

# Scalable Satisfiability Checking and Test Data Generation from Modeling Diagrams

Yannis Smaragdakis · Christoph Csallner ·  
Ranjith Subramanian

Received: date / Accepted: date

**Abstract** We explore the automatic generation of test data that respect constraints expressed in the Object-Role Modeling (ORM) language. ORM is a popular conceptual modeling language, primarily targeting database applications, with significant uses in practice. The general problem of even checking whether an ORM diagram is satisfiable is quite hard: restricted forms are easily NP-hard and the problem is undecidable for some expressive formulations of ORM. Brute-force mapping to input for constraint and SAT solvers does not scale: state-of-the-art solvers fail to find data to satisfy uniqueness and mandatory constraints in realistic time even for small examples. We instead define a restricted subset of ORM that allows efficient reasoning yet contains most constraints overwhelmingly used in practice. We show that the problem of deciding whether these constraints are consistent (i.e., whether we can generate appropriate test data) is solvable in polynomial time, and we produce a highly efficient (interactive speed) checker. Additionally, we analyze over 160 ORM diagrams that capture data models from industrial practice and demonstrate that our subset of ORM is expressive enough to handle their vast majority.

**Keywords** ORM · modeling · testing · databases · NP-hardness · ORM-

## 1 Introduction

Modeling languages offer a concise way to capture design decisions and express specifications at a level of abstraction higher than concrete code. The higher level of ab-

---

Y. Smaragdakis  
Computer Science  
University of Massachusetts, Amherst  
E-mail: yannis@cs.umass.edu

C. Csallner  
Computer Science and Engineering  
University of Texas at Arlington  
E-mail: csallner@uta.edu

R. Subramanian  
TheFind.com  
E-mail: ranjith.subramanian@gmail.com

straction often lends itself to automated reasoning. Nevertheless, a competing trend is that of adding more and more functionality to modeling languages, in order to bridge the gap between design and implementation. Several modern modeling languages (e.g., UML) have distinct sub-languages, some of which are highly abstract, while others are much closer to the implementation. Automatic reasoning is easy in the former case but becomes quite hard in the latter.

In this article, we analyze a modeling notation and, in particular, its tradeoffs between expressiveness and automatic reasoning ability. Our modeling language is *Object-Role Modeling (ORM)* [11, 24]: a modern, popular, and powerful data modeling language. ORM is a general language for conceptual modeling, although its primary application is in the database domain. For concreteness, we use database terminology in this article. The problem we want to solve is the automatic generation of test data that respect semantic constraints expressed in the model.

Producing valid test data is a problem with several applications in practice. Developers often want sample data both to verify that their specification corresponds to their understanding of the data, and to test programs as they are being developed. Generating unconstrained data is unlikely to be appropriate, however. Useful data typically needs to respect many semantic constraints: e.g., on a given table a certain field’s values may be unique (i.e., the field is a key); one table’s contents may be a subset of the contents of another; data values of a field may need to be in a specific range; etc. Such constraints are often concisely captured in data modeling languages. Thus, it is convenient and intuitively appealing to use a high-level model, expressed in a data modeling language, as a blueprint for creating large volumes of concrete well-formed data automatically.

Nevertheless, the problem of satisfying well-formedness constraints mechanically is often quite hard. In our case, producing sample data from ORM models (called *diagrams* in ORM parlance) is at best an intractable problem: even a simplified version of basic ORM constraints, with each table (*predicate* in ORM terminology) holding up to a single entry, makes the problem NP-hard. Checking realistic ORM constraints is typically even harder: constraints that describe relationships among different entries of a predicate typically result in a double-exponential runtime complexity. Even with small amounts of data to generate, a naive translation of the constraints into logic needs to model not just the predicate interconnections but also the contents of predicates. The result typically far exceeds the capabilities of modern constraint solving tools. In an experiment with state-of-the-art solvers we could not get a satisfying assignment for as few as 3 entities with 20 elements each.

Practical uses of ORM modeling, however, often only concentrate on a modest subset of the language. The main constraints on predicates encountered in practice are special forms of internal uniqueness constraints (“this field or set of fields is a primary key”), mandatory constraints (“every value of this entity needs to participate in this predicate”) and subtype constraints on objects. Thus, we define an interesting subset of ORM,  $\text{ORM}^-$  (pronounced “ORM-minus”), with such commonly used yet simple constraints. We define  $\text{ORM}^-$  precisely and present an algorithm that a) detects errors that make a diagram unsatisfiable in polynomial time (relative to the diagram size); and b) produces large quantities of data in time proportional to the size of the data produced. (Full ORM inputs can be used with our algorithm but constraints outside the  $\text{ORM}^-$  subset are ignored.)

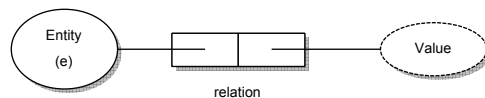
In overview, the novel elements of our work are as follows:

- We define  $\text{ORM}^-$ : a subset of the ORM modeling notation with desirable properties for automated reasoning. The satisfiability (i.e., consistency) of diagrams in  $\text{ORM}^-$  can be checked efficiently—both in theoretical terms (polynomial time) and in practice (sub-second run-time). Furthermore, test data can be generated efficiently for satisfiable  $\text{ORM}^-$  diagrams.  $\text{ORM}^-$  represents a sweet spot in the tradeoff between automation and expressibility: we show two small variations that make the problem intractable when added to the existing set of allowed constraints. Although there are similar results for other notations, we are not aware of any such result for ORM or modeling languages straightforwardly reducible to/from ORM.
- Our ORM subset is of practical interest: the vocabulary allowed was selected after consultation with database engineers who use the ORM notation. Our constraints are precisely those used heavily in practice, perhaps suggesting that database engineers avoid too-expressive constraints. We validate the applicability of  $\text{ORM}^-$  with an extensive study of ORM diagrams in real industrial use. We examine over 160 diagrams that contain some 1800 constraints in total. Of those constraints, only 24 are not expressible in  $\text{ORM}^-$ . The vast majority of diagrams are completely free of constraints outside the  $\text{ORM}^-$  subset.

## 2 Background: ORM

Object Role Modeling (ORM) is a data modeling language that attempts to model a system in terms of its objects and the roles that they play. ORM has a graphical notation and tries to capture common properties that are well-understood by database programmers. We next present briefly the most common elements of ORM. Some advanced elements are elided, in the interest of space, since they do not qualitatively affect our subsequent discussion.

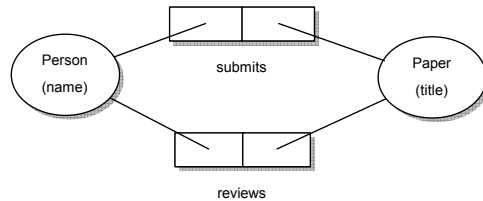
For the rest of this article, we follow the convention of calling a model specified in ORM a “diagram”. When we talk of a “model” (as a noun) of the diagram, we mean a set of data that satisfy the diagram’s constraints—analogously to the use of the term “model” in logic.



**Fig. 1** Basic notation

The basic components of the system being modeled are referred to as *objects*. Objects can be either *values* (i.e., well understood objects, such as a number or a string) or *entities* (derivative objects that are mapped to values). Object types are related to each other through *predicates*. Any number of object types can be related to each other, thus there can be predicates of any arity. If an object type  $e$  is related to another object type  $f$  through a predicate  $p$ , we say object types  $e$  and  $f$  play roles in the predicate  $p$ . The term *role* is used to refer to the relationship between an object type and a predicate, since an object type can be used in multiple predicates. In the relational world, a predicate can be thought of as a relation or table. The notation for value and entity types is shown in Figure 1. Value types are designated by a dotted-line ellipsis, while entity types are designated by a solid-line ellipsis. Typically each entity

is identified by a single value type (e.g., a name string), which is listed in parentheses under the entity name. In general, the difference between value types and entity types has little consequence in our examples, and we will occasionally use the term “entity” to mean “entity or value”.



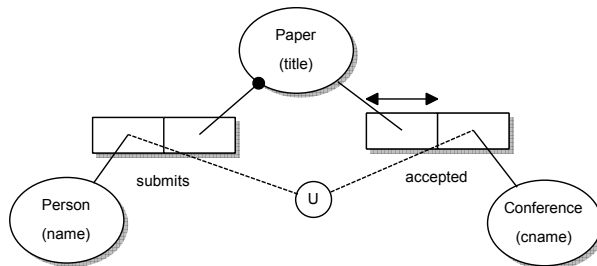
**Fig. 2** Simple ORM diagram

Consider the example in Figure 2. We are trying to model a simple system to store information about conference submissions. Two entity types, “Person” and “Paper” represent our objects. There are two predicates *reviews* and *submits*. Each of the entities play a role in both these predicates. A database that corresponds to this diagram will consist of two tables (assuming what is called the “standard mapping” of ORM diagrams to an implementation [12]): one with a list of people and the papers they submitted and another with a list of reviewers and papers they review.

There are no constraints specified on the ORM diagram in Figure 2. The modeler may want to specify a restriction that a person cannot both submit and review papers. Another restriction might be that a person can submit at most two papers. Or the modeler may require that a paper not have more than a set number of authors. Next, we list the main types of constraints that can be specified in an ORM diagram with examples. For a more detailed treatment of ORM, the reader is referred to Halpin’s book [11].

#### *Uniqueness constraints.*

This constraint specifies that the occurrence of a value in any column of a predicate is unique. In database terminology, a uniqueness constraint on a predicate identifies the keys of the corresponding table. Uniqueness can be represented in the ORM diagram by a double arrow over the roles that are unique in a predicate. For example, consider the diagram in Figure 3. The line over the role played by *Paper* in the predicate *accepted* indicates that a particular paper can occur in that table at most once—i.e., a paper can be accepted to at most one conference.



**Fig. 3** Uniqueness and mandatory constraints

Uniqueness can span multiple roles, implying that each tuple with values from the roles is unique. Furthermore, the roles can belong to different predicates—a constraint called *external uniqueness*. In this case, the tuple spanning multiple roles is considered on the (implicit) predicate resulting from joining the actual predicates on their common values. For instance, in Figure 3, the roles for *Person* in predicate *submits* and for *Conference* in predicate *accepted* are connected by an external uniqueness constraint, signified by the circled “U” in the diagram. This specifies that no person can have more than one paper in the same conference.

In our formulation of ORM, we enforce the common implicit uniqueness constraint spanning all roles of a predicate. That is, a predicate cannot contain two tuples that are the same across all roles.

#### *Mandatory constraints and independence.*

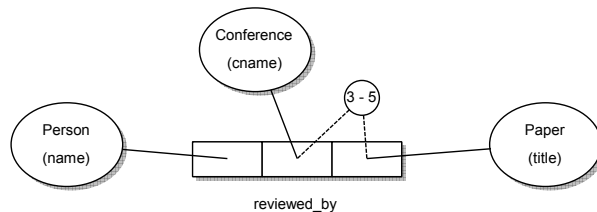
In Figure 3, the black dot connecting the entity *Paper* to the role it plays in the *submits* predicate denotes a *mandatory constraint*. This indicates that all papers in our universe are submitted.

A more complex form of the mandatory constraint is the *disjunctive mandatory constraint*. This links roles of the same entity in multiple predicates and signifies that all objects in the entity *have* to be a part of *some* (possibly all) of those predicates. This constraint is represented by connecting each of the roles that are mandatory to a black dot. We will not encounter an example of such a constraint in our discussion, except in its implicit form. In an ORM diagram, every object of an entity is implicitly assumed to participate in some of the predicates in which the entity has a role. This means that every entity in the system has an implicit disjunctive mandatory constraint over all its roles. (If an entity plays only one role, then this is equivalent to a regular mandatory constraint on that role. For instance, a black dot would be redundant on the links from entities *Person* or *Conference* in Figure 3.)

An exception to this rule is *independent* entity types, designated with an exclamation mark after their name. Objects of an independent type can exist without participating in predicates.

#### *Frequency constraints.*

Frequency constraints generalize uniqueness constraints, by allowing each tuple to occur a number of times, instead of just once. In Figure 4, the frequency constraint of (3-5) on the roles played by *Conference* and *Paper* states that each combination of conference and paper that appears in the predicate has to appear three to five times—i.e., each paper needs three to five reviewers for the same conference.

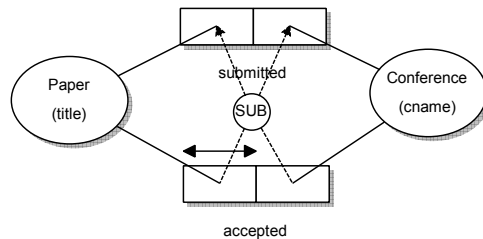


**Fig. 4** Frequency constraints

*Subset, Equality constraints.*

A *subset* constraint can be applied over two ranges of roles in two predicates. The constraint specifies that the tuples in the roles of the first predicate have to be among the tuples in the roles of the second. An important restriction here is that the types of the roles of the sub-predicate in a subset constraint have to match the corresponding roles of the super-predicate on a role by role basis. The *equality* constraint is a two-way subset: it forces the two sets of tuples to be equal, as the name suggests.

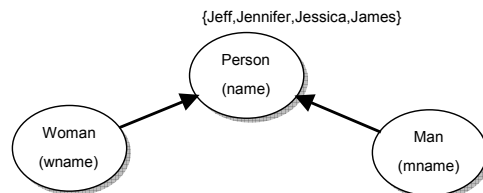
In the example of Figure 5, the subset constraint spans the entire *accepted* predicate, specifying that a paper needs to have been submitted to the conference, if it is to be accepted.



**Fig. 5** Subset constraint

*Subtype constraints.*

This constraint is used to model usual subtyping situations. For example, the *Man* and *Woman* entity types can be declared as subtypes of *Person*, as shown in Figure 6.



**Fig. 6** Subtype and value constraints

*Value and cardinality constraints.*

A value constraint is used to enumerate the range of values that an entity type may have. In the example of Figure 6, the allowed values for the *Person* entity type are given explicitly. Value constraints can also be used with a range notation—e.g., one can specify that the *Age* value type needs to be a number between 1 and 150. It is important to make the distinction between a value constraint, which specify the allowed values in a type, and *cardinality* constraints, which specify how many elements a type has. Clearly a value constraint implies an upper bound on the cardinality. Our graphical notation does not reflect cardinality constraints, although they do exist in the diagram. A cardinality constraint can be either a number or a range.

---

### Exclusion constraints.

This constraint is used to model a situation where a particular value can occur in only one table in any valid model of the diagram. For instance, in a system that stores the *current* state of papers, a paper cannot be both submitted and accepted. We will see the graphical notation for exclusion constraints in Section 3.1.

### Ring constraints.

A predicate represents a relation over the sets of values for each of the entities playing roles in the predicate. For predicates with all roles played by the same entity, we can specify various properties that the relation should possess. For example, in Figure 7, we state that the *works\_with* predicate is *symmetric*. That is, if a person *A* works with a person *B*, then *B* also works with *A*.

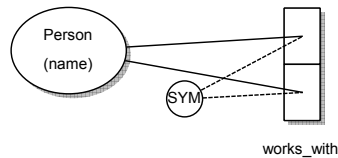


Fig. 7 Ring constraint

Other types of ring constraints include *acyclic*, *transitive*, *intransitive*, *reflexive*, *irreflexive*, and *asymmetric*.

## 3 Difficulty of Analyzing ORM

The problem we want to address is that of generating sample databases from ORM diagrams. We assume that we receive as input an ORM diagram and we want to determine its satisfiability (i.e., the existence of data that populate all types and predicates and satisfy all constraints). Some variants of ORM (including the one supported by the tool used in our experiments) allow arbitrary query-like code (in Datalog or SQL) to be used as a notation for constraints,<sup>1</sup> making the problem of determining satisfiability undecidable. (Although query languages guarantee termination, other common problems, such as query subsumption, are undecidable.) Even without allowing arbitrary queries, undecidability is not uncommon when dealing with highly expressive integrity constraints. For instance, the general class of “numeric dependencies” [8] in databases has no decidable inference process.

Even if we limit ourselves to “standard” ORM constraints, however, the problem is at least NP-hard. We next discuss the difficulties involved, and present a (failed) brute-force attempt at the problem, which serves to frame our later restriction of ORM to a simpler subset.

---

<sup>1</sup> The full set of ORM constraints is not standardized (“not all versions of ORM support all these symbols” [10]). The constraint that allows us to write queries as part of the diagram is the asterisk (“\*”) constraint—symbol 23 in [10].

### 3.1 NP-Hardness

Our problem of producing satisfying data for ORM diagrams is NP-hard. Bommel et al. have shown [28] the NP-hardness of some slightly different variants of the ORM satisfiability problem. Nevertheless, it is easy to produce a much simpler NP-hardness proof for our exact formulation of ORM and the satisfiability problem.

We reduce the well-known NP-complete boolean satisfiability problem, 3SAT, to our problem. A 3-CNF formula has the form  $(a_1 \vee a_2 \vee a_3) \wedge (b_1 \vee b_2 \vee b_3) \wedge \dots$ , with  $a_1, a_2, a_3, b_1, \dots$  representing expressions  $q_i$  or  $\neg q_i$  over a set of boolean variables  $q_1, q_2, \dots, q_n$ . A CNF formula can be translated to an ORM diagram using the following steps.

1. Introduce an object type for each variable  $q_i$  in the CNF formula
2. Introduce an object type for each clause  $c_i$  in the CNF formula
3. Impose a value constraint on each type introduced above. The value constraint restricts each type  $q_i$  to have just one value ‘qi’. Similarly, each clause is restricted to have one value ‘ci’.
4. For each variable  $q_i$ , let  $n$  be the number of times it occurs in the clauses  $c_i$  to  $c_j$ . If  $n > 0$ , create a  $(n + 1)$ -ary predicate  $q_i$ -true for the variable  $q_i$ .
5. Similarly, for each variable  $\neg q_i$  let  $m$  be the number of times it occurs in the clauses  $c_i$  to  $c_j$ . If  $m > 0$ , create a  $(m+1)$ -ary predicate  $q_i$ -false for the variable  $\neg q_i$ .
6. Connect type  $q_i$  to predicates  $q_i$ -true and  $q_i$ -false
7. Impose an *exclusion constraint* on the roles played by  $q_i$  in  $q_i$ -true and  $q_i$ -false
8. Connect each clause  $c_i$  to the predicate  $q_i$ -true if it contains  $q_i$  or  $q_i$ -false if it contains the variable  $\neg q_i$

Consider a CNF formula  $(q_1 \vee q_2 \vee \neg q_3) \wedge (q_1 \vee \neg q_2 \vee q_3) \wedge (\neg q_1 \vee q_2 \vee \neg q_3) \wedge \dots(\dots)$ . Figure 8 illustrates the relevant part of the formula’s transformation.

Based on the semantics of ORM constraints, it is easy to see that the 3SAT formula is satisfiable iff the ORM diagram is satisfiable. The exclusion constraints on the roles played by the  $q_i$ s force each value to appear in at most one of the  $q_i$ -true or  $q_i$ -false predicates. Recall that there is an implicit disjunctive mandatory constraint over all the roles played by the same type, since none of the types are “independent”. Therefore, each  $q_i$  value appears exactly once in the predicates. A similar implicit disjunctive mandatory constraint on the roles played by the  $c_i$ s ensures that the single value of each of these entities appears in at least one of the predicates, effectively implementing a disjunction.

Note that the ORM satisfiability problem is NP-hard, but not clearly NP-complete. Bommel et al. claimed an NP-completeness proof for their two satisfiability problems [28]. The proof is based on a brute-force nondeterministic guess of the contents of each predicate. The claim is that because all populations are bounded by  $m$ —the maximum frequency constraint upper bound—one can nondeterministically choose populations with up to  $m$  elements for each entity and predicate in polynomial time. Yet the populations will be of up to size  $m$ , which is *exponential* relative to the *representation* of  $m$ , which is the input to the problem. Thus the argument of Bommel et al. is wrong: even a non-deterministic Turing machine cannot guess all possible populations in time polynomial to the size of the input. (We reported this finding to the authors in May 2006.)

It is possible that when complex constraints are included (such as subset or ring constraints) the problem is undecidable (or with a double-exponential complexity, if



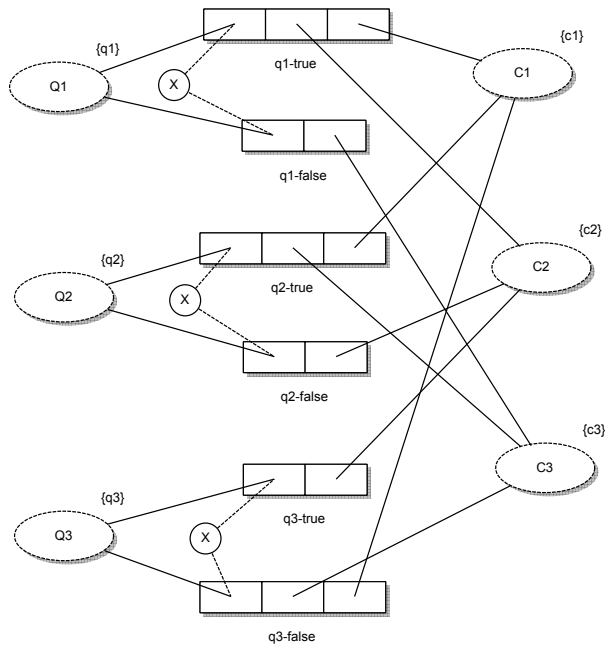


Fig. 8 Example translation of 3SAT to ORM

there is a known size bound). Yet, for a simpler subset of constraints the problem is just NP-complete, and for a yet simpler subset it is polynomial, as we discuss in Section 4.

### 3.2 Brute-Force Approach

As we saw, the ORM diagram satisfiability problem is fairly hard. There are, however, good reasons to try a combination of brute-force translation and heuristic approaches. There is a standard translation of ORM constraints to first-order logic (e.g., [9]) that directly captures the semantic intricacies of constraints. At the same time, reasoning tools have matured and state-of-the-art constraint solvers and SAT solvers often achieve impressive scalability. In practice, we found the brute-force + heuristic approach to be a bad fit for our problem. Nevertheless, the result yields insights that suggest a scalable solution.

We translated ORM diagrams to inputs for two different tools: the CoBaSa tool [20] and the Alloy Analyzer [14, 25]. This is a good choice of candidate tools: Alloy is well-known and mature, yet does not emphasize scalability, while CoBaSa offers an expressive front-end, an efficient translation, and a state-of-the-art pseudo-boolean constraint solver (PBS [1]) as its back-end. We used the standard translation of ORM to first-order logic, customized to the needs of the tool at hand. Generally we map every ORM predicate to a logic predicate, every value of an entity to a logical value, and specify constraints using a logic notation. For instance, if there is a mandatory constraint on the role played by the entity  $A$  in a predicate, this would translate to CoBaSa as:

```
For_all a in A { Sum i in id tableA(i, a) >= 1 }
```

**Table 1** Relation (A,B,C) with frequency constraint of (min=2, max=3) spanning (B,C). B and C are mandatory. PBS aborts SAT solving after 1500 seconds. The first four columns give cardinalities. “vars” and “clauses” measure the complexity of the SAT problem. SAT indicates a satisfiable problem. The last column is the CoBaSa processing time.

A	B	C	pred	vars	clauses	SAT	time [s]
3	2	2	10	150	200	yes	0.23
3	2	2	15	225	300	no	197.92
10	2	2	12	264	240	yes	0.12
10	2	2	13	286	260	no	119.51
10	2	2	15	330	300	no	260.67
10	4	4	40	2000	3200	yes	0.90
10	5	5	25			no	abort
10	5	5	70	4900	8750	yes	3.84
10	6	6	100	9400	18000	yes	4.76
10	8	8	150	23100	48000	yes	22.50
10	10	10	200	46000	100000	yes	68.14
10	10	10	250	57500	125000	yes	180.10

This states that every value in the domain of  $A$  has to occur at least once in *tableA*.

The results of one of our CoBaSa experiments are shown in Table 1. We use an input with three entity types: A, B, C. These are linked with a single ternary predicate (A,B,C) and a frequency constraint spanning (B,C). Furthermore, B and C are mandatory. This simple input should yield a lower bound for the cost of adding more complex constraints.

Overall, this approach does not scale very well. Even our simple, small examples take minutes to process. The problem size (in terms of variables and clauses) scales exponentially relative to the original input. (Note that the original input is the logarithmic representation of the number  $n$  of entities, while the input to the solver is  $O(n^c)$  clauses—in this case  $c = 4$ .) We certainly cannot answer the question of whether a satisfying configuration exists for realistic inputs (which will specify that entities contain millions of objects: recall that we want to produce test databases that satisfy the stated constraints). Our experience with the Alloy Analyzer was quite similar: we found it unable to yield solutions for an example domain of 3 entities with 5 to 40 values each.

The problem is that the brute-force approach is not a very good fit for our problem. It is much better for dealing with logically deep constraints, than with size constraints on a large space. Nevertheless, the approach has value for detecting some unsatisfiable (i.e., contradictory) configurations. It is often the case that a small model is sufficient for disproving the consistency of constraints. This has been a standard argument for practical uses of Alloy. For this to apply to our problem, the original ORM diagram must contain no explicit size constraints on object types or role frequencies that will make the search space unmanageable. A good example is a diagram that has no frequency, value, or cardinality constraints but has a predicate with two ring constraints, one irreflexive, one transitive, and two mandatory constraints on the same roles. This is a logically unsatisfiable constraint system: a transitive relation on a finite space will eventually include either a limit element or a cycle. The limit element cannot occur on the right side of the relation (so the mandatory constraint is violated), and the cycle (plus transitivity) contradicts the irreflexivity. This deep-but-not-size-sensitive reasoning is handled well by a brute-force translation into input for constraint solvers.

Although the brute-force approach fails, it suggests the direction to proceed. The problem with the brute-force approach is that it entails an inherent blowup in the

problem size, because it models the interrelations of the *members* of types. This makes the problem exponentially harder: for an input of  $N$  bits, the possible values are up to  $2^N$  and there are up to  $2^{2^N}$  possibilities for the contents of entities and predicates. For complex constraints, such as ring, subset, and exclusion, this modeling seems necessary. For other constraints, however, we only need to concern ourselves with the *sizes* of types. Thus, if we limit ourselves to simple constraints that only deal with the sizes of sets and not their contents, we can obtain a scalable solution. Fortunately, these simple constraints turn out to be exactly the ones used overwhelmingly in practice.

#### 4 ORM<sup>-</sup>: An Efficient and Expressive Subset

Given the difficulty of producing concrete data from full ORM specifications, we concentrated on a realistic subset of ORM. This subset was arrived upon through continuous refinement over the course of several months, in consultation with engineers of LogicBlox—a company that markets a database engine and ORM tools. The end result is a subset of ORM that captures the most commonly used constraints in practice, yet without sacrificing the ability to generate data very efficiently. We call our ORM subset ORM<sup>-</sup> (pronounced “ORM minus”).

##### 4.1 Definition of ORM<sup>-</sup>

ORM<sup>-</sup> is the ORM subset with the following constraints:

- *Uniqueness* constraints: only internal uniqueness (i.e., over roles in a single predicate) is supported. Uniqueness constraints can span multiple roles, but no two uniqueness constraints can overlap. E.g., it is not expressible that, in a single predicate, both roles A-and-B and roles B-and-C form unique keys.
- *Mandatory* constraints on a single role. Disjunctive mandatory constraints are not allowed, although the standard implicit disjunctive mandatory constraint holds, over all roles played by the same (non-independent) type. Independent types are also supported.
- *Frequency* constraints, which, just like uniqueness constraints, cannot overlap.
- *Value and cardinality*.
- *Subtype* constraints. We assume that a subtype’s value constraints specify a subset of those of the supertype.

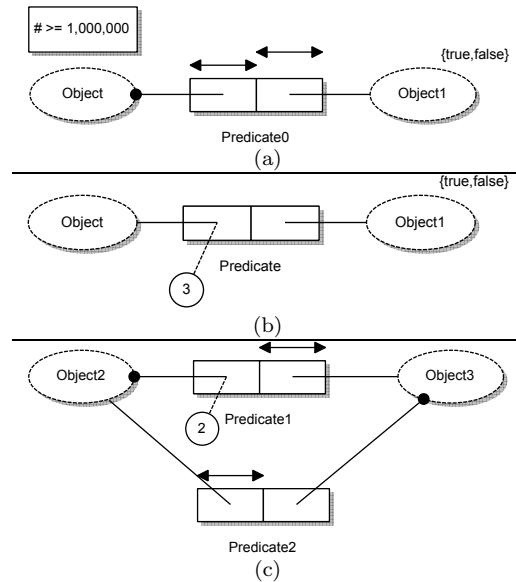
Notably absent are subset constraints, exclusion constraints, and ring constraints. It is possible that some of those can be added while maintaining the desirable properties of ORM<sup>-</sup> (i.e., sound and complete satisfiability decision in polynomial time).<sup>2</sup> Some additions (e.g., exclusion) immediately result in an NP-hard problem, however.

---

<sup>2</sup> For instance, we speculate that with a more complex translation we can support external uniqueness constraints, overlapping frequency constraints, as well as subset constraints under the restriction that all other constraints on the superset predicate are identical to the constraints of the subset predicate. We have not, however, attempted to prove that such a generalization works.

4.2 Testing Satisfiability for  $ORM^-$ 

A first question regarding satisfiability of  $ORM^-$  is “what are instances of *unsatisfiability*?” Diagrams in  $ORM^-$  can be unsatisfiable for a variety of reasons. Figure 9 shows three examples.



**Fig. 9** Three unsatisfiable diagrams.

The top example (a) represents a simple instance of contradictory cardinality constraints. (Recall that cardinality constraints are not depicted graphically in our diagrams, although they are maintained in the properties of diagram elements. In this case, we show the constraint in a text comment.) The diagram cannot be satisfied since it requires at least a million objects of one type (the designer wants a large test database) but the type is in a one-to-one mapping with a subset of a two-object type. The middle diagram (b) cannot be satisfied since the predicate needs to contain three tuples for each left-hand-side value, yet the right-hand-side type has at most two values. This requirement violates the implicit constraint that all tuples in a predicate be unique. The bottom diagram (c) is a bit more interesting: through a cycle of constraints the cardinality of type *Object2* is required to be at least two times itself.

We can test the satisfiability of  $ORM^-$  diagrams in polynomial time. The key property of  $ORM^-$  is that all its constraints can be expressed as *numeric constraints* (inequalities) on the sizes of sets, instead of their contents. The sets are polynomial in number, relative to the original input. Specifically, we translate all  $ORM^-$  constraints into numeric constraints of the form  $c \cdot x \leq y_1 \cdot y_2 \cdot \dots \cdot y_n$ , or  $x \leq y_1 + y_2 + \dots + y_n$ , where  $c$  is a positive integer constant and  $x, y_1, \dots, y_n$  are positive integer constants or variables. Note that, although these are constraints on the integers (suggesting that the problem is hard), we have no addition or multiplication of *variables* (i.e., unknowns) on the left hand side of the inequality.

Specifically, the translation rules are as follows (we list clarifying assumptions or explanations as second level bullets):

- For each type A, introduce a fresh variable  $a$ , representing its cardinality. (We follow this convention of upper and lower case letters in the following.) For each predicate R, introduce a variable  $r$ , representing its cardinality (number of tuples—distinct by definition). For each role in R played by type A, introduce a variable  $r_a$ , representing the number of unique elements from A in this role.
- Produce inequalities  $r_a \leq a$ , and  $r_a \leq r$ , for each role that type A plays in predicate R.
- For each *cardinality* constraint that limits the cardinality of a type A to be in a range  $min...max$ , produce inequalities  $min \leq a$  and  $a \leq max$ . Similarly for predicates.
  - We assume, without loss of generality, that each type and predicate has a maximum cardinality constraint. (They already have a minimum cardinality constraint of 1, as we do not want empty types or predicates in the solution.) In practice, for types or predicates that do not have an explicit cardinality constraint, we can produce the numeric constraint  $a \leq M$ , where  $M$  is an upper bound of the size of all entities or of the desired output (e.g., MaxInt).
- For each *value* constraint on a type A, produce inequality  $a \leq ds$ , where  $ds$  is the domain size defined by the value constraint.
  - E.g., for a value constraint that restricts a type to four distinct values  $ds = 4$ . For a value constraint of the form 5...30,  $ds = 26$ .
- For each *mandatory* constraint on the role played by type A in predicate R, produce inequality  $a \leq r_a$ .
  - Since we also have  $r_a \leq a$  for each role, the result is an equality, but all our reasoning is done in the inequality form.
- For each *frequency* constraint, with frequency range  $f_{min}...f_{max}$ , on the roles played in predicate R by types A, B, ..., K, introduce a variable  $r_{ab...k}$ , representing the number of unique tuples in these roles. Produce the following inequalities:
 
$$r \leq f_{max} \cdot r_{ab...k}$$

$$f_{min} \cdot r_{ab...k} \leq r$$

$$r_{ab...k} \leq r_a \cdot r_b \cdot \dots \cdot r_k$$

$$r_a \leq r_{ab...k}, r_b \leq r_{ab...k}, \dots, r_k \leq r_{ab...k}$$
  - Note that we are representing the number of unique tuples only for sets of roles that have a frequency constraint, not for all subsets of the roles of a predicate. Note also that uniqueness constraints and constant frequency constraints are a special case of the above. (A uniqueness constraint is a frequency constraint with a minimum and maximum frequency of 1.)
- For each *subtype* constraint between types A and S, produce the inequality  $a \leq s$ .
- For each predicate R with roles played by types A, B, ..., N, express the implicit uniqueness constraint over all roles as follows: Produce inequality  $r \leq r_{ab...k} \cdot r_{lm...p} \cdot \dots \cdot r_n$ , where the elements of the right hand side are all the variables corresponding to the ranges of roles that participate in some frequency constraint on the predicate, followed by all the variables corresponding to roles that are under no frequency constraint.
  - Recall that frequency constraints cannot overlap.
- For each non-independent type A that plays roles in predicates R, S, ..., V, introduce the numeric constraint  $a \leq r_a + s_a + \dots + v_a$ .

- This captures the implicit disjunctive mandatory constraint over all roles played by a type: each object of a non-independent type needs to appear in some predicate.

For illustration, consider the translation of the top example of Figure 9. If we call the left type in the diagram “A”, the predicate “R”, and the right type “B”, and follow our naming convention for numeric variables, we get the inequalities:  $1000000 \leq a$  (cardinality on A),  $a \leq r_a$  (mandatory on A),  $r \leq r_a$  (uniqueness on role of A),  $r \leq r_b$  (uniqueness on role of B),  $b \leq 2$  (value on B), as well as the implicit constraints on roles:  $r_a \leq a$ ,  $r_a \leq r$ ,  $r_b \leq r$ ,  $r_b \leq b$ , and the implicit cardinality constraints that make every variable at least 1 and smaller than a constant upper bound. (We omit the implicit uniqueness over all roles, and the implicit disjunctive mandatory constraint, as they are redundant for this example. The implicit mandatory constraint cannot be used, anyway, unless we know that the elements shown constitute the whole diagram and not just a part of it.) It is easy to see that these numeric constraints imply  $1000000 \leq 2$ —a contradiction.

The above translation supports the key properties of  $\text{ORM}^-$ : The translation is sound and complete, in that an ORM diagram is satisfiable if and only if the integer inequalities resulting from the translation admit a solution. Furthermore, solving the inequalities can be done in polynomial time. Note that the above does *not* mean that every dataset whose type, role, and predicate cardinalities satisfy the integer inequalities will satisfy the ORM constraints! Instead, it means that for every satisfying assignment of the integer inequalities, we can produce *some* dataset that will satisfy the ORM constraints.

We next present the key arguments establishing these properties: We first show how to satisfy the diagram constraints if the inequalities have a solution, and then discuss how we can check the inequalities in polynomial time.

#### *Sound and Complete Decision.*

If the ORM constraints are satisfiable, then the numeric constraints also hold: the numeric constraints just capture size properties of sets of values that satisfy the ORM constraints. The inverse direction is harder, but our process of showing this also yields a way to produce objects that satisfy the ORM constraints. Starting from an integer variable assignment satisfying all inequalities, we can show that there exists a data set that satisfies the original ORM constraints, as follows:

#### *Algorithm 1: Existence of Satisfying Data Assignment*

- For each type A without a subtype, create  $a$  new objects from the set specified by the type’s value constraint (if one exists). (We follow our lower/upper case convention—i.e.,  $a$  is type A’s cardinality in the solution of the size inequalities.) The crucial observation is that  $\text{ORM}^-$  constraints do not restrict *which* of these objects participate in a role, predicate, or type, except in three cases: a subtype should only contain objects from its supertype; the objects in a role  $R_A$  on predicate R played by type A should be a projection over  $R_A$  of all the tuples of any range of roles that includes  $R_A$ , including the entire R; the implicit disjunctive mandatory constraint should hold: each object of a type should participate in one of the roles. We can satisfy these requirements with an appropriate choice of objects. Note that in this section we discuss enumerations of objects without specifying in detail how the enumeration

is produced, as long as it is clear that one exists. The exact efficient algorithm for generating the desired data is described in the next section (4.3).

- For each supertype B of a type that has been populated with objects, pick  $b$  objects so that they include the objects in the subtype. (Recall that the value constraints of the subtype specify a subset of the value constraints of the supertype, per the definition of ORM<sup>-</sup>.) Repeat until all types are populated. The ORM subtype constraints of the diagram are now satisfied, as are any value and cardinality constraints on types.
- For each type A, populate all roles  $R_A, S_A, \dots, V_A$  that A plays with  $r_a, s_a$ , etc. elements. If A is not independent, ensure that all elements of A participate in some role. This is possible since  $a \leq r_a + s_a + \dots + v_a$ . After this step, the implicit disjunctive mandatory and any explicit mandatory constraints of the ORM diagram are satisfied, by straightforward properties of the construction.
- At this point all roles are populated with the unique objects they contain. (I.e., we have computed the projection of each predicate on each role.) For each range of roles  $R_A, R_B, \dots, R_K$  that participate in a frequency (or uniqueness, as a special case) constraint, pick  $r_{ab\dots k}$  distinct tuples so that their projections on roles are the computed population of  $R_A, R_B$ , etc. Such tuples are guaranteed to exist since  $r_{ab\dots k} \leq r_a \cdot r_b \cdot \dots \cdot r_k$ .
- Finally, for each predicate produce its entire contents by enumerating the unique combinations of the sub-tuples from each range of roles that participates in a frequency constraint, and the objects in the roles that are yet unfilled (i.e., roles that do not participate in frequency constraints) until producing  $r$  tuples. This is guaranteed to be possible since  $r \leq r_{ab\dots k} \cdot r_{lm\dots p} \cdot \dots \cdot r_n$ . The enumeration of tuples should be done so that the sub-tuples from roles in frequency constraints are covered evenly (i.e., two sub-tuples of the same role range are used either the same number of times or with a difference of one) and all objects from roles that are not under frequency constraints are used. This is guaranteed to produce the right frequencies for the contents of ranges of roles used in frequency constraints, since for each such range  $R_A, \dots, R_K$ , we have  $r \leq f_{max} \cdot r_{ab\dots k}$ , or  $r/r_{ab\dots k} \leq f_{max}$ , and similarly for  $f_{min}$ . That is, in producing  $r$  complete tuples, we are guaranteed to repeat each sub-tuple in role range  $R_A, \dots, R_K$  at least  $f_{min}$  and at most  $f_{max}$  times.

#### *Polynomial Testing.*

We can test the satisfiability of numeric inequalities produced by our translation through a fixpoint algorithm. Recall that all of our constraints are of the form  $c \cdot x \leq y_1 \cdot y_2 \cdot \dots \cdot y_n$ , or  $x \leq y_1 + y_2 + \dots + y_n$ , where  $c$  is a positive integer constant and  $x, y_1, \dots, y_n$  are positive integer constants or variables.

#### *Algorithm 2: Polynomial Satisfiability Check*

1. We first rewrite all inequalities of the form  $c \cdot x \leq y_1 \cdot y_2 \cdot \dots \cdot y_n$  into  $x \leq (y_1 \cdot y_2 \cdot \dots \cdot y_n) / c$ .
2. For each variable, we can maintain its current upper bound (recall that we assume cardinality constraints for all entities and predicates in the original input, although these can be MaxInt) and, on every step, we substitute the current upper bounds for all variables on the right hand side of each inequality and perform the arithmetic operations (the division is integer division, performed after all multiplications). This produces candidate upper bounds for variables. The minimum of the candidate upper bounds for a variable and its current upper bound becomes the variable's upper bound for the next step.

3. Step 2 repeats until either an upper bound crosses a lower bound (i.e., we get  $x \leq k$  and  $m \leq x$ , where  $m > k$  for some variable  $x$ ) or all upper bounds remain unchanged during a step. The former signifies an unsatisfiable case, the latter a satisfying assignment.

This algorithm is correct because it maintains a conservative upper bound on cardinalities on any step. Any satisfying assignment will need to have values at most equal to the respective upper bounds. When the algorithm reaches a fixpoint, all inequalities are satisfied by the current upper bounds, which, thus, constitute a solution.

It is easy to see that this algorithm runs in polynomial time relative to the size of the input. The fixpoint computation runs a polynomial number of times between successive applications of inequalities that contain division, i.e., of the form  $x \leq (y_1 \cdot y_2 \cdot \dots \cdot y_n)/c$ . (Inequalities of other forms do not shrink the strictest upper bound in the system, only propagate it to new variables. If there are  $n$  variables in the system, each inequality of a different form can yield a different upper bound at most  $n$  times: the variable on its right hand side with the smallest current bound can never get any smaller, as it satisfies all other inequalities in whose left hand side it appears—since the right hand sides can only have multiplication with positive integers. Similarly, the second-smallest-upper-bound variable will never get smaller after one round, the third will not shrink after two rounds, etc.) Hence, the only reason that candidate upper bounds may keep becoming smaller is inequalities that contain the division operator. Yet each of those can apply at most a number of times logarithmic relative to the value of the original upper bounds (since it reduces the upper bound by a multiplicative factor) which is polynomial relative to the input size.

#### *Discussion and Insights.*

A polynomial solution to this problem exists only because there is no addition or multiplication of variables on the left-hand side of an inequality (i.e., all our arithmetic operations are on upper bounds, yielding a monotonicity of the problem). Interestingly, simple additions to the  $\text{ORM}^-$  constraints result in addition or multiplication on the left-hand side, and make the problem NP-hard. For instance, adding exclusion constraints corresponds to variable addition on the left-hand side. Also, our constraint language does not support anything like an (imaginary) “external mandatory” constraint, that would require all combinations of objects (cross product) of two types to appear in a single predicate. Such a constraint would result into multiplication of variables on the left-hand side. For such richer languages, it is straightforward to show NP-hardness by taking techniques from the domain of integer inequalities and mapping them back to our constraints. E.g., with just multiplication of integers, we can encode an instance of 3SAT as a set of inequalities. For each propositional symbol  $q_i$  and its negation, we introduce two integer variables  $x_i$  and  $x'_i$ , both of which take values from 1 to 2 (for false and true, respectively). Then a clause like  $(q_1 \vee q_2 \vee \neg q_3)$  gets mapped into  $2 \leq x_1 \cdot x_2 \cdot x'_3$ . Multiplication on the left hand side of an inequality is used to establish that no symbol and its negation can be simultaneously true:  $x_i \cdot x'_i \leq 2$ . Our sound and complete translation of Algorithm 1 can then be used as the basis for deriving an inverse transformation where the inequalities get mapped back to diagrams that are only satisfiable if the respective inequalities are.



### 4.3 Generating Data

The previous section gave a sound and complete algorithm for determining the satisfiability of ORM<sup>-</sup> diagrams. The completeness depended on creating model data that satisfy the constraints of the diagram. This process is precisely how we generate test data. Nevertheless, it is not immediately obvious how to do this efficiently. Namely, although it is clear that producing the desired combinations of values is possible (because of the combinatorics entailed by the numeric constraints) it is not clear that doing so does not require a super-exponential search: The universe of potential output tuples is exponential relative to our input, and having to do a naive trial-and-error search for a set of tuples in this exponential space would yield prohibitive (double-exponential) complexity. We next discuss how the enumerations of the previous section can be produced efficiently.

To see a typical enumeration problem, consider an example set of cardinalities for a predicate P. If  $p = 15$  (i.e., we want to generate a table for P with 15 tuples in it),  $p_a = 4$ ,  $p_b = 6$ ,  $p_c = 2$  (i.e., 4, 6, and 2 distinct elements of types A, B, and C, respectively, participate in predicate P) how can we generate such tuples, so that they also satisfy the frequency constraint that every element in role  $P_A$  (i.e., each one of the 4 elements of type A that participates in P) appears between 3 and 5 times? (We assume all constraints are otherwise satisfiable, i.e., the above assignments were produced after a successful run of Algorithm 2.) This instance is typical of the enumeration problems that our implementation is called upon to solve.

Recall the observation from Algorithm 1: We can address any satisfiable frequency constraints by just ensuring that sub-tuples from roles in frequency constraints are covered evenly (i.e., two sub-tuples of the same role range are used either the same number of times or with a difference of one). In our example, each of the 4 elements of role  $P_A$  should be used either 3 or 4 times in the 15 total tuples.

The main workhorse of our data generation is an algorithm for addressing the above requirements for generating pairs of values from two sets. Specifically, for two sets of values, the enumeration algorithm produces combinations of the values with maximum and even coverage (i.e., all values from the two sets are used in tuples, and two values from the same set are either used the same number of times or with a difference of one) and the combinations also have maximum and even coverage (i.e., no combination repeats).

The enumeration algorithm, given input sets  $S$  and  $T$  and a desired number of output pairs  $n$  is below. The algorithm produces either  $n$  pairs or the maximum number of unique pairs, whichever is lower.

```

ia := 0, ib := 0, i := 0;
while i < n and i < |S| · |T| do
  produce( $S_{ia}$ ,  $T_{ib}$ );
  ia := (ia + 1) % |S|;
  ib := (ib + 1) % |T|;
  if (i + 1) % lcm(|S|, |T|) = 0 then
    end
  i := i + 1
end

```

( $S_x$  denotes the  $x$ -th element of the set, % is the modulo operator,  $|S|$  denotes the cardinality of the set, and  $lcm$  is the least common multiple of two numbers.)

In other words, the enumeration algorithm cyclically traverses both sets at the same time, matching their elements one-by-one to create pairs. When the pairs are about to start repeating (after every  $l$ -th iteration, where  $l$  is the least common multiple of the two set sizes) one set is *cyclically shifted* by one. This guarantees that all elements are covered equally, as the indices of both sets go over all their elements in a round robin fashion, and all possible combinations are exhausted. To show this rigorously, consider the case of the algorithm producing a pair that has also appeared before. The two instances of the pair are separated by a distance which is a multiple of  $|S|$ :  $c_1 \cdot |S|$  (since the first set is just always traversed cyclically) as well as equal to  $c_2 \cdot |T| + k$ , where  $k$  is the number of cyclic shifts in the second set that have occurred in the meantime. Thus,  $c_1 \cdot |S| = c_2 \cdot |T| + k$ . Consider now the greatest common divisor,  $g$ , of  $|S|$  and  $|T|$ . Since  $g$  divides the left hand side of the equation, it also divides the right hand side, and since it divides  $|T|$  it also needs to divide  $k$ . Thus, pairs repeat only after a multiple-of- $g$  number of shifts have occurred. Since, however, every shift happens after  $l$  pairs are produced, where  $l$  is the least common multiple of  $|S|$  and  $|T|$ , the distance between repeated pairs is (at least) a multiple of  $g \cdot l$ . From arithmetic, the product of the greatest common divisor and the least common multiple of two numbers is the numbers' product, hence pairs do not repeat until all  $|S| \cdot |T|$  possible pairs have been exhausted, which is prevented by the loop condition.

For tuples of higher arity (triples, quadruples, etc.), we can invoke this enumeration algorithm repeatedly with one of the input sets being the output of the previous step. For instance, in our example, we can first generate 15 pairs of the 4 elements of type A and the 6 elements of type B, and then combine these pairs with the 2 elements of type C. If set A contains elements  $I, II, III, IV$ , set B contains elements  $a, b, c, d, e, f$ , and set C contains elements 0,1, then the triples produced will be (in order):  $[I, a, 0]$ ,  $[II, b, 1]$ ,  $[III, c, 0]$ ,  $[IV, d, 1]$ ,  $[I, e, 0]$ ,  $[II, f, 1]$ ,  $[III, a, 0]$ ,  $[IV, b, 1]$ ,  $[I, c, 0]$ ,  $[II, d, 0]$ ,  $[III, e, 1]$ ,  $[IV, f, 0]$ ,  $[I, b, 1]$ ,  $[II, c, 0]$ ,  $[III, d, 1]$ .

The cost of this data generation algorithm is clearly linear in the amount of data generated (assuming constant-time arithmetic). In this sense, the algorithm is optimally efficient. An important note is that the notion of efficiency differs in this section from that in previous ones. Previous sections examined possible solutions from the perspective of the satisfiability of ORM diagrams. This section deals with producing test data that satisfy the diagrams. Which problem we choose affects our asymptotic analysis arguments. When producing test data, the output can be exponentially larger than the input, so it makes sense to express run-time complexity relative to the size of the output and not, as in previous sections, of the input.

We can now see in detail how to produce test data that satisfy all constraints, by combining Algorithm 1 of the previous section with our value enumeration approach. (Recall that Algorithm 1 just proves that there exist data that satisfy the constraints, whereas here we are interested in how exactly to generate such data without excessive search.)

- Data generation for the first two steps of Algorithm 1 is easy. We only need to “create  $a$  new objects” (possibly from a range specified by a value constraint) and “pick  $b$  objects so that they include the objects in the subtype”. These tasks are clearly of linear complexity in the amount of values generated.
- The next step of Algorithm 1 requires populating for each type A all roles  $R_A, S_A, \dots, V_A$ , so that, if type A is not independent, all its elements participate in some role. We can satisfy this by picking an ordering of the roles  $R_A, S_A, \dots, V_A$  that A

plays. We then populate  $R_A$  with the first  $r_a$  elements of A,  $S_A$  with the next  $s_a$  elements, etc. Since  $a \leq r_a + s_a + \dots + v_a$  (for a non-independent type) this process will exhaust the elements of A, possibly before exhausting the roles. Repeat from the first element of A in round-robin fashion, until all roles are populated. Again, this data generation process is linear in the amount of generated data.

- The next two steps of Algorithm 1 require producing combinations of existing sets of values, so that uniqueness and frequency constraints are satisfied. First, for each frequency constraint over roles  $R_A, R_B, \dots, R_K$  we need to generate  $r_{ab\dots k}$  unique tuples from combining the  $r_a, r_b, \dots, r_k$  values of the component sets. We do so by repeatedly employing our enumeration algorithm, above. (The algorithm enumerates pairs of values from two sets at a time, so conceptually it needs to be repeatedly applied to its previous output and a new set. In practice, the repeated application can be optimized into a single loop over all input sets.) Second, we need to generate  $r$  unique tuples for each predicate R in such a way that frequency constraints are satisfied. Again, we employ our enumeration algorithms on two sets at a time, where the sets either contain values from a role that is not under a frequency constraint, or sub-tuples from roles in frequency constraints. Our enumeration algorithm guarantees the satisfaction of the frequency constraints, because of its even coverage. Note that the final invocation of the enumeration algorithm is guaranteed to produce exactly the desired number,  $r$ , of tuples since for every predicate R we have the constraint  $r \leq r_{ab\dots k} \cdot r_{lm\dots p} \cdot \dots \cdot r_n$  (i.e.,  $r$  does not exceed the number of combinations of unique values).

In summary, our approach efficiently generates test data that satisfy a given set of  $\text{ORM}^-$  constraints.

## 5 ORM Diagrams in Practice

To validate the usability of  $\text{ORM}^-$ , we examined the use of ORM in practice by looking at example diagrams provided by our industrial partner, LogicBlox Inc. The diagrams model a range of consulting, internal use, and benchmarking projects: a retail prediction application, a cell bandwidth prediction application, a standard database benchmark (TPC-H), and various internal tools. These diagrams are the whole set of ORM data that LogicBlox has available: no selection or other filtering took place. There are, however, threats to the representativeness of the data: the majority of the diagrams were developed by a single engineer, and all diagrams model implementations that use the same database back-end. Nevertheless, given the variety of diagrams and domains, as well as the input we received in personal communication with multiple developers, we believe our study and findings to be highly valid.

Table 2 summarizes the elements of our ORM diagrams. There are 5 distinct projects, each with one or more documents, which in turn contain one or more diagrams. Each row represents a document and clusters of rows represent a project. For instance, the first project contains a single document with 7 diagrams, while the second project contains 9 documents, the largest of which has 46 diagrams. Diagrams in the same document are just different views: they can share entities and predicates, while omitting detail that is unnecessary for the aspect currently being modeled.

Overall, we analyzed 168 diagrams with about 1800 constraints. The vast majority of constraints are in the  $\text{ORM}^-$  subset. Only 24 constraints are not supported by

**Table 2** Elements found in ORM diagrams in practice. Dgr:# diagrams in the document. Ent:# entity types. Val ctr:# value constraints (on entity and value types, separately). IsA:# entity types that have a super type. Indp:# independent entity types. Val:# value types. Ent ref:# entity references. Rel cont:# relation container entities. Rel:# relations (predicates). Mnd:# mandatory constraints. Unq int:# internal uniqueness constraints. Frq:# frequency constraints. Unq ext:# external uniqueness constraints. Eq:# equality constraints. Ring AC:# acyclic constraints (ring). Usr:# user-defined constraints, not counting comment-only constraints. We did not see any other ORM predefined constraints such as irreflexive, intransitive, asymmetric, etc.

	Dgr	Ent	Val ctr	IsA	Indp	Val	Val ctr	Ent ref	Rel cont	Rel	Mnd	Unq int	Frq	Unq ext	Eq	Ring AC	Usr
co	7	21	0	0	0	17	2	0	0	57	0	56	0	0	0	0	4
la	11	19	3	2	3	0	0	8	0	71	7	64	0	0	0	0	0
lb	11	31	12	4	12	15	0	0	0	55	6	55	0	0	1	3	0
ld	7	6	0	1	0	1	0	15	1	76	0	73	0	0	0	0	0
le	1	3	0	1	0	0	0	2	0	5	4	4	0	1	0	0	0
ll	12	24	7	11	6	1	0	10	0	50	20	49	0	0	0	0	0
lol	12	32	6	10	5	0	0	12	0	112	0	111	0	0	0	0	0
lor	4	23	2	14	2	0	0	5	0	44	6	45	0	0	0	1	0
ls	1	4	0	0	0	0	0	5	0	17	0	17	0	0	0	0	0
lu	46	94	48	31	47	0	0	12	1	528	1	523	0	0	0	0	1
mc	1	2	0	0	0	3	0	0	0	4	0	4	0	0	0	0	0
md	1	1	0	0	0	0	0	2	0	4	0	4	0	0	0	0	0
mf	1	5	1	1	0	1	0	0	0	4	8	4	0	0	0	0	0
mm	3	9	0	2	0	5	0	0	0	16	8	16	1	1	1	2	1
m	1	2	0	0	0	1	0	0	0	3	0	3	0	0	0	0	1
mo	5	22	3	13	4	0	0	13	0	65	6	66	0	0	0	1	0
mp	2	2	1	0	1	0	0	5	0	20	0	15	0	0	0	0	0
mre	3	4	0	0	0	11	1	0	6	25	10	15	0	0	0	0	0
mru	2	2	0	0	0	0	0	7	0	18	0	14	0	0	0	0	0
mu	9	8	6	0	0	0	0	27	0	84	0	48	0	0	0	0	0
mv	3	9	7	0	7	0	0	10	0	45	0	35	0	0	0	0	0
ro	13	24	0	0	0	11	1	13	0	83	0	75	0	0	0	0	5
sb	1	5	4	0	4	2	0	0	1	3	0	3	0	0	0	0	1
sl	1	3	0	0	0	2	0	0	0	4	0	4	0	0	0	0	0
st	10	11	1	0	1	22	0	0	1	64	35	63	0	0	0	0	0
total	168	366	101	90	92	92	4	146	10	1457	111	1366	1	2	2	7	13

ORM<sup>-</sup>. Of these, 13 are *user-defined* constraints: they add arbitrary (query) code to the high-level diagram. This means that the developer adds consistency code (e.g., “this predicate is the join of two others over these roles”) to be generated together with normal consistency code that the ORM editor produces (e.g., for uniqueness, frequency, value, or mandatory constraints). Checking the satisfiability of such user-defined constraints seems unlikely, as this is an undecidable problem for the general query language. Thus, we do not believe that there is significantly more benefit to get over what ORM<sup>-</sup> already achieves. A promising direction may be acyclicity ring constraints, which comprise 7 of the remaining 11 constraints that ORM<sup>-</sup> does not handle.

Our analysis did not find interesting errors in the ORM documents. There were 19 cases of unsatisfiability, but all were relatively benign: an entity or value was not connected to any role or supertype. Nevertheless, finding no errors is hardly surprising for this dataset since the ORM tool generates consistency checking code (in Datalog) automatically. Therefore, mistakes in the specification are overwhelmingly likely to be caught when the database is populated with real data, and the diagrams we examined have produced databases that are in real use.

The run-time of our satisfiability check was negligible. Satisfiability checking is done on a per-document basis and on a 2 GHz AMD Athlon 64 X2 machine checking the largest document (with 528 predicates) takes about 35ms. Thus, satisfiability checking is fast enough to be done interactively inside the ORM editor, while the user is editing diagrams. Test data generation takes time proportional to the output, however, and thus needs to be done offline.

The tendency we observe is for developers to prefer to encode complex constraints in code, rather than using the semantically complex constraints of the modeling language. The modeling language is instead used for easier constraints that developers feel very comfortable with. One possible motivation may be exactly the lack of good tools for consistency checking. Since the ORM editor works as a code generator, errors in an ORM diagram stay undetected until later, when the database is populated. Going back and fixing the diagram is costly at this stage: it requires re-generating the database schema and the consistency code, re-importing data, and possibly dealing with the re-integration of hand-written code that was implemented against the faulty database. It is, thus, possible that tools like ours will encourage the use of more complex constraints.

## 6 Related Work

There are several earlier approaches to automatic example data generation. Compared to other work on checking the satisfiability of ORM, ours is distinguished by its completeness (for ORM<sup>-</sup>) and emphasis on practicality. Jarrar and Heymans [15] use the fact that “no complete satisfiability checker is known for ORM” to motivate the introduction of 9 *unsatisfiability* patterns that capture conceptual modeling mistakes. The presence of these patterns can be checked with a program search. Heymans [13] describes a translation of a subset of ORM to a formalism with an ExpTime decision procedure. Keet [17] reduces another subset of ORM (emphasizing ring constraints) to Description Logic languages. These reductions follow the general pattern of the standard translation of ORM to logic, and are not intended for efficient execution. We are working on the problem from the opposite end, and it may be interesting to try to

extend  $\text{ORM}^-$  when greater expressiveness is needed, rather than to try to restrict larger fragments.

Wilmore and Embury [30] propose an *intensional* approach to database testing, which is closely related in spirit to our test of constraint satisfiability. Nevertheless, they do not address the fundamental undecidability of the language they support, so their technique is heuristic, without a clear understanding of its expressiveness and limitations.

Other database testing techniques, such as that of Deng et al. [4], concentrate on producing test databases with arbitrary values or after pruning illegal values through heuristics. A more formal approach is followed by Neufeld et al. [23]: constraints are translated into logic and a generator specification is derived semi-automatically (user control may be needed) from the logical formula. This approach is, again, heuristic, in that it may not be able to produce appropriate data. Nevertheless, it offers a way to handle complex constraints in a unifying framework.

There is significant work in databases on the satisfiability of different kinds of constraints. Fan and Libkin [7] address XML document specifications with DTDs and integrity constraints. Calvanese and Lenzerini study the interaction between subtype and “cardinality” constraints [3]. (The latter correspond roughly to what we called “frequency” constraints in  $\text{ORM}$ .) Both of the above pieces of work are interesting because they use integer constraints for some of their formulations. Nevertheless, the kinds of constraints supported are significantly different from the  $\text{ORM}$  constraints we tackled. For instance, a key element in the Calvanese and Lenzerini work is that a subtype can refine the cardinality constraint of its supertype. E.g., we can specify that objects of type *Person* appear at least 3 times in some table, yet objects of subtype *Man* can appear at most 5 times.

In terms of consistency reasoning for modeling languages, UML is a common focus. UML is really a collection of specification languages that cover many aspects of object-oriented development—from problem specification over class diagrams, to module interactions, to software deployment. Egyed presents a fast technique for maintaining the consistency of different UML diagrams [5]. Recent work also proposes actions for fixing inconsistencies [22, 6]. Nevertheless, neither the core problem we target (inconsistency under *all* inputs) nor the kinds of constraints we support are closely related to that work.

Beradi et al. follow an approach similar to ours, by restricting the expressiveness of UML class diagrams to a decidable subset [2]. I.e., they also exclude user-constraints, which can contain arbitrarily complex expressions in full first-order logic. They prove their resulting language  $\text{ExpTime}$ -complete and present a mapping to description logic.  $\text{ORM}^-$ , in contrast, is decidable in polynomial time.

In addition to an  $\text{ExpTime}$  subset of UML class diagrams, Kaneiwa and Satoh also propose restrictions that enable consistency checking in  $\text{P}$ ,  $\text{NP}$ , and  $\text{PSPACE}$  [16]. Maraee and Balaban [21] discuss the efficient satisfiability of UML class diagrams by translating them into linear inequalities, much like our approach. But it remains unclear what proportion of UML class diagrams in practice can be handled with these subsets. It would be interesting to apply our algorithms to UML class diagrams and determine how many UML class diagrams our algorithms can handle in industrial practice.

Lenzerini and Nobili show that polynomial-time consistency checking of database specifications that include cardinality constraints is possible [19]. They work on Entity-Relationship Model (ERM) diagrams that include a form of frequency constraint but not the additional constraints we support, such as mandatory and subset. We extend

this work by giving fast algorithms for both consistency checking and test case generation. We also provide empirical results on which proportion of real-world database specifications our technique covers.

Alloy has been applied successfully in several other domains. Warren et al. compile software architecture specifications to Alloy to check their consistency before performing a proposed software architecture reconfiguration [29]. Khurshid and Jackson discovered bugs in the Intentional Naming System (INS) by modeling and checking it with Alloy [18]. Finally, several translation schemes from Java to Alloy have been proposed: Taghdiri's [27] is a good example.

## 7 Conclusions

Producing example data that meet a given specification enables interactive checking of a data model and testing of the implementation with large amounts of consistent, automatically produced data. In this article we presented an important advance in the direction of generating test data from modeling diagrams. We proposed a subset,  $\text{ORM}^-$ , of the popular ORM modeling language. For models written in this subset, we presented algorithms for determining the satisfiability of the model, as well as for producing data conforming to the model. An analysis of practical uses of ORM diagrams confirmed that the constraints we target are sufficient for covering the vast majority of real-world uses.

We believe that our results likely generalize to other settings, as the constraints in  $\text{ORM}^-$  are fairly standard. Much of the power of our approach comes from our support of multi-role frequency and uniqueness constraints. Interestingly, these constraints are also useful for providing input to the test data generation process. For instance, a user wishing to produce a set of random data is likely to use frequency constraints to specify his/her expectation as to how many times each value is used in a table. Thus, our constraint language can be viewed both as a way to specify a data schema (for later real use) and as a way to specify desired properties of an automatically generated sample dataset. This, combined with its other desirable properties, makes our approach a strong candidate for wide practical adoption.

**Acknowledgements** Molham Aref, Wes Hunter, and David Zook of LogicBlox gave valuable help on ORM constraint selection and ORM editor usage, as well as helpful comments on this article. Pete Manolios, Gayatri Subramanian, and Daron Vroon encouraged us to use CoBaSa and helped us with it. This work was funded by the National Science Foundation under grant NSF:CCR-0735267 and by LogicBlox Inc. This article is an extended version of an ASE'07 conference paper [26].

## References

1. F. A. Aloul, A. Ramani, I. L. Markov, and K. A. Sakallah. PBS: A backtrack search pseudo-boolean solver and optimizer. In *Proc. 5th International Symposium on the Theory and Applications of Satisfiability Testing (SAT)*, pages 346–353, May 2002.
2. D. Berardi, D. Calvanese, and G. D. Giacomo. Reasoning on UML class diagrams. *Artificial Intelligence*, 168:70–118, Oct. 2005.
3. D. Calvanese and M. Lenzerini. On the interaction between isa and cardinality constraints. In *Proc. 10th International Conference on Data Engineering (ICDE)*, pages 204–213. IEEE, Feb. 1994.

4. Y. Deng, P. G. Frankl, and D. Chays. Testing database transactions with AGENDA. In *Proc. 27th International Conference on Software Engineering (ICSE)*, pages 78–87. ACM, May 2005.
5. A. Egyed. Instant consistency checking for the UML. In *Proc. 28th International Conference on Software Engineering (ICSE)*, pages 381–390. ACM, May 2006.
6. A. Egyed. Fixing inconsistencies in UML design models. In *Proc. 29th International Conference on Software Engineering (ICSE)*, pages 292–301. IEEE, May 2007.
7. W. Fan and L. Libkin. On xml integrity constraints in the presence of dtlds. *J. ACM*, 49(3):368–406, 2002.
8. J. Grant and J. Minker. Inferences for numerical dependencies. *Theor. Comput. Sci.*, 41(2–3):271–287, 1985.
9. T. A. Halpin. A fact-oriented approach to schema transformation. In *Proc. 3rd Symposium on Mathematical Fundamentals of Database and Knowledge Base Systems (MFDBS)*, pages 342–356. Springer, 1991.
10. T. A. Halpin. Object-role modeling (ORM/NIAM). In *Handbook on Architectures of Information Systems*. Springer, 1998.
11. T. A. Halpin. *Information Modeling and Relational Databases*. Morgan Kaufmann, 2001.
12. T. A. Halpin and H. A. Proper. Database schema transformation and optimization. In *Proc. 14th International Conference on Object-Oriented and Entity-Relationship Modelling (OOER)*, pages 191–203. Springer, Dec. 1995.
13. S. Heynmans. *Decidable Open Answer Set Programming*. PhD thesis, Vrije Universiteit Brussel, Department of Computer Science, Feb. 2006.
14. D. Jackson. *Software abstractions: Logic, language, and analysis*. MIT Press, 2006.
15. M. Jarrar and S. Heymans. Unsatisfiability reasoning in ORM conceptual schemes. In *Proc. International Conference on Semantics of a Networked World (ICSNW)*. Springer, Mar. 2005.
16. K. Kaneiwa and K. Satoh. Consistency checking algorithms for restricted UML class diagrams. In *Proc. 4th International Symposium on the Foundations of Information and Knowledge Systems (FoIKS)*, pages 219–239, Feb. 2006.
17. C. M. Keet. Prospects for and issues with mapping the object-role modeling language into DL<sup>R</sup>ifd. In *Proc. 20th International Workshop on Description Logics*, June 2007.
18. S. Khurshid and D. Jackson. Exploring the design of an intentional naming scheme with an automatic constraint analyzer. In *Proc. 15th IEEE International Conference on Automated Software Engineering (ASE)*, pages 13–22. IEEE, Sept. 2000.
19. M. Lenzerini and P. Nobili. On the satisfiability of dependency constraints in entity-relationship schemata. In *Proc. 13th International Conference on Very Large Data Bases (VLDB)*, pages 147–154, Sept. 1987.
20. P. Manolios, G. Subramanian, and D. Vroon. Automating component-based system assembly. In *Proc. ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. ACM, July 2007.
21. A. Maraee and M. Balaban. Efficient reasoning about finite satisfiability of uml class diagrams with constrained generalization sets. In *Proc. 3rd European Conference on Model-Driven Architecture*, 2007.
22. C. Nentwich, W. Emmerich, and A. Finkelstein. Consistency management with repair actions. In *Proc. 25th International Conference on Software Engineering (ICSE)*, pages 455–464. IEEE, May 2003.
23. A. Neufeld, G. Moerkotte, and P. C. Lockemann. Generating consistent test data: Restricting the search space by a generator formula. *VLDB Journal*, 2:173–213, 1993.
24. G. M. Nijssen and T. A. Halpin. *Conceptual schema and relational database design: A fact oriented approach*. Prentice-Hall, 1989.
25. I. Shlyakhter, R. Seater, D. Jackson, M. Sridharan, and M. Taghdiri. Debugging overconstrained declarative models using unsatisfiable cores. In *Proc. 18th IEEE International Conference on Automated Software Engineering (ASE)*, pages 94–105. IEEE, Oct. 2003.
26. Y. Smaragdakis, C. Csallner, and R. Subramanian. Scalable automatic test data generation from modeling diagrams. In *Proc. 22nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 4–13. ACM, Nov. 2007.
27. M. Taghdiri. Inferring specifications to detect errors in code. In *Proc. 19th IEEE International Conference on Automated Software Engineering (ASE)*, pages 144–153. IEEE, Sept. 2004.
28. P. van Bommel, A. ter Hofstede, and T. van der Weide. Semantics and verification of object-role models. *Information Systems*, 16(5):471–495, Oct. 1991.



- 
29. I. Warren, J. Sun, S. Krishnamohan, and T. Weerasinghe. An automated formal approach to managing dynamic reconfiguration. In *Proc. 21st IEEE International Conference on Automated Software Engineering (ASE)*, pages 37–46. IEEE, Sept. 2006.
  30. D. Willmor and S. M. Embury. An intensional approach to the specification of test cases for database applications. In *Proc. 28th International Conference on Software Engineering (ICSE)*, pages 102–111. ACM, May 2006.