

Uncovering Environmental Change in the English Lake District: Using Computational Techniques to Trace Shifting Practice in the Historical Documentation of Flora

Journal:	<i>Digital Scholarship in the Humanities</i>
Manuscript ID	DSH-2019-0078.R3
Manuscript Type:	Full Paper
Date Submitted by the Author:	11-Jun-2020
Complete List of Authors:	Smail, Robert; Lancaster University, Lancaster Environment Centre Donaldson, Chris; Lancaster University, History Govaerts, Rafaël; Royal Botanic Gardens Kew, Biodiversity Informatics and Spatial Analysis Rayson, Paul; Lancaster University, Computing Stevens, Carly; Lancaster University, Lancaster Environment Centre
Keywords:	Floristic change, Historical flora, Geoparsing, Digital humanities, Corpus linguistics, Geographic information systems (GIS), English Lake District, Natural Language Processing (NLP), Plants of the World Online (POWO)

SCHOLARONE™
Manuscripts

Digital Scholarship in the Humanities

Uncovering Environmental Change in the English Lake District: Using Computational Techniques to Trace the Presence and Documentation of Historical Flora

1 Introduction

The English Lake District has long been regarded as a place of outstanding cultural and environmental importance. Situated in North West England, the region has been described and celebrated by naturalists, poets and painters since the late 1600s and ranks among one of the most iconic upland landscapes in Europe. In 2017 the Lake District became a UNESCO World Heritage Site (WHS) under the 'cultural landscape' category. UNESCO introduced its 'cultural landscape' category in the 1990s in order to establish guidelines for acknowledging places whose 'outstanding universal value' derives from 'the combined work of nature and mankind' (Denyer, 2016; UNESCO, 2017). The official designation of the Lake District as a cultural landscape has raised questions about how the character of the region has been shaped by natural processes and human industry. These questions are of interest for historians and environmental scientists, but they are also important for heritage and conservation organisations in the region, including the National Trust, Natural England and the Lake District National Park Authority, who are under mounting pressure to preserve, to protect and (in some cases) to restore the Lake District's historical environmental character. Answering these questions is not easy, not least because we lack a sufficiently coherent body of empirical evidence about environmental conditions in the Lake District before the mid-twentieth century (Lake District National Park, n.d.; UNESCO, 2017).

1
2
3
4 The root of the problem is not a lack of evidence, but instead the need for a
5
6 methodology that can enable researchers to compile and analyse evidence from a body of
7
8 otherwise disparate historical sources. In this article, we demonstrate the implementation of
9
10 such a methodology. Using a combination of computational techniques, we show how it is
11
12 possible to consolidate and interrogate a diverse assortment of texts that provide evidence
13
14 about the environmental character of the Lake District from the seventeenth century to the
15
16 twentieth century. These texts contain a wealth of information about the region's historical
17
18 environment, including accounts of flora, fauna and weather conditions. For the purposes of
19
20 this study, we have chosen to focus on information relating to the region's flora, and we have
21
22 done so for two reasons: firstly, because the Lake District's flora has received particular and
23
24 sustained interest historically, meaning that there is an extensive body of empirical evidence
25
26 dating back to the seventeenth century on which one can draw; secondly, we have selected
27
28 flora because it is a good indicator of broader environmental changes (Ellenberg 1974). By
29
30 tracing changes in plant species distribution and composition, one can also trace the
31
32 transformation of habitats and landscape characteristics more generally.
33
34
35
36
37
38

39 The method we present combines techniques from the digital and spatial humanities,
40
41 including Natural Language Processing (NLP), Named Entity Recognition (NER), corpus
42
43 linguistics (CL) and Geographic Information Systems (GIS). These techniques have been
44
45 shown to be effective in guiding the investigation and interrogation of geospatial themes
46
47 across historical textual corpora (Donaldson, Gregory, & Taylor, 2017; Gregory &
48
49 Donaldson, 2016). A few pioneering research projects have, moreover, indicated the
50
51 benefits of applying these sorts of techniques for environmental history (Cervera, Pino,
52
53 Marull, Padró, & Tello, 2019; Hinrichs et al., 2015), and there have recently been
54
55 comparable developments in the environmental sciences (Roll, Correia, & Berger-Tal, 2018).
56
57
58
59
60

1
2
3 As Kherai and Oosthoek note, however, this is still 'an emerging area' (Kheraj & Oosthoek,
4 2016, p. 245), and as yet few studies have attempted to use computational methods to
5
6 analyse source materials drawn from across the sort of broad span of time that we are
7
8 concerned with in this study. We are therefore keen to take this opportunity to advance
9
10 environmental research, and in the following pages we do so by demonstrating three things:
11
12 firstly, how information about the Lake District's flora can be identified across a range of
13
14 digitised historical sources on a large scale; secondly, how this information can be extracted
15
16 and transformed into a structured dataset; thirdly, how this dataset can be plotted and
17
18 analysed geospatially using digital tools including GIS. The method we outline here can be
19
20 adapted to support research into the environmental history of other locations.
21
22
23
24
25
26

27 The historical texts we shall consider constitute a corpus of 92 digitised works
28
29 published between 1682 and 1904 (a list of these works appears in the appendix).
30
31 Collectively, these texts total nearly 19 million words and include examples of genres
32
33 ranging from regional guides and travelogues to scientific journal articles and reports. The
34
35 historical breadth and variety of the corpus is deliberate, as working with both non-scientific
36
37 and scientific writing enables one to conduct more historically nuanced analyses. Selecting
38
39 texts published between the late seventeenth century and the early twentieth century was
40
41 also deliberate, as it enables us to study the composition of the Lake District's flora over a
42
43 much longer time frame than current datasets allow. We chose not to include works
44
45 published after 1904 due to UK copyright laws, which constrain our ability to work freely with
46
47 some more recent sources. All the texts in the corpus have been compiled from existing
48
49 open-source material, which collectively presents a uniquely rich and diverse body of
50
51 information about the Lake District's environmental history.
52
53
54
55
56
57
58
59
60

1
2
3
4 Our analysis will focus on how our computational methodology can enable
5
6 researchers to track changes in the language used to describe and document plant species
7
8 across the corpus and therefore across time. Plant taxonomy and nomenclature underwent
9
10 sustained change during the period represented in our corpus. To uncover broader themes
11
12 collectively registered across the corpus, it is essential that our computational methodology
13
14 can trace variations and changes in naming conventions. As part of our investigation, we
15
16 have drawn on the *Plants of the World Online* (POWO) database (POWO, 2018). Based on
17
18 the extensive plant collections of the Royal Botanical Gardens, Kew, POWO is one of the
19
20 world's most extensive databases of historical binomial plant name synonyms. Our use of
21
22 this database reveals the importance of adopting a temporally sensitive approach that allows
23
24 for synonyms to be evaluated in tandem with modern standardised plant names.
25
26
27
28
29

30 The methodology we outline and the dataset we derived from the corpus contains
31
32 information about 802 plant species, 510 (63.5%) of which are linked to locations within the
33
34 boundary of the Lake District National Park. This boundary, although not established until
35
36 1951, broadly corresponds with both the principal area of interest documented by the texts in
37
38 the corpus and the area in which the members of the Lake District World Heritage Site
39
40 Partnership (LDWHSP) have a vested interest. ~~In tracing the evolving language used to
41
42 describe plants we have been able to form a richer dataset that more accurately reflects the
43
44 historical observations recorded across our corpus. By uncovering changes in plant naming
45
46 conventions over time, we have been able to increase our understanding of where different
47
48 plant species were observed historically and also how they were documented in the past.~~ By
49
50 extracting information about plant species from the corpus and collating that information in a
51
52 structured and accessible digital form, this work stands to make an important contribution to
53
54 academic researchers investigating the historical distribution and composition of plant
55
56
57
58
59
60

1
2
3 species in the Lake District. We have also increased and enriched the historical knowledge
4
5 available to organisations, including members of the LDWHSP, who are directly involved in
6
7
8 landscape management and policy decisions in the region.
9

10 11 12 13 14 15 **2 Background**

16 17 18 *2.1 The natural environment as recorded in historical accounts*

19
20
21
22 The flora of the Lake District has been recorded with increasing enthusiasm and dedication
23
24 since the seventeenth century. Naturalists including John Ray (1627–1705), Archbishop
25
26 William Nicholson (1655–1727) and Thomas Lawson (c.1630–1691) are credited with
27
28 producing some of the earliest records of Lake District flora (Arber, 1943; E. Jean Whittaker,
29
30 1981; Hodgson & Goodchild, 1898, p. xxiii–xxiv). Over the course of the eighteenth and
31
32 nineteenth centuries, interest in uncovering and recording the Lake District's natural
33
34 environment grew as the region attained greater renown (Denyer, 2016; Lindop, 2005;
35
36 Nicholson, 1955). In addition to naturalists, the Lake District also attracted the attention of
37
38 tourists, travellers, writers and poets, many of whom published accounts that contain
39
40 information about the region's flora. ~~Consequently, there is an extensive and wide-ranging~~
41
42 ~~body of source material on which to draw.~~ This material provides a wealth of empirical
43
44 information about historical flora in the Lake District. The problem, however, is that this
45
46 material appears in a variety of forms, which makes it difficult to consult and to collate. To
47
48 understand how these sources can be combined we first need to understand the different
49
50 sorts of information being described and the varying ways it is presented.
51
52
53
54
55
56
57
58
59
60

1
2
3
4 To illustrate this point, consider briefly the species *Drosera anglica* (*English sundew*
5 or *great sundew*) and its description in two texts. *Drosera anglica* is a species commonly
6 found in the wetter parts of the region, including bogs and lake shores. The species has
7
8 been declining in England due to drainage, eutrophication and peat extraction (Stace, 2005,
9
10 p. 217). Describing the species in *Flora of The Lake District* (1885), John Baker writes:

11
12
13
14
15
16
17
18 *Drosera anglica*, Huds. (Great Sundew). Native.

19
20 Scottish type. Range i.

21
22
23 C[umberland]. Ullock Moss near Portinscale.—(W. Dickinson.) Helvellyn,—(J. Flintoft.) Moss
24
25 at Grange, abundant.

26
27 (J. C. Melvill.) Seathwaite in Borrowdale.—(Miss Edmunds.) Side of Crummock.—(W; B.
28
29 Waterfall.)

30
31
32 W[estmorland]. Foulshaw Moss and Brigstear Moss near Kendal, First
33
34 recorded by Wilson.

35
36
37 L[ancashire]. Stickle Pike, Donnerdale.—(W. F. Miller.)

38
39 (Baker, 1885, p. 44)

40
41
42
43
44 Baker (1834–1920) was a Fellow of the Royal Society and principal assistant to the
45
46 Keeper of the Kew Herbarium (he became Keeper in 1890; see Desmond, 1977, p. 36). His
47
48 *Flora* was the culmination of years of patient research, and it was intended for the botanical
49
50 community ~~(for the development of regional floras; see, David E Allen, 2003, p. 271–280),~~
51
52 providing a detailed (if rather dry) account of the different species found in the region.
53
54 Baker's *Flora* is essentially a reference work, it contains concise descriptions that feature
55
56 specialised abbreviations and technical terminology. In this format, *Drosera anglica* is
57
58
59
60

1
2
3 documented at eight different locations around the Lake District. In addition, the species is
4
5 also noted as being 'native' to the region and growing at *range i.*, which means it was
6
7 commonly observed to grow from sea level up to an altitude of 900 feet.
8
9

10 Contrast this with the description of *Drosera anglica* in Frederick Malleeson's (1819–
11
12 1897) *Holiday Studies of Wordsworth by Rivers, Woods and Alps*. Malleeson writes:
13
14

15
16
17 From the [Ulpha] bridge the main road leads along the foot of the Dunnerdale
18
19 Fells to Broughton and to Millom. The scenery, though grand and noble, is bare
20
21 and wild. In the boggy streams running off these fells is found the great sundew
22
23 (*Drosera anglica*), a rare plant. (Malleeson, 1890, p. 68)
24
25
26
27
28
29

30 This account is rather different than Baker's. A minister in the Church of England,
31
32 Malleeson served as vicar of Broughton-in-Furness between 1870 and 1897. The
33
34 observations just quoted were made during a two-day tour through the Duddon Valley in
35
36 1882. This tour, he tells us, was inspired by the writings of the Romantic poet William
37
38 Wordsworth. Specifically, Malleeson is referring to Wordsworth's River Duddon sonnets,
39
40 which were published in 1820. Malleeson's aim is to impart a vivid account of the valley's
41
42 scenery: the surroundings are described as 'surpassingly beautiful', as both 'grand' and
43
44 'noble' though 'bare and wild'. A competent botanist, Malleeson offers an account of *Drosera*
45
46 *anglica* that is knit into his poetically inspired account of this landscape; flora is intrinsically
47
48 linked to the emotional impressions the landscape makes upon him and which, in turn, he
49
50 attempts to impress upon the reader. Though the account is quite different to Baker's (both
51
52 in its purpose and its style), it too imparts some useful information about *Drosera anglica*.
53
54 Malleeson informs us of a locality in which the plant can be found. He also records its
55
56
57
58
59
60

1
2
3 conservation status, telling us it is a 'rare' species, and he gives us some details about its
4
5 habitat, noting that it can be found in 'boggy streams running off these fells'.
6
7

8 Despite their quite different outward appearances and intended audiences, both
9
10 commentators provide insightful accounts of Lake District flora. Moving past differences in
11
12 language, descriptive styles and layout, commonalities can be discerned regarding the flora
13
14 being described. In addition to recording a specific plant species, these accounts indicate
15
16 where those species can be found, as well as providing details about the plant's traits and
17
18 habitat. There is therefore a clear benefit in bringing these two sources together and
19
20 combining them with other accounts of *Drosera anglica*. Doing so enables us to build a more
21
22 comprehensive picture of the historical presence and distribution of plant species in the Lake
23
24 District.
25
26
27
28
29
30
31

32 *2.2 Correlating empirical evidence across historical sources*

33
34
35
36

37 As we have seen in the previous section, historical sources can provide a rich documentary
38
39 record of Lake District flora and its distribution. Taken in isolation, individual accounts give
40
41 us glimpses into the activities and interests of specific individuals, furnishing us with a
42
43 snapshot of flora seen at specific places and points in time. ~~In order to~~ To improve
44
45 further advance our understanding of the region's historical environmental character, and to
46
47 begin to examine changes, it is necessary to combine and correlate evidence from multiple
48
49 sources. Consequently, it is important that the sources collectively provide sufficient
50
51 geographic and temporal coverage. By geographic coverage, we refer to how observations
52
53 of a specific site can be iteratively combined with observations of other sites to enable one to
54
55 *construct* a picture of the whole Lake District. Similarly, by temporal coverage, we refer to
56
57
58
59
60

1
2
3 how accounts of specific sites made at different points in time can be compared to determine
4
5
6 if and when local environmental conditions changed.
7

8 As previously noted, Lake District flora has received interest from naturalists, plant
9
10 collectors, tourists, travellers, writers and poets alike. The individuals who constituted these
11
12 groups often had different motivations for exploring and recording the region's flora. These
13
14 motivations influenced the aspects of the natural environment they investigated and
15
16 recorded, including which plants they documented and which they overlooked. Even
17
18 observations collected by the most studious scientific researcher are still likely to contain
19
20 some degree of bias, reflecting the time spent observing a particular site, the skill level of the
21
22 observer and their specific interests, as well as the standards and conventions of scientific
23
24 practice at that time (this is a persisting issue in botanical field research; see, Rich, 1997).
25
26 Drawing from multiple 'witness groups', as they will be termed here, serves to elucidate the
27
28 interests of each group and to provide a sense of perspective. It is also essential if we are to
29
30 reveal a broader narrative of historical Lake District flora that is not limited to a single group.
31
32
33
34
35

36 How might we expect observational records to vary between different witness
37
38 groups? Genre aside, specialisation is one obvious cause of variation. Here, we do not use
39
40 the term in its more rigid sense to mean specialisation between scientific disciplines, as it
41
42 has been shown that the boundaries between disciplines in natural history were never to
43
44 become rigidly demarcated (David Elliston Allen, 1994; Kuklick & Kohler, 1996; Nicholas
45
46 Jardine, James A. Secord, 1996). Rather, specialisation refers here to the particular interests
47
48 of each observer: that which led them to favour the recording of specific taxonomic groups,
49
50 habitat types or geographical localities over others. One might turn to the species *Impatiens*
51
52 *noli-tangere* (*Touch-me-not-balsam*) to elaborate this point. *Impatiens noli-tangere* is rare in
53
54 the UK, but it is relatively common in the Lake District. Some observers might overlook this
55
56
57
58
59
60

1
2
3 species in favour of rarer or more distinctive ones. For the lepidopterologist, however,
4
5 *Impatiens noli-tangere* is a species of interest since it is the sole food source of the rare
6
7 moth *Eustroma reticulate* (*Netted carpet moth*). Within the discipline of Lepidoptero-
8
9 logy, then, places where *Impatiens noli-tangere* grew attracted particular attention.
10
11
12

13 Other witness groups including tourists, writers, travellers and poets approached the
14 region from yet different perspectives. Though the individuals affiliated with these groups
15 often made detailed observations, on the whole they were less interested in whether a
16 particular plant species belonged to a particular taxonomic group or habitat type than they
17 were in describing some facet of the landscape that surrounded them. Writers intent on
18 revealing the region's aesthetic qualities were, for example, more likely to focus their
19 attention on a plant or combination of plants that helped to define the landscape they were
20 trying to understand and describe. Consider Elizabeth Lynn Linton's *The Lake Country*.
21 Linton tells us her aim in writing this work was to 'illustrate and describe the most beautiful
22 places—both those popularly known, and those which only the residents ever find out.'
23 (Linton, 1864, p. x). Describing the walk up from Ambleside to the celebrated waterfall of
24 Stockghyll Force, it is the dramatic qualities of the locality that Linton attempts to impress
25 upon her reader. She writes:
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

following the wild path of rock and running water and twisted tree-roots—the rocks
below getting larger and more broken, the rift between them deeper and shaper—
the roar of the waters loader, and the rush more fierce and rapid ... there you come
upon the "loosening silver" of the fall (Linton, 1864, p. 18-19)

1
2
3
4 At this point, Linton uses a reference to the species *Pyrola media* to appeal to the
5
6 reader's imagination. She observes that, 'for those who have stout nerves ... the *Pyrola*
7
8 *media*, a rare growth of the Winter-green, [may be] found only among the rocks in the centre
9
10 of the fall' (Linton, 1864, p. 18-19). Linton's decision to include these observations indicates
11
12 how her interest in flora was intrinsically linked with her interest in the rugged landscape of
13
14 the region. The scarcity of the plant imbues the locality with a heightened sense of
15
16 uniqueness, while the difficulties involved in gaining a glimpse of it helps impress the
17
18 dramatic landscape of the waterfall.
19
20
21

22
23 Combining these sorts of discreet observations from texts that span the entire
24
25 timespan and geographic extent under investigation allows us to uncover a broader picture
26
27 of the region's flora as a whole. This picture becomes even richer and more nuanced when
28
29 observations from multiple witness groups, including guidebook writers, are included.
30
31
32
33

34 *2.3 Using digitised material and source selection*

35
36
37
38

39 The texts that form the corpus have been digitised using automated Optical Character
40
41 Recognition (OCR) software. This software uses NLP methodologies to recognise letters
42
43 and words from scanned images (Bennamoun & Mamic, 2012, p. 199-220). Automated OCR
44
45 digitisation is conventionally less expensive and time-consuming than manual digitisation,
46
47 however in digitising historical texts its accuracy has been shown to be variable, with a
48
49 percentage of characters and words being misidentified (Blanke, Bryant, & Hedges, 2012;
50
51 Tanner, Muñoz, & Ros, 2009). One way of dealing with OCR error is to use fuzzy matching.
52
53 Fuzzy matching relies upon the use of edit-distance to measure the similarity between two
54
55 words (Tanner et al., 2009). ~~This measure is frequently represented as a quantitative metric~~
56
57
58
59
60

1
2
3
4 ~~such as 0–100, where 100 represents an exact match.~~ Specifying a similarity threshold
5
6 allows for words that are not an exact match, due to OCR error, to be matched provided their
7
8 similarity falls above a specified quantitate threshold. However, for fuzzy matching to be
9
10 effective the threshold needs careful consideration. Setting a threshold that is too high is
11
12 unlikely to improve match rates significantly, while setting a threshold that is too low can
13
14 result in false positives. Achieving an effective balance is especially difficult when OCR
15
16 accuracy varies across a corpus, and for this reason it was not used in this study (Amelia,
17
18 2017; Gregory et al., 2016; Tanner et al., 2009).
19
20
21

22
23 Our decision to use OCR digitised texts, despite potential variations in accuracy, is
24
25 three-fold. Firstly, the manual digitisation of almost one hundred texts was beyond the
26
27 resources of this project. Secondly, a large amount of historical textual source material has
28
29 already been digitised at considerable expense to both the public and the private sector, and
30
31 this makes it desirable to explore the potential of this material and to assess its limitations
32
33 before investing in costly re-digitisation projects. Thirdly, we thought it important to determine
34
35 whether our developing methodology was robust enough to cope with ‘noisy’ corpora.
36
37
38

39
40 In the selection of texts, our intention was to form a corpus that reflected the broad
41
42 range of literature documenting Lake District flora. To achieve this several factors were
43
44 taken into consideration. These factors included the date of publication, geographical focus
45
46 and the text’s genre. Coverage of a broad time frame is essential if changes in plant
47
48 distribution are to be examined historically over time. The corpus therefore spans more than
49
50 two hundred years (1682–1904). Texts that focused on the Lake District or North West
51
52 England more generally were favoured, but texts covering a wider geographical area, such
53
54 as John Hull’s *The British Flora* (1799), were also admitted if they were deemed of special
55
56 significance to the recording of flora at a specific point in time.
57
58
59
60

1
2
3
4 Making selections based on genre proved especially challenging. As previously
5
6 noted, capturing the observations of different witness groups is essential for forming a more
7
8 comprehensive picture of plant distribution in the region. However, these observations were
9
10 often found to extend across several textual genres. Whereas the observations of the
11
12 national scientific elite are most commonly found in botanical floras, specialist botanical
13
14 journals and the reports and transactions of scientific societies, the observations of amateur
15
16 collectors and enthusiasts are more commonly found in published diaries and travelogues.
17
18 Consequently, we determined that it was necessary to include regional and national
19
20 botanical floras, botanical journals, scientific society reports and transactions as well as
21
22 periodicals, collected letters, diaries, botanical handbooks, travelogues and Lake District
23
24 tourist literature.
25
26
27
28
29
30
31

32 *2.4 Historic name variations: plant species synonyms*

33
34
35
36

37 In order to examine the distribution of flora recorded in a corpus of historical texts, it is
38
39 necessary to identify every plant name recorded in that corpus. This is made more
40
41 challenging as the scientific names of plants underwent substantial and sustained change
42
43 during our period of investigation. Many plant species have been recorded under different
44
45 names at different points in time and some plants have been recorded under several
46
47 different names at the same point in time. Changing plant taxonomy was one reason for
48
49 name variation. During the eighteenth and nineteenth centuries plant nomenclature was
50
51 considered intrinsically linked to taxonomic classification, with a plant's name serving to both
52
53 identify and distinguish it as a unique 'species', while at the same time linking that plant to
54
55 other plants with similar traits (Pickstone, 2001, p. 71; Sanderson, 2017). There was,
56
57
58
59
60

1
2
3 however, considerable debate over what traits should be considered the most important
4 basis for grouping plants together, and by the mid-eighteenth century several rival taxonomic
5 systems had been proposed, each setting out different naming principles as part of their
6 classification system (David Elliston Allen, 1994; Gledhill, 2002; Scharf, 2009).
7
8
9
10
11

12
13 Even within a single taxonomy, it was not uncommon for plant names to change over
14 time. The identification of new species collected from around the globe frequently required
15 plants to be re-grouped and re-named to incorporate new knowledge. ~~In these cases, plant~~
16 ~~names were sometimes also changed to reflect the re-grouping~~ (David Elliston Allen, 1994;
17 Scharf, 2009). A requirement of scientific naming systems is that each plant has a single
18 'accepted' name. However, it was frequently the case that plants collected and named in
19 different localities were later found to be in fact the same species. In these instances, one
20 name was selected as the accepted name and the other names were dropped. In both
21 situations the adoption of a single name was often a gradual process, and it was not
22 uncommon for naturalists in one social group or geographical region to continue using plant
23 names that had been rejected many years before by other groups or in other regions.
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

39 ~~Tracking historical changes in naming conventions across our corpus posed~~
40 ~~challenges for even state-of-the-art NLP methodologies. Conventionally, w~~When working
41 with modern corpora, one can employ NER techniques that identify and extract different
42 types of named entities by automatically comparing the contents of a corpus to a designated
43 search list. However, these lists normally use modern naming conventions and this makes
44 them less adept at detecting historic name variations (Butler, Donaldson, Taylor, & Gregory,
45 2017). In order to maintain the accuracy of NER techniques when working with historical
46 texts, it is necessary to develop naming inventories that map modern names to any known
47 historic variations. With this in mind, we compiled a list of historical plant synonym names
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 using the POWO database, which lists both currently accepted plant names alongside
4
5 documented alternative synonym names (POWO, 2018).
6
7

8 ~~In total POWO contains over 982,000 plant names. In spite its size though, the~~
9 ~~database is not an exhaustive list. Certain plant groups are still being evaluated and~~
10 ~~therefore are not yet available online (POWO, 2018). This limitation is not an issue for this~~
11 ~~study, however, as all the plants selected to form our search list were present in the POWO~~
12 ~~database. Furthermore, T~~the database primarily contains synonyms formed within the
13
14 binomial paradigm established by the Swedish naturalist Carl Linnaeus (1707–1778) during
15
16 the mid-eighteenth century. Before Linnaeus, naturalists had tended to use ‘diagnostic
17
18 phrase names’, which also served as a brief description of the plant and its traits (these
19
20 names could be very long, sometimes up to half a page in length; see, Gledhill, 2002;
21
22 Koerner, 1996, p. 149; Ogilvie, 2008; Reddy, 2007). As a result, most synonyms listed within
23
24 POWO date from the mid-eighteenth century onwards. Consequently, it is likely that POWO
25
26 will be less effective at identifying plant synonym names in texts published before the 1750s.
27
28 In order to determine the extent to which this was the case, we analysed the plant name
29
30 match rates per text to reveal if any decrease could be discerned across the corpus.
31
32
33
34
35
36
37
38
39
40
41
42
43

44 **3 Methodology**

45 *3.1 Forming the historical plant list*

46
47
48
49
50

51 In order to compile a historically sensitive plant search list, we first needed to draw together
52
53 a list of currently accepted plants relevant to the Lake District. We completed this task by
54
55 conducting a polygon query of the National Biodiversity Network (NBN Atlas) to find all plant
56
57 species listed as present in the Lake District (Atlas, 2019). The NBN Atlas is an open-source
58
59
60

1
2
3 online resource that provides access to extensive species records. It is linked to several
4
5 extensive databases, including the Botanical Society of Britain & Ireland (BSBI) and the
6
7 Biological Records Centre (BRC) (Atlas, 2019). We decided to focus on high-level species,
8
9 as these species are often more conspicuous in the landscape and are therefore more likely
10
11 to have been observed and recorded historically. As a result, all mosses, algae and liverwort
12
13 were filtered out of the downloaded dataset, leaving vascular plants and ferns. We then
14
15 compared this list with POWO, compiled the historical synonyms and ~~and~~ 'mapped' them to
16
17 their modern names. The final list used to search the corpus contains 952 currently accepted
18
19 plant species to which 9340 synonyms were linked (~~the plant search list, alongside the~~
20
21 ~~corpus and Python scripts have all be made available on GitHub and University online~~
22
23 ~~repository -- URL~~ reference anonymised for review). Within the compiled list a number of
24
25 historical synonyms were linked to two or more currently accepted names. As this could
26
27 cause potential ambiguity in the analysis of the results, duplicate names were filtered out if
28
29 matched in the corpus.
30
31
32
33
34
35
36
37
38
39

40 *3.2 Forming and formatting the corpus*

41
42
43
44 The corpus was formed entirely of existing digitised texts downloaded from open-source
45
46 online repositories including Archive.org, Google Books and BioDiversity.org (~~URL~~ reference
47
48 anonymised for review). Once all the texts had been downloaded, the first stage was to
49
50 standardise and format the corpus. Non UTF-8 characters were removed, and words split
51
52 over two lines were re-joined. Both processes were preformed using an automated approach
53
54 with regular expressions in Python 2.7 (~~URL~~ reference anonymised for review).
55
56
57
58
59
60

3.3 Extracting species and geographical locations from the corpus

To improve efficiency, we located plant species and place names within the corpus in two phases. Firstly, we used a keyword search method to identify plant species named in the corpus, including synonym variations and abbreviated forms. For certain search terms, particularly those where a word can have multiple meanings, keyword searching can lead to misleading results. However, the complex structure of plant names (which are formed of two Latin words) and their appearance in texts principally written in English, helped to reduce ambiguity between plant names and other word tokens. By manually checking a random sample of 500 match instances, we found that only 1.8% of the sample was misidentified. To support subsequent analyses of the data, all plant name match instances were then mapped to their current accepted name and a span of text (or co-text) adjacent to each match was extracted. The process was performed using an automated approach, with ~~the~~ scripts written and run in Python 2.7. The extracted text was then geoparsed using the Edinburgh Geoparser (Grover et al., 2010; Tobin, Grover, Byrne, Reid, & Walsh, 2010). This process enabled the automated identification and georeferencing of place name entities across text, which is a prerequisite for performing geospatial analysis using GIS software. Geoparsing the extracted co-text, and not the whole corpus, improves efficiency as only the relevant sections are examined (Rupp et al., 2015). This is especially valuable when working with the more general natural history texts, which are likely to contain numerous place names related to other research areas in the environmental sciences.

Place names can have a variety of potential meanings, and they are thus more likely to be affected by issues of ambiguity than plant names. As Rupp et. al. have noted, the place name Lancaster can be a town, but it can also be the name of a person (Stuart

1
2
3 Lancaster) or even an honorary title (the Duke of Lancaster). Furthermore, a single place
4 name can refer to multiple localities: there are settlements named Lancaster in England, the
5 USA and South Africa (~~Donaldson et al., 2017, p. 47~~; Rupp et al., 2015). To mitigate the
6 potential errors introduced by these ambiguities, the Edinburgh Geoparser includes two
7 interlinked components: a 'geo-tagger' and a 'geo-resolver'. The geo-tagger runs through a
8 sequence of NLP analysis steps (including tokenization, sentence splitting, POS tagging,
9 chunking and a rule-based named entity recogniser) to identify place names within the text
10 and to disambiguate them from other word token types (Grover et al., 2010; Tobin et al.,
11 2010). The geo-resolver then attempts to assign a pair of geographical coordinates to each
12 identified place name, using a ranking algorithm that considers population size and the
13 geospatial relation of the place to others in the document. This algorithm gives preference to
14 places that cluster with other locations in the same document (Grover et al., 2010; Tobin et
15 al., 2010). Assessing the impact of these heuristics on the output can be difficult. However,
16 manual checking of the output indicates that the ranking of locations in our corpus was more
17 frequently influenced by the geographical clustering of place names than by population size,
18 which reflects the regional focus of many of the texts in the corpus. In addition to these
19 metrics, it is possible to add a weighting to a particular geographical area when the
20 geographical locality under investigation is already known (Grover et al., 2010; Tobin et al.,
21 2010). For this study, the geoparser was used in conjunction with the Ordnance Survey
22 1:50,000 scale gazetteer, as this gazetteer provided the best coverage for the Lake District,
23 and a weighting of two was added for locations which fell inside the North West of England
24 ("OS Open Data," n.d.).

3.4 Determining plant and location collocates

1
2
3
4
5
6 Collocate analysis was used to trace the recorded localities of plants across the corpus. An
7
8 established analytical method within the fields of lexicography and corpus linguistics,
9
10 collocate analysis can be used to identify automatically every time a pair of co-occurring
11
12 search terms appear in close proximity to one another within a text. As a collocate is here
13
14 determined by the proximity between terms within a text, the span of the collocate 'window'
15
16 is of some significance: if the span is too narrow, place names associated with the search
17
18 term may be missed; if the span is too wide, place names are more likely to be erroneously
19
20 linked to the search term. Commonly, a collocate window is of a fixed size. In this study,
21
22 however, our initial experiments using a fixed size collocate window led to place names
23
24 being erroneously linked to plant names. Manual checking of these errors suggested that
25
26 they were principally a result of the very 'compact' format of the entries in many botanical
27
28 works, which frequently list information about flora and their geographical locations in a
29
30 compressed space (an example of which is shown in Fig. 1). In these instances, we found
31
32 that a fixed window would extend across multiple plant species match instances, leading to
33
34 false positives between plant and place. To overcome this problem, we adopted a 'dynamic'
35
36 collocate window. This was done by identifying all plant species matches and their
37
38 abbreviated forms across the whole corpus and setting a collocate boundary that extended
39
40 either up to 300 characters to the right of each match or up to the point where the next plant
41
42 species match occurred.
43
44
45
46
47
48
49

50
51 The direction of the span of the collocate window was also taken into consideration.
52
53 A span can be made to the left or right of each plant name in order to include the words that
54
55 co-occur immediately before and after each match. Close evaluation of the initial output of
56
57 the collocate analysis revealed that including a span to the left of the matched plant names
58
59
60

1
2
3 increased the number of false positives. Inspection of these false positives suggests that
4
5 they were a consequence of the compact format of many of the texts in the corpus and of
6
7 the tendency of the texts in the corpus to record the location where a species was observed
8
9 after naming the species itself. Accordingly, we decided to extend the collocate window only
10
11 to the right of each match instance, which improved the overall accuracy of the results. This
12
13 decision means that our results do not include instances where a location is mentioned
14
15 before a plant name. Given the early stage of this research we felt that prioritising the
16
17 accuracy of the results over the recall was justified, as it gives more weight to the potential
18
19 value of the generated dataset. We are confident that it should be possible to augment our
20
21 analyses once a methodology for accurately identifying collocate pairings with place names
22
23 mentioned before match instances has been developed.
24
25
26
27
28
29
30
31

32 **4 Results and Findings**

33 *4.1 The impact of historical synonyms on match instances*

34
35
36
37
38
39 Before we proceed to examine the geographies of historical Lake District flora that the
40
41 collocate analysis revealed, it is helpful first to consider the plant species matches and how
42
43 searching for historical synonyms influenced the results. Searching the corpus for modern
44
45 accepted plant names resulted in a total of 16216 match instances. The number of match
46
47 instances increased to 22659 when the search was expanded to include both accepted plant
48
49 names and their historical synonyms. We can account for this marked increase of 28% in
50
51 two ways. Firstly, the number of times a species was matched in the corpus can be seen to
52
53 have increased. When we searched the corpus for only modern accepted plants names, 673
54
55 species were matched at an average of 24.1 match instances per species. When synonyms
56
57
58
59
60

1
2
3 were included in the search, the matches for these 673 species increased to 28.8 match
4 instances per species. Secondly, the number of different species being matched also
5
6 increased.
7
8
9

10 As just noted, 673 of the 952 plant species listed in the search list were matched in
11 the corpus when we searched for modern accepted names. The total increased to 802
12 species matches when synonyms were added to the search list. This increase suggests that
13 many modern accepted species names were not in use during the timeframe under
14 investigation. Although further fine-grained analysis would be required to reveal if there are
15 any patterns regarding plant name changes, these results do demonstrate the improvements
16 in match instances that can be gained by using temporally sensitive search lists when
17 examining flora across historical source material. Through the results we can see an
18 increase not only in the number of times a plant species was matched across the corpus, but
19 also in the number of different unique species matched. More complete identification of plant
20 species across the whole corpus is critical if subsequent collocate analysis and geoparsing
21 are to provide an accurate picture of the geographies of historical Lake District flora.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 *4.2 The impact of plant recording practices on match rate*

43
44
45

46 The identification of historical synonyms can also assist in the investigation of observer
47 practices. ~~Plant synonyms account for 6494 of the total match instances.~~ Fig. 2, created
48 using the pyplotlib library in Python 2.7, shows the ratio between the number of times a plant
49 species was matched in the corpus under a historical synonym name and the number of
50 times it was matched under its modern accepted name. As is evident in Fig. 2, just over half
51 of the plants matched across the corpus were recorded under a single name: 367 were
52
53
54
55
56
57
58
59
60

1
2
3 recorded under only their modern accepted name, whereas 129 were recorded under a
4
5 single historical synonym name. This finding indicates that the use of many plant names was
6
7 relatively stable across the time period under investigation. This result surprises, especially
8
9 given the heterogeneity of the corpus and the span of time it represents.
10
11
12

13 For plants matched under two or more names, the transition from one name to
14
15 another can be used to explore evolving naming conventions and, more broadly, changes in
16
17 recorder habits. On the one hand, a smooth shift from one name variation to another might
18
19 tend to suggest a smooth transition in the modification of plant-naming habits, with a new
20
21 name being suggested and adopted by the botanical community over a relatively short span
22
23 of time. On the other hand, the use of several variant plant names over a sustained period
24
25 suggests a more staggered transition between names, either because some observers were
26
27 not aware a new name had been proposed or because they resisted adopting the new
28
29 name. Both trends can be detected in the 306 species recorded under two or more names.
30
31
32
33

34 Consider the species *Platanthera bifolia* (*lesser butterfly-orchid*) and *Ranunculus*
35
36 *aquatilis* (*Water Crowfoot*). Both are documented in modern floras as having taxonomies that
37
38 remained unresolved until the second half of the twentieth century (Preston, Pearman,
39
40 Dines, & others, 2002; Stace, 2010). Our findings extend this understanding, revealing
41
42 distinct differences in the trajectory of naming conventions for these two species over the
43
44 past two hundred years. *Platanthera bifolia* was matched 67 times across 34 texts, spanning
45
46 the entire timeframe of the corpus. The earliest text in which the species was recorded was
47
48 Volume 2 of Ray's *Historia Plantarum* (1688); the latest text was Crump and Crossland's
49
50 *Flora of the Parish of Halifax* (1904). In total, the species was recorded under five different
51
52 names: *Orchis bifolia*, *Orchis alba*, *Habenaria bifolia*, *Platanthera bifolia* and *Gymnadenia*
53
54 *bifolia*. Despite being recorded under a variety of different names, there is a discernible
55
56
57
58
59
60

1
2
3
4 pattern in the use of these names over time. From the 1680s until around the mid-eighteenth
5
6 century the plant was matched under the names *Orchis bifolia* and *Orchis alba*. After this
7
8 period, *Orchis alba* no longer appears in the matches and the species is instead most
9
10 commonly recorded under the name *Orchis bifolia*, as well as *Platanthera bifolia*,
11
12 *Gymnadenia bifolia* and *Habenaria bifolia*. This last name was first matched in Volume 4 of
13
14 Smith's *English Flora* (1828) and appears with increasing regularity until the mid-nineteenth
15
16 century, after which point it becomes the only name under which the species is recorded.
17
18
19

20
21 With *Ranunculus aquatilis* the situation is rather different. This species was matched
22
23 71 times across 31 texts in the corpus. As with *Platanthera bifolia*, the match instances for
24
25 *Ranunculus aquatilis* occurred across nearly the entire timeframe of the corpus, with the
26
27 earliest text being Volume 1 of Ray's *Historia Plantarum* (1686) and the latest text being an
28
29 edition of *The Naturalist* from 1893. Over the course of the 207 years that separate these
30
31 texts, *Ranunculus aquatilis* was recorded under six different synonyms: *Ranunculus*
32
33 *aquatilis*, *Ranunculus aquaticus*, *Batrachium heterophyllum*, *Ranunculus heterophyllus*,
34
35 *Ranunculus hydrocharis* and *Ranunculus diversifolius*. However, the transition between
36
37 these synonyms is not as smooth as with *Platanthera bifolia*. Instead, the synonyms under
38
39 which *Ranunculus aquatilis* is recorded appear to have been used interchangeably. The
40
41 name *Ranunculus aquatilis* was first matched in 1686 and last matched in 1890. Similarly,
42
43 the first match for the synonym *Ranunculus aquaticus* was in 1686 and the latest was in
44
45 1870. The other synonyms (~~*Batrachium heterophyllum*, *Ranunculus heterophyllus*,~~
46
47 ~~*Ranunculus hydrocharis* and *Ranunculus diversifolius*~~) were all matched over the course of
48
49 the nineteenth century.
50
51
52
53
54
55

56
57 These findings provide greater insight into the different ways naming conventions
58
59 evolved for these two species. Specifically, the findings indicate that new names proposed
60

1
2
3 for *Platanthera bifolia* were accepted and adopted swiftly, whereas newly proposed names
4
5 for *Ranunculus aquatilis* were taken up much more slowly, if at all. Identifying these trends
6
7 provides opportunities for further analysis into the factors that influenced plant-naming
8
9 conventions in the region. More immediately, these trends remind us that the texts in our
10
11 corpus reflect historical scientific conventions as well as biases, and that these conventions
12
13 and biases should be taken into account in the methods implemented in analysing the
14
15
16
17
18 corpus.

21 22 23 *4.3 Potential limitations of the computational methodology*

24
25
26
27 Despite the apparent improvement in match instances when using historically sensitive
28
29 search lists, one must still be cautious when analysing the results. Above all, one must be
30
31 mindful that some historical name variations may still be missed, either as a result of spelling
32
33 errors introduced during the OCR process or because the historical search list is still
34
35 incomplete. Assessing where computational methods have failed to identify plant names is
36
37 very challenging and requires further checks to be performed. One approach is to plot the
38
39 results across the corpus and assess changes to the number of match instances.
40
41
42

43
44 Fig. 3, created using the pyplotlib library in Python 2.7, plots the match instances of
45
46 accepted names (X) and accepted and synonym names combined (O). As already
47
48 established, match instances are noticeably increased when historical synonyms are
49
50 included in the search. However, Fig. 3 reveals a degree of variability in match rate per text
51
52 across the corpus, with a greater number of matches being discernible between 1820 and
53
54 1890. Before 1800 the number of match instances per text decreases sharply. Given that
55
56 these texts were selected on account of their relevance to the recording of Lake District flora,
57
58
59
60

1
2
3
4 it seems unlikely that they do not contain any plant names. A more plausible explanation is
5
6 that the plant names used in earlier texts in the corpus are not being detected. Using Fig. 3
7
8 as a guide, a detailed reading of the texts published around 1800 was performed to
9
10 understand why there was such a sharp drop-off in results. This evaluation revealed that
11
12 texts published before 1800 increasingly used non-Linnean taxonomic systems, such as
13
14 those of John Gerard (1545–1612), Caspar Bauhin (1560, 1624), Joseph Pitton de
15
16 Tournefort (1656–1708) and above all the English naturalist John Ray (1627–1705) (Charles,
17
18 1947; Ogilvie, 2008; Scharf, 2009). As a result, the names used to record plants in many of
19
20 the texts published before 1800 predate the names compiled from POWO and were
21
22 therefore being missed. Such omissions are noteworthy as they give us a better
23
24 understanding of the adoption of Linnaean plants names in Britain.
25
26
27
28
29

30 A further disparity in the results was uncovered when the distribution of match
31
32 instances, presented in Fig. 2, were considered in relation to witness groups. Here, match
33
34 rates across travel accounts were found to be frequently lower than scientific floras and
35
36 journals, as well as more general histories of the region. Closer inspection revealed that
37
38 tourists and travellers to the region frequently documented the plants they observed using
39
40 common names rather than scientific binomial names. Again, these names were not
41
42 compiled from POWO and had therefore been missed. These findings are helpful in
43
44 revealing not only how observing conventions varied over time and between social groups,
45
46 but also where and why our methodology has failed to identify plant names. Such issues
47
48 alert us to the allowances that need to be made in future analyses of the results. These
49
50 findings also provide insight into how our methodology can be refined and improved going
51
52 forwards. They highlight that the addition of pre-Linnean plant names and common plant
53
54 names into our plant species search list would likely capture further information.
55
56
57
58
59
60

4.4 Mapping extracted information

We shall now examine how searching for historical synonyms across the corpus affects the geographical distribution of the geoparsed results. The collocate results for modern accepted names returned a total of 1982 location matches, 515 of which were unique. This figure increased to 2569 total location matches and 576 unique locations when the modern accepted plant name collocates were combined with the synonym plant name collocates. Closer investigation of the geoparsed results revealed a discernible increase in both the composition of plant species that were geolocated and the locations to which they were linked. For the modern accepted plant name collocates, 400 individual plant species were linked to 515 different Lake District locations. This number increased to 510 plant species and 576 locations when the plant synonym collocates are added.

Assigning geographical coordinates to each location enables the results to be visualised. Fig. 4, created using ArcGIS 10, plots the results of the geoparsed modern plant name collocates against the combined modern and synonym collocates. This visualisation further exposes the geospatial differences that result from implementing the two search methodologies. Specifically, Fig. 4 helps reveal how both the *extent* (or, in other words, the geospatial distribution) and the *depth* (in other words, the number of plants linked to each location) improved when historical synonyms were included in the search list. Each dot on the map marks a location with which a plant species was collocated, with the size of the dot representing the number of times a plant was matched to that location. When one focuses on the extent of the geospatial dispersion of the results across the whole of the Lake District, the geoparsed modern name collocates and the geoparsed modern and synonym collocates

1
2
3 appears to be similar. However, subtle shifts can be discerned when one focuses in on
4
5 specific localities, such as the northwest of the Lake District. Here, place names including
6
7 Cardurnock, Blackdyke, Wampool, Thurstonfield, Howrigg and Harker can all be seen to
8
9 match with synonym names only. These place names stand out from other place name
10
11 collocates match instances as they all received comparatively less attention across the
12
13 corpus. Consequently, the impact of searching for historical synonyms alongside modern
14
15 plant names is accentuated for these place names when the findings are mapped
16
17 geospatially, as other place names with a higher number of collocates-match instances are
18
19 more likely to be collocated with modern plant names.
20
21
22
23
24

25 Linked with the extent of place names distribution, the depth of plant names
26
27 collocated to each location also improved when historical synonyms were included in the
28
29 search list. As can be seen in Fig. 4, the number of times plant names were linked to a
30
31 particular site visibly increases in several localities across the region. This increase is
32
33 especially evident around Buttermere and Derwent Water and around Coniston, Langdale
34
35 and Hawkshead. This provides a more complete picture regarding the species composition
36
37 in each locality; it also provides further insights into observer habits and the range of
38
39 different people making these observations.
40
41
42
43

44 Evaluation of the impact of searching for historical synonyms on the match instances
45
46 of individual plant species revealed that those plant species matched most frequently under
47
48 synonym names had the most pronounced shifts in the geographical locations to which they
49
50 were linked. For example, the species *Blechnum spicant* (*hard fern*) was matched 101 times
51
52 across the corpus and 90 of these instances were under synonym names. During
53
54 geoparsing, the species was collocated with 19 locations around the Lake District including
55
56 Keswick, Buttermere, Grisedale Pike, Wasdale, Scale Force, Kirkstone Pass, Mosser,
57
58
59
60

1
2
3 Ullock, Lamplugh, St. Bees Head, Glaramara, Styhead, Lingmell, Ponsonby, Bootle, Santon
4
5 Bridge, Birks Wood, Ease Gill and Millom. Of all these locations, only the collocations with
6
7 Bootle were matches with the modern accepted name of the species. The case of *Littorella*
8
9 *uniflora* (*Shoreweed*) is even more extreme. This species was matched in the corpus 62
10
11 times, but only under the synonym *Littorella lacustris*. As a result, its collocation with
12
13 Edenhall, Lodore, Blea Tarn, Styhead Tarn and Barrow-in-Furness would have been missed
14
15 if we had only searched for its modern name. As these examples demonstrate, using
16
17 modern plant species lists alone is likely to distort the results, giving an increased weighting
18
19 to plant species whose modern accepted names remained more consistent throughout the
20
21 period under investigation. Going forward, it will be necessary to assess the results of such
22
23 analyses in greater detail in order to determine where any plants species have been
24
25 incorrectly collocated to locations. This could be accomplished by reading a sample of the
26
27 collocate windows to establish instances where plant names and place names have been
28
29 misidentified. A further step would be to apply statistical measurements, such as Kulldorff's
30
31 spatial scan statistic, to distinguish the degree to which the collocate pairings have been
32
33 influenced by the underlying geography of the corpus (Kulldorff, 1997; Rupp et al., 2015).
34
35 Using this sort of approach might reveal which collocates form clusters even when the
36
37 underlying geographies of the text are taken into account (Rupp et al., 2014). This would
38
39 provide confidence in the results and pave the way for further analysis of the data though the
40
41 clustering of indicator species for different habitat types to assess changes in the observed
42
43 environment over time. These clusters could in turn be used to help determine whether the
44
45 discerned changes have been caused by human or natural events.
46
47
48
49
50
51
52
53
54
55
56
57

58 **5 Conclusions**

59
60

1
2
3
4
5
6 This article has introduced a methodology that uses computational techniques to extract the
7
8 geospatial information of Lake District flora from a corpus of disparate historical texts. Our
9
10 findings reveal the potential of using historical sources to provide detailed empirical evidence
11
12 regarding the historical distribution of different plant species across the Lake District. In total,
13
14 802 species were traced across the corpus and 510 of these species could be linked to
15
16 locations around the Lake District. Re-organising the extracted information into a structured
17
18 and searchable geo-temporal dataset enables further analysis; it also allows for the
19
20 formation of advanced queries that use resources such as GIS to visualise and examine
21
22 changes to plant species composition and distribution over time.
23
24
25
26

27 Our findings demonstrate the importance of developing temporally sensitive search
28
29 lists in order to trace plant names accurately across a corpus that comprises texts from an
30
31 array of genres and historical periods. Searching for historical synonym plant names
32
33 alongside modern accepted names improved the number of unique plant species found
34
35 across the corpus as well as the total number of finds. Focusing on the geoparsed results,
36
37 the inclusion of plant name synonyms increased both the number of species that could be
38
39 geolocated and the number of different locations to which they were geolocated.
40
41
42
43
44 Consequently, our temporally sensitive methodology can be seen to have resulted in a more
45
46 accurate reflection of the consulted corpus and to have reduced distortions that would have
47
48 arisen were we to have only searched for modern accepted names.
49
50

51 But the inclusion of historical plant synonyms does more than simply improve the
52
53 accuracy of the results. It also brings broader changes in plant classification systems and
54
55 naming conventions into focus. The results indicate that many plant names are more stable
56
57 than might initially be expected: 496 of the 802 plants found in the corpus matched under a
58
59
60

1
2
3
4 single name. It is important to bear in mind, however, that only 367 of the species that were
5
6 matched under a single name were recorded under their modern accepted name. If the
7
8 timeframe under investigation was extended farther into the twentieth century, it is likely that
9
10 greater shifts in plant names would be observed. Closer investigation of the results revealed
11
12 earlier texts were more likely to use pre-Linnean plant names, whereas travellers and
13
14 tourists were more likely to use common plant names. Including these names in our search
15
16 plant list would not only extend the dataset, but also reveal further shifts in naming
17
18 conventions.
19
20
21

22
23 For the 306 plants species matched under two or more names, the transition from
24
25 one name to the next across the corpus provides insight into the introduction and adoption of
26
27 new names by those making observations on Lake District flora. In this case, some plant
28
29 names were replaced by another relatively swiftly, which indicates that the new name was
30
31 accepted and adopted with little resistance. Other names, however, appear to have been
32
33 adopted less uniformly, with multiple plant names continuing to be used for a sustained
34
35 period. These findings open up new lines of research into changing naming practices of
36
37 Lake District flora and why some plant names were adopted more swiftly than others.
38
39
40
41

42
43 Even though searching for historical synonyms alongside modern accepted names
44
45 improved results, other factors should not be ignored as a decline in match instances is still
46
47 discernible in the results. In particular, the lower match rate before the year 1800 indicates
48
49 that the Linnaean binomial naming system was only adopted in Britain at the turn of the
50
51 nineteenth century. One option to overcome this limitation would be to extend the plant-
52
53 naming lists by adding pre-Linnaean plant names. Furthermore, errors introduced by OCR
54
55 also resulted in some plant and place names collocations being missed. A possible way to
56
57
58
59
60

1
2
3 address such errors in further investigations of the corpus would be to draw upon emerging
4
5
6 Machine Learning approaches to help identify plant and place names across the corpus.
7

8 Notwithstanding these limitations, the findings presented in this paper provides
9
10 researchers with a firmer empirical grounding for making accurate assessments of historical
11
12 flora. The structured digital dataset we have created from an otherwise disparate
13
14 assemblage of texts is a composite knowledge base which can be used to assess broader
15
16 floristic changes and to guide and complement more detailed analysis of specific sources. In
17
18 its digital form the dataset can be easily shared, manipulated, queried and visualised
19
20 geospatially; it can therefore be explored and exploited by a range of audiences to suit
21
22 different interests and needs. These audiences include: academic researchers such as
23
24 environmental historians and scientists, who have previously struggled to interrogate a wide
25
26 range of historical accounts; heritage and conservation organisations including members of
27
28 the LDWHSP, who are responsible for protecting and preserving the region; and members of
29
30 the general public, who want to know more about the history of the Lake District's
31
32 environment. It is hoped, moreover, that our computational methodology might be extended
33
34 beyond the Lake District, taking in an even broader range of digitised historical texts to
35
36 support research into the environmental history of other localities and natural features.
37
38 Collectively, such research could yield a much more detailed and comprehensive
39
40 understanding of the earth's environmental past.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

1
2
3 **Figure Captions**
4
5
6
7

8 Fig. 1. Excerpt from Winch's *Flora of the Lake District* (Winch, 1825, p. 27)
9

10 Fig. 2. Ratio (%) of plant species being matched in the corpus under synonym names
11
12

13 Fig. 3. Match rates across the corpus
14

15 Fig. 4 Geographical distribution of plant species matched from across the corpus
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

- Allen, David E. (2003). Four centuries of local flora-writing: some milestones. *WATSONIA-KINGS LYNN-BOTANICAL SOCIETY OF THE BRITISH ISLES*-, 24(3), 271–280.
- Allen, David Elliston. (1994). *The naturalist in Britain: a social history*. Princeton University Press.
- Amelia, T. J. (2017). *Corpus Linguistics for History: The Methodology of investigating place-name discourses in digitised nineteenth-century newspapers*. Lancaster University.
Retrieved from <https://doi.org/10.17635/lancaster/thesis/143>
- Arber, A. (1943). A seventeenth-century naturalist: John Ray. *Isis*, 34(4), 319–324.
- Atlas, N. B. N. (2019). No Title. Retrieved from <https://registry.nbnatlas.org/>
- Baker, J. G. (1885). *A flora of the English Lake District*. George Bell.
- Bennamoun, M., & Mamic, G. J. (2012). *Object recognition: fundamentals and case studies*. Springer Science & Business Media.

- 1
2
3
4 Blanke, T., Bryant, M., & Hedges, M. (2012). Open source optical character recognition for
5
6 historical research. *Journal of Documentation*, 68(5), 659–683.
7
8 <https://doi.org/10.1108/00220411211256021>
9
- 10 Butler, J. O., Donaldson, C. E., Taylor, J. E., & Gregory, I. N. (2017). Alts, Abbreviations, and
11
12 AKAs: historical onomastic variation and automated named entity recognition. *Journal*
13
14 *of Map & Geography Libraries*, 13(1), 58–81.
15
16
- 17 Cervera, T., Pino, J., Marull, J., Padró, R., & Tello, E. (2019). Understanding the long-term
18
19 dynamics of forest transition: From deforestation to afforestation in a Mediterranean
20
21 landscape (Catalonia, 1868–2005). *Land Use Policy*, 80, 318–331.
22
23
- 24 Charles, E. (1947). *Raven, English Naturalists from Neckam to Ray*. Cambridge: Cambridge
25
26 University Press.
27
28
- 29 Denyer, S. (2016). The Lake District Landscape: Cultural or Natural? In *The Making of a*
30
31 *Cultural Landscape* (pp. 19–46). Routledge.
32
33
- 34 Desmond, R. (1977). *Dictionary of British and Irish Botanists and Horticulturalists Including*
35
36 *Plant Collectors, Flower Painters and Garden Designers*. CRC Press.
37
38
- 39 Donaldson, C., Gregory, I. N., & Taylor, J. E. (2017). Locating the beautiful, picturesque, sublime and
40
41 majestic: spatially analysing the application of aesthetic terminology in descriptions of the
42
43 English Lake District. *Journal of Historical Geography*, 56, 43–60.
44
45 <https://doi.org/10.1016/j.jhg.2017.01.006>
46
- 47 Ellenberg, H. (1974) Zeigerwerte der Gefäßpflanzen Mitteleuropas. *Scripta Geobotanica*, 9.
48
49 Göttingen.
- 50 Gledhill, D. (2002). *The names of plants*. Cambridge University Press.
- 51
52 Gregory, I., Atkinson, P., Hardie, A., Joulain-Jay, A., Kershaw, D., Porter, C., ... Rupp, C.
53
54 (2016). From Digital Resources to Historical Scholarship with the British Library 19th
55
56 Century Newspaper Collection. *Journal of Siberian Federal University. Humanities &*
57
58 *Social Sciences*, 9(04), 994–1006. <https://doi.org/10.17516/1997-1370-2016-9-4-994->
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1006

- Gregory, I., & Donaldson, C. (2016). Geographical text analysis: Digital cartographies of Lake District literature. In *Literary mapping in the digital age* (pp. 85–105). Routledge.
- Grover, C., Tobin, R., Byrne, K., Woollard, M., Reid, J., Dunn, S., & Ball, J. (2010). Use of the {Edinburgh Geoparser} for Georeferencing Digitised Historical Collections. *Philosophical Transactions of the Royal Society A*, *368*(1925), 3875–3889.
- Hinrichs, U., Alex, B., Clifford, J., Watson, A., Quigley, A., Klein, E., & Coates, C. M. (2015). Trading consequences: A case study of combining text mining and visualization to facilitate document exploration. *Digital Scholarship in the Humanities*, *30*(suppl_1), i50–i75.
- Hodgson, W., & Goodchild, J. G. (1898). *Flora of Cumberland*. W. Meals and Company.
- Kheraj, S., & Oosthoek, K. J. (2016). 18 Online digital communication, networking, and environmental history. *Methodological Challenges in Nature-Culture and Environmental History Research*, 233.
- Koerner, L. (1996). Carl Linnaeus in his time and place. *Cultures of Natural History*, 145–162.
- Kuklick, H., & Kohler, R. E. (1996). Science in the field, Osiris, 2nd Series, Vol. 11.
- Kulldorff, M. (1997). A spatial scan statistic. *Communications in Statistics - Theory and Methods*, *26*(6), 1481–1496. <https://doi.org/10.1080/03610929708831995>
- Lake District National Park. (n.d.). Lake District Nomination and Partnership Plan. Retrieved from <http://www.lakedistrict.gov.uk/caringfor/projects/whs/lake-district-nomination>
- Lindop, G. (2005). *A literary guide to the Lake District*. Sigma Press.
- Linton, E. L. (1864). *The lake country*. Smith, Elder and Company.
- Malleson, F. A. (1890). *Holiday Studies of Wordsworth by Rivers, Woods, and Alps: The*

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Wharfe, the Duddon, and the Stelvio Pass. Cassell and Company.

Nicholas Jardine, James A. Secord, E. C. S. (Ed.). (1996). *Cultures of natural history*.
Cambridge University Press.

Nicholson, N. (1955). *The Lakers: the adventures of the first tourists*. Hale.

Ogilvie, B. W. (2008). *The science of describing: Natural history in Renaissance Europe*.
University of Chicago Press.

OS Open Data. (n.d.). Retrieved from

<https://www.ordnancesurvey.co.uk/opendatadownload/>

Pickstone, J. V. (2001). *Ways of knowing: A new history of science, technology, and
medicine*. University of Manchester Press.

POWO. (2018). Plants of the World Online. Retrieved from

<http://www.plantsoftheworldonline.org/>

Preston, C. D., Pearman, D., Dines, T. D., & others. (2002). *New atlas of the British & Irish
flora*. Oxford University Press.

Reddy, S. M. (2007). *University Botany-iii:(Plant Taxonomy, Plant Embryology, Plant
Physiology)* (Vol. 3). New Age International.

Rich, T. C. G. (1997). Is ad hoc good enough. *Transactions of the Suffolk Naturalists'
Society*, 33, 14–21.

Rupp, C. J., Rayson, P., Gregory, I., Hardie, A., Joulain, A., & Hartmann, D. (2014). Dealing
with heterogeneous big data when geoparsing historical corpora. In *2014 IEEE
International Conference on Big Data (Big Data)* (pp. 80–83). IEEE.

<https://doi.org/10.1109/BigData.2014.7004457>

Rupp, C. J., Rayson, P., Gregory, I., Hardie, A., Joulain, A., & Hartmann, D. (2015). Dealing
with heterogeneous big data when geoparsing historical corpora. *Proceedings - 2014*

1
2
3
4 *IEEE International Conference on Big Data, IEEE Big Data 2014*, (September 2013),
5
6 80–83. <https://doi.org/10.1109/BigData.2014.7004457>
7

8 Sanderson, J. (2017). *Plants, people and practices: the nature and history of the UPOV*
9
10 *convention*. Cambridge University Press.
11

12
13 Scharf, S. T. (2009). Identification keys, the “Natural Method,” and the development of plant
14
15 identification manuals. *Journal of the History of Biology*, 42(1), 73–117.
16

17
18 Stace, C. (2010). *New flora of the British Isles*. Cambridge University Press.
19

20
21 Tanner, S., Muñoz, T., & Ros, P. H. (2009). Measuring mass text digitization quality and
22
23 usefulness. *D-Lib Magazine*, 15(7/8), 1082–9873.
24

25
26 Tobin, R., Grover, C., Byrne, K., Reid, J., & Walsh, J. (2010). Evaluation of georeferencing.
27
28 In *Proceedings of the 6th Workshop on Geographic Information Retrieval - GIR '10* (p.
29
30 1). New York, New York, USA: ACM Press. <https://doi.org/10.1145/1722080.1722089>
31

32
33 UNESCO. (2017). The Operational Guidelines for the Implementation of the World Heritage
34
35 Convention. Retrieved August 20, 2009, from <https://whc.unesco.org/en/guidelines>
36

37
38 Whittaker, E. J. (Ed.). (1981). *A Seventeenth Century Flora of Cumbria: William Nicolson's*
39
40 *Catalogue of Plants 1690*. Publications of the Surtees Society,.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Acknowledgements

We are grateful to Mike Bennett and Dr. Beatrice Alex at Edinburgh for their support with the geoparser and the Unlock service. We are also grateful to John Iacona at The Royal Botanical Gardens Kew, for help and guidance in accessing POWO data. Further thanks are due to Prof. Ian Gregory and Dr. Patricia Murrieta-Flores for their insight and guidance.

For Peer Review

Funding

This work was supported by the EU-China Research and Innovations Partnership (ECRIP) [2014/348-010], as part of the project Addressing food Security, Environmental stress and Water by promoting multidisciplinary Research EU And China Partnerships in science and business (Sew reap); the Leverhulme Trust [RPG-2015-230], as part of the Geospatial Innovation in the Digital Humanities project; and the Arts and Humanities Research Council (AHRC) [AH/S012893/1], via the North West Consortium Doctoral Training Partnership (NWCDTP), as part of the project Augmented Humanity: Does the Human Enhance the Machine or the Machine Enhance the Human?

Appendix

List of texts in corpus

Author	Title	Year of Publication	Volume Number
Ray, John	<i>Historia Plantarum</i>	1686	1
Ray, John	<i>Historia Plantarum</i>	1688	2
Ray, John	<i>Synopsis methodica stirpium Britannicarum</i>	1690	
Ray, John	<i>Historia Plantarum</i>	1704	3
Wilson, John	<i>A synopsis of British Plants</i>	1744	
Hudson, William	<i>Flora Anglica</i>	1762	
Nicholson and Burns	<i>History and Antiquities of the Counties of Westmorland and Cumberland</i>	1777	1
Nicholson and Burns	<i>History and Antiquities of the Counties of Westmorland and Cumberland</i>	1777	2
	<i>Memoirs and proceedings of the Manchester Literary and Philosophical Society</i>	1789	1
Pulteney, Richard	<i>Historical and Biographical Sketches of the Progress of Botany in England</i>	1790	1
Pulteney, Richard	<i>Historical and Biographical Sketches of the Progress of Botany in England</i>	1790	2
	<i>Transactions of Linnaean Society of London</i>	1791	1
Hutchinson, W.	<i>History of the County of Cumberland</i>	1794	1
Hutchinson, W.	<i>History of the County of Cumberland</i>	1794	2
	<i>Transactions of Linnaean Society of London</i>	1797	3

1				
2				
3	Withering, William	<i>A Botanical Arrangement of British Plants</i>	1801	1
4				
5	Withering, William	<i>A Botanical Arrangement of British Plants</i>	1801	2
6				
7	Withering, William	<i>A Botanical Arrangement of British Plants</i>	1801	3
8				
9	Withering, William	<i>A Botanical Arrangement of British Plants</i>	1801	4
10				
11	Turner, D. and Dillwyn,	<i>Botanist's Guide through England and Wales</i>	1805	1
12				
13	L.			
14				
15	Turner, D. and Dillwyn,	<i>Botanist's Guide through England and Wales</i>	1805	2
16				
17	L.			
18				
19				
20				
21		<i>Memoirs and proceedings of the Manchester Literary and</i>	1805	1
22		<i>Philosophical Society, Second series</i>		
23				
24	Dugdale, J. T.	<i>A new British Traveller</i>	1819	1
25				
26	Dugdale, J. T.	<i>A new British Traveller</i>	1819	2
27				
28	Dugdale, J. T.	<i>A new British Traveller</i>	1819	3
29				
30	Dugdale, J. T.	<i>A new British Traveller</i>	1819	4
31				
32	Smith, James,	<i>English Flora</i>	1824	1
33				
34	Edward,			
35				
36	Smith, James,	<i>English Flora</i>	1824	2
37				
38	Edward,			
39				
40	Smith, James,	<i>English Flora</i>	1825	3
41				
42	Edward,			
43				
44	Winch, Nathaniel John	<i>Essays on the Geographical Distribution on Plants... of</i>	1825	
45		<i>Northumberland, Cumberland and Durham</i>		
46				
47	Smith, James,	<i>English Flora</i>	1828	4
48				
49	Edward,			
50				
51		<i>Magazine of natural history</i>	1829	1
52				
53	Otley, J.	<i>A concise description of the English lakes and adjacent</i>	1830	
54				
55				
56				
57				
58				
59				
60				

1				
2				
3				
4		<i>mountain</i>		
5	Stokes, J	<i>Botanical Commentaries</i>	1830	
6				
7	Hooker, William	<i>The British flora; comprising the Phaenogamous, or flowering</i>	1831	
8				
9				
10	Jackson, Sir,	<i>plants, and the ferns</i>		
11				
12	Smith P. eds	<i>Memoir and correspondence of the late Sir James Edward</i>	1832	1
13				
14		<i>Smith</i>		
15				
16	Smith P. eds	<i>Memoir and correspondence of the late Sir James Edward</i>	1832	2
17				
18		<i>Smith</i>		
19				
20				
21		<i>Magazine of Natural History</i>	1834	7
22				
23		<i>Companion of the Botanical Magazine</i>	1835	1
24				
25	Murray, A.	<i>Northern Flora: or a Descriptions of Wild Plants belonging to</i>	1836	1
26				
27		<i>the North and East of Scotland</i>		
28				
29	Watson, Hewett	<i>New Botanist's Guide to the Localities of the Rarer Plants of</i>	1837	2
30				
31	Cottrell	<i>Britain</i>		
32				
33		<i>Proceedings of the Botanical Society of London</i>	1839	1
34				
35		<i>Magazine of natural history, New Series</i>	1840	4
36				
37				
38	Otley, J.	<i>A descriptive guide to the English lakes and adjacent</i>	1842	
39				
40		<i>mountains</i>		
41				
42	Babington, Charles	<i>Manual of British Botany: Containing the Flowering Plants</i>	1843	
43				
44	Cardale	<i>and Ferns</i>		
45				
46	Hudson, J.	<i>Guide to Lakes</i>	1843	
47				
48				
49	Newman, Edward	<i>History of British Ferns and Allied Plants</i>	1844	
50				
51		<i>The Phytologist</i>	1844	1
52				
53		<i>Gardeners Chronicle</i>	1845	
54				
55		<i>The Phytologist</i>	1845	2
56				
57	Atkinson, T.	<i>Hand-Book to the English Lakes</i>	1847	
58				
59				
60				

1				
2				
3	Lancaster, E. ed.	<i>The Correspondence of John Ray</i>	1848	
4				
5	Baker, John Gilbert	<i>A supplement to Baines' Flora of Yorkshire,</i>	1854	
6				
7		<i>The Phytologist</i>	1854	5
8				
9				
10	Martineau, Harriet	<i>Complete Guide to the English Lakes</i>	1855	
11				
12	Walcott, M. E. C.	<i>Guide to the Mountains, Lakes and Northwest Coast of</i>	1860	
13				
14		<i>England</i>		
15				
16	Nicholson	<i>Annals of Kendal</i>	1861	
17				
18		<i>Phytologist, New Series</i>	1861	5
19				
20	Baker, John Gilbert	<i>North Yorkshire: studies of its botany, geology, climate, and</i>	1863	
21				
22		<i>physical geography</i>		
23				
24				
25	Linton, w. J.	<i>The Lake Country</i>	1864	
26				
27		<i>The Naturalist</i>	1864	1
28				
29	Bentham, George	<i>Handbook of the British Flora</i>	1865	1
30				
31	Bentham, George	<i>Handbook of the British Flora</i>	1865	2
32				
33				
34	Linton, W. J.	<i>The Ferns of the English Lake Country</i>	1865	
35				
36	Lowe. E. J.	<i>Our Native Ferns</i>	1865	1
37				
38	Lowe. E. J.	<i>Our Native Ferns</i>	1867	2
39				
40	Watson, Hewett	<i>Compendium of the Cybele Britannica</i>	1870	
41				
42	Cottrell			
43				
44		<i>Journal of Botany</i>	1871	9
45				
46	Watson, Hewett	<i>Topographical Botany</i>	1873	1
47				
48				
49	Cottrell			
50				
51		<i>Transactions of Cumberland Association for Advancement of</i>	1876	1
52				
53		<i>Literature and Science</i>		
54				
55		<i>Transactions of Cumberland Association for Advancement of</i>	1877	2
56				
57		<i>Literature and Science</i>		
58				
59				
60				

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

	<i>Transactions of Cumberland Association for Advancement of</i>	1878	3
	<i>Literature and Science</i>		
	<i>The Naturalist</i>	1878	4
	<i>Transactions of Cumberland Association for Advancement of</i>	1879	4
	<i>Literature and Science</i>		
	<i>Transactions of Cumberland Association for Advancement of</i>	1880	5
	<i>Literature and Science</i>		
	<i>Transactions of Cumberland Association for Advancement of</i>	1881	6
	<i>Literature and Science</i>		
	<i>Natural History Society of Glasgow</i>	1881	5
	<i>Transactions of Cumberland Association for Advancement of</i>	1882	7
	<i>Literature and Science</i>		
	<i>Transactions of Cumberland Association for Advancement of</i>	1883	8
	<i>Literature and Science</i>		
	<i>Transactions of Cumberland Association for Advancement of</i>	1884	9
	<i>Literature and Science</i>		
	<i>The Naturalist</i>	1883	8
	<i>Journal of Botany</i>	1883	21
Watson, Hewett	<i>Topographical Botany</i>	1883	2
Cottrell			
Baker, John Gilbert	<i>Flora of the English Lake District</i>	1885	
	<i>Journal of Botany</i>	1885	23
	<i>The Naturalist</i>	1886	
	<i>The Naturalist</i>	1888	
Malleson	<i>Holiday Studies of Wordsworth</i>	1890	
	<i>The Naturalist</i>	1890	

1			
2			
3		<i>The Naturalist</i>	1891
4			
5		<i>The Naturalist</i>	1893
6			
7			
8	Crump, W. B. and	<i>Flora of the Parish of Halifax</i>	1904
9			
10	Crossland, C.		
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			
26			
27			
28			
29			
30			
31			
32			
33			
34			
35			
36			
37			
38			
39			
40			
41			
42			
43			
44			
45			
46			
47			
48			
49			
50			
51			
52			
53			
54			
55			
56			
57			
58			
59			
60			

For Peer Review

Colchicum autumnale. Near Darlington, Egleston, and Butterby, Durham.

Convallaria majalis. In Scotswood, Denton, and Castle Eden Denes, in Gibside Woods, and near Winch Bridge, Teesdale; also near Warden Mill.

Ornithogalum luteum. By the Tees at Wycliffe, Barnard-Castle, and Egleston, and the Wear at Butterby.

Juncus subverticillatus—Bicheno. By the Lakes of Cumberland and Westmorland.

Epilobium alsinifolium. On Cheviot and Cronkley Fell, and on the highest ridge of Foalfoot, at the head of Longsedale. It is the *Epilobium alpinum* of Ray and Curtis.

Pyrola rotundifolia. In Castle Eden and Hawthorn Denes.

Pyrola media. In Scotswood, and East Common Wood, and by Roadley Lake, Northumberland. In Hounswood, and Blackstone Bank Wood, Durham, from 100 to 1,000 feet.

Pyrola minor. In Gibside Woods, and on Teesdale Forest, at Cocken, in Arngill, Cow Close, and Hindon Gills, and in Skullwood, near South Hamsterley, Durham. In East Common Wood, at Catcherside and Wallington, Northumberland.

Fig. 1. Excerpt from Winch's Flora of the Lake District (Winch, 1825, p. 27)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

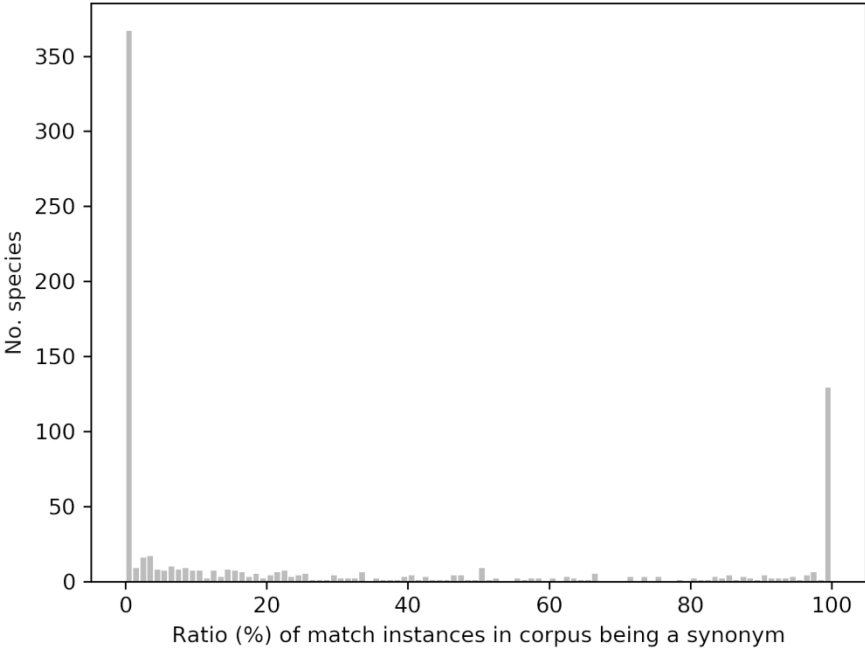


Fig. 2. Ratio (%) of plant species being matched in the corpus under synonym names

162x121mm (300 x 300 DPI)

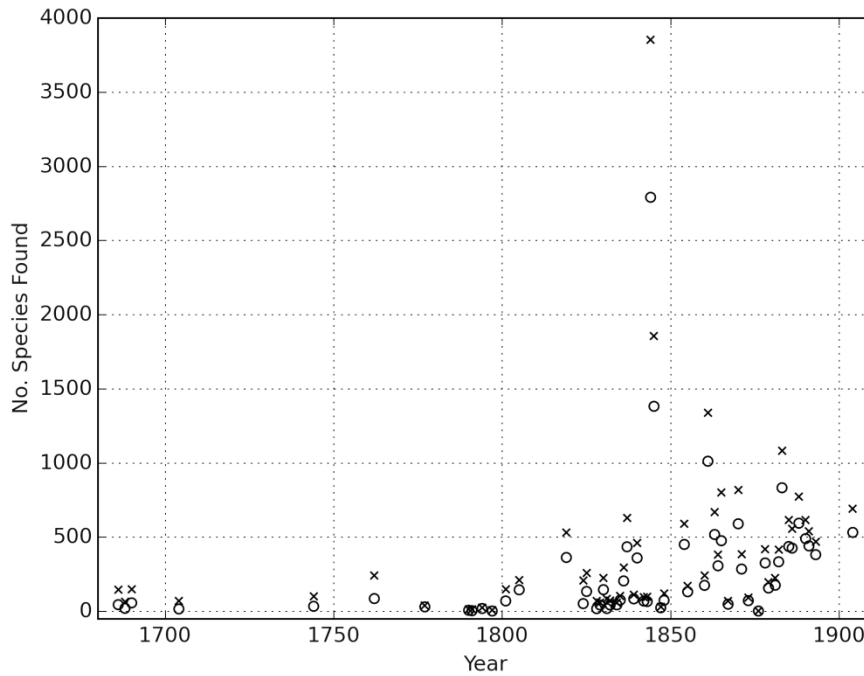


Fig. 3. Match rates across the corpus

203x152mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

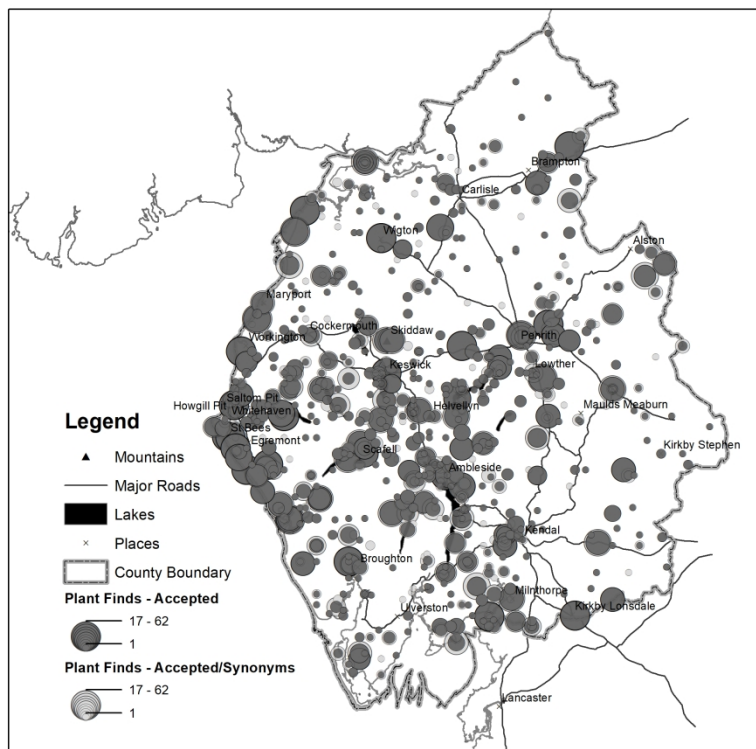


Fig. 4 Geographical distribution of plant species matched from across the corpus
210x297mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

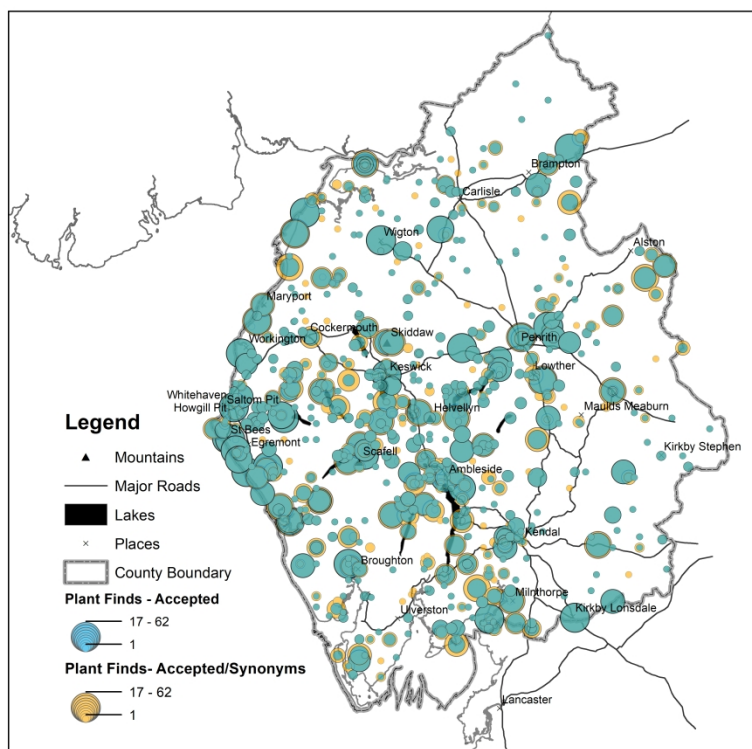


Fig. 4 Geographical distribution of plant species matched from across the corpus

210x297mm (600 x 600 DPI)