# GPT-Powered Data Insights: Advancing Clinical Data Analytics with Large Language Models (LLMs)

PHUSE US Connect 2025 - ML25

Mel Hullings, Director, Data Management & Analytics

March 19, 2025

**Formation Bio**

# Agenda

- Intro to AI in Analytics

- GPT-Powered Workflows

  1. Data Cleaning & Transformation

  2. Data Modeling & Analysis

  3. QC & Validation of AI-Generated Outputs

- Key Takeaways & Future Directions

# AI-Powered Clinical Data Geneticists

**Why LLMs in Clinical Data Analytics?**

Much like geneticists decode DNA to understand biological functions, LLMs can decode clinical data to extract meaningful insights while overcoming challenges such as:
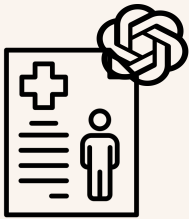
- **Data Growth & Complexity:** exponential growth of data (structured & unstructured)
- **Data Quality:** cleaning and standardizing data for analysis is slow and prone to error
- **Regulatory Requirements:** FDA & EMA require structured, traceable, and validated data reporting

💡 *Manual processes are no longer scalable.*
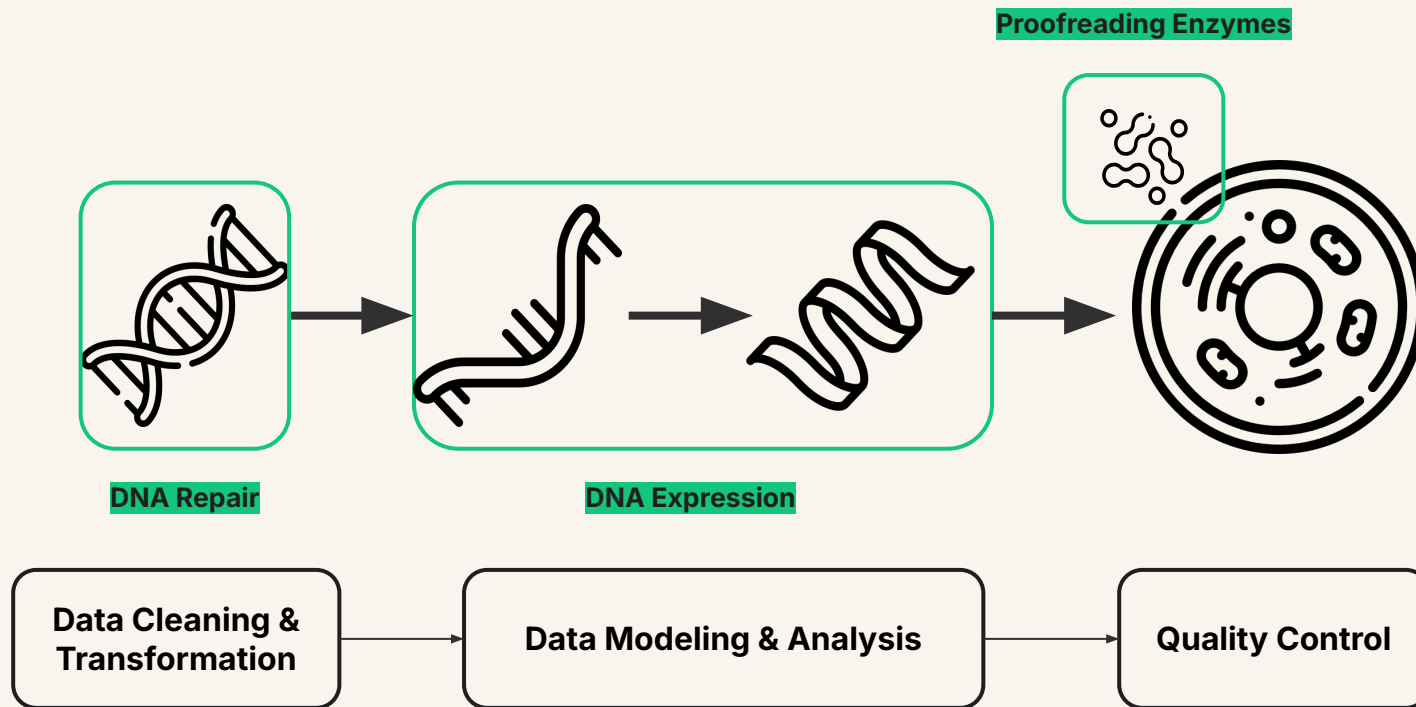
# Transforming Clinical Data Analytics with LLMs

✅ **Natural Language Understanding (NLU):** AI can put data into context to interpret study protocols, medical narratives, and investigator comments.

✅ **Pattern Recognition in Structured Data:** GPT models can flag inconsistencies in case report forms (CRFs) and clinical datasets.

✅ **Automated Report Generation:** AI can summarize safety trends, efficacy endpoints, and statistical analysis results.

📌 *Example: A recent study demonstrated that GPT-4 outperformed traditional rule-based NLP methods in summarizing unstructured EHR data and extracting clinical phenotypes*

# GPT-Powered Workflows

Proofreading Enzymes

DNA Repair

DNA Expression



| Data Cleaning & Transformation | → | Data Modeling & Analysis | → | Quality Control |

# Data Analysis Process



**Gather Inputs** → **Prepare Data** → **Analyze Data** → **Deliver Outputs**

Actions | Data Analysis | Data Analysis | Actions

- Write SQL
- Execute SQL

- Describe data
- Update columns
- Anomaly detection

- Brainstorm questions
- Turn a dashboard alive
- Advanced data science

- Save & download files
- Send an email
- File a ticket

*Data analysis is an iterative process where you continually explore questions, data, and results. There are opportunities for AI to support you throughout this process!*

# ChatGPT's "Data Analyst" Custom GPT

Collaborate with ChatGPT's AI Data Analyst Custom GPT
to get from data to insights faster.

**Capabilities:**

- Chat with your data
- Get proactive insights and suggestions
- Create customizable, presentation-ready charts

**Why ChatGPT is recommended for data analysis:**

- Automates repetitive analyses
- Advanced data exploration and insights
- Enterprise-grade data privacy and security
- No training the LLM based on your data, ever

**By ChatGPT**
GPTs created by the ChatGPT team

| | | Data Analyst |
|---|---|---|
| 1 | >_ | Drop in any files and I can help analyze and visualize your data. |
| | | By ChatGPT |

*Note: Reference your company's data governance guidelines before uploading data.*

# Case Study: Using AI to Analyze the Osteoarthritis Initiative (OAI) Public Database

Osteoarthritis Initiative. (2025). Osteoarthritis Initiative (OAI) Database. National Institute of Mental Health Data Archive.

# 1. Data Cleaning & Transformation



DNA Repair

| Data Cleaning & Transformation | → | Data Modeling & Analysis | → | Quality Control |

# Prompts to assess data quality and add columns

**Variable Guide:  by Data Collection Form**

**Enrollment Visit p 1**

**Variable Name:**     V00COHORT
**Label:**     EV:Subcohort assignment (calc)
**SAS Dataset:**     Enrollees
**Release Comments:**     Notes available: see individual SAS dataset documentation for details

| Category | SubCategory |
|---|---|
| Knee pain/OA status | Knee pain/OA status |
| Study eligibility | Other |

| Value | N | % | Cumulative N | Cumulative % |
|---|---|---|---|---|
| 1: Progression | 1,390 | 28.98 | 1,390 | 28.98 |
| 2: Incidence | 3,284 | 68.47 | 4,674 | 97.46 |
| 3: Non-exposed control group | 122 | 2.54 | 4,796 | 100.00 |

➜ **Data cleaning:** What are some anomaly checks to run in this dataset?

➜ **Data transformation:** Add a column after "V00COHORT" to the enrollees dataset labeled V00COHORT_Verbatim" and in that column for each subject's ID, enter the text "Progression" if "V00COHORT" is a 1, "Incidence" if "V00COHORT" is a 2, and "Control" if "V00COHORT" is a 3.

# 2. Data Modeling & Analysis



DNA Expression

Data Cleaning & Transformation → Data Modeling & Analysis → Quality Control

# Prompts to analyze and interact with data

➔ **Data summarization:** Will perform upon upload and then specify further, such as asking how many subjects in the "outcomes 99" table had a knee replacement

➔ **Data viz:** Create a bar chart to show the breakdown of race and ethnicity for each cohort.

➔ **Data interactivity:** Generate R code for a Shiny app that allows you to explore data.

| Custom Population Key Metrics | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Population | | N | 1-yr | 2-yr | 3-yr | 4-yr | 5-yr | Crude Risk |
| Custom | | 4796 | 99.5% | 98.6% | 97.7% | 96.4% | 95.2% | 1.1% |

# 3. QC & Validation of AI-Generated Outputs



Proofreading Enzymes

Data Cleaning & Transformation → Data Modeling & Analysis → Quality Control

# Quality Control

Regardless of the classification of data, critical thinking is required to evaluate risk based on the data, information, and decisions being made for each use case, including the ability to quality control (QC) outputs because AI can be inconsistent, factually incorrect, and hallucinate.

Tips to check quality of AI analyses:

- Logical checks
- Ask for sources
- Output verification (i.e. human-in-the-loop)
- Awareness of bias in the data

# Key Takeaways

1. **LLMs enhance data processing** by automating cleaning, transformation, and modeling, improving efficiency of data insights.

2. **AI-powered analytics facilitate data quality and exploration**, making complex clinical trial data more interpretable and actionable.

3. **Shiny apps with LLM integration** allow for real-time interaction with clinical datasets.

4. **Validation & QC remain critical**—AI-generated insights must be rigorously verified to ensure accuracy and compliance.

# Future Directions

🚀 **Advancements in Multi-Modal AI**

- Next-gen models (e.g., GPT-5, Med-PaLM 2) will integrate **text, images, and structured data**, enabling deeper insights from multi-source clinical datasets.

📜 **Regulatory Landscape Evolution**

- As AI adoption grows, **regulatory bodies (FDA, EMA, MHRA) will refine guidance** on AI-driven analytics, emphasizing transparency and validation.

🧠 **Ethical & Interpretability Challenges**

- While AI augments clinical decision-making, **explainability remains a challenge**—future AI models must improve interpretability and trustworthiness.

# Questions?