# Parsing Fortran-77 with proprietary extensions

Younoussa Sow
*DTIPD Framatome*
Paris, France
younoussa.sow@framatome.com

Larisa Safina
*Univ. Lille, Inria, CNRS, Centrale Lille,*
*UMR 9189 CRIStAL*
F-59000 Lille, France
larisa.safina@inria.fr

Léandre Brault
*DTIPD Framatome*
Paris, France
leandre.brault@framatome.com

Papa Ibou Diouf
*DTIPD Framatome*
Paris, France
papa-ibou.diouf@framatome.com

Stéphane Ducasse
*Univ. Lille, Inria, CNRS, Centrale Lille,*
*UMR 9189 CRIStAL*
F-59000 Lille, France
stephane.ducasse@inria.fr

Nicolas Anquetil
*Univ. Lille, Inria, CNRS, Centrale Lille,*
*UMR 9189 CRIStAL*
F-59000 Lille, France
nicolas.anquetil@inria.fr

*Abstract*—Far from the latest innovations in software development, many organizations still rely on old code written in "obsolete" programming languages. Because this source code is old and proven it often contributes significantly to the continuing success of these organizations. Yet to keep the applications relevant and running in an evolving environment, they sometimes need to be updated or migrated to new languages or new platforms. One difficulty of working with these "veteran languages" is being able to parse the source code to build a representation of it. Parsing can also allow modern software development tools and IDEs to offer better support to these veteran languages. We initiated a project between our group and the Framatome company to help migrate old Fortran-77 with proprietary extensions (called Esope) into more modern Fortran. In this paper, we explain how we parsed the Esope language with a combination of island grammar and regular parser to build an abstract syntax tree of the code.

*Index Terms*—program synthesis, refactoring and reengineering, code transformation, Fortran, Esope

## I. INTRODUCTION

A significant amount of the software used nowadays in critical sectors of the industry (*e.g.* nuclear, engineering, or banking sectors) is written in languages that are considered to be "veterans" in computer science, such as Fortran, COBOL, Ada, etc. Nowadays, Cobol remains dominating in the financial domain being used in 43% of US banking systems and 95% of ATMs, and estimated as having around 220 billion lines of code[1]. Fortran prevails as the number one language for scientific, engineering, and high-performance computing (HPC) domains due to its performance, numerical computation capabilities, and extensive legacy code base. For example, even the modern SciPy Python library is actually a wrap-up around old Fortran (and C) libraries.

But working with veteran languages raises specific problems:

- They tend to have limited support for modern software engineering practices (testing, code-commenting, docu-mentation, etc.) and do not benefit from such tools as code quality and security evaluation, code completion, code auto-generation, or refactoring.
- They have small abstraction power having been designed at a time when computers were much smaller and slower and developers needed fine control over the compiled code to optimize it (*e.g.* it was normal in Fortran to take memory allocation alignment into account when declaring variables[2]). They may also lack dynamic data allocation, forcing developers to produce complex and entangled solutions.
- Veteran languages are not taught anymore, thus a recent report[3] identified as a real threat the difficulty to find top-rate Fortran computer scientists in the future;
- Hardware platforms running this code are becoming more and more outdated. The code is not always ported on new platforms and does not profit from the new technologies with better performance.

Yet, organizations are forced to support irreplaceable old code-base, if only to adapt it to new needs or security standards. This prompted the creation of different extensions of these languages, for example for Fortran: Parametric Fortran, Fortran-S, Vienna Fortran, Esope. Because they have a more restricted distribution, these extensions present the same risks listed above to an even higher degree. This creates a strong demand for methods and tools for parsing, modeling, or migrating such veteran languages and/or their extensions.

We are experiencing these difficulties as we were called on a project to migrate Esope/Fortran-77 projects to Fortran-2003 by the Framatome company[4]. Framatome designs, builds, and services nuclear steam supply systems. The first roadblock of this project is to be able to parse the Esope/Fortran-77 source code (Esope being a proprietary extension of the Fortran language).

---

[1] http://fingfx.thomsonreuters.com/gfx/rngs/USA-BANKS-COBOL/010040KH18J/index.html

[2] https://fftw.org/doc/Allocating-aligned-memory-in-Fortran.html
[3] https://permalink.lanl.gov/object/tr?what=info:lanl-repo/lareport/LA-UR-23-23992
[4] https://www.framatome.com/en/

In this paper we explain what are the challenges linked to parsing veterans languages in general, and Fortran-77 in particular. In our case, these challenges are made harder by the fact that we are dealing with an extension of Fortran that has only one known parser available. We propose an approach that allowed us to parse this language at a low cost.

The paper is structured as follows: In Section II we present the background of the project and the specificities of the two languages we have to deal with: Esope and Fortran-77. In Section III we list the existing work on Fortran migration and some parsing technologies such as Island Grammars. In Section IV, we explain why parsing these veteran languages is difficult, we list our constraints for a Fortran parser, and review the existing solutions. In Section V we propose a solution to parse the Esope extension language without the need to go into the trouble of creating a whole parser for the language. We present the result of evaluating our solution on a small but representative and independent project in Section VI. Finally, we close the paper with a discussion of future work (Section VII) and the conclusion (Section VIII).

## II. BACKGROUND

### A. Fortran

For historical reasons, Fortran is one of the main programming languages for scientific computing, so much, that the SciPy Python library is actually a wrap-up around old Fortran code.

There are several dialects of the language, standardized at different times: early versions (Fortran-66, Fortran-77, Fortran-90, Fortran-95), and modern versions, adding for example concepts from Object-Oriented programming (Fortran-2003, Fortran-2008, Fortran-2018). There are also non-standard extensions:

- Parametric Fortran [1] enables the creation of Fortran extensions using parameter structures, which can be referenced in Fortran programs to indicate the dependence of program sections on these parameters.
- Fortran-S [2] enhanced Fortran-77 with directives to specify parallelism (shared data structures and parallel loops).
- Fortran-D [3] is an extension to Fortran90 enhanced with data decomposition specifications also used for writing parallel programs, that was later integrated to a newer version, Fortran 2008.
- Vienna Fortran 90 [4] is a language extension of Fortran 90 which enables the user to write programs for distributed memory multiprocessors using global data references only.
- Esope is a Fortran extension created by the CEA[5] to allow (i) complex data structures use, (ii) automatic memory management, and (iii) automatic swap management.

Having been the standard for 13 years, Fortran-77 is historically important and has a very large code base. This is the

standard on which Esope was based and that is much used at Framatome.

Fortran programs come in two forms: fixed-column (early versions) and free-form (since Fortran-90) formats. The free-form formats resemble current programming practices where the position of a character on the line has no formal significance and can be used for example for statement indentation.

The fixed-column inherits characteristics from the old punched card systems that were used to enter programs in a computer:

- Lines are limited to 72 characters (or columns);
- A "C" or "*" as the first character of the line (column 1) marks a comment line (the line is ignored by the compiler);
- Otherwise, the first 5 characters of a line (columns 1 to 5) are for a label field: a sequence of digits that can be referred to by GOTO statements;
- Column 6 is for a continuation field: a character other than a blank or a zero there causes the line to be taken as a continuation of the preceding line.
- Columns 7 to 72 serve as the statement field;
- Anything from column 73 is ignored and can be used for comment.

Many scientific or engineering programs today still use the fixed-column format because they were written a long time ago or by people who had learned to program in early versions of Fortran. We will see that it has consequences on migration projects.

### B. Esope

Before Fortran 90, there were some issues in data management:

- No user-defined data structures (other than arrays and matrices);
- No dynamic allocation of memory;
- No automatic memory swap (memory paging).

At the end of the seventies, the CEA created an extension to Fortran 77. The goal was to create a programming language allowing (1) Complex data structures, (2) Automatic memory management, (3) Automatic swap management. The objective of Esope is to facilitate the management of data and to allow the notion of an object by the structuring of data like struct in C. The motivation was to have a set of data within a single variable. This notion was unknown to Fortran-77. This is how an entity called SEGMENT was added, as well as primitives for manipulating it.

A segment is:

- A group of Fortran variables defined by the programmer;
- Referenced by a single variable called POINTER. Knowing the pointer is enough to access all the variables contained in the structure.

Esope is an extension to Fortran meaning that before learning how Esope works, one should know the basics of Fortran, especially control structures, subroutine calls, COMMON statements, etc. The functionality of Esope allows one

to define, initialize, copy, change, and suppress a segment, get and change a segment's dimension and declare multiple pointers attached to a segment.

We give in Listing 1 an example of Esope code for the management of books and users inside a library. Esope introduces:

- A new type of definition structure, the "segment". The example in Listing 1 declares a structure for a `user` with a `name` (line 5) and an array of borrowed books (`ubb`, line 6). The size of the array is dynamically defined by the variable ubbcnt. The "field" definitions have the same syntax as regular variable definitions in Fortran;
- A new data type, `pointer` for a pointer variable (ur) pointing to a `user` structure (line 8);
- New statements to manipulate the segments (*e.g.* an "instance" of `user` is created at line 11 with `segini`). There are six new statements: *segini*, *segact*, *segadj*, *segdes*, *segprt*, and *segsup*;
- New library functions (*e.g.* `actstr`, `ajpnt`, or `mypnt`, not shown in the listing);
- New notations to access the members of a segment (dot notation, line 12) and to get the size of an array (slash notation, line 13, gets the size of the first dimension of array `ur.ubb`).

We will explain in our solution how we deal with each of these extensions that generate potential problems for parsing.

```
1        SUBROUTINE NEWUSER(LIB,NAME)
2        IMPLICIT NONE
3        INTEGER UBBCNT
4        SEGMENT, USER
5         CHARACTER*40 UNAME
6         INTEGER UBB(UBBCNT)
7        END SEGMENT
8        POINTEUR UR.USER
9 C the user does not have a book yet
10       UBBCNT = 0
11       SEGINI, UR
12       UR.UNAME = NAME
13       WRITE(*,*) UR.UBB(/1)
14 [...]
```

Listing 1. Example of Esope code (see text for explanations)

*C. Aging tools*

The Esope transpiler takes an Esope source code (Esope new instructions inter-mixed with standard Fortran-77 source code) and generates pure Fortran-77 code with memory manipulation instructions. For this, it makes strong assumptions on how the Fortran-77 compiler (Intel Fortran compiler) allocates memory. The latest versions of the compiler are not able to optimize the source code generated, which is a huge drawback for Fortran applications. It is feared that the future version of the compiler will not be able to handle the code at all.

In the prevision of this, a research project was launched to offer an automatic translation solution from Esope sources to modern Fortran (Fortran-2003) allowing an object-oriented approach. A preliminary study showed that it is possible to perform this translation from Esope sources to Fortran-2003 manually. However, the large number of sources to be translated and their size do not allow manual translation

to be considered, which explains the need for an automatic approach. The objective is therefore to offer solutions for the automatic translation of codes from the Esope overlay into modern Fortran.

## III. RELATED WORK

Maintaining and providing an up-to-date infrastructure for proprietary programming languages, domain-specific programming languages, or language dialects remains a real problem in reverse engineering [5], [6], [7].

Jonge and Monajemi [5] suggest several options for resolving maintenance problems. They include outsourcing the development in the given (veteran) language, simplifying the development process, and migrating to a similar standardized programming language. They do mention that the last one (our case) can be very challenging. According to the same article and [8], the main effort in building a parser for migrating a language is in (re)defining a grammar, as opposed to building tools around it (building parse trees, pretty printing, etc.).

Evolving a DSL requires adapting the compiler of a hosting language which in turn requires corresponding skills and background from engineers, and time and money investments from companies. Lammel and Verhoef in [8] state that implementing a high-quality Cobol parser can take up to three years, and adapting an existing parser to deal with language changes/dialects can take three to five months.

Island grammars are a possible solution to the definition of a new grammar/parser. They are a partial definition of a language's grammar, a common way to identify and annotate statements of an embedded language [9], [10]. An island grammar is characterized by its production rules, which provide comprehensive descriptions of specific constructs known as "islands," alongside more lenient rules that encompass the remaining elements referred to as "water." Island grammars work where normal grammars do not (*e.g.* in handling incomplete and syntactically incorrect code, embedded code written in other programming language or dialect, code using preprocessors, etc.) and help to avoid the tedious and expensive writing of a complete language grammar and parser.

*A. Parsing and Migration of Fortran*

There exist numerous projects dedicated to parsing and converting code written in Fortran to other programming languages including Python [11], JVM Bytecode or Java [12], [13], [14], C/C++ [15], [16], [17], [18], Pascal [19], Basic [20], Ada [21], [22], Algol [23] and others.

Bysiek et al [11] present a semi-automatic transpilation for a subset of Fortran 77/90/95 to Python 3 to retain Fortran-level performance. The paper, however, does not give any details on the parsing process in transpilation.

The first automatic Fortran to C conversion, a program called f2c, was presented by Feldman *et al.* in [15], [16]. The creation of the tool was motivated by the need to run Fortran program on a machine that had a C compiler but no Fortran compiler. It allowed to better express certain functionalities that were easier to do in C than in Fortran (*e.g.* storage

management, some character operations, arrays of functions, heterogeneous data structures, and calls that depend on the operating system), to profit from the C tools like linters and verifiers for consistency and portability checks etc.[6]. The tool is based on Feldman's previous f77 compiler [24] that originally parsed Fortran to an intermediate representation, but was tweaked to produce a C parse tree used as an input for the second pass of the portable C compiler. We review f2c parsing capability in Section IV-C.

In [17] Grosse-Kunstleve et al. show the problems of integrating Fortran code to modern Object-Oriented Programming environments (*e.g.* extensive use of global variables, communicating using intermediate data files etc.) and present a Fortran to C++ source-to-source conversion tool FABLE. Comparing to [15] generated C++ code had better readability and the tool itself could be integrated in modern modular systems.

For Java, there exist two solutions, both called f2j and focused on porting of Fortran libraries to profit from Java's portability and proliferation. The first one, by Fox et al [12] is based on the Fortran to C conversion [25], [26] and parses an input Fortran program to an AST, converts Fortran AST to C AST, and then generates java source code from it. The second one, by Dongarra, Seymour, and Doolin [13], [14] focuses on parsing Fortran to a JVM-bytecode to facilitate translation of Fortran GOTO statements (that do not always have a direct translation to a corresponding Java code) and to explore possible code optimizations on the byte-code level. The tool does lexing/parsing to build a complete AST and to use it further for optimization, type assignments, and code generation purposes.

## IV. Modeling Fortran/Esope source code

Our migration project will use a model-driven approach, a proven technique for source code modification. The challenge of the project addressed in this paper is how to build a model of the Fortran/Esope source code.

### A. Challenge: Building a model of the code

We elected to work with the Moose platform [27] that offers tools for software analysis and manipulation. It has a specific AST (Abstract Syntax Tree) meta-model independent of the programming language. However, Moose does not have a Fortran importer. This is a bottleneck of the model-driven approaches that a model of the system must be built, which involves parsing the source code. Creating custom parsers for most programming languages is a non-trivial endeavor that can be even more difficult for veteran languages [8].

Esope source is currently handled by a pre-processor written in Esope itself. Although it could seem the most natural way to go, we choose not to use this pre-processor because it would have meant learning two new languages (Esope and Fortran-77) that lag behind modern programming techniques

---

and especially do not facilitate handling abstract concepts and data structures.

Another conceivable solution would have been to apply the Esope pre-processor to obtain pure Fortran-77 source code and parse it to find back the Esope instructions that led to the generated Fortran code. However, the pre-processor generates very low-level, memory manipulation, code that does not allow to recover the high-level instructions. For example, Listing 2 shows the preprocessed result of the Esope code of Listing 1 (note that Fortran is not case sensitive). Lines 4 to 16 in Listing 2 are the translation of the segment definition and pointer declaration (lines 4 to 8 in Listing 1). Lines 18 to 23 in Listing 2 are the translation of the `segini` instruction (line 11 in Listing 1). Lines 24 to 26 (Listing 2) are the translation of the assignment (line 12 in Listing 1). Note that line 26 (Listing 2) has a character in column 6, indicating it is the continuation of the preceding line. The character itself (a ampersand) is not part of the continuation, it is only a marker. Lines 27 and 28 (Listing 2) are the translation of the assignment (line 13 in Listing 1).

```
1       SUBROUTINE NEWUSER(LIB,NAME)
2       IMPLICIT NONE
3       INTEGER UBBCNT
4 C      segment, user
5 C      POINTEUR UR.USER
6 c the user does not have a book yet
7       COMMON/OOOCOM/OOT,OOV(2),OO_001,OO_002,OO_003,
8       OO_004
9       INTEGER*8OOW(1)
10      INTEGEROOV,OOO,OO1,OO2,OO3,OO4,USER,OO5,UR
11      INTEGEROOI(1)
12      INTEGER*8OOT
13      CHARACTER*4OOH(1)
14      EQUIVALENCE(OOV(1),OOW(1),OOI(1),OOH(1))
15      INTEGEROO_001(2),OO_002(2),OO_003(2)
16      CHARACTER*4OO_004(2)
17      UBBCNT = 0
18 C      SEGINI, UR
19      CALLOOOWIN(OO4,0,'NEWUSE_10_UR_',OO1,13+UBBCNT)
20      OO_001(-0002+OOW(OOT+OO1)+1)=40
21      OO_002(-0004+OOW(OOT+OO1)+2)=13
22      OO_003(-0006+OOW(OOT+OO1)+3)=UBBCNT
23      UR=OO1
24 C      UR.UNAME = NAME
25      OO_004(-0008+OOW(OOT+UR)+1)(OOV(2)+12+1:OOV(2)+12
26      &+OOV(OOW(OOT+UR)+1))=NAME
27 C      WRITE(*,*) ur.ubb(/1)
28      WRITE(*,*) OO_005(-0010+OOW(OOT+UR)+5)
29 [...]
```

Listing 2. Fortran-77 result of the preprocessed Esope code from Listing 1

We therefore looked for a Fortran-77 parser, open-source, that we could adapt to accept Esope source code (*i.e.* with Esope extensions).

Note that the first versions of Fortran (the early fifties, work lead by J.W. Backus[7]) slightly pre-dates the theory of formal languages (Chomksy, 1954) and as such did not have a formal grammar to describe the language. For example, Fortran compilers are expected to accept programs where some "tokens" do not need to be separated by white spaces which makes parsers cumbersome to write. There are several cases of this in Listing 2, for example line 11, " `INTEGEROOI(1)` " for " `INTEGER OOI(1)` ", or line 19, " `CALLOOOWIN(...)` "

---

for " `CALL OOOWIN(...)` ". This is one of the difficulties of Fortran parsing that had to be taken into account.

## B. Criteria for an Esope/Fortran parser

A good candidate able of transforming both Esope code and Fortran-77 code in the form of an AST must offer the following:

- Allow to easily convert its AST representation to Moose AST format. An acceptable solution for this is to output the parser AST in a standard data exchange format (eg: JSON or XML) and load it in Moose;
- Keep the comments of the analyzed source code. Migration between languages implies not only technical challenges of preserving the same functionality, performance, or security of the resulting code but also of preserving the original "knowledge" of the development team [28], [29]. Our goal must be to preserve the structure, identifiers, and comments of the current code to the greatest extent possible;
- Give access to the positions in the source file of the AST nodes as several tools of Moose use this information;
- Be open-source to be adaptable to the needs of the project.

We analyzed several tools that are able to parse Fortran and/or Esope source code and produce an intermediate representation in the form of AST: compilers, dedicated parsers, static analyzers, and code converters. We summarize our results in the next section.

## C. Fortran compilers

**gfortran**[8] is a free Fortran compiler included in the GNU Compiler Collection (GCC) suite. It fully implements the Fortran95/2003/2008/2018 standards and was supporting Fortran-77 until recently (the support ceased in gcc-4.0). It can parse Fortran-77 generated by Esope. After the parsing phase, gfortran can output an AST, however, in a "not standard" format which therefore requires additional processing and creating an additional grammar/parser for it. It does not give access to the comments and positions of the instructions in the source file.

**ifort**[9] is a proprietary compiler developed by Intel. It was included in the study as it is used by the Framatome engineers to compile their sources. However, there is no easy way to obtain an AST as well as other information we require.

**lfortran**[10] is an open-source Fortran compiler that is constructed on the foundation of LLVM. lfortran offers AST in a clear readable form that can also be exported to JSON. Unfortunately, it has full support only for Fortran2018 and cannot be used for this project[11].

**flang**[12] also referred to as "Classic Flang," is a Fortran compiler intended for use with LLVM. It represents an open-source iteration of pgfortran, which is a commercial Fortran compiler produced by PGI/NVIDIA. It produces AST but in a format requiring additional work to be imported.

**F2c [15]** is a free Fortran-77 compiler and Fortran-77 to C code translator. Developed in the 90s, it became a common mean to compile Fortran code and is used in many migration tools or compilers like gfortran. Before performing the translation between languages, f2c generates a parse tree, however, we did not find a way to export it. We considered making a direct translation of the Fortran-77/Esope code to Fortran2003 code without having to go through an AST (the process is described in [15]) but the changes we will have to do look too large to be handled without a model of the code. There existed a tool to improve the readability of the source code obtained by f2c, which is, unfortunately, no longer available.

## D. Fortran-specific parsers

**fortran-src**[13] [30] is a tool written in Haskell that provides lexing, parsing, and basic analysis for Fortran66/77/90/95/2003 as a part of a bigger CamFort (Cambridge Fortran infrastructure) project for reengineering and verification of scientific Fortran programs. Fortran-src outputs an AST expressed in Haskell algebraic data types which required parsing to obtain an easy-to-use format. However, Camfort recently provided JSON output of their AST, which makes it a good candidate for our project. Fortran-src also allows keeping comments and positions in source files.

**Open Fortran Parser**[14] (OFP) consists of a parser and associated tools for Fortran2008 based on the ANTLR (ANother Tool for Language Recognition) grammar for Fortran. It outputs AST in XML format, however, it is not suited to use with Fortran-77/90 and its output did not seem reliable for our purposes.

**ANTLR**[15](ANother Tool for Language Recognition) [31] is a generator of front-ends for compilers which, starting from a grammar in Backus–Naur form (BNF) produces lexical and syntactic analyzers. There exist community-written grammars for different versions of Fortran, but not for the code generated by Esope. Indeed, Esope does not impose to separate tokens in the code which ends up causing the lexical analyzer to fail (*e.g.* "INTEGERVAR" was not parsed into "INTEGER" and "VAR"). The solution would therefore be to find a way to modify expressions in a grammar used by a lexical analyzer which for the moment has not been conclusive.

**BNFC**[16] (or the BNF Converter) is an open-source tool used in compiler construction to generate the front end of a compiler based on a labeled BNF grammar. Initially designed for generating Haskell code, it can also generate other languages (*e.g.* Ada, C/C++, Java, etc). Based on the provided grammar, the tool produces several components, including an AST, and

---

[8]https://gcc.gnu.org/wiki/GFortran

[9]https://www.intel.com/content/www/us/en/developer/tools/oneapi/fortran-compiler.html#gs.0p7v9b

[10]https://lfortran.org

[11]After this survey, we discovered that the latest version of lfortran (version 0.19.0) should be able to parse Fortran-77 (https://lfortran.org/blog/2023/05/lfortran-breakthrough-now-building-legacy-and-modern-minpack/). We come back to this in Section VII

[12]https://github.com/flang-compiler/flang

[13]https://github.com/camfort/fortran-src

[14]https://github.com/OpenFortranProject/open-fortran-parser

[15]https://www.antlr.org

[16]https://hackage.haskell.org/package/BNFC

generates lexer and parser generator files that can be used with ANTLR and other tools. However, our attempts to compile the Fortran source code generated by Esope were not successful.

**ROSE**[17] [32] is a compiler infrastructure that is open-source, designed for creating tools that perform program transformation and analysis on a source-to-source basis. It is specifically developed for large-scale applications written in Fortran 77/95/2003, C, C++, OpenMP, and UPC. Unlike most other research compilers, ROSE aims to allow non-experts to take advantage of compiler technologies to build their own custom software analyzers and optimizers. At the time of this writing, we did not yet evaluate ROSE as a solution for our project. It seems a serious candidate.

*E. Other parsers*

**PetitParser2**[18] [33] PetitParser2 is an open-source framework for building parsers developed by Jan Kurš. It is based on a unique combination of four alternative parsing methodologies:

1) scannerless parsers;
2) parser combinators;
3) parsing expression grammars;
4) Packrat analyzers.

PetitParser2 allows one to define island grammars and would be capable of recognizing Esope instructions contained in an Esope source file.

**ESOPE** is the tool used by Framatome that allows translation from Esope programming language to Fortran-77. It is a preprocessor that is written in Esope itself. It generates Fortran-77 source code (as illustrated in Listing 2).

All the above-listed solutions have been analyzed according to our stated criteria. The results of this analysis are synthetized in Table I. The column titles in it correspond to:

1) **F77**: Idiomatic Fortran-77 support;
2) **ESP**: Support for Fortran-77 generated by Esope;
3) **AST**: Offers an AST as output;
4) **CMT**: Keeps comments;
5) **POS**: Gives access to the positions in the nodes file of the AST;
6) **OS**: Open-source tool.

Given the criteria mentioned and the tests carried out on all the sources provided, we distinguish two promising tools for parsing Fortran code: gfortran and fortran-src as they are capable of analyzing Fortran-77. We chose the latter due to its ability to keep both the comments and the positions of the instructions as well as due to having AST presented in easily-parsable JSON format.

## V. PROPOSED SOLUTION

Since finding an Esope parser fulfilling our needs is difficult, we chose to follow another path. We propose to work in three steps:

| Parseurs | F77 | ESP | AST | CMT | POS | OS |
|---|---|---|---|---|---|---|
| gfortran | X | X | X | | | X |
| ifort | X | X | | | | |
| lfortran | X | | X | X | X | X |
| flang | | | | | | |
| f2c | X | | | | | |
| Fortran-src | X | X | X | X | X | X |
| OFP | | | X | | | X |
| ANTLR | X | X | X | | | |
| BNFC | X | X | X | | | X |
| ESOPE | X | X | | | | |
| PetitParser2 | | X | X | X | X | X |

TABLE I: Tools under analysis

1) "de-Esopify" the Esope/Fortran code by creating pure Fortran-77 source with Esope "annotations" (in comments);
2) parse annotated Fortran-77 code to get its AST (with the comments);
3) process the generated AST to recover the Esope constructs/instructions.

*A. "De-Esopify" the Esope/Fortran code*

To remove Esope specific constructs and instructions, we defined an island grammar [10] parser. This parser rewrites the Esope source code into "annotated" Fortran-77 code. For this, it recognizes the Esope specific constructs and instructions (island grammar), and converts them either to Fortran comment or to valid Fortran code:

- Definition of a `segment` :
  - The first line of the segment (line 4 in Listing 1) is commented out so that we will be able to recover it later. We use a special marker " `c@_` " that is unlikely to be found in normal programs. Remember that a "c" in the first column indicates a comment line in Fortran;
  - the definition of the segment members ("attributes", lines 5 and 6 in Listing 1) are kept since they are valid variable definitions in Fortran;
  - the end line of the segment (line 7) is also commented out. It would be syntactically valid ("END" can be followed by anything), but would close something that was not opened in the Fortran code.
- The pointer definition is not valid Fortran as the `pointer` type is not known, therefore, the line is commented out
- The six new Esope statements (*segini* –on line 11–, *segact*, *segadj*, *segdes*, *segprt*, and *segsup*) are again commented out;
- The new Esope functions (*e.g.* `actstr`, `ajpnt`, or `mypnt` are kept because they are already valid Fortran functions. They will have to be treated when migrating the code, but they don't raise any parsing problems;
- The dot notation to access segment "attributes" is replaced by a special expression: `ur.uname` becomes

`D__(ur,uname)` ("D" for dot). `D__` is a valid Fortran identifier and unlikely to be found in real code. The generated expression can have two interpretations in Fortran, which we will discuss in the next section;

- Simlarly, the "slash notation" ( `array(/n)` ) to get the size of `array`'s *n*th dimension (line 13 of Listing 1) is replaced by another special expression: `ur.uname(/1)` becomes `S__(D__(ur,uname),1)` ("S" for "slash"). The result is a syntactically correct Fortran expression (with the same caveat as above).

The result of "De-Esopifying" the code example of Listing 1 is given in Listing 3. Note that this code is syntactically correct for Fortran-77, but semantically incorrect. This is not a problem as we only wish to get the AST from the code to be able to analyze it.

```
1          subroutine newuser(lib,name)
2          implicit none
3          integer ubbcnt
4  c@_  segment, user
5          character*40 uname
6          integer ubb(ubbcnt)
7  c@_  end segment
8  c@_  pointeur ur.user
9  c the user does not have a book yet
10         ubbcnt = 0
11 c@_  segini, ur
12         D__(ur,uname) = name
13         bor = S__(D__(ur,ubb),1)
14 [...]
```
Listing 3. "De-Esopification" of the code example of Listing 1. The result is syntactically valid Fortran code.

### B. Parsing the annotated Fortran code

There is no difficulty here once we have a Fortran-77 parser that matches our restrictions (see Section IV-B). The annotated Fortran code is valid and should result in a valid Fortran AST. We chose a parser that keeps comments in the AST and the annotations will be kept this way.

### C. Recover the Esope constructs/instructions

Once we have the AST of the pure Fortran code, we transform it into an Esope/Fortran AST. For this, we navigate the AST, looking for our annotations (line comment nodes). We recognize our annotations by the special comment marker that we created (" `c@_` "). Each case is handled appropriately:

- A " `c@_ segment...` " comment marks the start of a segment definition. Its sibling nodes in the AST should be variable declaration statements and a " `c@_ end segment` " comment. We remove all these nodes (the two comments and the declaration statements in-between) from their parent and create a new special node `EsopeSegmentDefinitionNode` and give it the variable declaration statements as children.
- A " `c@_ pointeur...` " comment marks the declaration of pointer variable. This node is replaced by an `EsopePointerVariableDeclaration` for the variable(s). This involves a trivial parsing of the comment line where the `pointeur` keyword is followed

by a word, a dot, and another word. The first word is the name of the variable and the second is the segment type pointed to. Optionally, several pointers may be declared on the same line (separated by commas), but this does not make the parsing much more difficult.

Note that a segment definition (previous item) may contain itself a pointer "attribute". In such a case, the two annotations are treated successively to return the correct result (a `EsopeSegmentDefinitionNode` containing a `EsopePointerVariableDeclaration` ).

- A " `c@_ segini...` " comment marks a special Esope instruction. It is replaced by a new node `EsopeSegmentIntruction` that takes the name of the Esope instruction (*segini*, *segact*, *segadj*, *segdes*, *segprt*, and *segsup*) and, as a child, the pointer variable(s) treated. Again this involves some trivial parsing of the comment line where the Esope command is followed by a comma-separated list of variable names.
- The new Esope functions (*e.g.* `actstr` , `ajpnt` , or `mypnt` ) are not modified for now, but we could check all function call nodes to see whether they invoke a special Esope function or not.
- The new dot and slash notations of Esope are found in the AST in two forms:
  - When the " `D__` " or " `S__` " is an expression (right-hand side of an assignment for example) it will be seen in the AST as a call to a " `D__` " or " `S__` " functions. So " `var = D__(ptr,attribute)` " will be seen in the AST as an assignment with a variable `var` on the left-hand side and a function call " `D__` " with two arguments ( `ptr` and `attribute` ) on the right-hand side. In this case, we replace the function call node by a node `EsopeAttributeAccess` with two children (the pointer and the attribute).
  - When the " `D__` " or " `S__` " is a variable (left-hand side of an assignment), it will be seen in the AST as a "statement function", which is a concise way in Fortran to define a function performing only one computation and returning its value. So " `D__(ptr,attribute) = value` " will be seen in the AST a statement function where the function `D__` has two parameters ( `ptr` and `attribute` ) and a body which is the `value` . In this case, we replace the statement function node by a node `EsopeAttributeAccess` with two children (the pointer and the attribute).

With these transformations of our special annotations in the original Fortran-77 AST, we obtain a new AST containing Esope specific nodes and representing accurately the initial Esope source code.

### VI. EVALUATION

We evaluated our solution on the BookLibrary application, a toy Esope program. The application was created independently

by Framatome developers long before this project. It contains 10 Esope source files. Its interest lies not in the size of the source code (around 500 LOC in total) but in the fact that it uses all the different Esope-specific instructions and constructs. Note that the running example of this paper (Listing 1 and following) comes from this application.

Table II lists some descriptive statistics about the BookLibrary. For the 10 files of the example, it gives the number of lines of code (column *LOC*), the number of segment definitions (column *segment*), number of pointer definitions (column *pointer*), number of uses of one of the six Esope instructions (segini,..., in column *instr.*), number of uses of the "dot notation" (column *dot*), and number of uses of the "slash notation" (column *slash*). The three segment definitions of the example (book, user, and library) appear in the "struc.inc". This file is included in all ".E" files that manipulate "instances" of the segments.

| file | LOC | segment | pointer | instr. | dot | slash |
|---|---|---|---|---|---|---|
| struc.inc | 31 | 3 | 0 | 0 | 0 | 0 |
| borbk.E | 65 | 0 | 3 | 5 | 5 | 3 |
| findbk.E | 54 | 0 | 3 | 4 | 4 | 2 |
| findur.E | 55 | 0 | 3 | 4 | 5 | 3 |
| libpnt.E | 73 | 0 | 4 | 10 | 16 | 3 |
| main.E | 35 | 0 | 1 | 0 | 0 | 0 |
| newbook.E | 52 | 0 | 4 | 6 | 7 | 2 |
| newlib.E | 39 | 0 | 3 | 3 | 0 | 0 |
| newuser.E | 53 | 0 | 5 | 6 | 4 | 2 |
| relbk.E | 80 | 0 | 4 | 5 | 7 | 3 |

TABLE II: Descriptive statistics on the files of our validation

To validate our model of the Esope/Fortran source code, we use a "loose round-trip" comparison (see Figure 1). The round-trip comparison consists in taking a source code, modeling it as an AST in Moose, then regenerate the source code from this AST and compare it with the original code. If the model is correct and complete, the generated source code will be the same as the initial source code.
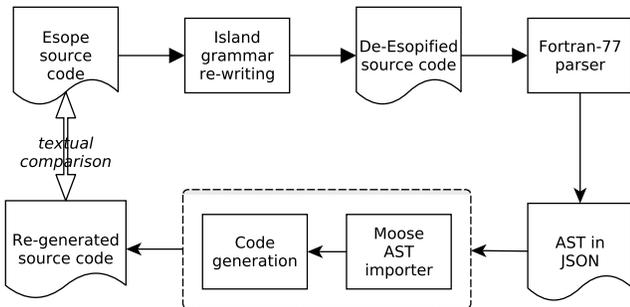


Fig. 1. Round-trip validation of our chain of tools to parse Esope/Fortran-77 source code and rebuild it from the AST

There are several problems linked to this validation:

- Whitespaces and empty lines in source code are often not significant (apart from the six spaces at the beginning of the lines in Fortran), for example " `( a .lt. 5 )` " is equivalent to " `(a.lt.5)` " despite the differences in whitespaces. Similarly, empty lines are not significant from the compiler's point of view, but they are from source code comparison.
- Also Fortran is not case-sensitive, so that " `if (a .lt. 5) then` " is equivalent to " `IF (A .lt. 5) THEN` ".
- Fortran is a somewhat flexible language offering several different syntaxes to express the same semantics. Comments can be expressed by a "C", "c", or "*" in the first column of a line, or a continuation line can be marked by any character in column 6.
- Still for flexibility, in Fortran, a subroutine declaration with no parameter can be written with or without parentheses: " `subroutine insgrp()` " can also be written " `subroutine insgrp` ".

A perfect regeneration of code should regenerate exactly the whitespace and empty lines that were in the original code although they rarely are part of the AST.

In our AST, we don't store the actual tokens found in the code because this is not relevant to our goal. For example, we don't store the actual "IF" token found in the code, but the fact that there is a "FortranIfStatement".

The rationale is that even if for the migration, we need to generate source code with the same structure (functions, instructions, flow of control) of the original one, and the same documentation (comments, identifiers) so that the developers will be able to recognize their source code and understand it. But the Esope constructs and specificities will disappear, so obviously some differences are expected. We will format the code in a way that makes it easily readable (source code pretty-printer or beautifier) and according to developers' wishes, but we will not aim to reproduce *exactly* the same format (which might actually not be correctly formatted in the first place).

Therefore, in the AST, we only indicate there is a comment line or a continuation line without keeping the exact character that marked it, or we just store a list of parameters that can be empty. When regenerating source code, we put a "C" in column 1 for comments, a "&" in column 6 for continuation lines, we write "IF" (not "if"), and "()" for an empty list of parameters. If this is not how the original code was formatted, a strict round-trip comparison will signal these differences.

For the evaluation, we, therefore, applied some modifications to the original source code:

- Apply a preprocessor to deal with the " `#include` " instruction;
- Ignore whitespace, some tools like the GNU diff program have an option for that (https://www.gnu.org/software/diffutils/manual/html_node/White-Space.html, `--ignore-all-space` );
- Convert both texts to compare to lowercase;

- Remove all "`()`" of subroutine declarations (note: it is easier to remove it from both texts than adding it to the original code when it is missing);
- Replace "`*`" at the beginning of a line (first column) with a "C".

This does not solve all problems, but greatly reduces the number of differences so that the remaining can be very quickly checked manually.

We applied our chain of tools (see Figure 1) to all the files of the BookLibrary example and the re-generated files were all equivalent to the original ones.

As an additional verification, we executed a small "test case" that uses the BookLibrary (creates books and users, loan and return books) and prints the state of the library after each action. We compared the printed information of this small test on the original and re-generated code and both outputs are strictly equal.

## VII. FUTURE WORK AND DISCUSSIONS

The markers for annotations ( `c@_` ), dot and slash notations ( `D__(...)` , and `S__(...)` ) were chosen to minimize the probability of finding them in real source code. We are currently just assuming they do not appear in the company's code. For more security, we should add a test on the source code before applying the "de-Esopification". This is a trivial task as a simple text search in the source code will tell us if these strings are found. In such a case, a warning should be raised and different markers should be chosen.

Because the Fortran parser that we are currently using (camfort) is permissive, we did not check for the length of the generated lines. When a line contains several cases of dot and slash notations, there is a high risk that the "de-Esopified" line will exceed 72 characters (we replace the single dot character by six characters: " `a.b` " becomes " `D__(a,b)` "). With more strict Fortran parsers, this should be checked and a continuation line (a line with a character in column six) should be introduced. We have so far resisted doing it because we wish to keep the "de-Esopified" code as close as possible to the Esope code.

Since our initial analysis of Fortran parsers (see Section IV), we discovered that lfortran evolved to handle Fortran-77, that ROSE looked like a good candidate, and that there was another parser (SYNTAX[19]) that could also be used. This is not a threat to our results since our approach accepts any Fortran parser. Changing the parser only requires adapting the Moose AST importer (see again Figure 1) which is not a difficult step.

## VIII. CONCLUSIONS

In industry, a vast amount of companies are still using software written in old programming languages that often do not conform to modern standards and best practices and are not handled by modern development environments and tools. However, companies are forced to keep such software since there does not exist a reliable solution to migrate it. The situation gets worse if it concerns the rarely used programming languages: dialects, proprietary languages, DSL, etc. In this paper, we present a part of the project made for the Framatome company in migration of the code written in such language called Esope which is an extension of Fortran-77. Here, we focus on parsing Esope: underlying challenges and possible solutions. We present an analysis of existing parsers for Esope/Fortran-77 and propose a solution consisting of: (1) using an island-grammar parsing phase that remove Esope constructs from the source code and creating pure Fortran-77 source with commented Esope "annotations", (2) parsing annotated ("de-esopified") Fortran-77 code to get its AST (with the comments), and (3) processing the generated AST to recover the Esope constructs. We evaluate our solution on the 500 LOC example provided by Framatome, that is, however, complete as it covers all possible Esope constructs.

The work described in this paper is only the very first step of the project that should lead to the migration of Esope code into Fortran-2003 code

## REFERENCES

[1] Z. F. Martin Erwig and B. Pflaum, "Parametric fortran: Program generation in scientific computing," *Journal of Software Maintenance and Evolution*, vol. 19, no. 3, pp. 155–182, 2007.

[2] F. Bodin, L. Kervella, and T. Priol, "Fortran-s: A fortran interface for shared virtual memory architectures," in *Proceedings of the 1993 ACM/IEEE Conference on Supercomputing*, ser. Supercomputing '93. New York, NY, USA: Association for Computing Machinery, 1993, pp. 274–283. [Online]. Available: https://doi.org/10.1145/169627.169732

[3] G. C. Fox, S. Hiranandani, K. Kennedy, C. Koelbel, U. Kremer, C.-W. Tseng, and M. L. C. Wu, "Fortran d language specification," Dept. of Computer Science, Rice University, Tech. Rep. CRPC#TR90–141, 1990.

[4] S. Benkner, B. Chapman, and H. Zima, "Vienna fortran 90," in *Proceedings Scalable High Performance Computing Conference SHPCC-92*, 1992, pp. 51–59.

[5] M. de Jonge and R. Monajemi, "Cost-effective maintenance tools for proprietary languages," in *Proceedings IEEE International Conference on Software Maintenance. ICSM 2001*, 2001, pp. 240–249.

[6] A. van Deursen and P. Klint, "Little languages: little maintenance?" *J. Softw. Maintenance Res. Pract.*, vol. 10, pp. 75–92, 1998.

[7] H. Sneed and C. Verhoef, "Re-implementing a legacy system," *Journal of Systems and Software*, vol. 155, pp. 162–184, 2019. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0164121219301050

[8] R. Lammel and C. Verhoef, "Cracking the 500-language problem," *IEEE Software*, vol. 18, no. 6, pp. 78–88, 2001.

[9] L. Moonen, "Generating robust parsers using island grammars," in *Proceedings Eight Working Conference on Reverse Engineering (WCRE 2001)*, E. Burd, P. Aiken, and R. Koschke, Eds. IEEE Computer Society, Oct. 2001, pp. 13–22.

[10] ——, "Lightweight impact analysis using island grammars," in *Proceedings 10th International Workshop on Program Comprehension*, 2002, pp. 219–228.

[11] M. Bysiek, A. Drozd, and S. Matsuoka, "Migrating legacy fortran to python while retaining fortran-level performance through transpilation and type hints," in *Proceedings of the 6th Workshop on Python for High-Performance and Scientific Computing*, ser. PyHPC '16. IEEE Press, 2016, p. 9–18.

[12] G. Fox, X. Li, Z. Qiang, and W. Zhigang, "A prototype of Fortran-to-Java converter," *Concurrency: Practice and Experience*, vol. 9, no. 11, pp. 1047–1061, Nov. 1997. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1002/%28SICI%291096-9128%28199711%299%3A11%3C1047%3A%3AAID-CPE348%3E3.0.CO%3B2-V

[13] K. Seymour and J. Dongarra, "Automatic translation of fortran to jvm bytecode," in *Proceedings of the 2001 Joint ACM-ISCOPE Conference on Java Grande*, ser. JGI '01. New York, NY, USA: Association for Computing Machinery, 2001, p. 126–133. [Online]. Available: https://doi.org/10.1145/376656.376833

[19]https://fr.wikipedia.org/wiki/SYNTAX

[14] D. Doolin, J. Dongarra, and K. Seymour, "Jlapack - compiling lapack fortran to java," *Scientific Programming*, vol. 7, no. 2, pp. 111–138, 1999.

[15] S. I. Feldman, "A fortran to c converter," *SIGPLAN Fortran Forum*, vol. 9, no. 2, p. 21–22, oct 1990. [Online]. Available: https://doi.org/10.1145/101363.101366

[16] S. I. Feldman, D. M. Gay, M. W. Maimone, and N. L. Schryer, "Availability of F2c-a Fortran to C Converter," *SIGPLAN Fortran Forum*, vol. 10, no. 2, pp. 14–15, Jul. 1991. [Online]. Available: http://doi.acm.org/10.1145/122006.122007

[17] R. W. Grosse-Kunstleve, T. C. Terwilliger, N. K. Sauter, and P. D. Adams, "Automatic fortran to c++ conversion with fable," *Source Code for Biology and Medicine*, vol. 7, no. 1, p. 5, may 2012. [Online]. Available: https://doi.org/10.1186/1751-0473-7-5

[18] "F2CPP: a Python script to convert Fortran 77 to C++ code," http://sourceforge.net/projects/f2cpp/, accessed: June 16, 2023.

[19] R. A. Freak, "A fortran to pascal translator," *Softw. Pract. Exp.*, vol. 11, no. 7, pp. 717–716, 1981. [Online]. Available: https://doi.org/10.1002/spe.4380110708

[20] R. B. Caringal and P. M. Dung, "A fortran iv to quickbasic translator," *SIGPLAN Not.*, vol. 27, no. 2, p. 75–87, feb 1992. [Online]. Available: https://doi.org/10.1145/130973.130979

[21] J. K. Slape and P. J. L. Wallis, "Conversion of fortran to ada using an intermediate tree representation," *Comput. J.*, vol. 26, no. 4, p. 344–353, nov 1983. [Online]. Available: https://doi.org/10.1093/comjnl/26.4.344

[22] M. Parsian, B. Basdell, Y. Bhayat, I. Caldwell, N. Garland, B. Jubanowsky, and J. Robinette, "Ada translation tools development: Automatic translation of fortran to ada," *Ada Lett.*, vol. VIII, no. 6, p. 57–71, nov 1988. [Online]. Available: https://doi.org/10.1145/51634.51637

[23] J. A. Prudom and M. A. Hennell, "Some problems concerning the automatic translation of fortran to algol 68," in *Proceedings of the Strathclyde ALGOL 68 Conference, Glasgow, Scotland, March 29-31, 1977*. ACM, 1977, pp. 138–143. [Online]. Available: https://doi.org/10.1145/800238.807153

[24] J. Feldman, "High level programming for distributed computing," *CACM*, vol. 22, no. 6, Jun. 1979.

[25] G. Zhang, B. Carpenter, G. C. Fox, X. Li, X. Li, and Y. Wen, "Pcrc-based HPF compilation," in *Languages and Compilers for Parallel Computing, 10th International Workshop, LCPC'97, Minneapolis, Minnesota, USA, August 7-9, 1997, Proceedings*, ser. Lecture Notes in Computer Science, Z. Li, P. Yew, S. Chatterjee, C. Huang, P. Sadayappan, and D. C. Sehr, Eds., vol. 1366. Springer, 1997, pp. 204–217. [Online]. Available: https://doi.org/10.1007/BFb0032693

[26] Q. Zheng, "A fortran to c translator based on sigma system," PACT, Technical Report, 1994.

[27] N. Anquetil, A. Etien, M. H. Houekpetodji, B. Verhaeghe, S. Ducasse, C. Toullec, F. Djareddir, J. Sudich, and M. Derras, "Modular moose: A new generation of software reengineering platform," in *International Conference on Software and Systems Reuse (ICSR'20)*, ser. LNCS, no. 12541, Dec. 2020.

[28] S. Bragagnolo, "A Holistic Approach to Migrate Industrial Legacy Systems," Theses, Universite de Lille ; Inria, May 2023. [Online]. Available: https://inria.hal.science/tel-04132315

[29] B. Verhaeghe, A. Etien, N. Anquetil, A. Seriai, L. Deruelle, S. Ducasse, and M. Derras, "GUI migration using MDE from GWT to Angular 6: An industrial case," in *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER'19)*, Hangzhou, China, 2019, pp. 579–583.

[30] D. A. Orchard and A. C. Rice, "Upgrading fortran source code using automatic refactoring," in *Proceedings of the 2013 ACM Workshop on Refactoring Tools, WRT@SPLASH 2013, Indianapolis, IN, USA, October 27, 2013*, E. R. Murphy-Hill and M. Schäfer, Eds. ACM, 2013, pp. 29–32. [Online]. Available: https://doi.org/10.1145/2541348.2541356

[31] T. Parr, *The Definitive ANTLR Reference: Building Domain-Specific Languages*. Pragmatic Programmers, May 2007.

[32] D. Quinlan and C. Liao, "The ROSE source-to-source compiler infrastructure," in *Cetus users and compiler infrastructure workshop, in conjunction with PACT*, vol. 2011. Citeseer, 2011, p. 1.

[33] J. Kurš, G. Larcheveque, L. Renggli, A. Bergel, D. Cassou, S. Ducasse, and J. Laval, "PetitParser: Building modular parsers," in *Deep Into Pharo*. Square Bracket Associates, Sep. 2013, p. 36.