



7. Page Layout Analysis

Finding text regions on pages from books, magazines, and newspapers.

Ray Smith, Google Inc.

Background

- Historically Tesseract had no page layout analysis, but did have text-line finding, assuming a single column of text.
- Cube relies on Tesseract's page-layout/line finding.
- Tesseract's existing text-line finding is also weak wrt diacritics, especially for Arabic and Thai.
- Past methods tend to be:
 - (Bottom-up) Insufficiently aware of page-layout rules or
 - (Top-down) Insufficiently general.

Past Methods: Bottom-Up

- Analyze groups of pixels or connected components to classify into text/image/graphic/blank/line.
- Spread/smear/anneal groups of pixels by some neighborhood voting scheme, morphology or voronoi/graph algorithms.
- Find connected components of labels to group pixels into typed regions.
- Box-up regions into rectangles where possible.
- Morphological approach is very similar.
- Hard to include knowledge like "Columns should usually be the same size."

Past Methods: Top Down

- Often starts with a (possibly pre-trained) model of layout, eg 2-column journal page.
- Attempts to cut the image into the required parts, either with recursive vertical/horizontal cuts, or finding rectangles of whitespace.
- Methods usually fail on non-rectangular regions.
- Methods can often only deal with pages that fit the model.

New Method: Hybrid

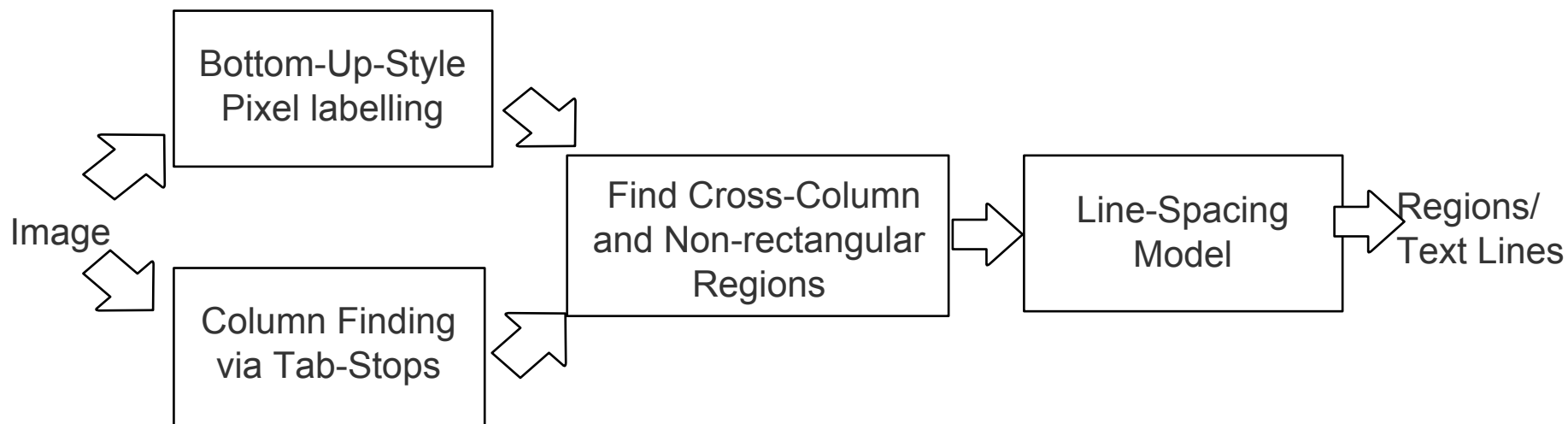
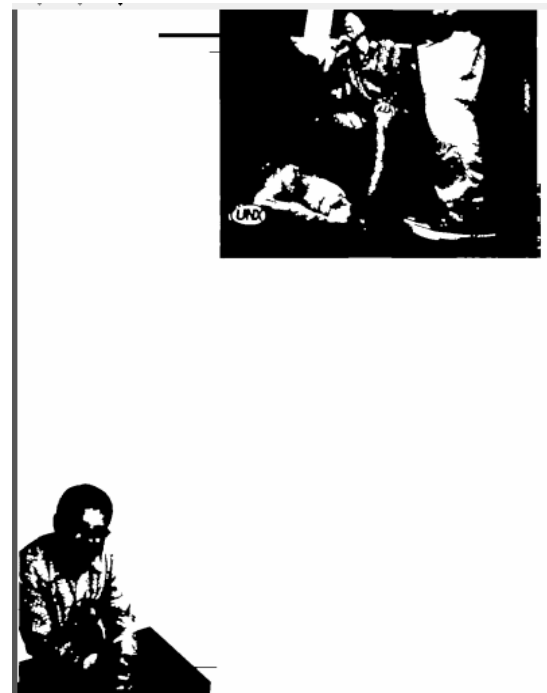
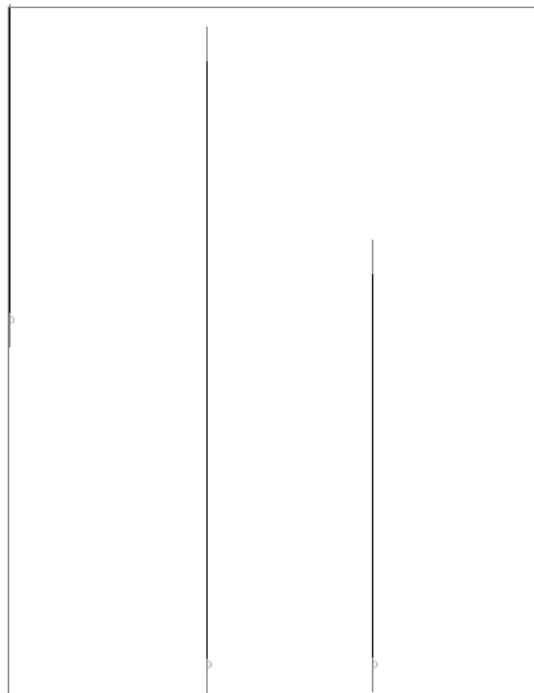


Image-Level Page Layout Analysis

Input Image

Detected Lines

Detected Images

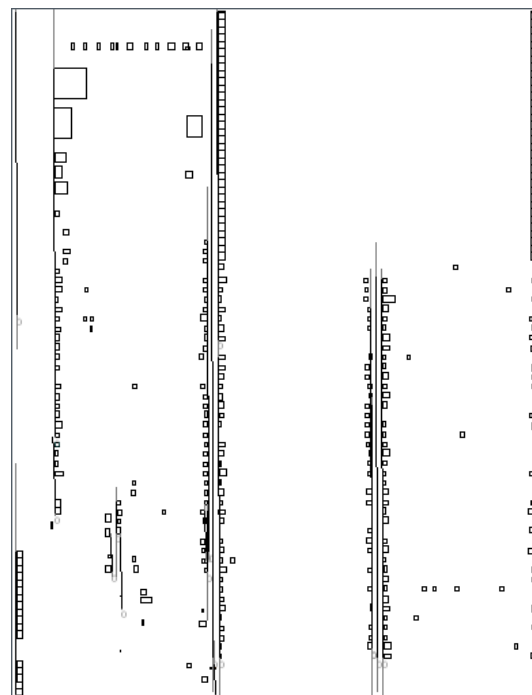
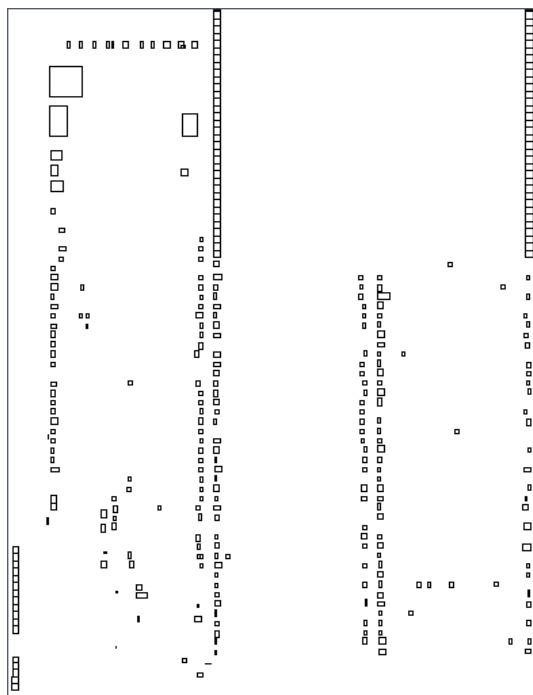
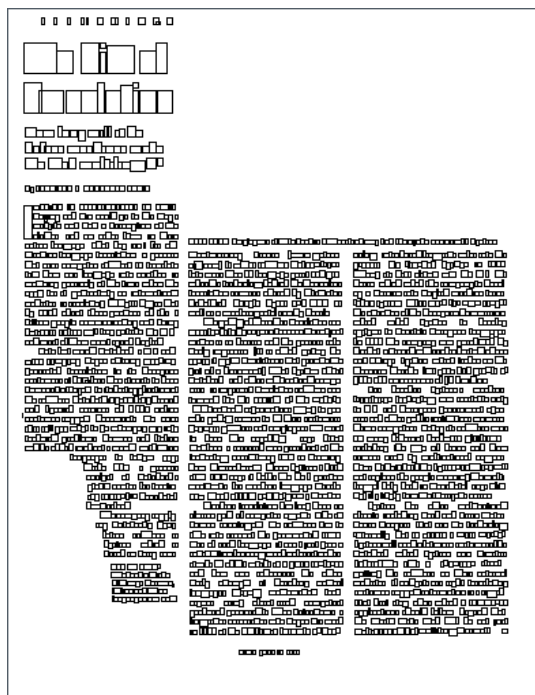


Connected Component Analysis

Text Components

Candidate Tab-Stop
Components

Detected Tab
Segments



Writing direction detection: Got Japanese?

Detect Local writing direction and do tab finding again...

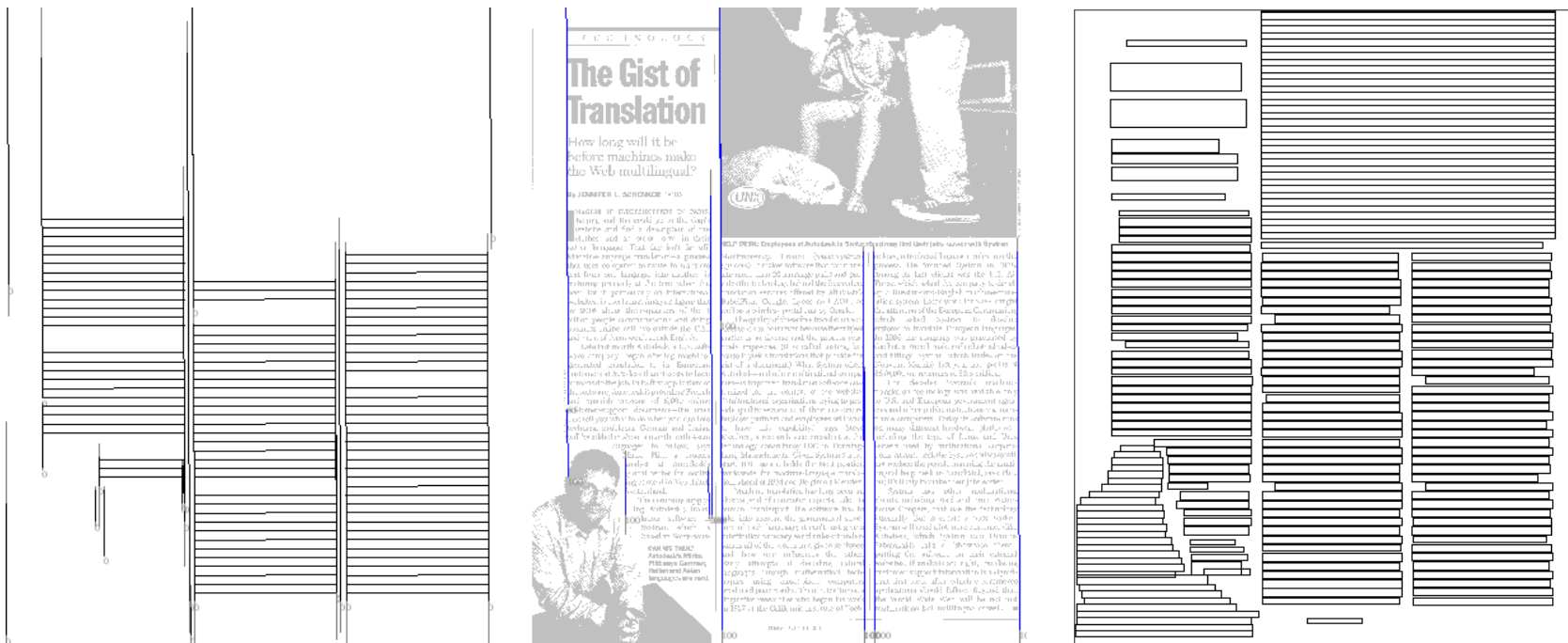
| | |
|---|-----|
| 申込書 | |
| ★季刊考古学 号一 号 20号まで各1345円 21～26号各1854円 27号以降各1860円 | 冊 |
| 注文書は書店さんに書き込んで下さい。お申し込みは郵便振替口座東京3-1685番をご利用下さい。 | |
| ★古墳時代の研究 全13巻 ⑥⑧⑫3800円 組 3500円 | |
| ★末永雅雄著作集 全5巻 各5,800円 組 | |
| ご住所 | |
| Tel. | ご職業 |
| ご氏名 | |

| | | | |
|---|--|---|--|
| 立注文を希望される本に○印をつけて下さい。 第30号 縄文土偶の世界 第31号 環濠集落とクニのおこり 第32号 古代の住居―縄文から古墳へ 第33号 古墳時代の日本と中国・朝鮮 | 古墳時代の研究全10巻 ①総論・研究史 ②集落と豪族居館 ③生活と祭祀 ④生産と流通 ⑤生産と流通Ⅱ ⑥土師器と須恵器 ⑦古墳Ⅰ墳丘と内部構造 ⑧古墳Ⅱ副葬品 ⑨古墳Ⅲ墳輪 ⑩地域の古墳Ⅰ西日本 ⑪地域の古墳Ⅱ東日本 ⑫古墳の造られた時代 ⑬東アジアの中の古墳文化 | 弥生文化の研究全10巻 ①巻三二〇四円 ②巻三二九六円 ③弥生人とその環境 ④生業 ⑤弥生土器Ⅱ ⑥弥生土器Ⅰ ⑦道具と技術Ⅱ ⑧弥生集落 ⑨弥生と瀬と技術Ⅰ ⑩弥生と技術Ⅱ ⑪弥生と技術Ⅲ ⑫弥生と技術Ⅳ ⑬弥生と技術Ⅴ ⑭弥生と技術Ⅵ ⑮弥生と技術Ⅶ ⑯弥生と技術Ⅷ ⑰弥生と技術Ⅷ ⑱弥生と技術Ⅸ ⑲弥生と技術Ⅹ ⑳弥生と技術Ⅺ ㉑弥生と技術Ⅻ ㉒弥生と技術Ⅼ ㉓弥生と技術Ⅽ ㉔弥生と技術Ⅾ ㉕弥生と技術Ⅿ ㉖弥生と技術ⅰ ㉗弥生と技術ⅱ ㉘弥生と技術ⅲ ㉙弥生と技術ⅴ ㉚弥生と技術ⅵ ㉛弥生と技術ⅶ ㉜弥生と技術ⅷ ㉝弥生と技術ⅸ ㉞弥生と技術ⅹ ㉟弥生と技術ⅺ ㊱弥生と技術ⅻ ㊲弥生と技術ⅼ ㊳弥生と技術ⅽ ㊴弥生と技術ⅾ ㊵弥生と技術ⅿ ㊶弥生と技術ⅿ ㊷弥生と技術ⅿ ㊸弥生と技術ⅿ ㊹弥生と技術ⅿ ㊺弥生と技術ⅿ | 論争学説 日本の考古学 別巻6巻 各二〇六〇円 巻一 各二〇六〇円 |
|---|--|---|--|

Column Finding Connected Tabs

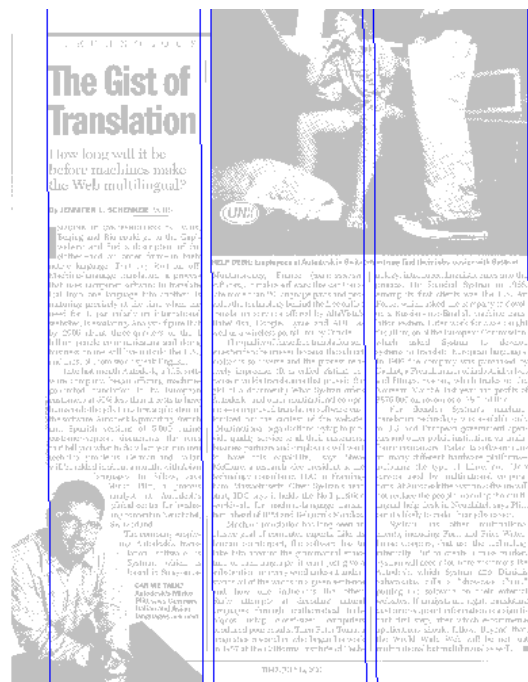
Validated Tab Segments

Candidate Column Partitions

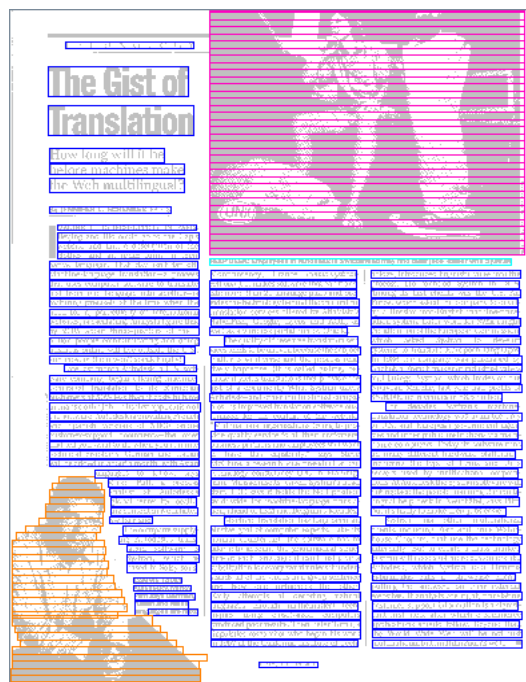


Block Finding

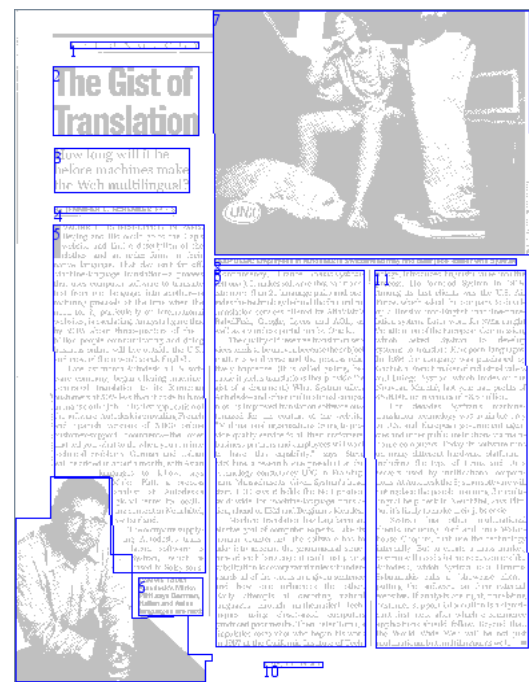
Detected Columns



Typed Column Partitions



Detected Blocks



Live demo of Page Layout Analysis

Simple magazine page with vertical text:

```
api/tesseract unlv/mag.3B/1/8001_044.3B.tif test1 inter segdemo strokewidth
```

Non-rectangular images:

```
api/tesseract unlv/mag.3B/0/8050_078.3B.tif test1 inter strokewidth
```

Multi-orientation Japanese:

```
api/tesseract -l jpn unlv/jpn.tif test1 inter segdemo strokewidth
```

Text on image:

```
api/tesseract unlv/mp00052bw.tif test1 inter strokewidth
```

Thanks for Listening!

Questions?