

Adult Content Consumption in Online Social Networks

Mauro Coletto · Luca Maria Aiello · Claudio
Lucchese · Fabrizio Silvestri

Final version available at <https://doi.org/10.1007/s13278-017-0449-y>

Abstract Users in online social networks naturally organize themselves into overlapping and inter-linked communities that are formed around common identity or shared topical interests. Some communities gather people around specific *deviant behaviors*, conducts that are commonly considered inappropriate with respect to the society's norms or moral standards such as drug use, eating disorders, and pornographic content consumption. From a network analysis perspective, the set of interactions between members of these communities form *deviant networks* that map how the deviant content is shared and consumed. It is commonly believed that deviant networks are small and isolated from the mainstream social-media life; accordingly, most research studies have considered them in isolation.

We focus on adult content consumption networks, which is one *deviant network* with a significant presence in on-line social media and in the Web in general. We investigate two large on-line social networks and discuss the following insights. *Deviant networks* are limited in size, tightly connected and structured in subgroups. Nevertheless, content originated in *deviant networks* spreads widely across the whole social graph possibly touching a large number of *unintentionally exposed* users, such that the average local perception is that neighboring users share more deviant content. Finally, we investigate how content production and consumption varies with age and

Mauro Coletto
IMT Lucca - Ca' Foscari University of Venice
E-mail: mauro.coletto@imtlucca.it

Luca Maria Aiello
Nokia Bell Labs
E-mail: luca.aiello@nokia.com

Claudio Lucchese
CNR Pisa
E-mail: claudio.lucchese@isti.cnr.it

Fabrizio Silvestri
CNR Pisa
E-mail: fabrizio.silvestri@isti.cnr.it

show that the consumption rate is very similar between male and female users up to the age of 25.

We conclude that deviant communities are deeply rooted into the relational fabric of a social network, and that a deeper understanding of how their activity impacts on every other user is required.

Keywords Deviant network · Deviant behaviour · Pornography · Adult content consumption · Sexual content production · Social Media · Online social network · Tumblr · Flickr

1 Introduction

The structure of online social networks is fundamentally related to the interests of their members. People assort spontaneously based on the topics that are relevant to them, forming social groups that revolve around different subjects. This tendency has been observed with quantitative studies in several online social media [4, 47]. In the past, researchers have explored the relationship between information diffusion and network structure [8], focusing on the structural and dynamical properties of specific topical communities such as groups supporting political parties [20], or discussion groups about rumors, hoaxes [60] and conspiracy theories [9].

Online social media are also favorable ecosystems for the formation of topical communities centered on matters that are not commonly taken up by the general public because of the embarrassment, discomfort, or shock they may cause. Those are communities that depict or discuss what are usually referred to as *deviant behaviors* [18]; these are conducts that are commonly considered inappropriate because they are somehow violative of society's norms or moral standards. Pornography consumption, drug use, excessive drinking, eating disorders, or any self-harming or addictive practices are all examples of deviant behaviors. Many of them are represented, to different extents, on social media [22, 35, 52]. However, since all these topics touch upon different societal taboos, the common-sense assumption is that they are embodied either in niche, isolated social groups or in communities that might be quite numerous but whose activity runs separately from the mainstream social media life. In line with this belief, research has mostly considered those groups in isolation, focusing predominantly on the patterns of communications among community members [70] or, from a sociological perspective, on the motivations to that make people join such groups [7].

In reality, people who are involved in deviant practices are not segregated outcasts, but are part of the fabric of the global society. As such, they can be members of multiple communities and interact with very diverse people, possibly exposing their deviant behavior to the public. We aim to go beyond previous studies that looked at deviant groups in isolation by observing them *in context*. In particular, we want to shed light on three matters that are relevant to both network science and social sciences: *i*) how much deviant groups are structurally secluded from the rest of the social network, and what are the characteristics of their sub-groups who build ties with the external world; *ii*) how the content produced by a deviant community spreads and what is the entity of the diffusion which reaches users outside the boundaries of the

36 deviant community who voluntarily or inadvertently access the adult content; and *iii*)
37 what is the demographic composition of producers and consumers of deviant content
38 and what is the potential risk that young boys and girls are exposed to it.

39 In this initial study we undertake to answer those questions focusing on the be-
40 havior of *adult content* consumption. Public depiction of pornographic material is
41 considered inappropriate in most cultures, yet the number of consumers is strikingly
42 high [63]. Despite that, we are not aware of any study about the interface between
43 adult content communities and the rest of the social network. The approach followed
44 in this study looks at a deviant network *in context*. The same kind of analysis could
45 be conducted on any other *deviant* topic.

46 We studied this phenomenon on two large dataset sampled from the Tumblr and
47 Flickr social networks. The Tumblr dataset contains more than 130 million users and
48 almost 7 billion directed dyadic interactions, while the Flickr dataset contains more
49 than 39 million users and almost 600 million directed dyadic interactions. In both
50 cases, we selected *deviant* users with a vocabulary-based approach. In Tumblr we
51 used a large sample of 146 million queries from a 7-month log of search queries from
52 Yahoo Search, from which we identified Tumblr pages clicked in response to *deviant*
53 queries. In Flickr we followed a similar approach by looking at image tags. The adult
54 dictionary is made publicly available [19].

55 Results show that:

- 56 – The deviant network is a tightly connected community structured in subgroups,
57 but it is linked with the rest of the network with a very high number of ties (Sec-
58 tion 4.1).
- 59 – The vastest amount of information originating in the deviant network is produced
60 from a limited core of nodes but spreads widely across the whole social graph,
61 potentially reaching a large audience of people who might see that type of content
62 unwillingly, depending on the sharing actions enabled by each social platform.
63 Although the consumption of deviant content remains a minority behavior, the
64 average local perception of users is that neighboring nodes reblog more deviant
65 content than they do (Section 4.2).
- 66 – There are clear differences in the age and gender distributions between producers
67 and consumers of adult content. The differences we found are compatible with
68 previous literature on adult material consumption: producers are older and more
69 predominantly male and age greatly affects the consumption habit, strengthening
70 it in males and weakening it in females (Section 4.3). Moreover if we look at age
71 distribution of the consumers we recognize a similar consumption pattern both in
72 Tumblr and in Flickr which is very different gender by gender: for male users the
73 consumption of adult content increases with age until a maximum around 40-55;
74 for female users the consumption increases only in youth to progressively sub-
75 stantially decrease. Relatively to the population of each social network analyzed
76 we verify a higher consumption of pornography by male users, but interestingly
77 male and female have similar consumption relative volumes up to the age of 25
78 on both networks.

79 **Summary of the contributions.** To the best of our knowledge, this is the first study
80 that analyzes the production and consumption of adult content in general-purpose

81 online social networks, at very large scale. The computer science and social science
82 literature in this area has focused so far on specific, small-scale deviant communities
83 of nodes without analyzing this type of behavior in the broader social context these
84 communities are immersed in. We study the phenomenon of adult content production
85 and consumption for the first time *in context* by considering the whole social network
86 of active users that surrounds the producers of pornographic material. The dataset we
87 study is original, includes two social networks different in scope and structure, and it
88 has no precedents for its size. This unique experimental setup allows us to study for
89 the first time to what extent adult content spills over the restricted groups and spreads
90 in the network; what we find provides strong evidence against the common assump-
91 tion that adult content is mostly confined inside groups dedicated to that specific
92 topic and provide a detailed analysis of the dynamics that characterize the spread-
93 ing of adult content. Furthermore, the scale of the data allowed us to answer some
94 longstanding social science questions around demographic patterns of pornographic
95 material consumers: this is the first study that provides a very large scale quantitative
96 measurement of adult content consumption across age and gender.

97 **2 Related Work**

98 Our contribution is related with previous work aimed at characterizing community
99 dynamics, deviant content consumption, and potential impact of deviant communities
100 on people's behavior.

101 **Groups in online social media.**

102 Computer science research has dealt extensively with the problem of classification
103 of groups along structural, temporal, behavioral, and topical dimensions [2, 33, 53].
104 The relationship between group connectivity and shape of information cascades has
105 also been explored, revealing an intertwinement between community boundaries and
106 cascade reach that is particularly tight in communities built upon a common theme
107 shared by all of their members [8, 27, 50, 61]. The degree of inter-community in-
108 teraction has been analyzed mostly in the context of heavily polarized networks,
109 the most classical example being online discussions between two opposing politi-
110 cal views [1, 20, 29]. These studies explored methods to quantify segregation [34],
111 but mainly focus on networks formed by two main divergent clusters.

112 **Deviant communities.**

113 A body of work has investigated the dynamics of misbehavior in online communities
114 including newsgroups [57], question-answering portals [40], chats [67], and multi-
115 player video games [10, 66], trying to quantify the negative impact of misconducts
116 on the community's health [21, 26, 72].

117 Studies about the depiction of drug and alcohol use in social media adopted
118 mainly the content perspective. Researchers aimed at identifying the elements that
119 boost content popularity, investigated the effect of gender on engagement, and stud-
120 ied the perceptions that deviant content arises in the young public [52]. Research
121 has been conducted around anorexia-centered online communities [12, 32, 59], also

122 on Tumblr [22], investigating a wide range of aspects including the construction and
123 management of member identities, the processes of social recognition, the emergence
124 of group norms, and the use of linguistic style markers. Similar studies have been pub-
125 lished over the years on communities of self-injurers and negative-enabling support
126 groups, in which members encourage negative or harmful behaviors [35]. Fewer stud-
127 ies touch upon network-related aspects. One notable example is the work by Gareth et
128 al. ([70]) that provides an overview of behavioral aspects of users in the PornHub so-
129 cial network, with particular focus on the role of sexuality and gender. More loosely
130 related are studies on the so-called *dark networks*, mostly motivated by the need of
131 finding effective methods to disrupt criminal or terroristic organizations [74]. The
132 study by Christakis et al. ([17]) about the communication network between smokers
133 and non-smokers is one of the few quantitative studies that addresses the interaction
134 between the social network and one of its sub-groups, but it strongly focuses on the
135 phenomenon of contagion.

136 **Adult content consumption.**

137 In the context of internet pornography consumption, computer science literature stud-
138 ied the categorization of content and frequency of use [38, 65, 69]. A wider corpus
139 of research has been produced by social and behavioral scientists by means of sur-
140 veys administered to relatively small groups. Special attention has been given to the
141 relationship between age/gender and the exposure (voluntary or unwanted) to the in-
142 ternet pornography [13, 15, 51, 63, 75], with particular interest to the age range of
143 young teens [15, 51, 73]. Numbers vary substantially between studies, but clearly
144 men are more exposed than women (approximately 75%-95% vs. 30%-60%), with
145 men exposed more frequently [36] and women more often involuntarily. It is esti-
146 mated that young teens that are often exposed accidentally (roughly 25% to 66% of
147 the times) and are also exposed to violent or degrading pornography (20% among
148 female, 60% among male) [62]. Researchers have also pointed out the potential harm
149 that adult material consumption through internet can cause, including addiction [42]
150 and increased chance of adopting aggressive behavior [5]. Exposition also correlates
151 with drug use [75] and with lack of egalitarian attitude towards the other sex [37]. Al-
152 though delving into the potential harm of pornography is far beyond the scope of our
153 work, this inherent risks provide an additional motivation to focus on this particular
154 type of deviant community.

155 **3 Deviant graph extraction**

156 Online social media platforms provide users the capability of *publishing* content, ac-
157 cessing content published by other users of the network, and interacting with them.
158 We model the social network as a graph whose nodes are the users and whose edges
159 are the observed *interactions* among them. We distinguish among three kinds of in-
160 teractions: *i) following* users enables a to receive updates from their content stream;
161 *ii) liking* any piece of content produced by others expresses explicit interest in the
162 activity of others; *iii) sharing* content produced by others increases the visibility of
163 an item by re-posting it on the sharer's feed.

164 Each specific social network may re-brand these actions to fit the goals of the
165 network; this study focuses on two online platforms: Tumblr and from Flickr. Tumblr
166 is a popular micro-blogging platform where users can publish content by *posting* new
167 entries on their blogs usually containing multimedia content, they can share content
168 by *reblogging* any other post on their blogs, and *follow* other users. Flickr is a photo-
169 sharing platform, including both amateur and professional users. Users can publish
170 content by *posting* new photos or videos in their *stream*. They cannot re-share content
171 but they can *like* photos and *follow* other users. Those platforms also provide some
172 other actions, e.g., commenting a photo, but the analysis of rich textual signals goes
173 beyond this work.

174 Tumblr and Flickr are ideal platforms for this type of study for three main rea-
175 sons. First, they are two general-purpose sites that contain a wide variety of topical
176 communities, they do not focus on few specific themes. Second, they do not enforce
177 any content restriction policy around pornographic imagery (unlike Instagram, for
178 example). Last, the two platforms are very different in scope, size, typical usage, and
179 demographics of the user base, which helps us drawing more general conclusions.

180 We note that, in both platforms, the same human user could, in principle, own
181 multiple user accounts (aka blogs). For the purpose of this study we consider blogs
182 as users, and we will use the two terms interchangeably. We believe blogs are a good
183 unit of analysis for the purpose of this study. All the dynamics of both social networks
184 under investigation happen at the blog level: the following relationships and all other
185 social actions (e.g., liking, reblogging) are done among blogs. Given that, it would be
186 undesirable to coalesce different blogs owned by the same person into a single node.
187 Furthermore, accurately matching blogs to actual users is hard even if one had the
188 full information about all users. In fact, users can register their multiple blogs under
189 different emails and using different nicknames. Historically, previous work on Flickr
190 and Tumblr adopted the same assumption as we do and studied the network of blogs,
191 not the network of users [22, 53].

192 We consider as *deviant nodes* those users who publish content about a given
193 *deviant topic*: in our study the deviant behavior under analysis is the pornography
194 production and consumption. To identify deviant nodes we resort to data from Web
195 search engine logs and tags. We first discuss the methodology adopted on the Tumblr
196 data, and then how this was adapted for the Flickr network.

197 As shown in previous studies analyzing pornography consumption in the context
198 of Web search [45], the higher the number of deviant queries *hitting* (i.e., leading to
199 the click of) a page, the higher the probability that page contains deviant content. In
200 Tumblr, to identify a pool of candidate *deviant nodes* we consider deviant queries
201 hitting the URLs of Tumblr blogs.

202 We used a seven-month long anonymized query log from Yahoo search engine,
203 from which we collected a random sample of 146M US query log entries whose
204 clicked URL belongs to the `tumblr.com` domain. After a simple query normaliza-
205 tion process we obtained about 26M unique queries that hit a total of 2.7M unique
206 Tumblr blogs. As expected, the distribution of number of queries hitting a blog is
207 very skewed, with most popular blogs being reached by hundreds of thousands of
208 clicks originating from search queries (Figure 1).

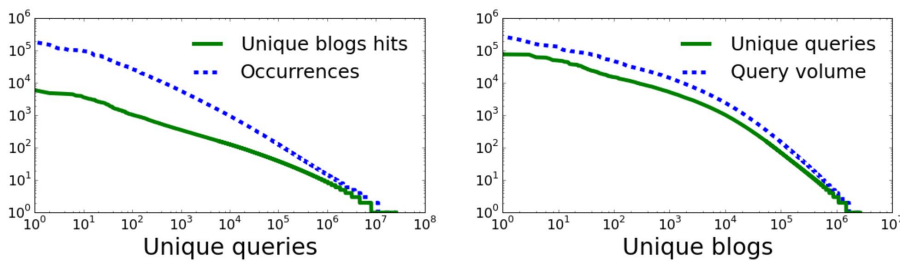


Fig. 1: Distributions of: (left) number of blogs hit by a query and number of occurrences of a query; (right) volume of (unique) queries hitting a blog.

209 To maximize the accuracy and coverage of the set of discovered deviant nodes,
 210 we devise an iterative semi-supervised *Deviant Graph Extraction* procedure [19].
 211 The procedure detects adult queries by mean of a dictionary of terms: adult-marked
 212 queries are used to detect adult blogs which are their landing pages and the process is
 213 iterated until convergence extending the dictionary with new terms extracted by fre-
 214 quent queries which point to highly adult blogs previously detected, i.e. blogs whose
 215 incoming query volume is almost entirely pornographic. Figure 2 shows that the *De-*
 216 *viant Graph Extraction* procedure converges quickly. We were able to identify 198K
 217 adult blogs, through 4.2M unique queries filtered by mean of a dictionary which was
 218 expanded at each iteration of the procedure including in the final step 7,361 terms.

219 Unlike in Tumblr, our Flickr datasets contains user-generated textual annotations
 220 that carry information about the content of published pictures, which allows us to
 221 directly estimate the pool of deviant content without the need of relying on external
 222 proxies¹.

223 To detect adult photos, we matched the vocabulary obtained through the *Deviant*
 224 *Graph Extraction* procedure with Flickr tags. Each picture in Flickr is labeled with
 225 manual tags. On average each picture has 6 tags and each user publishes around 122
 226 pictures. We slightly modified the adult vocabulary to remove misleading words in
 227 a photographic context (e.g. black and white) and we filtered photos labeled with at
 228 least one tag in the adult dictionary. We considered only users with at least two public
 229 adult photos identified as above in line with the query approach used for Tumblr
 230 where we marked a blog as adult only if it was reached through at least two unique
 231 adult queries.

232 This procedure resulted in about 6.5 million photos by about 73K deviant users.
 233 In Flickr users share content not only in their profile but the platform enables the
 234 creation of groups whose member share images according to the topic of the group.
 235 To improve the recall of the data collection, we also identified those groups which *i)*
 236 include at least one of the previously detected adult users and *ii)* with a group title
 237 matches at least one word in our adult vocabulary. All the users and photos in such
 238 groups were included in the adult cluster. In so doing, about than 10M photos and
 239 175K *deviant* users were detected in Flickr.

¹ Flickr users can mark their own photos as “adult”. We first attempted to use this self-reported information to detect adult photos. We found that this approach leads to many false positives, mainly because very often pictures are marked in big batches containing adult and non-adult pictures.

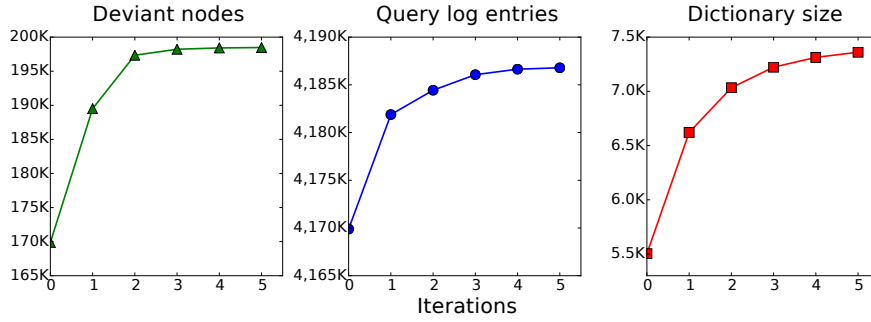


Fig. 2: Convergence of number of deviant nodes, query log entries and adult dictionary size during the Deviant Graph Extraction procedure.

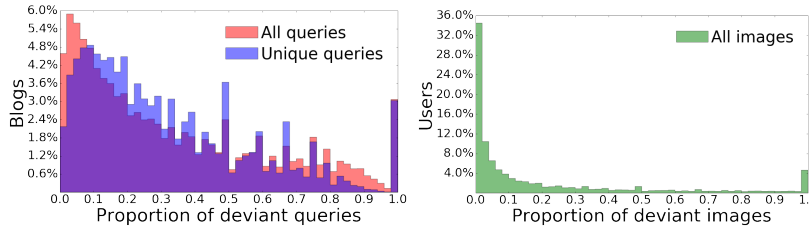


Fig. 3: (left) Distribution of deviant query volume ratio reaching deviant nodes in Tumblr, and (right) distribution of deviant photo volume ratio by adult user (publisher) in Flickr.

Table 1: Deviant users detected for Tumblr and Flickr: $|U|$, $|U^*|$ are users which are connected at least with one of the other users in $|U|$ through an incoming or outgoing reblog link (Tumblr) or favorite link (Flickr). $|U^*|$ in GCC are users in the giant component considering the reblog graph for Tumblr and the favorite graph for Flickr.

	$ U $	$ U^* $	$ U^* $ in GCC
Tumblr	105K	75.6K (72%)	75.2K
Flickr	171K	156.1K (91%)	156.1K

240 In Figure 3 we report the distribution of the *deviant query volume* ratio for Tumblr
 241 and the *deviant image volume* ratio for the deviant users detected in Flickr. For both
 242 cases we see that the distribution is skewed, showing a tail of blog/user which are hit
 243 by a majority of deviant queries (Tumblr) or who published a high portion of adult
 244 pictures over their uploads.

245 To study the interaction of deviant nodes with the rest of the social network, we
 246 extracted a subset of the Tumblr and Flickr social network with a snowball expansion
 247 starting from the identified deviant nodes up to 3-hops away. Deviant nodes detected
 248 are reported in Table 1. For Tumblr we considered the following and reblog actions
 249 while for Flickr since the reblog action is not present we collected information about
 250 favorites (or likes), which express a similar engagement even though they do not par-
 251 ticipate in the propagation of the content. The Tumblr follower network is a snapshot

Table 2: Network statistics for the reblog (R), follow (F), and favorite (L) networks of the full graph sample (*All*), the deviant graph (*Deviant*), and the communities that compose it (*Producers* and *Bridges*). All the statistics are about the giant weakly connected components and count only links whose both endpoints are in the considered node subset. $\langle k \rangle$ =average degree, D =density, ρ =reciprocity, C =clustering, \overline{spl} =average shortest path length, d =diameter.

		$ N $	$ E $	$\langle k \rangle$	D	ρ	C	\overline{spl}	d
<i>Tumblr</i>	All R	14M	472M	33	$2 \cdot 10^{-6}$	0.06	-	-	-
	All F	130M	6,892M	53	$4 \cdot 10^{-7}$	0.10	-	-	-
	Deviant R	105K	1.4M	13	$1 \cdot 10^{-4}$	0.04	0.10	3.73	11
	Deviant F	135K	24.6M	182	$1 \cdot 10^{-3}$	0.07	0.13	2.80	8
	Prod₁ R	48K	914K	19	$4 \cdot 10^{-4}$	0.04	0.09	3.44	9
	Prod₂ R	16K	305K	19	$1 \cdot 10^{-3}$	0.05	0.13	3.19	8
	Bridge₁ R	9K	36K	4	$5 \cdot 10^{-4}$	0.04	0.08	4.18	13
	Bridge₂ R	3K	32K	11	$4 \cdot 10^{-4}$	0.06	0.21	3.32	10
<i>Flickr</i>	All L	15M	553M	37	$2 \cdot 10^{-6}$	0.06	-	-	-
	All F	39M	566M	15	$4 \cdot 10^{-7}$	0.26	-	-	-
	Deviant L	171K	13.4M	79	$5 \cdot 10^{-4}$	0.03	0.17	3.06	9
	Deviant F	169K	37.9M	224	$1 \cdot 10^{-3}$	0.28	0.21	2.77	9
	Bridge₁ L	66K	2.7M	47	$6 \cdot 10^{-4}$	0.05	0.17	3.05	13
	Prod₁ L	53K	4.6M	99	$2 \cdot 10^{-3}$	0.03	0.18	2.83	13
	Prod₃ L	20K	1.5M	94	$4 \cdot 10^{-3}$	0.03	0.23	2.53	13
	Prod₂ L	16K	1.0M	83	$4 \cdot 10^{-3}$	0.04	0.28	2.52	14

252 of the graph done in December 2015; the reblog network was built from the reblog
 253 activity happened in the same month. The Flickr follower network is a snapshot of
 254 the graph done in March 2016; the favorite network was built from the like activity
 255 happened until the snapshot time. The resulting networks are discussed in the Sec-
 256 tion 4.

257 We also obtained information about self-declared age and gender for about $1.7M$
 258 Tumblr users and $12.3M$ Flickr users. The datasets include exclusively interactions
 259 between users who voluntarily opted-in for such studies. All the analysis we report
 260 next was performed in aggregate and on anonymized data.

261 4 Deviant graph in context

262 The availability of data about the interaction between deviant nodes and the social
 263 network that surrounds them provides the unique opportunity to study the structure
 264 and dynamics of a deviant network within its context. We first analyze the shape of the
 265 deviant network and measure its connectivity with the rest of the social graph (Sec-
 266 tion 4.1). We then look into how the information originating from deviant networks
 267 spreads across the boundaries of the deviant group (Section 4.2). Last, we study some
 268 demographic properties that characterize producers and consumers (Section 4.3).

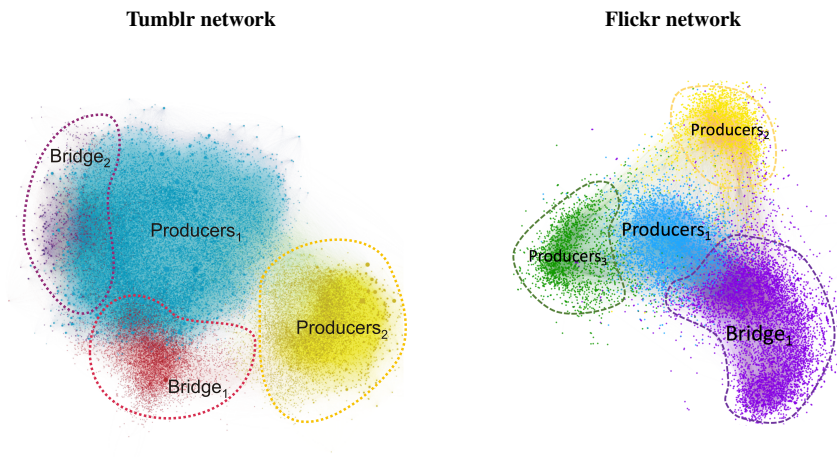


Fig. 4: Bird-eye view of the deviant network for Tumblr (reblog network, left) and Flickr (favorite network, right) with colors and labels denoting algorithmically-extracted communities.

269 4.1 Deviant network connectivity

270 The deviant network is a tiny portion of the whole graph, representing about 0.8%
 271 of all the nodes in the reblog graph in Tumblr, 1.1% of all the nodes in the favorite
 272 graph in Flickr and a even smaller portion in the follow network (0.1% in Tumblr and
 273 0.4% in Flickr). So few nodes could be scattered along the social network or clustered
 274 together. So we ask:

275 *Q1) Are deviant nodes organized in a community?*

276 We consider the *deviant networks* as the subgraphs of the follow and reblog/favorite
 277 networks induced by the *deviant nodes* in Tumblr and in Flickr. A directional link
 278 in the follow (reblog/favorite) network from node i to node j exists if i follows (or
 279 reblogs/likes the posts of) j , meaning that the information flows from j to i . Basic
 280 network statistics on such subgraphs reveal that the deviant networks are quite dense,
 281 yet they have a high diameter (Table 2). Similar statistics have been observed before
 282 in other social networks [3] and might be an indication of the presence of strong
 283 sub-groups patterns, as well as a signal of the absence of a community structure. To
 284 better determine the reason for such elongated shape, we run the Louvain community
 285 detection algorithm [11] on the deviant network². We considered the reblog network
 286 for Tumblr and the favorite network for Flickr as they are more representative than
 287 follow network of the dynamic interactions of the users in the OSN. Only the giant
 288 connected component of the network has been considered which corresponds to 72%
 289 of the network in Tumblr and 91% in Flickr (see Table1). The modularity measure
 290 for the clustering in Tumblr is 0.44 and 0.53 in Flickr; those are high values of modu-

² Louvain is a modularity-based graph clustering algorithm that shows very good performance across several benchmarks [31] and that is fast to compute even on large networks.

291 larity [54] that indicate the presence of structured, well-separated communities. Four
292 clusters emerge both in Tumblr and in Flickr, whose network statistics are summa-
293 rized in the bottom lines of Table 2. To determine their nature, we manually inspected
294 the content of 250 blogs in each of them.

295 In Tumblr more than 90% of all the blogs in the two largest clusters contain
296 blogs that *exclusively* produce explicit adult content, aimed at an heterosexual public
297 (*Producers₁*) or at a male homosexual public (*Producers₂*). The blogs in the two
298 remaining communities post less explicit adult content and more sporadically, often
299 by means of reblogging. They either focus on celebrities (*Bridge₁*), or function as
300 aggregator blogs with high content variety, including depiction of nudity (*Bridge₂*). In
301 Flickr, the composition is different with similar attributes. We found (*Producers₁*) and
302 (*Producers₂*) clusters with the same content characterization described for Tumblr
303 and in addition to them a new cluster has been identified (*Producers₃*) whose users
304 share mainly pictures representing transvestites or transsexuals. The same cluster is
305 likely present in Tumblr but its size is not large enough to be distinguished by other
306 producers. Producers clusters in Flickr are less than 58% of the deviant nodes with a
307 large bridge cluster (*Bridge₁*) which is characterized by a content less explicit (soft
308 porn, artistic nudity, manga).

309 From a bidimensional visualization of the network layout (Figure 4) it becomes
310 apparent that the *Producers₁* and *Producers₂* are two well-separated cores. In Flickr
311 *Producers₃* is very close to *Producers₁*. The remaining communities are peripheral
312 and arranged in a crown-like fashion in Tumblr (which explains their high diameter)
313 around the largest sub-cluster *Producers₁*; in Flickr instead the largest cluster is very
314 elongated (showing the highest diameter) indicating a strong presence of soft content
315 which from a network point of view is organized in long cluster connected with the
316 hard part only with one side. We named all the non-producers groups *bridge commu-*
317 *nities* as their main focus is often not on deviant content and, as we shall see next,
318 they act also as link towards the rest of the graph.

319 In short, we find that deviant nodes are not scattered in the social network but are
320 tightly organized in a structure of distinct communities. To find out about the nature
321 of their interaction with the rest of the social ecosystem, we proceed to answer the
322 next question.

323 *Q2) To what extent is the deviant graph connected to the rest of the social network?*

324 There are several ways to estimate the connectivity between two sets of nodes in a
325 graph. We use different metrics to measure it between the four communities of the
326 deviant network and the rest of Tumblr, as summarized by the matrices³ in Table 3;
327 rows represent the group of nodes from which the social tie originates, columns those
328 on which it lands.

329 The average volume of connections (left matrix) provides a first indication about
330 the difference in connectivity across different groups. For both Tumblr and Flickr,
331 the diagonal has the highest values because of the community structure of the deviant
332 network and of its sub-communities: members of a group have many more ties to-

³ Link directionality is considered: ties originate from groups listed on the rows and land on groups listed on the columns.

		Average volume					Density ($\cdot 10^{-2}$)					Null model comparison					
		P_1	P_2	B_1	B_2	O	P_1	P_2	B_1	B_2	O	P_1	P_2	B_1	B_2	O	
Tumblr	Follow	P_1	463	13	7.9	13	582	0.702	0.788	0.917	4.744	0.004	1165	94	110	569	0.5
		P_2	27	443	4.6	1.2	635	0.071	2.743	0.538	0.429	0.005	66	3199	62	50	0.6
		B_1	21	4.5	40	2.7	484	0.442	0.278	4.573	0.951	0.004	103	65	1074	223	0.9
		B_2	220	6.3	17	131	598	4.591	0.388	1.911	4.651	0.005	612	51	255	6205	0.6
		O	2.4	0.7	1.7	0.2	47	0.051	0.045	0.293	0.046	10^{-5}	125	112	487	165	0.9
Tumblr	Reblog	P_1	19	0.3	0.2	0.5	30	0.401	0.019	0.026	0.177	0.002	113	5.4	7.5	50	0.6
		P_2	0.6	19	0.2	0.01	39	0.012	0.167	0.017	0.014	0.003	3.0	281	4.2	3.4	0.7
		B_1	0.9	0.1	4.3	0.2	63	0.018	0.008	0.493	0.076	0.004	3.8	1.6	102	16	0.9
		B_2	7.0	0.2	0.8	11	44	0.147	0.013	0.008	4.018	0.003	33	2.8	22	897	0.7
		O	0.8	0.3	1.1	0.1	31	0.016	0.016	0.127	0.009	0.002	6.6	6.7	54	17	0.9
Flickr	Follow	B_1	99	33	8.1	5.2	388	1.499	1.634	0.153	0.216	0.01	109	119	11	23	0.7
		P_1	109	645	71	28	655	1.658	3.174	1.341	1.745	0.017	42	821	34	45	0.4
		P_2	9.8	27	89	4.1	91	0.149	0.155	1.685	0.248	0.002	26	237	295	43	0.4
		P_3	25	45	14	119	190	0.387	0.254	0.266	0.735	0.005	38	223	26	715	0.5
		O	1.0	0.7	0.3	0.1	16	0.016	0.031	0.005	0.007	-	49	104	15	23	0.8
Flickr	Favorite*	B_1	41	9.6	2.0	1.3	-	0.627	0.47	0.038	0.00	-	173	130	10	22	-
		P_1	34	225	21	10	-	0.515	1.111	0.396	0.609	-	26	574	20	31	-
		P_2	2.6	11	28	1.8	-	0.079	0.569	0.535	0.11	-	13	192	181	37	-
		P_3	7.3	16	5.0	61	-	0.11	0.828	0.091	3.738	-	18	137	15	621	-
		O	0.7	0.8	0.2	0.1	-	0.011	0.038	0.004	0.009	-	90	300	34	71	-
Flickr	Favorite	B_1	162	22	6.5	6.7	523	2.452	1.125	0.123	0.408	0.035	51	23	2.5	8.5	0.7
		P_1	94	356	43	19	357	1.43	1.755	0.819	1.191	0.024	24	302	14	20	0.4
		P_2	10	20	53	3.6	50	0.156	0.001	0.104	0.216	0.003	16	109	110	23	0.4
		P_3	32	28	0.5	113	147	0.489	1.387	0.175	6.871	0.01	22	63	8.2	313	0.5
		O	4.6	1.5	0.6	0.5	24	0.069	0.078	0.011	0.028	0.002	32	35	5.1	13	0.8

Table 3: (T=Tumblr, F=Flickr) Measures of connectivity between the communities in the deviant network in Tumblr (*Producers* P_1, P_2, P_3 (F) and *Bridges* B_1, B_2 (T)) and the rest of the social network O , for the follow, reblog (T), favorite to adult content or favorite* (F) and total favorite (F) relations.

wards other group members rather than to the outside. This is true in particular for *Producer* clusters. In Tumblr the difference between the diagonal values and the other cells is more prominent, however the volume of links incoming to the largest producer cluster is particularly high from the smallest bridge community (*Bridge₂*), which surrounds it. The average Tumblr user in our sample (see rows *O* in Table 3) follows around 51 users, between 2 or 3 of which are in the core of the deviant network and around 2 of them are in bridge communities; similarly, among the 33 users reblogged in one month by the average user, one is from a *Producer* cluster and one from a *Bridge* group. In Flickr, also the values not in the diagonal are generally high showing an significant connection among clusters: a user follows on average around 12 other users and among them 1 is from the bridge cluster which is the biggest in size and 1 is a member of a producer cluster; a user likes on average pictures from 31 users and around 5 of them are in the bridge cluster and between 2 or 3 are among producers. Since the classification of deviant content in Flickr has been done at picture level through tags, we can distinguish likes to adult content (forth row in Table 3: favorite*) and total likes (fifth row). Among the users liked on average among the bridge cluster only 15% are favorite links on adult content, while for the producers clusters the percentage rises to 42% confirming that the content of the bridge cluster is less explicit and often not adult.

When looking at raw volumes, the amount of links from the deviant network to the rest of the graph is very high, mainly due to the high dimensionality of the set of nodes that are not deviant. To partially account for dimensionality of the groups, we measure the connectivity with density computed as the ratio of edges between the two groups over the total number of possible edges between them (Table 3, center). Also in this case the overall patterns hold, but the connectivity towards the external graph drops significantly.

Values of density are still affected by size, though. It is known that in real networks there is a strong correlation between number of nodes and graph density [48]. To fix that, in the spirit of established work in complex systems [64] we resort to a comparison of the real network connectivity with a *null model* that randomly rewires the links while keeping the degree of each node unchanged. The values we report in Table 3 (right) indicate the ratio between the real observed connections and the null model. Also in this case, values on the diagonal are very high (except for the outer network, which has a value close to 1, as expected). For Tumblr, both the density matrix and the null model comparison show an high connectivity not only for the producers (in particular *Producer₂*), but also for bridge communities. This is particularly evident for *Bridge₂* both in the reblog and in the follow matrix. For Flickr instead the bridge community is characterized by the lower value both in favorite and follow in the diagonal among deviant clusters. Moreover if we look at *Producer₃* in Flickr we discover that is highly connected to *Producer₁*: members of *Producer₃* like and follow significantly members in *Producer₁* cluster while the contrary is much less evident. This confirms the assumption that a similar cluster *Producer₃* might be present for Tumblr but not identifiable since it is small or highly connected with *Producer₁* resulting indistinguishable from it. Also, this computation highlights that ordinary users have a tendency in Tumblr to reblog content from the core of the deviant network almost 7 times more than random and between 17 and 55 times more

379 than random from the bridge community members, in Flickr to like content from
 380 *Producer*₁ almost 35 times more than random (300 times if we look at adult content
 381 only) and 32 times more than random from the bridge community members (90 times
 382 if we look at adult content only).

383 Last, we complement the connectivity analysis with a measurement of the be-
 384 tweenness centrality of nodes in the different communities of adult content producers
 385 (Table 4). The betweenness centrality [71] quantifies the number of times a node acts
 386 as a bridge along the shortest path between any two other nodes in the whole net-
 387 work, so it is a strong indication of the brokerage potential of a node with respect to
 388 information flow. In both networks the largest betweenness is observed for the bridge
 389 clusters, with a peak for the *Bridge*₁ cluster in Tumblr (celebrities). This further find-
 390 ing confirms that bridge cluster play an important role in diffusing adult content to
 391 the rest of the network.

Table 4: Average betweenness centrality C_b of nodes belonging to different adult com-
 munities in Tumblr’s reblog network (R) and Flickr’s like network (L). Centrality
 values are max-min normalized.

<i>Tumblr</i>	C_b
Prod ₁ R	0.0038
Prod ₂ R	0.0061
Bridge ₁ R	0.0133
Bridge ₂ R	0.0061
<i>Flickr</i>	C_b
Bridge ₁ L	0.0074
Prod ₁ L	0.0060
Prod ₃ L	0.0053
Prod ₂ L	0.0049

392 In summary, the core of the deviant community is dense but it is far from being
 393 separated from the rest of the graph, which is connected to it both directly and even
 394 more tightly through bridge groups.

395 4.2 Deviant content reach

396 We found that, although the deviant network forms a tightly connected community, it
 397 is not isolated from the rest of the social graph. This calls for an investigation about
 398 the visibility that the deviant content has in the outer network and what are the main
 399 factors that determine its exposure. We do so by answering four research questions.

400 *Q3) How much deviant content spreads in the social graph and what are the main*
 401 *agents of diffusion?*

402 The exposure to deviant content goes beyond the members of the deviant net-
 403 work who are the *producers* of original adult material. Specifically, the *consumers* of
 404 deviant content can be categorized in three classes. The first is the class of *active con-*
 405 *sumers*: nodes who reblog (but not necessarily follow) adult posts, thus contributing

406 to its spreading along social ties. Posts can be re-blogged in chains and create diffu-
 407 sion trees that potentially spread many hops away from the original content producer,
 408 therefore active consumers could further be partitioned in those who spread the con-
 409 tent *directly* from the producers and those who do it from *indirect* reblogs. This class
 410 is present in Tumblr; in Flickr sharing actions are not made available but users can
 411 see the pictures liked by other people they follow in their feed. For this reason we
 412 called *active consumers* in Flickr all the users who like adult content. The favorite
 413 action is stronger than the following action which characterizes passive consumers.
 414 This influences the spreading dynamic which is limited in Flickr compared to Tum-
 415 blr also for these infrastructure and visibility limitations. The second is the class of
 416 *passive consumers*: nodes who do not contribute to the information diffusion process
 417 but are explicitly interested in adult content because they directly follow the producer
 418 nodes or the like their content. In Flickr we will distinguish passive consumers by
 419 looking at the actions: follow or both follow and favorite actions to the adult con-
 420 tent generated by producers. The last class is the one of *involuntary consumers* or
 421 *unintentionally exposed* users. In Tumblr, users in this class are the ones who do not
 422 follow any producer node and do not reblog their content, but happen to follow at
 423 least one active consumer who pushes adult content in their feed through reblogging.
 424 In Flickr, we consider unintentionally exposed all the users who follow people who
 425 liked adult content at least once; again this choice is motivated by the fact that the
 426 Flickr feed shows the pictures that neighbors recently liked. This way of estimating
 427 involuntary consumers users is a best-effort approach that provides an upper bound
 428 on the number of people who are actually exposed to adult content; as we do not have
 429 access to click logs, there is no way to know which pages users visited.

430 By drawing a quantitative description of the volume of deviant content reach-
 431 ing these three classes we can estimate how much the adult community is visible in
 432 the network at large. We adopt a conservative approach in which we consider the
 433 *Producers* communities as the only ones generating original explicit content. Given
 434 the results of the aforementioned manual inspection, we are very confident that their
 435 activity is mainly focused on the production of adult material.

436 **Tumblr.** We measure the size of the different consumer classes and the amount of
 437 content that flows through or to them by means of reblogging. The results are sum-
 438 marized by the schema in Figure 5. The network of deviant content producers is very
 439 small but receives a considerable amount of attention from direct observers. The au-
 440 dience of passive consumers counts almost 24M people. Around 2M users reblog
 441 directly from the deviant network, for a total of around 28M reblog actions in one
 442 month. A consistent part of the two *Bridge* communities within the deviant graph (a
 443 total of 3K users) are also direct consumers, and they reblog *Producers* 56K times per
 444 month. When looking at the set of 2.4M users who indirectly reblog deviant content,
 445 we see that only a small fraction of their monthly reblogs (less than 7%) is performed
 446 through bridge communities. However, in relative terms, bridge communities are con-
 447 siderably more efficient in spreading information than the average active consumer.
 448 If we consider efficiency η of a user set U as the ratio between reblogs done r_d and
 449 reblogs received r_r , weighted by the cardinality of the set $\left(\eta = \frac{r_r}{r_d \cdot |U|}\right)$, we discover
 450 that the bridge communities ($\eta = 1.5 \cdot 10^{-3}$) are several orders of magnitude more

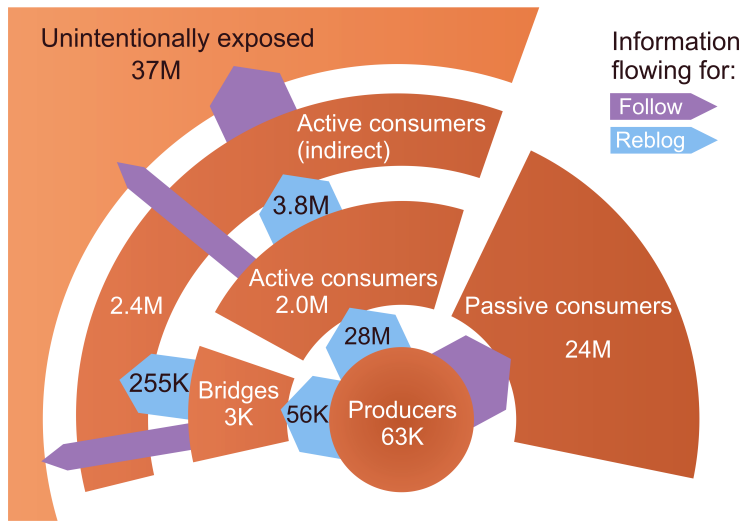


Fig. 5: Diffusion of deviant content from the core of Producers to the rest of the network in Tumblr. Sectors represent disjoint user classes and arrows encode the information flow between them. Reblog arrows report the total volume of reblogs between two classes.

451 effective in spreading the content farther away in the network than the rest of active
 452 consumers ($\eta = 6.7 \cdot 10^{-8}$). Last, the audience of users who are potentially exposed
 453 in an unintentional way to deviant content includes almost 40M people. This figure
 454 should be considered as an upper bound on the number of people who actually have
 455 been exposed, as a follower of an active consumer might not see the pieces of deviant
 456 content for a number of reasons (e.g., inactivity, amount of content in the feed). That
 457 said, the pool of people who are potentially exposed is still very wide.

458 **Flickr.** Similarly for Flickr we quantified the sizes of different classes. Because of the
 459 absence of resharing actions in Flickr we focused on *passive consumers* and *uninten-*
 460 *tionally exposed users* only. The results are summarized by the schema in Figure 6.
 461 The size of the producer clusters (90K) is almost 43% bigger than the case of Tum-
 462 blr but still very small compared to the whole network which is composed by 39M
 463 users. The size of the passive consumers instead is smaller in comparison with the
 464 same class in Tumblr: around 2M users, accessing deviant content by following the
 465 producers. The favorite action is quite limited: only 226K users exclusively like pro-
 466 ducers' content and 475 K like and follow producers at the same time. For Flickr the
 467 bridge cluster has a role different from the bridges communities in Tumblr: resharing
 468 actions are not enabled but they have an important role since their content is an entry
 469 point to access adult content by navigating the social network and the followers. In
 470 particular around 37% of them follow the producers and around 48% of them are in
 471 group of users who like and follow the producers. Last, the audience of users who
 472 are potentially exposed in an unintentional way to deviant content includes almost

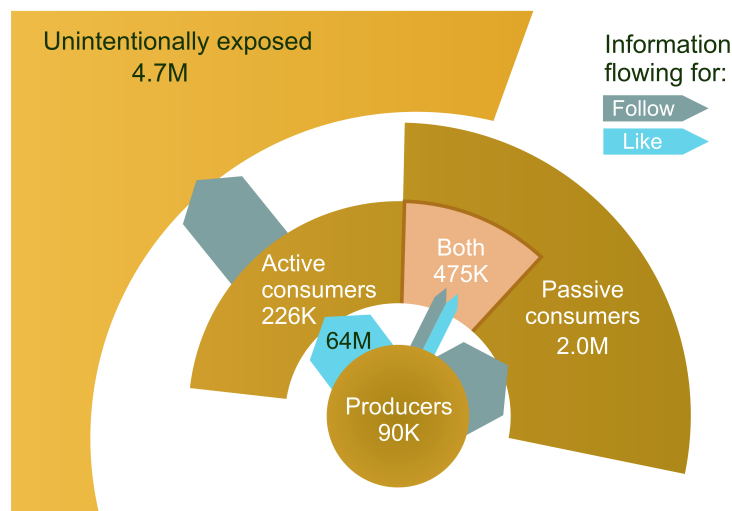


Fig. 6: Diffusion of deviant content from the core of Producers to the rest of the network in Flickr. Sectors represent disjoint user classes and arrows encode the information flow between them.

473 4.7M people. The members of this group like pictures from users who liked at least
 474 one picture shared by the deviant producers. The pool of people who are exposed to
 475 adult content in Flickr is significant but smaller than Tumblr's, in proportion, mainly
 476 because the platform enable less sharing tools (no reblog actions). Also, the two plat-
 477 forms have very different content-production dynamics. Manual inspection suggests
 478 that Tumblr deviant blogs tend to be more topically focused and contain multime-
 479 dia material taken from other Web sources. On the contrary, deviant Flickr users are
 480 mostly amateur who publish pictures that they take, which makes their content more
 481 appealing to their social circles rather than to the general audience, thus restraining
 482 wide diffusion.

483 **Content reach and node coreness.** Past literature has shown how efficient informa-
 484 tion spreaders do not necessarily correspond to the most highly connected nodes in
 485 the network [14]. Kitsak et al. [41] have found instead that the top spreaders are often
 486 those located within the core of the network as identified by a k -core decomposi-
 487 tion [71]. In network analysis, a k -core is a subgraph containing nodes of degree k or
 488 more *within* that subgraph. A k -shell is a set of nodes that belongs to the k -core but
 489 not to the $k + 1$ core. A node that belongs to a given k -shell is said to have a *shell*
 490 *index* equal to k . The shell index is a way to measure the *coreness* of a node in the
 491 graph.

492 We compute the *shell index* k_s on the follower networks for every adult node and
 493 plot it against the total number of (direct and indirect) reblogs the node's content has
 494 received in Tumblr and against the total number of likes received in Flickr (Figure 7).
 495 In Flickr, we have found a trend that is similar to what has been found in previous

496 work, with the node reach growing as the shell increases, with diminishing returns.
 497 In Tumblr, interestingly, the reach is more spread across several orders of magnitude
 498 for fixed values of shell index. This is mainly given by the behavior of bridge clusters
 499 (B1 and B2). Bridge nodes have a higher reach also at lower values of shell index,
 500 which confirms their role as brokers in the process of content spreading: bridges take
 501 advantage of their high-betweenness rather than their high coreness to act as effective
 spreaders.

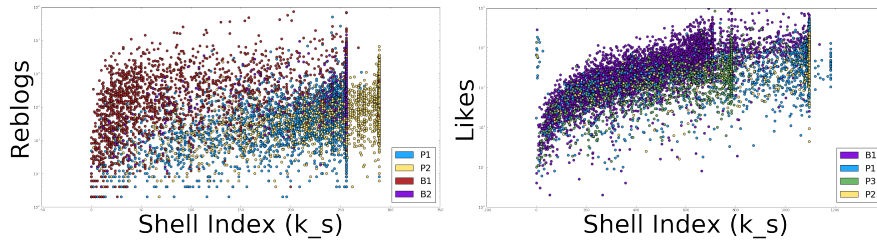


Fig. 7: Correlation between shell index k_s of the social graph (follower network) with the reblogs for Tumblr (left) and with the likes for Flickr (right). Each point represents a blog that produces adult content. Colors encode the different adult sub-communities.

502

503 *Q4) What is the perception of deviant content consumption from the perspective of*
 504 *individual nodes?*

505 Similar to real life, individuals in online social networks are most often aware
 506 of the activities of their direct social connections only but lack a global knowledge
 507 of the behavior of the rest of the population. In fact, the broad degree distribution
 508 of social networks may lead to the over-representation of rather rare nodal features
 509 when they observed in the local context of an ego-network. This phenomenon has
 510 been observed in the form of the so-called *friendship paradox* [28, 39], a statistical
 511 property of social networks for which on average people have fewer friends than their
 512 own friends. More recently the concept has been extended by the so-called *majority*
 513 *illusion* [46], which states that in a social network with binary node attributes there
 514 might be a systematic local perception that the majority of people (50% or more)
 515 possess that attribute even when it is globally rare. As an illustrative example, in
 516 a network where people drinking alcohol are a small minority, the local perception
 517 of most nodes can be that the majority of people are drinkers just because drinkers
 518 happen to be connected with many more neighbors than the average. In our case
 519 study, active deviant content consumption is definitely a minority behavior compared
 520 to millions users of our sample both in Tumblr and Flickr.

521 To estimate the presence of any skew in the local perception of deviant content
 522 consumption, we consider the nodes who are not producers and calculate the distri-
 523 bution of the proportion of their neighbors (in both the follow and reblog graphs)
 524 that either produce or reblog deviant material. The result is summarized in Figure 8.

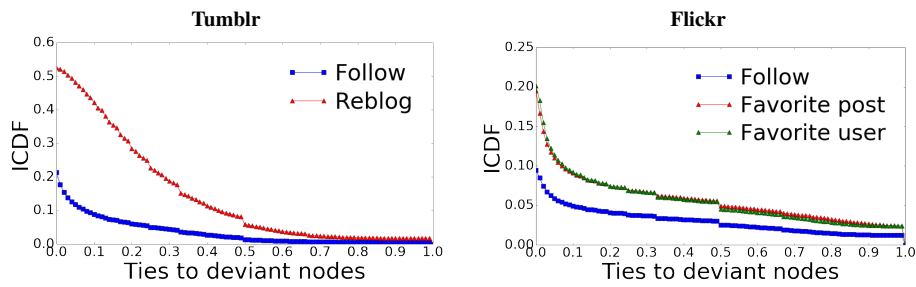


Fig. 8: Proportion of nodes with at least a given ratio of outlinks landing on deviant nodes (inverse cumulative density function) in Tumblr (left) and Flickr (right).

525 We observe that the follower network is nowhere close to exhibit the majority illu-
 526 sion phenomenon, with only the 10% of the population having 10% or more of their
 527 neighbors posting or reblogging deviant content in Tumblr and with only the 5% of
 528 the population having 10% or more of their neighbors posting or liking deviant con-
 529 tent in Flickr. The effect increases sensibly when considering the reblog network in
 530 Tumblr, with 40% of the population locally observing more than 10% of their contacts
 531 reblogging deviant content and almost 10% having more than half of their neighbors
 532 doing it. This happens partly because the size of the reblog network is one order of
 533 magnitude smaller than the one of the follower network, as we consider reblogging
 534 activity for one month only. In Flickr the majority illusion phenomenon is less promi-
 535 nent even though if we look at the favorite graph the effect is doubled compared to
 536 the follow graph; this means that when looking at recent activity only (reblogs or
 537 likes), local perception biases are much stronger (although not predominant) in the
 538 community than what can be inferred from the static follow graph.

539 Although strongly biased perceptions are not predominant when counting the
 540 number of neighbors, a stronger bias emerges when looking at the *volume* of de-
 541 viant content that is observed by a node from its neighbors in Tumblr. We calculated
 542 that more than 71% of nodes in Tumblr reblog less deviant content than the aver-
 543 age of their friends (considering friends who posted or reblogged at least once in the
 544 time frame we consider). This effect, that derives directly from the strong correla-
 545 tion between degree and number of posts and reblogs, suggests that the local users'
 546 perception of other people's behavior is skewed towards an image of pervasive con-
 547 sumption of deviant content in Tumblr that might be a driver to stimulate the diffusion
 548 of deviant content in this social network.

549 *Q5) How cascades of deviant content are characterized?*

550 We have found that the activity of a relatively small groups of producers can echo
 551 in the network and reach a very large audience. An open question remains about how
 552 an individual piece of deviant content spreads along the social ties in comparison to
 553 any other content type. To partly answer this question we focused on Tumblr, where
 554 a post can be reblogged by several users, possibly in long chains, thus generating
 555 *information cascades* [25, 49].

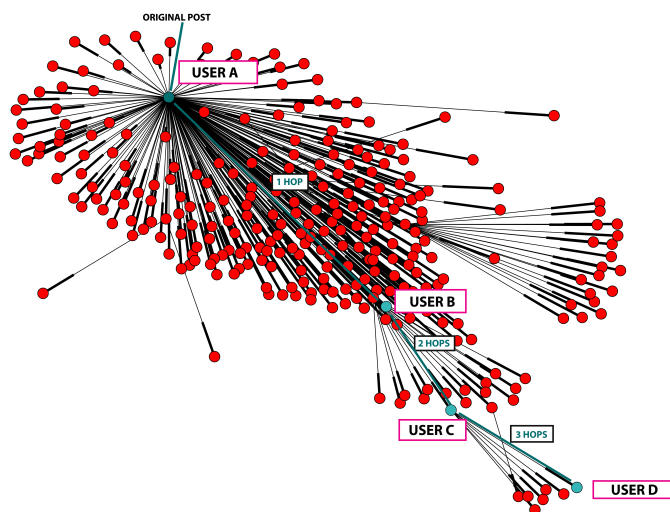


Fig. 9: Example of a reblog cascade generated by a post published by user A. In particular a chain is highlighted in light blue showing user D reblogging the post generated by user A through user C, who had reblogged the content from user B before. The chain is three-hops long.

556 Figure 9 show an example of a cascade generated by a post published by user A.
 557 All the nodes represent other users sharing that post directly from the original blog or
 558 through other users who previously had shared the content (e.g., user C from user B).
 559 The nodes without outgoing links are the leaves of the tree, final re-blogging actions
 560 whose content is not reshared anymore in the period of analysis. In the picture a chain
 561 of 3 hops is highlighted. We are interested in adult posts which extensively propagate,
 562 with long reblog chains (> 10 hops).

563 We selected 157K posts created in the first week of January 2016 by users in any
 564 of the *Producer* clusters and that are reblogged at least once. Figure 10 shows the
 565 cascade size and depth distribution for a sample of the selected adult posts and for
 566 non-adult posts. The distributions are in line with similar results obtained in previous
 567 studies related to cascades [16] with a weaker propagation effect for adult content
 568 cascades compared to regular type of content.

569 Most posts are reblogged just for a few days after they are published, even though
 570 there is a consistent tail of posts that gets reblogs for several weeks (Figure 11). We
 571 focus on the 529 posts that generate long cascades with 10 hops or more. We compare
 572 those with 657 non-adult posts (manually checked) with comparable virality (10 hops
 573 or more) published by users in the *Bridge* cluster.

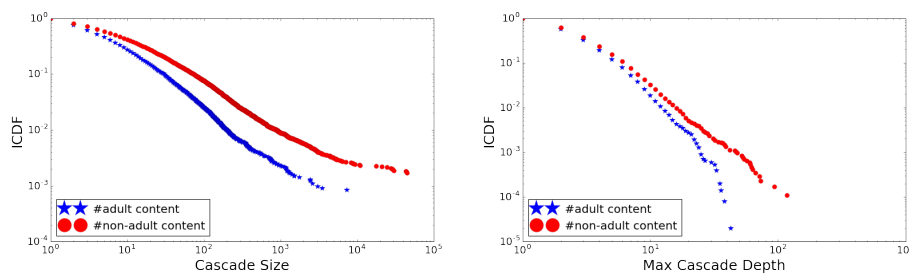


Fig. 10: The inverse cumulative distribution function (ICDF) of cascade size and depth for a sample of adult and non-adult posts.

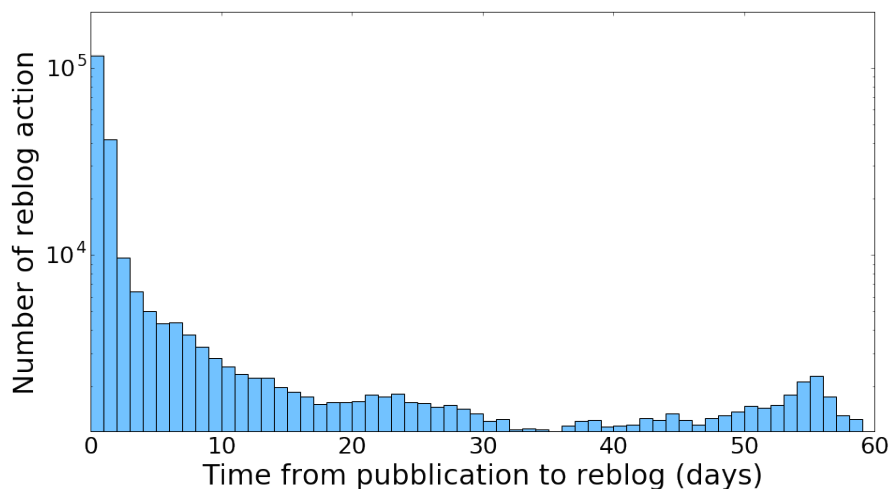


Fig. 11: Number of reblogs of adult posts created by Producers in the first week of January 2016. We considered all the reblog actions happening in January and February 2016 regarding the selected 157K posts created in the first week of January 2016. For each reblog action we plot the time in days from the publication date of the original root post to the reblog action time. The histogram shows the number of reblogs (y-axis) in relation to the amount of days passed by the original publication of the reblogged content.

574 We characterize each cascade by the maximum time of diffusion (time of the
 575 last reblog) and the number of leaves that the post generated (number of final users
 576 reblogging the content which is not reshared anymore from them). Figure 12 shows
 577 the relationship between the two dimensions, for maximum time of diffusion binned
 578 in 10-days long intervals.

579 The time of last reblog is an indicator of the persistence of a content item over
 580 time. We observe that the persistence in time is similar both for adult and non-adult
 581 content. We use the number of leaves to understand how the frontier of the cascade
 582 evolves over time. In this case, the frontier of adult content is much smaller than non-

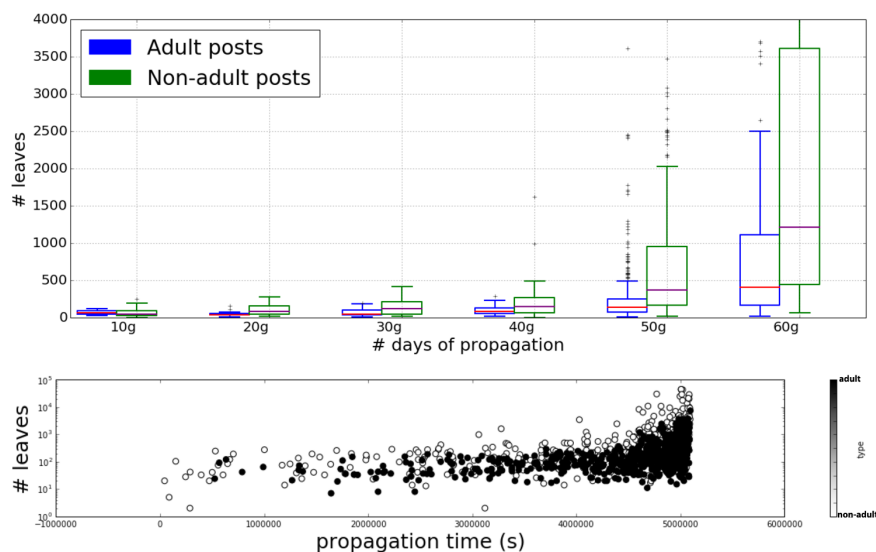


Fig. 12: Time length vs number of leaves for porn and not porn post cascades (> 10). Each boxplot (blue for adult posts and green for non-adult posts) shows how all the reblogs done in the temporal range (days from publication of the original post and the reblog action) are distributed (red line is the median and 50% of the data are contained in the box).

583 adult, meaning that adult content is less likely to be reblogged and that the cascade as
 584 a smaller width. This might indicate that adult content can become viral and it reaches
 585 users far in time and number of hops but the process involves less people than for non-
 586 adult content probably because of the embarrassment of sharing this kind of content
 587 which shows that drivers of the information are less but still they spread a lot the
 588 content. We verified statistically that the difference in number of leaves given a time
 589 diffusion for a post taken from the two samples (adult and non-adult) is significant
 590 through the Welch Two Sample t-test done on the ration leaves over time. The test is
 591 significant with a low p-value ($p = 867e - 07$).

592 *Q6) Is it possible to reduce the diffusion of deviant content in Tumblr with targeted*
 593 *interventions?*

594 Previous literature that investigated the properties of small-world networks indi-
 595 cates that information spreading or other phenomena of contagious nature can be
 596 drastically reduced by acting on a limited number of nodes in the graph [56]. Effec-
 597 tiveness of targeted interventions has been shown in a variety of domains, epidemics
 598 being the most prominent among them.

599 The intuition informed by previous work suggests that the wide diffusion of de-
 600 viant content can be reduced by properly marking the posts produced by a small set of
 601 core nodes and showing them only to people who explicitly declared their interest for
 602 that specific topic. In a simplified experimental scenario, we measure the proportion
 603 of active consumers reached by adult content in a setting where all the posts from a

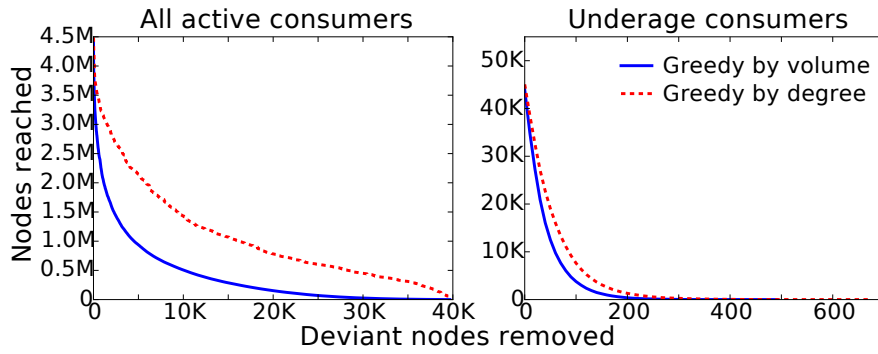


Fig. 13: Shrinkage of content diffusion after deviant nodes removal, using two different strategies.

604 set of core nodes C are erased. The question is how to select C and how big it needs
 605 to be to uproot the diffusion process.

606 The optimal selection of nodes is a set cover problem (NP-complete), but we test
 607 two common approximated strategies to solve it: *i) greedy by volume*, an algorithm
 608 that ranks nodes by the number of blogs that are reached by the content they pro-
 609 duce; and *ii) greedy by degree*, that takes into account the network structure only
 610 and ranks nodes by their in-degree in the reblog network. The effectiveness of the
 611 two approaches as $|C|$ increases is shown in Figure 13. Although using the indegree
 612 as proxy for the diffusion potential is not optimal, the removal of the 5,000 high-
 613 est indegree nodes curbs the diffusion by more than 50%. As expected, the strategy
 614 by volume is more effective (as it better approximates the optimal set cover), with a
 615 surprisingly sharp decay of the deviant content reach. The removal of the 5,000 top
 616 nodes reduces the information spreading by nearly 80%, which increases to almost
 617 100% when extending the block to 25,000 nodes. Furthermore, using our sample of
 618 demographic information, we find that to limit the exposure of underage users would
 619 be sufficient to remove the 200 top nodes, as identified by any of the two selection
 620 strategies. To monitor and control the capabilities a vertex may have on data flowing
 621 in the network other alternative approaches based on Routing Betweenness Central-
 622 ity (RBC) have been proposed [23], but we limited to simple strategies since the
 623 intervention strategies to limit the flow of this content are not the main scope of this
 624 work.

625 We used Tumblr as a case study to show the effect of targeted interventions be-
 626 cause in Flickr the spreading of the content is limited by the platform and the access
 627 of adult content is direct and not through cascades.

628 4.3 Demographics factors

629 The demographic composition of online adult content consumers has been measured
 630 by several sociological surveys (see Section 2), but none of them partitions the partic-
 631 ipants according to their type of consumption. Yet, we have shown that the categories

Table 5: Gender distribution among different categories in Tumblr and in Flickr.

	Group	Male	Female
Tumblr	All	29%	71%
	Producers	84%	16%
	Bridges	47%	53%
	Active consumers (reblog)	31%	69%
	Passive consumers (follow)	34%	66%
	Unintentionally exposed	20%	80%
Flickr	All	59%	41%
	Producers	76%	24%
	Bridges	72%	28%
	Active consumers (like)	87%	13%
	Passive consumers (follow)	68%	32%
	Unintentionally exposed	62%	38%

632 of people exposed to online deviant content range from the active content producers
633 to unintentional consumers. This calls for an investigation of the relationship between
634 type of consumption and demographic characterization.

635 *Q7) Is there a significant difference in the distribution of age and gender between*
636 *members of the deviant network and people with different levels of exposure to deviant*
637 *content?*

638 We report the distribution of age and gender of users with different levels of exposure
639 to adult content, computed on the sample of 1.7M Tumblr users and 12.3M Flickr
640 users who self-reported their demographic information. The two social networks are
641 very different in the age distribution of the users in addition to typology of content.
642 Tumblr is a social network targeted on young people. The average age in the sample
643 is slightly higher than 26, and female are the majority (71%). Flickr is more used by
644 adult people and professional photographers with an age on average of 41 and it is
645 more balanced in the gender of the users with 59% of male users (Table 5).

646 To partly validate the user-provided information, we first compare them with
647 third-party statistics. Our numbers are roughly compliant with several public reports
648 that rely on orthogonal methods for assessing the age and gender of users (e.g.,
649 surveys and clickstream monitoring [44, 58]). Also, we further validate the gender
650 data by assessing that the 95% of users in Tumblr and 95% of users in Flickr in the
651 *Producer*₂ cluster focused on male homosexual content are indeed male. The overall
652 age distribution of age by gender is shown in Figure 14 for Tumblr and for Flickr:
653 in both OSNs male tend to be older with an higher age variance in Flickr. In Tumblr
654 moreover the percentage on underage users is high: around 10% while in Flickr is
655 0.2%, this confirm why is important to study strategies to prevent visibility of adult
656 content to minors as we have discussed in the previous research question especially
657 for OSN popular among young people. Despite the spikes corresponding to birthdays
658 in round decades (1970, 1980, and 1990), probably due to misreporting, the distribu-
659 tions still tend to be Gaussian, as expected.

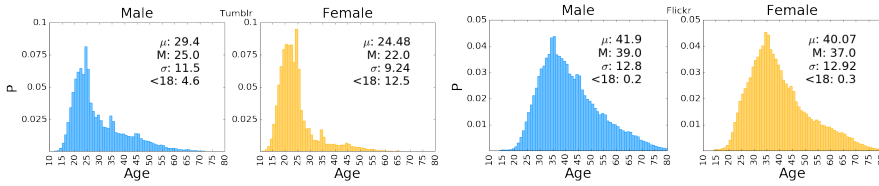


Fig. 14: Age distribution of users in our dataset for Tumblr (left) and Flickr (right). Mean μ , median M, standard deviation σ , and percentage of users under 18 years old are reported.

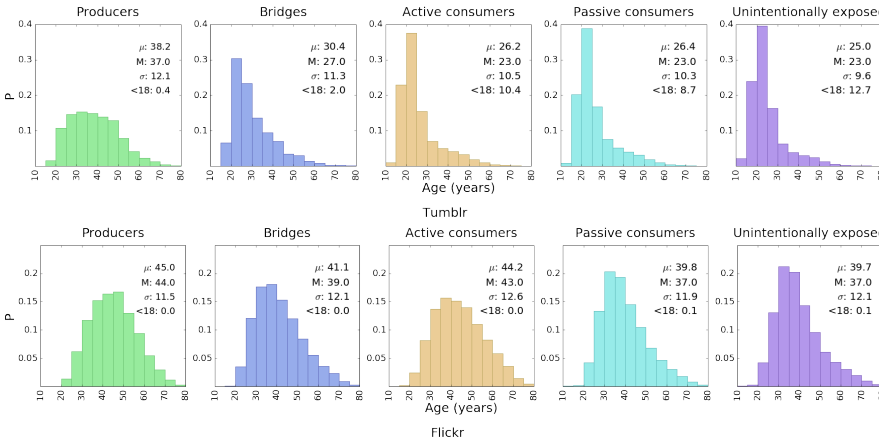


Fig. 15: Age distribution of different groups of producers and consumers of adult content in Tumblr (top) and in Flickr (bottom).

660 We then measure differences in age⁴ and gender distribution for the user classes
 661 of *producers*, *bridges*, *active consumers*, *passive consumers*, and *unintentionally ex-*
 662 *posed* users (Figure 15). In Flickr we have the same categories with a different char-
 663 acterization as reported in Figure 6: Likers (active consumers), Followers (passive
 664 consumers).

665 **Tumblr.** Figure 15 shows that producers are considerably older than the typical
 666 user, averaging around age 38 and with almost no underage users. Different from the
 667 overall distribution, they are mostly male (84%), in alignment with studies indicating
 668 that men are more involved in assiduous consumption of adult material (Table 5).
 669 Moreover we found that bridge groups are fairly gender-balanced (with more female
 670 –53%– in the celebrity-oriented community) and include younger people (30 years
 671 old on average). Consumers of deviant nodes who actively reblog or passively follow
 672 deviant blogs are covered by demographic data at 12%, proportion that drops to 4%
 673 among those who follow deviant nodes. In both classes, the age is quite representative
 674 of the overall Tumblr population in our sample (more than 66% female). A similar

⁴ The number of samples in each age distribution is high; therefore, as expected, all the differences between the average values are statistically significant ($p < 0.01$) under the Mann-Whitney test.

675 male-female proportion holds for people that are potentially exposed to deviant con-
676 tent in an unintentional way. This last class has the highest proportion of underage
677 people (13%), which reinforces the concern about young teens unwillingly seeing
678 inappropriate content.

679 **Flickr.** Figure 15 shows that producers are older than other categories (4 years
680 more than average) and unintentional exposed users are the youngest as it is in Tum-
681 blr. Moreover we found that producers in Flickr are for the 76% male, confirming a
682 higher presence of men in production of adult content. Among consumers we have
683 differences in the average age in we consider people who like adult content (avg. age
684 44) and users who follow adult content producers (avg. age 40). This shows a differ-
685 ent behavior in consumers: younger people tend to follow more than liking and male
686 users tend to like more than female. Indeed, if we look at gender we surprisingly see
687 that among likers 87% are male, while among followers only 68% (Table 5). The
688 bridge community has a lower percentage of male users compared to the produc-
689 ers (72%) probably showing indicating a more tendency to consume hard content by
690 male and soft content by female.

691 The fact that the gender distribution for consumers deviates only slightly from
692 the overall gender distribution is in partial disagreement with previous studies on
693 gender and sexual behavior [36, 43] which state that men are usually more exposed
694 than women to adult material. This is particular evident in Tumblr and we conjecture
695 that this might happen because of the tendency of female to have their peak of adult
696 content consumption in a much younger age than men (as shown by [30]), combined
697 with the predominance of young female among Tumblr users. To verify better the
698 relationship between age and gender in consumption of adult content in both Flickr
699 and Tumblr we aim to answer one more question.

700 *Q8) Does age have an effect on how different genders consume adult content?*

701 To find out, we measure the proportion of male and female actively exposed to deviant
702 content (by following deviant account) broken down by age (Figure 16).

703 Surprisingly both in Tumblr and in Flickr, despite the differences in the demo-
704 graphic of the users and the typology of social platform, we observe a similar trend.
705 The curve for men shows an increasing trend that plateaus at its maximum in the
706 range of age 40 to 55. In contrast, women, although less exposed than men at any
707 age, have their peak in their 20s, much earlier than men. This peak is longer in Flickr
708 up to 25 probably because the platform target a more adult audience. This observa-
709 tion supports previous findings [30] and explains the distributions we observed. In
710 absolute terms the volume of users consuming adult content (by following) in rela-
711 tion to the gender with different ages is interesting to compute and this is one of the
712 first large-scale study shedding light on that. In Figure 17 we report the ratio male
713 users over female users for different age ranges for adult content consumers both in
714 Tumblr and in Flickr. If we compare those values with the general ratio of the popula-
715 tion of the whole online social network we discover a similar trend. The consumption
716 of adult content is substantially equal between the two genders under 25 years old.
717 After that age the percentage of male users increases progressively compared to the
718 female users. In particular in Flickr which is a platform targeting more adult people,
719 the male consumption keep increasing until the age of 50.

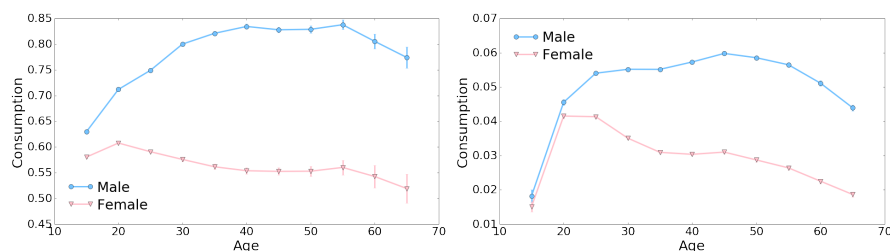


Fig. 16: Ratio of male and female consuming adult content for different age bands in Tumblr (left) and in Flickr (right).

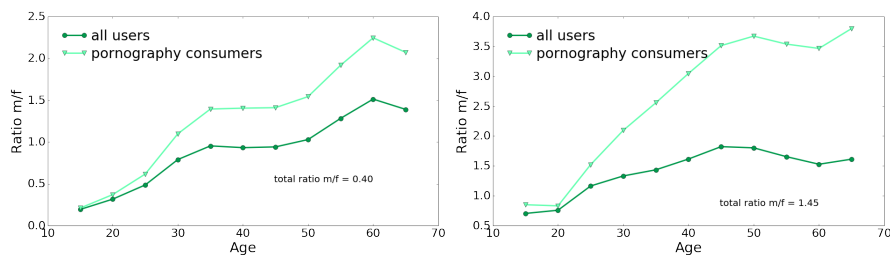


Fig. 17: Ratio of male and female users consuming adult content for different age bands compared to the ratio of the general population of the OSN in the sample for Tumblr (left) and Flickr (right).

720 Q9) What is the cross-platform behavior of deviant content consumers?

721 To better understand the dynamics of diffusion of deviant content, one could con-
 722 sider how adult material leaks from an online platform to another. It is very hard to
 723 track individual pieces of information moving across different social media layers.
 724 For example, a Flickr photo can be re-uploaded to Tumblr, which makes it difficult
 725 to trace it back to its original Flickr URL. To investigate how much the activity of
 726 people involved in the activity of production and consumption of adult content is
 727 cross-platform, we resort to an analysis at user-level instead. We do so by consider-
 728 ing those users who have subscribed to both Flickr and Tumblr using the same email
 729 address. We found 293K users with an account in both platform in our sample. As
 730 shown in Figure 18, approximately 32% of them are involved in pornography at least
 731 in one of two networks, the majority using Flickr instead of Tumblr.

732 There are two competing hypothesis regarding the involvement of users on mul-
 733 tiple social networks. On one hand, the literature in multiplex networks seems to
 734 point towards a scenario in which the activity indicators of nodes on multiple layers
 735 are strongly correlated [55]. This finding holds mostly for structural indicators (e.g.,
 736 node degree), but if the same principle applied also to the *type* of activity, people
 737 who actively engage in activities connected to deviant content on one platform would
 738 likely do the same also on the other one. On the other hand, social science studies sug-
 739 gest that the role people plays is strongly dependent from the context they are acting
 740 in. Social identification depends only partially by the characteristics of the individual

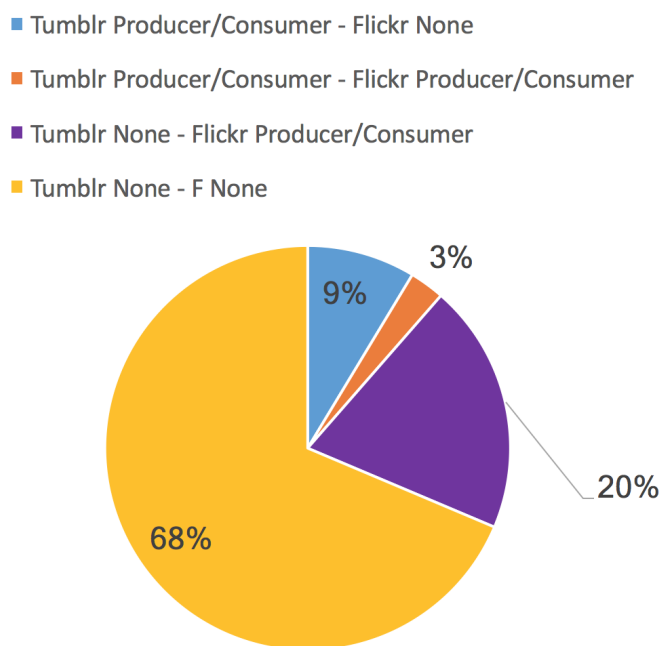


Fig. 18: Cross-platform behavior: matching users in Tumblr and in Flickr. Users are characterized by their role according to our classification (adult content producers or consumers. People not related to pornography even though they might be unintentionally exposed are classified as None).

741 and it is determined largely by factors that are distinctive of the group in which the
 742 individual is part [6, 68]. This holds when comparing the behavior of an individual
 743 across groups and, consequently, across social systems. The fact that both Tumblr and
 744 Flickr do not enforce strong user identification (pseudonyms are allowed), makes it
 745 even easier for a person to build different, possibly contrasting online personas. This
 746 leads us to hypothesize that deviant users on one platform do not have the same type
 747 of behavior, on average, on the other.

748 In agreement with this second hypothesis, we find that only 3% (see Figure 18)
 749 of the users with matching accounts consume or produce deviant content in both plat-
 750 forms: people tend to use different social media for different purposes. Moreover,
 751 even when people use both platforms to deal with porno-graphical content their roles
 752 are often different (producing or consuming). We further investigated the users pro-
 753 ducing deviant content and we found that less than 1% are producers in both net-
 754 works.

755 5 Conclusion

756 This work aims at motivating researchers who study all types of deviant communities
757 online as well as offline to explore in more depth the interaction between the agents in
758 such networks and the external social environment. Our contribution scratches only
759 the surface of the exploration space that underlies the many types of deviant networks
760 and the multitude of settings they are situated within. The study we have presented
761 is limited under many aspects, beginning from the focus on a single type of deviant
762 behavior –adult material consumption– that is much more pervasive than others (e.g.,
763 anorexia) and, in that, has unique characteristics that likely cannot generalize to other
764 deviant groups. In terms of methodology, alternative techniques (e.g., computer vi-
765 sion) could be used to identify adult content without a dedicated dictionary; those
766 could possibly lead to describe the same phenomenon from a slightly different angle,
767 for instance considering more exhaustively nodes that are not reached by search traf-
768 fic or by tags. Furthermore the notion of pornographic content is culture-dependent;
769 then it would be possible to study deviant behaviors under different cultural premises
770 and not only under a western country perspective that we adopted.

771 Our work has limitations that leave space for exploration in future work. We have
772 focused our study on Tumblr and Flickr, but more social media ecosystems could be
773 included; Twitter would be a good candidate as it does not enforce any strong restric-
774 tion on the content of tweets. We have studied the spread of pornographic content,
775 but a broader exploration of how other types of deviant content diffuses on general-
776 purpose social media would be very desirable. For what concerns the analysis tools
777 we have used to measure the importance that nodes have in spreading information,
778 one could investigate a plethora of alternative ways to assess the role of groups and
779 individuals in the process. In particular, we believe an interesting extension to our
780 work would be to look at more sophisticated centrality metrics such as routing be-
781 tweenness centrality [24].

782 Despite these limitations, we believe that our study has already important theo-
783 retical implications in revealing, for the first time on very large scale, that deviant
784 communities can be deeply rooted into the relational fabric of a social network, and
785 that the echo of their abnormal activity can reach a plenitude of ordinary users. Also,
786 from a practical point of view, learning the effect that a minority group can have on
787 a much larger audience is key to trigger mechanisms able to contain risky deviant
788 phenomena by means of targeted interventions on few nodes, as we have shown. We
789 believe that this work could set the basis for a line of study that could lead to a deeper
790 understanding of deviant networks and of their impact on everyone’s life.

791 References

- 792 1. Adamic, L. A. and Glance, N. (2005). The political blogosphere and the 2004 us
793 election: divided they blog. In *International workshop on Link discovery*. ACM.
- 794 2. Aiello, L. M. (2015). Group types in social media. In Paliouras, G., Papadopou-
795 los, S., Vogiatzis, D., and Kompatsiaris, Y., editors, *User Community Discovery*,

- 796 Human–Computer Interaction Series, pages 97–134. Springer International Pub-
797 lishing.
- 798 3. Aiello, L. M., Barrat, A., Cattuto, C., Ruffo, G., and Schifanella, R. (2010). Link
799 creation and profile alignment in the aNobii social network. In *SocialCom*.
- 800 4. Aiello, L. M., Barrat, A., Schifanella, R., Cattuto, C., Markines, B., and Menczer,
801 F. (2012). Friendship prediction and homophily in social media. *ACM Transactions*
802 *on the Web*, 6(2):9:1–9:33.
- 803 5. Allen, M., D’Alessio, D., and Brezgel, K. (1995). A meta-analysis summarizing
804 the effects of pornography ii aggression after exposure. *Human communication*
805 *research*, 22(2):258–283.
- 806 6. Ashforth, B. E. and Mael, F. (1989). Social identity theory and the organization.
807 *Academy of management review*, 14(1):20–39.
- 808 7. Attwood, F. (2005). What do people do with porn? qualitative research into the
809 consumption, use, and experience of pornography and other sexually explicit me-
810 dia. *Sexuality and culture*, 9(2).
- 811 8. Barbieri, N., Bonchi, F., and Manco, G. (2013). Cascade-based community detec-
812 tion. In *WSDM*. ACM.
- 813 9. Bessi, A., Coletto, M., Davidescu, G. A., Scala, A., Caldarelli, G., and Quattro-
814 ciocchi, W. (2015). Science vs conspiracy: collective narratives in the age of mis-
815 information. *PLoS one*, 10(2):02.
- 816 10. Blackburn, J., Simha, R., Kourtellis, N., Zuo, X., Ripeanu, M., Skvoretz, J., and
817 Iamnitchi, A. (2012). Branded with a scarlet c: Cheaters in a gaming social network.
818 In *WWW: Proceedings of the 21st International Conference on World Wide Web*,
819 pages 81–90, New York, NY, USA. ACM.
- 820 11. Blondel, V. D., Guillaume, J.-L., Lambiotte, R., and Lefebvre, E. (2008). Fast
821 unfolding of communities in large networks. *Journal of Statistical Mechanics:*
822 *Theory and Experiment*, 2008(10):P10008.
- 823 12. Boero, N. and Pascoe, C. J. (2012). Pro-anorexia communities and online inter-
824 action: Bringing the pro-ana body online. *Body & Society*, 18(2):27–57.
- 825 13. Buzzell, T. (2005). Demographic characteristics of persons using pornography
826 in three technological contexts. *Sexuality & Culture*, 9(1):28–48.
- 827 14. Cha, M., Haddadi, H., Benevenuto, F., and Gummadi, P. K. (2010). Measuring
828 user influence in twitter: The million follower fallacy. *ICWSM*, 10(10-17):30.
- 829 15. Chen, A.-S., Leung, M., Chen, C.-H., and Yang, S. C. (2013). Exposure to inter-
830 net pornography among taiwanese adolescents. *Social Behavior and Personality:*
831 *an international journal*, 41(1):157–164.
- 832 16. Cheng, J., Adamic, L., Dow, P. A., Kleinberg, J. M., and Leskovec, J. (2014).
833 Can cascades be predicted? In *Proceedings of the 23rd international conference on*
834 *World wide web*, pages 925–936. ACM.
- 835 17. Christakis, N. A. and Fowler, J. H. (2008). The collective dynamics of smoking
836 in a large social network. *New England journal of medicine*, 358(21):2249–2258.
- 837 18. Clinard, M. and Meier, R. (2015). *Sociology of deviant behavior*. Wadsworth
838 Cengage Learning.
- 839 19. Coletto, M., Aiello, L. M., Lucchese, C., and Silvestri, F. (2016). On the be-
840 haviour of deviant communities in online social networks. In *ICWSM 2016, Köln,*
841 *Germany*.

- 842 20. Conover, M., Ratkiewicz, J., Francisco, M., Gonçalves, B., Menczer, F., and
843 Flammini, A. (2011). Political polarization on twitter. In *ICWSM*.
- 844 21. Davis, J. P. (2002). The experience of ‘bad’ behavior in online social spaces: A
845 survey of online users. *Social Computing Group, Microsoft Research*.
- 846 22. De Choudhury, M. (2015). Anorexia on tumblr: A characterization study. In
847 *Digital Health*. ACM.
- 848 23. Dolev, S., Elovici, Y., and Puzis, R. (2010a). Routing betweenness centrality.
849 *Journal of the ACM (JACM)*, 57(4):25.
- 850 24. Dolev, S., Elovici, Y., and Puzis, R. (2010b). Routing betweenness centrality.
851 *Journal of ACM*, 57(4).
- 852 25. Dow, P. A., Adamic, L. A., and Friggeri, A. (2013). The anatomy of large face-
853 book cascades. In *ICWSM*.
- 854 26. Dunlop, P. D. and Lee, K. (2004). Workplace deviance, organizational citizen-
855 ship behavior, and business unit performance: The bad apples do spoil the whole
856 barrel. *Journal of organizational behavior*, 25(1):67–80.
- 857 27. Easley, D. and Kleinberg, J. (2010). *Networks, Crowds, and Markets: Reasoning*
858 *About a Highly Connected World*. Cambridge University Press.
- 859 28. Feld, S. L. (1991). Why your friends have more friends than you do. *American*
860 *Journal of Sociology*, pages 1464–1477.
- 861 29. Feller, A., Kuhnert, M., Sprenger, T. O., and Welp, I. M. (2011). Divided they
862 tweet: The network structure of political microbloggers and discussion topics. In
863 *ICWSM*.
- 864 30. Ferree, M. (2003). Women and the web: Cybersex activity and implications.
865 *Sexual and Relationship Therapy*, 18(3):385–393.
- 866 31. Fortunato, S. (2010). Community detection in graphs. *Physics Reports*,
867 486(3):75–174.
- 868 32. Gavin, J., Rodham, K., and Poyer, H. (2008). The presentation of “pro-anorexia”
869 in online group interactions. *Qualitative Health Research*, 18(3):325–333.
- 870 33. Grabowicz, P. A., Aiello, L. M., Eguiluz, V. M., and Jaimes, A. (2013). Distin-
871 guishing topical and social groups based on common identity and bond theory. In
872 *WSDM*. ACM.
- 873 34. Guerra, P. H. C., Meira Jr, W., Cardie, C., and Kleinberg, R. (2013). A measure of
874 polarization on social media networks based on community boundaries. In *ICWSM*.
- 875 35. Haas, S. M., Irr, M. E., Jennings, N. A., and Wagner, L. M. (2010). Online
876 negative enabling support groups. *New Media & Society*.
- 877 36. Hald, G. M. (2006). Gender differences in pornography consumption among
878 young heterosexual danish adults. *Archives of sexual behavior*, 35(5):577–585.
- 879 37. Hald, G. M., Malamuth, N. N., and Lange, T. (2013). Pornography and sexist
880 attitudes among heterosexuals. *Journal of Communication*, 63(4):638–660.
- 881 38. Hald, G. M. and Štulhofer, A. (2015). What types of pornography do people use
882 and do they cluster? assessing types and categories of pornography consumption in
883 a large-scale online sample. *The Journal of Sex Research*, pages 1–11.
- 884 39. Hodas, N. O., Kooti, F., and Lerman, K. (2013). Friendship paradox redux: Your
885 friends are more interesting than you. In *ICWSM*.
- 886 40. Kayes, I., Kourtellis, N., Quercia, D., Iamnitchi, A., and Bonchi, F. (2015). The
887 social world of content abusers in community question answering. In *WWW: Pro-*

- 888 *ceedings of the 24th International Conference on World Wide Web*, pages 570–580,
889 New York, NY, USA. ACM.
- 890 41. Kitsak, M., Gallos, L. K., Havlin, S., Liljeros, F., Muchnik, L., Stanley, H. E., and
891 Makse, H. A. (2010). Identification of influential spreaders in complex networks.
892 *Nature physics*, 6(11):888–893.
- 893 42. Kühn, S. and Gallinat, J. (2014). Brain structure and functional connectivity
894 associated with pornography consumption: the brain on porn. *JAMA psychiatry*,
895 71(7):827–834.
- 896 43. Kvaem, I. L., Træen, B., Lewin, B., and Štulhofer, A. (2014). Self-perceived ef-
897 fects of internet pornography use, genital appearance satisfaction, and sexual self-
898 esteem among young scandinavian adults. *Cyberpsychology: Journal of Psychoso-*
899 *cial Research on Cyberspace*, 8(4).
- 900 44. LaSala, R. (2012). The social makeup of social media. *Compete.com Tech Blog*.
- 901 45. Lee, L.-H. and Chen, H.-H. (2011). Collaborative cyberporn filtering with col-
902 lective intelligence. In *SIGIR*. ACM.
- 903 46. Lerman, K., Yan, X., and Wu, X.-Z. (2016). The “majority illusion” in social
904 networks. *PLoS ONE*, 11(2):1–13.
- 905 47. Leskovec, J. and Horvitz, E. (2008). Planetary-scale views on a large instant-
906 messaging network. In *Proceedings of the 17th international conference on World*
907 *Wide Web*, pages 915–924. ACM.
- 908 48. Leskovec, J., Kleinberg, J., and Faloutsos, C. (2005). Graphs over time: densi-
909 fication laws, shrinking diameters and possible explanations. In *ACM SIGKDD*.
910 ACM.
- 911 49. Leskovec, J., McGlohon, M., Faloutsos, C., Glance, N. S., and Hurst, M. (2007).
912 Patterns of cascading behavior in large blog graphs. In *SDM*, volume 7, pages
913 551–556. SIAM.
- 914 50. Martin-Borregon, D., Aiello, L. M., Grabowicz, P., Jaimes, A., and Baeza-Yates,
915 R. (2014). Characterization of online groups along space, time, and social dimen-
916 sions. *EPJ Data Science*, 3(1):8.
- 917 51. Mitchell, K. J., Finkelhor, D., and Wolak, J. (2003). The exposure of youth
918 to unwanted sexual material on the internet a national survey of risk, impact, and
919 prevention. *Youth & Society*, 34(3):330–358.
- 920 52. Morgan, E. M., Snelson, C., and Elison-Bowers, P. (2010). Image and video dis-
921 closure of substance use on social media websites. *Computers in Human Behavior*,
922 26(6):1405–1411.
- 923 53. Negoescu, R. A. and Gatica-Perez, D. (2008). Analyzing flickr groups. In *CIVR*,
924 New York, NY, USA. ACM.
- 925 54. Newman, M. E. (2006). Modularity and community structure in networks. *Pro-*
926 *ceedings of the national academy of sciences*, 103(23):8577–8582.
- 927 55. Nicosia, V. and Latora, V. (2015). Measuring and modeling correlations in mul-
928 tiplex networks. *Physical Review E*, 92(3):032805.
- 929 56. Pastor-Satorras, R. and Vespignani, A. (2005). Epidemics and immunization in
930 scale-free networks. *Handbook of graphs and networks: from the genome to the*
931 *internet*, pages 111–130.
- 932 57. Phillips, D. J. (1996). Defending the boundaries: Identifying and countering
933 threats in a usenet newsgroup. *The information society*, 12(1):39–62.

- 934 58. Pingdom (2012). Report: Social network demographics in 2012. *Pingdom.com*
935 *Tech Blog*.
- 936 59. Ramos, J. d. S., Pereira Neto, A. d. F., and Bagrichevsky, M. (2011). Pro-anorexia
937 cultural identity: characteristics of a lifestyle in a virtual community. *Interface*
938 (*Botucatu*), 15(37):447–460.
- 939 60. Ratkiewicz, J., Conover, M., Meiss, M., Gonçalves, B., Patil, S., Flammini, A.,
940 and Menczer, F. (2011). Detecting and tracking the spread of astroturf memes in
941 microblog streams. In *WWW*.
- 942 61. Romero, D. M., Tan, C., and Ugander, J. (2013). On the interplay between social
943 and topical structure. In *ICWSM*.
- 944 62. Romito, P. and Beltramini, L. (2015). Factors associated with exposure to violent
945 or degrading pornography among high school students. *The Journal of School*
946 *Nursing*, 31(4):280–90.
- 947 63. Sabina, C., Wolak, J., and Finkelhor, D. (2008). The nature and dynamics of
948 internet pornography exposure for youth. *CyberPsychology & Behavior*, 11(6).
- 949 64. Schifanella, R., Barrat, A., Cattuto, C., Markines, B., and Menczer, F. (2010).
950 Folks in folksonomies: Social link prediction from shared metadata. In *WSDM*.
951 ACM.
- 952 65. Schuhmacher, M., Zirn, C., and Völker, J. (2013). Exploring youporn categories,
953 tags, and nicknames for pleasant recommendations. In *Workshop on Search and*
954 *Exploration of X-Rated Information*. ACM.
- 955 66. Shores, K. B., He, Y., Swanenburg, K. L., Kraut, R., and Riedl, J. (2014). The
956 identification of deviance and its impact on retention in a multiplayer game. In
957 *CSCW: Proceedings of the 17th ACM Conference on Computer Supported Cooper-*
958 *ative Work & Social Computing*, pages 1356–1365, New York, NY, USA. ACM.
- 959 67. Suler, J. R. and Phillips, W. L. (1998). The bad boys of cyberspace: Deviant be-
960 havior in a multimedia chat community. *CyberPsychology & Behavior*, 1(3):275–
961 294.
- 962 68. Thibaut, J. W. and Kelley, H. H. (1959). The social psychology of groups.
963 *Database: PsycINFO*.
- 964 69. Tyson, G., Elkhatib, Y., Sastry, N., and Uhlig, S. (2013). Demystifying porn 2.0:
965 A look into a major adult video streaming website. In *IMC*. ACM.
- 966 70. Tyson, G., Elkhatib, Y., Sastry, N., and Uhlig, S. (2015). Are People Really
967 Social in Porn 2.0? In *ICWSM*.
- 968 71. Wasserman, S. and Faust, K. (1994). *Social network analysis: Methods and*
969 *applications*, volume 8. Cambridge university press.
- 970 72. Wellen, J. M. and Neale, M. (2006). Deviance, self-typicality, and group cohe-
971 sion the corrosive effects of the bad apples on the barrel. *Small Group Research*,
972 37(2):165–186.
- 973 73. Wolak, J., Mitchell, K., and Finkelhor, D. (2007). Unwanted and wanted expo-
974 sure to online pornography in a national sample of youth internet users. *Pediatrics*,
975 119(2):247–257.
- 976 74. Xu, J. and Chen, H. (2008). The topology of dark networks. *Communications of*
977 *the ACM*, 51(10):58–65.
- 978 75. Ybarra, M. L. and Mitchell, K. J. (2005). Exposure to internet pornography
979 among children and adolescents: A national survey. *CyberPsychology & Behavior*,

980 8(5):473–486.