



International Journal of Advanced Trends in Computer Applications

www.ijatca.com

Big Data Compression for Aadhaar Storage through Reduced Bit Level Ordering

¹Kulanthaivel. G, ²Ezhilarasu. P, ³Ulagamuthalvi. V

¹Professor, NITTTR, Chennai,

Tamil Nadu, India.

²Associate Professor,

Department of Computer Science and Engineering, SRMIST, Uttar Pradesh, India.

³Associate Professor,

Sathyabama Institute of Science and Technology,

Chennai, Tamil Nadu, India.

¹gkveldr@gmail.com, ²prof.p.ezhilarasu@gmail.com, ³ulagamv@gmail.com

Abstract: In this paper explores on compression of Aadhaar number storage through the concepts of reduced bit level ordering. The Aadhaar number is an unique number used in various government schemes. It is a twelve digit number. The population of our country is more than 135 crores. It means more than 135 crores of Aadhaar numbers will be available. The uniqueness of beneficiary is ensured by avoiding duplicate beneficiary. The uniqueness needs Aadhaar number comparison with minimal amount of time. Hence there is a need for reducing the Aadhaar storage for minimizing the search time. The numbers are represented by using various number representations with their needed storage. The compression of Aadhaar storage implemented by using the concept of bit level ordering which takes sorted integer as an input. The expected saved space shown in graphs and tables. The input from various bit level taken and output obtained as index bit and content bit. The obtained saved space depicted through table and graphs.

Keywords: Big data, Compression, Aadhaar, Bit level ordering, Space reduction.

I. INTRODUCTION

1.1 Big Data

The impact of social media leads to exponential growth of digital data. It is growing exponentially daily. This collection of digital data named as Big Data. The data are categorized as structured data, semi structured data and unstructured data. Incase of structured data sometime it is possible to process the data manually. But in case of semi structured and unstructured data we need an analytic tool for processing of data. The prime objective of big data analysis is to process data of high volume, velocity, variety, and veracity using various traditional and computational intelligent techniques [1].

1.2 Compression

The system of reducing the size of a data file referred to as data compression [2]. Data compression uses the trade-off in the form of space and time complexity. If the data stored in its original form, then there is no need to perform compression and decompression of the taken input data which require huge storage. By applying compression space can be reduces. This efficiently minimize space require and time for transfer. Data transfer rate is calculated by eq.1.

Data transfer time = Input data /

Data transfer rate (1)

Data transfer time required depend on the data rate and size of data. Depend on the storage of the destination files will be stored in storage. The space and time

complexity based on compression ratio. It can be calculated by using the following equation 2.

$$\text{Compression ratio} = \frac{\text{Actual uncompressed data}}{\text{Compressed Data}} \quad (2)$$

The data compression also has specific limitations for compression and decompression. If processing time is more then the need for performing compression practically not needed. Limitations can be Compression and Decompression cost, time and quality of data. Lossy and lossless compression can be used.

1.3 Aadhaar

The United Nations, Department of Economic and Social Affairs, Population Division released world population prospects 2019[3]. According to this report the population of India is 136.6 crores. Aadhaar is a verifiable 12-digit identification number issued by UIDAI to the resident of India for free of cost [4]. This number is similar to nine digit Social Security number issued by Social Security Administration in United States of America. The Aadhaar number used in many government schemes to identify beneficiary. In many bank Aadhaar enabled Payment System (AePS) is a bank led model which allows online interoperable financial inclusion transaction at PoS (MicroATM) through the Business correspondent of any bank using the Aadhaar authentication [5]. The fake beneficiary removed from the database using Aadhaar number in many government schemes. This helps to strengthen the economy of India. This Aadhaar number is based on population of India. The population of India is more than 135 crores means the available Aadhaar number is more than 135 crores. The Aadhaar number has many details like Enrolment number, Aadhaar number, Date of Birth, Address, QR code, Gender. The state and central government has many schemes based on Aadhaar number.

II. RELATED WORK

Claude Elwood Shannon who is known as father of Information Theory discussed about discrete random variable and its information content [6]. Inverted indexes, which is the most common data structure used by search engines to index their data, are made of lists of increasing integers corresponding to the documents of the collection[7]. Importance of security to save data in big data is done by various techniques [8,9]. If we compare Elias-Fano encoding space requirement with the theoretical lower bound we realize that this

structure is close to the bound, so it has been epithet quasi-succinct index [10]. The Elias-Fano integer encoding of monotone sequences of integers solved using succinct data structures [11, 12, 13]. Compression of inverted indices saves disk space, but more importantly also reduces disk and main memory accesses [14], resulting in faster evaluation. An inverted index is a collection of sorted sequences of integers [15, 16, 17]. Compressing such sequences is a crucial problem which has been studied since the 1950s; the literature presents several approaches, each of which introduces its own trade-off between space occupancy and decompression speed [20, 21, 22, 23]. Elias-Fano coding for various length of sorted integers shows that input with 3 bits length (1 highest bit) cannot produce compression. The 10 bit highest bit provides 73% of space savings. Hence for storing 100 bits we need only 27 bits. The increased number of highest bit (>10 bits) results in more amount of space savings (>73%)[21].

III. AADHAAR NUMBER REPRESENTATION

The Aadhaar number can be stored using the number representation methods such as

- ASCII CODE REPRESENTATION
- NORMAL BINARY REPRESENTATION

3.1 Number Representation

The ASCII code range is from 0 to 255. So we need eight digits ($2^8 = 256$) for representing each number. The numbers from 0 to 9 represented in ASCII using the numbers 48 to 57. In normal binary representation the maximum number 9 has the size 4. So we need four digits for representing each number. The number representation for both number representations is given in the table 1.

Table 1: ASCII code representation for numbers (0-9)

S.NO	NUMBERS	ASCII REPRESENTATION	NORMAL BINARY REPRESENTATION
1	0	00110000	0000
2	1	00110001	0001
3	2	00110010	0010

4	3	00110011	0011
5	4	00110100	0100
6	5	00110101	0101
7	6	00110110	0110
8	7	00110111	0111
9	8	00111000	1000
10	9	00111001	1001

3.2 Storage Need

The storage needed for storing 135 crores Aadhaar numbers using the above representation is given in Table 2 and Figure 1.

Table 2: Needed storage for 135 crores Aadhaar numbers

S.NO	REPRESENTATION	EXAMPLE AADHAAR NUMBERS	TOTAL DIGITS AFTER PAIRING	SIZE OF INDIVIDUAL PAIRS	SIZE OF SINGLE AADHAAR NUMBER	SIZE OF 135 CRORES AADHAAR NUMBERS
1	ASCII	9999 9999 9999	12	8	96	15.08 GB
2	NORMAL BINARY	5555 5555 5555	12	4	48	7.54 GB
3	2-COUPLE	99 99 99 99 99 99	6	7	42	6.60 GB
4	3-TRIPLE	453 453 453 453	4	10	40	6.28 GB
5	4-QUADRUPLE	1060 1060 1060	3	14	42	6.60 GB
6	6-SEXTUPLE	555555 555555	2	20	40	6.28 GB

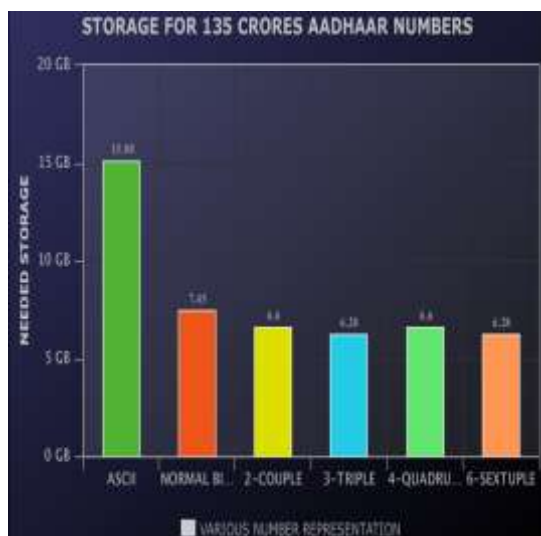


Figure 1: Needed storage for 135 corers Aadhaar numbers

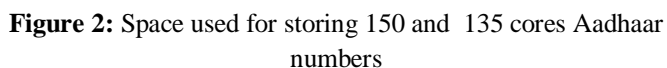
IV. AADHAAR COMPRESSION THROUGH BIT LEVEL ORDERING

The possible range of Aadhaar number is from 000000000000(12 Zero's) to 999999999999(12 Nine's). The available Aadhaar numbers are sorted in ascending order. Then they are compressed. The compression of sorted integers is implemented using eliasfano coding by splitting the numbers into high and low level bits. In this work only sorted integers are used. Instead of splitting as high level and low level, here bit level ordering is used. In bit level ordering the bit level and its occurrence are used. The twelve digit Aadhaar number is considered as decimal number. Then decimal to binary number conversion is performed. The size of binary number is termed as level. Ex. 123412341234 is converted into binary as

contents of index and storage bits is calculated and storage decided according to that. Space saved for storing all Aadhaar card is shown in table 3. The range from 000000000000-001349999999 denotes first 135 cores possible Aadhaar numbers. The needed bit for this range using reduced bit level ordering is 38352516324 bits. Average bits used to save are 28.41 and 39 bits after bit ordering.

If the total number of Aadhaar numbers are 150 and there consolidated count in reduced bit level are as given in the table 3, then it is possible to store and search Aadhaar numbers using index bit and storage bit with

S.NO	NUMBER REPRESENTATION	CHARACTER TYPE	CONSOLIDATED TOTAL BIT	SAVED SPACE
1	ASCII	Fixed - 8 bits	$150 * 96 = 14400$	0%
2	NORMAL BINARY	Fixed - maximum 4 bits	$150 * 48 = 7200$	50%
3	2-COUPLE	Fixed - maximum 7 bits for 2 characters	$150 * 42 = 6300$	56.25%
4	3-TRIPLE	Fixed - maximum 10 bits for 3 characters	$150 * 40 = 6000$	58.33%
5	4-QUADRUPLE	Fixed - maximum 14 bits for 4 characters	$150 * 42 = 6300$	56.25%



we need to search from 2470to 2506 then 2507 to 2543.
If match found then the number is present else the
given number is not present. Here the given number
137438953490(00000000000000000000000000000000
10100) matches the bit from 2507 to 2543. Hence the
given number is present. It is explained below.

10011000101111000000100101101111110000000010
11001000000000001000010000011001000011011111
0001110100000000000011000011000100100010100
00111011101000110100111100101110100111111100
0000000110010010011100101011001010110111111
1011111111000000000000010100001010000010110
1001011111010010001011000011011110101000000
0000000101010011011000101111011100111111111
0000000000000000110110111101000101111111111
0000000000000000000001100000000011100000110
1010001110001010010001010010000000000000111
11111111110000100110101000011001010010110101
1001011000000000000000000000011101010101011
011011010111111111111111000000000000000000
00011101001001010111110000100101111110010010
11111011110100100000000111010001001011111100
100010011111001101110101111111111011100100000

[illegible]

Information Sciences and Systems: 54-47, Princeton University, 1972.

[12] P. Elias, "Efficient storage and retrieval by content and address of static files", J. ACM 21(2): 246-260, 1974.

[13] R. M. Fano, "On the number of bits required to implement an associative memory", Memorandum 61, Computer Structures Group, Project MAC, Massachusetts Institute of Technology, 1971.

[14] S. Büttcher and C. L. A. Clarke. Index compression is good, especially for random access. In M. J. Silva, A. H. F.

[15] S. Buttcher, C. L. A. Clarke, and G. V. Cormack. "Information retrieval. "Implementing and evaluating search engines", MIT Press, Cambridge, Mass., 2010.

[16] C. D. Manning, P. Raghavan, and H. Schulze, "Introduction to Information Retrieval", Cambridge University Press, 2008.

[17] J. Zobel and A. Moffat, "Inverted files for text search engines", ACM Comput. Surv., 38(2), 2006.

[18] D. Lemire and L. Boytsov, "Decoding billions of integers per second through vectorization", Software: Practice & Experience, 2013.

[19] A. Moffat and L. Stuiver, "Binary interpolative coding for effective index compression", Inf. Retr., 3(1), 2000.

[20] D. Salomon, "Variable-length Codes for Data Compression", Springer, 2007.

[21] A. A. Stepanov, A. R. Gangolli, D. E. Rose, R. J. Ernst, and P. S. Oberoi. "Simd-based decoding of posting lists", In CIKM, pages 317–326, 2011.

[22] H. Yan, S. Ding, and T. Suel, "Inverted index compression and query processing with optimized document ordering", In WWW, pages 401–410, 2009.

[23] Ezhilarasu P, Mahapatra RP, Senthil R, "Analysis and Interpretation of Elias-Fano coding for Sorted Integers", LAP LAMBERT Academic Publishing, 2019.

Ph.D. degree in Information and Communication Engineering from Anna University, Chennai. He completed his Master's degree in Microwave and Optical Engineering from Madurai Kamaraj University and Bachelor's Degree in Electronics and Communication Engineering from University of Madras. He is having experience of more than 27 years out of which more than 24 years in training of technical teachers in India and abroad. His area of interest includes Telemedicine, Image Processing, Communication, IoT, Virtual Instrumentation, ICT applications in Teaching/Learning. He has published/presented many papers in the National/International Journals/Conferences.



Dr. P. Ezhilarasu received the B.E degree in Computer Science and Engineering from Bharathiar University Coimbatore. Then M.E degree in Computer Science and Engineering from Anna University, Chennai. He commended his Ph.D in Faculty of Information and Communication from Anna University Chennai. He is having two decades of academic experience. He guided so many UG and PG students in his career. He handled more than 30 subjects in his career and also having more than 80 international/national publication in conference/journal. His area of interest includes Image Processing, Theory of Computation and Data science.



Dr. V. Ulagamuthalvi, Associate Professor, department of Computer Science and Engineering, Sathyabama Institute of Science and Engineering received her Ph.D degrees in Medical Image Processing from sathyabama University. She is having experience of more than 20 years in teaching. Her area of interest are Data Mining, Image Processing, Big Data Analysis.

Author Profile



Dr. G. Kulanthaivel, Professor & Head of the Department of ECE, NITTTR, Chennai, India received his