703     van Zanten, B.T., Van Berkel, D.B., Meentemeyer, R.K., Smith, J.W., Tieskens, K.F. and Verburg,

704     P.H., 2016. Continental-scale quantification of landscape values using social media

705     data. *Proceedings of the National Academy of Sciences*, *113*(46), pp.12974-12979.

706     https://doi.org/10.1073/pnas.1614158113

707     Völske, M., Potthast, M., Syed, S. and Stein, B., 2017, September. TI; dr: Mining reddit to learn

708     automatic summarization. In Proceedings of the Workshop on New Frontiers in Summarization (pp.

709     59-63).

710     Wang, Z., Ye, X. and Tsou, M.H., 2016. Spatial, temporal, and content analysis of Twitter for wildfire

711     hazards. Natural Hazards, 83(1), pp.523-540. https://doi.org/10.1007/s11069-016-2329-6

712     Wood, S.A., Guerry, A.D., Silver, J.M., Lacayo, M., 2013. Using social media to quantify nature-based

713     tourism and recreation. Sci. Rep. 3, 1-7. https://doi.org/10.1038/srep02976

714     Zhang, S. and Zhou, W., 2018. Recreational visits to urban parks and factors affecting park visits:

715     Evidence from geotagged social media data. Landscape and urban planning, 180, pp.27-35.

716     https://doi.org/10.1016/j.landurbplan.2018.08.004

717

718     **Clean Version**

719     **Reddit: A novel data source for cultural ecosystem service studies**

720     **Abstract**

721     Social media sites have been gaining traction as a source of novel data for environmental research,
722     particularly for cultural ecosystem service (CES) assessments. However, Reddit, a discussion-based
723     site, has yet to establish itself as an important source of data for CES research, possibly due to
724     researchers not being aware of its potential applications or because Reddit posts lack georeferencing
725     information. Here, we demonstrate how researchers can search Reddit for CES datasets related to
726     recreation and how specific pages on Reddit may provide data for other CES such as aesthetics.
727     Using named-entity recognition, we developed an automated method of geocoding the approximate
728     location of where images in Reddit posts were taken. Furthermore, we compare posts from Reddit
729     and Flickr for a range of recreational activities and compare the content and textual metadata of
730     images relating to hiking. Though there is potential for Reddit data to be used in spatial analysis, we
731     highlight the limitations associated with georeferencing posts. We recommend that data from
732     Reddit is best suited to assessing general trends in CES, either for a given service or place. By
733     demonstrating the value of big data from Reddit we hope to encourage its inclusion in future CES
734     and environmental research.

735     **Key words:** Cultural ecosystem services, social media, Reddit, Flickr, aesthetic values, recreation

**1.0 Introduction**

Big data from social media sites has multiple benefits over conventional methods of data collection for environmental studies, providing access to large spatio-temporal scale datasets, through inexpensive and quick data collection methods (Barve 2014). The use of social media data is therefore becoming more prominent in environmental research, ranging from the use of Twitter to understand animal life-cycles (Hart et al. 2018) and prepare for natural hazards (Wang et al. 2016; Mendoza et al. 2019), to Flickr being used to assess species niches (Peña-Aguilera et al. 2019) and map invasive species (Allain 2019). One of the biggest applications for social media data in environmental research has been the assessment of cultural ecosystem services (CES) (Ghermandi and Sinclair 2019).

CES are the non-material benefits obtained through nature and are derived from the interaction of biodiversity (biotic nature) and geodiversity (abiotic nature) (Gray 2011; Fox et al. 2020a). CES include aesthetic value, recreational services and sense of place and can enhance physical and mental well-being (Haines-Young and Potschin 2010). They can deliver multiple benefits for both residents and tourists, supporting local and regional economies (Schirpke et al. 2016; King et al. 2017). However, the exploitation, destruction and consumption of natural landscapes by humans for activities such as intensive agriculture, urban development and recreational activities can damage ecosystems and reduce their capacity to provide CES (Figueroa-Alfaro and Tang 2017). Furthermore, our understanding of CES is more limited than that of provisioning and regulating ecosystem services (Milcu et al. 2013; Díaz et al. 2018), particularly because our interactions with CES are subjective and vary between individuals, which makes obtaining practical measurements of their benefits and values difficult (Daniel et al. 2012; Havinga et al. 2020). By developing a better understanding of the natural and social drivers of CES we can help inform policy and management strategies to alleviate the threats to their sustainable use (Clemente et al. 2019). Researchers, therefore, need to understand better the supply and demand of these services over relevant temporal and spatial scales (Langemyer et al. 2018). Here, social media datasets provide relatively quick and cost-effective data collection for assessing CES, versus traditional methods, and provide novel approaches to assessing how CES are generated as well as their perceived benefits and values over a range of spatial and temporal scales (Wood et al. 2013; Ghermandi and Sinclair 2019; Fox et al. 2020b).

Due to the vast quantity of data available on social media websites can be viewed as a source of big data and therefore benefit from the emergence of big data approaches to assessing human-nature relationships (Retka et al. 2019). Social media sites, including Twitter and Weibo (microblogging sites), Flickr, Instagram and Panaramio (image sharing sites), have already been widely used to assess a range of CES. Aesthetic value has been assessed through textual metadata from Twitter (Johnson et al. 2019), image and geographic distribution from Instagram (Guerrero et al. 2016; Chen et al. 2020), and image content and geographic distribution from Flickr (Figueroa-Alfaro and Tang 2017; Tieskens et al. 2018). Recreational preferences have been studied using Flickr (van Zanten et al. 2016; Graham and Eigenbrod 2019; Gosal et al. 2019) and Weibo (Zhang and Zhou 2019). Furthermore, Flickr has also been used to assess changes in cultural values over time (Thiagarajah et al. 2018) and identify trade-offs between multiple CES (Allan et al. 2015). However, some social media sites have either ceased operating (e.g. Paramio) or introduced restrictions to accessing data (e.g. Instagram) and therefore Flickr is becoming the main source of data for CES studies (Langemeyer et al. 2018; Retka et al. 2019).

Metadata available from Reddit, the social news aggregation and discussion orientated social media website, has been used in a vast array of scientific studies across a range of disciplines (Baumgartner et al. 2020), including health and psychology (e.g. Jamnik and Lane 2017; Park et al. 2018),

782 technological development (e.g. Derczynski et al. 2017; Volske et al. 2017) and political studies (e.g.
783 Guimaraes et al. 2019). Reddit, which is broken up into different forums or "subreddits'' themed
784 around different topics, allows for user to post a range of media such as images and text posts.
785 These posts, along with their associated metadata, draws parallels to the types of data from other
786 social media sites that are currently used in CES studies. However, there appears to have been little
787 to no application of big data from Reddit to assess any ecosystem service. A systematic review of the
788 applications of social media data in environmental research did not include any studies using Reddit
789 as a source of data (Ghermandi and Sinclair 2019). A search of the titles abstracts and keywords on
790 Web of Knowledge (https://wok.mimas.ac.uk) and Scopus (www.scopus.com) for "ecosystem
791 servic*" (the * denotes any end to the term e.g. service or services) AND "Reddit", carried out on
792 10th February, 2021, returned no results.

793 As there have been few studies comparing social media sources for CES, there is a need for a greater
794 understanding of the impacts of differences in data availability and biases among the various social
795 media sites used as data sources (Ostera-Roza et al. 2018). We therefore find it surprising that big
796 data from Reddit has yet to be explored in the context of CES, though we postulate that this may be
797 for two key reasons: researchers not being familiar with the website and its potential uses; and that
798 posts on the website are not georeferenced. In this paper, we aim to provide an overview of Reddit,
799 and to compare data from the site with that from another social media site, Flickr. We provide
800 examples of how data from Reddit can be used to assess recreational, aesthetic, spiritual and
801 cultural CES and address how Reddit can be a novel source of data for commonly used CES methods
802 such as assessing image contents and textual sentiment. We also provide an insight to the potential
803 uses and limitations of Reddit for spatial assessments.

804 **2.0 Methods**

805 Here we present multiple methods for searching Reddit for data suitable for CES assessments via its
806 Application Programming Interface (API), a computing interface that allows researchers to access a
807 platform via code. First, we searched the site for all posts containing a specific keyword and
808 compared these posts to those found using the same keyword search on Flickr. Second, we searched
809 for posts on subreddits that are based around topics of interest to CES research. Third, we
810 demonstrate a method for geocoding an estimated location for posts from Reddit as well as
811 combining a place keyword search with another keyword, or within a subreddit to find posts linked
812 to a particular location.

813 *2.1 Data sources*

814 *2.1.1 Reddit*

815 Reddit is a social media site consisting of over two million different communities called subreddits
816 (Table 1). Subreddits are built around a topic, each with their own rules on posts and comments. The
817 type of post is highly variable among subreddits. For example, the subreddit "r/EarthPorn" is limited
818 to photographs of landscapes, accompanied by a text title and a comment section, whereas the
819 subreddit "r/Culture" hosts primarily text-based posts with a title and a comment section.

820 Table 1. Selected examples of subreddits linked to cultural ecosystem services.

| Service | Subreddit | Extract of group description | Number of members (10th February 2021) | Primary metadata type |
| --- | --- | --- | --- | --- |

| Aesthetic views | r/EarthPorn | "EarthPorn is your community of landscape photographers and those who appreciate the natural beauty of our home planet." | 20.9m | Images |
|---|---|---|---|---|
| | r/BotanicalPorn | "High quality images of plants (fungi are allowed!)." | 167k | Images |
| Recreational activities | r/Outdoors | "Outdoors is for *all* outdoor experiences, not limited to any specific interest. Caving, mountain climbing, cycling, bushcraft, gardening, sailing, plants, birds, trees, going for a stroll -- it's all on topic here!" | 2.7m | Images |
| | r/Hiking | "The hikers' subreddit." | 1.3m | Images |
| Tourism | r/Travel | "r/travel is a community about exploring the world. Your pictures, questions, stories, or any good content is welcome." | 5.7m | Images and text |
| Spirituality and sense of place | r/Spirituality | "Here, we discuss such things as personal transformation, the meaning of life, death, and moments of clarity." | 190k | Text |
| | r/Culture | "A subreddit dedicated to sharing and discussing the many aspects of culture" | 6.3k | Text |

821

822 Posts and comments from Reddit can be searched and returned through the Reddit API, with text
823 and image posts, as well as their metadata including the, title, comments, date posted, how many
824 upvotes (the number of people that like a post) a post has, and the ratio of upvotes to downvotes
825 (the number of people that dislike a post). These data types are similar to data already being used in
826 CES and social media studies derived from Flickr, Instagram and Twitter.

827 Accessing data from Reddit has multiple benefits for researchers. First, data from Reddit is freely
828 available. Second, the data is accessible across multiple software tools and programming languages.
829 For example, the Pushshift tool (Baumgartner et al. 2020) provides researchers with an accessible
830 method for querying and retrieving data. The tool also benefits researchers by providing built-in
831 functionality which overcomes Reddit's 100 object limit per search. For researchers familiar with
832 writing scripts, functionality for searching the Reddit API is available in multiple programming
833 languages: packages "RedditExtractoR" (Rivera 2019) and "rreddit" (Kearney 2020) for the R
834 environment; "Python Reddit API Wrapper" (Boe 2020) within the Python environment; "jReddit"
835 (jReddit 2020) within the Java environment.

836 *2.1.2 Flickr*

837 Flickr is a popular social media site that hosts images and videos with up to 25 million uploads a day
838 (Ding and Fan 2019). Flickr has a broad user base, with a range of motivations for uploading
839 photographs (Oteros-Rozas et al. 2018), and therefore has potential as a source of data for a wide
840 range of CES. Posts on Flickr can have associated metadata that includes textual titles, description
841 and tags; spatial location in the form of latitude and longitude of where the image was taken; and
842 the time and date the image was taken. Flickr metadata is accessible through tools such as the

843    "photosearcher" package in the R environment (Fox et al. 2020b), and stand-alone software such as

844    the InVEST Recreational tool (Sharp et al. 2020).

845    *2.2 Data collection and analysis*

846    *A reproducible R file for the data collection methods has been included in the supplementary*

847    *material. To comply with API terms and privacy policies all data sets were anonymised, stored with*

848    *multiple layers of security and any unnecessary metadata was deleted.*

849    *2.2.1 Keyword search*

850    First, to find posts related to recreational activities, we searched the Pushshift tool (Baumgartner et

851    al. 2020) for any posts on Reddit containing a single keyword for four different activities; "hiking",

852    "camping", "skiing" and "kayaking", found in any textual metadata uploaded by the user e.g. the

853    posts title or description. We also constrained the search to any posts that were uploaded between

854    the 1st of January 2020 and the 1st of January 2021. We then repeated this query on Flickr, using the

855    photosearcher R package (Fox 2020b), ensuring that we made a comparable search using the same

856    keywords, again found in any of the posts textual metadata, and within the same uploaded date

857    range. We summarized the number of uploads per month as well as the mean character length of

858    the posts title and text, and the mean number of likes and comments on the images for each activity

859    across platforms. Furthermore, as posts on Reddit can be in a range of formats other than images

860    and text traditionally used in CES studies (e.g. links to other websites or videos), we calculate the

861    percentage of posts that were images or text.

862    To compare the contents of the images posted on the two sites, we took a random sample of 1,000

863    images related to hiking from both sites (images listed as adult material were not included in the

864    sample selection). The contents of the images were automatically tagged using the Google Cloud

865    Vision API (Google Cloud Vision 2020). The Google Cloud Vision API is a machine learning algorithm

866    that labels the content of images. The algorithm is based on a large pre-trained dataset and can label

867    image contents into millions of predefined categories including objects and expressions. Here, we

868    used the "imgrec" R package (Schwemmer 2019) to label each image with the 10 objects the

869    algorithm first detects. To ensure that the image contents were accurate without manual validation

870    we only kept labels that had a confidence score of > 0.6 (Gosal et al. 2019).

871    To compare the hiking images from Reddit and Flickr we used a chi-square test to compare the two

872    sources of data in terms of their image content (frequency of Google Cloud Vision API labels). As the

873    dataset is relatively large, some statistical tests may indicate statistical significance ($p < 0.05$)

874    irrespective of real-world significance in the data. Furthermore, statistical significance does not

875    provide information on the size of the effect (Kim 2017). Here, we primarily focus on the individual

876    contribution of features ($x^2{}_i$ eq. 1) to the total effect size $x^2 = \sum x^2{}_i$ , enabling us to understand

877    better the difference between the two datasets (Oakes and Farrow 2006).

878    $x^2{}_i = \frac{(\mathrm{obs}_i - \mathrm{exp}_i)}{\mathrm{exp}_i}$ (1),

879    where $\mathrm{obs}_i$ and $\mathrm{exp}_i$ are the observed and expected values of feature $i$, respectively.

880    As textual metadata can be useful for understanding characteristics of CES or eliciting emotion of

881    CES beneficiaries (Brindley et al. 2019; Hale et al. 2019), we also returned textual metadata for

882    analysis for the random sample of hiking images. As images uploaded to Reddit can only contain a

883    title, with no description text, the most comparable source of textual data for images from Flickr and

884    Reddit are the comment sections. We summarised the number of comments for these images as

885  well as the number of unique users interacting with the posts. The sentiment expressed in each
886  comment was calculated using the Afinn dictionary (Nielsen 2011), which has previously been used
887  to assess the sentiment value expressed in social media text posts (Koto and Adriani 2015). This
888  dictionary ranks words on a -5 (negative sentiment) to +5 (positive sentiment) scale. The sum
889  sentiment of each post was calculated, and the mean sentiment score of the posts on each site
890  calculated. We also filtered out automated messages, weblinks and commonly used words such as
891  "the" and "is" and calculated the most frequently used words in comments on the two sites.

892  *2.2.2 Subreddit Search*

893  A unique aspect to the Reddit API is the ability to search individual subreddits. Here, we searched
894  four subreddits that are themed around the aesthetic value of nature; "r/EarthPorn",
895  "r/BotanicalPorn", "r/WaterPorn" and "r/DesertPorn", as well as posts from two subreddit about
896  two recreational activities ("r/Birding" and "r/Scuba") and two subreddits that discuss spirituality
897  and culture ("r/Spirituality'' and "r/Culture"). The results were limited to posts uploaded in the year
898  2020. The aesthetic views subreddits have a set of rules that mean all posts on the subreddit are of
899  photographs pertaining to nature. Table 2 summarises the submission rules for the "r/EarthPorn"
900  subreddit, these rules are similar across the other aesthetic subreddits assessed, though the subject
901  of the photograph differs. The rules for the recreational and spiritual subreddits allow for both
902  images and discussion-based posts. To compare the contents of the images posted different
903  subreddits, we took a random sample of 1,000 images posted on "r/EarthPorn" and
904  "r/BotanicalPorn". These images were then automatically tagged using the Google Vision Cloud API
905  and the contents of the two sets of images were compared using a chi-square test.

906  Table 2. Selected rules for submissions to "r/EarthPorn" (as of the 10th February 2021).

| Rule | Description |
|---|---|
| A photograph | "No Paintings, illustrations, gifs, videos, or interactive images." |
| An image featuring a natural landscape | "Images must have visible land. Images with humans, machines, boats, roads, airplanes, farms, animals, buildings, or other man made objects in them will be removed." |
| A photograph you took (OC) | "Or one which you can provide and post the original source for. Do not rehost non OC images to reddit or imgur." |
| An unsilhouetted image | "Images where details in the landscape are not visible due to silhouetting will be removed." |
| The location of the area in the photo | "When it comes to location, the more specific the better. If you wish to not disclose the location you should at the very least name the state/country. Rule of thumb for naming only the location (e.g. a lake, mountain): if one can find the place immediately by searching it in google it's fine. For possibly ambiguous locations add state/country for safety." |

907

908  *2.2.3 Potential spatial uses for Reddit*

909  As Reddit posts are not geolocated, it is not possible to directly map the distribution of the CES
910  expressed in the posts. Instead, we developed an automated method for estimating the
911  approximate location of images posted to Reddit, following a similar method to Harrington (2018).
912  The subreddit "r/EarthPorn" requires that posts must contain the image location in the title. To
913  extract the location name, we used named-entity recognition, a technique that classifies words in a
914  text into predefined categories, one of which is a named location, on the 1,000 images randomly
915  sampled from "r/EarthPorn" (Alfred et al. 2014). Named-entity recognition was carried out using the

916 "entity" R package (Rinker 2015). A subset of 10% of the name-entities were manually validated by
917 comparing the returned name-entity with the post title. The extracted location names were then
918 geolocated using the Google Maps API through the "ggmap" R package (Kahle and Wickham 2013).
919 Based on the place name, the Google Maps API provides an estimated latitude and longitude. The
920 global distribution of both sets of data was mapped and the percentage of uploads from each
921 continent was calculated.

922 To assess whether Reddit posts can be used to assess general CES trends for a given location, we
923 also searched Reddit for posts containing given place names. We carried out two types of search;
924 first, we searched for posts containing a given place name as well as the term "hiking" and second,
925 we searched for posts containing a given place name within the subreddit "r/EarthPorn". The place
926 names were chosen to represent a range of scales; national ("USA" and "UK"), regional ("Wyoming"
927 and "Scotland") and National Park ("Yellowstone" and "Cairngorms"). The searches were carried out
928 for any post uploaded between the 1st of January 2010 and the 1st of January 2021. Total number of
929 posts was calculated.

930 **3.0 Results**
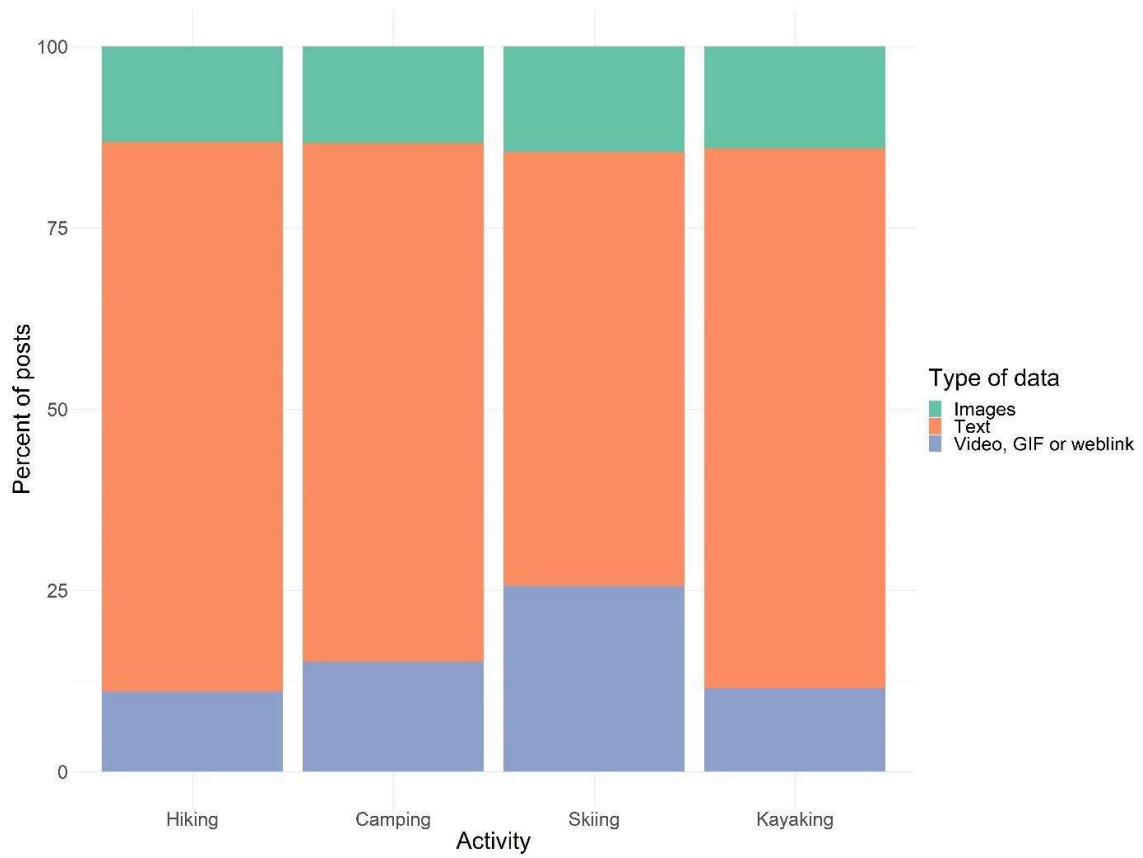
931 *3.1 Full datasets*

932 For each activity, the number of posts vary across each site. For hiking there were a similar number
933 of posts uploaded to each site in 2020 with 145,036 hiking posts on Reddit and 148,535 on Flickr.
934 There were also a similar number of posts relating to skiing across the two sites: 41,703 post on
935 Reddit and 59,455 posts on Flickr. For camping, more posts were uploaded to Reddit (143,446) than
936 Flickr (66,818); however for kayaking more posts were uploaded to Flickr (48,659) than Reddit
937 (15,107). The number of uploads fluctuate across the year for both websites (Fig. 1). For hiking and
938 skiing, even though there were a similar number of posts, Reddit had a greater quantity of unique
939 users generating the posts. For hiking, Reddit had 88,075 unique users posting whilst Flickr had
940 9,392, while for skiing Reddit had 20,934 unique users whilst Flickr had 4,309.
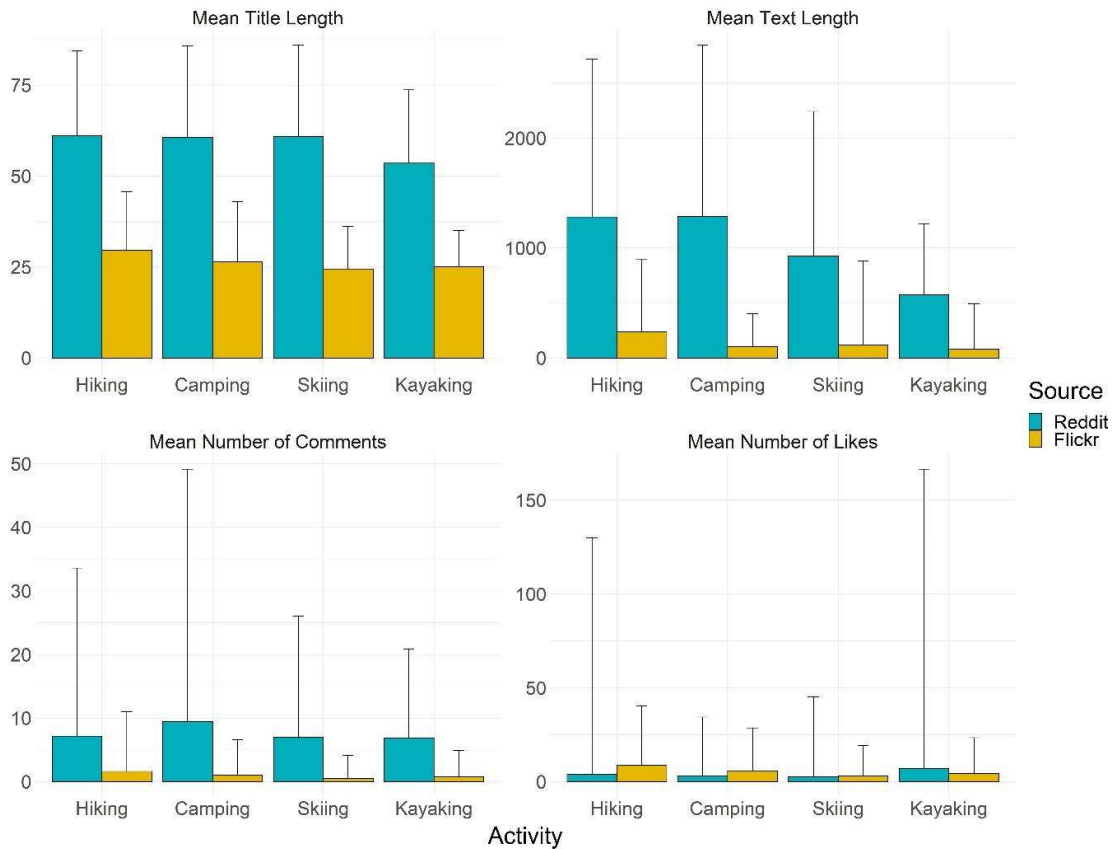
941

942
943   Fig.1: Uploads of posts including the words "hiking", "camping", "skiing" and "camping" to Reddit
944   and Flickr between the 1st of January 2020 and the 1st of January 2021.

945   For each activity that we searched, many of the posts uploaded to Reddit were text based (Fig. 2).
946   Only around 15% of the posts returned via a keyword search from Reddit were images. Compared to
947   posts uploaded to Flickr, posts on Reddit, in general, have longer titles and text descriptions as well
948   as a higher number of comments (Fig. 3). Posts relating to hiking, camping and skiing on Flickr have,
949   on average, more likes than posts on Reddit, though Kayaking posts on Reddit have a higher mean
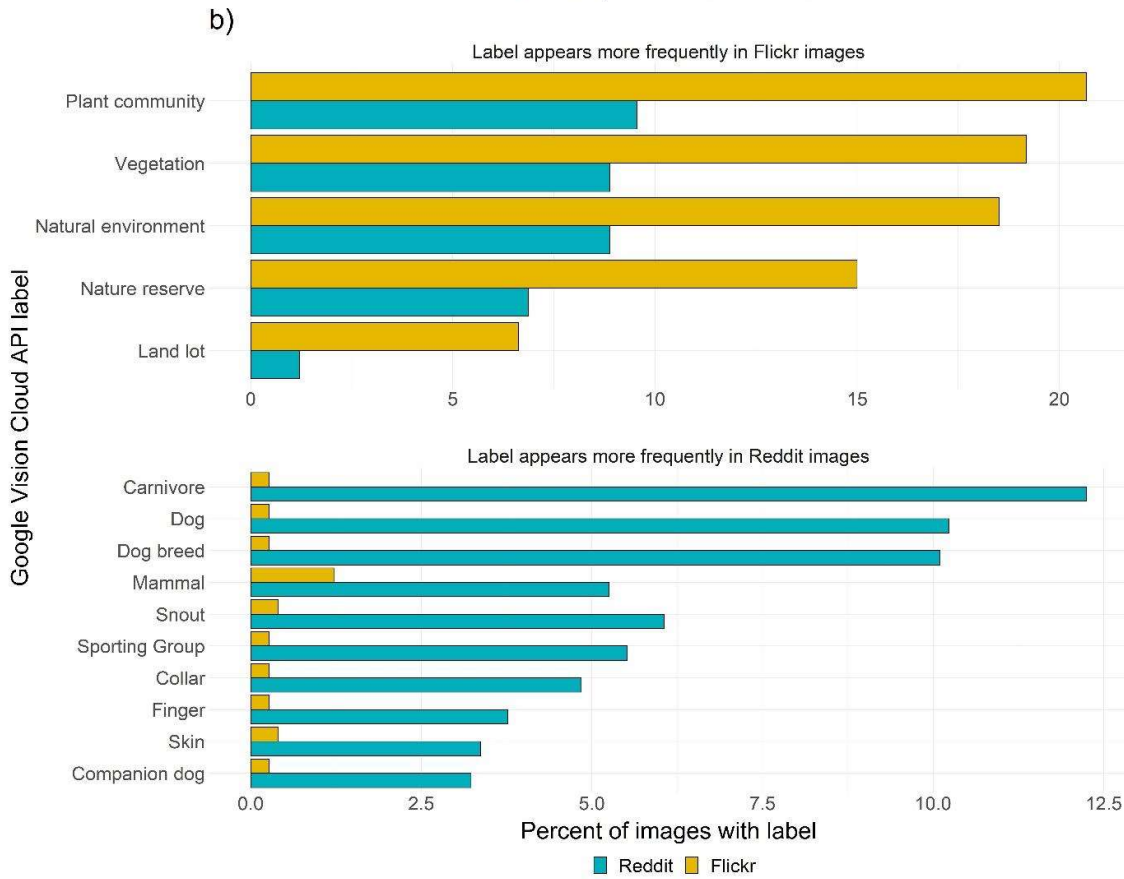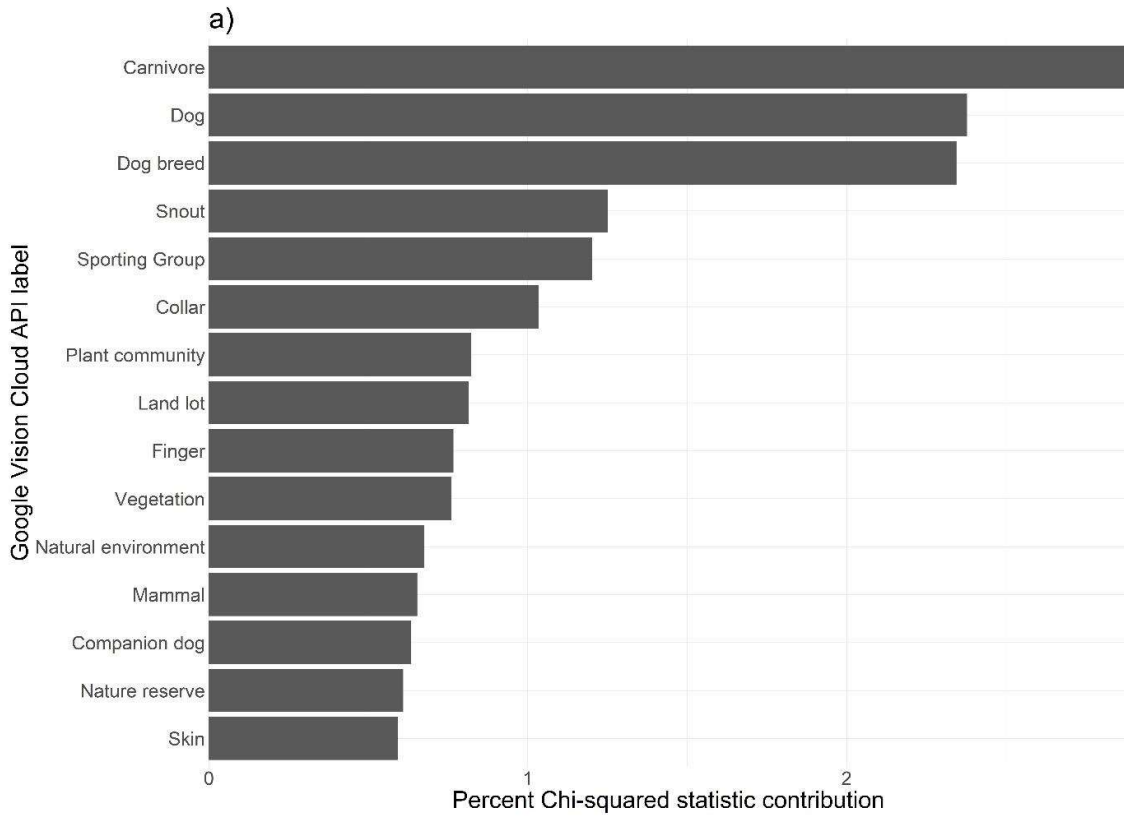950   number of likes than those on Flickr.

951

952     Fig. 2: Types of posts uploaded to Reddit.

953

Fig. 3. Summary of posts made on Reddit and Flickr (mean + 0.5 standard deviations).

While the majority of the most labelled objects were common between the two sets of images (e.g. tree and mountain), there was an overall significant difference in the contents of the two sets of photographs labelled by the Google Cloud Vision API, $(x^2 = 3,127.5, df = 1230, N = 13,582, p < 0.001)$ The 15 Google Cloud Vision API labels (1.22% of the total number of unique labels) that had the highest contribution to the total $x^2$ effect size contributed 17.42% of the total $x^2$ value (Figure 4a). Of these 15 labels, five ("plant community", "vegetation", "natural environment", "nature reserve" and "land lot") appeared more frequently in the images from Flickr (Fig. 4b). Though more frequent in Flickr images, the Google Cloud Vision API labels such as "plant community" and "natural environment" were present in 71 and 66 Reddit images, respectively. The other ten highest contributing labels, relating to dog walking, sports and people were more frequently photographed in Reddit images, with the labels such as "dog" and "dog breed" only being tagged in two of the Flickr images.

a)

b)

Label appears more frequently in Flickr images

Label appears more frequently in Reddit images

968      Figure 4: a) The 15 Google Cloud Vision API labels which had the greatest contribution to the overall
969      Chi-squared statistic (larger values indicate a larger difference between Reddit and Flickr); b) The
970      percentage of Reddit and Flickr images that the 15 labels appeared in.
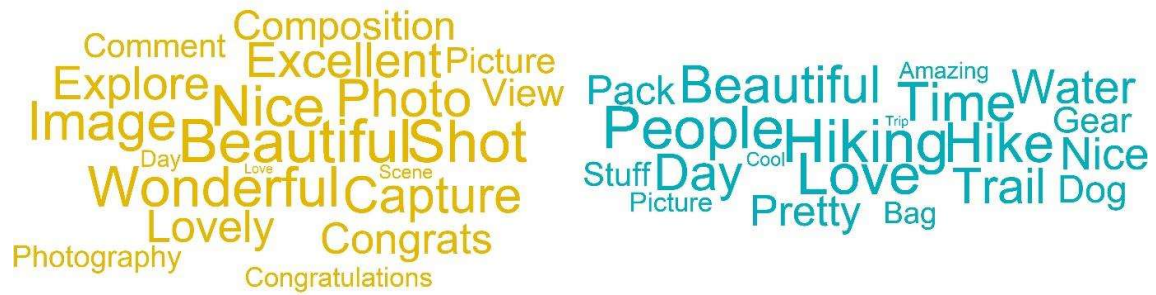
971      For the 1,000 hiking images from Reddit, 702 posts had comments, while for the 1,000 Flickr images
972      only 116 posts had comments. The 6,602 comments on the Reddit post were made by 4,142 unique
973      users, while the 1,702 Flickr comments were made by 1,119 unique users. A sentiment score could
974      be calculated for 642 Reddit comments and 108 Flickr comments, those where a score could not be
975      calculated did not contain any words in the AFINN dictionary. In general, the sentiment expressed in
976      Flickr comments was far higher than those on Reddit (Fig. 5). Only 1.90% of Flickr images expressed a
977      negative or neutral sentiment, whilst 11.66% of Reddit comments expressed a negative or neutral
978      sentiment. Many of the non-unique Flickr comments are "awards", a small sticker accompanied by a
979      text phrase, while on Reddit they were automatically generated messages from moderators of the
980      subreddit. After filtering, the most used words in Flickr and Reddit comments suggest that Flickr
981      users more frequently comment general positive comments regarding the picture, such as
982      "wonderful" and "excellent", while Reddit users more frequently comment regarding features of the
983      photograph, such as "trail", "water" and "dog" (Fig. 6).



984
985      Fig. 5: Mean +- 0.5 standard deviations for the AFINN sentiment score expressed in the comments of
986      hiking images on Reddit and Flickr.
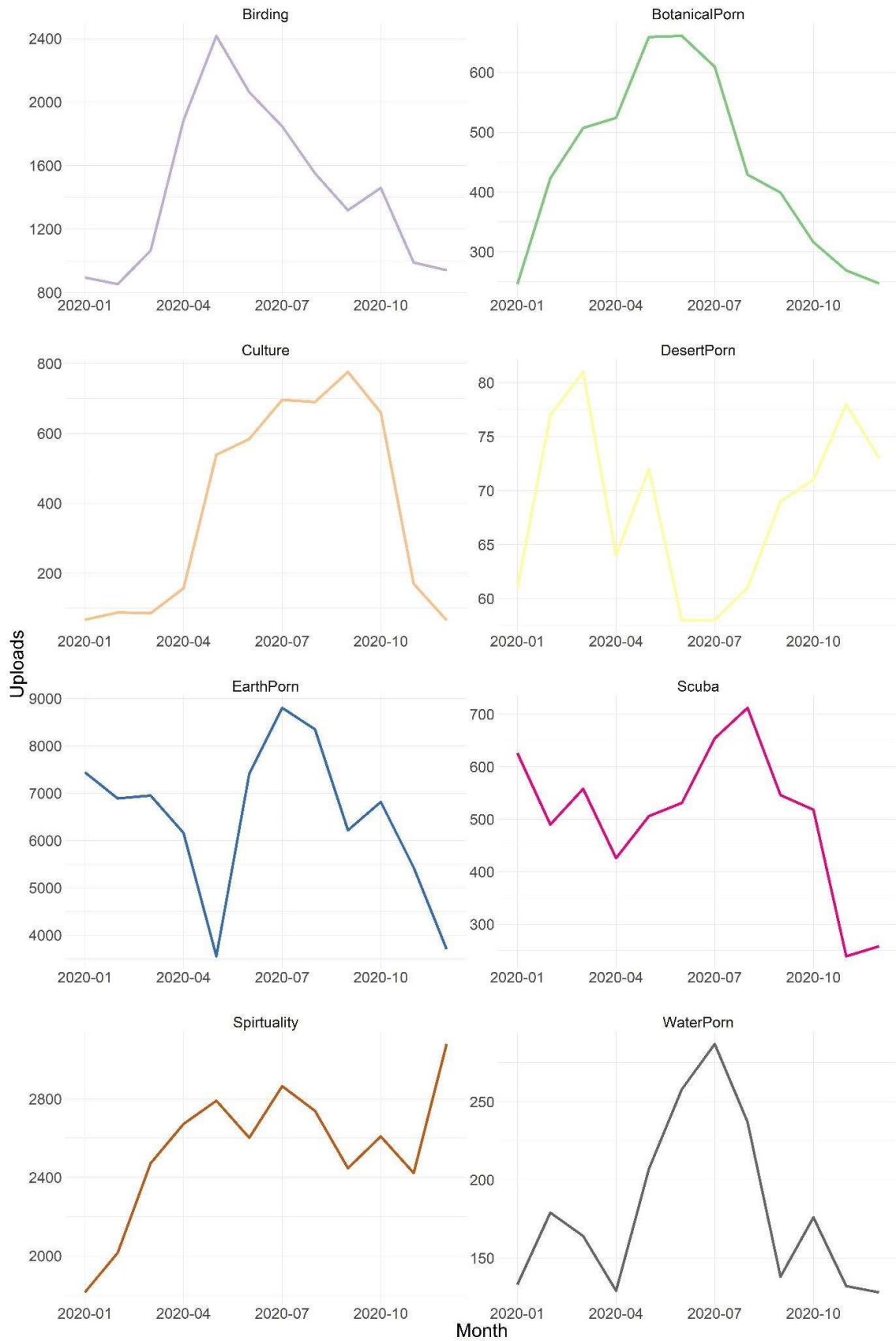
Fig. 6: The 20 most frequently used word in Flickr and Reddit comments after filtering.

*3.2 Subreddit search*

Of the subreddits relating to aesthetic values, "r/EarthPorn" was the most popular of the four we searched, with 77,717 photographs uploaded in 2020. The subreddit "r/BotanicalPorn" had 5,289 uploads, "r/WaterPorn" 2,168 and "r/DesertPorn" 823. The number of uploads to each subreddit varies by month (Fig. 7). The subreddits "r/Spirituality" and "r/Culture" also had a relatively large number of uploads during the year 2020 with 30,528 and 4,579 uploads, respectively. Furthermore the recreational based subreddits "r/Birding" had 17,280 post and "r/Scuba" had "6,064" posts.
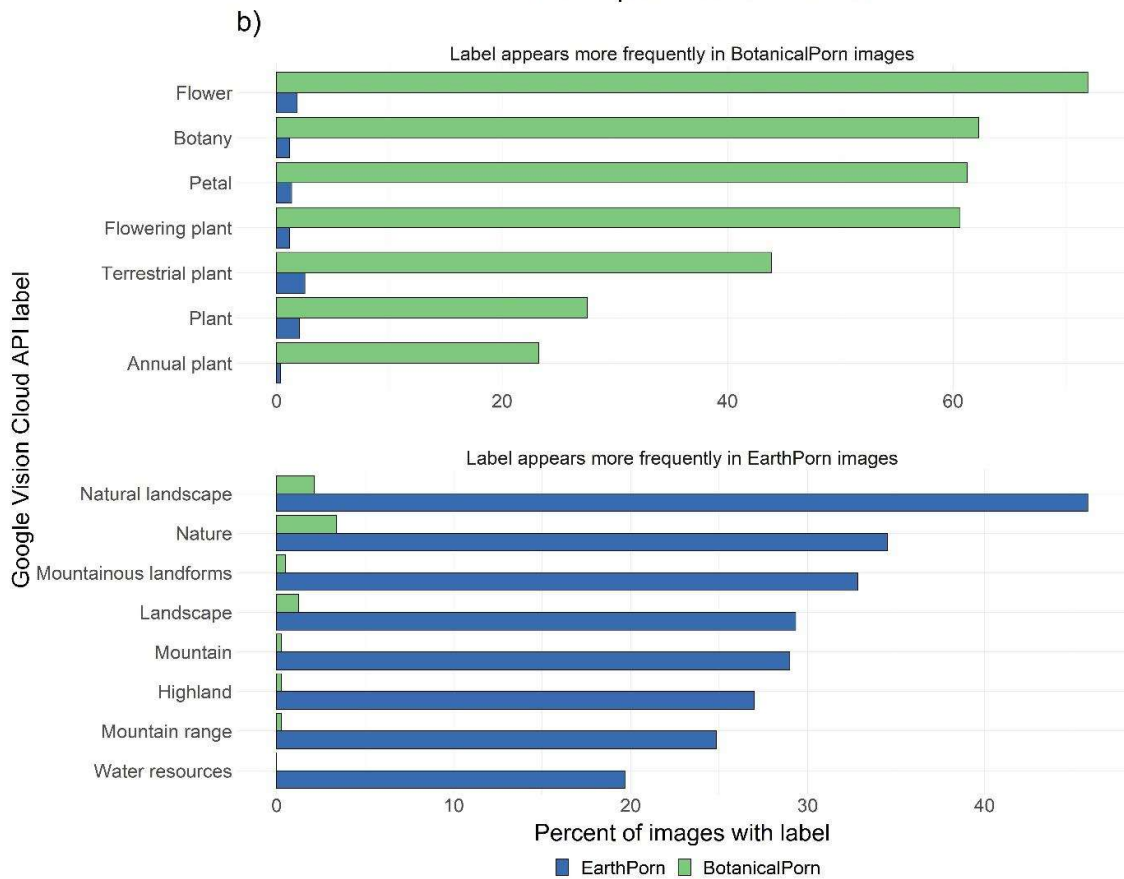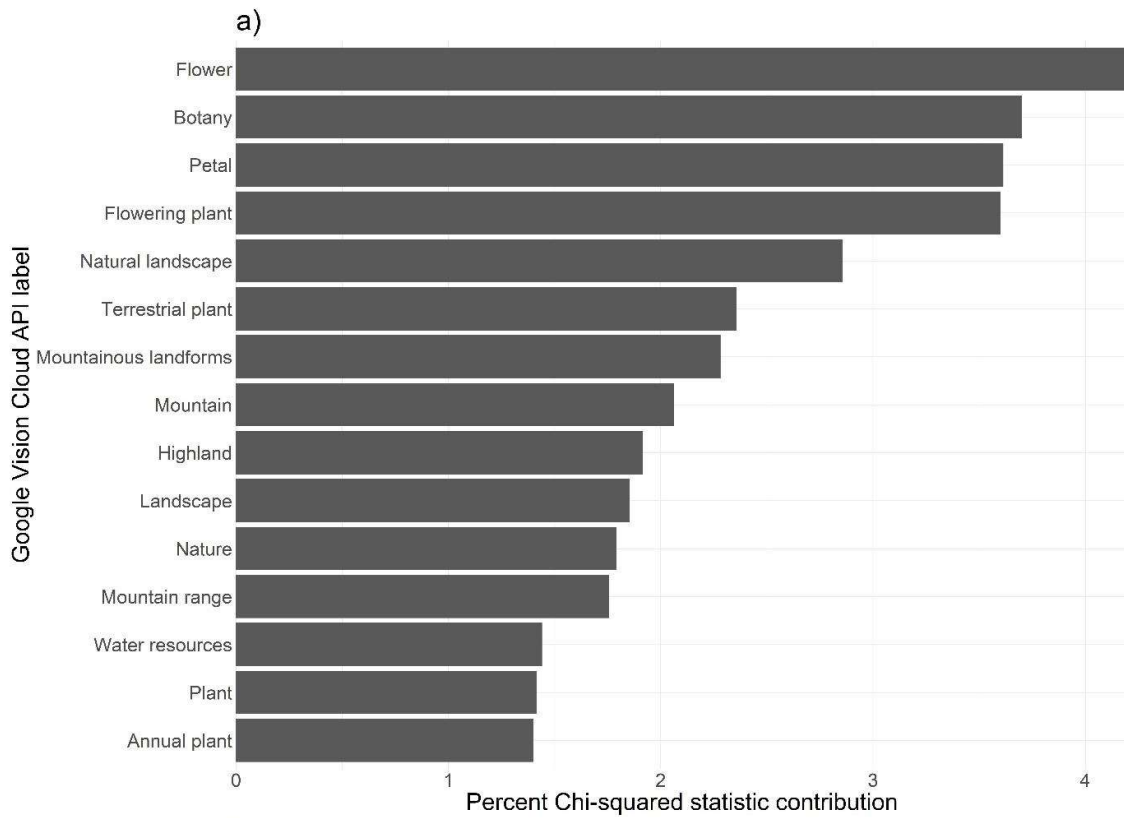
Figure 7: Uploads of posts to the subreddits "r/Birding", "r/BotanicalPorn", "r/Culture",

1000 "r/DesertPorn", "r/EarthPorn", "r/Scuba", "r/Spirituality" and "r/WaterPorn" between the 1st of
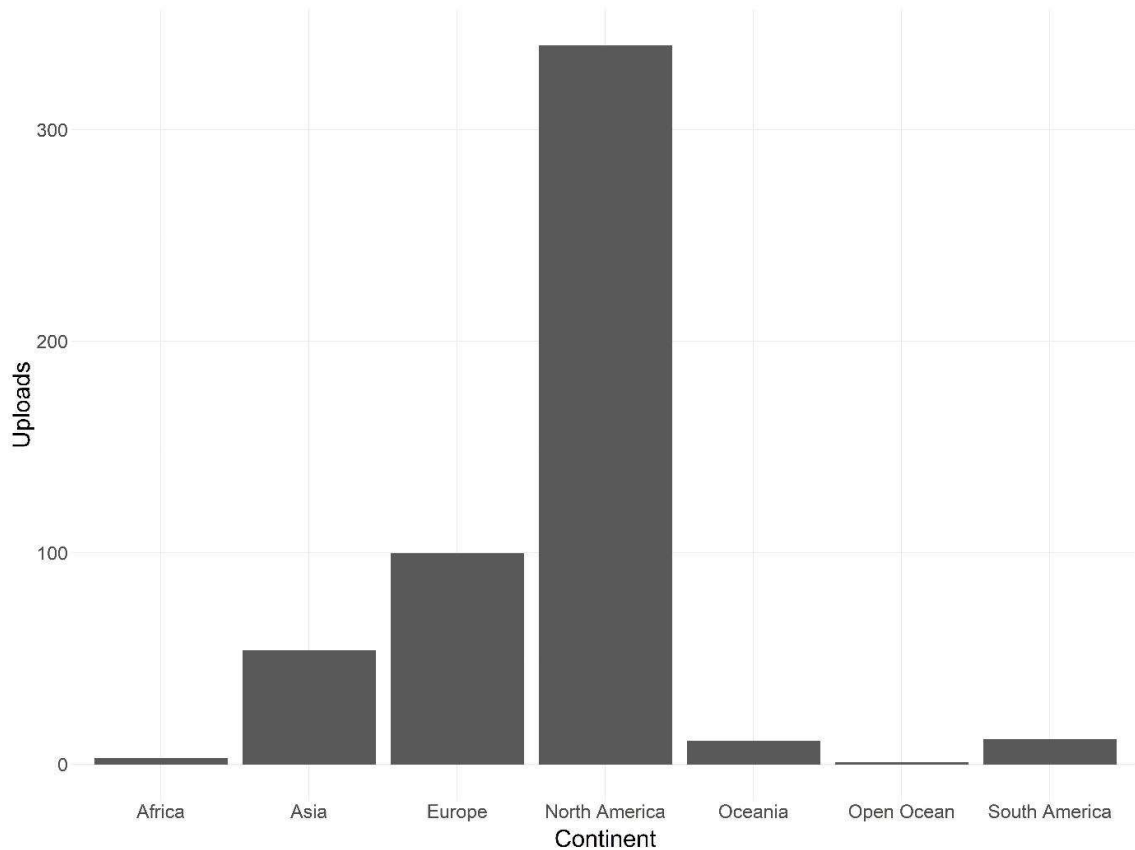1001 January 2020 and the 1st of January 2021.

1002 The was a large contrast between the labelled objects in images from the "r/EarthPorn" and
1003 "r/BotanicalPorn" subreddits, with an overall significant difference in the contents of the two sets of
1004 photographs labelled by the Google Cloud Vision API, $(x^2 = 10{,}205.5, df = 765, N = 13{,}196, p <$
1005 $0.001)$. The 15 Google Cloud Vision API labels (1.95% of the total number of unique labels) that had
1006 the highest contribution to the total $x^2$ effect size contributed 36.26% of the total $x^2$ value (Figure
1007 8a). Of these 15 labels, seven, all relating to plants and flowers, appeared more frequently in the
1008 images from "r/BotanicalPorn" (Fig. 8b). The other highest contributing labels, relating to
1009 landscapes, were more frequently photographed in "r/EarthPorn" images.

a)

b)

Label appears more frequently in BotanicalPorn images

Label appears more frequently in EarthPorn images

1011 Figure 8: a) The 15 Google Cloud Vision API labels which had the greatest contribution to the overall
1012 Chi-squared statistic (larger values indicate a larger difference between Reddit and Flickr); b) The
1013 percentage of "r/EarthPorn" and "r/BotanicalPorn" subreddit images that the 15 labels appeared in.

1014 *3.3 Potential spatial uses for Reddit*

1015 Our automated method for estimating image location returned a latitude and longitude for 574
1016 "r/EarthPorn" subreddit images (57.4%) (Fig. 9). The vast majority of images (65.26%) were
1017 distributed across North America. Overall, there were fewer images taken in the other continents,
1018 with Europe and Asia having relatively higher numbers of images than Oceania, South America and
1019 Africa.



1020

1021 Figure 9: Estimated locations of a subset of photographs from the "r/EarthPorn" subreddit.

1022 When searching the Reddit API for posts relating to a place name as a keyword the number of posts
1023 vary depending on the spatial scale and location (Table 3). For both searches containing a separate
1024 keyword ("hiking") and those from a specific subreddit ("r/EarthPorn") a large number of posts were
1025 returned.

1026 Table 3: Number of posts, when searching Reddit with a location name as a criterion.

| Country | Scale | Search Criteria | Number of Posts |
|---|---|---|---|
| USA | National | Text = "USA" AND "hiking" Subreddit = any | 13,148 |
| | | Text = "USA" Subreddit = "r/EarthPorn" | 12,336 |
| | Regional | Text = "Wyoming" AND "hiking" | 1,209 |

| | | Subreddit = any | |
|---|---|---|---|
| | | Text = "Wyoming" Subreddit = "r/EarthPorn" | 3,399 |
| | National park | Text = "Yellowstone" AND "hiking" Subreddit = any | 2,794 |
| | | Text = "Yellowstone" Subreddit = "r/EarthPorn" | 4,334 |
| UK | National | Text = "UK" AND "hiking" Subreddit = any | 8,196 |
| | | Text = "UK" Subreddit = "r/EarthPorn" | 5,539 |
| | Regional | Text = "Scotland" AND "hiking" Subreddit = any | 2,528 |
| | | Text = "Scotland" Subreddit = "r/EarthPorn" | 5,539 |
| | National park | Text = "Cairngorms" AND "hiking" Subreddit = any | 87 |
| | | Text = "Cairngorms" Subreddit = "r/EarthPorn" | 131 |

1027

## 4.0 Discussion

The main aim of this paper was to understand the potential applications for Reddit as a complementary or alternative source of CES data from social media sites. Here, we explored two methods of searching the Reddit API: a keyword search and searching specific subreddits. In general we were able to return a relatively large number of posts relating to a range of CES (recreation, aesthetic, spirituality and culture). Searches made via the keywords search showed that Reddit has a comparable number of available posts on recreational CES to Flickr. However, the posts returned via a keyword search on Reddit are primarily text based, which is unsurprising given that Reddit is marketed as a discussion-based social media site. The two sites had similar numbers of posts for hiking and skiing, though Reddit had more posts about camping and Flickr had more posts about kayaking. This suggests that the choice of site may depend on the activity of interest and thus the suitability for CES research is context dependent. Furthermore, even when the posts had a similar number of uploads between sites, the posts on Reddit were contributed by a far greater quantity of unique users. This gives rise to the potential for posts to be generated by a more diverse user base than Flickr. There are however socio-demographic biases associated with social media sites (Duggan and Smith 2018; Rekta et al. 2019), and these need to be explored fully before making generalisations about the wider population.

The biggest limitation of Reddit is that the posts do not have geotagged locations. Our automatic method for estimating the approximate location of a photograph calculated latitude and longitude for 57.4% of the Reddit posts. From our analysis of landscape photographs, the distribution of images uploaded to the "r/EarthPorn" subreddit are primarily concentrated in North America, though many images were also from Europe and Asia. Harrington (2018) estimated the distribution of the Reddit users base through geolocating statements in their comments and found that the demographic was primarily people living in North America, followed by Europe and Asia. Harrington (2018) also provides a potential method of establishing user origins, a key feature in understanding CES interaction from Flickr (Wood et al. 2013; Sinclair et al. 2020). The demographic of users and distribution of posts may have implications for studies that wish to assess CES across different

1055 continents, with previous studies assessing CES in North America potentially missing out on the
1056 wider range of photographs available from Reddit.

1057 A potential issue with Reddit, as well as other social media sites such as Flickr, is the potential biases
1058 introduced by the demographics of their users. For example, though Reddit has a large user base
1059 with high socio-demographic diversity, with an estimate that around 6% of internet users were
1060 active on Reddit, there is bias towards male users (8% of male internet users compared to 4%
1061 female) and a bias to younger users, with a higher percentage users aged 18-49 than those over 50
1062 (Duggan and Smith 2018). Furthermore, the users of both Reddit and Flickr are concentrated in
1063 western, developed countries. Where studies are at a global or super-continental scale, data from
1064 Reddit and Flickr should therefore be used in combination with each other and with other sources of
1065 data that are popular in other areas of the world. For example, in China where Flickr is banned and
1066 Reddit is not a popular social media site, alternative social media sites such as Weibo (Zhang and
1067 Zhou 2019), or travel comment portals websites such as Tuniu Travel (Dai et al. 2019), should be
1068 used to bridge the gap in CES data. At local and regional scales other sources of data may also help
1069 to complement social media data such as on-site survey data (Sinclair et al. 2020), online surveys
1070 (Moreno-Llorca et al. 2020) and national statistics (Graham and Eigenbrod 2019). Future work
1071 should begin to assess the respective biases in these alternative sources to ensure they are
1072 comparable. Furthermore, both Flickr and "r/EarthPorn" are related to images pertaining to high-
1073 end photography, which may restrict the demographics to only those with access to such technology
1074 (Chen et al. 2020). One possible source of data that we suggest needs exploring is other subreddits
1075 focused on natural landscapes, such as "r/AmatureEarthPorn", which do not restrict uploads to high-
1076 quality images and therefore may have greater representation of landscapes from a wider
1077 demographic.

1078 There are however several caveats to geocoding Reddit post locations. First, the extracted location
1079 name from the named-entity recognition may not be correct due to ambiguity in the text, spelling or
1080 language differences, or capitalizations (Goyal et al. 2018). Given that posts on "r/EarthPorn" are
1081 predominantly in the English language, this may not have been a significant issue in our analyses.
1082 The issue with multi-part names being extracted to a single word place name means that the finer
1083 spatial scale of the location is lost. The rules of the subreddit specifies that place names included in
1084 the post title should be as specific as possible. However, the named-entity recognition method often
1085 identified the location as the regional (i.e. state) or country part of the place name, losing the finer
1086 detail of the image's location. Though the named-entity recognition method can correctly recognise
1087 and extract places names with multiple parts (e.g. "Ocean Beach", "San Francisco" was correctly
1088 identified), for many multi-part place names the finer location detail can be lost. For example, "Mt.
1089 St. Helens, Washington" was extracted as "Washington". The automated extraction of the landscape
1090 image place name presented here may be best suited for generalising large-scale distributions.
1091 However, as the Reddit posts normally contain specific location details in their titles, studies that
1092 wish to assess spatial distribution on a finer scale may find success in manually extracting the place
1093 name.

1094 Second, the high number of available geocoding algorithms, as well as the potential for ambiguity in
1095 the named entity locations extracted from the Reddit comments, can introduce errors in the
1096 geocoded results (McDonald et al. 2018). For example, there are multiple locations globally named
1097 "Portland"; without more context the geocode algorithm may not correctly code the location. Third,
1098 though the geocoding method can provide a latitude and longitude with a high spatial accuracy,
1099 when geocoding is based on a general location name, the location will be plotted to a single point
1100 within that region. For example, multiple photographs taken in completely different areas of the

1101 Badlands National Park, US, all containing "Badlands National Park" in their title, will all be
1102 aggregated to the same point location. Furthermore, though this method was successful on posts to
1103 "r/EarthPorn", other subreddits may not stipulate that a location must be present in the text. We
1104 suggest that future studies using Reddit data for spatial analysis should consider methods for
1105 reducing geocoding inaccuracies (McDonald et al. 2018). Another possible source of geocoding a
1106 posts location is the Google Cloud Vision API which can estimate the location of an image; however
1107 this process is currently only capable of locating popular sites.

1108 Due to the limitations of geocoding Reddit posts, we do not recommend using posts from Reddit to
1109 assess the spatial variation of CES in a similar manner to those from Flickr, Twitter or Instagram (e.g.
1110 Graham and Eigenbrod 2019; Chen et al. 2020). Instead, one potential method for getting CES data
1111 for a location without the need for geocoding posts is searching for a given name place alongside
1112 other keywords or within a subreddit. This method has previously been used in CES studies from
1113 Flickr, for example Thiagarajah et al. (2015) searched Flickr for photographs based on the place
1114 names of four mangrove sites in Singapore, while Roberts (2017) queried Twitter posts for any
1115 containing the names of urban green spaces in Birmingham, UK. Here, we showed that searches for
1116 Reddit posts with a relevant study site as a key word provides a relatively large dataset across spatial
1117 scales and locations. Though we have demonstrated that Reddit data has the potential for spatial
1118 studies, we acknowledge these limitations do restrict the use of Reddit's data to assess spatial
1119 variations in CES and therefore suggest that Reddit posts are more suited to generalising CES based
1120 on a given search criteria e.g. a place name or specific activity. However, these limitations do not
1121 hinder the use of data for studies that assess CES through content analysis and textual analysis.

1122 We have shown that photographs associated with hiking from both Reddit and Flickr can both be
1123 used in the same image content analysis techniques, thus illustrating their potential for CES studies
1124 which use content analysis of images, without additional spatial analysis (e.g.Thiagarajah et al.
1125 2015). Oakes and Farrow (2006) demonstrated that words with the highest percentage contribution
1126 of the total $x^2$ value, relative to the other words in the set, best highlight the differences in two
1127 groups of words. Here, the small number of labels contributing to a high percentage of the total
1128 $x^2$ value indicates that, in general, many images contain similar scenes, but the difference between
1129 the two sites is driven by a small number of features identified with the Google Cloud Vision API. The
1130 differences between the two sites may be reflected in the user's motivations for undertaking hiking.
1131 As the reasons to undertake hiking are multifaceted (Wilcer et al. 2019), the difference in
1132 demographics between users of Reddit and Flickr suggests they may be undertaking hiking or
1133 uploading images to each site for different reasons. For example, results from our subset of images
1134 suggest that Reddit users are more likely to participate in hiking for physical activity and dog walking,
1135 whilst Flickr users are more likely to undertake hiking to access aesthetic views.

1136 We have demonstrated that, as the contents of images from Reddit and Flickr can provide
1137 essentially the same information about CES, Reddit may be a valuable additional source of data for
1138 assessing aesthetic landscape qualities (e.g. Oteros-Rozas et al. 2018) or recreational preferences
1139 (e.g. Gosal et al. 2019; Lee et al. 2019). The difference in contents may also be down to the
1140 motivations to upload to each platform. Kipp et al. (2017) found that Flickr users have multiple
1141 motivations for uploading photographs including wanting to get an opinion on their photographs
1142 and because they have an interest in a particular subject. However, as one of the main features of
1143 Reddit is the ranking of posts through user votes (Duggan and Smith 2013), further work should be
1144 undertaken to assess whether the relative motivations for uploading to Reddit are similar to other
1145 social media sites. Furthermore, our searches were only carried out in the English language, and

1146 therefore may introduce bias into the conclusion drawn about the motivations for undertaking
1147 different recreational activities.

1148 Comparison of image content from uploads to the subreddits "r/EarthPorn" and "r/BotanicalPorn",
1149 which focus on photographs of different aspects of nature, demonstrated distinctions between the
1150 two - and therefore provide unique sources of data for assessing the role of different aspects of
1151 nature to CES. Building on this, "r/WaterPorn" and "r/DesertPorn" may help to provide a robust
1152 dataset for untangling the contributions of geodiversity to CES by providing unique insights into
1153 peoples' opinions on abiotic features (Fox et al. 2020a). Furthermore, subreddits are not just useful
1154 for assessing aesthetic CES, but can also provide a large source of data for spirituality and recreation.
1155 There is a far larger range of subreddits available than accessed here, each with a unique theme that
1156 can help to understand CES, for example "r/Travel" (a discussion board for travel) could be a useful
1157 source of data for understanding the links between tourism and CES and "r/CityPorn" (images of
1158 cityscapes and urban areas) may help to investigate urban ecosystem services, although this may
1159 require some content filtering to remove purely architectural images. As our keyword searches
1160 return significantly more text-based posts than images, researchers should familiarise themselves
1161 with the different subreddit as potential sources of images, for example, titles of posts in
1162 "r/EarthPorn" generally do not contain words like "landscape" or "view" and would therefore not be
1163 returned through a keyword search looking for images relating to aesthetics. The results presented
1164 here demonstrate that Reddit has the potential to be a significant source of image data and may be
1165 beneficial to CES studies that incorporate content analysis.

1166 Studies can also use textual metadata to assess CES (e.g. Roberts 2017; Hale et al. 2019; Johnson et
1167 al. 2019). Flickr images tend to have description metadata that the uploader provides, which has
1168 been demonstrated to be useful in textual analysis such as sentiment analysis (Brindley et al. 2019)
1169 or eliciting information on CES from the text (Hale et al. 2019). A disadvantage of photographs
1170 uploaded to Reddit is that images do not have an equivalent description by the uploader, therefore
1171 we only compare the comment sections of the two websites. As many posts on Reddit have
1172 comments and because Reddit is a discussion-based platform, this large online database may help to
1173 understand the opinions of thousands of individuals. As the perception of the CES can only be drawn
1174 from those that comment (Dai et al. 2019), having a larger number of unique individuals interacting
1175 with CES related posts may enable the results to be generalised to the wider population and
1176 therefore better help to inform policy, planning and management (Dunkel 2015). Here, the text
1177 comments from the two sites vary regarding the sentiment expressed, with Flickr images having a
1178 more positive associated sentiment score, but also a large variability within the score. The subset
1179 analysed here also showed very few negative comments on Flickr, whilst on Reddit a negative
1180 sentiment was more frequently expressed. Moreover, the actual text contained within the
1181 comments differs between the two sources, with comments on Flickr tending to be more general
1182 appraisals of the photograph, while Reddit comments are more often a discussion around the image
1183 themselves, thus potentially providing richer information on the users' perspective of CES. Having
1184 access to a wider range of opinions, both positive and negative, may help to better generalise
1185 attitudes to CES.

1186 As Reddit is designed to be a discussion-based forum it may contribute to richer information on the
1187 users' perspective of CES. For example, the "r/Spirituality" subreddit encourages users to contribute
1188 to the discussion of any aspects of spirituality regardless of religion or ideology, thus providing the
1189 potential to assess the opinions of people from a wide range of backgrounds. Furthermore, Reddit
1190 comments can be longer than most other social media sites (e.g. Twitter has a 280-character limit
1191 and Instagram has a 300-character comment limit) and therefore a user can discuss their opinions in

1192 greater detail (Gkotsis et al. 2017). The discursive nature of Reddit provides researchers a unique
1193 opportunity to assess which aspects of a certain image or video people appreciate. There is also
1194 scope for this interactive and discussion-based platform to be used in experimental studies in which
1195 researchers post content and monitor feedback. Though as with all social media-based studies, we
1196 recommended that the ethics of these studies be discussed in further detail. We suggest that Reddit
1197 data is particularly useful for studies that wish to analyse users' comments in conjunction with the
1198 metadata available for each image for a more robust assessment of CES.

1199 For studies carrying out image content or textual analysis we suggest that combining Reddit data
1200 alongside other sources of data, would be useful in CES because (1) images and text from Reddit can
1201 provide comparable data used to assess aspects of CES; (2) Reddit potentially contains additional
1202 data previously overlooked; (3) they have different geographical biases (e.g. Reddit to North
1203 America, Flickr to Europe and Weibo to Asia). We therefore suggest that a more holistic approach of
1204 assessing CES would be to include cross-platform analysis including multiple sources (Retka et al.
1205 2019). However, we note that Reddit may not be suitable for integrating into studies assessing
1206 spatial variations in CES. Data integration, the bringing together of data from multiple sources, could
1207 be implemented to allow data from social media sites to be analysed as a complete unit. Data
1208 integration methods, which control for differing biases and sizes of datasets, have been successfully
1209 used in other scientific fields such as species distribution modelling (Issac et al. 2019) and those
1210 using satellite imagery (Aires 2014). As accessing data from Reddit requires a similar skill level as
1211 accessing datasets from other social media websites, data integration of these multiple sources is
1212 feasible. The tools and software used in this manuscript make these datasets more accessible and
1213 reproducible for non-data scientists and enable us to start to bridge the gap in integrating multiple
1214 sources. We therefore recommend that CES and wider environmental science studies make use of
1215 these tools to include the vast amount of data from Reddit alongside other social media data sources
1216 in their future studies.

## 5.0 Conclusion

1218 We have demonstrated that posts from Reddit can be used in commonly applied CES assessment
1219 methods, such as image content analysis and textual analysis, which leverage the power of big data
1220 from social media sites. The results from this study show that Reddit can provide a large source of
1221 data similar to Flickr. However, the posts available on Reddit are not geolocated and the geocoding
1222 of a post's location has several limitations meaning that Reddit is not as suited to assessing the
1223 spatial variation of CES as other social media sites. The large quantity of data available on Reddit is
1224 most appropriate for assessing general trends in CES through image content analysis and textual
1225 analysis. The discursive nature of Reddit provides a unique opportunity to assess a wide range of CES
1226 including recreational activities, aesthetic views, spirituality and culture. We argue that Reddit
1227 should be more widely considered as a useful source of data for CES studies and we hope that this
1228 paper sets a precedent for including big datasets from Reddit in future studies.

1229