

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

2

3 Oliver C. Stringham^{1,2}, Adam Toomes¹, Aurelie M. Kanishka^{1,3}, Lewis Mitchell², Sarah Heinrich¹, Joshua V.

4 Ross², Phillip Cassey¹

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

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 systematic method to automate data collection from relevant websites. This guide is adaptable to the
20 multitude of trade-based contexts including different focal taxa or derived parts, and locations of
21 interest. We provide recommendations for working with the diversity of websites that trade wildlife,
22 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, pet stores, and
31 digital platforms on the Internet (e.g., Alfino & Roberts 2018). Reliable data on the quantity and
32 composition of the wildlife trade (legal and illegal) is vital for informing decisions about conservation,
33 biosecurity and law enforcement, and developing human behavior change campaigns; yet this data is
34 not collected often, or is difficult to obtain (Regueira & Bernard 2012; Eskew et al. 2020). In recent
35 years, the Internet has played an increasingly important role in facilitating the wildlife trade (Siriwat &
36 Nijman 2020).

37

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

50

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

58

59 Here, we present an accessible guide to using the Internet to gather data on the wildlife trade. We
60 developed the methodology through our collective knowledge of working with web data and the wildlife
61 trade, combined with the methods used in prior published studies. Our goal is for this guide to be used
62 by scientists, NGOs, government agencies, and other interested parties, who wish to utilize the Internet
63 as a source of data on the wildlife trade. We do not intend this paper to be adopted as a strict protocol,
64 as the Internet is highly transient and there needs to be the flexibility to adapt to changing contexts and
65 technology.

66

67 *Structure of the Internet*

68 The Internet (i.e., the World Wide Web or simply the Web) is categorized into three distinct 'layers': the
69 surface web, the deep web, and the dark web (Figure 1; Bergman 2001). Each layer differs in one of two
70 factors: whether it's accessible without logging in or invitation (i.e., is publicly viewable) and whether it
71 is indexed by a search engine (i.e., will appear as a result in a search engine). The surface web includes
72 any website that is publicly viewable and is indexed by search engines (e.g., e-commerce websites). The
73 deep web includes websites or online content that require either logging in or an invitation to view (e.g.,

74 social media, private messaging apps). Some deep web sites may be indexed by a search engine (e.g.,
75 Facebook, Twitter), while others may not (e.g., WhatsApp). The dark web contains purposefully hidden
76 content that requires special software to access, requires either logging in, or an invitation to view, and
77 is not indexed by any search engine (Chen 2011; CRS 2017). The degree to which researchers can find
78 relevant wildlife trade content on the Internet will be influenced by how 'findable' a website/content is
79 (e.g., can a search engine find it). Further, the ethical considerations of collecting data will in part
80 depend on how 'accessible' the website is (e.g., is deceit required to gain access to the content/data;
81 Section 4.1).

82

83 *What data is available?*

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

95 **Guide to using the Internet to monitor and quantify the wildlife trade**

96 We specify our guide in six steps (Figure 2): (1) defining the scope and purpose of the project; (2) finding
97 candidate websites; (3) selecting target websites to monitor; (4) collecting and storing data from
98 websites; (5) cleaning data; and (6) analysis. Here, we detail the process of each step leading up to
99 analysis. We generalize this guide for websites found in any layer of the Internet (including social media)
100 and discuss how to adapt this guide to different languages and countries. For more generalized
101 frameworks on working with social media and online news outlets, refer to Toivonen et al. (2019) and
102 Sonricker Hansen et al. (2012), respectively.

103

104 *1. Defining the scope and purpose of the project*

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

114

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

116 Here, we detail a method to find candidate websites (e.g., e-commerce sites, forums) by using search
117 engines. Finding relevant social media content requires special considerations, which we detail in
118 Section 2.4. Outside of the search engines, other ways of finding candidate websites and choosing target

119 websites (Section 3) include: interviewing a specific community of practice (e.g., reptile keepers and
120 traders); or collaborating with other researchers actively engaged in online wildlife-trade monitoring
121 (e.g., governmental agencies or NGOs). It is important to note that the Internet is transient: traders go
122 out of business and new ones emerge. Thus, websites found at one point in time can differ in
123 composition and function if surveyed later. If the goal is long-term monitoring, we suggest revising the
124 list of current relevant websites at regularly timed intervals.

125

126 2.1 Surface Web

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

133

134 2.1.1 Defining search phrases

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

143 names, or names of breeds, morphs, and mutations (e.g., Lyons and Natusch 2013), which are not
144 captured in the initial formulation of search phrases.

145

146 2.1.2 Using search engines to perform searches

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

162

163 2.2 Deep web and dark web (non-social media)

164 Websites on the deep web indexed on search engines will be findable using the same approach outlined
165 for the surface web (e.g., private forums). Currently, there are no generalizable or automated methods
166 of finding deep web content or websites not indexed by search engines (e.g. WhatsApp, other private

167 messaging apps) nor the dark web content, outside of expert consultation or interviewing communities
168 of practice. While some algorithms exist for querying deep websites (e.g., Liakos et al. 2016), the actual
169 implementation of these algorithms as web crawlers must be tailored for each individual instance, and
170 require unique login details. This severely limits any large-scale monitoring effort such as that which we
171 are discussing here.

172

173 2.3 Classifying search engine results

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

182

183 2.4 Social media

184 2.4.1 Types of social media content

185 Social media websites vary in structure and format (Appendix S1; Toivonen et al. 2019). For our
186 purposes, we categorize content found on social media websites into two categories: consolidated and
187 unconsolidated. The differences between each category will influence how researchers find relevant
188 social media content related to wildlife trade. Consolidated social media content includes ‘groups’
189 dedicated to a particular purpose (e.g., ornamental orchid traders) where users can share content that is
190 only viewable by other group members (e.g., Facebook groups). Social media ‘groups’ function similarly

191 to forum websites. Unconsolidated social media content consists of users posting to the social media
192 platform at large or to a group of followers. Twitter, for example, is mostly public, where all ‘Tweets’
193 (i.e., posts) are viewable by all users. Some social media websites, such as Facebook, may have both
194 consolidated and unconsolidated content.

195

196 2.4.2 Finding social media content

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

209

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

215 posts. Importantly, content with privacy restrictions may be hidden from the internal search engine or
216 API results.

217

218 *3. Selecting target sites to monitor*

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

232

233 *4. Collecting data from websites*

234 Data collection can occur in one of two main ways: manual or automated. Manual data collection
235 involves visiting the website and recording what taxa/product is being traded along with desired
236 associated attributes (e.g., price, location). Automated data collection involves constructing “web
237 scrapers” to visit the website and extract desired relevant information (Figure 3; Singrodia et al. 2019).
238 Web scrapers organize the contents of a website into a structured tabular format (for more information

239 on web scrapers and data storage see Appendix S5). Since each website differs in its underlying
240 structure, custom web scrapers need to be coded for each website individually. A few highly visited
241 websites may have APIs that allow for easy collection of data; this is more likely to be the case for social
242 media websites (e.g., Twitter; Toivonen et al. 2019). Choosing manual or automatic data collection will
243 depend on how long and how often data is being collected, as it takes technical expertise and time to
244 build web scrapers, which may not be necessary if the number of target websites is small and the data
245 collection window is short (e.g., Heinrich et al. 2020). These methods of data collection apply to
246 websites and content on the deep web (including social media) and dark web, as long as researchers
247 have access to the website/content (e.g. Cunliff et al. 2019).

248

249 4.1 Ethical considerations of collecting data from the Internet

250 Ethics approval may be required to collect information from the Internet, especially when personally
251 identifiable material is collected, including, but not limited to, social media sites (Zimmer 2010). Care
252 should be taken to ensure de-identified information is used for analyses and subsequent publication
253 (Harriman & Patel 2014; Sula 2016). Furthermore, we recommend acquiring ethics approval for
254 collecting data from any deep or dark web sites since deceit may be required to obtain a login or
255 approval to join (Tai et al. 2012). Also, automated data collection processes (i.e., web scraping) currently
256 encompasses a legal gray area (Zamora 2019), and thus we encourage researchers to acquire ethics
257 approval prior to using them.

258

259 *5. Data Cleaning*

260 Data cleaning involves curating each listing (i.e., post or advertisement) for attributes that could not be
261 automatically extracted, but are required for the analysis, such as species name, quantity, price, or
262 location. Data cleaning is often a tedious and time-consuming task (Freitas & Curry 2016) and could

263 possibly be the most time-consuming part of the entire project. The amount of cleaning required will
264 depend on the structure of the website and will vary by individual website (Appendix S1). For instance, a
265 website may have a separate field for species names, while another may just have one free form open
266 text box where the user can type anything. Our experience with websites involving the wildlife trade is
267 with the latter, which takes substantially more time to clean. If collecting data manually, simultaneously
268 cleaning data during collection is possible and likely desirable.

269

270 5.1 Resolving species names

271 Resolving the species name in a listing or post is one of the most important aspects of data cleaning.
272 Some pet stores and specialist classifieds websites explicitly state the scientific name while other sites
273 may mention common names, trade names, or simply supply a photo. For all practical purposes,
274 identifications down to the rank of species are needed for effective action on conservation, biosecurity,
275 and crime (Rhyne et al. 2012). Therefore, we recommend identifying the taxa to the most specific
276 taxonomic level as possible. If pictures are provided, taxonomic experts can aid in species identification.
277 Yet, pictures may be too poor in quality to properly identify species. In some instances, online traders
278 may simply not provide enough information in the listing to identify to species level, which is an
279 unfortunate limitation of this data.

280

281 If monitoring many species, we recommend relating the species/taxa name to a taxonomic database
282 (e.g., GBIF 2020). Doing so will facilitate conformation to taxonomic names by avoiding synonyms and
283 misspellings (Gallagher et al. 2020). In addition, it will enable the researcher to easily acquire upstream
284 taxonomy (e.g., Family and Order; R package *taxize*: Chamberlain & Szocs 2013) for analysis.

285 **Advantages and caveats of web data**

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

300 **Considerations for the Deep and Dark Web**

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

308 potential for trade to be occurring elsewhere on the dark web (Roberts & Hernandez-Castro 2017).
309 Further, future policies enacted in response to conservation, and the criminological and welfare
310 concerns of wildlife trade, may shift the balance of where wildlife trade occurs on the Internet (Roe et
311 al. 2020). Specifically, new regulations or improved enforcement of illegal trade can unintentionally
312 drive trade away from the open and indexed deep web to the unindexed deep web and dark web
313 (Nijman 2020; Appendix S6), ultimately making it more difficult for researchers to locate wildlife trade
314 online.

315

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

325

326 **Automated data cleaning**

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

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

341 **Conclusions**

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

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

355

356

357 **Acknowledgements**

358 We thank Talia Wittmann for graphic design of the figures and Stephanie Moncayo for data curation in a
359 previous version of this paper. This work was supported by funding from the Centre for Invasive Species
360 Solutions (PO1-I-002: 'Understanding and intervening in illegal trade in non-native species').

361 **References**

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

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

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

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

390 Harriman S, Patel J. 2014. The ethics and editorial challenges of internet-based research. *BMC Medicine*
391 **12**:124.

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

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

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

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

400 Holmberg RJ, Tlusty MF, Futoma E, Kaufman L, Morris JA, Rhyne AL. 2015. The 800-Pound Grouper in the
401 Room: Asymptotic Body Size and Invasiveness of Marine Aquarium Fishes. *Marine Policy* **53**:7–
402 12.

403 Humair F, Humair L, Kuhn F, Kueffer C. 2015. E-commerce trade in invasive plants. *Conservation Biology*
404 **29**:1658–1665.

405 International Fund for Animal Welfare. 2018. “Disrupt: Wildlife Cybercrime”. Available on
406 [https://d1jyxxz9imt9yb.cloudfront.net/resource/218/attachment/original/IFAW_-](https://d1jyxxz9imt9yb.cloudfront.net/resource/218/attachment/original/IFAW_-_Disrupt_Wildlife_Cybercrime_-_FINAL_English_-_new_logo.pdf)
407 [_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)

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

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

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

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

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

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

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

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

425 Liakos P, Ntoulas A, Labrinidis A, Delis A. 2016. Focused crawling for the hidden web. *World Wide Web*
426 **19**:605–631.

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

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

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

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

435 Morgan J, Chng S. 2018. Rising internet-based trade in the Critically Endangered ploughshare tortoise
436 *Astrochelys yniphora* in Indonesia highlights need for improved enforcement of CITES. *Oryx*
437 **52**:744–750. Cambridge University Press.

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

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

443 Pew Research Center, February, 2016, “Smartphone Ownership and Internet Usage Continues to Climb
444 in Emerging Economies”. Available on
445 [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/)
446 [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)

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

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

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

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

457 Rowley JLL, Shepherd CR, Stuart BL, Nguyen TQ, Hoang HD, Cutajar TP, Wogan GOU, Phimmachak S.
458 2016. Estimating the global trade in Southeast Asian newts. *Biological Conservation* 199:96–100.

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

462 Setiawan A, Iqbal M, Halim A, Saputra RF, Setiawan D, Yustian I. 2019. First description of an immature
463 Sumatran striped rabbit (*Nesolagus netscheri*), with special reference to the wildlife trade in
464 South Sumatra. *Mammalia* 1. De Gruyter. Available from
465 [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)
466 [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).

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

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

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

473 Siritwat P, Nijman V. 2020. Wildlife trade shifts from brick-and-mortar markets to virtual marketplaces: A
474 case study of birds of prey trade in Thailand. *Journal of Asia-Pacific Biodiversity*. Available from
475 <http://www.sciencedirect.com/science/article/pii/S2287884X2030042X> (accessed April 24,
476 2020).

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

480 Stringham OC, Lockwood JL. 2018. Pet problems: Biological and economic factors that influence the
481 release of alien reptiles and amphibians by pet owners. *Journal of Applied Ecology* **55**:2632–
482 2640.

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

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

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

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

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

493 Toivonen T, Heikinheimo V, Fink C, Hausmann A, Hiippala T, Järvi O, Tenkanen H, Di Minin E. 2019. Social
494 media data for conservation science: A methodological overview. *Biological Conservation*
495 **233**:298–315.

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

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

500 <https://www.frontiersin.org/articles/10.3389/fdata.2019.00028/full> (accessed February 14,
501 2020).

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

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

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

509



511

512 *Figure 1.*

513 Where wildlife trade occurs on the Internet. Within the Internet, there are three ‘layers’ of where
514 websites can exist: the surface web, the deep web, and the dark web. As wildlife trade moves to
515 websites on the deep and dark web, it becomes increasingly obfuscated (denoted by darkening gray
516 background), making it more difficult for researchers to detect and monitor (Appendix S6).

517

USING THE INTERNET TO MONITOR WILDLIFE TRADE



518

519 *Figure 2.*

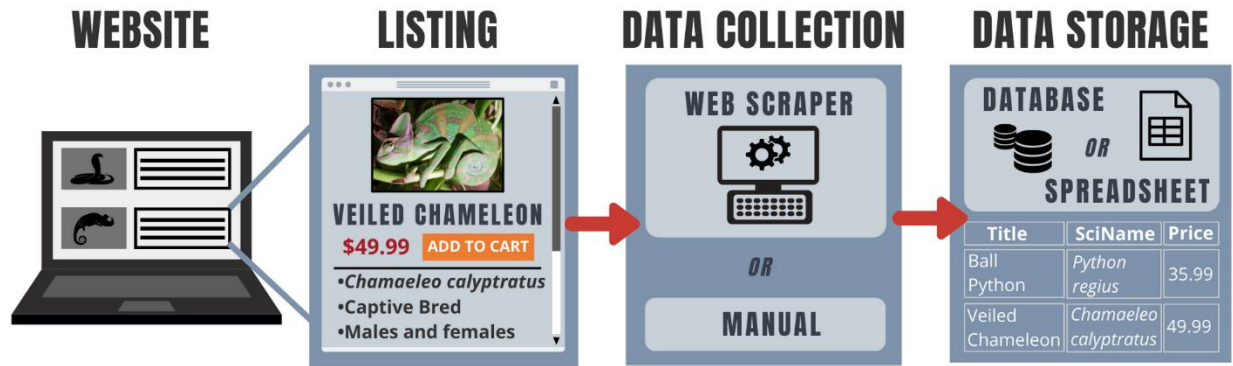
520 Flowchart of our guide to using the Internet to monitor and quantify the wildlife trade. Adjusting our

521 guide to a specific taxa/product, language, or location will occur in Step 2 (Section 2), where search

522 phrases will be tailored to a specific context (and language) and search engines will be restricted to a

523 particular country.

524



525

526 *Figure 3.*

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

535 **Supporting Information**

536 **Table of Contents**

537

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

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

540 Appendix S3: Further information on search engines

541 Appendix S4: Further information on web traffic statistics

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

543 Appendix S6: Consequences of regulations

544

545

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

547 We categorize eight types of websites/platforms that can be relevant to the wildlife trade and provide
548 select example references that use each type of website as a data source.

549 1. **Pet stores** are often physical storefronts that have websites where they list what species they
550 are selling. Sometimes pet stores will specify whether they can ship, and to where. Other times
551 the online pet stores do not have a physical storefront and is exclusively an online store that
552 ships directly to consumers. Pet stores reliably give either scientific names and/or common
553 names of the species they are selling and their price. In our experience, most pet store websites
554 are not updated daily and therefore collecting data weekly or fortnightly can be appropriate.

555 Select references

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

559 *Nelufule T, Robertson MP, Wilson JRU, Faulkner KT, Sole C, Kumschick S. 2020. The*
560 *threats posed by the pet trade in alien terrestrial invertebrates in South Africa.*
561 *Journal for Nature Conservation:125831.*

562 2. **Classifieds**, including **e-commerce** websites, are websites where individuals or companies can
563 post their animal/wildlife/products they wish to trade. They usually appear on screen in reverse
564 chronological order where the most recent listings appear first. Some classified websites are
565 exclusive to particular taxa (e.g., only reptiles), while others have separate categories for
566 multiple taxa (e.g., a bird section and a reptile section), and other sites may not have a distinct
567 category for any taxa. Classified listings often contain some form of the taxa or product name:
568 scientific, common, trade names. However, this will vary by website, by taxa, and by individual
569 traders. Most classified websites remove listings once they are “sold”. Price is usually provided

570 by the user and therefore a distribution of prices for a given species or products can be derived.
571 The location of the sale is usually given as well. For popular classifieds, data collection will likely
572 be daily or every two to three days. Some classified websites make listings ‘inaccessible’ once
573 the seller finds a buyer. Therefore, for these types of websites, it’s important to collect data
574 more frequently in order to capture listings before they are removed. Further, most dark web
575 marketplaces (i.e., crypto-markets) function nearly identically to surface web classifieds/e-
576 commerce websites. In the US, a popular classifieds website is craigslist
577 (<https://www.craigslist.org>) and in Australia is Gumtree (<https://www.gumtree.com.au/>).

578 Select references

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

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

584 3. **Forums** are specialist websites where enthusiasts discuss various aspects of the taxa of interest.
585 Many forums have a dedicated marketplace subforum where trading occurs. The marketplace
586 subforums are structurally similar to classified websites. One key difference is that users can
587 comment below the initial post asking clarifying questions. From these questions it may be
588 possible to determine if the transaction/sale took place. Another difference is that users of
589 forums do not typically remove “sold” listings. Either the common, scientific, or trade name is
590 provided. The location and price are usually provided. Most forums keep an archive of all posts
591 and don’t remove old posts, therefore regular data collection efforts are less essential compared
592 to classifieds.

593 Select reference

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

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

601 Select references

602 *Kikillus KH, Hare KM, Hartley S. 2012. Online trading tools as a method of estimating*
603 *propagule pressure via the pet-release pathway. Biological Invasions 14:2657–*
604 *2664.*

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

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

609 Select references

610 *None known to authors*

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

614 Select references

615 *Indraswari K, Friedman RS, Noske R, Shepherd CR, Biggs D, Susilawati C, Wilson C. 2020.*

616 *It's in the news: Characterising Indonesia's wild bird trade network from media-*
617 *reported seizure incidents. Biological Conservation 243:108431.*

618 *TRAFFIC International (2020) Wildlife Trade Portal. Available at*
619 www.wildlifetradeportal.org.

620 7. **Social media** websites vary drastically in their structure and content relating the wildlife trade.

621 Broadly, content on social media websites can be separated into: (i) 'groups' with a particular

622 purpose where people can join and (ii) users that post to the social media platform at large, or

623 to a group of followers (Main text, Section 2.4). Some 'groups' focus on trading particular taxa or

624 products. The posts are similar in structure to forums, except with usually less organization.

625 Some 'groups' are open to the public (e.g., users with a login to the social media platform can

626 view) and others require an invitation or approval to join (e.g., 'private' groups). In addition,

627 individual stores, breeders, or traders may maintain social media accounts where they advertise

628 wildlife. Data collection frequency will be similar to that of classifieds. Social media websites are

629 among the most popularly used websites. Examples of social media websites investigated for

630 wildlife trade in the primary literature include Facebook, Twitter, Instagram, and YouTube.

631 Select references

632 *Jensen TJ, Auliya M, Burgess ND, Aust PW, Pertoldi C, Strand J. 2019. Exploring the*

633 *international trade in African snakes not listed on CITES: highlighting the role of*

634 *the internet and social media. Biodiversity and Conservation 28:1–19.*

635 *Kitson, H., & Nekaris, K. A. I. (2017). Instagram-fuelled illegal slow loris trade uncovered*

636 *in Marmaris, Turkey. Oryx, 51(3), 394.*

637 *Measey J, Basson A, Rebelo AD, Nunes AL, Vimercati G, Louw M, Mohanty NP. 2019.*
638 *Why Have a Pet Amphibian? Insights From YouTube. Frontiers in Ecology and*
639 *Evolution 7. Frontiers. Available from*
640 *<https://www.frontiersin.org/articles/10.3389/fevo.2019.00052/full>*

641 *Van TP, Luu VQ, Tien TV, Leprince B, Khanh LTT, Luiselli L. 2019. Longitudinal monitoring*
642 *of turtle trade through Facebook in Vietnam. The Herpetological Journal 29:48–*
643 *56.*

644 *Xu Q, Li J, Cai M, Mackey TK. 2019. Use of Machine Learning to Detect Wildlife Product*
645 *Promotion and Sales on Twitter. Frontiers in Big Data 2. Available from*
646 *<https://www.frontiersin.org/articles/10.3389/fdata.2019.00028/full> (accessed*
647 *February 14, 2020).*

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

655 Select references

656 *Sánchez-Mercado A, Cardozo-Urdaneta A, Moran L, Ovalle L, Arvelo MÁ, Morales-*
657 *Campos J, Coyle B, Braun MJ, Rodríguez-Clark KM. 2020. Social network analysis*
658 *reveals specialized trade in an Endangered songbird. Animal Conservation*
659 *23:132–144.*

660

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

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

- 667 1. Buy {taxa} {location}
- 668 2. {taxa} for sale OR purchase {location}
- 669 3. {taxa} {website type} {location}

670

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

671

672

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

674 *Application Programming Interfaces (APIs)*

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

679

680 *Some websites do not appear in search engines*

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

685

686 *Choosing a cutoff point*

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

693

694 **References**

- 695 Carl Drott M. 2002. Indexing aids at corporate websites: the use of robots.txt and META tags.
696 Information Processing & Management 38:209–219.
697
698 Thelwall M, Sud P. 2012. Webometric research with the Bing Search API 2.0. Journal of Informetrics
699 6:44–52.

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

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

715

716

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

718 Background on web scrapers

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

733

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

741

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

748

749 Data storage

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

756

757 Duplicated listings

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

765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787

References

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

Harrison J. 2020. RSelenium: R Bindings for 'Selenium WebDriver'. R package version 1.7.7.
<https://CRAN.R-project.org/package=RSelenium>

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

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

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

Wickham H. 2019. rvest: Easily Harvest (Scrape) Web Pages. R package version 0.3.5. <https://CRAN.R-project.org/package=rvest>

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

788 **Appendix S6: Consequences of regulations**

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

807

808 References

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

812 Roe D, Dickman A, Kock R, Milner-Gulland EJ, Rihoy E, 't Sas-Rolfes M. 2020. Beyond banning wildlife
813 trade: COVID-19, conservation and development. *World Development* 136:105121.
814
815