

An Empirical Study on GitHub Repositories - Exploring the Types of Repositories

COMP 4520 - Undergraduate Honours Project Report
Supervisor: Shaowei Wang

Darshan Pandhi
Computer Science
University of Manitoba
Winnipeg, Manitoba, Canada
pandhid@myumanitoba.ca

1 Introduction

GitHub¹ is the largest and most advanced development platform in the world. Millions of developers and companies build, ship and maintain their software on GitHub. According to the company, it contains over 200 million repositories with over 65 million developers utilizing the platform, making it the largest source code host as of April 2020 [10]. As a result, it lends itself to data mining applications for several research studies, leveraging the invaluable open-source projects available in the website's massive database.

Being home to a large chunk of code hosted online does not take away from the fact that it is incredibly diverse. Hence, the goal of this study was to take a deep dive into all the source code available on GitHub in order to infer trends. This shall help quell the curiosity surrounding the proportions of various types of software and the type that enjoys the most support for open-source contributions. In addition, this would enable and act as a catalyst for future studies regarding software data collection in general.

2 Background and Related Work

GitHub Data Mining has been a growing matter of interest over the years. Mathew Russell breaks down the entire process skillfully in their book "Mining the Social Web: Data Mining Facebook, Twitter, LinkedIn, Google+, Github, and More" to the point where it is almost akin to serving as a tutorial [6]. No wonder the methodology for this study is inspired by it.

Another paper, "The Promises and Perils of Mining GitHub" sheds light on how to approach the data in GitHub based on an initial qualitative assessment. It highlights the common mistakes researchers often make when going about this route. However,

since the paper was published in 2014 it will be interesting to see if some of their conclusions still hold.

Valerio et al. take the reverse gear by examining a number of research papers that use GitHub data. Their paper "Findings from GitHub: methods, datasets and limitations" talks about i) the empirical methods employed, ii) the datasets they used and iii) the limitations reported [7]. Lately, Michael Färber follows a similar approach to conduct a thorough analysis of the code repositories linked in scientific papers using the Microsoft Academic Graph as a data source [5].

While there are numerous (previous) studies that utilize GitHub data, no one has really recently explored what kind of projects does GitHub contain i.e., the quality of the dataset. This is what will be the primary focus of this study. With technology advancing at an unprecedented pace, the common approach needs to be constantly evolving in order to keep up. The results for this study shall help in honing GitHub data collection and analysis.

3 Research Questions

The following are the 2 research questions that this paper aims to answer:

- **RQ1: What are the different types of repositories on GitHub?**

While GitHub is commonly used to host open-source projects, understanding them in greater depth will go a long way in dealing with unintentional bias. The assumption is that most of the repositories on the version control platform are engineered software projects.

Software can be divided into 4 main types: Application Software, System Software, Programming Software and Driver Software [8]. The goal is to break them down further using manual classification to gain a granular perception of the current state of things.

¹ <https://github.com/>

Moreover, such detailed introspection will also allow for understanding the language mix.

- **RQ2: Can we build a classifier for GitHub repositories?**

Why do we even need such a classifier in the first place? It will be an understatement to claim that GitHub contains a lot of software code. However, more often than not research projects focus on a narrow subsection, involving a particular type of software. For example, one of my recent projects was “What Are the Characteristics of Most-Popular Apps? - A Case Study on Android Applications”². Undoubtedly, I examined quite a few mobile applications and not just any mobile applications but Android applications specifically. Therefore, studies like these would benefit immensely from a GitHub classifier that helps them collect the right subset of the available data quickly and efficaciously.

4 Data Collection and Experimental Setting

For this study, I needed a large number of software projects but more importantly supplementary data like the languages used in the codebase, ratio and lines of code, number of contributors, etc. in order to compare and contrast and infer trends.

I started with the GitHub Activity Data hosted on Google BigQuery. With over 3 million repositories³, it is the largest released source of GitHub activity to date. However, the schema not only missed several key attributes but also was devoid of an easy way to download the repository contents inexpensively and efficiently. Consequent attempts included experimenting with various tools and datasets like GHTorrent⁴, SEART⁵Tool [1], Kaggle⁶, etc. Unfortunately, none of these attempts proved to be fruitful due to a plethora of reasons including but not limited to inconsistencies, obsolete and incomplete data, broken URLs, data skewness, lack of documentation, etc. Ultimately, GitHub Search API⁷ was the one that outsmarts the competition since it offered the maximum information and was easy to use, thanks to the massive community support it enjoys. To see this in action, the Python programming language is used for implementing data mining and data processing. PyGitHub⁸ complements this choice beautifully by making it very convenient to access the GitHub

REST API. Finally, the code for the entire project is rightfully hosted on GitHub⁹ itself.

The pertinent task of selecting the right data for further analysis was taken care of by the API itself as it returned a random assortment of the said number of repositories. For this experiment, 100 such repositories were taken a closer look at. Despite that assurance, the first batch of repositories wasn’t the most favourable for the task at hand. This is because a majority of those were outdated, unknown, incomplete, poorly documented and at times even empty / moved to new repositories. This made manual analysis hard and time-consuming. Accurately labelling them often involved a quick examination of the source code. Even after all this, there were ample repositories in the end that didn’t have enough information to clearly decipher their intent. Hence, I started playing around with a variety of filters in order to obtain a higher-quality dataset. After several attempts, filtering out repositories with less than 50 stars and less than 50 people watching seemed to do the job. In addition, only repositories that had a README file made the final cut. With these additional criteria in effect, getting the same 100 repositories involved programmatically scanning over 87,500 of them. The large number also meant data collection spanned over multiple sessions due to the rate-limiting restrictions imposed by the API. Be that as it may, the increased popularity of the repositories meant that there were more contributors to ensure the repositories are well maintained. Simply scanning the project description and /or the README file was sufficient enough this time around and far worth it honestly to confidently label the given software.

4.1 Labels

While on one hand. I was finetuning data selection, I made sure to simultaneously refine the types of software (labels) we started with to comprehensively represent all the data. After careful thought and multiple iterations, the following labels were finalized for software classification. Since the terms are often loosely defined and open to common interpretation, I have used the given connotations for the purposes of this study. This does not take away from the fact that a particular software can belong to multiple categories, i.e., multi-label classification.

1. Library/Framework

To put it simply, a library works within the software while in the case of a framework, the software works within the framework. Initially, both library and framework were separate labels. However, one cannot deny the fact that these are perhaps the most confusing terms, often used interchangeably. Hence for the sake of simplicity, I have currently classified both libraries

² GitHub Repository containing project documentation and associated codebase: <https://github.com/darshanpandhi/What-Are-the-Characteristics-of-Most-Popular-Apps->

³ Link to the Dataset: <https://console.cloud.google.com/marketplace/product/github/github-repos>

⁴ <https://ghtorrent.org/>

⁵ <https://seart-ghs.si.usi.ch/>

⁶ <https://www.kaggle.com/>

⁷ <https://docs.github.com/en/rest/reference/search>

⁸ <https://github.com/PyGithub/PyGithub>

⁹ <https://github.com/darshanpandhi/simplifying-development>

and frameworks into a single category for now. Future studies can involve splitting these into 2 separate ones.

2. Language

This includes all types of languages like general-purpose programming languages, domain-specific languages (DSLs), Scripting Languages, etc.

3. Plugin

A plugin is a software add-on that is installed on a program, enhancing its capability. For example, the Adobe Flash Player plugin that is used to stream and view video, audio and multimedia [4].

4. Web Application

A web application is a client-side and server-side software application in which the client runs or request in a web browser. Popular web applications include eCommerce stores, email, online banking among others.

5. Docs/Tutorials

Any repositories that simply have textual files reside here. This also includes repositories that contain instructional documentation intended to impart knowledge to fellow developers.

6. Moved/Empty/Archived

Just as the name suggests, these include all those that are sans code, moved to a new repository or have been made read-only by archiving.

7. Utility

Utility software is designed to aid in analyzing, optimizing, configuring and maintaining a computer system. It is often the one that supports the computer infrastructure.

8. Protocol

Protocol can be thought of as a set of rules or procedures for transmitting data/information between electronic devices, such as computers.

9. Database/Storage

Anything related to database/persistence falls in this category. Naturally, this also includes the likes of an Object-Relational-Mapper (ORM).

10. Mobile Applications

A mobile application (app) is a type of application software designed to run on a mobile device, such as smartphones or tablets. Most common apps are for the ones for the Android and iOS platforms.

11. Testing

This includes any software or tools that aid in ensuring code works as expected by writing test cases.

12. APIs

An Application Programming Interface is simply a software intermediary that facilitates data exchange between two applications.

13. Parsers

A parser is a compiler or interpreter component that breaks data into smaller elements for easier translation into another language.

14. Other

This includes all the remaining categories of software that had very little representation on its own. Some prominent ones are games, user interface (code), interpreters and operating systems.

5 Results

This section describes the results obtained to answer each of the 2 research questions outlined in section 3.

• RQ1: What are the different types of repositories on GitHub?

Table 1 contains a snapshot of 25 repositories and their types. The complete table is hosted on GitHub¹⁰. Alternately, Figure 1 is a graphical representation highlighting the distribution of the types of repositories found.

Libraries/Frameworks are the clear winner with 20% whereas plugins occupy the runner-up spot with 14% to their name. 11% of GitHub is languages while docs/tutorials take up another 11%. Together, these four categories make up for more than half of GitHub.

Next up, 9% repositories are empty now, have been moved or are archived. And this is what is the most surprising revelation of all. The fact that a significant part of GitHub is not pure software projects as one would expect. As a matter of fact, there were countless instances of only text files, abandoned projects, tutorials including class notes, half-baked academic assignments, archived repos, resumes hosted online, lists of popular movies, books, software companies, food and so on and so forth. Furthermore, the majority of the projects are personal and inactive. As you can see in the chart these two (docs/tutorials and moved/empty/archived) together accumulate for 20% of the lot i.e., 1 out of every 5 repos. To make things

¹⁰ <https://github.com/darshanpandhi/simplifying-development>

worse, many repositories are forks or forks of forks and most of the time the original project would be of greater interest for research purposes rather than their duplicated counterparts.

Web applications aren't too far behind with 7% and lastly, we have database/storage, testing, APIs, parsers, protocol, utility and mobile applications grabbing negligible shares between 2 – 4% each.

Table 1: Snapshot of the Repository List (along with their labels)

Sr. No.	Repo Name	Download URL	Type
1	mojombo/grit	https://github.com/mojombo/grit	Library/Framework
2	rubinius/rubinius	https://github.com/rubinius/rubinius	Other
3	mojombo/god	https://github.com/mojombo/god	Library/Framework Utility
4	mojombo/chronic	https://github.com/mojombo/chronic	Plugin
5	engineyard/eycap	https://github.com/engineyard/eycap	Other
6	defunkt/facebox	https://github.com/defunkt/facebox	Library/Framework
7	haml/haml	https://github.com/haml/haml	Language
8	sferik/twitter	https://github.com/sferik/twitter	Other
9	IoLanguage/io	https://github.com/IoLanguage/io	Language
10	seven1m/onebody	https://github.com/seven1m/onebody	Web Application
11	mislav/will_paginate	https://github.com/mislav/will_paginate	Library/Framework
12	aasm/aasm	https://github.com/aasm/aasm	Library/Framework
13	dustin/java-memcached-client	https://github.com/dustin/java-memcached-client	Database/Storage
14	programming-nu/nu	https://github.com/programming-nu/nu	Language
15	preservim/nerdtree	https://github.com/preservim/nerdtree	Plugin
16	preservim/nerdcommenter	https://github.com/preservim/nerdcommenter	Plugin
17	spree/spree	https://github.com/spree/spree	Web Application
18	arclanguage/anarki	https://github.com/arclanguage/anarki	Language
19	psychs/limechat	https://github.com/psychs/limechat	Other
20	DotNetOpenAuth/DotNetOpenAuth	https://github.com/DotNetOpenAuth/DotNetOpenAuth	Library/Framework Protocol
21	mbleigh/acts-as-taggable-on	https://github.com/mbleigh/acts-as-taggable-on	Plugin
22	bborn/communityengine	https://github.com/bborn/communityengine	Plugin
23	prototypejs/prototype	https://github.com/prototypejs/prototype	Library/Framework
24	sstephenson/sprockets	https://github.com/sstephenson/sprockets	Moved
25	phusion/passenger	https://github.com/phusion/passenger	Web Application

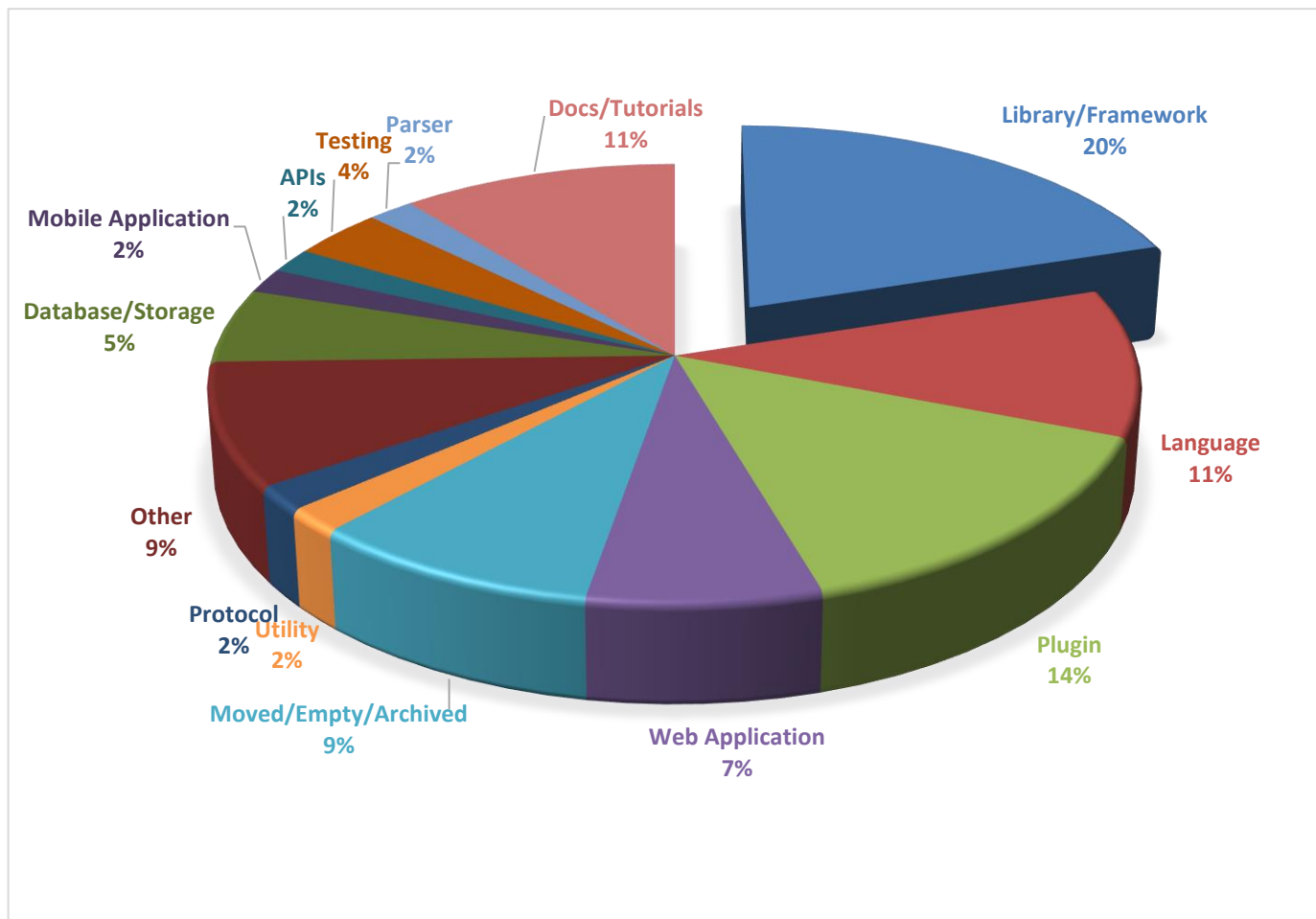


Figure 1: Types of GitHub Repositories

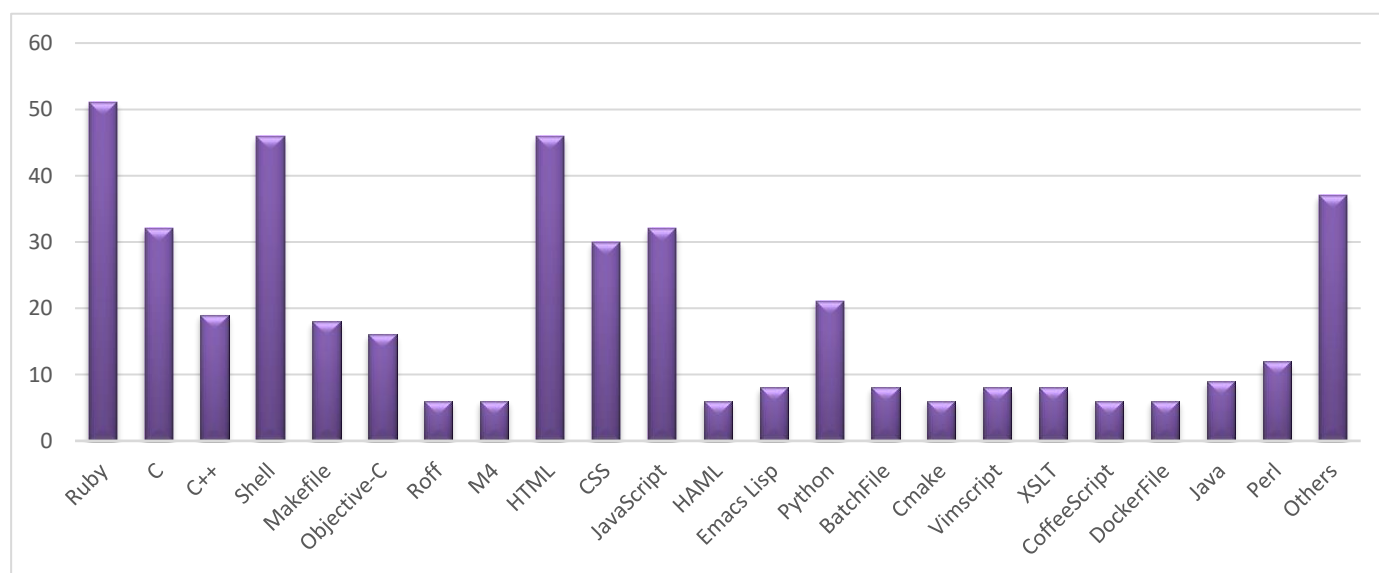


Figure 2: Language Mix

Lastly, remaining software like games, drivers, operating systems, firmware collectively forms 9% of all repositories (as Others).

All in all, the distribution isn't completely asymmetrical. There isn't one type of software that greatly dominates. Although there is a preference for libraries/frameworks, this is again due to the fact that these two were combined as one category.

The GitHub Search API also made it possible to compare the various languages used across the board as seen in Figure 2. Ruby is used the maximum with languages like HTML and Shell following closely behind. C is used in 32 out of the 100 repositories and so is JavaScript. CSS is next with 30 occurrences. Python is gaining popularity with while Java seems to be diminishing over the years with only 9 repositories actively using it. Languages like Objective-C and Makefile occupy the middle spot whereas BatchFiles, Cmake scripts, Vimscripts find themselves in single digits. The Others category combines languages that have nearly negligible representation in the grand scheme of things but together constitutes about 37%. It includes the likes of TeX, R, VBScript, Dtrace, Pascal, Ada, Assembly, etc.

RQ2: Can we build a classifier for GitHub repositories?

These results give more vigour to develop a robust classifier that could separate various types of software. It shall save lots of time and effort otherwise spent on preprocessing and filtering the data set. Features that can very well be used to develop the model are ratio/lines of code, programming languages (if any), the number of contributors, stars, forks, people watching, last commit, pull request comments, issues, task board, etc. Alongside, care should also be taken to exclude prehistoric data points by going over the repository activity.

As it happens, manual analysis of the repositories during the course of this study allowed me to form many preliminary associations between the available statistics and the type of software. One obvious one is to look for any presence of code. GitHub conveniently provides information regarding the languages used in a project. So, if a particular repository does not use any language, it has to be either documentation or an empty project or a project that has been moved to a new place.

6 Implications

The biggest takeaway from this research study hands down is the fact that GitHub is not only used for software but also merely as a means to store something for private use in the future or to host and share something with the entire world.

Therefore, this study should serve as a cautionary tale for subsequent research activities that rely on data from GitHub. After all, an accurate dataset is what makes up the foundation for any investigation. Corroding that will not only lead to invalid and biased results but also affect future dependent studies.

Just to secure a dataset of 100, easy enough for precise labelling saw me going over 87,500 repositories. That is nearly ~0.11%. Clearly, cleaning up the dataset and filtering out all the noise is an absolute must.

7 Future Work

This experiment helped solidify the significance of a GitHub repository classifier. Such a model would have limitless applications. From recommendation systems to even search engine optimization (based on the type of software), the list is practically endless. And incorporating the language trends acquired here would be the icing on the cake. Speaking of the language spread, the comparison done in this study is at an elementary level. There is room for future expansion to understand the rational behind the numbers we observe in the results.

Further, one of the threats to validity for this project was the sample size. GitHub rate limits to 1000 requests per hour for its API¹¹. This combined with insanely high noise levels and time constraints made it difficult to manually analyze large samples of data, leaving room for future explorations.

Acknowledgement

I would like to thank Shaowei Wang under whose supervision I have been conducting this study. Without his constant support and guidance, this paper would simply not have been possible, to say the least. Further, I would also like to express my gratitude to Mathew Russell whose work on Social Web Mining paved the path forward with minimum obstacles while formulating the methodology for this study.

¹¹ This is when `GITHUB_TOKEN` is GitHub Actions is used.

July 9, 2014. During recent years, GITHUB (2008) has become the largest code host in the world.

References

- [1] Dabic, O., Aghajani, E. and Bavota, G., 2021. Sampling Projects in GitHub for {MSR} Studies. IEEE, pp.560--564.
- [2] squareboat, n.d. Different Types of Software with Examples. Available at: <<https://squareboat.com/blog/different-types-of-software-with-examples>> [Accessed 31 August 2021].
- [3] Hope, C., 2021. What is a Plugin?. [online] Computerhope.com. Available at: <<https://www.computerhope.com/jargon/p/plugin.htm>> [Accessed 31 August 2021].
- [4] McLaughlin, E., 2021. What is Adobe Flash Player?. [online] SearchCIO. Available at: <<https://searchcio.techtarget.com/definition/Adobe-Flash-Player#:~:text=Adobe%20Flash%20Player%20is%20software,c omputer%20or%20supported%20mobile%20device.>>> [Accessed 31 August 2021].
- [5] Michael Färber. 2020. Analyzing the GitHub Repositories of Research Papers. In Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020 (JCDL '20). Association for Computing Machinery, New York, NY, USA, 491–492. DOI:<https://doi.org/10.1145/3383583.3398578>
- [6] Russel, M., 2014. Mining the Social Web: Data Mining Facebook, Twitter, LinkedIn, Google+ 2nd ed. Sebastopol, CA: O'Reilly.
- [7] Valerio Cosentino, Javier Luis, and Jordi Cabot. 2016. Findings from GitHub: methods, datasets and limitations. In Proceedings of the 13th International Conference on Mining Software Repositories (MSR '16). Association for Computing Machinery, New York, NY, USA, 137–141. DOI:<https://doi.org/10.1145/2901739.2901776>
- [8] Wilcox, L., 2021. [Blog] The 4 Main Types of Software, Available at: <<https://www.leadwithprimitive.com/blog/the-4-main-types-of-software>> [Accessed 1 September 2021].
- [9] Y. Tian, M. Nagappan, D. Lo and A. E. Hassan, "What are the characteristics of high-rated apps? A case study on free Android Applications," 2015 IEEE International Conference on Software Maintenance and Evolution (ICSME), Bremen, 2015, pp. 301-310, DOI: 10.1109/ICSM.2015.7332476. ISBN:978-1-4503-0000-0/18/06
- [10] Zaidman, Andy; Gousios, Georgios; Vasilescu, Bogdan; Serebrenik, Alexander. "Lean GHTorrent: GitHub Data on Demand" (PDF). The Netherlands: Delft University of Technology & †Eindhoven University of Technology: 1. Archived (PDF) from the original on April 25, 2020. Retrieved