# Ask Captain Cyber
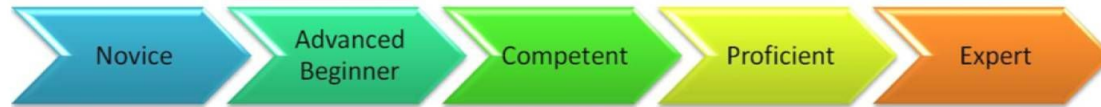
Team: SDMay25-07

Client & Advisor: Dr. Doug Jacobson

# Problem Statement

- Not easy to convey cybersecurity topics to general population

- Information on the web is not always accurate

- Everyone should have basic cyber knowledge and a convenient, reliable way to learn it

| Novice | Advanced Beginner | Competent | Proficient | Expert |

# Solution: Ask Captain Cyber

- AI-integrated chatbot to ask cybersecurity-related questions

  - Hosted on Iowa State University's servers

- Answers are proofread by cybersecurity professionals

- Contains database table of previously asked questions

  - Ensures reliable information

- Free to use

# Users

Cybersecurity Amateurs

- Little cybersecurity knowledge
- Asks basic questions

Cybersecurity Enthusiasts

- Decent cybersecurity knowledge
- Asks complex questions

Cybersecurity Experts

- Substantial knowledge
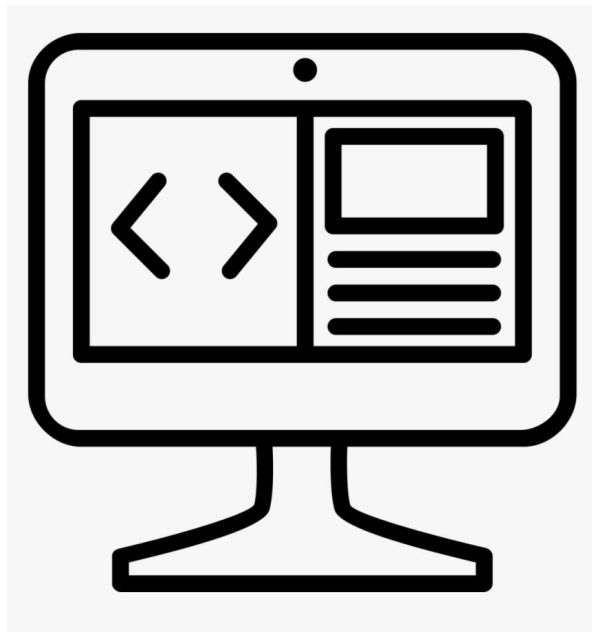- Answers questions from previous two user groups

# Design Requirements

- Must accurately interpret and respond to user queries

- Experts must be able to easily vet initial AI response

- Store vetted answers in database for future usage

- Maintain overall security posture for login, user information, data storage/transfer

# Our Design

- Frontend
    - WordPress - Content management
    - PHP, HTML, JS - Frontend design

- Backend
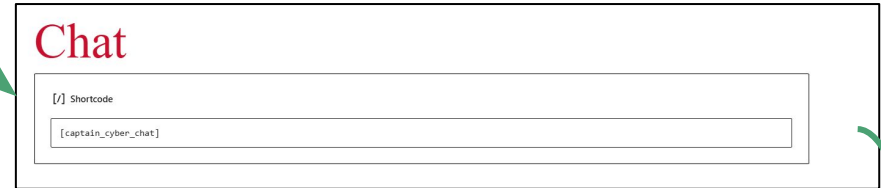    - Flask Rest API - Middleware
    - OpenAI API - LLM
    - MySQL - Database

- Version control
    - Github

# Design - Frontend

- WPCode plugin implements "code snippets" directly into site

- PHP, HTML, and JS connects with MySQL and Flask to populate pages

- Promotes easy updating and troubleshooting as needed in the future

# Design - Accessing Captain Cyber



Ask Captain Cyber

Search for existing answers.

Search question

Can't find the answer: Ask me the question!



**Ask Captain Cyber!**

Click here to chat with Captain Cyber to learn about any cybersecurity-related topics!

**Ask Captain Cyber**

# Design - Chat



**IOWA STATE UNIVERSITY**

Ask Captain Cyber

## Chat with Captain Cyber to learn more about Cybersecurity!

what is phishing?

Phishing is a type of cyber attack where attackers impersonate legitimate organizations or individuals to trick people into revealing sensitive information, such as usernames, passwords, or credit card numbers. This is usually done via deceptive emails, websites, or messages. This answer was vetted by a professional.

Ask your question...    Send

# Design - Vetting Dashboard

# Design - Account Management

# Design - Account Management

# Design - Backend implementation

- Flask Rest API
    - Variety of API endpoints
    - Efficient user, database, and AI workflow
    - Modular design for scalability
- Database
    - Questions, answers, category, status, tags, updated_at, created_at
    - Wide variety of querying possibilities
    - Scalable for expansion
    - Seeded the database

# Design - Backend Dataflow



## Data Flow if Query is Valid

Open API — User Input → Flask RestAPI — User Input → WordPress Frontend

Flask RestAPI — Database Query from OpenAI → Open API

Flask RestAPI — Response from Database → WordPress Frontend

Database — Response from Database → Flask RestAPI

Flask RestAPI — Database Query from OpenAI → Database

Blue: Backend Communication
Green: User Input
Backend Elements: Red
Frontend Elements: Yellow

## Data Flow if Query is Invalid

Custom AI Response if Database Query is Invalid

Open API — User Input → Flask RestAPI — User Input → WordPress Frontend

Flask RestAPI — Database Query from OpenAI → Open API

Flask RestAPI — Response from OpenAI → WordPress Frontend

User Input resent for AI Response

Invalid Indicator Response

Database Query from OpenAI

Database

Blue: Backend Communication
Green: User Input
Backend Elements: Red
Frontend Elements: Yellow

# Design - AI

- OpenAI Assistant
  - GPT-4o
  - Temperature: 0.30
  - Top P: 1.00

- Prompt Engineering
  - Role assignment
  - Zero-shot prompting
  - Behavioral conditioning
  - Iterative reframing



OpenAI API Playground

# Results

- Project effectively meets client's requirements
    - Successfully responds to user queries
    - User-friendly
    - Ethical responses
    - Provides accurate cybersecurity information

- Technical functionality
    - Accessible user interface
    - Full communication between all components
    - Rapid response time
    - Log user input in OpenAI

# Project Demo

IOWA STATE UNIVERSITY

**Ask Captain Cyber**

Resources      Adversarial News Network      Ask Captain Cyber

| Ask your question... | Send |

*Disclaimer: Captain Cyber is designed to promote cybersecurity awareness and education; any attempt to misuse, manipulate, or exploit this system for malicious purposes is strictly prohibited. Users are reminded that engaging in unauthorized cyber activities can lead to severe penalties, including criminal charges, financial liabilities, and reputational damage.*

IOWA STATE UNIVERSITY
OF SCIENCE AND TECHNOLOGY

# Value provided

- Educates users on cybersecurity concepts

- Promotes digital literacy

- Provides resources

- Raises awareness

# Next Steps

- Hand off to our client

- Additional features
  - Enhancing user interface
  - Security improvements (e.g. MFA)

- Unlikely to be part of another project
  - Not a component; ACC is the product
  - Cyber House Rock website redesign

# Conclusion