### The If Works

by James Coglan

**Op**en **sou**rce
**Boo**ks
**Conf**erence **tal**ks
**Podc**ast **inter**views

[**Bu**y **m**y **bo**ok, *Build*ing *Gi*t](#)

# Myers diff in linear space: theory

**If y**ou **enj**oy **th**is **arti**cle, I **ha**ve **publi**shed a **bo**ok **explain**ing **t**he **inter**nals **of G**it **thro**ugh **impleme**ntation: [**Build**ing **Gi**t](#).

—

**T**he **interest**ing **th**ing **abo**ut **di**ff **algor**ithms **i**s **th**at **th**ey're a **m**ix **of comp**uter **sci**ence **a**nd **hum**an **fact**ors. **The**re **a**re **ma**ny **equa**lly **go**od **dif**fs **betw**een **t**wo **fil**es, **judg**ing **th**em **b**y **t**he **len**gth **o**f **t**he **ed**it **sequ**ence, **a**nd **choos**ing **betw**een **th**em **req**uires **a**n **algor**ithm **th**at **c**an **be**st **mat**ch **peo**ple's **intui**tion **abo**ut **h**ow **the**ir **co**de **h**as **chan**ged.

**A**s **lu**ck **wou**ld **ha**ve **i**t, **a**s **so**on **a**s **I h**ad **fini**shed **u**p **m**y **prev**ious **arti**cles **o**n [**My**ers **dif**fs](#) **a**nd **go**ne **ba**ck **t**o **ma**ke **prog**ress **o**n **ano**ther **proj**ect, **I stum**bled **in**to a **ca**se **whe**re **G**it **prod**uced a **confus**ing **di**ff **f**or a **fi**le **I'**d **ju**st **chan**ged, **a**nd **I h**ad **t**o **kn**ow **wh**y. **He**re's **t**he **port**ion **of co**de **I h**ad **be**en **work**ing **o**n. **It'**s a **cou**ple **o**f C **funct**ions **th**at **co**py **byt**es **f**rom **o**ne **buf**fer **t**o **anot**her, **check**ing **t**he **siz**es **o**f **t**he **reque**sted **regi**ons **t**o **ma**ke **su**re **th**ey're **wit**hin **t**he **buf**fer. (**Th**is **i**s **n**ot **liter**ally **t**he **co**de **I w**as **work**ing **o**n; **I'v**e **remo**ved a **f**ew **th**ings **t**o **ma**ke **t**he **exam**ple **smal**ler.)

```c
void Chunk_copy(Chunk *src, size_t src_start, Chunk *dst, size_t dst_start, siz
{
    if (!Chunk_bounds_check(src, src_start, n)) return;
    if (!Chunk_bounds_check(dst, dst_start, n)) return;

    memcpy(dst->data + dst_start, src->data + src_start, n);
}

int Chunk_bounds_check(Chunk *chunk, size_t start, size_t n)
{
    if (chunk == NULL) return 0;

    return start <= chunk->length && n <= chunk->length - start;
}
```

**T**he **cha**nge **I ma**de **invo**lved **swapp**ing **t**he **ord**er **o**f **the**se **t**wo **funct**ions, **th**at **i**s:

```c
int Chunk_bounds_check(Chunk *chunk, size_t start, size_t n)
{
    if (chunk == NULL) return 0;

    return start <= chunk->length && n <= chunk->length - start;
}

void Chunk_copy(Chunk *src, size_t src_start, Chunk *dst, size_t dst_start, siz
{
    if (!Chunk_bounds_check(src, src_start, n)) return;
    if (!Chunk_bounds_check(dst, dst_start, n)) return;

    memcpy(dst->data + dst_start, src->data + src_start, n);
}
```

**Look**ing **a**t **th**is **cha**nge, **o**ne **wou**ld **expe**ct **th**at a **di**ff **f**or **i**t **sho**uld **sh**ow **t**he **shor**ter **func**tion Chunk_bounds_check **be**ing **mov**ed **t**o **app**ear **bef**ore Chunk_copy, i.e.:

```
+int Chunk_bounds_check(Chunk *chunk, size_t start, size_t n)
+{
+    if (chunk == NULL) return 0;
+
+    return start <= chunk->length && n <= chunk->length - start;
+}
+
 void Chunk_copy(Chunk *src, size_t src_start, Chunk *dst, size_t dst_start, si
 {
     if (!Chunk_bounds_check(src, src_start, n)) return;
     if (!Chunk_bounds_check(dst, dst_start, n)) return;

     memcpy(dst->data + dst_start, src->data + src_start, n);
 }
-
-int Chunk_bounds_check(Chunk *chunk, size_t start, size_t n)
-{
-    if (chunk == NULL) return 0;
-
-    return start <= chunk->length && n <= chunk->length - start;
-}
```

And **indeed**, **if** **y**ou **t**ry **runn**ing **the**se **t**wo **co**de **snip**pets **thro**ugh **o**ur **prev**ious **Mye**rs **impleme**ntation, **th**is is **exac**tly **wh**at **y**ou **g**et. **Howe**ver, **wh**en **y**ou **a**sk **G**it to **comp**are **the**se **vers**ions, **he**re's **wh**at **happ**ens:

```
-void Chunk_copy(Chunk *src, size_t src_start, Chunk *dst, size_t dst_start, si
+int Chunk_bounds_check(Chunk *chunk, size_t start, size_t n)
 {
-    if (!Chunk_bounds_check(src, src_start, n)) return;
-    if (!Chunk_bounds_check(dst, dst_start, n)) return;
+    if (chunk == NULL) return 0;

-    memcpy(dst->data + dst_start, src->data + src_start, n);
+    return start <= chunk->length && n <= chunk->length - start;
 }

-int Chunk_bounds_check(Chunk *chunk, size_t start, size_t n)
+void Chunk_copy(Chunk *src, size_t src_start, Chunk *dst, size_t dst_start, si
 {
-    if (chunk == NULL) return 0;
+    if (!Chunk_bounds_check(src, src_start, n)) return;
+    if (!Chunk_bounds_check(dst, dst_start, n)) return;

-    return start <= chunk->length && n <= chunk->length - start;
+    memcpy(dst->data + dst_start, src->data + src_start, n);
 }
```

**Th**is **di**ff is *correct*; **th**at **i**s, **i**t is a **leg**al **a**nd **mini**mal **ed**it **sequ**ence **th**at **trans**forms **t**he **fir**st **vers**ion **in**to **t**he **sec**ond, **a**nd **i**t **cont**ains **exac**tly **a**s **ma**ny **chan**ges **a**s **t**he "**expe**cted" **di**ff. **Howe**ver, **i**t is a *poor quality* **di**ff **i**n **th**at **i**t **fai**ls to **con**vey **t**he **hi**gh-**lev**el **mean**ing **o**f **t**he **cha**nge **t**o a **progr**ammer. **I**t **pres**ents **t**he **cha**nge **a**s a **ser**ies **o**f **arbitr**arily **interl**eaved **delet**ions **a**nd **inser**tions, **wher**eas **t**he **expe**cted **di**ff **clea**rly **sho**ws **t**he **t**wo **funct**ions **switch**ing **pla**ces.

**S**o, **w**hy **i**s **i**t **th**at **G**it **prod**uces **th**is **di**ff **ra**ther **th**an **t**he **o**ne **w**e'd **exp**ect? **W**e **c**an **beg**in **t**o **ans**wer **th**is **b**y **draw**ing **o**ut **t**he **ed**it **gra**ph **f**or **th**is **cha**nge, **ju**st **a**s **w**e **d**id **wh**en **w**e **fir**st **invest**igated **t**he **Mye**rs **algo**rithm.

```
     0   1   2   3   4   5   6   7   8   9   10  11  12  13  14
 0   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   |   |   |   |   | \ |   |   |   |   |   |
 1   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   | \ |   |   |   |   |   |   |   | \ |   |   |   |   |
 2   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   |   |   |   |   |   |   |   | \ |   |   |
```

```
 3   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   | \ |   |   | \ |   |   | \ |   |   |
 4   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   |   |   |   |   |   |   |   | \ |   |
 5   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   |   |   | \ |   |   |   |   |   | \ |
 6   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   | \ |   | \ |   |   |   | \ |   |   |
 7   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     | \ |   |   |   |   |   |   |   |   |   |   |   |   |
 8   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   | \ |   |   |   |   |   |   |   | \ |   |   |   |
 9   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   | \ |   |   |   |   |   |   |   |   |   |   |
10   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   | \ |   |   |   |   |   |   |   |   |   |
11   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   | \ |   | \ |   |   |   | \ |   |   |
12   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   |   | \ |   |   |   |   |   |   |
13   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
     |   |   |   |   |   | \ |   |   |   |   |   | \ |
14   o---o---o---o---o---o---o---o---o---o---o---o---o---o---o
```

You'll notice this graph contains two major diagonals, one of length 7 beginning at (0,7) and ending at (7,14), and one of length 6 beginning at (8,0) and ending at (14,6). These diagonals represent the two functions in the code that are not modified by the change, they simply trade places. Apart from these, there are many scattered short diagonals that correspond to braces and blank lines that appear in both versions. These types of lines appear multiple times in the text, whereas all the actual lines of code are unique and are matched just once between versions.

Now, the path corresponding to the expected diff is this one, where we first travel from (0,0) to (0,7), inserting the Chunk_bounds_check function, then from (0,7) to (7,14) to display the Chunk_copy function unmodified, then from (7,14) to (14,14) to delete the Chunk_bounds_check from its previous location.

```
     0   1   2   3   4   5   6   7   8   9   10  11  12  13  14
 0   o
     |
 1   o
     |
 2   o
     |
 3   o
     |
 4   o
     |
 5   o
     |
 6   o
     |
 7   o
       \
 8       o
           \
 9           o
               \
10               o
                   \
11                   o
                       \
12                       o
                           \
13                           o
```

```
                                                          \
          14                                   o---o---o---o---o---o---o---o
```

This path contains 7 **diagonals**, **and rec**all **th**at **t**he **a**im **of** **t**he **Myers**
**algo**rithm **is** **t**o **maxi**mise **t**he **num**ber **of** **diago**nals **tak**en – **th**at **is**, **t**o **fi**nd **t**he
**long**est **com**mon **subseq**uence – **s**o **a**s **t**o **ma**ke **a**s **f**ew **chan**ges **a**s **poss**ible.
**Howe**ver, **it** **tur**ns **o**ut **th**at **th**is **pa**th **is** **n**ot **uni**que **in** **contain**ing 7 **diago**nals,
**f**or **exam**ple, **he**re **is** **anot**her **pa**th **contain**ing **t**he **sa**me **num**ber:

```
           0   1   2   3   4   5   6   7   8   9   10  11  12  13  14
      0    o---o
           |
      1        o
                \
      2            o---o---o
                       |
      3                o
                        \
      4                    o---o
                               |
      5                        o
                                \
      6                            o
                                    \
      7                                o---o
                                           |
      8                                    o
                                            \
      9                                        o---o
                                                   |
     10                                            o
                                                   |
     11                                            o
                                                    \
     12                                                o---o
                                                           |
     13                                                    o
                                                            \
     14                                                        o
```

**If** **y**ou **comp**are **th**is **wi**th **t**he **diff**s **abo**ve, **y**ou'll **s**ee **th**at **th**is **is** **t**he **pa**th **tak**en
**b**y **G**it. (**The**re **a**re **ma**ny **oth**er **poss**ible **pa**ths **thro**ugh **th**ese **diago**nals **th**at
**y**ou **c**an **g**et **if** **y**ou **va**ry **t**he **ord**er **of** **right**ward **a**nd **down**ward **mov**es, **th**is **is**
**ju**st **t**he **o**ne **G**it **happ**ens **t**o **ta**ke.) **Beca**use **t**he **num**ber **of** **matc**hes **betw**een
**bra**ces **a**nd **bla**nk **lin**es – **tok**ens **th**at **hap**pen **t**o **b**e **t**he **sa**me **b**ut **ar**en't
**signif**icant **lin**es **of** **prog**ram **te**xt – **is** **equ**al **t**o **t**he **num**ber **of** **lin**es **in** **t**he
**ent**ire **matc**hed Chunk_copy **func**tion, **bo**th **pa**ths **a**re **equa**lly **legit**imate **in**
**ter**ms **of** **minimis**ing **t**he **num**ber **of** **edi**ts. **Howe**ver, **it**'s **n**ot **a** **use**ful **diff** **t**o **a**
**hum**an **rea**der, **a**nd **it** **wi**ll **prod**uce **mo**re **mer**ge **confl**icts **th**an **t**he **expe**cted
**diff** **beca**use **it** **doe**sn't **d**o **a**s **go**od **a** **j**ob **of** **separat**ing **sect**ions **of** **chan**ges
**fr**om **ea**ch **oth**er.

**B**ut **bef**ore **w**e **c**an **lo**ok **a**t **h**ow **t**o **fi**x **the**se **prob**lems, **we** **ne**ed **t**o **under**stand
**wh**at **cau**ses **th**em. **Giv**en **t**he **Mye**rs **algor**ithm **w**e've **s**een **s**o **f**ar **do**es **n**ot
**cho**ose **th**is **pa**th, **w**hy **do**es **G**it **cho**ose **i**t?

**T**he **ans**wer **is** **th**at **G**it **doe**sn't **u**se **t**he **algor**ithm **we** **implem**ented, **it** **us**es **a**
**varia**tion **of** **it**. **T**he **orig**inal **Mye**rs **algor**ithm **tak**es **a** **quadr**atic **amo**unt **of**
**spa**ce **t**o **calcu**late **t**he **ed**it **sequ**ence. **T**he **cent**ral **tri**ck **in** **t**he **Mye**rs
**algor**ithm **is** **th**at **rat**her **th**an **requir**ing **a** **gene**ral-**purp**ose **gra**ph **sea**rch
**algor**ithm, **wh**ich **wou**ld **expl**ore **a**nd **sto**re **a**ll **poss**ible **bran**ches, **we** **c**an find
the length of the longest common subsequence **us**ing **on**ly **a** **sin**gle **arr**ay **of**
**si**ze **propor**tional **t**o N + M, **t**he **s**um **of** **t**he **leng**ths **of** **t**he **t**wo **stri**ngs.

However if we want not just the length but the actual edit sequence, this requires storing a copy of the array as we move through the graph so we can backtrack at the end. In the worst case, it will take N + M moves to traverse the graph, and so the algorithm requires space proportional to (N + M)².

Requiring quadratic space is generally not a good property for an algorithm to have, since it rapidly becomes unfeasible to run the algorithm in memory as the input size increases. But fortunately, in the same paper where the original algorithm appears, Myers presents a modified version of it that runs in linear space, that is, space proportional to the sum of the string lengths. It is explained briefly in section 4b, but I'll give a worked example for our case here.

Whereas the original algorithm deterministically walks from the top-left to the bottom-right of the graph in a single pass, the linear space version uses a divide and conquer approach. The edit sequence is made up of a series of what Myers calls *snakes*, that is a rightward or downward step followed by zero or more diagonal ones. The linear-space version works by finding the *middle snake* of a possible edit path, that is a snake that crosses the halfway distance from top-left to bottom-right, and using the endpoints of that to divide the original graph into two smaller regions. It then works recursively on these regions until they're so small that no further work is required.

Here's an example: suppose that you were somehow able to determine that there exists a snake from (6,5) to (9,7), halfway between (0,0) and (14,14) in our example graph. (We will examine how to find this middle snake later.) This divides the original problem of getting from (0,0) to (14,14) into two smaller problems: getting from (0,0) to (6,5), and from (9,7) to (14,14). We can visualise this split like so:

```
        0   1   2   3   4   5   6   7   8   9   10  11  12  13  14
   0    o---o---o---o---o---o---o
        |   |   |   |   |   |   |
   1    o---o---o---o---o---o---o
        |   |   \   |   |   |   |
   2    o---o---o---o---o---o---o
        |   |   |   |   |   |   |
   3    o---o---o---o---o---o---o
        |   |   |   |   | \ |   |
   4    o---o---o---o---o---o---o
        |   |   |   |   |   |   |
   5    o---o---o---o---o---o---o
                                 \
   6                              @
                                   \
   7                                @---o---o---o---o---o---o
                                    |   |   |   |   |   |
   8                                o---o---o---o---o---o
                                    | \ |   |   |   |   |
   9                                o---o---o---o---o---o
                                    |   |   |   |   |   |
  10                                o---o---o---o---o---o
                                    |   |   |   |   |   |
  11                                o---o---o---o---o---o
                                    |   |   | \ |   |   |
  12                                o---o---o---o---o---o
                                    |   |   |   |   |   |
  13                                o---o---o---o---o---o
                                    |   |   |   |   | \ |
  14                                o---o---o---o---o---o
```

I have marked with a @ symbol those points that have been removed from the problem and are not part of any remaining sub-region of the graph.
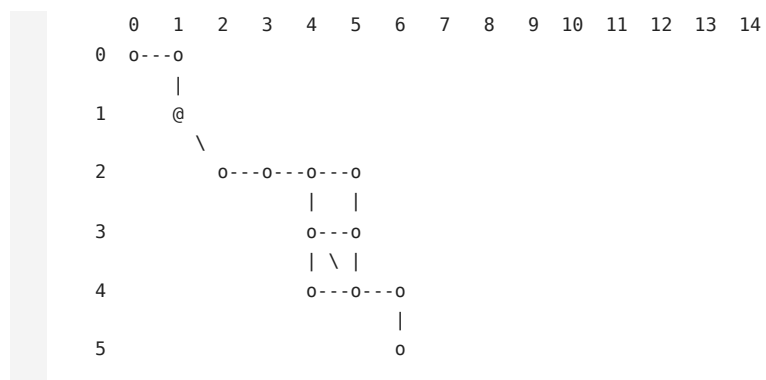
Having split the problem in two like this, we can apply the same technique to the sub-regions. In the (0,0)–(6,5) region there is a middle snake from (3,2) to (4,2), splitting the box into (0,0)–(3,2) and (4,2)–(6,5). In the (9,7)–(14,14) region the middle snake is from (11,10) to (11,11), splitting the box into (9,7)–(11,10) and (11,11)–(14,14).

```
        0   1   2   3   4   5   6   7   8   9   10  11  12  13  14
  0   o---o---o---o
        |   |   |   |
  1   o---o---o---o
        |   | \ |   |
  2   o---o---o---o---o---o---o
                    |   |   |
  3               o---o---o
                    | \ |   |
  4               o---o---o
                    |   |   |
  5               o---o---o
                                \
  6                                   @
                                        \
  7                                   @---o---o---o
                                        |   |   |
  8                                   o---o---o
                                        | \ |   |
  9                                   o---o---o
                                        |   |   |
  10                                  o---o---o
                                        |
  11                                          o---o---o---o
                                                | \ |   |   |
  12                                          o---o---o---o
                                                |   |   |   |
  13                                          o---o---o---o
                                                |   |   | \ |
  14                                          o---o---o---o
```
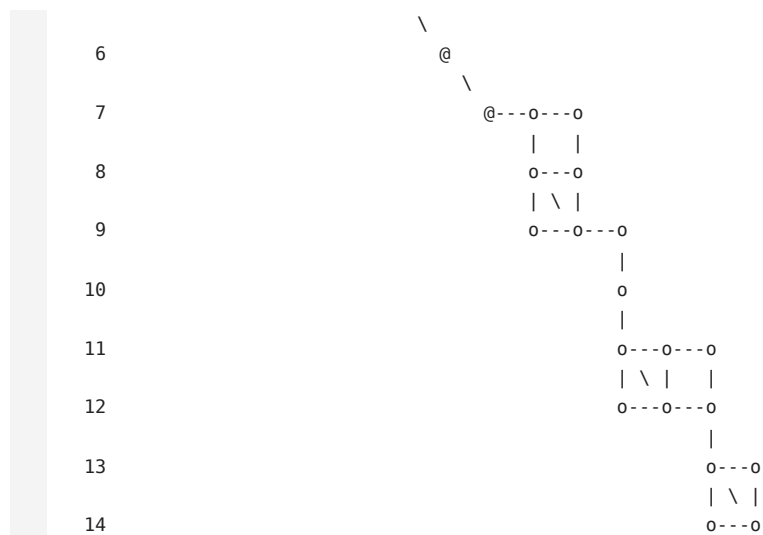
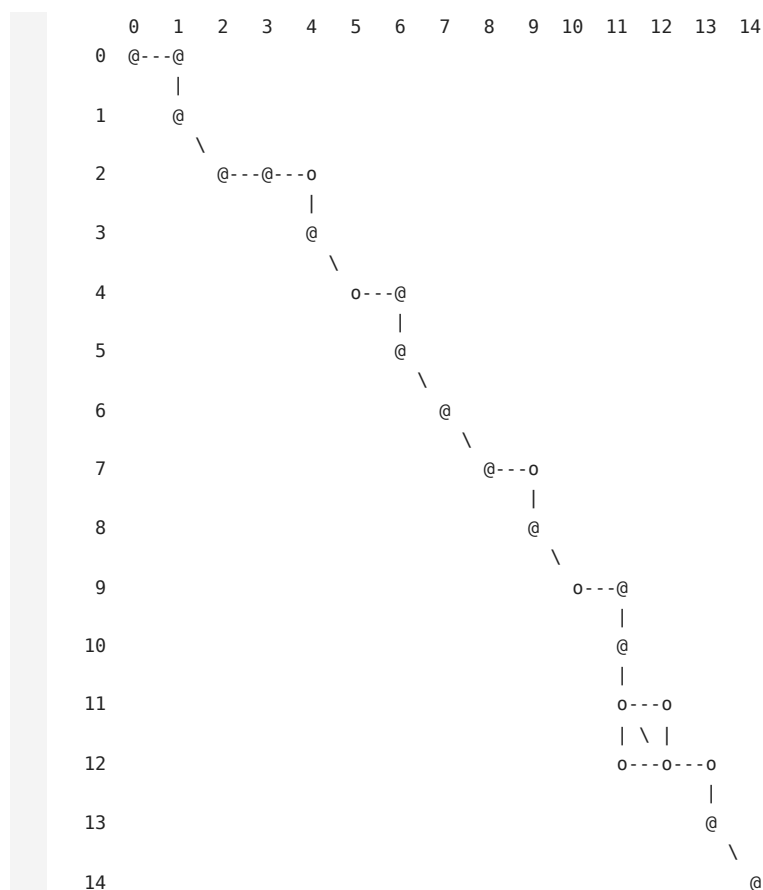Applying the process yet again to these four boxes gives the following splits:

- (0,0)–(3,2) → (0,0)–(1,0), (2,2)–(3,2)
- (4,2)–(6,5) → (4,2)–(5,4), (6,4)–(6,5)
- (9,7)–(11,10) → (9,7)–(10,9), (11,9)–(11,10)
- (11,11)–(14,14) → (11,11)–(13,12), (13,13)–(14,14)

Notice that some of the resulting boxes have zero area, and some have corner points in common. This is perfectly normal, for example in the space from (4,2) to (6,5) below there are two resulting sub-regions: a 1-by-2 box (4,2)–(5,4), and a 0-by-1 box (6,4)–(6,5).

```
        0   1   2   3   4   5   6   7   8   9   10  11  12  13  14
  0   o---o
            |
  1       @
              \
  2               o---o---o---o
                    |   |
  3               o---o
                    | \ |
  4               o---o---o
                        |
  5                   o
```

```
              \
   6           @
                \
   7           @---o---o
               |   |
   8           o---o
               | \ |
   9           o---o---o
               |
  10           o
               |
  11           o---o---o
               | \ |   |
  12           o---o---o
               |
  13               o---o
                   | \ |
  14               o---o
```

One **final** **pa**ss **a**nd **y**ou **c**an **s**ee **the**re **a**re **ve**ry **few** **po**ints **le**ft **t**o **vis**it, **a**nd **all** **o**f **th**em **a**re **p**art **of** **ze**ro-**ar**ea **o**r **un**it **box**es **whe**re **t**he **pa**th **f**rom **o**ne **e**nd **t**o **t**he **oth**er **i**s **obvi**ous. **Wh**at **w**e're **le**ft **w**ith **i**s **t**he **pa**th **foll**owed **b**y **G**it's **di**ff.

```
          0   1   2   3   4   5   6   7   8   9  10  11  12  13  14
   0  @---@
          |
   1      @
           \
   2          @---@---o
                  |
   3              @
                   \
   4                  o---@
                          |
   5                      @
                           \
   6                          @
                               \
   7                          @---o
                                  |
   8                              @
                                   \
   9                              o---@
                                      |
  10                                  @
                                      |
  11                              o---o
                                  | \ |
  12                              o---o---o
                                          |
  13                                      @
                                           \
  14                                          @
```

**Th**is **i**s **a** **ne**at **tri**ck, **a**nd **i**t's **actu**ally **t**he **rea**son **th**is **vari**ant **o**f **t**he **algor**ithm **c**an **disc**over **mult**iple **pat**hs **thro**ugh **t**he **gra**ph, **a**s **w**e'll **s**ee **short**ly. **B**ut **t**o **g**et **i**t **t**o **wo**rk, **w**e **ne**ed **t**o **kn**ow **h**ow **t**o **f**ind **tho**se **mid**dle **sna**kes. **Th**is **see**ms **l**ike **a** **chic**ken-**a**nd-**e**gg **prob**lem: **i**n **ord**er **t**o **ma**ke **a** **work**ing **di**ff **algor**ithm, **w**e **f**irst **ne**ed **someth**ing **th**at **c**an **f**ind **opti**mal **ed**it **pat**hs, **whic**h **i**s **ju**st **wh**at **t**he **orig**inal **Mye**rs **algor**ithm **do**es!

**I**n **fa**ct, **th**is **i**s **exac**tly **wh**at **happ**ens. **T**o **f**ind **t**he **mid**dle **sna**ke, **w**e **r**un **t**he **orig**inal **Mye**rs **algor**ithm **b**oth *forward* **f**rom **t**he **t**op-**le**ft, **a**nd *backward* **f**rom **t**he **b**ottom-**righ**t, **unt**il **w**e **f**ind **a** **po**int **a**t **whic**h **val**ues **o**n **t**he **sa**me **diag**onal **me**et **ea**ch **oth**er. **Th**at **prob**ably **sou**nds **a** **b**it **cryp**tic **s**o **l**et's **lo**ok **a**t **a** **sma**ll **exam**ple.
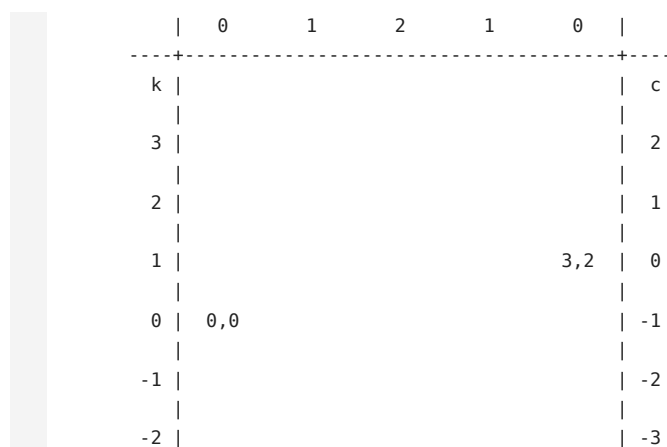
Say we're working on the (0,0)–(3,2) box. Let's reproduce it here, labelled with its (x, y) values, and also with each diagonal labelled with its value of *k*. Recall that $k = x - y$ and it either increases or decreases by 1 each time we make a move. Similarly, when working on the backward scan, we will label diagonals with a value *c* which is zero at the bottom-right corner, increases with each step upward, and decreases with each step leftward. The relationship between *k* and *c* is:

$k = c + delta$, where *delta* = *width* − *height*, the difference between the box's dimensions.

```
          x                                k
                                     0     1     2     3
          0     1     2     3        \     \     \     \
    y  0  o-----o-----o-----o          o-----o-----o-----o
          |     |     |     |     -1   |     |     |     | \
          |     |     |     |        \ |     |     |     |   2
       1  o-----o-----o-----o          o-----o-----o-----o
          |     | \   |     |     -2   |     | \   |     | \
          |     |   \ |     |        \ |     |   \ |     |   1
       2  o-----o-----o-----o          o-----o-----o-----o
                                          \     \     \     \
                                         -3    -2    -1     0
                                                              c
```

We run the original algorithm forward from (0,0) and backward from (3,2). Let's begin exploring the graph, marking our progress as we go. Just like our previous graph-walking explorations, these diagrams place *d*, the number of moves taken, on the horizontal axis, and *k*, the number of the current diagonal, on the vertical axis, and record (x, y) points as we walk the graph. But unlike before, we also have the vertical axis labelled with *c*, which I've displayed on the right hand side, and we'll include the points in the backward scan beginning from this side. Points on the same horizontal line in the forward and backward scan are on the same diagonal in the edit graph, and the forward and backward scans will overlap when a point in the forward scan has a greater *x* value than a point on the same diagonal in the backward scan.

We begin by marking our starting positions at opposite ends of a grid. (0,0) appears at (d, k) = (0,0), and (3,2) at (d, c) = (0,0).

```
       |  0      1      2      1      0   |
   ----+---------------------------------+----
    k  |                                  |  c
       |                                  |
    3  |                                  |  2
       |                                  |
    2  |                                  |  1
       |                                  |
    1  |                            3,2   |  0
       |                                  |
    0  |  0,0                             |  -1
       |                                  |
   -1  |                                  |  -2
       |                                  |
   -2  |                                  |  -3
```

Next we explore the graph to one level deep, that is we move to (d, k) = (1,−1) and (1,1) and similarly for (d, c).

Moving downward from (0,0) lands us at (0,1), and moving rightward at (1,0). Likewise, moving leftward from (3,2) takes us to (2,2), which has a diagonal leading to (1,1), and moving upward takes us to (3,1).

```
       |  0      1      2      1      0   |
  ----+-----------------------------------+----
   k |                                    |  c
     |                                    |
   3 |                                    |  2
     |                                    |
   2 |                            3,1     |  1
     |                               \    |
   1 |           1,0                3,2  |  0
     |           /                  /     |
   0 |  0,0                  1,1          | -1
     |     \                              |
  -1 |           0,1                      | -2
     |                                    |
  -2 |                                    | -3
```

**S**o far, **the**re is **n**o **opport**unity **f**or the **forw**ard **a**nd **back**ward **par**ts **t**o **me**et; **t**he **forw**ard **pa**th **is o**n *k* **diago**nals −1 **a**nd 1 **whi**le **t**he **back**ward **pa**th **is o**n *k* **diago**nals 0 **a**nd 2.

**Let**'s **ta**ke **o**ne **mo**re **mo**ve **i**n the **forw**ard **dire**ction. **W**e **c**an **mo**ve **down**ward **fr**om (0,1) **t**o **rea**ch (0,2), **a**nd **fr**om (1,0) **t**o (1,1), **tak**ing **a diag**onal **mo**ve **t**o (2,2). **W**e **c**an **mo**ve **right**ward **fr**om (1,0) **t**o (2,0).

```
       |  0      1      2      1      0   |
  ----+-----------------------------------+----
   k |                                    |  c
     |                                    |
   3 |                                    |  2
     |                                    |
   2 |                   2,0    3,1       |  1
     |                   /         \      |
   1 |           1,0                3,2  |  0
     |           /      \           /     |
   0 |  0,0             2,2    1,1        | -1
     |     \                              |
  -1 |           0,1                      | -2
     |                  \                 |
  -2 |                   0,2              | -3
```

**N**ow, **we**'ve **fou**nd **a**n **over**lap: the **poi**nt (2,2) **w**as **fou**nd **a**t (*d*, *k*) = (2,0), **a**nd the **poi**nt (1,1) **a**t (*d*, *c*) = (1,−1), **whi**ch **giv**en **t**he **relati**onship **betw**een *k* **a**nd *c* **mea**ns **t**hey **liv**e **o**n **t**he **sa**me **diag**onal. **A**nd, (2,2) **o**n **t**he **forw**ard **pa**th **i**n (*x*, *y*) **spa**ce **is t**o **t**he **rig**ht **of** (1,1) **o**n **t**he **back**ward **pa**th, **a**nd **s**o **the**se **pa**ths **h**ave **over**lapped. **Th**is **mea**ns **we**'ve **fou**nd **o**ur **mid**dle **sna**ke.

**Wh**en *delta* **is ev**en, **t**he **opti**mal **pa**th **len**gth **will al**so **b**e **ev**en, **a**nd **s**o **we will f**ind **a**n **over**lap **aft**er **tak**ing **bo**th **a forw**ard **a**nd **a back**ward **mo**ve **o**n **ea**ch **itera**tion **of** *d*. **B**ut **wh**en **it is o**dd, **w**e **will f**ind **t**he **over**lap **aft**er **ju**st **tak**ing **a forw**ard **mo**ve **o**n **t**he **fin**al **itera**tion, **aft**er **seve**ral **tur**ns **of bo**th **forw**ard **a**nd **back**ward **mo**ves. **A**nd **s**o **b**y **conve**ntion, **wh**en *delta* **is o**dd, **w**e **cons**ider **t**he **forw**ard **mo**ve **w**e **ju**st **ma**de **t**o **b**e **t**he **mid**dle **sna**ke, **wher**eas **wh**en **it**'s **ev**en, **w**e **u**se **t**he **la**st **back**ward **mo**ve.

**I**n **t**his **ca**se, *delta* **is o**dd **a**nd **s**o **t**he **mid**dle **sna**ke **is t**he **la**st **forw**ard **mo**ve, **fr**om (1,0) **t**o (2,2). **There**fore **we c**an **spl**it **o**ur **orig**inal **b**ox **in**to (0,0)–(1,0) **a**nd (2,2)–(3,2).

```
          x

          0    1    2    3
    y  0  o-----o
```

```
                1


                2                       o-----o
```

Notice that, although I've drawn the lines of motion on the traces above, that's purely to help visualise what's going on. Since we're only doing the forward scan of the original algorithm and we're not going to backtrack at the end to calculate the full path, we don't need to keep copies of the history of the walk as we go - we only need to keep the current state. That's what allows this variant to work in linear space; we only need to allocate two arrays of size *width* + *height* + 1, half the size of the array in the original version since we're only going to scan halfway. We use one array for the forward scan and one for the backward scan.

Let's do one more example; this will be in the middle of the grid and also shows up a situation where the algorithm can make choices between equally good paths. We'll take the (4,2)–(6,5) box.

```
           x                       k
                                   0     1     2
           4     5     6            \     \     \
   y   2   o-----o-----o              o-----o-----o
           |     |     |         -1   |     |     | \
           |     |     |              \ |     |     |   3
       3   o-----o-----o              o-----o-----o
           | \   |     |         -2   | \   |     | \
           |   \ |     |              \ |   \ |     |   2
       4   o-----o-----o              o-----o-----o
           |     |     |         -3   |     |     | \
           |     |     |              \ |     |     |   1
       5   o-----o-----o              o-----o-----o
                                       \     \     \
                                       -2    -1     0
                                                    c
```

Just as before, we mark our starting points in appropriate positions in (*d*, *k*) space.

```
          |   0      1      2      1      0   |
      ----+---------------------------------+----
      k |                                  |  c
        |                                  |
      2 |                                  |  3
        |                                  |
      1 |                                  |  2
        |                                  |
      0 |   4,2                            |  1
        |                                  |
     -1 |                           6,5 |  0
        |                                  |
     -2 |                                  | -1
        |                                  |
     -3 |                                  | -2
```

We follow the graph one move. These moves are all a single step, except for the move downward from (4,2) which can take a diagonal to (5,4). So far, no overlaps are possible.

```
          |   0      1      2      1      0   |
      ----+---------------------------------+----
      k |                                  |  c
        |                                  |
      2 |                           |  3
        |                                  |
      1 |          5,2                     |  2
```

```
      |           /                        |
    0 |    4,2                    6,4       |   1
      |           \                    \    |
   -1 |             5,4                  6,5 |   0
      |                                /    |
   -2 |                         5,5         |  -1
      |                                     |
   -3 |                                     |  -2
```

Now one more move. After this round, there are two overlaps: on the $k=-2$ diagonal both paths have found (5,5), and on the $k=0$ diagonal both have found (6,4). We now need to make a choice about which move we're going to consider our middle snake.

```
      |   0       1       2       1       0   |
  ----+---------------------------------------+----
    k |                                       |   c
      |                                       |
    2 |                   6,2                 |   3
      |                 /                     |
    1 |           5,2                         |   2
      |         /                             |
    0 |   4,2             6,4     6,4         |   1
      |         \       /               \     |
   -1 |           5,4                     6,5 |   0
      |                 \               /     |
   -2 |                   5,5     5,5         |  -1
      |                                       |
   -3 |                                       |  -2
```

We know from the eventual outcome that Git picks (6,4) as the midpoint here, but why? To find out, I went looking in [xdiffi.c](#) in the Git source tree. This is actually part of [LibXDiff](#), originally developed by Davide Libenzi and then imported and modified by the Git authors.

It actually boils down to the order in which Git visits graph locations: it iterates over $k$ (or $c$ in the backward direction) [in descending order, from $d$ to $-d$](#). That means in this final round it generates point (6,2) in the forward direction, and then (6,4), and then it stops since it's found an overlap. It never finds point (5,5) in the forward scan. If you change it to iterate from $-d$ to $d$ in ascending order, then it generates the same diff as the original Myers method. It's not making a decision between options as such, it's just an accident of implementation which one it picks, although this iteration order may have been chosen to prefer paths on the upper diagonals.

Having chosen (6,4) as the split point, that gives the following boxes after removing the middle snake.

```
          x

          4     5     6
  y   2   o-----o
          |     |
          |     |
      3   o-----o
          | \   |
          |   \ |
      4   o-----o     o
                      |
                      |
      5                o
```

A similar choice presents itself on the top-level scan of the whole graph. As well as the snake from (6,5) to (9,7), there is one from (0,7) to (7,13), which

is found via a path that took the diagonal from (14,14) to (13,13), but Git choses the first simply because it iterates over diagonals in descending order.

This process shows why it is that Git can choose a different path to the obvious one: rather than the original algorithm which is searching for a single point: the bottom-right, this algorithm might find many possible candidates for a midpoint and faces the problem (or opportunity!) of choosing between them. A lot of subtle implementation details affect the decisions it makes: what order we iterate diagonals in, whether we optimise for *x* or *y* in the forward and backward scans, which direction we pick when the two previous points have the same *x* value. We'll see some of these choices in the next article when we come to implement this algorithm.

---