

# A Generic Model Decomposition Technique and its Application to the Eclipse Modeling Framework

Qin Ma<sup>12</sup>, Pierre Kelsen<sup>12</sup>, Christian Glodt<sup>1</sup>

<sup>1</sup> FSTC, University of Luxembourg

<sup>2</sup> SnT, University of Luxembourg

Received: date / Revised version: date

**Abstract** Model-driven software development aims at easing the process of software development by using models as primary artifacts. Although less complex than the real systems they are based on, models tend to be complex nevertheless, thus making the task of handling them non-trivial in many cases. In this paper we propose a generic model decomposition technique to facilitate model management by decomposing complex models into smaller sub-models that conform to the same metamodel as the original model. The technique is based upon a formal foundation that consists of a formal capturing of the concepts of models, metamodels, and model conformance; a formal constraint language based on EssentialOCL; and a set of formally proved properties of the technique. We organize the decomposed sub-models in a mathematical structure as a lattice, and design a linear-time algorithm for constructing this decomposition.

The generic model decomposition technique is applied to the Eclipse Modeling Framework (EMF) and the result is used to build a solution to a specific model comprehension problem of Ecore models based upon model pruning. We report two case studies of the model comprehension method: one in BPMN and the other in fUML.

---

**Key words** MDE, EMF, Model decomposition, Model comprehension, Linear-time algorithm, Sub-model lattice, OCL, EssentialOCL, BPMN, fUML.

## 1 Introduction

In model-driven software development models are the primary artifacts. Typically several models are used to describe the different concerns of a system. One of the main motivations for using models is the problem of dealing with the complexity of real systems: because models

represent abstractions of a system, they are typically less complex than the systems they represent.

Nevertheless models for real systems can be complex themselves and thus may require aids for facilitating human comprehension. The problem of understanding complex models is at the heart of this paper. We propose to rely on a model decomposition technique that subdivides models into smaller relevant sub-models to aid in their comprehension.

An example of a concrete application scenario is the following: when trying to understand a large model, one starts with a subset of model elements that one is interested in (such as the concept of Class in the UML metamodel). Our method allows to construct a small sub-model of the initial model that contains all entities of interest and that conforms to the original metamodel (in the case of UML the original metamodel would be MOF). The latter condition ensures that the sub-model can be viewed in the same way as the original model and that it has a well-defined semantics. The smaller size (compared to the original model) should facilitate comprehension.

Instead of providing a particular solution to the specific comprehension problem, we first study a more general model decomposition problem in an abstract setting. More specifically, the decomposition problem deals with the following: given a metamodel and a model conforming to the metamodel, how does one derive all the conformant sub-models and how are the sub-models related to

each other? Based upon a well established formal foundation that consists of a formal capturing of the concepts of models, metamodels, and model conformance, we propose a linear time algorithm to build the decomposition hierarchy of a model from which all the conformant sub-models can be constructed in a straightforward manner. We prove formally the correctness of the algorithm, and present the mathematical structure of these conformant sub-models as a lattice. A lattice (drawn upon order theory) is a partially ordered set in which any two elements have a least upper bound and a greatest lower bound. The original model is the greatest element in the lattice at the top and the empty sub-model is the least element in the lattice at the bottom.

The solution to the model decomposition problem is generic in two ways: first, the technique can be employed to decompose any model conforming to any metamodel; second, it considers the collection of relevant sub-models in its totality rather than a single sub-model.

The generic solution to the model decomposition problem is then customized to target a solution to the specific model comprehension problem we discussed earlier. This is achieved in two steps. First, we instantiate the decomposition solution to work within a concrete environment: Eclipse Modeling Framework (EMF) and EssentialOCL [27]. Second, the smallest sub-model in the decomposition hierarchy that contains all the concepts of interest is selected as an answer to the comprehension problem.

We implemented the model comprehension approach in an Ecore model comprehension tool [1] and carried out two case studies for validation. The first case study takes place in the context of BPMN [26] (business process model and notation), where we try to comprehend the *Gateway* concept. The second case study is about understanding the *Class* concept and *Namespace* concept in fUML [25] (describing a subset of executable UML models). In the first case, the size of the model for comprehension decreases by 93% (in terms of the number of model elements). In the second case, the size decreases by 62% and 89% (respectively for *Class* and *Namespace*).

*Roadmap* Metamodels and models are specified in practice with concrete tools. An abstraction layer is derived on top of all these concrete environments to retain only concepts that are relevant for model decomposition. We present the abstraction layer in Section 2 with formal definitions of models, metamodels and model conformance, and describe a technique for model decomposition in Section 3. We prove that the sub-models are conformant to the same metamodel as the original model and can be organized in a mathematical structure called the lattice of sub-models. Being an abstract metamodeling and modeling environment, the model decomposition technique defined at this level of abstraction is generally usable in any concrete environment. In this paper, we demonstrate its usage in the Eclipse Modeling Framework (EMF): Section 4 lays the ground for the usage of

the model decomposition technique by providing a short summary of EMF; Section 5 presents a concrete yet formal constraint language, called CoreOCL (a core of EssentialOCL), for the specification of invariants attached to metamodels; and Section 6 discusses the actual steps involved, namely first establishing an alignment between EMF and the abstraction layer, then customizing the general model decomposition algorithm for EMF model decomposition, and finally discussing the soundness of the customized algorithm. We report the application of the model decomposition technique in a concrete scenario for Ecore model comprehension in Section 7 and present the results of two case studies. We evaluate our approach in Section 8 and discuss related work in Section 9. Finally, we present concluding remarks and point out some interesting directions for future work in the last section.

*Extension statement* We have presented a previous version of the model decomposition technique in [22] at the Fourteenth International Conference on Fundamental Approaches to Software Engineering (FASE 2011). Thanks to the valuable comments and suggestions from the anonymous referees of FASE 2011, we here are happy to present an extended and improved version of the model decomposition technique completed with all its formal treatments. Notably, three new pieces of work have been carried out in this extension, namely, a new formal constraint language called CoreOCL (reported in Section 5),

the usage of the model decomposition technique in the Eclipse Modeling Framework in general (reported in Section 4 and 6), and in the Ecore model comprehension example in particular (reported in Section 7.1).

## 2 Models and Metamodels

OMG puts forward MOF [24] as a platform independent framework for defining, manipulating and integrating metamodels and models. The MOF standard has been implemented in various concrete environments to support (meta-)modeling in real life. Examples of such concrete environments include EMF (the Eclipse Modeling Framework) <sup>1</sup>, Kermeta (Kernel Metamodeling) <sup>2</sup>, AMMS (ATLAS Model Management Architecture) <sup>3</sup>, and MOFLON <sup>4</sup>. Inspired by MOF, we abstract from the metamodeling and modeling concepts present in mainstream concrete tools and formalize a set of definitions of models, metamodels, and model conformance, where only model decomposition relevant information is kept.

The following notational conventions will be used:

1. For any tuple  $p$ , we use  $\text{fst}(p)$  to denote its first element,  $\text{snd}(p)$  to denote its second element, and  $\text{trd}(p)$  to denote its third element.
2. For any set  $s$ , we use  $\#s$  to denote its cardinality.

<sup>1</sup> <http://www.eclipse.org/modeling/emf/>

<sup>2</sup> <http://www.kermeta.org/>

<sup>3</sup> <http://wiki.eclipse.org/AMMA>

<sup>4</sup> <http://www.moflon.org/>

### 2.1 Metamodels

A metamodel defines (the abstract syntax of) a language for expressing models to be decomposed. A metamodel consists of the following parts: a finite set of metaclasses; a finite set of associations between metaclasses; and a finite inheritance relation between metaclasses. Moreover, a set of invariants may be specified in the contexts of metaclasses as additional well-formedness rules imposed upon models.

**Definition 1 (Metamodel)** *A metamodel is defined by a tuple  $\mathbb{M} = (\mathbb{N}, \mathbb{A}, \mathbb{H}, \mathbb{I}nv, \mathbf{s}, \mathbf{t}, \mu_s, \mu_t, \text{ctx})$  where:*

- $\mathbb{N}$  is the set of metaclasses, and  $n \in \mathbb{N}$  ranges over it.
- $\mathbb{A}$  is the set of (directed) associations between metaclasses and  $a \in \mathbb{A}$  ranges over it.
- $\mathbb{H} \subseteq \mathbb{N} \times \mathbb{N}$  denotes the inheritance relation among metaclasses. The transitive closure of  $\mathbb{H}$  should be irreflexive. Namely, a metaclass cannot inherit, directly or indirectly, from itself. The subtyping relation ( $\preceq$ ) between metaclasses is defined as the reflexive and transitive closure of  $\mathbb{H}$ .
- $\mathbb{I}nv$  is the set of invariants and  $inv \in \mathbb{I}nv$  ranges over it.
- $\mathbf{s} : \mathbb{A} \rightarrow \mathbb{N}$  and  $\mathbf{t} : \mathbb{A} \rightarrow \mathbb{N}$  are two functions from associations to metaclasses. They specify respectively the types of the source and target ends of an association.
- $\mu_s : \mathbb{A} \rightarrow \mathbb{Z}$  and  $\mu_t : \mathbb{A} \rightarrow \mathbb{Z}$  are two functions from associations to well-formed multiplicities. They specify

respectively the multiplicities of the source and target ends of an association.  $\Sigma \subseteq \text{Nat} \times \{\text{Int}^+ \cup \{*\}\}$  captures the set of multiplicities, where  $\text{Nat}$  is the set of natural numbers (i.e. non-negative integers) and  $\text{Int}^+$  is the set of positive integers. Multiplicities are ranged over by  $\sigma \in \Sigma$ . We assume that all the multiplicities in  $\Sigma$  are well-formed, namely  $\text{snd}(\sigma) = *$  or  $\text{fst}(\sigma) \leq \text{snd}(\sigma)$ . A well-formed multiplicity defines an inclusive interval from a lower bound to an upper bound. An asterisk  $*$  is used for denoting an unlimited upper bound.

- $\text{ctx} : \text{Inv} \rightarrow \mathbb{N}$  is a function from invariants to metaclasses. It specifies the context metaclass of an invariant.

In addition to the context, an invariant also comes with a body which, in practice, is either documented in natural language, or specified as expressions of a constraint language such as EssentialOCL [27], or implemented in terms of code in a programming language such as Java.

Note that we omit the following concepts in meta-models on purpose:

- *Attributes* are treated as a special kind of associations, where one end of the connected metaclasses is a data type. Consequently, data types (including both primitive types e.g., integers and booleans and enumeration types) are all metaclasses.

- *Composition* or *containment* relations are treated as a special kind of associations with extra constraints:
  1. the source multiplicity of a composition association is (0..1);
  2. a metaclass instance cannot be the target of more than one composition link (i.e., instantiations of composition associations);
  3. and a metaclass instance cannot be connected to itself via a sequence of composition links.

Moreover, we have also omitted the concept of operations in metaclasses because differently from associations (or attributes or references), operations only exist at the level of the metamodel. Side effect free helper operations that are used for invariant specifications can be defined directly with the corresponding invariants when needed.

## 2.2 Models

A model is built by instantiating the metaclasses and associations of the metamodel.

**Definition 2 (Model)** *A model is defined by a tuple  $M = (\mathbb{M}, \mathbb{N}, \mathbb{A}, \tau, \text{src}, \text{tgt})$  where:*

- $\mathbb{M}$  is the metamodel in which the model is expressed.
- $\mathbb{N}$  is the set of metaclass instantiations, and  $n \in \mathbb{N}$  ranges over it. They are often referred to as *instances*.
- $\mathbb{A}$  is the set of association instantiations, and  $a \in \mathbb{A}$  ranges over it. They are often referred to as *links*.

–  $\tau : (\mathbf{N} \rightarrow \mathbf{N}) \cup (\mathbf{A} \rightarrow \mathbf{A})$  is the typing function.

It records the type information of the instances and links in the model, i.e., from which metaclasses or associations of the metamodel  $\mathbb{M}$  they are instantiated.

–  $\text{src} : \mathbf{A} \rightarrow \mathbf{N}$  and  $\text{tgt} : \mathbf{A} \rightarrow \mathbf{N}$  are two functions from links to instances. They specify respectively the source and target ends of a link.

### 2.3 Model Conformance

Not all models following the definition above are valid, or “conform to” the metamodel. A valid model should satisfy all the typing, multiplicity, and extra well-formedness invariants captured in the metamodel.

**Definition 3 (Model conformance)** We say a model  $\mathbf{M} = (\mathbb{M}, \mathbf{N}, \mathbf{A}, \tau, \text{src}, \text{tgt})$  conforms to its metamodel  $\mathbb{M}$  or is valid when the following conditions are met:

1. *Typing condition: links only connect instances whose types are compatible with (i.e., subtypes of) the metaclasses specified for the corresponding association ends.*

Namely,  $\forall a \in \mathbf{A}$ , we have both  $\tau(\text{src}(a)) \preceq s(\tau(a))$  and  $\tau(\text{tgt}(a)) \preceq t(\tau(a))$ .

2. *Multiplicity condition: the numbers of links must fall in the ranges specified by the multiplicities of the corresponding associations. We shall consider both the source and target multiplicities.*

(a)  $\forall n \in \mathbf{N}, a \in \mathbf{A}$ , if  $\tau(n) \preceq s(a)$ , then let  $k = \#\{a \in \mathbf{A} \mid \tau(a) = a \text{ and } \text{tgt}(a) = n\}$  (i.e., the number of  $a$ -typed links ending at the instance  $n$  in

model  $\mathbf{M}$ ), we must have  $k \geq \text{fst}(\mu_s(a))$  (i.e., the lower bound of the source multiplicity of  $a$ ) and  $k \leq \text{snd}(\mu_s(a))$  (i.e., the upper bound of the source multiplicity of  $a$ ) in case the latter is not  $*$ ;

(b)  $\forall n \in \mathbf{N}, a \in \mathbf{A}$ , if  $\tau(n) \preceq t(a)$ , then let  $k = \#\{a \in \mathbf{A} \mid \tau(a) = a \text{ and } \text{src}(a) = n\}$  (i.e., the number of  $a$ -typed links leaving the instance  $n$  in model  $\mathbf{M}$ ), we must have  $k \geq \text{fst}(\mu_t(a))$  (i.e., the lower bound of the target multiplicity of  $a$ ) and  $k \leq \text{snd}(\mu_t(a))$  (i.e., the upper bound of the target multiplicity of  $a$ ) in case the latter is not  $*$ .

3. *Invariant condition: all invariants should hold.*  $\forall \text{inv} \in \mathbb{I}\text{nv}$ ,  $\forall n \in \mathbf{N}$  where  $\tau(n) \preceq \text{ctx}(\text{inv})$ , (i.e., the type of  $n$  is compatible with the context metaclass,) the invariant  $\text{inv}$  should evaluate to true in model  $\mathbf{M}$  for the contextual instance  $n$ . (Section 5 defines the evaluation semantics of invariants written in CoreOCL.)

## 3 Model Decomposition

### 3.1 Criteria

Model decomposition starts from a model that conforms to a metamodel, and decomposes it into smaller parts. Our model decomposition technique is designed using the following as main criterion: the derived parts should be valid models conforming to the original metamodel. Achieving this goal has two main advantages:

1. the derived parts, being themselves valid models, can be comprehended on their own according to the fa-

miliar abstract syntax and semantics (if defined) of the modeling language;

2. the derived parts can be wrapped up into modules and reused in the construction of other system models, following a modular model composition paradigm such as [21].

In general, a decomposed smaller part of a model  $M$  is a sub-model whose instance set is a subset of that of  $M$ , and whose link set is a subset of the links of  $M$  restricted to the instance subset. In our model decomposition technique we consider a particular kind of sub-models, called *instance induced sub-models*, as we view models as essentially sets of instances, augmented with links among these instances. As a consequence, induced sub-models include all the links involving the instances included in the sub-models. In graph-theoretic terms this corresponds to induced subgraphs [12] (if we view models as graphs).

**Definition 4 (Instance induced sub-model)** *We say a model  $M' = (\mathbb{M}, N', A', \tau', \text{src}', \text{tgt}')$  is an instance induced sub-model of another model  $M = (\mathbb{M}, N, A, \tau, \text{src}, \text{tgt})$  if and only if:*

1.  $N' \subseteq N$ ;
2.  $A' = \{a \mid a \in A \text{ and } \text{src}(a) \in N' \text{ and } \text{tgt}(a) \in N'\}$ ;
3.  $\tau'$  is the restriction of  $\tau$  to  $N'$  and  $\text{src}'$  and  $\text{tgt}'$  are the restrictions of  $\text{src}$  and  $\text{tgt}$  to  $A'$ .

From now on, all sub-models we talk about are instance induced unless mentioned otherwise explicitly. Namely,

when constructing a sub-model, we shall only discuss the inclusion of instances and let the inclusion of links be induced by the instances included in the sub-model.

In order to make the sub-model  $M'$  also conform to  $\mathbb{M}$ , we will propose three conditions - one for the meta-model (Condition 2 below, regarding the nature of the invariants) and two conditions for the sub-model (Conditions 1 and 3). Altogether these three conditions will be sufficient to ensure conformance of the sub-model.

The starting point of our investigation is the definition of conformance (Definition 3). Three conditions must be met in order for sub-model  $M'$  to conform to metamodel  $\mathbb{M}$ .

We first tackle the most sophisticated one: the invariant condition in the conformance definition. It requires that all metamodel invariants are satisfied in sub-model  $M'$ . These invariants are known to be satisfied in model  $M$  because  $M$  conforms to the metamodel. Therefore, it is sufficient to maintain the same evaluation of these invariants in  $M'$ .

In order to achieve this goal, let us first introduce some knowledge about invariant evaluation in models. Invariants are evaluated on so called well-formed evaluation points defined as follows:

**Definition 5 (Well-formed evaluation point)** *We call a triplet of the following form  $(\text{inv}, M, n)$  an evaluation point of an invariant  $\text{inv}$  in a model  $M$  on an instance  $n$ . An evaluation point is well-formed if  $\text{inv}$  is*

defined for the metamodel of  $M$ ,  $n$  is an instance of  $M$ , and the type of  $n$  is a subtype of the context metaclass of  $\text{inv}$ .

The result of evaluating a well-formed evaluation point is determined by its *scope*, defined below:

**Definition 6 (Scope of invariant evaluation)** *The scope of a well-formed evaluation point  $(\text{inv}, M, n)$  is a sub-model of  $M$  induced by all the instances that are referenced during the evaluation of the invariant  $\text{inv}$  in model  $M$  on the contextual instance  $n$ .*

For example, CoreOCL provides three ways to reference instances in invariants:

1. referencing the contextual instance via keyword `self`;
2. following links to reference target instances;
3. using the `AllInstances` operation to reference all the instances of a given metaclass.

Section 5 discusses how scopes are constructed along the evaluation of invariants. Moreover, a formal proof is also presented in Section 5 to demonstrate that letting a sub-model contain the scopes of invariant evaluations suffices to preserve the invariants holding in the original model (see Theorem 5). More specifically, given an evaluation point  $(\text{inv}, M, n)$  and a sub-model  $M'$  of  $M$ , letting  $M'$  include all the instances of the scope of  $(\text{inv}, M, n)$  suffices to have the same result for evaluating both  $(\text{inv}, M, n)$  and  $(\text{inv}, M', n)$ , i.e., the invariant  $\text{inv}$  has the same evaluation in both  $M$  and  $M'$ .

In theory, a metamodel can specify an invariant, where an evaluation point of this invariant in a model on a contextual instance has a scope that spans the whole instance set of the model. Following the discussion above, if the model decomposition technique is expected to always preserve the evaluation of such “global” invariants in a sub-model that contains the contextual instance, the sub-model should effectively include all the instances of the original model (i.e., equivalent to the original model), hence leaves no room for effective sub-modeling.

Fortunately, most of the scopes implied by invariants in practice involves only a portion of the original model that is reachable from the contextual instance. We refer to such kind of invariants as *forward invariants*, precisely defined below:

**Definition 7 (Reachability)** *Given a model  $M$ , we say an instance  $n' \in N$  is reachable from an instance  $n \in N$ , if and only if there exists  $a_1, a_2, \dots, a_k$  in  $A$  such that  $\text{src}(a_1) = n$ ,  $\text{tgt}(a_i) = \text{src}(a_{i+1})$ , for all  $i$ ,  $1 \leq i \leq k-1$ , and  $\text{tgt}(a_k) = n'$ .*

**Definition 8 (Forward invariant)** *An invariant  $\text{inv}$  is forward if for any well-formed evaluation point  $(\text{inv}, M, n)$ , all the instances in the corresponding scope  $N_S$  are reachable from  $n$ .*

Following the discussion above, we impose a condition on the sub-model to include all the instances that are reachable from an instance that is already included, formally expressed below:

**Condition 1**  $\forall a \in A, \text{src}(a) \in N' \text{ implies } \text{tgt}(a) \in N'$ .

Meanwhile, we only allow forward invariants, as expressed in the following condition over the metamodel:

**Condition 2** *All invariants of metamodel  $\mathbb{M}$  are forward invariants.*

It is not difficult to see that Conditions 1 and 2 together guarantee that an invariant satisfied on contextual instance  $n$  in  $M$  is also satisfied in the instance induced sub-model  $M'$  on the same contextual instance, since the evaluations of the two points are indeed determined by the same scope (and all sub-models are instance induced). Consequently, the invariant condition in the conformance definition is ensured.

The multiplicity condition for conformance concerns the numbers of links in models. The number of links of a given association may decrease from the original model  $M$  to the sub-model  $M'$  in cases where an instance connected at an end of a link is not included in  $M'$ . We need to examine both the number of links ending at the instance in  $M'$ , which must agree with the source multiplicity of the corresponding association in the metamodel, and the number of links leaving the instance in  $M'$ , which must agree with the target multiplicity of the corresponding association.

The multiplicity condition for links leaving an instance of the sub-model  $M'$  is ensured by Condition 1 and the fact that  $M'$  is instance induced. The two together guarantee that the number of links leaving an

instance in  $M'$  is exactly the same as the number in  $M$ , hence remains within the range allowed by the target multiplicity.

To ensure the multiplicity condition for links ending at an instance of the sub-model, we will introduce the notion of *fragmentable links*, whose type (i.e., the corresponding association) has an unconstrained (i.e., being 0) lower bound for the source multiplicity.

**Definition 9 (Fragmentable link)** *Given a model  $M = (\mathbb{M}, N, A, \tau, \text{src}, \text{tgt})$ , a link  $a \in A$  is fragmentable if the lower bound of the source multiplicity of the corresponding association is 0, i.e.,  $\mu_s(\tau(a)) = (0, -)$ , where  $-$  represents an arbitrary upper bound.*

Fragmentable incoming links (together with the source instances) of  $M$  to instances in  $M'$  are safe to exclude but this is not the case for non-fragmentable links, which should all be included. We thus obtain the following condition on sub-model  $M'$ :

**Condition 3**  $\forall a \in A$  where  $a$  is non-fragmentable,  $\text{tgt}(a) \in N'$  implies  $\text{src}(a) \in N'$ .

Finally, the typing condition for conformance follows directly from the fact that  $M'$  is an instance induced sub-model of  $M$  and  $M$  conforms to  $\mathbb{M}$ .

Summarizing the discussion above, we thus obtain the following result:

**Theorem 1** *Given a metamodel  $\mathbb{M}$ , a model  $M$ , and an instance induced sub-model  $M'$  of  $M$ , suppose that:*

1. *model  $M$  conforms to  $\mathbb{M}$ ;*

2. model  $M'$  satisfies Conditions 1 and 3;

3. metamodel  $M$  satisfies Condition 2;

then model  $M'$  also conforms to the metamodel  $M$ .

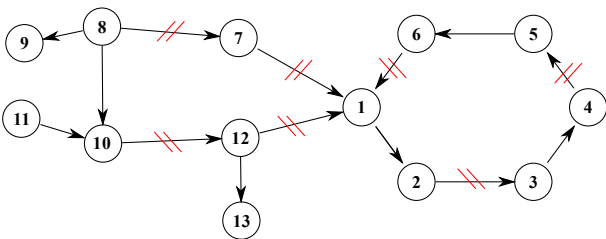
*Proof* The result follows from the discussion above.  $\square$

### 3.2 Algorithm

From hereon we shall assume that the metamodel under consideration satisfies Condition 2. In this subsection we describe an algorithm that finds, for a given model  $M$ , a partition of  $M$  such that any sub-model of  $M$  that satisfies both Conditions 1 and 3 can be derived from the partition by uniting some components of the partition.

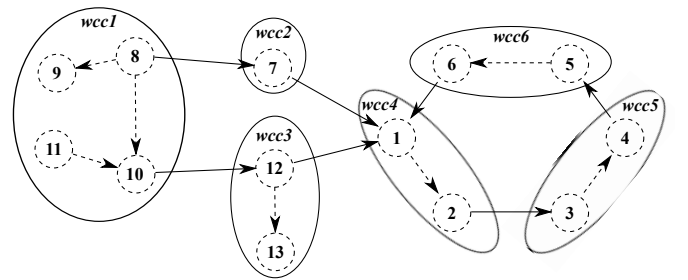
We reach the goal in two steps: (1) ensure Condition 1 with respect to only non-fragmentable links and Condition 3; (2) ensure Condition 1 with respect to fragmentable links. Details of each step are discussed below.

Treating instances as nodes and links as edges, models are just graphs. For illustration purpose, consider an example model as presented in Figure 1 where all fragmentable links are indicated by two short parallel lines crossing the links.



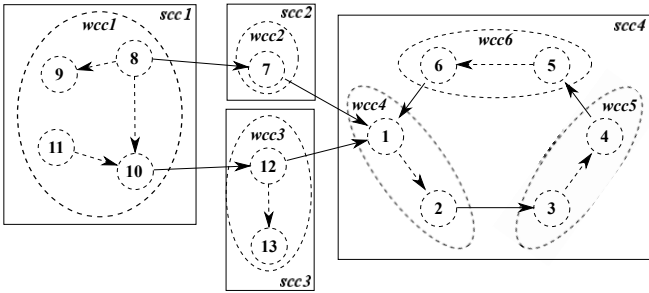
**Fig. 1** An example model

Let  $G$  be the graph derived by removing the fragmentable links from  $M$ . Because all the links in  $G$  are non-fragmentable, for a sub-graph of  $G$  to satisfy both Conditions 1 and 3, an instance is included in the sub-graph if and only if all its ancestor and descendant instances are also included. An instance together with its ancestors and descendants, from the point of view of graph theory, constitute a weakly connected component (wcc) of graph  $G$  (i.e., a connected component if we ignore edge directions). The first step of the model decomposition computes all such wcc's of  $G$ , which disjointly cover all the instances in model  $M$ , then puts back the fragmentable links. We collapse all the instances that belong to one wcc into one node, (referred to as a *wcc-node* in contrast to the original nodes), and refer to the result as graph  $W$ . After the first step, the corresponding graph  $W$  of the example model contains six wcc-nodes inter-connected by fragmentable links. We show  $W$  in Figure 2, where the instances grouped in each wcc-node remain visible (in dashed border style) for traceability purpose.



**Fig. 2** The corresponding  $W$  graph of the example model after step 1

A sub-model of  $M$  that is induced by the instances grouped in one wcc-node in  $W$  satisfies Condition 1 and Condition 3, but only with respect to non-fragmentable links for Condition 1, because wcc's are computed in the context of  $G$  where fragmentable links are removed. The second step of the model decomposition starts from graph  $W$  and tries to satisfy Condition 1 with respect to fragmentable links, i.e., following outgoing fragmentable links. More specifically, we compute all the strongly connected components (scc's) in  $W$  (see [33] for a definition of strongly connected components) and collapse all the wcc-nodes that belong to one scc into one node, (referred to as *an scc-node*), and refer to the result as graph  $D$ . After the second step, the corresponding graph  $D$  of the example model looks as depicted in Figure 3, where the original model instances and previous wcc-nodes that are grouped in each scc-node are also shown (in dashed border style) for traceability purpose. The three wcc-nodes



**Fig. 3** The corresponding  $D$  graph of the example model after step 2

$wcc4$ ,  $wcc5$  and  $wcc6$  of graph  $W$  are collapsed into one scc-node  $scc4$  because they lie on a (directed) cycle.

Note that we only collapse wcc-nodes of a strongly connected component in the second step instead of any reachable ones following outgoing fragmentable links in  $W$ , because we do not want to lose any potential sub-model of  $M$  satisfying both Conditions 1 and 3 on the way. More precisely, a set of instances is grouped in an scc-node only if for every instance induced sub-model  $M'$  of  $M$  satisfying both Conditions 1 and 3, it is either completely contained in  $M'$  or disjoint with  $M'$ , i.e., no such  $M'$  can tell the instances in the set apart.

The computational complexity of the above algorithm is dominated by the complexity of computing weakly and strongly connected components in the model graph. Computing weakly connected components amounts to computing connected components if we ignore the direction of the edges. We can compute connected components and strongly connected components in linear time using depth-first search [33]. Thus the overall complexity is linear in the size of the model graph.

### 3.3 Correctness

Graph  $D$  obtained at the end of the algorithm is a DAG (Directed Acyclic Graph) with all the edges being fragmentable links. Graph  $D$  represents a partition of the original model  $M$  where all the instances that are grouped in an scc-node in  $D$  constitute a component in the partition. We call graph  $D$  the *decomposition hierarchy* of model  $M$ .

To relate the decomposition hierarchy to the sub-models, we introduce the concept of an *antichain-node*. An antichain-node is derived by collapsing a (possibly empty) antichain of scc-nodes (i.e., a set of scc-nodes that are neither descendants nor ancestors of one another, the concept of antichain being borrowed from order theory) plus their descendants (briefly an antichain plus descendants) in the decomposition hierarchy. For example, in the decomposition hierarchy of the example model given in Figure 3, the two scc-nodes: *scc2* and *scc3* constitute an antichain, and collapsing them with their descendant *scc4*, gives rise to an antichain-node, which groups the nine instances that are previously grouped in the collapsed scc-nodes. Sub-models are then induced by the instances grouped in antichain-nodes.

To demonstrate the correctness of the algorithm, we prove both the soundness, i.e., a sub-model induced by the instances grouped in an antichain-node satisfies both Conditions 1 and 3, and the completeness, i.e., any sub-model satisfying both Conditions 1 and 3 can be induced by the instances grouped in an antichain-node. We formally capture the correctness by the following theorem:

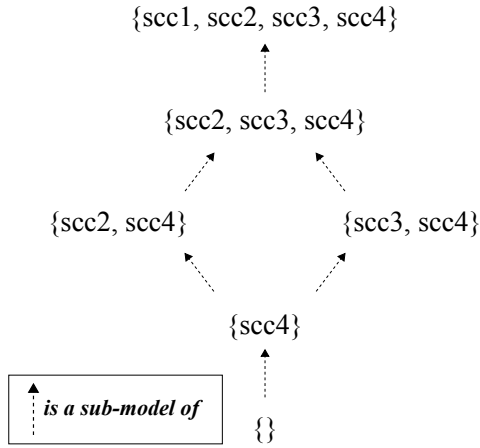
**Theorem 2** *Given a model  $M = (\mathbb{M}, \mathbb{N}, A, \tau, \text{src}, \text{tgt})$  and an instance induced sub-model  $M' = (\mathbb{M}, \mathbb{N}', A', \tau', \text{src}', \text{tgt}')$  of  $M$ ,  $M'$  satisfies both Conditions 1 and 3 if and only if there exists a corresponding antichain-node of the decomposition hierarchy of  $M$  where  $M'$  is induced by all the instances grouped in this antichain-node.*

*Proof* See Appendix A.1. □

### 3.4 The lattice of sub-models

Recall that a lattice is a partially-ordered set in which every pair of elements has a least upper bound and a greatest lower bound. Thanks to Theorem 2, we can now refer to an instance induced sub-model  $M'$  of model  $M$  that satisfies both Conditions 1 and 3 by the corresponding antichain-node  $\alpha$  of the decomposition hierarchy of  $M$ . Given a model  $M$ , all the instance induced sub-models that satisfy both Conditions 1 and 3 constitute a lattice ordered by the relation “is a sub-model of”, referred to as *the sub-model lattice* of  $M$ . Let  $\alpha_1$  and  $\alpha_2$  denote two such sub-models. The least upper bound  $(\alpha_1 \vee \alpha_2)$  and the greatest lower bound  $(\alpha_1 \wedge \alpha_2)$  are computed in the following way:

- $\alpha_1 \vee \alpha_2$  collapses all the scc-nodes that are collapsed in either  $\alpha_1$  or  $\alpha_2$ . If an scc-node is collapsed in either  $\alpha_1$  and  $\alpha_2$ , so are all its descendants because  $\alpha_1$  and  $\alpha_2$  are antichain-nodes. Therefore,  $\alpha_1 \vee \alpha_2$  is an antichain-node, i.e., an instance induced sub-model of  $M$ . Moreover, it is the least one of which both  $\alpha_1$  and  $\alpha_2$  are sub-models.
- $\alpha_1 \wedge \alpha_2$  collapses all the scc-nodes that are both collapsed in  $\alpha_1$  and  $\alpha_2$ . If an scc-node is collapsed in both  $\alpha_1$  and  $\alpha_2$ , so are all its descendants because  $\alpha_1$  and  $\alpha_2$  are antichain-nodes. Therefore,  $\alpha_1 \wedge \alpha_2$  is an antichain-node, i.e., an instance induced sub-



**Fig. 4** The sub-model lattice of the example model in Figure 1 whose decomposition hierarchy is given in Figure 2

model of  $M$ . Moreover, it is the greatest one that is a sub-model of both  $\alpha_1$  and  $\alpha_2$ .

The top of the sub-model lattice is  $M$  itself (whose corresponding antichain-node collapses all the scc-nodes), and the bottom is the empty sub-model (whose corresponding antichain-node collapses none scc-node).

For the example model discussed in Section 3.2 whose decomposition hierarchy is given in Figure 3, six possible antichain-nodes can be derived from the decomposition hierarchy, denoted by the set of scc-nodes that are collapsed. They are ordered in a lattice as shown in Figure 4.

### 3.5 Implementation

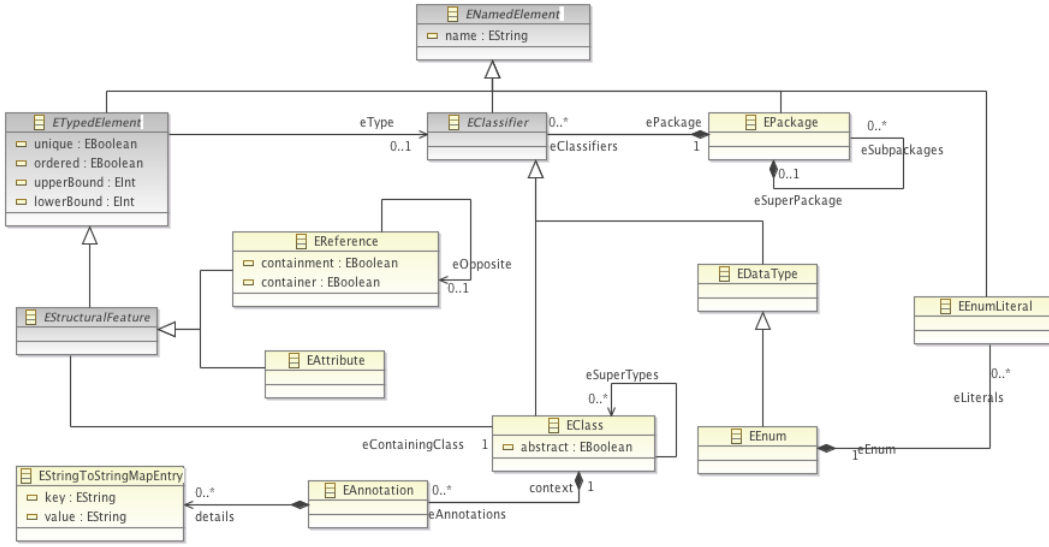
We have implemented the model decomposition technique [1] based on linear graph algorithms to compute connected components and strongly connected components [33]. The implementation takes a model of any

metamodel that follows Definition 1 as input, and computes the decomposition hierarchy of it from which the sub-model lattice can be constructed by enumerating all the antichain-nodes of the decomposition hierarchy. Note that in the worst case where the decomposition hierarchy contains no edges, the size of the sub-model lattice equals the size of the power-set of the decomposition hierarchy, which is exponential. However, this does not impact the complexity of the implementation, which is linear in the size of the model graph, because we do not require an explicit representation of the sub-model lattice. Instead, the lattice and the sub-models in it are all derived from the decomposition hierarchy.

## 4 EMF in A Nutshell

The definitions of Section 2 capture a minimal abstraction of the metamodeling and modeling concepts that are relevant for model decomposition and common to mainstream tools. As a consequence, the model decomposition technique defined for this abstraction layer can be used to decompose models conforming to metamodels defined in any of these concrete environments. We demonstrate in this paper the usage of the model decomposition technique in EMF (the Eclipse Modeling Framework) — arguably the most popular modeling framework used by the model-driven software development community.

EMF [32] is a framework for metamodel specification and provides a reflective editor for model creation. A



**Fig. 5** Class diagram of a core part of Ecore focusing on metamodel specification and oriented to model decomposition.

metamodel in EMF is instantiated from Ecore (`Ecore.ecore`) — a meta metamodel aligned with the OMG EMOF [24] standard, and is accompanied with a constraint model that captures extra well-formedness invariants. This section lays the ground for the usage of our model decomposition technique in EMF by providing a short summary of metamodels, models and conformance in EMF. It by no means attempts to give a formal capture of EMF neither of Ecore. Works on this topic deserve their own contributions, such as [6] and [3].

The focus is put on the metamodeling capability of EMF and the applicability of the model decomposition technique in EMF. More specifically, the class diagram presented in Figure 5 represents a core part of Ecore that is oriented to this focus in the following sense:

1. Implementation relevant elements and generics are omitted.

2. The notion of operations is left out of metamodels.

Side effect free helper operations that are used for constraint specifications will be defined directly in the accompanying constraint model when needed.

#### 4.1 Metamodels in EMF

A metamodel  $\mathbb{M}$  in EMF (represented by a `.ecore` file) is constructed by instantiating the class diagram given in Figure 5. Briefly speaking, an EMF metamodel consists of sets of elements instantiated from the non-abstract classes (boxes with yellow background) in the diagram, and these elements are connected by instantiating the relations (lines) in the diagram, where all the conditions posed upon the relations with adornments (such as a black diamond) and multiplicities on the corresponding line ends must be respected. We summarize the main sets of elements of an EMF metamodel by a tuple  $\mathbb{M}$

$= (\mathbb{C}, \mathbb{DT}, \mathbb{DV}, \mathbb{Rf}, \mathbb{At})$ , and explain the sets and their relations below.

1. A finite set of classes  $\mathbb{C}$  ranged over by  $c \in \mathbb{C}$  (instances of `EClass`).
  - A class  $c$  may be abstract, i.e.,  $c.\text{abstract}=\text{true}$ .
  - We can connect a class to another class within the relation `eSuperTypes` to say that the former inherits from the latter. The transitive closure of `eSuperTypes` should be irreflexive. We define the subtyping relation ( $\preceq$ ) between classes as the reflexive and transitive closure of `eSuperTypes`.
  - We can associate an annotation (instance of `EAnnotation`) to a class with a specific key equal to `''constraints''` so as to enumerate the names of the invariants that are imposed on the class in the value of the annotation. Enumerated invariants are to be defined in the accompanying constraint model of the metamodel, with exactly the same context classes and invariant names.
2. A finite set of data types  $\mathbb{DT}$  ranged over by  $dt \in \mathbb{DT}$  (instances of `EDataType` and `EEnum`) to capture both primitive types, i.e., `Integer`, `Boolean`, `String`, and `Real`, predefined for any metamodel, and enumeration types, specific to the metamodel.
  - The subtyping relation ( $\preceq$ ) between primitive types follows the conventional definition, i.e., `Integer` is a subtype of `Real`. No subtyping relation holds between enumeration types.

3. A finite set of data values  $\mathbb{DV}$  ranged over by  $dv \in \mathbb{DV}$ .
  - Values of metamodel specific enumeration types are explicitly specified by enumeration literals (instances of `EEnumLiteral`) to extensionally define the types.
  - Values of predefined primitive types are implicitly given following the conventional definitions.
  - A typing function from data values to data types:  $\tau_d : \mathbb{DV} \rightarrow \mathbb{DT}$  tells the type of a data value. Note that if the data value  $dv$  is an enumeration literal,  $\tau_d(dv)$  should coincide with  $dv.\text{eEnum}$ , i.e.,  $dv$  is connected to  $\tau_d(dv)$  within the relation `eEnum`.
4. A finite set of references  $\mathbb{Rf}$  ranged over by  $rf \in \mathbb{Rf}$  (instances of `EReference`). A reference  $rf$  has an owning class (aka. the source of the reference), specified by  $rf.\text{eContainingClass}$ .
5. A finite set of attributes  $\mathbb{At}$  ranged over by  $at \in \mathbb{At}$  (instances of `EAttribute`). An attribute  $at$  has an owning class (aka. the source of the attribute), specified by  $at.\text{eContainingClass}$ .

References and attributes are both known as the structural features of a class. A structural feature, written  $sf$ , has itself the following features.

- A type (aka. the target of the structural feature), specified by  $sf.\text{eType}$ . Note that the type of a reference must be a class and the type of an attribute must be a data type.

- A multiplicity, specified by `sf.lowerBound` for its lower bound, which is a natural number (i.e., zero or positive integer), and `sf.upperBound` for its upper bound, which is either a positive integer or  $-1$ . Note that EMF implements an unlimited upper bound, i.e.,  $*$ , by  $-1$ . In case the upper bound of the multiplicity is not  $-1$ , it must be greater than or equal to the lower bound.
- When `sf.ordered = true`, we call the structural feature `sf` ordered.
- When `sf.unique = true`, we call the structural feature `sf` unique.

In addition to the common features above, references have the following extra features.

- A reference `rf` is a containment, when `rf.containment = true` (depicted by a line with a filled diamond at the source end in class diagrams).
- A reference `rf` can have an opposite reference, specified by `rf.eOpposite`, where the relation `eOpposite` should be functional, symmetric, and irreflexive.

The accompanying constraint model of a metamodel in EMF consists of a finite set of invariants  $\mathbb{E}Inv$ , ranged over by  $eInv$ , to provide specifications to the invariants enumerated in the metamodel. The context of an invariant  $eInv$  is the same as the class to which the corresponding annotation (i.e., the annotation in the value of which the name of the invariant is enumerated) is associated. The body of a invariant is given by an expression. We

keep the notion of expression abstract to accommodate all potential languages for expression specification. In Section 5, we define a formal constraint language called CoreOCL for this purpose. In addition, a finite set of helper operations can be defined in the constraint model.

#### 4.2 Models in EMF

A model  $\mathbf{EM}$  in EMF is constructed by instantiating the classes of an EMF metamodel  $\mathbb{EM}$  and assigning values to attributes and references owned by the classes. We summarize the components of a model in EMF by a tuple  $\mathbf{EM} = (\mathbb{EM}, \mathbf{O}, \tau_{\mathbf{O}}, \mathbf{RA}, \tau_{\mathbf{ra}}, \mathbf{ra}_s, \mathbf{ra}_t, \mathbf{AA}, \tau_{\mathbf{aa}}, \mathbf{aa}_s, \mathbf{aa}_t)$  and explain below.

1. A finite set of objects  $\mathbf{O}$  ranged over by  $\mathbf{o}$ .
2. An object typing function  $\tau_{\mathbf{O}} : \mathbf{O} \rightarrow \mathbb{C}$  from objects to classes, to tell the type of an object, i.e., from which class the object is instantiated.
3. A finite set of reference assignments  $\mathbf{RA}$  ranged over by  $\mathbf{ra}$ .
4. Reference assignments involve three functions:  $\tau_{\mathbf{ra}} : \mathbf{RA} \rightarrow \mathbb{Rf}$ ,  $\mathbf{ra}_s : \mathbf{RA} \rightarrow \mathbf{O}$ , and  $\mathbf{ra}_t : \mathbf{RA} \rightarrow \mathbf{O}$  that specify respectively the type, the source, and the target of a reference assignment. The semantics of a reference assignment is  $\mathbf{ra}_s(\mathbf{ra}).\tau_{\mathbf{ra}}(\mathbf{ra}) := \mathbf{ra}_t(\mathbf{ra})$ .
5. The set of reference assignments  $\mathbf{RA}$  is partially ordered, where for each ordered reference `rf` of the metamodel  $\mathbb{EM}$  and an object  $\mathbf{o}$  of the model  $\mathbf{EM}$ ,  $\{\mathbf{ra} \mid \mathbf{ra} \in \mathbf{RA} \text{ and } \tau_{\mathbf{ra}}(\mathbf{ra}) = \mathbf{rf} \text{ and } \mathbf{ra}_s(\mathbf{ra}) = \mathbf{o}\}$  con-

stitutes a totally ordered subset of  $RA$ , (i.e., a chain in the partial order on  $RA$ ).

6. A finite set of attribute assignments  $AA$  ranged over by  $aa$ .
7. Attribute assignments involve three functions:  $\tau_{aa} : AA \rightarrow \mathbb{A}t$ ,  $aa_s : AA \rightarrow O$ , and  $aa_t : AA \rightarrow DV$  that specify respectively the type, the source, and the target of an attribute assignment. The semantics of an attribute assignment is  $aa_s(aa).\tau_{aa}(aa) := aa_t(aa)$ .
8. The set of attribute assignments  $AA$  is partially ordered, where for each ordered attribute  $at$  of the metamodel  $\mathbb{EM}$  and an object  $o$  of the model  $EM$ ,  $\{aa \mid aa \in AA \text{ and } \tau_{aa}(aa) = at \text{ and } aa_s(aa) = o\}$  constitutes a totally ordered subset of  $AA$ , (i.e., a chain in the partial order on  $AA$ ).

#### 4.3 Conformance in EMF

In EMF, the conformance of a model  $EM$  to a metamodel  $\mathbb{EM}$  requires the following conditions to hold:

1. Typing condition: all reference assignments are valid.  
For a reference assignment  $ra$  to be valid,
  - (a)  $\tau_o(ra_s(ra))$  is a subtype of  $\tau_{ra}(ra).eContainingClass$ , i.e.,  $\tau_{ra}(ra)$  should be a reference defined for the class of  $ra_s(ra)$ ;
  - (b) the class of  $ra_t(ra)$  is compatible with the type of  $\tau_{ra}(ra)$ , i.e.,  $\tau_o(ra_t(ra))$  is a subtype of  $\tau_{ra}(ra).eType$ .
2. Attribute typing condition: all attribute assignments are valid. For an attribute assignment  $aa$  to be valid,

- (a)  $\tau_o(aa_s(aa))$  is a subtype of  $\tau_{aa}(aa).eContainingClass$ , i.e.,  $\tau_{aa}(aa)$  should be an attribute defined for the class of  $aa_s(aa)$ ;
- (b) the type of  $aa_t(aa)$  is compatible with the type of  $\tau_{aa}(aa)$ , i.e.,  $\tau_d(aa_t(aa))$  is a subtype of  $\tau_{aa}(aa).eType$ .

3. Reference multiplicity condition: the total number of assignments to a reference  $rf$  in an object should fall in the range specified by the multiplicity of  $rf$ . Namely,  $\forall o \in O, rf \in \mathbb{R}f$ , let  $k = \#\{ra \in RA \mid \tau_{ra}(ra) = rf \text{ and } ra_s(ra) = o\}$  (i.e., the number of assignments to the reference  $rf$  in object  $o$ ), we must have  $k \geq rf.lowerBound$  and  $k \leq rf.upperBound$  in case  $rf.upperBound \neq -1$ .
4. Attribute multiplicity condition: the total number of assignments to an attribute  $at$  in an object should fall in the range specified by the multiplicity of  $at$ . Namely,  $\forall o \in O, at \in \mathbb{A}t$ , let  $k = \#\{aa \in AA \mid \tau_{aa}(aa) = at \text{ and } aa_s(aa) = o\}$  (i.e., the number of assignments to the attribute  $at$  in object  $o$ ), we must have  $k \geq at.lowerBound$  and  $k \leq at.upperBound$  in case  $at.upperBound \neq -1$ .
5. Opposite reference condition: opposite references are assigned in a pair-wise way. Namely, given two references  $rf_1$  and  $rf_2$  where  $rf_1.eOpposite = rf_2$ , for any reference assignment  $ra_1 \in RA$  where  $\tau_{ra}(ra_1) = rf_1$ , there exists a reference assignment  $ra_2 \in RA$  such that  $\tau_{ra}(ra_2) = rf_2$ ,  $ra_s(ra_1) = ra_t(ra_2)$ , and  $ra_t(ra_1) = ra_s(ra_2)$ .

6. Containment condition: an object cannot be contained in more than one object, neither can it be contained in itself.
7. Abstract class condition: no objects can be instantiated from an abstract class.
8. Uniqueness condition: if a reference (resp. an attribute) is flagged as unique, the assigned values to the reference (resp. the attribute) in an object must be distinct from one to another.
9. Invariant condition: all the invariants should hold.

## 5 CoreOCL: A Formal Capture of Core EssentialOCL

“EssentialOCL is the package exposing the minimal OCL required to work with EMOF.” (OMG OCL specification [27], p. 187). As a consequence, concepts such as messages, association classes, and states are not part of EssentialOCL.

In this section, we formalize a core of EssentialOCL, called CoreOCL. CoreOCL provides a theoretical framework for defining invariants and discussing their evaluation semantics. CoreOCL is the formal abstraction and mathematical foundation of EssentialOCL. CoreOCL consists of the core constructs of EssentialOCL and exposes the same expressive power as the latter. The relationship between EssentialOCL and CoreOCL can be characterized by the following keywords: abstraction, normalization, simplification, and formalization.

*Abstraction* Pre-defined operations in the OCL standard library that are relevant for EssentialOCL are abstracted into kappa operations, i.e., we use  $\kappa$  to range over their names in the syntax. Kappa operations include pre-defined operations for primitive types such as arithmetic operations and logical operations; pre-defined operations on collection types such as collection size and collection union; and pre-defined operations for all values such as equality-check operations and operations whose names come with an ocl prefix (e.g., `oclIsUndefined()`). Note that kappa operations take only values as parameters. In case a pre-defined operation in the OCL standard library takes also types as parameters, such as the `oclAsType` and `allInstances` operations, they are captured in CoreOCL explicitly.

*Normalization* A uniform prefix notation is used for all kappa operation calls. For example,  $x \wedge y$  is expressed as `AND( $x, y$ )`,  $x > y$  is expressed as `>( $x, y$ )`, and `Set{ $x$ } → including( $x$ )` is expressed as `including(Set{ $x$ },  $x$ )`, where `AND`, `>`, and `including` are all instances of  $\kappa$ .

Moreover, for each language concept, CoreOCL supports only one canonical syntactic representation. Inessential “syntactic sugar”, i.e., syntax constructs that can be removed without any effect on the expressive power, are dismissed in CoreOCL. Applications of syntactic sugar can be systematically replaced with equivalent constructs from the core subset. Here are some examples of dis-

missed syntactic sugar and the corresponding encoding in CoreOCL:

- Tuple: the type information in a tuple is required in legal CoreOCL syntax, and all tuples without types specified for their parts, which is legal syntax in EssentialOCL, should be augmented with the inferred type information. For example, the following tuple: `Tuple(name="Smith Johns", age=30)` will be expressed as `Tuple(name:String="Smith Johns", age: Integer=30)`.
- Collection iterator: only the `iterate` operation is legal syntax in CoreOCL (see CoreOCL syntax below), and all other collection iterators need to be described in terms of `iterate`. For example, the following iterator: `forAll(x | > (x, 0))`, checks if a set has only positive integers. This will be expressed as `iterate(x; y = true | if > (x, 1) then AND(y, true) else AND(y, false))`.
- Non-empty collection literals: only empty collection literals are legal syntax in CoreOCL (see CoreOCL syntax below), and all non-empty collection literals such as `Set{1, 2}` will be expressed as nested calls to the pre-defined “including” operation on collections, i.e., `including(including(Set{}, 1), 2)`.

*Simplification* CoreOCL omits the following in the syntax: `OclAny`, `OclVoid`, `OclInvalid`, `null`, and `invalid`, because they are not supposed to be used directly in expressions. The first three types are used in the type system of CoreOCL, and the last two values are used in the evaluation semantics of CoreOCL.

*Formalization* CoreOCL lays a formal foundation upon which we can precisely define the evaluation semantics of invariants written in CoreOCL and the corresponding scope that has been introduced abstractly in Definition 6. Moreover, we are also able to demonstrate (in Theorem 5) formally that letting a sub-model include all the instances of the evaluation scope of an invariant as they are in the original model suffices to maintain the same evaluation, whereas in Section 3, this statement can only be discussed informally given the abstract setting. Finally, we explore a sufficient condition for identifying forward invariants as defined in Definition 8 in a practical way, by proving that CoreOCL invariants that do not call `allInstances` are all forward.

We present CoreOCL in the context of EMF. Namely, the concrete syntax of CoreOCL adopts EMF notations for the specification of invariants for EMF metamodels and the semantics of CoreOCL is defined in terms of how these invariants are evaluated on EMF models.

### 5.1 CoreOCL expressions: syntax

The detailed definition of CoreOCL expressions is given below in EBNF format.

$\mathbb{t} ::=$	<b>Types</b>
$\mathsf{dt}$	Data types
$\mathsf{c}$	Classes
$\mathsf{Collection}(\mathbb{t})$	Collection types
$\mathsf{Bag}(\mathbb{t})$	Bag types
$\mathsf{Set}(\mathbb{t})$	Set types
$\mathsf{Seq}(\mathbb{t})$	Sequence types
$\mathsf{OSet}(\mathbb{t})$	Ordered set types
$\mathsf{Tuple}(\overline{\mathbb{tn}} : \bar{\mathbb{t}})$	Tuple types
$\mathsf{e} ::=$	<b>Expressions</b>
$\mathsf{dv}$	Data value
$x \mid \mathsf{self}$	Variable
$\mathsf{e.sf}$	Structural feature call
$\mathsf{Tuple}\{\overline{\mathbb{tn}} : \bar{\mathbb{t}} = \bar{\mathsf{e}}\}$	Tuple
$\mathsf{e.tn}$	Tuple part call
$\mathsf{let } x = \mathsf{e}_1 \text{ in } \mathsf{e}_2$	Let binding
$\mathsf{def } \mathsf{c} :: \mathsf{op} = \lambda(\overline{x} : \bar{\mathbb{t}}).\mathsf{e}_1 \text{ in } \mathsf{e}_2$	Def expression
$\mathsf{e.op}(\bar{\mathsf{e}})$	Object operation call
$\kappa(\bar{\mathsf{e}})$	Kappa operation call
$\mathsf{if } \mathsf{e}_1 \text{ then } \mathsf{e}_2 \text{ else } \mathsf{e}_3$	Conditional expression
$\mathsf{e.asInstanceOf } \mathsf{c}$	Downcast
$\mathsf{e.isInstanceOf } \mathsf{c}$	Is Instance Of
$\mathsf{e.isOfKind } \mathsf{c}$	Is Kind Of
$\mathsf{Bag}\{\} \mid \mathsf{Set}\{\} \mid \mathsf{Seq}\{\} \mid \mathsf{OSet}\{\}$	Empty collections
$\mathsf{e}_1 \rightarrow \mathsf{iterate}(x; y = \mathsf{e}_2 \mid \mathsf{e}_3)$	Collection iteration
$\mathsf{allInstances}(\mathsf{c})$	All instances

The following notational convention is adopted for sequences: we write  $\overline{X}$  as the shorthand for a sequence

of the form  $X_1, \dots, X_k$ , when the members of the sequence are not necessarily required to be accessible explicitly. We also abbreviate handling of paired sequences in an obvious way. For example, we write  $(\overline{x} : \bar{\mathbb{t}})$  as the shorthand for  $(x_1 : \mathbb{t}_1, \dots, x_k : \mathbb{t}_k)$ ,  $(\overline{\mathbb{tn}} : \bar{\mathbb{t}} = \bar{\mathsf{e}})$  as the shorthand for  $(\mathbb{tn}_1 : \mathbb{t}_1 = \mathsf{e}_1, \dots, \mathbb{tn}_k : \mathbb{t}_k = \mathsf{e}_k)$ , and  $(\overline{x} := \bar{\mathsf{v}})$  as the short hand for  $(x_1 := \mathsf{v}_1, \dots, x_k := \mathsf{v}_k)$ . The first is used in the syntax to specify the parameters and their types of an object operation. The second is used in the syntax to specify the parts of a tuple. And the third is used in the semantic rule for evaluating object operation calls to assign arguments to parameters (See rule OBJECTOPERATIONCALL in Figure 6).

In addition to data types and classes, OCL introduces the notion of collections and tuples. Collections are used to handle navigations of multi-valued structural features (i.e., references, attributes, etc.). A structural feature is *single-valued* when its upper bound is 1, otherwise *multi-valued*. Five collection types are provided, where four of them, i.e., the **Bag**, **Set**, **Seq**, and **OSet** types are concrete and the fifth, i.e., *Collection* is the abstract supertype of the other four. Moreover, both **Set** and **Seq** are subtypes of **Bag**, and **OSet** is a subtype of both **Set** and **Seq**.

Tuple types (resp. values) provide a way to compose several types (resp. values) together to form compound types (resp. values). The parts of a tuple can be accessed by their names (ranged over by  $\mathbb{tn}$ ) using the same dot notation as is used for accessing structural features.

Variables are ranged over by identifiers such as  $x, y$ , etc. The keyword **self** is a reserved identifier that represents a special variable referring to the current contextual object. We use  $\mathfrak{s}$  to range over both attributes and references, which are commonly referred to as structural features in EMF (and properties in EMOF). We can call a structural feature on an object to retrieve the assigned value(s).

Let bindings and def expressions define variables and object operations (ranged over by  $\mathfrak{op}$ ) reusable in the nested expressions. Note that according to [27], all the object operations used in OCL are side effect free. The parameters of an operation are treated as variables local to the corresponding body expression of the operation. We call an operation on an object by passing concrete arguments to parameters.

Finally, iterator operations are a special type of operations on collection types that enable one to iterate over the elements in a collection. We only capture the most fundamental and generic collection iterator called **iterate** in CoreOCL to keep the language small, as all the other collection iterators can be derived from it [34].

## 5.2 CoreOCL expressions: evaluation

### 5.2.1 Evaluation judgments

The evaluation of a CoreOCL expression  $e$  takes place in a model  $\mathbf{EM}$  and an evaluation environment  $\Lambda$  and results in a pair  $(v, O_S)$ , where  $v$  is the value of the ex-

pression and  $O_S$  consists of the objects of  $\mathbf{EM}$  that are referenced by  $e$  during the evaluation. The model  $\mathbf{EM}$  is exploited during the course of the evaluation for relevant model information such as the types of instances and the values assigned to attributes and references.

Thereafter, we make use of the following judgment to denote the evaluation semantics of CoreOCL expressions:

$$\Lambda \stackrel{\mathbf{EM}}{\models} e \Downarrow (v, O_S)$$

which we read as: within the context  $\Lambda$  and the model  $\mathbf{EM}$ , the value of  $e$  is  $(v, O_S)$ .

### 5.2.2 Evaluation environments

An evaluation environment  $\Lambda$  is a function that relates qualified object operation names ( $c :: \mathfrak{op}$ ) to operation definitions ( $\lambda(\bar{x}).e$ ) and variables to values. In CoreOCL, we consider the following types of values:

$v ::=$	<b>Values</b>
$dv$	Data value
$o$	Object value
$\text{Tuple}\{\bar{t}_n : \bar{t} = \bar{v}\}$	Tuple value
$\text{null}$	Unknown value
$\text{Bag}\{\bar{v}\} \mid \text{Set}\{\bar{v}\} \mid \text{Seq}\{\bar{v}\} \mid \text{OSet}\{\bar{v}\}$	Collection value

Evaluation environments can be specified extensionally by listing the pairs related in it. The overriding of an evaluation environment  $\Lambda_1$  by another evaluation environment  $\Lambda_2$  (also called overriding union) is an extension of  $\Lambda_2$  denoted as  $\Lambda_1 \oplus \Lambda_2$ . The domain of  $\Lambda_1 \oplus \Lambda_2$  is the union of the domain of  $\Lambda_1$  and the domain of  $\Lambda_2$ . In  $\Lambda_1 \oplus \Lambda_2$ , any element of the domain of  $\Lambda_2$  is related to

its image under  $A_2$ , and any other element exclusive to the domain of  $A_1$  is related to its image under  $A_1$ .

### 5.2.3 Expression evaluation rules

We say an evaluation judgment holds if a derivation tree of the evaluation judgment can be constructed by instantiating the evaluation rules presented in Figure 6. Note that in the presentation of the rules, wherever the model information is required for a computation, **EM** appears explicitly on top of the corresponding computation sign. Given an evaluation environment  $A$ , an expression  $e$  and a model **EM**, if no evaluation judgment holds for them, we say  $e$  is *invalid* within context  $A$  and model **EM**.

*Some words about typing* Without making the type system explicit, the evaluation rules assign only semantics to well-typed expressions. Notably, we draw the attention of the reader to the following points.

1. Due to subtype polymorphism, an expression of type  $\mathbb{t}$  has also all the super types of  $\mathbb{t}$ .
2. Variables are first defined via let bindings then referenced in nested variable expressions. Note that being a pre-defined variable for the reference of the current contextual object, no explicit let binding is needed for **self**.
3. Operations defined via def expressions carry distinct parameters.
4. Structural feature calls and object operation calls are only addressed on objects and the callee (i.e., a reference or an attribute or an object operation) is actually defined for the type of the object. Subsequently, calling structural features or object operations on a collection of objects should be explicitly implemented by using the **iterate** operation that iterates the structural feature or object operation call on all the element objects of the collection.
5. Tuple part calls are only addressed on tuple values and the callee (i.e., the called tuple part name  $\mathbb{tn}$ ) is actually a defined part of the tuple. Similarly, calling parts on a collection of tuples should be explicitly implemented by using the **iterate** operation.
6. Given a tuple  $\text{Tuple}\{\overline{\mathbb{tn}} : \overline{\mathbb{t}} = \overline{\mathbb{e}}\}$ , the type of the expression  $e_i$  that is used to indicate the value of the part  $\mathbb{tn}_i$  in the tuple is compatible with the type that is specified for the part in the tuple, i.e.,  $\mathbb{t}_i$ .
7. Both object operations and kappa operations are called with type and cardinality compatible arguments.
8. Given a conditional expression **if**  $e_1$  **then**  $e_2$  **else**  $e_3$ ,  $e_1$  has type **Boolean**, and  $e_2$  and  $e_3$  have the same type.
9. Downcast expressions (**AsInstanceOf**), type testing expressions (**isInstanceOf**), and subtype testing expressions (**isKindOf**) are only applied to objects.
10. Finally, given a collection iteration expression  $e_1 \rightarrow \text{iterate}(x; y = e_2 \mid e_3)$ ,  $e_1$  has a collection type. More-

<p><b>DATAVALUE</b></p> $\frac{}{\Lambda \stackrel{\text{EM}}{\Vdash} \text{dv} \Downarrow (\text{dv}, \emptyset)}$	<p><b>VARIABLE</b></p> $\frac{\Lambda(x) = \mathbf{v} \quad \mathbf{O}_S(\mathbf{v}) \subseteq \mathbf{O}}{\Lambda \stackrel{\text{EM}}{\Vdash} x \Downarrow (\mathbf{v}, \mathbf{O}_S(\mathbf{v}))}$	<p><b>STRUCTURALFEATURECALL</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S^1) \quad \mathbf{o}.\text{sf} \stackrel{\text{EM}}{=} \mathbf{v}}{\Lambda \stackrel{\text{EM}}{\Vdash} e.\text{sf} \Downarrow (\mathbf{v}, \mathbf{O}_S^1 \cup \mathbf{O}_S(\mathbf{v}))}$
<p><b>TUPLE</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} \bar{e} \Downarrow (\bar{\mathbf{v}}, \overline{\mathbf{O}_S})}{\Lambda \stackrel{\text{EM}}{\Vdash} \text{Tuple}\{\bar{\text{tn}} : \bar{\ell} = \bar{e}\} \Downarrow (\text{Tuple}\{\bar{\text{tn}} : \bar{\ell} = \bar{\mathbf{v}}\}, \bigcup \overline{\mathbf{O}_S})}$	<p><b>TUPLEPARTCALL</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\text{Tuple}\{\bar{\text{tn}} : \bar{\ell} = \bar{\mathbf{v}}\}, \mathbf{O}_S) \quad \text{tn}_i = \text{tn}}{\Lambda \stackrel{\text{EM}}{\Vdash} e.\text{tn} \Downarrow (\mathbf{v}_i, \mathbf{O}_S)}$	
<p><b>LETBINDING</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e_1 \Downarrow (\mathbf{v}_1, \mathbf{O}_S^1) \quad \Lambda \oplus \{(x : \mathbf{v}_1)\} \stackrel{\text{EM}}{\Vdash} e_2 \Downarrow (\mathbf{v}_2, \mathbf{O}_S^2)}{\Lambda \stackrel{\text{EM}}{\Vdash} \text{let } x = e_1 \text{ in } e_2 \Downarrow (\mathbf{v}_2, \mathbf{O}_S^1 \cup \mathbf{O}_S^2)}$	<p><b>DEFEXP</b></p> $\frac{\Lambda \oplus \{(\text{op} : \lambda(\bar{x}).e_1)\} \stackrel{\text{EM}}{\Vdash} e_2 \Downarrow (\mathbf{v}, \mathbf{O}_S)}{\Lambda \stackrel{\text{EM}}{\Vdash} \text{def } c :: \text{op} = \lambda(\bar{x} : \bar{\ell}).e_1 \text{ in } e_2 \Downarrow (\mathbf{v}, \mathbf{O}_S)}$	
<p><b>OBJECTOPERATIONCALL</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S^{\mathbf{o}}) \quad \Lambda \stackrel{\text{EM}}{\Vdash} \bar{e} \Downarrow (\bar{\mathbf{v}}, \overline{\mathbf{O}_S}) \quad \Lambda(\text{op}) = \lambda(\bar{x}).e_{\text{op}} \quad \Lambda \oplus \{(\text{self} : \mathbf{o})\} \oplus \{(\bar{x} : \bar{\mathbf{v}})\} \stackrel{\text{EM}}{\Vdash} e_{\text{op}} \Downarrow (\mathbf{v}, \mathbf{O}_S^{\text{op}})}{\Lambda \stackrel{\text{EM}}{\Vdash} e.\text{op}(\bar{x} := \bar{e}) \Downarrow (\mathbf{v}, \mathbf{O}_S^{\mathbf{o}} \cup (\bigcup \overline{\mathbf{O}_S}) \cup \mathbf{O}_S^{\text{op}})}$		
<p><b>KAPPACALL</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} \bar{e} \Downarrow (\bar{\mathbf{v}}, \overline{\mathbf{O}_S}) \quad \kappa(\bar{\mathbf{v}}) = \mathbf{v}}{\Lambda \stackrel{\text{EM}}{\Vdash} \kappa(\bar{e}) \Downarrow (\mathbf{v}, \bigcup \overline{\mathbf{O}_S})}$	<p><b>IFTRUE</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e_1 \Downarrow (\text{true}, \mathbf{O}_S^1) \quad \Lambda \stackrel{\text{EM}}{\Vdash} e_2 \Downarrow (\mathbf{v}, \mathbf{O}_S^2)}{\Lambda \stackrel{\text{EM}}{\Vdash} \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \Downarrow (\mathbf{v}, \mathbf{O}_S^1 \cup \mathbf{O}_S^2)}$	
<p><b>IFFALSE</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e_1 \Downarrow (\text{false}, \mathbf{O}_S^1) \quad \Lambda \stackrel{\text{EM}}{\Vdash} e_3 \Downarrow (\mathbf{v}, \mathbf{O}_S^2)}{\Lambda \stackrel{\text{EM}}{\Vdash} \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \Downarrow (\mathbf{v}, \mathbf{O}_S^1 \cup \mathbf{O}_S^2)}$	<p><b>DOWNCASTOK</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S) \quad \tau_{\mathbf{o}}(\mathbf{o}) \preceq \mathbf{c}}{\Lambda \stackrel{\text{EM}}{\Vdash} e \text{ asInstanceOf } \mathbf{c} \Downarrow (\mathbf{o}, \mathbf{O}_S)}$	<p><b>DOWNCASTNOTOK</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S) \quad \tau_{\mathbf{o}}(\mathbf{o}) \not\preceq \mathbf{c}}{\Lambda \stackrel{\text{EM}}{\Vdash} e \text{ asInstanceOf } \mathbf{c} \Downarrow (\text{null}, \mathbf{O}_S)}$
<p><b>ISINSTANCEOFTRUE</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S) \quad \tau_{\mathbf{o}}(\mathbf{o}) \stackrel{\text{EM}}{=} \mathbf{c}}{\Lambda \stackrel{\text{EM}}{\Vdash} e \text{ isInstanceOf } \mathbf{c} \Downarrow (\text{true}, \mathbf{O}_S)}$	<p><b>ISINSTANCEOFFALSE</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S) \quad \tau_{\mathbf{o}}(\mathbf{o}) \stackrel{\text{EM}}{\neq} \mathbf{c}}{\Lambda \stackrel{\text{EM}}{\Vdash} e \text{ isInstanceOf } \mathbf{c} \Downarrow (\text{false}, \mathbf{O}_S)}$	<p><b>ISKINDOFTRUE</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S) \quad \tau_{\mathbf{o}}(\mathbf{o}) \stackrel{\text{EM}}{=} \mathbf{c}' \quad \mathbf{c}' \preceq \mathbf{c}}{\Lambda \stackrel{\text{EM}}{\Vdash} e \text{ isOfKind } \mathbf{c} \Downarrow (\text{true}, \mathbf{O}_S)}$
<p><b>ISKINDOFFALSE</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e \Downarrow (\mathbf{o}, \mathbf{O}_S) \quad \tau_{\mathbf{o}}(\mathbf{o}) \stackrel{\text{EM}}{=} \mathbf{c}' \quad \mathbf{c}' \not\preceq \mathbf{c}}{\Lambda \stackrel{\text{EM}}{\Vdash} e \text{ isOfKind } \mathbf{c} \Downarrow (\text{false}, \mathbf{O}_S)}$	<p><b>EMPTYBAG</b></p> $\Lambda \stackrel{\text{EM}}{\Vdash} \text{Bag}\{\} \Downarrow (\text{Bag}\{\}, \emptyset)$	<p><b>EMPTYSET</b></p> $\Lambda \stackrel{\text{EM}}{\Vdash} \text{Set}\{\} \Downarrow (\text{Set}\{\}, \emptyset)$
<p><b>EMPTYSEQ</b></p> $\Lambda \stackrel{\text{EM}}{\Vdash} \text{Seq}\{\} \Downarrow (\text{Seq}\{\}, \emptyset)$	<p><b>EMPTYOSET</b></p> $\Lambda \stackrel{\text{EM}}{\Vdash} \text{OSet}\{\} \Downarrow (\text{OSet}\{\}, \emptyset)$	<p><b>COLLECTIONITERATION</b></p> $\frac{\Lambda \stackrel{\text{EM}}{\Vdash} e_1 \Downarrow (\text{Collection}\{\mathbf{v}'_1, \dots, \mathbf{v}'_k\}, \mathbf{O}_S^0) \quad \Lambda \stackrel{\text{EM}}{\Vdash} e_2 \Downarrow (\mathbf{v}_1, \mathbf{O}_S^1) \quad (\Lambda \oplus \{(x : \mathbf{v}'_i)\} \oplus \{(y : \mathbf{v}_i)\} \stackrel{\text{EM}}{\Vdash} e_3 \Downarrow (\mathbf{v}_{i+1}, \mathbf{O}_S^{i+1}))^{i=1, \dots, k}}{\Lambda \stackrel{\text{EM}}{\Vdash} e_1 \rightarrow \text{iterate}(x; y = e_2 \mid e_3) \Downarrow (\mathbf{v}_{k+1}, \bigcup_{0 \leq i \leq k+1} \mathbf{O}_S^i)}$
<p><b>ALLINSTANCES</b></p> $\Lambda \stackrel{\text{EM}}{\Vdash} \text{allInstances}(\mathbf{c}) \Downarrow (\text{Set}\{\mathbf{o} \mid \mathbf{o} \in \mathbf{O} \text{ and } \tau_{\mathbf{o}}(\mathbf{o}) \stackrel{\text{EM}}{=} \mathbf{c}\}, \{\mathbf{o} \mid \mathbf{o} \in \mathbf{O} \text{ and } \tau_{\mathbf{o}}(\mathbf{o}) \stackrel{\text{EM}}{\neq} \mathbf{c}\})$		

Fig. 6 Evaluation rules for CoreOCL expressions

over, the expression  $e_2$  that assigns initial value to the accumulator  $y$  and the body expression  $e_3$  have the same type, in order for the iteration to take place properly.

*Digesting the rules* Rules **STRUCTURALFEATURECALL** and **VARIABLE** in Figure 6 both rely on an auxiliary definition  $O_S(v)$  that computes the set of objects directly referenced by a value  $v$ . The computation is defined recursively for different kinds of values as follows:

$$\begin{aligned}
O_S(dv) &= \emptyset \\
O_S(o) &= \{o\} \\
O_S(\text{Tuple}\{\overline{tn} : \overline{t} = \overline{v}\}) &= \bigcup_{v_i \in \overline{v}} O_S(v_i) \\
O_S(\text{null}) &= \emptyset \\
O_S(\text{Bag}\{\overline{v}\}) &= \bigcup_{v_i \in \overline{v}} O_S(v_i) \\
O_S(\text{Set}\{\overline{v}\}) &= \bigcup_{v_i \in \overline{v}} O_S(v_i) \\
O_S(\text{Seq}\{\overline{v}\}) &= \bigcup_{v_i \in \overline{v}} O_S(v_i) \\
O_S(\text{OSet}\{\overline{v}\}) &= \bigcup_{v_i \in \overline{v}} O_S(v_i)
\end{aligned}$$

The second premise of rule **VARIABLE** requires that the value  $v$  associated to the variable  $x$  in the evaluation environment should only reference objects of **EM** i.e.,  $O$ , in order for the value to be meaningful.

The computation of a structural feature on an object in a model is present as the second premises of rule **STRUCTURALFEATURECALL**. A structural feature is either a reference or an attribute. We write  $n.sf \stackrel{EM}{=} v$  to denote the computation of a well-typed structural feature call in a model, whose result is either a collection, a data value, an object, or **null**, depending on the target multiplicity of the structural feature. More specifically,

if the structural feature is single-valued, i.e., its upper bound is 1, the result is the data value assigned to the attribute or the object assigned to the reference, or **null** if no value is assigned to the structural feature at all. Otherwise, the result is a collection of the values assigned to the structural feature in the model. Then depending on the two flags: ordered and unique, the collection is either a set (not ordered but unique), a bag (not ordered and not unique), a sequence (ordered but not unique), or a ordered set (ordered and unique).

We summarize the computation of a structural feature call on an object in a model in Table 1, in which, we use **sfa** to range over  $RA \cup AA$ , i.e., assignments to structural features in a model **EM**, and use **sfa.t** to denote the target of the assignment, i.e., if **sfa** is a reference assignment then  $ra_t(\text{sfa})$  and if **sfa** is an attribute assignment then  $aa_t(\text{sfa})$ .

In rule **OBJECTOPERATIONCALL**, calling an operation  $op$  on an object  $o$  evaluates the body expression of the operation  $e_{op}$  with the parameters  $\overline{x}$  replaced by proper arguments.

The evaluation of  $\kappa(\overline{v})$  in rule **KAPPACALL** is defined by the semantics of the pre-defined operation  $\kappa$ . For example, we have  $isDefined(\text{null}) = \text{false}$ ,  $+(1, 2) = 3$ , and  $size(\text{Set}\{\}) = 0$ , etc.

Given a call to the **iterate** operation  $e_1 \rightarrow \text{iterate}(x; y = e_2 \mid e_3)$ , expression  $e_1$  should have a collection value of either the four collection types: **Bag**, **Set**, **Seq**, and **OSet**. In rule **COLLECTIONITERATION**, we write the abstract value

$\mathbf{o.sf} \stackrel{\mathbf{EM}}{=} ?$	$\mathbf{sf}$ multi-valued (i.e., $\mathbf{sf.upperBound} > 1$ or $= -1$ )				$\mathbf{sf}$ single-valued (i.e., $\mathbf{sf.upperBound} = 1$ )
$\exists$ assignments to $\mathbf{sf}$ for $\mathbf{o}$ in EM	$\mathbf{sf.ordered} = \mathbf{true}$  Let $(\mathbf{sfa}_1, \dots, \mathbf{sfa}_n)$ be the chain of assignments to $\mathbf{sf}$ for $\mathbf{o}$ in EM (see Section 4.2)		$\mathbf{sf.ordered} = \mathbf{false}$  Let $\{\mathbf{sfa}_1, \dots, \mathbf{sfa}_n\}$ be the set of assignments to $\mathbf{sf}$ for $\mathbf{o}$ in EM (see Section 4.2)		The value  assigned to $\mathbf{sf}$  for $\mathbf{o}$ in EM
	$\mathbf{sf.unique} = \mathbf{true}$	$\mathbf{sf.unique} = \mathbf{false}$	$\mathbf{sf.unique} = \mathbf{true}$	$\mathbf{sf.unique} = \mathbf{false}$	
	$\mathbf{OSet}\{\mathbf{sfa}_1.t, \dots,$ $\mathbf{sfa}_n.t\}$	$\mathbf{Seq}\{\mathbf{sfa}_1.t, \dots,$ $\mathbf{sfa}_n.t\}$	$\mathbf{Set}\{\mathbf{sfa}_1.t, \dots,$ $\mathbf{sfa}_n.t\}$	$\mathbf{Bag}\{\mathbf{sfa}_1.t, \dots,$ $\mathbf{sfa}_n.t\}$	
	otherwise	$\mathbf{OSet}\{\}$	$\mathbf{Seq}\{\}$	$\mathbf{Set}\{\}$	
					null

**Table 1** Calling structural features (i.e., attributes or references) on an object in a model

$Collection\{\bar{v}\}$  to stand for any of the four kinds of collection values. The result of calling the *iterate* operation on  $Collection\{\bar{v}\}$  is a value obtained by iterating over all elements in the collection. The variable  $x$  is the *iterator*. It goes through the elements in the source collection for each round of the iteration. The variable  $y$  is the *accumulator*. It gets an initial value given by expression  $e_2$  and is used to accumulate results during the iteration. For each element in the source collection referenced by the iterator  $x$ , the body expression  $e_3$  is calculated using the current value of the accumulator  $y$ , and the result is assigned back to the accumulator for the next iteration until the last element of the collection. A simple example of collection iteration is given below to calculate the sum of the elements of a set of integers:

$$\mathbf{Set}\{1, 2, 3\} \rightarrow \mathbf{iterate}(x; y = 0 \mid x + y)$$

### 5.3 CoreOCL invariants

An invariant specified in CoreOCL, written  $e_{inv} = (c, e)$ , comprises a context class  $c$  and a boolean CoreOCL expression  $e$  being the body of the invariant.

To check if a model EM satisfies all the invariants specified for its metamodel, we need to evaluate all the well-formed evaluation points of the form  $(e_{inv}, EM, o)$  (see Definition 5), where  $e_{inv}$  is an invariant defined for the metamodel of EM,  $o$  is an object of EM, and the type of  $o$  is a subtype of the context of  $e_{inv}$ . Evaluating such a well-formed evaluation point amounts to evaluating the boolean body expression of  $e_{inv}$  within the model EM and an *initial evaluation environment*. The domain of the initial evaluation environment consists of the keyword *self* associated with the current contextual object  $o$ . We write  $\mathbf{eval}(e_{inv}, EM, o)$  to denote the result of evaluating  $(e_{inv}, EM, o)$ .

The first part of the result, which is a boolean value, tells whether the invariant  $\text{eInv}$  holds for the contextual object  $\mathbf{o}$  in model  $\text{EM}$ , and the second part of the result is used for the construction of the scope of  $(\text{eInv}, \text{EM}, \mathbf{o})$ , which is an object induced sub-model of  $\text{EM}$  following the definition below.

**Definition 10 (Object induced EMF sub-model)**

We say an EMF model  $\text{EM}' = (\mathbb{M}, \mathbf{O}', \tau_{\mathbf{O}}', \text{RA}', \tau_{\text{ra}}', \text{ra}_t', \text{ra}_s', \text{AA}', \tau_{\text{aa}}', \text{aa}_s', \text{aa}_t')$  is an object induced sub-model (shortly sub-model) of an EMF model  $\text{EM} = (\mathbb{M}, \mathbf{O}, \tau_{\mathbf{O}}, \text{RA}, \tau_{\text{ra}}, \text{ra}_t, \text{ra}_s, \text{AA}, \tau_{\text{aa}}, \text{aa}_s, \text{aa}_t)$ , if and only if:

1.  $\mathbf{O}' \subseteq \mathbf{O}$ ;
2.  $\tau_{\mathbf{O}}'$  is the restriction of  $\tau_{\mathbf{O}}$  to  $\mathbf{O}'$ ;
3.  $\text{RA}' = \{\text{ra} \mid \text{ra} \in \text{RA}, \text{ra}_s(\text{ra}) \in \mathbf{O}', \text{ra}_t(\text{ra}) \in \mathbf{O}'\}$ ;
4.  $\tau_{\text{ra}}', \text{ra}_s'$  and  $\text{ra}_t'$  are respectively the restriction of  $\tau_{\text{ra}}, \text{ra}_s$ , and  $\text{ra}_t$  to  $\text{RA}'$ ;
5. The partial order over  $\text{RA}$  in  $\text{EM}$  is restricted to  $\text{RA}'$ ;
6.  $\text{AA}' = \{\text{aa} \mid \text{aa} \in \text{AA}, \text{aa}_s(\text{aa}) \in \mathbf{O}'\}$ ;
7.  $\tau_{\text{aa}}', \text{aa}_s'$  and  $\text{aa}_t'$  are respectively the restriction of  $\tau_{\text{aa}}, \text{aa}_s$ , and  $\text{aa}_t$  to  $\text{AA}'$ ;
8. The partial order over  $\text{AA}$  in  $\text{EM}$  is restricted to  $\text{AA}'$ .

**5.4 Scopes determine invariant evaluations**

We first prove in the following lemma that the evaluation of a CoreOCL expression only involves the objects of the context model in which the expression is being evaluated.

**Lemma 3** *If an evaluation judgment  $\Lambda \stackrel{\text{EM}}{\models} \mathbf{e} \Downarrow (\mathbf{v}, \mathbf{O}_S)$  holds, then  $\mathbf{O}_S(\mathbf{v}) \subseteq \mathbf{O}_S$ .*

*Proof* Straightforward by induction on the structure of expression  $\mathbf{e}$ . □

**Lemma 4** *Given a model  $\text{EM}$ , an evaluation environment  $\Lambda$ , and an expression of CoreOCL  $\mathbf{e}$ , if  $\Lambda \stackrel{\text{EM}}{\models} \mathbf{e} \Downarrow (\mathbf{v}, \mathbf{O}_S)$  holds, then for any object induced sub-model  $\text{EM}'$  of  $\text{EM}$  where  $\mathbf{O}_S \subseteq \mathbf{O}'$ ,  $\Lambda \stackrel{\text{EM}'}{\models} \mathbf{e} \Downarrow (\mathbf{v}, \mathbf{O}_S)$  also holds.*

*Proof* See Appendix A.1. □

**Theorem 5** *Given two well formed evaluation points:  $(\text{eInv}, \text{EM}, \mathbf{o})$  and  $(\text{eInv}, \text{EM}', \mathbf{o})$ , where  $\text{EM}'$  an object induced sub-model of  $\text{EM}$  and  $\text{EM}'$  includes all the objects of the scope of  $(\text{eInv}, \text{EM}, \mathbf{o})$ ,  $\text{eval}(\text{eInv}, \text{EM}, \mathbf{o}) = \text{eval}(\text{eInv}, \text{EM}', \mathbf{o})$ .*

*Proof* Follows as a corollary of Lemma 4. □

**5.5 CoreOCL invariants without allInstances are forward**

**Theorem 6** *Without calling allInstances, CoreOCL invariants are all forward following Definition 8.*

*Proof* See Appendix A.1. □

**Lemma 7** *Given a model  $\text{EM}$ , an object  $\mathbf{o}$  in it, a CoreOCL expression  $\mathbf{e}$  with no calls to allInstances (i.e.,  $\mathbf{e}$  is an expression written in the sub-language of CoreOCL excluding allInstances), and an evaluation environment*

$\Lambda$  in which values bound to variables only refer to objects reachable from  $\mathbf{o}$  if any, i.e.,  $\forall x$  in the domain of  $\Lambda$ ,  $\forall \mathbf{o}' \in \mathcal{O}_S(\Lambda(x))$ ,  $\mathbf{o}'$  is reachable from  $\mathbf{o}$  (see Definition 7). Suppose  $\Lambda \stackrel{\text{EM}}{\models} e \Downarrow (v, \mathcal{O}_S)$ . We have the following two statements hold:

1.  $\forall \mathbf{o}' \in \mathcal{O}_S$ ,  $\mathbf{o}'$  is reachable from  $\mathbf{o}$ ;
2.  $\forall \mathbf{o}' \in \mathcal{O}_S(v)$ ,  $\mathbf{o}'$  is reachable from  $\mathbf{o}$ .

*Proof* See Appendix A.1.  $\square$

## 6 Model Decomposition in EMF

### 6.1 Abstract EMF for Model Decomposition

In order to exploit the model decomposition technique defined in Section 3 for EMF models, EMF metamodels and EMF models are first abstracted into Definition 1 metamodels and Definition 2 models. Briefly speaking, at the metamodeling level, both classes and data types are mapped to metaclasses in the abstract setting; both references and attributes are mapped to associations; operations are not considered; and invariants are discarded, because invariants are not used in the model decomposition algorithm and hence do not need a counterpart in the abstract setting following the mapping. After the decomposition, we examine whether the invariants keep holding in the sub-models on the EMF side (see Theorem 10 below). At the modeling level, both objects and data values are mapped to instances; and both reference assignments and attribute assignments are mapped to links.

**EMF metamodel abstraction** Given an EMF metamodel  $\text{EM} = (\mathbb{C}, \mathbb{DT}, \mathbb{DV}, \mathbb{Rf}, \mathbb{At})$ , mapping  $\mathbb{F} = (\mathbb{F}_\mathbb{N}, \mathbb{F}_\mathbb{A})$  abstracts it into a metamodel, written as  $\mathbb{F}(\text{EM}) = \mathbb{M} = (\mathbb{N}, \mathbb{A}, \mathbb{H}, \emptyset, \mathbf{s}, \mathbf{t}, \mu_\mathbf{s}, \mu_\mathbf{t}, \emptyset)$ , if the following properties hold.

1.  $\mathbb{C} \cup \mathbb{DT}$  is bijectively mapped to  $\mathbb{N}$  by  $\mathbb{F}_\mathbb{N}$ .
2. Subtyping is preserved, i.e., given two classes  $\mathbf{c}_1$  and  $\mathbf{c}_2$  in  $\mathbb{C}$ ,  $\mathbb{F}_\mathbb{N}(\mathbf{c}_1) \preceq \mathbb{F}_\mathbb{N}(\mathbf{c}_2)$  if and only if  $\mathbf{c}_1 \preceq \mathbf{c}_2$ ; given two data types  $\mathbf{dt}_1$  and  $\mathbf{dt}_2$  in  $\mathbb{DT}$ ,  $\mathbb{F}_\mathbb{N}(\mathbf{dt}_1) \preceq \mathbb{F}_\mathbb{N}(\mathbf{dt}_2)$  if and only if  $\mathbf{dt}_1 \preceq \mathbf{dt}_2$ .
3.  $\mathbb{Rf} \cup \mathbb{At}$  is bijectively mapped to  $\mathbb{A}$  by  $\mathbb{F}_\mathbb{A}$ .
4. Target types and source types are preserved.

Given a reference  $\mathbf{rf} \in \mathbb{Rf}$ ,

- (a)  $\mathbf{t}(\mathbb{F}_\mathbb{A}(\mathbf{rf})) = \mathbb{F}_\mathbb{N}(\mathbf{rf.eType})$ ;
- (b)  $\mathbf{s}(\mathbb{F}_\mathbb{A}(\mathbf{rf})) = \mathbb{F}_\mathbb{N}(\mathbf{rf.eContainingClass})$ .

Given an attribute  $\mathbf{at} \in \mathbb{At}$ ,

- (a)  $\mathbf{t}(\mathbb{F}_\mathbb{A}(\mathbf{at})) = \mathbb{F}_\mathbb{N}(\mathbf{at.eType})$ ;
- (b)  $\mathbf{s}(\mathbb{F}_\mathbb{A}(\mathbf{at})) = \mathbb{F}_\mathbb{N}(\mathbf{at.eContainingClass})$ .

5. Multiplicities are preserved.

Given a reference  $\mathbf{rf} \in \mathbb{Rf}$ ,

- (a) the target multiplicity of  $\mathbb{F}_\mathbb{A}(\mathbf{rf})$  is the same as the multiplicity of  $\mathbf{rf}$ , i.e.,  $\mu_\mathbf{t}(\mathbb{F}_\mathbb{A}(\mathbf{rf})) = (\mathbf{rf.lowerBound}, \mathbf{rf.upperBound})$  in case where  $\mathbf{rf.upperBound} \neq -1$ , otherwise  $\mu_\mathbf{t}(\mathbb{F}_\mathbb{A}(\mathbf{rf})) = (\mathbf{rf.lowerBound}, *)$ ;
- (b) the source multiplicity of  $\mathbb{F}_\mathbb{A}(\mathbf{rf})$  (i.e.  $\mu_\mathbf{s}(\mathbb{F}_\mathbb{A}(\mathbf{rf}))$ ) is  $(0, 1)$  if  $\mathbf{rf}$  is a containment, otherwise  $(0, *)$ .

Given an attribute  $\mathbf{at} \in \mathbb{At}$ ,

- (a) the target multiplicity of  $\mathbb{F}_\mathbb{A}(\mathbf{at})$  is the same as the multiplicity of  $\mathbf{at}$ , i.e.,  $\mu_\mathbf{t}(\mathbb{F}_\mathbb{A}(\mathbf{at})) = (\mathbf{at.lowerBound},$

$\mathfrak{ol}.\text{upperBound}$ ) in case where  $\mathfrak{ol}.\text{upperBound} \neq -1$ , otherwise  $\mu_t(\mathbb{F}_A(\mathfrak{ol})) = (\mathfrak{ol}.\text{lowerBound}, *)$ ;

(b) the source multiplicity of  $\mathbb{F}_A(\mathfrak{ol})$  (i.e.  $\mu_s(\mathbb{F}_A(\mathfrak{ol}))$ ) is  $(0, *)$ .

6. The set of invariants and consequently the invariant context function are empty. Although invariants play a role in determining model conformance (Section 4.3), they are not referenced by the model decomposition algorithm (Section 6.2) hence can be abstracted away.

*EMF models abstraction* Given an EMF model  $\mathbb{EM} = (\mathbb{EM}, \mathbb{O}, \tau_o, \mathbb{RA}, \tau_{ra}, \mathbf{ra}_t, \mathbf{ra}_s, \mathbb{AA}, \tau_{aa}, \mathbf{aa}_s, \mathbf{aa}_t)$ , mapping  $F = (F_N, F_A)$  abstracts it into a model, written as  $F(\mathbb{EM}) = \mathbb{M} = (\mathbb{M}, \mathbb{N}, \mathbb{A}, \tau, \text{src}, \text{tgt})$ , if the following properties hold.

1.  $\mathbb{M} = F(\mathbb{EM})$  is the abstract counterpart of  $\mathbb{EM}$ .
2. The union of  $\mathbb{O}$  (from  $\mathbb{EM}$ ) and  $DV$  (from  $\mathbb{EM}$ ) is bijectively mapped to  $\mathbb{N}$  by  $F_N$ .
3.  $\mathbb{RA} \cup \mathbb{AA}$  is bijectively mapped to  $\mathbb{A}$  by  $F_A$
4. Typing is preserved, i.e.,
  - (a) given an object  $\mathfrak{o} \in \mathbb{O}$ ,  $\tau(F_N(\mathfrak{o})) = F_N(\tau_o(\mathfrak{o}))$ ;
  - (b) given a data value  $\mathfrak{dv} \in DV$ ,  $\tau(F_N(\mathfrak{dv})) = F_N(\tau_d(\mathfrak{dv}))$ ;
  - (c) given a reference assignment  $\mathbf{ra} \in \mathbb{RA}$ ,  $\tau(F_A(\mathbf{ra})) = F_A(\tau_{ra}(\mathbf{ra}))$ ;
  - (d) given an attribute assignment  $\mathbf{aa} \in \mathbb{AA}$ ,  $\tau(F_A(\mathbf{aa})) = F_A(\tau_{aa}(\mathbf{aa}))$ .
5. Source ends and target ends are preserved.

Given a reference assignment  $\mathbf{ra} \in \mathbb{RA}$ ,

(a)  $\text{src}(F_A(\mathbf{ra})) = F_N(\mathbf{ra}_s(\mathbf{ra}))$ ;

(b)  $\text{tgt}(F_A(\mathbf{ra})) = F_N(\mathbf{ra}_t(\mathbf{ra}))$ .

Given an attribute assignment  $\mathbf{aa} \in \mathbb{AA}$ ,

(a)  $\text{src}(F_A(\mathbf{aa})) = F_N(\mathbf{aa}_s(\mathbf{aa}))$ ;

(b)  $\text{tgt}(F_A(\mathbf{aa})) = F_N(\mathbf{aa}_t(\mathbf{aa}))$ .

Note that in EMF, neither references nor attributes come with a definition for their source multiplicities. The way we give source multiplicities to their images in  $\mathbb{M}$  following  $\mathbb{F}_A$  (i.e.,  $(0, *)$  by default and  $(0, 1)$  in case of containment references) reflects the least constrained interpretation of the containment relation following the MOF specification [24] (p. 38). As a consequence, all reference assignments and all attribute assignments are fragmentable (in the sense of Definition 9) from the point of view of the assigned values. This allows a fine-degree decomposition wherever possible. However, if two objects are connected by a pair of opposite links, they become indecomposable (see Condition 4 below).

## 6.2 EMF model decomposition algorithm

Given an EMF model  $\mathbb{EM} = (\mathbb{EM}, \mathbb{O}, \tau_o, \mathbb{RA}, \tau_{ra}, \mathbf{ra}_t, \mathbf{ra}_s, \mathbb{AA}, \tau_{aa}, \mathbf{aa}_s, \mathbf{aa}_t)$ , the following algorithm decomposes  $\mathbb{EM}$  into (object induced) sub-models as defined in Definition 10.

1. Let  $\mathbb{M} = F(\mathbb{EM})$ .
2. Apply the model decomposition algorithm as defined in Section 3.2 to  $\mathbb{M}$  and we get the decomposition hierarchy of  $\mathbb{M}$ .

3. For any sub-model of  $M$ , written  $M' = (\mathbb{M}, N', A', \tau', \text{src}', \text{tgt}')$ , which is induced by the instances grouped in an antichain-node of the decomposition hierarchy of  $M$ , return the corresponding  $EM'$ , where  $EM'$  is a sub-model of  $EM$  induced by all the objects of  $EM$  whose images following  $F_N$  are included in  $M$ , namely  $\{o \mid o \in O, F(o) \in N'\}$ .

### 6.3 Soundness of EMF model decomposition

Given a model  $EM$  conforming to a metamodel  $\mathbb{EM}$  in EMF, given a sub-model  $EM'$  of  $EM$  derived by the EMF model decomposition algorithm of Section 6.2, we demonstrate the soundness of the EMF model decomposition algorithm by proving the conformance of  $EM'$  to the original metamodel  $\mathbb{EM}$ . More precisely, we will achieve the goal in two steps: Firstly, similar to the sufficient conditions elicited in the abstract setting, i.e., Conditions 1, 2, and 3, we propose two conditions in the EMF setting, i.e., Conditions 4 and 5, and prove that these two conditions will be sufficient to ensure conformance of the EMF sub-model  $EM'$ ; Secondly, we prove that the sub-model  $EM'$  derived by the EMF model decomposition algorithm satisfies Condition 4. As a consequence, if the EMF metamodel  $\mathbb{EM}$  satisfies Condition 5, the conformance of  $EM'$  to  $\mathbb{EM}$  can be concluded. Note that we do not need to add a condition corresponding to Condition 3 in the EMF setting because following the definition of EMF abstraction described above in Section 6.1

and the way source multiplicities are given, all links are fragmentable.

The first sufficient condition (Condition 4 defined below and corresponding to Condition 1) guarantees target reachability in sub-models.

**Condition 4**  $\forall ra \in RA, ra_s(ra) \in O' \text{ implies } ra_t(ra) \in O'$ .

The second sufficient condition (Condition 5 defined below and corresponding to Condition 2) restricts the nature of the invariants associated to a metamodel to all be forward.

**Condition 5** *All invariants of metamodel  $\mathbb{EM}$  are forward invariants.*

We prove in the following lemma that Conditions 5 and 4 together suffice to ensure the conformance of sub-models in EMF.

**Lemma 8** *Given an EMF metamodel  $\mathbb{EM} = (\mathbb{C}, \mathbb{DT}, DV, \mathbb{Rf}, \mathbb{At})$ , an EMF model  $EM = (EM, O, \tau_o, RA, \tau_{ra}, ra_t, ra_s, AA, \tau_{aa}, aa_s, aa_t)$ , and an object induced sub-model of  $EM$ , written  $EM' = (EM, O', \tau_o', RA', \tau_{ra}', ra_t', ra_s', AA', \tau_{aa}', aa_s', aa_t')$ , suppose that:*

1.  $EM$  conforms to  $\mathbb{EM}$ ;
2.  $EM'$  satisfies Condition 4;
3.  $\mathbb{EM}$  satisfies Condition 5;

*then  $EM'$  also conforms to  $\mathbb{EM}$ .*

*Proof* See Appendix A.1. □

Moreover, we prove in the following lemma that any sub-model derived by the algorithm satisfies Condition 4.

**Lemma 9** *Given an EMF model  $EM = (\mathbb{EM}, O, \tau_o, RA, \tau_{ra}, ra_t, ra_s, AA, \tau_{aa}, aa_s, aa_t)$  conforming to an EMF metamodel  $\mathbb{EM}$ , and any object induced sub-model of  $EM$  derived by the EMF model decomposition algorithm, written  $EM' = (\mathbb{EM}, O', \tau_o', RA', \tau_{ra}', ra_t', ra_s', AA', \tau_{aa}', aa_s', aa_t')$ ,  $EM'$  satisfies Condition 4.*

*Proof* See Appendix A.1.  $\square$

Summarizing the discussion above, we thus reach the following conclusion about the soundness of the EMF model decomposition algorithm.

**Theorem 10** *Applying our model decomposition technique to a model  $EM$  conforming to a metamodel  $\mathbb{EM}$  in EMF following the algorithm of Section 6.2 produces only conformant sub-models if all the invariants of the metamodel are forward.*

*Proof* Follows from Lemma 8 and Lemma 9.  $\square$

## 7 Application Example: Pruning Based Ecore Model Comprehension

In this section, we demonstrate the power of our generic model decomposition technique by reporting one of its applications in a pruning-based model comprehension method. A typical comprehension question one would like to have answered for a large model is:

*“Given a set of instances of interest in the model, how does one construct a substantially smaller sub-model that is relevant for the comprehension of these instances?”*

Model readers, when confronted with such a problem, would typically start from the interesting instances and browse through the whole model attempting to manually identify the relevant parts. Even with the best model documentation and the support of model browsing tools, such a task may still be too complicated to solve by hand, especially when the complexity of the original model is high. Moreover, guaranteeing by construction that the identified parts (together with the interesting instances) indeed constitute a valid model further complicates the problem.

Our model decomposition technique can be exploited to provide a linear time automated solution to the problem above. The general idea is to simply take the union of all the scc-nodes, each of which contains at least one interesting instance, and their descendant scc-nodes in the decomposition hierarchy of the original model.

We have implemented the idea in an Ecore model comprehension tool [1] based on the implementation of the EMF model decomposition algorithm defined in Section 6.2. In this context, models to be decomposed are Ecore models (i.e. models conforming to `Ecore.ecore`). Interestingly `Ecore.ecore` conforms to itself, in particular to the core part of Ecore as depicted in Figure 5, hence is a metamodel in EMF as defined in Section 4.1.

**Listing 1** Constraint model of Ecore implemented in EssentialOCL (part 1).

---

```

import 'CoreEcore.ecore'

package ecore

context ENamedElement
/*
 * The name can not be empty string
 */
inv WellFormedName:
    self.name.size() > 0

context ETypedElement
/*
 * The lower bound must be greater or equal to 0
 */
inv ValidLowerBound:
    self.lowerBound >= 0
/*
 * The upper bound must be greater than 0 or unlimited (i.e., -1)
 */
inv ValidUpperBound:
    let ub = self.upperBound in
        ub = -1 or ub > 0
/*
 * The lower bound must be less than or equal to the upper bound, unless the upper bound is
    -1
 */
inv ConsistentBounds:
    let lb = self.lowerBound in
    let ub = self.upperBound in
        ub = -1 or lb <= ub
/*
 * Only operations can be without a type to represent void.
 */
inv ValidType:
    self.eType -> isEmpty() implies self.oclIsTypeOf(EOperation)

context EPackage
/*
 * No two sub-packages can have the same name.
 */
inv UniqueSubpackageNames:
    self.eSubpackages -> isUnique(name)
/*
 * No two classifiers can have the same name.
 */
inv UniqueClassifierNames:
    self.eClassifiers -> isUnique(name)

context EOperation
/*
 * No two parameters can have the same name.
 */
inv UniqueParameterNames:
    self.eParameters -> isUnique(name)
/*
 * An operation without a type, which represents void, must have an upper bound of 1.
 */
inv NoRepeatingVoid:
    self.eType -> isEmpty() implies self.upperBound = 1
/*
 * The types of the exceptions raised by an operation are limited to be classes.
 */
inv ValidException:
    self.eExceptions -> forAll(oclIsTypeOf(EClass))

context EClass
def: allFeatures() : Set(EStructuralFeature) =
    self.eSuperTypes -> iterate(
        s: EClass;
        r: Set(EStructuralFeature) = self.eStructuralFeatures -> asSet() |
        r -> union(s.allFeatures())
    )
def: allOperations() : Set(EOperation) =
    self.eSuperTypes -> iterate(

```

---

**Listing 2** Constraint model of Ecore implemented in EssentialOCL (part 2).

---

```

        s: EClass;
        r: Set(EOperation) = self.eOperations -> asSet() |
        r -> union(s.allOperations())
    )
def: allSuperTypes(dejaCalled: Set(EClass)) : Set(EClass) =
    self.eSuperTypes -> iterate(
        s: EClass;
        r: Set(EClass) = self.eSuperTypes -> asSet() |
        if dejaCalled->includes(s)
        then r
        else r->union(s.allSuperTypes(dejaCalled->including(s))
        )
    )
endif
)
/*
 * No two features (attributes or references) can be the same name.
 */
inv UniqueFeatureNames:
    self.allFeatures() -> isUnique(name)
/*
 * No two operations can have the same signature, which is the name of the operation, the
 * number of parameters and their types.
 */
inv UniqueOperationSignatures:
    self.allOperations() -> iterate(
        op1: EOperation;
        r1: Boolean = true |
        r1 and (self.allOperations() -> iterate(
            op2: EOperation;
            r2: Boolean = true |
            r2 and (op1 = op2
                or op1.name <> op2.name
                or let op1Parameters = op1.eParameters in
                    let op2Parameters = op2.eParameters in
                        op1Parameters -> size() <> op2Parameters ->
                            size()
                or op1Parameters.eType <> op2Parameters.eType
            )
        )
    )
)
/*
 * A class can not be a super type of itself.
 */
inv NoCircularSuperTypes:
    self.allSuperTypes(Set{self}) -> excludes(self)

context EEnum
/*
 * No two literals can have the same name.
 */
inv UniqueEnumeratorNames:
    self.eLiterals -> isUnique(name)

context EAttribute
/*
 * The type of an attribute must be a data type
 */
inv ValidType:
    self.eType.ocIsTypeOf(EDDataType)

context EReference::container: EBoolean
derive: if eOpposite.ocIsUndefined()
    then false
    else eOpposite.containment
endif

context EReference
/*
 * The type of a reference must be a class.
 */
inv ValidType:
    self.eType.ocIsTypeOf(EClass)
/*
 * If opposite exists,
 *** it must be a feature of this references's type,
 *** it must have this reference as its opposite,

```

---

**Listing 3** Constraint model of Ecore implemented in EssentialOCL (part 3).

---

```

*** this reference and its opposite can not both be containments
*/
inv ConsistentOpposite:
    (self.eOpposite -> isEmpty())
    or (let oppoRef = self.eOpposite in
        oppoRef.eContainingClass = self.eType
        and oppoRef.eOpposite = self
        and (self.containment implies not oppoRef.containment)
    )
/*
 * A container reference (i.e., whose opposite is containment) must have upperbound = 1
 */
inv SingleContainer:
    self.eOpposite -> isEmpty()
    or self.containment implies self.eOpposite.upperBound = 1
/*
 * A reference that may have multiple values must be unique if it is containment or has
   opposite
 */
inv ConsistentUnique:
    ((self.upperBound = -1 or (self.upperBound - self.lowerBound > 1))
    and (self.containment or not self.eOpposite -> isEmpty()))
    implies self.unique
/*
 * If a container reference is required (i.e., lower bound >=1), the contained class can not
   have other container references
 */
inv ConsistentContainer:
    (self.container and self.lowerBound >=1) implies
        self.eContainingClass.allFeatures() -> forAll (f |
            f.oclIsTypeOf(EReference) implies not f.oclAsType(EReference).container
        )
endpackage

```

---

Consequently, Ecore models are all EMF models in the sense of Section 4.2. Briefly speaking, the tool takes as input an Ecore model and a set of objects in the model marked as interesting for a comprehension task. It runs the EMF model decomposition algorithm, where in the third step, the minimal sub-model in the sub-model lattice which contains all the instances corresponding to a marked object is selected, and returns the result.

### 7.1 Ecore invariants are all forward

Invariants pertinent to Ecore are identified in `Ecore.ecore` as annotations nested in the context classes and directly implemented in `EcoreValidator.java`. For example, the following snippet extracted from `Ecore.ecore` summarizes

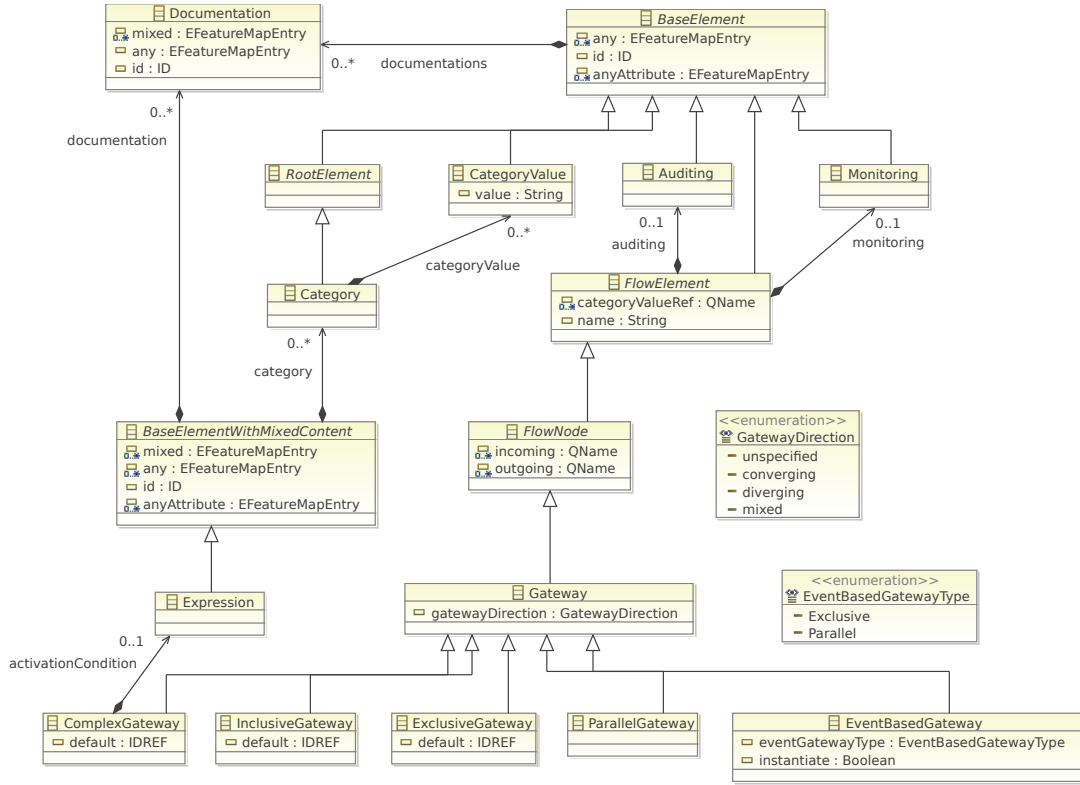
the invariants of class `ETypedElement`, with the following names: `ValidLowerBound`, `ValidUpperBound`, `ConsistentBounds`, and `ValidType`.

```

<eClassifiers xsi:type="ecore:EClass" name="
ETypedElement" abstract="true"
eSuperTypes="#//ENamedElement">
  <eAnnotations source="http://www.eclipse.
    org/emf/2002/Ecore">
    <details key="constraints" value="
      ValidLowerBound ValidUpperBound
ConsistentBounds ValidType" />
  </eAnnotations>
  .....
</eClassifiers>

```

We re-implemented the constraint model of Ecore in EssentialOCL. Please refer to Listings 1, 2, and 3 for the details of invariant specifications (started with keyword `inv`) and helper operation specifications (started with keyword `def`) in the constraint model, and to Section 5 for the formal counterpart of EssentialOCL. Thanks to



**Fig. 7** Pruned class diagram for understanding the Gateway concept in BPMN.

the formal treatment of EssentialOCL in terms of Core-OCL, a precise and formal assessment of the forwardness of the Ecore invariants becomes straightforward, i.e., by checking if `allInstances` is called in the body of the invariant expression.

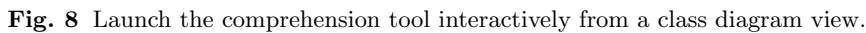
We notice that none of the Ecore invariants in Listings 1, 2, 3 makes use of `allInstances`. Therefore, according to Theorem 6, they are all forward invariants. As a consequence of this fact and by Theorem 10, applying the EMF model decomposition technique to models conforming to Ecore following the steps discussed in Section 6.2 only returns conformant sub-models.

## 7.2 Case study 1: BPMN

As the first case study we have chosen to comprehend the *Gateway* concept in BPMN [26]: business process model and notation. Gateways are modeling elements in BPMN used to control how sequence flows interact as they converge and diverge within a business process. Five types of gateways are identified in order to cater to different types of sequence flow control semantics: exclusive, inclusive, parallel, complex, and event-based.

Inputs to the comprehension tool for the case study are the following:

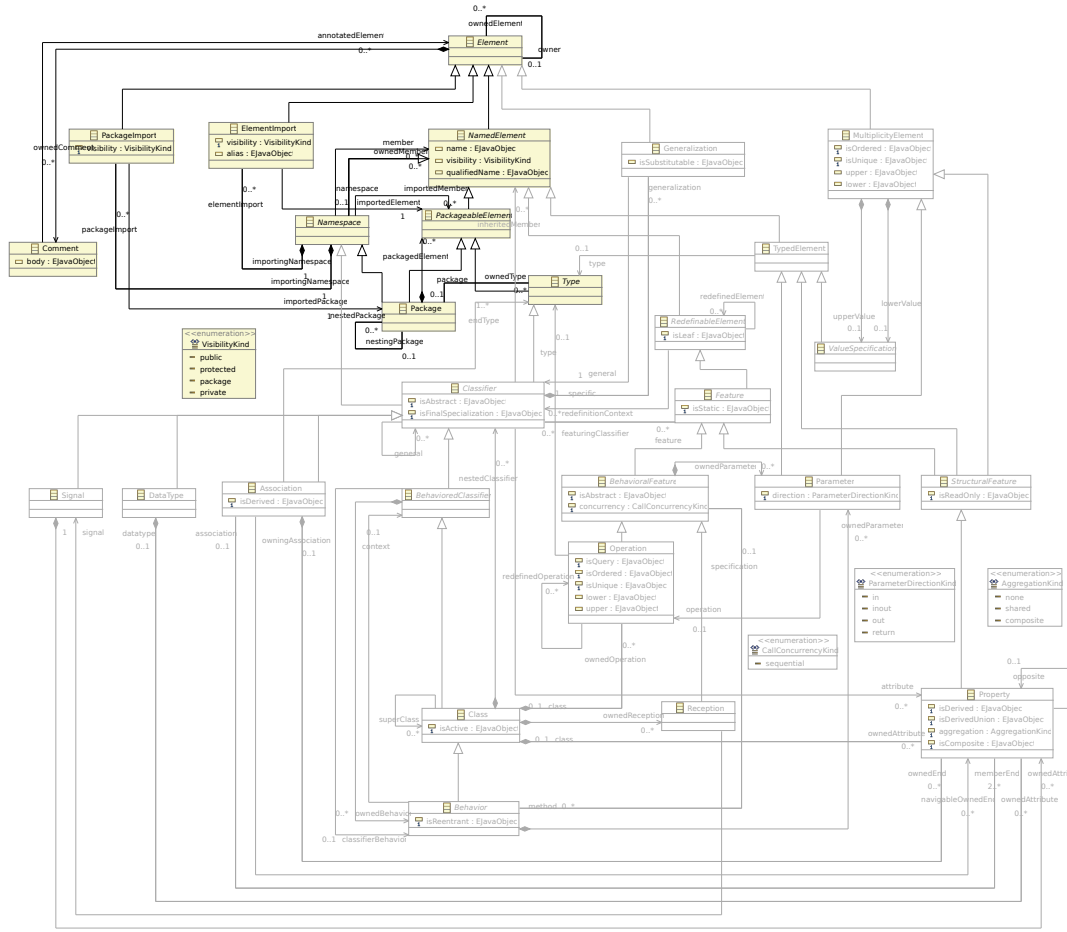
- The BPMN Ecore model (`bpmn.ecore`) containing 156 classifiers (of which 134 are objects of `EClass`, 11 are objects of `EEnum`, and 11 are objects of `EDataType`),



- a set of EClass objects capturing the key notions of the design of gateways in BPMN: Gateway, ExclusiveGateway, InclusiveGateway, ParallelGateway, ComplexGateway, and EventBasedGateway.

dent concepts of BPMN such as *Activity*, *Event*, *Connector*, and *Artifact*, are pruned out. The class diagram view of the pruned BPMN model is shown in Figure 7. Note that it corresponds well to the class diagram that is sketched in the chapter for describing gateways in the BPMN specification [26]. We have also verified that the pruned BPMN model is indeed a valid Ecore model by calling the Ecore Validator in EMF.

As the second case study we have chosen to study fUML [25]: a subset of executable UML models. We take two steps. In the first step, the key concept `Class` is selected as the seed to produce a smaller sub-model of fUML. Inputs to



**Fig. 9** Pruned class diagram for understanding the Namespace concept in fUML with all irrelevant parts faded out.

the comprehension tool for the first step are the following:

- The fUML Ecore model (fuml.ecore) containing 109 classifiers (in which 104 are objects of EClass and 5 are objects of EEnum), 160 references (objects of EReference), and 56 attributes (objects of EAttribute). Altogether, it adds up to a class diagram that is too large to fit on the screen.
- the EClass object with name Class.

The sub-model returned by the tool contains 28 classes, 4 enumerations, 60 references, and 32 attributes. The computation of the first step takes around 0.04 seconds.

Although the number of nodes does not decrease dramatically after this step, the result becomes a class diagram that fits on the screen (see Figure 8).

We then call the comprehension tool one more time on the concept **Namespace**. Figure 8 illustrates how the tool allows a user to launch the pruning functionality interactively from the Sample Ecore Model Editor in EMF, and Figure 9 shows the result where all irrelevant parts for comprehending namespace are faded out. The highlighted part in Figure 9 contains 9 classes, 1 enumerations, 19 references and 7 attributes, and corresponds well to the class diagram appearing in the fUML speci-

cation [25] (p. 25). The computation of the second step takes around 0.3 seconds including updating the class diagram view in Sample Ecore Model Editor.

## 8 Discussion

### 8.1 About the inclusion conditions

When an instance is included in a sub-model, we decide if a neighbour instance (i.e., an instance connected by either an incoming or an outgoing link) should also be included based on two conditions: Conditions 1 and 3, for the purpose of preserving the conformance of the sub-model. Recall that Condition 1 requires to always follow outgoing links and include also the target instances, and Condition 3 requires to always follow non-fragmentable incoming links and include also the source instances. However, these two conditions are not always necessary. For example, regarding Condition 1, if an outgoing link is *target optional*, i.e., its type has 0 as the lower bound for its target multiplicity, and is not referred to in any invariant evaluation scope, then it is not necessary to also include the target instance to ensure conformance.

The generic nature of our technique and the attempt to avoid sophisticated analysis of invariants account for the existence of such a gap. On one hand, we target a general solution that should work for models expressed in any metamodels. On the other hand, arbitrary invariants can be specified for a metamodel. We opt for a lightweight strategy to automatically guarantee the preserva-

tion of invariants in sub-models by simply restricting the invariants inside a property, i.e., being forward, which can be checked easily at the syntax level as proved by Theorem 6. (See Section 8.2 below for discussion on the implication of this restriction.) The evaluation scopes of forward invariants are always reachable from the context instances, and Condition 1 ensures the inclusion of all reachable instances from the context. These two together guarantee a forward invariant always has the same evaluation scope in a sub-model as in the original model, hence preserved.

Egyed reports in [14] an efficient technique to automatically compute evaluation scopes through model profiling. This technique could be combined with our model decomposition technique to reduce the gap from being sufficient to being necessary. More specifically, for all target optional outgoing links from an instance already in a sub-model, the corresponding target instances must be included as well if and only if the latter appear in at least one evaluation scope of the former.

### 8.2 About `allInstances`

Theorem 6 allows a user to simply check the occurrence of calls to `allInstances` in the invariant body in order to decide if an invariant written in CoreOCL (and eventually EssentialOCL) is forward (and eventually to decide if the decomposition technique can be applied). However, the applicability of our model decomposition technique is not limited to only `allInstances`-free invariants.

Firstly, not all calls to `allInstances` are necessary. A typical example of abuse is shown below where `P` stands for some property:

```
context A
inv: A.allInstances() -> forAll(a | P(a))
```

This invariant can be expressed in an equivalent way without using `allInstances` as follows:

```
context A
inv: P(self)
```

*Normalization* is a process that removes unnecessary calls to `allInstances` in a set of invariants and results in an equivalent set of invariants. Two sets of invariants are equivalent if they accept/reject the same set of models. After normalization, the power of forward invariants is the same as the power of invariants excluding `allInstances`.

Secondly, invariants with necessary calls to `allInstances` do not always stop holding after sub-modeling. Actually, a big number of them are *monotonic*, i.e., if an invariant holds for a model then it holds for any sub-models. A typical example of a monotonic invariant checks for identity uniqueness of all instances of a type in a model.

Only in case when invariants with calls to `allInstances` are both necessary and non-monotonic, the conformance of sub-models is not automatically guaranteed. However, we can still apply our technique but need to re-validate the concerned invariants on the sub-models produced by the decomposition.

We have carried out an empirical study on a pool of 707 invariants written in OCL collected from various

sources ranging from OMG (e.g., metamodel specifications of UML, MOF, and OCL), to academic community (e.g., metamodel of RBAC for role-based access control defined at University of Bremen and metamodel of B language defined at IMAG), to industrial community (e.g., the SAM metamodel from the Topcased open source software project). Out of these 707 invariants studied, only 49 call `allInstances`. And among the 49, 5 can be eliminated by normalization, 39 are monotonic, and only 5 left to be re-validated, which amounts to an overall percentage of 0.7%. Table 2 in the appendix gives a detailed break-down of these numbers.

### 8.3 EMF opposite links

The model decomposition technique works with associations that have both a source end and a target end. However, in EMF, the notion of a source end is indirectly captured by an opposite reference. This brings consequences to the model decomposition technique. For example, when applying our model decomposition technique in the Ecore model comprehension tool, due to the universal existence of opposite references to all containment references in Ecore (see Figure 5), any Ecore model is a strongly connected graph and consequently, the inclusion of any object in an Ecore model eventually requires the inclusion of the whole model. (Note that it is not a problem caused by containment references themselves, whose links are fragmentable according to our definition in Section 6.1.) A solution to this is to

introduce extra pre-/post-processing before/after model decomposition.

In our Ecore model comprehension application, we omitted two references in the pre-processing: `eClassifiers` (from `EPackage` to `EClassifier`) and `eSubPackages` (from `EPackage` to `EPackage`). This omission breaks the strong bi-directional connection between packages and their members, being either sub-packages or classifiers. It opens up opportunities for a desired degree of decomposition because the inclusion of a member in a package does not imply the inclusion of all the other members in the same package any more (as the outgoing links from packages to their members are omitted). Let us refer to the abridged metamodel by  $Ecore'$ . For each Ecore model  $M$ , a corresponding  $Ecore'$  model  $M'$  removes from  $M$  all links of the two aforementioned references. Then  $M'$  conforms to  $Ecore'$  if  $M$  conforms to  $Ecore$ . The actual input for the decomposition is  $Ecore'$  (as the metamodel) and  $M'$  (as the model). Note that we have carefully chosen to omit `eClassifiers` and `eSubPackages` instead of their opposites, i.e., `ePackage` and `eSuperPackage`. Omission of the latter could also bring a good degree of decomposition but would lose the package hierarchy in the result sub-models. Namely, the inclusion of a member does not require the inclusion of its containing package because containment references are fragmentable.

After decomposition, in the post-processing, links of the omitted references can be restored. For each sub-model of  $M'$  in the sub-model lattice, called  $M'$ -sub, a

corresponding model  $M$ -sub adds for each link of `ePackage` (from `EClassifier` to `EPackage`) an opposite link instantiating `eClassifiers` (from `EPackage` to `EClassifier`) and for each link of `eSuperPackage` (from `EPackage` to `EPackage`) an opposite link instantiating `eSubPackages` (from `EPackage` to `EPackage`). Then  $M$ -sub conforms to  $Ecore$  if  $M'$ -sub conforms to  $Ecore'$ .

Pre-/post-processings customize our general purpose solution to model decomposition to work in a specific situation and/or for a specific need. They are application specific, and therefore cannot be captured in a general sense but need to be specified by the user case by case. Moreover, extra effort may become necessary to show that the pre-/post-processings do not influence any properties established in the general settings. In the Ecore model comprehension application discussed above, it amounts to demonstrating the following two properties: (1)  $M'$  conforms to  $Ecore'$  if  $M$  conforms to  $Ecore$ , (2)  $M$ -sub conforms to  $Ecore$  if  $M'$ -sub conforms to  $Ecore'$ , which hold straightforwardly by simply checking against the conformance conditions listed in Section 4.3.

## 9 Related Work

*Model slicing* In general model slicing consists in identifying sub-models of a model that satisfy certain properties. As such it is a generalization of the work on program slicing to the domain of models. The slicing criterion for model slicing depends on the purpose of the

slicing process. In our case the primary purpose is model comprehension.

Other authors have investigated other uses of model slicing. For instance in [31] the goal is to check satisfiability of a UML class diagram equipped with a set of OCL constraints, i.e., to check the existence of an instance of the class diagram that satisfies all the integrity constraints. Each slice is a valid UML class diagram with constraints belonging to the same slice if they constrain the same model element. The work of [23] strives to establish traceability links between safety requirements and software design elements. The slicing criterion is the inclusion of elements from a SysML model that are relevant for a specific safety requirement.

Another difference between our technique and existing approaches is its genericity: it is not restricted to a particular metamodel. Existing work on model slicing generally applies only to a particular modeling language. Two examples of UML-based model slicing approaches are [20] where model slicing of UML class diagrams is investigated and [5] which considers the problem of slicing the UML metamodel into metamodels corresponding to the different diagram types in UML. With the emergence of an increasing number of domain specific modeling languages genericity becomes an important issue. The need for a generic model slicing technique was identified in [7] which proposed a language for modeling model slicers, allowing for automatic generation of model slicers for any given metamodel.

Note that some work in model slicing allows rewriting of the original model when looking for “sub-models”. Such slices are called amorphous slices [4]. These approaches clearly fall outside the scope of our paper since they do not necessarily produce a proper sub-model of the original model.

*Metamodel pruning* In a similar line of work some authors have investigated the possibility of pruning metamodels in order to make them more manageable. The idea is to remove elements from a metamodel to obtain a minimal set of modeling elements containing a given subset of elements of interest. Such an approach is described in [30]. This work differs from our work in several respects: first, just like model slicing it focuses on a single model rather than considering the collection of relevant sub-models in its totality; second, it is less generic in the sense that it restricts its attention to Ecore metamodels (and the pruning algorithm they present is very dependent on the structure of Ecore), and lastly their goal is not just to get a conformant sub-model but rather to find a sub-metamodel that is a supertype of the original model. This added constraint is due to the main use of the sub-metamodel in model transformation testing.

*Model abstraction* The general idea of simplifying models (which can be seen as a generalization of model slicing and pruning) has also been investigated in the area of model abstraction (see [18] for an overview). In the area of simulation model abstraction is a method for reducing

the complexity of a simulation model while maintaining the validity of the simulation results with respect to the question that the simulation is being used to address. Work in this area differs from ours in two ways: first, model abstraction techniques generally transform models and do not necessarily result in sub-models; second, conformance of the resulting model with a metamodel is not the main concern but rather validity of simulation results.

*Constraint evaluation scope* A profiling based technique was implemented in [14] for the computation of evaluation scopes of invariants, where a model profiler monitors what model elements the invariant evaluation engine accesses and logs the accessed model elements in a scope database. Treating all invariants as black boxes, the profiling based technique is independent of the language in which the invariants are expressed. Complementarily, we demonstrate how scopes are computed formally when invariants are all white boxes, i.e., both the syntactic representations and the formal semantics of invariant evaluation are accessible.

In [14], scopes were used to detecting inconsistencies (i.e., breaking of invariants) that occurs during model changes. In our context, we exploit scopes of invariant evaluation to determine if a sub-model automatically preserves invariants satisfied by the original model. More specifically, Theorem 5 demonstrates that letting a sub-

model contain the scopes of invariant evaluations suffices to preserve the invariants hold in the original model.

*Metamodel/model formalization* There has been many attempts at full formalization of metamodeling/modeling. Here we list a few representatives: [11] formally defined the Meta-Modeling Language (a core of UML 2.0) based on the  $\zeta$ -calculus of Cardelli and Abadi; [19] and [6] followed a graph-based approach to formalizing metamodels/models; [2] provided a rich set-theoretic setting for metamodeling; and [8] presented an algebraic semantics of MOF using membership equational logic and term rewriting as a formal foundation.

Our formalization of metamodels and models in Section 2 can be considered as a simplified version of graph-based formalization exploited by work on model transformations such as [6], where a metamodel is expressed as an attributed type graph with inheritance, composition, and multiplicities and a model is captured as a graph typed by the corresponding type graph. In comparison, we omit attributes, containment relations, and operations in metamodels.

*Formalizing OCL* CoreOCL provides a formal account of EssentialOCL to allow detailed elaborations of evaluation semantics and evaluation scopes. There are various work in the literature formalizing OCL from different aspects: A set theory based denotational semantics for OCL was first defined in [28] then extended in [29]. A type inference system and a big-step operational se-

antics for OCL 1.4 were defined in [10], [17] and [15] completed existing work by formalizing the semantics of OCL messages. In [16] a formal semantics is provided for state-based temporal OCL expressions. A graph-based semantics for OCL was developed by translating OCL constraints into expressions over graph rules in [9], and [13] provided a formal semantics of OCL by defining a mapping from OCL to a temporal logic.

## 10 Conclusion and Future Work

A generic model decomposition technique is put forward in this paper to work with metamodels and models that can be abstracted into a simplified graph-based formalism. A formal foundation is established to demonstrate properties of the technique including correctness and conformance of derived sub-models. We position this generic model decomposition technique as an infrastructural service that can be exploited in many concrete application domains in model-driven software development. A detailed instantiation of the technique to EMF and EssentialOCL and an application of the instantiated model decomposition technique for Ecore model comprehension are presented in this paper. We provide tool implementation and report the results of two concrete case studies on BPMN and fUML comprehension.

We plan to validate the Ecore model comprehension method with more examples from model zoos such as the Repository for Model Driven Development (ReMoDD) <sup>5</sup>

and the AtlanMod metamodel zoos <sup>6</sup> to have a better idea of the applicability of the method and the efficiency of the tool implementation, and to discover potential limitations.

Beyond the domain of model comprehension described in this paper, the generic model decomposition technique, being a fundamental facilitating technique, can have more applications in other kinds of tasks in the life cycle of model-driven software development. We mention in the following some examples. (1) Given a software that takes models conforming to a metamodel as input (e.g., a model transformation), an important part in testing the software is test case generation. Our model decomposition technique could help with the generation of new test cases by using one existing test case as the seed. New test cases of various complexity degree could be automatically generated following the sub-model lattice of the seed test case. (2) Our model decomposition could also help with the debugging activity when a failure of the software is observed on a test case. The idea is that we will find a sub-model of the original test case which is responsible for triggering the bug. Although both the reduced test case and the original one are relevant, the smaller test case is easier to understand and investigate. (3) A major obstacle to the massive model reuse in model-based software engineering is the cost of building a repository of reusable model components. A more effective alternative to creating those reusable model com-

<sup>5</sup> <http://www.cs.colostate.edu/remodd/v1/>

<sup>6</sup> <http://www.emn.fr/z-info/atlanmod/index.php/Zoos>

ponents from scratch is to discover them from existing system models. Sub-models of a system extracted by following our model decomposition technique are all guaranteed to be valid models hence can be wrapped up into modules and reused in the construction of other systems following our modular model composition paradigm [21].

(4) Finally, applying the model decomposition technique to system models in multi-view modeling so that distinct and separate sub-models capturing different aspects of a system can be extracted is a potential avenue for future work.

Implementation at this stage is only for proof of concept. As a consequence, neither user-friendliness nor efficiency is of high priority in the current version (although the performance is still agreeable, i.e., within a second for the case studies). Optimization is also in our future work agenda. We plan to equip the comprehension tool with a better graphical user interface and provide the tool as an Eclipse plugin to reach a wider audience.

## References

1. Democles tool. <http://democles.lassy.uni.lu/>.
2. Marcus Alanen and Ivan Porres. A metamodeling language supporting subset and union properties. *Software and System Modeling (SoSyM)*, 7(1):103–124, 2008.
3. Moussa Amrani. A Formal Semantics of Kermeta. In Marjan Mernik, editor, *Formal and Practical Aspects of Domain-Specific Languages: Recent Developments*, chapter 10. IGI Global, 2012.
4. Kelly Androutsopoulos, David Clark, Mark Harman, Robert M. Hierons, Zheng Li, and Laurence Tratt. Amorphous slicing of extended finite state machines. *IEEE Transactions on Software Engineering*, 99(PrePrints):1, 2012.
5. Jung Ho Bae, KwangMin Lee, and Heung Seok Chae. Modularization of the UML metamodel using model slicing. *Fifth International Conference on Information Technology: New Generations*, 0:1253–1254, 2008.
6. Enrico Biermann, Claudia Ermel, and Gabriele Taentzer. Formal foundation of consistent emf model transformations by algebraic graph transformation. *Software and Systems Modeling (SoSyM)*, 11(2):227–250, 2012.
7. Arnaud Blouin, Benoît Combemale, Benoit Baudry, and Olivier Beaudoux. Modeling model slicers. In *the Proceedings of ACM/IEEE 14th International Conference on Model Driven Engineering Languages and Systems (MoDELS 2011)*, pages 62–76, 2011.
8. Artur Boronat and José Meseguer. An algebraic semantics for MOF. *Formal Aspects of Computing*, 22(3-4):269–296, 2010.
9. Paolo Bottoni, Manuel Koch, Francesco Parisi-Presicce, and Gabriele Taentzer. Consistency checking and visualization of OCL constraints. In *the Proceedings of 3rd International Conference on The Unified Modeling Language: Advancing the Standard (UML 2001)*, LNCS 1393, pages 294–308, 2000.
10. María Victoria Cengarle and Alexander Knapp. A formal semantics for OCL 1.4. In *the Proceedings of 4th International Conference on the Unified Modeling Language: Modeling Languages, Concepts, and Tools (UML*

- 2001), LNCS 2185, pages 118–133, 2001.
11. Tony Clark, Andy Evans, and Stuart Kent. The meta-modelling language calculus: Foundation semantics for UML. In *the Proceedings of 4th International Conference on 4th International Conference (FASE 2001)*, LNCS 2029, pages 17–31, 2001.
  12. Reinhard Diestel. *Graph Theory Fourth Edition*. Springer-Verlag, 2010.
  13. Dino Distefano, Joost-Pieter Katoen, and Arend Rensink. On a temporal logic for object-based systems. In *the Proceedings of 4th International Conference on Formal Methods for Open Object-Based Distributed Systems (FMOODS 2000)*, pages 305–325, 2000.
  14. Alexander Egyed. Automatically detecting and tracking inconsistencies in software design models. *IEEE Transactions on Software Engineering*, 37(2):188–204, 2011.
  15. Stephan Flake. Towards the completion of the formal semantics of OCL 2.0. In *the Proceedings of 27th Australasian Conference on Computer science (ACSC) '04*, pages 73–82, 2004.
  16. Stephan Flake and Wolfgang Mueller. Formal semantics of static and temporal state-oriented OCL constraints. *Software and Systems Modeling (SoSyM)*, 2:164–186, 2003.
  17. Stephan Flake and Wolfgang Müller. Formal semantics of OCL messages. *Electronic Notes in Theoretical Computer Science*, 102:77–97, 2004.
  18. F.K. Frantz. A taxonomy of model abstraction techniques. *Winter Simulation Conference*, 0:1413–1420, 1995.
  19. Frédéric Jouault and Jean Bézivin. KM3: A DSL for metamodel specification. In *the Proceedings of 8th International Conference on Formal Methods for Open Object-Based Distributed Systems (FMOODS 2006)*, LNCS 4037, pages 171–185, 2006.
  20. Huzefa Kagdi, Jonathan I. Maletic, and Andrew Sutton. Context-free slicing of UML class models. In *ICSM '05: Proceedings of the 21st IEEE International Conference on Software Maintenance*, pages 635–638, Washington, DC, USA, 2005. IEEE Computer Society.
  21. Pierre Kelsen and Qin Ma. A modular model composition technique. In *the Proceedings of 13th International Conference on Fundamental Approaches to Software Engineering, (FASE 2010)*, volume LNCS 6013, pages 173–187, 2010.
  22. Pierre Kelsen, Qin Ma, and Christian Glodt. Models within models: Taming model complexity using the sub-model lattice. In *the Proceedings of 14th International Conference on Fundamental Approaches to Software Engineering, (FASE 2011)*, volume LNCS 6603, pages 171–185, 2011.
  23. Shiva Nejati, Mehrdad Sabetzadeh, Davide Falessi, Lionel C. Briand, and Thierry Coq. A sysml-based approach to traceability management and design slicing in support of safety certification: Framework, tool support, and case studies. *Information & Software Technology*, 54(6):569–590, 2012.
  24. OMG. Meta Object Facility (MOF) Core Specification Version 2.4.1, 2011.
  25. OMG. Semantics of a Foundational Subset for Executable UML Models (fUML), v1.0, 2011.
  26. OMG. Business Process Model and Notation (BPMN) Version 2.0, 2012.
  27. OMG. Object Constraint Language Version 2.3.1, 2012.

28. Mark Richters and Martin Gogolla. On formalizing the UML object constraint language OCL. In *the Proceedings of 17th International Conference on Conceptual Modeling (ER '98)*, pages 449–464, 1998.
29. Mark Richters and Martin Gogolla. OCL: Syntax, semantics, and tools. In *the Proceedings of Object Modeling with the OCL, The Rationale behind the Object Constraint Language (OCL 2002)*, LNCS 2263, pages 42–68, 2002.
30. Sagar Sen, Naouel Moha, Benoit Baudry, and Jean-Marc Jézéquel. Meta-model pruning. In *the Proceedings of the 12th international conference on Model Driven Engineering Languages and Systems (MoDELS 2009)*, pages 32–46, 2009.
31. Asadullah Shaikh, Robert Clarisó, Uffe Kock Wiil, and Nasrullah Memon. Verification-driven slicing of UML/OCL models. In *the Proceedings of 25th IEEE/ACM International Conference on Automated Software Engineering (ASE 2010)*, pages 185–194, 2010.
32. Dave Steinberg, Frank Budinsky, Marcelo Paternostro, and Ed Merks. *EMF: Eclipse Modeling Framework, 2nd Edition*. Addison-Wesley Professional, 2008.
33. Robert Endre Tarjan. Depth-first search and linear graph algorithms. *SIAM J. Comput.*, 1(2):146–160, 1972.
34. Jos Warmer and Anneke Kleppe. *The Object Constraint Language: Getting Your Models Ready for MDA*. Addison-Wesley Professional, 2003.

## A Appendix

### A.1 Proofs

**Theorem 2** *Given a model  $M = (\mathbb{M}, N, A, \tau, \text{src}, \text{tgt})$  and an instance induced sub-model  $M' = (\mathbb{M}, N', A', \tau', \text{src}', \text{tgt}')$  of  $M$ ,  $M'$  satisfies both Conditions 1 and 3 if and only if there exists a corresponding antichain-node of the decomposition hierarchy of  $M$  where  $M'$  is induced by all the instances grouped in this antichain-node.*

*Proof* We first demonstrate that if  $M'$  is induced by the set of instances that are grouped in an antichain-node  $\alpha$  of the decomposition hierarchy of  $M$ , then  $M'$  satisfies both Conditions 1 and 3.

- Check  $M'$  against Condition 1: following the algorithm in Section 3.2 to compute the decomposition hierarchy, given a link  $a \in A$ ,  $\text{src}(a)$  and  $\text{tgt}(a)$  are either grouped in one scc-node, or the scc-node  $s_1$  that groups  $\text{src}(a)$  is different from the scc-node  $s_2$  that groups  $\text{tgt}(a)$  and  $s_2$  is a descendant of  $s_1$  because of the presence of  $a$  (which must be a fragmentable link). In both cases, grouping  $\text{src}(a)$  in an antichain-node implies grouping also  $\text{tgt}(a)$  in the same antichain-node. Therefore, if  $\text{src}(a) \in N'$ , i.e.,  $\text{src}(a)$  is grouped in the antichain-node  $\alpha$ , so is  $\text{tgt}(a)$ , i.e.,  $\text{tgt}(a) \in N'$ .
- Check  $M'$  against Condition 3: following the algorithm in Section 3.2 to compute the decomposition hierarchy, given a link  $a \in A$  that is non-fragmentable,

$\text{src}(a)$  and  $\text{tgt}(a)$  are grouped in one wcc-node, hence are grouped in one scc-node, hence will be grouped always together in any antichain-nodes. Therefore, if  $\text{tgt}(a) \in N'$ , i.e.,  $\text{tgt}(a)$  is grouped in the antichain-node  $\alpha$ , so is  $\text{src}(a)$ , i.e.,  $\text{src}(a) \in N'$ .

We now demonstrate the other direction of the theorem, namely, if  $M'$  satisfies both Conditions 1 and 3, then there exists an antichain-node of the decomposition hierarchy of  $M$  such that  $M'$  is induced by the set of instances that are grouped in this antichain-node.

Let  $S$  refer to the set of scc-nodes in the decomposition hierarchy, each of which includes at least one instance of  $M'$ .

1. All the instances that are grouped in an scc-node in  $S$  belong to  $M'$ . Indeed given an scc-node  $s \in S$ , there must exist an instance  $n$  grouped in  $s$  and  $n \in N'$  in order for  $s$  to be included in  $S$ . Let  $n'$  be another instance grouped in  $s$ . Following the algorithm in Section 3.2 to compute the decomposition hierarchy,  $n$  and  $n'$  are grouped in one scc-node either in the first or the second step.
  - (a) If they are grouped in the first step, that means the two instances are weakly connected by non-fragmentable links, and because  $M'$  satisfies Conditions 1 and 3,  $n'$  should also be in  $M'$ .
  - (b) If they are grouped in the second step but not in the first step, that means  $n$  and  $n'$  are grouped in two separate wcc-nodes in the first step, called  $w$

and  $w'$ , which are strongly connected by a path of fragmentable links. Referring to the other wcc-nodes on the path by  $w_1, \dots, w_k$ , there exists a set of instances  $n_0 \in w$ ,  $n'_0 \in w'$ , and  $n_i, n'_i \in w_i$  for  $1 \leq i \leq k$ , such that the following fragmentable links exist: from  $n_0$  to  $n_1$ , from  $n'_i$  to  $n_{i+1}$  ( $\forall i. 1 \leq i < k$ ), and from  $n'_k$  to  $n'_0$ . Since  $M'$  satisfies Condition 1, if the source instance of these fragmentable links belongs to  $M'$ , so does the target instance. Because the following pairs of instances:  $n$  and  $n_0$ ,  $n_i$  and  $n'_i$  ( $\forall i. 1 \leq i \leq k$ ), and  $n'_0$  and  $n'$ , are respectively grouped in a wcc-node, if one instance in a pair belongs to  $M'$  then the other instance in the pair must belong to  $M'$  as well following Conditions 1 and 3. From the above discussion and by applying mathematical induction,  $n$  belonging to  $M'$  implies that  $n'$  belongs to  $M'$  as well.

2.  $S$  constitutes an antichain plus descendant. We partition  $S$  into two subsets:  $S_1$  contains all the scc-nodes in  $S$  that do not have another scc-node also in  $S$  as ancestor;  $S_2$  contains the rest, i.e.,  $S_2 = S \setminus S_1$ . Clearly  $S_1$  constitutes an antichain, and any scc-node in  $S_2$  is a descendant of an scc-node in  $S_1$  because otherwise the former scc-node should belong to  $S_1$  instead of  $S_2$ . Moreover,  $S_2$  contains all the descendants of scc-nodes in  $S_1$ . Indeed given a child  $s_2$  of an scc-node  $s_1 \in S_1$ , the fragmentable link from  $s_1$  to  $s_2$  connects an instance  $n_1$  grouped in  $s_1$  to an

instance  $n_2$  grouped in  $s_2$ . Because  $s_1 \in S_1$ , following the demonstrated item 1 above, we have  $n_1 \in N'$ . Because of the outgoing fragmentable link from  $n_1$  to  $n_2$  and since  $M'$  satisfies Condition 1, we also have  $n_2 \in N'$ . Therefore we have  $s_2 \in S$ . Furthermore,  $s_2 \notin S_1$  because it has  $s_1 \in S$  as its ancestor. Hence we have  $s_2 \in S_2$ . Inductively we conclude that any descendant of  $s_1$  belongs to  $S_2$  and hence  $S$  constitutes an antichain plus descendant.

3. Collapse all the scc-nodes in  $S$  into an antichain-node called  $\alpha$ . We demonstrate that  $M'$  is induced by the instances grouped in  $\alpha$ , i.e.,  $N'$  equals to the set of instances grouped in  $\alpha$ : any instance of  $M'$  is grouped in  $\alpha$  because of the selection criteria of  $S$ , and any instance grouped in  $\alpha$  is an instance of  $M'$  following the demonstrated item 1 above.

□

**Lemma 4** *Given a model  $EM$ , an evaluation environment  $\Lambda$ , and an expression of CoreOCL  $e$ , if  $\Lambda \stackrel{EM}{\models} e \Downarrow (v, O_S)$  holds, then for any object induced sub-model  $EM'$  of  $EM$  where  $O_S \subseteq O'$ ,  $\Lambda \stackrel{EM'}{\models} e \Downarrow (v, O_S)$  also holds.*

*Proof* We prove by induction on the depth of the structure of the expression  $e$ .

*Base cases: expressions of depth 1*

$(e \equiv dv)$  Trivial.

$(e \equiv x)$  Because of the same evaluation environment  $\Lambda$  and  $O_S \subseteq O'$ .

$(e \equiv \text{Bag}\{\} \mid \text{Set}\{\} \mid \text{Seq}\{\} \mid \text{OSet}\{\})$  Trivial.

$(e \equiv \text{allInstances}(\mathbb{C}))$  Because of  $\{o \mid o \in O \text{ and } \tau_o(o) \stackrel{\text{EM}}{=} \mathbb{C}\} \subseteq O'$ , and the fact that  $\text{EM}'$  is an object induced sub-model of  $\text{EM}$ , we have  $\{o \mid o \in O' \text{ and } \tau_o'(o) \stackrel{\text{EM}'}{=} \mathbb{C}\} = \{o \mid o \in O \text{ and } \tau_o(o) \stackrel{\text{EM}}{=} \mathbb{C}\}$ .

*Induction cases: expressions of depth  $n$*  We suppose that for any expression of depth  $< n$ , the lemma holds. We prove that it is also the case for expressions of depth  $n$ . We distinguish the top most structure the expression may have:

$(e \equiv e_1.\text{sf})$  Assume  $\Lambda \stackrel{\text{EM}}{\models} e_1.\text{sf} \Downarrow (v, O_S)$  holds. Following `STRUCTURALFEATURECALL`, we have both  $\Lambda \stackrel{\text{EM}}{\models} e_1 \Downarrow (o, O_S^1)$  and  $o.\text{sf} \stackrel{\text{EM}}{=} v$  hold, where  $O_S = O_S^1 \cup O_S(v)$ . By induction hypothesis, we have also  $\Lambda \stackrel{\text{EM}'}{\models} e_1 \Downarrow (o, O_S^1)$  holds. According to Lemma 3 and the lemma assumption,  $o \in O_S^1 \subseteq O_S \subseteq O'$  and  $O_S(v) \subseteq O_S \subseteq O'$ . Because  $\text{EM}'$  is an object induced sub-model of  $\text{EM}$ ,  $o.\text{sf}$  computes to the same value  $v$  in  $\text{EM}'$  as in  $\text{EM}$  according to Table 1. Namely, we also have  $o.\text{sf} \stackrel{\text{EM}'}{=} v$ . As a consequence, by applying rule `STRUCTURALFEATURECALL`, we have  $\Lambda \stackrel{\text{EM}'}{\models} e_1.\text{sf} \Downarrow (v, O_S^1 \cup O_S(v))$  holds, i.e.,  $\Lambda \stackrel{\text{EM}'}{\models} e \Downarrow (v, O_S)$  holds.

$(e \equiv \text{Tuple}\{\overline{\text{t}} : \overline{\text{t}} = \overline{v}\})$  Following `TUPLE` and by induction on the sub-expressions  $\overline{e}$ .

$(e \equiv e_1.\text{tn})$  Following `TUPLEPARTCALL` and by induction on the sub-expression  $e_1$ .

$(e \equiv \text{let } x = e_1 \text{ in } e_2)$  Following `LETBINDING` and by induction on the two sub-expressions  $e_1$  and  $e_2$ .

$(e \equiv \text{def } c :: op = \lambda(\overline{x} : \overline{\text{t}}).e_1 \text{ in } e_2)$  Following `DEFEXP` and by induction on the sub-expression  $e_2$ .

$(e \equiv e_1.op(\overline{x} := \overline{e}))$  Following `OBJECTOPERATIONCALL` and by induction on the sub-expressions:  $e_1$ ,  $\overline{e}$ , and body expression of the operation  $e_{op}$ .

$(e \equiv \kappa(\overline{e}))$  Following `KAPPCALL` and by induction on the sub-expressions  $\overline{e}$ .

$(e \equiv \text{if } e_1 \text{ then } e_2 \text{ else } e_3)$  Following `IFTRUE` (respectively `IFFALSE`) and by induction on the sub-expressions  $e_1$  and  $e_2$  (respectively  $e_3$ ).

$(e \equiv e_1 \text{ asInstanceOf } c)$  Following `DOWNCASTOK` (respectively `DOWNCASTNOTOK`), by induction on the sub-expression  $e_1$ , and because  $\text{EM}'$  is an object induced sub-model of  $\text{EM}$ , i.e., for an object  $o$  in both  $\text{EM}$  and  $\text{EM}'$ ,  $\tau_o(o) \preceq c$  iff  $\tau_o'(o) \preceq [\text{EM}']c$  (respectively  $\tau_o(o) \not\preceq c$  iff  $\tau_o'(o) \not\preceq [\text{EM}']c$ ).

$(e \equiv e_1 \text{ isInstanceOf } c)$  Following `ISINSTANCEOFTTRUE` (respectively `ISINSTANCEOFFALSE`), by induction on the sub-expression  $e_1$ , and because  $\text{EM}'$  is an object induced sub-model of  $\text{EM}$ , i.e., for an object  $o$  in both  $\text{EM}$  and  $\text{EM}'$ ,  $\tau_o(o) \stackrel{\text{EM}}{=} c$  iff  $\tau_o'(o) \stackrel{\text{EM}'}{=} c$ .

$(e \equiv e_1 \text{ isOfKind } c)$  Following `ISKINDOFTTRUE` (respectively `ISKINDOFFALSE`), by induction on the sub-expression  $e_1$ , and because  $\text{EM}'$  is an object induced sub-model of  $\text{EM}$ , i.e., for an object  $o$  in both  $\text{EM}$  and  $\text{EM}'$ ,  $\tau_o(o) \stackrel{\text{EM}}{=} c'$  iff  $\tau_o'(o) \stackrel{\text{EM}'}{=} c'$ ; and  $\text{EM}$  and  $\text{EM}'$  have the same metamodel where  $c' \preceq c$  (respectively  $c' \not\preceq c$ ).

$(e \equiv e_1 \rightarrow \text{iterate}(x; y = e_2 \mid e_3))$  Following evaluation rule

COLLECTIONITERATION and by induction on the sub-expressions  $e_1$ ,  $e_2$  and  $e_3$ .

□

**Theorem 6** *Without calling allInstances, CoreOCL invariants are all forward following Definition 8.*

*Proof* Let  $e_{\text{inv}} = (c, e)$  be a CoreOCL invariant defined for the metamodel of a model EM where  $e$  is an expression written in the sub-language of CoreOCL excluding allInstances. Let  $o$  be an object of EM where  $\tau_o(o) \preceq c$ . Let  $\Lambda$  be the initial evaluation environment for  $\text{eval}(e_{\text{inv}}, \text{EM}, o)$ . Suppose  $\Lambda \stackrel{\text{EM}}{\models} e \Downarrow (v, O_S)$ . The set of objects in the scope of  $(e_{\text{inv}}, \text{EM}, o)$ , i.e.,  $O_S$ , is reachable from  $o$  following the first statement of Lemma 7. □

**Lemma 7** *Given a model EM, an object o in it, a CoreOCL expression e with no calls to allInstances (i.e., e is an expression written in the sub-language of CoreOCL excluding allInstances), and an evaluation environment  $\Lambda$  in which values bound to variables only refer to objects reachable from o if any, i.e.,  $\forall x$  in the domain of  $\Lambda$ ,  $\forall o' \in O_S(\Lambda(x))$ ,  $o'$  is reachable from o (see Definition 7). Suppose  $\Lambda \stackrel{\text{EM}}{\models} e \Downarrow (v, O_S)$ . We have the following two statements hold:*

1.  $\forall o' \in O_S$ ,  $o'$  is reachable from o;
2.  $\forall o' \in O_S(v)$ ,  $o'$  is reachable from o.

*Proof* We prove by induction on the depth of the structure of the expression  $e$ .

*Base cases: expressions of depth 1*

$(e \equiv dv)$  Trivial.

$(e \equiv x)$  Because values bound in the evaluation environment  $\Lambda$  only refer to objects reachable from  $o$ .

$(e \equiv \text{Bag}\{\} \mid \text{Set}\{\} \mid \text{Seq}\{\} \mid \text{OSet}\{\})$  Trivial.

*Induction cases: expressions of depth n* We suppose that for any expression of depth  $< n$ , the two statements hold. We prove that it is also the case for expressions of depth  $n$ . We distinguish the top most structure the expression may have:

$(e \equiv e_1.\text{sf})$  Assume  $\Lambda \stackrel{\text{EM}}{\models} e_1.\text{sf} \Downarrow (v, O_S)$  holds. Following STRUCTURALFEATURECALL, we have  $\Lambda \stackrel{\text{EM}}{\models} e_1 \Downarrow (o_1, O_S^1)$  and  $o_1.\text{sf} \stackrel{\text{EM}}{=} v$  hold, where  $O_S = O_S^1 \cup O_S(v)$ . By induction hypothesis, we have  $o_1$  reachable from  $o$ . Therefore we have also objects in  $O_S(v)$  all reachable from  $o$  via  $o_1$  and  $\text{sf}$  following Table 1 and Definition 7. As a consequence, objects in  $O_S$  are all reachable from  $o$  because it is the case for all the objects in  $O_S^1$  (by induction hypothesis) and  $O_S(v)$ .

$(e \equiv \text{Tuple}\{\overline{\text{tn}} : \overline{\text{te}} = \overline{\text{e}}\})$  Assume  $\Lambda \stackrel{\text{EM}}{\models} \text{Tuple}\{\overline{\text{tn}} : \overline{\text{te}} = \overline{\text{e}}\} \Downarrow (v, O_S)$  holds. Following TUPLE, we have  $\Lambda \stackrel{\text{EM}}{\models} e_i \Downarrow (v_i, O_S^i)$  holds for all  $e_i \in \overline{\text{e}}$  where  $1 \leq i \leq k$ ,  $v = \text{Tuple}\{\text{tn}_1 : \text{te}_1 = v_1, \dots, \text{tn}_k : \text{te}_k = v_k\}$ , and  $O_S = \bigcup_{1 \leq i \leq k} O_S^i$ . By induction hypothesis, objects in  $O_S^i$  and  $O_S(v_i)$  are all reachable from  $o$  for all  $1 \leq i \leq k$ . Therefore, objects in  $O_S$  and objects in  $O_S(v) = O_S(\text{Tuple}\{\text{tn}_1 : \text{te}_1 = v_1, \dots, \text{tn}_k : \text{te}_k = v_k\}) = \bigcup_{1 \leq i \leq k} O_S(v_i)$  are all be reachable from  $o$ .

$(e \equiv e_1.\text{!n})$  Assume  $\Lambda \stackrel{\text{EM}}{\models} e_1.\text{!n} \Downarrow (v, O_S)$  holds. Following `TUPLEPARTCALL`, we have  $\Lambda \stackrel{\text{EM}}{\models} e_1 \Downarrow (\text{Tuple}\{\overline{\text{!n}} : \overline{\text{!}} = \overline{v}\}, O_S)$  holds. By induction, objects in  $O_S$  and  $O_S(\text{Tuple}\{\overline{\text{!n}} : \overline{\text{!}} = \overline{v}\}) = \bigcup_{v_i \in \overline{v}} O_S(v_i)$  are all reachable from  $o$ . Therefore, objects in  $O_S(v_i)$  are all reachable from  $o$ .

$(e \equiv \text{let } x = e_1 \text{ in } e_2)$  Following `LETBINDING`, and by induction on the two sub-expressions  $e_1$  and  $e_2$ .

$(e \equiv \text{def } c :: \text{op} = \lambda(\overline{x} : \overline{\text{!}}).e_1 \text{ in } e_2)$  Following `DEFEXP` and by induction on the sub-expression  $e_2$ .

$(e \equiv e_1.\text{op}(\overline{x} := \overline{e}))$  Assume  $\Lambda \stackrel{\text{EM}}{\models} e_1.\text{op}(\overline{x} := \overline{e}) \Downarrow (v, O_S)$ . Following `OBJECTOPERATIONCALL`, we have  $\Lambda \stackrel{\text{EM}}{\models} e_1 \Downarrow (o_1, O_S^{o_1})$ ,  $\Lambda \stackrel{\text{EM}}{\models} \overline{e} \Downarrow (\overline{v}, \overline{O_S})$ , and  $\Lambda \oplus \{(\text{self} : o_1)\} \oplus \{(\overline{x} : \overline{v})\} \stackrel{\text{EM}}{\models} e_{\text{op}} \Downarrow (v, O_S^{\text{op}})$ , where  $O_S = O_S^{o_1} \cup (\bigcup \overline{O_S}) \cup O_S^{\text{op}}$ . By induction hypothesis, we have  $o_1$ , objects in  $O_S^{o_1}$ , objects in  $O_S(\overline{v})$ , and objects in  $\overline{O_S}$  all reachable from  $o$ . Therefore, the overridden environment  $\Lambda \oplus \{(\text{self} : o_1)\} \oplus \{(\overline{x} : \overline{v})\}$  binds variables to values that only refer to objects reachable from  $o$  if any. Hence by induction hypothesis, objects in  $O_S(v)$  and  $O_S^{\text{op}}$  are all reachable from  $o$ . As a consequence, objects in  $O_S$  are all reachable from  $o$ .

$(e \equiv \kappa(e_1, \dots, e_k))$  Assume  $\Lambda \stackrel{\text{EM}}{\models} \kappa(e_1, \dots, e_k) \Downarrow (v, O_S)$  holds. Following `KAPPACALL`, we have  $\Lambda \stackrel{\text{EM}}{\models} e_i \Downarrow (v_i, O_S^i)$  holds for  $1 \leq i \leq k$ ,  $\kappa(v_1, \dots, v_k) = v$ , and  $O_S = \bigcup_{1 \leq i \leq k} O_S^i$ . By induction hypothesis, objects in  $O_S^i$  and  $O_S(v_i)$  are all reachable from  $o$  for

$1 \leq i \leq k$ . Therefore, objects in  $O_S$  are all reachable from  $o$ . Moreover,  $\kappa$  being the pre-defined operations, does not introduce any new objects other than those are referenced in its operands, i.e.,  $v_i$  for  $1 \leq i \leq k$ . Hence objects in  $O_S(v)$  should all be reachable from  $o$  as well.

$(e \equiv \text{if } e_1 \text{ then } e_2 \text{ else } e_3)$  Following `IFTRUE` (respectively `IFFALSE`) and by induction on the sub-expressions  $e_1$  and  $e_2$  (respectively  $e_3$ ).

$(e \equiv e_1 \text{ asInstanceOf } c)$  Following `DOWNCASTOK` (respectively `DOWNCASTNOTOK`) and by induction on the sub-expression  $e_1$ .

$(e \equiv e_1 \text{ isInstanceOf } c)$  Following `ISINSTANCEOFTTRUE` (respectively `ISINSTANCEOFFALSE`) and by induction on the sub-expression  $e_1$ .

$(e \equiv e_1 \text{ isOfKind } c)$  Following `ISKINDOFTTRUE` (respectively `ISKINDOFFALSE`) and by induction on the sub-expression  $e_1$ .

$(e \equiv e_1 \rightarrow \text{iterate}(x; y = e_2 \mid e_3))$  Following evaluation rule `COLLECTIONITERATION` and by induction on the sub-expressions  $e_1$ ,  $e_2$  and  $e_3$ .

□

**Lemma 8** *Given an EMF metamodel  $\mathbb{EM} = (\mathbb{C}, \mathbb{DT}, \text{DV}, \mathbb{Rf}, \mathbb{At})$ , an EMF model  $\text{EM} = (\text{EM}, O, \tau_o, \text{RA}, \tau_{ra}, \text{ra}_t, \text{ra}_s, \text{AA}, \tau_{aa}, \text{aa}_s, \text{aa}_t)$ , and an object induced sub-model of  $\text{EM}$ , written  $\text{EM}' = (\text{EM}, O', \tau_o', \text{RA}', \tau_{ra}', \text{ra}_t', \text{ra}_s', \text{AA}', \tau_{aa}', \text{aa}_s', \text{aa}_t')$ , suppose that:*

1.  $\text{EM}$  conforms to  $\mathbb{EM}$ ;

2.  $EM'$  satisfies Condition 4;
3.  $EM$  satisfies Condition 5;

then  $EM'$  also conforms to  $EM$ .

*Proof* To check the conformance of  $EM'$  to  $EM$ , it amounts to checking against the ten conformance conditions in EMF summarized in Section 4.3.

1. Reference typing condition holds in  $EM'$  because  $EM'$  is an object induced sub-model of  $EM$ , i.e., a reference assignment in  $EM'$  is also in  $EM$ , and  $EM$  is conformant to  $EM$ , i.e., all reference assignments in  $EM$  are type valid.
2. Attribute typing condition holds in  $EM'$  because  $EM'$  is an object induced sub-model of  $EM$ , i.e., an attribute assignment in  $EM'$  is also in  $EM$ , and  $EM$  is conformant to  $EM$ , i.e., all attribute assignments in  $EM$  are type valid.
3. Reference multiplicity condition holds in  $EM'$  because  $EM'$  is an object induced sub-model of  $EM$  and  $EM'$  satisfies Condition 4, i.e., for any object in  $EM'$ , the set of reference assignments whose source is the object is exactly the same as in  $EM$ , and  $EM$  is conformant to  $EM$ , i.e., the number of reference assignments is in the correct range.
4. Attribute multiplicity condition holds in  $EM'$  because  $EM'$  is an object induced sub-model of  $EM$ , i.e., any object in  $EM'$  has exactly the same set of attribute assignments as in  $EM$ , and  $EM$  is conformant to  $EM$ , i.e., the number of attribute assignments is in the correct range.
5. Opposite reference condition: given two references of  $EM$  called  $rf_1$  and  $rf_2$  where  $rf_1.eOpposite = rf_2$ , two objects  $o_1, o_2 \in O'$ , and a reference assignment  $ra_1 \in RA'$  where  $\tau_{ra}'(ra_1) = rf_1$ ,  $ra_s'(ra_1) = o_1$ , and  $ra_t'(ra_1) = o_2$ , because  $EM'$  is an object induced sub-model of  $EM$ ,  $o_1, o_2 \in O$  and  $ra_1 \in RA$ . Because  $EM$  is conformant to  $EM$ , there exists a reference assignment  $ra_2 \in RA$  such that  $\tau_{ra}(ra_2) = rf_2$ ,  $ra_s(ra_2) = o_1 = ra_t(ra_1)$ , and  $ra_t(ra_2) = o_2 = ra_s(ra_1)$ . Again because  $EM'$  is an object induced sub-model of  $EM$ , we have  $ra_2 \in RA'$  where  $\tau_{ra}'(ra_2) = rf_2$ ,  $ra_s'(ra_2) = o_1$ , and  $ra_t'(ra_2) = o_2$ .
6. Containment condition: on one hand, because  $EM'$  is an object induced sub-model of  $EM$ , an object  $o_2$  is contained in another object  $o_1$  in  $EM'$  if and only if it is also the case in  $EM$ ; on the other hand,  $EM$  is conformant to  $EM$ , i.e., no object is contained in more than one object, neither is it contained in itself. Therefore, it is also the case in  $EM'$ .
7. Abstract class condition holds in  $EM'$  because  $EM'$  is an object induced sub-model of  $EM$ , i.e., all the objects in  $EM'$  are also in  $EM$ , and  $EM$  is conformant to  $EM$ , i.e., none of the objects in  $EM$  are instantiated from an abstract class.
8. Uniqueness condition holds in  $EM'$  because  $EM'$  is an object induced sub-model of  $EM$ , i.e., any object in  $EM'$  has the subset of reference and attribute assign-

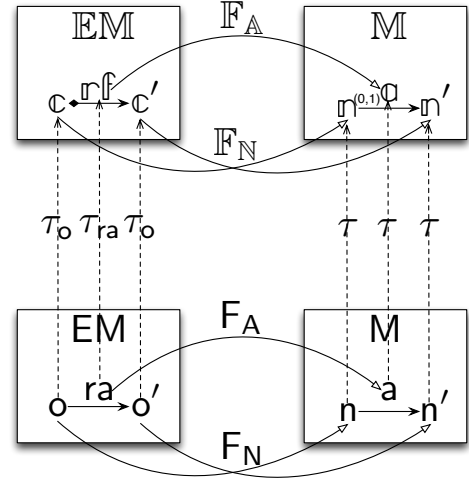
ments as in  $\mathbb{EM}$ , and  $\mathbb{EM}$  is conformant to  $\mathbb{EM}$ , i.e., reference and attribute assignments are never duplicated for unique references and attributes. A subset of a non-multiset is also a non-multiset.

9. Invariant condition: because  $\mathbb{EM}$  satisfies Condition 5, all invariants to be checked are forward. Because  $\mathbb{EM}'$  satisfies Condition 4, for any forward invariant  $\mathbf{einv}$  of  $\mathbb{EM}$ , an object  $\mathbf{o}$  in  $\mathbb{EM}$  whose type is compatible with the context of  $\mathbf{einv}$ , the scope of  $(\mathbf{einv}, \mathbb{EM}, \mathbf{o})$  is included in  $\mathbb{EM}'$  as long as  $\mathbf{o}$  is. Following Theorem 5, we thus have  $\text{eval}(\mathbf{einv}, \mathbb{EM}, \mathbf{o}) = \text{eval}(\mathbf{einv}, \mathbb{EM}', \mathbf{o})$ . We already know that  $\mathbf{einv}$  holds in  $\mathbb{EM}$  because it conforms to  $\mathbb{EM}$ . As a consequence,  $\mathbf{einv}$  also holds in  $\mathbb{EM}'$ .

□

**Lemma 9** *Given an EMF model  $\mathbb{EM} = (\mathbb{EM}, \mathbf{O}, \tau_{\mathbf{o}}, \mathbf{RA}, \tau_{\mathbf{ra}}, \mathbf{ra}_{\mathbf{t}}, \mathbf{ra}_{\mathbf{s}}, \mathbf{AA}, \tau_{\mathbf{aa}}, \mathbf{aa}_{\mathbf{s}}, \mathbf{aa}_{\mathbf{t}})$  conforming to an EMF metamodel  $\mathbb{EM}$ , and any object induced sub-model of  $\mathbb{EM}$  derived by the EMF model decomposition algorithm, written  $\mathbb{EM}' = (\mathbb{EM}, \mathbf{O}', \tau_{\mathbf{o}}', \mathbf{RA}', \tau_{\mathbf{ra}}', \mathbf{ra}_{\mathbf{t}}', \mathbf{ra}_{\mathbf{s}}', \mathbf{AA}', \tau_{\mathbf{aa}}', \mathbf{aa}_{\mathbf{s}}', \mathbf{aa}_{\mathbf{t}}')$ ,  $\mathbb{EM}'$  satisfies Condition 4.*

*Proof* Let  $\mathbf{M} = (\mathbb{M}, \mathbf{N}, \mathbf{A}, \tau, \mathbf{src}, \mathbf{tgt})$  be the corresponding abstract model of  $\mathbb{EM}$  (i.e.,  $\mathbf{M} = \mathbf{F}(\mathbb{EM})$ ) and  $\mathbf{M}' = (\mathbb{M}, \mathbf{N}', \mathbf{A}', \tau', \mathbf{src}', \mathbf{tgt}')$  be the sub-model of  $\mathbf{M}$ , from which  $\mathbb{EM}'$  is constructed. Assume the setting as illustrated in Figure 10 holds. We demonstrate that  $\mathbb{EM}'$  satisfies Condition 4 in three steps:



**Fig. 10** Setting assumed for the proof of Lemma 9

1. according to step 3 of the EMF model decomposition algorithm,  $\mathbf{o} \in \mathbf{O}' \implies \mathbf{n} \in \mathbf{N}'$ ;
2. according to Theorem 2, model  $\mathbf{M}'$  satisfies Condition 1, i.e.,  $\mathbf{n} \in \mathbf{N}' \implies \mathbf{n}' \in \mathbf{N}'$ ;
3. according to step 3 of the EMF model decomposition algorithm,  $\mathbf{n}' \in \mathbf{N}' \implies \mathbf{o}' \in \mathbf{O}'$ .

□

## A.2 Statistical Table of the Usage of allInstances

Metamodel	Source	# Invariants	# Invariants with <code>allInstances</code> originally	# Invariants with <code>allInstances</code> after normalization	
				monotonic	non-monotonic
UML	OMG	172	0	0	0
MOF	OMG	85	0	0	0
OCL	OMG	46	13	13	0
CORBA	OMG	40	1	1	0
CWM	OMG	90	0	0	0
DiagramDefinition	OMG	16	0	0	0
BLanguage	Academia	9	0	0	0
SAD3	Academia	9	0	0	0
CPFSTool	Academia	27	0	0	0
DeclarativeWorkflow	Academia	23	0	0	0
ER2RE	Academia	59	11	11	0
RBAC	Academia	33	6	0	1
DBLP	Academia	26	7	7	0
SAM	Industry	73	11	7	4
Total	N/A	707	49	39	5

**Table 2** Statistics on the Usage of `allInstances` in Metamodel Specifications