

---

# Examine Dynamic Data Structures and Complex Algorithms

---

**Yiqiao Yin**  
Ph.D. Student

## Abstract

1        This report examines the famous file sharing algorithms used in dis-  
2        tribution systems. We investigate the architecture of these algorithms  
3        and compares how they are deployed. In addition, this report also  
4        exposes some of the potential issues with the algorithms. Examples  
5        and applications are presented to illustrate the survey of the algorithms  
6        in this report.

## 7    1    Introduction

8        The search on similarity of two documents is a well established and import subject in  
9        computer science (Wang et al., 2014; Enbody and Du, 1988; Chi and Zhu, 2017; Gui  
10       et al., 2017). Amongst many searching algorithms, hashing becomes famous and rise to  
11       popularity because of the level of simplicity at deployment and the low cost in regards  
12       of time and space analysis. Throughout the years, hashing techniques have been evolved  
13       overtime and many different upgrades have been developed (Chi and Zhu, 2017). The  
14       hashing techniques are desired to retrieve information from large volume of documents  
15       and texts for post process analysis and it is widely used in day-to-day work in computer  
16       science. When deep learning was developed in the late 1980s, hashing algorithms  
17       have also been used to encourage promising work in machine learning since its design  
18       provide crucial guideline for scholars to understand the how searching works from  
19       programming perspective. In large-scale data analysis and today's BIGDATA concept,  
20       hashing provides the most fundamental understanding of information representation and  
21       mapping system between data.

22       Clustering technique is another important field of practice in computer science (Saxena  
23       et al., 2017; Rai and Singh, 2010). In this report, different clustering techniques are  
24       discussed and a variety of grouping techniques are covered. There is a clear motivation  
25       in today's computer science field work to understand and also develop the advancement  
26       of different clustering for a variety of different purposes in mathematics, engineering,  
27       computer science, statistics, and so on. Taking a hospital scenario as an example.  
28       Patients with different diseases are, evidently, treated differently. This is because  
29       different diseases have different symptoms and hence requires different medicine. A  
30       good clustering technique can help us navigate the unsupervised world with orders  
31       instead of chaos.

## 32 2 Description of Different Architectures in Hashing

33 This section describes the different architectures of hashing techniques.

34 The first fundamental technique to introduce is the hashing function. It is a mathematical  
35 function that maps information and text to a list of dictionary of keys. To express the  
36 idea better, we will address the input of the hash functions as keys and the output of  
37 the hash functions as hash values. In other words, an abstract formula can be expressed  
38 using a hash function  $h(\cdot)$

$$h : \text{keys} \rightarrow \text{hash values} \quad (1)$$

39 An important property for this setup is that it utilizes statistical concept of how functions  
40 interact with each other. In addition, the key functionality of a hash function is the  
41 capability of searching for keys and values even when the information is scrambled  
42 because the hash function ensures that the resulting values from the function remains  
43 uniformly distributed when generating the output. This is to avoid collision. A collision  
44 in a hash function refers to a scenario where two different keys are mapped to one single  
45 hash value. Hence, the important properties of a good hash function can be summarized  
46 below

- 47 1. the hash function is efficient and easy to compute;
- 48 2. the hash function is capable of minimizing duplicated scenarios.

### 49 2.1 Properties of Hashing

50 From earlier section, it is discussed that a good hashing function has expected inputs to  
51 be spread over the output range as evenly as possible. This is the uniformity property.  
52 Like the uniform distribution phenomenon in statistics, every hash value ideally has the  
53 same probability. A uniform probability distribution states that a random variable  $x$  can  
54 take values from a range  $[a, b]$  while each value has the same probability. In other words,  
55 we can state this formally in the following. Assume that there is support of a random  
56 variable  $x$  to be  $\mathbb{R}_x = [a, b]$ . The probability density function of a uniform random  
57 variable takes the following form

$$f_X(x) = \frac{1}{b-a} \text{ if } x \in \mathbb{R}_x \text{ and } f_X(x) = 0 \text{ elsewhere} \quad (2)$$

58 This property can be tested and the test for uniformity is called Chi-square test (Rao,  
59 1972; Inglot and Janic-Wróblewska, 2003). The goal is to produce a goodness-of-fit  
60 measure using Chi-square test. The key here is to measure whether the actual distribution  
61 is close to the expected distribution. The closer they are, the higher confidence computer  
62 scientists have to claim that the actual distribution is uniform. To conduct this test, the  
63 following formula is used

$$\frac{n/2m}{n+2m-1} \sum_{j=0}^{m-1} (b_j)(b_j+1)/2 \quad (3)$$

64 where  $n$  is the number of keys,  $m$  is the number of buckets,  $b_j$  is the number of items  
65 in bucket  $j$ . Commonly it is desirable to have the test value to be 0.95-1.00 so that it  
66 gives computer scientists confidence that the actual distribution is uniform (Castro et al.,  
67 2005).

## 68 2.2 Discussion of Architecture, Coding, and Implementation

69 This section provides some demonstration of the architecture of hashing. Since hashing  
70 function is a technique to solve a particular data structure question, there is no universal  
71 hashing algorithm to solve all problems. Hence, this subsection will use example-based  
72 approach to explain the architecture of the hashing algorithm, the coding style, and the  
73 implementation.

74 The first example is to search if one list of items is inside another list of items. It is one  
75 of the most common data structure problems that need to be solved. Suppose there are  
76 two lists and each of them has certain number of digits. Denote them list 1 and list 2.  
77 One simple goal is to check if one list has all of its items inside the other. This type of  
78 task is called "sub-list search" and it is a common task abstracted to solve different kinds  
79 of computer programming problems.

80 A sample syntax is presented below. The syntax creates a function called "is\_this\_subset"  
81 and the function takes two inputs. The two inputs can be two lists or two arrays. The  
82 algorithm presents a nested for loop or double for loop. The first for loop searches  
83 through each item in the second array while the nested loop or the inside loop searches  
84 for each item in the first array. Each of the inside loop checks if the second array is  
85 inside the first and breaks if the condition does trigger. Then the algorithm checks if the  
86 running index  $j$  is the same with the length of first array. The algorithm returns negative  
87 or "0" value if the condition triggers. In the end, if everything checks, the algorithm  
88 returns positive or "1".

```
89
90 def is_this_subset(the_first_array, the_second_array):
91     i = 0
92     j = 0
93
94     l1 = len(the_first_array)
95     l2 = len(the_second_array)
96
97     for i in range(l2):
98         for j in range(l1):
99             print(i, j)
100             if(the_second_array[i] == the_first_array[j]):
101                 print(i, j, the_second_array[i], the_first_array[j]
102                     ↪ ] )
103                 break
104
105         # the above nested loop can break
106         # when the condition triggers
107         # else it runs the following
108         # then returns 0 if j is the
109         # same as l1
110         if (j == l1):
111             return 0
112
113     # when the nested for loops
114     # finishes above
115     # we can wrap things up here
116     return 1
117
118 the_first_array = [11, 4, 452, 6, 213, 3]
119 the_second_array = [11, 213, 4, 3]
```

```
120 |  
121 | print(is_this_subset(the_first_array, the_second_array))
```

### 123 3 Clustering Techniques

124 This section covers the fundamentals of clustering techniques. The basic purpose for  
125 clustering techniques is to divide data into different groups based on similarity. For each  
126 group, there are certain number of items from the original data and these items form a  
127 cluster. Inside of this cluster, the items are similar to each other more than they are to the  
128 rest of the items. There can be many building blocks inside of the clustering algorithm  
129 and many places require their own tuning parameters range. The key idea for most  
130 clustering techniques is to use a distance function to measure how far one data point  
131 is from the other. If the distance function produces a value that is high, then it implies  
132 that the two data points are far away. This distance value can be used as a comparison  
133 purposes amongst all data points. The most basic distance function is the Euclidean  
134 distance and it takes the following form

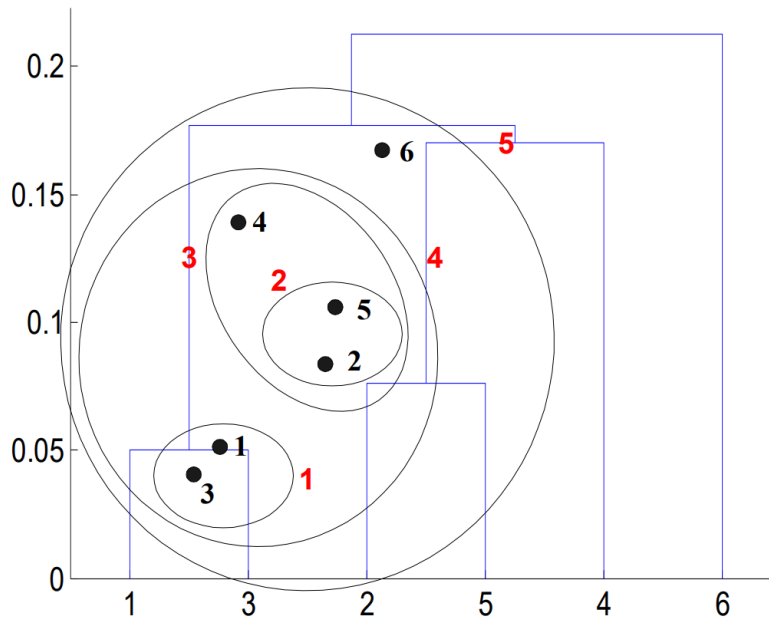
$$d(x, y) = \sqrt{\sum_i (x_i - y_i)^2} \quad (4)$$

135 where  $x$  and  $y$  are two attributes in the data and the running index  $i$  tracks each instance  
136 (or sample) in the data. There are a few more famous distance functions desired for most  
137 clustering techniques and we list them below

$$\begin{array}{ll} \text{Euclidean:} & d(x, y) = \sqrt{\sum_i (x_i - y_i)^2} \\ \text{Squared Euclidean:} & d(x, y) = \sum_i (x_i - y_i)^2 \\ \text{Manhattan:} & d(x, y) = \sum_i |x_i - y_i| \end{array} \quad (5)$$

138 There are many algorithms surveyed in clustering techniques (Rai and Singh, 2010).  
139 They come down to two original ideas. The first one is nearest neighbor method and the  
140 second one is the hierarchical method. Both methods can be demonstrated in the Figure  
141 1. The nearest neighbor method uses distance functions to measure how far a data point  
142 is from the rest of the data points. The data points with closer distances are grouped  
143 together. The second method is hierarchical method. A famous hierarchical method is  
144 dendrogram. The algorithm investigates the pair-wise distances between data points and  
145 assign them to different branches accordingly. The branch accelerates in levels when  
146 the distances get past certain thresholds where these thresholds are numerical tuning  
147 parameter.

Figure 1: **Clustering Techniques.** Famous clustering techniques can be seen in this diagram. By using two attributes, the data points can be grouped into different clusters using nearest neighbors or hierarchical methods. The graph is an adaption from Figure 4 of this research (Rai and Singh, 2010).



#### 148 **4 Deployment**

149 This section briefly discusses the crucial steps required for deploying a hashing algorithm  
 150 in practice. The stages and levels of deployment depends on the existing system and  
 151 pipeline in place.

152 First, it is likely that the existing system in place consists of some of hashing techniques  
 153 and the system works smoothly and efficiently to a certain level. It is up to the computer  
 154 scientists to evaluate the system and the business requirement to decide whether the  
 155 current system needs to be replaced or upgraded.

156 Second, the scenario can be that there is no existing system in place. In this case, there  
 157 is more freedom of what hashing algorithm to be put in place. A new algorithm will be  
 158 proposed and certain measurements need to be taken.

159 In regards to measurements, our previous work investigated the Big-O analysis on  
 160 different algorithms. This can be taken as the major measurement to assess and evaluate  
 161 the performance of different hashing algorithms. It can certainly be used to provide  
 162 crucial guidance whether the current system needs to be replaced. It is recommended  
 163 and certainly desired to investigate the time and space consumption using Big-O notation  
 164 before any upgrade or novel proposal of new hashing algorithm.

#### 165 **5 Potential Issues**

166 Potential issues can occur when the current or the proposed hashing algorithms do not  
 167 satisfy the business moat. Every business has an economic moat and the range of such

168 moat expands as the business scales up in size. Different consumer groups will arise and  
169 the business will have to adjust its operation activities accordingly. This is where data  
170 type and data structure changes.

171 Whenever there is a change on data type or structure, it is absolutely necessary to  
172 evaluate the time and space consumption of the current hashing algorithms. If the current  
173 benchmark does not meet the requirements anymore, it would have to be replaced or  
174 upgraded.

## 175 **References**

176 Castro, J. C. H., Sierra, J. M., Sez nec, A., Izquierdo, A., and Ribagorda, A. (2005). The  
177 strict avalanche criterion randomness test. *Mathematics and Computers in Simulation*,  
178 68(1):1–7.

179 Chi, L. and Zhu, X. (2017). Hashing techniques: A survey and taxonomy. *ACM*  
180 *Computing Surveys (CSUR)*, 50(1):1–36.

181 Enbody, R. J. and Du, H.-C. (1988). Dynamic hashing schemes. *ACM Computing*  
182 *Surveys (CSUR)*, 20(2):850–113.

183 Gui, J., Liu, T., Sun, Z., Tao, D., and Tan, T. (2017). Fast supervised discrete hashing.  
184 *IEEE transactions on pattern analysis and machine intelligence*, 40(2):490–496.

185 Inglot, T. and Janic-Wróblewska, A. (2003). Data driven chi-square test for uniformity  
186 with unequal cells. *Journal of Statistical Computation and Simulation*, 73(8):545–561.

187 Rai, P. and Singh, S. (2010). A survey of clustering techniques. *International Journal of*  
188 *Computer Applications*, 7(12):1–5.

189 Rao, J. (1972). Some variants of chi-square for testing uniformity on the circle. *Zeitschrift*  
190 *für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 22(1):33–44.

191 Saxena, A., Prasad, M., Gupta, A., Bharill, N., Patel, O. P., Tiwari, A., Er, M. J., Ding,  
192 W., and Lin, C.-T. (2017). A review of clustering techniques and developments.  
193 *Neurocomputing*, 267:664–681.

194 Wang, Q., Si, L., Zhang, Z., and Zhang, N. (2014). Active hashing with joint data  
195 example and tag selection. In *Proceedings of the 37th international ACM SIGIR*  
196 *conference on Research & development in information retrieval*, pages 405–414.