# MICHIGAN STATE UNIVERSITY

# Project Plan
## Open Source Intel

## The Capstone Experience

### Team GM

Ben Buscarino
Will Crecelius
Igli Ndoj
Qiming Ren
Taylor Zachar

Department of Computer Science and Engineering
Michigan State University

Spring 2020

*From Students…*
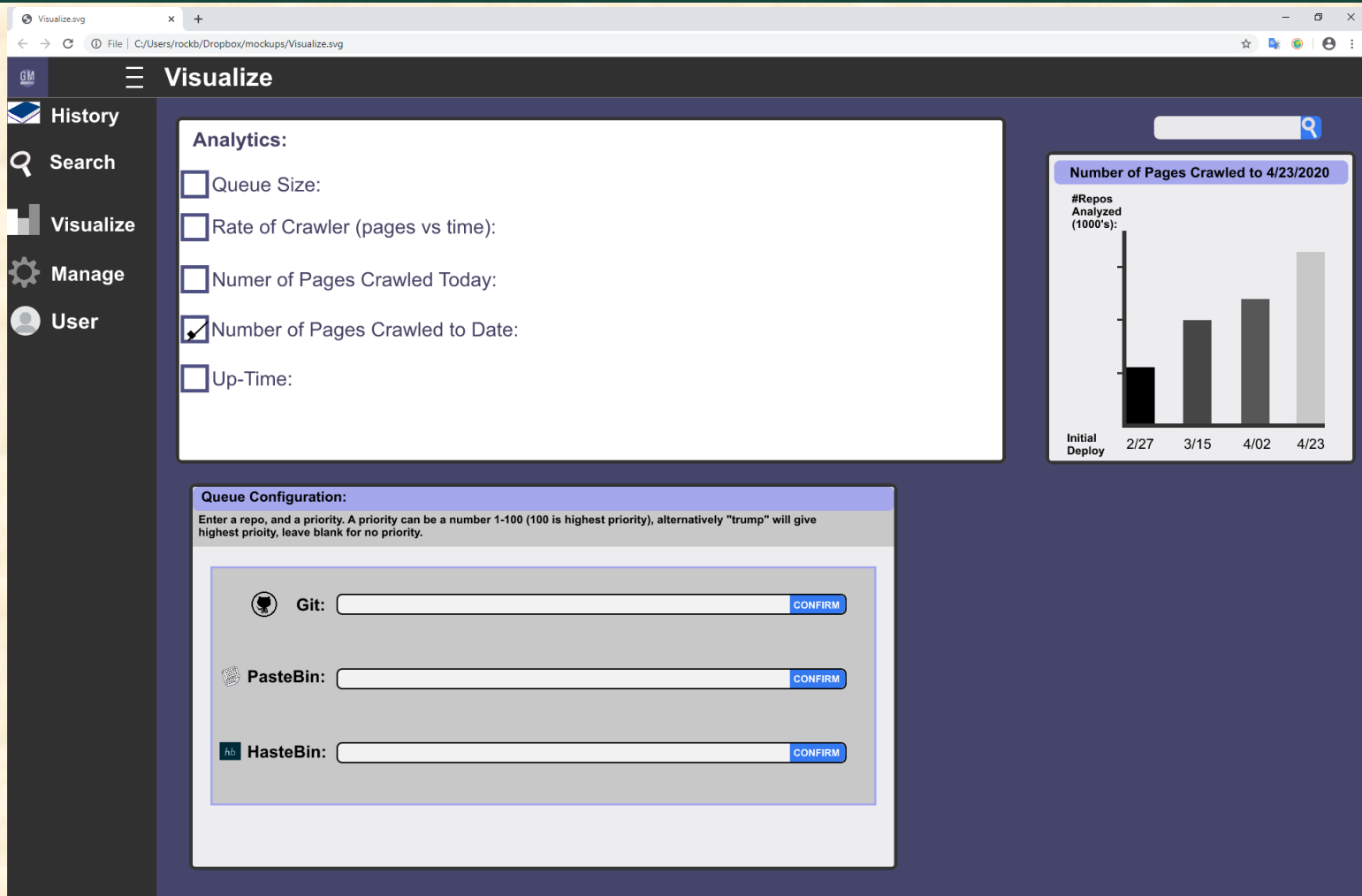   *…to Professionals*

# Functional Specifications

- GM is so large that it is at risk of having intellectual property leaked on public repositories.

- Open Source Intel will be able to find any GM owned IP on public repositories and then alert the GM team of findings

- A confidence rating will be given of the assurance that something is found on the site.
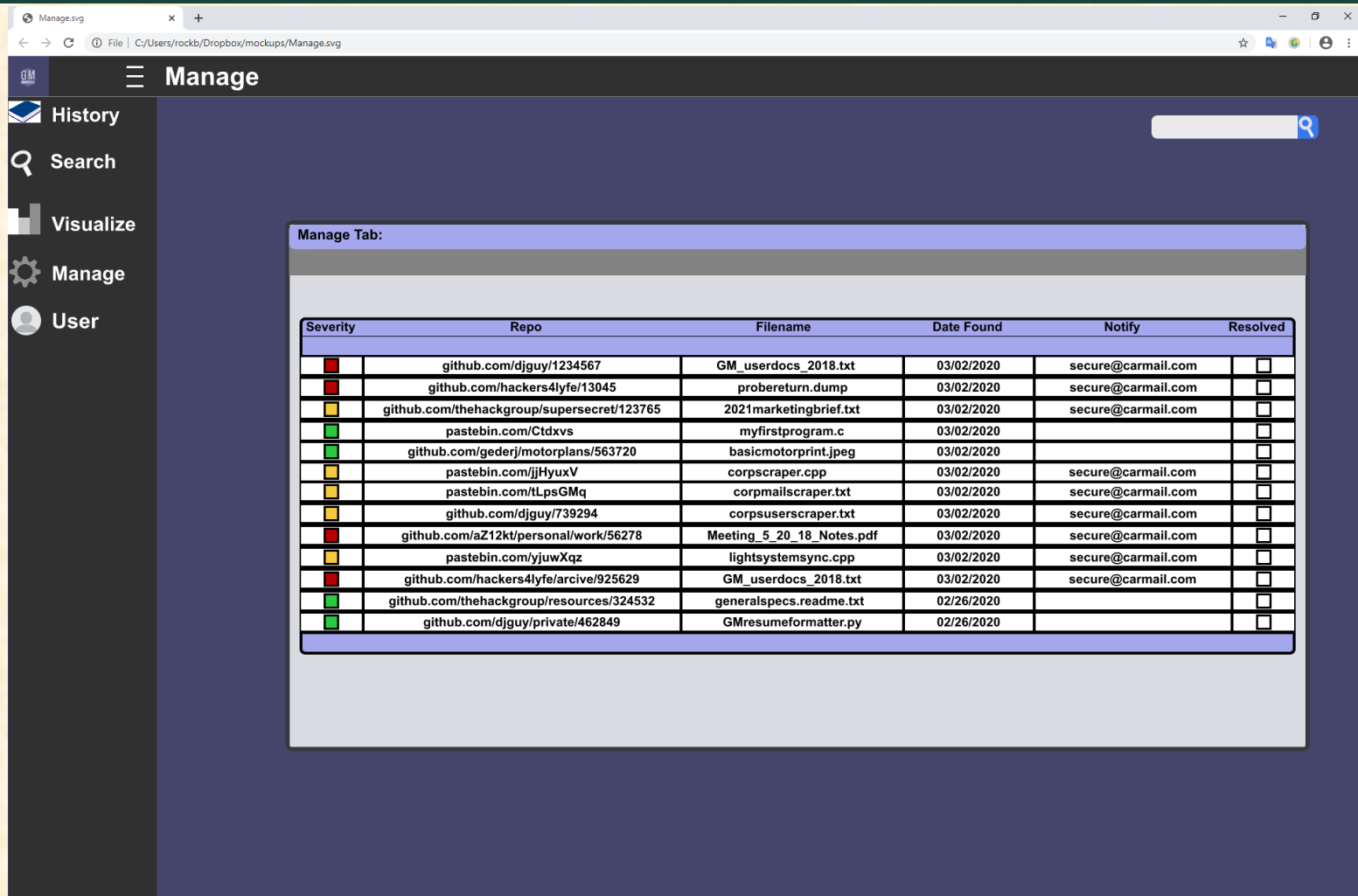
# Design Specifications

- Scrape various public repositories to identify GM property and security threats such as: usernames, API keys, and code snippets.

- User friendly web frontend to display the data collected by the backend Python scrapers.

- Once the threat has been identified and verified, GM can proceed with legal action through the website.
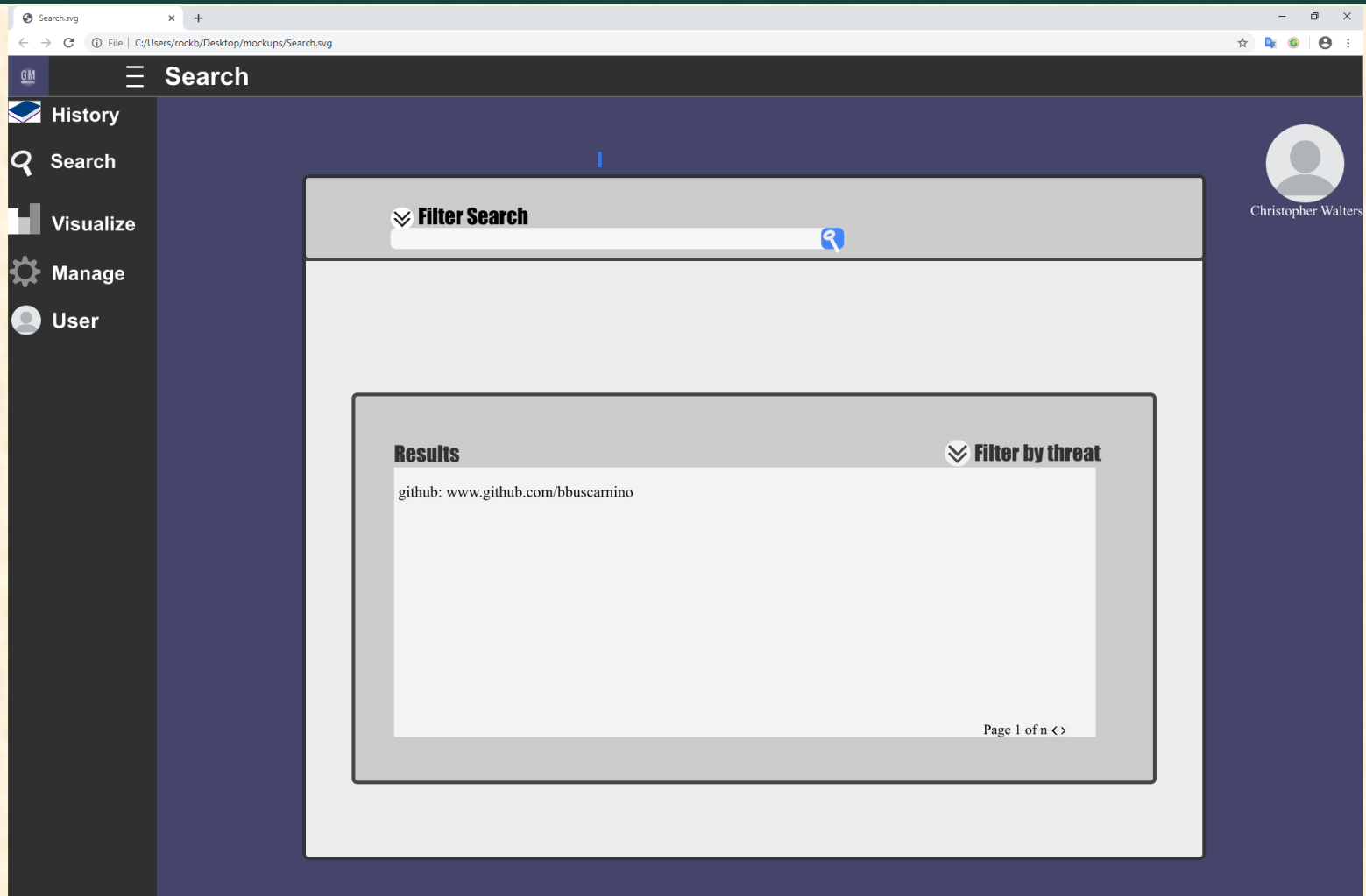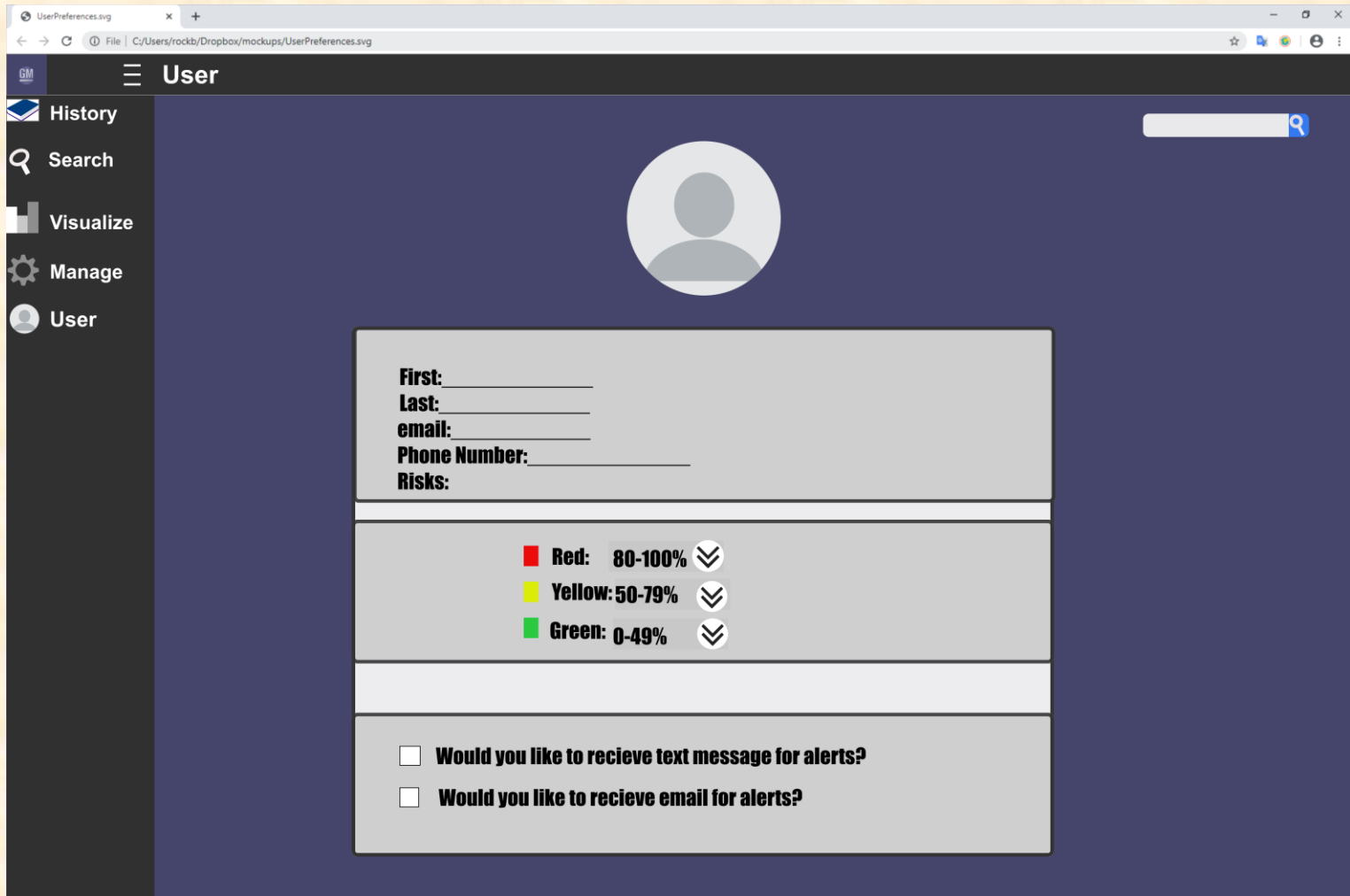
# Screen Mockups:Visualize
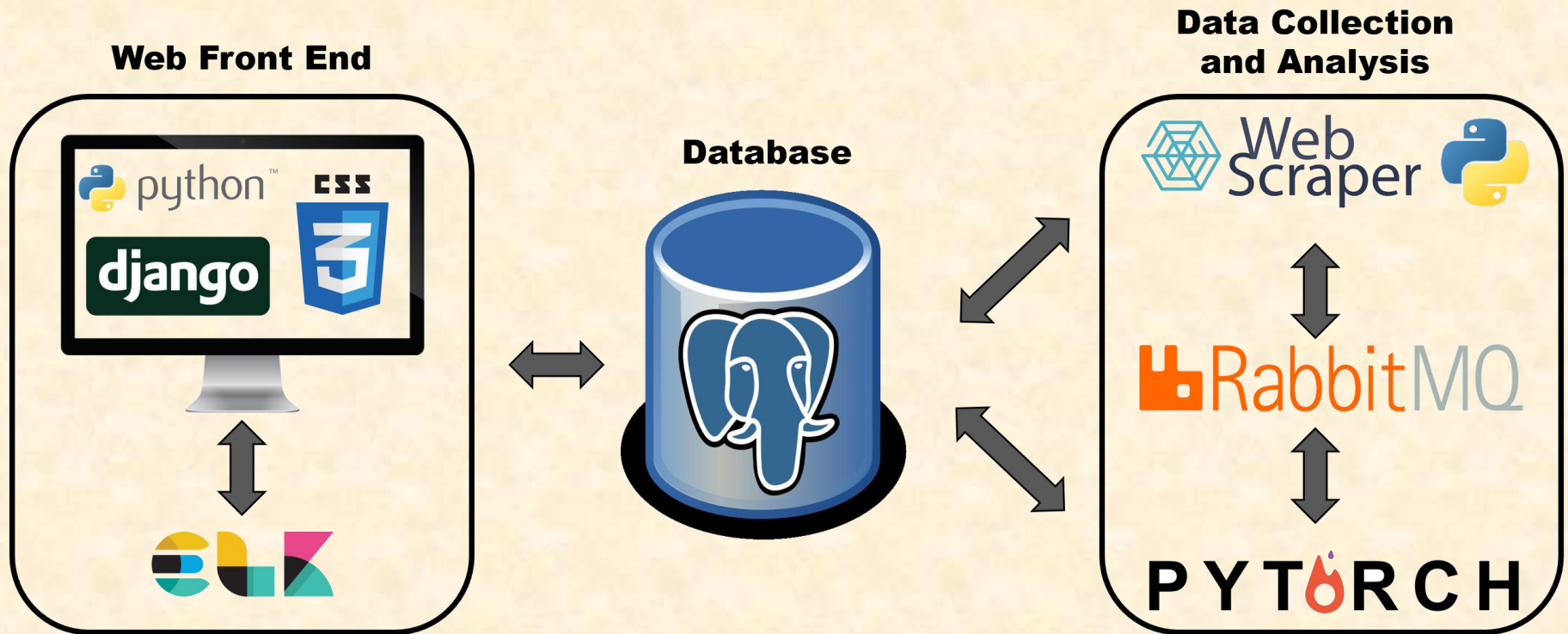
# Screen Mockups: Manage

# Screen Mockups: Search

# Screen Mockups: User-Preferences

# System Architecture



**Web Front End**

**Database**

**Data Collection and Analysis**

# Technical Specifications

- GM employees interact with a Python web application.
- The front end uses Django, a high-level Python web framework, styled with the Semantic UI CSS Framework.
- RabbitMQ serving as a work-queue between a suite of web crawlers and PyTorch machine learning service.
- PostgreSQL used for data warehousing of confidence rating, crawled websites, and details of identified leaks.

# System Components

- Hardware Platforms
  - Rack-mounted server running NixOS.
  - MacOS Catalina.
- Software Platforms / Technologies
  - Python, Django, Semantic UI.
  - PostgreSQL and RabbitMQ.
  - Web scraper and ML services using Beautiful Soup, Requests, and PyTorch.
  - Kubernetes and Azure Kubernetes Service.

# Risks

- Risk 1: Computing constraints
  - Scanning the entirety of GitHub/GitLab/Bitbucket/Pastebin plus their revision histories.
  - Mitigation: Leveraging bounding and caching much as possible to reduce the amount of computationally expensive work done.

- Risk 2: Identifying problematic content
  - GM is such a large company that they don't know all types of IP the enterprise may have exposed.
  - Mitigation:  The crawlers will need to have a large breadth for their search,. Treating findings as metaphorical "breadcrumbs" and showing them to the user in a meaningful way so that the user can finish the investigation is an encouraged solution offered by our client.

- Risk 3: Bringing machine learning to production
  - Nobody on the team has experience with machine learning.
  - Mitigation: Leverage outside resources such as MSU CSE faculty and industry contacts.

# Questions?