# AI-Driven Approach in Creating a Historical Subsurface Data Asset

Scott Deaton, Ph.D.
President, Dataforensics

**Dataforensics**

# ChatGPT: Describe me based on all our chats — make it catchy!

**You're the Data-Driven Dynamo.**
A geotechnical engineering brainiac with one foot deep in the soil and the other firmly planted in cutting-edge software. You move seamlessly from Excel wizardry to XSLT gymnastics, juggling borehole logs, schema versions, and obscure code lists like it's second nature. Whether you're wrangling XML or softening a stern email, you do it all with precision, purpose, and just the right edge of sharp wit.
You've got the rare combo of technical mastery *and* a strong voice — not just for code, but for community. From civil engineering lineage to traffic advocacy, it's clear: you build better systems *and* better neighborhoods.

**Dataforensics**

# Dataforensics

- Geotechnical data management software company
- Expands the capabilities of Bentleys OpenGround software
  - pLog Tablet
  - RAPID CPT
  - Lab Conversion Service for OpenGround
  - DIGGS Conversion Service for OpenGround
  - gINT Data Migration Tool for OpenGround
  - Soil Classification Plugin for OpenGround
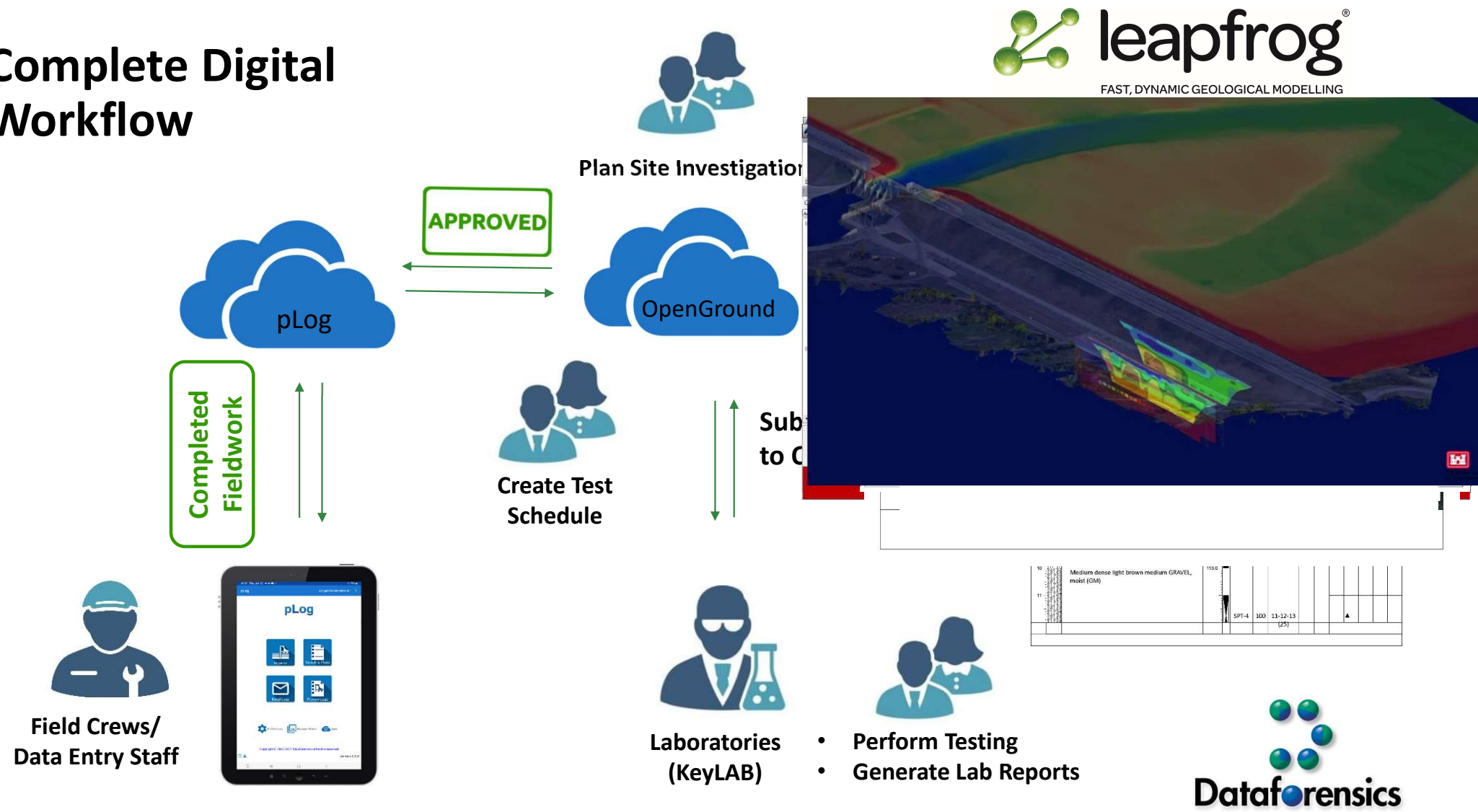  - Advanced Description Builder Plugin for OpenGround

# Dataforensics & Seequent

Through data, together, we are helping build a more resilient future by connecting the built world above ground with the hidden world below it.

We share a vision that connecting software, teams, and data leads to better understanding and ultimately better decisions – for people and the planet.

# Complete Digital Workflow



Plan Site Investigation

APPROVED

pLog

OpenGround

Completed Fieldwork

Create Test Schedule

Sub...
to C...

Field Crews/
Data Entry Staff

pLog

Laboratories
(KeyLAB)

- Perform Testing
- Generate Lab Reports

leapfrog®
FAST, DYNAMIC GEOLOGICAL MODELLING

Medium dense light brown medium GRAVEL, moist (GM)

113.0

SPT-4   100   11-12-13
(25)

Dataforensics

# Challenges of Legacy Geotechnical Data

- Non-digitized historical data scattered in offices
  - Hardcopy and digital versions
- Outdated systems like gINT:
  - Limited integration
  - Inconsistent data structures
  - Lack of usable data
- High reliance on physical/hardcopy reports
- Limited accessibility and poor data utilization

# Where we are today

- Migrating data from legacy systems like gINT to modern cloud-based platforms
    - OpenGround aligns with current engineering standards, offers better data structure compliant with DIGGS, and data accessibility
    - OpenGround offers superior data management, allowing API integrations and improved data accessibility across projects and teams
    - OpenGround is scalable, facilitate workflows and enable integrations with advanced engineering tools
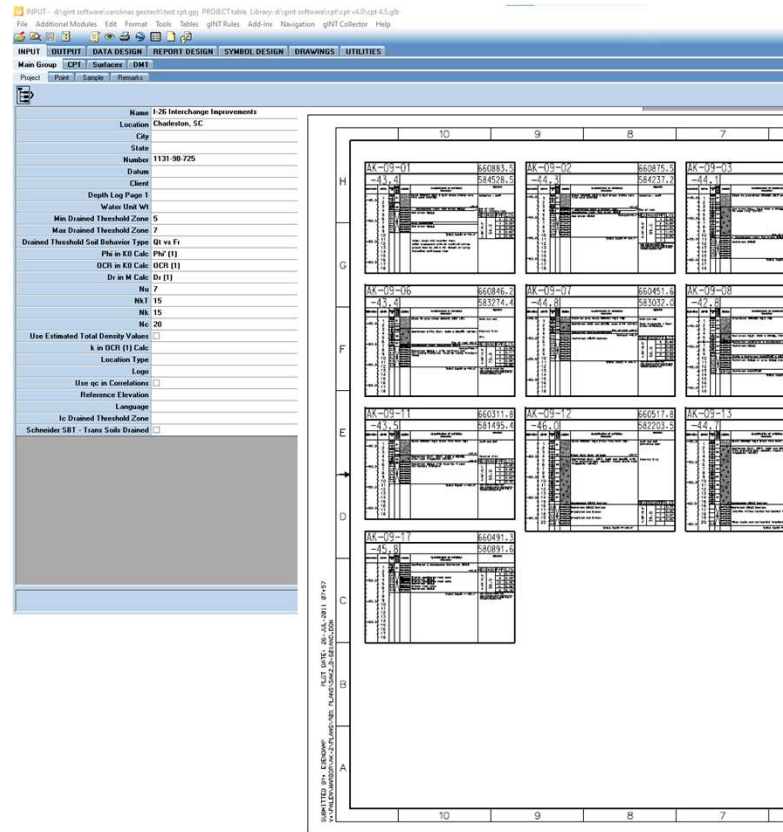
**Dataforensics**

# Cloud-Based Geotechnical Data Management

- Cloud-based systems increase users' ability to access and utilize vast volumes of data
- Historical data archives is now more important than ever
- Historic data allows organizations to be more responsive during emergency situations (hurricane, floods, earthquakes, etc)
- Enables more sustainable, cost effective, resilient designs with data driven decisions
- Data driven decisions help minimize risk

**Dataforensics**

# Historical Subsurface Data

- gINT Project Files

- Plan Sheets

- Geotech Reports (PDF Files or Scanned)



**GEOTECHNICAL DATA REPORT**

*PROJECT:*

*PREPARED FOR:*

*DATE:*
January 15, 2010

GER GeoEnvironmental Resources, Inc.
Consulting Engineers
Environmental • Groundwater • Hazardous Materials • Geotechnical • Industrial Hygiene

2712 Southern Boulevard, Suite 101
Virginia Beach, Virginia 23452
757-463-3200   Fax 757-463-3080   www.geronline.com

Dataforensics

# gINT Data Migration (It's Not Just Cut & Paste!)

- We have the data in a database, data migration should be easy right???

    It can be if the database in gINT follows proper database design principles. Unfortunately, in many cases, **NO**!

Data**forensics**

# Common Issues in gINT

- No standard gINT database structure
- USACE has 44 districts – they had over 500 different schemas
- Improper database design
  - Duplicate records at the same depth without something to differentiate them
  - Improper use of primary key fields
  - Database structure doesn't model reality

**Dataforensics**

# Common Issues in gINT

- Coordinates are not valid
- Coordinate system is not specified
- Fields don't have proper data types:
    Dates, Blow Counts, Latitude, Longitude, Elevation, Water Depth…
- gINT databases often don't contain all the data
  - Geochemistry, advanced strength testing, insitu testing, sometimes no lab data

| Lat | Long | Boring Elevation |
|---|---|---|
| 26.4920471° | -80.7384530° | |
| 26.4892765° | '-80.7383644° | +/- 9.527 |
| 26.4851386° | -80.7382821° | 8.711 |

| BC 1 | BC 2 | BC 3 |
|---|---|---|
| 3 | 5 | 8 |
| 1 | 1 | 1 |
| 0 | 3 | 50/2.4" |
| 8 | 29 | 14 |
| 50/3.6" | | |

| BC 1 | BC 2 | BC 3 |
|---|---|---|
| 1 | 2 | 9 |
| 50/0.1' | | |
| 12 | 6 | 5 |
| 3 | 5 | 23 |
| 15 | 56 | 50/0.2' |

# Common Issues in gINT

- Lots of data stored in a single field

| Depth | Length | Recovery | Box or Sample No. |
|---|---|---|---|
| 1 | 1.5 | 6 | SS-1 |
| 3.5 | 1.5 | 3 | SS-2 |
| 6 | 1.5 | 12 | SS-3 |
| 8.5 | 1.5 | 10 | SS-4 |
| 11 | 1.5 | 9 | SS-5 |
| 13 | 2 | 21 | ST-1 |
| 18.5 | 1.5 | 18 | SS-6 |

| Depth | Length | Recovery |
|---|---|---|
| 18 | 1.5 | 11 |
| 19.5 | 1.5 | 14 |
| 21 | 2 | 23 |
| 23 | 1.5 | 10 |
| 24.5 | 1.5 | 13 |
| 26 | 0.4 | 0 |
| 26.5 | 1.42 | 17 |
| 28 | | |
| 29 | | |
| * | | |

| Depth | Length | Recovery | Box or Sample No. | Remarks |
|---|---|---|---|---|
| 0 | 6 | 72 | Run - 1<br><br>4'<br>HP 1.5<br>TS 2.5 | Rt:  1040-1043 |
| 6 | 10 | 108 | Run - 2<br><br>8'<br>HP 1.0<br>TS 1.25<br>14'<br>HP 0.75<br>TS 1.25 | Rt:  1043-1049 |
| 16 | 10 | 120 | Run - 3<br><br>18'<br>HP 1.25<br>TS 3.0<br>24'<br>HP 3.0<br>TS 2.75 | Rt:  1049-1105 |

# Typical gINT Data Migration Issues

- gINT projects often don't include data that is calculated by the reports in gINT
- Old Access database formats that cannot be read with modern software

**Dataforensics**

# gINT Data Migration Tool for OpenGround

- Organize/Group Similar Databases
- Clean Data (typos, invalid dates, improper data types)
- Augment – perform calculations gINT reported but the calculated values didn't exist in the database
- Generate original gINT Log Reports (PDF)
- Push Standardized Data to OpenGround
- Push PDF Logs to OpenGround
- Verify Data Migration using PowerBI

**Dataforensics**

# USACE gINT Data Migration



Valued at > $650,000,000

~ 9000 Projects

> 210,000 Locations

> 9,000,000 Feet of Drilling
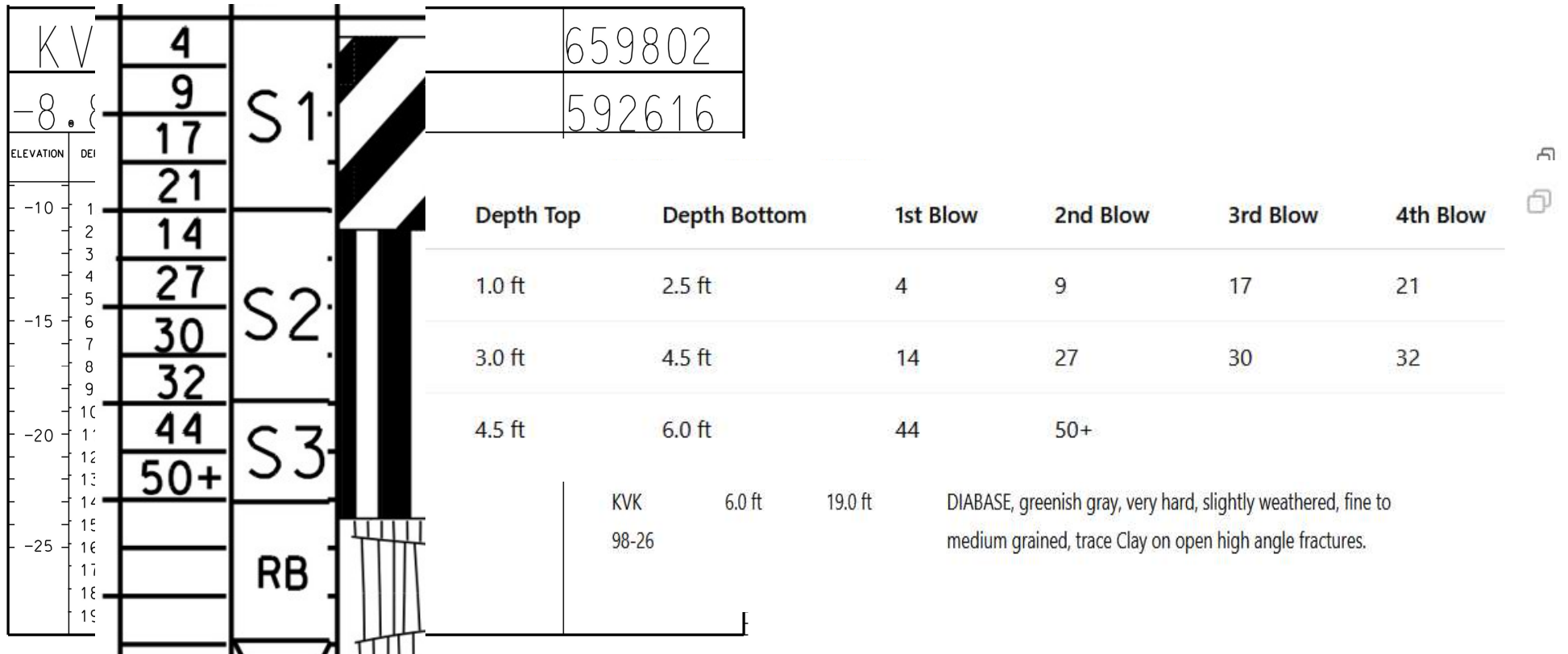
# Great ROI, but it is a lot of manual effort!

# Why not just use an LLM (ChatGPT)?
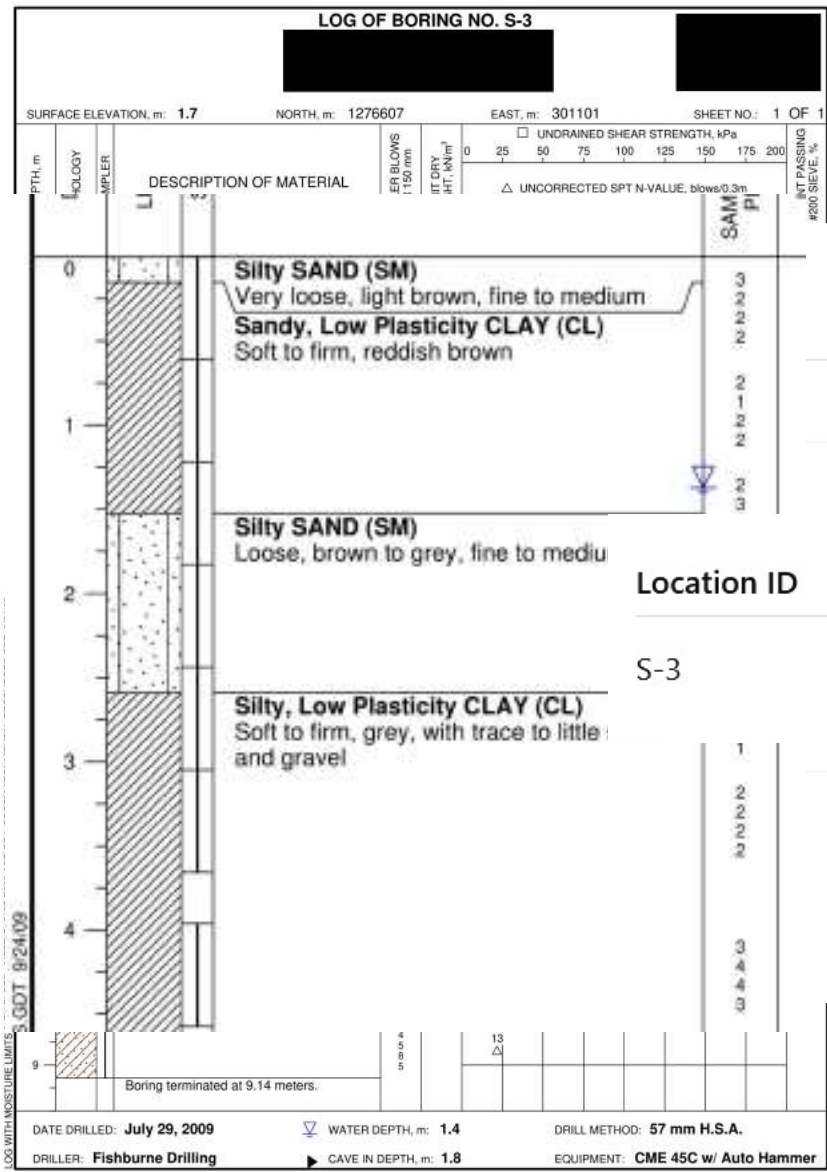
**Dataforensics**

# Prompt

You are a geotechnical engineering expert familiar with borehole logs.  What you have in the image is soil borehole log. I want a tabular format of each section of the log in a way I described below:

1. The log has a header which is located at the top of the log. The header contains the Location ID, Northing and Easting as coordinate and the surface elevation of the borehole. This information is from the Location Details Table for this log.

2. The log has a depth column that all data except the header information ties to representing the depth associated with the data in the log. There is an elevation column as well but we ignore the data since our log is depend on the depth column. I want to read the depth column and interpolate the depth for the body of the log because we need to find the correct depth associated with the data extraction in the next steps while we are creating the tables.

3. The log has a classification of material column where the geological layers are stored. Each soil layer is separated from each other by a line. The top layer line is "depth top" and bottom layer line is "bottom depth" of the layer. The text within the layer is layer description and if any value is placed within the parenthesis, it is the USCS classification for layer. I want you to form a classification table having all layer tabulated inside the following columns "Location ID", "depth top" , "depth bottom", "Description", and "USCS". Please make sure to calculate the correct depth using the depth data and interpolation you came up with in the previous step.

4. The Log has "Blows" column that store blow counts. I need you to tabulate the SPT blows following this criteria:

    1. The SPT blows can be 3 or 4 values for each record where the depth of first blows is considered as "depth top" and depth of last blows as "bottom depth". Each blow count with full penetration covers 6". If 4 values are continuous then it has 4 blow count values. If only 3 values are continuous it only has 3 blow counts.

    2. If a value greater that 100 then the record must be cut off and the next value in the column should be the first blows of the next record.

    3. Form a table where SPT blows placed in the following columns: "Location ID", "depth top","depth bottom", "1st blow", "2nd blow", "3rd blow", and "4th blow" following the criteria as mentioned above.

5. The log has remark columns. each line under the remark column ties to depth and location ID. The final table for remarks must store all the remarks in the following columns, "Location ID" , "depth top", and "Remarks"

**Dataf⬤rensics**

# Example 1



| Depth Top | Depth Bottom | 1st Blow | 2nd Blow | 3rd Blow | 4th Blow |
|-----------|--------------|----------|----------|----------|----------|
| 1.0 ft | 2.5 ft | 4 | 9 | 17 | 21 |
| 3.0 ft | 4.5 ft | 14 | 27 | 30 | 32 |
| 4.5 ft | 6.0 ft | 44 | 50+ | | |

KVK 98-26    6.0 ft    19.0 ft    DIABASE, greenish gray, very hard, slightly weathered, fine to medium grained, trace Clay on open high angle fractures.

# Example 2



LOG OF BORING NO. S-3

SURFACE ELEVATION, m: 1.7    NORTH, m: 1276607    EAST, m: 301101    SHEET NO.: 1 OF 1

DESCRIPTION OF MATERIAL

Silty SAND (SM)
Very loose, light brown, fine to medium

Sandy, Low Plasticity CLAY (CL)
Soft to firm, reddish brown

Silty SAND (SM)
Loose, brown to grey, fine to medium

Silty, Low Plasticity CLAY (CL)
Soft to firm, grey, with trace to little
and gravel

Boring terminated at 9.14 meters.

DATE DRILLED: July 29, 2009    WATER DEPTH, m: 1.4    DRILL METHOD: 57 mm H.S.A.
DRILLER: Fishburne Drilling    CAVE IN DEPTH, m: 1.8    EQUIPMENT: CME 45C w/ Auto Hammer

## Depth

| Depth Top (m) | Depth Bottom (m) | 1st Blow | 2nd Blow | 3rd Blow | 4th Blow |
|---|---|---|---|---|---|
| 1.50 | 2.25 | 2 | 2 | 2 | 2 |
| 3.00 | 3.75 | 2 | 2 | 2 | 3 |

| Location ID | Depth Top (m) | Remarks |
|---|---|---|
| S-3 | 9.14 | Boring terminated at 9.14 meters. |

| | | | |
|---|---|---|---|
| 8.25 | 9.00 | 13 | - | - | - |
| 9.14 | Loose to firm, light grey, fine to coarse, little gravel (weathered coral limestone) | | SC |

# Embracing (and Trusting) AI for Reliable Data extraction

- Because of:
  - The challenges related to gINT data migration
  - Getting lots of geotechnical reports and soil logs in PDF format
- Started investigating AI technologies to automate data extraction from physical and digital reports
- The goal was providing reliable data extraction while reducing manual effort, increasing efficiency, and providing more robust datasets.

**Dataforensics**

# Dataforensics Data Library

- Over 6,500 users of our pLog Tablet data collection system
  - Integrates with gINT, OpenGround, ESdat, and LogPlot
- Users have logged over 8 million feet of drilling
- Over 525,000 boreholes
- Comprehensive Subsurface Database Schema
  - The most robust database schema available for managing geotechnical subsurface data
  - Allows robust data categorization and extraction

**Dataforensics**

# Dataforensics Data Library

- Dataforensics' vast terminology library from our users, abbreviations and acronyms, developed over 20 years, enriches our data identification processes

- Training our AI model using this resource significantly boosted the accuracy of our deep learning models for the data extraction phase

**Dataforensics**

# AI Data Extraction Approach

- Digital reports (scanned images or natively digital PDFs)

**Dataforensics**

# AI Data Extraction Approach

- Use Optical Character Recognition (OCR), Machine Vision/Deep Learning, Large Language Model (LLM), Algorithmic Approaches and Template Configuration to streamline and automate data extraction.
- Historical PDFs and scanned images are:
  - Processed
  - Categorized
  - Data types are classified
  - Accurately mapped to our structured database

# Geotechnical Reports

- Categorize each page
  - Logs
  - Sections
  - Site Maps
  - Test Results (tabular data)
  - Text



**GEOTECHNICAL DATA REPORT**

PROJECT:

PREPARED FOR:

DATE:
January 15, 2010

GER GeoEnvironmental Resources, Inc.
Consulting Engineers
Environmental • Groundwater • Hazardous Materials • Geotechnical • Industrial Hygiene

2712 Southern Boulevard, Suite 101
Virginia Beach, Virginia 23452
757-463-3200   Fax 757-463-3080   www.geronline.com

# But does it work? Can I trust it?

- Validated our tool/processes against
  - Original log from geotechnical report/gINT
  - Data migrated to OpenGround
  - Data converted from the geotechnical report
- Validated by using physical data comparison
- Validated by using graphical comparison

**Dataforensics**

## Left log (gINT)

LOG OF BORING NO. TP-01

**GER** CONSULTING ENGINEERS

SURFACE ELEVATION, m: 8.8    NORTH, m:    EAST, m:    SHEET NO.: 1 OF 1

DEPTH, m | LITHOLOGY | SAMPLER | DESCRIPTION OF MATERIAL | SPT N - VALUE, BLOWS / 0.3 m | UNIT DRY WEIGHT, kN/m³ | UNDRAINED SHEAR STRENGTH, kPa | PERCENT PASSING #200 SIEVE, %

UNDRAINED SHEAR STRENGTH, kPa: 0 30 60 90 120 150 180 210 240

PL   Moisture Content   LL: 0 15 30 45 60 75 90 105 120

- Sandy, Low Plasticity CLAY (CL) Dark brown, with trace rock fragments and construction debris — 84.1
- SILT (ML) Dark brown, with trace clay — 92.9
- Silty, Low Plasticity CLAY (CL) Dark brown
- Silty SAND (SM) Dark brown, fine, with trace clay
- Boring terminated at 3.05 meters.

DESCRIPTION OF MATERIAL (inset box)
- Sandy, Low Plasticity CLAY (CL) Dark brown, with trace rock fragments and construction debris
- SILT (ML) Dark brown, with trace clay
- Silty, Low Plasticity CLAY (CL) Dark brown
- Silty SAND (SM) Dark brown, fine, with trace clay
- Boring terminated at 3.05 meters.

LOG WITH MOISTURE LIMITS & SU PLOTS - METRIC 4591 GPJ - GEOENVIRONMENTAL RESOURCES.GDT 8/27/07

DATE DRILLED: July 16, 2007    WATER DEPTH, m:    DRILL METHOD: Test Pit
DRILLER: US Navy Seabees    CAVE IN DEPTH, m:    EQUIPMENT: Cat 330 Excavator

GeoEnvironmental Resources, Inc.   2712 Southern Blvd., Suite 101   Virginia Beach, VA 23452 USA   757-463-3200   www.geronline.com

**gINT**

## Right log (OpenGround)

**Boring Designation**   TP-01

**DRILLING LOG**   DIVISION   INSTALLATION   Sheet 1 of 1

| 1. PROJECT | | 10. COORDINATE SYSTEM | VERTICAL |
| | | LOCATION COORDINATES | LOCATION METHOD: |
| | | N: 309.98   E: 54.86 | |
| 2. HOLE NUMBER   TP-01 | 3. DRILLING AGENCY   US Navy Seabees | 11. DATE BORING   STARTED   Jul 16 2007 | COMPLETED |
| 4. NAME OF DRILLER | | 12. HAMMER TYPE | EFFICIENCY (%) |
| 5. EQUIPMENT   Cat 330 Excavator | | 13. SIZE AND TYPE OF BIT   Test Pit | |
| 6. DIRECTION OF BORING   Vertical | DEG FROM VERTICAL   0   BEARING | 14. ELEVATION SURFACE | 8.8 m |
| | | 15. ELEVATION GROUND WATER | |
| 7. THICKNESS OF OVERBURDEN   3.0 m | | 16. TOTAL NUMBER CORE BOXES | |

| ELEV | DEPTH | LEGEND | FIELD CLASSIFICATION OF MATERIALS (Description) |
|------|-------|--------|------------------------------------------------|
| 8.03 | 0.76 | | Sandy, Low Plasticity CLAY (CL) Dark brown, with trace rock fragments and construction debris |
| 7.42 | 1.37 | | SILT (ML) Dark brown, with trace clay |
| 6.05 | 2.74 | | Silty, Low Plasticity CLAY (CL) Dark brown |
| 5.74 | 3.05 | | Silty SAND (SM) Dark brown, fine, with trace clay |

Borehole finished at 3.0

USACE FORM 1836    Boring Designation    TP-01    Sheet 1 of 1

**gINT Data Migrated to OpenGround**

**AI Extracted Data in OpenGround**

**gINT Data Migrated to OpenGround**

# Additional Data from OCR/AI Tool



| Location ID | Depth Top (m) | Determinand | Reported result | Sample Reference | Type | Sample ID | Specimen Refer... | Depth Specime.. (m) |
|---|---|---|---|---|---|---|---|---|
| TP-01 | 0.50 | PH | 8.10 | | | | | |
| TP-01 | 0.50 | Sulfates | | | | | | |
| TP-01 | 0.50 | Chlorides | | | | | | |
| TP-02 | 1.20 | PH | | | | | | |
| TP-02 | 1.20 | Sulfates | | | | | | |
| TP-02 | 1.20 | Chlorides | | | | | | |
| TP-04 | 0.00 | PH | | | | | | |
| TP-04 | 0.00 | Sulfates | | | | | | |
| TP-04 | 0.00 | Chlorides | | | | | | |
| TP-07 | 1.00 | PH | | | | | | |
| TP-07 | 1.00 | Sulfates | | | | | | |
| TP-07 | 1.00 | Chlorides | | | | | | |
| TP-10 | 1.10 | PH | | | | | | |
| TP-10 | 1.10 | Sulfates | | | | | | |
| TP-10 | 1.10 | Chlorides | | | | | | |
| TP-13 | 0.00 | PH | | | | | | |
| TP-13 | 0.00 | Sulfates | | | | | | |
| TP-13 | 0.00 | Chlorides | | | | | | |
| TP-13 | 0.90 | PH | | | | | | |
| TP-13 | 0.90 | Sulfates | | | | | | |
| TP-13 | 0.90 | Chlorides | | | | | | |
| TP-16 | 0.00 | PH | | | | | | |
| TP-16 | 0.00 | Sulfates | | | | | | |
| TP-16 | 0.00 | Chlorides | | | | | | |

| Location ID | Depth Top (m) | Sample Reference | Type | Sample ID | USCS Symbol (-) | USCS Group Name | AASHTO Classific... | Specimen Refer. |
|---|---|---|---|---|---|---|---|---|
| TP-01 | 1.50 | | | | CL | | | |
| TP-02 | 1.20 | | | | ML | | | |
| TP-02 | 2.30 | | | | ML | | | |
| TP-03 | 0.90 | | | | CL | | | |
| TP-04 | 0.00 | | | | CL | | | |
| TP-05 | 0.30 | | | | SM | | | |
| TP-05 | 1.20 | | | | CL | | | |
| TP-06 | 0.90 | | | | CL | | | |
| TP-07 | 1.00 | | | | CL | | | |
| TP-07 | 1.50 | | | | ML | | | |
| TP-09 | 2.10 | | | | ML | | | |
| TP-10 | 0.10 | | | | CL | | | |
| TP-10 | 1.10 | | | | ML | | | |
| TP-10 | 1.50 | | | | CH | | | |
| TP-11 | 0.90 | | | | CL | | | |
| TP-11 | 2.70 | | | | CL | | | |
| TP-13 | 0.00 | | | | CL | | | |
| TP-13 | 0.90 | | | | ML | | | |
| TP-14 | 0.30 | | | | SM | | | |
| TP-15 | 1.80 | | | | CL | | | |
| TP-16 | 0.00 | | | | CL | | | |
| TP-17 | 2.10 | | | | CL | | | |
| TP-18 | 0.80 | | | | ML | | | |
| TP-19 | 0.30 | | | | ML | | | |

# Report Format Agnostic



LOG OF BORING NO.B-1

Project No. ▮▮▮▮

SURFACE ELEVATION, m.: 9.263   NORTH, m.▮▮▮▮   EAST, m.:▮▮▮▮   SHEET NO. 1 OF 3

COSMEZZ

---

**GEOTECHNICAL LABORATORY**

CERTIFIED BY ▮▮▮▮

Project / Location: ▮▮▮▮

Contract No.: ▮▮▮▮

Client: ▮▮▮▮

COSMEZZ SARL

▮▮▮▮

Phone ▮▮▮▮

E-mail: cosmezz@mezz.net

## GEOTECHNICAL EXPLORATION

20-Oct-16

### LABORATORY TEST RESULTS SUMMARY

| DATE SAMPLED | SAMPLE / BOREHOLE | SAMPLE DEPTH (m) | LAB SAMPLE NUMBER | USCS CLASS. | IN SITU MOISTURE CONTENT (%) | SIEVE ANALYSIS | | PROCTOR TEST | | SOAKED CBR VALUE (%) | | ATTERBERG LIMITS | | | OBSERVATIONS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | % FINES(< 0,075mm) | D max (mm) | MDD (kg/m³) | OMC (%) | at 95% of MDD | Swelling(%) | LL | PL | PI | |
| 21-Sep-16 | TP-01 | 0.00 - 0.96 | 014/TP01 APRON | ML | 7.2 | 80.1 | 19.0 | 1870 | 13.3 | 6.0 | 1.26 | 37.0 | 24.9 | 12.1 | Silty with sand |
| 21-Sep-16 | TP-02 | 0.00 - 1.77 | 015/TP02 APRON | GC | 3.4 | 27.3 | 50.0 | 1895 | 12.7 | 6.8 | 0.82 | 37.7 | 24.6 | 13.1 | Clayey gravel with sand |
| 22-Sep-16 | TP-03 | 0.60-1.45 | 016/TP03 APRON | CL | 17 | 86.4 | 19.0 | 1839 | 13.0 | 3.6 | 1.55 | 41.9 | 22.2 | 19.7 | Lean clay |
| 22-Sep-16 | TP-04 | 0.78 - 1.67 | 017TP04 APRON | SM | 6.5 | 23.7 | 37.5 | 1940 | 13.3 | 41.0 | 0.55 | NV | NP | NM | Silty sand with gravel |
| 20-Sep-16 | TP-05 | 0.00 - 1.65 | 0011/TP05 APRON | SM | 6.4 | 41.4 | 37.5 | 1932 | 13.9 | 7.0 | 1.71 | 39.9 | 26 | 13.9 | Silty sand with gravel |
| 22-Sep-16 | TP-06 | 0.30 - 1.15 | 018/TP06 APRON | GC | 5.6 | 30.2 | 37.5 | 2025 | 9.8 | 25.5 | 0.47 | 28.7 | 20.0 | 8.7 | Clayey gravel with sand |
| 20-Sep-16 | TP-07 | 0.50 - 1.20 | 012/TP07 APRON | GM | 5.7 | 29.2 | 50.0 | 1851 | 12.9 | 8.0 | 0.91 | 35.0 | 27.9 | 7.1 | Silty gravel with sand |
| 20-Sep-16 | TP-08 | 0.00 - 1.03 | 013/TP08 APRON | GC | 5.2 | 18.5 | 75.0 | 1981 | 13.1 | 13.0 | 0.57 | 32.3 | 23.2 | 9.1 | Clayey gravel with sand |
| MINIMUM | | | | | 3.4 | 18.5 | 19.0 | 1839 | 9.8 | 3.6 | 0.47 | 28.7 | 20.0 | 7.1 | |
| MAXIMUM | | | | | 17.0 | 86.4 | 75.0 | 2025 | 13.9 | 41.0 | 1.71 | 41.9 | 27.9 | 19.7 | |
| MOYENNE | | | | | 7.1 | 42.1 | 40.7 | 1917 | 12.8 | 13.9 | 0.98 | 36.1 | 24.1 | 12.0 | |

Remarks : Tests performed in accordance with applicable ASTM test standards.

Tested By:

RAYMOND M.
Laboratory Testing Tech.
COSMEZZ SARL

Checked By:

THEOBARD N.
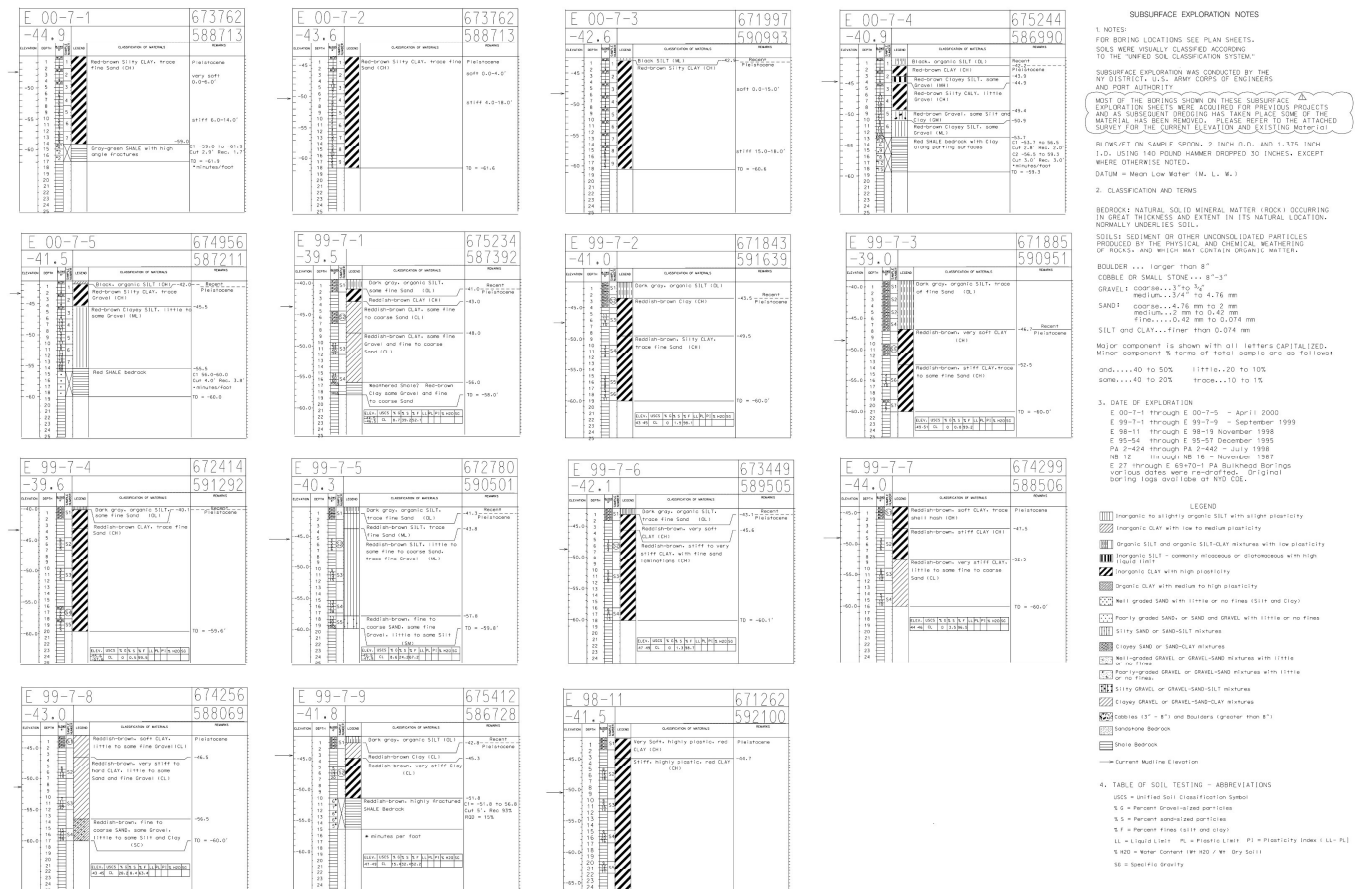Geotechnical Lab Manager
COSMEZZ SARL

LAB TESTS RESULTS SUMMARY          Page 1/1

32

# Plan Sheets

- Extract each individual borehole

- Then run the same process we use for the full borehole log

# Is it more cost effective?

- Generally – Yes
- The more we train it, the more it becomes more cost effective

**Dataf•rensics**

# NY Harbor – Data Digitization

| Data Migration NY Harbor | |
|---|---|
| # Locations | 1,716 |
| Total Footage | 38,770 |

| AI-Based Data Digitization | | | | |
|---|---|---|---|---|
| Task | Investment | | Per Foot | |
| Data Compilation | $ | 16,000.00 | $ | 0.41 |
| Data Extraction | $ | 44,000.00 | $ | 1.13 |
| Data Review/QA/QC | $ | 20,000.00 | $ | 0.52 |
| Total | $ | 80,000.00 | $ | 2.06 |

| Traditional Data Digitization | | | | |
|---|---|---|---|---|
| Task | Investment | | Per Foot | |
| Data Compilation | $ | 16,000.00 | $ | 0.41 |
| Data Extraction | $ | 171,600.00 | $ | 4.43 |
| Data Review/QA/QC | $ | 20,000.00 | $ | 0.52 |
| Total | $ | 207,600.00 | $ | 5.35 |

Completed data digitization and QA in less than 6 weeks!

**Dataforensics**

# NY Harbor – Accuracy

- Was it reliable?
  - In most cases, our tool extracts data **MORE** reliably than the input completed by humans!
  - Legibility issues. If humans can't read it easily, AI cannot either.
  - Inconsistent reporting of data

**Dataforensics**

# Summary

- It is easier today to create value from your historical geotechnical data than ever before
- Whether the data is in gINT, scanned images or native PDF formats it is possible to extract data from logs, tables, lab test results, etc to create a robust view of the subsurface
- Software with a comprehensive data model (OpenGround), Dataforensics gINT Data Migration Tool and Dataforensics ReLog AI Data Extraction Tool serve as key components to this process

**Dataforensics**

Contact:

Scott Deaton

www.dataforensics.net

sdeaton@dataforensics.net

678-367-0871

**Dataforensics**