# Algorithms in Bioinformatics

## Exact pattern matching

# Recap

- HW up, due next Tuesday
- No midterm - longer reading assignment instead
- Tomorrow: no office hours 2-3pm (sorry!)
  Instead: tomorrow 11:30-noon
  Friday: 11-12
  ↳ but please email to set up!

# Today: Exact matching & repeat finding

- 50% of human genome is repeats

- However, repeats are also important!
  (Go read Section 9.1 - associated with disease, evolution, etc.)

In particular: long, maximal repeats.

Different from motifs:
## pattern is known

# First tool: Hashing

Hash tables, dictionaries or associative arrays, are built into most languages these days.

## Side note:

**Don't ever implement these yourself!**

## Hashing's goal:

- Book says duplicate removal.

This is one goal — but not the biggest one from a CSV or BCB perspective!

# Hashing: Fast data storage

Given key/value pairs, want to be able to retrieve value ~~quickly~~ given the key. $O(1)$

(As well as store/update)

## Examples:

- Course # & schedule info
- URL and html page
- Flight # & arrival info
- Color and BMP
- Directors & movies
- l-mers & repetition locations in a sequence

# Dictionary

A data structure which supports:

3 supported funs
- insert (key, data)
- find (key)
- remove (key)

$\left.\right\}$ O(1)

## Note: An array is a kind of dictionary!

key: index/position

data: stored value



Other implementations:

Linked List:



Vectors

# Hashing

Assume $m \gg n$, so

possible keys ↑     ↑ #'of entries

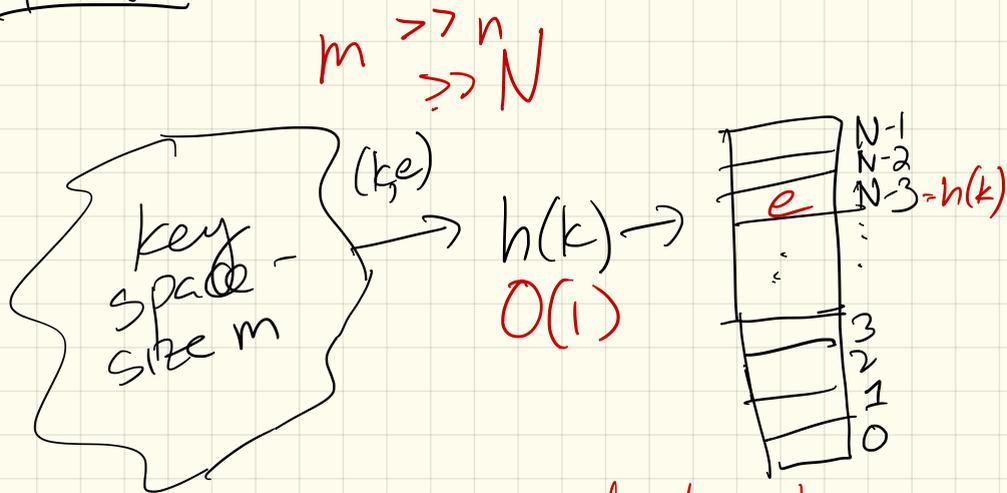array takes too much space.

<u>Goal</u>: $O(n)$ space

fast lookup/insert/
remove

A <u>hash function</u> $h$ maps
each key to an integer
in range $[0 .. N-1]$

Goal: $N$ is bigger than $n$,
but much smaller
than $m$.

Then: Given $(k, e)$, store
it in $A[h(k)]$ (in an
array).

# Picture:

$m >> n$
$>> N$



key space - size $m$

$(k,e)$

$h(k) \rightarrow$
$O(1)$

N-1
N-2
$e$  N-3 = $h(k)$
:
:
3
2
1
0

$n$ actual values stored in array

## Good hash functions:

- are fast $O(1)$!

- avoid collisions

↳ if $k \neq k'$
want $h(k) \neq h(k')$
with high probability

So, how to do this?

① Make the key a #

② Compress # to $[0,...,N-1]$

③ Handle Collisions

①&② : often combined,
& saw some of it
in data structure

We'll recap a bit...

# First idea

For something like ASCII,
can break into pieces & treat
as bits!

$$E \quad r \quad i \quad n$$

$$69 + 114 + 105 + 110 = \# \ddot\smile$$

Then what?

## Problem: this can backfire
w/ words:

$$h(temp01) = h(temp10)$$
$$\neq h(pm0te1)$$

Want to avoid collisions.

So...

# Polynomial Hash Codes

Split date to 32-bit pieces.

$$x = (x_0, \dots, x_{k-1})$$

Pick $a \neq 1$.

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

$\uparrow t$     $\uparrow e$     $\uparrow m$     $\uparrow p$

## Ex: Erin (or $\underline{69}, \underline{105}, \underline{114}, \underline{110}$)
and $a = \underline{37}$:
$$p(x) = 69 \cdot 37^3 + 105 \cdot 37^2$$
$$+ 114 \cdot 37 + 110$$

## Why?
- relatively fast
- avoids collisions!

(more tricks like this)

<u>Next</u>:  Compress:

$h(k) \longrightarrow$

#       b/t
0 ↔ N-1

N-1
⋮
1
0

<u>Idea</u>:  Take $h(k) \bmod N$

Recall:  $3 \bmod 10 = 3$       (Python    % in C)

$50 \bmod 10 = 0$

$14 \bmod 10 = 4$

# Example: $h(k) = k \bmod \underset{N}{11}$

A:

$(12,E)$ … $(37,I)$ $(16,N)$ … $(4,R)$

0 1 2 3 4 5 6 7 8 9 10

Insert: key

$(12, E)$ : $h(12) =$ 12 mod 11 = 1

$(21, R)$ : $h(21) = 10$

$(37, I)$ $h(37) = 4$

$(16, N)$ $h(16) = 5$

$(26, C)$ $h(26) = 4$  X

$(5, H)$

Comment: Works best if #s are prime.

^ relatively

Why?

go take number theory

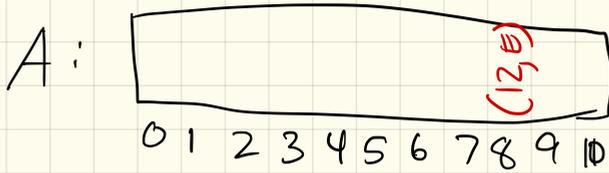Another way: M.A.D
  Instead of $h(k)$ mod N,
  do $\underline{h(k)} = |\underline{a}k + \underline{b}|$ mod N
  where a + b are:
      - relatively prime
      - less than N

Why? go take NT

# Example:   $h(k) = 3k+5 \mod 11$

A:

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

Insert:

$h(12) =$

(12, E) $= 3 \cdot 12 + 5 \mod 11 = 8$

(21, R)

(37, I)

(16, N)

(26, C)
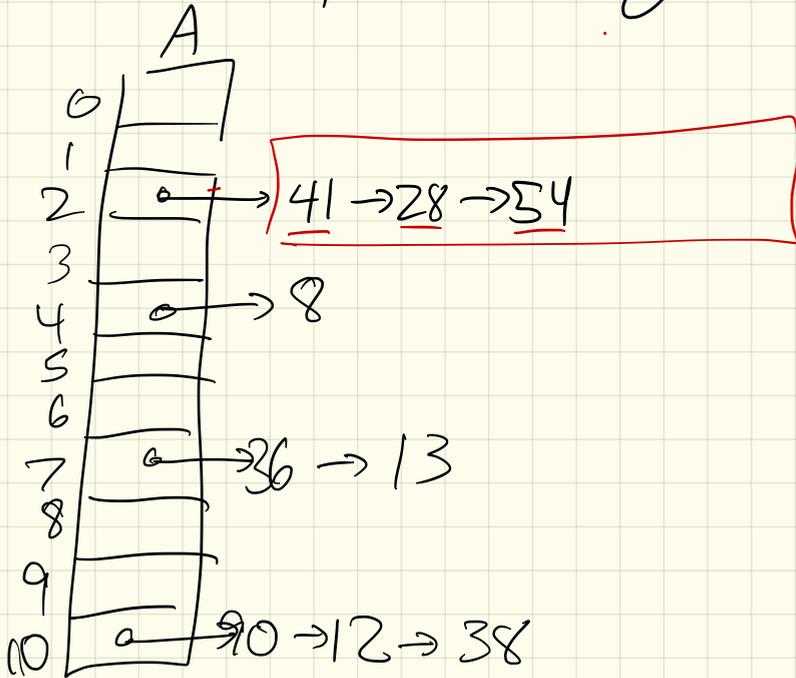
(5, H)

(Collisions may still happen)

Why bother?

MUCH better in practice

# Step 3: Handle Collisions
(Hint: What data structures can store more than 1 thing??)

## Ex: Simple Chaining:

A

```
 0
 1
 2  •───────→ 41 →28 →54
 3
 4  •──→ 8
 5
 6
 7  •───→ 36 → 13
 8
 9
10  •──→ 10 →12→ 38
```

Run times:

<span style="color:red">Worst case, bad hash function.
↳ insert/lookup list time</span>

# Other techniques:

- linear probing
- quadratic probing
- re-hashing

# Takeaway:

Handle collisions
On most data, all of
these work well in
practice.
(No theoretical guarantee)

# Load Factors

Whatever method you use, usually starts to do badly if

n gets close to N:

Want $\frac{n}{N} < .5$

more than half full

Re hashing:

When more than half full, most implementations double the array size + choose a new hash function

(Hence, don't write these yourself!)

# Back to pattern matching

## Naive Pattern matching:

pattern
⌐input string (longer)

```python
def naive(p, t):
    occurrences = []
    for i in range(len(t) - len(p) + 1):   # loop over alignments
        match = True
        for j in range(len(p)):            # loop over characters
            if t[i+j] != p[j]:             # compare characters
                match = False              # mismatch; reject alignment
                break
        if match:
            occurrences.append(i)          # all chars matched; record
    return occurrences
```

P: **word**
T: **There would have been a time for such a word**
       ------ **word** ------- **word** --------------------→ **word**
          ---→     -→                   ------→

## Runtime:

$$|P| = n$$
$$|T| = m$$
$$n(m-n+1) = O(mn)$$

# [Boyer -Moore]

How to improve?

① Skip pointless alignments:
("Bad character rule")

Align P at start of T:

ⓐ Look at position of the
last occurence of a
mismatching character

If this character exists in
pattern, realign to last
P(prior) occurence

Step 1:
```
T:  G C T T Ⓒ T G C T A C C T T T T G C G C G C G C G C G G A A
P:  C Ⓒ T T T T G C                                    Case (a)
```

ⓑ If that character isn't in
pattern, just go past entirely

Step 1:
```
T:  G C T T Ⓒ T G C T A C C T T T T G C G C G C G C G C G G A A
P:  C Ⓒ T T T T G C                                    Case (a)
```

Step 2:
```
T:  G C T T C T G C T Ⓐ C C T T T T G C G C G C G C G C G G A A
P:        C C T T T T G C                              Case (b)
```

© Of course, if you don't find this 'character', you've hit a match!

**Step 1:**
T: **G C T T C T G C T A C C T T T T G C G C G C G C G C G G A A**
P: **C C T T T T G C**
*Case (a)*

**Step 2:**
T: **G C T T C T G C T A C C T T T T G C G C G C G C G C G G A A**
P: **C C T T T T G C**
*Case (b)*

**Step 3:**
T: **G C T T C T G C T A C C T T T T G C G C G C G C G C G G A A**
P: **C C T T T T G C**
*Case (c)*

Run-time of this:

Still O(mn), since could get all the same character:

AAA A
AA AA ---- A

# ② Good suffix rule

Let $t$ = substring matched
   by the inner loop

Look at suffixes:

Skip until either

ⓐ no mismatches
   between P & $t$

ⓑ P moves past $t$

**Step 1:**

T: C G T G C **C** `TAC` T T A C T T A C T T A C T T A C G C G A A
P: C T `TAC` **T** T A C

**Step 2:**

T: C G T G C **C** `TACTTAC` T T A C T T A C T T A C G C G A A
P: `CTTAC` T T A C

**Step 3:**

T: C G T G C C T A C T T A C T T A C T T A C T T A C G C G A A
P: C T T A C T T A C

# Note: Can break @ down to cases:

Step 1:
T: C G T G C **C TAC** T T A C T T A C T T A C T T A C G C G A A
P: C T **TAC** **T TAC**

*t* occurs *in its entirety* to the left within P

Step 2:
T: C G T G C **C TAC T TAC** T T A C T T A C T T A C G C G A A
P: **C T TAC** **T TAC**

*prefix* of P matches a *suffix* of *t*

Step 3:
T: C G T G C C T A C T T A C T T A C T T A C T T A C G C G A A
P: C T T A C T T A C

# Algorithm:
## Tradeoff b/t these rules
## just suffix rule: O(mn)