

Gangam: A Transformation Modeling Framework*

Kajal T. Claypool
Department of Computer Science
University of Massachusetts, Lowell
kajal@cs.uml.edu

Elke A. Rundensteiner
Department of Computer Science
Worcester Polytechnic Institute
rundenst@cs.wpi.edu

Abstract

Integration of multiple heterogeneous data sources continues to be a critical problem for many application domains and a challenge for researchers world-wide. One aspect of integration is the translation of schema and data across data model boundaries. Researchers in the past have looked at both customized algorithmic approaches as well as generic meta-modeling approaches as viable solutions. We now take the meta-modeling approach the next step forward. In this paper, we propose a flexible, extensible and re-usable transformation modeling framework which allows users to (1) model their transformations; (2) to choose from a set of possible execution strategies to translate the underlying schema and data; and (3) to access and re-use a library of transformation generators. In this paper, we present the core of our modeling framework - a set of cross algebra operators that covers the class of linear transformations, and two different techniques of composing these operators into larger transformation expressions. We also present an evaluation strategy to execute the modeled transformation, and thereby transform the input schema and data into the target schema and data assuming that data model wrappers are provided for each data model. The proposed framework has been implemented, and we give an overview of this prototype system.

Keywords: Cross Model Mapping Algebra, Heterogeneous System Integration, Schema Transformation

1. Introduction

Integration of multiple heterogeneous data sources continues to be a critical problem for many application domains and a challenge for researchers world-wide [5]. Each database brings with it its own concepts, semantics, data formats, and access methods. Currently, the burden falls on the human to manually resolve conflicts, integrate the data,

*This work was supported in part by the NSF NYI grant #IRI 97-96264.

and interpret the results. More often than not, this barrier proves too difficult or too time-consuming to overcome and data hence often is under-exploited.

Data integration as a research field looks at automating as many of the tasks related to the above process, and hence aims to provide better and painless access to data no matter what data source or format it is stored in. One aspect of data integration is schema matching. Schema matching is the task of finding semantic correspondences between elements of two schemas [16]. Many researchers have addressed the schema matching problem either for a specific domain [6, 3, 4] or in a generic domain-independent way [14, 16, 5, 2, 12, 17, 20]. Another aspect of the integration problem with respect to the heterogeneity of information, i.e., the different data models, is the *translation* of schema (and data) from one data model to another. Solutions for this include customized algorithmic approaches [26, 11, 22, 7, 23] and meta-modeling approaches [2, 12, 20, 5, 17]. A customized algorithmic approach provides fixed translation algorithms that convert schema and data between a given pair of data models. The meta-modeling approach provides a more general technique that goes beyond translations for a given pair of data models. The translations themselves are generally expressed either via rules [17, 2] or via fragments of code [5, 12].

In our work we focus on the meta-modeling approach for the translation of schemas across data model boundaries. In particular, we focus on the explicit *modeling* and the subsequent execution of the transformations themselves, an aspect not addressed by previous research [2, 12, 20, 5, 17]. While models, such as the UML model, for static concepts like schemata or data models have been much studied, models for the more dynamic aspects such as for transformations have been largely overlooked. The goal of our work thus is to provide a *flexible, extensible* and *re-usable* translation modeling framework wherein users can (1) explicitly model the translations between schemas; (2) choose automated execution strategies to execute the modeled translations that would transform the source schema and data to the desired target schema and data; and (3) choose and compose trans-

formations from an existing library of translations. Such a framework offers many advantages over previous translation approaches. In particular it allows for (1) optimized execution strategies as each modeled transformation can be reasoned over to determine the optimal execution plan similar to the algebraic query optimization; (2) development of a *generic* tool set for facilitating activities such as maintenance; and (3) query merging and translation in a multi-tier environment.

To enable this translation modeling framework, we identify (1) the fundamental operations required to express and model the translation process (Section 3); and (2) the flexible techniques necessary for composing these core operations into larger meaningful translations (Section 4). A modeled transformation can be executed using any one of the many possible execution strategies, to perform the requisite data transformation process. These execution strategies range from the mapping of the cross algebra expressions to full-fledged query languages such as SQL [1] and XQuery [10] to having customized execution algorithms. In this paper, we briefly sketch out a customized algorithm (Section 5) for executing the cross algebra expressions to illustrate the simplicity of this task.

2. Background: Sangam Graph Model

```

<!ELEMENT item (location,
                mailbox, name)>
<!ATTLIST item id ID #REQUIRED
              featured CDATA #IMPLIED>
<!ELEMENT location (#PCDATA)>
<!ELEMENT mailbox (mail*)>
<!ATTLIST mailbox id CDATA>
<!ELEMENT mail (from, to, date)>
<!ATTLIST mail text CDATA>
<!ELEMENT from (#PCDATA)>
<!ELEMENT to (#PCDATA)>
<!ELEMENT date (#PCDATA)>
<!ELEMENT name (firstName,
                lastName)>
<!ELEMENT firstName (#PCDATA)>
<!ELEMENT lastName (#PCDATA)>

```

Figure 1. A Fragment of the XMark Benchmark DTD.

We assume, as in previous modeling approaches [2, 12, 20, 5, 17], that schemas from different data models are first represented in one common data model. In our work we assume that all schemas are represented by a simple graph called a *Sangam graph*, an instance of the Sangam graph

model [8]. A Sangam graph $G = (N, E, \lambda)$ is a directed graph of nodes N and edges E , and a set of labels λ . Each node has an associated type *complex* (\square) or *atomic* (\circ); and each edge is either a *containment* (\rightarrow) or a *property* ($--\rightarrow$) edge. A *containment* edge is an edge between two complex nodes, while a *property* edge exists between a complex node and an atomic node. Each node n has associated with it a set of objects, its *extent* denoted as $I(n)$. Each object $o \in I(n)$ is a pair $\langle id, v \rangle$ where id is a globally unique identifier and v is its data value. Each edge $e: \langle n_1, n_2 \rangle$ also has associated with it a set of objects, termed its *extent* $R(e)$. Each object $o_e \in R$ is a triple $\langle id, o_1, o_2 \rangle$ where id is a system-generated identifier, object $o_1 \in I(n_1)$ and object $o_2 \in I(n_2)$. There may be zero to multiple edges between the same two nodes. In addition, each edge e is annotated with a set of properties ζ , possibly empty. This set of properties includes a *local order*, denoted by ρ , and *quantifier* annotation, denoted by Ω . The local order ρ gives the relative local ordering for all outgoing edges from a given node n in the Sangam graph. A *quantifier* is a pair of integers $[\min:\max]$, with $0 \leq \min \leq \max < \infty$ where \min specifies the minimum and \max the maximum occurrences of objects of a node n_2 for a given object o of node n_1 associated via the relationship (edge) e .

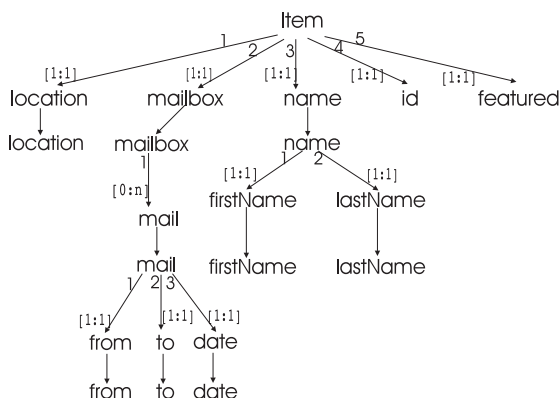


Figure 2. A Fragment of the XMark Benchmark DTD as shown in Figure 1 depicted as a Sangam Graph.

Figure 2 shows the Sangam graph for the XMark benchmark schema of Figure 1. Here each element and attribute is represented by a Sangam graph node. For example, the edge e_1 between the node labeled *item* and *location* represents a relationship between the *item* element and its sub-element *location*. The edge e_1 has an order annotation of 1. The quantifier annotation $[1:1]$ on edge e_1 denotes a functional edge, i.e., that there must be exactly one object of *location* that can participate in a binary relationship with one object of *item*.

3. The Bricks: Cross Algebra Operators

The key factors that influence the achievement of a flexible and extensible translation framework are the building blocks that would enable users to model different translations in order to transform a schema. To enable the modeling of such translations we provide two main building blocks: (1) the bricks: the cross algebra operators which allow the user to express a variety of linear transformations; and (2) the mortar: different techniques that allow users to compose the operators together to represent larger translation units. In this section we present the first building blocks, i.e., the cross algebra operators.

In our work we have identified four basic transformation operators, `cross`, `connect`, `smooth`, and `subdivide`. These operators, termed the *cross algebra operators*, represent the primitive set of operations in the class of linear graph transformations [13] on the basis of which larger more complex linear transformations can be defined. In this section we briefly describe the semantics of these operators. For more details refer to [8].

3.1. Cross Operator

The `cross` algebra operator \otimes takes as input a node n in G and produces as output a node n' in G' . The cross operator is a total mapping, i.e., the objects in the extent of n given by $I(n)$ are mapped one-to-one to the objects in the extent of n' given by $I(n')$ in the output Sangam graph. Figure 3 (a) depicts the cross operator. We use the notation $\otimes_{n'}(n)$ to depict a cross operator with input n and output n' .

3.2. Connect Operator

A `connect` algebra operator (\ominus) corresponds to an edge creation in G' . It takes as input an edge e between two nodes n_1 and n_2 in G and produces an edge e' between two nodes n_1' and n_2' in G' . All objects $o \in I(e)$ are also copied as part of this process. The connect operation succeeds if and only if nodes n_1 and n_2 have already been mapped to the nodes n_1' and n_2' respectively using two cross operators. The connect operator preserves the annotations of the edge e , i.e., the output edge e' will have the same quantifier and local ordering annotation as the input edge e . Figure 3 (b) gives an example of the connect operator. We use the notation $\ominus_{e'}(e)$ to depict a connect operator that maps the edge $e: \langle n_1, n_2 \rangle$ to edge $e': \langle n_1', n_2' \rangle$.

3.3. Smooth Operator

A `smooth` operator (\odot) models the combination of two relationships in G to form one relationship in G' . Let G

be a Sangam graph with three nodes n_1, n_2 , and n_3 , and two relationships represented by edges $e1: \langle n_1, n_2 \rangle$ and $e2: \langle n_2, n_3 \rangle$. The `smooth` (\odot) operator replaces the relationships represented by edges $e1$ and $e2$ in G with a new relationship represented by edge $e': \langle n_1', n_3' \rangle$ in G' . The smooth operator can only be applied when $\otimes(n_1) = n_1'$ and $\otimes(n_3) = n_3'$. The *local order* annotation on the edge e' is set to the *local order* annotation of the edge $e1$. However, as the edge e' has a larger information capacity than the edges $e1$ and $e2$, the *quantifier* annotation of the edge e' is given as: $\rho(e') = \rho(e1) * \rho(e2)$. Figure 4 gives an example of the smooth operator. We use the notation $\odot_{e'}(e1, e2)$ to depict a smooth operator that maps edges $e1: \langle n_1, n_2 \rangle$ and $e2: \langle n_2, n_3 \rangle$ to edge $e': \langle n_1', n_3' \rangle$ in the output.

3.4. Subdivide Operator

A `subdivide` operator \oslash intuitively performs the inverse operation of the `smooth` operator, i.e., it splits a given relationship into two relationships connected via a node. Let G have two nodes n_1 and n_3 and edge $e: \langle n_1, n_3 \rangle$. The `subdivide` operator introduces a new node n_2' in G' such that the edge e in G is replaced by two edges $e1': \langle n_1', n_2' \rangle$ and $e2': \langle n_2', n_3' \rangle$ in G' . The `subdivide` operator is only valid if $\otimes(n_1) = n_1'$ and $\otimes(n_3) = n_3'$. The local order annotation for the edge $e1': \langle n_1', n_2' \rangle$ is the same as the local order annotation of the edge e as $\otimes(n_1) = n_1'$. The edge $e2'$ is the only edge added for the node n_2' and thus has a local order annotation of 1. To preserve the extent $I(e)$, the edges $e1'$ and $e2'$ are assigned quantifier annotations as follows. If $\min(\rho(e)) = 0$, then the quantifier range for $e1'$ is given as $[0 : 1]$, else it is always set to $[1 : 1]$. The quantifier of edge $e2'$ is set equal to the quantifier of edge e . We use the notation $\oslash_{e1', e2', n_2'}(e)$ to depict a subdivide node that maps edge $e: \langle n_1, n_3 \rangle$ to $e1': \langle n_1', n_2' \rangle$ and $e2': \langle n_2', n_3' \rangle$ in the output.

4. The Mortar: Composition Techniques

Cross algebra operators can be composed into larger transformations using two techniques: (1) context dependency; and (2) derivation. The context dependency composition enables several algebra operators to collaborate and jointly operate on sub-graphs to produce one combined output graph. The derivation composition enables the nesting of several algebra operators wherein output of one or more operators becomes the input of another operator. Derivation and context dependency can in turn be combined together to produce larger, more complex transformations. In this section, we present the rules governing their combination.



Figure 3. (a) Example of Cross Algebra Operator; (b) Example of Connect Algebra Operator.

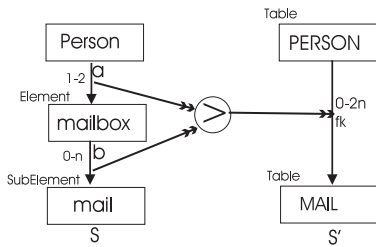


Figure 4. Example of Smooth Operator.

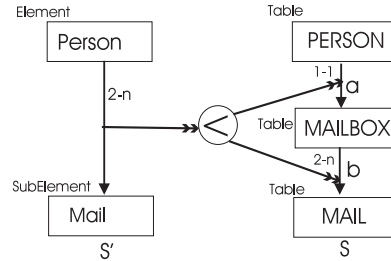


Figure 5. Example of Subdivide Operator.

4.1. Context Dependency Composition

The first composition technique, called the context dependency composition, enables several algebra operators to collaborate and jointly operate on sub-graphs to produce one combined output graph. Figure 7 denotes such a context dependency composition CT of three cross algebra operators. Here, the algebra operators $op1_{A'}(A)$ and $op2_{B'}(B)$ are cross operators that map the nodes A and B in G to nodes A' and B' respectively in the output Sangam graph G'. The algebra operator $op3_{e'}(e)$, a \ominus operator is the root of CT and maps the edge $e: \langle A, B \rangle$ between the nodes A and B in the input Sangam graph G to the edge $e': \langle A', B' \rangle$ between the nodes A' and B' in the output Sangam graph G'. Here the outputs of all operators $op1$, $op2$, and $op3$ together produce G'.

Definition 1 Given an input Sangam graph G, a context dependency expression CT_o is specified as:

$$CT_{o(out_o)}(in_o) = \begin{cases} op_{i(out_i)}(in_i) \\ op_{i(out_i)}(in_i), \\ (CT_{k(out_k)}(in_k)) [\circ \\ (CT_{l(out_l)}(in_l))]^+ \end{cases}$$

where op_i is the parent operator of CT_k and CT_l denoting that op_i must be executed after CT_k and CT_l . op_i uses outputs, inputs and mapping of CT_k and CT_l and $out_o = out_j \cup out_k \cup out_l$. The context dependency expression CT_o operates on nodes n_i and edges $e_i \in G$, and produces

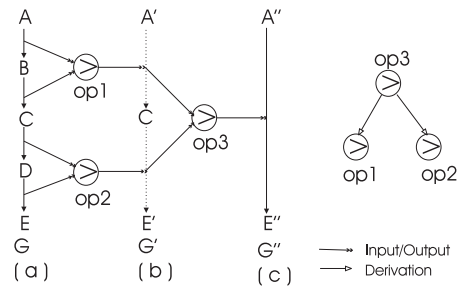


Figure 6. Derivation Composition.

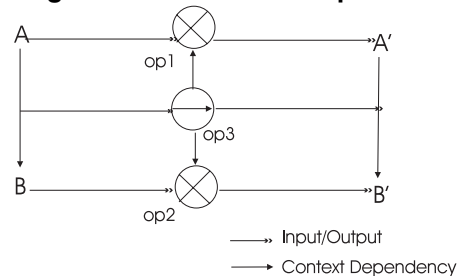


Figure 7. Context Dependency Composition Example.

as output a Sangam graph G' such that all nodes n_i' and/or edges e_i' produced as output by any of the individual operators $op_i \in CT$ are in G'. Here the symbol "+" is part of the BNF grammar syntax to indicate that the expression contained in "[]" may occur one or more times.

As an example, the expression for the context dependency composition in Figure 7 is given as: $CT_{out}(in)$

$= \text{op3}_{e'}(e), (\text{op1}_{A'}(A) \circ \text{op2}_{B'}(B))$. Here the operator op3 is the root of the composition, and op1 and op2 are the children operators. Here $e: \langle A, B \rangle$ and $e': \langle A', B' \rangle$. The final output $\text{out} = e' \cup A' \cup B'$ corresponds to the final output G' .

The context dependency expression CT_o is evaluated from right to left and from inside out. That is all children operators are evaluated prior to the evaluation of their parent operator. The order of evaluation between the sibling operations is immaterial. Beyond the order of evaluation, the context dependency relation between two operators $\text{op3} \rightarrow \text{op1}$ (Figure 7) implies that the operator op3 uses the following three pieces of information in its calculation: (1) the input of op1 ; (2) the output of op1 ; and (3) the mapping ϕ of op1 as established by the type $\otimes, \ominus, \otimes$ and \otimes of op1 .

4.2. Derivation Composition Technique

The second composition technique is the *derivation composition*. This technique enables the nesting of several algebra operators wherein output of one or more operators becomes the input of another operator. Figure 6 gives an example of the modeling of a derivation composition that transforms the path in the Sangam graph G shown in Figure 6 (a) to the edge in the Sangam graph G' given in Figure 6 (c) by applying three smooth nodes \otimes . Let $e_1: \langle A, B \rangle, e_2: \langle B, C \rangle, e_3: \langle C, D \rangle$ and $e_4: \langle D, E \rangle$ be edges in G . Operators $\text{op1}_{e_1'}(e_1, e_2)$ and $\text{op2}_{e_2'}(e_3, e_4)$ are applied to the input edges e_1 and e_2 , and e_3 and e_4 respectively to first produce the intermediate edges $e_1': \langle A', C' \rangle$ and $e_2': \langle C', E' \rangle$ as shown in Figure 6 (b)¹. The operator $\text{op3}_{e_3''}(e_1', e_2')$ operates on these intermediate edges e_1' and e_2' and produces the desired output edge $e_3'': \langle A'', E'' \rangle$ as shown in Figure 6 (c). One approach to achieving this is to first produce the intermediate edges and then the final output edge e_3'' . Equivalently we can express this by nesting. Thus, the output of the algebra expression $\text{op3}_{e_3''}(\text{op1}_{e_1'}(e_1, e_2), \text{op2}_{e_2'}(e_3, e_4))$ is the output edge e_3'' . The output of the operator op3 is said to be *derived* from the outputs of operators op1 and op2 , or put differently, the output of operators op1 and op2 are the *inputs* of operator op3 and are consumed by op3 .

Definition 2 Given an input Sangam graph G , a derivation expression DT_o is given as:

¹We do not include here the context dependency composition that would be needed to map the nodes of the graph.

$$\text{DT}_{o(\text{out}_o)}(\text{in}_o) = \begin{cases} \text{op}_{i(\text{out}_i)}(\text{in}_i) \\ \text{op}_{i(\text{out}_i)}(\text{DT}_{j(\text{out}_j)}(\text{in}_j)) \\ \text{op}_{i(\text{out}_i)}(\text{DT}_{k(\text{out}_k)}(\text{in}_k)), \\ (\text{DT}_{l(\text{out}_l)}(\text{in}_l)) \end{cases}$$

where DT_k and DT_l are derivation trees and $\text{in}_i = \{\text{out}_j\}$ or $\{\text{out}_k, \text{out}_l\}$. The expression DT_o produces as output out_o node and edge elements for an output Sangam graph G' , such that out_o is the output of the root operator op_i and thus also of the complete derivation tree DT_o . Here, “(” and “)” pairs denote nesting and the symbol “,” separates input arguments of an operator op_i .

4.3. Cross Algebra Graphs (CAG): Combining Context Dependency and Derivation Compositions

Derivation and context dependency compositions can be combined in one *cross algebra tree* (CAT) or a *cross algebra graph* (CAG) to model a complex transformation of a Sangam graph G into into a graph G' . Figure 8 represents an example of a CAT. Here let $e_1: \langle A, B \rangle, e_2: \langle B, C \rangle, e_3: \langle C, D \rangle$ and $e_4: \langle D, E \rangle$ represent edges in G , and let $e': \langle A', E' \rangle$ represent an edge in G' . The expression for Figure 8 is $\text{CAT} = \text{CAT3}$.

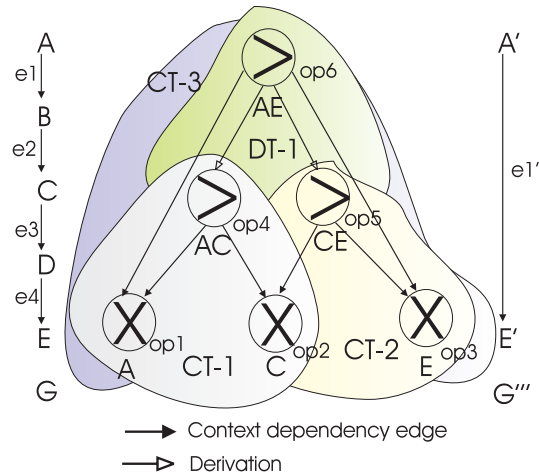


Figure 8. A Cross Algebra Graph.

The output of CAT3 ($\text{CAT3} = \text{DT1}, (\text{op1}_{A'}(A) \circ \text{op3}_{E'}(E))$) is produced by the evaluation of its three inputs, the derivation composition DT1 , and the two cross algebra operators op1 and op3 . The evaluation of the two primitive operators op1 and op3 produces the nodes A' and E'

respectively. The expression DT1 produces two intermediate Sangam graphs G' and G'' (outputs of CT1 and CT2 respectively) that smooth the edges e_1 and e_2 to produce edge e_{Temp1} and smooth edges e_3 and e_4 to produce edge e_{Temp2} respectively. The operator op_6 then gets its inputs from CT1 and CT2 and produces the edge $e' : \langle A', E' \rangle$. The operator op_6 also participates in the composition CAT3, the output of which is the final Sangam graph G''' with nodes A', E' and the edge $e' : \langle A', E' \rangle$. Thus, the CAT in Figure 8 operates on the input Sangam graph G and produces as output the Sangam graph G''' . Formally, we define a CAT as follows.

Definition 3 (CAT) A CAT is an expression that operates on one or more input Sangam graphs G and produces one or more output Sangam graph G' such that:

$$CAT_{o(out_o)}(in_o) = \begin{cases} DT_{i(out_i)}(in_i) \\ CT_{i(out_i)}(in_i) \\ (CAT_{j(out_j)}(in_j))^\dagger, \\ ((CAT_{k(out_k)}(in_k)))^\ddagger \\ op_{j(out_j)} \\ (CAT_{k(out_k)}(in_k))^\ddagger, \\ op_{j(out_j)}((CAT_{k(out_k)}(in_k))), \\ (CAT_{l(out_l)}(in_l))^\gamma. \end{cases}$$

$\dagger =$ A context dependency edge is added from the root op_j of CAT_j to the root op_k of CAT_k .

$\ddagger =$ A derivation edge is added from op_j to op_k , the root of CAT_k .

$\gamma =$ A derivation edge is added from op_j to op_k and op_l , the roots of CAT_k and CAT_l respectively.

Based on the definition of a CAT (Definition 3), we now define a cross algebra graph (CAG). Intuitively, a CAG is a collection of cross algebra trees that may operate on possibly disjoint input graphs to produce possibly disjoint output graphs.

Definition 4 (CAG) A cross algebra graph (CAG) is an expression that operates on one or more input Sangam graphs G and produces one or more output Sangam graph G' such that:

$$CAG_{o(out_o)}(in_o) = \begin{matrix} (CAT_{j(out_j)}(in_j)) \\ (CAT_{k(out_k)}(in_k)) \end{matrix} \begin{matrix} [o \\]^+ \end{matrix}$$

where CAT_j and CAT_k are sibling CATs such that $out_o = out_j \cup out_k$.

5. Execution Strategies for Modeled Transformations

In Sections 3 and 4 we have introduced the bricks and the mortar of our translation modeling framework, and shown how large complex transformations can be modeled in the same. In this section, we briefly describe how the cross algebra expressions representing the modeled transformations can be executed.

function EvaluateCAG (**input:** CAG cag , Sangam graph G ,
output: Sangam graph G')

```
{
  List roots ← cag.getRoots()
  while (roots != null) {
    operator op ← roots.getNext()
    EvaluateCAT (op, G, G')
  }
}
```

function EvaluateCAT (**input:** Operator op , Sangam graph G ,
output: Sangam graph G')

```
{
  if (!op.hasChildren())
    Sangam graph  $G' \leftarrow op.evaluate(G, G')$ 
    op.markDone()
    Sangam graph out ←  $G'$  // cache the local output
    return local $G'$ 
}
```

```
while (op.hasChildren()) {
  operator opC ← op.getNextChild()
  if (e : <op, opC> = derivation)
    SAG  $G_{local} \leftarrow EvaluateCAT (opC, G, G')$ 
    SAG  $G' \leftarrow op.evaluate(G_{local}, G')$ 
    op.markDone()
    SAG out ←  $G'$  // local cached output
    return  $G'$ 
}
```

```
elseif (e : <op, opC> = context dependency)
  SAG  $G_{local} \leftarrow EvaluateCAT (opC, G, G')$ 
  SAG  $G'_{local} \leftarrow op.evaluate(G, G')$ 
  SAG  $G' \leftarrow G_{local} \cup G'_{local}$ 
  op.markDone()
  SAG out ←  $G'_{local}$  // local cached output
  return  $G'$ 
}
```

```
}
```

Figure 9. The Evaluation Algorithm for a Cross Algebra Graph.

The execution of the cross algebra expression here implies the transformation of the source schema and data as per the modeled transformation to produce the target

schema and data. As stated in Section 1 there are many possible strategies for executing these modeled transformations. These strategies range from mapping the modeled transformation into query languages such as SQL [1] or XQuery [10] to applying execution algorithms directly using the transformation model. To keep the discussion simple we now sketch out a customized algorithm for executing the cross algebra expressions as shown in Figure 9.

A cross algebra graph (CAG) can be viewed as a forest of nested context dependency and derivation compositions. Each composition CAT_i of the CAG can be evaluated independently of any other composition CAT_j of the CAG. To evaluate each individual CAT , we use post-order evaluation, i.e., all children operators op_j of an operator op_k are evaluated prior to the evaluation of op_k .

Figure 9 gives the algorithm for evaluating the cross algebra graph. Here, to facilitate evaluation of shared operators, each operator op_j is marked "visited" the first time it is evaluated, and its local output is cached. If the operator op_j is re-visited, no further evaluation of op_j is done, instead its cached output is returned to the invoking parent operator op_i . Evaluation of the tree terminates with the evaluation of the root operator².

6. Architectural Overview of Sangam

The Sangam system incorporating all the techniques discussed in the paper has been developed using Java technology and a variety of tools such as the JAXP [24] for parsing XML documents and the DTD-Parser [25] for parsing the DTDs. Figure 10 gives an architectural overview of the system. Here **SAG-Loader** translates XML and relational schema and data into Sangam graphs; **CAG-Builder** houses the library of transformations and builds transformation models based on the chosen generators for the given input Sangam graph; the **CAG-Evaluator** evaluates the CAG; and lastly the **CAG-Generator** is able to translate the given output Sangam graph into a relational or XML schema and data.

7. Related Work

There is extensive literature under the umbrella of schema transformation and integration [15, 19, 11, 18, 9, 21]. However, this work is typically specific to either an application domain or to a particular data model and does not deal with meta-modeling [2, 5, 20, 21]. Recent work related to ours are Clio [15] a research project at IBM's Almaden Research Center and work by Milo and Zohar [19]. Clio, a tool for creating mappings between two data representations semi-automatically with user inputs, focuses on

²Proof of termination can be found in [8].

supporting querying of data in either the source or the target representation and on just in time cleansing and transformation of data. Milo et al. [19] have looked at the problem of data translations based on schema-matching. They follow an approach similar to Atzeni et al.[2] and Papazoglou et al. [20], but not at the meta-level, in that they define a set of translation logic rules to enable discovery of relationships between two application schemas. Bernstein et al. [5] have also proposed a meta-modeling framework to represent schemas in a common data model to facilitate a host of tools such as the *match* operator which produces a set of matches between two given schemas. While we share this vision, our focus is more on the representation and the execution of transformations that may be produced either by such matches or by a user via a GUI. In fact, based on the discussion with the authors of [5], their meta-model can be extended with a *makeMap* meta-operator to capture the semantics of our work.

We can directly make use of translation algorithms from the literature, such as the algorithms for translating between an XML-DTD and relational schema [11] or mapping rules [19] or output from match operators [5]. That is, we can develop generators that capture such specific algorithms, and then generate separate transformation models that can be executed. Work on equivalence of the translations between models [18] is of particular importance as such properties could also be established for the cross algebra.

8. Conclusions

In this paper, we have presented our *flexible, extensible* and *re-usable* transformation modeling framework. This framework has the advantage of allowing users to model complex translations using the basic building blocks that we have presented in this paper, and then with a click of a button generating and executing code that would perform the data translation. While in this paper, we make substantial contributions towards providing a transformation framework, we believe this is just the tip of the iceberg. Some possible extensions to this work are handling of non-linear transformations and optimization of evaluation strategies.

References

- [1] ANSI Standard. The SQL 92 Standard. <http://www.ansi.org/>, 1992.
- [2] P. Atzeni and R. Torlone. Management of Multiple Models in an Extensible Database Design Tool. In P. M. G. Apers and et al., editors, *Advances in Database Technology - EDBT'96*, Avignon, France, March 25-29, LNCS. Springer, 1996.
- [3] S. Bergamaschi, S. Castano, M. Vincini, and D. Beneventano. Semantic integration of heterogeneous information sources. *Data and Knowledge Engineering*, 36(3):215-249, 2001.

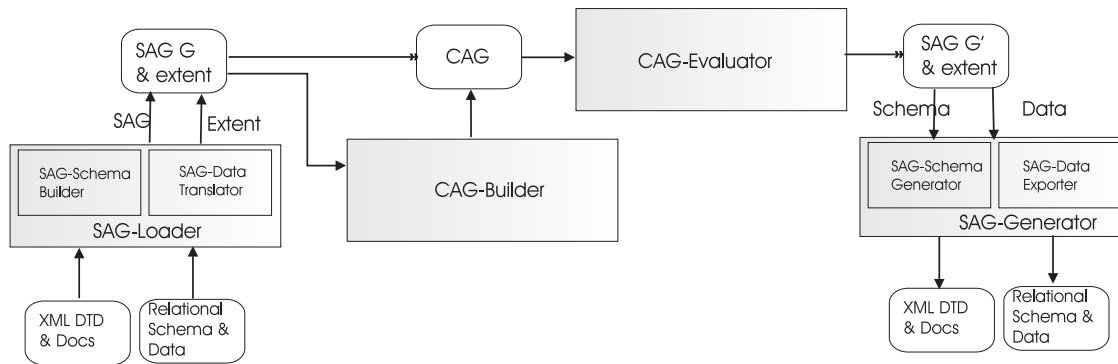


Figure 10. Architecture of \mathcal{G} angan

- [4] J. Berlin and A. Motro. AutoPlex: Automated Discovery of Content for Virtual Databases. In *CoopIS*, pages 108–122, 2001.
- [5] P. A. Bernstein and E. Rahm. Data Warehouse Scenarios for Model Management. In *International Conference on Conceptual Modeling*, 2000.
- [6] M. Bright, A. Hurson, and S. H. Pakzad. Automated Resolution of Semantic Heterogeneity in Multidatabases. *TODS*, 19(2):212–253, 1994.
- [7] B. Catania, E. Ferrari, A. Levy, and A. Meldelzon. XML and Object Technology. In *ECOOOP Workshop on XML and Object Technology, LNCS 1964*, pages 191–202, 2000.
- [8] K. Claypool. *Managing Change in Databases*. PhD thesis, Worcester Polytechnic Institute, May 2002.
- [9] K. Claypool, J. Jin, and E. Rundensteiner. SERF: Schema Evolution through an Extensible, Re-usable and Flexible Framework. In *Int. Conf. on Information and Knowledge Management*, pages 314–321, November 1998.
- [10] P. Fankhauser, M. Fernandez, A. Malhotra, M. a. Rys, J. Simeon, and P. Wadler. XQuery 1.0 Formal Semantics. <http://www.w3c.org/TR/2001/query-semantics>, June 2001.
- [11] D. Florescu and D. Kossmann. Storing and Querying XML Data Using an RDBMS. In *Bulletin of the Technical Committee on Data Engineering*, pages 27–34, Sept. 1999.
- [12] S. Göbel and K. Lutze. Development of Meta Databases for Geospatial Data in the WWW. In *ACM-GIS*, pages 94–99, 1998.
- [13] J. Gross and J. Yellen. *Graph Theory and its Applications*. CRC Press, 1998.
- [14] L. Haas, R. Miller, B. Niswonger, M. Roth, P. Schwarz, and E. Wimmers. Transforming Heterogeneous Data with Database Middleware: Beyond Integration. *IEEE Data Engineering Bulletin*, 22(1):31–36, 1999.
- [15] L. Haas, R. Miller, B. Niswonger, M. Roth, P. Schwarz, and E. Wimmers. Transforming Heterogeneous Data with Database Middleware: Beyond Integration. *IEEE Data Engineering Bulletin*, 22(1):31–36, 1999.
- [16] J. Madhavan, P. Bernstein, and E. Rahm. Generic Schema Matching with Cupid. In *vldb*, pages 49–58, 2001.
- [17] L. Mark and N. Roussopoulos. Integration of Data, Schema and Meta-Schema in the Context of Self-Documenting Data Models. In C. G. Davis, S. Jajodia, P. A. Ng, and R. T. Yeh, editors, *Proceedings of the 3rd Int. Conf. on Entity-Relationship Approach (ER'83)*, pages 585–602. North Holland, 1983.
- [18] R. J. Miller, Y. E. Ioannidis, and R. Ramakrishnan. The Use of Information Capacity in Schema Integration and Translation. In *Int. Conference on Very Large Data Bases*, pages 120–133, 1993.
- [19] T. Milo and S. Zohar. Using Schema Matching to Simplify Heterogeneous Data Translation. In A. Gupta, O. Shmueli, and J. Widom, editors, *VLDB'98, Proceedings of 24th International Conference on Very Large Data Bases, August 24-27, 1998, New York City, New York, USA*, pages 122–133. Morgan Kaufmann, 1998.
- [20] M. Papazoglou and N. Russell. A Semantic Meta-Modeling Approach to Schema Transformation. In *CIKM '95*, pages 113–121. ACM, 1995.
- [21] A. Rosenthal and D. Reiner. Theoretically Sound Transformations for Practical Database Design. In S. T. March, editor, *Entity-Relationship Approach, Proceedings of the Sixth International Conference on Entity-Relationship Approach, New York, USA, November 9-11, 1987*, pages 115–131, 1987.
- [22] J. Shanmugasundaram, G. He, K. Tufte, C. Zhang, D. DeWitt, and J. Naughton. Relational Databases for Querying XML Documents: Limitations and Opportunities. In *Proceedings of 25th International Conference on Very Large Data Bases (VLDB'99)*, pages 302–314, 1999.
- [23] T. Shimura, M. Yoshikawa, and S. Uemura. Storage and Retrieval of XML Documents using Object-Relational Databases. In *Int. Conference and Workshop on Database and Expert Systems Applications*, pages 206–217. Springer-Verlag, 1999.
- [24] J. Systems. The jaxp 1.1.1parser. <http://java.sun.com>, November 2001.
- [25] M. Wutka. The dtd parser. <http://www.wutka.com>, March 2001.
- [26] X. Zhang, W.-C. Lee, G. Mitchell, and E. A. Rundensteiner. Clock: Synchronizing Internal Relational Storage with External XML Documents. In *ICDE-RIDE 2001*, 2001.