# The Invoice Whisperer

Liad Magen

919.5

1015

32,000,000,000

## Invoice 1 — Correct Care Solutions, LLC

**Correct Care Solutions, LLC**
1283 Murfreesboro Road
Suite 500
Nashville, TN 37217
1-800-592-2974, ext 5775

Invoice No. 081116-03

**INVOICE**

**Customer**

GRAHAM COUNTY SHERIFF'S OFFICE
ATTN: Sue Taylor
410 N POMEROY SUITE 8
Hill City    KS    67642

| Date | 081116 |
| PO No. | |
| Rep | |
| FOB | |

| | Description | Unit Price | TOTAL |
|---|---|---|---|
| 3 | Inmate healthcare services repricing 2015 -2016 | $18.00 | $54.00 |
| 0 | Inmate healthcare services repricing 2014 | $17.50 | $0.00 |

Inmate Names:
John Doe    Control No. 12345678-00
Mary Ann    Control No. 45678999-00
Sam Cook    Control No. 89765555-00

---

## Invoice 2 — KNCCI

**KNCCI** — Growing your Business together

**INVOICE**

Account Name: XYZ
Account No: 070617130546

Invoice Date: 06/07/2017
Invoice No: 10370617329516
Invoice Expiry:

| No. | Item | Item Description | Item Amount |
|---|---|---|---|
| 1. | Flowers | Qty: 100 Wgt: 1000.0 | 10,000.00 |

| Total Item Amount | EUR | 10000.0 |

**Certificate of Origin Fee Calculation**

| Block Range | 9001.0 - 18000.0 |
| Block % Fee | 0.0 |
| Block Fixed Charges | 4.4 |
| Block Category | Enterprise Member |

| Total Invoice Amount | EUR | 4.0 |

**Payment Instructions:**
Please read below information on how to deposit your payment against this invoice.

1. Mobile Money
Pay via Equitel or Mpesa using PayBill No 150090. Provide the Invoice No as the Account No. Note that only KES payments are supported for this option.

2. Bank Branch or Agent
Pay to any Equity Bank Ltd Branch or Agent using the account information provided below. Provide the Invoice No as the main reference for your transaction.

3. Online, EFT, RTGS, Cheque
Pay using any of the above channels using the account information provided below. Provide the Invoice No as main reference for your transaction.

**Account Information**
Account Name: KENYA NATIONAL CHAMBER OF COMMERCE AND INDUSTRY
Accounts:
1. KES 1510264669132
2. USD 1510265711317
3. EUR 1510265711361
4. GBP 1510265711408

Bank Name: Equity Bank (K) Limited
Branch Name: Mayfair Supreme Centre
Bank Code: 68
Branch Code: 151
SWIFT CODE: EQBLKENA

---

## Invoice 3 — NIKE

**NIKE**
GSTIN 27ABCCT2727Q1ZX
Olympia Building, 66/1
Bagmane TechPark
Bengaluru, 29-KARNATAKA, 560093
Mobile None

**TAX INVOICE**
ORIGINAL FOR RECIPIENT

Bill To:
**Rohan**
Ph: 9999999999
Jeevantara building
Sansad Maarg
Delhi, 07-DELHI, 110001

Ship to:
Jeevantara building
Sansad Maarg
Delhi, 07-DELHI, 110001

Invoice #: INV-20
Invoice Date: 06 May 2022
Place of Supply: 07-DELHI

| # | Item | HSN/SAC | Rate/Item | Qty | Amount |
|---|---|---|---|---|---|
| 1 | Nike Air Force 1 '07 | 640110 | 6,944.92 | 1 | 6,944.92 |
| 2 | Nike Dri-FIT Sport Clash Men's Training T-shirt | 61091000 | 1,424.11 | 1 | 1,424.11 |

| Taxable Amount | ₹ 8,369.02 |
| IGST 12.0% | ₹ 170.89 |
| IGST 18.0% | ₹ 1,250.08 |
| **Total** | **₹ 9,790.00** |

Total Items / Qty : 2 / 2.00
Total amount (in words): INR Nine Thousand, Seven Hundred And Ninety Rupees Only.
Amount Payable: ₹ 9,790.00

Pay using UPI:
Bank Details:
Bank: Yes Bank
Account #: 999999999999999
IFSC: YES99999
Branch: Kodihalli

For NIKE

Notes:
Thank you for the Business

Terms and Conditions:
1. Goods once sold cannot be taken back or exchanged.
2. We are not the manufacturers, company will stand for warranty as per their terms and conditions.
3. Interest @24% p.a. will be charged for uncleared bills beyond 15 days.
4. Subject to local Jurisdiction.

---

## Invoice 4 — GitLab

**GitLab** AUTHORISED RESELLER

1233 Howard Street, Suite 2F
San Francisco, California 94103
United States
about.gitlab.com
sales@gitlab.com

| Bill To: | For resale only to: |
|---|---|
| {Reseller Accounting Contact} | {End User Name} |
| {Reseller Name} | {End User Company} |
| {Reseller Address} | {End User Address} |
| Quote Number: {UID} | Quote Date: {YYYY.MM.DD} |
| | Currency: USD |

| Item | Quantity | Per item | Total |
|---|---|---|---|
| {Product Name} | {nn} | {$$} | {$$} |
| Reseller Discount | {%%} | | {$$} |
| | | Total | {$$} |

**TERMS OF QUOTATION**

1. **Valid for:** 30 days from the quote date
2. **Warranty:** 45 days money back.
3. **Payment:** Net 30 days via international wire transfer, credit card, or check (see below)

**ACCEPTANCE**

To accept this quote your customer must accept the terms and conditions of the GitLab subscription agreement found at:
https://about.gitlab.com/terms/print/gitlab_subscription_terms.pdf
. They can do this be either:
1. Order the subscription online with a credit card at http://www.gitlab.com/subscription/
2. Signing the GitLab subscription agreement and emailing a scanned copy to sales@gitlab.com.
3. Accepting the GitLab Click Through EULA at time of license key download

I hereby agree to all of the terms and conditions the price and terms specified above to the exclusion of all other terms (including, without limitation, those included on any purchase order).

I hereby represent that I have the authority to bind the organization set forth at the top of this quote.

Signature          Full name          Title          Date

| Bank information | Beneficiary information |
|---|---|
| Rabobank | GitLab BV |
| Croesselaan 18, 3500HG | 108 Ondiep, 3552 EK |
| Utrecht, The Netherlands | Utrecht, The Netherlands |
| Bank account number: 153952644 | VAT #: NL853740343B01 |
| SWIFT Code: RABONL2U | DO NOT SEND REMITTANCE TO THIS ADDRESS |
| Routing number: 121137522 | |
| IBAN: NL48RABO0153952644 | |

---

## Invoice 5 — Amazon

**TAX INVOICE**                    ORIGINAL FOR RECIPIENT

**Amazon**
GSTIN 26ADCDE3836R1ZQ
Q-city, 2nd Floor-Block A & Block B Survey Number-109,110,111/2,
Nanakramguda Village Serlingamplayy Mandal, Ranga Reddy Dist.
Hyderabad, 36-TELANGANA, 500032
Mobile 9999999999

Invoice #: **INV-13**
Invoice Date: 06 Apr 2022

Customer Details:
**Gaurav Gupta**

Billing address:
Babuganj, Hasanganj
Lucknow, 09-UTTARPRADESH, 226007

Shipping address:
Babuganj, Hasanganj
Lucknow, 09-UTTARPRADESH, 226007

Place of Supply: **09-UTTARPRADESH**

| # | Item | Rate/Item | Qty | Taxable Value | Tax Amount | Amount |
|---|---|---|---|---|---|---|
| 1 | Samsung Galaxy F23 HSN: 8517 Color - Aqua Green Storage - 128 GB Ram - 6 GB | 15,677.10 | 1 | 15,677.12 | 2,821.88 (18%) | 18,499.00 |
| 2 | Samsung 45 Watt Travel Adapter HSN: 8504 EP-TA845XBNGIN Color - Black | 2,541.53 | 1 | 2,541.52 | 457.48 (18%) | 2,999.00 |

| Taxable Amount | ₹ 18,218.63 |
| IGST 18.0% | ₹ 3,279.36 |
| **Total** | **₹ 21,498.00** |

Total Items / Qty : 2 / 2.00
Total amount (in words): INR Twenty-One Thousand, Four Hundred And Ninety-Eight Rupees Only.
Amount Payable: ₹ 21,498.00

Pay using UPI:
Bank Details:
Bank: Yes Bank
Account #: 999999999999999
IFSC: YES99999
Branch: Kodihalli

For Amazon

Notes:
Thank you for the Business

Terms and Conditions:
1. Goods once sold cannot be taken back or exchanged.
2. We are not the manufacturers, company will stand for warranty as per their terms and conditions.
3. Interest @24% p.a. will be charged for uncleared bills beyond 15 days.
4. Subject to local Jurisdiction.

---

## Invoice 6 — Instagram

**Instagram**
**TAX INVOICE**
GSTIN 36ABCCS2942R1ZR
Block 3 , Indira Nagar, Gachibowli
Hyderabad, 36-TELANGANA, 500032
Mobile 9999999999

Bill To:
**Taniya**
GSTIN: 27ABCCS2942R1ZR
Ph: 9999999999
Survey 115/1, ISB Rd, Financial District,
Gachibowli, Nanakramguda, Telangana
36-TELANGANA, 500032

Ship to:
Survey 115/1, ISB Rd, Financial District,
Gachibowli, Nanakramguda, Telangana
36-TELANGANA, 500032

Invoice #: INV-5
Invoice Date: 02 Mar 2022
Place of Supply: 36-TELANGANA

| # | Description | HSN/SAC | Amount |
|---|---|---|---|
| 1 | Digital Marketing Facebook Ads: for 1 Month - 5 Paid campaigns - 20 Ads. Instagram Ad - Organic + Reels and Story Ads | 45346 | 21,186.40 |
| 2 | Logo Designing | - | 5,000.00 |
| 3 | Professional fee | - | 1,694.92 |

| Subtotal | ₹ 27,880.32 |
| Taxable Amount | ₹ 27,881.32 |
| CGST 9.0% | ₹ 2,509.32 |
| SGST 9.0% | ₹ 2,509.32 |
| **Total** | **₹ 32,900.00** |

Total Items / Qty : 3 / 3.00
Total amount (in words): INR Thirty-Two Thousand, Nine Hundred Rupees Only.

Pay using UPI:
Bank Details:
Bank: Yes B
Account #: 999999999999999
IFSC: YES99999
Branch: Somajiguda

For Instagram

Diversity

## Receipt 1 (Credit Card / Visa Sale)

```
MID: 000000003391861          TID: 051528??
367325695883

CREDIT CARD

VISA SALE

                              XXXXXXXXXXX█
CARD #                        A0000000031010
Chip Card AID:
ATC:                                    0002
TC:             ED7724DAB5DC11E2
                                        0003
INVOICE                               000109
Batch #:                              16503D
Approval Code:
Entry Method:                      Chip Read
Mode:                                 Issuer
```

## Receipt 2 (Invoice — Menen Hotel Inc)

Menen Hotel Inc

| | Invoice Number | MISR/228/19-20 |
| | Reference | Aiwo Boat Harbour, |
| | Hire of Equipment | Aiwo District, |
| | | PO Box No. 99 |
| | | Republic Of Nauru |

| Description | Quantity | Unit Price | Amount AUD |
| --- | --- | --- | --- |
| ██████ Hire of Equipment (2 ways) | 1.00 | 450.00 | 450.00 |
| | | Subtotal | 450.00 |
| | | **TOTAL AUD** | **450.00** |

**Due Date: 9 Jun 2020**
Please Pay by Direct Credit to:
Bendigo & Adelaide Bank
Account Name: Port Authority of Nauru
Account No: 168 861 599
BSB: 633 000
Swift Code: BENDAU3B

* Please mention Invoice number as reference in your payment receipt.

* If you have any questions about this Invoice, Please contact; Mr.Naveen Kumar,
Phone: +674 557 2996, Email:naveenkannothnpa@gmail.com

------------
**PAYMENT ADVICE**                    Customer ██████
                                      Invoice Number ██████

## Receipt 3 (Costco)

```
COSTCO
WHOLESALE

MISSION VALLEY 488
2345 FENTON PKWY
SAN DIEGO, CA  92108
LW-G ET 90-102181
MEMBER #111783833925

2751   ACTIVE YEAST        3.95
261105 KS PEPPER GR        3.99
261104 KS MED SALT         3.49
32911  KS VANILLA          9.99
15306  MRM ANTIPSTO        4.97
227731 GREEK YOGURT        6.99
30669  BANANAS             1.32
135283 OMEGA-3 EGGS        5.49
285471 STRWBRY WFFL       10.99
330482 NEPRICA             7.99  A
128912 ORGANIC BUTR        6.89
188227 BLUE CHEESE         7.29
15446  OLIVES              6.49
18267  FIRM TOFU           2.15
83505  SM  RED POTA        5.79
47019  AVOCADOS            5.29
23003  MOZZ/OG 2/1#        7.49

       SUBTOTAL           96.56
A  7.75 % TAX RATE          .62
VF     TOTAL              ██████
       American Express   97.18

XXXXXXXXXXXX              SWIPED
                         589551
American Express
```

## Receipt 4 (Netto — German)

```
              0,45 B
SB S-Medaill.natur300gQS   2,99 B
RABATT        30,00%      -0,90
SB Schinkenspeck 300gQS    2,29 B
SB MG Prosciutto Par.70g   1,69 B
SUMME [6]                 11,50
EC-Karte EUR              11,50

-K-U-N-D-E-N-B-E-L-E-G-

Terminal-ID:  65434052
TA-Nr 031864     BWr 9895
Kartenzahlung
Online Lastschrift
EUR 11,50
           ###.###2705
Karte 1 gültig bis 12/15
Datum 05.11.13 18:30 Uhr
*** Zahlung erfolgt ***

           MWSt   BRUTTO   NETTO
B  7,0%    0,75   11,50    10,75

Bei Bezahlung mit cardNmore erhalten Sie
           12 Punkt(e)
Informationen unter www.netto-online.de
*************************************
      *Stars.Style.Sparen*
Ab jetzt bei Netto:gold-Des Star-Magazin
```

## Receipt 5 (Ministry of Finance — Bhutan)

```
དཔལ་ལྡན་གཞུང་གི་གཞུང་ཚབ།

MINISTRY OF FINANCE
DEPARTMENT OF REVENUE & CUSTOMS
དངུལ་རྩིས་ལས་ཁུངས།
Revenue Money Receipt
```

Serial No:- **1596603**

| | | System Receipt | RC449555 |
| Receipt Date : 23-Nov-2015 | | Security No | A429527D5ED496 |
| Generated On : 23-Nov-2015 | | Page No | |
| Book No : | | | |
| Pre Print No : 1596603 | | | |

Received from Regional Director, Regional Trade & Industry Office (RTIO), Thimphu ( RAY00605 ) a sum of BTN 1600.00 ( ONE THOUSA
HUNDRED NGULTRUM ONLY ) by CASH towards the following accounts.

| Sl. No. | A/C Head Code | Particulars of head of Accounts | Amounts |
| --- | --- | --- | --- |
| 1 | 113423905 | Service License Registration Fees - Cottage | |
| 2 | 113423906 | Service License Renewal Fees - Cottage | |
| 3 | 131130012 | License Booklet Fee | |
| | | Total | BTN |

Received By: b.k.sham, Cash Officer of Regional Director, Regional Trade & Industry Office (RTIO), Thimphu
Name - Designation & Office Seal                                    (Signature)

Validity of the receipts is subject to realisation of the amount.

## Receipt 6 (Tesco — UK)

```
COOKING █████           0.31
SOUP                    0.42
PEAR QUARTRS            0.31
SOUP                    2.15
CAT FOOD           *    0.38
ORANGE JUICE       *    2.34
OPTIONS DRINK           0.69
WHEAT BISCUITS          0.79
TESCO CRISPS       *    0.21
INSTANT NOODLE          0.76
CAKE SLICES
T W/SOME ROLLS          0.35
REDUCED PRICE
EGGS                    0.84
KM SQUARE MED           0.83

SUB-TOTAL              35.42
---------------------------
MULTIBUY SAVINGS
FROZEN CHICKEN 2 FOR £3 -0.48
CHICKEN TONIGHT BOGOF  -1.16
CANTALOUPE MELON BOGOF -1.47
---------------------------
TOTAL SAVINGS          -3.11
```

# Document Processing Pipeline

# Traditional Document Processing Pipeline

| Document | OCR | Information Extraction | Response |
|:---:|:---:|:---:|:---:|

➢Image extraction

➢PDF Text extraction

➢Image Pre-process

➢Text Post-process

➢NER Models

➢Results Post Process

Format (JSON)

# Traditional Document Processing Pipeline

**Document**
- Image extraction
- PDF Text extraction

**OCR**
- Image Pre-process
- Text Post-process

**Information Extraction**
- NER Models
- Results Post Process

**Response**
- Format (JSON)

OCR

Often requires preprocessing

# OCR: Preprocessing Steps



Normalization

Skew Correction

Image Scaling

Noise Removal

Morphological Operations: Thinning, Erosion and Dilation

Gray scaling

Thresholding / Binarization

# OCR: Preprocessing Steps

JN 1891
CAP 8 doc (146)

28 de Agosto

Meu caro Barão,

D'aqui lhe mando as ultimas
despedidas cheias de toda philo-
sophia que semelhante momento
ainda uma vez mais me suggere
sobre as viagens, a vida, a amisa-
de, eo nosso paiz. Espero tornar a
vel-o, mas como isto só pode ser
d'este lado receio muito que se
passe d'esta vez bastante tempo
eo parenthesis seja o mais longo
que tem havido em nossa
velha convivencia. O imprevisto
porem representa um papel tão

# Case Study: Thresholding

- Global:
    - Two-peaks, Otsu, Entropy, …
- Global Multi-Threshold
    - Liao, Kapur, …
- Local:
    - Local-Mean, Gaussian, Sauvola, Wolf, …

OCR Engine: **Tesseract**

Metric: **Character Error Rate (CER)**

$$CER = \frac{i_c + s_c + d_c}{n_c}$$

Best CER Results Per Thresholding Method

- Kapur Multithreshold + LAB Contrast Enhacement, 8.75%
- Sauvola Threshold window_size=15, 8.55%
- Nick Threshold, 8.46%
- NiBlack Threshold, 7.22%
- Kapur Threshold, 10.52%
- Kapur Multithreshold, 17.10%
- LAB Contrast Enhancement, 20.72%
- Min-Err Threshold, 18.69%

How to Choose?

Kapur Multithreshold + LAB Contrast Enhacement, 8.75%

Sauvola Threshold window_size=15, 8.55%

Kapur Threshold, 10.52%

Nick Threshold, 8.46%

NiBlack Threshold, 7.22%

Kapur Multithreshold, 17.10%

LAB Contrast Enhancement, 20.72%

Min-Err Threshold, 18.69%

# How to Choose?

"Most of the submitted algorithms employed machine learning techniques and performed best on the most complex images.

==Traditional algorithms== provided ==very good results== at a ==fraction of the time=="

ICDAR 2019 Time-Quality Binarization Competition | IEEE Conference Publication | IEEE Xplore

# How to Choose?



J. Imaging | Free Full-Text | Using Paper Texture for Choosing a Suitable Algorithm for Scanned Document Image Binarization (mdpi.com)

# How to Choose?

Still an open question.

# OCR Post-Processing



Figure 1. sOCRates and the pipeline for OCR extraction and correction.

sOCRates - a post-OCR text correction method | ScienceGate

# OCR: Post-processing

## Error correction

- Lexical - 40% non-words
- Context-dependent (Language Models)
- Seq2Seq - Translation

Survey of Post-OCR Processing Approaches | ACM Computing Surveys

# Traditional Document Processing Pipeline

| Document | OCR | Information Extraction | Response |
|---|---|---|---|

➢Image extraction

➢PDF Text extraction

➢Image Pre-process

➢Text Post-process

➢NER Models

➢Results Post Process

Format (JSON)

# Traditional Document Processing Pipeline

# Named Entity Recognition (NER)

# Named Entity Recognition (NER)

- Input: List of tokens (+ Position/Bounding Boxes)
- Output: List of Labels:
  - Invoice Date
  - Brut amount
  - Net amount
  - Tax percentage
  - Tax value
  - Recipient
  - Sender
  …

# Named Entity Recognition (NER) - BIO

**BIO Format:**

- BILLA          B-Issuer
- Sagt           O
- Mein           O
- Hausverstand   O
- 1060           B-Zip
- Wien           B-City
- Webgasse       B-Address
- 37             I-Address

# NER Post-Processing

Filter and join predicted labels

**What is the amount on this receipt?**



The amount on the receipt in the image you sent is 10.80 euros.

The line that says MONTANT: is the amount of the purchase, and it is followed by the number 10,80 EUR.

The receipt is also in French, and the line at the bottom that says TICKET CLIENT A CONSERVER means "Cust

I hope this is helpful!

# Large Language Models

**Pros:**

- Out of the box
- Can perform OCR
- Can analyze PDFs
- Multilingual support

**Cons:**

- Prompt-dependent
- Requires examples for more than basics
- Security risk:
  - Attacks
  - Prompt injections
- No on-premises option
- Reproducibility issues
  - Model Updates
  - Lack of quality control
- Liability issues

# Trustworthy AI

- [leiwand.ai - fair and trustworthy AI](#)
- Consulting & Training

# Position-based NER/IE Models: **GNNs**



An Invoice Reading System Using a Graph Convolutional Network | SpringerLink

# Position-based NER/IE Models: **LayoutLM**



**Figure 2: An example of LayoutLM, where 2-D layout and image embeddings are integrated into the original BERT architecture. The LayoutLM embeddings and image embeddings from Faster R-CNN work together for downstream tasks.**

[1912.13318] LayoutLM: Pre-training of Text and Layout for Document Image Understanding (arxiv.org)

# Position-based NER/IE Models: **LayoutLMv2**



Figure 1: An illustration of the model architecture and pre-training strategies for LayoutLMv2

[2012.14740] LayoutLMv2: Multi-modal Pre-training for Visually-Rich Document Understanding (arxiv.org)

# Position-based NER/IE Models: **LayoutLMv3**



Figure 3: The architecture and pre-training objectives of LayoutLMv3. LayoutLMv3 is a pre-trained multimodal Transformer for Document AI with unified text and image masking objectives. Given an input document image and its corresponding text and layout position information, the model takes the linear projection of patches and word tokens as inputs and encodes them into contextualized vector representations. LayoutLMv3 is pre-trained with discrete token reconstructive objectives of Masked Language Modeling (MLM) and Masked Image Modeling (MIM). Additionally, LayoutLMv3 is pre-trained with a Word-Patch Alignment (WPA) objective to learn cross-modal alignment by predicting whether the corresponding image patch of a text word is masked. "Seg" denotes segment-level positions. "[CLS]", "[MASK]", "[SEP]" and "[SPE]" are special tokens.

[2204.08387] LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking (arxiv.org)

# Position-based NER/IE Models: BROS*

Figure 2: An overview of BROS. The tokens in the document image are masked through token- and area-masking strategy. The position difference between text blocks is encoded directly to the attention mechanism in Transformer. The output token representations are used in both pre-training and fine-tuning.

[2108.04539] BROS: A Pre-trained Language Model Focusing on Text and Layout for Better Key Information Extraction from Documents (arxiv.org)

# Position-based NER/IE Models: BROS*

**SPADE Decoder**



(a) Initial token classification  (b) Subsequent token classification  (c) Entity linking (EL) task

[2108.04539] BROS: A Pre-trained Language Model Focusing on Text and Layout for Better Key Information Extraction from Documents (arxiv.org)

# Position-based NER/IE Models: **GeoLayoutLM**



Figure 2. An overview of GeoLayoutLM.

[2304.10759] GeoLayoutLM: Geometric Pre-training for Visual Information Extraction (arxiv.org)

# Position-based NER/IE Models: **Donut**



**Fig. 3. The pipeline of Donut.** The encoder maps a given document image into embeddings. With the encoded embeddings, the decoder generates a sequence of tokens that can be converted into a target type of information in a structured form

# Position-based NER/IE Models: Donut

| | Fine-tuning set | OCR | #Params$^\dagger$ | Time (ms) | ANLS test set | ANLS* handwritten |
|---|---|---|---|---|---|---|
| BERT [64] | train set | ✓ | $110M + \alpha^\ddagger$ | 1517 | 63.5 | n/a |
| LayoutLM[65] | train set | ✓ | $113M + \alpha^\ddagger$ | 1519 | 69.8 | n/a |
| LayoutLMv2[64] | train set | ✓ | $200M + \alpha^\ddagger$ | 1610 | 78.1 | n/a |
| Donut | train set | | 176M | **782** | 67.5 | **72.1** |
| LayoutLMv2-Large-QG[64] | train + dev + QG | ✓ | $390M + \alpha^\ddagger$ | 1698 | **86.7** | 67.3 |

- Pretrain + Finetune
  - SWIN Transformer

- Handwriting

- Synthetic training data

Q: What is the phone number given?
Answer: 336-723-6100
Donut: 336-723-6100
LayoutLMv2-Large-QG: **336-723- 4100**

Q: What is the name of the passenger?
Answer: DR. William J. Darby
Donut: DR. William J. Darby
LayoutLMv2-Large-QG: **DR. William J. Jarry**

Q: What is the Publication No.?
Answer: 540
Donut: **943** (another number in the image is extracted)
LayoutLMv2-Large-QG: 540

# Position-based NER/IE Models: **Nougat**



Figure 1: Our simple end-to-end architecture followin Donut [28]. The Swin Transformer encoder takes a document image and converts it into latent embeddings, which are subsequently converted to a sequence of tokens in a auto-regressive manner

Nougat (facebookresearch.github.io)

# Position-based NER/IE Models: **Nougat**



Figure 2: List of the different image augmentation methods used during training on an example snippet form a sample document.

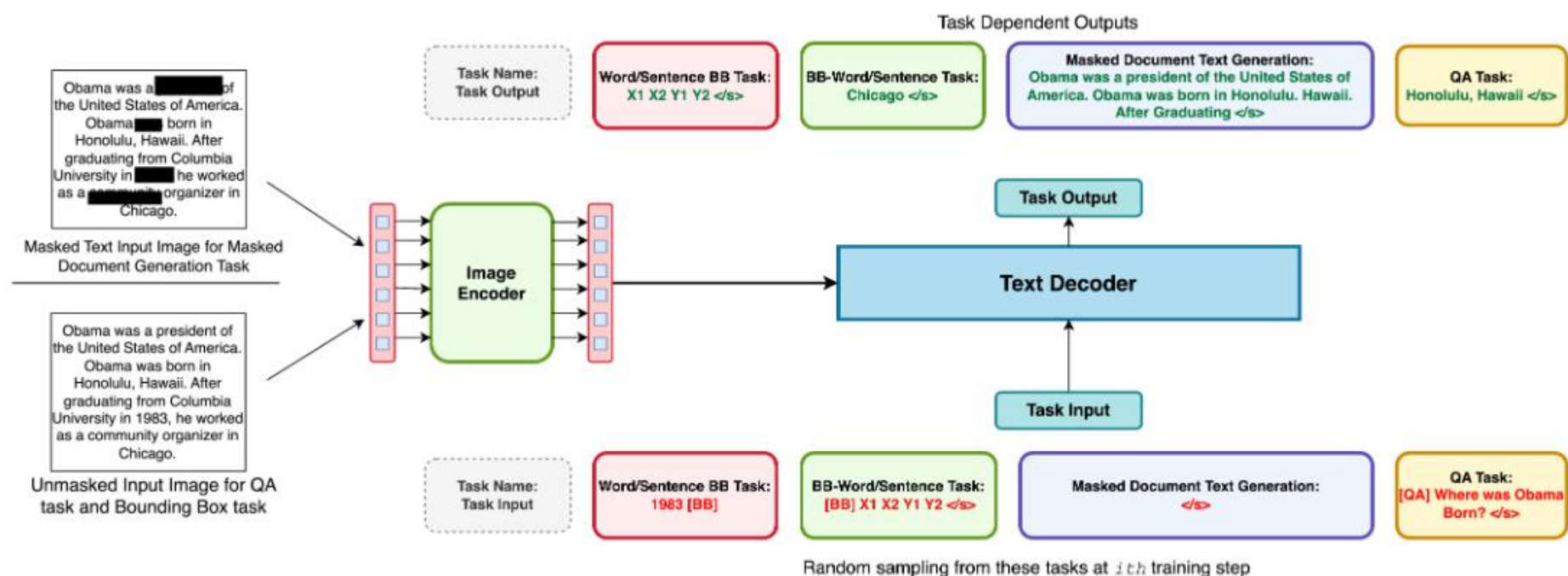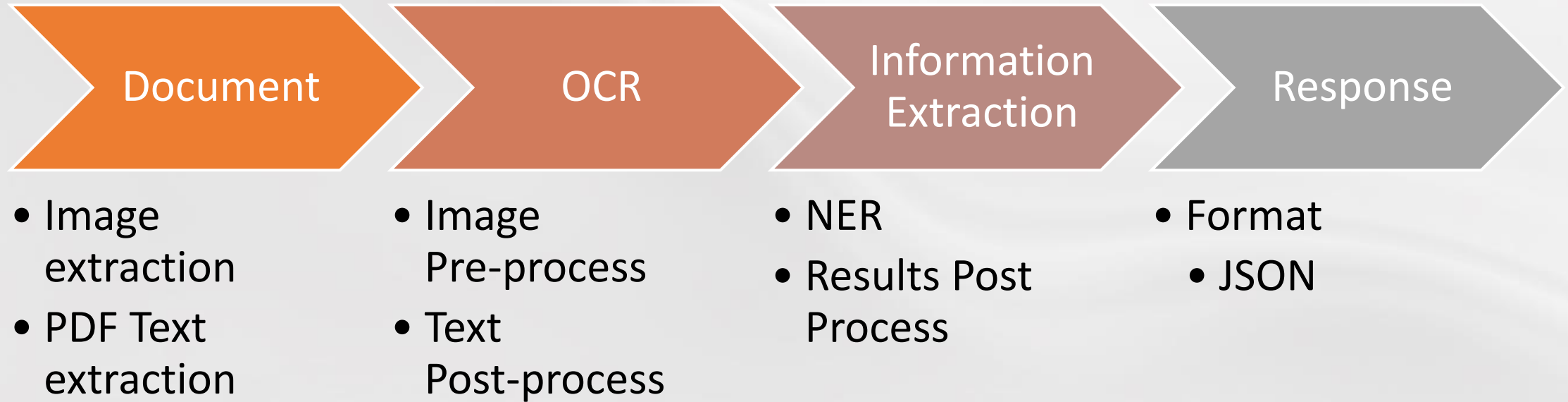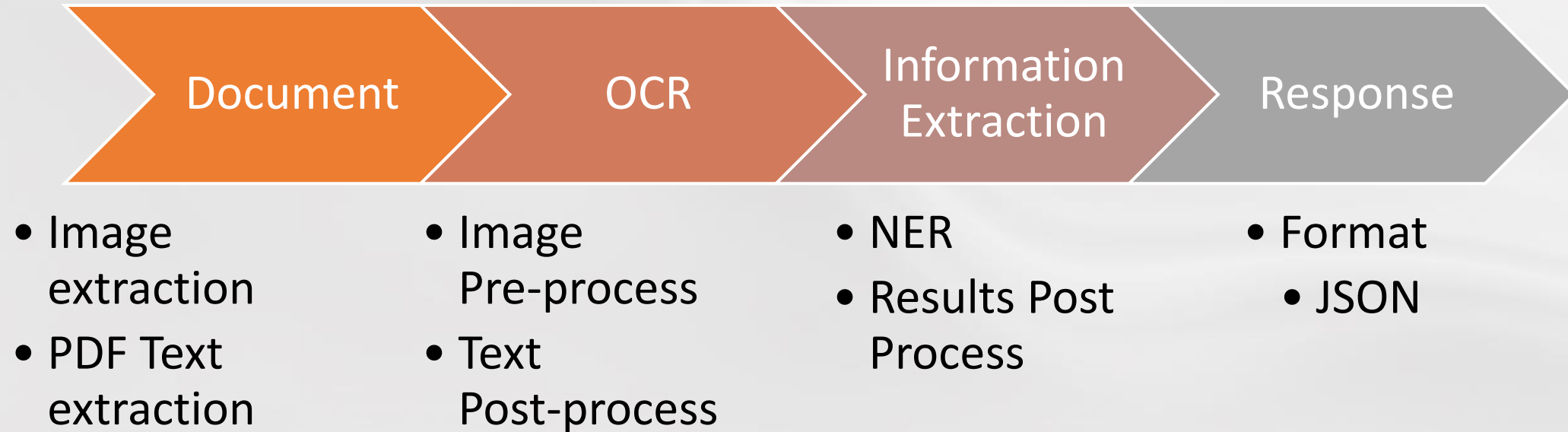# Position-based NER/IE Models: DUBLIN



Figure 2: Illustration of three tasks in the DUBLIN pretraining framework: Bounding Box, Rendered QA, and Masked Document Text Generation.

[2305.14218] DUBLIN -- Document Understanding By Language-Image Network (arxiv.org)

# Traditional Document Processing Pipeline

**Document**
- Image extraction
- PDF Text extraction

**OCR**
- Image Pre-process
- Text Post-process

**Information Extraction**
- NER
- Results Post Process

**Response**
- Format
- JSON

# Modern Document Processing Pipeline



**Document**
- Image extraction
- PDF Text extraction

**OCR**
- Image Pre-process
- Text Post-process

**Information Extraction**
- NER
- Results Post Process

**Response**
- Format
- JSON

Take-Aways

Diversity and Collaboration

Doubt fuels science

No one (method) is perfect

Thank you for your attention!

# References

- The EU invoice volume and the scalability of our blockchain-based invoice reporting system · summitto blog

- sOCRates - a post-OCR text correction method | Anais do Simpósio Brasileiro de Banco de Dados (SBBD) (sbc.org.br)

- [2108.04539] BROS: A Pre-trained Language Model Focusing on Text and Layout for Better Key Information Extraction from Documents (arxiv.org)

- [2111.15664] OCR-free Document Understanding Transformer (arxiv.org)

- [2204.08387] LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking (arxiv.org)

- [2308.13418] Nougat: Neural Optical Understanding for Academic Documents (arxiv.org)

- [2203.16618] End-to-end Document Recognition and Understanding with Dessurt (arxiv.org)

- An Invoice Reading System Using a Graph Convolutional Network | SpringerLink

# Diversity

Diversity can make your teams, your data, and your models – **stronger**!

International Conference on Document Analysis and Recognition (ICDAR)