

通过Hybrid-Top检测进行混合页面布局分析 Hybrid Page Layout Analysis via Tab-Stop Detection

贾史密斯

Google Inc. 1600 Amphitheatre Parkway, Mountain View, CA 94043, USA.

theraysmith@gmail.com

Abstract

一种新的混合页面布局分析算法是提出的。该方法使用自下而上的方法来形成初始数据类型的假设，并找到格式化页面时使用的制表位。检测到的制表位用于推断页面的列布局。然后以自上而下的方式应用列布局，将结构和阅读顺序强加到检测到的区域上。

完整的C++源代码实现 + source code implementation is available as part of the Tesseract open source OCR engine at <http://code.google.com/p/tesseract-ocr/>.

1. Introduction

物理页面布局分析，第一步是OCR的功能是将图像分为文本和非文本区域，以及将多文本分成多个列。本文不涉及逻辑布局分析，而是检测页面，检测标题、编号列表和文章分段、body text, numbered lists, and segmentation into articles.

物理页面布局分析对于实现OCR引擎以处理任意页面，例如书籍、杂志、报纸、信件和报告。物理布局分析的方法大致分为两类：分析分析 fall roughly into two categories:

自下而上的方法是最早的方法，甚至更多。最近发表的方法[1]，他们将图像的一小部分（像素、像素组或连接的组件）分类，并像类型一样聚集在一起以形成区域。自下而上的方法的主要优点是它们可以轻松处理任意形状的区域。关键缺点是它们努力考虑图像中的高级结构，例如列。这通常导致过度碎片化的区域。

另一种方法以递归方式切割图像。沿着空白的垂直和水平方向，这些空白应该是列边界或段落边界。尽管自上而下的方法具有

优点是它们从查看页面上最大的结构开始，因此无法处理许多杂志页面上出现的各种格式，例如非矩形区域和交叉列标题无缝地融合到下面的列中。 headings that blend seamlessly into the columns below.

第三种方法[5-7]是基于分析图像中的空白。通过对间隙进行自底向上的分析，在列之间查找间隙，从而明确查找白色矩形。从而解决了递归自顶向下方法中的某些缺陷。这些算法 algorithms mostly still suffer from the problem of being unable to handle non-rectangular regions.

2. Page layout via tab-stop detection

有些页面是由某些人 laid out, either by a professional 发布系统或公共文字处理程序。页面区域由制表符 regions of a page are bounded by tab-stops. The margins, column edges, indentation, and columns of a 限制、边距、列边缘、缩进和可表格都放置在固定的x位置。在这些位置上文本行的边缘或中心垂直对齐。制表位停止将表格与正文分离。他们也限制了矩形的排列 from body text, and they also bound rectangular non-column elements, such as inset images and pull-out quotes.

制表符在 tab-stops in the 图1的示例是带有用于页面缩进的附加制表符的列 column boundaries with an additional tab-stop for the 边界。这对于查找页面布局不是必需的。非矩形插图通常会偏离列边界。 that is not required for finding the page layout. The non-rectangular inset image, typically, strays outside of the column boundaries.

从某种意义上说，白色矩形与制表符匹配，但背景矩形或背景图像可能会干扰白色矩形。也是白色的两端 or background images. Also the ends of white 矩形与制表位限制的区域的不端不匹配。因为白色短形一直延伸到垂直空白中。 into the perpendicular whitespace.



Fig. 1 Input image.

The proposed algorithm is similar to the whitespace rectangle method, it uses bottom-up and top-down methods to find a top-down structure, but instead of finding the space between columns, it looks for the 'tab-stops' that mark its boundary cells, and combines bottom-up and top-down methods. The proposed algorithm copes easily with non-rectangular regions.

There are two main phases: preprocessing, in which a bottom-up shape analysis and component analysis formed the local data types, initial hypotheses over the local data types, bottom-up tab-stop detections, finding the tab layout, finally applying tab-stops to create an ordered set of typed regions. These phases will be detailed in sections 3-6.

3. 预处理

The aim of the preprocessing step is to identify line separators, image regions, and will separate the likely connected components into likely text components and a smaller number of uncertain type.



Fig.2. (a) Vertical lines, (b) Image elements.

8. 如图 2(a) 所示, 将图 1 中的图像 Fig. 1, 的形态学处理检测得到图 2(a) 所示的垂直线, 如图 2(b) 所示的图像蒙版。在将清洗后的图像传递到连接的分量分析之前, 从输入图像中减去这些检测到的元素, 得到图 2(c) 所示的图像。图 2(c) 所示的图像是图 2(a) 所示的图像减去图 2(b) 所示的图像蒙版后的结果。图 2(c) 所示的图像是图 2(a) 所示的图像减去图 2(b) 所示的图像蒙版后的结果。图 2(c) 所示的图像是图 2(a) 所示的图像减去图 2(b) 所示的图像蒙版后的结果。

The connected components (CCs) are filtered by width, height and area. In this paper, the CCs with size smaller than 300 pixels are removed. The CCs with height $h \leq 7$ (at 300 dpi) are small. The 75th percentile of the heights of the remainder, h_{75} , is computed, and CCs with $h < 2h_{75}$ are small, $h \geq 2h_{75}$ are medium, and the rest are large. The same logic is applied to width and area.

此小写字母是重要的, 因为小写字母 (noise 或变音符号) 和较大的非文本标题 (线条图、徽标或框架, logos, or frames) 是 likely 来混淆的。text-line algorithms, 但是较大的文本标题对于阅读 可能会混淆文本线算法, 但是较大的文本标题对于阅读 顺序检测很重要。如果左右或邻居的笔触宽度相似, 则在 reading order detection. Large CCs are considered text 此阶段将大型 CC 视为文本。在“强调”字体上, 笔触 has a similar stroke width. On “stressed” fonts, the stroke

垂直线上的宽度大于水平线上的宽度。因此笔划宽度是在两个方向上分别计算的。笔划宽度是根据 width 进制图像上距离最近的水平和垂直局部最大值 (local maxima) 的距离值来计算的。图 3 显示将 CC 过滤为中文或英文文本。CCs are filtered as medium or large text.



Fig.3. 过滤的炒送 CCs

4. Finding tab positions as line segments

寻找制表位的过程包括以下几个主要步骤：

1. 找到看起来像在文本区域边缘的候选制表位。它们可能位于文本区域的边缘，或者在文本区域内部。然后，将它们分组为制表位线，然后找到制表位线之间的连接，从而消除误报。

4.1.1. Finding candidate tab-stop components

图4(a)展示了候选制表位CC。在候选制表位CC中，可以识别出以下表格找到。初始搜索从预处理中的每个已过滤CC开始。假设CC位于表格处理位置。Assuming that the CC is at a tab-stop, 则搜索将查找对齐的邻居和邻居在排水沟的搜索 looks for aligned neighbors and neighbors in the gutter where there should be a space. Each CC is 每个CC独立处理，并根据其是否为候选在选项卡或选项卡或两者都不标记而进行标记。图4(a) 示出了候选制表位CC。Fig. 4(a) illustrates the candidate tab-stop CCs.

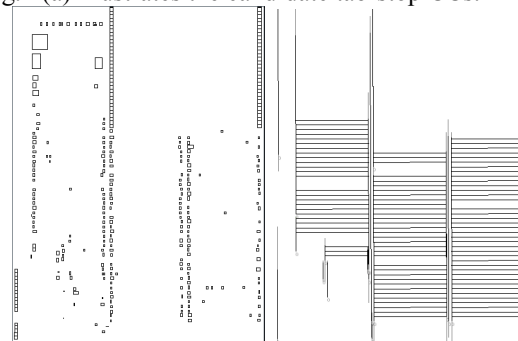


Fig. 4. (a) Candidate tab-stop components
(b) Fitted tab lines and traces connections.

4.2. Grouping candidate tab components

如果组中有足够多的描述, 则保留它们。最小二乘平方中值算法用于将线拟合到组中每个CC的适当(左或右)边缘。找到并拟合一条线到适当的(左或右)边缘 of each CC in a group. 所有线都重新调整至页面均值方向, all the lines are refitted to the page-mean direction.

然后,迭代过程会标记最长的段列候选对象之一(解释的连续失败) page y-coordinates that is explained by one of the column candidates. Fig. 7 shows the result of this process.

6. 寻找地区

找到这些后,为CP赋予一个类型。根据它们跨越多少列,具有单列的CP正在流动,within a single column are *flowing*, partitions that touch more than one column, but do not span to the outer edges of the page are *pull-out*, and partitions that completely span more than one column are *heading*.

6.1. 创建CP流

每个CP选择其最佳匹配的上下伙伴,即垂直最接近的CP重叠在一起。由于每个CP向其水平方向注册自己,每个CP注册其自己,其选择的伙伴注册自己,因此每个CP可能具有零个或多个注册的上下伙伴。

注册伙伴列表的大小限制为注册伙伴的数目。依次使用以下规则将上下限分别设为零或1,类型如下:

1. Type. If there are multiple types, text can only stay with its own (text) type, whereas image can stay with its own (image) type. 2. Transitive partner shortcuts are broken. If A has 2 partners B and C, and also B has C as a partner in the same direction, delete C as a partner of A. 3. (仅文本) 如果A仍然有2个伙伴B和C,请删除B和C的伙伴,看看哪个拥有最长的链。从A删除拥有最短链的伙伴。然后,将最短链的类型转换为拉出。

4. (仅图像) 选择水平重叠最大的伙伴CP, with the largest horizontal overlap.

现在所有CP都具有0或1个伙伴。Even so, (re)run the rules. 这会将文本的所有链简化为单一类型,并将文本链与图像链分开。通过将所有链中的所有CP设置为链中最通用的类型,可以纯化图像链。图8显示了最终建立的CP。其链的初始类型:CPs, 其中流动文本为蓝色,标题文本为青色,标题图像为洋红色,拉出图像为橙色。

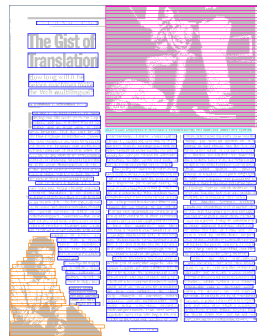


Fig. 8 Typed partition chains.

CPs of text CPs are further divided into groups of uniform line spacing. 使文本块,现在每个CP链代表一个块。每个链的CPs代表一个候选区域,但必须对这些区域进行排序。

6.2. 阅读顺序确定

回忆,图像和文本分区的类型为:1. 流动。流动块在它们接触的列之间。2. 拉出。拉出块在它们接触的列之间。3. 标题。标题跨越多个列,并位于所跨越的列中或其之间的上方的任何内容之后。

4. 列布局的更改就像标题一样,任何更改的列内容之前发生。5. 在标题之间,列的内容从左到右排序。

6. 列布局的更改就像标题一样,任何更改的列内容之前发生。

7. 在标题之间,列的内容从左到右排序。

8. 在标题之间,列的内容从左到右排序。

9. 在标题之间,列的内容从左到右排序。

10. 在标题之间,列的内容从左到右排序。

6.3. 找到每个区域的多边形边界

为了简化实现,多边形边界是等规的,即边在水平和垂直方向上。多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。

多边形边界是等规的,即边在水平和垂直方向上。



Fig. 9 Final blocks.

7. 测试与结果

本文描述的算法在以下位置实现: C++, 并且源代码作为Tesseract开源OCR系统的一部分提供[9,10]。它在3.4 GHz Pentium 4上运行,在典型8MPixel图像上大约需要1秒,在非典型8MPixel图像上运行。

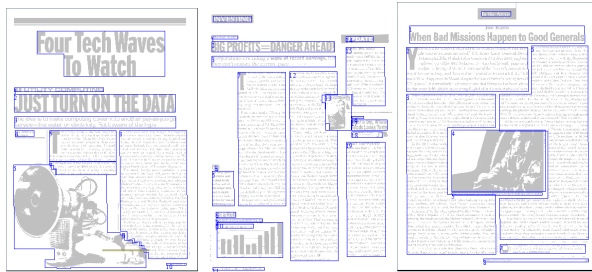


Fig. 10. Results on some of the ICDAR2007 set.

Property testing in page layout analysis is a difficult problem [11]. For complex magazine pages, the UNLV test set [12] only measures text regions, and counts errors unless figure captions are placed after all the body text.

The ICDAR page layout analysis competitions provided a better overall precision. The results of this algorithm appear in the 2009 competition [13]. Some graphical results are shown in Fig. 10. The results in Table 1 are computed on only the 2007 test set, and the author would like to thank Apostolos Antonacopoulos for providing these results. For more information on the testing methodology, see references [11] and [13].

Table 1. Results on the ICDAR 2007 set.

Method	Noise	Sep	Text	Image	Overall
PRIMA Metric					
2007-Besús	86.8%	76.9%	37.4%	42.5%	35.9%
2007-HITP	79.7%	68.0%	79.7%	46.2%	67.6%
2007-HF	79.6%	79.6%	79.6%	48.4%	65.7%
Tesseract	74.1%	65.6%	74.1%	55.3%	68.4%
F-measure					
2007-Besús	62.9%	76.2%	95.8%	57.2%	90.2%
2007-HITP	80.7%	79.9%	80.7%	72.1%	88.2%
2007-HF	80.7%	80.6%	80.6%	72.4%	88.6%
Tesseract	70.9%	70.9%	70.9%	82.0%	91.3%
Recall					
2007-Besús	65.7%	71.7%	94.9%	67.0%	88.2%
2007-HITP	79.5%	69.5%	94.9%	66.4%	89.8%
2007-HF	79.5%	69.5%	94.9%	66.9%	90.2%
Tesseract	65.6%	65.6%	94.9%	76.5%	93.8%
Precision					
2007-Besús	60.4%	81.3%	96.7%	50.0%	92.2%
2007-HITP	81.9%	81.9%	87.4%	79.0%	86.7%
2007-HF	81.7%	81.7%	87.9%	79.0%	87.0%
Tesseract	62.8%	82.8%	89.0%	88.3%	88.9%

10. Conclusion and further work

Table stops make an interesting and useful alternative to white rectangles for finding the column structure of a

page. Combining the top-down concept of column layout analysis with the bottom-up method of column layout analysis to easily handle the complex non-rectangular layouts of modern magazine pages without using bottom-up methods alone.

The algorithm described has no table detection or analysis, but table detection is a useful feature for both, so table analysis will be added in the future.

11. References

- [1] F. Wan, W. Wang, W. Gao, R. Casey, "Block segmentation and text extraction in mixed-text documents," *Computer Graphics and Image Processing*, vol. 62, no. 3, pp. 359-370, 1992.
- [2] J. W. F. Cheng, "A new method for document layout analysis," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 33, no. 4, pp. 401-408, 2003.
- [3] S. P. Chowdhury, S. Mandal, A. K. Das, B. Chanda, "Segmentation of Text and Graphics from Document Images," *Proc. 9th Int. Conf. on Document Analysis and Recognition*, vol. 1, pp. 619-623, 2007.
- [4] H. S. Baird, S. E. Jones, S. J. Fortune, "Image Segmentation by Shape- and Size-based Clustering," *Proc. 5th Int. Conf. on Document Analysis and Recognition*, vol. 1, pp. 347-349, 1984.
- [5] H. S. Baird, S. E. Jones, S. J. Fortune, "Image Segmentation by Shape- and Size-based Clustering," *Proc. 5th Int. Conf. on Document Analysis and Recognition*, vol. 1, pp. 347-349, 1984.
- [6] J. Zhou, "Page Segmentation and Classification," *Proc. 9th Int. Conf. on Document Analysis and Recognition*, vol. 1, pp. 619-623, 2007.
- [7] J. Zhou, "Page Segmentation and Classification," *Proc. 9th Int. Conf. on Document Analysis and Recognition*, vol. 1, pp. 619-623, 2007.
- [8] Leptonica image processing and analysis library. <http://www.leptonica.com/>
- [9] R. Smith, "A new method for document layout analysis," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 33, no. 4, pp. 401-408, 2003.
- [10] The Tesseract open source OCR engine. <http://code.google.com/p/tesseract-ocr/>
- [11] A. Antonacopoulos, B. Gatos, D. Bridson, "ICDAR2007 Page Segmentation Competition," *Proc. 9th Int. Conf. on Document Analysis and Recognition*, vol. 1, pp. 619-623, 2007.
- [12] UNLV-ISRI OCR testing toolkit and database. <http://www.isri.unlv.edu/ISRI/OCR/>
- [13] A. Antonacopoulos et al., "ICDAR2009 Page Segmentation Competition," *Proc. 10th Int. Conf. on Document Analysis and Recognition*, vol. 1, pp. 619-623, 2009.