

Aurorasaurus database of real-time, crowd-sourced aurora data for space weather research

B. C. Kosar^{1,2,3*}, Elizabeth A. MacDonald^{2,3}, Nathan A. Case⁴, and Matthew Heavner⁵

¹The Catholic University of America, Washington D.C., USA

²NASA Goddard Space Flight Center, Greenbelt, MD, USA

³New Mexico Consortium, Los Alamos, NM, USA

⁴Lancaster University, Lancaster, UK

⁵Los Alamos National Laboratory, Los Alamos, NM, USA

Key Points:

- Newly available Aurorasaurus database offers quality-controlled, citizen science and social media reports of aurora
- Demonstration of the breadth of scientific utility of Aurorasaurus data
- Crowd-sourced aurora data for model validation efforts of the space weather research community

*Greenbelt, MD, USA

Corresponding author: B. C. Kosar, bkosar@my.fit.edu

Abstract

This technical report documents the details of Aurorasaurus citizen science data for the period spanning 2015 and 2016 as well as its routine data filtering protocols. Aurorasaurus citizen science data is a collection of auroral sightings submitted to the project via its website or apps and mined from social media. It is a robust data set and particularly abundant during strong geomagnetic storms when auroral precipitation models have the highest uncertainty. This data is offered to the scientific community for use through an open-access database in its raw and scientific formats, each of which is described in detail in this technical report. Furthermore, by demonstrating its scientific utility, we aim to encourage its integration into auroral research.

1 Introduction

Knowing the accurate location of the auroral oval with the progression of a geomagnetic storm is important for auroral research. Auroral oval predictions are generally based on the incorporation of data collected by various space-based particle detectors or imagers into empirical models [Hardy *et al.*, 1985, 1989; Evans, 1987; Newell *et al.*, 2009, 2010a, 2014], however, the extent of their real-time prediction accuracy is unclear. Generally, they do not take into account contributions from substorms (explosive energy release within Earth's magnetic field) that can cause the auroral oval to expand and contract significantly within a few minutes. The time scale of dynamic auroral processes is faster than current operational models can predict. Auroral oval images obtained by space- and ground-based instruments provide more morphological detail in comparison to empirical model predictions. These observations are limited by coverage and typically the data are not readily available in real-time due to image processing time requirements.

Aurorasaurus [MacDonald *et al.*, 2015] is an innovative citizen science project focused on two fundamental scientific objectives: (1) collect real-time, ground-based aurora data from citizen scientists whose personal devices act as a form of soft-sensor and (2) incorporate this new type of data into scientific investigations related to aurora. Such citizen science and crowdsourcing data are becoming more common and important within space science [Cushley and Noël, 2014; Frisell *et al.*, 2014].

2 Overview of Aurorasaurus Data

Aurorasaurus data are composed of direct reports submitted to the project via its website (*aurorasaurus.org*) and iOS and Android apps and tweets that are mined from Twitter via keyword searching and geotagging [Case *et al.*, 2016b]. Direct reports can either be positive or negative, corresponding to whether or not the observer saw the aurora. The project has been live since September 2014. During the period of 2015-2016, the database compiled a total of 9,519 raw observations. The distribution of direct reports is shown in Figure 1[a]. The gray frame corresponds to the total number of direct reports collected by the project in 2015 (bar filled with diagonal lines) and 2016 (bar filled with dots). The green and the red frames show the number of positive and negative direct reports, respectively for each year. Figure 1[b] shows the distribution of tweets that are mined from the Twitter social media platform. Twitter offers public access to its Application Programming Interface (API) through which interested communities can interact with their data. The pink frame corresponds to the total number of aurora-related tweets

Figure 1. A distribution of 2015-2016 raw [a] direct reports collected via the projects's website and apps and [b] data mined from the Twitter search API.

scraped from the Twitter search API. About 15% of these tweets, shown by the purple frame, contained geographical information (or location) with them. The geolocated tweets were presented to the Aurorasaurus community to vote on. The Aurorasaurus project engages its community in tweet verification efforts by asking them to up or down vote the tweets presented on the Aurorasaurus platforms (website and apps). Tweets that are up-voted to be real-time auroral sightings are classified as "positive verified" tweets highlighted by the blue frame. The orange frame shows the number of "negatively verified" tweets indicating that they were not real-time auroral sightings or not actual auroral sightings at all and therefore, down-voted by the community [Case *et al.*, 2016b]. The total number of negatively verified tweets for both years are significantly larger compared to positive verified tweets, reflecting the noise levels inherent in the Twitter data. The black frame shows that approximately 70% of the tweets were "unverified". An earlier study by Case *et al.* [2015a] showed that the number of reports submitted to Aurorasaurus scales with the strength of the geomagnetic activity. Even though 2015 was more active in terms

74 of geomagnetic storms, the total number of reports submitted to the project increased by
75 40% in 2016. This demonstrates that the number of submissions is affected by other fac-
76 tors as well such as the growth of the size of the Aurorasaurus community, which grew
77 from ~3500 in 2015 to ~5,000 in 2016. A large number of the direct reports submitted
78 during 2016 are negative which is expected and clearly emphasized by the 50% increase
79 in number compared to 2015. The data mined from Twitter is consistently smaller in num-
80 ber during 2016 compared to 2015, likely due to declining geomagnetic activity. Even
81 though the data scraped from the Twitter API are more numerous, only a small fraction of
82 it is considered to be scientifically useful. Twitter is a unique source for robustly picking
83 out relevant data during strong geomagnetic storm conditions [*Case et al.*, 2015a].

84 Aurorasaurus uses Postgres relational databases to store its data securely and orga-
85 nize it structurally (into rows and columns) for easy access via Structured Query Lan-
86 guage (SQL) query operations. Full database access is currently limited to project team
87 members as well as the admin staff responsible for managing and maintaining it. Monthly
88 data dumps from the database track data statistics and content. These files are stored at
89 the New Mexico Consortium (NMC) servers and are maintained by the technical staff
90 of the institution. Recently, the Aurorasaurus database has increased its functionality by
91 providing access to its data through an API for research and re-serving purposes. Before
92 making this dataset open access on Zenodo repository, interested research communities
93 were granted limited access to Aurorasaurus dataset upon request. Per our privacy pol-
94 icy, access to sensitive information such as the account details of the community members
95 through the API is not permitted. Protecting the privacy of our community is a high prior-
96 ity of the project.

97 **2.1 Description of the content of Aurorasaurus data files**

99 The hierarchical tree structure of the Aurorasaurus data files is shown in Figure 2.
100 This dataset is currently open access at Zenodo data repository (zenodo.org) [*Kosar et al.*,
2018a].

98 **Figure 2.** The hierarchical tree structure of the Aurorasaurus data files.

102 The two years (2015-2016) of shared data are either in their raw or scientific for-
 103 mats. Scientific data is the cleaned version of raw data by the processes described later in
 104 this section. For the raw data, three files are shared: Tweets (yyyy_tweets_mm_raw.csv or
 105 T-file), Positive Verified Tweets (yyyy_pos_verified_mm_raw.csv or PVT-file), and Web
 106 Observations (yyyy_web_observations_mm_raw.csv or WO-file). The yyyy and mm corre-
 107 spond to year and month of each year (i.e. 01 is January), respectively. WO-files contain
 108 reports submitted directly to the project via Aurorasaurus platforms. T-files contain all
 109 the aurora-related tweets that are mined from the Twitter search API via keyword search-
 110 ing such as "aurora" or "northern lights". The Aurorasaurus server primarily filters this
 111 data by removing retweets, tweets containing spam terms, and Twitter users with "aurora"
 112 in their username. The content of the raw T- and PVT-files as well as cleaned PVT-files
 113 (yyyy_pos_verified_tweets_cleaned.csv) are described in Table 1.

Table 1: Description of data attributes found in raw T- and PVT-files.
 Raw and cleaned version of PVT-file headers are identical to each other and they are a subset of column headers found in T-file with four additional fields. The distinction between T- and PVT-files is demonstrated in the last column.

Column Header	Description	Tweets (T) or Positive Verified Tweets (PVT)?
id	Unique for each tweet	T, PVT
user_screen_name	Screen name of the community member who posted the tweet on Twitter	T, PVT
created_at	Posting time of the tweet	T, PVT
text	140 character text (frequently includes a link to the tweet window)	T, PVT
location	Well-known text (WKT) format describing the location of the community member	T
geotagged	Boolean (true or false) flag indicating if the tweets had a location embedded within them	T

Continued on next page

Table 1 – continued from previous page

Column Header	Description	Tweets (T) or Positive Verified Tweets (PVT)?
location_full_name	Full location where the tweet was originated from	T, PVT
location_country	Country where the tweet was originated from	T, PVT
clavin_enriched	Boolean (true [t] or false [f]) flag indicating if CLAVIN software was used to extract the location information of the community member through the text of the tweet.	T
verified	Time when the tweet was verified	T
verified_type	If the tweet was verified, this field indicated the verification type (positive or negative)	T, PVT
st_y ($\pm 0-90^\circ$)	Latitude of the observation location	PVT
st_x ($\pm 0-180^\circ$)	Longitude of the observation location	PVT
total_votes	Number of votes cast on the tweet	PVT
score	Final score of the tweet (positive vote = +1 and negative vote = -1)	PVT

114

115 Most of the data attributes found in T- and PVT-files are self-explanatory, however,
116 it's worth giving a more detailed explanation of a few of them than what is given in Table
117 1. The allowed number of characters per tweet has traditionally been 140, as noted under
118 "text" column, however, this has been updated to 280 characters per tweet starting late
119 2017. Therefore, Aurorasaurus data collected after 2017 will contain longer tweet texts.
120 The location information (under "location" column) of the community member is saved as
121 Well-Known Text (WKT) format that is an alphanumeric representation of geometry on a
122 map. This alphanumeric string can be converted to more readable geographic coordinates
123 (latitude, st_y, and longitude, st_x) via query operations. If the location information is
124 available, this means that the tweet has an embedded native geotag therefore the geotagged

125 column will be true ("t"). The geotagged tweet may also include location information in
126 the textual format (e.g., Quincy, MA - United States) which is consecutively saved under
127 location_full_name and location_country columns. In this scenario, the clavin_enriched
128 column will show false ("f"). However, for tweets that do not come with a native geotag
129 or a place name, we utilize an open source geoparsing software CLAVIN (Cartographic
130 Location And Vicinity INdexer) [Greenbacker and Pinney, 2012-2014] to extract location
131 information from the tweet text. In this scenario, the clavin_enriched column will be true
132 ("t").

133 PVT-files are subsets of T-files containing only the tweets that are positively verified
134 as real-time aurora sightings by the members of the Aurorasaurus community. There are a
135 total of 10 header fields in PVT-files and seven of them overlap with the content of T-files
136 already described in Table 1. The four additional fields are: st_y, st_x, total_votes, and
137 score, two of which (st_y and st_x) are described earlier. Total_votes and score represent
138 the number of votes cast on the tweet and the final score of the tweet (positive vote = +1
139 and negative vote = -1), respectively. The final score of a tweet must be greater than or
140 equal to the threshold value set by the Aurorasaurus team to be classified as a positively
141 verified tweet. Currently, this value is set to 2.

142 The Aurorasaurus project presents the citizen science community with a simple form
143 to fill out for reporting their auroral sightings. The observer is asked to fill out the infor-
144 mation on the location where the aurora was seen, and the observation period (start and
145 end time of the observation). These geolocated and time-stamped records of auroral vis-
146 ibility are frequently accompanied by optional, additional data describing the observed
147 aurora and local environmental conditions (such as color, strength of the activity, location
148 of the aurora in the night sky, and auroral type). Raw WO-file have 24 data attributes that
149 are identical to headers found in the cleaned version of this file (yyyy_web_observations_cleaned.csv)
150 and they are described in Table 2. Web observations have the latitude and longitude infor-
151 mation systematically obscured by a random amount of a kilometer or less, introducing an
152 error of ± 1 km, for privacy reasons.

Table 2: Description of data attributes in raw and cleaned WO-files.

Column Header	Description
id	Unique for each observation
activities_id	Option for choosing the level of auroral activity (Quiet, Active, or Very Active)
height_id	Option for choosing the auroral height in the sky (Overhead, Northern Horizon, 45°N, 45°S, or Whole Sky)
sky_id (N/A for positive reports)	Option for choosing the sky condition during the observation (Cloudy, Clear, or Bright)
observer_id	Unique for each community member (blank for anonymous submissions)
timestamp (yyyy-mm-dd hh:mm:ss UT)	Observation submission time into Aurorasaurus platforms
address_country	Country of the observation
address_state	State of the observation (Effective for U.S. and Canada)
location	Well-known text (WKT) format describing the location of the community member
see_aurora	Boolean (true [t] or false [f]) flag indicating if the observer saw the aurora or not
sky_other	"Other" field allows observers to manually input description of the sky condition other than the options provided (see sky_id)
time_start (yyyy-mm-dd hh:mm UT)	Beginning time of the observation (15-min resolution)
time_end (yyyy-mm-dd hh:mm UT)	Ending time of the observation (15-min resolution)
on_going	Boolean (true [t] or false [f]) flag indicating if the auroral activity is continuing at the time of the report submission
height_other	"Other" field allows observers to manually input description of the auroral height in the sky other than the options provided (see height_id)
activities_other	"Other" field allows observers to manually input description of the level of auroral activity other than the options provided (see activities_id)

Continued on next page

Table 2 – continued from previous page

Column Header	Description
colors_other	"Other" field allows observers to manually input auroral colors observed other than the options provided (see colors_id)
types_other	"Other" field allows observers to manually input auroral types observed other than the options provided (see types)
comment	Allows observers to provide additional comments
image	If an auroral image captured by the observer was submitted to the server - yes [y] otherwise no [n]
st_y ($\pm 0-90^\circ$)	Latitude of the observation location ($\sim 1\text{km}$ accuracy)
st_x ($\pm 0-180^\circ$)	Longitude of the observation location ($\sim 1\text{km}$ accuracy)
colors	Option for choosing auroral colors (Red, Green, White, or Pink - community members can pick multiple colors)
types	Option for choosing auroral types (Discrete Arcs, Diffuse Glows, or Patches - community members can pick multiple types)

153

154 The scientific data is the processed version of the raw data and, maintains the same
 155 header fields. For ease of use, scientific data for all months for each year are combined
 156 into one file for positive verified tweets (yyyy_pos_verified_tweets_cleaned.csv) and one
 157 file for web observations (yyyy_web_observations_cleaned.csv).

158 Aurorasaurus, like any other citizen science project, exercises high data quality stan-
 159 dards essential to the success of the project. Data are subject to thorough inspection for
 160 quality and integrity. Duplicate reports that are posted due to technical issues encountered
 161 during submission are filtered. Of interest to our primary scientific investigations are the
 162 negative reports with an indication of clear, unobscured view of the night sky. Therefore,
 163 negative reports that specify the sky condition to be "cloudy" or "bright" are removed
 164 from the dataset. However, negative reports that come with no indication of the sky con-
 165 dition (i.e., community member skips sky_id field), are counted as scientifically valuable
 166 data because the sky condition being clear is equally likely as being bright or cloudy.

167 Twitter data is also subject to rigorous processing for data quality by means of a
 168 three-step system: filtering, verification, and validation. As mentioned earlier, aurora-

169 related tweets mined from Twitter are subject to filtering before being presented to the
170 community on the Aurorasaurus platforms. Besides filtering, extracting meaningful sig-
171 nals from Twitter data requires verification and manual validation. Filtered tweets with
172 location information are initially presented to the community members on Aurorasaurus
173 platforms to verify if they are real-time aurora sightings. After exceeding a certain thresh-
174 old (the final vote score should be greater than or equal to 2) a tweet is classified as a
175 "positive verified tweet". Verified tweets are checked annually following a pre-determined
176 set of rules to ensure their validity for detailed scientific analysis. The verification is a
177 time consuming and labor intensive task that is primarily done by the Aurorasaurus team
178 members and/or volunteers recruited under a standard protocol. Team members are the
179 core group of scientists that are/were affiliated with the project. Volunteers are usually
180 recruited from high school/undergraduate students through education and outreach ac-
181 tivities of the project by the team members. Team members or volunteers involved in
182 manual validation are required to read and understand the privacy policy of the project
183 (<http://aurorasaurus.org/privacy>) prior to any sort of data handling or database access. Au-
184 rorasaurus community members are protected by our privacy policy. Personally identifi-
185 able information and data that requires proper crediting to their owner (such as images)
186 are excluded from the public access.

187 The details of manual tweet verification are discussed in an earlier study [*Case et al.*,
188 2016b] based on the analysis of tweets collected during March and April 2015 that in-
189 cludes the period of St. Patrick's Day storm [*Case et al.*, 2015b]. The raw positively veri-
190 fied tweets are sifted through one at a time and they are divided into two major categories,
191 valid or invalid. The valid category represents tweets that were identified correctly as
192 real-time auroral sightings while the invalid category is a collection of tweets that were
193 misidentified as real-time auroral sightings by the Aurorasaurus community. The invalid
194 category is further broken down into subcategories i.e., not real-time (red), not original
195 (yellow), overlap (orange), wrong location (blue), not a positive sighting (gray), and junk
196 (purple). The distribution of these categories for 2015 and 2016 data is shown in Figure
197 3. The description of each category can be found in the work of *Case et al.* [2016b]. True
198 and false positives (TP and FP) refer to positively verified tweets that are valid and in-
199 valid, respectively. By utilizing the number of TP and FP, the positive predictive value
200 (PPV) for the tweet verification system was found to be 20% and 31% for 2015 and 2016,
201 respectively. In other words, 20% and 31% of the tweets identified as positively verified

202 in 2015 and 2016 were actually valid. There is an increase in this value for 2016, how-
 203 ever, the source of this variance is not well understood. The increase is not attributable to
 204 sample size because although 2015 was more active (hence higher number of positively
 205 verified tweets) in comparison to 2016, the number of valid tweets is fewer.

207 Figure 3 also shows that the percentage of the "not real-time" subcategory of invalid
 tweets is reduced in 2016. Identifying a tweet as real-time or not requires detailed investi-

206 **Figure 3.** The distribution of positively verified tweets collected during 2015 and 2016.

208
 209 gation of many aspects of that particular tweet. The procedure is a set of rules developed
 210 by the Aurorasaurus team members. For data quality assurance, team members and volun-
 211 teers are trained on the same set of hundred tweets that were used during the project's first
 212 validation efforts. Because validating a large dataset tends to be a time-consuming pro-
 213 cess, alternative techniques (such as machine learning algorithms) to speed up or eliminate
 214 manual validation efforts are being explored. The project currently has two years of data
 215 (2015-2016) validated for quality and readily available for scientific use. This data can be
 216 utilized for evaluation of existing models [Newell *et al.*, 2009, 2014; Zhang and Paxton,
 217 2008] and used as a new data source complementing the data-sparse field of Heliophysics.

218 **2.2 Citizen scientist descriptions of auroral observations in 2015-2016**

220 Of the 1,740 and 2,435 raw reports submitted in 2015 and 2016, 19.8% and 19.7%
 221 of them included an image of the observed aurora. Submitted auroral images are com-
 222 posed of smartphone photos of the back screen of a Digital Single-Lens Reflex (DSLR)
 223 camera, lower-quality smartphone images taken of the aurora directly, and high-quality
 224 post-processed images. On average 52% of the reports also contain descriptive informa-
 225 tion about the observed aurora. If a community member skips one question on the form
 226 (e.g. color), they often skip the rest (i.e. type, sky location, activity). This is apparent in
 227 the percentages of each data attribute skipped being very similar. Figure 4 shows how citi-
 228 zen scientists described their observations during 2015-2016. Most of the observed aurora
 229 were either typical green auroral emission or multicolor (combination of green with other
 230 colors). The observed types are dominated by discrete arcs and diffuse glows or multi-
 231 ple types (combination of arcs, glows, and pulsating patches). Most observers described

219 **Figure 4.** Description of the observed aurora by the citizen scientists during 2015-2016.

232 aurora being on the northern horizon or 45° above the horizon. The whole sky obser-
 233 vations are sparse, which is likely due to the limited number of inhabitants at latitudes
 234 likely to see overhead aurora. Aurora was reported to be more active in 2015 (please see
 235 <http://blog.aurorasaurus.org/?p=356>) in comparison to 2016.

236 **3 Scientific utility of Aurorasaurus database**

242 The cleaned positive verified tweets and direct reports are subject to two more fil-
 243 ters that are implemented in IDL codes. The plots shown in Figure 5 are produced for
 244 the time period of 2015-01-01 00:00:00 UT to 2016-12-31 23:59:59 UT. The first filter
 245 applied to the cleaned data files further checks to make sure the report times fall within
 246 this range. This filter removes only a few reports from the total (2 positive verified tweets
 247 and 12 positive reports). During submission, community members occasionally pick an
 248 incorrect time period (the difference between the end_time and the start_time) for their
 249 observations. The second filter removes positive/negative reports with an observation time
 250 period exceeding 3-hrs, as they may contain an error or not be specific enough for anal-
 ysis. In total, 214 positive and 18 negative reports are removed by filter two. Figure 5[a]

237 **Figure 5.** Distribution of validated [a] positive verified tweets and [b] web reports over the globe and the
 238 distribution of validated and filtered data as a function of [c] absolute magnetic latitude and [d] magnetic
 239 local time. Green and red filled circles correspond to positive and negative web reports, and blue filled circles
 240 correspond to positively verified tweets. The color code used for making the stacked bars refer to the same
 241 data types.

251 and 5[b] are distributions of positive verified tweets and direct reports on a world map.
 252 This data is a collection of geolocated and timestamped signals of auroral visibility ob-
 253 tained from soft-sensors. These signals exhibit a sparse spatial organization with isolated
 254 regions of high signal density nested within low signal density distribution over the globe.
 255 Data coverage over land is reasonable, particularly around populated sectors of the high
 256 latitude regions of the northern hemisphere where aurora is typically visible. This scenario
 257

258 reverses to no data over the ocean and only a few points on the southern hemisphere due
259 to the limited land area from which an aurora might be visible. With our systematic out-
260 reach efforts, particularly during strong geomagnetic activity, the Aurorasaurus community
261 and contributed observations will continue to grow in the near future. In the world map
262 shown in Figure 5[b], there are a few data points (positive and negative reports) coming
263 from very low latitude regions. While positive sightings at very low latitudes are highly
264 unlikely, negative reports are still reasonable. Positive reports are most likely submitted
265 by mistake or could be spam members submitting anonymously since there was no geo-
266 magnetic storm large enough to cause the auroral oval to expand that far south. This rep-
267 represents a minor caveat in positive reports.

268 Figure 5[c] and [d] show the distribution of Aurorasaurus reports submitted during
269 2015-2016, grouped by absolute magnetic latitude in 0.5° bins and magnetic local time in
270 30 min bins, respectively. The stacked green, red, and blue bars indicate the number of
271 positive reports, negative reports, and verified tweets that fall into each bin. The distribu-
272 tion of this data as a function of absolute magnetic latitude indicates that the number of
273 reports peak around $\sim 58^\circ$ latitude and span a wide range between 40 to 75° latitude. Au-
274 rorasaurus report submission hours span a range between 18:00 to 06:00 MLT with a peak
275 around midnight. Most auroral models typically have the highest uncertainty during large
276 geomagnetic storms when Aurorasaurus data is the most abundant. This unique data set
277 can potentially help reduce this uncertainty.

278 **3.1 Example scientific application**

279 The scientific utility of this innovative and robust citizen science data collected by
280 the Aurorasaurus project has been demonstrated in numerous publications across multiple
281 disciplines. *Case et al.* [2015a] is the first study showing the effectiveness of social me-
282 dia (Twitter) in detecting real-time auroral activity, specifically during strong geomagnetic
283 disturbances. The large number of initial reports collected during the St. Patrick's Day
284 storm of 2015 [*Case et al.*, 2015b] by the Aurorasaurus platform, were evaluated against
285 the "view-line" - an aurora forecast product of NOAA's Space Weather Prediction Center
286 (SWPC) that is obtained using the predictions of Oval Variation, Assessment, Tracking,
287 Intensity, and Online Nowcasting (OVATION) Prime 2010 auroral precipitation model and
288 demonstrates the most southern latitude of the visible aurora. The results indicated that
289 the latitudes of the majority of the citizen science reports were significantly equatorward

290 of the view-line latitudes predicted by the SWPC [*Case et al.*, 2016a]. We note that the
 291 latitude of the citizen science reports solely represent the location of the observer submit-
 292 ting the report. The latitude is not derived using the location of the aurora in the sky. A
 293 recent case study [*Kosar et al.*, 2018b] compared a subset of this data with the equatorial
 294 boundaries of the auroral oval at a fixed flux level obtained from the solar wind driven
 295 OVATION Prime 2013 (OP-13) model [*Newell et al.*, 2014] and the Kp-dependent Zhang-
 296 Paxton model [*Zhang and Paxton*, 2008]. It was found that the OP-13 boundary is slightly
 297 more consistent with the citizen science data.

305 Global auroral particle precipitation is a result of coupling between the magnetosphere-
 306 ionosphere system that is driven by the external solar wind plasma input. The OVATION
 307 Prime 2013 (OP-13) auroral precipitation model uses a solar wind-magnetosphere cou-
 308 pling function to produce its high-resolution electron energy flux maps for the aurora. As
 309 described in *Case et al.* [2016a], this electron energy flux can be converted to a probability
 310 of visible aurora by scaling the summed precipitation energy flux (j) and adding an offset
 311 to it (i.e. $P(A) = 10 + 8\sum j$). In addition to this empirical conversion, NOAA's SWPC has
 312 a coarse estimate of a view line to account for the auroral height in the sky. The SWPC
 313 view line represents the lowest latitude where aurora should be visible. Aurorasaurus data
 314 is mostly clustered around the equatorial edge of the auroral oval hence offering useful
 data for assessing the accuracy of the view line. Following the earlier work [*Case et al.*,

298 **Figure 6.** The differences in latitude between Aurorasaurus reports collected in 2015 and the SWPC view
 299 line at the same longitude are grouped into 0.5° bins. Stacked bars indicate number of each type of report
 300 falling into each interval. The color code used for the data types is the same as earlier. Approximately $\sim 50\%$
 301 of the observations are reported from latitudes that are further equatorward of the view line estimated by
 302 the NOAA SWPC. The accuracy is calculated using true positive (TP) reports that include positively veri-
 303 fied tweets (blue) and positive web reports (green) and true negative (TN) reports that include negative web
 304 reports (red).

315
 316 2016a], outputs of the OVATION Prime 2013 model with a 15-minute cadence were pro-
 317 duced and the energy flux outputs were converted to percent probability of visible aurora.
 318 Figure 6 shows the distribution of Aurorasaurus data collected in 2015, grouped by lati-
 319 tude differences between Aurorasaurus data ($|\phi_{obs}|$) and SWPC view lines ($|\phi_{VL}^{SWPC}|$) into

320 0.5° bins. The accuracy is calculated using a statistical technique suggested by *Machol*
 321 *et al.* [2012], $ACC = (\sum TP + \sum TN) / \sum R$ where $\sum TP$ is the total number of true positive
 322 reports that fall within, $\sum TN$ is the total number of true negative reports that fall out-
 323 side of the view-line, $\sum R$ is the total number of reports. This equation yields an accuracy
 324 (ACC) of approximately 50.3% for the SWPC view line.

325 **3.2 Aurorasaurus database of optical, geotagged auroral imagery**

330 Recent technological advancements have equipped citizen scientists with devices
 331 (smartphones, DSLR cameras) that are capable of capturing high-quality image data. In
 332 the two year period of 2015-2016, a total of 823 auroral images have been submitted to
 333 the Aurorasaurus project accompanying the auroral sighting reports. We note that the im-
 334 age data are not shared on Zenodo due to the terms and conditions of the Aurorasaurus
 335 privacy policy. This database has permission for research use offering a unique collec-
 336 tion of geotagged and optical auroral imagery as well as time lapse. Even though image
 337 sequences captured by the citizen scientists are rare, they are particularly useful in visu-
 338 alizing temporal and spatial dynamics of auroral arcs during geomagnetic storms. One
 339 example are auroral beads that are repeating patterns or structures within the auroral arcs.
 Typically, scientific instruments such as imagers on-board satellites or all-sky cameras cap-



326 **Figure 7.** [a] Side view image of auroral beads observed during a geomagnetic storm from Saskatoon,
 327 Canada using a DSLR camera. The beads have a 20 km spacing based on star-tracking and analysis. [b] Im-
 328 age of STEVE (Strong Thermal Emission Velocity Enhancement) and its accompanying green picket fence
 329 features forming south of the traditional auroral oval.

340 ture them from above or below and may not have the resolution for fine scale structures.
341 Citizen science images, such as the one shown in Figure 7[a], provide scientists with a
342 new set of data obtained from ground but from a different perspective and resolution. This
343 particular side profile image of auroral beads allowed us to determine dimensions of an
344 individual upright ray (width ~ 5 km and length ~ 15 km), the separation between two arbi-
345 trarily selected rays (~ 20 km), and the approximate total arc size within the field of view
346 (~ 500 km) using star field analysis. The image sequence of this particular event allowed
347 us to observe the direction of motion of individual rays. Citizen scientists collecting im-
348 ages of auroral arcs such as these provide new pieces of information about aurora that
349 contribute to research interests of the space weather community. The Aurorasaurus blog
350 has posted an article (<http://blog.aurorasaurus.org/?p=398>) on auroral beads featuring this
351 particular image and discussing it relative to images of auroral beads captured by all-sky
352 imagers and instruments on-board Earth-orbiting satellites [*Henderson, 2008; Kalmoni*
353 *et al., 2015*].

354 A collaborative research opportunity between the Aurorasaurus citizen science net-
355 work and auroral researchers has recently led to the discovery of an optical signature of a
356 new sub-auroral phenomena (see Figure 7[b]) - STEVE (Strong Thermal Emission Veloc-
357 ity Enhancement) [*MacDonald et al., 2018; Gallardo-Lacourt et al., 2018*]. This transient
358 structure forms equatorward of the traditional auroral oval and displays a purplish color
359 that is not typical of an auroral emission. In the declining period of solar maximum, these
360 phenomena have been frequently caught on citizen scientists cameras and submitted to the
361 Aurorasaurus project. With an expanding Aurorasaurus community, this image database
362 will continuously grow to allow opportunities for detailed analysis of STEVE in the near
363 future.

364 **4 Conclusions**

365 The Aurorasaurus project provides curated citizen science aurora data, particularly
366 abundant during strong geomagnetic storms, as a useful resource for the space weather
367 research community. Currently, two years (2015-2016) of data are available for scientific
368 use due to data validation challenges. Alternative solutions for automating this effort is
369 a work in progress and an important future step for the Aurorasaurus project. The newly
370 emerging fields of artificial intelligence and machine learning offers algorithms (natural

371 language processing, classification, etc.) that may be well-suited for the tweet validation
372 efforts of the project.

373 To demonstrate the scientific utility of this dataset, Aurorasaurus reports are com-
374 pared with the OVATION-driven view line predictions of NOAA SWPC for 2015. Au-
375 rorasaurus reports are mostly clustered around the equatorial edge of the auroral oval
376 hence offering a useful dataset for assessing accuracy. We find that ~50% of the obser-
377 vations are reported from the latitudes that are further equatorward of the view line es-
378 timated by NOAA SWPC. This unique dataset has a great potential for validating, im-
379 proving, and complementing existing models for auroral oval predictions and specifica-
380 tions. Emerging computational methods based on data-model integration offer new in-
381 sights that could potentially improve real-time assessment and space weather prediction
382 when citizen science data are combined with traditional sources. A future study will fo-
383 cus on developing a state-of-the-art auroral assimilative model that combines observational
384 data (citizen science reports) with existing empirical models. Once developed, this as-
385 similative model will provide feedback to model validation and ionospheric conductance
386 challenges introduced by the NASA Community Coordinated Modeling Center (CCMC)
387 (<https://ccmc.gsfc.nasa.gov>).

388 The Aurorasaurus database also offers high quality images and time-lapse sequences
389 of aurora captured by the community members. This geotagged image database contains a
390 new set of data obtained from the ground but from a different perspective in comparison
391 to ground- and space-based scientific equipment. This image database is a valuable com-
392 plement to current scientific research and also provides opportunities for new discoveries
393 advancing our understanding of the night sky.

394 **Acknowledgments**

395 This material is based upon work supported, in part, by the National Science Foundation
396 (NSF) under grant 1344296 and NASA CAN award. Any opinions, findings, and conclu-
397 sions or recommendations expressed in this material are those of the author(s) and do not
398 necessarily reflect the views of NSF. NAC was supported during this study by STFC grant
399 ST/M001059/1. The OVATION Prime output and associated view line were kindly sup-
400 plied by the Space Weather Prediction Center, Boulder, CO, National Oceanic and Atmo-
401 spheric Administration (NOAA), U.S. Department of Commerce. The output can be freely
402 downloaded from the NOAA SWPC product pages (<http://www.swpc.noaa.gov/products/aurora->

403 30-minute-forecast). Reid Priedhorsky provided the initial software engineering and contri-
404 butions to the Aurorasaurus vision. Ideum and David Kingman coded the database, web-
405 site, and apps described here. We would like to thank Sean McCloat, Mike Cook, and
406 Madison Smith for their valuable help with tweet validation efforts of the Aurorasaurus
407 project.

408 **References**

- 409 Case, N., E. MacDonald, M. Heavner, A. Tapia, and N. Lalone (2015a), Mapping auroral
410 activity with twitter, *Geophysical Research Letters*, *42*(10), 3668–3676.
- 411 Case, N., E. MacDonald, and K. Patel (2015b), Aurorasaurus and the st. patrick’s day
412 storm, *Astronomy and Geophysics*, *56*(3), 3–13.
- 413 Case, N., E. A. MacDonald, and R. Viereck (2016a), Using citizen science reports to de-
414 fine the equatorial extent of auroral visibility, *Space Weather*, *14*(3), 198–209.
- 415 Case, N. A., E. A. MacDonald, S. McCloat, N. Lalone, and A. Tapia (2016b), Determin-
416 ing the accuracy of crowdsourced tweet verification for auroral research, *Citizen Sci-*
417 *ence: Theory and Practice*, 2016.
- 418 Cushley, A., and J.-M. Noël (2014), Ionospheric tomography using ads-b signals, *Radio*
419 *Science*, *49*(7), 549–563.
- 420 Evans, D. (1987), Global statistical patterns of auroral phenomena, pp. 325–330.
- 421 Frissell, N., E. Miller, S. Kaeppler, F. Ceglia, D. Pascoe, N. Sinanis, P. Smith,
422 R. Williams, and A. Shovkoplyas (2014), Ionospheric sounding using real-time amateur
423 radio reporting networks, *Space Weather*, *12*(12), 651–656.
- 424 Gallardo-Lacourt, B., J. Liang, Y. Nishimura, and E. Donovan (2018), On the origin of
425 steve: Particle precipitation or ionospheric skyglow?, *Geophysical Research Letters*,
426 *45*(16), 7968–7973.
- 427 Greenbacker, C., and T. Pinney (2012-2014), Clavin [software], *Berico Technologies*.
- 428 Hardy, D. A., M. Gussenhoven, and E. Holeman (1985), A statistical model of auroral
429 electron precipitation, *Journal of Geophysical Research: Space Physics*, *90*(A5), 4229–
430 4248.
- 431 Hardy, D. A., M. Gussenhoven, and D. Brautigam (1989), A statistical model of auroral
432 ion precipitation, *Journal of Geophysical Research: Space Physics*, *94*(A1), 370–392.
- 433 Henderson, M. G. (2008), Observational evidence for an inside-out substorm onset sce-
434 nario, in *Annales Geophysicae*, LA-UR-08-07251; LA-UR-08-7251, Los Alamos Na-

- 435 tional Lab.(LANL), Los Alamos, NM (United States).
- 436 Kalmoni, N. M., I. J. Rae, C. E. Watt, K. R. Murphy, C. Forsyth, and C. J. Owen (2015),
437 Statistical characterization of the growth and spatial scales of the substorm onset arc,
438 *Journal of Geophysical Research: Space Physics*, *120*(10), 8503–8516.
- 439 Kosar, B. C., E. A. MacDonald, N. A. Case, and M. Heavner (2018a), Aurorasaurus real-
440 time citizen science aurora data [data set], *Zenodo*, doi:10.5281/zenodo.1255196.
- 441 Kosar, B. C., E. A. MacDonald, N. A. Case, Y. Zhang, E. J. Mitchell, and R. Viereck
442 (2018b), A case study comparing citizen science aurora data with global auroral bound-
443 aries derived from satellite imagery and empirical models, *Journal of Atmospheric and*
444 *Solar-Terrestrial Physics*.
- 445 MacDonald, E. A., N. A. Case, J. H. Clayton, M. K. Hall, M. Heavner, N. Lalone, K. G.
446 Patel, and A. Tapia (2015), Aurorasaurus: A citizen science platform for viewing and
447 reporting the aurora, *Space Weather*, *13*(9), 548–559.
- 448 MacDonald, E. A., E. Donovan, Y. Nishimura, N. A. Case, D. M. Gillies, B. Gallardo-
449 Lacourt, W. E. Archer, E. L. Spanswick, N. Bourassa, M. Connors, M. Heavner,
450 B. Jackel, B. Kosar, D. J. Knudsen, C. Ratzlaff, and I. Schofield (2018), New science
451 in plain sight: Citizen scientists lead to the discovery of optical structure in the upper
452 atmosphere, *Science Advances*, *4*(3), doi:10.1126/sciadv.aag0030.
- 453 Machol, J. L., J. C. Green, R. J. Redmon, R. A. Viereck, and P. T. Newell (2012), Evalua-
454 tion of OVATION Prime as a forecast model for visible aurorae, *Space Weather*, *10*(3).
- 455 Newell, P., T. Sotirelis, and S. Wing (2009), Diffuse, monoenergetic, and broadband au-
456 rora: The global precipitation budget, *Journal of Geophysical Research: Space Physics*,
457 *114*(A9).
- 458 Newell, P., K. Liou, Y. Zhang, T. Sotirelis, L. Paxton, and E. Mitchell (2014), OVATION
459 Prime-2013: Extension of auroral precipitation model to higher disturbance levels,
460 *Space Weather*, *12*(6), 368–379.
- 461 Newell, P. T., T. Sotirelis, and S. Wing (2010a), Seasonal variations in diffuse, monoener-
462 getic, and broadband aurora, *Journal of Geophysical Research: Space Physics*, *115*(A3).
- 463 Zhang, Y., and L. Paxton (2008), An empirical Kp-dependent global auroral model based
464 on TIMED/GUVI FUV data, *Journal of Atmospheric and Solar-Terrestrial Physics*,
465 *70*(8), 1231–1242.

Figure 1.

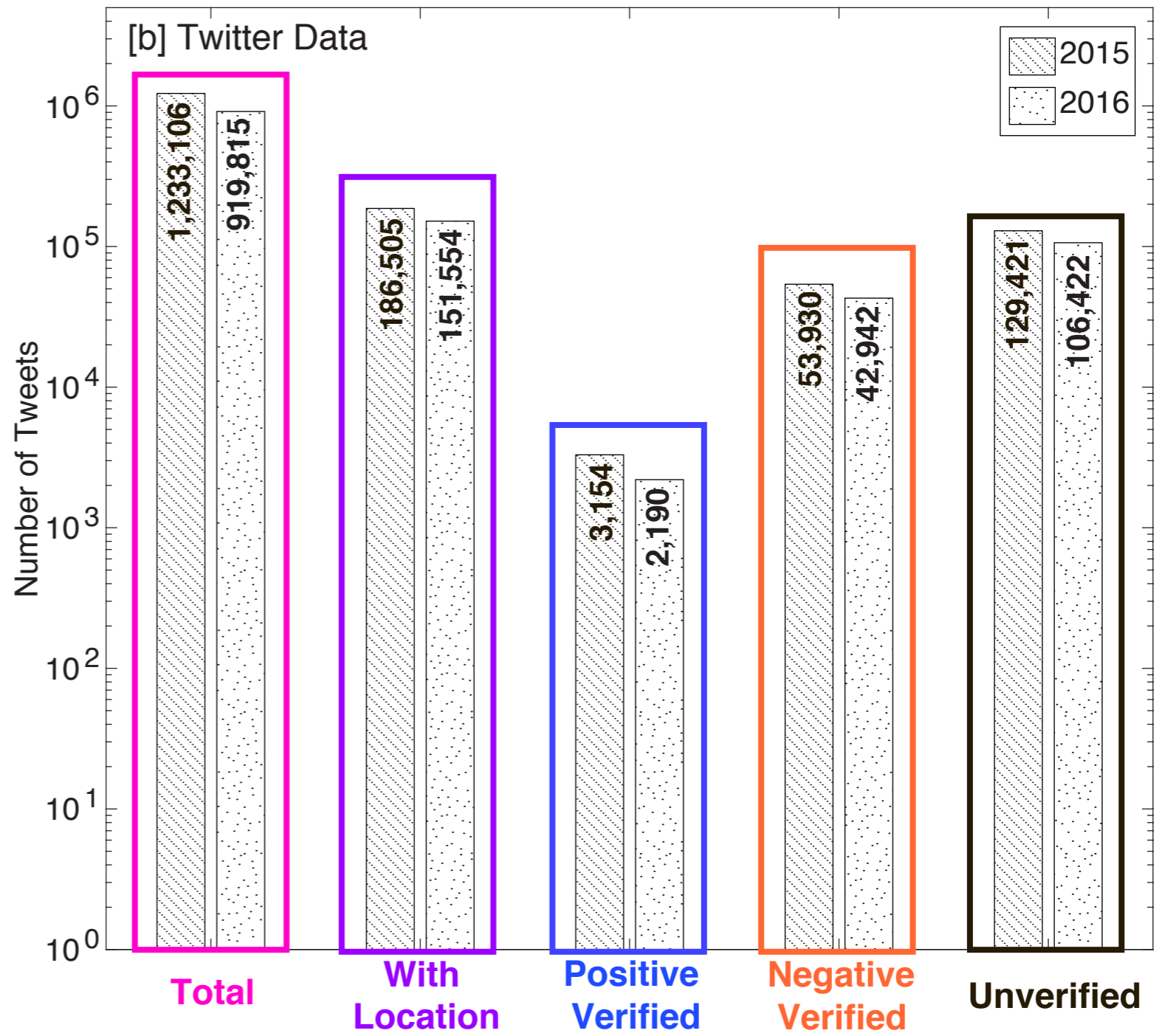
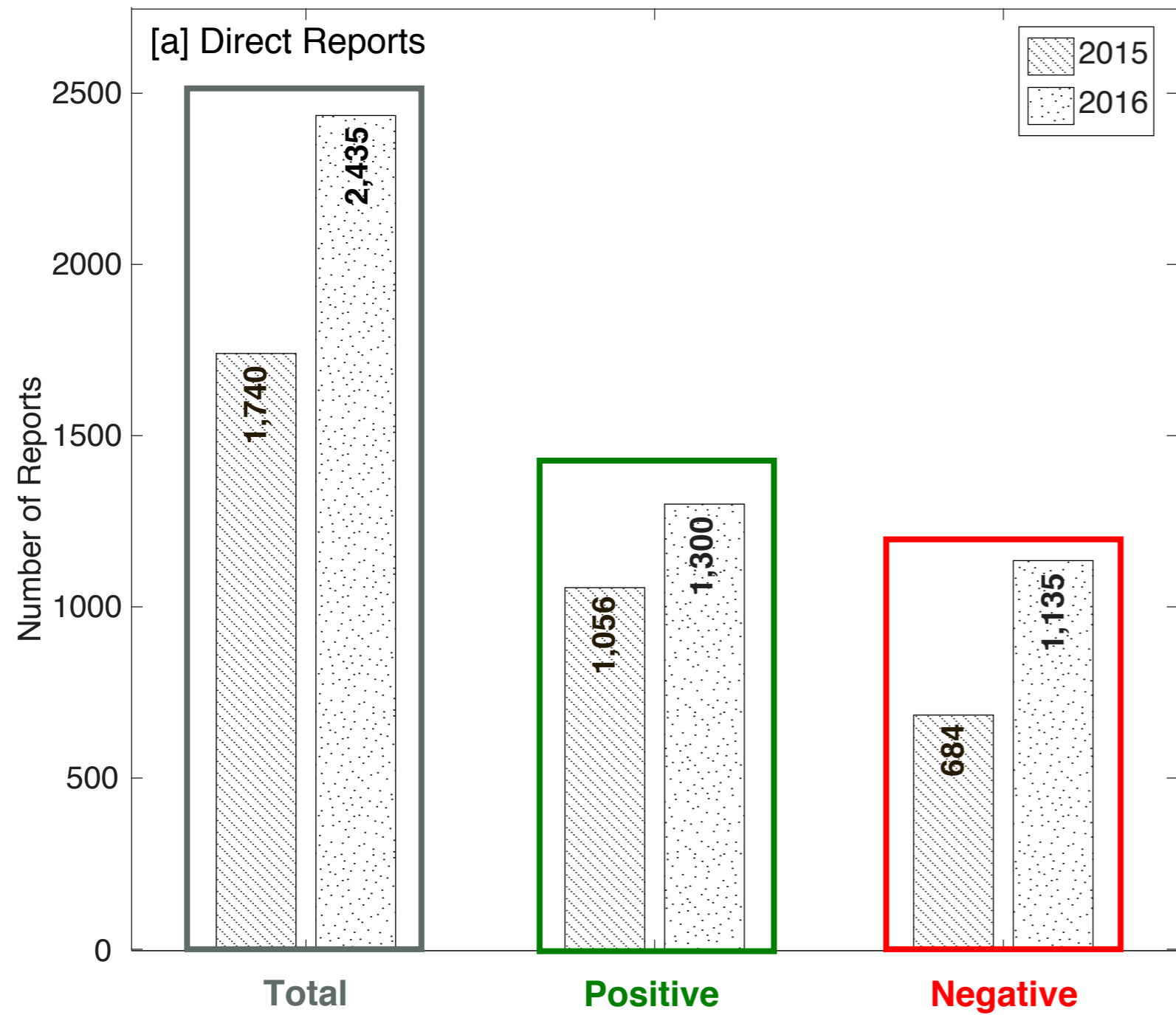


Figure 2.

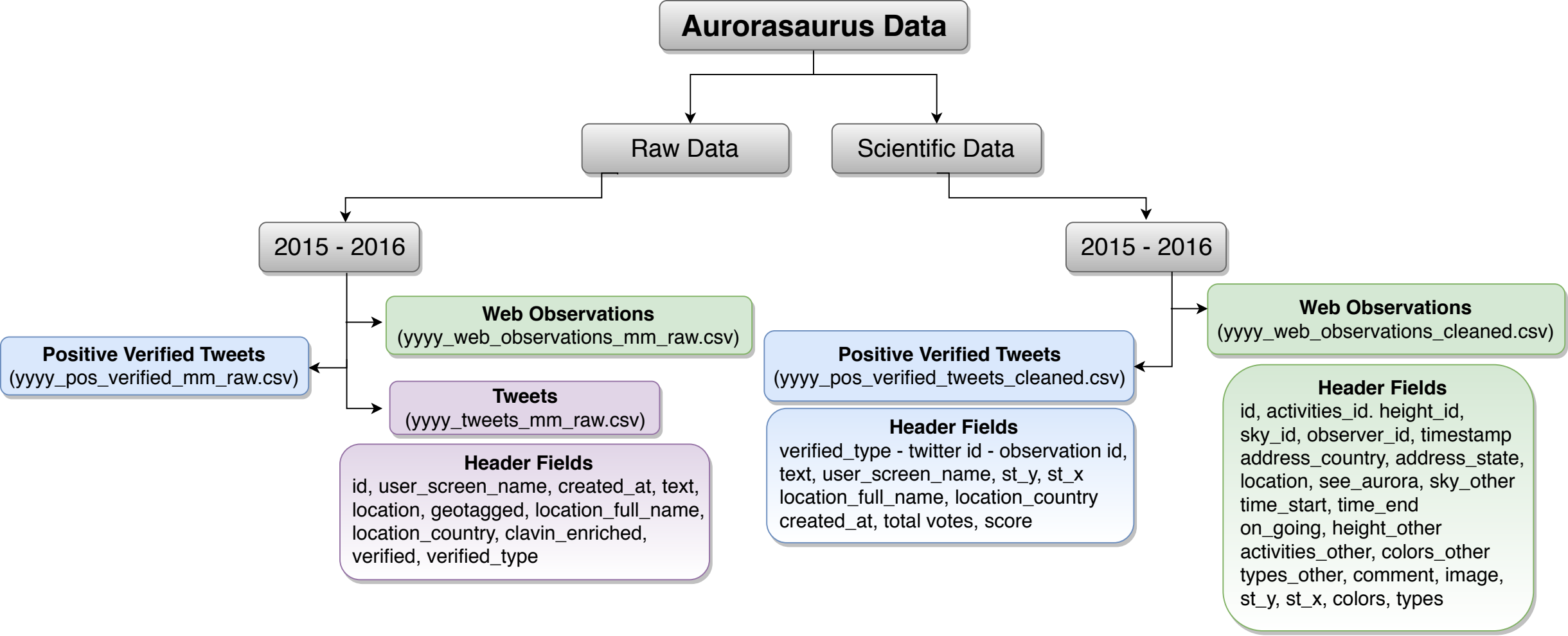


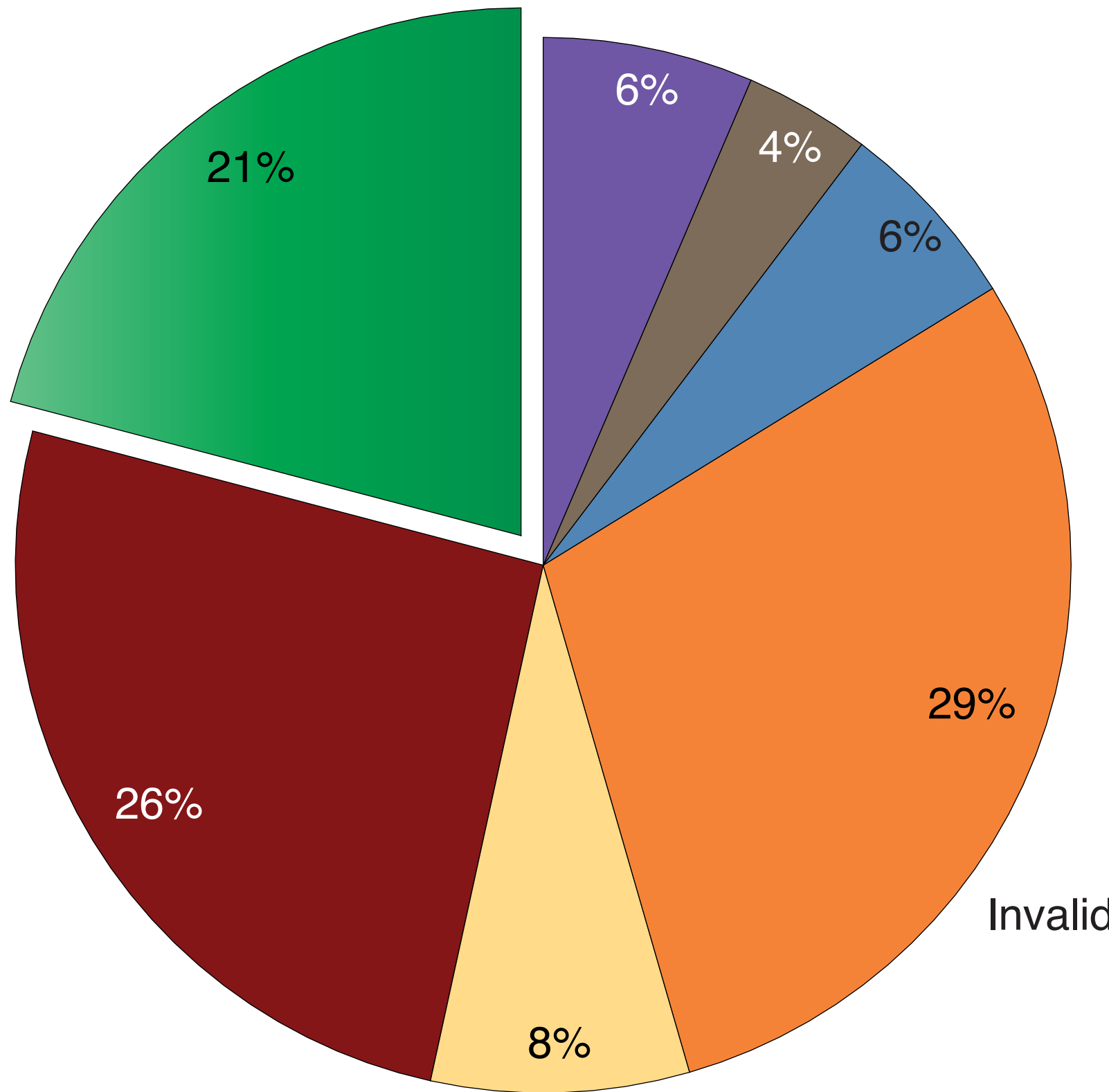
Figure 3.

Positively verified tweets (raw): **3154**
After manual validation (cleaned): **659**

2015

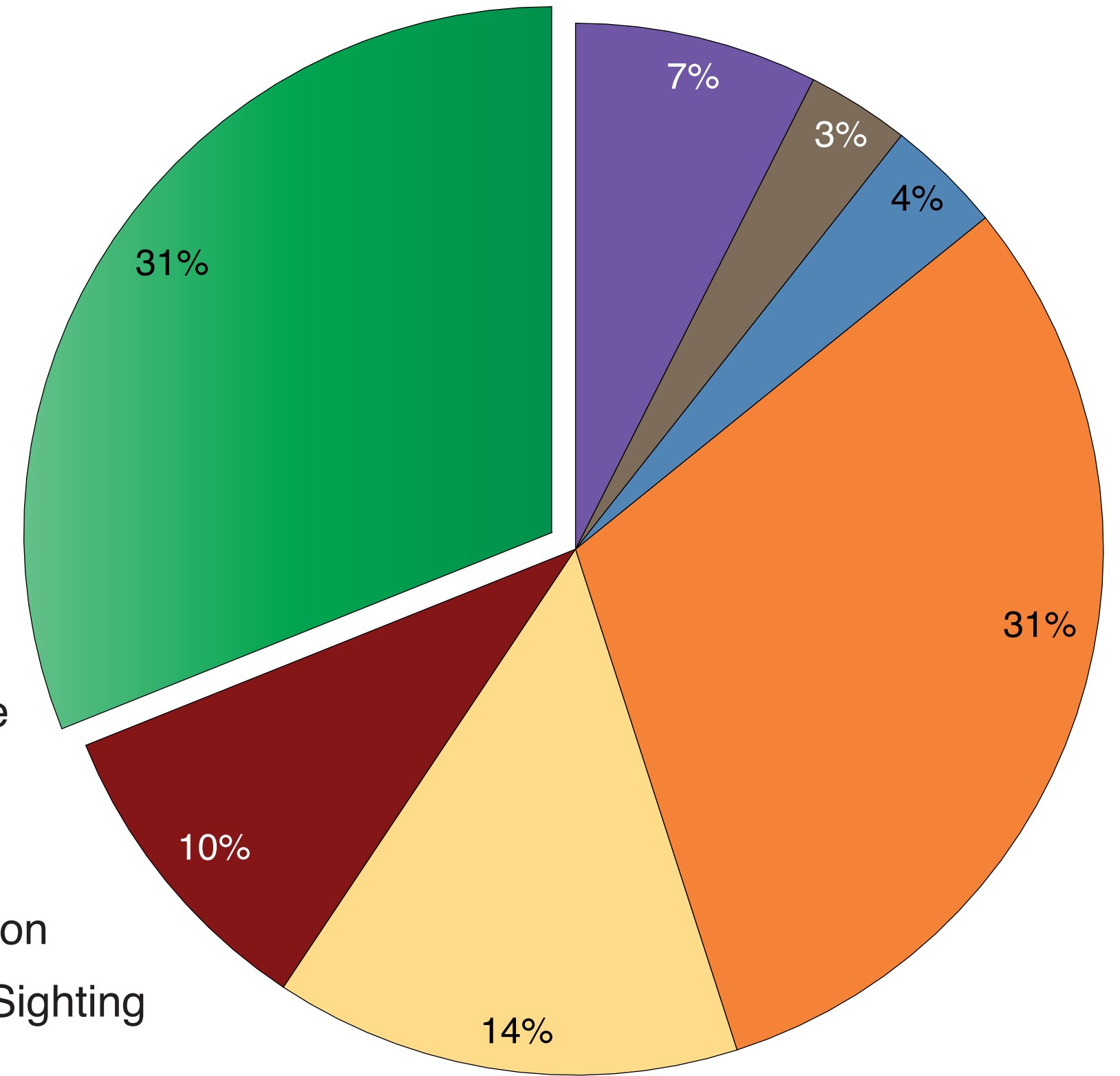
Positively verified tweets (raw): **2190**
After manual validation (cleaned): **681**

2016



- Valid
- Not Real-time
- Not Original
- Overlap
- Wrong Location
- Not Positive Sighting
- Junk

Positive Predictive Ability (PPV)
 $[\sum TP / (\sum TP + \sum FP)] = [659 / (659 + 2495)] = 20\%$



Positive Predictive Ability (PPV)
 $[\sum TP / (\sum TP + \sum FP)] = [681 / (681 + 1509)] = 31\%$

Figure 4.

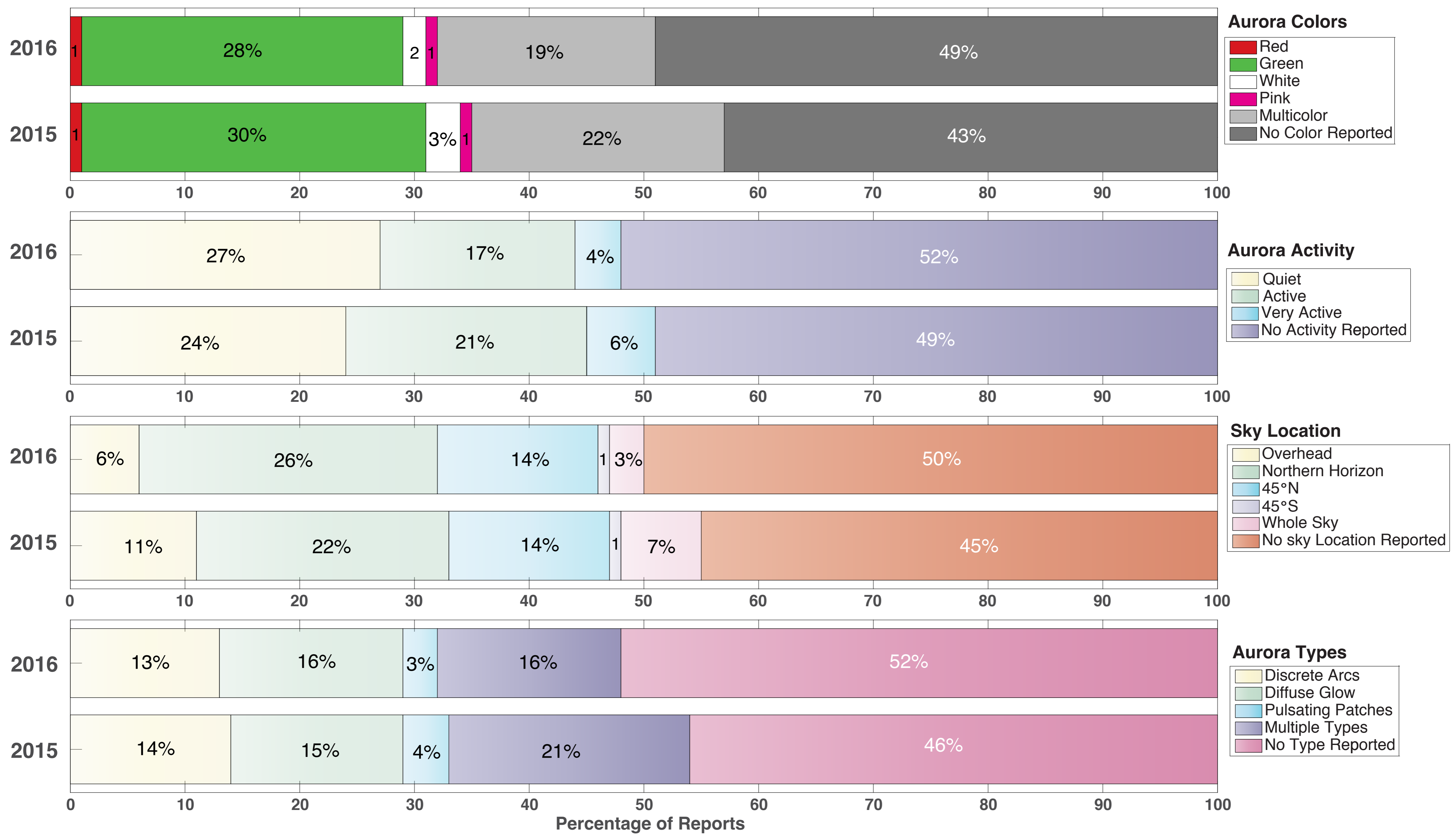


Figure 5.

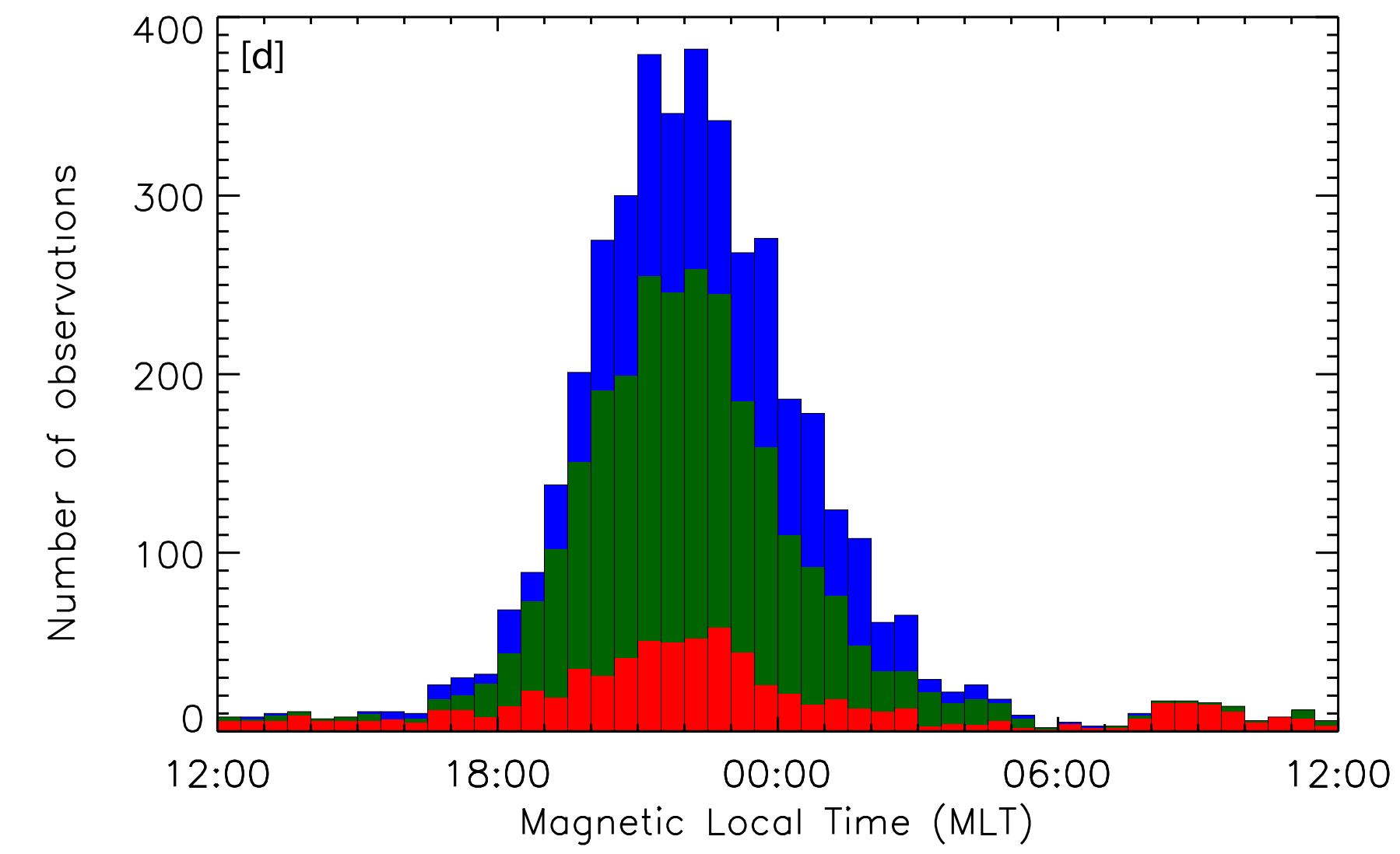
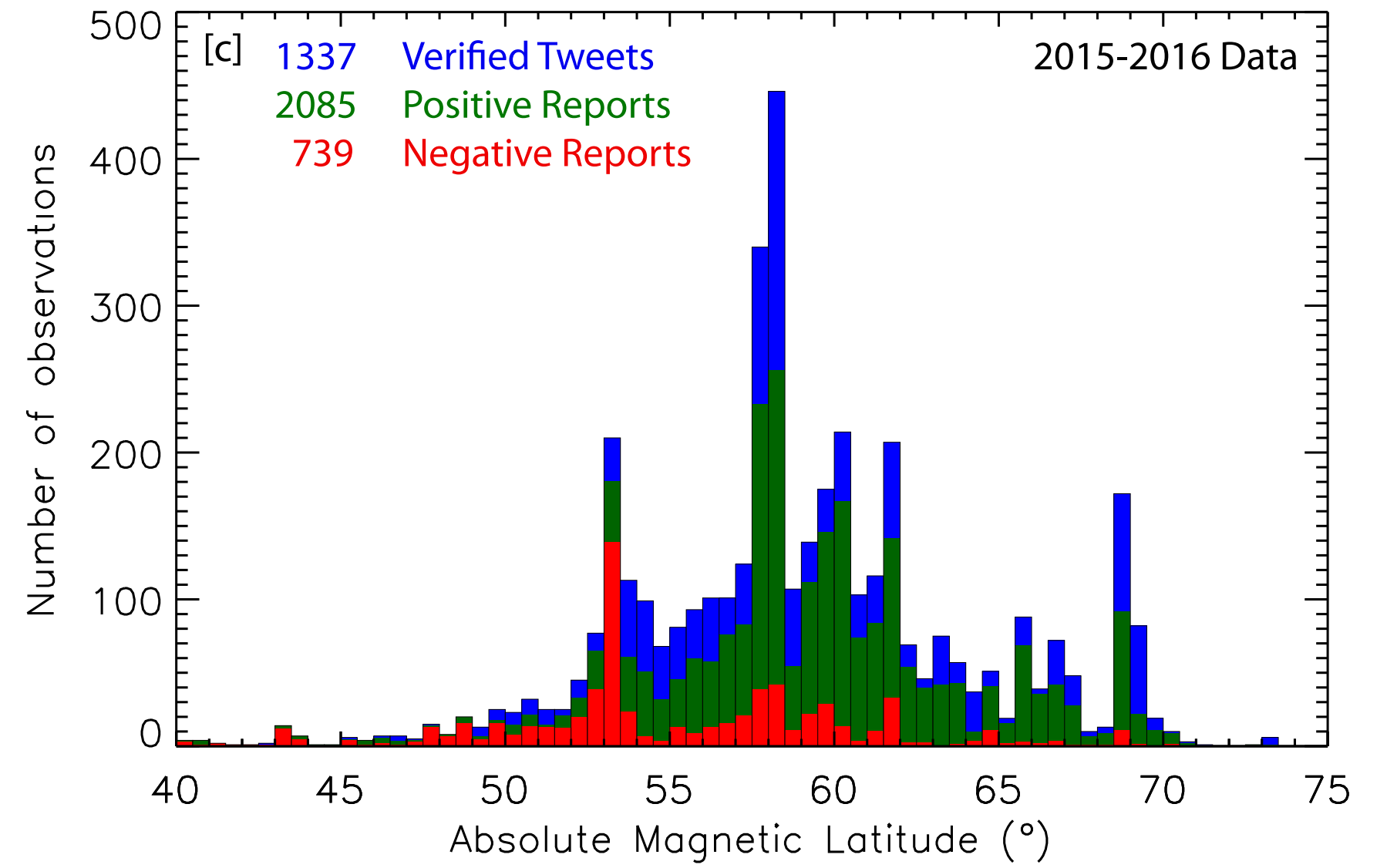
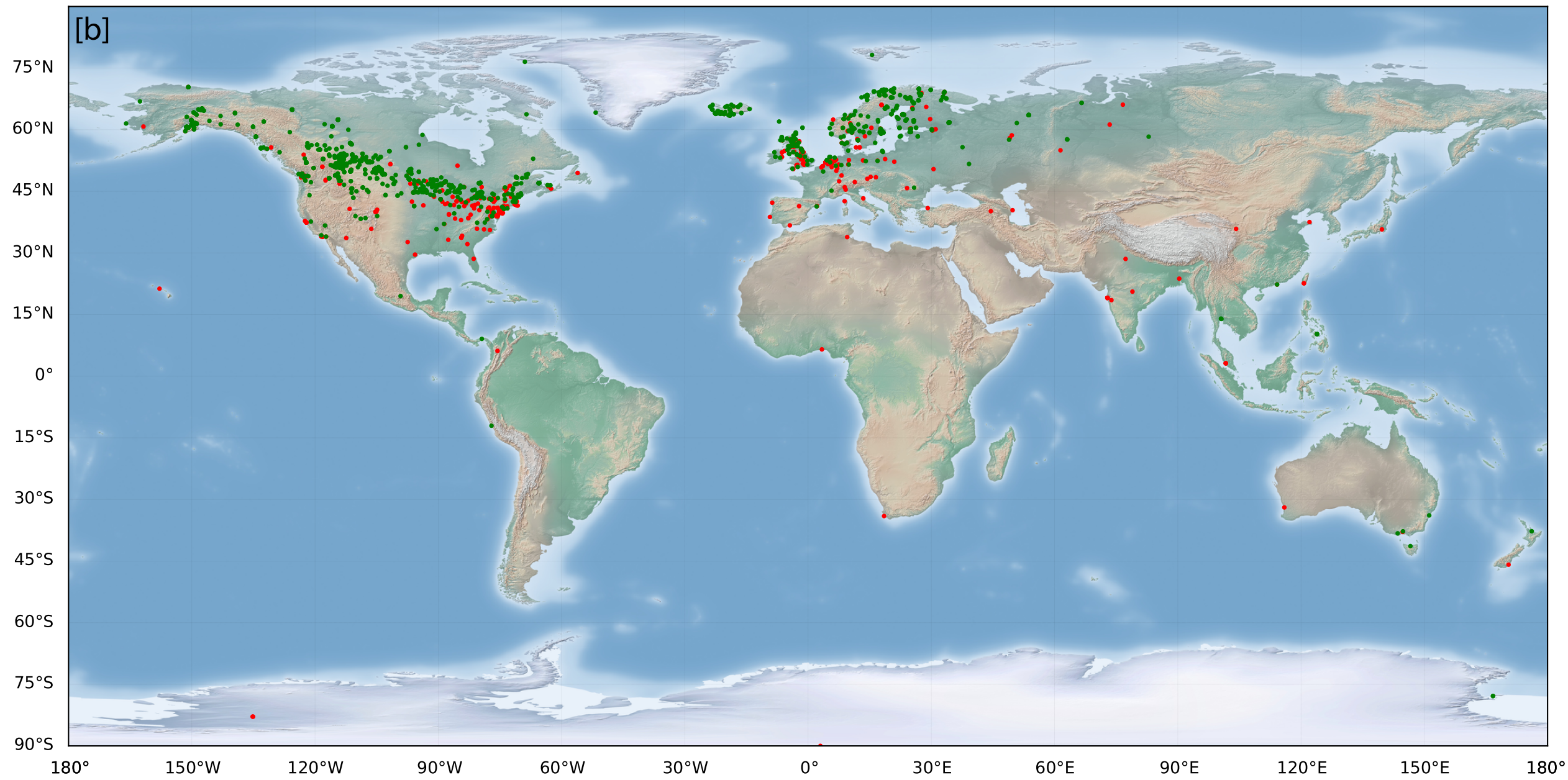
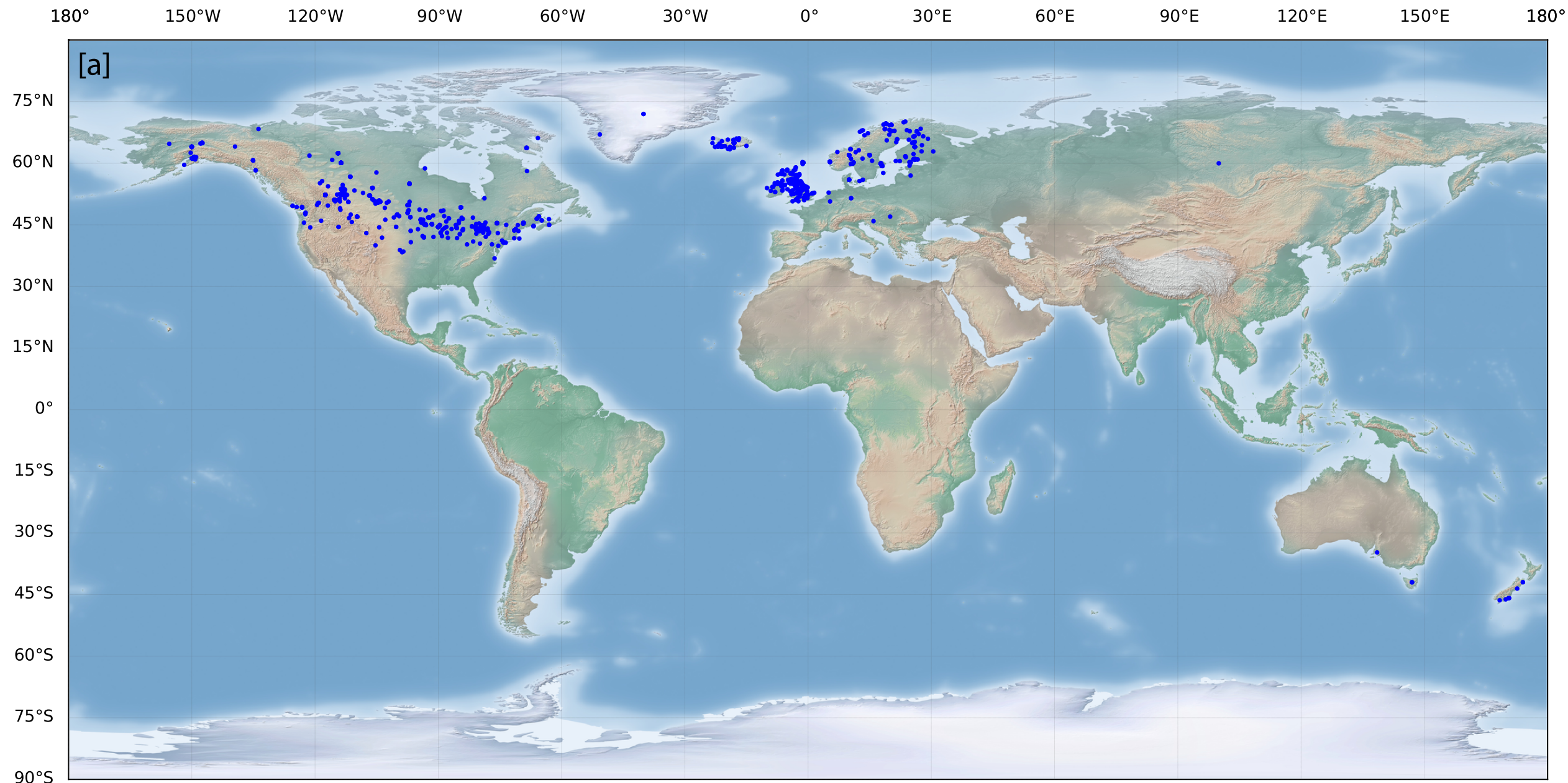


Figure 6.

SWPC View Line

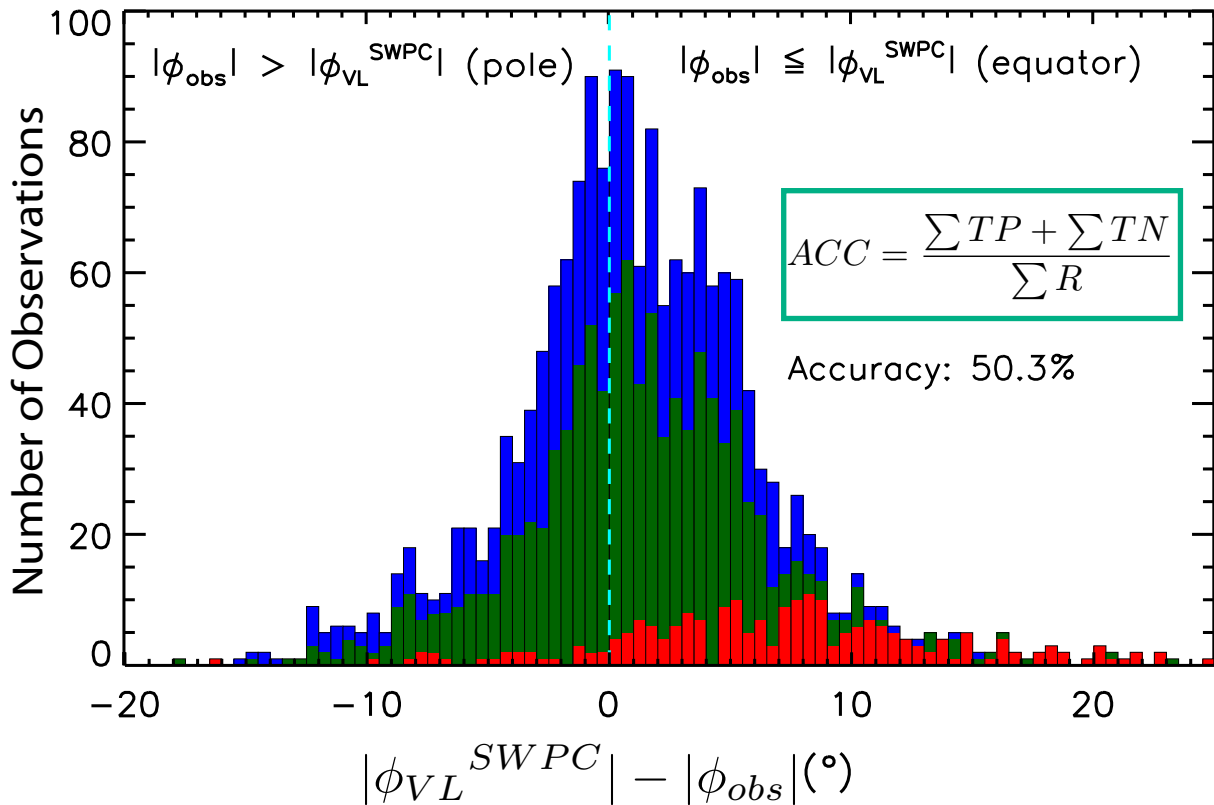
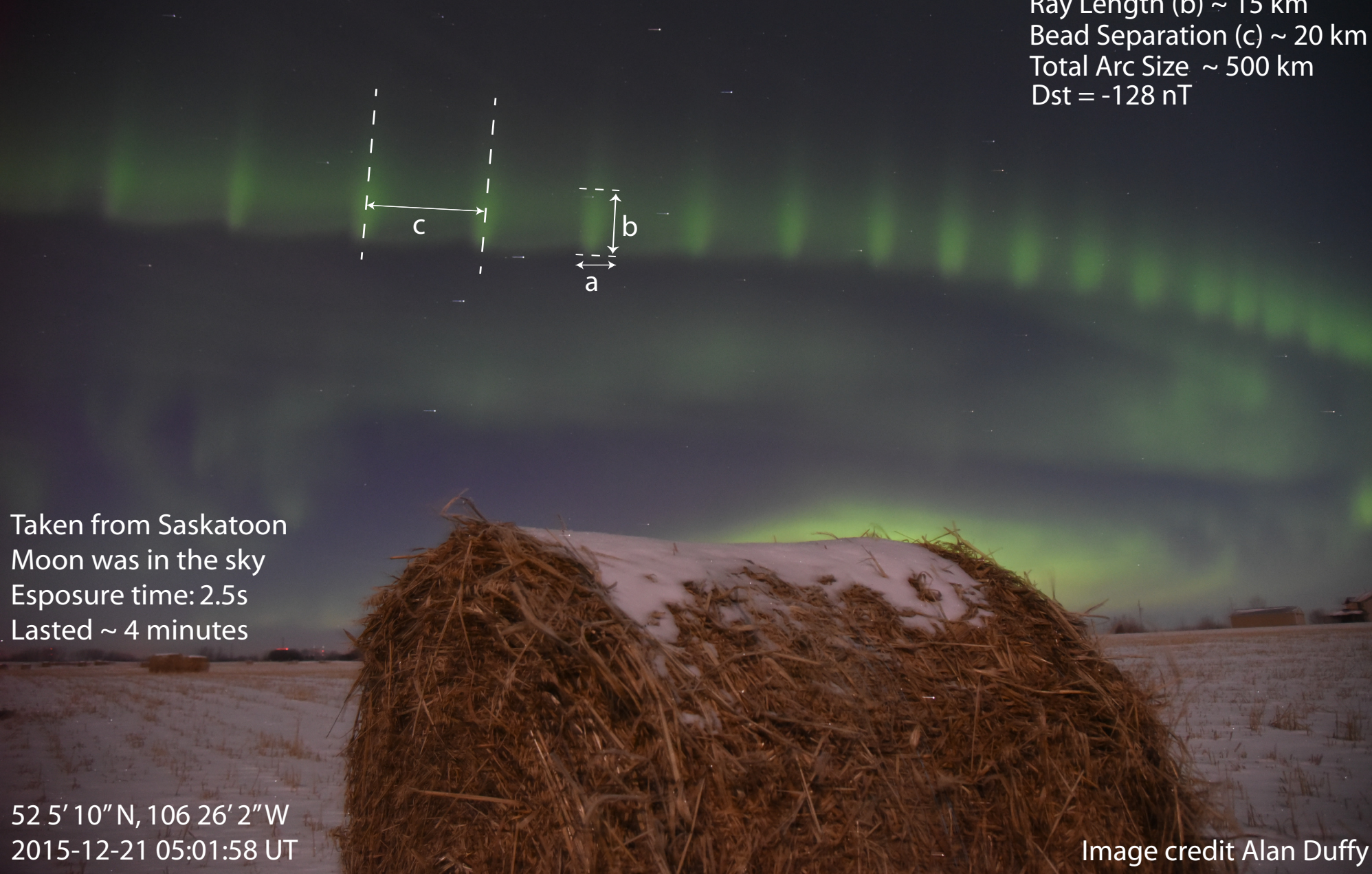


Figure 7.

(a) Citizen Science Image of Unusual Auroral Beads



(b) STEVE - the subauroral arc

