

CSE 373: Data Structures and Algorithms

Lecture 7: Hash Table Collisions

Instructor: Lilian de Greef
Quarter: Summer 2017

Today

- Announcements
- Hash Table Collisions
- Collision Resolution Schemes
 - Separate Chaining
 - Open Addressing / Probing
 - Linear Probing
 - Quadratic Probing
 - Double Hashing
- Rehashing

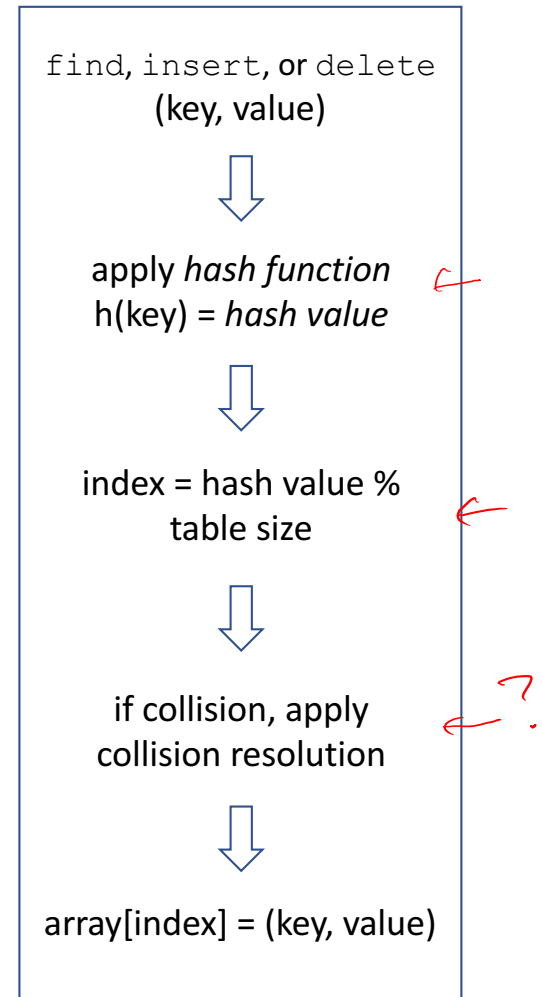
Announcements

- Reminder: homework 2 due tomorrow
- Homework 3: Hash Tables
 - Will be out tomorrow night
 - Pair-programming opportunity! (work with a partner)
 - Ideas for finding partner: before/after class, section, Piazza
- Pair-programming: write code together
 - 2 people, 1 keyboard
 - One is the “navigator,” the other the “driver”
 - Regularly switch off to spend equal time in both roles
 - Side note: our brains tend to edit out when we make typos
 - Need to be in same physical space for entire assignment, so partner and plan accordingly!

Review: Hash Tables & Collisions

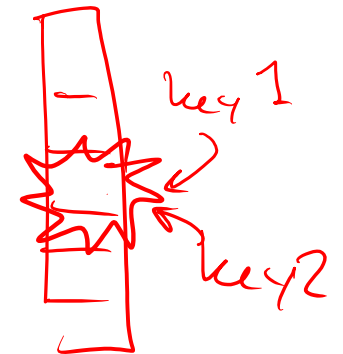
Hash Tables: Review

- A data-structure for the dictionary ADT
- *Average case* $O(1)$ find, insert, and delete (when under some often-reasonable *assumptions*)
- An array storing (key, value) pairs
- Use *hash value* and table size to calculate array index
- Hash value calculated from key using *hash function*

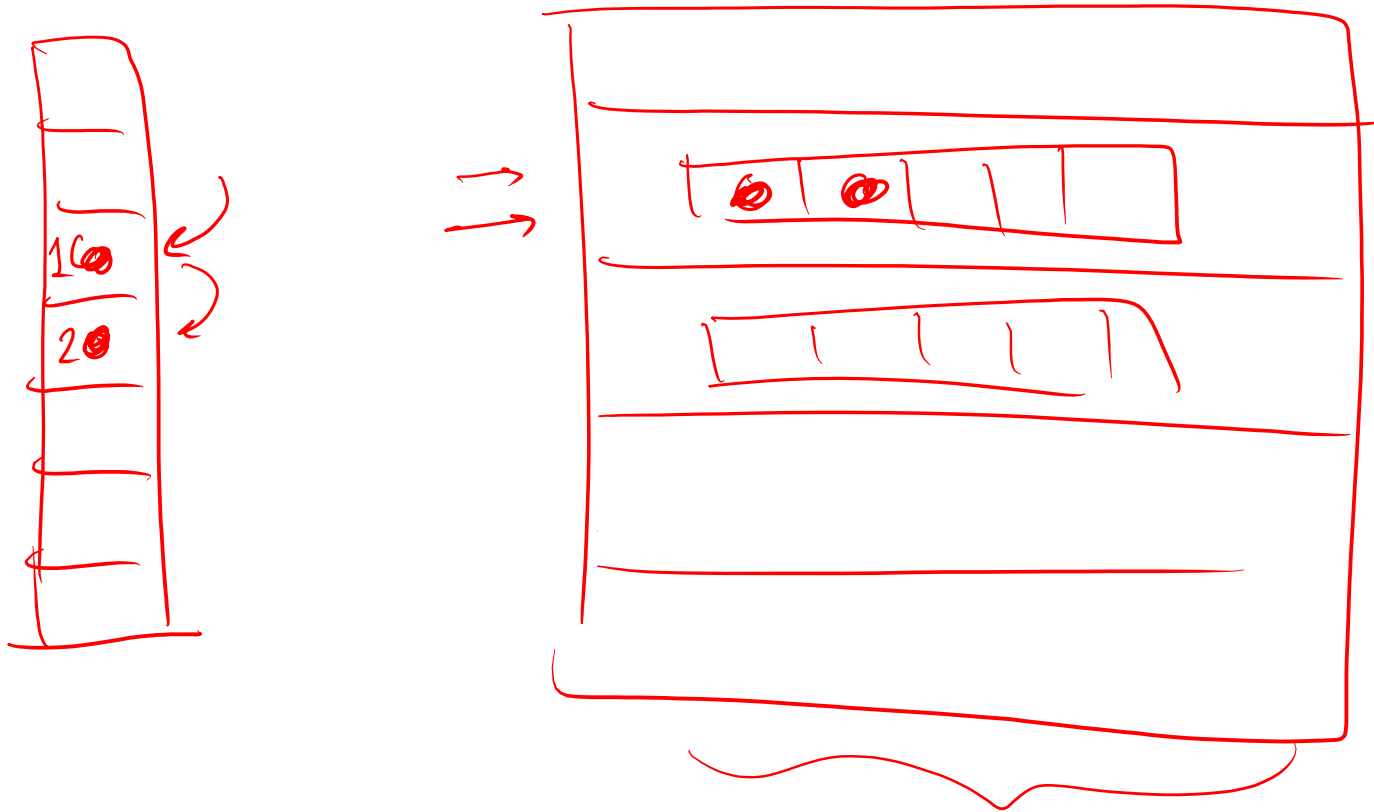


Hash Table Collisions: Review

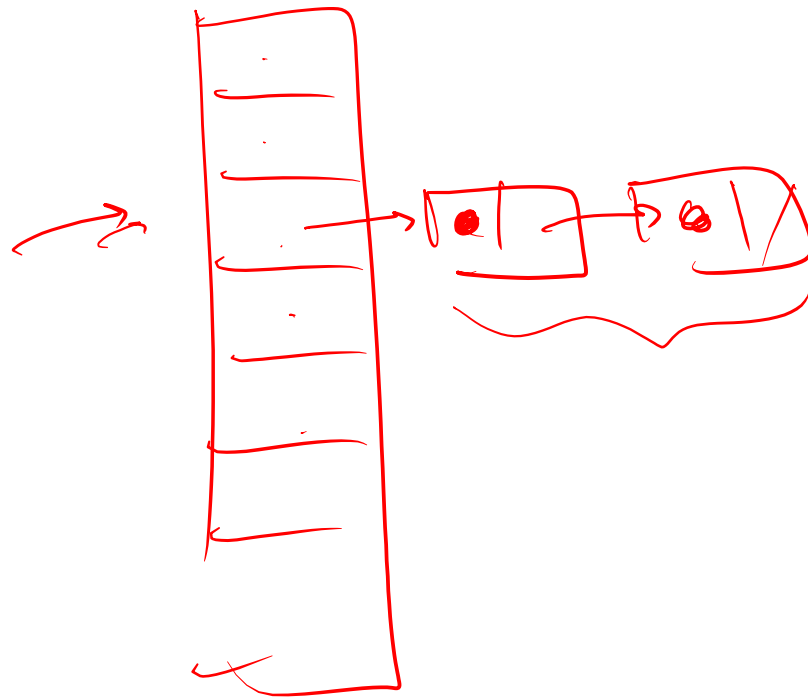
- Collision: when two keys map to the same location in the hash table
- We try to *avoid* them by having a good hash function (unique indexes)
- Unfortunately, collisions are unavoidable in practice
 - Number of possible keys \gg table size
 - No perfect hash function & table-index combo



Collision Resolution Schemes: your ideas



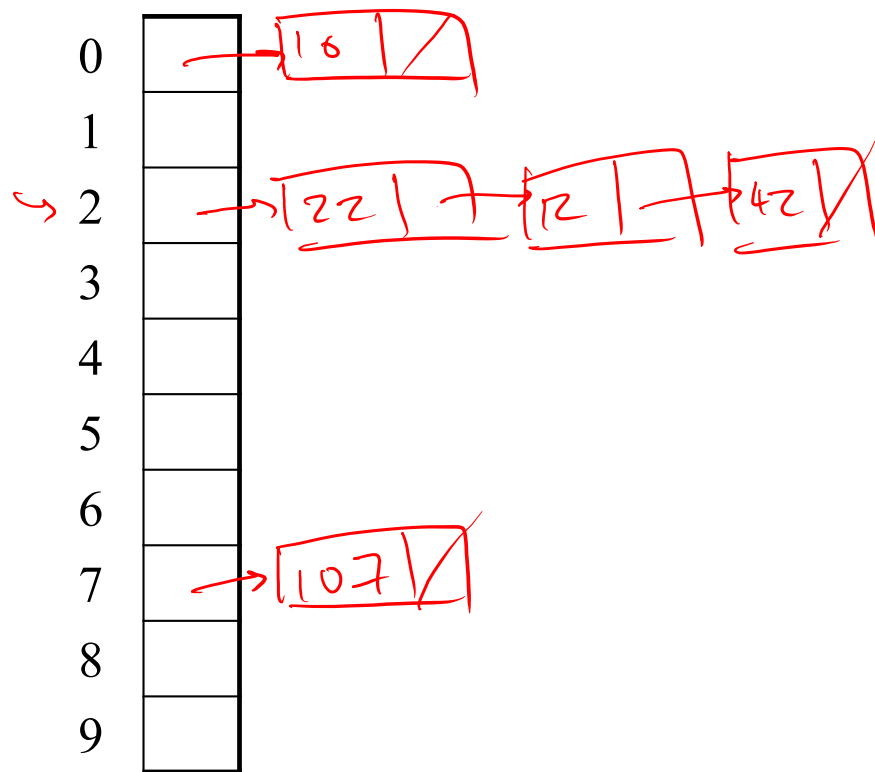
Collision Resolution Schemes: your ideas



Separate Chaining

One of several collision resolution schemes

Separate Chaining



All keys that map to the same table location (aka “bucket”) are kept in a list (“chain”).

Example:

insert 10, 22, 107, 12, 42

and **TableSize** = 10

(for illustrative purposes,
we’re inserting hash values)

Separate Chaining: Worst-Case

What's the worst-case scenario for `find`?

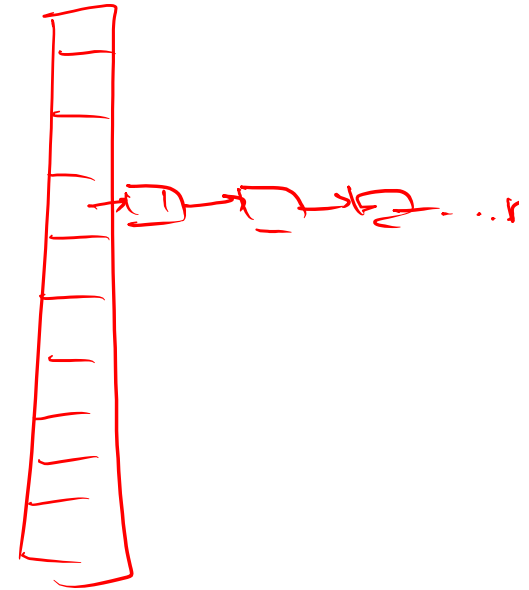
all keys indexed to the
same bucket

What's the worst-case running time for `find`?

$O(n)$ linear

But only with really bad luck or really bad hash function

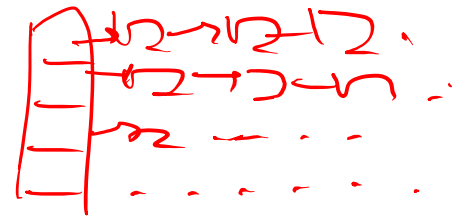
↳ not worth avoiding
worst-case



Separate Chaining: Further Analysis

- How can find become slow when we have a good hash function?

way bigger
elements \gg table size
mean long chains



- How can we reduce its likelihood?

Maintain a good ratio of
elements to the table size
(resize the table as needed)

Rigorous Analysis: Load Factor

Definition: The **load factor** (λ) of a hash table with N elements is

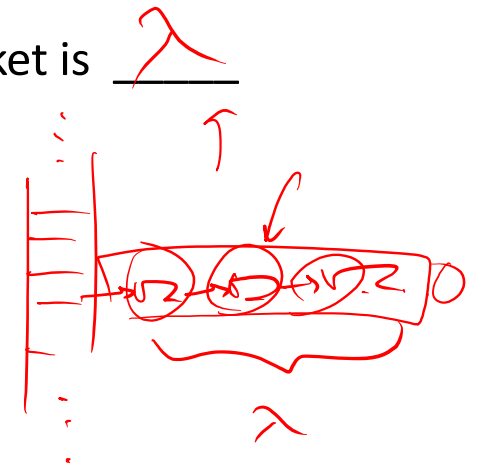
$$\lambda = \frac{N}{\text{table size}}$$

*$N = \text{table size}$
average: 1 element/bucket*

Under separate chaining, the average number of elements per bucket is λ

For a *random* find, on average

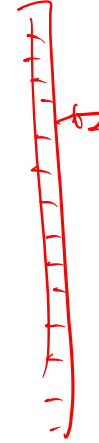
- an unsuccessful find compares against λ items
- a successful find compares against $\lambda/2$ items



Rigorous Analysis: Load Factor

Definition: The load factor (λ) of a hash table with N elements is

$$\lambda = \frac{N}{\text{table size}}$$



To choose a good load factor, what are our goals?

- short chains (not too high)
- efficient use of table space (not too low)

So for separate chaining, a good load factor is 1, 1.5, or 2

Open Addressing / Probing

Another family of collision resolution schemes

Idea: use empty space in the table

0	8
1	109
2	10
3	
4	
5	
6	
7	
8	38
9	19

- If $h(\text{key})$ is already full,
 - try $(h(\text{key}) + 1) \% \text{TableSize}$. If full,
 - try $(h(\text{key}) + 2) \% \text{TableSize}$. If full,
 - try $(h(\text{key}) + 3) \% \text{TableSize}$. If full...

- Example: insert 38, 19, 8, 109, 10

$$38 \% 10 = 8$$

↓
index

Open Addressing Terminology

Trying the next spot is called *probing* (also called *open addressing*)

- We just did *linear probing*
 i^{th} probe was $(h(\text{key}) + i) \% \text{TableSize}$
- In general have some *probe function* f and use
 $(h(\text{key}) + f(i)) \% \text{TableSize}$

Dictionary Operations with Open Addressing

`insert` finds an open table position using a probe function

What about `find`?

- must use same probe function to "retrace the trail"
- unsuccessful search when reach empty bucket

What about `delete`?

- use "lazy" deletion
 - ↳ replace element with marker/flag to say "no data here, but keep probing"
- Note: delete with separate chaining is plain-old list-remove



Practice:

The keys 12, 18, 13, 2, 3, 23, 5 and 15 are inserted into an initially empty hash table of length 10 using open addressing with hash function $h(k) = k \bmod 10$ and linear probing. What is the resultant hash table?

0	
1	
2	2
3	23
4	
5	15
6	
7	
8	18
9	

(A)

0	
1	
2	12
3	13
4	
5	5
6	
7	
8	18
9	

(B)

0	
1	
2	12
3	13
4	2
5	3
6	23
7	5
8	18
9	15

(C)

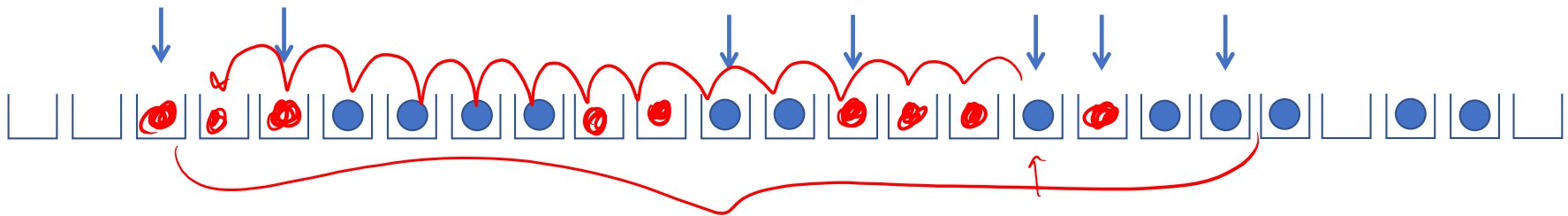
0	
1	
2	<u>12, 2</u>
3	<u>13, 3, 23</u>
4	
5	5, 15
6	
7	
8	18
9	

(D)

Open Addressing: Linear Probing

- Quick to compute! 😊
- But mostly a *bad idea*. Why?

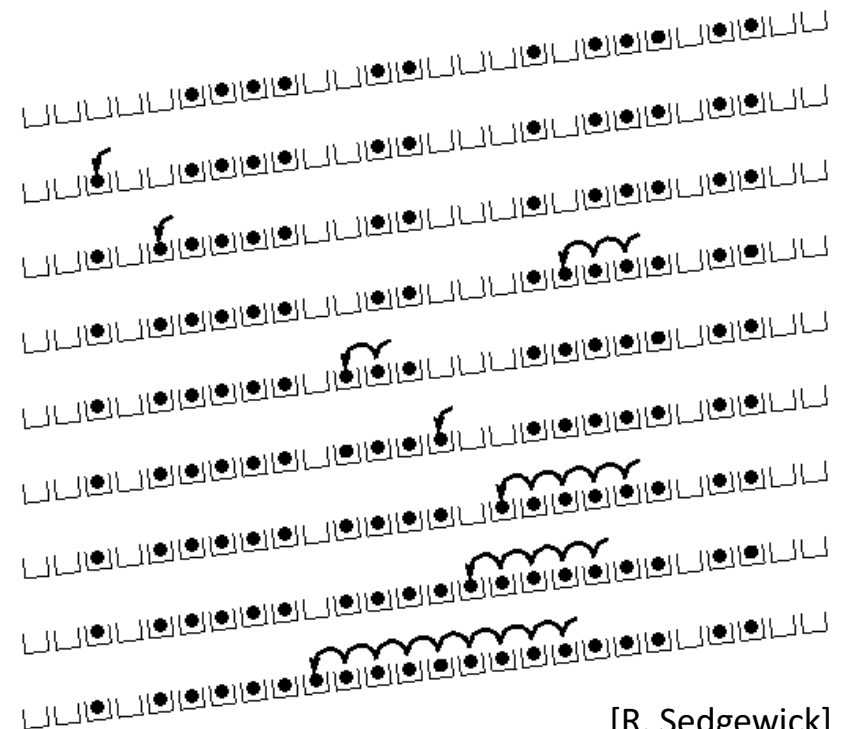
$\lambda > 1$ (#elements > table size)



(Primary) Clustering

Linear probing tends to produce ***clusters***, which lead to long probing sequences

- Called *primary clustering*
- Saw this starting in our example



[R. Sedgewick]

Analysis of Linear Probing

- For any $\lambda < 1$, linear probing will find an empty slot
 - It is “safe” in this sense: no infinite loop unless table is full

- Non-trivial facts we won’t prove:

Average # of probes given λ (in the limit as **TableSize** $\rightarrow \infty$)

- Unsuccessful search: $\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)^2} \right)$

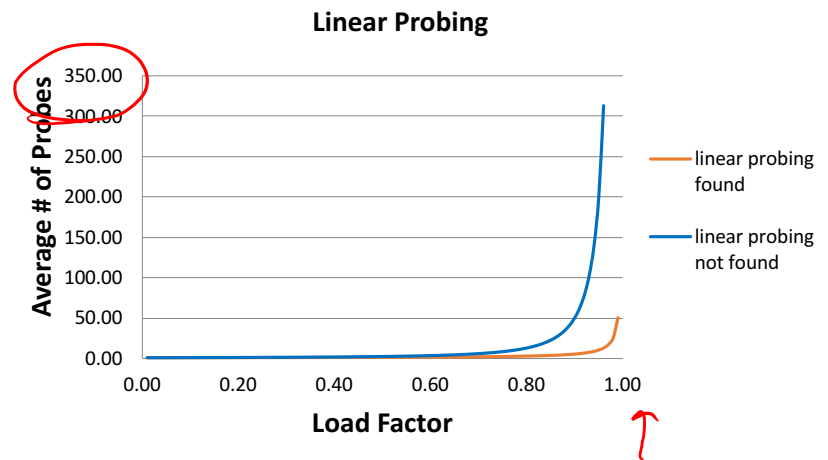
- Successful search: $\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)} \right)$

As lambda increases (gets closer to 1), so does the number of probes

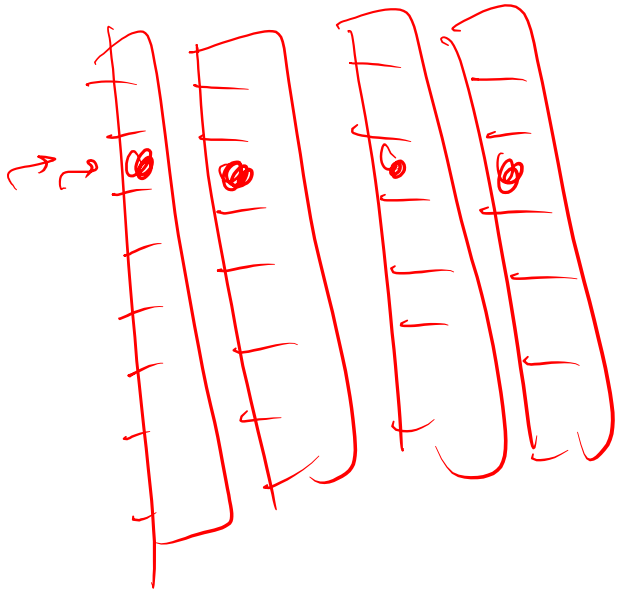
$\left\{ \begin{array}{l} \lambda = 0.75 \rightarrow \text{expect } \sim 8.5 \text{ probes} \\ \lambda = 0.9 \rightarrow \text{expect } \sim 50 \text{ probes} \end{array} \right.$

- This is pretty bad: need to leave sufficient empty space in the table to get decent performance (see chart)

- By comparison, chaining performance is linear in λ and has no trouble with $\lambda > 1$



Any ideas for alternatives?



Different
probe
function!

Open Addressing: Quadratic Probing

- We can avoid primary clustering by changing the probe function

$$(h(\text{key}) + f(i)) \% \text{TableSize}$$

(linear probing:
 $f(i) = i$)

- A common technique is quadratic probing: $f(i) = i^2$

- So probe sequence is:

- 0th probe: $h(\text{key}) \% \text{TableSize}$

- 1st probe: $(h(\text{key}) + 1) \% \text{table size}$

- 2nd probe: $(h(\text{key}) + 4) \% \text{table size}$

- 3rd probe: $(h(\text{key}) + 9) \% \text{table}$

- ...

- i^{th} probe: $(h(\text{key}) + i^2) \% \text{TableSize}$

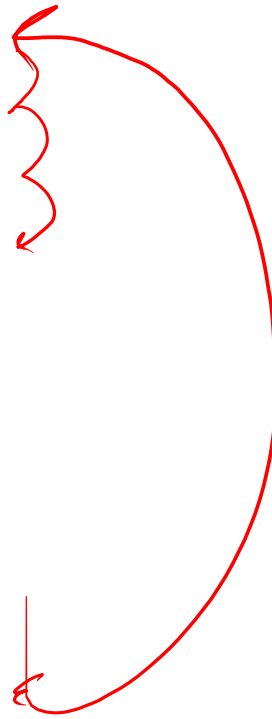
$$f(1) = 1^2 = 1$$

$$f(2) = 2^2 = 4$$

- Intuition: Probes quickly “leave the neighborhood”

Quadratic Probing Example #1

0	49
1	
2	58
3	79
4	
5	
6	
7	
8	18
9	89



TableSize = 10

Insert:

89

18

49

58

79

$9 + 1 \% 10 = 0$
 $(8 + 0 \% 10 = 9$
 $(8 + 4 \% 10 =$

$i^{\text{th}} \text{ probe: } (h(\text{key}) + i^2) \% \text{ TableSize}$

Quadratic Probing Example #2

0	48
1	
2	5
3	55
4	
5	40
6	76

TableSize = 7

Insert:

<u>76</u>	(76 % 7 = 6)
<u>40</u>	(40 % 7 = 5)
<u>48</u>	(48 % 7 = 6)
<u>5</u>	(5 % 7 = 5)
<u>55</u>	(55 % 7 = 6)
<u>47</u>	(47 % 7 = 5)

Yikes! For all n , $(n^2 + 5) \% 7$
is 0, 2, 5, or 6!
47

i^{th} probe: $(h(\text{key}) + i^2) \% \text{TableSize}$