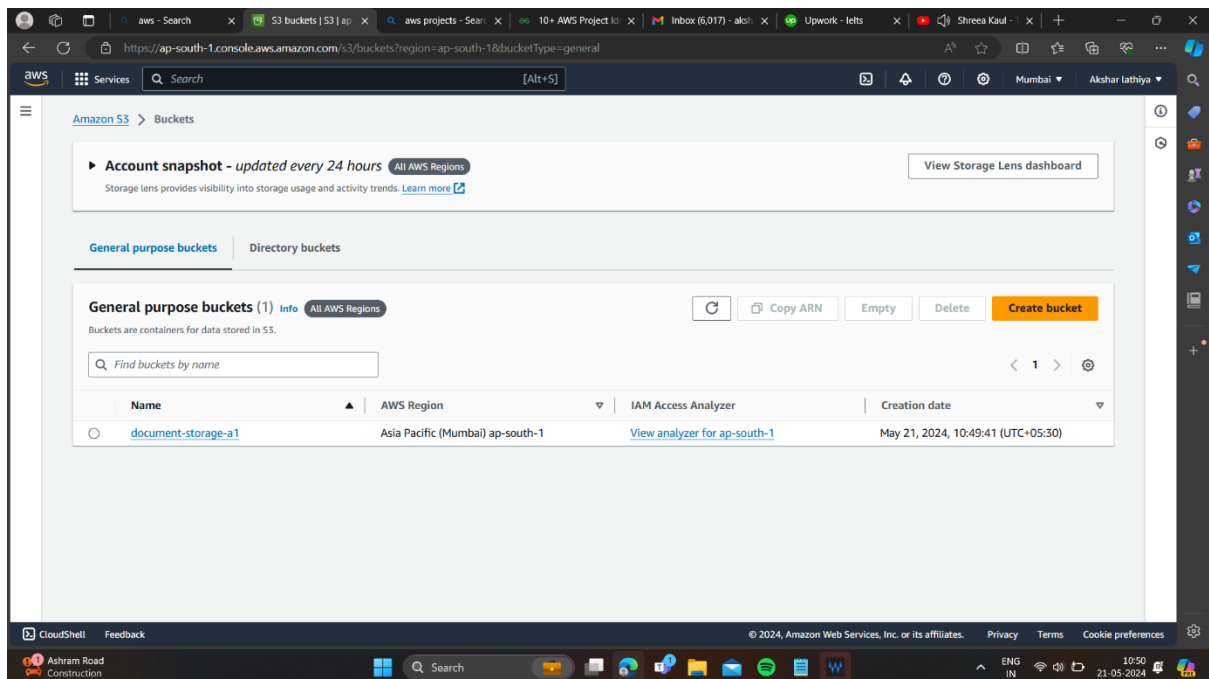
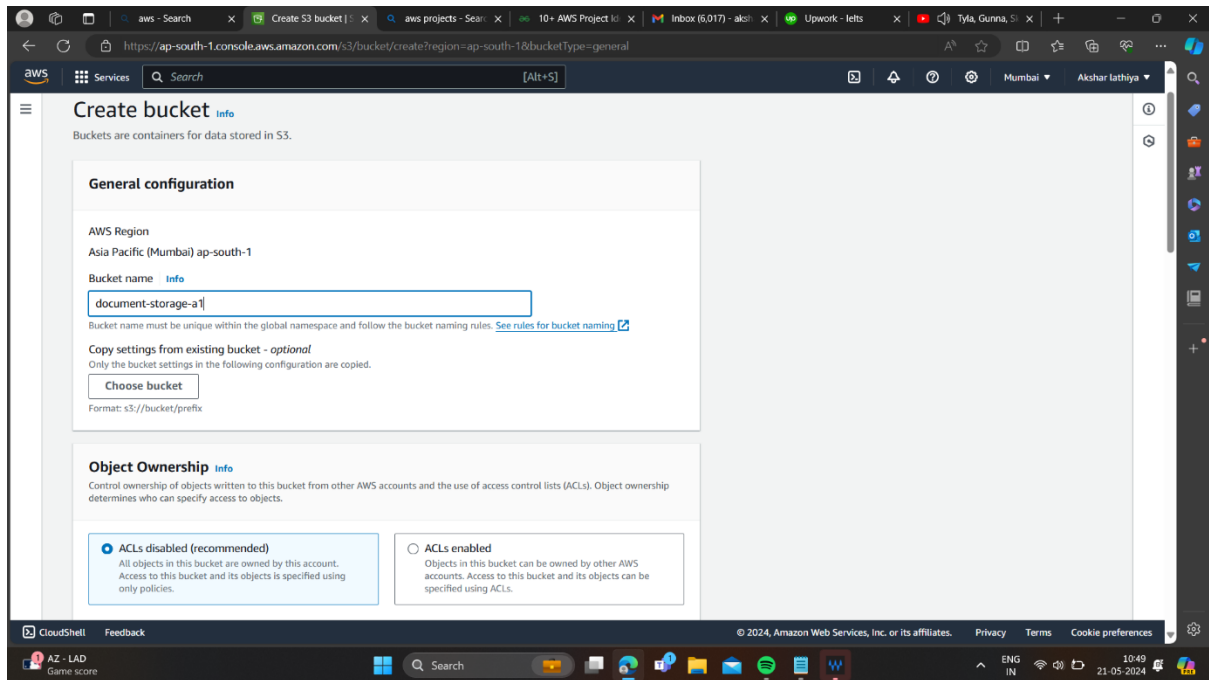


PDF conversion of any document uploaded in S3 bucket

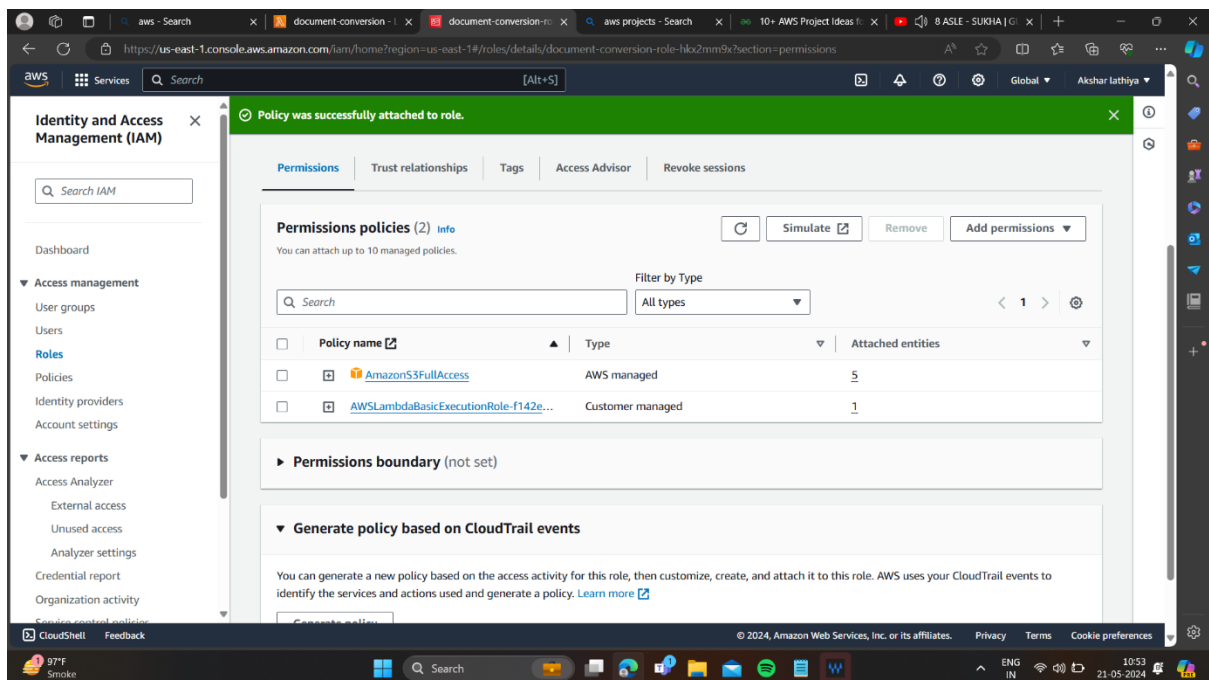
- **S3 Bucket setup:-**

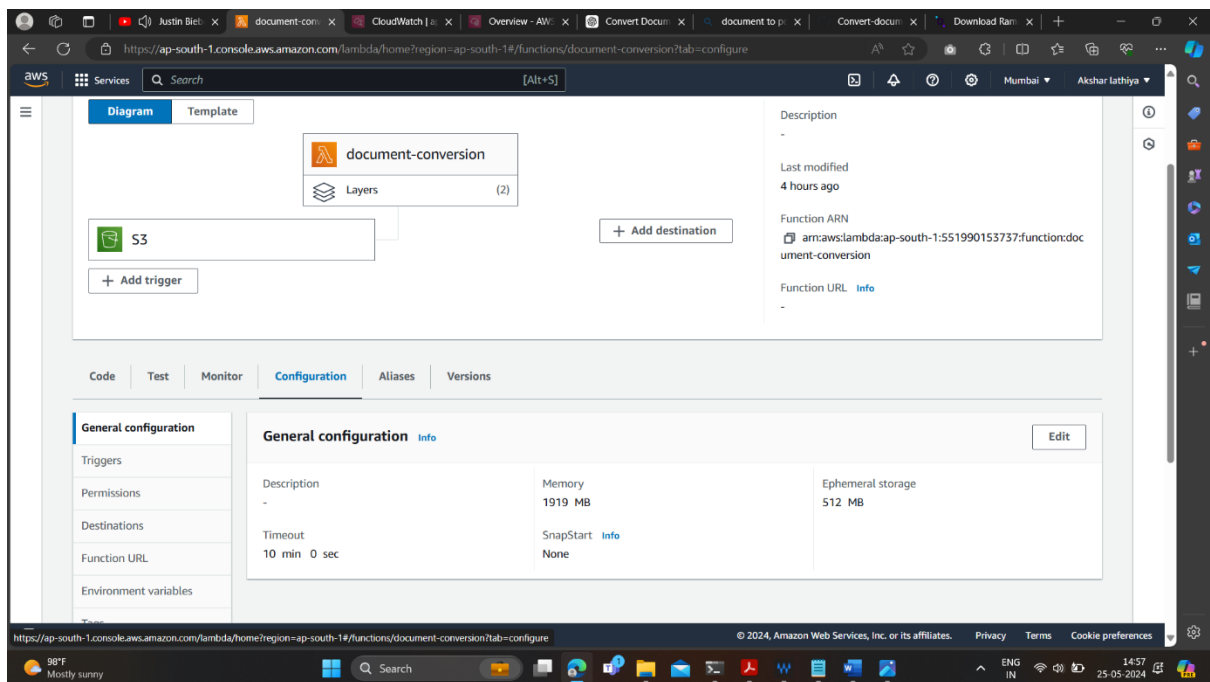
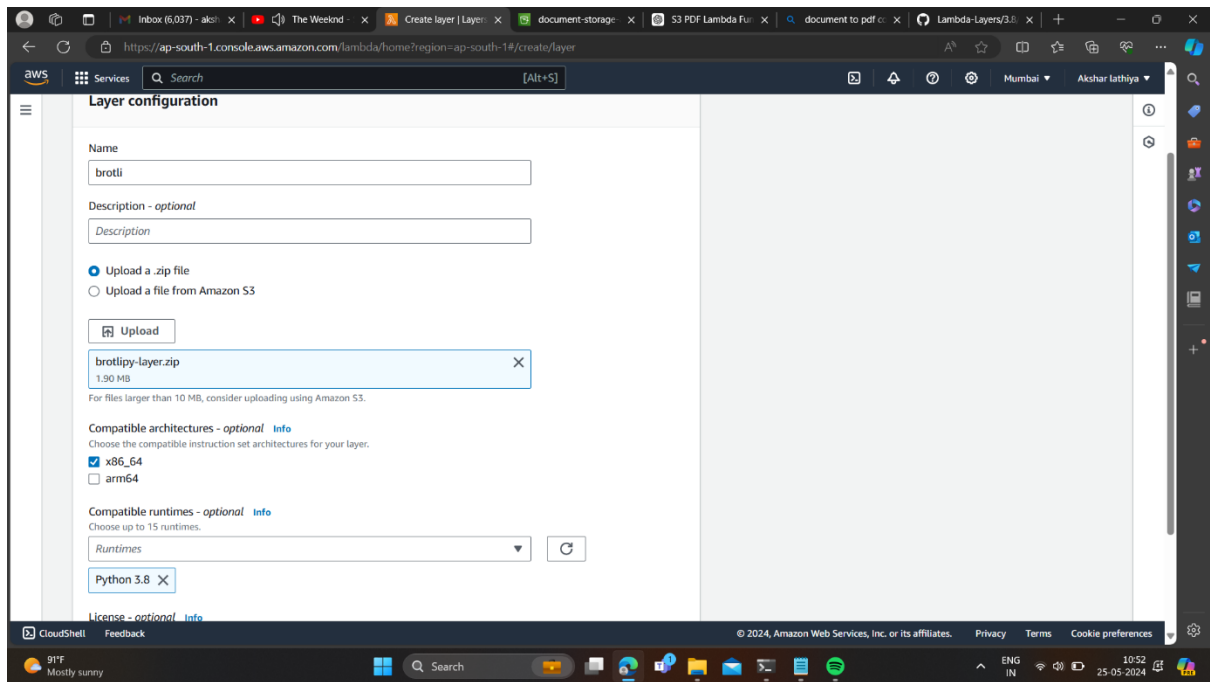
- Firstly , setup two buckets for source and destination.
- Set both buckets as per default buckets configurations and make sure that both buckets are in same region.



- **Lambda Function setup:-**

- Now configure lambda function and add python 3.8 as runtime .
- The function contains layers , so add layer in section and configure by uploading brotli-layer.zip file to it.
- Layers will add dependency required while running lambda function .
- Add layer created in function and give aws arn of layer created .
- Add aws arn of libreoffice of the Mumbai region.
- Now insert python code for file conversion.
- Configure memory and timeout for function.
- Configure IAM role of lambda and give full S3 access to upload and access files in bucket.





```

import subprocess
import tarfile
import sys
from io import BytesIO
import os
import time
import boto3
import urllib.parse
import json

sys.path.append("/opt/brotli")
import brotli

ACCESS_KEY = "AKIAYBBJNKIEYBERKW4H"
SECRET_ACCESS_KEY = "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"

def extract_libre_office():
    buffer = BytesIO()
    with open('/opt/lo.tar.br', mode='rb') as fout:
        file = fout.read()
        buffer.write(brotli.decompress(file))
        buffer.seek(0)
        with tarfile.open(fileobj=buffer) as tar:
            tar.extractall('/tmp')

def lambda_handler(event, context):
    if os.path.exists("/tmp/instdir/program/soffice.bin"):
        pass
    else :
        # load libre
        extract_libre_office()

    # Get Trigger event
    if 'Records' in event.keys():
        input_bucket = event['Records'][0]['s3']['bucket']['name']
        output_bucket = "document-pdf-a2"
        input_key = urllib.parse.unquote_plus(event['Records'][0]['s3']['object']['key'])
        s3 = boto3.client('s3', aws_access_key_id=ACCESS_KEY, aws_secret_access_key=SECRET_ACCESS_KEY)
    else :
        return 'test finished'

    # get S3 Object
    file_path = '/tmp/'+input_key
    s3.download_file(input_bucket, input_key, file_path)

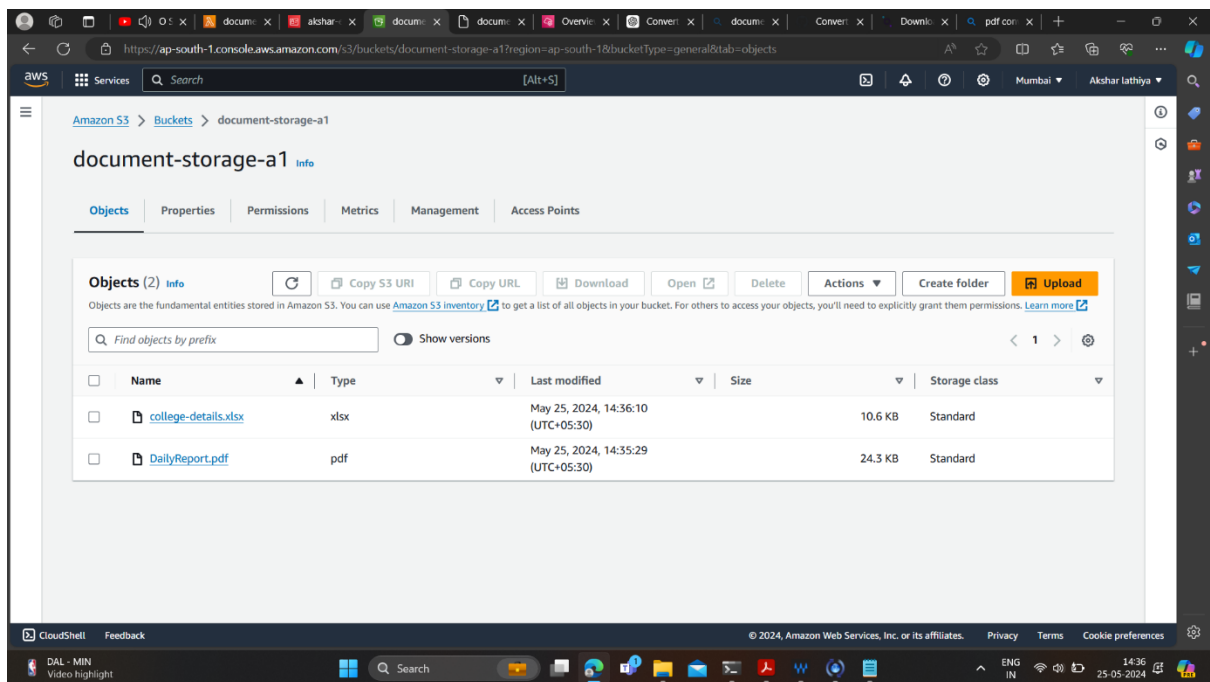
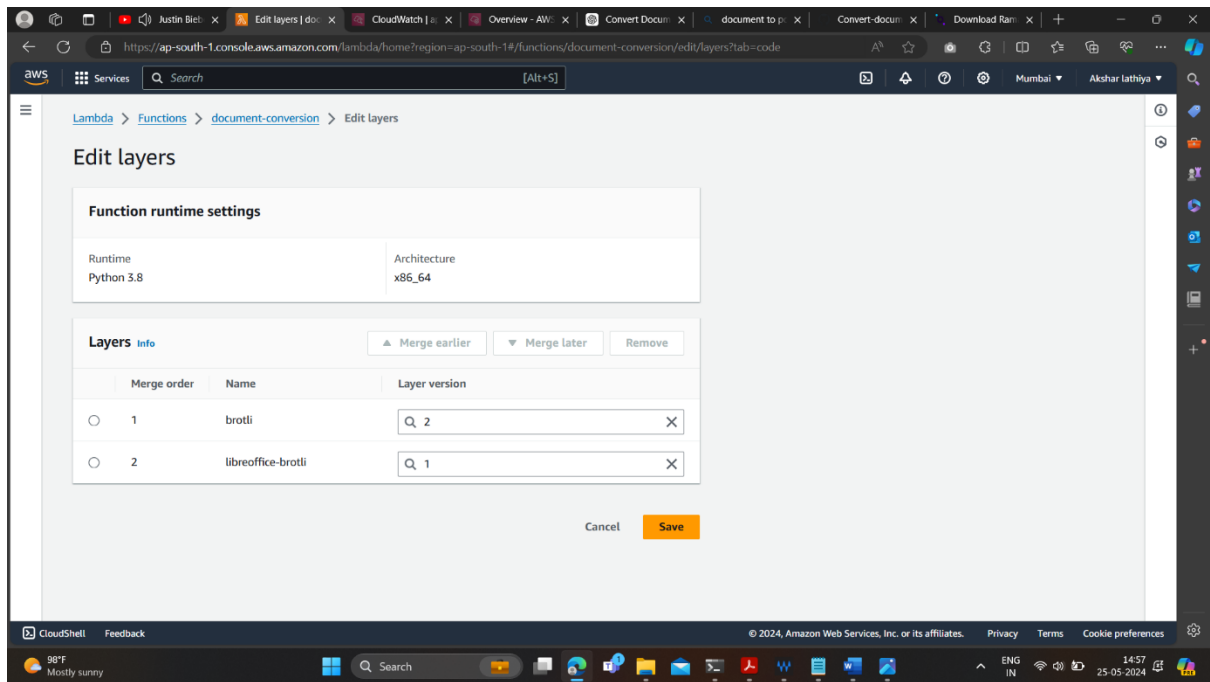
    # Document -> pdf
    proc = subprocess.run("/tmp/instdir/program/soffice.bin --headless --norestore --inv
    print('STDOUT: {}'.format(proc.stdout))
    print('STDERR: {}'.format(proc.stderr))

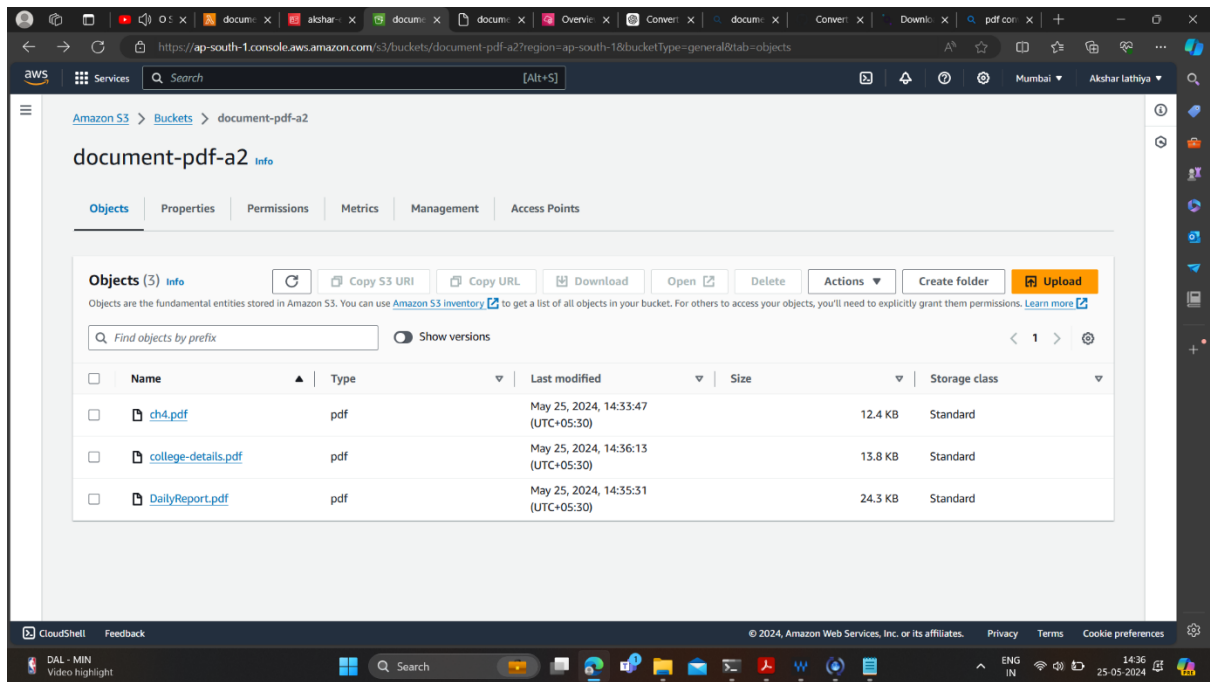
    # get pdf path
    key_list = input_key.split('.')
    pdf_path = "/tmp/"+input_key.replace(key_list[-1], 'pdf')

    # put S3 Object
    if os.path.exists(pdf_path):
        print('PDF: {}'.format(pdf_path.replace("/tmp/", "")))
        print('Size: {}'.format(os.path.getsize(pdf_path)))
        data = open(pdf_path, 'rb')
        s3.put_object(Bucket=output_bucket, Key=pdf_path.replace("/tmp/", ""), Body=data)
        data.close()
    else :
        print("The PDF file({}) cannot be found".format(pdf_path))

    return ''

```





Code:-

```
import subprocess
```

```
import tarfile
```

```
import sys
```

```
from io import BytesIO
```

```
import os
```

```
import time
```

```
import boto3
```

```
import urllib.parse
```

```
import json
```

```
sys.path.append("/opt/brotli")
```

```
import brotli
```

```
ACCESS_KEY = "AKIAYBBJNKIEYBERKW4H"
```

```
SECRET_ACCESS_KEY = "c17ipWwmATWOq+nkH9e+icMjE9/EnsjLYMyHW9dW"
```

```
def extract_libre_office():
```

```
    buffer = BytesIO()
```

```
    with open('/opt/lo.tar.br', mode='rb') as fout:
```

```
        file = fout.read()
```

```
        buffer.write(brotli.decompress(file))
```

```
        buffer.seek(0)
```

```
        with tarfile.open(fileobj=buffer) as tar:
```

```
            tar.extractall('/tmp')
```

```
def lambda_handler(event, context):
```

```
    if os.path.exists("/tmp/instdir/program/soffice.bin"):
```

```
        pass
```

```

else :

    # load libre

    extract_libre_office()

# Get Trigger event
if 'Records' in event.keys():

    input_bucket = event['Records'][0]['s3']['bucket']['name']

    output_bucket = "document-pdf-a2"

    input_key = urllib.parse.unquote_plus(event['Records'][0]['s3']['object']['key'],
encoding='utf-8')

    s3 = boto3.client('s3', aws_access_key_id=ACCESS_KEY,
aws_secret_access_key=SECRET_ACCESS_KEY)

else :

    return 'test finished'

# get S3 Object
file_path = '/tmp/'+input_key

s3.download_file(input_bucket, input_key, file_path)

# Document -> pdf

proc = subprocess.run("/tmp/instdir/program/soffice.bin --headless --norestore --invisible -
-nodetach --nofirststartwizard --nolockcheck --nologo --convert-to pdf:writer_pdf_Export --
outdir /tmp {}".format("/tmp/"+input_key), shell=True, stdout=subprocess.PIPE,
stderr=subprocess.PIPE)

print('STDOUT: {}'.format(proc.stdout))

print('STDERR: {}'.format(proc.stderr))

# get pdf path
key_list = input_key.split('.')

pdf_path = "/tmp/"+input_key.replace(key_list[-1], 'pdf')

```



```
# put S3 Object
if os.path.exists(pdf_path):
    print('PDF: {}'.format(pdf_path.replace("/tmp/", "")))
    print('Size: {}'.format(os.path.getsize(pdf_path)))
    data = open(pdf_path, 'rb')
    s3.put_object(Bucket=output_bucket, Key=pdf_path.replace("/tmp/", ""),Body=data)
    data.close()
else :
    print("The PDF file({}) cannot be found".format(pdf_path))

return "
```

libreoffice aws arn :- arn:aws:lambda:ap-south-1:764866452798:layer:libreoffice-brotli:1