# Identifying and Classifying Harassment in Arabic Wikipedia: A "Netnography"

Prepared for the Wikimedia Foundation

by

Michael Raish

Feb 2019

# Executive summary

This study investigated the extent and scope of harassment, trolling, insults, and other threatening communications within the context of Arabic Wikipedia using both observational and participatory data collection. This study found that:

- Although threatening and harassing communications can be frequently observed within this community, the community mechanisms intended to mitigate the effects of these behaviors, such as the *block request* mechanism, are robust and see the participation of numerous individuals.
- Many of the negative behaviors identified in this study occurred in disputes between established and long-standing community members. Participants in these disputes use insults to undermine their opponents' legitimacy within the Arabic Wikipedia community.
- Dispute participants and bystanders frame conflicts using the language of human rights or in terms of neutrality and bias, indicating the importance placed on these values in this community.
- There is an observable vein of conflict in the Arabic Wikipedia community related to administrator–editor relationships.
- Bystander intervention is repeatedly cited by interview participants as a crucial factor in mitigating the negative outcomes experienced by victims of abuse.
- Transparently conducted research efforts are an important way in which WMF can signal engagement and interest to non-English Wikipedia communities.
- Further study is needed to explore the experiences of new editors and to poll the views and experiences of administrators in a more structured way.

# 1. Community profile

Arabic Wikipedia (ARWP) is a growing, moderately sized community composed of 27 administrators and 4,760 active users. This community can be conceived of as a series of embedded circles in which the center represents a core group of active administrators, surrounded by a larger group of active editors, who are in turn surrounded by a much larger and more diffuse group of occasional editors. Within this matrix, however, there are a number of sub-groupings that span multiple levels of membership, and there are likewise a number of observable ongoing conflicts within the community that some participants have sought to describe as conflicts between editors and the moderation team. Although many active community members know each other in offline contexts and participate in regular meetups, many others only know each other virtually. Regardless, active community members maintain a number of non-public methods of communication, including the Wikimessage system, private email, participation in private and public Facebook groups, etc.

The ARWP community has several observable fault lines related to issues of community importance such as elections and awarding of administrator privileges. Although a core group of dedicated administrators has been managing the mechanisms of community for several years, most of these individuals are male[1], and recent years have likewise seen the growing participation of a new generation of younger editors. These new community members, who are often college-age or younger, bring with them new norms of communication and politeness in online contexts. In one instance, for example, a administrator admonished an editor that "what you learned on the forum groups doesn't apply here," i.e., attempting to reinforce the fact that Arabic Wikipedia has higher standards for expression and interaction than the other online spaces this user was presumably active in. While new users may be inconsistently capable or desirous of communicating in ways that conform with long-established norms in the Arabic Wikipedia community, the older generation may in turn be losing touch with the ways in which younger community members are accustomed to communicating online.

---

[1] This observation is based both on a dearth of Arabic administrator and user accounts that self-identify as female, as well as on input from knowledgeable WMF staff and knowledge from CE insights.

# 2. Methodology

This study employed a two-pronged approach to data collection consisting of (a) observational identification and classification of problematic and threatening communications on Arabic Wikipedia, and (b) asynchronous (email) interviews with knowledgeable dispute participants in order to explore their views and experiences.

## 2.1 Establishment of extent of threatening communications

In order to establish the extent and scope of the behaviors under investigation in Arabic Wikipedia, the block request page was initially selected as a jumping-off point for analysis. The block request page represents a centralized location in which problematic communicative acts, including the behaviors under investigation in this study, have been identified by community members as worthy of some form of sanction. Block requests within the two-year period leading up to the time of data collection (December, 2018) were grouped into one of nine categories based on the **first justification presented** by the reporting user in each request (many requests contain reports of multiple types of problematic behavior).

Second, individual instances of threatening communications were collected, translated, and categorized in the context of an annotated list of disputed interactions. Communications were primarily identified via having been flagged by community members in the block request page, although other tokens emerged via keyword searches for terms such as "harassment," "bullying," "insult," etc., as well as examination of community discussion forums, personal talk pages, and the talk pages of controversial articles. It should be noted that keyword searches are limited in their ability to identify certain practices. Multiple Arabic terms may be used to describe phenomena that might be described with a single word in English—in this study, searches for "bullying" تنمر tanammur were inconclusive, and searches for the "troll" ترول revealed two only results.

## 2.2 Email interview recruitment and administration

Following ongoing analysis of the *block request page*, a constellation of editors and administrators who are actively involved in reporting and mediating disputes emerged. A list

including members of both groups was compiled, as well as a WMF-approved recruitment script. Prospective participants were contacted via Wikipedia's private internal message system. Ultimately 18 private solicitation emails were sent, resulting in seven responses, of which six respondents agreed to participate. Three individuals eventually responded to posed questions, resulting in a 16.7% response rate.[2] See Appendix A for lessons learned and suggestions for future email-based interview studies.

# 3. Findings

## 3.1 Analysis of block requests

Initially 782 relevant block requests were identified within the two years leading up to the time of data collection. Of these, 577 were sustained by administrators (78.3% block rate). These requests were in-turn categorized into one of nine categories, based on the first violation listed in each request. The most common violation was *inappropriate username* (315 requests, 87.6% sustained), followed by *vandalism* (147 requests, 74.6% sustained), *non-encyclopedic* content (86 requests, 53.4% sustained), and **personal attacks** (67 requests, 40.3% sustained).[3] See Appendix B for a presentation of all analyzed requests. This analysis reveals that (a) personal attacks, which include threats, insults, harassment, and other behaviors under investigation in this study, are indeed a widely reported phenomenon on Arabic Wikipedia, and that (b) the *personal attack* category is the **least likely category** to result in a block. That is, the *personal attack* category is the only one of the nine identified categories in which less than 50% of block requests result in a block being levied against the accused.

## 3.2 Categorization of threatening behaviors

The initial analysis of block requests revealed a number of interactions between Arabic Wikipedia users that contained instances of threatening communications such as insults,

---

[2] Given the difficulties inherent to unsolicited recruitment of participants for participation in online survey studies, this rate of participation is in line with expectations.

[3] The phrase "personal attack" was chosen for this category due to its common use in reports of violations (تهجم شخصي *tahajjum shakhṣiyy*)

harassment and threats. Keyword searches for certain terms, such as تحرش *taḥarrush* and مضايقة *muḍāyaqa* "harassment" turned up additional examples of these practices.

A.  **Insults** consist of pejorative or derogatory remarks directed at a victim's person, race, religion, or contributions to Wikipedia. As with other phenomena such as harassment and trolling, insults can be difficult to identify—while one party to a dispute may perceive having been insulted, the aggressor may just as likely insist aftering being reported that the comment in question was not an insult. Identification of insults thus requires interpretation of both audience perception and speaker intent.

    In this community, **Person-oriented insults** aim to undermine an opponent's legitimacy by attacking their person. This is achieved through an author's affectation of a dismissive or derogatory stance toward the victim, such as through accusation of racism, the use of pejorative labels such as "idiot," or "Nazi," etc.

*As you can see, **[he is] an impolite person**, with such an **underdeveloped and arrogant tone**; how can you discuss with him—he doesn't have any culture of disagreement, and he insults and offends whomever disagrees with him.*

– block request discussion

> Notably, individuals accused of insulting their interlocutors often reject the label "insult" being applied to their comments. This may be even more frequent in disputes between experienced editors who are adept at using Arabic Wikipedia's various reporting mechanisms to flag and ultimately sanction their opponents. In the case of an editor who was reported for rhetorically asking their opponent, "are you a Nazi?", the editor clarified in the subsequent discussion of the possible block that "Nazi is a philosophy like many ideas or orientations. Are you a secularist? Are you an Islamist? Are you a capitalist? Are you a communist? Do these questions appear to be insults?"

*…[Y]our head has only learned how to oppress others, and I don't blame you for that personally, rather **I blame your country** which didn't teach you freedom of responsible opinion.*

– personal talk page

**Content-oriented insults**, on the other hand, aim to undermine an opponent's legitimacy by casting their contributions to Wikipedia as outside the bounds of acceptable authorship. In these cases, the insulting party attempts to cast their opponents as having failed to correctly interpret or adhere to Wikipedia's content, style, and tone guidelines. Content-directed insults thus seek to cast the victim as and incomplete or deficient member of the community. Typically, disputes between experienced contributors see the concurrent use of several overlapping communicative strategies, such as in accusations that one's opponent is a "**fabricator**," i.e., a dishonest individual who seeks to vandalize or otherwise publish false material on Wikipedia, thereby placing themselves outside the bounds of community norms.

*I didn't use any weak sources—all of the sources that I brought are **correct and online**, while the source that you brought (and **you didn't even list the page number**) are offline and it's basically a book in which the name of the mountain isn't even mentioned (I looked it up).*

– block request discussion

B. **Threats and commands** aim to dictate or control the actions of others by implicitly implying or explicitly articulating a series of future consequences. These actions represent claims of power over the recipient of the threat. The threatening party seeks to undermine the victim's autonomy and agency within the community. Threats and commands are often embedded within or accompanied by politeness markers.

*In general you have nothing to do with this issue (you're not a moderator) **my dear**…rather, you're the one who made this request to stir up trouble between me and Faisel, nothing more, and the attempt ended in failure, **and when you go hunting for my mistakes I'll go hunting for yours too in the future.***

– block request discussion

C. **Trolling** is likely a frequent occurrence within Arabic Wikipedia. Individuals viewed as trolls will often engage in **circular argumentation** for the sake of argument rather than resolution, insisting that they are acting in good faith while accusing interlocutors of acting in bad faith. Trolling occurs via disruptions of ongoing conversations and prolongation of arguments, and can lead to extended, aggravated disputes. Although

this behavior occurs—likely frequently, as evidenced by a **administrator listing trolling as the most significant threat** to building trust in the community—it is methodologically difficult to identify. Trolling is only successful if it remains unidentified, and its identification requires subjective interpretation of audience perception and speaker intention. Indeed, the term "**troll**" as a description of an opponent or their contributions in Arabic Wikipedia disputes is relatively rare, making keyword searches for trolling accusations problematic. The community's preference for academic, formal Arabic may discourage the use of certain English loanwords such as "troll," indicating that alternative terms or framings may be more common to describe this behavior, and that more research may be required to understand its full extent within Arabic Wikipedia.

*He has become a **troll**; he purposely tries to **obstruct the documentation** of Wikipedia and its editors and moderators and other people who work to develop it.*

– block request discussion

> D. **Harassment**, in the form of repeated, unwanted, negative attention directed by one or more aggressors against the victim and their contributions, occurs frequently on Arabic Wikipedia, and indeed is a frequent accusation leveled against individuals reported via the *Block Request* mechanism. Individuals' definitions and perceptions of harassment tend to vary greatly, however, as evidenced in the fact that multiple overlapping terms are used to refer to it, including تحرش *taḥarrush* and مضايقة *muḍāyaqa*, which in turn overlap with Wikihounding and bullying.

*It's clear that you are **harassing the user** by the way that you're following his edits.*

– block request discussion

> E. **Sexual Harassment** is somewhat less visible as a practice in the community forums and talk pages of Arabic Wikipedia that were accessed pursuant to this study. Actions such as **flirting** and **gendered insults** exist and are flagged by the community, however their full extent in Arabic Wikipedia remains nebulous. In a conversation on the Facebook page of an active Arabic Wikipedia group, a prominent administrator admonished a new editor, "please don't treat the Encyclopedia like a forum for dating or

entertainment," indicating that flirting is a frequently encountered communicative practice. This study uncovered evidence that deeper structural problems, such as the **collective overruling and marginalization of female users' contributions** may occur, but this phenomenon requires further study.

Finally, although **vandalism** lies somewhat outside the scope of the current study, this effort has produced several findings related to this practice in Arabic Wikipedia. First, vandalism is clearly a frequently engaged-in practice (25% of analyzed block requests), however interview respondents did not identify it as a significant barrier to trust within the community. Articles about the **Arab tribes**[4] and various Gulf and Jordanian **royal families** are both mentioned and observed to be frequently subjected to coordinated and individual vandalism and non-encyclopedic contributions by new or rarely-active users. Conflicts on the talk page of more contemporary topics (e.g., Jamal Khashogji, the Egyptian Armed Forces, Israel-related articles, etc.), on the other hand, often appear in block request discussions due to **disputes between established editors**.

# 3.3 Talking about disputes: Maintaining the frame

Dispute participants who narrate and justify their actions in community forums such as the *block request* page, as well as email interview participants recalling previous instances of abuse, rely heavily on two overlapping **discursive frames** to characterize interpersonal conflicts within Arabic Wikipedia. Discursive frames are **culturally-bound templates** used to represent and interpret situations, events, and relationships. In the current study, dispute participants recalling or describing past abuse frequently cast these interactions in either the **human rights** frame or the **objectivity** frame, using the language of each to represent actions and events.

### 3.3.1 Human rights frame

Dispute participants employing the **Human Rights Frame** cast their interactions with perceived abusers as akin to the relationship between the dictatorial state and the powerless citizen. Incidents of abuse cast in this frame are described as violations of "**rights**" and

---

[4] See, for example, see this audit request of a user reported for repeatedly changing the surnames of historical figures to ascribe them to the Anazzah tribe.

"**dignity**," for which aggrieved parties are entitled to demand "**redress**" and "**restitution**" in the form of sanctions imposed by the moderation team on the aggressor. This is especially true when participants see themselves as the losing or aggrieved party in the dispute. Democracy, as in the "**peaceful democracy**" used to describe the nomination process for administrators by an interview participant, is framed as a positive aspect of the community, although abusers are able to "**exploit the democratic climate**" of Wikipedia for their own ends. The Human Rights frame is often invoked when dispute participants perceive an imbalance in community power between participating parties. Individual administrators who weigh into disputes as arbiters or who directly participate in them, for example, may be described as wielding their administrator privileges in a "**dictatorial**" fashion. One interview participant characterizes the perceived ongoing conflict between some editors and some administrators as one "**of the powerful over the weak**."

*…he deleted a number of old and new articles because of "his personal problem with the name of the article," using his moderator privileges in a **dictatorial** way, without returning to the deletion discussion that Wikipedia insists on and **without heeding the policies** of Wikipedia. This is a very serious issue.*

– email interview excerpt

### 3.3.2 Neutrality frame

Dispute participants likewise frame conflict in terms of **neutrality**, appealing to ideals of objectivity, and invoking the values of "**the Encyclopedia**." In this frame, the Encyclopedia and neutrality are aspired-to concepts, while bias and ideology are devalued. Dispute participants invoking this frame do not cast themselves as disagreeing on the fundamental values of the community, but rather as disagreeing over each other's ability to interpret and apply these values. Claims that an opponent is biased, relies on weak sources, or otherwise engages in "non-encyclopedic" contributions represent an attempt to undermine the opponent's legitimacy within the community.

*…there are a number of **ideological possibilities**, also, as well as similarities and overlaps between the Encyclopedia's administrators, and I mean here specifically **adopting the same opinions and ideas** and so on…*

– block request discussion

In the statement above, an active editor expresses their view that the moderation team is **ideologically motivated** in their moderation activities, i.e., that by possessing ideological motivations, the administrators have placed themselves outside the bounds of acceptable behaviors within Wikipedia and thus are inherently illegitimate arbiters. This statement is in fact used as evidence against the speaker in a discussion of their possible block.

## 3.4 Edit wars: A continuation of politics by other means

Edit wars and block requests between established and experienced users are particularly common in the context of articles that have contemporary political significance. This research revealed a number of wide-ranging and enduring disputes on articles including the [Egyptian Armed Forces](#), Egyptian president [Abdel Fattah el-Sisi](#), and articles related to Israel. In these cases, ongoing interpersonal disputes characterized by political or ideological disagreements can play out via community mechanisms for mediating conflict, such as through repeated block requests exchanged between the same editors for contributions to the same article over the course of years.

## 3.5 Politeness as a weapon

Underscoring the difficulties inherent in identifying instances of harassment, trolling, threats, etc. is the fact that these speech acts are often embedded within or accompanied by overt markers of politeness intended to either soften the impact of a face-threatening utterance,[5] or possibly to reduce the likelihood that such an utterance might be reported. In these cases, in which politeness is used as a vehicle for criticism, the **form** of the utterance can be face-affirming, while its **function** remains functionally face threatening. In the excerpt below,

---

[5] In interactional sociolinguistics, **politeness theory** assumes that speakers are concerned with maintaining **face**, i.e., their desired public image (Brown & Levinson, 1987). While positive face refers to the desire to be accepted and approved of by others, negative face refers to the desire for autonomy and freedom from constraint. Although specific desired attributes vary, face needs operate in all cultures and affect both the sender and the receiver in an interaction, for example in the formulation of speech acts such as requests or apologies. **Face-threatening acts** are contrary to the face needs of senders and receivers, some speech acts may be face-threatening by definition, including insults, apologies, and demands. For example, these acts may threaten receivers' face by impinging upon their autonomy (negative face), while an apology may conversely threaten the senders' face because admission of wrongdoing may elicit disapproval from others (positive face).

for example, the author responds to user Alaa's claim to have an upcoming test that is interfering with editing duties:

*My dear Alaa, you don't have any test, and you don't study at all. From what I've seen your contributions [are] very large, goodbye.*

– personal talk page

After employing a polite term of address ("my dear Alaa"), the author continues to deny Alaa's claims and to furthermore indicate ongoing surveillance or observation of Alaa, an implicit claim to power on the basis of access to information that Alaa may have wished to conceal. In another example, a dispute participant admonishes their opponent by writing that, "your problem, my dear, is that you think the goal of Wikipedia is to serve your opinions, and that you're the only one who has the right to edit." Although **embedded in an affectionate term indicating solidarity** with the audience, the criticism contained in this utterance directly undermines the audience party's legitimacy within the ARWP community, in which correct interpretation and application of Wikipedia's values are implicitly understood to be paramount.

Apologies and self-deprecation are likewise used as **amplifiers or vehicles of criticism** of opponents. In the excerpt below, for example, the author initiates a series of overlapping threats and insults with the self-deprecating admission of being "ignorant and illiterate and stupid," before proceeding to threaten the opponent and tarnish them as a "fabricator" guilty of of disingenuous participation in Wikipedia:

*Really **I'm ignorant and illiterate and stupid**, among other things, but I won't allow you to fabricate any information on Wikipedia (your account has been put under close observation), **if you really wanted to participate in Wikipedia,** try starting articles about your country (your IP address has been identified) which everyone associates with terrorism!*

– personal talk page

On the other hand, "genuine" apologies likewise surfaced during this research effort, affirming the presence of community norms that value politeness and civility. Acts of apology are always

**face-threatening for the speaker**. After being admonished for an aggravated response to criticism, for example a user responds:

*__I'm sorry for my tone in speaking__, I got emotional when I didn't find the edits, but __thank you for helping me understand the problem__. I'll try to write useful things in my sandbox about this and I'll show you when I finish.*

– personal talk page

This research has revealed both apologies that appear to be genuine, as well as apologies that are so effusive as to appear disingenuous.

## 3.6 Importance of bystander intervention

Bystanders play a crucial role in acts of bullying, harassment, and aggression in both online and offline contexts. Although bystander intervention has the potential to **de-escalate interpersonal disputes** and improve outcomes for victims, rates of bystander intervention on behalf of victims are typically low in a variety of contexts. In the context of ARWP, in which a number of aggravated disputes among established community members can be observed, bystander intervention has the potential to mitigate negative outcomes for victims. Victims finding a "defender" within the community may in turn experience less anxiety or desire to leave ARWP. In the current study, bystander intervention was observed to have occurred in the moderation of block request discussions, as in the excerpt below in which a administrator places the disputing parties on equal footing by reminding both of the community's norms of discussion and thanking them for their participation:

*But the important thing here is that __"the language of the discussion above" is unacceptable__, especially from two experienced editors like you. Therefore I urge you to exercise wisdom and to discuss the issue in the article's talk page, and you can request the intervention of any colleagues from the team wherever you disagree. Thank you to both of you. – [administrator]*

– block request discussion

In the current study, interview participants consistently stressed the importance of the intervention of sympathetic administrators in the navigation of disputes. Just as the "dictatorial" wielding of power by some administrators is seen by some in the community as a key impediment to the development of trust, the intervention of a sympathetic administrator is likewise cited by interview participants as a crucial link in the chain of editor retention. One participant recalls that, after being subjected perceived as abuse directed against their national identity (Kurdish), they made the decision to leave Wikipedia and delete their account. Ultimately, however, the ongoing support of a sympathetic administrator aided in the decision to remain an active member in ARWP, in spite of the perceived ongoing injustice resulting from this incident:

*...finding some **stakeholders [i.e., administrators] standing by my side** also prevented me from carrying out this negative decision of mine, but so far my rights have yet to be restituted.*

– email interview excerpt

A second participant indicates that, in the event they are subjected to similar abuse, their first course of action would be to "**request the intervention of one of the administrators** in whose neutrality I trust and I will be comforted by their presence." Bystanders can play a number of socially- and culturally-bound roles as interpersonal disputes play out in public spaces within Wikipedia, including supporters/defenders of the victim, disengaged onlookers/outsiders, or reinforcers/assistants of the aggressor. Further research may be able to shed light on the specific forms and functions associated with bystander roles within the context of this community, although this research has established that bystander support can be a crucial factor in mitigating the negative outcomes of abuse.

## 3.7 "Leaving Wikipedia": Harassment is highly personalized

Rather than representing a singular event, **harassment is an ongoing process** that is both targeted and personalized. Even experienced and dedicated Wikipedia editors may be driven by the trauma associated with harassment to leave the community. In the current study, an interview participant recalled having been harassed by another user who had created a series

of sock puppet accounts. The victim reported that this harassment took place across multiple sites within ARWP, including the victim's contributions, talk page, and personal page, and that this highly individualized attack ultimately led them to "**think about leaving Wikipedia**", an experience reported by other harassment victims. Although the victim in this case attributes their decision to remain in the community to a sense of personal perseverance, a second interviewed harassment victim cites the intervention of a sympathetic administrator as crucial to their continuation in the community. These experiences highlight one of the major risks associated with harassment in the context of Wikipedia, namely that otherwise active and dedicated community members may be driven away, weakening the community as a whole.

## 3.8 Ongoing concerns with administrator neutrality

This research effort revealed an apparent ongoing rift within the Arabic Wikipedia community in which a number of active editors see themselves as opposed by a small clique of administrators, who are in turn perceived to be biased, ideologically driven, and closed to new ideas and modes of communication within the community. Although this phenomenon is somewhat beyond the scope of the current research, it was nevertheless the subject of several instances of observed insults and threats, and it was likewise cited as a main factor impeding the development of trust within the community by two non-administrator interview participants. In one case, for example, an editor publicly referred on a administrator's *talk* page to "**a big gang of sugar and honey**," indicating that the moderation team speaks kindly but in reality operates as a gang. In other instances, editors speculate in public spaces on Wikipedia that the moderation team is collectively biased, or complain about them on in non-Wikipedia spaces such as Facebook. Further research—such as data collection via survey instruments—will likely be able to shed more and useful light on this phenomenon, however at this time it is evident that this perceived editor–administrator rift is a sometimes-observable feature of community discussion forums and likewise featured prominently in the interviews conducted pursuant to this study.

# 4. Recommendations for future research

This exploratory study has revealed a number of ongoing phenomena of relevance to the Anti-Harassment tools team in the context of Arabic Wikipedia. Although this study addresses the *how* of harassment and threats among some users of ARWP, the *why* begs further investigation. Why do experienced users frequently appear trading insults and threats in the *block request* page, for example? Why do new editors appear to be absent from this page, unless they are being reported for abuse? How might future efforts leverage education and awareness building to promote bystander intervention during interpersonal disputes? To address these and other questions, this report concludes with some suggestions for future avenues of research by the WMF.

**Suggested avenues for future research**

- **Initiate a survey study** to provide more context and nuance to the phenomena discussed in this report.

A formal survey of the **27 administrators of Arabic Wikipedia** would allow community exports a standardized format through which to rate, assess, and provide suggestions for the improvement of community functions. A **larger survey of active editors**, on the other hand would provide more a more robust framework for conceptualizing problems facing the community. A survey of either type would be an excellent source of recruitment of participants for further email interviews given the inclusion of an option for respondents to indicate interest in this activity. Finally, both surveys could be combined in the same instrument, directing respondents to versions of the survey following their self-identification as "administrator" or "editor." Furthermore, it is possible that a "survey" is favored within this community, and that recruitment for a survey may be greater than recruitment for participation in an "email interview."

- **Conduct a longitudinal study of new editors' experiences**.

Rather than focusing recruitment on the *block request* page, where potential participants may be motivated to participate due to an ongoing grievance, a longitudinal study of new editors would recruit from among active new editors and could make use either of interview methods, survey instruments, or both. For example, such a study might ask new users to periodically respond to a unified survey or questionnaire as their participation in Arabic Wikipedia develops, allowing their changes in stance toward various entities to be tracked over time. By providing new users with an avenue for recording and reporting experiences, this effort might uncover previously unreported challenges to the new editor experience, and may shed light on phenomena of harassment or abuse that new editors might experience but choose not to report.

- **Expand research to other-language Wikis** in order to gain a clearer view of the editing and moderation landscape with respect to issues of community trust across Wikipedia as a whole.

Investigating the forms that harassment and trolling take in the **Spanish** (73 administrators) and **German** (191 administrators) Wikipedias, for example, will provide more context to the current findings and help indicate whether the phenomena described here are unique to Arabic. Like Arabic, Spanish Wikipedia unites editors from multiple countries, while German Wikipedia may provide a counterpoint by exhibiting more geographical unity. Similar research in these and other-language Wikis likewise has the potential to reveal successful innovative approaches to mitigating barriers to community trust that might be adapted to Arabic.

- **Compile a case study profile of one or more abusers and victims**.

This effort would seek to explore harassment from multiple perspectives in order to determine, for example, at what point in a devolving interaction would the participants be most susceptible to bystander intervention. Given that a large proportion of the insults and threats collected in this study were produced by active community members, there is ample

opportunity to explore interpersonal conflicts from the perspective of both the insulting and insulted parties.

**Suggested product innovations**

- **Create a dedicated research page** that clarifies, explains, and otherwise publicizes ongoing research efforts in this vein.

A dedicated research page will help lay the groundwork for relationships with future participants by communicating research objectives and possibly providing a forum for feedback, for example through the submission of "experience reports" in which respondents describe or report incidents of the type under study. Such a space will also be helpful in legitimizing and providing visibility for research efforts. This page should be made available in all languages of active research and could link to products or tools produced as a result of these efforts. Finally, such a page would positively serve WMF by improving transparency and indicating WMF engagement in the effort to improve the functioning of its disparate communities. An email interview participant noted that this type of research by WMF, "if it doesn't solve the problem [of harassment], might at least disrupt it."

# 5. Conclusion

This research has established the fact that harassment and other types of threatening communications occur at a non-trivial rate in the context of Arabic Wikipedia. These phenomena may be one-off events in which pseudonymous aggressors take advantage of their relative anonymity to abuse others, or they may likewise be indicative of deeper rifts that continue to divide the community. Indeed, many of the instances catalogued by this effort consist of ongoing disputes between active community members who adeptly use the mechanisms of community moderation (e.g., the block request page) as strategic resources in their contest for power and legitimacy within the community. The discursive frames used by community members to talk about disputes are indicative of the value that community members place on issues such as human rights, neutrality, and objectivity, and this research has likewise revealed that certain positive behaviors, such as bystander intervention, are crucial to mitigating the negative outcomes experienced by abuse victims. Moving forward,

more research is needed to contextualize and complexify the phenomena discussed here, however it is clear that the WMF's engagement in this type of research promotes transparency and is well received by the community.

## References

Brown, P., & Levinson, S. C. (1987). *Politeness: Some universals in language use.* New York: Cambridge University Press.

Meho, L. I. (2006). E-mail interviewing in qualitative research: A methodological discussion. *Journal of the American Society for Information Science and Technology, 57*(10), 1284–1295. https://doi.org/10.1002/asi.20416

Petersen, L. N. (2017). "The florals": Female fans over 50 in the Sherlock fandom. *Transformative Works and Cultures, 23.* https://doi.org/10.3983/twc.2017.0956

Hollway, W., & Jefferson, T. (2008). Researching Defended Subjects with the Free Association Narrative Interview Method. In L. M. Given (Ed.), *The SAGE Encyclopedia of Qualitative Research Methods* (pp. 296–315). Sevenoaks, CA: Sage.

# Appendix A: Lessons learned for future asynchronous (email) interview studies

This research effort resulted in a number of lessons learned concerning best practices in conducting asynchronous (email) interviews with members of Wikipedia's global community:

A. **Transparency is paramount**—potential participants may be suspicious of recruitment efforts by non-WMF-associated individuals. This may be even more notable in contexts such as Arabic Wikipedia, which sees the participation of individuals from countries in which national intelligence services are known to closely monitor their citizens' online activities. To alleviate any confusion during the recruitment phase, recruiting calls should be associated with a recognized name or team in the community.

B. **Response rates may be low**, and study sponsors may wish to consider what an "acceptable" number of responses represents before recruitment efforts begin, in the expectation that several recruitment emails will likely have to be sent for each successful interview (a 6:1 ratio of recruitment emails to successful interviews in the current study). In the absence of examples of similar previous studies, the expected recruitment rate was unclear. In general, prospective participants in online surveys or email interview studies often delete invitations before they are read or change their email addresses. However, reminders can significantly increase participation rates. Recruitment pools can also be expanded in the event that initial recruitment is limited.

C. **Question sequencing determines the trajectory of interviews**, and should therefore be carefully considered. Petersen (2017) recommends submitting seven introductory questions to each participant, followed by a series of one to five follow-up emails containing two to four questions each. This method allows both for multiple interviews to be conducted simultaneously, as well as for the researcher to have access to multiple participants' initial responses before selecting or composing follow-up questions for each. In the current study, the initial sevel questions contained in the second email to participants were shortened to four questions in order to lower the barrier to response.

D. **Different participants may have different timelines for participation**. Meho (2006) notes that typical time to completion of data collection in studies relying on email interview protocols vary widely, with researchers completing data collection in anywhere from one week to several months. The length of the data collection period may depend on several factors, including the number of participants, the number of questions posed, participants' motivation to contribute, participants' engagement with Wikipedia at the time that questions are posed, and the time that both participants and interviewers are able to devote to the interviews.

E. **Questions should be carefully worded,** and if possible should be pilot-tested. Due to the lack of opportunity for clarification, lack of nonverbal cues, and the ability of participants to take their time in responding for questions, Meho (2006) notes that email interview questions must be "much more self-explanatory than those posed face-to-face, with a clear indication given of the responses required" (p. 1290). Hollway and Jefferson (2008) recommend using open-ended questions to elicit stories, and to avoid "why" questions. The wording of individual questions in this study was amended after the initial round of responses due to confusion on the part of a single respondent, for example.

# Appendix B: Categorization of block requests

This table presents a categorization of block requests spanning the period from January 2017 to December 20th, 2018.

| Complaint type | Instances | Description |
|---|---|---|
| **Username** | 315 | Inappropriate username, often due to containing the name of a prominent person, company, or containing profanity. |
| Blocks granted | 276 (87.6%) | |
| **Vandalism** | 197 | Specific accusation of "vandalism". |
| Blocks granted | 147 (74.6%) | |
| **Non-encyclopedic content** | 86 | "Non-encyclopedic edits," e.g., repeatedly re-creating rejected articles, self-promotion, edits via poor machine translation, inappropriate articles and contributions, etc. |
| Blocks granted | 46 (53.4%) | |
| **Personal attack** | 67 | Engaging in "personal attacks," issuing threats, insults, "harassment," or other abusive behavior. |
| Blocks granted | 27 (40.3%) | |
| **Propaganda** | 43 | Political or commercial propagandizing—e.g., using contributions to promote a particular company or political point of view. |
| Blocks granted | 33 (76.7%) | |
| **Sock puppet** | 29 | Accusations of being sock puppet or otherwise maintaining multiple accounts. |
| Blocks granted | 20 (68.9%) | |
| **Trivial account** | 25 | Accounts accused of existing "purely for annoying," socializing, flirting, etc. |

| | | |
|---|---|---|
| Blocks granted | 18 (72%) | |
| **Use of profanity** | 11 | Use of inappropriate language, curses, etc. The contents of these accounts' reported posts are quickly removed. |
| Blocks granted | 10 (90.9%) | |
| **Engaging in edit war** | 9 | Users accused of systematically, deliberately, and/or maliciously reverting the edits of another user. |
| Blocks granted | 0 (0%) | |
| **Total** | 782 (577 blocks, 73.8%) | |