

Synthesizing Number Transformations from Input-Output Examples

Rishabh Singh^{1*} and Sumit Gulwani²

¹ MIT CSAIL, Cambridge, MA, USA

² Microsoft Research, Redmond, WA, USA

Abstract. Numbers are one of the most widely used data type in programming languages. Number transformations like formatting and rounding present a challenge even for experienced programmers as they find it difficult to remember different number format strings supported by different programming languages. These transformations present an even bigger challenge for end-users of spreadsheet systems like Microsoft Excel where providing such custom format strings is beyond their expertise. In our extensive case study of help forums of many programming languages and Excel, we found that both programmers and end-users struggle with these number transformations, but are able to easily express their intent using input-output examples.

In this paper, we present a framework that can learn such number transformations from very few input-output examples. We first describe an expressive number transformation language that can model these transformations, and then present an inductive synthesis algorithm that can learn all expressions in this language that are consistent with a given set of examples. We also present a ranking scheme of these expressions that enables efficient learning of the desired transformation from very few examples. By combining our inductive synthesis algorithm for number transformations with an inductive synthesis algorithm for syntactic string transformations, we are able to obtain an inductive synthesis algorithm for manipulating data types that have numbers as a constituent sub-type such as date, unit, and time. We have implemented our algorithms as an Excel add-in and have evaluated it successfully over several benchmarks obtained from the help forums and the Excel product team.

1 Introduction

Numbers represent one of the most widely used data type in programming languages. Number transformations like formatting and rounding present a challenge even for experienced programmers. First, the custom number format strings for formatting numbers are complex and take some time to get accustomed to, and second, different programming languages support different format strings, which makes it difficult for programmers to remember each variant.

* Work done during an internship at Microsoft Research.

Number transformations present an even bigger challenge for end-users: the large class of users who do not have a programming background but want to create small, *often one-off*, applications to support business functions [4]. Spreadsheet systems like Microsoft Excel support a finite set of commonly used number formats and also let users write their own custom formats using a number formatting language similar to that of .NET. This hard-coded set of number formats is often insufficient for the user’s needs and providing custom number formats is typically beyond their expertise. This leads them to solicit help on various online help forums, where experts typically respond with the desired formulas (or scripts) after few rounds of interaction, which spans over a few days.

In an extensive case study of help forums of many programming languages and Excel, we found that even though both programmers and end-users struggled while performing these transformations, they were able to easily express their intent using input-output examples. In fact, in some cases the initial English description of the task provided by the users on forums was inaccurate and only after they provided a few input-output examples, the forum experts could provide the desired code snippet.

In this paper, we present a framework to learn number formatting and rounding transformations from a given set of input-output examples. We first describe a domain-specific language for performing number transformations and an inductive synthesis algorithm to learn the set of all expressions that are consistent with the user-provided examples. The key idea in the algorithm is to use the interval abstract domain [2] to represent a large collection of consistent format expressions symbolically, which also allows for efficient intersection, enumeration, and execution of these expressions. We also present a ranking mechanism to rank these expressions that enables efficient learning of the desired transformation from very few examples.

We then combine the number transformation language with a syntactic string transformation language [6] and present an inductive synthesis algorithm for the combined language. The combined language lets us model transformations on strings that represent data types consisting of number as a constituent subtype such as date, unit, time, and currency. The key idea in the algorithm is to succinctly represent an exponential number of consistent expressions in the combined language using a **Dag** data structure, which is similar to the BDD [1] representation of Boolean formulas. The **Dag** data structure consists of program expressions on the edges (as opposed to Boolean values on BDD edges). Similar to the BDDs, our data structure does not create a quadratic blowup after intersection in practice.

We have implemented our algorithms both as a stand-alone binary and as an Excel add-in. We have evaluated it successfully on over 50 representative benchmark problems obtained from help forums and the Excel product team.

This paper makes the following key contributions:

- We develop an expressive number transformation language for performing number formatting and rounding transformations, and an inductive synthesis algorithm for learning expressions in it.

- We combine the number transformation language with a syntactic string transformation language to manipulate richer data types.
- We describe an experimental prototype of our system with an attractive user interface that is ready to be deployed. We present the evaluation of our system over a large number of benchmark examples.

2 Motivating Examples

We motivate our framework with the help of a few examples taken from Excel help forums.

Example 1 (Date Manipulation). An Excel user stated that, as an unavoidable outcome of data extraction from a software package, she ended up with a series of dates in the input column v_1 as shown in the table. She wanted to convert them into a consistent date format as shown in the output column such that both month and day in the date are of two digits.

Input v_1	Output
1112011	01/11/2011
12012011	12/01/2011
1252010	01/25/2010
11152011	11/15/2011

It turns out that no Excel date format string matches the string in input column v_1 . The user struggled to format the date as desired and posted the problem on a help forum. After a few rounds of interactions (in which the user provided additional examples), the user managed to obtain the following formula for performing the transformation:

```
TEXT(IF(LEN(A1)=8,DATE(RIGHT(A1,4),MID(A1,3,2),LEFT(A1,2)),
DATE(RIGHT(A1,4),MID(A1,2,2),LEFT(A1,1))), "mm/dd/yyyy")
```

In our tool, the user has to provide only the first two input-output examples from which the tool learns the desired transformation, and executes the synthesized transformation on the remaining strings in the input column to produce the corresponding outputs (shown in bold font for emphasis).

We now briefly describe some of the challenges involved in learning this transformation. We first require a way to extract different substrings of the input date for extracting the day, month, and year parts of the date, which can be performed using the syntactic string transformation language [6]. We then require a number transformation language that can map 1 to 01, i.e. format a number to two digits. Consider the first input-output example 1112011 \rightarrow 01/11/2011. The first two characters in the output string can be obtained by extracting 1 from the input string from any of the five locations where it occurs, and formatting it to 01 using a number format expression. Alternatively, the first 0 in the output string can also be a constant string or can be obtained from the 3rd last character in the input. All these different choices for each substring of the output string leads to an exponential number of choices for the complete transformation. We use an efficient data structure for succinctly representing such exponential number of consistent expressions in polynomial space.

Example 2 (Duration Manipulation). An Excel user wanted to convert the “raw data” in the input column to the lower range of the corresponding “30-min interval” as shown in the output column. An expert responded by providing the following macro, which is quite unreadable and error-prone.

Input v_1	Output
0d 5h 26m	5:00
0d 4h 57m	4:30
0d 4h 27m	4:00
0d 3h 57m	3:30

```
FLOOR(TIME(MID(C1,FIND(" ",C1)+1,FIND("h",C1)- FIND(" ",C1)-1)+0,
MID(C1,FIND("h",C1)+2,FIND("m",C1)-FIND("h",C1)-2)+0,0)*24,0.5)/24
```

Our tool learns the desired transformation using only the first two examples. In this case, we first need to be able to extract the hour and minute components of the duration in input column v_1 , and then perform a rounding operation on the minute part of the input to round it to the lower 30-min interval.

3 Overview of the Synthesis Approach

In this section, we define the formalism that we use in the paper for developing inductive synthesizers [8].

Domain-specific language: We develop a domain-specific language L that is expressive enough to capture the desired tasks and, at the same time, is concise for enabling efficient learning from examples.

Data structure for representing a set of expressions: The number of expressions that are consistent with a given input-output example can potentially be very large. We, therefore, develop an efficient data structure D that can succinctly represent a large number of expressions in L .

Synthesis algorithm: The synthesis algorithm `Synthesize` consists of the following two procedures:

- **GenerateStr:** The `GenerateStr` procedure learns the set of all expressions in the language L (represented using the data structure D) that are consistent with a given input-output example (σ_i, s_i) . An input state σ holds values for m string variables v_1, \dots, v_m (denoting m input columns in a spreadsheet).
- **Intersect:** The `Intersect` procedure intersects two sets of expressions to compute the common set of expressions.

The synthesis algorithm `Synthesize` takes as input a set of input-output examples and generates a set of expressions in L that are consistent with them. It uses `GenerateStr` procedure to generate a set of expressions for each individual input-output example and then uses the `Intersect` procedure to intersect the corresponding sets to compute the common set of expressions.

```
Synthesize(( $\sigma_1, s_1$ ), ..., ( $\sigma_n, s_n$ ))
  P := GenerateStr( $\sigma_1, s_1$ );
  for i = 2 to n:
    P' := GenerateStr( $\sigma_i, s_i$ );
    P := Intersect(P, P');
  return P;
```

Ranking: Since there are typically a large number of consistent expressions for each input-output example, we rank them using the Occam’s razor principle that

states that smaller and simpler explanations are usually the correct ones. This enables users to provide only a few input-output examples for quick convergence to the desired transformation.

4 Number Transformations

In this section, we first describe the number transformation language L_n that can perform formatting and rounding transformations on numbers. We then describe an efficient data structure to succinctly represent a large number of expressions in L_n , and present an inductive synthesis algorithm to learn all expressions in the language that are consistent with a given set of input-output examples.

4.1 Number Transformation Language L_n

<div style="display: flex; flex-direction: column; gap: 5px;"> <div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <p>Expr. $e_n := \text{Dec}(u, \eta_1, f)$</p> <p style="padding-left: 20px;"> $\text{Exp}(u, \eta_1, f, \eta_2)$</p> <p style="padding-left: 20px;"> $\text{Ord}(u)$</p> <p style="padding-left: 20px;"> $\text{Word}(u)$</p> <p style="padding-left: 20px;"> u</p> </div> <div style="width: 50%;"> <p>Dec. Fmt. $f := (\odot, \eta) \mid \perp$</p> <p>Number $u := v_i$</p> <p style="padding-left: 20px;"> $\text{Round}(v_i, r)$</p> </div> </div> <div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <p>Round Fmt. $r := (z, \delta, m)$</p> <p style="padding-left: 20px;">Mode $m := \downarrow \mid \uparrow \mid \updownarrow$</p> <p>Num. Fmt. $\eta := (\alpha, \beta, \gamma)$</p> </div> <div style="width: 50%;"></div> </div> </div>	<div style="display: flex; flex-direction: column; gap: 5px;"> <p>$\llbracket \text{Dec}(u, \eta_1, f) \rrbracket \sigma = \llbracket (\text{Int}(\llbracket u \rrbracket^R), \eta_1) \rrbracket^R \sigma \star \llbracket f \rrbracket \sigma$</p> <p>$\llbracket \text{Exp}(u, \eta_1, f, \eta_2) \rrbracket \sigma = \llbracket (\text{Int}(\llbracket u \rrbracket^R), \eta_1) \rrbracket^R \sigma \star \llbracket f \rrbracket \sigma \star \llbracket (\text{E}(\llbracket u \rrbracket^R), \eta_2) \rrbracket^R \sigma$</p> <p>$\llbracket \text{Ord}(u) \rrbracket \sigma = \text{numToOrd}(\llbracket u \rrbracket \sigma)$</p> <p>$\llbracket \text{Word}(u) \rrbracket \sigma = \text{numToWord}(\llbracket u \rrbracket \sigma)$</p> <p>$\llbracket (\odot, \eta) \rrbracket \sigma = \llbracket \odot \rrbracket \sigma \star \llbracket (\text{Frac}(\llbracket u \rrbracket), \eta) \rrbracket \sigma$</p> <p>$\llbracket \perp \rrbracket \sigma = \epsilon$</p> <p>$\llbracket v_i \rrbracket \sigma = \sigma(v_i)$</p> <p>$\llbracket \text{Round}(v_i, r) \rrbracket \sigma = \text{RoundNumber}(\sigma(v_i), z, \delta, m)$ where $r = (z, \delta, m)$</p> <p>$\llbracket (d, \eta) \rrbracket \sigma = \text{FormatDigits}(d, \alpha, \beta, \gamma)$ where $\eta = (\alpha, \beta, \gamma)$</p> </div>
(a)	(b)

Fig. 1. The (a) syntax and (b) semantics of the number transformation language L_n . The variable v_i denotes an input number variable, $z, \delta, \alpha, \beta,$ and γ are integer constants, and \star denotes the concatenation operation.

The syntax of the number transformation language L_n is shown in Figure 1(a). The top-level expression e_n of the language denotes a number formatting expression of one of the following forms:

- $\text{Dec}(u, \eta_1, f)$: formats the number u in decimal form (e.g. 1.23), where η_1 denotes the number format for the integer part of u ($\text{Int}(u)$), and f represents the optional format consisting of the decimal separator and the number format for the fractional part ($\text{Frac}(u)$).
- $\text{Exp}(u, \eta_1, f, \eta_2)$: formats the number u in exponential form (e.g. 1.23E+2). It consists of an additional number format η_2 as compared to the decimal format expression, which denotes the number format of the exponent digits of u .

<pre> <u>RoundNumber</u>(n, z, δ, m) 1 $n' := \left\lfloor \frac{n - z}{\delta} \right\rfloor \times \delta + z;$ 2 if ($n = n'$) return n; 3 if ($m = \uparrow$) return $n' + \delta$; 4 if ($m = \downarrow$) return n'; 5 if ($m = \updownarrow \wedge (n - n') \times 2 < \delta$) return n'; 6 if ($m = \updownarrow \wedge (n - n') \times 2 \geq \delta$) return $n' + \delta$; (a) </pre>	<pre> <u>FormatDigits</u>(d, α, β, γ) 1 if ($\text{len}(d) \geq \beta$) 2 return significant(d, β); 3 else if ($\text{len}(d) \geq \alpha$) 4 $\{z := 0; s := 0\};$ 5 else $\{s := \text{Min}(\gamma, \alpha - \text{len}(d));$ 6 $z := \alpha - \text{len}(d) - s\};$ 7 return concat($d, 0_z, ' 's$); (b) </pre>
---	---

Fig. 2. The functions (a) `RoundNumber` for rounding numbers and (b) `FormatDigits` for formatting a digit string

- `Ord`(u): formats the number u in *ordinal* form, e.g. it formats the number 4 to its ordinal form 4th.
- `Word`(u): formats the number u in *word* form, e.g. it formats the number 4 to its word form **four**.

The number u can either be an input number variable v_i or a number obtained after performing a rounding transformation on an input number. A rounding transformation `Round`(v_i, z, δ, m) performs the rounding of number present in v_i based on its rounding format (z, δ, m), where z denotes the *zero* of the rounding interval, δ denotes the interval size of the rounding interval, and m denotes one of the rounding mode from the set of modes {upper(\uparrow), lower(\downarrow), nearest(\updownarrow)}.

We define a *digit string* d to be a sequence of digits with trailing whitespaces. A number format η of a digit string d is defined by a 3-tuple (α, β, γ) , where α denotes the minimum number of significant digits and trailing whitespaces of d in the output string, β denotes the maximum number of significant digits of d in the output string, and γ denotes the maximum number of trailing whitespaces in the output string. A number format, thus, maintains the invariant: $\gamma \leq \alpha \leq \beta$.

The semantics of language L_n is shown in Figure 1(b). A digit string d is formatted with a number format (α, β, γ) using the `FormatDigits` function shown in Figure 2(b). The `FormatDigits` function returns the first β digits of the digit string d (with appropriate rounding) if the length of d is greater than the maximum number of significant digits β to be printed. If the length of d is lesser than β but greater than the minimum number of significant digits α to be printed, it returns the digits itself. Finally, if the length of d is less than α , it appends the digit string with appropriate number of zeros (z) and whitespaces (s) as computed in Lines 5 and 6. The semantics of the rounding transformation is to perform the appropriate rounding of number denoted by v_i using the `RoundNumber` function shown in Figure 2(a). The function computes a number n' which lies on the number line defined by zero z with unit separation δ as shown in Figure 3. It returns the value n' or $(n' + \delta)$ based on the rounding mode m and the distance between n and n' as described in Figure 2(a).

The semantics of a decimal form formatting expression on a number u is to concatenate the reverse of the string obtained by formatting the reverse of

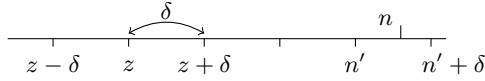


Fig. 3. The RoundNumber function rounding-off number n to n' or $n' + \delta$

integral part $\text{Int}(u)$ with the string obtained from the decimal format f . Since the `FormatDigits` function adds only trailing zeros and whitespaces to format a digit string, the formatting of the integer part of u is performed on its reverse digit string and the resulting formatted string is reversed again before performing the concatenation. The semantics of decimal format f is to concatenate the decimal separator \odot with the string obtained by formatting the fractional part $\text{Frac}(u)$. The semantics of exponential form formatting expression is similar to that of the decimal form formatting expression and the semantics of ordinal form and word form formatting expressions is to simply convert the number u into its corresponding ordinal form and word form respectively.

We now present some examples taken from various help forums that can be represented in the number transformation language L_n .

Example 3. A python programmer posted a query on the `StackOverflow` forum after struggling to print double values from an array of doubles (of different lengths) such that the decimal point for each value is aligned consistently across different columns. He posted an example of the desired formatting as shown on the right. He also wanted to print a single 0 after the decimal if the double value had no decimal part.

Input v_1	Output
3264.28	3264.28
53.5645	53.5645
235	235.0
5.23	5.23
345.213	345.213
3857.82	3857.82
536	536.0

The programmer started the post saying “*This should be easy*”. An expert replied that after a thorough investigation, he couldn’t find a way to perform this task without some post-processing. The expert provided the following python snippet that pads spaces to the left and zeros to the right of the decimal, and then removes trailing zeros:

```
ut0 = re.compile(r'(\d)0+$')
thelist = textwrap.dedent(
    '\n'.join(ut0.sub(r'\1', "%20f" % x) for x in a)).splitlines()
print '\n'.join(thelist)
```

This formatting transformation can be represented in L_n as $\text{Dec}(v_1, \eta_1, (".", \eta_2))$, where $\eta_1 \equiv (4, \infty, 4)$ and $\eta_2 \equiv (4, \infty, 3)$.

Example 4. This is an interesting post taken from a help forum where the user initially posted that she wanted to round numbers in an excel column to nearest 45 or 95, but the examples later showed that she actually wanted to round it to *upper* 45 or 95.

Input v_1	Output
11	45
32	45
46	95
1865	1895

Some of the solutions suggested by experts were:

```
=Min(Roundup(A1/45,0)*45,Roundup(A1/95,0)*95)
=CEILING(A1+5,50)-5
=A1-MOD(A1,100)+IF(MOD(A1,100)>45,95,45)
```

This rounding transformation can be expressed in our language as:

```
Dec(Round( $v_1$ , (45, 50,  $\uparrow$ )), (0,  $\infty$ , 0),  $\perp$ ).
```

4.2 Data structure for a set of expressions in L_n

Figure 4 describes the syntax and semantics of the data structure for succinctly representing a set of expressions from language L_n . The expressions \tilde{e}_n are now associated with a set of numbers \tilde{u} and a set of number formats $\tilde{\eta}$. We represent the set of numbers obtained after performing rounding transformation in two ways: $\text{Round}(v_i, \tilde{r})$ and $\text{Round}(v_i, n_p)$, which we describe in more detail in section 4.3. The set of number formats $\tilde{\eta}$ are represented using a 3-tuple (i_1, i_2, i_3) , where i_1 , i_2 and i_3 denote a set of values of α , β and γ respectively using an interval domain. This representation lets us represent $O(n^3)$ number of number format expressions in $O(1)$ space, where n denotes the length of each interval.

The semantics of evaluating the set of rounding transformations $\text{Round}(v_i, \tilde{r})$ is to return the set of results of performing rounding transformation on v_i for all rounding formats in the set \tilde{r} . The expression $\text{Round}(v_i, (n_1, n'_1))$ represents an infinite number of rounding transformations (as there exists an infinite number of rounding formats that conform to the rounding transformation $n_1 \rightarrow n'_1$). For evaluating this expression, we select one conforming rounding format with $z = 0$, $\delta = n'_1$ and an appropriate m as shown in the figure. The evaluation of a set of format strings $\tilde{\eta} = (i_1, i_2, i_3)$ on a digit string d returns a set of values, one for each possible combination of $\alpha \in i_1$, $\beta \in i_2$ and $\gamma \in i_3$. Similarly, we obtain a set of values from the evaluation of expression \tilde{e}_n .

4.3 Synthesis Algorithm

Procedure GenerateStr_n: The algorithm **GenDFmt** in Figure 5 takes as input two digit sequences d_1 and d_2 , and computes the set of all number formats $\tilde{\eta}$ that are consistent for formatting d_1 to d_2 . The algorithm first converts the digit sequence d_1 to its canonical form d'_1 by removing trailing zeros and whitespaces from d_1 . It then compares the lengths \mathbf{l}_1 of d'_1 and \mathbf{l}_2 of d_2 . If \mathbf{l}_1 is greater than \mathbf{l}_2 , then we can be sure that the digits got truncated and can therefore set the interval for i_2 (the maximum number of significant digits) to be $[\mathbf{l}_2, \mathbf{l}_2]$. The intervals for α and γ are set to $[0, \mathbf{l}_2]$ because of the number format invariant. On the other hand if \mathbf{l}_1 is smaller than \mathbf{l}_2 , we can be sure that the least number of significant digits need to be \mathbf{l}_2 , *i.e.* we can set the interval i_1 to be $[\mathbf{l}_2, \mathbf{l}_2]$. Also, we can set the interval i_2 to $[\mathbf{l}_2, \infty]$ because of the number format invariant. For interval i_3 , we either set it to $[\xi, \xi]$ (when $\mathbf{l}_2 - \xi \neq \mathbf{l}_1$) or $[\xi, \mathbf{l}_2]$ (when $\mathbf{l}_2 - \xi = \mathbf{l}_1$) where ξ denotes the number of trailing spaces in d_2 . In the former case, we can be sure about the exact number of trailing whitespaces to be printed.

$\begin{aligned} \tilde{e}_n &:= \text{Dec}(\tilde{u}, \tilde{\eta}_1, \tilde{f}) \\ & \text{Exp}(\tilde{u}, \tilde{\eta}_1, \tilde{f}, \tilde{\eta}_2) \\ & \text{Ord}(\tilde{u}) \\ & \text{Word}(\tilde{u}) \\ & \tilde{u} \\ \tilde{f} &:= (\odot, \tilde{\eta}) \mid \perp \\ \tilde{u} &:= v_i \\ & \text{Round}(v_i, \tilde{r}) \\ & \text{Round}(v_i, n_p) \\ \text{Pair } n_p &:= (n_1, n'_1) \\ \tilde{\eta} &:= (i_1, i_2, i_3) \\ \text{Interval } i &:= (l, h) \end{aligned}$ <p style="text-align: center;">(a)</p>	$\begin{aligned} \llbracket \text{Dec}(\tilde{u}, \tilde{\eta}_1, \tilde{f}) \rrbracket &= \{\text{Dec}(u, \eta_1, f) \mid u \in \tilde{u}, \eta_1 \in \tilde{\eta}_1, f \in \tilde{f}\} \\ \llbracket \text{Exp}(\tilde{u}, \tilde{\eta}_1, \tilde{f}, \tilde{\eta}_2) \rrbracket &= \{\text{Exp}(u, \eta_1, f, \eta_2) \mid u \in \tilde{u}, \eta_1 \in \tilde{\eta}_1, \\ &f \in \tilde{f}, \eta_2 \in \tilde{\eta}_2\} \\ \llbracket \text{Ord}(\tilde{u}) \rrbracket &= \{\text{Ord}(u) \mid u \in \tilde{u}\} \\ \llbracket \text{Word}(\tilde{u}) \rrbracket &= \{\text{Word}(u) \mid u \in \tilde{u}\} \\ \llbracket (\odot, \tilde{f}) \rrbracket &= \{(\odot, f) \mid f \in \tilde{f}\} \\ \llbracket \perp \rrbracket &= \epsilon \\ \llbracket v_i \rrbracket &= \{v_i\} \\ \llbracket \text{Round}(v_i, \tilde{r}) \rrbracket &= \{\text{Round}(v_i, (z, \delta, m)) \mid (z, \delta, m) \in \tilde{r}\} \\ \llbracket \text{Round}(v_i, n_p) \rrbracket &= \{\text{Round}(v_i, (0, n'_1, m)) \mid n_p \equiv (n_1, n'_1), \\ &\text{if } (n_1 \leq n'_1) m \equiv \uparrow \text{ else } m \equiv \downarrow\} \\ \llbracket (d, (i_1, i_2, i_3)) \rrbracket &= \{(d, \alpha, \beta, \gamma) \mid \alpha \in i_1, \beta \in i_2, \gamma \in i_3\} \end{aligned}$ <p style="text-align: center;">(b)</p>
---	---

Fig. 4. The (a) syntax and (b) semantics of a data structure for succinctly representing a set of expressions from language L_n .

The **GenerateStr_n** algorithm in Figure 5 learns the set of all expressions in L_n that are consistent with a given input-output example. The algorithm searches over all input variables v_i to find the inputs from which the output number n' can be obtained. It first converts the numbers $\sigma(v_i)$ and n' to their *canonical forms* n_C and n'_C respectively in Line 3. We define canonical form of a number to be its decimal value. If the two canonical forms n_C and n'_C are not equal, the algorithm tries to learn a rounding transformation such that n_C can be rounded to n'_C . We note that there is not enough information present in one input-output example pair to learn the exact rounding format as there exists an infinite family of such formats that are consistent. Therefore, we represent such rounding formats symbolically using the input-output example pair (n_C, n'_C) , which gets concretized by the **Intersect** method in Figure 6. The algorithm then normalizes the number $\sigma(u)$ with respect to n' using the **Normalize** method in Line 6 to obtain $n = (n_i, n_f, n_e)$ such that both n and n' are of the same form. For decimal and exponential forms, it learns a set of number formats $\tilde{\eta}$ for each of its constituent digit strings from the pairs (n_i^R, n_i^R) , (n_f, n'_f) , and (n_e^R, n_e^R) where n_i^R denotes the reverse of digit string n_i . As noted earlier, we need to learn the number format on the reversed digit strings for integer and exponential parts. For ordinal and word type numbers, it simply returns the expressions to compute ordinal and word forms of the corresponding input number respectively.

Procedure Intersect_n: The **Intersect_n** procedure for intersecting two sets of L_n expressions is described as a set of rules in Figure 6. The procedure computes the intersection of sets of expressions by recursively computing the intersection of their corresponding sets of sub-expressions. We describe below the four cases of

```

GenerateStrn(σ: inp state, n': out number)
1 Sn := ∅;
2 foreach input variable vi:
3   nc = Canonical(σ(vi)); n'c = Canonical(n');
4   if (nc ≠ n'c) u := Round(vi, (nc, n'c));
5   else u := vi;
6   (ni, nf, ne) := Normalize(σ(u), n');
7   match n' with
8     DecNum(n'i, n'f, ∘) →
9       η1 := GenDFmt(n'iR, n'iR);
10      if (∘ = ε) Sn := Sn ∪ Dec(u, η1, ⊥);
11      else { η2 := GenDFmt(n'f, n'f);
12             Sn := Sn ∪ Dec(u, η1, ∘, η2); }
13     ExpNum(n'i, n'f, n'e, ∘) →
14       η1 := GenDFmt(n'iR, n'iR);
15       η3 := GenDFmt(n'eR, n'eR);
16       if (∘ = ε) Sn := Sn ∪ Exp(u, η1, ⊥, η3);
17       else { η2 := GenDFmt(n'f, n'f);
18             Sn := Sn ∪ Exp(u, η1, ∘, η2, η3); }
19     OrdNum(n'i) →
20       Sn := Sn ∪ Ord(u);
21     WordNum(n'i) →
22       Sn := Sn ∪ Word(u);
23 return Sn;

GenDFmt(d1: inp digits, d2: out digits)
1 d'1 := RemoveTrailingZerosSpaces(d1);
2 l1 := len(d'1); l2 := len(d2);
3 ξ := numTrailingSpaces(d2);
4 if (l1 > l2)
5   (i1, i2, i3) := ([0, l2], [l2, l2], [0, l2]);
6 else if (l1 < l2) {
7   i1 := [l2, l2]; i2 := [l2, ∞];
8   if (l2 - ξ = l1) i3 := [ξ, l2];
9   else i3 := [ξ, ξ]; }
10 else (i1, i2, i3) := ([0, l2], [l2, ∞], [0, l2]);
11 return η(i1, i2, i3);

Normalize(n: inp number, n': out number)
n1 = n = (ni, nf, ne);
if (Type(n) = ExpNum ∧ Type(n') ≠ ExpNum)
  n1 := n × 10ne;
if (Type(n) ≠ ExpNum ∧ Type(n') = ExpNum)
  { n' = (n'i, n'f, n'e); n1 := n/10n'e; }
return n1;

```

Fig. 5. The GenerateStr_n procedure for generating the set of all expressions in language L_N that are consistent with the given set of input-output examples

intersecting rounding transformation expressions. The first case is of intersecting a finite rounding format set \tilde{r} with another finite set \tilde{r}' . The other two cases intersect a finite set \tilde{r} with an input-output pair n_p , which is performed by selecting a subset of the finite set of rounding formats that are consistent with the pair n_p . The final case of intersecting two input-output pairs to obtain a finite set of rounding formats is performed using the IntersectPair algorithm shown in Figure 7.

Consider the example of rounding numbers to nearest 45 or 95 for which we have the following two examples: $32 \rightarrow 45$ and $81 \rightarrow 95$. Our goal is to learn the rounding format (z, δ, m) that can perform the desired rounding transformation. We represent the infinite family of formats that satisfy the rounding constraint for each example as individual pairs $(32, 45)$ and $(81, 95)$ respectively. When we intersect these pairs, we can assign z to be 45 without loss of generality. We then compute all divisors $\tilde{\delta}$ of $95 - 45 = 50$. With the constraint that $\delta \geq (\text{Max}(45 - 32, 95 - 81) = 14)$, we finally arrive at the set $\tilde{\delta} = \{25, 50\}$. The rounding modes m are appropriately learned as shown in Figure 7. For decimal numbers, we compute the divisors by first scaling them appropriately and then re-scaling them back for learning the rounding formats. In our data structure, we do not store all divisors

```

IntersectPair((n1, n'1), (n2, n'2))
z := n'1;
δ̃ := Divisors(∥n'2 - n'1∥);
S := ∅;
foreach δ ∈ δ̃:
  if (δ ≥ Max(∥n1 - n'1∥, ∥n2 - n'2∥))
    if (2 × Max(∥n1 - n'1∥, ∥n2 - n'2∥) ≤ δ)
      S := S ∪ (z, δ, ↓);
    if (n1 > n'1 ∧ n2 > n'2)
      S := S ∪ (z, δ, ↓);
    if (n1 < n'1 ∧ n2 < n'2)
      S := S ∪ (z, δ, ↑);
return S;

```

Fig. 7. Intersection of Round expressions

$$\begin{aligned}
\text{Intersect}_n(\text{Dec}(\tilde{u}, \tilde{\eta}_1, \tilde{f}), \text{Dec}(\tilde{u}', \tilde{\eta}'_1, \tilde{f}')) &= \text{Dec}(\text{Intersect}_n(\tilde{u}, \tilde{u}'), \text{Intersect}_n(\tilde{\eta}_1, \tilde{\eta}'_1), \\
&\quad \text{Intersect}_n(\tilde{f}, \tilde{f}')) \\
\text{Intersect}_n(\text{Exp}(\tilde{u}, \tilde{\eta}_1, \tilde{f}, \tilde{\eta}_2), \text{Exp}(\tilde{u}', \tilde{\eta}'_1, \tilde{f}', \tilde{\eta}'_2)) &= \text{Exp}(\text{Intersect}_n(\tilde{u}, \tilde{u}'), \text{Intersect}_n(\tilde{\eta}_1, \tilde{\eta}'_1), \\
&\quad \text{Intersect}_n(\tilde{f}, \tilde{f}'), \text{Intersect}_n(\tilde{\eta}_2, \tilde{\eta}'_2)) \\
\text{Intersect}_n(\text{Ord}(\tilde{u}), \text{Ord}(\tilde{u}')) &= \text{Ord}(\text{Intersect}_n(\tilde{u}, \tilde{u}')) \\
\text{Intersect}_n(\text{Word}(\tilde{u}), \text{Word}(\tilde{u}')) &= \text{Word}(\text{Intersect}_n(\tilde{u}, \tilde{u}')) \\
\text{Intersect}_n(v_i, v_i) &= v_i \\
\text{Intersect}_n((\odot, \tilde{\eta}), (\odot', \tilde{\eta}')) &= (\text{Intersect}_n(\odot, \odot'), \text{Intersect}_n(\tilde{\eta}, \tilde{\eta}')) \\
\text{Intersect}_n(\text{Round}(v_i, \tilde{r}), \text{Round}(v_i, \tilde{r}')) &= \text{Round}(v_i, \text{Intersect}_n(\tilde{r}, \tilde{r}')) \\
\text{Intersect}_n(\text{Round}(v_i, \tilde{r}), \text{Round}(v_i, n_p)) &= \text{Round}(v_i, \text{Intersect}_n(\tilde{r}, n_p)) \\
\text{Intersect}_n(\text{Round}(v_i, n_p), \text{Round}(v_i, \tilde{r})) &= \text{Round}(v_i, \text{Intersect}_n(n_p, \tilde{r})) \\
\text{Intersect}_n(\text{Round}(v_i, n_p), \text{Round}(v_i, n'_p)) &= \text{Round}(v_i, \text{IntersectPair}(n_p, n'_p)) \\
\text{Intersect}_n((i_1, i_2, i_3), (i'_1, i'_2, i'_3)) &= (\text{Intersect}_n(i_1, i'_1), \text{Intersect}_n(i_2, i'_2), \\
&\quad \text{Intersect}_n(i_3, i'_3)) \\
\text{Intersect}_n((l, h), (l', h')) &= (\text{Max}(l, l'), \text{Min}(h, h'))
\end{aligned}$$

Fig. 6. The Intersect_n function for intersecting sets of expressions from language L_n . The Intersect_n function returns ϕ in all other case not covered above.

explicitly as this set might become too large for big numbers. We observe that we only need to store the greatest and least divisors amongst them, and then we can intersect two such sets efficiently by computing the `gcd` of the two corresponding greatest divisors and the `max` of the two corresponding least divisors.

Ranking: We rank higher the lower value for α in the interval i_1 (to prefer lesser trailing zeros and whitespaces), the higher value of β in i_2 (to minimize un-necessary number truncation), the lower value of γ in i_3 (to prefer trailing zeros more than trailing whitespaces), and the greatest divisor in the set of divisors $\tilde{\delta}$ of the rounding format (to minimize the length of rounding intervals). We rank expressions consisting of rounding transformations lower than the ones that consist of only number formatting expressions.

Theorem 1 (Correctness of Learning Algorithm for L_n).

- (a) The procedure GenerateStr_n is sound and complete. The complexity of GenerateStr_n is $O(|s|)$, where $|s|$ denotes the length of the output string.
- (b) The procedure Intersect_n is sound and complete.

Proof. The proof of (a) follows from the invariant of function GenDFmt which is that it maintains the set of all possible number formats in the language L_n that can format the input digits to the output digits. The case in which the two canonical forms are not equal, the round-off function takes care of the truncation case from Theroem 3. The rounding transformation is also sound and complete

as we maintain all possible choices for δ by maintaining a set of all divisors. From Lemma 1, we have that the value of z can be soundly chosen to be any output value n' . The proof of (b) follows from the semantics of the **Intersect** method, as it never loses any possible interpretations when intersecting and maintains the set of all common expressions.

Example 5. Figure 8 shows a range of number formatting transformations and presents the format strings that are required to be provided in Excel, .NET, Python and C, as well as the format expressions that are synthesized by our algorithm. An N.A. entry denotes that the corresponding formatting task cannot be done in the corresponding language.

Formatting of Doubles				
Input String	Output String	Excel/C# Format String	Python/C Format String	Synthesized format $\text{Dec}(u, \eta_1, (".", \eta_2))$ or $\text{Exp}(u, \eta_1, (".", \eta_2), \eta_3)$
123.4567 123.4	123.46 123.40	#.00	.2f	$\eta_1 \equiv ([0, 3], [3, \infty], [0, 3])$ $\eta_2 \equiv ([2, 2], [2, 2], [0, 0])$
123.4567 123.4	123.46 123.4	###	N.A.	$\eta_1 \equiv ([0, 3], [3, \infty], [0, 3])$ $\eta_2 \equiv ([0, 1], [2, 2], [0, 1])$
123.4567 3.4	123.46 03.40	00.00	05.2f	$\eta_1 \equiv ([2, 2], [3, \infty], [0, 0])$ $\eta_2 \equiv ([2, 2], [2, 2], [0, 0])$
123.4567 3.4	123.46 03.4	00.##	N.A.	$\eta_1 \equiv ([2, 2], [3, \infty], [0, 0])$ $\eta_2 \equiv ([0, 1], [2, 2], [0, 1])$
9723.00 0.823	9.723E+03 8.23E-01	##### E 00	N.A.	$\eta_1 \equiv ([0, 1], [1, \infty], [0, 1])$ $\eta_2 \equiv ([0, 3], [3, \infty], [0, 3])$ $\eta_3 \equiv ([2, 2], [2, \infty], [0, 0])$
243 12	00243 00012	00000	05d	$\eta_1 \equiv ([5, 5], [5, \infty], [0, 0])$
1.2 18	1.2_ 18.---	#.??	N.A.	$\eta_1 \equiv ([0, 1], [2, \infty], [0, 1])$ $\eta_2 \equiv ([2, 2], [2, \infty], [2, 2])$
1.2 18	_.1.2__ _.18.---	???.???	N.A.	$\eta_1 \equiv ([3, 3], [3, \infty], [2, 3])$ $\eta_2 \equiv ([3, 3], [3, \infty], [3, 3])$
1.2 18	_.1.20_ _.18.00_	???.00?	N.A.	$\eta_1 \equiv ([3, 3], [3, \infty], [2, 3])$ $\eta_2 \equiv ([3, 3], [3, \infty], [1, 1])$

Fig. 8. We compare the custom number format strings required to perform formatting of doubles in Excel/C# and Python/C languages. An N.A. entry in a format string denotes that the corresponding formatting is not possible using format strings only. The last column presents the corresponding L_n expressions (_ denotes whitespaces).

5 Combining Number Transformations with Syntactic String Transformations

In this section, we present the combination of number transformation language L_n with the syntactic string transformation language L_s [6] to obtain the combined language L_c , which can model transformations on strings that contain numbers as substrings. We first present a brief background description of the syntactic string transformation language and then present the combined language L_c . We also present an inductive synthesis algorithm for L_c obtained by combining the inductive synthesis algorithms for L_n and L_s respectively.

Syntactic String Transformation Language L_s (Background) Gulwani [6] introduced an expression language for performing syntactic string transformations. We reproduce here a small subset of (the rules of) that language and call it L_s (with e_s being the top-level symbol) as shown in Figure 9. The formal se-

$$\begin{aligned}
 e_s &:= \text{Concatenate}(f_1, \dots, f_n) \mid f \\
 \text{Atomic expr } f &:= \text{ConstStr}(s) \mid v_i \mid \text{SubStr}(v_i, p_1, p_2) \\
 \text{Position } p &:= k \mid \text{pos}(r_1, r_2, c) \\
 \text{Integer expr } c &:= k \mid k_1 w + k_2 \\
 \text{Regular expr } r &:= \epsilon \mid T \mid \text{TokenSeq}(T_1, \dots, T_n)
 \end{aligned}$$

Fig. 9. The syntax of syntactic string transformation language L_s .

mantics of L_s can be found in [6]. For completeness, we briefly describe some key aspects of this language. The top-level expression e_s is either an atomic expression f or is obtained by concatenating atomic expressions f_1, \dots, f_n using the **Concatenate** constructor. Each atomic expression f can either be a constant string **ConstStr**(s), an input string variable v_i , or a substring of some input string v_i . The substring expression **SubStr**(v_i, p_1, p_2) is defined partly by two *position expressions* p_1 and p_2 , each of which implicitly refers to the (subject) string v_i and must evaluate to a position within the string v_i . (A string with ℓ characters has $\ell + 1$ positions, numbered from 0 to ℓ starting from left.) **SubStr**(v_i, p_1, p_2) is the substring of string v_i in between positions p_1 and p_2 . A position expression represented by a non-negative constant k denotes the k^{th} position in the string. For a negative constant k , it denotes the $(\ell + 1 + k)^{\text{th}}$ position in the string, where $\ell = \text{Length}(s)$. **pos**(r_1, r_2, c) is another position expression, where r_1 and r_2 are regular expressions and integer expression c evaluates to a non-zero integer. **pos**(r_1, r_2, c) evaluates to a position t in the subject string s such that r_1 matches some suffix of $s[0 : t]$, and r_2 matches some prefix of $s[t : \ell]$, where $\ell = \text{Length}(s)$. Furthermore, if c is positive (negative), then t is the $|c|^{\text{th}}$ such match starting from the left side (right side). We use the expression $s[t_1 : t_2]$ to denote the substring of s between positions t_1 and t_2 . We use

the notation $\text{SubStr2}(v_i, r, c)$ as an abbreviation to denote the c^{th} occurrence of regular expression r in v_i , i.e., $\text{SubStr}(v_i, \text{pos}(\epsilon, r, c), \text{pos}(r, \epsilon, c))$.

A regular expression r is either ϵ (which matches the empty string, and therefore can match at any position of any string), a token T , or a token sequence $\text{TokenSeq}(T_1, \dots, T_n)$. The tokens T range over a finite extensible set and typically correspond to character classes and special characters. For example, tokens `CapitalTok`, `NumTok`, and `WordTok` match a nonempty sequence of uppercase alphabetic characters, numeric digits, and alphanumeric characters respectively.

A `Dag` based data structure is used to succinctly represent a set of L_s expressions. The `Dag` structure consists of a node corresponding to each position in the output string s , and a map W maps an edge between node i and node j to the set of all L_c expressions that can compute the substring $s[i..j]$. This representation enables sharing of common subexpressions amongst the set of expressions and represents an exponential number of expressions using polynomial space.

Example 6. An Excel user wanted to modify the delimiter in dates present in a column from “/” to “-”, and gave the following input-output example “08/15/2010” \rightarrow “08-15-2010”. An expression in L_s that can perform this transformation is: $\text{Concatenate}(f_1, \text{ConstStr}(\text{“-”}), f_2, \text{ConstStr}(\text{“-”}), f_3)$, where $f_1 \equiv \text{SubStr2}(v_1, \text{NumTok}, 1)$, $f_2 \equiv \text{SubStr2}(v_1, \text{NumTok}, 2)$, and $f_3 \equiv \text{SubStr2}(v_1, \text{NumTok}, 3)$. This expression constructs the output string by concatenating the first, second, and third numbers of input string with constant strings “-”.

5.1 The Combination Language L_c

The grammar rules R_c for the combined language L_c are obtained by taking the union of the rules for the two languages R_n and R_s with the top-level rule e_s . The modified rules are shown in the figure on the right. The combined language consists of an additional expression rule g that corresponds to

$$\begin{aligned} f &:= \text{ConstStr}(s) \mid v_i \\ &\quad \mid \text{SubStr}(v_i, p_1, p_2) \mid e_n \\ u &:= g \mid \text{Round}(g, r) \\ g &:= v_i \mid \text{SubStr}(v_i, p_1, p_2) \end{aligned}$$

either some input column v_i or a substring of some input column. This expression g is then passed over to the number variable expression u for performing number transformations on it. This rule enables the combined language to perform number transformations on substrings of input strings. The top-level expression of the number language e_n is added to the atomic expr f of the string language. This enables number transformation expressions to be present on the `Dag` edges together with the syntactic string transformation expressions.

The transformation in Example 1 is represented in L_c as: $\text{Concatenate}(f_1, \text{ConstStr}(\text{“/”}), f_2, \text{ConstStr}(\text{“/”}), f_3)$, where $f_1 \equiv \text{Dec}(g_1, (2, \infty, 0), \perp)$, $g_1 \equiv \text{SubStr}(v_1, 1, -7)$, $f_2 \equiv \text{SubStr}(v_1, -7, -5)$, and $f_3 \equiv \text{SubStr}(v_1, -5, -1)$. The transformation in Example 2 is represented as: $\text{Concatenate}(f_1, \text{“:”}, f_2)$, where $f_1 \equiv \text{SubStr2}(v_1, \text{NumTok}, 2)$, $f_2 \equiv \text{Dec}(u_1, (2, \infty, 0), \perp)$, and $u_1 \equiv \text{Round}(\text{SubStr2}(v_1, \text{NumTok}, 3), (0, 30, \downarrow))$.

5.2 Data structure for representing a set of expressions in L_c

Let \tilde{R}_n and \tilde{R}_s denote the set of grammar rules for the data structures that represent a set of expressions in L_n and L_s respectively. We obtain the grammar rules \tilde{R}_c for succinctly representing a set of expressions of L_c by taking the union of the two rule sets \tilde{R}_n and \tilde{R}_s with the updated rules as shown in Figure 10(a). The updated rules have expected semantics and can be defined as in Figure 4(b).

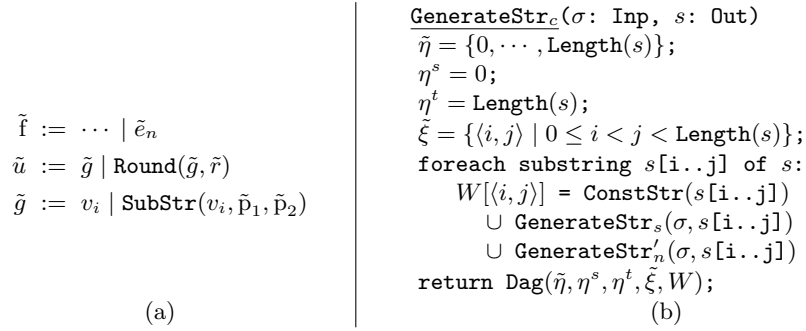


Fig. 10. (a) The data structure and (b) the GenerateStr_c procedure for L_c expressions.

5.3 Synthesis Algorithm

Procedure GenerateStr_c :

We first make the following two modifications in the GenerateStr_n procedure to obtain $\text{GenerateStr}'_n$ procedure. The first modification is that we now search over all substrings of input string variables v_i instead of just v_i in Line 2 in Figure 5. This lets us model transformations where number transformations are required to be performed on substrings of input strings. The second modification is that we replace each occurrence of v_i by $\text{GenerateStr}_s(\sigma, v_i)$ inside the loop body. This lets us learn the syntactic string program to extract the corresponding substring from the input string variables. The GenerateStr_c procedure for the combined language is shown in the Figure 10(b). The procedure first creates a Dag of $(\text{Length}(s) + 1)$ number of nodes with start node $\eta^s = 0$ and target node $\eta^t = \text{Length}(s)$. The procedure iterates over all substrings $\mathbf{s}[i..j]$ of the output string \mathbf{s} , and adds a constant string expression, a set of substring expressions (GenerateStr_s) and a set of number transformation expressions ($\text{GenerateStr}'_n$) that can generate the substring $\mathbf{s}[i..j]$ from the input state σ . These expressions are then added to a map $W[\langle i, j \rangle]$, where W maps each edge $\langle i, j \rangle$ of the dag to a set of expressions in L_c that can generate the corresponding substring $\mathbf{s}[i..j]$.

Procedure Intersect_c : The rules for Intersect_c procedure for intersecting sets of expressions in L_c are obtained by taking the union of intersection rules

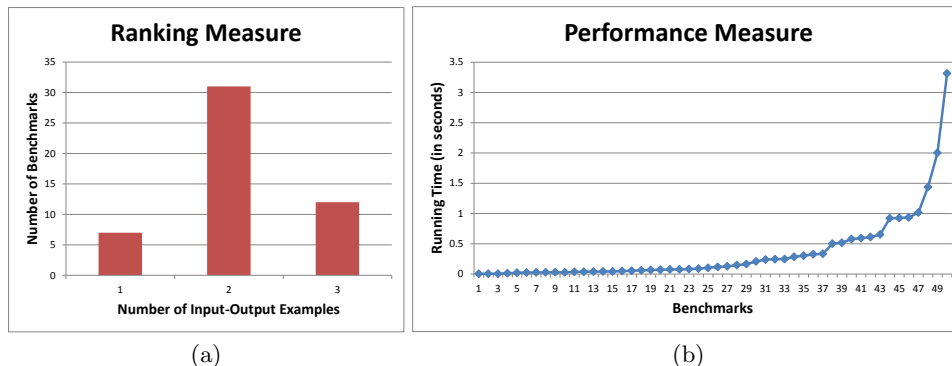


Fig. 11. (a) Number of examples required and (b) the running time of algorithm (in seconds) to learn the desired transformation

of Intersect_n and Intersect_s procedures together with corresponding intersection rules for the updated and new rules.

Ranking: The ranking scheme of the combined language L_c is obtained by combining the ranking schemes of languages L_n and L_s . In addition, we prefer substring expressions corresponding to longer input substrings that can be formatted or rounded to obtain the output number string.

Theorem 2 (Correctness of Learning Algorithm for combined language).

(a) The procedure GenerateStr_c is sound and complete with complexity $O(|s|^{3l^2})$, where $|s|$ denotes the length of the output string and l denotes the length of the longest input string.

(b) The procedure Intersect_c is sound and complete.

6 Experiments

We have implemented our algorithms in C# as an add-in to the Microsoft Excel spreadsheet system. The user provides input-output examples using an Excel table with a set of input and output columns. Our tool learns the expressions in L_c for each output column separately and executes the learned set of expressions on the remaining entries in the input columns to generate their corresponding outputs. We have evaluated our implementation on over 50 benchmarks obtained from various help forums, mailing lists, books and the Excel product team. More details about the benchmark problems can be found in [22].

The results of our evaluation are shown in Figure 11. The experiments were run on an Intel Core-i7 1.87 Ghz CPU with 4GB of RAM. We evaluate our algorithm on the following two dimensions:

Ranking: Figure 11(a) shows the number of input-output examples required by our tool to learn the desired transformation. All benchmarks required at most 3 examples, with majority (76%) taking only 2 examples to learn the desired transformation. We ran this experiment in an automated counter-example

guided manner such that given a set of input-output examples, we learned the transformations using a subset of the examples (training set). The tool iteratively added the failing test examples to the training set until the synthesized transformation conformed to all the remaining examples.

Performance: The running time of our tool on the benchmarks is shown in Figure 11(b). Our tool took at most 3.5 seconds each to learn the desired transformation for the benchmarks, with majority (94%) taking less than a second.

7 Related Work

The closest related work to ours is our previous work on synthesizing syntactic string transformations [6]. The algorithm presented in that work assumes strings to be a sequence of characters and can only perform concatenation of input substrings and constant strings to generate the desired output string. None of our benchmarks presented in this paper can be synthesized by that algorithm as it lacks reasoning about the semantics of numbers present in the input string.

There has been a lot of work in the HCI community for automating end-user tasks. Topes [20] system lets users create abstractions (called topes) for different data present in the spreadsheet. It involves defining constraints on the data to generate a context free grammar using a GUI and then this grammar is used to validate and reformat the data. There are several *programming by demonstration* [3] (PBD) systems that have been developed for data validation, cleaning and formatting, which requires the user to specify a complete demonstration or trace visualization on a representative data instead of code. Some of such systems include Simultaneous Editing [18] for string manipulation, SMARTedit [17] for text manipulation and Wrangler [15] for table transformations. In contrast to these systems, our system is based on programming by example (PBE) – it requires the user to provide only the input and output examples without providing the intermediate configurations which renders our system more usable [16], although at the expense of making the learning problem harder. Our expression languages also learns more sophisticated transformations involving conditionals. The by-example interface [7] has also been developed for synthesizing bit-vector algorithms [14], spreadsheet macros [8] (including semantic string manipulation [21] and table layout manipulation [12]), and even some intelligent tutoring scenarios (such as geometry constructions [10] and algebra problems [23]).

Programming by example can be seen as an instantiation of the general program synthesis problem, where the provided input-output examples constitutes the specification. Program synthesis has been used recently to synthesize many classes of non-trivial algorithms, e.g. graph algorithms [13], bit-streaming programs [26, 9], program inverses [27], interactive code snippets [11, 19], and data-structures [24, 25]. There are a range of techniques used in these systems including exhaustive search, constraint-based reasoning, probabilistic inference, type-based search, theorem proving and version-space algebra. A recent survey [5] explains them in more details. Lau et al. used the version-space algebra

based technique for learning functions in a PBD setting [17], our system uses it for learning expressions in a PBE setting.

8 Conclusions

We have presented a number transformation language that can model number formatting and rounding transformations, and an inductive synthesis algorithm that can learn transformations in this language from a few input-output examples. We also showed how to combine our system for number transformations with the one for syntactic string transformations [6] to enable manipulation of data types that contain numbers as substrings (such as date and time). In addition to helping end-users who lack programming expertise, we believe that our system is also useful for programmers since it can provide a consistent number formatting interface across all programming languages.

References

1. R. E. Bryant. Graph-based algorithms for boolean function manipulation. *IEEE Trans. Computers*, 35(8):677–691, 1986.
2. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL*, 1977.
3. A. Cypher, editor. *Watch What I Do – Programming by Demonstration*. MIT Press, 1993.
4. M. Gualtieri. Deputize end-user developers to deliver business agility and reduce costs. In *Forrester Report for Application Development and Program Management Professionals*, April 2009.
5. S. Gulwani. Dimensions in program synthesis. In *PPDP*, 2010.
6. S. Gulwani. Automating string processing in spreadsheets using input-output examples. In *POPL*, 2011.
7. S. Gulwani. Synthesis from examples. *WAMBSE (Workshop on Advances in Model-Based Software Engineering) Special Issue, Infosys Labs Briefings*, 10(2), 2012.
8. S. Gulwani, W. R. Harris, and R. Singh. Spreadsheet data manipulation using examples. In *Communications of the ACM*, 2012. To Appear.
9. S. Gulwani, S. Jha, A. Tiwari, and R. Venkatesan. Synthesis of loop-free programs. In *PLDI*, 2011.
10. S. Gulwani, V. A. Korthikanti, and A. Tiwari. Synthesizing geometry constructions. In *PLDI*, pages 50–61, 2011.
11. T. Gvero, V. Kuncak, and R. Piskac. Interactive synthesis of code snippets. In *CAV*, pages 418–423, 2011.
12. W. R. Harris and S. Gulwani. Spreadsheet table transformations from examples. In *PLDI*, pages 317–328, 2011.
13. S. Itzhaky, S. Gulwani, N. Immerman, and M. Sagiv. A simple inductive synthesis methodology and its applications. In *OOPSLA*, 2010.
14. S. Jha, S. Gulwani, S. Seshia, and A. Tiwari. Oracle-guided component-based program synthesis. In *ICSE*, 2010.
15. S. Kandel, A. Paepcke, J. Hellerstein, and J. Heer. Wrangler: Interactive visual specification of data transformation scripts. In *CHI*, 2011.

16. T. Lau. Why PBD systems fail: Lessons learned for usable AI. In *CHI Workshop on Usable AI*, 2008.
17. T. Lau, S. Wolfman, P. Domingos, and D. Weld. Programming by demonstration using version space algebra. *Machine Learning*, 53(1-2):111–156, 2003.
18. R. C. Miller and B. A. Myers. Interactive simultaneous editing of multiple text regions. In *USENIX Annual Technical Conference*, 2001.
19. D. Perelman, S. Gulwani, T. Ball, and D. Grossman. Type-directed completion of partial expressions. In *PLDI*, 2012.
20. C. Scaffidi, B. A. Myers, and M. Shaw. Topes: reusable abstractions for validating data. In *ICSE*, pages 1–10, 2008.
21. R. Singh and S. Gulwani. Learning semantic string transformations from examples. *PVLDB*, 5, 2012. (To appear).
22. R. Singh and S. Gulwani. Synthesizing number transformations from input-output examples. Technical Report MSR-TR-2012-42, Apr 2012.
23. R. Singh, S. Gulwani, and S. Rajamani. Automatically generating algebra problems. In *AAAI*, 2012. (To appear).
24. R. Singh and A. Solar-Lezama. Synthesizing data structure manipulations from storyboards. In *SIGSOFT FSE*, pages 289–299, 2011.
25. A. Solar-Lezama, C. G. Jones, and R. Bodík. Sketching concurrent data structures. In *PLDI*, pages 136–148, 2008.
26. A. Solar-Lezama, R. M. Rabbah, R. Bodík, and K. Ebcioğlu. Programming by sketching for bit-streaming programs. In *PLDI*, pages 281–294, 2005.
27. S. Srivastava, S. Gulwani, S. Chaudhuri, and J. S. Foster. Path-based inductive synthesis for program inversion. In *PLDI*, 2011.

A Appendix

A.1 Additional Examples

Example 7. A user wanted to round off values in a column to the nearest 10

Input v_1	Output
112	110
117	120
11112	11110
11119	11120

value and gave the following examples on a help forum.

Experts on the help-forum suggested the following solutions:
`=IF(RIGHT(A1,1)<5,A1-RIGHT(A1,1),A1-RIGHT(A1,1)+10)`
`=ROUND(D3*0,1;0)*10`
`=round(number,-1)`

Our framework synthesizes the following expression:
 $L_{\mathbb{N}}$ expression: $\text{Dec}(\text{Round}(v_1, 110, 10, \downarrow), ([0, 3], [6, \infty], [0, 3]), \perp)$

Example 8. The goal of this problem, taken from a help forum, is to round a decimal number to the nearest 100th place. Again, examples give a better

Input v_1	Output
19.57	19.55
19.58	19.60

understanding of the desired intention.

Suggested Solutions:

=ROUND(A1/5,2)*5
 =ROUND(A1/0.05,0)*0.05

$L_{\mathbb{N}}$ expression: $\text{Dec}(\text{Round}(v_1, 19.55, 0.05, \downarrow), ([0, 2], [2, \infty], [0, 2]),$
 “.”, $([2, 2][2, \infty][0, 0])$)

Example 9. The goal of this problem, taken from a help forum, is to round the

Input v_1	Output
1.63	1.99
1.24	1.49

numbers to the upper 0.49 or 0.99 number.

$L_{\mathbb{N}}$ expression: $\text{Dec}(\text{Round}(v_1, 1.49, 0.50, \uparrow), ([0, 1], [1, \infty], [0, 1]),$
 “.”, $([0, 2][2, \infty][0, 2])$)

Example 10. The goal of this problem, taken from a help forum, is to round the

Input v_1	Output
1.36	1.36
1.43	1.44
1.37	1.38

numbers to the upper even numbers.

$L_{\mathbb{N}}$ expression: $\text{Dec}(\text{Round}(v_1, 1.36, 0.02, \uparrow), ([0, 1], [1, \infty], [0, 1]),$
 “.”, $([0, 2][2, \infty][0, 2])$)

Example 11. An end-user was trying to generate an Excel report in Jasper but was unsuccessful in formatting numbers appropriately. He gave the examples shown in the table for specifying his intention.

Input v_1	Output
234.1	234.100
0	0.000

The task can be represented in our language as: $\text{Dec}(v_1, \eta_1, \text{“.”}, \eta_2)$ where $\eta_1 \equiv$
 $([0, 1], [3, \infty], [0, 1]), \eta_2 \equiv ([3, 3], [3, \infty], [0, 0])$

Example 12. An Excel user wanted to round a number lesser than 1.25 to 1.0, a number between 1.25 to 1.75 to 1.5 and numbers greater than 1.75 to 2.0 and similarly for other numbers in the columns.

Input v_1	Output
1.14	1.0
1.45	1.5
1.82	2.0
3.65	3.5

An expert on the help forum suggested the following formula:

=IF(A1="", "", IF(A1<1.25, 1, IF(AND (A1>=1.25, A1<=1.75), 1.5, IF(A1>1.75, 2))))

The desired transformation can be expressed in our language as:

$\text{Dec}(\text{Round}(v_1, 1.00, 0.50, \uparrow), ([0, 1], [1, \infty], [0, 1]), \text{“.”}, ([1, 1][1, \infty][0, 0]))$

A.2 Correctness and Soundness Theorems

Theorem 3. *For every truncation transformation, there exists an equivalent round-off transformation that generates the same output digits.*

Proof. A truncate transformation $\text{FormatDigits}(d, \alpha, \beta, \gamma)$ where $\beta > \text{Length}(d)$ that produces digits d' can be represented equivalently as round-off transformation $\text{Round}(v_i, (z, \delta, m))$ where $z = 0$, $\delta = 10^{-\beta}$ with appropriate mode m .

Lemma 1. *In a rounding transformation, any output value n' can be chosen soundly to be the zero z of the rounding interval.*

Theorem 4 (Correctness of Learning Algorithm.).

- (a) *The procedure GenerateStr_n is sound and complete.*
- (b) *The procedure Intersect_n is sound and complete.*

Proof. The proof of (a) follows from the invariant of function GenDFmt which is that it maintains the set of all possible number formats in the language L_n that can format the input digits to the output digits. The case in which the two canonical forms are not equal, the round-off function takes care of the truncation case from Theorem 3. The rounding transformation is also sound and complete as we maintain all possible choices for δ by maintaining a set of all divisors. From Lemma 1, we have that the value of z can be soundly chosen to be any output value n' . The proof of (b) follows from the semantics of the Intersect method, as it never loses any possible interpretations when intersecting and maintains the set of all common expressions.

Theorem 5 (Correctness of Learning Algorithm for combined language).

- (a) *The procedure GenerateStr_c is sound and complete.*
- (b) *The procedure Intersect_c is sound and complete.*

Proof. The soundness of GenerateStr_c comes directly from the soundness of GenerateStr_n and GenerateStr_s . The GenerateStr_c procedure associates each edge of the Dag with all corresponding number and syntactic string expressions and hence is complete. The soundness and completeness of Intersect_c also comes from soundness and completeness of Intersect_n and Intersect_s respectively.

A.3 Detailed results

Benchmark	# I-O examples	Running Time (in s)
benchmark1	2	0.5920339
benchmark2	3	1.0150580
benchmark3	2	0.0350020
benchmark4	2	0.0040003
benchmark5	2	0.1270073
benchmark6	3	0.2060117
benchmark7	2	0.0580033
benchmark8	2	0.2450140
benchmark9	1	0.0750042
benchmark10	3	00.5170296
benchmark11	2	00.1630093
benchmark12	3	00.3260187
benchmark13	2	00.1460084
benchmark14	1	00.0760043
benchmark15	2	00.3330191
benchmark16	4	03.3161897
benchmark17	2	00.2850163
benchmark18	2	00.6100349
benchmark19	2	00.6520373
benchmark20	2	00.5780330
benchmark21	2	00.5010286
benchmark22	2	00.0410023
benchmark23	1	00.0370021
benchmark24	2	00.0250015
benchmark25	2	00.0270016
benchmark26	2	00.0890051
benchmark27	2	00.0230013
benchmark28	2	00.0530031
benchmark29	3	00.0660038
benchmark30	2	00.0270015
benchmark31	3	00.0120006
benchmark32	2	00.0040002
benchmark33	3	00.0210012
benchmark34	2	00.0030002
benchmark35	2	00.0400023
benchmark36	2	00.0280016
benchmark37	2	00.2470141
benchmark38	2	00.2380137
benchmark39	1	00.0470027
benchmark40	3	00.3030173
benchmark41	2	02.0021145
benchmark42	1	00.0810047
benchmark43	3	00.1170067
benchmark44	3	00.9260530
benchmark45	1	00.0690039
benchmark46	2	00.0410023
benchmark47	3	00.9210527
benchmark48	1	00.1020059
benchmark49	3	01.4370822
benchmark50	2	00.9350534

Table 1. Detailed results of the experiments

benchmark1

Input v1	Output
243	00243
12.5	00012.5
2345.23292	
10	
1202.3433	
23224.1	

at least 5 digits before decimal

benchmark2

Input v1	Output
243.1	243.100
12.5	12.500
2345.23292	2345.233
10	
1202.3433	
23224.1	

exactly 3 digits after decimal

benchmark3

Input v1	Output
112	110
117	120
11112	
11119	
548	
23224	

round to nearest 10

benchmark4

Input v1	Output
11	45
46	95
32	
1865	
105	
546	

round to upper 45/95

benchmark5

Input v1	Output
19.57	19.55
19.58	19.60
21.48	
43.32	
16.42	
102.12	

round to nearest 0.05

benchmark6

Input v1	Output
1.36	1.36
1.43	1.44
1.37	1.38
43.32	
16.42	
102.11	

round to upper 0.02

benchmark7

Input v1	Output
1.14	1.0
1.45	1.5
1.82	
43.32	
16.42	
102.11	

round to nearest 0.5

benchmark8

Input v1	Output
123.4567	123.46
3.4	3.40
123.4	
1.82	
43.3235	
102	

exactly 2 decimal places

benchmark9

Input v1	Output
123.4567	123.46
3.4	3.4
123.4	
1.82	
43.3235	
102	

atmost 2 decimal places

benchmark10

Input v1	Output
123.4567	123.46
3.4	03.40
123.4	123.40
1.82	
43.3235	
102	

at least 2 digits before decimal
and exactly 2 decimal places

benchmark11

Input v1	Output
123.4567	123.46
3.4	03.4
123.4	
1.82	
43.3235	
102	

atleast 2 digits before decimal
atmost 2 decimal places

benchmark12

Input v1	Output
1.2	1.2
18	18.
3.4	3.4
1.82	
43.3235	
102	

Two spaces after decimal if no decimal digits

benchmark13

Input v1	Output
1.2	1.2
18	18.
1.82	1.82
3.4	
43.3235	
102	

Three spaces before and three spaces after the decimal

benchmark14

Input v1	Output
1.2	1.20
18	18.00
1.82	1.82
3.4	
43.3235	
102	

Three spaces before decimal, two zeros after decimal
and one space after that

benchmark15

Input v1	Output
5.23	5.23
325.213	325.213
53.5645	53.5645
3246.28	
235	
102	

at least 3 places (with spaces) before decimal
atleast 4 decimal places (with spaces)

benchmark16

Input v1	Output
1112011	01/11/2011
12012011	12/01/2011
1252010	01/25/2011
11152011	11/15/2011
6112011	
12062011	
11082010	
6012010	
10192011	
8211998	
11162010	

converting to mm/dd/yyyy format

benchmark17

Input v1	Output
20040717	2004/7/17
19991108	1999/11/8
19991108	
20080615	
20010918	
19960822	
19960725	
20000309	
20081217	
20031116	
20020402	
20010719	
20070317	
20000217	
20060709	
20020109	
20060620	
20091230	
20080805	
20000228	
20041126	
20010831	
19980716	

converting dates to yyyy/(m | mm)/(d | dd)

benchmark20

Input v1	Output
2004/07/17	17-7-2004
1999/11/08	8-11-1998
1999/11/08	
2008/06/15	
2001/09/18	
1996/08/22	
1996/07/25	
2000/03/09	
2008/12/17	
2003/11/16	
2002/04/02	
2001/07/19	
2007/03/17	
2000/02/17	
2006/07/09	
2002/01/09	
2006/06/20	
2009/12/30	
2008/08/05	
2000/02/28	
2004/11/26	
2001/08/31	
1998/07/16	

converting dates to (d|dd) - (m|mm) - yyyy

benchmark18

Input v1	Output
2004/07/17	2004717
1999/11/08	1999118
1999/11/08	
2008/06/15	
2001/09/18	
1996/08/22	
1996/07/25	
2000/03/09	
2008/12/17	
2003/11/16	
2002/04/02	
2001/07/19	
2007/03/17	
2000/02/17	
2006/07/09	
2002/01/09	
2006/06/20	
2009/12/30	
2008/08/05	
2000/02/28	
2004/11/26	
2001/08/31	
1998/07/16	

converting to yyyy(m|mm)(d|dd) format

benchmark19

Input v1	Output
2004/07/17	2004717
1999/11/08	19991108
1999/11/08	
2008/06/15	
2001/09/18	
1996/08/22	
1996/07/25	
2000/03/09	
2008/12/17	
2003/11/16	
2002/04/02	
2001/07/19	
2007/03/17	
2000/02/17	
2006/07/09	
2002/01/09	
2006/06/20	
2009/12/30	
2008/08/05	
2000/02/28	
2004/11/26	
2001/08/31	
1998/07/16	

converting to yyyy(m|mm)dd format

benchmark21

Input v1	Output
2004/07/17	17-07-2004
1999/11/08	8-11-1998
1999/11/08	
2008/06/15	
2001/09/18	
1996/08/22	
1996/07/25	
2000/03/09	
2008/12/17	
2003/11/16	
2002/04/02	
2001/07/19	
2007/03/17	
2000/02/17	
2006/07/09	
2002/01/09	
2006/06/20	
2009/12/30	
2008/08/05	
2000/02/28	
2004/11/26	
2001/08/31	
1998/07/16	

converting to (d|dd)-mm-yyyy format

benchmark22

Input v1	Output
26/4	04-26
5/11	11-05
23/9	
8/12	
14/5	
21/9	
14/8	
16/11	
24/9	
18/12	
1/2	
4/8	
15/5	
21/12	
18/2	
16/4	
22/9	
11/11	
12/11	
19/2	
25/4	
26/1	
15/7	

converting to mm-dd format

benchmark23

Input v1	Output
0930	930
1520	1520
1648	
0830	
1015	
2010	
1002	
1425	
2345	
1247	
0957	
1036	
0452	
1308	
1123	

converting to (h | hh)mm format

benchmark24

Input v1	Output
3.48	3.5
3.89	4.0
2342.35	
10.76	
1284.42	
23224.98	
1024.21	
14.98	

rounding to upper 0.5

benchmark25

Input v1	Output
3.48	3.50
3.89	4.00
2342.35	
10.76	
1284.42	
23224.98	
1024.21	
14.98	

rounding to upper 0.5
exactly 2 decimal places

benchmark26

Input v1	Output
249.60	250.00
247.10	245.00
2342.35	
10.76	
1284.42	
23224.98	
1024.21	
14.98	

rounding to nearest 1.00

benchmark27

Input v1	Output
4.56	4.99
7.23	9.99
2.45	
3.14	
8.56	
7.24	
9.94	
8.76	

rounding to upper 4.99/9.99

benchmark28

Input v1	Output
71.8	72.0
71.2	71.0
82.4	
103.8	
96.7	
74.8	
99.9	
78.8	

rounding to nearest 1.0

benchmark29

Input v1	Output
92.00	91.99
81.00	80.99
91.00	90.99
3.00	
8.00	
75.00	
99.00	
78.00	

rounding to lower 0.99

benchmark30

Input v1	Output
23.87	24.95
25.00	29.95
91.00	
3.00	
8.00	
75.00	
99.00	
78.00	

rounding to upper 4.95/9.95

benchmark31

Input v1	Output
542	500
954	1000
234	200
1321	
8330	
6265	
9812	
458	

rounding to nearest 100

benchmark32

Input v1	Output
542	1000
1954	2000
234	
1321	
8330	
6265	
9812	
458	

rounding to upper 1000

benchmark33

Input v1	Output
90	80
131	120
234	240
65	
124	
605	
842	
964	

rounding to nearest 40

benchmark34

Input v1	Output
64	100
158	200
556	
6265	
234	
605	
842	
964	

rounding to upper 100

benchmark35

Input v1	Output
234	235
232	230
238	
65	
124	
605	
842	
964	

rounding to nearest 5

benchmark36

Input v1	Output
423531	400000
324223	300000
234432	
763827	
283872	
234282	
932828	
832727	

rounding to nearest 100000

benchmark37

Input v1	Output
22666622	22666622
2321	00002321
2381	
65645424	
1244	
5321	
12541253	
9828	

atleast 8 digits before decimal

benchmark38

Input v1	Output
1.1	1.10
1.2345678	1.2345678
1	
1.234	
1.23456	
1.235	
4	
1.32256	

at least 2 decimal places

benchmark39

Input v1	Output
029.3	29.3
030.4	
028.2	
031.0	
13.24	
200.0	
4.62	
1.32256	

no leading zeros

benchmark40

Input v1	Output
1.243	1.2 pts
372.32	372.3 pts
1.25238	
5	
200.0	
8.238	
4.18	
19.252	

at most 1 decimal place with pts added at the end

benchmark41

Input v1	Output
100.34	100.3400
0.000347	0.0003
1.25238	
5	
200.083913	
8.238	
4.18	
19.252	

exactly 4 decimal places

benchmark42

Input v1	Output
9:30 5/14	09 30 05 14
14:25 11/23	
6:25 4/8	
9:05 12/5	
21:32 6/21	
22:48 8/16	
8:55 9/21	
7:30 4/18	

converting to format hh mm dd mm

benchmark43

Input v1	Output
+91	+0091
+617	+0617
+3523	+3523
+1	
+626	
+125	
+43	
+9	

exactly 4 digits after +

benchmark44

Input v1	Input v2	Output
1.23	2.54	1.2 * 2.5
32.624	5.216	32.6 * 5.2
11.26	2	11.3 * 2.0
5.21	2.15	
2.62	5.2	
92.252	7.15	
1.252	8.2	
8.25	2.165	

rounding to 1 decimal place
before multiplication

benchmark45

Input v1	Input v2	Input v3	Output
5	6	2001	05/06/2001
11	21	2001	
4	20	2002	
12	5	2002	
8	16	2001	
2	21	2002	
9	21	2002	
10	5	2001	

converting to
mm/dd/yyyy

benchmark46

Input v1	Input v2	Output
9	45	09:45
6	5	06:05
11	32	
21	30	
15	15	
8	26	
11	48	
23	16	

converting to
hh:mm format

benchmark47

Input v1	Input v2	Input v3	Output
521	632	2642	521-632-2642
91	24	2834	091-024-2531
14	52	728	014-052-0728
617	95	9217	
33	285	928423	
925	982	1892	
1	425	28520	
512	272	8167	

converting to
ddd-ddd-dddd format

benchmark48

Input v1	Input v2	Input v3	Output
1.23	2.54	12.421	1.2 + 2.5 + 12.4
32.624	5.216	5	
11.26	2	1.52	
5.21	2.15	9.2162	
1.58	8.26	8.28	
21.282	16.21	18.278	
1.43	1.2	1.82	
64.262	15.81	38.28	

converting to 1 decimal
place before adding

benchmark49

Input v1	Output
3:14:23	03h 14m 23s
11:5:25	11h 05m 25s
4:23:1	04h 23m 01s
21:7:6	
14:5:25	
12:8:29	
15:29:9	
11:23:45	

converting to hh~~h~~ mm~~m~~ sss format

benchmark50

Input v1	Output
0d 5h 26m	5:00
0d 4h 57m	4:30
0d 4h 27m	
0d 3h 57m	
0d 6h 23m	
0d 8h 42m	
0d 2h 56m	
0d 4h 35m	

rounding to lower 30 minutes interval