

ES-SQL: Visually Querying Spreadsheets

Jácome Cunha^{*†}, João Paulo Fernandes^{*‡}, Jorge Mendes^{*}, Rui Pereira^{*}, and João Saraiva^{*}

^{*} HASLab/INESC TEC & Universidade do Minho, Portugal

[†] CIICESI, ESTGF, Instituto Politécnico do Porto, Portugal

[‡] RELEASE, Universidade da Beira Interior, Portugal

{jacome,jpaulo,jorgemendes,ruipereira,jas}@di.uminho.pt

Abstract—This paper presents ES-SQL, an embedded tool for visually constructing queries over spreadsheets. This tool provides an expressive query environment which has knowledge on the business logic of spreadsheets, and by this knowledge it assists the user in defining the intended queries.

I. INTRODUCTION

Spreadsheets are one of the most successful and widely used software systems targeted specially at end users. The uses of spreadsheets range from simple single user applications to large business oriented decision making calculations.

Although conceived to be simple, easy, visual, and human-friendly, the concrete uses of spreadsheets in the real world tend to evolve into large and complex data-centric software systems. In this scenario spreadsheets often grow bigger than thousands of lines per thousands of columns and it becomes difficult to extract, query, and reason about their data.

To simplify the spreadsheet querying process, we have proposed an Embedded Spreadsheet-Structured Query Language (ES-SQL) approach for end users [1], which relies on our previous work on model-driven spreadsheets where a concise model abstracts the structure and logic of a potentially large spreadsheet [2]. This abstraction allows queries to be expressed by names, instead of column letters, referencing entities.

ES-SQL, which itself has been built on top of *MDSheet* [2], provides the powerful and expressive information extraction setting of [3], [4], [5], while maintaining a simple to read, write, and understand querying system, all in the user's spreadsheet. Moreover, the visual query building environment is synchronized with the spreadsheet model, maintaining full consistency, after model/instance evolution, between these two.

II. USING ES-SQL

Before we present ES-SQL, let us introduce a running example to use throughout the paper. Figure 1 shows a spreadsheet model written in the *ClassSheet* language [6] containing information about the *Budget* of a research team.

	A	B	C	D	E	F	G
1	Budget		Year			...	
2			year=2005			...	
3	Category	Name	Qty	Cost	Total	...	
4		name="abc"	qty=0	cost=0	total=qty*cost	...	
5						...	
6						...	

Fig. 1. A Budget *ClassSheet* model

This **Budget** model contains a **Category** class (with a **Name** attribute), and a **Year** class (with a **Year** attribute), expanding vertically and horizontally, respectively (expressed by the ellipsis). An instance of each class gives us information on the **Quantity**, the **Cost**, and the **Total**, of a **Category** in a given **Year**.

Using ES-SQL, users can write queries in their familiar spreadsheet environment, without the need of learning textual (SQL-like) notation. This is achieved by guiding users in query construction, and is achieved using drop-down boxes to select attributes, filter conditions, and other querying conditions. This eliminates both syntactic and semantic errors.

Let us now address the following question under ES-SQL:

In a budget instance of the model in Figure 1, *What was the total per year, in decreasing order, from 2010 onwards?*

In ES-SQL, we display all the information from our original model-driven query language in a human-friendly way: along with the *ClassSheet* model and instance, the ES-SQL query is also shown in its own worksheet. In Figure 2 we show how to construct a query to answer our previous question.

	A	B	C	D	E	F	G	H	I
1	Attributes	✓	Formula	Sort		Unique rows			
2	Budget					Limit rows	0	Run	
3	Year								
4	year	✓							
5	Category								
6	name								
7	(Year,Category)								
8	qty					Attribute	Op	Value	
9	cost					Year.year	>=	2010	-
10	total	✓	Sum	Z↘A			+		

Fig. 2. ES-SQL representing the answer to the previous question

The steps to construct this embedded query are as follows:

- 1) Using the drop-down boxes under the *Selected* column (B), click on cell B4 and B10 to select the *Check Mark*, selecting both *year* and *total* to be used, respectively.
- 2) Under the *Formula* column, click on cell C10 and choose *Sum* from the drop-down box list.
- 3) Click on the “+” button to add a new condition row.
- 4) Select the *Year.year* attribute, and *>=* operation using the drop-down boxes in the new condition row. Afterwards, fill in *2010* in the value cell.
- 5) Click “Run”.

III. MAPPING ES-SQL TO PLAIN SQL

This section presents the graphical to textual query translation in the ES-SQL system. Looking at Table I we can see

The first, and fourth author were funded by FCT: SFRH/BPD/73358/2010, B13-2013_PTDC/EIA-CCO/116796/2010_UMINHO, respectively.

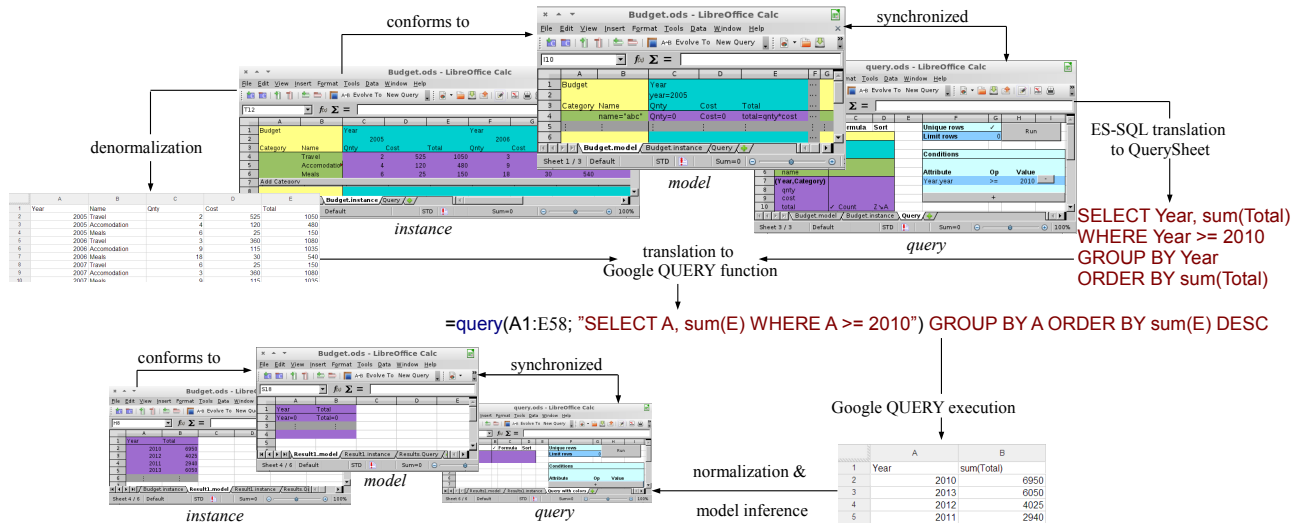


Fig. 3. Architecture of the embedded model-driven query system

the visual query on the left, with its translated textual SQL equivalent on the right. Looking at the right column, we can see bold text, which represents what part of the SQL clause we are generating, and italic text, representing the pseudocode for that translation. The remaining text is part of the query the system is currently translating. Our system automatically calculates where a *group by* is needed, a transformation which is applied to the textual translation before running the query.

IV. ARCHITECTURE OF ES-SQL

ES-SQL builds upon *QuerySheet* [4], and we follow here a similar approach to the one we followed in the past. Indeed, we translate the embedded language to our model-driven *QuerySheet* language, while the rest of the process conceptually remains the same. When the user clicks “Run”, the visual language is translated to our model-driven language. Afterwards, the data (spreadsheet instance) is automatically denormalized, and using this data alongside the model-driven query, we translate it to Google’s QUERY function using advanced compiler techniques (generalized top-down parsing). Finally, both the denormalized data and Google QUERY function are sent to be executed. The results are returned with a corresponding model/instance and a new embedded query sheet. This architecture can be seen in Figure 3.

REFERENCES

- [1] J. Cunha, J. P. Fernandes, J. Mendes, R. Pereira, and J. Saraiva, “Embedding model-driven spreadsheet queries in spreadsheet systems,” in *VLHCC’14*. IEEE, 2014, to appear.
- [2] J. Cunha, J. P. Fernandes, J. Mendes, and J. Saraiva, “MDSheet: A framework for model-driven spreadsheet engineering,” in *ICSE’12*. IEEE Press, 2012, pp. 1395–1398.
- [3] J. Cunha, J. P. Fernandes, J. Mendes, R. Pereira, and J. Saraiva, “Querying model-driven spreadsheets,” in *VLHCC’13*. IEEE, 2013, pp. 83–86.
- [4] O. Belo, J. Cunha, J. P. Fernandes, J. Mendes, R. Pereira, and J. Saraiva, “Querysheet: A bidirectional query environment for model-driven spreadsheets,” in *VLHCC’13*. IEEE, 2013, pp. 199–200.
- [5] R. Pereira, “Querying for model-driven spreadsheets,” Master’s thesis, University of Minho, 2013.

TABLE I. MAPPING THE EMBEDDED QUERY TO STANDARD SQL

Graphical	Textual																																			
<table border="1"> <thead> <tr> <th></th> <th>A</th> <th>B</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>Attributes</td> <td>✓</td> </tr> <tr> <td>2</td> <td>Class</td> <td></td> </tr> <tr> <td>3</td> <td>attr1</td> <td></td> </tr> <tr> <td>4</td> <td>attr2</td> <td>✓</td> </tr> <tr> <td>5</td> <td>⋮</td> <td>⋮</td> </tr> <tr> <td>6</td> <td>attrN</td> <td>✓</td> </tr> </tbody> </table>		A	B	1	Attributes	✓	2	Class		3	attr1		4	attr2	✓	5	⋮	⋮	6	attrN	✓	SELECT Class.attr2, ..., Class.attrN														
	A	B																																		
1	Attributes	✓																																		
2	Class																																			
3	attr1																																			
4	attr2	✓																																		
5	⋮	⋮																																		
6	attrN	✓																																		
<table border="1"> <thead> <tr> <th></th> <th>A</th> <th>C</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>Attributes</td> <td>Formula</td> </tr> <tr> <td>2</td> <td>Class</td> <td></td> </tr> <tr> <td>3</td> <td>attr1</td> <td>Sum</td> </tr> <tr> <td>4</td> <td>attr2</td> <td>Sum</td> </tr> <tr> <td>5</td> <td>⋮</td> <td>Count</td> </tr> <tr> <td>6</td> <td>attrN</td> <td>Average</td> </tr> <tr> <td>7</td> <td></td> <td>Min</td> </tr> <tr> <td></td> <td></td> <td>Max</td> </tr> </tbody> </table>		A	C	1	Attributes	Formula	2	Class		3	attr1	Sum	4	attr2	Sum	5	⋮	Count	6	attrN	Average	7		Min			Max	SELECT <i>SWITCH</i> C3 Sum: Sum (Class.attr1) Count: Count (Class.attr1) Avg: Avg (Class.attr1) Min: Min (Class.attr1) Max: Max (Class.attr1)								
	A	C																																		
1	Attributes	Formula																																		
2	Class																																			
3	attr1	Sum																																		
4	attr2	Sum																																		
5	⋮	Count																																		
6	attrN	Average																																		
7		Min																																		
		Max																																		
<table border="1"> <thead> <tr> <th></th> <th>A</th> <th>D</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>Attributes</td> <td>Sort</td> </tr> <tr> <td>2</td> <td>Class</td> <td></td> </tr> <tr> <td>3</td> <td>attr1</td> <td></td> </tr> <tr> <td>4</td> <td>attr2</td> <td>Z\A</td> </tr> <tr> <td>5</td> <td>⋮</td> <td>A\Z</td> </tr> <tr> <td>6</td> <td>attrN</td> <td>Z\A</td> </tr> </tbody> </table>		A	D	1	Attributes	Sort	2	Class		3	attr1		4	attr2	Z\A	5	⋮	A\Z	6	attrN	Z\A	SELECT ... <i>if</i> (D4 == Z\A) ORDER BY Class.attr2 DESC <i>elsif</i> (D4 == A\Z) ORDER BY Class.attr2 ASC														
	A	D																																		
1	Attributes	Sort																																		
2	Class																																			
3	attr1																																			
4	attr2	Z\A																																		
5	⋮	A\Z																																		
6	attrN	Z\A																																		
<table border="1"> <thead> <tr> <th>E</th> <th>F</th> <th>G</th> <th>H</th> <th>I</th> </tr> </thead> <tbody> <tr> <td>6</td> <td>Attribute</td> <td>Op</td> <td>Value</td> <td></td> </tr> <tr> <td>7</td> <td>Class.attr2</td> <td>=</td> <td>value</td> <td>-</td> </tr> <tr> <td>8</td> <td></td> <td>=</td> <td></td> <td></td> </tr> <tr> <td>9</td> <td></td> <td>></td> <td></td> <td></td> </tr> <tr> <td>10</td> <td></td> <td>></td> <td></td> <td></td> </tr> <tr> <td>11</td> <td></td> <td>!</td> <td></td> <td></td> </tr> </tbody> </table>	E	F	G	H	I	6	Attribute	Op	Value		7	Class.attr2	=	value	-	8		=			9		>			10		>			11		!			SELECT ... WHERE <i>SWITCH</i> C3 =: Class.attr2 = value >: Class.attr2 > value >=: Class.attr2 >= value <: Class.attr2 < value <=: Class.attr2 <= value !=: Class.attr2 != value
E	F	G	H	I																																
6	Attribute	Op	Value																																	
7	Class.attr2	=	value	-																																
8		=																																		
9		>																																		
10		>																																		
11		!																																		
<table border="1"> <thead> <tr> <th></th> <th>F</th> <th>G</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>Unique rows</td> <td>✓</td> </tr> </tbody> </table>		F	G	1	Unique rows	✓	<i>if</i> (G1 == ✓) SELECT DISTINCT ... else SELECT ...																													
	F	G																																		
1	Unique rows	✓																																		
<table border="1"> <thead> <tr> <th></th> <th>F</th> <th>G</th> </tr> </thead> <tbody> <tr> <td>2</td> <td>Limit rows</td> <td>10</td> </tr> </tbody> </table>		F	G	2	Limit rows	10	SELECT ... <i>if</i> (G5 > 0) LIMIT G5																													
	F	G																																		
2	Limit rows	10																																		

- [6] G. Engels and M. Erwig, “ClassSheets: automatic generation of spreadsheet applications from object-oriented specifications,” in *ASE’2005*. ACM, 2005, pp. 124–133.