

# Static Analysis for Data Science

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



# Data Science is Everywhere

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



mobile devices



data science is **revolutionizing industries**



- personalized recommendations
- targeted marketing



- predictive models
- patient selection



- equipment failure predictions
- internet of things



- predictive models
- customized product offerings



- exploration and discovery
- accident prevention



- personalized treatments
- preventive care

# DATA SCIENCE SOFTWARE



data



data preparation



model training



model deployment



predictions



**TODAY**

**NEXT WEEK**

# 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

American Economic Review: Papers & Proceedings 100 (May 2010): 573–578  
<http://www.aeaweb.org/articles.php?doi=10.1257/aer.100.2.573>

## Growth in a Time of Debt

By CARMEN M. REINHART AND KENNETH S. ROGOFF\*

In this paper, we exploit a new multi-country historical dataset on public (government) debt to search for a systemic relationship between high public debt levels, growth and inflation.<sup>1</sup> Our main result is that whereas the link between growth and debt seems relatively weak at “normal” debt levels, median growth rates for countries with public debt over roughly 90 percent of GDP are about one percent lower than otherwise; average (mean) growth rates are several percent lower. Surprisingly, the relationship between public debt and growth is remarkably similar across emerging markets and advanced economies. This is not the case for inflation. We find no systematic relationship between high debt levels and inflation for advanced economies as a group (albeit with individual country exceptions including the United States). By contrast, in emerging market countries, high public debt levels coincide with higher inflation.

Our topic would seem to be a timely one. Public debt has been soaring in the wake of the recent global financial maelstrom, especially in the epicenter countries. This should not be surprising, given the experience of earlier severe financial crises.<sup>2</sup> Outsized deficits and epic bank bailouts may be useful in fighting a downturn, but what is the long-run macroeconomic impact,

especially against the backdrop of graying populations and rising social insurance costs? Are sharply elevated public debts ultimately a manageable policy challenge?

Our approach here is decidedly empirical, taking advantage of a broad new historical dataset on public debt (in particular, central government debt) first presented in Carmen M. Reinhart and Kenneth S. Rogoff (2008, 2009b). Prior to this dataset, it was exceedingly difficult to get more than two or three decades of public debt data even for many rich countries, and virtually impossible for most emerging markets. Our results incorporate data on 44 countries spanning about 200 years. Taken together, the data incorporate over 3,700 annual observations covering a wide range of political systems, institutions, exchange rate and monetary arrangements, and historic circumstances.

We also employ more recent data on external debt, including debt owed both by governments and by private entities. For emerging markets, we find that there exists a significantly more severe threshold for total gross external debt (public and private)—which is almost exclusively denominated in a foreign currency—than for total public debt (the domestically issued component of which is largely denominated in home currency). When gross external debt reaches 60 percent of GDP, annual growth declines by about two percent; for levels of external debt in excess of 90 percent of GDP, growth rates are roughly cut in half. We are not in a position to calculate separate total external debt thresholds (as opposed to public debt thresholds) for advanced countries. The available time-series is too recent, beginning only in 2000. We do note, however, that external debt levels in advanced countries now average nearly 200 percent of GDP, with external debt levels being particularly high across Europe.

The focus of this paper is on the longer term macroeconomic implications of much higher public and external debt. The final section, how-

	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
28	Maximum		5.4	4.9	10.2	3.6	13.3
29							
30	US	1946-2009	n.a.	3.4	3.3	-2.0	n.a.
31	UK	1946-2009	n.a.	2.4	2.5	2.4	n.a.
32	Sweden	1946-2009	3.6	2.9	2.7	n.a.	6.3
33	Spain	1946-2009	1.5	3.4	4.2	n.a.	9.9
34	Portugal	1952-2009	4.8	2.5	0.3	n.a.	7.9
35	New Zealand	1948-2009	2.5	2.9	3.9	-7.9	2.6
36	Netherlands	1956-2009	4.1	2.7	1.1	n.a.	6.4
37	Norway	1947-2009	3.4	5.1	n.a.	n.a.	5.4
38	Japan	1946-2009	7.0	4.0	1.0	0.7	7.0
39	Italy	1951-2009	5.4	2.1	1.8	1.0	5.6
40	Ireland	1948-2009	4.4	4.5	4.0	2.4	2.9
41	Greece	1970-2009	4.0	0.3	2.7	2.9	13.3
42	Germany	1946-2009	3.9	0.9	n.a.	n.a.	3.2
43	France	1949-2009	4.9	2.7	3.0	n.a.	5.2
44	Finland	1946-2009	3.8	2.4	5.5	n.a.	7.0
45	Denmark	1950-2009	3.5	1.7	2.4	n.a.	5.6
46	Canada	1951-2009	1.9	3.6	4.1	n.a.	2.2
47	Belgium	1947-2009	n.a.	4.2	3.1	2.6	n.a.
48	Austria	1948-2009	5.2	3.3	-3.8	n.a.	5.7
49	Australia	1951-2009	3.2	4.9	4.0	n.a.	5.9
50							
51			4.1	2.8	2.8	-AVERAGE(L30:L44)	

data excluded from the analysis

\*Reinhart: Department of Economics, 4115 Tydings Hall, University of Maryland, College Park, MD 20742 (e-mail: creinhar@umd.edu); Rogoff: Economics Department, 216 Littauer Center, Harvard University, Cambridge MA 02138-3001 (e-mail: krogoff@harvard.edu). The authors would like to thank Olivier Jeanne and Vincent R. Reinhart for helpful comments.

<sup>1</sup> In this paper “public debt” refers to gross central government debt. “Domestic public debt” is government debt issued under domestic legal jurisdiction. Public debt does not include debts carrying a government guarantee. Total gross external debt includes the external debts of all branches of government as well as private debt that is issued by domestic private entities under a foreign jurisdiction.

<sup>2</sup> Reinhart and Rogoff (2009a, b) demonstrate that the aftermath of a deep financial crisis typically involves a protracted period of macroeconomic adjustment, particularly in employment and housing prices. On average, public

# 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

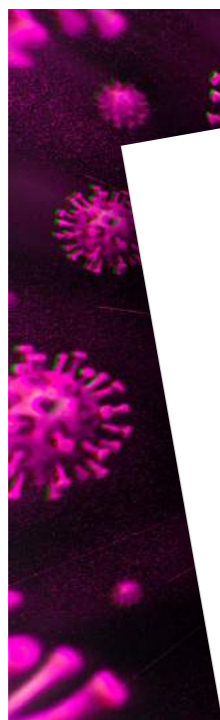
# England Covid-19 Cases Error

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



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/> on 06 October 2020 at 10:13:36 AM GMT+01:00.



# 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

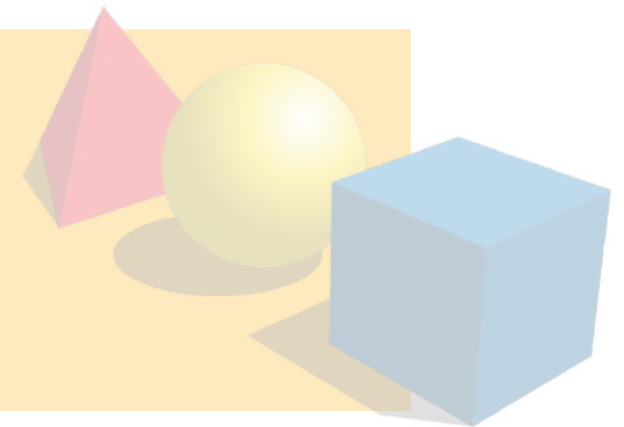
## practical tools

targeting specific programs



## algorithmic approaches

to decide program properties

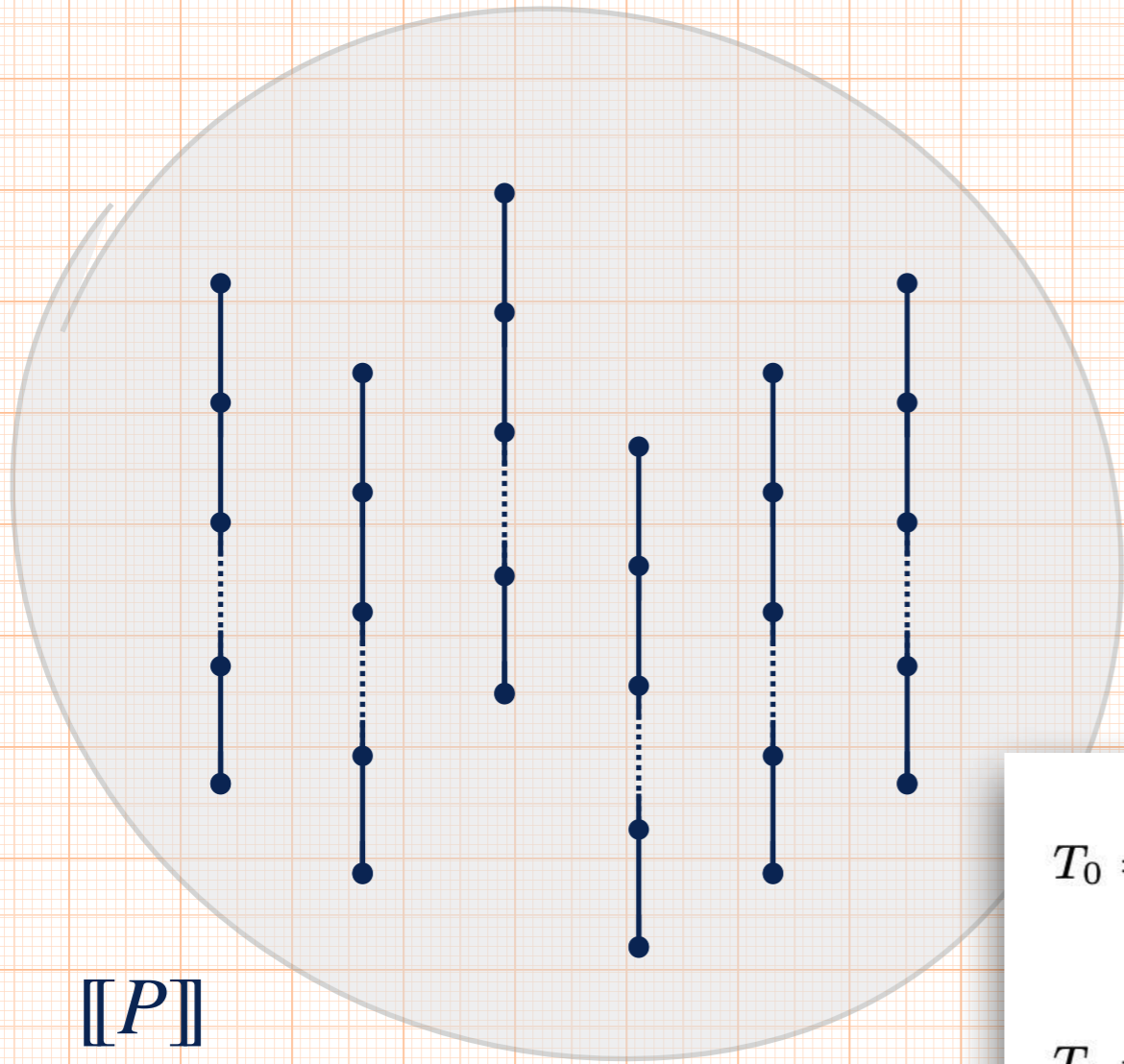


## mathematical models

of the program behavior



# Maximal Trace Semantics



## Least fixpoint formulation of maximal traces

Idea: To get a **least fixpoint** formulation for whole  $\mathcal{M}_\infty$ , 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



[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  $J \subseteq I_P$

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

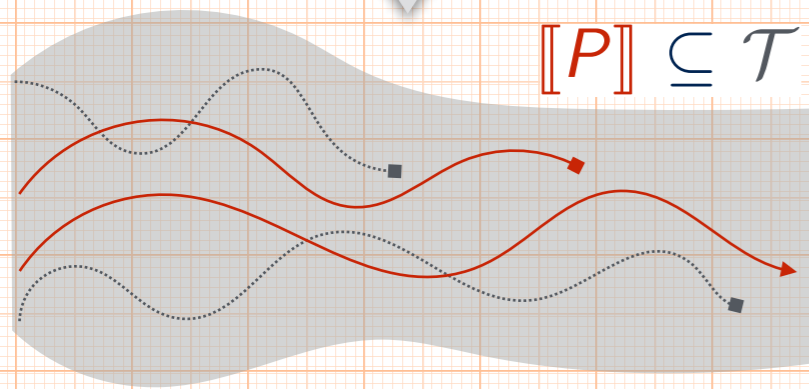
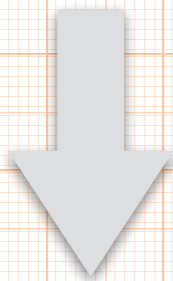
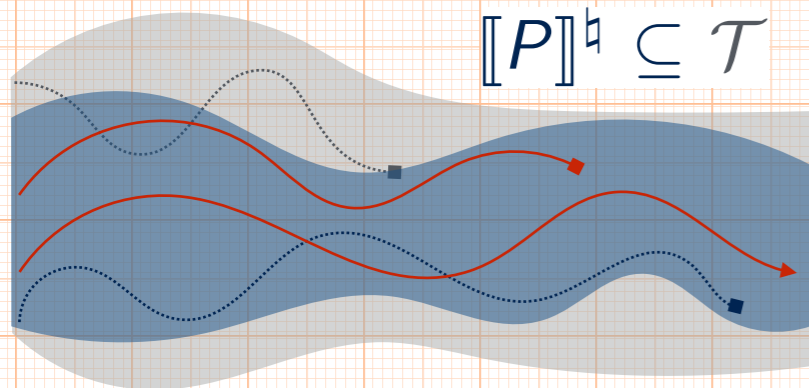
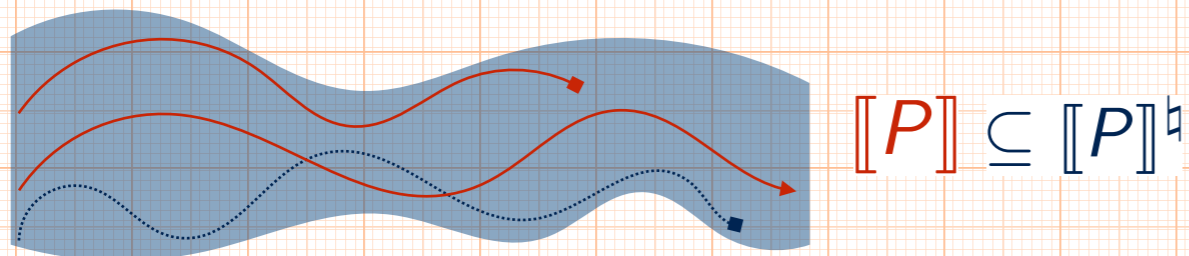
Intuitively: **any possible program outcome** is possible from any value of the input variable  $x_{0,i}$

## 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



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**:  $Col(prog) \stackrel{def}{=} [[prog]]$
- the verification problem remains an inclusion check:  $[[prog]] \subseteq P$
- abstractions will over-approximate the set of traces  $[[prog]]$

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  $b$ )

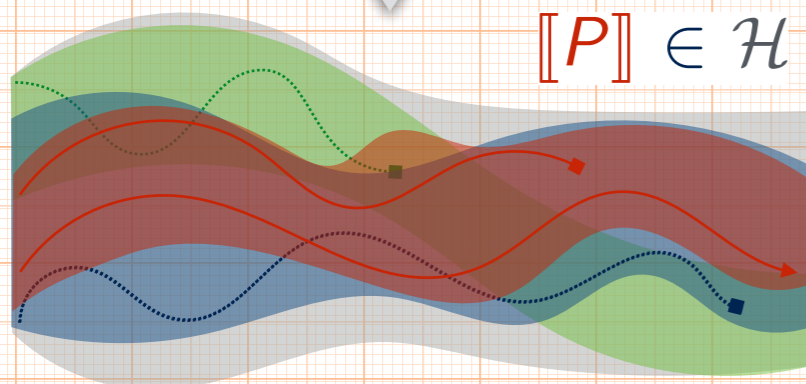
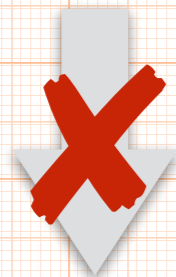
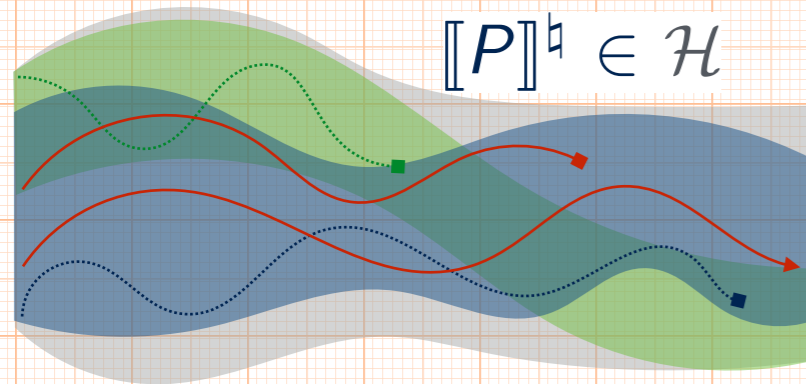
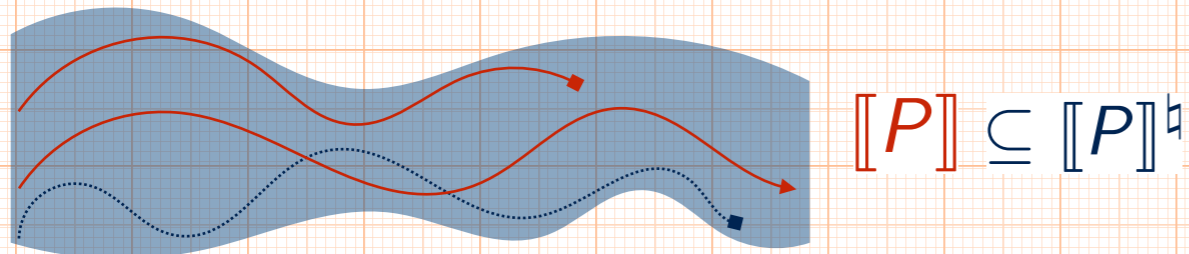
Course 2

Program Semantics and Properties

Antoine Miné

p. 25 / 98

# Program Properties



Collecting semantics and properties

## General properties

General setting:

- given a program  $prog \in Prog$
- its **semantics**:  $[\cdot] : Prog \rightarrow \mathcal{P}(\Sigma^*)$  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)

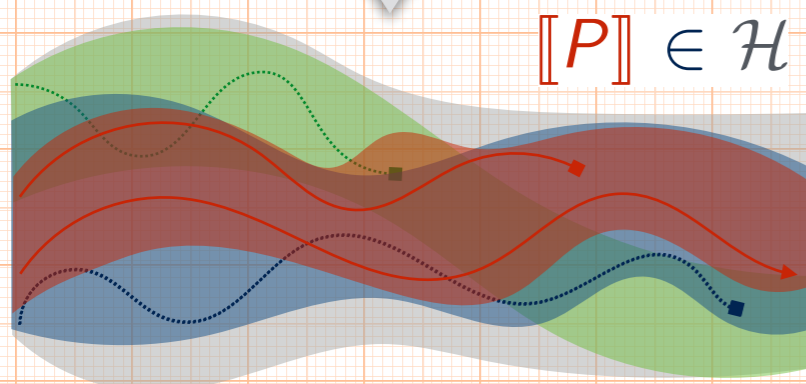
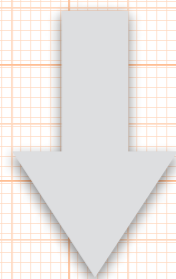
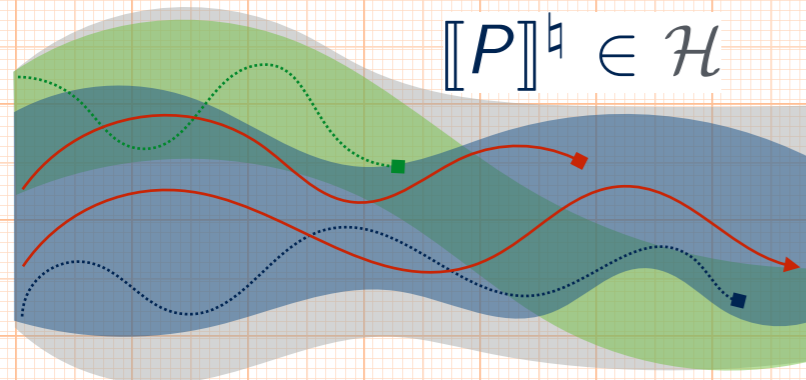
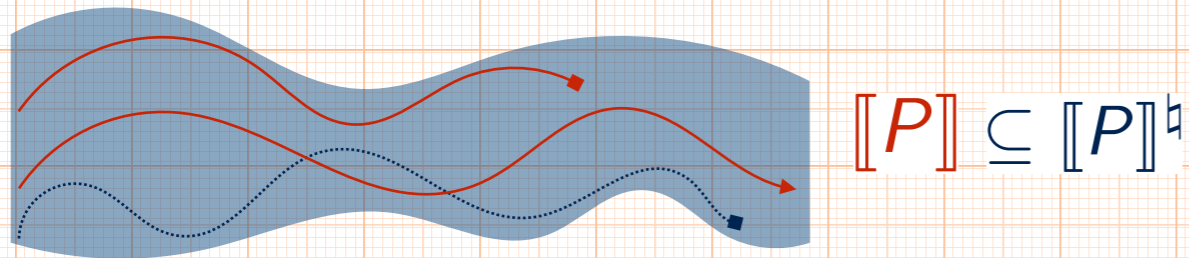
Course 2

Program Semantics and Properties

Antoine Miné

p. 23 / 98

# Subset-Closed Properties



**Program Properties**

$[P] \subseteq [P]^\sharp$

$[P]^\sharp \in \mathcal{H}$

$[P] \in \mathcal{H}$

**General properties**

Collecting semantics and properties

**General setting:**

- given a program  $prog \in Prog$
- its semantics:  $[\cdot] : Prog \rightarrow \mathcal{P}(\Sigma^*)$  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 7

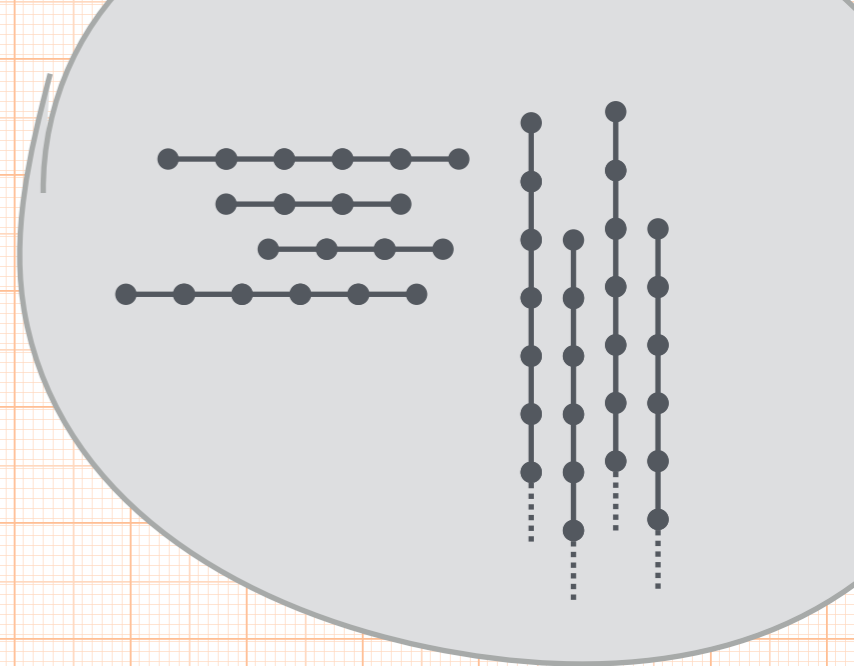
Static Analysis for Data Science

Caterina Urban

15



# 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

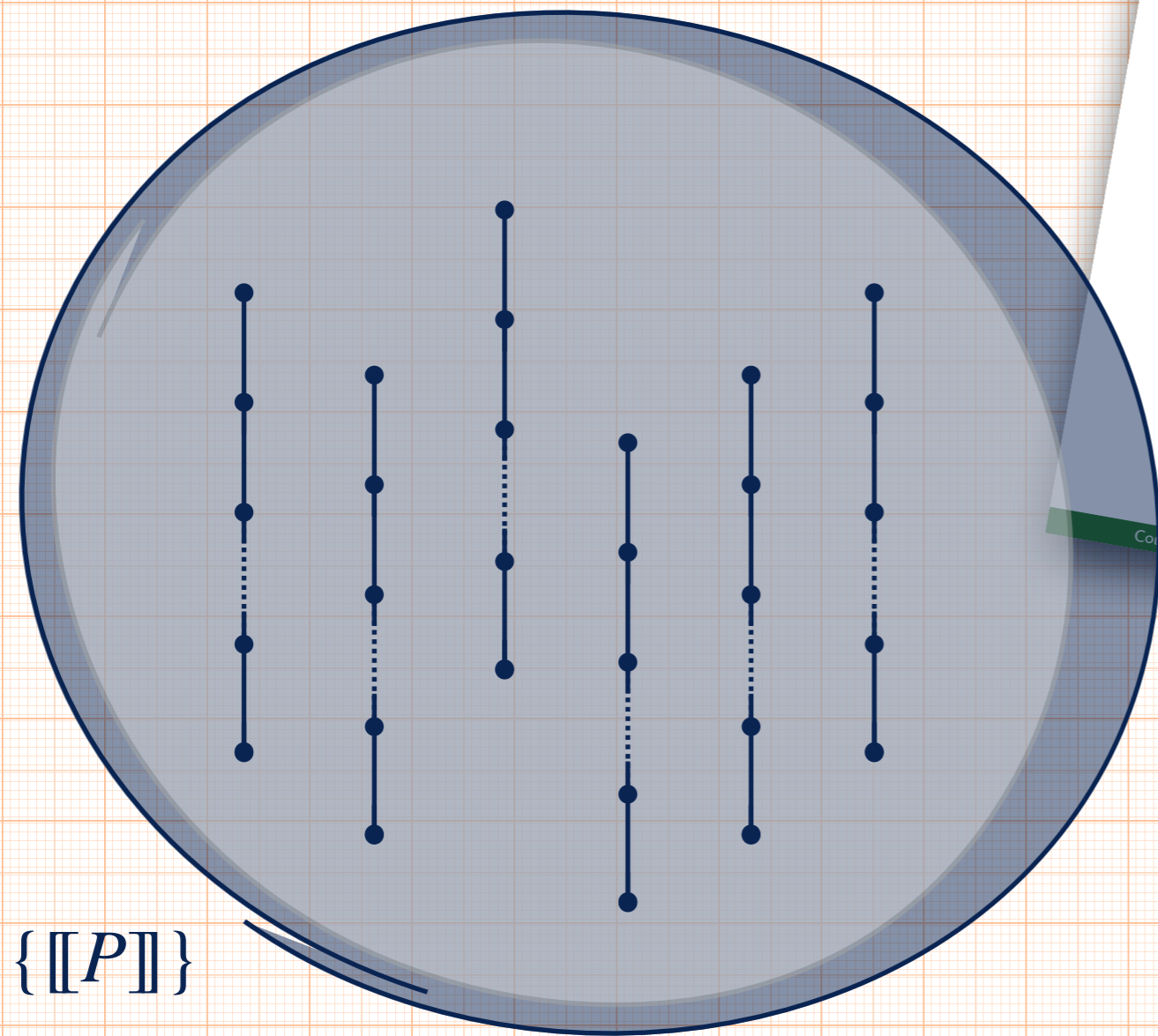
$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \} \subseteq \mathcal{N}_J$$

### 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:  $S^\# \subseteq P$

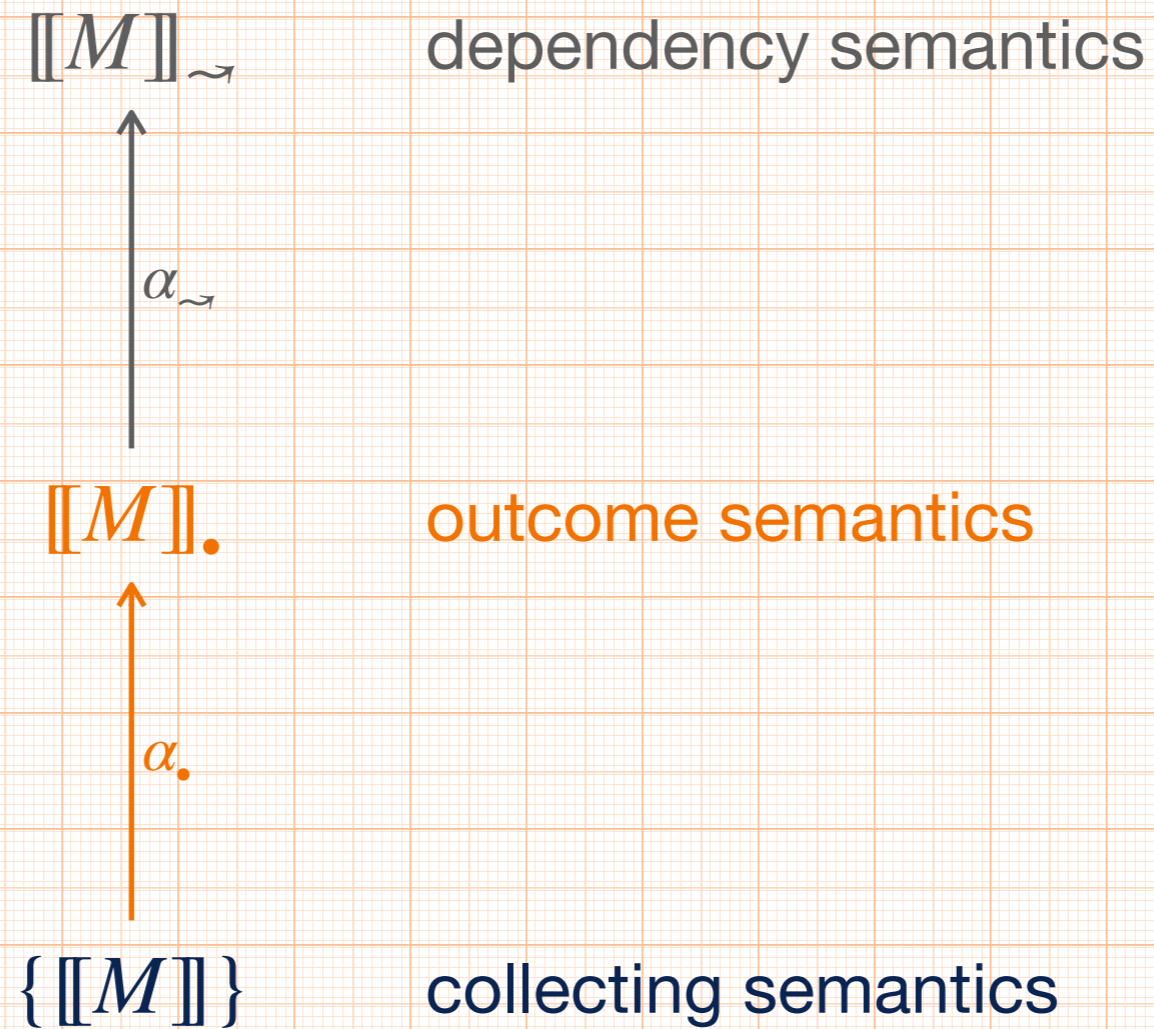
Course 2

Program Semantics and Properties

Antoine Miné

p. 24 / 98

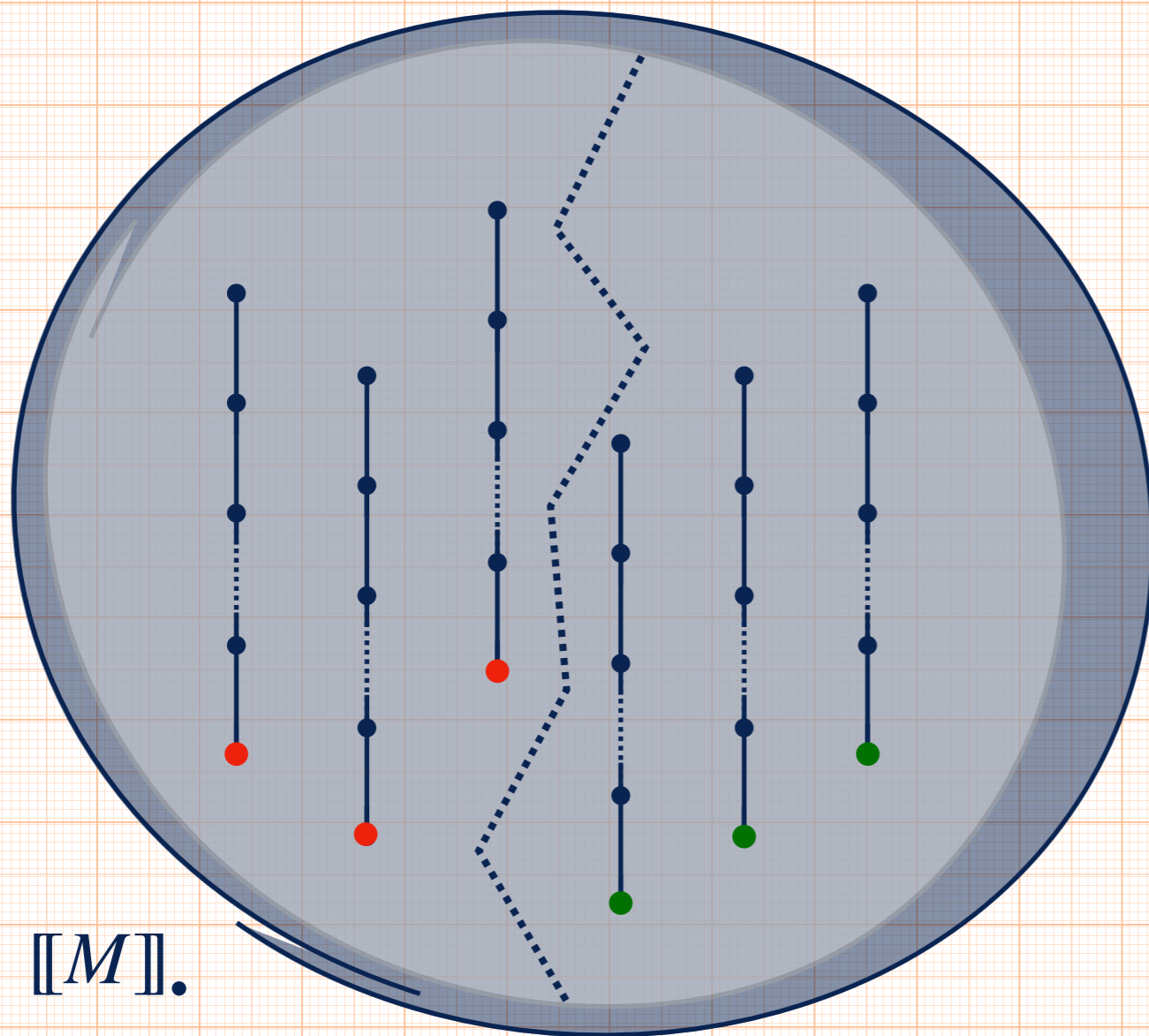
# (Another) Hierarchy of Semantics



# 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



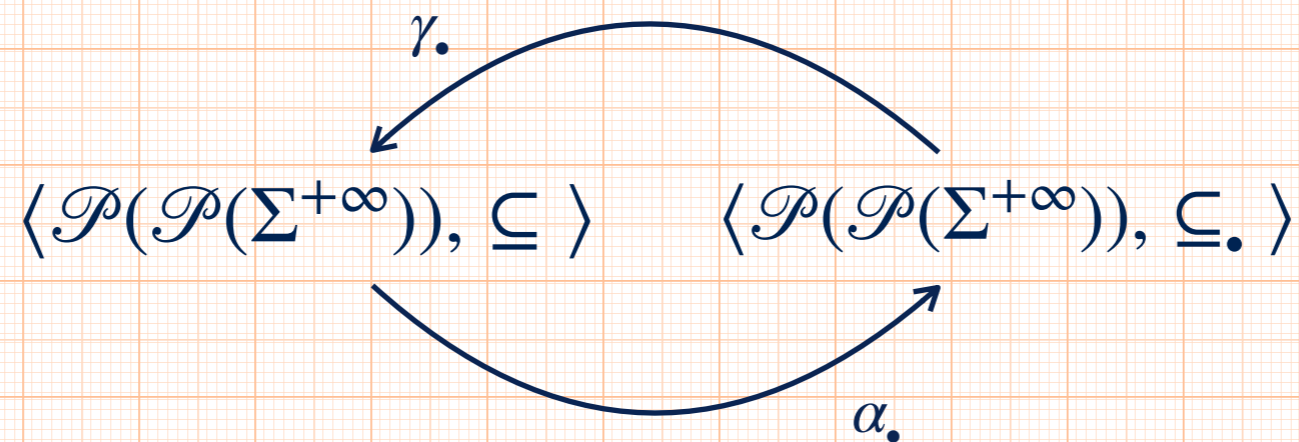
# Outcome Semantics

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

outcomes

Lemma

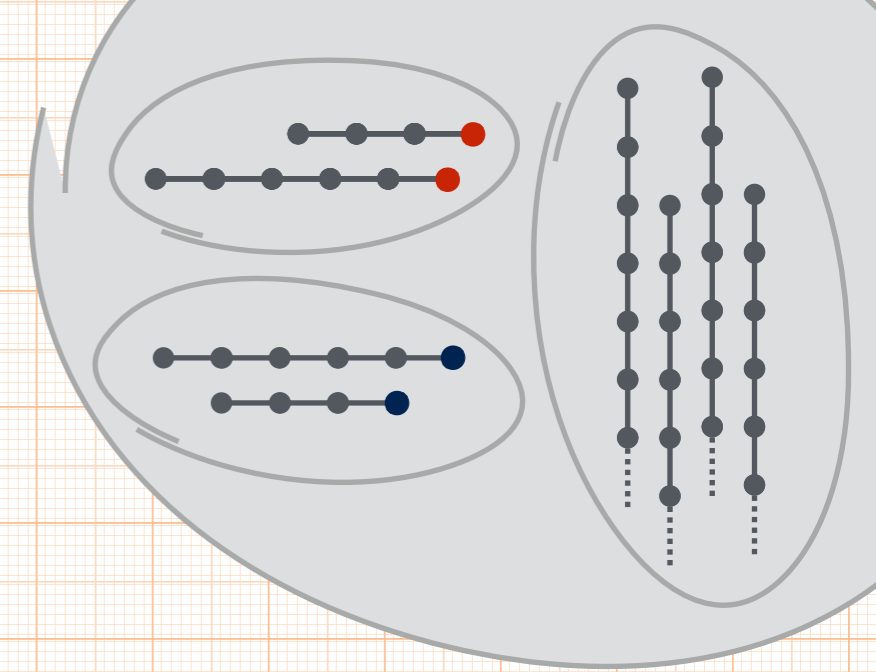
$$P \models \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \cap O \mid O \in \mathbb{O} \} \subseteq \mathcal{N}_J$$



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

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 = <b>True</b>	
english → _	english → _
math → <b>T</b>	math → <b>T</b>
science → _	science → _
bonus → _	bonus → _
passing → ?	passing → <b>T</b>

passing = <b>True</b>	passing = <b>False or bonus</b>	passing = <b>False or bonus</b>	passing = <b>False or bonus</b>
english → _	english → _	english → _	english → _
math → <b>F</b>	math → <b>F</b>	math → <b>F</b>	math → <b>F</b>
science → _	science → _	science → _	science → _
bonus → <b>T</b>	bonus → <b>T</b>	bonus → <b>T</b>	bonus → <b>T</b>
passing → ?	passing → <b>T</b>	passing → <b>T</b>	passing → <b>T</b>

passing = <b>True</b>	passing = <b>False or bonus</b>	passing = <b>False or bonus</b>	passing = <b>False or bonus</b>
english → _	english → _	english → _	english → _
math → <b>F</b>	math → <b>F</b>	math → <b>F</b>	math → <b>F</b>
science → _	science → _	science → _	science → _
bonus → <b>F</b>	bonus → <b>F</b>	bonus → <b>F</b>	bonus → <b>F</b>
passing → ?	passing → <b>T</b>	passing → <b>F</b>	passing → <b>F</b>



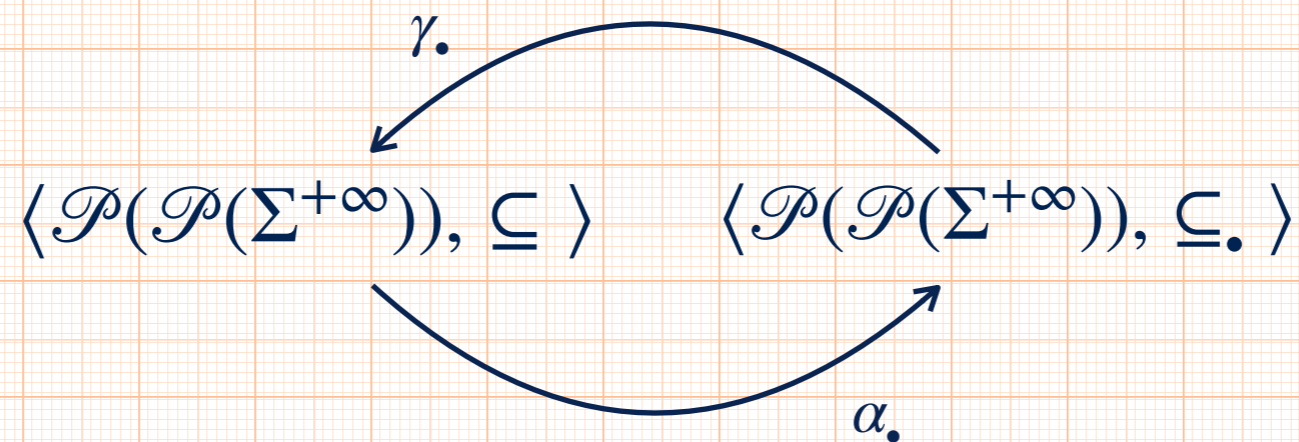
# Outcome Semantics

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

outcomes

Lemma

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

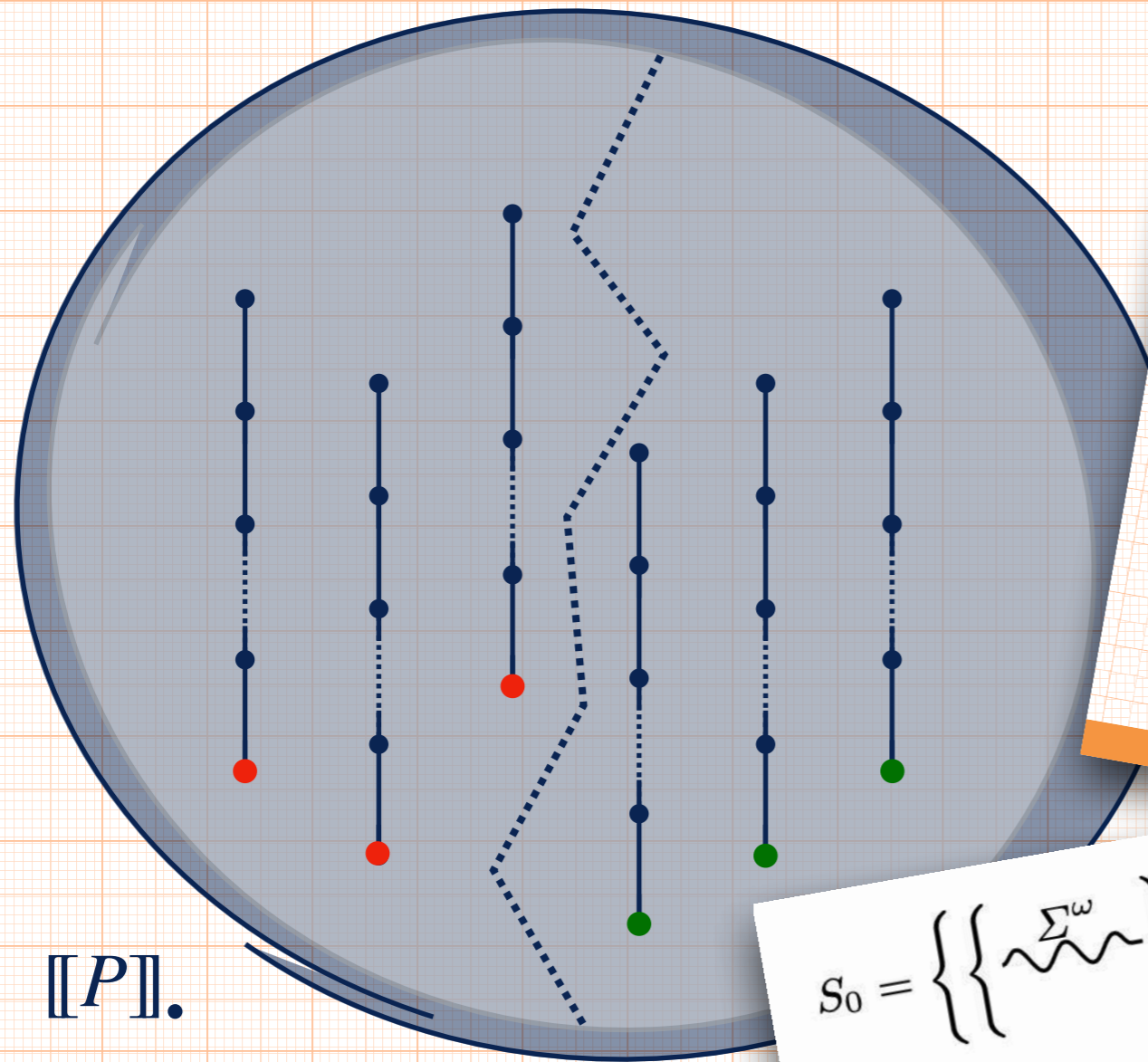


$$\alpha_\bullet(S) \stackrel{\text{def}}{=} \{ T \cap O \mid T \in S \wedge O \in \mathbb{O} \}$$

outcome abstraction

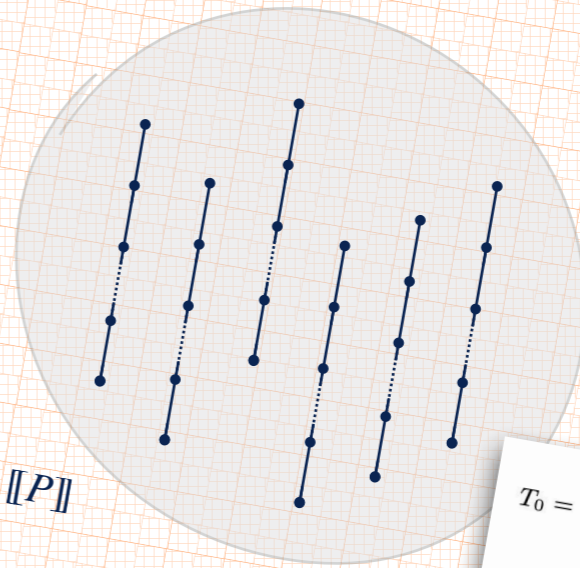
$$\llbracket P \rrbracket_\bullet \stackrel{\text{def}}{=} \alpha_\bullet(\{ \llbracket P \rrbracket \}) = \{ \llbracket P \rrbracket \cap O \mid O \in \mathbb{O} \}$$

# Outcome Semantics



[[P]].

## Maximal Trace Semantics



[[P]]

Least fixpoint formulation of maximal trace semantics

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  $(\mathcal{P}(\Sigma^\omega), \subseteq, \cup, \cap, \emptyset, \Sigma^\omega)$ .  
 (we transform the greatest fixpoint into a least fixpoint!)

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

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

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

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

$$T_1 = \left\{ \begin{array}{c} \Omega \\ \bullet \end{array} \right\} \cup \left\{ \begin{array}{c} \tau \\ \bullet \end{array} \right\} \rightarrow \begin{array}{c} \Sigma^\omega \\ \sim \end{array}$$

$$T_2 = \left\{ \begin{array}{c} \Omega \\ \bullet \end{array} \right\} \cup \left\{ \begin{array}{c} \tau \\ \bullet \end{array} \right\} \rightarrow \begin{array}{c} \Omega \\ \bullet \end{array} \right\} \cup \left\{ \begin{array}{c} \tau \\ \bullet \end{array} \right\} \rightarrow \begin{array}{c} \tau \\ \bullet \end{array} \rightarrow \begin{array}{c} \Sigma^\omega \\ \sim \end{array}$$

$$S_0 = \left\{ \left\{ \begin{array}{c} \Sigma^\omega \\ \sim \end{array} \right\}, \emptyset \right\}$$

$$S_1 = \left\{ \left\{ \begin{array}{c} \Omega_{o=v} \\ \bullet \end{array} \right\} \mid v \in V \right\} \cup \left\{ \begin{array}{c} \tau \\ \bullet \end{array} \right\} \rightarrow \begin{array}{c} \Sigma^\omega \\ \sim \end{array} \right\}$$

$$S_2 = \left\{ \left\{ \begin{array}{c} \Omega_{o=v} \\ \bullet \end{array} \right\} \cup \left\{ \begin{array}{c} \tau \\ \bullet \end{array} \right\} \rightarrow \begin{array}{c} \Omega_{o=v} \\ \bullet \end{array} \right\} \mid v \in V \right\} \cup \left\{ \begin{array}{c} \tau \\ \bullet \end{array} \right\} \rightarrow \begin{array}{c} \tau \\ \bullet \end{array} \rightarrow \begin{array}{c} \Sigma^\omega \\ \sim \end{array} \right\}$$

# 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:

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

$$\Theta_\bullet(S) \stackrel{\text{def}}{=} \{\Omega_{o_1=v_1, \dots, o_k=v_k} \mid v_1, \dots, v_k \in V\} \cup \{\tau; T \mid T \in S\} \quad (9)$$

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

(proof by Kleenian Fixpoint Transfer [Urban18])

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

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

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

# 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  $J \subseteq I_P$

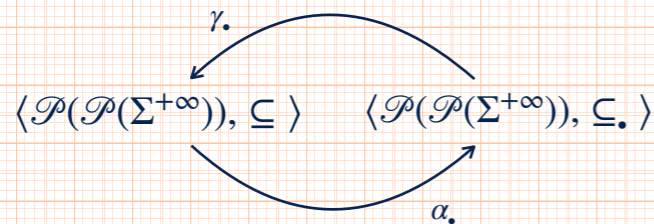
## Outcome Semantics

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

outcomes

Lemma

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



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

outcome abstraction

$$\llbracket P \rrbracket . \stackrel{\text{def}}{=} \alpha.(\{ \llbracket P \rrbracket \}) = \{ \llbracket P \rrbracket \cap \mathbb{O} \mid \mathbb{O} \in \mathbb{O} \}$$

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

## 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  $J \subseteq I_P$

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

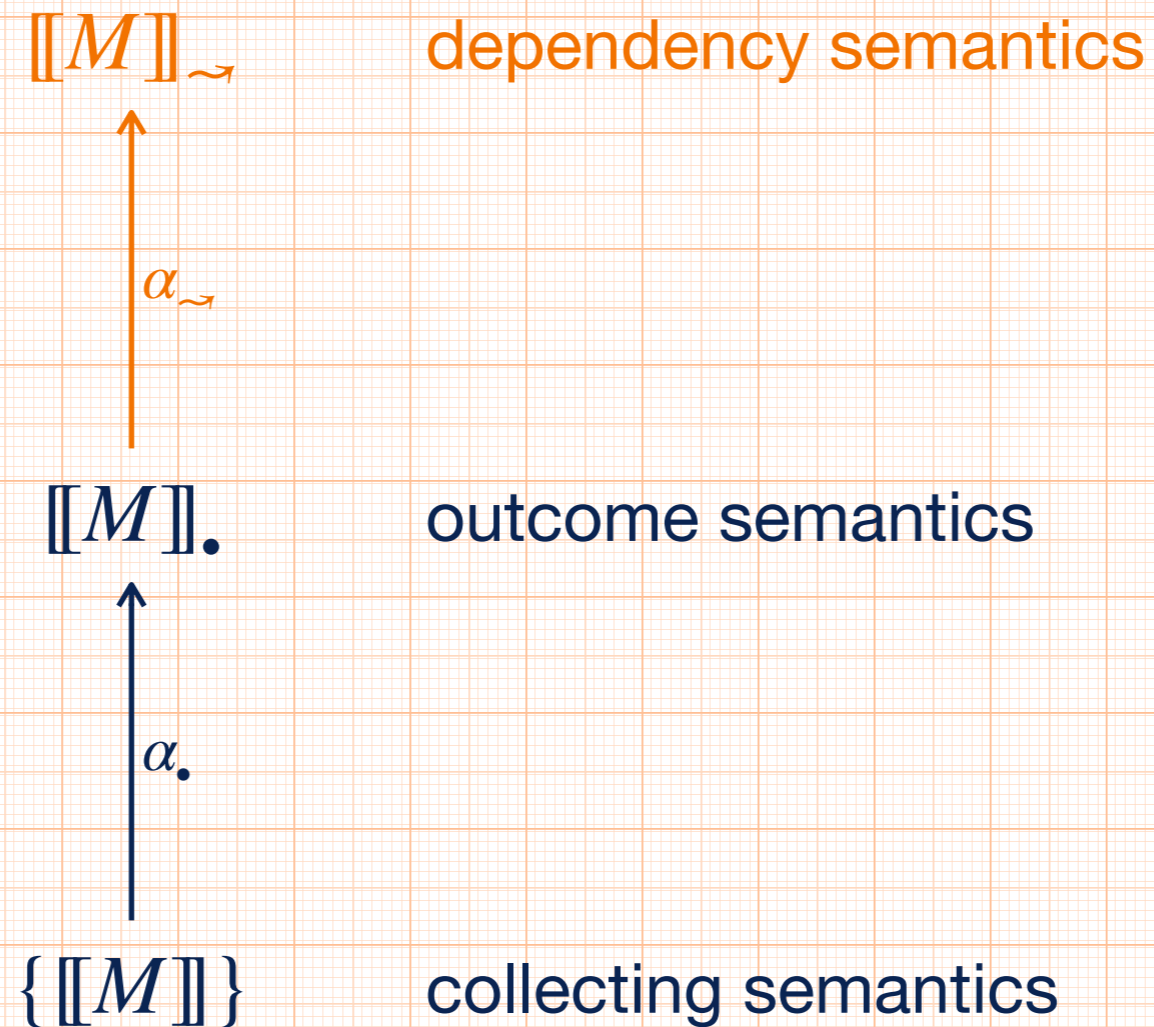
## Theorem

$$P \vDash \mathcal{N}_J \Leftrightarrow \{ \llbracket P \rrbracket \} \subseteq \mathcal{N}_J \Leftrightarrow \alpha.(\{ \llbracket P \rrbracket \}) \subseteq \mathcal{N}_J \Leftrightarrow \llbracket P \rrbracket . \subseteq \mathcal{N}_J$$

Theorem

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

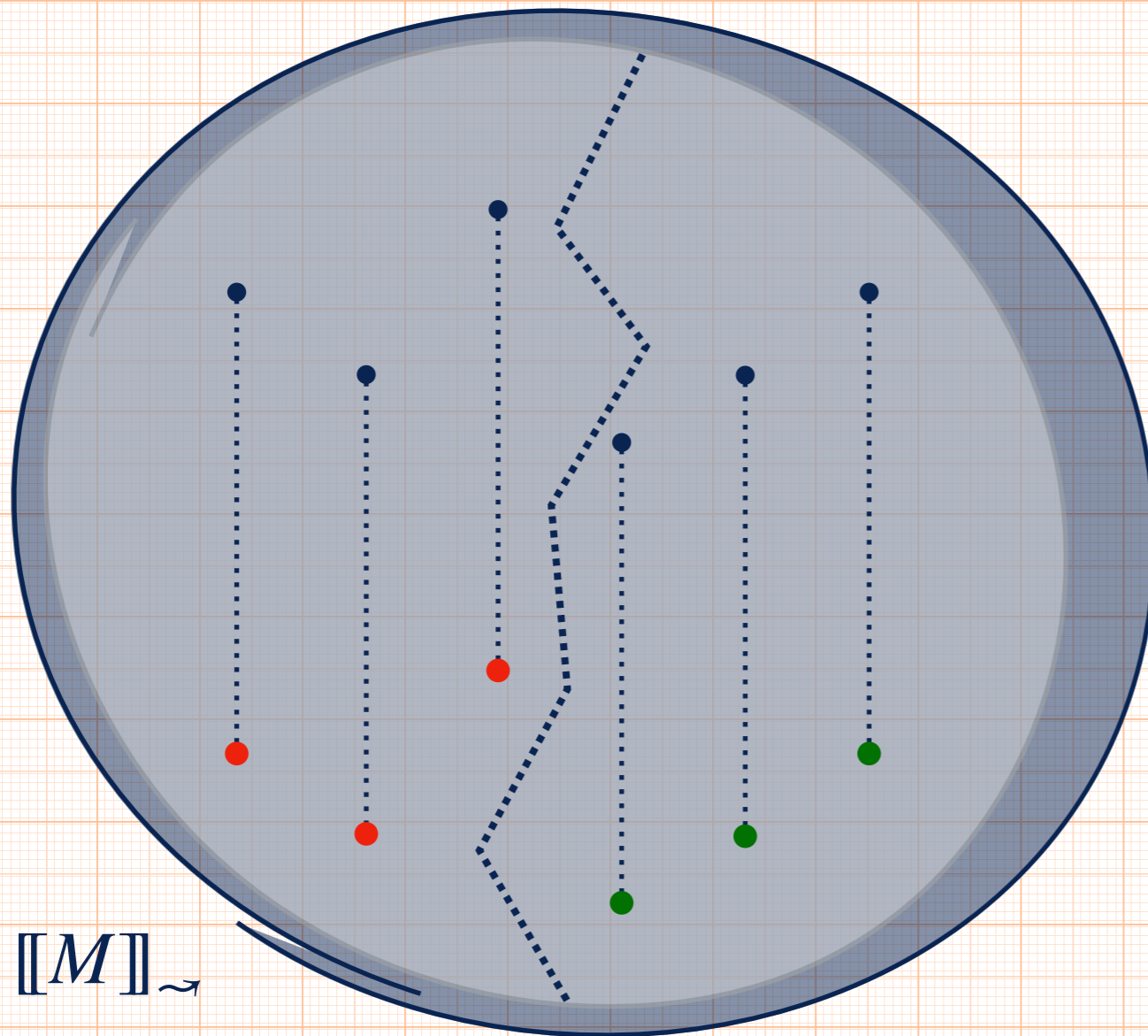
# (Another) Hierarchy of Semantics



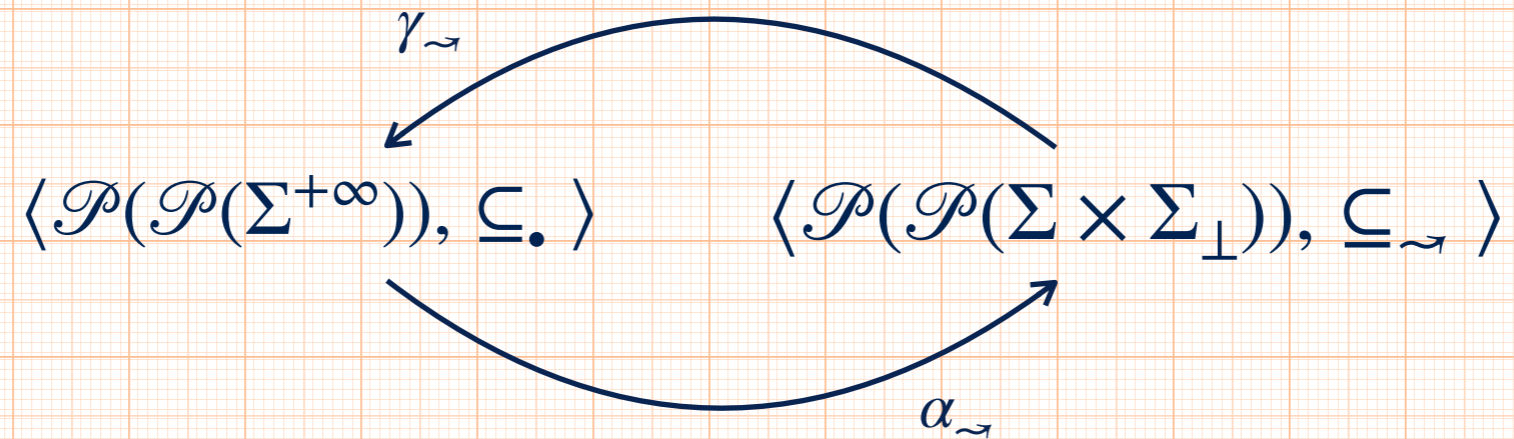
# Dependency Semantics



💡 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)

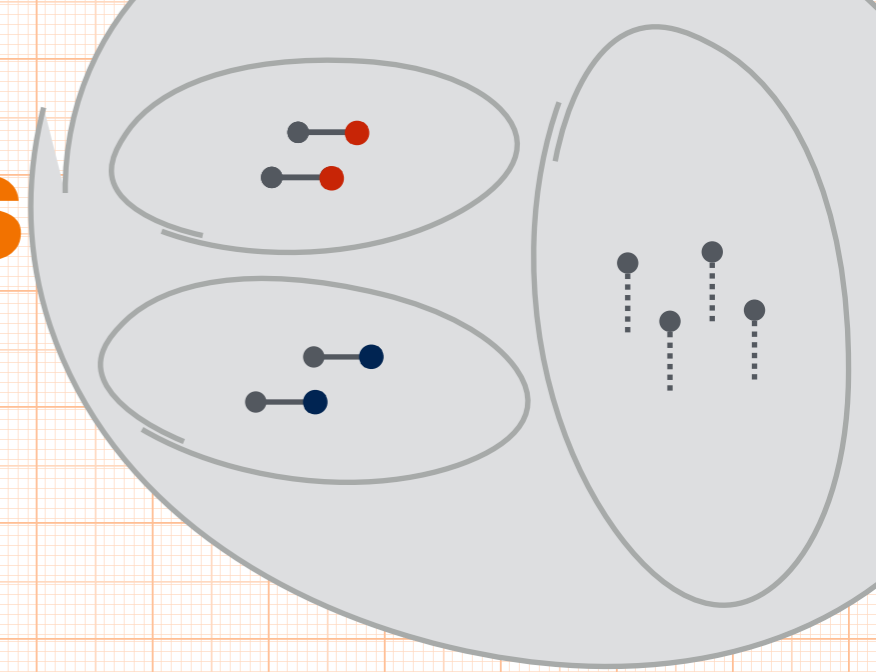


# Dependency Semantics



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

# Dependency Semantics



$[P]$   $\rightsquigarrow$

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

english → _	english → _
math → <b>T</b>	math → <b>T</b>
science → _	science → _
bonus → _	bonus → _
passing → ?	passing → <b>T</b>

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

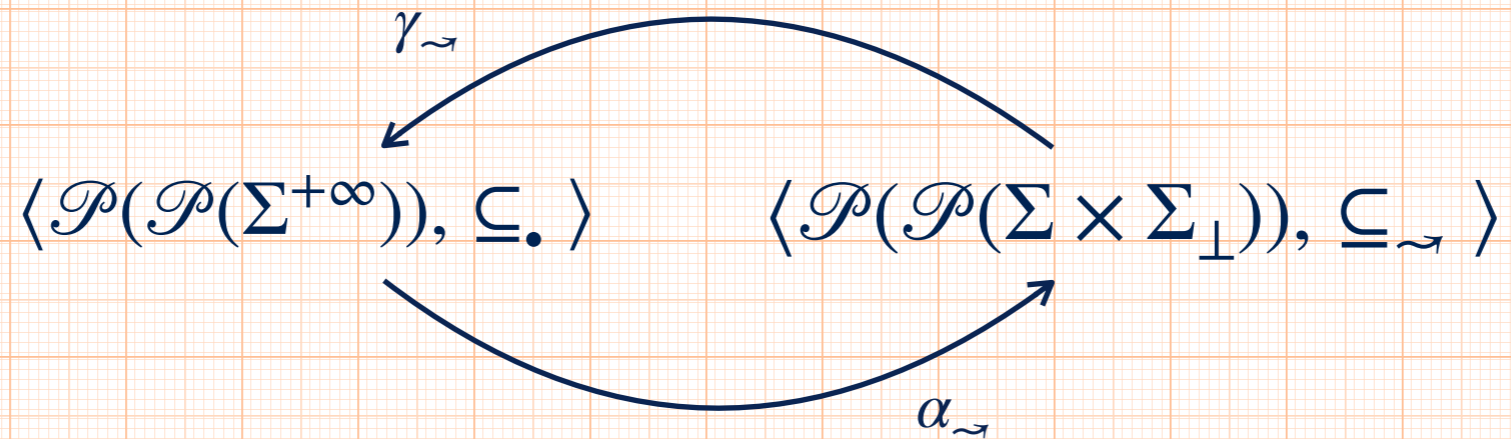
english → \_  
 math → **F**  
 science → \_  
 bonus → **T**  
 passing → **T**

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

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



# Dependency Semantics



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

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

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

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

the outcome semantics  $\Lambda_\bullet$  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

$$\Lambda_{\rightsquigarrow}^{\omega} = \text{lfp}_{\{\Sigma \times \{\perp\}\}}^{\sqsubseteq} \Theta_{\rightsquigarrow}^{\omega} \tag{15}$$

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

*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:

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

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

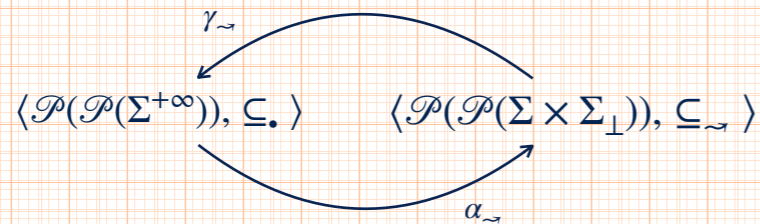
*Proof (Sketch).* The proof follows immediately from Lemma 2 and Lemma 3.  $\square$

# 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



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

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

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

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

## 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

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

## Theorem

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

# Unused Data Analysis

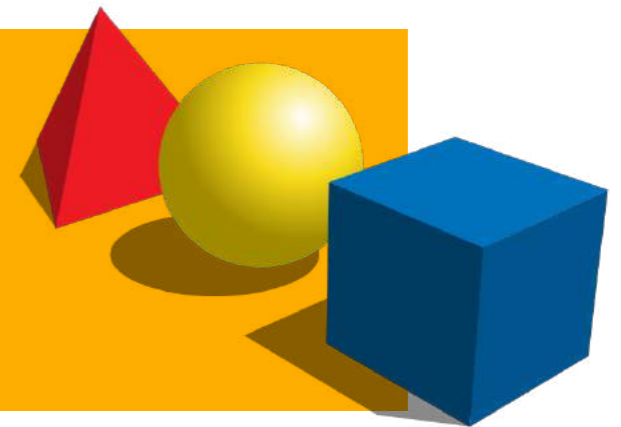
## practical tools

targeting specific programs



## algorithmic approaches

to decide program properties



## mathematical models

of the program behavior



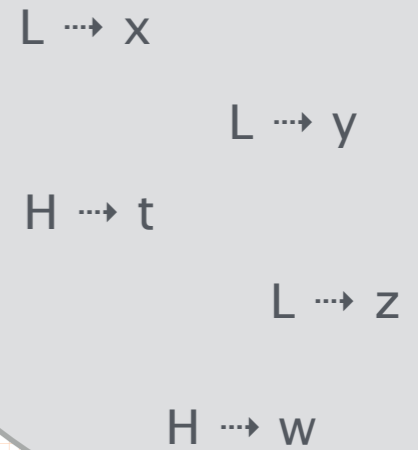
# Input Data (Non-)Usage Abstractions

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



**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)

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

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

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

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

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

## Hypercollecting Semantics and Its Application to Static Analysis of Information Flow

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

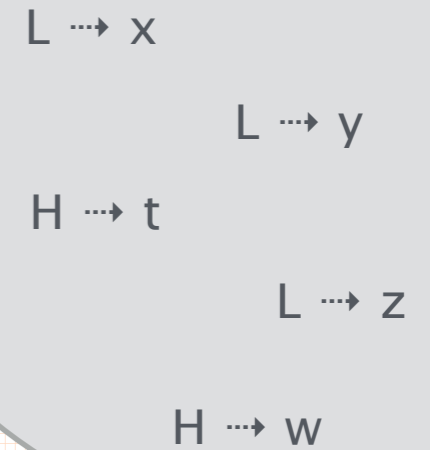
David A. Naumann  
Stevens Institute of Technology,  
Hoboken, US  
first.last@stevens.edu

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

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

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

# 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)

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

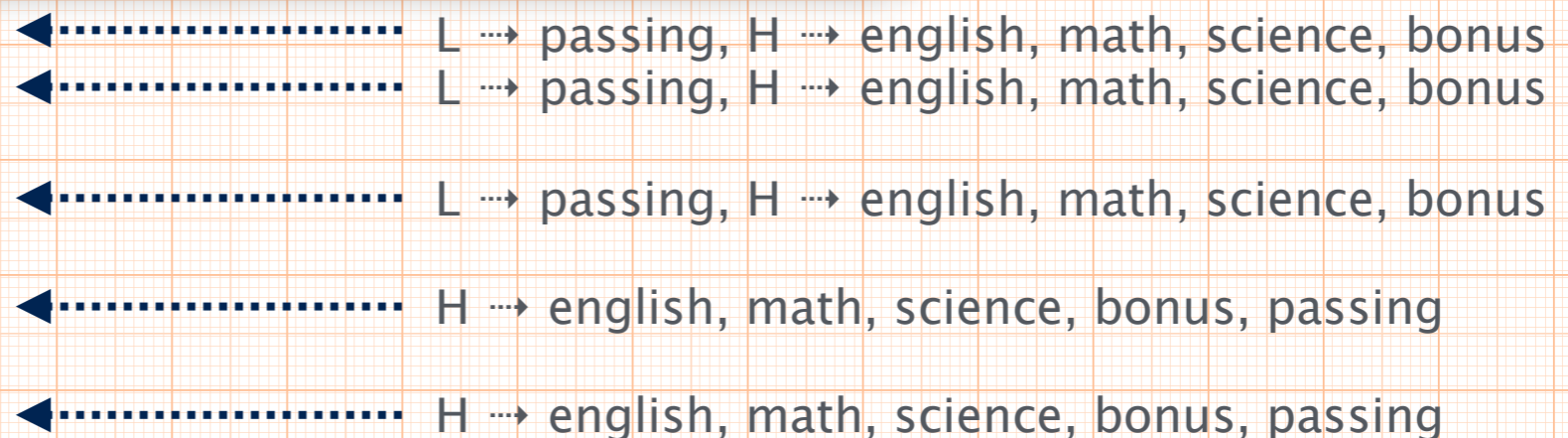
$$\Theta_F[x = e](S) \stackrel{\text{def}}{=} \{L \rightsquigarrow y \in S \mid y \neq x\} \cup \{L \rightsquigarrow x \mid \mathcal{V}_F[e]S\}$$

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

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

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

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





# Secure Information Flow

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

**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**
- 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)

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

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

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

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

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

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

$\leftarrow \dots \dots \dots L \rightsquigarrow \text{passing}, H \rightsquigarrow \text{english, math, science, bonus}$   
 $\leftarrow \dots \dots \dots L \rightsquigarrow \text{passing}, H \rightsquigarrow \text{english, math, science, bonus}$   
 $\leftarrow \dots \dots \dots L \rightsquigarrow \text{passing}, H \rightsquigarrow \text{english, math, science, bonus}$

## 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)

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

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

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

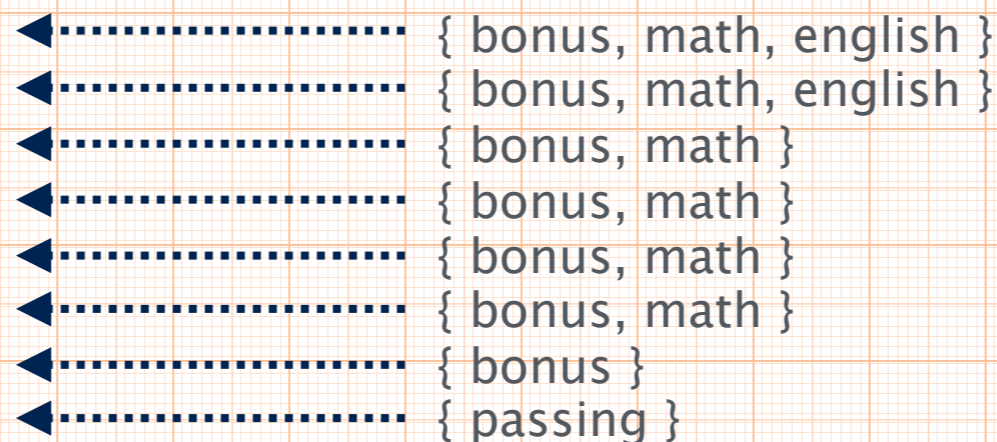
$$\Theta_X[\text{while } b: s](S) \stackrel{\text{def}}{=} \text{VARS}(b) \cup \Theta_X[s](S)$$

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

## Theorem

$$P \vDash \mathcal{N}_J \Leftarrow \gamma_{\sim}(\gamma_X(\llbracket P \rrbracket_X)) \subseteq \mathcal{N}_J$$

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



# Syntactic (Non-)Usage

$x \mapsto U$

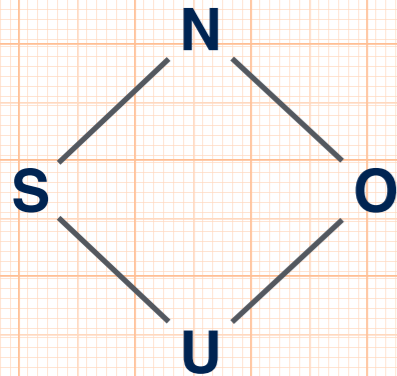
$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

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

$\llbracket P \rrbracket_U$



- **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

$x \mapsto U$

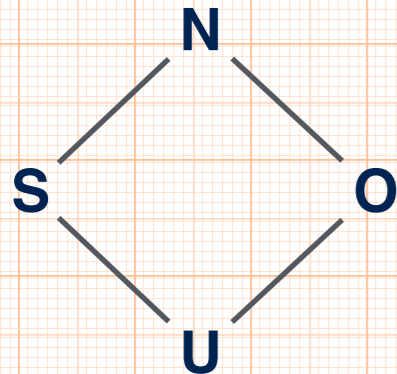
$y \mapsto S \mid y \mapsto U$

$t \mapsto N$

$z \mapsto N$

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

$\llbracket P \rrbracket_U$



- **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

• ..... math, bonus  $\mapsto U$ , passing  $\mapsto O$   
 passing = **True**

• ..... math, bonus, passing  $\mapsto U$   
**if not** english:

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

• ..... math, bonus, passing  $\mapsto S \mid$  math, bonus, passing  $\mapsto U$   
**if not** math:

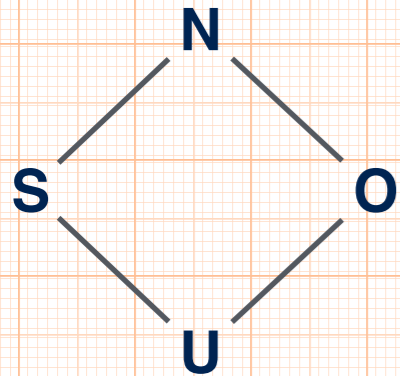
• ..... math, honus, passing  $\mapsto U$   
 passing = **False or** bonus

• .....  
**if not** math:

• .....  
 passing = **False or** bonus

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

# 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

$x \rightsquigarrow U$

$y \rightsquigarrow S \mid y \rightsquigarrow U$

$t \rightsquigarrow N$

$z \rightsquigarrow N$

$w \rightsquigarrow O \mid w \rightsquigarrow U$

$\llbracket P \rrbracket_U$

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

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

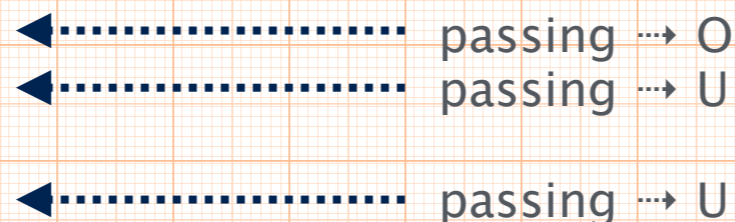
$$\Theta_Q[x = e](q) \stackrel{\text{def}}{=} \text{ASSIGN}[x = e](q)$$

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

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

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

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



## Theorem

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

# Piecewise Syntactic (Non-)Usage

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

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

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

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

$[P]_s$

```
grades = list(map(int, input().split()))
```

```
count = 0
```

● ..... grades  $\rightarrow \{0\} N \{1\}^? U \{2\}^? U \{\text{len}(\text{grades})\}^?$

**i = 1** ← ..... **ERROR: 1 SHOULD BE 0**

● ..... grades  $\rightarrow \{0\} N \{i\}^? U \{i+1\}^? U \{\text{len}(\text{grades})\}^?$

```
while i < len(grades):
```

● ..... grades  $\rightarrow \{0\} N \{i\}^? U \{i+1\}^? S \{i+2\}^? S \{\text{len}(\text{grades})\}^? \mid \dots$

```
  if grades[i] < 4:
```

```
    count = count + 1
```

● ..... grades  $\rightarrow \{0\} N \{i+1\}^? S \{i+2\}^? S \{\text{len}(\text{grades})\}^? \mid \dots \mid \dots$

● ..... grades  $\rightarrow \{0\} N \{i+1\}^? S \{i+2\}^? S \{\text{len}(\text{grades})\}^? \mid \dots$

```
    i = i + 1
```

● ..... grades  $\rightarrow \{0\} N \{i\}^? S \{i+1\}^? S \{\text{len}(\text{grades})\}^? \mid \dots$

● ..... grades  $\rightarrow \{0\} N \{\text{len}(\text{grades})\}^?$

```
if 2 * count < len(grades):
```

```
  passing = True
```

```
else:
```

```
  passing = False
```

● ..... grades  $\rightarrow \{0\} N \{\text{len}(\text{grades})\}^?$

```
print(passing)
```

# Unused Data Analysis

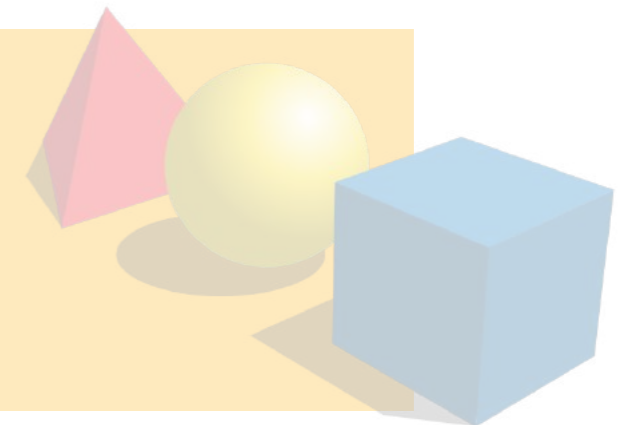
## practical tools

targeting specific programs



## algorithmic approaches

to decide program properties



## mathematical models

of the program behavior



# Jupyter Notebooks

```
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

## Beyond Interactive: Notebook Innovation at Netflix

By [Michelle Ufford](#), [M Pacer](#), [Matthew Seal](#), and [Kyle Kelley](#)

Notebooks have rapidly grown in popularity among data scientists as a de facto standard for quick prototyping and experimentation, pushing the boundaries even further. We'll share our motivations for using it, and what they can do for you. This vision a reality.

In this post, we'll share our motivations for using it, and what they can do for you. We'll also introduce you to some of the novel ways we're using it to explore some of the novel ways we're using it.

If you're short on time, we suggest you read this post.

### Motivations

Data powered by notebooks

## Databricks for Data Science

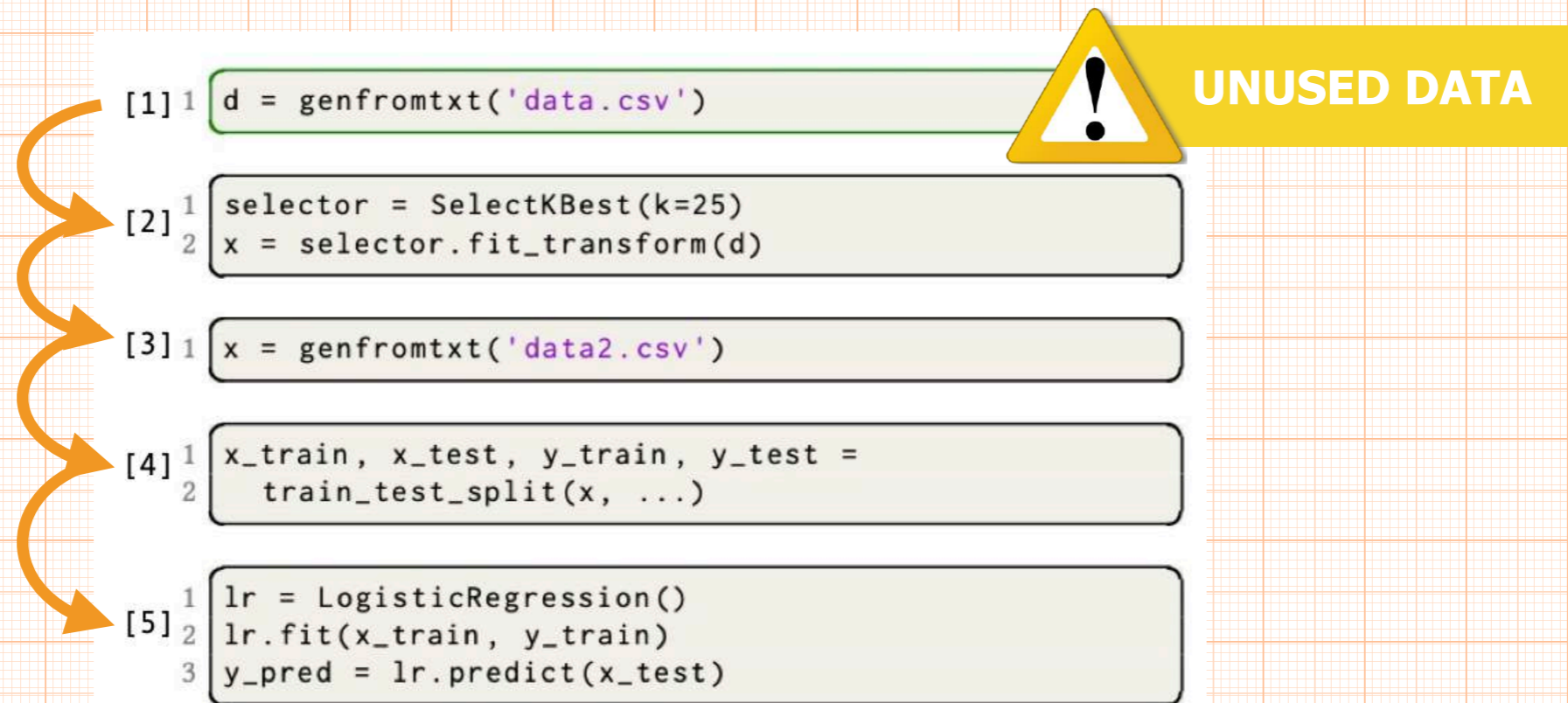
An open and unified platform to collaboratively run all types of analytics workloads, from data preparation to exploratory analysis and predictive analytics, at scale.

```
SELECT *
FROM silver_loan_stats
WHERE int_rate > 6.24
ORDER BY int_rate ASC
```

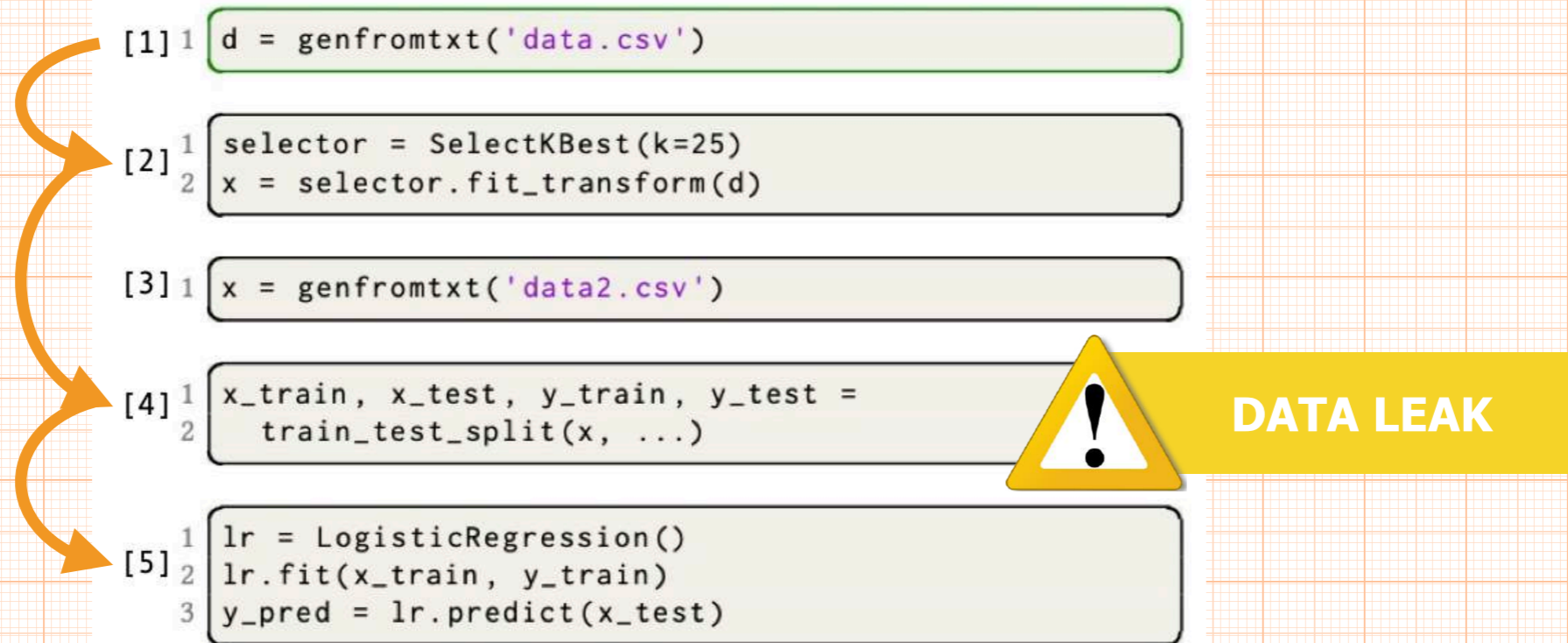
loan_status	int_rate	revolt_date	issue_date	verified_status	emp_length	verification_status	total_payment	loan_amount	grade	annual_inc	dti	addr_state	term
Fully Paid	6.39	55	Apr-2015	Dec-1989	8	Not Verified	8902.76	8400	A	120000	14.33	NY	36 mo
Fully Paid	6.39	15.4	Mar-2015	Sep-1989	7	Not Verified	12058.08	12000	A	75000	4.66	NY	36 mo



# Jupyter Notebooks



# Jupyter Notebooks



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

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

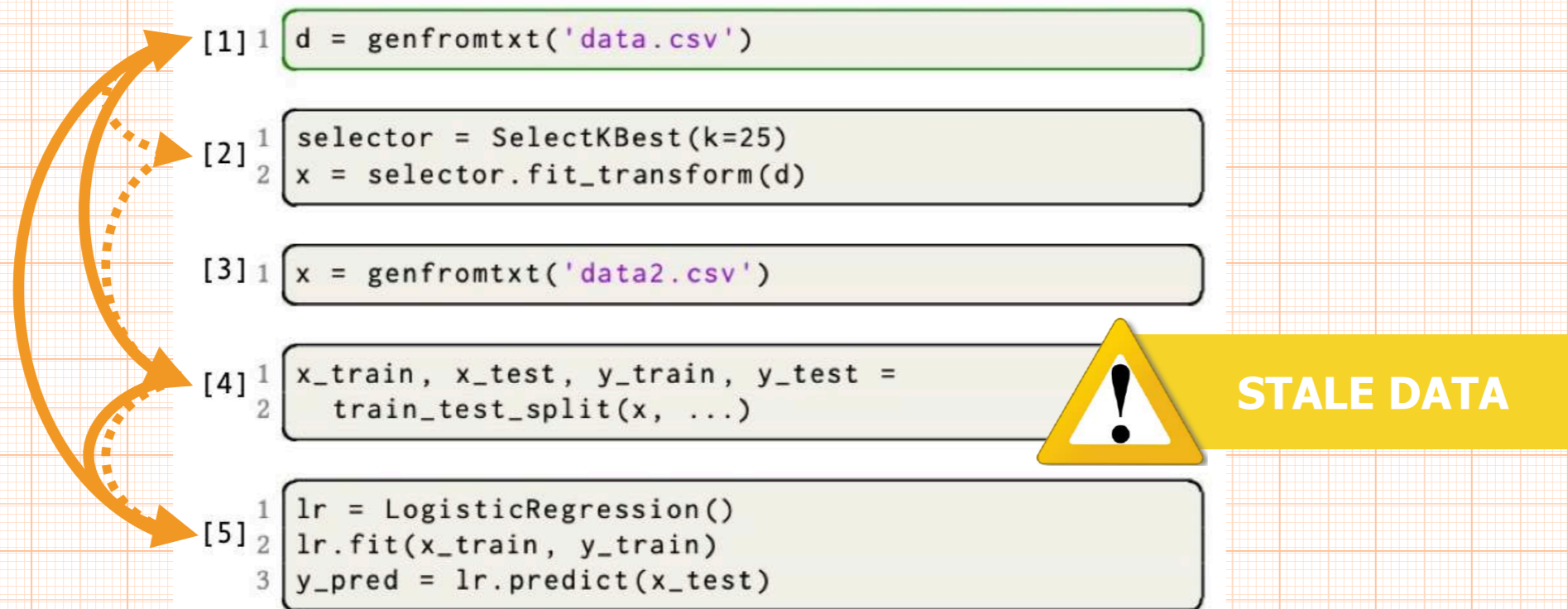
[3] 1 x = genfromtxt('data2.csv')

[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)
```

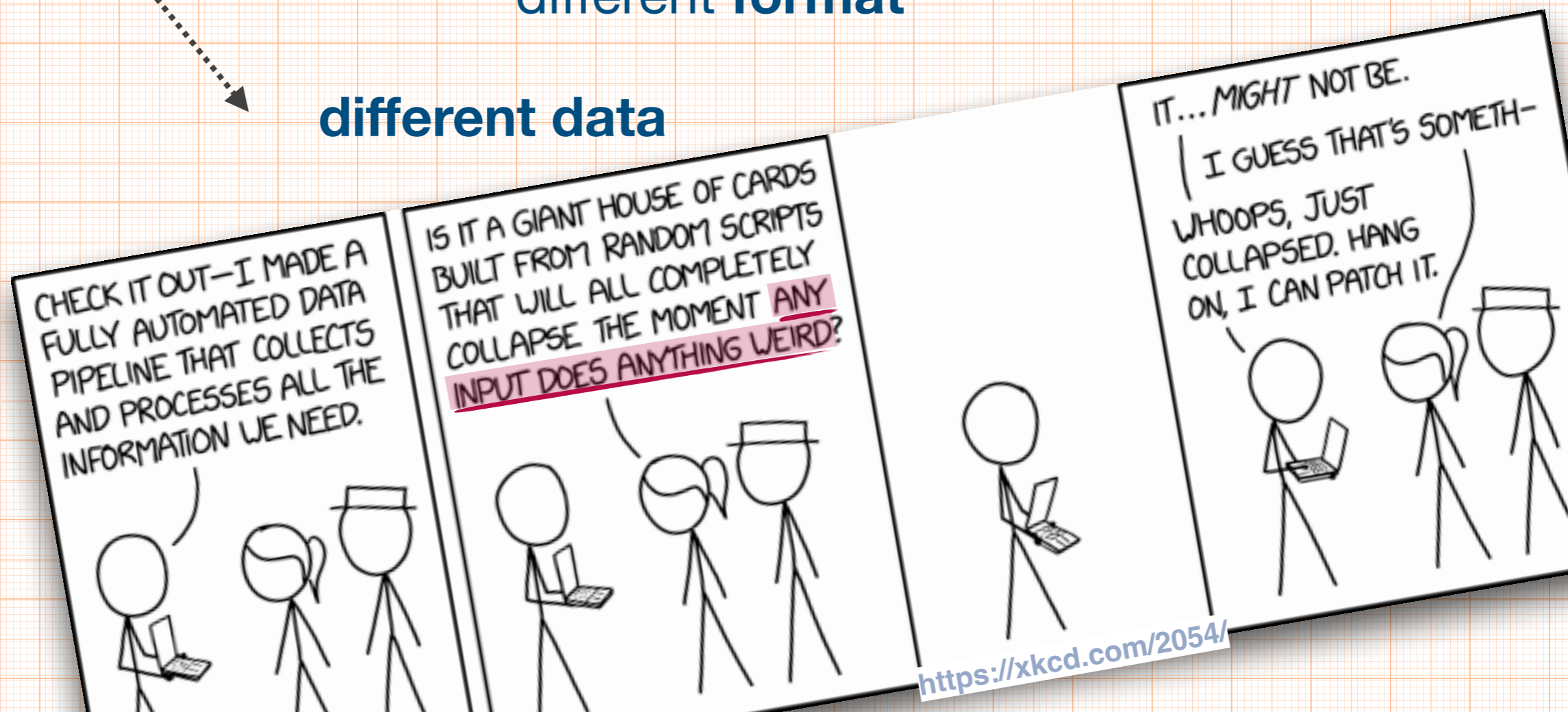
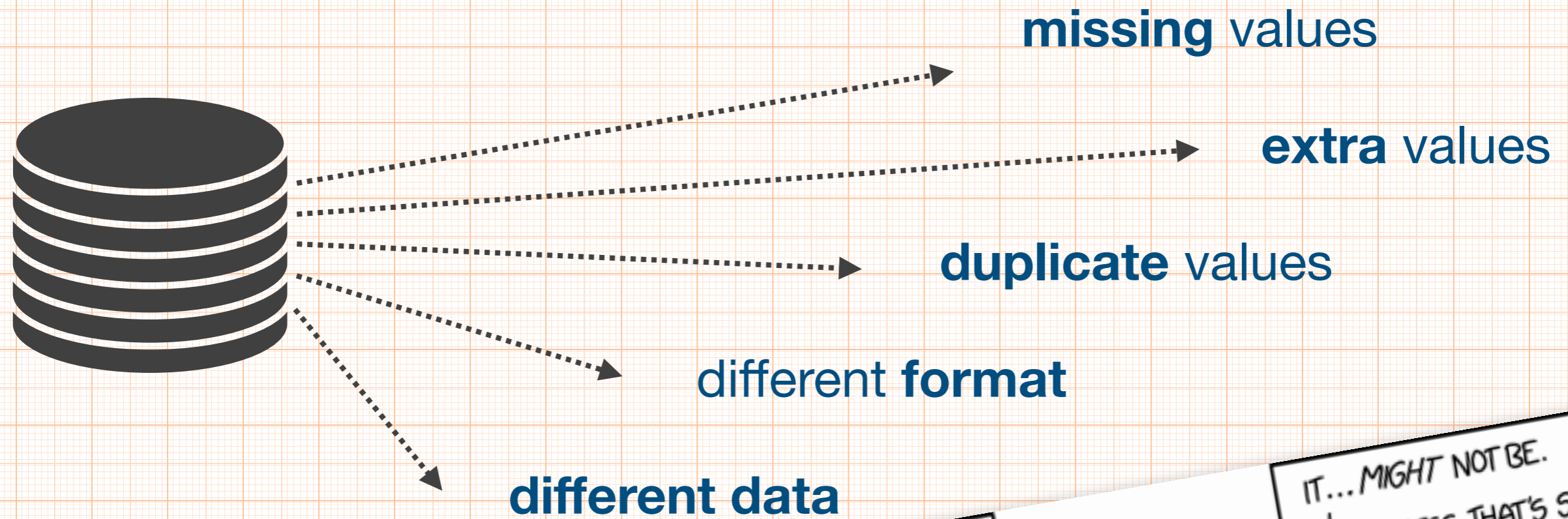
**DATA LEAK**

# 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

**Left Notebook (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

**Right Notebook (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]:
```

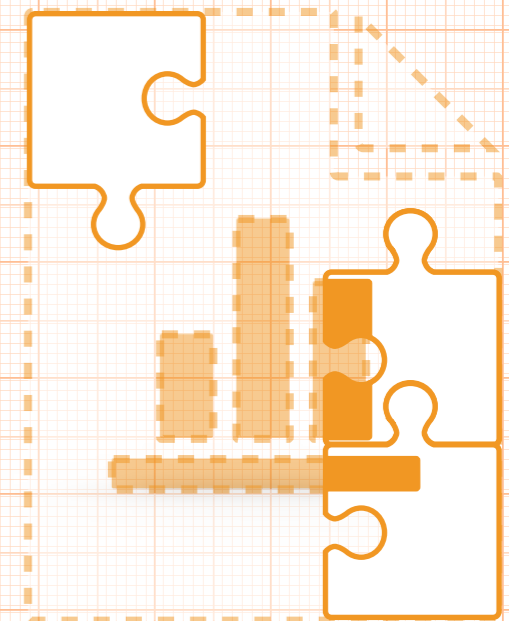
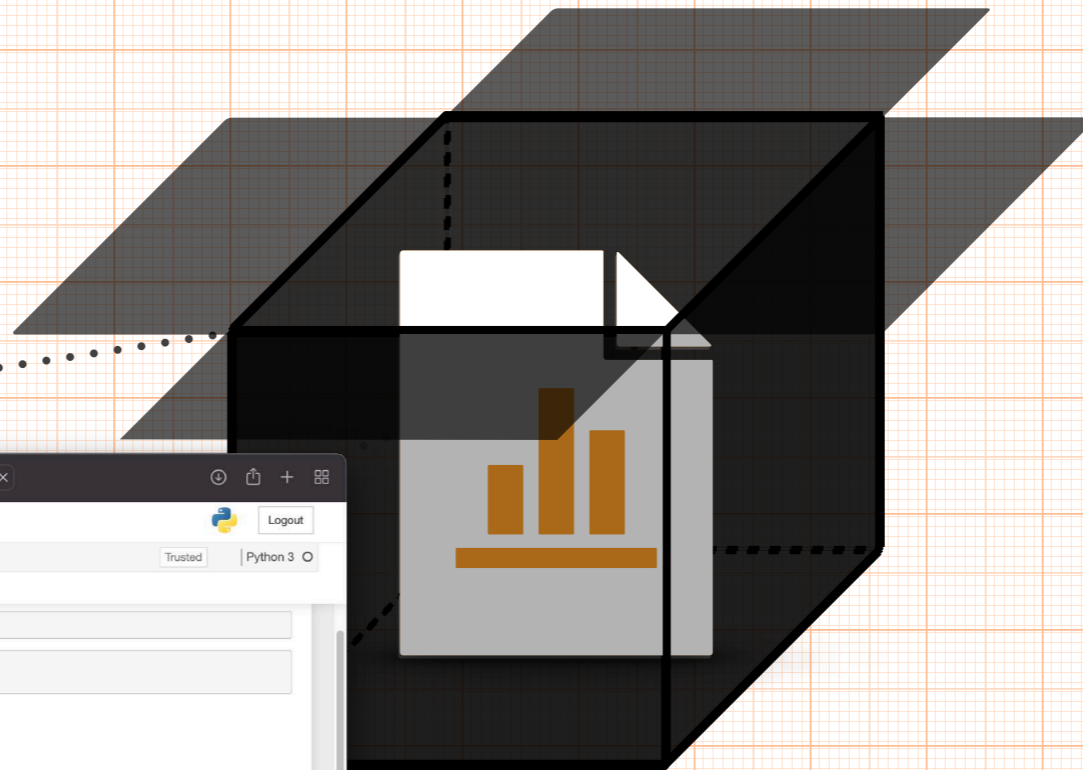
	Name	Q1	Q2	Q3
ID				
2394	Alice	A	A	A
4583	Bob	F	B+	B
3956	Carol	F	A	C
9578	David	D	F	C

```
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

# Data Expectations Analysis

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



# Data Expectations Analysis

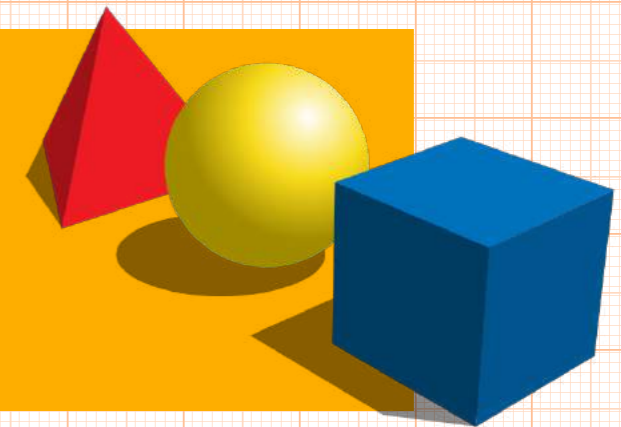
## practical tools

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.