

# Draft for \*Typed Markup Revisited\*

JONAS U. BENN

This paper is mostly a case study [for/on] the final tagless approach—which is a solution for the *Expression Problem*. We describe how the final tagless encoding can help to create truly extensible representations for markup documents. This is not a new finding and was first used in the Haskell project HSXML.

We provide a comparison of the FT-encoding with the ADT-encoding that Pandoc is using and describe the essential implementation techniques that HSXML is using to create a context aware encoding. This context aware FT-encoding has great potential for creating a representation that is not only extensible but also provides a lot of guarantees in regards to the well-formedness of the resulting abstract syntax trees.

## ACM Reference Format:

Jonas U. Benn. 2018. Draft for \*Typed Markup Revisited\*. 1, 1 (March 2018), 11 pages.

---

Author's address: Jonas U. Benn.

---

2018. XXXX-XXXX/2018/3-ART  
<https://doi.org/>

50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98

[7]

1 INTRODUCTION

1.1 Motivation

In the age of digital documents, an author of content is confronted with the question which document format to choose. Since every document format has its advantages, one might not want to commit to a specific format to soon.

A series of blog posts might turn into a book or at least a pretty typeset *pdf*. An author might want to give the reader the freedom to read their text on differently sized displays – if the reader has ever tried to read a paper in *pdf*-format on an e-book reader, no further motivation might be needed.

Luckily the problem of decoupling initial text from output seems to be solved by the rise of markup languages such as Markdown and alike. These type of documents can be easily compiled into all sorts of output formats by programs as *Pandoc* [4].

If the reader has no objections to such a publishing system, they might read no further and write away their next *format-agnostic* document. But if they are interested in how they can easily extend the representation of their document and let a type-checker reason about the *well-formedness* of it, they may find the findings gathered in this paper worth while.

1.2 Type-safe extensibility

This paper mostly outlines the ideas of the work on *HSXML: Typed SXML* [2] and the underlying idea of *Finally Tagless Interpreters* [1, 3].

The *final tagless encoding* is a solution to the expression problem [6]. It is closely related to the problem at hand, in that it is concerned with the simultaneous extension of syntactic variants and interpretations of them. This can be seen as an extensibility in two dimensions [add schema like in Object-oriented programming versus abstract data types].

1.3 ?

In short a *final tagless encoded* representation of documents like *HSXML* has in our opinion two major advantages over markup languages such as Pandoc’s internal one:

- (1) Guarantee the well-formedness of the document by construction
- (2) Easy extensibility without losing the guarantees of 1.

While having these two advantages we still do not want to loose perspective and be true to our initial goal:

- (1) Writing documents that are format agnostic—i.e. observe our source in different ways or as described in the Wikipedia-article on *Markup Languages*

Descriptive markup

Markup is used to label parts of the document rather than to provide specific instructions as to how they should be processed. Well-known examples include  $\text{\LaTeX}$ , HTML, and XML. The objective is to decouple the inherent structure of the document from any particular treatment or rendition of it. Such markup is often described as "semantic".

## 2 EXTENSIBILITY OF MARKUP REPRESENTATIONS

### 2.1 Extensible Observers

Pandoc achieves the separation of input and output format by choosing an algebraic data type (ADT) as its intermediate representation. We will quickly sketch why such an encoding leads to an easy extensibility of constructors by looking at a subset of Pandoc's Abstract Syntax Tree and writing some *observers* for it.

Given the representation (Figure 1) we can write observers that interpret this data in different ways (Figure 2). So in the dimension of observers an ADT encoding is obviously extensible.

Now we can construct a tree in the host language and interpret it in two different ways:

```
groceryList :: [Block]
groceryList
  = [ Heading 1 [ Str "Grocery list" ]
      , BulletList [ Paragraph [ Str "1 Banana" ]
                      , Paragraph [ Str "2 "
                                    , Emph [Str "organic"]
                                    , Str " Apples" ] ] ]

groceryListCM :: Markdown
groceryListCM = mconcatMap docToCMark groceryList

groceryListLaTeX :: LaTeX
groceryListLaTeX = mconcatMap docToLaTeX groceryList
```

We can make our life a bit easier by adding an instance for `IsString` for our representation. This injects `String` automatically into our data types by applying `fromString` to it.

```
instance IsString Inline where
  fromString = Str
```

Our initial definition is now even more concise:

```
groceryListShort :: [Block]
groceryListShort
  = [ Heading 1 [ "Grocery list" ]
      , BulletList [ Paragraph ["1 Banana"]
                      , Paragraph ["2 ", Emph ["organic"], " Apples" ] ] ]
```

```

197 data Block
198   = Paragraph    [Inline] -- ^ Paragraph
199   | BulletList   [Block]  -- ^ Bullet list (list of items, each a block)
200   | Heading Int [Inline] -- ^ Heading - level (int) and text (inlines)
201
202 data Inline
203   = Str String      -- ^ Text (string)
204   | EmDash           -- ^ em dash
205   | Emph    [Inline] -- ^ Emphasized text (list of inlines)
206   | Strong  [Inline] -- ^ Strongly emphasized text (list of inlines)
207

```

Fig. 1. This is part of Pandoc's ADT-encoded AST modulo EmDash

```

210 docToCMark :: Block -> Markdown
211 docToCMark (Paragraph text)      = mconcatMap inlineToCMark text
212 docToCMark (BulletList docs)     = addLineBreak $ mconcatMap (mappend "-"
213   " . docToCMark) docs
214 docToCMark (Heading level text) = addLineBreak $ headingPrefix `mappend`
215   mconcatMap inlineToCMark text
216 where
217   headingPrefix = mconcat $ replicate level "#"
218
219 addLineBreak :: Markdown -> Markdown
220 addLineBreak text = text `mappend` "\n"
221
222 inlineToCMark :: Inline -> Markdown
223 inlineToCMark (Str content)      = fromString content
224 inlineToCMark (Emph contents)    = "*" `mappend` mconcatMap inlineToCMark
225   contents `mappend` "*"
226 inlineToCMark (Strong contents) = "**" `mappend` mconcatMap
227   inlineToCMark contents `mappend` "**"
228 inlineToCMark EmDash             = "---"
229
230 docToLaTeX :: Block -> LaTeX
231 ...
232
233 inlineToLaTeX :: Inline -> LaTeX
234 ...
235

```

Fig. 2. Observers of ADT encoding

## 2.2 Extensible Variants

The simple ADT encoding works very well, as long as we have foreseen every variant we might want to create. But as soon as we want to add a new kind of variant (e.g. a node representing the em dash) we are out of luck. Even if we have access to the original ADT-definition and we could add this new variant, this would break all existing observers that were written for the original set of variants.

## 2.3 Relationship to the Expression Problem

To be extensible in the dimension of observers as well as the dimension of the variants—while still guaranteeing statically their compatibility—is quite a challenge and one that is common when writing software. It was coined as the *Expression Problem* by Wadler [6] and many solutions have been proposed.

The most prominent solutions—that are widely used the Haskell-ecosystem—are described in *Data-types a la carte* [5] and in *Finally Tagless, Partially Evaluated* [1]. Kiselyov’s et al. solution to this is, in our opinion, both easy to use and the notation for constructing AST is extremely similar to the ADT-encoded one.

## 3 SIMPLE FINAL TAGLESS ENCODING

Our first attempt to encode our document in the final tagless encoding will not have the distinction between Doc and Inline—which was enforced by the Pandoc-encoding. But later we will see that we are able to recover that property quite easily with great extensibility properties.

The basic idea of the final tagless encoding is as follows:

- Create a type class that specifies all our constructors in Church [add footnote with Böhm Berarducci citation] encoding (Figure 3)
- Parametrize over the return-type and recursive fields of those constructors (Figure 4)

The type classes look basically like a GADT-encoding where all recursive occurrences and the return-type are parametrized over.

The observers will now be instances of these type classes. The reader might notice that we cannot use the same carrier type for different interpretations of our AST — otherwise we would get overlapping instances. This can be quite easily solved by wrapping the carrier type into a *newtype* and add or derive the needed instances for it. In our case Markdown is simply a *newtype* of String. Therefore the instances for IsString and Monoid are straightforward to implement.

Figure 5 shows the implementation of an observer in the final tagless encoding. The implementation is really similar the one in the ADT encoding. But if we have close look, we can see that in Church encoding our observers do not need to be called recursive explicitly. This makes both our code simpler and extensible.

```

295 data Doc = Doc String
296
297 instance Monoid (Doc doc) where
298     mappend (Doc doc1) (Doc doc2) = Doc $ doc1 `mappend` doc2
299     mempty = Doc mempty
300
301 -- Constructors
302
303 class Block where
304     paragraph  ::          [Doc] -> Doc
305     bulletList ::          [Doc] -> Doc
306     heading    :: Int -> [Doc] -> Doc
307
308 class Inline a where
309     emDash ::          Doc
310     str     :: String -> Doc
311     str = Doc
312

```

Fig. 3. First Step FT-encoding

```

315 -- DocConstraint defined using ConstraintKinds
316 type DocConstraint doc = (Monoid doc, IsString doc)
317
318 newtype Doc doc = Doc doc
319
320 instance DocConstraint doc => -- Have to restrict for the use of 'mempty'
321     Monoid (Doc doc) where
322         mappend (Doc doc1) (Doc doc2) = Doc $ doc1 `mappend` doc2
323         mempty = Doc mempty
324
325 -- Constructors
326
327 class Block a where
328     paragraph  ::          [Doc a] -> Doc a
329     bulletList ::          [Doc a] -> Doc a
330     heading    :: Int -> [Doc a] -> Doc a
331
332 class DocConstraint a =>
333     Inline a where
334         emDash ::          Doc a
335         str     :: String -> Doc a
336         str = Doc . fromString
337

```

Fig. 4. Second Step FT-encoding

```

344  -- Implement Markdown observer
345
346  instance Block Markdown where
347    paragraph = fromInline . mconcat
348    bulletList = addLineBreak . mconcat . map (mappend (fromInline "\n- "))
349    heading level = addLineBreak . fromInline . mappend (mconcat $
350      replicate level "#") . mconcat
351
352  addLineBreak :: DocAtts doc => DocWithCtx ctx doc -> DocWithCtx ctx doc
353  addLineBreak (DocWithCtx doc) = DocWithCtx $ doc `mappend` "\n"
354
355  instance Inline Markdown where
356    emDash = "---"
357
358  instance Styles Markdown where
359    emph texts = "*" `mappend` mconcat texts `mappend` "*"
360    strong texts = "**" `mappend` mconcat texts `mappend` "**"
361
362
363  -- Implement LaTeX observer
364
365  instance Block LaTeX where
366    ...
367
368  instance Inline LaTeX where
369    ...
370
371  instance Styles LaTeX where
372    ...
373

```

Fig. 5. Observer implementation in the final tagless encoding



Let's see how our example from above looks in our new encoding:

```
groceryList
= [ heading 1 [str "Grocery list"]
  , bulletList [ paragraph [ str "1 Banana"]
                  , paragraph [ str "2 "
                               , emph [str "fresh"]
                               , str " Apples" ] ] ]
```

As before, we can automate the injection of `String` into our encoding by using the `OverloadedStrings` language pragma. We do this by adding a constraint on the type classes, so every output format (i.e. carrier type) must have an `IsString` instance.

Interestingly `Doc` has now no dependency on `Inline` anymore and we are now allowed to create the following AST:

```
badHeading = [ heading 1 [ heading 2 [str "Headingception!!"] ] ]
```

As noted above, we lost the distinction between `Doc` and `Inline`. But we also gained something—`Doc` can now be used without `Inline` and we can now also add new variants without changing our original variant definitions:

```
class Styles doc where
  emph    :: [DocWithCtx InlineCtx doc] -> DocWithCtx InlineCtx doc
  strong  :: [DocWithCtx InlineCtx doc] -> DocWithCtx InlineCtx doc
```

Not only can we now mix those node types at will, but the type of an expression will reflect which type classes (i.e. algebras) we used for constructing it:

```
stylishNote :: (Inline a, Styles a) => a
stylishNote = strong [str "Green Tea keeps me awake"]
```

That is why the type system can now statically tell us whether we can evaluate `stylishNote` to a particular type.

If we wanted to evaluate an expression, that uses constructors that belong to a type class `X` and we would want to evaluate the expression to some carrier type `C`, `C` has to be instance of `X`. Since this is a static property, it can be decided at compile time.

### 3.1 A short note on GHC's Type Inference

When we define an AST like `stylishNote` GHC's type inference might come in our way. If no type signature for `stylishNote` is supplied GHC will try to infer a concrete type for this definition and not the most generalized type.

We can avoid this by either supplying the generalized type signature—as done above—or using the language pragma `NoMonomorphismRestriction`.

## 4 RECOVER CONTEXT AWARENESS

To regain the context awareness of the Pandoc encoding, we add another field named `ctx` to our `Doc` wrapper (Figure 6). `ctx` is a phantom type and with its help we can specify in which context a constructor can be used. Since phantom types are not materialized on the value level, we are simply using empty data declarations as context types.

```

442 newtype DocWithCtx ctx doc = DocWithCtx doc
443

```

Fig. 6. Context-aware wrapper

```

446
447 -- Context definitions

```

```

448 data InlineCtx

```

```

449 data BlockCtx

```

As shown before, the first Final Tagless encoding had the disadvantage, that we could construct a heading inside another heading. To prohibit this, the heading constructor has the following context-aware definition:

```

454 class Block doc where

```

```

455   heading :: Int -> [DocWithCtx InlineCtx doc] -> DocWithCtx BlockCtx doc
456   ...

```

The type signature states, that the function expects a DocWithCtx-wrapper in the InlineCtx-context and returns a wrapper in the BlockCtx-context. With this refined signature a heading inside a heading will be rejected by the type system.

To convince Haskell's type system that a conversion from InlineCtx to BlockCtx is possible, we can use the following type class:

```

463 class FromInline ctx where

```

```

464   fromInline :: DocWithCtx InlineCtx doc -> DocWithCtx ctx doc
465   fromInline (DocWithCtx doc) = DocWithCtx doc

```

```

467 instance FromInline BlockCtx

```

The set of available contexts should be defined generously, since all independent extensions of the AST should agree on them. This is obviously a restriction—but one that is intended.

It is also possible to create context independent constructors. This can be achieved by parametrizing over the context:

```

474 class Math doc where

```

```

475   qed :: DocWithCtx ctx doc

```

## 5 CONCLUSION

We presented two different encodings that we can choose from for representing a markup language. While the ADT encoding might look like the tool for the job, we have seen that it has some serious limitations. If we do not want to commit to a set of variants forever and do not want to break other people's observers by changing the ADT definition—the final tagless approach might be a good solution for our instance of the *Expression Problem*.

For those who want to study this approach in more depth—the lecture notes on *Typed Tagless Final Interpreters* [3] are a great resource.

REFERENCES

[1] Jacques Carette, Oleg Kiselyov, and Chung-chieh Shan. “Finally Tagless, Partially Evaluated”. In: *Programming Languages and Systems*. Springer Berlin Heidelberg, pp. 222–238. DOI: 10.1007/978-3-540-76637-7\_15. URL: [https://doi.org/10.1007/978-3-540-76637-7\\_15](https://doi.org/10.1007/978-3-540-76637-7_15).

[2] Oleg Kiselyov. *HSXML: Typed SXML*. 2014.

[3] Oleg Kiselyov. “Typed Tagless Final Interpreters”. In: *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2012, pp. 130–174. DOI: 10.1007/978-3-642-32202-0\_3. URL: [https://doi.org/10.1007/978-3-642-32202-0\\_3](https://doi.org/10.1007/978-3-642-32202-0_3).

[4] John MacFarlane. “Pandoc: a universal document converter”. In: URL: <http://pandoc.org> (2018).

[5] Wouter Swierstra. “Data types à la carte”. In: *Journal of functional programming* 18.4 (2008), pp. 423–436.

[6] Philip Wadler. “The expression problem”. In: *Java-genericity mailing list* (1998).

[7] Glen Wright et al. “This Study is Intentionally Left Blank - A systematic literature review of blank pages in academic publishing”. In: (Nov. 2014). DOI: 10.6084/m9.figshare.1230110.v6. URL: [https://figshare.com/articles/This\\_Study\\_is\\_Intentionally\\_Left\\_Blank\\_A\\_systematic\\_literature\\_review\\_of\\_blank\\_pages\\_in\\_academic\\_publishing/1230110](https://figshare.com/articles/This_Study_is_Intentionally_Left_Blank_A_systematic_literature_review_of_blank_pages_in_academic_publishing/1230110).