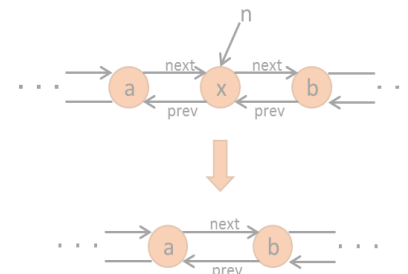
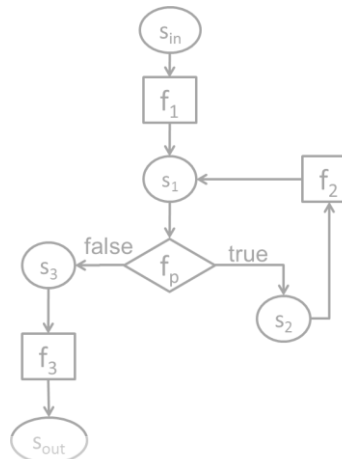


$$\exists c \forall in \ Q(c, in)$$

```

/* Average of x and y without using x+y (avoid overflow)*/
int avg(int x, int y){
  int t = expr({x/2, y/2, x%2, y%2, 2 }, {PLUS, DIV});
  assert t == (x+y)/2;
  return t;
}

```

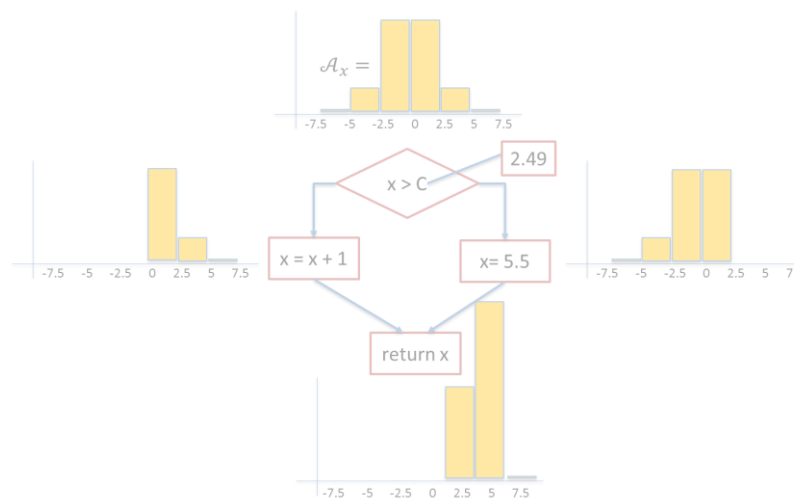
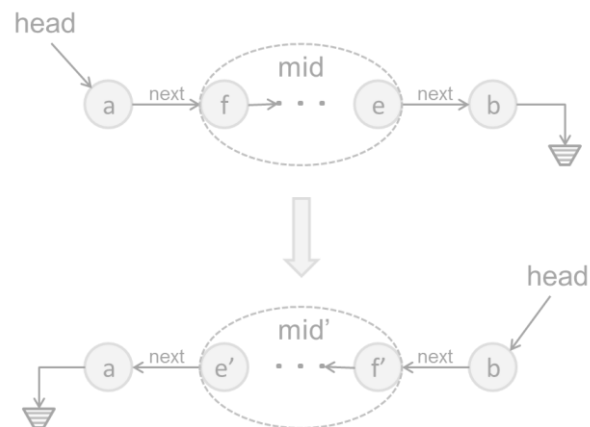


```

{
  s = n.succ;
  p = n.pred;
  p.succ = s;
  s.pred = p;
}

```

# Module III: Applications of Synthesis



$$\varphi(p)$$

$$Sk[c](in)$$

# Logistics

---

## Project presentations

- Tuesday Mar 17, 3-6pm; CSE 2154
- 20 min per team (15 min presentation + questions)
- Structure: motivation, **demo**, technique, evaluation
- Talk to me if you can't make it

## Project reports

- Due on Mar 20 (start working on them now!)
- Format: see course organization page (3-5 pages, SIGPLAN format)

# Lecture 15

## Overview of Applications

*Nadia Polikarpova*

# Applications of synthesis

---

We have seen:

- End-user spreadsheet programming [FlashFill, BlinkFill]
- Superoptimization [Stoke]

Today:

- Custom data structures
- Data extraction and data wrangling
- Databases

Thursday: synthesis as AI

Next week: synthesis for programmers

# Applications of synthesis

---

- Custom data structures
  - Loncaric, Torlak, Ernst: Fast synthesis of fast collections. PLDI'16

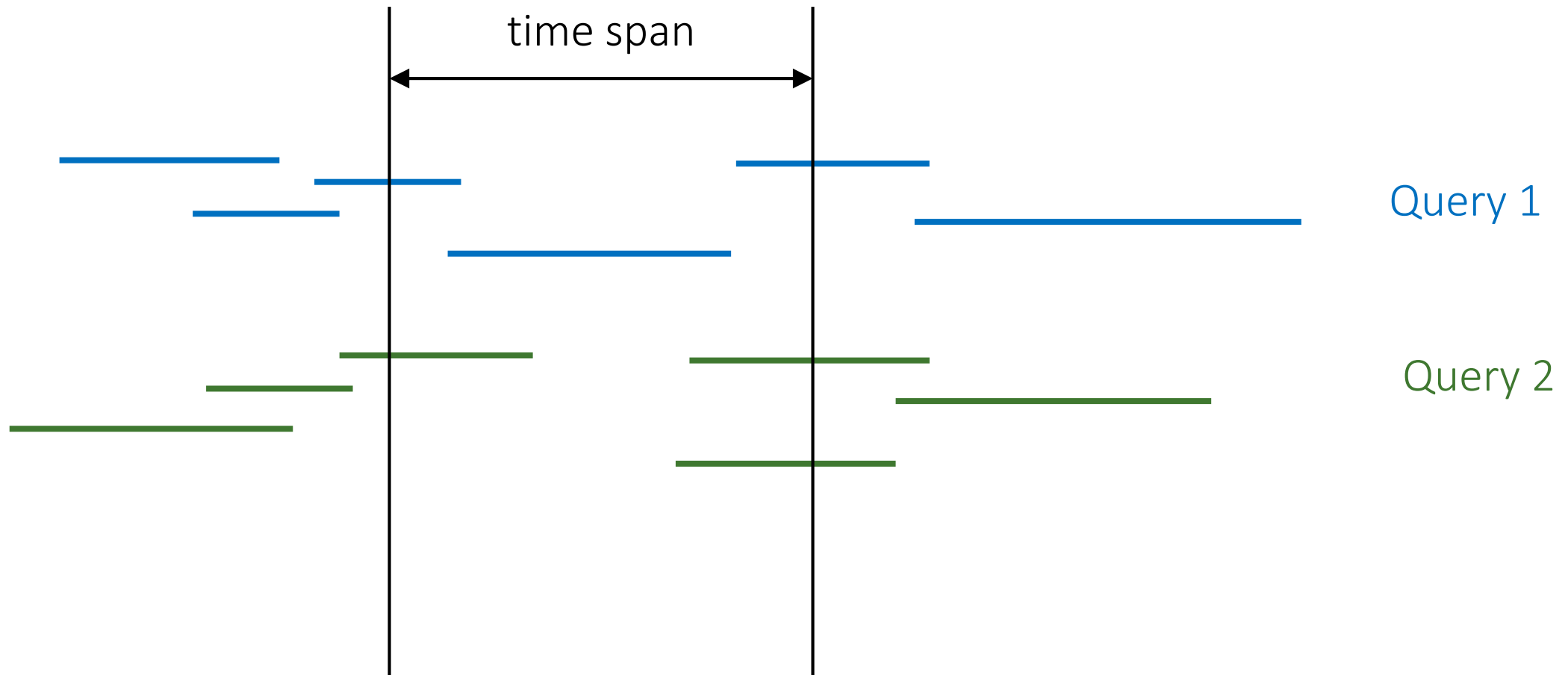
Data extraction and data wrangling

Databases

# Fast Synthesis of Fast Collections

[Loncaric et al. PLDI'16]

Myria distributed database: needs to retrieve all queries in a given timespan



# Specification

---

## fields

queryId:long, subqueryId:long,  
fragmentId:int, opId:int,  
startTime:long, endTime:long,

**assume** startTime <= endTime

**query** getAnalyticsInTimespan(  
    v\_queryId:long, v\_subqueryId:long,  
    v\_fragmentId:int,  
    v\_start:long, v\_end:long)

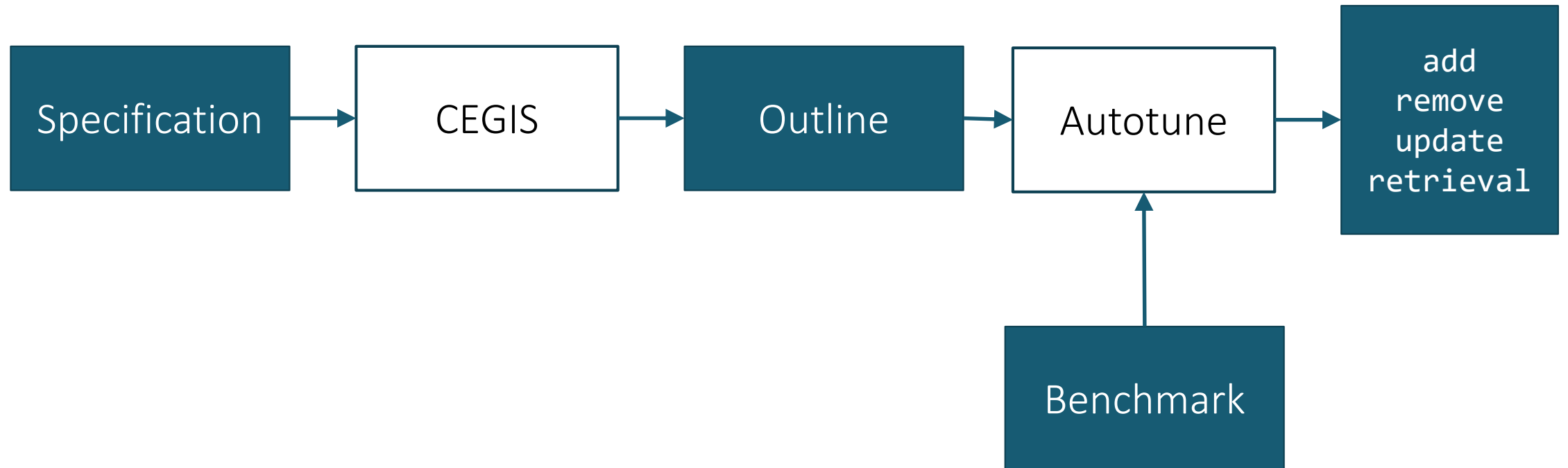
**assume** v\_start <= v\_end

queryId == v\_queryId and  
subqueryId == v\_subqueryId and  
fragmentId == v\_fragmentId and  
startTime < v\_end and  
endTime >= v\_start

costmodel myria-cost.java

# The Cozy workflow

---





# Results

---

Match performance of hand-written code on four real-world applications:

- Myria (a distributed data-base)
- Bullet (a physics simulation library)
- ZTopo (a topographic map viewer)
- Sat4J (a SAT solver)

# Applications of synthesis

---

Custom data structures

→ Data extraction and data wrangling

- Vu Le, Sumit Gulwani: FlashExtract: a framework for data extraction by examples. PLDI'14
- Inala, Singh: WebRelate: integrating web data with spreadsheets using examples. POPL'17
- Feng, Martins, Van Geffen, Dillig, Chaudhuri: Component-based synthesis of table consolidation and transformation tasks from examples. PLDI'17

Databases

# FlashExtract

[Le, Gulwani. PLDI'14]

**Problem:** extract data from semi-structured sources (e.g. log file) into a list of records

**User input:**

- output schema
- highlights examples of fields

**Search strategy:** VSA

```
DLZ - Summary Report
"Sample ID:,""5007-01""
"Sample Date/Time:,""Wednesday, May 30, 2006 00:43:51""
Intensities
"I/S,""Analyte"",""Mass"",""Conc. Mean"",""Unit"",""Conc. SD"",""RSD"",""Mean""
"|,""Be""",9,0.070073,""ug/L""",0.009,12.542,121.334"
"|>,""Sc""",45,""ug/L""",,,404615.043"
"|,""Ti""",48,10.653153,""ug/L""",0.847,7.949,181379.200"
"|,""Se""",82,1.009204,""ug/L""",0.026,2.613,457.487"
"|,""Sr""",88,20.163079,""ug/L""",2.005,9.943,718014.023"
"|>,""Rh""",103,""ug/L""",,,438976.176"

DLZ - Summary Report
"Sample ID:,""5007-02""
"Sample Date/Time:,""Wednesday, May 30, 2006 01:02:38""
Intensities
"I/S,""Analyte"",""Mass"",""Conc. Mean"",""Unit"",""Conc. SD"",""RSD"",""Mean""
"|,""Mn""",55,71.705740,""ug/L""",0.350,0.489,2428667.736"
"|,""Co""",59,0.131132,""ug/L""",0.004,3.315,3606.816"
"|,""Ba""",138,129.339264,""ug/L""",3.088,2.387,4648771.382"
"|,""Hf""",178,""ug/L""",,,338359.496"
"|,""Ti""",205,2.876992,""ug/L""",0.730,25.380,129217.588"
"|,""Pb""",208,3.671043,""ug/L""",0.026,0.702,228830.402"
```

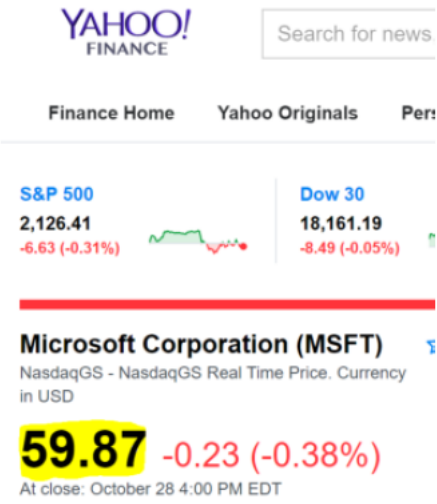
# WebRelate

[Inala, Singh. POPL'17]

**Problem:** extract data from web pages into spreadsheets

**User input:** navigate to a webpage and select content

	Company	URL	Stock price
1	MSFT	https://finance.yahoo.com/q?s=msft	59.87
2	AMZN	https://finance.yahoo.com/q?s=amzn	775.88
3	AAPL	https://finance.yahoo.com/q?s=aapl	113.69
4	TWTR	https://finance.yahoo.com/q?s=twtr	17.66
5	T	https://finance.yahoo.com/q?s=t	36.51
6	S	https://finance.yahoo.com/q?s=s	6.31

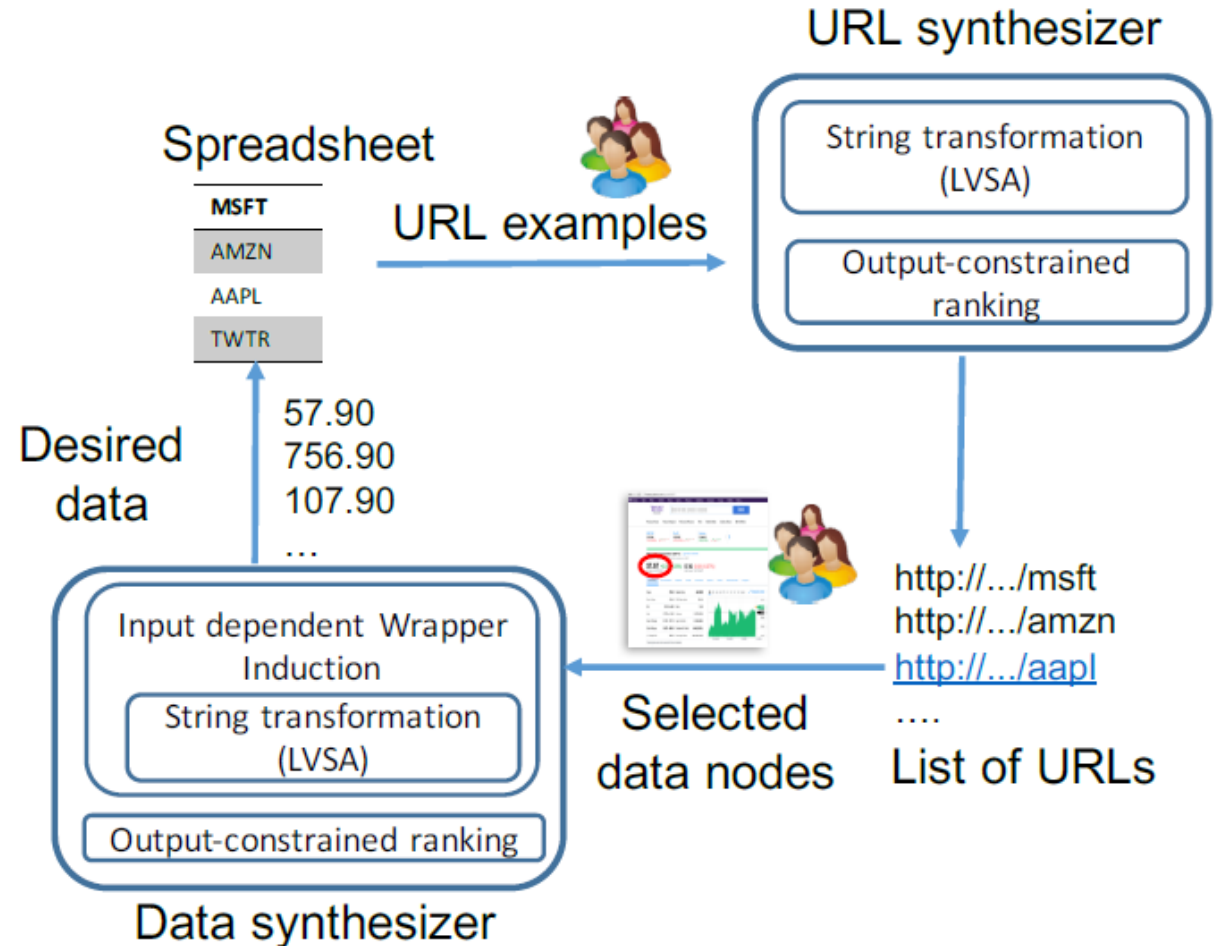


# WebRelate

Search strategy: VSA

Optimizations:

- Layered VSA (URLs are too long for FlashFill-style VSAs)
- Output-constrained synthesis: we know the space of possible outputs



# Morpheus


[Feng et al. PLDI'17]

**Problem:** table data wrangling

**User input:** input-output examples (small tables)

**Search strategy:** enumerative search with deduction

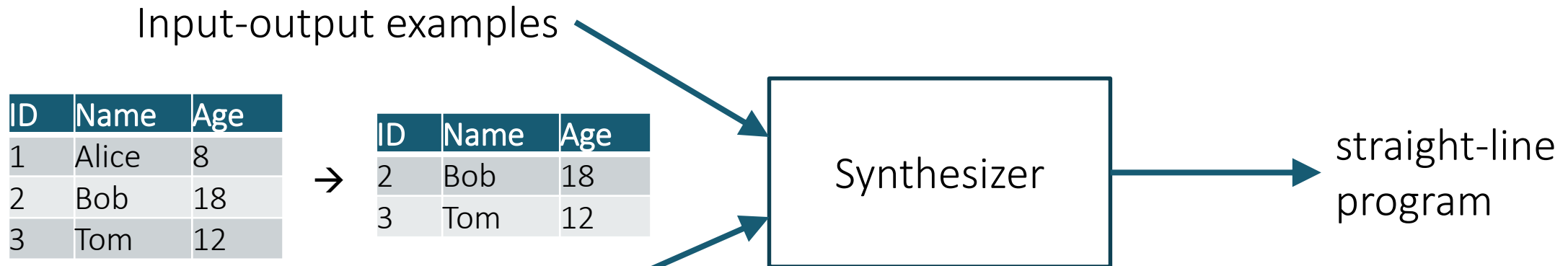
id	year	A	B
1	2007	5	10
2	2009	3	50
1	2007	5	17
2	2009	6	17



<i>id</i>	<i>A_2007</i>	<i>B_2007</i>	<i>A_2009</i>	<i>B_2009</i>
<i>1</i>	<i>5</i>	<i>10</i>	<i>5</i>	<i>17</i>
<i>2</i>	<i>3</i>	<i>50</i>	<i>6</i>	<i>17</i>

# Morpheus: TDP with deduction

[Feng et al'17]



Components

`select : Table → [Col] → Table`

`filter : Table → (Row → Bool) → Table`

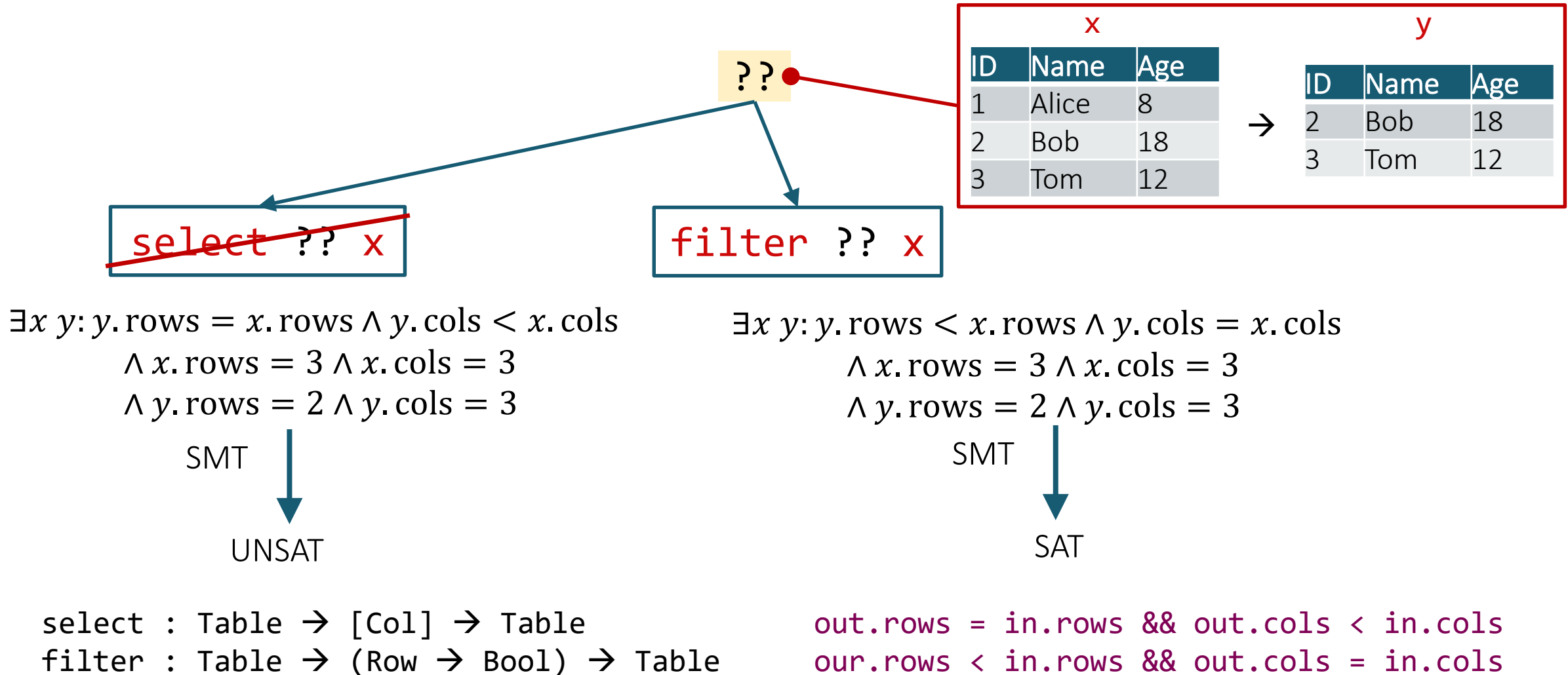
with partial specifications!

`out.rows = in.rows  
&& out.cols < in.cols`

`our.rows < in.rows  
&& out.cols = in.cols`

# Morpheus: TDP with deduction

[Feng et al'17]





# Applications of synthesis

---

Custom data structures

Data extraction and data wrangling

→ Databases

- Wang, Cheung, Bodík: Synthesizing highly expressive SQL queries from input-output examples. PLDI'17
- Yaghmazadeh, Wang, Dillig, Dillig: SQLizer: Query Synthesis from Natural Language. OOPSLA'17
- Singh, Meduri, Elmagarmid, Madden, Papotti, Quiané-Ruiz, Solar-Lezama, Tang: Synthesizing Entity Matching Rules by Examples. VLDB'17

Problem: SQL query synthesis

User input: input-output examples (small tables)

$T_1$			$T_2$		$T_{out}$				
id	date	uid	oid	val	$c_0$	$c_1$	$c_2$	$c_3$	$c_4$
1	12/25	1	1	30	1	12/25	1	1	30
2	11/21	1	1	10	4	12/24	2	2	10
4	12/24	2	2	50					
			2	10					

Output:

```
Select *
From (Select *
      From T1
      Where T1.date = 12/24
            Or T1.date = 12/25) T3
Join (Select oid, Max(val)
      From (Select *
            From T2
            Where T2.val < 50) T4
      Group By oid) T5
On T3.uid = T5.oid
```

# Scythe: technique

---

[Wang et al. PLDI'17]

Sketch generation via bottom-up enumeration with OE

- sketches are SQL queries with holes for predicates
- similar idea to Cosi

Predicate synthesis via bottom-up enumeration with optimizations

# SQLizer

---

[Yaghmazadeh et al. 2017]

**Problem:** SQL query synthesis

**User input:** natural language + DB schema

“Find the number of papers in OOPSLA 2010”

**Output:**

```
SELECT count(Publication.pid)
FROM Publication JOIN Conference ON Publication.cid = Conference.cid
WHERE Conference.name = "OOPSLA" AND Publication.year = 2010
```

# SQLizer: techniques

---

[Yaghmazadeh et al. OOPSLA'17]

Sketch generation via semantic parsing

- similar idea to Codi and Scythe

Quantitative type inhabitation

- deductive synthesis that also deduces weights

Sketch refinement

- most similar to program repair

# Entity matching

[Singh et al. VLDB'17]

**Entity matching:** which rows correspond to the same person?

**Goal:** more interpretable results than existing approaches (e.g. decision trees)

**Search strategy:** Sketch + techniques for handling noise

(a)  $D_1$ : an instance of schema  $R$

	name	address	email	nation	gender
$r_1$	Catherine Zeta-Jones	9601 Wilshire Blvd., Beverly Hills, CA 90210-5213	c.jones@gmail.com	Wales	F
$r_2$	C. Zeta-Jones	3rd Floor, Beverly Hills, CA 90210	c.jones@gmail.com	US	F
$r_3$	Michael Jordan	676 North Michigan Avenue, Suite 293, Chicago		US	M
$r_4$	Bob Dylan	1230 Avenue of the Americas, NY 10020		US	M

(b)  $D_2$ : An instance of the schema  $S$

	name	apt	email	country	sex
$s_1$	Catherine Zeta-Jones	9601 Wilshire, 3rd Floor, Beverly Hills, CA 90210	c.jones@gmail.com	Wales	F
$s_2$	B. Dylan	1230 Avenue of the Americas, NY 10020	bob.dylan@gmail.com	US	M
$s_3$	Micheal Jordan	427 Evans Hall #3860, Berkeley, CA 94720	jordan@cs.berkeley.edu	US	M

# Entity matching

[Singh et al. VLDB'17]

