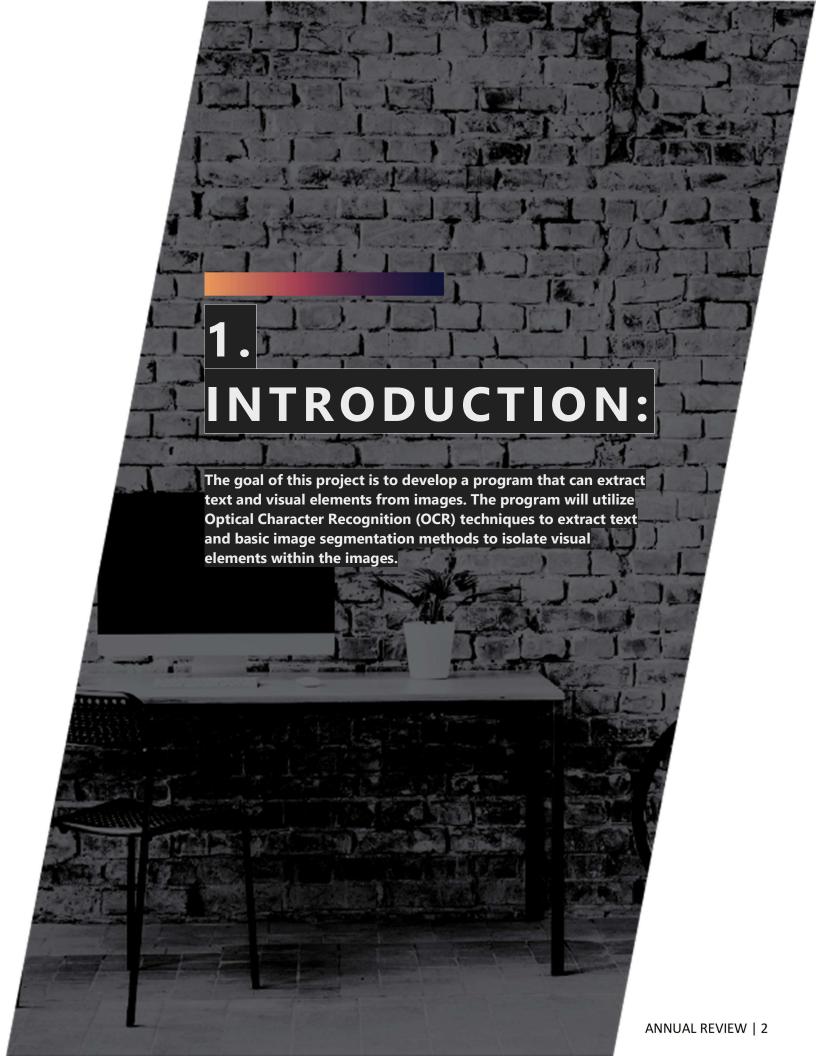
# DETAILED REPORT VISUAL CONTENT EXTRACTOR



# 2. APPROACH

## TECHNOLOGIES USED:

- a. Python: Programming language used for development.
- b. OpenCV (cv2): Library for computer vision and image processing tasks.
- c. Tesseract OCR (pytesseract): OCR engine used for text extraction.
- d. BeautifulSoup: Library for parsing HTML and XML documents.
- e. Requests: Library for making HTTP requests.
- f. Base64: Encoding library for handling image data.
- g. Heroku API: External API used for listing Heroku apps.

# IMPLEMENTATION DETAILS:

### TEXT EXTRACTION: **A**.

Utilized the pytesseract library, which is a wrapper for Tesseract OCR, to extract text from images.

Loaded the image using OpenCV and passed it to pytesseract for OCR processing.

Extracted text content from the image.

# B. VISUAL ELEMENT SEGMENTATION:

Implemented basic image segmentation techniques using OpenCV to isolate visual elements.

Converted the image to grayscale and applied thresholding to create a binary image.

Detected contours in the binary image and extracted individual visual elements by bounding rectangle.

### C. HTML GENERATION:

Used BeautifulSoup to generate HTML with extracted text content and visual elements.

Styled the HTML output with CSS for improved presentation. Saved the generated HTML to a file named 'output.html'.

### D. HEROKU APP LISTING:

Implemented a function to list Heroku apps using the Heroku API. Used the Requests library to send HTTP requests with proper authorization headers.

Processed the response to retrieve and display the list of Heroku apps.

# 4. CHALLENGES **ENCOUNTERED:**

- a. OCR Accuracy: One of the challenges was achieving high accuracy in text extraction, especially for images with complex backgrounds or low-quality text.
- b. Visual Element Segmentation: Implementing effective image segmentation techniques to accurately isolate visual elements posed a challenge, especially for images with complex layouts.
- c. Integration with Heroku API: Ensuring proper authentication and handling of API responses posed challenges during integration with the Heroku API.

# 5. CONCLUSION:

In conclusion, the developed program successfully achieves the goal of extracting text and visual elements from images. By leveraging OCR techniques and basic image segmentation, the program can accurately extract textual content and identify visual elements within images. Additionally, integration with the Heroku API provides additional functionality for listing Heroku apps. Despite some challenges encountered during development, the program provides a solid foundation for further refinement and extension.

