# The Myers diff algorithm: part 2

If you enjoy this article, I have published a book explaining the internals of Git through implementation: **Building Git**.
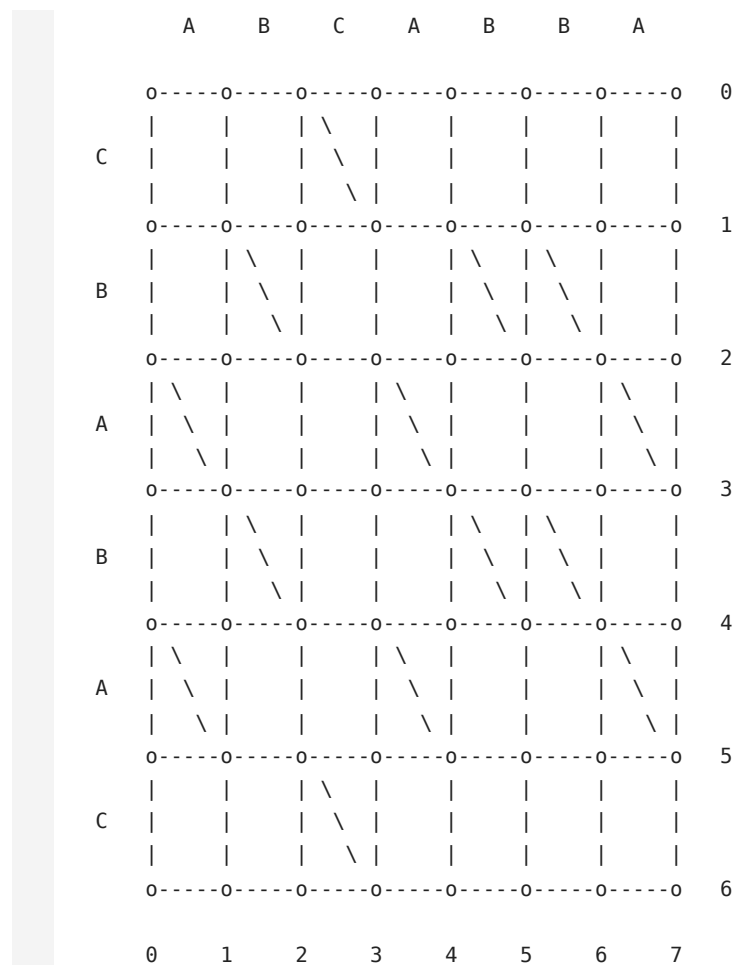
—

In part 1 of this series, we saw how the ==diff between two strings is modelled as a graph search problem==. We worked through the ==shortest edit script== between two strings:
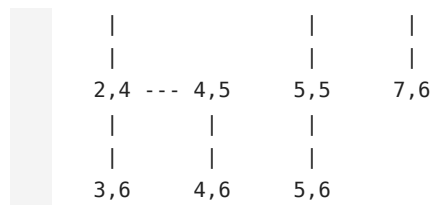
- *a* = ABCABBA
- *b* = CBABAC

We saw that the edit graph for these strings looks like this:

```
        A       B       C       A       B       B       A

    o-----o-----o-----o-----o-----o-----o-----o    0
    |     |     | \   |     |     |     |     |
  C |     |     |  \  |     |     |     |     |
    |     |     |   \ |     |     |     |     |
    o-----o-----o-----o-----o-----o-----o-----o    1
    |     | \   |     |     | \   | \   |     |
  B |     |  \  |     |     |  \  |  \  |     |
    |     |   \ |     |     |   \ |   \ |     |
    o-----o-----o-----o-----o-----o-----o-----o    2
    | \   |     |     | \   |     |     | \   |
  A |  \  |     |     |  \  |     |     |  \  |
    |   \ |     |     |   \ |     |     |   \ |
    o-----o-----o-----o-----o-----o-----o-----o    3
    |     | \   |     |     | \   | \   |     |
  B |     |  \  |     |     |  \  |  \  |     |
    |     |   \ |     |     |   \ |   \ |     |
    o-----o-----o-----o-----o-----o-----o-----o    4
    | \   |     |     | \   |     |     | \   |
  A |  \  |     |     |  \  |     |     |  \  |
    |   \ |     |     |   \ |     |     |   \ |
    o-----o-----o-----o-----o-----o-----o-----o    5
    |     |     | \   |     |     |     |     |
  C |     |     |  \  |     |     |     |     |
    |     |     |   \ |     |     |     |     |
    o-----o-----o-----o-----o-----o-----o-----o    6

    0     1     2     3     4     5     6     7
```
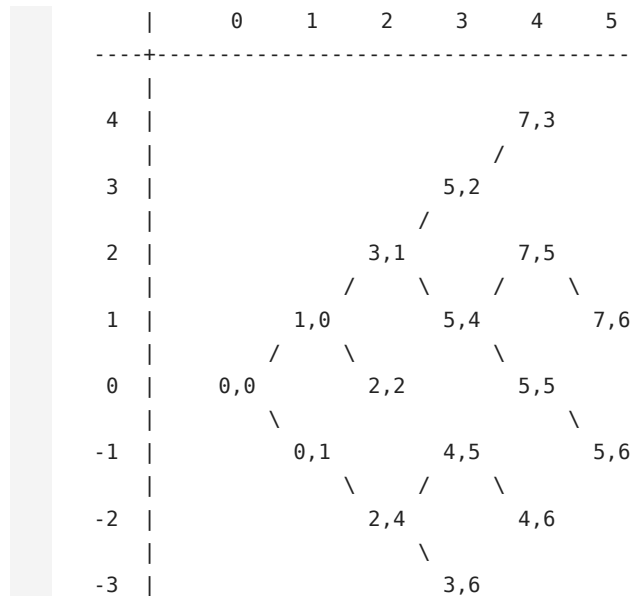
And, we recorded a trace through the graph to find the shortest path from (0,0) to the bottom-right corner (7,6).

```
0,0 --- 1,0 --- 3,1 --- 5,2 --- 7,3
 |       |       |
 |       |       |
0,1     2,2     5,4 --- 7,5
```

```
       |               |       |
       |               |       |
      2,4 --- 4,5     5,5     7,6
       |       |       |
       |       |       |
      3,6     4,6     5,6
```

Now, having seen how the graph search works, we're going to change the representation slightly to get us toward how the Myers algorithm actually works. Imagine that we take the above graph walk and render it rotated by 45 degrees.
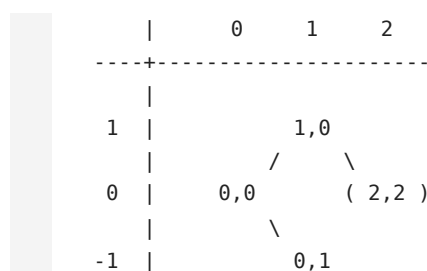
```
     |    0    1    2    3    4    5
  ---+----------------------------------
     |
  4  |                         7,3
     |                        /
  3  |                   5,2
     |                  /
  2  |              3,1       7,5
     |             /    \    /    \
  1  |         1,0       5,4       7,6
     |        /    \         \
  0  |    0,0       2,2       5,5
     |        \              \
 -1  |         0,1       4,5       5,6
     |             \    /    \
 -2  |              2,4       4,6
     |                  \
 -3  |                   3,6
```

The number along the horizontal axis, $d$, is the depth we've reached in the graph, i.e. how many moves we've made so far, remembering that diagonal moves are free. The number along the vertical axis we call $k$, and notice that for every move in the graph, $k = x - y$ for each move on that row.

Moving rightward increases $x$, and so increases $k$ by 1. Moving downward increases $y$ and so decreases $k$ by 1. Moving diagonally increases both $x$ and $y$, and so it keeps $k$ the same. So, each time we make a rightward or downward move followed by a chain of diagonals, we either increment or decrement $k$ by 1. What we are recording is the furthest through the edit graph we can reach for each value of $k$, at each step.

Now, here's how the algorithm proceeds. <mark>For each $d$ beginning with 0, we fill in each move for $k$ from $-d$ to $+d$ in steps of 2.</mark> Our aim at each $(d, k)$ position is to determine the best move we can make from the previous position. The <mark>best move is the one that gives us the highest $x$ value; maximising $x$ rather than $y$ means we prioritise deletion over insertion.</mark>
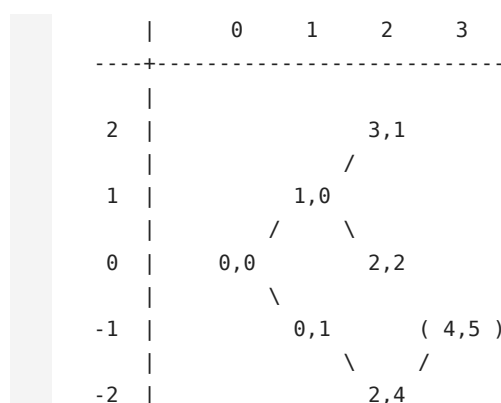
To discover the best move, we need to decide whether we should pick a downward move from $(d - 1, k + 1)$, or a rightward move from $(d - 1, k - 1)$. If $k$ is $-d$ then the move must be downward, likewise if $k$ is $+d$ then we must move rightward. For all other values of $k$, we pick the position with the highest $x$ from the two adjacent $k$ values in the previous column, and determine where that move leads us.

For **exam**ple, **cons**ider **t**he **mo**ve **a**t $(d, k) = (2,0)$. **We c**an **eit**her **mo**ve **right**ward **fr**om $(d, k) = (1,−1)$ **whe**re $(x, y) = (0,1)$, **or down**ward **fr**om $(d, k) = (1,1)$ **whe**re $(x, y) = (1,0)$.

```
      |    0    1    2
  ----+---------------------
      |
  1   |         1,0
      |        /    \
  0   |    0,0      ( 2,2 )
      |        \
 -1   |         0,1
```

(1,0) **h**as a **hig**her $x$ **val**ue **th**an (0,1), **s**o **we pi**ck a **mo**ve **down**ward **fr**om (1,0) **to** (1,1), **whi**ch **lea**ds **us to** (2,2) **diago**nally. **There**fore **we rec**ord $(x, y) = (2,2)$ **for** $(d, k) = (2,0)$. **Th**is **expl**ains **wh**y **we reco**rded **t**he **mo**ve **vi**a **th**is **pa**th **wh**en (2,0) **is al**so **reach**able **b**y **go**ing **right**ward **fr**om (0,1); **pick**ing **t**he **prev**ious **posi**tion **wi**th **t**he **high**est $x$ **val**ue **mea**ns **we t**ry **to maxi**mise **t**he **num**ber **o**f **dele**tions **we ma**ke **bef**ore **try**ing **inser**tions.

**In so**me **situa**tions, **t**he **t**wo **prev**ious **posi**tions **wi**ll **ha**ve **t**he **sa**me $x$ **val**ue. **F**or **exam**ple, **cons**ider **t**he **mo**ve **a**t $(d, k) = (3,−1)$, **whe**re **we c**an **mo**ve **down**ward **fr**om $(x, y) = (2,2)$ **or right**ward **fr**om $(x, y) = (2,4)$. **Mov**ing **right**ward **wi**ll **incr**ease $x$, **s**o **we mo**ve **fr**om (2,4) **t**o (3,4) **a**nd **th**en **diago**nally **t**o (4,5).

```
      |    0    1    2    3
  ----+----------------------------
      |
  2   |                  3,1
      |                 /
  1   |            1,0
      |           /    \
  0   |    0,0         2,2
      |        \
 -1   |            0,1      ( 4,5 )
      |               \    /
 -2   |                  2,4
```

**The**re **a**re a **f**ew **fin**al **simplifi**cations **th**at **g**et **us t**o **t**he **algo**rithm **a**s **prese**nted **i**n **t**he **pap**er. **T**he **fir**st **is th**at, **sin**ce **we**'re **stor**ing **ea**ch $(x, y)$ **posi**tion **inde**xed **aga**inst $k$, **a**nd $k = x − y$, **we d**on't **ne**ed **t**o **sto**re $y$ **sin**ce **i**t **c**an **b**e **calcu**lated **fr**om **t**he **val**ues **o**f $k$ **a**nd $x$. **T**he **sec**ond **is th**at **we d**on't **ne**ed **t**o **sto**re **t**he **direc**tion **o**f **t**he **mo**ve **tak**en **a**t **ea**ch **st**ep, **we ju**st **sto**re **t**he **be**st $x$ **val**ue **we c**an **achi**eve **a**t **ea**ch **po**int. **T**he **pa**th **wi**ll **b**e **deri**ved *after* **we**'ve **comp**leted **th**is **proc**ess **t**o **f**ind **t**he **smal**lest $d$ **th**at **g**ets **us t**o **t**he **bot**tom-**rig**ht **posi**tion; **on**ce **we kn**ow **whe**re **t**he **fin**al **posi**tion **sho**ws **u**p **we c**an **backt**rack **t**o **f**ind **whi**ch **sin**gle **pa**th **o**ut **o**f **t**he **ma**ny **we**'ve **expl**ored **wi**ll **le**ad **us the**re.

**Remov**ing **tho**se **det**ails **lea**ves **u**s **wi**th **th**is **inform**ation:

```
      |    0    1    2    3    4    5
  ----+---------------------------------------
      |
```

```
 4  |                             7
    |
 3  |                       5
    |
 2  |                 3           7
    |
 1  |           1           5           7
    |
 0  |      0           2           5
    |
-1  |           0           4           5
    |
-2  |                 2           4
    |
-3  |                       3
```

The **final simplification is that the** *x* **values in the** *d*th **round depend only on those in the** (*d* − 1)th **round, and because each round alternately** modifies **either the odd or the even** *k* **positions, each round does not modify the values it depends on from the previous round. Therefore the** *x* **values can be stored in a single flat array, indexed by** *k*. **In our example, this array would evolve as follows with each value of** *d*:

```
      k |   -3    -2    -1    0     1     2     3     4
--------+-----------------------------------------------
        |
 d = 0  |                      0
        |
 d = 1  |                0     0     1
        |
 d = 2  |          2     0     2     1     3
        |
 d = 3  |    3     2     4     2     5     3     5
        |
 d = 4  |    3     4     4     5     5     7     5     7
        |
 d = 5  |    3     4     5     5     7     7     5     7
```

The **iteration stops when we discover we can reach** (*x*, *y*) = (7,6) **at** (*d*, *k*) = (5,1).

**We've now arrived at the representation of the problem used in the algorithm, and we can translate this into working code. We create a function that takes two lists, a and b, which will contain** `Diff::Line` **objects that represent lines in a file:**

```ruby
module Diff
  Line = Struct.new(:number, :text)

  def self.lines(document)
    document = document.lines if document.is_a?(String)
    document.map.with_index { |text, i| Line.new(i + 1, text) }
  end
end
```

**Our diff code will mostly only rely on the** `text` **field of these objects, but it's useful to store the original line numbers as it will make printing and other operations easier later. Since all the algorithms we'll look at will need to**

break a string into lines, let's make a utility function for doing that and then running the resulting lists through a diff algorithm of our choice.

```
module Diff
  def self.diff(a, b, differ: Myers)
    differ.diff(lines(a), lines(b))
  end
end
```

Now we'll begin writing the Myers class that will implement the algorithm we've been discussing. To start with, we'll just make some boilerplate for storing the two strings as instance variables on an object that implements the diff method, which we'll define later once all the building blocks are in place.

```
class Myers
  def self.diff(a, b)
    new(a, b).diff
  end

  def initialize(a, b)
    @a, @b = a, b
  end

  def diff
    # TODO
  end
end
```

To return a diff, we need to find the shortest edit path. We begin by storing n as the size of @a and m as the size of @b, and max as the sum of those; that's the most number of moves we might need to make.

```
def shortest_edit
  n, m = @a.size, @b.size
  max  = n + m
```

Then we set up an array to store the latest value of x for each k. k can take values from -max to max, and in Ruby a negative array index is interpreted as reading from the end of the array. The actual order of the elements doesn't matter, we just need the array to be big enough so that there's space for the positive and negative k values.

We set v[1] = 0 so that the iteration for $d = 0$ picks $x = 0$. We need to treat the $d = 0$ iteration just the same as the later iterations since we might be allowed to move diagonally immediately. Setting v[1] = 0 makes the algorithm behave as though it begins with a virtual move downwards from $(x, y) = (0, -1)$.

```
v    = Array.new(2 * max + 1)
v[1] = 0
```

Next, we begin a nested loop: we iterate d from 0 to max in the outer loop, and k from -d to d in steps of 2 in the inner loop.

```
        (0 .. max).step do |d|
          (-d .. d).step(2) do |k|
```

Within the loop, we begin by choosing whether to move downward or rightward from the previous round. If k is -d, or if it's not d and the k + 1 value is greater than the k - 1 value, then we move downward, i.e. we take the x value as being equal to the k + 1 value in the previous round. Otherwise we move rightward and take x as one greater than the previous k - 1 value. We calculate y from this chosen x value and the current k.

```
        if k == -d or (k != d and v[k - 1] < v[k + 1])
          x = v[k + 1]
        else
          x = v[k - 1] + 1
        end

        y = x - k
```

Having taken a single step rightward or downward, we see if we can take any diagonal steps. As long as we've not deleted the entire @a string or added the entire @b string, and the elements of each string at the current position are the same, we can increment both x and y. Once we finish moving, we store off the value of x we reached for the current k.

```
        while x < n and y < m and @a[x].text == @b[y].text
          x, y = x + 1, y + 1
        end

        v[k] = x
```

Finally, we return the current value of d if we've reached the bottom-right position, telling the caller the minimum number of edits required to convert from a to b.

```
          return d if x >= n and y >= m
        end
      end
    end
```

This minimal version of the algorithm calculates the smallest number of edits we need to make, but not the what those edits are. To do that, we need to back-track through the history of *x* values generated by the algorithm. We'll see how to do this in the next and final article in this series.