

The structure of a compiler

- Trees
- Running a compiler
- Trees and ATerms
- Pretty printing
- Trees and ATerms



Compiler Construction

WWW: <http://www.cs.uu.nl/wiki/Cco>

Edition 2011/2012

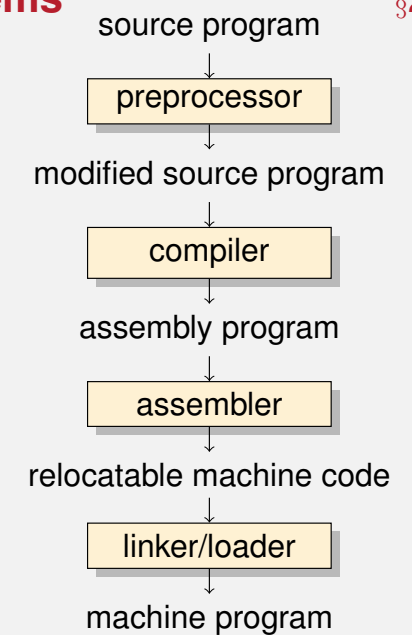


2. The structure of a compiler

Language-processing systems §2

Typically, compilers are parts of larger systems.

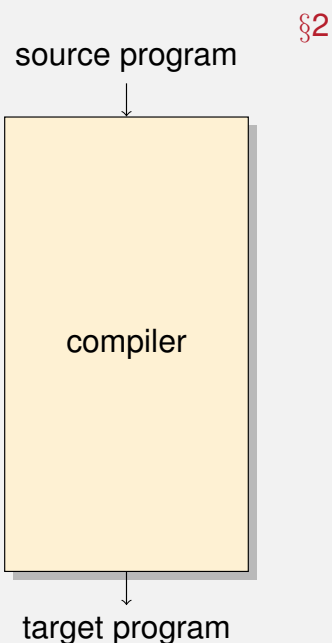
For example, for the generation of machine-executable code from a source program written in a high-level programming language, besides a compiler, several other programs may be involved: e.g., a preprocessor, an assembler, and/or a linker.



Black-box view

As a whole, a compiler performs a meaning-preserving translation from its source language into its target language.

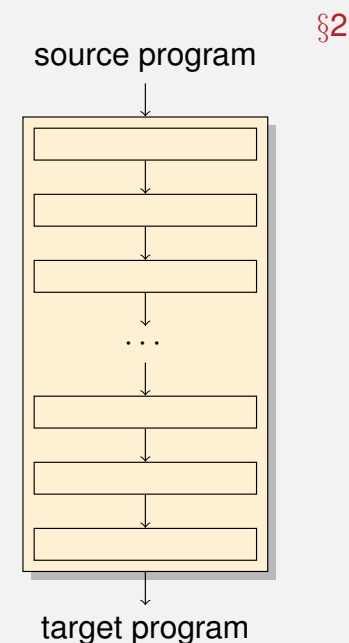
Optionally, it also validates the source program, reporting (possible) errors to the user, and/or produces a log of its activities.



White-box view

Internally, a compiler is typically implemented as a **pipeline** of components.

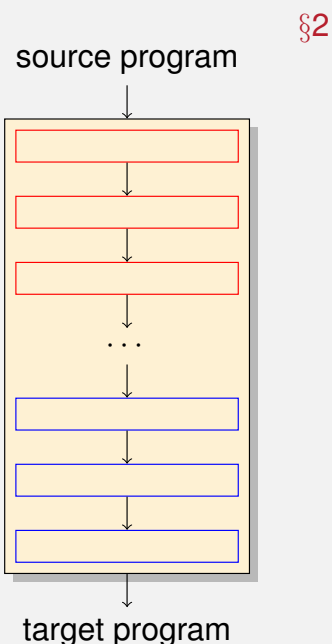
Each component in the pipeline takes as input the output of its predecessor and produces as output the input for its successor.



Analysis and synthesis

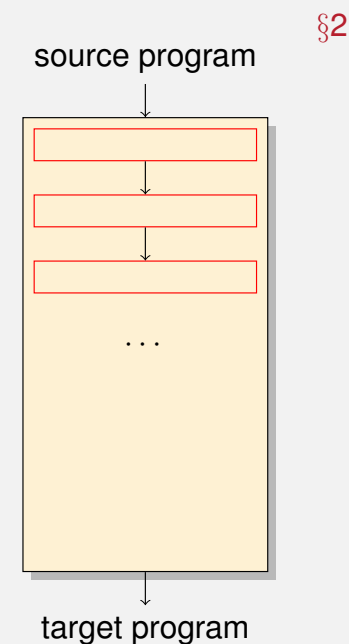
The first components in the pipeline form the **front end** of the compiler and perform the **analysis** of the source program.

The last components form the **back end** and take care of the **synthesis** of the target program.



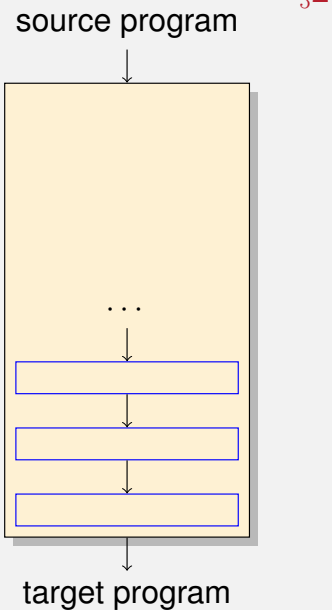
Front end

- ▶ Breaks up the source program into its constituent pieces.
- ▶ Imposes a grammatical structure on it.
- ▶ Checks for syntactic and semantic errors.
- ▶ Produces informative warnings and error messages.
- ▶ Constructs a symbol table.
- ▶ Creates an intermediate representation.



Back end

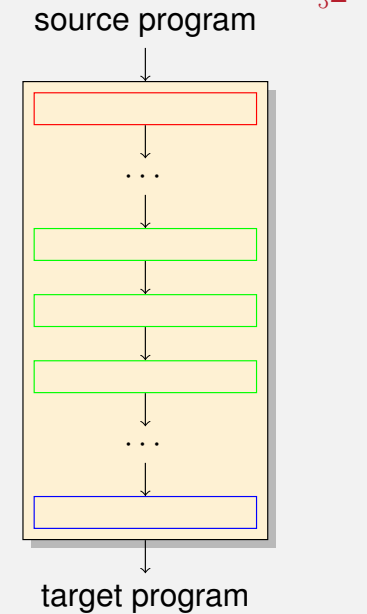
- ▶ Constructs the target program from the intermediate representation and the information stored in the symbol table.



Middle end

Sometimes we separate out a **middle end** in which optimisations on intermediate representations of the program have place.

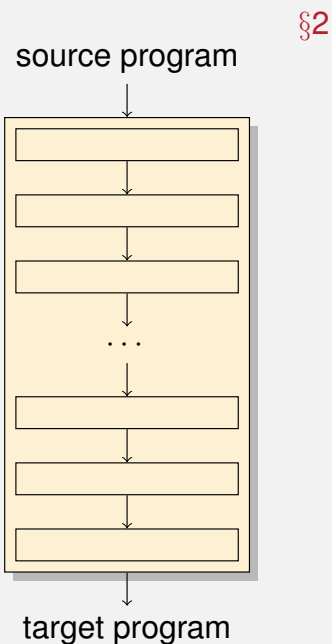
Ideally, the middle end operates on an intermediate representation that is completely independent from the source and target languages. This way, we can easily combine different front and back ends.



Phases

Logically, the compilation pipeline consists of a sequence of **phases**. In each phase, one representation of the source program is transformed into another.

The logical division of the compilation process into phases is not necessarily matched by the actual implementation. In practice, several phases are grouped together or split into subphases. Furthermore, some intermediate representations may not be constructed explicitly.



Phases: example

Typical phases in the compilation of a high-level program into machine-executable code are:

- ▶ Lexical analysis.
- ▶ Syntactic analysis.
- ▶ Semantic analysis.
- ▶ Intermediate-code generation.
- ▶ Intermediate-code optimisation.
- ▶ Code generation.



Symbol tables are data structures that contain incrementally obtained information about the source program. This information is typically produced in the front end of the compiler and consumed in the back end.

For example, the symbol table may contain information about identifiers such as their position in the source program and their type.

The symbol table may be shared by the different components that make up a compiler. Alternatively, it can be passed around by the compiler components.

In practice, the symbol table is not really a single table, but a compound data structure consisting of multiple tables, that may each contain subtables and so on.



Internally, the intermediate representations of the source program typically take the form of **trees**.

Interacting with their environments, however, compilers consume and produce “flat” **data streams**, i.e., strings of characters or bytes.

Implementing a compiler, we therefore have to be able to convert between flat data and trees, and to transform one type of trees into another.



From a software engineer’s point of view it may be advantageous to construct the components of a compiler as loosely coupled as possible.

Such modularity keeps the overall design comprehensible and makes it possible to debug and test components in isolation.

Ultimately, each component can exist as a stand-alone executable program, taking its inputs and producing its outputs from and to the command line or a file. Components can then be composed at the command line to form (sub)compilers.



2.1 Trees



Internally, a compiler passes around tree-shaped representations of the source program.

In Haskell, such tree-structured data is typically represented in terms of algebraic data types.

For example:

```
data Exp = Const Int
        | Add Exp Exp
        | Mul Exp Exp
```

☞ Recall: a declaration of an algebraic data type introduces both a type constructor (*Exp*) and a family of data constructors (*Const*, *Add*, and *Mul*).



One possibility is to rely on implementations of Haskell's *Read* and *Show* classes:

```
class Read α where
  read :: String → α

class Show α where
  show :: α → String
```

☞ Type classes are groups of types that share some common functionality.



But if we want to encapsulate each compiler component in its own executable program, we need to pass tree-shaped data between programs.

Hence, we need to be able to convert between tree-shaped and “flat”, textual representations of trees, so that we can read and write tree-structured data from and to files and terminals.



Read and *Show* are so-called derivable type classes: for a large set of programmer-defined data types, a Haskell compiler can automatically derive instances of these classes:

```
data Exp = Const Int
        | Add Exp Exp
        | Mul Exp Exp
        deriving (Read, Show)
```



```
*Main> Const 2 'Add' Const 3
Add (Const 2) (Const 3)

*Main> show (Const 2 'Add' Const 3)
"Add (Const 2) (Const 3)"

*Main> read ("Add (Const 2) (Const 3)") :: Exp
Add (Const 2) (Const 3)
```

☞ Why do we need the explicit type annotation in `read ("Add (Const 2) (Const 3)") :: Exp`?



Alternatives

As an alternative to *Read* and *Show* we may consider a more widespread format for representing tree-shaped structures, such as XML or ATerm.

XML is ubiquitous and has excellent tool support. However, the format is quite verbose.

The ATerm format is perhaps less known, but it is specifically targeted at compilers and has fairly good tool support.



Pro:

- ▶ **Easy implementable:** for most data types, *Read* and *Show* can be derived by the Haskell compiler.

Cons:

- ▶ **Haskell-centric:** the format on which *Read* and *Show* operate is essentially the Haskell syntax for constructor application. If we consider exchange between components with different implementation languages, there may be formats that are better supported across different programming languages.
- ▶ **Single-line output:** derived implementations of *show* produce their output on a single line and, hence, the display of large trees appears quite chaotic. Inspecting the output of a compiler component, it may be hard to recognise the tree structure.



The ATerm format

Annotated Term format: a structured representation of arbitrary tree-shaped data.

See:

- ▶ Mark van den Brand, Hayco de Jong, Paul Klint, and Pieter A. Olivier. Efficient annotated terms. *Software—Practice and Experience (SPE)*, 30(3):259–291, 2000.
- ▶ <http://www.meta-environment.org/Meta-Environment/ATerms>.

As a means of representing trees, ATerms are, in a sense, comparable to both XML and algebraic data types.



An ATerm \mathcal{A} can be:

- ▶ An **integer** constant: 2, 3, ...
- ▶ A **floating-point** constant: 3.14, ...
- ▶ A **string** constant: "x", "abc", "", ...
- ▶ A **constructor application** $C(\mathcal{A}_1, \dots, \mathcal{A}_n)$:
Const(2), *Ident*("x"), *Pos*(1, 1), *EOF*, ...
- ▶ A **tuple** $(\mathcal{A}_1, \dots, \mathcal{A}_n)$: (*Ident*("x"), *Op*("*")), ...
- ▶ A **list** $[\mathcal{A}_1, \dots, \mathcal{A}_n]$: [*Const*(2), *Op*("+)], ...
- ▶ An **annotated term** $\mathcal{A} \{ \mathcal{A}_1, \dots, \mathcal{A}_n \}$:
Ident("x") { *Pos*(1, 9) }, ...



Representing ATerms in Haskell

The CCO library exposes a module *CCO.Tree* that contains an algebraic data type for representing ATerms:

```
data ATerm = Integer Integer
           | Float Double
           | String String
           | App Con [ATerm]
           | Tuple [ATerm]
           | List [ATerm]
           | Ann ATerm [ATerm]
           deriving (Eq, Read, Show)

type Con = String
```



ATerm for *Const* 2:

```
Const(2)
```

ATerm for *Const* 2 'Add' *Const* 3:

```
Add(Const(2), Const(3))
```

ATerm for *Const* 2 'Add' (*Const* 3 'Mul' *Const* 5):

```
Add(Const(2)
     , Mul(Const(3), Const(5))
     ) {size(5), depth(3), value(17)}
```



A class for tree types

The class *Tree* contains types that can be represented as ATerms:

```
class Tree α where
  fromTree :: α → ATerm
  toTree   :: ATerm → α
```

The methods *fromTree* and *toTree* convert between trees and ATerms.



To make *Exp* a member of the *Tree* class, we need to provide an instance declaration:

```
instance Tree Exp where
```

```
  fromTree = ...
```

```
  toTree   = ...
```



```
toTree (App "Const" [Integer n]) = Const (fromInteger n)
```

```
toTree (App "Add" [a1, a2]) = Add (toTree a1) (toTree a2)
```

```
toTree (App "Mul" [a1, a2]) = Mul (toTree a1) (toTree a2)
```

```
toTree (Ann a _) = toTree a
```

```
toTree _ = error "toTree: ..."
```

Converting from ATerms to tree types constitutes a partial function.

In this particular example, annotations are ignored.



```
fromTree (Const n)
  = App "Const" [Integer (toInteger n)]
fromTree (Add e1 e2)
  = App "Add" [fromTree e1, fromTree e2]
fromTree (Mul e1 e2)
  = App "Mul" [fromTree e1, fromTree e2]
```

Haskell integers are converted to ATerm integers.

Haskell constructor applications are converted to ATerm constructor applications.



2.2 Running a compiler



For dealing with feedback it is relevant to virtually all parts of a compilation pipeline, we introduce a monad for feedback management:

```
data Feedback  $\alpha$  -- abstract, instance of Monad
```

☞ A monad is a data type μ for representing and constructing effectful computations, supporting the operations $return :: \alpha \rightarrow \mu \alpha$ and $(\gg=) :: \mu \alpha \rightarrow (\alpha \rightarrow \mu \beta) \rightarrow \mu \beta$ (“bind”).



Running a feedback computation

We can “run” a feedback computation and, as a side effect, have its messages written to some file-system object:

```
runFeedback :: Feedback  $\alpha$   $\rightarrow$  -- the computation
             Int  $\rightarrow$  -- verbosity level
             Int  $\rightarrow$  -- severity level
             Handle  $\rightarrow$  -- object to write messages to
             IO (Maybe  $\alpha$ )
```

Feedback computations fail if an error message has been issued.

☞ *Maybe* is the *Prelude* type of computations that may fail:

```
data Maybe  $\alpha$  = Nothing -- for failure
             | Just  $\alpha$  -- for succes
```



Feedback is constructed from three types of messages: log messages, warning messages, and error messages.

Log and warning messages have, respectively, a verbosity and a severity level associated, that allows control over the amount of feedback that is generated.

```
trace :: Int  $\rightarrow$  String  $\rightarrow$  Feedback () -- for log messages
warn :: Int  $\rightarrow$  String  $\rightarrow$  Feedback () -- for warning messages
fail :: String  $\rightarrow$  Feedback  $\alpha$  -- for error messages
```

☞ Actually, *fail* is just $fail :: Monad \mu \Rightarrow String \rightarrow \mu \alpha$ from the *Monad* class.



Example of a feedback computation

```
divider :: Int  $\rightarrow$  Int  $\rightarrow$  Feedback Int
divider m n = do trace 2 "Start divider ..."
                 if n  $\equiv$  0 then fail "Error: division by zero!"
                 else do trace 1 "Dividing ..."
                        return (m 'div' n)

runDivider :: Int  $\rightarrow$  Int  $\rightarrow$  Int  $\rightarrow$  IO ()
runDivider v m n = do let d = divider m n
                      result  $\leftarrow$  runFeedback d v 1 stderr
                      case result of
                        Nothing  $\rightarrow$  return ()
                        Just k  $\rightarrow$  print k
```



```
*Divider> runDivider 2 8 4
Start divider ...
Dividing ...
2

*Divider> runDivider 1 8 4
Dividing ...
2

*Divider> runDivider 2 8 0
Start divider ...
Error: division by zero!
```



Arrows

Arrows are a generalisation of monads and provide a common interface to effectful computations:

```
class Arrow  $\varphi$  where
  pure   :: ( $\alpha \rightarrow \beta$ )  $\rightarrow \varphi \alpha \beta$ 
  ( $\gg$ )  ::  $\varphi \alpha \beta \rightarrow \varphi \beta \gamma \rightarrow \varphi \alpha \gamma$ 
  first  ::  $\varphi \alpha \beta \rightarrow \varphi (\alpha, \gamma) (\beta, \gamma)$ 
  second ::  $\varphi \alpha \beta \rightarrow \varphi (\gamma, \alpha) (\gamma, \beta)$ 
```

See: John Hughes. Generalising monads to arrows. *Science of Computer Programming*, 37(1–3):67–111, 2000.

For *Component*, we have:

```
pure   :: ( $\alpha \rightarrow \beta$ )  $\rightarrow$  Component  $\alpha \beta$ 
( $\gg$ )  :: Component  $\alpha \beta \rightarrow$  Component  $\beta \gamma \rightarrow$  Component  $\alpha \gamma$ 
first  :: Component  $\alpha \beta \rightarrow$  Component  $(\alpha, \gamma) (\beta, \gamma)$ 
second :: Component  $\alpha \beta \rightarrow$  Component  $(\gamma, \alpha) (\gamma, \beta)$ 
```



Components

Recall that a compiler is typically implemented as a pipeline of components.

Components inside the pipeline take as input the output of their predecessor and produce as output the input to their successor. The component at the beginning of a pipeline reads its input from a file or the command line, while the component at the end of a pipeline writes its output to a file or the command line.

To represent components that take inputs of type α and produce outputs of type β , the CCO library provides an abstract type constructor *Component*:

```
data Component  $\alpha \beta$  -- abstract, instance of Arrow
```



Creating components

Primitive components can be created either from pure computations,

```
pure :: ( $\alpha \rightarrow \beta$ )  $\rightarrow$  Component  $\alpha \beta$ 
```

or from computations that involve *Feedback*:

```
component :: ( $\alpha \rightarrow$  Feedback  $\beta$ )  $\rightarrow$  Component  $\alpha \beta$ 
```



Assume we have to construct a compilation pipeline that

- ▶ Reads a floating-point value from the standard input.
- ▶ Checks that the value is nonnegative.
- ▶ Calculates the square root of the value.
- ▶ Writes the square root to the standard output.



Next, we define a *Component* that checks whether its *Double*-input is nonnegative. If the check fails, a warning message is emitted. The component always returns its input unmodified.

```
validator :: Component Double Double
validator = component $ \r → do
    trace_ "Validating ..."
    when (r < 0) (warn_ "Warning: negative input!")
    return r
```

☞ The function *warn_* is defined as *warn_ = warn 1*.



To read, we create a *Component* that consumes a *String* and produces a *Double*. If the *String* cannot be parsed into a *Double*, we issue an error message.

```
parser :: Component String Double
parser = component $ \input → do
    trace_ "Parsing ..."
    case [r | (r, _) ← reads input] of
        r : _ → return r
        _     → fail "Parse error!"
```

☞ The function *trace_* is defined as *trace_ = trace 1*.

☞ For the actual parsing, we instantiate the list-of-successes parser *reads* :: *Read α ⇒ String → [(α, String)]* for *Double*.



The actual calculation is performed by a *Component* that simply applies the *Prelude* function *sqrt* to its input:

```
calculator :: Component Double Double
calculator = component $ \r → do
    trace_ "Calculating ..."
    return (sqrt r)
```



The final *Component* in the pipeline uses the function *show* to turn a *Double* into a *String*.

```
printer :: Component Double String
printer = component $ \r → do
    trace_ "Printing ..."
    return (show r)
```



Wrapping components

Recall that a compilation pipeline, at its source and sink sides, is supposed to exchange flat data with file-system objects.

Indeed, if a *Component* consumes and produces flat data, i.e., *String*-values, it can be turned into a stand-alone program that reads from the standard input channel and writes to the standard output channel:

```
ioWrap :: Component String String → IO ()
```



To assemble a pipeline, we compose individual *Component*-values by means of the *Arrow*-combinator (\gg):

```
pipeline :: Component String String
pipeline = parser >>> validator >>> calculator >>> printer
```

The pipeline itself constitutes a *Component* that takes *String*-inputs to *String*-outputs.



Wrapping components: example

For our square-root calculator with *pipeline* :: *Component String String*, we have:

```
main :: IO ()
main = ioWrap pipeline
```

☞ The *main* function is where a Haskell program starts its execution.



At the command line:

```
% ls
Sqrt.hs

% ghc --make -o sqrt Sqrt.hs
[1 of 1] Compiling Main          ( Sqrt.hs, Sqrt.o )
Linking sqrt ...

% ls
Sqrt.hi Sqrt.hs Sqrt.o sqrt

% echo 25.0 | ./sqrt
Parsing ...
Validating ...
Calculating ...
Printing ...
5.0
```



```
% cat > in.db1
49.0
^D

% ls
Sqrt.hi Sqrt.hs Sqrt.o in.db1 sqrt

% cat in.db1 | ./sqrt > out.db1
Parsing ...
Validating ...
Calculating ...
Printing ...

% ls
Sqrt.hi Sqrt.hs Sqrt.o in.db1 out.db1 sqrt

% cat out.db1
7.0
```



```
% echo xyz | ./sqrt
Parsing ...
Parse error!

% echo -3.14 | ./sqrt
Parsing ...
Validating ...
Warning: negative input!
Calculating ...
Printing ...
NaN
```



2.3 Trees and ATerms



Recall that the ATerm format provides a generic representation for the tree-like structures that typically appear in the internals of the compilation pipeline.

In Haskell, we represented ATerms by

```
data ATerm = Integer Integer
          | Float Double
          | String String
          | App Con [ATerm]
          | Tuple [ATerm]
          | List [ATerm]
          | Ann ATerm [ATerm]
          deriving (Eq, Read, Show)
```



Converting between trees and ATerms

The class *Tree* contains types that can be represented as ATerms:

```
class Tree α where
  fromTree :: α → ATerm
  toTree   :: ATerm → α
```

toTree is a partial function: not every ATerm can be converted into a tree of the appropriate type.

Hence, we let *toTree* produce its result in the *Feedback* monad:

```
class Tree α where
  fromTree :: α → ATerm
  toTree   :: ATerm → Feedback α
```



As an example of an *ATerm*, consider

```
Mul(
  Add(Const(2), Mul(Const(3), Const(5)))
, Mul(Add(Const(7), Const(9)), Const(11))
)
```

represented in Haskell by

```
App "Mul" [
  App "Add" [
    App "Const" [Integer 2]
  , App "Mul" [App "Const" [Integer 3], App "Const" [Integer 5]]
  ]
, App "Mul" [
  App "Add" [App "Const" [Integer 7], App "Const" [Integer 9]]
, App "Const" [Integer 11]
  ]
]
```



2.4 Pretty printing



Often it is necessary to display trees in a human-readable format; for instance, for testing or debugging purposes.

Rather than relying on a derived *Show* instance,

```
App "Mul" [App "Add" [App "Const" [Integer 2], App
"Mul" [App "Const" [Integer 3], App "Const" [In\
teger 5]]], App "Mul" [App "Add" [App "Const" [In\
teger 7], App "Const" [Integer 9]], App "Const" [I\
n\teger 11]]]
```

we typically want to present the user with a representation in concrete syntax:

```
Mul(Add(Const(2), Mul(Const(3), Const(5))), Mul(\
Add(Const(7), Const(9)), Const(11)))
```



Pretty-printer combinators

The CCO library exposes a module *CCO.Printing* that provides a family of **pretty-printer combinators**.

These combinators can be used to construct and combine values of the abstract data type *Doc* of printable documents:

```
data Doc -- abstract
```



Ideally, the concrete-syntax representation makes the tree structure explicit:

```
Mul(
  Add(Const(2), Mul(Const(3), Const(5)))
, Mul(Add(Const(7), Const(9)), Const(11))
)
```

To display a tree in a human-readable form so that the structure of the tree is easily perceivable, we employ a so-called **pretty printer**.



Primitive document constructors

The empty document:

```
empty :: Doc
```

A document containing a specified text:

```
text :: String -> Doc
```

For example:

```
text "pretty"
```

yields

```
pretty
```



The combinator ($>-<$) is used to place one document on top of another:

```
infixr 2 >-<
(>-<) :: Doc → Doc → Doc
```

For example:

```
text "pretty" >-< text "printing"
```

yields

```
pretty
printing
```



The combinator ($>|<$) is used to concatenate two documents:

```
infixr 3 >|<
(>|<) :: Doc → Doc → Doc
```

For example:

```
text "pretty" >|< text "printing"
```

yields

```
prettyprinting
```



The function *indent* is used to indent a document by a given amount of spaces:

```
indent :: Int → Doc → Doc
```

For example:

```
text "pretty" >-< indent 2 (text "printing")
```

yields

```
pretty
  printing
```



If its first operand is a multiline document, ($>|<$) performs what is known as “dovetailing”.

For example:

```
(text "combinators" >-< text "for ") >|<
(text "pretty" >-< text "printing")
```

yields

```
combinators
for pretty
  printing
```



The combinator ($>/// $) is used to introduce a choice between two alternative formattings. When a document is printed, the most space-efficient one is chosen.$$

```
infixr 1 >///  
(>/// $) :: Doc \to Doc \to Doc$ 
```

Parallelisation should only be used with care: nested uses of ($>/// $) can easily cause an explosion in the number of alternatives to consider. Therefore, a local choice can be enforced by means of the function $join$:$$

```
join :: Doc \to Doc
```



2.5 Trees and ATerms



Documents are rendered by means of the function $render$:

```
render :: Int \to Doc \to Maybe String
```

$render$ takes as its first argument the amount of horizontal space that is available for printing. It then constructs the most space-efficient rendering specified by its Doc argument that still fits the available space. If the document cannot be rendered within the available space, $Nothing$ is returned.

As an alternative, $render_*$ always produces a rendering, enlarging the given amount of space as necessary:

```
render_* :: Int \to Doc \to String
```



Exercise: utilities for turning ATerms into trees

§2.5

```
class Tree \alpha where  
  fromTree :: \alpha \to ATerm  
  toTree   :: ATerm \to Feedback \alpha
```

```
data Either \alpha \beta = Left \alpha | Right \beta
```

```
instance (Tree \alpha, Tree \beta) \Rightarrow Tree (Either \alpha \beta) where  
  fromTree (Left x) = App "Left" [fromTree x]  
  fromTree (Right y) = App "Right" [fromTree y]  
  toTree = parseTree [app "Left" (Left <$> arg)  
                    , app "Right" (Right <$> arg)  
                    ]
```



```
test :: String → IO ()
test input = case runFeedback feedback 1 1 stderr of
    Nothing → return ()
    Just tree → print tree
```

where

```
feedback :: Feedback (Either Bool Int)
feedback = toTree (parseATerm input)
```



```
*Test> test "Left()"
Error in ATerm: Left takes 1 argument, but none
were given.
*** In term : Left()

*Test> test "Left(False, True)"
Error in ATerm: Left takes 1 argument, but 2
were given.
*** In term : Left(False, True)
```



```
*Test> test "Left(False)"
Left False

*Test> test "InBetween(True)"
Error in ATerm.
*** Unexpected : InBetween
*** Expected  : Left or Right
*** In term   : InBetween(True)

*Test> test "Left(Perhaps)"
Error in ATerm.
*** Unexpected : Perhaps
*** Expected   : False or True
*** In term    : Perhaps
```



318 lines of Haskell code:

- ▶ comment and whitespace: 161 lines
- ▶ generating error messages: 146 lines

Available from the CCO library through the
CCO.Tree.Parser module.

