# Fine-tuning question answering model with layoutLM v2

Welcome to our exploration of fine-tuning LayoutLMv2 for Document Question Answering. In this article, we'll delve into the practical aspects of refining LayoutLMv2, a versatile model known for its ability to understand both text and visual layout in documents. Document Question Answering has practical applications, and we'll guide you through the steps of fine-tuning, from data preparation to model configuration. Specifically we will cover :

1. Document Question Answering
2. DocVQA dataset
3. Environment
4. Data loading
5. Data processing
6. Training
7. Inference
8. Conclusion

**LayoutLMv2: Multi-modal Pre-training for Visually-rich Document Understanding**

Yang Xu[1*], Yiheng Xu[2*], Tengchao Lv[2*], Lei Cui[2], Furu Wei[2], Guoxin Wang[3], Yijuan Lu[3], Dinei Florencio[3], Cha Zhang[3], Wanxiang Che[1], Min Zhang[4], Lidong Zhou[2]
[1]Research Center for Social Computing and Information Retrieval, Harbin Institute of Technology
[2]Microsoft Research Asia    [3]Microsoft Azure AI    [4]Soochow University
[1]{yxu,car}@ir.hit.edu.cn,
[2]{v-yixu,v-telv,lecu,fuwei,lidongz}@microsoft.com,
[3]{guow,yijlu,dinei,chazhang}@microsoft.com [4]minzhang@suda.edu.cn

Source : LayoutLMV2

## 1- Document Question Answering:

Document question answering (DQA) is a cutting-edge natural language processing task that focuses on developing models capable of extracting

information from large textual sources to answer user queries. Unlike traditional question answering, which relies on predefined knowledge bases or structured databases, DQA involves comprehending unstructured documents, such as articles, reports, or books, to generate accurate and contextually relevant responses.

## 2- DocVQA dataset:

The DocVQA dataset serves as a crucial resource in the realm of Document Visual Question Answering (DocVQA), offering a standardized benchmark for evaluating the capabilities of models designed to comprehend textual and visual information within documents. Comprising images of diverse documents, accompanied by corresponding questions and answers, this dataset challenges researchers and practitioners in the fields of computer vision and natural language processing to develop robust algorithms capable of simultaneously processing visual content and answering questions related to the document's context.
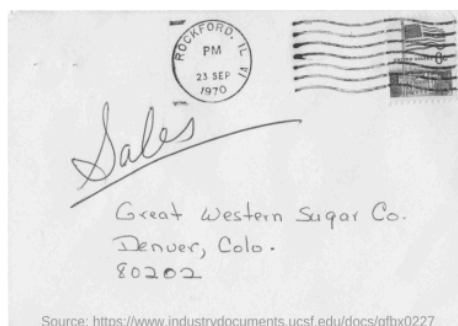
**DocVQA: A Dataset for VQA on Document Images**

Minesh Mathew[1]    Dimosthenis Karatzas[2]    C.V. Jawahar[1]
[1]CVIT, IIIT Hyderabad, India    [2]Computer Vision Center, UAB, Spain
minesh.mathew@research.iiit.ac.in, dimos@cvc.uab.es, jawahar@iiit.ac.in

**Abstract**

*We present a new dataset for Visual Question Answering (VQA) on document images called DocVQA. The dataset consists of 50,000 questions defined on 12,000+ document images. Detailed analysis of the dataset in comparison with similar datasets for VQA and reading comprehension is presented. We report several baseline results by adopting existing VQA and reading comprehension models. Although the existing models perform reasonably well on certain types of questions, there is large performance gap compared to human performance (94.36% accuracy). The models need to improve specifically on questions where understanding structure of the document is crucial. The dataset, code and leaderboard are available at docvqa.org*

**1. Introduction**

Source: https://www.industrydocuments.ucsf.edu/docs/gfbx0227

**Q:** Mention the ZIP code written?
**A:** 80202
**Q:** What date is seen on the seal at the top of the letter?
**A:** 23 sep 1970
**Q:** Which company address is mentioned on the letter?
**A:** Great western sugar Co.

Source : DocVQA

The tasks involved in DocVQA demand a fusion of image processing techniques to interpret the visual aspects of documents and sophisticated natural language understanding to extract relevant information for accurate responses. As the dataset evolves and expands, it continues to play a pivotal role in advancing the state-of-the-art in document understanding, fostering innovation in AI systems tailored for complex document-based tasks. Researchers frequently leverage the DocVQA dataset to push the boundaries of what is achievable in terms of visual and textual comprehension, with the ultimate goal of enhancing the efficiency and accuracy of question answering systems in real-world document analysis scenarios.

## 3-Environment:

Prior to starting, ensure that you have installed all the required libraries. LayoutLMv2 relies on detectron2, torchvision, and tesseract.

```
!pip install -q transformers datasets
```

```
507.1/507.1 kB 4.6 MB/s eta 0:00:00
115.3/115.3 kB 6.9 MB/s eta 0:00:00
134.8/134.8 kB 7.0 MB/s eta 0:00:00
```

Source : Author

```
!pip install 'git+https://github.com/facebookresearch/detectron2.git'
!pip install torchvision
```

```
Collecting git+https://github.com/facebookresearch/detectron2.git
  Cloning https://github.com/facebookresearch/detectron2.git to /tmp/pip-req-build-0vcghgje
  Running command git clone --filter=blob:none --quiet https://github.com/facebookresearch/detectron2.git /tmp/pip-req-build-0vcghgje
  Resolved https://github.com/facebookresearch/detectron2.git to commit 864913f0e57e87a75c8cc0c7d79ecbd774fc669b
  Preparing metadata (setup.py) ... done
Requirement already satisfied: Pillow>=7.1 in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (9.4.0)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (3.7.1)
Requirement already satisfied: pycocotools>=2.0.2 in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (2.0.7)
Requirement already satisfied: termcolor>=1.1 in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (2.4.0)
Collecting yacs>=0.1.8 (from detectron2==0.6)
  Downloading yacs-0.1.8-py3-none-any.whl (14 kB)
Requirement already satisfied: tabulate in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (0.9.0)
Requirement already satisfied: cloudpickle in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (2.2.1)
Requirement already satisfied: tqdm>4.29.0 in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (4.66.1)
Requirement already satisfied: tensorboard in /usr/local/lib/python3.10/dist-packages (from detectron2==0.6) (2.15.1)
Collecting fvcore<0.1.6,>=0.1.5 (from detectron2==0.6)
  Downloading fvcore-0.1.5.post20221221.tar.gz (50 kB)
```

Source : Author

```
!sudo apt install tesseract-ocr
!pip install -q pytesseract
```

```
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following additional packages will be installed:
  tesseract-ocr-eng tesseract-ocr-osd
The following NEW packages will be installed:
  tesseract-ocr tesseract-ocr-eng tesseract-ocr-osd
0 upgraded, 3 newly installed, 0 to remove and 30 not upgraded.
Need to get 4,816 kB of archives.
After this operation, 15.6 MB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu jammy/universe amd64 tesseract-ocr-eng all 1:4.00~git30-7274cfa-1.1 [1,591 kB]
Get:2 http://archive.ubuntu.com/ubuntu jammy/universe amd64 tesseract-ocr-osd all 1:4.00~git30-7274cfa-1.1 [2,990 kB]
Get:3 http://archive.ubuntu.com/ubuntu jammy/universe amd64 tesseract-ocr amd64 4.1.1-2.1build1 [236 kB]
Fetched 4,816 kB in 1s (5,677 kB/s)
```

Source : Author

Now let's define some global variables:

```
model_checkpoint = "microsoft/layoutlmv2-base-uncased"
batch_size = 4
```

## 4- Data Loading:

In this tutorial, we utilize a compact sample of preprocessed DocVQA. If you prefer working with the complete DocVQA dataset, you can register and obtain it from the DocVQA homepage.

```
from datasets import load_dataset
```

```
dataset = load_dataset("nielsr/docvqa_1200_examples")
dataset
```

The dataset has already been divided into training and testing sets.

```
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:88: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret in your
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
  warnings.warn(
Downloading metadata: 100%    1.74k/1.74k [00:00<00:00, 90.5kB/s]
Downloading data: 100%        123M/123M [00:05<00:00, 29.2MB/s]
Downloading data: 100%        25.2M/25.2M [00:01<00:00, 21.1MB/s]
Generating train split: 100%  1000/1000 [00:01<00:00, 518.80 examples/s]
Generating test split: 100%   200/200 [00:00<00:00, 394.48 examples/s]
DatasetDict({
    train: Dataset({
        features: ['id', 'image', 'query', 'answers', 'words', 'bounding_boxes', 'answer'],
        num_rows: 1000
    })
    test: Dataset({
        features: ['id', 'image', 'query', 'answers', 'words', 'bounding_boxes', 'answer'],
        num_rows: 200
    })
})
```

```
dataset["train"].features
```

In the dataset, each example is characterized by several fields. The "id" serves as a unique identifier for the example, while the "image" field contains a PIL.Image.Image object representing the document image. The "query" field encapsulates the natural language question posed in various languages, and the corresponding correct answers provided by human annotators are stored in the "answers" field. The "words" and "bounding_boxes" fields hold the results of Optical Character Recognition (OCR), which we won't be utilizing in this context. The "answer" field contains responses matched by a different model, and for our purposes, we will exclude this feature.

```
{'id': Value(dtype='string', id=None),
 'image': Image(decode=True, id=None),
 'query': {'de': Value(dtype='string', id=None),
  'en': Value(dtype='string', id=None),
  'es': Value(dtype='string', id=None),
  'fr': Value(dtype='string', id=None),
  'it': Value(dtype='string', id=None)},
 'answers': Sequence(feature=Value(dtype='string', id=None), length=-1, id=None),
 'words': Sequence(feature=Value(dtype='string', id=None), length=-1, id=None),
 'bounding_boxes': Sequence(feature=Sequence(feature=Value(dtype='float32', id=None), length=4, id=No
 'answer': {'match_score': Value(dtype='float64', id=None),
  'matched_text': Value(dtype='string', id=None),
  'start': Value(dtype='int64', id=None),
  'text': Value(dtype='string', id=None)}}
```

To streamline the dataset, we'll filter for only English questions and discard the "answer" feature. Additionally, we'll retain only the first answer from the set provided by annotators, or you can opt to randomly sample it for further analysis.

```
updated_dataset = dataset.map(lambda example: {"question":
example["query"]["en"]}, remove_columns=["query"])
updated_dataset = updated_dataset.map(
    lambda example: {"answer": example["answers"][0]},
remove_columns=["answer", "answers"]
)
```

```
Map: 100%  ████████████████  1000/1000 [00:04<00:00, 214.30 examples/s]
Map: 100%  ████████████████  200/200 [00:00<00:00, 1089.48 examples/s]
Map: 100%  ████████████████  1000/1000 [00:05<00:00, 174.78 examples/s]
Map: 100%  ████████████████  200/200 [00:00<00:00, 1040.44 examples/s]
```

## 5- Data Processing:

It's important to note that the LayoutLMv2 checkpoint utilized in this guide has been trained with a maximum position embedding value of 512, as indicated in the checkpoint's config.json file. While we can truncate

examples, it's crucial to avoid situations where the answer might be positioned at the end of a lengthy document and end up being truncated. For that purpose we only keep those that are most likely to not exceed 512.

```
updated_dataset = updated_dataset.filter(lambda x: len(x["words"]) +
len(x["question"].split()) < 512)
```

```
Filter: 100%    ████████████████████    1000/1000 [00:12<00:00, 79.22 examples/s]
Filter: 100%    ████████████████████    200/200 [00:01<00:00, 105.82 examples/s]
```

Source : Author

At this stage, we will remove the Optical Character Recognition (OCR) features from the dataset. These features were generated during the OCR process for fine-tuning a different model. However, they require additional processing to align with the input requirements of the model used in this guide.

```
updated_dataset = updated_dataset.remove_columns("words")
updated_dataset = updated_dataset.remove_columns("bounding_boxes")
```

Here's an example of an how the data looks like now :

PHILIP MORRIS U. S. A.

INTER-OFFICE CORRESPONDENCE

Richmond, Virginia

**To:** Mr. James L. Myracle      **Date:** April 27, 1990

**From:** H. L. Spielberg

**Subject:** Flavor Development Monthly Summary for April, 1990    *C90-03267*

PROJECT ART:

Samples of the GMC, 8842-5 aftercut and POL 0704 models were sent to ARD for method development and quantification of GMC on filler and in solution. Preliminary analytical data indicate approximately 75% of the target level of GMC is on the filler. Additional models were also submitted for analyses. These included the 20% and 40% reduced level, control level and 20% increased level of GMC on filler.

The lowest acceptable OV study has been initiated on the 9 mg non-menthol salesman samples. Models are complete at 10.8, 11.3, 11.9, 12.7, 13.0 and 13.6% OV. These models will be subjectively evaluated against a control product at 13.0% OV.

Production of Next KS and 100's for the Tampa, Florida test market began C Shift April 23, 1990 at the MC. Menthol production is complete and non-menthol production will be complete by the end of the week.

Evaluations of citric acid stems from single versus double-batched trials have shown promise at the 2% substitution level in Pilot RL. After the Flavor Development panel showed the single to be comparable to double-batched stems, results of the MC panel showed no differences between the control and double-batched test RL in Marlboro.

Similar comparisons at the 5% substitution level in production RCB are in progress. Combinations of test RL and RCB will be selected to determine utilization levels. Qualification of double-batched stems at specific levels should allow usage of single or half-batched stems being tested at Bermuda Hundred.

Nippon ART concept Danchi test to evaluate tar and tipping color variations has been shipped. A second Danchi test to evaluate ART menthol levels has been shipped. Small-scale samples were made to evaluate PM Super Lights and "Castor" type flavors.

TAR/NICOTINE INTERACTION:

Subjective Studies:
A Tar/Nicotine interaction study to assess 100 mm models from 3 to 15 mg tar and from 0.3 to 1.5% filler nicotine is in progress.

2022156205

Source : Author

For the Document Question Answering task, which involves multiple modalities, it is crucial to preprocess inputs from each modality according to the model's expectations. To initiate this process, we'll begin by loading the LayoutLMv2Processor. This processor internally integrates an image processor capable of handling image data and a tokenizer designed to encode text data. This combined functionality allows for comprehensive preprocessing of both image and text inputs to meet the model's requirements.

```
from transformers import AutoProcessor
```

```
processor = AutoProcessor.from_pretrained(model_checkpoint)
```

- First let's process the documents  :

```
image_processor = processor.image_processor


def get_ocr_words_and_boxes(examples):
    images = [image.convert("RGB") for image in examples["image"]]
    encoded_inputs = image_processor(images)

    examples["image"] = encoded_inputs.pixel_values
    examples["words"] = encoded_inputs.words
    examples["boxes"] = encoded_inputs.boxes

    return examples
```

And apply the processing to the whole dataset

```
dataset_with_ocr =
updated_dataset.map(get_ocr_words_and_boxes, batched=True,
batch_size=2)
```

Source : Author

- Now let's process the Text using the tokenizer:

```
tokenizer = processor.tokenizer
def subfinder(words_list, answer_list):
    matches = []
    start_indices = []
    end_indices = []
    for idx, i in enumerate(range(len(words_list))):
        if words_list[i] == answer_list[0] and words_list[i : i +
len(answer_list)] == answer_list:
            matches.append(answer_list)
            start_indices.append(idx)
            end_indices.append(idx + len(answer_list) - 1)
    if matches:
        return matches[0], start_indices[0], end_indices[0]
    else:
        return None, 0, 0
```

This function takes two input lists, words_list and answer_list. It iterates over words_list, checks for a match with the first word of answer_list, and verifies if the sublist of words_list starting from the current word and of the same length as answer_list is equal to answer_list. If a match is found, it records the match's starting index (idx) and ending index (end_idx). If multiple matches occur, the function returns only the first one. If no match is found, it returns (None, 0, 0).

To illustrate how this function finds the position of the answer, let's use it on an example:

```python
example = dataset_with_ocr[1]
words = [word.lower() for word in example["words"]]
match, word_idx_start, word_idx_end = subfinder(words,
example["answer"].lower().split())
print("Question: ", example["question"])
print("Words:", words)
print("Answer: ", example["answer"])
print("start_index", word_idx_start)
print("end_index", word_idx_end)
```

```
Question:  Who is in  cc in this letter?
Words: ['wie', 'baw', 'brown', '&', 'williamson', 'tobacco', 'corporation', 'research',
Answer:  T.F. Riehl
start_index 17
end_index 18
```

Source : Author

After the encoding process, the examples will have the following appearance:

```python
encoding = tokenizer(example["question"], example["words"],
example["boxes"])
tokenizer.decode(encoding["input_ids"])
```

```
'[CLS] who is in cc in this letter? [SEP] wie baw brown & williamson tobacco corporation resear
c. j. cook date : may 8, 1995 subject : review of existing brainstorming ideas / 483 the major
uld be profitable to manufacture and sell. novel is defined as : of a new kind, or different fr
introduced ; act of innovating ; introduction of new things or methods. the products may incorp
taste or look. the first task of the product innovation group was to assemble, review and categ
es labeled appearance and taste / aroma. these categories are used for novel products t...' 
```

Source : Author

Upon completion of the encoding process, a pivotal step is to pinpoint the position of the answer within the encoded input. This involves leveraging crucial pieces of information: Firstly, the token_type_ids offer insights into the classification of tokens, distinguishing those relevant to the question and those associated with the document's words.

 This distinction is instrumental in comprehending the structural composition of the encoded input. Secondly, the presence of a special token at the input's inception, identified by tokenizer.cls_token_id, aids in recognizing the overall organization of the encoded input. Lastly, the word_ids play a crucial role in establishing a connection between the answer in the original words and its counterpart within the fully encoded input.

This linkage is paramount for determining the precise start and end positions of the answer within the encoded input. In summary, these elements collectively contribute to the effective localization of the answer in the encoded representation.

```python
def encode_dataset(examples, max_length=512):
    questions = examples["question"]
    words = examples["words"]
    boxes = examples["boxes"]
    answers = examples["answer"]

    # encode the batch of examples and initialize the start_positions and end_positions
    encoding = tokenizer(questions, words, boxes, max_length=max_length, padding="max_length", truncation=True)
    start_positions = []
    end_positions = []

    # loop through the examples in the batch
    for i in range(len(questions)):
```

```python
        cls_index = encoding["input_ids"][i].index(tokenizer.cls_token_id)

        # find the position of the answer in example's words
        words_example = [word.lower() for word in words[i]]
        answer = answers[i]
        match, word_idx_start, word_idx_end = subfinder(words_example,
answer.lower().split())

        if match:
            # if match is found, use `token_type_ids` to find where words start
in the encoding
            token_type_ids = encoding["token_type_ids"][i]
            token_start_index = 0
            while token_type_ids[token_start_index] != 1:
                token_start_index += 1

            token_end_index = len(encoding["input_ids"][i]) - 1
            while token_type_ids[token_end_index] != 1:
                token_end_index -= 1

            word_ids = encoding.word_ids(i)[token_start_index :
token_end_index + 1]
            start_position = cls_index
            end_position = cls_index

            # loop over word_ids and increase `token_start_index` until it
matches the answer position in words
            # once it matches, save the `token_start_index` as the
`start_position` of the answer in the encoding
```

```python
        for id in word_ids:
            if id == word_idx_start:
                start_position = token_start_index
            else:
                token_start_index += 1

        # similarly loop over `word_ids` starting from the end to find the
`end_position` of the answer
        for id in word_ids[::-1]:
            if id == word_idx_end:
                end_position = token_end_index
            else:
                token_end_index -= 1

        start_positions.append(start_position)
        end_positions.append(end_position)

    else:
        start_positions.append(cls_index)
        end_positions.append(cls_index)

encoding["image"] = examples["image"]
encoding["start_positions"] = start_positions
encoding["end_positions"] = end_positions

return encoding
```

With the completion of this preprocessing function, we are now poised to encode the entire dataset.

```python
encoded_train_dataset = dataset_with_ocr.map(
```

```
    encode_dataset, batched=True, batch_size=2,
remove_columns=dataset_with_ocr.column_names
)
encoded_test_dataset = dataset_with_ocr.map(
    encode_dataset, batched=True, batch_size=2,
remove_columns=dataset_with_ocr.column_names
)
```

Map: 100% ████████████████████████ 10/10 [00:02<00:00, 4.86 examples/s]
Map: 100% ████████████████████████ 10/10 [00:01<00:00, 5.41 examples/s]

Source : Author

encoded dataset should look like:

```
{'image': Sequence(feature=Sequence(feature=Sequenc
 'input_ids': Sequence(feature=Value(dtype='int32',
 'token_type_ids': Sequence(feature=Value(dtype='in
 'attention_mask': Sequence(feature=Value(dtype='in
 'bbox': Sequence(feature=Sequence(feature=Value(dt
 'start_positions': Value(dtype='int64', id=None),
 'end_positions': Value(dtype='int64', id=None)}
```

Source : Author

Another Way to process the data is to use UBIAI's Annotation tool. So head to UBIAI's website. Create new projects, upload your data and start annotating. You can find the annotation steps dedicated to the last article about fine tuning GPT models.

The idea here is to start by creating Entities for questions and their answers. Then create a relationship between the entities that way we specify that The answer to each question.

Once the entities are specified you can annotate your documents.

Once you finish annotating. Export the annotated data as JSON file. And it is ready for use.

# 6- Training:

Load the model using AutoModelForDocumentQuestionAnswering with the same checkpoint utilized in the preprocessing step. Specify training hyperparameters with TrainingArguments. Create a function to batch examples, where the DefaultDataCollator is suitable. Provide the training arguments, model, dataset, and data collator to the Trainer. Execute the train() method to commence fine-tuning your model.

```python
from transformers import AutoModelForDocumentQuestionAnswering


model =
AutoModelForDocumentQuestionAnswering.from_pretrained(model_check
point)
```

```
pytorch_model.bin: 100% |████████████████████████| 802M/802M [00:15<00:00, 47.0MB/s]
Some weights of LayoutLMv2ForQuestionAnswering were not initialized from the model checkpoint at microsoft/layout
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
```

Source : Author

This code snippet reflects the guidance you provided, using the output_dir parameter in TrainingArguments to specify where to save your model, and the push_to_hub parameter set to True if you want to share your model with the Hugging Face community. The output_dir will also serve as the name of the repository where your model checkpoint will be pushed if you choose to upload it.

```python
from transformers import TrainingArguments

# REPLACE THIS WITH YOUR REPO ID
repo_id = "MariaK/layoutlmv2-base-uncased_finetuned_docvqa"
```

```python
training_args = TrainingArguments(
    output_dir=repo_id,
    per_device_train_batch_size=4,
    num_train_epochs=20,
    save_steps=200,
    logging_steps=50,
    evaluation_strategy="steps",
    learning_rate=5e-5,
    save_total_limit=2,
    remove_unused_columns=False,
    push_to_hub=True,
)
```

Create a straightforward data collator to group examples into batches. Utilize the DefaultDataCollator from the Transformers library to achieve this. Once the data collator is instantiated, combine all components, including the model, training arguments, datasets, and tokenizer. Conclude by invoking the train() function to initiate the training process.

```python
from transformers import DefaultDataCollator

data_collator = DefaultDataCollator()
```

```python
from transformers import Trainer
import transformers

# Set your Hugging Face API token
transformers.hf_hub_token = "My_ID"

# Assuming 'model', 'training_args', 'data_collator', 'encoded_train_dataset',
# 'encoded_test_dataset', 'processor' are defined
```

```
trainer = Trainer(
    model=model,
    args=training_args,
    data_collator=data_collator,
    train_dataset=encoded_train_dataset,
    eval_dataset=encoded_test_dataset,
    tokenizer=processor,
)
trainer.train()
```

[60/60 01:06, Epoch 20/20]

| Step | Training Loss | Validation Loss |
|------|---------------|-----------------|
| 50   | 0.000000      | nan             |

```
TrainOutput(global_step=60, training_loss=0.0, metrics={'train_runtime':
'train_loss': 0.0, 'epoch': 20.0})
```

Source : Author

## 7- Inference:

Example given this image from the testing set. and  the question :
**Who is 'presiding' TRRF GENERAL SESSION**
We get as an answer :

```
example = dataset["test"][2]
question = example["query"]["en"]
image = example["image"]
print(question)
```

```
print(example["answers"])
```

Who is 'presiding' TRRF GENERAL SESSION (PART 1)?
['TRRF Vice President', 'lee a. waller']

| | |
|---|---|
| 11:14 to 11:39 a.m. | Coffee Break<br>Coffee will be served for men and women in the lobby adjacent to exhibit area. Please move into exhibit area. **(Exhibits Open)** |
| 11:39 a.m. | TRRF GENERAL SESSION (PART I)<br>Presiding: Lee A. Waller<br>TRRF Vice President |
| 11:39 to 11:44 a.m. | **"Introductory Remarks"**<br>Lee A. Waller, TRRF Vice President |
| 11:44 a.m. to 12:25 p.m. | **Individual Interviews with TRRF Public Board Members and Scientific Advisory Council Members**<br>Conducted by TRRF Treasurer Philip G. Kuehn to get answers which the public refrigerated warehousing industry is looking for. Plus questions from the floor. Dr. Emil M. Mrak, University of California, Chairman, TRRF Board; Sam R. Cecil, University of Georgia College of Agriculture; Dr. Stanley Charm, Tufts University School of Medicine; Dr. Robert H. Cotton, ITT Continental Baking Company; Dr. Owen Fennema, University of Wisconsin; Dr. Robert E. Hardenburg, USDA. |
| 12:25 to 12:58 p.m. | Questions and Answers |
| 12:58 to 4:00 p.m. | **Exhibits Open**<br>Capt. Jack Stoney Room |
| 2:00 to 5:00 p.m. | **TRRF Scientific Advisory Council Meeting**<br>Ballroom Foyer |

Source : Author

## 8- Conclusion:

In summary, fine-tuning LayoutLMv2 for Document Question Answering marks a significant step forward in advancing natural language processing

for document analysis. This article has covered key aspects of the fine-tuning process, emphasizing the model's dual proficiency in understanding text and layout. By integrating document layout insights, we enhance question answering systems, promising improved information retrieval across various applications.