

Spreadsheet Programming using Examples



Sumit Gulwani



Keynote at SEMS

July 2016

Motivation



99% of computer users cannot program!
They struggle with simple repetitive tasks.

Programming by examples (PBE)
can revolutionize this landscape!

Spreadsheet help forums



Excel Forum  Download Excel forum android application

MS Office Application Help

Forum	What's New?	Members List	Calendar	Forum Rules	Winner Board
--------------	-------------	--------------	----------	-------------	--------------

Today's Posts FAQ Calendar Community ▾ Forum Actions ▾ Quick Links ▾

 **ExcelExperts.com** Rectangular Site

Excel Consultancy, VBA Consultancy, Training and Tips Call:+442081234832

MREXCEL.COM 

Your One Stop for Excel Tips & Solutions

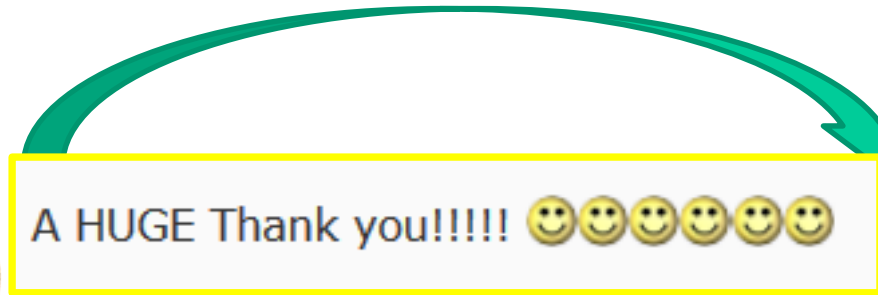
Forum	Search New Posts	Zero Reply Posts	Subscribed Threads
--------------	------------------	------------------	--------------------

FAQ Forum Actions ▾ Quick Links ▾

Typical help-forum interaction

300_w30_aniSh_c1_b → w30

300_w5_aniSh_c1_b → w5



=MID(B1,5,2)

=MID(B1,FIND("_",\$B:\$B)+1,
FIND("_",REPLACE(\$B:\$B,1,FIND("_",\$B:\$B),""))-1)

Flash Fill (Excel 2013 feature) demo



Excel 2013's coolest new feature that should have been available years ago

The screenshot shows the Microsoft Excel 2013 interface. The ribbon is set to 'DESIGN' under 'TABLE TOOLS'. The active cell is C6, containing the text 'Andrew'. The spreadsheet data is as follows:

	A	B	C	D
4		Email	First Name	
5		Nancy.FreeHafer@fourthcoffee.com	Nancy	
6		Andrew.Cencini@northwindtraders.com	Andrew	
7		Jan.Kotas@litwareinc.com	Jan	
8		Mariya.Sergienko@graphicdesigninstitute.com	Mariya	
9		Steven.Thorpe@northwindtraders.com	Steven	
10		Michael.Neipper@northwindtraders.com	Michael	
11		Robert.Zare@northwindtraders.com	Robert	
12		Laura.Giussani@adventure-works.com	Laura	
13		Anne.HL@northwindtraders.com	Anne	
14		Alexander.David@contoso.com	Alexander	
15		Kim.Shane@northwindtraders.com	Kim	
16		Manish.Chopra@northwindtraders.com	Manish	
17		Gerwald.Oberleitner@northwindtraders.com	Gerwald	
18		Amr.Zaki@northwindtraders.com	Amr	

At the bottom of the window, a taskbar shows the steps: Start, 1. Fill, 2. Analyze, 3. Chart, and Learn More. An 'ENTER' button is visible at the bottom left of the spreadsheet area.

“Automating string processing in spreadsheets using input-output examples”;
POPL 2011; Sumit Gulwani

Number Transformations

Input	Output (Round to 2 decimal places)
123.4567	123.46
123.4	123.40
78.234	78.23

Excel/C#: `#.00`

Python/C: `.2f`

Java: `###`

Input	Output (Nearest lower half hour)
0d 5h 26m	5:00
0d 4h 57m	4:30
0d 4h 27m	4:00
0d 3h 57m	3:30

Semantic String Transformations

MarkupRec Table

Id	Name	Markup
S33	Stroller	30%
B56	Bib	45%
D32	Diapers	35%
W98	Wipes	40%
A46	Aspirator	30%

CostRec Table

Id	Date	Price
S33	12/2010	\$145.67
S33	11/2010	\$142.38
B56	12/2010	\$3.56
D32	1/2011	\$21.45
W98	4/2009	\$5.12

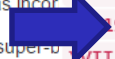
Input v₁	Input v₂	Output (Price + Markup*Price)
Stroller	10/12/2010	$\$145.67 + 0.30 * 145.67$
Bib	23/12/2010	$\$3.56 + 0.45 * 3.56$
Diapers	21/1/2011	
Wipes	2/4/2009	
Aspirator	23/2/2010	

Data Science Class Assignment



```
cat superbowl.txt | awk '$1=$1' ORS=' ' | sed 's/|-|/\n|/g' | grep "^| style=\"text-align: center;\"" | grep -v "Championship"
```

To get Started!



FlashExtract Demo

Ships inside two Microsoft products:



ConvertFrom-String cmdlet



Custom Log,
Custom Field

“FlashExtract: A Framework for data extraction by examples”;
PLDI 2014; Vu Le, Sumit Gulwani

Layout Transformations

Input Table

Bureau of I.A.	
Regional Director	Numbers
Niles C.	Tel: (800)645-8397
	Fax: (907)586-7252
Jean H.	Tel: (918)781-4600
	Fax: (918)781-4604
Frank K.	Tel: (615)564-6500
	Fax: (615)564-6701



Output Table

	Tel	Fax
Niles C.	(800)645-8397	(907)586-7252
Jean H.	(918)781-4600	(918)781-4604
Frank K.	(615)564-6500	(615)564-6701

	A	B	C	D	E
1		value	year	value	year
2	Albania	1000	1950	930	1981
3	Austria	3139	1951	3177	1955
4	Belgium	541	1947	601	1950
5	Bulgaria	2964	1947	1959	1958
6	Czech ...	2416	1950	2503	1960

...

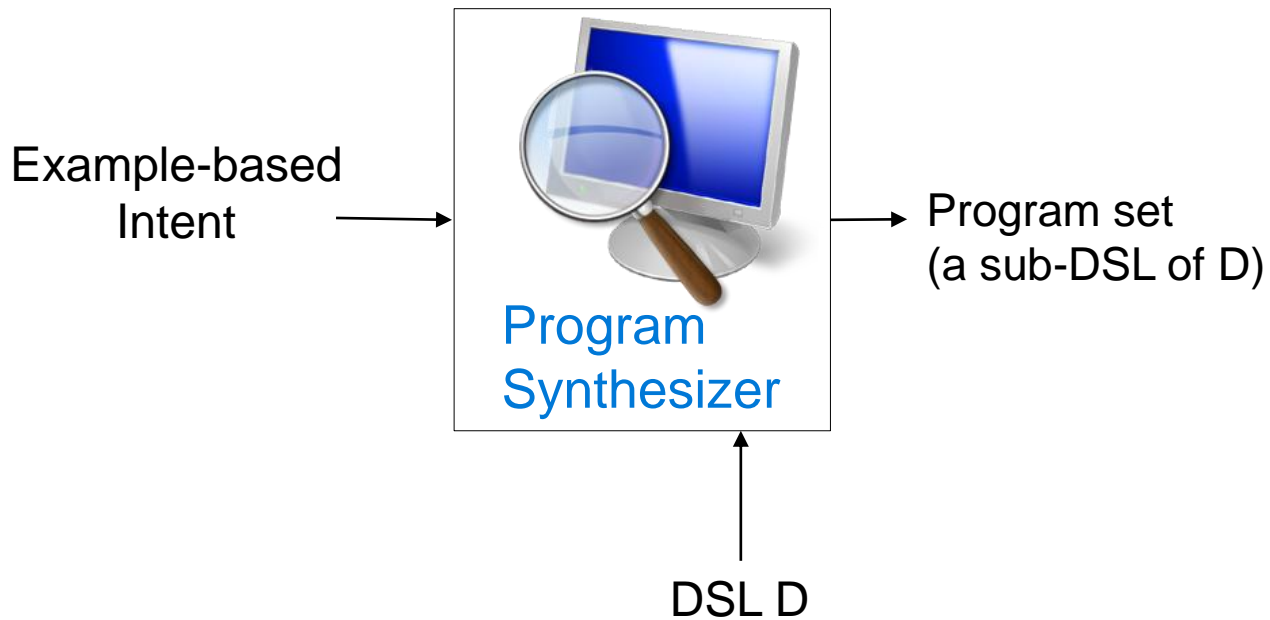


Country	Harvest	Date
Albania	1000	1950
Albania	930	1981
...		
Austria	3139	1951
Austria	3177	1955
...		
Belgium	541	1947
Belgium	601	1950

...

PBE allows creation of output table from couple of example tuples.

Programming-by-Examples Architecture



Domain-specific Language (DSL)

- **Balanced Expressiveness**
 - Expressive enough to cover wide range of tasks
 - Restricted enough to enable efficient search
 - Restricted set of operators
 - those with small inverse sets
 - Restricted syntactic composition of those operators
- **Natural computation patterns**
 - Increased user understanding/confidence
 - Enables selection between programs, editing

Flash Fill DSL (String Transformations)

$Tuple(String\ x_1, \dots, String\ x_n) \rightarrow String$

top-level expr $T :=$ `if-then-else`(B,C,T)
| C

condition-free expr $C :=$ `Concatenate`(A,C)
| A

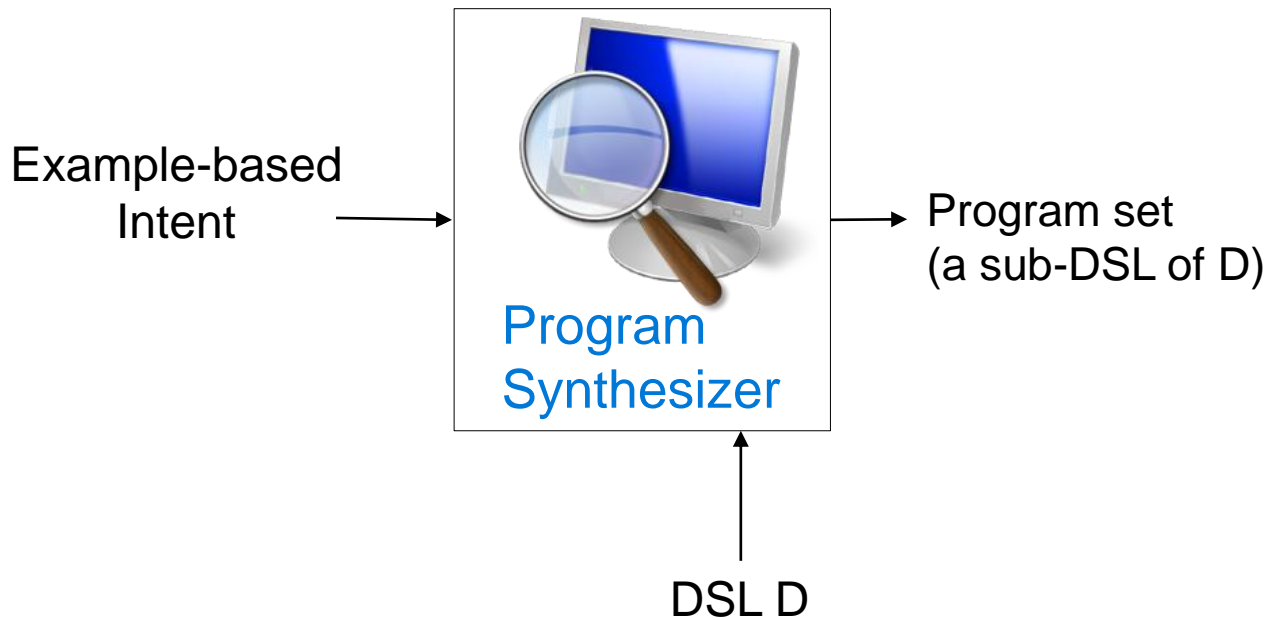
atomic expression $A :=$ `SubStr`(X,P,P)
| `ConstantString`

input string $X :=$ x_1 | x_2 | ...

position expression $P :=$ K
| `Pos`(X, R_1 , R_2 , K)
Kth position in X whose left/right side
matches with R_1/R_2 .

Boolean expression $B :=$...

Programming-by-Examples Architecture



Search Methodology

Goal: Set of program expr of kind e that satisfies spec ϕ
[denoted $[e \models \phi]$]

e : DSL (top-level) expression

ϕ : Conjunction of (input state $\sigma \rightsquigarrow$ output value v)

Methodology: Based on **divide-and-conquer** style problem decomposition.

- $[e \models \phi]$ is reduced to **simpler problems** (over sub-expressions of e or sub-constraints of ϕ).
- Top-down (as opposed to bottom-up enumerative search).

*"FlashMeta: A Framework for Inductive Program Synthesis";
OOPSLA 2015; Alex Polozov, Sumit Gulwani*

Problem Reduction Rules

$$[e \models \phi_1 \wedge \phi_2] = \text{Intersect}([e \models \phi_1], [e \models \phi_2])$$

An alternative strategy:

$$[e \models \phi_1 \wedge \phi_2] = \text{Filter}([e \models \phi_1], \phi_2)$$

Let e be a non-terminal defined as $e ::= e_1 \mid e_2$

$$[e \models \phi] = \text{Union}([e_1 \models \phi], [e_2 \models \phi])$$

Problem Reduction Rules

$$[F(e_1, \dots, e_n) \models \sigma \rightsquigarrow v] = \text{Union}(\{F([e_1 \models \sigma \rightsquigarrow u_1], \dots, [e_n \models \sigma \rightsquigarrow u_n]) \mid (u_1, \dots, u_n) \in F^{-1}(v)\})$$

$F(S_1, \dots, S_n)$ denotes $\{F(e_1, \dots, e_n) \mid e_1 \in S_1, \dots, e_n \in S_n\}$

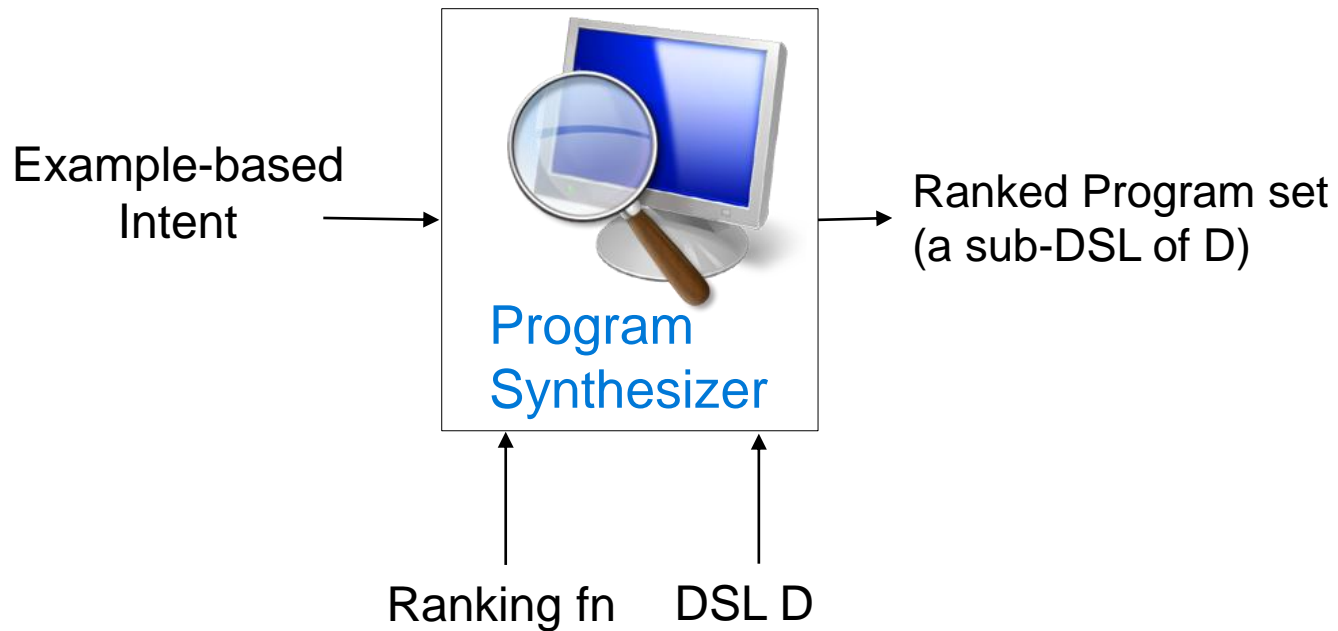
Inverse Set: Let F be an n -ary operator.

$$F^{-1}(v) = \{(u_1, \dots, u_n) \mid F(u_1, \dots, u_n) = v\}$$

$$\text{Concat}^{-1}(\text{"Abc"}) = \{(\text{"Abc"}, \epsilon), (\text{"Ab"}, \text{"c"}), (\text{"A"}, \text{"bc"}), (\epsilon, \text{"Abc"})\}$$

$$[\text{Concat}(X, Y) \models (\sigma \rightsquigarrow \text{"Abc"})] = \text{Union}(\{ \text{Concat}([X \models (\sigma \rightsquigarrow \text{"Abc"})], [Y \models (\sigma \rightsquigarrow \epsilon)]), \text{Concat}([X \models (\sigma \rightsquigarrow \text{"Ab"})], [Y \models (\sigma \rightsquigarrow \text{"c"})]), \text{Concat}([X \models (\sigma \rightsquigarrow \text{"A"})], [Y \models (\sigma \rightsquigarrow \text{"bc"})]), \text{Concat}([X \models (\sigma \rightsquigarrow \epsilon)], [Y \models (\sigma \rightsquigarrow \text{"Abc"})]) \})$$

Programming-by-Examples Architecture



Ranking scheme: Program features

Prefer simpler programs

- Fewer constants.
- Smaller constants.

Input	Output
Rishabh Singh	Rishabh
Ben Zorn	Ben



- 1st Word
- If (input = “Rishabh Singh”) then “Rishabh” else “Ben”
- “Rishabh”

“Predicting a correct program in Programming by Example”;
[CAV 2015] Rishabh Singh, Sumit Gulwani

Ranking scheme: Data features

Prefer simpler programs

- Fewer constants.
- Smaller constants.



Input	Output
Missing page numbers, 1993 64-67, 1995	1993 1995

- 1st Number from the beginning
- 1st Number from the end

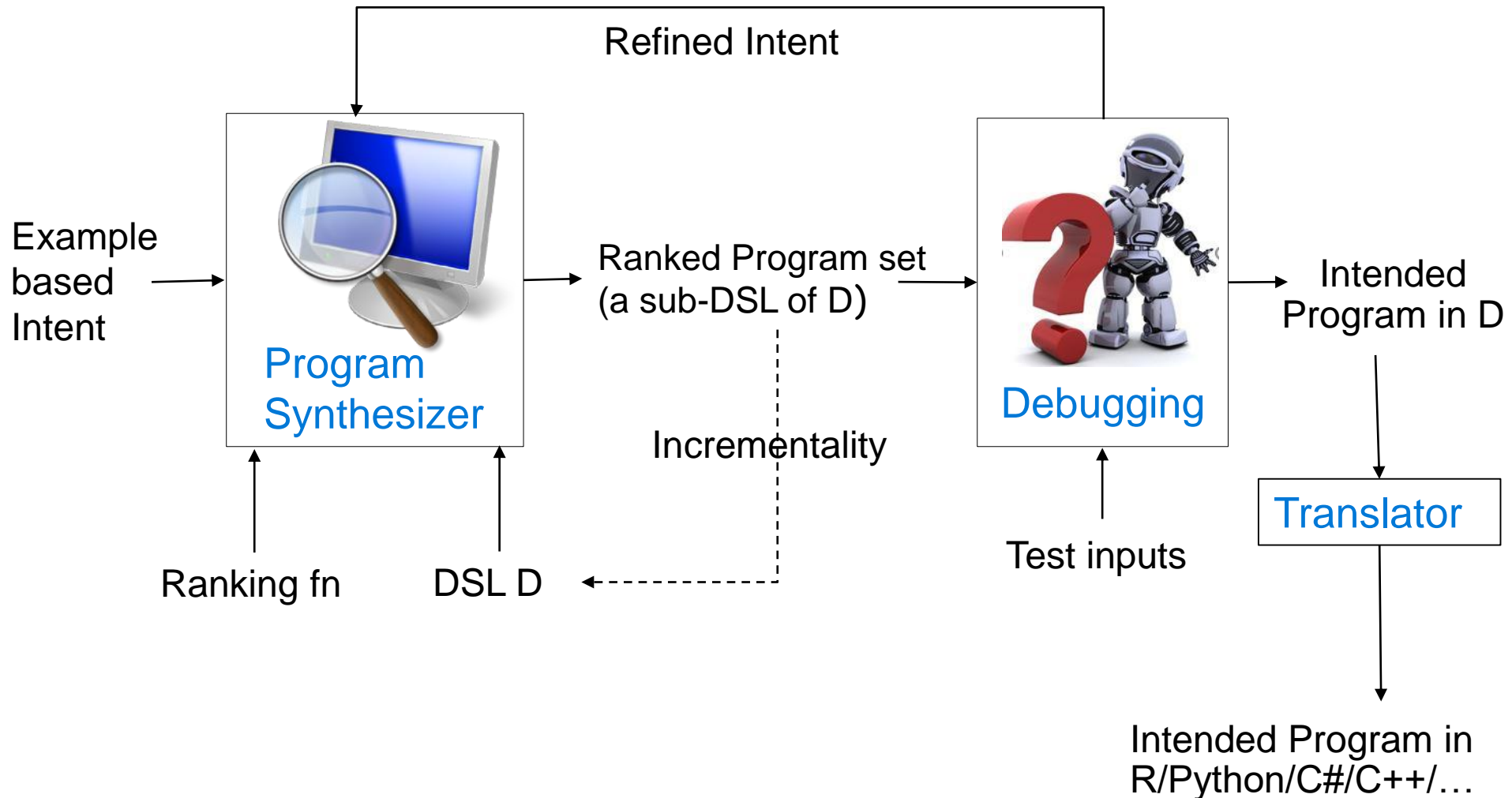
How to select between programs with
same number of same-sized constants?

Prefer programs that generate more uniform output.

Outline

- Core Synthesis Architecture
 - Domain-specific Language
 - Search methodology
 - Ranking function
- Next generation Synthesis
 - Interactive
 - Predictive
 - Adaptive

Programming-by-Examples Architecture



Interactive Debugging

- Sampling inputs
- Asking questions
- Visual explanations of the synthesized program

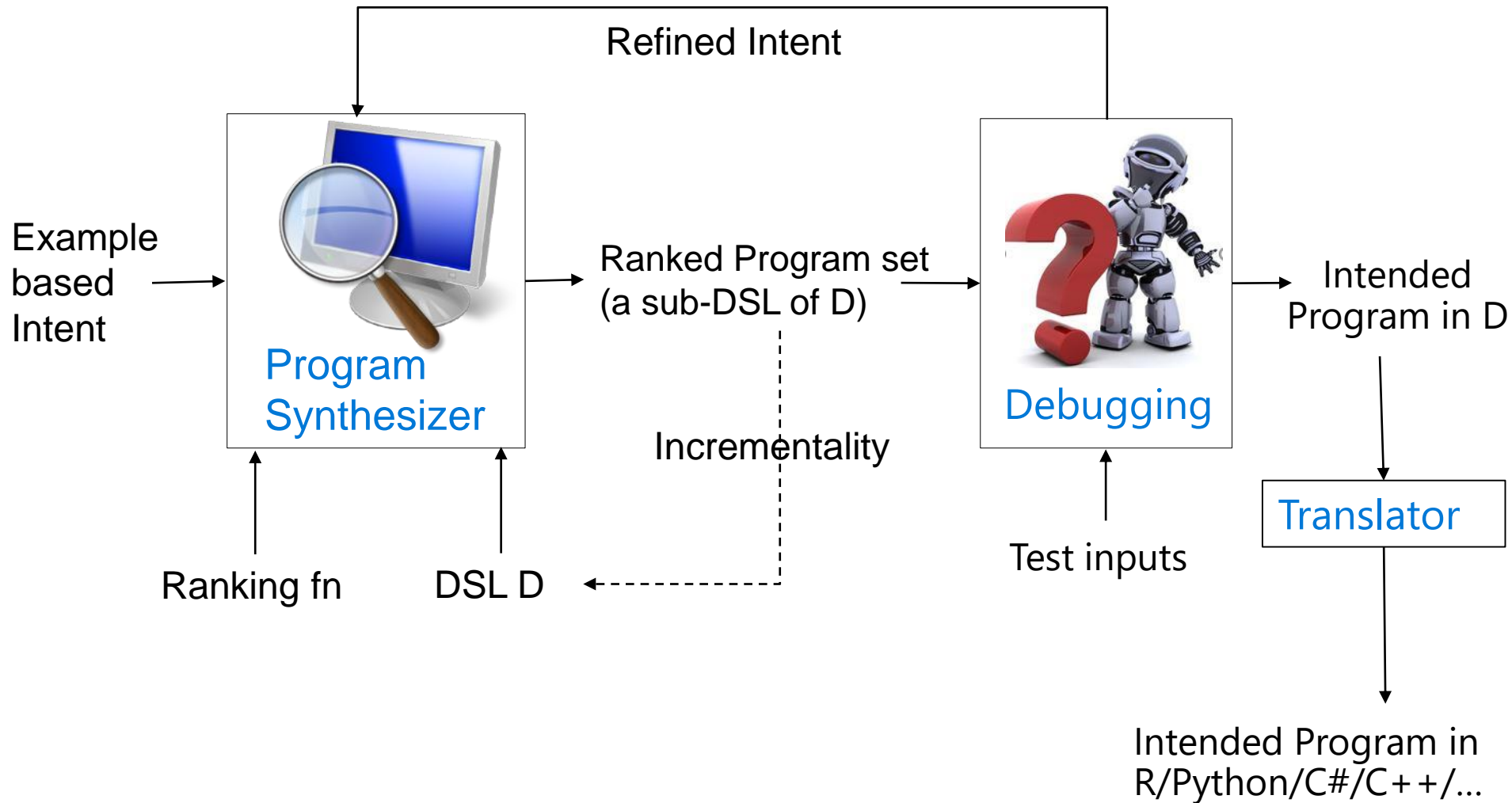


Predictive



- Intended programs can sometimes be synthesized from just the input.
 - Tabular data extraction, Sort, Join
- Can save large amount of user effort.
 - User need not provide examples for each of tens of columns.

Programming-by-Examples Architecture

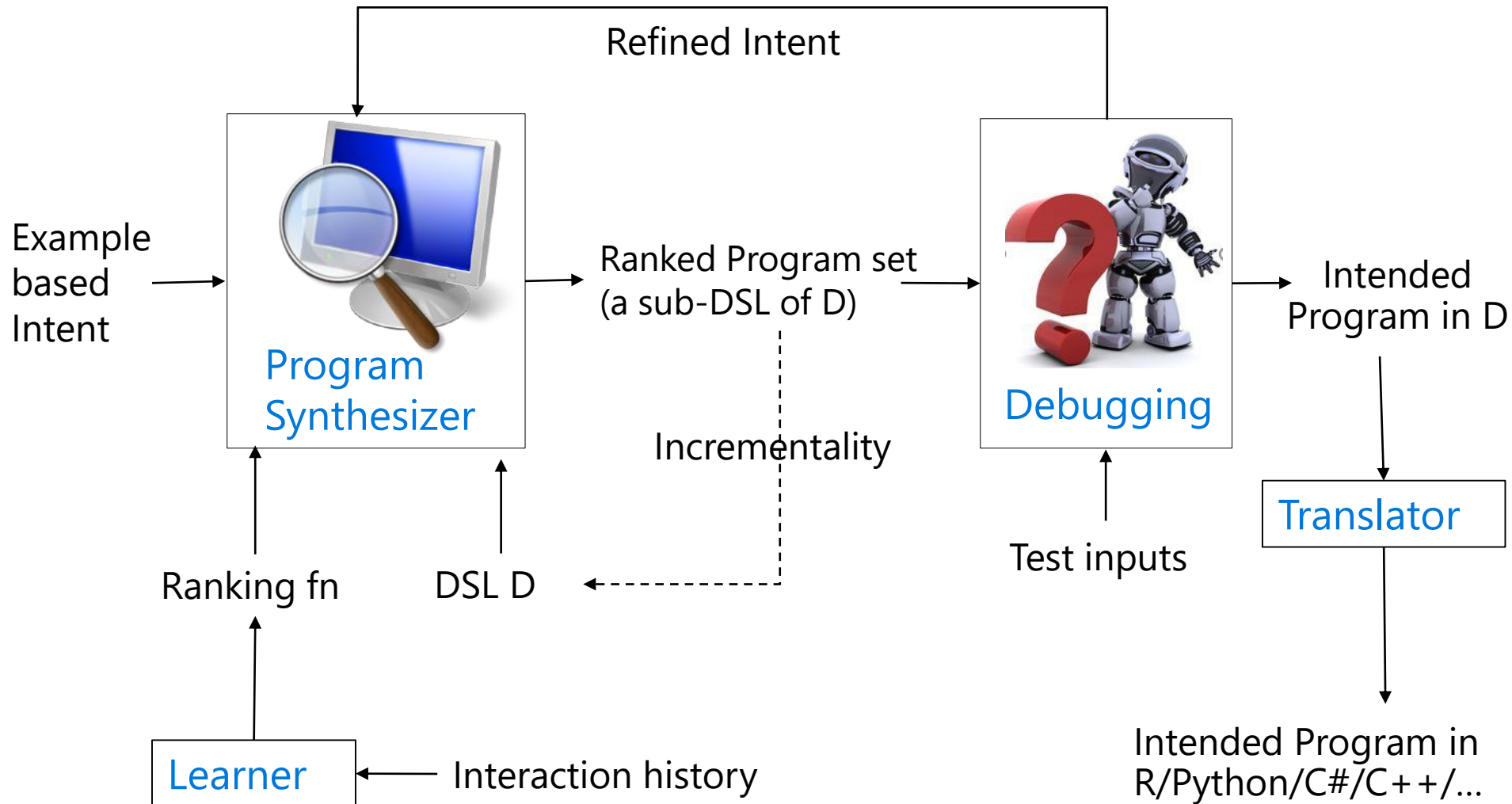


Adaptive



- Learn from past interactions
 - of the same user (personalized experience).
 - of other users in the enterprise/cloud.
- The synthesis sessions now require less interaction.

Programming-by-Examples Architecture



PROSE Framework

<https://microsoft.github.io/prose>

- Efficient implementation of the generic search methodology.
- Provides a library of reduction rules.

Role of synthesis designer

- Implement a DSL and provide reduction rules for new operators.
- Provide ranking strategy.
- Can specify tactics to resolve non-determinism in search.

The PROSE Team



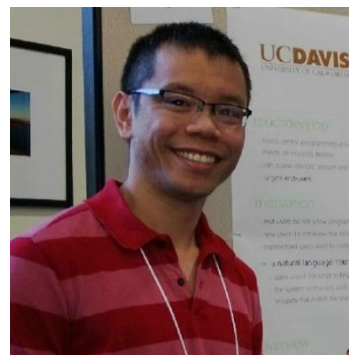
Allen
Cypher



Sumit
Gulwani



Ranvijay
Kumar



Vu Le



Daniel
Perelman

We are hiring interns/full-time!



Alex
Polozov



Mohammad
Raza



Danny
Simmons



Adam
Smith



Abhishek
Udupa

Future Directions

- Learn from usage data
- Probabilistic noise handling
- Programming using natural language
- Application to robotics

Conclusion

- PBE can enable easier & faster data wrangling.
 - 99% of computer users are non-programmers.
 - Data scientists spend 80% time cleaning data.
- Algorithmic search
 - Domain-specific language
 - Deductive methodology based on back-propagation
- Ambiguity resolution
 - Ranking
 - Interactivity



Reference: “*Programming by Examples (and its applications in Data Wrangling)*”,
In Verification and Synthesis of Correct and Secure Systems; IOS Press; 2016
[based on Marktoberdorf Summer School 2015 Lecture Notes]