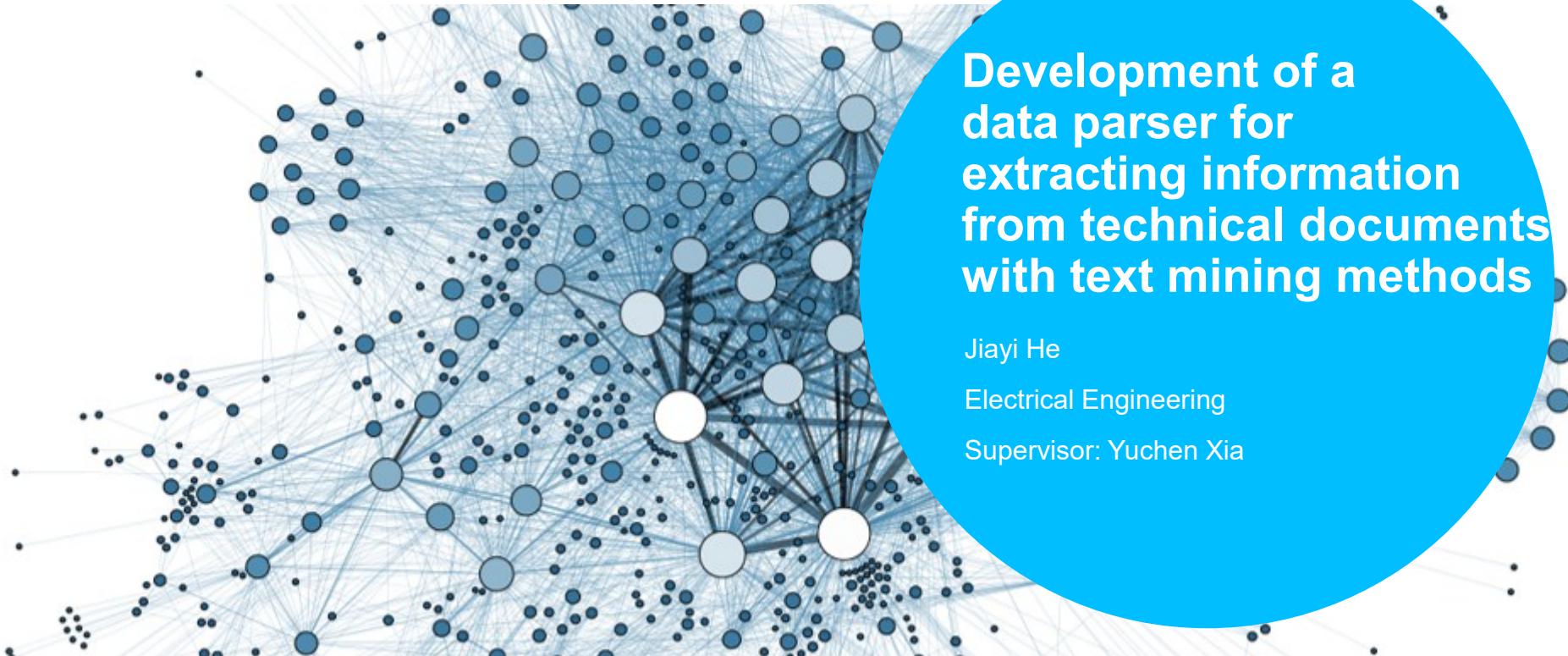




**University of Stuttgart**  
Institute of Industrial Automation  
and Software Engineering



## **Development of a data parser for extracting information from technical documents with text mining methods**

Jiayi He

Electrical Engineering

Supervisor: Yuchen Xia



# Agenda

- Motivation and Difficulties
- Basis
- Conception
- Implementation
- Evaluation
- Summary and Outlook



# Motivation

## Technical Documents in Modern Automation Industry



[2]



[3]

- Equipment maintenance
- Quality control
- Supply chain management
- Regulatory compliance

### Defect of Manual Data Extraction & Entry:

- Time-Consuming: large number of documents
- Error-Prone: manual copying and typing
- Inconsistent: inconsistent formats from different staffs in one project

# Difficulties

- PDF file format variations: text-based, scanned documents.
- Layout complexity: multi-column, complex tables, embedded images.
- Large size: multiple tables in pages, time-consuming.

## Accuracy

Measurement  
Linearizer  
Display range  
CJC accuracy

## Electrical input

| Input type | M |
|------------|---|
| mV         |   |
| V          |   |
| mA         |   |

**TSSP530..**  
Vishay Semiconductors

**IR Sensor Module for Reflective Sensor, Light Barrier, and Fast Proximity Applications**



**FEATURES**

- Up to 2 m for presence and proximity sensing
- Uses modulated bursts of infrared light
- PIN diode and sensor IC in one package
- Low supply current
- Shielding against EMI
- Visible light is suppressed by IR filter
- Insensitive to supply voltage ripple and noise
- Supply voltage: 2.5 V to 5.5 V
- Material categorization: for definitions of compliance please see [www.vishay.com/doc/99912](http://www.vishay.com/doc/99912)

**MECHANICAL DATA**

Pinning:  
1 = OUT, 2 = GND, 3 = V<sub>G</sub>

**DESCRIPTION**

The TSSP530.. series are compact infrared detector modules for presence and fast proximity sensing applications. They provide an active low output in response to infrared bursts at 940 nm. The frequency of the burst should correspond to the carrier frequency shown in the parts table.

This component has not been qualified according to automotive specifications.

**LINKS TO ADDITIONAL RESOURCES**

[STB](#) [ICMark](#)

**APPLICATIONS**

- Reflective sensors for hand dryers, towel or soap dispensers, water faucets, toilet flush
- Vending machine fall detection
- Security and pet gates
- Person or object vicinity activation
- Fast proximity sensors for toys, robotics, drones, and other consumer and industrial uses

**PARTS TABLE**

| Carrier frequency | 38 kHz | TSSP53038                                |
|-------------------|--------|--|
|                   | 56 kHz | TSSP53056                                |
| Package           |        | Minimold                                 |
| Pinning           |        | 1 = OUT, 2 = GND, 3 = V <sub>G</sub>     |
| Dimensions (mm)   |        | 6.0 W x 6.95 H x 5.6 D                   |
| Mounting          |        | Leaded                                   |
| Application       |        | Presence sensors, fast proximity sensors |

| °C    | °F     |
|-------|--------|
| Max.  | Min.   |
| 1650  | 32     |
| 1649  | 32     |
| 205.4 | 32     |
| 450   | 32     |
| 761   | 32     |
| 262   | -328   |
| 260.6 | 32     |
| 760   | -328   |
| 1373  | -328   |
| 205.7 | 32     |
| 450   | 32     |
| 762   | 32     |
| 1842  | 211    |
| 1399  | 32     |
| 800   | 32.0   |
| 100   | -149.7 |
| 206   | -328   |
| 537.3 | -149.7 |
| 100.9 | 32     |
| 300   | 32     |
| 800   | 32.0   |

**TSSP530..**  
Vishay Semiconductors

**IR Sensor Module for Reflective Sensor, Light Barrier, and Fast Proximity Applications**



**FEATURES**

- Up to 2 m for presence and proximity sensing
- Uses modulated bursts of infrared light
- PIN diode and sensor IC in one package
- Low supply current
- Shielding against EMI
- Visible light is suppressed by IR filter
- Insensitive to supply voltage ripple and noise
- Supply voltage: 2.5 V to 5.5 V
- Material categorization: for definitions of compliance please see [www.vishay.com/doc/99912](http://www.vishay.com/doc/99912)

**MECHANICAL DATA**

Pinning:  
1 = OUT, 2 = GND, 3 = V<sub>G</sub>

**DESCRIPTION**

The TSSP530.. series are compact infrared detector modules for presence and fast proximity sensing applications. They provide an active low output in response to infrared bursts at 940 nm. The frequency of the burst should correspond to the carrier frequency shown in the parts table.

This component has not been qualified according to automotive specifications.

**LINKS TO ADDITIONAL RESOURCES**

[STB](#) [ICMark](#)

**APPLICATIONS**

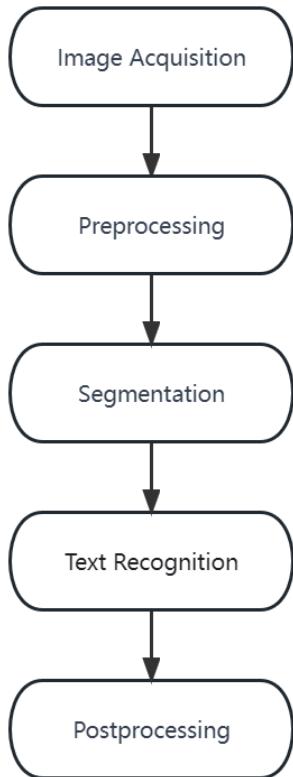
- Reflective sensors for hand dryers, towel or soap dispensers, water faucets, toilet flush
- Vending machine fall detection
- Security and pet gates
- Person or object vicinity activation
- Fast proximity sensors for toys, robotics, drones, and other consumer and industrial uses

**PARTS TABLE**

| Carrier frequency | 38 kHz | TSSP53038                                |
|-------------------|--------|--|
|                   | 56 kHz | TSSP53056                                |
| Package           |        | Minimold                                 |
| Pinning           |        | 1 = OUT, 2 = GND, 3 = V <sub>G</sub>     |
| Dimensions (mm)   |        | 6.0 W x 6.95 H x 5.6 D                   |
| Mounting          |        | Leaded                                   |
| Application       |        | Presence sensors, fast proximity sensors |

# Basis

## Optical Character Recognition (OCR) Method



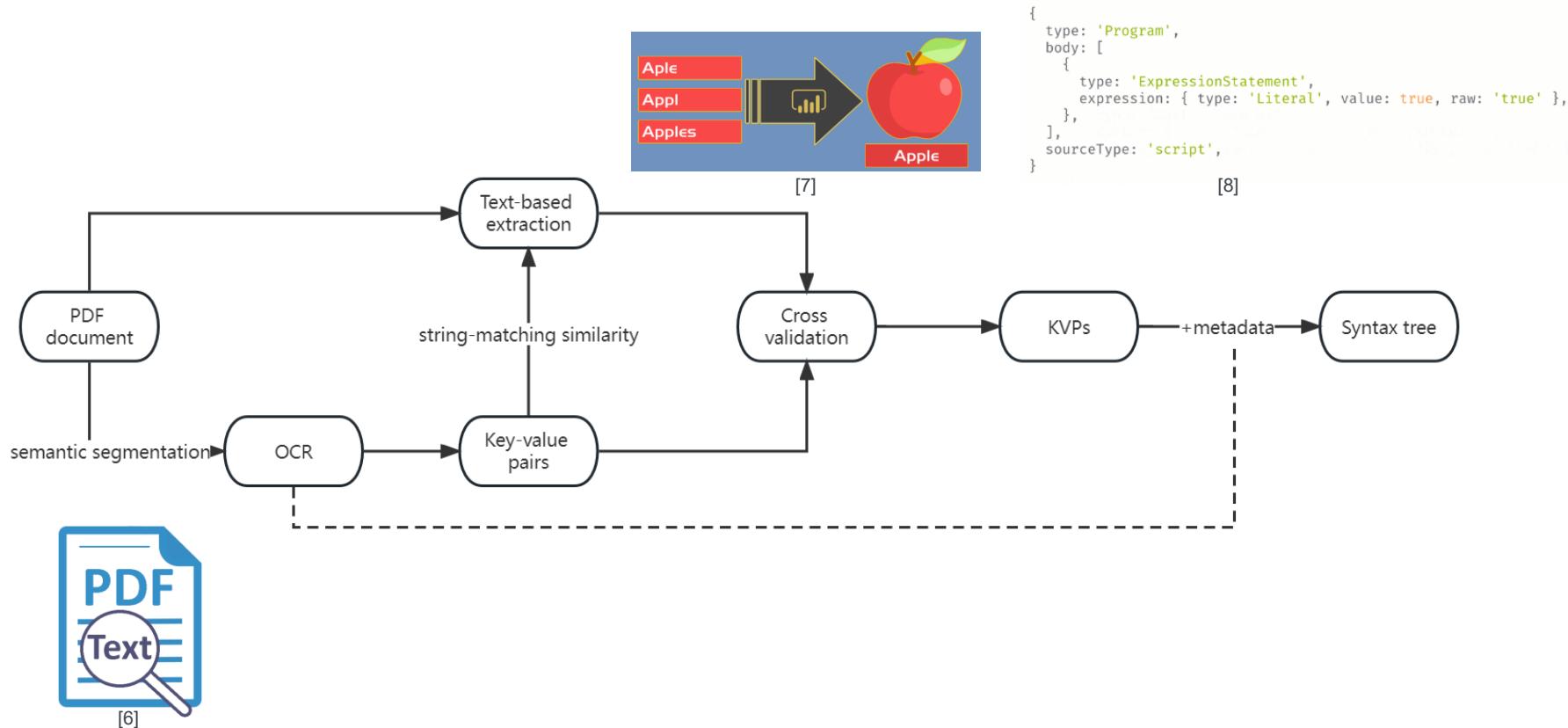
### Pros:

- All types of PDF (text-based and image-based).
- High accuracy with pre-trained ML models.
- Time- and effort-saving.

### Cons:

- Errors in special symbols.
- Limited supporting languages.
- Costly.

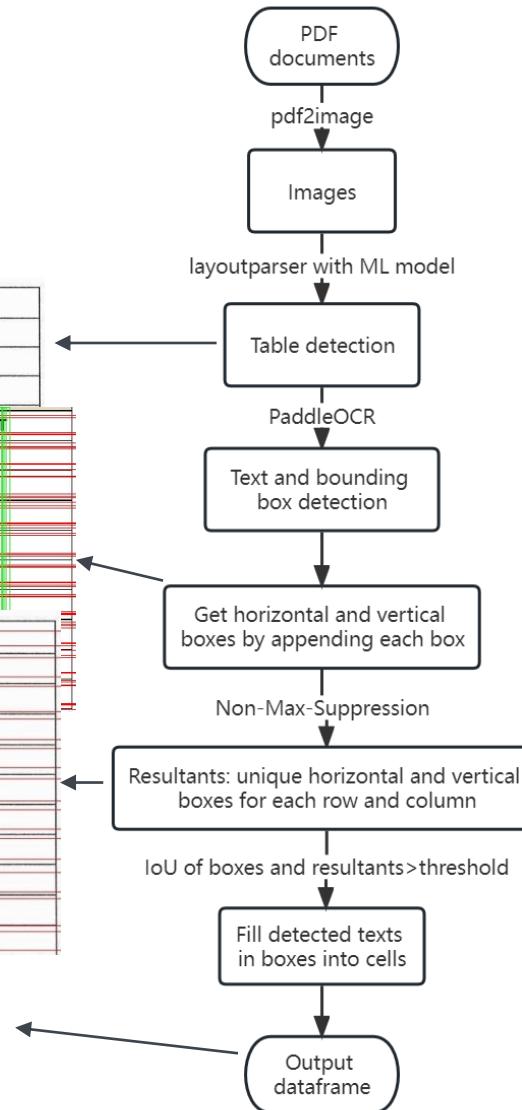
# Data Parser Conception



# Implementation

## OCR with Semantic Segmentation

| ABSOLUTE MAXIMUM RATINGS       |                      |                                 |                                |      |
|--------------------------------|----------------------|---------------------------------|--------------------------------|------|
| PARAMETER                      | TEST CONDITION       | SYMBOL                          | VALUE                          | UNIT |
| Supply voltage(pin 3)          |                      | V <sub>s</sub>                  | -0.3 to +6.0                   | V    |
| Supply voltage (pin 3)         |                      | I <sub>s</sub>                  | 5                              | mA   |
| Supply current (pin 3)         |                      |                                 | 0.3 to 5.5                     |      |
| Junction temperature (pin 1)   |                      | T <sub>j</sub>                  | 100                            | °C   |
| ABSOLUTE MAXIMUM RATINGS       |                      |                                 |                                |      |
| PARAMETER                      | TEST CONDITION       | SYMBOL                          | VALUE                          | UNIT |
| Supply voltage (pin 3)         |                      | V <sub>s</sub>                  | -0.3 to +6.0                   | V    |
| Supply current (pin 3)         |                      | I <sub>s</sub>                  | 5                              | mA   |
| Output voltage (pin 1)         |                      | V <sub>o</sub>                  | -0.3 to 5.5                    | V    |
| Voltage at output to supply    |                      | V <sub>s</sub> - V <sub>o</sub> | -0.3 to (V <sub>s</sub> + 0.3) | V    |
| Output current (pin 1)         |                      | I <sub>o</sub>                  | 5                              | mA   |
| Junction temperature           |                      | T <sub>j</sub>                  | 100                            | °C   |
| 1 ABSOLUTE MAXIMUM RATINGS     |                      |                                 |                                |      |
| 2 PARAMETER                    | TEST CONDITION       | SYMBOL                          | VALUE                          | UNIT |
| 3 Supply voltage(pin 3)        |                      | V <sub>s</sub>                  | -0.3 to +6.0                   | V    |
| 4 Supply current (pin 3)       |                      | I <sub>s</sub>                  | 5                              | mA   |
| 5 Output voltage (pin 1)       |                      | V <sub>o</sub>                  | -0.3 to 5.5                    | V    |
| 6 Voltage at output to supply  |                      | V <sub>s</sub> -V <sub>o</sub>  | -0.3 to V <sub>s</sub> +0.3    | V    |
| 7 Output current (pin 1)       |                      | I <sub>o</sub>                  | 5                              | mA   |
| 8 Junction temperature         |                      | T <sub>j</sub>                  | 100                            | °C   |
| 9 Storage temperature range    |                      | T <sub>stg</sub>                | -25 to +85                     |      |
| 10 Operating temperature range |                      | T <sub>amb</sub>                | -25 to +85                     |      |
| 11 Power consumption           | T <sub>amb</sub> 85C | P <sub>tot</sub>                | 10 mW                          |      |



# Implementation

## Text-Based Extraction with Camelot

-  output\_Cam-page-1-table-1.csv
-  output\_Cam-page-2-table-1.csv
-  output\_Cam-page-2-table-2.csv
-  output\_Cam-page-3-table-1.csv
-  output\_Cam-page-3-table-2.csv

|                            | A                           | B              | C       | D                  | E                     |
|----------------------------|-----------------------------|----------------|---------|--------------------|-----------------------|
| 1                          | Column1                     | Column2        | Column3 | Column4            | Column5               |
| 2                          |                             |                |         |                    | Vishay Semiconductors |
| 3 ABSOLUTE MAXIMUM RATINGS |                             |                |         |                    |                       |
| 4                          | PARAMETER                   | TEST CONDITION | SYMBOL  | VALUE              | UNIT                  |
| 5                          | Supply voltage (pin 3)      |                | VS      | -0.3 to +6.0       | V                     |
| 6                          | Supply current (pin 3)      |                | IS      |                    | 5 mA                  |
| 7                          | Output voltage (pin 1)      |                | VO      | -0.3 to 5.5        | V                     |
| 8                          | Voltage at output to supply |                | VS - VO | -0.3 to (VS + 0.3) | V                     |
| 9                          | Output current (pin 1)      |                | IO      |                    | 5 mA                  |
| 10                         | Junction temperature        |                | Tj      |                    | 100 °C                |
| 11                         | Storage temperature range   |                | Tstg    | -25 to +85         | °C                    |
| 12                         | Operating temperature range |                | Tamb    | -25 to +85         | °C                    |
| 13                         | Power consumption           | Tamb ≤ 85 °C   | Ptot    |                    | 10 mW                 |
| 14                         | Note                        |                |         |                    |                       |

Misrecognizes other objects as tables

Performances better in recognizing special symbols

# Implementation

## Cross Validation using Fuzzy Match

|    | A            |   | A  | B               |   |
|----|--------------|---|----|-----------------|---|
| 1  |              | 0 | 1  |                 | 0 |
| 2  | VPF44.50F15  |   | 2  | 0               |   |
| 3  |              |   | 3  | 1               |   |
| 4  | VPF54.50F15  |   | 4  | 2 VPF44.50F15   |   |
| 5  | VPF44.50F25  |   | 5  | 3               |   |
| 6  |              |   | 6  | 4 VPF44.50F15   |   |
| 7  | VPF54.50F25  |   | 7  | 5 VPF44.50F25   |   |
| 8  | VPF44.65F25  |   | 8  | 6               |   |
| 9  | VPF54.65F25  |   | 9  | 7 VPF54.50F25   |   |
| 10 | VPF44.65F35  |   | 10 | 8 VPF44.65F25   |   |
| 11 | VPF54.65F35  |   | 11 | 9 VPF54.65F25   |   |
| 12 | VPF44.80F35  |   | 12 | 10 VPF44.65F35  |   |
| 13 |              |   | 13 | 11 VPF54.65F35  |   |
| 14 | VPF54.80F35  |   | 14 | 12 VPF44.80F35  |   |
| 15 | VPF44.80F45  |   | 15 | 13              |   |
| 16 |              |   | 16 | 14 VPF54.80F35  |   |
| 17 | VPF54.80F45  |   | 17 | 15 VPF44.80F45  |   |
| 18 | VPF44.100F70 |   | 18 | 16              |   |
| 19 |              |   | 19 | 17 VPF54.80F45  |   |
| 20 | VPF54.100F70 |   | 20 | 18 VPF44.100F70 |   |
| 21 | VPF44.100F90 |   | 21 | 19              |   |
| 22 |              |   | 22 | 20 VPF54.100F70 |   |
| 23 | VPF54.100F90 |   | 23 | 21 VPF44.100F90 |   |
| 24 |              |   | 24 | 22              |   |
| 25 |              |   | 25 | 26 VPF54.100F90 |   |
| 26 |              |   | 26 |                 |   |

Camelot

OCR

|    | 0               | fuzzy match        | similarity score |
|----|-----------------|--------------------|------------------|
| 1  | VPF44.50F15     | VPF44.50F15        | 100              |
| 2  | VPF54.50F15     | VPF54.50F15        | 100              |
| 3  | VPF44.50F25     | VPF44.50F25        | 100              |
|    | A               | B                  | C                |
| 1  | OCR Column4     | fuzzy match        | similarity score |
| 2  | VALUE           | VALUE              | 100              |
| 3  | -0.3 to +6.0    | -0.3 to +6.0       | 100              |
| 4  | 5               | 5                  | 100              |
| 5  | -0.3 to 5.5     | 5                  | 100              |
| 6  | -0.3 to Vs+0.3) | -0.3 to (Vs + 0.3) | 94               |
| 7  | 5               | 5                  | 100              |
| 8  | 100             | 100                | 100              |
| 9  | -25 to +85      | -25 to +85         | 93               |
| 10 | -25 to +85      | -25 to +85         | 100              |
| 11 | 10              | 10                 | 100              |
|    | VPF54.100F90    | VPF54.100F90       | 92               |

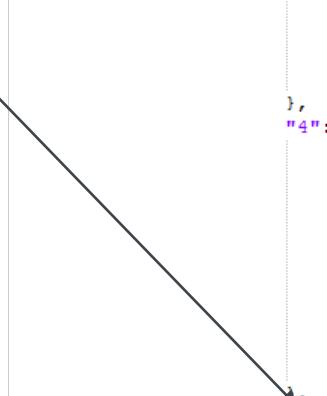
Fuzzy match correction

- ✓ After fuzzy match, 'O' in OCR result is corrected to '0'.
- ✓ White spaces and some symbols can be corrected.

# Implementation

## Integration into Syntax Tree

A hierarchical data structure that represents the syntactic structure of text in a tree-like form.



**VISHAY** www.vishay.com Vishay Semiconductors

**ABSOLUTE MAXIMUM RATINGS**

| PARAMETER                   | TEST CONDITION                  | SYMBOL                  | VALUE        | UNIT |
|-----------------------------|---------------------------------|-------------------------|--------------|------|
| Supply voltage (pin 3)      |                                 | $V_S$                   | -0.3 to +6.0 | V    |
| Supply current (pin 3)      |                                 | $I_S$                   | 5            | mA   |
| Output voltage (pin 1)      |                                 | $V_O$                   | -0.3 to 5.5  | V    |
| Voltage at output to supply | $V_S = V_O$                     | -0.3 to ( $V_S + 0.3$ ) | V            |      |
| Output current (pin 1)      |                                 | $I_O$                   | 5            | mA   |
| Junction temperature        |                                 | $T_J$                   | 100          | °C   |
| Storage temperature range   |                                 | $T_{STG}$               | -25 to +85   | °C   |
| Operating temperature range |                                 | $T_{OPR}$               | -25 to +85   | °C   |
| Power consumption           | $T_{AMB} \leq 85^\circ\text{C}$ | $P_{DET}$               | 10           | mW   |

**Note**

- Stresses beyond those listed under "Absolute Maximum Ratings" may cause permanent damage to the device. This is a stress rating only and functional operation of the device at these or any other conditions beyond those indicated in the operational sections of this specification is not implied. Exposure to absolute maximum rating conditions for extended periods may affect the device reliability.

**ELECTRICAL AND OPTICAL CHARACTERISTICS ( $T_{AMB} = 25^\circ\text{C}$ , unless otherwise specified)**

| PARAMETER                  | TEST CONDITION   | SYMBOL       | MIN. | TYP.     | MAX. | UNIT              |
|----------------------------|--|--------------|------|----------|------|-------------------|
| Supply current (pin 3)     | $E_V = 0$ , $V_S = 5\text{V}$  | $I_{SO}$     | 0.55 | 0.7      | 0.9  | mA                |
|                            | $E_V = 40\text{ kV}$ , sunlight  | $I_{SR}$     | -    | 0.8      | -    | mA                |
| Supply voltage             |  | $V_S$        | 2.5  | -        | 5.5  | V                 |
| Transmission distance      | $E_V = 0$ , test signal see Fig. 1,<br>IR diode TSAL6200,<br>$I_R = 50\text{ mA}$  | $d$          | -    | 12       | -    | m                 |
| Output voltage low (pin 1) | $I_{OLH} = 0.5\text{ mA}$ , $E_V = 2\text{ mV}/\text{mW}$ ,<br>$I_{OLL} = 0.5\text{ mA}$ , $E_V = 2\text{ mV}/\text{mW}$ ,<br>$T_{AMB} = 25^\circ\text{C}$ | $V_{O(L)}$   | -    | -        | 100  | mV                |
| Minimum irradiance         | Pulse width tolerance:<br>$t_{p1} = 5t_0$ , $t_{p2} < t_{p1} + 6t_0$ ,<br>test signal see Fig. 1   | $E_{I,min.}$ | -    | 0.4      | 0.7  | mW/m <sup>2</sup> |
| Maximum irradiance         | $t_{p1} = 5t_0$ , $t_{p2} > t_{p1} + 6t_0$ ,<br>test signal see Fig. 1   | $E_{I,max.}$ | 50   | -        | -    | W/m <sup>2</sup>  |
| Directivity                | Angle of half transmission<br>distance   | $\phi_{1/2}$ | -    | $\pm 45$ | -    | deg               |

Rev. 1.6, 09-Jul-2021 2 Document Number: 82780

THIS DOCUMENT IS SUBJECT TO CHANGE WITHOUT NOTICE. THE PRODUCTS DESCRIBED HEREIN AND THIS DOCUMENT ARE SUBJECT TO SPECIFIC DISCLAIMERS, SET FORTH AT [www.vishay.com/doc/20704](http://www.vishay.com/doc/20704)

```
"3": {
    "0": "VALUE",
    "1": "-0.3 to +6.0",
    "2": "5",
    "3": "-0.3 to 5.5",
    "4": "-0.3 to (Vs+0.3)",
    "5": "5",
    "6": "100",
    "7": "-25 to +85",
    "8": "-25 to +85",
    "9": "10"
},
"4": {
    "0": "UNIT",
    "1": "mA",
    "2": "mA",
    "3": "V",
    "4": "V",
    "5": "mA",
    "6": "°C",
    "7": "°C",
    "8": "°C",
    "9": "mW"
},
"location": {
    "0": "0.22",
    "1": 1,
    "2": null,
    "3": null,
    "4": null,
    "5": null,
    "6": null,
    "7": null,
    "8": null,
    "9": null
},
```

Data frame format: index

Data frame format: columns

# Evaluation

## Table Detection

|                          | Number of detected tables | Number of real tables | Extraction rate $\varphi$ | Completeness $\chi$ |
|--------------------------|---------------------------|-----------------------|---------------------------|---------------------|
| <b>Sensor</b>            |                           |                       |                           |                     |
| Document 1 (Temperature) | 4                         | 4                     | 1.00                      | 0.94                |
| Document 2 (Reflective)  | 2                         | 3                     | 0.67                      | 0.97                |
| <b>Actuator</b>          |                           |                       |                           |                     |
| Document 3 (Camera)      | 1                         | 1                     | 1.00                      | 1.00                |
| Document 4 (Motor)       | 0                         | 1                     | 0.00                      | 0.00                |
| <b>Controller</b>        |                           |                       |                           |                     |
| Document 5 (Nebula)      | 1                         | 1                     | 1.00                      | 1.00                |
| Document 6 (DIN)         | 2                         | 2                     | 1.00                      | 1.00                |

- Most tables are detected with high completeness.
- Some missing lines is tolerable.

### Failed cases:

- No ruling lines at all.
- Objects around the table have similar structure.

# Evaluation

## Table Structure

|                   | Total adjacency relations | Detected adjacency relations | Correct adjacency relations | Recall β | Precision α | F1-score |
|-------------------|---------------------------|------------------------------|-----------------------------|----------|-------------|----------|
| <b>Document 1</b> |                           |                              |                             |          |             |          |
| Table 1           | 94                        | 92                           | 68                          | 0.723    | 0.739       | 0.731    |
| Table 2           | 52                        | 52                           | 52                          | 1.000    | 1.000       | 1.000    |
| Table 3           | 31                        | 39                           | 26                          | 0.839    | 0.667       | 0.743    |
| Table 4           | 108                       | 110                          | 100                         | 0.926    | 0.909       | 0.918    |
| <b>Document 2</b> |                           |                              |                             |          |             |          |
| Table 1           | -                         | -                            | -                           | -        | -           | -        |
| Table 2           | 85                        | 85                           | 85                          | 1.000    | 1.000       | 1.000    |
| Table 3           | 109                       | 127                          | 94                          | 0.862    | 0.740       | 0.797    |
| <b>Document 3</b> |                           |                              |                             |          |             |          |
| Table 1           | 58                        | 58                           | 58                          | 1.000    | 1.000       | 1.000    |
| <b>Document 4</b> |                           |                              |                             |          |             |          |
| -                 | -                         | -                            | -                           | -        | -           | -        |
| <b>Document 5</b> |                           |                              |                             |          |             |          |
| Table 1           | 31                        | 33                           | 29                          | 0.935    | 0.879       | 0.906    |
| <b>Document 6</b> |                           |                              |                             |          |             |          |
| Table 1           | 186                       | 167                          | 137                         | 0.737    | 0.820       | 0.776    |
| Table 2           | 40                        | 40                           | 40                          | 1.000    | 1.000       | 1.000    |

- Almost no error with a neat structure.

### Failed cases:

- Merging rows/columns.
- Multiple lines in one cell.

ns in Predicted Table

|                        |  |
|------------------------|--|
| Supply current (pin 3) | $E_v = 0, V_S = 5 \text{ V}$   |
|                        | $E_v = 40 \text{ klx, sunlight}$   |
| Supply voltage         |  |
| Transmission distance  | $E_v = 0, \text{test signal see Fig. 1,}$<br>IR diode TSAL6200,<br>$I_F = 50 \text{ mA}$ |

# Evaluation

## Table Content

|                             | Levenshtein distance | Number of true characters | Error rate | Number of corrected strings | Error rate after correction |
|-----------------------------|----------------------|---------------------------|------------|-----------------------------|-----------------------------|
| <b>Sensor</b>               |                      |                           |            |                             |                             |
| Document 1<br>(Temperature) | 74                   | 1488                      | 5.0%       | 11                          | 4.2%                        |
| Document 2<br>(Reflective)  | 71                   | 852                       | 8.3%       | 9                           | 7.3%                        |
| <b>Actuator</b>             |                      |                           |            |                             |                             |
| Document 3<br>(Camera)      | 103                  | 720                       | 14.3%      | 10                          | 12.9%                       |
| Document 4<br>(Motor)       | -                    | -                         | -          | -                           | -                           |
| <b>Controller</b>           |                      |                           |            |                             |                             |
| Document 5<br>(Nebula)      | 5                    | 637                       | 0.8%       | 2                           | 0.3%                        |
| Document 6<br>(DIN)         | 3                    | 449                       | 0.7%       | 2                           | 0.2%                        |

## Main Errors:

- Missing information.
- Special symbols such as  $\circ$ ,  $\leq$ .

|                    |  |
|--------------------|--|
| Maximum irradiance | $t_{pi} - 5/f_0 < t_{po} < t_{pi} + 6/f_0$ ,<br>test signal see Fig. 1 |
|--------------------|--|



|                    |  |
|--------------------|--|
| Maximum irradiance | $t_{pi} - 5/f_0 < t_{po} < t_{pi} + 6/f_0$ |
|--------------------|--|

# Summary and Outlook

## Functions:

- Data parsing from text-based and image-based PDFs
- Automatically key-value-pairs extraction from tables
- Loop through each page of PDF
- Well-structured multiple output formats available

## Deficiencies:

- Relatively low accuracy for complex table
- Manual error correction
- Limited metadata provided

## Outlook:

- NER and semantic analysis
- Store in database for user's query

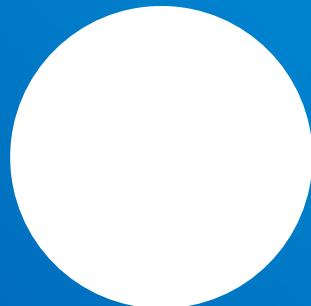
# Quelle

1. <https://www.tibco.com/reference-center/what-is-data-parsing>
2. <https://www.boschmanufacturingsolutions.com/>
3. <https://www.dvz.de/rubriken/logistik/detail/news/industrie-40-bosch-schafft-einheit-fuer-software-und-services.html>
4. G. Endignoux, O. Levillain, J.Y. Migeon, “Caradoc: a pragmatic approach to PDF parsing and validation” in IEEE Security and Privacy Workshops, 2016, pp.126-139.
5. <https://products.aspose.app/ocr/de/pdf-ocr>
6. <https://www.imranabdullah.com/2021-09-17/Fuzzy-word-replace-from-string-in-Python>
7. <https://www.twilio.com/blog/abstract-syntax-trees>



**University of Stuttgart**  
Institut of Industrial Automation  
and Software Engineering

# Thank you!



**Jiayi He**

e-mail [st176731@stud.uni-stuttgart.de](mailto:st176731@stud.uni-stuttgart.de)

phone +49 (0) 711 685-

fax +49 (0) 711 685-

University of Stuttgart



# Basis

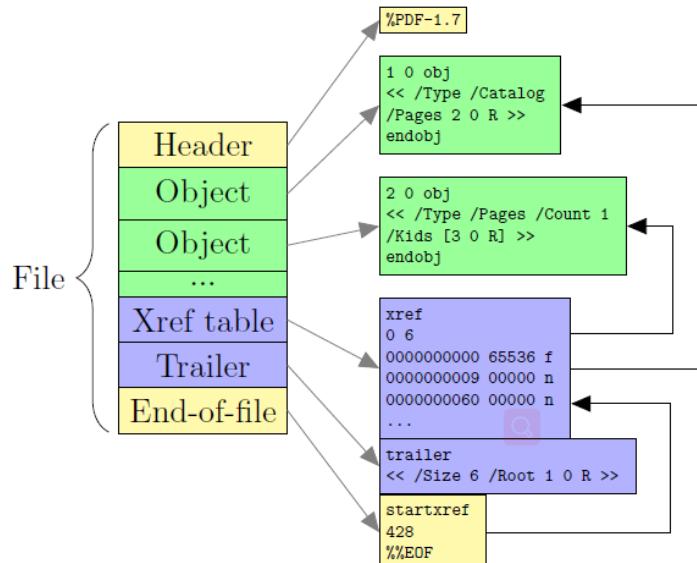
## Text Mining



- Preprocessing: filtering, stemming, lemmatizing.
- Text mining methods: NLP, Named Entity Recognition, information extraction.
- Text analysis: semantic analysis.
- Discovery of knowledge: stored in knowledge database.

## PDF Syntax-Postscript

A page description language, provides internal structure of documents.



[5]

### Pros:

- Preserves the file formatting.
- Metadata accessible.

### Cons:

- Requires technical expertise to work with.
- Limited supporting tools and researches.

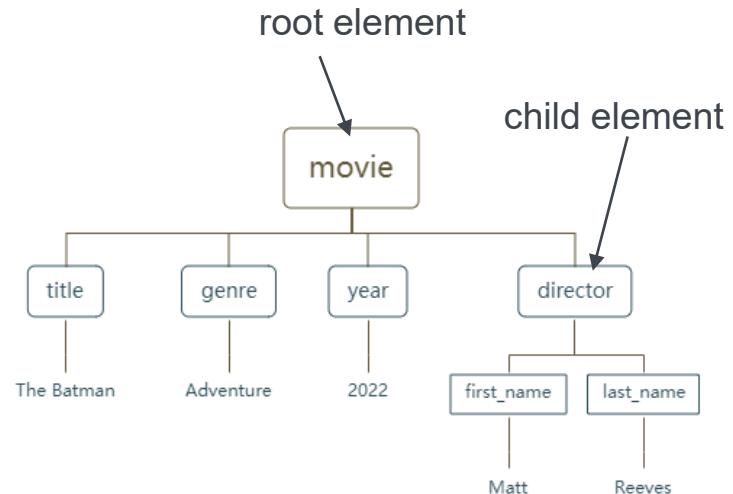
# Makeup Language Conversion

## Pros:

- Structured data, helps to identify different elements of the PDF document.
- Use tags and attributes to define elements, accurately extract symbols and units.

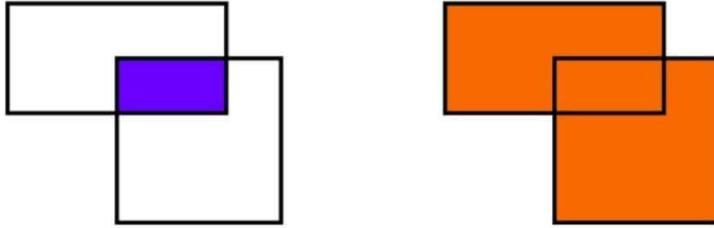
## Cons:

- Conversion can be difficult.  
(especially for image-based PDF)
- Increased system complexity.
- Data loss.



```
<movie lang="English">  
  <title>The Batman</title>  
  <genre>Adventure</genre>  
  ...  
</movie>
```

# NMS



- Firstly, the boxes whose confidence are smaller than a threshold should be removed.
- Secondly, the box with a lower confidence among the boxes which overlap too much with each other ( $\text{IoU} > 0.1$ ) should also be removed.
- At last, the process is repeated until only one unique box is left for each class, which represents the final prediction.