# Cloud Computing Redaction Project

Michail Angelos Karvelas

# Project Overview

**The aim of this project is to establish a secure system that anonymizes sensitive documents, similar to a redaction service.**

## 1. The Goal

Create a trusted system that protects privacy by anonymizing documents.

Automatically clear names, emails, and addresses from text.

## 2. The Solution

A simple service that receives text and sends back a fully cleaned PDF.

Powered by language models that find sensitive details with context.

## 3. The Architecture

Built for the cloud with containers and Kubernetes.

Ready for high demand with smart routing and traffic control.

Sensitive content is erased before any file is generated.

# The Problem & Solution

- The Challenge: Handling sensitive documents (PII) securely is difficult. Manual redaction is slow and error-prone.
- The Solution: An automated Microservice Architecture that detects and permanently removes PII.
- Key Capabilities:
  - Intelligent Detection: Combines Regex (Email/Phone) and AI (Names/Locations).
  - Visual Redaction: Physically draws black boxes on PDFs (not just covering text).

  Scalable: Runs on a self-healing Cloud Cluster (AKS).

# System Architecture

- Frontend (Consumer):
User Interface for uploading files.

- Backend (Producer):
Python FastAPI service for OCR and NLP processing.

- Service Discovery:
Consul allows the Frontend to find the Backend dynamically.

# System Architecture

## Backend Explanation



Regex: Used for deterministic patterns (Emails, Phones). High speed, 100% accuracy for fixed formats.

NLP (spaCy): Used for context-aware entities (Names, Cities). Capable of distinguishing "Park" (Name) from "Park" (Location).

**Why Hybrid?**

Combining both covers the weaknesses of each, ensuring maximum accuracy.

# System Architecture

## Backend Code  Screenshot

# System Architecture

## FrontEnd

# Kubernetes Infrastructure (AKS)

- Key Specs:
  - Cluster: Azure Kubernetes Service (AKS)
  - Nodes: 2 Nodes (Standard_B2s) for hardware redundancy
- Networking:
  - LoadBalancer: Exposes the Frontend to the public internet

ClusterIP: Keeps Backend communication private and secure

# CI/CD Pipeline (DevOps)

- Push: Code committed to GitHub triggers the pipeline.
- Build: Docker Image created and pushed to Azure Container Registry (ACR).
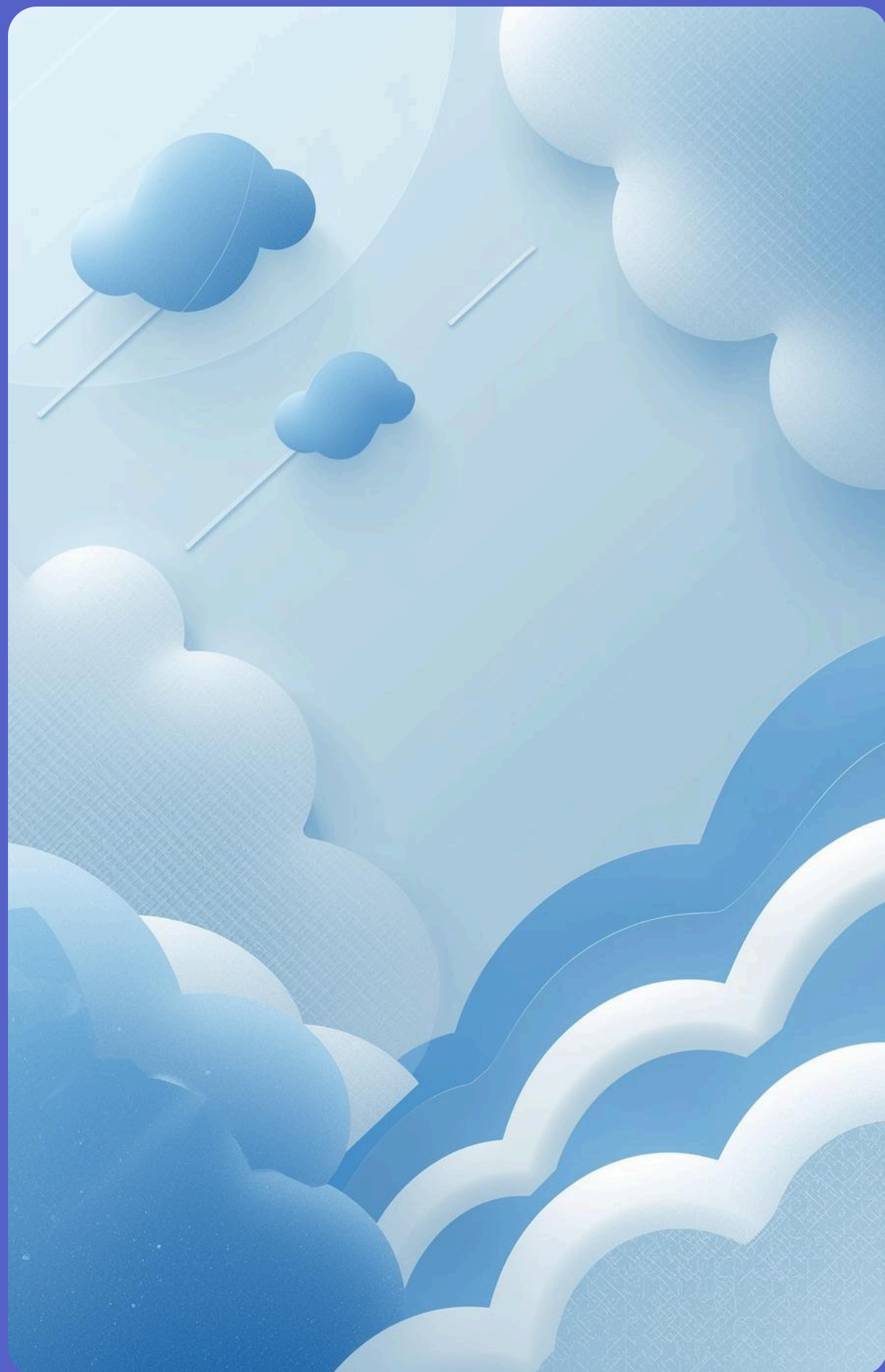- Deploy: AKS cluster pulls the new image and updates the deployment.

# Implementation Steps and Phases

Development began locally, where I built and tested the backend and frontend microservices to ensure stability. After confirming that all components functioned correctly in isolation, I deployed the full solution to Microsoft Azure, making the web application accessible via the public cloud."

# Future Improvements

- Advanced Visual Privacy (Scene Text Recognition):
  - Current State: The system works on scanned documents (white background, black text).
  - Future: Implement Object Detection (e.g., using YOLO or CRAFT) to identify and blur sensitive information in natural photos.
  - Use Case: Automatically blurring house numbers, license plates, and street signs in uploaded photos to protect location privacy.
- Multi-Language Support:
  - Future: Integrate additional spaCy models (e.g., es_core_news_sm for Spanish) to support international document redaction.
- User Authentication (RBAC):
- Future: Implement OAuth2 (Login with Google/Microsoft) so that only authorized employees can redact documents, with audit logs tracking who redacted what.

# Thank you for your attention