CSE 143

Highlights of Tables and Hashing

Tables: Ch. 11, pp.515-522 Hashing: Ch. 12 pp.598-604

12/8/00 Y-1











- friend could be instantly located in the array by their social security number.
- •What's wrong with the above scheme?

12/8/00 Y-7



Hashing and Tables

- Hashing gives us another implementation of Table ADT
- •Hash the key; this gives an index; use it to find the value stored in the table
- If this scheme worked, it would be O(1)
- Great improvement over Log N.
- Main problems
- Finding a good hash function
- Collisions
- ·Wasted space in the table

12/8/00 Y-9

An Apparent Sidetrack										
Problem to solve: Given a list of n integers, determine if there is a pair of duplicate values										
37591	31576	64085	42782	25475	70900	79953	76186			
67887	84848	81309	30822	77867	45852	65289	8322			
79367	40520	58053	16030	34723	22116	41073	60522			
34399	31616	85965	82102	73707	38316	153	11282			
7623	61416	10741	46686	73123	69780	65105	21866			
75567	5760	66525	80214	63835	48652	49593	42066			
20055	16248	12213	35758	12147	13828	7729	2266			
13263	57984	73181	34246	51755	58053	31817	52754			
23863	9160	56677	62462	65715	68404	48097	66762			
37519	52480	28045	68294	71131	6252	81689	51570			
72119	71944	9797	77822	56563	67348	51553	86986			
88303	10656	6925	89654	63099	25036	84393	47426			
						12/8/	⁰⁰ Y-10			



12/8/00 Y-11

Element Uniqueness (2)								
Step 1: Assign to buckets, based on value mod								
375 <mark>91</mark>	315 <mark>76</mark>	64085	427 <mark>82</mark>	25475	709 <mark>00</mark>	799 <mark>53</mark>	761 <mark>86</mark>	
678 <mark>87</mark>	848 <mark>48</mark>	813 <mark>09</mark>	308 <mark>22</mark>	778 <mark>67</mark>	458 <mark>52</mark>	652 <mark>89</mark>	83 <mark>22</mark>	
793 <mark>67</mark>	405 <mark>20</mark>	580 <mark>53</mark>	160 <mark>30</mark>	347 <mark>23</mark>	221 <mark>16</mark>	410 <mark>73</mark>	605 <mark>22</mark>	
343 <mark>99</mark>	316 <mark>16</mark>	859 <mark>65</mark>	821 <mark>02</mark>	737 <mark>07</mark>	383 <mark>16</mark>	153	112 <mark>82</mark>	
76 <mark>23</mark>	614 <mark>16</mark>	107 <mark>41</mark>	466 <mark>86</mark>	731 <mark>23</mark>	697 <mark>80</mark>	651 <mark>05</mark>	218 <mark>66</mark>	
755 <mark>67</mark>	57 <mark>60</mark>	665 <mark>25</mark>	802 <mark>14</mark>	638 <mark>35</mark>	486 <mark>52</mark>	495 <mark>93</mark>	420 <mark>66</mark>	
200 <mark>55</mark>	162 <mark>48</mark>	122 <mark>13</mark>	357 <mark>58</mark>	121 <mark>47</mark>	138 <mark>28</mark>	77 <mark>29</mark>	22 <mark>66</mark>	
132 <mark>63</mark>	579 <mark>84</mark>	731 <mark>81</mark>	342 <mark>46</mark>	517 <mark>55</mark>	580 <mark>53</mark>	318 <mark>17</mark>	527 <mark>54</mark>	
238 <mark>63</mark>	91 <mark>60</mark>	566 <mark>77</mark>	624 <mark>62</mark>	657 <mark>15</mark>	684 <mark>04</mark>	480 <mark>97</mark>	667 <mark>62</mark>	
375 <mark>19</mark>	524 <mark>80</mark>	280 <mark>45</mark>	682 <mark>94</mark>	711 <mark>31</mark>	62 <mark>52</mark>	816 <mark>89</mark>	515 <mark>70</mark>	
721 <mark>19</mark>	719 <mark>44</mark>	97 <mark>97</mark>	778 <mark>22</mark>	565 <mark>63</mark>	673 <mark>48</mark>	515 <mark>53</mark>	869 <mark>86</mark>	
883 <mark>03</mark>	106 <mark>56</mark>	69 <mark>25</mark>	896 <mark>54</mark>	630 <mark>99</mark>	250 <mark>36</mark>	843 <mark>93</mark>	474 <mark>26</mark>	
						12/8/0)0 V 12	

Ele	men	t Ur	niqu	ene	SS	(3)	
Step 2	2: Look	inside	each	bucke	t for du	uplicat	es
375 <mark>9</mark>	1 315 <mark>76</mark>	640 <mark>85</mark>	427 <mark>82</mark>	254 <mark>75</mark>	709 <mark>00</mark>	799 <mark>53</mark>	761 <mark>86</mark>
678 <mark>8</mark>	7 848 <mark>48</mark>	813 <mark>09</mark>	308 <mark>22</mark>	778 <mark>67</mark>	458 <mark>52</mark>	652 <mark>89</mark>	83 <mark>22</mark>
793 <mark>6</mark>	7 405 <mark>20</mark>	580 <mark>53</mark>	160 <mark>30</mark>	347 <mark>23</mark>	221 <mark>16</mark>	410 <mark>73</mark>	605 <mark>22</mark>
343 <mark>9</mark>	9 316 <mark>16</mark>	859 <mark>65</mark>	821 <mark>02</mark>	737 <mark>07</mark>	383 <mark>16</mark>	153	112 <mark>82</mark>
76 <mark>2</mark>	3 614 <mark>16</mark>	107 <mark>41</mark>	466 <mark>86</mark>	731 <mark>23</mark>	697 <mark>80</mark>	651 <mark>05</mark>	218 <mark>66</mark>
755 <mark>6</mark>	7 57 <mark>60</mark>	665 <mark>25</mark>	802 <mark>14</mark>	638 <mark>35</mark>	486 <mark>52</mark>	495 <mark>93</mark>	420 <mark>66</mark>
200 <mark>5</mark>	5 162 <mark>48</mark>	122 <mark>13</mark>	357 <mark>58</mark>	121 <mark>47</mark>	138 <mark>28</mark>	77 <mark>29</mark>	22 <mark>66</mark>
132 <mark>6</mark>	3 579 <mark>84</mark>	731 <mark>81</mark>	342 <mark>46</mark>	517 <mark>55</mark>	580 <mark>53</mark>	318 <mark>17</mark>	527 <mark>54</mark>
238 <mark>6</mark>	<mark>3 9160</mark>	566 <mark>77</mark>	624 <mark>62</mark>	657 <mark>15</mark>	684 <mark>04</mark>	480 <mark>97</mark>	667 <mark>62</mark>
375 <mark>1</mark>	9 524 <mark>80</mark>	280 <mark>45</mark>	682 <mark>94</mark>	711 <mark>31</mark>	62 <mark>52</mark>	816 <mark>89</mark>	515 <mark>70</mark>
721 <mark>1</mark>	9 719 <mark>44</mark>	97 <mark>97</mark>	778 <mark>22</mark>	565 <mark>63</mark>	673 <mark>48</mark>	515 <mark>53</mark>	869 <mark>86</mark>
883 <mark>0</mark>	3 106 <mark>56</mark>	69 <mark>25</mark>	896 <mark>54</mark>	630 <mark>99</mark>	250 <mark>36</mark>	843 <mark>93</mark>	474 <mark>26</mark>
						12/8/0	^{IO} Y-13



Collisions

 Collisions occur when multiple items are mapped to same cell
 h(idNumber) = idNumber % B

h(678921) = 21h(354521) = 21

Issues

- Relative size of table to number of data items
 Choice of hash function
- •With a bad choice of hash function we can have lots of collisions
- Even with a good choice of hash functions there may be some collisions



Analysis of hash table ops

- Insert is easy to analyze:
 - It is just the cost of calculating the hash value O(1), plus the cost of inserting into the front of a linked list O(1)
- *Retrieve* and *Delete* are harder. To do the analysis, we need to know:
- The number of elements in the table (N)
- The number of buckets (B)
- The quality of the hash function

^{12/8/00} Y-17



• Note that this means growing (rehashing) the hash table as more items are inserted.

^{12/8/00} Y-18





Hashing and Files

•We've spoken of the hashed data as being stored in an array (in memory)

Hashing is also very appropriate for disk files

 Efficient look-up techniques for disk data are essential

Disks are thousands of times slower than memoryEven a LogN look-up algorithm is too slow for a

database application! Many structures we have studied (linked lists, trees, etc.)

do not scale well to large disk files

12/8/00 Y-21

Drawbacks to Hashing

- Finding a good hash function
 Small risk of bad behavior
- Dealing with collisions
 Simplest method to use linked list for buckets
- Wasted space in the array
 Not a big deal if memory is cheap
- Doesn't support ordering queries (such as we would want for a real dictionary)

12/8/00 Y-22

Summary

- Hash tables are specialized for dictionary operations: Insert, Delete, Lookup
- Principle: Turn the key field of the record into a number, which we use as an index for locating the item in an array.
- •O(1) in the ideal case; less in practice
- Problems: collisions, wasted space
- Implementations: open hashing, closed hashing, dynamic hashing
- •Highly suitable for database files, too

12/8/00 Y-23