

# Optical Character Recognition systems for Document Understanding

**Mikuláš Bankovič**

**456421@mail.muni.cz**

Faculty of Informatics, Masaryk University

October 10, 2022

## Research sources

- A Survey of Deep Learning Approaches for OCR and Document Understanding [5]
- ICDAR 2019 - Scanned Receipts OCR and Information Extraction (SROIE) [2]

# Computer Vision problems



Object Detection



Semantic Segmentation

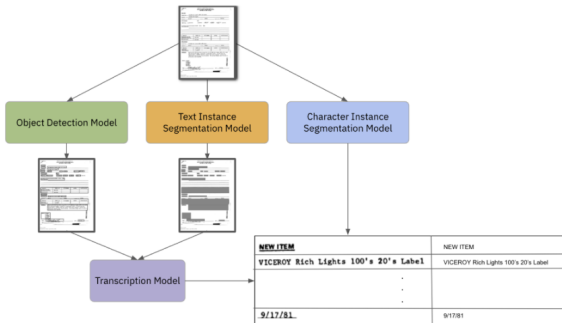


Instance Segmentation

*Source – Ref. 8*

# OCR pipeline

Figure 1 given by Subramani et al. [5] shows different approaches to OCR systems.



**Figure:** Object detection and segmentation need transcription(text recognition model)

# Text Detection

- CRAFT based on CNN, specifically FCN [1]
- Differentiable Binarization Network (DBNet) [3]

# Comparisons

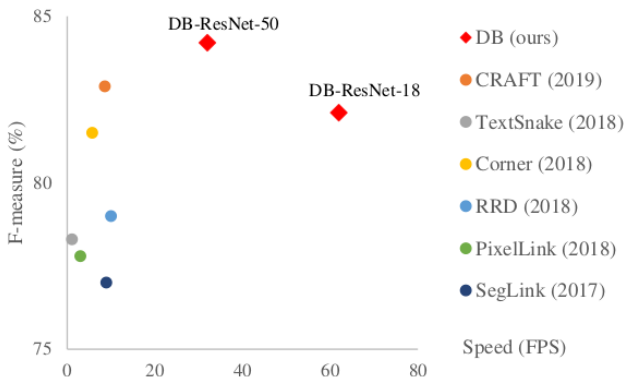


Figure: Speed and metric comparison

# Text Recognition

- CRNN with different backbones (mobilenet, vgg)
- Master - CNN + Transformer [4]

# Master

Method	Input	Accuracy	Inference Time (ms)	Training Time (h)
SAR [8]	$48 \times 160$	91.5	16.1	51
MASTER (original)	$48 \times 160$	95.0	9.2	36
MASTER (improved)	$48 \times 160$	95.0	4.3	36

Figure: Speed comparison to previous SOTA



# Master

- Connectionist Temporal Classification (CTC) loss
- does not require character-level annotations but word-level annotations
- high training parallelization compared to RNN

## OCR frameworks

OCR systems allow users to use text detection and text recognition models without having to create own pipeline and visualization tools.

- EasyOCR - many languages - only a few models
- DocTR - newer interface and better variety of models including selected ones

Both have poor documentation and are in their code infancy.



Figure: We do not want to reinvent the wheel

# DocTR

- DBNet (pretrained text detection)
- Master (we want to train for our domain)



Figure: Doctr

# EasyOCR model

```
1 import torch.nn as nn
2 from .modules import ResNet_FeatureExtractor, BidirectionalLSTM
3
4 class Model(nn.Module):
5
6     def __init__(self, input_channel, output_channel, hidden_size, num_class):
7         super(Model, self).__init__()
8         """ FeatureExtraction """
9         self.FeatureExtraction = ResNet_FeatureExtractor(input_channel, output_channel)
10        self.FeatureExtraction_output = output_channel # int(imgH/16-1) * 512
11        self.AdaptiveAvgPool = nn.AdaptiveAvgPool2d((None, 1)) # Transform final (imgH/16-1) -> 1
12
13        """ Sequence modeling"""
14        self.SequenceModeling = nn.Sequential(
15            BidirectionalLSTM(self.FeatureExtraction_output, hidden_size, hidden_size),
16            BidirectionalLSTM(hidden_size, hidden_size, hidden_size))
17        self.SequenceModeling_output = hidden_size
18
19        """ Prediction """
20        self.Prediction = nn.Linear(self.SequenceModeling_output, num_class)
21
22
23    def forward(self, input, text):
24        """ Feature extraction stage """
25        visual_feature = self.FeatureExtraction(input)
26        visual_feature = self.AdaptiveAvgPool(visual_feature.permute(0, 3, 1, 2)) # [b, c, h, w] -> [b, w, c, h]
27        visual_feature = visual_feature.squeeze(3)
28
29        """ Sequence modeling stage """
30        contextual_feature = self.SequenceModeling(visual_feature)
31
32        """ Prediction stage """
33        prediction = self.Prediction(contextual_feature.contiguous())
34
35        return prediction
```

Figure: Source code of so far unnamed model implemented by EasyOCR

Mikuláš Bankovič 456421@mail.muni.cz • Optical Character Recognition systems for Document Understanding • Octo

# Invoice

<b>CZC.CZ</b> naše online výměna elektronice	<b>FAKTURA</b> - DAN
<b>Dodavatel:</b> CZC.cz s.r.o. L. Maše s.r.o. 103, Moravská Ostrava, 70300, Ostrava IČ: 26855701, DIČ: CZ26855701 Banka: Raiffeisenbank, účet: 427293001/5500 Var. úč. 427293001/5500 Dobročinný příspěvek: 100 Kč, číslo účtu: 427293001/5500	
<b>Datum vystavení:</b>	16.8.2021
<b>Datum zdanění / plnění:</b>	16.8.2021
<b>Společnost:</b>	23.8.2021
<b>Způsob platby:</b>	Online bankovním účtem ČSOB
<b>Pracovní:</b>	Česká pošta - Do ruky (bez sob.) Rudná
<b>Vystavitel:</b>	CZC.cz
<b>Objednávka č.:</b>	42104/1723

# Reconstruction

CZC.CZ

vám

FAKTURA

IDA

Dodavate:

CZC.cz s.r.o.

1. máje 3236/103, Moravská Ostrava 70300 Ostrava

IC 25655701, DIČ: CZ25655701

Banka: Raiffeisen učet 327293001/55

Vars.: 1121455944 K.s. 0008

Obchodní rejstřík Městský úřad v Praze oddíl C, vložka 58549

Datum vystavení 16.8.2021

Datum zdanění 16.8.2021

Splatnost 23.8.2021

Způsob platby Online bankovním účtem CSOB

Doprava Česká pošta Do rukou (bezdobrá) Rudná

Vystavil CZC.cz

Objednávka číslo: 4210471721

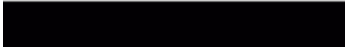
# Born-Digital Dataset

- Scrape pdf invoices from the internet (provided for us)
- Process pdf files with python libraries (pdfminer, etc) to extract bounding boxes and texts
- Upload on HuggingFace Hub

## Born-Digital Dataset

We have 699585 pairs cropped image: extracted text.

7.11.2018-8.11.2018



**Figure:** Example of born-digital cropped word. Custom model predicted:  
7.1..2018.8.1..2018



## Compare OCR on Born Digital dataset

OCR engine	Exact	Partial	No text	FPS
easyocr_generation2	0.88	0.89	0.00	56.13
easyocr_generation1	0.81	0.82	0.01	3.99
doctr_vgg16_bn*	0.78	0.78	0.00	18.93
doctr_mobilenet_v3_large*	0.77	0.78	0.00	40.70
doctr_mobilenet_v3_small*	0.77	0.77	0.00	51.36
tesseract	0.72	0.72	0.24	2.36
doctr_master	0.82	0.82	X	X

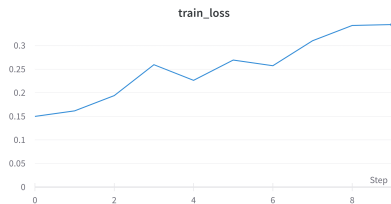
**Table:** Comparison of different recognition networks on born-digital dataset

1

---

<sup>1</sup>\*These are all CRNN architectures with different backbones

# Training and validation loss

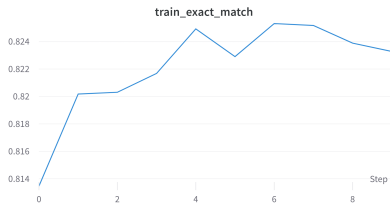


(a) Doctr custom model training loss

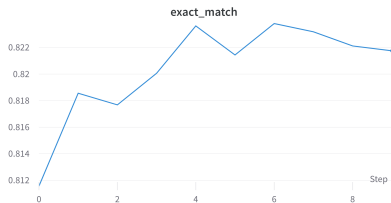


(b) Doctr custom model validation loss

# Training and validation exact match



(a) Doctr custom model training exact match



(b) Doctr custom model validation exact match

# Bibliography I

- [1] Youngmin Baek et al. *Character Region Awareness for Text Detection*. 2019. arXiv: 1904.01941 [cs.CV].
- [2] *CDAR 2019 Robust Reading Challenge on Scanned Receipts OCR and Information Extraction*. URL: <https://rrc.cvc.uab.es/?ch=13&com=introduction> (visited on 07/21/2022).
- [3] Minghui Liao et al. "Real-time Scene Text Detection with Differentiable Binarization". In: (). URL: <http://arxiv.org/abs/1911.08947> (visited on 07/21/2022).

## Bibliography II

- [4] Ning Lu et al. “MASTER: Multi-aspect non-local network for scene text recognition”. In: (). DOI: 10.1016/j.patcog.2021.107980. URL: <https://doi.org/10.1016%5C%2Fj.patcog.2021.107980>.
- [5] Nishant Subramani et al. “A Survey of Deep Learning Approaches for OCR and Document Understanding”. In: (). URL: <https://arxiv.org/abs/2011.13534> (visited on 07/21/2022).

1. We created Czech recognition dataset from invoices
2. We trained prototype custom model architecture
3. We compared pretrained solution and evaluated our prototype

**MUNI**

FACULTY

OF INFORMATICS