

Modules, Signatures and Functors^{★★, ★★★}

Topics

- Modules, Signatures, Functors
- First class modules, Dynlink
- Symbolic Document Processing, Abstract Syntax Trees

Introduction

The goal of this series of exercises is to design the architecture of a command line tool for converting documents between rich text formats, and apply processing passes on them.

We want to end up with command line invocations such as:

```
rtftool -from latex -to text:out.txt -pass flatten-maths  
rtftool -from latex:f.tex -to html:f.html -pass retrieve-urls
```

General design

In the end, the tool will be made of the following components:

- The Rtf intermediate language module.
- A series of input rich text format modules with their own internal representations, either ad hoc or decorator modules around existing libraries.
- Their associated converters to Rtf that may fail if the input document is too complex for the common language.
- The same for output formats.
- A functor Converter.Make that takes an input and an output format and instantiates a dedicated conversion function.
- A registry associating names to first class format and conversion modules and to processing passes.
- A main program analyzing the command line to load the appropriate modules from the registry or by dynlinking, and launching the conversion and processing.

Exercise 1 – Warm-up

Question 1.1 – Write a module Rtf. This module will serve as the core interchange rich text format. It defines a type `t` representing the abstract syntax tree of documents.

The required features are:

- Document structure (paragraphs, sections, chapters, etc.)
- Lists, tables.
- In-line formatting (bold, italic, etc.)
- Internal and external references.
- Equations and pictures.

We may not implement every one of them in every handled format, but we want them to be available in the intermediate representation.

Question 1.2 – Write a proper interface, with `ocaml`doc formatted comments.

Question 1.3 – Write two processing passes: in-line formatting stripping and equation flattening, that will allow them to be exported to formats that do not handle them natively.

Question 1.4 – Add a processing pass registration mechanism in a Registry module: a registration function taking a name and a pass and storing it in an internal global table, and a lookup function that may raise `Not_found`. Pre-register the two previous passes.

Exercise 2 – Formats

The signature for rich text file formats is:

```
1 module type FORMAT = sig
2   type t
3   type input_params
4   type output_params
5   val read : input_params -> t
6   val write : output_params -> t -> unit
7   val name : string
8   val parse_input_params : string -> input_params
9   val parse_output_params : string -> output_params
10 end
```

Defined in a file `converter.ml`

The type `t` is the internal representation specific to this format. We will try and use existing libraries, this is likely to be the types they define.

The types `*_params` are the parameters required for reading and writing files in this format. It will probably be a file name, but may include things such as encoding, author, format version, etc.

Other fields are for registering the format in the system, when we introduce the command line / plugin system.

Question 2.1 – We may not want to provide both input and output for all formats. Rewrite `FORMAT` as a composition of two new interfaces `INPUT_FORMAT` and `OUTPUT_FORMAT`.

Question 2.2 – Implement an `INPUT_FORMAT Text` whose only smart behaviour will be to split the text in paragraphs at double line breaks.

Question 2.3 – Implement an `OUTPUT_FORMAT Latex` or `Html`, whichever you feel most comfortable with.

Question 2.4 – The previous will be enough for testing, but feel free to implement any other format. You can try using `omd` for Markdown, `xmlm` for XHTML.

Exercise 3 – Rtf projections

Question 3.1 – Write two interfaces `TO_RTF` and `OF_RTF` that, for a specific type `t` provide the `to_rtf` and `of_rtf` functions, converting to and from the format specific type to the intermediate `Rtf.t` format.

Question 3.2 – Write `RTF_CONVERTERS` that combine both `TO_RTF` and `OF_RTF`, in a way similar to `FORMAT`.

Question 3.3 – Write instances `Text_to_rtf`, `Latex_of_rtf` or `Html_of_rtf` and for any other format you defined in the previous exercise. Your converters may fail with `"unsupported feature"`.

Exercise 4 – Converter program

Question 4.1 – Write a functor `Converter.Make` taking the input and output formats and their `Rtf` projections and returning a module of signature:

```
1 module type S = sig
2   type input_params
3   type output_params
4   val convert : input_params -> output_params -> (Rtf.t -> Rtf.t) list -> unit
5 end
```

The `convert` function calls the input reader with its parameters, the input to `Rtf.t` converter, applies a series of passes, calls the `Rtf.t` to output converter, and finally the output writer with its parameters.

Question 4.2 – Statically apply the functor and provide a first basic line interface to the instantiated conversion function.

Question 4.3 – Use the `Arg` module to provide a `-pass` option that will apply passes from their registered names on the intermediate document. It will also take input and output arguments that will be passed to the parameter parsers of the input and output formats.

Exercise 5 – Run-time selection of formats

Question 5.1 – Add to the `Registry` functions to register by name input and output formats as first class modules (with their associated converters). Pre-register the formats and converters you already wrote.

Question 5.2 – Update the main so that it supports the `-input` and `-output` arguments as demonstrated in the header: a format name and its arguments separated by a colon.

The main will obtain the format and converter modules by name from the `Registry` and locally apply the `Converter.Make` functor and call the instantiated `convert` function.

Question 5.3 – If the format specified on the command line is `format.cmxs`, `dynlink` the module and then try another search for `format` in the registry. The idea is that the `format` `dynlinkable` module will register a format at toplevel under its own name.

Exercise 6 – To go further

Question 6.1 – Find a better way to handle features unsupported by a processing pass or converter.

Question 6.2 – Design a stream version of the tool that does not maintain the document in memory.

Question 6.3 – Implement document format autodetection.

Question 6.4 – Using objects, implement customizable map and fold over the type `Rtf.t`.

Modules, Signatures and Functors

SOLUTIONS – SOLUTIONS – SOLUTIONS – SOLUTIONS

Solution to exercise 1

Makefile:

```
1 rtftool: rtf.cmx registry.cmx
2     ocamlOPT $^ -o $@
3 %.cmx: %.ml
4     ocamlOPT -c $<
5 %.cmi: %.mli
6     ocamlOPT -c $<
7 rtf.cmx: rtf.cmi
8 registry.cmx: rtf.cmi
```

File rtf.ml:

```
1 type structure =
2   (anchor * structure_item) list
3 and structure_item =
4   | Section of inline option * structure
5   | Picture of inline option * url
6   | Equation of math
7   | Paragraph of inline
8 and math =
9   | Variable of string
10  | Op of string * math list
11  | Seq of math list
12  | Number of float
13  | Inline of inline
14 and inline =
15   inline_item list
16 and inline_item =
17   | Text of string
18   | External of url
19   | Internal of anchor
20   | Style of attribute list * inline
21 and attribute =
22   | Bold | Italic | Monospace
23   | Color of string
24 and url = string
25 and anchor = string
26
27 type t = structure
28
29 let strip_formatting l =
30   let rec structure l =
31     List.map (fun (a, item) -> (a, structure_item item)) l
32   and structure_item = function
```

```

33 | Section (None, s) -> Section (None, structure s)
34 | Section (Some i, s) -> Section (Some (inline i), structure s)
35 | Picture (None, u) -> Picture (None, u)
36 | Picture (Some i, u) -> Picture (Some (inline i), u)
37 | Equation _ as eq -> eq
38 | Paragraph i -> Paragraph (inline i)
39 and inline l =
40   List.flatten (List.map inline_item l)
41 and inline_item = function
42   | Style (_, i) -> inline i
43   | Text _ | External _ | Internal _ as untouched -> [ untouched ]
44 in structure l
45
46 let flatten_maths l =
47   let rec structure l =
48     List.map (fun (a, item) -> (a, structure_item item)) l
49   and structure_item = function
50     | Section (t, s) -> Section (t, structure s)
51     | Equation m -> Paragraph (math m)
52     | Picture _ | Paragraph _ as untouched -> untouched
53   and math = function
54     | Variable n -> [ Text n ]
55     | Op (n, []) -> [ Text n ; Text "()" ]
56     | Op (n, m :: ms) ->
57       [ Text n ; Text "(" ]
58       @ math m
59       @ List.fold_right
60         (fun m acc -> Text "," :: math m @ acc)
61         ms [ Text ")" ]
62     | Seq ms -> List.flatten (List.map math ms)
63     | Number n -> [ Text (string_of_float n) ]
64     | Inline i -> i
65 in structure l

```

File rtf.mli:

```

1 type structure =
2   (anchor * structure_item) list
3 and structure_item =
4   | Section of inline option * structure
5   | Picture of inline option * url
6   | Equation of math
7   | Paragraph of inline
8 and math =
9   | Variable of string
10  | Op of string * math list
11  | Seq of math list
12  | Number of float
13  | Inline of inline
14 and inline =
15   inline_item list
16 and inline_item =
17   | Text of string
18   | External of url
19   | Internal of anchor

```

```

20 | Style of attribute list * inline
21 and attribute =
22 | Bold | Italic | Monospace
23 | Color of string
24 and url = string
25 and anchor = string
26
27 type t = structure
28
29 val strip_formatting : t -> t
30 val flatten_maths : t -> t

```

File registry.ml:

```

1 let all_passes = ref []
2 let register_pass name pass =
3   all_passes := (name, pass) :: !all_passes
4 let lookup_pass name =
5   List.assoc name !all_passes
6 let () =
7   register_pass "strip-formatting" Rtf.strip_formatting ;
8   register_pass "flatten-maths" Rtf.flatten_maths

```

File registry.mli:

```

1 val register_pass : string -> (Rtf.t -> Rtf.t) -> unit
2 val lookup_pass : string -> (Rtf.t -> Rtf.t)

```

Solution to question 2.1

```

1 module type INPUT_FORMAT = sig
2   type t
3   type input_params
4   val read : input_params -> t
5   val name : string
6   val parse_input_params : string -> input_params
7 end
8
9 module type OUTPUT_FORMAT = sig
10  type t
11  type output_params
12  val write : output_params -> t -> unit
13  val name : string
14  val parse_output_params : string -> output_params
15 end
16
17 module type FORMAT = sig
18  type t
19  include INPUT_FORMAT with type t := t
20  include OUTPUT_FORMAT with type t := t
21 end

```

Solution to question 2.2

```

1 type t = string list
2 type input_params = in_channel
3 let name = "text"

```

```

4 let parse_input_params = function
5   | "" | "stdin" | "-" -> stdin
6   | fn -> open_in fn
7 let read fp =
8   let rec read lacc pacc =
9     match input_line fp with
10    | "" -> read [] (pack lacc pacc)
11    | line -> read (line :: lacc) pacc
12    | exception End_of_file -> pack lacc pacc
13   and pack lacc pacc =
14     if lacc = [] then pacc else String.concat "_" (List.rev lacc) :: pacc in
15   let pars = read [] [] in
16   close_in fp ; List.rev pars

```

Solution to question 3.1

```

1 module type TO_RTF = sig
2   type t
3   val to_rtf : t -> Rtf.t
4 end
5
6 module type OF_RTF = sig
7   type t
8   val of_rtf : Rtf.t -> t
9 end

```

Solution to question 3.2

```

1 module type RTF_CONVERTERS = sig
2   type t
3   include TO_RTF with type t := t
4   include OF_RTF with type t := t
5 end

```

Solution to question 3.3

```

1 type t = string list
2 let to_rtf l = List.map (fun t -> Rtf.(None, Paragraph [ Text t ])) l

```

Solution to question 4.1

```

1 module Make
2   (IF : INPUT_FORMAT)
3   (IC : TO_RTF with type t = IF.t)
4   (OF : OUTPUT_FORMAT)
5   (OC : OF_RTF with type t = OF.t)
6 : S with type input_params = IF.input_params
7   and type output_params = OF.output_params = struct
8 type input_params = IF.input_params
9 type output_params = OF.output_params
10 let convert input_params output_params passes =
11   let input = IF.read input_params in
12   let rtf = IC.to_rtf input in
13   let rtf = List.fold_left (fun rtf pass -> pass rtf) rtf passes in
14   let output = OC.of_rtf rtf in
15   OF.write output_params output

```


16 **end**

Solution to question 4.2

```
1 module TextToHtml = Converter.Make (Text) (Text) (Html) (Html)
2 let () =
3   let input_params = Text.parse_input_params Sys.argv.(1) in
4   let output_params = Html.parse_output_params Sys.argv.(2) in
5   TextToHtml.convert input_params output_params []
```

Solution to question 5.1

```
1 module type INPUT = sig
2   type t
3   include Converter.INPUT_FORMAT with type t := t
4   include Converter.TO_RTF with type t := t
5 end
6
7 module type OUTPUT = sig
8   type t
9   include Converter.OUTPUT_FORMAT with type t := t
10  include Converter.OF_RTF with type t := t
11 end
12
13 let all_input_formats : (string * (module INPUT)) list ref = ref []
14 let register_input_format name m =
15   all_input_formats := (name, m) :: !all_input_formats
16 let lookup_input_format name =
17   List.assoc name !all_input_formats
18
19 let all_output_formats : (string * (module OUTPUT)) list ref = ref []
20 let register_output_format name m =
21   all_output_formats := (name, m) :: !all_output_formats
22 let lookup_output_format name =
23   List.assoc name !all_output_formats
24 let () =
25   register_input_format "text" (module Text : INPUT) ;
26   register_output_format "html" (module Html : OUTPUT)
```

Solution to question 5.2

```
1 let () =
2   let input = ref ((module Text : Registry.INPUT), "") in
3   let output = ref ((module Html : Registry.OUTPUT), "") in
4   let passes = ref [] in
5   Arg.parse
6     [ "-input", Arg.String (fun s -> input := parse_input_arg s),
7       "set_the_input_format_and_params" ;
8       "-output", Arg.String (fun s -> output := parse_output_arg s),
9       "set_the_output_format_and_params" ;
10      "-pass", Arg.String (fun pass -> passes := pass :: !passes),
11      "add_a_processing_pass"]
12     (fun _ -> raise (Arg.Bad "anonymous_arguments_unsupported"))
13     "Usage: _rtftool_-input_format:_params_-output_format:params_pass1_pass2_..." ;
14   let (module Input), input_params = !input in
```

```
15 let (module Output), output_params = !output in
16 let input_params = Input.parse_input_params input_params in
17 let output_params = Output.parse_output_params output_params in
18 let module Converter =
19   Converter.Make (Input) (Input) (Output) (Output) in
20 Converter.convert
21   input_params output_params
22   (List.map Registry.lookup_pass (List.rev !passes))
```