# CS180 - Hash Tables (part 2)

## Announcements

- Checkpoint is tomorrow

- Review session: Dec. 10
    NOT: 8am, 10am, 2-6
    
$\rightarrow$ | 10 am - noon |

# Hashing

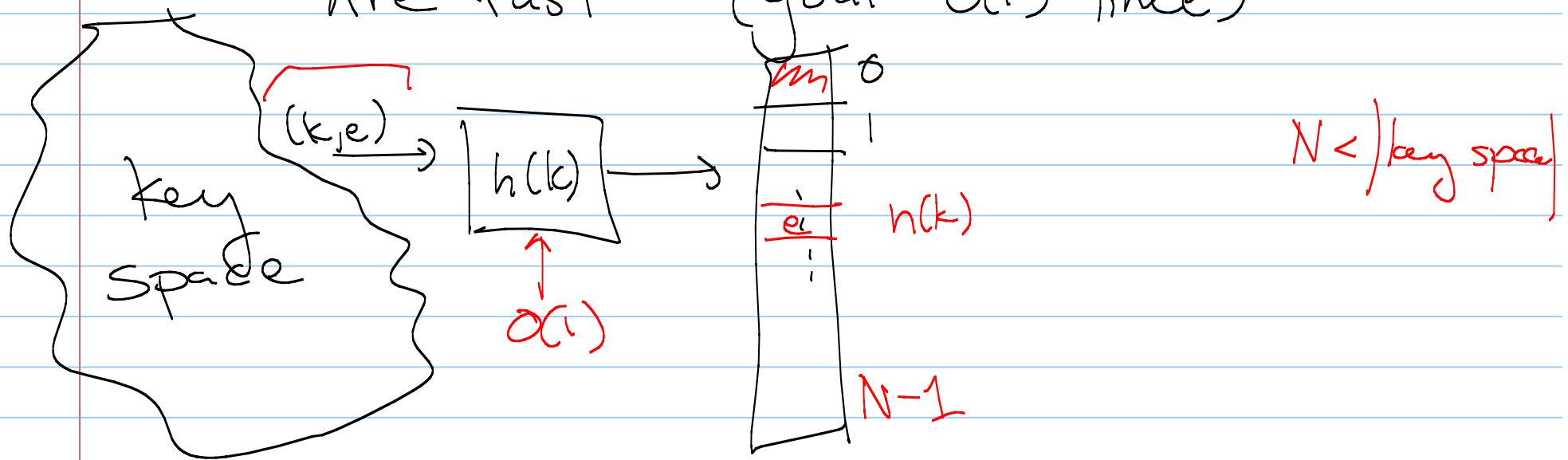An array is not very space efficient. We would like to take the key & make is smaller.

A hash function h maps each key in our dictionary to an integer in the range $[0, N-1]$.

(N should be much smaller than the # of keys.)

Then we store $(k, e)$ in $A[h(k)]$

# Good hash functions:

— Are fast    (goal: $O(1)$ time)



Key space

$(k, e) \rightarrow$ $h(k) \rightarrow$

$O(1)$

0
1
$e_i$    $h(k)$
$N-1$

$N < |$key space$|$

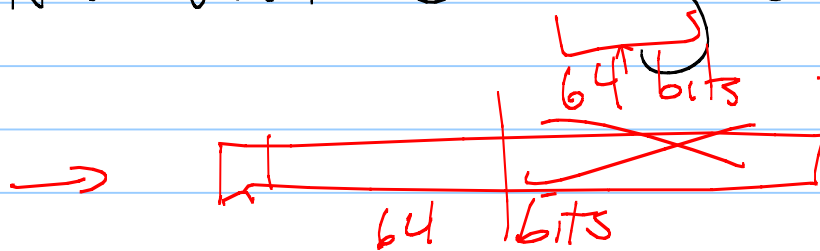— Don't have collisions.

Collisions are unavoidable.

when $k_1 \neq k_2$
but $h(k_1) = h(k_2)$

(1) First: map key to a number

32 bits

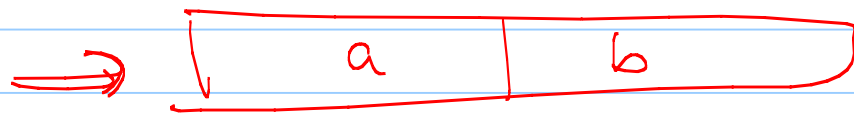Say we want keys to fit in an int.

What can we do for int, char, & short types?

32 bits   32 bits   32 bits

Now what about long or float?

64 bits

64 bits

$a \gg 32$

| a | b |

$a + b$ ← simplist way to hash

This can backfire. Remember ASCII?

128-bits (full newest version)

$$q_0 \, q_1 \, q_2 \, q_3 \, q_4 \, q_5 = \overbrace{\qquad} = \overbrace{\qquad}$$

temp01 and temp10 & pm0te1

$$\underbrace{q_0 + q_1 + q_2 + q_3 + q_4 + q_5}$$

all will go to same #

Goal:

Better way to avoid collisions between "similar" keys.

A better idea: Polynomial Hash codes

Pick $\boxed{a \neq 1}$ and split data into k 32-bit parts
$$(x_0, x_1, \ldots, x_{k-1}) = x$$

Let $h(x) = x_0 a^{k-1} + x_1 a^{k-2} + \cdots + x_{k-2}a + x_{k-1}$

$\overset{x_0}{"t"}\ \overset{x_1}{"e"}$

temp 01

$a = 37$

$"t" \cdot 37^5 + "e" \cdot 37^4 + "m" \cdot 37^3 + "p" \cdot 37^2 +$
$\qquad\qquad "0" \cdot 37 + "1" \cdot 37^0$

temp10

$"1"$

$\left[\ "1" \cdot 37 + "0" \cdot 37^0 \right.$

Aside: Efficiency

$$\sum_{i=1}^{n} i = O(n^2)$$

$$x_0 \cdot a^4 + x_1 \cdot a^3 + x_2 \cdot a^2 + x_3 \cdot a + x_4$$

Horner's rule: $x_{k-1} + a(x_{k-2} + a(x_{k-3} + \cdots)))$

$$x_4 + a(x_3 + a(x_2 + a(x_1 + a x_0)))$$

$$O(n)$$

This strategy makes it less likely that "similar" words/data will collide.

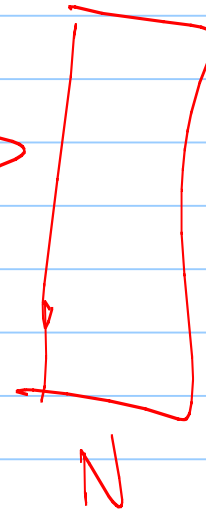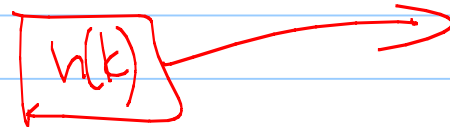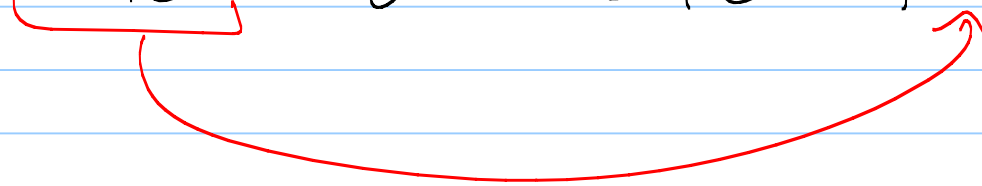What about overflow? (Remember, we want only 32-bits in key.)

chop it at 32-bits

(not so great)

# Cyclic Shift Hash codes

shift bits in representation somehow

$$1010001\,0\cdots 0\,1\,000\,1$$

$h(k)$

N

# Compression Map:

**2:** Once we have an integer key representation;
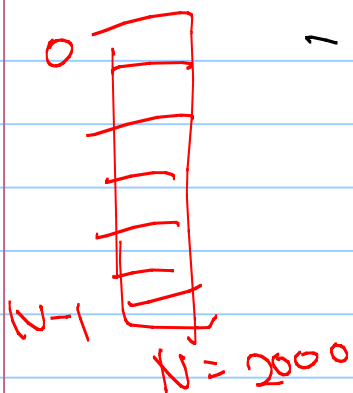
Need to make sure it is between $0$ & $N-1$, so is in our array.

Ideas?  map everything to $0$  lots of collisions!

Want to spread things out evenly

— modular arithmetic

32-bit integer

$$h(k) \mod N$$

$h(0)$

0

$N-1$

$N = 2000$

Compressing number down to something
between 0 and N-1

Ideas.

Modular arithmetic

$h(x) \bmod N$

lut

$x \bmod n$

# Example



A:

| | (12,E) | | | | | | | | | (21,R) |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

- Map is $h(k) = k \mod \boxed{11}$

insert( 12 , E )    $h(12) = 1$    E goes in A[1]
insert( 21 , R )    $h(21) = 10$
insert( 37 , I )    $h(37) = 4$
insert( 26 , N )    $h(26) = 4$
insert( 16 , C )    $h(16) = 5$
insert( 5 , H )    $h(5) = 5$

key & data

Strategy 1: go to next open spot

strategy 2: Secondary hash

Strategy 3: chaining

Some comments:

This works better if size of table is a prime number.

Why?

Go take number theory

(check book)

The MAD (multiply add + divide) method
is a bit better:

$$h(x) = |ax + b| \mod N$$

where a + b are

- not equal
- relatively prime
$$gcd(a, b) = 1$$

(if not)

Collisions
get more
frequent

21
3, 7

20
2, 4, 5, 10

← even better)
2 prime
numbers

- less than N

Goal: Simple Uniform Hashing Assumption:

For all $k \in$ keyspace,

$$\Pr[h(k) = i] = \frac{1}{N}$$

(Essentially, elements are "thrown randomly" into the buckets)

## Collisions

Can we ever totally avoid collisions?

No — goal was to minimize them,

How to deal with them?

(had 3 strategies)

How can we handle collisions?

(Do we have data structures to store
more than one thing??)

- list ← bad search

- vectors ← insert can be bad

- trees ← $O(\log n)$
  (in between lists & vectors)