

Reasoning About Correctness: Binary Search

```
// Pre: assumes data is ordered
// Post: returns index position (low) such that all
// elements indexed less than low are smaller than val, and
// all elements indexed greater than low are greater than
// or equal to val
```

```
int _binarySearch(TYPE * data, int size, TYPE val) {
    int low = 0;
    int high = size;
    int mid;
    //(A)
    while (low < high) {
        mid = (low + high) / 2;
        //mid less than val looking for
        if (LT(data[mid],val))
            low = mid + 1;
        else high = mid;
    }
    //(B)
    }
    //(C)
    return low;
}
```

Three important questions to ask about an algorithm...

- A. Is the result correct?
- B. Does the algorithm terminate?
- C. What is the execution time?

A. For `_BinarySearch`, is the resulting index correct?

To be correct, the resulting index returned from Binary Search must satisfy the following:

1. The result index should be ≥ 0 and \leq the number of elements in the array (ie. it can go anywhere inside the array or at the very end)
2. All values in positions $<$ the result index must be strictly $<$ value at index
3. All values in positions \geq result index must be \geq the value at index

To argue that it's correct, we will:

- **Establish assertions that correspond to 1-3** Before, Inside, After the loop (A,B, and C respectively in the code above)
- Argue paths from each assertion to the next [Note: An assertion in a loop is an invariant]

Assertion#1

The elements at index less than low are themselves strictly less than the argument (argument is the value we're looking for): For all $j < \text{low}$, $a[j] < \text{val}$ (for #2 above)

Assertion #2

The elements at index $\geq \text{high}$ are themselves $\geq \text{value}$: For all $j \geq \text{high}$, $a[j] \geq \text{val}$ (for #3 above)

Do Assertions Hold at A,B, and C?

1. Before the Loop (A)

Low is initialized to 0 and high is initialized to one larger than the very highest legal index ($n+1$)

Therefore, the sets described by the invariants are EMPTY so the invariants must hold!

2. In the loop (B) (they're called invariants in the loop!)

Assume invariants 1&2 hold at the beginning of an iteration and we execute the body of the loop.

If the $a[\text{mid}] < \text{val}$ (ie. condition is true), then all elements at positions **mid** and below must have values $< \text{val}$ so we can safely move **low** to **mid + 1** and still preserve invariant#1.

If the condition is false, then $a[\text{mid}] \geq \text{val}$. Likewise, since these are sorted, the values at positions higher than **mid** must also be $\geq \text{val}$, so we can assign **high** to **mid** and still preserve invariant #2

3. At the end of the Loop (C)

If the loop executes 0 times, the invariant was true before so must be true after the loop

If the loop executes some number of times, then the invariants must have been true at the last iteration of the loop so they must still be true when we're done since we have not done anything additional after the loop. Thus, invariants 1&2 are still preserved.

So, We've established our invariants and demonstrated that they hold from (A) to (B) [ie. before loop to inside loop], from (B) to (B) [inside the loop] and from (B) to (C) when loop is done.

Finally, let's show that the resulting integer index has the desired properties (e.g. it's ≥ 0 and $\leq \text{size}$). To do so, we combine our invariants with the observation that when the loop terminates, **high** \leq **low**, however, high can't possibly be less than low, so **high = low**.

Why can't high ever be less than low???

1. Remember that mid must be strictly $< \text{high}$ because of our integer arithmetic with truncation (but could be = low)

2. At worst, we have only two elements, low at 0 and high at 1. In this case, mid is 0, so either low gets set to mid + 1 (or high), or high gets set to mid (or low) so high= low.

Now, let's tie it all together: If Low == High and 1, 2, 3, are true, we have found the index where either the element is...or where it should be inserted to maintain the sorted order of the array, and our post condition is held.

B. To Demonstrate Termination, we need to find a quantity that...

1. Is an integer value
2. Changes each iteration through the loop
3. Is decreasing
4. Is strictly positive

In this algorithm, we can't use **high** or **low** because only one of them is changed in each iteration. However, we can consider **high-low**

- It changes each iteration
- It is integer since we're using integer arithmetic
- Is always decreasing because we compute mid as **low+high/2** with integer arithmetic so it must be strictly less than high and we set low either to mid+1 or high to mid so we move low up or high down so high-low MUST decrease.

We've found a quantity that meets all criteria so it MUST terminate.

C. Execution Time

We divide the array roughly in half at each step. We can do this at most $\log_2 n$ times. So, `_binarySearch` is $O(\log n)$ [Much, much better than $O(n)$ linear searching....what is $\log_2(1,000,000)$?]

Bug?

In 2006, Google found a bug in the `_binarySearch` algorithm. Can you find it? Can you fix our argument such that it uncovers this bug? Give up??? [<http://googleresearch.blogspot.com/2006/06/extra-extra-read-all-about-it-nearly.html>]

Efficiency Feature

Finally, we can add one feature to potentially speed up our algorithm. That is, if the value is found, we can simply return the value, terminating the algorithm. Can you make this change to the algorithm? Does it require any changes to the correctness argument?