

Documentation MyPandoc

Document Structure Overview

Our Pandoc system uses an Abstract Syntax Tree (AST) to represent structured documents. The AST handles Markdown, JSON, and XML inputs/outputs through a unified representation.

Core Data Types

The AST consists of these primary types:

Document

```
data Document = Document Meta [Block]
```

A document consists of metadata (`Meta`) and a list of content blocks (`[Block]`).

Meta

```
data Meta = Meta
  { metaTitle   :: [Inline]    -- Document title
  , metaAuthors :: [[Inline]]  -- List of authors
  , metaDate    :: [Inline]    -- Date of the document
  }
```

Metadata includes title, authors (multiple authors possible), and date.

Block

```
data Block
  = Plain [Inline]           -- Plain text without paragraph
  | Para [Inline]            -- Paragraph
  | CodeBlock String         -- Code block
  | RawBlock Format String   -- Raw block (HTML, LaTeX, etc.)
  | OrderedList ListAttributes [[Block]] -- Numbered list
  | BulletList [[Block]]    -- Bullet list
  | DefinitionList [(Inline), [[Block]]] -- Definition list
  | Header Int [Inline]     -- Section header
  | Section Int [Inline] [Block] -- Section with nested blocks
  | Null                    -- Empty block
```

Block elements represent structural components of a document.

Inline

```
data Inline
  = Str String           -- Plain text
  | Emph [Inline]        -- Emphasized text (italic)
  | Strong [Inline]      -- Strong emphasis (bold)
  | Code String          -- Inline code
  | RawInline Format String -- Raw inline content
  | Link [Inline] Target -- Hyperlink
  | Image [Inline] Target -- Image
  | Note [Block]         -- Footnote
  | Span [Inline]        -- Generic inline container
```

Inline elements represent text-level formatting and entities.

Supporting Types

```
type Target = (String, String)           -- URL and title
newtype Format = Format String            -- Format specifier
type ListAttributes = (Int, ListNumberStyle, ListNumberDelim)
data ListNumberStyle = DefaultStyle | Example | Decimal | LowerRoman | UpperRoman
                      | LowerAlpha | UpperAlpha
data ListNumberDelim = DefaultDelim | Period | OneParen | TwoParens
```

These types provide additional details for specialized elements.

Implementing a New Renderer

To implement a new renderer for the AST, you need to create functions that traverse the AST and generate output in your target format. Here's how to approach it:

1. Create a Main Renderer Function

Start by creating a main function that takes a `Document` and returns a `String`:

```
myFormatRender :: Document -> String
myFormatRender (Document meta blocks) =
    renderMeta meta ++ renderBlocks blocks
```

2. Implement Metadata Rendering

Create a function to render document metadata:

```
renderMeta :: Meta -> String
renderMeta (Meta title authors date) =
    -- Format-specific metadata rendering
```

3. Implement Block Rendering

Create functions to render Block elements:

```
renderBlocks :: [Block] -> String
renderBlocks blocks = concatMap renderBlock blocks
```

```
renderBlock :: Block -> String
renderBlock (Para inlines) =
    -- Paragraph rendering
renderBlock (Header level inlines) =
    -- Header rendering
renderBlock (BulletList items) =
    -- Bullet list rendering
-- ... and so on for other block types
```

4. Implement Inline Rendering

Create functions to render Inline elements:

```
renderInlines :: [Inline] -> String
renderInlines inlines = concatMap renderInline inlines
```

```
renderInline :: Inline -> String
renderInline (Str text) =
    -- Plain text rendering
```

```
renderInline (Emph inlines) =
    -- Emphasized text rendering
renderInline (Strong inlines) =
    -- Bold text rendering
-- ... and so on for other inline types
```

5. Handle Special Cases

Implement special cases for your format:

```
escapeSpecialChars :: String -> String
escapeSpecialChars = -- Format-specific character escaping
```

Example: JSON Renderer Implementation

The provided `JsonRenderer` module demonstrates how to implement a renderer:

1. **Main entry point:** `jsonRender` function converts a full Document to JSON
2. **Block rendering:** `jsonRenderBlock` handles each type of Block
3. **Inline rendering:** `jsonRenderInline` and `jsonRenderInlines` handle Inline elements
4. **Helper functions:** For escaping characters, formatting lists, etc.

Key Rendering Patterns

Several patterns can be observed in the JSON renderer that apply to new renderers:

Recursive Rendering

Block elements may contain inline elements, and inline elements may contain other inline elements. Use recursive functions to handle this.

String Escaping

Always handle escaping of special characters required by your output format.

Element Mapping

Map each AST element type to its corresponding representation in your target format.

Composition

Build complex output by composing the output of smaller elements.

Implementation Tips

1. Start with a simple, functional renderer that handles basic elements
2. Add support for more complex elements incrementally
3. Test with simple documents before moving to complex ones
4. Pay attention to whitespace and formatting in your output
5. Consider using helper functions for repeated patterns

Common Challenges

- **Nested formatting:** Elements can be deeply nested, requiring careful recursive handling
- **Special characters:** Different output formats have different character escaping needs
- **Whitespace handling:** Ensuring correct spacing between elements
- **Context-sensitive rendering:** Some elements may need to be rendered differently based on context

Adding New Formats to MyPandoc

To extend MyPandoc with new input or output formats, you need to modify several components of the system. Here's a step-by-step guide:

Adding a New Output Format

1. Create a new **renderer module** following the renderer implementation guidelines above.
2. Update the supported output formats list:

```
getSupportedOutput :: [String]
getSupportedOutput = ["html", "json", "markdown", "xml", "your_new_format"]
```

3. Register your new renderer in the `getValidDoc` function:

```
getValidDoc :: String -> Document -> String
getValidDoc format doc
  | format == "html" = htmlRender doc
  | format == "xml" = xmlRender doc
  | format == "json" = jsonRender doc
  | format == "markdown" = markdownRender doc
  | format == "your_new_format" = yourFormatRender doc
  | otherwise = "Not valid format"
```

4. Export your new renderer function in the appropriate module file.

Adding a New Input Format

1. Create a new **parser module** following the parser implementation guidelines below.
2. Update the `getinputfile` function to recognize your format:

```
getinputfile :: String -> Opts -> Maybe Document
getinputfile content opts
  | fromJust (inputFormat opts) == "json" = parseJson content newdoc
  | fromJust (inputFormat opts) == "markdown" = parsemd content newdoc
  | fromJust (inputFormat opts) == "xml" = parsexml content newdoc
  | fromJust (inputFormat opts) == "your_new_format" = parseYourFormat content newdoc
  | otherwise = Nothing
```

3. Export your parser function in the appropriate module file.

Example: Adding LaTeX Support

Here's a simplified example of adding LaTeX as both an input and output format:

Output Format (LaTeX Renderer)

1. Create a LaTeX renderer module:

```
module LaTeXRender (latexRender) where

import Types

latexRender :: Document -> String
latexRender (Document meta blocks) =
  "\\documentclass{article}\\n\\begin{document}\\n" ++
  renderLatexMeta meta ++
  renderLatexBlocks blocks ++
```

```

"\end{document}"

renderLatexMeta :: Meta -> String
renderLatexMeta (Meta title authors date) =
  "\\title{" ++ renderLatexInlines title ++ "}\n" ++
  "\\author{" ++ renderLatexAuthors authors ++ "}\n" ++
  "\\date{" ++ renderLatexInlines date ++ "}\n" ++
  "\\maketitle\n\n"

renderLatexAuthors :: [[Inline]] -> String
renderLatexAuthors = -- Implementation

renderLatexBlocks :: [Block] -> String
renderLatexBlocks = -- Implementation

renderLatexInlines :: [Inline] -> String
renderLatexInlines = -- Implementation

```

2. Update the supported formats list:

```

getSupportedOutput :: [String]
getSupportedOutput = ["html", "json", "markdown", "xml", "latex"]

```

3. Register the renderer:

```

getValidDoc :: String -> Document -> String
getValidDoc format doc
  | format == "html" = htmlRender doc
  | format == "xml" = xmlRender doc
  | format == "json" = jsonRender doc
  | format == "markdown" = markdownRender doc
  | format == "latex" = latexRender doc
  | otherwise = "Not valid format"

```

Input Format (LaTeX Parser)

1. Create a LaTeX parser module:

```

module LaTeXParser (parseLatex) where

import Types
import Lib

-- Define intermediate representation
data LaTeXNode
  = LaTeXCommand String [LaTeXNode]
  | LaTeXEnvironment String [LaTeXNode]
  | LaTeXText String

-- Create parser functions
parseLatexNode :: Parser LaTeXNode
parseLatexNode = -- Implementation

-- Create conversion functions
latexToDocument :: LaTeXNode -> Document
latexToDocument = -- Implementation

```

```

-- Create entry point
parseLatex :: String -> Document -> Maybe Document
parseLatex input _ = do
  (node, _) <- parseLatexNode input
  return (latexToDocument node)

```

2. Update the input file handler:

```

getinputfile :: String -> Opts -> Maybe Document
getinputfile content opts
  | fromJust (inputFormat opts) == "json" = parseJson content newdoc
  | fromJust (inputFormat opts) == "markdown" = parsemd content newdoc
  | fromJust (inputFormat opts) == "xml" = parsexml content newdoc
  | fromJust (inputFormat opts) == "latex" = parseLatex content newdoc
  | otherwise = Nothing

```

Parsing Input Formats

In addition to rendering AST to different output formats, Mini Pandoc also needs to parse various input formats (Markdown, XML, JSON) into the AST.

Parser Fundamentals

Our parsing system is built on a basic parser combinator approach:

```

-- Basic parser type from Lib
type Parser a = String -> Maybe (a, String)

```

A parser takes a string input and returns either: - Nothing (if parsing fails) - Just (result, remaining) where result is the parsed value and remaining is the unparsed portion of the input

Example: XML Parser Implementation

The provided XML parser demonstrates the approach:

1. **Parser structure definition:**

```

data XmlValue
  = XmlElement String Int [(String, String)] [XmlValue]
  | XmlText String

```

2. **Basic parsers** for elements, text, and attributes:

```

parseXmlValue :: Int -> Parser XmlValue
parseXmlText :: Parser XmlValue
parseXmlElement :: Int -> Parser XmlValue
parseXmlAttribute :: Parser (String, String)

```

3. **Conversion functions** that transform the parsed structure to AST:

```

xmlToDocument :: XmlValue -> Document
getmeta :: [XmlValue] -> Meta
getblocks :: [XmlValue] -> [Block]
getblock :: XmlValue -> [Block]
getInline :: XmlValue -> [Inline]

```

4. **Helper functions** for finding elements and extracting content:

```

findElement :: String -> [XmlValue] -> Maybe XmlValue
extractText :: String -> [XmlValue] -> Maybe String

```

5. **Entry point function** that handles the entire parsing process:

```
parsexml :: String -> Document -> Maybe Document
```

Implementing a New Parser

To create a new parser for an input format, follow these steps:

1. **Define an intermediate representation** Create data types that represent the structure of your input format:

```
data MyFormatValue
  = MyElement String [MyFormatValue]
  | MyText String
  -- Other format-specific elements
```

2. **Create basic parsers** Implement parsers for each element in your format:

```
parseMyFormat :: Parser MyFormatValue
parseMyFormat = parseMyElement <|> parseMyText
```

```
parseMyElement :: Parser MyFormatValue
parseMyElement = do
  -- Format-specific parsing logic
```

```
parseMyText :: Parser MyFormatValue
parseMyText = do
  -- Text parsing logic
```

3. **Implement conversion to AST** Create functions to convert your parsed structure to the AST:

```
myFormatToDocument :: MyFormatValue -> Document
myFormatToDocument value =
  Document meta blocks
  where
    meta = extractMeta value
    blocks = extractBlocks value
```

```
extractMeta :: MyFormatValue -> Meta
extractMeta = -- Logic to extract metadata
```

```
extractBlocks :: MyFormatValue -> [Block]
extractBlocks = -- Logic to extract blocks
```

4. **Create helper functions** Implement any helper functions needed for extraction and conversion:

```
findMyElement :: String -> [MyFormatValue] -> Maybe MyFormatValue
extractMyText :: MyFormatValue -> String
```

5. **Create the main entry point** Provide a function that parses input text and returns a Document:

```
parseMyFormat :: String -> Document -> Maybe Document
parseMyFormat input _ = do
  value <- runMyFormatParser input
  return (myFormatToDocument value)
```

Parser Combinator Patterns

Several key patterns are used in the parsing system:

1. **Alternative (<|>):** Try one parser, and if it fails, try another

```
parseXmlValue level = parseXmlElement level <|> parseXmlText
```

2. **Sequential composition (do notation):** Run parsers in sequence

```
parseXmlElement level = do
  _ <- parseChar '<'
  name <- parseSome (parseAnyChar (['a'..'z'] ++ ['A'..'Z'] ++ ['0'..'9']))
  -- other parsing steps
```

3. **Repetition (parseMany, parseSome):** Parse repeated patterns

```
attribute <- parseMany parseXmlAttribute
```

4. **Character and string matching:**

```
_ <- parseChar '<'
_ <- parseString "</"
```

Implementation Tips for Parsers

1. **Start small:** Begin with basic elements and add complexity gradually
2. **Test frequently:** Test each parser function individually before combining
3. **Handle edge cases:** Consider empty documents, malformed input, etc.
4. **Use meaningful error messages:** When possible, provide context for parsing failures
5. **Optimize for readability:** Parser code should be clear and maintainable

By following these patterns, you can implement parsers for various input formats that convert content into your AST for further processing and rendering.