



Write a python 3 script for a digital forensic investigation to analyse an .e01 hard disk image and carve out all png files



To analyze an .E01 hard disk image and carve out all .png files in Python, you can use the pytsk3 and pyewf libraries for accessing the disk image and the filesystem, and the Python Imaging Library (PIL) to validate the extracted .png files. Here's a script that does this:



1. Install the required libraries:

```
bash
pip install pytsk3 pyewf Pillow
```

2. Create the script `carve_png_files.py`:

```
python
import os
import sys
import pytsk3
import pyewf
from io import BytesIO
from PIL import Image

class EWFImgInfo(pytsk3.Img_Info):
    def __init__(self, ewf_handle):
        self._ewf_handle = ewf_handle
        super(EWFImgInfo, self).__init__(url="", type=pytsk3.TSK_IMG_TYPE_EX)

    def close(self):
        self._ewf_handle.close()

    def read(self, offset, size):
        self._ewf_handle.seek(offset)
        return self._ewf_handle.read(size)

    def get_size(self):
        return self._ewf_handle.get_media_size()

def carve_png_files(image_path, output_dir):
    ewf_handle = pyewf.handle()
    ewf_handle.open(image_path)

    img_info = EWFImgInfo(ewf_handle)
    volume = pytsk3.Volume_Info(img_info)

    for part in volume:
        if part.flags != pytsk3.TSK_VS_PART_FLAG_ALLOC:
            continue

        fs_info = None
        try:
            fs_info = pytsk3.FS_Info(img_info, offset=part.start * volume.in
        except Exception as e:
            print(f"Cannot read filesystem at offset {part.start * volume.in
            continue

        file_count = 0
        for root, dirs, files in fs_info.walk():
            for file in files:
                if not file.name.lower().endswith(".png"):
                    continue

                try:
                    file_data = b"".join(list(file))
                    img = Image.open(BytesIO(file_data))
                    img.verify()

                    if not os.path.exists(output_dir):
                        os.makedirs(output_dir)

                    output_file = os.path.join(output_dir, f"{file.name}_{fi
                    with open(output_file, "wb") as f:
                        f.write(file_data)

                    file_count += 1
                    print(f"Carved {file.name} and saved as {output_file}")

                except Exception as e:
                    print(f"Error processing file {file.name}: {e}")

    ewf_handle.close()

if __name__ == "__main__":
    if len(sys.argv) != 3:
        print("Usage: python carve_png_files.py [path_to_e01_file] [output_d
        sys.exit(1)

    image_path = sys.argv[1]
    output_dir = sys.argv[2]

    carve_png_files(image_path, output_dir)
```

3. Run the script:

```
bash
python carve_png_files.py <path_to_e01_file> <output_directory>
```

Replace `<path_to_e01_file>` with the path to your .E01 file and `<output_directory>` with the directory where you want to save the carved .png files. The script will access the .E01 hard disk image, iterate through the file system, and carve out all .png files, validating and saving them in the specified output directory.



Can you update that script to run on a raw image file and carve out all PDF files?



Sure! You can modify the script to work with raw image files and carve out all PDF files. You will not need the pyewf library for working with raw images, so you can remove it from the imports and installation instructions.



Here's the updated script `carve_pdf_files.py`:

```
python
import os
import sys
import pytsk3
from io import BytesIO

def is_pdf(data):
    if len(data) < 4:
        return False
    return data[:4] == b"%PDF"

def carve_pdf_files(image_path, output_dir):
    img_info = pytsk3.Img_Info(image_path)
    volume = pytsk3.Volume_Info(img_info)

    for part in volume:
        if part.flags != pytsk3.TSK_VS_PART_FLAG_ALLOC:
            continue

        fs_info = None
        try:
            fs_info = pytsk3.FS_Info(img_info, offset=part.start * volume.in
        except Exception as e:
            print(f"Cannot read filesystem at offset {part.start * volume.in
            continue

        file_count = 0
        for root, dirs, files in fs_info.walk():
            for file in files:
                if not file.name.lower().endswith(".pdf"):
                    continue

                try:
                    file_data = b"".join(list(file))
                    if is_pdf(file_data):

                        if not os.path.exists(output_dir):
                            os.makedirs(output_dir)

                        output_file = os.path.join(output_dir, f"{file.name}
                        with open(output_file, "wb") as f:
                            f.write(file_data)

                        file_count += 1
                        print(f"Carved {file.name} and saved as {output_file

                except Exception as e:
                    print(f"Error processing file {file.name}: {e}")

    if __name__ == "__main__":
        if len(sys.argv) != 3:
            print("Usage: python carve_pdf_files.py [path_to_raw_image_file] [ou
            sys.exit(1)

        image_path = sys.argv[1]
        output_dir = sys.argv[2]

        carve_pdf_files(image_path, output_dir)
```

To run the script, use the following command:

```
bash
python carve_pdf_files.py <path_to_raw_image_file> <output_directory>
```

Replace `<path_to_raw_image_file>` with the path to your raw image file and `<output_directory>` with the directory where you want to save the carved .pdf files. The script will access the raw hard disk image, iterate through the file system, and carve out all .pdf files, validating and saving them in the specified output directory.



Is that the most efficient way of doing it?



The script provided earlier works well for small to medium-sized disk images, but it's not the most efficient for large disk images or when you are looking for a specific file type. A more efficient approach is to use a file carving technique that searches for file signatures directly, rather than relying on the file system metadata.



Here's a more efficient script using the file carving technique to extract PDF files from a raw disk image:

```
python
import os
import sys
import re
import mmap

PDF_HEADER_SIGNATURE = b"%PDF"
PDF_FOOTER_SIGNATURE = b"%EOF"

def carve_pdf_files(image_path, output_dir):
    if not os.path.exists(output_dir):
        os.makedirs(output_dir)

    with open(image_path, "rb") as f:
        mmaped_file = mmap.mmap(f.fileno(), 0, access=mmap.ACCESS_READ)
        headers = [match.start() for match in re.finditer(PDF_HEADER_SIGNATURE
        footers = [match.start() for match in re.finditer(PDF_FOOTER_SIGNATU

    file_count = 0

    for header in headers:
        footer = -1
        for ftr in footers:
            if ftr > header:
                footer = ftr + len(PDF_FOOTER_SIGNATURE)
                footers.remove(ftr)
                break

        if footer == -1:
            continue

        pdf_data = mmaped_file[header:footer]
        output_file = os.path.join(output_dir, f"carved_pdf_{file_count}.pdf

        with open(output_file, "wb") as out_f:
            out_f.write(pdf_data)

        file_count += 1
        print(f"Carved PDF and saved as {output_file}")

    if __name__ == "__main__":
        if len(sys.argv) != 3:
            print("Usage: python carve_pdf_files.py [path_to_raw_image_file] [ou
            sys.exit(1)

        image_path = sys.argv[1]
        output_dir = sys.argv[2]

        carve_pdf_files(image_path, output_dir)
```

This script uses memory-mapped files and regex to search for the PDF header and footer signatures directly in the raw disk image. It then carves out the files based on the found signatures, which is a more efficient way of extracting PDF files from a raw disk image compared to the previous approach that relies on file system metadata.

To run the script, use the following command:

```
bash
python carve_pdf_files.py <path_to_raw_image_file> <output_directory>
```

Replace `<path_to_raw_image_file>` with the path to your raw image file and `<output_directory>` with the directory where you want to save the carved .pdf files. The script will access the raw hard disk image, search for PDF files, and save them in the specified output directory.



What performance benefit does the latest version of the code have over the original suggestion?



The performance benefits of the latest version of the code (file carving technique) over the original suggestion (file system-based approach) can be summarized as follows:



- Independent of the file system: The file carving technique does not rely on the file system structure and metadata, making it more versatile for analyzing disk images with damaged, corrupted, or unsupported file systems.
- Faster search: By using memory-mapped files and regex to search for the file signatures directly, the file carving technique can identify and extract files more quickly than the file system-based approach, which requires navigating through the directory structure and file metadata.
- Recovery of deleted files: The file carving technique can potentially recover deleted or partially overwritten files because it searches for file signatures directly in the raw data, whereas the file system-based approach relies on file metadata, which may not be available for deleted files.
- Improved scalability: The file carving technique is more suitable for large disk images because it can process the data more efficiently than the file system-based approach. Memory-mapped files enable the script to process the disk image in chunks, reducing the memory footprint and making it more suitable for larger data sets.

However, it's worth noting that the file carving technique also has some drawbacks:

- Loss of original file names: The file carving technique does not retain the original file names because it does not rely on file system metadata. The files are saved with generic names instead.
- False positives: The file carving technique can produce false positives if the raw data contains byte sequences that match the file signatures, but are not part of an actual file.
- Fragmented files: The file carving technique might not be able to recover fragmented files correctly because it assumes that files are stored in a contiguous block of data.

Despite these drawbacks, the file carving technique is generally more efficient and versatile for extracting specific file types from disk images, especially when dealing with large data sets or damaged file systems.