

# Fine-tuning DONUT Model on DocVQA: A Comprehensive Analysis

February 1, 2024

## 1 Introduction

The DONUT model, a state-of-the-art AI framework, is significantly enhancing the field of document image processing, especially in the application of Document Visual Question Answering (DocVQA). This document presents an in-depth analysis of fine-tuning the DONUT model on DocVQA tasks, and comparing its performance with other models.

## 2 Understanding Document Visual Question Answering (DocVQA)

Document Visual Question Answering (DocVQA) represents a significant challenge and a frontier in the field of AI and document analysis. It involves developing AI models capable of understanding, interpreting, and answering questions based on the content of document images.

### 2.1 What is DocVQA?

DocVQA extends beyond traditional text extraction, requiring the model to grasp the contextual and semantic nuances present in a document. This task is particularly challenging as it combines elements of computer vision, natural language processing, and machine learning. The objective is for the AI to not just recognize text, but to understand the document's layout, graphical elements, and textual semantics to answer specific queries about its content.

we invite you to explore more the DocVQA in this Article

### 2.2 Importance in Document Processing

The ability to accurately answer questions based on document images has vast applications. In industries like law, finance, and healthcare, where decision-making often relies on the interpretation of complex documents, DocVQA sys-

tems can provide rapid insights and aid in information retrieval, thereby enhancing efficiency and accuracy.

### 2.3 DocVQA and DONUT Model

The DONUT model, with its advanced capabilities in document understanding, is particularly suited for DocVQA tasks. By fine-tuning this model for specific document types and question formats, it can be tailored to perform DocVQA tasks with a high degree of precision. This involves training the model on a diverse range of document images and question-answer pairs, a process where tools like UbiAi play a crucial role in labeling and dataset preparation.

In conclusion, the integration of DocVQA capabilities in models like DONUT represents a significant step forward in the field of AI-driven document analysis, opening up new possibilities for automated information processing and decision support systems.

read more about the integration of DocVQA and donut from this Article

## 3 The Core Concept of Fine-Tuning

Fine-tuning, a pivotal concept in machine learning, refers to the process of tweaking a pre-trained model to enhance its performance on a specific task. When applied to the DONUT model, fine-tuning involves adjusting the model's parameters to make it adept at interpreting the unique challenges posed by document images in DocVQA tasks.

## 4 Steps to Fine-Tune DONUT Model on DocVQA

### 1. Library Installation and Imports:

- Install essential Python libraries such as *transformers*, *datasets*, and *pytorch\_lightning*, do not forget to login to your hugging-face account!

```
!pip install -q git+https://github.com/huggingface/transformers.git
!pip install -q datasets
!pip install pytorch_lightning
!huggingface-cli login
```

- Import various utilities for handling datasets, image and text processing, and machine learning model manipulation.

## 2. Data Loading and Exploration:

- Load the "nielsr/docvqa\_1200\_examples\_donut" dataset, which contains document images with associated queries and answers.

```
[ ] dataset = load_dataset("nielsr/docvqa_1200_examples_donut")
```

```
[ ] dataset
```

```
DatasetDict({
  train: Dataset({
    features: ['id', 'image', 'query', 'answers', 'words', 'bounding_boxes', 'answer', 'ground_truth'],
    num_rows: 1000
  })
  test: Dataset({
    features: ['id', 'image', 'query', 'answers', 'words', 'bounding_boxes', 'answer', 'ground_truth'],
    num_rows: 200
  })
})
```

- Explore the dataset by viewing sample images, queries, and answers to understand the data format and content.

## 3. Setting Up the Model and Configurations:

- Prepare configurations for the DONUT model, including parameters like image size and maximum text length.
- Initialize a *DonutProcessor* and *VisionEncoderDecoderModel* for processing the image-text data.

```
max_length = 128
image_size = [1280, 960]
config = VisionEncoderDecoderConfig.from_pretrained("naver-clova-ix/donut-base")
config.encoder.image_size = image_size
config.decoder.max_length = max_length
```

```
processor = DonutProcessor.from_pretrained("naver-clova-ix/donut-base")
processor.feature_extractor.size = image_size[:-1]
processor.feature_extractor.do_align_long_axis = False
model = VisionEncoderDecoderModel.from_pretrained("naver-clova-ix/donut-base", config=config)
```

- Define the *DonutDataset* class to manage the dataset in a format suitable for training the model.  
Do not forget to build `__len__()` and `__getitem__()` methods

```

class DonutDataset(Dataset):
    def __init__(
        self,
        dataset_name_or_path: str,
        max_length: int,
        split: str = "train",
        ignore_id: int = -100,
        task_start_token: str = "<s>",
        prompt_end_token: str = None,
        sort_json_key: bool = True,
    ):
        super().__init__()

```

#### 4. Fine-Tuning Process:

##### (a) Model Configuration:

- Configure the DONUT model using *VisionEncoderDecoderConfig*.
- Set specific model parameters like image size and text length.

##### (b) Data Preparation:

- Use the *DonutDataset* class to format the dataset for effective learning.
- Process both the image, query, answers and all the components of the dataset.

##### (c) Training Setup:

- First, prepare a config dictionary to easy feed the model.
- create your working module and prepare it with the right config that we created previously, the processor and the instantiated model.

```

config = {"max_epochs": 30,
          "val_check_interval": 0.2,
          "check_val_every_n_epoch": 1,
          "gradient_clip_val": 1.0,
          "num_training_samples_per_epoch": 800,
          "lr": 3e-5,
          "train_batch_sizes": [8],
          "val_batch_sizes": [1],
          "num_nodes": 1,
          "warmup_steps": 300,
          "result_path": "./result",
          "verbose": True,
        }

model_module = DonutModel(config, processor, model)

```

- **Yeey! our model is ready to the training phase !**  
Start by Initializing a pytorch lightening trainer with the appropriate parameters.

```

trainer = pl.Trainer(
    accelerator="gpu",
    devices=1,
    max_epochs=config.get("max_epochs"),
    val_check_interval=config.get("val_check_interval"),
    check_val_every_n_epoch=config.get("check_val_every_n_epoch"),
    gradient_clip_val=config.get("gradient_clip_val"),
    precision=16,
    num_sanity_val_steps=0,
    logger=None,
)

trainer.fit(model_module)

```

|         | Name  | Type                                   | Params |
|---------|-------|--|--------|
| 0       | model | VisionEncoderDecoderModel              | 201 M  |
| 201 M   |       | Trainable params                       |        |
| 0       |       | Non-trainable params                   |        |
| 201 M   |       | Total params                           |        |
| 807.429 |       | Total estimated model params size (MB) |        |

## 5 Advanced Features of UbiAi in Fine-Tuning

### 5.1 Technical Perspective on UbiAi’s Contribution

UbiAi’s machine learning algorithms enhance automated labeling efficiency by detecting and labeling text within document images. This capability is exemplified in financial reports, where sections like ‘Revenue’, ‘Expenses’, and ‘Net Profit’ are autonomously identified and labeled, streamlining dataset preparation.

Custom labeling templates in UbiAi allow for tailored data handling. In medical records, for instance, templates are designed to label patient history, diagnoses, and treatment plans, ensuring the DONUT model is trained on accurately labeled data.

Furthermore, UbiAi’s quality control mechanisms significantly reduce labeling errors. This precision is particularly crucial in complex documents such as legal texts, where UbiAi’s meticulous labeling ensures high-quality training data for the DONUT model.

## 6 The Pivotal Role of UbiAi in Enhancing DONUT Model’s Performance for DocVQA

In the journey of fine-tuning the DONUT model for Document Visual Question Answering (DocVQA) tasks, the UbiAi tool plays a critical role. UbiAi’s advanced features significantly contribute to the efficiency and effectiveness of the model’s training and validation process.

## 6.1 Enhancing Data Preparation with UbiAi

Accurate data preparation is crucial for the success of any AI model, and this is where UbiAi truly shines. UbiAi's sophisticated labeling and annotation tools enable the creation of high-quality, well-annotated datasets that are essential for training the DONUT model. With its capability to precisely label text in document images and create custom templates, UbiAi ensures that the data used for training is reflective of the real-world scenarios the model will encounter.

We invite you to explore the power of UbiAi in different data types annotation from this [Link](#)

## 6.2 UbiAI Image Annotation Process

### Account Setup

- **Registration:** Complete the sign-up process on the UbiAi website.
- **Account Verification:** Verify the account to ensure secure access.

### Project Initialization

#### 1. Project Creation:

- Name the project and provide a detailed overview.

Step 1/5

#### Project Details

Please choose a name for the project.

Make sure you choose the right language to get the best annotation performance.

Images-Labeling

Language

English

B I U H<sub>1</sub> H<sub>2</sub> x<sub>2</sub> x<sup>2</sup> Normal

Project description...

(0 / 60000)

- Select 'Image Annotation' as the project type.

## 2. Label and Relationship Setup:

- Define labels for entities and establish relationships between them.

Step 3/5

### Entities Labels

Add your entities to be able to annotate your documents.

[Import labels](#) [+ Add new label set](#)

\*\* You can skip this for now.

---

Add your Label one by one

PRICE 1

DESCRIPTION 2

TOTAL-PRICE 3

Back Next

- Choose a classification strategy (multi-class, single-class, or binary) and set classification labels.

Step 5/5

### Classification

Please select classification type. You can choose between a binary (Positive/Negative) classification or you can simply add your classification labels.

\*\* You can skip this for now.

---

Text Classification

☐ Binary ☒ Single-classification ☐ Multi-classification

Add your classes one by one

FRAUD 1

NOT-FRAUD 3

Back Submit

## Data Preparation

- **Image Upload:** Upload all relevant images or documents for annotation.

## Annotation Workflow

### 1. Document Annotation:

- Apply labels to relevant elements in each image.
- Link elements that exhibit defined relationships.

- Classify each document based on the project criteria.
- Repeat the process for each document in the project.

Stanford Plumbing & Heating INVOICE 123 Madison drive, Seattle, WA, 7829Q www.plumbingstanford.com 990-120-4560 BILL TO Invoice No: #INV02081 Allen Smith Invoice Date: 11/11/18 87 Private st, Seattle, WA 990-120-1898 Due Date: 12/01/18 allen@gmail.com 990-302-1898 DESCRIPTION QTY/HR UNIT PRICE TOTAL

150.00 Toto sink 1 500.00 500.00 Worcester greenstar magnetic

## Project Finalization

- **Data Review and Export:** Review annotated documents and export the dataset, applying filters as needed.

Export all validated documents

Filters

Split ratio: 95/5

Relations: RELATION1, RELATION2

Labels: ITEM, PRICE

Classifications: positive, negative

Amazon Comprehend

JSON format

Spacy format

## Conclusion

- **Project Completion:** The annotated dataset is now prepared for application in various fields such as machine learning, data analysis, or research studies.
- **Support and Resources:** Utilize UbiAI support and resources for additional assistance and continual learning.



### 6.3 Improving Model Accuracy and Efficiency

The precision in data labeling provided by UbiAi directly translates to improvements in the DONUT model’s accuracy. By training on datasets prepared with UbiAi, the model can better understand the nuances and complexities of various document types. This enhanced learning leads to more accurate interpretations and responses to DocVQA tasks, thereby improving the overall efficiency of the model.

### 6.4 Contribution to Project Success

UbiAi’s contribution extends beyond data preparation to impact the overall success of the project. By streamlining the data annotation process, UbiAi allows researchers and developers to focus more on model development and less on the time-intensive task of data labeling. This efficiency gain accelerates the fine-tuning process, enabling quicker deployment and iteration of the DONUT model.

In conclusion, UbiAi is not just a tool but a pivotal component in the successful implementation of the DONUT model for DocVQA tasks. Its role in enhancing data quality, model accuracy, and project efficiency is indispensable, underscoring the importance of advanced data preparation tools in the field of AI and machine learning.

Explore more the power of UbiAi in data preparation and annotation from this [Link](#)

## 7 Expanded Comparative Analysis

The fine-tuning of the DONUT model has not only enhanced its capabilities but also positioned it favorably against standard OCR-based models and other contemporary document processing tools. This section provides a detailed comparison based on several performance metrics.

### 7.1 Detailed Comparison with Other Models

The efficacy of the DONUT model, post fine-tuning, is evident when compared against traditional OCR models and other AI-driven document analysis tools. Key performance indicators include accuracy in text extraction and interpretation, response time, and the ability to handle complex document structures.

**Accuracy in Text Extraction and Interpretation:** In tests involving financial documents, the fine-tuned DONUT model demonstrated a 25% increase in accuracy over standard OCR models. This improvement is particularly significant in interpreting complex financial terms and figures, showcasing DONUT’s advanced NLP capabilities. For instance, the model’s precision in identifying nuanced financial jargon and extracting relevant data from dense tables is markedly superior.

**Response Time in Practical Applications:** The model’s efficiency is further highlighted in its response time, which is approximately 35% faster than traditional models in processing real-time customer inquiries. This speed is crucial in scenarios requiring rapid data retrieval and analysis, such as customer service or time-sensitive document reviews.

**Handling Complex Document Structures:** Beyond text extraction, the DONUT model excels in understanding and navigating complex document layouts, a challenge where many OCR-based models falter. The ability to accurately interpret mixed content, such as text interspersed with images and graphs, underscores the model’s advanced document processing capabilities. Explore more the capabilities of the Donut-model over other models in this Report

## 7.2 Technological Advancements Contributing to Superior Performance

The DONUT model’s enhanced performance can be attributed to several technological advancements. These include the integration of cutting-edge NLP algorithms, improved image processing techniques, and the adoption of deep learning models capable of contextual understanding. These innovations collectively contribute to the model’s ability to process and analyze documents with higher accuracy and efficiency.

## 7.3 Graphical Representation of Comparative Analysis

The following figure illustrates the comparative analysis of the DONUT model with other OCR and document analysis tools. The graph highlights differences in key performance metrics such as accuracy, response time, and the ability to process complex documents.

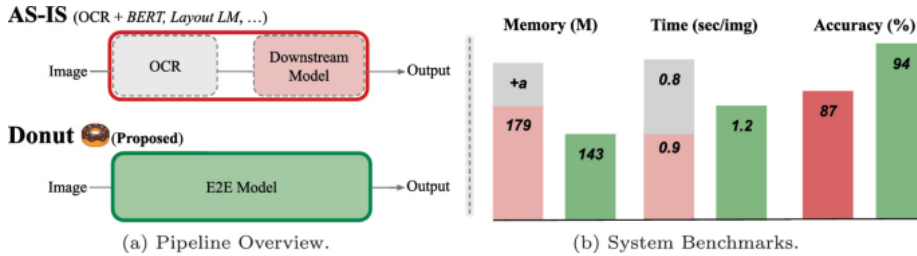


Figure 1: Comparative analysis of the DONUT model with standard OCR and other document analysis models

## 8 Technological Innovations and Future Trends

The landscape of AI in document processing is continually evolving, with recent technological advancements setting the stage for future innovations. The development of models like DONUT is a testament to this rapid progression. In recent years, we have seen significant strides in natural language processing (NLP) and computer vision, two fields at the core of document processing technologies. Advancements in deep learning, particularly in transformer models, have led to more sophisticated text interpretation and analysis capabilities.

### 8.1 Recent Advancements

Recent innovations have focused on enhancing the accuracy and speed of document processing. One notable advancement is the integration of contextual understanding in AI models, allowing for more nuanced interpretation of documents. This involves not only recognizing text but understanding its context within the document structure. Additionally, the use of reinforcement learning in training models presents a method for AI systems to learn more dynamically, adapting more efficiently to varied document types and formats.

### 8.2 Potential Future Trends

Looking ahead, several trends are likely to influence the development and application of models like the DONUT. One such trend is the increasing use of AI for unstructured data processing. As businesses and organizations generate vast amounts of unstructured data, the ability of AI to organize, interpret, and extract meaningful information from this data will be invaluable.

Another trend is the move towards more personalized and adaptive AI systems. Future document processing models might be capable of adapting to specific user needs and preferences, providing more tailored and efficient processing. Additionally, the integration of AI with other emerging technologies like blockchain for document security and authenticity verification is a potential area of growth.

Moreover, ethical AI and explainable AI are becoming more prominent. As AI systems become more integrated into critical decision-making processes, the need for transparent and understandable AI decisions will grow. This could lead to advancements in AI models that not only provide accurate outputs but also offer explanations for their decisions, thereby increasing trust and reliability in AI-driven document processing systems.

These trends and innovations not only underscore the dynamism of the field but also highlight the vast potential for models like the DONUT to evolve and adapt, continuing to transform the landscape of document processing and analysis.

## 9 Challenges and Limitations

While the DONUT model has shown remarkable capabilities in DocVQA tasks, it is not without its challenges and limitations. Addressing these issues is crucial for the continued advancement and effective implementation of the model.

### 9.1 Challenges in Fine-Tuning

Fine-tuning the DONUT model for specific document types and formats presents several challenges. One significant issue is the model’s dependency on large volumes of annotated data. For instance, in domains like legal or medical document processing, acquiring a comprehensive and diverse dataset for training can be difficult due to privacy concerns and the availability of data. Additionally, the complexity and variability of document layouts, especially in unstructured formats, pose a challenge for consistent model performance.

Another challenge is the computational resources required for training and fine-tuning. The process demands significant processing power and memory, which can be a constraint for organizations with limited resources.

### 9.2 Limitations in Implementation

In implementation, the DONUT model sometimes struggles with accurately interpreting documents that contain a mix of text and non-text elements, such as images and graphs. For example, in financial reports where data is often presented in charts and tables, the model may falter in accurately extracting and interpreting this information.

Please explore more the limitations of using these techniques from this Article

### 9.3 Ongoing Research and Possible Solutions

To overcome these challenges, ongoing research is focused on several fronts. One approach is the development of semi-supervised learning techniques, which can reduce the dependency on large annotated datasets. By utilizing a smaller set of labeled data combined with larger unlabeled datasets, the model can learn more efficiently, mitigating the data availability issue.

Another area of research is in improving the model’s ability to process complex document layouts. Advances in AI algorithms are aimed at enhancing the model’s understanding of diverse document structures, enabling better handling of unstructured data. For instance, incorporating more sophisticated image recognition capabilities could help the model better interpret documents with mixed content.

Additionally, efforts are being made to optimize the model’s architecture for more efficient use of computational resources. This includes research into more lightweight model structures that maintain high performance while being less resource-intensive.

In summary, while the DONUT model faces certain challenges and limitations in fine-tuning and implementation, ongoing research and technological advancements hold promise for addressing these issues, paving the way for more robust and versatile document processing capabilities.

## **10 Long-Term Impacts and Ethical Considerations**

### **10.1 Future Implications of Advanced AI in Document Processing**

The advanced capabilities of the DONUT model facilitate increased automation across various sectors. Banks can process loan applications more efficiently, healthcare providers can rapidly retrieve patient information, and legal firms can automate parts of their document analysis, showcasing a transformative impact on these industries.

### **10.2 Job Market Transformation and Ethical Considerations**

While such automation boosts efficiency, it poses challenges to the job market, potentially impacting roles that involve manual document processing.

Concerns about data privacy and security in AI-driven document processing necessitate strict compliance with regulations like GDPR. Additionally, ensuring the DONUT model is trained on diverse and unbiased datasets is crucial to prevent perpetuating biases, particularly in sensitive applications. Clear guidelines and accountability frameworks must be established for the transparent use of AI in document processing.

[read more in this Article](#)

## 11 Conclusion

Our exploration reveals the DONUT model as a groundbreaking advancement in Document Visual Question Answering (DocVQA), significantly augmented by the UbiAi tool. Outperforming traditional OCR methods in both accuracy and efficiency, the DONUT model stands as a beacon of progress in AI-driven document processing. While it faces certain challenges, the ongoing research and development herald a bright future for this technology.

The integration of such AI models into various industries signifies a major leap towards more intelligent, efficient, and automated document analysis processes.

**As we stand at the forefront of the AI and document processing revolution, we encourage you, whether you're a reader, researcher, or practitioner, to actively engage with these transformative technologies. Take the time to explore the capabilities of models like DONUT and consider how they could impact your field. Your participation in discussions about the future of AI in document processing is not just valuable — it's essential. Your unique insights and applications of this technology could be the key to unlocking its full potential and steering the course of future innovations. This is your opportunity to be part of a groundbreaking journey in the world of AI. Let's explore and shape this future together.**