A detailed stained glass-style illustration of a sun rising over a series of ocean waves. The sun is large and yellow with rays extending to the right. The waves are depicted in shades of blue, white, and gold, with highlights suggesting sunlight reflecting off the water. Small glowing particles or stars are scattered throughout the scene.

COMPOSITION BOOK

OCR STEP : CODE EXPLANAT^o &

STRUCTURE

100 SHEETS
9.75 X 7.5 IN / 24.76 X 19.0 CM

PROJECT STRUCTURE

C:\CR\NewAttempt\

main.py (main orchestrator)

config.py (path & settings common to all programs)

ModelTraining.py

init.py

fix_labels.py (fix labels when they're names don't match the image's).

training_classifier.py (training model w/ our dataset)

YOLO3DBatcher.py (YOLO Detection Class)

processors/ => OCR processing pipeline.

init.py

contentDetector.py (detect content type: text, image, ID, audio, ...)

image_processor.py (image enhancement)

ocr_processor.py (text extraction)

pdf_loader.py (PDF to image).

text_processor.py (text formatting)

visualizer.py

STEP BY STEP

① MAIN.PY

workflow: Start → load document → for each page → [Process Pipeline] → Save Results → Finish

• INPUT: folder path containing PDF Document

OUTPUT: folder w/ extracted text + visualization for each PDF

• Key functions → process_pdf_file() - processes single PDF

→ process_pdf_folder() - processes all PDFs in folder

→ OutputManager class - manage file organization

② CONFIG.PY

• Control location for all files & settings. Contains

model paths
input/output folders
processing parameters
Model Training Settings

• key variables → TRAINED_MODEL_PATH = "/path/to/trained-model.pt" (our trained model, for now only hosted on my PC)
→ INPUT_FOLDER = "/path/to/input/PDFs"
→ OUTPUT_FOLDER = "/path/to/output/results"

③ ROBOTEZ TRAINING / YOLO3DDETECTOR.py

- workflow: initialize → load Model → Detect Objects → Return Results

Input: B&W data (images here)

Output: List of Detected objects w/ Bounding Boxes, classes, confidence.

- key components → prepare_dataset() - Prepare training data in YOLO format (YAML) & split dataset for training: 80/20
 - train_model() - Train model w/ parameters & data augmentation
 - detect_3d_cards() - Detect 3D cards & pictures in images.

DON'T RUN IT !

④ TRAINING - CLASSIFIER.py

- workflow: fix labels → Prepare Dataset → Train Model → Save Model

Input: Raw Images + Label file

Output: Trained Model weights (best.pt)] will call YOLO3Ddetector.py !

⑤ FIX_LABELS.py

Fix mismatched label file names & class IDs.

- workflow: Find labels → Rename to match images → convert Class IDs → Save.

e.g. 0e2f75c6-img-154.txt → img-154.txt, class 15 → 0

⑥ PROCESSORS / OCR PROCESSING PIPELINE

- workflow: PDF → Page Image → Preprocess → Detect Content → OCR → format ↓
visualize.

⑥.1 PDF LOADER.py ⇒ convert PDF pages in 2 images

- workflow: load PDF → Extract Page → Convert to OpenCV Image → Return

Input: PDF file path

Output: OpenCV image (numpy array)

⑥.2 IMAGE_PROCESSOR.py ⇒ Enhance images (* contrast) for ⚡ efficient OCR

- workflow: Grayscale → Denoise → Enhance Contrast → Sharpen → Return

Input: Raw Image

Output: Enhanced grayscale image.

(3) CONTENT_DETECTOR.PY \Rightarrow Determine Content type (3D Pic / Pic/text)

Workflow: Load YOLO model \rightarrow Detect Objects \rightarrow Classify Content \rightarrow Return Type

Outputs: "3d_card"; "picture"; "text"

(4) OCR_PROCESSOR.PY \Rightarrow Extract text from images using EasyOCR

Workflow: Load EasyOCR \rightarrow Process Image \rightarrow Extract Text \rightarrow Return Results

Input: Enhanced Grayscale image

Output: List of (Bounding Box, text, confidence) tuples.

(5) TEXT_PROCESSOR.PY \Rightarrow Format OCR results based on content type

Workflow: Check Content Type \rightarrow Format Accordingly \rightarrow Return Structured Text

E.g. Outputs \rightarrow 3D card : [3D CARD] In Extracted text... \rightarrow Text: Raw extracted text...
 \rightarrow Picture: [PICTURE] In Extracted text...

(6) VISUALIZE.PY \Rightarrow Create visualization for OCR results

Workflow: Draw Bounding Boxes \rightarrow Create Side Panel \rightarrow Combine \rightarrow Save Image

INPUT: Image + OCR result + content type

OUTPUT: Combined visualized image (original + text extract panel)

RUNNING CODE:

L \rightarrow training phase (one-time setup): fix_labels.py \rightarrow train_model.py \rightarrow YOLODetector.py \rightarrow trained_model.pt

L \rightarrow processing phase (daily use): main.py

pdf-loader.py
img_processor.py
content_detector.py
ocr_processor.py
text_processor.py
visualizer.py

DATA FLOW!

: the only ones you actually click 'run' for, they call the other ones mentioned!

PDF FOLDER
main.py selects PDF

PDF FILE
pdf-loader.py converts

PAGE IMAGES
pdf-loader.py
converts

CONTENT TYPE + BOUNDING BOXES
content-detector.py
classifies

TEXT + CONFIDENCE SCORE
ocr_processor.py
extracts

text_processor.py formats

STRUCTURED TEXT

visualizer.py creates.

VISUALIZATI+ TEXT FILE

L \rightarrow key technology * YOLOv8-Object Detec*

* NumPy - Numerical Operat*

* EasyOCR-text extract*

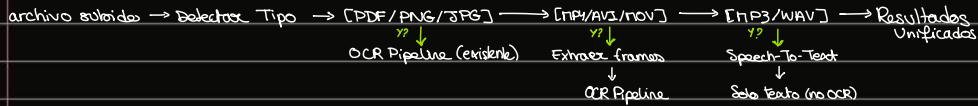
* PIL/Pillow - Image Manipulat*

* PyPDF2 (fitz) - PDF processing

* OpenCV - image processing

NOW INCLUDING SPEECH TO TEXT

→ flujo:



↳ ESTRUCTURA ACTUALIZADA

New Attempt /

main.py	(Actualizar el Orchestrator)
config.py	(Config Actualizado)
ModelTraining/	
...	

processors /

file-type-detector.py	NUEVO: Detector q tipos de Archivo es
audio-processor.py	NUEVO: Speech - To - Text
video-processor.py	NUEVO: Extract frames
pdf-reader.py	Existente
contentDetector.py	
...	

requirements.txt ⇒ AÑADIR DEPENDENCIAS AUDIO/VIDEO.

Passo 1: Actualizar config.py

```

# Anadir al config.py
SUPPORTED_FILE_TYPES = [
    'pdf': ['application/pdf', '.pdf'],
    'image': ['image/png', 'image/jpeg', '.png', '.jpe', '.jpeg'],
    'video': ['video/mp4', 'video/avi', 'video/mov', '.mp4', '.avi', '.mov'],
    'audio': ['audio/mpeg', 'audio/wav', 'audio/mka', '.mp3', '.wav', '.mka']
]

# Speech-to-text settings
WHISPER_MODEL = 'base' # 'tiny', 'base', 'small', 'medium', 'large'
LANGUAGE = 'es' # 'en', 'auto' para detección automática
    
```

Passo 2: file-type-detector.py

```

# Importar librerías
import os
from typing import List
from config import SUPPORTED_FILE_TYPES

class FileTypeDetector:
    def __init__(self):
        self._supported_types = SUPPORTED_FILE_TYPES

    def detect(self, file_path):
        """Detector tipo de archivo"""
        if not os.path.exists(file_path):
            return None

        # Por extensión
        ext = os.path.splitext(file_path)[1].lower()

        # Por MIME type
        mime_type, _ =MimeType.mime_type(ext[1:], ext)

        # Determinar categoría
        for category in SUPPORTED_FILE_TYPES.keys():
            for ext in SUPPORTED_FILE_TYPES[category]:
                if ext[1:] == ext or mime_type == ext[1:]:
                    return category

        return None

    def get_processor(self, file_type):
        """Retorna el procesador que maneja el tipo de archivo"""
        processor_map = {
            'pdf': PdfReader(),
            'image': ImageReader(),
            'video': VideoReader(),
            'audio': AudioReader()
        }
        return processor_map.get(file_type, None)
    
```

Passo 3: audio-processor.py

```

# Importar librerías
import whisper
import os
import time
from config import WHISPER_MODEL, LANGUAGE

class AudioProcessor:
    def __init__(self):
        self._model = whisper.load_model(WHISPER_MODEL)
        self._language = whisper.load_model(WHISPER_MODEL)

    def transcribe(self, audio_path):
        """Convertir audio a texto usando Whisper"""
        try:
            result = self._model.transcribe(
                audio_path,
                language=LANGUAGE,
                whisper_kwargs={
                    'langs': LANGUAGE,
                    'hybrid': True
                }
            )
            return result['text'], result['tokens'], result['logprobs'], result['language'], result['duration']
        except Exception as e:
            print(f"Error transcribing audio: {e}")
            return None

    def extract_from_video(self, video_path):
        """Extraer audio de video y transcribir"""
        try:
            import moviepy.editor as mp
            video = mp.VideoFileClip(video_path)
            audio = video.audio
            audio.write_audiofile(audio_path)
            logprobs = self.transcribe(audio_path)
            return logprobs
        except Exception as e:
            print(f"Error extracting audio from video: {e}")

    # Función para procesar audio
    with tempfile.NamedTemporaryFile(suffix=".mp3", delete=True) as temp:
        temp.write(open(audio_path, "rb").read())
        video = AudioFileClip(temp.name)
        video.write_audiofile(audio_path, verbose=False)

    def transcribe(self, audio_path):
        """Transcribir"""
        result = self._model.transcribe(audio_path)
        return result['text']

    # Limpiar
    os.unlink(audio_path)
    video.close()

    return result
    
```

These are just an idea of how to incorporate speech-to-text.
Not THE final result.

↑
WAIT FOR YASSIN'S CODE, NOT THIS ONE ☺

Phase 4: video_processor.py

Paso 5: Actualizar main.py

```

    result = process_attributes_file(file_path, output_base_folder)
    if result == 0:
        processed_no += 1
        print("Process (%d) completed: (%s)" % (processed_no, file_path))
    else:
        exception = Exception("Process (%d) failed: (%s)" % (processed_no, file_path))
        print(exception)
        print("Process (%d) failed processing (%s)" % (processed_no, file_path))

print("Process (%d) Processed (%s) files successfully!" % (processed_no, file_path))

```

add in requirement.txt:

openai-whisper = 2023/11/7

Speech Recognition 3.10.0

`pydub >= 0.25.1`