

1           **A guide to using the Internet to monitor and quantify the wildlife trade**

2

3   Oliver C. Stringham<sup>1,2</sup>, Adam Toomes<sup>1</sup>, Aurelie M. Kanishka<sup>1,3</sup>, Lewis Mitchell<sup>2</sup>, Sarah Heinrich<sup>1</sup>, Joshua V.

4   Ross<sup>2</sup>, Phillip Cassey<sup>1</sup>

5           1. School Biological Sciences, University of Adelaide, SA 5005, Australia

6           2. School of Mathematical Sciences, University of Adelaide, SA 5005, Australia

7           3. Fenner School of Environment and Society, The Australian National University, Canberra, ACT

8           2601, Australia

9

10   Corresponding author: Oliver Stringham, [oliverstringham@gmail.com](mailto:oliverstringham@gmail.com)

11

12   **Keywords:** big data, e-commerce; Internet; pet trade; social media; web-scraping

13 **Abstract**

14 The unrivalled growth in e-commerce of animals and plants presents an unprecedented opportunity to  
15 monitor wildlife trade to inform conservation, biosecurity, and law enforcement. Using the Internet to  
16 quantify the scale of the wildlife trade (volume, frequency) is a relatively recent and rapidly developing  
17 approach, which currently lacks an accessible framework for locating relevant websites and collecting  
18 data. Here, we present an accessible guide for Internet-based wildlife trade surveillance, which uses a  
19 repeatable and systematic method to automate data collection from relevant websites. Our guide is  
20 adaptable to the multitude of trade-based contexts including different focal taxa or derived parts, and  
21 locations of interest. We provide information for working with the diversity of websites that trade  
22 wildlife, including social media platforms. Finally, we discuss the advantages and limitations of web data,  
23 including the challenges presented by trade occurring on clandestine sections of the Internet (e.g., deep  
24 and dark web).

25

26 **Background**

27 *Introduction*

28 The wildlife trade is an influential driver of species endangerment; as well as spreading invasive species  
29 and diseases, and provisioning criminal activity (‘t Sas-Rolfes et al. 2019). Wildlife trade occurs across a  
30 variety of physical and virtual settings, including ‘brick and mortar’ stores, wet markets, and digital  
31 platforms on the Internet (e.g., Alfino & Roberts 2018). Reliable data on the quantity and composition of  
32 the wildlife trade (legal and illegal) is vital for informing decisions about conservation, biosecurity, and  
33 law enforcement; and developing human behavior change campaigns. Yet this data is rarely collected, or  
34 is difficult to obtain (Regueira & Bernard 2012; Eskew et al. 2020). In recent years, the Internet has  
35 played an increasingly important role in facilitating the wildlife trade (Siriwat & Nijman 2020).

36

37 To date, researchers have used data from the Internet in various ways to inform wildlife trade research,  
38 and assist in practical management; including law enforcement. These studies have generally been small  
39 in scale (i.e., monitoring one or few websites), but have nonetheless revealed the utility of the Internet  
40 to describe different aspects of the wildlife trade. In the context of conservation, classified websites  
41 have been used to estimate intensity of trade, and support increases in the legal protection of high-risk  
42 species (Rowley et al. 2016). For biological invasions, online pet stores have been used to inventory non-  
43 native species (Stringham & Lockwood 2018). Lost and Found websites have been used to estimate  
44 propagule pressure, a major determinant of non-native establishment probability (Cassey et al. 2018),  
45 for commonly held exotic pets (e.g., turtles: Kikillus et al. 2012). In terms of assisting law enforcement,  
46 listings from online classifieds have been used to quantify the illegal trade (Ye et al. 2020), and social  
47 media websites have been used to track the intensity of legal and illegal trade (Jensen et al. 2019).

48

49 As the volume and frequency of wildlife trade increases over the Internet, having a unified method for  
50 using the Internet to obtain data on the wildlife trade becomes more critical for researchers. However,  
51 such a methodology, or 'guide', does not currently exist. By outlining a guide with repeatable steps we  
52 hope to facilitate reproducible methods for using the Internet as a data source (including finding  
53 websites, data collection and curation). Further, a guide can serve as a primer for investigating  
54 unexplored contexts of the trade, including new locations and different focal taxa, or emerging trade in  
55 derived parts and commodities.

56

57 Here, we present an accessible guide to using the Internet to gather data on the wildlife trade. We  
58 developed the methodology through our collective knowledge of working with web data and the wildlife  
59 trade, combined with the methods used in prior published studies. Our goal is for this guide to be used  
60 by scientists, NGOs, government agencies, and other interested parties, who wish to utilize the Internet  
61 as a source of data on the wildlife trade.

62

### 63 *Structure of the Internet*

64 The Internet (i.e., the World Wide Web or simply the Web) is categorized into three distinct 'layers': the  
65 surface web, the deep web, and the dark web (Figure 1; Bergman 2001). Each layer differs in one of two  
66 factors: whether it's accessible without logging in or invitation (i.e., is publicly viewable) and whether it  
67 is indexed by a search engine (i.e., will appear as a result in a search engine). The surface web includes  
68 any website that is publicly viewable and is indexed by search engines (e.g., e-commerce websites). The  
69 deep web includes websites or online content that require either logging in or an invitation to view (e.g.,  
70 social media, private messaging apps). Some deep web sites may be indexed by a search engine (e.g.,  
71 public Facebook or Twitter posts), while others may not (e.g., WhatsApp). The dark web contains  
72 purposefully hidden content that requires specialized software to access, requires either logging in, or

73 an invitation to view, and is not indexed by any search engine (Chen 2011; CRS 2017). The degree to  
74 which researchers can find relevant wildlife trade content on the Internet will be influenced by how  
75 'findable' a website/content is (e.g., can a search engine find it). Further, the ethical considerations of  
76 collecting data will in part depend on how 'accessible' the website is (e.g., is deceit or limited disclosure  
77 required to gain access to the content; Section 4.1).

78

### 79 *What data is available?*

80 Data availability on wildlife trade varies by website and even within a website (Appendix S1; Toivonen et  
81 al. 2019). On a basic level, online advertisements (i.e., listings or posts) are provided in the form of text,  
82 pictures and/or videos. Foremost, the name of the species, taxa, or derived product traded is usually  
83 stated. Characteristics of the traded taxa or product can include quantity (number, size, volume), age,  
84 sex, size, color, morph, and provenance (domestic bred or wild caught/harvested). The physical location  
85 of the advertisement (i.e., city) and metadata on the advertisement itself, such as the number of page  
86 views and username of the trader, may be provided. Further, the current purpose for which the wildlife  
87 is being used (pet, medicinal, food, etc.) along with the rationale for trading the wildlife (e.g., profit,  
88 lifestyle change) can sometimes be ascertained from advertisements with open text fields. These  
89 attributes may aid in understanding motives around wildlife trade participation or consumption (i.e.,  
90 Conservation Culturonomics: Ladle et al. 2016).

### 91 **Guide to using the Internet to monitor and quantify the wildlife trade**

92 We specify our guide in six steps (Figure 2): (1) defining the scope and purpose of the project; (2) finding  
93 candidate websites; (3) selecting target websites to monitor; (4) collecting and storing data from  
94 websites; (5) cleaning data; and (6) analysis. Here, we detail the process of each step leading up to  
95 analysis. We generalize this guide for websites found in any layer of the Internet (including social media)

96 and discuss how to adapt this guide to different languages and countries. Further, in Figure 3, we  
97 provide two hypothetical case studies to accompany and contextualize each step of the guide. For more  
98 generalized frameworks on working with social media and online news outlets, refer to Toivonen et al.  
99 (2019) and Sonricker Hansen et al. (2012), respectively.

100

### 101 *1. Defining the scope and purpose of the project*

102 At a minimum, it is essential to decide which species, taxa, and/or derived products are of interest, the  
103 location(s) of interest, and the timeframe for data collection (i.e., one-time snapshot, versus ongoing  
104 monitoring for months to years). Further, considering what type of website (Appendix S1) or layer of the  
105 Internet may be appropriate. On a practical note, the research questions that can be answered will be  
106 influenced by the data available on the Internet. Thus, there will likely need to be some exploration of  
107 the websites and the kind of data they provide (Steps 2 – 3). Examples of project aims include:  
108 quantifying the trade in parrots in different regions of China (Ye et al. 2020); investigating the sale of  
109 pangolin-leather boots in the US (Heinrich et al. 2019); exploring the social network structure of sellers  
110 of horticultural orchids (Hinsley et al. 2016).

111

### 112 *2. Finding candidate websites where specific taxa and wildlife products are traded*

113 Here, we detail a method to find candidate websites (e.g., e-commerce sites, forums) by using search  
114 engines. Finding relevant social media content requires special considerations, which we detail in  
115 Section 2.4. Outside of the search engines, other approaches to finding candidate websites and choosing  
116 target websites (Section 3) include: interviewing a specific community of practice (e.g., reptile keepers  
117 and traders); or collaborating with other researchers actively engaged in online wildlife-trade  
118 monitoring (e.g., governmental agencies or NGOs). It is important to note that the Internet is transient:  
119 traders go out of business and new ones emerge. Thus, websites found at one point in time can differ in

120 composition and function if surveyed later. If the goal is long-term monitoring, we suggest revising the  
121 list of current relevant websites at regularly timed intervals.

122

## 123 2.1 Surface Web

124 For the surface web, finding candidate websites involves three steps: (1) defining keyword phrases to  
125 search; (2) using a search engine to perform searches; and (3) classifying the relevance of each search  
126 result. This part of the methodology is akin to the process of finding relevant scientific papers in a  
127 systematic review or meta-analysis (i.e., PRISMA methodology: Koricheva et al. 2013). However, instead  
128 of searching the scientific literature, the Internet is searched (via search engines), and not all candidate  
129 results will be used for data collection (Section 3).

130

### 131 2.1.1 Defining search phrases

132 Search phrases are composed of a combination of relevant keywords. We recommend developing a  
133 suite of keywords for each target taxa (e.g., species name, common name, product name), type of  
134 websites (Appendix S1), and location of interest. Other useful keywords include adding the terms “for  
135 sale” or “buy”. Example search phrases may be: “snakes for sale Australia”, “marine fish forum USA”, or  
136 “orchid store UK” (See Appendix S2 for a detailed example). These search phrases should be in the  
137 language(s) written in the location of interest. There may be a need to refine keywords after exploratory  
138 investigation of search engine results. In particular, there may be trade names (i.e., names for species or  
139 taxa used in the wildlife trade community, but not commonly used among scientists), local/regional  
140 names, or names of breeds, morphs, and mutations (e.g., Lyons and Natusch 2013), which are not  
141 captured in the initial formulation of search phrases.

142

143 2.1.2 Using search engines to perform searches

144 Search engines (e.g., Google) use proprietary algorithms to return a list of URLs (i.e., website addresses)  
145 when a search phrase is input. Search engine algorithms consider the relevance of the keywords, the  
146 popularity of the website (i.e., the number of page views), and, increasingly, the location of where the  
147 search occurs (Langville and Meyer 2011). The results from a search engine are expected to change at  
148 any point in time for a number of reasons including: changes to the search engine algorithm, changes to  
149 website popularity metrics, the emergence of new websites, and a change in the location of where the  
150 search is performed. Once a keyword phrase is searched, the search engine will likely return millions of  
151 URLs per phrase. We recommend choosing a cutoff point that balances the quality of search results with  
152 search effort (Appendix S3). Because search engines can use the user's location to provide personalized  
153 results (e.g., Google: <https://policies.google.com/technologies/location-data>), extra steps must be taken  
154 to ensure that the search engine provides location- and language-relevant results. One way to control  
155 the location is to use advanced search features (e.g., [https://www.google.com/advanced\\_search](https://www.google.com/advanced_search)), which  
156 allows the researcher to specify which country and languages to restrict a search to. In addition, using a  
157 VPN (virtual private network) may alleviate location issues. For more information on search engines, see  
158 Appendix S3.

159

160 2.2 Deep web and dark web (non-social media)

161 Websites on the deep web indexed on search engines will be findable using the same approach outlined  
162 for the surface web (e.g., private forums). Currently, there are no generalizable or automated methods  
163 for locating deep web content or websites not indexed by search engines (e.g. WhatsApp, WeChat,  
164 other private messaging apps) nor dark web content, outside of expert consultation or interviewing  
165 communities of practice. While some algorithms exist for querying deep websites (e.g., Liakos et al.  
166 2016), the actual implementation of these algorithms as web crawlers must be tailored for each



167 individual instance, and require unique login details. This severely limits any large-scale monitoring  
168 efforts.

169

## 170 2.3 Classifying search engine results

171 After obtaining URLs from search results, each will need to be categorized as relevant or irrelevant.

172 Relevance is subjective and we recommend defining inclusion/exclusion criteria depending on the scope  
173 and purpose of the study. One obvious inclusion criterion is whether the target taxa is traded on the  
174 website. Another criterion can be the type of transaction that occurs on the website. Specifically, on the  
175 Internet, there are varying levels of ‘directness’ of trade. For instance, some e-commerce companies will  
176 ship live animals or products to a customer’s doorstep (e.g., pet stores: Holmberg et al. 2015) and there  
177 are less direct websites that facilitate the transaction of selling wildlife online, but leave it up to the  
178 individuals in the transaction to conduct the exchange (e.g., classifieds: Sung & Fong 2019).

179

## 180 2.4 Social media

### 181 2.4.1 Types of social media content

182 Social media websites vary in structure and format (Appendix S1; Toivonen et al. 2019). For our  
183 purposes, we categorize content found on social media websites into two categories: consolidated and  
184 unconsolidated. The differences between each category will influence how researchers find relevant  
185 social media content related to wildlife trade. Consolidated social media content includes ‘groups’  
186 dedicated to a particular purpose (e.g., ornamental orchid traders) where users can share content that is  
187 only viewable by other group members (e.g., Facebook groups). Social media ‘groups’ function similarly  
188 to forum websites. Unconsolidated social media content consists of users posting to the social media  
189 platform at large or to a group of followers. Twitter, for example, is mostly public, where all ‘Tweets’

190 (i.e., posts) are viewable by all users. Some social media websites, such as Facebook, may have both  
191 consolidated and unconsolidated content.

192

#### 193 2.4.2 Finding social media content

194 Social media websites have their own internal search engine, which searches through content of the  
195 specific social media site. Thus, for consolidated social media content, we recommend adapting our  
196 approach outlined for the surface web (i.e., using search phrases; Section 2.1) for internal search  
197 engines to find relevant social media ‘groups’ (e.g., Siriwat & Nijmans 2020). These ‘groups’ can then be  
198 classified by their relevance (Section 2.3) and considered for monitoring (Section 3). For unconsolidated  
199 social media content, we recommend simply using the internal search engine to search for relevant  
200 posts. The posts returned by the search engine become the data itself (e.g., Xu et al. 2019), where the  
201 ‘classification’ and ‘selection’ steps of this guide are skipped (Section 2.3; Section 3). Importantly, many  
202 social media users utilize hashtags (denoted by the number sign: #), which are user-generated tags  
203 relating to the post’s content (e.g., #ivory). Thus, for social media sites, determining what hashtags are  
204 used for a specific context of wildlife trade may yield more relevant search results than keyword phrases  
205 for both consolidated and unconsolidated social media content (e.g., Morgan & Chng 2018).

206

207 APIs (Application Programming Interfaces) may be available for some social media websites, which may  
208 allow for ‘bulk’ searches (i.e., more than one search at once) and streamlined data collection (Section 4).  
209 Filters may be available in advanced search options of internal search engines or in APIs to restrict  
210 search results to certain countries and languages. Finally, social media companies have allowed users to  
211 adjust their privacy settings, so that only their ‘followers’ or a pre-selected group of users can view their  
212 posts. Importantly, content with privacy restrictions may be hidden from the internal search engine or  
213 API results.

214

215 *3. Selecting target sites to monitor*

216 After obtaining the list of candidate websites, the next step is to select which websites to collect data  
217 from (i.e., target websites). This step of the framework is the most subjective and therefore some level  
218 of justification and transparency should be provided when choosing target websites. To make informed  
219 decisions on selecting target websites, metadata on candidate websites can be collected. For surface  
220 websites, one metadata attribute is web traffic statistics, which includes information such as the  
221 number of page views per month (see Appendix S4 for more information). In addition, for any website,  
222 researchers can calculate the average number of posts or listings per day and use this as a proxy for  
223 popularity. Ultimately, researcher discretion is needed to choose target websites, because measures of  
224 website metadata are not available for all candidate websites and project relevance is not always  
225 straightforward to quantify. The number of target websites chosen will vary based on the project aim(s)  
226 and the resources available to collect and clean data (Sections 4 and 5). Again, expert opinion and  
227 communities of practice can provide opinions on what websites are most relevant.

228

229 *4. Collecting data from websites*

230 Data collection can occur in one of two main ways: manual or automated. Manual data collection  
231 involves visiting the website and recording what taxa/product is being traded along with desired  
232 associated attributes (e.g., price, location). Automated data collection involves constructing “web  
233 scrapers” to visit the website and extract desired relevant information (Figure 4; Singrodia et al. 2019).  
234 Web scrapers organize the contents of a website into a structured tabular format (for more information  
235 on web scrapers and data storage see Appendix S5). Since each website differs in its underlying  
236 structure, custom web scrapers need to be coded for each website individually. A few highly visited  
237 websites may have APIs that allow for easy collection of data; this is more likely to be the case for social

238 media websites (e.g., Twitter; Toivonen et al. 2019). Choosing manual or automated data collection will  
239 depend on how long and how often data is being collected, as it takes technical expertise and time to  
240 build web scrapers, which may not be necessary if the number of target websites is small and the data  
241 collection window is short (e.g., Heinrich et al. 2020). These methods of data collection apply to  
242 websites and content on the deep web (including social media) and dark web, as long as researchers  
243 have access to the website/content (e.g. Cunliffe et al. 2019).

244

#### 245 4.1 Ethical considerations of collecting data from the Internet

246 Ethics approval is required to collect information from the Internet, especially when personally  
247 identifiable material is collected, including, but not limited to, social media sites (Zimmer 2010). Care  
248 should be taken to ensure de-identified information is used for analyses and subsequent publication  
249 (Harriman & Patel 2014; Sula 2016). Furthermore, ethics approval for collecting data from any deep or  
250 dark web sites will include obtaining a login or approval to join (Tai et al. 2012); since deceit or limited  
251 disclosure of research aims may be required. Also, automated data collection processes (i.e., web  
252 scraping) currently encompasses a legal gray area (Zamora 2019), and thus we encourage researchers to  
253 acquire ethics approval prior to using them. For specific recommendations of ethical practice, refer to  
254 Appendix S6.

255

#### 256 *5. Data Cleaning*

257 Data cleaning involves curating each listing (i.e., post or advertisement) for attributes that could not be  
258 automatically extracted, but are required for the analysis, such as species name, quantity, price, or  
259 location. Data cleaning is often a tedious and time-consuming task (Freitas & Curry 2016) and could  
260 possibly be the most time-consuming part of the entire project. The amount of cleaning required will  
261 depend on the structure of the website and will vary by individual website (Appendix S1). For instance, a

262 website may have a separate field for species names, while another may just have one free form open  
263 text box where the user can type anything. Our experience with websites involving the wildlife trade is  
264 with the latter, which takes substantially more time to clean. If collecting data manually, simultaneously  
265 cleaning data during collection is possible and likely desirable.

266

### 267 5.1 Resolving species names

268 Resolving the species name in a listing or post is one of the most important aspects of data cleaning.

269 Some pet stores and specialist classifieds websites explicitly state the scientific name while other sites

270 may mention common names, trade names, or simply supply a photo. For all practical purposes,

271 identifications down to the rank of species are needed for effective action on conservation, biosecurity,

272 and crime (Rhyne et al. 2012). Therefore, we recommend identifying the taxa to the most specific

273 taxonomic level as possible. If pictures are provided, taxonomic experts can aid in species identification.

274 Yet, pictures may be too poor in quality to properly identify species. In some instances, online traders

275 may simply not provide enough information in the listing to identify to species level.

276

277 If monitoring many species, we recommend relating the species/taxa name to a taxonomic database

278 (e.g., GBIF 2020). Doing so will facilitate conformation to taxonomic names by avoiding synonyms and

279 misspellings (Gallagher et al. 2020). In addition, it will enable the researcher to easily acquire upstream

280 taxonomy (e.g., Family and Order). We recommend the R package *taxize*, which automates the

281 gathering of upstream taxonomy if supplied a scientific name or database identifier for the taxa of

282 interest (Chamberlain & Szocs 2013).

### 283 **Advantages and caveats of web data**

284 The ease of gathering data from the Internet is the main advantage compared to surveying physical  
285 markets or stores, especially if using automated data collection techniques (i.e., web scrapers).  
286 Furthermore, using the Internet could potentially allow for a more complete picture of the trade both  
287 spatially and temporally than would normally be possible for researchers or organizations who have  
288 limited resources for traditional surveys. However, the Internet is not a panacea for monitoring the  
289 wildlife trade and relying on the Internet for data on the wildlife trade has several disadvantages. First,  
290 not all trade occurs nor is observable online (e.g., bushmeat trade; McNamara et al 2019). The degree to  
291 which trade occurs online will depend on the type of trade (i.e., pet, derived products, food, etc.), the  
292 taxa, the country or culture in question (i.e., Internet use varies by country; Pew Research Center 2016),  
293 and possibly by target/consumer group. To the best of our knowledge, there are no estimates of the  
294 ratio of physical versus online trade for any context. Another downside is that it is difficult, if not  
295 impossible, to verify the validity of online listings of wildlife (i.e., fake or scam versus genuine  
296 advertisements). Supplementing data collected online with physical surveys is a more holistic approach  
297 that may be more impactful when considering applied outcomes (e.g., Rowley et al. 2016).

### 298 **Considerations for the Deep and Dark Web**

299 Currently, wildlife trade on the surface web and indexed deep web (e.g., social media) is extremely  
300 abundant (Sung & Fong 2018; Xu et al. 2020; IFAW 2018). The unindexed deep web, such as private text  
301 messaging apps (e.g., WhatsApp; Facebook Messenger), has remained relatively unexplored until  
302 recently (e.g., Sanchez-Mercado et al. 2020; Setiawan et al. 2019), thus the extent of trade is unknown.  
303 Given the ease of access of private messaging apps and the anonymity they provide, we hypothesize  
304 that trade is also abundant on the unindexed deep web. The dark web remains elusive. While there is  
305 evidence that wildlife is not traded on common dark web marketplaces, this does not discount the

306 potential for trade to be occurring elsewhere on the dark web (Roberts & Hernandez-Castro 2017).  
307 Further, future policies enacted in response to conservation, and the criminological and welfare  
308 concerns of wildlife trade, may shift the balance of where wildlife trade occurs on the Internet (Roe et  
309 al. 2020). Specifically, new regulations or improved enforcement of illegal trade can unintentionally  
310 drive trade away from the open and indexed deep web to the unindexed deep web and dark web  
311 (Nijman 2020; Appendix S7), ultimately making it more difficult for researchers to locate wildlife trade  
312 online.

313

314 Websites and content on the deep and dark web present several challenges for researchers. First,  
315 finding websites that trade wildlife on the unindexed deep and dark web is difficult because they are not  
316 accessible by search engines. This is an unfortunate reality for researchers, but reflects an intentional  
317 design to keep this information private. Further, obtaining access to deep and dark websites often  
318 requires researchers to use deceit for successful infiltration. Using deceit requires ethics approval  
319 (Section 4.1) and infiltration requires skills and training that conservation researchers may not have  
320 (e.g., remaining anonymous). Thus, interdisciplinary collaborations with criminologists, sociologists,  
321 computer scientists, and agencies that specialize in infiltrating and tracking cybercrime (e.g., law  
322 enforcement) will be beneficial.

323

#### 324 **Automated data cleaning**

325 Automated data cleaning of wildlife trade web data has not been attempted, however, there is potential  
326 from computer science subfields, such as machine learning, to help with cleaning messy data (Lamda et  
327 al. 2019; Norouzzadeh et al. 2020). Tools relevant to wildlife trade websites are image classification and  
328 text classification (e.g., Deep learning and Natural Language Processing: Di Minin et al. 2018; Silge &

329 Robinson 2020), which can potentially use images or text to identify certain attributes of a given listing,  
330 such as the species being traded. However, there is a paucity of applications of these tools/fields to web  
331 data of the wildlife trade specifically (Xu et al. 2019). Importantly, underlying all of these machine-  
332 learning tools are training sets, which are a representative sample of listings that have been manually  
333 classified by a human(s) for the machine-learning algorithm to use (Lamda et al. 2019). The larger the  
334 training set, the more likely the machine-learning model will perform better (Norouzzadeh et al. 2020).  
335 Importantly, there will always be the need for human data cleaning and labelling. One major barrier to  
336 successful implementation of automated data cleaning tools for wildlife trade data is the number of  
337 species involved in the trade, where research contexts can encompass hundreds to thousands of species  
338 and wildlife parts/derivatives (e.g., Humair et al. 2015).

### 339 **Conclusions**

340 As more of the global human population shifts to using the Internet, and as ethical and disease concerns  
341 of physical markets arise (Roe et al. 2020), the online trade of wildlife is poised to increase. Thus, the  
342 Internet is, and will continue to be, an invaluable source of data (Lavorgna 2014). Despite the limitations  
343 of data collected from the Internet, there are vast opportunities to inform conservation, biosecurity, and  
344 law enforcement objectives. Current strategies of researchers using small-scale monitoring (i.e., one or  
345 few websites) should continue to provide insight into specific taxa/products contexts (Sung & Fong  
346 2018). With the development of machine learning tools to clean 'messy' web data, there will be the  
347 possibility of creating large-scale (i.e., for many websites) automated systems to detect illegal trade to  
348 help inform law enforcement and conservation efforts. Likewise, early risk-screening and rapid response  
349 systems may be possible for invasive species (e.g., Marshall Meyers 2020; Suiter & Sferrazza 2007),  
350 especially for 'exotic' animals and ornamental plants whose online trade is commonplace (Lockwood et



351 al. 2019; Lenda et al. 2014). Regardless of the ultimate application, our guide can serve as a primer and  
352 starting point to establishing research agendas related to wildlife trade occurring on the Internet.  
353

354 **References**

- 355 Alfino S, Roberts DL. 2018. Code word usage in the online ivory trade across four European Union  
356 member states. *Oryx*:1–5. Cambridge University Press.
- 357 Bergman MK. 2001. White Paper: The Deep Web: Surfacing Hidden Value. *Journal of Electronic*  
358 *Publishing* 7. Available from <http://hdl.handle.net/2027/spo.3336451.0007.104>.
- 359 Cassey P, Delean S, Lockwood JL, Sadowski JS, Blackburn TM. 2018. Dissecting the null model for  
360 biological invasions: A meta-analysis of the propagule pressure effect. *PLOS Biology*  
361 16:e2005987. Public Library of Science.
- 362 Chamberlain S, Szocs E. 2013. “taxize - taxonomic search and retrieval in R.” *F1000Research*.  
363 <http://f1000research.com/articles/2-191/v2>.
- 364 Chen H. 2011. *Dark Web: Exploring and Data Mining the Dark Side of the Web*. Springer Science &  
365 Business Media.
- 366 Congressional Research Service, March 10, 2017 , “Dark Web”. Available on  
367 <https://crsreports.congress.gov/product/pdf/R/R44101> (accessed May 2020)
- 368 Cunliffe J, Décary-Hêtu D, Pollak TA. 2019. Nonmedical prescription psychiatric drug use and the  
369 darknet: A cryptomarket analysis. *International Journal of Drug Policy* 73:263–272.
- 370 Di Minin E, Fink C, Tenkanen H, Hiippala T. 2018. Machine learning for tracking illegal wildlife trade on  
371 social media. *Nature Ecology & Evolution* 2:406–407. Nature Publishing Group.
- 372 Eskew EA, White AM, Ross N, Smith KM, Smith KF, Rodríguez JP, Zambrana-Torrel C, Karesh WB,  
373 Daszak P. 2020. United States wildlife and wildlife product imports from 2000–2014. *Scientific*  
374 *Data* 7:1–8.
- 375 Freitas A, Curry E. 2016. Big Data Curation. Pages 87–118 in J. M. Cavanillas, E. Curry, and W. Wahlster,  
376 editors. *New Horizons for a Data-Driven Economy: A Roadmap for Usage and Exploitation of Big*

377 Data in Europe. Springer International Publishing, Cham. Available from  
378 [https://doi.org/10.1007/978-3-319-21569-3\\_6](https://doi.org/10.1007/978-3-319-21569-3_6) (accessed May 13, 2020).

379 Gallagher RV et al. 2020. Open Science principles for accelerating trait-based science across the Tree of  
380 Life. *Nature Ecology & Evolution* **4**:294–303. Nature Publishing Group.

381 GBIF: The Global Biodiversity Information Facility (2020) “What is GBIF?”. Available from  
382 <https://www.gbif.org/what-is-gbif> (accessed May 2020)

383 Harriman S, Patel J. 2014. The ethics and editorial challenges of internet-based research. *BMC Medicine*  
384 **12**:124.

385 Harrison JR, Roberts DL, Hernandez-Castro J. 2016. Assessing the extent and nature of wildlife trade on  
386 the dark web. *Conservation Biology* **30**:900–904.

387 Heinrich S, Ross JV, Cassey P. 2019. Of cowboys, fish, and pangolins: US trade in exotic leather.  
388 *Conservation Science and Practice* **1**:e75.

389 Heinrich S, Toomes A, Gomez L. 2020. Valuable Stones: The Trade in Porcupine Bezoars. *Global Ecology*  
390 *and Conservation*:e01204.

391 Hinsley A, Lee TE, Harrison JR, Roberts DL. 2016. Estimating the extent and structure of trade in  
392 horticultural orchids via social media. *Conservation Biology* **30**:1038–1047.

393 Holmberg RJ, Tlusty MF, Futoma E, Kaufman L, Morris JA, Rhyne AL. 2015. The 800-Pound Grouper in the  
394 Room: Asymptotic Body Size and Invasiveness of Marine Aquarium Fishes. *Marine Policy* **53**:7–  
395 12.

396 Humair F, Humair L, Kuhn F, Kueffer C. 2015. E-commerce trade in invasive plants. *Conservation Biology*  
397 **29**:1658–1665.

398 International Fund for Animal Welfare. 2018. “Disrupt: Wildlife Cybercrime”. Available on  
399 [https://d1jyxxz9imt9yb.cloudfront.net/resource/218/attachment/original/IFAW\\_-  
400 \[\\\_Disrupt\\\_Wildlife\\\_Cybercrime\\\_-\\\_FINAL\\\_English\\\_-\\\_new\\\_logo.pdf\]\(https://d1jyxxz9imt9yb.cloudfront.net/resource/218/attachment/original/IFAW\_-\_Disrupt\_Wildlife\_Cybercrime\_-\_FINAL\_English\_-\_new\_logo.pdf\) \(Accessed May 2020\)](https://d1jyxxz9imt9yb.cloudfront.net/resource/218/attachment/original/IFAW_-_Disrupt_Wildlife_Cybercrime_-_FINAL_English_-_new_logo.pdf)

401 Jensen TJ, Auliya M, Burgess ND, Aust PW, Pertoldi C, Strand J. 2019. Exploring the international trade in  
402 African snakes not listed on CITES: highlighting the role of the internet and social media.  
403 *Biodiversity and Conservation* **28**:1–19.

404 Kikillus KH, Hare KM, Hartley S. 2012. Online trading tools as a method of estimating propagule pressure  
405 via the pet-release pathway. *Biological Invasions* **14**:2657–2664.

406 Koricheva J, Gurevitch J, Mengersen K. 2013. *Handbook of Meta-analysis in Ecology and Evolution*.  
407 Princeton University Press.

408 Ladle RJ, Correia RA, Do Y, Joo G-J, Malhado AC, Proulx R, Roberge J-M, Jepson P. 2016. Conservation  
409 culturomics. *Frontiers in Ecology and the Environment* **14**:269–275.

410 Lamba A, Cassey P, Segaran RR, Koh LP. 2019. Deep learning for environmental conservation. *Current*  
411 *Biology* **29**:R977–R982.

412 Langville AN, Meyer CD. 2011. *Google’s PageRank and Beyond: The Science of Search Engine Rankings*.  
413 Princeton University Press.

414 Lavorgna A. 2014. Wildlife trafficking in the Internet age. *Crime Science* **3**:5.

415 Lenda M, Skórka P, Knops JMH, Moroń D, Sutherland WJ, Kuszewska K, Woyciechowski M. 2014. Effect  
416 of the Internet Commerce on Dispersal Modes of Invasive Alien Species. *PLOS ONE* **9**:e99786.  
417 Public Library of Science.

418 Liakos P, Ntoulas A, Labrinidis A, Delis A. 2016. Focused crawling for the hidden web. *World Wide Web*  
419 **19**:605–631.

420 Lockwood JL et al. 2019. When pets become pests: the role of the exotic pet trade in producing invasive  
421 vertebrate animals. *Frontiers in Ecology and the Environment* **17**:323–330.

422 Lyons JA, Natusch DJD. 2013. Effects of consumer preferences for rarity on the harvest of wild  
423 populations within a species. *Ecological Economics* **93**:278–283.

424 Marshall Meyers N, Reaser JK, Hoff MH. 2020. Instituting a national early detection and rapid response  
425 program: needs for building federal risk screening capacity. *Biological Invasions* **22**:53–65.

426 McNamara J, Fa JE, Ntiamoa-Baidu Y. 2019. Understanding drivers of urban bushmeat demand in a  
427 Ghanaian market. *Biological Conservation* **239**:108291.

428 Morgan J, Chng S. 2018. Rising internet-based trade in the Critically Endangered ploughshare tortoise  
429 *Astrochelys yniphora* in Indonesia highlights need for improved enforcement of CITES. *Oryx*  
430 **52**:744–750. Cambridge University Press.

431 Nijman V. 2020. Illegal trade in Indonesia’s National Rare Animal has moved online. *Oryx* **54**:12–13.

432 Norouzzadeh MS, Nguyen A, Kosmala M, Swanson A, Palmer MS, Packer C, Clune J. 2018. Automatically  
433 identifying, counting, and describing wild animals in camera-trap images with deep learning.  
434 *Proceedings of the National Academy of Sciences* **115**:E5716–E5725. National Academy of  
435 Sciences.

436 Pew Research Center, February, 2016, “Smartphone Ownership and Internet Usage Continues to Climb  
437 in Emerging Economies”. Available on  
438 [https://www.pewresearch.org/global/2016/02/22/smartphone-ownership-and-internet-usage-](https://www.pewresearch.org/global/2016/02/22/smartphone-ownership-and-internet-usage-continues-to-climb-in-emerging-economies/)  
439 [continues-to-climb-in-emerging-economies/](https://www.pewresearch.org/global/2016/02/22/smartphone-ownership-and-internet-usage-continues-to-climb-in-emerging-economies/) (Accessed May 2020)

440 Regueira RFS, Bernard E. 2012. Wildlife sinks: Quantifying the impact of illegal bird trade in street  
441 markets in Brazil. *Biological Conservation* **149**:16–22.

442 Rhyne AL, Tlusty MF, Schofield PJ, Kaufman L, Morris JA, Bruckner AW. 2012. Revealing the Appetite of  
443 the Marine Aquarium Fish Trade: The Volume and Biodiversity of Fish Imported into the United  
444 States. *PLoS ONE* **7**. Available from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3357433/>  
445 (accessed April 2, 2020).

446 Roberts DL, Hernandez-Castro J. 2017. Bycatch and illegal wildlife trade on the dark web. *Oryx* **51**:393–  
447 394. Cambridge University Press.

448 Roe D, Dickman A, Kock R, Milner-Gulland EJ, Rihoy E, 't Sas-Rolfes M. 2020. Beyond banning wildlife  
449 trade: COVID-19, conservation and development. *World Development* 136:105121.

450 Rowley JLL, Shepherd CR, Stuart BL, Nguyen TQ, Hoang HD, Cutajar TP, Wogan GOU, Phimmachak S.  
451 2016. Estimating the global trade in Southeast Asian newts. *Biological Conservation* **199**:96–100.

452 Sánchez-Mercado A, Cardozo-Urdaneta A, Moran L, Ovalle L, Arvelo MÁ, Morales-Campos J, Coyle B,  
453 Braun MJ, Rodríguez-Clark KM. 2020. Social network analysis reveals specialized trade in an  
454 Endangered songbird. *Animal Conservation* 23:132–144.

455 Setiawan A, Iqbal M, Halim A, Saputra RF, Setiawan D, Yustian I. 2019. First description of an immature  
456 Sumatran striped rabbit (*Nesolagus netscheri*), with special reference to the wildlife trade in  
457 South Sumatra. *Mammalia* **1**. De Gruyter. Available from  
458 [https://www.degruyter.com/view/journals/mamm/ahead-of-print/article-10.1515-mammalia-](https://www.degruyter.com/view/journals/mamm/ahead-of-print/article-10.1515-mammalia-2018-0217/article-10.1515-mammalia-2018-0217.xml)  
459 [2018-0217/article-10.1515-mammalia-2018-0217.xml](https://www.degruyter.com/view/journals/mamm/ahead-of-print/article-10.1515-mammalia-2018-0217/article-10.1515-mammalia-2018-0217.xml) (accessed April 3, 2020).

460 Silge J, Robinson D. (n.d.). Text Mining with R. Available from <https://www.tidytextmining.com/>  
461 (accessed May 14, 2020).

462 Singrodia V, Mitra A, Paul S. 2019. A Review on Web Scrapping and its Applications. Pages 1–6 2019  
463 International Conference on Computer Communication and Informatics (ICCCI).

464 Siritwat P, Nijman V. 2018. Illegal pet trade on social media as an emerging impediment to the  
465 conservation of Asian otters species. *Journal of Asia-Pacific Biodiversity* **11**:469–475.

466 Siritwat P, Nijman V. 2020. Wildlife trade shifts from brick-and-mortar markets to virtual marketplaces: A  
467 case study of birds of prey trade in Thailand. *Journal of Asia-Pacific Biodiversity*. Available from  
468 <http://www.sciencedirect.com/science/article/pii/S2287884X2030042X> (accessed April 24,  
469 2020).

470 Sonricker Hansen AL, Li A, Joly D, Mearu S, Brownstein JS. 2012. Digital Surveillance: A Novel Approach  
471 to Monitoring the Illegal Wildlife Trade. PLoS ONE **7**. Available from  
472 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3517447/> (accessed February 14, 2020).

473 Stringham OC, Lockwood JL. 2018. Pet problems: Biological and economic factors that influence the  
474 release of alien reptiles and amphibians by pet owners. *Journal of Applied Ecology* **55**:2632–  
475 2640.

476 Suiter K, Sferrazza S. 2007. MONITORING THE SALE AND TRAFFICKING OF INVASIVE VERTEBRATE SPECIES  
477 USING AUTOMATED INTERNET SEARCH AND SURVEILLANCE TOOLS:5.

478 Sula CA. 2016. Research Ethics in an Age of Big Data. *Bulletin of the Association for Information Science  
479 and Technology* **42**:17–21.

480 Sung Y-H, Fong JJ. 2018. Assessing consumer trends and illegal activity by monitoring the online wildlife  
481 trade. *Biological Conservation* **227**:219–225.

482 't Sas-Rolfes M, Challender DWS, Hinsley A, Veríssimo D, Milner-Gulland EJ. 2019. Illegal Wildlife Trade:  
483 Scale, Processes, and Governance. *Annual Review of Environment and Resources* **44**:201–228.

484 Tai MC-T. 2012. Deception and informed consent in social, behavioral, and educational research (SBER).  
485 *Tzu Chi Medical Journal* **24**:218–222.

486 Toivonen T, Heikinheimo V, Fink C, Hausmann A, Hiippala T, Järv O, Tenkanen H, Di Minin E. 2019. Social  
487 media data for conservation science: A methodological overview. *Biological Conservation*  
488 **233**:298–315.

489 Xu Q, Cai M, Mackey TK. 2020. The illegal wildlife digital market: an analysis of Chinese wildlife  
490 marketing and sale on Facebook. *Environmental Conservation*:1–7. Cambridge University Press.

491 Xu Q, Li J, Cai M, Mackey TK. 2019. Use of Machine Learning to Detect Wildlife Product Promotion and  
492 Sales on Twitter. *Frontiers in Big Data* **2**. Available from

493 <https://www.frontiersin.org/articles/10.3389/fdata.2019.00028/full> (accessed February 14,  
494 2020).

495 Ye Y-C, Yu W-H, Newman C, Buesching CD, Xu Y, Xiao X, Macdonald DW, Zhou Z-M. 2020. Effects of  
496 regional economics on the online sale of protected parrots and turtles in China. *Conservation*  
497 *Science and Practice* n/a:e161.

498 Zamora A. 2019. Making Room for Big Data: Web Scraping and an Affirmative Right to Access Publicly  
499 Available Information Online. *Journal of Business, Entrepreneurship and the Law* 12:203–228.

500 Zimmer M. 2010. “But the data is already public”: on the ethics of research in Facebook. *Ethics and*  
501 *Information Technology* 12:313–325.

502





504

505 *Figure 1.*

506 Where wildlife trade occurs on the Internet. Within the Internet, there are three 'layers' of where  
507 websites can exist: the surface web, the deep web, and the dark web. As wildlife trade moves to  
508 websites on the deep and dark web, it becomes increasingly obfuscated (denoted by darkening gray  
509 background), making it more difficult for researchers to detect and monitor (Appendix S7).

510

# USING THE INTERNET TO MONITOR WILDLIFE TRADE



511

512 *Figure 2.*

513 Flowchart of our guide to using the Internet to monitor and quantify the wildlife trade. Each number

514 corresponds to a step (and subheading title) described in text. Adjusting our guide to a specific

515 taxa/product, language, or location will occur in step 2 (Section 2), where search phrases will be tailored

516 to a specific context (and language) and search engines will be restricted to a particular country. Search

517 and selection of website/content will vary if exploring social media and other deep web content (Section

518 2.4). Collecting data is a similar process for all websites regardless of where on the web the website

519 exists (Section 4; Figure 1). Cleaning data involves processing the collected data so that it can be  
520 analyzed (e.g., verifying the species traded). Machine learning and natural language processing tools  
521 have the potential to help speed up the data cleaning process.

522



524

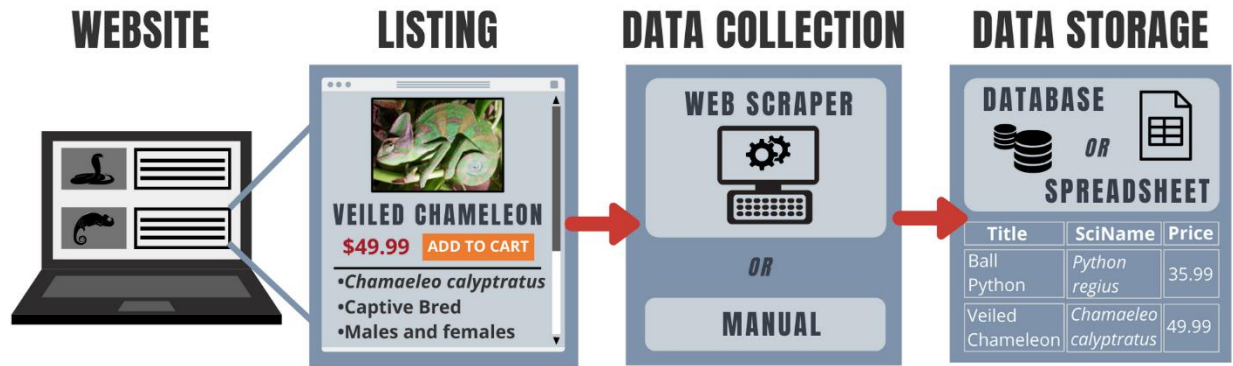
525 *Figure 3.*

526 Two hypothetical case studies using the Internet to study wildlife trade. The first (left column) explored  
 527 the trade of alien ornamental plants found in online plant shops/nurseries in Australia (i.e., open web;  
 528 sensu Lenda et al. 2014). For this study, keywords were generated based on the alien plant species of  
 529 interest (along with their scientific and trade names). In addition, qualifiers such as “for sale” or “store”  
 530 were added to the keywords to create search phrases for the search engines (See appendix S2 for  
 531 details). Next, the search engines (Google and DuckDuckGo) provided a list of candidate websites, from  
 532 which a subset were chosen based on inclusion criteria (Section 3) for data collection (designated by

533 green checkmark). For this hypothetical, the inclusion criteria were: (1) if the store sold one or more of  
534 the species of interest; and (2) if the store offered to ship plants or seeds interstate. Next, web scrapers  
535 were constructed for each website and data collected on a bi-weekly basis for one year (indicated by  
536 green check marks on calendar). Finally, data cleaning was undertaken. Since data collected from  
537 individual stores/shops tend to be more organized than other types of websites (Appendix S1) data  
538 cleaning was less intensive compared to the next case study (denoted by one hourglass icon). Further,  
539 since this study explored many species, linking each traded taxa to a taxonomic database (e.g., Global  
540 Biodiversity Information Facility: GBIF) facilitated data analysis at the species level (Section 5.1). The  
541 second hypothetical case study (right column) explored the trade of exotic leather boots made from  
542 pangolin skins occurring on social media in the United States (sensu Heinrich et al. 2019). Preliminary  
543 investigation revealed several hashtags were used in sale of pangolin-leather boots. These hashtags  
544 were supplied to the internal search engines of the social media sites (Facebook and Instagram). For this  
545 study, all posts returned from the search engine become the data itself (i.e., unconsolidated social  
546 media content; Section 2.4.2). Data collection occurred every other day since social media content tends  
547 to be updated frequently. Finally, data cleaning took longer for this study because: (i) more listings were  
548 collected compared to stores, and (ii) the listings were structured as open-text boxes, which must be  
549 read and parsed by humans to verify what is being advertised. Natural language processing and  
550 associated tools (i.e., fuzzy string matching) can be used to narrow down the number of listings needed  
551 to be cleaned (e.g., using text classification models to identify and remove irrelevant posts). GBIF logo  
552 credits: <https://www.gbif.org>.

553

554



555

556 *Figure 4.*

557 Data collection and data storage procedure for websites trading wildlife. Websites have underlying  
 558 HTML code that web scrapers can parse to extract relevant information, which can then be stored in a  
 559 database or spreadsheet. This process can be repeated for different websites using custom web scraper  
 560 code (see Appendix S5 for more information). The frequency with which to collect data will depend on  
 561 the nature of the website, including how often the website is updated. If data collection is to occur  
 562 frequently, we recommend using automated data collection because manual data collection is laborious  
 563 and time consuming. However, there is a trade-off between the resources invested in creating web  
 564 scrapers and the quantity of data that will be collected. Chameleon photo credits: Chris Kade.

565

566 **Supporting Information**

567 **Table of Contents**

568

569 Appendix S1: Types of websites relevant to wildlife trade research

570 Appendix S2: Table of search phrases from an example study

571 Appendix S3: Further information on search engines

572 Appendix S4: Further information on web traffic statistics

573 Appendix S5: Further information on web scrapers, data storage, and marking duplicates

574 Appendix S6: Recommendations for ethical practice of wildlife e-commerce surveillance

575 Appendix S7: Consequences of regulations

576 **Appendix S1: Types of websites relevant to wildlife trade research**

577 We categorize eight types of websites/platforms, which can be relevant to the wildlife trade, and we  
578 provide select examples (with references) that use each type of website as a data source.

579 1. **Shops/stores** are often physical storefronts that have websites where they list the species or  
580 wildlife products they are selling (e.g., pet stores, ornamental plant nurseries, medicinal shops).  
581 Sometimes stores will specify whether they can ship, and to where. Other times the online  
582 stores do not have a physical storefront and are exclusively an online store that ships directly to  
583 consumers. Pet stores and ornamental plant shops reliably provide either scientific names  
584 and/or common names of the species they are selling and their price. In our experience, most  
585 store websites are not updated daily and therefore collecting data weekly or fortnightly can be  
586 appropriate.

587 Select references

588 *Holmberg RJ, Tlusty MF, Futoma E, Kaufman L, Morris JA, Rhyne AL. 2015. The 800-*  
589 *Pound Grouper in the Room: Asymptotic Body Size and Invasiveness of Marine*  
590 *Aquarium Fishes. Marine Policy 53:7–12.*  
591 *Nelufule T, Robertson MP, Wilson JRU, Faulkner KT, Sole C, Kumschick S. 2020. The*  
592 *threats posed by the pet trade in alien terrestrial invertebrates in South Africa.*  
593 *Journal for Nature Conservation:125831.*

594 2. **Classifieds**, including **e-commerce** websites, are websites where individuals or companies can  
595 post their animal/wildlife/products they wish to trade. They usually appear on screen in reverse  
596 chronological order where the most recent listings appear first. Some classified websites are  
597 exclusive to particular taxa (e.g., only reptiles), while others have separate categories for  
598 multiple taxa (e.g., a bird section and a reptile section), and other sites may not have a distinct  
599 category for any taxa. Classified listings often contain some form of the taxa or product name:



600 scientific, common, trade names. However, this will vary by website, by taxa, and by individual  
601 traders. Most classified websites remove listings once they are “sold”. Price is usually provided  
602 by the user and therefore a distribution of prices for a given species or products can be derived.  
603 The location of the sale is usually given as well. For popular classifieds, data collection will likely  
604 be daily or every two to three days. Some classified websites make listings ‘inaccessible’ once  
605 the seller finds a buyer. Therefore, for these types of websites, it’s important to collect data  
606 more frequently in order to capture listings before they are removed. Further, most dark web  
607 marketplaces (i.e., crypto-markets) function nearly identically to surface web classifieds/e-  
608 commerce websites. In the US, a popular open-web classifieds website is craigslist  
609 (<https://www.craigslist.org>) and in Australia there is Gumtree (<https://www.gumtree.com.au/>).

#### 610 Select references

611 *Heinrich S, Ross JV, Cassey P. 2019. Of cowboys, fish, and pangolins: US trade in exotic*  
612 *leather. Conservation Science and Practice 1:e75.*

613 *Ye Y-C, Yu W-H, Newman C, Buesching CD, Xu Y, Xiao X, Macdonald DW, Zhou Z-M. 2020.*  
614 *Effects of regional economics on the online sale of protected parrots and turtles*  
615 *in China. Conservation Science and Practice n/a:e161.*

616 3. **Forums** are specialist websites where enthusiasts discuss various aspects of the taxa of interest.  
617 Many forums have a dedicated marketplace subforum where trading occurs. The marketplace  
618 subforums are structurally similar to classified websites. One key difference is that users can  
619 comment below the initial post asking clarifying questions. From these questions it may be  
620 possible to determine if the transaction/sale took place. Another difference is that users of  
621 forums do not typically remove “sold” listings. Either the common, scientific, or trade name is  
622 provided. The location and price are usually provided. Most forums keep an archive of all posts

623 and don't remove old posts, therefore regular data collection efforts are less essential compared  
624 to classifieds.

625 Select reference

626 *Sung Y-H, Fong JJ. 2018. Assessing consumer trends and illegal activity by monitoring the*  
627 *online wildlife trade. Biological Conservation 227:219–225.*

628 4. **Lost and Found** websites allow users to report a lost or found pet. They are structurally similar  
629 to classifieds websites. They may provide useful information if exploring invasive species risks.  
630 They are usually only available for highly visible species such as domesticated mammals (cats,  
631 dogs, rabbits), turtles, and birds. The species name (scientific, common, or trade name) as well  
632 as the location and date is usually provided.

633 Select references

634 *Kikillus KH, Hare KM, Hartley S. 2012. Online trading tools as a method of estimating*  
635 *propagule pressure via the pet-release pathway. Biological Invasions 14:2657–*  
636 *2664.*

637 *Vall-Ilosera M, Cassey P. 2017. Leaky doors: Private captivity as a prominent source of*  
638 *bird introductions in Australia. PLOS ONE 12:e0172851.*

639 5. **Adoption** websites post pet animals that are available for adoption. This is considered the  
640 secondary market for pets. They are structurally similar to classifieds websites.

641 6. **News** websites contain news from either print or electronic news companies. For the wildlife  
642 trade, many seizures of illegal wildlife are often reported in the news and may be used as a  
643 source of data.

644 Select references

645 *Indraswari K, Friedman RS, Noske R, Shepherd CR, Biggs D, Susilawati C, Wilson C. 2020.*  
646 *It's in the news: Characterising Indonesia's wild bird trade network from media-*  
647 *reported seizure incidents. Biological Conservation 243:108431.*  
648 *TRAFFIC International (2020) Wildlife Trade Portal. Available at*  
649 [www.wildlifetradeportal.org](http://www.wildlifetradeportal.org).

650 7. **Social media** websites vary drastically in their structure and content relating to the wildlife  
651 trade. Broadly, content on social media websites can be separated into: (i) 'groups' with a  
652 particular purpose where people can join and (ii) users that post to the social media platform at  
653 large, or to a group of followers (Main text, Section 2.4). Some 'groups' focus on trading  
654 particular taxa or products. The posts are similar in structure to forums, except with usually less  
655 organization. Some 'groups' are open to the public (e.g., users with a login to the social media  
656 platform can view) and others require an invitation or approval to join (e.g., 'private' groups). In  
657 addition, individual stores, breeders, or traders may maintain social media accounts where they  
658 advertise wildlife. Data collection frequency will be similar to that of classifieds. Social media  
659 websites are among the most popularly used websites. Examples of social media websites  
660 investigated for wildlife trade in the primary literature include Facebook, Twitter, Instagram,  
661 and YouTube.

662 Select references

663 *Jensen TJ, Auliya M, Burgess ND, Aust PW, Pertoldi C, Strand J. 2019. Exploring the*  
664 *international trade in African snakes not listed on CITES: highlighting the role of*  
665 *the internet and social media. Biodiversity and Conservation 28:1–19.*  
666 *Kitson, H., & Nekaris, K. A. I. (2017). Instagram-fuelled illegal slow loris trade uncovered*  
667 *in Marmaris, Turkey. Oryx, 51(3), 394.*

668 *Measey J, Basson A, Rebelo AD, Nunes AL, Vimercati G, Louw M, Mohanty NP. 2019.*  
669 *Why Have a Pet Amphibian? Insights From YouTube. Frontiers in Ecology and*  
670 *Evolution 7. Frontiers. Available from*  
671 *<https://www.frontiersin.org/articles/10.3389/fevo.2019.00052/full>*  
672 *Van TP, Luu VQ, Tien TV, Leprince B, Khanh LTT, Luiselli L. 2019. Longitudinal monitoring*  
673 *of turtle trade through Facebook in Vietnam. The Herpetological Journal 29:48–*  
674 *56.*  
675 *Xu Q, Li J, Cai M, Mackey TK. 2019. Use of Machine Learning to Detect Wildlife Product*  
676 *Promotion and Sales on Twitter. Frontiers in Big Data 2. Available from*  
677 *<https://www.frontiersin.org/articles/10.3389/fdata.2019.00028/full> (accessed*  
678 *February 14, 2020).*

679 8. **Private messaging apps** including WhatsApp and Facebook messenger (among others) are  
680 dedicated apps/platform for instant messaging between two or more people. Search engines do  
681 not index private messaging apps, so it is not possible to find individual chats using the Internet.  
682 The type of information provided about traded taxa in private messaging is likely similar to  
683 classifieds and forums. Once access is granted to a private messaging group, exporting a 'log' of  
684 the entire chat is a commonly available feature, thus negating the need for web scrapers.

685 Select references

686 *Sánchez-Mercado A, Cardozo-Urdaneta A, Moran L, Ovalle L, Arvelo MÁ, Morales-*  
687 *Campos J, Coyle B, Braun MJ, Rodríguez-Clark KM. 2020. Social network analysis*  
688 *reveals specialized trade in an Endangered songbird. Animal Conservation*  
689 *23:132–144.*

690                    *Stoner, S & Shepherd, C. 2020. Using intelligence to tackle the criminal elements of the*  
691                    *illegal trade in Indian Star Tortoises *Geochelone elegans* in Asia. *Global Ecology**  
692                    *and Conservation. 23. e01097. 10.1016/j.gecco.2020.e01097.*  
693

694 **Appendix S2: Table of search phrases from an example study**

695 We provide a table of keywords used to generate search phrases for an example study quantifying the  
696 exotic pet trade in three countries (United States, United Kingdom, and Australia). The “taxa” column  
697 refers to the taxa of interest; the “location” refers to our target locations, and “website type” refers to  
698 the website types of interest. We obtained the search phrases by performing all combinations of “taxa”,  
699 “location”, and “website type”, using the follow search phrase templates:

- 700 1. Buy {taxa} {location}  
701 2. {taxa} for sale OR purchase {location}  
702 3. {taxa} {website type} {location}

703

<b>Taxa</b>	<b>Location</b>	<b>Website Type</b>
freshwater aquarium fish	United States	Forum
marine aquarium fish	United Kingdom	Store
pet birds	Australia	Breeder
exotic pet reptiles		Adoption
exotic pet amphibians		Classifieds

704

705

706 **Appendix S3: Further information on search engines**

707 *Application Programming Interfaces (APIs)*

708 Certain search engines offer APIs, which can automate the search process by iterating over each search  
709 phrase using computer programming (e.g., Bing: Thelwall and Sud 2012). Currently, Microsoft's search  
710 engine, Bing, offers an API ([https://azure.microsoft.com/en-au/services/cognitive-services/bing-web-](https://azure.microsoft.com/en-au/services/cognitive-services/bing-web-search-api/)  
711 [search-api/](https://azure.microsoft.com/en-au/services/cognitive-services/bing-web-search-api/)) while Google does not.

712

713 *Some websites do not appear in search engines*

714 Some surface websites can opt out of appearing on search engines (Carl Drott et al. 2002), so if a  
715 website is known to be important, but does not appear in the search engine results it may still be worth  
716 considering it as a candidate website. Checking a website's "robots.txt" file will reveal if they have opted  
717 out of appearing on search engines (<https://technicalseo.com/tools/robots-txt/>).

718

719 *Choosing a cutoff point*

720 Search engines return millions of URLs per search. Thus, choosing a cutoff point to stop recording  
721 resultant URLs is important to optimize search effort. While there can be various methods of choosing a  
722 cutoff point, it is important that the chosen method is transparent and repeatable. One semi-  
723 quantitative method to decide a cutoff point can be to explore the cumulative proportion of relevant  
724 results as a function of cutoff point. The point at which the curve flattens, or begins to flatten, can be  
725 considered an optimal cutoff point.

726

727 **References**

728 Carl Drott M. 2002. Indexing aids at corporate websites: the use of robots.txt and META tags.  
729 Information Processing & Management 38:209–219.  
730  
731 Thelwall M, Sud P. 2012. Webometric research with the Bing Search API 2.0. Journal of Informetrics  
732 6:44–52.

733 **Appendix S4: Further information on web traffic statistics**

734 Many websites have web traffic statistics (i.e., metadata) that have been recorded by third party  
735 companies. For a given website, these traffic statistics can include: the number of page views per  
736 month, the rank/popularity, the country where the website is most popular, and more. One provider of  
737 website metadata is Amazon Alexa Web Information Services (<https://www.alexa.com/siteinfo>), which  
738 also has an API ([https://aws.amazon.com/marketplace/pp/B07Q71HJ3H?ref=srh\\_res\\_product\\_title](https://aws.amazon.com/marketplace/pp/B07Q71HJ3H?ref=srh_res_product_title)).  
739 There are a couple of caveats to using web traffic statistics. First is that traffic statistics are calculated for  
740 the entire website (i.e., website domain). If the website's only purpose is to trade the target taxa, then  
741 this will not be an issue (i.e., online pet store). However, for many websites, there are other reasons  
742 people visit the website than to trade the target taxa. For example, the web traffic statistics for eBay, a  
743 popular American e-commerce marketplace, would pertain to all trade on eBay and would therefore be  
744 unrepresentative of the specific trade. This makes it difficult to compare traffic statistics between  
745 websites. In addition, it's important to note that web traffic statistics are not available for all websites.  
746 Given these caveats, we recommend using web traffic data as only one line of evidence in choosing a  
747 target website.

748

749



## 750 **Appendix S5: Detailed information on web scrapers and data storage**

### 751 Background on web scrapers

752 Web scrapers are computer code that convert unstructured web data into a structured data format (i.e.,  
753 tabular data format; Singrodia et al. 2019). Coding web scrapers involves technical expertise (Mitchell  
754 2018). Outside of learning to code their own web scrapers, researchers may hire data scientists or  
755 contractors to code web scrapers. There are several open-source programming languages that can be  
756 used to code web scrapers. Some examples include the language Python with libraries bs4  
757 (<https://www.crummy.com/software/BeautifulSoup/>), requests (Chandra & Varanasi 2015), and  
758 selenium (<https://selenium-python.readthedocs.io/>). Web scraping is possible in other programming  
759 languages including R with the packages RSelenium (Harrison 2020) or rvest (Wickham 2019). In  
760 addition, there are “no code” web scrapers, which is “point and click” software that facilitates building  
761 of web scrapers without knowledge of programming (de S Sirisuriya 2015). Since web scrapers rely on  
762 the underlying HTML of a website, if a website changes its HTML structure (i.e., an update in the website  
763 layout), the web scraper may ‘break’ and will need to be updated. There must be a separate custom web  
764 scraper coded for each target website (Mitchell 2018; Holmberg et al 2015). In addition to tabular text  
765 data, web scrapers can also be programmed to download images.

766

767 Web scrapers can cause ‘harm’ to the targeted website because they take up bandwidth on the  
768 website’s server (Zamora 2019). Care should be taken not to overwhelm the targeted website with the  
769 web scraper by spacing out visits to the website (i.e., a few seconds between navigating pages). Some  
770 websites specify the amount of time to ‘wait’ in between visits in their “robots.txt” file (called crawl-  
771 delay). Spacing out visits is especially important in web scraper development. Some websites may have  
772 an auto block feature, where they will block an IP address if too many visits occur in a short amount of  
773 time.

774

775 Running web scrapers takes computing resources, however, most modern computers can handle  
776 running several web scrapers simultaneously without issues. Alternatively, setting up web scrapers to  
777 run on a cloud server or a separate dedicated computer may be desirable. If the data collection is  
778 recurrent, then establishing a system to schedule web scrapers to run at regular intervals is possible  
779 through built-in software available on all popular computer operating systems (Windows: Task  
780 Scheduler, Mac/Linux: cron).

781

782 Data storage

783 Data collected by web scrapers must be stored in a way that is retrievable for cleaning and subsequent  
784 analysis. Data storage can be achieved by using spreadsheets or databases (i.e., Database Management  
785 Systems such as MySQL). The choice is dependent on the researcher's familiarity with either, and the  
786 frequency or total number of data collection events to be stored. Regardless of the data storage  
787 technique, since the fields or columns will likely differ between websites, the researcher will need to  
788 organize and collate data for each website separately.

789

790 Duplicated listings

791 Determining and marking duplicated listings is an important post data-collection step. Detecting  
792 duplicates can be achieved by selecting a column(s) to search for duplicates. If more than one row  
793 contains the exact value for the selected column(s) then it can be labelled as a duplicate. For instance,  
794 for a pet store, a researcher may decide that if two or more listings share the exact title and exact text  
795 description, they are duplicates. Other rules/assumptions can be made depending on the specific  
796 website. Labelling unique listings with a unique identifier can help to integrate the raw data with the  
797 data cleaning.

798

799 References

800 Chandra, RV, Varanasi, BS. 2015. Python requests essentials. Packt Publishing Ltd.

801

802 Harrison J. 2020. RSelenium: R Bindings for 'Selenium WebDriver'. R package version 1.7.7.

803 <https://CRAN.R-project.org/package=RSelenium>

804

805 Holmberg RJ, Tlusty MF, Futoma E, Kaufman L, Morris JA, Rhyne AL. 2015. The 800-Pound Grouper in the  
806 Room: Asymptotic Body Size and Invasiveness of Marine Aquarium Fishes. Marine Policy 53:7–  
807 12.

808

809 Mitchell R. 2018. Web Scraping with Python: Collecting More Data from the Modern Web. O'Reilly  
810 Media, Inc.

811

812 De S Sirisuriya SCM. 2015. A Comparative Study on Web Scraping. Available from  
813 <http://ir.kdu.ac.lk/handle/345/1051> (accessed May 13, 2020).

814

815 Wickham H. 2019. rvest: Easily Harvest (Scrape) Web Pages. R package version

816 0.3.5. <https://CRAN.R-project.org/package=rvest>

817

818 Zamora A. 2019. Making Room for Big Data: Web Scraping and an Affirmative Right to Access Publicly  
819 Available Information Online. Journal of Business, Entrepreneurship and the Law 12:203–228.

820

## 821 **Appendix S6: Recommendations for ethical practice of wildlife e-commerce surveillance**

822 Minimizing harm to research participants (i.e., Internet users) is a key ethical consideration when  
823 conducting online surveillance (Buchanan 2010). Additionally, it may also be important to protect the  
824 identity of researchers, as wildlife trade activity may at times involve criminal behavior (e.g., Décarry-  
825 Héту & Aldridge 2015). Thus, we recommend researchers follow the following ethical practices:

- 826 • **Compliance with legislation and acquiring ethics approval.** Approval from relevant ethics  
827 committees should always be obtained prior to conducting research and projects should be  
828 planned with consideration of relevant legislation (e.g., Australian Government 2020).
- 829 • **De-identification of identifiable/re-identifiable data.** Some data collected by researchers can  
830 be used to identify an individual person (i.e., identifiable data). This data will require de-  
831 identification in order to maintain anonymity of the participants and thus minimize any  
832 potential harm to them. Types of data that require de-identification include (but are not limited  
833 to): names of participants, user names of participants, age of participants, and locations of  
834 online posts (e.g., street address or even postcode). The ease with which data can be de-  
835 identified will be largely dependent on the structure of a given website. For example, some sites  
836 may have a specific field for usernames, which will be straightforward to de-identify after data  
837 collection. In other instances, identifiable information may be present in an unstructured  
838 manner in free-form text boxes. This presents a much greater logistical challenge, especially if  
839 researchers are collecting large quantities of data that cannot feasibly be manually processed.  
840 Finally, in some instances, it is worth considering de-identification of the name of the website(s),  
841 group(s), or platform(s) where researchers are collecting data from (e.g., Hinsley et al. 2016).  
842 Doing so will preserve researcher anonymity and may prevent future behavioral changes of  
843 participants (e.g., switching to trading on a more secure website: Appendix S7).

- 844       • **Secure storage of research data.** Data should be stored in a secure manner with transparency  
845       around which researchers, institutions or third parties will be granted access to data with  
846       different levels of de-identification (Samuel and Buchanan 2020). It is important to consider  
847       optimal ways to store ‘big data’ while maintaining security (e.g., cloud-based storage; Buchanan  
848       and Ess 2016). Consideration of data storage should also extend beyond the anticipated lifespan  
849       of any given research project.
- 850       • **Consideration of limited disclosure of research.** Whether or not researchers can ethically justify  
851       the decision to withhold or partially withhold disclosure of their identity or research aims and  
852       methodology is highly context dependent. Given the number and potential anonymity of  
853       website users whose wildlife trade activity may be monitored, it would be logistically unfeasible  
854       to contact everyone in a manner that preserves their identity. However, the same cannot  
855       necessarily be said for contacting website administrators or subsite moderators. In fact, for deep  
856       web platforms, contact may be necessary in order to gain access to a particular subsection of a  
857       website (such as a subforum). It should be noted that such contact may induce changes in  
858       wildlife trade behavior that would reduce the value of surveillance research (e.g., individuals  
859       who would otherwise participate in illegal trade may choose to do so in an undetectable manner  
860       if they are aware that research is taking place). Therefore, the importance of outcomes likely to  
861       result from research must be weighed against the importance of gaining consent from research  
862       participants.

863

864       Australian Government, 2020. National Statement On Ethical Conduct In Human Research. Available at:  
865       [http://www.nhmrc.gov.au/\\_files\\_nhmrc/publications/attachments/e72.pdf](http://www.nhmrc.gov.au/_files_nhmrc/publications/attachments/e72.pdf) [Accessed 09  
866       November 2020].

867

868       Buchanan, E., 2010, “Internet Research Ethics: Past, Present, Future,” in *The Blackwell Handbook of*  
869       *Internet Studies*, C. Ess and M. Consalvo, (eds.), Oxford: Oxford University Press.

870

871 Buchanan, E. and C. Ess. 2016. "Ethics in Digital Research," in *The Handbook of Social Practices and*  
872 *Digital Everyday Worlds*, M. Nolden., G. Rebane, M. Schreiter (eds.), Springer.  
873

874 Décary-Héту, D, Aldridge, J. 2015. Sifting through the net: Monitoring of online offenders by researchers.  
875 *European Review of Organised Crime* 2, 122-141.  
876

877 Hinsley A, Lee TE, Harrison JR, Roberts DL. 2016. Estimating the extent and structure of trade in  
878 horticultural orchids via social media. *Conservation Biology* 30:1038–1047.  
879

880 Samuel, G, Buchanan, E, 2020. Guest Editorial: Ethical Issues in Social Media Research. SAGE Publications  
881 Sage CA: Los Angeles, CA.  
882

## 883 **Appendix S7: Consequences of regulations**

884 As more concerns around the wildlife trade emerge (e.g., ethical, disease risk, etc.), governments may  
885 impose stricter laws around the online trade of wildlife and/or companies may impose stricter self-  
886 regulations of their own websites (Roe et al. 2020). Previous research on trade bans suggest that stricter  
887 regulations can have unintended consequences, such as increasing or redirecting trade instead of  
888 eliminating it (Challender et al. 2015). Thus, we hypothesize that as regulations around wildlife trade  
889 become stricter for the open and indexed deep web, traders will likely move to the either the unindexed  
890 deep web (e.g., private messaging apps) or potentially the dark web to avoid detection. One recent  
891 incident highlights this possibility. Facebook recently implemented a ban on the sale of animals (live and  
892 derived parts) on its website  
893 ([https://www.facebook.com/policies/commerce/prohibited\\_content/animals](https://www.facebook.com/policies/commerce/prohibited_content/animals)). The efficacy of the ban  
894 in reducing trade on Facebook has not been evaluated; however, trade did not stop on Facebook  
895 because of the ban (Nijman 2020). Instead, users adjusted how they advertised wildlife, presumably to  
896 avoid detection. Users started using code names or acronyms for species and stopped including asking  
897 prices (Nijman 2020; author's personal observations). Further, users directed all questions about the  
898 advertisement to a private chat (Facebook messenger: Nijman 2020; author's personal observations). In  
899 this case, the result of stricter regulations served to decrease the amount of available information for  
900 researchers. We note that while stricter regulations may decrease information available to researchers,  
901 regulations may have the intended consequence of reducing trade overall.

902

## 903 References

904 Challender DWS, Harrop SR, MacMillan DC. 2015. Towards informed and multi-faceted wildlife trade  
905 interventions. *Global Ecology and Conservation* 3:129–148.  
906 Nijman V. 2020. Illegal trade in Indonesia's National Rare Animal has moved online. *Oryx* 54:12–13.

907 Roe D, Dickman A, Kock R, Milner-Gulland EJ, Rihoy E, 't Sas-Rolfes M. 2020. Beyond banning wildlife  
908 trade: COVID-19, conservation and development. *World Development* 136:105121.  
909  
910