# Data Collection Documentation (Draft)

Steven Deutekom

June 2019

# Contents

# 1 Introduction

The goal of this project is to collect source code samples from individual authors from various online sources. It requires finding samples that have certain sociolinguistic characteristics such as gender, region, and experience. In addition to sociolinguistic characteristics, it is necessary that the samples are the work of only a single author. Data that meets these needs is being collected and used to create a dataset with information on authors and their source code.

The collected data will be used for sociolinguistic research into how people use programming languages. Current and future University of Lethbridge students will use this research to learn how sociolinguistic characteristics affect how programmers write code. Previous research was conducted with a small dataset of student programs. This new dataset will contain many samples from a larger set of programmers.

The document is broken into three main sections. Each section details one of the sources that was used to collect data. First, the collection methods used to gather source code from GitHub. Then, the methods used to collect source code from Codeforces. Lastly, the methods that were used to add gender data to the samples collected.

Each section introduces the sources and collection methods. Then it discusses the pros and cons of the source. Next an overview of the process of collecting data is given, with diagrams to help visualize it. Following this a more in depth technical examination of the collection process takes place. Finally, some reflection on the source is given (needed?).

# 2 Collecting Source Code From GitHub

# 3 Collecting Source Code From Codeforces

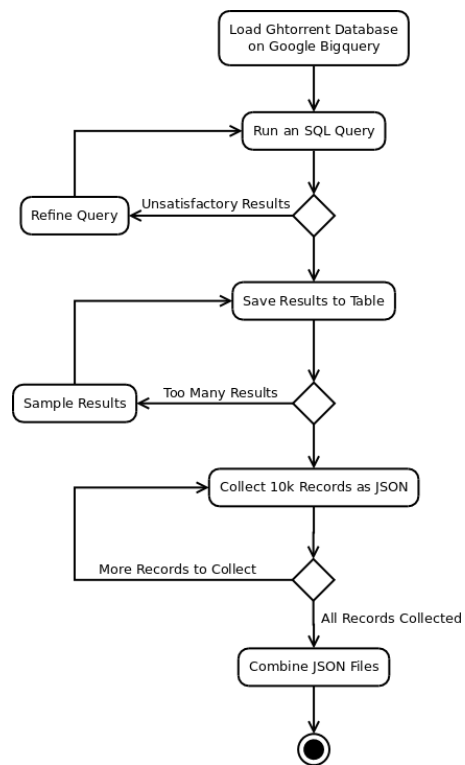# 4 Adding Gender Labels to Authors
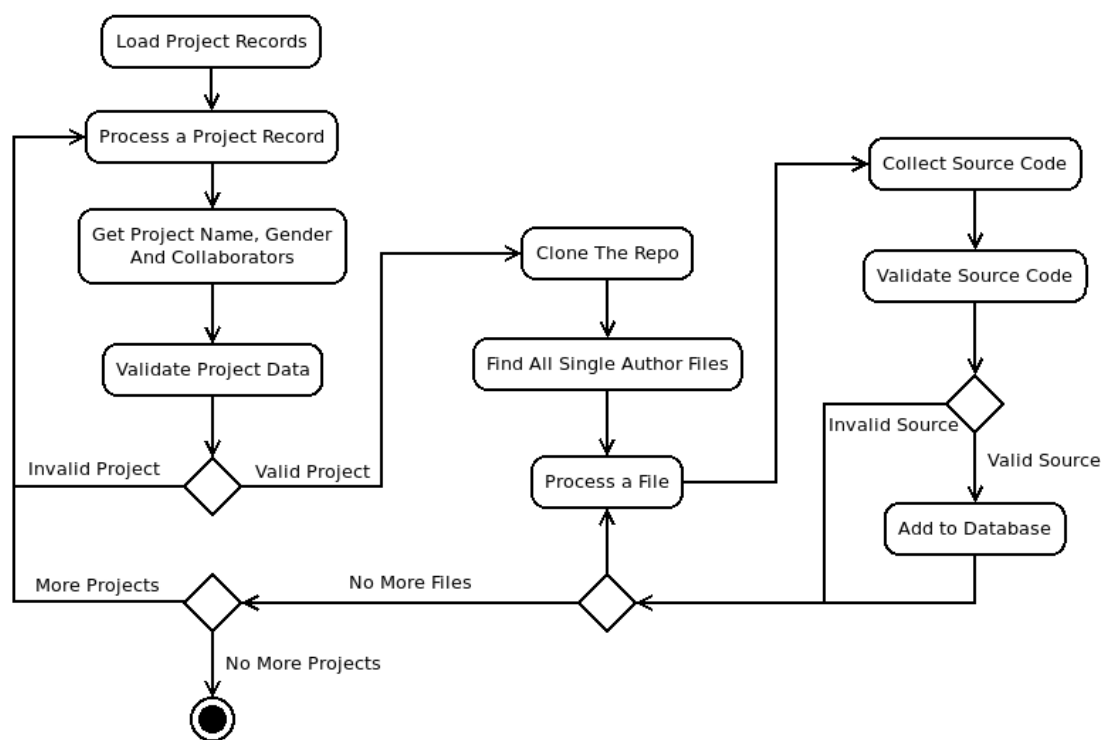
# 5 Conclusion

Figure 1: Getting Data From Ghtorrent

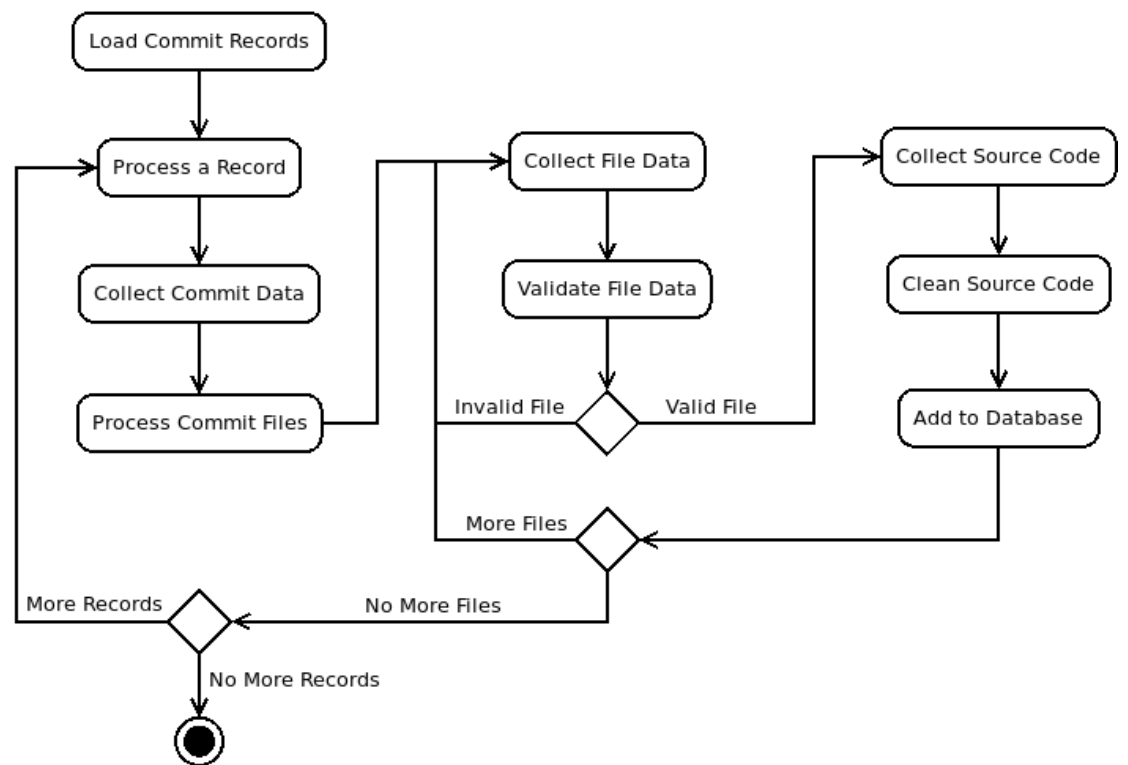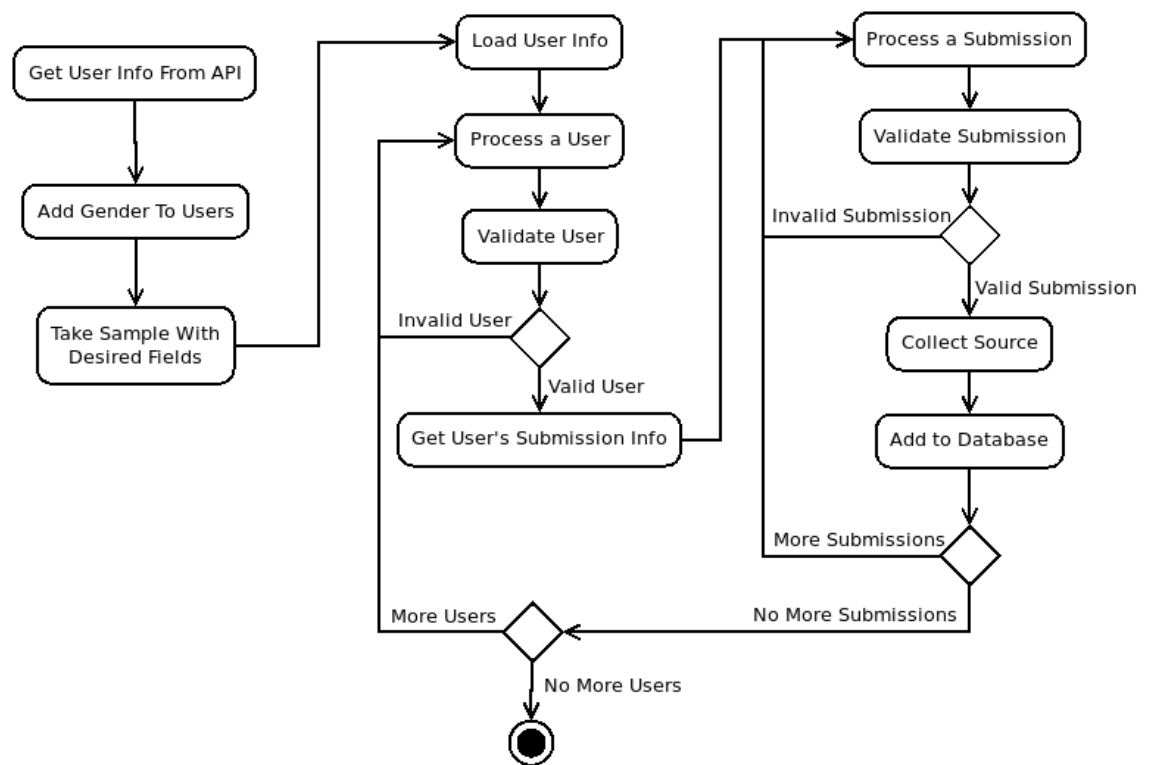Figure 2: Getting GitHub Project Source Code

Figure 3: Getting GitHub Commit Source Code

Figure 4: Getting Data From Codeforces