

In most general applications, we cannot know exactly what set of key values will need to be hashed until the hash function and table have been designed and put to use.

At that point, changing the hash function or changing the size of the table will be extremely expensive since either would require re-hashing every key.

A *perfect hash function* is one that maps the set of actual key values to the table without any collisions.

A *minimal perfect hash function* does so using a table that has only as many slots as there are key values to be hashed.

If the set of keys **IS** known in advance, it is possible to construct a specialized hash function that is perfect, perhaps even minimal perfect.

Algorithms for constructing perfect hash functions tend to be tedious, but a number are known.

This is used primarily when it is necessary to hash a relatively small collection of keys, such as the set of reserved words for a programming language.

The basic formula is:

$$h(S) = S.length() + g(S[0]) + g(S[S.length() - 1])$$

where $g()$ is constructed using Cichelli's algorithm so that $h()$ will return a different hash value for each word in the set.

The algorithm has three phases:

- computation of the letter frequencies in the words
- ordering the words
- searching

Suppose we need to hash the words in the list below:

```
calliope
clio
erato
euterpe
melpomene
polyhymnia
terpsichore
thalia
urania
```

Determine the frequency with which each first and last letter occurs:

letter:	e	a	c	o	t	m	p	u
freq:	6	3	2	2	2	1	1	1

Score the words by summing the frequencies of their first and last letters, and then sort them in descending order:

calliope	8
clio	4
erato	8
euterpe	12
melpomene	7
polyhymnia	4
terpsichore	8
thalia	5
urania	4

euterpe
calliope
erato
terpsichore
melpomene
thalia
clio
polyhymnia
urania

Finally, consider the words in order and define $g(x)$ for each possible first and last letter in such a way that each of the words will have a distinct hash value:

word	g_value assigned	h(word)	table slot
euterpe	e-->0	7	7 ok
calliope	c-->0	8	8 ok
erato	o-->0	5	5 ok
terpsichore	t-->0	11	2 ok
melpomene	m-->0	9	0 ok
thalia	a-->0	6	6 ok
clio	none	4	4 ok
polyhymnia	p-->0	10	1 ok
urania	u-->0	6	6 reject
	u-->1	7	7 reject
	u-->2	8	8 reject
	u-->3	9	0 reject
	u-->4	10	1 reject

Cichelli's method imposes a limit on the search at this point (we're assuming it's 5 steps), and so we back up to the previous word and redefine the mapping there:

word	g_value assigned	h(word)	table slot
polyhymnia	p-->0	10	1 reject
	p-->1	11	2 reject
	p-->2	12	3
urania	u-->0	6	6 reject
	u-->1	7	7 reject
	u-->2	8	8 reject
	u-->3	9	0 reject
	u-->4	10	1 ok

So, if we define $g()$ as determined above, then $h()$ will be a minimal perfect hash function on the given set of words.

The primary difficulty is the cost, because the search phase can degenerate to exponential performance, and so it is only practical for small sets of words.

A duplex hash strategy to achieve worst case searching of $\Theta(1)$ using the power of choice.

Two tables (size = # items) are employed each with a separate hash function. A key is hashed by the two different functions and will always be located in the first or second table. Thus only two lookups are ever required to find an item.

Hash Fn table 1: $h(k) = k \bmod m$
 where $m = \text{table size}, (8)$

Hash Fn table 2: $g(k) = 1 + k \bmod (m-1)$,

<i>T1</i>	36	22	84	68	17	14	13	65
0								
1								
2								
3								
4	36		84					
5								
6		22						
7								

Standard Cuckoo hashing does not automatically check table 2 when a collision occurs in a table 1 insert.

<i>T2</i>	36	22	84	68	17	14	13	65
0								
1			84					
2								
3								
4								
5								
6								
7								

On the third insert, 84 collides with 36 in table 1. Thus it is inserted in table 2 where no collision occurs.

The fourth & fifth inserts cause no displacements. On the sixth insert, 14 collides with 22 in table 1. When 14 is inserted in table 2 it displaces 84 sending it back to table 1 which results in 36 being displaced from table 1. The displacements end when 36 is inserted into table 2.

Hash Fn table 1: $h(k) = k \bmod m$
 where $m = \text{table size}, (8)$

Hash Fn table 2: $g(k) = 1 + k \bmod (m-1)$,

<i>T1</i>	36	22	84	68	17	14	13	65
0								
1					17			
2								
3								
4	36		84	68				
5								
6		22				14		
7								

<i>T2</i>	36	22	84	68	17	14	13	65
0								
1			84			14		
2	36							
3								
4								
5								
6				68				
7								

The last two inserts result in no displacements.

Hash Fn table 1: $h(k) = k \bmod m$
 where $m = \text{table size, } (8)$

Hash Fn table 2: $g(k) = 1 + k \bmod (m-1)$,

$T1$	36	22	84	68	17	14	13	65
0								
1					17			65
2								
3								
4	36		84	68				
5							13	
6		22				14		
7								

$T2$	36	22	84	68	17	14	13	65
0								
1			84			14		
2	36							
3								65
4								
5								
6				68				
7								

Table displacements may result in a cycle which must be detected.

If each table load factor is $< 50\%$ then the probability of a cycle is quite small, with a small constant of displacements and most insertions will need only $O(\log N)$ displacements.

Thus tables may need to be rebuilt with different hash functions if a number of displacements ($\log N + \text{very small constant}$) occurs during insertion.

key	$h(k)$	$g(k)$
36	4	2
22	6	2
84	4	1
68	4	6
17	1	4
14	6	1
13	5	7
65	1	3

Cuckoo hashing requires a set of hash Fns, (many standard hash Fns perform poorly in cuckoo hashing). Cuckoo hashing guarantees worst-case constant lookup, trivial deletion and constant insertion if the load factor, $\lambda < 50\%$. The expected insertion cost bound is:

$$\frac{1}{\left(1 - (4\lambda^2)^{1/3}\right)}$$