

KHEPERA: A System for Rapid Implementation of Domain Specific Languages

Rickard E. Faith

Lars S. Nyland

Jan F. Prins

*Department of Computer Science
University of North Carolina
CB #3175, Sitterson Hall
Chapel Hill NC 27599-3175
{faith,nyland,prins}@cs.unc.edu*

Abstract

The KHEPERA system is a toolkit for the rapid implementation and long-term maintenance of domain specific languages (DSLs). Our viewpoint is that DSLs are most easily implemented via source-to-source translation from the DSL into another language and that this translation should be based on simple parsing, sophisticated tree-based analysis and manipulation, and source generation using pretty-printing techniques. KHEPERA emphasizes the use of familiar, pre-existing tools and provides support for transformation replay and debugging for the DSL processor and end-user programs. In this paper, we present an overview of our approach, including implementation details and a short example.

1 Introduction

Domain specific languages (DSL) can often be implemented as a source-to-source translator composed with a processor for another language. For example, PIC [8], a classic “little language” for typesetting figures, is translated into `troff`, a general-purpose typesetting language. Language composition can be extended in either direction: the CHEM language [1], a DSL used for drawing chemical structures, is translated into PIC, while `troff` is commonly translated into PostScript.

Other DSLs translate into general-purpose high-level programming languages. For example, ControlH, a DSL for the domain of real-time Guidance, Navigation, and Control (GN&C) software, translates into Ada [5]; and RISLA, a DSL for financial engineering, translates into COBOL [18].

The composition of a DSL processor with (for ex-

ample) a C compiler is attractive because it provides portability over a large class of architectures, while achieving performance through the near universal availability of architecture-specific optimizing C compilers.

Yet there are some drawbacks to this approach. While DSLs are often simpler than general purpose programming languages, the domain-specific information available may result in a generated program that can be much larger and substantially different in structure than the original code written in the DSL. This can make debugging very difficult: an exception raised on some line of an incomprehensible C program generated by the DSL processor is a long way removed in abstraction from the DSL input program.

Since the DSL processor is composed with a native high-level compiler, and does not have to perform machine-code generation or optimization, we believe that there are some basic differences between the construction of a compiler for a general purpose programming language and the construction of a translator for a DSL. Our view is that DSL translation is most simply expressed as

1. simple parsing of input into an abstract syntax tree (AST),
2. translation via sophisticated tree-based analysis and manipulation, and
3. output source generation using versatile pretty-printing techniques.

We add the additional caveat that the translation process retain enough information to support the inverse mapping problem, i.e., given a locus in the output source, determine the tree manipulations and input source elements that are responsible for it. This

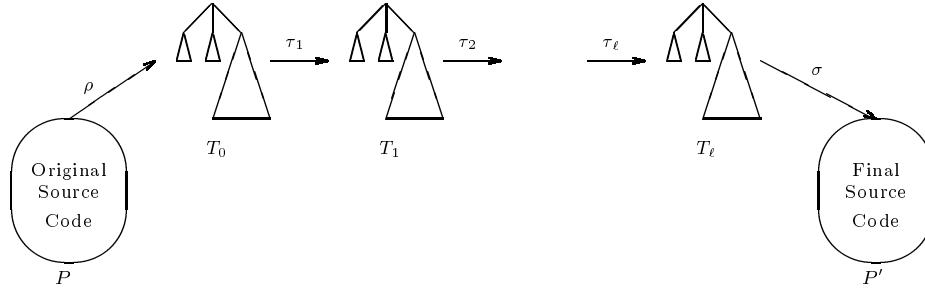


Figure 1: Transformation Process

facility would be useful both for the DSL developer to trace erroneous translation and for the DSL user to trace (run-time) errors back to the input source.

For the translation step we advocate the use of arbitrary AST traversals and transformations. We believe that this approach is simpler for source-to-source translation than the use of attribute grammars, since it decouples the AST analysis and program synthesis from the grammar of the input and output languages. Further, this approach minimizes the need for parsing “heroics”, since simple grammars, close or identical to the natural specification of the DSL syntax, can be used to generate an AST that is specialized in subsequent analysis. By decoupling the input grammar, translation process, and output grammar, this approach is better able to accommodate changes during the evolution of the DSL syntax and semantics.

Throughout this paper, we will use “AST” to refer to abstract syntax tree derived from parsing the input file, and to any intermediate tree-based representations derived from this original AST, even if those representations do not strictly represent an “abstract syntax”.

In our own work we use the DSL paradigm in the compilation of parallel programs. We are particularly interested in the translation into HPF of irregular computations expressed in the PROTEUS [12] language, a DSL providing specialized notation. Our observation was that we were spending a disproportionate amount of effort working on a custom translator implementation to incorporate changes in PROTEUS syntax and improvements in the translation scheme—thus we were motivated to investigate general tool support for DSL translation to simplify this process.

1.1 Goals for a DSL Implementation Toolkit

The implementation of a DSL translator can require considerable overhead, both for the initial implementation and as the DSL evolves. A toolkit should leverage existing, familiar tools as much as possible. Use of such tools takes advantage of previous implementor knowledge and the availability of comprehensive resources explaining these tools (which may not be widely available for a DSL toolkit).

A transformational model for DSL design fits in well with these high-level goals. Consider the problem of translating a program, P , written in the domain specific language, L . In Figure 1, T_0 is an AST which represents P after the parsing phase, ρ . T_ℓ is the final transformed AST, and P' is a valid program in the output language, L' , constructed from T_ℓ during the pretty-printing phase, σ . The transformation process is viewed as a sequential application of various transformations functions, $\tau_{k+1}(T_k) = T_{k+1}$, to the AST. The determination of which transformation function to apply next may require extensive analysis of the AST. Once the transformation functions are determined, however, they can be rapidly applied for replay or debugging.

Within a transformational model, a DSL-building toolkit can simplify the implementation process by providing specialized tools where pre-existing tools are not already available, and to transparently integrate support for debugging within this framework.

The KHEPERA system facilitates both the problem of rapid DSL prototyping and the problem of long-term DSL maintenance through the following specific design goals:

Familiar, modularized parsing components. KHEPERA supports the use of familiar scanning and parsing tools (e.g., the traditional `lex` and `yacc`, or the newer PCCTS [11]) for implementation of a DSL

processor. Because KHEPERA concentrates on providing the “missing pieces” that help with rapid implementation of DSLs, previous knowledge can be utilized, thereby decreasing the slope of the learning curve necessary for the rapid implementation of a DSL.

Familiar, flexible, and efficient semantic analysis. KHEPERA uses the source-to-source transformational model outlined in Figure 1. This model uses tree-pattern matching for AST manipulation, analysis, and attribute calculation. For tedious but common tasks, such as tree-pattern matching, sub-tree creation, and sub-tree replacement, KHEPERA provides a little language for describing tree matches and for building trees. For unpredictable or language-specific tasks, such as attribute manipulation or analysis, the KHEPERA little language provides an escape to a familiar general-purpose programming language (C). Standard tree traversal algorithms are supported (e.g., bottom up, top down), as well as arbitrarily complicated syntax-directed sequencing. Rapid pattern matching is provided via data-structure maintenance, which can perform rapid pattern matches in a standard tree traversal order for many commonly-used patterns.

Familiar output mechanism. A pretty-printing facility is provided that can output the AST in an easily readable format at any time. One strong advantage of this pretty-printer when compared with other systems is that it will always be able to print the AST, regardless of how much of the transformation has been performed. If the AST is in the original input format or the original output format, then the pretty-printed program will probably be executable in the input language, L , or the output language, L' . However, if the AST being printed is one of the T_n intermediate trees, then the output will use some combination of the syntax of L and L' , with a fallback to simple S-expressions for AST constructs which do not have well-defined concrete syntax. While the program printed may not be executable, it does use a familiar syntax which may be helpful for the human when replaying transformations while debugging.

Debugging support for DSL translation. KHEPERA tracks transformation application and AST modifications, can replay the transformation sequence, and has support for answering questions about which transformations were applied at which points on the AST. This is helpful when writing and debugging the DSL processor, as well as when implementing a debugger for the DSL itself.

Transformations are either written in the high-level KHEPERA language and are transformed by KHEPERA into executable C with calls to the

KHEPERA library (as discussed in Section 4.6 and shown in Figure 8 and Figure 9); or the transformations are written using explicit calls to the KHEPERA library tree manipulation functions. In either case, low-level hooks in the KHEPERA library track debugging information when nodes or subtrees are created, destroyed, copied, or replaced. This low-level information can be analyzed to provide the ability to navigate through intermediate versions of the transformed program, and the ability to answer specific queries that support the debugging of the final transformed output:

- setting breakpoints
- determining current execution location (e.g., in response to a breakpoint or program exception)
- reporting a procedure traceback
- displaying values of variables

These tracking and debugging capabilities are the subject of Faith’s forthcoming dissertation and will be not be discussed in detail in this paper. An example of setting a breakpoint will be shown in Section 4.

2 Related Work

KHEPERA is similar to some compiler construction kits. However, these systems usually restrict the scanning and parsing tools used [6]; specify AST transformations using a low-level language, such as C [17] (instead of a high-level transformation-oriented language); or require that the AST always conforms to a single grammar specification, making translation from one language to another difficult [4, 3, 14]. Further, some systems rely on an attribute grammars for all AST transformations, without providing for a more general-purpose scheme for tree-pattern matching and replacement.

SORCERER, from the PCCTS toolkit [11], is the most similar, since it does not require the use of specific scanning and parsing tools, and since it provides a “little language” in the style of `lex` and `yacc` with embedded procedures written in another general-purpose programming language (e.g., C). SORCERER and KHEPERA share abilities to describe tree structures, perform syntax-directed translations, and support the writing of AST-based interpreters. In contrast, KHEPERA also supports rule-based translations that do not require a complete grammar specification; KHEPERA rules are well suited for the construction of “use-def” chains, data-flow dependency graphs, and other compiler-required analyses; and writing pretty-printer rules in KHEPERA does not

require a complete tree-grammar specification. This allows pretty-printing to easily take place during grammar evolution.

None of the previous systems, including SORCERER, contain built-in support for “replay” of transformations, or for automatic and transparent tracking of debugging information. When translating programs from one language to another, the “discovery” of the best order for transformation application is often difficult, involving considerable AST analysis. The code to perform this analysis is often difficult to verify or is undergoing constant change during the implementation phase of a DSL. However, after the transformations are discovered and recorded in a database, a much simpler program (i.e., one that is easier to verify) could be written that applies all of the discovered transformations in the specified order, thereby proving, by construction, that the translation preserves semantics. In this case, only the semantics preserving characteristics of the transformations themselves must be proven—not the code which performs analysis and discovery. While we have not yet implemented such a prover, we have utilized the transformation discovery and replay capabilities of KHEPERA to implement a browser that presents intermediate views of the transformation process, and which can answer typical queries posed by a debugger (see Section 4.6).

3 Overview of KHEPERA

The KHEPERA library provides low-level support for:

- building an AST
- applying transformation rules to the AST (tree traversal, matching, and replacement)
- “pretty-printing” the P' source code from the T_ℓ AST (pretty-printing is actually the σ “transformation”)

An overview of the KHEPERA system is shown in Figure 2. KHEPERA encapsulates low-level details of the DSL implementation: AST manipulation, symbol and type table management, and management of line-number and lexical information. On a higher level, library routines are available to support pretty-printing (currently, with a small language to describe how to print each node type in the AST), type inference, and the tracking functions for debugging information. Further, a “little language” has been implemented to support a high-level description of the transformation rules. If transformation rules are written in the KHEPERA language, or if

they are written in an ad hoc manner using the underlying KHEPERA AST manipulation library, then the debugging tracking and transformation replay support will be automatically provided.

An overview of how the KHEPERA system fits into a complete DSL implementation solution is shown in Figure 3. In the example shown in the next section, we explain how the scanner and parser specifications are simplified by using calls to the KHEPERA library and will provide an example showing how other important input files are specified.

In Figure 4, the “DSL Processor” from Figure 3 is expanded, showing the basic blocks that are created from the source code and showing how the DSL processor is used during the compilation of a program written in the DSL.

4 Example

A simple language translation problem based on [12] will be used to illustrate the KHEPERA system. The DSL is a subset of Fortran 90 with the addition of a *sequence comprehension* construct that can be used to construct (nested) sequences. The translation problem is to remove all sequence comprehension constructs and replace them with simple data-parallel operations, yielding a program suitable for compilation with a standard Fortran 90 compiler.

4.1 Example DSL Syntax

The lexical elements of the DSL are:

Id Num (/ /) + , : = in

A program is described by the following context-free grammar (CFG):

```

program ::= statement-list
statement ::= Id = expression
statement-list ::= statement
                | statement-list statement
expr ::= Id
        | Num
        | expr + expr
        | length( expr )
        | range( expr )

```

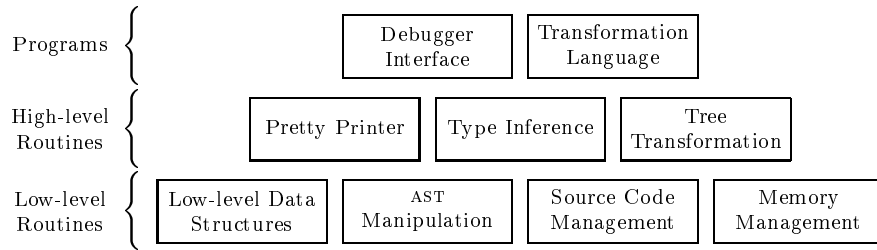


Figure 2: The KHEPERA Transformation System

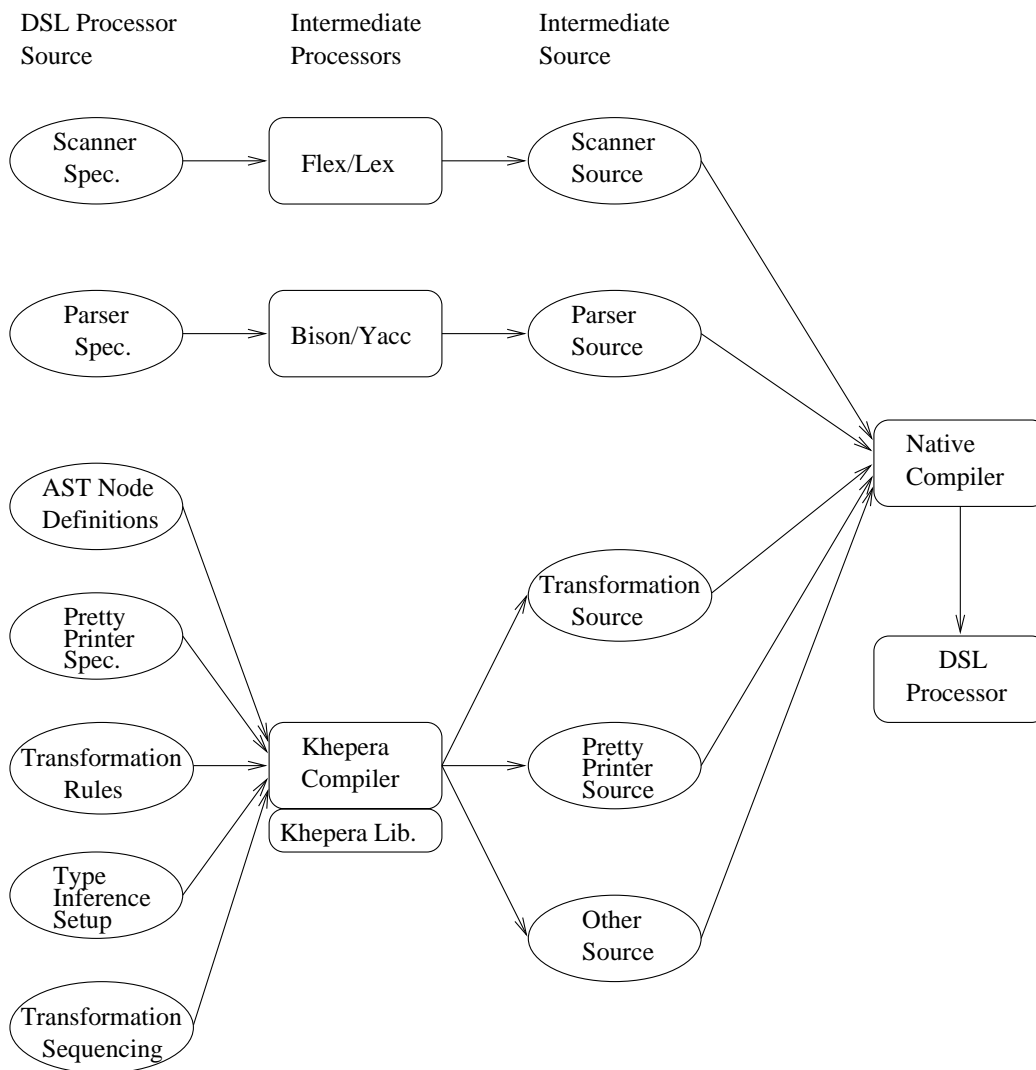


Figure 3: Using the KHEPERA Transformation System

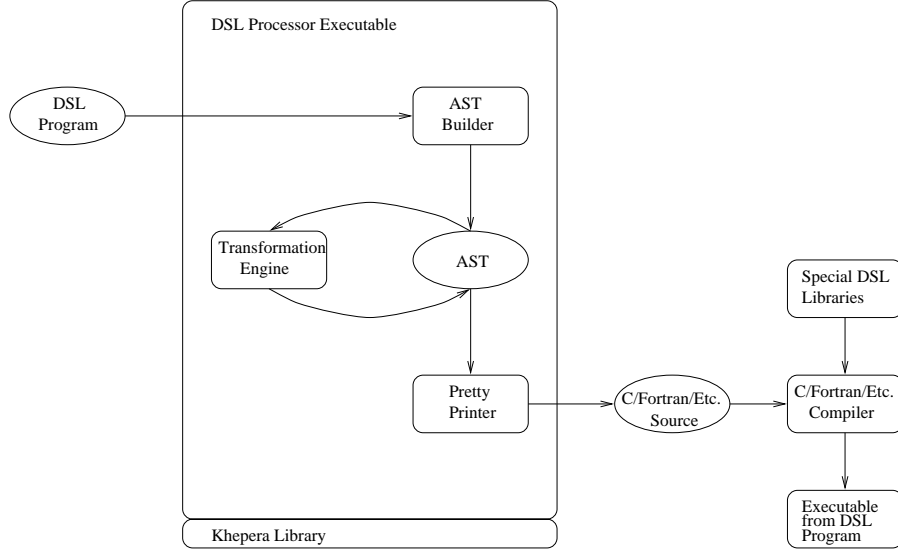


Figure 4: Using the DSL Processor

```

|  dist( expr , expr )
|  (/ expr-list /)
|  (/ Id in expr : expr /)

```

```

A = range(3);
B = (/ i in A: i + i /);
C = (/ i in A:
      (/ j in range(i): i /) /)

```

For this example, we use the array constructor notation from Fortran 90 to specify literal sequences and a similar notation to specify the sequence comprehension construct. However, the sequence comprehension construct creates arbitrarily nested, *irregular* sequences. (In contrast, the array constructor from Fortran 90 can only generate vectors or rectangular arrays.)

yields:

```

A = (/ 1, 2, 3 /)
B = (/ 2, 4, 6 /)
C = (/ (/ 1 /),
      (/ 2, 2 /),
      (/ 3, 3, 3 /) /)

```

4.2 Example DSL Semantics

DSL values have types drawn from $D = \mathbf{Int}|\mathbf{Seq}(D)$. We define, $\forall n \in \mathbf{Int}, c \in D$:

```

range(n) = (/ 1, 2, ..., n /)
dist(c, n) = (/ c, c, ..., c /)

```

with $\text{length}(\text{dist}(c, n)) = \text{length}(\text{range}(n)) = n$. For an expression, e , the sequence comprehension

```

(/ i in A : e(i) /)

```

yields the sequence of successive values of e obtained when i is bound to successive values in A .

For example, the sample program:

We omit here a collection of type (inference) rules for the language that define a well-typed program.

4.3 Example Translation

We view a program in terms of the natural AST corresponding to the CFG of Section 4.1. In the AST, an application of one of the four basic operations is written as a function application node with the operation to be applied in the *name* attribute and a *depth* attribute that is 0. The children of the node are expression(s) for each of the arguments.

The following 3 rules can be used to eliminate all sequence comprehension constructs from the AST:

Rule 1

```

(/ x1 in e1 : x1 /) → e1

```

Rule 2 Provided e_2 is an **Id** or **Num**, and $e_2 \neq x_1$,

```
(/ x1 in e1 : e2 /)
→ dist(e2, length(e1))
```

Rule 3

```
(/ x1 in e0 :
  fn_app( name = f,
          depth = d,
          args = n,
          e1, ..., en ) /)
→ fn_app( name = f,
          depth = d + 1,
          args = n,
          (/ x1 in e0 : e1 /),
          ...,
          (/ x1 in e0 : en /) )
```

The resultant AST can be written out in as Fortran 90 with the depth attribute supplied as an extra argument to the basic functions (**add**, **length**, **range**, **dist**). Given an appropriate implementation of these basic four functions, the resultant program specifies fully parallel execution of each sequence comprehension construct, regardless of the degree of nesting and sequence sizes.

For example, using these rules, the program from Section 4.2 would be transformed as follows (using $f(\dots)$ as a shorthand for $\text{fn_app}(\text{name} = f, \dots)$):

```
A = range(depth=0, 3)
B = add(depth=1, A, A)
C = dist(depth=1,
        A,
        length(depth=1,
              range(depth=1, A)))
```

Note that functions with `depth = 0` operate on scalar arguments, whereas functions with `depth = 1` operate on vector arguments.

The rules shown for this example are terminating and confluent. When the source language is more expressive and optimization becomes an issue, the rules used are not necessarily terminating, hence additional sequencing rules must be added to control rule application [10].

In the following sections, we shall show how KHEPERA can be used to implement translations, such as the one specified above, in an efficient manner.

4.4 Parsing and AST Construction

The AST is constructed using a scanner and parser generator of the implementor's choice with calls to

the KHEPERA library AST construction routines. At the level of the scanner, KHEPERA provides support for source code line number and token offset tracking. This support is optional, but is very helpful for debugging. If the implementor desires line number and token offset tracking, the scanner must interact with KHEPERA in three ways: first, each line of source code must be registered. In versions of `lex` that support states, providing this information is trivial (although inefficient), as show in Figure 5. For other scanner generators, or if scanning efficiency is of great concern, other techniques can be used. The routine `src_line` stores a copy of the line using low-level string-handling support. While the routines used in these examples are tailored for `lex` semantics, the routines are generally wrapper routines for lower-level KHEPERA functions and would, therefore, be easy to implement for other front-end tools.

KHEPERA also handles interpretation of line number information generated by the C preprocessor. This requires a simple `lex` action:

```
^#\ .*      src_cpp_line(yytext, yleng);
```

Finally, every scanner action must advance a pointer to the current position on the current line. This is accomplished by having every action make a call to `src_get(yytext)`, a minor inconvenience that can be encapsulated in a macro.

The productions in the parser need only call KHEPERA tree-building routines—all other work can be reserved for later tree walking. This tends to simplify the parser description file, and allows the implementor to concentrate on parsing issues during this phase of development. A few example `yacc` productions are shown in Figure 6. The second argument to `tre_mk` is a pointer to the (optional) source position information obtained during scanning. The abstract representation of the constructed AST is that of an n -ary tree, and routines are available to walk the tree using this viewpoint.¹

Immediately after the parsing phase, the AST is available for printing. Without any pretty-printer description, the AST is printed as a nested S-expression, as shown in Figure 7.

4.5 Pretty-printing

For pretty-printing, KHEPERA uses a modification of the algorithm presented by [9]. This algorithm

¹Physically, the tree is stored as a rotated binary tree, although other underlying representations would also be possible.

```

NL          \n
...
%%
<INITIAL>{
    .*{NL}      src_line(yytext,yy leng); yless(0); BEGIN(OTHER);
    .*          src_line(yytext,yy leng); yless(0); BEGIN(OTHER);
}
...
{NL}        BEGIN(INITIAL);

```

Figure 5: Storing Lines While Scanning

```

Statement: Identifier '=' Expression
          { $$ = tre_mk(N_Assign, $2.src,
                      $1, $3, 0); }
          ;

StatementList: Statement
              {
                $$ = tre_mk(N_StatementList,
                          tre_src($1),
                          $1, 0);
              }
          | StatementList Statement
          {
            $$ = tre_append($1, $2);
          }
          ;

```

Figure 6: Building the AST While Parsing

is linear in space and time, and does not backtrack when printing. The implementation was straightforward, with modifications added to support source line tracking and formatted pretty-printing. Other algorithms for pretty printing, some of which support a finer-grain control over the formatting, are presented in [7, 2, 15, 16].

For each node type in the AST, a short description, using `printf`-like syntax, tells how to print that node and its children. If the node can have several different numbers of children, several descriptions may be present, one for each variation. List nodes may have an unknown number of children. Multiple descriptions may be present for multiple languages, with “fallback” from one language to another specified at printing time (so, Fortran may be printed for all of those nodes that have Fortran-specific descriptions, with initial fallback to unlabeled nodes (perhaps for C or for the original DSL), and with final fallback to generic S-expressions). This fallback scheme provides usable pretty-printing during devel-

opment, even before the complete pretty-printer description is finished and debugged.

For printing which requires local analysis, implementor-defined functions can be used to return pre-formatted information or to force a line break. These functions are passed a pointer to the current node, so they have access to the complete AST from the locus being printed. While the pretty-printer is source-language independent and is unaware of the specific application-defined attributes present on the AST, the implementor-defined functions have access to all of this information. We typically use these functions to format type information or to add comments to the generated source codes.

Additional pretty-printer description syntax allows line breaks to be declared as “inconsistent” or “consistent”²; allows for forced line breaks; and permits indentation adjustment after breaks.

²See [9] for details. Each group may have several places where a break is possible. An inconsistent break will select one of those possible places to break the line, whereas a consistent break will select *all* of these places if a break is needed anywhere in the group. This allows the following formatting to be realized (assuming breaks are possible before +):

```

Inconsistent
( x = a + b + c
  + d + e + f)
Consistent
( x = a
  + b
  + c
  + d
  + e
  + f)

```


4.6 The KHEPERA Transformation Language

KHEPERA transformations are specified in a special “little language” that is compiled into C code for tree-pattern matching and replacement. A simple transformation rule conditionally matches a tree, builds a new tree, and performs a replacement. The rule that implements the first sequence comprehension elimination transformation (Rule 1 from Section 4.3) is shown in Figure 8.

Original Program:

```
A = range(depth=0, 3)
B = (/ i in A : i + i /)
C = (/ i in A :
      (/ j in range(depth=0, i) :
        i /) /)
```

Initial AST (with attribute values shown after the slash):

```
(N_StatementList
(N_Assign
(N_Identifier/"A")
(N_Call
(N_Identifier/"range")
(N_ExpressionList
(N_Integer/3))))
(N_Assign
(N_Identifier/"B")
(N_SequenceBuilder
(N_Iterator
(N_Identifier/"i")
(N_Identifier/"A"))
(N_Add
(N_Identifier/"i")
(N_Identifier/"i"))))
(N_Assign
(N_Identifier/"C")
(N_SequenceBuilder
(N_Iterator
(N_Identifier/"i")
(N_Identifier/"A"))
(N_SequenceBuilder
(N_Iterator
(N_Identifier/"j")
(N_Call
(N_Identifier/"range")
(N_ExpressionList
(N_Identifier/"i"))))
(N_Identifier/"i"))))
```

Figure 7: Example Input and Initial AST

```
rule eliminate_iterator1
{
  match (N_SequenceBuilder
         (N_Iterator id1:N_Identifier D:.)
         id2:N_Identifier)
  when (tre_symbol(id1)
        == tre_symbol(id2))
  build new with D
  replace with new
}
```

Figure 8: Simple Transformation Rule

In Figure 8, a tree pattern follows the **match** keyword. Tree patterns are written as S-expressions for convenience. The tree pattern in this example is compiled to the pattern matching code shown in the first part of Figure 9 (code for sections of the rule follow the comment containing that section).

The **when** expression, which contains arbitrary C code, guards the match, preventing the rest of the rule from being executed unless the expression evaluates to true. The **build** statement creates a new subtree, taking care to copy subtrees from the matched tree, since those subtrees are likely to be deleted by a **replace** command.

The tracking necessary for debugging and transformation replay is performed at a low-level in the KHEPERA library. However, the KHEPERA language translator automatically adds functions (with names starting with `trk_`) to the generated rules. These functions add high-level descriptive information which allows fine-grain navigation during transformation replay, but which is not necessary for answering debugger queries.

A more complicated KHEPERA rule is shown in Figure 10. This rule implements the third sequence comprehension elimination transformation (Rule 3 from Section 4.3).

The example in Figure 10 uses the **children** statement to iterate over the children of the

```

int rule_eliminate_iterator1( int *_kh_flag, tre_Node _kh_node )
{
    const char *_kh_rule = "rule_eliminate_iterator1";
    Node _kh_pt;
    Node this = NULL; /* sym */
    Node id1 = NULL; /* sym */
    Node D = NULL; /* sym */
    Node id2 = NULL; /* sym */
    Node new = NULL;

    /* match (this:N_SequenceBuilder
              (N_Iterator id1:N_Identifier D:.) id2:N_Identifier) */

    _kh_pt = _kh_node;
    if (_kh_pt && tre_id( this = _kh_pt ) == N_SequenceBuilder) {
        _kh_pt = tre_child( _kh_pt ); /* N_Node */
        if (_kh_pt && tre_id( _kh_pt ) == N_Iterator) {
            _kh_pt = tre_child( _kh_pt ); /* N_Node */
            if (_kh_pt && tre_id( id1 = _kh_pt ) == N_Identifier) {
                _kh_pt = tre_right( _kh_pt );
                if (_kh_pt) {
                    D = _kh_pt;
                    _kh_pt = tre_parent( _kh_pt );
                    _kh_pt = tre_right( _kh_pt );
                    if (_kh_pt && tre_id( id2 = _kh_pt ) == N_Identifier) {
                        _kh_pt = tre_parent( _kh_pt );
                        assert( _kh_pt == _kh_node );

                        /* when (tre_symbol(id1) == tre_symbol(id2)) */

                        if (tre_symbol(id1) == tre_symbol(id2)) {
                            trk_application( _kh_rule, _kh_node );

                            /* build new with D */

                            new = tre_copy(D);

                            /* replace with new */

                            ++*_kh_flag;
                            trk_work( _kh_rule, _kh_node );
                            tre_replace( _kh_node, new );
                        }
                    }
                }
            }
        }
    }
}
return 0;
}

```

Figure 9: Generated Tree-Pattern Matching Code

```

rule dp_func_call
{
  match (this:N_SequenceBuilder
        iter:N_Iterator
        (f:N_Call
         fn:N_Identifier
         plist:N_ExpressionList))

  build newPlist with (N_ExpressionList)
  children plist {
    match (p:.)
    build next with (N_SequenceBuilder
                    iter p)
    do { tre_append(newPlist, next); }
  }

  build call with (N_Call fn newPlist)
  delete newPlist
  do { call->prime = f->prime + 1; }
  replace with call
}

```

Figure 10: Iterator Distributing Transformation Rule

`N_ExpressionList` node, and uses the `do` statement as a general-purpose escape to C. This escape mechanism is used to build up a new list with the `tre_append` function, and to modify an implementor-defined attribute (`prime`).

KHEPERA language features not discussed here include the use of a conditional **if-then-else** statement in place of a **when** statement, the ability to break out of a **children** loop, and the ability to perform tree traversals of matched subtree sections (this is useful when an expression must be examined to determine if it is independent of some variable under consideration).

4.7 Debugging with KHEPERA

The KHEPERA library tracks changes to the AST throughout the transformation process. The tracking is performed, automatically, at the lowest levels of AST manipulation: creation, destruction, copying, and replacement of individual nodes and subtrees. This tracking is transparent, assuming that the programmer always uses the KHEPERA AST-manipulation library, either via direct calls or via the KHEPERA transformation language, to perform all AST transformations. This assumption is reasonable because use of the KHEPERA library is required to maintain AST integrity through the transformation process. Since the programmer does not have to

remember to add tracking capabilities to his transformations, the overhead of implementing debugging support in a DSL processor is greatly reduced.

The tracking algorithms associate the tree being transformed (T_i in Figure 1), the transformation rule (τ) being applied, and the specific changes made to the AST. This information can then be analyzed to answer queries about the transformation process. For example, the DSL implementor may have identified two intermediate ASTs, T_i and T_{i+1} , and may ask for a summary of the changes between these two ASTs.

On a more sophisticated level, the user may identify a node in the DSL program and request that a breakpoint be placed in the program output. An example of this is shown in Figure 11. Here, the user clicked on the scalar `+` node in the left window. In the right window, the generated program, after 13 transformations have been applied, is displayed, showing that the breakpoint should be set on the call to the vector `add` function.

At this point, the user could navigate backward and forward among the transformations, viewing the particular intermediate ASTs which were involved in transforming the original `+` into the call to `add`. The ability to navigate among these views is unique to the KHEPERA system and helps the user to understand how the transformations changed the original program. This is especially useful when many transformations are composed.

The tracking algorithms can also be used to understand relationships between variables in the original and transformed programs. For example, in Figure 12, the user has selected an iterator variable `i` which was removed from the final transformed output. In this case, both occurrences of `A` are marked in the final output, showing that these vectors correspond, in some way, to the use of the scalar `i` in the original input.

In addition to the “forward” tracking, described here, KHEPERA also supports reverse tracking, which can be used to determine the current execution point in source terms, or to map a compile or run-time error back to the input source.

5 Conclusion and Future Work

In this paper, we have presented an overview of our transformation-based approach to DSL processor implementation, with emphasis on how this approach provides increased ease of implementation and more flexibility during the DSL lifetime when compared with more traditional compiler implement-

Ra: The Khepera Viewer	
Original Program	Program after 13 iterations
<pre>A = range(depth=0, 3) B = (/ i in A : i + i /) C = (/ i in A : (/ j in range(depth=0, i) : i /) /)</pre>	<pre>A = range(depth=0, 3) B = add(depth=1, A, A) C = dist(depth=1, A, length(depth=1, range(depth=1, A)))</pre>

Figure 11: Debugging with KHEPERA (Example 1)

Ra: The Khepera Viewer	
Original Program	Program after 13 iterations
<pre>A = range(depth=0, 3) B = (/ i in A : i + i /) C = (/ i in A : (/ j in range(depth=0, i) : i /) /)</pre>	<pre>A = range(depth=0, 3) B = add(depth=1, A, A) C = dist(depth=1, A, length(depth=1, range(depth=1, A)))</pre>

Figure 12: Debugging with KHEPERA (Example 2)

tation methods.

In the previous section we have provided an overview of the KHEPERA system using a small example. We have shown how the KHEPERA library supports AST construction and pretty-printing, and have demonstrated some of the capabilities of the KHEPERA transformation language and debugging system. Many additional features of the KHEPERA system are difficult to demonstrate in a short paper. These features include low-level support for common compiler-related data structures such as hash tables, skip lists, string pools, and symbol tables and for high-level functionality such as type inference and type checking. The availability of these commonly-used features in the KHEPERA library can shorten the time needed to implement a DSL processor.

Further, we have found that keeping lists of nodes, by type, can dramatically improve transformation speed. Instead of traversing the whole AST, we traverse only those node types which will yield a match for the current rule. However, since some transformations may assume a pre-order or post-order traversal of the AST, the “fast tree walk” problem is more difficult than simply keeping node lists: the lists must be ordered and the data structure holding the lists must be updateable during the tree traversal (this eliminates many balanced binary trees from consideration for the underlying data structure). We have found that an implementation based on skip lists [13] was viable—preliminary empirical results demonstrate a significant transformation speed compared with pattern matching over the whole AST.

More details on this work will be presented in a future paper.

Another advantage of KHEPERA is the support for debugging via transformation replay. When the transformations are applied to the AST using the KHEPERA library support (with or without using the KHEPERA transformation language), then those transformations are tracked and can be replayed at a later time. KHEPERA includes support for arranging the transformations in an abstract hierarchy, thereby facilitating meaningful viewing by a DSL implementor. As part of a complete debugging system, KHEPERA also provides mappings which allow loci in the output source to be mapped back through the AST transformations to the input source (written in the DSL). These debugging capabilities are the subject of Faith’s forthcoming dissertation.

6 Availability

Snapshots of the KHEPERA library, including working examples similar to those discussed in this paper, are available from <ftp://ftp.cs.unc.edu/pub/projects/proteus/src/>.

References

- [1] Jon L. Bentley, Lynn W. Jelinski, and Brian W.

- Kernighan. CHEM—a program for phototypesetting chemical structure diagrams. *Computers and Chemistry*, **11**(4):281–97, 1987.
- [2] Robert D. Cameron. An abstract pretty printer. *IEEE Softw.*, **5**(6):61–7, Nov. 1988.
- [3] James R. Cordy and Ian H. Carmichael. *The TXL programming language syntax and informal semantics, version 7*. Software Technology Laboratory, Department of Computing and Information Science, Queen’s University at Kingston, June 1993.
- [4] James R. Cordy, Charles D. Halpern-Hamu, and Eric Promislow. TXL: a rapid prototyping system for programming language dialects. *Comp. Lang.*, **16**(1):97–107, Jan. 1991.
- [5] Matt Englehart and Mike Jackson. ControlH: A Fourth Generation Language for Real-time GN&C Applications. *Symp. on Computer-Aided Control System Design* (Tucson, Arizona, Mar. 1994), Mar. 1994.
- [6] J. Grosch and H. Emmelmann. *A tool box for compiler construction*, Compiler Generation Report No. 20. GMD Forschungsstelle an der Universität Karlsruhe, 21 Jan. 1990.
- [7] Matti O. Jokinen. A language-independent prettyprinter. *Softw.—Practice and Experience*, **19**(9):839–56, Sep. 1989.
- [8] Brian W. Kernighan. PIC—a language for typesetting graphics. *Softw.—Practice and Experience*, **12**:1–21, 1982. Published as AT&T Bell Laboratories (Murray Hill, New Jersey) Computing Science Technical Report No. 116: *PIC—a graphics language for typesetting (user manual)*, available as <http://cm.bell-labs.com/cm/cs/cstr/116.ps.gz>.
- [9] Derek C. Oppen. Prettyprinting. *ACM Trans. on Prog. Lang. and Sys.*, **2**(4):465–83, Oct. 1980.
- [10] Daniel William Palmer. *Efficient execution of nested data-parallel programs*. PhD thesis, published as Technical report TR97-015. University of North Carolina at Chapel Hill, 1996.
- [11] Terence John Parr. *Language Translation Using PCCTS and C++: A Reference Guide*. San Jose: Automata Publishing Co., 1997.
- [12] Jan F. Prins and Daniel W. Palmer. Transforming high-level data-parallel programs into vector operations. *Proc. 4th Annual Symp. on Princ. and Practice of Parallel Prog.* (San Diego, CA, 19–22 May 1993). Published as *SIGPLAN Notices*, **28**(7):119–28. ACM, July 1993.
- [13] William Pugh. Skip lists: a probabilistic alternative to balanced trees. *Communications of the ACM*, **33**(6):668–76, June 1990.
- [14] Reasoning Systems. *REFINE user’s guide*, 25 May 1990.
- [15] Lisa F. Rubin. Syntax-directed pretty printing—a first step towards a syntax-directed editor. *IEEE Trans. on Softw. Eng.*, **SE-9**(2):119–27, Mar. 1983.
- [16] Martin Ruckert. Conservative pretty printing. *SIGPLAN Notices*, **32**(2):39–44, Feb. 1997.
- [17] S. Tjiang, M. Wolf, M. Lam, K. Pieper, and J. Hennessy. Integrating Scalar Optimization and Parallelization. *Languages and Compilers for Parallel Computing (Fourth International Workshop)* (Santa Clara, California, 7–9 Aug. 1991). Published as U. Banerjee, D. Gelernter, A. Nicolau, and D. Padua, editors, *Lecture Notes in Computer Science*, **589**:137–51. Springer-Verlag, 1992. An overview of a more recent version of SUIF is available as Robert P. Wilson, Robert S. French, Christopher S. Wilson, Saman P. Amarasinghe, Jennifer M. Anderson, Steve W. K. Tjiang, Shih-Wei Liao, Chau-Wen Tseng, Mary W. Hall, Monica S. Lam, and John L. Hennessy, *SUIF: An Infrastructure for Research on Parallelizing and Optimizing Compilers*, available at <http://suif.stanford.edu/suif/suif-overview/suif.html>.
- [18] Arie van Deursen and Paul Klint. Little languages: little maintenance? *Proceedings of DSL ’97 (First ACM SIGPLAN Workshop on Domain-Specific Languages)* (Paris, France, 18 Jan. 1997). Published as *University of Illinois Computer Science Report*, <http://www-sal.cs.uiuc.edu/~kamin/dsl/109-27>, Jan. 1997.