#### CSC148 winter 2018

tracing, hash table week 11

```
Danny Heap
heap@cs.toronto.edu / BA4270 (behind elevators)
http://www.teach.cs.toronto.edu/~csc148h/winter/
416-978-5899
```

March 23, 2018





#### Outline

debugger delving

hash tables

### call stack

### recursion redux



# resolving hierarchy

#### missed reference



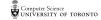
## why hash

lists are contiguous (adjacent) sequences of references to objects, so access to a list position is fast (just arithmetic)

```
[ , , , , , ..... , , , , , ]
```

what if we could convert — hash — other data to a suitable integer for a list index, we'd want:

- ▶ fast
- ▶ deterministic: the same (or equivalent values) gets hashed to the same integer each time.
- well-distributed: We'd like a typical set of values to get hashed pretty uniformly over the available list positions.





### you can't hash everything!

```
>>> list1 = [0]

>>> id(list1)

3069263116

>>> list2 = [0, 1]

>>> id(list2)

3069528300

>>> list1.append(1)

>>> id(list1)

3069263116
```

oops!



# hash to hash table (dictionary)...

Once you have hashed an object to a number, you can easily use part of that number as an index into a list to store the object, or something related to that object. If the list is of length n, you might store information about object o at index hash(o) % n.

#### collisions

even a well-distributed hash function will have a surprising number of collisions...

how many people do you need to poll before you find two with the same birthday (out of 366 possibilities, including leap-year)?

the mathematics is a bit counter-intuitive... the probability of a non-collision for 23 birthdays is:

$$p = \frac{366}{366} \times \frac{365}{366} \times \cdots \times \frac{344}{366} \approx 0.493$$





## chaining or probing

a couple of tactics for dealing with two different keys ending up at the same index

chaining: keep a small (one hopes) list at that index (sometimes called bucket)

probing: explore, in a systematic way, until the next open index

either tactic has costs, so keep collisions to a minimum by keeping the list partly empty





Python dictionaries are implemented<sup>1</sup> using hash tables and probing. The cost of collisions is kept small by enlarging the underlying list when necessary, and the cost of enlarging is amortized over many dictionary accesses.

The result is that access to a dictionary element is  $\Theta(1)$ , essentially the time it takes to access a list element.

One downside is that extra work is required to order the keys or values of a dictionary. What is their "natural" order?



<sup>&</sup>lt;sup>1</sup>Google "Tim Peters" but beware of obnoxious culture → ⟨ ႃ → ⟨ ႃ → ト | ■