

Program abstraction by transformation: Abstraction of Visual Basic to UML

K. Lano
kevin.lano@kcl.ac.uk
King's College London
London, UK

S. Kolahdouz-Rahimi
shekoufeh.rahimi@roehampton.ac.uk
University of Roehampton
London, UK

ABSTRACT

Program abstraction is a key step in the extraction of information from executable code, in order to understand legacy code, produce documentation in the form of models, or to perform re-engineering to an alternative program platform/language. Several special-purpose model transformation languages have been developed to perform program abstraction, however it remains an open research question what kinds of transformation facilities and techniques are most appropriate for the problem. In this case, we define a task for abstracting a subset of VB6/VBA to UML and OCL, this task can be used to perform comparative evaluation of different transformation approaches for the abstraction problem.

KEYWORDS

Program abstraction; Model-driven engineering; Reverse-engineering; Re-engineering

ACM Reference Format:

K. Lano and S. Kolahdouz-Rahimi. 2023. Program abstraction by transformation: Abstraction of Visual Basic to UML. In *Proceedings of The 15th Transformation Tool Contest (TTC 2023)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Program abstraction is the process of extracting formalised information from executable program code. The input could be either source code [15] or object code [17], and the outputs can include data flow or control flow information. The purpose could be for program comprehension [5] or for refactoring or other quality improvement of the source [4]. Here we will focus on the task of abstracting software models from source code, for the purpose of re-engineering, in particular for translating the source code to a different programming language [10].

An important property in this situation is *semantic preservation*: the abstraction should accurately capture the semantics of the source code, in order that the organisation which owns the code can have confidence that the re-engineered version still performs the same functionality as the original. Thus the abstraction needs to be expressed in a language which supports detailed specification

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
TTC 2023, 20th July 2023, Leicester, UK

© 2023 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

of behaviour. Here we propose the use of OCL [1, 13] together with UML class specifications, however other appropriate formalisms, such as activity diagrams or state machines, could also be used.

Program abstraction involves the co-use of parsing/grammar-based technologies with model-based technologies such as transformations. This is a similar situation to the combined use of grammar-based and model-based techniques for DSL tooling [2].

The specific re-engineering task is translation from Visual Basic version 6 (VB6) to Python version 3.9. To make the task practical, only a small subset of VB6 will be considered here, essentially modules with a top-level linear sequence of variable declarations and assignments.

The research questions we wish to investigate are:

- RQ1** What form of transformation language and transformation language facilities are particularly effective for program abstraction?
- RQ2** What are the specific challenges of defining program abstraction transformations?
- RQ3** Are there any transformation design patterns or idioms which are particularly relevant for this domain?
- RQ4** How should parsing and grammar-based technologies be integrated with transformations for program abstraction?

The case materials are available at: zenodo.org/record/7801436.

2 VISUAL BASIC

BASIC¹ was intended, as its name suggests, as a language for inexperienced programmers to use for relatively simple programming problems. It became popular with the advent of home PCs in the 1970s, and as Visual Basic (VB) and Visual Basic for Applications (VBA) became the main language for defining auxiliary code modules within MS applications such as Excel [12]. Visual Basic 6.0 (VB6), released in 1998, was the last version of VB prior to VB.NET, and is still supported on Windows platforms.

The principal challenges for software modernisation and re-engineering of VB/VBA are:

- The use of implicit typing for data items
- GOTO statements
- The large number of kinds of statements (67 in VB6)
- The complexity of MS applications such as Excel, with complex spreadsheet data and hundreds of application functions, which can be called from VBA code.

We will restrict the considered subset of VB6 to those programs written using variable declarations (DIM statements), assignment statements (including LSET, RSET and REDIM) and sequencing. We recommend the use of the ANTLR VB6 grammar, which is

¹Beginner's All Purpose Symbolic Instruction Code

supplied with the case materials, together with an executable parser generated from the grammar.

In terms of this grammar, the case will concern programs parsed according to the grammar/lexical rules for *module*, *moduleBody*, *moduleBodyElement*, *moduleBlock*, *block*, *blockStmt*, *letStmt* (restricted to the form *LHS = RHS*), *variableStmt*, *implicitCallStmt_InStmt*, *valueStmt*, *variableListStmt*, *variableSubStmt*, *subscripts*, *asTypeClause*, *iCS_S_VariableOrProcedureCall*, *iCS_S_ProcedureOrArrayCall*, *ambiguousIdentifier*, *baseType*, *complexType*, *argsCall*, *argCall*, *type_*, *subscript_*, *IDENTIFIER*, *literal*, *doubleLiteral*, *integerLiteral*, *STRINGLITERAL*, *TRUE*, *FALSE*, *rsetStmt*, *lsetStmt*, *redimStmt*, *redimSubStmt*, *iCS_B_MemberProcedureCall*, *iCS_S_MemberCall*. The VB6 grammar *VisualBasic6.g4* is included in the case materials *grammar* directory.

Figure 1 shows the metamodel of the considered subset of VB6. This is available as an EMF metamodel in the case materials. LSET and RSET statements are combined with LET statements in this metamodel.

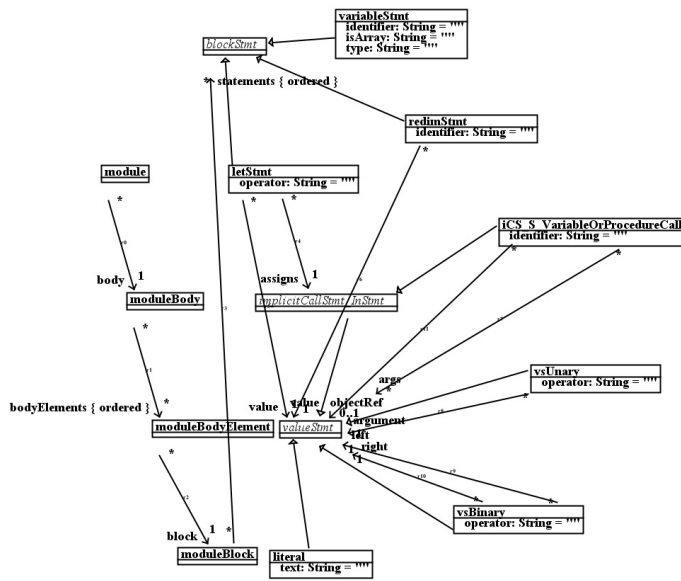


Figure 1: VB6 subset metamodel

For example, the statement

$X(2) = Y$

would be expressed in terms of this metamodel as an instance

$s : letStmt$

where

$s.operator = "="$
 $s.assigns = lhs$
 $s.value = rhs$
 $lhs : iCS_S_VariableOrProcedureCall$
 $lhs.identifier = "X"$
 $lit2 : literal$
 $lit2 : lhs.args$
 $lit2.text = "2"$
 $rhs : iCS_S_VariableOrProcedureCall$
 $rhs.identifier = "Y"$

3 CASE SPECIFICATION

The intended mapping from the VB6 subset to UML/OCL is as follows. The mapping of types is shown in Table 1.

VB6 type t	UML/OCL translation t'
Boolean, Integer, Long	Boolean, Integer, Integer
String	String
Float, Double	Real
Array type $t()$	Sequence(t')
Collection	Sequence(Map(String, OclAny))

Table 1: Mapping of VB6 types to OCL types

The VB6 data types Boolean, Integer (16-bit integers), Long (32-bit integers) and String translate directly to OCL types. However the VB6 Double is a semantically distinct subset (IEEE 754 64-bit floating point range) of OCL Real. A specific computational type *double* with the necessary properties could be used to abstract VB6 Double. The VB6 Collection type is conceptually an ordered map type, whereby elements can be accessed by index as well as by key. One way to model this is as the OCL type

$Sequence(Map(String, OclAny))$

For each basic VB6 type t , there is a default OCL value $default_t$ for the type: 0 for integer types, 0.0 for floating-point types, *false* for booleans, and the empty string "" for strings.

The mapping of expressions is shown in Tables 2 and 3.

Table 4 describes the mapping of VB6 statements to procedural OCL. Note that when elements are added to a collection using an explicit key, then the key must not already be used in the collection keyset [12].

4 SOLUTION CRITERIA

The case tasks are:

- Abstraction:** Implement the specified mapping from the VB6 subset to UML/OCL, using your chosen parsing technology and transformation language/languages.
- Validation:** Check that the test cases are correctly abstracted, by inspection or by execution of the abstracted specification.
- Translation (optional):** translate the abstracted UML/OCL into Python 3.9 and test that the result satisfies the expected semantics. An existing MDE toolset or code generator can be used for this part.

VB6 source expression e	UML/OCL translation e'
Numeric literal v	v
String literal "s"	"s"
True, False	true, false
Nothing	null
Identifier Id	Id
Array access $e(v)$	$e' \rightarrow at(v')$
Bracketed expression (e)	(e')
Floor(x)	$(x').floor()$
Max(s)	$Set\{s'\} \rightarrow max()$
Min(s)	$Set\{s'\} \rightarrow min()$
Pow(x,y)	$(x').pow(y')$
Len(s)	$(s') \rightarrow size()$
Mid(s,i,j)	$(s').substring(i', i' + j' - 1)$
Unary expressions $+e, -e$	$e', -e'$
NOT(e)	$not(e')$
Binary expressions $e1 + e2, e1 - e2,$ $e1 * e2, e1 / e2,$ $e1 \setminus e2, e1 \wedge e2,$ $e1 < e2, e1 <= e2,$ $e1 <> e2, e1 = e2,$ $e1 > e2, e1 >= e2,$ $e1 \& e2$ String $e1$ $e1 \& e2$ integer $e1$ $e1 \text{ MOD } e2$ $e1 \text{ AND } e2, e1 \text{ OR } e2$ $e1 \text{ XOR } e2$ $e1 \text{ IMP } e2, e1 \text{ EQV } e2$ $e1 \text{ LIKE } e2$	$e1' + e2', e1' - e2',$ $e1' * e2', e1' / e2',$ $e1' \text{ div } e2', (e1').pow(e2'),$ $e1' < e2', e1' <= e2',$ $e1' / = e2', e1' = e2',$ $e1' > e2', e1' >= e2',$ $e1 + e2$ $MathLib.bitwiseAnd(e1', e2')$ $e1' \text{ mod } e2'$ $e1' \text{ and } e2', e1' \text{ or } e2'$ $e1' \text{ xor } e2'$ $e1' \text{ implies } e2', (e1' = e2')$ $(e1') \rightarrow isMatch(e2')$

Table 2: Mapping of VB6 expressions to OCL expressions

VB6 expression e	UML/OCL translation e'
NEW Collection	$Sequence\{\}$
$id.Item(v)$	$id \rightarrow select(m \mid$ $m \rightarrow keys() \rightarrow includes(v') \rightarrow any() \rightarrow at(v')$
$id(v)$	
$id.Count$	$id \rightarrow size()$
$id.Items$	$id \rightarrow collect(m \mid m \rightarrow values() \rightarrow any())$
$id.Keys$	$id \rightarrow collect(m \mid m \rightarrow keys() \rightarrow any())$
$id.RemoveAll$	$Sequence\{\}$

Table 3: Mapping of VB6 collection expressions to OCL expressions

The specific criteria to be evaluated are:

- Coverage and completeness:** the abstraction transformation should be able to process the given 10 example programs. This includes the abstraction of the VB6 types *Double, Integer, Long, Boolean, String, Collection* to appropriate UML/OCL types.
- Correctness:** the abstracted specifications should be correct wrt the mapping of Tables 1, 2, 3, 4.
- Efficiency:** the abstraction process should be of practical efficiency (ie., execution time less than 1 minute for examples of 500 LOC, and a linear time complexity).

VB6 statement s	UML/OCL translation s'
DIM id AS t	$var\ id : t' := default_t$
DIM id	$var\ id : OclAny := null$
DIM $id()$ AS t	$var\ id : Sequence(t') := Sequence\{default_t\}$
DIM $id()$	$var\ id : Sequence(OclAny) := Sequence\{null\}$
$id = e$	$id := e'$
$e(v) = value$	$e' := e'.setAt(v', value')$
$id.Add\ v$	$id := id \rightarrow including(Map\{null \mapsto v'\})$
$id.Add\ Key := k,$ $Item := v$	$id := id \rightarrow including(Map\{k' \mapsto v'\})$
$id.Remove\ v$	$id := id \rightarrow excludingAt(v')$ when v integer
$id.Remove\ v$	$id := id \rightarrow select(m \mid$ $m \rightarrow keys() \rightarrow excludes(v')$ when v string
LSET $id = e$	$id := StringLib.leftAlignInto(e', id.size)$
REDIM $id(val)$	$id := Sequence\{1..(val')\} \rightarrow collect(id \rightarrow any())$
RSET $id = e2$	$id := StringLib.rightAlignInto(e', id.size)$

Table 4: Mapping of VB6 statements to OCL statements

10 small VB6 examples are provided in the *examples* directory, both in source code form and as parse trees generated by the ANTLR VB6 parser. Your solution should correctly abstract these examples and optionally translate them to correct Python code. Compute the percentage of cases which are correctly abstracted/translated.

Test cases for each program are specified in the directory *tests*. There are 21 test cases in total. Compute the overall percentage of test cases which have the same result as the source in (1) their UML/OCL representation; (2, optional) the Python target code.

Five larger examples for testing performance are given in the *performance* directory. Compute the execution time of your approach on these examples, as an average of three executions. Also provide a specification of your execution environment.

Desirable characteristics of solutions are (1) clear and modular expression of abstraction rules, for example, that the abstraction of each source language construct is defined by a specific transformation rule for that construct; (2) efficient processing of program source data and generation of target text; (3) preservation of source code structure in the abstraction and target, in order to enhance traceability; (4) adaptable and extensible transformations, which could be extended to process larger subsets of VB using the same transformation approach.

4.1 Scores for solutions

Solutions will be evaluated according to these measures:

- Corr1*: The percentage of the 10 example programs which are correctly abstracted to UML/OCL (also optionally: the percentage correctly translated to Python)
- Corr2*: The percentage of the 21 tests which have the same result in the source and abstraction (also optionally: in the source and the translation to Python)
- Perf*: Percentage of performance examples for which your approach has the same or better performance than the reference solution, on similar hardware.

5 JOURNAL PUBLICATION

Case solutions which meet a threshold standard of capabilities and scores will be selected for incorporation into a JOT article. JOT is an appropriate venue as it is concerned with the application of MDE technologies in practical software development contexts. Re-engineering of legacy systems into modernised and object-oriented platforms/languages is of high concern to businesses that utilise software [7].

6 REFERENCE SOLUTION

A solution to the abstraction part of the case is provided in the *solution* directory, using the CGTL/CSTL text-to-text transformation language [9, 11]. The VB2UML.cstl script, together with vbDeclarations.cstl and vbFunctions.cstl, defines abstraction rules for each grammar clause of the VB6 grammar. This covers almost the entire VisualBasic6.g4 grammar definition.

For example, the VB6 grammar definition for the *valueStmt* non-terminal includes the productions:

```
valueStmt :
  ...
  | valueStmt WS? AMPERSAND WS? valueStmt
  | valueStmt WS? (EQ | NEQ |
    LT | GT | LEQ |
    GEQ | LIKE | IS) WS? valueStmt
```

Thus the corresponding abstraction ruleset *valueStmt::* has rules for each of the 9 binary operators of these cases:

```
valueStmt::
  ...
  _1 & _2 |-->(_1 + _2)
  _1 = _2 |-->_1 = _2
  _1 <> _2 |-->_1 /= _2
  _1 < _2 |-->_1 < _2
  _1 > _2 |-->_1 > _2
  _1 <= _2 |-->_1 <= _2
  _1 >= _2 |-->_1 >= _2
  _1 LIKE _2 |-->(_1)->isMatch(_2)
  _1 IS _2 |-->_1 <=> _2
```

Likewise for other forms of expression and statement. A CGTL/CSTL rule

```
LHS |-->RHS
```

of ruleset *tg::* matches against AST terms with tag *tg* which correspond element-by-element to the LHS tokens. Eg., a term *t* of form (*valueStmt* *t1* & *t2*) will match against the LHS of the *valueStmt* rule

```
_1 & _2 |-->(_1 + _2)
```

with *t1* bound to *_1* and *t2* bound to *_2*.

The subterms *t1* and *t2* are then recursively mapped to strings *s1* and *s2*, and the result of the rule formed as the substitution RHS[*s1*/*_1*, *s2*/*_2*], in this case this is (*s1* + *s2*).

User-defined functions *f* can also be applied to terms by the notation *_if*, where *f* is defined by a ruleset *f::*. This enables processing of subterms of a term bound to *_i*, ie., source terms can be inspected to any depth using this technique.

For forward engineering to Python, the Python code generator of AgileUML² is used.

Figure 2 shows the overall execution time for abstraction of the five performance examples. The time is computed as the average of 3 executions, on a Windows 10 quad-core laptop (Intel i5 2.8GHz processor). Table 5 gives the solution evaluation scores for the reference solution.

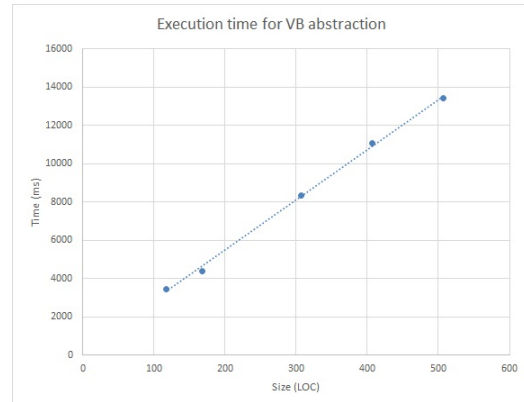


Figure 2: Performance of reference solution

Criteria	Score
Corr1	100%
Corr2	100%
Perf	100%

Table 5: Solution evaluation of reference solution

7 RELATED WORK

Related TTC cases are (1) [5] and (2) [4]. These concern (1) the extraction of state machines from Java code, and (2) the refactoring of Java code. An earlier case at GraBaTs '09 also concerned reverse engineering of Java for program comprehension [18]. This concerned the production of control flow and program dependence graphs.

The present case differs from these previous cases by (i) focussing on the fine-grained semantic modelling of program variables and data types, and (ii) by addressing a legacy source language (VB6) instead of Java. It also concerns program translation rather than comprehension or refactoring.

Specialised transformation approaches and languages have been utilised for program abstraction and re-engineering tasks: the TGraph concept and GReQL/GReTL languages are used for software migration in [3], and Gra2MoL for extracting models from code in [6]. These approaches have in common the need to effectively search and extract information from large graph or tree-structured program representations, which is a key task also in the present case. The present case however extends the scope of the abstraction task by requiring that a detailed semantic (mathematical) model

²github.com/eclipse/agileuml

is produced by abstraction, rather than specific search results or a syntactic (structural) model.

REFERENCES

- [1] Eclipse, *Eclipse OCL*, <https://wiki.eclipse.org/OCL/OCLinEcore>, 2023.
- [2] M. Eysholdt, H. Behrens, *Xtext: implement your language faster than the quick and dirty way*, OOPSLA 2010, pp. 307–309.
- [3] A. Fuhr, T. Horn, V. Riediger, A. Winter, *Model-driven software migration into service-oriented architectures*, *Comput. Sci. Res. Dev.*, vol. 28, 2013, pp. 35–84.
- [4] M. Geza Kulcsar, S. Peldszus, M. Lochau, *Case Study: object-oriented refactoring of Java programs using graph transformation*, TTC 2015.
- [5] T. Horn, *Program Understanding: a reengineering case for the Transformation Tool Contest*, TTC 2011, EPTCS.
- [6] J. Izquierdo, J. Molina, *Extracting models from source code in software modernisation*, *SoSyM* vol. 13, 2014, pp. 713–734.
- [7] R. Khadka et al., *How do professionals perceive legacy systems and software modernization?*, ICSE 2014, ACM Press, 2014.
- [8] S. Kolahdouz-Rahimi, K. Lano, et al., *A comparison of quality flaws and technical debt in model transformation specifications*, *JSS*, vol. 169, 2020.
- [9] K. Lano, Q. Xue, S. Kolahdouz-Rahimi, *Agile specification of code generators for model-driven engineering*, ICSEA 2020.
- [10] K. Lano, *Program translation using model-driven engineering*, short paper, ICSE 2022.
- [11] K. Lano, Q. Xue, *Lightweight software language processing using Antr and CGTL*, *Modelsward* 2023.
- [12] Microsoft Com, *Office VBA Reference*, <https://learn.microsoft.com/en-us/office/vba/api/overview>, Oct. 2022.
- [13] OMG, *Object Constraint Language 2.4 Specification*, OMG document formal/2014-02-03, 2014.
- [14] J. Perez et al., *Data reverse engineering of legacy databases to OO conceptual schemas*, *ENTCS* 72, no. 4, 2003, pp. 7–19.
- [15] R. Perez-Castillo, I. Garcia-Rodriguez de Guzman, M. Piattini, *Implementing business process recovery patterns through QVT transformations*, *ICMT* 2010.
- [16] R. Perez-Castillo, I. Garcia-Rodriguez de Guzman, M. Piattini, *Knowledge discovery metamodel ISO/IEC 19506: A standard to modernize legacy systems*, *Computer Standards and Interfaces*, vol. 33, 2011, pp. 519–532.
- [17] T. Sen, R. Mall, *Extracting finite-state representation of Java programs*, *SoSyM*, vol. 15 (2), 2016, pp. 497–511.
- [18] J-S. Sottet, F. Jouault, *Program Comprehension Case*, *GraBaTs* 2009.