# EXAFLUENCE
## Data Driven Influence

# ExfSurveyBuddy - GenAI Powered Platform for Survey Creation, Administration & Analytics

**Malaya Rout**
February 2024

## Introduction

In today's world, conducting a survey for a specific objective is marred with many challenges. Designing a questionnaire that brings out the most honest response, selecting the right respondents, and analysing the survey responses are just a few of them. Leave aside the number of months it takes to complete one survey cycle and plan the improvements to be incorporated in the next survey design. The less we talk about how we feel the moment we see a set of thirty survey questions to be answered, the better it is. Four of the thirty questions would need us to respond in a descriptive manner.

## What was the motivation behind developing the Exafluence Survey Buddy Platform?

At Exafluence, we have built ExfSurveyBuddy, a companion for survey design and analytics. It is a Generative AI-based platform.

Going through an end-to-end survey cycle is entirely manual and time-consuming. The various phases involved are 1) Designing the questionnaire, 2) Reaching out to respondents, 3) Getting their responses, and 4) Analysing the responses. The failure points in each phase could be multiple.

### Designing questionnaire
- Are we asking the right questions?
- Are we covering all aspects of the requirement?
- Are the questions interesting and engaging enough to evoke a response?

### Reaching out to respondents

- Have we reached out to a sample of respondents' representative of the organisation's total base?

### Getting responses

- The response rate is meagre in most cases (less than 10%)
- Often, it takes months to get a response

### Analysing responses

- Going through the responses (both ratings and free text comments) manually takes time and leads to human errors

The ExfSurveyBuddy addresses all the failure points above by using Large Language Models (LLMs), the backbone of Generative AI. The platform can be used to get customer feedback on products and services, understand respondents' behaviour, determine respondents' health status, carry out market research, and any other usage of surveys that we can think of. The user's flexibility in entering a specific context for the survey makes it possible to go as specific and broad as possible.

## What are the capabilities of ExfSurveyBuddy?

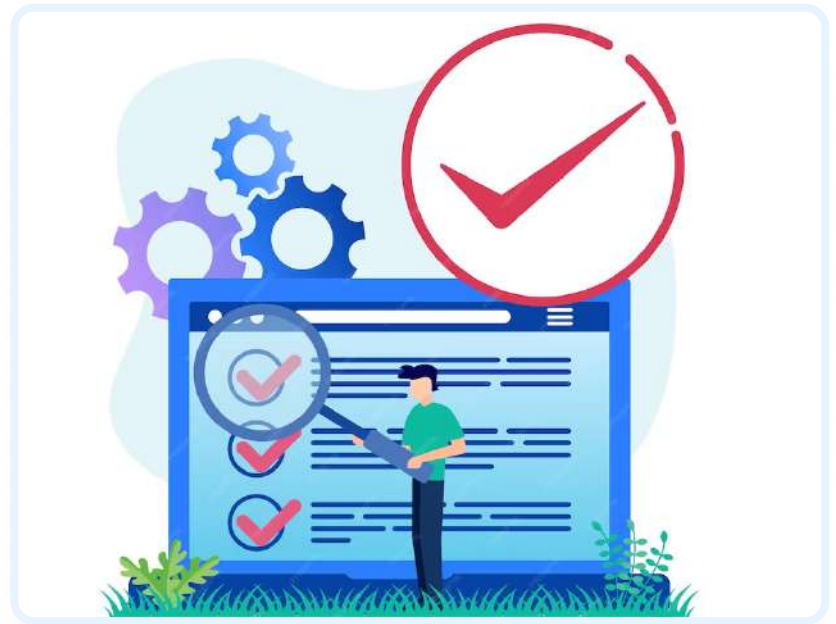The following is a list of capabilities of the platform

1. Generate questions with interesting-ness and likelihood of response scores
2. Generate synthetic respondents based on sample respondents
3. Generate responses based on questions and respondent profiles
4. Generate the top three reasons justifying the response
5. Generate insights from responses
6. Set up configuration properties of LLM for
   - Generating questions
   - Generating Python code used to create synthetic respondents
   - Generating responses
   - Generating Python code used to generate insights
7. Context setting by the user to generate questions
8. Upload a file containing sample respondents in CSV or XLS format
9. Upload a file containing questions and another containing synthetic respondents
10. Upload a file containing responses

11. Generate Python codes for
   - Creating synthetic respondents
   - Analysing categorical variables
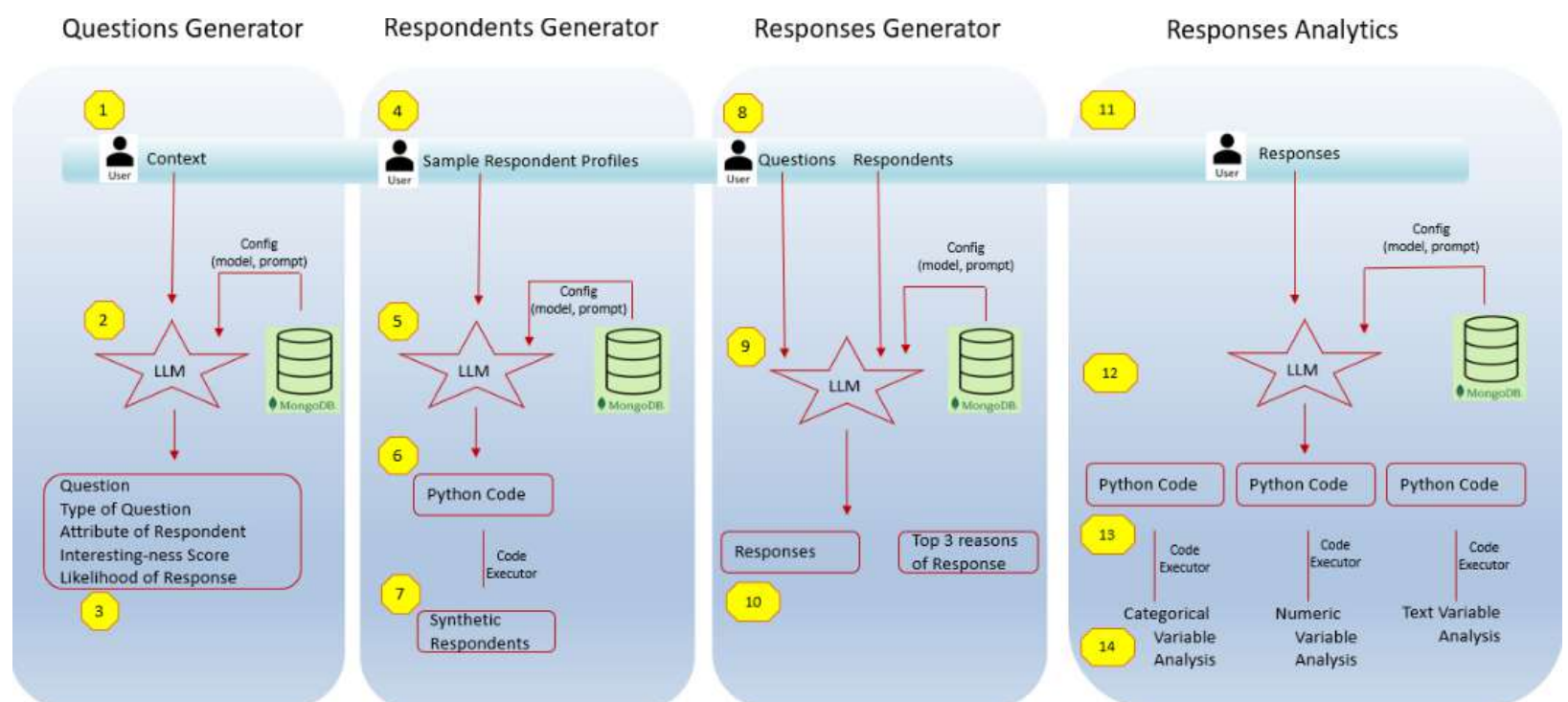   - Analysing numeric variables
   - Analysing textual variables
12. Download the following files as CSV
   - Questions
   - Synthetic Respondents
   - Responses

## Which architecture and technology stack does it follow?

The platform has four modules. The following is an architecture diagram of the platform.



## Would you mind giving us a brief description of each of the four modules?

### Generation of Questions (module 1)

The user gets to enter a brief context in which the LLM will generate the survey questions. The user also controls the number of questions, the scale of responses, and the number of free-text questions to be generated. The questions can be of three types.

- **NPS (Net Promoter Score):** There can be a maximum of one question related to the NPS calculation

- **Free Text:** Certain questions require the user to enter text to express his response in a descriptive manner

- **Regular:** All other questions are categorised as regular questions, providing a scale for rating

The output from the LLM is extracted in a structured JSON format as key/value pairs. It is converted into a Pandas data frame before displaying it as a table on the user interface. Each question is limited to 30 words. The LLM also identifies which attribute of the respondent the question reflects. The LLM also assigns an interestingness score out of 100 to each question. It also gives a likelihood of response to each question based on its historical knowledge. The questions are finally arranged in a decreasing order of interestingness score. This should allow the user to focus on questions towards the top and ignore questions at the bottom if necessary. The user can download the list of questions as a CSV file.

**Factors considered by the LLM for arriving at an interestingness score are as follows.**
- Relevance to Context
- Engagement Potential
- Intrusiveness
- Creativity and Variety
- Appeal to Common Experiences

**Factors considered for arriving at a likelihood of response are as follows.**
- Positive Tone
- Flexibility
- Relatability
- Everyday Relevance

These factors are weighted and combined to assign scores to each question. The goal is to craft questions that are not only relevant to the context but also intriguing, relatable, and designed to evoke a response from respondents.

## Generation of Synthetic Respondents (module 2)

The user uploads a sample file containing profiles of a few real-life respondents. The user can view the contents of the uploaded file. The user also decides and enters the number of respondents the LLM should create synthetically. The LLM generates

the Python code for creating the synthetic respondents, whereas the actual generation happens outside the LLM by executing the Python code. This approach helps keep the tokens exchanged with the LLM to a minimum. There is no upper limit to the number of respondents generated. The respondents can be downloaded as a CSV file.

The LLM generates a user-defined Python function having a pre-defined name. The input to the Python function is a Pandas data frame, as is the output. The synthetic respondents are created within certain statistical constraints. For every categorical variable, the frequency distribution of values in the sample file matches the frequency distribution of values in the synthetic dataset. The generated numeric values are all positive. For every numeric variable, the mean, median, mode, and standard deviation in the synthetic dataset match the mean, median, mode, and standard deviation in the sample file. The synthetic dataset has the exact columns (or variables) as the sample dataset. A respondent ID is generated against each profile. The respondent ID is the first column. The rest of the columns of the synthetic dataset are in precisely the same order as the sample dataset.

## Generation of Responses (module 3)

The user uploads the questions downloaded from the first module, and the respondents downloaded from the second module. They can view the contents of both the files uploaded. The LLM generates the responses corresponding to each question and each respondent profile. The response contains a rating / free text. The response also includes the top three reasons why the LLM thinks the particular rating or textual response is the most likely rating or response, respectively. The user can download the responses as a CSV file.

## Generation of Insights (module 4)

The user uploads the file containing the responses downloaded from the third module. They can view the contents of the uploaded file. There are two kinds of analysis. One is in a tabular form, and the other is in a graphical format. Each of these forms of analysis is carried out and displayed for categorical, numeric, and textual variables. Following is a matrix of the various types of analysis the platform provides.

| Type of Variable | Tabular Analysis | Graphical Analysis |
|---|---|---|
| Categorical Variable | Average rating against each class | Multiple box plot |
| Numeric Variable | Correlation coefficient | Scatter plot with regression line |
| Textual Variable – Feedback | Sentiment determination of each free-text response<br>Frequent patterns in positive sentiments<br>Frequent patterns in negative sentiment | Bar chart for positive and negative respondents |
| Textual Variable – Reasons of Rating | Topmost Reason – Summary of all Respondents<br>Second Reason – Summary of all respondents<br>Third Reason – Summary of all respondents | Bar chart for reasons corresponding to positive and negative respondents |
| Calculation of NPS | % promoters minus % detractors<br>Promoter – Rates 9 or 10 to NPS question<br>Detractor – Rates 1 to 6 to NPS question<br>Neutral – Rates 7 or 8 to NPS question | Bar chart for promoters and detractors |

If you would like to learn more, write to us at marketing@exafluence.com

For interesting videos about our solutions subscribe to our YouTube channel-
https://tinyurl.com/YouTubeExf

For regular updates on our solutions follow us on LinkedIn
https://tinyurl.com/LinkedInExf