

HTML の文書構造と表示制御を分離するコンバータの開発

松岡 幸典 志田 晃一郎 横山 孝典
 武蔵工業大学

1 はじめに

Web 上に公開する文書の標準形式となっている HTML (HyperText Markup Language) [1] は, その策定組織である W3C (World Wide Web Consortium) によって仕様改定が続けられてきた. その経過の中で表示上の見栄えを設定する記述が策定されたが, 現在 W3C はその機能には問題があるとして使用を非推奨とし, 徐々に言語仕様から削除していく方向へ向かっている.

HTML の表示制御機能の問題は, ユーザがその制御を自由に行うことが出来ない点にある. 使いようによってアクセシビリティを大きく低下させる要因となる.

上記の問題の解決のため W3C は, HTML の記法を用いずに表示制御を行うための言語として CSS (Cascading StyleSheet) [2] を勧告した. 外部 CSS を用いて表示制御を行えば, 図 1 に示すようにユーザが任意に CSS を切り替えることによって表示の変更が容易になる.

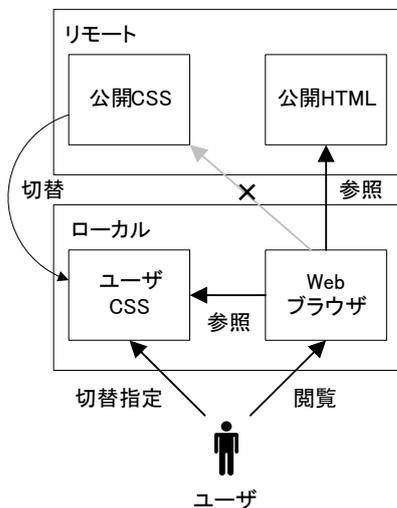


図 1 ユーザによる CSS の切替

続いて W3C は, CSS の有効な使い方や HTML の記述についてのガイドライン (WCAG : Web Content Accessibility Guidelines 等) を策定し, それを基準とする HTML や CSS

HTML-CSS Converter for dividing document structure from presentation control

Yukinori Matsuoka, Koichiro Shida and Takanori Yokoyama,
 Musashi Institute of Technology

の評価ツールを公開した. このような HTML の検査ツールは現在さまざまな組織によって開発, 公開されている.

また, HTML の文法やアクセシビリティの検査を行い, 部分的な修正を行うツール Tidy [3] がある. Tidy は HTML の文法的な誤りや推奨されない記述のうち, 修正方法が自明なものについての自動修正を行うツールである. しかし, HTML の表示制御の記述を CSS に置き換える機能は限定的であり, 一部のものにしか対応していない.

そこで本研究では, 表示制御に用いられる W3C 非推奨タグ・属性について, Tidy 未対応部分を含め CSS の代替記述に自動で置き換えるコンバータの作成を行う.

2 ツール概要

本研究で試作するコンバータでは, 表 1 に示すタグ・属性を CSS の代替記述に置き換える. これらは, W3C が非推奨とするタグ・属性のうち, 同等の表示をするための代替記述が一意に特定できるものである. 非推奨であるタグや属性の種類は, HTML の仕様を定義する DTD [4] に指定されている.

本研究では, W3C 非推奨のタグや属性であっても, その代替となる記述方法に主要な Web ブラウザが対応していない場合は, 変換を行わない仕様とした.

表 1 主要な対応タグ・属性

タグ	属性	Tidy の対応
BODY, TABLE	bgcolor	
BODY	background	
BODY	text	
BODY	link	
BODY	alink	
BODY	vlink	
CENTER	(タグ全体)	○
FONT	(タグ全体)	○
TABLE, TH, TD,	align	
IMG, HR, P	align	
TH, TD	height	
TH, TD	width	
TH, TD	bgcolor	
IMG	border	
IMG	vspace	
IMG	hspace	

また、本研究では変換対象は文法に違反のない HTML に限る。

3 実装

本ツールの処理手順は

- 1) HTML 構文解析
- 2) 非推奨記述検知
- 3) 非推奨記述を HTML から CSS へ変換
- 4) 非推奨記述を HTML から削除または置換
- 5) 変換後 HTML・CSS 出力

となる。

1) の機能については、本ツールでは HTML の構造に関し行う処理は多くないので簡易的な木構造を構築するまでに留めた。また、この処理を行う過程で同時に 2) の処理を行うことが可能である。

3) の処理では、HTML と CSS の文法が異なるために単純なファイル間の文字列移動では不足である。変換の際、読み込んだタグの名前や属性を一度 CSS を概念的に表したクラスとして表現し、そこから CSS の文字列を生成する方式をとった。

4) の機能について、本ツールでは変換元の HTML に対して無用の編集を行ってしまうことを避けるために、HTML 要素のデータ構造から文字列を生成しなおすのではなく、元の HTML 文字列に対してテキスト処理を施すことでタグや属性の削除・置換を行う。

5) の機能については、本ツールでは CSS の記述形態の中で、HTML の外部ファイルとして CSS を参照するようにしている。そのため、HTML の HEAD 要素内に、参照する CSS ファイルを指定するためのリンク記述を追加している。Tidy では、変換した CSS は検査した HTML の内部に埋め込む形で出力しているが、複数の HTML に同一の表示制御を適用させる場合などにおいては、外部ファイルに記述するほうが利便性が高い。

4 評価実験

サンプルデータは、実際に公開されている Web サイトのトップページ数点と、極端な例として HTML の表示制御のみを用いてレイアウトを行った HTML を用いた。これらの HTML から、変換を行う前のもの、Tidy による変換を行ったもの、本ツールによる変換を行ったもの、Tidy と本ツールを組み合わせ変換したものの複数のパターンを既存の検査ツールにかけて、その結果から本ツールの行う変換の有効性を確認した。

また、変換を行う前と後の HTML をそれぞれ Web ブラウ

ザで表示した際に表示崩れが起こらないことについても確認した。

5 今後の課題

現在は置き換えた後の CSS の個々のスタイルの名前（セレクタ）に意味の無い固定の文字列を当てている。これでは CSS の利便性を完全に生かしているとは言えない。この点についての改良を進めたい。

また、HTML の次世代フォーマットとして、XML に準拠した XHTML が勧告され徐々に浸透している。それに伴い、スタイルシートも CSS に限らず XSL が使用できるようになった。これらの新しい規格への対応も考慮していく。

参考文献

- [1] "HTML 4.01 Specification", W3C, <http://www.w3.org/TR/html401/>
- [2] "Cascading Style Sheets, level 2 CSS2 Specification", W3C, <http://www.w3.org/TR/REC-CSS2/>
- [3] "Clean up your Web pages with HTML TIDY", Dave Raggett, <http://www.w3.org/People/Raggett/tidy/>
- [4] "HTML 4 Document Type Definition", W3C, <http://www.w3.org/TR/html401/loose.dtd>