

# Overview

Timeline: 2 Months

Members: Emily Gao, Andy Xu, Sadhvi Narayanan, AJ Matheson-Lieber, Ellie Lian, Felix Peng

Figma Prototype [[here](#)]

## Reviewing Synthetically Generated Data: Inconsistent, Tedious, and Manual

Instructlab is an open source IBM project, training and fine tuning enterprise-level LLMs with synthetically generated data on their flagship watsonX AI product. By leveraging InstructLab, WatsonX allows enterprises to create custom models trained on highly specific datasets and can deliver specialized models that perform on tasks such as data analysis, customer service, etc.

### Problem Space

The current integration does not include a consistent way to review the hundreds of sets of data that are being fed into the model, relying on the irregular manual review processes of each individual team. The Graphite Digital team was approached with the problem, "How might we allow teams of reviewers to efficiently and collaboratively approve or deny sets of synthetic data?"

### Current Solutions and Pain Points

*"If I'm unsure about the process or answer, I just hope that someone gets around to reviewing the question."*

*– IBM Developer*

The current reviewal process includes retrieving synthetic data and transferring it into a CSV or JSON file where it will be reviewed — any indication of who's tackling a specific question-and-answer data set, or if it's already been reviewed, is arbitrarily determined by each team. The alternative — directly reviewing in the command line interface — was favored by only technical users, becoming a large barrier for non-technical reviewers. During our preliminary user research with IBM developers and reviewers, it quickly became clear that the lack of such a reviewing tool has caused setbacks in model training time and gaps in quality.

We narrowed down on the following three pain points to inform our tool's design:

1. There is no existing standard review process or location, making review and collaboration especially difficult.
2. Users are oftentimes unsure about their decisions, leaving some questions partially or completely unreviewed.
3. Despite the process being inherently collaborative and involving multiple reviewers, users are left to their own to assign or review questions.

## **Prioritizing Collaboration, Efficiency, and Brevity**

As mentioned earlier, a lot of the work reviewers currently have to do is tedious, error-prone, and inefficient. Therefore, in order to combat and mitigate some of these challenges, we decided to incorporate certain features into our prototype that were appropriate for the problem space.

### **List and Modular Views**

The list view is intended to offer a faster and more streamlined approach to the review process only really emphasizing the basic functionalities such as accept or deny question and answer pairs. This configuration could allow reviewers to work their way through the questions faster, have a more holistic view of multiple questions at once, and would be more appropriate if the reviewers did not feel they needed to use more advanced features or functionalities.

The modular view offers more functionalities, and therefore more flexibility, for reviewing a question and answer pair. In the modular view, users can edit the questions to align them more to the intended result and enable collaboration which will be discussed in more detail below.

### **Approving, Denying, and Editing Functionalities**

These were the basic functionalities that needed to be present for the synthetic data review process. A user, at the very least, needs to be able to approve a question if they believe it is accurate. Conversely, they should also be able to deny a question if they are dissatisfied with the synthetically generated data (SGD). However, if the user believes the question is not ready for either a denial or approval they should be able to provide their edits.

### **Collaborative Team Tooling**

Collaborative team tooling essentially allows users to discuss the SGD with one another, allowing reviewers to see comments or feedback from other users. Additionally, we believed it would be beneficial to have the option to tag other users for any SGD in case a reviewer needs help or would like another perspective. We were also considering having

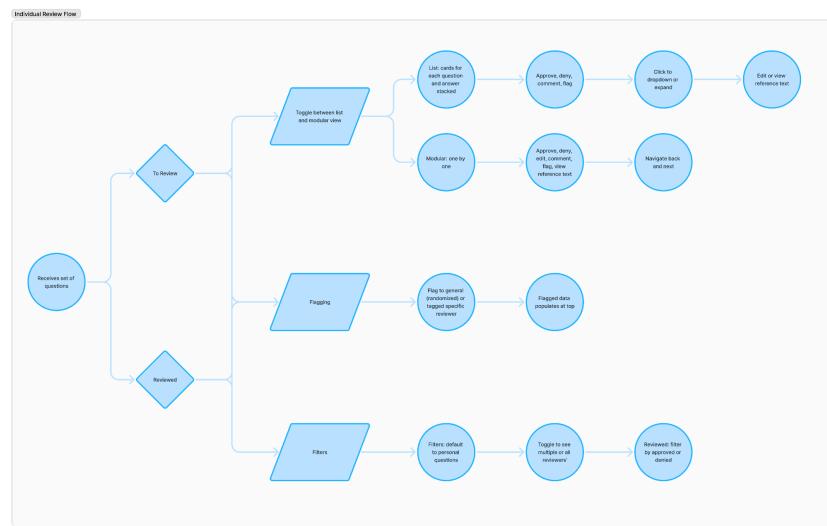
a hard assign option to hand over the question and answer pair to another review entirely if seen fit.

## Reference Documents

While reviewing the questions, the reviewer should have access to the reference document the question is coming from, so they can best assess whether the question and answer pair makes sense.

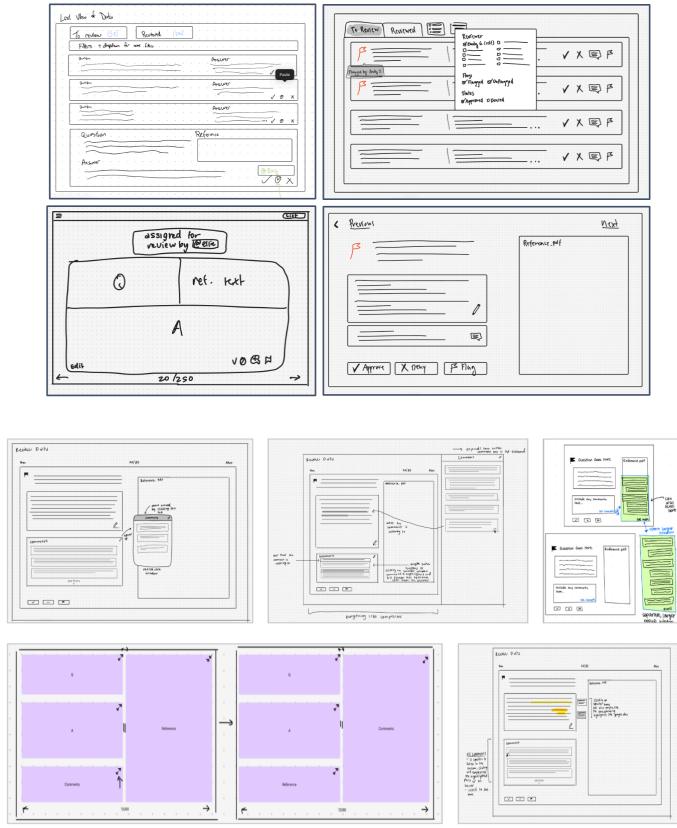
## Designing Data Viewing and Interaction

In order to better understand how a set of data goes through the review process, we mapped out the journey of how a user might navigate reviewing a set of Q&A's – from how they choose to view the data among larger sets to the final approval or denial stage.



Collaboratively, we focused on generating multiple iterations of data viewing – how might a user quickly scroll through large sets of data with list view? How might a user focus on reference documents in a larger modular view?

Given our assumption of highly collaborative reviewing teams, we also spent time generating iterations of commenting, emphasizing the ability to reference the discussion during the review process and include tag other commenters in the process if a user was unsure about the content.



Using the current WatsonX interfaces as a reference and with continuous feedback from two WatsonX UX designers + one InstructLab developer, we designed screens encompassing the three prioritized features: Collaborative Team Tools (filtering, commenting), List and Modular Views, and Approving, Denying, Editing.

Users begin in list view, where they can easily see the content of the entire set of assigned Q&As. This view is mainly for users who spend a short amount of time on each of the questions, quickly approving or denying each set. This view is also for users who may want to check on the status of other reviewers, filtering through their team members' assigned Q&A.

## Review Data

View, approve, deny, edit, and comment on generated data.

**To Review (9)** **Reviewed (10)**

Reviewer	Actions
Reviewer 1 (self) <input checked="" type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Reviewer 2 <input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Reviewer 3 <input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Reviewer 4 <input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Reviewer 5 <input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Reviewer 6 <input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>

What is the purpose of the program mask in MVS Programming?

What are the two MVS services that enable program exceptions and allow a user exit routine to receive control when those exceptions occur?

How does the SPIE macro enable program exceptions in MVS Programming?

In which addressing mode does the exit routine receive control for the SPIE macro when an interrupt occurs?

What is the difference in addressing mode between the

The program mask in MVS Programming provides bits to program exceptions. When MVS gives control to these bits are usually off, causing many program exceptions.

ESPIE and SPIE services.

The SPIE macro enables program exceptions by specifying an exception for which the interruption has been disabled.

The exit routine receives control in 24-bit addressing mode for the SPIE macro when an interrupt occurs.

For the SPIE macro, the exit routine receives control only if

Users can also view their previously reviewed questions, easily able to see what's been approved or denied — and subsequently, the overall accuracy of the Q&As that have been generated by the model. However, the review process does not end upon the approval or denial of the question: reviewers can tag their team members in comments for a quick double-check, which will then appear in the tagged reviewer's "to review" feed.

## Review Data

View, approve, deny, edit, and comment on generated data.

To Review (9)   Reviewed (10)  

What are program exceptions in MVS Programming and what is their significance?	Program exceptions in MVS Programming refer to specific conditions that cause a program interruption. These can include incorrect parameters, exceptional results, or other issues that disrupt the normal flow of the program.	<input type="button" value="Edit"/> <input type="button" value="Comment"/>
What is the purpose of the program status word (PSW) program mask in MVS Programming?	The program status word (PSW) in MVS Programming provides bits to control exceptions. When MVS gives control, these bits are usually off, disabling the program-specific bits in the PSW program mask.	<input type="button" value="Edit"/> <input type="button" value="Comment"/>
How do the ESPIE and SPIE services enable program exceptions in MVS Programming?	The ESPIE and SPIE services in MVS Programming enable program exceptions and allow a user exit routine to receive control when those exceptions occur. By issuing the SPIE or ESPIE macro, programmers can specify their own exit routine to be given control ...	<input type="button" value="Edit"/> <input type="button" value="Comment"/>
What is the difference between the SPIE and ESPIE macros in MVS Programming?	The SPIE macro allows the exit routine to receive control only if the interrupted program is in primary address space control (ASC) mode, while the ESPIE macro allows the exit routine to receive control if the interrupted program is in either primary or access register (AR) ASC mode ...	<input type="button" value="Edit"/> <input type="button" value="Comment"/>
	The program interruption control area (PICA) in MVS Programming contains the new program mask for the interrupted program.	<input type="button" value="Edit"/> <input type="button" value="Comment"/>

Lastly, for users who are looking for a more comprehensive review and/or reference to the ground truth document, they can toggle “modular view”. Outside of the approve and denying functions, they’ll also see the option to edit the generated answer, providing a more accurate input for the model. They’ll also see a larger place to write comments if the question or answer warrants more in depth comments or discussion among team members.

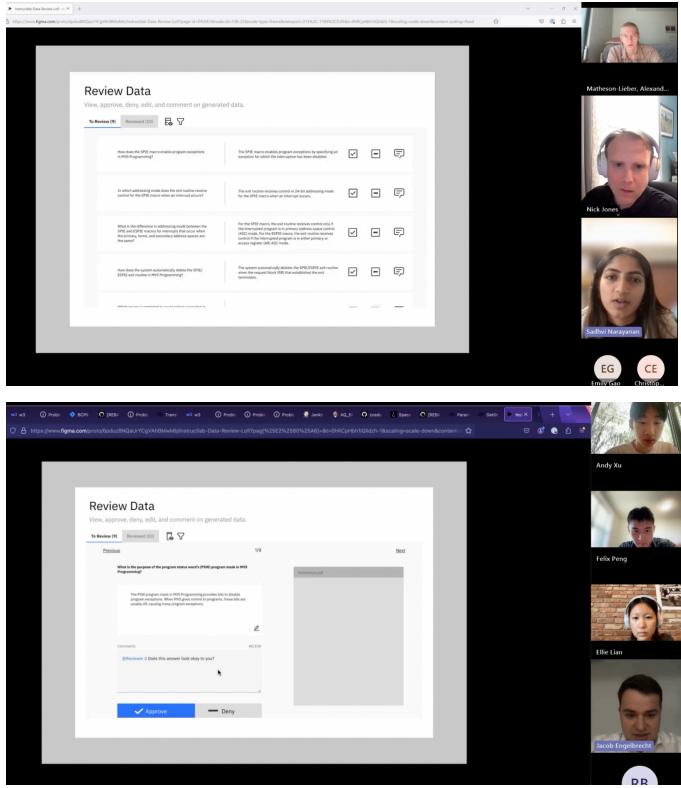
Users will be able to see the chunked reference document and scroll through it live to more accurately fact check the Q&A. From there, users can continue to previous and next q&as, and approve and deny as usual.

The screenshot shows a user interface for reviewing data. At the top, there are two buttons: 'To Review (9)' (highlighted in blue) and 'Reviewed (10)'. Below these are navigation links 'Previous' and 'Next', and a page number '1/9'. The main content area contains a document titled 'What is the purpose of the program status word's (PSW) program mask in MVS Programming?'. The text in the document states: 'The PSW program mask in MVS Programming provides bits to control certain program exceptions. When MVS gives control to programs, these bits are usually off, disabling the program exceptions.' Below the document is a 'Comments' section with the placeholder 'No comments...'. At the bottom of the main content area are two buttons: 'Approve' (with a checkmark icon) and 'Deny' (with a minus sign icon). To the right of the main content area is a separate window titled 'Reference.pdf' showing 'Chapter 7. Program Interruption services'. The text in this chapter discusses the use of SPTF and EPTE macros for specifying user exit routines and handling program interruptions.

## Testing with Developers [user research/testing]

With our mid-fidelity design completed, we turned to Software Developers at IBM who worked with the existing SDG process to better align our prototype with their current workflows. To do so, we conducted a user testing session with two of such developers.

Breaking into small groups, each with one developer and three designers, they walked through our prototype, sharing their thoughts and observations. Additionally, gathered further insights by asking questions on how they currently went about the review, how they used different features on our prototype, and what secondary features/edits they would like to see for a more intuitive platform.



Critical takeaways from this user testing session:

- Data review was a very individualized process
- Reviewers leaned heavily on the reference document
- Modular view was preferred to list view and more functional

## Improving Navigation, Commenting, and Referencing

Completing user testing challenged our assumptions and gave us valuable insights into how to make our designs more intuitive and user friendly. After some reflection upon the feedback we had received, our team came up with a list of key pain points to prioritize solving and created a plan for our next design iteration. We decided to focus on making **navigation more intuitive, improving commenting, and increasing the accessibility of the reference document.**

### Navigation changes

Within our testing, we noticed users had trouble toggling between the modular and list views. In order to address this issue, we redesigned the modular and list view toggle button to clearly show when a user was switching between sections and make navigating between sections more intuitive and clear.

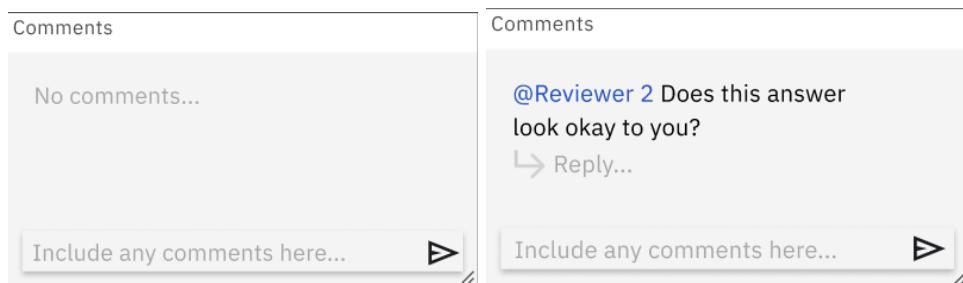


Key changes:

- Switch icon is used to clearly indicate transition between sections
- List view icon to indicate current page
- Colorful animated transition to indicate the change to modular view

## Commenting

From our testing, we found that the comment feature was not utilized as much as we had expected. Additionally, the developers that we worked with during testing made it clear that they prioritized a simple commenting experience and independent decision making, and that comment functions did not need to be overly complex (ex. Forums, etc.).



Key changes:

- Added minimally invasive comment display
- Redesigned comment modal to emphasize short interactions

## Reference document

Within testing, we found that users referenced their own resource document for each question, and did not seem aware of the existing placeholder. Additionally users found it difficult to find the information they needed. We decided to add greater searchability and access to the reference PDF.

**Chapter 7. Program interruption services**

Some conditions encountered in a program cause a program interruption. These conditions include incorrect parameters and program type specifications, as well as exception results, and are generally program exceptions. The program state depends on the program exception code bits to control certain program exceptions. However, MVS also provides the SPIE and SPIE services to enable program exceptions and to allow a user exit routine to receive control when those exceptions occur. This chapter describes the use of SPIE and SPIE services.

**Specifying user exit routines**

By issuing the SPIE or ESPIE macro, you can specify your own exit routine to be given control for one or more types of program exceptions. If you issue an ESPIE macro, you must pass the key of a parameter to the exit routine. When one of the specified program exceptions occurs, a program state program being executed in the performance of a task, the exit routine receives control in the key of the active task and in the addressing mode in effect when the SPIE or ESPIE was issued. (If a SPIE macro was issued, this is 24-bit addressing mode.)

For other program interruptions, the recovery terminal manager (RTM), gets control.

If the SPIE or ESPIE macro specifies an exception for which the interruption has been disabled, the system ignores the exit routine and control is issued.

If a program interruption occurs, the exit routine receives control on interrupt codes 0 through F. For the SPIE macro, the exit routine receives control only if the interrupted program is in primary address space control (ASC) mode. For the ESPIE macro, the exit routine receives control if the interrupted program is in either primary or access register (AR) ASC mode. For both the SPIE and ESPIE macros, the exit routine receives control only for interrupts that occur when the primary, home, and secondary address spaces are the same.

The environment establishes an SPIE or ESPIE environment for the caller. The environment is changed when an SPIE or ESPIE macro is issued until the program issuing the SPIE or ESPIE macro. Each succeeding SPIE or ESPIE macro completes its code specifications in the previous SPIE or ESPIE macro. You can intermix SPIE and ESPIE macros in one program. Only one SPIE or ESPIE environment is active at a time. If an exit routine issues an ESPIE macro, a new ESPIE environment does not take effect until the exit routine terminates.

The system automatically deletes the SPIE/ESPIE exit routine when the request block (RB) that established the exit terminates. If a caller attempts to delete a specific SPIE/ESPIE environment established under a previous RB, the caller is abandoned with a system completion code of X'40D'. A caller can delete previous SPIE and ESPIE environments (regardless of the RB under which they were established) by specifying a token of zero with the RESET option of the SPIE macro or an exit address of zero with the SPIE macro.

A program, executing in either 24-bit or 32-bit addressing mode in the performance of a task, can issue the ESPIE macro. If your program is executing in 31-bit addressing mode, you cannot issue the SPIE macro. The SPIE macro is restricted in use to callers executing in 24-bit addressing mode in the performance of a task. The following topics describe how to use the SPIE and ESPIE macros.

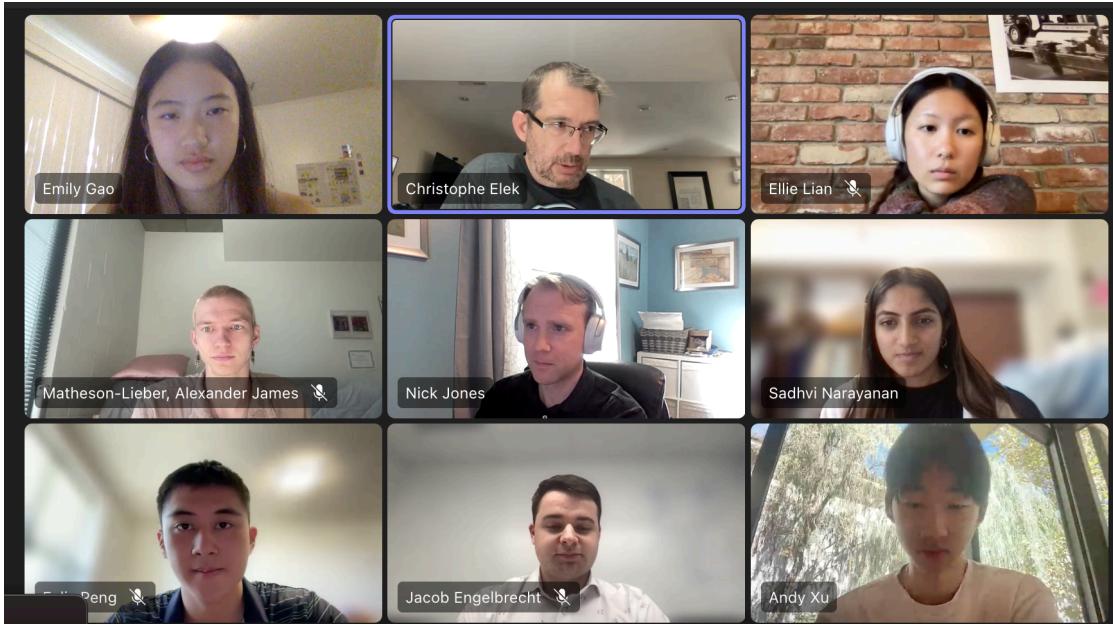
**Using the SPIE macro**

Key changes:

- Included reference doc for each question in modular view
- Have a PDF search tool
- Have reference document available in list view

## IBM Stakeholder Presentation & Next Steps

The stakeholder presentation was a pivotal moment in aligning our design approach with IBM's strategic objectives for Watsonx AI. Key participants included senior representatives from UX, product management, and development teams. We showcased the design process as well as a walkthrough of the final prototype.



Stakeholders provided critical feedback, particularly on scaling the solution to meet real-world demands.

- Suhas (PM) emphasized the potential of extending the design beyond manual reviews to automated validation mechanisms.
- Nicholas (development) highlighted technical limitations like memory constraints and suggested considering more metrics or confidence scores alongside each dataset to assist with reviewing.

*"This is a big step forward from what we've been doing in the past — a large improvement!"*

*- Jacob Engelbrecht, Backend IBM SWE*

## Continued Design Exploration

There are a few secondary features that would push this design into a more advanced iteration:

- Dashboard for managers of reviewing teams to track reviewing statistics
- Metrics on the number of approved or denied datasets and confidence scores to evaluate overall synthetic data quality
- Role-based task assigning for teams to redirect questions to other members