# Business Concept

## Overview

Our business goal is to get a better understanding of current technology trends. In this project the main focus will be on the emerging data startups, especially companies related to the modern data stack.

We want to understand which are the hot projects and see how the monitored projects perform compared to others.

The objective was to take a real-life problem and provide a complete, automated and parameterizable solution at the end, where all the stakeholders are represented.

The project will help to provide end-users enabling them to make business decisions, get a real picture of these companies and they can then perform business actions on this basis.

The main objective is to better understand the development of the tech market. To help with this, there are many data sources available and some of which have been selected to help answer these complex questions.

To make this data meaningful for end users, a data visualization tool will be used in the final phase of the project. We want to make the data prepared, accessible and meaningful through dashboards, charts, KPIs and other metrics.

Possible questions the business side may seek answers:

- The active community of open-source and other solutions within a domain
- Rate and popularity of new players compared to old players
- What activity a new release has generated within a given product
- Comparing two projects from the same field by certain attributes
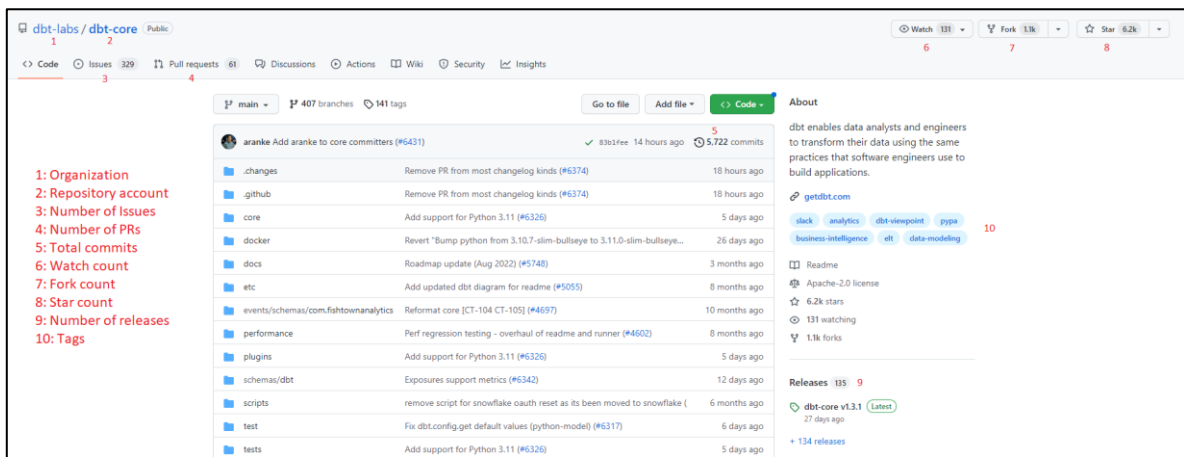
## Datasets

To answer these questions, we will use the following freely available datasets provided by BigQuery:

- GitHub
- Stack Overflow
- Hacker News

### GitHub

GitHub is a code hosting platform for version control and collaboration. It lets people work together on projects from anywhere.

The BigQuery dataset comprises the largest released source of GitHub activities. It contains a full snapshot of the content of repositories including PRs, commits, forks and other activities that can be carried out on the site.
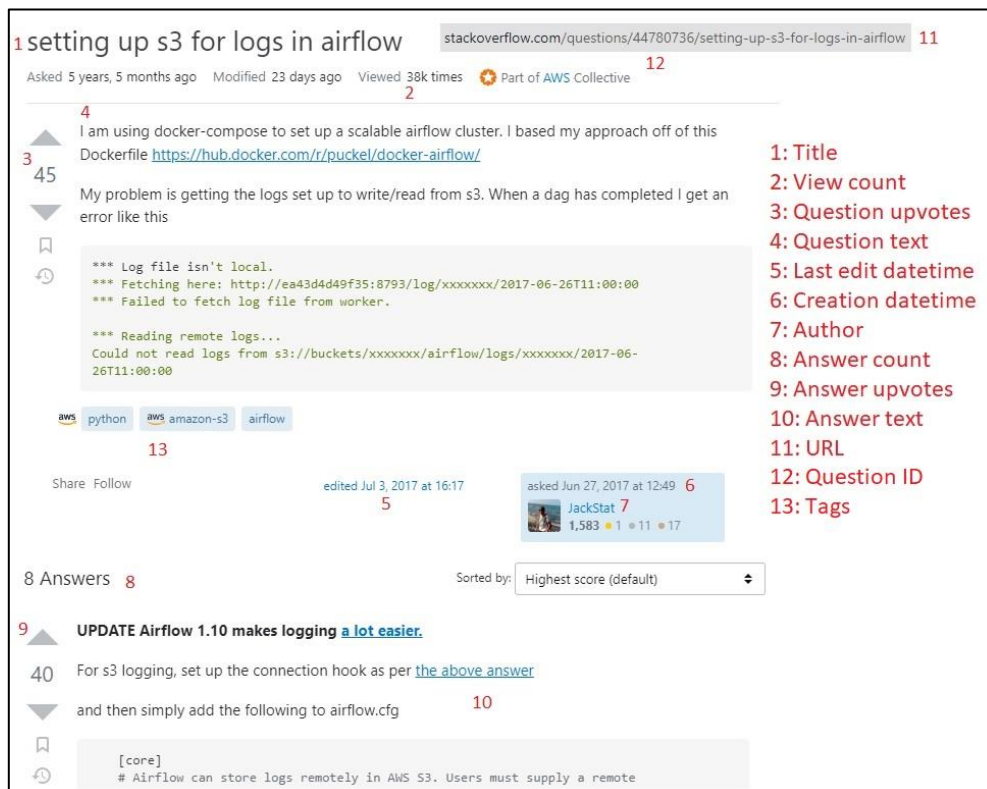
*GitHub screenshot with comments*

This screenshot shows the home page of a random GitHub repository. The numbers in the screenshot indicate the attributes visible and readable on the web page that are also stored in the BigQuery data source, but in a different format.

For the GitHub data source, the organizing principle for separating the projects will be the organization and possibly an exact repository within an organization.

**Stack Overflow**

Stack Overflow is a question and answer website for professional and enthusiast programmers. It features questions and answers on a wide range of topics in computer programming. It is the largest online community for programmers to learn, share their knowledge and advance their careers.

The BigQuery dataset includes an archive of all the Stack Overflow content, including posts, votes, tags and badges.



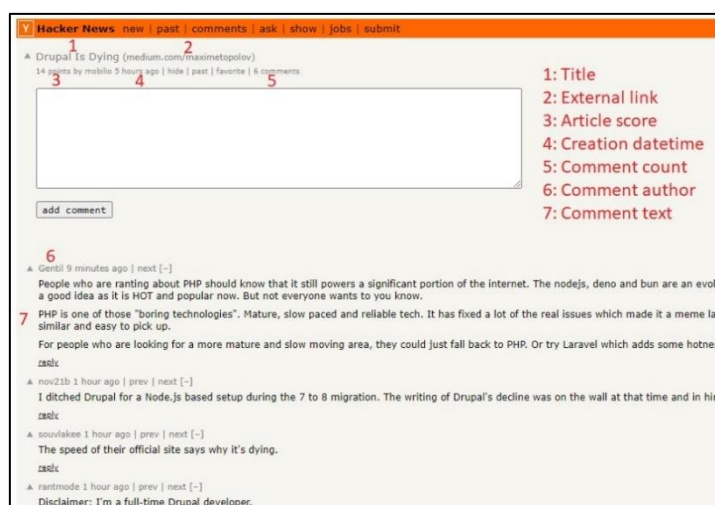*Stack Overflow screenshot with comments*

This screenshot above shows a random Stack Overflow question. The numbers in the screenshot indicate the attributes visible and readable on the web page that are also stored in the BigQuery data source, but in a different format.

For the Stack Overflow data source, the organizing principle for separating the projects will be the tags in the questions or possibly a labeling process derived from the question text.

**Hacker News**

Hacker News is a social news website focusing on computer science and entrepreneurship.

The BigQuery dataset contains all stories and comments from Hacker News from its launch in 2006 to present. Each story contains the author that made the post, when it was written, the number of points the story received and all the related comments.



*Hacker News screenshot with comments*

**Configuration**

From a business point of view, only certain projects are of interest and it is necessary to separate these from the complete data sources.

This requires an input document that contains a list of the companies currently being used, with additional dimensions that are not part of the BigQuery data sources but are of high relevance to the project. For example, what category a company falls into, or whether it has an open-source solution or not.

In this project, it will be stored in a google sheet, which will be ingested and added to our data warehouse.

To simulate real business operations, this google sheet may change from time to time as companies are added or removed along the way. We will also manage these changes as part of the project.

| Organization | Repository account | Repository name | L1 type | L2 type | L3 type | Tags | Open source available |
|---|---|---|---|---|---|---|---|
| Astronomer | astronomer | | Modern data stack | Orchestration | | astronomer | No |
| Monte Carlo | monte-carlo-data | | Modern data stack | Monitoring | | montecarlo | Yes |
| Lightdash | lightdash | lightdash | Modern data stack | Analytics | | lightdash | Yes |
| Superset | apache | superset | Modern data stack | Analytics | | apache-superset, superset, apache superset | Yes |
| Snowplow | snowplow | snowplow | Modern data stack | Web tracking | | snowplow | Yes |
| Firebolt | firebolt-db | | Modern data stack | Data warehouse | | firebolt | No |

*Company details example*