

BACHELOR THESIS
COMPUTING SCIENCE



RADBOUD UNIVERSITY

Formalizing the C preprocessor

Author:
Alexander Wolters
s1022959

Daily supervisor/assessor:
Dr. Freek Wiedijk
freek@cs.ru.nl

Second supervisor/assessor:
Dr. Robbert Krebbers
mail@robbertkrebbers.nl

January 14, 2022

Abstract

There have been a number of formalizations which model the C programming language, as defined by its ISO standard documents, but in general, such models simply take the C preprocessor as given.

In this thesis, we present a model of a simplified version of the C preprocessor and investigate how such a model can be used to prove the correctness of statements involving both preprocessing and language-level requirements, using a minimal toy language in place of C.

Contents

1	Introduction	3
1.1	Related Work	4
1.2	Overview	4
2	Behaviour of the C preprocessor	6
2.1	Examples	6
3	A toy language	10
4	Notation	12
4.1	Sets and Types	12
4.2	Lambda functions	13
4.3	Partial functions	13
4.4	Updating	13
5	Formal grammars	14
5.1	Common non-terminals	14
5.2	Before expansion	15
5.3	After expansion	15
6	Preprocessing function	17
6.1	Wrapper and main function	19
6.2	Object Macro	20
6.3	Function-Like Macro	20
6.4	Text Line	20
6.5	Auxiliary functions	23
6.6	Example: Simple macro derivation	23
7	Interpretation	26
7.1	Unparsing function	26
7.2	Processing Function	28
7.3	Evaluation	28
8	Correctness of the expansion process	31

9	Predicates over the state(s)	35
10	Discussion	38
10.1	Expansion to other directives	38
10.2	Expansion to C	39
A	Common.hs	43
B	PreProcess.hs	44
C	PostProcess.hs	47

Chapter 1

Introduction

Ever since it became a major language (see [1]), there have been several formalizations of the behaviour of the C language (or at least parts of it) as defined in the C standard ([2],[3],[4],[5],[6], among others).

However, those papers generally just assume the C preprocessor and the standard header files as a given, or relegate it to future work.

As such, the questions that we want to consider in this thesis are these:

1. How can we model the behaviour of the C preprocessor, specifically of the macro replacement directives?
 - And how can we do this with as little state as possible to simplify proofs?
2. (How) can we use this model to prove the correctness of program and header files?

In this thesis, we only consider files without condition macros, and instead focus on presenting a mathematical model of a simplified version of the C preprocessor. This simplified version only includes macro replacement directives and omits conditionals, include directives and anything else the C preprocessor usually does. While this is only one aspect, it serves as a basis to build upon, especially with regard to verification of C header files. The model uses as its core idea a transformation of lines to other lines, with as little state as possible being carried between lines. The output of this model is a list of tokens which form a program. In this thesis, the program is not a C program, but rather is written in a toy language, but this is irrelevant, as the preprocessor is largely language-agnostic (except for the definition of what is a token, and the model can cope with adding to that definition just fine).

This is a first step towards building a model of the full preprocessor, which is important as it is required to prove the formal correctness of C programs that include preprocessor directives, which in practice is true of most large-scale programs, but to the best of my knowledge, no such model exists at

this point in time. Such a model would also allow the formal verification of C header files, which form an essential part of the C ecosystem, but are typically merely assumed to be correct.

1.1 Related Work

As already mentioned, there has been little interest in the C preprocessor in academic work so far. Indeed, there is only one work that I am aware of that investigates the difficulties that arise from this additional step in compiling a C program.

Kastner et al. [7] have investigated the challenges of proving statements over code for all values of conditional macros, and while they do acknowledge the problems that arise from the combination of conditionals with macro expansion, their solution builds on a pre-existing compiler to handle the macro expansion itself, and they do not consider this a focus of their work. Consequently, their work, while interesting in itself, is no help if one wishes to have a model that can be written in a proof assistant, as this thesis aims to achieve.

1.2 Overview

We first go over some preliminary knowledge, which consists of select examples of the behaviour of the C preprocessor in chapter 2, an explanation of the toy language in chapter 3 and some notation used in this thesis in chapter 4. Then we give the formal grammar of our languages before and after the preprocessing step in chapter 5. After that, we give the functions for the preprocessing (chapter 6) and the evaluation (chapter 7) of our language, the former of which is the answer to the first research question. Finally, we consider how to show the correctness of our model in chapter 8 and answer the second research question in chapter 9, before we finish up with our overall conclusions and a discussion on expandability of the model and its applicability to C in chapter 10. A large-scale overview of this thesis in a picture is given in fig. 1.1.

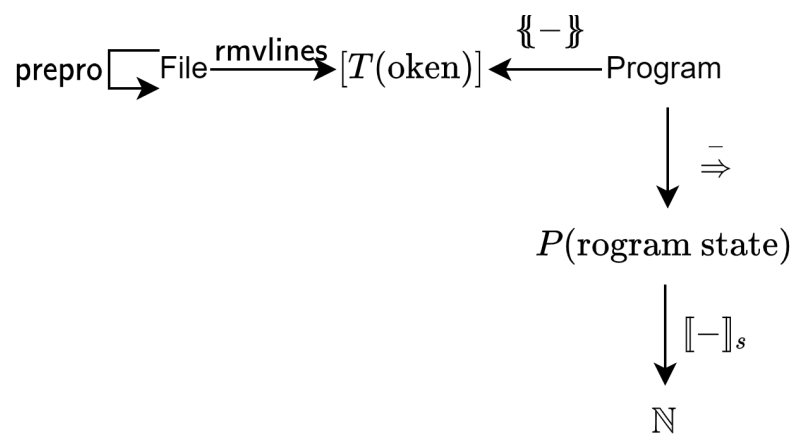


Figure 1.1: Large-scale overview of the thesis

Chapter 2

Behaviour of the C preprocessor

For a full description of the C preprocessor and all its directives, please refer to [8], with particular attention to sections 5.1.1.2 and 6.10. Here, I merely give some examples to give the reader a general understanding of the directives from the standard that are formalized in this thesis.

2.1 Examples

As the preprocessor is primarily a transformation over lines, the examples are given in a similar style, with the left side being the input line and the right side the output resulting thereof. As some lines do not give any output, the corresponding lines on the right side will be left blank.

Horizontal lines are used to separate examples from each other, consider each block between horizontal lines to be entirely unrelated to any other such block.

#define A B A	B
#define A B #define B C A B	C C
#define B C #define A B A B	C C
#define A(X) X A(2)	2
#define A(X) X + 2 A(3)	3 + 2
#define A(X) 3 A(2)	3
#define B(X) 2 + X #define A B A(2) B	B(2) B
#define A(X) X, Z declare f(A(Z)); Note that this is not legal in the toy language, but that is irrelevant to the preprocessor.	declare f(Z, Z);
#define ZERO(X) const X = 0; ZERO(z)	const z = 0;
A #define A B	A
#define A B #define B A A B	A B

#define A B #define B C #define C A A	A
#define A(X,Y) X + Y A(2, 3)	2 + 3
#define A(X,Y) X + Y A(2, #define B C 3)	2 + 3
A(2, #define A(X,Y) X + Y 3)	A(2, 3)

As not all macro definitions are legal, I also give some examples of legal and illegal macro definitions, but they can all essentially be summed by this quote from the C standard: "An identifier currently defined as an object-like macro shall not be redefined by another #define preprocessing directive unless the second definition is an object-like macro definition and the two replacement lists are identical" ([8], 6.10.3.2), and similarly for function-like macros.

Legal:

#define A B
#define A B
#define A(X) X
#define A(X) X
#define A(X) X
#define A(Y) Y
#define A(X, Y) X + Y
#define A(Y, X) Y + X

Note that the last two cases are technically illegal according to the standard as well. This is the one corner case where this thesis does not quite follow the standard and is slightly more permissive instead.

I also give some examples of illegal definitions where macros are redefined differently than they were:

```
=====
#define A B
#define A C
-----
#define A B
#define A(X) X
-----
#define A(X) X
#define A B
-----
#define A(X, Y) X + Y
#define A(Y, X) X + Y
-----
#define A B
#define C A
#define C B
-----
#define A(X, X) X
-----
```

Chapter 3

A toy language

The toy language that is used in this thesis as the target of the preprocessor is a very simple language. The full grammar is given in section 5.3, but in short, a program is a list of statements, and there are only three types of statements: Function declarations, definitions, and constant definitions.

The following example illustrates each of these three in turn:

```
declare f(y);  
define f(x) = ((5 + x) * 3);  
const c = f(2);
```

In this example, the value of the constant c is 21. As the language is not very complex, it is expected that the meaning of any expression is easy enough to determine for anyone with experience with any C-like language. However, some things are not permitted, so the next two examples would be invalid programs:

```
define f(x) = ((5 + x) * 3);  
define f(x) = ((5 + x) * 3);
```

```
const c = 12;  
const c = 13;
```

Note that it does not matter if the second definition is the same as the first or different, no function or constant may be defined twice.

It is allowed, however, to declare a function multiple times, even with different parameter names (But parameter names must be unique, even for declarations):

```
declare f(x);  
declare f(y);
```

Of course, the number of parameters must be the same, so this would be forbidden:

```
declare f(x);  
declare f(x, y);
```

And a function must be defined if it is to be used, though there is no restriction on the order of the two, so while this is not forbidden, the value of c is not defined until a definition for f is given (note that this example would be valid even without the function declaration):

```
declare f(x);  
const c = f(5);
```

The following example is not forbidden, but it will never be possible to evaluate the constant c :

```
const c = (c + 1);
```

Indeed, any constant or function definition that refers back to itself is legal, but in general, the value of such a constant or function cannot be used, unless there is a multiplication with zero involved:

```
const c = (0 * d); define f(x) = (x * f((x - 1)));
```

Note that even though the constant d does not exist yet, the evaluation of c is possible, as what is written on the right hand side does not matter. Also, since $f(0) = 0$, the function f always returns 0, and not the factorial of its input, as one may assume.

This example also demonstrates that it is legal for an expression to contain a constant name that is not yet introduced, as long as the name is defined when it is evaluated.

Chapter 4

Notation

4.1 Sets and Types

Wherever this thesis uses symbolic notation to represent the type for some function or value, the following symbols always have the same meaning:

- Given some type X , $[X]$ stands for a list of elements of type X , which may have repetitions, is well ordered and has variable length (unlike tuples).
There are also several constructors for and functions on lists, which have the same meaning as the versions in the Haskell standard library if nothing else is explicitly mentioned.
- Given some type X , we use $X?$ to represent a type that could be called "possibly X ", that is, a type which is either really a X or nothing at all, which we denote as $?$. (This type is analogous to Haskell's `Maybe` datatype [9] and Java's `Optional` class [10].)
Given an element $x \in X?$, we can write $x\#$ to express that $x \in X$, and we can write x'' to express that it is $?$.
- B stands for the set of Boolean values, that is $B = \{T, F\}$, where the symbols T and F represent the Boolean values true and false, respectively.
- id stands for the set of valid identifiers, which is a letter or underscore followed by any number of letters, digits and underscores. Identifiers are printed in monospace.
- T stands for the set of tokens, as defined in the grammar in section 3.1.1.

4.2 Lambda functions

As there are several higher-order functions in this thesis, and it is not useful to give a name to each function, we may wish to pass in as arguments to them, we make use of lambda functions as a shorthand notation. In this document, they are written as such: $\lambda x, y. x + y$, which denotes a function that takes two inputs and returns the sum of those inputs.

Note that this notation does not include types, they can be inferred from the context where the lambda function is given.

4.3 Partial functions

In this thesis, partial functions are sometimes used where no reasonable definition would exist for some input to a function, or to represent some kind of error state (For example, trying to redefine a macro constant with a different value). In this case, the type of the function is written as $f : A \rightarrow B$, where A and B are types.

In such a case, we write $f(x) \#$ to represent that f is defined for input x , and we write $f(x) "$ to represent that f is not defined for input x .

We also define the special function σ , which is not defined for any input, and therefore can match any function type. (Much as the empty set can be a set of any type.)

4.4 Updating

We use the notation of $f[a \nabla b]$ for some (possibly partial) function $f : A \rightarrow B$ and some values $a \in A, b \in B$ to represent a new (possibly partial) function, which gives the same result as f for all inputs, except for input a , for which it gives b .

We use the same notation for lists as well, but here, a is a natural number, and the meaning is that the a 'th element of the list is set to b , with all the other elements unaltered.

Chapter 5

Formal grammars

We give formal grammars for our languages, both for before and after the macro expansion. Terminals are given in typewriter font, whereas non-terminals are given in *italics*. Subsequent indented lines indicate alternative definitions. λ represents the empty string.

5.1 Common non-terminals

ArgumentList:

Identi er
Identi er , ArgumentList

Token:

Identi er
IntegerLiteral
(one of)
+ - * / () ; = ,

Identi er:

Letter
Identi er Letter
Identi er Digit

Letter:

(one of)
a b c d e f g h i j k l m n o p q r s t u v w x y z
A B C D E F G H I J K L M N O P Q R S T U V W X Y Z _

IntegerLiteral:

IntegerLiteral Digit
Digit

Digit:

(one of)

0 1 2 3 4 5 6 7 8 9

Note that *Identifier* is effectively the same as *identifier* is in the C standard, though, for simplicity, we do not use universal characters.

5.2 Before expansion

File:

λ

Line m File

Line:

TokenList

MacroLine

TokenList:

λ

Token TokenList

MacroLine:

#define Identifier TokenList

#define Identifier(ArgumentList) TokenList

5.3 After expansion

Program:

λ

Statement Program

Statement:

FuncDecl

FuncDefn

ConstDefn

FuncDecl:

declare Identifier (ArgumentList) ;

FuncDefn:

define Identifier (ArgumentList) = Expression ;

ConstDefn:

const Identifier = Expression ;

Expression:

(Expression BinOp Expression)

IntegerLiteral

Identifier

Identifier (ExpressionList)

ExpressionList:

Expression

Expression, ExpressionList

BinOp:

(one of)

*+ - * /*

Chapter 6

Preprocessing function

The preprocessing function turns a *File* with macros into a *File* that only consists of *TokenLists*. To define this function, we create several sub-functions, one for each type of possible line, which are then combined into one overarching function. This function is only a partial function, it is undefined if the *File* would not be legal to write according to the restrictions laid out in the C standard (e.g. redefining a macro with a different value). Symbolically, the `\main` function for this step has the type $\text{prepro} : \text{File} \rightarrow \text{File}$.

For brevity, we define two types that we use throughout this section to avoid giving the expanded version each time:

- M stands for the set of macro mappings, which store the macro definitions.

Such a mapping is a partial function that maps macro names to information about the number of arguments (if any) and their replacements. If the macro takes arguments (i.e. the first element of the tuple is greater than 0), it is a function like macro. In this case, the replacement may contain non-negative integers, which represent the places where the arguments should be substituted. Symbolically:

$$M = \text{File} \rightarrow (\mathbb{N} \times [T \mid \mathbb{N}])$$

- S stands for the state related to the expansion of function-like macros. This state must be preserved between lines, as macro arguments may span multiple lines.

Such a mapping is a 6-tuple of multiple different types, namely

$$S = (f, 0, 1, 2, g, \mathbb{N} \times [[T]] \times P(T) \times T^? \times \mathbb{N})$$

Some explanation is in order as to the meaning of those tuple members: The first member is always 0, 1 or 2, and it represents the state

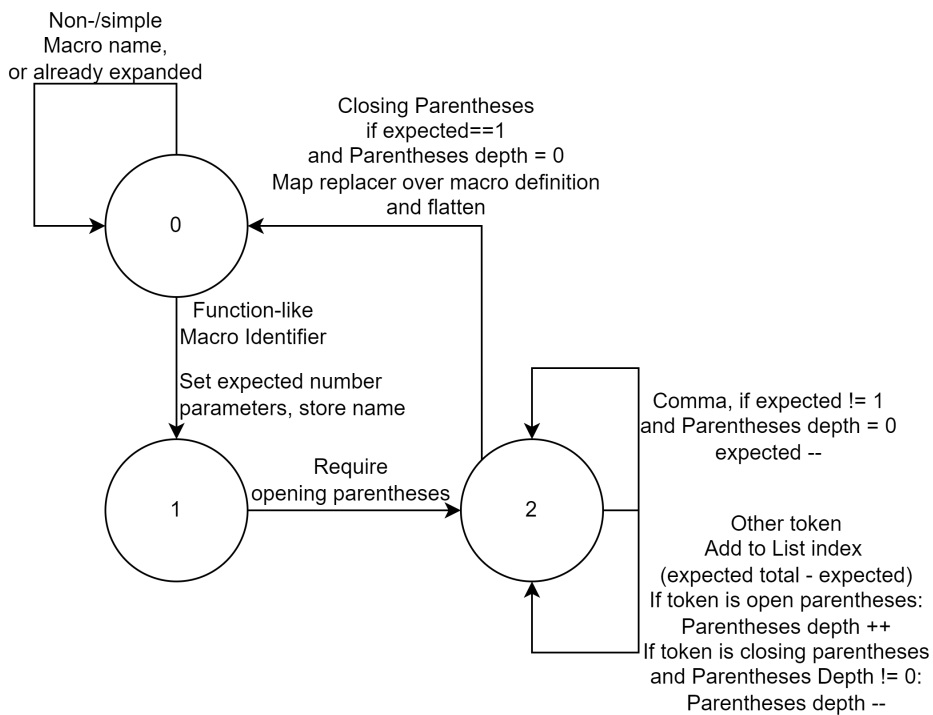


Figure 6.1: State diagram for S

that the preprocessor is in, which is used in gathering the arguments for function-like macros. Refer to g. 6.1 for a reference of what each state means.

The next integer is the number of arguments still remaining to be read in the current function-like macro, and each sub-list of the double list in the third member represents what is gathered so far for those arguments.

The fourth member, the set of *Tokens*, is there to store the macros that were already expanded in recursive expansion to avoid infinite loops.

The fifth member is the name of the function-like macro we are currently working on. The name is required to do the macro replacement, and is not present if we are not currently gathering any macro arguments.

The sixth member keeps track of the parentheses depth, as macro arguments may themselves contain opening and closing parentheses, and we would otherwise be unable to determine whether a closing parentheses is intended to be part of the parameter or to end the parameter list.

Since giving all the definitions in one big expression would be difficult to read, we split the input up into different cases: The two kinds of *MacroLines* and the case for *TokenList*, and we dispatch based on the type of line in `prepro`, together with some minor input preprocessing.

A Haskell version of all the functions in this section is also given in annexes A and B.

6.1 Wrapper and main function

The function `prepro`, with the type as given above, is only a wrapper function that sets up the arguments to the real main function, `dispatch`. This is necessary to keep the type of the function clean, but the additional information is needed for the actual preprocessing work.

We therefore define:

$$\begin{aligned} \text{prepro} &: \text{File} \rightarrow \text{File} \\ \text{prepro}(f) &= \text{dispatch}(f, \sigma, (0, 0, [], \cdot, \cdot, \cdot, 0)) \end{aligned}$$

$$\begin{aligned} \text{dispatch} &: \text{File} \quad M \quad S \rightarrow \text{File} \\ \text{dispatch}(\lambda, m, s) &= \lambda \end{aligned}$$

$$\begin{aligned} \text{dispatch}(\# \text{define } i \text{ t } m \ f, m, s) &= r^\theta \ m \ \text{dispatch}(f, m^\theta, s^\theta) \\ &\quad \mathbf{where} \ r^\theta := \text{unlistify}(r) \\ &\quad \mathbf{where} \ (m^\theta, s^\theta, r) := \text{om}(m, s, i, \text{tlisify}(t)) \end{aligned}$$

$$\begin{aligned} \text{dispatch}(\# \text{define } i \ (a) \ t \ m \ f, m, s) &= r^\theta \ m \ \text{dispatch}(f, m^\theta, s^\theta) \\ &\quad \mathbf{where} \ r^\theta := \text{unlistify}(r) \\ &\quad \mathbf{where} \ (m^\theta, s^\theta, r) := \text{fm}(m, s, i, \text{alisify}(a), \text{tlisify}(t)) \end{aligned}$$

$$\begin{aligned} \text{dispatch}(t \ m \ f, m, s) &= r^\theta \ m \ \text{dispatch}(f, m^\theta, s^\theta) \\ &\quad \mathbf{where} \ r^\theta := \text{unlistify}(r) \\ &\quad \mathbf{where} \ (m^\theta, s^\theta, r) := \text{tl}(m, s, \text{tlisify}(t)) \end{aligned}$$

$$\begin{aligned} \text{unlistify} &: [T] \rightarrow \text{TokenList} \\ \text{unlistify}([]) &= \lambda \\ \text{unlistify}(t : ts) &= t \ \text{unlistify}(ts) \end{aligned}$$

We also give a function which turns a file into list of tokens, which is a necessary step to parsing it as a program. This function is only partial,

because it is impossible to turn a *MacroLine* into a list of *Tokens*, and the function is only intended to be called on the output of *prepro*, anyways.

$$\begin{aligned} \text{rmvlines} &: \text{File} \rightarrow [T] \\ \text{rmvlines}(\lambda) &= [] \end{aligned}$$

$$\text{rmvlines}(t \text{ on } f) = \text{tlitfy}(t) + \text{rmvlines}(f) \quad \mathbf{iff} \ t \in \text{TokenList}$$

6.2 Object Macro

This case is relatively simple, as we only need to transform the *TokenList* into a proper list of tokens and store that in the state, making sure not to overwrite a previous macro.

$$\begin{aligned} \text{om} &: M \ S \ [T] \rightarrow M \ S \ [T] \\ \text{om}(m, s, i, r) &= (m[i \notin (0, r)], s, []) \quad \mathbf{iff} \ m(i) \neq m(i) \notin (0, r) \end{aligned}$$

6.3 Function-Like Macro

This case is not quite so simple, so we introduce an intermediate variable to keep the final equation legible:

$$\begin{aligned} \text{fm} &: M \ S \ [] \ [T] \rightarrow M \ S \ [T] \\ \text{fm}(m, s, i, a, r) &= (m[i \notin \text{res}], s, []) \quad \mathbf{iff} \ m(i) \neq m(i) \notin \text{res} \\ &\quad \mathbf{where} \ \text{res} := (\text{len}(a), \text{map}(\lambda t. \text{repind}(a, t), r)) \end{aligned}$$

6.4 Text Line

We handle a text line by creating a simple state machine (e.g. 6.1) and then transforming that state machine into a series of equations that describe the behaviour of the macro expansion process.

We first remove the requirement to always carry along the mapping function, as it is never changed and would unnecessarily complicate the definition even further.

$$\begin{aligned} \text{tl} &: M \ S \ [T] \rightarrow M \ S \ [T] \\ \text{tl}(m, s, t) &= (m, s^\theta, t^\theta) \quad \mathbf{where} \ (s^\theta, t^\theta) := \text{tlr}(m, s, t) \end{aligned}$$

We first define the tlr function, and then give the auxiliary functions, followed by explanations for each equation in English afterwards.

$$\begin{aligned} \text{tlr} &: M \times S \times [T] \rightarrow S \times [T] \\ \text{tlr}(m, (0, 0, [], u, ?, 0), []) &= ((0, 0, [], u, ?, 0), []) \end{aligned} \quad (1)$$

$$\begin{aligned} \text{tlr}(m, (0, 0, [], u, ?, 0), t : ts) &= (s, t : r) \\ \mathbf{where} \ (s, r) &:= \text{tlr}(m, (0, 0, [], u, ?, 0), ts) \\ &\mathbf{iff} \ t \geq u _ m(t) \end{aligned} \quad (2)$$

$$\begin{aligned} \text{tlr}(m, (0, 0, [], u, ?), t : ts) &= (s, \text{tlr}(m, (0, 0, [], \text{ftg} [u, ?, 0), e) ++ r) \\ \mathbf{where} \ (s, r) &:= \text{tlr}(m, (0, 0, [], u, ?, 0), ts) \\ &\mathbf{iff} \ t \not\geq u \wedge m(t) \# \wedge m(t) = (0, e) \end{aligned} \quad (3)$$

$$\begin{aligned} \text{tlr}(m, (0, 0, [], u, ?, 0), t : ts) &= \text{tlr}(m, (1, n, \text{repeat}([], n), u, t, 0), ts) \\ \mathbf{iff} \ t \not\geq u \wedge s(t) \# \wedge s(t) = (n, r) \wedge n \notin 0 \end{aligned} \quad (4)$$

$$\text{tlr}(m, (1, n, ls, u, t, 0), (: ts) = \text{tlr}(m, (2, n, ls, u, t, 0), ts) \quad (5)$$

$$\text{tlr}(m, (2, n, ls, u, t, 0), (: ts) = \text{tlr}(m, (2, n - 1, ls, u, t, 0), ts) \mathbf{iff} \ n \notin 1 \quad (6)$$

$$\begin{aligned} \text{tlr}(m, (2, 1, ls, u, t, 0), (: ts) &= (s, \text{expand}(s, ls, u, t) ++ r) \\ \mathbf{where} \ (s, r) &:= \text{tlr}(m, (0, 0, [], u, ?, 0), ts) \end{aligned} \quad (7)$$

$$\begin{aligned} \text{tlr}(m, (2, n, ls, u, t, d), ct : ts) &= \text{tlr}(m, (2, n, lsm, u, t, d + dd), ts) \\ \mathbf{where} \ lsm &:= ls[nt - n \mathcal{V} \text{append}(ls[nt - n], ct)], \end{aligned} \quad (8)$$

$$(nt, r) = m(t), dd = \begin{cases} 1 & \mathbf{iff} \ ct = (\\ 0 & \mathbf{otherwise} \end{cases}$$

$$\begin{aligned} \text{expand} &: M \times [[T]] \times P(T) \times T \rightarrow [T] \\ \text{expand}(m, ls, u, t) &= \text{tlr}(m, (0, 0, [], \text{ftg} [u, ?, 0], \text{atten}(\text{map}(\text{repl_int}(ls), r))) \\ &\mathbf{where} \ (n, r) := s(t) \end{aligned}$$

$$\begin{aligned} \text{atten} &: [[T]] \rightarrow [T] \\ \text{atten}([]) &= [] \\ \text{atten}([: bs) &= \text{atten}(bs) \\ \text{atten}((a : as) : bs) &= a : \text{atten}(as : bs) \end{aligned}$$

$$\begin{aligned} \text{repl_int} &: [T] \rightarrow (T \rightarrow Z) \rightarrow T \\ \text{repl_int}(ls)(i) &= ls[i] && \text{iff } i \geq N \\ \text{repl_int}(ls)(t) &= t && \text{iff } t \geq T \end{aligned}$$

Explanations

- (1) This is the easiest case: An empty line remains just that, no expansion can take place.
- (2) If we already expanded the current token we leave it alone to avoid infinite replacement loops, and if the token is not defined to have any replacement, then it also stays verbatim.
- (3) If the token is a known macro, and the first entry of the tuple is 0, then it refers to a constant macro, we then expand that macro, and possibly any further macro names that result from the expansion, and prepend that expansion to the expansion of the rest of the line.
- (4) If we encounter a macro which we know to refer to a function macro, we set up the parameters accordingly, as is mentioned in the state diagram: Store the name of the macro that we are collecting arguments for, keep track of how many arguments we still expect to find and set the state to 1 to await the opening parenthesis that must follow this macro name.
- (5) This case is also rather straight-forward, in state 1, we only want to see an opening parenthesis, anything else is not allowed, and we then process the next token(s) in state 2.
- (6) If we encounter a comma in state 2 and with no more pending open parentheses pairs, this means that we have reached the end of one argument and are about to encounter the start of another. This is not permitted if we only expected one more argument, as this would mean that we gave too many arguments to a macro.
- (7) A closing parentheses is only allowed to close the argument list once we gathered all the necessary arguments, and if there is no pending open parenthesis in the current argument. We then do the expansion, including any recursive macro expansion resulting therefrom, using the aptly named expand function, and continue on with the remainder of the line.
- (8) If we encounter any other token than a comma or a closing parentheses, then we consider it part of the macro arguments, and consequently add it to the appropriate sublist of the 5th argument. In addition, if the token is an opening parenthesis, then we keep track of that so we

know when we find its closing counterpart rather than the end of the argument list.

6.5 Auxiliary functions

Following are the definitions of the functions that were used above, but not defined.

tlistify

The function `tlistify` turns the syntactical element *TokenList* into the corresponding list of tokens.

$$\begin{aligned} \text{tlistify} &: \text{TokenList} \rightarrow [T] \\ \text{tlistify}(\lambda) &= [] \\ \text{tlistify}(t \text{ tl}) &= t : \text{tlistify}(\text{tl}) \end{aligned}$$

alistify

The function `alistify` turns an *ArgumentList* into a list of *Identifiers*, without allowing any duplicates.

$$\begin{aligned} \text{alistify} &: \text{ArgumentList} \rightarrow [] \\ \text{alistify}(a) &= \text{alistify}^\circ(a, \cdot) \\ \text{alistify}^\circ &: \text{ArgumentList} \rightarrow P(\cdot) \rightarrow [] \\ \text{alistify}^\circ(a, s) &= a : [] && \text{iff } a \notin s \\ \text{alistify}^\circ((a, al), s) &= a : \text{alistify}^\circ(al, \text{flag } [s]) && \text{iff } a \notin s \end{aligned}$$

repind

This function is mapped over a *Token* list, replacing each element that is in the input list by its index and otherwise leaves it alone.

$$\begin{aligned} \text{repind} &: [] \rightarrow T \rightarrow T [Z] \\ \text{repind}(a, b) &= \text{index}(a, b) && \text{iff } \text{index}(a, b) \neq \\ \text{repind}(a, b) &= b && \text{iff } \text{index}(a, b) = \end{aligned}$$

6.6 Example: Simple macro derivation

To aid the understanding of the process described by the functions above, we go through a simple example involving a function-like macro step-by-step. The example that we use is

```
#define X(A, B) A + B
X(2, 3)
```

or `#define X(A, B) A + B` `m X(2, 3) m` `λ`.

We start by applying `prepro` to this, and then evaluate each function in turn.

$$\begin{aligned}
 & \text{prepro}(\#define\ X(A, B)\ A + B\ m\ X(2, 3)\ m\ \lambda) \\
 = & \text{dispatch}(\#define\ X(A, B)\ A + B\ m\ X(2, 3)\ m\ \lambda, \sigma, (0, 0, [], :, ?, 0)) \\
 = & \lambda m\ \text{dispatch}(X(2, 3)\ m\ \lambda, \sigma[X\ \mathcal{D}\ (2, [1, +, 2])], (0, 0, [], :, ?, 0)) \quad (1) \\
 = & \lambda m\ 2 + 3\ m\ \text{dispatch}(\lambda, \sigma[X\ \mathcal{D}\ (2, [1, +, 2])], (0, 0, [], :, ?, 0)) \quad (2) \\
 = & \lambda m\ 2 + 3\ m\ \lambda
 \end{aligned}$$

The equal signs in the two numbered equations need some justification to show that they satisfy the equations for `dispatch`.

Regarding the first one, that is according to the third equation for `dispatch`, we need to show that

$$\text{fm}(\sigma, (0, 0, [], :, ?, 0), X, [A, B], [A, +, B]) = (\sigma[X\ \mathcal{D}\ (2, [1, +, 2])], (0, 0, [], :, ?, 0), \lambda)$$

To show this, we first calculate the value of `res` from the equation for `fm`:

$$\begin{aligned}
 \text{res} &= (\text{length}([A, B]), \text{map}(\lambda t. \text{repind}([A, B], t), [A, +B])) \\
 &= (2, [1, +, 2])
 \end{aligned}$$

(We collapse the calculation of the length and the index replacement, as that would add a lot of highly trivial steps, without aiding in understanding.)

We then see that `m(res)` " holds, so the condition is fulfilled and we can apply the equation for `fm`, and get the result written above.

We now need to justify the second numbered equals sign, which involves the fourth equation of `dispatch`. Therefore we need to show that

$$\begin{aligned}
 & \text{tl}(\sigma[X\ \mathcal{D}\ (2, [1, +, 2])], (0, 0, [], :, ?, 0), [X, (, 2, , 3,)]) = \\
 & (\sigma[X\ \mathcal{D}\ (2, [1, +, 2])], (0, 0, [], :, ?, 0), [2, +, 3])
 \end{aligned}$$

We start with $\text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (0, 0, [], :, ?, 0), [X, (, 2, , , 3,)])$, and go through the steps from there.

$$\begin{aligned}
& \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (0, 0, [], :, ?, 0), [X, (, 2, , , 3,)]) \\
= & \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (1, 2, \text{repeat}([], 2), :, X, 0), [(, 2, , , 3,)]) \\
= & \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (1, 2, [], [], :, X, 0), [(, 2, , , 3,)]) \\
= & \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (2, 2, [], [], :, X, 0), [2, , , 3,)]) \\
= & \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (2, 2, [[2], []], :, X, 0), [, , 3,)]) \\
= & \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (2, 1, [[2], []], :, X, 0), [3,)]) \\
= & \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], (2, 1, [[2], [3]], :, X, 0), []) \\
= & (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], \text{expand}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [[2], [3]], :, X)+[]) \\
= & (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [[2], [3]], (0, 0, [], fXg, ?, 0), \\
& \quad \text{atten}(\text{map}(\text{repl_int}([[2], [3]]), [1, +, 2]))) + []) \\
= & (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [[2], [3]], (0, 0, [], fXg, ?, 0), \\
& \quad \text{atten}([[2], [+], [3]])) + []) \\
= & (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [[2], [3]], (0, 0, [], fXg, ?, 0), [2, +, 3]) + []) \\
= & (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], 2 : \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [[2], [3]], (0, 0, [], fXg, ?, 0), [+ , 3]) + []) \\
= & (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], 2 : + : \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [[2], [3]], (0, 0, [], fXg, ?, 0), [3]) + []) \\
= & (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], 2 : + : 3 : \text{tlr}(\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [[2], [3]], (0, 0, [], fXg, ?, 0), []) + []) \\
& = (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], 2 : + : 3 : [] + []) \\
& = (\sigma[X \ \mathcal{D} \ (2, [1, +, 2])], [2, +, 3])
\end{aligned}$$

And from the definition of tl , which adds the m back to the tuple, we can see that we got where we needed to get.

However, this derivation also shows that the macro expansion process requires a lot of small step function applications, and so any proof involving it are best only done by automatic proof assistant systems. (You may notice that the above list of equations is not entirely precise, as I do several minor function expansions and tuple member extractions implicitly between steps, but otherwise, the derivation would be even longer, and no easier to read.)

Chapter 7

Interpretation

Once the pre-processing step has concluded, we are left with a list of tokens, which may conform, and be parseable according to, the grammar set out in section 1.3.

The parser for this is not defined explicitly, but rather, we define a function that turns a *Program* into a list of tokens, and define the parsing function as the inverse of that function, and leave it undefined if there is no *Program* that could turn into a given list of tokens, which means that the output of the pre-processing step is not a valid program.

We also give a processing and an evaluation function which together determine the value of any constant or function in the program.

A Haskell version of all the functions in this section is available in Appendix C.

7.1 Unparsing function

We write the unparsing function as $\mathit{ff}x\mathit{gg}$, where x is some syntactical element that is being turned into a list of tokens. Note that we omit the concatenation and cons operators for clarity, they are present in the Haskell

source files if one wishes to see them.

$$ff_s p \mathcal{G} = ff_s \mathcal{G} ff_p \mathcal{G}$$

$$\begin{aligned} ff_declare\ i\ (a); \mathcal{G} &= declare\ ff_i \mathcal{G}(ff_a \mathcal{G}); \\ ff_define\ i\ (a) = e; \mathcal{G} &= define\ ff_i \mathcal{G}(ff_a \mathcal{G}) = ff_e \mathcal{G}; \\ ff_const\ i = e; \mathcal{G} &= const\ ff_i \mathcal{G} = ff_e \mathcal{G}; \end{aligned}$$

$$\begin{aligned} ff(e_1 + e_2) \mathcal{G} &= (ff_{e_1} \mathcal{G} + ff_{e_2} \mathcal{G}) \\ ff(e_1 - e_2) \mathcal{G} &= (ff_{e_1} \mathcal{G} - ff_{e_2} \mathcal{G}) \\ ff(e_1 * e_2) \mathcal{G} &= (ff_{e_1} \mathcal{G} * ff_{e_2} \mathcal{G}) \\ ff(e_1 / e_2) \mathcal{G} &= (ff_{e_1} \mathcal{G} / ff_{e_2} \mathcal{G}) \end{aligned}$$

$$ff_l \mathcal{G} = l$$

$$ff_i \mathcal{G} = i$$

$$ff_i(l) \mathcal{G} = ff_i \mathcal{G}(ff_l \mathcal{G})$$

$$ff_{e,l} \mathcal{G} = ff_e \mathcal{G}, ff_l \mathcal{G}$$

$$ff_{i,a} \mathcal{G} = ff_i \mathcal{G}, ff_a \mathcal{G}$$

where the letters used for the variables indicate the type of syntactical element:

- *a* is used for *ArgumentLists*
- *e* is used for *Expressions*
- *p* is used for *Programs*
- *i* is used for *Identifiers*
- *l* is used for *ExpressionLists*
- *s* is used for *Statements*

7.2 Processing Function

We write the processing function for a single statement as $s \xrightarrow{t} s^0$ to mean that reading the statement t in state s produces a new state s^0 .

The type of this state is henceforth denoted as P , which is a shorthand notation for

$$P = \rightarrow ([], Expression?)$$

Expression is defined in section 1.1.3.

The meaning of this function is as follows: If the value of the function for an *Identifier* is not defined, it means that no constant or function declaration by that name is known.

If it is defined, and the list is empty, then the second member of the tuple must be an *Expression*, and this represents a constant definition.

Otherwise, the members of the list represent the names of function parameters, and the second member is either the *Expression* which forms the body of the function in its definition, or it is $?$, in which case the function was only declared, but not defined.

Symbolically, we therefore write the type of this function as:

$$\begin{array}{l}
 P \xrightarrow{Statement} P \\
 s \xrightarrow{\text{declare } i(a)} s \quad \text{iff } s(i) \neq \wedge \text{len}(\text{alistify}(a)) = \text{len}(s(i)_1) \\
 s \xrightarrow{\text{declare } i(a)} s[i \nabla (\text{alistify}(a), 0)] \quad \text{iff } s(i) = \\
 s \xrightarrow{\text{declare } i(a) = e} s[i \nabla (\text{alistify}(a), e)] \quad \text{iff } s(i) = _ (\text{len}(\text{alistify}(a)) = \text{len}(s(i)_0) \wedge s(i)_1 = 0) \\
 s \xrightarrow{\text{const } i = e} s[i \nabla ([], e)] \quad \text{iff } s(i) =
 \end{array}$$

Given this definition, we also define the function that takes an entire *Program* and a state to mutate and gives back the state after reading that entire *Program*. Symbolically, we write the type of this function as:

$$P \xrightarrow{Program} P$$

$$s \xrightarrow{\lambda} s$$

$$s \xrightarrow{t} s^0 \text{ where } s \xrightarrow{t} s^0 \text{ and } s^0 \xrightarrow{B} s^0$$

7.3 Evaluation

We give three distinct functions for evaluation, all of which return a natural number:

1. The first of the three functions evaluates constants, and it is denoted as $\llbracket c \rrbracket_s$, for some constant name c and state s .
2. The second function is used to evaluate functions. It is denoted as $\llbracket f \rrbracket_s^l$, for some function name f , state s and a list of natural numbers l , which represent the function parameters.
3. The third function evaluates expression, and it is denoted as $\llbracket e \rrbracket_s^p$, for some expression e , state s and a function from *Identifiers* to natural numbers to represent the parameter context (when the expression is part of the definition of a function).

These functions are only partial functions, as there is no sensible result to be given in certain instances, such as a constant that is not defined or a function with the wrong name of parameters.

$$\begin{aligned}
\llbracket c \rrbracket_s &= \llbracket s(c)_1 \rrbracket_s^\sigma \text{ iff } s(c) \neq \# \wedge s(c)_0 = [] \\
\llbracket f \rrbracket_s^l &= \llbracket e \rrbracket_s^{l^0 \otimes l} \text{ iff } s(f) \neq \# \wedge \text{len}(l^0) = \text{len}(l) \\
&\quad \text{where } (l^0, e) := s(f) \\
\llbracket (l + r) \rrbracket_s^p &= \llbracket l \rrbracket_s^p + \llbracket r \rrbracket_s^p \\
\llbracket (l - r) \rrbracket_s^p &= 0 \text{ iff } \llbracket l \rrbracket_s^p \leq \llbracket r \rrbracket_s^p \\
\llbracket (l - r) \rrbracket_s^p &= \llbracket l \rrbracket_s^p - \llbracket r \rrbracket_s^p \text{ iff } \llbracket l \rrbracket_s^p > \llbracket r \rrbracket_s^p \\
\llbracket (l \cdot r) \rrbracket_s^p &= \llbracket l \rrbracket_s^p \cdot \llbracket r \rrbracket_s^p \\
\llbracket (l/r) \rrbracket_s^p &= b\llbracket l \rrbracket_s^p / \llbracket r \rrbracket_s^p c \text{ iff } \llbracket r \rrbracket_s^p \neq 0 \\
\llbracket (l/r) \rrbracket_s^p &'' \text{ iff } \llbracket r \rrbracket_s^p = 0 \\
\llbracket il \ d \rrbracket_s^p &= 10 \cdot \llbracket il \rrbracket_s^p + \llbracket d \rrbracket_s^p \\
\llbracket 0 \rrbracket_s^p &= 0 \text{ and similarly for the other digits} \\
\llbracket i \rrbracket_s^p &= p(i) \text{ iff } p(i) \neq \# \\
\llbracket i \rrbracket_s^p &= \llbracket i \rrbracket_s \text{ iff } p(i) = \# \\
\llbracket f(l) \rrbracket_s^p &= \llbracket f \rrbracket_s^{\text{map}(\lambda a. \llbracket a \rrbracket_s^p, \bar{l})} \text{ where } \bar{e} = [e] \text{ and } \overline{e, \bar{l}} = e : \bar{l}
\end{aligned}$$

We also give the definition of an auxiliary function that is used in the definition of the evaluation function. This function takes two lists and gives back another function, which goes through the two lists in lockstep, returning the current element of the second list if the argument matches the current element of the first list, and is undefined when it reaches the end of either list.

As that is a rather abstract description, consider the example

$$f = ([1, 2, 3] \quad [1, 4, 9])$$

Then, $f(1) = 1, f(2) = 4, f(3) = 9$, and for all other inputs, f is undefined.

$$\begin{aligned}
 & : [a] \quad [b] \quad ! \quad (a \rightarrow b) \\
 (a : as) \quad (b : bs) = \lambda c. & \begin{cases} b \text{ iff } a = c \\ (as \quad bs)(f) \text{ iff } a \notin c \end{cases} \\
 & [] \quad bs = \sigma \\
 as \quad [] = \sigma
 \end{aligned}$$

Chapter 8

Correctness of the expansion process

There are some properties that we would like for the expansion process itself to have: For one, we would like all the examples in section 2.1 to give the same outputs in our model. Of course, the identifiers in those examples were arbitrary, so we would wish to generalise, e.g. the first example as:

$$\forall x, y \in T, \text{prepro}(\#define\ x\ y\ m\ x) = \lambda\ m\ y$$

While it would certainly be possible to write this down in some proof checking system, such as Coq, proving that this holds is merely a matter of repeatedly applying function definitions. As such, this kind of proof is not particularly interesting, although it does serve as a good sanity check, especially to diagnose any errors if such are present. Also, doing this for all the examples would add at least 15 pages of dense function derivations, and the chance for error is far too high to do so by hand.

Far more interesting are properties of the functions themselves. For example, we would expect that applying `prepro` to a file more than once should not change anything, i.e. `prepro` is idempotent, as there are no macros left in the output. Formally, we would write this as

$$\forall f \in File, \text{prepro}(\text{prepro}(f)) = \text{prepro}(f)$$

Lemma 8.0.1. *Every Line in the output of the prepro function is a TokenList.*

Proof.

1. As the `prepro` function is defined in terms of dispatch, we prove that every *Line* in the output of the dispatch function is a *TokenList*.
2. We prove this by induction on the first argument of dispatch.

3. Base case: The first argument of the dispatch function is λ .
4. Then the output is also λ . Therefore, there are no lines in the output, and the lemma holds vacuously.
5. Inductive case 1: The first argument of the dispatch function is a *TokenList* l followed by a *File* f .
6. Inductive Hypothesis: Every *Line* in $\text{dispatch}(f, m, s)$ is a *TokenList*, for all m and s .
7. To show: Every *Line* in $\text{dispatch}(l \text{ } m \text{ } f, m, s)$ is a *TokenList*, for all m and s .
8. $\text{dispatch}(l \text{ } m \text{ } f, m, s) = r^\theta \text{ } m \text{ } \text{dispatch}(f, m^\theta, s^\theta)$
where $r^\theta := \text{unlistify}(r)$
where $(r, m^\theta, s^\theta) := \text{tl}(m, s, \text{tlisify}(t))$
9. According to the definition of unlistify , r^θ is a *TokenList*.
10. According to step 6, every line in $\text{dispatch}(f, m^\theta, s^\theta)$ is a *TokenList*.
11. From 9 and 10 it follows that 7 holds.
12. Inductive cases 2 and 3: If the first argument is either type of *Macro-Line* followed by a file, steps 6-11 hold by the same construction.
13. From 4 and 11, it follows by induction that the claim in 1 holds.
14. Therefore, as the output of dispatch is precisely the output of prepro , it follows that the lemma holds.

□

Lemma 8.0.2. $\forall t \in [T], \text{tlr}(\sigma, (0, 0, [], ;, 1, 0), t) = (((0, 0, [], ;, 1, 0)), t)$

Proof.

1. Proof by induction on t .
2. Base case: $t = []$.
3. By definition of tlr in equation 5.1, $\text{tlr}(\sigma, (0, 0, [], ;, 1, 0), []) = ((0, 0, [], ;, 1, 0), [])$. Therefore, the lemma holds in the base case.
4. Inductive case: $t = e : l$ for some *Token* e and some *TokenList* l .
5. Inductive hypothesis: $\text{tlr}(\sigma, (0, 0, [], ;, 1, 0), l) = ((0, 0, [], ;, 1, 0), l)$
6. Then, $\text{tlr}(\sigma, (0, 0, [], ;, 1, 0), e : l) = ((0, 0, [], ;, 1, 0), e : l)$ by using the induction hypothesis in equation 5.2, which applies because $\sigma(e) = \text{}$. Therefore, the inductive case holds.

7. From 3 and 6 it follows that the lemma holds.

□

Lemma 8.0.3. *8t 2 TokenList, unlistify(tlistify(t)) = t*

Proof.

1. Proof by induction on tl .
2. Base case: $tl = \lambda$.
3. By de nition of tlistify and unlistify, $unlistify(tlistify(\lambda)) = \lambda$, so the lemma holds in this case.
4. Inductive case: $tl = tr$ for some *Token* t and some *TokenList* l .
5. Induction hypothesis: $unlistify(tlistify(r)) = r$
6. Then, $unlistify(tlistify(tr)) = unlistify(t : tlistify(r))$ by de nition of tlistify.
7. We apply the de nition of unlistify, and get $unlistify(t : tlistify(r)) = tunlistify(tlistify(r))$.
8. We apply the induction hypothesis to that and get $tunlistify(tlistify(r)) = tr$.
9. Therefore, the lemma holds in the recursive case.
10. From 3 and 9, the lemma holds in all cases by induction.

□

Theorem 8.0.4. *8f 2 File, prepro(prepro(f)) = prepro(f)*

Proof. For brevity, let g refer to $prepro(f)$ for the rest of this proof.

1. Proof by induction on the structure of g .
2. Base case: $g = \lambda$. Then, by the de nitions of prepro and the rst de nition of dispatch, $prepro(g) = \lambda$, and so the theorem holds.
3. Induction hypothesis: Let $g = l \ m \ h$. Then, $prepro(h) = h$
4. Inductive case 1: g is a *MacroLine* (either type) followed by another *File*.
5. In this case, as g is an output of prepro, we have a contradiction with lemma 8.0.1, therefore, the Theorem holds in this case.
6. Inductive case 2: g is a *TokenList* followed by another *File*.

7. Then, by the definitions of prepro and the last definition of dispatch , and by lemma 8.0.2, we have that $\text{prepro}(g) = \text{dispatch}(l \text{ } m \text{ } h, \sigma, (0, 0, [], ;, 1, 0)) = \text{unlistify}(\text{tlistify}(l)) \text{ } m \text{ } \text{dispatch}(h, \sigma, (0, 0, [], ;, 1, 0))$.
8. By applying the definition of prepro backwards, we rewrite that as $\text{prepro}(g) = \text{unlistify}(\text{tlistify}(l)) \text{ } m \text{ } \text{prepro}(h)$.
9. We then apply lemma 8.0.3 and get $\text{prepro}(g) = l \text{ } m \text{ } h = g$, so the Theorem holds in the inductive case.
10. From 2, 6 and 10, we conclude by induction that the theorem holds in all cases.

□

Result: Each line stays the same when applying prepro again, by induction on lines. Note: IH must include both the equality of the lines and that the state and mapping are default/empty. Just the equality would not do.

Chapter 9

Predicates over the state(s)

We can now come back to our second research question of proving the correctness of C files with preprocessing macros.

We express the correctness requirements on a file with predicates over that file, such as "\The file shall define a constant or constant macro called CHAR_BIT with a value of at least 8." (cf. [8], 5.2.4.2.1) This cannot simply be a predicate over any one state, for two reasons: For one, the macro and the program constants are kept in two separate states. And secondly, we wish to make sure that the definition is only added in the file, not previously present.

To be able to define such predicates properly, we need to define a function that gives us access to the preprocessing mappings after the file has been preprocessed. We therefore define this function `preprol` as follows. Note that this is exactly the same as `prepro` in section 6.1, except that it passes the mapping back as well.

$$\begin{aligned} \text{prepro}^0 &: File \rightarrow File \quad M \\ \text{prepro}^0(f) &= \text{dispatch}^0(f, \sigma, (0, 0, [], :, 1, 0)) \end{aligned}$$

$$\begin{aligned} \text{dispatch}^0 &: File \quad M \quad S \rightarrow File \quad M \\ \text{dispatch}^0(\lambda, m, s) &= (\lambda, m) \end{aligned}$$

$$\begin{aligned} \text{dispatch}^0(\#define \ i \ t \ m \ f, m, s) &= (r^0 \ m \ r^{00}, m^{00}) \\ \text{where } (r^{00}, m^{00}) &:= \text{dispatch}^0(f, m^0, s^0) \\ \text{where } r^0 &:= \text{unlistify}(r) \\ \text{where } (m^0, s^0, r) &:= \text{om}(m, s, i, \text{tlisify}(t)) \end{aligned}$$

$$\begin{aligned} \text{dispatch}^0(\#define \ i \ (\ a \) \ t \ m \ f, m, s) &= (r^0 \ m \ r^{00}, m^{00}) \\ \text{where } (r^{00}, m^{00}) &:= \text{dispatch}^0(f, m^0, s^0) \\ \text{where } r^0 &:= \text{unlistify}(r) \\ \text{where } (m^0, s^0, r) &:= \text{fm}(m, s, i, \text{alisify}(a), \text{tlisify}(t)) \end{aligned}$$

$$\begin{aligned} \text{dispatch}^0(t \ m \ f, m, s) &= (r^0 \ m \ r^{00}, m^{00}) \\ \text{where } (r^{00}, m^{00}) &:= \text{dispatch}^0(f, m^0, s^0) \\ \text{where } r^0 &:= \text{unlistify}(r) \\ \text{where } (m^0, s^0, r) &:= \text{tl}(m, s, \text{tlisify}(t)) \end{aligned}$$

Then, we need to consider four separate states, the *Program* that would result from it's expansion, and the intermediate list of *Tokens* to express the above predicate on some given *File* f :

$$\begin{aligned} \exists m, m^0 \in M, \exists s, s^0 \in P, \exists p \in Program, \exists t \in [T], \\ ((\text{prepro}'(f, m)_0 = t \wedge \text{ffpgg} = t \wedge \text{prepro}'(f, m)_1 = m^0 \wedge s \xrightarrow{R} s^0 \\ \wedge \text{expand}(m, [[\text{CHAR_BI T}]], :, \text{CHAR_BI T}) \wedge s(\text{CHAR_BI T}) =) \\ ([\text{expand}(m^0, [[\text{CHAR_BI T}]], :, \text{CHAR_BI T})]_{s^0}^{\sigma} \ \& \\ _ [[\text{CHAR_BI T}]]_{s^0} \ \&)) \end{aligned}$$

As it can be seen that even this relatively simple predicate leads to a rather complex formula, and writing such a formula out fully should be done only as rarely as possible to avoid errors, I therefore introduce various higher-order functions that produce predicates over a *File*, such as the one above.

We then write this higher-order function like this:
(The name is an abbreviation for Macro Constant or Constant Predicate.)

$$\begin{aligned} \text{MCCP} &: ((N \rightarrow B)) ! (File \rightarrow B) \\ \text{MCCP}(i, f) &= \delta m, m^\theta \in M, \delta s, s^\theta \in P, \delta p \in Program, \delta t \in [T], \\ ((\text{prepro}'(f, m)_0 = t \wedge \text{ffp} \text{gg} = t \wedge \text{prepro}'(f, m)_1 = m^\theta \wedge s \xrightarrow{\theta} s^\theta \\ &\wedge \text{expand}(m, [[i]], ;, i) " \wedge s(i) ") =) \\ &(f(\llbracket \text{expand}(m^\theta, [[i]], ;, i) \rrbracket_{s^\theta}^\sigma) - f(\llbracket i \rrbracket_{s^\theta})) \end{aligned}$$

And then we can write the above example as

$$\text{MCCP}(\text{CHAR_BIT}, \lambda v. v \quad 8)$$

However, other requirements may set multiple constants in relation to each other. For example, it may be required that $\text{UCHAR_MAX} = 2^{\text{CHAR_BIT}} - 1$ must hold (cf. [8], 5.2.4.2.1). For this, we define another higher-order function, which also provides a function that maps *Identifiers* to natural numbers:

$$\begin{aligned} \text{MCCAP} &: ((N \rightarrow (N \rightarrow N) \rightarrow B)) ! (File \rightarrow B) \\ \text{MCCAP}(i, f) &= \delta m, m^\theta \in M, \delta s, s^\theta \in P, \delta p \in Program, \delta t \in [T], \\ ((\text{prepro}'(f, m)_0 = t \wedge \text{ffp} \text{gg} = t \wedge \text{prepro}'(f, m)_1 = m^\theta \wedge s \xrightarrow{\theta} s^\theta \\ &\wedge \text{expand}(m, [[i]], ;, i) " \wedge s(i) ") =) \\ &(f(\llbracket \text{expand}(m^\theta, [[i]], ;, i) \rrbracket_{s^\theta}^\sigma, v(m^\theta, s^\theta)) - f(\llbracket i \rrbracket_{s^\theta}, v(m^\theta, s^\theta)))) \end{aligned}$$

$$\begin{aligned} &v : M \rightarrow P \rightarrow (N \rightarrow N) \\ v(m, s) \lambda i. &\begin{cases} \llbracket \text{expand}(m, [[i]], ;, i) \rrbracket_s^\sigma & \text{iff } m(i) \# \\ \llbracket i \rrbracket_s & \text{iff } m(i) " \end{cases} \end{aligned}$$

And we can then write the above requirement as

$$\text{MCCAP}(\text{UCHAR_MAX}, \lambda v. f. v = 2^{f(\text{CHAR_BIT})} - 1)$$

Chapter 10

Discussion

10.1 Expansion to other directives

The C standard specifies a number of other preprocessing directives besides macro expansion, such as conditional directives (`#if`, `#ifdef`, ...) and source file inclusion, among others. (For a full list see [8], 6.10.) In addition, there are also the `#` and `##` operators in macro replacement lists, which were also not considered in this thesis. However, to model the full C preprocessor, all of this needs to be included as well. Therefore, I will now briefly discuss how each of the directives and features could be added to the model, without giving any explicit equations for this.

- For the conditional directives, one needs to define the evaluation function for the conditions, which is trivial if one has the syntax tree corresponding to the constant expression. In addition, one must also keep track of how many nested if groups one is in, both those that should be skipped and those that should not.
This could be achieved by adding two more integers to the state that is carried between lines, one for each kind of nesting (which I shall call s and c), such that if one is currently in a group that should be skipped ($s > 0$), all `#if` and `#ifdef` directives just increase s by one, and all `#endif` directives decrease it by one.
And if the current group should not be skipped ($s = 0$), then any `#if` and `#ifdef` directives increase either s or c by one, depending on whether the condition is true, and all `#endif` directives decrease c by one. This is only a rough outline, but it should serve as a starting point for anyone who wishes to expand this model.
- Regarding file inclusion, the standard only prescribes that a file name might need to undergo macro expansion before it can be included, but leaves everything else implementation-defined. This means that in a

model, we can only fall back to a "magic function" that maps the filename onto a *File*, and act as though the lines of that *File* had been in our original *File* instead of the inclusion directive.

- The `#` operator is easily added to the model by modifying the type of the state M to also allow negative numbers to represent parameters that should be stringified, and adapting the `repl.int` function to perform this process, but that is not a particularly large effort, though it requires string literals to be a meaningful concept, which is why it was not done in this thesis. The `##` operator is also not difficult to model, the largest difficulty here will be in determining whether something is a valid preprocessing token.
- The `#undef` directive is trivial to implement as a modification of the macro mapping state, it was not done in this thesis solely to keep the model as simple as possible.
- Keeping track of line numbers is possible by adding another integer to the state passed between lines that is incremented each time, and a special case should be added to the definition of `tlr` to make sure that the macro `__LINE__` always expands to the current line number. Then, the implementation of the line control directive is merely to change this additional integer as indicated by the argument to the directive.
- As regards the `error`, `pragma` and `null` directives, they all likely should simply be treated as no-ops, as we are interested in the behaviour of the result of preprocessing, and these directives do not alter it, except for the floating point directives, which must be somehow preserved into the C program (as they follow its block structure), presumably using some sort of non-syntactical marker that is left in place of the `pragma` directive.

10.2 Expansion to C

As mentioned in the introduction, the preprocessor is agnostic with regard to its target language, up to the definition of what constitutes a token. As such, expanding the model of the preprocessor to target C merely requires expanding this part of the grammar and does not require any change to the definitions.

On the other hand, the toy language used in this thesis is rather far removed from C, and so one may expect the predicates to test file correctness are not very useful anymore after the transition. However this is only somewhat correct: While the evaluation function will certainly become far more complex, and will also need to contend with the type system, the general model of those predicates is still sound as long as the outputs of a function depend

only on its inputs.

This is not the case for all standard functions, such as `strtok`, but it is true of most functions in the standard library, and in any case, the result type of the evaluation function can be modified to also include a new program state if an evaluation does change the state.

In conclusion, while the transition from our toy language to C as a target language does certainly require some changes, I am confident that the basic structure and model presented in this thesis will form a solid foundation to build on for any future work in this direction.

Bibliography

- [1] TIOBE Software popularity index. <https://www.tiobe.com/tiobe-index/>.
- [2] Sandrine Blazy and Xavier Leroy. Mechanized semantics for the clight subset of the c language. *Journal of Automated Reasoning*, 43(3):263{288, 2009.
- [3] Andrew W. Appel and Sandrine Blazy. Separation logic for small-step cminor. In Klaus Schneider and Jens Brandt, editors, *Theorem Proving in Higher Order Logics*, pages 5{21, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg.
- [4] Valery A Nepomniaschy, Igor S Anureev, and AV Promskii. Towards verification of c programs: Axiomatic semantics of the c-kernel language. *Programming and Computer Software*, 29(6):338{350, 2003.
- [5] Xavier Leroy and Sandrine Blazy. Formal verification of a c-like memory model and its uses for verifying program transformations. *Journal of Automated Reasoning*, 41(1):1{31, 2008.
- [6] Robbert Jan Krebbers. *The C standard formalized in Coq*. PhD thesis, Radboud Universiteit Nijmegen, 2015.
- [7] Christian Kästner, Paolo G. Giarrusso, Tillmann Rendel, Sebastian Erdweg, Klaus Ostermann, and Thorsten Berger. Variability-Aware Parsing in the Presence of Lexical Macros and Conditional Compilation. *SIGPLAN Not.*, 46(10):805{824, oct 2011.
- [8] ISO/IEC. C standard working draft N2731 (likely to be formalised as ISO/IEC 9899:2023). <http://open-std.org/jtc1/sc22/wg14/www/docs/n2731.pdf>.
- [9] University of Glasgow. *Haskell2010 Reference, Data.Maybe*, 2010. <https://hackage.haskell.org/package/base-4.16.0.0/docs/Data-Maybe.html>.

- [10] Oracle Corporation. *Java SE 17 & JDK 17 Reference*, *java.util.Optional<T>*, 2021. <https://docs.oracle.com/en/java/javase/17/docs/api/java.base/java/util/Optional.html>.

Appendix A

Common.hs

```
1  module Common where
   import qualified Data.Set as Set
   5  type Token = [Char]
      type Identifier = [Char]
      data TokOrInt = Token Token | Int Int deriving (Eq, Show)
      data TokenList = Lambda | C Token TokenList
      data ArgumentList = Arg Identifier | CA Identifier ArgumentList
   10
      instance Show ArgumentList where
          show (Arg a) = a
          show (CA a al) = a ++ ", " ++ show al
   15
      alistify :: ArgumentList -> [Token]
      alistify a = alistify' a Set.empty
          where
              alistify' (Arg a) s = [a]
              alistify' (CA a al) s =
   20                 if Set.member a s
                    then error "Duplicate Argument name"
                    else a : alistify' al (Set.insert a s)
```

Appendix B

PreProcess.hs

```
1  module PreProcess (Line (OM, FM, TL), emptyState) where
    import Common
    import Data.List (elemIndex)
5   import Data.Maybe (fromJust, isJust, isNothing)
    import qualified Data.Set as Set

    data Line = OM Token TokenList | FM Token ArgumentList TokenList |
              TL TokenList

10  data File = FileEmpty | FileCons Line File

    data StateInt = ZERO | ONE | TWO

    type Mapping = Token -> Maybe (Int, [TokOrInt])

15  type State = (StateInt, Int, [[Token]], Set.Set Token, Token, Int)

    prepro :: File -> File
    prepro f = dispatch f emptyState (ZERO, 0, [], Set.empty, "1", 0)

20  rmvlines :: File -> [Token]
    rmvlines FileEmpty = []
    rmvlines (FileCons (TL t) f) = (tlistify t) ++ (rmvlines f)

25  dispatch :: File -> Mapping -> State -> File
    dispatch FileEmpty m s = FileEmpty
    dispatch (FileCons (OM i t) f) m s = FileCons (TL (foldr C Lambda
        r)) (dispatch f m' s')
        where
            (m', s', r) = om m s i (tlistify t)
30  dispatch (FileCons (FM i a t) f) m s = FileCons (TL (foldr C
        Lambda r)) (dispatch f m' s')
        where
```

```

    (m', s', r) = fm m s i (alistify a) (tlistify t)
dispatch (FileCons (TL t) f) m s = FileCons (TL (foldr C Lambda r)
) (dispatch f m' s')
  where
35   (m', s', r) = tl m s (tlistify t)

om :: Mapping -> State -> Token -> [Token] -> (Mapping, State, [
  Token])
om m s i l = (modMapping m i (0, map Token (l)), s, [])

40 fm :: Mapping -> State -> Token -> [Identifier] -> [Token] -> (
  Mapping, State, [Token])
fm m s i a l = (modMapping m i res, s, [])
  where
    res = (length a, map (repind a) (l))

45 tl :: Mapping -> State -> [Token] -> (Mapping, State, [Token])
tl m s t = (m, s', t') where (s', t') = tlr m s t

tlr :: Mapping -> State -> [Token] -> (State, [Token])
tlr m s@(ZERO, 0, [], u, "1", 0) [] = (s, [])
50 tlr m s@(ZERO, 0, [], u, "1", 0) (t : ts) =
  if Set.member t u || isNothing (m t)
  then (s, t : snd (tlr m s ts))
  else
    if n == 0
55     then (fst (tlr m s ts), snd (tlr m (ZERO, 0, [], Set.
insert t u, "1", 0) r) ++ snd (tlr m s ts))
    else tlr m (ONE, n, replicate n [], u, t, 0) ts
  where
    r = filterTokens r' --if n==0, this will be the same list, but
      this satisfies the types.
    Just (n, r') = m t
60 tlr m (ONE, n, ls, u, t, 0) ("(" : ts) = tlr m (TWO, n, ls, u, t,
0) ts
tlr m (ONE, n, ls, u, t, 0) ts = (s', t : ts') where (s', ts') =
  tlr m (ZERO, 0, [], Set.empty, "1", 0) ts
tlr m (TWO, 1, ls, u, t, 0) ("," : ts) = error "Comma to end
argument list?"
tlr m (TWO, n, ls, u, t, 0) ("," : ts) = tlr m (TWO, n - 1, ls, u,
t, 0) ts
tlr m (TWO, 1, ls, u, t, 0) (")" : ts) = (s', (expand m ls u t) ++
  ts') where (s', ts') = (tlr m (ZERO, 0, [], u, "1", 0) ts)
65 tlr m (TWO, n, ls, u, t, d) (ct : ts) = tlr m (TWO, n, lsm, u, t,
  if ct == "(" then d + 1 else d) ts
  where
    lsm = replaceIndex ls (nt - n) ((ls !! (nt - n)) ++ [ct])
    Just (nt, r) = m t

```

```

70 flatten :: [[a]] -> [a]
   flatten [] = []
   flatten ([] : bs) = flatten bs
   flatten ((a : as) : bs) = a : flatten (as : bs)

75 replInt :: [[Token]] -> TokOrInt -> [Token]
   replInt ls (Int i) = ls !! i
   replInt ls (Token t) = [t]

   expand :: Mapping -> [[Token]] -> Set.Set Token -> Token -> [Token]
   ]
80 expand m ls u t = snd (tlr m (ZERO, 0, [], Set.insert t u, "1", 0)
   (flatten (map (replInt ls) r)))
   where
   Just (n, r) = m t

85 replacelndex :: [a] -> Int -> a -> [a]
   replacelndex (a : as) 0 r = r : as
   replacelndex (a : as) n r = a : replacelndex as (n - 1) r

   filterTokens :: [TokOrInt] -> [Token]
   filterTokens [] = []
90 filterTokens ((Token t) : ts) = t : filterTokens ts

   tlistify :: TokenList -> [Token]
   tlistify Lambda = []
   tlistify (C t tl) = t : tlistify tl

95 repind :: [Token] -> Token -> TokOrInt
   repind a b =
   if isJust (elemIndex b a)
   then Int (fromJust (elemIndex b a))
100 else Token b

   emptyState :: Mapping
   emptyState t = Nothing

105 modMapping :: Mapping -> Token -> (Int, [TokOrInt]) -> Mapping
   modMapping m t r t' =
   if isJust (m t) && m t /= Just r
   then error "Redefining existing macro!"
   else if t == t' then Just r else m t'

```


Appendix C

PostProcess.hs

```
1  module PostProcess (eval, File (..), Statement (..), Expression
    (..), ExpressionList (..), BinOp (..)) where

import Common
import Data.Map ((!))
5  import qualified Data.Map as Map
import Data.Maybe (isJust)

type PState = Map.Map Identifier ([Identifier], Maybe Expression)

10 data File = FS Statement | FC Statement File

data Statement = FDecl Identifier ArgumentList | FDefn Identifier
    ArgumentList Expression | CDefn Identifier Expression

data Expression = LIT Integer | IDENT Identifier | FUNC Identifier
    ExpressionList | BI Expression BinOp Expression
15 data ExpressionList = ELS Expression | ELC Expression
    ExpressionList

data BinOp = PLUS | MINUS | TIMES | OVER

20 instance Show File where
    show (FS s) = show s
    show (FC s f) = show s ++ "\n" ++ show f

instance Show Statement where
25 show (FDecl i al) = "declare " ++ i ++ "(" ++ show al ++ ";"
    show (FDefn i al e) = "define " ++ i ++ "(" ++ show al ++ ") = "
        ++ show e ++ ";"
    show (CDefn i e) = "const " ++ i ++ " = " ++ show e ++ ";"

instance Show Expression where
```

```

30 show (LIT i) = show i
   show (IDENT i) = i
   show (FUNC i el) = i ++ "(" ++ show el ++ ")"
   show (BI | PLUS r) = "(" ++ show l ++ "+" ++ show r ++ ")"
   show (BI | MINUS r) = "(" ++ show l ++ "-" ++ show r ++ ")"
35 show (BI | TIMES r) = "(" ++ show l ++ "*" ++ show r ++ ")"
   show (BI | OVER r) = "(" ++ show l ++ "/" ++ show r ++ ")"

instance Show ExpressionList where
   show (ELS e) = show e
40   show (ELC e el) = show e ++ ", " ++ show el

class Unparseable u where
   unparse :: u -> [Token]

45 instance Unparseable File where
   unparse (FS s) = unparse s
   unparse (FC s f) = unparse s ++ unparse f

instance Unparseable Statement where
50   unparse (FDecl i al) = ["declare", i, "(" ++ unparse al ++ "]"
      , ";" ]
   unparse (FDefn i al e) = ["define", i, "(" ++ unparse al ++ "]"
      , "=" ++ unparse e ++ ";" ]
   unparse (CDefn i e) = ["const", i, "=" ++ unparse e ++ ";" ]

instance Unparseable Expression where
55   unparse (LIT i) = [show i]
   unparse (IDENT i) = [i]
   unparse (FUNC i el) = [i, "(" ++ unparse el ++ "]" ]
   unparse (BI | o r) = "(" : unparse l ++ unparse o ++ unparse r
      ++ "]"

60 instance Unparseable ExpressionList where
   unparse (ELS e) = unparse e
   unparse (ELC e el) = unparse e ++ [", "] ++ unparse el

instance Unparseable BinOp where
65   unparse PLUS = ["+"]
   unparse MINUS = ["-"]
   unparse TIMES = ["*"]
   unparse OVER = ["/"]

70 instance Unparseable ArgumentList where
   unparse (Arg a) = [a]
   unparse (CA a al) = [a, ", "] ++ unparse al

eval :: File -> Integer
75 eval f = evalC "main" (processF f Map.empty)

```

```

processF :: File -> PState -> PState
processF (FC s f) os = processF f (processS s os)
processF (FS s) os = processS s os
80
processS :: Statement -> PState -> PState
processS (CDefn i e) os =
  if Map.member i os
    then error ("Attempted to redefine constant: " ++ i ++ "!")
85    else Map.insert i ([], Just e) os
processS (FDecl i a) os =
  if Map.member i os
    then
90      if length al == length (fst (os ! i))
        then os
        else error ("Function declaration with wrong number of
parameters: " ++ i ++ "!")
      else Map.insert i (al, Nothing) os
  where
    al = alistify a
95 processS (FDefn i a e) os =
  if Map.member i os
    then
100      if length al == length (fst (os ! i))
        then
          if isJust (snd (os ! i))
            then error ("Tried to redefine function" ++ i ++ "")
            else Map.insert i (al, Just e) os
          else error ("Function definition with wrong number of
parameters: " ++ i ++ "!")
        else Map.insert i (al, Just e) os
105      where
        al = alistify a

evalC :: Identifier -> PState -> Integer
evalC i s =
110  if Map.member i s
    then
      if null is
        then evalE e s Map.empty
        else error ("Tried to evaluate constant as function: " ++
i ++ "!")
115    else error ("Tried to evaluate unknown constant: " ++ i ++ "!")
  )
  where
    (is, Just e) = s ! i

evalF :: Identifier -> PState -> [Integer] -> Integer
120 evalF i s a =

```

```

    if Map.member i s
    then
      if length is == length a
      then evalE e s (Map.fromList (zip is a))
125      else error ("Tried to evaluate function with wrong number
of arguments: " ++ i ++ "!")
      else error ("Tried to evaluate unknown function: " ++ i ++ "!")
    )
  where
    (is, Just e) = s ! i

130 evalE :: Expression -> PState -> Map.Map Identifier Integer ->
Integer
evalE (BI | PLUS r) s m = evalE l s m + evalE r s m
evalE (BI | MINUS r) s m = max 0 (evalE l s m - evalE r s m)
evalE (BI | TIMES r) s m = evalE l s m * evalE r s m
evalE (BI | OVER r) s m = evalE l s m `div` evalE r s m
135 evalE (LIT i) _ _ = i
evalE (IDENT i) s m =
  if Map.member i m
  then m ! i
  else evalC i s
140 evalE (FUNC i es) s m = evalF i s (map (\e -> evalE e s m) (
  elistify es))

elistify :: ExpressionList -> [Expression]
elistify (ELS e) = [e]
elistify (ELC e es) = e : elistify es

```