# Document Understanding in the Age of AI
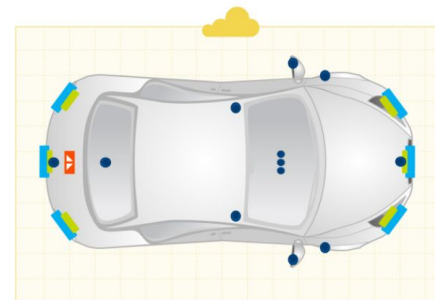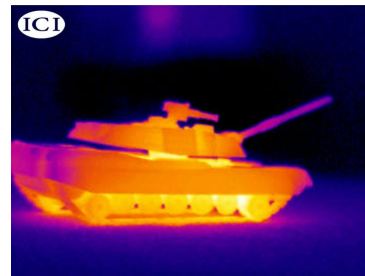
Kshitij Agrawal

28 June 2022

# About Me

Seasoned computer vision scientist

Worked across applications in

- Document understanding
- Autonomous Driving
- In game advertising
- Defence Imaging

Multiple publications in conferences

# Overview

1. Need for Document Understanding
2. Stages of Document Processing
3. Text Extraction - In detail
4. Doc classification - In detail
5. Information Extraction - In detail
6. Conclusion

# Why do we need Document Understanding

Documents are part of our everyday way of interacting and transacting with the world

- Many types of docs - Forms, texts, handwritten, certificates, legal documents

- Multiple formats - images, pdfs, raw text, books

- Varied content - tables, images, stylized text

- Challenges of lighting, shape, size and orientation

# Information in Documents



What are the **titles** of the **books**?



Who is the **invoice to**?

What is the **date**?

What is the **total**?

# Information in Documents



Which restaurant's menu?

What is the price of a burger? ice cream?



What is the content of the email?

What is the transaction amount?

Is it fraud/alert? Should the account be blocked?

# If we can extract information, we can search and automate processes

# Journey of Document Digitization

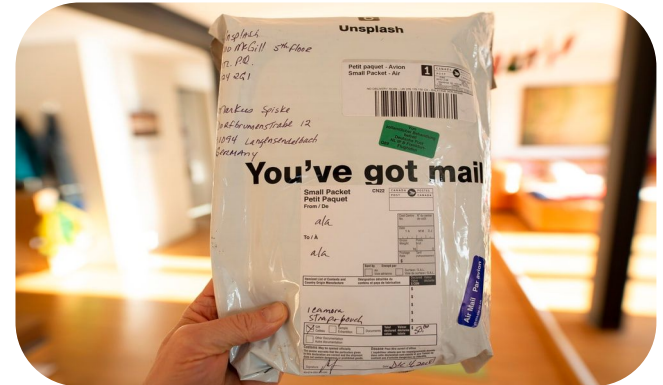| Document is input to the system | ▶ | Operator <u>views</u> and <u>enters</u> fields manually | ▶ | Another agent <u>verifies</u> and <u>approves</u> | ▶ | Digitized information stored in database for later use |
|---|---|---|---|---|---|---|

1. Time consuming
2. Error prone - if a business decision relies on it, can lead to mishaps!
3. High amount of variability -
   a. Varied formats - pdfs, jpegs, pngs etc
   b. Varied layouts - single column, multi column, tabular, forms
   c. Multiple pages

# Applications

If we can *automate* information extraction then we can automate business processes.

1. Banking, Insurance,credit card
   Customer verification, Loan document processing, email classification

2. Ecommerce and logistics
   Delivery verification, shipment processing, payment processing
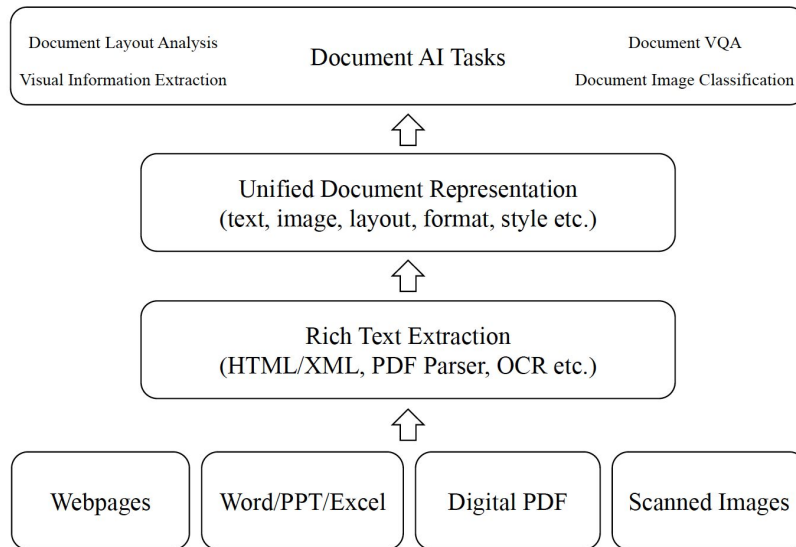   invoice processing

# Document Processing Stages

Document understanding

Visual understanding / Text
understanding / High level embeddings
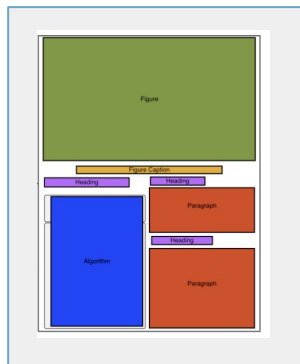
Text Extraction

Input doc

| Document Layout Analysis | Document AI Tasks | Document VQA |
| --- | --- | --- |
| Visual Information Extraction | | Document Image Classification |

↑

Unified Document Representation
(text, image, layout, format, style etc.)

↑

Rich Text Extraction
(HTML/XML, PDF Parser, OCR etc.)

↑

| Webpages | Word/PPT/Excel | Digital PDF | Scanned Images |

# Document Understanding Tasks

Information Extraction

## Image Classification



```
{
  "type": "form",
}
```

## Layout Understanding



Layout segmentation
Table Detection
Image Detection

## Form Extraction



```
{
  "company": "Brown &
Williamson",
  "date": "19/01/1982",
  "product": "viceroy"
}
```
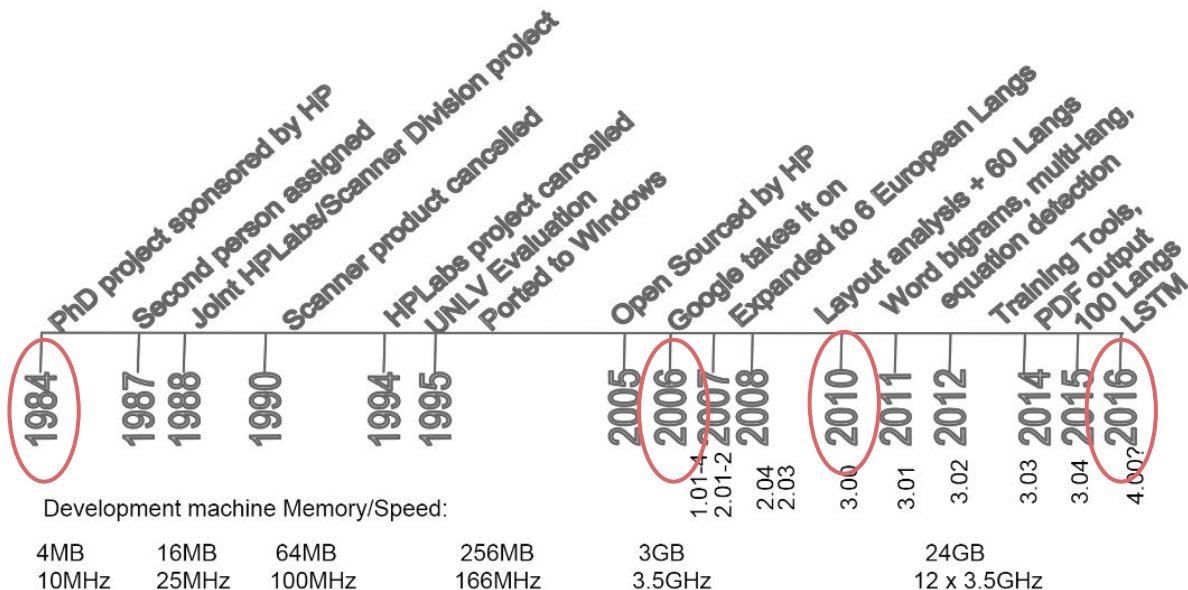
## Receipt Extraction



```
{
  "company": "Uroko
Japanese Cusine SDN BHD",
  "date": "20/03/2018",
  "total": "53.00"
}
```

## Form Understanding



Entity detection
Entity Linking
Visual Question Answering

# From Images to Text - Optical Character Recognition



Tesseract Timeline
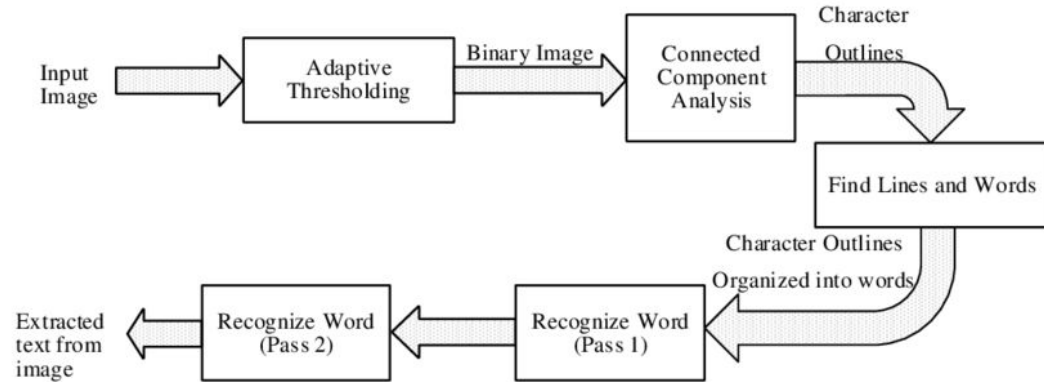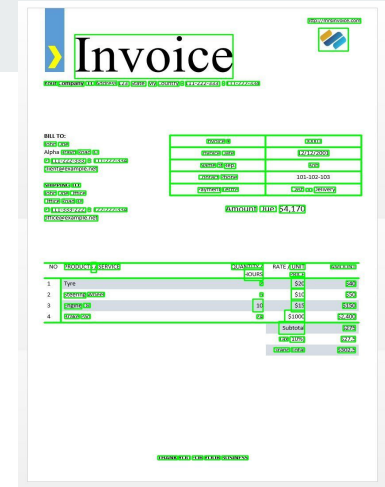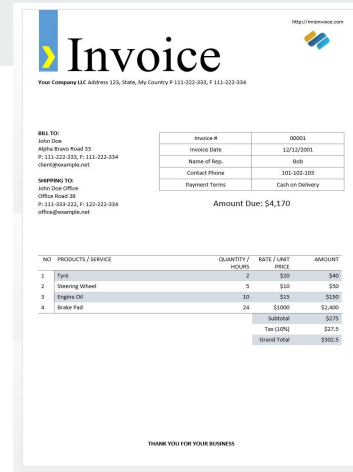
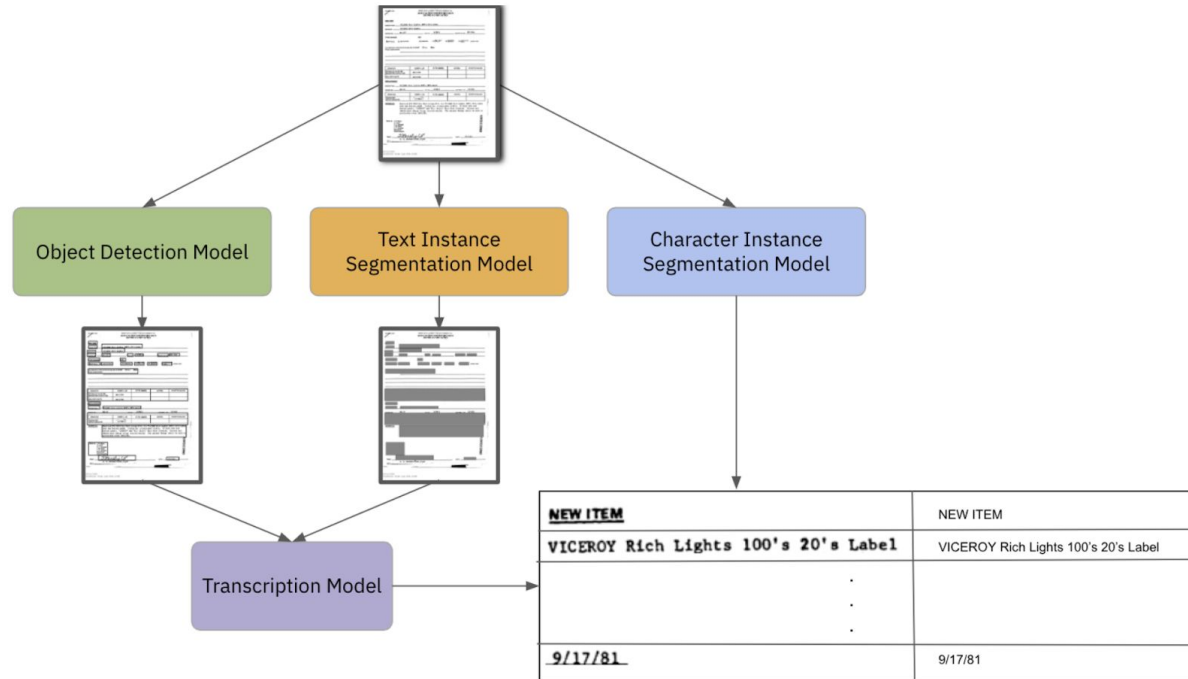# OCR - Process

**Tesseract 4.0**

- 100+ languages supported
- Full layout analysis
- Table detection
- Equation detection
- Better (Deeper) language models
- Improved segmentation search
- Word bigrams
- Training tools for custom datasets
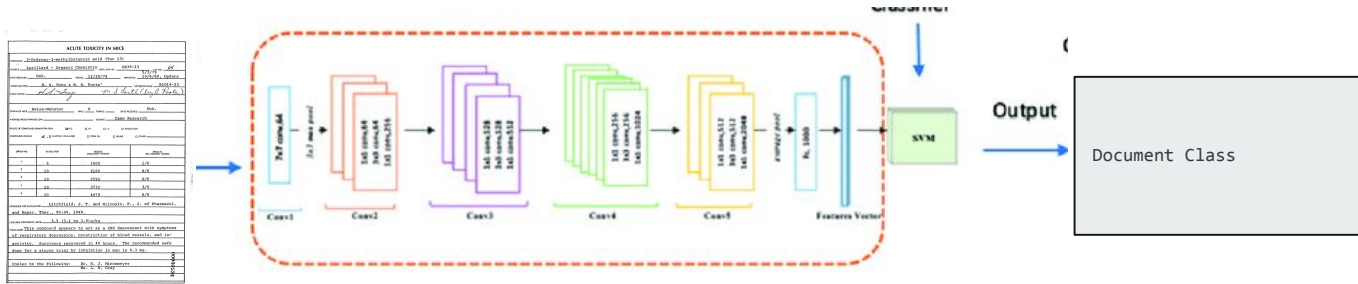
apt install tesseract-ocr



source: Architecture of Tesseract

# Optical Character Recognition - Modern architectures

1. Detection Based
   a. EAST
   b. Faster RCNN/SSD
2. Text Instance Segmentation Model
   a. Borrowed from segmentation - FCN
3. Character Instance Segmentation Model

# Document Image Classification

1. Dataset - RVLCPID
   a. 400k images
   b. 16 categories - letter, email, form, handwritten, invoice, advertisement etc.

2. Classification network - Resnet50

3. Accuracy - 90%

# Vision + Language Models

Can we leverage vision and text to improve accuracy ?

1. Vision based image features with language models - LayoutLM

2. Leverages the design of BERT models

3. Document classification Accuracy - 94.5%



"LayoutLM: Pre-training of Text and Layout for Document Image Understanding", Y. Xu et al, KDD2020

# LayoutLM

1. First time that text and layout are jointly learned in a single network for document level pre-training

2. Two types of input embeddings in BERT Model :
   a. 2D position embedding that denotes the relative position of a token within a document
   b. image embedding for scanned token text within a document

3. Architecture
   a. Multi-task learning objective -
      Masked Visual-Language Model (MVLM) loss and
      a Multi-label Document Classification (MDC) loss

Transformers ++ : Word + position embedding

Tasks

- form understanding (from 70.72 to 79.27),
- receipt understanding (from 94.02 to 95.24) and
- document image classification (from 93.07 to 94.42).

# Where to get started?

Information extraction - track finances from purchase receipts

Tesseract + LayoutLM

Result - Key value pairs of the extracted information

SROIE dataset

Information extraction from receipts

OCR text with bounding boxes

Four fields for extraction - company, date, address, total



Extracted Information

```json
{
  "company": "STARBUCKS STORE #10208",
  "address": "11302 EUCLID AVENUE, CLEVELAND, OH (216) 229-0749",
  "date": "14/03/2015",
  "total": "4.95"
}
```

Notebook: https://github.com/NielsRogge/Transformers-Tutorials/blob/master/LayoutLM/Fine_tuning_LayoutLMForSequenceClassification_on_RVL_CDIP.ipynb

# Open Research Areas in Document Understanding

Self supervised learning from large scale data

Few shot learning from scarce data

Improving image capture and preprocessing

Dealing with noise in data

# Fin.

Questions?