# A Math Formula Extraction and Evaluation Framework for PDF Documents

**Ayush Kumar Shah,** Rochester Institute of Technology, USA

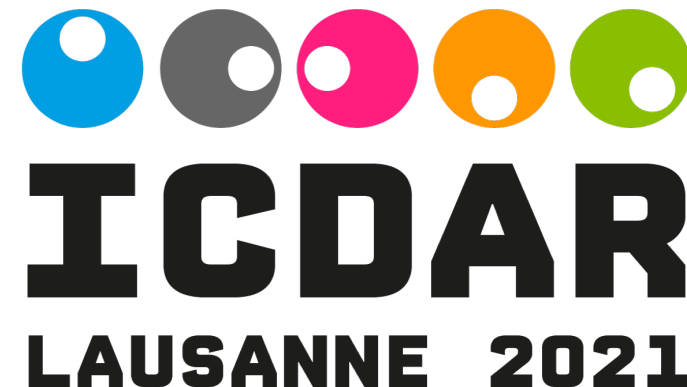**Abhisek Dey,** Rochester Institute of Technology, USA

**Richard Zanibbi,** Rochester Institute of Technology, USA
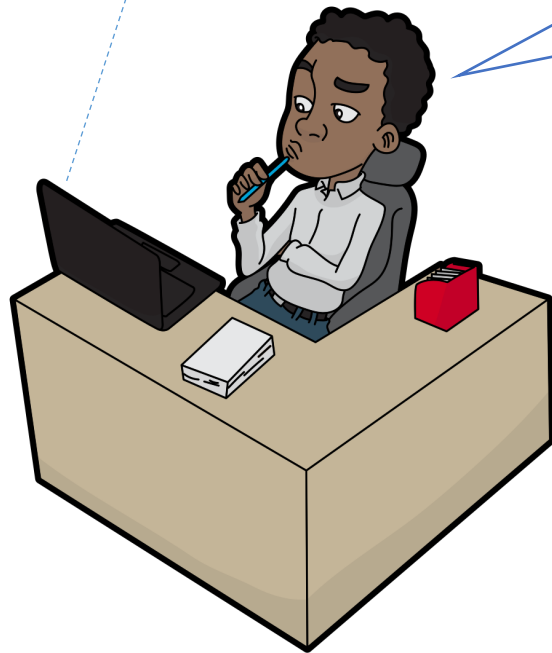
{as1211, ad4529, rxzvcs}@rit.edu

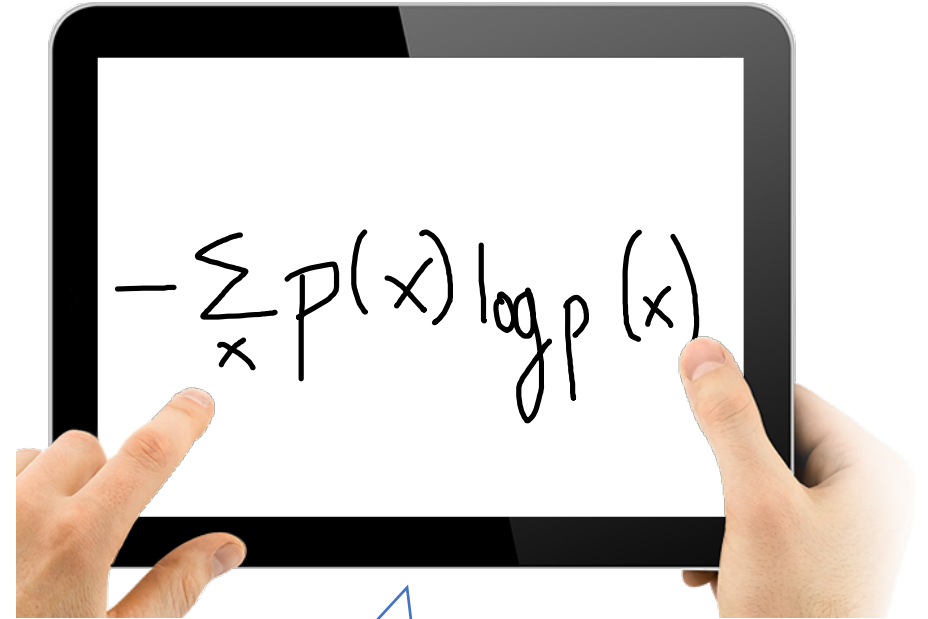Document and Pattern Recognition Lab

RIT | Rochester Institute of Technology

ICDAR LAUSANNE 2021

# Motivation: Information Retrieval

$$-\sum_x p(x) \log p(x)$$

What is the name of this formula?

$$-\sum_x p(x) \log p(x)$$

Let's find documents describing this formula.

# Mathematical Formula Extraction: Overview

**SymbolScraper**

Extract character BBs and labels using pdf info (no OCR)

**ScanSSD**

Locate formula regions in document images

Example Document Page (Input)

ScanSSD + SymbolScraper Output

# Mathematical Formula Extraction: Overview

ScanSSD Output

ScanSSD + SymbolScraper Output

**CC Boxes**

**SS Boxes + labels**

CC Extraction

$(\delta(u,m) = 1$

**Formula Box**

Combine

$( \quad \delta \quad ( \quad u \quad , \quad m \quad ) \quad =$

$( \quad \delta \quad ( \quad u \quad , \quad m \quad ) \quad = \quad \_$

**QD-GGA**

Parse formula structure

**Symbol Layout Tree (Output)**

# SymbolScraper: Extracting Symbols in PDF

- Based on Apache PDFBox
- Avoids OCR in **born-digital** PDF documents and instead uses vector drawing commands in PDF
- Unicode, writing line position and attributes derived from PDF encoding
- *'em box'* or underlying character outlines (glyphs) represent symbol outlines in a font as boxes

| PDF Miner | PyMuPDF | PDFBox | SymbolScraper |

Unlike other methods, SymbolScraper uses glyphs to fine-tune bounding box locations

# SymbolScraper: Extracting Symbols in PDF

- Glyphs and font scaling information used to obtain precise bounding box locations
- Compound characters (large braces, square roots, etc.) are formed of 2 or more characters



Compound Symbol

Correcting glyph origins

# ScanSSD: Locating Formula Regions

- **Scanning Single-Shot Detector,** CNN which locates formula bounding boxes using a sliding window

- 600 dpi images broken into windows of 1200 x 1200 pixels, SSD applied in each window at 10% stride

- Non-Maximal Suppression selects the highest confidence regions from overlapping detections

- Wider default boxes sizes used with aspect ratios of 5, 7, and 10 -> increased recall

A window (in Blue) slides across the grid (in Red)
(50% stride)

Default boxes around a grid point

- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: **SSD: Single shot multibox detector**. In: European conference on computer vision. pp. 21–37. Springer (2016)
- Mali, P., Kukkadapu, P., Mahdavi, M., Zanibbi, R.: **ScanSSD: Scanning Single Shot Detector for Mathematical Formulas in PDF Document Images**. arXiv:2003.08005 [cs] (2020)

# ScanSSD: Locating Formula Regions

- A sliding window divides the page into windows which are processed by ScanSSD
- The partial predictions at the window-level are pooled together and the final regions are identified using pixel-wise voting (stitching)



Input Page

Sliding Window

SSD

Stitch Patches

- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: **SSD: Single shot multibox detector**. In: European conference on computer vision. pp. 21–37. Springer (2016)
- Mali, P., Kukkadapu, P., Mahdavi, M., Zanibbi, R.: **ScanSSD: Scanning Single Shot Detector for Mathematical Formulas in PDF Document Images**. arXiv:2003.08005 [cs] (2020)

# ScanSSD: Locating Formula Regions



**Bottom:** Window-level Predictions
**Top:** Confidence masks

Pixel-wise voting

Likely Regions

Thresholding

Binary masks

CC Extraction

Final Predictions

• Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: **SSD: Single shot multibox detector**. In: European conference on computer vision. pp. 21–37. Springer (2016)
• Mali, P., Kukkadapu, P., Mahdavi, M., Zanibbi, R.: **ScanSSD: Scanning Single Shot Detector for Mathematical Formulas in PDF Document Images**. arXiv:2003.08005 [cs] (2020)

# QD-GGA: Recognizing Formula Structure (Parsing)

1. **Construct** graph over CCs

2. **Prune:** Convert to LOS graph

3. **Classify** edges as merge/split and relationships, nodes as symbols

4. **New LOS graph:** detected symbols

5. **Extract MST** using Edmond's arborescence algorithm



1. Complete graph

2. LOS graph

3. CNN outputs — Symbol classification, Relation classification, Symbol detection

4. Symbol-level LOS graph

5. Symbol Layout Tree

Mahdavi, M.; Sun, L.; Zanibbi, R.: **Visual Parsing with Query-Driven Global Graph Attention (QD-GGA)**. In Conference on Computer Vision and Pattern Recognition Workshops (2020)

# Inputs



Input Formula with
CCs or Symbols

# Outputs



Symbol Layout Tree (SLT)

```
<math xmlns="http://www.w3.org/1998/Math/MathML">
<mrow>
    <mi xml:id="0:">ζ</mi>
    <mrow>
        <mo xml:id="1:">(</mo>
        <mrow>
            <mrow>
                <mi xml:id="2:">T</mi>
                <mo xml:id="3:">,</mo>
            </mrow>
            <mrow>
                <mi xml:id="4:">k</mi>
                <mo xml:id="5:">)</mo>
            </mrow>
        </mrow>
    </mrow>
</mrow>
</math>
```
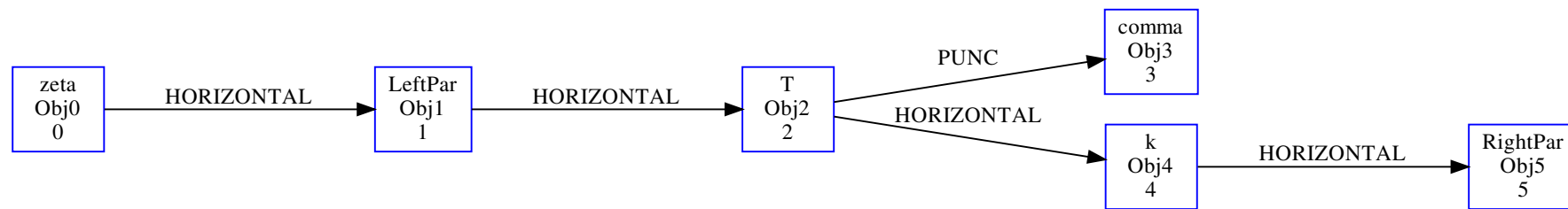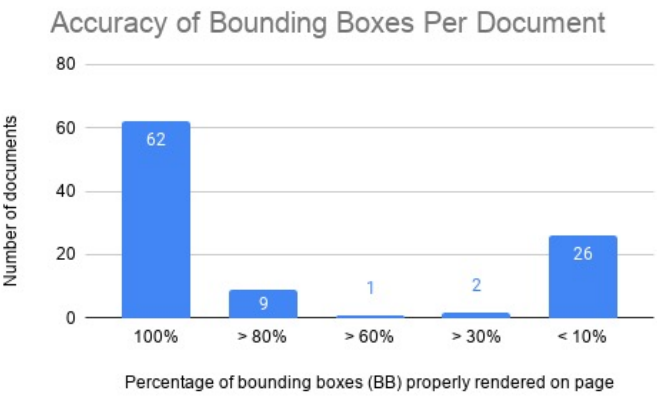
SLT in MathML

```
\(\zeta\left( {{T,}\left. k \right)} \right.\)
```

SLT in LATEX

Mahdavi, M.; Sun, L.; Zanibbi, R.: **Visual Parsing with Query-Driven Global Graph Attention (QD-GGA)**. In Conference on Computer Vision and Pattern Recognition Workshops (2020)

# SymbolScraper Results

## ScanSSD Results

## QD-GGA Results

Summary of SymbolScraper Accuracy

Formula Detection Results for TFD-ICDAR2019

Formula Recognition Results for InftyMCCDB-2[1] Test set



Accuracy of Bounding Boxes Per Document

| | IOU $\geq$ 0.75 | | | IOU $\geq$ 0.5 | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F-score | Precision | Recall | F-score |
| ScanSSD | **0.774** | **0.690** | **0.730** | **0.851** | **0.759** | **0.802** |
| RIT 2 | 0.753 | 0.625 | 0.683 | 0.831 | 0.670 | 0.754 |
| RIT 1 | 0.632 | 0.582 | 0.606 | 0.744 | 0.685 | 0.713 |
| Mitchiking | 0.191 | 0.139 | 0.161 | 0.369 | 0.270 | 0.312 |
| Samsung* | 0.941 | 0.927 | 0.934 | 0.944 | 0.929 | 0.936 |

*Used character information

| Metrics | Value |
|---|---|
| Structure rate | 92.56 |
| Structure + Classification rate | 85.94 |

| Hardware Specifications and Speed | |
|---|---|
| Storage | HDD |
| Dataset Size | 100 pages |
| Total time | 28 mins, 19 secs |
| Average time | 1.7 secs/page |

| Hardware Specifications and Speed | |
|---|---|
| Storage | HDD |
| RAM | 32 GB |
| Graphics | Nvidia RTX 2080 Ti |
| Processor | AMD Ryzen 7 2700 |
| Dataset Size | 233 pages |
| Total time | 4 hrs, 33 mins, 31 secs |
| Average time | 70.4 secs/page |

| Hardware Specifications and Speed | |
|---|---|
| Storage | HDD |
| RAM | 32 GB |
| Graphics | Nvidia GTX 1080 |
| Processor | Intel(R) Core(TM) i7-9700KF |
| Dataset Size | 6830 images |
| Total time | 26 mins, 25 secs |
| Average time | 232 ms/formula |

- [1] https://zenodo.org/record/3483048#.XaCwmOdKjVo
- Mali,P.,Kukkadapu,P.,Mahdavi,M.,Zanibbi,R.:**ScanSSD:ScanningSingleShot Detector for Mathematical Formulas in PDF Document Images**. arXiv:2003.08005 [cs] (2020)
- Mahdavi, M.; Sun, L.; Zanibbi, R.: **Visual Parsing with Query-Driven Global Graph Attention (QD-GGA)**. In Conference on Computer Vision and Pattern Recognition Workshops (2020)

# Recognition Results Visualization (HTML)



**MathSeer Pipeline Results Visualization**

Pdf name: K15-1002

Page: 4

[Previous page] [Home] [Next page]

Page image

# LgEval Extension: Error Visualization



Errors organized by decreasing frequency

Specific instances where 'z' is misclassified as '2,' seen after clicking on the '22 errors' link

# LgEval Extension: Error Visualization



**Zoomed in:** Specific instances where 'z' is misclassified as '2,' seen after clicking on the '22 errors' link

# Conclusion and Future Work

- Open-source formula extraction pipeline for PDF documents
  - **https://www.cs.rit.edu/~dprl/software.html**
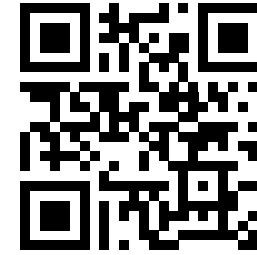- PDF symbol extractor that identifies precise bounding box locations in born-digital PDFs
- A simple and effective algorithm for detection of math expressions using visual features alone
- Extended tools for visualizing recognition results and formula parsing errors
- **ScanSSD-XYc:** Unified page and window level merging using recursive XY Cuts avoiding NMS speeding up detection by 300 times approximately (included in the repository)

**Future work**

- **SymbolScraper:** Handle Type 3 Fonts and faster system for symbol extraction, better handling of compound characters
- **Pipeline:** End-to-End trainable system for detection and parsing

- Dey, A.; Zanibbi, R.: ScanSSD-XYc: Faster Detection of Math Formulas. In The 14th IAPR International Workshop on Graphics Recognition (GREC 2021), to appear; 2021.

# Thank You

**Alfred P. Sloan FOUNDATION**

**dprl**

**Document and Pattern Recognition Lab**