# MICHIGAN STATE
# U N I V E R S I T Y

# Project Plan

## Machine Learning Document Classification and Redaction

## The Capstone Experience

### Team Technology Services Group

Lazaro Cruz
Genya Dobrev
Will Giger
Jacob Harris
Xiaokuan Zhang

Department of Computer Science and Engineering
Michigan State University

Spring 2020

*From Students…*
*…to Professionals*

# Functional Specifications

- Removes sensitive personal information from documents.

- In doing so, private information will not be viewed by who is not supposed to see it.

- This is important in medical records especially.
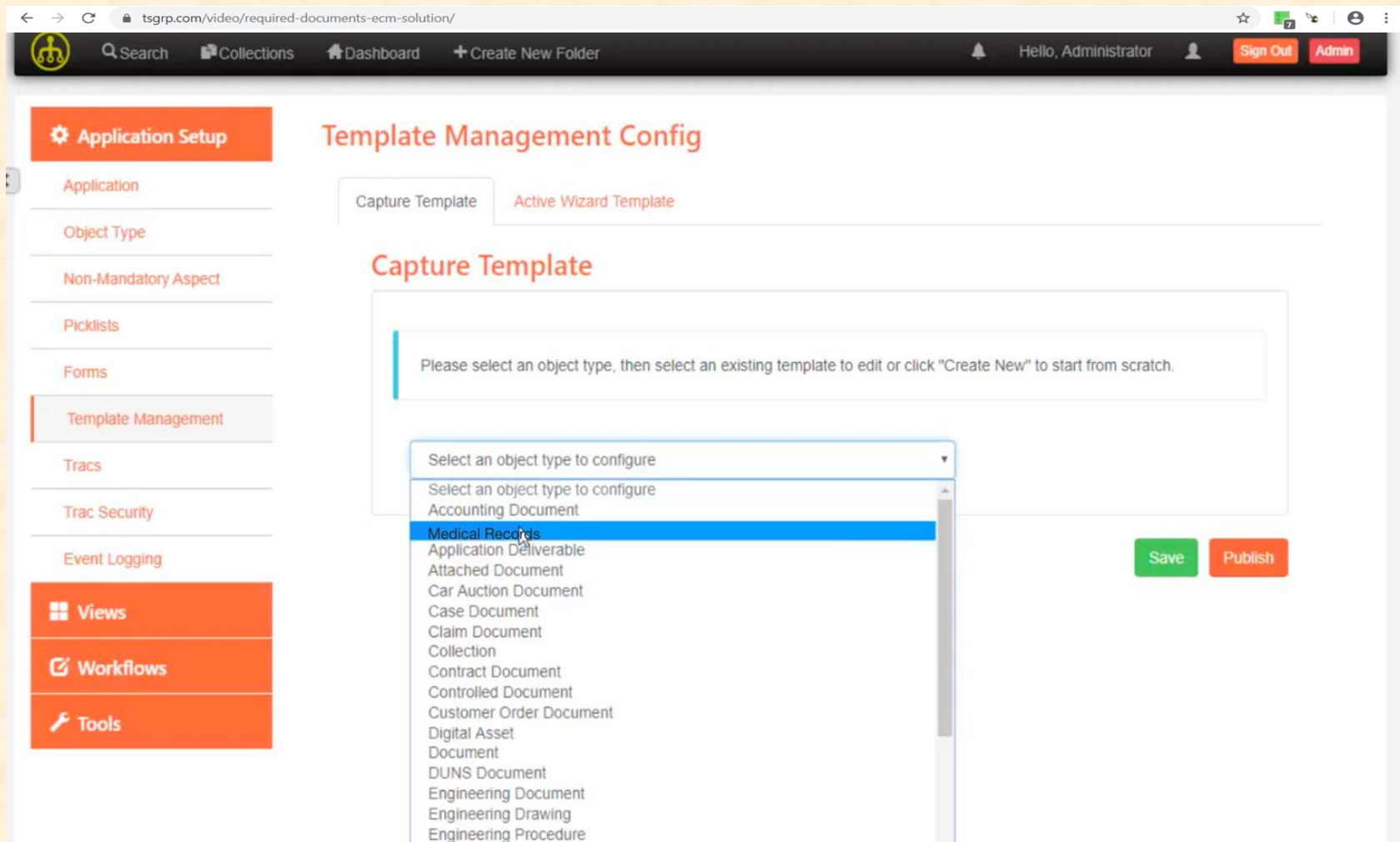
# Design Specifications

- A person should be able to upload a document through a computer.

- PII can be redacted during template configuration and fields can be un-redacted as needed.

- User is displayed redacted version during indexing of the document.

# Screen Mockup: Selecting Template

# Screen Mockup: Identifying Values

# Screen Mockup: Redacted Indexing

# Screen Mockup: Redaction Editing

# Technical Specifications

- Using Apache TomCat server to host the program.

- Uploading documents through front-end with a JavaScript application (OCMS).

- Managing back-end with Java and HBase database.

- Using Azure machine learning to recognize information and redact

# System Architecture

# System Components

- Hardware Platforms
  - Linux
  - Ubuntu
- Software Platforms / Technologies
  - Apache Tomcat Server
  - Hadoop cluster
  - OpenContent from client
  - Azure Machine Learning

# Risks

- Which Azure Machine Learning environment would work best, if any would?
  - Time consuming process of testing each environment to find the best solution.
  - Mitigation: Upfront research for the most integratable environment.
- Redaction confidence level
  - Current client software can find metadata with strong confidence.  In dealing with PII redaction the confidence level will need to be much higher.
  - Mitigation: Making sure to benchmark our Machine Learning continuously throughout development.
- What is PII exactly and how to measure it accurately
  - What information is PII exactly and how to train model to recognize it?
  - Mitigation: Worst case is doing it manually by going through documents and running it by the person redacting manually in the client's company.
- Client storage platform is still unclear
  - In last call with client it was unclear which storage platform the client would like to use.
  - Mitigation: Upfront research on Azure to fall back on.

# Questions?

? ? ? ?

? ? ? ?

? ?