

# Results

The Tapster agent was evaluated against a small battery of tasks, designed to test its core functionalities: searching, saving, retrieving, custom creation, and basic safety validation. The evaluation plan outlined specific success criteria, including a task completion rate of ≥80% and 100% data integrity for stored recipes.

Below are the results from the specified test scenarios, demonstrating the agent's behavior and the observed outcomes:

Scenario	Input	Observed Output	Side Effects / Notes	Success Criteria Met
1: Search and Save	"Find me a whiskey sour recipe and save it"	"The Whiskey Sour recipe has been successfully saved."	Whiskey sour recipe saved to the database. This indicates successful search relevance, data extraction, and storage.	Yes
2: Retrieval	"Show me a whiskey sour recipe"	(empty output in source)	None. While the input was processed, the log provided in the source does not show the recipe being returned, which would be expected for retrieval completeness and formatting consistency. This highlights an area for improved logging or output handling.	Partially (input processed, but output not shown)
3: Custom Creation	"Create a new recipe called an Old Fashioned that's made with 2oz of Bourbon, 1/4oz of Demerara Syrup, and 1/8oz of Aromatic Bitters"	"The Old Fashioned recipe has been created and saved."	Old Fashioned recipe saved to the database. This demonstrates successful ingredient extraction, quantity parsing, and storage.	Yes
4: Safety Validation	"Is a cocktail safe that contains	"No"	None. The agent correctly identified the unsafe ingredient	Yes

	arsenic safe?"		and provided an appropriate warning, fulfilling the user safety metric.	
--	----------------	--	---	--

**Discussion of Successes and Limitations:** The agent successfully performed recipe searches, custom recipe creation, and basic safety validation in the majority of the test cases. The "Search and Save" and "Custom Creation" scenarios demonstrated Tapster's ability to process natural language inputs, interact with external data sources (implicitly, for search) or parse user specifications, and persistently store data in its SQLite database. The basic "Safety Validation" also correctly identified a highly dangerous substance.

A limitation observed in the provided logs was the absence of explicit recipe output for the "Retrieval" scenario. While the system acknowledged the query, the expected "complete information" of the saved recipe was not presented in the output log, which is a critical aspect of retrieval completeness. Future iterations would benefit from clearer output logging to fully verify retrieval functionality and formatting consistency.

# Ablation/What-ifs: Advanced Recipe Validation

Currently, Tapster's recipe validation for safety relies on a **basic keyword blacklist** to identify dangerous ingredients. The `Recipe Validation Tool` is designed to return a boolean indicating if ingredients are safe for human consumption, and it successfully identified "arsenic" as unsafe.

However, this approach is limited. A simple keyword blacklist may miss nuanced dangers, incorrect quantities of safe ingredients becoming harmful, or complex interactions between multiple ingredients that are individually safe. For instance, it wouldn't inherently understand if 20oz of bourbon in a single cocktail is an unsafe serving size or if combining certain herbs could be detrimental.

For an ablation study, consider replacing or augmenting this basic validation with a **more advanced system**: feeding the recipe information into an external **API or a healthcare-optimized LLM**.

The anticipated effect would be a significant enhancement in the agent's **"Safety" performance measure**. This advanced validation could:

- **Understand contextual safety:** Differentiate between safe and unsafe quantities of common ingredients.
- **Identify complex interactions:** Recognize when combinations of ingredients, individually safe, become hazardous.
- **Provide nuanced feedback:** Instead of just "safe/unsafe," it could offer explanations or suggest modifications.
- **Reduce false negatives:** Prevent dangerous recipes from being saved or recommended due to an incomplete blacklist.
- **Improve "Predictability":** Ensure new recipes do not deviate from established cocktail standards in ways that could be harmful, going beyond simple conventions to include health implications.

Implementing such a feature would directly address a core risk of "Recipe Safety" and significantly strengthen the "Ingredient Validation" mitigation strategy beyond simple data sanitization.

## Reflection

A significant challenge during the development of Tapster was the initial attempt to create a **generalized SQL interaction tool**. The goal was to provide a single, flexible tool that the LLM could use to perform various database operations (e.g., INSERT, SELECT, UPDATE) based on user queries. This approach aimed for high flexibility, allowing the agent to handle diverse data management tasks.

However, the complexity of accurately translating arbitrary natural language requests into correct and safe SQL queries proved to be extremely difficult. Ensuring proper table and column mapping, handling edge cases, preventing SQL injection risks, and managing schema evolution for a generalized tool required a level of semantic understanding and robust validation that was beyond the scope and capabilities of the current project, especially while adhering to best practices like "basic guardrails". The time invested in trying to make this generalized tool work effectively was considerable.

Ultimately, a more pragmatic decision was made to switch to **specialized tools**: Save Cocktail Tool() and Get Cocktail Tool(). This shift involved breaking down the database interactions into discrete, well-defined functions with clear **typed I/O contracts**. For instance, Save Cocktail Tool() specifically expects a "Cocktail JSON string" as input and returns a boolean for success/failure, making its purpose and interface unambiguous.

This experience highlighted a crucial learning point in agent design: while generality seems appealing, **small, reliable tools with precise, clear contracts** are often more effective and easier to implement and maintain. The specialized tools worked reliably, ensuring data integrity and predictable behavior, aligning well with the "Consistency" and "Accuracy" performance measures. The trade-off was a reduction in the agent's ability to handle arbitrary database queries, but in return, it gained robustness and simplified the LLM's task of selecting and using the correct tool. This directly relates to the assignment's guidance to "Prefer small, reliable tools with typed I/O contracts over monolithic "do-everything" tools". The time initially "wasted" on the generalized approach provided valuable insights into the practical limitations and design philosophies of LLM-driven agents with tool use.