

Draft for *Typed Final Markup Revisited*

JONAS U. BENN

This paper is mostly a case study of the tagless-final style—which is a solution to the *Expression Problem*. We describe how the tagless-final 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 tagless-final encoding with the *algebraic data type* encoding—that *Pandoc* is using—and describe the essential implementation techniques that HSXML’s implementation is based on to create a context aware encoding. This *context aware tagless-final encoding* has great potential for creating a representation, that is

- truly extensible—i.e. in the dimension of *constructors* and the dimension of *observations*
- provides strong guarantees in regards to the well-formedness of the created abstract syntax trees

ACM Reference Format:

Jonas U. Benn. 2018. Draft for *Typed Final Markup Revisited*. 1, 1 (March 2018), 14 pages.

Author’s address: Jonas U. Benn, mail@benn.in.

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

Blank

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

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 also might want to give the reader the freedom to read their text on different digital devices—e.g. mobile phones, tablets and e-readers.

Luckily the problem of decoupling the initial document from output seems to be solved by the rise of markup languages such as Markdown and the like. These types of documents can be easily compiled into all sorts of output formats by programs such as *Pandoc* [5].

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* [3] and its underlying approach of *tagless-final style* [2, 4].

The *tagless-final style* is a solution to the expression problem [8]. It is closely related to the problem at hand, in that it is concerned with the simultaneous extension of syntactic *constructors* and interpretations of them—we call those *observations*. This can be seen as an extensibility in two dimensions. For those unfamiliar with the the expression problem—section 2 provides a short introduction.

1.3 ?

In short a *tagless-final 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 and full extensibility without loosing the guarantees of 1.

While having these two advantages we still do not want to loose perspective and solve 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 \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 BACKGROUND: THE EXPRESSION PROBLEM

If the reader is already familiar with the expression problem, they might just skip this section and read our case study in section .

2.1 Definition

The following description of the expression problem is short and precise and stems from Zenger's and Odersky's paper *Independently Extensible Solutions to the Expression Problem* [9].

Since software evolves over time, it is essential for software systems to be extensible. But the development of extensible software poses many design and implementation problems, especially if extensions cannot be anticipated. The *expression problem* is probably the most fundamental one among these problems. It arises when recursively defined datatypes and operations on these types have to be extended simultaneously.

In this paper we call those “datatypes” “data variants” or short “variants” and the operations on them are called “observations”. This is inspired by *Extensibility for the Masses* [6].

2.2 Expression Problem by Example

To get a better intuition on what the expression problem is really concerned with, we introduce a small example and will explain the meaning of the *mystical dimensions* with its help. We will represent *algebraic expressions*, like e.g. $2 + 4 \cdot 3$, with two different encodings and will show what kind of extensibility they allow.

2.2.1 ADT encoding. When encoding algebraic expressions in *algebraic data type (ADT) encoding*, we could write this definitions:

```
data Expr
  = Lit Int
  | Add Expr Expr
```

With the above code we defined two data variants—*Lit* and *Add*—and now we can write multiple observations on those variants easily:

```
eval :: Expr -> Int
eval (Lit i)    = i
eval (Add l r) = eval l + eval r

pretty :: Expr -> String
pretty (Lit i)  = show i
pretty (Add l r) =
  "(" ++ pretty l ++ "+"
  ++ "(" ++ pretty r ++ ")"
```

We assess that the ADT encoding is extensible in the dimension of observations.

If we wanted to add another variant—e.g. one for negation—we would have to change not only the ADT definition but also all observations. This might be feasible as long as we feel comfortable with changing the original code. But as soon as someone else wrote observations depending on the original set of variants, we risk breaking compatibility.

2.2.2 *OO encoding*. If we wanted to ensure that our representation is extensible in the dimension of data variants, we could choose the *object oriented (OO) encoding*.

```

data Expr00 = Expr { evalThis :: Int
                      , prettyThis :: String}

newLit :: Int -> Expr00
newLit i = Expr i (show i)

newAdd :: Expr00 -> Expr00 -> Expr00
newAdd l r = Expr evalResult prettyResult
  where
    evalResult = evalThis l + evalThis r
    prettyResult =
      "(" ++ prettyThis l ++ " "
      ++ "+"
      ++ " (" ++ prettyThis r ++ " "
      ++ ")"

ex00 :: Expr00
ex00 = newAdd (newLit 4) (newLit 2)

evalEx00 :: Int
evalEx00 = evalThis ex00

newNeg :: Expr00 -> Expr00
newNeg e = Expr (- evalThis e) (" - " ++ prettyThis e)

```

2.2.3 *Church/Boehm-Berarducci encoding*. Finally we will choose *Boehm-Berarducci (BB) encoding* for our representation.

Boehm and Berarducci used a technique, that is similar to Church encoding, to show that ADTs can be represented by using solely using function application and abstraction in *System F* (i.e. polymorphic lambda-calculus) [1].

(Figure 1)

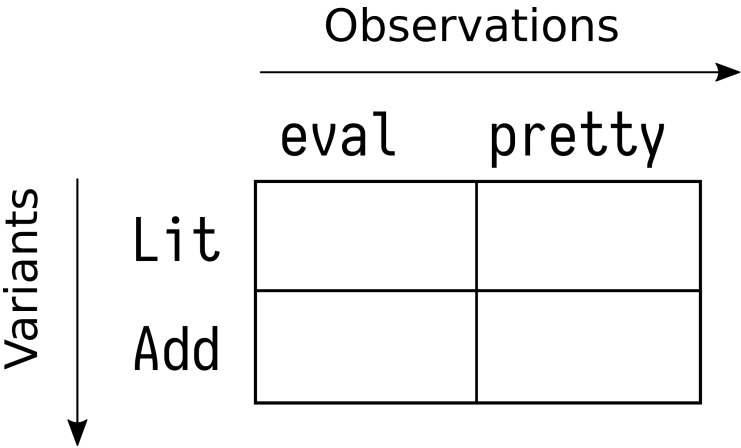


Fig. 1. Dimensions of the *Expression Problem*

3 EXTENSIBILITY OF MARKUP REPRESENTATIONS

3.1 Extensible Observations

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 (AST) and writing some *observations* for it.

Given the representation (Figure 2) we can write observations that interpret this data in different ways (Figure 3). So in the dimension of observations 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" ] ] ]
```

```

344 data Block
345   = Paragraph    [Inline] -- ^ Paragraph
346   | BulletList   [Block]  -- ^ Bullet list (list of items, each a block)
347   | Heading Int  [Inline] -- ^ Heading - level (int) and text (inlines)
348
349 data Inline
350   = Str String    -- ^ Text (string)
351   | EmDash        -- ^ em dash
352   | Emph  [Inline] -- ^ Emphasized text (list of inlines)
353   | Strong [Inline] -- ^ Strongly emphasized text (list of inlines)
354

```

Fig. 2. This is part of Pandoc's ADT-encoded AST modulo EmDash

```

357 docToCMark :: Block -> Markdown
358 docToCMark (Paragraph text)      = mconcatMap inlineToCMark text
359 docToCMark (BulletList docs)     = addLineBreak $ mconcatMap (mappend "-"
360   " . docToCMark) docs
361 docToCMark (Heading level text) = addLineBreak $ headingPrefix `mappend`
362   mconcatMap inlineToCMark text
363 where
364   headingPrefix = mconcat $ replicate level "#"
365
366 addLineBreak :: Markdown -> Markdown
367 addLineBreak text = text `mappend` "\n"
368
369 inlineToCMark :: Inline -> Markdown
370 inlineToCMark (Str content)      = fromString content
371 inlineToCMark (Emph contents)    = "*" `mappend` mconcatMap inlineToCMark
372   contents `mappend` "*"
373 inlineToCMark (Strong contents) = "**" `mappend` mconcatMap
374   inlineToCMark contents `mappend` "**"
375 inlineToCMark EmDash             = "---"
376
377 docToLaTeX :: Block -> LaTeX
378 ...
379
380 inlineToLaTeX :: Inline -> LaTeX
381 ...
382
383 deleteme$
384

```

Fig. 3. Observations of ADT encoding

3.2 Extensible constructors

The simple ADT encoding works very well, as long as we have foreseen every constructor we might want to create. But as soon as we want to add a new kind of constructor—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 constructor, this would break all existing observations that were written for the original set of constructors.

3.3 Relationship to the Expression Problem

To be extensible in the dimension of observations as well as the dimension of the constructors—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 [8] 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* [7] and in *Finally Tagless, Partially Evaluated* [2]. 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.

4 SIMPLE TAGLESS-FINAL ENCODING

Our first attempt to encode our document in the tagless-final 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 tagless-final encoding is as follows:

- Create a type class that specifies all our constructors in Church [add footnote with Böhm Berarducci citation] encoding (Figure 4)
- Parametrize over the return-type and recursive fields of those constructors (Figure 5)

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

The observations 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 6 shows the implementation of an observer in the tagless-final encoding. The implementation is really similar the one in the ADT encoding. But if we have close look, we can see that—since our data type is Church encoded—the observations do not need to be called recursive explicitly. This makes both our code simpler and is essential for extensibility.

```

442 newtype Doc = Doc String
443
444 instance Monoid (Doc doc) where
445     mappend (Doc doc1) (Doc doc2) = Doc $ doc1 `mappend` doc2
446     mempty = Doc mempty
447
448 -- Constructors
449
450 class Block where
451     paragraph :: [Doc] -> Doc
452     bulletList :: [Doc] -> Doc
453     heading :: Int -> [Doc] -> Doc
454
455 class Inline a where
456     emDash :: Doc
457     str :: String -> Doc
458     str = Doc
459 
```

Fig. 4. First Step FT-encoding

```

462 newtype Doc doc = Doc doc
463
464 -- DocConstraint defined using ConstraintKinds
465 type DocConstraint doc = (Monoid doc, IsString doc)
466
467 instance DocConstraint doc => -- Have to restrict for the use of 'mempty'
468     Monoid (Doc doc) where
469     mappend (Doc doc1) (Doc doc2) = Doc $ doc1 `mappend` doc2
470     mempty = Doc mempty
471
472 -- Constructors
473
474 class Block a where
475     paragraph :: [Doc a] -> Doc a
476     bulletList :: [Doc a] -> Doc a
477     heading :: Int -> [Doc a] -> Doc a
478
479 class DocConstraint a =>
480     Inline a where
481     emDash :: Doc a
482     str :: String -> Doc a
483     str = Doc . fromString
484 
```

Fig. 5. Second Step FT-encoding

```

491  -- Implement Markdown observer
492
493  instance Block Markdown where
494    paragraph = fromInline . mconcat
495    bulletList = addLineBreak . mconcat . map (mappend (fromInline "\n- "))
496    heading level = addLineBreak . fromInline . mappend (mconcat $
497      replicate level "#") . mconcat
498
499  addLineBreak :: DocAtts doc => DocWithCtx ctx doc -> DocWithCtx ctx doc
500  addLineBreak (DocWithCtx doc) = DocWithCtx $ doc `mappend` "\n"
501
502  instance Inline Markdown where
503    emDash = "---"
504
505  instance Styles Markdown where
506    emph texts = "*" `mappend` mconcat texts `mappend` "*"
507    strong texts = "**" `mappend` mconcat texts `mappend` "**"
508
509
510  -- Implement LaTeX observer
511
512  instance Block LaTeX where
513    ...
514
515  instance Inline LaTeX where
516    ...
517
518  instance Styles LaTeX where
519    ...
520
521

```

Fig. 6. Observer implementation in the tagless-final 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 constructors without changing our original constructor 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 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.

4.1 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`.

5 RECOVER CONTEXT AWARENESS

To regain the context awareness of the Pandoc encoding, we add another field named `ctx` to our `Doc` wrapper (Figure 7). `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.

```

589 newtype DocWithCtx ctx doc = DocWithCtx doc
590

```

Fig. 7. Context-aware wrapper

```

593 -- Context definitions
594

```

```

595 data InlineCtx
596

```

```

596 data BlockCtx
597

```

As shown before, the first tagless-final 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:

```

601 class Block doc where

```

```

602   heading :: Int -> [DocWithCtx InlineCtx doc] -> DocWithCtx BlockCtx doc
603   ...
604

```

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:

```

610 class FromInline ctx where

```

```

611   fromInline :: DocWithCtx InlineCtx doc -> DocWithCtx ctx doc
612   fromInline (DocWithCtx doc) = DocWithCtx doc
613

```

```

614 instance FromInline BlockCtx
615

```

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:

```

621 class Math doc where

```

```

622   qed :: DocWithCtx ctx doc
623

```

624 6 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. Especially if our set of constructors might scale up and we would do not want to break other people's observations by changing the ADT definition—the tagless-final approach might be a good solution also for this instance of the *Expression Problem*.

For those who want to study this approach in more depth—the lecture notes on *Typed Tagless Final Interpreters* [4] are a great resource.

REFERENCES

- [1] Corrado Böhm and Alessandro Berarducci. “Automatic Synthesis of Typed Lambda-Programs on Term Algebras”. In: *Theor. Comput. Sci.* 39 (1985), pp. 135–154. DOI: [10.1016/0304-3975\(85\)90135-5](https://doi.org/10.1016/0304-3975(85)90135-5). URL: [https://doi.org/10.1016/0304-3975\(85\)90135-5](https://doi.org/10.1016/0304-3975(85)90135-5).
- [2] 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](https://doi.org/10.1007/978-3-540-76637-7_15). URL: https://doi.org/10.1007/978-3-540-76637-7_15.
- [3] Oleg Kiselyov. *HSXML: Typed SXML*. 2014.
- [4] 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](https://doi.org/10.1007/978-3-642-32202-0_3). URL: https://doi.org/10.1007/978-3-642-32202-0_3.
- [5] John MacFarlane. “Pandoc: a universal document converter”. In: URL: <http://pandoc.org> (2018).
- [6] Bruno C. d. S. Oliveira and William R. Cook. “Extensibility for the Masses”. In: *ECOOP 2012 – Object-Oriented Programming*. Ed. by James Noble. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 2–27. ISBN: 978-3-642-31057-7.
- [7] Wouter Swierstra. “Data types à la carte”. In: *Journal of functional programming* 18.4 (2008), pp. 423–436.
- [8] Philip Wadler. “The expression problem”. In: *Java-genericity mailing list* (1998).
- [9] Matthias Zenger and Martin Odersky. *Independently extensible solutions to the expression problem*. Tech. rep. 2004.