



# County of Santa Cruz



## HEALTH SERVICES AGENCY

### Environmental Health Division

701 Ocean St. Room 312, Santa Cruz, CA 95060

(831) 454-2022 TDD/ TTY: Call 711

EnvironmentalHealth@santacruzcountyca.gov

[www.scceh.org](http://www.scceh.org)

# AI-OCR Well Layer Update Project

## Executive Summary:

### Project Overview

The AI-OCR Well Layer Update Project was undertaken to modernize the County of Santa Cruz's GIS well layer—a foundational dataset for groundwater and water-resource management. Funded through the 2023 County Drought Resilience Planning Grant, the project applied AI-powered optical character recognition (OCR), custom Python scripts, GIS workflows, and structured quality-assurance (QA) methods to extract, clean, and integrate data from over 17,000 historical well completion files (representing approximately 10,000 unique wells). Many of these reports were decades-old, scanned PDFs with handwritten or inconsistent entries. The project's automated workflows brought them into a standardized, GIS-ready format to support more accurate analysis, planning, and long-term resilience.

### Objectives

- Improve completeness and accuracy of historical well records.
- Standardize key fields including well ID, well type, depth, water levels, locations, and screen intervals.
- Align County, OCR derived, and Department of Water Resources (DWR) datasets for GIS integration.
- Reduce manual data-entry workload and significantly improve geocoding accuracy.

### Approach

The project combined a custom-trained Microsoft Azure Document Intelligence model, six purpose-built Python scripts, fuzzy matching and aliasing logic, and multi-layered QA workflows to digitize, clean, classify, and integrate over 17,000 scanned well completion reports. ArcGIS-based workflows then geocoded wells to the highest achievable accuracy, with confidence flags for transparency.

## **Key Achievements**

- **Records Processed:** 17,672 well completion reports (7,046 County-held; 10,626 DWR) processed in under one day—saving thousands of hours of manual effort.
- **Unique Wells Identified:** Created the most complete countywide inventory to date: 10,007 unique, deduplicated wells, including 713 previously unrecorded wells discovered in County archives.
- **Dataset Alignment:** Reconstructed 95% of DWR entries directly from scanned originals.
- **Attribute Completeness:**
  - OCR workflows filled 78% of missing drill depths.
  - Improved extraction of static water level, depths, yield, screen intervals, drilling methods, and locations.
  - Classified wells into production (6,077), destruction (2,185), monitoring, and additional types.
- **Location Geocoding Improvements:**
  - 91% of wells geocoded beyond PLSS centroid placement.
  - 5,799 used precise County-verified coordinates;
  - 3,288 mapped to parcel centroids;
  - 898 mapped using PLSS or manual fallback.

- New “Exact” vs. “Approximate” vs “PLSS” confidence flags improve spatial transparency.
- **Duplicate Resolution:** Resolved distinct 66 duplicate IDs and more than 6,000 redundant scans, improving dataset quality.

### **Quality Assurance (QA) Highlights**

- Approximately 15% of AI supplemented records were flagged for review based on issues such as:
  - Conflicting depth measurements
  - Well type inconsistencies
  - Low-confidence APN/address matches
- Red-highlighted AI-derived fields, confidence scoring, and structured review queues ensured efficient and reliable manual verification.

### **Outcomes & Impact**

The project produced the most accurate, complete, and spatially reliable well layer ever assembled for Santa Cruz County. The dataset, which is published to the County’s public portal, GISWeb, now supports:

- Drought resilience planning
- Groundwater monitoring and basin assessment
- Well permitting and oversight
- Long-term water resource protection

The workflow is reproducible, scalable, and statewide-ready, offering a model for other counties and agencies seeking to modernize archival water datasets.

## **Conclusion**

By integrating AI-OCR, Python automation, QA, and GIS, the County transformed decades of handwritten, inconsistent, and poorly geocoded records into a unified, high-confidence dataset. The project dramatically reduced manual workload, increased spatial and attribute completeness, and demonstrated how modern AI tools can convert legacy water records into actionable information for sustainable groundwater management.

*The remainder of this report expands on these findings, outlining the methods, data processing steps, QA approach, and final outcomes in greater detail.*

# Contents

- Executive Summary: ..... 1
- Introduction ..... 6
- Background and Challenge..... 6
- Methodology..... 7
- Quality Assurance..... 13
- GIS Workflow ..... 17
- Results and Discussion..... 19
- Potential Enhancements ..... 24
- Conclusion ..... 28
- Appendix..... 29

## Introduction

The AI-OCR Well Layer Update Project modernized the County of Santa Cruz's GIS well layer, a critical tool for water resource management within the County, by automating data extraction and matching from well completion reports. Funded through the 2023 County Drought Resilience Planning Grant, the project combines AI-powered OCR, custom Python scripting, QA procedures, and GIS integration to significantly improve data accuracy and completeness. Initiated in 2024, this report outlines the project's methodology, including model training, six core data processing scripts, QA workflows, and GIS updates, as well as the resulting improvements.

## Background and Challenge

Since its creation in 2005, through the merger of County and partner agency datasets, the County of Santa Cruz's GIS well layer has faced persistent issues with incomplete and inconsistent information. A 2009 integration with the DWR well layer provided some improvements; however, more than 75% of the 11,700 wells still lacked well type classifications, and a significant number of records were missing altogether. In addition, several wells were entered into the dataset more than once. These gaps and inconsistencies have limited the Water Resources Division and other agencies' ability to accurately identify production wells and assess their impacts.

The current project was designed to address these data gaps and more, by applying AI-powered OCR to scanned well completion reports. Using custom scripts, the project extracted, cleaned, and standardized key fields—such as well ID, depth, water levels, APN, screen intervals, and address—and employed fuzzy matching to align these records with the existing DWR's well database. This process enhanced data completeness and consistency while preserving the County's existing spatial accuracy and retaining the county-specific fields

we still find valuable, such as whether an electric log is associated with the well record.

## Methodology

The project employed Microsoft Azure Document Intelligence for OCR, supported by six custom Python scripts for data processing, continuous QA to validate outputs, and GIS integration to update and visualize the results. The overall workflow is summarized below, with a visual overview of the script sequence provided in the attached diagram (*Script Sequence Workflow*, Appendix Page 1).

### **Model Training**

A custom Azure Document Intelligence model was trained using nine standardized DWR well completion report templates, dating back to the late 1940s to early 1950s. The model was designed to extract key fields such as well ID, total depth, water levels, completion dates, yields, well type, seals, lithology, screen intervals, APNs and other locational information, and more. It was optimized to handle a wide range of report formats—including handwritten entries and low-resolution scans—ensuring consistent and accurate data extraction across standard California DWR reports.

An example of the model's digitized output, displaying color-coded bounding boxes for extracted data and a right-hand panel listing some of the parsed fields, is provided in the following Azure Document Intelligence screenshot:

STATE OF CALIFORNIA  
WELL COMPLETION REPORT

10.51/341-2

Owner's Well No. 313.01  
Date Work Began 3.13.01  
Local Permit No. Santa Cruz County Health  
Permit No. 00-277  
Permit Date 11.29.00

WELL OWNER  
Name Tom Zingale / Marjorie Cameron  
Mailing Address 4415 Bonny Doon Rd  
Santa Cruz, CA 95060

LOCATION  
Address Same  
City Santa Cruz  
County Santa Cruz  
APN Book 103 Page 081 Parcel 10  
Township Range Section  
Latitude Longitude

WELL DEPTH  
DEPTH FROM SURFACE  
DEPTH TO FIRST WATER  
DEPTH OF SHAFT  
WATER LEVEL  
ESTIMATED YIELD  
TEST LENGTH  
TOTAL DEPTH OF BORING  
TOTAL DEPTH OF COMPLETED WELL

WATER LEVEL & YIELD OF COMPLETED WELL  
DEPTH TO FIRST WATER 175  
DEPTH OF SHAFT 313.01  
WATER LEVEL 176  
ESTIMATED YIELD 176  
TEST LENGTH 2  
TOTAL DEPTH OF COMPLETED WELL 250

ANNULAR MATERIAL

DEPTH FROM SURFACE	DEPTH TO FIRST WATER	DEPTH OF SHAFT	WATER LEVEL	ESTIMATED YIELD	TEST LENGTH	TOTAL DEPTH OF COMPLETED WELL
0	175	313.01	176	176	2	250

ATTACHMENTS

CERTIFICATION STATEMENT

Earth Flow Drilling Company  
21000 Smith Grade Santa Cruz, CA 95060  
Bookkeeper 4-1-01 679020

Item	Value	Percentage
template_1_0004F44D (1).pdf		
other_well #1		9.80%
unselected		
test_well #1		9.80%
unselected		
Log (11) #1		29.70%
template1_screen_row5 #1		35.60%
unselected		
template1_screen_row6 #1		35.60%
unselected		
Casings (3) #1		49.40%
owner_address #1		52.20%
4415 Bonny Doon Rd Santa Cruz, CA 95060		
contractor #1		58.80%
Earth Flow Drilling Company		
total_depth #1		59.60%
250.		

## Data Processing Scripts

Six custom Python scripts were developed to process, clean, classify, and enhance the data extracted through OCR. The scripts leverage AI-assisted coding tools (e.g., ChatGPT, Copilot) to accelerate prototyping and debugging. They utilize libraries including pandas, fitz, tempfile, numpy, os, json, logging, re, collections, rapidfuzz, azure-ai-formrecognizer, and openpyxl.

The scripts work sequentially to extract structured data from scanned well reports, clean and standardize fields, classify well types, match records

across multiple datasets, integrate county records, and highlight AI-derived fields for quality control. The workflow for each script is summarized below.

## **1. OCR Extraction Script**

Purpose: Convert raw PDF well completion reports into structured JSON data suitable for downstream processing.

### Overview:

Processes scanned PDF well reports using Azure Document Intelligence for OCR, with fitz handling PDF manipulation and compression.

### Key Functions:

- Compresses large PDFs (>4 MB) to ensure compatibility with Azure OCR.
- Extracts fields such as well\_id, depth, APN, and Casings using the custom-trained Azure model.
- Captures page- and line-level content in nested JSON dictionaries.
- Logs extraction errors and skips failed PDFs.
- Outputs JSON files to specified directory.

## **2. Data Extraction and Classification Script**

Purpose: Transform JSON outputs from OCR into structured Excel tables, with keyword detection to support well type classification.

### Overview:

Processes JSON outputs from the OCR Extraction Script, extracting key fields and identifying destruction-related keywords. Outputs two Excel sheets: *Main Data* and *Casings*.

### Key Functions:

- Parses JSON dictionaries, extracting fields like well\_id, destruction\_well, address, and APN.
- Detects destruction-related keywords (e.g., "abandonment," "destroy," "decommission") using fuzzywuzzy (token\_sort\_ratio threshold 85), including common misspellings.

- Converts Casings lists into DataFrame rows with traceability via source\_file and row\_index.
- Skips invalid JSON structures and logs errors.
- Generates Excel file output to specified directory.

### **3. Data Cleaning and Matching Script**

Purpose: Produce a cleaned, deduplicated dataset with standardized fields, enhanced classifications, and matched geospatial data.

Overview:

Processes the Excel output from extraction to clean, standardize, and enrich data. Handles OCR errors, refines well classifications, extracts perforation intervals, and prepares addresses/APNs for matching.

Key Functions:

- Cleans well\_id using regex, handling OCR misreads.
- Classifies templates (1–9 or Unknown) via revision-based signatures.
- Refines destruction classifications using baseline keyword lists.
- Extracts perforation intervals and validates top/bottom values.
- Validates the values of Total Completed Depth and Total Drill Depth, ensuring that Total Completed Depth is always greater than Total Drill Depth.
- Detects duplicates using weighted scoring and retains the most complete record.
- Cleans addresses and other fields (e.g. City, Driller Name, Fluid Type, etc.) via dictionary-based aliases and regex;
- Searches for and extracts APNs from the document (which may be located in multiple places within the document).
- Derives Type of Work hierarchically from fields like keyword\_destruction\_well or new\_well.
- Handles invalid numeric fields (NaN) and missing data.
- Generates Excel file output to specified directory.

### **4. Data Matching and Merging Script**

Purpose: Align OCR-derived well records with the DWR database, enhance APN and address accuracy, and produce a comprehensive matched dataset.

Overview:

Integrates cleaned OCR data with DWR records and optionally checks Monterey County or other County adjacent datasets, generating multiple output files for matched and unmatched records.

Key Functions:

- Cleans well IDs and normalizes APNs (handling 8- and 9-digit formats).
- Performs primary matching on well IDs, validated by depth or water level.
- Multiple matches are resolved with a weighted scoring algorithm
- Refines APNs and addresses using parcel validation and fuzzy matching techniques.
- Validates DWR APNs against a dictionary mapping to parcel addresses. If an APN is invalid, it extracts APNs from the Well Location or other fields if present, or employs rapidfuzz for address matching. Final APNs are formatted to align with the parcel database standards used for GIS integration.
- Checks unmatched OCR records against adjacent County datasets.
- Outputs:
  - unmatched\_ocr\_output.xlsx
  - unmatched\_dwr\_output.xlsx
  - matched\_output.xlsx
- Maintains source tracking and logs matching details for QA.

## **5. County Database Matching Scripts**

Purpose: Integrate county database records with DWR or unmatched OCR data, validating and enriching well records.

Overview:

Two scripts align county records with either DWR data or OCR data, producing matched and unmatched outputs.

### Key Functions:

- Normalizes county and source IDs using regex (handling OCR misreads and variations).
- Matches county Log\_Number to DWR Legacy Log Number/WCR Number or OCR well\_id\_cleaned.
- Prevents duplicate matches with sets and logs metadata.
- Produces three Excel outputs per script: matched records, unmatched county records, and unmatched source records.

## **6. Field Highlighting Script**

Purpose: Visually highlight OCR-derived or non-validated fields in the final supplemented dataset for QA based on the data source.

### Overview:

Applies conditional formatting to the supplemented dataset Excel file using openpyxl, marking AI-sourced fields or unmatched values in red font.

### Key Functions:

- Loads Excel workbook and maps columns dynamically.
- Highlights final\_apn and final\_address if match\_source indicates non-validated or OCR-derived data.
- Highlights other fields based on a predefined field\_map when source = "AI."
- Skips missing columns and logs errors.
- Saves the modified Excel file, overwriting the input file.

The following image shows an example output of the field-highlighting script, with color-coded AI-generated values:

	A	B	E	H	I	N	P	AI	AK	AL	AM	AN	AS	BE
1	WCR Number	Legacy Log Number	City	Permit Date	Permit Number	Driller Name	Driller License Number	Total Drill Depth	Ital Completed Depp	Of Perforated Interm	of Perforated In	Casing Diameter	Static Water Level	B118WellUse
92	WCR0153482	E0263622	Scotts Valley	3/19/15	SM-15-002	EXPLORATION GEOSERVICES INC (score=9	484288		30					Monitoring Well
93	WCR1992-018220	361482			92-136	POLLOCK'S DRILLING & PUMP SERVICE	379086	780	760	720	860		350	Domestic
94	WCR0274630	18373				C & N PUMP & WELL CO	68648							Domestic/Public Supply Well
95	WCR0290064	NN				FREEDOM DRILLING	unknown	150	145	66	150	10		
96	WCR0117565	NN	SANTA CRUZ	4/25/06	05-319	PIERCE & SON	unknown	580	220	370	510	12	99	
97	WCR0276070	NN	Watsonville	11/6/15	15-072	LITCHFIELD	unknown	415	362	0	470		60	
98	WCR0251594	327077	FELTON		8910189A	AQUA SCIENCE ENGINEERS INC	487000	83	81	79	81	2	79.4	Monitoring
99	WCR0177731	926182	LOS GATOS	2010-09-09 00:00:00	10-092	GREGG DRILLING & TESTING INC	485165	30	30	20	30	2	23	Monitoring Well
100	WCR0118327	261906	SANTA CRUZ		1928-28 AND 88-261	MAGGIORA BROS DRILLING INC	249957	7	17	7	17	2	9	Monitoring
101	WCR0029130	298476	FELTON			AQUA SCIENCE ENGINEERS INC	487000	60	59.7	57.7	59.7	2		Monitoring
102	WCR0046082	324356	FELTON		8910189A	JEFFREY D MORGAN WATER DEVELOPMEN	283326	334	334	329	334	2	314.2	Monitoring
103	WCR0025311	324357	FELTON		8910189A	PITCHER DRILLING CO	263085	326	305	295	305	1	267.05	Monitoring
104	WCR0277831	200034	S CRUZ			MAGGIORA BROS DRILLING INC	249957	15	14.5	4	14	4	6.5	Monitoring
105	WCR0123199	200037	S CRUZ			MAGGIORA BROS DRILLING INC	249957	14.8	14.8	4.3	14.3	4	6.1	Monitoring
106	WCR0021275	425649	WATSONVILLE	Na	325-25.01	R BARRICK FOR BAYLANDS	374152							Monitoring
107	WCR0315111	381860	WAT			G E WEBER & ASSOCIATES INC	unknown	20	20	10	20	2	4	Monitoring
108	WCR0323923	55603				ROMAN WELL DRILLING	265064	415	365	325	355	12	10	Public Supply
109	WCR0032522	313021			1988-11-01 00:00:00	EARTH FLOW DRILLING	320331		300	220	300	5	240	Domestic
110	WCR0092975	261740			88-08	MAGGIORA BROS DRILLING INC	249957	16	15.1	4.7	14.7	4.5	12	Monitoring
111	WCR2010-011978	E0136967	CAPITOLA			R S I DRILLING	802334	5.5	5.5	4.5	5			Other Well
112	WCR0118334	550821	WATSONVILLE			EXPLORATION GEOSERVICES INC	484288	25	25	15	25	4	16.68	Monitoring
113	WCR0088754	365674		May 22, 1991	1991-05-20 00:00:00	PC EXPLORATION INC	265556							Monitoring
114	WCR2015-007806	E0285609	SANTA CRUZ	9/5/14	SM14-188	ALL WELL ABANDONMENT	848359							
115	WCR0276740	305149	CAPITOLA		88-000295	BAYLAND DRILLING CO	374152	30.5	27.8	8.3	27.8	2	13.7	Monitoring
116	WCR0305688	192402	SANTA CRUZ			ROGERS JOHNSON & ASSOCIATES	unknown	15	14.5	4.5	14.5	2	5.6	Monitoring
117	WCR0188616	192403	SANTA CRUZ			ROGERS JOHNSON & ASSOCIATES	unknown	16	15.5	5.5	15.5	2	5.8	Monitoring
118	WCR0248871	550823	WATSONVILLE			EXPLORATION GEOSERVICES INC	484288	25	25	15	25	2	17.44	Monitoring
119	WCR0304159	313867	WATSONVILLE		88-234	EXPLORATION GEOSERVICES INC	484288	33	30	9	30	2	16	Monitoring
120	WCR2015-007799	E0285622	SANTA CRUZ	9/5/14	SM14-192	ALL WELL ABANDONMENT	848359							
121	WCR0123198	423358	SOQUEL	2/5/91	91-034	WOODWARD DRILLING COMPANY	581639	36	36	21	36	4		Monitoring
122	WCR2010-011979	E0136966	CAPITOLA			R S I DRILLING	802334	10.5	10.5	9.5	10			Other Well
123	WCR2024-003460		Santa Cruz	2023-12-05 00:00:00	SR0015594	CONFLUENCE TECHNICAL SERVICES INC	1035255		8					Other

Note: See the Appendix for the code and detailed function descriptions.

## Quality Assurance

The Well Layer Update data processing workflow incorporates multiple layered quality assurance mechanisms to improve accuracy, reliability, and traceability. These mechanisms operate at the field-cleaning, matching, and validation stages and generate flags, scores, source tracking, and visual cues that guide manual review. Outputs from the Data Cleaning and Matching Script, Data Matching and Merging Script, County Database Matching Scripts, and Field Highlighting Script include dedicated QA columns, confidence metrics, and conditional formatting designed to highlight low-confidence or potentially erroneous values.

The key QA mechanisms are summarized below.

### 1. Location Matching Confidence and Source Tracking

Purpose: Provide transparent, quantifiable confidence in the final APN and address used for geocoding, enabling targeted manual review of lower-quality locations.

Implementation:

- Matching follows a strict hierarchy: exact validated APN → exact validated address → fuzzy validated APN → fuzzy validated address → fallback methods.
- Every record receives:
  - match\_source / OCR\_match\_source: a human-readable string describing the match method and score (e.g., “ID match → APN validated”, “Fuzzy address match (score 78.2%)”)
  - Numeric fuzzy scores (0–100) from rapidfuzz when fuzzy address or APN parcel is used.
  - Alias/regex confidence scores when street-name or city dictionaries are applied during address cleaning.

QA Role: Records with fuzzy scores <85, non-APN-based matches, or low alias/regex scores (<80%) are prioritized for manual review using GIS tools and source PDFs.

## 2. **AI Source Indicators for OCR-Derived Data**

Purpose: Distinguish OCR-derived data from DWR or county-validated data.

Implementation: Scripts track OCR-derived fields using source columns (e.g., apn\_source, depth\_source) marked “AI.” The Field Highlighting Script applies red font to these cells (e.g., final\_apn, final\_address, seal\_cleaned, template\_id).

QA Role: Red highlights visually flags potentially error-prone data, prompting user to verify against the original PDFs.

### 3. **Type Mismatches for Record Consistency**

Purpose: Detect discrepancies between DWR and OCR-derived well type classifications.

Implementation: The Record\_Type\_Mismatch field flags conflicts between DWR B118WellUse and OCR-derived types.

QA Role: Prioritize mismatched records and cross-reference source documents to resolve classification errors.

### 4. **Destruction Well Cross-Referencing**

Purpose: Improve destruction well accuracy, especially when handwritten notes override checked boxes (common when a proposed well was abandoned after drilling began).

Implementation: keyword\_destruction\_well refines the initial checkbox detection by comparing additional keyword counts against template-specific baselines.

QA Role: Records where OCR and keyword indicators disagree and are flagged for manual confirmation.

### 5. **Total Depth Validation Flag (depth\_flag / review\_depths)**

Purpose: Ensure consistency between Total Drill Depth and Total Completed Depth and catch common OCR or data-entry errors.

Implementation:

- Validates that Total Drill Depth  $\geq$  Total Completed Depth (with minor tolerance).
- Sets depth\_flag = 'review depths' when the relationship is violated or values appear swapped.

## 6. Dictionary-Based Alias Cleaning with Regex Confidence Scoring

Purpose: Standardize inconsistent or misspelled categorical/text fields (City, Driller Name, Fluid, Test Type, Drilling Method, etc.) while quantifying confidence.

Implementation:

- Large alias dictionaries map historical and variant names to standardized values.
- Two-stage cleaning: (1) dictionary lookup with rapidfuzz token-set ratio, (2) regex pattern fallback if no strong dictionary hit.
- Each cleaned value receives an alias\_match\_score (0–100)

### Manual Review Process

The QA outputs feeds into a structured review workflow:

- Red-highlighted AI-sourced fields (Field Highlighting Script)
- depth\_flag = True or conflicting depth measurements
- Low dictionary/regex alias match scores (<80%) on City, Driller Name, Fluid Type, Drilling Method, etc.
- Record\_Type\_Mismatch or conflicting destruction indicators
- Fuzzy address/APN match scores <85 (especially non-APN or direct address matches)
- OCR\_match\_source indicating fallback or low-confidence location methods
- B118WellUse anomalies possibly caused by 'Confidential' stamp interferences

Verification Tools: ArcGIS Pro, original scanned PDFs, county parcel records and historical lists are utilized to update flagged record values.

# GIS Workflow

The GIS workflow for the Well Layer Update emphasizes maximizing spatial accuracy of well locations by integrating updated database outputs with county-maintained spatial layers. It combines automated geocoding, join operations, and QA field value comparison checks to produce a high-confidence, GIS-ready dataset. The workflow prioritizes county-verified locations, leverages parcel-based centroids for unmatched (county) APN records, and applies PLSS grid centroids or manual location sketch plotting as fallbacks for records with missing or low-confidence location data.

## **Geocoding Process**

### 1. Matching with County-Mapped Well Records:

Purpose: Preserve the accuracy of well locations already verified and mapped by county staff.

Implementation:

- Updated database records that match existing county-mapped wells are spatially joined based on Log\_Number and well\_id\_cleaned.
- The county well layer geometry, established over years from DWR report records and maps, is preserved.
- County unique fields such as county\_ELOG are appended to matched records.

QA Role:

- High-confidence matches inherit county-verified locations, reducing the need for further location verification.

### 2. APN and Address-Based Geocoding with Parcel Layer

Purpose: Assign the most precise possible locations to wells using validated APNs or addresses, while distinguishing exact from approximate placements.

Implementation:

- Records with a validated APN or address (cross-referenced against the County's Parcel database) are spatially joined to the centroid of the corresponding parcel polygon.
- Each record receives a Location\_Accuracy field with one of two values:
  - "Exact" – assigned when the match is based on an exact, validated APN or address match, or inheritance from an existing county-mapped well.
  - "Approximate" – assigned when the match relies on fuzzy APN or address matching with a fuzz score <100 but ≥80–85

3. Fallback to PLSS Grid or Manual Review

Purpose: Provide approximate locations for records that are lacking or have low confidence APN/address data.

Implementation:

- PLSS Grid: Wells are mapped to the centroid of their corresponding PLSS Section grid in the county.
- Manual Review with Location Sketches: Extract well positions (ongoing) from sketches in DWR reports into ArcGIS.
- Records with low (<80) OCR\_match\_source fuzzy scores (e.g., "Address matched (score 63.64)") are flagged for fallback or manual review processing before mapping to that address or PLSS.

QA Role:

- Wells without confident location data are mapped via PLSS centroids and are subject to future cross-checks with the source file to improve location data manually.

### GIS Summary:

This GIS workflow ensures well locations in the Well Layer Update are geocoded with maximum accuracy by:

1. Prioritizing county-verified coordinates.
2. Leveraging parcel-based centroids for direct match or high-confidence APN records.
3. Applying existing DWR PLSS grid mapping for low-confidence or missing data.

## Results and Discussion

The project processed 17,672 well completion reports (7,046 county, 10,626 DWR) in batches using the developed Python scripts, achieving improvements in data integration, accuracy, and efficiency. The following results and metrics highlight improvements across data processing, matching, geocoding, field supplementation, QA, and GIS integration.

### **Data Processing**

- Total Records Processed: Successfully extracted data from 17,672 files, with nearly all files processed without errors, despite challenges with low-quality scans and illegible handwritten data. Note some of these files ended up not being well completion reports (i.e. e-logs, site plans, boring logs, etc.)

### **Matching Accuracy**

- County-to-DWR Matching: 5,683 of the 6,698 county-scanned well completion reports (85%) were successfully matched to corresponding DWR records using a combination of exact well ID matching and fuzzy reconciliation techniques.
- The remaining ~15% (1,015 records) were unmatched due to either the record not existing in the DWR database or poor scan quality or illegible handwriting that prevented accurate well ID extraction.
- As of April 28, 2025, the DWR online system for WCR contained 9,406 well records geographically associated with Santa Cruz County. Around the same date, 10,651 PDF files were downloaded from the DWR website for processing.
- After the initial county-to-DWR matching, 3,723 DWR records remained unsupplemented by county scans (9,406 total DWR records – 5,683 county-matched). Of these, 3,070 were subsequently matched and enriched using OCR data extracted directly from the downloaded DWR-sourced PDFs (following deduplication in the Data Cleaning and Matching Script). This left only 653 DWR records without OCR-derived supplemental data.
- Unique and Non-Well Records: 713 county-scanned records had no corresponding entry in the DWR database (true new records).
- ~530 county files were identified during QA review as non-well completion reports (e.g., boring logs, site plans, or electric logs) or were illegible and excluded from final counts.
- The OCR extraction and processing pipeline, when applied to the 10,651 PDFs downloaded from DWR's OSWCR website, successfully recreated 95% of the 9,406 official DWR records for Santa Cruz County using only the digitized documents themselves. When combined with the 713 previously unrecorded wells identified from county-held scans, the project achieved 103% relative completeness compared to the official Santa Cruz County DWR database as of April 2025.
- The updated master database contains 10,007 unique well records. This final dataset is a synthesis of DWR records enriched with AI-OCR-

extracted data, 713 previously unrecorded wells identified from county-held scans, and county records that were supplemented and validated using OCR.

### **Location Accuracy and Improvements**

- Geocoding Success: Of the 10,007 wells in the updated dataset, 9,135 (91%) now have geocode locations ranging from exact well locations to approximate well locations—a significant improvement over the previous layer.
  - County-Mapped Wells: 5,799 matched records inherited precise county coordinates through GIS joins.
  - APN-Based Geocoding: 3,288 wells were mapped to parcel centroids, with 80% (2,644) validated as high-confidence (e.g. “APN matched”) and the remaining 644 as approximate locations.
  - PLSS/Manual Fallback: 898 wells with low fuzzy-match scores (<80%) or missing location data were mapped to PLSS centroids, flagged for later QA review.
  - The remaining 48 inherited locations from DWR associated with GPS, Other, or Surveyed to Benchmark.

### **Well Type Improvements**

The previous County well database had limited and inconsistently defined well type information. It often wasn’t clear whether a well was a production well—or what type of production it served (domestic, irrigation, industrial, or municipal)—and it also did not reliably distinguish destruction wells from ‘other’ categories such as test, monitoring, or boring wells.

The updated database now classifies each well type based on the information available in the well completion reports, resulting in a more accurate and consistent countywide dataset. The updated counts include 6,077 production wells, 2,185 destruction wells, and 1,747 wells classified as other types.

**Field Supplementation**

- Critical Fields Filled: The Data Matching and Merging Script supplemented previously blank or unknown DWR fields using OCR and county data:

Field	Records Filled	Total Matched	Percent Filled
Well completion date	513	8,752	5.9%
Total drill depth	6,852	8,754	78.3%
Final well depth	405	8,754	4.6%
Static water level	545	8,754	6.2%
Well yield	1,074	8,754	12.3%
Well use/source	1,788	8,754	20.4%
Drilling method	896	8,754	10.2%
Screen intervals	670	8,754	7.7%

- Unique Fields: Additional fields (e.g., well seals, first water levels encountered, geologic materials) were extracted. Note: geologic materials were not processed or incorporated in any of the output datasets.

## Duplicate Cleanup

- Duplicate Resolution: The Data Cleaning and Matching Script reconciled 66 distinct duplicate well IDs—cases where the same well log ID referred to different wells—and eliminated more than 6,000 duplicate PDF scans, significantly reducing redundancy in the dataset.

## Time Savings

- Processing Efficiency: 17,672 reports were processed in under one day, saving significant amounts of time compared to manual extraction and geocoding, equivalent to thousands of hours saved if that would have been required to manually supplement or input the data.
- Script Runtime Breakdown:
  - OCR Extraction Script: ~70% of runtime
  - Post-Processing Scripts: ~30% hours, optimized by batch processing and indexing

## Quality Assurance Metrics

QA-Flagged Records: Approximately 15% of all records (2,651 of 17,672) were flagged for manual review. Key flag categories included:

- ~2000 records affected by “*Confidential*” stamp interference (misread numbers, checkbox errors).
- 393 records with conflicting depth measurements (review\_depths flag).
- 181 records with well type mismatches (often caused by “*Confidential*” stamp interference).

- 1,857 records with fuzzy match scores below 80% (from OCR\_match\_source).
- Roughly 5% of supplemented fields using the street-name alias dictionary generated low alias match scores (<80%), requiring manual refinement.

### **Challenges and Solutions:**

- Poor scan quality / illegible handwriting affected numerous reports, mitigated through data cleaning, OCR error logging, and QA.
- Template inconsistencies: Handled by template-specific cleaning logic
- OCR errors: Minimized using a combination of fuzzy matching, dictionary-based aliasing, standardized or pattern-based cleaning routines, and QA-flagged review.

### **Scalability**

- Statewide-Ready Framework:

The California DWR template-based processing, APN integration, and GIS workflows create a scalable framework that can be extended to other counties or applied to statewide DWR datasets.

## Potential Enhancements

### **Model Training & Template Improvements**

- Train additional documents to improve value extraction, especially for Template 2.
- Add training for all selection boxes across templates (equipment-type checkboxes, cathodic protection wells, injection wells, etc.) to prevent lesser-used boxes from interfering (common) with more critical selections located nearby.
- Additional training across all known templates to address missed or unrecognized templates.
- Retrain the AI-OCR model to better detect monitoring wells and other specialized well types.
- Refinement of template classifications.
- Improved template-specific logic for screen interval extraction, especially for Template 1.
- Template-specific well ID cleaning, such as enforcing expected ID ranges (e.g., Template 2).

### **OCR & Document Handling**

- Improved OCR performance for low-quality or complex scans.
- Automated detection and removal of “Confidential” stamps.
- Enhanced parsing and handling of PDFs containing multiple well completion reports.
- Add logic to detect and classify documents that are not well completion reports (borings, e-logs, site plans, etc.).

### **Data Cleaning & Field Logic Enhancements**

- Strengthen cleaning algorithms for Total Depth, including:
  - cross-checking Total Drill Depth against the bottom depth in the geologic log,
  - handling “Same” entries,
  - applying fallback logic to populate final depth from geologic log values when needed.
- Add logic and fields for detailed screen interval depths.
- Add lithology post-processing.
- Enhance screen interval logic (e.g., ensure the screen checkbox aligns with slot-size entries in the same row).
- Add a new field to classify wells that were attempted but ultimately dry (Files that contain: “Dry”, “Dry Hole”, “Dry Well”).
- Add a new field to identify wells with accompanying electric logs (e.g., files containing terms such as “Electric Log,” “E-Log,” or “Elog,” geophysical logs marked as checked, or forms where “Was electric log made?” is checked).
- Logic to classify wells with no “Proposed Use” selected as Production if a screen interval exists, with a QA flag to ensure they aren’t monitoring wells.
- More robust handling of duplicate well IDs.

### **Quality Assurance Enhancements**

- Increased QA checks comparing AI-extracted data with DWR records to identify and correct errors

- Add built in QA checks (rather than relying later on GIS checks which are more time intensive) to flag impossible or inconsistent values, such as:
  - static water level or first water deeper than total drill depth,
  - bottom of perforated interval deeper than total drill depth.
- Several DWR data entry errors were identified during GIS QA, including:
  - “Top of perforations” often recorded as the bottom-most screening interval rather than the top-most
  - Static water level and other fields off by up to hundreds of feet
  - Incorporate these types of checks directly in the extraction scripts to reduce reliance on later GIS QA which can be time intensive
  - When conflicts are detected, use AI-extracted data as a replacement if it falls within expected value ranges relative to other related fields
- Additional QA checks to ensure fuzzy-matched low-score addresses remain within previously identified PLSS grids.
- Incorporate confidence scores in JSON outputs to guide QA and validation.

### **APN, Address, and Location Improvements**

- Integrate historical APN datasets to account for subdivisions or APNs that have changed over time.
- Expand location-improvement logic, including using owner name information to fuzzy match against parcel databases and verify whether a corresponding parcel address aligns with the WCR-reported street/location.

- Add logic to split extracted TRS values into Township, Range, and Section fields for new or supplemented records.

### **Additional Functional Enhancements**

- Utilize unredacted reports (when available) to access owner information, which is essential when WCRs list the well location as “same” as the owner’s address.

## Conclusion

The Well Layer Update project demonstrates the effectiveness of automated data processing and GIS integration in modernizing well record management, as well as the transformative potential of AI-powered OCR software. Recent advancements in AI OCR have significantly improved accuracy—particularly when handling complex layouts, mixed formatting, and handwritten text—making it possible to rapidly digitize and analyze large volumes of historical documentation.

By processing 17,672 well completion reports in under a day, identifying missing records, supplementing key fields, and achieving 91% geocoding accuracy (tied to locations other than PLSS centroids), the project illustrates major gains in efficiency, completeness, and spatial reliability. These outcomes highlight the value of applying AI OCR to historical records across other County divisions, as well as in other counties or organizations with large archival datasets.

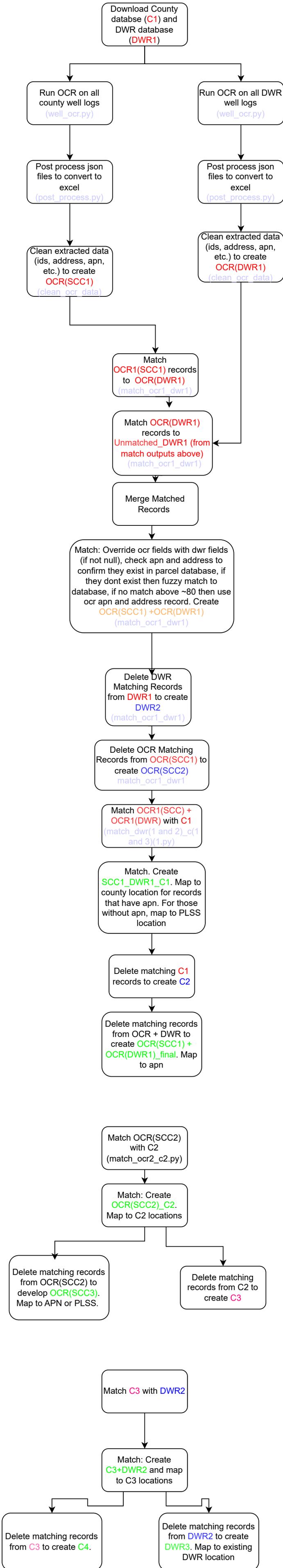
As AI technology continues to advance, extraction accuracy will further improve, expanding opportunities to modernize legacy datasets and integrate them into contemporary water management and planning workflows.

The updated well layer is scheduled to be integrated into GISWeb in January 2026, supporting drought resilience and groundwater sustainability efforts, and enhancing water resource decision-making. The project demonstrates how modern tools can transform previously inaccessible or labor-intensive data into a reliable, comprehensive resource for staff, analysts, and policymakers.

## Appendix

Below is the script flowchart workflow as well as the primary scripts, with comments explaining the purpose of key functions.

# Script Sequence Workflow



## C:\Well Layer Update\Python\Azure\OCR Extraction Script secure.py

```
1 import os
2 import json
3 import fitz
4 import tempfile
5 from azure.core.credentials import AzureKeyCredential
6 from azure.ai.formrecognizer import DocumentAnalysisClient
7
8 # Azure config
9 # Set these in your environment before running the script:
10 # AZURE_FORMRECOGNIZER_ENDPOINT = https://your-resource.cognitiveservices.azure.com/
11 # AZURE_FORMRECOGNIZER_KEY      = your-key-here
12
13
14 endpoint = os.getenv("AZURE_FORMRECOGNIZER_ENDPOINT")
15 key      = os.getenv("AZURE_FORMRECOGNIZER_KEY")
16 model_id = os.getenv("AZURE_FORMRECOGNIZER_MODEL_ID", "wcr11")
17
18 # Validate that required variables are present
19 if not endpoint or not key:
20     raise EnvironmentError(
21         "Please set the environment variables before running this script."
22     )
23
24 # Trim any accidental whitespace or trailing slashes
25 endpoint = endpoint.rstrip("/")
26
27 BASE_DIR = r"<internal_project_directory>"
28
29 input_folder = f"{BASE_DIR}\WCR_Files"
30 output_folder = f"{BASE_DIR}\json files"
31 os.makedirs(output_folder, exist_ok=True)
32
33 # Azure client
34 client = DocumentAnalysisClient(endpoint=endpoint, credential=AzureKeyCredential(key))
35
36 # Compress large image only PDFs
37 def compress_pdf_with_pymupdf(file_path, max_size_mb=4):
38     file_size_mb = os.path.getsize(file_path) / (1024 * 1024)
39     is_image_only = True
40
41     try:
42         doc = fitz.open(file_path)
43         for page in doc:
44             if page.get_text().strip():
45                 is_image_only = False
46                 break
47
48     if file_size_mb <= max_size_mb or not is_image_only:
```

```
49         return file_path # Skip compression
50
51     print(f"🔗 Compressing large image-only PDF: {os.path.basename(file_path)}")
52     temp_file = tempfile.NamedTemporaryFile(delete=False, suffix=".pdf")
53     temp_file.close() # Close so fitz can write to it
54     new_doc = fitz.open()
55
56     for page in doc:
57         pix = page.get_pixmap(dpi=100)
58         img_bytes = pix.tobytes("jpeg")
59         img_doc = fitz.open("jpeg", img_bytes)
60         rect = img_doc[0].rect
61         pdf_page = new_doc.new_page(width=rect.width, height=rect.height)
62         pdf_page.insert_image(rect, stream=img_bytes)
63
64     new_doc.save(temp_file.name)
65     new_doc.close()
66     return temp_file.name
67
68     except Exception as e:
69         print(f"Could not compress PDF ({file_path}): {e}")
70         return file_path
71
72 # Process each PDF
73 for dirpath, dirnames, filenames in os.walk(input_folder):
74     for filename in filenames: #iterate over the filenames list
75         if filename.lower().endswith(".pdf"):
76             original_path = os.path.join(dirpath, filename) # Use dirpath to construct full
path
77             pdf_to_process = compress_pdf_with_pymupdf(original_path)
78
79             print(f"\n📄 Processing file: {filename}")
80
81             with open(pdf_to_process, "rb") as f:
82                 try:
83                     poller = client.begin_analyze_document(
84                         model_id=model_id,
85                         document=f
86                     )
87                     result = poller.result()
88
89                     result_dict = result.to_dict()
90                     json_filename = os.path.splitext(filename)[0] + ".json"
91                     json_path = os.path.join(output_folder, json_filename)
92
93                     with open(json_path, "w", encoding="utf-8") as json_file:
94                         json.dump(result_dict, json_file, indent=2)
95
96                     print(f"Saved JSON: {json_path}")
97
```

```
98         except Exception as e:
99             print(f"Error processing {filename}: {e}")
100
101     # Clean up temp file if we compressed
102     if pdf_to_process != original_path:
103         try:
104             os.remove(pdf_to_process)
105             print(f"Deleted temporary file: {os.path.basename(pdf_to_process)}")
106         except Exception as e:
107             print(f"Failed to delete temp file: {e}")
```

## C:\Well Layer Update\Python\Azure\Data Extraction and Classification Script.py

```
1 import os
2 import json
3 import pandas as pd
4 from fuzzywuzzy import fuzz
5 import re
6
7 BASE_DIR = r"<internal_project_directory>"
8
9 # Paths for input JSON files and output Excel file
10 input_folder = f"{BASE_DIR}\json files"
11 output_excel = f"{BASE_DIR}\extracted_data_dwr.xlsx"
12
13 # Destruction related keywords (including common misspellings)
14 destruction_keywords = [
15     "abandonment", "abandon", "abandone", "abandoment",
16     "destroy", "destruction", "distruction", "demolish", "decommision", "decommission",
17     "destroying"
18 ]
19 destruction_keywords = [kw.lower() for kw in destruction_keywords]
20
21 # Extract raw text from JSON content for keyword search
22 def extract_raw_text(json_content):
23     # Try to get the main content field
24     raw_text = json_content.get("content", "").lower()
25     if not raw_text: # Fallback to concatenating line content from pages
26         raw_text = " ".join(
27             line["content"].lower()
28             for page in json_content.get("pages", [])
29             for line in page.get("lines", [])
30             if "content" in line and isinstance(line["content"], str)
31         )
32     return raw_text
33
34 #Count destruction keywords using fuzzy matching to handle OCR errors
35 def count_destruction_keywords(text, keywords, threshold=85):
36     words = re.findall(r'\b\w+\b', text) # Split into words
37     matched = []
38     for word in words:
39         scores = [fuzz.ratio(word, kw) for kw in keywords]
40         max_score = max(scores, default=0)
41         if max_score >= threshold:
42             best_index = scores.index(max_score)
43             best_match = keywords[best_index]
44             matched.append((word, best_match))
45     # Get unique original words
46     unique_matched = sorted(set(word for word, _ in matched))
47     return unique_matched
```

```
48 # Initialize lists for data storage
49 main_data = []
50 casings_data = []
51
52 # Process JSON files to extract fields, generate destruction keywords, and structure data
53 for filename in os.listdir(input_folder):
54     if filename.endswith(".json"):
55         path = os.path.join(input_folder, filename)
56         with open(path, "r", encoding="utf-8") as file:
57             try:
58                 json_content = json.load(file)
59                 doc = json_content.get("documents", [{}])[0]
60                 fields = doc.get("fields", {})
61
62                 # Extract raw text for keyword search
63                 raw_text = extract_raw_text(json_content)
64                 matched_keywords = count_destruction_keywords(raw_text, destruction_keywords)
65
66                 # Build keyword string
67                 keywords_str = ', '.join(matched_keywords) if matched_keywords else ''
68
69                 # Main data: all top-level fields except Casings and Log, preserve existing
destruction_well
70                 record = {
71                     k: (v.get("content") if isinstance(v, dict) else v)
72                     for k, v in fields.items()
73                     if k not in ["Log", "Casings"]}
74                 }
75                 record["source_file"] = filename
76                 record["additional_destruction_keywords"] = keywords_str
77
78                 main_data.append(record)
79
80                 # Casings data: parse each entry in the list
81                 if "Casings" in fields and fields["Casings"].get("value_type") == "list":
82                     for idx, casing in enumerate(fields["Casings"]["value"], start=1):
83                         row = {
84                             "source_file": filename,
85                             "row_index": idx
86                         }
87                         if casing.get("value_type") == "dictionary":
88                             for key, val in casing["value"].items():
89                                 row[key] = val.get("content") if isinstance(val, dict) else
val
90                                 casings_data.append(row)
91
92             except Exception as e:
93                 print(f"Error processing {filename}: {e}")
94
95 # write extracted and classified data to Excel with separate sheets for Main Data and Casings
```

```
96 with pd.ExcelWriter(output_excel, engine="openpyxl") as writer:  
97     pd.DataFrame(main_data).to_excel(writer, index=False, sheet_name="Main Data")  
98     pd.DataFrame(casings_data).to_excel(writer, index=False, sheet_name="Casings")  
99  
100 print(f"Excel file created with 'Main Data' and 'Casings' at: {output_excel}")
```

## C:\Well Layer Update\Python\Azure\Data Cleaning and Matching Script.py

```

1 import pandas as pd
2 import re
3 from rapidfuzz import fuzz, process
4
5 BASE_DIR = r"<internal_project_directory>"
6 # === File Paths ===
7 input_excel = f"{BASE_DIR}\extracted_data_SCC_3.xlsx"
8 parcel_excel = f"{BASE_DIR}\Parcels.xlsx"
9 output_excel = f"{BASE_DIR}\OCR(SCC1)_3.xlsx"
10 unique_excel = f"{BASE_DIR}\unique_values.xlsx"
11
12 def clean_well_id(raw):
13     if pd.isna(raw):
14         return ""
15
16     original = str(raw).strip()
17     # Convert 'c' or 'C' to 'e' at the start
18     if original.lower().startswith('c'):
19         original = 'e' + original[1:]
20     original = re.sub(r'[oO]', '0', original)
21     cleaned = original.replace(" ", "").upper().strip()
22
23     # Handle well IDs starting with E/e
24     if original.lower().startswith('e'):
25         temp = original.lower().replace('o', '0').upper()
26         match = re.match(r"^(E)(\d{7,8})(.*)$", temp)
27         if match:
28             prefix = match.group(1)
29             digits = match.group(2)
30             return prefix + digits
31         digits_after_e = re.sub(r"\D", "", temp[1:])
32         if 6 <= len(digits_after_e) <= 8:
33             return "E" + digits_after_e
34         digits_match = re.match(r"^(E)(\d+)", temp)
35         if digits_match:
36             return digits_match.group(1) + digits_match.group(2)
37         return temp
38
39     # Preserve WCR-formatted IDs
40     if re.match(r"^\WCR\d{4}-\d{6}$", cleaned):
41         return cleaned
42
43     # Fallback cleanup for non-WCR values
44     text = original.lower().replace(", ", "").replace("no.", "").replace("/", " ")
45     text = re.sub(r"[\s\n\r]+", " ", text).strip()
46
47     # Check for date-like prefix (e.g., MM-DD-YYYY or MM.DD.YYYY)
48     match = re.match(r"^\d{1,2}[\.-](\d{1,2})[\.-](\d{2,4})(\d{5,7})$", text)

```

```

49     if match:
50         year_str = match.group(3)
51         try:
52             year = int(year_str)
53             if year < 100:
54                 year += 1900 if year >= 35 else 2000
55             if 1935 <= year <= 2025:
56                 return match.group(4)
57         except ValueError:
58             pass
59
60     # Remove date-like prefix if present
61     prefix_removed = re.sub(r"^\d{1,2}[\.-]\d{1,2}[\.-]\d{2,4}", "", text).strip()
62     segments = re.split(r"[\s]+", prefix_removed)
63
64     #Collect all digit sequences (5-8 digits) from segments
65     all_segment_digits = []
66     for segment in segments:
67         segment_digits = re.findall(r"\d{5,8}", segment)
68         all_segment_digits.extend(segment_digits)
69
70     if all_segment_digits:
71         # If multiple sequences, prefer the last one unless an earlier one is longer
72         longest = max(all_segment_digits, key=len)
73         last_segment_digits = [d for d in re.findall(r"\d{5,8}", segments[-1])] if segments
else []
74         if last_segment_digits and len(max(last_segment_digits, key=len, default="")) >=
len(longest) - 1:
75             return max(last_segment_digits, key=len)
76         return longest
77
78     # Global fallback for any remaining digit sequences
79     digits = re.findall(r"\d{5,8}", text)
80     if digits:
81         return max(digits, key=len)
82
83     return ""
84
85 # Load OCR Extracted Data
86 df = pd.read_excel(input_excel, dtype=str).fillna("")
87
88 # Load unique values for fuzzy matching
89 unique_df = pd.read_excel(unique_excel, header=None)
90 field_names = ['city', 'driller_name', 'drilling_method', 'fluid', 'test_type']
91 uniques_dict = {}
92 for idx, field in enumerate(field_names):
93     col_data = unique_df[idx].dropna().tolist()
94     if col_data:
95         uniques_dict[field] = col_data[1:] # Skip the first row (label)
96

```

```
97 #Map to df column names (lowercase as per input)
98 col_map = {
99     'city': 'city',
100    'driller_name': 'contractor',
101    'drilling_method': 'drilling_method',
102    'fluid': 'fluid',
103    'test_type': 'test_type'
104 }
105
106 # Clean each field using fuzzy matching to unique values, appending score for non-exact
    matches or '(no match)' if no match found
107 for field, col in col_map.items():
108     if col in df.columns:
109         def clean_val(val):
110             if pd.isna(val) or not str(val).strip():
111                 return val
112             val_str = str(val).strip().lower()
113             choices = uniques_dict.get(field, [])
114             if not choices:
115                 return val_str
116             choices_lower = [str(c).strip().lower() for c in choices]
117             choices_map = dict(zip(choices_lower, choices))
118             if val_str in choices_lower:
119                 return choices_map[val_str]
120             result = process.extractOne(val_str, choices_lower, scorer=fuzz.partial_ratio)
121             if result:
122                 match_lower, score, _ = result
123                 if score >= 70:
124                     return f"{choices_map.get(match_lower, match_lower)} (score={score})"
125             return f"{val_str} (no match)"
126         df[col] = df[col].apply(clean_val)
127
128 # Column names
129 APN_COL = "apn"
130 ADDR_COL = "address"
131 FALLBACK_ADDR_COL = "owner_address"
132 LOC_DESC_COL = "location_details"
133 TEMPLATE_ID_COL = "revision_cleaned"
134
135 # Define template signatures
136 TEMPLATE_SIGNATURES = {
137     "Template 1": ['1197', '790', '0503'],
138     "Template 2": ['965', '905', '968', '958', '969', '960', '246354', '354', '554', '154',
139     '334', '384', '954', '304', '254', '454', '568', '934', '87649354', '87028354'],
140     "Template 3": ['776', '1286', '1288', '775', '778', '7741'],
141     "Template 4": ['46170', '46370'],
142     "Template 5": ['115', '118', '113', '119'],
143     "Template 6": ['246', '2467'],
144     "Template 7": ['263', '253', '2634764'],
145     "Template 8": ['188'],
```

```
145     "Template 9": ['1200', '12008', '112006', '12006', '1200', '12000', '12009', '10008',
146     '40008']
147 }
148 DATE_SIGNATURES = {
149     "Template 8": ['12/19/2017'],
150     "Template 9": ['1/2006', '1/2008', '1/2000', '1/2005']
151 }
152
153 flat_signature_map = [(val, template) for template, vals in TEMPLATE_SIGNATURES.items() for
154 val in vals]
155 flat_date_signature_map = [(val, template) for template, vals in DATE_SIGNATURES.items() for
156 val in vals]
157 def preprocess_apn(apn):
158     apn = str(apn).strip()
159     if '\n' in apn:
160         parts = [p.strip() for p in apn.split('\n') if p.strip()]
161         if len(parts) == 3 and all(p.isdigit() for p in parts):
162             first, second, third = parts
163             if len(third) > len(second):
164                 second, third = third, second
165             return '-'.join([first, second, third])
166     return apn
167
168 def clean_apn_format(apn):
169     if not isinstance(apn, str):
170         return ""
171
172     apn = preprocess_apn(apn)
173
174     # Remove all non-digit characters
175     digits = re.sub(r"\D", "", apn)
176
177     # Expecting 8 digits (3+3+2 format)
178     if len(digits) == 8:
179         return f"{digits[:3]}-{digits[3:6]}-{digits[6:]}"
180
181     # Sometimes 7-digit (e.g., missing leading 0 in first section)
182     elif len(digits) == 7:
183         return f"{digits[:2].zfill(3)}-{digits[2:5]}-{digits[5:].zfill(2)}"
184
185     # Other cases invalid
186     return ""
187
188 def clean_revision(value):
189     if pd.isna(value) or not str(value).strip():
190         return ""
191     value = str(value).strip()
192     # Check for date-like patterns
```

```

193     date_patterns = [
194         r"^\d{1,2}/\d{1,2}/\d{4}$", # MM/DD/YYYY or M/D/YYYY
195         r"^\d{1,2}-\d{1,2}-\d{4}$", # MM-DD-YYYY or M-D-YYYY
196         r"^\d{1,2}/\d{4}$", # M/YYYY or MM/YYYY
197         r"^\d{1,2}-\d{4}$" # M-YYYY or MM-YYYY
198     ]
199     for pattern in date_patterns:
200         if re.match(pattern, value):
201             return value # Return original string if it's a date
202 # Otherwise, strip all non-digit characters and newlines
203 digits = re.sub(r"\D", "", value)
204 return digits
205
206 def identify_template(row):
207     text = row.get(TEMPLATE_ID_COL, "")
208     if not isinstance(text, str) or not text.strip():
209         return "Unknown"
210 # First check for date signatures (exact match)
211 for val, template in flat_date_signature_map:
212     if text == val:
213         return template
214 # Then check for numeric signatures (fuzzy match)
215 match, score, _ = process.extractOne(
216     text.strip(), [val for val, _ in flat_signature_map], scorer=fuzz.token_sort_ratio
217 )
218 if score >= 80:
219     for val, template in flat_signature_map:
220         if match == val:
221             return template
222 return "Unknown"
223
224 #Create revision_cleaned column before template identification
225 df["revision_cleaned"] = df["revision"].apply(clean_revision)
226 df["template_id"] = df.apply(identify_template, axis=1)
227
228 # Known baseline occurrences per template
229 known_baselines = {
230     "Template 1": 1, # DESTROY- (1)
231     "Template 2": 2, # destruction- (1) + Destroying- (1)
232     "Template 3": 2, # Destruction (1) + destruction- (1)
233     "Template 6": 1, # abandoned (1)
234     "Template 9": 1, # DESTROY- (1)
235     "Template 4": 0,
236     "Template 5": 0,
237     "Template 7": 0,
238     "Template 8": 0,
239     "Unknown": 0
240 }
241
242 # Baseline keywords per template (lowercase)

```

```
243 baseline_keywords = {
244     "Template 1": ['destroy'],
245     "Template 2": ['destruction', 'destroying'],
246     "Template 3": ['destruction'],
247     "Template 6": ['abandoned'],
248     "Template 9": ['destroy'],
249     "Template 4": [],
250     "Template 5": [],
251     "Template 7": [],
252     "Template 8": [],
253     "Unknown": []
254 }
255
256 # Function to process additional_destruction_keywords and set keyword_destruction_well
257 def enhance_destruction_detection(row):
258     template = row['template_id']
259     keywords_str = row.get('additional_destruction_keywords', '')
260     baseline_set = set(baseline_keywords.get(template, []))
261     destruction_well = row.get('destruction_well', ':unselected:').strip()
262
263     if not keywords_str:
264         row['additional_destruction_keywords'] = ''
265         row['keyword_destruction_well'] = destruction_well
266         return row
267
268     # Split into list, strip spaces and lowercase for comparison
269     keyword_list = [kw.strip().lower() for kw in keywords_str.split(',')]
270
271     # Filter out 'instruction' or items starting with 'instruction'
272     filtered_list = [kw for kw in keyword_list if not kw.startswith('instruction')]
273
274     # Mark keywords based on template
275     marked_list = []
276     if template == "Unknown":
277         marked_list = [f"{kw} (unknown)" for kw in filtered_list]
278     else:
279         for kw in filtered_list:
280             if kw in baseline_set:
281                 marked_list.append(f"{kw} (baseline)")
282             else:
283                 marked_list.append(kw)
284
285     # Update the column to cleaned/marked string
286     row['additional_destruction_keywords'] = ', '.join(marked_list)
287
288     # Set keyword_destruction_well: default to destruction_well
289     row['keyword_destruction_well'] = destruction_well
290
291     # For known templates, override to :selected: if keyword count exceeds baseline and
    destruction_well is :unselected:
```

```

292     if template != "Unknown":
293         total_keyword_count = len(filtered_list)
294         baseline_count = known_baselines.get(template, 0)
295         if total_keyword_count > baseline_count and destruction_well == ":unselected:":
296             row['keyword_destruction_well'] = ':selected:'
297             # Unset other fields
298             for field in ['deepen_well', 'recondition_other.', 'domestic_well',
'irrigation_well',
299                             'public_well', 'monitoring_well', 'industrial_well', 'other_well',
300                             'new_well', 'test_well']:
301                 row[field] = ':unselected:'
302
303     return row
304
305 # Apply the enhancement after template_id is set
306 df = df.apply(enhance_destruction_detection, axis=1)
307
308 # Normalize other_well_no for Template 6 and merge duplicates
309 def normalize_other_well_no(val):
310     if not isinstance(val, str):
311         return ""
312     return re.sub(r"^\w+", "", val).upper() # Remove all non-alphanumeric characters and
uppercase
313
314 # Only process Template 6
315 template6_df = df[df["template_id"] == "Template 6"].copy()
316 template6_df["other_well_no_normalized"] =
template6_df["other_well_no"].apply(normalize_other_well_no)
317
318 # Group by normalized other_well_no and merge fields
319 merged_template6 = (
320     template6_df.groupby("other_well_no_normalized", as_index=False)
321     .agg(lambda x: next((i for i in x if i.strip()), "") if x.dtype == "O" else x.iloc[0])
322 )
323
324 # Drop duplicates in original df
325 df = df[df["template_id"] != "Template 6"]
326
327 # Append merged template6 records back
328 df = pd.concat([df, merged_template6], ignore_index=True)
329
330 # Load Casings Sheet
331 casings_df = pd.read_excel(input_excel, sheet_name="Casings", dtype=str).fillna("")
332 casings_df["row_index"] = casings_df["row_index"].astype(int)
333 casings_df["top"] = casings_df["top"].astype(str)
334 casings_df["bottom"] = casings_df["bottom"].astype(str)
335
336 def clean_to_digits(value, is_top=False):
337     if pd.isna(value) or not str(value).strip():
338         return None
339     # Extract the first sequence of digits before a non-digit character

```

```

340     match = re.match(r"^\d+", str(value).strip())
341     if match:
342         cleaned = match.group(0)
343         if not cleaned.isdigit():
344             return None
345         result = int(cleaned)
346         # For top values > 999 ending with '0', strip the trailing '0'
347         if is_top and result > 999 and cleaned.endswith('0'):
348             result = int(cleaned[:-1]) # Remove the last digit ('0')
349         return result
350     return None
351
352 def extract_perforation_fields(row):
353     template = row.get("template_id")
354     source_file = row.get("source_file", "").strip()
355
356     if template == "Template 1":
357         selected_rows = []
358         for i in range(1, 7):
359             colname = f"template1_screen_row{i}"
360             if row.get(colname, "").strip() == ":selected:":
361                 selected_rows.append(i)
362
363         if not selected_rows:
364             return pd.Series({"perforation_interval": "", "top_perforations": "",
"bottom_perforations": ""})
365
366         matches = casings_df[
367             (casings_df["source_file"] == source_file) &
368             (casings_df["row_index"].isin(selected_rows))
369         ].copy()
370
371         if matches.empty:
372             return pd.Series({"perforation_interval": "", "top_perforations": "",
"bottom_perforations": ""})
373
374         matches["top_cleaned"] = matches["top"].apply(lambda x: clean_to_digits(x,
is_top=True))
375         matches["bottom_cleaned"] = matches["bottom"].apply(clean_to_digits)
376
377         top = matches["top_cleaned"].min()
378         bottom = matches["bottom_cleaned"].max()
379
380         if pd.isna(top) or pd.isna(bottom):
381             return pd.Series({"perforation_interval": "", "top_perforations": "",
"bottom_perforations": ""})
382
383         # Check if bottom is less than top
384         if bottom < top:
385             return pd.Series({
386                 "perforation_interval": "",

```

```

387         "top_perforations": top,
388         "bottom_perforations": ""
389     })
390
391     return pd.Series({
392         "perforation_interval": f"{top} - {bottom}",
393         "top_perforations": top,
394         "bottom_perforations": bottom
395     })
396
397     elif template in {"Template 2", "Template 3", "Template 4", "Template 6", "Template 8",
"Template 9", "Unknown"}:
398         filtered = casings_df[casings_df["source_file"] == source_file].copy()
399
400         if template in {"Template 8", "Template 9"}:
401             if template == "Template 8":
402                 has_empty_type = (filtered["type"] == "").any()
403                 if has_empty_type:
404                     # Misidentified Template 8: process all rows (like Templates 2, 3, 4, 6,
Unknown)
405                     pass
406                 else:
407                     # True Template 8: only "Screen" rows
408                     filtered = filtered[filtered["type"] == "Screen"]
409             elif template == "Template 9":
410                 filtered = filtered[filtered["type"] == "Screen"]
411
412         if filtered.empty:
413             return pd.Series({"perforation_interval": "", "top_perforations": "",
"bottom_perforations": ""})
414
415         filtered["top_cleaned"] = filtered["top"].apply(lambda x: clean_to_digits(x,
is_top=True))
416         filtered["bottom_cleaned"] = filtered["bottom"].apply(clean_to_digits)
417
418         filtered = filtered[
419             (filtered["top_cleaned"].notnull()) &
420             (filtered["bottom_cleaned"].notnull()) &
421             (filtered["bottom_cleaned"] < 2000)
422         ]
423
424         if filtered.empty:
425             return pd.Series({"perforation_interval": "", "top_perforations": "",
"bottom_perforations": ""})
426
427         top_row = filtered.loc[filtered["row_index"].idxmin()]
428         bottom_row = filtered.loc[filtered["row_index"].idxmax()]
429
430         top = top_row["top_cleaned"]
431         bottom = bottom_row["bottom_cleaned"]
432

```

```
433     if top is None or bottom is None:
434         return pd.Series({"perforation_interval": "", "top_perforations": "",
"bottom_perforations": ""})
435
436     # Check if bottom is less than top
437     if bottom < top:
438         return pd.Series({
439             "perforation_interval": "",
440             "top_perforations": top,
441             "bottom_perforations": ""
442         })
443
444     return pd.Series({
445         "perforation_interval": f"{top} - {bottom}",
446         "top_perforations": top,
447         "bottom_perforations": bottom
448     })
449
450     elif template == "Template 5":
451         raw_text = row.get("perforations_template5", "")
452         if not raw_text.strip():
453             return pd.Series({"perforation_interval": "", "top_perforations": "",
"bottom_perforations": ""})
454
455         ranges = re.findall(r"(\d{2,4})\s*[--]\s*(\d{2,4})", raw_text)
456         numbers = []
457
458         for start, end in ranges:
459             try:
460                 start_int = int(start)
461                 end_int = int(end)
462                 if 10 < start_int < 2000 and 10 < end_int < 2000:
463                     numbers.extend([start_int, end_int])
464             except ValueError:
465                 continue
466
467         singles = re.findall(r"\b(\d{2,4})\b", raw_text)
468         for val in singles:
469             try:
470                 num = int(val)
471                 if 10 < num < 2000:
472                     numbers.append(num)
473             except ValueError:
474                 continue
475
476         if numbers:
477             top = min(numbers)
478             bottom = max(numbers)
479             # Check if bottom is less than top
480             if bottom < top:
```

```
481         return pd.Series({
482             "perforation_interval": "",
483             "top_perforations": top,
484             "bottom_perforations": ""
485         })
486     return pd.Series({
487         "perforation_interval": f"{top} - {bottom}",
488         "top_perforations": top,
489         "bottom_perforations": bottom
490     })
491
492     return pd.Series({"perforation_interval": "", "top_perforations": "",
493 "bottom_perforations": ""})
494
495     elif template == "Template 7":
496         raw_text = row.get("perforations_template7", "")
497         if not raw_text.strip():
498             return pd.Series({"perforation_interval": "", "top_perforations": "",
499 "bottom_perforations": ""})
500
501         ranges = re.findall(r"(\d{2,4})\s*[-]\s*(\d{2,4})", raw_text)
502         singles = re.findall(r"\b(\d{2,4})\b", raw_text)
503         numbers = []
504
505         for start, end in ranges:
506             try:
507                 start_int = int(start)
508                 end_int = int(end)
509                 if 10 < start_int < 2000 and 10 < end_int < 2000:
510                     numbers.extend([start_int, end_int])
511             except ValueError:
512                 continue
513
514         for val in singles:
515             try:
516                 num = int(val)
517                 if 10 < num < 2000:
518                     numbers.append(num)
519             except ValueError:
520                 continue
521
522         if numbers:
523             top = min(numbers)
524             bottom = max(numbers)
525             # Check if bottom is less than top
526             if bottom < top:
527                 return pd.Series({
528                     "perforation_interval": "",
529                     "top_perforations": top,
530                     "bottom_perforations": ""
```

```
529         })
530     return pd.Series({
531         "perforation_interval": f"{top} - {bottom}",
532         "top_perforations": top,
533         "bottom_perforations": bottom
534     })
535
536     return pd.Series({"perforation_interval": "", "top_perforations": "",
537 "bottom_perforations": ""})
538
539     return pd.Series({"perforation_interval": "", "top_perforations": "",
540 "bottom_perforations": ""})
541
542 df[["perforation_interval", "top_perforations", "bottom_perforations"]] =
543 df.apply(extract_perforation_fields, axis=1)
544
545 # Clean Well ID
546 df["well_id_cleaned"] = df.iloc[:, 0].apply(clean_well_id)
547
548 # Extract potential APN pattern from a text string
549 def extract_apn_from_text(text):
550     if not isinstance(text, str):
551         return ""
552     text = text.upper().strip()
553
554     text = text.replace("\n", " ").replace("APN#", "APN").replace("A.P.N.", "APN")
555
556     match = re.search(r"\b(\d{2,3}-\d{3}-\d{2}[A-Z]?)\b", text)
557     if match:
558         return match.group(1)
559     return ""
560
561 def fallback_apn_with_source(row):
562     apn_raw = str(row.get(APN_COL, "")).strip()
563     apn_digits = clean_apn_format(apn_raw)
564
565     if apn_digits:
566         return apn_digits.zfill(7), "Original APN field"
567
568     apn_candidate = extract_apn_from_text(row.get("location_details", ""))
569     source = "location_details"
570
571     if not apn_candidate:
572         apn_candidate = extract_apn_from_text(row.get("owner_well_id", ""))
573         source = "owner_well_id"
574
575     if apn_candidate:
576         digits_only = clean_apn_format(apn_candidate)
577         return digits_only.zfill(7), f"Extracted from {source}"
578
579     return "", "No APN found"
```

```
577
578 df[["apn_cleaned", "apn_cleaned_source"]] = df.apply(fallback_apn_with_source, axis=1,
579 result_type="expand")
580 # Define cleaning function for pump_test_length
581 def clean_pump_test_length(value):
582     if pd.isna(value) or not str(value).strip():
583         return ""
584
585     value = str(value).strip()
586
587     # Return empty string if slash is present
588     if '/' in value:
589         return ""
590
591     # Remove any digits (and subsequent characters) after a non-digit, non-period character
592     match = re.match(r'^[\d.]+', value)
593     if not match:
594         return ""
595
596     cleaned = match.group(0)
597
598     # Keep only digits and periods
599     cleaned = re.sub(r'^[\0-9.]+', '', cleaned)
600
601     # Ensure valid number format (remove multiple periods, trailing/leading periods)
602     if cleaned.count('.') > 1:
603         parts = cleaned.split('.')
604         cleaned = parts[0] + '.' + ''.join(parts[1:])
605     cleaned = cleaned.strip('.')
606
607     # Return empty string if result is empty or invalid
608     return cleaned if cleaned and cleaned != '.' else ""
609
610 numeric_columns = {
611     "water_static": "water_static_cleaned",
612     "depth_well": "depth_well_cleaned",
613     "yield": "yield_cleaned",
614     "seal": "seal_cleaned",
615     "first_water": "first_water_cleaned",
616     "total_depth": "total_depth_cleaned",
617     "lat": "lat_cleaned",
618     "long": "long_cleaned",
619     "pump_test_length": "pump_test_length_cleaned"
620 }
621
622 def clean_numeric_field(value):
623     if not isinstance(value, str):
624         value = str(value)
625
```

```

626     value = value.strip().lower()
627
628     # Handle newline: keep only the first number before the newline
629     if '\n' in value:
630         value = value.split('\n')[0].strip()
631
632     # Remove apostrophe and any digits following it
633     value = re.sub(r"'s*\d*(\"|'|\')?", "", value)
634
635     value = re.sub(r'^[^\d]+', '', value)      #Remove non-digit prefix
636     value = re.sub(r'^\d+$$', '', value)
637
638     range_pattern = r"^\s*(\d+(?:\.\d+)?)\s*(?:[--]|to|t)\s*(\d+(?:\.\d+)?)\s*$"
639     range_match = re.match(range_pattern, value)
640     if range_match:
641         try:
642             num1 = float(range_match.group(1))
643             num2 = float(range_match.group(2))
644             avg = (num1 + num2) / 2
645             return f"{avg:.2f}".rstrip('0').rstrip('.')
646         except ValueError:
647             pass
648
649     value = re.sub(r'(?<=\d),(?=\d)', '.', value)
650     value = re.sub(r'^0-9.', '', value)
651
652     if value.count('.') > 1:
653         parts = value.split('.')
654         value = parts[0] + '.' + ''.join(parts[1:])
655
656     return value
657
658 for colname, new_col in numeric_columns.items():
659     if colname == "pump_test_length":
660         df[new_col] = df[colname].apply(clean_pump_test_length)
661     else:
662         df[new_col] = df[colname].apply(clean_numeric_field)
663
664
665 # Initialize qa_flag column
666 df['depths_flag'] = ''
667
668 # Only check rows where both fields are non-empty
669 mask = (df['total_depth_cleaned'] != '') & (df['depth_well_cleaned'] != '')
670
671 # Apply logic: if total_depth < depth_well → flag 'review depths'
672 df.loc[mask, 'depths_flag'] = df.loc[mask].apply(
673     lambda row: 'review depths' if float(row['total_depth_cleaned']) <
674     float(row['depth_well_cleaned']) else '',

```

```
675 )
676
677
678 # Load and structure alias dictionary
679 alias_df = pd.read_excel(r"C:\Well Layer Update\GIS\data_tables\street_names_with_aliases.xlsx").fillna("")
680 alias_dict = {}
681 for _, row in alias_df.iterrows():
682     canonical = row['canonical_name'].strip().upper()
683     for alias in str(row['alias']).split(";"):
684         alias = alias.strip().upper()
685         if alias:
686             alias_dict[alias] = canonical
687
688 for canonical in alias_df["canonical_name"].unique():
689     canonical = canonical.strip().upper()
690     alias_dict[canonical] = canonical
691
692 street_keywords = set(alias_dict.keys())
693
694 def extract_address_from_location(desc):
695     if not isinstance(desc, str):
696         return ""
697     desc = desc.upper()
698     desc = re.sub(r"^[^w\s]", " ", desc)
699     desc = re.sub(r"\s+", " ", desc).strip()
700
701     if re.search(r"\b(UTM|TOPO|SHEET|1[;:,\\.]?62[,\.\.]?500|GRID|QUAD|COORDINATE|LAT|LONG|MAP)\b", desc):
702         return ""
703
704     if re.search(r"\b\d{1,5}\s*(FEET|FT|YARDS|MILES?)\b", desc):
705         return ""
706
707     tokens = desc.split()
708     num_tokens = len(tokens)
709
710     for i in range(num_tokens):
711         for j in range(1, 5):
712             if i + j > num_tokens:
713                 continue
714             candidate_tokens = tokens[i:i + j]
715             candidate = " ".join(candidate_tokens)
716             if candidate in street_keywords:
717                 canonical = alias_dict.get(candidate, candidate)
718                 if i > 0 and re.fullmatch(r"\d{2,6}", tokens[i - 1]):
719                     return f"{tokens[i - 1]} {canonical}"
720             else:
721                 return canonical
722
```

```

723     for i in range(num_tokens):
724         for j in range(1, 5):
725             if i + j > num_tokens:
726                 continue
727             candidate_tokens = tokens[i:i + j]
728             candidate = " ".join(candidate_tokens)
729             best_match, score, _ = process.extractOne(
730                 candidate,
731                 list(street_keywords),
732                 scorer=fuzz.token_sort_ratio
733             )
734             if score >= 85:
735                 canonical = alias_dict.get(best_match, best_match)
736                 if i > 0 and re.fullmatch(r"\d{2,6}", tokens[i - 1]):
737                     return f"{tokens[i - 1]} {canonical}"
738                 else:
739                     return canonical
740
741     return ""
742
743 def resolve_address_with_source(row):
744     addr = str(row.get(ADDR_COL, "")).strip()
745     fallback = str(row.get(FALLBACK_ADDR_COL, "")).strip()
746     loc_extracted = extract_address_from_location(row.get(LOC_DESC_COL, ""))
747
748     fallback_keywords = ["same", "dame", "come", "came", "some", "sama", "same as above"]
749     if any(fuzz.ratio(addr.lower(), word) >= 80 for word in fallback_keywords):
750         return fallback, "Fallback address (approximate 'same')"
751
752     if addr:
753         return addr, "Direct address"
754     if loc_extracted:
755         return loc_extracted, "Extracted from location description"
756     if fallback:
757         return fallback, "Fallback address (owners)"
758     return "", "No address available"
759
760 def strip_city_and_zip(address):
761     if not address:
762         return address
763
764     cities = [
765         "APTOS", "BEN LOMOND", "BOULDER CREEK", "CAPITOLA", "DAVENPORT",
766         "LOS GATOS", "SANTA CRUZ", "SCOTTS VALLEY", "SOQUEL", "WATSONVILLE"
767     ]
768     address = address.upper()
769
770     house_number_match = re.match(r"^\d+\s+", address)
771     if not house_number_match:
772         return address

```

```

773
774     start = house_number_match.end()
775     post_number = address[start:]
776     tokens = post_number.split()
777     retained_tokens = tokens[:1]
778     remaining_tokens = tokens[1:]
779
780     for idx, token in enumerate(remaining_tokens):
781         joined = " ".join([*retained_tokens, *remaining_tokens[:idx + 1]])
782         for city in cities:
783             if city in joined:
784                 return address[:address.upper().find(city)].strip()
785
786     return address
787
788 def clean_for_matching(addr):
789     if not addr:
790         return ""
791     addr = addr.upper()
792     addr = re.sub(r"\s*\n\s*", " ", addr)
793     addr = re.sub(r"\s+", " ", addr).strip()
794
795
796     remove_keywords = [
797         "P.O. BOX", "PO BOX", "O. BOX", "P.O.", "PO", "BOX", "PABOX", "DRAWER",
798         "ABOVE", "OFF", "MAP", "MAPSHEET", "APN", "NO ADDRESS", "END OF", "ENDOF",
799         "PART OF", "P/R", "SUITE", "C/O", "SHOPS", "NOTA", "WIR", "EX", "FEET NORTH OF"
800     ]
801     for keyword in remove_keywords:
802         addr = re.sub(rf"\b{re.escape(keyword)}\b", "", addr)
803
804     #dictionary and spell-checker for common OCR misspellings
805     ocr_fixes = {
806         r"\bARE\b": "AVE", r"\bSTH\b": "ST", r"\bDRI\b": "DR", r"5T": "ST",
807         r"\bOAK STREE\b": "OAK ST", r"(\d)O(\d)": r"\1 0 \2",
808         r"\bMOUNT\b": "MT", r"\bM\b": "MT", r"\bR\b": "RD",
809         r"\bDRIE\b": "DR", r"\bQUINET\b": "QUINNETTE", r"\bJANE\b": "LANE",
810         r"\bCEUK\b": "CREEEK", r"\bTHNY\b": "TWINY", r"\bELTON\b": "FELTON",
811         r"\bBENLOMOND\b": "BEN LOMOND", r"\bBENLAMIND\b": "BEN LOMOND",
812         r"\bBEN SOMAND\b": "BEN LOMOND", r"\bBEN HOWARD\b": "BEN LOMOND",
813         r"\bBONNY DOON\b": "BONNY DOON", r"\bBUNNY DOON\b": "BONNY DOON",
814         r"\bBROOKDALE\b": "", r"\bSANTACRUZ\b": "SANTA CRUZ",
815         r"\bSANTACURZ\b": "SANTA CRUZ", r"\bCUPERTING\b": "CUPERTINO",
816         r"\bROALND\b": "ROLAND", r"\bZIMT\b": "", r"\bHERMOURA\b": "HERMOSA",
817         r"\bDALTON\b": "", r"\bSA\b": "CA", r"\bCALIFÓRNIA\b": "CA",
818         r"\bCALLIL\b": "CA", r"\bCALIF\b": "CA", r"\bCALL\b": "CA", r"\bCAO\b": "CA",
819         r"\bCAMILLA\b": "CAMPBELL", r"\bLETT?S VALLY\b": "LETTS VALLEY",
820         r"\bENGINE GRADE\b": "EMPIRE GRADE", r"\bFORMAND\b": "LOMOND",
821         r"\bSTREEL\b": "STREET", r"\bARE\b": "AVE", r"\bMª\b": "MA", r"\bª\b": "A"
822     }

```

```

823     for pattern, replacement in ocr_fixes.items():
824         addr = re.sub(pattern, replacement, addr)
825
826     # Remove irrelevant leading words found in well reports
827     addr = re.sub(r"^(WELL( SITE| HEAD| SA| LOCATION| HOUSE)?|ADDRESS|LOCATED AT|NEAR|NEXT
TO|ADJACENT TO|BEHIND|ADJOINING)\s+", "", addr)
828
829     # Normalize street suffixes for consistent matching
830     replacements = {
831         " STREET": " ST", " AVENUE": " AVE", " ROAD": " RD", " DRIVE": " DR",
832         " BOULEVARD": " BLVD", " MOUNT ": " MT ", " HIGHWAY": " HWY",
833         " LANE": " LN", " COURT": " CT", " PLACE": " PL", " TERRACE": " TER",
834         " CIRCLE": " CIR", " TRAIL": " TRL", " CREEK": " CRK", " WAY": " WAY"
835     }
836     for k, v in replacements.items():
837         addr = addr.replace(k, v)
838
839     addr = re.sub(r"\s+[A-Z\s]*\sCA\s\d{5,6}$", "", addr)
840     addr = re.sub(r"^[^w\s]+$", "", addr)
841     addr = re.sub(r"\b[w\b]$", "", addr).strip()
842     addr = strip_city_and_zip(addr)
843     return addr
844
845 df[["address_raw", "address_source"]] = df.apply(lambda row: pd.Series(resolve_address_with_
_source(row)), axis=1)
846 df["address_cleaned"] = df["address_raw"].astype(str).apply(lambda x: x.splitlines()[0] if
x.strip() else "").str.strip().str.upper()
847 df["address_cleaned_for_matching"] = df["address_cleaned"].apply(clean_for_matching)
848
849 # Load Parcel Data and Prepare Matching Columns
850 parcels_df = pd.read_excel(parcel_excel, dtype=str).fillna("")
851 parcels_df.columns = parcels_df.columns.str.lower()
852 parcels_df["clean_apn"] = parcels_df.iloc[:, 0]
853 parcels_df["clean_address_parcel"] = parcels_df.iloc[:,
13].astype(str).str.strip().str.upper()
854 parcels_df["match_ready_address"] = parcels_df["clean_address_parcel"]
855
856 apn_lookup = parcels_df.set_index("clean_apn")[["clean_address_parcel"]].to_dict("index")
857
858 def fuzzy_address_lookup(addr):
859     cleaned_input = clean_for_matching(addr)
860     if not cleaned_input:
861         return "", "", [], 0
862
863     input_house_num = re.match(r"^(\\d+)", cleaned_input)
864     input_house_num = input_house_num.group(1) if input_house_num else ""
865
866     matches = process.extract(
867         cleaned_input,
868         parcels_df["match_ready_address"].tolist(),
869         scorer=fuzz.token_sort_ratio,

```

```

870     limit=30
871 )
872
873 candidates_log = []
874 best_row, best_score, best_dist = None, 0, float("inf")
875
876 for match_text, score, idx in matches:
877     row = parcels_df.iloc[idx]
878     match_addr = row["clean_address_parcel"]
879     match_num = re.match(r"^\d+", match_addr)
880     match_num = match_num.group(1) if match_num else ""
881     dist = abs(int(input_house_num) - int(match_num)) if input_house_num and match_num
else float("inf")
882     candidates_log.append(f"{match_addr} (score={score}, dist={dist})")
883
884     if score > best_score or (score == best_score and dist < best_dist):
885         best_score, best_dist = score, dist
886         best_row = row
887
888     if best_row is not None:
889         return best_row["clean_apn"], best_row["clean_address_parcel"], candidates_log,
best_score
890
891     return "", "", candidates_log, 0
892
893 # Normalize well_id_cleaned for duplicate comparison (strip leading zeros)
894 def normalize_well_id_for_comparison(well_id):
895     if not isinstance(well_id, str) or not well_id.strip():
896         return ""
897     # Remove leading zeros, but preserve the rest of the ID (including 'E' prefix if
present)
898     if well_id.startswith('E'):
899         return 'E' + well_id[1:].lstrip('0')
900     return well_id.lstrip('0')
901
902 # Remove Duplicates Based on well_id_cleaned and multiple fields, including all-empty
matches
903 def count_non_empty(row):
904     """Count non-empty, non-whitespace fields in a row."""
905     return sum(1 for val in row if isinstance(val, str) and val.strip() != "" or not
pd.isna(val))
906
907 # Add duplicate details column
908 df['duplicate_exists'] = "No duplicates"
909
910 # Add a column to count non-empty fields for each row
911 df["non_empty_count"] = df.apply(count_non_empty, axis=1)
912
913 # Add a column for normalized well_id_cleaned for duplicate comparison
914 df["well_id_normalized"] = df["well_id_cleaned"].apply(normalize_well_id_for_comparison)
915

```

```
916 # Identify duplicates based on normalized well_id_cleaned and any field match or all fields
    empty
917 fields_to_check = [
918     "total_depth_cleaned",
919     "depth_well_cleaned",
920     "water_static_cleaned",
921     "yield_cleaned",
922     "seal_cleaned",
923     "first_water_cleaned",
924     "permit_date",
925     "permit_no",
926     "date",
927     "apn_cleaned"
928 ]
929
930 # Filter for non-empty well_id_cleaned
931 non_empty_well_id_mask = df["well_id_cleaned"].notna() & (df["well_id_cleaned"] != "")
932 duplicates = pd.Series([False] * len(df), index=df.index)
933 for field in fields_to_check:
934     # Mark duplicates where normalized well_id_cleaned and field match, and both are non-
    empty
935     mask = (
936         df.duplicated(subset=["well_id_normalized", field], keep=False) &
937         (df[field].notna() & (df[field] != "")) &
938         non_empty_well_id_mask
939     )
940     duplicates |= mask
941
942 # Add duplicates where all fields are empty for a given normalized well_id_cleaned
943 all_empty_mask = df[fields_to_check].apply(lambda x: x.isna() | (x == "")).all(axis=1)
944 all_empty_duplicates = (
945     df[all_empty_mask & non_empty_well_id_mask].duplicated(subset=["well_id_normalized"],
    keep=False)
946 )
947 duplicates |= all_empty_duplicates
948
949 # Identify duplicate groups for logging
950 duplicate_groups = df[duplicates & non_empty_well_id_mask][["well_id_normalized",
    "well_id_cleaned", "total_depth_cleaned", "depth_well_cleaned"]].drop_duplicates()
951
952 print(f"Found {len(duplicate_groups)} groups of duplicates based on normalized
    well_id_cleaned and multiple fields or all fields empty")
953
954 # Process each duplicate group
955 records_to_keep = set()
956 for _, group in duplicate_groups.iterrows():
957     well_id_norm = group["well_id_normalized"]
958     total_depth = group["total_depth_cleaned"]
959     depth_well = group["depth_well_cleaned"]
960
```

```

961 # Get all records in this duplicate group (matching normalized well_id_cleaned and any
    field or all fields empty)
962 group_records = df[
963     (df["well_id_normalized"] == well_id_norm) & (
964         (df["total_depth_cleaned"] == total_depth) |
965         (df["depth_well_cleaned"] == depth_well) |
966         (df["water_static_cleaned"] == df[df["well_id_normalized"] == well_id_norm]
["water_static_cleaned"].iloc[0])) |
967         (df["yield_cleaned"] == df[df["well_id_normalized"] == well_id_norm]
["yield_cleaned"].iloc[0])) |
968         (df["seal_cleaned"] == df[df["well_id_normalized"] == well_id_norm]
["seal_cleaned"].iloc[0])) |
969         (df["first_water_cleaned"] == df[df["well_id_normalized"] == well_id_norm]
["first_water_cleaned"].iloc[0])) |
970         (df["permit_date"] == df[df["well_id_normalized"] == well_id_norm]
["permit_date"].iloc[0])) |
971         (df["permit_no"] == df[df["well_id_normalized"] == well_id_norm]
["permit_no"].iloc[0])) |
972         (df["date"] == df[df["well_id_normalized"] == well_id_norm]["date"].iloc[0])) |
973         (df["apn_cleaned"] == df[df["well_id_normalized"] == well_id_norm]
["apn_cleaned"].iloc[0])) |
974         (
975             df[fields_to_check].apply(lambda x: x.isna() | (x == "")).all(axis=1) &
976             (df["well_id_normalized"] == well_id_norm)
977         )
978     ) & non_empty_well_id_mask
979 ]
980
981 if len(group_records) > 1:
982     # Sort by non_empty_count (descending), then by source_file (prefer newer or cleaner
files)
983     most_complete = group_records.sort_values(
984         by=["non_empty_count", "source_file"], ascending=[False, True]
985     ).iloc[0]
986     df.at[most_complete.name, 'duplicate_exists'] = "Duplicates deleted (matched depths,
water_static, yield, seal, first_water, permit_date, permit_no, date, apn_cleaned, or all
fields empty, including leading zero variations)"
987     records_to_keep.add(most_complete.name)
988     else:
989         records_to_keep.add(group_records.index[0])
990
991 # Keep non-duplicate records and the most complete record from each duplicate group
992 non_duplicates = df[~duplicates]
993 duplicates_to_keep = df.loc[list(records_to_keep)]
994 df = pd.concat([non_duplicates, duplicates_to_keep], ignore_index=True)
995
996 # Ensure records with empty well_id_cleaned are marked as "No duplicates"
997 df.loc[df["well_id_cleaned"].isna() | (df["well_id_cleaned"] == ""), "duplicate_exists"] =
"No duplicates"
998
999 # Flag distinct duplicates in the final DataFrame, excluding empty well_id_cleaned

```

```

1000 well_id_counts = df[df["well_id_cleaned"].notna() & (df["well_id_cleaned"] != "")]
1001     ["well_id_normalized"].value_counts()
1002
1003 duplicate_well_ids = well_id_counts[well_id_counts > 1].index
1004
1005 for well_id_norm in duplicate_well_ids:
1006     mask = df["well_id_normalized"] == well_id_norm
1007     for idx in df[mask].index:
1008         current = df.at[idx, 'duplicate_exists']
1009         if current == "No duplicates":
1010             df.at[idx, 'duplicate_exists'] = "Distinct duplicates exist (different depths,
water_static, yield, seal, first_water, permit_date, permit_no, date, apn_cleaned, and not
all empty, including leading zero variations)"
1011         elif current == "Duplicates deleted (matched depths, water_static, yield, seal,
first_water, permit_date, permit_no, date, apn_cleaned, or all fields empty, including
leading zero variations)":
1012             df.at[idx, 'duplicate_exists'] = "Duplicates deleted (matched depths,
water_static, yield, seal, first_water, permit_date, permit_no, date, apn_cleaned, or all
fields empty, including leading zero variations) and distinct duplicates exist"
1013
1014 # Drop temporary columns
1015 df = df.drop(columns=["non_empty_count", "well_id_normalized"])
1016
1017 # Verify no duplicates remain
1018 remaining_duplicates = df[
1019     (df["well_id_cleaned"].notna() & (df["well_id_cleaned"] != "")) & (
1020         df.duplicated(subset=["well_id_cleaned", "total_depth_cleaned"], keep=False) |
1021         df.duplicated(subset=["well_id_cleaned", "depth_well_cleaned"], keep=False) |
1022         df.duplicated(subset=["well_id_cleaned", "water_static_cleaned"], keep=False) |
1023         df.duplicated(subset=["well_id_cleaned", "yield_cleaned"], keep=False) |
1024         df.duplicated(subset=["well_id_cleaned", "seal_cleaned"], keep=False) |
1025         df.duplicated(subset=["well_id_cleaned", "first_water_cleaned"], keep=False) |
1026         df.duplicated(subset=["well_id_cleaned", "permit_date"], keep=False) |
1027         df.duplicated(subset=["well_id_cleaned", "permit_no"], keep=False) |
1028         df.duplicated(subset=["well_id_cleaned", "date"], keep=False) |
1029         df.duplicated(subset=["well_id_cleaned", "apn_cleaned"], keep=False) |
1030         (
1031             df[fields_to_check].apply(lambda x: x.isna() | (x == "")).all(axis=1) &
1032             df.duplicated(subset=["well_id_cleaned"], keep=False)
1033         )
1034     )
1035 ]
1036 if not remaining_duplicates.empty:
1037     print("flag: some duplicates remain")
1038     print(remaining_duplicates[["well_id_cleaned", "total_depth_cleaned",
"depth_well_cleaned", "water_static_cleaned", "yield_cleaned", "seal_cleaned",
"first_water_cleaned", "permit_date", "permit_no", "date", "apn_cleaned", "source_file"]])
1039 else:
1040     print("No duplicates remain.")
1041
1042 # Fuzzy Matching Loop
1043 final_apn, final_address, match_source, fuzzy_logs = [], [], [], []

```

```
1042
1043 for i, (_, row) in enumerate(df.iterrows()):
1044     print(f"Processing record {i + 1} of {len(df)}")
1045     cleaned_apn = row["apn_cleaned"]
1046     cleaned_addr = row["address_cleaned"]
1047     source = row["address_source"]
1048
1049     if cleaned_apn in apn_lookup:
1050         final_apn.append(cleaned_apn)
1051         final_address.append(apn_lookup[cleaned_apn]["clean_address_parcel"])
1052         match_source.append("APN matched")
1053         fuzzy_logs.append("")
1054     else:
1055         if source == "Extracted from location description":
1056             fuzzy_threshold = 75
1057         elif "Fallback address" in source:
1058             fuzzy_threshold = 85
1059         else:
1060             fuzzy_threshold = 62
1061
1062         fuzzy_apn, fuzzy_addr, log, score = fuzzy_address_lookup(cleaned_addr)
1063         if fuzzy_apn and score >= fuzzy_threshold:
1064             final_apn.append(fuzzy_apn)
1065             final_address.append(fuzzy_addr)
1066             if "Fallback address" in source:
1067                 match_source.append(f"Owners address (fallback) matched (score {score})")
1068             else:
1069                 match_source.append(f"Address matched (score {score})")
1070         else:
1071             final_apn.append(cleaned_apn)
1072             final_address.append("")
1073             match_source.append("No match")
1074             fuzzy_logs.append("; ".join(log))
1075
1076     planned_use_columns = {
1077         "domestic_well": "Domestic Well",
1078         "irrigation_well": "Irrigation Well",
1079         "public_well": "Public Supply Well",
1080         "monitoring_well": "Monitoring Well",
1081         "industrial_well": "Industrial Well",
1082         "other_well": "Other Well",
1083         "test_well": "Test Well"
1084     }
1085
1086     def determine_planned_use(row):
1087         selected = [
1088             label for col, label in planned_use_columns.items()
1089             if row.get(col, "").strip() == ":selected:"
1090         ]
1091         if len(selected) == 1:
```

```
1092     return selected[0]
1093     elif len(selected) > 1:
1094         return "Two or more uses/Possible Error from 'Confidential' Stamp"
1095     else:
1096         return "Other"
1097
1098 df["planned_use_former_use"] = df.apply(determine_planned_use, axis=1)
1099
1100 type_of_work_columns = {
1101     "new_well": "New Well",
1102     "keyword_destruction_well": "Well Destruction",
1103     "deepen_well": "Deepen Well",
1104     "replacement": "Replacement Well",
1105     "recondition_other.": "Well Reconditioning/Other"
1106 }
1107
1108 def determine_type_of_work(row):
1109     selected = [
1110         label for col, label in type_of_work_columns.items()
1111         if row.get(col, "").strip() == ":selected:"
1112     ]
1113     if len(selected) == 1:
1114         return selected[0]
1115     elif len(selected) > 1:
1116         return "Two or more types"
1117     else:
1118         return "Unknown"
1119
1120 df["Type of Work"] = df.apply(determine_type_of_work, axis=1)
1121
1122 def map_record_type(type_of_work):
1123     if type_of_work == "New Well":
1124         return "WellCompletion/New/Production or Monitoring/NA"
1125     elif type_of_work == "Unknown":
1126         return "Unknown"
1127     elif type_of_work in ["Drill and Destroy", "Destroy"]:
1128         return "WellCompletion/Drill and Destroy/NA/NA"
1129     elif type_of_work == "Well Destruction":
1130         return "WellCompletion/Destruction/NA/NA"
1131     elif type_of_work in ["Deepen Well", "Replacement Well", "Well Reconditioning/Other"]:
1132         return "WellCompletion/Modification or Repair/Production or Monitoring/NA"
1133     else:
1134         return "Unknown"
1135
1136 df["Record Type"] = df["Type of Work"].apply(map_record_type)
1137
1138 df["final_apn"] = final_apn
1139 df["final_address"] = final_address
1140 df["match_source"] = match_source
1141 df["fuzzy_candidates_log"] = fuzzy_logs
```

```
1142
1143 df.loc[df["template_id"] == "Template 8", "planned_use_former_use"] =
df.loc[df["template_id"] == "Template 8", "Planned Use"]
1144 df.loc[df["template_id"] == "Template 8", "Type of Work"] = df.loc[df["template_id"] ==
"Template 8", "Activity"]
1145 df["Record Type"] = df["Type of Work"].apply(map_record_type)
1146
1147 df.to_excel(output_excel, index=False)
1148 print(f"Cleaned data written to: {output_excel}")
```

## C:\Well Layer Update\Python\Azure\Data Matching and Merging Script.py

```
1 import pandas as pd
2 import re
3 from rapidfuzz import process, fuzz
4 import logging
5 import numpy as np
6 from collections import defaultdict
7
8 BASE_DIR = r"<internal_project_directory>"
9
10 # === File paths ===
11 file1_path = f"{BASE_DIR}\OCR(SCC1).xlsx"
12 file2_path = f"{BASE_DIR}\DWR1.xlsx"
13 file_monterey = f"{BASE_DIR}\Well Completion Reports Monterey.xlsx"
14 matched_output = f"{BASE_DIR}\matched_output.xlsx"
15 unmatched_output = f"{BASE_DIR}\unmatched_output.xlsx"
16 unmatched_ocr_output = f"{BASE_DIR}\unmatched_ocr_output.xlsx"
17 unmatched_dwr_output = f"{BASE_DIR}\unmatched_dwr_output.xlsx"
18 supplemented_output_path = f"{BASE_DIR}\matched_output_supplemented.xlsx"
19 parcel_excel = f"{BASE_DIR}\Parcels.xlsx"
20
21
22
23
24
25 #Load input datasets (OCR, DWR, Monterey, and parcel data) and preserve original match_source
26 df1 = pd.read_excel(file1_path)
27 df2 = pd.read_excel(file2_path)
28 df3 = pd.read_excel(file_monterey)
29 parcels_df = pd.read_excel(parcel_excel)
30
31 if 'match_source' in df1.columns:
32     df1['ocr_original_match_source'] = df1['match_source']
33
34 # Define helper function to clean and normalize well IDs
35 def clean_id(raw):
36     if pd.isna(raw):
37         return ""
38     original = str(raw).strip()
39     cleaned = original.replace(" ", "").upper().strip()
40     # Convert 'c' to 'E' for IDs starting with 'c' or 'C'
41     if cleaned.lower().startswith('c'):
42         cleaned = 'E' + cleaned[1:]
43     if re.match(r"^\WCR\d{4}-\d{6}$", cleaned):
44         return cleaned
45     text = original.lower().replace(", ", "")
46     text = re.sub(r"no\.", "", text)
47     text = text.replace("/", " ")
48     text = re.sub(r"[\s\n]+", " ", text)
```

```

49     match = re.match(r"^\d{1,2}[\.-](\d{1,2})[\.-](\d{2,4})(\d{3,8})$", text)
50     if match:
51         year_str = match.group(3)
52         try:
53             year = int(year_str)
54             if year < 100:
55                 year += 1900 if year >= 35 else 2000
56             if 1935 <= year <= 2025:
57                 return match.group(4)
58         except ValueError:
59             pass
60     prefix_removed = re.sub(r"^\d{1,2}[\.-]\d{1,2}[\.-]\d{2,4}", "", text).strip()
61     segments = re.split(r"[s]+", prefix_removed)
62     for segment in reversed(segments):
63         segment_digits = re.findall(r"\d{3,8}", segment)
64         if segment_digits:
65             return segment_digits[-1].rstrip('0')
66     digits = re.findall(r"\d{3,8}", text)
67     if digits:
68         return digits[-1].rstrip('0')
69     # Return cleaned ID if it starts with 'E' and has digits
70     if cleaned.startswith('E') and any(c.isdigit() for c in cleaned[1:]):
71         return cleaned
72     return ""
73
74 #Format APN strings to xxx-yyy-zz pattern for 8-digit APNs, consistent with SCC Parcel db
75 def format_apn_with_dashes(apn_str):
76     apn_str = str(apn_str).strip()
77     clean_apn = re.sub(r'\D', '', apn_str)
78     if len(clean_apn) == 8:
79         return f"{clean_apn[:3]}-{clean_apn[3:6]}-{clean_apn[6:]}"
80     else:
81         return apn_str
82
83 # Normalize parcel APNs by removing non-digits and zero-padding to 8 digits
84 def normalize_parcel_apn(apn_str):
85     return re.sub(r"\D", "", str(apn_str)).zfill(8)
86
87 # Generate APN signature by keeping non-zero digits
88 def apn_signature(apn):
89     return ''.join([ch for ch in str(apn) if ch.isdigit() and ch != '0'])
90
91 #Clean address strings for consistent matching
92 def clean_for_matching(addr):
93     if not isinstance(addr, str):
94         addr = str(addr) if not pd.isna(addr) else ""
95     addr = addr.upper()
96     addr = addr.replace("(", "").replace(")", "")
97     addr = re.sub(r"\bARE\b", "AVE", addr)
98     addr = re.sub(r"\bSTH\b", "ST", addr)

```

```

99     addr = re.sub(r"\bDRI\b", "DR", addr)
100    addr = re.sub(r"\bMOUNT\b", "MT", addr)
101    addr = re.sub(r"\bR\b", "RD", addr)
102    addr = re.sub(r" STREET", " ST", addr)
103    addr = re.sub(r" AVENUE", " AVE", addr)
104    addr = re.sub(r" ROAD", " RD", addr)
105    addr = re.sub(r" DRIVE", " DR", addr)
106    addr = re.sub(r" LANE", " LN", addr)
107    addr = re.sub(r" COURT", " CT", addr)
108    addr = re.sub(r" PLACE", " PL", addr)
109    addr = re.sub(r" TERRACE", " TER", addr)
110    addr = re.sub(r" CIRCLE", " CIR", addr)
111    addr = re.sub(r" TRAIL", " TRL", addr)
112    addr = re.sub(r" CREEK", " CRK", addr)
113    addr = re.sub(r" WAY", " WAY", addr)
114    addr = re.sub(r"\s+[A-Z\s]*\sCA\s\d{5}$", "", addr)
115    addr = re.sub(r"^\w\s+$", "", addr)
116    addr = re.sub(r"\s+", " ", addr).strip()
117    addr = re.sub(r"\b\w\b$", "", addr).strip()
118    return addr
119
120 # Normalize APN strings by removing non-digits and handling specific cases
121 def normalize_apn(apn_str):
122     digits = re.sub(r"\D", "", str(apn_str))
123     if len(digits) == 9 and digits[-3] == '0':
124         digits = digits[:-3] + digits[-2:]
125     return digits.zfill(8)
126
127 # Generate APN signature for matching by keeping non-zero digits
128 def get_apn_signature(apn_str):
129     return "".join(d for d in str(apn_str) if d.isdigit() and d != "0")
130
131 # Extract APN from location text using regex
132 def extract_apn_from_location(location_str):
133     match = re.search(r'\bAPN\s*(\d{2,3}[- ]?\d{2,3}[- ]?\d{2,3})\b', str(location_str),
134 re.IGNORECASE)
135     if match:
136         return normalize_apn(match.group(1))
137     return None
138
139 #Normalize year values for comparisons
140 def normalize_year(v):
141     if pd.isna(v) or not str(v).strip():
142         return None
143     s = str(v).strip()
144     if '/' in s:
145         parts = s.split('/')
146         if len(parts) >= 3 and parts[2].isdigit():
147             y = int(parts[2])
148             if len(parts[2]) == 2:

```

```

148         y += 1900 if y >= 35 else 2000
149     return y
150     s = re.sub(r'\s', '', s)
151     match = re.search(r'(\d{2,4})$', s)
152     if match:
153         y = int(match.group(1))
154         if 0 <= y <= 99:
155             y += 1900 if y >= 35 else 2000
156         if 1935 <= y <= 2030:
157             return y
158     return None
159
160 # Find best matching DWR record based on well ID and field comparisons
161 def get_best_match(well_id, records, df1_row=None, fields_to_check=None, field_map=None,
162                   used_indices=None, target_df=None):
163     if used_indices is None:
164         used_indices = set()
165
166     if not well_id.strip():
167         return None, None, None, None, None, 0
168
169     well_id_clean = str(well_id).strip().rstrip('M').replace('.0', '')
170     if well_id_clean.startswith('E'):
171         well_id_normalized = 'E' + well_id_clean[1:].lstrip('0')
172     else:
173         well_id_normalized = well_id_clean.lstrip('0')
174
175     exact_matches = []
176     for idx, (target_id, wcr_number, apn, addr, depth, water) in enumerate(records):
177         if idx in used_indices:
178             continue
179         target_clean = str(target_id).strip().rstrip('M').replace('.0', '')
180         if target_clean.startswith('E'):
181             target_normalized = 'E' + target_clean[1:].lstrip('0')
182         else:
183             target_normalized = target_clean.lstrip('0')
184
185         if well_id_normalized == target_normalized:
186             exact_matches.append((target_id, wcr_number, apn, addr, depth, water, idx))
187
188     if exact_matches:
189         if len(exact_matches) == 1:
190             target_id, wcr_number, apn, addr, depth, water, idx = exact_matches[0]
191             logging.debug(f"Exact match: well_id_normalized='{well_id_normalized}',
192 target='{target_id}' (single match)")
193             return target_id, 100, apn, addr, idx, 0
194         else:
195             logging.debug(f"Multiple matches found for
196 well_id_normalized='{well_id_normalized}': {len(exact_matches)} candidates")
197             best_match = None
198             best_match_score = -1

```

```

196     best_match_idx = None
197     best_matched_fields = []
198
199     for match in exact_matches:
200         target_id, wcr_number, apn, addr, depth, water, idx = match
201         df2_row = df2.iloc[idx] if target_df is None else target_df.iloc[idx]
202         match_score = 0
203         matched_fields = []
204
205         for df1_field in ["date", "depth_well_cleaned"] + [f for f in fields_to_check
if f not in ["date", "depth_well_cleaned"]]:
206             df2_field = None
207             for key, (df1_col, _) in field_map.items():
208                 if df1_col == df1_field:
209                     df2_field = key
210                     break
211             if not df2_field:
212                 continue
213
214             df1_value = df1_row.get(df1_field) if df1_row is not None else np.nan
215             df2_value = df2_row.get(df2_field)
216
217             if isinstance(df1_value, str) and df1_value.strip() == '':
218                 df1_value = np.nan
219             if isinstance(df2_value, str) and df2_value.strip() == '':
220                 df2_value = np.nan
221
222             if df1_field == "date":
223                 y1 = normalize_year(df1_value)
224                 y2 = normalize_year(df2_value)
225                 if y1 == y2 and y1 is not None:
226                     match_score += 2
227                     matched_fields.append(df1_field)
228                     continue
229             elif df1_field == "depth_well_cleaned":
230                 if not pd.isna(df1_value) and not pd.isna(df2_value):
231                     if abs(df1_value - df2_value) <= 5:
232                         match_score += 2
233                         matched_fields.append(df1_field)
234                         continue
235             elif df1_field in ["water_static_cleaned", "total_depth_cleaned",
"yield_cleaned"]:
236                 if not pd.isna(df1_value) and not pd.isna(df2_value):
237                     if abs(df1_value - df2_value) <= 5:
238                         match_score += 1
239                         matched_fields.append(df1_field)
240                         continue
241             elif df1_field == "apn_cleaned":
242                 v1 = normalize_apn(df1_value) if df1_value else ''
243                 v2 = normalize_apn(df2_value) if df2_value else ''

```

```

244         if v1 == v2 and v1:
245             match_score += 1
246             matched_fields.append(df1_field)
247             continue
248
249         if pd.isna(df1_value) and pd.isna(df2_value):
250             continue
251         elif isinstance(df1_value, str) and isinstance(df2_value, str):
252             if df1_value.strip().upper() == df2_value.strip().upper():
253                 match_score += 1
254                 matched_fields.append(df1_field)
255         elif df1_value == df2_value:
256             match_score += 1
257             matched_fields.append(df1_field)
258
259         if match_score > best_match_score:
260             best_match = match
261             best_match_score = match_score
262             best_match_idx = idx
263             best_matched_fields = matched_fields
264
265         if best_match:
266             target_id, wcr_number, apn, addr, depth, water, idx = best_match
267             logging.debug(f"Selected match for well_id_normalized='{well_id_normalized}'
based on field matching (score={best_match_score}, matched_fields={best_matched_fields}")
268             return target_id, 100, apn, addr, idx, best_match_score
269
270         logging.debug(f"no match from field scoring for
well_id_normalized='{well_id_normalized}' - falling back to depth/water")
271         best_match = exact_matches[0]
272         best_depth_match = False
273         df1_depth = df1_row.get("depth_well_cleaned", np.nan) if df1_row is not None else
np.nan
274         df1_water = df1_row.get("water_static_cleaned", np.nan) if df1_row is not None
else np.nan
275         for match in exact_matches:
276             target_id, wcr_number, apn, addr, depth, water, idx = match
277             depth_match = (not pd.isna(df1_depth) and not pd.isna(depth) and
abs(df1_depth - depth) <= 5)
278             water_match = (not pd.isna(df1_water) and not pd.isna(water) and
abs(df1_water - water) <= 5)
279             if depth_match:
280                 logging.debug(f"Selected match for
well_id_normalized='{well_id_normalized}' due to depth match (df1_depth={df1_depth},
candidate_depth={depth}")
281                 return target_id, 100, apn, addr, idx, 0
282             elif water_match and not best_depth_match:
283                 best_match = match
284                 logging.debug(f"Selected match for
well_id_normalized='{well_id_normalized}' due to water match (df1_water={df1_water},
candidate_water={water}")

```

```

285         best_depth_match = False
286         target_id, wcr_number, apn, addr, depth, water, idx = best_match
287         logging.debug(f"Fallback match: well_id_normalized='{well_id_normalized}',
target='{target_id}' (index={idx})")
288         return target_id, 100, apn, addr, idx, 0
289
290     df1_depth = df1_row.get("depth_well_cleaned", np.nan) if df1_row is not None else np.nan
291     df1_water = df1_row.get("water_static_cleaned", np.nan) if df1_row is not None else
np.nan
292     for idx, (target_id, wcr_number, apn, addr, depth, water) in enumerate(records):
293         if idx in used_indices:
294             continue
295         target_clean = str(target_id).strip().rstrip('M').replace('.0', '')
296         if target_clean.startswith('E'):
297             target_normalized = 'E' + target_clean[1:].lstrip('0')
298         else:
299             target_normalized = target_clean.lstrip('0')
300
301         is_substring_match = False
302         if well_id_normalized in target_normalized and len(target_normalized) >
len(well_id_normalized):
303             is_substring_match = True
304             match_type = f"Substring match (well_id in target):
well_id_normalized='{well_id_normalized}', target='{target_clean}'"
305             elif target_normalized in well_id_normalized and len(well_id_normalized) >
len(target_normalized):
306                 is_substring_match = True
307                 match_type = f"Substring match (target in well_id):
well_id_normalized='{well_id_normalized}', target='{target_clean}'"
308
309             if is_substring_match:
310                 depth_match = (not pd.isna(df1_depth) and not pd.isna(depth) and abs(df1_depth -
depth) <= 5)
311                 water_match = (not pd.isna(df1_water) and not pd.isna(water) and abs(df1_water -
water) <= 5)
312                 if not (depth_match or water_match):
313                     logging.debug(f"Skipping substring match for
well_id_normalized='{well_id_normalized}', target='{target_clean}' due to no depth/water
match")
314                 continue
315
316                 logging.debug(match_type)
317                 return target_id, 80, apn, addr, idx, 0
318
319             logging.debug(f"No match found for well_id_normalized='{well_id_normalized}'")
320             return None, None, None, None, None, 0
321
322 # Perform fuzzy matching on addresses against parcel data
323 def fuzzy_address_lookup(addr):
324     cleaned_input = clean_for_matching(addr)
325     if not cleaned_input:

```

```

326         return None, None, 0
327
328     matches = process.extract(
329         cleaned_input,
330         parcels_df["parcel_address"].tolist(),
331         scorer=fuzz.token_sort_ratio,
332         limit=10
333     )
334
335     for match_text, score, idx in matches:
336         if score >= 85:
337             best_row = parcels_df.iloc[idx]
338             return best_row["parcel_apn"], best_row["parcel_address"], score
339
340     return None, None, 0
341
342 # Prepare DWR dataset (df2) for matching by cleaning IDs, APNs, and addresses
343 df2["clean_apn"] = df2.iloc[:, 27].apply(normalize_apn)
344 df2["clean_address"] = df2.iloc[:, 3].apply(clean_for_matching)
345 df2_depth = df2["Total Completed Depth"].astype(float).fillna(np.nan)
346 df2_water = df2["Static Water Level"].astype(float).fillna(np.nan)
347
348 df2_ids = []
349 for i in range(len(df2)):
350     id1_raw = df2.iloc[i, 1] # Legacy Log Number
351     id0_raw = df2.iloc[i, 0] # WCR Number
352     id1 = clean_id(id1_raw)
353     id0 = clean_id(id0_raw)
354     if id1 and re.match(r"^(WCR\d{4}-\d{6}|\d{3,8}|E\d+)$", id1):
355         df2_ids.append(id1)
356     else:
357         df2_ids.append(id0)
358 df2_ids = pd.Series(df2_ids, index=df2.index)
359
360 df2_apn = df2.iloc[:, 27].astype(str).fillna('')
361 df2_addr = df2.iloc[:, 3].astype(str).fillna('')
362 df2_wcr = df2.iloc[:, 0].astype(str).fillna('')
363 df2_records = list(zip(df2_ids, df2_wcr, df2_apn, df2_addr, df2_depth, df2_water))
364
365 # Prepare Monterey dataset (df3) for secondary matching
366 df3.columns = df3.columns.str.strip()
367 df3_ids = []
368 for i in range(len(df3)):
369     id1_raw = df3.iloc[i, 1] # Legacy Log Number
370     id0_raw = df3.iloc[i, 0] # WCR Number
371     id1 = clean_id(id1_raw)
372     id0 = clean_id(id0_raw)
373     if id1 and re.match(r"^(WCR\d{4}-\d{6}|\d{3,8}|E\d+)$", id1):
374         df3_ids.append(id1)
375     else:

```

```
376     df3_ids.append(id0)
377 df3_ids = pd.Series(df3_ids, index=df3.index)
378 df3_wcr = df3.iloc[:, 0].astype(str).fillna('')
379 df3_apn = df3.iloc[:, 27].astype(str).fillna('')
380 df3_addr = df3.iloc[:, 3].astype(str).fillna('')
381 df3_depth = df3["Total Completed Depth"].astype(float).fillna(np.nan)
382 df3_water = df3["Static Water Level"].astype(float).fillna(np.nan)
383 df3_records = list(zip(df3_ids, df3_wcr, df3_apn, df3_addr, df3_depth, df3_water))
384
385 # Prepare parcel data for APN and address lookups
386 parcels_df["parcel_address"] = parcels_df.iloc[:, 13].astype(str).apply(clean_for_matching)
387 parcels_df["parcel_apn"] = parcels_df.iloc[:, 1].apply(normalize_parcel_apn)
388 parcel_apn_to_address = dict(zip(parcels_df["parcel_apn"], parcels_df["parcel_address"]))
389 signature_to_apns = defaultdict(list)
390 for apn in parcels_df["parcel_apn"]:
391     sig = get_apn_signature(apn)
392     signature_to_apns[sig].append(apn)
393 apn_lookup = parcels_df.set_index("parcel_apn")["parcel_address"].to_dict()
394
395 # Set up logging for debugging and tracking matches
396 logging.basicConfig(
397     filename=r"C:\Well Layer Update\Python\Azure\Results\wc11\excel 3\fuzzy_match_debug.log",
398     filemode="a",
399     format="%(asctime)s [%(levelname)s] %(message)s",
400     level=logging.DEBUG
401 )
402
403 # Log sample data for debugging
404 logging.debug(f"df1 well_id_cleaned sample: {df1['well_id_cleaned'].head(10).tolist()}")
405 logging.debug(f"df2_ids sample: {df2_ids.head(10).tolist()}")
406
407
408 # Define fields for validating duplicate matches
409 fields_to_check = [
410     "total_depth_cleaned",
411     "depth_well_cleaned",
412     "water_static_cleaned",
413     "yield_cleaned",
414     "permit_date",
415     "permit_no",
416     "date",
417     "apn_cleaned"
418 ]
419
420 #Define mapping of DWR fields to OCR fields and their source columns
421 field_map = {
422     "Owner Assigned Well Number": ("owner_well_id", "owner_assigned_source"),
423     "Date Work Ended": ("date", "date_source"),
424     "Permit Number": ("permit_no", "permit_number_source"),
425     "Driller Name": ("contractor", "driller_name_source"),
```

```

426     "Total Drill Depth": ("total_depth_cleaned", "total_drill_depth_source"),
427     "Total Completed Depth": ("depth_well_cleaned", "total_completed_depth_source"),
428     "Static Water Level": ("water_static_cleaned", "static_water_level_source"),
429     "Well Yield": ("yield_cleaned", "well_yield_source"),
430     "B118WellUse": ("planned_use_former_use", "b118welluse_source"),
431     "Record Type": ("Type of Work", "record_type_source"),
432     "Permit Date": ("permit_date", "permit_date_source"),
433     "Driller License Number": ("contractor_license", "driller_license_number_source"),
434     "Drilling Method": ("drilling_method", "drilling_method_source"),
435     "Fluid": ("fluid", "fluid_source"),
436     "Total Draw Down": ("total_drawdown", "total_drawdown_source"),
437     "Test Type": ("test_type", "test_type_source"),
438     "Pump Test Length": ("pump_test_length_cleaned", "pump_test_length_source"),
439     "Other Observations": ("other_observations", "other_observations_source"),
440     "City": ("city", "city_source"),
441     "Ground Surface Elevation": ("surface_elevation", "surface_elevation_source"),
442     "Elevation Determination Method": ("elevation_method", "elevation_method_source"),
443     "Top Of Perforated Interval": ("top_perforations", "top_perforations_source"),
444     "Bottom of Perforated Interval": ("bottom_perforations", "bottom_perforations_source"),
445     "Decimal Latitude": ("lat_cleaned", "latitude_source"),
446     "Decimal Longitude": ("long_cleaned", "longitude_source"),
447     "APN": ("apn_cleaned", "apn_source")
448 }
449
450 #Define additional OCR fields to carry over to supplemented output
451 additional_fields = {
452     "seal_cleaned": ("seal_cleaned", "seal_cleaned_source"),
453     "first_water_cleaned": ("first_water_cleaned", "first_water_cleaned_source"),
454     "template_id": ("template_id", "template_id_source"),
455     "OCR_match_source": ("ocr_original_match_source", "OCR_match_source_source")
456 }
457
458 #Initialize set to track used DWR indices for matching
459 used_dwr_indices = set()
460
461 # Match OCR records to DWR records and update APN/address
462 def match_apn_and_address(row):
463     well_id = clean_id(row["well_id_cleaned"]).strip()
464     if not well_id:
465         return row["final_address"], str(row["final_apn"]).zfill(7), "No match", "No valid
well ID", 0
466
467     well_id = str(well_id).strip().rstrip('0').replace('.0', '')
468     logging.debug(f"Processing well_id: {well_id}")
469
470     match_id, score, apn_dwr, addr_dwr, dwr_idx, match_score = get_best_match(
471         well_id, df2_records, df1_row=row, fields_to_check=fields_to_check,
field_map=field_map, used_indices=used_dwr_indices, target_df=df2
472     )
473     if match_id and dwr_idx is not None:

```

```

474     dwr_row = df2.iloc[dwr_idx]
475     wcr_number = dwr_row["WCR Number"]
476     dwr_apn = dwr_row["clean_apn"]
477     dwr_address = dwr_row["clean_address"]
478     logging.debug(f"Match found for well_id='{well_id}' to df2_id='{match_id}' (score=
{score}, duplicate_match_score={match_score})")
479     if not dwr_apn or dwr_apn == '00000000':
480         potential_apn = extract_apn_from_location(dwr_row["Well Location"])
481         if potential_apn:
482             dwr_apn = potential_apn
483             logging.debug(f"Extracted APN '{dwr_apn}' from Well Location")
484
485     if dwr_apn in parcel_apn_to_address:
486         return parcel_apn_to_address[dwr_apn], dwr_apn, f"ID match (score {score}) → APN
match", f"OCR ID '{well_id}' → DWR ID '{match_id}' | WCR Number: {wcr_number}", match_score
487
488     matched_apn, matched_address, addr_score = fuzzy_address_lookup(dwr_address)
489     if matched_apn and matched_address:
490         return (
491             matched_address,
492             matched_apn,
493             f"ID match (score {score}) → fuzzy address",
494             f"OCR ID '{well_id}' → DWR ID '{match_id}' → Address (score: {addr_score}) |
WCR Number: {wcr_number}",
495             match_score
496         )
497
498     logging.debug(f"ID match for '{well_id}' to '{match_id}' (score: {score}) but no
parcel match")
499     return row["final_address"], str(row["final_apn"]).zfill(7), f"ID match (score
{score}) with no parcel match", f"OCR ID '{well_id}' → DWR ID '{match_id}' | WCR Number:
{wcr_number}", match_score
500
501     location_details_raw = row.get("location_details", "")
502     location_details_clean = clean_for_matching(location_details_raw)
503
504     if location_details_clean:
505         dwr_locations_raw = df2["Well Location"].fillna("").astype(str)
506         dwr_locations_cleaned = dwr_locations_raw.apply(clean_for_matching).tolist()
507
508         match_result = process.extractOne(location_details_clean, dwr_locations_cleaned,
scorer=fuzz.token_sort_ratio)
509
510         if match_result:
511             best_match_text, score, match_index = match_result
512             matched_row = df2.iloc[match_index]
513             original_dwr_location = dwr_locations_raw.iloc[match_index]
514
515             if score >= 80:
516                 legacy_log = str(matched_row.get("Legacy Log Number", "")).strip().upper()
517                 if legacy_log != "NN":

```

```

518         ocr_apn = str(row["final_apn"]).zfill(7)
519         ocr_match_source = str(row.get("match_source", ""))
520         if ocr_match_source == "No match" or ocr_apn not in parcel_apn_to_addresses:
521             ocr_apn = "0000000"
522             return row["final_address"], ocr_apn, "Fallback to OCR values", "Fuzzy
match blocked: Legacy Log Number not 'NN'", 0
523
524         dwr_apn = matched_row["clean_apn"]
525         dwr_address = matched_row["clean_address"]
526
527         if not dwr_apn or dwr_apn == '0000000':
528             potential_apn = extract_apn_from_location(matched_row["Well Location"])
529             if potential_apn:
530                 dwr_apn = potential_apn
531
532         if dwr_apn in parcel_apn_to_address:
533             wcr_number = matched_row["WCR Number"]
534             return parcel_apn_to_address[dwr_apn], dwr_apn, "Fuzzy location → APN
match", f"Location '{location_details_raw}' → '{original_dwr_location}' | WCR Number:
{wcr_number}", 0
535
536         matched_apn, matched_address, addr_score = fuzzy_address_lookup(dwr_address)
537         if matched_apn and matched_address:
538             wcr_number = matched_row["WCR Number"]
539             return (
540                 matched_address,
541                 matched_apn,
542                 "Fuzzy location → fuzzy address",
543                 f"Location '{location_details_raw}' → '{original_dwr_location}' →
Address: {dwr_address} (score: {score}) → Parcel (score: {addr_score}) | WCR Number:
{wcr_number}",
544                 0
545             )
546
547         wcr_number = matched_row["WCR Number"]
548         return row["final_address"], str(row["final_apn"]).zfill(7), "Fuzzy location
match with no parcel match", f"Location '{location_details_raw}' → '{original_dwr_location}'
(score: {score}) | WCR Number: {wcr_number}", 0
549
550         ocr_apn = str(row["final_apn"]).zfill(7)
551         ocr_match_source = str(row.get("match_source", ""))
552         if ocr_match_source == "No match":
553             ocr_apn = "0000000"
554         return row["final_address"], ocr_apn, "Fallback to OCR values", "No match found", 0
555
556 # Apply matching logic to update APN and address fields in OCR dataset
557 df1_ids = df1["well_id_cleaned"].apply(clean_id)
558 df1_apn = df1["final_apn"].astype(str).fillna('')
559 df1_addr = df1["final_address"].astype(str).fillna('')
560

```

```

561 df1[["final_address", "final_apn", "match_source", "match_log", "duplicate_match_score"]] =
df1.apply(
562     lambda row: pd.Series(match_apn_and_address(row)),
563     axis=1
564 )
565
566 # Process matched and unmatched records, including Monterey dataset checks
567 matched_records = []
568 matched_dwr_ids = set()
569 unmatched_ocr_records = []
570
571 for i in range(len(df1)):
572     source_id = df1_ids.iloc[i]
573     match_source = df1.at[i, "match_source"]
574
575     if match_source == "Fallback to OCR values":
576         matched_in_monterey = "No"
577         match_id_mon, score_mon, apn_mon, addr_mon, mon_idx, _ = get_best_match(source_id,
df3_records, df1_row=df1.iloc[i], fields_to_check=fields_to_check, field_map=field_map,
target_df=df3)
578         print(f"NO MATCH: df1_id raw='{df1.at[i, 'well_id']}', cleaned='{source_id}'")
579         if match_id_mon:
580             matched_in_monterey = "Yes"
581             unmatched_record = df1.iloc[i].copy()
582             unmatched_record["Matched_in_Monterey"] = matched_in_monterey
583             unmatched_ocr_records.append(unmatched_record)
584
585         else:
586             df1_depth = df1.at[i, "depth_well_cleaned"] if "depth_well_cleaned" in df1.columns
else np.nan
587             df1_water = df1.at[i, "water_static_cleaned"] if "water_static_cleaned" in
df1.columns else np.nan
588             match_id, score, apn_dwr, addr_dwr, dwr_index, match_score =
get_best_match(source_id, df2_records, df1_row=df1.iloc[i], fields_to_check=fields_to_check,
field_map=field_map, used_indices=used_dwr_indices, target_df=df2)
589             if match_id and dwr_index is not None and dwr_index not in used_dwr_indices:
590                 logging.debug(f"[Row {i}] Match found for df1_id='{source_id}' to
df2_id='{match_id}' (score={score}, duplicate_match_score={match_score})")
591                 record = df2.iloc[dwr_index].copy()
592                 record["final_address"] = df1.at[i, "final_address"]
593                 record["final_apn"] = df1.at[i, "final_apn"]
594                 record["match_source"] = match_source if match_source else f"ID match (score
{score})"
595                 record["matched_df2_well_id"] = match_id
596                 record["matched_df1_well_id"] = source_id
597                 record["duplicate_match_score"] = match_score
598                 record["df1_row_index"] = i
599                 matched_records.append(record)
600                 matched_dwr_ids.add(match_id)
601                 used_dwr_indices.add(dwr_index)
602             else:

```

```

603     logging.debug(f"[Row {i}] No df2 match for df1_id='{source_id}' or df2 index
{dwr_index} already used")
604     unmatched_record = df1.iloc[i].copy()
605     unmatched_record["Matched_in_Monterey"] = "No"
606     unmatched_ocr_records.append(unmatched_record)
607
608     if match_source in [
609         "Fuzzy location → APN match",
610         "Fuzzy location → fuzzy address",
611         "Fuzzy location match with no parcel match"
612     ]:
613         location_raw = df1.at[i, "location_details"]
614         location_clean = clean_for_matching(location_raw)
615         dwr_locations_raw = df2["Well Location"].fillna("").astype(str)
616         dwr_locations_clean = dwr_locations_raw.apply(clean_for_matching).tolist()
617
618         match_result = process.extractOne(
619             location_clean,
620             dwr_locations_clean,
621             scorer=fuzz.token_sort_ratio
622         )
623
624         if match_result:
625             match_value, score, match_index = match_result
626             if match_index not in used_dwr_indices:
627                 logging.debug(f"[Row {i}] Best fuzzy match='{match_value}' with score=
{score}")
628                 matched_row = df2.iloc[match_index].copy()
629                 matched_row["final_address"] = df1.at[i, "final_address"]
630                 matched_row["final_apn"] = df1.at[i, "final_apn"]
631                 matched_row["match_source"] = match_source
632                 matched_row["matched_df2_well_id"] = df2_ids.iloc[match_index]
633                 matched_row["matched_df1_well_id"] = source_id
634                 matched_row["duplicate_match_score"] = 0
635                 matched_row["df1_row_index"] = i
636                 matched_records.append(matched_row)
637                 matched_dwr_ids.add(df2_ids.iloc[match_index])
638                 used_dwr_indices.add(match_index)
639             else:
640                 logging.debug(f"[Row {i}] Fuzzy match index {match_index} already used")
641                 unmatched_record = df1.iloc[i].copy()
642                 unmatched_record["Matched_in_Monterey"] = "No"
643                 unmatched_ocr_records.append(unmatched_record)
644         else:
645             logging.warning(f"[Row {i}] Fuzzy match expected but no location match found for
location='{location_clean}'")
646             unmatched_record = df1.iloc[i].copy()
647             unmatched_record["Matched_in_Monterey"] = "No"
648             unmatched_ocr_records.append(unmatched_record)
649
650 # Generate unmatched DWR records (those not matched to OCR)

```

```
651 unmatched_dwr_records = df2[~df2_ids.isin(matched_dwr_ids)]
652
653 # Save matched and unmatched records to Excel files
654 if matched_records:
655     matched_df = pd.DataFrame(matched_records)
656     matched_df["final_apn"] = matched_df["final_apn"].fillna("").astype(str)
657     mask = (matched_df["match_source"] == "No match") & (
658         matched_df["final_address"].isna() | (matched_df["final_address"].str.strip() == "")
659     )
660     matched_df.loc[mask, "final_apn"] = "0000000"
661
662     fallback_lookup = df1.set_index("well_id_cleaned")["match_source"].to_dict()
663     matched_df["fallback_match_source"] = ""
664
665     for idx, row in matched_df.iterrows():
666         if row.get("match_source") == "ID match with no parcel match":
667             well_id = row.get("matched_df1_well_id")
668             if well_id in fallback_lookup:
669                 matched_df.at[idx, "fallback_match_source"] = fallback_lookup[well_id]
670
671     matched_df.to_excel(matched_output, index=False)
672     print(f"Matched DWR records written to: {matched_output}")
673
674 if unmatched_ocr_records:
675     unmatched_ocr_df = pd.DataFrame(unmatched_ocr_records)
676     unmatched_ocr_df.to_excel(unmatched_ocr_output, index=False)
677     print(f"Unmatched OCR records written to: {unmatched_ocr_output}")
678
679 if not unmatched_dwr_records.empty:
680     unmatched_dwr_records.to_excel(unmatched_dwr_output, index=False)
681     print(f"Unmatched DWR records written to: {unmatched_dwr_output}")
682
683 # Format APN fields in matched and unmatched OCR datasets
684 matched_df = pd.DataFrame(matched_records)
685 unmatched_ocr_df = pd.DataFrame(unmatched_ocr_records)
686
687 if not matched_df.empty:
688     matched_df["final_apn"] = matched_df["final_apn"].apply(format_apn_with_dashes)
689
690 if not unmatched_ocr_df.empty:
691     unmatched_ocr_df["final_apn"] = unmatched_ocr_df["final_apn"].apply(format_apn_with_dashes)
692
693 # Supplement matched dataset with OCR fields and source tracking
694 supplemented_df = matched_df.copy().reset_index(drop=True)
695
696 supplemented_df["duplicate_exists"] = ""
697 for _, source_col in list(field_map.values()) + list(additional_fields.values()):
698     supplemented_df[source_col] = ""
699 for final_col, _ in additional_fields.items():
```

```

700     supplemented_df[final_col] = np.nan
701
702 # create column for the OCR depths flag
703 supplemented_df["depths_flag"] = ""
704
705 columns_to_supplement = [
706     (final_col, ocr_col, source_col)
707     for final_col, (ocr_col, source_col) in field_map.items()
708 ]
709 additional_columns = [
710     (final_col, ocr_col, source_col)
711     for final_col, (ocr_col, source_col) in additional_fields.items()
712 ]
713
714 df1["well_id_cleaned"] = df1["well_id_cleaned"].apply(lambda x: str(x).strip().rstrip(".0")
if not pd.isna(x) else "")
715
716 for idx, row in supplemented_df.iterrows():
717     df1_idx = row.get("df1_row_index")
718     if pd.isna(df1_idx):
719         continue
720     ocr_row = df1.iloc[int(df1_idx)]
721     supplemented_df.at[idx, "duplicate_exists"] = ocr_row.get("duplicate_exists", "")
722     for match_col, ocr_col, source_col in columns_to_supplement:
723         original_val = supplemented_df.at[idx, match_col]
724         ocr_val = ocr_row.get(ocr_col, np.nan)
725         if match_col == "B118WellUse":
726             dwr_val_str = str(original_val).strip().upper() if isinstance(original_val, str)
else ""
727             ocr_val_str = str(ocr_val).strip()
728             if dwr_val_str == "UNKNOWN" and ocr_val_str.upper() not in {"", "NAN"}:
729                 supplemented_df.at[idx, match_col] = ocr_val
730                 supplemented_df.at[idx, source_col] = "AI"
731             continue
732         if pd.isna(original_val) or (isinstance(original_val, str) and original_val.strip()
== ""):
733             if not pd.isna(ocr_val) and str(ocr_val).strip() != "":
734                 supplemented_df.at[idx, match_col] = ocr_val
735                 supplemented_df.at[idx, source_col] = "AI"
736             else:
737                 supplemented_df.at[idx, source_col] = "DWR"
738         else:
739             supplemented_df.at[idx, source_col] = "DWR"
740     for match_col, ocr_col, source_col in additional_columns:
741         ocr_val = ocr_row.get(ocr_col, np.nan)
742         if not pd.isna(ocr_val) and str(ocr_val).strip() != "":
743             supplemented_df.at[idx, match_col] = ocr_val
744             supplemented_df.at[idx, source_col] = "AI"
745         else:
746             supplemented_df.at[idx, match_col] = np.nan
747             supplemented_df.at[idx, source_col] = "AI"

```

```
748
749
750 # Pull the two source columns that belong to the depth fields
751 completed_src = supplemented_df.at[idx, "total_completed_depth_source"]
752 drill_src     = supplemented_df.at[idx, "total_drill_depth_source"]
753
754 # If either source contains "AI" then copy OCR flag
755 if (pd.notna(completed_src) and "AI" in str(completed_src).upper()) or \
756     (pd.notna(drill_src)     and "AI" in str(drill_src).upper()):
757     ocr_flag = ocr_row.get("depths_flag", "")
758     supplemented_df.at[idx, "depths_flag"] = ocr_flag
759 else:
760     supplemented_df.at[idx, "depths_flag"] = "" # blank when both are DWR (or empty)
761
762 # Flag record type mismatches between DWR and OCR datasets
763 for idx, row in supplemented_df.iterrows():
764     df1_idx = row.get("df1_row_index")
765     if pd.isna(df1_idx):
766         continue
767     ocr_row = df1.iloc[int(df1_idx)]
768     dwr_record_type = str(row.get("Record Type", "")).strip().upper()
769     ocr_record_type = str(ocr_row.get("Record Type", "")).strip().upper()
770     if ocr_record_type in {"", "UNKNOWN"}:
771         continue
772     if dwr_record_type != ocr_record_type:
773         supplemented_df.at[idx, "Record_Type_Mismatch"] = f"DWR: {dwr_record_type} ≠ OCR:
{ocr_record_type}"
774
775 #Add fallback match source for records with no parcel match
776 supplemented_df["fallback_match_source"] = ""
777 for idx, row in supplemented_df.iterrows():
778     if row.get("match_source") == "ID match with no parcel match":
779         df1_idx = row.get("df1_row_index")
780         if not pd.isna(df1_idx):
781             fallback_val = df1.iloc[int(df1_idx)].get("match_source", "")
782             supplemented_df.at[idx, "fallback_match_source"] = fallback_val
783
784 # Save supplemented dataset to Excel
785 supplemented_df.to_excel(supplemented_output_path, index=False)
786 print("supplemented file written to:", supplemented_output_path)
```

## C:\Well Layer Update\Python\Azure\County Database Matching Script.py

```

1 import pandas as pd
2 import re
3 import logging
4
5
6
7 BASE_DIR = r"<internal_project_directory>"
8 # df4_path = f"{BASE_DIR}\C1.xlsx"
9 # df5_path = f"{BASE_DIR}\OCR(SCC1)_OCR(DWR1)2.xlsx"
10 # matched_output = f"{BASE_DIR}\SCC1_DWR1_C1_2.xlsx"
11 # unmatched_county_output = f"{BASE_DIR}\C2.xlsx"
12 # unmatched_source_output = f"{BASE_DIR}\OCR(SCC1)_OCR(DWR1)_Final.xlsx"
13
14 #c3 and dwr2
15 df4_path = f"{BASE_DIR}\C3.xlsx"
16 df5_path = f"{BASE_DIR}\DWR2.xlsx"
17 matched_output = f"{BASE_DIR}\C3_DWR2.xlsx"
18 unmatched_county_output = f"{BASE_DIR}\C3.xlsx"
19 unmatched_source_output = f"{BASE_DIR}\C3_DWR2\DWR3.xlsx"
20
21 df4 = pd.read_excel(df4_path)
22 df5 = pd.read_excel(df5_path)
23
24
25 def clean_id(raw):
26     if pd.isna(raw):
27         return ""
28     original = str(raw).strip()
29     # Replace '0' with 'Ø' and 'e' with 'E'
30     cleaned = original.replace("0", "Ø").replace("o", "O").replace("e", "E").upper().strip()
31     if cleaned.lower().startswith('c'):
32         cleaned = 'E' + cleaned[1:]
33     if re.match(r"^\WCR\d{4}-\d{6}$", cleaned):
34         return cleaned
35     text = original.lower().replace("0", "Ø").replace("o", "O").replace("e",
36 "E").replace(",", "").replace("no.", "").replace("/", " ")
37     text = re.sub(r"[\s\n]+", " ", text)
38     match = re.match(r"^\d{1,2}[\.-](\d{1,2})[\.-](\d{2,4})(\d{3,8})$", text)
39     if match:
40         year_str = match.group(3)
41         try:
42             year = int(year_str)
43             if year < 100:
44                 year += 1900 if year >= 35 else 2000
45             if 1935 <= year <= 2025:
46                 return match.group(4)
47         except ValueError:
48             pass

```

```

48     prefix_removed = re.sub(r"^\d{1,2}[\.-]\d{1,2}[\.-]\d{2,4}", "", text).strip()
49     segments = re.split(r"[\s]+", prefix_removed)
50     for segment in reversed(segments):
51         segment_digits = re.findall(r"\d{3,8}", segment)
52         if segment_digits:
53             return segment_digits[-1].lstrip('0')
54     digits = re.findall(r"\d{3,8}", text)
55     if digits:
56         return digits[-1].lstrip('0')
57     if cleaned.startswith('E') and any(c.isdigit() for c in cleaned[1:]):
58         return cleaned
59     return ""
60
61
62 logging.basicConfig(
63     filename=r"C:\Well Layer Update\GIS\master files\county_match_debug.log",
64     filemode="a",
65     format="%(asctime)s %(levelname)s %(message)s",
66     level=logging.DEBUG
67 )
68
69
70 df4_ids = df4["Log_Number"].apply(clean_id)
71 df5_legacy_ids = df5["Legacy Log Number"].apply(clean_id) # Apply clean_id for consistency
72 df5_wcr_ids = df5["WCR Number"].apply(clean_id) if "WCR Number" in df5.columns else
pd.Series([""] * len(df5))
73
74
75 matched_records = []
76 unmatched_county_records = []
77 unmatched_source_records = []
78 used_df5_indices = set()
79
80 for i, county_id in enumerate(df4_ids):
81     if not county_id.strip():
82         continue # Skip empty Log_Number
83     county_id_clean = str(county_id).strip().rstrip('M').replace('.0', '')
84     county_id_normalized = 'E' + county_id_clean[1:].lstrip('0') if
county_id_clean.startswith('E') else county_id_clean.lstrip('0')
85
86     # Debug specific case
87     if county_id_clean.lower() == 'e0139534' or county_id_clean == 'E0139534':
88         logging.debug(f"Debugging e0139534: county_id_clean='{county_id_clean}',
county_id_normalized='{county_id_normalized}'")
89
90     match_found = False
91     matched_field = ""
92     df5_row = None
93     df4_row = df4.iloc[i]
94
95     # try matching with Legacy Log Number

```

```

96     for j, legacy_id in enumerate(df5_legacy_ids):
97         if j in used_df5_indices:
98             continue
99         legacy_id_clean = str(legacy_id).strip().rstrip('M').replace('.0', '')
100        legacy_id_normalized = 'E' + legacy_id_clean[1:].lstrip('0') if
legacy_id_clean.startswith('E') else legacy_id_clean.lstrip('0')
101
102
103        if county_id_normalized == legacy_id_normalized:
104            logging.debug(f"Match found: county_id='{county_id}' to
source_legacy_id='{legacy_id}'")
105            df5_row = df5.iloc[j].copy()
106            matched_field = "Legacy Log Number"
107            used_df5_indices.add(j)
108            match_found = True
109            break
110
111        #if no match on Legacy Log Number, try WCR Number
112        if not match_found and "WCR Number" in df5.columns:
113            for j, wcr_id in enumerate(df5_wcr_ids):
114                if j in used_df5_indices:
115                    continue
116                wcr_id_clean = str(wcr_id).strip().rstrip('M').replace('.0', '')
117                wcr_id_normalized = 'E' + wcr_id_clean[1:].lstrip('0') if
wcr_id_clean.startswith('E') else wcr_id_clean.lstrip('0')
118
119                if county_id_normalized == wcr_id_normalized:
120                    logging.debug(f"Match found: county_id='{county_id}' to
source_wcr_id='{wcr_id}'")
121                    df5_row = df5.iloc[j].copy()
122                    matched_field = "WCR Number"
123                    used_df5_indices.add(j)
124                    match_found = True
125                    break
126
127        if match_found:
128            # Append all df4 columns with 'county_' prefix
129            df4_row_prefixed = df4_row.rename(lambda x: f"county_{x}")
130            # Create record with df5 columns, df4 columns, and match metadata
131            record = pd.concat([df5_row, df4_row_prefixed])
132            record["matched_field"] = matched_field
133            record["cleaned_county_id"] = county_id
134            record["source_id"] = legacy_id if matched_field == "Legacy Log Number" else wcr_id
135            matched_records.append(record)
136        else:
137            unmatched_county_records.append(df4_row)
138
139
140    for j in range(len(df5)):
141        if j not in used_df5_indices:
142            df5_row = df5.iloc[j].copy()

```

```
143     # Collect only unmatched df5 records
144     unmatched_source_records.append(df5_row)
145     logging.debug(f"Unmatched df5 record added: index={j},
source_id='{df5_row.get('Legacy Log Number', 'N/A') if df5_row.get('Legacy Log Number') else
df5_row.get('WCR Number', 'N/A')}'")
146
147 # save matched
148 if matched_records:
149     matched_df = pd.DataFrame(matched_records)
150     matched_df.to_excel(matched_output, index=False)
151     print(f"Matched records written to: {matched_output}")
152 else:
153     print("No matched records to write to matched output.")
154
155 # save unmatched
156 if unmatched_county_records:
157     unmatched_county_df = pd.DataFrame(unmatched_county_records)
158
159     if 'IFILE' not in unmatched_county_df.columns:
160         unmatched_county_df['IFILE'] = pd.NA
161     unmatched_county_df.to_excel(unmatched_county_output, index=False)
162     print(f"Unmatched county records written to: {unmatched_county_output}")
163 else:
164     print("No unmatched county records with non-empty Log_Number.")
165
166 # save sources
167 if unmatched_source_records:
168     unmatched_source_df = pd.DataFrame(unmatched_source_records)
169     unmatched_source_df.to_excel(unmatched_source_output, index=False)
170     print(f"Unmatched source records written to: {unmatched_source_output}")
171 else:
172     print("No unmatched source records.")
```

## C:\Well Layer Update\Python\Azure\Field Highlighting Script.py

```
1 from openpyxl import load_workbook
2 from openpyxl.styles import Font
3
4 BASE_DIR = r"<internal_project_directory>"
5 supplemented_path = f"{BASE_DIR}\SCC_WELLS_FINAL_table.xlsx"
6
7 # Updated field_map to include additional fields
8 # field_map = {
9 #     "Owner_Assigned_Well_Number": "owner_assigned_source",
10 #     "Permit Number": "permit_number_source",
11 #     "Driller Name": "driller_name_source",
12 #     "Total Drill Depth": "total_drill_depth_source",
13 #     "Total Completed Depth": "total_completed_depth_source",
14 #     "Static Water Level": "static_water_level_source",
15 #     "Well Yield": "well_yield_source",
16 #     "B118WellUse": "b118welluse_source",
17 #     "Record Type": "record_type_source",
18 #     "Permit Date": "permit_date_source",
19 #     "Driller License Number": "driller_license_number_source",
20 #     "Drilling Method": "drilling_method_source",
21 #     "Fluid": "fluid_source",
22 #     "Total Draw Down": "total_drawdown_source",
23 #     "Test Type": "test_type_source",
24 #     "Pump Test Length": "pump_test_length_source",
25 #     "Other Observations": "other_observations_source",
26 #     "City": "city_source",
27 #     "Ground Surface Elevation": "surface_elevation_source",
28 #     "Elevation Determination Method": "elevation_method_source",
29 #     "Top Of Perforated Interval": "top_perforations_source",
30 #     "Bottom of Perforated Interval": "bottom_perforations_source",
31 #     "Decimal Latitude": "latitude_source",
32 #     "Decimal Longitude": "longitude_source",
33 #     "APN": "apn_source",
34 #     "seal_cleaned": "seal_cleaned_source",
35 #     "first_water_cleaned": "first_water_cleaned_source",
36 #     "template_id": "template_id_source",
37 #     "OCR_match_source": "OCR_match_source_source"
38 # }
39
40 #after gis export
41 field_map = {
42     "Owner_Assigned_Well_Number": "owner_assigned_source",
43     "Permit_Number": "permit_number_source",
44     "Date_Work_Ended": "date_source",
45     "Driller_Name": "driller_name_source",
46     "Total_Drill_Depth": "total_drill_depth_source",
47     "Total_Completed_Depth": "total_completed_depth_source",
48     "Static_Water_Level": "static_water_level_source",
```

```
49     "Well_Yield": "well_yield_source",
50     "B118WellUse": "b118welluse_source",
51     "Well_Type": "b118welluse_source",
52     "Record_Type": "record_type_source",
53     "Permit_Date": "permit_date_source",
54     "Driller_License_Number": "driller_license_number_source",
55     "Drilling_Method": "drilling_method_source",
56     "Fluid": "fluid_source",
57     "Total_Draw_Down": "total_drawdown_source",
58     "Test_Type": "test_type_source",
59     "Pump_Test_Length": "pump_test_length_source",
60     "City": "city_source",
61     "Ground_Surface_Elevation": "surface_elevation_source",
62     "Top_Of_Perforated_Interval": "top_perforations_source",
63     "Bottom_of_Perforated_Interval": "bottom_perforations_source",
64     "APN_mapped": "apn_source",
65     "Seal_Depth": "seal_cleaned_source",
66     "First_Water": "first_water_cleaned_source"
67 }
68
69
70 # Load workbook and sheet
71 wb = load_workbook(supplemented_path)
72 ws = wb.active
73
74
75 headers = {cell.value: idx for idx, cell in enumerate(ws[1], 1)} # 1-based indexing
76
77 red_font = Font(color="FF0000")
78
79 # Get indexes for special fields
80 final_apn_idx = headers.get("final_apn")
81 final_address_idx = headers.get("final_address")
82 match_source_idx = headers.get("match_source")
83
84 # Loop through rows starting from row 2 (excluding header)
85 for row in ws.iter_rows(min_row=2, max_row=ws.max_row):
86     # Highlight special logic for final_apn and final_address
87     match_source_val = row[match_source_idx - 1].value if match_source_idx else None
88
89     if final_apn_idx:
90         apn_cell = row[final_apn_idx - 1]
91         if match_source_val not in ["DWR APN matched to parcel database", "Fuzzy DWR address
92 match", "No match"]:
93             apn_cell.font = red_font # Mark as red if not in accepted match sources
94
95     if final_address_idx:
96         addr_cell = row[final_address_idx - 1]
97         if match_source_val not in ["DWR APN matched to parcel database", "Fuzzy DWR address
98 match"]:
```

```
97         addr_cell.font = red_font
98
99     #Apply red font to "AI"-sourced fields
100     for field, source_col in field_map.items():
101         col_idx = headers.get(field)
102         source_idx = headers.get(source_col)
103
104         if not col_idx or not source_idx:
105             continue
106
107         source_val = row[source_idx - 1].value
108         if source_val == "AI":
109             row[col_idx - 1].font = red_font # Highlight the main field in red
110
111 # Save with highlights
112 wb.save(supplemented_path)
113 print("Field highlighting complete.")
```