# Federated Data Management Aggregation Framework by using NLP

**[1]Prof. Swapnil Wani, [2]Miss.Jyoti Gupta, [3]Miss.Namrata Kasturi, [4]Mr.Suraj Yadav**

**[1]Asst.Professor, [2,3,4]UG Student, [1,2,3,4]Computer Engg. Dept. Shivajirao S. Jondhle College of Engineering & Technology, Asangaon, Maharashtra, India. [1]swapnilwani27@gmail.com,[2] jyoti23147@gmail.com, [3]namratakasturi30@gmail.com, [4]sy929611@gmail.com.**

**Abstract- Given the increasing range of heterogeneous information hold on in relative databases, file systems or cloud environments, it has to be simply accessed and semantically connected for any information analysis. The potential of information federation is essentially unrealized, this paper provides associate interactive information federation system by applying large-scale techniques as well as heterogeneous information federation, association rules and linguistics internet to perform information retrieval and social network analytics information. Providing an easy framework to make virtual databases so mapping them over to the actual information instance implementing "Eventual Consistency". The project contains a javascript framework that has associate object relative plotter, configuration management scripting setting, associate API to attach to databases in a happening driven fashion and a separate daemon which is able to act because the human service for the API.**

**Keywords:-- federation, information analysis, small services, software system design**

## I. INTRODUCTION

Given the growing amount a collection of disparate data stored in file systems, relational databases, and cloud settings, it must be retrieved quickly and semantically connected for further data analysis. Progress in bar-code technology has created it double for retail organizations to gather and store large amounts of sales information, mentioned because the basket information. A record in such information generally contain the dealings date and therefore the things bought within the dealings. No-hit organizations read such databases as vital items of the selling infrastructure. They're fascinated by instituting information-driven selling processes, managed by information technology, that change marketers to develop and implement tailor-made selling programs and techniques. Semantic net techniques (RDF, SPARQL) are wide used for information federation and linkage. However, the first issue of the linguistics net is short integrated resolution. To solve this issue, this paper tends to applied information federation, linguistics net, and data processing technologies to develop this technique, which might complement with one another. The system permits users to pick information sources, act with visualized graphs, and run custom-built queries across united information to fulfill specific wants for information analytics. Range aggregation may be a fundamental process in abstraction information applications, and there is an increasing desire for such processes to be supported over a knowledge base federation, wherever the whole abstraction information are individually command by multiple information suppliers

(a.k.a., information silos). Information federations notably increase the number of knowledge offered for data-intensive applications like good quality coming up with and public health emergency responses. However they conjointly challenge the traditional implementation of vary aggregation queries as a result of the data can't be shared at intervals the federation and also the information partition at every data silo is fastened throughout question process.

## II. AIMS AND OBJECTIVE

### a) Aim

Providing a simple framework to create virtual databases and then mapping them over to the database itself instance implementing "Eventual Consistency". A javascript framework and a predictive query parser for one shot multi-data layered query, in declarative style. Data federation will employ association rule analysis to establish patterns and identify the most important co-occurrences of variables across many data sources. This paper have a tendency to created virtual info (VDB) for knowledge federation, wherever knowledge is accessed and nearly integrated in time period across distributed knowledge sources while not repeating or otherwise moving knowledge from its system of record.

### b) Objective

The project's major objective is to make aggregator microservices and API gateways more adaptable and simple to set up.. Less computational power or resources like RAM, CPU, etc. In less bandwidth of data, can achieve

more network efficiency. Microservices are the most relevant applications of backend software engineering.

## III. LITERATURE SURVEY

### Paper 1: Fast Algorithms for Mining Association Rules:

The two proposed methods for sizing scale-up can be combined to create a hybrid algorithm, called AprioriHybrid. Additionally, As the quantity of objects inside the information grows, the execution time lowers somewhat. Because the average group action size will increase (while keeping the info size constant), the execution time will increase solely step by step. These demonstrations show that practicability of victimization AprioriHybrid in real applications involving terribly giant databases. The algorithms given during this paper are enforced on many knowledge repositories, together with the bird genus classification system, DB2/6000, and DBS/MVS.

### Paper 2: Efficient Approximate Range Aggregation over Large-scale Spatial Data Federation:

Data federations significantly enhance the quantity of data available for data-intensive applications like good quality designing and medical emergency responses. Nevertheless they conjointly challenge the standard implementation of vary aggregation queries as a result of the data cannot be distributed among the federation and therefore the information partition at every data silo is fastened throughout question process. These control limits the planning house of distributed vary aggregation question process. During this work, the paper proposes approximate algorithms for economical vary aggregation over spatial information federation.

### Paper 3: Graph-based Interactive Data Federation System for Heterogeneous Data Retrieval and Analytics:

For heterogeneous data analytics, the article presented a graph-based data federation architecture. It's a Query builder and result set visualizer in SPARQL for diverse data sources that's free source, which permits users to effortlessly generate and examine data over a variety of data sources. The system is scalable because it allows users, particularly academics, to simply add additional information from multiple sources and customize data format and analytics by users, especially researchers based on their own interests. To digitally integrate data from numerous data sources, the system generates a Virtual Database (VDB). Then, using an RDF generator and SPARQL queries, a data model search over the generated textual data is supported by NLP.

## IV. EXISTING SYSTEM

There exist some popular IBM Info Server for the Federation of Spheres (https://ibm.co/2qWQbom) and Oracle Data Service Integrator (https://goo.gl/6MKXkF) are two examples of enterprise data virtualization products. Similar efforts in data federation have been seen from academia such as Bio Mart (www.ensembl.org/biomart) and Maelstrom (www.maelstrom-research.org). However, considering advanced information analytics across united information is unnoticed. In this paper, have a tendency to planned RDF-supported information visual image framework over united data-processing-enhanced databases (For instance, consider association rules) and information science techniques (e.g., sentiment analysis). It will with efficiency federate and analyze giant heterogeneous information sources for general or specific analysis desires.

## V. COMPARTIVE STUDY

| SR NO. | PAPER TITLE | AUTHOR NAME | PUBLICATION | TECHNOLOGY | PURPOSE |
|--------|-------------|-------------|-------------|------------|---------|
| 1. | Fast Algorithms for Mining Association Rules | Rakesh Agarwal, Ramakrishnan Srikant | IBM Almaden Research Center, 1998 | AprioriHybrid | Identifying all relevant item-to-item relationship rules in a large transaction database |
| 2. | Efficient Approximate Range Aggregation over Large-Scale Spatial Data Federation | Yexuan Sh, Y Tong, Yuxiang Z, Zimu Zhou, Bolin Ding, Lei Chen | IEEE 2021 | Data Federation (Single Silo) | Single-Silo sampling methods that process queries in parallel, as well as a level sampling-based algorithm |
| 3. | Graph-based Interactive Data Federation System for Heterogeneous Data Retrieval and Analytics | Xuan-Son Vu, Addi A, Erik E,Lili Jiang | IEEE 2021 | Custom SPARQL query builder | Users may quickly create and view information or data from a variety of sources. |

## VI. PROBLEM STATEMENT

Given the growing volume of heterogeneous information in relational database systems, file systems, and cloud settings, it must be easily accessible and semantically connected for further data analysis. The potential of knowledge federation is essentially unutilized, this paper is

going to discuss associate interactive information federation system by applying large-scale techniques as well as heterogeneous information federation, association rules and linguistics net to perform information extraction and analysis on social network information.

## VII. SCOPE

Providing a simple framework to create virtual databases and then mapping them over to the actual database instance implementing "Eventual Consistency". An additional object relational mapper shall be provided to create an abstracted experience for the end user. The project includes a javascript framework which has an object relational mapper, configuration management scripting environment, an API to connect to databases in an event-driven manner and a separate daemon which will act as the aggregator service for the API. A framework is to be made with three phases: predict token, parse query and execute, return multilayer result. Extensive experiments area unit nonetheless to be apply to illustrate the efficiency of methodology.

## VIII. PROPOSED SYSTEM

A JSON file composed of a query with nested queries including database refs. A daemon service running in a loop to recursively parse the json received through file system. A database connector service in adapter fashion with observers looking for change in the file system. A JSON file as soon as placed in the specified path in file system would be observed by the daemon. The service will automatically recursively predict and parse the tokens and build multiple queries ready to fire to their respective databases. All the queries get executed parallely using worker threads. The results are aggregated using the parser specified NLP analysis then returned to an exposed TCP port.

## IX. ALGORITHM

### GUN Graph Synchronization Protocol

**Step 1:** If the string doesn't contain any control or quotation characters, as well as any backslash characters, it's safe to quote it. Otherwise, needs to substitute safe escape sequences for the problematic characters.

```
rx_escapable.test(string)
?"\"" + string.replace(rx_escapable, function (a) {
var c = meta[a];
return typeof c === "string"
? c
:"\\u"+("0000"+ a.charCodeAt(0).toString(16)).slice(-4);
}) + "\""
: "\"" + string + "\""
```

**Step 2:** Make a string out of holder[key].

```
var value = holder[key]
```

**Step 3:** To get a replacement value, If a JSON method exists for the value, use it.

```
if (
value
&& typeof value === "object"
&& typeof value.toJSON === "function"
){
value = value.toJSON(key);
```

**Step 4:** Numbers in JSON must be finite. Non-finite numbers should be encoded as null.

```
return (isFinite(value))
? String(value)
: "null"
```

**Step 5:** Convert the value to a string if it is boolean or null.

```
String(value)
```

**Step 6:** Make a list to retain the partial stringification results for the value of this object.

```
gap += indent;
partial = []
```

**Step 7:** The value will be regarded as an array if it is one. Every element should be stringified. As a placeholder, use null.

```
length = value.length;
for (i = 0; i < length; i += 1) {
partial[i] = str(i, value) || "null";
```

**Step 8:** Wrap all of the pieces together, with commas between them.

```
v = partial.length === 0
? "[]"
: gap
? (
"[\n"
+ gap
+ partial.join(",\n" + gap)
+ "\n"
+ mind
+ "]"
)
: "[" + partial.join(",") + "]";
gap = mind;
return v
```

**Step 9:** Use the members to be replaced if the replacer is a collection (reviver function).

```
if (rep && typeof rep === "object") {
length = rep.length;
for (i = 0; i < length; i += 1) {
if (typeof rep[i] === "string") {
k = rep[i];
v = str(k, value);
if (v) {
partial.push(quote(k) + (
(gap)
? ": "
: ":"
) + v);
}
}
}
```

**Step 10:** Co-relate nested id to timestamp and collection by splitting the token by colon ':'

```
Let uniqueCollectionName =`$
{req.body.collectionName}:$
{req.body.document.id}`          await
gun.get(uniqueCollectionName).put(req.body.document)
```

**Step 11:** Update uuid by incrementing last bit (edge weight) of the found id.

```
await gun.get(`${collectionName}:${id}`).once(node => {
    let result = JSON.parse(JSON.stringify(node))
delete result._;
data = result;
});
data = {
...data,
...req.body.modifiedData,
}
Await gun.get(`${collectionName}:${id}`).put(data);
```

**Step 12:** Update collection record file

```
data = {
...data,
```

```
...req.body.modifiedData,
}
await gun.get(`${collectionName}:${id}`).put(data)
```

**Step 13:** Walk the new structure in a recursive fashion[step 2], giving each name/value pair to a reviver function for possible change.

**Step 14:** A SyntaxError is generated if the text cannot be parsed as JSON.

## X. MATHEMATICAL MODEL

**Handshaking Theorem :**

What is the result of adding the sum of the degrees of all the vertices of a graph? When dealing with an undirected graph,, the contribution of every edge is two times, first to the beginning vertex and again to the termination vertex. As a result, the total number of degrees corresponds to the number of edges multiplied by two. In the Handshaking Theorem, this fact is stated.

A graph that is not directed. 'e' edges is G = (V, E). After that, $2e = \sum_{u\in V} deg(u)$.

In case G is a directed graph,

$\sum_{u\in V} deg^-(u) = \sum_{u\in V} deg^+(u) = |E|$.

An significant implication of the handshaking principle for undirected graphs is that an undirected graph contains an even number of odd-degree vertices. Proof : Let $V_{1}$ and $V_{2}$ be set of an even odd order vertex, respectively. In handshaking theorem, $2e = \sum_{u\in V} deg(u)$.

So, $2e = \sum_{u\in V} deg(u) = \sum_{u\in V_{1}} deg(u) + \sum_{u\in V_{2}} deg(u)$.

The total number of degrees of vertices with even degrees is even. The LHS is likewise even, implying because the total of vertices with odd degrees must be even, resulting in an even number of odd-degree vertices.
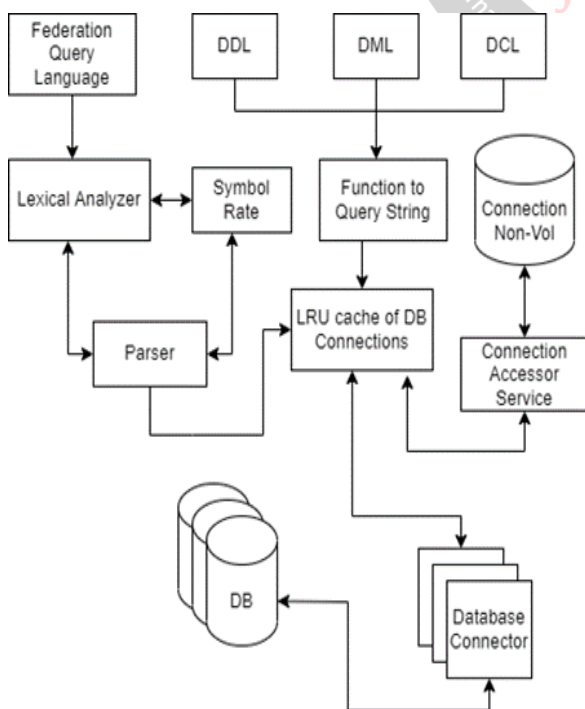
## XI. SYSTEM ARCHITECTURE



Fig.1: System Architecture

A json oriented query language is parsed and executed on multiple nodes to access data. A caching mechanism creates a pseudo connection pool for all database connection nodes. This enables faster connector selection and query speeds. User may also manually write commands directly to a preselected node. The custom json query is parsed using the lexical analyzer and the symbols to correspond them to. A nonvolatile file contains all the frozen connections that can be brought alive during the runtime if a cache miss occurs.

## XII. ADVANATGES

(1) Aggregator microservices and API gateways are more configurable and adaptable.

(2) Less computational power or resources like RAM, CPU, etc.

(3) In less bandwidth of data, it can achieve more network efficiency.

(4) Microservices are the most relevant applications of backend software engineering.

(5) Less computation as generators can generate tokens procedurally.
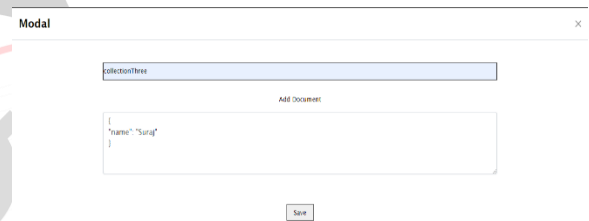
## XIII. OUTPUT



Fig. 2

Web modal to add a document to a particular collection. Documents can be normal text, or special syntactical json for graphical functions for the db. Same modal is used for editing existing documents.
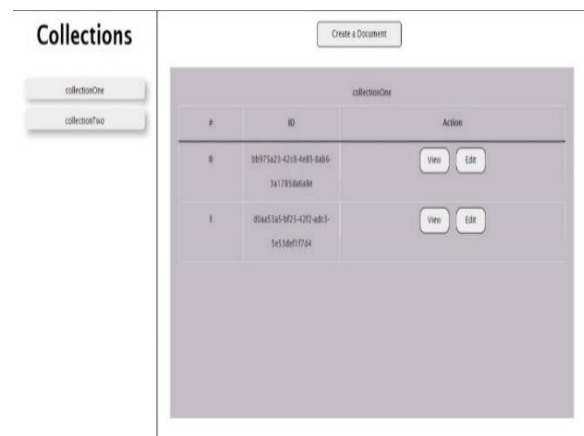


Fig. 3

Birds eye view for all documents in the selected collection. It is a single page all inclusive dashboard for selecting a

collection and viewing, creating and editing documents for the selected collection.

## XIV. CONCLUSION

Thus we have tried to implement the paper "Yexuan Shi, Yongxin Tong, Yuxiang Zeng, Zimu Zhou, Bolin Ding, Lei Chen", "Efficient Approximate Range Aggregation over Large-scale Spatial Data Federation", IEEE 2021 and according to the implementation the conclusion is, a JSON parser to parse document references. A daemon service running in a loop to recursively parse the json received through file system (GUN). A database connector service in adapter fashion with observers looking for change in the file system. (API). As a result, we may modify the weights of relations among nodes (documents) via. Latent Dirichlet allocation for speedier searches. Allowing cycles also makes sure that documents are inherently in BCNF as only 1-1 cardinality is present between 2 distinct documents across the database related using id. Thus an API was created in nodejs for storing interrelated heterogeneous data in a federated manner. Graph based interactive data federation system by "Xuan-Son Vu" was referenced while building the aggregation getter of the library. For data federation efficient approximate range aggregation (eara) was implemented for json object based documents.

## REFERENCES

[1] Yexuan Shi, Yongxin Tong, Yuxiang Zeng, Zimu Zhou, Bolin Ding, Lei Chen, "Efficient Approximate Range Aggregation over Large-scale Spatial Data Federation", IEEE 2021

[2] Xuan-Son Vu, Addi Ait-Mlouk, Erik Elmroth, Lili Jiang, "Graph-based Interactive Data Federation System for Heterogeneous Data Retrieval and Analytics", Association for Computing Machinery, New York, NY, United States.

[3] Rakesh Agrawal and Ramakrishnan Srikant. 1994. Fast Algorithms for Mining Association Rules in Large Databases. In Proceedings of the 20th International Conference on Very Large Data Bases(VLDB '94). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 487-499.

[4] Miika Alonen, Tomi Kauppinen, Osma Suominen, and Eero Hyvönen. 2013. Exploring the Linked University Data with Visualization Tools. In The Semantic Web: ESWC 2013 Satellite Events, Philipp Cimiano, Miriam Fernández, Vanessa Lopez, Stefan Schlobach, and Johanna Völker (Eds.). Springer Berlin Heidelberg, 204-208.

[5] Josep Maria Brunetti, Sören Auer, and Roberto García. 2012. The Linked Data Visualization Model. In Proceedings of the 2012th International Conference on Posters & Demonstrations Track - Volume 914(ISWC-PD'12). Germany, 5-8.

[6] Marija Djokic-Petrovic, Vladimir Cvjetkovic, Jeremy Yang, Marko Zivanovic, and David J. Wild. 2017. PIBAS FedSPARQL: a web-based platform for integration and exploration of bioinformatics datasets. Journal of Biomedical Semantics 8, 1 (20 Sep 2017), 42.

[7] Michael Hahsler and Radoslaw Karpienko. 2017. Visualizing association rules in hierarchical groups. Journal of Business Economics 87, 3 (01 Apr 2017), 317-335.

[8] Martin G. Skjæveland. 2015. Sgvizler: A JavaScript Wrapper for Easy Visualization of SPARQL Result Sets. In The Semantic Web: ESWC 2012 Satellite Events, Elena Simperl, Barry Norton, Dunja Mladenic, Emanuele Della Valle, Irini Fundulaki, Alexandre Passant, and Raphaël Troncy (Eds.). Springer Berlin, 361-365.

[9] Magnus Stuhr, Dumitru Roman, and David Norheim. 2010. LODWheel - JavaScript-based Visualization of RDF Data -. In Proceedings of the Second International Conference on Consuming Linked Data - Volume 782(COLD'11). Germany, 73-84.

[10] Xuan-Son Vu, Lucie Flekova, Lili Jiang, and Iryna Gurevych. 2018. Lexical-semantic resources: yet powerful resources for automatic personality classification. In Proceedings of the 9th Global WordNet Conference. 173-182.

[11] Xuan-Son Vu and Lili Jiang. 2018. Self-adaptive Privacy Concern Detection for User-generated Content. In Proceedings of the 19th International Conference on Computational Linguistics and Intelligent Text Processing, Long papers, p., March 2018.

[12] Xuan-Son Vu, Lili Jiang, Anders Brändström, and Erik Elmroth. 2017. Personality-based Knowledge Extraction for Privacy-preserving Data Analysis. In Proceedings of the Knowledge Capture Conference(K-CAP 2017). ACM, USA, Article 45, 4 pages.