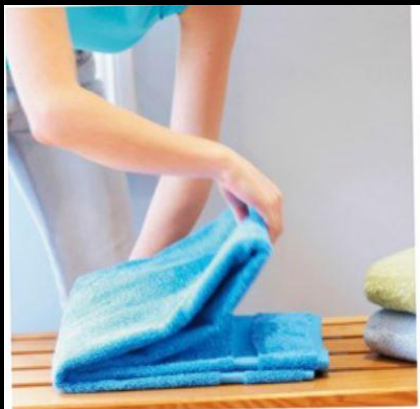# Programming by Examples
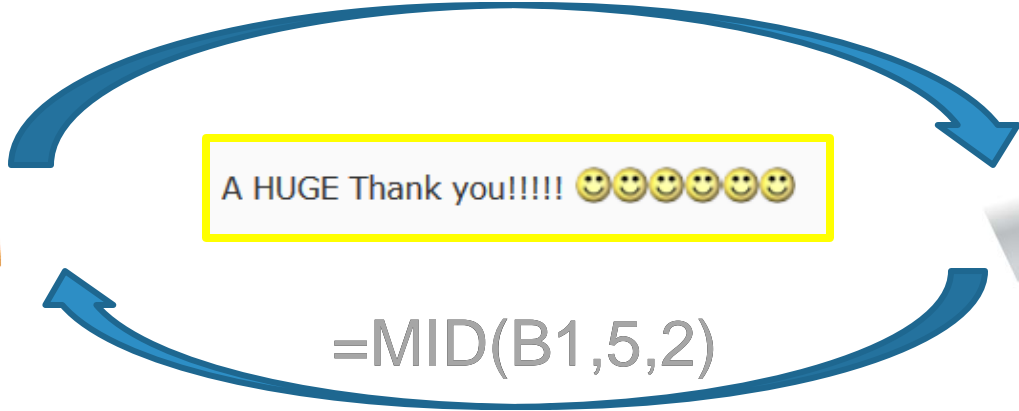
Sumit Gulwani

Microsoft

ECML/PKDD Conference

Sep 2019

# Example-based help-forum interaction

300_w30_aniSh_c1_b → w30
300_w5_aniSh_c1_b → w5

A HUGE Thank you!!!!! 😊😊😊😊😊

=MID(B1,5,2)

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

# Flash Fill (Excel feature)



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

| | A | B |
|---|---|---|
| 1 | **Email** | **Column 2** |
| 2 | Nancy.FreeHafer@fourthcoffee.com | freehafer |
| 3 | Andrew.Cencici@northwindtraders.com | cencici |
| 4 | Jan.Kotas@litwareinc.com | kotas |
| 5 | Mariya.Sergienko@gradicdesigninstitute.com | sergienko |
| 6 | Steven.Thorpe@northwindtraders.com | thorpe |
| 7 | Michael.Neipper@northwindtraders.com | neipper |
| 8 | Robert.Zare@northwindtraders.com | zare |
| 9 | Laura.Giussani@adventure-works.com | giussani |
| 10 | Anne.HL@northwindtraders.com | hl |
| 11 | Alexander.David@contoso.com | david |
| 12 | Kim.Shane@northwindtraders.com | shane |

*"Automating string processing in spreadsheets using input-output examples"*
[POPL 2011] Sumit Gulwani

3

Darin
@crushspread

AI is going to take over the world... and this is what Excel auto-populated today.

| K | L | M | N |
|---|---|---|---|
| | DEC | December | |
| | NOV | November | |
| | OCT | Octember | |
| | APR | Aprember | |
| | AUG | Augember | |
| | FEB | Febember | |
| | JAN | Janember | |
| | JUL | Julember | |
| | JUN | Junember | |
| | MAR | Marember | |
| | MAY | Mayember | |
| | SEP | Sepember | |

5:00 AM · Oct 23, 2018 · Twitter for iPhone

**12.5K** Retweets   **32K** Likes

|   | A | B | C |
|---|---|---|---|
| 1 | DEC | December | |
| 2 | NOV | November | |
| 3 | OCT | Octember | |
| 4 | APR | Aprember | |
| 5 | AUG | Augember | |
| 6 | FEB | Febember | |
| 7 | JAN | Janember | |
| 8 | JUL | Julember | |
| 9 | JUN | Junember | |
| 10 | MAR | Marember | |
| 11 | MAY | Mayember | |
| 12 | SEP | Sepember | |
| 13 | | | |

# Number, DateTime 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 (3-hour weekday bucket) |
|-------|--------------------------------|
| CEDAR AVE & COTTAGE AVE; HORSHAM; 2015-12-11 @ 13:34:52; | Fri, 12PM - 3PM |
| RT202 PKWY; MONTGOMERY; 2016-01-13 @ 09:05:41-Station:STA18; | Wed, 9AM - 12PM |
| ; UPPER GWYNEDD; 2015-12-11 @ 21:11:18; | Fri, 9PM - 12AM |

[CAV 2012] *"Synthesizing Number Transformations from Input-Output Examples"*; Singh, Gulwani
[POPL 2015] *"Transforming Spreadsheet data types using Examples"*; Singh, Gulwani

# Table Extraction

```
cat superbowl.txt | awk '$1=$1' ORS=' ' | sed 's/|- |/\n|/g' | grep "^| style=\"t
ext-align: center;\"" | grep -v "Championship"
```

*"FlashExtract: A Framework for data extraction by examples"*
[PLDI 2014] Vu Le, Sumit Gulwani

# Table Reshaping

| Bureau of I.A. | |
|---|---|
| Regional Dir. | 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 |

FlashRelate

From few examples of rows in 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 |

50% spreadsheets are semi-structured.
KPMG, Deloitte budget millions of dollars for normalization.

# PBE Architecture



Examples →
DSL D →
Search Engine

Examples

Program set

Program Ranker

Ranked Program set

Disambiguator

Test inputs

Intended Program (in D)

Huge search space
- Prune using Logical reasoning
- Guide using Machine learning

Under-specification
- Guess using Ranking (PL features, ML models)
- Interact: leverage extra inputs (clustering) and programs (execution)

# Flash Fill DSL

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

top-level expr $T := C \mid ifThenElse(B, C, T)$

condition-free expr $C := A \mid Concat(A, C)$

atomic expression $A := SubStr(X, P, P) \mid ConstantString$

input string $X := x_1 \mid x_2 \mid \ldots$

position expression $P := K \mid Pos(X, R_1, R_2, K)$

$K^{th}$ position in X whose left/right side matches with $R_1$/$R_2$.

# Search Idea 1: Deduction

Let $[G \vDash \phi]$ denote programs in grammar G that satisfy spec $\phi$

$\phi$ is a Boolean constraint over (input state $i \rightsquigarrow$ output value $o$)

**Divide-and-conquer style problem reduction**

$$[G \vDash \phi_1 \wedge \phi_2] = Intersect([G \vDash \phi_1], [G \vDash \phi_2])$$
$$= [G_1 \vDash \phi_2] \text{ where } G_1 = [G \vDash \phi_1]$$

Let $\text{G} \coloneqq G_1 \mid G_2$

$$[G \vDash \phi] = [G_1 \vDash \phi] \mid [G_2 \vDash \phi]$$

# Search Idea 1: Deduction

Inverse Set: $F^{-1}(o) \stackrel{\text{def}}{=} \{ (u, v) \mid F(u, v) = o \}$

E.g. $Concat^{-1}("Abc") = \{ ("A", "bc"), ("Ab", "c"), \dots \}$

Let $G := F(G_1, G_2)$

Let $F^{-1}(o)$ be $\{ (u, v), (u', v') \}$

$[G \vDash (i \rightsquigarrow o)] = F([G_1 \vDash (i \rightsquigarrow u)], [G_2 \vDash (i \rightsquigarrow v)])$
$\qquad\qquad\qquad | \; F([G_1 \vDash (i \rightsquigarrow u')], [G_2 \vDash (i \rightsquigarrow v')])$

# Search Idea 2: Learning

Machine Learning for ordering search
- Which grammar production to try first?
- Which sub-goal resulting from inverse semantics to try first?

Prediction based on supervised training
- standard LSTM architecture
- Training: 100s of tasks, 1 task yields 1000s of sub-problems.
- Results: Up to 20x speedup with average speedup of 1.67

# Ranking Idea 1: Program Features

| Input | Output |
|---|---|
| Vasu Singh | v.s. |
| Stuart Russell | s.r. |

P1:  Lower($1^{st}$ char) + ".s."
P2:  Lower($1^{st}$ char) + "." + $3^{rd}$ char + "."
P3:  Lower($1^{st}$ char) + "." + Lower($1^{st}$ char after space) + "."

Prefer programs (P3) with simpler Kolmogorov complexity
- Fewer constants
- Smaller constants

# Ranking Idea 2: Output Features

| Input | Output | Output of P1 |
|-------|--------|--------------|
| [CPT-123 | [CPT-123] | [CPT-123] |
| [CPT-456] | [CPT-456] | [CPT-456]] |

P1:  Input + "]"

P2:  Prefix of input upto 1st number + "]"

Examine features of outputs of a program on extra inputs:

- IsYear, Numeric Deviation, # of characters, IsPerson

# Disambiguation

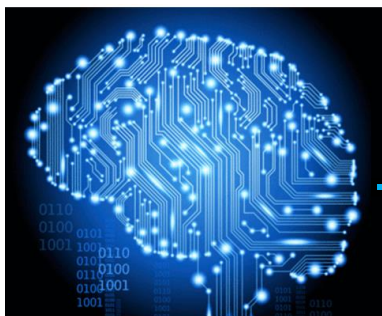Communicate actionable information back to user.

## Program-based disambiguation

- Enable effective navigation between top-ranked programs.
- Highlight ambiguity based on *distinguishing inputs*.

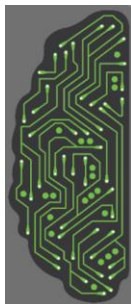## Heuristics that can be machine learned

- Highlight ambiguity based on clustering of inputs/outputs.
- When to stop highlighting ambiguity?

# ML in PBE



PBE Component  →  Logical strategies + Creative heuristics  →  Features + Model

Written by developers

Can be learned and maintained by ML-backed runtime

Advantages
- Better models
- Less time to author
- Online adaptation, personalization

# Mode-less Synthesis

Non-intrusively watch, learn, and make suggestions

Advantages: Usability, Avoids Discoverability

Applications: Document Editing, Code Refactoring, Robotic Process Automation

Key Idea: Identify related examples within noisy action traces

# **Predictive Synthesis**

Synthesis of intended programs from just the input.

Predictive Synthesis **:** PBE **::** Unsupervised **:** Supervised ML

Applications: Tabular data extraction, Join, Sort, Split

Key Idea: Structure inference over inputs

# **Synthesis of Readable Code**

Synthesis in target language of choice.
- Python, R, Scala, PySpark

Advantages:
- Transparency
- Education
- Integration with existing workflows in IDEs, Notebooks

Challenges: Quantify readability, Quantitative PBE

Key Idea: Observationally-equivalent (but non-semantic preserving) transformation of an intended program

# Program Synthesis meets Notebooks

A match made in heaven!

PS can synthesize small code fragments. Sufficient for notebook cell-based programming.

PS can synthesize code in different languages.
A good solution for polyglot challenge in notebooks.

PS needs interactivity. Notebooks provide that.

# Other Topics in Program Synthesis

- Search methodology: Code repositories [Murali et.al., ICLR 2018]

- Language: Neural program induction
  - [Graves et al., 2014; Reed & De Freitas, 2016; Zaremba et al., 2016]

- Intent specification:
  - Natural language [Huang et.al., NAACL-HLT 2018; Gulwani, Marron SIGMOD 2014, Shin et al. NeurIPS 2019]
  - Conversational pair programming

- Applications:
  - Super-optimization for model training/inference
  - Personalized Learning [Gulwani; CACM 2014]

# Conclusion

*Program Synthesis:* key to next-generational programming
- Future: Multi-modal programming with Examples and NL
- 100x more programmers
- 10-100x productivity increase in several domains.

Next-generational AI techniques under the hood
- Logical Reasoning + Machine Learning

Questions/Feedback: Contact me at **sumitg@microsoft.com**