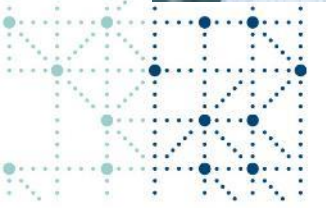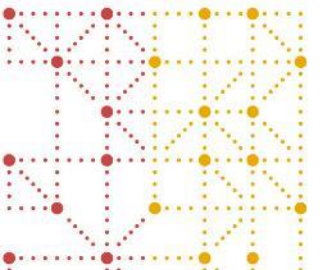2024 CDISC + TMF
EUROPE INTERCHANGE

BERLIN

24-25 APRIL: CONFERENCE & EXPO | 22, 23, 26 APRIL: TRAININGS

**AI agents as Co-Scientists: Pioneering the Next Frontier in Pharma with Advanced Clinical Data Insights**

Presented by Sam Khalil, PhD
VP Data Insights and FounData
Novo Nordisk

# Meet the Speaker

## Sam khalil

Title: VP Data Insights and FounData

Organization: Novo Nordisk

Sam Khalil brings over 15 years of experience in the biopharmaceutical industry, where he has contributed to the integration of life sciences, healthcare, and advanced technologies to support patient care. As the Vice-President of Data Insights at Novo Nordisk in Denmark, a leading biopharmaceutical company, he plays a key role in harnessing data to inform decision-making processes.

His expertise covers a broad range of areas, including fundamental research, pre-clinical studies, and clinical trials, providing him with a well-rounded understanding of the drug development process. Currently, he is leading the development of FounData, a comprehensive clinical data lakehouse at Novo Nordisk. This platform is designed to consolidate clinical trials data, along with proteomics, genomics, and imaging data, to facilitate more informed research and development strategies.

# Disclaimer and Disclosures

- *The views and opinions expressed in this presentation are those of the author(s) and do not necessarily reflect the official policy or position of CDISC.*

# Agenda

# **FounData |** Find & Access securely all data sources



FOUNDATA
BRING · DATA · TO · LIFE

**~1,700** Clinical Trials

**>900K** Patients

**MULTI - MODAL**

Text
Tabular
Imaging
Genomic
Proteomics
...

# FounData | Key components



**FIND**　　　**ACCESS**　　　　　　　　　　　**CONNECT**　　**SOLVE**

**For more information on:**
Data Quality considerations when providing access to pooled clinical trial data for secondary Use

*Sid Prabhu, Jonathan Dean*

# Imagine if... |

we could **reduce further the barriers to data** and **accelerate** **complex insights generation** from all our **clinical trials data**

**In the past:**

Only professional photographers or serious hobbyist could capture images.

From **complex process**, requiring **specialized knowledge of chemistry** and **optics**, as well as access to **cumbersome equipment**.

The art of photography is **now accessible to all.**

Without worrying about the technical intricacies behind the lens.

# Key areas | Data – Governance - Technology -

**FOUNDATIONAL**

**Data Platform**

**A seamless but controled**

**GOVERNANCE**

**ABSTRACTING**

**complex analytics tasks with**

**RESPONSIBLE AI**

# Can we just add Generative AI on top ? |
## Not yet…. For complex or multi-steps scientific analytics



**You**
Given the 4 numbers 4 9 10 13 what would be a simple arithmetic operation (*-+/) to obtain 24.

$$9 * (13 - 10) + 4 = 24$$

Here's the breakdown:

1. Subtract 10 from 13, which equals 3.
2. Multiply 9 by the result (3), which equals 27.
3. Add 4 to get 24.

# Current LLMs are fast thinkers |
## We need PLANNING, RESONING & CRITIC



THE NEW YORK TIMES BESTSELLER

THINKING, FAST AND SLOW

DANIEL KAHNEMAN

WINNER OF THE NOBEL PRIZE IN ECONOMICS

"[A] masterpiece... This is one of the greatest and most engaging collections of insights into the human mind I have read." —WILLIAM EASTERLY, Financial Times

## SYSTEM 1



**Fast & Automatic,**

*But by design error prone and instinctual*

$$2+2 = 4 \quad \checkmark$$

$$9 * (13 - 10) + 4 = 24 \quad \times$$

## SYSTEM 2



**Slow & Deliberate**

*but requires conscious effort / planning*

In patients where the disease is still Progressing, can we identify any patterns in our data that could guide better outcomes?

cdisc

# We need Tools | Essential for complex tasks



M. Hilbert, Online Course *Digital Technology & Social Change*, University of California: https://canvas.instructure.com/courses/949415

**Tool use is a distinguishing characteristic of human beings.**
We create, modify and utilize external objects to do things that go beyond our physical and cognitive limits. Equipping LLMs with external tools can significantly extend the model capabilities.

**Some example of tools:**

Big Data Analytics

Database querying

Genomic & Proteomic analysis

Internet & intranet access

..... Many more use case specific

# We need Memory | Essential for complex tasks

| | System Messages / Context | |
|---|---|---|
| **Short term** <br> Context and prompting | **Fixed** | **Variable** <br> ***Programaticaly*** |
| **Long term** <br><br> Clinical and regulatory Documents <br> Past interactions and insights <br> Molecules interactions | **RAG** | **RAG** <br> **+ Knowledge Graph** |

# FOUNDATA AGENTS | Advanced capabilities with focused purpose

**LLM**

**AGENT**

| PLANNING | **Decomposition of the task in small steps**, reflection on the plan and self-critics (CoT, ToT etc…) |
|---|---|
| TOOLS | A **list of all the tools** to which the agent has access, with instructions on how to access them |
| MEMORY | R**elevant memory** items (based on the question) from **past conversations** with the user<br>New "memories" can be generated and stored if required |

# AGENTS | They work better in teams with each having a focused expertise and set of tools



DATA ANALYSIS PLANNER

DATABASE EXPERT

CLINICAL DATA EXPERT

MANAGER

GOVERNANCE MANAGER

BIOSTATS / CODING EXPERT

SCIENTIFIC REVIEWER / CRITIC

cdisc

# FounData AI agents | Building Blocks



**Data Engineer**

**FounData, Agent Team**

I work as a data engineer in the FounData agents unit. I specialize on loading data from different systems including databases, blob storage & delta tables.

**+ Follow**   **Message**   **More**

**System Prompt**

You are a senior data engineer. You help a team of autonomous agents fulfill user queries. You should always provide the exact location of the requested data along with the metadata and column names.

**Tools**

☐ **query_clinical_trial_db:** query the secure clinical trial database for SDTM or ADAM data.

☐ **crawl_blob_storage_container:** list all the files in a given blob storage container

**SDTM Expert**
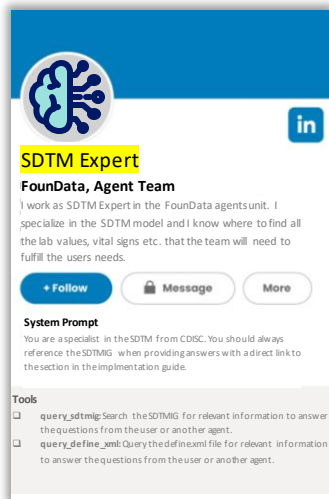
**FounData, Agent Team**

I work as SDTM Expert in the FounData agents unit. I specialize in the SDTM model and I know where to find all the lab values, vital signs etc. that the team will need to fulfill the users needs.

**+ Follow**   **Message**   **More**

**System Prompt**

You are a specialist in the SDTM from CDISC. You should always reference the SDTMIG when providing answers with a direct link to the section in the implmentation guide.

**Tools**

☐ **query_sdtmig:** Search the SDTMIG for relevant information to answer the questions from the user or another agent.

☐ **query_define_xml:** Query the define.xml file for relevant information to answer the questions from the user or another agent.

## CHARACTERISTICS OF AGENTS

Specialized with key attributes:

- **NAME:** agent identifier
- **DESCRIPTION:** agent specialty and skills
- **SYSTEM PROMPT:** Agent personality and directiv (e.g., programmer)
- **LLM Config:** Which LLM to use
- **TOOLS:** Annotated user defined functions that th agent can utilize (e.g., read data from secure database)

# Group chats | How agents interact with each other to solve complex queries

**Group Chat Manager**



## What is a GROUP CHATS?

Group chats define how agents can interact with each other.

- **Simple group chats** can have a linear interaction like "round robin" or a "directed graph", where the user has predetermined the flow of the conversation

- **Complex group chats** like the "generate plan", will have an agent compare the user query with a list of all available agents to produce a plan with a set of tagged steps (task list). This provides additional flexibility.

# **Workgroups |** Multiple agents iterating over steps in the plan
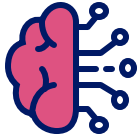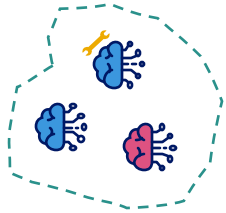
**Group Chat Manager**



Visualization Workgroup

Statistical Programming Workgroup

## What is a WORKGROUP?

To the group chat manager, a workgroup looks like a single agent, but in reality, the workgroup works as a small group chat on its own. This can be highly beneficial in multiple scenarios:

- **Programming tasks:** Iterating over coding tasks with a set of agents like a *programmer*, a *code_executioner* and a *qa_specialist,* can give better results than a single programmer.

- **Visualisation tasks:** Producing good figures or graphs can be improved by having "visual_critique" agent iterate with the programmer.

# Workflows | Joining it all together gives a highly customizable framework for LLM applications

**Workflow**

Group Chat Manager



## What is a WORKFLOW?

Workflows combine all the building blocks into a single configuration. A workflow can be general to serve many needs, or highly specialized for a narrow LLM application. Workflow contain:

❑ **Agents:** List of available agents

❑ **Planner:** Used for generating plans for group chats

❑ **Group Chats:** How the agents interact

❑ **Manager:** Agent for managing the group chat

**You** 12:39 PM

Create a pie chart of the race distribution of all female patients in the CDISCPILOT01 trial. Use data from the Demographics domain, include names only in the labels and title the graph 'Race Distribution - F'.

**F** **FounData Assistant** 12:39 PM

Using data_engineer_lead ⌄



DATA ANALYSIS

FounData Co-Scientists agents working in synergy

FOUNDATA
BRING·DATA·TO·LIFE

successfully fulfilled. The Demographics domain CSV file is located at the following full path: /app/data_csv/dm.csv. The necessary data has been extracted for all female patients, focusing on their race information.

**Lead Programmer** 12:39 PM

The necessary Python libraries have been successfully installed, and the code to create the pie chart of the race distribution for the extracted female patient data has been executed. The pie chart has been saved as a JSON file named 'race_distribution_f.json' in the local working directory. The next step is to retrieve the 'race_distribution_f.json' file and present it to the user with the title 'Race Distribution - F'.

**Data Analyst Lead** 12:39 PM

The task has been successfully fulfilled. The race_distribution_f.png has been retrieved from the local working directory and presented with the title "Race Distribution - F." The plot was successfully displayed and the task is now complete.
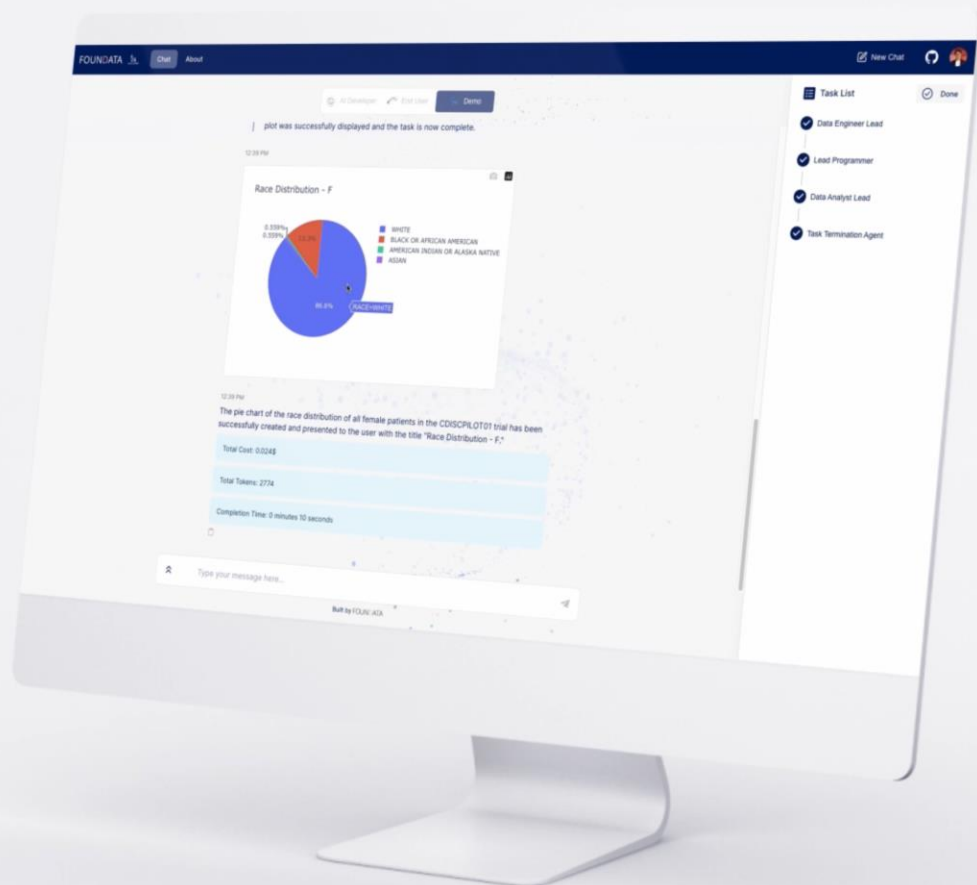
12:39 PM

Race Distribution - F

## Analytical output
Required multiple steps:
1- Planning
2- Data Retrieval
3- Understanding of data tables
4- Writing SQL queries and reviewing code
5- Creating visualisation

FOUNDATA
BRING · DATA · TO · LIFE

# Built on a Scalable Framework

Agents, workgroups, and workflows are fully configurable and most importantly can be tested in units during deployment.

# THANK YOU!

HKLL@NOVONORDISK.COM

cdisc