ORPALIS
A PSPDFKit Company

working together!

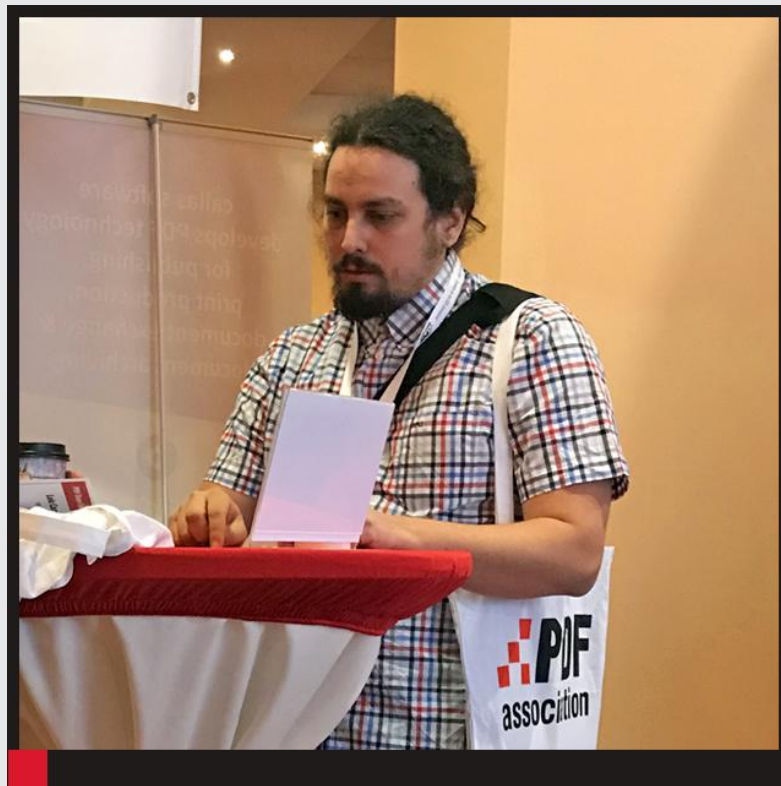PDF Days Europe 2022 | Berlin

# How document understanding can leverage your PDF workflow

How to overcome the different challenges - fields of application

# Content

- PDF, an Unstructured Format
- Challenges of PDF
- What is Document Understanding?
- Document Layout Analysis
- Optical Character Recognition

- Key-Value Pair
- Natural Language Processing
- Named-Entity Recognition
- Fields of Application

# Speakers

## Matúš Pizúr
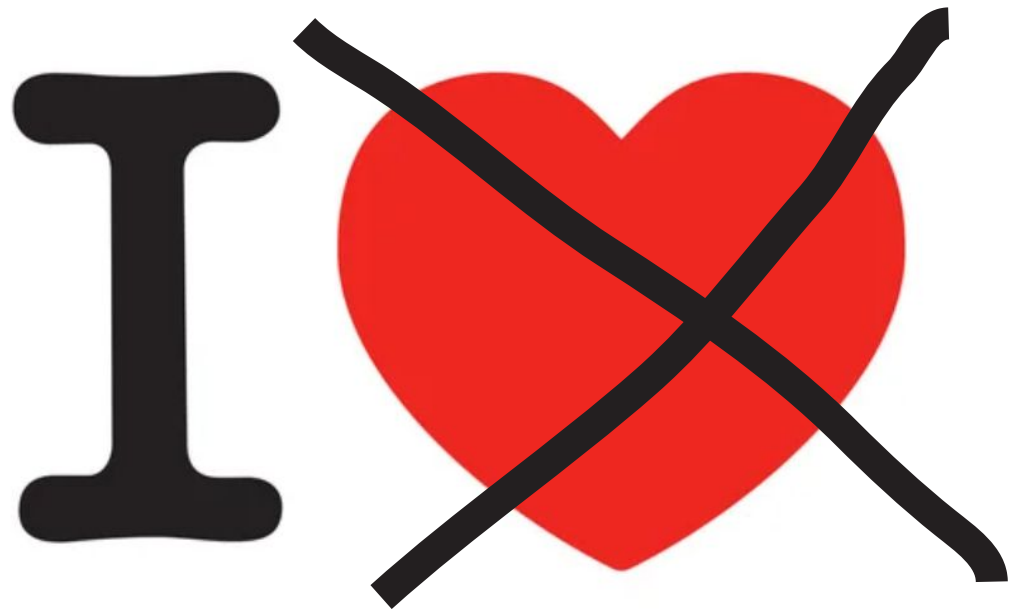
Senior developer and PDF specialist

## Elodie Tellier

Copywriter and PDF Association Board Member

I PDF

## File format for digital distribution of final form documents

Created with focus on:

- Consistency
- Fidelity

- Convenience
- Security

**Product of generation from data or conversion of other digital documents**

- Contain unstructured text and images as individual graphic objects
- Makes use of all available features of PDF specification
- No structure on the content stream level

## Basic example



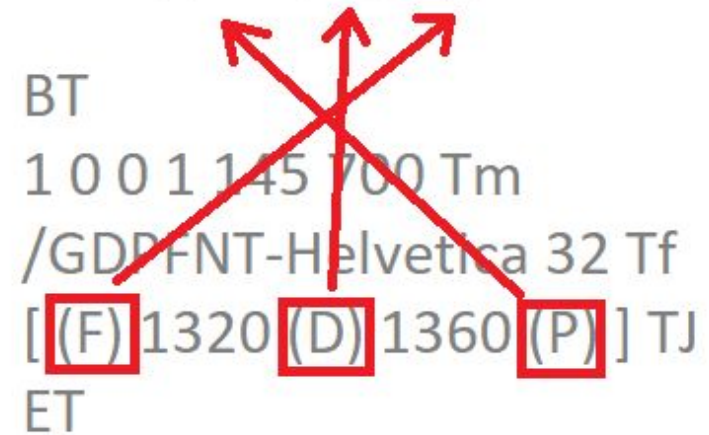```
BT
1 0 0 1 100 500 Tm
/FNT1 32 Tf
(PDF) Tj
ET
```
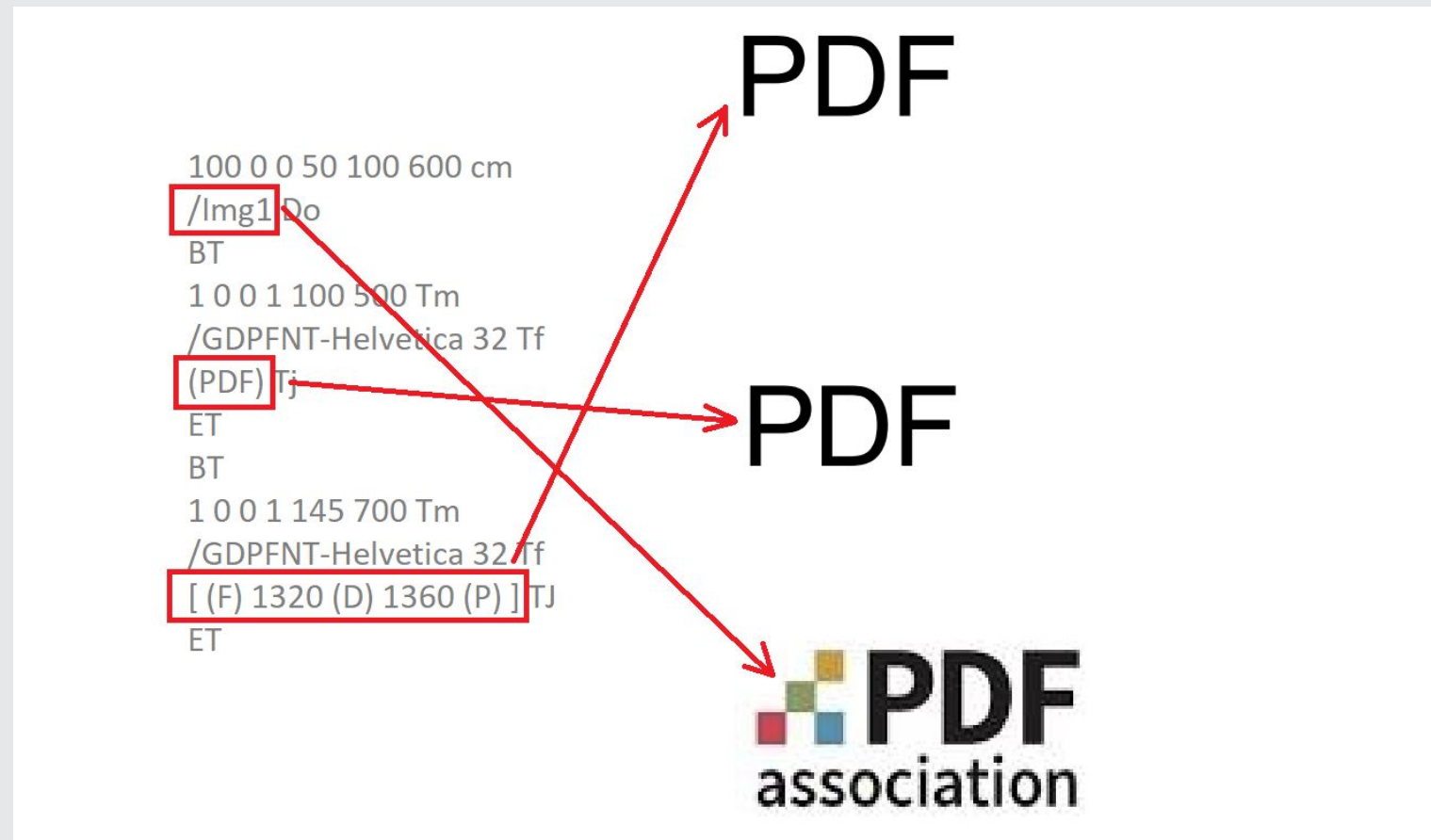
## Not so basic



PDF

```
BT
1 0 0 1 145 700 Tm
/GDPFNT-Helvetica 32 Tf
[ (F) 1320 (D) 1360 (P) ] TJ
ET
```

## Object level
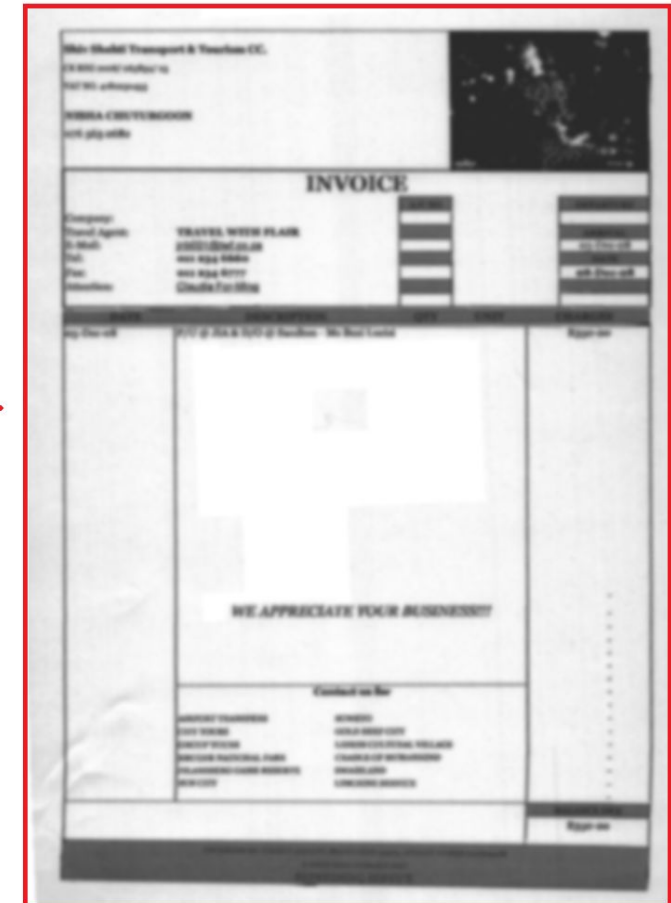
## Tagging and logical structure

- Tagged PDF introduced in PDF 1.4 (2001)
- PDF/UA - ISO 14289 (2012)
- Still only small portion of consumed files are tagged ( 14% based on 2018 data)

# Challenges of PDF: Digitized PDF

- Single image per page
- No structure
- Constitutes a considerable portion of all PDF documents consumed

395.75999 0 0 611.75999 0 0 cm
/Im0 Do

**Unstructured documents represent 90% of all documents generated**

These documents can be:

- Not searchable and discoverable
- Problematic in day to day work
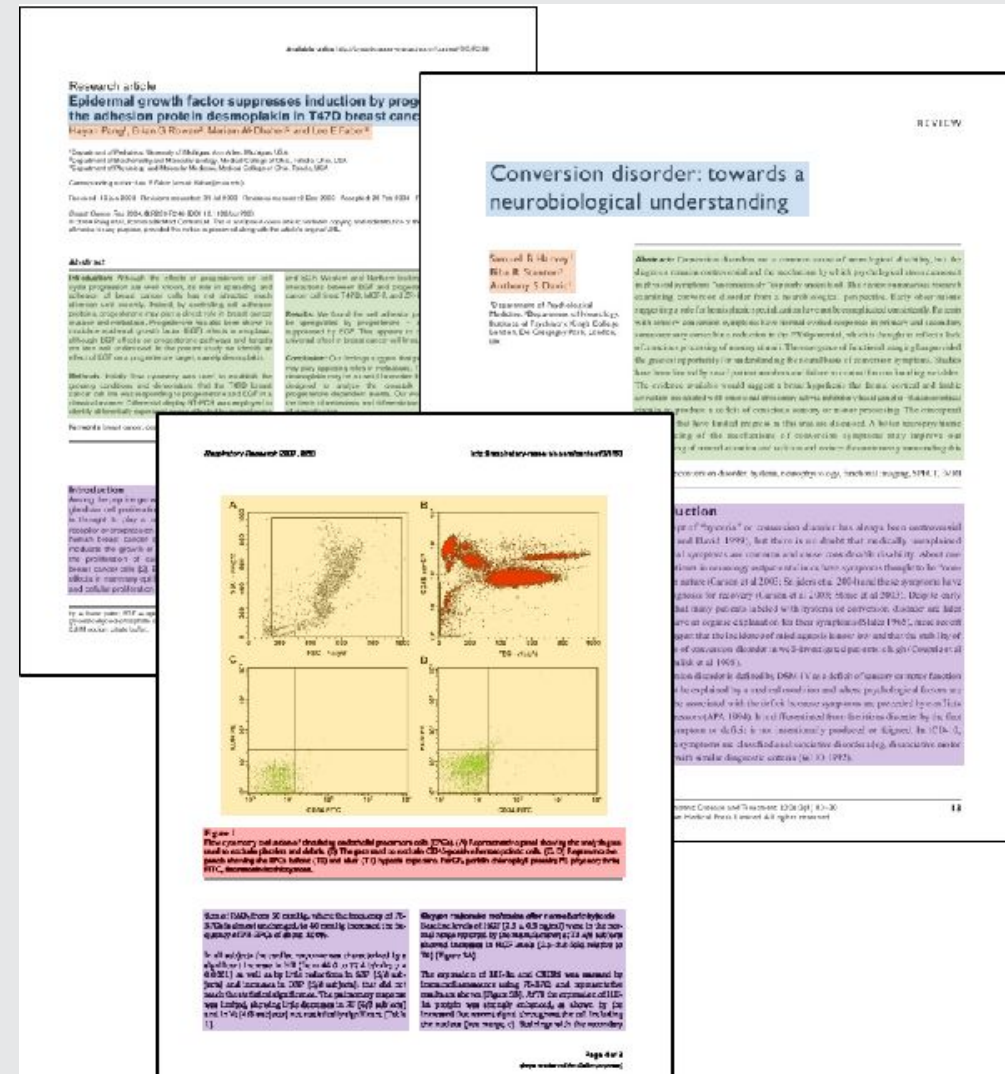- Unable to be processed in automated workflows
- Not accessible

**Document understanding is at the core of any Intelligent Document Processing solution**

- Layout analysis + recognition → First level of document understanding

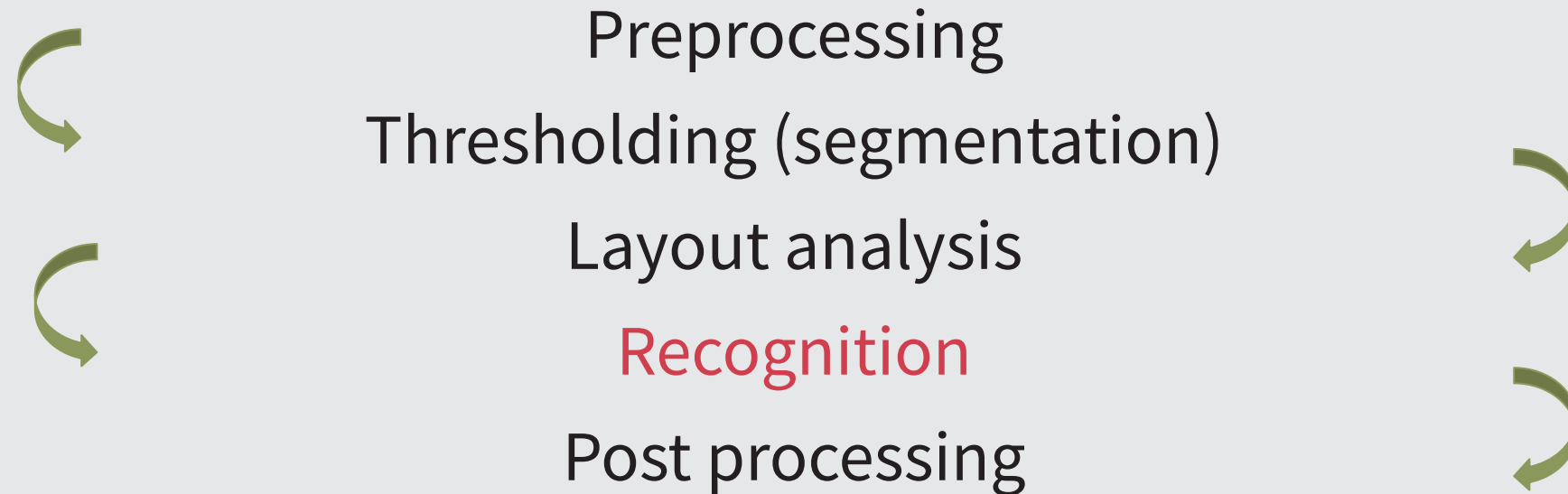- First level of document understanding + NLP → Reinforced document understanding

- DLA is the identification and categorization of regions
- DLA implies a geometric analysis of tables, pictures, equations, and barcodes and a logical layout analysis (paragraphs, lines, words, characters) of the document



Soto, C. and Shinjae Yoo. "Visual Detection with Context for Document Layout Analysis." EMNLP (2019).

## What's in an OCR engine?

Preprocessing

Thresholding (segmentation)

Layout analysis

Recognition

Post processing

## A traditional/standard OCR is not enough

- colored backgrounds
- glaring
- skew
- text in tables and graphics
- handwritten text
- same font in different size
- underlined text
- hard to scale

## Key-Value Pair (KVP)

KVPs are two related data items, a key, and a value. The key defines the data and is fixed, and the value is variable and describes the key.

Example: BIC (key) : 12345678901 (value)

## Data extraction

1.  Extracting unstructured information or text with OCR or OCR + DL

2.  Making sense of this unstructured information with DL + NLP

- NER is a form of Natural Language Processing (NLP), a subfield of artificial intelligence

- NER helps with KVP extraction: faster & more accurate

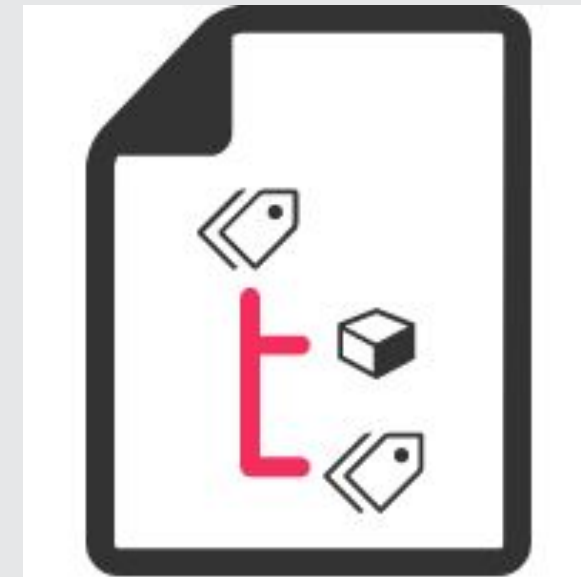- Building an algorithm with NER is not sufficient to address all the challenges

## Benefits of a performant document understanding system on OCR

- improves the OCR recognition step, especially with LSTM-based recognizers
- permits to move beyond simple OCR processes and towards document understanding

## Benefits of a performant document understanding system on tagging

- automatic tagging
- better accessibility support by making the conversion to PDF/UA easier

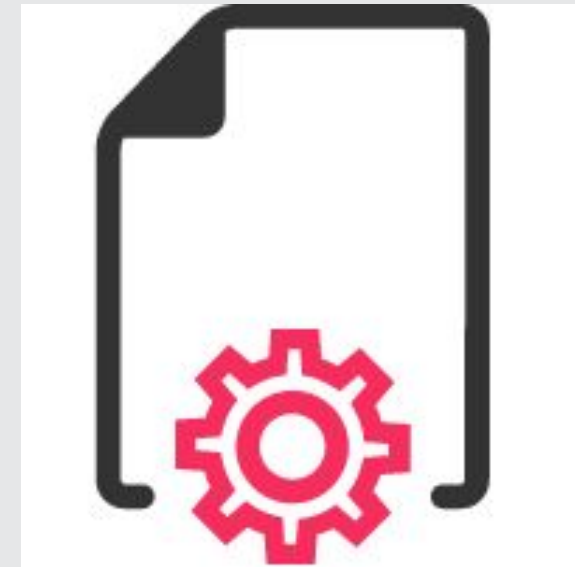## Benefits of a performant document understanding system on conversion

- structured layout conversion

- improves the conversion of a fixed into an editable layout (PDF to Office formats)

**Benefits of a performant document understanding system on automation processes**

- automatic indexing
- automatic labeling

## Benefits of a performant document understanding system on redaction

- automatic redaction

## Benefits of a performant document understanding system on compression

- better MRC compression
- better color detection compression
- better PDF optimization

**?**

# QUESTIONS

Keep in touch with us! m.pizur@orpalis.com / e.tellier@orpalis.com

www.orpalis.com

PDF association