# Math Spotting in Technical Documents Using Handwritten Queries
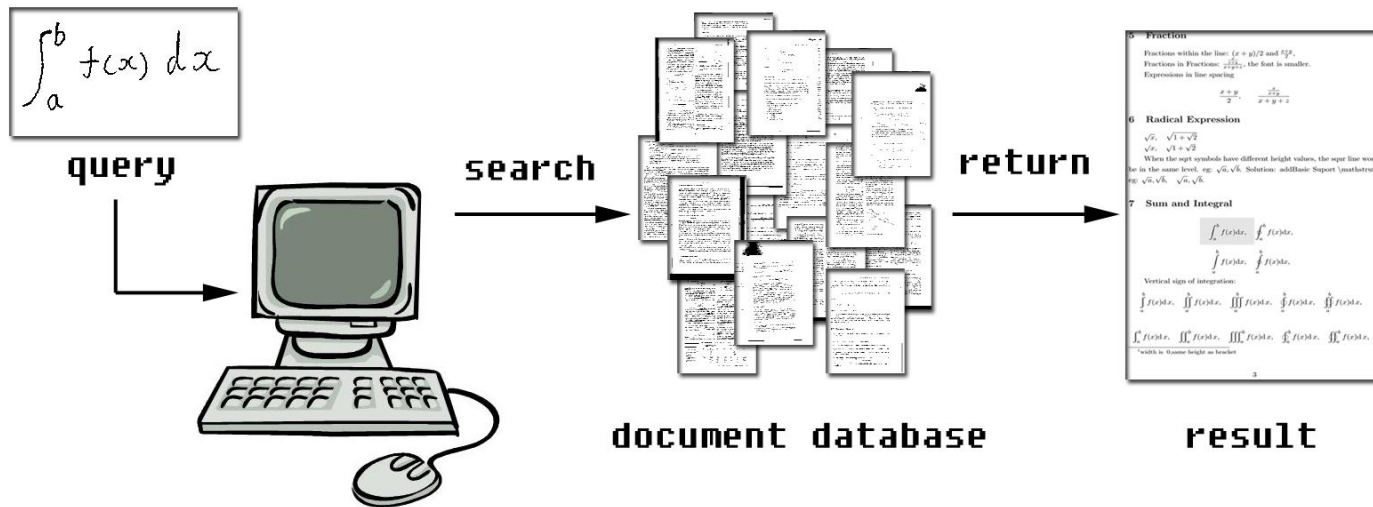
Li Yu and Richard Zanibbi

Document and Pattern Recognition Lab

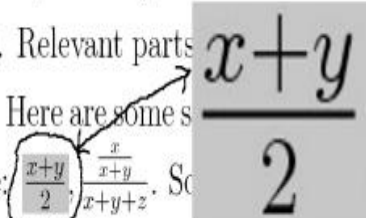Rochester Institute of Technology, NY, USA

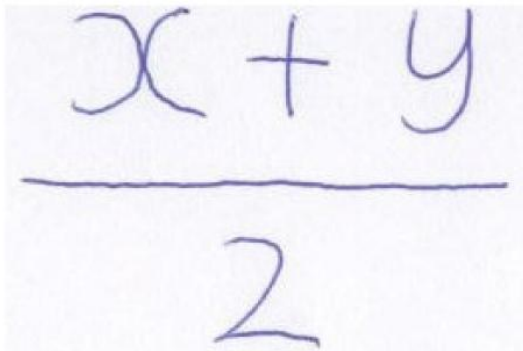lxy1783@rit.edu, rlaz@cs.rit.edu

# Math spotting



query     search     return

document database     result

- **OCR (optical character recognition) avoided**
- **Structure feature & Visual feature**

# Document image and query image

This file is prepared by running latex A.tex and cutting the pictures out of the resulting preview. Relevant parts $x+y$ ode are reproduced under each of the pictures. Here are some s $\dfrac{x+y}{2}$ $, x_{11}^{22}, x_m^{(k)}, {}^*x^*, x^{m^n}, x^{x^{x^x}}$ Other examples include: $\boxed{\frac{x+y}{2}}$ $\frac{\frac{x}{x+y}}{x+y+z}$. So mbols have an explanatory text. $\int_a^b f(x)\,dx$, $\oint_a^b f(x)\,dx$, This text is found in the latex code, mostly stating that they are parts of some spacial setup and cannot be used in standard LaTeX.
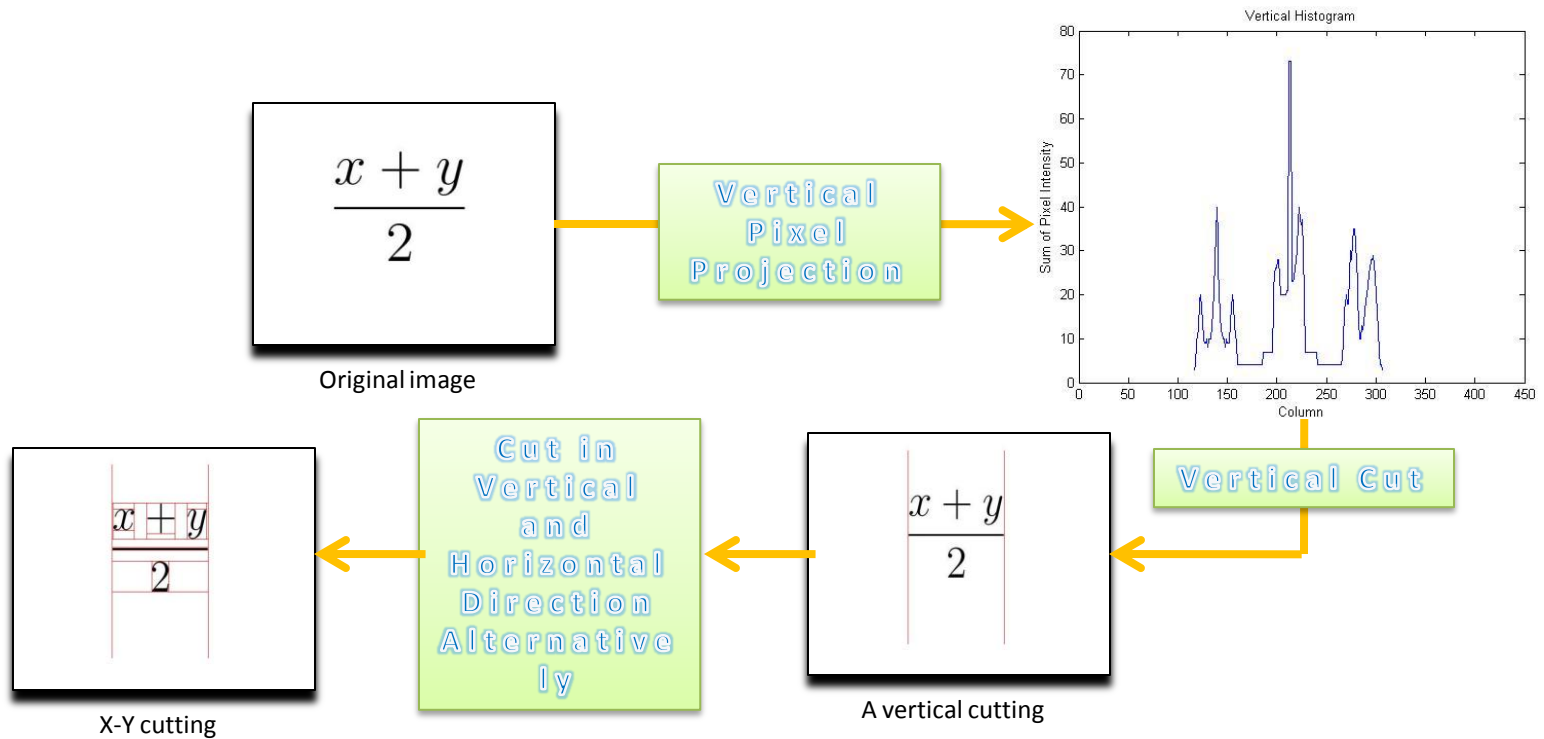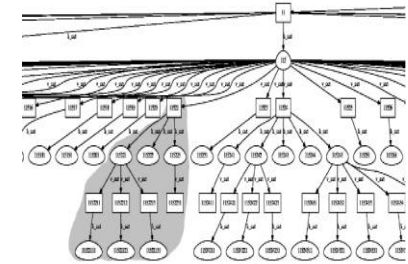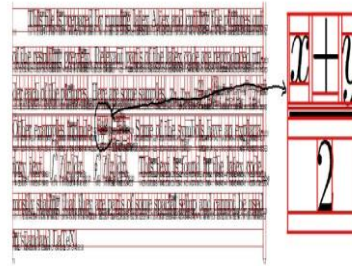
Document image

Query image

# X-Y cutting



Original image

Vertical Pixel Projection

Vertical Histogram

Cut in Vertical and Horizontal Direction Alternatively

Vertical Cut

X-Y cutting

A vertical cutting

$$\frac{x+y}{2}$$

G. Nagy and S. Seth, "Hierarchical representation of optically scanned documents," *Proc. of ICPR, (1984) 347-349.*

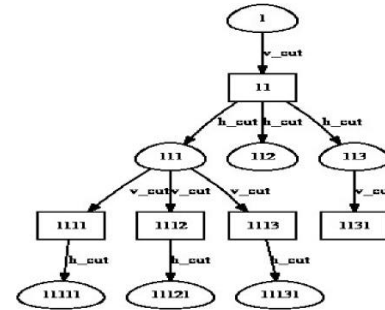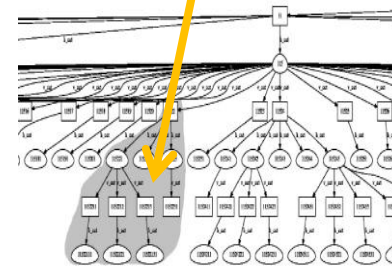# X-Y cut and X-Y tree



Page

Query

Images

X-Y cutting

X-Y trees

# Sub-tree matching

- **What if we can find a matched sub-tree in the page tree?**

- **What we want?**
  **Speed & Accuracy**

- **Problems?**
  **Inexact matching**



X-Y tree for query
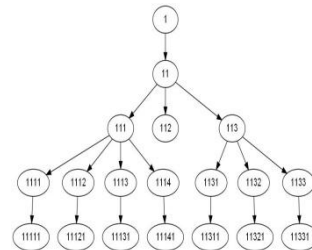
X-Y tree for page

# Noise and "Bad Division"



Cutting in Query

Cutting in Page

Noise

Bad Division

- **Avoid noise**
- **Control the way in which regions are cut**
- **Rectangles whose size smaller than thresholds will be ignored**

# Thresholds

One Line In Document

Width of Peaks

Horizontal Projection

- **Dominant height/width of characters**
- $C_h = Mode(h_1,h_2,...h_n)$, **where $h_n$ represents the heights of lines in one page**
- $W_h = Mode(W_1,W_2,...W_n)$, **where $W_n$ represents the widths of blank spaces in one line**
- **Scaled linearly based on the current region's height and width**

# Equivalency Class

"((()())())"  (3)

"(()())"  (2)    (1)  "()"

(1)    (1)

"()"    "()"

(4)  "(()(()))"

"()"  (1)    (2)  "(()())"

(1)    (1)  — Equivalency Class Num

"()"    "()"  — Code

- **Two trees have same code (equivalence class number) if and only if they are isomorphic**
- **Bottom-up algorithm with linear time in the size of the trees**

A.V. Aho, J.E. Hopcroft, and J.D. Ullman. The design and analysis of computer algorithms. Addison Wesley, Reading, Mass, 1974.

# Ranking by Equivalency Class



Equivalency Class Num

Equivalency Class Num for Query Root (Ni)

Decreasing Similarity Order

E_NUM:N $_i$

E_NUM:N $_{i-1}$

E_NUM:N $_{i-2}$

E_NUM:N $_{i-3}$

E_NUM:N $_{i-4}$

1

Order Nodes by Equivalency Class Num

X-Y tree for query image

X-Y tree for page image

# Ranking by Equivalency Class

$$p = \frac{P}{(R_0\, m^0 + R_1\, m^1)\, T}$$

11

Query

physical equations were solved together with the dynamical equations. Transformation rates are given as functions of partial densities $\rho^k$, cloud droplet concentration at cloud base $N_{cl,0}$, and spectral width $C_f$. If mass fractions $m^k$ are required, $\rho^k$ may simply be replaced by $\rho^k = \rho m^k$ with $\rho$ determined by

$$p = \frac{\mu}{(R_0\, m^0 + R_1\, m^1)\, T} \qquad (23)$$

Rank2

## Spectral distributions

Cloud drops are assumed to follow a non-normalized log-normal density distribution:

$$f_{cl}(\ln m) = \frac{N_{cl}}{\sigma_{cl}\sqrt{2\pi}}\exp\left[-\frac{(\ln m - \mu_{cl})^2}{2\sigma_{cl}^2}\right] \qquad (24)$$

Rank3     Rank4

with

$m$ – drop mass

$N_{cl}$ – total number concentration of cloud drops in $m^{-3}$

$\sigma_{cl}^2$ – variance of $f_{cl}(\ln m)$

- **The query are included in the page**

Page

# Ranking by Equivalency Class



Query
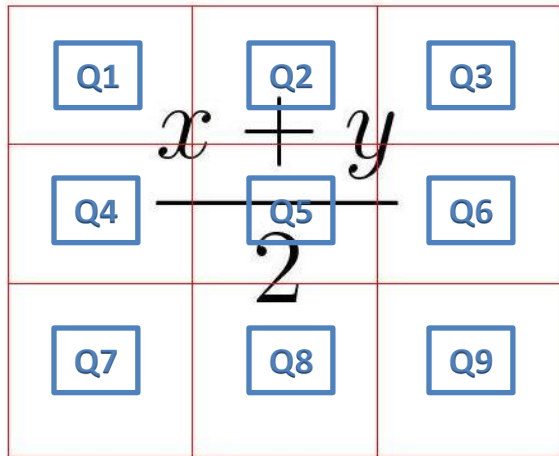
Page

Rank5

Rank3

Rank2

Rank5

Rank5

Rank3

- **The query are not included in the page**

# Other Rankings

- Ranking by Number of Nodes:

  - Divide the page nodes into bins based on their size.
  - Start with the size of the query root.
  - Search for the page nodes in decreasing size order.

- Ranking by both equivalence class number and number of nodes:

  - Generate the equivalence class number for both query and page.
  - Start with the query root and by decreasing order.
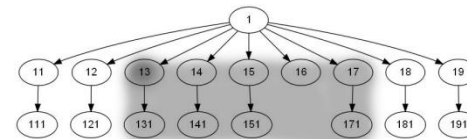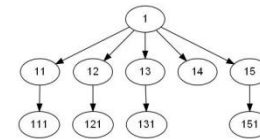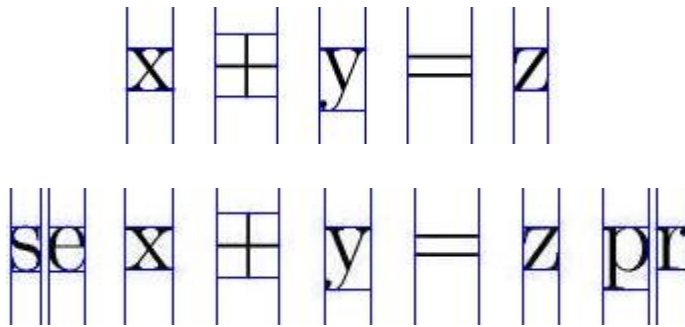  - Find all the exact sub-matches in the page tree.

# Visual Feature

| Q1 | Q2 | Q3 |
|----|----|----|
| Q4 | Q5 | Q6 |
| Q7 | Q8 | Q9 |

$$\frac{x+y}{2}$$

$$Dis\tan ce = \sum_{i=1}^{9} \frac{(Q_i - P_i)^2}{P_i}$$

**Where Qi and Pi represents the sum of pixel intensity in the sub-region in query and candidate respectively**

- **Dividing the region into nine sub-regions and computing sum of pixel intensity respectively**
- **Ranking the candidates by decreasing visual similarity**

# Problems and Future work

- **The situation where the target is "scattered" in the page.**



- **q03vp03.htm**

# Problems and Future work

- **Different Rankings**

- **More visual features && comparison**

- **Document image indexing**

# Thanks

**Question?**