



**Martin is getting
the projector
to work with
his laptop.**

Hashed Collections

**You Can
(and sometimes should)
Build Your Own**

Hashed Collections

It slices, it dices...

Martin McClure



Andrés Valloud

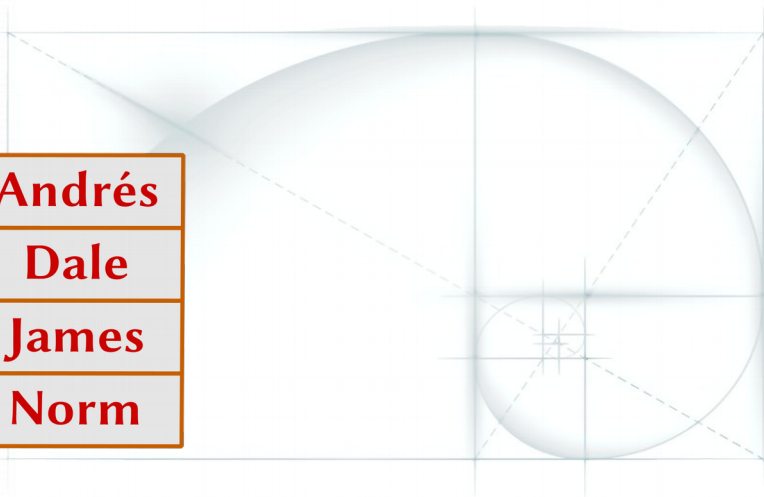


Andrés

Dale

James

Norm





Andrés

Dale

James

Norm

Linear search $O(n)$



Andrés

Dale

James

Norm

Linear search $O(n)$
Binary search $O(\lg n)$



Andrés
Dale
James
Norm

Linear search $O(n)$
Binary search $O(\lg n)$
 $O(1)$?



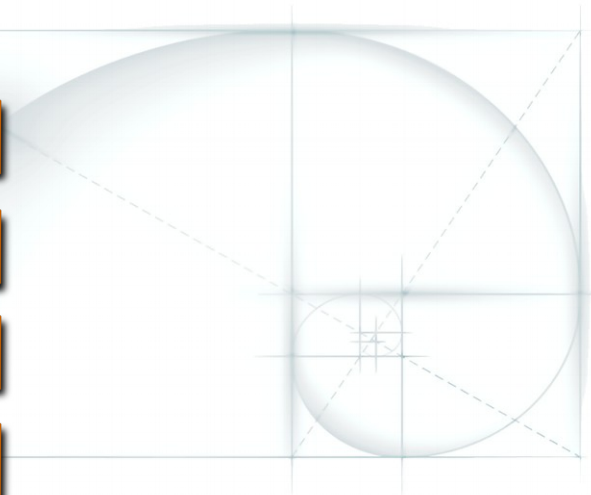
**What
is a
Hash
Table?**

Andrés

Dale

James

Norm



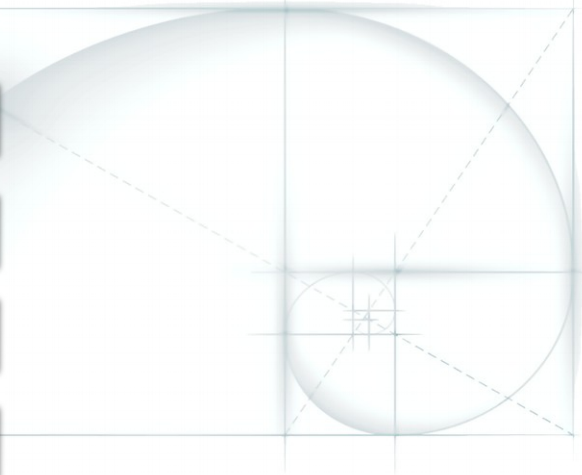
Keys

Andrés

Dale

James

Norm



Keys

**Hash
Function**

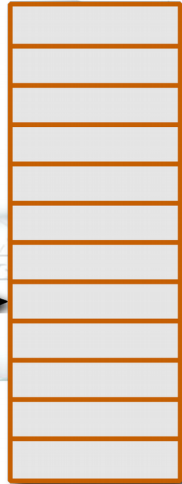
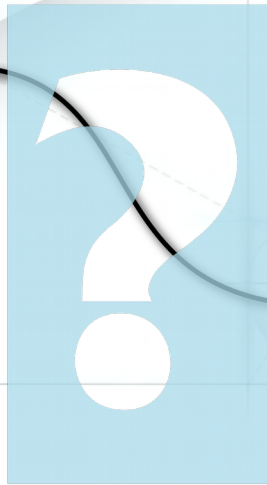
**Hash
Table**

Andrés

Dale

James

Norm



Keys

**Hash
Function**

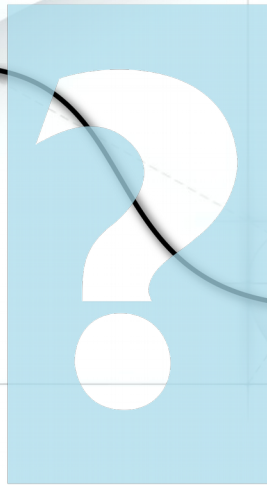
**Hash
Table**

Andrés

Dale

James

Norm



Andrés

Keys

**Hash
Function**

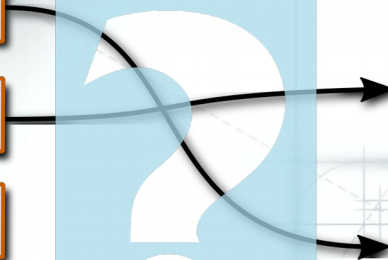
**Hash
Table**

Andrés

Dale

James

Norm



Andrés

Keys

**Hash
Function**

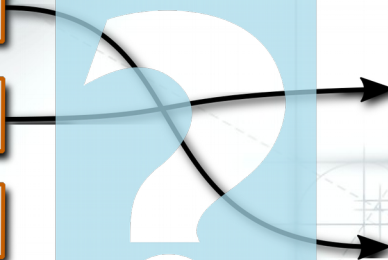
**Hash
Table**

Andrés

Dale

James

Norm



Keys

Hash
Function

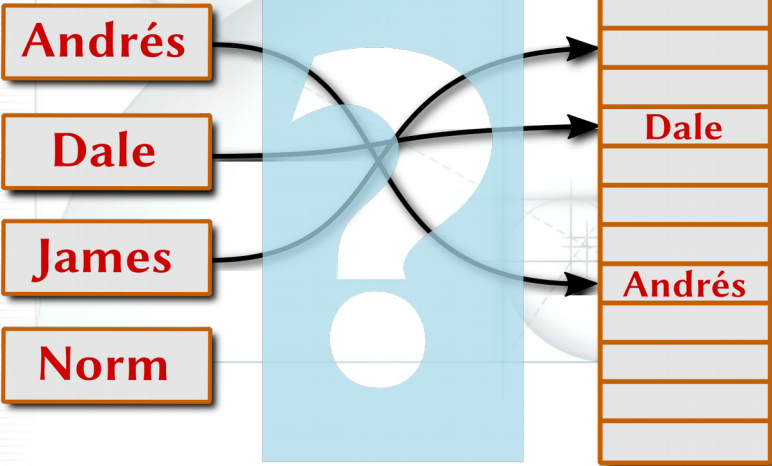
Hash
Table

Andrés

Dale

James

Norm



Keys

Hash
Function

Hash
Table

Andrés

Dale

James

Norm



James
Dale
Andrés



Keys

Hash
Function

Hash
Table

Andrés

Dale

James

Norm



James
Dale
Andrés



Keys

Hash
Function

Hash
Table

Andrés

Dale

James

Norm



James

Dale

Andrés

Norm



Keys

Hash
Function

Hash
Table

Andrés

Dale

James

Norm



James

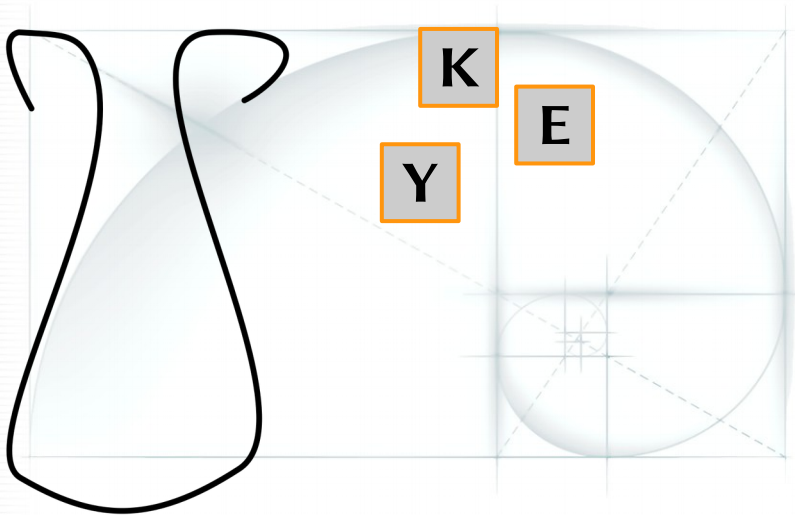
Dale

Andrés

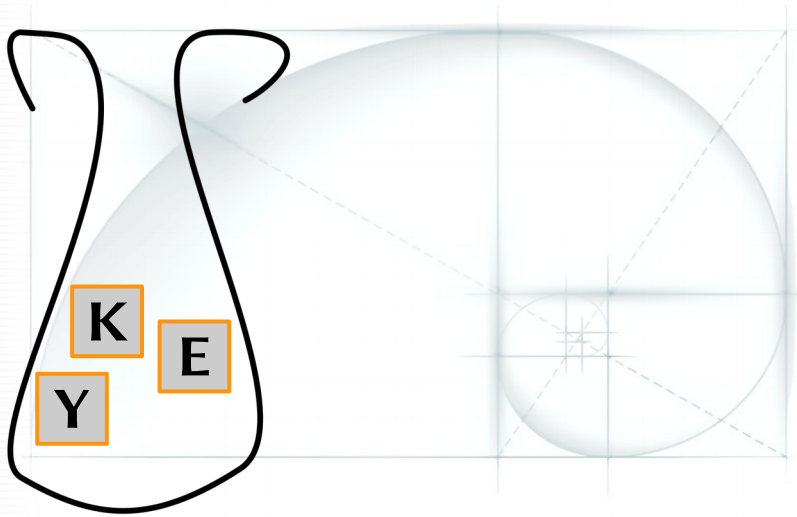
Norm



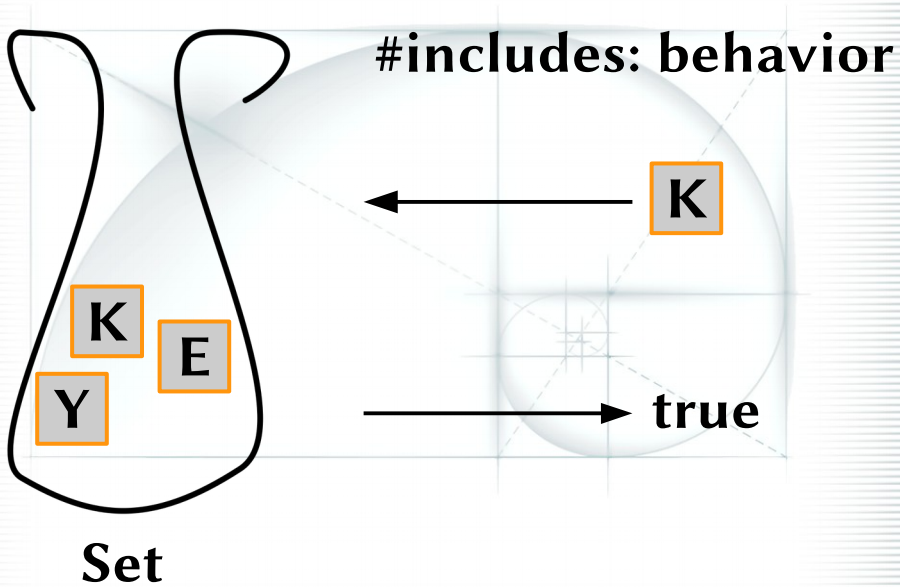
Hash Tables in Smalltalk

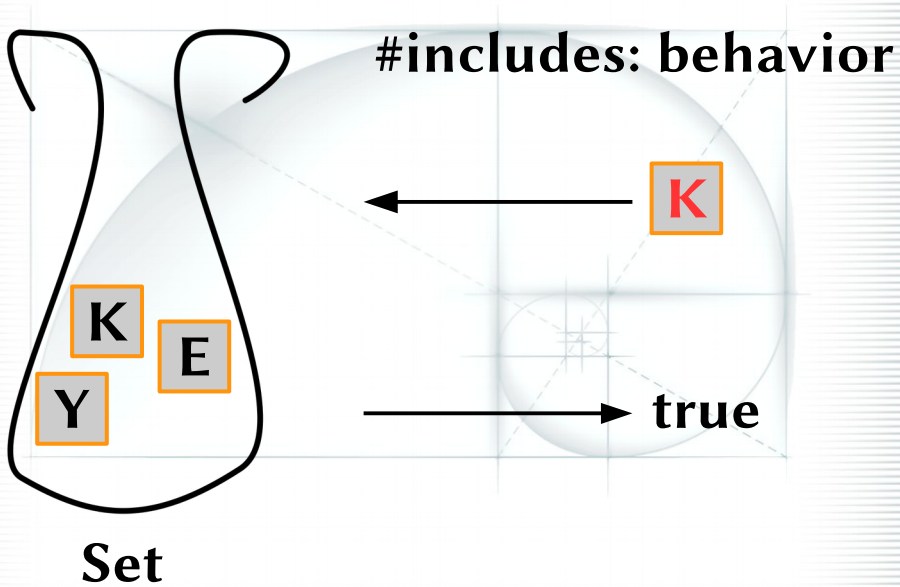


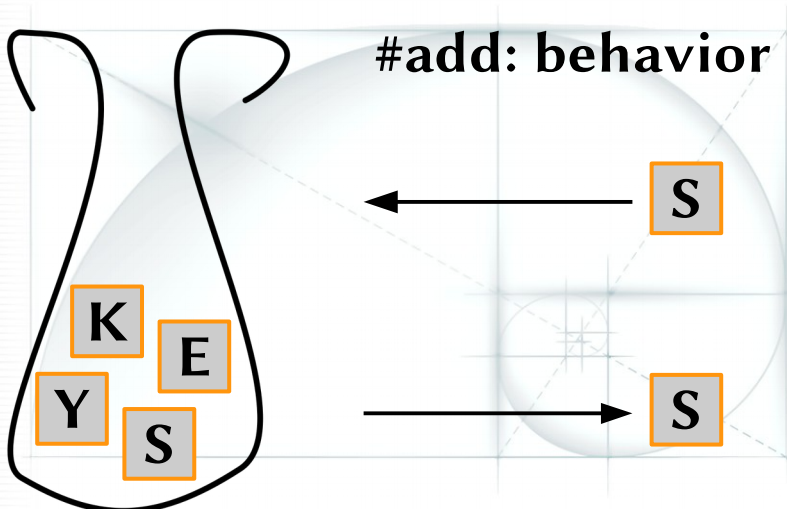
Set



Set







#add: behavior

S

K

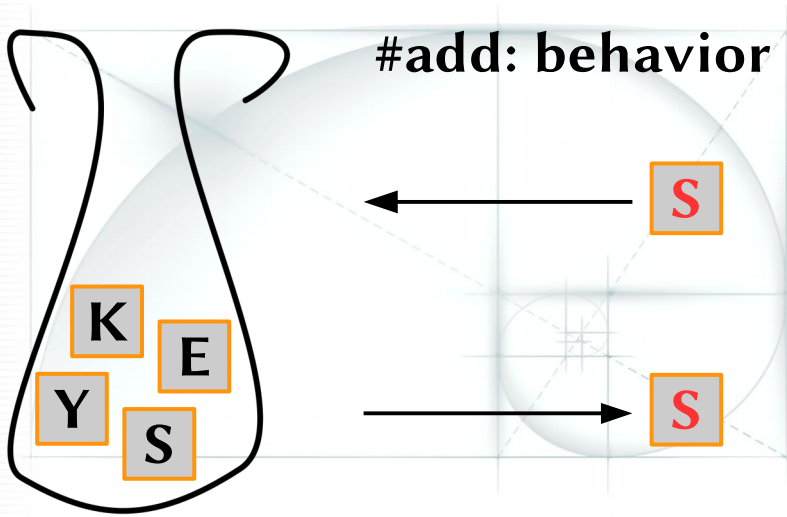
E

Y

S

S

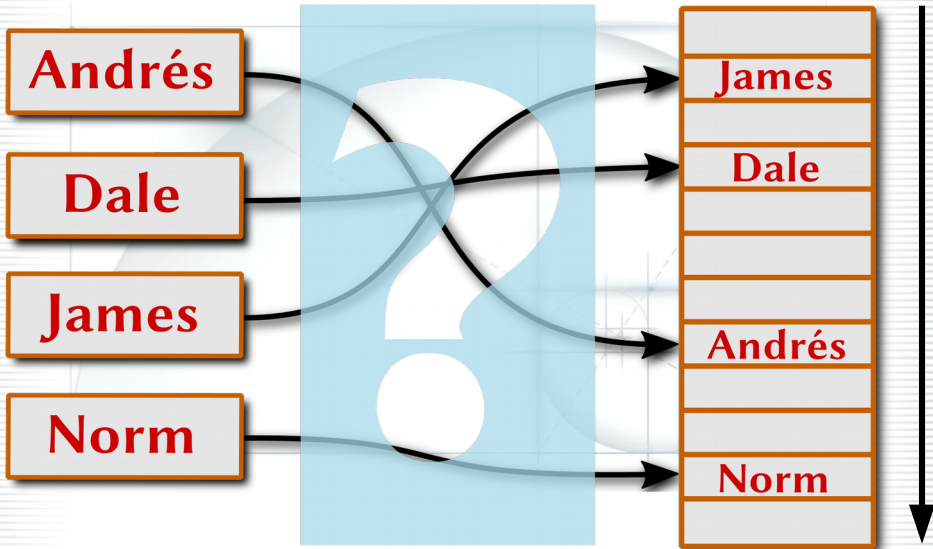
Set

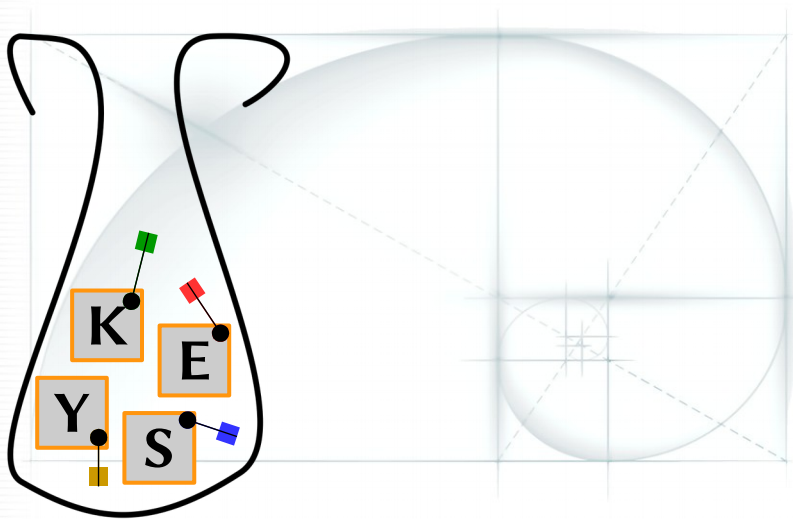


#add: behavior

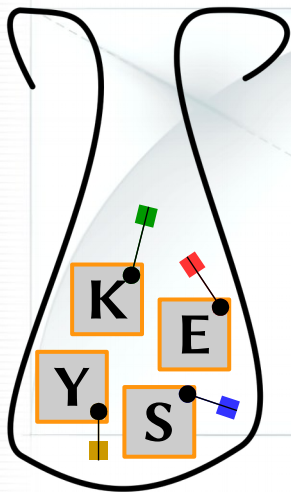
Set

#do: aBlock



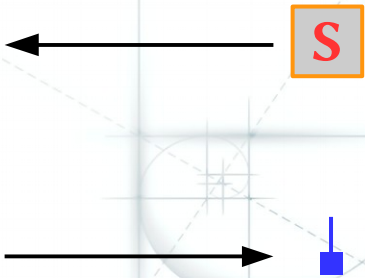


Dictionary

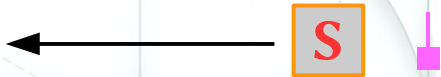
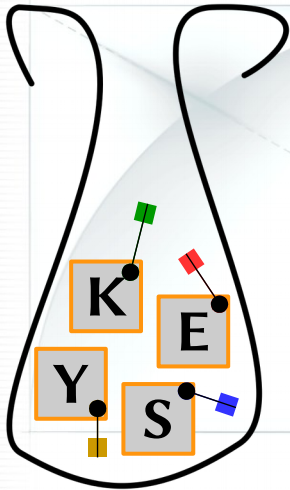


Dictionary

#at: behavior

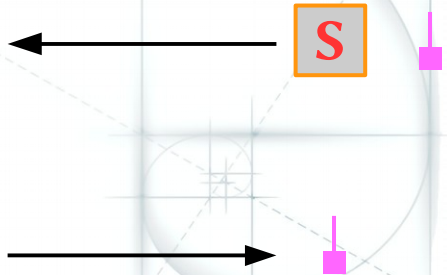
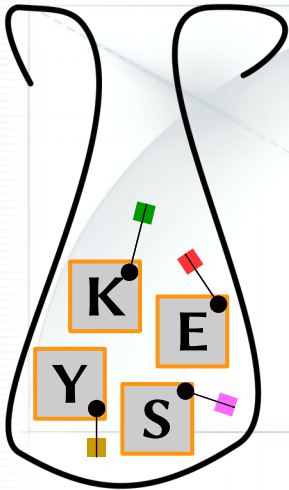


#at:put: behavior



Dictionary

#at:put: behavior



Dictionary



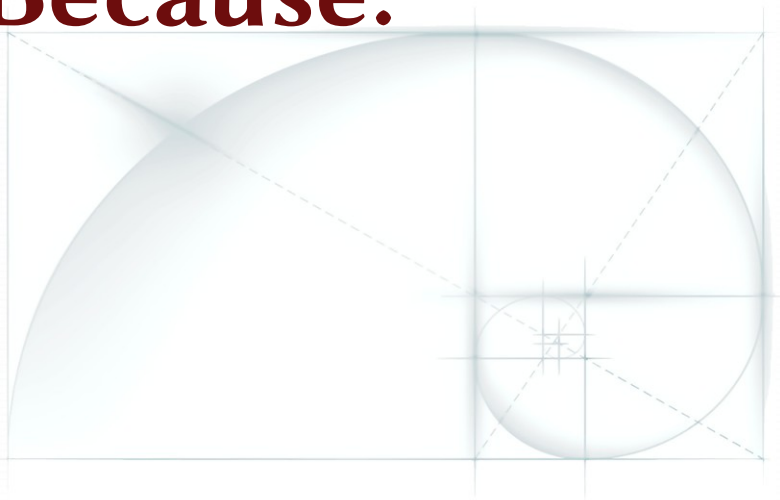
**Build
Your Own
Hashed
Collections**

**Build
Your Own
Hashed
Collections**

Why?



Because:



Because:

- **Special Requirements**



Because:

- **Special Requirements**
- **Security**



Because:

- **Special Requirements**
- **Security**
- **Performance**



Because:


- **Special Requirements**
- **Security**
- **Performance**
 - **Space**



Because:

- **Special Requirements**
- **Security**
- **Performance**
 - **Space**
 - **Speed**





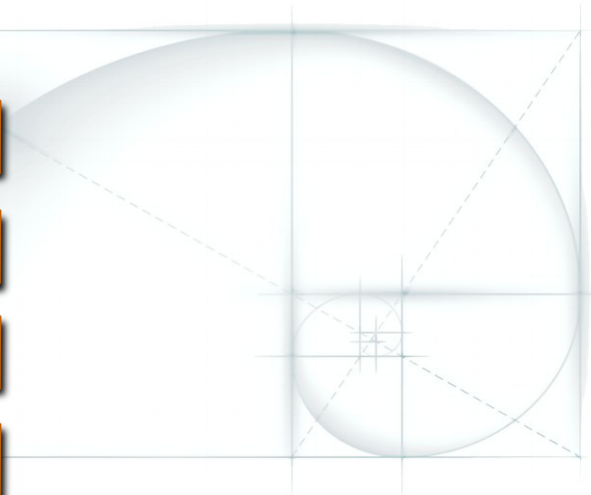
Hash Functions and Domains

Andrés

Dale

James

Norm



Perfect hash function

Andrés

$$f(\text{Andrés}) = 10$$

Dale

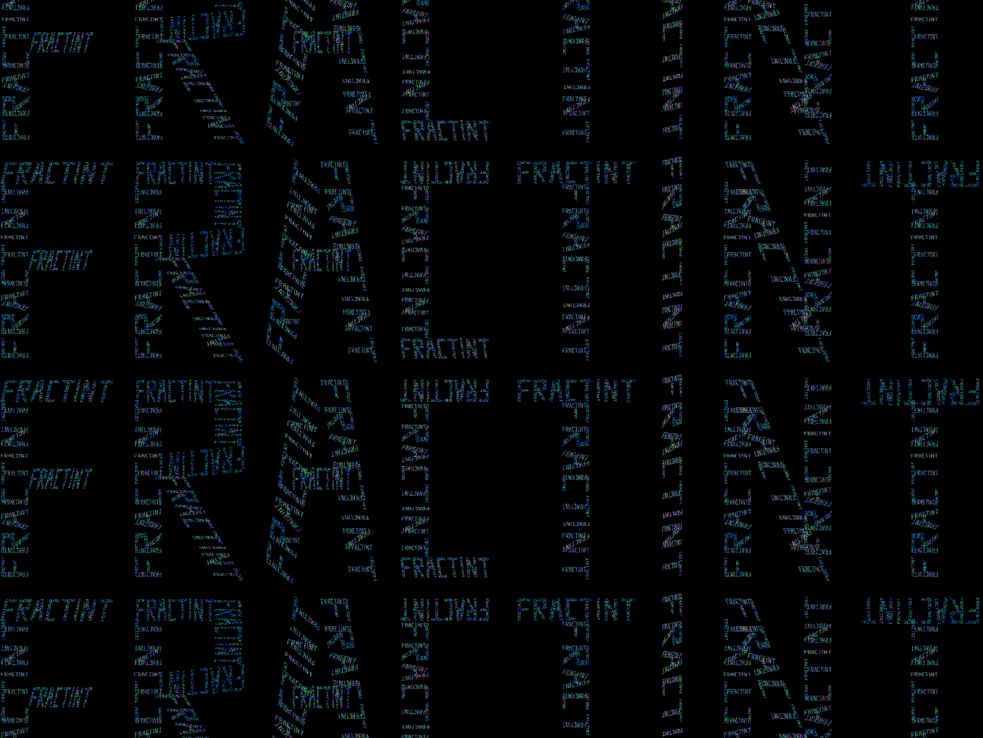
$$f(\text{Dale}) = 4$$

James

$$f(\text{James}) = 7$$

Norm

$$f(\text{Norm}) = 1$$



Pragmatic hash functions

Andrés

Dale

James

Norm

$$f(\text{string}) = k \bmod s$$

#hash

k

k mod s

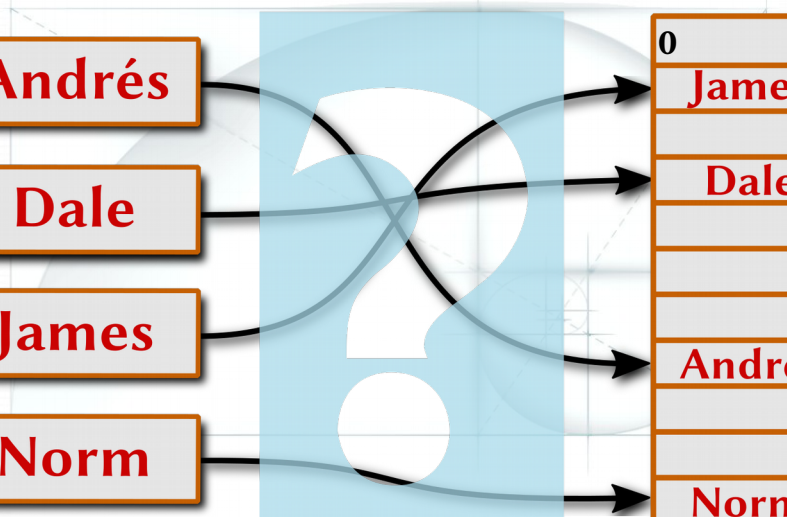
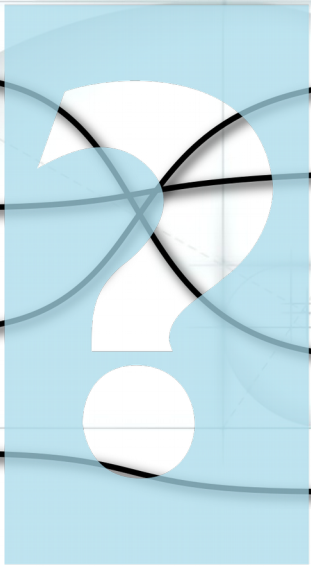
Andrés

Dale

James

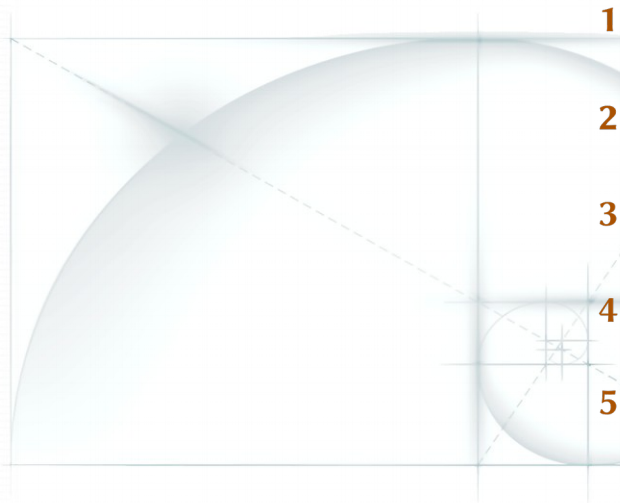
Norm

0
James
Dale
Andrés
Norm
s - 1





Collisions



1

2

3

4

5

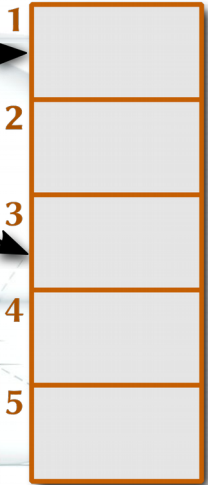
K_1



1	
2	
3	
4	
5	

K_1

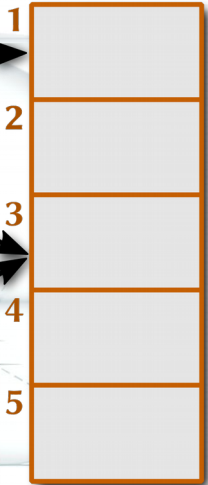
K_2

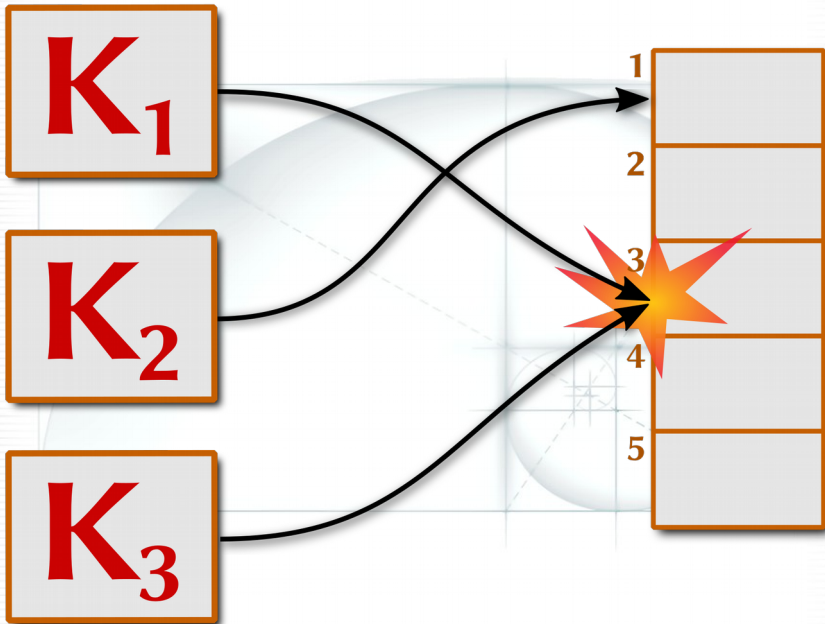


K_1

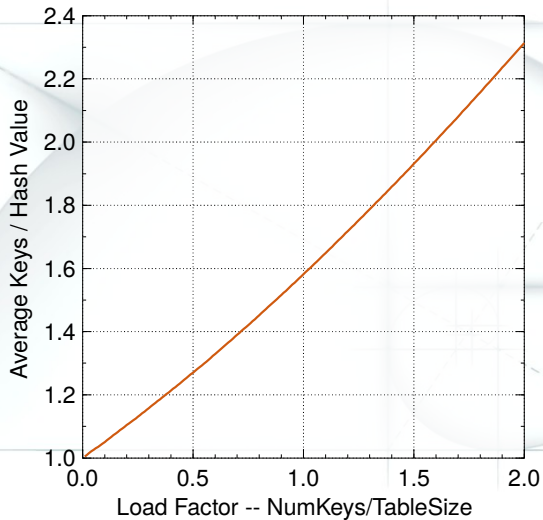
K_2

K_3





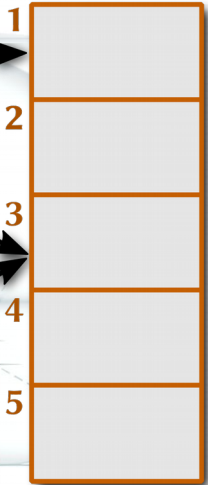
Random Hash Collision Rate



K_1

K_2

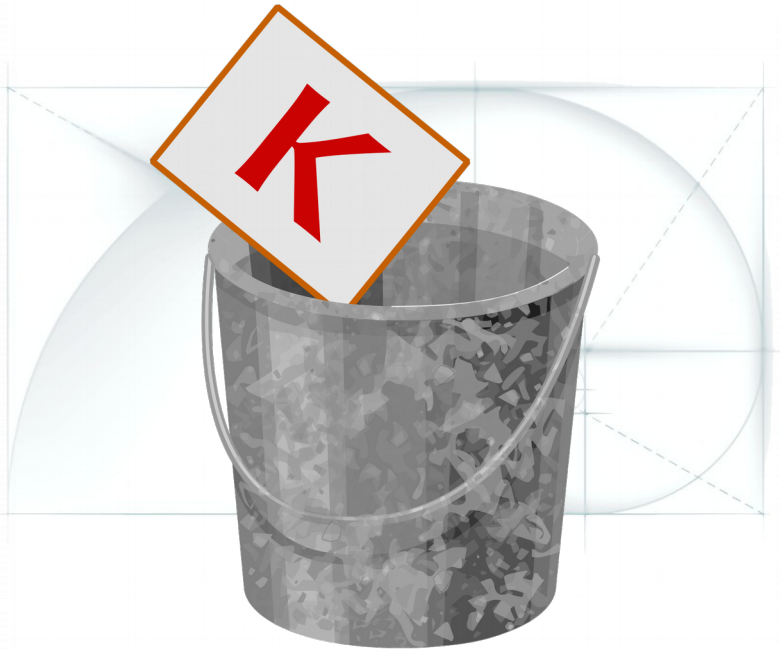
K_3



A technical drawing of a sphere, showing its construction with a grid and dashed lines. The sphere is rendered in a light blue color with a gradient. The word "Buckets" is written in a large, bold, dark red serif font across the center of the sphere. The background is white with a light blue grid and dashed lines indicating the sphere's geometry.

Buckets











Linked Buckets

1



2



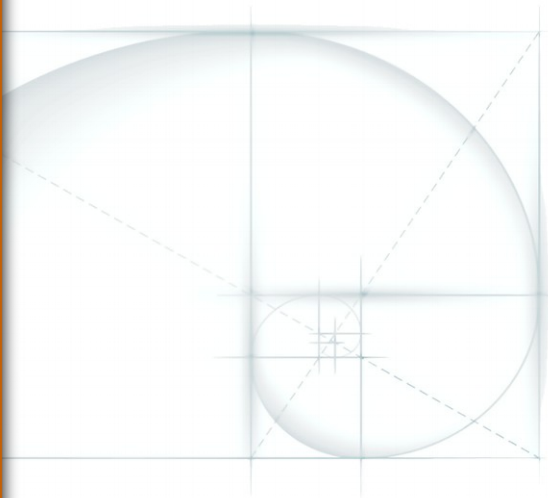
3

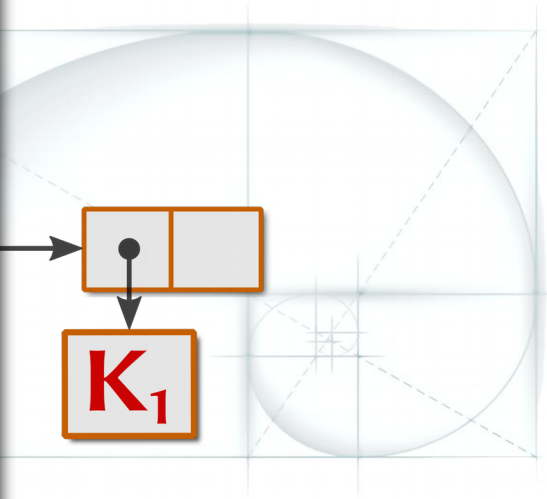
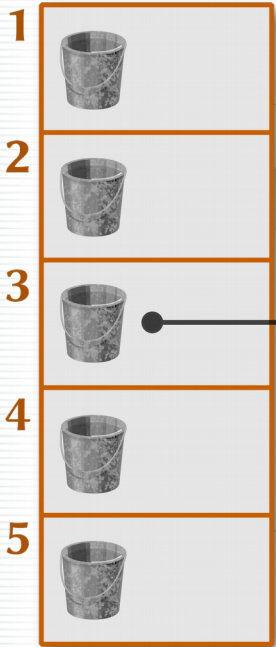


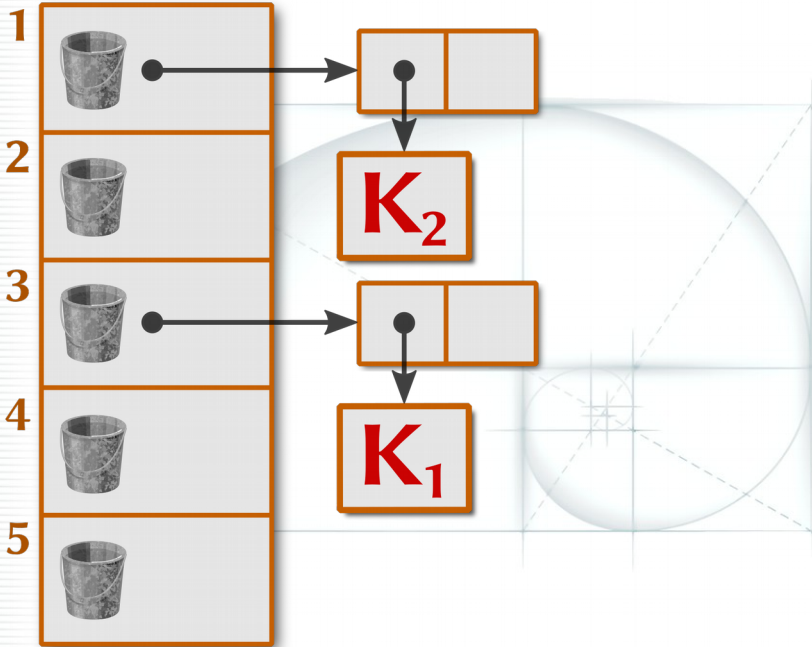
4

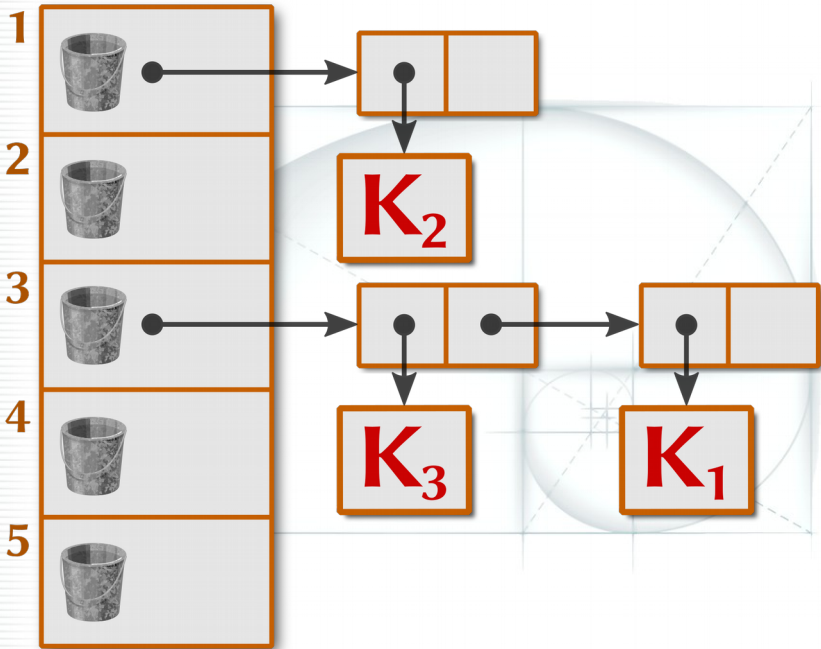


5



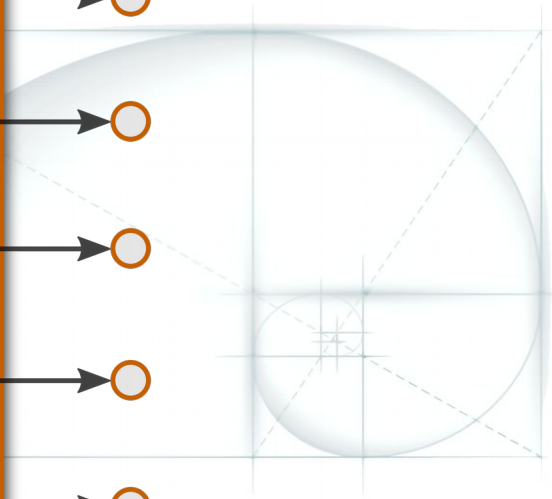
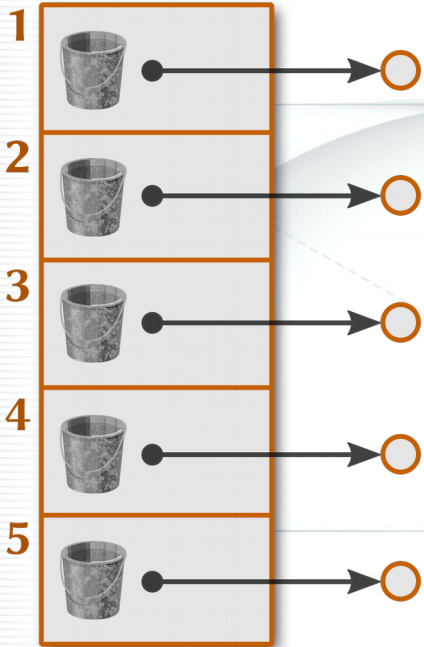


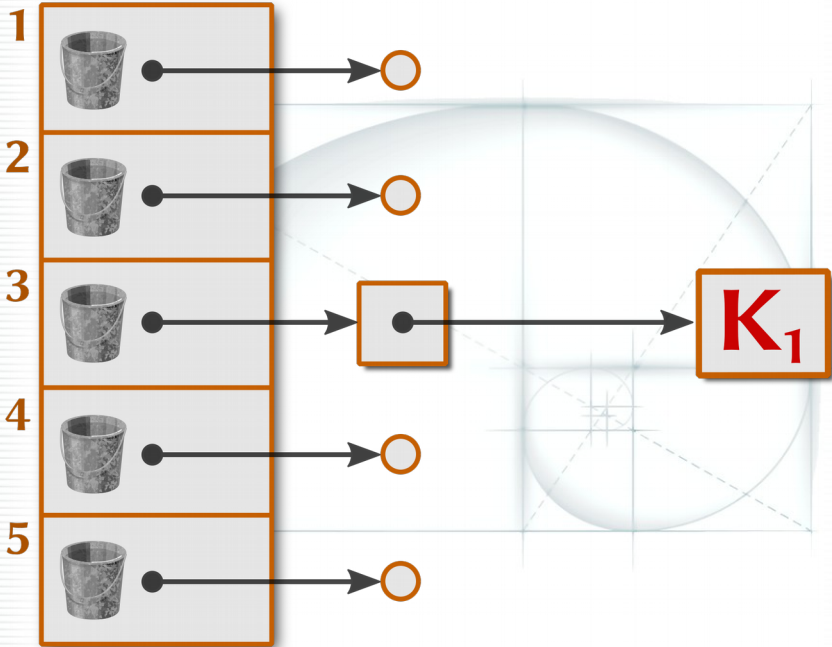


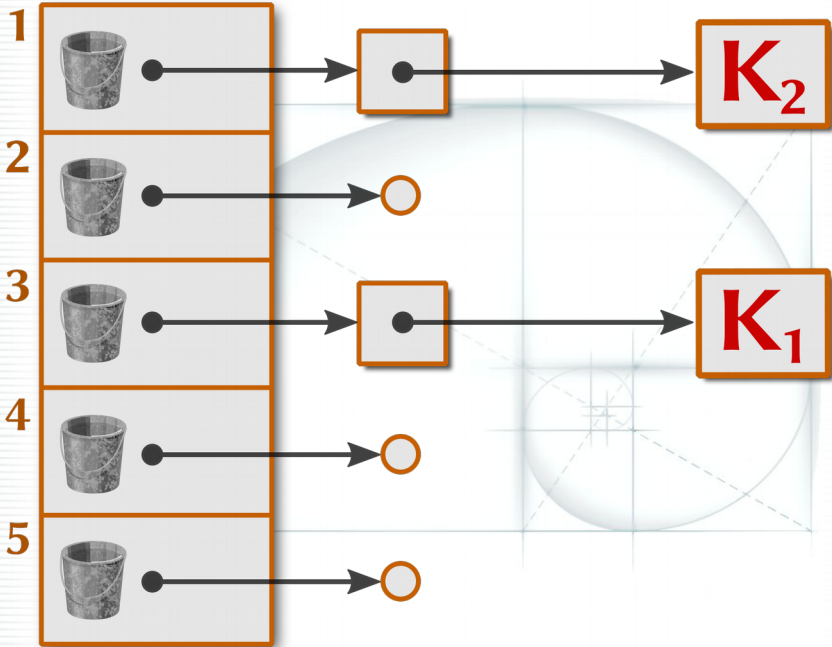


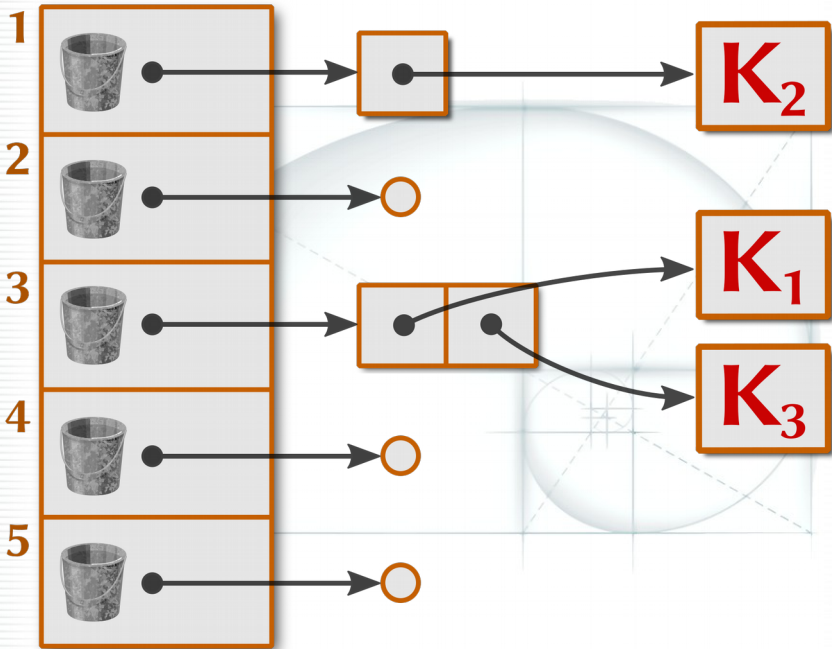



Array Buckets



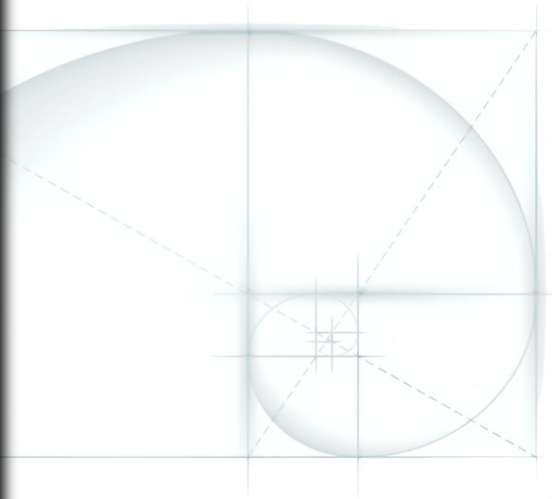


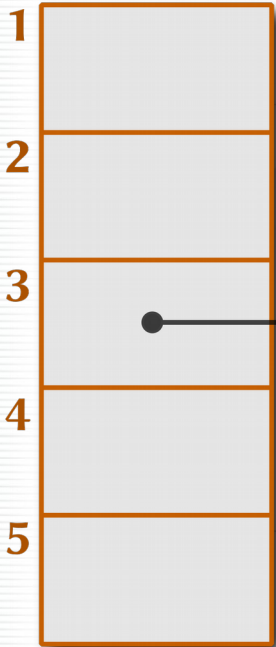




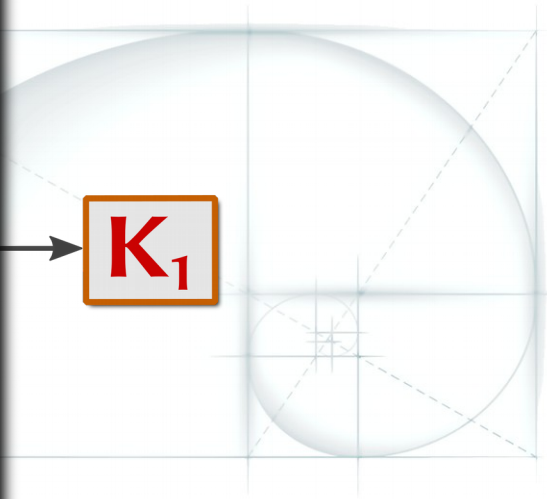


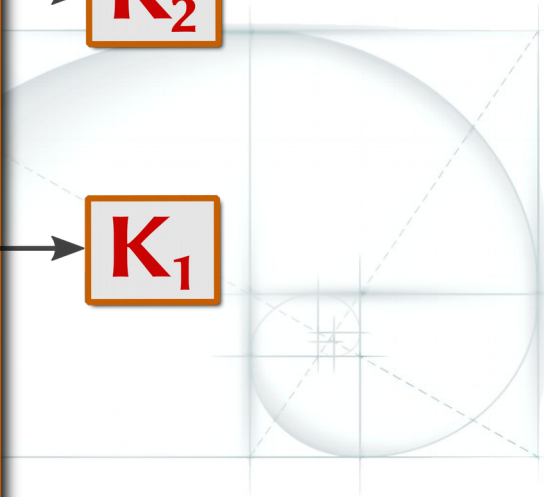
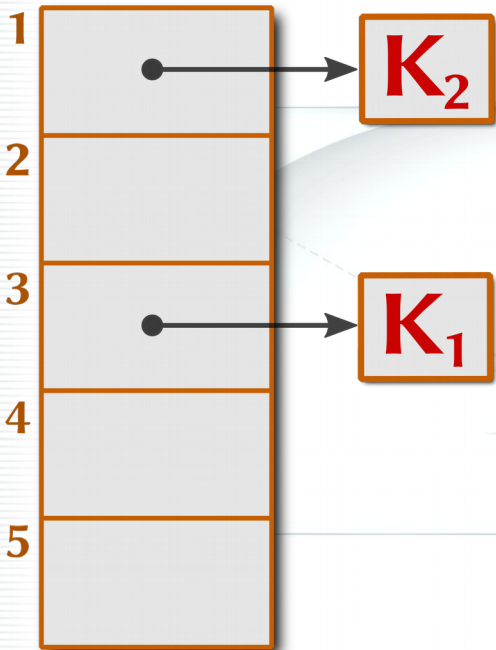
Linear Probing

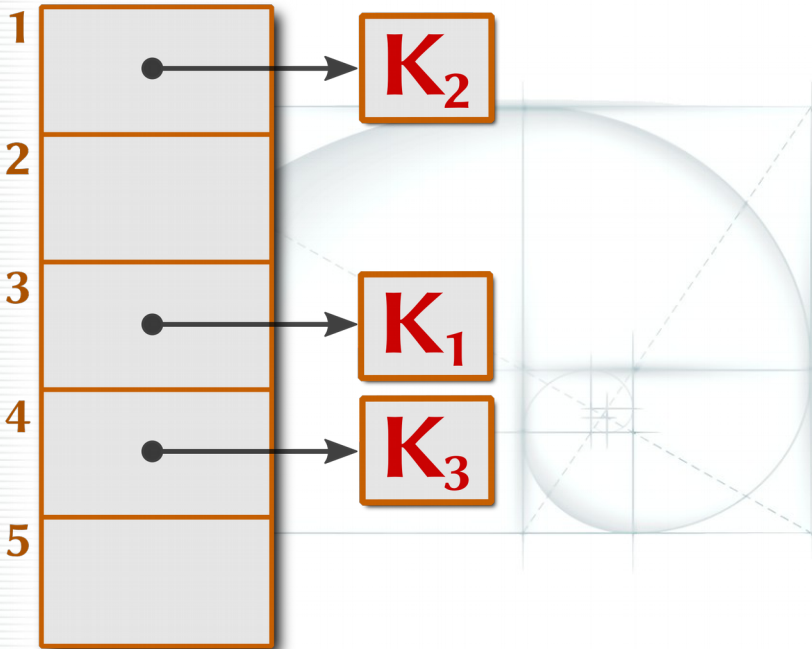




K_1









**“Lazy”
Buckets**

1



2



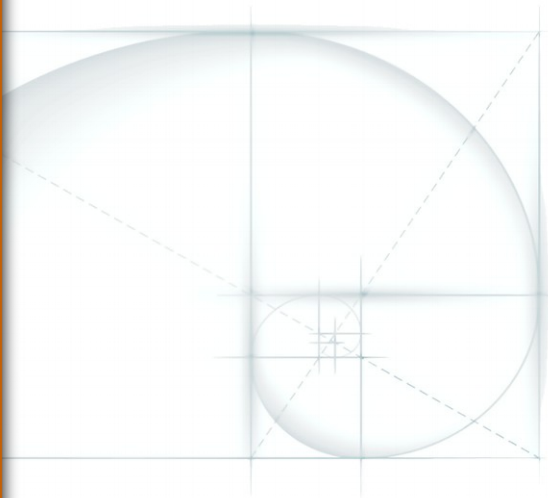
3



4



5



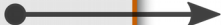
1



2



3

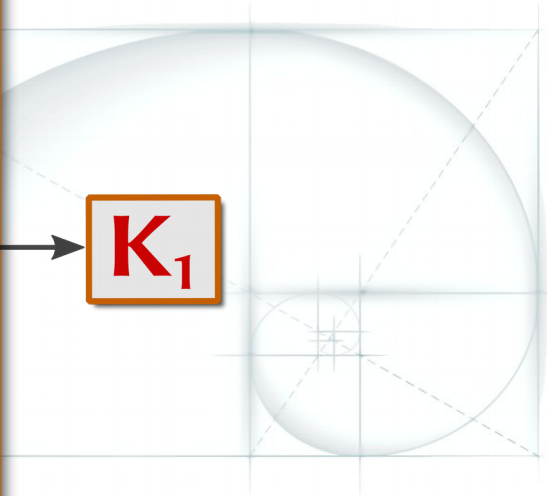


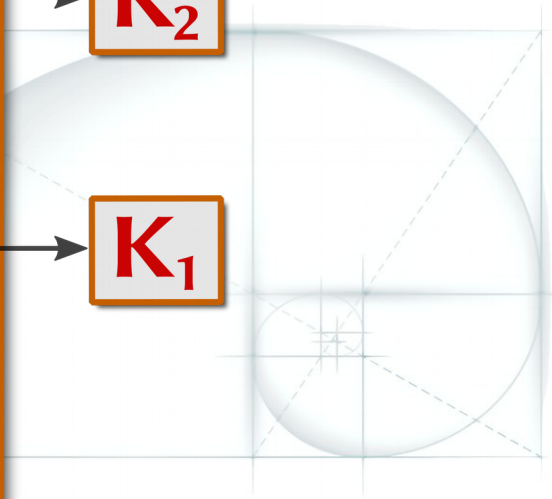
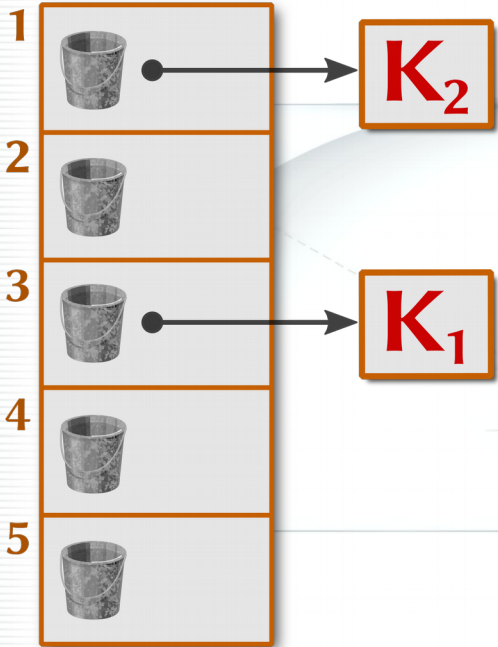
K₁

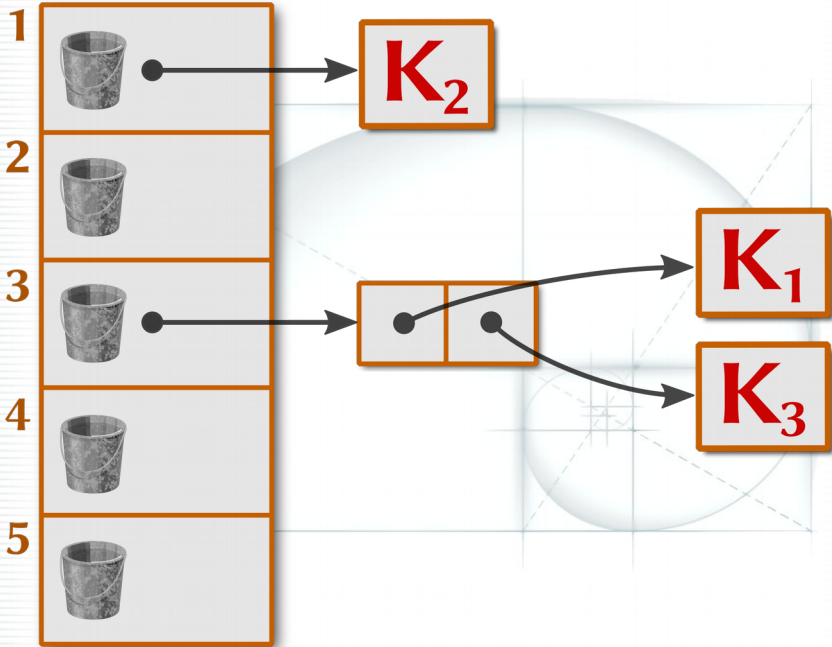
4



5









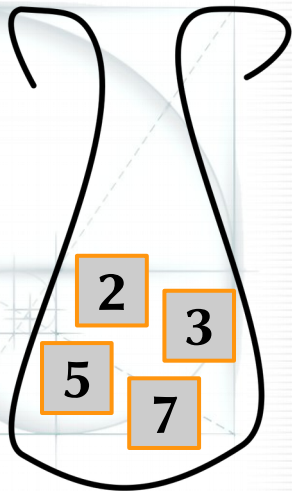
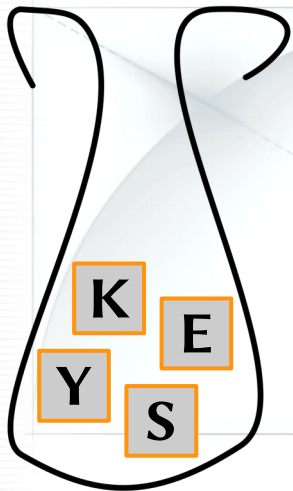
**Some
Advice**



Choosing a Hash Function

D

f(D)

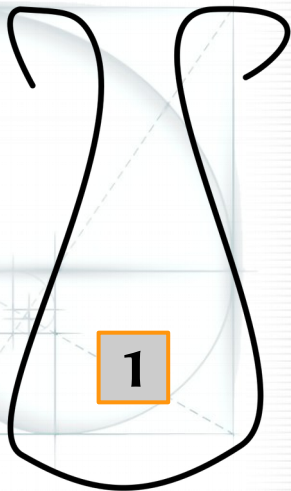
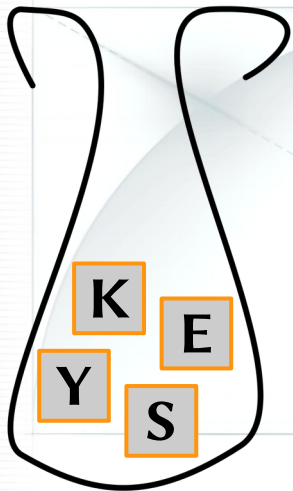


Data Set

#hash Values

D

f(D)

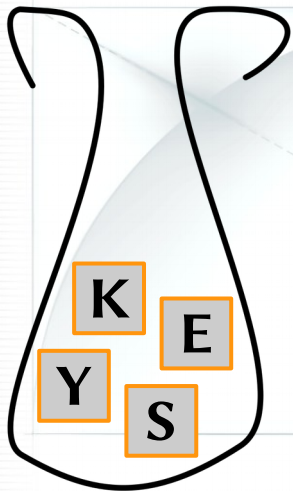


Data Set

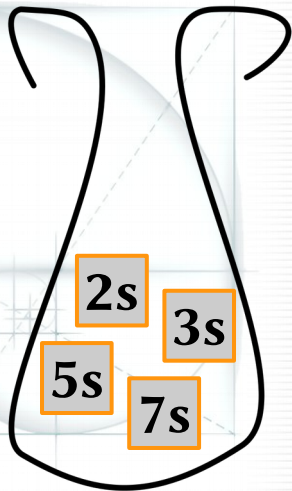
#hash Values

D

f(D)



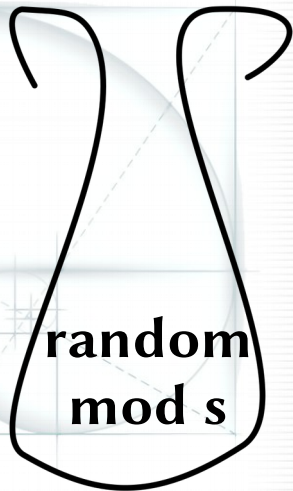
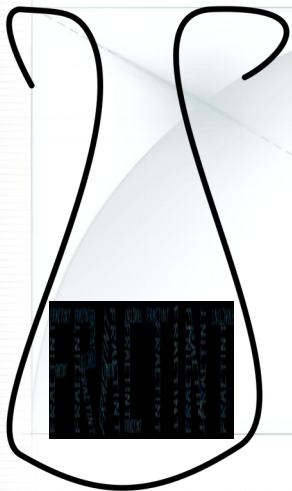
Data Set



#hash Values

D

f(D)

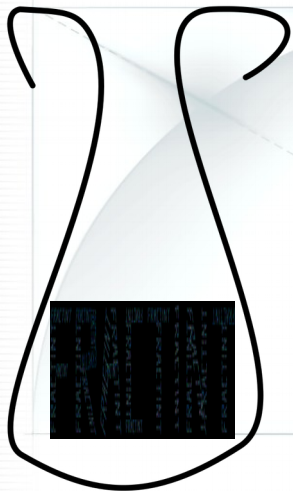


Data Set

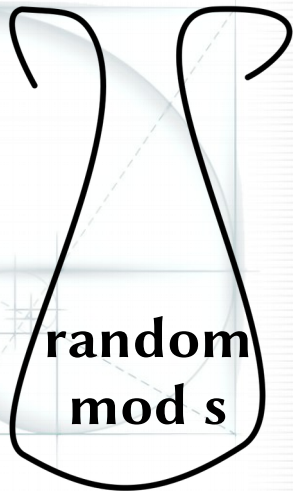
#hash Values

D

f(D)



**crypto
hash**

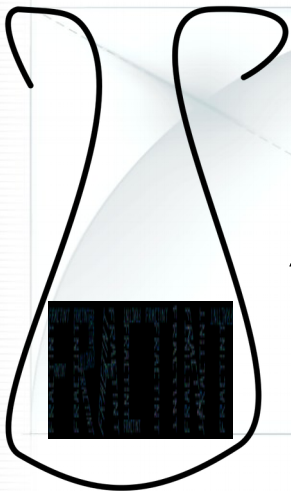


Data Set

#hash Values

D

f(D)



**“hash
function”**

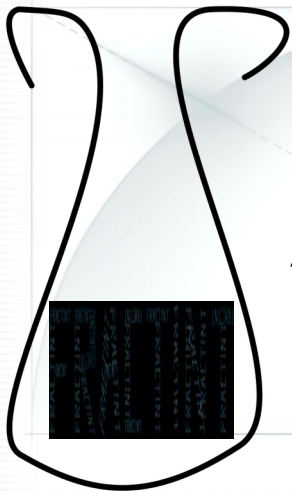


Data Set

#hash Values

D

f(D)



**“hash
function”**

prime s



**maybe
random
mod s**

Data Set

#hash Values

Collisions

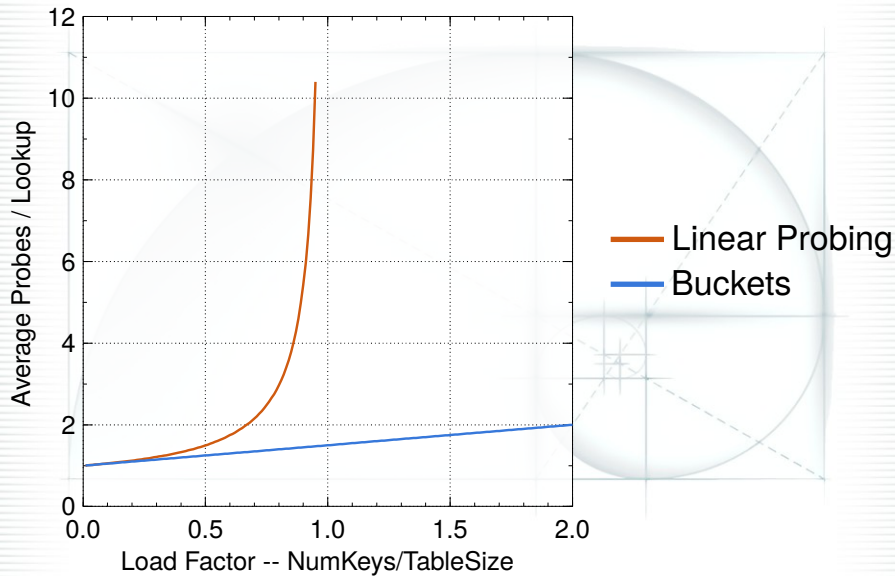
The background of the slide features a light blue, semi-transparent sphere. Overlaid on the sphere is a golden spiral, which is a mathematical curve that starts from a central point and spirals outwards, getting closer to the curve as it moves. The spiral is composed of several overlapping squares of decreasing size. A grid of dashed lines is also overlaid on the sphere, consisting of several vertical and horizontal lines that intersect to form a grid pattern. The word "Collisions" is written in a large, bold, dark red serif font across the center of the sphere.



Collisions

...again

Number of probes





€/probe

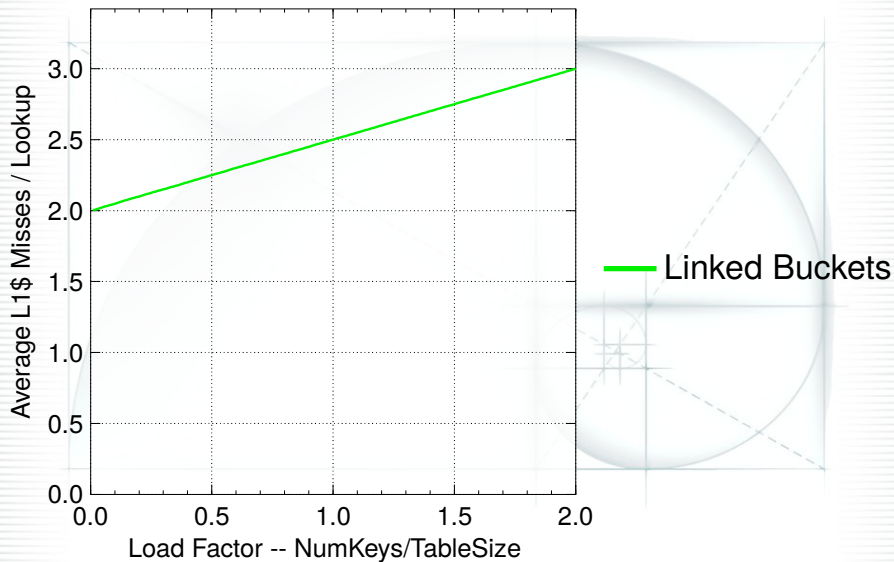
≠

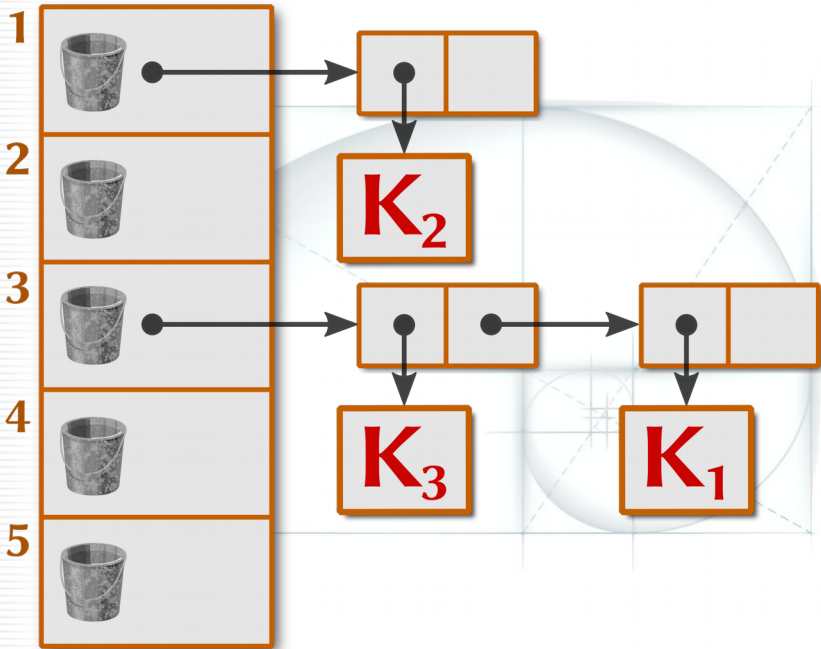
Hypothesis:

$$\epsilon \propto$$

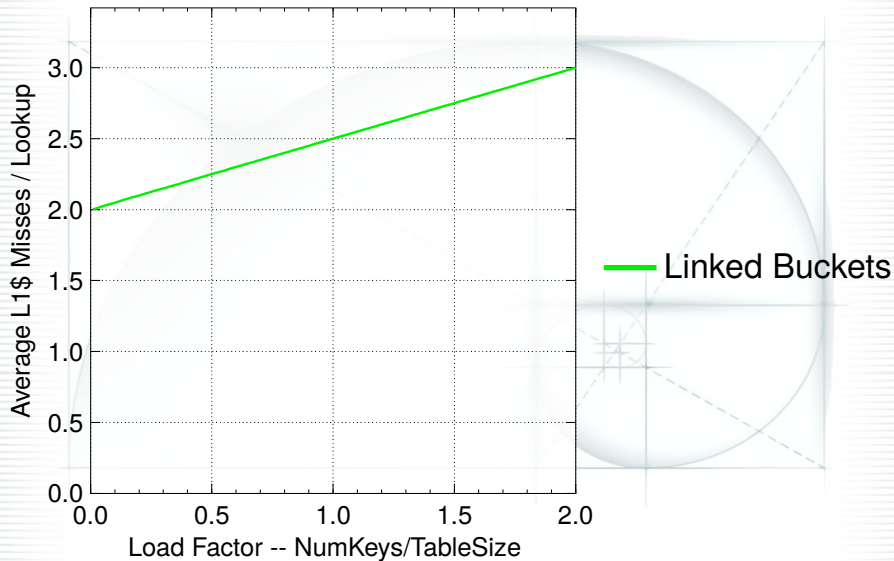
L1\$ misses

Predicted Cache Misses

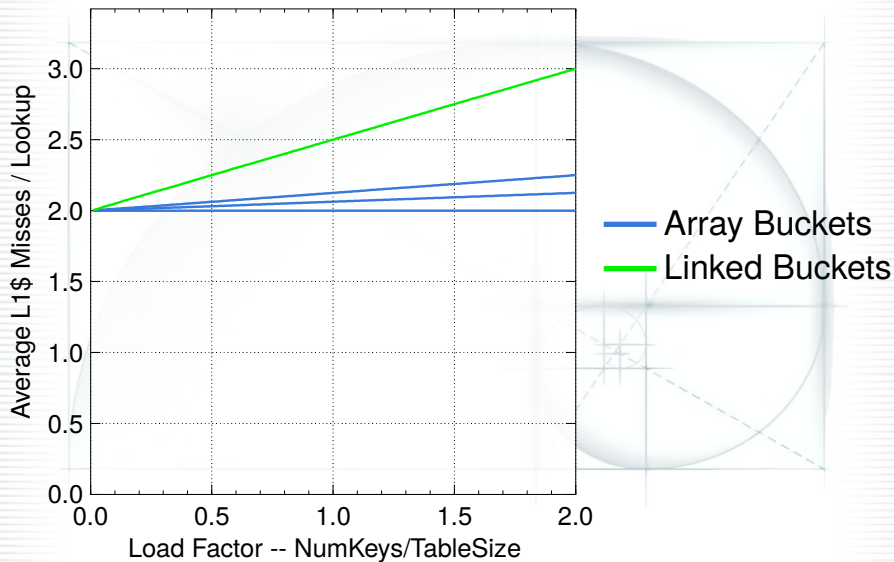


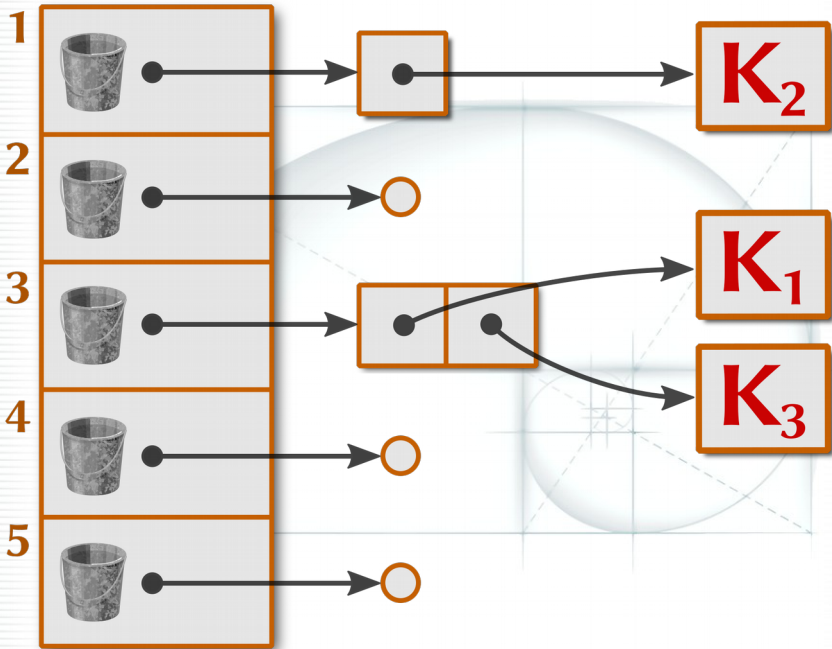


Predicted Cache Misses

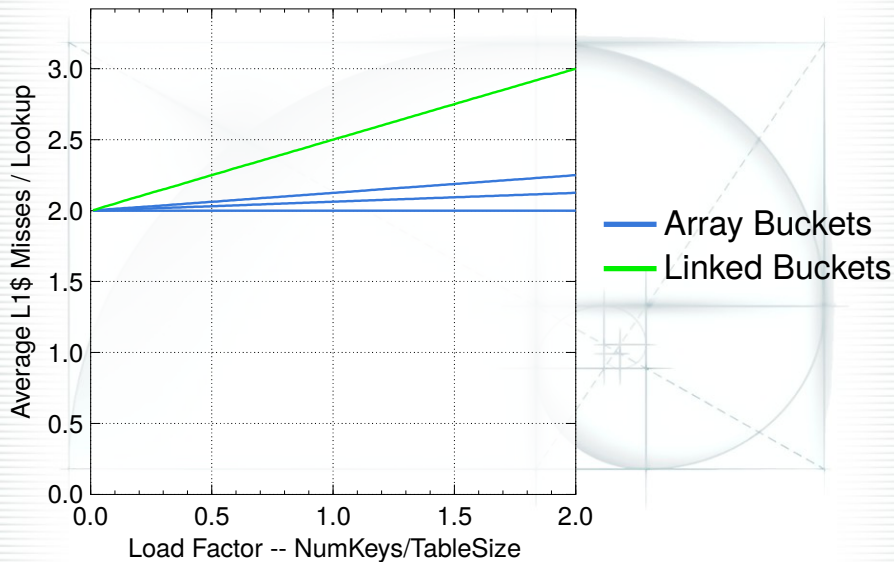


Predicted Cache Misses

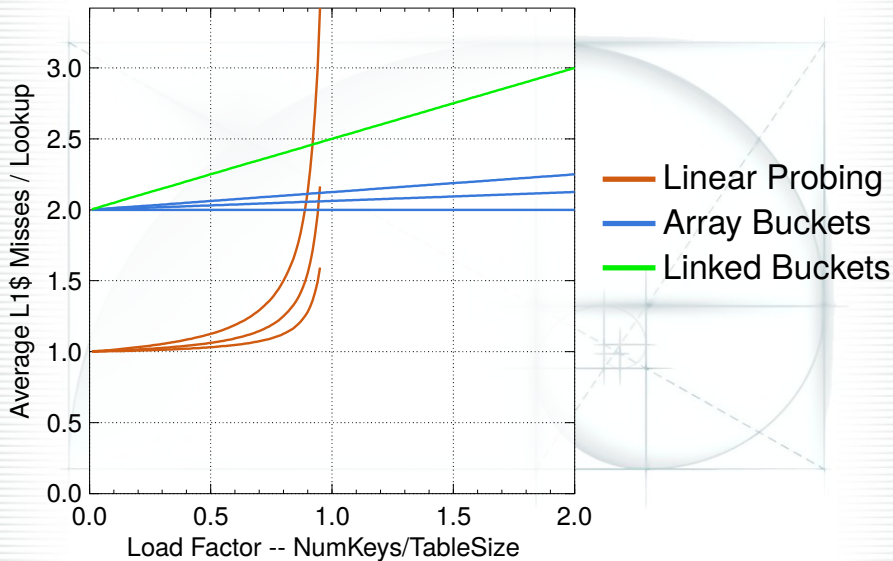


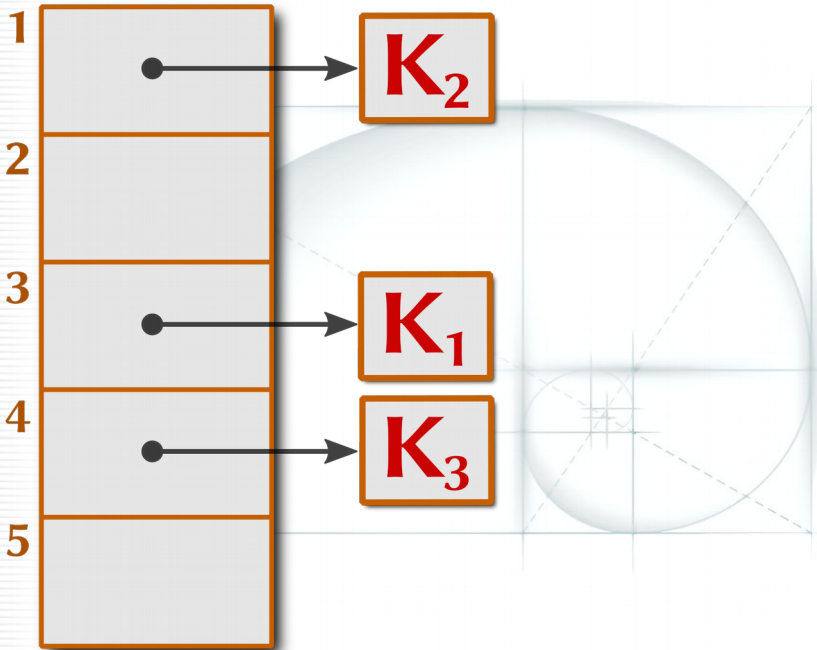


Predicted Cache Misses

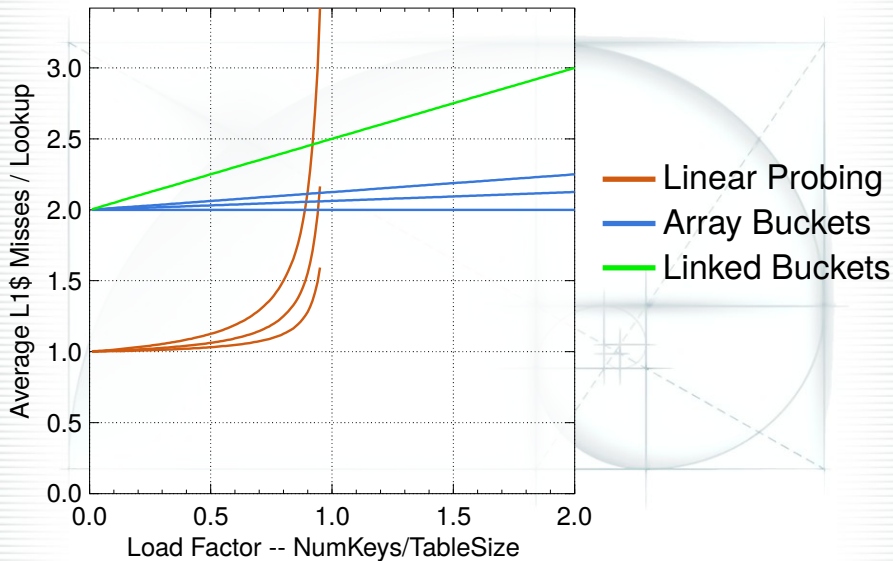


Predicted Cache Misses

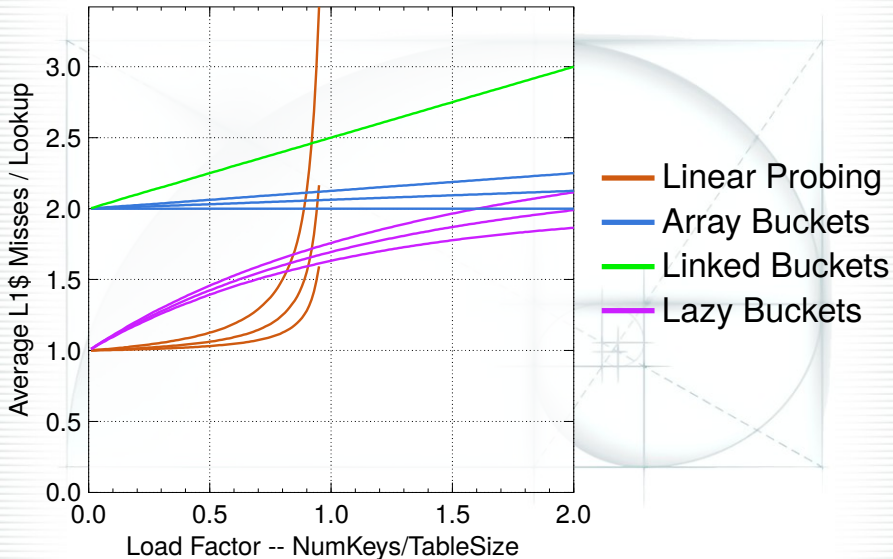


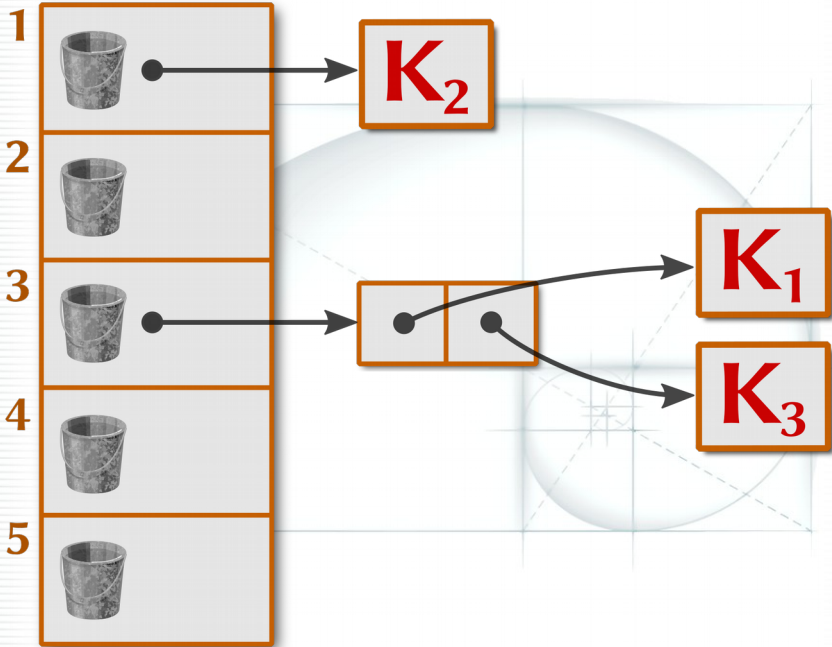


Predicted Cache Misses

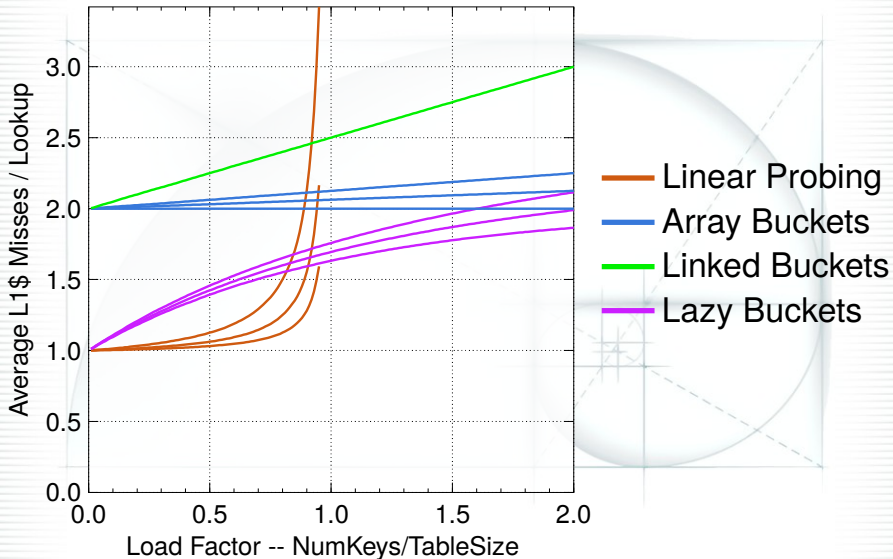


Predicted Cache Misses

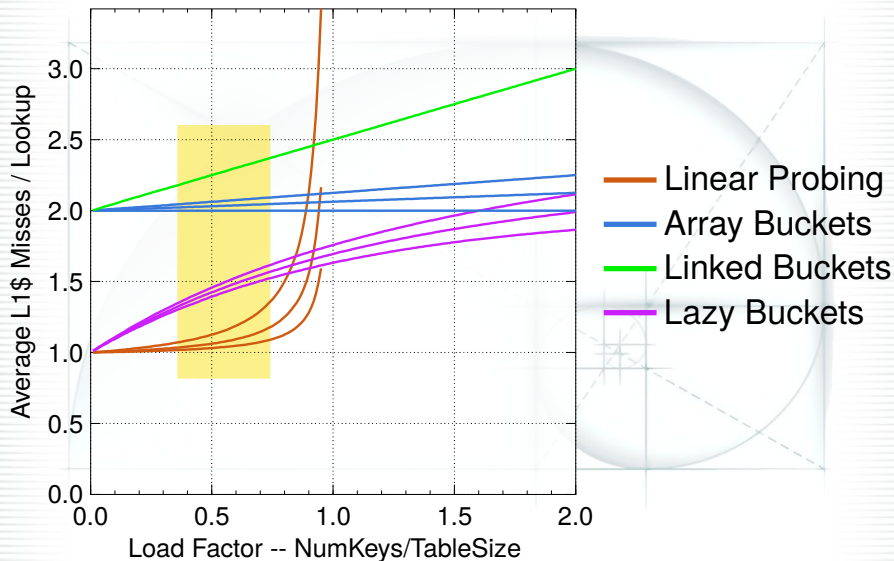




Predicted Cache Misses



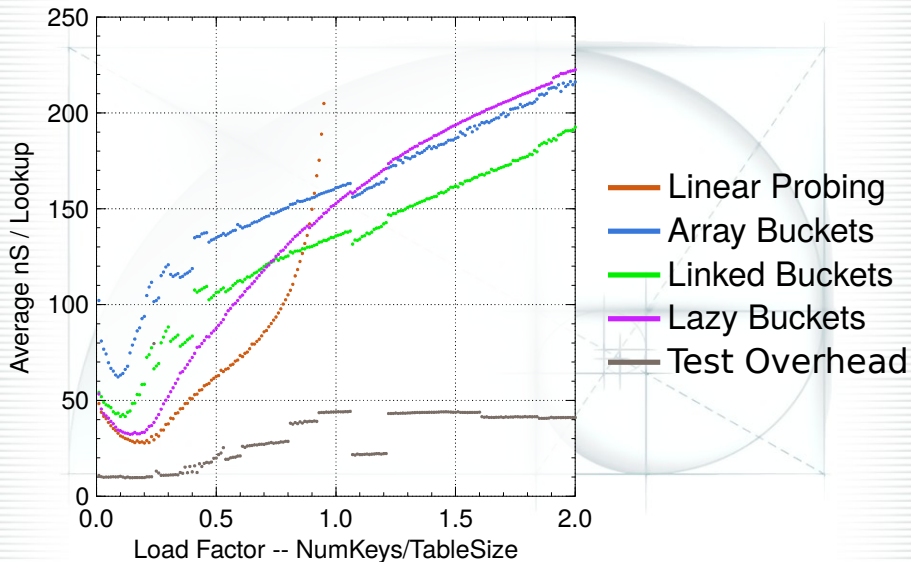
Predicted Cache Misses



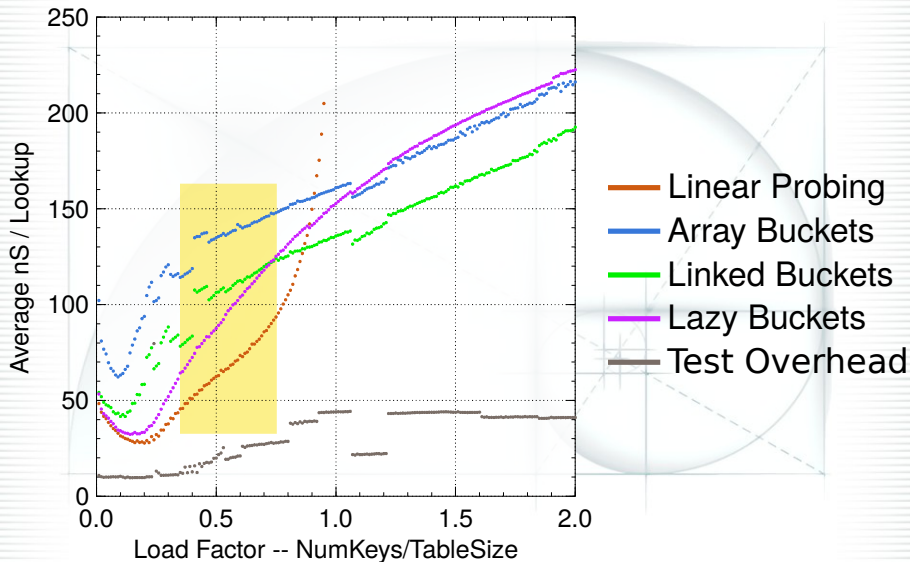


Test the Hypothesis

Average Lookup Time



Average Lookup Time

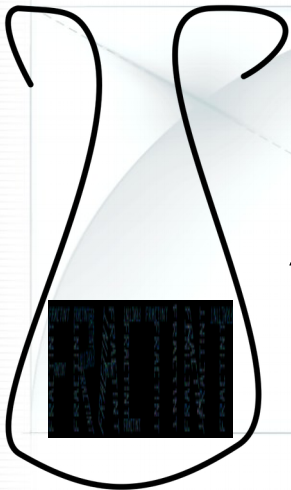




Security

D

f(D)



**“hash
function”**



Data Set

#hash Values

D

f(D)



**“hash
function”**



Data Set

#hash Values

A technical drawing of a dome or semi-circle. The dome is shaded in light blue and is centered within a rectangular frame. A golden spiral is drawn on the right side of the dome, starting from a small square and expanding outwards. The drawing includes a grid of solid lines and dashed lines, suggesting a geometric construction or architectural plan. The text "Why Test?" is overlaid in a bold, dark red font across the center of the dome.

Why Test?

JDK 1.5

```
/**  
 * Returns a hash value for the specified object. In addition to  
 * the object's own hashCode, this method applies a "supplemental  
 * hash function," which defends against poor quality hash functions.  
 * This is critical because HashMap uses power-of-two length  
 * hash tables.<p>  
 *  
 * The shift distances in this function were chosen as the result  
 * of an automated search over the entire four-dimensional search space.  
 */  
static int hash(Object x) {  
    int h = x.hashCode();  
  
    h += ~(h << 9);  
    h ^= (h >>> 14);  
    h += (h << 4);  
    h ^= (h >>> 10);  
    return h;  
}
```

Multiplicative hash function

String>>hash

```
| answer |  
answer := self size.  
self do:  
  [:each |  
    answer := answer * magic bitAnd: 2^k.  
    answer := answer + each asInteger  
  ].  
^answer
```

Multiplicative hash function

Magic	Hash function
31, 131, 1313, 13131, ...	K&R's digital method, Java
5	C++ STL
1664525	#hashMultiply
19	OCaml
2, 8	Bjarne Stroustrup
13, 53	C# 2005 Professional
32.25	Boost, Justin Sobel
33	Bernstein, Standard ML

Funny hash methods

SomeClass>>hash

^12345

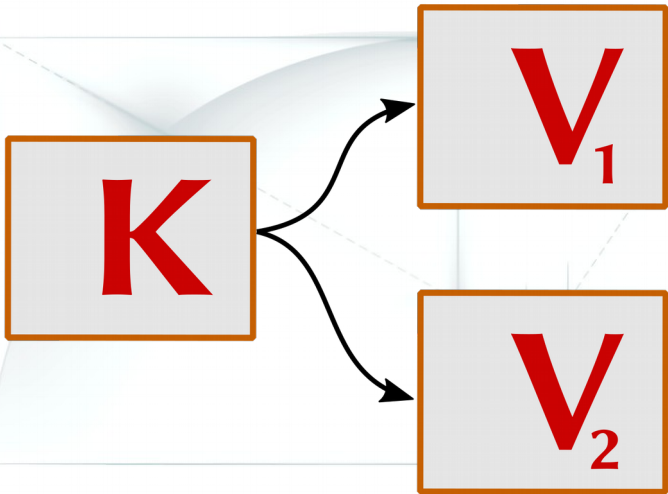
```
bitXor: self a hash;  
bitXor: self b hash;  
bitXor: self c hash;  
yourself
```

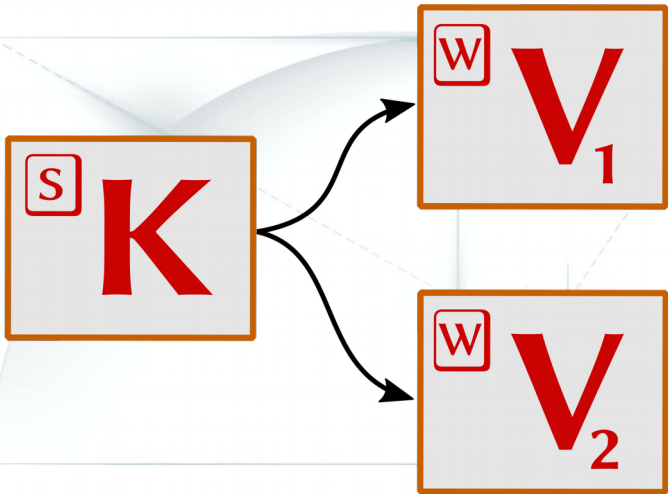


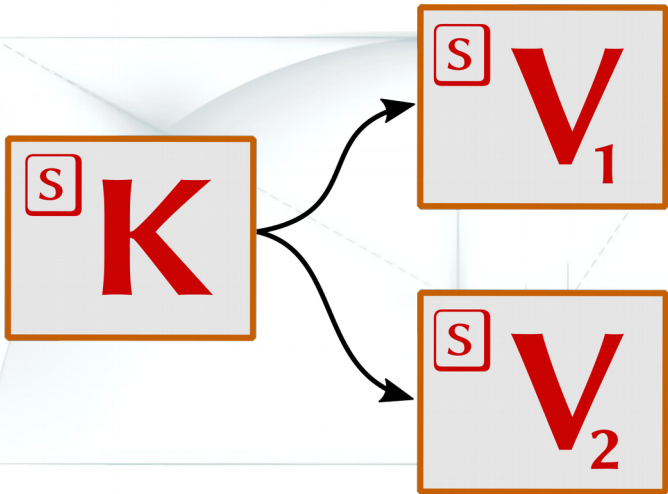

Real-World Examples



GBS ServerMap

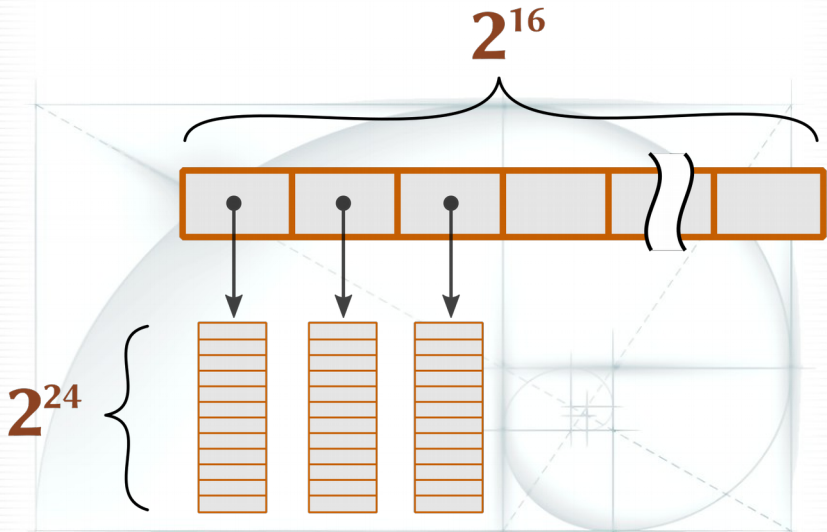






Also

- Integer keys, $1-2^{40}$
(typically smallish)
- ~2M keys
- ~100K updates/sec
- ~ $(1-10) \times 100K$ lookups/sec



2^{24}

V_2

K

V_1

V_2

K

2^{24}

W

V_2

K

V_1

V_2

K

2^{24}

S

V_2

K

V_1

V_2

K



#hashMultiply



SymmetricPoint

Symmetric points

SymmetricPoint>>= aPoint

```
(self x = aPoint x
  and: [self y = aPoint y]) ifTrue: [^true].
(self y = aPoint x
  and: [self x = aPoint y]) ifTrue: [^true].
^false
```

Symmetric points

#hash

1x

```
SymmetricPoint>>hash
```

```
^self x + self y
```

```
SymmetricPoint>>hash
```

```
| xSquared ySquared product |
```

```
xSquared := self x * self x.
```

```
ySquared := self y * self y.
```

```
product := (xSquared + 1) * (ySquared + 1).
```

```
^self x * self y + product
```

5x

Symmetric points

#hash

Use case

1x

```
SymmetricPoint>>hash
```

```
^self x + self y
```

20x

```
SymmetricPoint>>hash
```

```
| xSquared ySquared product |  
xSquared := self x * self x.  
ySquared := self y * self y.  
product := (xSquared + 1) * (ySquared + 1).  
^self x * self y + product
```

5x

1x

Number crunching

1.5gb

2^{30} paths join detection



Number crunching

1.5gb

2^{30} paths join detection

6.4mb

2^{33} paths join detection

Cheney scavenging



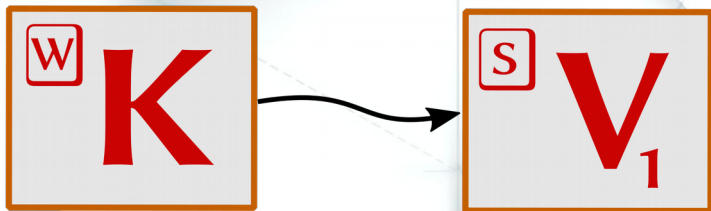
GBS
ClientMap

ClientMap

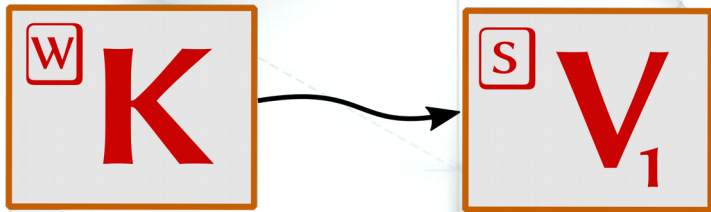
- Identity Dictionary
- ~2M keys
- ~100K updates/sec
- ~ $(1-10) \times 100K$ lookups/sec



32-bit VW

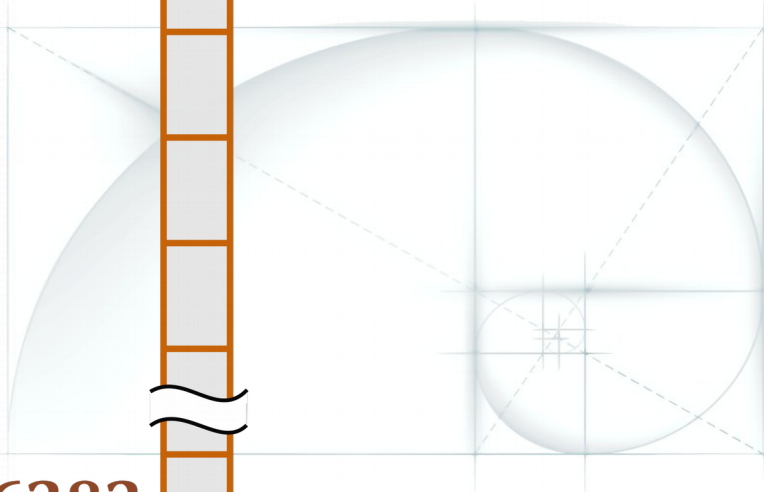


32-bit VW



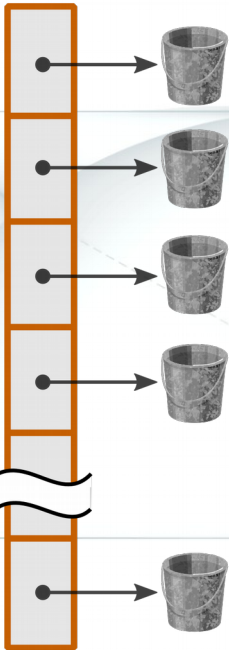
2^{14} hash values

1

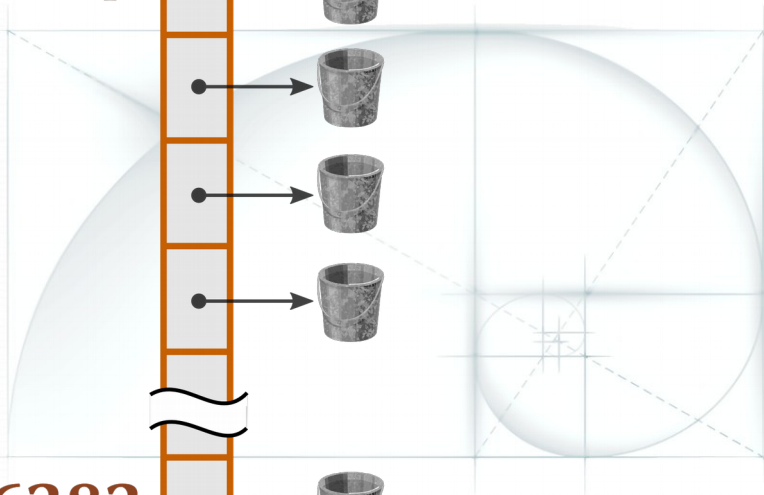


16383

1



16383





W

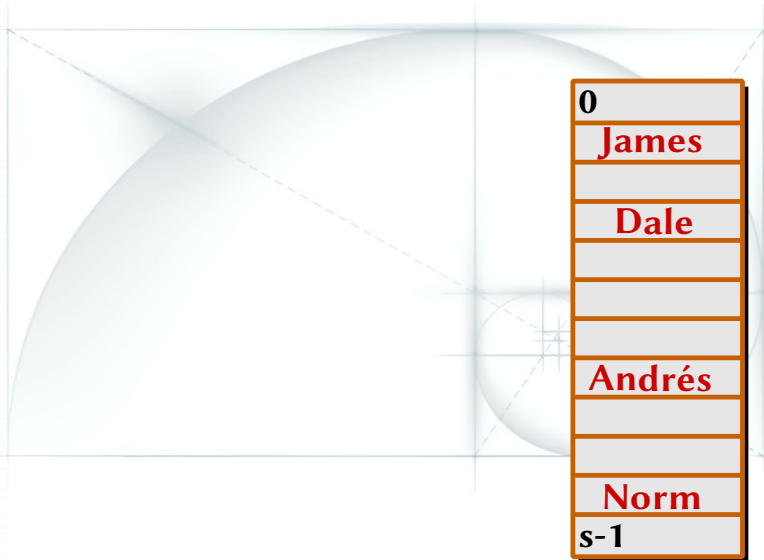




Ways to Fail

(a very incomplete list)

Effects of become:



Effects of become:

'Norm' oneWayBecome: 'Dale'

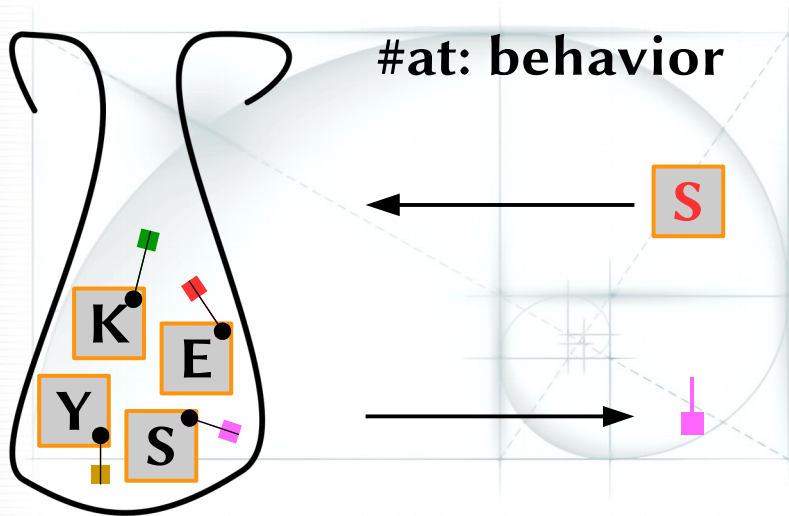
0
James
Dale
Andrés
Dale
s-1

Effects of become:

'Norm' become: 'Dale'

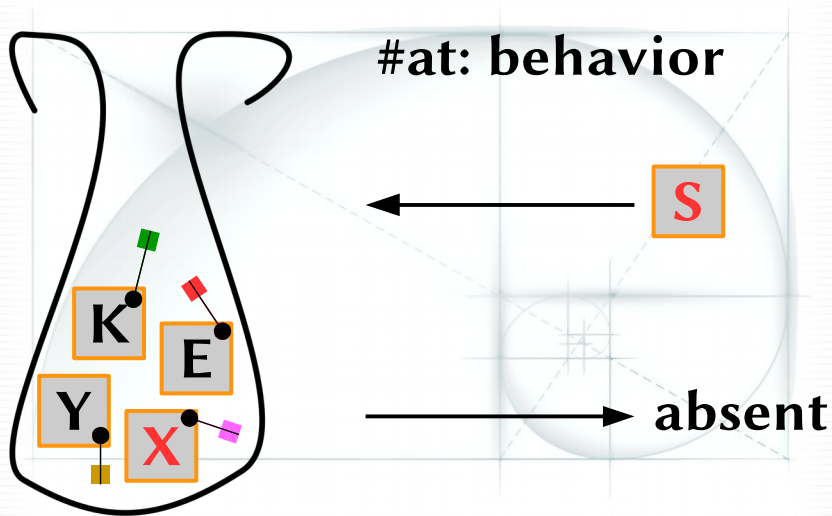
0
James
Norm
Andrés
Dale
s-1

Key mutation



Dictionary

Key mutation



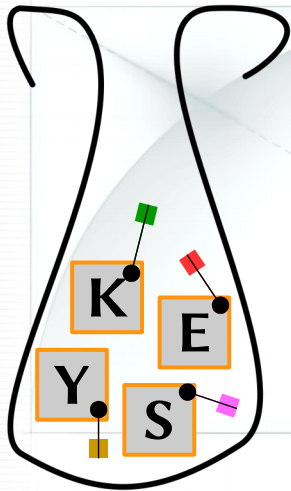
#at: behavior

S

absent

Dictionary

Threading



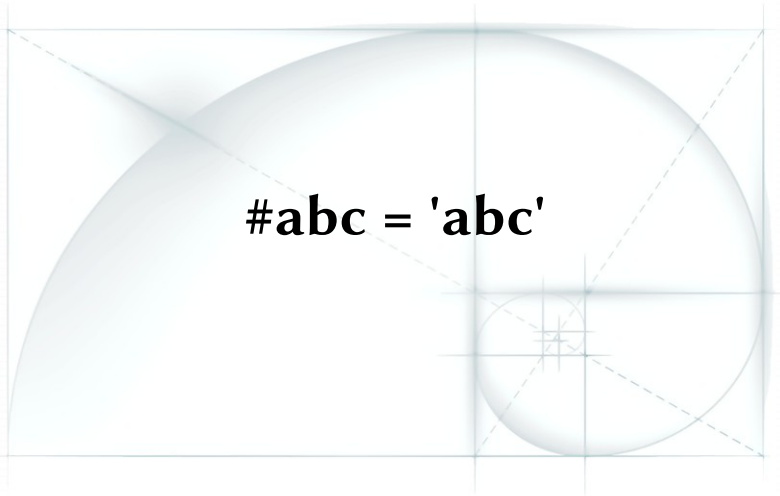
#at:put: ...

<process switch>

#at:put: ...

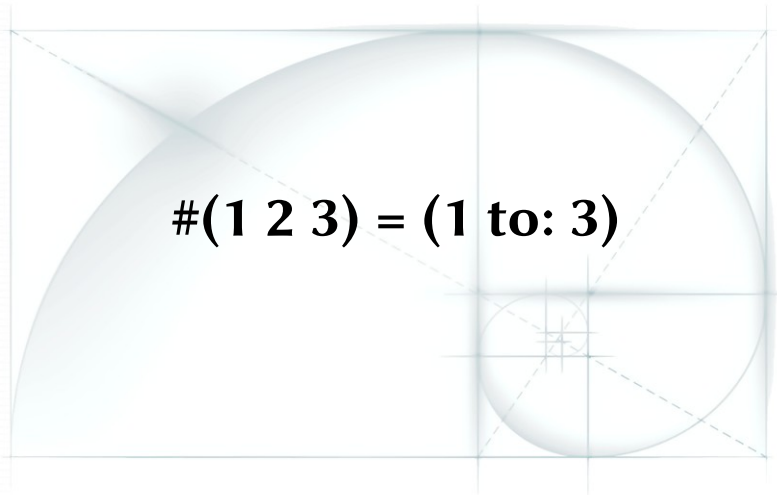
Dictionary

Funny equal methods

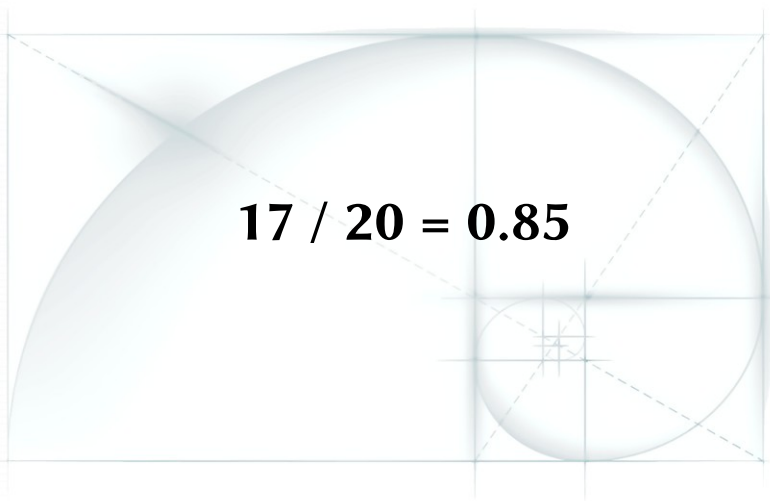


Funny equal methods

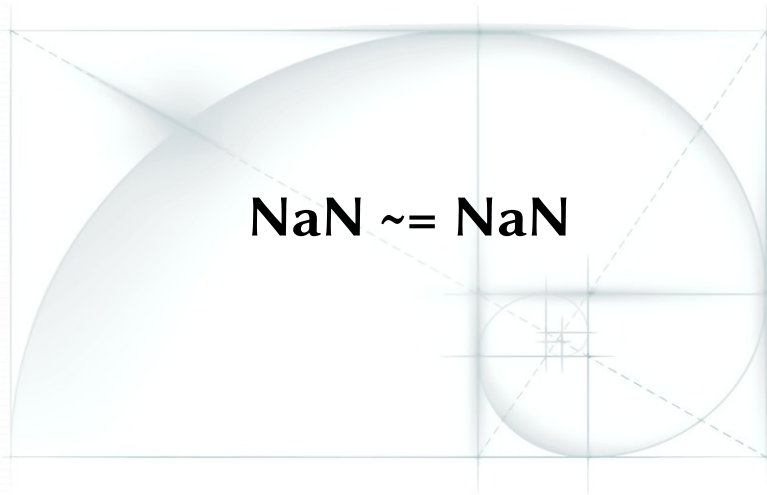
$\#(1\ 2\ 3) = (1\ \text{to:}\ 3)$



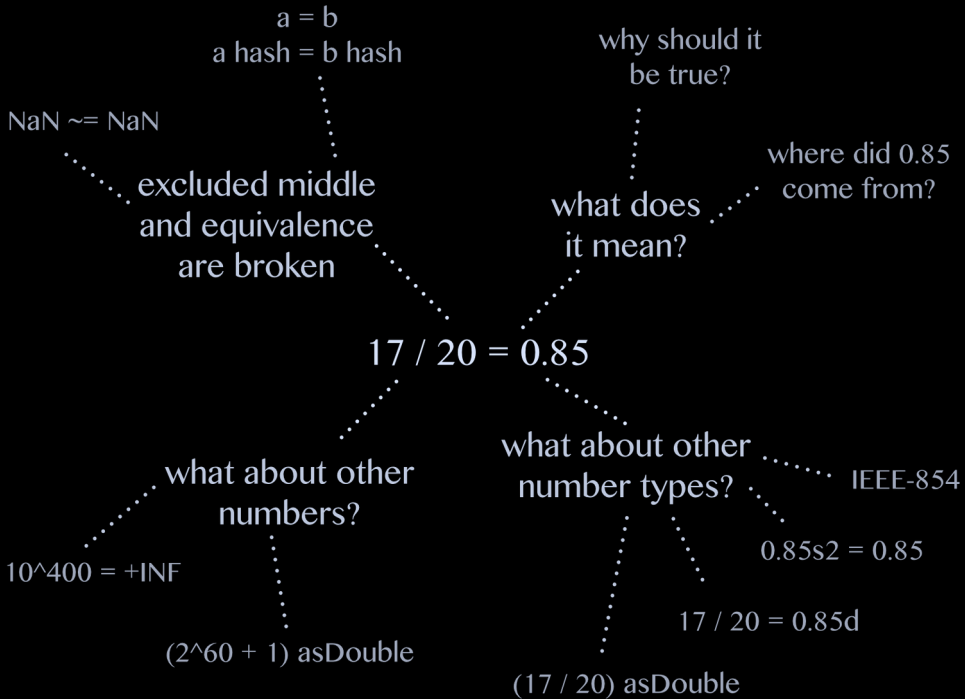
Funny equal methods



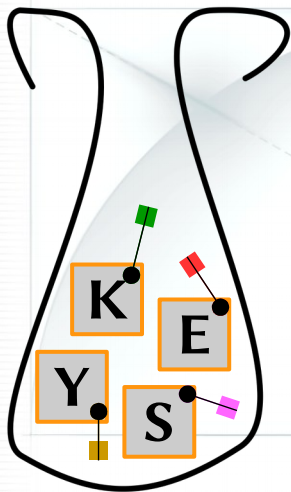
Funny equal methods



NaN \approx NaN




Reverse hashing



keyAtValue: 

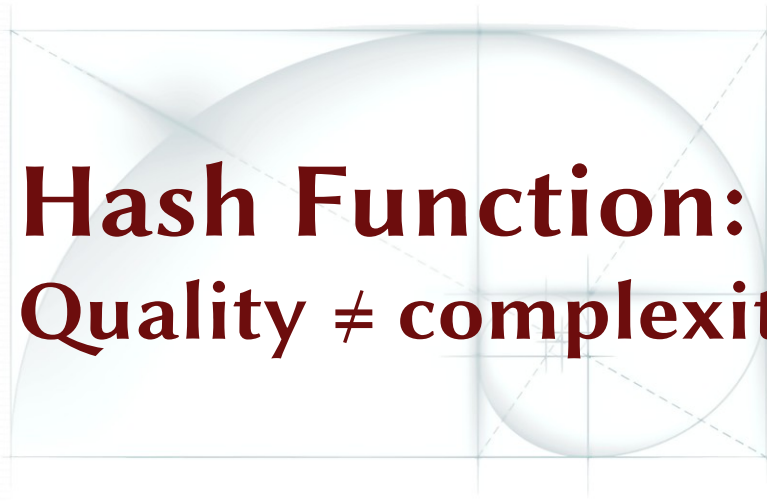
Dictionary

A light blue sphere is centered on a white background. A grid of thin blue lines is overlaid on the sphere. A dashed blue square is drawn around the sphere, and a dashed blue circle is drawn around the sphere's equator. The text "Final Thoughts" is written in a bold, dark red serif font across the center of the sphere.


Final Thoughts



**Hash Function:
Quality over
speed**



**Hash Function:
Quality \neq complexity**



**Hash Method
is only half of
Hash Function**



**Test
Profile
TEST**

Hashed Collections

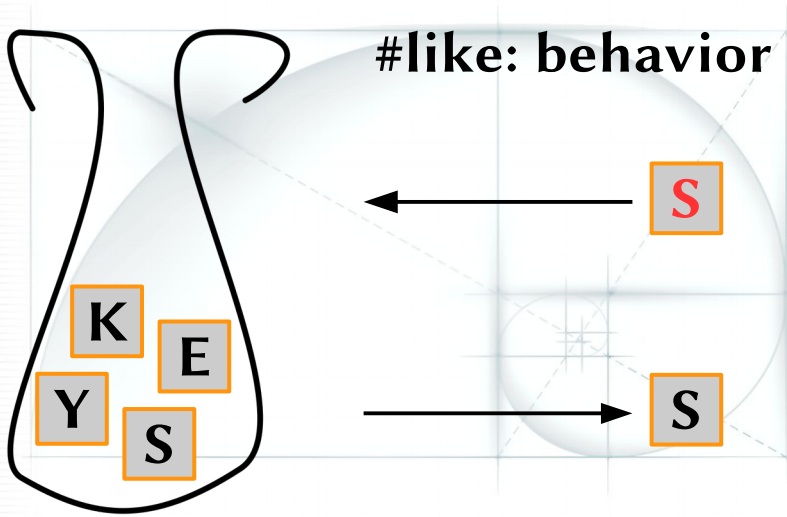
...it julienne fries!

Martin McClure



Andrés Valloud





#like: behavior

S

S

Set

JDK 1.6

```
/**  
 * Applies a supplemental hash function to a given hashCode, which  
 * defends against poor quality hash functions. This is critical  
 * because HashMap uses power-of-two length hash tables, that  
 * otherwise encounter collisions for hashCodes that do not differ  
 * in lower bits. Note: Null keys always map to hash 0, thus index 0.  
 */  
static int hash(int h) {  
    // This function ensures that hashCodes that differ only by  
    // constant multiples at each bit position have a bounded  
    // number of collisions (approximately 8 at default load factor).  
    h ^= (h >>> 20) ^ (h >>> 12);  
    return h ^ (h >>> 7) ^ (h >>> 4);  
}
```

Multiplicative hash function

Magic

16,777,619

1, 099, 511, 628, 211

309, 485, 009, 821, 345, 068, 724, 781, 401

374, 144, 419, 156, 711, 147, 060, 143, 317, 175,
368, 453, 031, 918, 731, 002, 211 35, 835, 915,
874, 844, 867, 368, 919, 076, 489, 095, 108, 449,
946, 327, 955, 754, 392 558, 399, 825, 615, 420,
669, 938, 882, 575, 126, 094, 039, 892, 345, 713,
852, 759

5, 016, 456, 510, 113, 118, 655, 434, 598, 811,
035, 278, 955, 030, 765, 345, 404, 790 744, 303,
017, 523, 831, 112, 055, 108, 147, 451, 509, 157,
692, 220, 295, 382, 716 162, 651, 878, 526, 895,
249, 385, 292, 291, 816, 524, 375, 083, 746, 691,
371, 804 094, 271, 873, 160, 484, 737, 966, 720,
260, 389, 217, 684, 476, 157, 468, 082, 573

Hash function

Fowler, Noll, Vo

Multiplicative hash function

Magic	Hash function
0	Black hole
1	Identity
3.14	Full circle
256	Byte oriented
2018	ESUG 2018 Conference
1234567890	Diversity
9876543210	Diversity the Great
$10^{10^{100}}$	googleplex