

Bijoux: Data Generator for Evaluating ETL Process Quality

Vasileios Theodorou^a, Petar Jovanovic^a, Alberto Abelló^a,
Emona Nakuçi^a

^a*Universitat Politècnica de Catalunya, BarcelonaTech
Barcelona, Spain ({vasileios|petar|aabello}@essi.upc.edu), emona.nakuci@est.fib.upc.edu*

Abstract

Obtaining the right set of data for evaluating the fulfillment of different quality factors in the extract-transform-load (ETL) process design is rather challenging. First, the real data might be out of reach due to different privacy constraints, while manually providing a synthetic set of data is known as a labor-intensive task that needs to take various combinations of process parameters into account. More importantly, having a single dataset usually does not represent the evolution of data throughout the complete process lifespan, hence missing the plethora of possible test cases. To facilitate such demanding task, in this paper we propose an automatic data generator (i.e., *Bijoux*). Starting from a given ETL process model, *Bijoux* extracts the semantics of data transformations, analyzes the constraints they imply over input data, and automatically generates testing datasets. *Bijoux* is highly modular and configurable to enable end-users to generate datasets for a variety of interesting test scenarios (e.g., evaluating specific parts of an input ETL process design, with different input dataset sizes, different distributions of data, and different operation selectivities). We have developed a running prototype that implements the functionality of our data generation framework and here we report our experimental findings showing the effectiveness and scalability of our approach.

Keywords: Data generator, ETL, process quality

1. Introduction

2 Data-intensive processes constitute a crucial part of complex business
3 intelligence (BI) systems responsible for delivering information to satisfy the
4 needs of different end users. Besides delivering the right information to end

5 users, data-intensive processes must also satisfy various quality standards
6 to ensure that the data delivery is done in the most efficient way, whilst the
7 delivered data are of certain quality level. The quality level is usually agreed
8 beforehand in the form of service-level agreements (SLAs) or business-level
9 objects (BLOs).

10 In order to guarantee the fulfillment of the agreed quality standards (e.g.,
11 data quality, performance, reliability, recoverability; see [1, 2, 3]), an exten-
12 sive set of experiments over the designed process must be performed to test
13 the behavior of the process in a plethora of possible execution scenarios.
14 Essentially, the properties of input data (e.g., value distribution, cleanness,
15 consistency) play a major role in evaluating the resulting quality character-
16 istics of a data-intensive process. Furthermore, to obtain the finest level of
17 granularity of process metrics, quantitative analysis techniques for business
18 processes (e.g., [4]) propose analyzing the quality characteristics at the level
19 of individual activities and resources. Moreover, one of the most popular
20 techniques for quantitative analysis of process models is process simulation
21 [4], which assumes creating large number of hypothetical process instances
22 that will simulate the execution of the process flow for different scenarios.
23 In the case of data-intensive processes, the simulation should be additionally
24 accompanied by a sample of input data (i.e., *work item* in the language of
25 [4]) created for simulating a specific scenario.

26 Nonetheless, obtaining input data for performing such experiments is
27 rather challenging. Sometimes, easy access to the real source data is hard,
28 either due to data confidentiality or high data transfer costs. However, in
29 most cases the complexity comes from the fact that a single instance of avail-
30 able data, usually does not represent the evolution of data throughout the
31 complete process lifespan, and hence it cannot cover the variety of possible
32 test scenarios. At the same time, providing synthetic sets of data is known
33 as a labor intensive task that needs to take various combinations of process
34 parameters into account.

35 In the field of software testing, many approaches (e.g., [5]) have tackled
36 the problem of synthetic test data generation. However, the main focus was
37 on testing the correctness of the developed systems, rather than evaluat-
38 ing different data quality characteristics, which are critical when designing
39 data-intensive processes. Moreover, since the execution of data-intensive
40 processes is typically fully automated and time-critical, ensuring their cor-
41 rect, efficient and reliable execution, as well as certain levels of data quality
42 of their produced output is pivotal.

43 In the data warehousing (DW) context, an example of a complex, data in-
44 tensive and often error-prone data-intensive process is the extract-transform-

45 load (ETL) process, responsible for periodically populating a data warehouse
46 from the available data sources. Gartner has reported in [6] that the correct
47 ETL implementation may take up to 80% of the entire DW project. More-
48 over, the ETL design tools available in the market [7] do not provide any
49 automated support for ensuring the fulfillment of different quality param-
50 eters of the process, and still a considerable manual effort is expected from
51 the designer. Thus, we identified the real need for facilitating the task of
52 testing and evaluating ETL processes in a configurable manner.

53 In this paper, we revisit the problem of synthetic data generation for the
54 context of ETL processes, for evaluating different quality characteristics of
55 the process design. To this end, we propose an automated data generation
56 framework for evaluating ETL processes (i.e., *Bijoux*). Growing amounts
57 of data represent hidden treasury assets of an enterprise. However, due
58 to dynamic business environments, data quickly and unpredictably evolve,
59 possibly making the software that processes them (e.g., ETL) inefficient and
60 obsolete. Therefore, we need to generate delicately crafted sets of data (i.e.,
61 *bijoux*) to test different execution scenarios of an ETL process and detect
62 its behavior (e.g., *performance*) over a variety of changing parameters (e.g.,
63 *dataset size, process complexity, input data quality*) .

64 For overcoming the complexity and heterogeneity of typical ETL pro-
65 cesses, we tackle the problem of formalizing the semantics of ETL operations
66 and classifying the operations based on the part of input data they access for
67 processing. This largely facilitates *Bijoux* during data generation processes
68 both for identifying the constraints that specific operation semantics imply
69 over input data, as well as for deciding at which level the data should be
70 generated (e.g., single field, single tuple, complete dataset).

71 Furthermore, *Bijoux* offers data generation capabilities in a modular and
72 configurable manner. Instead of relying on the default data generation func-
73 tionality provided by the tool, more experienced users may also select specific
74 parts of an input ETL process, as well as desired quality characteristics to
75 be evaluated using generated datasets.

76 To illustrate the functionality of our data generation framework, we
77 introduce the running toy example that shows an ETL process (see Fig-
78 ure 1), which is a simplified implementation of the process defined in the
79 *TPC-DI benchmark*¹ for loading the *DimSecurity* table during the *Histor-*
80 *ical Load* phase². The ETL process extracts data from a file with fixed-

¹<http://www.tpc.org/tpcdi/>

²Full implementation available at: <https://github.com/AKartashoff/TPCDI-PDI/>

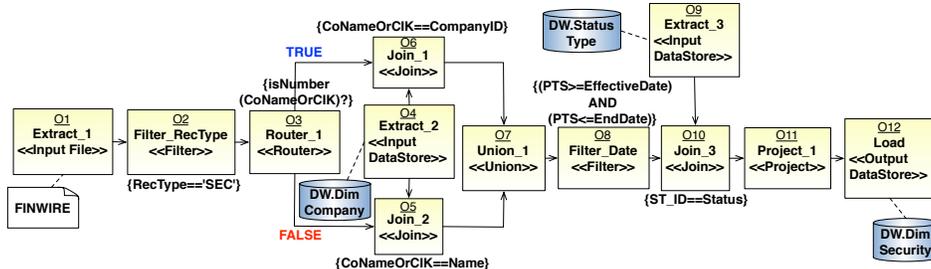


Figure 1: ETL flow example: TPC-DI DimSecurity population

81 width fields (flat file in the *Staging Area*), which is a merged collection
 82 of financial information about companies and securities coming from a fi-
 83 nancial newswire (FINWIRE) service. The input set is filtered to keep
 84 only records about Securities ($RecType == 'SEC'$) and then rows are split
 85 to two different routes, based on whether or not their values for the field
 86 $CoNameOrCIK$ are numbers ($isNumber(CoNameOrCIK)$) or not. For the
 87 first case, data are matched with data about companies through an equi-join
 88 on the company ID number ($CoNameOrCIK == CompanyID$). On the other
 89 hand, for the second case, data are matched with data about companies
 90 through an equi-join on the company name ($CoNameOrCIK == Name$). In
 91 both cases, data about companies are extracted from the *DimCompany*
 92 table of the data warehouse. Subsequently, after both routes are merged,
 93 data are filtered to keep only records for which the posting date and time (PTS)
 94 correspond to company data that are current ($(PTS \geq EffectiveDate)$ AND
 95 $(PTS \leq EndDate)$). Lastly, after data are matched with an equi-join to the
 96 data from the *StatusType* table, to get the corresponding status type for
 97 each status id ($ST_ID == Status$), only the fields of interest are maintained
 98 through a projection and then data are loaded to the *DimSecurity* table of
 99 the DW.

100 For the sake of simplicity, in what follows we will refer to the operators
 101 of our example ETL, using the label noted for each operator in Figure 1
 102 (i.e., O1 for *Extract_1*, O2 for *Filter_RecType*, etc.). Given that an ETL
 103 process model can be seen as a directed acyclic graph (DAG), *Bijoux* follows
 104 a topological order of its nodes, i.e., operations (e.g., O1, O2, O3, O4, O5, O6,
 105 O7, O8, O9, O10, O11, and O12), and extracts the found flow constraints (e.g.,
 106 $RecType == 'SEC'$ or $CoNameOrCIK == Name$). Finally, *Bijoux* generates
 107 the data that satisfy the given constraints and can be used to simulate the
 108 execution of the given ETL process.

109 Our framework, *Bijoux*, is useful during the early phases of the ETL

110 process design, when the typical time-consuming evaluation tasks are facil-
111 itated with automated data generation. Moreover, *Bijoux* can also assist
112 the complete process lifecycle, enabling easier re-evaluation of an ETL pro-
113 cess redesigned for new or changed information and quality requirements
114 (e.g., *adding new data sources*, *adding mechanisms for improving data con-*
115 *sistency*). Finally, the *Bijoux*'s functionality for automated generation of
116 synthetic data is also relevant during the ETL process deployment. It pro-
117 vides users with the valuable benchmarking support (i.e., synthetic datasets)
118 when selecting the right execution platform for their processes.

119 **Outline.** The rest of the paper is structured as follows. Section 2 for-
120 malizes the notation of ETL processes in the context of data generation and
121 presents a general overview of our approach using an example ETL pro-
122 cess. Section 3 formally presents *Bijoux*, our framework and its algorithms
123 for the automatic data generation. Section 4 introduces modified versions
124 of our example ETL process and showcases the benefits of *Bijoux* for re-
125 evaluating flow changes. In Section 5, we introduce the architecture of the
126 prototype system that implements the functionality of the *Bijoux* framework
127 and further report our experimental results. Finally, Section 6 discusses the
128 related work, while Section 7 concludes the paper and discusses possible
129 future directions.

130 2. Overview of our approach

131 In this section, we present the overview of our data generation frame-
132 work. We classify the ETL process operations and formalize the ETL process
133 elements in the context of data generation and subsequently, in a nutshell,
134 we present the overview of the data generation process of the *Bijoux* frame-
135 work.

136 2.1. ETL operation classification

137 To ensure applicability of our approach to ETL processes coming from
138 major ETL design tools and their typical operations, we performed a com-
139 parative study of these tools with the goal of producing a common subset
140 of supported ETL operations. To this end, we considered and analyzed four
141 major ETL tools in the market; two commercial, i.e., Microsoft SQL Server
142 Integration Services (SSIS) and Oracle Warehouse Builder (OWB); and two
143 open source tools, i.e., Pentaho Data Integration (PDI) and Talend Open
144 Studio for Data Integration.

145 We noticed that some of these tools have a very broad palette of specific
146 operations (e.g., PDI has a support for invoking external web services for

Table 1: Comparison of ETL operations through selected ETL tools - Part 1

Operation Level	Operation Type	Pentaho PDI	Talend Data Integration	SSIS	Oracle Warehouse Builder
Field	Field Value Alteration	Add constant	tMap	Character Map	Constant Operator
		Formula	tConvertType	Derived Column	Expression Operator
Dataset	Duplicate Removal	Number ranges	tReplaceList	Copy Column	Data Generator
		Add sequence	tUniqueRow	Data Conversion	Transformation
		Calculator	tSortRow	Fuzzy Grouping	Mapping Sequence
		Add a checksum	Sort Rows	Sort	Deduplicator
		Unique Rows	Reservoir Sampling	Percentage Sampling	Sorter
		Unique Rows (HashSet)	Sample Rows	Row Sampling	
		Sort Rows	Group by	AggregateRow	Aggregator
		Reservoir Sampling	Memory Group by	tAggregateSortedRow	
		Sample Rows	Clone Row	tReplicate	Multicast
		Aggregation	Filter Rows	tRowGenerator	
Row	Dataset Copy	Dataset Copy	tFilterRow	Conditional Split	Filter
		Duplicate Row	tMap		
		Filter	tSchemaComplianceCheck		
		Join	Merge Join	Merge Join	Joiner
		Router	Stream Lookup	Fuzzy Lookup	Key Lookup Operator
		Set Operation - Intersect	Database lookup		
		Set Operation - Difference	Merge Rows		
		Set Operation - Union	Multway Merge Join		
			Fuzzy Match		
			Switch/Case		
	Merge Rows (diff)	Conditional Split	Splitter		
	Merge Rows (diff)	Merge Join	Set Operation		
	Sorted MergeAppend streams	tMap	Set Operation		
		tMap	Set Operation		
		tUnite	Set Operation		
			Merge		
			Union All		

Table 2: Comparison of ETL operations through selected ETL tools - Part 2

Operation Level	Operation Type	Pentaho PDI	Talend Data Integration	SSIS	Oracle Warehouse Builder
Schema	Field Addition	Set field value to a constant	tMap tExtractRegexFields tAddCRCRow	Derived Column Character Map Row Count Audit Transformation	Constant Operator Expression Operator Data Generator Mapping Input/Output parameter
		String operations			
		Strings cut			
		Replace in string			
Table	Datatype Conversion Field Renaming Projection	Formula	iConvertType tMap tFilterColumns	Data Conversion Derived Column	Anydata Cast Operator
		Concat Fields			
		Add value fields changing sequence			
		Sample rows			
Value	Single Value Alteration	Select Values	tDenormalize tDenormalizeSortedRow tNormalize tSplitRow	Pivot Unpivot	Unpivot Pivot
		Select Values			
		Row Denormalizer			
		Row Normalizer			
Source Operation	Extraction	Split field to rows	tMap tReplace	Derived Column	Constant Operator Expression Operator Match-Merge Operator Mapping Input/Output parameter
		If field value is null			
		Null if			
		Modified Java Script Value			
Target Operation	Loading	SQL Execute	tFileInputDelimited tDBInput tFileInputExcel	ADO .NET / DataReader Source Excel Source Flat File Source OLE DB Source XML Source	Table Operator Flat File Operator Dimension Operator Cube Operator
		GSV file input			
		Microsoft Excel Input			
		Table input			
Target Operation	Loading	Text file input	tFileOutput tDelimited tDBOutput tFileOutputExcel	Dimension Processing Excel Destination Flat File Destination OLE DB Destination SQL Server Destination	Table Operator Flat File Operator Dimension Operator Cube Operator
		XML Input			
		Text file output			
		Microsoft Excel Output			

Table 3: List of operations considered in the framework

Considered ETL Operations	
Aggregation	Intersect
Cross Join	Join (Outer)
Dataset Copy	Pivoting
Datatype Conversion	Projection
Difference	Router
Duplicate Removal	Single Value Alteration
Duplicate Row	Sampling
Field Addition	Sort
Field Alteration	Union
Field Renaming	Unpivoting
Filter	

147 performing the computations specified by these services). Moreover, some
 148 operations can be parametrized to perform different kinds of transformation
 149 (e.g., tMap in Talend), while others can have overlapping functionalities, or
 150 different implementations for the same functionality (e.g., *FilterRows* and
 151 *JavaFilter* in PDI). Tables 1 and 2 show the resulting classification of the
 152 ETL operations from the considered tools.

153 To generalize such a heterogeneous set of ETL operations from different
 154 ETL tools, we considered the common functionalities that are supported by
 155 all the analyzed tools. As a result, we produced an extensible list of ETL
 156 operations considered by our approach (see Table 3). Notice that this list
 157 covers all operations of our running example in Figure 1, except extraction
 158 and loading ones, which are not assumed to carry any specific semantics
 159 over input data and thus are not considered operations by our classification.

160 A similar study of typical ETL operations inside several ETL tools has
 161 been performed before in [8]. However, this study classifies ETL opera-
 162 tions based on the relationship of their input and output (e.g., *unary*, *n-ary*
 163 operations). Such operation classification is useful for processing ETL opera-
 164 tions (e.g., in the context of ETL process optimization). In this paper, we
 165 further complement such taxonomy for the data generation context. There-
 166 fore, we classify ETL operations based on the part of the input table they
 167 access when processing the data (i.e., *table*, *dataset*, *row*, *schema*, *field*, or
 168 *field value*; see the first column of Table 1 and Table 2) in order to assist
 169 *Bijoux* when deciding at which level data should be generated. In Figure
 170 2, we conceptually depict the relationships between different parts of input

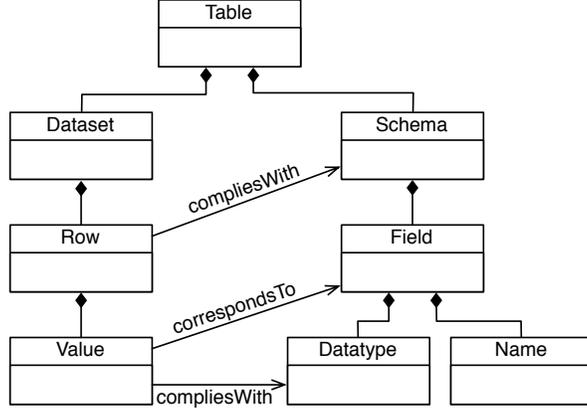


Figure 2: Table-access based classification, UML notation

171 data, which forms the basis for our ETL operation classification. In our
 172 approach, we consider the *Name* of a *Field* to act as its identifier.

173 *2.2. Formalizing ETL processes*

174 The modeling and design of ETL processes is a thoroughly studied area,
 175 both in the academia [9, 10, 11, 12] and industry, where many tools avail-
 176 able in the market often provide overlapping functionalities for the design
 177 and execution of ETL processes [7]. Still, however, no particular standard
 178 for the modeling and design of ETL processes has been defined, while ETL
 179 tools usually use proprietary (platform-specific) languages to represent an
 180 ETL process model. To overcome such heterogeneity, *Bijoux* uses a logical
 181 (platform-independent) representation of an ETL process, which in the lit-
 182 erature is usually represented as a directed acyclic graph (DAG) [12, 13].
 183 We thus formalize an ETL process as a DAG consisting of a set of nodes
 184 (\mathbf{V}), which are either source or target data stores ($\mathbf{DS} = \mathbf{DS}_S \cup \mathbf{DS}_T$) or
 185 operations (\mathbf{O}), while the graph edges (\mathbf{E}) represent the directed data flow
 186 among the nodes of the graph ($v_1 \prec v_2$). Formally:

187 $ETL = (\mathbf{V}, \mathbf{E})$, such that:
 188 $\mathbf{V} = \mathbf{DS} \cup \mathbf{O}$ and $\forall e \in \mathbf{E} : \exists(v_1, v_2), v_1 \in \mathbf{V} \wedge v_2 \in \mathbf{V} \wedge v_1 \prec v_2$
 189

190 *Data store* nodes (\mathbf{DS}) in an ETL flow graph are defined by a schema
 191 (i.e., finite list of fields) and a connection to a source (\mathbf{DS}_S) or a target
 192 (\mathbf{DS}_T) storage for respectively extracting or loading the data processed by
 193 the flow.

194 On the other side, we assume an ETL *operation* to be an atomic process-
 195 ing unit responsible for a single transformation over the input data. Notice
 196 that we model input and output data of an ETL process in terms of one or
 197 more *tables* (see Figure 2).

198 We formally define an ETL flow *operation* as a quintuple:

199

200 $o = (\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})$, where:

201

- 202 • $\mathbb{I} = \{I_1, \dots, I_n\}$ is a finite set of input tables.
- 203 • $\mathbb{O} = \{O_1, \dots, O_m\}$ is a finite set of output tables.
- 204 • \mathbf{X} ($\mathbf{X} \subseteq Attr(\mathbb{I})$) is a subset of fields of the input tables \mathbb{I} required by
 205 the operation. Notice that the function *Attr* for a given set of input
 206 or output tables, returns a set of fields (i.e., attributes) that builds the
 207 schema of these tables.
- 208 • $\mathbf{S} = (\mathbb{P}, \mathbb{F})$ represents ETL operation semantics in terms of:
 - 209 – $\mathbb{P} = \{P_1(X_1), \dots, P_p(X_p)\}$: a set of conjunctive predicates over
 210 subsets of fields in \mathbf{X} (e.g., *Age > 25*).
 - 211 – $\mathbb{F} = \{F_1(X_1), \dots, F_f(X_f)\}$: a set of functions applied over subsets
 212 of fields in \mathbf{X} (e.g., *Substr(Name, 0, 1)*). The results of these
 213 functions are used either to alter the existing fields or to generate
 214 new fields in the output table.
- 215 • \mathbf{A} is the subset of fields from the output tables, added or altered during
 216 the operation.

217 Intuitively, the above ETL notation defines a transformation of the in-
 218 put tables (\mathbb{I}) into the result tables (\mathbb{O}) by evaluating the predicate(s) and
 219 function(s) of semantics \mathbf{S} over the functionality schema \mathbf{X} and potentially
 220 generating or altering fields in \mathbf{A} .

221 An ETL operation processes input tables \mathbb{I} , hence based on the clas-
 222 sification in Figure 2, the semantics of an ETL operation should express
 223 transformations at (1) the *schema* (i.e., generated/projected-out schema),
 224 (2) the *row* (i.e., passed/modified/generated/removed rows), and (3) the
 225 *dataset* level (i.e., output cardinality).

226 In Table 4, we formalize the semantics of ETL operations considered
 227 by the framework (i.e., operations previously listed in Table 3). Notice
 228 that some operations are missing from Table 4, as they can be derived

Table 4: Table of ETL operations semantics

Op. Level	Op. Type	Op. Semantics
Value	Single Value Alteration	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{I}) = Attr(\mathbb{O}) \wedge \mathbb{I} = \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}(P_i(t_{in}[\mathbf{X}]) \rightarrow \exists t_{out} \in \mathbb{O}(t_{out}[Attr(\mathbb{O}) \setminus \mathbf{A}] = t_{in}[Attr(\mathbb{I}) \setminus \mathbf{A}] \wedge t_{out}(\mathbf{A}) = F_j(t_{in}[\mathbf{X}]))$
Field	Field Alteration	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{I}) = Attr(\mathbb{O}) \wedge \mathbb{I} = \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}, \exists t_{out} \in \mathbb{O}(t_{out}[Attr(\mathbb{O}) \setminus \mathbf{A}] = t_{in}[Attr(\mathbb{I}) \setminus \mathbf{A}] \wedge t_{out}(\mathbf{A}) = F_j(t_{in}[\mathbf{X}]))$
Row	Duplicate Row	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{I}) = Attr(\mathbb{O}) \wedge \mathbb{I} < \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}, \exists \mathbb{O}' \subseteq \mathbb{O}, \mathbb{O}' = n^{-3} \wedge \forall t_{out} \in \mathbb{O}', t_{out} = t_{in}$
	Router	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow \forall j(Attr(O_j) = Attr(\mathbb{I}) \wedge \mathbb{I} \geq O_j))$ $\forall j, \forall t_{in} \in \mathbb{I}(P_j(t_{in}[x_j]) \rightarrow \exists t_{out} \in O_j(t_{out} = t_{in}))$
	Filter	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{O}) = Attr(\mathbb{I}) \wedge \mathbb{I} \geq \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}(P_j(t_{in}[\mathbf{X}]) \rightarrow \exists t_{out} \in \mathbb{O}, (t_{out} = t_{in}))$
	Join	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{O}) = Attr(I_1) \cup Attr(I_2) \wedge \mathbb{O} \leq I_1 \times I_2))$ $\forall t_{in_1} \in I_1, t_{in_2} \in I_2, (P(t_{in_1}[x_1], t_{in_2}[x_2]) \rightarrow \exists t_{out} \in \mathbb{O}(t_{out} = t_{in_1} \bullet t_{in_2}))$
	Union	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(I_1) = Attr(I_2) \wedge Attr(\mathbb{O}) = Attr(I_1) \wedge \mathbb{O} = I_1 + I_2))$ $\forall t_{in} \in (I_1 \cup I_2) \rightarrow \exists t_{out} \in \mathbb{O}(t_{out} = t_{in})$
	Difference	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(I_1) = Attr(I_2) \wedge Attr(\mathbb{O}) = Attr(I_1) \wedge \mathbb{O} \leq I_1))$ $\forall t_{in}(t_{in} \in I_1 \wedge t_{in} \notin I_2) \rightarrow \exists t_{out} \in \mathbb{O}(t_{out} = t_{in})$
	Dataset	Aggregation
Sort		$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{I}) = Attr(\mathbb{O}) \wedge \mathbb{I} = \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}, \exists t_{out} \in \mathbb{O}(t_{out} = t_{in})$ $\forall t_{out}, t_{out}' \in \mathbb{O}(t_{out}[\mathbf{X}] < t_{out}'[\mathbf{X}] \rightarrow t_{out} < t_{out}')$
Duplicate Removal		$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{I}) = Attr(\mathbb{O}) \wedge \mathbb{I} \geq \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}, \exists t_{out} \in \mathbb{O}(t_{out} = t_{in})$
Dataset Copy		$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow \forall j(Attr(O_j) = Attr(\mathbb{I}) \wedge \mathbb{I} = O_j))$ $\forall j, \forall t_{in} \in \mathbb{I}, \exists t_{out} \in O_j, (t_{out} = t_{in})$
Schema	Projection	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{O}) = Attr(\mathbb{I}) \setminus \mathbf{X} \wedge \mathbb{I} = \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}, \exists t_{out} \in \mathbb{O}(t_{out}[Attr(\mathbb{O})] = t_{in}[Attr(\mathbb{I}) \setminus \mathbf{X}]))$
	Field Renaming	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{O}) = (Attr(\mathbb{I}) \setminus \mathbf{X}) \cup \mathbf{A} \wedge \mathbb{I} = \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}, \exists t_{out} \in \mathbb{O}(t_{out}[Attr(\mathbb{O}) \setminus \mathbf{A}] = t_{in}[Attr(\mathbb{I})] \wedge t_{out}[\mathbf{A}] = F(t_{in}[\mathbf{X}]))$
	Field Addition	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{O}) = Attr(\mathbb{I}) \cup \mathbf{A} \wedge \mathbb{I} = \mathbb{O}))$ $\forall t_{in} \in \mathbb{I}, \exists t_{out} \in \mathbb{O}(t_{out}[Attr(\mathbb{O}) \setminus \mathbf{A}] = t_{in}[Attr(\mathbb{I})] \wedge t_{out}[\mathbf{A}] = F(t_{in}[\mathbf{X}]))$
Table	Pivoting	$\forall(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A})(F(\mathbb{I}, \mathbb{O}, \mathbf{X}, \mathbf{S}, \mathbf{A}) \rightarrow (Attr(\mathbb{O}) = (Attr(\mathbb{I}) \setminus \mathbf{X}) \cup \mathbf{A} \wedge \mathbb{O} = \mathbb{I} _a \wedge \mathbb{I} = \mathbb{O} _a))$ $\forall t_{in} \in \mathbb{I}, \forall a \in Attr(\mathbb{I}), \exists t_{out} \in \mathbb{O}, \exists b \in Attr(\mathbb{O})(t_{out}[b] = t_{in}[a]))$

229 from the semantics of other listed operations (e.g., Intersection as a special
230 case of Join, Unpivoting as an inverse operation to Pivoting, and Datatype
231 Conversion as a special case of Field Alteration using a specific conversion
232 function).

233 In our approach, we use such formalization of operation semantics to
234 automatically extract the constraints that an operation implies over the
235 input data, hence to further generate the input data for covering such
236 operations. However, notice that some operations in Table 4 may imply
237 specific semantics over input data that are not explicitly expressed in the
238 given formalizations (e.g., *Field Addition/Alteration*, *Single Value Alter-*
239 *ation*). Such semantics may span from simple arithmetic expressions (e.g.,
240 $yield = dividend \div DM_CLOSE$), to complex user defined functions ex-
241 pressed in terms of an ad hoc script or code snippets. While the former case
242 can be easily tackled by powerful expression parsers [13], in the later case
243 the operation’s semantics must be carefully analyzed to extract the con-
244 straints implied over input data (e.g., by means of the static code analysis,
245 as suggested in [14]).

246 2.3. *Bijoux* overview

247 Intuitively, starting from a logical model of an ETL process and the se-
248 mantics of ETL operations, *Bijoux* analyzes how the fields of input data
249 stores are restricted by the semantics of the ETL process operations (e.g.,
250 *filter* or *join* predicates) in order to generate the data that satisfy these
251 restrictions. To this end, *Bijoux* moves iteratively through the topologi-
252 cal order of the nodes inside the DAG of an ETL process and extracts the
253 semantics of each ETL operation to analyze the constraints that the opera-
254 tions imply over the input fields. At the same time, *Bijoux* also follows the
255 constraints' dependencies among the operations to simultaneously collect
256 the necessary parameters for generating data for the correlated fields (i.e.,
257 *value ranges*, *datatypes*, and *the sizes of generated data*). Using the collected
258 parameters, *Bijoux* then generates input datasets to satisfy all found con-
259 strains, i.e., to simulate the execution of selected parts of the data flow. The
260 algorithm can be additionally parametrized to support data generation for
261 different execution scenarios.

262 Typically, an ETL process should be tested for different sizes of input
263 datasets (i.e., different *scale factors*) to examine its scalability in terms of
264 growing data. Importantly, *Bijoux* is extensible to support data generation
265 for different characteristics of input datasets (e.g., *size*), fields (e.g., *value*
266 *distribution*) or ETL operations (e.g., *operation selectivity*). We present
267 in more detail the functionality of our data generation algorithm in the
268 following section.

269 3. *Bijoux* data generation framework

270 The data generation process includes four main stages (i.e., 1 - *path*
271 *enumeration*, 2 - *constraints extraction*, 3 - *constraints analysis*, and 4 - *data*
272 *generation*).

273 3.1. *Preliminaries and Challenges*

274 We first discuss some of the important challenges of generating data for
275 evaluating general ETL flows, as well as the main structures maintained
276 during the data generation process.

277 The workflow-graph structure of the ETL logical model that we adopt
278 for our analysis consists of ETL operations as graph *nodes*, input data stores
279 as graph *sources* and output data stores as graph *sinks*. In particular, input

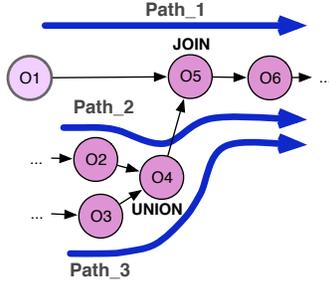
³_n is the number of replicas in the Replicate Row operation semantics

280 data stores, as well as routing operations (e.g., Routers) that direct rows
281 to different outputs based on specified conditions, introduce alternative di-
282 rected paths of the input graph (in the rest of the paper referred to as *paths*),
283 which can be followed by input data. Hence, there are two properties of the
284 generated input data that can be defined:

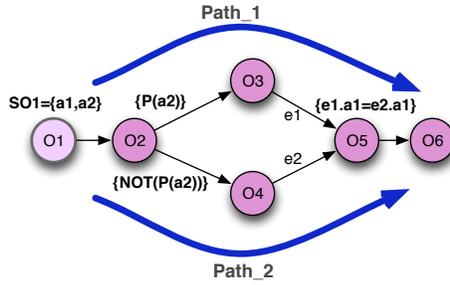
- 285 • *Path Coverage*: Input data are sufficient to “cover” a specific path,
286 i.e., each and every edge (or node) that is on this path is visited by at
287 least one row of data.

- 288 • *Flow Coverage*: Input data are sufficient to “cover” the complete flow
289 graph, i.e., each and every edge (or node) of the flow graph is visited
290 by at least one row of data.

291 The apparently simple case of *Path Coverage* hides an inherent complex-
292 ity, deriving from the fact that some joining operations (i.e., *joining nodes*;
293 e.g., Join, Intersection) require the involvement of multiple paths in order
294 to direct data to their output. In addition, new fields are introduced to the
295 flow either through input data stores or Field Addition operations (see Table
296 4), while the fields from different paths are *fused*/joined together through
297 joining operations. This in turn implies two facts: i) *Path Coverage* is not
298 guaranteed by generating the right input data only for the input data store
299 that is involved in a specific path; instead, data generation should be con-
300 ducted for a combination of paths (i.e., their included input data stores),
301 and ii) during the *Path Coverage* analysis, referring to a field solely by its
302 name is not sufficient; the same field might participate in multiple paths
303 from a combination of paths, in each path holding different properties com-
304 ing from extracted constraints of different operations. Thus, the *name* of a
305 field should be combined with a *pathid* to identify one distinct entity with
306 specific properties.



(a) Alternative path combinations for coverage of the same path



(b) Multiple rows from same input source required for coverage

Figure 3: Notable cases of graph patterns

307 In Figure 3, we show some notable cases of graph patterns that require
 308 special attention during the coverage analysis, as described above.

309 In Figure 3a, we can see how the coverage of **Path_1** ($O1 \rightarrow O5 \rightarrow O6 \dots$)
 310 needs multiple paths to be considered for data generation, because of the
 311 joining operation $O5$ that requires multiple inputs (e.g., a *Join* operation).
 312 Thus, coverage can be ensured by using alternative combinations, either
 313 **Path_1** in combination with **Path_2** ($\dots O2 \rightarrow O4 \rightarrow O5 \rightarrow O6 \dots$), or **Path_1** in
 314 combination with **Path_3** ($\dots O2 \rightarrow O4 \rightarrow O5 \rightarrow O6 \dots$). It should be mentioned
 315 that operation $O4$ is of a merging type that does not require both of its
 316 incoming edges to be crossed in order to pass data to its output (i.e., a
 317 *Union* operation) and thus **Path_2** and **Path_3** can be used interchangeably
 318 for coverage.

319 In Figure 3b, we show how the coverage of one path might require the
 320 generation of multiple rows for the same input source. For example, for the
 321 Path Coverage of **Path_1** ($O1 \rightarrow O2 \rightarrow O3 \rightarrow O5 \rightarrow O6$) it is required to additionally
 322 generate data for **Path_2** ($O1 \rightarrow O2 \rightarrow O4 \rightarrow O5 \rightarrow O6$), because of the
 323 existence of the joining operation $O5$. It should be noticed here that fields

324 $a1$ and $a2$ in `Path_1` belong to a different instance than in `Path_2`, since
 325 the condition of the routing operator `O2` imposes different predicates over
 326 $a2$ for different paths (i.e., $P(a2)$ and $\text{NOT}(P(a2))$), respectively). Hence,
 327 at least two different rows from the same input data store are required for
 328 *Path Coverage of Path_1*.

329 **Example.** For illustrating the functionality of our algorithm, we will
 330 use the running example introduced in Section 1 (see Figure 1). For the sake
 331 of simplicity, we will not use the complete schemata of the input data stores
 332 as specified in the *TPC-DI benchmark*, but instead we assume simplified
 333 versions, where the only fields present are the ones that are used in the
 334 ETL flow, i.e., taking part in predicates or functions. In this manner, input
 335 data stores of the example ETL flow are: $\mathbb{I} = \{O1, O4, O9\}$, with schemata
 336 $SO1 = \{PTS, RecType, Status, CoNameOrCIK\}$, $SO4 = \{CompanyID,$
 337 $Name, EffectiveDate, EndDate\}$ and $SO9 = \{ST_ID, ST_NAME\}$; whilst
 338 a topological order of its nodes is: $\{O1, O2, O3, O4, O5, O6, O7, O8, O9,$
 339 $O10, O11, O12\}$. Besides this running example, we will also use the auxiliary
 340 example graph from Figure 4a to support the description of the complete
 341 functionality of *Bijoux* □

342 3.2. Data structures

343 Before going into the details of algorithms 1 and 2 in Section 3.4, we
 344 present the main structures maintained by these algorithms.

345 While analyzing a given ETL graph, in Algorithm 1, *Bijoux* builds the
 346 following structures that partially or completely record the path structures
 347 of the input ETL graph (i.e., path traces):

- 348 • *Path Traces* (\mathbb{PT}) collection keeps traces of operations and edges that
 349 have been visited, when following a specific path up to a specific node
 350 in the ETL graph. Traces of individual paths PT ($PT \in \mathbb{PT}$) are built
 351 incrementally and thus, following a specific path on the graph, if a
 352 Path Trace $PT1$ is generated at an earlier point than the generation of
 353 a Path Trace $PT2$, then $PT1$ will include a subset of the trace of $PT2$
 354 (i.e., $PT1 \subseteq PT2$). From an implementation point of view, each PT
 355 holds a *Signature* as a property, which can be a string concatenation
 356 of graph elements that shows which route has been followed in the
 357 case of alternative paths. This enables very efficient PT analysis and
 358 comparisons by simply applying string operations.

359 **Example.** Referring to our running example in Section 1 we can have
 360 the following signature of a Path Trace $PT1$:

361 $Sig(PT1) = "I[O1].S[O2, true].S[O3, true].J[O6, e1]"$

362 From this signature we can conclude that PT1 starts from I (i.e., Input
363 Source): $O1$; passes through S (i.e., Splitting Operation): $O2$ coming
364 from its outgoing edge that corresponds to the evaluation: $true$ of
365 its condition; passes through S (i.e., Splitting Operation): $O3$ coming
366 from its outgoing edge that corresponds to the evaluation: $true$; passes
367 through J (i.e., Joining Operation): $O6$ coming from its incoming edge:
368 $e1$; and so on. For some operations (e.g., Joins) it makes sense to keep
369 track of the incoming edge through which they have been reached in
370 the specific path and for some others (e.g., Routers), it makes sense
371 to keep track of the outgoing edge that was followed for the path.
372 Looking at the following signature of Path Trace PT2:
373 $Sig(PT2) = "I[O1].S[O2, true].S[O3, true]"$, we can infer that PT1
374 and PT2 are on the same path of the ETL graph, PT2 being generated
375 at an “earlier” point, since the signature of PT2 is a substring of the
376 signature of PT1. \square

- 377 • *Tagged Nodes* (\mathbb{TN}) structure records, for each node, the set of paths
378 (i.e., operations and edges) reaching that node from the input data
379 store nodes (i.e., source nodes). Thus, each node is “tagged” with a
380 set of Path Traces (\mathbb{PT}) which are being built incrementally, as ex-
381 plained above.

382 **Example.** Referring to our running example, within \mathbb{TN} the $O7$ op-
383 eration node will be “tagged” with four different path traces, $PT1$,
384 $PT2$, $PT3$ and $PT4$ with the following signatures:

- 385 - $Sig(PT1) = "I[O1].S[O2, true].S[O3, true].J[O6, e1].J[O7, e1]"$
- 386 - $Sig(PT2) = "I[O1].S[O2, true].S[O3, false].J[O5, e1].J[O7, e2]"$
- 387 - $Sig(PT3) = "I[O4].J[O6, e2].J[O7, e1]"$
- 388 - $Sig(PT4) = "I[O4].J[O5, e2].J[O7, e2]"$ \square

- 389 • *Final path traces* (\mathbb{FP}) structure records all the *complete* (i.e., *source-*
390 *to-sink*) paths from the input ETL graph, by maintaining all source-
391 to-sink Path Traces (i.e., the union of all Path Traces that tag sink
392 nodes).

393 When it comes to formalizing the main structure that is being built by
394 Algorithm 2 (i.e., data generation pattern), we define its structure as follows:

- 395 • A data generation pattern (*Pattern*) consists of a set of path con-
396 straints (i.e., *pathConstr*), where each path constraint is a set of
397 constraints over the input fields introduced by the operations of an
398 individual path. Formally:

399 $Pattern = \{pathConstr_i | i = 1, \dots, pathNum\}$

400

401 **Example.** In our running example (Figure 1), so as to cover the
402 path $Path1=(O1 \rightarrow O2 \rightarrow O3 \rightarrow O6 \rightarrow O7 \rightarrow O8 \rightarrow O10 \rightarrow O11 \rightarrow O12)$, addi-
403 tionally, the path $Path2=(O4 \rightarrow O6 \rightarrow O7 \rightarrow O8 \rightarrow O10 \rightarrow O11 \rightarrow O12)$ and
404 the path $Path3=(O9 \rightarrow O10 \rightarrow O11 \rightarrow O12)$ need to be covered as well,
405 because of the equi-join operators $O6$ and $O10$. The *Pattern* would
406 then consist of three constraints sets ($pathConstr1$, $pathConstr2$ and
407 $pathConstr3$), one for each (source-to-sink) path of the flow that has
408 to be covered. \square

409 • A path constraint (i.e., $pathConstr_i$) consists of a set of constraints
410 over individual fields of the given path (i.e., $fieldConstr$). Formally:
411 $pathConstr_i = \{fieldConstr_j | j = 1, \dots, pathFieldNum\}$

412 **Example.** Each constraints set in our example will contain a set of
413 constraints for any of the fields that are involved in imposed predi-
414 cates of operations on the related path. For example, $pathConstr1$
415 will contain constraints over the fields: $Path1.PTS$, $Path1.RecType$,
416 $Path1.Status$, $Path1.CoNameOrCIK$, $Path1.CompanyID$, $Path1.Name$,
417 $Path1.EffectiveDate$, $Path1.EndDate$, $Path1.ST_ID$, $Path1.ST_name$.
418 Notice that each field is also defined by the related path. Respec-
419 tively, $pathConstr2$ and $pathConstr3$ will contain constraints over
420 the same fields as $pathConstr1$, but with the corresponding path as
421 identifier (e.g., $Path2.PTS$, $Path2.RecType$ and so on for $pathConstr2$
422 and $Path3.PTS$, $Path3.RecType$ and so on for $pathConstr3$). In our
423 example, it does not make any difference maintaining constraints com-
424 ing from fields of $O4$ for $Path1$ (for e.g., $CompanyId$ for $Path1$), since
425 the flow is not split after it merges, but in the general case they are
426 necessary for cases of indirect implications over fields from one path
427 and for determining the number of rows that need to be generated. \square

428 • A field constraint (i.e., $fieldConstr_j$) is defined as a pair of an input
429 field and an ordered list of constraint predicates over this field. For-
430 mally:

431 $fieldConstr_j = [field_j, \mathbb{S}_j]$

432 **Example.** An example field constraint that can be found in our run-
433 ning scenario within $pathConstr1$, is:

434 $fieldConstr_1 = [Path1.RecType, \{(RecType == 'SEC')\}]$ \square

435

436 • Finally, a constraint predicates list defines the logical predicates over

437 the given field in the topological order they are applied over the field
438 in the given path. Formally:
439 $\mathbb{S}_j = \langle P_1(\text{field}_j), \dots, P_{\text{constrNum}}(\text{field}_j) \rangle$
440 The list needs to be ordered to respect the order of operations, since
441 in the general case:
442 $f_1(f_2(\text{field}_x)) \neq f_2(f_1(\text{field}_x))$
443

444 After processing the input ETL graph in Algorithm 1, Algorithm 2 uses
445 the previously generated collection of final path traces (i.e., \mathbb{FP}) for travers-
446 ing a selected *complete* path (i.e., $PT \in \mathbb{FP}$) and constructing a *data genera-*
447 *tion pattern* used finally for generating data that will guarantee its *coverage*.
448 Thus, Algorithm 2 implements the construction of a data generation pattern
449 for *path coverage* of one specific path. For *flow coverage* we can repeat Al-
450 gorithm 2, starting every time with a different PT from the set of final path
451 traces \mathbb{FP} , until each node of the ETL graph has been visited at least once.
452 We should notice here that an alternative to presenting two algorithms —
453 one for path enumeration and one for pattern construction — would be to
454 present a merged algorithm, which traverses the ETL graph and at the same
455 time extracts constraints and constructs the data generation pattern. How-
456 ever, we decided to keep Algorithm 1 separate for two reasons: i) this way
457 the space complexity is reduced while computational complexity remains
458 the same and ii) we believe that the path enumeration algorithm extends
459 beyond the scope of ETL flows and can be reused in a general case for imple-
460 menting a directed path enumeration in polynomial time, while constructing
461 efficient structures for comparison and analysis (i.e., *Path Traces*). A similar
462 approach of using a compact and efficient way to represent ETL workflows
463 using string signatures has been previously introduced in [15].

464 3.3. Path Enumeration Stage

465 In what follows, we present the path enumeration stage, carried out by
466 Algorithm 1.

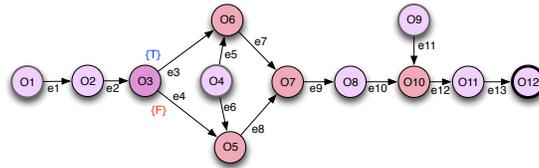
467 In the initial stage of our data generation process, *Bijoux* processes the
468 input ETL process graph in a topological order (Step 2) and for each source
469 node starts a new path trace (Step 5), initialized with the operation rep-
470 resented by a given source node. At the same time, the source node is
471 tagged by the created path trace (Step 6). For other (non-source) nodes,
472 *Bijoux* gathers the path traces from all the previously tagged predecessor
473 nodes (Step 8), extends these path traces with the current operation o_i (Step
474 9), while o_i is tagged with these updated path traces (\mathbb{PT}). Finally, if the

Algorithm 1 Enumerate Paths and Generate Path Traces

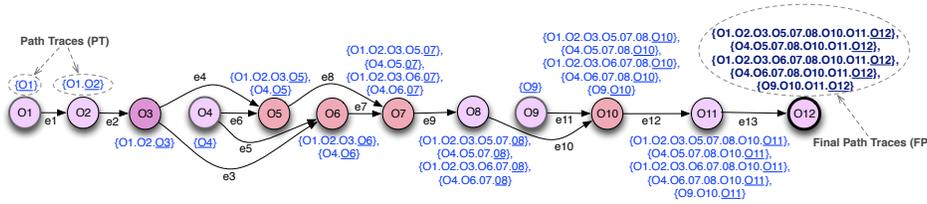
Input: ETL**Output:** FP

```
1: TN  $\leftarrow$  new Tagged Nodes; FP  $\leftarrow$   $\emptyset$ ;  
2: for each operation  $o_i \in \text{TopOrder}(\text{ETL})$  do  
3:   if ( $o_i$  is source) then  
4:     PT  $\leftarrow$   $\emptyset$ ;  
5:     PT.addElement(new Path Trace( $o_i$ ));  
6:     TN.addTag(PT,  $o_i$ );  
7:   else  
8:     PT  $\leftarrow$  TN.UnionOfAll_PTs_forAllPredecessorNodesOf( $o_i$ );  
9:     PT.updateBasedOnOperation( $o_i$ );  
10:    if ( $o_i$  is sink) then  
11:      FP.addAllElementsFrom(PT);  
12:    else  
13:      TN.addTag(PT,  $o_i$ );  
14:    end if  
15:  end if  
16: end for  
17: return FP;
```

475 visited operation is a sink node, the traces of the paths that reach this node
 476 are added to the list of final path traces (i.e., \mathbb{FP}). Processing the input
 477 ETL process graph in this manner, Algorithm 1 gathers the complete set
 478 of final path traces, that potentially can be covered by the generated input
 479 data. An example of the execution of Algorithm 1 applied on our running
 480 example and the 5 resulting final path traces are shown in Figure 4.



(a) DAG representation of our running example



(b) Execution of Algorithm 1 for the topological order of the DAG representation of our running example

Figure 4: Example of execution of Algorithm 1

481 *3.4. Constraints Extraction and Analysis Stage*

482 In what follows, we discuss in detail the constraints extraction and anal-
 483 ysis stages of our data generation process, carried out by Algorithm 2.

484 After all possible *final* paths of input ETL graph are processed and
 485 their traces recorded in \mathbb{FP} , an end-user may select an individual path she
 486 wants to *cover*. To this end, *Bijoux* runs Algorithm 2, with a selected path
 487 $PT \in \mathbb{FP}$, and builds a data generation *Pattern* to cover (at least) the
 488 given path. Algorithm 2 iterates over all the operation nodes of the selected
 489 path (Step 2), and for each *joining node* (i.e., node with multiple incoming
 490 edges), it searches in \mathbb{FP} for all paths that reach the same *joining node*,
 491 from now on, *incident paths* (Steps 5 - 11). As discussed in Section 3.2,
 492 routing operations (e.g., *Router*) introduce such paths, and they need to be
 493 considered separately when generating data for their coverage (see Figure
 494 3). In general, there may be several *joining nodes* on the selected path,
 495 hence Algorithm 2 must take into account all possible combinations of the

Algorithm 2 Construct Data Generation Pattern for one Path

Input: ETL, PT, \mathbb{FP} **Output:** Pattern

```
1:  $\mathbb{AP} \leftarrow \emptyset$ ;  
2: for each operation  $o_i$  crossedBy PT do  
3:   if ( $o_i$  is of type joining_node) then  
4:      $\mathbb{AP}_i \leftarrow \emptyset$   
5:     for each Path Trace  $PT_j \in \mathbb{TN}.getAllPathTracesFor(o_i)$  do  
6:       if ( $PT_j.PredecessorOf(o_i) \neq PT.PredecessorOf(o_i)$ ) then  
7:          $\mathbb{AP}_i.add(PT_j)$ ;  
8:       end if  
9:     end for  
10:     $\mathbb{AP}.add(\mathbb{AP}_i)$ ;  
11:   end if  
12: end for  
13:  $\mathbb{C} \leftarrow allCombinations(PT, \mathbb{AP})$ ;  
14: for each Combination  $C \in \mathbb{C}$  do  
15:   Pattern  $\leftarrow \emptyset$ ;  
16:   for each Path Trace  $PT_i \in C$  do  
17:     for each operation  $o_j$  crossedBy  $PT_i$  do  
18:       Pattern.addConstraints( $o_j$ );  
19:       if ( $\neg$ Pattern.isFeasible) then  
20:         abortPatternSearchForC();  
21:       end if  
22:     end for  
23:   end for  
24:   return Pattern;  
25: end for  
26: return  $\emptyset$ ;
```

496 alternative incident paths that reach these nodes (Step 13).

497 **Example.** Referring to the DAG of Figure 4a, if the path to be covered is
498 $(O9 \rightarrow O10 \rightarrow O11 \rightarrow O12)$, it would require the coverage of additional path(s)
499 because of the equi-join operator $O10$. In other words, data would also need
500 to be coming from edge $e10$ in order to be matched with data from edge $e11$.
501 However, because of the existence of a Union operator ($O7$), there are differ-
502 ent alternative combinations of paths that can meet this requirement. The
503 reason is that data coming from either of the incoming edges of a *Union* oper-
504 ator reach its outgoing edge. Hence, data reaching $O10$ from edge $e10$ could
505 pass through path $(O1 \rightarrow O2 \rightarrow O3 \rightarrow O6 \rightarrow O7 \rightarrow O8 \dots)$ combined with path
506 $(O4 \rightarrow O6 \rightarrow O7 \rightarrow O8 \dots)$ or through path $(O1 \rightarrow O2 \rightarrow O3 \rightarrow O5 \rightarrow O7 \rightarrow O8 \dots)$
507 combined with $(O4 \rightarrow O6 \rightarrow O7 \rightarrow O8 \dots)$. Thus, we see how two alternative
508 combinations of paths, each containing three different paths, can be used
509 for the coverage of one single path. \square

510 For each combination, Algorithm 2 attempts to build a data generation
511 pattern, as explained above. However, some combination of paths may raise
512 a contradiction between the constraints over an input field, which in fact
513 results in disjoint value ranges for this field and thus makes it unfeasible to
514 cover the combination of these paths using a single instance of the input
515 field (Step 20). In such cases, Algorithm 2 aborts pattern creation for a
516 given combination and tries with the next one.

517 **Example.** Referring to the DAG of Figure 4a, we can imagine field $f1$,
518 being present in the schema of operation $O6$ and field $f2$ being present in
519 the schema of operation $O9$. We can also imagine that the datatype of $f1$
520 is *integer* and the datatype of $f2$ is *positive integer*. Then, if the joining
521 condition of operation $O10$ is $(f1 = f2)$ and at the same time, there is
522 a constraint (e.g., in operation $O6$) that $(f1 < 0)$, the algorithm will fail
523 to create a feasible data generation pattern for the combination of paths
524 $(O1 \rightarrow O2 \rightarrow O3 \rightarrow O5 \dots \rightarrow O12)$ and $(O9 \rightarrow O10 \rightarrow O11 \rightarrow O12)$. \square

525 Otherwise, the algorithm updates currently built *Pattern* with the con-
526 straints of the next operation (o_j) found on the path trace.

527 As soon as it finds a combination that does not raise any contradiction
528 and builds a complete feasible *Pattern*, Algorithm 2 finishes and returns the
529 created data generation pattern (Step 24). Notice that by covering at least
530 one combination (i.e., for each *joining node*, each and every incoming edge
531 is crossed by one selected path), Algorithm 2 can guarantee the coverage of
532 the selected input path PT .

533 Importantly, if Algorithm 2 does not find a feasible data generation pat-
534 tern for any of the alternative combinations, it returns an empty pattern
535 (Step 26). This further indicates that the input ETL process model is not

	Datatype	Distribution Type	Modification
PTS	Integer	Triangular	-
RecType	String	Uniform Discrete	deform 2%
CoNameOrCIK	String	Complex	-
CompanyID	Long	Uniform	addNullValues 1%
EffectiveDate	Integer	Uniform	setToMaxInt 1%

Field Parameters (FP)

	Selectivity
O2 (Filter_RecType)	0.3
O3 (Router_1)	0.7
O6 (Join_1)	1
O5 (Join_2)	0.95
O8 (Filter_Date)	0.6

Operation Parameters (OP)

Figure 5: Data generation parameters (FP and OP)

536 correct, i.e., that some of the path branches are not reachable for any com-
 537 bination of input data.

538 The above description has covered the general case of data generation
 539 without considering other generation parameters. However, given that our
 540 data generator aims at generating data to satisfy other configurable param-
 541 eters, we illustrate here as an example the adaptability of our algorithm to
 542 the problem of generating data to additionally satisfy *operation selectivity*.
 543 To this end, the algorithm now also analyzes the parameters at the oper-
 544 ation level (OP) (see Figure 5:right). Notice that such parameters can be
 545 either obtained by analyzing the input ETL process for a set of previous
 546 real executions, or simply provided by the user, for example, for analyzing
 547 the flow for a specific set of operation selectivities.

548 Selectivity of an operation o expresses the ratio of the size of the dataset
 549 at the output (i.e., $card(o)$), to the size at the input of an operation (i.e.,
 550 $input(o)$). Intuitively, for filtering operations, we express selectivity as the
 551 percentage of data satisfying the filtering predicate (i.e., $sel(o) = \frac{card(o)}{input(o)}$),
 552 while for n-ary (join) operations, for each input e_i , we express it as the
 553 percentage of the data coming from this input that will match with other
 554 inputs of an operation (i.e., $sel(o, e_i) = \frac{card(o)}{input(o, e_i)}$).

555 From the OP (see Figure 5:right), *Bijoux* finds that operation *O2* (Fil-
 556 ter_RecType) has a selectivity of 0.3 . While processing a selected path
 557 starting from the operation *O1*, *Bijoux* extracts operation semantics for *O2*
 558 and finds that it uses the field *RecType* ($RecType == 'SEC'$). With the selec-
 559 tivity factor of 0.3 from OP, *Bijoux* infers that out of all incoming rows for
 560 the Filter, 30% should satisfy the constraint that *RecType* should be equal
 561 to SEC, while 70% should not. We analyze the selectivity as follows:

- 562 • To determine the total number of incoming rows for operation *O8*
 563 (Filter_Date), we consider predecessor operations, which in our case
 564 come from multiple paths.

- 565 • As mentioned above, operation O_2 will allow only 30% of incoming
566 rows to pass. Assuming that the input load size from *FINWIRE* is
567 1000, this means that in total $0.3 * 1000 = 300$ rows pass the filter
568 condition.
- 569 • From these 300 rows only 70%, based on the O_3 (*Router_1*) selec-
570 tivity, (i.e., 210 rows) will successfully pass both the filtering (*Rec-*
571 *Type*==‘*SEC*’) and the router condition (*isNumber(CoNameOrCIK)*)
572 and hence will be routed to the route that evaluates to *true*. The rest
573 ((i.e., $300 - 210 = 90$ rows)) will be routed to the route that evaluates
574 to *false*.
- 575 • The 210 rows that pass both previous conditions, will be matched
576 with rows coming from operation O_4 through the join operation O_6
577 (*Join_1*). Since the selectivity of operation O_6 is 1, all 210 tuples will
578 be matched with tuples coming from O_4 and meeting the condition
579 *CoNameOrCIK*==*CompanyID* and hence will pass the join condition.
580 On the other hand, the selectivity of operation O_5 (*Join_2*), for the
581 input coming from O_3 (*Router_1*), is 0.95, which means that from the
582 90 rows that evaluated to *false* for the routing condition, only 85 will
583 be matched with tuples coming from O_4 and meeting the condition
584 *CoNameOrCIK*==*Name*. Thus, $210 + 85 = 295$ tuples will reach the
585 union operation O_6 and pass it.
- 586 • Finally, from the 295 rows that will reach operation O_8 (*Filter_Date*)
587 coming from the preceding union operation, only $0.6 * 295 = 177$ will
588 successfully pass the condition (*PTS*>=*EffectiveDate*) AND (*PTS*<=
589 *EndDate*), as the selectivity of OP_8 is 0.6.

590 In order to generate the data that do not pass a specific operation of the
591 flow, a data generate pattern inverse to the initially generated *Pattern*
592 in Algorithm 2 needs to be created to guarantee the percentage of data
593 that will fail the given predicate.

594 Similarly, other parameters can be set for the generated input data to
595 evaluate different quality characteristics of the flow, (see Figure 5:left). As
596 an example, the percentage of null values or incorrect values (e.g., wrong
597 size of telephone numbers or negative age) can be set for the input data,
598 to evaluate the measured data quality of the flow output, regarding *data*
599 *completeness* and *data accuracy*, respectively. Other quality characteristics
600 like *reliability* and *recoverability* can be examined as well, by adjusting the
601 distribution of input data that result to exceptions and the selectivity of

602 exception handling operations. Examples of the above will be presented in
 603 Section 4.

604 3.5. Data Generation Stage

605 Lastly, after the previous stage builds data generation patterns for cov-
 606 ering either a single path, combination of paths, or a complete flow, the last
 607 (data generation) stage proceeds with generating data for each input field
 608 f . Data are generated within the ranges (i.e., R) defined by the constraints
 609 of the provided pattern, using either random numerical values within the
 610 interval or dictionaries for selecting correct values for other (textual) fields.

611 For each field f , data generation starts from the complete domain of the
 612 field's datatype $dt(f)$.

Each constraint P , when applied over the an input field f , generates a
 set of disjoint ranges of values $R_i^{f,init}$ in which the data should be gener-
 ated, and each range being inside the domain of the field's datatype $dt(f)$.
 Formally:

$$P(f) = R^{f,init} = \left\{ r^{f,init} \mid r^{f,init} \subseteq dt(f) \right\} \quad (1)$$

613 For example, depending on the field's datatype, a value range for numeric
 614 datatypes is an interval of values (i.e., $[x, y]$), while for other (textual) fields
 615 it is a set of possible values a field can take (e.g., *personal names*, *geographical*
 616 *names*).

617 After applying the first constraint P_1 , *Bijoux* generates a set of disjoint,
 618 non-empty value ranges R_1^f , each range being an intersection with the do-
 619 main of the field's datatype.

$$R_1^f = \left\{ r_1^f \mid \forall r_1^{f,init} \in R_1^{f,init}, \exists r_1^f, s.t. : \right. \\ \left. (r_1^f = r_1^{f,init} \cap dt(f) \wedge r_1^f \neq \emptyset) \right\} \quad (2)$$

620 Iteratively, the data generation stage proceeds through all the constraints
 621 of the generation pattern. For each constraint P_i it updates the resulting
 622 value ranges as an intersection with the ranges produced in the previous
 623 step, and produces a new set of ranges R_i^f .

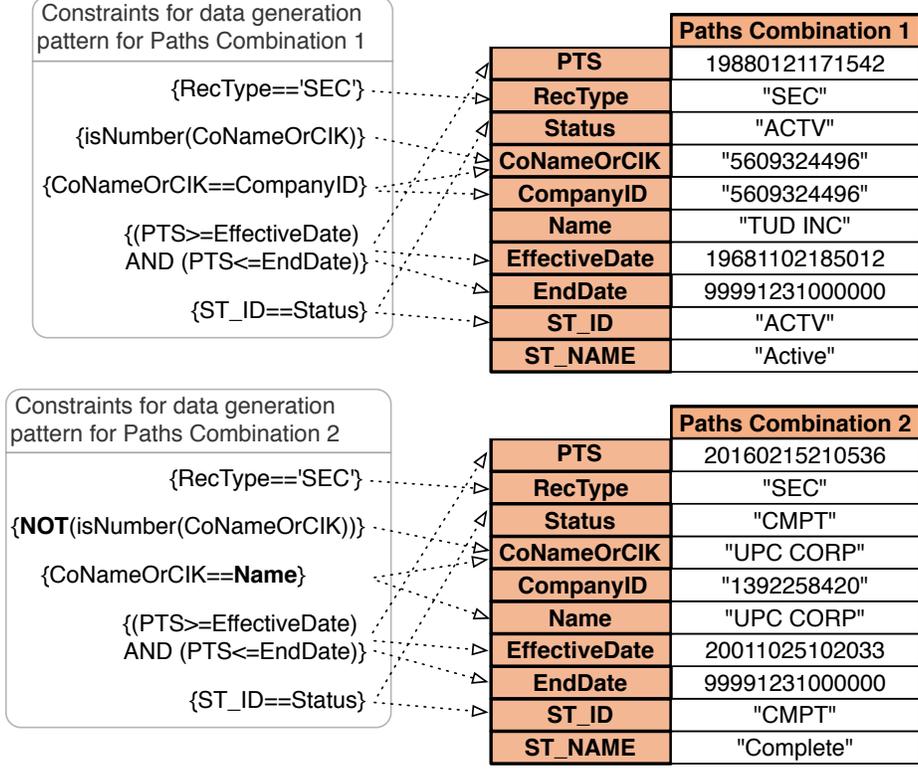


Figure 6: Data generated after analyzing all ETL operations

$$R_i^f = \left\{ r_i^f \mid \forall r_i^{f,init} \in R_i^{f,init}, \forall r_{i-1}^f \in R_{i-1}^f, \exists r_i^f, s.t. : \right. \quad (3)$$

$$\left. (r_i^f = r_i^{f,init} \cap r_{i-1}^f \wedge r_i^f \neq \emptyset) \right\}$$

624 Finally, following the above formalization, for each input field f *Bi-*
625 *joux* produces a final set of disjoint, non-empty value ranges ($R_i^{f,final}$) and
626 for each range it generates an instance of f inside that interval.

627 See for example, in Figure 6 and Figure 7, the generated data sets for
628 covering the ETL process flow of our running example. We should mention
629 at this point, that non conflicting constraints for the same field that is
630 present in different paths and/or path combinations, can be merged and
631 determine a single range (i.e., the intersection of all the ranges resulting

FINWIRE				StatusType	
PTS	RecType	Status	CoNameOrCIK	ST_ID	ST_NAME
19880121171542	"SEC"	"ACTV"	"5609324496"	"ACTV"	"Active"
20160215210536	"SEC"	"CMPT"	"UPC CORP"	"CMPT"	"Complete"

DimCompany			
CompanyID	Name	EffectiveDate	EndDate
"5609324496"	"TUD INC"	19681102185012	99991231000000
"1392258420"	"UPC CORP"	20011025102033	99991231000000

Figure 7: Generated datasets corresponding to the generated data

632 from the different paths). This way, under some conditions, the same value
633 within that interval can be used for the coverage of different paths. As
634 an example, in Figure 6, the fields *Status* and *ST_ID* that exist in both
635 path combinations, all hold a constraint ($ST_ID == Status$). These can be
636 merged into one single constraint, allowing for the generation of only one
637 row for the table *StatusType* that can be used for the coverage of both path
638 combinations, as long as both generated values for the field *Status* equal the
639 generated value for the field *ST_ID* (e.g., "ACTV").

640 Following this idea, it can easily be shown that under specific conditions,
641 the resulting constraints for the different path combinations from the appli-
642 cation of our algorithm, can be further reduced, until they can produce a
643 minimal set of datasets for the coverage of the ETL flow.

644 Data generation patterns must be further combined with other user-
645 defined data generation parameters (e.g., selectivities, value distribution,
646 etc.). We provide more details regarding this within our test case in Section
647 4.

648 3.6. Theoretical validation

649 We further provide a theoretical validation of our data generation pro-
650 cess in terms of: the *correctness* of generated data sets (i.e., *path and flow*
651 *coverage*).

652 A theoretical proof of the correctness of the *Bijoux* data generation pro-
653 cess is divided into the three following components.

1. *Completeness of path traces*. Following from Algorithm 1, for each ETL graph node (i.e., datastores and operations, see Section 2.2) *Bijoux* builds path traces of all the paths reaching that node (e.g., see

Figure 4b). Formally, given that an ETL graph node can represent either an operation (\mathbf{O}), a source (\mathbf{DS}_S), or a target data store (\mathbf{DS}_T), we recursively formalize the existence of path traces as follows:

$$\forall v_i \in \mathbf{O} \cup \mathbf{DS}_T, \mathbb{P}\mathbb{T}_{v_i} = \bigcup_{j=1}^{|\{v_j | v_j \prec v_i\}|} \left\{ PT_{v_j}^1 \cdot v_i, \dots, PT_{v_j}^{|\mathbb{P}\mathbb{T}_{v_j}|} \cdot v_i \right\}. \quad (4)$$

$$\forall v_i \in \mathbf{DS}_S, \mathbb{P}\mathbb{T}_{v_i} = \{PT_{v_i}\}, PT_{v_i} = v_i. \quad (5)$$

654 Considering that ETL graph nodes are visited in a topological order
 655 (see Step 2 in Algorithm 1), the path traces of each ETL graph node
 656 are built after visiting all its predeceasing sub-paths. This guarantees
 657 that path traces of each node v_i are complete with regard to all its
 658 predecessors (i.e., $\{v_j | v_j \prec v_i\}$), hence the final path traces $\mathbb{P}\mathbb{T}$ (i.e.,
 659 path traces of target data store nodes) are also complete.

660 2. *Path coverage.* Having the complete path traces recorded in Algorithm
 661 1, Algorithm 2 traverses a selected path (i.e., PT), with all its alter-
 662 native *incidence paths*, and builds a data generation *Patern* including
 663 a list of constraints over the input fields. Following from 1, this list
 664 of constraints is complete. Moreover, as explained in Section 3.5, *Bi-*
 665 *joux* iteratively applies given constraints, and for each input field f
 666 produces a set of value ranges ($R^{f,final}$), within which the field values
 667 should be generated.

668 Given the statements 1 - 3 in Section 3.5, *Bijoux* guarantees that the
 669 data generation stage applies all the constraints over the input fields
 670 when generating $R^{f,final}$, thus guaranteeing that the complete selected
 671 path will be covered.

672 On the other side, if at any step of the data generation stage a result of
 673 applying a new constraint P_i leads to an empty set of value ranges, the
 674 collected list of constraints must be contradictory. Formally (following
 675 from statement 3 in Section 3.5):

$$676 (\exists R_i^{f,init}, R_{i-1}^f | R_i^f = \emptyset) \rightarrow \perp.$$

677 This further implies that the input ETL graph has contradictory path
 678 constraints that would lead to an unreachable sub-path, which could
 679 never be executed. As an additional functionality, *Bijoux* detects such
 680 behavior and accordingly warns the user that the input ETL flow is
 681 not correct.

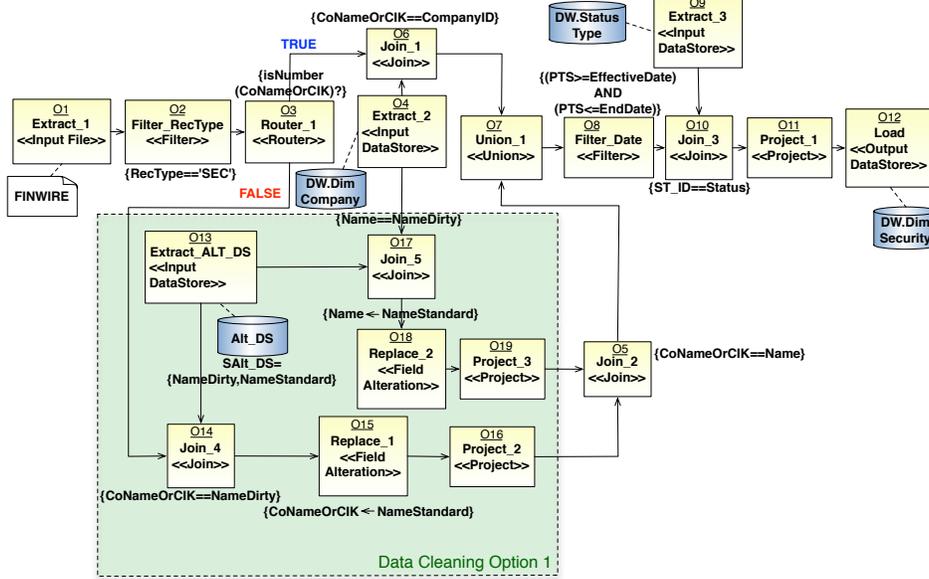


Figure 8: ETL flow for data cleaning, using a dictionary

682 3. *Flow coverage.* Following from 2, Algorithm 2 generates data that
 683 guarantee the coverage of a single path from \mathbb{FP} . In addition, if Algo-
 684 rithm 2 is executed for each final path $PT_i \in \mathbb{FP}$, it is straightforward
 685 that *Bijoux* will produce data that guarantee the coverage of the complete
 686 ETL flow (i.e., ETL graph), unless a constraints contradiction
 687 for an individual path has been detected.

688 4. Test case

689 The running example of the ETL flow that we have used so far is expres-
 690 sive enough to illustrate the functionality of our framework, but it appears
 691 too simple to showcase the benefits of our approach regarding the evalua-
 692 tion of the quality of the flow. In this respect, we present in this section
 693 representative examples of how our framework can generate data, not only
 694 to enact specific parts of the ETL flow, but also to evaluate the performance
 695 and the data quality of these flow parts.

696 Going back to our running example (Figure 1), from now on referred
 697 to as *Flow_A*, we can identify a part of the flow that can be the source
 698 of data quality issues. That is, rows whose values for the field *CoName-*
 699 *OrCIK* are not numbers are matched with data about companies from the

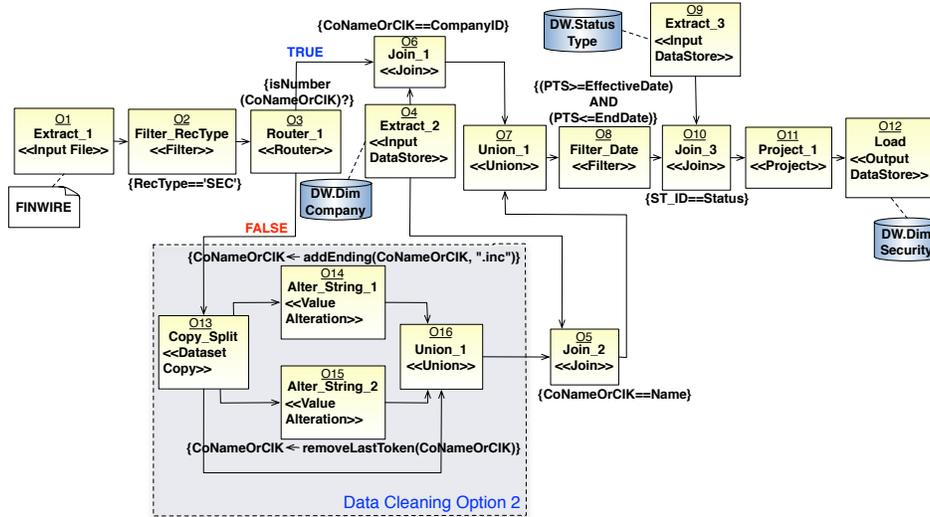


Figure 9: ETL flow for data cleaning, trying different string variations for the join key

700 *DimCompany* table, through an equi-join on the company name (CoName-
 701 OrCIK==Name). However, company names are typical cases of attributes
 702 that can take multiple values in different systems or even within the same
 703 system. For example, for a company *Abcd Efgh*, its name might be stored
 704 as “Abcd Efgh”, or followed by a word indicating its type of business en-
 705 tity (e.g., “Abcd Efgh Incorporated”) or its abbreviation with or without
 706 a comma (e.g., “Abcd Efgh Inc.” or “Abcd Efgh, Inc.”). It is also possi-
 707 ble that it might be stored using its acronym (e.g., “ABEF”) or with a
 708 different reordering of the words in its name, especially when the two first
 709 words are name and surname of a person (e.g., “Efgh Abcd”). Moreover,
 710 there can be different uppercase and lowercase variations of the same string,
 711 combinations of the above-mentioned variations or even misspelled values.

712 Hence, there are many cases that the equi-join (CoNameOrCIK==Name)
 713 will fail to match the incoming data from the *FINWIRE* source with the
 714 rows from the *DimCompany* table, because they might simply be using a
 715 different variation of the company name value. This will have an impact
 716 on *data completeness*, since it will result in fewer rows being output to the
 717 *DimSecurity* than there should be.

718 To this end, we introduce here two more complex ETL flows (Figure 8
 719 and Figure 9), which perform the same task as the running example, but
 720 include additional operations in order to improve the data quality of the out-

721 put data. The ETL flow in Figure 8, from now on referred to as *Flow_B*,
722 uses a dictionary (*Alt_DS*) as an alternative data source. This dictionary is
723 assumed to have a very simple schema of two fields — *NameDirty* and *Name-*
724 *Standard*, to maintain a correspondence between different dirty variations
725 of a company name and its standard name. For simplicity, we assume that
726 for each company name, there is also one row in the dictionary containing
727 the standard name, both as value for the *NameDirty* and the *NameStandard*
728 fields. Operations *O14* and *O17* are used to match both the company names
729 from the *FINWIRE* and the table, to the corresponding dictionary entries
730 and subsequently, rows are matched with the standard name value being the
731 join key, since the values for the join keys are replaced by the standard name
732 values (*(Name←NameStandard)* and *(CoNameOrCIK←NameStandard)*).

733 Another alternative option for data cleaning is to try different variations
734 of the company name value, by adding to the flow various string operations
735 that alter the value of *CoNameOrCIK*. The ETL flow in Figure 9, from
736 now on referred to as *Flow_C*, generates different variations of the value
737 for *CoNameOrCIK* with operations *O14* and *O15*, who concatenate the
738 abbreviation “inc.” at the end of the word and remove the last token of
739 the string, respectively. After the rows from these operations are merged
740 through a Union operation (*O16*), together with the original *CoNameOrCIK*
741 value, all these different variations are tried out to match with rows coming
742 from *DimCompany*.

743 4.1. Evaluating the performance overhead of alternative ETL flows

744 In the first set of experiments, we implemented the three different ETL
745 flows (*Flow_A*, *Flow_B* and *Flow_C*) using Pentaho Data Integration⁴ and
746 we measured their time performance by executing them on Kettle Engine,
747 running on Mac OS X, 1.7 GHz Intel Core i5, 4GB DDR3 and keeping
748 average values from 10 executions.

749 For each flow, we used *Bijoux* to generate data to cover only the part of
750 the flow that was of interest, i.e., to cover the paths from Operations *O1* to
751 *O12* who are covered by the rows that are evaluated as False by operation
752 *O3*. Hence, one important advantage of our tool is that it can generate data
753 to evaluate specific part of the flow, as opposed to random data generators
754 (e.g., the TPC-DI data generator provided on the official website) who can
755 only generate data agnostically of which part of the flow is being covered.
756 This gives *Bijoux* not only a quality advantage, being able to evaluate the

⁴<http://www.pentaho.com/product/data-integration>

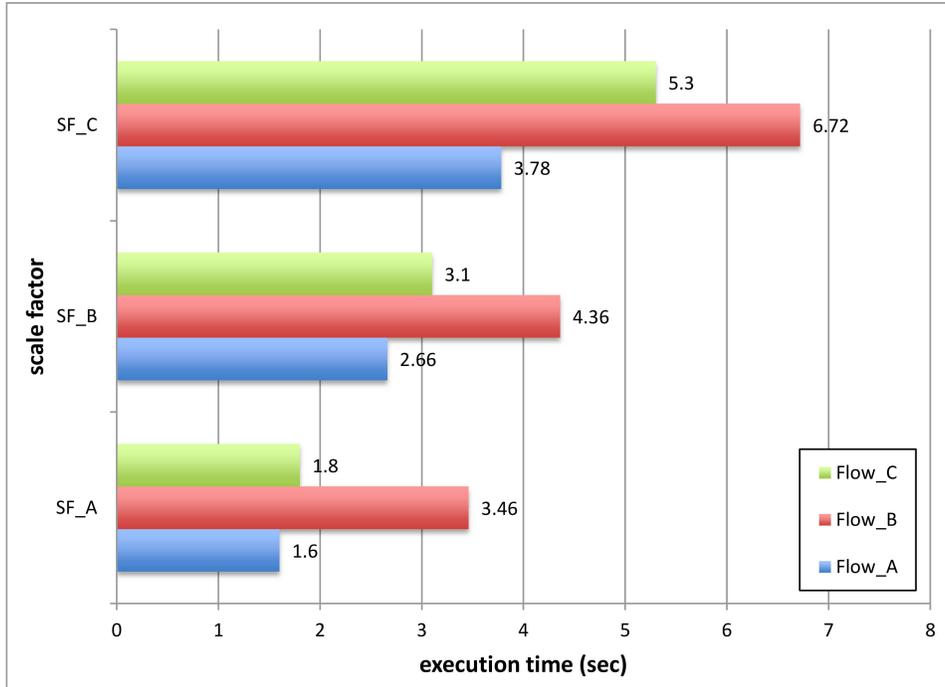


Figure 10: Performance evaluation of the flows using different scale factors

757 flow in greater granularity, but also a practical advantage, since the size of
 758 data that need to be generated can be significantly smaller. For instance, the
 759 TPC-DI data generator generates data for the FINWIRE file, only around
 760 $\frac{1}{3}$ of which are evaluated as *true* by the filter *RecType==SEC* and from
 761 them only around $\frac{1}{3}$ contains a company name instead of a number.

762 In order to generate realistic values for the company name fields, we used
 763 a catalog of company names that we found online ⁵ and we used *Bijoux* to
 764 generate data not only for the attributes that have been mentioned above,
 765 but for all of the attributes of the schemata of the involved data sources as
 766 defined in the TPC-DI documentation, so as to measure more accurate time
 767 results.

768 For each flow, we generated data of different size in order to evaluate
 769 how their performance can scale with respect to input data size, as shown in
 770 the below table, where we can see the number of rows for each data source
 771 for the three different scale factors (*SF*).

⁵<https://www.sec.gov/rules/other/4-460list.htm>

Data source →	FINWIRE	DimCompany	Alt_DS (for <i>Flow_B</i>)
SF_A	4000	4000	60000
SF_B	8000	8000	60000
SF_C	16000	16000	60000

772 For these experiments, for each flow we assumed selectivities that would
773 guarantee the matching of all the rows in *FINWIRE* with rows in *DimCom-*
774 *pany* and the results can be seen in Figure 10 For *Flow_C*.

775 As we expected, the results show an overhead in performance imposed by
776 the data cleaning operations. It was also intuitive to expect that the lookup
777 in the dictionary (*Flow_B*) would impose greater overhead than the string
778 alterations (*Flow_C*). Nevertheless, some interesting finding that was not
779 obvious is that as input data scale in size, the overhed of *Flow_B* appears
780 to come closer and closer to the overhed of (*Flow_C*), which appears to
781 become greater as input data size grows. We should notice at this point
782 that our results regard the performance and scalability of a specific part of
783 the flow – not the complete flow in general – which is a unique advantage
784 of our approach, especially in cases of dealing with bottlenecks.

785 Consequently, we conducted experiments assuming different levels of in-
786 put data dirtiness, by setting the selectivity of the different join operations
787 for the different flows. The scenario we intended to simulate was a pre-
788 defined percentage of different types of data dirtiness. In this respect, we
789 considered four different types of dirtiness:

- 790 1. Missing the abbreviation “inc.” at the end of the company name
791 (Type_I)
- 792 2. A word (e.g., company type abbreviation) exists at the end of the
793 name when it should not (Type_II)
- 794 3. The ending of the company name is mistakenly in an extended format
795 (e.g., “incorporated’ ’ instead of “inc.”) (Type_III)
- 796 4. Miscellaneous that cannot be predicted (e.g., “corp.” instead of “inc.”
797 or misspelled names) (Type_IV)

798 We assumed that *Flow_A* cannot handle any of these cases (i.e., dirty
799 names as an input for the *FINWIRE* source will fail to be matched to data
800 coming from *DimCompany*); that *Flow_B* can solve all the cases for Type_I

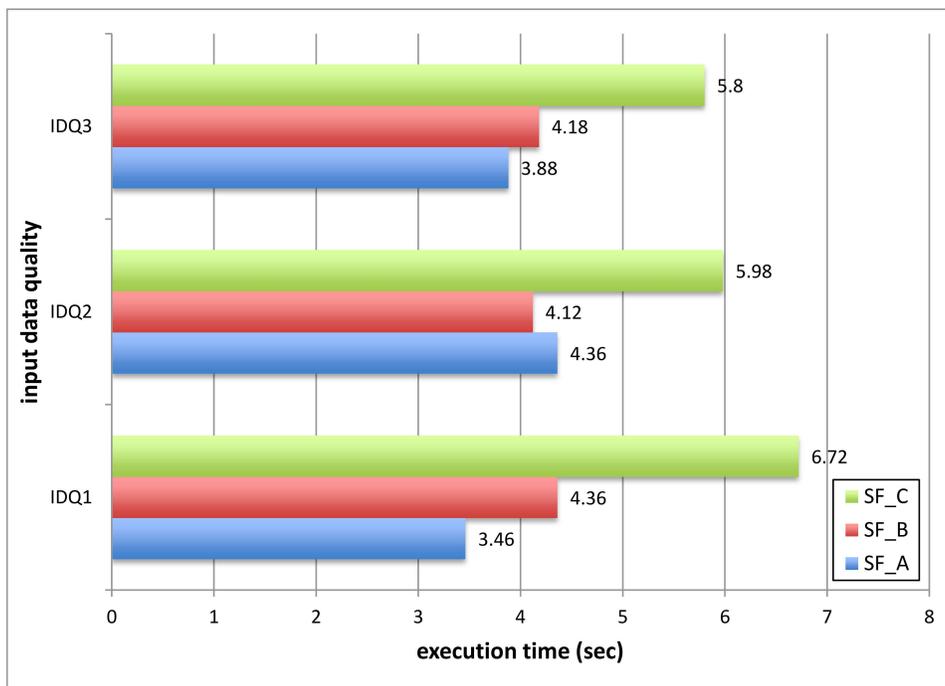


Figure 11: Performance evaluation of *Flow_B* using different levels of input data quality

801 and *Type_III* (i.e., there will be entries in the dictionary covering both of
 802 these types of dirtiness); and *Flow_C* can cover all the cases for *Type_I*
 803 and *Type_II*, because of the operation that it performs.

804 Thus, we generated data that were using real company names from the
 805 online catalog; we considered those names as the standard company names
 806 versions to generate data for the *DimCompany* source; and we indirectly
 807 introduced specified percentages of the different types of dirtiness, by set-
 808 ting a) the selectivities of the join operators and b) by manually generating
 809 entries in our dictionary (*Alt_DS*) that included all the names from the
 810 catalog together with their corresponding names manually transformed to
 811 *Type_I* and *Type_II*. The percentages of input data quality (IDQ) that
 812 were used for our experiments can be seen in the following table.

813 In Figure 11, we show how the performance of *Flow_B* scales with re-
 814 spect to different scale factors and data quality of input data. What is
 815 interesting about those results, is that the flow appears to be performing
 816 better when the levels of dirtiness of the input data are higher. This might
 817 appear counter-intuitive, but a possible explanation could be that less data

Dirtness Type →	Type_I	Type_II	Type_III	Type_IV
IDQ1	0%	0%	0%	0%
IDQ2	1%	1%	3%	1%
IDQ3	2%	2%	6%	2%

818 (i.e., fewer rows) actually reach the extraction operation, keeping in mind
819 that read/write operations are very costly for ETL flows.

820 4.2. Evaluating the data quality of alternative ETL flows

821 In the above-mentioned experiments, we evaluated the time performance
822 of different flows, assuming that both data quality levels and data dirtiness
823 characterization were a given. However, in order to evaluate an ETL flow
824 with respect to the quality of the data cleaning that it can provide, it is not
825 sufficient to only evaluate the time performance of different data cleaning
826 options. To this end, in the second set of experiments, our goal was to
827 evaluate which data cleaning option would produce the lowest levels of data
828 incompleteness in the output data of the flow (*DimSecurity* table), using
829 realistic datasets. In this respect, we used the company names from our
830 catalog and for each of them we prepared a query to scrap the Freebase online
831 database⁶ and retrieve data about the company name and the known aliases
832 of those names. Consequently, starting from 940 unique company names of
833 our catalog, we were able to construct a dictionary that contained 2520
834 entries, each containing an alias of a company name and its corresponding
835 standard name. We then used this dictionary as our *Alt_DS* dictionary; the
836 standard names to populate the *DimCompany* table; and the names as they
837 were on the catalog to populate the *FINWIRE* file.

838 Using *Bijoux*, we generated data that used *Flow_A* semantics in order to
839 pass through the part of the flow that was of our interest and the dictionaries
840 as mentioned above to generate realistic data. Despite the fact that it might
841 appear as if the use of dictionaries devalues the use of our algorithm, in fact
842 this is one strength of our approach — that it can be configured to generate
843 data with different degrees of freedom, based on the constraints defined both
844 by the flow semantics and the user. Therefore, it is possible to conduct such
845 analysis, using a hybrid approach and evaluating the flows based on realistic

⁶<https://www.freebase.com/>

846 data. The contribution of our algorithm in this case is to generate, on one
847 hand all the data for the different fields of the schemata that are required for
848 the flow execution and to make sure, on the other hand that the generated
849 rows will cover specific parts of the flow.

850 After executing Flow_B and Flow_C with these input data, we used
851 the following measure for data completeness:

$$852 \quad DI = \%_of_missing_entities_from_their_appropriate_storage \text{ [16]}$$

853 The results for the two flows were the following:

$$854 \quad DI_{Flow_B} = \frac{56}{940} * 100 \approx 6\%$$

855

$$856 \quad DI_{Flow_C} = \frac{726}{940} * 100 \approx 77\%$$

857

858 According to these results, we can see a clear advantage of Flow_B
859 regarding the data quality that it provides, suggesting that the performance
860 overhead that it introduces, combined with potential cost of obtaining and
861 maintaining a dictionary, might be worth undertaking, if data completeness
862 is a goal of high priority.

863 We have explained above how the parametrization of our input data
864 generation enables the evaluation of an ETL process and various design al-
865 terations over it, with respect to data quality and performance. Essentially,
866 alternative implementations for the same ETL can be simulated using dif-
867 ferent variations of the data generation properties and the measured quality
868 characteristics will indicate the best models, as well as how they can scale
869 with respect not only to data size but also to data quality of the input data.
870 Similarly, other quality characteristics can be considered, like *reliability* and
871 *recoverability*, by adjusting the percentage of input data that result to excep-
872 tions and the selectivity of exception handling operations. In addition, we
873 have shown through our examples how data properties in the input sources
874 can guide the selection between alternative ETL flows during design time.

875 5. *Bijoux* performance evaluation

876 In this section, we report the experimental findings, after scrutinizing
877 different performance parameters of *Bijoux*, by using the prototype that
878 implements its functionalities.

879 We first introduce the architecture of a prototype system that imple-
880 ments the functionality of the *Bijoux* algorithm.

881 **Input.** The main input of the *Bijoux* framework is an ETL process.
882 As we previously discussed, we consider that ETL processes are provided

883 in the logical (platform-independent) form, following previously defined for-
 884 malization (see Section 2.2). Users can also provide various parameters (see
 885 Figure 5) that can lead the process of data generation, which can refer to
 886 specific fields (e.g., *field distribution*), operations (e.g., *operation selectivity*)
 887 or general data generation parameters (e.g., *scale factors*).

888 **Output.** The output of our framework is the collection of datasets
 889 generated for each input data store of the ETL process. These datasets
 890 are generated to satisfy the constraints extracted from the flow, as well
 891 as the parameters provided by the users for the process description (i.e.,
 892 distribution, operation selectivity, load size).

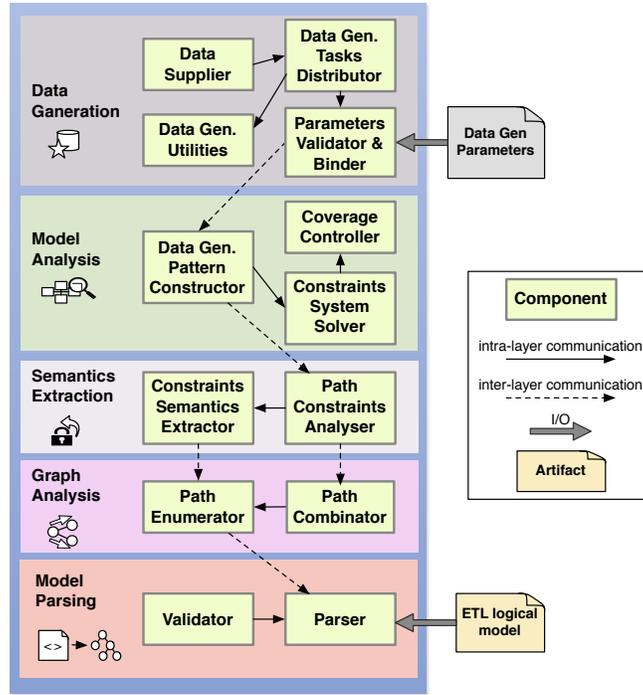


Figure 12: Bijoux prototype architecture

893 ***Bijoux's* architecture.** The *Bijoux's* prototype is modular and based
 894 on a layered architecture, as shown in Figure 12. The four main layers
 895 implement the core functionality of the *Bijoux* algorithm (i.e., *graph anal-*
 896 *ysis*, *semantics extraction*, *model analysis*, and *data generation*), while the
 897 additional bottom layer is responsible for importing ETL flows from corre-
 898 sponding files and can be externally provided and plugged to our framework
 899 (e.g., flow import plugin [13]). We further discuss all the layers in more

900 detail.

- 901 • The bottom layer (*Model Parsing*) of the framework is responsible for
902 parsing the model of the ETL process (*Parser* component) from the
903 given logical representation of the flow (e.g., XML), and importing a
904 DAG representation for the process inside the framework. In general,
905 the *Model Parsing* layer can be extended with external parser plugins
906 for handling different logical representations of an ETL process (e.g.,
907 [12, 13]). This layer also includes a *Validator* component to ensure
908 syntactic, schematic and logical (e.g., cycle detection) correctness of
909 the imported models.
- 910 • The *Graph Analysis* layer analyses the DAG representation of the ETL
911 flow model. Thus, it is responsible for identifying and modeling all the
912 ETL flow paths (*Path Enumerator* component; see Algorithm 1), as
913 well as constructing all their possible combinations (*Path Combinator*
914 component).
- 915 • The *Semantics Extraction* layer extracts relevant information needed
916 to process the ETL flow. The information extracted in this layer (from
917 the *Constraints Semantics Extractor* component) includes information
918 about input datasets, operation semantics, order of operations,
919 schema changes, and other parameters for data generation. This layer
920 is also responsible for modeling constraints grouped by path (*Path*
921 *Constraints Analyzer*; see Algorithm 2) to provide the required con-
922 structs for feasibility analysis and the construction of a data generation
923 pattern to the layer above (*Model Analysis*).
- 924 • *Model Analysis* layer realizes the construction of a data generation
925 pattern (*Data Gen. Pattern Constructor* component) that computes
926 for each field (i.e., attribute), in each table, the ranges of values ac-
927 cording to the extracted semantics of operations and their positioning
928 within paths and path combinations. To this end, this layer includes
929 the *Coverage Controller* component for implementing such analysis
930 according to the set coverage goal (i.e., path coverage, flow cover-
931 age). In addition, it includes the *Constraints System Solver* compo-
932 nent, which solves the systems of gathered constraints (e.g., system of
933 logical predicates and equations over specified attributes) and returns
934 the computed restrictions over the ranges.
- 935 • *Data Generation* layer controls the data generation stage according to
936 the constraints (i.e., data generation patterns) extracted and analyzed

937 in the previous layer, as well as the Data Gen. Parameters provided
938 externally (e.g., distribution, selectivity). The *Parameters Validator*
939 *& Binder* component binds the externally provided parameters to the
940 ETL model and ensures their compliance with the data generation pat-
941 terns, if it is possible. The *Data Gen. Tasks Distributor* component
942 is responsible for managing the generation of data in a distributed
943 fashion, where different threads can handle the data generation for
944 different (pairs of) attributes, taking as input the computed ranges
945 and properties (e.g., *generate 1000 values of normally distributed in-*
946 *tegers where 80% of them are lower than “10”*). For that purpose, it
947 utilizes the *Data Gen. Utilities* component, that exploits dictionaries
948 and random number generation methods. Finally, the *Data Supplier*
949 component outputs generated data in the form of files (e.g., CSV files).

950 5.1. Experimental setup

951 Here, we focused on testing both the functionality and correctness of
952 the *Bijoux* algorithm discussed in Section 3, and different quality aspects,
953 i.e., data generation overhead (*performance*) wrt. the growing complexity of
954 the ETL model. The reason that we do not additionally test those quality
955 aspects wrt. input load sizes is that such analysis is irrelevant according to
956 the *Bijoux* algorithm. The output of the analysis phase is a set of ranges
957 and data generation parameters for each attribute. Hence, the actual data
958 generation phase does not depend on the efficiency of the proposed algo-
959 rithm, but instead can be realized in an obvious and distributed fashion.
960 Thus, we present our results from experiments that span across the phases
961 of the algorithm up until the generation of ranges for each attribute. We
962 performed the performance testing considering several ETL test cases, which
963 we describe in what follows.

964 Our experiments were carried under an OS X 64-bit machine, Processor
965 Intel Core i5, 1.7 GHz and 4GB of DDR3 RAM. The test cases consider a
966 subset of ETL operations, i.e., *Input DataStore*, *Join*, *Filter*, *Router*, *UDF*,
967 *Aggregation* and *Output DataStore*. Based on the TPC-H benchmark⁷, our
968 basic scenario is an ETL process, which extracts data from a source re-
969 lational database (TPC-H DB) and after processing, loads data to a data
970 warehouse (DW) and can be described by the following query: *Load in the*
971 *DW all the suppliers in Europe together with their information (phones, ad-*
972 *dresses etc.), sorted on their revenue and separated by their account balance*

⁷<http://www.tpc.org/tpch/>

973 (either low or high), as can be seen in Fig. 13.

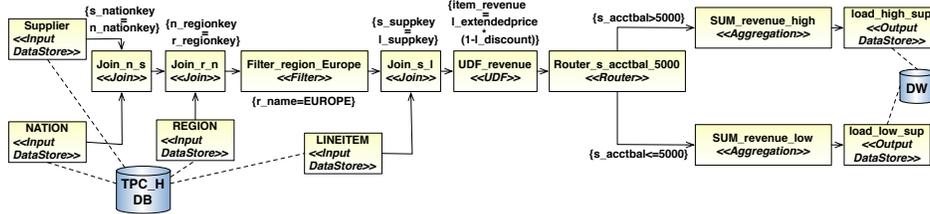


Figure 13: Basic scenario ETL process for experiments

974 The tables that are used from the source database are *Supplier*, *Nation*,
 975 *Region* and *Lineitem*. After *Supplier* entries have been filtered to keep
 976 only suppliers in Europe, the revenue for each supplier is calculated based
 977 on the supplied lineitems and subsequently, they are sorted on revenue,
 978 separated by their account balance and loaded to different tables in the
 979 DW. Starting from the basic scenario, we use *POIESIS* [17], a tool for ETL
 980 Process redesign that allows for the automatic addition of flow patterns on
 981 an ETL model. Thus, we create other, more complex, synthetic ETL flows.
 982 The motivation for using tools for automatic ETL flow generation stems from
 983 the fact that obtaining real world ETL flows covering different scenarios with
 984 different complexity and load sizes is hard and often impossible.

985 **Scenarios creation.** Starting from this basic scenario, we create more
 986 complex ETL flows by adding additional operations, i.e., *Join*, *Filter*, *Input*
 987 *DataStore*, *Project* in various (random) positions on the original flow. We
 988 add two different Flow Component Patterns (FCP) [17] on the initial ETL
 989 flow in different cardinalities and combinations. The first pattern — *Join* —
 990 adds 3 operations every time it is applied on a flow: one Input DataStore, one
 991 Join and one Project operation in order to guarantee matching schemata;
 992 the second pattern — *Filter* — adds one Filter operation with a random
 993 (inequality) condition on a random numerical field (i.e., attribute).

994 We iteratively create 5 cases of different ETL flow complexities and ob-
 995 serve the *Bijoux*'s execution time for these cases, starting from the basic
 996 ETL flow:

- 997 • *Case 1.* Basic ETL scenario, consisting of twenty-two (22) operations,
 998 as described above (before each join operation there exists also one
 999 joining key sorting operation which is not shown in Fig. 13, so that
 1000 the flow is executable by most popular ETL engines).
- 1001 • *Case 2.* ETL scenario consisting of 27 operations, starting from the

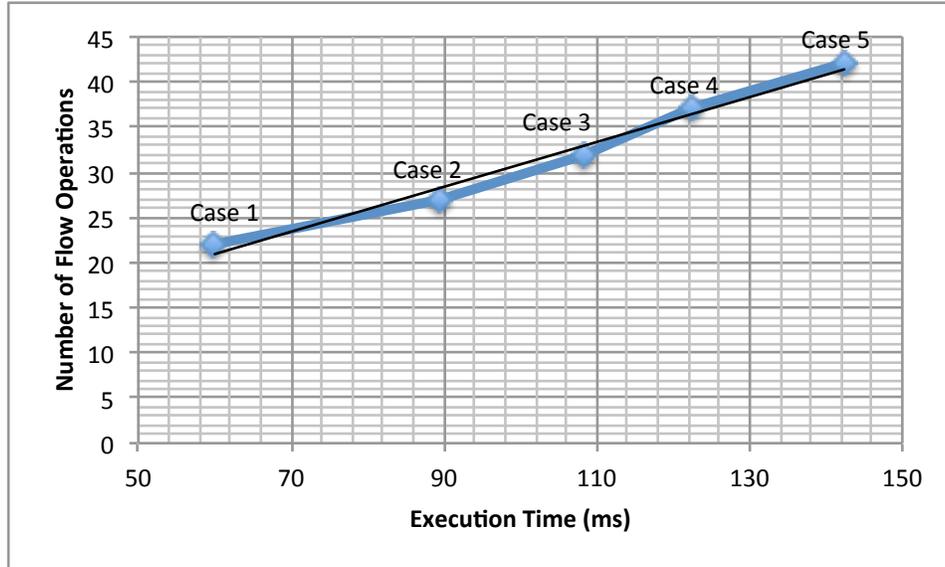


Figure 14: Linear trend of constraints extraction time wrt. the increasing number of operations (ETL flow complexity)

1002 basic one and adding an additional *Join* FCP and 2 *Filter* FCP to the
 1003 flow.

1004 • *Case 3.* ETL scenario consisting of 32 operations, starting from the
 1005 basic one and adding 2 additional *Join* FCP and 4 *Filter* FCP to the
 1006 flow.

1007 • *Case 4.* ETL scenario consisting of 37 operations, starting from the
 1008 basic one and adding 3 additional *Join* FCP and 6 *Filter* FCP to the
 1009 flow.

1010 • *Case 5.* ETL scenario consisting of 42 operations, starting from the
 1011 basic one and adding 4 additional *Join* FCP and 8 *Filter* FCP to the
 1012 flow.

1013 5.2. Experimental results

1014 We measure the average execution time of the path enumeration, ex-
 1015 traction and analysis phase for the above 5 scenarios covering different ETL
 1016 flow complexities.

1017 Figure 14 illustrates the increase of execution time when moving from the
 1018 simplest ETL scenario to a more complex one. As can be observed, execution

1019 time appears to follow a linear trend wrt. the number of operations of the
1020 ETL flow (i.e., flow complexity). This can be justified by the efficiency of our
1021 graph analysis algorithms and by the extensive use of indexing techniques
1022 (e.g., hash tables) to store computed properties for each operation and field,
1023 perhaps with a small overhead on memory usage. This result might appear
1024 contradictory, regarding the combinatorial part of our algorithm, computing
1025 and dealing with all possible path combinations. Despite the fact that it
1026 imposes factorial complexity, it is apparent that it does not constitute a
1027 performance issue for ETL flows of such complexity. To this end, the solution
1028 space is significantly reduced by i) our proposed greedy evaluation of the
1029 feasibility of a pattern every time it is updated and ii) by disregarding path
1030 combinations that do not comply to specific rules, e.g., *when considering*
1031 *path coverage, every input of a joining operation involved in any path of a*
1032 *path combination must be flowed (crossed by) at least one other path of that*
1033 *combination.*

1034 6. Related Work

1035 *Software development and testing.* In the software engineering field, *test-*
1036 *driven development* has studied the problem of software development by
1037 creating tests cases in advance for each newly added feature in the cur-
1038 rent software configuration [18]. However, in our work, we do not focus on
1039 the design (i.e., development) of ETL processes per se, but on automat-
1040 ing the evaluation of quality features of the existing designs. We analyze
1041 how the semantics of ETL processes entail the constraints over the input
1042 data, and then consequently create the testing data. Similarly, the problem
1043 of constraint-guided generation of synthetic data has been also previously
1044 studied in the field of software testing [5]. The context of this work is the mu-
1045 tation analysis of software programs, where for a program, there are several
1046 “mutants” (i.e., program instances created with small, incorrect modifica-
1047 tions from the initial system). The approach analyzes the constraints that
1048 “mutants” impose to the program execution and generates data to ensure
1049 the incorrectness of modified programs (i.e., “to kill the mutants”). This
1050 problem resembles our work in a way that it analyzes both the constraints
1051 when the program executes and when it fails to generate data to cover both
1052 scenarios. However, this work mostly considered generating data to test
1053 the correctness of the program executions and not its quality criteria (e.g.,
1054 performance, recoverability, reliability, etc.).

1055 *Data generation for relational databases.* Moving toward the database
1056 world, [19] presents a fault-based approach to the generation of database in-

1057 stances for application programs, specifically aiming to the data generation
1058 problem in support of white-box testing of embedded SQL programs. Given
1059 an SQL statement, the database schema definition and tester requirements,
1060 the approach generates a set of constraints, which can be given to existing
1061 constraints solvers. If the constraints are satisfiable, a desired database in-
1062 stances are obtained. Similarly, for testing the correctness of relational DB
1063 systems, a study in [20] proposes a semi-automatic approach for populating
1064 the database with meaningful data that satisfy database constraints. Work
1065 in [21] focuses on a specific set of constraints (i.e., cardinality constraints)
1066 and introduces efficient algorithms for generating synthetic databases that
1067 satisfy them. Unlike the previous attempts, in [21], the authors generate
1068 synthetic database instance from scratch, rather than by modifying the ex-
1069 isting one. Furthermore, [22] proposes a query-aware test database genera-
1070 tor called *QAGen*. The generated database satisfies not only constraints of
1071 database schemata, table semantics, but also the query along with the set of
1072 user-defined constraints on each query operator. Other work [23] presents a
1073 generic graph-based data generation approach, arguing that the graph rep-
1074 resentation supports the customizable data generation for databases with
1075 more complex attribute dependencies. The approach most similar to ours
1076 [24] proposes a multi-objective test set creation. They tackle the problem of
1077 generating "branch-adequate" test sets, which aims at creating test sets to
1078 guarantee the execution of each of the *reachable* branches of the program.
1079 Moreover, they model the data generation problem as a multi-objective
1080 search problem, focusing not only on covering the branch execution, but also
1081 on additional goals the tester might require, e.g., memory consumption cri-
1082 terion. However, the above works focus solely on relational data generation
1083 by resolving the constraints of the existing database systems. Our approach
1084 follows this line, but in a broader way, given that *Bijoux* is not restricted
1085 to relational schema and is able to tackle more complex constraint types,
1086 not supported by the SQL semantics (e.g., complex user defined functions,
1087 pivot/unpivot). In addition, we do not generate a single database instance,
1088 but rather the heterogeneous datasets based on different information (e.g.,
1089 input schema, data types, distribution, etc.) extracted from the ETL flow.

1090 *Benchmarking data integration processes.* In a more general context,
1091 both research and industry are particularly interested in benchmarking ETL
1092 and data integration processes in order to evaluate process designs and com-
1093 pare different integration tools (e.g., [25, 26]). Both these works note the lack
1094 of a widely accepted standard for evaluating data integration processes. The
1095 former work focuses on defining a benchmark at the logical level of data inte-
1096 gration processes, meanwhile assessing optimization criteria as configuration

1097 parameters. Whereas, the later works at the physical level by providing a
1098 multi-layered benchmarking platform called *DIPBench* used for evaluating
1099 the performance of data integration systems. These works also note that
1100 an important factor in benchmarking data integration systems is defining
1101 similar workloads while testing different scenarios to evaluate the process
1102 design and measure satisfaction of different quality objectives. These ap-
1103 proaches do not provide any automatable means for generating benchmark
1104 data loads, while their conclusions do motivate our work in this direction.

1105 *General data generators.* Other approaches have been working on pro-
1106 viding data generators that are able to simulate real-world data sets for
1107 the purpose of benchmarking and evaluation. [27] presents one of the first
1108 attempts of how to generate synthetic data used as input for workloads
1109 when testing the performance of database systems. They mainly focus on
1110 the challenges of how to scale up and speed up the data generation process
1111 using parallel computer architectures. In [28], the authors present a tool
1112 called Big Data Generator Suite (BDGS) for generating Big Data mean-
1113 while preserving the 4V characteristics of Big Data ⁸. BDGS is part of the
1114 BigDataBench benchmark [29] and it is used to generate textual, graph and
1115 table structured datasets. BDGS uses samples of real world data, analyzes
1116 and extracts the characteristics of the existing data to generate loads of “self-
1117 similar” datasets. In [30], the parallel data generation framework (PDGF) is
1118 presented. PDGF generator uses XML configuration files for data descrip-
1119 tion and distribution and generates large-scale data loads. Thus its data
1120 generation functionalities can be used for benchmarking standard DBMSs as
1121 well as the large scale platforms (e.g., MapReduce platforms). Other pro-
1122 totypes (e.g., [31]) offer similar data generation functionalities. In general,
1123 this prototype allows *inter-rows*, *intra-rows*, and *inter-table* dependencies
1124 which are important when generating data for ETL processes as they must
1125 ensure the multidimensional integrity constraints of the target data stores.
1126 The above mentioned data generators provide powerful capabilities to ad-
1127 dress the issue of generating data for testing and benchmarking purposes
1128 for database systems. However, the data generation is not led by the con-
1129 straints that the operations entail over the input data, hence they cannot be
1130 customized for evaluating different quality features of ETL-like processes.

1131 *Process simulation.* Lastly, given that the simulation is a technique that
1132 imitates the behavior of real-life processes, and hence represents an impor-
1133 tant means for evaluating processes for different execution scenarios [32], we

⁸*volume, variety, velocity and veracity*

1134 discuss several works in the field of simulating business processes. Simula-
1135 tion models are usually expected to provide a qualitative and quantitative
1136 analysis that are useful during the re-engineering phase and generally for un-
1137 derstanding the process behavior and reaction due to changes in the process
1138 [33]. [34] further discusses several quality criteria that should be considered
1139 for the successful design of business processes (i.e., *correctness, relevance,*
1140 *economic efficiency, clarity, comparability, systematic design*). However, as
1141 shown in [35] most of the business process modeling tools do not provide
1142 full support for simulating business process execution and the analysis of the
1143 relevant quality objectives. We take the lessons learned from the simulation
1144 approaches in the general field of business processes and go a step further fo-
1145 cusing our work to data-centric (i.e., ETL) processes and the quality criteria
1146 for the design of this kind of processes [36, 3].

1147 7. Conclusions and Future Work

1148 In this paper, we study the problem of synthetic data generation in
1149 the context of multi-objective evaluation of ETL processes. We propose an
1150 ETL data generation framework (*Bijoux*), which aims at automating the
1151 parametrized data generation for evaluating different quality factors of ETL
1152 process models (e.g., data completeness, reliability, freshness, etc.), ensuring
1153 both accurate and efficient data delivery. Thus, beside the semantics of ETL
1154 operations and the constraints they imply over input data, *Bijoux* takes
1155 into account different quality-related parameters, extracted or configured by
1156 an end-user, and guarantees that generated datasets fulfill the restrictions
1157 implied by these parameters (e.g., operation selectivity).

1158 We have evaluated the feasibility and scalability of our approach by pro-
1159 totyping our data generation framework. The experimental results have
1160 shown a linear (but increasing) behavior of *Bijoux*'s overhead, which sug-
1161 gests that the algorithm is potentially scalable to accommodate more in-
1162 tensive tasks. At the same time, we have observed different optimization
1163 opportunities to scale up the performance of *Bijoux*, especially considering
1164 larger volumes of generated data.

1165 As an immediate future step, we plan on additionally validating and
1166 exploiting the functionality of this approach in the context of quality-driven
1167 ETL process design and tuning, as explained in our test case scenario.

1168 **8. Acknowledgements**

1169 This research has been funded by the European Commission through the
1170 Erasmus Mundus Joint Doctorate “Information Technologies for Business
1171 Intelligence - Doctoral College” (IT4BI-DC).

1172 **References**

- 1173 1. Barbacci, M., Klein, M.H., Longstaff, T.A., Weinstock, C.B.. Quality
1174 attributes. Tech. Rep.; CMU SEI; 1995.
- 1175 2. Simitsis, A., Wilkinson, K., Castellanos, M., Dayal, U.. Qox-driven etl
1176 design: reducing the cost of etl consulting engagements. In: *SIGMOD Con-*
1177 *ference*. 2009, p. 953–960.
- 1178 3. Theodorou, V., Abelló, A., Lehner, W., Thiele, M.. Quality measures for
1179 ETL processes: from goals to implementation. *Concurrency and Computation*
1180 *Practice and Experience: Version of record online* 2016;.
- 1181 4. Dumas, M., Rosa, M.L., Mendling, J., Reijers, H.A.. *Fundamentals of*
1182 *Business Process Management*. Springer; 2013. ISBN 978-3-642-33142-8.
- 1183 5. DeMillo, R.A., Offutt, A.J.. Constraint-based automatic test data generation.
1184 *IEEE Trans Software Eng* 1991;**17**(9):900–910.
- 1185 6. Strange, K.H.. ETL Was the Key to This Data Warehouse’s Success. March
1186 2002. Gartner Research, CS-15-3143.
- 1187 7. Pall, A.S., Khaira, J.S.. A comparative review of extraction, transformation
1188 and loading tools. *Database Systems Journal* 2013;**4**(2):42–51.
- 1189 8. Vassiliadis, P., Simitsis, A., Baikousi, E.. A taxonomy of etl activities. In:
1190 *DOLAP*. 2009, p. 25–32.
- 1191 9. Vassiliadis, P., Simitsis, A., Skiadopoulos, S.. Conceptual modeling for ETL
1192 processes. In: *DOLAP*. 2002, p. 14–21.
- 1193 10. Muñoz, L., Mazón, J.N., Pardillo, J., Trujillo, J.. Modelling etl processes
1194 of data warehouses with uml activity diagrams. In: *OTM Workshops*. 2008, p.
1195 44–53.
- 1196 11. Akkaoui, Z.E., Zimányi, E., Mazón, J.N., Trujillo, J.. A bpmn-based design
1197 and maintenance framework for etl processes. *IJDWM* 2013;**9**(3):46–72.
- 1198 12. Wilkinson, K., Simitsis, A., Castellanos, M., Dayal, U.. Leveraging business
1199 process models for etl design. In: *ER*. 2010, p. 15–30.
- 1200 13. Jovanovic, P., Simitsis, A., Wilkinson, K.. Engine independence for logical
1201 analytic flows. In: *ICDE*. 2014, p. 1060–1071.
- 1202 14. Hueske, F., Peters, M., Sax, M., Rheinländer, A., Bergmann, R., Krettek,
1203 A., et al. Opening the black boxes in data flow optimization. *PVLDB* 2012;
1204 **5**(11):1256–1267.

- 1205 15. Tziouvara, V., Vassiliadis, P., Simitsis, A.. Deciding the physical imple-
1206 mentation of etl workflows. In: *Proceedings of the ACM Tenth International*
1207 *Workshop on Data Warehousing and OLAP*; DOLAP '07. 2007, p. 49–56.
- 1208 16. Simitsis, A., Vassiliadis, P., Dayal, U., Karagiannis, A., Tziouvara, V..
1209 Benchmarking etl workflows. In: *Performance Evaluation and Benchmarking:*
1210 *First TPC Technology Conference, TPCTC 2009, Lyon, France, August 24-28,*
1211 *2009, Revised Selected Papers*. 2009, p. 199–220.
- 1212 17. Theodorou, V., Abelló, A., Thiele, M., Lehner, W.. POIESIS: a tool
1213 for quality-aware ETL process redesign. In: *Proceedings of the 18th Interna-*
1214 *tional Conference on Extending Database Technology, EDBT 2015, Brussels,*
1215 *Belgium, March 23-27, 2015*. 2015, p. 545–548.
- 1216 18. Beck, K.. *Test-driven development: by example*. Addison-Wesley Professional;
1217 2003.
- 1218 19. Zhang, J., Xu, C., Cheung, S.C.. Automatic generation of database instances
1219 for white-box testing. In: *COMPSAC*. 2001, p. 161–165.
- 1220 20. Chays, D., Dan, S., Frankl, P.G., Vokolos, F.I., Weber, E.J.. A framework
1221 for testing database applications. In: *ISSTA*. 2000, p. 147–157.
- 1222 21. Arasu, A., Kaushik, R., Li, J.. Data generation using declarative constraints.
1223 In: *SIGMOD Conference*. 2011, p. 685–696.
- 1224 22. Binnig, C., Kossmann, D., Lo, E., Özsu, M.T.. Qagen: generating query-
1225 aware test databases. In: *SIGMOD Conference*. 2007, p. 341–352.
- 1226 23. Houkjær, K., Torp, K., Wind, R.. Simple and realistic data generation. In:
1227 *VLDB*. 2006, p. 1243–1246.
- 1228 24. Lakhota, K., Harman, M., McMinn, P.. A multi-objective approach to
1229 search-based test data generation. In: *GECCO*. 2007, p. 1098–1105.
- 1230 25. Simitsis, A., Vassiliadis, P., Dayal, U., Karagiannis, A., Tziouvara, V..
1231 Benchmarking etl workflows. In: *TPCTC*. 2009, p. 199–220.
- 1232 26. Böhm, M., Habich, D., Lehner, W., Wloka, U.. Dipbench toolsuite: A
1233 framework for benchmarking integration systems. In: *ICDE*. 2008, p. 1596–
1234 1599.
- 1235 27. Gray, J., Sundaresan, P., Englert, S., Baclawski, K., Weinberger, P.J..
1236 Quickly Generating Billion-Record Synthetic Databases. In: *SIGMOD Con-*
1237 *ference*. 1994, p. 243–252.
- 1238 28. Ming, Z., Luo, C., Gao, W., et al. Bdgs: A scalable big data generator suite
1239 in big data benchmarking. *CoRR* 2014;[abs/1401.5465](https://arxiv.org/abs/1401.5465).

- 1240 29. Luo, C., Gao, W., Jia, Z., Han, R., Li, J., Lin, X., et al. Hand-
1241 book of BigDataBench (Version 3.1) - A Big Data Benchmark Suite. 2014.
1242 Last accessed: 13/05/2015; URL [http://prof.ict.ac.cn/BigDataBench/
1243 wp-content/uploads/2014/12/BigDataBench-handbook-6-12-16.pdf](http://prof.ict.ac.cn/BigDataBench/wp-content/uploads/2014/12/BigDataBench-handbook-6-12-16.pdf).
- 1244 30. Rabl, T., Frank, M., Sergieh, H.M., Kosch, H.. A data generator for
1245 cloud-scale benchmarking. In: *TPCTC*. 2010, p. 41–56.
- 1246 31. Hoag, J.E., Thompson, C.W.. A parallel general-purpose synthetic data
1247 generator. *SIGMOD Record* 2007;**36**(1):19–24.
- 1248 32. Paul, R.J., Hlupic, V., Giaglis, G.M.. Simulation modelling of business
1249 processes. In: *Proceedings of the 3 rd U.K. Academy of Information Systems
1250 Conference, McGraw-Hill*. McGraw-Hill; 1998, p. 311–320.
- 1251 33. Law, A.M., Kelton, W.D., Kelton, W.D.. *Simulation modeling and analysis*;
1252 vol. 2. McGraw-Hill; 1991.
- 1253 34. Becker, J., Kugeler, M., Rosemann, M.. *Process Management: a guide for
1254 the design of business processes: with 83 figures and 34 tables*. Springer; 2003.
- 1255 35. Jansen-Vullers, M., Netjes, M.. Business process simulation—a tool survey.
1256 In: *Workshop and Tutorial on Practical Use of Coloured Petri Nets and the
1257 CPN Tools*; vol. 38. 2006, p. –.
- 1258 36. Simitsis, A., Wilkinson, K., Castellanos, M., Dayal, U.. QoX-driven ETL
1259 design: reducing the cost of ETL consulting engagements. In: *SIGMOD*. 2009,
1260 p. 953–960.
1261