

A Grammar for Spreadsheet Formulas Evaluated on Two Large Datasets

Efthimia Aivaloglou, David Hoepelman, Felienne Hermans

Software Engineering Research Group

Delft University of Technology

Mekelweg 4, 2628 CD Delft, the Netherlands

e.aivaloglou@tudelft.nl, d.j.hoepelman@student.tudelft.nl, f.f.j.hermans@tudelft.nl

Abstract—Spreadsheets are ubiquitous in the industrial world and often perform a role similar to other computer programs, which makes them interesting research targets. However, there does not exist a reliable grammar that is concise enough to facilitate formula parsing and analysis and to support research on spreadsheet codebases.

This paper presents a grammar for spreadsheet formulas that is compatible with the spreadsheet formula language, is compact enough to feasibly implement with a parser generator, and produces parse trees aimed at further manipulation and analysis. We evaluate the grammar against more than one million unique formulas extracted from the well known EUSES and Enron spreadsheet datasets, successfully parsing 99.99%. Additionally, we utilize the grammar to analyze these datasets and measure the frequency of usage of language features in spreadsheet formulas. Finally, we identify smelly constructs and uncommon cases in the syntax of formulas.

I. INTRODUCTION

Spreadsheets are widely used in industry: Winston [1] estimates that 90% of all analysts in industry perform calculations in spreadsheets. Their use is diverse, ranging from inventory administration to educational applications and from scientific modeling to financial systems. It is estimated that 90% of desktops have Excel installed [2] and that the number of spreadsheet programmers is bigger than that of software programmers [3].

Because of their widespread use, spreadsheets have been the topic of research since the nineties [4]. Recent research has often focused on analyzing and visualizing spreadsheets [5], [6]. More recently, researchers have attempted to define *spreadsheet smells*: applications of Fowler’s code smells to spreadsheets [7], [8], followed by approaches to refactor spreadsheets [9], [10]. These research works analyze the formulas within spreadsheets, and therefore often involve formula parsing. This is done either by using simple grammars which have not been evaluated ([10]), or through implied, undefined grammars ([5], [7]–[9]).

The above analyses are our main motivation towards a defined grammar. Having such a grammar will enable parsing spreadsheet formulas into processable parse trees which can in turn be used to analyze cell references, extract metrics, find code smells and explore the structure of spreadsheets. Essentially, a reliable and consistent grammar and its parser implementation, available to the spreadsheet research commu-

nity, can support research on spreadsheet formula codebases and can enhance the understanding and usability of research results.

To make a grammar suitable for this goal, the requirements that we set for it are (1) to be compatible with the official Excel formula language, (2) to produce parse trees suited for further manipulation and analysis, and (3) to be compact enough to feasibly implement with a parser generator. The approach that we took towards developing the grammar was gradual enrichment through trial-and-error: we started from a simple grammar containing only the most common and well known formula structures and implemented it using a parser generator. Then we repeatedly tested it against formulas extracted from spreadsheet datasets, leading to further enrichments and refinements, until all common and rare cases found in the datasets were supported. We used two major datasets that are available in the spreadsheet research community: The EUSES dataset [11] and the Enron corpus [12], jointly containing over 20,000 spreadsheets.

The contributions of this paper are (1) a concise grammar for spreadsheet formulas, (2) the evaluation of the compatibility of the grammar using two major datasets, and (3) an analysis of the common formula characteristics and of the rare grammatical cases in the datasets.

II. BACKGROUND

Spreadsheets are cell-oriented dataflow programs which are Turing complete [13]. A single spreadsheet *file* corresponds to a single (*work*)*book*. A workbook can contain any number of (*work*)*sheets*. A sheet consists of a two-dimensional grid of *cells*. The grid consists of horizontal *rows* and vertical *columns*. Rows are numbered sequentially top-to-bottom starting at 1, while columns are numbered left-to-right alphabetically, i.e. base-26 using A to Z as digits, starting at ‘A’, making column 27 ‘AA’.

A cell can be empty or contain a *constant value*, a *formula* or an *array formula*. Formulas consist of expressions which can contain constant values, arithmetic operators and *function calls* such as `SUM(. . .)` and, most importantly, *references* to other cells. Functions can be built-in or user-defined (UDFs).

A. References

References are the core component of spreadsheets. The value of any cell can be used in a formula by concatenating its column and row number, producing a reference like B5. If the value of a cell changes this new value will be propagated to all formulas that use it.

When copying a cell to another cell by default references will be adjusted by the offset, for example copying =A1 from cell B1 to C2 will cause the copied formula to become =B2. This can be prevented by prepending a \$ to the column index, row index or both. The formula =\$A\$1 will remain the same on copy while =A1 will still have its row number adjusted.

References can also be *ranges*, which are collections of cells. Ranges can be constructed by three operators: the range operator :, the union operator , (a comma) and the intersection operator \sqcap (a single whitespace). The range operator creates a rectangular range with the two cells as top-left and bottom-right corners, so =SUM(A1:B10) will sum all cells in columns A and B with row number 1 through 10. The range operator is also used to construct ranges of whole rows or columns, for example 3:5 is the range of the complete rows three through five. The union operator, which is different from the mathematical union as duplicates are allowed, combines two references, so A1,C5 will be a range of two cells, A1 and C5. Lastly the intersection operator returns only the cells which are occurring in both ranges, =A:A 5:5 will thus be equivalent to =A5.

A user can also give a name to any collection of cells, thus creating a *named range* which can be referenced in formulas by name.

B. Sheet and External References and DDE

By default references point to cells or ranges in the same sheet as the formula, but this can be modified with a prefix. A prefix consists of an identifier, followed by an exclamation mark, followed by the actual reference.

A reference to another sheet in the same workbook is indicated using the sheetname as prefix: =Sheetname!A1. References to external spreadsheet files are defined by prepending the file name in between square brackets: =[Filename]Sheetname!A1. A peculiar type of prefix are those that indicate multiple sheets: =Sheet1:Sheet10!A1 means cell A1 in Sheet1 through Sheet10. Sheet names are enclosed in single quotes if they contain special characters or spaces, e.g. ='Sheetname with space'!A1.

C. Array Formulas and Arrays

In spreadsheet programs it is possible to work with one- or two-dimensional matrices. When constructed from constant values they are called *array constants*, e.g. {1,2;3,4}. They are surrounded by curly brackets, columns are separated by commas, and rows by semicolons. Several matrix operations are available, for example =SUM({1,2,3}*10) will evaluate to 60.

Array Formulas use the same syntax as normal formulas, except that the user must signal that it is an array formula, usually by pressing *Ctrl + Shift + Enter*. Marking a formula as an array formula will enable one- or two-dimensional ranges to be treated as arrays. For example, if A1,A3,A3 contain the values 1,2,3, the array formula {=SUM(A1:A3*10)} will evaluate to 60.

III. SPREADSHEET FORMULA GRAMMAR

For previous and ongoing research the authors needed a grammar for Microsoft Excel spreadsheet formulas with the following requirements:

- 1) Compatibility with the official language
- 2) Produce parse trees suited for further manipulation and analysis with minimal post-processing required
- 3) Be compact enough to feasibly implement with a parser generator

While an official grammar for Excel formulas is published [14], it does not meet the above requirements for two reasons: first, it is over 30 pages long and contains hundreds of production rules and thus fails requirement 3. Second, because of the detail of the grammar and the large number of production rules, the resulting parse trees are very complex and thus fail requirement 2. An example is given in Figure 1(a): the relatively simple formula SUM(B2,5) results in a 37-node tree with a depth of 18 nodes.

For these reasons the authors decided to construct a new grammar with the above requirements as design goals.

A. Grammar Class

While the class of this grammar is not strictly LALR(1) due to the ambiguity discussed in Section III-F, we implemented this grammar using a LALR(1) parser generator. The present ambiguity can be solved by defining operator precedence (section III-D) and manually resolving conflicts (Section III-F). These two features are supported by most LALR(1) parser generators.

B. Lexical Analysis

Table I contains the lexical tokens of the grammar, along with their identification patterns in the regular expression language. All tokens are case-insensitive. Characters are defined as unicode characters x9,xA,xD and x20 and upwards.

Lexical analysis requires the scanner to support token priorities. Removing the necessity for token priorities is possible by altering the tokens and production rules, but makes the grammar more complicated and the resulting tree harder to use, thus being detrimental to design goals 2 and 3.

Some simple tokens (e.g. '%', '!') are directly defined in the production rules in Figure 2 in between quotes for readability and compactness.

1) *Dates*: The appearance of date and time values in spreadsheets depends on the presentation settings of cells. Internally, date and time values are stored as positive floating point numbers with the integer portion representing the

TABLE I: Lexical tokens used in the grammar

Token Name	Description	Contents	Priority
BOOL	Boolean literal	TRUE FALSE	0
CELL	Cell reference	\$? [A-Z]+ \$? [0-9]+	2
DDECALL	Dynamic Data Exchange link	' ([^ '] ")+'	0
ERROR	Error literal	#NULL! #DIV/0! #VALUE! #NAME? #NUM! #N/A	0
ERROR-REF	Reference error literal	#REF!	0
EXCEL-FUNCTION	Excel built-in function	(Any entry from the function list ¹) \ (5
FILE	External file reference	\ [0-9]+ \ (5
HORIZONTAL-RANGE	Range of rows	\$? [0-9]+ : \$? [0-9]+	0
MULTIPLE-SHEETS	Multiple sheet references	((\ [2+ : \ [2+]) (\ ' (\ [3 ")+" : (\ [3 ")+" ')) !	1
NR	Named range	[A-Z_\][A-Z0-9_\ \ [1]*	-2
NR-PREFIXED	Named range which starts with a string that could be another token	(TRUE FALSE [A-Z]+[0-9]+) [A-Z0-9_\ [1]+	3
NUMBER	An integer, floating point or scientific notation number literal	[0-9]+ ,? [0-9]* (e [0-9]+)?	0
QUOTED-FILE-SHEET	A file reference within single quotes	' \ [0-9]+ \ (\ [3 ")+" ' !	5
REF-FUNCTION	Excel built-in reference function	(INDEX OFFSET INDIRECT)\ (5
REF-FUNCTION-COND	Excel built-in conditional reference function	(IF CHOOSE)\ (5
RESERVED-NAME	An Excel reserved name	_xlnm\ [A-Z_\]+	-1
SHEET	The name of a worksheet	(\ [2+ ' (\ [3 ")+" ') !	5
STRING	String literal	" ([^ "] ")*" "	0
UDF	User Defined Function	(_xll\)? [A-Z_\][A-Z0-9_\ \ [1]* (4
VERTICAL-RANGE	Range of columns	\$? [A-Z]+ : \$? [A-Z]+	0
Placeholder character	Placeholder for	Specification	
\ [1]	Extended characters	Non-control Unicode characters x80 and up	
\ [2]	Sheet characters	Any character except ' * [] \ : / ? () ; { } # " = < > & + - * / ^ % , _	
\ [3]	Enclosed sheet characters	Any character except ' * [] \ : / ?	

¹ A function list is available as part of the reference implementation. Lists provided by Microsoft are also available in [15] and [14].

number of days since a Jan 0 1900 epoch and the fractional portion representing the portion of the day passed.

For this reason, the grammar only parses numeric dates and times and these are not distinguishable from other numbers.

2) **External References:** The file names of external references in formulas, both to external files and DDE, are not stored as part of the formula in the Microsoft Excel storage format, but instead are replaced by a numeric index. This index is then stored in a file level dictionary of external references. A formula that is presented to the user as `= [C:\Path\Filename.xlsx]Sheet1!A1` is internally stored as `[X]Sheet1!A1`, where X can be any number.

For this reason the presented grammar supports only numeric file names in external references. Adding support for full filenames can be achieved by introducing an additional token or altering the FILE token. Note that external filenames can be presented to and entered by the user in a number of different formats, depending on conditions such as whether or not the file is open in the spreadsheet program.

C. Syntactical Analysis

The complete production rules of our grammar in Extended BNF syntax are listed in Figure 2. Patterns inside { and } can be repeated zero or more times. The start symbol is *Start*.

An example parse tree produced using this grammar is drawn in Figure 1(b).

Formula and *Reference* are the two most important production rules. These are also illustrated as syntax diagrams, with most production rules expanded, in Figures 3 and 4.

The *Formula* rule covers all types of spreadsheet formula expressions: they can be constants (=5), references (=A3), function calls (=SUM(A1:A3)), array constants (= {1, 2; 3, 4}, explained in Section II-C), or reserved names (= _xlnm.Print_Area). Function calls invoke actual named (built-in or user defined) functions or operators applied to one or more formulas.

The *Reference* rule covers all types of referencing expressions, which are diverse. The simple case of a reference to a cell range can be expressed in any of the following ways:

$$\begin{aligned}
 & \text{SUM}(A1:A2) & = & \text{SUM}(A1, A2) \\
 = & \text{SUM}(\text{Sheet}!A1:A2) & = & \text{SUM}((A1, A2)) \\
 = & \text{SUM}(\text{Sheet}!A1:(A2)) & = & \text{SUM}(A1:A2:A1) \\
 = & \text{SUM}(' \text{Sheet}' !A1:A2) & = & \text{SUM}(A1:A2 A:A) \\
 = & \text{SUM}(\text{namedRange}A1A2) & &
 \end{aligned}$$

The *Reference* rule, as shown in Figure 4, supports internal (in the same or in different sheets), or external single cell

references, cell ranges, horizontal and vertical ranges, named ranges and reference-returning, built-in or user-defined, functions.

D. Operator Precedence

All operators in Excel are left-associative, including the exponentiation operator, which in most other languages is right-associative. In order to resolve ambiguities, a LALR parser generator needs the operator precedence to be defined as listed in Table II.

E. Intersection Operator

The intersection binary operator in Excel formulas is a single space. While this is straightforward to define in EBNF, it can be challenging to implement using a parser generator.

The parser generator we used for implementing the grammar supports a feature called implicit operators which was used to implement this operator. Implicit operators are operators which are left out and only implied, for example in calculus the multiplication operator is often omitted: $5a$ is equivalent to $5 \cdot a$.

F. Ambiguity

Due to trade-offs on parsing references (see section III-G1) and on parsing unions (see section III-G2) our grammar is not fully unambiguous. Ambiguity exists between the following production rules:

- 1) $\langle Reference \rangle ::= '(\langle Reference \rangle)'$
- 2) $\langle Union \rangle ::= '(\langle Reference \rangle \{ ', ' \langle Reference \rangle \})'$
- 3) $\langle Formula \rangle ::= '(\langle Formula \rangle)'$

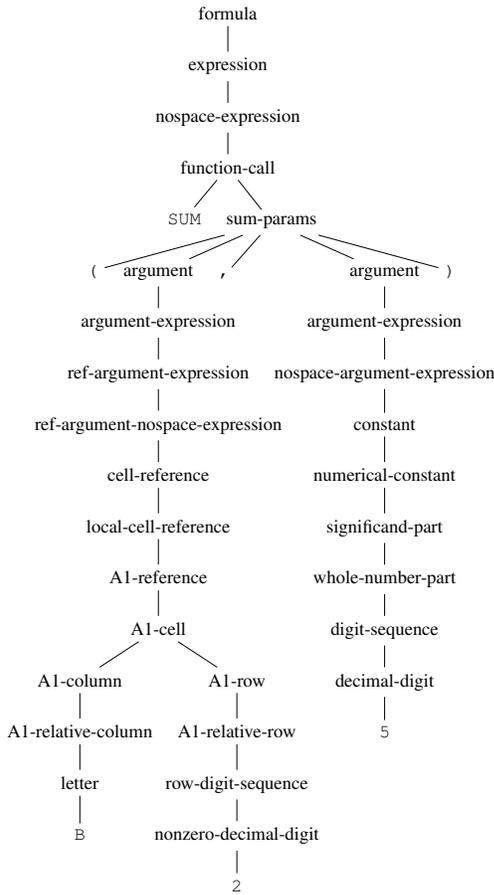
A formula like $=(A1)$ can be interpreted as either a bracketed reference, a union of one reference, or a reference within a bracketed formula.

In a LALR(1) parser the ambiguity manifests in a state where, on a $') '$ token, shifting on rule 1 and reducing on either rule 2 or 3 are possibilities, causing a shift-reduce conflict. This was solved by instructing the parser generator to shift on Rule 1 (bracketed $\langle Reference \rangle$) in case of this conflict, because this always is a correct interpretation and thus results in correct parse trees.

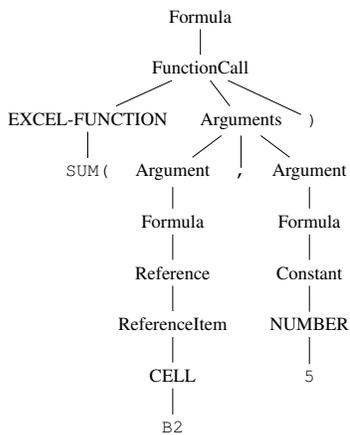
G. Trade-offs

1) **References:** References are of great importance in spreadsheet formulas, and thus of interest for analysis. To support easier analysis (Design goal 2) references have different production rules than other expressions. This causes references to be easily identified and isolated, but has the downside of increasing ambiguity, as explained in Section III-F.

Another approach would be to parse all formulas similarly and implement a type system, however this would be detrimental to ease of implementation (Design goal 3).



(a) Microsoft Excel parse tree, constructed based on reference [14]



(b) Parse tree produced using the grammar defined in this paper

Fig. 1: Parse trees for formula $SUM(B2, 5)$

```

<Start> ::= <Constant>
| '=' <Formula>
| '{=' <Formula> '}'

<Formula> ::= <Constant>
| <Reference>
| <FunctionCall>
| '(' <Formula> ')',
| <ConstantArray>
| RESERVED-NAME

<Constant> ::= NUMBER | STRING | BOOL | ERROR

<FunctionCall> ::= EXCEL-FUNCTION <Arguments> ')',
| <UnOpPrefix> <Formula>
| <Formula> '%',
| <Formula> <BinOp> <Formula>

<UnOpPrefix> ::= '+' | '-'

<BinOp> ::= '+' | '-' | '*' | '/' | '^'
| '<' | '>' | '=' | '<=' | '>=' | '<>'

<Arguments> ::= <Argument> { ',' <Argument> } | ε

<Argument> ::= <Formula> | ε

<Reference> ::= <ReferenceItem>
| <RefFunctionCall>
| '(' <Reference> ')',
| <Prefix> <ReferenceItem>
| FILE '!' DDECALL

<RefFunctionCall> ::= <Union>
| <RefFunctionName> <Arguments> ')',
| <Reference> ':' <Reference>
| <Reference> '␣' <Reference>

<ReferenceItem> ::= CELL
| <NamedRange>
| VERTICAL-RANGE
| HORIZONTAL-RANGE
| UDF <Arguments> ')',
| ERROR-REF

<Prefix> ::= SHEET
| FILE SHEET
| FILE '!'
| QUOTED-FILE-SHEET
| MULTIPLE-SHEETS
| FILE MULTIPLE-SHEETS

<RefFunctionName> ::= REF-FUNCTION
| REF-FUNCTION-COND

<NamedRange> ::= NR | NR-PREFIXED

<Union> ::= '(' <Reference> { ',' <Reference> } ')',
<ConstantArray> ::= '{' <ArrayColumns> '}'
<ArrayColumns> ::= <ArrayRows> { ';' <ArrayRows> }
<ArrayRows> ::= <ArrayConst> { ',' <ArrayConst> }
<ArrayConst> ::= <Constant>
| <UnOpPrefix> NUMBER
| ERROR-REF

```

Fig. 2: Production rules

TABLE II: Operator precedence in formulas

Precedence (higher is greater)	Operator(s)
1	= < > <= >= <>
2	&
3	+ - (binary)
4	* /
5	^
6	%
7	+ - (unary)
8	,
9	␣
10	:

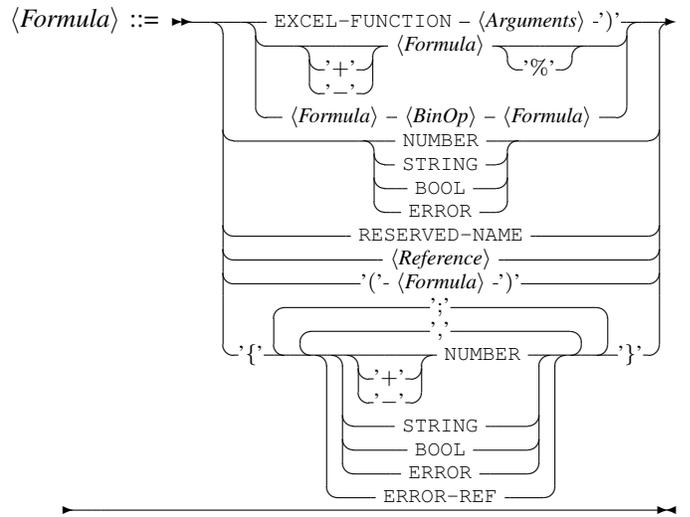


Fig. 3: Syntax diagram of the <Formula> production rule with most production rules expanded

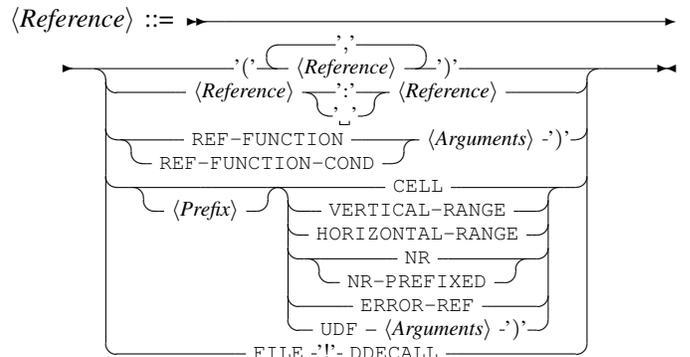


Fig. 4: Syntax diagram of the <Reference> production rule with most production rules expanded

2) *Unions*: The comma serves both as an union operator and a function argument separator. This proves challenging to correctly implement in a LALR(1) grammar.

A straightforward implementation would use production rules similar to this:

$$\langle \text{Union} \rangle ::= \langle \text{Reference} \rangle \{ ', ' \langle \text{Reference} \rangle \}$$

$$\langle \text{Arguments} \rangle ::= \langle \text{Argument} \rangle \{ ', ' \langle \text{Argument} \rangle \} | \epsilon$$

However, this will cause a reduce-reduce conflict because the parser will have a state wherein it can reduce to both a $\langle \text{Union} \rangle$ or an $\langle \text{Argument} \rangle$ on a $,$ token. Unfortunately there is no correct choice: in a formula like $=\text{SUM}(A1, 1)$ the parser must reduce on the $\langle \text{Argument} \rangle$ nonterminal, while in a formula like $=A1, A1$ the parser must reduce to the $\langle \text{Union} \rangle$ nonterminal. With the above production rules a LALR(1) parser could not correctly parse the language.

The presented grammar only parses unions in between parentheses, e.g. $=\text{SMALL}(A1, A2), 1$. This is a trade-off between a lower compatibility (Design goal 1) and an easier implementation (Design goal 3). We deem this decreased compatibility to be acceptable since unions are very rare (discussed in Section IV-B) and, in the datasets we used, all but two were with parentheses (Section IV-A).

IV. EVALUATION

In this section we explain how we implemented and evaluated the grammar using two large datasets and we discuss the obtained results and formula parse failures. In the grammar analysis in Section IV-B we examine how frequently language features occur in the formulas of the datasets.

The grammar is implemented using the Irony parser generator framework² and the resulting parser, named XLParse, is available for download³. An online demo is also available.⁴

To extract unique formulas from spreadsheets and use them as input to the parser we built a tool that opens spreadsheets using a third-party library called Gembox⁵. This tool reads all cells and identifies the formulas that are unique when adjusted for cell location (R1C1 representation), thus rejecting the formulas with adjusted references due to cell copying (e.g. formulas $=C1$ and $=C2$ are considered the same if contained in cells A1 and A2 respectively). The tool then uses each unique formula string as input to the parser.

To evaluate the grammar we use it for parsing a total of 1,035,586 unique formulas. These originate from the two major datasets available in the spreadsheet research community: The EUSES dataset [11], comprising of 4,498 spreadsheets and the Enron email corpus [12], which became available after the Enron company declared bankruptcy, comprising of 16,190 spreadsheets. We were not able to process 1087 (5.25%) of these spreadsheets, either because they are password protected, or because of read failures in the Gembox library. In total, the 19,601 spreadsheets that were processed from the two

datasets include 22,310,406 formula cells with 1,035,586 unique formulas—89,266 from the EUSES and 946,320 from the Enron dataset.

To give a rough indication, processing these two datasets and extracting these results takes around 4 hours on a computer with an Intel Core i7 processor, 16GB of RAM and a SSD.

Out of the 1,035,586 unique formulas from the two datasets that were used as input to the parser, 1,035,558 (99.99%) were parsed successfully. This satisfies our first design goal of compatibility with the official language. Regarding the second and third design goals, the implementation of the parser proved feasible and the resulting parse trees are suited for analysis and manipulation, having only 19 types of non-terminal nodes.

A. Unparsable Formulas

The 28 formulas that were not parsed using the grammar defined in Section III are:

- $=\text{NOX}, \text{Regi}$ and $=\text{SO}_2, \text{Regi}$, found in two different workbooks in the Enron dataset. These are cases of an union operations without parentheses that the grammar does not parse as explained in Section III-G2.
- $=+\text{Ë}\%$ was included in an Enron file that we assume to be either corrupt or another type of binary file, as the file is indecipherable.
- 25 formulas that are not returned correctly from the Gembox library. For example our tool reads and attempts to parse the formula $\text{IF}(=7, \text{AVERAGE}(C4:C11), 0)$ and fails, but in reality the formula is $\text{IF}(B8=7, \text{AVERAGE}(C4:C11), 0)$ which can be parsed. All these 25 cases are parsed successfully when we manually provide them as input to the parser.

B. Grammar Analysis

In this section we describe an analysis of the formulas in the datasets and measure the frequency of their characteristics. We also identify potentially smelly grammatical constructs and rare syntactical cases.

1) *Formulas and Functions*: Table III shows how frequently each of the production rules occurs in the formulas of the two datasets. Jointly, 84.91% of the formulas include a function call. Built-in value-returning functions are invoked by 35.82% of the formulas. A significant amount of formulas (286,210 or 1.28%) invoke user-defined functions—e.g., $=[1]!\text{erUserEmail}(\text{User_Id})$. A special case of user defined functions are the ones created using an Excel add-in. These are invoked as $_xll.\text{functionName}$ in 0.57% of the formulas.

Operators are used in 66.74% of the formulas, with binary operators being the most common ones, appearing in 59.77% of the formulas. Analyzing the utilization of constants, we find that 39.14% of the formulas contain at least one; more than one third (35.18%) of the formulas contain a number and 11.60% are formulas that contain text. Reserved names are uncommon, with 271 occurrences of the $_xlnm.\text{Print_Area}$ and 5 occurrences of $_xlnm.\text{Database}$.

²<https://irony.codeplex.com/>

³<https://github.com/PerfectXL/XLParse>

⁴<http://xlparser.perfectxl.nl/demo>

⁵<http://www.gemboxsoftware.com/>

TABLE III: Frequency of spreadsheet formulas with specific grammatical structures in the combined EUSES and Enron datasets

Syntax	Example	Unique formulas	Total formulas		
$\langle Formula \rangle$	=1+2	1,035,586		22,310,406	
$\langle Reference \rangle$	=E9/E10	962,783	92.97%	22,131,002	99.20%
CELL	=A5	951,521	91.88%	22,021,833	98.71%
$\langle FunctionCall \rangle$	=SUM(A5:A22)	701,626	67.75%	18,944,204	84.91%
$\langle BinOp \rangle$	=H10-H8	397,580	38.39%	13,333,844	59.77%
$\langle Constant \rangle$	=A5+134	271,585	26.23%	8,731,489	39.14%
EXCEL-FUNCTION	=SUM(A5:A22)	264,673	25.56%	7,991,329	35.82%
NUMBER	=(B8/48)*15	250,085	24.15%	7,849,495	35.18%
$\langle Prefix \rangle$	=Sheet1!B1	337,727	32.61%	5,599,011	25.10%
$\langle RefFunctionName \rangle$	=SUM(J9:INDEX(J9:J41,B43))	55,680	5.38%	5,349,237	23.98%
SHEET	=Sheet1!B1	303,981	29.35%	5,282,386	23.68%
REF-FUNCTION-COND	=IF(A1<0,0,1)	50,171	4.84%	4,872,661	21.84%
$\langle Reference \rangle$ ':' $\langle Reference \rangle$	=SUM(A5:A22)	184,451	17.81%	3,735,005	16.74%
$\langle UnOpPrefix \rangle$	=+B11+1	218,397	21.09%	3,289,326	14.74%
STRING	=IF(AD3<0,"buy","sell")	56,635	5.47%	2,587,971	11.60%
$\langle NamedRange \rangle$	=SUM(freq)	20,686	2.00%	1,645,120	7.37%
BOOL	=IF(AND(R11=1,R14=TRUE),G19,0)	7,532	0.73%	1,183,798	5.31%
FILE	=[1]Sheet1!C5	104,892	10.13%	1,135,185	5.09%
REF-FUNCTION	=SUM(J9:INDEX(J9:J41,B43))	9,907	0.96%	778,056	3.49%
QUOTED-FILE-SHEET	=' [2]Detail I&E'D62)/1000	33,781	3.26%	325,498	1.46%
UDF	=SQRT(_eoq2(C5,C4,C6,C7))	21,352	2.06%	286,210	1.28%
' $\langle Reference \rangle$ ''	=(2*(B29))/(1+B29)	6,394	0.62%	266,420	1.19%
_xl.	=_xl.RiskTriang(F9,F7,F8)	11,922	1.15%	127,348	0.57%
ERROR-REF	=AVERAGE(#REF!)	3,477	0.34%	123,447	0.55%
VERTICAL-RANGE	=COUNT(A:A)	851	0.08%	55,254	0.25%
FILE '!'	=[1]!today	2,040	0.20%	28,448	0.13%
ERROR	=IF(R14=TRUE,G19,#N/A)	379	0.04%	27,237	0.12%
'%'	=IF(E5>I8,3%,0%)	858	0.08%	16,606	0.07%
Empty argument	=DCOUNT(Lettergrades,,I80:I81)	1,343	0.13%	10,512	0.05%
Complex ranges	=SUM(I8:K8:M8)	367	0.04%	8,581	0.04%
DDECALL	=TWINDDE RSFRec!'NGH2 NET.CHNG'	3,276	0.32%	3,686	0.02%
Intersection	=Ending_Inventory_Jan	304	0.03%	2,343	0.01%
MULTIPLE-SHEETS	=SUM(Sheet1:Sheet20!I29)	173	0.02%	1,986	0.01%
Prefixed right ref. limit	=SUM('Tot-1'!\$B8:'Tot-1'!B8)	147	0.01%	1,501	0.01%
UDF reference	=[1]!wbname()	332	0.03%	855	0.00%
HORIZONTAL-RANGE	=MATCH(F3,Prices!2:2,0)	11	0.00%	836	0.00%
$\langle Union \rangle$	=LARGE(F38,C38),1)	10	0.00%	385	0.00%
RESERVED_NAME	=C23/_xlnm.Print_Area	70	0.01%	276	0.00%
FILE MULTIPLE-SHEETS	=SUM([2]Section3A:formulas!B11)	4	0.00%	189	0.00%
$\langle ConstantArray \rangle$	=FVSCHEDULE(1,0.09;0.11;0.1)	15	0.00%	19	0.00%

Regarding function arguments, spreadsheet systems allow empty arguments (e.g. =SUM(,E35,E37)) but this is rarely done—in only 0.05% of the formulas. Unions are found in only 385 formulas, e.g. =LARGE(F38,C38),1). All occurrences were arguments of the LARGE and SMALL functions—these two functions require a range of cells to be declared as a single argument, necessitating a union if the cells are not in a single range. In the EUSES dataset we also found 19 cases of constant arrays used as arguments, e.g. =FVSCHEDULE(1,{0.09;0.11;0.1}).

The array formulas rule, covering $\langle Formula \rangle$ s surrounded by brackets, is the only part of the grammar that is not evaluated. The Gembox library that we use for reading spreadsheets does not support array formulas—it reads them as regular formulas,

without the surrounding brackets. For this reason, we cannot we extract information on their frequency in the two datasets.

2) *References*: 99.2% of the formulas in the two datasets contain at least one $\langle Reference \rangle$, and 25.10% of these contain a reference that is not local, since it includes a $\langle Prefix \rangle$. External file references exist in almost 5.09% of the formulas. 16.74% of the formulas include a reference to a ':' separated cell range. Named ranges exist in 7.37% of formulas. Interestingly, horizontal and vertical ranges are rarely used (jointly, in 0.25% of the formulas). 0.55% of formulas include references to errors, e.g. =#REF!E3. These reference errors are more than four times as common as all other types of errors combined—the ERROR token exists in 0.12% of the formulas.

Moving to the edge cases of the grammar, the structures that are less common in the datasets include:

File-only external references

External references are normally in the form [File]Sheet!Cell. In 28,448 formulas (0.13%), however, the sheet is not specified, e.g. =[2]!LastTrade. These are cases of references to either external named ranges or external UDFs.

Multiple sheet references

1,986 formulas (0.01%) contain this complex case of reference, which spans across multiple sheets. An example formula is =SUM(Sheet1:Sheet10!A5), evaluated by summing all cells in position A5 from Sheet1 to Sheet10. In 189 formulas, the reference is to external files.

References to external UDFs

855 formulas (0.004%) contain references to external UDFs, for example =[1]!SheetName().

Prefixed right limits

1,501 formulas (0.01%) include a reference with a prefix in the right limit, e.g. =SUM('Deals'!F9:'Deals'!F16). In all cases this prefix is identical to the first one, as continuous ranges spanning across multiple sheets are not supported by Excel. Still, this syntax is supported.

A special case in the grammar are the functions that always return references, namely the INDEX, OFFSET and INDIRECT functions and the conditional functions that sometimes return references, namely IF and CHOOSE. For example, INDEX returns the reference of the cell at the intersection of a particular row and, optionally, column, so INDEX(B1:B10, 3) returns a reference to cell B3 and can be used in a formula as =SUM(A1:INDEX(B1:B10, 3)) being equivalent to =SUM(A1:B3). These reference returning functions are relatively common: they are found in 3.49% of the formulas, with the most common one being the INDEX (in 2.51% of formulas) and the least common one being the INDIRECT (in 0.2%). While the IF and CHOOSE functions can be part of reference expressions, there were no formulas in the datasets using them as such. An example of using those functions like this would be =SUM(IF(A1=1,A2,A5):A10), which is equivalent to =SUM(A2:A10) if A1 is 1 and to =SUM(A5:A10) otherwise.

Another rare case of references are the dynamic data exchange links, which were found in 3,686 formulas. These take the form of =Program|Topic!Arguments, e.g. =Database|TableA!Column1, and are used in Windows versions of Microsoft Excel to receive data from other applications.

3) *Smelly Grammar Constructs*: There are two constructs in the spreadsheet formula grammar that we consider to be smelly, i.e. counterintuitive or inconsistent to the rest of the grammar and error-prone: complex ranges and the intersect operator.

By *complex ranges* we mean $\langle Reference \rangle$ s that include more than two or different types of ' : ' separated $\langle ReferenceItem \rangle$ s.

	A	B	C	D	E
1					
2					
3					
4					
5					
6					

(a) A range with four limits B2:D4:C1:C5, equivalent to the area marked gray B1:D5

	A	B	C	D	E
1					
2					
3					
4					
5					
6					

(b) A range with a named range rangeC2D3:B1, equivalent to the area marked gray A1:C3

Fig. 5: Examples of references to complex ranges

An example is range B2:D4:C1:C5, illustrated in Figure 5a. The smelly aspect of complex ranges is their evaluation. Simple cell ranges are in the form top-left:bottom-right, including all cells in between the two limits. However, the limits in complex ranges are not the ones specified in the formula: they are calculated as the upper leftmost and lower rightmost cell in the square that includes all defined cells. For example, range B2:D4:C1:C5 is equivalent to B1:D5. Understanding the limits of the range becomes even less intuitive when vertical or horizontal ranges or named ranges are used, like in Figure 5b. This syntax does not add to expressiveness of the grammar: each range is still calculated as the cells within a single square, but without clearly user-defined limits. Complex ranges are rare: 8,581 formulas (0.04%) include complex ranges in the Enron dataset, and they are all defined using three cell locations.

The *intersect operator* is included in this discussion because it is `_`, a single whitespace. An advantage of this approach is that it enables more natural language definition of intersections, e.g. =SUM((Total_Cost Jan):(Total_Cost Apr.)). However, we find this representation inconsistent to the grammar, because whitespace does not carry meaning in the rest of the language. Other spreadsheet systems do not use whitespace for this operator, either by using an alternative like LibreOffice which uses ! or simply not supporting it. In the two datasets intersection operations are rare, as they are found in only 2,343 formulas (0.01%).

V. DISCUSSION AND LIMITATIONS

The currently defined formula grammar is able to parse 99,99% of the 1,035,586 unique formulas in the EUSES and the Enron datasets. In this section, we discuss a variety of issues that affect its applicability and suitability.

	A6		=Product A Store 2	
	A	B	C	D
1		Store 1	Store 2	
2	Product A	100	50	
3	Product B	110	60	
4	Product C	120	70	
5				
6			50	
7				

Fig. 6: A natural language formula in Excel 2003

A. Dialects

While other spreadsheet programs (e.g. Numbers, LibreOffice, Google Sheets) have generally adopted the Excel formula syntax, there are slight differences between programs and even Excel versions. The grammar has been designed as a generically as possible and has been enriched to include all syntactical features found in the two datasets. Both datasets, however, contain spreadsheets created in, or converted to, the Excel 2007 format. This limits the grammar support for language elements that are spreadsheet system-dependent or even version-dependent. The built-in functions list for example might change across versions, which would make the parser mistakenly recognize built-in functions as user-defined functions. Another example is found in LibreOffice, which uses `~` as the union operator instead of `,`. The presented grammar will need to be modified to account for these differences before it can be used on other dialects.

Syntactical features have also been deprecated between Excel versions. An example is regular expressions in formulas. Excel allows defining formulas that include regular expressions, for example `=SUM('S*'!A1)` or `=SUM('Sheet?'!A1)`. However, in Excel 2010 and up, regular expressions are instantly resolved—in the example, to `=SUM(Sheet2:Sheet3!A1)`, summing up all A1 cells between Sheet2 and Sheet3, where the sheets are all sheets matching the regular expression, except the one that the formula is on. This way, in Excel versions 2010 and up, saved spreadsheets never contain regular expressions.

The use of labels in formulas (referred to as natural language formulas) is another feature that was discontinued in Excel 2007. Labels were the headings that were typed above columns and before rows, and they could be used in formulas instead of defined names or cell ranges. Figure 6 shows an example in Excel 2003, where formula `=Product A Store 2` returns the intersection between the cell range with heading `Product A` and the one with heading `Store 2`. This feature is replaced in newer versions of Excel with the less error-prone named ranges feature. When processing spreadsheets with newer versions of Excel, the references that include labels are automatically converted to cell-only references—in the example, the formula is converted to `=C2`. The grammar does not support labels, and it would mistakenly parse them as named ranges.

B. Internationalization

Excel formulas differ depending on the language of the software. For example function arguments are separated by a semicolon instead of a comma in locales that use the comma as a decimal separator: the formula `=SUM(1.5,A1)` in the English version would be shown as `=SOM(1,5;A1)` in the Dutch version. Our grammar supports only the English locale. Grammars for other locales can be derived by replacing delimiters, error values and function names with their localized versions.

It is worth noting that Excel will always save formulas in either a locale-independent form (Excel 2003 and earlier format) or in its English version (Excel 2007 and later format). When interacting with Excel through its API two versions of the formula can be read or written: the English version and the version in the current locale. This makes a grammar for the English version useful, since the parser can process all spreadsheets as long as their formulas are read using the always available English locale.

C. Rejection of Invalid Formulas

As stated in the design goals in Section III, the goal of this grammar is to facilitate analysis of formulas, which means correctly parsing valid spreadsheet formulas. Rejecting invalid formulas is not among the primary goals of this grammar, as the parser will normally not encounter invalid formulas in Excel files. Furthermore, while there exist two big datasets of valid formulas, no such datasets of invalid formulas exist. As such we expect that the presented grammar will parse formulas which are not valid. Using this grammar to parse possibly-invalid formulas like user-input might thus require additional safeguards.

On one point we know the grammar to be too broad: Excel places several limitation on formulas like the number of arguments of a function (255), nested function calls (64), row number (2^{20}), column number (2^{14}) and total formula length (2^{13}), with lower numbers in older file formats. Our grammar does not incorporate any of these limits.

D. Parse Tree Correctness

While we have empirically shown a high compatibility in terms of successful parse rate, we do not have as much evidence that the produced parse trees are correct as this is only tested by usage and unit tests in the reference implementation. We have manually sampled numerous parse trees and we have found them to be correct. We believe it is unlikely that a formula parsed with the presented grammar would be interpreted differently by Excel, but we do seek additional feedback on possible erroneous parse trees from the research community.

VI. RELATED WORK

Efforts to reverse-engineer language characteristics based on existing artifacts have been successful for other languages, including COBOL [16] and C, C++, C# and Java [17].

Most related to our research on the spreadsheet formula language is the work of Badame and Dig [10] who, as part of their proposed spreadsheet refactoring approach, presented a grammar for spreadsheet formulas. However, they do not evaluate their grammar, and upon inspection one can see that key ingredients are missing: e.g. external references, intersections, unions, named ranges and operator precedence. An extension of the same grammar was used to refactor formulas by Hermans and Dig [9].

There is a large body of related work that relies on parsing spreadsheet formulas to analyze spreadsheets. This includes our own work in which we have created an algorithm to visualize spreadsheets as dataflow diagrams [5], and subsequently on detecting smells in spreadsheets [7], [8]. Related approaches exist, for example the work of Cunha that have worked on code smells [18] and smell-based fault localization [19]. These papers also analyze spreadsheet formulas but do not detail which grammar or analysis method they use.

VII. CONCLUSION

In this paper we (1) present a grammar for spreadsheet formulas, (2) evaluate it against over one million unique formulas, successfully parsing 99.99%, and (3) use it to analyze the formulas in the dataset, measure the frequency of their characteristics and find uncommon grammatical cases.

The grammar is compact and produces processable parse trees, suited for further manipulation and analysis. We believe that the grammar is reliable and concise enough to facilitate further research on spreadsheet formula codebases. It has already been applied in other works for analyzing formula characteristics, calculation chains and code smells and for applying formula transformations. The XLParse⁶ is published as open-source software and an online demo is also available.⁷

A point of improvement for the grammar is that its exact compatibility with the official Excel grammar is unknown. A comparison to the official specification could lead to either improving compatibility, or extending the number of known limitations. In general, the problem of determining whether two context-free grammars are equivalent is undecidable, but in practice several techniques have been successfully used for this purpose [20], [21].

REFERENCES

- [1] W. Winston, "Executive education opportunities," *OR/MS Today*, vol. 28, no. 4, pp. 8–10, 2001.
- [2] L. Bradley and K. McDaid, "Using bayesian statistical methods to determine the level of error in large spreadsheets," in *Proc. of ICSE '09, Companion Volume*, 2009, pp. 351–354.
- [3] C. Scaffidi, M. Shaw, and B. A. Myers, "Estimating the numbers of end users and end user programmers," in *Proc. of VL/HCC '05*, 2005, pp. 207–214.
- [4] D. Bell and M. Parr, "Spreadsheets: A research agenda," *SIGPLAN Notices*, vol. 28, no. 9, pp. 26–28, 1993.
- [5] F. Hermans, M. Pinzger, and A. van Deursen, "Supporting professional spreadsheet users by generating leveled dataflow diagrams," in *Proc. of ICSE '11*, 2011, pp. 451–460.
- [6] K. Shiozawa, K. Okada, and Y. Matsushita, "3d interactive visualization for inter-cell dependencies of spreadsheets," in *Proc. of INFOVIS*. IEEE, 1999, pp. 79–83.
- [7] F. Hermans, M. Pinzger, and A. van Deursen, "Detecting and visualizing inter-worksheet smells in spreadsheets," in *Proc. of ICSE '12*, 2012, pp. 441–451.
- [8] —, "Detecting code smells in spreadsheet formulas," in *Proc. of ICSM '12*, 2012.
- [9] F. Hermans and D. Dig, "Bumblebee: A refactoring environment for spreadsheet formulas," in *FSE 2014*, 2014, pp. 747–750.
- [10] S. Badame and D. Dig, "Refactoring meets spreadsheet formulas," in *Proc. of ICSM 2012*. IEEE, 2012, pp. 399–409.
- [11] M. Fisher and G. Rothermel, "The euses spreadsheet corpus: A shared resource for supporting experimentation with spreadsheet dependability mechanisms," *SIGSOFT Softw. Eng. Notes*, vol. 30, no. 4, pp. 1–5, May 2005.
- [12] B. Klimt and Y. Yang, "The enron corpus: A new dataset for email classification research," in *Machine Learning: ECML 2004*. Springer Berlin Heidelberg, 2004, vol. 3201, pp. 217–226.
- [13] F. Hermans, "Excel turing machine." [Online]. Available: <http://www.felienne.com/archives/2974>
- [14] Microsoft, "Excel (.xlsx) extensions to the office open xml spreadsheetml file format." [Online]. Available: [https://msdn.microsoft.com/en-us/library/dd922181\(v=office.12\).aspx](https://msdn.microsoft.com/en-us/library/dd922181(v=office.12).aspx)
- [15] —, "Excel functions (alphabetical)." [Online]. Available: <https://support.office.com/en-in/article/Excel-functions-alphabetical-b3944572-255d-4efb-bb96-c6d90033e188>
- [16] M. Van Den Brand, M. Sellink, and C. Verhoef, "Obtaining a cobol grammar from legacy code for reengineering purposes," in *Proc. of the 2nd International Workshop on the Theory and Practice of Algebraic Specifications*. Springer verlag, 1997.
- [17] V. V. Zaytsev, "Recovery, convergence and documentation of languages," Ph.D. dissertation, Vrije Universiteit, 2010.
- [18] J. Cunha, J. P. Fernandes, J. Mendes, and J. S. Hugo Pacheco, "Towards a Catalog of Spreadsheet Smells," in *Proc. of ICCSA'12*, vol. 7336. LNCS, 2012, pp. 202–216.
- [19] R. Abreu, J. Cunha, J. a. P. Fernandes, P. Martins, A. Perez, and J. a. Saraiva, "Smelling faults in spreadsheets," in *Proc. of ICSME'14*. IEEE Computer Society, 2014, pp. 111–120.
- [20] R. Lämmel and V. Zaytsev, "An introduction to grammar convergence," in *Integrated formal methods*. Springer, 2009, pp. 246–260.
- [21] B. Fischer, R. Lämmel, and V. Zaytsev, "Comparison of context-free grammars based on parsing generated test data," in *Software Language Engineering*. Springer, 2012, pp. 324–343.

⁶<https://github.com/PerfectXL/XLParse>

⁷<http://xlparser.perfectxl.nl/demo>