



# 5. 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.

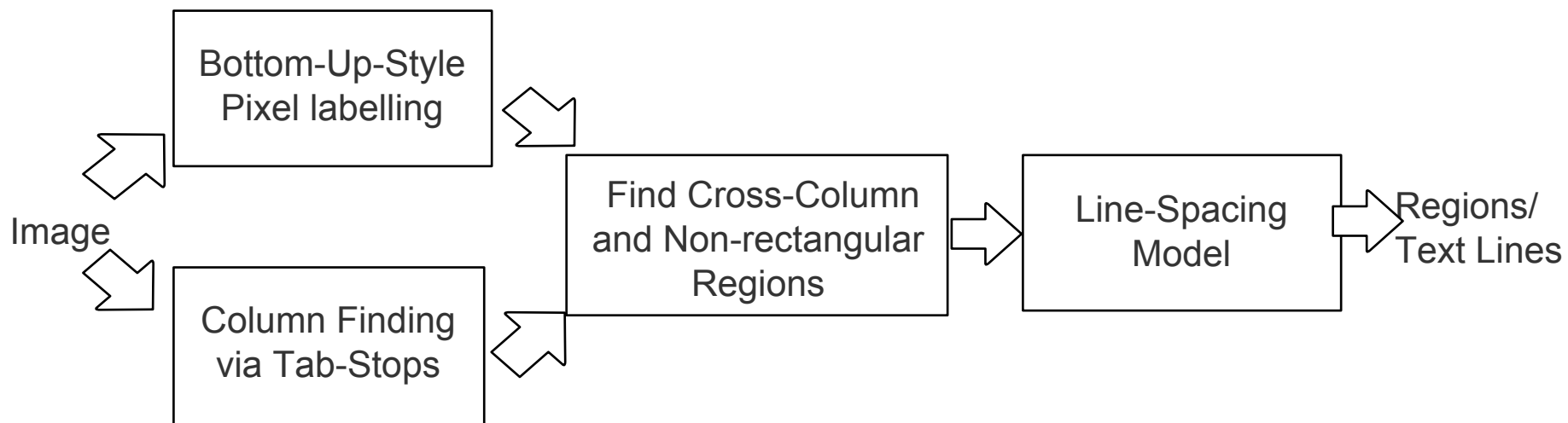
# Past Methods: Bottom-Up

- 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

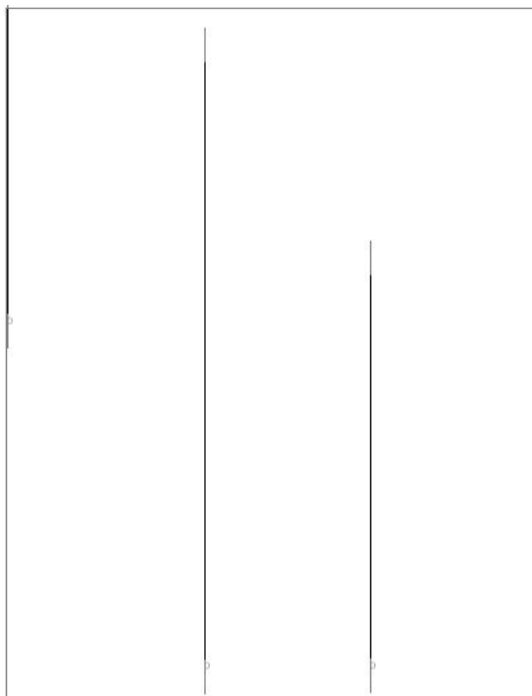


# Image-Level Page Layout Analysis

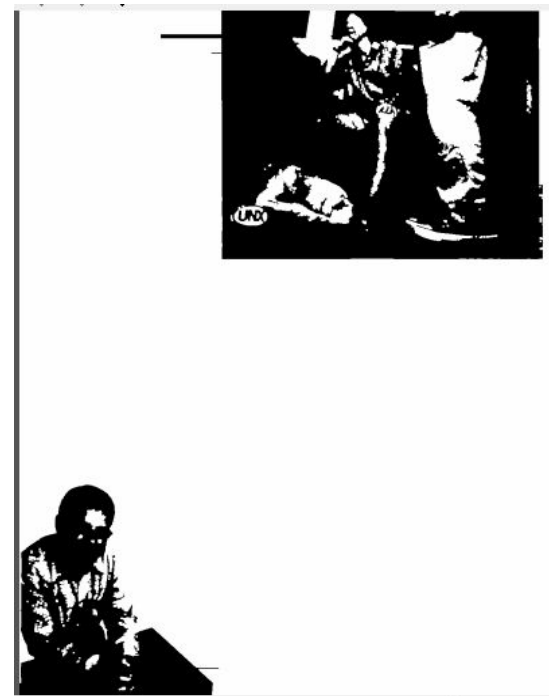
## Input Image



## Detected Lines



## Detected Images

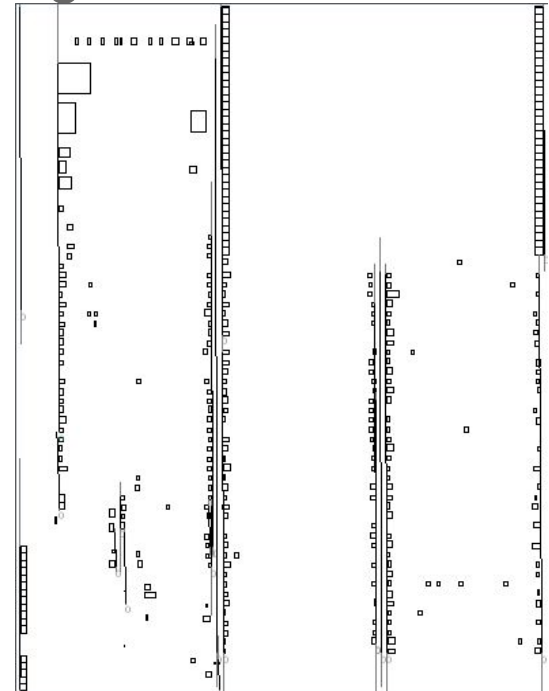
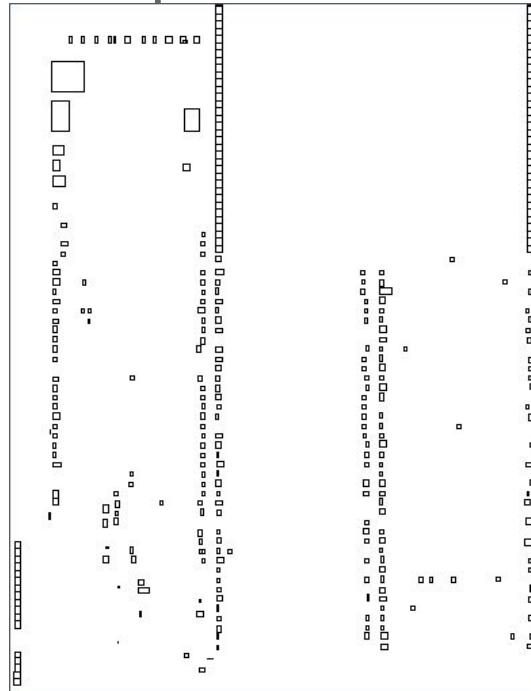
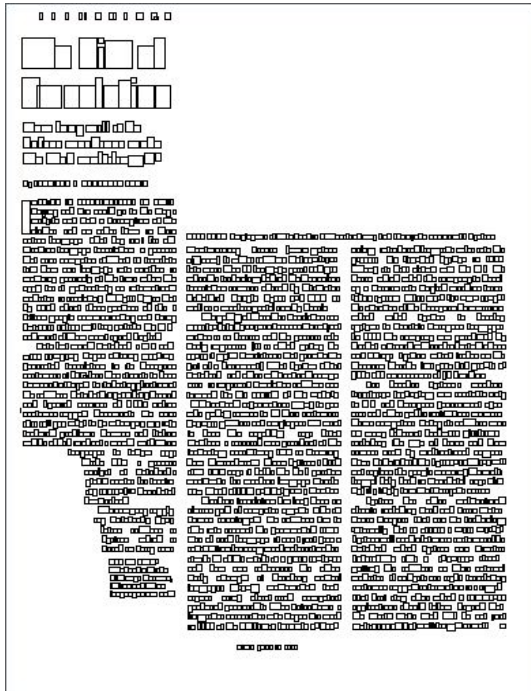


# Connected Component Analysis

Text Components

Candidate Tab-Stop  
Components

Detected Tab  
Segments



# Writing direction detection: Got Japanese?

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

在日朝鮮人の北朝鮮帰還をめぐる  
韓国内閣は「在日朝鮮人の大  
半は戦時中、日本政府が強制労  
働をさせるためにつれてきたもの

## 大半、自由意思で居住

### 外務省、在日朝鮮人で発表

### 戦時徴用は245人

で、いまだに不要になつたため逃  
避するのだ」との趣旨の中傷を行  
なつてゐるのに対し、外務省はこ  
のほど「在日朝鮮人の引揚に關す  
るいきさつ」について発表した。  
これによれば在日朝鮮人の総数は  
約六十一万人だが、このうち戦時  
中に徴用労働者として日本に來た  
者は二百四十五人にすぎないとさ  
れてゐる。主な内容は次の通り。

一、戦前（昭和十四年）に日本内  
地に住んでゐた朝鮮人は約百万  
人で、終戦直前（昭和二十年）  
には約百万人となつた。増加  
した百万人のうち、七十万人は  
自分から進んで内地に職を求め  
てきた個別徴用者で、その間の  
出生によるものである。残りの  
三十万人は大部分、工礦業、土  
木事業の募集に応じてきた者  
で、戦時中の國民徴用令による  
徴用労働者はごく少数である。  
また、國民留用令は日本内地で  
は昭和十四年七月に實施された  
が、朝鮮への適用はさしつかえ  
昭和十九年九月に実施されてお

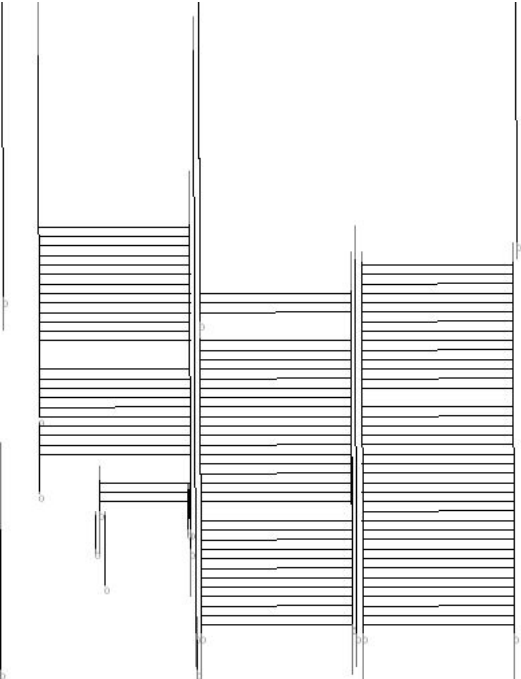
り、朝鮮人徴用労働者が導入  
されたのは、翌年三月の下関一  
釜山間の運輸が止るまでのわず  
か七日間であった。

一、終戦後、昭和二十年八月から  
翌年三月まで、捕虜者が政府の  
船舶、個別引揚げで合計百四十  
万人が帰還したほか、北朝鮮へ  
は昭和二十一年三月、赤台団の  
指令に基づく北朝鮮引揚計画で三  
百五十人が帰還するなど、終戦  
時まで在日してゐた者のうち  
七五％が帰還してゐる。戦時中  
に來日した労働者、復讐軍人、  
軍捕などは日本内地になじみが  
薄いため終戦後、残留した者は  
ごく少数である。現在、登録さ  
れてゐる在日朝鮮人は総計六十  
一万人で、関係各省で來日の事  
情を調査した結果、戦時中に徴  
用労働者としてきた者は二百四  
十五人にすぎず、現在、日本に  
居住してゐる者は犯罪者を除  
き、自由意思によつて在留した  
者である。

By English: Asahi Shimbun 日本語: 朝日新聞 [Public domain], via Wikimedia Commons  
[https://commons.wikimedia.org/wiki/File%3AAsahi\\_Shimbun\\_newspaper\\_clipping\\_\(13\\_July\\_1959\\_issue\).jpg](https://commons.wikimedia.org/wiki/File%3AAsahi_Shimbun_newspaper_clipping_(13_July_1959_issue).jpg)



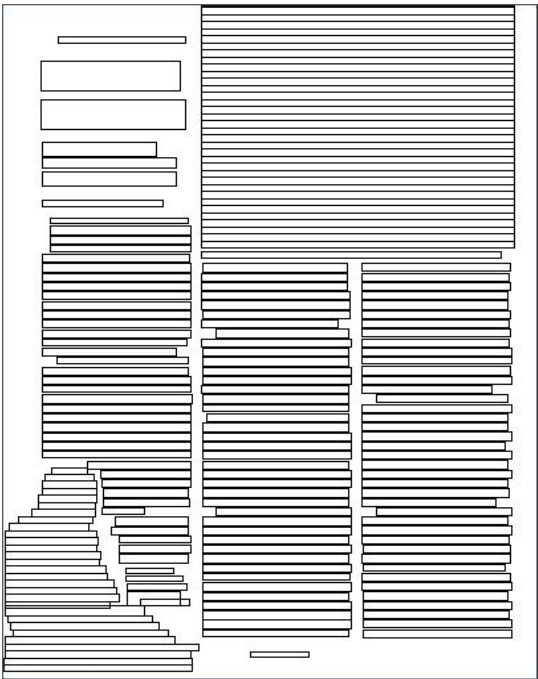
# Column Finding Connected Tabs



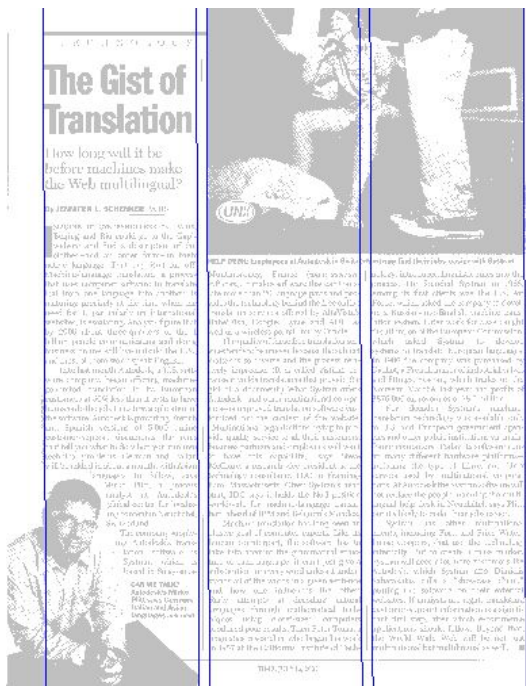
# Validated Tab Segments

This image shows a scanned newspaper page with a grid overlay. The grid is used to identify and validate tab segments. The page contains text, a photograph of a person, and a small logo. The grid lines are drawn over the page, highlighting specific segments and their connections. The text on the page includes the title "The Gist of Translation" and a sub-headline "How long will it be before machines make the Web multilingual?". The grid helps in identifying the structure and relationships between different parts of the page.

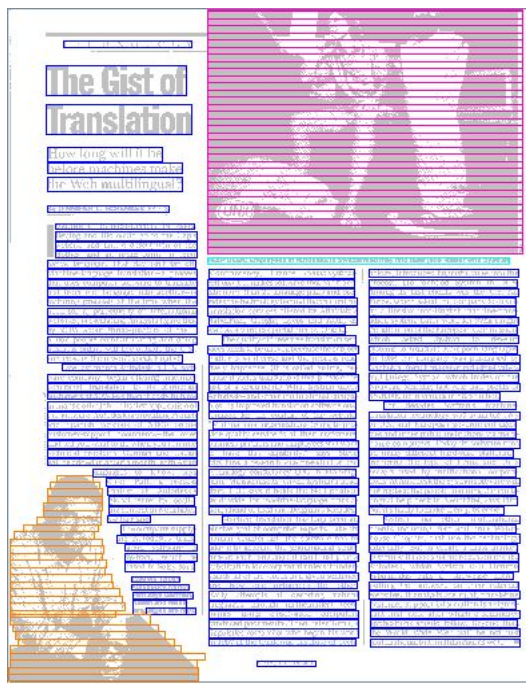
# Candidate Column Partitions



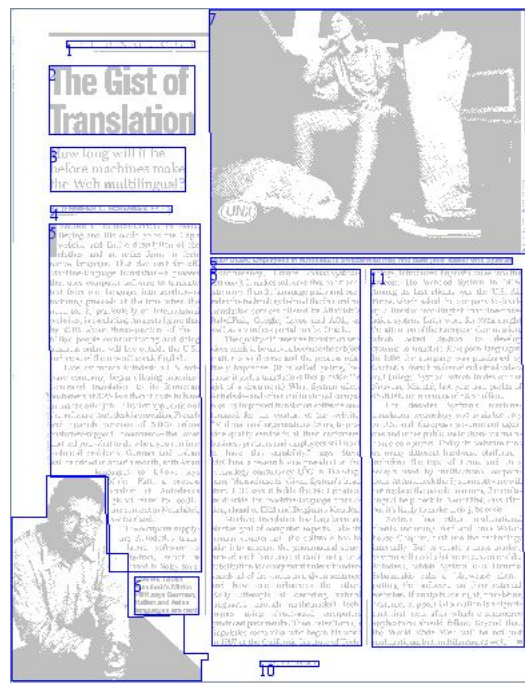
# Block Finding Detected Columns



# Typed Column Partitions



# Detected Blocks



# Thanks for Listening!

## Questions?