

Normalisierung von Dokumentformaten für einfache Beschreibungstexte

Frieda Josi

12. Juni 2017

Aufgabenstellung

Funktionsweise von Pandoc

Filter in Pandoc

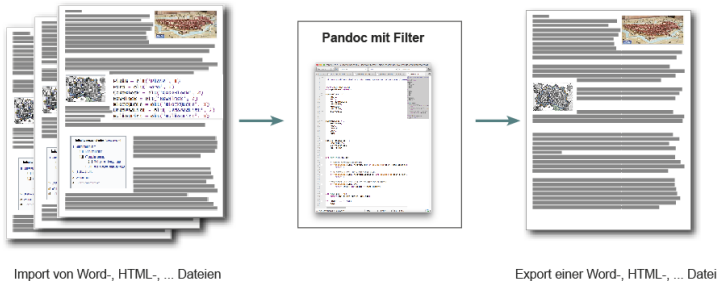
Eigener Filter

Quellen

Weitere Möglichkeiten

Aufgabenstellung: Normalisierung mit Pandoc

- ▶ Filter für Kommandozeilentool Pandoc entwickeln:
 - ▶ Aus Dokumentformate (Word, OpenOffice, HTML, LaTeX, Markdown,...) werden nur Fotos, Abbildungen, und deren Beschreibungstexte in das Ausgabeformat übernommen.
 - ▶ Der Filter wird mit Python umgesetzt

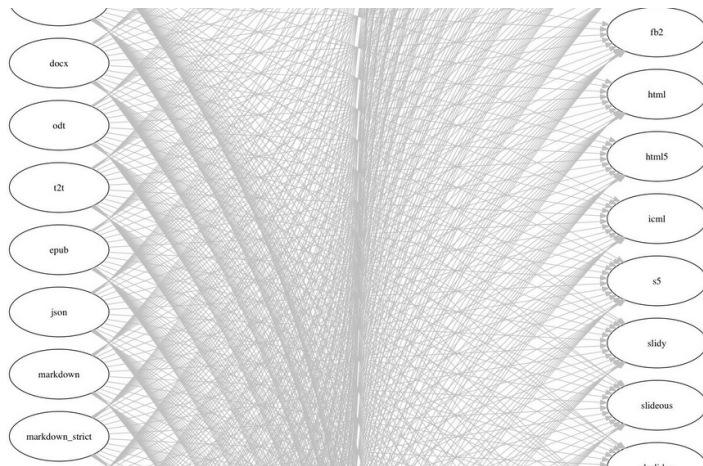


Funktionsweise von Pandoc: Mögliche Dokumentformate

Pandoc ist ein Programm, das viele verschiedene Dokumentformate ineinander umwandeln kann.

Lesen/Importieren:	Schreiben/Exportieren:
Markdown, CommonMark, reStructuredText, textile, MediaWiki, HTML, DocBook, LaTeX, Microsoft Word DOCX (OOXML), LibreOffice/OpenOffice ODT, EPUB, TWiki, Org-mode	HTML-Formate, Textverarbeitungsformate (Microsoft Word, LibreOffice/OpenOffice ODT, ...), Ebooks, Dokumentationsformate (DocBook, ...), Seitenlayoutformate (z.B. InDesign), TeX-Formate (z.B. LaTeX), PDF via LaTeX, Vereinfachte Auszeichnungssprachen (z.B. CommonMark, Markdown, reStructuredText, ...)

Funktionsweise von Pandoc: Ausschnitt Dokumentformate



Quelle: <http://pandoc.org/>

Funktionsweise von Pandoc: Markdown

- ▶ Pandoc basiert auf einer Erweiterung der Auszeichnungssprache Markdown.
- ▶ Bei der Umwandlung sollen Strukturelemente eines Dokuments erhalten bleiben, nicht ihre Formatierung. Problematisch sind aber komplexe Tabellenstrukturen.

```
1 # Überschrift in Ebene 1
2 ### Überschrift in Ebene 4
3
4 > Dies ist ein Zitat mit zwei Absätzen. Lorem ipsum dolor sit amet,
5 > consectetur adipiscing elit.
6 >
7 > Donec sit amet nisl. Aliquam semper ipsum sit amet velit.
8
9 ![nur ein Beispiel](https://commons.wikimedia.org/wiki/File:Example_de.jpg "Beispielbild")
10
11 1. Ein Punkt in einer geordneten Liste
12 1. Hund
13   1. Hund im Hund 1
14   2. Die Katze im Hund 2
15   3. Die Katze im Hund 3
16
17 Jetzt eine horizontale Linie
18 -----
19 Ein normaler Absatz, mit wenig Text. Normaler Text wird so dargestellt wie eingegeben.
20
21 Eine Leerzeile erzeugt einen Absatz.
22
23     Ein Code-Block
24     durch Einrückung mit vier Leerzeichen
```

Funktionsweise von Pandoc: Umwandlung

- ▶ Pandoc besteht aus einer Haskellbibliothek und einem eigenständigen Kommandozeilenprogramm.
- ▶ Diese Bibliothek besteht aus Modulen für jedes Eingabe- und Ausgabeformat. Es können neue Dokumentformate über neue Module hinzugefügt werden.
- ▶ Es gibt verschiedene Kommandozeilen Lese- und Schreiboptionen (hier `-s` für das Lesen der Importdatei (smart) und `-o` für die Exportdatei und auch das Ausgabeformat (output)).

Pandoc verwenden über Terminaleingabe:

```
pandoc -s ausstellung.md -o ausstellung.html
```

Quellen: <http://pandoc.org>

Pandoc Filter Funktionsweise

Pandoc bietet Schnittstelle für Einfluss auf Syntaxbaum der Zwischendarstellung. Dafür werden Filter (z.B. in Python, Perl, Haskell, ... geschrieben) eingesetzt.

```
INPUT --reader--> AST --filter--> AST --writer--> OUTPUT
```

- ▶ Das Eingabedokument wird von Pandoc in eine Zwischendarstellung (AST) gelesen,
- ▶ der Filter nimmt Einfluss darauf (z.B. Alle Überschriften werden mit Grossbuchstaben geschrieben),
- ▶ dann wird modifizierter Syntaxbaum erstellt,
- ▶ danach erfolgt das Schreiben des Zieldokumentformates.

Terminaleingabe: `pandoc -s ausstellung.md -filter remove.py -o ausgabe.html`

Quelle: <http://pandoc.org/scripting.html>

Pandoc-Elemente

Erstellte Übersicht als Vorbereitung für den Filter.

Elemente, die gelöscht werden:	Elemente, die bleiben:	Elemente, die teilweise bleiben:
CodeBlock, RawBlock, BlockQuote, HorizontalRule, Table, Null, Strikeout, Quoted, Code, Math, RawInline, Note	Plain, Para, LineBlock, Header, Div, Str, Emph, Strong, Superscript, Subscrip, SmallCaps, Cite, Space, SoftBreak, LineBreak, Link, Image, Span	OrderedList, DefinitionList, BulletList

Quelle: <http://hackage.haskell.org/packagepandoc-types-1.17.0.5/docs/Text-Pandoc-Definition.htm>

Eigener Filter Teil 1

Der erstellte Python-Filter entfernt alle nicht gewollten Elemente.

Datei: remove.py

```
import panflute as pf
from panflute import *
```

```
removeBlock = (
    Table ,
    Null ,
    HorizontalRule ,
    CodeBlock ,
    RawBlock ,
    RawInline ,
    BlockQuote )
```

```
removeInline = (
    Strikeout ,
    Note ,
    Quoted ,
    Code ,
    Math )
```

```
removeNestedLists = (
    OrderedList ,
    DefinitionList ,
    BulletList )
```

Eigener Filter Teil 2

```
def action(elem, doc):  
    # remove forbidden elements  
    if isinstance(elem, removeBlock) or isinstance(elem, removeInline):  
        return []  
  
    # no nested lists, only the first level  
    if isinstance(elem, removeNestedLists) and isinstance(elem.parent,  
                                                             pf.ListItem):  
        return []  
  
    # transform Div into list of its content elements  
    if isinstance(elem, Div):  
        return [block for block in elem.content]  
  
def main(doc=None):  
    return run_filter(action, doc=doc)  
  
if __name__ == '__main__':
```

Durchführung

Terminaleingabe:

`pandoc -s ausstellung.md -filter remove.py -o ausstellung.html`

Import Markdowndatei -> Filter Einsatz -> Export HTML-Datei

Import Markdown

```
1  ---
2  title: Hinweistafeln die in der Ausstellung **Mi
3  ...
4  Weit hinten, hinter den Wortbergen, fern der Länd
5
6  Weit hinten, hinter den Wortbergen
7
8  Weit hinten, hinter den Wortbergen
9  Weit hinten, hinter den Wortbergen
10
11 | was | wo | wie | warum | wieso |
12 |---|---|---|---|---|
13 | Nichts | Noch weniger | Gar nichts | Gar nicht
14 | Nichts | Noch weniger | Gar nichts | Gar nicht
15
16 > Dies ist ein Zitat mit zwei Absätzen. Lorem ipsum
17 > consectetur adipiscing elit.
18 >
19 > Donec sit amet nisl. Aliquam semper ipsum sit amet
20 > suspendisse id sem consectetur libero luctus a
21
22 1. Hund
23   1. Hund im Hund
24   2. Katze in Katze
25   3. Katze in Katze
26 1. Katze
```

Export HTML

```
1  <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML
2  <html xmlns="http://www.w3.org/1999/x
3  <head>
4  <meta http-equiv="Content-Type" con
5  <meta http-equiv="Content-Style-Typ
6  <meta name="generator" content="pan
7  <title>Hinweistafeln die in der Aus
8  <style type="text/css">code{white-s
9  </head>
10 <body>
11 <div id="header">
12 <h1 class="title">Hinweistafeln die i
13 </div>
14 <p>Weit hinten, hinter den Wortbergen
15 <p>Weit hinten, hinter den Wortbergen
16 <p>Weit hinten, hinter den Wortbergen
17 <ol style="list-style-type: decimal">
18 <li>Hund</li>
19 <li>Katze</li>
20 <li>Maus</li>
21 </ol>
22 <p>Noch mehr einfacher Text. Noch meh
23 <ul>
24 <li>Lorem ipsum dolor sit amet, conse
25 <li>Donec sit amet nisl. Aliquam semp
26 </ul>
```

Quellen

Pandoc: <http://pandoc.org/>

Pandoc Github: <https://github.com/jgm/pandoc>

Pandoc Filters:

<https://github.com/jgm/pandoc/wiki/Pandoc-Filters>

Pandoc Panflute: <https://github.com/sergiocorreia/panflute>

Pandoc Optionen: <http://pandoc.org/MANUAL.html#options>

Pandoc Elemente:

<http://hackage.haskell.org/packagepandoc-types-1.17.0.5/docs/Text-Pandoc-Definition.html>

Pandoc online testen: <http://pandoc.org/try/>

Pandoc Filter schreiben: <http://pandoc.org/scripting.html>

Ergänzung

Umwandeln einer PDF-Datei zu Word, HTML,... geht mit Pandoc nicht.

Alternative: Acrobat Pro DC

- ▶ PDF -> Word, Excel (Sehr guter Export)
- ▶ PDF -> HTML (Sehr guter Export)
- ▶ PDF -> RTF (Formatierung kaum übernommen)
- ▶ PDF -> Text (Sehr guter Export, aber ohne Formatierung)
- ▶ PDF -> Einzelne Bilder in jeweils eine separate Bilddatei exportieren

Quelle: https://helpx.adobe.com/de/acrobat/using/file-format-options-pdf-export.html#file_format_options_for_pdf_export, 27. März 2017