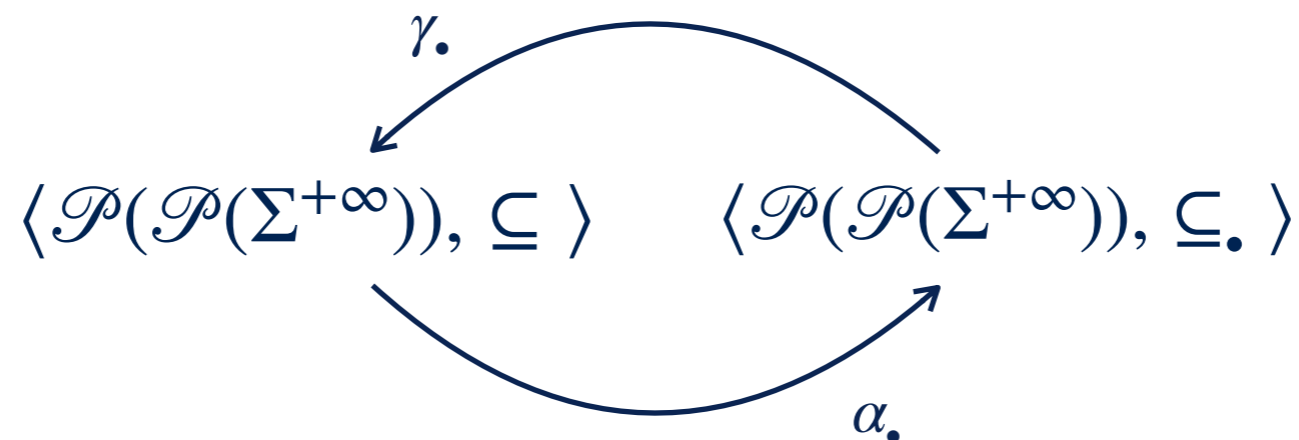
Outcome Semantics

$$\mathbb{O} \stackrel{\text{def}}{=} \{ \Sigma_{o_1=v_1, \dots, o_k=v_k}^+ \mid v_1, \dots, v_k \in \mathcal{V} \} \cup \{ \Sigma^\omega \}$$

outcomes

Lemma

$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \cap O \mid O \in \mathbb{O} \} \subseteq \mathcal{N}_J$$

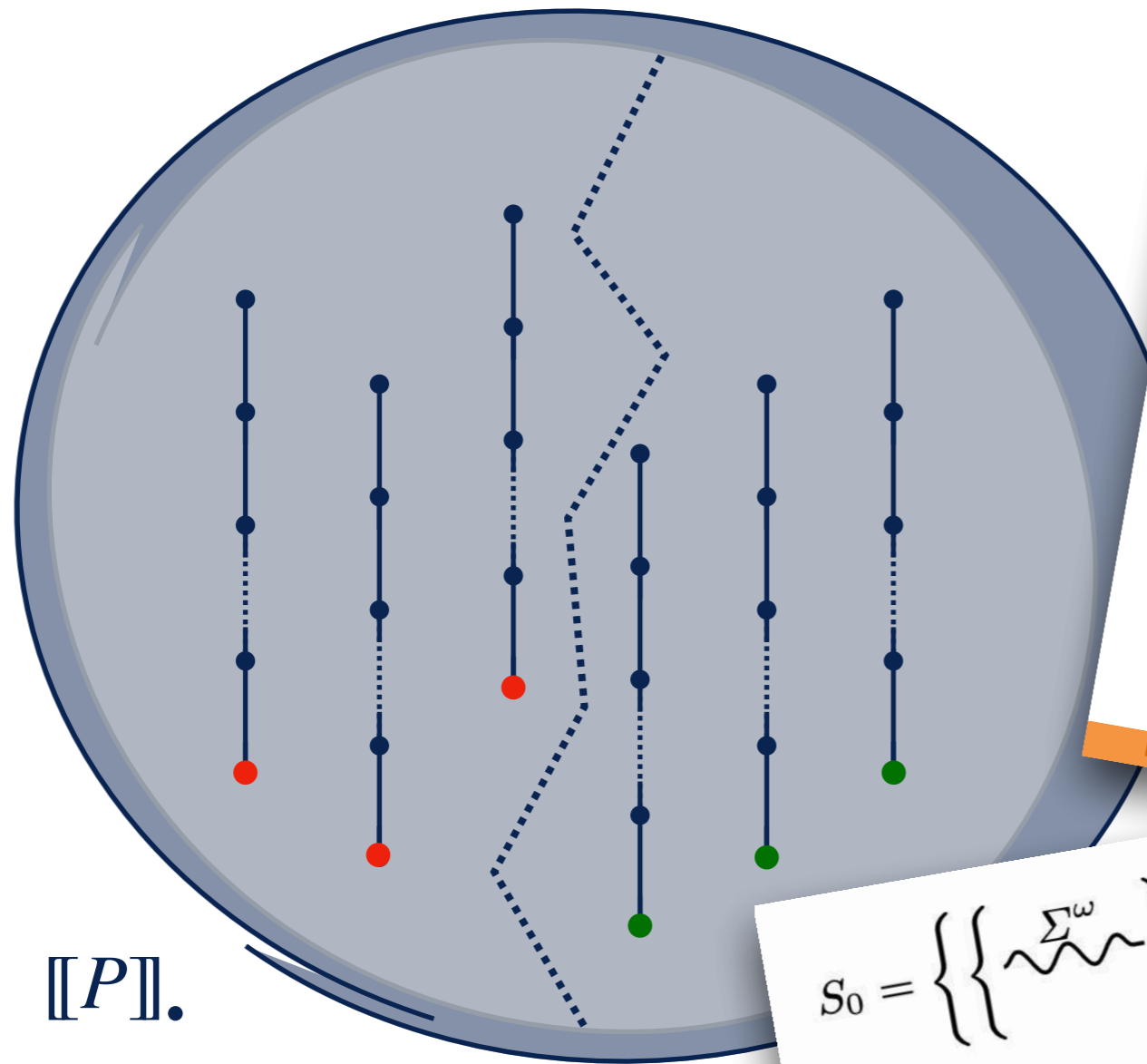


$$\alpha_\bullet(S) \stackrel{\text{def}}{=} \{ T \cap O \mid T \in S \wedge O \in \mathbb{O} \}$$

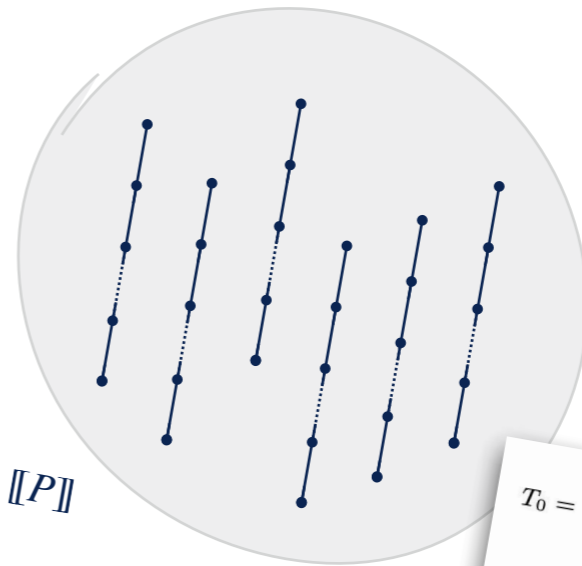
outcome abstraction

$$\llbracket P \rrbracket_\bullet \stackrel{\text{def}}{=} \alpha_\bullet(\{ \llbracket P \rrbracket \}) = \{ \llbracket P \rrbracket \cap O \mid O \in \mathbb{O} \}$$

Outcome Semantics



Maximal Trace Semantics



Least fixpoint formulation of maximal trace semantics

we merge finite and infinite maximal trace least fixpoint formulation for whole \mathcal{M}_∞

Fixpoint fusion:

$\mathcal{M}_\infty \cap \Sigma^*$ is best defined on $(\mathcal{P}(\Sigma^*), \subseteq, \cup, \cap, \emptyset, \Sigma^*)$.

$\mathcal{M}_\infty \cap \Sigma^\omega$ is best defined on $(\mathcal{P}(\Sigma^\omega), \subseteq, \cup, \cap, \emptyset)$, the dual lattice. (we transform the greatest fixpoint into a least fixpoint!)

We mix them into a new complete lattice $(\mathcal{P}(\Sigma^\infty), \subseteq, \cup, \cap, \perp, \top)$

- $A \subseteq B \stackrel{\text{def}}{=} (A \cap \Sigma^*) \subseteq (B \cap \Sigma^*) \wedge (A \cap \Sigma^\omega) \supseteq (B \cap \Sigma^\omega)$
- $A \cup B \stackrel{\text{def}}{=} ((A \cap \Sigma^*) \cup (B \cap \Sigma^*)) \cup ((A \cap \Sigma^\omega) \cap (B \cap \Sigma^\omega))$
- $A \cap B \stackrel{\text{def}}{=} ((A \cap \Sigma^*) \cap (B \cap \Sigma^*)) \cup ((A \cap \Sigma^\omega) \cap (B \cap \Sigma^\omega))$
- $\perp \stackrel{\text{def}}{=} \Sigma^\omega$
- $\top \stackrel{\text{def}}{=} \Sigma^*$

In this lattice, $\mathcal{M}_\infty = \text{lfp } F_s$ where $F_s(T) \stackrel{\text{def}}{=} B \cup T \cap T$

Lesson 10

Static Analysis for Data Science

Caterina Urban

10

$$T_0 = \left\{ \Sigma^\omega \right\}$$

$$T_1 = \left\{ \Omega \right\} \cup \left\{ \tau \cdot \Sigma^\omega \right\}$$

$$T_2 = \left\{ \Omega \right\} \cup \left\{ \tau \cdot \Omega \right\} \cup \left\{ \tau \cdot \tau \cdot \Sigma^\omega \right\}$$

$$S_0 = \left\{ \left\{ \Sigma^\omega \right\}, \emptyset \right\}$$

$$S_1 = \left\{ \left\{ \Omega_{o=v} \right\} \mid v \in V \right\} \cup \left\{ \left\{ \tau \cdot \Sigma^\omega \right\} \right\}$$

$$S_2 = \left\{ \left\{ \Omega_{o=v} \right\} \cup \left\{ \tau \cdot \Omega_{o=v} \right\} \mid v \in V \right\} \cup \left\{ \left\{ \tau \cdot \tau \cdot \Sigma^\omega \right\} \right\}$$

Outcome Semantics

$$S_1 \sqsubseteq S_2 \stackrel{\text{def}}{=} \bigwedge_{v_1, \dots, v_k \in V} S_{1, o_1=v_1, \dots, o_k=v_k}^+ \subseteq S_{2, o_1=v_1, \dots, o_k=v_k}^+ \wedge S_1^\omega \supseteq S_2^\omega$$

Theorem 1. The outcome semantics $\Lambda_\bullet \in \mathcal{P}(\mathcal{P}(\Sigma^{+\infty}))$ can be expressed as a least fixpoint in $\langle \mathcal{P}(\mathcal{P}(\Sigma^{+\infty})), \sqsubseteq, \sqcup, \sqcap, \{\Sigma^\omega, \emptyset\}, \{\emptyset, \Sigma^+\} \rangle$ as:

$$\Lambda_\bullet = \text{lfp}^{\sqsubseteq} \Theta_\bullet \tag{9}$$

$$\Theta_\bullet(S) \stackrel{\text{def}}{=} \{ \Omega_{o_1=v_1, \dots, o_k=v_k} \mid v_1, \dots, v_k \in V \} \cup \{ \tau ; T \mid T \in S \}$$

where $S_1 \sqcup S_2 \stackrel{\text{def}}{=} \{ S_{1, o_1=v_1, \dots, o_k=v_k}^+ \cup S_{2, o_1=v_1, \dots, o_k=v_k}^+ \mid v_1, \dots, v_k \in V \} \cup S_1^\omega \cup S_2^\omega$.

(proof by Kleenian Fixpoint Transfer [Urban18])

$$\begin{aligned}
 S_0 &= \left\{ \left\{ \Sigma^\omega \right\}, \emptyset \right\} \\
 S_1 &= \left\{ \left\{ \Omega_{o=v} \right\} \mid v \in V \right\} \cup \left\{ \left\{ \tau \cdot \Sigma^\omega \right\} \right\} \\
 S_2 &= \left\{ \left\{ \Omega_{o=v} \right\} \right\} \cup \left\{ \left\{ \tau \cdot \Omega_{o=v} \right\} \mid v \in V \right\} \cup \left\{ \left\{ \tau \cdot \tau \cdot \Sigma^\omega \right\} \right\}
 \end{aligned}$$

Input Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

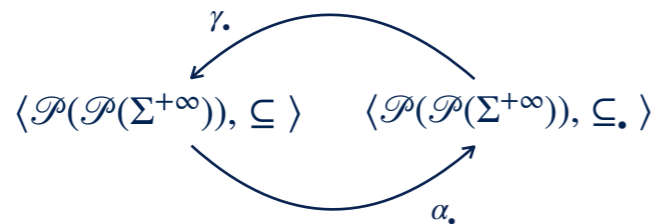
\mathcal{N}_J is the set of all programs P (or, rather, their semantics $\llbracket P \rrbracket$) that **do not use** the value of the input variables in $J \subseteq I_P$

Outcome Semantics

$$\mathbb{O} \stackrel{\text{def}}{=} \{ \Sigma_{o_1=v_1, \dots, o_k=v_k}^+ \mid v_1, \dots, v_k \in \mathcal{V} \} \cup \{ \Sigma^\omega \} \quad \text{outcomes}$$

Lemma

$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \cap O \mid O \in \mathbb{O} \} \subseteq \mathcal{N}_J$$



$$\alpha_\cdot(S) \stackrel{\text{def}}{=} \{ T \cap O \mid T \in S \wedge O \in \mathbb{O} \} \quad \text{outcome abstraction}$$

$$\llbracket P \rrbracket_\cdot \stackrel{\text{def}}{=} \alpha_\cdot(\{ \llbracket P \rrbracket \}) = \{ \llbracket P \rrbracket \cap O \mid O \in \mathbb{O} \}$$

$$t_0(J) \neq V \Rightarrow \exists t' \in \llbracket P \rrbracket : J \Rightarrow t_0(i) = t'_0(i)$$

Input Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

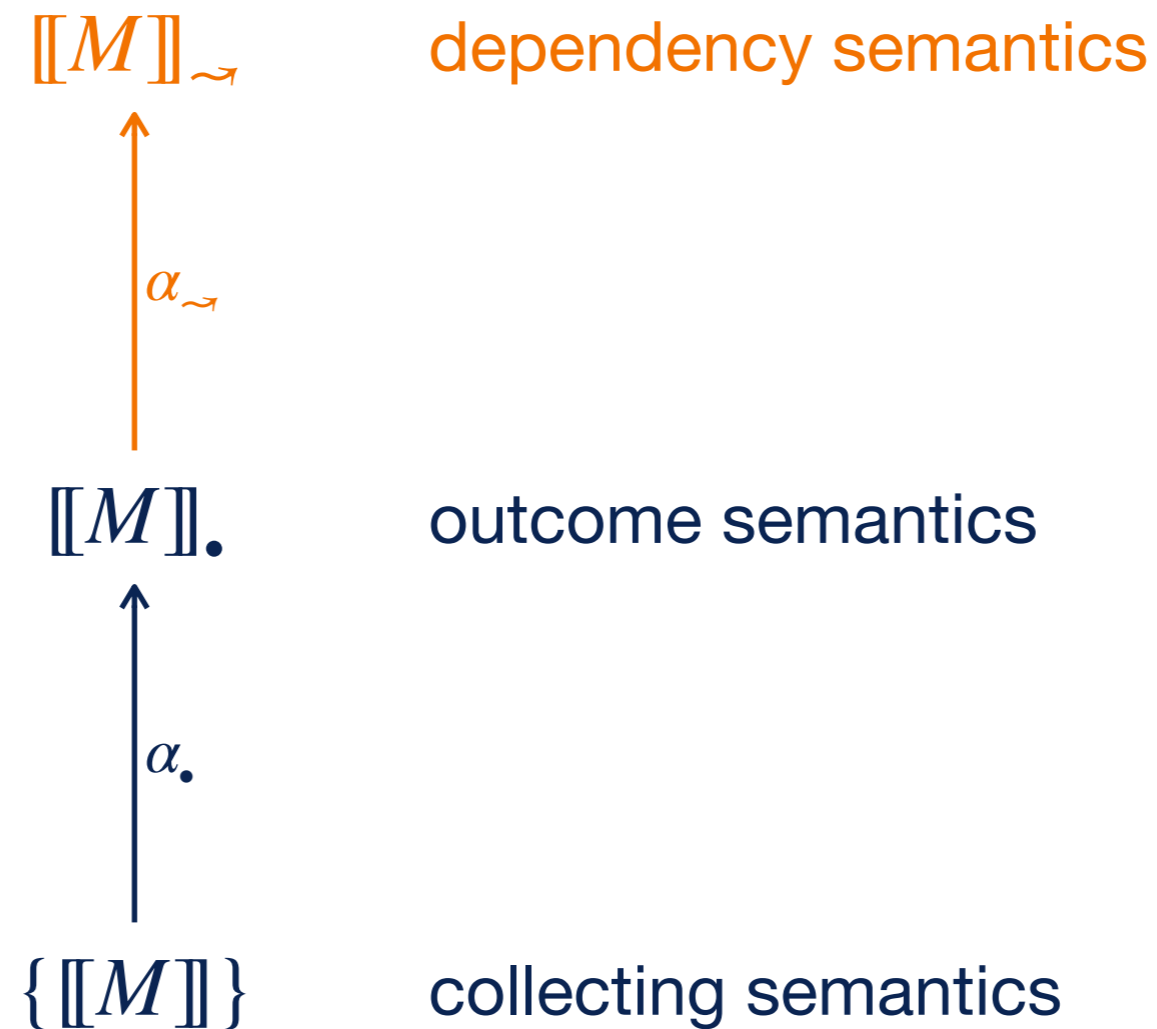
\mathcal{N}_J is the set of all programs P (or, rather, their semantics $\llbracket P \rrbracket$) that **do not use** the value of the input variables in $J \subseteq I_P$

$$\text{UNUSED}_J(\llbracket M \rrbracket) \stackrel{\text{def}}{=} \forall t \in \llbracket M \rrbracket, V \in \mathcal{V} : t_0(J) \neq V \Rightarrow \exists t' \in \llbracket M \rrbracket : (\forall 0 \leq i \leq |I_P| : i \notin J \Rightarrow t_0(i) = t'_0(i))$$

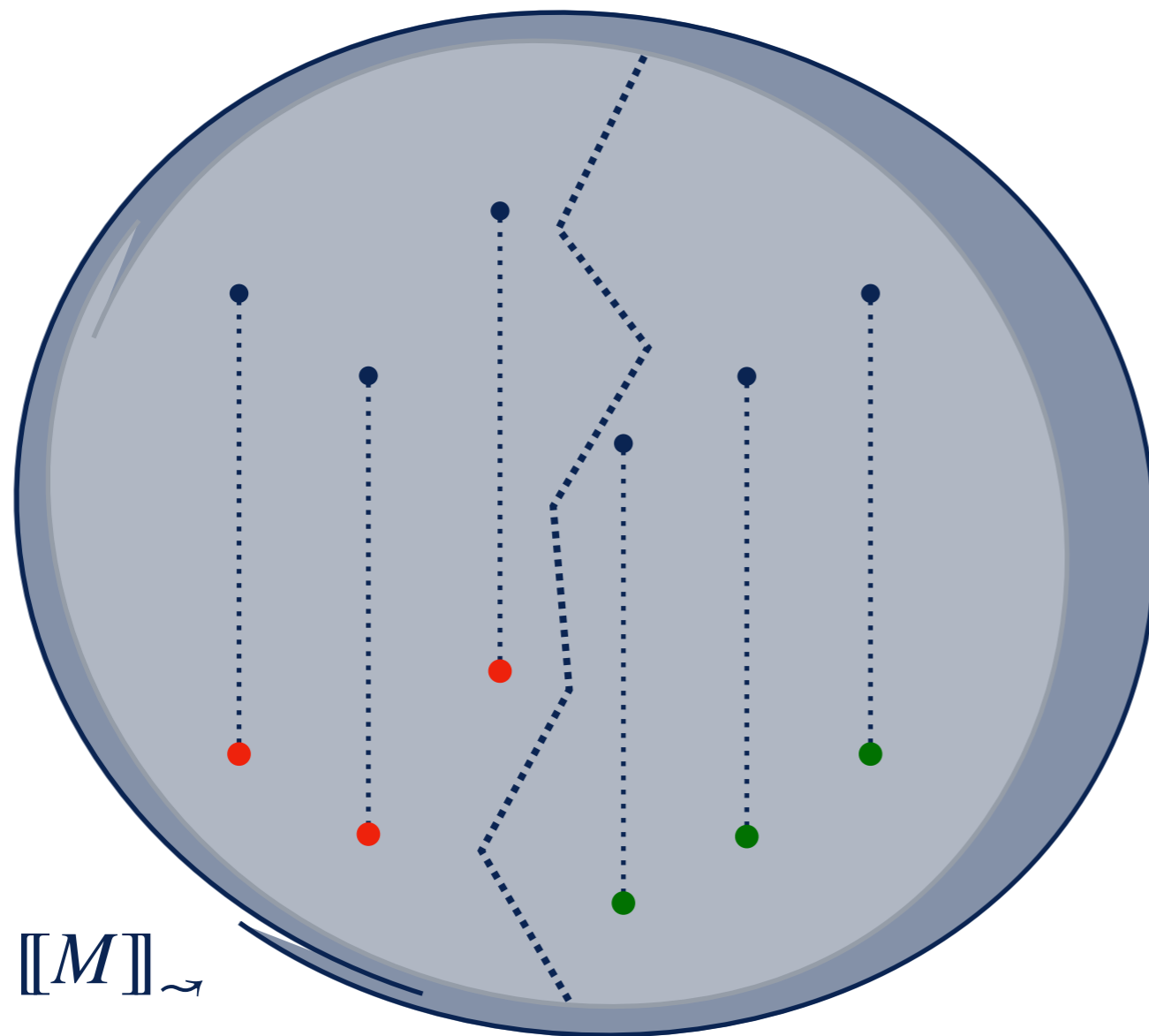
Theorem

$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \} \subseteq \mathcal{N}_J \Leftrightarrow \alpha_\cdot(\{ \llbracket P \rrbracket \}) \subseteq \mathcal{N}_J \Leftrightarrow \llbracket P \rrbracket_\cdot \subseteq \mathcal{N}_J$$

(Another) Hierarchy of Semantics

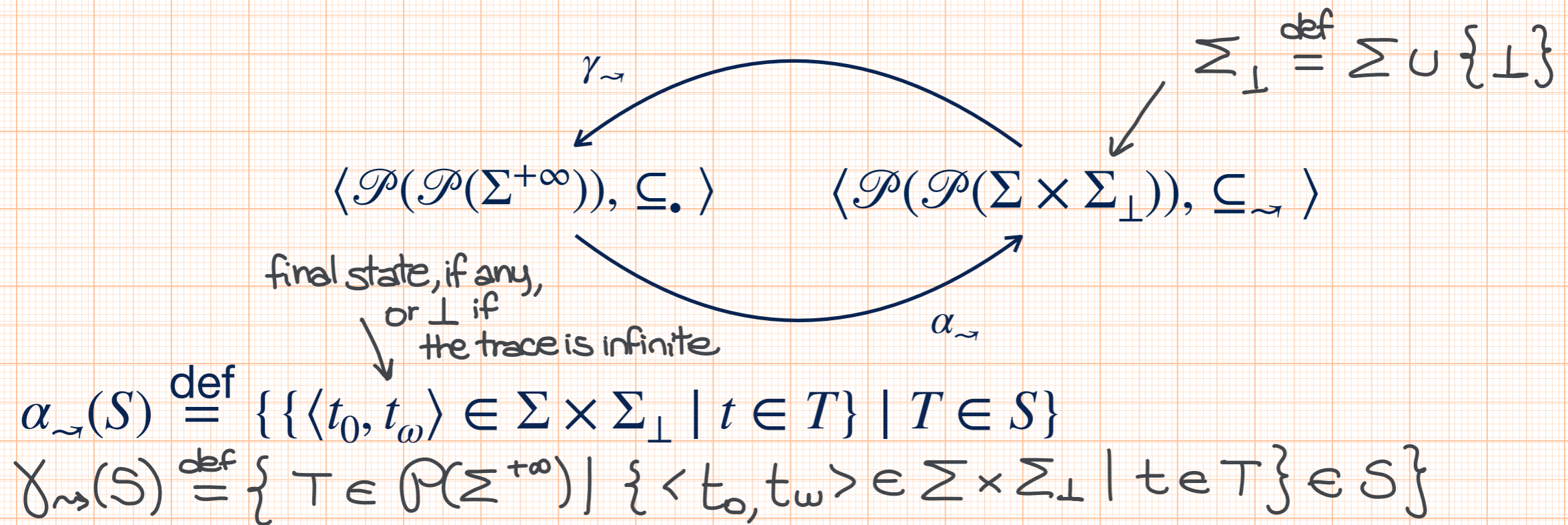


Dependency Semantics

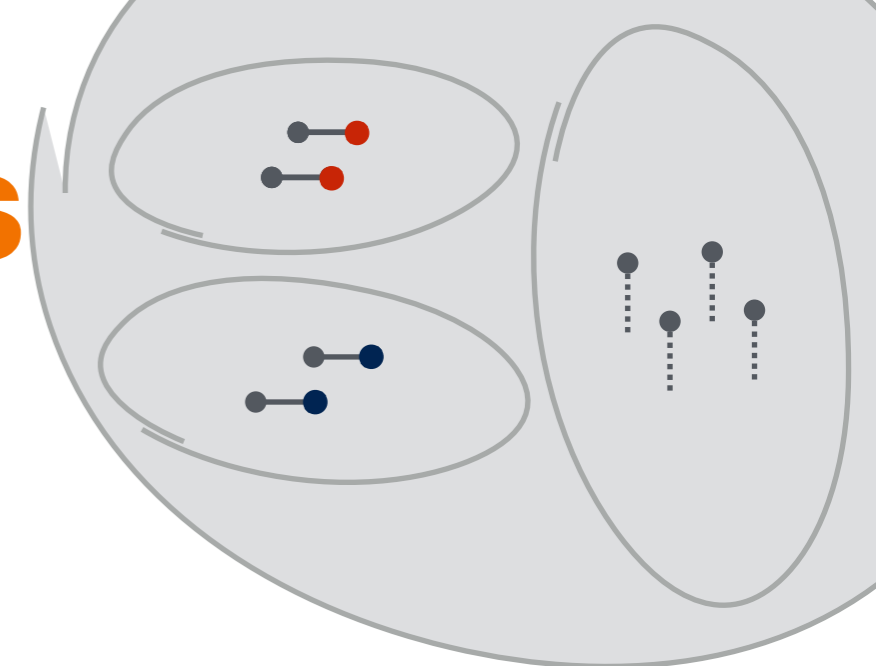


💡 to reason about input data (non-)usage **we do not need to consider all intermediate computations** between the initial and final states of a trace (if any)

Dependency Semantics



Dependency Semantics



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

english → _	english → _
math → T	math → T
science → _	science → _
bonus → _	bonus → _
passing → ?	passing → T

$[[P]] \rightsquigarrow$

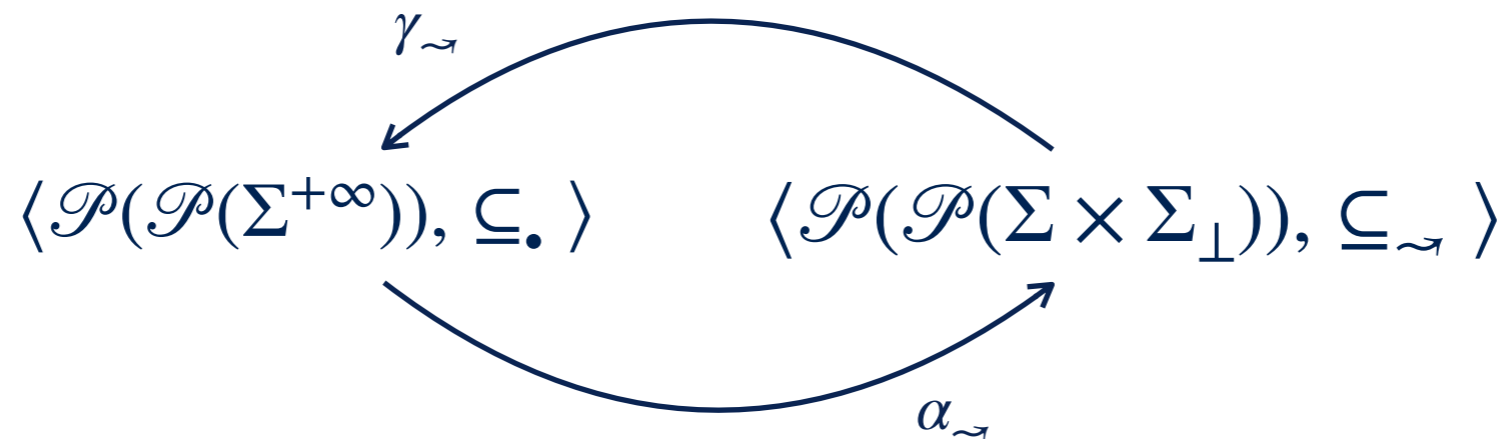
english → _
 math → **F**
 science → _
 bonus → **T**
 passing → ?

english → _
 math → **F**
 science → _
 bonus → **T**
 passing → **T**

english → _
 math → **F**
 science → _
 bonus → **F**
 passing → ?

english → _
 math → **F**
 science → _
 bonus → **F**
 passing → **F**

Dependency Semantics



$$\alpha_{\sim}(S) \stackrel{\text{def}}{=} \{ \{ \langle t_0, t_{\omega} \rangle \in \Sigma \times \Sigma_{\perp} \mid t \in T \} \mid T \in S \}$$

$$\gamma_{\sim}(S) \stackrel{\text{def}}{=} \{ T \in \mathcal{P}(\Sigma^{+\infty}) \mid \{ \langle t_0, t_{\omega} \rangle \in \Sigma \times \Sigma_{\perp} \mid t \in T \} \in S \}$$

$$\llbracket P \rrbracket_{\sim} \stackrel{\text{def}}{=} \alpha_{\sim}(\llbracket P \rrbracket \cdot) = \{ \{ \langle t_0, t_{\omega} \rangle \in \Sigma \times \Sigma \mid t \in \llbracket P \rrbracket \cap O \} \mid O \in \mathbb{O} \}$$

Dependency Semantics

Outcome Semantics

Theorem 1. The outcome semantics $\Lambda_\bullet \in \mathcal{P}(\mathcal{P}(\Sigma^{+\infty}))$ can be expressed as a least fixpoint in $\langle \mathcal{P}(\mathcal{P}(\Sigma^{+\infty})), \sqsubseteq, \sqcup, \sqcap, \{\Sigma^\omega, \emptyset\}, \{\emptyset, \Sigma^+\} \rangle$ as:

$$\Lambda_\bullet = \text{lfp}_{\emptyset}^{\sqsubseteq} \Theta_\bullet$$

$$\Theta_\bullet(S) \stackrel{\text{def}}{=} \{\Omega_{o_1=v_1, \dots, o_k=v_k} \mid v_1, \dots, v_k \in V\} \cup \{\tau; T \mid T \in S\}$$

where $S_1 \sqcup S_2 \stackrel{\text{def}}{=} \{S_{1_{o_1=v_1, \dots, o_k=v_k}} \mid v_1, \dots, v_k \in V\} \cup \{\tau; T \mid T \in S\}$ (9)

the outcome semantics Λ_\bullet can be equivalently expressed as follows:

$$\Lambda_\bullet = \Lambda_\bullet^+ \cup \Lambda_\bullet^\omega = \text{lfp}_{\emptyset}^{\sqsubseteq} \Theta_\bullet^+ \cup \text{lfp}_{\{\Sigma^\omega\}}^{\sqsubseteq} \Theta_\bullet^\omega$$

$$\Theta_\bullet^+(S) \stackrel{\text{def}}{=} \{\Omega_{o_1=v_1, \dots, o_k=v_k} \mid v_1, \dots, v_k \in V\} \cup \{\tau; T \mid T \in S\} \quad (13)$$

$$\Theta_\bullet^\omega(S) \stackrel{\text{def}}{=} \{\tau; T \mid T \in S\}$$

Lemma 2. The abstraction $\Lambda_{\rightsquigarrow}^+ \stackrel{\text{def}}{=} \alpha_{\rightsquigarrow}(\Lambda_\bullet^+) \in \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma))$ can be expressed as a least fixpoint in $\langle \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma_\perp)), \sqsubseteq, \sqcup, \sqcap, \{\Sigma \times \{\perp\}, \emptyset\}, \{\emptyset, \Sigma \times \Sigma\} \rangle$ as:

$$\Lambda_{\rightsquigarrow}^+ = \text{lfp}_{\{\emptyset\}}^{\sqsubseteq} \Theta_{\rightsquigarrow}^+$$

$$\Theta_{\rightsquigarrow}^+(S) \stackrel{\text{def}}{=} \{\Omega_{o_1=v_1, \dots, o_k=v_k} \times \Omega_{o_1=v_1, \dots, o_k=v_k} \mid v_1, \dots, v_k \in V\} \cup \{\tau \circ R \mid R \in S\}$$

Lemma 3. The abstraction $\Lambda_{\rightsquigarrow}^\omega \stackrel{\text{def}}{=} \alpha_{\rightsquigarrow}(\Lambda_\bullet^\omega) \in \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma))$ can be expressed as a least fixpoint in $\langle \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma_\perp)), \sqsubseteq, \sqcup, \sqcap, \{\Sigma \times \{\perp\}, \emptyset\}, \{\emptyset, \Sigma \times \Sigma\} \rangle$ as:

Proof (Sketch). By

$$\Lambda_{\rightsquigarrow}^\omega = \text{lfp}_{\{\Sigma \times \{\perp\}\}}^{\sqsubseteq} \Theta_{\rightsquigarrow}^\omega$$

$$\Theta_{\rightsquigarrow}^\omega(S) \stackrel{\text{def}}{=} \{\tau \circ R \mid R \in S\} \quad (15)$$

Proof (Sketch). By Tarskian fixpoint transfer (cf. Theorem 18 in [12]). \square

Dependency Semantics

Lemma 2. The abstraction $\Lambda_{\rightsquigarrow}^+ \stackrel{\text{def}}{=} \alpha_{\rightsquigarrow}(\Lambda_{\bullet}^+) \in \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma))$ can be expressed as a least fixpoint in $\langle \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma_{\perp})), \sqsubseteq, \sqcup, \sqcap, \{\Sigma \times \{\perp\}, \emptyset\}, \{\emptyset, \Sigma \times \Sigma\} \rangle$ as:

$$\Lambda_{\rightsquigarrow}^+ = \text{lfp}_{\{\emptyset\}}^{\sqsubseteq} \Theta_{\rightsquigarrow}^+$$

$$\Theta_{\rightsquigarrow}^+(S) \stackrel{\text{def}}{=} \{\Omega_{o_1=v_1, \dots, o_k=v_k} \times \Omega_{o_1=v_1, \dots, o_k=v_k} \mid v_1, \dots, v_k \in V\} \cup \{\tau \circ R \mid R \in S\}$$

Lemma 3. The abstraction $\Lambda_{\rightsquigarrow}^{\omega} \stackrel{\text{def}}{=} \alpha_{\rightsquigarrow}(\Lambda_{\bullet}^{\omega}) \in \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma))$ can be expressed as a least fixpoint in $\langle \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma_{\perp})), \sqsubseteq, \sqcup, \sqcap, \{\Sigma \times \{\perp\}, \emptyset\}, \{\emptyset, \Sigma \times \Sigma\} \rangle$ as:

Proof (Sketch). By

$$\begin{aligned} \Lambda_{\rightsquigarrow}^{\omega} &= \text{lfp}_{\{\Sigma \times \{\perp\}\}}^{\sqsubseteq} \Theta_{\rightsquigarrow}^{\omega} \\ \Theta_{\rightsquigarrow}^{\omega}(S) &\stackrel{\text{def}}{=} \{\tau \circ R \mid R \in S\} \end{aligned} \tag{15}$$

Proof (Sketch). By Tarskian fixpoint transfer (cf. Theorem 18 in [12]). \square

Theorem 3. The dependency semantics $\Lambda_{\rightsquigarrow} \in \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma_{\perp}))$ can be expressed as a least fixpoint in $\langle \mathcal{P}(\mathcal{P}(\Sigma \times \Sigma_{\perp})), \sqsubseteq, \sqcup, \sqcap, \{\Sigma \times \{\perp\}, \emptyset\}, \{\emptyset, \Sigma \times \Sigma\} \rangle$ as:

$$\Lambda_{\rightsquigarrow} = \Lambda_{\rightsquigarrow}^+ \cup \Lambda_{\rightsquigarrow}^{\omega} = \text{lfp}_{\{\Sigma \times \{\perp\}, \emptyset\}}^{\sqsubseteq} \Theta_{\rightsquigarrow}$$

$$\Theta_{\rightsquigarrow}(S) \stackrel{\text{def}}{=} \{\Omega_{o_1=v_1, \dots, o_k=v_k} \times \Omega_{o_1=v_1, \dots, o_k=v_k} \mid v_1, \dots, v_k \in V\} \cup \{\tau \circ R \mid R \in S\} \tag{16}$$

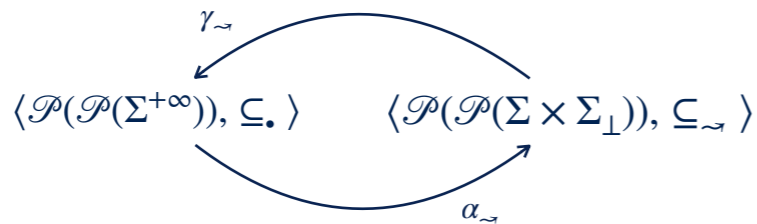
Proof (Sketch). The proof follows immediately from Lemma 2 and Lemma 3. \square

Input Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

\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 $J \subseteq I_P$

Dependency Semantics



$$\alpha_{\sim}(S) \stackrel{\text{def}}{=} \{ \{ \langle t_0, t_w \rangle \in \Sigma \times \Sigma_{\perp} \mid t \in T \} \mid T \in S \}$$

$$\gamma_{\sim}(S) \stackrel{\text{def}}{=} \{ T \in \mathcal{P}(\Sigma^{+\infty}) \mid \{ \langle t_0, t_w \rangle \in \Sigma \times \Sigma_{\perp} \mid t \in T \} \in S \}$$

$$\llbracket P \rrbracket_{\sim} \stackrel{\text{def}}{=} \alpha_{\sim}(\llbracket P \rrbracket \cdot) = \{ \{ \langle t_0, t_w \rangle \in \Sigma \times \Sigma \mid t \in \llbracket P \rrbracket \cap O \} \mid O \in \mathbb{O} \}$$

$$t_0(J) \neq V \Rightarrow \exists t' \in \llbracket P \rrbracket : J \Rightarrow t_0(i) = t'_0(i)$$

Input Data (Non-)Usage

$$\mathcal{N}_J \stackrel{\text{def}}{=} \{ \llbracket P \rrbracket \in \mathcal{P}(\Sigma^{+\infty}) \mid \text{UNUSED}_J(\llbracket P \rrbracket) \}$$

\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 $J \subseteq I_P$

$$t_0(J) \neq V \Rightarrow \exists t' \in \llbracket P \rrbracket : J \Rightarrow t_0(i) = t'_0(i)$$

Theorem

$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \} \subseteq \mathcal{N}_J \Leftrightarrow \llbracket P \rrbracket \cdot \subseteq \mathcal{N}_J \Leftrightarrow \gamma_{\sim}(\llbracket P \rrbracket_{\sim}) \subseteq \mathcal{N}_J$$

Unused Data Analysis

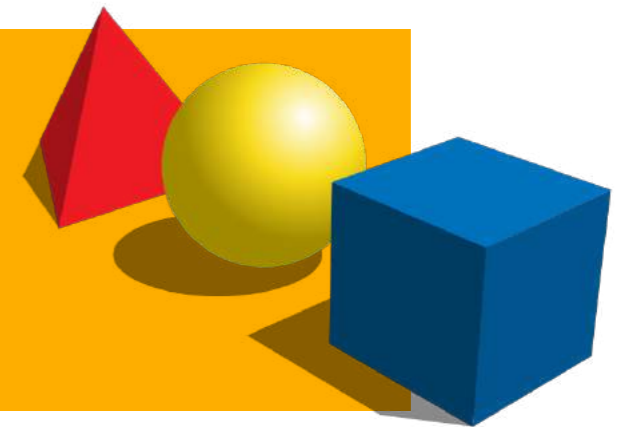
practical tools

targeting specific programs



algorithmic approaches

to decide program properties



mathematical models

of the program behavior

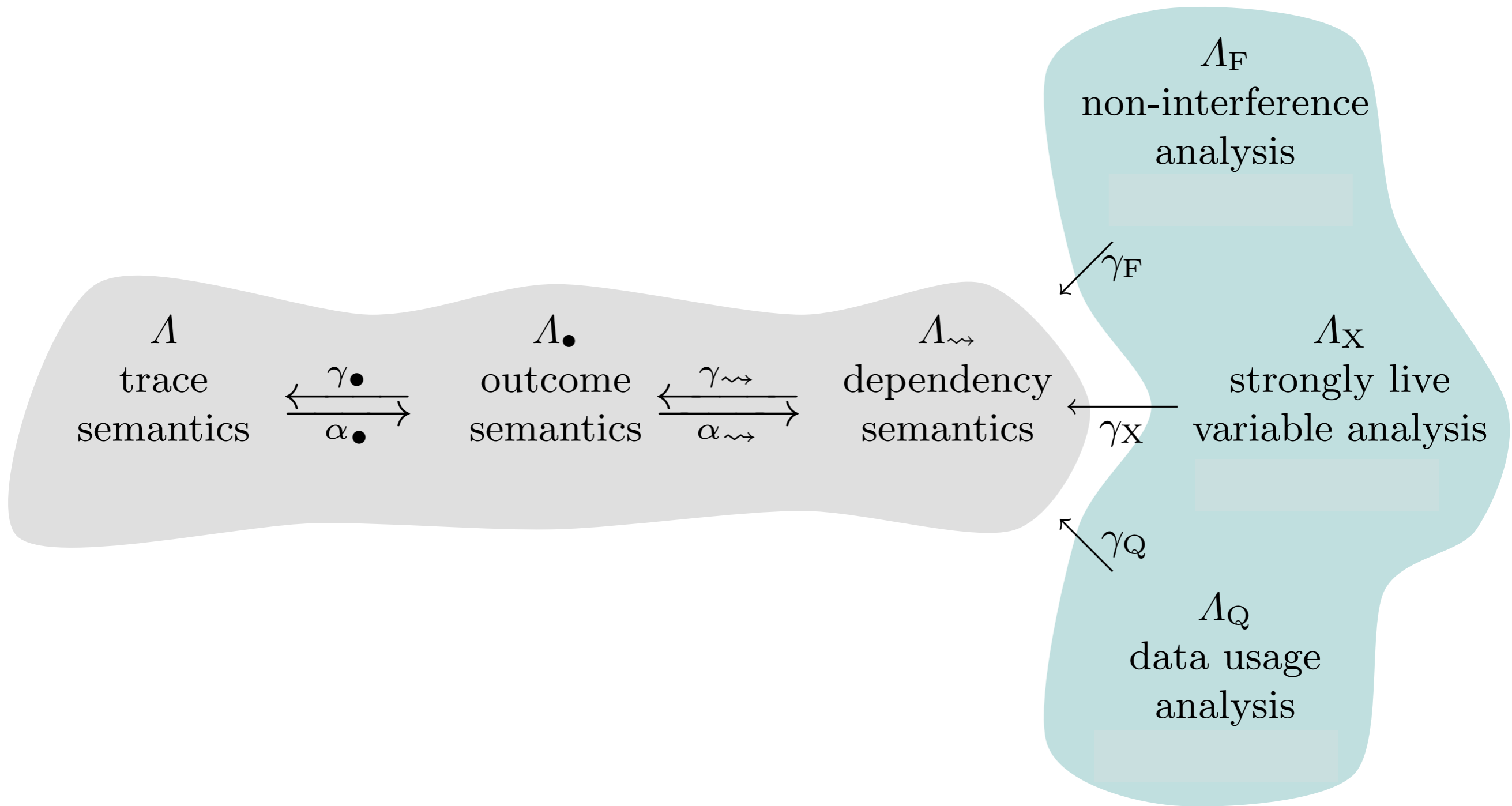


Input Data (Non-)Usage Abstractions

Over-Approximation of the Used Input Data

⇒ **Under-Approximation** of the Unused Input Data

$$P \models \mathcal{N}_{J^{\#} \subseteq J} \Leftarrow \gamma_{\sim}(\gamma_A(\llbracket P \rrbracket_A)) \subseteq \mathcal{N}_{J^{\#} \subseteq J}$$



Secure Information Flow

$L \rightsquigarrow x$
 $L \rightsquigarrow y$
 $H \rightsquigarrow t$
 $L \rightsquigarrow z$
 $H \rightsquigarrow w$

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

implicit usage flows

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

(expressions)
(statements)

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

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

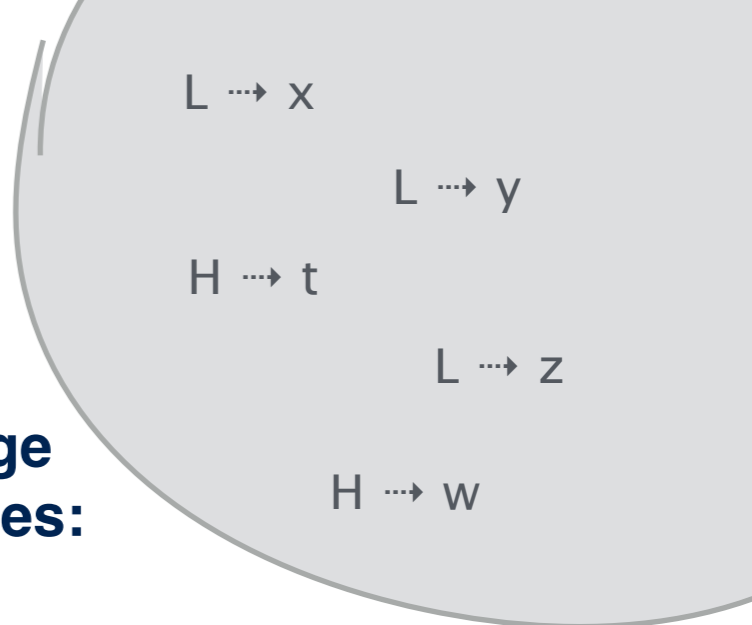
Éric Total
CIDRE, CentraleSupélec, Rennes, FR
first.last@centralesupelec.fr

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

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

program is correct if all its traces satisfy the predicate. By such *trace properties*, extensional definitions of dependencies with more than one trace. To express that the final value of a variable depend only on the initial value of a variable, we use the predicate no_interference .

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**

$\llbracket P \rrbracket_F$

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

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

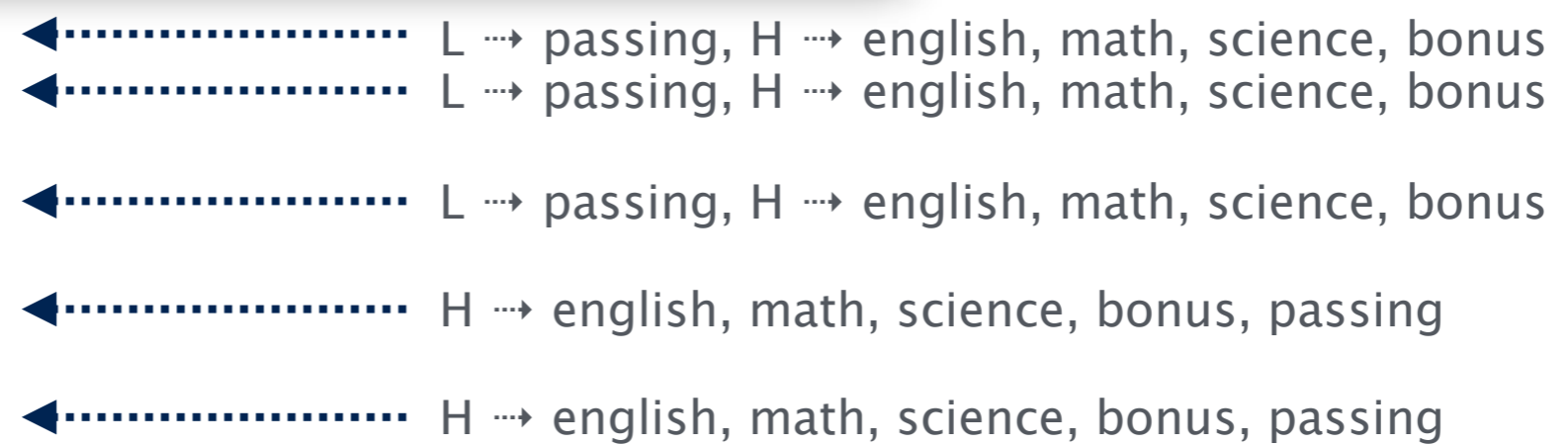
$$\Theta_F[x = e](S) \stackrel{\text{def}}{=} \{L \rightsquigarrow y \in S \mid y \neq x\} \cup \{L \rightsquigarrow x \mid \mathcal{V}_F[e]S\}$$

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

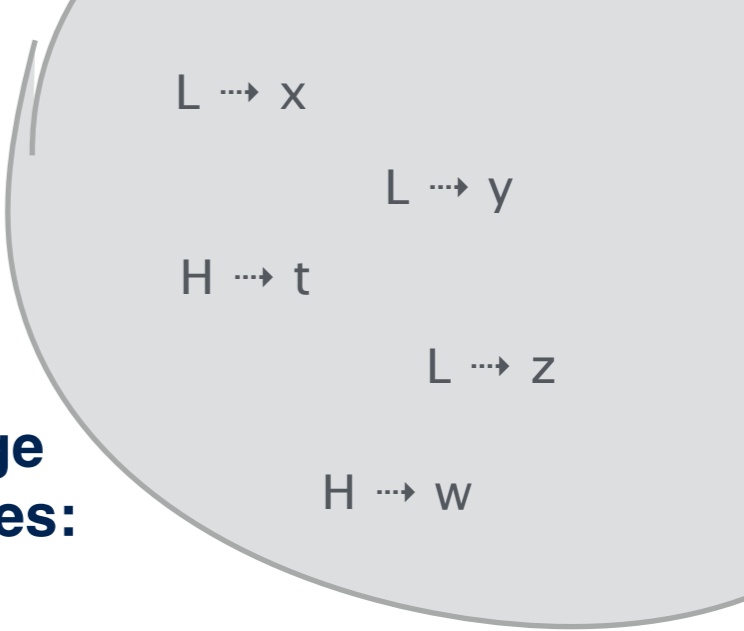
$$\Theta_F[\text{while } e: s](S) \stackrel{\text{def}}{=} \text{lfp}_{S^F}^{\sqsubseteq_F} \Theta_F[\text{if } e: s \text{ else: } \text{skip}]$$

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

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



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)

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

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

$$P \models \mathcal{N}_{I_P}^* \Leftarrow \gamma_{\rightsquigarrow}(\gamma_F(\llbracket P \rrbracket_F)) \subseteq \mathcal{N}_{I_P}^*$$

Strong-Liveness



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

$\llbracket P \rrbracket_X$

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

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

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

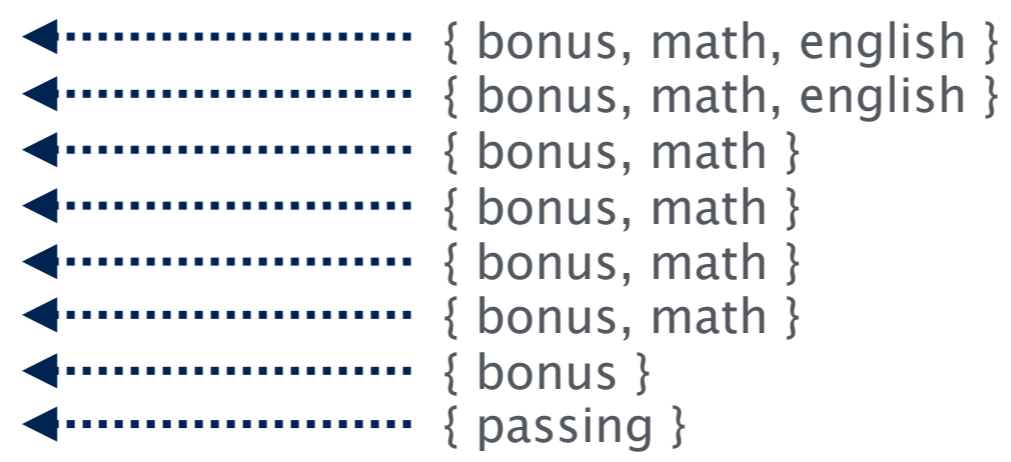
$$\Theta_X[\text{if } b: s_1 \text{ else: } s_2](S) \stackrel{\text{def}}{=} \text{VARS}(b) \cup \Theta_X[s_1](S) \cup \Theta_X[s_2](S)$$

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

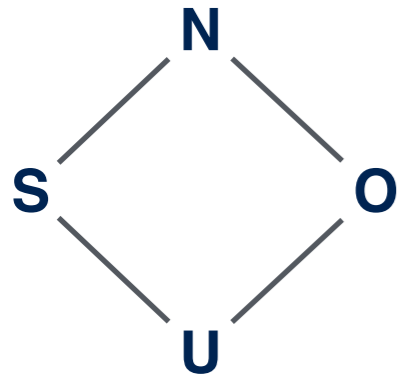
$$\Theta_X[s_1 \ s_2](S) \stackrel{\text{def}}{=} \Theta_X[s_1] \circ \Theta_X[s_2](S)$$

Theorem
 $P \models \mathcal{N}_J \Leftarrow \gamma_{\rightarrow}(\gamma_X(\llbracket P \rrbracket_X)) \subseteq \mathcal{N}_J$

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



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

$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

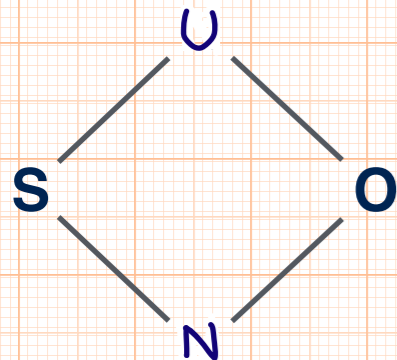
$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$[[P]]_U$

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

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

$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

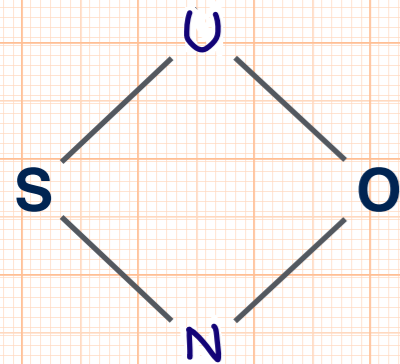
if not math:

passing = **False or** bonus

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

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

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

$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

if not math:

passing = **False or** bonus

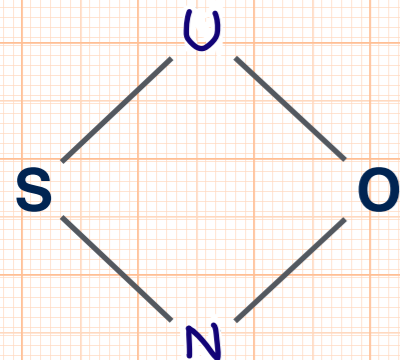


$$\begin{aligned} \Theta_Q[\text{skip}](q) &\stackrel{\text{def}}{=} q \\ \Theta_Q[x = e](q) &\stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q) \\ \Theta_Q[\text{if } b: s_1 \text{ else: } s_2](q) &\stackrel{\text{def}}{=} \text{POP} \circ \text{FILTER}[b] \circ \Theta_Q[s_1] \circ \text{PUSH}(q) \\ &\quad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

$x \mapsto U \Rightarrow x \mapsto S$
 $x \mapsto O \Rightarrow x \mapsto N$

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

$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$

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

• passing = **False or** bonus

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

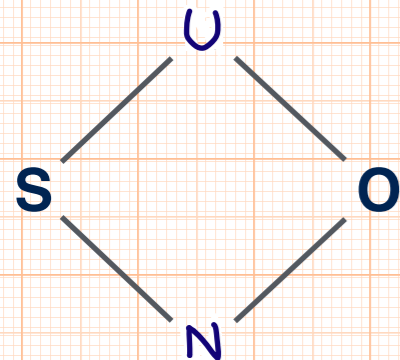
if the assigned variable was used (U or S) it becomes overwritten (O) if not also used in the expression being assigned; otherwise it becomes freshly used (U)

bonus $\mapsto U$, passing $\mapsto O \mid$ passing $\mapsto U$

passing $\mapsto S \mid$ 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

$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

if not math:

passing = **False or** bonus

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

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 \mapsto U

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

passing \mapsto S | passing \mapsto U

passing \mapsto U

Syntactic (Non-)Usage

$x \mapsto U$

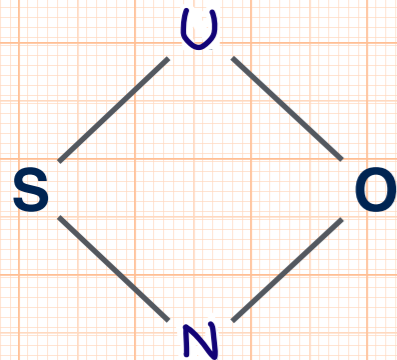
$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$



- **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**

if not english:

english = **False**

if not math:

passing = **False or** bonus

if not math:

passing = **False or** bonus

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

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

math, bonus, passing $\mapsto U$

math, bonus $\mapsto U$, passing $\mapsto O$

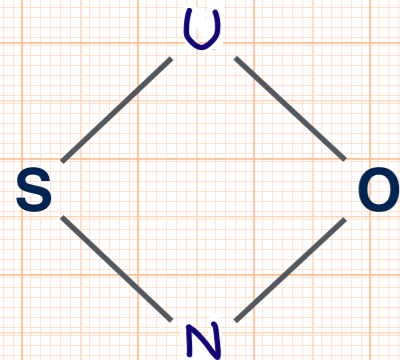
bonus $\mapsto U$, passing $\mapsto O \mid$ passing $\mapsto U$

passing $\mapsto S \mid$ passing $\mapsto U$

passing $\mapsto U$

Syntactic (Non-)Usage

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



- **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

$\llbracket P \rrbracket_U$

passing = **True**

if not english:

english = **False**

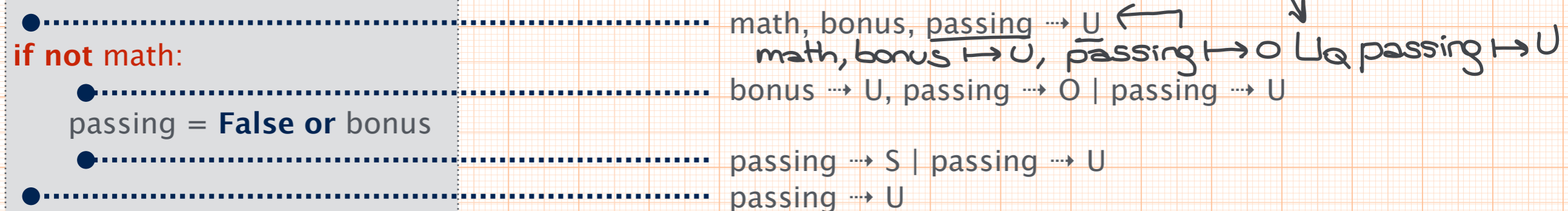
if not math:

passing = **False or** bonus

if not math:

passing = **False or** bonus

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



Syntactic (Non-)Usage

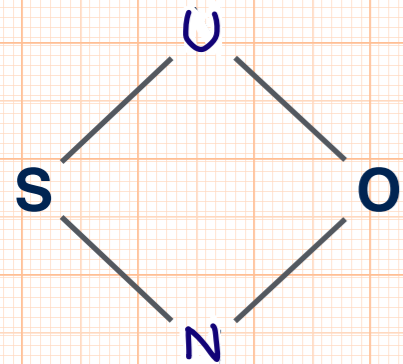
$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$



- **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

$\llbracket P \rrbracket_U$

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

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

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

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

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

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

passing = **True**

if not english:

english = **False**

if not math:

passing = **False or** bonus

if not math:

passing = **False or** bonus

math, bonus, passing $\mapsto S \mid$ math, bonus, passing $\mapsto U$

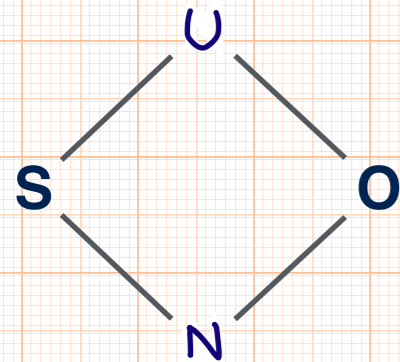
math, bonus, passing $\mapsto U$

bonus $\mapsto U$, passing $\mapsto O \mid$ passing $\mapsto U$

passing $\mapsto S \mid$ 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

$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$

passing = **True**

if not english:

english = **False**

if not math:

• math $\mapsto S$, bonus $\mapsto U$, passing $\mapsto O \mid \dots$
 passing = **False or** bonus

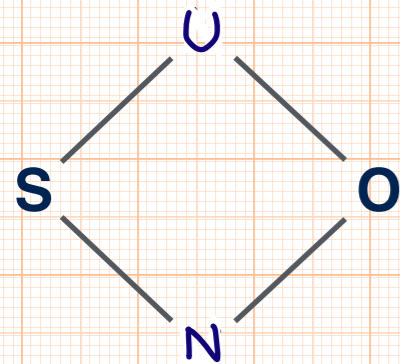
if not math:

• math, bonus, passing $\mapsto S \mid$ math, bonus, passing $\mapsto U$
 • math, bonus, passing $\mapsto U$
 passing = **False or** bonus

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

• bonus $\mapsto U$, passing $\mapsto O \mid$ passing $\mapsto U$
 • passing $\mapsto S \mid$ 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

$x \mapsto U$

$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$

passing = **True**

if not english:

english = **False**

•
if not math:

•
passing = **False or** bonus

•
if not math:

•
passing = **False or** bonus

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

math, bonus, passing $\mapsto U$

math $\mapsto S$, bonus $\mapsto U$, passing $\mapsto O \mid \dots$

math, bonus, passing $\mapsto S \mid$ math, bonus, passing $\mapsto U$

math, bonus, passing $\mapsto U$

bonus $\mapsto U$, passing $\mapsto O \mid$ passing $\mapsto U$

passing $\mapsto S \mid$ passing $\mapsto U$

passing $\mapsto U$

Syntactic (Non-)Usage

$x \mapsto U$

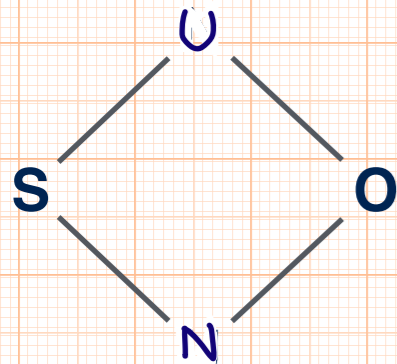
$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$



- **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

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

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

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

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

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

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

passing = **True**

if not english:

english = False



math, bonus, passing $\mapsto S \mid$ math, bonus, passing $\mapsto U$



math, bonus, passing $\mapsto U$

if not math:



math $\mapsto S$, bonus $\mapsto U$, passing $\mapsto O \mid \dots$

passing = **False or** bonus



math, bonus, passing $\mapsto S \mid$ math, bonus, passing $\mapsto U$



math, bonus, passing $\mapsto U$

if not math:



bonus $\mapsto U$, passing $\mapsto O \mid$ passing $\mapsto U$

passing = **False or** bonus



passing $\mapsto S \mid$ passing $\mapsto U$



passing $\mapsto U$

Syntactic (Non-)Usage

$x \mapsto U$

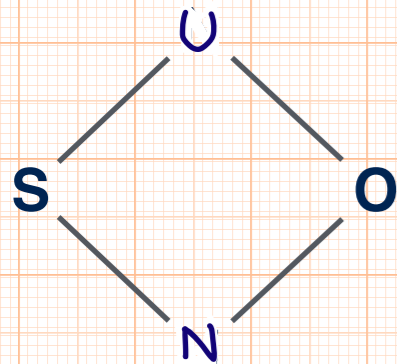
$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

$w \mapsto O \mid w \mapsto U$

$\llbracket P \rrbracket_U$



- **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**

•
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
 $\boxed{\text{math, bonus} \mapsto U, \text{passing} \mapsto O}$

math, bonus, passing $\mapsto U$

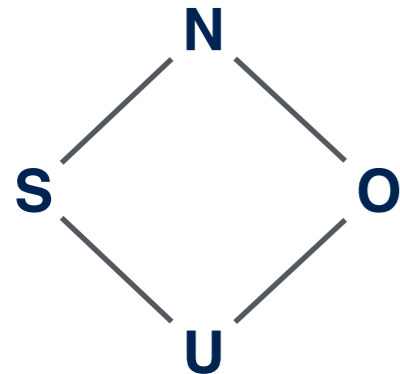
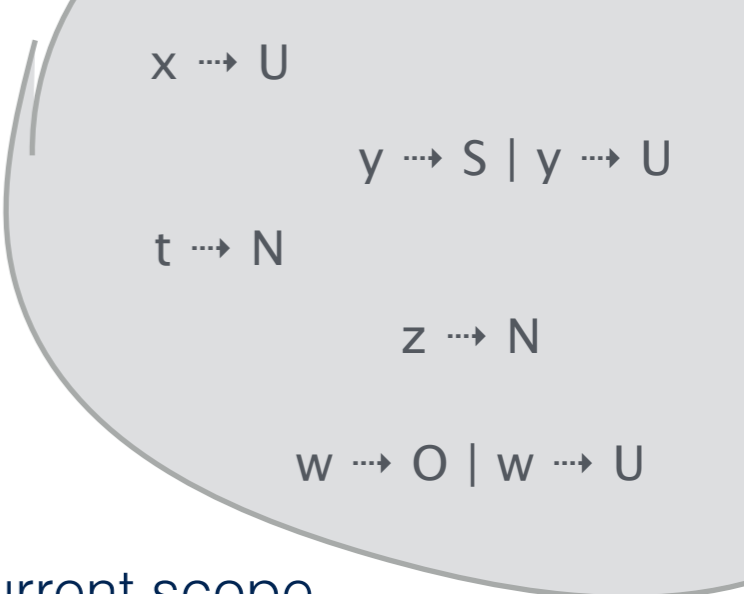
math, bonus, passing $\mapsto S \mid$ math, bonus, passing $\mapsto U$

math, bonus, passing $\mapsto S \mid$ math, bonus, passing $\mapsto U$

math, bonus, passing $\mapsto U$

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

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

$\llbracket P \rrbracket_U$

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

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

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



Theorem

$$P \models \mathcal{N}_J^* \Leftarrow \gamma_{\rightsquigarrow}(\gamma_Q(\llbracket P \rrbracket_Q)) \subseteq \mathcal{N}_J^*$$

Piecewise Syntactic (Non-)Usage

$\{0\} N \{i\} U \{i+1\} O \{len\}$

$\{0\} N \{len\}$

$\{0\} S \{len\} | \{0\} U \{len\}$

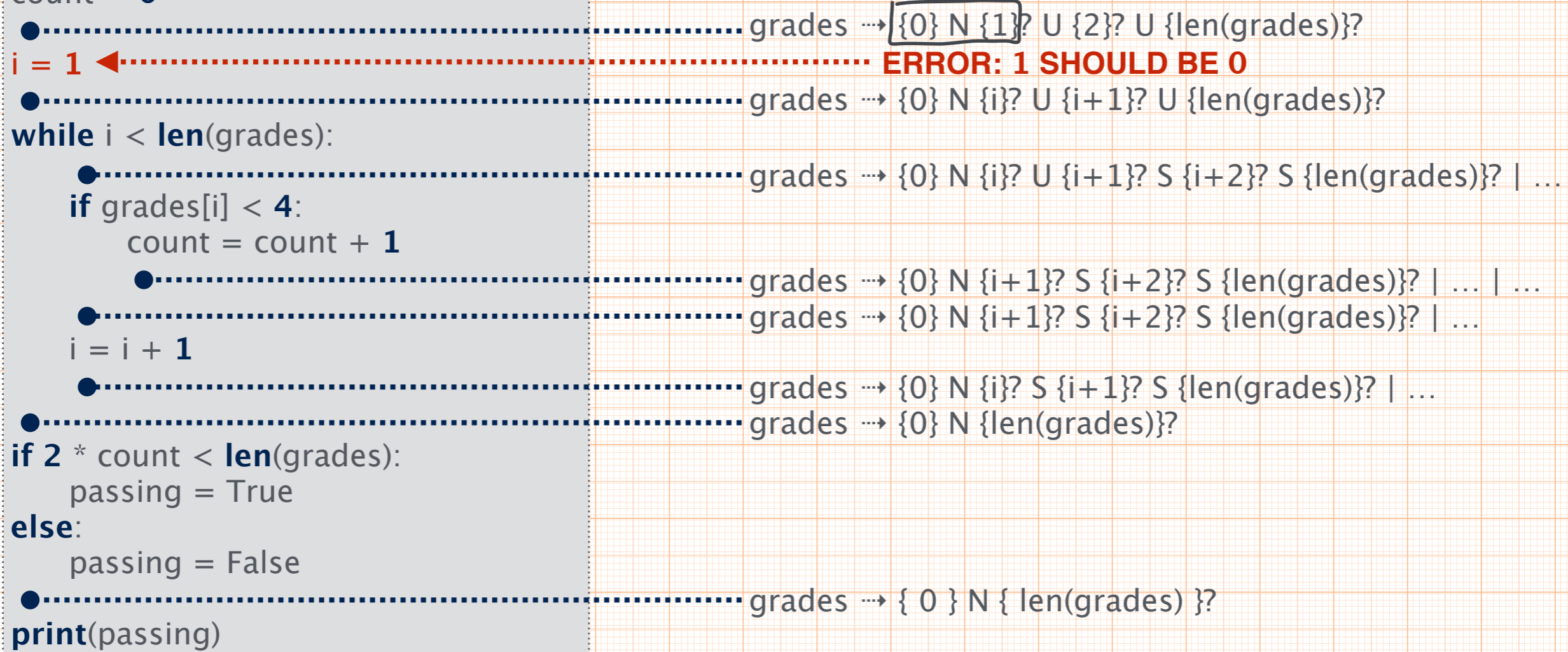
$e ::= v | x | \text{not } e | e \text{ and } e | e \text{ or } e$ $| a[e] | \text{len}(a) | e \oplus e | e \otimes e$ (expressions)
 $s ::= \text{skip} | x = e | \text{if } e: s \text{ else: } s | \text{while } e: s | s s$ $| a[e] = e$ (statements)

$[P]_S$

```

grades = list(map(int, input().split()))
count = 0
i = 1
while i < len(grades):
    if grades[i] < 4:
        count = count + 1
    i = i + 1
if 2 * count < len(grades):
    passing = True
else:
    passing = False
print(passing)
    
```

grades[0] is definitely not used by the program



Unused Data Analysis

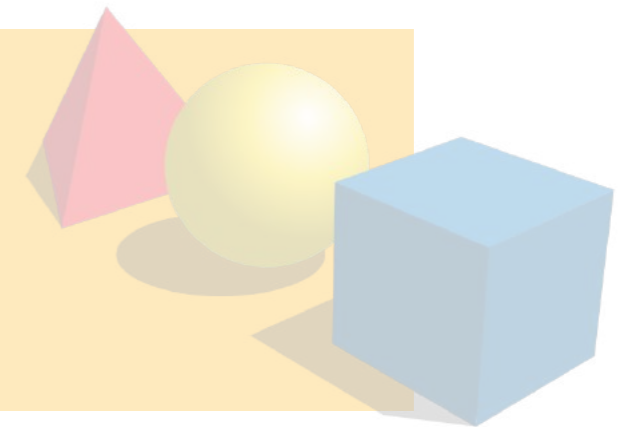
practical tools

targeting specific programs



algorithmic approaches

to decide program properties



mathematical models

of the program behavior



github.com

caterinaurban / Lyra

Type to search

Code Issues 1 Pull requests 1 Actions Projects Wiki Security 6 Insights Settings

Lyra Public

Unpin Unwatch 4 Fork 9 Star 25

master 1 branch 0 tags


Go to file Add file Code

caterinaurban update for Python 3.9 ✓ e37b228 on Nov 7 1,144 commits

docs	documentation	5 years ago
src	update for Python 3.9	last month
.gitignore	[wip] adding .DS_Store mac file	9 months ago
.travis.yml	added fulara unittests to travis	5 years ago
LICENSE	Initial commit	6 years ago
README.md	Merge pull request #78 from caterinaurban/build-status	5 years ago
icon.png	various	6 years ago
lyra.png	logo	6 years ago
requirements.txt	list creation	5 years ago
setup.py	main file	6 years ago

README.md

Lyra - Static Analyzer for Data Science Applications



About

No description or website provided.

python data-science static-analysis abstract-interpretation

Readme MPL-2.0 license Activity 25 stars 4 watching 9 forks

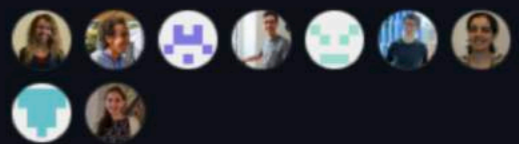
Releases

No releases published [Create a new release](#)

Packages

No packages published [Publish your first package](#)

Contributors 9



Deployments 188

Jupyter Notebooks

The top 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
------	----	----	----

The middle image shows a 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 use of notebooks in production environments.

The bottom image shows the Databricks landing page, which describes it as "An open and unified platform to collaboratively run all types of analytics workloads, from data preparation to exploratory analysis and predictive analytics, at scale." Below the landing page is a screenshot of a Jupyter Notebook interface showing a SQL query and its results.

Jupyter Notebooks



The diagram illustrates a sequence of five Jupyter Notebook cells. The first cell is highlighted with a yellow warning triangle and a yellow banner labeled 'UNUSED DATA'. Orange arrows point from the left to each cell. The code in each cell is as follows:

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

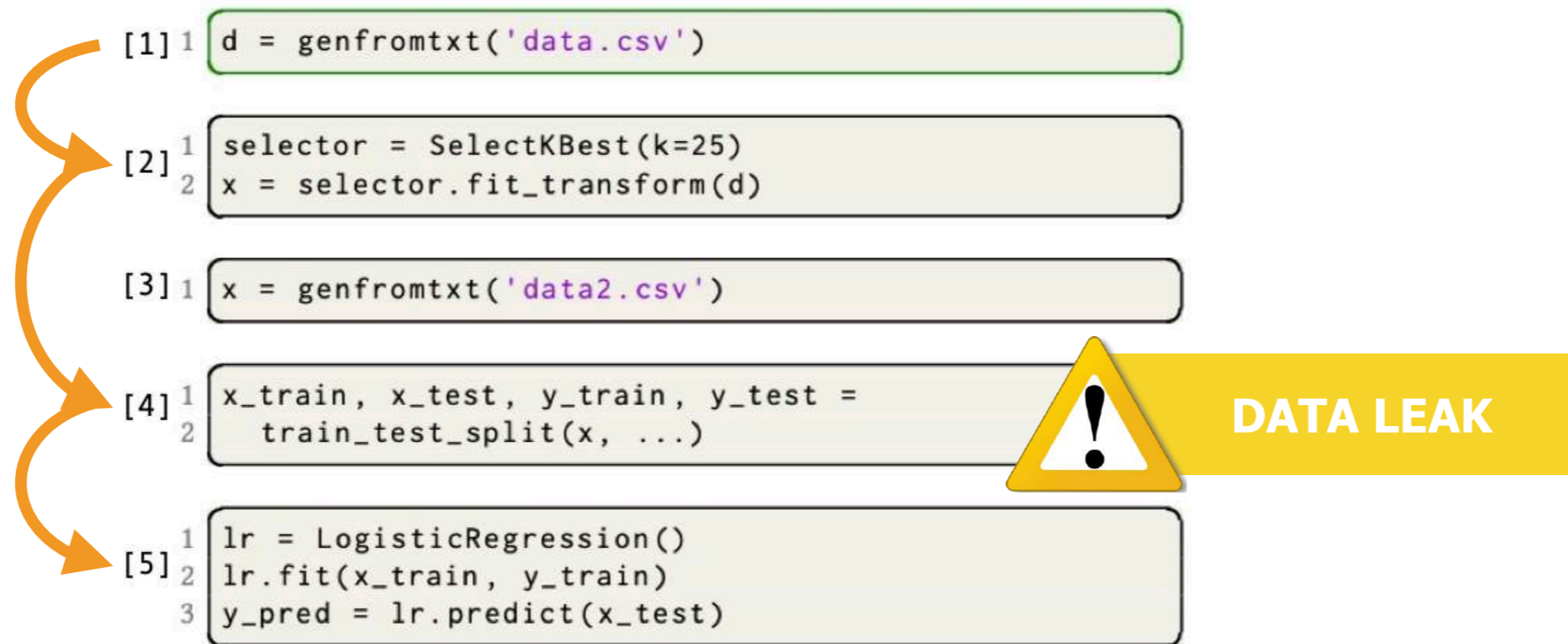
```
[2] 1 selector = SelectKBest(k=25)  
2 x = selector.fit_transform(d)
```

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

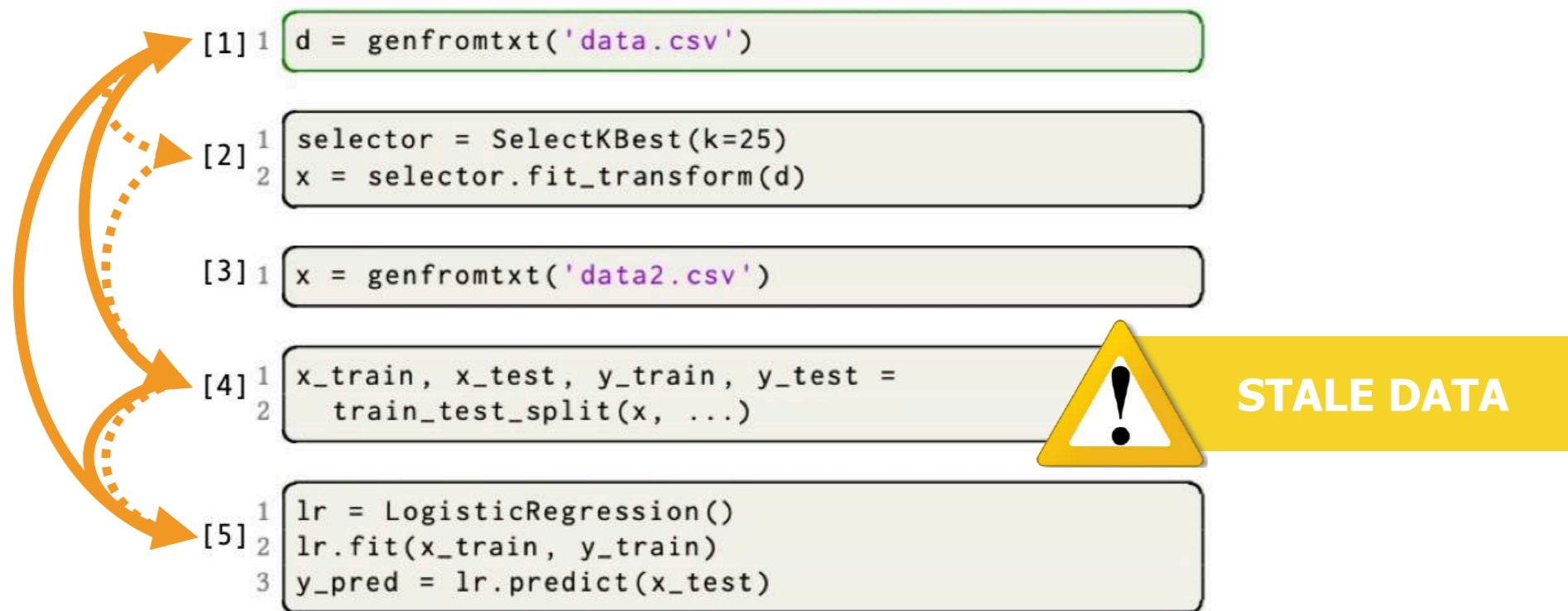
```
[4] 1 x_train, x_test, y_train, y_test =  
2 train_test_split(x, ...)
```

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

Jupyter Notebooks

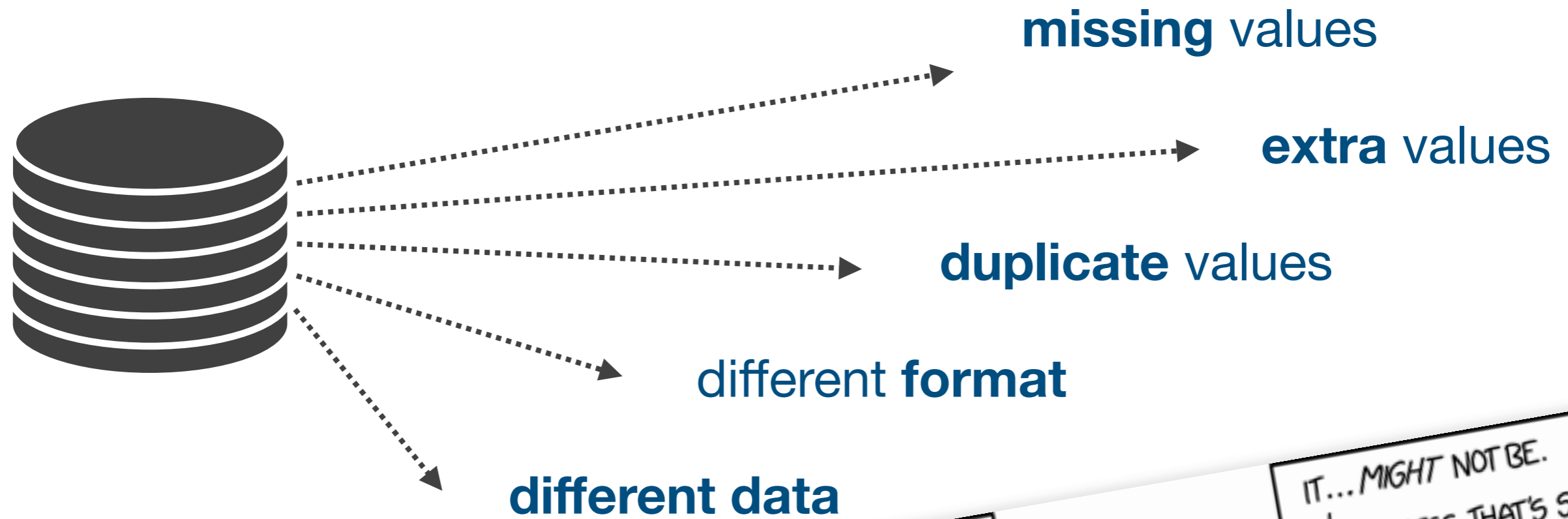


Jupyter Notebooks



Unexpected Data

Unexpected Data



localhost

jupyter Gradebook Last Checkpoint: a few seconds ago (autosaved) Logout

File Edit View Insert Cell Kernel Help Trusted Python 3

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 (autosaved)
```

File Edit View Insert Cell Kernel

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

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

```
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

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

File Edit View Insert Cell Kernel Help

```
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)
```

```
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 (autosaved)

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

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	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)
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 Formatted Values

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

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

	Name	Q1	Q2	Q3
ID				
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)
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]:
```

	Email	Mean
ID		
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: 14 minutes ago (unsaved changes)

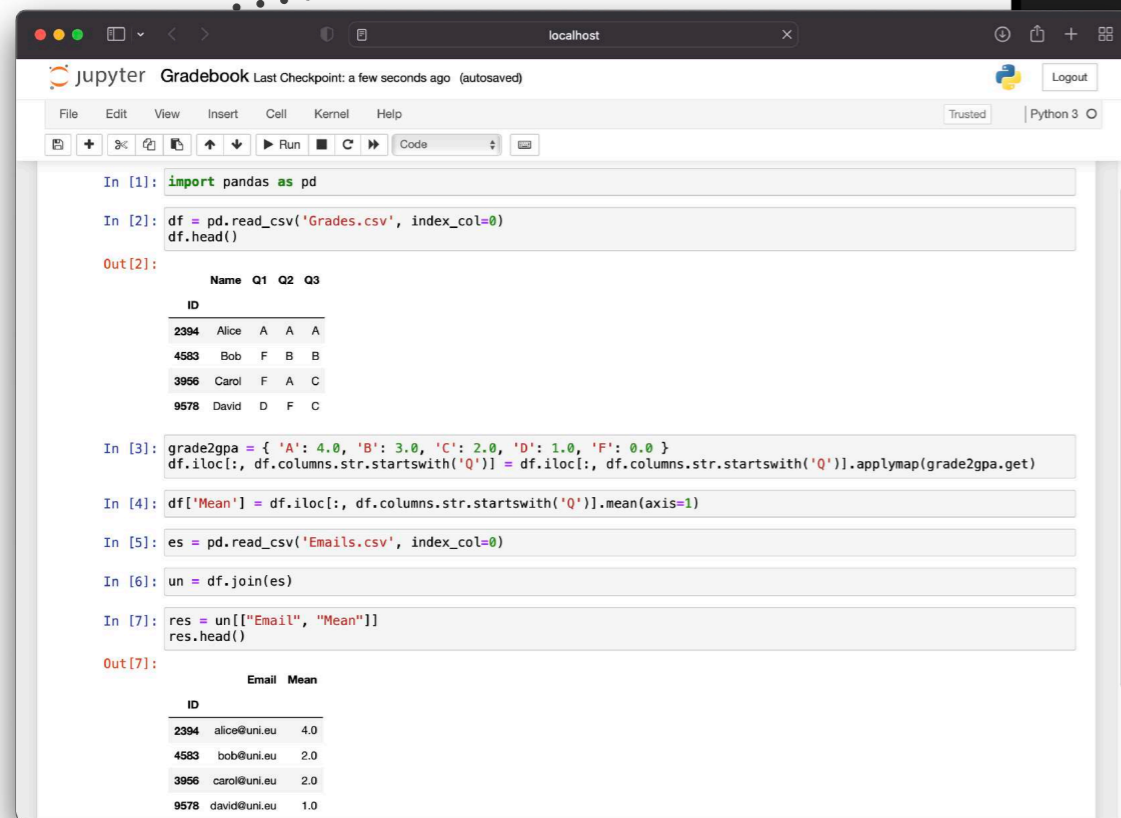
```
In [1]: import pandas as pd
In [2]: df = pd.read_csv('Grades.csv', index_col=0)
df.head()
Out[2]:
```

	Name	Q1	Q2	Q3
ID				
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)
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]:
```

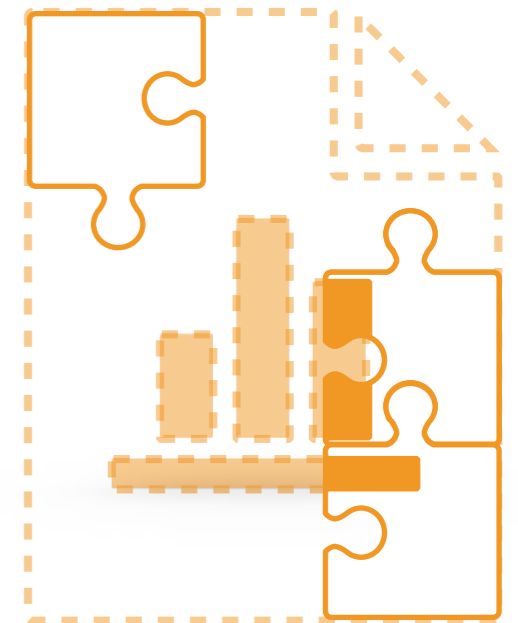
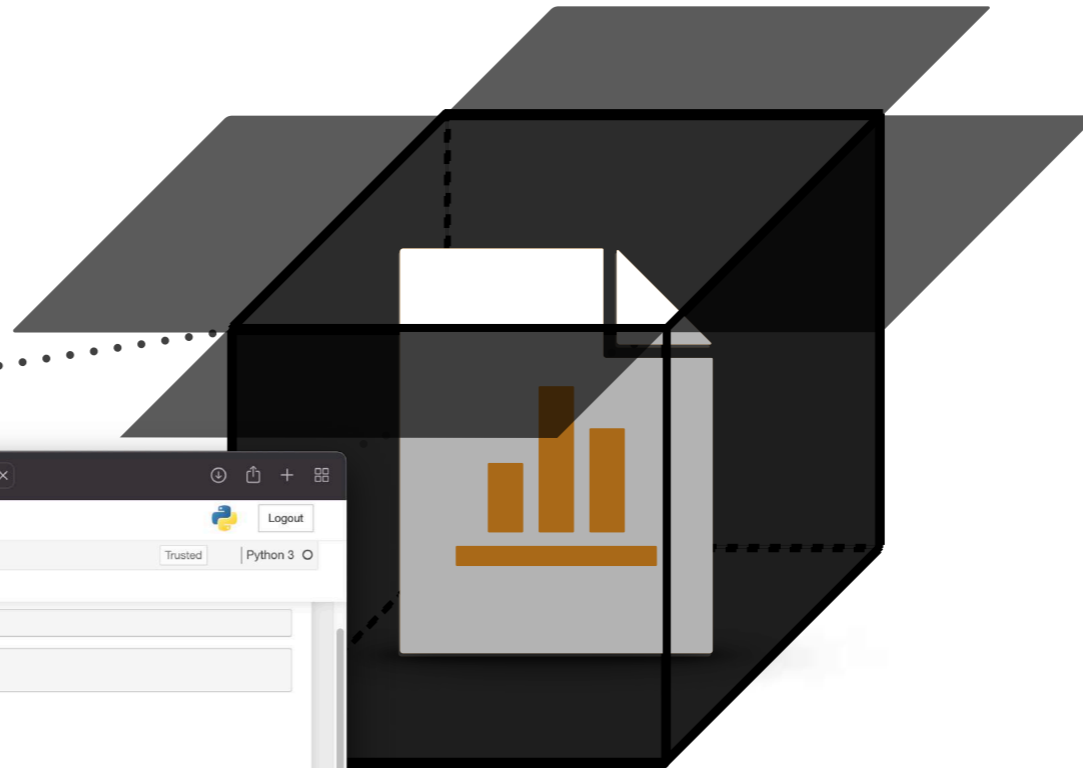
	Email	Mean
ID		
2394	alice@uni.eu	4.0
4583	bob@uni.eu	1.5
3956	carol@uni.eu	2.0
9578	david@uni.eu	1.0

Expected Data Analysis



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



Expected Unused Data Analysis

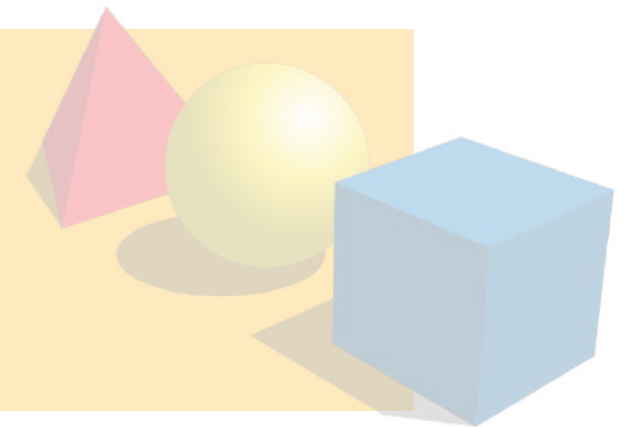
practical tools

targeting specific programs



algorithmic approaches

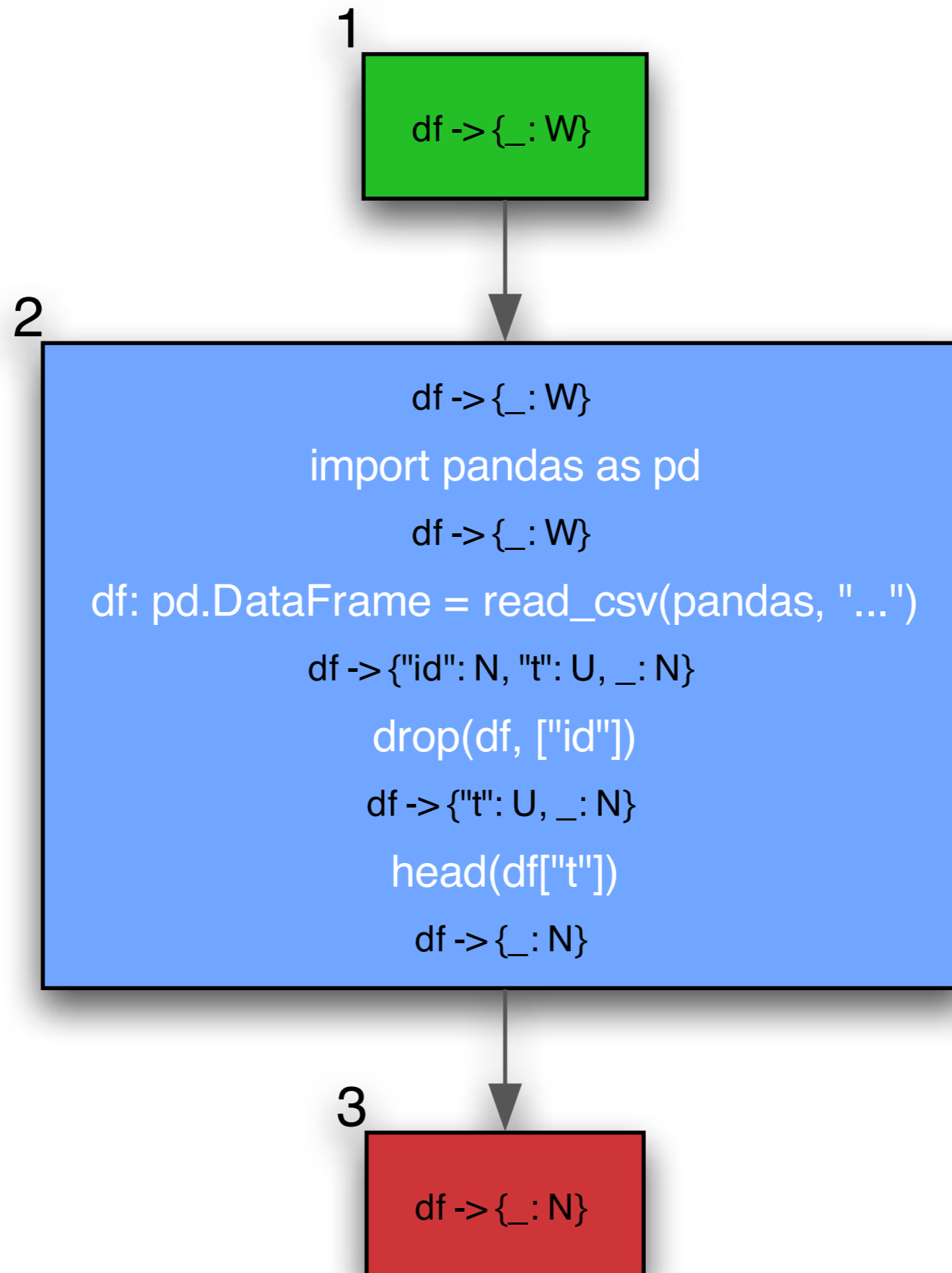
to decide program properties

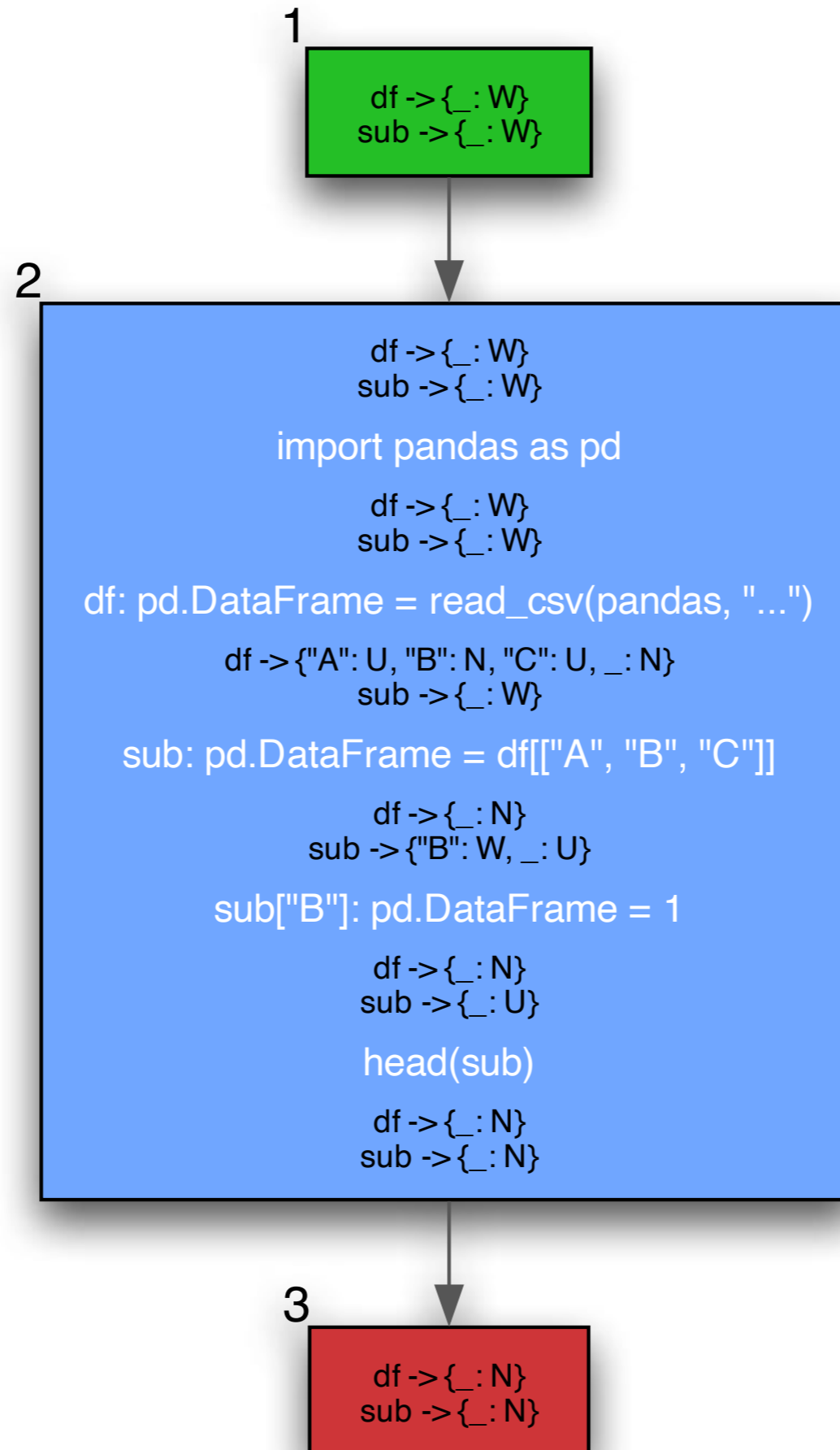


mathematical models

of the program behavior







Bibliography

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

[Assaf17] **Mounir Assaf, David A. Naumann, Julien Signoles, Éric Totel, and Frédéric Tronel.** Hypercollecting Semantics and Its Application to Static Analysis of Information Flow. In POPL, 2017.

[Giegerich81] **Robert Giegerich, Ulrich Möncke, and Reinhard Wilhelm.** Invariance of Approximate Semantics with Respect to Program Transformations. 1981.

[Urban19] **Caterina Urban.** Static Analysis of Data Science Software. In SAS, 2019.

[Urban20] **Caterina Urban.** What Programs Want: Automatic Inference of Input Data Specifications. <https://arxiv.org/abs/2007.10688>, 2020.