

The hOCR Microformat for OCR Workflow and Results

Thomas M. Breuel
U. Kaiserslautern and DFKI
Trippstadter Str. 122
67663 Kaiserslautern
Germany
tmb@iupr.net

Abstract

Large scale scanning and document conversion efforts have led to a renewed interest in OCR systems and workflows. This paper describes a new format for representing both intermediate and final OCR results, developed in response to the needs of a newly developed OCR system and ground truth data release. The format embeds OCR information invisibly inside the HTML and CSS standards and therefore can represent a wide range of linguistic and typographic phenomena with already well-defined, widely understood markup and can be processed using widely available and known tools. The format is based on a new, multi-level abstraction of OCR results based on logical markup, common typesetting models, and OCR engine-specific markup, making it suitable both for the support of existing workflows and the development of future model-based OCR engines.

1 Introduction

Over the history of OCR systems, a significant number of formats have been proposed for representing the output of OCR systems. We can distinguish three major classes of OCR output formats: *logical formats*, suitable for direct use of OCR results by end users (RTF, HTML, LaTeX, and Microsoft Word), *OCR engine-specific formats* [5, 9], and *benchmarking formats* [11, 10] proposed for benchmarking various aspects of OCR systems. Many OCR engines support multiple output formats and represent information (like text) redundantly within those formats, and commercial and research engines each have their own, mutually incompatible formats.

We are currently developing a multi-lingual OCR system that is targeting the major writing systems and languages of the world, and is intended to deal

with a wide variety of layouts. In addition, our system will have to integrate into existing workflows and text databases. In choosing an output format for our system, we evaluated many of the existing formats and encountered numerous limitations. For example, existing formats usually have limited multi-lingual capabilities, lack support for many common typographic phenomena, lack a well-defined underlying page and typesetting model, and often force the use of separate formats for final output and intermediate results.

Faced with these and other limitations of existing formats, it was important to find an alternative solution. To this end, we defined a number of requirements:

- the ability to cope with common typographic and layout phenomena in the major languages and writing systems
- the ability to support all processing stages in an OCR system with a single format
- the ability to represent layout information in terms of a generative typesetting model
- the ability to encapsulate existing OCR engine output formats
- the ability to represent and associate information from different stages of OCR processing

Furthermore, ideally, we wanted to reuse as much of existing formats and tools as possible, and we wanted a format that permitted future extensibility. As we will see, our format fulfills all these requirements.

2 A Microformat

Our approach to addressing the various requirements without creating a complex new format from scratch has been to adopt HTML/XHTML [2] as

the basis format, together with CSS (cascading style sheets) [4, 8] for representing typographic markup, and to enhance this format by embedding additional information using facilities of standard HTML. One of the key benefits of using HTML+CSS as the foundation of is the hundreds of man-years that has been invested into those formats in identifying typographic and linguistic phenomena across a wide variety of languages and scripts, and defining markup to represent those phenomena. For example, HTML and CSS provide support for representing fonts, styles, hyphenation, flexible spacing, justification, kashida (flexible Arabic characters), Urdu ligatures, Japanese ruby, mixed horizontal/vertical layout, inline changes in writing direction, and many others. None of the existing OCR formats address this wealth of linguistic and typographic phenomena.

Existing standards-conforming HTML output from an OCR system is already minimally compliant OCR output in our system, albeit without any OCR specific information (e.g., geometric information). Adding OCR specific information to this format then means incorporating a few tags to the HTML output that indicate geometric and OCR-specific information. These additional tags do not alter the appearance of the output, are fully compliant with the HTML standard, and are processed and preserved when the HTML is processed using standard tools.

This is accomplished by using the `DIV` and `SPAN` tags; these tags have no specific meaning in HTML, but they may be used to associate style and other information with regions of text inside an HTML document. Both tags allow a small number of standard attributes, including the `class`, `style` and `title` attribute. The `class` attribute is used for identifying a tag as belonging to a particular class or application; we use this attribute to identify tags representing OCR information. In particular, all OCR-related tags have a class of `ocr_` or `ocrx_`.

The `style` attribute is used for associating style information with regions of text; we are using it directly to associated typographic and layout information with text in the standard way defined by the CSS standard. The `title` attribute may contain arbitrary text, and we use it to encode tag specific information.

For the hOCR format, we will refer to `DIV` or `SPAN` tags that contain OCR-related information as *elements*, and we refer to the information encoded in its `title` attribute as *properties*. In fact, the hOCR format is equivalent to (and can be automatically translated to and from) an XML format, in which the hOCR elements and properties correspond to XML tags and attributes (the size of the two representations is ap-

hOCR format:

```
<div class="ocr_block"
  title="bounds 112 17 213 53">text...</div>
```

XML equivalent:

```
<ocr_block
  bounds="112 17 213 53">text...</ocr_block>
```

Figure 1. hOCR markup is fully equivalent to well-formed XML and can naturally be transformed into XML automatically and without loss of information. (Note that the reverse is not true: arbitrary XML cannot easily be embedded as a format.)

proximately the same). In that sense, none of the actual design of hOCR is tied to its representation as a format, and all the considerations and design going into hOCR apply equally well to a possibly future XML-based format. Note, however, that the reverse is not true: an XML format cannot be automatically and naturally converted to a format, since many XML constructs are difficult to encode naturally and readably inside a format.

3 Supported Markup

Above, we have described the hOCR format in terms of its embedding into HTML. Just as important is what hOCR is actually representing—the content elements. In fact, the content elements could easily form the basis of a non-HTML OCR markup language as well.

3.1 Logical Markup

Most document types have a tree structure, representing nested large scale to small scale divisions of the document. Within that tree structure, the text itself is represented linearly in reading order. Not only is this structure shared, but most document types and markup languages have common, recurring section types, found in document types like memoranda, articles, etc. The hOCR logical markup elements define and represent these most commonly found divisions. Like all hOCR markup, this markup is optional; if it does not fit the particular needs of a document type, OCR systems are free to omit these, or use other embedded formats to represent their specific document structures.

Of course, HTML itself has its own set of hierarchical document structuring elements (H1 through H4,

```

ocr_document
  ocr_linear
    ocr_title
    ocr_author
    ocr_abstract
    ocr_part
      ocr_chapter [H1]
      ocr_section [H2]
      ocr_sub*section [H3,H4]
      ocr_display
      ocr_blockquote [BLOCKQUOTE]
      ocr_par [P]

```

Figure 2. Logical markup available in hOCR. This is markup for the logical hierarchy of a document, independent of where or how on the page it is rendered. This particular markup is usable both for individual documents (memos, articles, etc.), as well as compound documents consisting of multiple, possibly interleaved, texts (newspapers, magazines, collections).

P, etc.). Those structuring elements are recognized by many HTML processing applications and browsers. hOCR recommends a mapping of hOCR structuring elements onto HTML elements that OCR systems can use if there is no other reason to prefer alternative mappings. Note that when a mapping is used, the hOCR elements can be encoded directly in the HTML tags, as in `<p class="ocr_par" title="...">`.

One important element is `ocr_linear`, used for representing output for document types like newspapers. Newspapers generally consist of multiple stories, each with a fairly regular, linear document structure (title, author, sections, paragraphs, etc.) and with text in reading order. However, the different articles of the newspaper themselves do not have a unique reading order; at the top level, the document is a collection of unordered articles, each of which is a self-contained linear document (possibly with cross-references to other documents). This relationship is represented by the `ocr_document` and `ocr_linear` elements.

3.2 Typesetting Markup

In the previous section, we have discussed markup for representing common linear document structure, found in many typesetting and word processing systems. In order to generate paginated, printed output, typesetting systems and word processors use a nearly universal typesetting model that divides each page

```

ocr_page
  ocr_carea ("content area")
    ocr_line [SPAN]
    ocr_float
      (subclasses for images, formulas, etc.)
    ocr_separator
    ocr_noise

```

Figure 3. Typesetting-related elements represent blocks and floats that are filled with the document content in reading order in standard typesetting models.

into blocks and floats ([1], [6], Microsoft Word, Adobe Framemaker). Typesetting is performed by flowing the linear, logically marked up text into these boxes in reading order. The combination of logical markup and typesetting markup permits us to use hOCR as an intermediate format for performing OCR as model-based *reverse typesetting*, an approach advocated, for example by Kopeck and Chou [7].

3.3 Engine-Specific Markup

At the lowest level, the hOCR format represents OCR engine-specific, physical layout, like text blocks, images, and other page content. This kind of physical markup is the most commonly found output of existing OCR systems. Common elements include “text blocks”, “figures”, “lines”, and “words”.

However, unlike typesetting markup, which is generally well-defined in terms of typesetting models, the kind of physical markup produced by OCR engines is implementation dependent. For example, a “text block” in an engine may be defined in terms of the existence of whitespace separators of minimal size, or the alignment of individual characters. Likewise, “blocks” are often also style-dependent; for example, a document rendered in a style with vertical inter-paragraph spacing may be represented with a single block for each paragraph, while in the same document rendered in a different style, an entire column of multiple paragraphs may be returned as a single block in the OCR system. In contrast, in a typesetting model of page layout, these two styles would be represented in the same way as a flowable content area, which would also correspond to the underlying page layout in the source document in any of the standard typesetting systems.

Another source of implementation dependencies is the definition of words. OCR systems sometimes define “words” in terms of spacing, at other times linguisti-

cally.

In order to indicate that engine-specific markup is not portable between OCR engines, its elements are prefixed with `ocrx_` instead of `ocr_`. Common physical markup tags are `ocrx_block`, `ocrx_line`, and `ocrx_word`. Authors are free to extend the set of `ocrx_` elements as needed.

Given that physical markup is OCR engine specific, the question arises why one would want to represent it at all. The reason is that many existing workflows already rely on this kind of markup. For example, a specific workflow may run the Abbyy or Textbridge OCR engine and then perform rule-based post-processing of the physical layout information produced by those systems. By having a portable, OCR-engine independent representation of these features, we can carry forward engine-specific information and yet use portable tools to manipulate the content. Furthermore, although this kind of markup is engine-specific, there are many commonalities so that code that works for one engine can, after some testing and modification, be used with another engine.

3.4 Character Recognition

We have already mentioned above that one advantage of the hOCR format over other OCR formats is that it can reuse the expertise that has gone into the development of textual representations for HTML. Generally speaking, all common style-, font-, script-, language-, and typesetting-specific phenomena (hyphenation, spacing, ruby, kashida, etc.) are to be represented using their HTML, CSS, and Unicode representations.

One area where existing HTML does not provide an obvious representation is in the representation of segmentation and recognition alternatives and weights. Generally, segmentation alternatives are represented using the HTML `INS` and `DEL` tags, with associated classes and additional information encoded in the elements. In terms of recognition confidences and geometric segmentation information, hOCR provides two sets of properties: engine-specific properties (commonly, “confidence scores” and “bounding boxes”), and engine-independent properties (posterior probabilities and character cutting paths).

3.5 Meta Information

An important requirement for hOCR files is that their headers contain information about what kind of processing was used in order to obtain the output. This information consists not only of the usual OCR

```
import libxml2, re, os, string

doc = libxml2.parseFile('doc.xhtml')
lines = doc.xpathEval("//*[ @class='ocr_line' ]")
for line in lines:
    textnodes = line.xpathEval("./text()")
    for node in textnodes:
        print node.getContent()
```

Figure 4. Example of processing hOCR markup in Python. This outputs the text contained in the document line-by-line. Simple processing of hOCR documents is enabled by relying on standard DOM functions, available for HTML processing.

meta information (software version numbers, scan resolution, number of pages, processing host, etc.), but also includes headers that indicate the capabilities of the OCR engine. In particular, the hOCR meta information must indicate which kind of markup the engine is capable of generating, regardless of whether the actual document contains that markup. Without this meta information, it is impossible to tell whether the lack of a specific markup in an OCR output file (e.g., paragraphs, language identification) is due to the absence of the corresponding markup in the source document, or due to the inability of the OCR engine to recognize that form of markup. OCR system capabilities are indicated using HTML meta tags simply by listing the element names and attributes the OCR system is, in principle, capable of generating. Inclusion of this meta information is therefore an important part of an OCR system independent OCR format.

3.6 Other Content

Let us note briefly that other common typesetting objects already have common representations for their embedding into HTML: MathML for mathematics, ChemML for chemical formulas, the `IMG` tag for bitmapped images, and SVG for vector graphics.

In addition, there is a wealth of additional embedded formats already available for encoding information such as resume structures, addresses, identity, personal relationships, and bibliographic references. Furthermore, the hOCR format itself may be extended, or separate, additional OCR-related embedded formats may be defined for handling special needs such as indexing, table of contents, etc.

4 Tool Support

The hOCR format has been developed over the last 12 months and we have gathered significant hands-on experience with it. It has been used by different groups as a workflow format representing the output of commercial OCR systems, as the output format of a newly developed OCR system (to be released in open source format), and as the format for the release of OCR ground truth for research in OCR. Both integration into existing HTML output routines of OCR engines, as well as development of new hOCR generators from scratch has turned out to be easy.

The ability to use tools intended for standard HTML with hOCR has turned out to be very useful. Output in hOCR format can be viewed and edited using standard HTML viewers and HTML editors. It can also be processed using standard DOM processing tools (see Figure 4). Such processing generally preserves the hOCR information and presents all levels of OCR output (content, physical, typesetting, logical) within a single, consistent document.

The specific properties of hOCR have also made it possible to create general-purpose tools that work for many different OCR engines. For example, post-processing of the output of commercial OCR systems (e.g., deriving logical layout information from engine specific output, replacing engine-specific fonts with standard HTML fonts) can be implemented using tools that transform hOCR input into hOCR output and only require small amounts of adaptation to different OCR engines. We have also begun creating a standard suite of evaluation tools that measure OCR accuracy and the performance of different levels of layout analysis directly on hOCR-encoded OCR output.

5 Conclusions

The hOCR format introduces several distinct concepts in OCR output formats and workflows: the use of HTML-based embedded formats for the representation of OCR outputs, the use of typesetting models in the design of OCR formats, and the use of a single format for representing results of all levels of OCR within a single output file. The format has proven itself for about a year in actual use now, and it will be the basis of the OCRopus [3] open source OCR system and an upcoming release of ground truth data from scanned books and magazines. We hope that a universal, extensible standards-based OCR format, together with ground truth data, OCR software, and various tools will be persuasive enough to induce both commercial vendors and researchers to standardize on a single pro-

cessing and interchange format. A full specification of the hOCR format and sample documents are available at xxxxxx.

References

- [1] Extensible stylesheet language (xsl) version 1.1. <http://www.w3.org/TR/xsl/>, 2001.
- [2] HTML/XHTML Specifications. <http://www.w3.org/MarkUp>, 2006.
- [3] The ocropus ocr system home page, 2007.
- [4] Bert Bos, Hakon Wium Lie, Chris Lilley, and Ian Jacobs. Cascading Style Sheets, level 2: CSS2 Specification. <http://www.w3.org/TR/REC-CSS2/>, 1998.
- [5] Daniel S. Connelley and Beth Paddock. Xdoc data format: Technical specification. Xerox Imaging Systems part no. 00-07571-00, 2000.
- [6] Donald Knuth. *TeXbook*. Addison-Wesley, 1984.
- [7] G. E. Kopec and P. A. Chou. Document image decoding using markov source models. *IEEE PAMI*, 16(6), 1994.
- [8] Eric A. Meyer and Bert Bos. Introduction to CSS3. <http://www.w3.org/TR/css3-roadmap/>, 2001.
- [9] Nuance, Inc. SSDOC-SCHEMA2 (XML Schema). <http://www.scansoft.com/omnipage/xml/SSDOC-SCHEMA2.xml>, 2006.
- [10] RAF Technologies, Inc. (Redmond, WA). Illuminator user’s manual. <http://documents.cfar.umd.edu/resources/source/illuminator.html>, 1995.
- [11] B. Yanikoglu and L. Vincent. Pink panther: a complete environment for ground-truthing and benchmarking document page segmentation. *Pattern Recognition*, 31:1191–1204, 1998.