

```
1 Start coding or generate with AI.
```

```
1 Start coding or generate with AI.
```

```
1 Prerequisites:
```

```
2 # pip install onnxruntime
3 # pip install onnxruntime-gpu
```

```
1 !pip install onnxruntime
```

→ Collecting onnxruntime

```
  Downloading onnxruntime-1.21.0-cp311-cp311-manylinux_2_27_x86_64.manylinux_2_28_x86_64.whl.metadata (4.5 kB)
Collecting coloredlogs (from onnxruntime)
  Downloading coloredlogs-15.0.1-py2.py3-none-any.whl.metadata (12 kB)
Requirement already satisfied: flatbuffers in /usr/local/lib/python3.11/dist-packages (from onnxruntime) (25.2.10)
Requirement already satisfied: numpy>=1.21.6 in /usr/local/lib/python3.11/dist-packages (from onnxruntime) (2.0.2)
Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from onnxruntime) (24.2)
Requirement already satisfied: protobuf in /usr/local/lib/python3.11/dist-packages (from onnxruntime) (5.29.3)
Requirement already satisfied: sympy in /usr/local/lib/python3.11/dist-packages (from onnxruntime) (1.13.1)
Collecting humanfriendly>=9.1 (from coloredlogs->onnxruntime)
  Downloading humanfriendly-10.0-py2.py3-none-any.whl.metadata (9.2 kB)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy->onnxruntime) (1.3.0)
  Downloading onnxruntime-1.21.0-cp311-cp311-manylinux_2_27_x86_64.manylinux_2_28_x86_64.whl (16.0 MB)
  16.0/16.0 MB 36.4 MB/s eta 0:00:00
```

```
  Downloading coloredlogs-15.0.1-py2.py3-none-any.whl (46 kB)
```

```
  46.0/46.0 kB 3.1 MB/s eta 0:00:00
```

```
  Downloading humanfriendly-10.0-py2.py3-none-any.whl (86 kB)
```

```
  86.8/86.8 kB 6.0 MB/s eta 0:00:00
```

```
Installing collected packages: humanfriendly, coloredlogs, onnxruntime
```

```
Successfully installed coloredlogs-15.0.1 humanfriendly-10.0 onnxruntime-1.21.0
```

```
1 !pip install onnxruntime-gpu
```

→ Collecting onnxruntime-gpu

```
  Downloading onnxruntime_gpu-1.21.0-cp311-cp311-manylinux_2_27_x86_64.manylinux_2_28_x86_64.whl.metadata (4.8 kB)
Requirement already satisfied: coloredlogs in /usr/local/lib/python3.11/dist-packages (from onnxruntime-gpu) (15.0.1)
Requirement already satisfied: flatbuffers in /usr/local/lib/python3.11/dist-packages (from onnxruntime-gpu) (25.2.10)
Requirement already satisfied: numpy>=1.21.6 in /usr/local/lib/python3.11/dist-packages (from onnxruntime-gpu) (2.0.2)
Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from onnxruntime-gpu) (24.2)
Requirement already satisfied: protobuf in /usr/local/lib/python3.11/dist-packages (from onnxruntime-gpu) (5.29.3)
Requirement already satisfied: sympy in /usr/local/lib/python3.11/dist-packages (from onnxruntime-gpu) (1.13.1)
Requirement already satisfied: humanfriendly>=9.1 in /usr/local/lib/python3.11/dist-packages (from coloredlogs->onnxruntime-gpu) (10.0)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy->onnxruntime-gpu) (1.3.0)
  Downloading onnxruntime_gpu-1.21.0-cp311-cp311-manylinux_2_27_x86_64.manylinux_2_28_x86_64.whl (280.8 MB)
  280.8/280.8 MB 3.2 MB/s eta 0:00:00
```

```
Installing collected packages: onnxruntime-gpu
```

```
Successfully installed onnxruntime-gpu-1.21.0
```

```
1 !pip install docling_core
```

→ Collecting docling_core

```
  Downloading docling_core-2.23.3-py3-none-any.whl.metadata (5.8 kB)
Collecting jsonref<2.0.0,>=1.1.0 (from docling_core)
  Downloading jsonref-1.1.0-py3-none-any.whl.metadata (2.7 kB)
Requirement already satisfied: jsonschema<5.0.0,>=4.16.0 in /usr/local/lib/python3.11/dist-packages (from docling_core) (4.23.0)
Collecting latex2mathml<4.0.0,>=3.77.0 (from docling_core)
  Downloading latex2mathml-3.77.0-py3-none-any.whl.metadata (14 kB)
Requirement already satisfied: pandas<3.0.0,>=2.1.4 in /usr/local/lib/python3.11/dist-packages (from docling_core) (2.2.2)
Requirement already satisfied: pillow<12.0.0,>=10.0.0 in /usr/local/lib/python3.11/dist-packages (from docling_core) (11.1.0)
Requirement already satisfied: pydantic!=2.10.0,!=2.10.1,!=2.10.2,<3.0.0,>=2.6.0 in /usr/local/lib/python3.11/dist-packages (from docling_core)
Requirement already satisfied: pyyaml<7.0.0,>=5.1 in /usr/local/lib/python3.11/dist-packages (from docling_core) (6.0.2)
Requirement already satisfied: tabulate<0.10.0,>=0.9.0 in /usr/local/lib/python3.11/dist-packages (from docling_core) (0.9.0)
Collecting typer<0.13.0,>=0.12.5 (from docling_core)
  Downloading typer-0.12.5-py3-none-any.whl.metadata (15 kB)
Requirement already satisfied: typing-extensions<5.0.0,>=4.12.2 in /usr/local/lib/python3.11/dist-packages (from docling_core) (4.12.2)
Requirement already satisfied: attrs=>22.2.0 in /usr/local/lib/python3.11/dist-packages (from jsonschema<5.0.0,>=4.16.0->docling_core) (22.2.0)
Requirement already satisfied: jsonschema-specifications==2023.03.6 in /usr/local/lib/python3.11/dist-packages (from jsonschema<5.0.0,>=4.16.0->docling_core)
Requirement already satisfied: referencing>=0.28.4 in /usr/local/lib/python3.11/dist-packages (from jsonschema<5.0.0,>=4.16.0->docling_core)
Requirement already satisfied: rpds-py>=0.7.1 in /usr/local/lib/python3.11/dist-packages (from jsonschema<5.0.0,>=4.16.0->docling_core)
Requirement already satisfied: numpy>=1.23.2 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0.0,>=2.1.4->docling_core) (2.0.2)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0.0,>=2.1.4->docling_core)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0.0,>=2.1.4->docling_core) (2025.1)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0.0,>=2.1.4->docling_core) (2025.1)
Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.11/dist-packages (from pydantic!=2.10.0,!=2.10.1,!=2.10.2)
Requirement already satisfied: pydantic-core==2.27.2 in /usr/local/lib/python3.11/dist-packages (from pydantic!=2.10.0,!=2.10.1,!=2.10.2)
Requirement already satisfied: click>=8.0.0 in /usr/local/lib/python3.11/dist-packages (from typer<0.13.0,>=0.12.5->docling_core) (8.1.8)
Requirement already satisfied: shellingham>=1.3.0 in /usr/local/lib/python3.11/dist-packages (from typer<0.13.0,>=0.12.5->docling_core)
Requirement already satisfied: rich>=10.11.0 in /usr/local/lib/python3.11/dist-packages (from typer<0.13.0,>=0.12.5->docling_core) (13.9.0)
```

```

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas<3.0.0,>=2.1.4->d
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.11/dist-packages (from rich>=10.11.0->typer<0.13.0,>=0.12
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from rich>=10.11.0->typer<0.13.0,>=0.
Requirement already satisfied: mdurl~=0.1 in /usr/local/lib/python3.11/dist-packages (from markdown-it-py>=2.2.0->rich>=10.11.0->typer<0
Downloading docling_core-2.23.3-py3-none-any.whl (115 kB)
  115.9/115.9 kB 5.2 MB/s eta 0:00:00
Downloading jsonref-1.1.0-py3-none-any.whl (9.4 kB)
Downloading latex2mathml-3.77.0-py3-none-any.whl (73 kB)
  73.7/73.7 kB 4.9 MB/s eta 0:00:00
Downloading typer-0.12.5-py3-none-any.whl (47 kB)
  47.3/47.3 kB 2.6 MB/s eta 0:00:00
Installing collected packages: latex2mathml, jsonref, typer, docling_core
Attempting uninstall: typer
  Found existing installation: typer 0.15.2
  Uninstalling typer-0.15.2:
    Successfully uninstalled typer-0.15.2
Successfully installed docling_core-2.23.3 jsonref-1.1.0 latex2mathml-3.77.0 typer-0.12.5

```

```

1
2 from transformers import AutoConfig, AutoProcessor
3 from transformers.image_utils import load_image
4 import onnxruntime
5 import numpy as np
6 import os
7 from docling_core.types.doc import DoclingDocument
8 from docling_core.types.doc.document import DocTagsDocument
9
10 os.environ["OMP_NUM_THREADS"] = "1"
11 # cuda
12 os.environ["ORT_CUDA_USE_MAX_WORKSPACE"] = "1"
13
14 # 1. Load models
15 ## Load config and processor
16 model_id = "ds4sd/SmolDocling-256M-preview"
17 config = AutoConfig.from_pretrained(model_id)
18 processor = AutoProcessor.from_pretrained(model_id)
19
20 ## Load sessions
21 # !wget https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/vision_encoder.onnx
22 # !wget https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/embed_tokens.onnx
23 # !wget https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/decoder_model_merged.onnx
24 # cpu
25 vision_session = onnxruntime.InferenceSession("vision_encoder.onnx")
26 embed_session = onnxruntime.InferenceSession("embed_tokens.onnx")
27 decoder_session = onnxruntime.InferenceSession("decoder_model_merged.onnx")
28
29 # cuda
30 vision_session = onnxruntime.InferenceSession("vision_encoder.onnx", providers=["CUDAExecutionProvider"])
31 embed_session = onnxruntime.InferenceSession("embed_tokens.onnx", providers=["CUDAExecutionProvider"])
32 decoder_session = onnxruntime.InferenceSession("decoder_model_merged.onnx", providers=["CUDAExecutionProvider"])
33
34 ## Set config values
35 num_key_value_heads = config.text_config.num_key_value_heads
36 head_dim = config.text_config.head_dim
37 num_hidden_layers = config.text_config.num_hidden_layers
38 eos_token_id = config.text_config.eos_token_id
39 image_token_id = config.image_token_id
40 end_of_utterance_id = processor.tokenizer.convert_tokens_to_ids("<end_of_utterance>")
41
42 # 2. Prepare inputs
43 ## Create input messages
44 messages = [
45     {
46         "role": "user",
47         "content": [
48             {"type": "image"},
49             {"type": "text", "text": "Convert this page to docling."}
50         ]
51     },
52 ]
53
54 ## Load image and apply processor
55 image = load_image("https://ibm.biz/docling-page-with-table")
56 prompt = processor.apply_chat_template(messages, add_generation_prompt=True)
57 inputs = processor(text=prompt, images=[image], return_tensors="np")
58
59 ## Prepare decoder inputs

```

```

60 batch_size = inputs['input_ids'].shape[0]
61 past_key_values = {
62     f'past_key_values.{layer}.\{kv\}': np.zeros([batch_size, num_key_value_heads, 0, head_dim], dtype=np.float32)
63     for layer in range(num_hidden_layers)
64     for kv in ('key', 'value')
65 }
66 image_features = None
67 input_ids = inputs['input_ids']
68 attention_mask = inputs['attention_mask']
69 position_ids = np.cumsum(inputs['attention_mask'], axis=-1)
70
71
72 # 3. Generation loop
73 max_new_tokens = 8192
74 generated_tokens = np.array([[]], dtype=np.int64)
75 for i in range(max_new_tokens):
76     inputs_embeds = embed_session.run(None, {'input_ids': input_ids})[0]
77
78     if image_features is None:
79         ## Only compute vision features if not already computed
80         image_features = vision_session.run(
81             ['image_features'], # List of output names or indices
82             {
83                 'pixel_values': inputs['pixel_values'],
84                 'pixel_attention_mask': inputs['pixel_attention_mask'].astype(np.bool_),
85             },
86         )[0]
87
88     ## Merge text and vision embeddings
89     inputs_embeds[inputs['input_ids'] == image_token_id] = image_features.reshape(-1, image_features.shape[-1])
90
91     logits, *present_key_values = decoder_session.run(None, dict(
92         inputs_embeds=inputs_embeds,
93         attention_mask=attention_mask,
94         position_ids=position_ids,
95         **past_key_values,
96     ))
97
98     ## Update values for next generation loop
99     input_ids = logits[:, -1].argmax(-1, keepdims=True)
100    attention_mask = np.ones_like(input_ids)
101    position_ids = position_ids[:, -1:] + 1
102    for j, key in enumerate(past_key_values):
103        past_key_values[key] = present_key_values[j]
104
105    generated_tokens = np.concatenate([generated_tokens, input_ids], axis=-1)
106    if (input_ids == eos_token_id).all() or (input_ids == end_of_utterance_id).all():
107        break # Stop predicting
108
109 doctags = processor.batch_decode(
110     generated_tokens,
111     skip_special_tokens=False,
112 )[0].lstrip()
113
114 print(doctags)
115
116 doctags_doc = DocTagsDocument.from_doctags_and_image_pairs([doctags], [image])
117 print(doctags)
118 # create a docling document
119 doc = DoclingDocument(name="Document")
120 doc.load_from_doctags(doctags_doc)
121
122 print(doc.export_to_markdown())

```

```

↳ /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: 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 secre
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(
config.json: 100%                                         3.90k/3.90k [00:00<00:00, 314kB/s]

processor_config.json: 100%                                     68.0/68.0 [00:00<00:00, 4.17kB/s]

chat_template.json: 100%                                     430/430 [00:00<00:00, 31.6kB/s]

preprocessor_config.json: 100%                                486/486 [00:00<00:00, 28.7kB/s]

tokenizer_config.json: 100%                                27.4k/27.4k [00:00<00:00, 1.13MB/s]

vocab.json: 100%                                         801k/801k [00:00<00:00, 9.12MB/s]

merges.txt: 100%                                         466k/466k [00:00<00:00, 6.28MB/s]

tokenizer.json: 100%                                     3.55M/3.55M [00:00<00:00, 9.68MB/s]

added_tokens.json: 100%                                 3.67k/3.67k [00:00<00:00, 314kB/s]

special_tokens_map.json: 100%                            1.07k/1.07k [00:00<00:00, 117kB/s]

-----
NoSuchFile                                                 Traceback (most recent call last)
<ipython-input-5-251b49be7d7b> in <cell line: 0>()
      27
      28 # cuda
--> 29 vision_session = onnxruntime.InferenceSession("vision_encoder.onnx", providers=["CUDAExecutionProvider"])
  30 embed_session = onnxruntime.InferenceSession("embed_tokens.onnx", providers=["CUDAExecutionProvider"])
  31 decoder_session = onnxruntime.InferenceSession("decoder_model_merged.onnx", providers=["CUDAExecutionProvider"])

  1 frames
/usr/local/lib/python3.11/dist-packages/onnxruntime/capi/onnxruntime_inference_collection.py in _create_inference_session(self, providers, provider_options, disabled_optimizers)
  528
  529     if self._model_path:
--> 530         sess = C.InferenceSession(session_options, self._model_path, True, self._read_config_from_model)
  531     else:
  532         sess = C.InferenceSession(session_options, self._model_bytes, False, self._read_config_from_model)

NoSuchFile: [ONNXRuntimeError] : 3 : NO SUCHFILE : Load model from vision_encoder.onnx failed:Load model vision_encoder.onnx failed.
File doesn't exist

```

Next steps: [Explain error](#)

```

1 from transformers import AutoConfig, AutoProcessor
2 from transformers.image_utils import load_image
3 import onnxruntime
4 import numpy as np
5 import os
6 from doceling_core.types.doc import DocelingDocument
7 from doceling_core.types.doc.document import DocTagsDocument
8
9 os.environ["OMP_NUM_THREADS"] = "1"
10 # cuda
11 os.environ["ORT_CUDA_USE_MAX_WORKSPACE"] = "1"
12
13 # 1. Load models
14 ## Load config and processor
15 model_id = "ds4sd/SmolDocling-256M-preview"
16 config = AutoConfig.from_pretrained(model_id)
17 processor = AutoProcessor.from_pretrained(model_id)
18
19 ## Load sessions
20 # Download the ONNX models if they are not present
21 !wget https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/vision_encoder.onnx
22 !wget https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/embed_tokens.onnx
23 !wget https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/decoder_model_merged.onnx
24 # cpu
25 # vision_session = onnxruntime.InferenceSession("vision_encoder.onnx")
26 # embed_session = onnxruntime.InferenceSession("embed_tokens.onnx")
27 # decoder_session = onnxruntime.InferenceSession("decoder_model_merged.onnx")
28
29 # cuda

```

```

30 vision_session = onnxruntime.InferenceSession("vision_encoder.onnx", providers=["CUDAExecutionProvider"])
31 embed_session = onnxruntime.InferenceSession("embed_tokens.onnx", providers=["CUDAExecutionProvider"])
32 decoder_session = onnxruntime.InferenceSession("decoder_model_merged.onnx", providers=["CUDAExecutionProvider"])
33
34 # ... rest of the code ...

✉ /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: 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.
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(
--2025-03-23 00:18:28-- https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/vision\_encoder.onnx
Resolving huggingface.co (huggingface.co)... 18.164.174.17, 18.164.174.118, 18.164.174.55, ...
Connecting to huggingface.co (huggingface.co)|18.164.174.17|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://cdn-lfs-us-1.hf.co/repos/b8/6e/b86e88e6db456aab3c09c424eee0c6e3c444c4122486d3f9ec088b13c389b0f9/9c97aaaf0543632c77893
--2025-03-23 00:18:28-- https://cdn-lfs-us-1.hf.co/repos/b8/6e/b86e88e6db456aab3c09c424eee0c6e3c444c4122486d3f9ec088b13c389b0f9/9c97aaaf0543632c77893
Resolving cdn-lfs-us-1.hf.co (cdn-lfs-us-1.hf.co)... 3.167.192.98, 3.167.192.6, 3.167.192.118, ...
Connecting to cdn-lfs-us-1.hf.co (cdn-lfs-us-1.hf.co)|3.167.192.98|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 374316454 (357M) [binary/octet-stream]
Saving to: 'vision_encoder.onnx'

vision_encoder.onnx 100%[=====] 356.98M 149MB/s in 2.4s

2025-03-23 00:18:31 (149 MB/s) - 'vision_encoder.onnx' saved [374316454/374316454]

--2025-03-23 00:18:31-- https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/embed\_tokens.onnx
Resolving huggingface.co (huggingface.co)... 18.164.174.17, 18.164.174.118, 18.164.174.55, ...
Connecting to huggingface.co (huggingface.co)|18.164.174.17|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://cdn-lfs-us-1.hf.co/repos/b8/6e/b86e88e6db456aab3c09c424eee0c6e3c444c4122486d3f9ec088b13c389b0f9/816939a02a48a540330655
--2025-03-23 00:18:31-- https://cdn-lfs-us-1.hf.co/repos/b8/6e/b86e88e6db456aab3c09c424eee0c6e3c444c4122486d3f9ec088b13c389b0f9/816939a02a48a540330655
Resolving cdn-lfs-us-1.hf.co (cdn-lfs-us-1.hf.co)... 3.167.192.98, 3.167.192.6, 3.167.192.118, ...
Connecting to cdn-lfs-us-1.hf.co (cdn-lfs-us-1.hf.co)|3.167.192.98|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 113541419 (108M) [binary/octet-stream]
Saving to: 'embed_tokens.onnx'

embed_tokens.onnx 100%[=====] 108.28M 121MB/s in 0.9s

2025-03-23 00:18:32 (121 MB/s) - 'embed_tokens.onnx' saved [113541419/113541419]

--2025-03-23 00:18:32-- https://huggingface.co/ds4sd/SmolDocling-256M-preview/resolve/main/onnx/decoder\_model\_merged.onnx
Resolving huggingface.co (huggingface.co)... 18.164.174.17, 18.164.174.118, 18.164.174.55, ...
Connecting to huggingface.co (huggingface.co)|18.164.174.17|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://cdn-lfs-us-1.hf.co/repos/b8/6e/b86e88e6db456aab3c09c424eee0c6e3c444c4122486d3f9ec088b13c389b0f9/c75155e5e121b9d89bf543
--2025-03-23 00:18:32-- https://cdn-lfs-us-1.hf.co/repos/b8/6e/b86e88e6db456aab3c09c424eee0c6e3c444c4122486d3f9ec088b13c389b0f9/c75155e5e121b9d89bf543
Resolving cdn-lfs-us-1.hf.co (cdn-lfs-us-1.hf.co)... 3.167.192.98, 3.167.192.6, 3.167.192.118, ...
Connecting to cdn-lfs-us-1.hf.co (cdn-lfs-us-1.hf.co)|3.167.192.98|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 540972236 (516M) [binary/octet-stream]
Saving to: 'decoder_model_merged.onnx'

decoder_model_merge 100%[=====] 515.91M 165MB/s in 3.2s

2025-03-23 00:18:35 (160 MB/s) - 'decoder_model_merged.onnx' saved [540972236/540972236]

```

```

1 ## Set config values
2 num_key_value_heads = config.text_config.num_key_value_heads
3 head_dim = config.text_config.head_dim
4 num_hidden_layers = config.text_config.num_hidden_layers
5 eos_token_id = config.text_config.eos_token_id
6 image_token_id = config.image_token_id
7 end_of_utterance_id = processor.tokenizer.convert_tokens_to_ids("<end_of_utterance>")
8
9 # 2. Prepare inputs
10 ## Create input messages
11 messages = [
12     {
13         "role": "user",
14         "content": [
15             {"type": "image"},
16             {"type": "text", "text": "Convert this page to docling."}
17         ]
18     },
19 ]

```

```

21 ## Load image and apply processor
22 image = load_image("https://ibm.biz/docling-page-with-table")
23 prompt = processor.apply_chat_template(messages, add_generation_prompt=True)
24 inputs = processor(text=prompt, images=[image], return_tensors="np")
25
26 ## Prepare decoder inputs
27 batch_size = inputs['input_ids'].shape[0]
28 past_key_values = {
29     f'past_key_values.{layer}.{kv}': np.zeros([batch_size, num_key_value_heads, 0, head_dim], dtype=np.float32)
30     for layer in range(num_hidden_layers)
31     for kv in ('key', 'value')
32 }
33 image_features = None
34 input_ids = inputs['input_ids']
35 attention_mask = inputs['attention_mask']
36 position_ids = np.cumsum(inputs['attention_mask'], axis=-1)
37
38
39 # 3. Generation loop
40 max_new_tokens = 8192
41 generated_tokens = np.array([[[]], dtype=np.int64)
42 for i in range(max_new_tokens):
43     inputs_embeds = embed_session.run(None, {'input_ids': input_ids})[0]
44
45 if image_features is None:
46     ## Only compute vision features if not already computed
47     image_features = vision_session.run(
48         ['image_features'], # List of output names or indices
49         {
50             'pixel_values': inputs['pixel_values'],
51             'pixel_attention_mask': inputs['pixel_attention_mask'].astype(np.bool_),
52         }
53     )[0]
54
55     ## Merge text and vision embeddings
56     inputs_embeds[input_ids == image_token_id] = image_features.reshape(-1, image_features.shape[-1])
57
58 logits, *present_key_values = decoder_session.run(None, dict(
59     inputs_embeds=inputs_embeds,
60     attention_mask=attention_mask,
61     position_ids=position_ids,
62     **past_key_values,
63 ))
64
65 ## Update values for next generation loop
66 input_ids = logits[:, -1].argmax(-1, keepdims=True)
67 attention_mask = np.ones_like(input_ids)
68 position_ids = position_ids[:, -1:] + 1
69 for j, key in enumerate(past_key_values):
70     past_key_values[key] = present_key_values[j]
71
72 generated_tokens = np.concatenate([generated_tokens, input_ids], axis=-1)
73 if (input_ids == eos_token_id).all() or (input_ids == end_of_utterance_id).all():
74     break # Stop predicting
75
76 doctags = processor.batch_decode(
77     generated_tokens,
78     skip_special_tokens=False,
79 )[0].lstrip()
80
81 print(doctags)
82
83 doctags_doc = DocTagsDocument.from_doctags_and_image_pairs([doctags], [image])
84 print(doctags)
85 # create a docling document
86 doc = DoclingDocument(name="Document")
87 doc.load_from_doctags(doctags_doc)
88
89 print(doc.export_to_markdown())

```

Truncation was not explicitly activated but `max_length` is provided a specific value, please use `truncation=True` to explicitly truncate the input sequence. The following table summarizes the results of the experiments. Inference timing results for all experiments were obtained from the Tabular Inference API. We have chosen the PubTabNet data set to perform HPO, since it includes a highly diverse set of documents.

```
<section_header_level_1><loc_57><loc_343><loc_209><loc_351>5.2 Quantitative Results</section_header_level_1>
<text><loc_57><loc_356><loc_457><loc_427>We picked the model parameter configuration that produced the best prediction quality (enc=6, d
<text><loc_57><loc_429><loc_457><loc_464>Additionally, the results show that OTSL has an advantage over HTML when applied on a bigger da
</doctag><end_of_utterance>
<doctag><page_header><loc_127><loc_27><loc_420><loc_34>Optimized Table Tokenization for Table Structure Recognition</page_header>
<page_header><loc_451><loc_27><loc_457><loc_34>9</page_header>
<text><loc_57><loc_46><loc_457><loc_71>order to compute the TED score. Inference timing results for all experiments were obtained from t
<section_header_level_1><loc_57><loc_84><loc_271><loc_92>5.1 Hyper Parameter Optimization</section_header_level_1>
<text><loc_57><loc_97><loc_457><loc_149>We have chosen the PubTabNet data set to perform HPO, since it includes a highly diverse set of
<otsl><loc_62><loc_212><loc_451><loc_315><ched>#<ched>#<ched>#<ched>TEDs<lcel><ched>mAP<ched>Inference time (secs)<n1><rhd>enc-layers<
<section_header_level_1><loc_57><loc_343><loc_209><loc_351>5.2 Quantitative Results</section_header_level_1>
<text><loc_57><loc_356><loc_457><loc_427>We picked the model parameter configuration that produced the best prediction quality (enc=6, d
<text><loc_57><loc_429><loc_457><loc_464>Additionally, the results show that OTSL has an advantage over HTML when applied on a bigger da
</doctag><end_of_utterance>
Optimized Table Tokenization for Table Structure Recognition
```

9

order to compute the TED score. Inference timing results for all experiments were obtained from the same machine on a single core with A

5.1 Hyper Parameter Optimization

We have chosen the PubTabNet data set to perform HPO, since it includes a highly diverse set of tables. Also we report TED scores separa

Table 1. HPO performed in OTSL and HTML representation on the same transformer-based TableFormer [9] architecture, trained only on PubTa

#	#	#	TEDs	TEDs	mAP	Inference time (secs)	
enc-layers	dec-layers	Language	simple	complex	all	(0.75)	
6	6	OTSL	0.965	0.934	0.955	0.88	2.73
		HTML	0.969	0.927	0.955	0.857	5.39
4	4	OTSL	0.938	0.904	0.927	0.853	1.97
		HTML	0.952	0.909	0.938	0.843	3.77
2	4	OTSL	0.923	0.897	0.915	0.859	1.91
		HTML	0.945	0.901	0.931	0.834	3.81
4	2	OTSL	0.952	0.92	0.942	0.857	1.22
		HTML	0.944	0.903	0.931	0.824	2

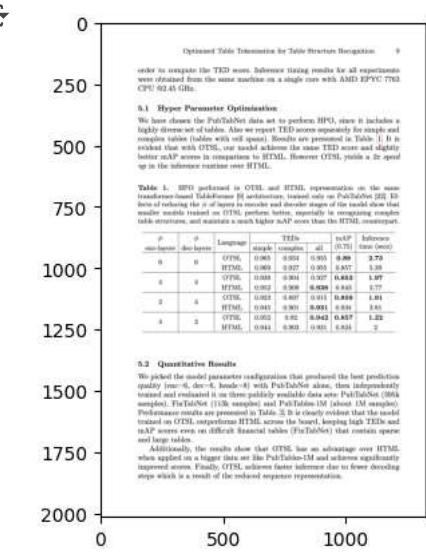
5.2 Quantitative Results

We picked the model parameter configuration that produced the best prediction quality (enc=6, dec=6, heads=8) with PubTabNet alone, then

Additionally, the results show that OTSL has an advantage over HTML when applied on a bigger data set like PubTables-1M and achieves sig

1 Start coding or generate with AI.

```
1 import matplotlib.pyplot as plt
2 import matplotlib.image as mpimg
3
4 # مسار الصورة
5 image_path = "/content/page_with_table.png"
6
7 # قراءة الصورة
8 img = mpimg.imread(image_path)
9
10 # عرض الصورة
11 imgplot = plt.imshow(img)
12 plt.show()
```



1 Start coding or generate with AI.

```

1 # Prerequisites:
2 # pip install torch
3 # pip install docling_core
4 # pip install transformers
5
6 import torch
7 from docling_core.types.doc import DoclingDocument
8 from docling_core.types.doc.document import DocTagsDocument
9 from transformers import AutoProcessor, AutoModelForVision2Seq
10 from transformers.image_utils import load_image
11 from pathlib import Path
12
13 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
14
15 # Load images
16 image = load_image("https://upload.wikimedia.org/wikipedia/commons/7/76/GazetteerdeFrance.jpg")
17
18 # Initialize processor and model
19 processor = AutoProcessor.from_pretrained("ds4sd/SmolDocling-256M-preview")
20 model = AutoModelForVision2Seq.from_pretrained(
21     "ds4sd/SmolDocling-256M-preview",
22     torch_dtype=torch.bfloat16,
23     _attn_implementation="flash_attention_2" if DEVICE == "cuda" else "eager",
24 ).to(DEVICE)
25
26 # Create input messages
27 messages = [
28     {
29         "role": "user",
30         "content": [
31             {"type": "image"},
32             {"type": "text", "text": "Convert this page to docling."}
33         ]
34     },
35 ]
36
37 # Prepare inputs

```

```
38 prompt = processor.apply_chat_template(messages, add_generation_prompt=True)
39 inputs = processor(text=prompt, images=[image], return_tensors="pt")
40 inputs = inputs.to(DEVICE)
41
42 # Generate outputs
43 generated_ids = model.generate(**inputs, max_new_tokens=8192)
44 prompt_length = inputs.input_ids.shape[1]
45 trimmed_generated_ids = generated_ids[:, prompt_length:]
46 doctags = processor.batch_decode(
47     trimmed_generated_ids,
48     skip_special_tokens=False,
49 )[0].lstrip()
50
51 # Populate document
```