

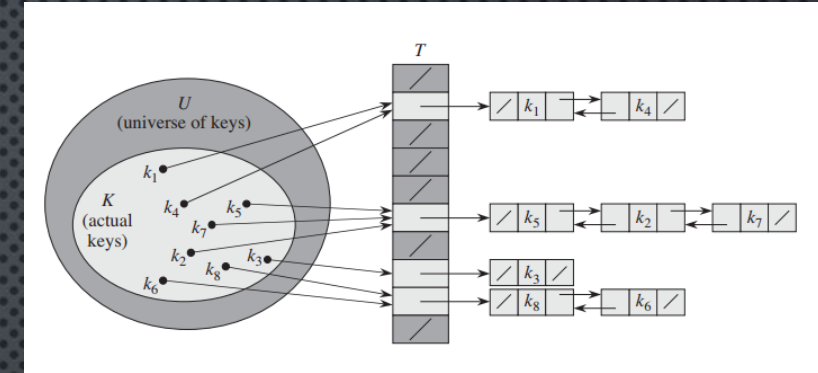
# PROBABILISTIC ANALYSIS AND HASHING

RECITATION WEEK 4



# HASH TABLES

- MANY APPLICATIONS REQUIRE SOME CONCEPT OF A **DICTIONARY**
  - A DICTIONARY STORES ELEMENTS AND CAN GROW (OR SHRINK) OVER TIME
    - **SEARCH**
    - **INSERT**
    - **DELETE**
- SOME DICTIONARIES, SUCH AS THE **HASH TABLE** USE ADDRESSES TO EXPEDITE THIS PROCESS
  - A HASH FUNCTION  $h$  TAKES THE **UNIVERSE OF POSSIBLE KEYS** AND ASSOCIATES THEM WITH A VALUE USING SOME FUNCTION
  - A KEY  $k$  IS THEN STORED AT LOCATION  $h(k)$  IN SOME AUXILIARY ARRAY
    - THE ARRAY IS USUALLY MUCH SMALLER THAN THE UNIVERSE OF KEYS – **COLLISIONS**
  - IF TWO ELEMENTS WOULD BE STORED AT THE SAME LOCATION, WE **CHAIN**
    - SO EACH INDEX OF THE AUXILIARY ARRAY POINTS TO A LINKED LIST
  - THE PERFORMANCE OF AN ALGORITHM USING A HASH TABLE IS DIRECTLY RELATED TO THE LENGTH OF THE CHAINS





# WARM-UP – HASHING AND RANDOMNESS

- TWO PROBLEMS TO THINK ABOUT
- WHAT IF OUR HASHING FUNCTION GOES WRONG?
  - SUPPOSE  $n$  ITEMS HAVE BEEN INSERTED INTO A HASH TABLE WITH  $m$  SLOTS. WHAT IS THE WORST-CASE TOTAL RUNNING TIME OF SEARCHING FOR EACH OF THESE ITEMS?
- “FINDING A 1”
  - INPUT: A HALF FULL HASH TABLE WITH UNKNOWN HASHING FUNCTION
  - OUTPUT: THE INDEX OF A USED HASH TABLE ENTRY

PROPOSE A DETERMINISTIC ALGORITHM AND ANALYZE ITS WORST-CASE RUNNING TIME AND THEN PROPOSE A RANDOMIZED ALGORITHM THAT YOU THINK WOULD BE BETTER



# RANDOMIZED ANALYSIS AND HASHING

- IN HASHING, WE ARE CONCERNED WITH COLLISIONS AS IT SLOWS DOWN HASH TABLE OPERATIONS AND FUNCTIONALITY
- SAY WE HAVE A UNIFORM HASHING FUNCTION AND WE INSERT UNIFORMLY AND INDEPENDENTLY GENERATED ITEMS INTO OUR HASH TABLE WITH  $m$  SLOTS. HOW MANY INSERTIONS ARE REQUIRED BEFORE WE EXPECT A COLLISION?
  - GUESS?
- CAN WE MODEL THIS USING INDICATOR RANDOM VARIABLES AND COMBINE THEM IN SOME WAY?
- $X_{ij} = I\{\text{ITEMS } i \text{ AND } j \text{ HASH TO THE SAME LOCATION}\}$
- LET  $X$  BE THE SUM OVER  $i$  AND  $j$ . WHAT DOES  $X$  REPRESENT?

$X$  is the number of pairs of elements with the same hash

$$X = \sum_{i=1}^k \sum_{j=i+1}^k X_{ij}$$

How can we use this random variable to solve our problem?

# EXPECTING A COLLISION

- $X_{ij} = I\{\text{ITEMS } i \text{ AND } j \text{ HASH TO THE SAME LOCATION}\}$
- $X = \{\text{THE NUMBER OF PAIRS OF ELEMENTS WITH THE SAME HASH}\}$

When this exceeds 1, the expected number of pairs is greater than 1.

$$k = \sqrt{2m} + 1$$

$$\begin{aligned} E[X] &= E\left[\sum_{i=1}^k \sum_{j=i+1}^k X_{ij}\right] \\ &= \sum_{i=1}^k \sum_{j=i+1}^k E[X_{ij}] \quad \bullet \quad \bullet \\ &= \binom{k}{2} \cdot \frac{1}{m} \\ &= \frac{k(k-1)}{2m} \end{aligned}$$

What does  $E[X_{ij}]$  equal?



# HOW MANY ELEMENTS WITH A GIVEN HASH?

- RECALL THE CHAINING PARADIGM USES A LINKED LIST AT EACH ELEMENT OF THE HASH TABLE. COLLISIONS ARE RESOLVED BY APPENDING THE ELEMENT TO THE END OF THE LINKED LIST.
- WE WOULD LIKE TO HAVE VERY SHORT LINKED LISTS FOR CHAINING TO WORK WELL
- SUPPOSE HAVE A UNIFORM HASHING FUNCTION FOR A HASH TABLE WITH  $m$  SLOTS. WE USE CHAINING TO RESOLVE COLLISIONS. IF  $n$  ELEMENTS ARE GENERATED INDEPENDENTLY AND UNIFORMLY AT RANDOM AND INSERTED INTO OUR HASH TABLE, WHAT IS THE EXPECTED LENGTH OF THE LINKED LIST AT TABLE ENTRY 3?
  - WHAT ARE OUR INDICATOR VARIABLES? WHAT ARE WE LOOKING TO SOLVE?
  - $X_j = I\{\text{THE } j\text{-TH ELEMENT HASHES TO ENTRY 3}\}$
  - NOW WE CAN COMBINE OUR INDICATOR VARIABLES USING LINEARITY OF EXPECTATION
  - $E\left[\sum_{j=1}^n X_j\right] = \sum_{j=1}^n E[X_j] = \frac{n}{m}$



# CHAIN HASHING AND PROBABILISTIC ANALYSIS

- SAY WE HAVE A UNIFORM HASHING FUNCTION AND WE INSERT UNIFORMLY AND INDEPENDENTLY GENERATED ITEMS INTO OUR HASH TABLE WITH  $m$  SLOTS. WE RESOLVE COLLISIONS WITH CHAINING.
- HOW MANY ELEMENTS DO WE NEED TO HASH, IN EXPECTATION, SO THAT THERE IS AT LEAST ONE ELEMENT IN HASH POSITION 3?
- WE KNOW THAT, IN EXPECTATION,  $\frac{n}{m}$  ELEMENTS ARE HASHED TO POSITION 3 WHEN WE INSERT  $n$  ELEMENTS. CAN WE USE THIS RESULT TO SOLVE OUR DESIRED PROBLEM?
- $\frac{n}{m} \geq 1 \Rightarrow n \geq m$
- **FOLLOW-UP QUESTION:** HOW MANY ELEMENTS DO YOU NEED TO INSERT, IN EXPECTATION, SO THAT THERE IS AT LEAST ONE ELEMENT AT EVERY HASH POSITION?



# FULL HASH TABLES USING CHAINING

- **MAIN QUESTION:** IF WE USE CHAIN HASHING, HOW MANY ELEMENTS DO YOU NEED TO INSERT, IN EXPECTATION, SO THAT EVERY HASH TABLE LOCATION IS USED?
  - HINTS
    - WHEN SELECTING A RANDOM VARIABLE(S), THINK: HOW CAN I **COMBINE** THEM? WHERE IS THERE ANY **RANDOMNESS**?
    - WHAT'S THE PROBABILITY OF HASHING THE FIRST INPUT TO A "NEW" LOCATION?
    - ONCE THE FIRST DISTINCT LOCATION HAS BEEN USED, WHAT'S THE PROBABILITY OF HASHING TO A NEW LOCATION ON THE FIRST SUBSEQUENT PULL?
- $X_i = I\{\text{THE } i\text{-TH DISTINCT HASH TABLE INDEX IS FILLED ON A THROW (GIVEN THAT } i - 1 \text{ ENTRIES WERE FILLED PRIOR)}\}$
- $E[X_i] = \frac{m-i+1}{m}$
- WHAT DOES THIS MEAN ABOUT THE EXPECTED NUMBER OF INSERTS NEEDED TO USE THE  $i$ -TH DISTINCT HASH ENTRY AFTER HAVING USED  $i - 1$  ENTRIES PRIOR?



# FULL HASH TABLES USING CHAINING

- $X_{ik} = I\{\text{A NEW HASH TABLE ENTRY IS USED ON THE } k\text{-TH INPUT AFTER HAVING USED } i-1 \text{ DISTINCT ENTRIES}\}$
- $E[X_{ik}] = \frac{m-i+1}{m}$
- LET  $Y_i$  BE THE NUMBER OF INSERTS REQUIRED TO USE THE  $i$ -TH DISTINCT ENTRY (AFTER HAVING USED  $i - 1$  DISTINCT ENTRIES PRIOR). HOW DOES  $Y_i$  RELATE TO  $X_{ik}$ ?
  - $E[Y_i] = \frac{1}{E[X_{ik}]}$
- TO FIND THE TOTAL INSERTIONS NEEDED TO FILL ALL  $m$  DISTINCT HASHES, WE CAN SUM OVER THE  $Y_i$ .

$$E \left[ \sum_{i=1}^m Y_i \right] =$$

$$\sum_{i=1}^m E[Y_i] =$$

$$\sum_{i=1}^m \frac{m}{m-i+1} =$$

$$\Theta(m \cdot \log(m))$$