

# Lecture 8: Speed-up Synthesis with Abstract Semantics

Yu Feng  
Spring 2021

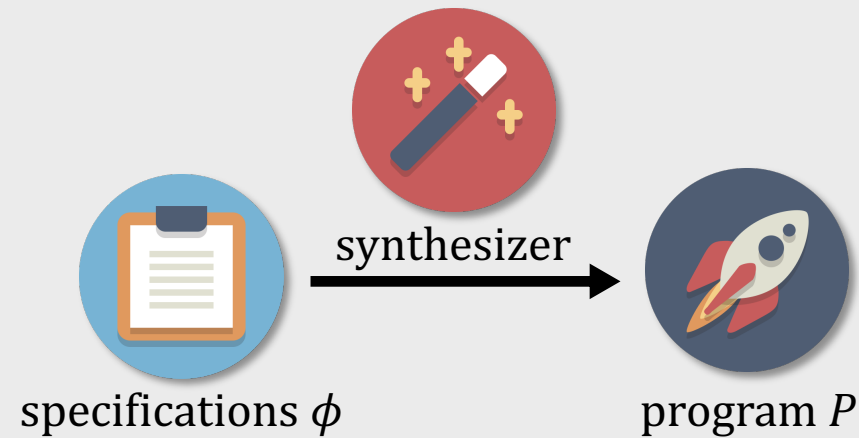


# Program Synthesis in a Nutshell

- Problem Formalization
- A Data Wrangling Example & DSL

Program Synthesis in a Nutshell

# Problem Formalization



**Find a program  $P$  that satisfies specifications  $\phi$ .**

# Problem Formalization

multi-modal

sample_ID	site	coll_date	species	TOT	inf_status
382870	site1	27/10/2007	SpeciesB	1	positive
382872	site2	27/10/2007	SpeciesB	1	negative
487405	site3	28/10/2007	SpeciesA	1	positive
487405	site3	28/10/2007	SpeciesA	1	positive

site	coll_date	SP_A_pos	SP_A_neg	SP_B_pos	SP_B_neg
site1	27/10/2007	0	0	1	0
site2	27/10/2007	0	0	0	1
site3	28/10/2007	2	0	0	0

examples

$\text{occurs}(\text{unite}) \wedge$   
 $\text{occurs}(\text{group\_by}) \wedge$   
 $\text{hasChild}(\text{group\_by}, \text{unite}) \wedge$

... logical constraints

I need to reformat the data so that there is just one row per site visit (i.e. in a given site name and date combo) with columns for total found by species and the fish status (i.e. speciesA\_pos, SpeciesA\_neg, Sp\_B\_pos.. etc).

natural languages

specifications  $\phi$

synthesizer

program  $P$

Find a program  $P$  that satisfies specifications  $\phi$ .

# Problem Formalization

## multi-modal

sample_ID	site	coll_date	species	TOT	inf_status
382870	site1	27/10/2007	SpeciesB	1	positive
382872	site2	27/10/2007	SpeciesB	1	negative
487405	site3	28/10/2007	SpeciesA	1	positive
487405	site3	28/10/2007	SpeciesA	1	positive

site	coll_date	SP_A_pos	SP_A_neg	SP_B_pos	SP_B_neg
site1	27/10/2007	0	0	1	0
site2	27/10/2007	0	0	0	1
site3	28/10/2007	2	0	0	0

examples

$\text{occurs}(\text{unite}) \wedge$   
 $\text{occurs}(\text{group\_by}) \wedge$   
 $\text{hasChild}(\text{group\_by}, \text{unite}) \wedge$

... logical constraints

I need to reformat the data so that there is just one row per site visit (i.e. in a given site name and date combo) with columns for total found by species and the fish status (i.e. speciesA\_pos, SpeciesA\_neg, Sp\_B\_pos.. etc).

natural languages

## multi-paradigm



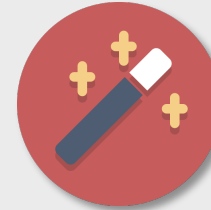
neural



deductive



...



synthesizer



specifications  $\phi$



program  $P$

**Find a program  $P$  that satisfies specifications  $\phi$ .**

# A Running Example from StackOverflow<sup>[1]</sup>

**[Title]** r script to count columns within dataset

**[Example]**

sample_ID	site	coll_date	species	TOT	inf_status
382870	site1	27/10/2007	SpeciesB	1	positive
382872	site2	27/10/2007	SpeciesB	1	negative
487405	site3	28/10/2007	SpeciesA	1	positive
487405	site3	28/10/2007	SpeciesA	1	positive



site	cat	sts
site1	SpeciesB_positive	1
site2	SpeciesB_negative	1
site3	SpeciesA_positive	2

**[Description]**

I need to reformat the data so that there is just one row per site visit (i.e. in a given site name and date combo) with columns for total found by species and the fish status (i.e. speciesA\_pos, SpeciesA\_neg, Sp\_B\_pos.. etc).

figured I need to sum within site. My thoughts were to use split/apply/aggregate/for loops etc but tried various combinations and not getting anywhere. apologies I'm not familiar with R. any comments appreciated!

[1] Example adapted from <https://stackoverflow.com/questions/39369502/r-script-to-reshape-and-count-columns-within-dataset>

# A Running DSL for Data Wrangling<sup>[1]</sup>

$t \rightarrow x_i$  (input table)  
`select( $t, \vec{c}_{arg}$ )` (column projection)  
`unite( $t, c_{tgt}, \vec{c}_{arg}$ )` (column merging)  
`separate( $t, \vec{c}_{tgt}, c_{arg}$ )` (column splitting)  
`mutate( $t, c_{tgt}, op, \vec{c}_{arg}$ )` (column arithmetic)  
`group_by( $t, \vec{c}_{arg}$ )` (row grouping)  
`summarise( $t, c_{tgt}, a, \vec{c}_{arg}$ )` (row aggregation)  
`filter( $t, f, \vec{c}_{arg}$ )` (row filtering)  
 $op \rightarrow + \mid - \mid \times \mid \div$   
 $a \rightarrow \min \mid \max \mid \text{sum} \mid \text{count} \mid \text{avg}$

$x_i$ : the  $i$ -th input table

$t$ : table

$c, \vec{c}$ : column(s) of table

$op$ : arithmetic operation

$a$ : aggregation function

$f$ : higher-order boolean function

A	B	C	D
A1	B1	1	5
A2	B2	2	6
A3	B3	3	7
A4	B4	4	8

select



A	C
A1	1
A2	2
A3	3
A4	4

A	B	C	D
A1	B1	1	5
A2	B2	2	6
A3	B3	3	7
A4	B4	4	8

unite



A_B	C	D
A1_B1	1	5
A2_B2	2	6
A3_B3	3	7
A4_B4	4	8

separate



A	B	C	D
A1	B1	1	5
A2	B2	2	6
A3	B3	3	7
A4	B4	4	8

mutate



A	B	C	D	C+D
A1	B1	1	5	6
A2	B2	2	6	8
A3	B3	3	7	10
A4	B4	4	8	12

A	B	C	D
X	B1	1	5
X	B2	2	6
Y	B3	3	7
Y	B4	4	8

group\_by



A	avg.D
X	5.5
Y	7.5

summarise

A	B	C	D
A1	B1	1	5
A2	B2	2	6
A3	B3	3	7
A4	B4	4	8

filter



A	B	C	D
A1	B1	1	5
A2	B2	2	6
A3	B3	3	7

[1] DSL adapted from Wang. C. et al. Visualization by Example. POPL'20

# A Running Example from StackOverflow

[Example]

sample_ID	site	coll_date	species	TOT	inf_status
382870	site1	27/10/2007	SpeciesB	1	positive
382872	site2	27/10/2007	SpeciesB	1	negative
487405	site3	28/10/2007	SpeciesA	1	positive
487405	site3	28/10/2007	SpeciesA	1	positive

unite

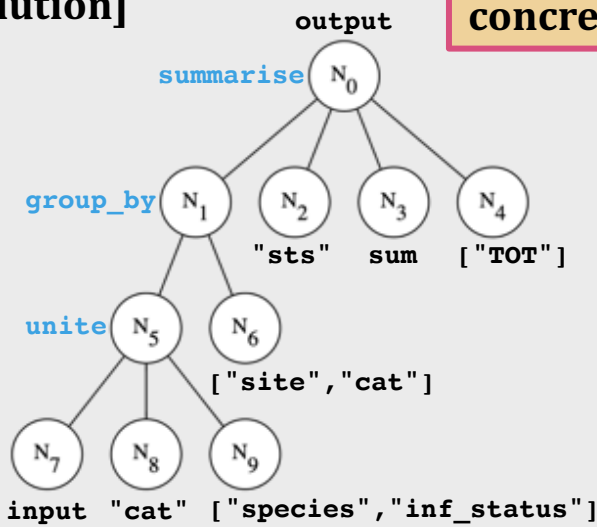
sample_ID	site	coll_date	cat	TOT
382870	site1	27/10/2007	SpeciesB_positive	1
382872	site2	27/10/2007	SpeciesB_negative	1
487405	site3	28/10/2007	SpeciesA_positive	1
487405	site3	28/10/2007	SpeciesA_positive	1

group\_by  
summarise

site	cat	sts
site1	SpeciesB_positive	1
site2	SpeciesB_negative	1
site3	SpeciesA_positive	2

[Solution]

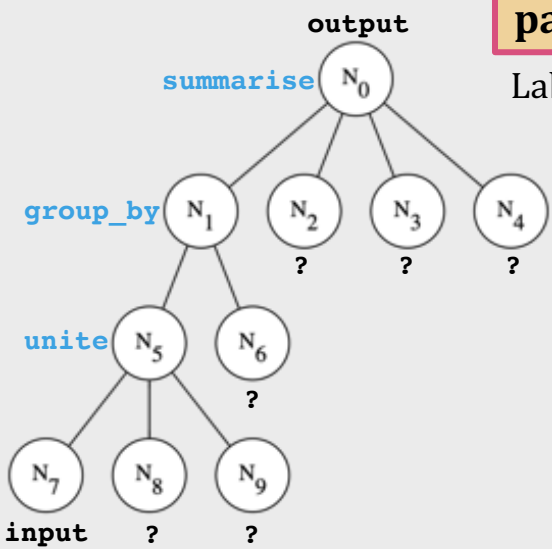
concrete program



```
T0 = unite( input, "cat", ["species", "inf_status"] )
T1 = group_by( T0, ["site", "cat"] )
output = summarise( T1, "sts", sum, ["TOT"] )
```

partial program / sketch

Labels of some AST nodes are yet to be determined.



```
T0 = unite( input, ?, ? )
T1 = group_by( T0, ? )
output = summarise( T1, ?, ?, ? )
```



# Terminologies

site	cat	sts
site1	SpeciesB_positive	1
site2	SpeciesB_negative	1
site3	SpeciesA_positive	2

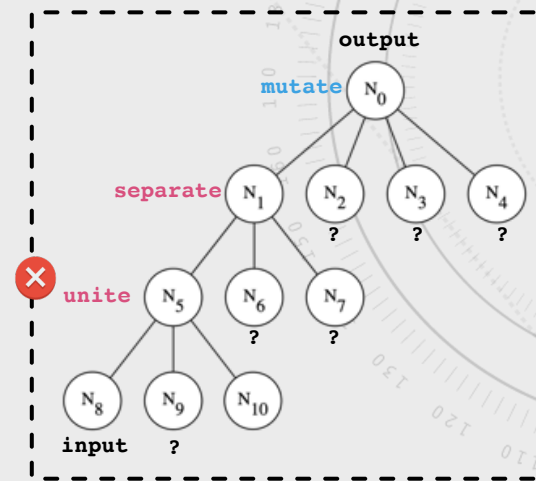
`output.row == 3`  
`output.col == 3`

sample_ID	site	coll_date	species	TOT	inf_status
382870	site1	27/10/2007	SpeciesB	1	positive
382872	site2	27/10/2007	SpeciesB	1	negative
487405	site3	28/10/2007	SpeciesA	1	positive
487405	site3	28/10/2007	SpeciesA	1	positive

`input.row == 4`  
`input.col == 6`



synthesizer



A Partial Program

$t \rightarrow x_i$

- | `select`( $t, \vec{c}_{arg}$ )
- | `unite`( $t, c_{tgt}, \vec{c}_{arg}$ )
- | `separate`( $t, \vec{c}_{tgt}, c_{arg}$ )
- | `mutate`( $t, c_{tgt}, op, \vec{c}_{arg}$ )
- | `group_by`( $t, \vec{c}_{arg}$ )
- | `summarise`( $t, c_{tgt}, a, \vec{c}_{arg}$ )
- | `filter`( $t, f, \vec{c}_{arg}$ )

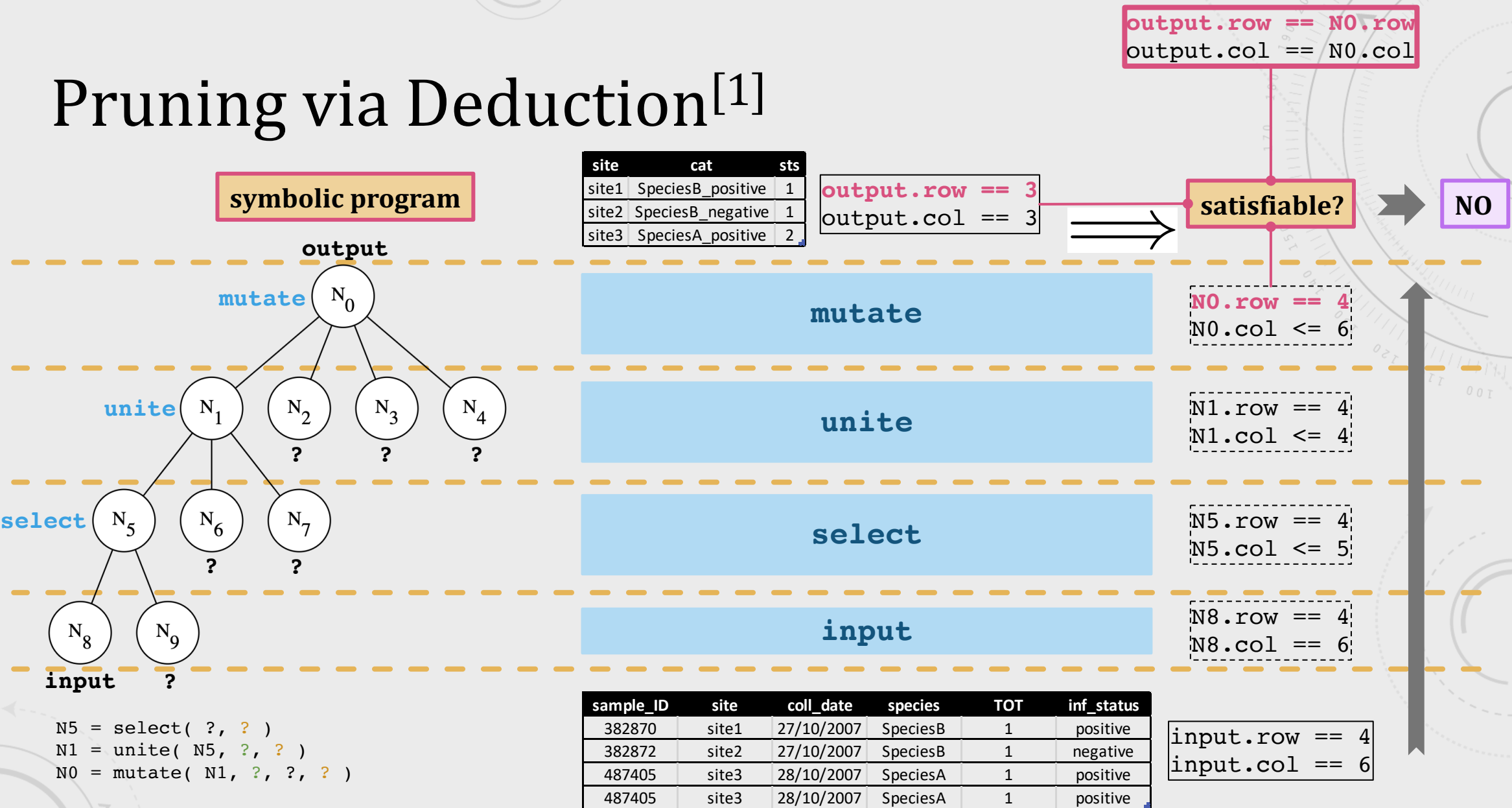
$op \rightarrow + \mid - \mid \times \mid \div$

$a \rightarrow \min \mid \max \mid \text{sum} \mid \text{count} \mid \text{avg}$

<b>select</b>	<code>out.row==in.row <math>\wedge</math> out.col&lt;=in.col-1</code>
<b>unite</b>	<code>out.row==in.row <math>\wedge</math> out.col==in.col-1</code>
<b>separate</b>	<code>out.row==in.row <math>\wedge</math> out.col==in.col+1</code>
<b>mutate</b>	<code>out.row==in.row <math>\wedge</math> out.col==in.col+1</code>
<b>group_by</b>	<code>out.row==in.row <math>\wedge</math> out.col==in.col</code>
<b>summarise</b>	<code>out.row&lt;=in.row <math>\wedge</math> out.col&lt;=in.col+1</code>
<b>filter</b>	<code>out.row&lt;=in.row-1 <math>\wedge</math> out.col==in.col</code>

Abstract Semantics of Components

# Pruning via Deduction<sup>[1]</sup>



[1] Feng, Y. et al.. Component-based Synthesis of Table Consolidation and Transformation Tasks from Examples. PLDI'17

# Learning from Past Mistakes

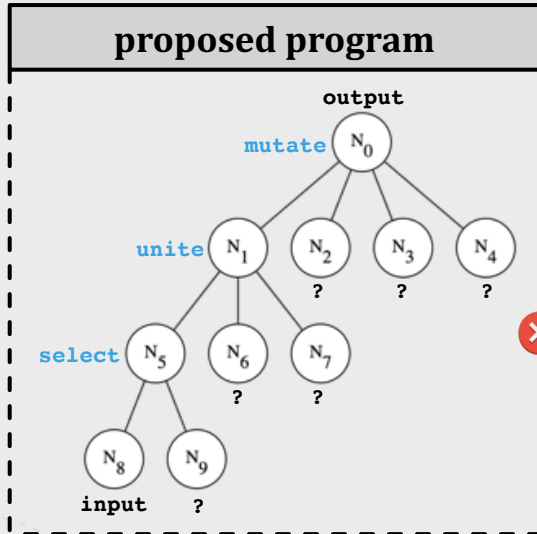
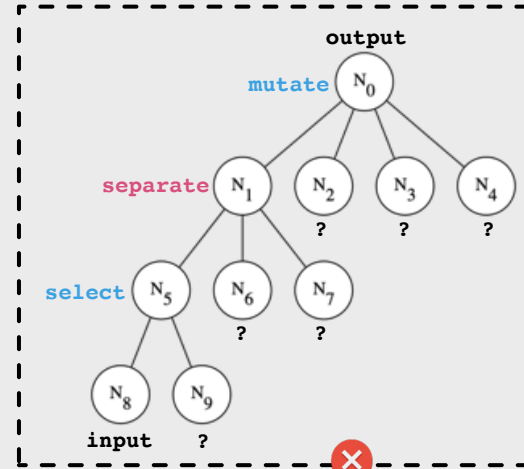
## Equivalent Modulo Conflict (EMC)<sup>[1]</sup>

<b>select</b>	$\text{out.row} == \text{in.row} \wedge \text{out.col} \leq \text{in.col} - 1$
<b>unite</b>	$\text{out.row} == \text{in.row} \wedge \text{out.col} == \text{in.col} - 1$
<b>separate</b>	$\text{out.row} == \text{in.row} \wedge \text{out.col} == \text{in.col} + 1$
<b>mutate</b>	$\text{out.row} == \text{in.row} \wedge \text{out.col} == \text{in.col} + 1$
<b>group_by</b>	$\text{out.row} == \text{in.row} \wedge \text{out.col} == \text{in.col}$
<b>summarise</b>	$\text{out.row} \leq \text{in.row} \wedge \text{out.col} \leq \text{in.col} + 1$
<b>filter</b>	$\text{out.row} \leq \text{in.row} - 1 \wedge \text{out.col} == \text{in.col}$

Abstract Semantics of Components

site	cat	sts
site1	SpeciesB_positive	1
site2	SpeciesB_negative	1
site3	SpeciesA_positive	2

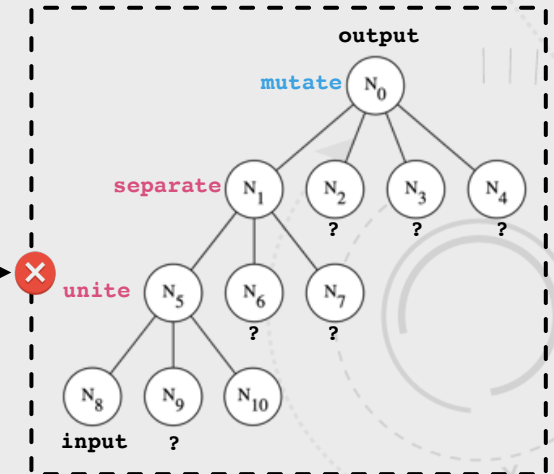
$\text{output.row} == 3$   
 $\text{output.col} == 3$



## Reasoning via Abstract Semantics

$\text{input.row} == 4 \wedge \text{input.col} == 6 \wedge$   
 $\text{N8.row} == \text{input.row} \wedge \text{N8.col} == \text{input.col} \wedge$   
 $\text{N5.row} == \text{N8.row} \wedge \text{N5.col} \leq \text{N8.col} - 1 \wedge$   
 $\text{N1.row} == \text{N5.row} \wedge \text{N1.col} == \text{N5.col} - 1 \wedge$   
 $\text{N0.row} == \text{N1.row} \wedge \text{N0.col} == \text{N1.col} + 1 \wedge$   
 $\text{output.row} == \text{N0.row} \wedge \text{output.col} == \text{N0.col} \wedge$   
 $\text{output.row} == 3 \wedge \text{output.col} == 3$

(input example)  
 (input alignment)  
 (select semantics)  
 (unite semantics)  
 (mutate semantics)  
 (output alignment)  
 (output example)

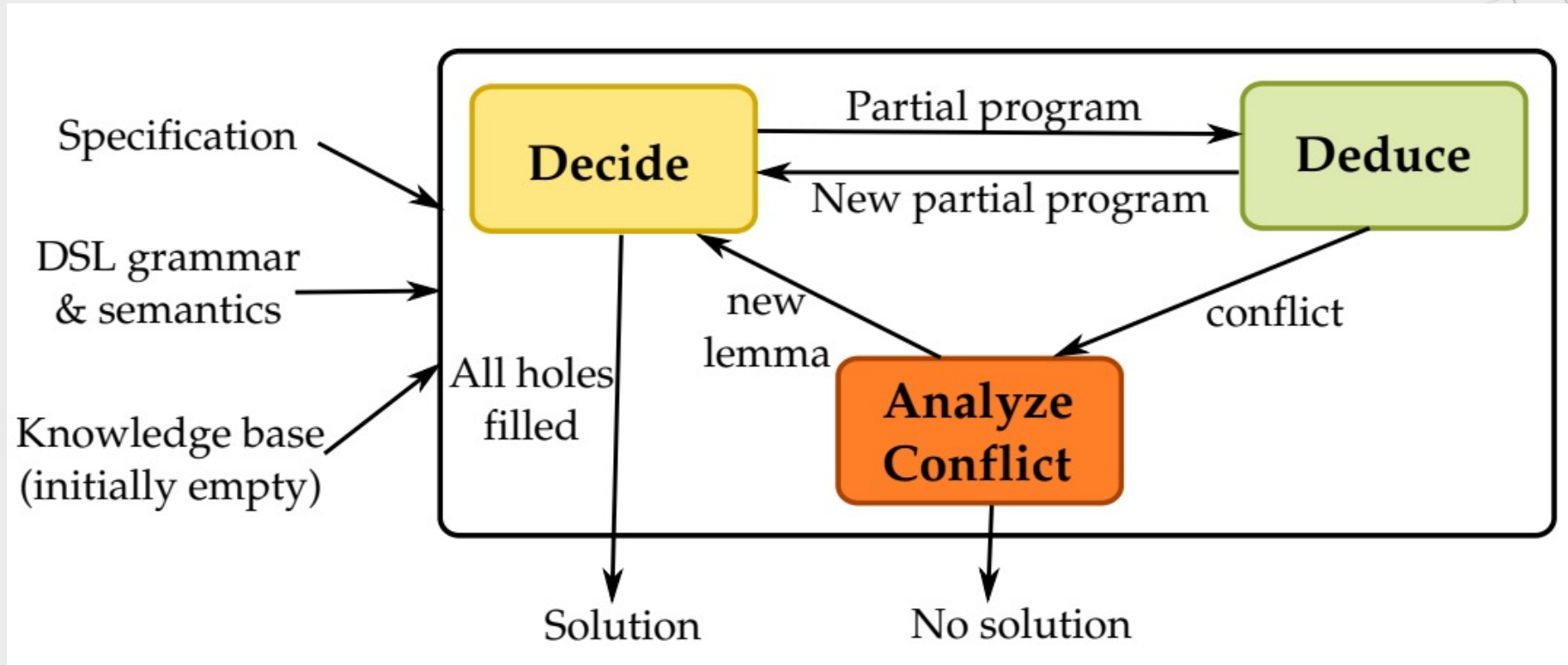


sample_ID	site	coll_date	species	TOT	inf_status
382870	site1	27/10/2007	SpeciesB	1	positive
382872	site2	27/10/2007	SpeciesB	1	negative
487405	site3	28/10/2007	SpeciesA	1	positive
487405	site3	28/10/2007	SpeciesA	1	positive

$\text{input.row} == 4$   
 $\text{input.col} == 6$

UNSAT Core

# Overview of NEO

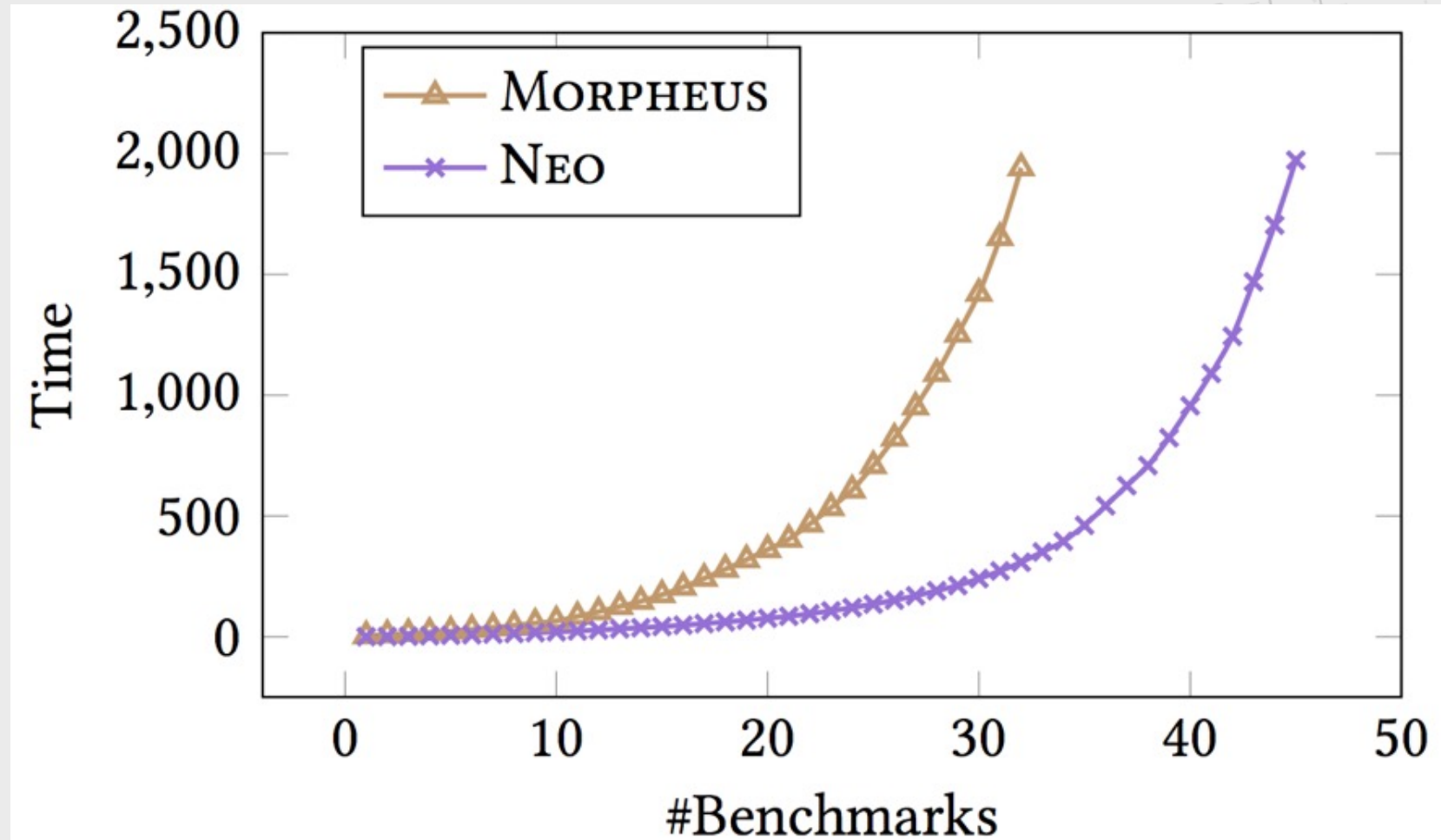


# Evaluation Setup

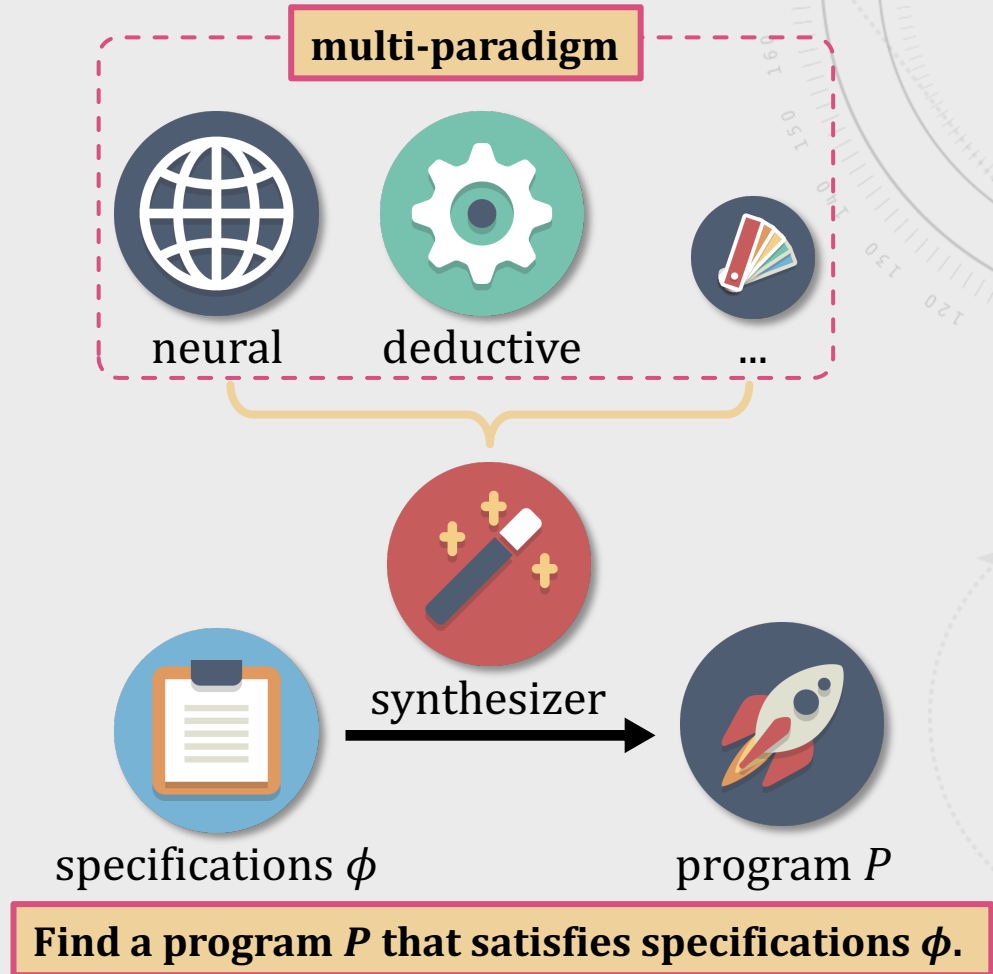
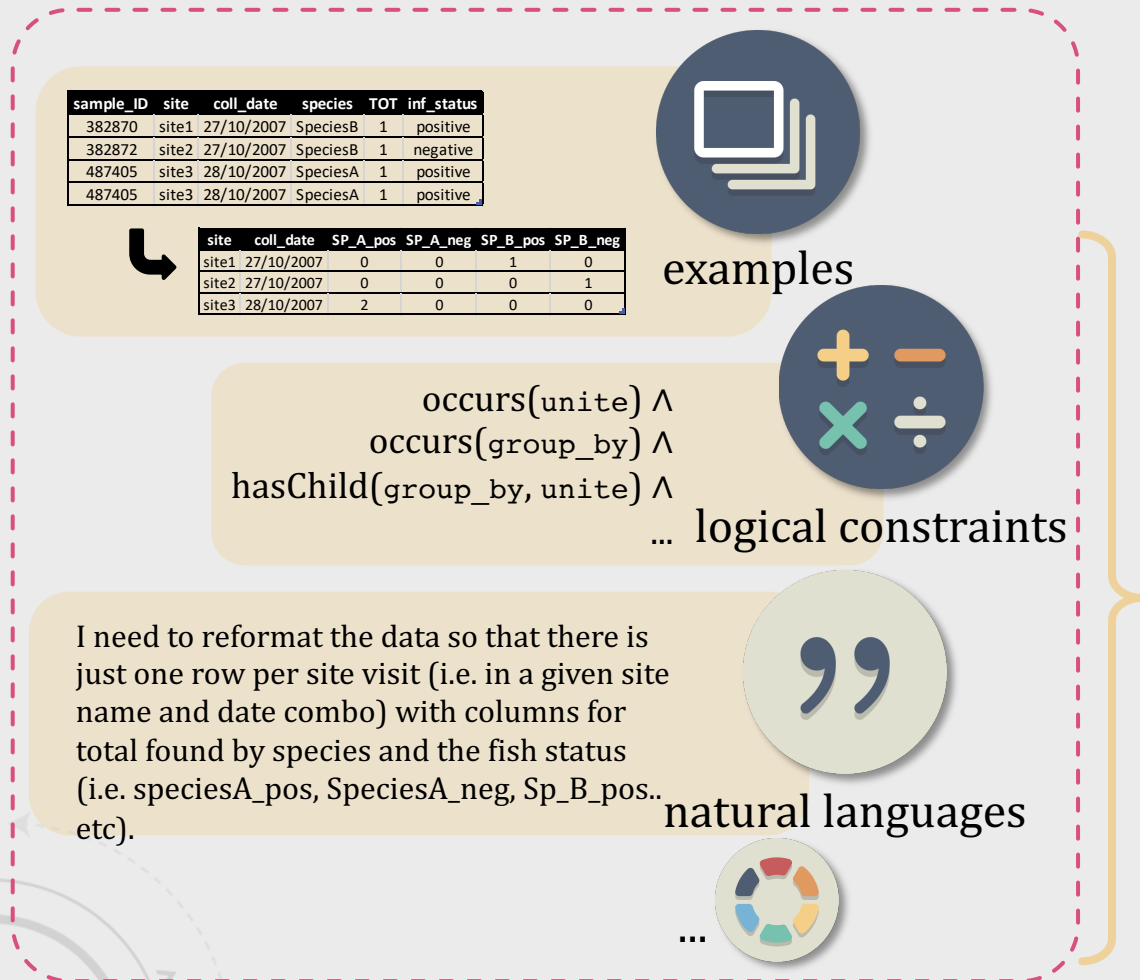
- Experiment Setup
  - Benchmarks: 50 Real-World Challenging Data Wrangling Tasks
  - Uses 2-gram model to guide the search
  - Comparison to MORPHEUS<sup>[1]</sup>
  - For more evaluation details, please refer to our paper

# Evaluation Results

- Timeout: 5 mins



# Caveats: Scalability and Over-fitting



# Challenges, Conclusions & Future Works

DEEPCODER (Balog et al. 2017); EXEC (Chen et al. 2018); NEO (Feng et al. 2018); SQLIZER (Yaghmazadeh et al. 2018); AutoPandas (Bavishi et al. 2019); METAL (Si et al. 2019); SKETCHADAPT (Nye et al. 2018); PROBE (Barke et al. 2020); CONCORD (Chen et al. 2020); ...

