# Learning to Extract Semantic Structure from Documents Using Multimodal Fully Convolutional Neural Networks

## SUPPLEMENTAL MATERIAL

Xiao Yang[‡], Ersin Yumer[†], Paul Asente[†], Mike Kraley[†], Daniel Kifer[‡], C. Lee Giles[‡]
[‡]The Pennsylvania State University    [†]Adobe Research
xuy111@psu.edu  {yumer, asente, mkraley}@adobe.com  dkifer@cse.psu.edu  giles@ist.psu.edu

## 1. Synthetic Document Data

We introduced two methods to generate documents. In the first method, we generate LaTeX source files in which elements like paragraphs, figures, tables, captions, section headings and lists are randomly arranged using the "textblock" environment from the "textpos" package. Compiling these LaTeX files gives single, double, or triple-column PDFs. The generation process is summarized in Algorithm 1.

---

**Algorithm 1** Synthetic Document Generation

---

1: $s \leftarrow$ a string containing preamble and necessary packages of a LaTeX source file
2: Select a LaTeX source file type $T \in \{$single-column, double-column, triple-column$\}$
3: **while** space remains on the page **do**
4:     Select an element type $E \in \{$figure, table, caption, section heading, list, paragraph$\}$
5:     Select an example $e$ of type $E$
6:     $s_e \leftarrow$ a string of LaTeX code that generates $e$ using the "textblock" environment
7:     $s \leftarrow s + s_e$
8: **end while**
**Output:** $s$
**Output:** A PDF document after compiling $s$

---

Elements in a document are carefully selected following the guidelines below. Figure 1 shows several examples of the figures and tables used in the synthetic data generation.

- Candidate figures include natural images from MS COCO [3], academic-style figures and graphic drawings downloaded using web image search.

- Candidates tables include table images downloaded using web image search. Various queries are used to increase the diversity of downloaded tables.

- For paragraphs, we randomly sample sentences from a 2016 English Wikipedia dump [2].

- For section headings, we sample sentences and phrases that are section or subsection headings in the "Contents" block in a Wikipedia page.

- For lists, we sample list items from Wikipedia pages, ensuring that all items in a list come from the same Wikipedia page.

- For captions, we either use the associated caption (for images from MS COCO) or the title of the image in web image search, which can be found in the span with class name "irc_pt".

In the second document generation method, we collected and labeled 271 documents with varied, complicated layouts. We then randomly replaced each element with a standalone paragraph, figure, table, caption, section heading or list generated as stated above. Figure 2 shows several examples from the 271 documents.

## 2. Visualizing the Segmentation Results

Each pixel $p$ in the model's output layer is assigned the color of the most likely class label $l$. The RGB value of that color is then weighted by the probability $P(l)$.

## 3. Post-processing

We apply an optional post-processing step to clean up segment masks for documents in PDF format. First, we obtain candidate bounding boxes by using the auto-tagging capabilities of Adobe Acrobat [1] and parsing the results. Boxes are stored in a tree structure, and each node's box can be a TextRun (a sequence of characters), TextLine (potentially a text line), Paragraph (potentially a paragraph) or Container (potentially figures or tables). Note that we

Figure 1: Sample figures and tables used in synthetic documents generation. (1) Natural images from MS COCO dataset. (2) Academic-style figures from web image search. (3) Symbols and graphic drawings from web image search. (4) Tables from web image search.
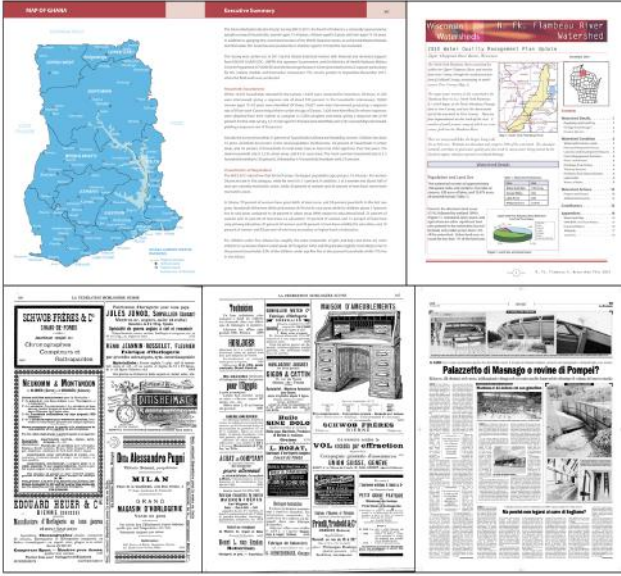
Figure 2: Examples of documents with complicated layout. We labeled regions in each document and then randomly replaced them with a standalone paragraph, figure, table, caption, section heading or list, as described in Sec. 1

ignore the semantic meanings associated with these boxes and only use the boxes as candidate bounding boxes in post-processing. Figure 3 (2) and 4 (2) illustrate candidate bounding boxes for each document.

---

**Algorithm 2** Segmentation Post-processing

**Input:** $P \leftarrow$ probability map, $P(u, v) \in R^{|C|}$ is a vector containing the probability of each class $c \in C$ at location $(u, v)$
**Input:** $Boxes \leftarrow$ candidate bounding boxes
 1: $S \leftarrow$ segmentation to be generated
 2: **for** each location $(x, y) \in S$ **do**
 3:     $S(x, y) \leftarrow$ background
 4: **end for**
 5: **for** each $b \in Boxes$ **do**     ▷ parent box comes before child boxes
 6:     $\bar{p} \leftarrow \sum_{(u,v) \in b} P(u, v)$
 7:     $l \leftarrow \arg\max \bar{p}$
 8:     **for** each location $(u, v) \in b$ **do**
 9:        **if** S(u, v) is background **then**
10:          $S(u, v) \leftarrow l$
11:        **end if**
12:     **end for**
13: **end for**
**Output:** $S$

---

Using these bounding box candidates, we refine the segmentation masks by first calculating the average class probability for pixels belonging to the same box, followed by assigning the most likely label to these pixels. The process is summarized in Algorithm 2.

## 4. Additional Visualization Results

Figures 3 and 4 show additional visualization examples of synthetic documents, and Figure 5 shows additional examples of real documents.

## References

[1] Adobe Acrobat. http://www.adobe.com/accessibility/products/acrobat.html. 1

[2] Wikipedia. https://dumps.wikimedia.org/. 1

[3] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 1
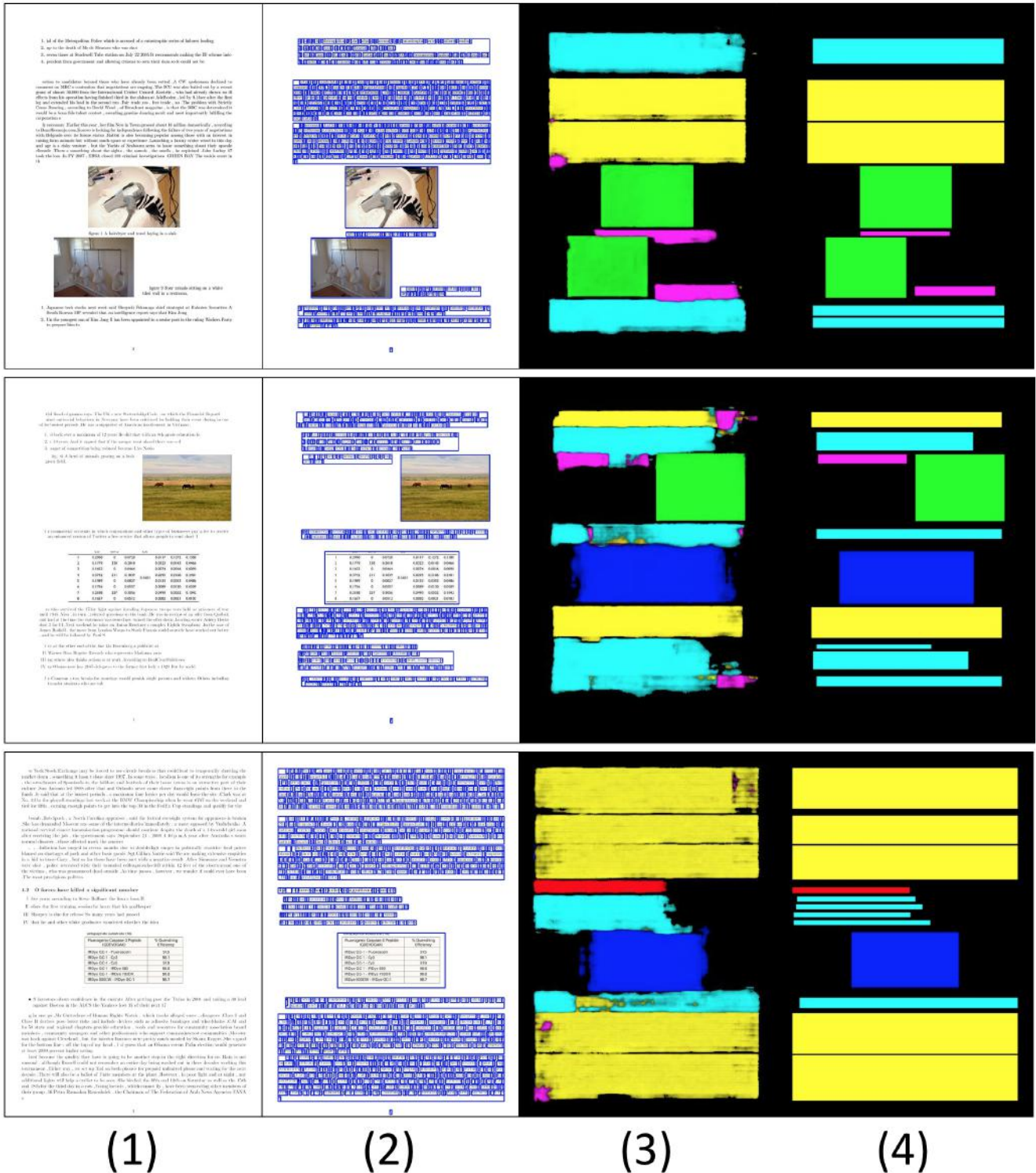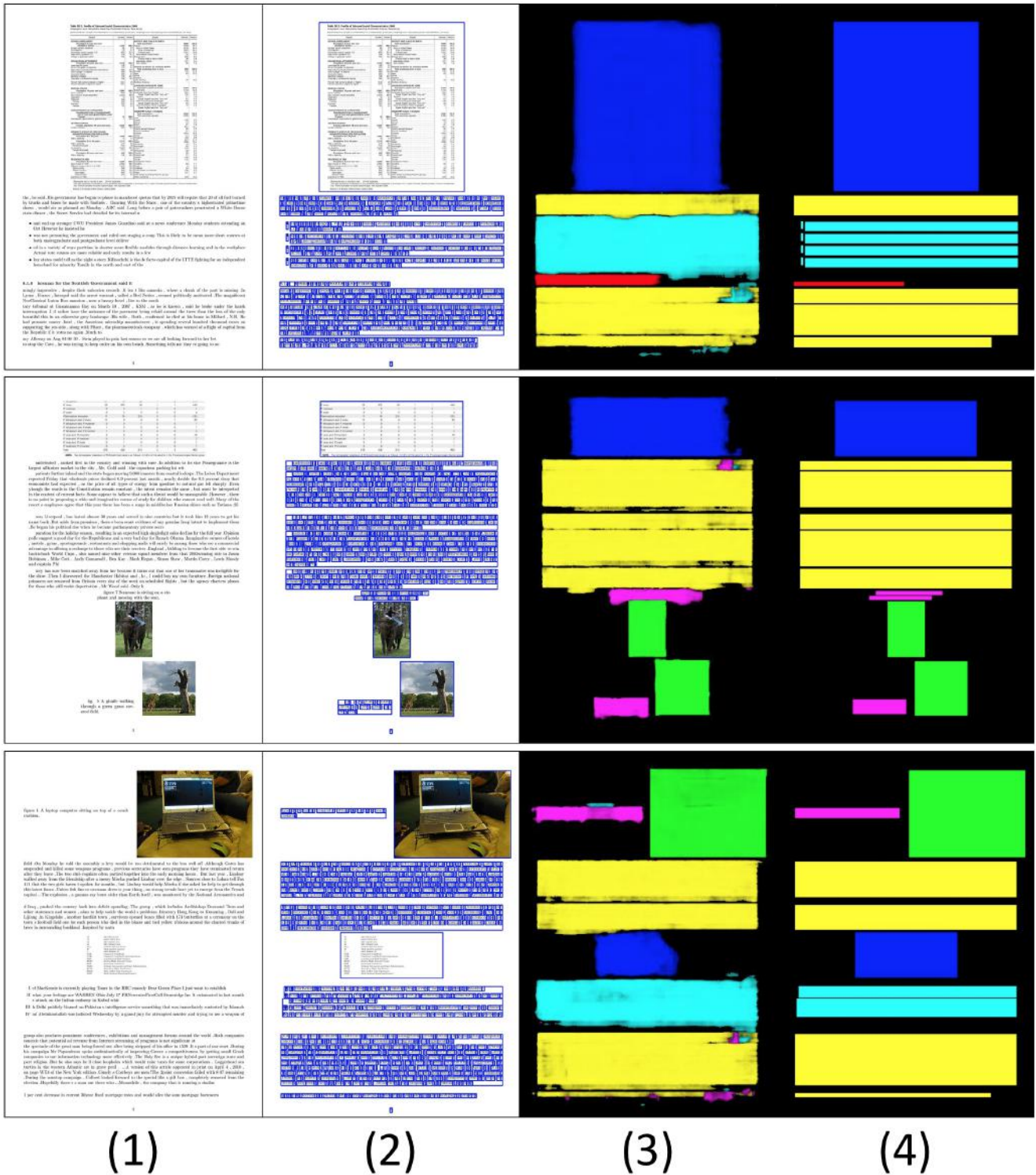
Figure 3: Synthetic documents and the corresponding segmentations. (1) Input synthetic documents. (2) Candidate bounding boxes obtained by parsing the PDF rendering commands. (3) Raw segmentation outputs. (4) Segmentations after post-processing. Segmentation label colors are: **figure** , **table** , **section heading** , **caption** , **list** and **paragraph** .

Figure 4: Synthetic documents and the corresponding segmentations. (1) Input synthetic documents. (2) Candidate bounding boxes obtained by parsing the PDF rendering commands. (3) Raw segmentation outputs. (4) Segmentations after post-processing. Segmentation label colors are: **figure**, **table**, **section heading**, **caption**, **list** and **paragraph**.
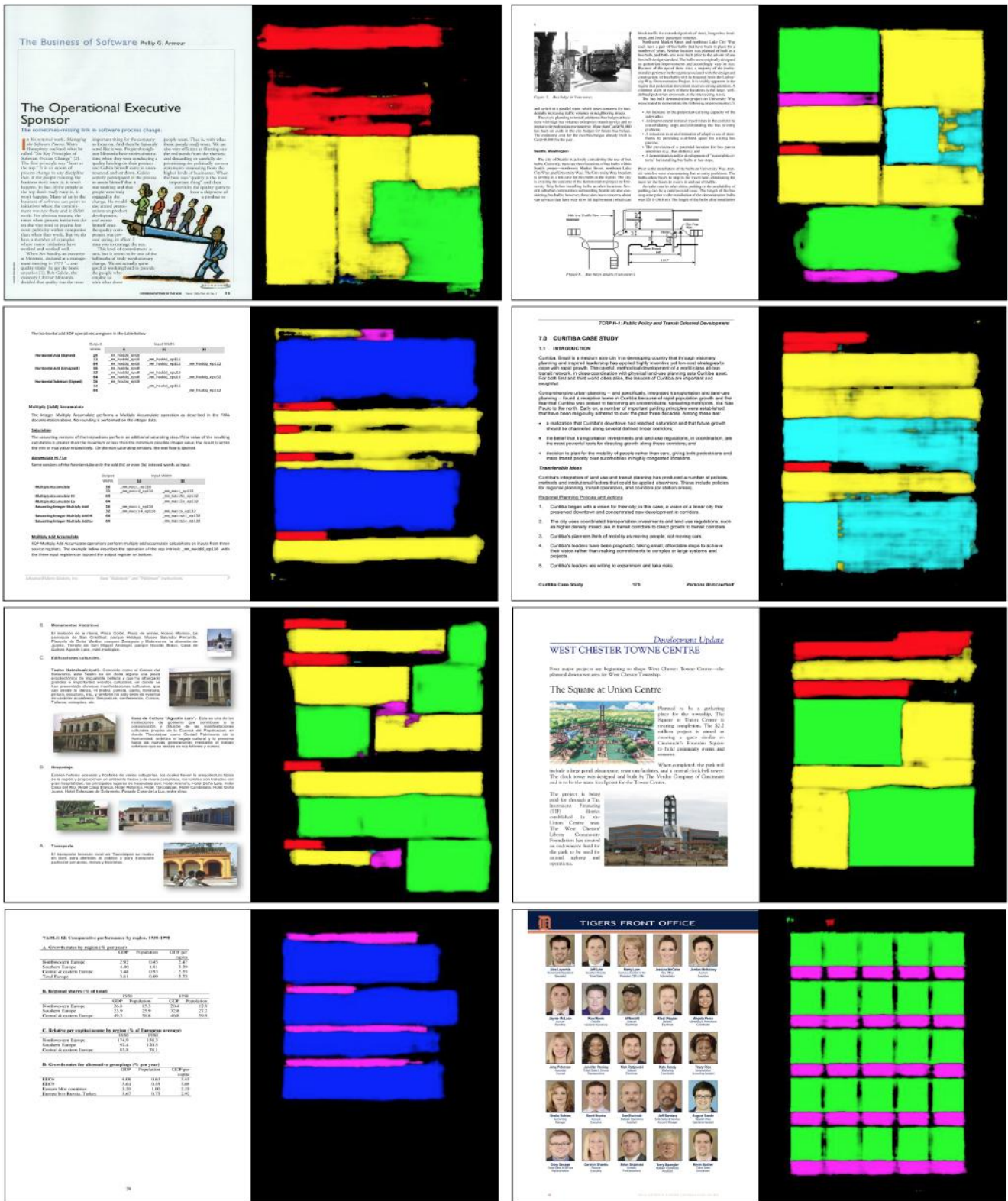
Figure 5: Real documents and the corresponding segmentations. Segmentation label colors are: **figure**, **table**, **section heading**, **caption**, **list** and **paragraph**.