

# Performance Analysis of Elastic Search Technique in Identification and Removal of Duplicate Data

Subhani Shaik, Nallamothu Naga Malleswara Rao

**Abstract:** Elastic search is a way to organize the data and make it easily accessible. It is a server based search on Lucene. It is a highly scalable, distributed and full-text search engine. Elastic search is developed in Java. It is published as open source under the terms of the Apache License. Elastic search is the most popular enterprise search engine. Elastic search includes all advances in speed, security, scalability, and hardware efficiency. Elastic search is a tool for querying written words. It can perform some other smart tasks, but its principal is returning text similar to a given query and statistical analyses of a quantity of text. Elasticsearch is a standalone database server, which is written in Java and using HTTP/JSON protocol, it takes data and optimized the data according to language based searches and stores it in a sophisticated format. Elastic search is very convenient, supporting clustering and leader selection out of the box. Whether it's searching a database of trade products by description, finding similar text in a body of crawled web pages. In this manuscript elastic search capability of copied data identification and its removing techniques performance are analyzed.

**Keywords:** Elastic search, duplicate data identification, data removing techniques, performance analysis.

## I. INTRODUCTION

Assumption is additionally a significant issue when Big Data is acquired utilizing web sources. Despite the fact that searching gives more command over the gathering procedure, there are numerous questions about the connection between the data accessible on the site and data that the site proprietor does not give. Likewise, server issues, arrange load, site refresh strategies, poor page plan, and the non-arbitrary nature of query items are additionally only a portion of the variables that lead to inspecting mistakes when gathering Big Data. Versatile searching is another database worked to deal with tremendous measures of information volume with high accessibility and to convey itself crosswise over numerous machines to be blame tolerant and adaptable, at the same time keeping up a basic yet incredible API that permits applications from any language or system access to the database. Elastic Search was at first written in Java to help free and open source data recovery programming library. Be that as it may, as it turns out, it was very hard to use since it's only a library and expects Java to work with it. Henceforth in ahead of schedule 2000s, a designer names Shay Banon began to work on a deliberation layer over Lucene that made working with scan applications for Java software engineers simpler and named it Compass5.

After a few years, Compass libraries were changed to give ongoing, circulated and superior web search tools. The independent server was discharged with name of Elastic search. The issue of identifying and expelling these copied shards from an archive or database is recognized as shard de-duplication. It is additionally alluded as shard linkage [1], information cleaning [2]. Information de-duplication can be utilized to enhance information quality and honesty, which serves to re-utilization of existing information hotspots for new investigations, and to decrease expenses and endeavors in getting information. In the de-duplication procedure, interesting lumps of information, or byte designs, are distinguished and put away amid a procedure of investigation. As the assessment proceeds, different pieces are contrasted with the duplicate shards and at whatever point a match happens, the repetitive shard is replaced with a little reference that focuses to the original shard. De-duplication is a key activity in coordinating information from numerous sources. The importance of elastic search is confidence web crawler that is fundamentally another product extend called Lucene. It is perhaps most straightforward to comprehend elastic search as a part of base constructed adjacent Lucene's Java libraries. In elastic search, everything is identified with the genuine calculations for coordinating content and putting away advanced files of question terms is executed by Lucene. Elastic search itself gives a more practical and conservative API, versatility, and operational apparatuses overhead Lucene's search usage. Lucene is antiquated in web years, seeing back to 1999. Lucene is illustrated, tried, and is broadly considered best of breed in open-source seek programming. The greater part of the balanced exertion clients of elastic search allot to the assignment of search will be identified with utilizing the Lucene APIs elastic search technique. A data repository is essentially a database which is having accidental replication of shards made among a huge number of information from different sources can scarcely be maintained a strategic distance from others. In the information distribution center group [5], the undertaking of finding copied shards inside information distribution center has for some time been a persistent issue and has turned into a range of dynamic investigation. There are many research endeavors to tackle these issues of information replications caused by copy ruination of information.

## II. LITERATURE SURVEY

The initial step is in charge of the estimation of the area covariance lattices and it is performed utilizing two MapReduce calculations.

Manuscript published on 30 August 2019.

\*Correspondence Author(s)

Nallamothu Naga Malleswara Rao, Professor, Department of IT, RVR & JC College of Engineering, Chowdavaram, Guntur, A.P, India.

Subhani Shaik, Research scholar, Department of CSE, Acharya Nagarjuna University, Nagarjuna Nagar, Guntur, A.P, India.

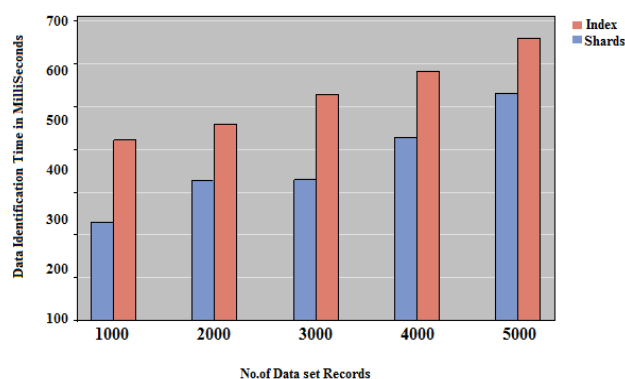
© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an open access article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

The second step plays out the calculation of ellipsoidal  $\tau$ -neighborhoods and it is performed utilizing two other MapReduce calculations. The third step finds center clusters, which is finished by a solitary MapReduce calculation. At long last, the last advance is in charge of the converge of center groups and it is performed with a solitary MapReduce algorithm. In this heuristic, bunches are acquired by sending downright expansive informational collections as indicated by the social investigation approach. The social examination approach gives a scientific formalism where the issue of clustering appears as a direct program with  $n^2$  whole number traits. Heuristics are the most advantageous answer for produce edible clustering results in the quickest time, especially with regards to Big Data, where the quantity of occurrences is vast and the reaction time is a basic factor. Since the first heuristic is consecutive, it should be changed in accordance with the MapReduce display. Flexible search has its very own questioning utilizing JSON called Query DSL. This pursuit can be performed in flexible search in two different ways: Using ansearch or utilizing aprocess forsearch of channel. The basic qualification between them is that asearch finds out and distributes each returned report with the hugeness score while a channel does not do as such. Thus, looking for by methods for channel is speedier than by methods for search. The flexible search documentation recommends using addresses similarly as a piece of two conditions: For full-content endeavors or when the noteworthiness of each result in the interest is basic. For straightforwardness, we will use term request to delineate the two channels and questions; in any case, our contribution with versatile search is compelled to filling in so to speak with channels; in this way, we don't report about usage of inquiries.

### III. PERFORMANCE ANALYSIS

The proposed method is developed in PHP which considers election commission data in which elastic search technique is applied to identify semi and completely copied data and eliminate them from the dataset which reduces the memory wastage.

Elasticsearch is a distributed information holding framework. It can store and bring complex information structures serialized as JSON records progressively [4]. The example in which an archive has been recorded in Elasticsearch, it tends to be recovered from any hub in the cluster. When performing data extraction from the dataset, JSON objects are given to Elasticsearch and result acquired is additionally in JSON position. Elasticsearch is predominantly founded on Apache Lucene. The littlest unit of Elasticsearch is shard and biggest unit is Index. A list can have numerous shards. A shard in Elasticsearch is a Lucene Index. Search process is more accurate when the dataset is divided as shards instead of index. The search process comparison is illustrated in below figure.



**Figure-1: Data Identification comparison levels**

The greatest number of records in a Lucene file is fixed and relies upon the adaptation of Lucene utilized. As of Lucene 5843 the farthest point is 2147483519. When we make an ES record, we can give the quantity of shards it can have. Every shard in itself is an autonomous file and can be facilitated on any hub of the Elasticsearch Cluster. In Elasticsearch, all information in each field is filed as a matter of course. That is, each field has a committed rearranged record for quick recovery. What's more, not normal for most different databases, it can utilize those transformed lists in a similar inquiry, to return results at stunning pace. The information is put away as a reversed list that is advanced for content hunts and in this manner proficient. For instance, in the event that we scan for the word "computer" inside the setting all things considered in Elasticsearch, the time taken is 4060 ms(4.06 s) to locate a sum of 192118 records where the "computer" word is available in record content. The word identification and duplicate data identification is quickly performed when the entire dataset is divided as shards, where as if the dataset is divided as index, then the time for searching and duplicate data identification will be increased.

Aftereffect of catchphrase "computer" from all shards from database

```
{
  "took": 4060,
  "coordinated out": false,
  "shards":
  {
    "absolute": 106,
    "similar": 106,
    "little matched": 0,
    "notmached": 321451
  }
  "hits":
  {
    "absolute": 192118,
    "max_score": 15.110959,
  }
}
```

There are two different ways of executing an essential full-content question: utilizing the Search Lite API, which anticipates that all the pursuit parameters should be passed in as a component of the URL, or utilizing the full JSON body which permits you utilize the full Elasticsearch query. Here is an essential match question that looks for the string "engineering" in every one of the fields: When the data set is divided as shards the query utilized is  
GET/bookdb\_index/engineering/\_search?q=engineering

[Results]

"hits": [

```
{
  "_code": "branch_code",
  "_type": "engineering",
  "_id": "5",
  "_score": 5.12,
  "_source": {
    "title": "computer science and engineering",
  },
  "stream": "engineering",
  "starting date": "20012-04-05",
  "num_students": 120,
  "college": "Nagarjuna University"
},
{
  "_code": "branch_code",
  "_type": "engineering",
  "_id": "4",
  "_score": 4.43,
  "_source": {
    "title": "Electronics and communication engineering",
  },
}
```

Hits :[

```
{
  "_code": "branch_code",
  "_type": "engineering",
  "_id": "4",
  "_score": 4.43,
  "_source": {
    "title": "Electronics and communication engineering",
  },
}
```

Where as if the dataset is divided into index the query utilized is

```
{
  "inquiry":
  {
    "multi_match" :
    {
      "inquiry" : "engineering",
      "fields" : ["code", "type", "id", "score", "source"]
    }
  }
}
```

If the index identification mechanism is initiated then the results are not more accurate when compared to searching of shards.

Elastic search supports fuzzy linked queries also which are applied can be applied on shards along with remaining structures. The fuzzy queries when applied with shards are very accurate and speed is also improved in terms of deduplication identification and removal. The data removing time for different elastic search structures are depicted in below figure.

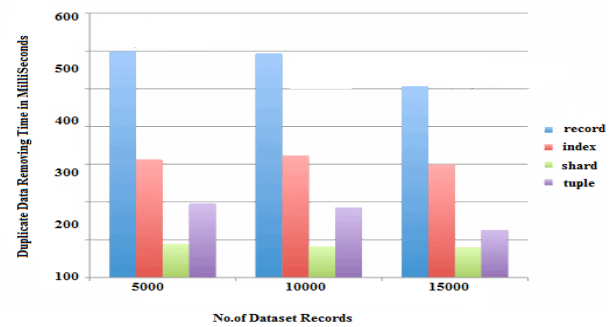


Figure-2: Duplicate data removing time comparison.

Fuzzy query coordinating can be empowered on Match and Multi-Match inquiries to find spelling mistakes, duplicate data and to remove data duplications. The level of fuzziness is determined dependent on the first word, for example the quantity of one-character changes that should be made to one string to make it equivalent to another string.

```
{
  "inquiry": {
    "multi_match" : {
      "inquiry" : "elastic search",
      "fields": ["title", "summary"],

```

```
"fuzzyness": "AUTO"
}
```

```
},

```

```
"_source": ["subject", "topics", "ranking"],
"estimate": 1
}
```

[Results]

"hits": [

```
{
  "_subject": "elastic search",
  "_topic": "computer science",
  "_ranking": "6",
  "_score": 9.54,
}
```

Different parameters are examined with general search techniques and the proposed elastic method exhibits better performance in all aspects. The Grapher software is used to plot the graphs and the values are calculated using the below equations.

The similarity index of shards in the dataset is identified using the equation

$$SIM_{ID} = T(i) \parallel T(i+1)$$

Where T(i) is the probability of matching i<sup>th</sup> shard.

The semi-copied shards are identified by the equation.

$$\frac{1}{2} * SIM_{ID}(I1 \parallel I1+1) + \frac{1}{2} * SIM_{ID}(J1 \parallel J1+1).$$

SC<sub>ID</sub> =

# Performance Analysis of Elastic Search Technique in Identification and Removal of Duplicate Data

Total Shard Count- Count(SIM<sub>ID</sub>)

The fully copied shards are identified by the equation.

$$FC_{ID} = \frac{SC_{ID} \in SC_{ID}(i) \&\& SC_{ID}(n)}{\text{Total Shard Count} - \text{Count}(SIM_{ID})}$$

The CPU utilization is calculated as

$$U=R/C$$

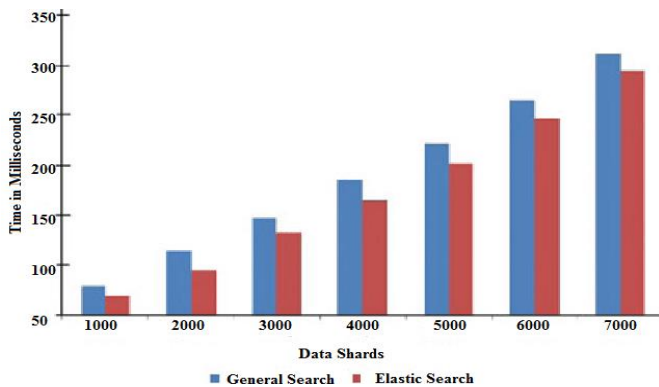
U= Utilization

R= Requirements which in simple terms is the BUSY TIME

C= Capacity which in simple terms is BUSY TIME + IDLE TIME

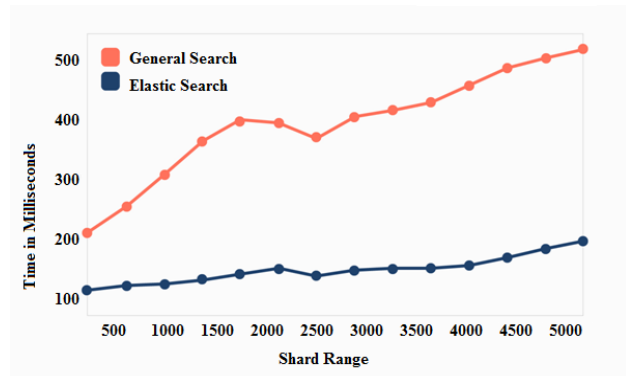
$U = (SC_{ID} / \text{Total shards}) * 1000$ . It is the CPU utilization time for semi-copied shard identification.

$U = (FC_{ID} / \text{Total shards}) - SC_{ID} * 1000$ . It is the CPU utilization time for Fully-copied shard identification.



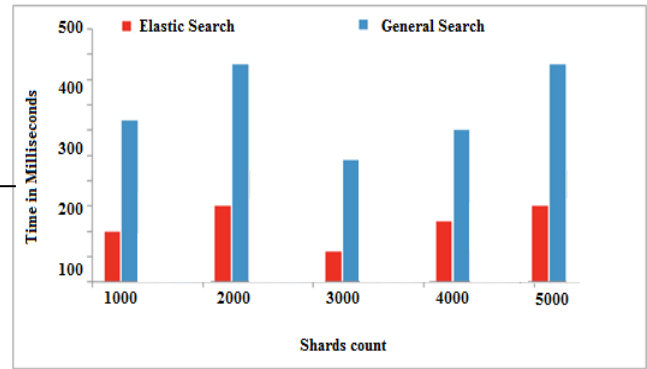
**Figure-3: Shard Loading time Comparison**

The above figure illustrates the shard loading time on proposed and existing methods and the results demonstrate that the shards will be loading much quickly than the traditional method.



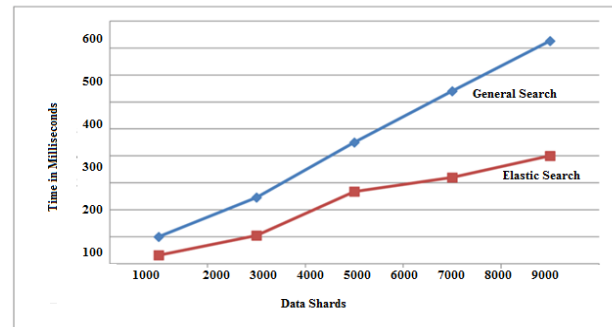
**Figure-4: Semi-copied data identification time levels.**

The above figure demonstrates the process of identification of semi-copied data from the dataset. The semi-copied data is identified in less time in elastic search method.



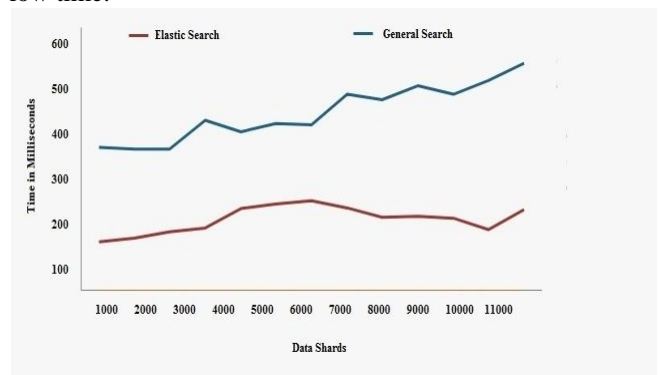
**Figure-5: Fully-copied data identification time levels.**

The above figure demonstrates the process of identification of Fully-copied data from the dataset. The Fully-copied data is identified in less time in elastic search method. The result state that the fully copied data can be identified in a very low time.



**Figure-6: Semi-copied data removing time levels**

The above figure demonstrates the process of removal of semi-copied data from the dataset. The semi-copied data is removed in less time in elastic search method. The result state that the semi- copied data can be removed in a very low time.



**Figure-7: Fully-copied data removing time levels**

The above figure demonstrates the process of removal of Fully-copied data from the dataset. The Fully-copied data is removed in less time in elastic search method.



The result state that the fully copied data can be removed in a very low time. The proposed methodology exhibits better performance and the accuracy levels are also high.

#### IV. CONCLUSION

Elastic search strategy and quality assortment methods are considered for the distinguishing proof of copy information from tremendous information. Shard idea which is utilized to isolate a group into numerous segments and clustering is performed on them and the connected imperatives removes the semi, or complete copy block of information which gives high productivity of the execution of the framework. This manuscript performs analysis on different record structures like shards, indexes for performing elastic search and data duplication removal. Numeric kind of attributes are used to make bundles which are valuable in reducing the amount of examinations. we have proposed replicated shards evacuation procedure, which expels incomplete duplicated information and completely duplicated information from the predetermined dataset. This paper illustrates different parameters and it is identified that the proposed method performs better in all aspects. The proposed strategy utilizes Hadoop technique which is extremely effective in information cleaning process in the information base which reduces memory wastage. The proposed manuscript performs comparison on various structures and the results show that dividing the dataset into shards will improve the elastic search performance and time for removing duplicate data is also reduced when compared to remaining record structures.

#### REFERENCES

1. Verstrepn K, Goethals B. Top-n suggestion for shared records. In: Proceedings of the ninth ACM Conference on Recommender Systems, RecSys'15, New York, NY, 2015, ACM, pp. 59–66.
2. Bender S, Jarmin R, Kreuter F, Lane J. Protection and classification. In: Big Data and Social Science Research: Theory and Practical Approaches. CRC Press, 2016.
3. Hauge M, Stevenson M, Rossmo D, Le Comber S. Labeling banksy: Using geographic profiling to explore a cutting edge craftsmanship riddle. J Spat Sci. 2016;61:185–190.
4. Shmueli G., Yahav I., 2014, "Handling Simpson's Catch 22 in Big Data with grouping and relapse trees", Proceedings of the European Conference on Information Systems (ECIS) 2014, Tel Aviv, Israel, June 9-11, 2014, ISBN 978-0-9915567-0-0
5. Slaoui SC, Lamari Y. Grouping of substantial information dependent on the social investigation. In: Proceedings of the universal meeting on canny frameworks and PC vision. 25–26 March 2015; Fez. Washington: IEEE Computer Society. 2015. p. 1–7.
6. Popov G., Magomedov Sh. Similar investigation of different strategies treatment master evaluations. Universal Journal of Advanced Computer Science and Applications. 2017. T. 8. № 5. C. 35-39. DOI: 10.14569/IJACSA.2017.080505 Elastic Search: Distributed, Lucene-based Search Engine" May 2010, Sematext Blog.
7. O. Baysal, R. Holmes, and M. W. Godfrey. Designer dashboards: The requirement for subjective investigation. IEEE Software, 30(4):46{52, 2013.
8. Long B, Chapelle JB, Zhang Y. Positioning through anticipated misfortune advancement. IEEE Trans Knowl Data Eng 2015;27(5):267-74.
9. Lin J, Ryaboy D, Weil K. Full-content ordering for advancing choice activities in substantial scale information examination. San Jose, California, New York, USA: ACM; 2011. Available from: <http://www.acm.org/publications.Unpluggable Similarity>. Accessible from: <https://www.elastic.co/manage/en/versatile search/direct/current/pluggablesimilarities.htm>.
10. Butler MH, Rutherford J. Appropriated Lucene: A Distributed Free Text Index for Hadoop. HP Laboratories, HPL; 2008. Elastic search invigorate interim versus ordering execution. Accessible from:

<https://sematext.com/preparing/elastic-search/>. Date Accessed: 8/07/2013.8. You know for Search. Accessible from: <https://www.elas-tic.co/blog/you-know-for-see inc>. Date Accessed: 13/07/2016

11. Sai Divya M, Goyal SK. ElasticSearch: A progressed and snappy search strategy to deal with voluminous information. COM-PUSOFT. An International Journal of Advanced Computer Technology. 2013 Jun; 2(6):171–5.
12. Holovaty A, Kaplan-Moss J. The Definitive Guide to Django: Web Development Done Right. Apress, (second edn). 2009 Jul. Django Documentation, Available from: <http://stackoverflow.com/questions/29957604/django-1-8-static-records-doesnt-work>. Date Accessed: 30/04/2015.
13. Django-Haystack Documentation. Accessible from: <https://pi.python.org/pypi/django-bundle/2.4.0>. Date Accessed: 09/06/2015. Understanding Information Retrieval by Using Apache Lucene and Tika Part-1. Accessible from: <https://dzone.com/articles/understanding-data>. Date Accessed: 22/08/2014.
14. Senthil Kumar NK, Kishore Kumar K, Rajkumar N, Amsa-valli K. Website streamlining by Fuzzy Classification and Prediction. Indian Journal of Science and Technology. 2016 Jan; 9(2):1–5.
15. Kausar A, Dhaka VS, Singh SK. Plan of Web Crawler for the Client – Server Technology. Indian Journal of Science and Technology. 2015 Dec; 8(36):1–7.
16. Stephen Mock, "Designsafe: Using Elasticsearch to share and Search Data on Science Web Portal", Practice and Experience in Advanced Research Computing 2017 on Sustainability, Success and Impact, ISBN: 978-1-4503-5272-7, July 13, 2017.
17. Joaquin Delgado, "Scalable Recommender Systems: Where Machine Learning Meets Search", 9th ACM Conference on Recommender Systems, ISBN: 978-1-4503-3692-5, September 20, 2015.
18. Anton Firsov, "Traditional IR meets Ontology Engineering in look for Data", 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, ISBN: 978-1-4503-5022-8, August 2017.
19. Andreas Both, "An administration situated look structure for full content, geospatial and semantic search", 10th International Conference on Semantic Systems, ISBN: 978-1-4503-2927-9, September 2014.
20. Shelly Garion, "Big Data Analysis of distributed storage logs utilizing sparkle", tenth ACM International Systems and Storage Conference, ISBN: 978-1-4503-5035-8, May 2017.
21. Bela Gipp, "Citolytics: A Link-based Recommender System for Wikipedia", 11th ACM Conference on Recommender Systems, ISBN: 978-1-4503-4652-8, August 2017.

#### ABOUT THE AUTHORS



**Mr. Subhani Shaik** is working as Assistant professor in Department of computer science and Engineering at St. Mary's group of institutions Guntur, he has 12 years of Teaching Experience in the academics.



**Dr. Nallamothu Naga Malleswara Rao** is working as Professor in the Department of Information Technology at RVR & JC College of Engineering with 25 years of Teaching Experience in the academics.