



How Machine Learning Can Automate Data Anomaly Detection?

Streamlining Query Management to Accelerate Database Lock and Regulatory Submission

John Hall, Senior Vice President

4th March 2025

Agenda

1

Data management review today – challenges and opportunities.

2

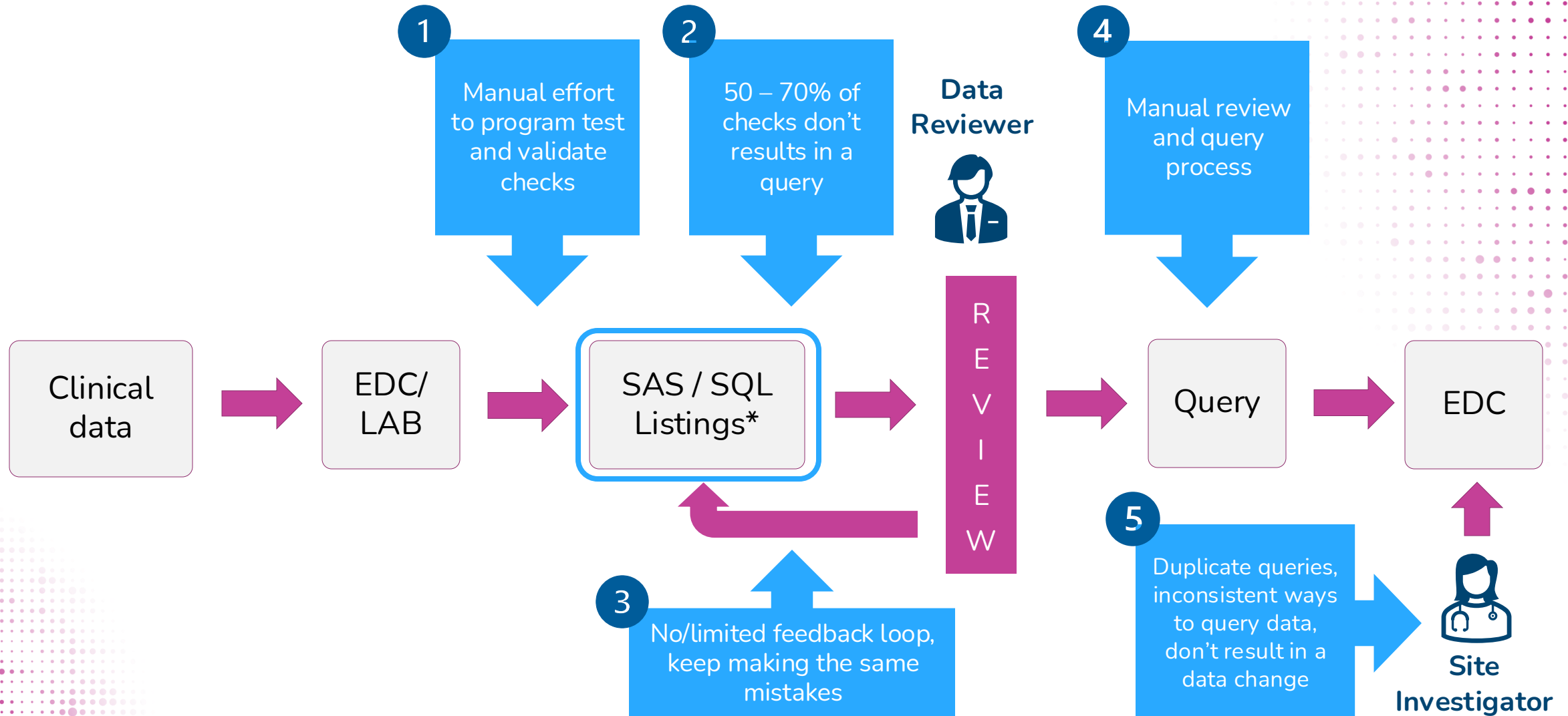
AI/ML to the rescue – or is it? Ensuring performance, scalability and regulatory compliance.

3

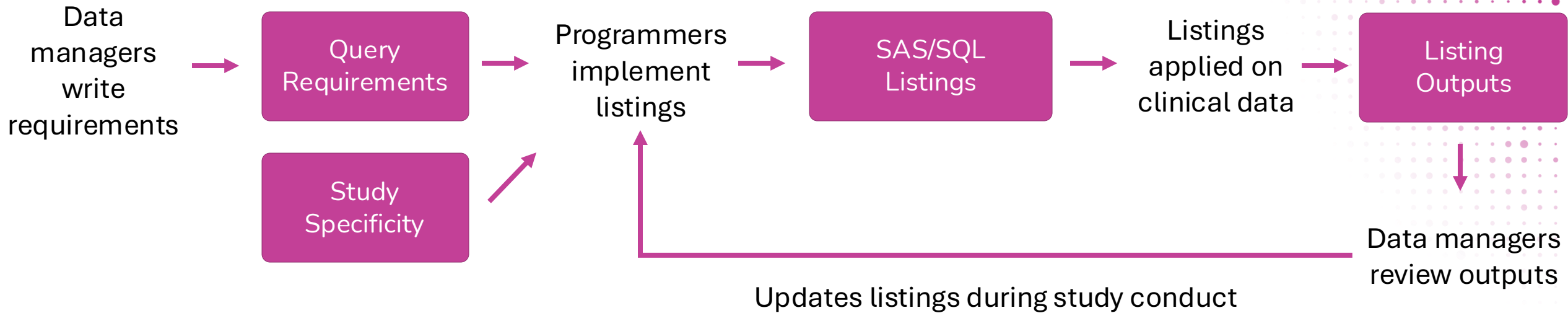
What could the future be like?

Data Cleaning Today

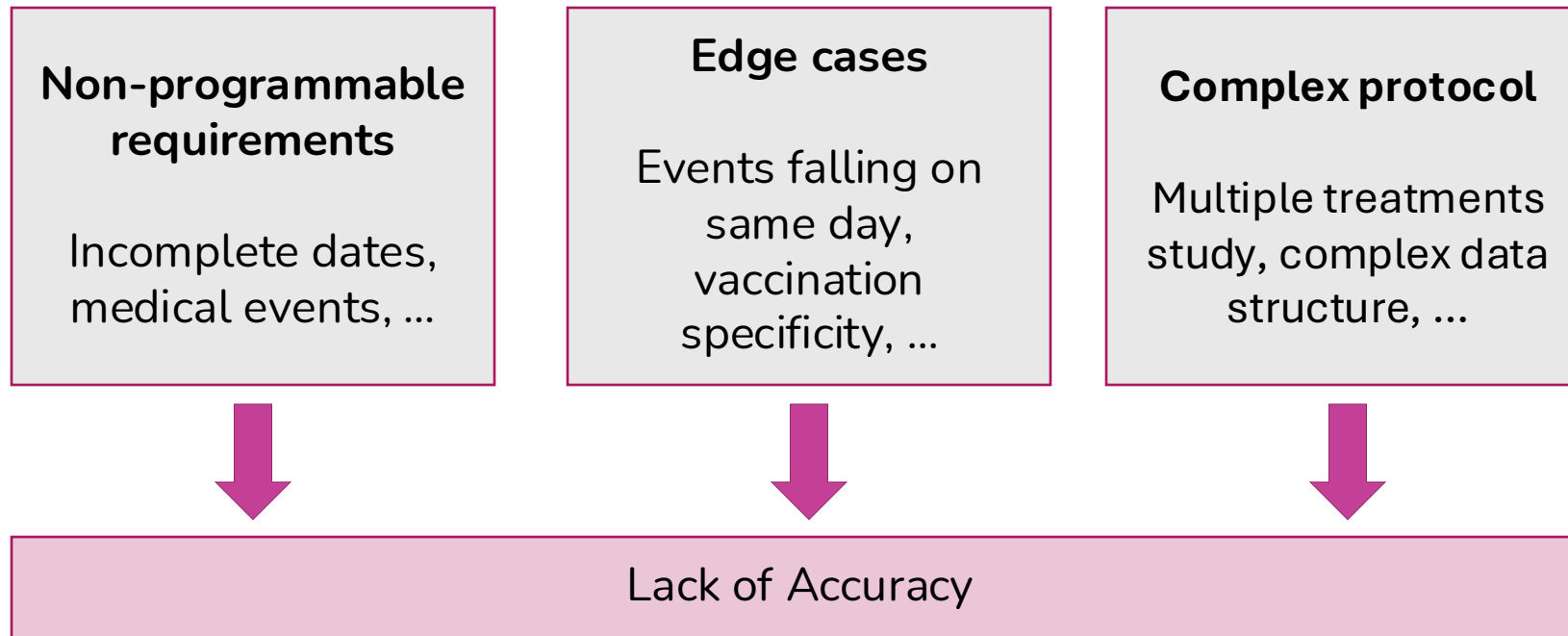
Data Queries: Current Challenges



SAS/SQL Listings: Current Process



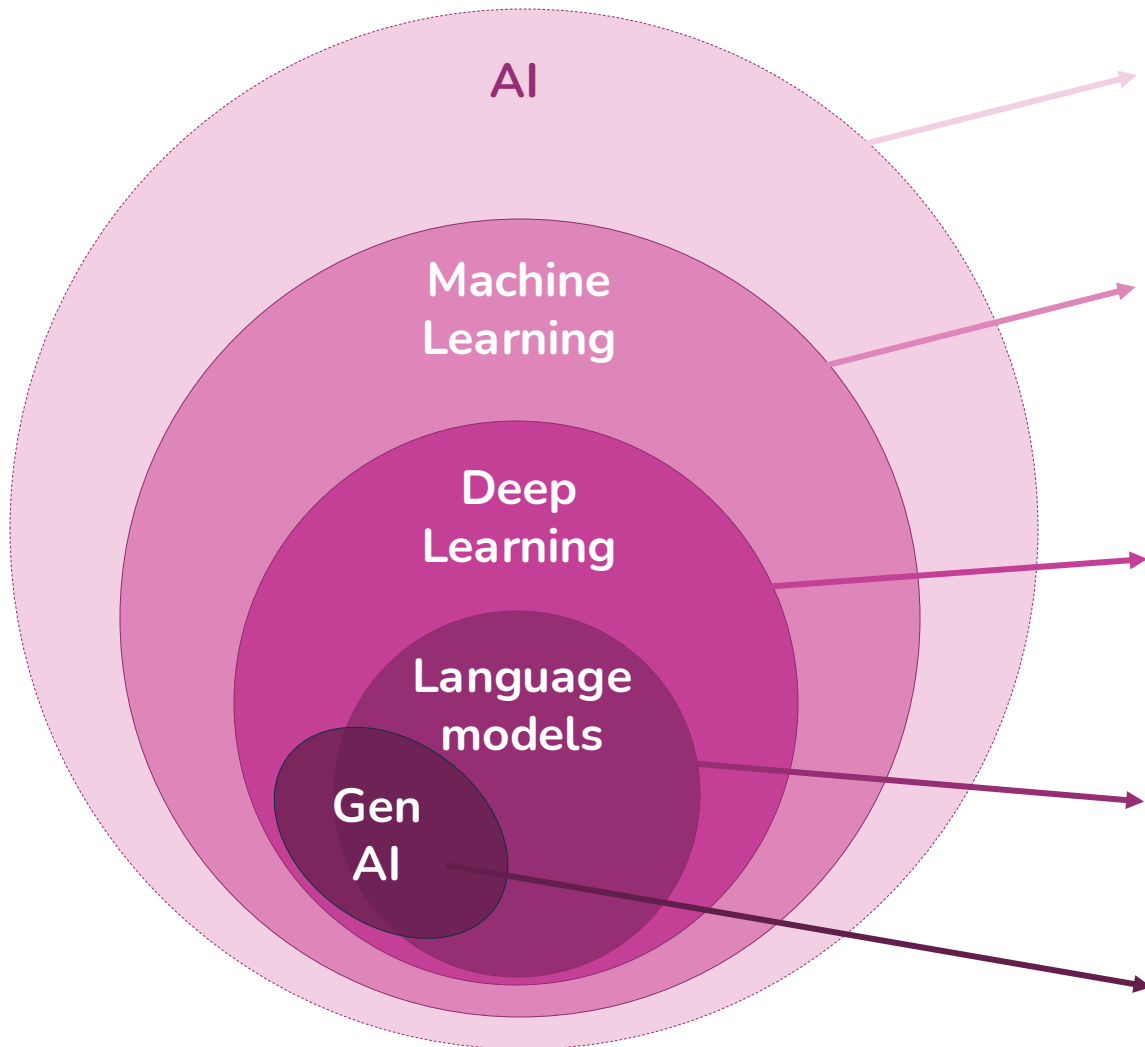
SAS/SQL Listings: Current Challenges



AI/ML to the rescue – or is it?

Ensuring performance, scalability and regulatory compliance.

The terminology is changing....

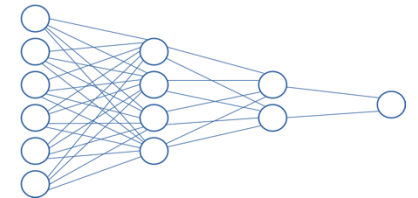


AI: a set of techniques? A solution? A goal?

Machine learning models learn patterns from data



Deep learning models are deep neural networks able to process complex data (images, text, etc.)

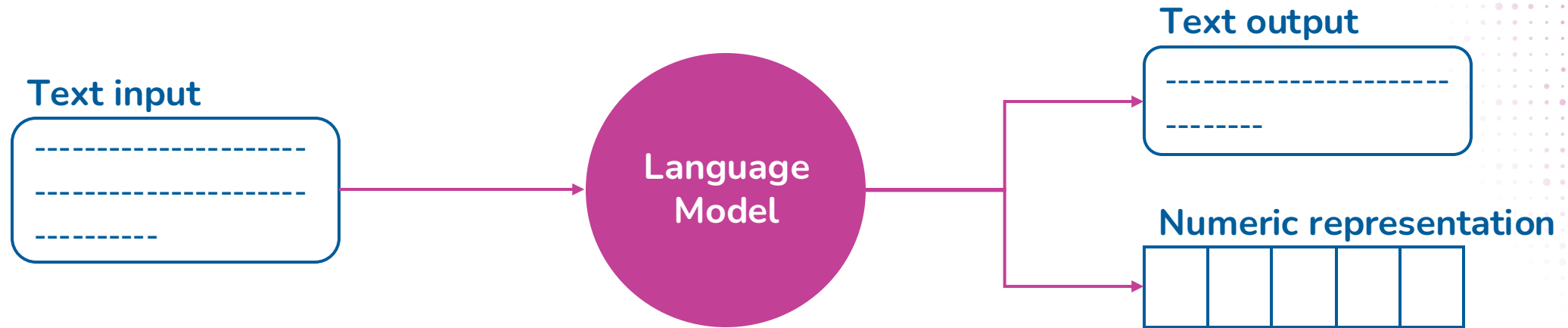


Language models learn probabilistic relationships in text

$$p(\text{"The dog barks"}) > p(\text{"The cat barks"})$$

Gen AI refers to deep learning models (language models or others) able to generate data (text, images, etc.)

Language models have disrupted the field



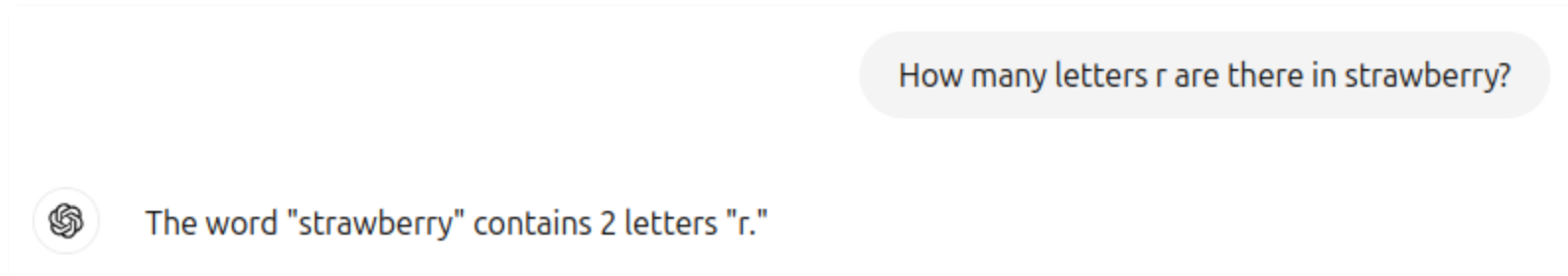
Why using language models for clinical data review?



But be careful....



Language models make mistakes, sometimes unexpected ones



GPT-4o (Sep 2024)



Needs process to avoid this behavior

Requirements for language models



Specialized in
clinical data

- Understand the concepts related to clinical data
- Understand clinical data structure



Understand
the task

- Know what action is expected
- Know how to respond



Actionable
answers

- Possible to parse the answer
- Take decision from the answer







Reliable and
consistent

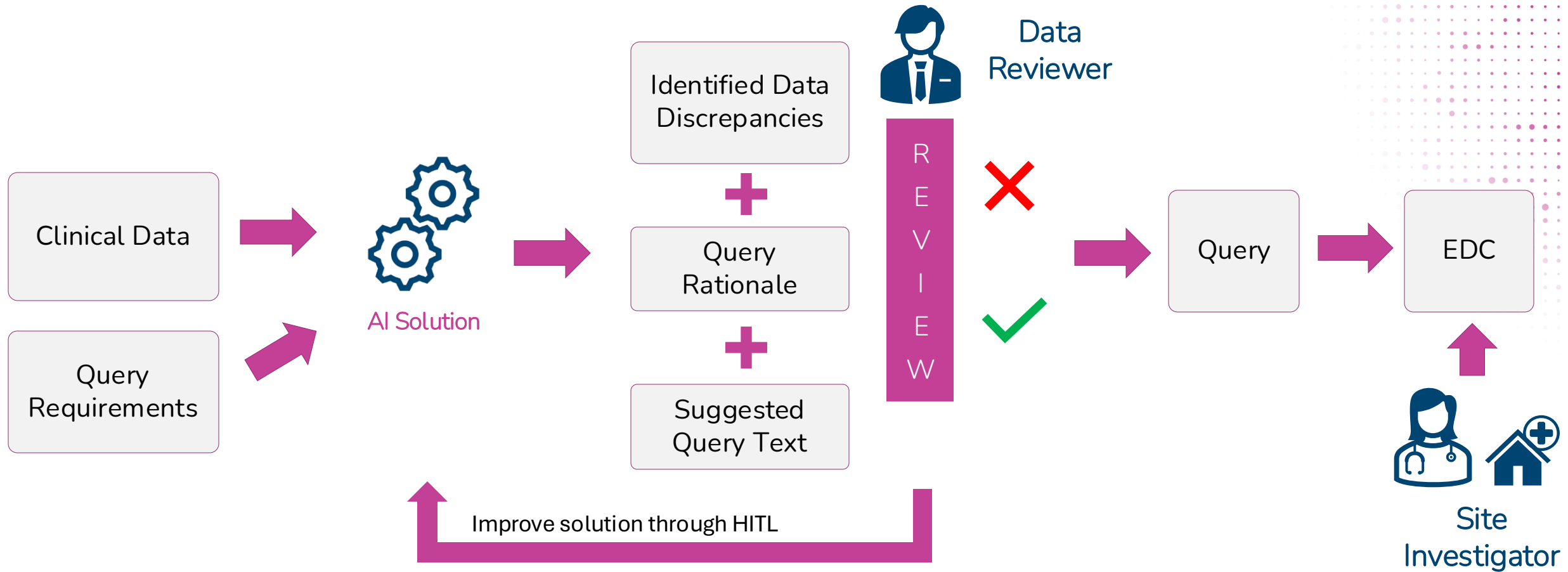
- Consistent in the process
- High accuracy and no unexpected mistakes

How to achieve these requirements?

»»» Fine-tuning open-source language models!

 Open-source model	 Clinical data for training	 Examples to learn from	 SME support
<ul style="list-style-type: none">• Existing powerful open-source models• Full control on these models	<ul style="list-style-type: none">• Access to diverse clinical studies	<ul style="list-style-type: none">• Examples to understand the task• Consistency in the ground truth	<ul style="list-style-type: none">• Feedback from solution's output

AI Driven Detection of Data Queries: How?



Example Query Scenario

Scenario

Identify Duplicate/Overlapping Concomitant Medication/Therapy records

Study specific exceptions
(wouldn't raise a query)

1. Except for such overlapping Concomitant Medications where the **end date** of 1st CM and **start date** of the overlapping CM are the **same day** and where the **frequency has changed**.
2. Except where the CMs have **different start dates** but **missing end dates**, but that the type of CM is a **vaccination**.
3.

E.g. 1

Con Med Treatment

Con Med Treatment		Dose	Frequency
Con Med 1	20/07/2023	7.5mg	BID
Con Med 2		7.5 mg	TID

E.g. 2

Con Med Treatment

Con Med Treatment		CM Category
Con Med 1	01/07/2023	Vaccination
Con Med 2	08/07/2024	Vaccination

Example Query Scenario

Scenario

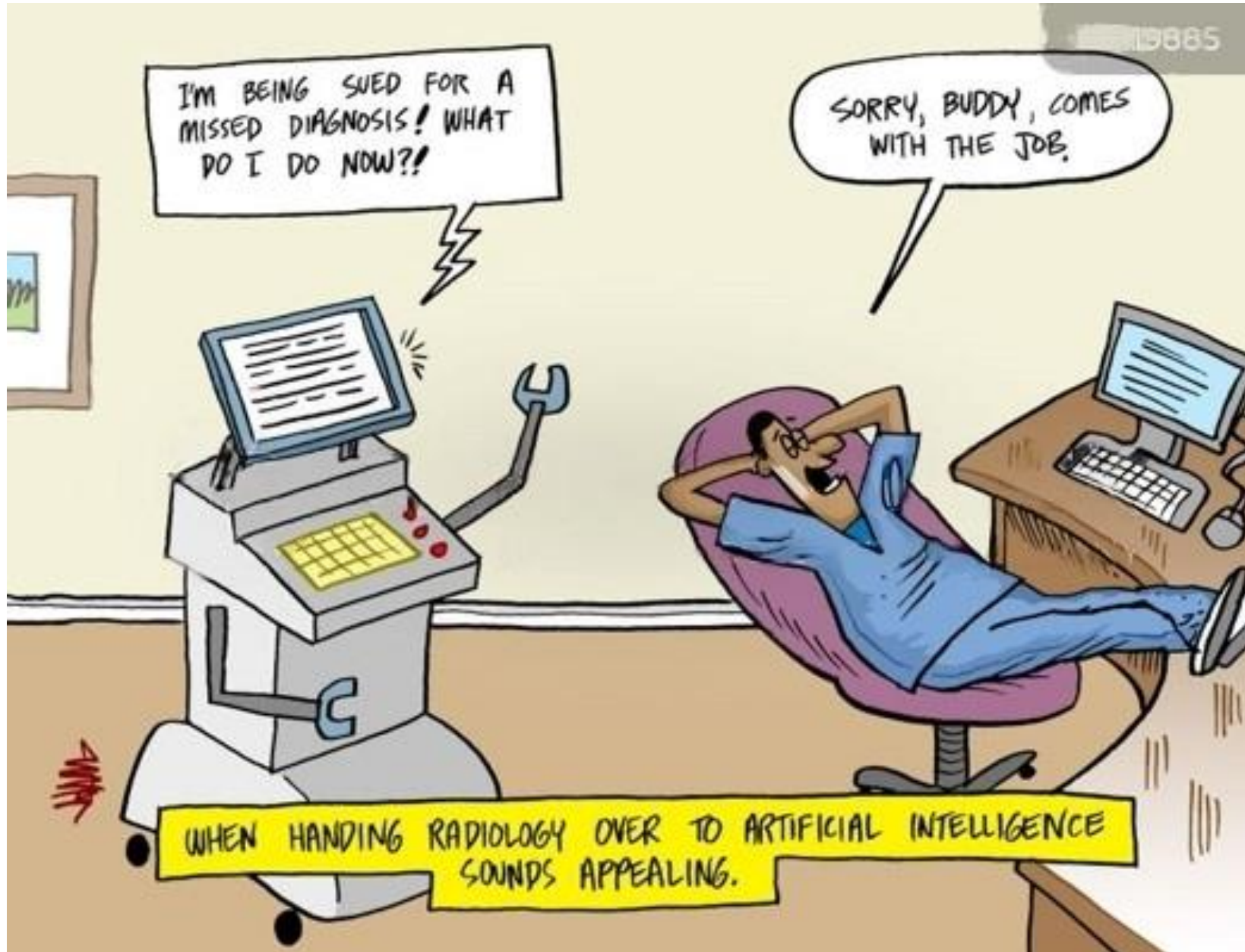
Identify Duplicate/Overlapping Concomitant Medication/Therapy records

Study specific conditions

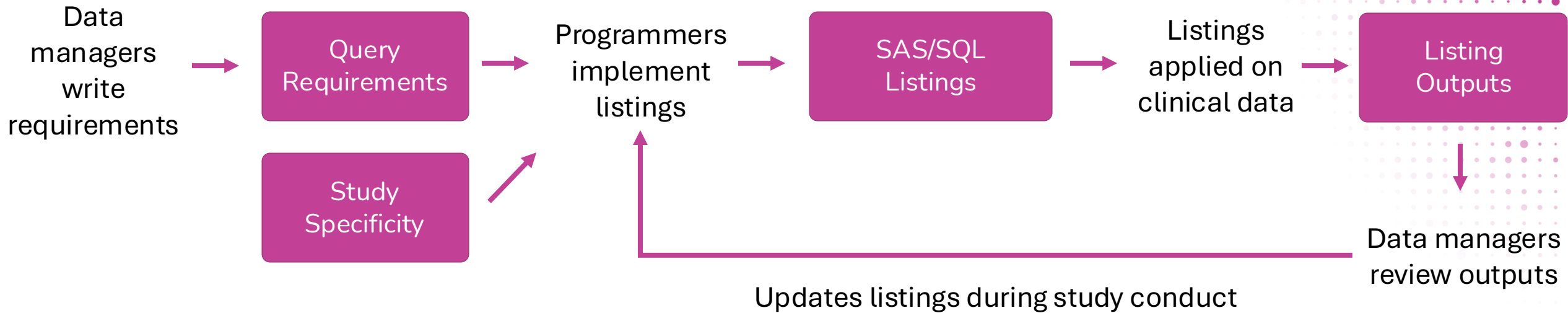
1. Except for such overlapping Concomitant Medications where the **end date** of 1st CM and **start date** of the overlapping CM are the **same day** and where the **frequency has changed**.
2. Except where the CM is a vaccination and the end date are missing, as long as the start dates are different.
3.

Model Benefit

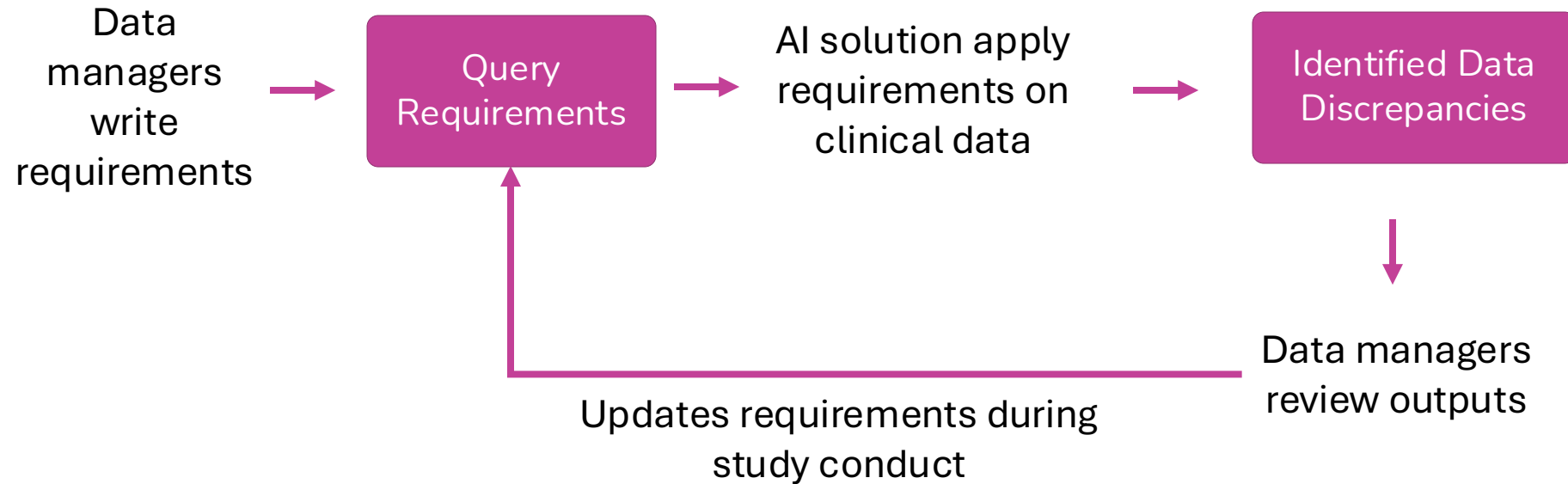
- **Reduced the burden** to program the scenario.
- Model is able to understand all the **relevant conditions** (indication, route, dosage, frequency, missing dates)
- Significantly reduce false positive rate (model achieves **95-100%** true positive rate)
- Improved **speed to review** query scenarios



SAS/SQL Listings: Current Process



AI Driven Detection of Data Queries: How?



- Query requirements are part of the implementation and are an input to the AI solution
- AI solution automatically adapts to study specificity (data structure, study protocol, ...)
- Link between requirements and quality of queries created
- Improve consistency in how requirements are applied and queries are raised

AI Driven Detection of Data Queries: Improve regulatory compliance

Traceability of Query Requirements

Keep track of what exact logic has been applied to clean data

Keep track on how requirements led to database change

Traceability of Model Decisions and Reasoning

Plain English text description of the models reasoning for each query scenario that can be accepted or rejected by a DM.

Used to for re-training.

Consistency on Raising Queries

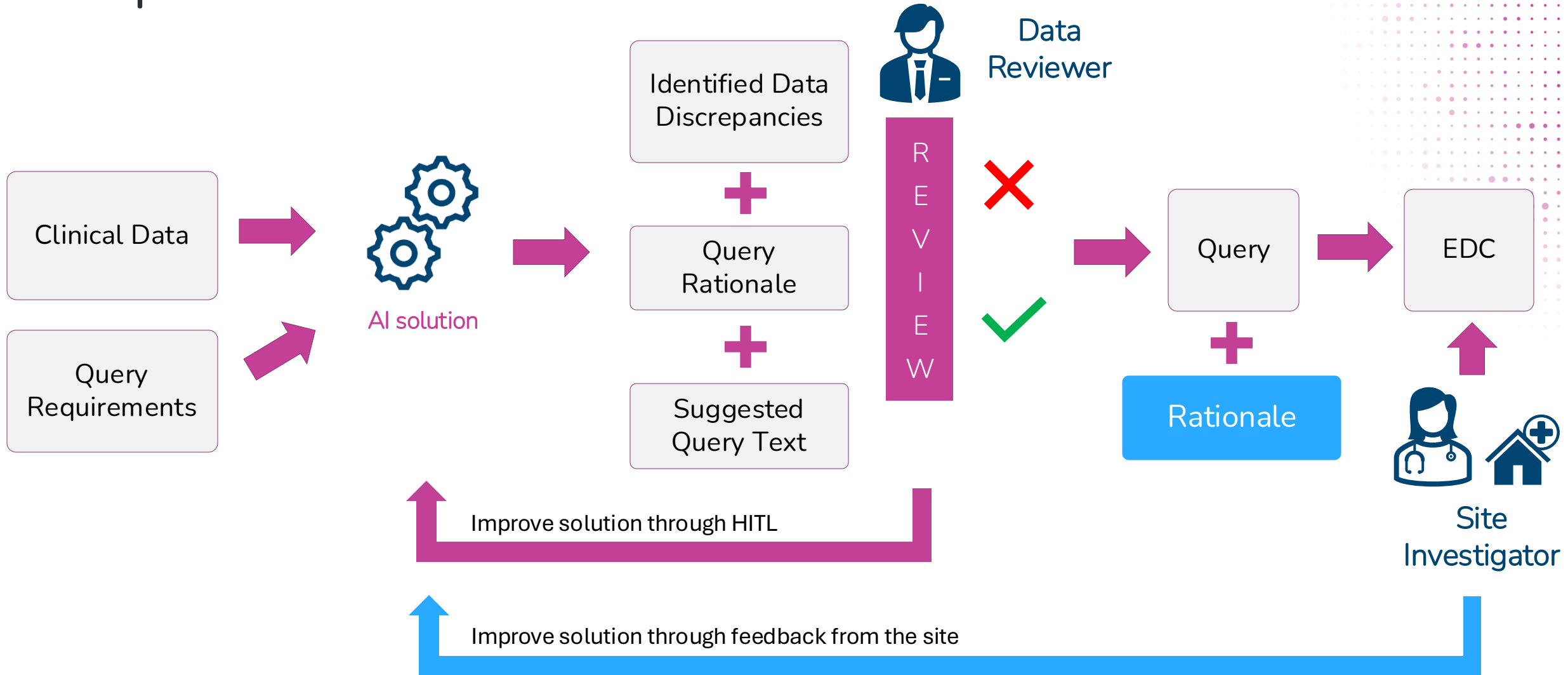
Thinking is done upfront instead of real time

Logic to raise queries is applied globally and comprehensively and not dependent on experience of DM



What could the
future be like?

Collect feedback from the site: Improve impact of queries created



Vision for the future

- Data managers can define **any** query requirements “**using plain English**” without programming expertise
- Raise queries in a uniform way across **datasources** and **studies** (EDC, ECOA, EPRO, Lab)
- Improves **consistency** of query text for site staff, improve % of queries that result in a data change, **reducing site burden**
- Solution suggests relevant **query requirements tailored** to a given study

Thank you

Any Questions?