# Data Science Clinic – International Rescue Committee (IRC)

IRC aims to address the educational challenges faced by crisis-affected communities through their AI-driven chatbot, aprendIA. They have developed the AudioClass System (ACS) to support Venezuelan migrants in Colombia who are at risk of dropping out of school. To help IRC continue to improve educational access, the team analyzed 349,000 conversational messages to provide product analytics on ACS performance and the user experience. The team also developed a pipeline that cleans and processes conversational text data to assist IRC's future product development.

The team analyzed product usage patterns in chatbot conversations, finding that the activity volume increased after July 2023 and users primarily engaged with the chatbot in the morning. More than half of the conversations were under 45 minutes, with an average duration of 30 minutes. Identifying these activity patterns can help IRC scale resources according to high-demand periods. The team also assessed user engagement by measuring setup time and content engagement duration. The findings revealed that approximately 70% of the students were engaged with the educational content. However, the accuracy of the data may be misleading as missing messages prevented correct identification of the beginning and end of some conversations.

The team identified three areas of improvement for the chatbot. Firstly, the chatbot inconsistently detects user disconnections, with detection times spanning from 10 minutes to one hour. Enhancing the chatbot's ability to detect disconnections will allow for more targeted messages encouraging engagement. Secondly, the team investigated the instances where the chatbot failed to understand user inputs and revealed that misunderstanding caused some students to stop engaging with the chatbot. Thus, improving the chatbot's language parsing ability will help increase user interest and engagement. Finally, the team also proposed improvements to IRC's data management in response to difficulties with missing, mislabeled, and testing data.