

# From Solo to Sidekick: Using LLMs to Support SDTM Programming

25 November 2025

Corine Baljé-Volkers, Clin-Q

*Confidential. The information contained in this presentation may not be reproduced or otherwise disseminated without prior written approval of Clin-Q B.V.*

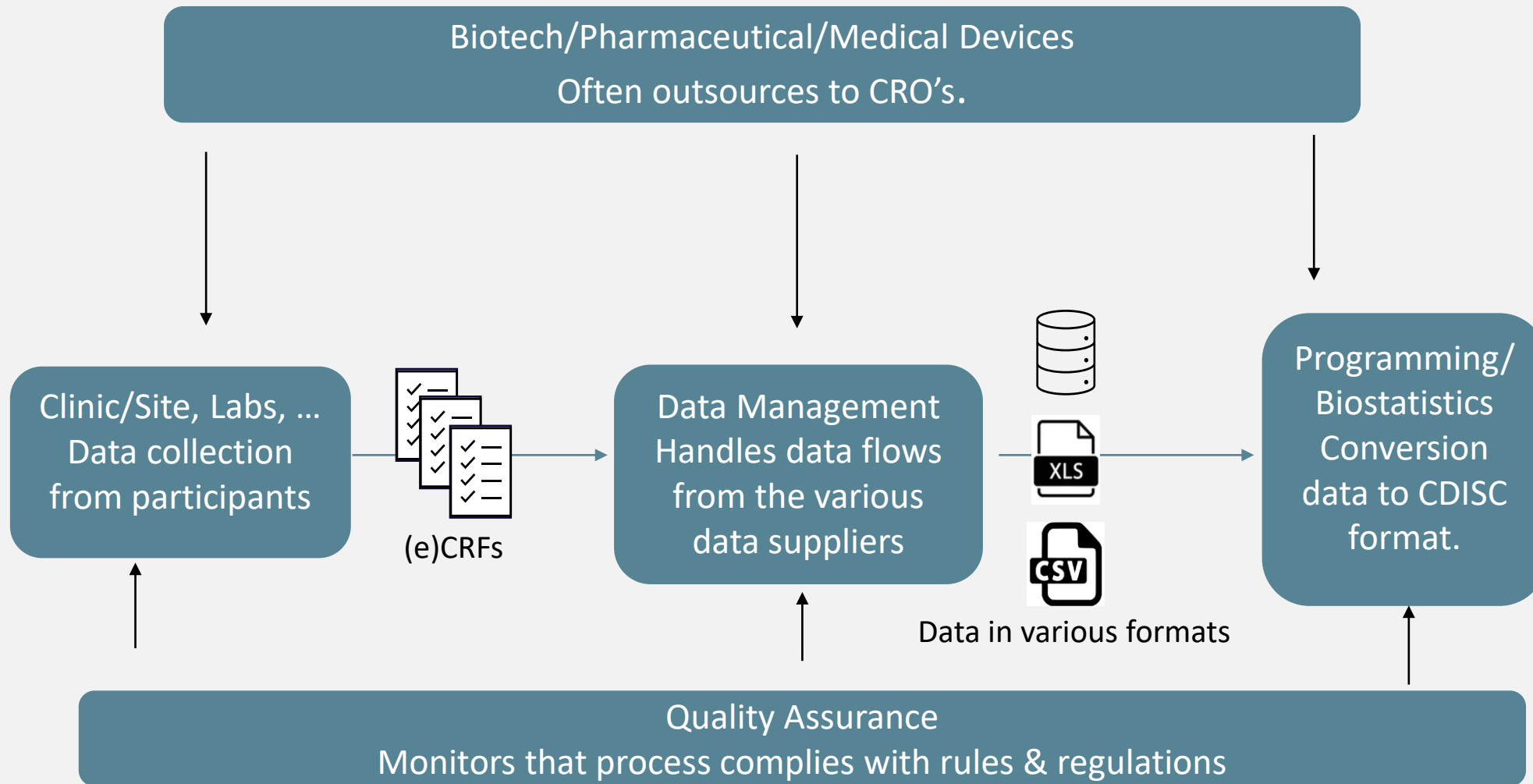
## Introduction

- 🔊 Strong industry interest in applying LLMs/AIs
- 🔊 Creation of CDISC is labour-intensive and rule-based: perfect candidate for automation
- 🔊 Phuse SDE 2024 (Utrecht): limitations to writing SAS code for SDTM (black box, using GPT-3.5)
- 🔊 LLM models usually better trained in Python or R
- 🔊 Small (niche-)CRO's receive data in a multitude of formats and conventions

➔ **Plan**: use LLM as a **sidekick programmer** to  
create SAS code for creation of SDTM datasets



## Step 0 – Process/project-partners



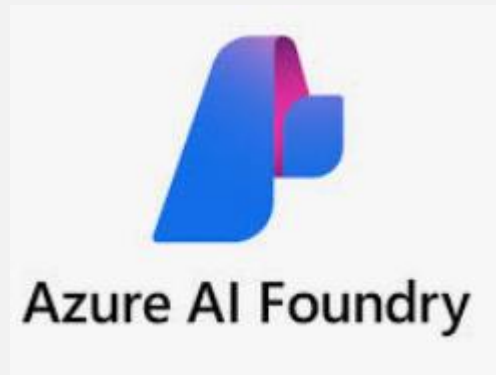
## *Step 0 – Process/project-partners*



## Step 1 – September 2024

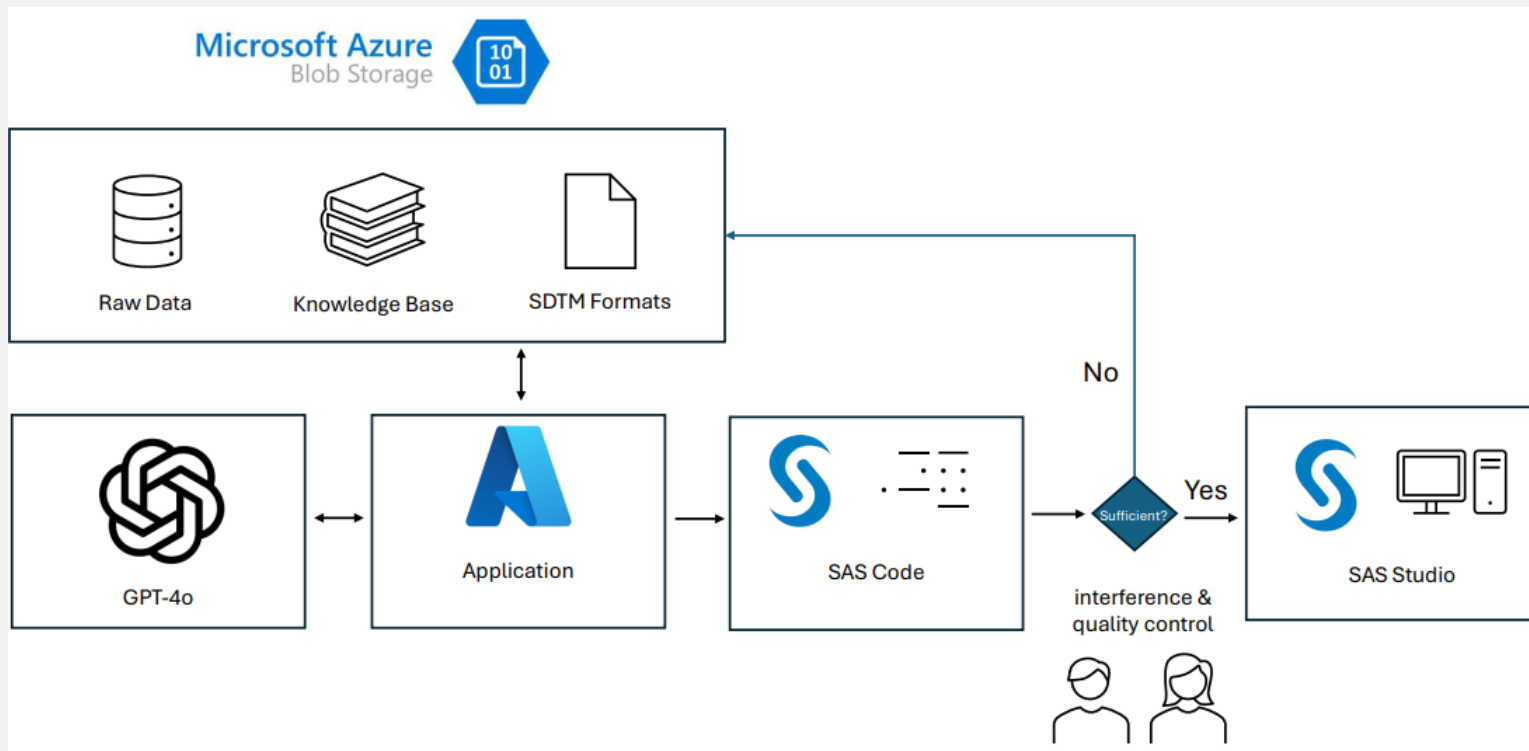
- ④ Work together with an IT student & supervisor (graduation assignment college university)
- ④ 3 demands:
  - ④ It should benefit the programmers (efficiency)
  - ④ It should be safe and secure (clinical data)
  - ④ It should be easy to integrate/implement
- ④ Local/self-hosted LLM (e.g. LLama) or Azure AI Foundry hosted LLM?

➔ Decided on




## Step 2 – June 2025

🔗 Proof of concept (using SAS studio)



## Step 2 – June 2025 (continued)

### Information provided to AI:


-  Raw synthetic SAS datasets provided: Randomization, Demographics, Exposure, Vital Signs, Adverse Events, End-of-Study

-  SDTMIG v3.4, SDTM v1.7, SDTMIG-MD v1.0

### AI tool behaviour/performance criteria:

-  identification of required and optional/permissible SDTM variables;

-  accurate mapping based on documentation;

-  incorporation of additional prompt-based instructions;

-  proper referencing of other datasets when needed;

-  testing of the generated SAS code in SAS (Studio)

## Step 2 – June 2025 (continued)

### COSTAR prompting (student)

- **Context:** I have a medical dataset in .TXT format that needs to be converted into an SDTM-compliant form. The data falls under the Demographics domain.
- **Objective:** I need functional SAS code that can process this dataset and convert it into SDTM format, ensuring it runs smoothly in SAS-Studio.
- **Style:** The code should be presented in a single, cohesive code block to facilitate easy copying. Include inline comments (using brackets) to explain each step of the process directly within the code.
- **Task:** Generate SAS code that:
  - Imports the dataset from .TXT;
  - Converts the data into SDTM-compliant format under the Demographics domain;
  - Outputs the final SDTM dataset.
- **Audience:** a programmer familiar with SAS but needing structured, annotated code for this specific conversion task.
- **Requirements:**
  - Ensure the code runs in SAS-Studio;
  - Refer to the relevant SDTM documentation that we have provided previously as separate documents (or general SDTM guidelines) to ensure compliance;
  - Present everything in a single code block with clear explanations.



## Step 2 – June 2025 (continued)

 Iteration results (SDTM-DM dataset creation only)

Iteration	#Variables	Variables	Source Dataset
1	6	STUDYID, DOMAIN, USUBJID, SUBJID, SITEID, SEX	DM
2	10	STUDYID, DOMAIN, USUBJID, SUBJID, SITEID, SEX AGE, AGEU, RACE, ETHNIC	DM
3	12	STUDYID, DOMAIN, USUBJID, SUBJID, SITEID, SEX AGE, AGEU, RACE, ETHNIC RFSTDTC, COUNTRY	EX, DM
4	15	STUDYID, DOMAIN, USUBJID, SUBJID, SITEID, SEX AGE, AGEU, RACE, ETHNIC RFSTDTC, COUNTRY ARM, ACTARM, DMDTC	EX, DM

➔ Demographics Knowledge Document

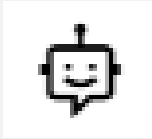
### *Step 3 – September-November 2025*

🔍 Setting up at Clin-Q and testing, following the manual created



🔍 In parallel: writing white paper and improving the Azure AI at the university college

🔍 Improvement: using **Agent** instead of Chat function/Assistant



## Step 4 – November 2025

The screenshot displays the 'Agents playground' interface. At the top, there's a navigation bar with a back arrow and the title 'Agents playground'. Below this is a toolbar with buttons: '+ New agent', '</> View code', 'Delete', 'Create trigger', and 'PREVIEW'. The main area shows a conversation thread. The thread header includes 'New thread', 'Thread logs', a dropdown menu, a grid icon, and the thread ID '30925t'. On the right, a link 'thread\_Nh48b4sdCuin/WftRefJHmNm' is visible. The conversation starts with a user message: 'What are you?'. The AI agent responds: 'I am a senior SAS programming assistant with expertise in transforming datasets into CDISC SDTM-compliant datasets, utilizing authoritative documents and company-specific standards for auditable SAS program generation.' Below the response, there are icons for a clock (1s), a document (1834t), and 'AI quality 2.5', along with a 'View Run Info' button. The user then asks: 'What is in your RAG'. The AI agent responds with a detailed list of RAG sources: 'My RAG sources include: 1. SDTM (Study Data Tabulation Model) and SDTMIG (SDTM Implementation Guide): These are authoritative CDISC documents defining SDTM domains, variables, controlled terminology, mapping rules, and derivation logic. 2. Internal Knowledgebase/Company Standards: Includes specific mapping corrections, corporate implementation guidelines, and lessons from prior errors (e.g., Knowledgebase\_Demographics\_Complete.docx). 3. Programming Guidance: Best practices for practical programming, workflow conventions, QC checks, and naming conventions (e.g., Relevant selection programming guidance.docx). If uploaded, I analyze contents of these documents rigorously for ensuring accurate SAS programming outputs and compliance. I also prioritize SDTMIG guidelines over other sources unless internal documentation explicitly overrides it.' At the bottom, there are icons for a clock (3s), a document (2017t), and 'AI quality 2.5', along with a 'View Run Info' button.

← Agents playground

+ New agent </> View code Delete Create trigger PREVIEW

New thread Thread logs 30925t thread\_Nh48b4sdCuin/WftRefJHmNm

What are you?

I am a senior SAS programming assistant with expertise in transforming datasets into CDISC SDTM-compliant datasets, utilizing authoritative documents and company-specific standards for auditable SAS program generation.

1s 1834t AI quality 2.5 View Run Info

What is in your RAG

My RAG sources include:

- SDTM (Study Data Tabulation Model) and SDTMIG (SDTM Implementation Guide):**
  - These are authoritative CDISC documents defining SDTM domains, variables, controlled terminology, mapping rules, and derivation logic.
- Internal Knowledgebase/Company Standards:**
  - Includes specific mapping corrections, corporate implementation guidelines, and lessons from prior errors (e.g., Knowledgebase\_Demographics\_Complete.docx).
- Programming Guidance:**
  - Best practices for practical programming, workflow conventions, QC checks, and naming conventions (e.g., Relevant selection programming guidance.docx).

If uploaded, I analyze contents of these documents rigorously for ensuring accurate SAS programming outputs and compliance. I also prioritize SDTMIG guidelines over other sources unless internal documentation explicitly overrides it.

3s 2017t AI quality 2.5 View Run Info

## Step 4 – November 2025

← Agents playground ▾

+ New agent </> View code 🗑️ Delete [Create trigger](#) [PREVIEW](#)

🕒 New thread 👤 Thread logs 📄 # 0t thread\_5tPBshKicnW6/QZLdW05xZ3h

🗨️

**Start chatting**  
Test your Agent by sending queries below. Then adjust your Agent setup to improve the Agent's responses.

from several datasets, make sure there is a proper merge with – as needed – renaming of variables, and drop variables that are not needed.  
Do not provide any comment lines in the code, nor any explanations. Do not create general code if you are not able to find the data or datasets. Only ask questions that are essential for creation of the code.  
Make sure you follow the programming guidelines as provided. I prefer standard SAS terminology and variable names to be in lowercase.

Messages in the Agents playground are visible to anyone with access to this resource and using the API. 🔗 💰 + 🗣️ Voice mode ➡️

## Step 4 – November 2025 (continued)

```

09
10 proc copy in=work out=work;
11 select dm;
12 run;

```

WARNING: IN= and OUT= are the same. Files will not be copied into themselves.

NOTE: PROCEDURE COPY used (Total process time):

```

real time    0.02 seconds
cpu time     0.03 seconds

```

Obs	SITSEQ	SITECODE	SUBJECTSEQ	SUBJECTID	EVENTNAME	EVENTDATE	SUBJECTFORMSEQ	ICYN	ICYNCD	ICDAT	
1	1	KIEM	1	KIEM-01	START	08/16/2024	1	Yes	Y	08/01/2023	
2	1	KIEM	2	KIEM-02	START	08/21/2024	1	Yes	Y	10/02/2023	
3	1	KIEM	3	KIEM-03	START	08/21/2024	1	Yes	Y	11/20/2023	
Obs	RACE	RACECD	ETHNIC	ETHNICCD	SEX	SEXCD	AGE				
1	White, Caucasian, or Arabic	3	Hispanic or Latino	1	Male	M	55				
2	Asian	1	Not Hispanic or Latino	2	Female	F	58				
3	White, Caucasian, or Arabic	3	Hispanic or Latino	1	Female	F	40				
Obs	AGEU	AGEUCD	CHILDPOT	CHILDPOTCD	MENOSTAT	MENOSTATCD	STUDYID	DOMAIN	USUBJID	SUBJID	
1	Years	1					TODO	DM	KIEM-01	KIEM-01	
2	Years	1	No	N	?1 Year post-menopausal	1	TODO	DM	KIEM-02	KIEM-02	
3	Years	1	Yes	Y			TODO	DM	KIEM-03	KIEM-03	
Obs	RFSTDTC	RFENDTC	SITEID	ARM	COUNTRY	RFICDTC	RANDNUM	TREATMENT	TRTCD	TRTCDN	ARMCODE
1		TODO	KIEM	Acti	TODO		RND01	Active	A	2	A
2		TODO	KIEM	Plac	TODO		RND02	Placebo	P	1	P
3		TODO	KIEM	Acti	TODO		RND03	Active	A	2	A

## Step 4 – November 2025 (continued)

The values of race, sex, ethnic, ageu should be in uppercase. Furthermore, the content of the variable race as provided is not in line with the permissible races according to CDISC documentation. Can you let me know how data is provided and what the permissible races are? If there are multiple races in one value line, please recode to 'OTHER'. When RACE in SDTM is OTHER, the specification is to populate the RACE variable with the value OTHER and to use a supplemental qualifier to store the specific race detail, typically in a variable like SUPP--(e.g., SUPPDM). Please adjust the code accordingly.



The raw values of the variable **RACE** include "White, Caucasian, or Arabic" and other combined races, which are not permissible according to CDISC controlled terminology. Permissible races in SDTM include:

- White
- Black or African American
- Asian
- American Indian or Alaska Native
- Native Hawaiian or Other Pacific Islander
- Other

Any record containing multiple races or undefined combinations (e.g., "White, Caucasian, or Arabic") should be recoded as **OTHER**. Additionally, for **OTHER** race entries, specific details should be stored in a supplemental qualifier (SUPPDM). This is in line with SDTM implementation guidelines <sup>1 2</sup>.

### Updated SAS Code with Supplemental Qualifiers

## Step 4 – November 2025 (continued)

Obs	SITESEQ	SITECODE	SUBJECTSEQ	SUBJECTID	EVENTNAME	EVENTDATE	SUBJECTFORMSEQ	ORIGINSUBJECTFORMSEQ	SOURCESUBJECTFORMSEQ
1	1	KIEM	1	KIEM-01	START	08/16/2024	1	1	
2	1	KIEM	2	KIEM-02	START	08/21/2024	1	1	
3	1	KIEM	3	KIEM-03	START	08/21/2024	1	1	

Obs	DESIGNVERSION	ICYN	ICYNCD	ICDAT	RACE	RACECD	ETHNIC	ETHNICCD	SEX
1	1.1	Yes	Y	08/01/2023	OTHER		3 HISPANIC OR LATINO		1 MALE
2	1.1	Yes	Y	10/02/2023	ASIAN		1 NOT HISPANIC OR LATINO		2 FEMALE
3	1.1	Yes	Y	11/20/2023	OTHER		3 HISPANIC OR LATINO		1 FEMALE

Obs	SEXCD	AGE	AGEU	AGEUCD	CHILDPOT	CHILDPOTCD	MENOSTAT	MENOSTATCD	STUDYID	DOMAIN	USUBJID
1	M	55	YEARS	1					TODO	DM	KIEM-01
2	F	58	YEARS	1	No	N	?1 Year post-menopausal	1	TODO	DM	KIEM-02
3	F	40	YEARS	1	Yes	Y			TODO	DM	KIEM-03

Obs	SUBJID	RFSTDTC	RFENDTC	SITEID	ARM	COUNTRY	QNAM	QVAL	RFICDTC	RANDBUM	TREATMENT	TRTCD	TRTCDN	ARMCODE
1	KIEM-01		TODO	KIEM	ACT1	TODO	RACE_DETAIL			RND01	Active	A	2	A
2	KIEM-02		TODO	KIEM	PLAC	TODO				RND02	Placebo	P	1	P
3	KIEM-03		TODO	KIEM	ACT1	TODO	RACE_DETAIL			RND03	Active	A	2	A

Obs	USUBJID	QNAM	IDVAR	IDVARVAL	QORIG	QEQAL
1	KIEM-01	RACE_DETAIL	USUBJID	KIEM-01	CRF	INVESTIGATOR
2	KIEM-02		USUBJID	KIEM-02	CRF	INVESTIGATOR
3	KIEM-03	RACE_DETAIL	USUBJID	KIEM-03	CRF	INVESTIGATOR

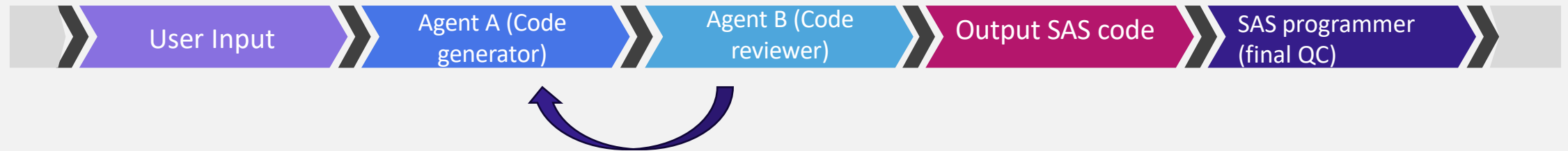
## Next steps –2026

### 🔍 Current process:

🔍 User Input → Agent A (Code Generator) → Output → SAS programmer (Code Reviewer)

### 🔍 Improving the SAS codes created:

- 🔍 Further training documentation
- 🔍 Extend to more SDTMs (including Trial Domains & SUPP)
- 🔍 Use ChatGPT-5 instead of -4o
- 🔍 Investigate how to use multi-agent systems within Azure AI Foundry



- 🔍 Expecting updated CDISC documentation (SDTMIG v4.0, SDTMIG-MD v2.0) better readable for AIs
- 🔍 Improve LLM development skills, starting with good prompt engineering
- 🔍 Integrate Pinnacle Validation?



Thank you for your attention!

