# Structure Recognition

**Joris Schellekens**
**Software Engineer - iText**

Joris Schellekens
Software Engineer
iText

# The iText R&D team

Structure Recognition

- Various **standards**: PDF/A-1, PDF/A-2, PDF/A-3, PDF/UA, PDF/E, PAdES, etc.

- Level of **descriptiveness** (e.g. metadata) varies.

- **Basic** level:
  instructions for how a viewer (e.g. Adobe Reader) should render a document.

```
[a, -28.7356, p, 27.2652, p, 27.2652, e, -27.2652, a, -28.7356, r, 64.6889, a, -28.7356, n, 27.2652, c, -38.7594,
e, 444] TJ
/R10 10.44 Tf
68.16 0.24 Td
[", 17.1965, P, -18.7118, i, -9.35592, l, -9.35592, o, -17.2414, t, -9.35636, ", 17.1965,  , 250] TJ
```

Joris Schellekens
Software Engineer
iText

# Expectations of end-users

- People are used to programs like Microsoft® Word™.

- Reflow (e.g. flowing content around an image).

- Structure is part of the document.

- Export to various formats.

- Extracting data.

- Change appearance of a logical unit (e.g. word, line, or paragraph).

Joris Schellekens
Software Engineer
iText

- **iText Group nv is an open source company:**
  - **Close to the community.**
  - **Code should be something you want other people to see - accessible (rather than a dirty hack) .**
  - **Code should be something the user can change - contribution.**
  - Submitting pull requests.
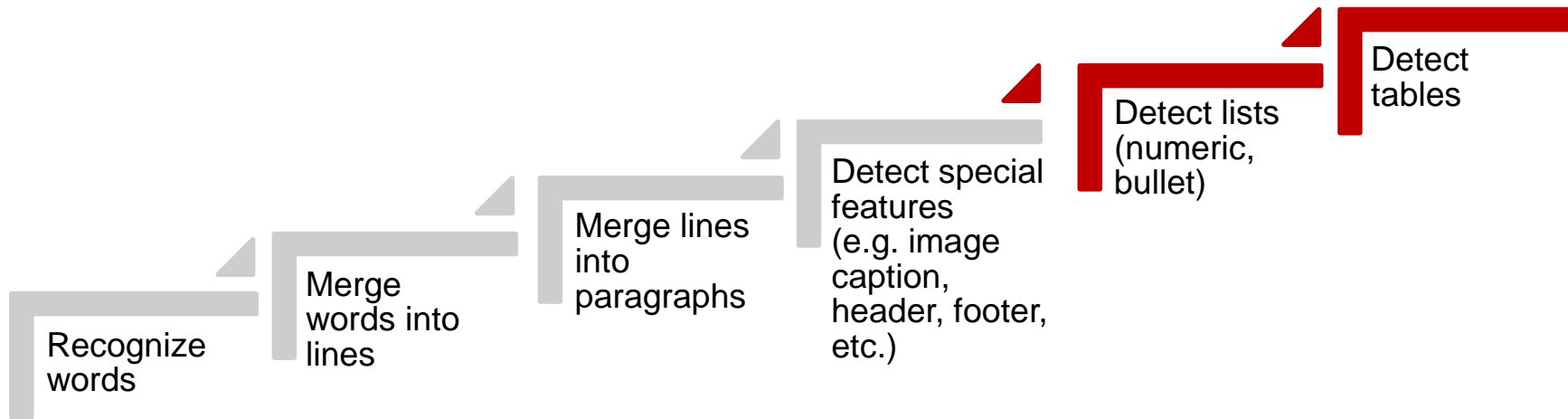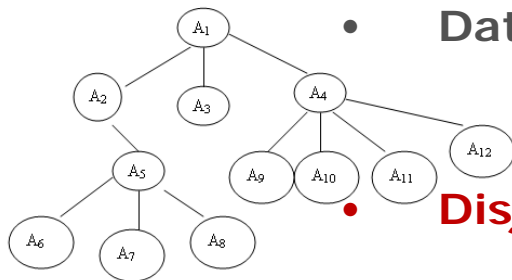  - Tweak and fine-tune to match experience and expertise.

**AGPL**

Joris Schellekens
Software Engineer
iText

Technical

Recognize words

Merge words into lines

Merge lines into paragraphs

Detect special features (e.g. image caption, header, footer, etc.)

Detect lists (numeric, bullet)

Detect tables

Joris Schellekens
Software Engineer
iText

- Data structure – evaluating equivalence.

- **Disjoint set** algorithm:

  - Space, find, merge all in *O(a(n))*.

  - Maps nicely to PDF ideas.

Joris Schellekens
Software Engineer
iText

- **Disjoint set tells us how the merging happens, not when.**

- **When depends on what is being merged - 2 approaches:**

1. **(Human) logic**
   - Chunks are merged into words based on distance.
   - Words are merged into lines based on distance and global page layout.
   - Lines into paragraphs based on distance and visual cues.

Joris Schellekens
Software Engineer
iText

# When to merge (alt)

- **Alternative (bends the constraints) – use Artificial Intelligence.**

- **Tackle the problem like an image recognition problem**

**PDF**
- PDF
- Preferably tagged (to enable supervised learning)

**Convert to graphics**

**Deep neural network**
- Mark certain sections
- Table, list, paragraph, header, footer, etc.

- Deep neural network coefficients
- Error rate
- Not retractable

**Deliverable**

Joris Schellekens
Software Engineer
iText

# In Depth

# Building a training set

- **Use of NN requires training data**
  - **Large volume of (perfectly) tagged PDF documents**

- **Not readily available**

- **Build our own?**
  - **pdfHTML (convert random crawled HTML to PDF)**
  - **Gutenberg**
  - **Crawl the internet for PDFs**
  - **Industry contacts**

Joris Schellekens
Software Engineer
iText

# Convert training set

# Convert training set (2)

- **Feature selection**
  - **X**
  - **Y**
  - **Width**
  - **Height**
  - **Fontsize**
  - **Bold?**
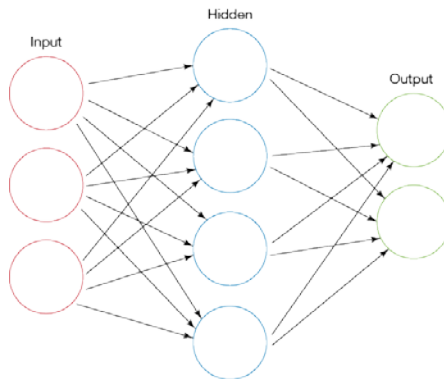  - **Italic?**
  - **Underline?**
  - **ΔX**
  - **ΔY**

Joris Schellekens
Software Engineer
iText

# Feed the beast

| 41 | 56 | 10 | 26 | 16 | 52 | 56 | 8 | 26 | 16 | 1 | 0 | 1 |
| 52 | 56 | 8 | 26 | 16 | 41 | 56 | 10 | 26 | 16 | 1 | 0 | 1 |
| 52 | 56 | 8 | 26 | 16 | 60 | 56 | 9 | 26 | 16 | 0 | 0 | 1 |
| 60 | 56 | 9 | 26 | 16 | 52 | 56 | 8 | 26 | 16 | 0 | 0 | 1 |
| 60 | 56 | 9 | 26 | 16 | 70 | 56 | 5 | 26 | 16 | 1 | 0 | 1 |
| 70 | 56 | 5 | 26 | 16 | 60 | 56 | 9 | 26 | 16 | 1 | 0 | 1 |
| 70 | 56 | 5 | 26 | 16 | 76 | 56 | 8 | 26 | 16 | 1 | 0 | 1 |
| 76 | 56 | 8 | 26 | 16 | 70 | 56 | 5 | 26 | 16 | 1 | 0 | 1 |
| 76 | 56 | 8 | 26 | 16 | 84 | 56 | 14 | 26 | 16 | 0 | 0 | 1 |
| 84 | 56 | 14 | 26 | 16 | 76 | 56 | 8 | 26 | 16 | 0 | 0 | 1 |
| 84 | 56 | 14 | 26 | 16 | 99 | 56 | 9 | 26 | 16 | 1 | 0 | 1 |
| 84 | 56 | 14 | 26 | 16 | 276 | 190 | 17 | 865 | 16 | 178 | 134 | 0 |



Joris Schellekens
Software Engineer
iText

```java
private void tagUsingAI(InputStream wekaModel, int confidence, List<StructureNodeImpl> nodes, boolean converge){
    IMerger nn = new WekaModelMerger()
            .load(wekaModel)
            .setThreshold(confidence / 100.0);

    // converge
    if(converge) {
        IMerger cv = new Convergence(nn);
        cv.apply(nodes);
    }
    else{
        nn.apply(nodes);
    }
}


public interface IMerger {

    void apply(List<StructureNodeImpl> l);

}
```

Joris Schellekens
Software Engineer
iText

# Benefits

# Tagging an untagged document

```
[...]
        <MCID text="e" x="409.0" y="686.0" width="11.0" height="25.0" />
        <MCID text="r" x="421.0" y="686.0" width="8.0" height="25.0" />
        <MCID text="l" x="429.0" y="686.0" width="5.0" height="25.0" />
        <MCID text="a" x="435.0" y="686.0" width="11.0" height="25.0" />
        <MCID text="n" x="446.0" y="686.0" width="12.0" height="25.0" />
        <MCID text="d" x="459.0" y="686.0" width="12.0" height="25.0" />
        <MCID text=" " x="472.0" y="686.0" width="5.0" height="25.0" />
    </Span>
</P>
<P lang="nl" x="72.0" y="295.0" width="468.0" height="109.0">
    <Span x="72.0" y="295.0" width="182.0" height="14.0">
        <MCID text="m" x="72.0" y="295.0" width="9.0" height="14.0" />
        <MCID text="e" x="81.0" y="295.0" width="5.0" height="14.0" />
        <MCID text="t" x="87.0" y="295.0" width="3.0" height="14.0" />
        <MCID text=" " x="91.0" y="295.0" width="2.0" height="14.0" />
    [...]
```
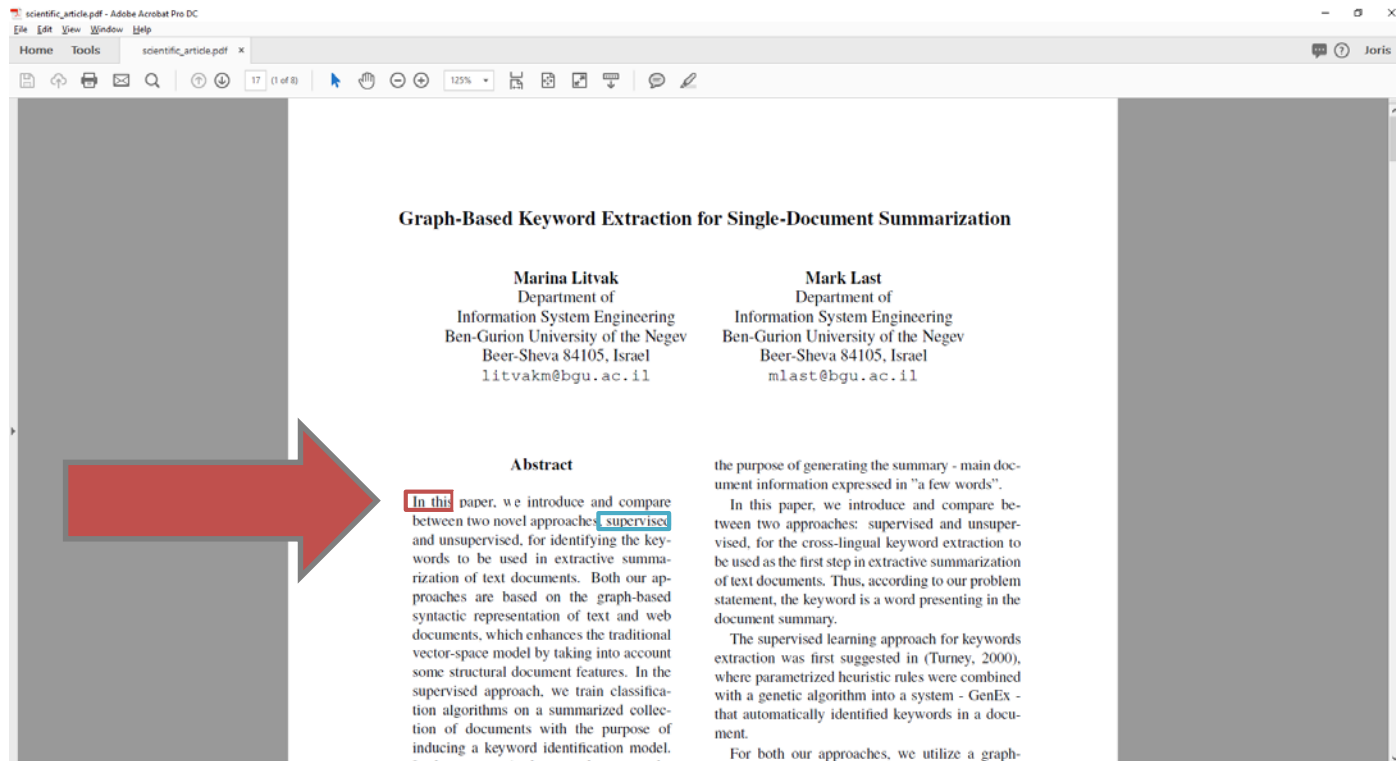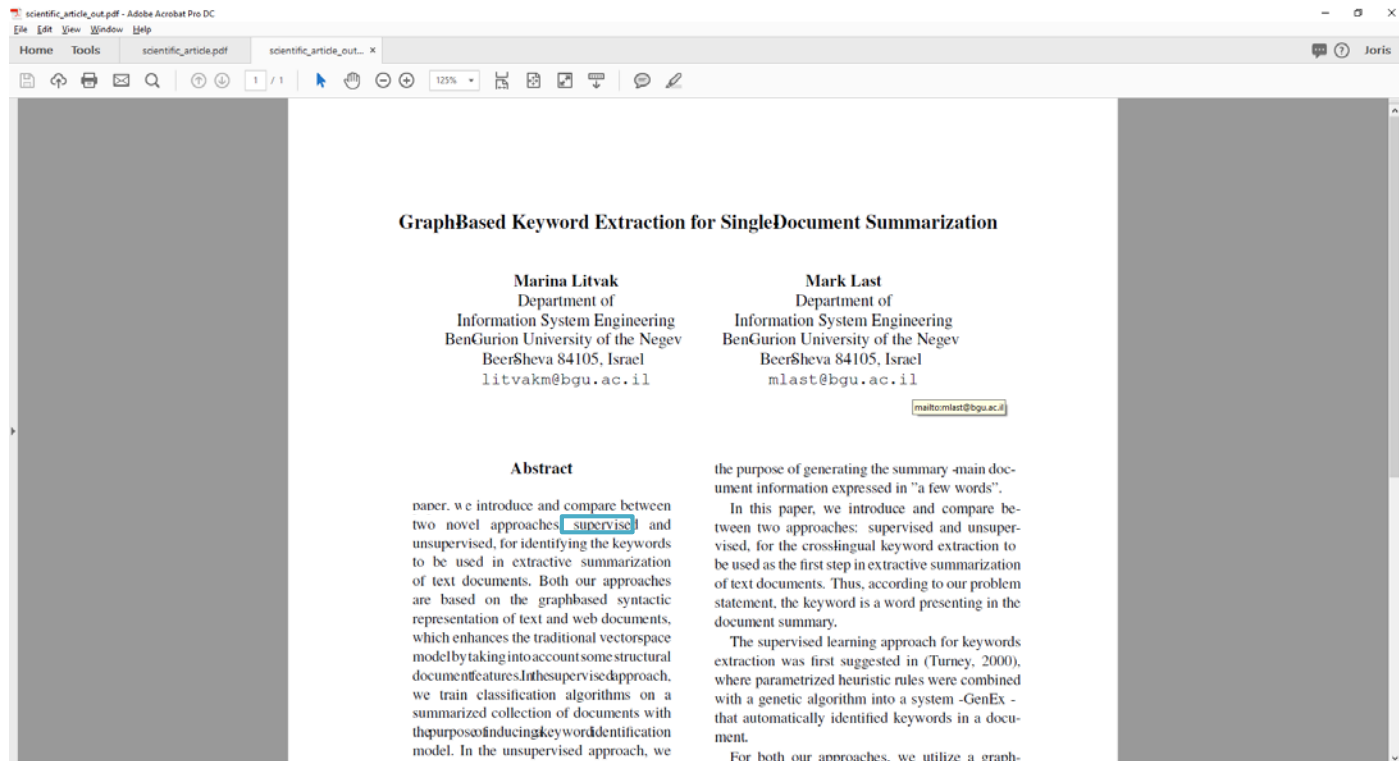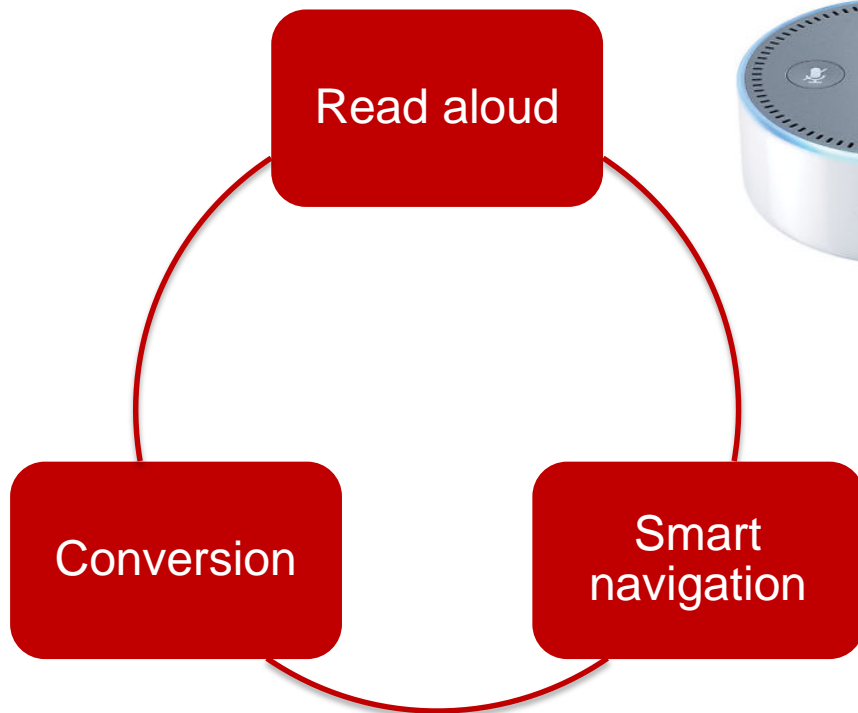
Joris Schellekens
Software Engineer
iText

Structure Recognition

19

# Reflow

# Reflow

Joris Schellekens
Software Engineer
iText

2018-05-14

- **Extreme case:**

  - **Layout content of A4 PDF on A3 canvas.**

Read aloud

Conversion

Smart navigation

Joris Schellekens
Software Engineer
iText

# Data extraction

Joris Schellekens
Software Engineer
iText

```
┌─────────────┐  ┌─────────────┐  ┌─────────────┐  ┌─────────────┐  ┌─────────────┐
│  Recognize  │  │    Merge    │  │ Merge lines │  │   Detect    │  │   Detect    │
│    words    │  │ words into  │  │    into     │  │   spatial   │  │ tables and  │
│             │  │    lines    │  │ paragraphs  │  │  features   │  │    lists    │
└─────────────┘  └─────────────┘  └─────────────┘  └─────────────┘  └─────────────┘
```

- Data structures
  - Disjoint set: (Human) logic | Global level | Artificial intelligence
- Benefits:
  - still configurable.
  - able to train network to work well for documents you typically process.
  - 'soft fail' mode.

ITEXT

Joris Schellekens
Software Engineer
iText

www.pdfa.org

# Thank you! Any questions?

Joris Schellekens
Software Engineer
iText

Get in touch:     joris.schellekens@itextpdf.com
Web site:          www.itextpdf.com
Twitter:           @itext