

OCR-free Document Understanding Transformer

- This paper introduces a model called **Donut** made by NAVER CLOVA AI Lab in 2022.
- I chose it because it solves document reading without using OCR.
- It's useful for reading receipts, invoices, and forms.
- OCR makes mistakes when the text is small, blurry, or handwritten.
- Donut removes OCR and reads the document directly from the image.
- Donut is the first model that reads documents **without OCR**.
- It can take an image and produce structured text or JSON directly.
- The team created a dataset called **SynthDoG** for training.
- The model is open-source and works faster than older systems.

Research: Shifts NLP for documents from OCR-dependency to **vision-language integration**.

Industry:

Automates invoice, receipt, and ID card processing in finance, healthcare, and government.

Reduces infrastructure costs (no separate OCR pipeline).

Society: Improves accessibility (visually impaired support via text-to-speech), digitization of public records

- Donut has two main parts: an **encoder** that understands the image and a **decoder** that writes text.
- The encoder uses a Swin Transformer to find details in the image.
- The decoder, similar to BART, generates the final text or output.
- It learns by predicting the next word in a sequence during training.

Donut (OCR-Free Transformer for Documents)

Architecture:

Vision Encoder: **Swin Transformer** for image feature extraction.

Text Decoder: **BART Transformer** for autoregressive text generation.

Innovations:

Eliminates OCR engines (like Tesseract).

End-to-end learning → converts input document images directly into structured JSON/text.

Advantage: Robust against noisy scans, handwriting, and multi-lingual text, while being deployable on edge devices.

- The model takes a full document image as input.
- The encoder turns image parts into small visual features.

- The decoder creates a structured answer in text or JSON.
- It has about **220 million parameters** and supports multiple languages.
- Donut was tested on many datasets like **CORD**, **DocVQA**, and **RVL-CDIP**.
- It was compared with models like **LayoutLMv2** and **BERT**.
- The tests used accuracy and F1 scores to measure performance.
- Training used large GPUs and included multi-language documents.

1. Vision Encoder (Swin Transformer)

Learns visual features directly from the document image.

Parameters depend on the variant:

Swin Base encoder → ~88 million parameters.

Extracts patch embeddings and builds hierarchical feature maps.



```
*** RESULT ***
What does the document say? 82,000</s_price><sep></s_nm> ICE BLACKCOFFE</s_nm><s_cnt> 2</s_cnt><s_price> 61,000</s_price><sep></s_nm> AVOCAD COFFEE</s_nm><s_cnt> 1</s_cnt><s_price> 61,000</s_price><sep></s_nm> Ome CHIKEN KATSU FF</s_nm><s_cnt> 1</s_cnt><s_price> 51,000</s_price></s_menu><s_sub_total><s_subtotal_price> 194,000</s_subtotal_price><s_discount_price> 19,400</s_discount_price></s_sub_total><s_total><s_total_price> 174,600</s_total_price><s_cashprice> 200,000</s_cashprice><s_changeprice> 25,400</s_changeprice></s_total>
C:\Users\draj0.CHANDRIKAV\Downloads>python run.py
Loading model: naver-clova-ix/donut-base-finetuned-cord-v2
The following generation flags are not valid and may be ignored: ['early_stopping']. Set 'TRANSFORMERS_VERTOSITY=info' for more details.

*** RESULT ***
<s_menu></s_nm> ICE BLACKCOFFE</s_nm><s_cnt> 2</s_cnt><s_price> 82,000</s_price><sep></s_nm> AVOCAD COFFEE</s_nm><s_cnt> 1</s_cnt><s_price> 61,000</s_price><sep></s_nm> Ome CHIKEN KATSU FF</s_nm><s_cnt> 1</s_cnt><s_price> 51,000</s_price></s_menu><s_sub_total><s_subtotal_price> 194,000</s_subtotal_price><s_discount_price> 19,400</s_discount_price></s_sub_total><s_total><s_total_price> 174,600</s_total_price><s_cashprice> 200,000</s_cashprice><s_changeprice> 25,400</s_changeprice></s_total>
C:\Users\draj0.CHANDRIKAV\Downloads>
```

2. Text Decoder (BART-like Transformer)

Autoregressively generates structured text (e.g., JSON, answers).

Parameters:

BART Base decoder → ~139 million parameters.

3. Total Parameters

For **Donut-base**: around **220–230 million parameters** (Vision encoder + Text decoder).

For finetuned versions (CORD, DocVQA, SROIE, etc.), the **parameter count stays the same**—only weights are updated with domain-specific data.

The model is **lighter than very large LLMs** (like GPT-3/4), making it practical for **document AI on GPUs with limited memory**.

Despite its size, Donut achieves strong performance because it removes the OCR step, learning **directly from pixels to text**

- Donut performed better and faster than OCR-based models.
- It gave higher accuracy on tasks like document classification and QA.
- The model works well across English, Korean, and Chinese text.
- It's simpler to run since it doesn't depend on OCR software.

For **finetuned versions** (like donut-base-finetuned-cord-v2 for receipts):

You just give the **task prompt** (e.g., <s_cord-v2>), and it generates structured JSON. This works almost like **one-shot**, since no extra examples are required at inference.

For **general models (docvqa, pretrained)**:

You can add **question prompts** (e.g., “*What is the total amount?*”), and Donut generates the answer. This is closer to **few-shot prompting**, where the model relies on its training examples to generalize.

- Donut outputs clean and organized text from images.
- Example: it can turn a receipt into a JSON file with store name and total price.
- It works even with handwriting or noisy images.
- The model automatically focuses on important areas in the document.
- Donut is more accurate than models like LayoutLMv2.

Model Type

Multimodal Transformer (Vision + Language).

It is an **encoder-decoder** model:

Encoder → Swin Transformer (for images).

Decoder → BART-style Transformer (for text generation).

```
MODEL_ID = "naver-clova-ix/donut-base-finetuned-docvqa"
```

```
task_prompt = "<s_docvqa><s_question>What does this document say?</s_question><s_answer>"
```

CORD V2 (receipts, English dataset)

```
MODEL_ID = "naver-clova-ix/donut-base-finetuned-cord-v2"
```

```
task_prompt = "<s_cord-v2>"
```

SROIE

```
MODEL_ID = "naver-clova-ix/donut-base-finetuned-sroie"
```

```
task_prompt = "<s_sroie>"
```

RVL-CDIP (document classification: letter, invoice, form, email, etc.)

```
MODEL_ID = "naver-clova-ix/donut-base-finetuned-rvlcdip"
```

```
task_prompt = "<s_rvlcdip>"
```

Example swap for receipts (SROIE):

```
MODEL_ID = "naver-clova-ix/donut-base-finetuned-sroie"
```

```
task_prompt = "<s_sroie>"
```

Category

Donut is **not a pure LLM** like GPT, nor a diffusion model.

It belongs to the class of **OCR-free vision-language models**, meaning it directly learns the mapping from **document pixels → text/JSON output**.

Key Idea

Instead of detecting text with OCR and then passing it through NLP, Donut **generates text tokens directly** from the image features, making it an **end-to-end generative document parser**.

- **Strengths:** Works without OCR, handles many languages, and is open-source.
- It's easy to reproduce and gives fast results.
- **Limitations:** Needs strong GPUs for large images.
- Can make mistakes if the JSON output structure breaks.
- New versions like Donut v2 and LayoutLMv3 build on its ideas.
- It's part of a growing trend in **vision + language AI**.
- Donut changed document AI by removing the need for OCR.