# Citizen Friendly Report of diversitydatakids.org

Anil Gubbala, Ali Hussain Ladiwala, Urja Naik and Priyanka Cornelius

Department of Software Engineering, San Jose State University´

San Jose, CA´

anil.gubbala@sjsu.edu, alihussain.ladiwala@sjsu.edu, urjamitulkumar.naik@sjsu.edu,
priyanka.cornelius@sjsu.edu

*Abstract*—Using Divrsitykids datasets, create a dashboard and use NLG to generate a citizen friendly report.About Diversitydatakids.org: It is a research project that examines who our children are, whether they have what they need to grow up healthy and achieve their full potential, whether social policies are well designed to improve children's lives and how to make them better to improve equity.Data collected by diversitydatakids is in tabular format, which makes it difficult for a general citizen to understand the content.Our goal is to convert few topics of this tabular content to a citizen friendly report so that it becomes more readable & helps increase awareness among citizens.

## I. INTRODUCTION

Early childhood is a period of rapid development for children. Early childhood experiences at home, in school and in their neighborhoods shape their ability to thrive in childhood and beyond. diversitydatakids.org examines equity (or inequity) in children's early experiences and in policies that support their wellbeing in early childhood, such as early care and education so that policymakers and practitioners can build on the work as they consider expansions of existing early childhood programs and to inform new federal, state and local policies. However users from non-engineering backgrounds might find it cumbersome to interpret these extensive datasets as numerical or tabular data is not intuitive to the brain. Our attempt is to obtain the data in citizen friendly textual format so as to augment its impact.

## II. ARCHITECTURE

There are 5 major components:
Dataset on diversitydatakids.org which was fetched using CKAN's API , Processing & filtering this data in node.js & react.js, Natural Language Generation library RosaeNLG based on pug templates, Display tabular data and generated text in frontend, Deploy application on AWS server.

CKAN API is used to fetch datasets that are used in the diversitydatakids.org website and then do data filtering & data extraction. The filtered data is sent to the Node backend server which pushes the data to RosaeNLG that uses Pug templates for NLG structuring.
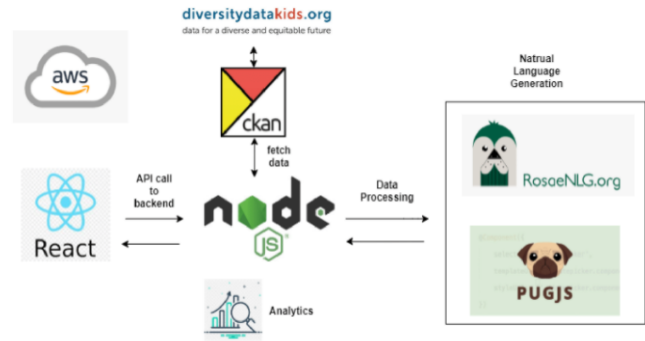


Fig. 1. Project Architectures

## III. DATA ANALYSIS

Established in 2014 with support from the W.K. Kellogg Foundation and the Robert Wood Johnson Foundation, diversitydatakids.org has 323 datasets, subdivided into Topics such as Child Opportunity Index, Demographics, Income, Education, Family, Health, Housing, Early Childhood, Neighborhood and Work. The datasets can further be filtered on the basis of 'Availability by Race and Ethnicity', 'Nativity', 'Age Group', 'Gender', Geography' and 'Time Scale'. The dataset is publicly accessible from the author's website [1].

### A. Data Distribution

The dataset available is in the form of csv files & JSON, categorised into several topics and subtopics like 'Early Childhood', 'Income, Work and Family' and 'Neighbourhoods' etc. In total there are over 6000+ datasets available in the database and can be fetched using CKAN data API. Majority of the datasets contain data examining how racial/ethnic and income segregation in neighborhoods and early childhood programs contribute to inequities in children's early experiences and opportunities. Column titles are tags that have been described in a separate table and include geographic ID, region and numerical values for opportunity indices termed as raw indicator values. Sample Data for Child Opportunity Index is shown in Fig 2.

Showing 1 to 10 of 144,408 entries

| geoid | year | in100 | msaid15 | msaname15 | countyfips |
|-------|------|-------|---------|-----------|------------|
| 01001020100 | 2015 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |
| 01001020100 | 2010 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |
| 01001020200 | 2015 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |
| 01001020200 | 2010 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |
| 01001020300 | 2010 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |
| 01001020300 | 2015 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |
| 01001020400 | 2010 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |
| 01001020400 | 2015 | 0 | 33860 | Montgomery, AL Metro Area | 01001 |

.

Fig.2. COI 2.0 Index Data

### B. Fetching Dataset

Fetching a dataset is not a straightforward process. Inorder to fetch the dataset first we have to select a subtopic and then selection geographic region to get the dataset. In datastore at present there are over 320 subtopics available. Only alternative is to enter the exact resource ID and package ID of the dataset to access that particular dataset and view the data. Our implementation includes both these approaches.

### C. Dataset Column titles

The column titles in the datasets on diversitydatakids.org are column tags or labels that have been separately described in a table. In order to use these column headers for natural language generation, they have been converted into a more descriptive format that clearly explains the contents of the column.

### D. Dataset statistics

The CKAN data API supports several forms of querying including  Javascript , python and URI. In addition, it also supports inclusion of direct SQL query & Apache Solr querying for a few API calls. So we make use of these features to extract summary of the datasets  using aggregate functions in SQL.

### E. Text  Decoding

Even though there are over 6000 datasets available in the datastore, the structural format of all the datasets is the same for all the datasets.  Every dataset also maintains some keywords specifying to which topics, subtopics or categories the current dataset represents. There are around 10 major categories which are further categorized to several other sub categories. This structural uniformity made it easy to extract data from the dataset required for NLG for almost all the datasets available in the datastore.

### F.  Keywords in NLG

The filtered and extracted data is sent to RosaeNLG where the actual text generation is implemented. The keywords & the metadata information in each dataset is used here to generate the text in RosaeNLG. This data is used in the major 5 stages of Natural language generation to generate the text[3].  Few of the text used in the NLG is static but most of the text is dynamically imported. So the text generation supports almost all the datasets in datastore.

## IV. ARCHITECTURAL COMPONENTS

Natural Language Generation tools and libraries require structured data in order to convert it to textual format. We obtain the required structured data in the form of JSON using CKAN data API. We fetch the data from the website using CKAN API employing SQL, Apache Solr queries and then format the JSON output in a Pug template that RosaeNLG uses to generate natural language text.

### A. CKAN data API

CKAN is an open-source DMS (data management system) for storage and distribution of open data. Instead of a database, CKAN API has been used to fetch data from the datastore. We performed API calls from both react as well as node.js based on functionality being implemented. [4][6]

### B. PostgreSQL & Apache Solr

CKAN maintains information in PostgreSQL databases since PostgreSQL is open source and widely available. Ckan has an Apache Solr component that enables search implementation with results in JSON format. Once we get the JSON content from CKAN we filtered the data and updated the content based on the filters. [5]

## C. Pug

Pug is a templating engine designed for node or browser execution. It is fast, robust and feature-rich. RosaeNLG templates are basically pug templates where you use RosaeNLG structures and mixins to complete the standard pug syntax.[2][8]



Fig3.Pug template

## D. RosaeNLG

RosaeNLG is a Natural Language Generation library for node.js or client side (browser) execution, based on the Pug template engine[8]. The filtered data obtained from CKAN is sent to the Node backend server which in turn sends the data as an input to the RosaeNLG library in order to perform NLG conversion.[7]

## E. React.js & Node.js

In our application we used react.js for frontend development and Node.js for backend development. Node.js is used to connect to RosaeNLG whereas both react.js & node.js are used to connect to CKAN API to fetch required data.

## V. IMPLEMENTATION

The Citizen Friendly Report has been presented on a dashboard that lets the user select datasets from a wide range of topics, search for datasets of their choice using relevant search keywords or choose from a list of available categories and further select datasets according to the region and tenure for which the data was recorded. The selected dataset can then be viewed in the dashboard and the textual report is generated summarising the range of the entire dataset, range of the region wise filtered data and the selected rows. Each selected row has also been drawn as a bar chart for visualization and comparison.

## A. Find relevant SubTopic using Keyword Search

Considering the huge count of datasets on the website and the vast research topics they cover, a search keyword filtering mechanism has been implemented that allows the user to enter keywords and search datasets relevant to their choice. The user can also select relevant search keys from a list of 'Available Categories' that are found to be the most common in the dataset descriptions. The search results appear as a list of clickable links to the dataset's region-wise categorization. A description of each dataset has been provided alongside the topic, fetched from the topic descriptions on diversitydatakids.org .



Fig. 4.Find relevant datasets using keyword search

## B. Region wise filtering of data

Each Subtopic listed as a search result in the Search Page can be further segregated according to region and tenure represented by the dataset. This segregated list of datasets can further enhance the search for relevant datasets. The user may choose a dataset belonging to a particular City, State, Nation, Zip Code, Census Division, Unified School Districts, Metro Areas, Counties etc, and further pick a 1 year dataset or a 5 year dataset in most cases.

Fig. 5. Region wise filtering of data

## C. Additional Filters

Once a user accesses a particular dataset , by default a region is selected and the data corresponding to that region alone is displayed to the user. Users can search the region manually or select from the available dropdown which gets lazily updated on input. Without this filter there can be thousands of records per dataset which makes it difficult for users to process the information.



Fig.7. Default region selected at top left

## C. Citizen Friendly Report

The Citizen friendly report generated by the rosaenlg library has been presented under three subsections that can provide a holistic view of the dataset as well as represent each row of the data in textual format.

Firstly, running an analysis on the selected dataset, a textual representation of the range of data is represented which includes observations such as the mean, topmost and the least values cataloged in the dataset for the selected region and tenure.

Following this, a textual report has been generated that presents an analysis of the data belonging to the selected city/state. A dropdown list of each distinct region recorded in the dataset has been provided so as to enable the user to generalise region specific data. This report shows mean and

the highest and lowest values recorded in the dataset for the selected region and the average distribution of factors for ethnicities, races, income groups etc.



Fig. 6. Range of cataloged data in textual format

Finally, upon selecting each row in the dataset, a textual representation of the selected row is presented which gives the stats and observations recorded in the row in a structured format that is more intuitive for users.



Fi. 8. Selected row data in textual format

Since each row in the dataset contains numerical stats for various races, ethnicities or such factors specific to the dataset topic, a graphical representation has been included that aids users in visual comparison. In most cases a bar chart has been used to represent the data values for each ethnicity/racial group.



Fig. 9. Graphical representation of selected row

## X. Conclusion

We were able to develop a citizen friendly version of diversitydatakids.org website with the help of nlg. This will help US citizens to get a proper knowledge about the data and help in making the rules beneficial for the children and helps in research of the needs for kids to get proper development.

## XI. Future Improvements

While we are proud that we were able to develop a site which is citizen friendly and can create a greater knowledge for the society in one semester, we know the idea still has room to grow.

A. Adding more graphical Data to the website for better visualization of the data
B. Add more analytics in NLG like what's the trend in past few years or trend of a particular ethnic group.
C. Data visualization for the overview of the whole dataset.

## XII. Deliverables

A. *GitHub Repository*
https://github.com/SJSUSpring21/Citizen-Friendly-Report-of-diversitydatakids.org

B. *Presentation*
https://github.com/SJSUSpring21/Citizen-Friendly-Report-of-diversitydatakids.org/blob/main/Citizen%20Friendly%20Report%20of%20DiversityDataKids.org.pdf

C. *Link of Running Application (Deployed in AWS)*
http://3.16.112.6:3000

## References

[1] https://data.diversitydatakids.org/dataset
[2] https://rosaenlg.org/rosaenlg/3.0.0/advanced/pug.html
[3] https://en.wikipedia.org/wiki/Natural-language_generation
[4] https://docs.ckan.org/en/latest/api/index.html
[5] https://solr.apache.org/guide/6_6/common-query-parameters.html
[6] https://docs.ckan.org/en/2.8/maintaining/datastore.html
[7] https://rosaenlg.org/rosaenlg/3.0.0/tutorials/tutorial_en_US.html
[8] https://pugjs.org/api/getting-started.html