

# Extraction of Solar Cell Data from PDF Datasheets

## Master Thesis Presentation

Muhammad Moez Malik

Albert-Ludwigs-Universität Freiburg

*moezmalik13@gmail.com*

April 14, 2023

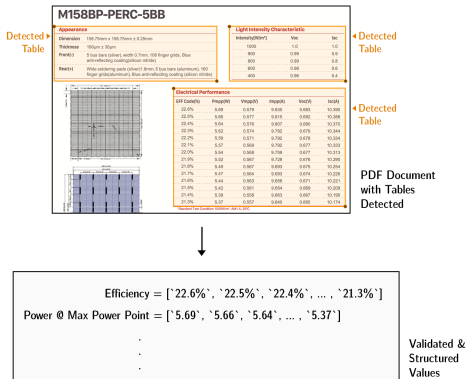
# Content

- ① Introduction
- ② Approach
- ③ Experiments
- ④ Conclusions
- ⑤ References
- ⑥ Appendix

# Introduction

# Introduction

What is the thesis about?



- Automatic extraction of solar cell data from data sheets

# Introduction: Motivation

## The Need for an Automated Solution

- Automatic generation of indexed and searchable database
- Process a large number of documents quickly
- Prediction of market trends using the data
- Uncovering possible research areas by analysis of data

# Introduction: PDF Documents

A Standard for Digital Document Distribution. Why?

- Pros

- Compact Files
- Industry standard
- Same rendering of documents regardless of hardware and software.
- Ensures information reproducibility

- Cons

- Document is constructed using only positional information, the row and column information for tables is lost



PDF  
Document

# Introduction: Solar Cell Data Sheets

Distributed as PDF Documents

## CONSOIT®

Mono 5BB Solar Cell 158.75 Bifacial

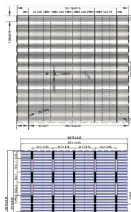
### M158BP-PERC-5BB

#### Appearance

Dimension	158.75mm x 158.75mm $\pm$ 0.25mm
Thickness	190 $\mu$ m $\pm$ 30 $\mu$ m
Front(-)	5 bus bars (silver), width 0.7mm, 106 finger grids, Blue anti-reflecting coating(silicon nitride)
Rear(+)	Wide soldering pads (silver)1.8mm, 5 bus bars (aluminum), 160 finger grids(aluminum), Blue anti-reflecting coating (silicon nitride)

#### Light Intensity Characteristic

Intensity(W/m <sup>2</sup> )	Voc	Isc
1000	1.0	1.0
900	0.99	0.9
800	0.99	0.8
600	0.98	0.6
400	0.96	0.4



#### Electrical Performance

EFF Code(%)	Pmpp(W)	Vmpp(V)	Ippp(A)	Voc(V)	Isc(A)
22.6%	5.69	0.579	9.835	0.683	10.390
22.5%	5.66	0.577	9.815	0.682	10.388
22.4%	5.64	0.576	9.807	0.680	10.370
22.3%	5.62	0.574	9.792	0.679	10.344
22.2%	5.59	0.571	9.792	0.678	10.334
22.1%	5.57	0.569	9.792	0.677	10.333
22.0%	5.54	0.568	9.759	0.677	10.313
21.9%	5.52	0.567	9.728	0.676	10.290
21.8%	5.49	0.567	9.693	0.675	10.254
21.7%	5.47	0.564	9.693	0.674	10.226
21.6%	5.44	0.563	9.658	0.671	10.221
21.5%	5.42	0.561	9.654	0.669	10.209
21.4%	5.39	0.558	9.663	0.667	10.195
21.3%	5.37	0.557	9.640	0.665	10.174

\*Standard Test Condition 1000W/m<sup>2</sup>, AM1.5, 25°C

Figure: An example of solar cell data sheet

# Introduction: Solar Cell Data Sheets

Distributed as PDF Documents

## CONSOIT®

Mono 5BB Solar Cell 158.75 Bifacial

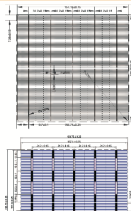
### M158BP-PERC-5BB

#### Appearance

Dimension	158.75mm x 158.75mm ± 0.25mm
Thickness	190µm ± 30µm
Front(-)	5 bus bars (silver), width 0.7mm, 106 finger grids, Blue anti-reflecting coating(silicon nitride)
Rear(+)	Wide soldering pads (silver)1.8mm, 5 bus bars (aluminum), 160 finger grids(aluminum), Blue anti-reflecting coating (silicon nitride)

#### Light Intensity Characteristic

Intensity(W/m²)	Voc	Isc
1000	1.0	1.0
900	0.99	0.9
800	0.99	0.8
600	0.98	0.6
400	0.96	0.4



#### Electrical Performance

EFF Code(%)	Pmpp(W)	Vmpp(V)	Ippp(A)	Voc(V)	Isc(A)
22.6%	5.69	0.579	9.835	0.683	10.390
22.5%	5.66	0.577	9.815	0.682	10.388
22.4%	5.64	0.576	9.807	0.680	10.370
22.3%	5.62	0.574	9.792	0.679	10.344
22.2%	5.59	0.571	9.792	0.678	10.334
22.1%	5.57	0.569	9.792	0.677	10.333
22.0%	5.54	0.568	9.759	0.677	10.313
21.9%	5.52	0.567	9.728	0.676	10.290
21.8%	5.49	0.567	9.693	0.675	10.254
21.7%	5.47	0.564	9.693	0.674	10.226
21.6%	5.44	0.563	9.658	0.671	10.221
21.5%	5.42	0.561	9.654	0.669	10.209
21.4%	5.39	0.558	9.663	0.667	10.195
21.3%	5.37	0.557	9.640	0.665	10.174

\*Standard Test Condition 1000W/m², AM1.5, 25°C

Figure: Critical information is in tables



# Introduction: Challenges

## Variability in Data Sheet Design

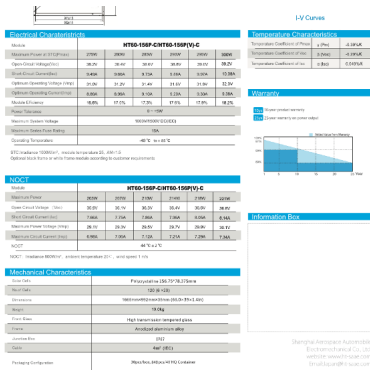


Figure: Data sheet design variation 1

# Introduction: Challenges

## Variability in Data Sheet Design

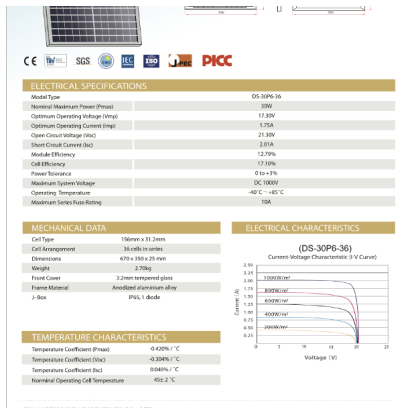


Figure: Data sheet design variation 2



# Introduction: Challenges

## Variability in Data Sheet Design

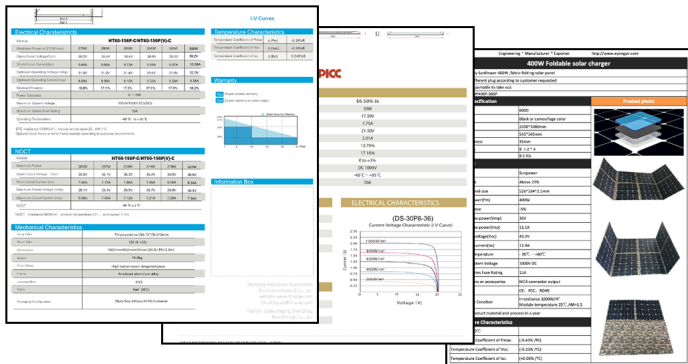


Figure: Data sheet design variations

# Introduction: Existing Solutions

## Shortcomings of Existing Off-the-Shelf Solutions

- Existing rule-based solutions
  - Tabula [1]
  - Camelot [2]
- Problems uncovered [3]
  - Failure in detection of tables
  - Overlap of detected table areas
  - Failure to extract required values

# Introduction: Summary

- Existing off-the-shelf solutions fail on solar cell data sheets
- The variability in the data sheet design presents a challenge
- Manual labour not feasible for large quantities of documents
- **Goal:** An end-to-end solution for automatic extraction of solar cell data from tables in the PDF data sheets

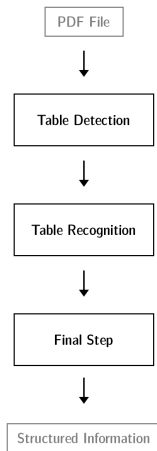
Questions?

# Approach



# Approach: Overview

## Major Steps Involved in Approach



### ① Table Detection

Locating where the tables are

### ② Table Recognition

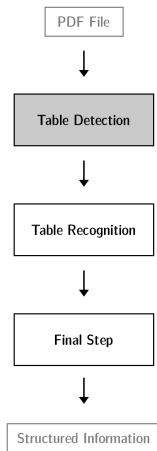
Extracting the raw values from the tables

### ③ Final Step

Validating and structuring extracted raw values

# Approach: Overview

## Major Steps Involved in Approach



### ① Table Detection

Locating where the tables are

### ② Table Recognition

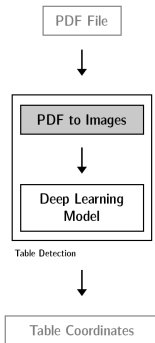
Extracting the raw values from the tables

### ③ Final Step

Validating and structuring extracted raw values

# Approach: Table Detection

Locating tables in the documents

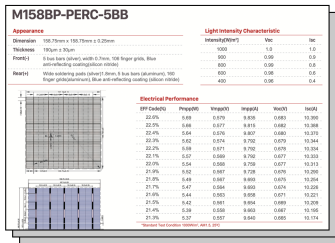


- Pages in PDF are converted to images



PDF Document

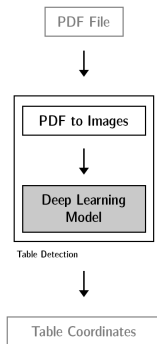
Convert Pages  
of PDF to Images



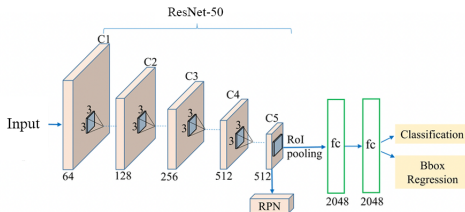
Images  
from PDF

# Approach: Table Detection

Locating tables in the documents

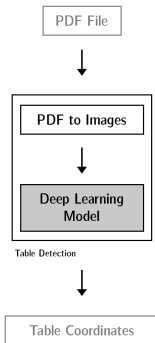


- An example Object Detector architecture (FasterRCNN) [4]



# Approach: Table Detection

Locating tables in the documents



- Detected Tables are saved as CSV

The screenshot shows a PDF document titled 'M158BP-PERC-5BB'. It contains several tables and a figure. The tables are highlighted with orange borders and labels. The labels are: 'Detected Table' (pointing to the 'Appearance' table), 'Detected Table' (pointing to the 'Light Intensity Characteristics' table), 'Detected Table' (pointing to the 'Electrical Performance' table), and 'PDF Document with Tables Detected' (pointing to the entire document area).

**Appearance**

Dimension	Value
158.75mm x 100.75mm x 0.25mm	158.75

**Light Intensity Characteristics**

Intensity (W/m²)	Watt	Inc.
1000	1.0	1.4
900	0.9	0.6
800	0.8	0.6
400	0.4	0.6
200	0.2	0.6

**Electrical Performance**

Eff. Coeff (%)	Power (%)	Watt (%)	Inc (%)	Watt (%)	Inc (%)
22.7%	0.88	0.576	0.888	0.888	10.288
22.4%	0.88	0.577	0.875	0.882	10.288
22.4%	0.88	0.576	0.867	0.888	10.275
22.7%	0.88	0.574	0.792	0.878	10.284
22.2%	0.88	0.571	0.792	0.878	10.284
22.7%	0.87	0.568	0.792	0.877	10.284
22.2%	0.84	0.568	0.792	0.877	10.284
21.4%	0.82	0.568	0.792	0.877	10.284
21.4%	0.48	0.567	0.588	0.875	10.274
21.7%	0.47	0.584	0.588	0.874	10.274
21.4%	0.44	0.583	0.588	0.874	10.274
21.2%	0.42	0.581	0.584	0.888	10.288
21.4%	0.38	0.588	0.588	0.887	10.287
21.2%	0.37	0.587	0.588	0.888	10.274

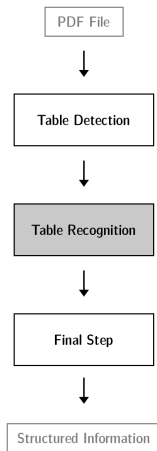
**Detection using Object Detectors**

filename, pageno, x1, y1, x2, y2  
1-4..pdf, 1, 252, 289, 575, 215  
1-4..pdf, 1, 254, 538, 576, 325  
1-4..pdf, 1, 354, 647, 576, 556  
1-4..pdf, 1, 587, 237, 212, 198  
1-4..pdf, 1, 38, 647, 335, 557

CSV File containing Table Coordinates

# Approach: Overview

## Major Steps Involved in Approach



### ① Table Detection

Locating where the tables are

### ② Table Recognition

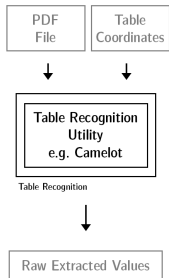
Extracting the raw values from the tables

### ③ Final Step

Validating and structuring extracted raw values

# Approach: Table Recognition

Extracting raw values from table



PDF Document

+

```
filename,pageNo,x1,y1,x2,y2
1~4,.pdf,1,252,289,575,215
1~4,.pdf,1,254,530,576,325
1~4,.pdf,1,359,647,576,556
1~4,.pdf,1,50,237,212,190
1~4,.pdf,1,38,647,335,557
```

CSV File containing  
Table Coordinates

Extract Raw  
Table Values

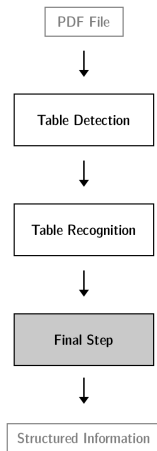
	A	B	C	D	E	F
1	EFF Code(% Pmpp(W)	Vmpp(V)	Impp(A)	Voc(V)	Isc(A)	
2	22.6%	5.69	0.579	9.835	0.683	10.390
3	22.5%	5.66	0.577	9.815	0.682	10.388
4	22.4%	5.64	0.576	9.807	0.680	10.370
5	22.3%	5.62	0.574	9.792	0.679	10.344
6	22.2%	5.59	0.571	9.792	0.678	10.334
7	22.1%	5.57	0.569	9.792	0.677	10.333
8	22.0%	5.54	0.568	9.759	0.677	10.313
9	21.9%	5.52	0.567	9.728	0.676	10.290
10	21.8%	5.49	0.567	9.693	0.675	10.254
11	21.7%	5.47	0.564	9.693	0.674	10.226
12	21.6%	5.44	0.563	9.658	0.671	10.221
13	21.5%	5.42	0.561	9.654	0.669	10.209
14	21.4%	5.39	0.558	9.663	0.667	10.195
15	21.3%	5.37	0.557	9.640	0.665	10.174

Raw Values  
Extracted  
from Tables

Figure: Raw values are extracted from tables

# Approach: Overview

## Major Steps Involved in Approach



### ① Table Detection

Locating where the tables are

### ② Table Recognition

Extracting the raw values from the tables

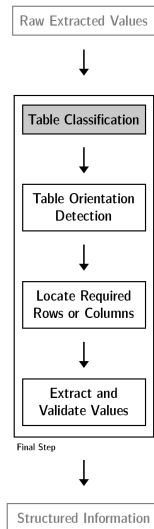
### ③ Final Step

Validating and structuring extracted raw values



# Approach: Final Step

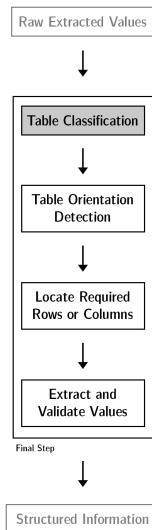
## Table Classification



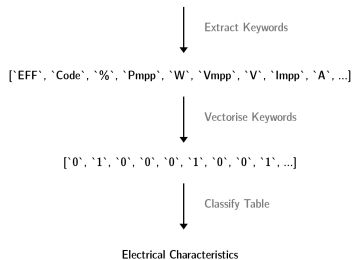
- Table type needs to be identified before the values can be validated
- The table title might be missing or not detected
- Can use the table content to identify the table type

# Approach: Final Step

## Table Classification

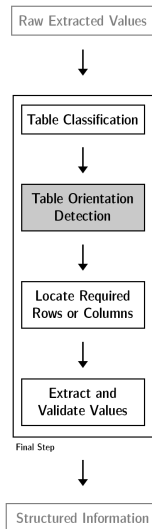


EFF Code(%)	Pmpp(W)	Vmpp(V)	Impp(A)	Voc(V)	Isc(A)
22.6%	5.69	0.579	9.835	0.683	10.390
22.5%	5.66	0.577	9.815	0.682	10.388
22.4%	5.64	0.576	9.807	0.680	10.370
22.3%	5.62	0.574	9.792	0.679	10.344
22.2%	5.59	0.571	9.792	0.678	10.334
22.1%	5.57	0.569	9.792	0.677	10.333
22.0%	5.54	0.568	9.759	0.677	10.313
21.9%	5.52	0.567	9.728	0.676	10.290
21.8%	5.49	0.567	9.693	0.675	10.254
21.7%	5.47	0.564	9.693	0.674	10.226
21.6%	5.44	0.563	9.658	0.671	10.221
21.5%	5.42	0.561	9.654	0.669	10.209
21.4%	5.39	0.558	9.663	0.667	10.195
21.3%	5.37	0.557	9.640	0.665	10.174



# Approach: Final Step

## Table Orientation Detection



Horizontal Table

Efficiency Code		196	195	194	193	192	191	190	189
Efficiency	Eff(%)	19.60	19.50	19.40	19.30	19.20	19.10	19.00	18.90
Power	P <sub>mp</sub> (W)	4.68	4.66	4.64	4.61	4.59	4.56	4.54	4.52
Max. Power Current	I <sub>pm</sub> (A)	8.64	8.61	8.58	8.55	8.53	8.51	8.49	8.48
Short Circuit Current	I <sub>sc</sub> (A)	9.14	9.11	9.08	9.05	9.03	9.02	9.01	9.01
Max. Power Voltage	V <sub>mp</sub> (V)	0.542	0.541	0.541	0.539	0.538	0.536	0.535	0.533
Open Circuit Voltage	V <sub>oc</sub> (V)	0.643	0.643	0.642	0.641	0.641	0.640	0.639	0.638

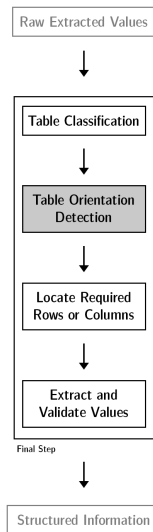
Vertical Table

ELECTRICAL PARAMETERS					
Efficiency Ncell(%)	Maximum power P <sub>mp</sub> (w)	Peak voltage V <sub>mp</sub> (V)	Peak current I <sub>mp</sub> (A)	Open-circuit voltage V <sub>oc</sub> (V)	Short-circuit current I <sub>sc</sub> (A)
≥18.6	4.53	0.540	8.389	0.641	8.889
18.4-18.6	4.48	0.539	8.311	0.641	8.802
18.2-18.4	4.43	0.538	8.234	0.640	8.729
18.0-18.2	4.38	0.536	8.172	0.639	8.655
17.8-18.0	4.33	0.534	8.109	0.638	8.580
17.6-17.8	4.28	0.531	8.060	0.636	8.518
17.4-17.6	4.23	0.528	8.011	0.634	8.456

- Table can be vertical or horizontal
- Need to detect before the next steps

# Approach: Final Step

## Table Orientation Detection



Eff (%)	17.9	17.8	17.7	17.6	17.5	17.4	17.3	17.2
Pmpp (W)	4.359	4.334	4.308	4.283	4.259	4.234	4.210	4.186
Voc(V)	0.630	0.629	0.629	0.628	0.627	0.626	0.626	0.624
Isc(A)	8.760	8.747	8.709	8.680	8.656	8.648	8.626	8.573
Ump(V)	0.530	0.529	0.527	0.524	0.524	0.524	0.523	0.520
Iimp(A)	8.233	8.194	8.168	8.173	8.133	8.088	8.053	8.052
FF(%)	78.973	78.785	78.639	78.567	78.441	78.240	77.970	78.290

Extract Values

17.9	17.8	17.7	17.6	17.5	17.4	17.3	17.2
4.359	4.334	4.308	4.283	4.259	4.234	4.210	4.186
0.630	0.629	0.629	0.628	0.627	0.626	0.626	0.624
8.760	8.747	8.709	8.680	8.656	8.648	8.626	8.573
0.530	0.529	0.527	0.524	0.524	0.524	0.523	0.520
8.233	8.194	8.168	8.173	8.133	8.088	8.053	8.052
78.973	78.785	78.639	78.567	78.441	78.240	77.970	78.290

Calculate Variances

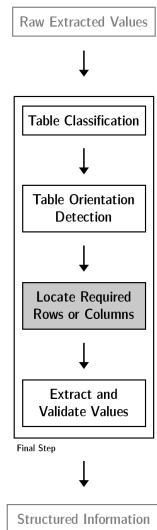
row-variance = 0.1755  
column-variance = 6173.95

Decide Orientation

Horizontal Table

# Approach: Final Step

## Locating Required Rows or Columns



- Locate the required rows or columns using Regex pattern matching

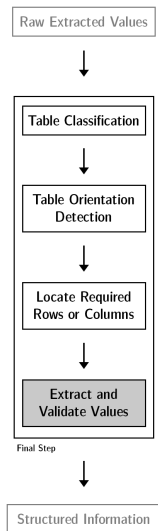
Eff (%)	17.9	17.8	17.7	17.6	17.5	17.4	17.3	17.2
Pmpp (W)	4.359	4.334	4.308	4.283	4.259	4.234	4.210	4.186
Uoc(V)	0.630	0.629	0.629	0.628	0.627	0.626	0.626	0.624
Isc(A)	8.760	8.747	8.709	8.680	8.656	8.648	8.626	8.573
Umpp(V)	0.530	0.529	0.527	0.524	0.524	0.524	0.523	0.520
Impp(A)	8.233	8.194	8.168	8.173	8.133	8.088	8.053	8.052
FF(%)	78.973	78.785	78.639	78.567	78.441	78.240	77.970	78.290

Locate Open Circuit Voltage Row

Uoc(V)	0.630	0.629	0.629	0.628	0.627	0.626	0.626	0.624
--------	-------	-------	-------	-------	-------	-------	-------	-------

# Approach: Final Step

## Validating and structuring raw values



- Extract the valid values from the located row

$V_{oc}(V)$	0.630	0.629	0.629	0.628	0.627	0.626	0.626	0.624
-------------	-------	-------	-------	-------	-------	-------	-------	-------

Extract Values

Open Circuit Voltage = ['0.63', '0.629', '0.629', '0.628', '0.627', ...]

Questions?

# Experiments



# Experiments: General Evaluation Metrics

		Actual	
		Positive	Negative
Predicted	Positive	True Positive (TP)	False Positive (FP)
	Negative	False Negative (FN)	True Negative (TN)

- Confusion matrix values for each class will be calculated

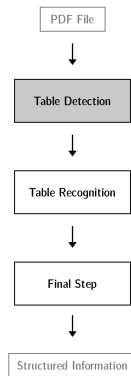
# Experiments: General Evaluation Metrics

- Confusion matrix values for each class will be calculated
  - **TP**: Actually positive, detected as positive
  - **FP**: Actually negative, detected as positive
  - **FN**: Actually positive, detected as negative
- Precision =  $\frac{TP}{TP+FP}$
- Recall =  $\frac{TP}{TP+FN}$
- F1-Score =  $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

# Experiments: Table Detection Evaluation

## Overview

This experiment evaluates the approach described for the table detection step.



# Experiments: Table Detection Evaluation

## Evaluation Metrics

- IoU Overlap Threshold
  - 75% IoU Overlap Threshold
  - 90% IoU Overlap Threshold

# Experiments: Table Detection Evaluation

## Evaluation Metrics

- IoU Overlap Threshold
  - 75% IoU Overlap Threshold
  - 90% IoU Overlap Threshold
- Precision =  $\frac{TP}{TP+FP}$ 

How many guesses made by model were actually tables?
- Recall =  $\frac{TP}{TP+FN}$ 

How many actual tables were detected by the model?

# Experiments: Table Detection Evaluation

## Evaluation Metrics

- IoU Overlap Threshold
  - 75% IoU Overlap Threshold
  - 90% IoU Overlap Threshold
- Precision =  $\frac{TP}{TP+FP}$ 

How many guesses made by model were actually tables?
- Recall =  $\frac{TP}{TP+FN}$ 

How many actual tables were detected by the model?

### Note

For selecting the best model Recall performance at 90% IoU threshold was considered

# Experiments: Table Detection Evaluation

## Dataset

- Randomly selected PDF documents
- Images: 2675 in total
- Tables: 5896 in total
- The tables were manually labelled

# Experiments: Table Detection Evaluation

## Setup

- Models
  - Single Stage Detectors
    - RetinaNet [5]
    - RetinaNet v2 [5]
  - Two Stage Detectors
    - FasterRCNN [6]
    - FasterRCNN v2 [6]



# Experiments: Table Detection Evaluation

## Setup

- Models
  - Single Stage Detectors
    - RetinaNet [5]
    - RetinaNet v2 [5]
  - Two Stage Detectors
    - FasterRCNN [6]
    - FasterRCNN v2 [6]
- Method
  - Trained for 50 epochs
  - Trained on NVidia P100 GPUs on Kaggle
  - Initial weights: Microsoft COCO [7]
  - Hyper-parameter Tuning
    - Batch Size
    - Learning Rate
  - Training-Evaluation Split
    - 80% for training
    - 20% for evaluation

# Experiments: Table Detection Evaluation

## Results

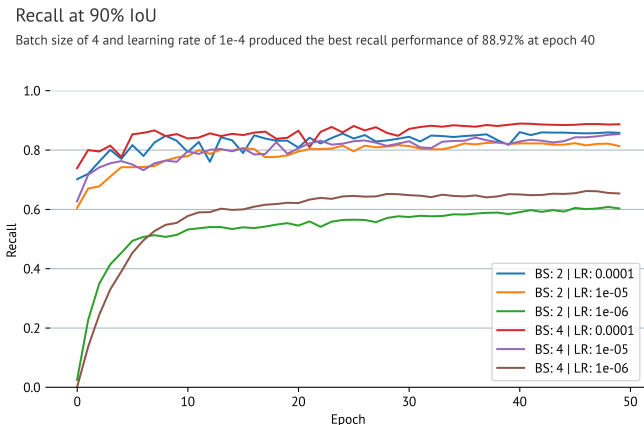


Figure: FasterRCNN V2 Hyperparameter Tuning

# Experiments: Table Detection Evaluation

## Results

Architecture	Configuration		Epoch	Recall		Precision	
	BS	LR		90% IoU	75% IoU	90% IoU	75% IoU
RetinaNet	4	1e-04	47	0.8007	0.9247	0.8155	0.9419
RetinaNet v2	4	1e-04	34	0.8007	0.9280	0.8259	0.9573
FasterRCNN	2	1e-04	23	0.8817	0.9520	0.8825	0.9528
FasterRCNN v2	4	<b>1e-04</b>	<b>40</b>	<b>0.8892</b>	<b>0.9512</b>	<b>0.9018</b>	<b>0.9648</b>

**Table:** Object Detector Model Comparison

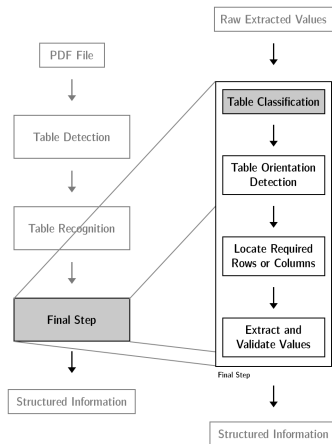
- For the required task:
  - Two-stage detectors perform better than single-stage detectors
  - FasterRCNN v2 performs the best

Questions?

# Experiments: Table Classification Evaluation

## Overview

This experiment evaluates the approach described for table classification as part of the final step.



# Experiments: Table Classification Evaluation

## Evaluation Metrics

- Precision =  $\frac{TP}{TP+FP}$   
How many guesses made by model belonged to the correct class?
- Recall =  $\frac{TP}{TP+FN}$   
Per class, how many actual tables were identified by the model?
- F1-Score =  $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

### Note

F1-Score was used to determine the best table classification approach

# Experiments: Table Classification Evaluation

## Dataset

- Dataset
  - 60 excel files
  - 215 tables
- Classes
  - Electrical Characteristics (EC)
  - Thermal Characteristics (TC)
  - Mechanical Characteristics (MC)
  - Other (O)

# Experiments: Table Classification Evaluation

## Setup

- Word vectorisers
  - Count Vectoriser [8]
  - TF-IDF [9]
- Classifiers
  - K Nearest Neighbours [10]
  - Naive Bayes [11]
- Cross compared using 5-Fold Cross Validation



# Experiments: Table Classification Evaluation

## Results

Approach	F1-Score for Classes			
	EC	TC	MC	O
K Nearest Neighbours & Count Vectoriser	0.89	0.81	0.89	0.83
K Nearest Neighbours & TF-IDF	0.72	0.87	0.77	0.74
Naive Bayes & Count Vectoriser	0.95	0.80	0.96	0.85
Naive Bayes & TF-IDF	<b>0.95</b>	<b>0.95</b>	<b>0.94</b>	<b>0.93</b>

**Table: Comparison of Table Classification Approaches** The F1-Scores for each class are averaged over the 5-fold cross validation run

# Experiments: Table Classification Evaluation

## Results

Approach	F1-Score for Classes			
	EC	TC	MC	O
K Nearest Neighbours & Count Vectoriser	0.89	0.81	0.89	0.83
K Nearest Neighbours & TF-IDF	0.72	0.87	0.77	0.74
Naive Bayes & Count Vectoriser	0.95	0.80	0.96	0.85
Naive Bayes & TF-IDF	<b>0.95</b>	<b>0.95</b>	<b>0.94</b>	<b>0.93</b>

**Table: Comparison of Table Classification Approaches** The F1-Scores for each class are averaged over the 5-fold cross validation run

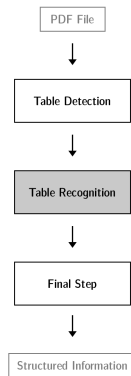
- For the required task:
  - Naive Bayes is a better classifier than K Nearest Neighbours
  - The combination of Naive Bayes with TF-IDF performs the best

Questions?

# Experiments: Complete Pipeline Evaluation

## Overview

This experiment will evaluate the complete pipeline as well as the table recognition utilities.



# Experiments: Complete Pipeline Evaluation

## Evaluation Metrics

- Precision =  $\frac{TP}{TP+FP}$   
How many of the values extracted were the actual values?
- Recall =  $\frac{TP}{TP+FN}$   
How many of the actual values were extracted?
- F1-Score =  $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

### Note

F1-Score was used to compare table recognition utilities

# Experiments: Complete Pipeline Evaluation

## Dataset

- 10 PDF documents
- Diverse selection
- Ground-truth
  - Electrical Characteristics
  - Thermal Characteristics
  - Manually extracted and saved in a YAML file

# Experiments: Complete Pipeline Evaluation

## Table Recognition Approaches Compared

- **Baseline**  
Custom utility created using raw text and positional information from the PDF.
- **Tabula [1]**  
Open-source Python utility for extracting tabular data from PDF documents.
- **Camelot [2]**  
Open-source Python utility for extracting tabular data from PDF documents.

# Experiments: Complete Pipeline Evaluation

## Setup

- 1 10 PDFs are fed into the pipeline and electrical and thermal characteristics are extracted
- 2 The extracted values are compared against ground truth values and precision, recall and f1-scores are calculated.
- 3 The process is repeated for all the table recognition utilities.



# Experiments: Complete Pipeline Evaluation

## Results

Approach	Precision	Recall	F1-Score
Baseline	0.7714	0.7066	0.7376
<b>Camelot</b>	<b>0.9986</b>	<b>0.9490</b>	<b>0.9731</b>
Tabula	0.8405	0.7769	0.8074

**Table:** Effect of different table recognition approaches on performance of complete pipeline

Questions?

## Key Takeaways

- Deep Learning based Object Detectors are well suited for Table Detection.
- Text Classification is a suitable technique for table type identification.
- Rule-based Table Recognition utilities are good enough for simple tables.
- The complete pipeline is well suited for extracting tabular information from solar cell data sheets and it would be interesting to test its feasibility on other domains as well.

Thank You!  
Questions?

- [1] Tabula, “Tabula: Extract Tables from PDFs.”  
<https://tabula.technology>, accessed on 2023-04-11.
- [2] Camelot, “Camelot: PDF Table Extraction for Humans.”  
<https://camelot-py.readthedocs.io/en/master/>,  
accessed on 2023-04-11.
- [3] B. D. D. Artha, “Data mining im bereich von photovoltaik-technologie,” 2021.
- [4] R. Ding, L. Dai, G. Li, and H. Liu, “Tdd-net: a tiny defect detection network for printed circuit boards,” *CAAI Transactions on Intelligence Technology*, vol. 4, no. 2, pp. 110–116, 2019.
- [5] PyTorch, “RetinaNet.”  
<https://pytorch.org/vision/main/models/retinanet.html>,  
accessed on 2023-04-10.

[6] PyTorch, “FasterRCNN.”

[https://pytorch.org/vision/master/models/faster\\_rcnn.html](https://pytorch.org/vision/master/models/faster_rcnn.html), accessed on 2023-04-10.

[7] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755, Springer, 2014.

[8] Sci-Kit-Learn, “CountVectorizer.”

[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.CountVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html), accessed on 2023-04-11.

[9] Sci-Kit-Learn, “TfidfVectorizer.”

[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.TfidfVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html), accessed on 2023-04-11.

[10] Sci-Kit-Learn, “Nearest Neighbors.”

<https://scikit-learn.org/stable/modules/neighbors.html>, accessed on 2023-04-11.

[11] Sci-Kit-Learn, “Naive Bayes.”

[https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html), accessed on 2023-04-11.

# Appendix: RetinaNet Training Run Graph

## Recall at 90% IoU

Batch size of 4 and learning rate of  $1e-4$  produced the best recall performance of 80.07% at epoch 47

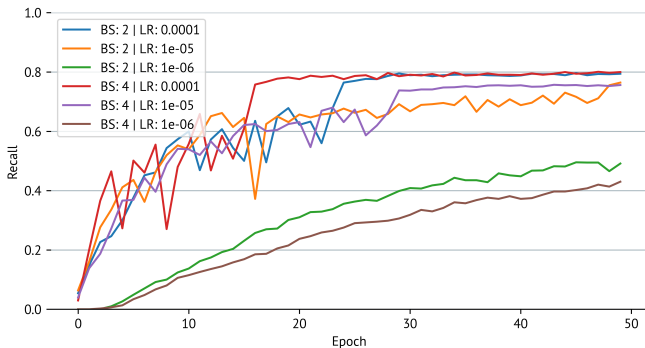


Figure: RetinaNet Hyperparameter Tuning



# Appendix: RetinaNet Training Run Results

Configuration		Epoch	Recall		Precision	
Batch Size	Learning Rate		90% IoU	75% IoU	90% IoU	75% IoU
2	1e-04	45	0.7965	0.9322	0.8099	0.9479
2	1e-05	49	0.7643	0.9388	0.7732	0.9498
2	1e-06	45	0.4955	0.8859	0.5008	0.8955
<b>4</b>	<b>1e-04</b>	<b>47</b>	<b>0.8007</b>	<b>0.9247</b>	<b>0.8155</b>	<b>0.9419</b>
4	1e-05	43	0.7568	0.9247	0.7663	0.9363
4	1e-06	49	0.4301	0.8354	0.4758	0.9241

Table: RetinaNet Results

# Appendix: RetinaNet v2 Training Run Graph

## Recall at 90% IoU

Batch size of 2 and learning rate of  $1e-4$  produced the best recall performance of 80.07% at epoch 34

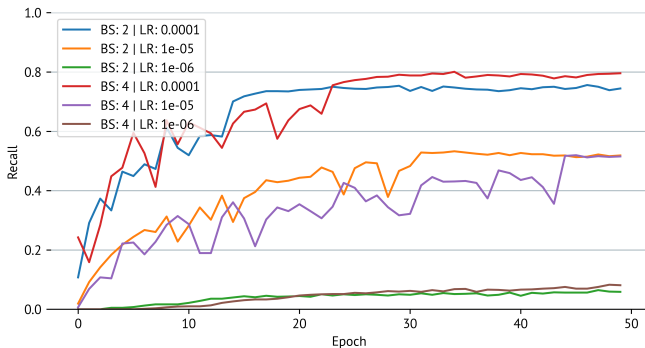


Figure: RetinaNet V2 Hyperparameter Tuning

# Appendix: RetinaNet v2 Training Run Results

Configuration		Epoch	Recall		Precision	
Batch Size	Learning Rate		90% IoU	75% IoU	90% IoU	75% IoU
2	1e-04	46	0.756	0.9156	0.7455	0.9029
2	1e-05	34	0.5327	0.8519	0.5389	0.8619
2	1e-06	47	0.0645	0.5285	0.0707	0.5788
4	<b>1e-04</b>	<b>34</b>	<b>0.8007</b>	<b>0.928</b>	<b>0.8259</b>	<b>0.9573</b>
4	1e-05	45	0.5194	0.8594	0.5182	0.8573
4	1e-06	48	0.0827	0.5476	0.104	0.6881

Table: RetinaNet V2 Results

# Appendix: FasterRCNN Training Run Graph

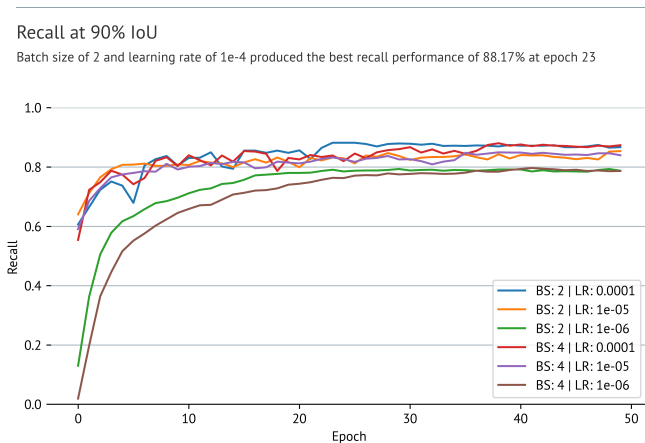


Figure: FasterRCNN Hyperparameter Tuning

# Appendix: FasterRCNN Training Run Results

Configuration		Epoch	Recall		Precision	
Batch Size	Learning Rate		90% IoU	75% IoU	90% IoU	75% IoU
2	1e-04	23	0.8817	0.952	0.8825	0.9528
2	1e-05	49	0.8536	0.9545	0.8411	0.9405
2	1e-06	29	0.7932	0.957	0.7551	0.911
4	1e-04	38	0.8801	0.9495	0.8837	0.9535
4	1e-05	38	0.8495	0.9529	0.8397	0.9419
4	1e-06	41	0.7965	0.9504	0.7631	0.9105

Table: FasterRCNN Results

# Appendix: FasterRCNN v2 Training Run Graph

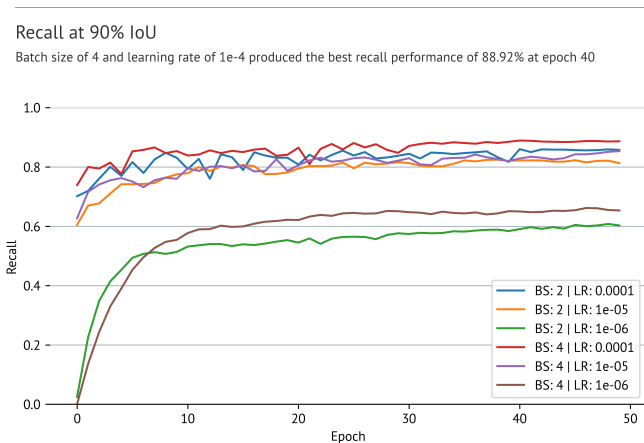


Figure: FasterRCNN V2 Hyperparameter Tuning

# Appendix: FasterRCNN v2 Training Run Results

Configuration		Epoch	Recall		Precision	
Batch Size	Learning Rate		90% IoU	75% IoU	90% IoU	75% IoU
2	1e-04	40	0.8602	0.9322	0.8784	0.9519
2	1e-05	38	0.8246	0.9313	0.8295	0.9368
2	1e-06	48	0.6079	0.8759	0.5829	0.8398
<b>4</b>	<b>1e-04</b>	<b>40</b>	<b>0.8892</b>	<b>0.9512</b>	<b>0.9018</b>	<b>0.9648</b>
4	1e-05	49	0.8536	0.9429	0.8445	0.9329
4	1e-06	46	0.6617	0.8941	0.6385	0.8627

**Table: FasterRCNN V2 Results**

## Appendix: Table Class - KNN with CV

Class	Precision	Recall	F1-Score
Electrical Characteristics	0.92	0.87	0.89
Thermal Characteristics	0.97	0.71	0.81
Mechanical Characteristics	0.95	0.87	0.89
Other	0.75	0.93	0.83
Overall Accuracy			0.86

**Table: K-Nearest Neighbours with Count Vectoriser** The table shows the averaged evaluation results of 5 Fold Cross Validation performed on the dataset.



## Appendix: Table Class - KNN with TF-IDF

Class	Precision	Recall	F1-Score
Electrical Characteristics	0.95	0.6	0.72
Thermal Characteristics	1.00	0.78	0.87
Mechanical Characteristics	1.00	0.65	0.77
Other	0.60	0.98	0.74
Overall Accuracy			0.77

**Table: K-Nearest Neighbours with TF-IDF Vectoriser** The table shows the averaged evaluation results of 5 Fold Cross Validation performed on the dataset.

## Appendix: Table Class - Naive Bayes with CV

Class	Precision	Recall	F1-Score
Electrical Characteristics	0.93	0.96	0.95
Thermal Characteristics	0.78	0.91	0.80
Mechanical Characteristics	0.93	1.0	0.96
Other	0.96	0.80	0.85
Overall Accuracy			0.89

**Table: Naive Bayes with Count Vectoriser** The table shows the averaged evaluation results of 5 Fold Cross Validation performed on the dataset.

## Appendix: Table Class - Naive Bayes with TF-IDF

Class	Precision	Recall	F1-Score
Electrical Characteristics	0.92	0.98	0.95
Thermal Characteristics	1.00	0.91	0.95
Mechanical Characteristics	0.93	0.97	0.94
Other	0.94	0.92	0.93
Overall Accuracy			0.94

**Table: Naive Bayes with TF-IDF Vectoriser** The table shows the averaged evaluation results of 5 Fold Cross Validation performed on the dataset.