

2004

# Scaling a Dataflow Testing Methodology to the Multiparadigm World of Commercial Spreadsheets

Marc Randall Fisher II

*University of Nebraska at Lincoln*, [fisherii@google.com](mailto:fisherii@google.com)

Gregg Rothermel

*University of Nebraska - Lincoln*, [grothermel2@unl.edu](mailto:grothermel2@unl.edu)

Tyler Creelan

*Oregon State University*

Margaret Burnett

*Oregon State University*, [burnett@eecs.oregonstate.edu](mailto:burnett@eecs.oregonstate.edu)

Follow this and additional works at: <http://digitalcommons.unl.edu/csetechreports>



Part of the [Computer Sciences Commons](#)

---

Fisher II, Marc Randall; Rothermel, Gregg; Creelan, Tyler; and Burnett, Margaret, "Scaling a Dataflow Testing Methodology to the Multiparadigm World of Commercial Spreadsheets" (2004). *CSE Technical reports*. 36.

<http://digitalcommons.unl.edu/csetechreports/36>

This Article is brought to you for free and open access by the Computer Science and Engineering, Department of at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in CSE Technical reports by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

# Scaling a Dataflow Testing Methodology to the Multiparadigm World of Commercial Spreadsheets

Marc Fisher II, Gregg Rothermel  
University of Nebraska-Lincoln  
{mfisher,grother}@cse.unl.edu

Tyler Creelan, Margaret Burnett  
Oregon State University  
{creelan,burnett}@cs.orst.edu

## ABSTRACT

Spreadsheet languages are widely used by end users to perform a broad range of important tasks. Evidence shows, however, that spreadsheets often contain faults. Thus, in prior work we presented a dataflow testing methodology for use with spreadsheets, that provides feedback about the coverage of cells in spreadsheets via visual devices. Studies have shown that this methodology, which we call WYSIWYT (What You See Is What You Test), can be used cost-effectively by end-user programmers. To date, however, the methodology has been investigated across a limited set of spreadsheet language features. Commercial spreadsheet environments are multiparadigm languages, utilizing features often associated with dataflow, functional, imperative, and database query languages, and these features are not accommodated by prior approaches. In addition, most spreadsheets contain large numbers of replicated formulas that differ only in the cells they reference, and these severely limit the efficiency of dataflow testing approaches. We show how to handle these two aspects of commercial spreadsheet environments through a new dataflow adequacy criteria and automated detection of areas of replicated formulas. We report results of a controlled experiment investigating several factors important to the feasibility of our approach.

## 1. INTRODUCTION

Spreadsheets are used by a wide range of end users to perform a variety of important tasks, such as managing retirement funds, performing tax calculations, and forecasting revenues. Evidence shows, however, that spreadsheets often contain faults, and that these faults can have severe consequences. For example, an erroneous spreadsheet formula inflated the University of Toledo's projected annual revenue by 2.4 million dollars, requiring budget cuts [24]. Another formula error caused the stocks of Shurgard Inc. to be devalued after employees were overpaid by \$700,000 [23], and a cut-and-paste error in a bidding spreadsheet cost Transalta Corporation 24 million dollars through overbidding [10].

Researchers have been responding to these problems by creating approaches that address dependability issues for spreadsheets, including unit inference and checking systems [1, 2, 3], visual-

ization approaches [7, 8, 22], interval analysis techniques [4, 5], and approaches for automatic generation of spreadsheets from a model [11]. Commercial spreadsheet systems such as Microsoft Excel have also incorporated several tools for assisting with spreadsheet dependability, including dataflow arrows, anomaly detection heuristics, and data validation facilities.

In our own prior research, we have presented an integrated family of approaches meant to help end users improve the dependability of their spreadsheets. At the core of these approaches is a dataflow test adequacy criterion, and a dataflow testing methodology that helps spreadsheet users address potential problems in interactions between cell formulas – a prevalent source of spreadsheet errors [15, 16]. This WYSIWYT (What You See Is What You Test) methodology [19] uses visual devices to provide feedback about test coverage of the spreadsheet obtained relative to a dataflow adequacy criterion. We have augmented this methodology with techniques for automated test case generation [12], fault localization [21], and test reuse and replay mechanisms [13]. Our studies of the WYSIWYT methodology itself [17, 20] suggest that it can be effective, and can be applied by end users with no specific training in the underlying testing theories.

However, commercial spreadsheet environments are multiparadigm languages with features such as higher-order functions (functional paradigm), table query constructs (database query languages), user-defined functions<sup>1</sup> (implemented in an imperative sublanguage), meta-program constructs, and pointers, and these features are not accommodated by prior approaches. In addition, most spreadsheets have large areas of replicated formulas which require some form of aggregation and abstraction to allow our methodologies to scale reasonably (i.e., operate sufficiently efficiently). The only previous approach to testing methodologies for spreadsheet regions [6] has required a form of “declaring” regions, and thus does not provide *unassisted* discovery of the testing needs of the informal regions that exist in commercial spreadsheets.

In this paper, we address these two problems. To direct our efforts, we analyzed the features found in commercial spreadsheets that have not yet been supported by our WYSIWYT methodology.<sup>2</sup> From these we determined the features that were most integral to supporting the methodology in commercial spreadsheets, including the particular multiparadigmatic spreadsheet features that need to be accommodated, and the degree to which replicated formulas

<sup>1</sup>User-defined functions are a subset of macros. Macros are also used for automation and adding event-driven behaviors to spreadsheets. Since these other applications are not related to formula correctness, we don't consider them here.

<sup>2</sup>Our platform for this work has been Excel, the de-facto standard commercial spreadsheet environment. Our extended methodology also supports other environments mimicking Excel's feature set, such as OpenOffice/StarOffice and Gnumeric.

	A	B	C	D	E	F	G	H	I	J	K
1		Type	Points		Percent		Points		Percent		
2			HW1	HW2	HW1	HW2	Midterm	Final	Midterm	Final	Grade
3	Student	Max	65	55	100.00%	100.00%	75	45	100.00%	100.00%	
4	Beatrice	G	60	42	92.31%	76.36%	65	45	86.67%	100.00%	88.83%
5	Benedick	G	58	54	89.23%	98.18%	68	44	90.67%	97.78%	93.96%
6	Claudio	U	63		96.92%		68	35	90.67%	77.78%	88.46%
7	Dogberry	U	27		41.54%		55	10	73.33%	22.22%	45.70%
8	Hero	U	48		73.85%		73	40	97.33%	88.89%	86.69%
9	Verges	U	29		44.62%		39	45	52.00%	100.00%	65.54%

Figure 1: An Excel spreadsheet. The numbered rectangles are used later to facilitate discussion of region inference.

require support. To support the multiparadigmatic features, we devised a generalization of our prior dataflow test adequacy criterion that considers functions in the formulas to determine their patterns of execution. For replicated formulas, we implemented a family of techniques for combining them into regions.

To assess the resulting new methodology we performed an experiment within a prototype Excel-based WYSIWYT system on a set of non-trivial Excel spreadsheets. This experiment evaluates the costs of our methodology along several dimensions, and also compares the different techniques we have devised for finding regions to a baseline (no-regions) approach. Our results suggest that our algorithms can support the use of WYSIWYT on commercial spreadsheets; they also reveal important tradeoffs among the region inference algorithms.

## 2. BACKGROUND: WYSIWYT

The WYSIWYT methodology [5, 12, 13, 19, 21] provides several techniques and mechanisms with which end-user programmers can increase the dependability of their spreadsheets. Underlying these approaches is a dataflow test adequacy criterion that helps end users incrementally check the correctness of their spreadsheet formulas as they create or modify a spreadsheet. End-user support for this approach is provided via visual devices that are tightly integrated into the spreadsheet environment, and let users communicate testing decisions and track the adequacy of their validation efforts.

The basic computational unit of a spreadsheet is a cell’s formula. As such, our adequacy criterion is developed at the granularity of cells. Since many of the errors in spreadsheets are reference errors, we focus on the data and control dependencies between cells. This allows us to catch a wide range of faults, including reference, operator, and logic faults.

The test adequacy criterion underlying WYSIWYT is based on a model of a spreadsheet’s formulas called the Cell Relation Graph (CRG). Figure 1 shows an Excel spreadsheet, *Grades*, and Figure 2 shows a portion of the CRG corresponding to row 4 of that spreadsheet. In the CRG, nodes correspond to the cells in the spreadsheet. Within each CRG node there is a cell formula graph (CFG) that uses nodes to represent subexpressions in formulas, and edges to represent the flow of control between subexpressions. The CFG has two types of nodes, predicate nodes such as node 29 in *R4C11*, and computation nodes such as node 30 in *R4C11*.

The edges between CFGs in the CRG represent *du-associations*, which link definitions of cell values to their uses. A *definition* is an assignment to a cell of some value; each computa-

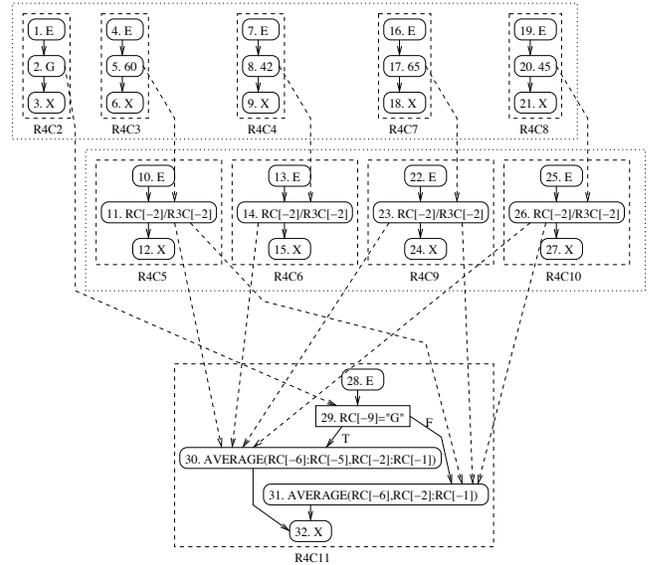


Figure 2: A CRG for the *Grades* spreadsheet

tion node provides a definition of the cell in which it resides. A *use* of a cell *C* is a reference to *C*, either in a predicate node (a p-use) or in a computation node (a c-use) of some other cell. For each use *U* of cell *C*, a *du-association* connects each definition of *C* to *U*.

Based on the CRG model, we defined the *output influencing definition-use adequacy criterion (du-adequacy)* for spreadsheets. Under this criterion, *du-associations* are classified into two categories, those ending in c-uses, and those ending in p-uses. A *du-association* ending in a c-use is considered exercised if, given the current inputs, both the definition and the use node are executed, and the cell containing the c-use or some cell downstream in dataflow from it is explicitly marked by a user as containing a value that is valid given the current assignment of values to other cells. *Du-associations* ending in a p-use are split into two separate associations, one corresponding to the p-use evaluating to true, and one corresponding to the p-use evaluating to false. A *du-association* ending in a p-use, corresponding to value *v* (either true or false) is considered exercised if the definition node is executed, the use node is executed, the use node evaluates to *v*, and the cell containing the p-use or some cell downstream in dataflow from it is marked by a user as containing a value that is valid, given the current assignment of values to other cells. A test suite is considered adequate if

all feasible (executable) du-associations in the CRG are exercised.

Spreadsheets often contain many duplicated formulas. In such cases it is impractical to require a tester to make separate decisions about each cell containing one of these duplicated formulas. Thus, in prior work [6], we extended WYSIWYT to handle regions of duplicate formulas. In that approach, a *region* is a set of cells explicitly identified by the user as sharing the same formula. (It is also possible that regions could be identified based on copy/paste actions). Figure 2 shows the regions identified by this process as boxes drawn using dotted lines around the cells included in each region. Cells *R4C5*, *R4C6*, *R4C9*, and *R4C10* form a region and all of the input cells (cells containing constants rather than formulas) form another region.

To extend the du-adequacy criterion to spreadsheets containing such regions we grouped nodes and du-associations. Within a given region, two CFG nodes are *corresponding* if they are in the same location in their respective CFGs. In Figure 2, CFG nodes 11, 14, 23, and 26 are corresponding nodes. We defined an equivalence class relationship over du-associations such that two du-associations are in the same class if and only if their definition nodes are corresponding and their use nodes are corresponding. In Figure 2, du-associations (11, 30), (14, 30), (23, 30), and (26, 30) are in the same equivalence class. Our modified adequacy criterion stated that if any du-association in an equivalence class is tested, then all of the du-associations in that class are tested.

### 3. ASSESSING PROBLEMS FOR SCALING

To better understand the problems for scaling our dataflow testing methodology to commercial spreadsheet environments and to direct our subsequent efforts, we began by assessing the ways in which such spreadsheet environments are not supported by our current methodology. We considered Excel’s documentation and user interface, and examined a large collection of Excel spreadsheets to determine what features were commonly found there and needed to be supported. We classified these features into two categories: the multiparadigmatic nature of cell formulas and formula replication. In addition, we identified a third category: spreadsheet features not directly related to cell formulas.

Where the first category of features, *the multiparadigmatic nature of cell formulas*, is concerned, there are several things to consider. First, there are combinations of simple functions such as IF which, though still purely declarative, are not supported by du-adequacy. These combinations of functions prevent cell formula graphs (as previously defined) from being constructed. In addition, Excel includes a large library (over 300) of functions, many of which represent fundamental extensions to the first-order functional paradigm represented by standard spreadsheet constructs. For example, functions like SUMIF provide limited support for constructing higher-order functions, while INDIRECT is similar to pointer arithmetic operations in other languages. Users can also implement their own “user-defined” functions (UDFs) in an imperative language (Visual Basic for Applications). These features require fundamental extensions to our du-adequacy criterion. Section 4 describes our techniques for addressing this category of features.

Where the second category of features related to *formula replication* is concerned, spreadsheets are often constructed through replication of formulas across many cells. A survey of 4432 spreadsheets from various sources found that of the 1918 that had formulas, 1714 had duplicate formulas [14]. As explained in Section 2, large numbers of replicated cells increase the cost of calculating the information required for dataflow analysis and the effort required for validation. One way to reduce cost and effort is to group similar formulas into regions that can be analyzed and tested

together. In Forms/3 there is a mechanism whereby users explicitly create regions by grouping cells together and assigning them a shared formula. However, most commercial spreadsheets, including Excel, do not allow users to explicitly create regions of cells with shared formulas, and requiring users to create regions for use by WYSIWYT places too great a burden on them. Additionally, since we wanted to be able to find regions in existing spreadsheets, the method suggested in [6] of tracking copy/paste actions is also not appropriate in the commercial spreadsheet context. Therefore, to support the use of regions in commercial spreadsheet environments, we must provide some new method for identifying regions. We provide a family of such methods in Section 5.

Where the third category of *features not directly related to cell formulas* is concerned, we determined that because WYSIWYT is primarily designed to assist users with testing cell formulas and interactions among them, these features are of only tertiary interest. However, some of these features could be accommodated by WYSIWYT. Examples of non-formula features in Excel include charts, macros for automating editing tasks or adding event-driven behaviors to spreadsheets, the ability to link to external data-sources such as databases, and the ability to insert external objects such as form fields into spreadsheets. In Section 7 we sketch some possible approaches for accommodating some of these features.

## 4. SUPPORTING THE MULTIPARADIGMATIC NATURE OF CELL FORMULAS

### 4.1 A New Adequacy Criterion

In Section 2, we presented the du-adequacy criterion that has been used in WYSIWYT research to date based on the CRG model of spreadsheets, but as outlined in Section 3, there are formula constructs in commercial spreadsheet languages that this du-adequacy criterion does not support. For example, consider cell *A3* in Figure 3. With two IF expressions added together it is no longer possible to directly convert this formula into a cell formula graph using just predicate and definition nodes. We develop our new adequacy criterion by first considering how to handle this (still purely declarative) subtlety, and then demonstrate the criterion’s ability to scale to multiparadigmatic aspects of spreadsheets.

One obvious method for dealing with the formula in cell *A3* in Figure 3 is to convert it into equivalent imperative code such as in Figure 4. But consider what happens with the formula in cell *B1* that references *A3*: the definitions in lines 2, 3, 5, and 6 do not form du-associations with the uses in *B1*. Since it is possible in WYSIWYT to validate any cell, it is not necessary for a user to consider how these definitions interact with the use in *B1* to achieve adequacy.

To address this issue, we decompose the problem of handling these formulas into two steps. The first step involves identifying *sources*, a generalized form of definitions that represent part of a cell’s computation, and *destinations*, a generalized form of uses. The second step involves connecting sources to destinations to define interactions between cells that need to be tested. To show how this process works, we walk through the process using Figure 3.

To determine the cell interactions for this example, we need to determine the sets of sources and destinations for each of the cells. Cells *A1*, *A2* and *B1* are simple cases that can be handled in the same fashion as in previous versions of WYSIWYT. Any formula that does not include conditional functions, functions that operate on or return references, or user-defined functions has only a single source. Any references in such a formula become destinations.

Cell *A3* is more interesting. To facilitate discussion of its han-

	A	B
1	5	=A3/2
2	-3	
3	=IF(A1>0,A1,0) + IF(A2>0,A2,0)	
4		

Figure 3: Spreadsheet fragment with problematic formula

1. if  $A1 > 0$  then
2.  $t1 := A1$
- else
3.  $t1 := 0$
- end if
4. if  $A2 > 0$  then
5.  $t2 := A2$
- else
6.  $t2 := 0$
- end if
7.  $A3 := t1 + t2$

Figure 4: Imperative code for cell A3 in Figure 3.

ding we use the AST in Figure 5. To determine sources for complex formulas such as this, we follow two steps. The first step is to identify the *source components* that represent different patterns of computations that can be performed by functions in the formula. The second is to combine these source components into the sources that represent the patterns of computation for the formula.

The formula for Cell A3 contains two function calls that need to be considered; namely each of the IF subexpressions. All IF’s have two possible patterns of evaluation, one that corresponds to the predicate evaluating to true, and one that corresponds to the predicate evaluating to false. We would like to capture these differing patterns of evaluation in the definition of our source components. One approach we considered was to convert all Excel functions into an equivalent UDF, and use the technique described later in Section 4.3 to determine source components and destinations. However, because this requires at least as much effort as considering the functions individually (since we do not have access to source code for the built-in functions, we would have to reverse-engineer UDF code for each of them), and because of imprecisions involved in the handling of UDFs, we chose to consider them individually. Consider the first IF (node 2 and its children in the AST); for this IF, we recognize two source components, (2, true) and (2, false). (The 2 indicates the AST node, and true or false indicates which “behavior” we are interested in). Similarly, for node 3 and its children we create the source components (3, true) and (3, false).

The source components are combined to form sources for cell A3. We consider two methods for doing this. One method is to consider sets of feasible combinations of source components. For cell A3, these combinations are  $\{(2, \text{true}), (3, \text{true})\}$ ,  $\{(2, \text{true}), (3, \text{false})\}$ ,  $\{(2, \text{false}), (3, \text{true})\}$  and  $\{(2, \text{false}), (3, \text{false})\}$ . For the current input assignment, the source  $\{(2, \text{true}), (3, \text{false})\}$  is exercised. This method captures all of the possible computation patterns for the formula and could be used when particularly rigorous testing is needed, but generates a number of sources exponential in the number of function calls in the formula.

The second method is to create a source for each source component in the formula. This creates fewer sources (in general), and on any given execution, allows multiple sources to be exercised. In our example, for the given inputs, sources (2, true) and (3, false) would be exercised. For the rest of the discussion, we assume we are using this simpler method.

Destinations for A3 are defined in the same way as uses were for

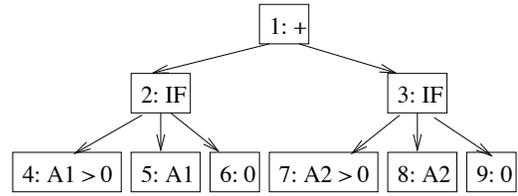


Figure 5: AST for formula in A3

du-adequacy. The destinations are (4, A1, true), (4, A1, false), (5, A1), (7, A2, true), (7, A2, false), and (8, A2).

Next we build a set of interactions that we wish to test. As we did with du-adequacy, we consider all source-destination pairs. For the example we have been considering these are  $\{(A1, (4, A1, \text{true})), (A1, (4, A1, \text{false})), (A1, (6, A1)), (A2, (7, A2, \text{true})), (A2, (7, A2, \text{false})), (A2, (9, A2)), ((2, \text{true}), B1), ((2, \text{false}), B1), ((3, \text{true}), B1), ((3, \text{false}), B1)\}$ .

Since the process of generating source components, sources, and destinations is syntax-driven, it can be readily automated using standard parsed representations (such as ASTs) of cell formulas. In addition, determining which source components and destinations are exercised requires only execution traces of the formulas, something that is easy to gather in a spreadsheet engine [19].

An additional question involves the interaction of our new du-adequacy criterion with the region mechanism for handling duplicated formulas described in Section 2. In that description we defined corresponding definitions and uses, and used those to define corresponding du-associations. For our new du-adequacy criterion we can use a similar process, defining corresponding source components and destinations based on the locations of the constructs in the cell formulas. Then two sources,  $S_1$  and  $S_2$ , are corresponding if for each source component  $C_i$  in  $S_1$  there is at least one corresponding source component in  $S_2$ , and for each source component  $C_j$  in  $S_2$  there is at least one corresponding source component in  $S_1$ . Interactions are considered corresponding if their sources and destinations are corresponding.

## 4.2 Handling Built-in Excel Functions

The previous section described our new adequacy criterion, but we still have to demonstrate how it can be applied to the built-in Excel functions that give rise to the multiparadigmatic nature of the language. To facilitate consideration of this, we partition the built-in functions into a small number of classes according to language features to which they relate: higher-order functions, meta-programming constructs, pointers, querying, and matrix operations.

### 4.2.1 Handling higher-order functions

Although higher-order functions are often considered to be an advanced programming language feature commonly associated with functional programming languages, there is support for a form of higher-order functions in Excel formulas. More precisely, Excel has a small number of functions that allow the dynamic construction of predicate expressions used for simple iterative computations, including SUMIF, COUNT, COUNTA, COUNTBLANK, and COUNTIF. To show how our approach handles these, we consider the two-parameter version of SUMIF as used in Example 1 in Table 1. The first parameter of SUMIF is a reference to a range of cells. The second parameter is a predicate to be applied to each of the cells referred to by the first parameter, which, if it evaluates to true, causes that cell’s value to be added to the running total. We can convert the SUMIF into a corresponding formula us-

1	= SUMIF(A1 : A2, "> 0")
2	= ROW(A1)
3	= OFFSET(A1, B1, C1)
4	= HLOOKUP(A1 : B5, 6)
5	{= MMULT(A1 : B2, A3 : B4)}
6	= SUMGREATERTHAN(A1 : A2, 0)

**Table 1: Examples of previously unsupported formulas.**

ing addition and IF. For Example 1, this would be = IF(A1 > 0, A1, 0) + IF(A2 > 0, A2, 0). Notice that this transformed version is the same as the formula in cell A3 of Figure 3, and the source components and destinations are the same.

One issue with this method is that it generates sources and destinations for each of the IF functions, without consideration for the symmetry between the IF expressions. To address this, we can exploit the symmetry in a fashion similar to that used for regions. By defining sets of corresponding source components and destinations, and applying the modified du-adequacy criterion, we can greatly reduce the number of interactions. In the above example, (2, true) and (3, true) are one set of corresponding source components, and (4, A1, false) and (7, A2, false) are one set of corresponding destinations.<sup>3</sup>

#### 4.2.2 Handling meta-programming constructs

Excel includes a class of functions that allow meta-programming constructs. Meta-programming constructs allow programming logic based on attributes of the source code rather than attributes of the data. These include ISBLANK, CELL, AREAS, COLUMN, COLUMNS, ROW, and ROWS. ISBLANK is a predicate that returns true if and only if the referenced cell's formula is blank. CELL allows a user to query for cell formatting, protection, and address information. AREAS, COLUMNS, and ROWS return information about the number or areas, columns, or rows included in a cell reference. COLUMN and ROW return the position (column or row) of the first cell in a cell reference. For all of these functions, the important thing to note is that they don't operate on values, and instead operate on features of the spreadsheet akin to the source code of most other languages. Consider Example 2 (Table 1). This formula returns the value 1, regardless of the value in cell A1. Therefore we do not create destinations for the references in parameters to these functions or propagate testing decisions to the referenced cells.

#### 4.2.3 Handling pointer constructs

Excel has three functions that are similar to pointer arithmetic as found in some imperative languages such as C: INDIRECT, OFFSET, and INDEX. Consider Example 3 (Table 1). Assume that cells B1 and C1 have values 1 and 2 respectively. In this case, the call to OFFSET in Example 3 returns a reference to cell C2 (1 row down and 2 columns right from cell A1). There are two potential issues with these functions. First, they can use references in their arguments. For INDEX and OFFSET, the first argument is a reference to a cell or range that is used as a starting point, and the additional arguments provide an offset relative to the original cell or range. Since the value in the range referred to in the first argument is not used (A1 in the example), we do not create any destinations for this reference or propagate testing decisions back

<sup>3</sup>An analysis similar to that used for finding regions in Section 5 could be applied to the process of finding corresponding sub-expressions in arbitrary formulas. At this time we have chosen not to do this, and to use corresponding subexpressions only in special cases such as with these functions.

to the referencing cells. However, any references used in the other arguments (B1 and C1) are dereferenced, and the corresponding values (1 and 2) are used in the calculation, therefore we can create destinations for these references and propagate testing decisions to the referenced cells as usual for computational functions.

The second issue with these functions is the handling of the returned reference (C2 in the example). For purposes of *propagating* testing decisions, it makes sense to treat the returned reference as we would a regular reference. The issue of *generating* destinations for the returned reference is more complicated. In general, these functions allow a reference to any cell in any spreadsheet ever created, although in practice their use will be much more limited (for INDEX we know the returned reference will be in the range provided in the first parameter, and for OFFSET we know the returned reference will be in the worksheet referenced in the first parameter). Since in many cases it may be intractable to calculate all of the references that can be returned by these functions, we require an approximation to determine which destinations to create.

There are several approaches that can be used for this. We could create no destinations for the returned reference; this minimizes the effort required of both the system and the user testing the spreadsheet, but may cause some interactions to be untested. We could generate the set of destinations based on the history of the spreadsheet by keeping track of the returned references of these functions and creating a new destination any time a cell that had not been used before is referenced. This method forces the user to make testing decisions that are influenced by each of the interactions seen by the system, but could still miss possible interactions. It also has the undesirable effect of having input cell changes potentially change the testedness of the spreadsheet (by creating new, necessarily untested, interactions and thereby decreasing the testedness percentage). A third possibility is to create destinations for any cells that could be referenced by the function (in the case of INDIRECT, we would limit this to cells in the workbook containing the function call). This would prevent the methodology from missing any interactions, but could create a large number of infeasible interactions. Further experimentation is needed to determine which of these possibilities is best, but for now our prototype does not create any destinations for the returned references.

#### 4.2.4 Handling query constructs

Excel has four functions, LOOKUP, HLOOKUP, VLOOKUP, and MATCH, that search for values in a range or array and return either a corresponding value or position. These are similar to standard query operations found in database query languages. Consider Example 4 (Table 1): the function searches through the cells in the top row of the range A1 : B5, in order from left to right, until it finds a cell with a value greater than or equal to 6, and returns the value in the corresponding cell from the bottom row of the range. Since sequencing seems to be an inherent property of these functions, we use the method suggested earlier of creating an equivalent UDF definition for these functions and using the testing approach described in Section 4.3.

#### 4.2.5 Handling matrix constructs

Excel has several matrix processing functions (Excel uses the term arrays) such as MMULT as used in Example 5. Formulas using these functions are typically assigned to a range of cells. Although there is some similarity between these ranges and the regions with shared formulas as used in Forms/3, it is primarily superficial. A matrix formula in Excel computes a single value (that happens to be a matrix), and "distributes" the value over a range of cells. In our new methodology, cells that participate in a matrix formula are

```

function SUMGREATERTHAN( $R, V$ )
input   $R$  : a range
         $V$  : a number
1.   $total = 0$ 
2.  for each cell  $C$  in  $R$ 
3.    if  $C > V$  then
4.       $total = total + C$ 
5.    end if
6.  end for
7.  return  $total$ 

```

**Figure 6: Implementation of UDF, SUMGREATERTHAN**

treated as an aggregate cell with a single decision box to validate the value of that cell. When validated, the testing decision propagates backwards through all referenced cells. References to cells involved in a matrix formula are treated as normal destinations with a single source, but unless the reference is a range reference that includes all cells involved in the formula, testing decisions are not propagated backwards through the formula. In all other respects, matrix functions are treated as simple computational functions.

### 4.3 Handling Imperative Code in Spreadsheets

Excel allows imperative code to be added to spreadsheets for a variety of tasks. One of the most common uses is for creating user-defined functions (UDFs). To integrate UDFs into this scheme, we need to be able to statically determine the source components and destinations relevant to those UDFs, and dynamically determine which source components and destinations are exercised when tests are applied. We use program analysis techniques on the UDFs to determine the source components and destinations.<sup>4</sup>

To determine the destinations in the UDF, we consider references in the parameters of the UDF. For each reference, we create a destination. To determine which destinations are executed, we use dynamic slicing on the return value of the UDF. In the case of a range being passed in as a parameter to the UDF (first parameter of Example 6 and  $R$  of the corresponding UDF definition in Figure 6), we create a destination for each cell in the range, and classify these destinations as corresponding destinations (similar to the corresponding destinations created for SUMIF). Therefore, for SUMGREATERTHAN as used in Example 6, the destinations are  $\{A1, A2\}$ , and they are corresponding destinations. If the formula in Example 6 was replaced with the functionally equivalent formula in A3 in Figure 3, both destinations would be considered exercised (and would in fact be considered exercised regardless of the inputs). This difference is one of the reasons we have chosen to handle the built-in functions on a case-by-case basis rather than by converting them into equivalent UDFs.

Determining the source components of the UDF is more complicated. Since source components represent subcomputations of formulas, one approach is to consider the subcomputations, or statements (which can be generalized to flow graph nodes), of the function. Then we have a source component for each statement. For SUMGREATERTHAN, the source components are  $\{1, 2, 3, 4, 7\}$ , and if the formula in Example 6 was substituted for the formula in cell A3 in Figure 3, all of these source components would be considered exercised (if A1 was changed to a value less than 0, then 4 would not be exercised); again this is weaker than the generated source components for the equivalent IF or SUMIF expressions.

<sup>4</sup>We assume that UDFs are implemented in an imperative sub-language, access spreadsheet state through parameters, and do not update global state. Examining the spreadsheets we surveyed (Section 3) revealed that all the UDFs defined there met these criteria.

## 5. SCALING IN THE PRESENCE OF REPLICATED FORMULAS

The notion of aggregating cells into regions of similar cells in spreadsheets is not new. For example, Sajaniemi defines a number of methods for doing so [22], and others have further extended his definitions [7, 8]. However, their work has focused on using these regions for visualization and auditing tasks. To use regions for our testing methodology we require that it be possible to define corresponding source components and destinations between the cells in the region, and that it be possible to efficiently update regions as formulas change, neither of these requirements is met by the approaches of [7, 8, 22].

We divide the task of inferring regions into two sub-tasks. The first subtask involves determining whether cells are similar, and the second involves grouping similar cells into regions.

### 5.1 Determining Whether Cells are Similar

The first step in developing a region inference algorithm is to define a criterion for determining whether two cells belong in the same region. Work by Sajaniemi [22] defines a number of equivalence relationships over cells in spreadsheets. For our purposes, we will consider his *formula equivalence* and *similarity* relationships, and define a variation on these that we will call *formula similarity*.

Two cells are formula equivalent if and only if one cell’s formula could have directly resulted from a copy action from the other cell’s formula. Sajaniemi also goes on to show that, under a certain referencing scheme, formula equivalence can be determined by textual comparison of the formulas. Most commercial spreadsheets, including Excel, include support for the necessary referencing scheme. In Excel it is called R1C1-style.

Sajaniemi defines two cells as being similar if and only if they are formula equivalent and *format equivalent* or neither contains any references to other cells and they are format equivalent. For our purposes, we choose to ignore format equivalence. Therefore we define two cells as formula similar if and only if they are formula equivalent or neither contains any references to other cells.

### 5.2 Finding Regions

The second issue we considered when defining our region inference techniques is the spatial relationships between the cells. Prior work has focused on rectangular areas. However, it is not necessary that regions be rectangular, and by allowing non-rectangular regions we allow larger regions to be found, thereby decreasing testing and computational effort (as well as avoiding problems with updating rectangular regions). Therefore, we consider three different candidate spatial relationships for inferring regions: discontinuous, contiguous, and rectangular. For each relationship, we present our algorithm for finding regions in an existing spreadsheet, and we then discuss mechanisms for incrementally updating regions as the spreadsheet is updated.

#### 5.2.1 Discontiguous regions

Using formula similarity and no additional constraints yields the most general concept of what constitutes a region: all cells in a worksheet that are formula similar are in the same region. Under this concept, regions can be discontinuous, containing cells that are not neighbors.

Discontiguous regions can be identified by iterating through the cells in a spreadsheet and looking up region identifiers in a hashtable indexed by cell formula. This process is linear in the number of cells in the spreadsheet.

Figure 7 presents an efficient algorithm, D-Regions, for finding the discontinuous regions in a worksheet. The algorithm relies

```

algorithm D-Regions(WS)
input   WS : a Worksheet
1. for each cell C in WS.Cells
2.   R = WS.Regions.Find(C.Formula)
3.   R.Cells.Add(C)
4.   C.RegionID = R.RegionID
5. end for

```

**Figure 7: Calculating discontinuous regions.**

Information kept per cell	
<i>Cell.RegionID</i>	An integer specifying the region a cell belongs to
<i>.Formula</i>	The R1C1-style formula for this cell
Information kept per region	
<i>Region.RegionID</i>	An integer unique to each region
<i>.Cells</i>	A list of cells belonging to the region
<i>.Formula</i>	The R1C1-style formula for this region
Information kept per worksheet	
<i>Sheet.Regions</i>	A hash table of the regions in the worksheet indexed by R1C1-style formulas
<i>.Cells</i>	An array of cells in the worksheet indexed by row and column

**Table 2: Data structures used by D-Regions.**

on the data structures shown in Table 2. D-Regions considers each cell *C* in worksheet *WS*, and finds the corresponding region *R* in *WS.Regions* (creating *R* if necessary) using *C.Formula*. It then updates *R* and *C* to reflect that *C* is part of region *R*. Since D-Regions considers each cell in the spreadsheet only once, and all of the operations within the loop can be completed in constant time, the algorithm runs in time  $O(\text{sizeof}(WS.Cells))$ .

This technique finds two regions in the *Grades* spreadsheet in Figure 1, one with the cells in the areas labeled 1, 2 and 3, and one with the cells in the area labeled 4.

To incrementally update regions there are several operations to consider. A cell’s formula could be changed (through user entry or a copy/paste operation), a cell could be inserted into the spreadsheet, or a cell could be deleted from the spreadsheet. First suppose cell *C*’s formula is changed. In this case, *C* is removed from the region it is in, and if *C* is the only cell in its region, that region is deleted. Next the technique finds the region to which *C* should be added; this is done by looking up the new region in the hashtable used to find the regions initially. This is a constant time operation.

When a cell (or cells) is (are) added to a spreadsheet, all of the cells below (or to the right of, at the user’s discretion) the inserted cell are shifted down (or to the right). This also causes references to the shifted cells to update to reflect the cells’ new locations. Each cell that references a cell that is shifted must have its region information updated. References change in a similar manner when cells are deleted from the spreadsheet, and are treated similarly.

### 5.2.2 Contiguous regions

The discontinuous algorithm is simple and efficient; however, it is important to consider what kinds of regions end users will be able to make use of. Allowing discontinuous regions requires the creation of some device to indicate the relationship between the disconnected areas that comprise regions, which could be difficult to do in a fashion that users can understand and use. Therefore, it may be useful to require regions to be contiguous.

To find contiguous regions, our technique iterates through the cells in a spreadsheet, comparing their formulas to their neighbor-

```

algorithm C-Regions(WS)
input   WS : a Worksheet
1. for row = 1 to WS.LastRow
2.   for col = 1 to WS.LastCol
3.     cell1 = WS.Cells(row, col)
4.     cell2 = WS.Cells(row + 1, col)
5.     cell3 = WS.Cells(row, col + 1)
6.     if cell1.Formula = cell2.Formula then
7.       C-Merge(WS, cell1, cell2)
8.     end if
9.     if cell1.Formula = cell3.Formula then
10.      C-Merge(WS, cell1, cell3)
11.    end if
12.  end for
13. end for
14. C-RegionsPP(WS)

```

**Figure 8: Calculating contiguous regions.**

ing cells, and merging formula similar cells into regions. With an efficiently implemented merge operation, this cost is linear in the number of cells in the spreadsheet.

Figure 8 presents algorithm C-Regions for finding contiguous regions. C-Regions iterates through each cell in a worksheet, and checks the cells to the right of and below it to see whether they belong in the same region. If they do, then the regions including the cells are merged. The cost of this algorithm is  $O(\text{number of cells} \times \text{cost of merge})$ .<sup>5</sup>

To represent contiguous regions, we use the same data structures as for D-Regions with the additional slots shown in Table 3. These data structures allow us to implement the constant time C-Merge algorithm shown in Figure 9. For each region, an additional slot *EqRegion* indicates a lower numbered region that this region is part of. C-Merge works by updating this slot. A post-processing step, C-RegionsPP (Figure 10), completes the merge when the main portion of C-Regions concludes. C-RegionsPP adds an additional  $O(\text{number of cells})$  time to the cost of C-Regions, bringing the cost of C-Regions to  $O(\text{number of cells})$ .

This technique finds three regions in *Grades* (Figure 1): (1) the cells in the areas labeled 1 and 2, (2) the cells in the area labeled 3, and (3) the cells in the area labeled 4.

Information kept per region	
<i>Region.EqRegion</i>	An integer specifying an region that this region belongs to
Information kept per worksheet	
<i>Sheet.LastRow</i>	The last non-empty row in the worksheet
<i>.LastCol</i>	The last non-empty column in the worksheet

**Table 3: Additional data for C-Regions.**

With contiguous regions, to update regions when a formula in cell *C* in region *R* is changed, there are two factors to consider. First, *C* is removed from *R*, but then it must be determined whether *C* is required to keep two or more areas of *R* connected. This can occur only if two or more of the cells adjacent to *C* were in *R*. To determine whether *R* should be split, a search is performed on the cells in *R* starting with one of the cells adjacent to *C*. If all cells in *R* that were adjacent to *C* can be reached, it is not necessary to split the region. If any adjacent cells are not reached in the search,

<sup>5</sup>An alternative approach would be a depth or breadth first search starting from cell in the worksheet and growing the regions. The cost of this approach is  $O(\text{number of cells})$ , but it may have larger constants. The cost of the merge operation depends on the data structures being used.

```

algorithm C-Merge(WS, cell1, cell2)
input   WS : a worksheet
input   cell1, cell2 : cells
1. id1 = cell1.RegionID
2. id2 = cell2.RegionID
3. if id1 = 0 and id2 = 0 then
4.   WS.Regions.AddRegion(cell1,cell2)
5. else if id1 = 0 then
6.   WS.Regions(id2).AddCell(cell1)
7. else if id2 = 0 then
8.   WS.Regions(id1).AddCell(cell2)
9. else if WS.Regions(id1).EqRegion <
   WS.Regions(id2).EqRegion then
10.  WS.Regions(id2).EqRegion =
     WS.Regions(id1).EqRegion
11. else
12.  WS.Regions(id1).EqRegion =
     WS.Regions(id2).EqRegion
13. end if

```

**Figure 9: Region merge.**

```

algorithm C-RegionsPP(WS)
input   WS : a worksheet
1. for i = 1 to WS.Regions.size
2.   WS.Regions(i).EqRegion =
     WS.Regions(WS.Regions(i).EqRegion).EqRegion
3. end for
4. for each cell in WS.cells
5.   cell.RegionID =
     WS.regions(cell.RegionID).EqRegion
6. end for

```

**Figure 10: Region post-processing.**

then the cells traversed in the search must be split off from the rest of the region. If two or more adjacent cells are not reached, the search process is repeated with another adjacent cell. In addition, it is also possible that changing the formula allows two neighbor regions to be merged. If the changed cell now has the same formula as two of its neighbors and those cells are in different regions, they need to be merged. Because of the need to potentially split or merge regions, this operation is linear in the size of  $R$  and of any other regions adjacent to  $C$ . A similar procedure is performed when a cell is deleted or inserted, taking into account changing references as in Section 5.2.1.

### 5.2.3 Rectangular regions

Forms/3 required regions to be rectangular, and Excel users may tend to think of their spreadsheets in rectangular blocks. Thus we also consider an algorithm that creates rectangular regions. To find rectangular regions, our technique first iterates through the cells, comparing their formulas to those of the cells directly above or below, creating all regions one cell wide of maximum height. It then iterates through these regions, comparing them to the regions on either side of them, and merging the adjacent regions with formula similar formulas with the same height.<sup>6</sup> Again, assuming an efficient region merge algorithm, this technique is linear in the number of cells. Figure 11 presents the algorithm, `R-Regions`; Table 4 depicts the additional data structures it uses. In lines 1 to 9, `R-Regions` creates regions one cell wide and of maximum height in  $WS$ . In lines 10 to 20, `R-Regions` finds regions created in the first pass that can merged into rectangular regions of width greater than 1. `R-Merge`, called in lines 6 and 17 is similar to the

<sup>6</sup>As described, the algorithm creates regions that favor height over width. It could instead create regions that favor width over height by altering the two passes.

`C-Merge` procedure used by `C-Regions`, but also updates region *Top*, *Left*, *Width*, and *Height* attributes. `R-RegionsPP` is the same as `C-RegionsPP`. This technique finds four regions in the Grades spreadsheet in Figure 1, one for each of the labeled areas.

```

algorithm R-Regions(WS)
input   WS : a Worksheet
1. for col = 1 to WS.LastCol
2.   for row = 1 to WS.LastRow
3.     cell1 = WS.Cells(row, col)
4.     cell2 = WS.Cells(row + 1, col)
5.     if cell1.Formula = cell2.Formula then
6.       R-Merge(WS,cell1,cell2)
7.     end if
8.   end for
9. end for
10. for row = 1 to WS.LastRow
11.  for col = 1 to WS.LastCol
12.    cell1 = WS.Cells(row, col)
13.    cell2 = WS.Cells(row, col + 1)
14.    region1 = WS.Regions(cell1.RegionID)
15.    region2 = WS.Regions(cell2.RegionID)
16.    if cell1.Formula = cell2.Formula and
       region1.Top = region2.Top and
       region1.Height = region2.Height then
17.      R-Merge(WS,cell1,cell2)
18.    end if
19.  end for
20. end for
21. R-RegionsPP(WS)

```

**Figure 11: Calculating rectangular regions.**

Information kept per region	
<i>Region.Top</i>	The top row of the region
<i>.Left</i>	The left column of the region
<i>.Width</i>	The width (in cells) of the region
<i>.Height</i>	The height (in cells) of the region

**Table 4: Additional data for R-Regions**

When a formula in cell  $C$  in region  $R$  is changed the region is split into five regions. This can be done in many ways, but to be consistent with our algorithm for finding regions it proceeds as follows: one region includes all cells in  $R$  to the left of  $C$ , one region includes all cells in  $R$  to the right of  $C$ , one includes the cells in  $R$  directly above  $C$ , one includes the cells in  $R$  directly below  $C$ , and the last includes only  $C$  (depending on where the modified cell is located in the original region, one or more of these regions may include no cells). Each of these regions is then compared with its neighbor regions to determine whether they should be merged. The total cost of this operation depends on the number of cells in the region that is broken up and its neighboring regions.

There is one important thing to note about this approach: it does not guarantee that the regions created are the same as they would be if we re-ran the batch operation. For example, in Figure 1, if the formula of cell  $I9$  was changed to match the formulas in area 3, it would be assigned to its own region. However, if it had been the same as the formulas in area 3 when the batch operation was performed, area 3 would have been divided into two regions (one for column  $I$  with  $I9$  and one for column  $J$ ). (Any update algorithm that attempted to recreate the regions that were inferred by the batch rectangular regions algorithm could potentially have wide-ranging effects on the structure of the updated regions that could be confusing to the user.) A similar procedure is performed when a cell is deleted or inserted, taking into account changing references as mentioned in Section 5.2.1.

## 6. ASSESSMENT

Ultimately, our techniques must be empirically studied in the hands of end users, to address questions about their usability and effectiveness. Such studies, however, are expensive, and before undertaking them, it is worth first assessing the more fundamental questions of whether our techniques for handling formulas and regions scale cost-effectively to real world spreadsheets, and how our different region inference algorithms perform when applied to real spreadsheets. If such assessments prove negative, they obviate the need for human studies; if they prove positive, they provide insights into the issues and factors that should be considered in designing and conducting human studies.

More formally we consider the following research questions:

**RQ1:** How much does the use of WYSIWYT as extended slow down commercial spreadsheets, and how does this vary with region inference algorithms?

**RQ2:** How much savings in testing effort can be gained by each of the region inference algorithms?

**RQ3:** How do the different region inference algorithms differ in terms of the regions they identify?

To investigate these questions, we implemented a prototype in Excel using Java and VBA. The Java component performs the underlying analysis required for determining du-associations and tracking coverage, while the VBA component evaluates formulas and expressions and displays our visual devices. The prototype version used for this study provides support for most of the functions described in Section 4, treating unsupported functions as a simple computational function for purposes of testing. (It does not yet support imperative code in spreadsheets.)

### 6.1 Objects of Analysis

As objects of analysis, we drew a sample of the spreadsheets collected in our survey (Section 3), working with just the 1826 of those spreadsheets that contained formulas and did not use macros. The 176 selected spreadsheets ranged in size from 41 to 12,121 non-empty cells, with a mean of 1,235 non-empty cells.

### 6.2 Variables and measures

#### 6.2.1 Independent variables

Our experiment involved two independent variables: region inference algorithm and spreadsheet size.

We used all three region inference algorithms described in Section 5: D-Regions, C-Regions, and R-Regions. As a baseline we also used a version without region inference, No-Regions.

To measure spreadsheet size we used the number of non-empty cells in the spreadsheet.

#### 6.2.2 Dependent variables and measures

We explored four dependent variables: total time to load the spreadsheet, time required for analysis on load, number of interactions in the spreadsheet, and number of regions found.

To measure total time to load, we used the time from when an open call is made for a spreadsheet until all of its user interface devices have been displayed. We chose this measure because it is during the loading of the spreadsheet that the most work in calculating regions and interactions must occur and because previous work has demonstrated that reasonable bounds hold on the time required to respond to other user actions within the WYSIWYT methodology.

To measure time for analysis on load, we measured the time that was spent in the analysis portion of loading the spreadsheet. The

total time to load measure lets us examine RQ1 for a specific implementation, but since Excel provides only limited options for programmability, it is useful to also track the effort required when the overhead associated with the creation of visual devices in Excel is not included. This measure includes the time required to infer regions and find all interactions in the spreadsheet.

To approximate the testing effort required by the different region algorithms we use the number of interactions in the spreadsheet. This works as an upper-bound on the amount of testing required, since if there exists a test,  $t$ , in the test suite such that all interactions that  $t$  validates are validated by other tests, then  $t$  can be removed from the test suite without sacrificing coverage.

We know that if two of our region inference algorithms find the same number of regions in a spreadsheet, they have found identical regions. Thus, measuring the number of regions found lets us quickly determine whether two algorithms act identically, and we can then further inspect interesting cases when this metric differs. For the No-Regions algorithm, the number of regions is equal to the number of cells in the spreadsheet.

### 6.3 Experiment Methodology

For each spreadsheet, we ran five different executions that each sequentially opened a spreadsheet, collected our measures, and then closed the spreadsheet. The first execution was without the WYSIWYT framework, and was used to gather base-level timings of Excel opening spreadsheets. We then did an execution for each of the four region inference algorithms utilizing the prototype Excel interface and Java analysis engine described in [9].

### 6.4 Threats to Validity

Our study involves the first use of WYSIWYT on large-scale, commercial spreadsheets and as such, addresses several threats to external validity present in earlier studies of WYSIWYT.

Where internal validity is concerned, our load times may be affected by characteristics of our implementation in Excel; however, our separate measurement of analysis time factors out most potential sources of such effects.

Where construct validity is concerned, our measures represent several factors important to the efficiency and effectiveness of WYSIWYT, but several other factors could also be considered. Our use of numbers of interactions to assess testing effort is an approximation. We have also not yet considered the costs of modification operations such as formula changes and cell addition/deletion. Finally, we have not yet considered measures involving user interaction with the system, such as the support that the various region algorithms provide for testing and for uncovering failures.

### 6.5 Data and Analysis

#### 6.5.1 RQ1: Slowdown

To determine whether our algorithms slow Excel down, we measured the total time required to load each spreadsheet in “normal” Excel. There was little variance in this time across spreadsheet sizes: all times fell between 0 and 3 seconds. For the purpose of considering slowdown we treat this time as a small constant.

Figures 12(a) and 12(b) display boxplots showing the distributions of the total time and time for analysis on loads, respectively. In both plots, times cluster at lower values. The median total time to load is between 14 and 18 seconds depending on the methodology used. Times for analysis on load, for techniques that identified regions, are considerably lower than total time to load, or analysis times for the No-Regions case. Each of the techniques had a few extreme outliers (cases that are more than three box-lengths from

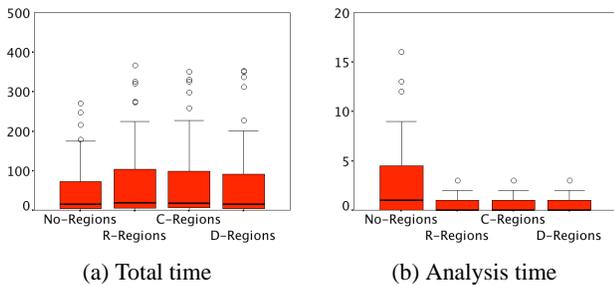
	No-Regions	R-Regions	C-Regions
R-Regions	20.93		
C-Regions	22.02	1.09	
D-Regions	23.02	2.09	1.01

**Table 5: Mean differences, total time to load.**

	No-Regions	R-Regions	C-Regions
R-Regions	9.68		
C-Regions	9.69	.01	
D-Regions	9.24	-.44	-.45

**Table 6: Mean differences, analysis time on load.**

the end of the box) including one spreadsheet that required between 25,000 and 31,000 seconds to load; these have been omitted from the figures to allow the range of typical values to be viewed.

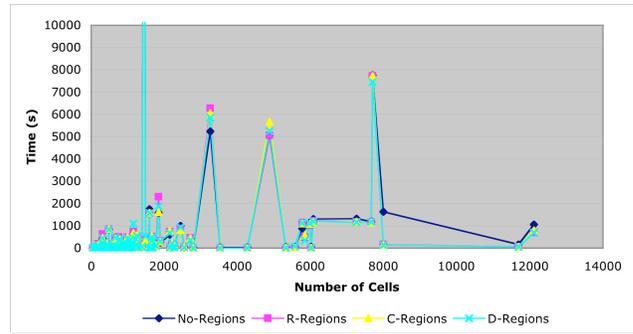


**Figure 12: Load times (seconds) per algorithm**

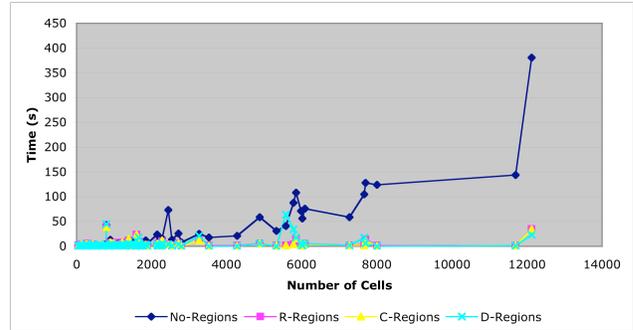
Figure 13 plots load time against spreadsheet size, for each region inference algorithm. As the figure shows, load time varies greatly, and there is no obvious trend relating total time to spreadsheet size. Also of note, the extreme outlier mentioned above is not a particularly large spreadsheet, which we thought to be unusual. Further examination of this spreadsheet showed that it had many errors that were indicated by Excel’s built-in error-checking system, and we believe that these errors interacted with the drawing of our visual devices to slow the system down.

Figure 14 provides a similar view of analysis time on load; this omits the messaging and display overhead associated with our particular implementation. For the techniques with regions, there again does not appear to be any correlation between size of spreadsheet and time; however, with the No-Regions approach it appears that such a correlation might exist. A bivariate linear correlation analysis of the data resulted in a Pearson value of .863 significant with a p-value of less than .01, indicating a reasonably strong correlation between analysis time for No-Regions and spreadsheet size.

To further examine differences in load time across region algorithms, Tables 5 and 6 show the mean differences in times (seconds) between techniques, with significant differences (paired t-test, p-value < .05) italicized. There were no significant differences in the total time between any of the techniques, but No-Regions was significantly slower for analysis time than any of the techniques.



**Figure 13: Total time to load vs. size**



**Figure 14: Analysis time on load vs. size**

### 6.5.2 RQ2: Testing effort

Table 7 shows the total number of interactions found by each techniques. No-Regions has more than 14 times as many interactions as any of the other techniques on average (significant, paired t-test, p-value < .05). Both C-Regions and R-Regions had a slightly larger number of interactions than D-Regions (significant, paired t-test, p-value < .05), and approximately the same number of interactions as each other. These results suggest that testing effort could be reduced dramatically through the use of our region inference algorithms.

### 6.5.3 RQ3: Differences in regions

Examination of the number of regions found by the different techniques shows that for 172 of the spreadsheets R-Regions found the same number of regions as C-Regions. Due to characteristics of the algorithms, this implies that R-Regions found regions identical to those found by C-Regions in these cases.

D-Regions found fewer regions than R-Regions and fewer than C-Regions, as indicated in Table 7 (significant, paired t-test, p-value < .05). Further examination shows that D-Regions found the same set of regions as C-Regions on only 36 spreadsheets.

	No-Regions	R-Regions	C-Regions	D-Regions
interactions	1162.90	81.30	81.23	50.26
regions	1234.99	39.60	39.46	20.16

**Table 7: Mean number of interactions and regions**

## 6.6 Discussion

Our analysis timings show that it is feasible to perform WYSIWYT analysis on real spreadsheets, and that with region inference and our formula extensions, WYSIWYT seems to scale quite well to larger spreadsheets. Larger load times and differences in load and analysis times suggest, however, that our current implementation incurs slowdown with respect to overhead in visual devices, a problem that we believe an implementation integrated more tightly into the Excel code could avoid. In addition, from the point of view of timing, it does not seem to make much difference which region inference algorithm is used.

As expected, D-Regions found significantly fewer (therefore larger) regions than the other techniques, which led to fewer interactions in the spreadsheet, implying less testing effort. The lack of difference between C-Regions and R-Regions, however, was somewhat surprising, although useful. As mentioned in Section 5, R-Regions are difficult to update and efficient updating algorithms could lead to an inconsistent state, a problem that C-Regions does not suffer from. Since it appears that the vast majority of contiguous regions created by users are inherently rectangular in nature, there seems to be little reason to use R-Regions. However, since there is a significant difference between the regions identified by D-Regions and R-Regions, user studies are needed to determine which of these methodologies provides the best balance between usability and efficiency for users.

## 7. CONCLUSIONS

In this paper we have presented a new adequacy criterion, aimed at supporting not only the usual dataflow relationships between formulas, but also the more challenging multiparadigmatic features of commercial spreadsheets. We show how the adequacy criterion can be applied to Excel's support of higher-order functions, meta-programming constructs, pointer constructs, query language mechanisms, and matrix constructs. We also present algorithms to support the high degree of formula replication common in commercial spreadsheets. Finally, we report on the first studies of WYSIWYT to ever be conducted within a commercial spreadsheet environment.

In our continuing work, we are considering approaches for handling other features of commercial spreadsheets. Charts could be handled as a special form of cell that have targets for each cell whose value is used to generate the chart. External data sources are a form of input cell into the system. For purposes of testing the spreadsheet, replacing them with temporary user-settable input cells would allow the user to test the logic of the spreadsheet. Using an anomaly detection mechanism on the data feeds themselves similar to that proposed in [18] could help to ensure that the data feeds are reliable.

Also, our WYSIWYT methodologies in Forms/3 include support for several dependability features other than testing, including fault localization, automated test case generation, regression test selection and replay, and assertions, and we intend to extend these features to support commercial spreadsheets. Ultimately, however, a goal of this work is to provide a system that can be used to further evaluate these methodologies with end users and large-scale spreadsheets, and in particular, that can be used in long-term ethnographic studies.

### Acknowledgements

This work was supported in part by the EUSES Consortium via NSF Grant ITR-0325273. Much of this work was performed while the first two authors were at Oregon State University.

## 8. REFERENCES

- [1] R. Abraham and M. Erwig. Header and unit inference for spreadsheets through spatial analyses. In *Symposium on Visual Languages and Human-Centric Computing*, pages 165–172, 2004.
- [2] Y. Ahmad, T. Antoniu, S. Goldwater, and S. Krishnamurthi. A type system for statically detecting spreadsheet errors. In *International Conference on Automated Software Engineering*, Oct. 2003.
- [3] T. Antoniu, P. A. Steckler, S. Krishnamurthi, E. Neuwirth, and M. Felleisen. Validating the unit correctness of spreadsheet programs. In *International Conference on Software Engineering*, May 2004.
- [4] Y. Ayalew and R. Mittermeir. Interval-based testing for spreadsheets. In *International Arab Conference on Information Technology*, Dec. 2002.
- [5] M. Burnett, C. Cook, O. Pendse, G. Rothermel, J. Summet, and C. Wallace. End-user software engineering with assertions in the spreadsheet paradigm. In *International Conference on Software Engineering*, pages 93–103, May 2003.
- [6] M. Burnett, A. Sheretov, B. Ren, and G. Rothermel. Testing homogeneous spreadsheet grids with the “What You See Is What You Test” methodology. *IEEE Transactions on Software Engineering*, 28(6):576–594, June 2002.
- [7] M. Clermont. Analyzing large spreadsheet programs. In *Working Conference on Reverse Engineering*, pages 306–315, Nov. 2003.
- [8] M. Clermont and R. Mittermeir. Auditing large spreadsheet programs. In *International Conference on Information Systems Implementation and Modelling*, Apr. 2003.
- [9] T. Creelan and M. Fisher II. Scaling up an end-user dependability framework for spreadsheets. Technical Report 04-60-09, Oregon State University, Aug. 2004.
- [10] D. Cullen. Excel snafu costs firm \$24m. *The Register*, June 19 2003.
- [11] M. Erwig, R. Abraham, I. Cooperstein, and S. Kollmansberger. Automatic generation and maintenance of correct spreadsheets. In *International Conference on Software Engineering*, pages 136–145, May 2005.
- [12] M. Fisher II, M. Cao, G. Rothermel, C. Cook, and M. Burnett. Automated test case generation for spreadsheets. In *International Conference on Software Engineering*, May 2002.
- [13] M. Fisher II, D. Jin, G. Rothermel, and M. Burnett. Test reuse in the spreadsheet paradigm. In *International Symposium on Software Reliability Engineering*, 2002.
- [14] M. Fisher II and G. Rothermel. The EUSES Spreadsheet Corpus: A shared resource for supporting experimentation with spreadsheet dependability mechanisms. In *Workshop on End-user Software Engineering*, May 2005.
- [15] R. Panko. What we know about spreadsheet errors. *Journal of End User Computing*, pages 15–21, Spring 1998.
- [16] R. Panko and R. Halverson. Spreadsheets on trial: A survey of research on spreadsheet risks. In *Hawaii International Conference on System Sciences*, Jan. 1996.
- [17] S. Prabhakararao, C. Cook, J. Ruthruff, E. Creswick, M. Main, M. Durham, and M. Burnett. Strategies and behaviors of end-user programmers with interactive fault localization. In *IEEE Symposium on Human-Centric Languages and Environments*, 2003.
- [18] O. Raz, P. Koopman, and M. Shaw. Semantic anomaly

- detection in online data sources. In *International Conference on Software Engineering*, May 2002.
- [19] G. Rothermel, M. Burnett, L. Li, C. DuPuis, and A. Sheretov. A methodology for testing spreadsheets. *ACM Transactions on Software Engineering*, pages 110–147, Jan. 2001.
- [20] K. Rothermel, C. Cook, M. Burnett, J. Schonfeld, T. Green, and G. Rothermel. WYSIWYT testing in the spreadsheet paradigm: An empirical evaluation. In *International Conference on Software Engineering*, June 2000.
- [21] J. Ruthruff, M. Burnett, and G. Rothermel. An empirical study of fault localization for end-user programmers. In *International Conference on Software Engineering*, pages 352–361, May 2005.
- [22] J. Sajaniemi. Modeling spreadsheet audit: A rigorous approach to automatic visualization. *Journal of Visual Languages and Computing*, 11(1):49–82, 2000.
- [23] A. Scott. Shurgard stock dives after auditor quits over company's accounting. *The Seattle Times*, Nov 18 2003.
- [24] R. Smith. University of Toledo loses \$2.4m in projected revenue. *The Toledo Blade*, May 1 2004.