

Harnessing LLMs for AppSec

Seth Law & Ken Johnson

Introductions - Seth Law



Founder of **Redpoint Security**

Based in Salt Lake City, Utah

Consultant, Programmer, Trainer,
Speaker

@sethlaw

Introductions - Ken Johnson

CTO @ DryRun Security // Lives in Northern VA

Building AI powered AppSec Automation

AppSec Background:

- GitHub
- CTO @ “a consultancy”
- LivingSocial
- FishNet Security (now Optiv)
- US Pentagon

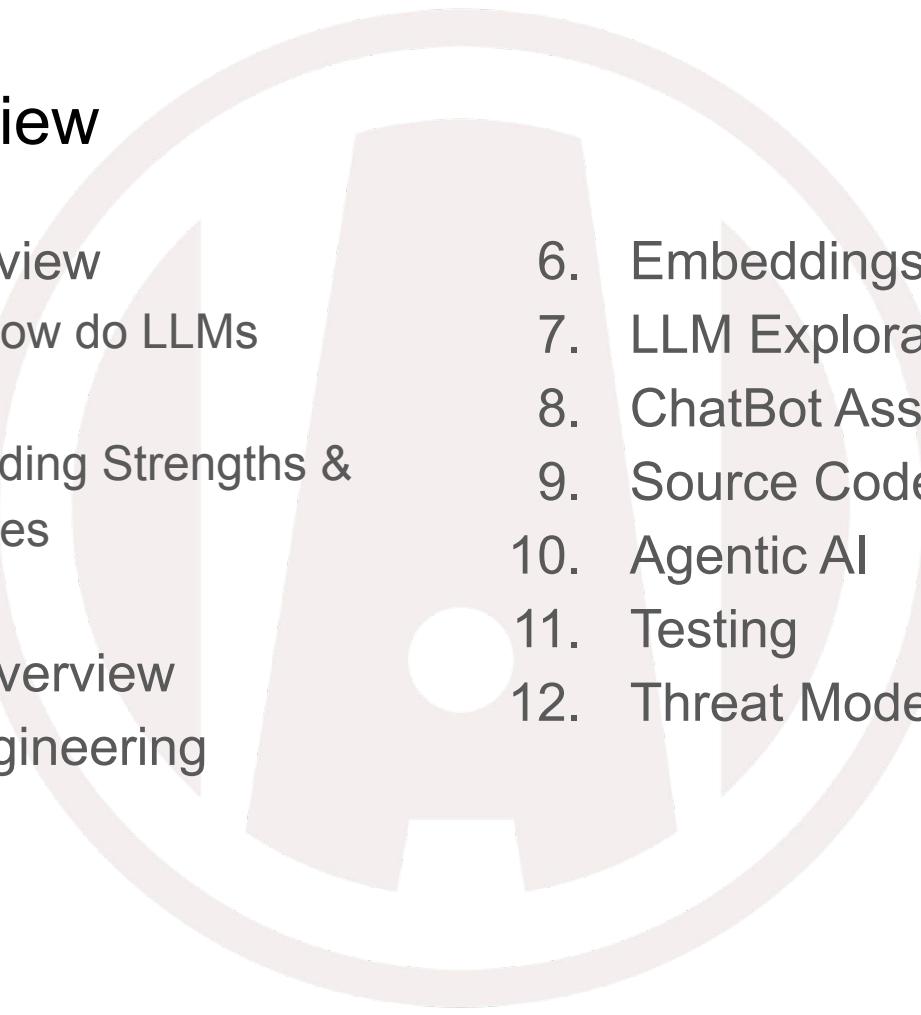


Introductions

- Name
- Current job/title/interest/career path
- Favorite LLM (or programming language)



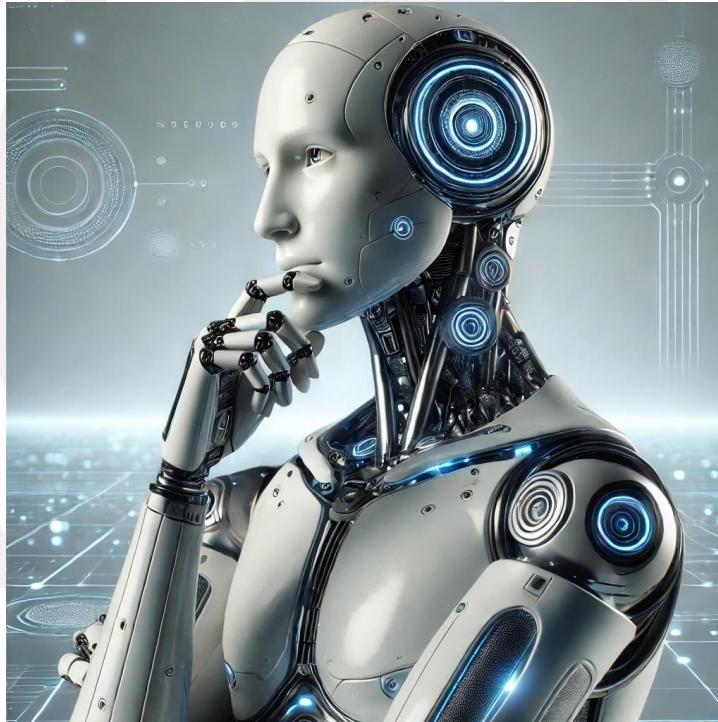
Course Overview

- 
1. Intros & Overview
 - a. Gen AI - How do LLMs work?
 - b. Understanding Strengths & Weaknesses
 2. Lab Setup
 3. LangChain Overview
 4. Prompts / Engineering
 5. Context
 6. Embeddings & Vector Stores
 7. LLM Exploration
 8. ChatBot Assistant
 9. Source Code Analysis
 10. Agentic AI
 11. Testing
 12. Threat Modeling

Background behind the Course

- We are using this technology today and seeing massive benefits
 - Explanations of Code Changes - Pull Request Summaries
 - Helps easily categorize changes in an organization like new endpoints, authentication mechanisms, new infrastructure, etc.
 - Deconstructing code bases and providing compositional analysis
 - Detection of anomalous patterns (Surveyor)
 - Creating risk scores and applying them to various security related analysis
 - Chat - extending your team by answering first level questions
 - ... really just uncovering new use cases all the time
- We want the Security industry to embrace this new technology! Not get left behind

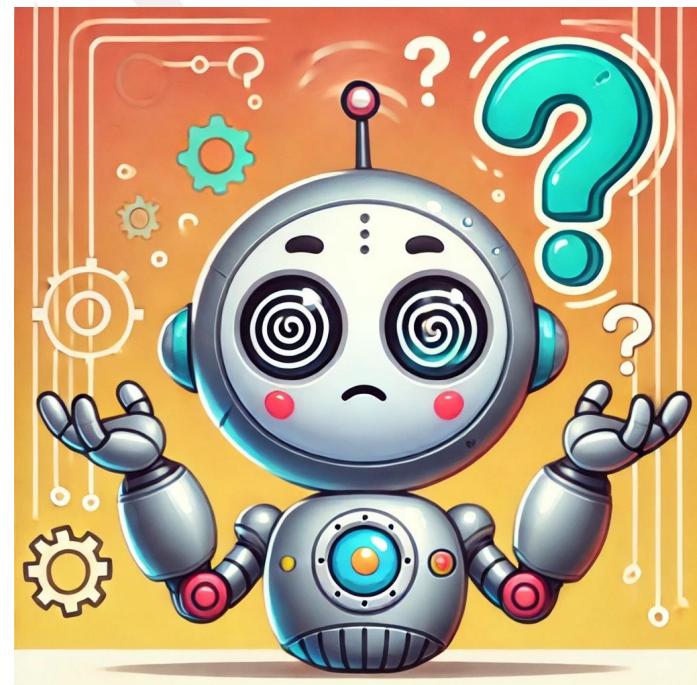
Gen AI Concepts - How do LLMs work?



Gen AI Concepts - How do LLMs work?

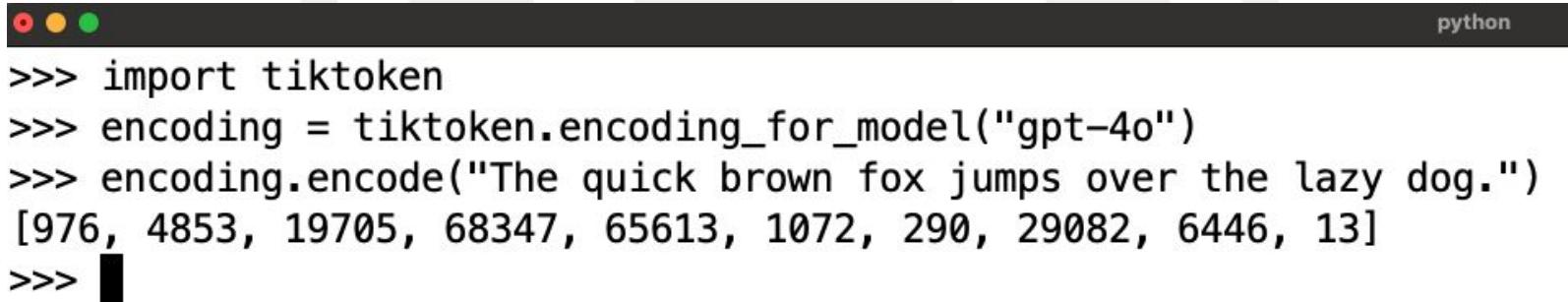
Large Language Models / Small Language Models

- How do they work
 - Tokens & Prediction
 - Inference
 - Prompting
 - Training vs RAG
- Common Misconceptions
- Notable events



Gen AI Concepts - How do LLMs work?

BPE - Byte Pair Encoding creates a vocabulary for the LLM mapping characters to tokens



A terminal window interface with three colored status indicators (red, yellow, green) at the top left. The word "python" is at the top right. The main area contains the following Python code:

```
>>> import tiktoken  
>>> encoding = tiktoken.encoding_for_model("gpt-4o")  
>>> encoding.encode("The quick brown fox jumps over the lazy dog.")  
[976, 4853, 19705, 68347, 65613, 1072, 290, 29082, 6446, 13]  
>>> █
```

Gen AI Concepts - How do LLMs work?

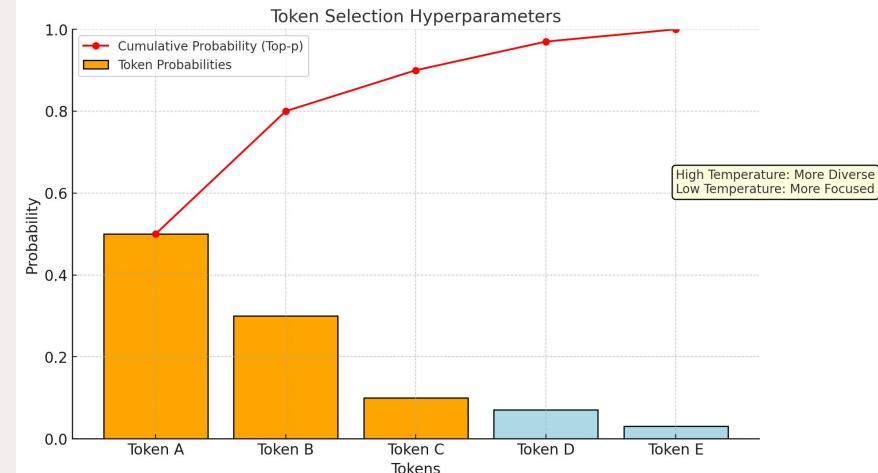
Tokens & Token Prediction

```
▶▶▶ encoding.encode("Payment-Gateway")
[17036, 13112, 17182]
▶▶▶ encoding.decode([17036])
'Payment'
▶▶▶ encoding.decode([13112])
'-G'
▶▶▶ encoding.decode([17182])
'ateway'
▶▶▶ █
```

Gen AI Concepts - How do LLMs work?

Token Prediction in LLMs:

1. **Input:** Receive tokenized inputs from the user.
2. **Prediction:** Generate a ranked list of the most probable next tokens based on the input sequence.
3. **Selection:** Choose the next token using either:
 - **Greedy:** Select the highest-probability token.
 - **Probabilistic:** Random selection weighted by the assigned probabilities.
4. **Control:** Probabilities are influenced by hyperparameters like:
 - **Temperature:** Adjusts randomness (higher = more random).
 - **Top-k:** Limits choices to the top-k highest probabilities.
 - **Top-p (nucleus sampling):** Limits choices to a cumulative probability threshold.



Gen AI Concepts - How do LLMs work?

Hyperparameters for Token Selection

1. **Temperature**
 - Adjusts randomness in token selection.
 - **Low values (e.g., 0.2):** Focused on high-probability tokens (deterministic).
 - **High values (e.g., 1.5):** More diverse and creative outputs (random).
2. **Top-p (Nucleus Sampling)**
 - Chooses from tokens that cumulatively cover a probability threshold (e.g., 90%).
 - Ensures only the most relevant tokens (based on context) are considered.
3. **Top-k**
 - Restricts token choices to the **k most probable options**.
 - Example: Top-3 only considers the three tokens with the highest probabilities.

Gen AI Concepts - How do LLMs work?

Helpful Articles:

<https://writings.stephenwolfram.com/2023/02/what-is-chatgpt-doing-and-why-does-it-work/>

<https://blog.miquelgrinberg.com/post/how-langs-work-explained-without-math>

Gen AI Concepts - How do LLMs work?

Training:

- I love AppSec
- You love AppSec
- We love AppSec

-	I	you	we	love	AppSec
I				1	
you				1	
we				1	
love					3
AppSec					

Gen AI Concepts - How do LLMs work?

Training:

- I love AppSec
- You love AppSec
- We love AppSec

-	I	you	we	love	AppSec
I					100%
you					100%
we					100%
love	33.3%	33.3%	33.3%		100%
*AppSec	25%	25%	25%	25%	

* Denotes a token that has no “next” tokens (ie - “Hole”)

LLMs Strengths & Weaknesses



LLMs Strengths & Weaknesses

Strengths

1. Summarization of Code Changes
2. Approximation of Intent and Behavior
3. Pattern Detection and Prediction
4. Automating Repetitive Tasks
5. Contextual Awareness
6. Support for Multiple Languages and Frameworks
7. Efficiency and Scalability

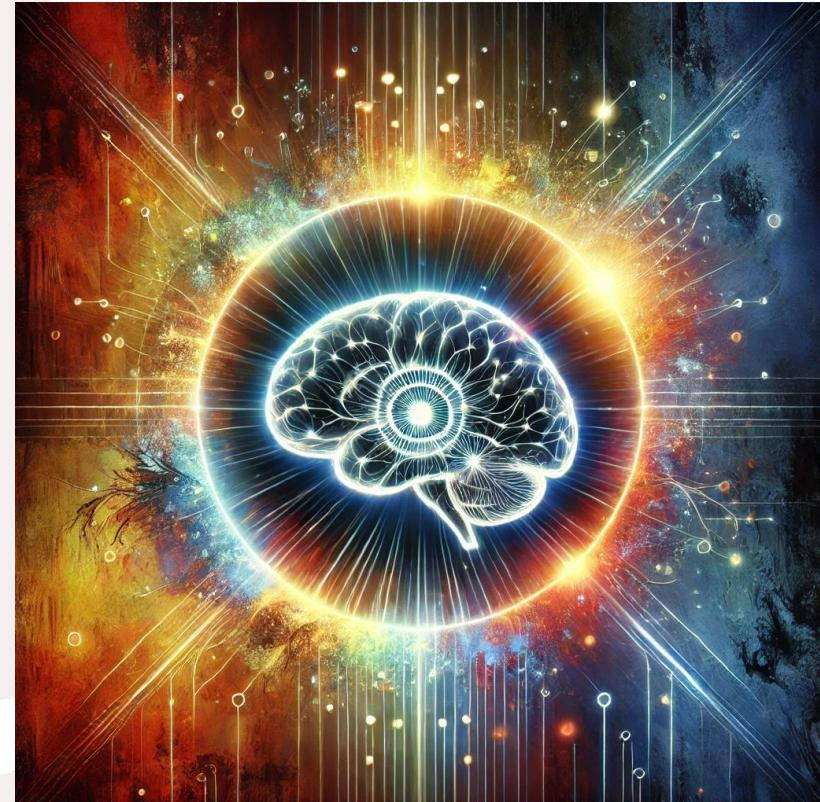
Weaknesses

1. Non-Deterministic Outputs
2. Hallucination of Risks
3. Resource Intensity
4. Lack of True Causal Reasoning
5. Dependency on Training Data
6. Difficulty Handling Complex Codebases
7. False Confidence in Results
8. Ambiguity in Prompts

LLMs Strengths & Weaknesses

We can mitigate almost all LLM deficiencies via:

- Context
- Prompt Engineering
- Technique
- Tuning
- Agents & Custom Tools
- Proper Testing



Lab Setup



Lab Setup

1. Clone the following repo:

https://github.com/absoluteappsec/harnessing_llms

2. Navigate to cloned repo, type:

```
$ python -m venv venv
```

```
$ source venv/bin/activate
```

```
$ pip install -r requirements.txt
```

3. Download your .env file (Gist link provided to you in an email) and place in the root of the cloned repo

LangChain



LangChain

Simplifies process of creating advanced AI applications

Provides abstractions for many things such as:

1. Chains - Combines multiple utilities into a single workflow
2. Memory - Adds state to chains, retaining context across multiple interactions
3. Agents - Enable dynamic decision-making, where the system determines what actions to take
4. Retrieval Augmented Generation (RAG) - Integrate with external knowledge bases to enhance outcomes
5. Much more!

LangChain - Equivalent Frameworks

Haystack -

<https://haystack.deepset.ai/>

Semantic Kernel -

<https://github.com/microsoft/semantic-kernel>

LLamaIndex

https://github.com/run-llama/llama_index

Comparison Table

Feature	LangChain	LlamaIndex	Haystack	Semantic Kernel	Hugging Face Transformers	OpenAI Function Calling
Workflow Management	✓	Limited	✓	✓	✗	✗
Memory Support	✓	✗	✓	✓	✗	✗
RAG Support	✓	✓	✓	✓	✗	✗
Indexing	✓	✓	✓	✗	✗	✗
Ease of Use	High	Moderate	Moderate	High	Low	High

LangChain

LangChain has 3 primary products:

LangChain (OSS) - <https://github.com/langchain-ai/langchain>

LangSmith (Paid) - <https://smith.langchain.com/>

LangGraph (Paid) - <https://www.langchain.com/langgraph>

LangChain

Things to know before getting started:

- Frequent changes means even minor version changes can break production systems
- Documentation is improving but not necessarily the best (**sorted by Python & JS versions**) -
<https://python.langchain.com/docs/introduction>

Primary Components:

langchain-core: Base abstractions and LangChain Expression Language (LCEL).

Langchain: Chains, agents, and retrieval strategies

Langchain-community: Third party integrations, community-maintained

LangChain - Navigating Docs

Introduction: Good way to get a lay of the land

Tutorials: Dig in right away to some examples of building different AI applications

How-to Guides: Guides that are goal oriented/concrete

Conceptual Guide: Provides some high level information on the various components that you will utilize within LangChain

API Reference:

The screenshot shows the LangChain Python API Reference documentation. The left sidebar has a dark theme with white text. It features a logo at the top, followed by 'Section Navigation' and two main categories: 'Base packages' and 'Integrations'. Under 'Base packages', there are links for Core, Langchain, Text Splitters, Community, and Experimental. Under 'Integrations', there are links for AI21, Anthropic, Astradb, AWS, Azure Dynamic Sessions, Box, Cerebras, Chroma, Cohere, Couchbase, Databricks, Elasticsearch, Exa, Fireworks, Google Community, Google GenAI, Groq, Huggingface, IBM, Milvus, MistralAI, MongoDB, and Neo4J. The main content area has a light background. At the top, it says 'LangChain Python API Reference'. Below that is a header 'LangChain Python API Reference #'. It includes a welcome message, user guides, and legacy API reference information. There are two main sections: 'Base packages #' and 'Integrations #'. Each section contains a table with two columns. The 'Base packages #' table has three rows: Core (langchain-core: 0.3.21), Langchain (langchain: 0.3.9), and Text Splitters (langchain-text-splitters: 0.3.2). The 'Community' row in the 'Base packages' sidebar also links to the 'Community' row in the table. The 'Integrations #' table has four rows: OpenAI (langchain-openai 0.2.10), Anthropic (langchain-anthropic 0.3.0), AWS (langchain-aws 0.2.7), Huggingface (langchain-huggingface 0.1.2), MistralAI (langchain-mistralai 0.2.2), and AI21 (langchain-ai21 1.0.1). A note at the bottom of the integrations table says 'See the full list of integrations in the Section Navigation.'

LangSmith

The screenshot shows the LangSmith application interface. At the top, there's a navigation bar with icons for home, back, and search, followed by the path: DryRun Security > Tracing projects > pr-tldr-dev-analyzer. Below the path is the project name "pr-tldr-dev-analyzer". A horizontal menu bar includes "Runs" (which is underlined), "Threads", "Monitor", and "Setup". Under "Runs", there are buttons for "2 filters", "Last 7 days", "Root Runs", "LLM Calls", and "All Runs". The main content area is titled "Runs" and lists 12 entries, each with a checkbox, a green checkmark icon, the name of the run, and its input file. The entries are:

	Name	Input
checkbox	Generate TLDR	[{ "filename": "dryrun-api.yaml", "description": "The code change in the 'dryrun-api.yaml' file is a minor update to the API endpoint descriptions." }]
checkbox	Summarize File	serverless.yml
checkbox	Summarize File	pkg/usecases/users_test.go
checkbox	Summarize File	pkg/usecases/users.go
checkbox	Summarize File	pkg/usecases/user_creator_test.go
checkbox	Summarize File	pkg/usecases/account_integrations_test.go
checkbox	Summarize File	pkg/domain/user.go
checkbox	Summarize File	pkg/api/user.go
checkbox	Summarize File	handler/verify-email/main.go
checkbox	Summarize File	handler/resend-verification-email/main.go
checkbox	Summarize File	handler/user-update/main.go

LangSmith

The screenshot shows the LangSmith application interface. On the left, there's a sidebar titled "TRACE" with buttons for "Collapse", "Stats", "Filter", and "Show All". Below this are two items: "Summarize File" (status: 5.45s, seq: 11,888) and "PromptTemplate" (status: 0.00s, seq: step:1). Further down is "ChatBedrock anthropic.claude-3-haiku-20240307-v1:0" (status: 5.44s, seq: 11,888, seq: step:2). The main area has tabs for "Run", "Feedback", and "Metadata", with "Run" selected. Under "Input", there's a code editor showing a YAML configuration file. The code is as follows:

```
1 filename: serverless.yml
2 patch: |-
3   @@ -47,6 +47,7 @@ provider:
4     FRONTEND_URL: ${self:custom.frontendUrl.${sls:stage}}
5     HARD_REPO_LIMIT: ${self:custom.hardRepoLimit}
6     FEATURE_GATE_MODEL: ${self:custom.featureGateModel.${sls:stage}}
7     + ALLOWED_USERS: ${self:custom.allowedUsers.${sls:stage}}
8       iam:
9         role:
10           statements:
11 contents: |
12   # Welcome to Serverless!
13   #
14   # This file is the main config file for your service.
15   # It's very minimal at this point and uses default values.
16   # You can always add more config options for more control.
17   # We've included some commented out config examples here.
18   # Just uncomment any of them to get that config option.
19   #
20   # For full config options, check the docs:
21   #   docs.serverless.com
22   #
23   # Happy Coding!
24
25 service: dryrun-api
26
```

The code editor has a "YAML" dropdown at the bottom. At the very bottom of the screen, there's a "Rendered Output" dropdown.

LangGraph (Agentic development)

<https://langchain-ai.github.io/langgraph/>

Server: Opinionated API / Arch for building agentic apps

Studio: IDE that connects to a server for visualization/debugging

CLI: Interface that enables local interactions with LangGraph

SDK: Provides a programmatic way to interact with deployed LangGraph Applications

Remote Graph: Allows you to interact with any deployed LangGraph application as though it were running locally

Prompt Engineering



Prompt Engineering

Prompts are the input instructions or text provided to a large language model (LLM) to guide its response. They set the context, define the task, and influence the quality and relevance of the output.

💡 Be mindful of context length windows however... they're fairly large these days



<https://github.com/dair-ai/Prompt-Engineering-Guide>

Context: Background information on task

Instruction: Specify what the LLM should do

Input Data: Dynamic/user input and usually defined as a Human or User prompt

Output Format: Defines how the response should be structured

Constraints or Rules: If you need to define how the response should look

Examples: Usually done with few shot prompting

Tone/Style: Describes the style in which the LLM should response (ie - professional, pirate, yoda, etc.)

Prompt Engineering

Some popular frameworks:

CO-STAR: Context, Objective, Scale, Time, Actor, Resources

RACE: Role, Action, Context, Explanation

CARE: Context, Action, Result, Example

APE: Action, Purpose, Execution

RICE: Role, Input, Steps, Execution

Prompt Engineering - Co-Star Example

Context: This code is part of a login system for a financial application.

Objective: Identify potential authentication security vulnerabilities in this code.

Scale: Focus only on the `authenticate_user` function

Time: This analysis must be completed in under 1 minute

Actor: You are an application security engineer.

Resources: Refer to OWASP's top 10 vulnerabilities for guidance.



Exercise: Craft Prompt

Exercise: Craft Prompt

🎯 Objective: Find out as much information about the Juice Shop application as possible

Instructions:

1. Open `scripts/prompt_engineering.py`
2. Modify both “`system_prompt_template`” and “`user_question`” fields
3. Have fun! This is free form, use one of the prompt engineering frameworks to guide you
4. ⚠️ Only modify user questions and system prompts for now

```
system_prompt_template = """  
You are a helpful code review assistant who is proficient  
in both security as well as functional review. You will be  
provided source code of a web application and tasked with  
answering questions about it.
```

```
<context>  
{context}  
</context>  
....
```

```
user_question = """Tell me the following information about the code base I am  
providing you:
```

- Purpose of the application
- Web technologies used in the application
- Templating language used in the application
- Database used in the application
- Authentication mechanisms used in the application
- Authorization mechanisms used in the application

```
List libraries by their name, purpose, and version that are  
used in the vtm application for the following categories:
```

- Security
- Testing
- Documentation
- Build
- Database
- Authentication / Authorization
- HTML Templating (ex: pug, handlebars)
- CSS Frameworks (ex: bootstrap, tailwind)
- widgets / UI components

```
....
```

Exercise: Craft Prompt

Retrospective:

- What worked and what did not?
- What kind of experiments did we run?
- What are your thoughts on how modification of prompts affect the outcome?
- Anything else interesting?



Context



Context

Context is the background information or situational details provided in a prompt to help the LLM understand the task environment and relevance.

💡 Context can be a number of things just remember it's included in the “context window” length calculation

Can be provided in a number of ways:

- Loaded from a static source such as a yaml file, database entry, or markdown file
- Chained using a vector store
- Statically defined inside of the prompt

Context Examples

- Technical documentation
- Large portions of source code
- Call Graphs / Abstract Syntax Tree (AST)
- Directory listings
- Framework or language information
- System or application logs
- Output of Scanner Findings



Context

Loading the contents of a markdown file to provide information about node.js web applications.

Typically, Node.js web application libraries will be located in a `yarn` file or `package.json`

When it comes to Node.js web applications, the following is a list of security-specific libraries:

Input Sanitization/Validation

==|

- Joi
- Validator
- Express-Validator

Authentication

==|

- jsonwebtoken
- passport.js
- express-jwt

Authorization

==|

- casbin
- accesscontrol

Secure Session Management

Contents of “juice_shop_knowldegebase.md”

Context

Now the code is “code” and the file contents are a “knowledgebase”

```
system_prompt_template = """  
You are a highly skilled and detail-oriented code review assistant with expertise in  
both application security and functional code analysis. Your role is to assist  
developers and security professionals by providing accurate, concise, and actionable  
insights.
```

```
Your task is to analyze the provided source code of a web application and answer  
specific questions about its functionality, security, and technologies. Always  
maintain a professional tone and prioritize clarity in your responses.
```

```
In your analysis:  
- Clearly identify and explain the purpose and technologies used in the codebase.  
- Highlight critical security mechanisms such as authentication and authorization.  
- Provide details on libraries, tools, and frameworks, organized by their categories  
and roles in the application.  
- When relevant, make recommendations for improving security or functionality.
```

```
Search the following codebase to answer the questions:
```

```
<code>  
{code}  
</code>
```

```
Use the following context to help answer questions:  
<context>  
{context}  
</context>
```

```
chain = (  
    {  
        "code": retriever,  
        "question": RunnablePassthrough(),  
        "context": RunnablePassthrough()  
    }  
    | prompt  
    | llm  
    | StrOutputParser()  
)
```

```
print('performing analysis....')  
response = chain.invoke(input=user_question, context=context)  
print(response)
```

We pass the context as
“RunnablePassThrough()” and use
.invoke() to pass in the context

Exercise: Craft Context



Exercise: Craft Context

🎯 **Objective:** Experiment with various types of context to improve the results of Juice Shop analysis

Instructions:

1. Open `scripts/building_with_context.py`
2. Open `data/juice_shop_knowledgebase.md`
3. Experiment with providing different types of context to improve the results of our LLMs analysis
4. ⚠️ Only modify context and prompt

Exercise: Craft Context

Retrospective:

- What types of context did you experiment with?
- Were there any combinations of prompt + context that worked better than others?

Embeddings & Vector Stores



Embeddings & Vector Stores

Embeddings:

Dense, numerical representations of data (e.g., words, sentences, or code) in a continuous vector space. They capture semantic meaning, enabling similar items (like related words or similar code snippets) to have closer representations in this space.

- Dense Representation - Fixed Size
- Semantic Encoding - Relationships between things
- Domain Specific - Can be pre-trained

Vector Stores:

Vector stores are specialized databases that store and manage embeddings. They enable efficient similarity searches by finding vectors close to a given query, often used in applications like document retrieval, recommendation systems, or question answering.

 For local development, we often use FAISS. For production purposes, Pinecone works very well. Many options are available!

Embeddings & Vector Stores

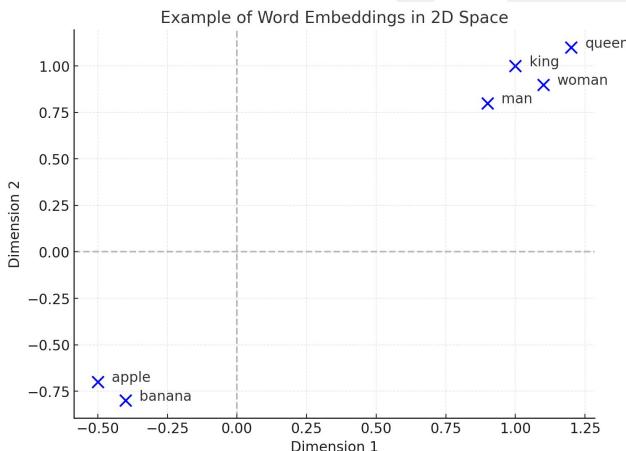
Embedding Format Example

For the word "quick", a 3-dimensional embedding might look like this:

plaintext

Copy code

```
[0.345, -0.221, 0.891]
```



BROWSER METRICS NAMESPACES (6) IMPORTS

Index records

Namespace: (Default) Query by: dense vector value Vector: 0.43,0.17,0.5,0.78,0.14,0.97,0.91,0.71,0.87,0.82,0.21,0.16,0.39,0.97,0.61,0.1,0.31,0.92,0.49,0.83,0.37,0.41,0.12, ... Top K: 10 Query

Matches: 1

ID	VALUES
d3a27ef4-48...	0.75390625, -0.34375, -0.37890625, -0.51171875, -0.2734375, -0.143554688, -0.247070312, -0.00076675415, 0.34765625, -0.2734375, 0...

SCORE: -0.0119 METADATA: text: "This is a content of the document"

Embeddings & Vector Stores

 Not every LLM supports Embeddings, this is why you will notice that we utilize one LLM for inference tasks and another for performing Embeddings.

```
embeddings = BedrockEmbeddings(model_id='amazon.titan-embed-text-v1')
```

```
llm = ChatBedrock(  
    model_id='anthropic.claude-3-haiku-20240307-v1:0',  
    model_kwargs={"temperature": 0.6},  
)
```

Embeddings & Vector Stores

Embeddings:

Comprehensive list of vector databases supported by LangChain

<https://python.langchain.com/v0.1/docs/integrations/vectorstores/>

For our local purposes, you'll notice that we're utilizing FAISS

 Pro-Tip, save the index database so that you do not have to create embeddings and then store them every  single  time 

```
db = FAISS.from_documents(texts, embeddings)  
db.save_local("juice_shop.faiss")
```

Embeddings & Vector Stores

Considerations:

- Embeddings & Vector stores must align dimension-wise
 - If the embedding LLM only supports 768 dimensions but you're using a 1536 dimension db... issues shall arise
- Cost
- Production or for local dev?
- How much control do you want?
 - Example: PGVector vs Pinecone

Embeddings & Vector Stores

Text Splitting:

- Size limitations - remember embedding is similar to an inference operation, it has a set context window length
- When splitting code, LangChain offers a utility which separates code by language
- To see the list of supported languages:

```
>>> from langchain_text_splitters import ( Language)
>>> [e.value for e in Language]
['cpp', 'go', 'java', 'kotlin', 'js', 'ts', 'php', 'proto', 'python',
 'rst', 'ruby', 'rust', 'scala', 'swift', 'markdown', 'latex', 'html',
 'sol', 'csharp', 'cobol', 'c', 'lua', 'perl', 'haskell', 'elixir',
 'powershell']
>>> █
```

```
loader = GenericLoader.from_filesystem(
    local_path,
    glob="**/*",
    suffixes=['.py'],
    parser=LanguageParser(language=Language.PYTHON),
    show_progress=True
)
```

```
splitter = RecursiveCharacterTextSplitter.from_language(language=Language.PYTHON,
    chunk_size=8000,
    chunk_overlap=20
)
texts = splitter.split_documents(documents)

db = FAISS.from_documents(texts, embeddings)
```

Exercise: Embedding & Vectorstores

Exercise: Embed(dings) & (Vector) Store

🎯 **Objective:** Load a repository of your choosing and store in our vector database FAISS, use previous code from “`build_context.py`” to analyze the database

Instructions:

1. Open `scripts/embed_and_store.py`
2. Modify the “RELEVANT” portions (specified in comments) to support the repo of your choice
3. Incorporate `building_with_context.py` code to build a feature complete script

Exercise: Embed(dings) & (Vector) Store

Retrospective:

- What type of repo(s) did you experiment with?
- What challenges did you face?
- Did you experiment with any of the splitting options?
- Are we feeling more comfortable with loading and analyzing code?

LLM Exploration



LLM Exploration

Considerations:

- Commercial vs OSS
- Transactional vs Dedicated
- Training vs Agentic / RAG (*can be all 3)
- Context-Window Size
- Purpose of the LLM

Commercial Vs OSS

Commercial

- Supports larger context windows
- Support usually is an option
- Good for transactional costs
- Bad option for training (though “fine tuning” options exist in some cases)

OSS

- Free to use, not to run
- No support
- Hosted on AI’s equivalent of “npm” (HuggingFace) so typical supply chain woes exist
- Low context window
- Best for training
- Fixed cost options exist

LLM Exploration

LLM	Context Window	Input Modalities	Accuracy for AppSec Tasks
GPT-4o	128k	Text/Image	High
Claude 3/3.5	200k	Text	High
Gemini 1.5	1 million	Text/Image	*High
CodeLLaMa	4k	Text	Low
LLaMa	4k	Text	Low
Mixtral/Mistral	4k	Text	Moderate

* denotes not personally used but research performed

Chatbot Assistant



Chatbot Assistant

Common AppSec Use Cases:

1. Feedback mechanism for LLM analysis
2. Assistant on Slack to help answer first level triage questions
3. Answer coding related questions directly from GitHub/GitLab/etc.



Chatbot Assistant

Overview of Steps:

1. Load PDF into LangChain's "Document" format
2. Create embeddings and store in vector database
3. Build Prompts
4. Chain components together and reference docs to answer questions

Chatbot Assistant

Supported Loaders: PDFs, Markdown, HTML, File Directories, JSON, and more...

Comprehensive list -

https://python.langchain.com/docs/integrations/document_loaders/

*You've actually already used a loader already 😊



Exercise: Chatbot Assistant

Exercise: Chatbot Assistant

⌚Objective: Load and save document, use as the source of knowledge for a chatbot, and answer questions while improving the existing prompt

Instructions:

1. Use `for_instructors/load_guide.py` as inspiration and langchain “Document loaders” documentation for guidance
2. Create or download a document you would like to answer questions based on
3. Store in `../vector_databases` utilizing a name of your choice
4. Utilize and modify `scripts/chatbot.py` for:
 - a. Improving the prompt
 - b. Load the vector database containing the document you wish to reference
5. Your chatbot, as much as is possible, should always respond with “Contact the security team” if the answer is not directly in the document

100 BONUS POINTS: 100

Try to eliminate prompt injection

Exercise: Chatbot Assistant

Retrospective:

- What types of documents did you utilize?
- What challenges did you face trying to limit your LLM to referencing only what is in the document?
- What prompt techniques did you utilize to improve the outcome?
- Anything else to share?

*Did anyone try to create a chat history?

Source Code Analysis



Source Code Analysis

Lessons Learned:

- Broad questions get fuzzy answers
- Pull Requests/Git Patches: Incorporate both the git patch and full file contents into analysis
- Few Shot Prompting helps guide the LLM
- Mixing deterministic with probabilistic is the best
- Context!
 - Tell the LLM what type of language, library, or web framework (relevant to the type of analysis)
 - When the LLM is blind to a specific tech stack, RAG can get you the rest of the way (so long as the language is supported)

Source Code Analysis

Types of Source Analysis

- Free form: Allow the LLM to reason about what type of vulnerabilities it thinks it is detecting, provide guidance, potential for tuning
- Deterministic: Look for a specific vulnerability, ask discrete code specific questions, use tuning methods

Exercise:

Source Code Review

(Free Form)

Exercise: Source Code Review (Free Form)

🎯 **Objective:** Iterate over all relevant source code files, store results, and then parse thru results to identify highest impact vulnerabilities

Instructions:

1. Either use vtm or choose an open source project from GitHub
2. Open `scripts/sca_repo.py` and modify the necessary bits like the database name and types of file extensions (ours defaults to python)
3. Run the script, allow the results to populate into your FAISS db
4. Now, modify `scripts/sca_repo_analysis.py` so that it loads the database of scan results
5. Experiment with prompts and providing contextual clues to improve the results

Exercise: Source Code Review (Free Form)

Retrospective:

- What worked?
- What did not?
- What changes did you make that assisted in improving the results?

Exercise:

Source Code Review

(Deterministic-esque)

Exercise: Source Code Review (Deterministic-esque)

🎯 **Objective:** Iterate over all relevant source code files looking specifically for a specific category of a vulnerability (default is SQL Injection) using Few Shot Prompting to improve accuracy of results

Instructions:

1. Either use vtm or choose an open source project from GitHub
2. Open **scripts/sca_deterministic_few_shot.py**
3. Look for another category of vulnerability, modify few shot prompt examples, question, and prompt to discover the type of vulnerability you are looking for

💡 Try without few shot prompts first, then, use few shot prompting to improve the results

Exercise: Source Code Review (Deterministic-esqu)

Retrospective:

- What worked?
- What did not?
- What changes did you make that assisted in improving the results?
- Anything else interesting to share?

Agentic AI



Agentic AI

- Refers to providing the LLM with tools it can use as well as a plan to use those tools (ie - “agents”)
- Overcomes:
 - Accuracy issues
 - Context Window Lengths
- You can use the built-in agent options with LangChain
- Platforms exist for production applications to create complex reasoning structures/paths

Agentic AI

Langchain built-in Agents (VERSION 1.0!)

https://python.langchain.com/v0.1/docs/modules/agents/agent_types

Agent Type	Intended Model Type	Supports Chat History	Supports Multi-Input Tools	Supports Parallel Function Calling	Required Model Params	When to Use	API
Tool Calling	Chat	✓	✓	✓	tools	If you are using a tool-calling model	Ref
OpenAI Tools	Chat	✓	✓	✓	tools	[Legacy] If you are using a recent OpenAI model (1106 onwards). Generic Tool Calling agent recommended instead.	Ref
OpenAI Functions	Chat	✓	✓		functions	[Legacy] If you are using an OpenAI model, or an open-source model that has been finetuned for function calling and exposes the same functions parameters as OpenAI. Generic Tool Calling agent recommended instead	Ref

Agentic AI

Langchain built-in Agents (VERSION 1.0!)

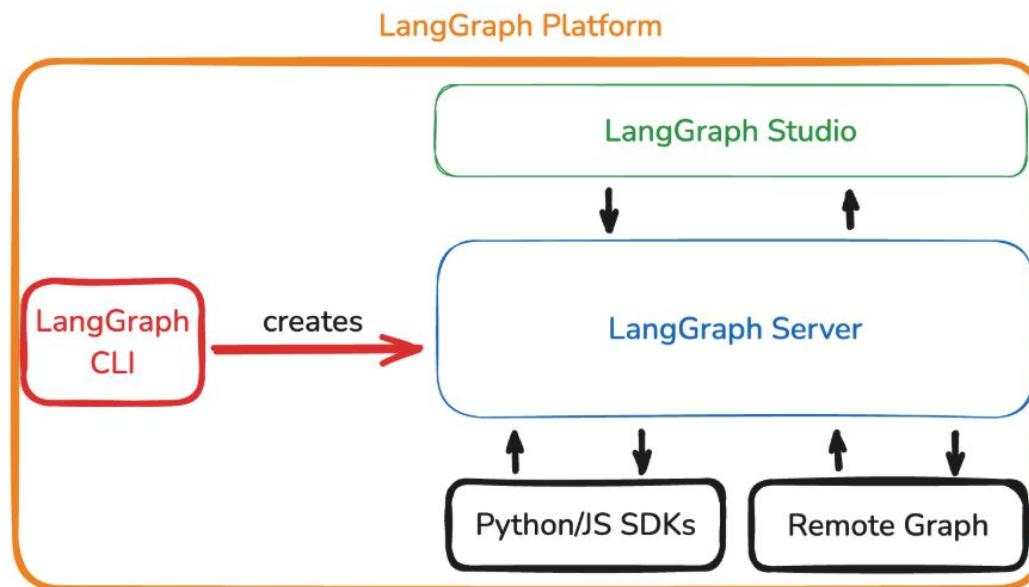
https://python.langchain.com/v0.1/docs/modules/agents/agent_types

XML	LLM	<input checked="" type="checkbox"/>			If you are using Anthropic models, or other models good at XML	Ref
Structured Chat	Chat	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		If you need to support tools with multiple inputs	Ref
JSON Chat	Chat	<input checked="" type="checkbox"/>			If you are using a model good at JSON	Ref
ReAct	LLM	<input checked="" type="checkbox"/>			If you are using a simple model	Ref
Self Ask With Search	LLM				If you are using a simple model and only have one search tool	Ref

Agentic AI

⚠️ Langchain is trying to move folks to LangGraph

ℹ️ We utilize ReAct



Agentic AI

LangGraph Benefits ->

Here are some common issues that arise in complex deployments, which LangGraph Platform addresses:

- **Streaming Support:** As agents grow more sophisticated, they often benefit from streaming both token outputs and intermediate states back to the user. Without this, users are left waiting for potentially long operations with no feedback. LangGraph Server provides [multiple streaming modes](#) optimized for various application needs.
- **Background Runs:** For agents that take longer to process (e.g., hours), maintaining an open connection can be impractical. The LangGraph Server supports launching agent runs in the background and provides both polling endpoints and webhooks to monitor run status effectively.
- **Support for long runs:** Vanilla server setups often encounter timeouts or disruptions when handling requests that take a long time to complete. LangGraph Server's API provides robust support for these tasks by sending regular heartbeat signals, preventing unexpected connection closures during prolonged processes.
- **Handling Burstiness:** Certain applications, especially those with real-time user interaction, may experience "bursty" request loads where numerous requests hit the server simultaneously. LangGraph Server includes a task queue, ensuring requests are handled consistently without loss, even under heavy loads.
- **Double Texting:** In user-driven applications, it's common for users to send multiple messages rapidly. This "double texting" can disrupt agent flows if not handled properly. LangGraph Server offers built-in strategies to address and manage such interactions.
- **Checkpointers and Memory Management:** For agents needing persistence (e.g., conversation memory), deploying a robust storage solution can be complex. LangGraph Platform includes optimized [checkpointers](#) and a [memory store](#), managing state across sessions without the need for custom solutions.
- **Human-in-the-loop Support:** In many applications, users require a way to intervene in agent processes. LangGraph Server provides specialized endpoints for human-in-the-loop scenarios, simplifying the integration of manual oversight into agent workflows.

Agentic AI

(Reason + Act) = ReAct

Explanation <https://react-lm.github.io/>

Research material

<https://research.google/blog/react-synergizing-reasoning-and-acting-in-language-models/>

ReAct (Reason + Act):

ReAct is a framework for task-solving in LLMs that combines reasoning and action. It encourages the model to first "reason" by explaining its thought process, then "act" by performing an action or generating a response. This iterative process helps break down complex problems, verify assumptions, and refine outputs for accurate and logical results.

Agentic AI

```
class SearchInput(BaseModel):
    query: str = Field(description="should be a search query")

class CustomSearchTool(BaseTool):
    name: str = "custom_search" # Explicitly type all attributes
    description: str = "Useful for when you need to answer questions about code"
    args_schema: Type[SearchInput] = SearchInput # Explicitly set schema

    def _run(
        self, query: str, run_manager: Optional[CallbackManagerForToolRun] = None
    ) -> str:
        """Use the tool."""
        # Repo details to clone
        faiss_db_path = "../vector_databases/vtm_faiss"
        db = FAISS.load_local(
            faiss_db_path,
            BedrockEmbeddings(model_id='amazon.titan-embed-text-v1'),
            allow_dangerous_deserialization=True
        )
        val = db.similarity_search(query)
        return val
    #return "LangChain"

async def _arun(
    self, query: str, run_manager: Optional[AsyncCallbackManagerForToolRun] = None
) -> str:
    """Use the tool asynchronously."""
    raise NotImplementedError("custom_search does not support async")
```

```
tools = [CustomSearchTool()]
```

Agentic AI

```
You have access to the following tools:
```

```
{tools}
```

```
To use a tool, please use the following format:
```

```
...
```

```
Thought: Do I need to use a tool? Yes
```

```
Action: the action to take, should be one of [{tool_names}]
```

```
Action Input: the input to the action
```

```
Observation: the result of the action
```

```
...
```

```
When you have a response to say to the Human, or if you do not
```

```
...
```

```
Thought: Do I need to use a tool? No
```

```
Final Answer: [your response here]
```

```
...
```

```
Begin!
```

```
Previous conversation history:
```

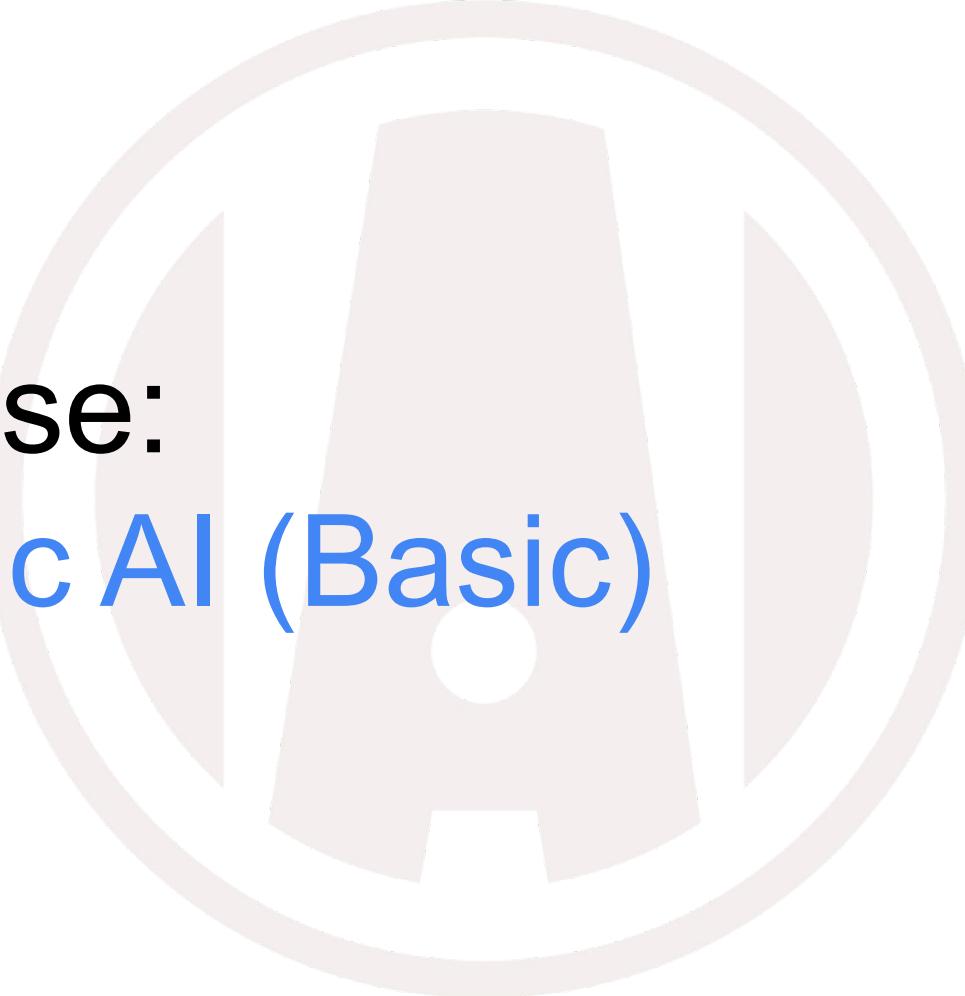
```
{chat_history}
```

```
New input: {input}
```

```
{agent_scratchpad}
```

```
agent_executor.invoke({"input": input, "chat_history": ""})
```

The `{agent_scratchpad}` placeholder in the prompt is a dynamic field that gets replaced with the intermediate reasoning or action steps that the agent generates during its thought process.



Exercise: Agentic AI (Basic)

Exercise: Agentic AI (Basic)

🎯 Objective:

Instructions:

1. Open `scripts/agentic_basic.py`
2. Run `scripts/agentic_basic.py` and note any inconsistencies
3. Make the minimum following modifications until consistent results:
 - a. Context
 - b. Prompt
 - c. Few Shot Prompt Examples

Exercise: Agentic AI (Basic)

Retrospective:

- What worked?
- What did not?
- What changes did you make that assisted in improving the results?
- Anything else interesting to share?

Testing LLMs



Testing LLMs - Background

- LLMs used to build out and run integration and unit tests for tuning
- Create repeatable results if and when you change variables that impact the LLM results.
- Even small changes in context, prompts, etc. can have big impacts
- Important if you want to change/swap/upgrade LLMs

Testing Framework

1. Use PyTest
2. Other options exist including:
 - a. <https://python.langchain.com/docs/concepts/testing/>
 - b. <https://github.com/openai/evals>
 - c. <https://github.com/confident-ai/deepeval>

Testing Framework

4 important components:

1. Makefile (specify what “make test” does)
2. pytest.ini (specify test markers)
3. Test file (ex: `agenetic_test.py`) ends in
“_test.py”
4. Use the marker `@pytest.mark.integration`
for integration tests

```
@pytest.mark.integration
def test_analyze_code_insecure():
```

*Note that we use unit-test to test structures
(json, yaml), and integration for behaviors

Exercise: Testing

Exercise: Testing

🎯 Objective:

Instructions:

1. Open `scripts/agentic_test.py`
2. Add an integration test (hint: you can make one that checks for `insecure == false`)
3. Bonus points - Create a unit-test that checks JSON or other structure for validity

Exercise: Agentic AI (Basic)

Retrospective:

- Any issues?
- Did anyone create a unit-test?
- Anything else to share/ask?

Threat Modeling



Threat Modeling with LLMs

- Creativity built in to LLMs is useful for Threat Modeling
- Initially just used Chat GPT to ask questions about threat, utilizing knowledge gleaned during initial Source Code Review process.
- Fell back on STRIDE and DREAD for classification
- Then it was codified.

Threat Modeling - Stride GPT

“STRIDE GPT is an AI-powered threat modelling tool that leverages Large Language Models (LLMs) to generate threat models and attack trees for a given application based on the STRIDE methodology.”

<https://github.com/mrwadam/s/stride-gpt>



Threat Modeling - Stride GPT



STRIDE GPT

AI-powered threat modeling

How to use STRIDE GPT

Select your preferred model provider: [?](#)

Ollama

Select the model you would like to use:

llama3.1:latest

Enter your GitHub API key (optional): [?](#)

..... 

Threat Model [Attack Tree](#) [Mitigations](#) [DREAD](#) [Test Cases](#)

A threat model helps identify and evaluate potential security threats to applications / systems. It provides a systematic approach to understanding possible vulnerabilities and attack vectors. Use this tab to generate a threat model using the STRIDE methodology.

Enter GitHub repository URL (optional)

https://github.com/owner/repo

② Select the application type

Web application

Describe the application to be modelled

Enter your application details...

② What is the highest sensitivity level of the data processed by the application?

Top Secret

Is the application internet-facing?

Yes

What authentication methods are supported by the application?

Choose an option

Deploy

Threat Modeling - Stride GPT - BHIMA

Enter GitHub repository URL (optional)

<https://github.com/owner/repo>

Describe the application to be modelled

BHIMA

=====

BHIMA is a free, open source accounting and hospital information management system

(HIMS) tailored for rural hospitals in the Democratic Republic of the Congo. We are an international team based all over the world.

BHIMA is an acronym for _basic hospital information management application_. It was originally developed by [IMA World Health](<https://imaworldhealth.org/>) with funding from the Foreign Commonwealth and Development Office (FCDO).

⑦ Select the application type

Web application

⑦ What is the highest sensitivity level of the data processed by the application?

Confidential

Is the application internet-facing?

Yes

What authentication methods are supported by the application?

Basic x

Generate Threat Model

Threat Modeling - Stride GPT - BHIMA Threats

Threat Model

Threat Type	Scenario	Potential Impact
Spoofing	Attacker creates a fake login page and steals users' credentials by manipulating the Basic authentication method used in BHIMA.	Unauthorized access to sensitive data, including patient and financial information.
Spoofing	An attacker sends a phishing email to hospital administrators with a link to a fake version of the BHIMA login page, tricking them into entering their credentials.	Compromise of sensitive data and potential malware infection.
Tampering	An attacker intercepts and modifies patient billing information while it's being transmitted between the client and server, causing incorrect charges to be applied.	Financial loss for patients and hospitals due to tampered data.
Tampering	A malicious user hacks into the Redis session management system and alters the login credentials of other users, allowing them to access sensitive areas of the application.	Unauthorized access to sensitive areas of the application and potential data tampering.
Reputation	A user makes changes to their account settings without actually performing the actions, then denies responsibility for the alterations when confronted about it.	Loss of trust in the system and potentially, legal issues due to lack of accountability.
Reputation	An attacker uses BHIMA's reporting plugins to fabricate fake reports, then denies responsibility for the malicious activity.	Loss of trust in the system and potentially, legal issues due to lack of accountability.
Information Disclosure	BHIMA is hacked, and sensitive data such as patient records or financial information is leaked online.	Compromise of sensitive data and potential reputational damage for the hospital and BHIMA.

Threat Modeling - Stride GPT - BHIMA Threats

Information Disclosure	BHIMA is hacked, and sensitive data such as patient records or financial information is leaked online.	Compromise of sensitive data and potential reputational damage for the hospital and BHIMA.
Information Disclosure	An attacker exploits a vulnerability in BHIMA's code to gain access to server logs, revealing sensitive information about user activity.	Compromise of sensitive data and potential reputational damage for the hospital and BHIMA.
Denial of Service	A large number of requests flood the server, causing it to become unresponsive and preventing legitimate users from accessing their accounts.	Loss of productivity due to system downtime and potential financial loss for hospitals.
Denial of Service	An attacker uses a botnet to send numerous requests to BHIMA's server, overwhelming it with traffic and causing it to crash.	System downtime, financial loss due to lost productivity, and potential reputational damage.
Elevation of Privilege	A user gains elevated privileges within BHIMA by exploiting a vulnerability in the NodeJS server-side code, allowing them to access sensitive areas of the application.	Unauthorized access to sensitive areas of the application and potential data tampering.
Elevation of Privilege	An attacker discovers a backdoor in BHIMA's Redis session management system, allowing them to gain elevated privileges within the application.	Unauthorized access to sensitive areas of the application and potential data tampering.

Improvement Suggestions

- Implement additional security measures for user authentication, such as multi-factor authentication or more secure password storage.
- Regularly update dependencies in the NodeJS code to prevent vulnerabilities from being exploited by attackers.
- Use a Web Application Firewall (WAF) to monitor and filter incoming traffic to the application, preventing common web attacks.

Threat Modeling - Stride GPT - BHIMA Attack Trees

Generate Attack Tree

Attack Tree Code:

```
graph TD
    A["Compromise Web Application (BHIMA)"] --> B("Spoofing")
    B --> C["Impersonate User Credentials"]
    B --> D["Bypass Basic Authentication"]

    A --> E("Tampering")
    E --> F["Modify Sensitive Data in Transit"]
    E --> G["Alter Data in MySQL Database"]

    A --> H("Repudiation")
    H --> I["Log Tampering in NodeJS Server"]
    H --> J["Modify Redis Session Data"]

    A --> K("Information Disclosure")
    K --> L("Data Leak via HTTP")
    K --> M("Sensitive Data Exposure in MySQL")

    A --> N("Denial of Service (DoS)")
    N --> O["Overload Redis with Sessions"]
```

Threat Modeling - Stride GPT - Mitigations

Threat Mitigation Strategies

Threat Type	Scenario	Suggested Mitigation(s)
Spoofing	Attacker creates a fake login page and steals users' credentials by manipulating the Basic authentication method used in BHIMA.	Implement a robust authentication mechanism, such as OAuth or OpenID Connect, to reduce reliance on Basic authentication. Utilize Secure Sockets Layer/Transport Layer Security (SSL/TLS) certificates for encryption and identity verification.
Spoofing	An attacker sends a phishing email to hospital administrators with a link to a fake version of the BHIMA login page, tricking them into entering their credentials.	Educate users on phishing techniques and provide training on identifying suspicious emails and links. Implement email authentication mechanisms like SPF, DKIM, or DMARC to prevent spoofed emails.
Tampering	An attacker intercepts and modifies patient billing information while it's being transmitted between the client and server, causing incorrect charges to be applied.	Use Transport Layer Security (TLS) encryption for data in transit and verify certificates of all parties involved. Implement a digital signature scheme, such as HMAC or ECDSA, to ensure data integrity during transmission.
Tampering	A malicious user hacks into the Redis session management system and alters the login credentials of other users, allowing them to access sensitive areas of the application.	Regularly update Redis sessions with secure random values, ensuring each user has unique credentials. Implement Redis security features like password authentication, ACLs (Access Control Lists), or REDIS_CLUSTER to prevent unauthorized access.

Threat Modeling - Stride GPT - DREAD Risk Assessment

Threat Type	Scenario	Damage Potential	Reproducibility	Exploitability	Affected Users	Discoverability	Risk Score
Spoofing	Attacker creates a fake login page and steals users' credentials by manipulating the Basic authentication method used in BHIMA.	9	8	7	10	6	8.00
Spoofing	An attacker sends a phishing email to hospital administrators with a link to a fake version of the BHIMA login page, tricking them into entering their credentials.	8	7	6	9	5	7.00
Tampering	An attacker intercepts and modifies patient billing information while it's being transmitted between the client and server, causing incorrect charges to be applied.	8	6	5	9	7	7.00
Tampering	A malicious user hacks into the Redis session management system and alters the login credentials of other users, allowing them to access sensitive areas of the application.	9	8	7	10	6	8.00
Repudiation	A user makes changes to their account settings without actually performing the actions, then denies responsibility for the alterations when confronted about it.	7	5	4	8	6	6.00
Repudiation	An attacker uses BHIMA's reporting plugins to fabricate fake reports, then denies responsibility for the malicious activity.	8	7	6	9	5	7.00

Threat Modeling - Stride GPT - Test Cases

Spoofing - An attacker sends a phishing email to hospital administrators with a link to a fake version of the BHIMA login page, tricking them into entering their credentials

****Test Case: Phishing Email****

Given an attacker has sent a phishing email with a link to a fake BHIMA login page to a hospital administrator
And the administrator clicks on the link and enters their credentials
Then the attacker should have obtained the administrator's credentials

Tampering - An attacker intercepts and modifies patient billing information while it's being transmitted between the client and server, causing incorrect charges to be applied

****Test Case: Modified Patient Billing Information****

Given an attacker has intercepted patient billing information being transmitted between the client and server
And the attacker modifies the billing information to reflect incorrect charges
Then the system should apply the incorrect charges to the patient's account

Tampering - A malicious user hacks into the Redis session management system and alters the login credentials of other users, allowing them to access sensitive areas of the application



Exercise: Threat Model

Exercise: Threat Model

⌚ Objective: Create a Threat Model of an Open Source Project

Instructions:

1. Clone and follow instructions for running Stride GPT (<https://github.com/mrwadams/stride-gpt>)
2. Choose an open source project from Github (e.g. https://github.com/railsbridge/bridge_troll)
3. Copy and paste the raw README file from the chosen project into Stride GPT.
4. ***** Note, if you navigate away from a tab it clears out any generated content. Copy/Paste if you want to save it.**
5. Utilize Ollama + llama3.1 to generate Threats, Attack Trees (sometimes), Mitigations, Risk Assessment and Test Cases
6. Utilize Ollama + mistral to generate Threats, Attack Trees (sometimes), Mitigations, Risk Assessment and Test Cases
7. Review output for appropriateness, analysis, note good and bad

Exercise: Threat Model

Retrospective:

- What worked?
- What did not?
- Which model worked better for Threats, Attack Trees, Mitigations, DREAD, and Test Cases?