

Algorithms and Data Structures for Data Science

Stacks and Queues (and Sets)

CS 277

February 13, 2023

Brad Solomon



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science

Learning Objectives

Explore tradeoffs in data structures

Introduce the stack and queue

Introduce sets

Tradeoffs in data structures

As we progress in the class, we will see that $O(n)$ isn't very good.

Take searching for a specific list value:

2	7	5	9	7	14	1	0	8	3
---	---	---	---	---	----	---	---	---	---

0	1	2	3	5	7	7	8	9	14
---	---	---	---	---	---	---	---	---	----

Tradeoffs in data structures

Getting the size of a linked list has a Big O of:



Tradeoffs in data structures: Bag of Words

Genome assembly databases are growing rapidly. The sequence content in each new assembly can be largely redundant with previous ones, but this is neither conceptually nor algorithmically easy to measure. We propose new methods and a new tool called DandD that addresses the question of how much new sequence is gained when a sequence collection grows. DandD can describe how much human structural variation is being discovered in each new human genome assembly and when discoveries will level off in the future. DandD uses a measure called δ ("delta"), developed initially for data compression. Computing δ directly requires counting k-mers, but DandD can rapidly estimate it using genomic sketches. We also propose δ as an alternative to k-mer-specific cardinalities when computing the Jaccard coefficient, avoiding the pitfalls of a poor choice of k. We demonstrate the utility of DandD's functions for estimating δ , characterizing the rate of pangenome growth, and computing allpairs similarities using k-independent Jaccard. DandD is open source software available at: <https://github.com/jessicabonnie/dandd>.



Tradeoffs in data structures

I want a list that can add and remove in $O(1)$

I am willing to make random access impossible to do so

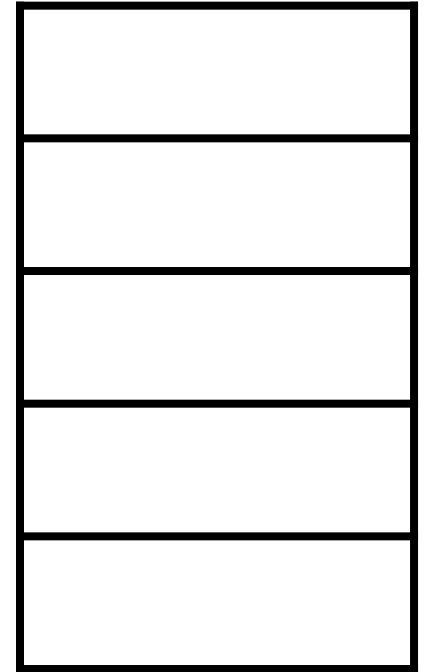
Stack Data Structure

A **stack** stores an ordered collection of objects (like a list)

However you can only do two operations:

Push: Put an item on top of the stack

Pop: Remove the top item of the stack (and return it)



```
push (3) ; push (5) ; pop () ; push (2)
```

Stack Data Structure

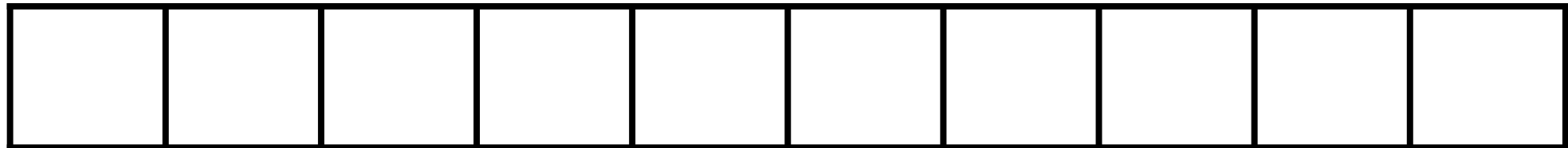
The stack is a **last in — first out** data structure (LIFO)



```
1 def reverse(inList):  
2     s = stack()  
3     for v in inList:  
4         s.push(v)  
5  
6     out = []  
7     while not s.empty():  
8         out.append(s.pop())  
9     return out
```

Not a Python built-in!

`reverse([3, 4, 5, 6, 7, 8])`

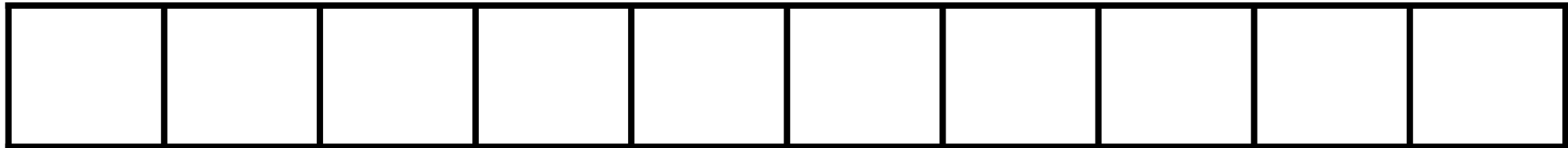


Stack Data Structure

The Python list has all the necessary stack operations!

```
1 def push(self, val):
2     self.push_count += 1
3     self.data.append(val)
4
5 def pop(self):
6     self.pop_count += 1
7     if self.__len__() > 0:
8         return self.data.pop()
9
```

Stack s
s.push(3)
s.push(8)
s.push(4)
s.pop()
s.push(7)
s.pop()
s.pop()
s.push(2)
s.push(1)
s.push(3)
s.push(5)
s.pop()
s.push(9)



Stack Data Structure

push (X)

The stack is also easily implemented as a linked list



Stack Data Structure



The stack is also easily implemented as a linked list



Programming w/ the Call Stack

```
1 def Happy():
2     return "Happy"
3
4 def Little():
5     return "Little"
6
7 def Trees():
8     return "Trees"
9
10 def LittleTrees():
11     return Little() + Trees()
12
13 def BobRoss():
14     return Happy() + LittleTrees()
15
16 if __name__ == '__main__':
17     print(BobRoss())
18
19
20
21
22
23
```

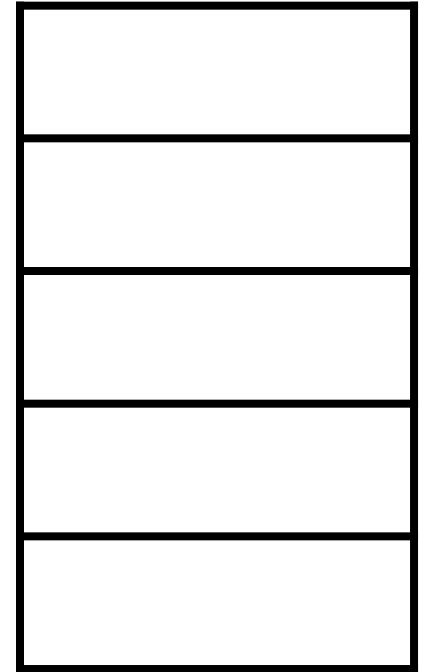
Queue Data Structure

A **queue** stores an ordered collection of objects (like a list)

However you can only do two operations:

Enqueue: Put an item at the back of the queue

Dequeue: Remove and return the front item of the queue



```
enqueue (3) ; enqueue (5) ; dequeue () ; enqueue (2)
```

Queue Data Structure

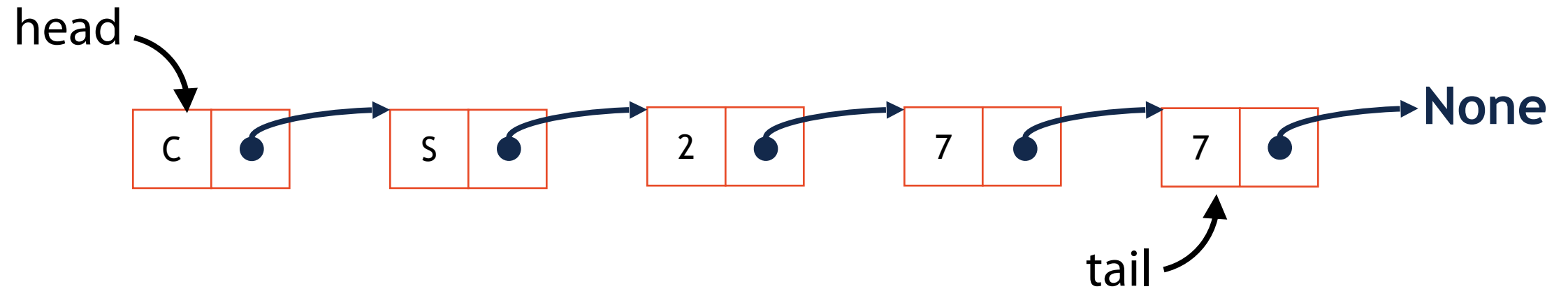
The queue is a **first in — first out** data structure (FIFO)

What data structure excels at removing from the front?

Can we make that same data structure good at inserting at the end?

Queue Data Structure

The queue is a **first in — first out** data structure (FIFO)



Queue Data Structure

The queue is a **first in — first out** data structure (FIFO)

```
Queue q  
q.enqueue(3)  
q.enqueue(8)  
q.enqueue(4)  
q.dequeue()  
q.enqueue(7)  
q.dequeue()  
q.dequeue()  
q.enqueue(2)  
q.enqueue(1)  
q.enqueue(3)  
q.enqueue(5)  
q.dequeue()  
q.enqueue(9)
```


Queue Data Structure

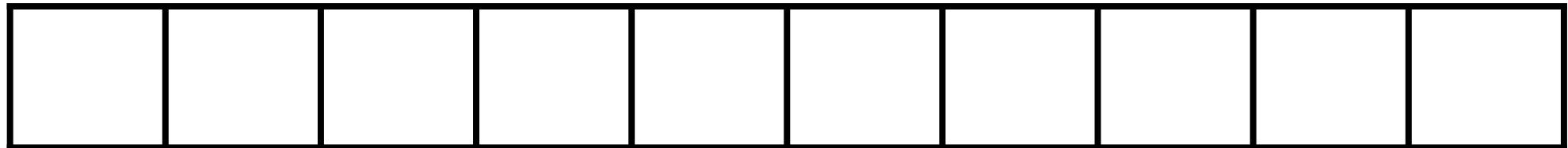
An array can implement a queue as well!

We just need to track three numbers:

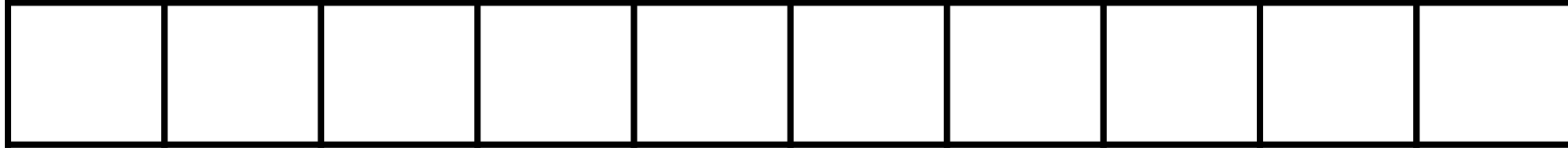
Front

Capacity

Size



Queue Data Structure



Front:

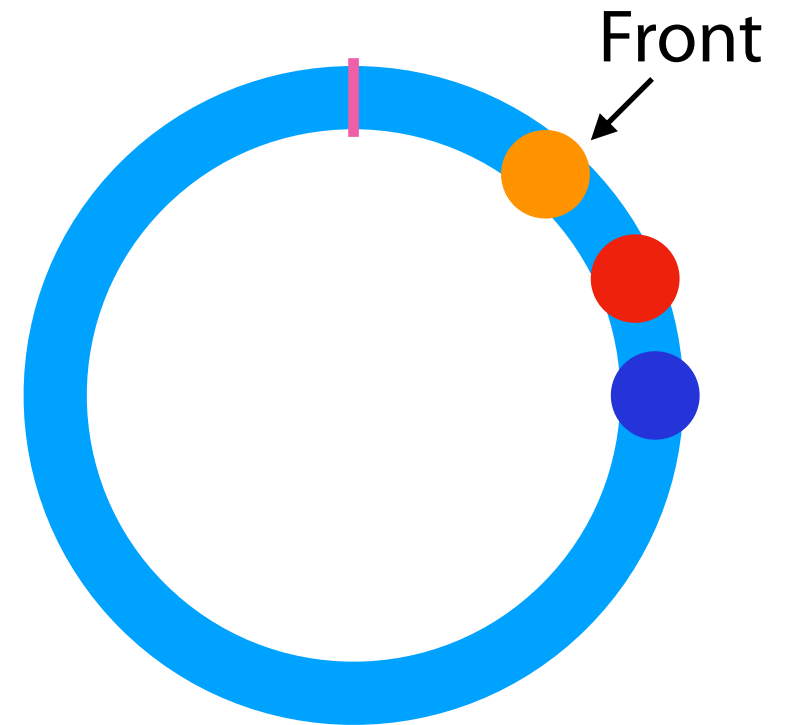
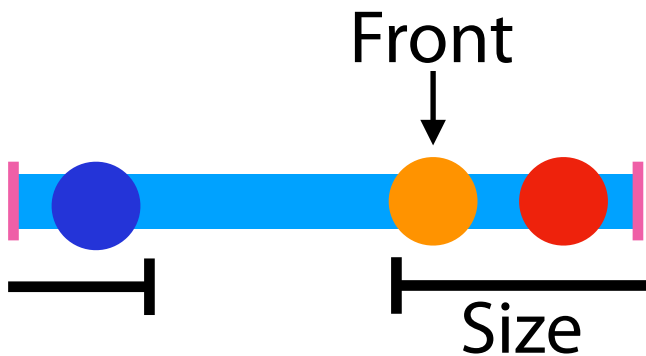
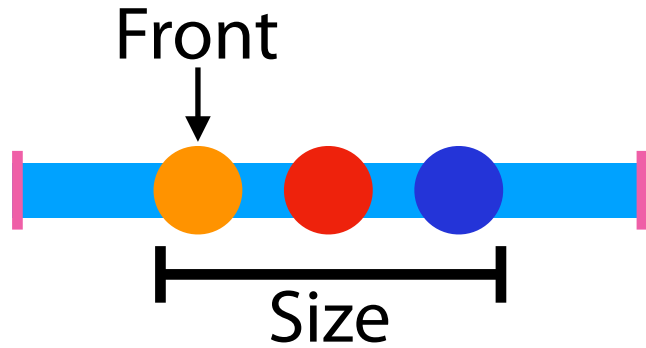
Capacity:

Size:

Queue q
q.enqueue(3)
q.enqueue(8)
q.enqueue(4)
q.dequeue()
q.enqueue(7)
q.dequeue()
q.dequeue()
q.enqueue(2)
q.enqueue(1)
q.enqueue(3)
q.enqueue(5)
q.dequeue()
q.enqueue(9)

Queue Data Structure

The array implementation treats the allocated memory as a circle



Queue Array Resizing

6	7	0	1	2	3	4	5
---	---	---	---	---	---	---	---

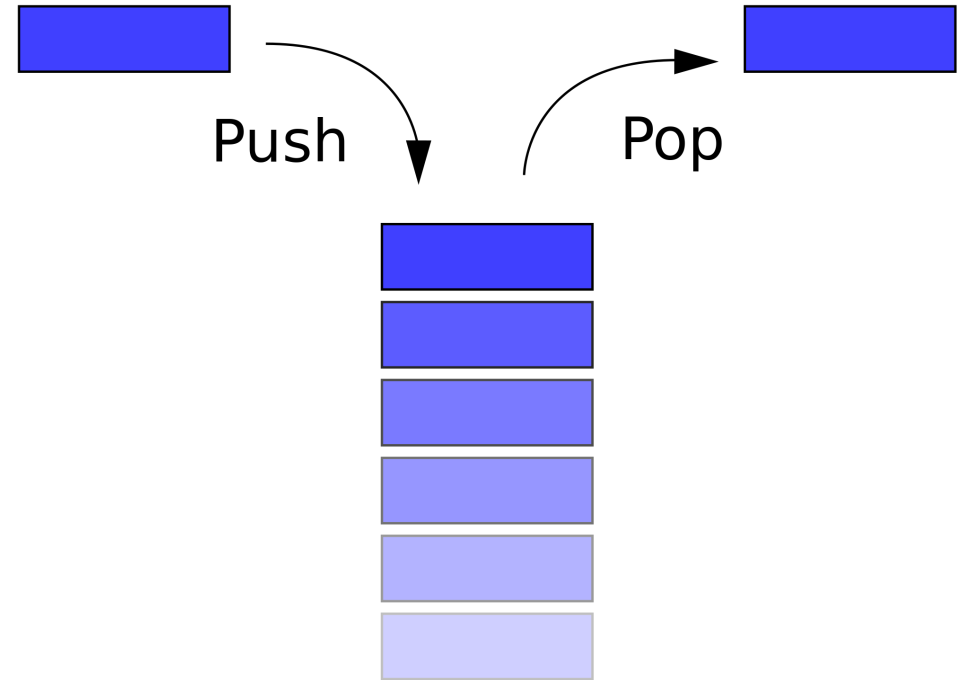
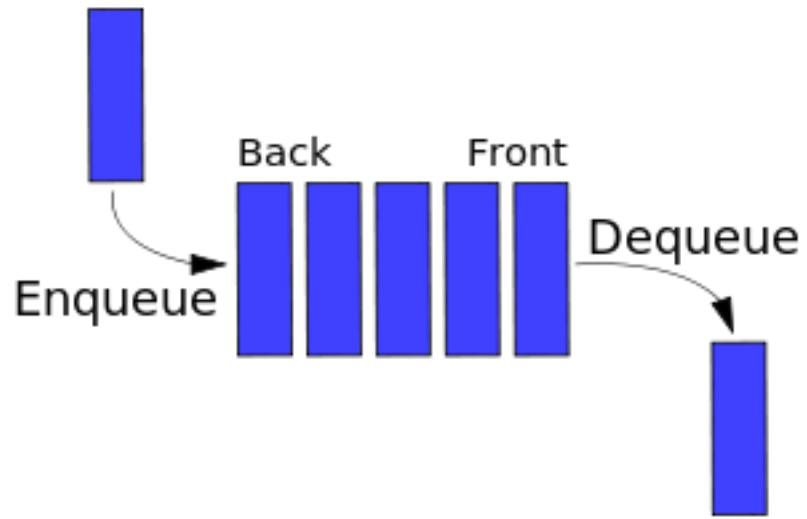
Queue q
...
Fill queue
...
`q.enqueue(8)`

Front: 2














Capacity: 8

Size: 8

Stacks and Queues



Engineering vs Theory Efficiency

	Time x1 billion	Like
L1 cache reference	0.5 seconds	Heartbeat 
Branch mispredict	5 seconds	Yawn 
L2 cache reference	7 seconds	Long yawn   
Mutex lock/unlock	25 seconds	Make coffee 
Main memory reference	100 seconds	Brush teeth
Compress 1K bytes	50 minutes	TV show 
Send 2K bytes over 1 Gbps network	5.5 hours	(Brief) Night's sleep 
SSD random read	1.7 days	Weekend
Read 1 MB sequentially from memory	2.9 days	Long weekend
Read 1 MB sequentially from SSD	11.6 days	2 weeks for delivery 
Disk seek	16.5 weeks	Semester
Read 1 MB sequentially from disk	7.8 months	Human gestation 
Above two together	1 year	 
Send packet CA->Netherlands->CA	4.8 years	Ph.D. 

(Care of <https://gist.github.com/hellerbarde/2843375>)

Engineering vs Theory Efficiency

```
1 class queue:
2     def __init__(self):
3         self.data = []
4
5
6     def enqueue(self, val):
7         self.data.append(val)
8
9
10    def dequeue(self):
11        if len(self.data) > 0:
12            return self.data.pop(0)
13
14
```

```
1 def enqueue(self, val):
2
3     i = (self.front + self.size) % self.capacity
4
5     self.data[i]=val
6     self.size += 1
7
8
9     if self.size == self.capacity:
10        temp = [None] * self.capacity * 2
11
12        for i in range(self.size):
13            pos = (self.front + i ) % self.capacity
14            temp[i]=self.data[pos]
15
16        self.front = 0
17        self.data = temp
18        self.capacity = len(temp)
19
20
21
22
23
```

Tradeoffs in data structures

I want a data structure that can add, remove, and find items in $O(1)$ *

I am willing to remove the 'ordered' property of my collection to do this

I am willing to remove the ability to store duplicate elements to do this

Set “Data Structure”

A **set** stores an unordered collection of objects with no duplicates

Genome assembly databases are growing rapidly. The sequence content in each new assembly can be largely redundant with previous ones, but this is neither conceptually nor algorithmically easy to measure. We propose new methods and a new tool called DandD that addresses the question of how much new sequence is gained when a sequence collection grows. DandD can describe how much human structural variation is being discovered in each new human genome assembly and when discoveries will level off in the future. DandD uses a measure called δ (“delta”), developed initially for data compression. Computing δ directly requires counting k-mers, but DandD can rapidly estimate it using genomic sketches. We also propose δ as an alternative to k-mer-specific cardinalities when computing the Jaccard coefficient, avoiding the pitfalls of a poor choice of k. We demonstrate the utility of DandD’s functions for estimating δ , characterizing the rate of pangenome growth, and computing allpairs similarities using k-independent Jaccard. DandD is open source software available at: <https://github.com/jessicabonnie/dandd>.



Sets in Python: Constructor

The set constructor **set()** takes a list or tuple as input

```
1 s1 = set([1,2,3,4])
2
3 s2 = set((3,4,5,6))
4
5
6
7
8
9
10
11
12
13
14
15
16
```

Sets in Python: Add

Add(x) adds object x to the set; it does nothing if x is already present

```
1 mySet = set()
2
3
4 mySet.add(1)
5
6
7 mySet.add(1)
8
9
10 mySet.add(3)
11
12
13
14
15
16
```

Sets in Python: Remove

Remove(x) removes the object *x* from the set if it is present.

If *x* does not exist, it will crash.

```
1 mySet = set([1,2,3,4,5])
2
3 mySet.remove(3)
4
5 print(mySet)
6
7 mySet.remove(10)
8
9
10
11
12
13
14
15
16
```

Sets in Python: Data Access

Sets have no indices (no order of objects). We can only access by:

1. Looping through a set for each element
2. Looking up a specific element in our set

```
1 mySet = set([1,2,3,4,5])
2 for obj in mySet:
3     print(obj)
4
5 print(10 in mySet)
6
7
8
9
10
11
12
13
```

Sets in Python: Implementation



A Python set is implemented using a **hash table**.

Claim: A hash table has direct lookup, add, and remove in $O(1)$ *

Set Operations

$$A = \{1, 2, 3, 4\} \quad B = \{3, 4, 5, 6, 7\}$$

Union $A \cup B$

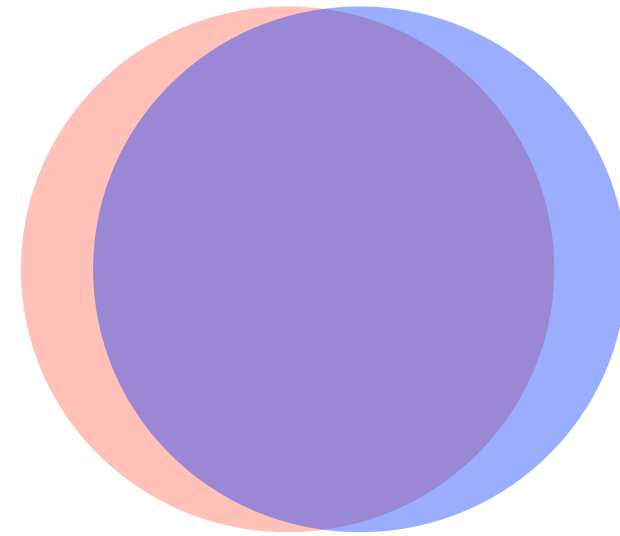
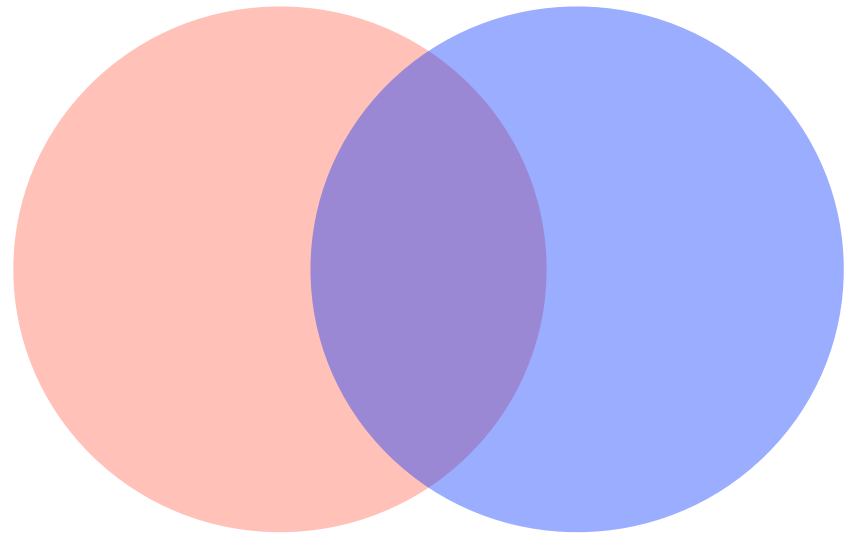
Intersection $A \cap B$

Difference A / B

Symmetric difference $A \triangle B$

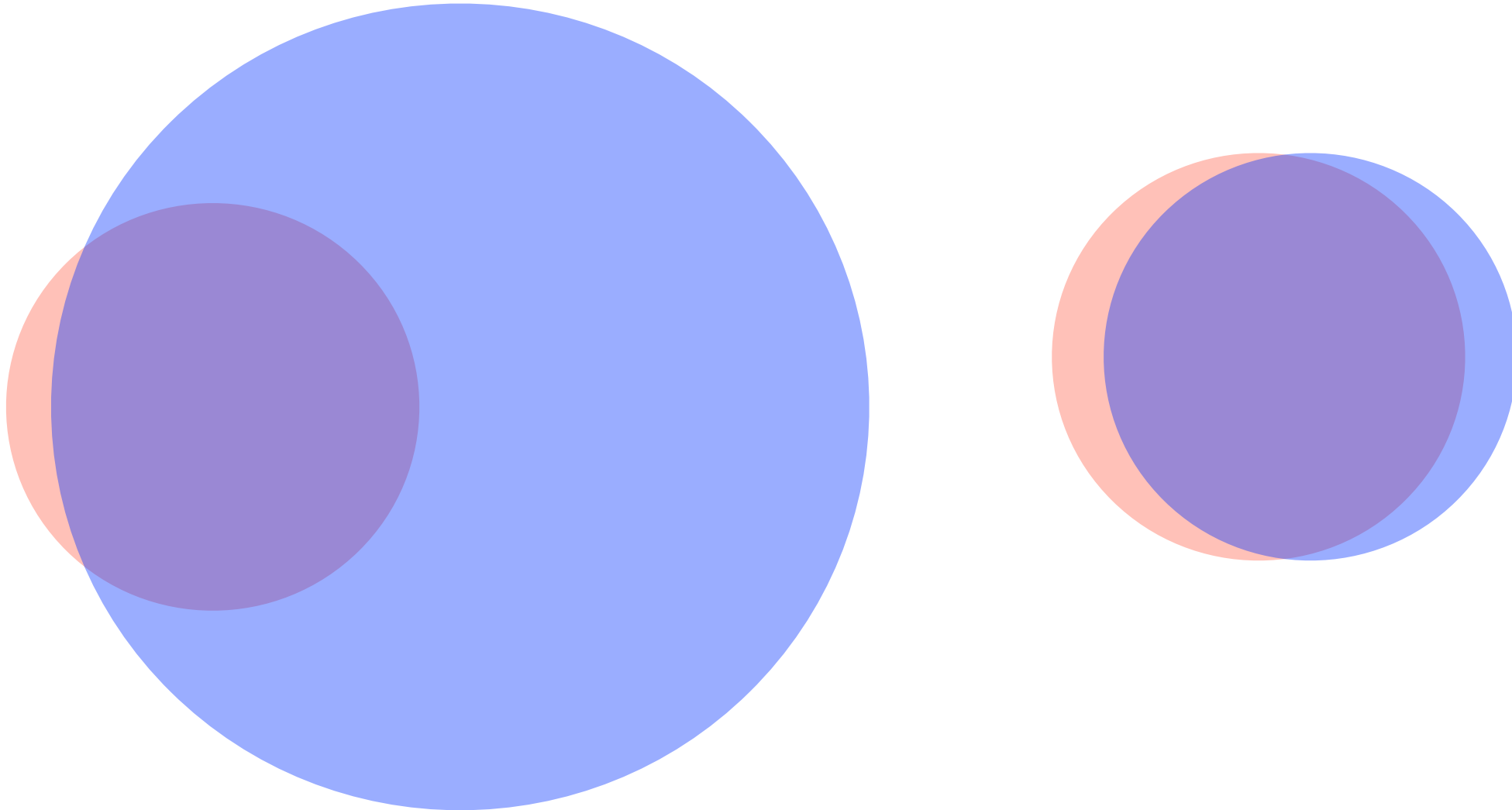
Set Similarity

How can we describe how *similar* two sets are?



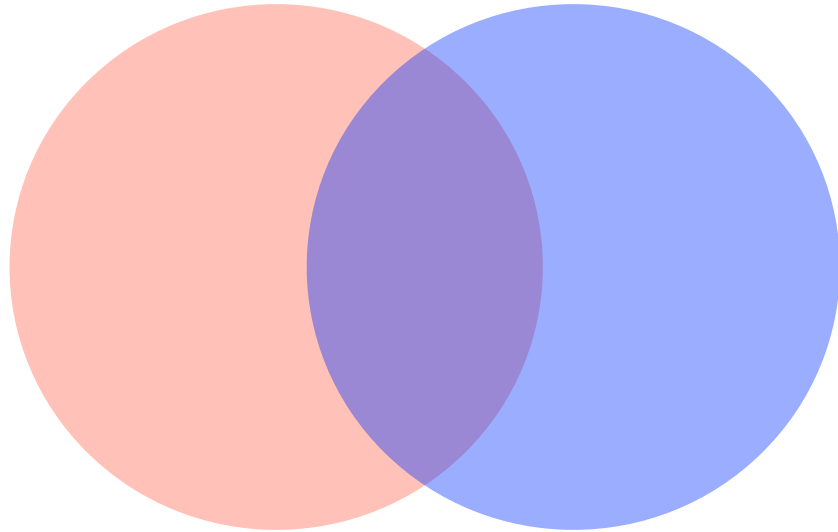
Set Similarity

How can we describe how *similar* two sets are?



Set Similarity

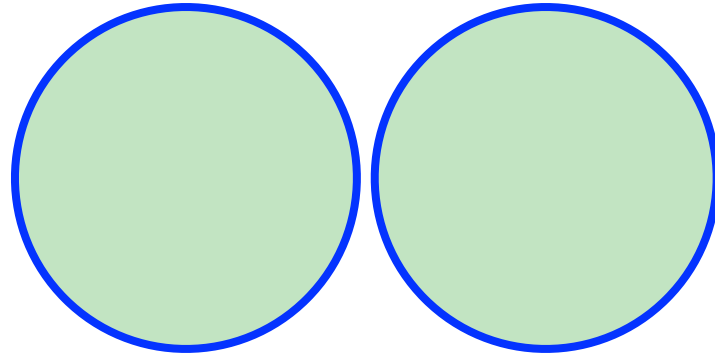
To measure **similarity** of A & B , we need both a measure of how similar the sets are but also the total size of both sets.



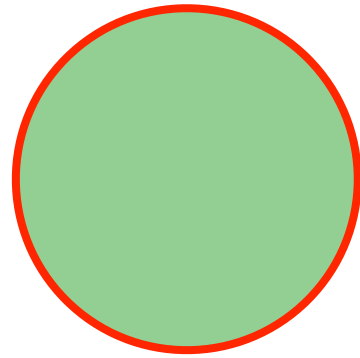
$$J = \frac{|A \cap B|}{|A \cup B|}$$

J is the **Jaccard coefficient**

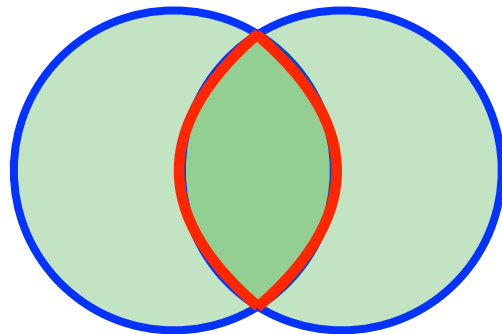
Set Similarity



$$\frac{|A \cap B|}{|A \cup B|} = 0$$



$$\frac{|A \cap B|}{|A \cup B|} = 1$$



$$0 < \frac{|A \cap B|}{|A \cup B|} < 1$$

Set Similarity

$$A = \{1, 2, 3, 4\} \quad B = \{3, 4, 5, 6, 7\}$$

$$J = \frac{|A \cap B|}{|A \cup B|} =$$

Set Similarity

$$A = \{1, 2, 3, 4\} \quad B = \{3, 4, 5, 6, 7\}$$

$$J = \frac{|A \cap B|}{|A \cup B|} = \frac{|\{3, 4\}|}{|\{1, 2, 3, 4, 5, 6, 7\}|} = \frac{2}{7}$$