# Political Polarization on the Digital Sphere: A Cross-platform, Over-time Analysis of Interactional, Positional, and Affective Polarization on Social Media

Moran Yarchi, Christian Baden & Neta Kligler-Vilenchik

|  |  |
|---|---|
| View supplementary material ⎘ | Published online: 14 Jul 2020. |
| Submit your article to this journal ⎘ | Article views: 10451 |
| View related articles ⎘ | View Crossmark data ⎘ |
| Citing articles: 77 View citing articles ⎘ |  |

Routledge
Taylor & Francis Group

RESEARCH ARTICLE

Check for updates

# Political Polarization on the Digital Sphere: A Cross-platform, Over-time Analysis of Interactional, Positional, and Affective Polarization on Social Media

Moran Yarchi [a], Christian Baden[b], and Neta Kligler-Vilenchik[b]

aThe Sammy Ofer School of Communications, Interdisciplinary Center Herzliya, Herzliya, Israel; bNoah Mozes Department of Communication and Journalism, Hebrew University of Jerusalem, Jerusalem, Israel

**ABSTRACT**

Political polarization on the digital sphere poses a real challenge to many democracies around the world. Although the issue has received some scholarly attention, there is a need to improve the conceptual precision in the increasingly blurry debate. The use of computational communication science approaches allows us to track political conversations in a fine-grained manner within their natural settings – the realm of interactive social media. The present study combines different algorithmic approaches to studying social media data in order to capture both the interactional structure and content of dynamic political talk online. We conducted an analysis of political polarization across social media platforms (analyzing Facebook, Twitter, and WhatsApp) over 16 months, with close to a quarter million online contributions regarding a political controversy in Israel. Our comprehensive measurement of interactive political talk enables us to address three key aspects of political polarization: (1) interactional polarization – homophilic versus heterophilic user interactions; (2) positional polarization – the positions expressed, and (3) affective polarization – the emotions and attitudes expressed. Our findings indicate that political polarization on social media cannot be conceptualized as a unified phenomenon, as there are significant cross-platform differences. While interactions on Twitter largely conform to established expectations (homophilic interaction patterns, aggravating positional polarization, pronounced inter-group hostility), on WhatsApp, de-polarization occurred over time. Surprisingly, Facebook was found to be the least homophilic platform in terms of interactions, positions, and emotions expressed. Our analysis points to key conceptual distinctions and raises important questions about the drivers and dynamics of political polarization online.

Over the last decades, processes of digitalization and mediatization have changed how we communicate and interact in political talk. Due to the rise of digital social media, political discussions that might have previously occurred during lunch breaks or over the dinner table have entered the realm of public, computer-mediated communication (Lelkes, 2016; Stroud, 2010; Van Aelst et al., 2017). Consequently, political communication research has seized the opportunities offered by the abundant availability of public political talk to study political interactions on social media, develop new digital methods, advance theory,

**CONTACT** Moran Yarchi ✉ moran.yarchi@gmail.com 💬 Moran Yarchi, the Sammy Ofer School of Communications, Interdisciplinary Center (IDC) Herzliya Kanfei Nesharim St. P.O.Box 167 Herzliya 46150, Israel

and raise new questions. Research into political polarization on social media is at the core of this enterprise. Observing a structural transformation of the political public sphere, many have voiced concerns that the ways in which social media shape political talk may contribute to a progressing fragmentation of political debates, increasing polarization, and growing hostility among political camps. Related to recent global occurrences, such as the 2016 US elections, the "Brexit" referendum, or the debate surrounding the refugee crisis in European countries, it has been argued that political polarization on social media poses a key challenge to contemporary democracies.

In the present paper, we use a novel computational approach to the study of political talk on social media in order to investigate the hypothesized linkages between the changing structure of digital interactions and political polarization. Using a computational communication science approach allows us a better understanding of public political discussions in their natural setting: interactive social media. Specifically, our computational strategy enables us to merge existing tools in our analysis in order to simultaneously capture expressed positions, interpretations, and sentiments as well as the structure of interactions, offering a more comprehensive view of political talk on social media. By linking content classification with relation extraction, we create a rich interaction network and focus on its dyadic interaction structures.

In doing so, we address three key challenges that have limited the ability of political communication researchers to fully capture the new dynamics of political talk in digital networks. First, in an effort to improve conceptual precision in the increasingly blurry debate (Baylis, 2012), we review the theoretical arguments tying political polarization to the specific changes brought about by social media. Specifically, we distinguish three aspects of political polarization: interactional polarization, positional polarization, and affective polarization. While all three aspects have been used to operationalize political polarization in recent political communication scholarship, each is related in different ways to the reconfiguration of online political talk.

Second, researchers' heavy reliance on Twitter's open API and the comparative difficulty of studying other platforms have obscured the important differences that exist between social media platforms, which structure political talk in different ways. By studying political polarization on three different social media platforms – Twitter, Facebook, and WhatsApp – our study's comparative design contributes to a better understanding of how specific algorithmic and social features of social media contribute to different polarization dynamics.

Third, due to limited interoperability between available digital methods, most studies have focused either on networked interactions or the textual contents of social media interactions, but few have combined these perspectives – even less so in a diachronic perspective. The present study combines the extraction of discourse contents with an analysis of networked interactions, tracking changes in both contents and interactions over time. This integration is enabled by an innovative methodological approach, which links algorithmically extracted relational data to the computer-assisted deductive classification of textual contents. Relying on refined conceptual distinctions, we examine different aspects of polarization in political talk online in its dynamic, natural environment (Garcia et al., 2012; Iyengar et al., 2012) to understand how different social media platforms affect evolving interaction patterns, positions, and relational attitudes in contrasting ways.

Our analysis focuses on social media interactions surrounding a political controversy in Israel, within the context of the Israeli–Palestinian conflict. Unlike the US, whose predominantly binary party system has been the focus of most existing research on polarization, Israeli society is characterized by a variety of cross-cutting conflict lines, enabling different configurations of opposing camps depending on the issue and focus of the debate. Using close to a quarter million social media posts on Facebook (publicly open discussions), Twitter and WhatsApp, collected over a period of 16 months, we study how each platform contributes to interactional, positional, and affective polarization. Capturing both content characteristics and interaction structures of social media political talk, over time and across different platforms, we use computational methods to offer a deeper understanding of political polarization.

## *Political Polarization and Social Media*

Much of the current concern about the rise of polarized political discourse is tied to the growing role that digital media play in structuring public communication (Dvir-Gvirsman et al., 2016; Van Aelst et al., 2017). From a traditional focus on polarization as the outcome of political-ideological and partisan identities, the rising visibility of interactive political talk on social media has shifted scholarly attention toward the potentially polarizing effects of homophilic communication patterns (Settle, 2018). Scholars have introduced new approaches to the study of polarization in digital media, typically employing computational methods to extract social network structures, semantic contents, sentiments, or other properties from rich social media data.

However, the methodological diversification was accompanied by a conceptual blurring: Owing to the difficulty of transferring traditional, survey-based measures to the realm of social media, scholars have focused on different proxies as measures of political polarization. Garcia et al. (2012) extracted specific issue positions from digital text, relying on computational dictionaries and natural language processing tools. Barberá et al. (2015) analyzed homophilic interaction patterns, operationalizing political polarization as the fragmentation of public spheres. Iyengar et al. (2012) studied polarization as the expression of hostile attitudes toward political rivals. In our own work (Kligler-Vilenchik et al., 2019) we have argued that the interpretation of contentious issues in incompatible terms can also be viewed as an aspect of political polarization. Polarization thus may refer to different kinds of attitudes (e.g., toward issues, in- and outgroups) and beliefs (e.g., about the nature of contested issues), and may be manifested in different kinds of behavior (e.g., verbal expression, interaction patterns, political choices).

At the same time, little attention has been paid to the possibility that these aspects may be affected in different ways by the structure of social media conversations. In the following, we will distinguish three key aspects of political polarization – interactional polarization, positional polarization, and affective polarization – and discuss their interrelation with communication practices specific to different social media. In particular, we (1) examine why social media are believed to contribute to each aspect of polarization; (2) discuss whether each aspect can give rise to possible self-reinforcing tendencies in polarization; and (3) consider the ways in which polarization can be expected to take different shapes on different social media platforms.

### Interactional Polarization

The first aspect of political polarization, which we call "interactional polarization," focuses on a process whereby participants in a debate increasingly interact with like-minded individuals, while disengaging from interactions with others who hold opposing viewpoints. As a consequence, social media publics may become fragmented into two or more opposing camps (Bodrunova et al., 2019), with only little political talk cutting across political conflict lines. Interactional polarization threatens democratic governance by increasing the number of relays needed to exchange viewpoints and demands between competing societal groups, obstructing their negotiation and constructive processing.

Despite the recent salience of theories regarding fragmented "echo chambers" or "filter bubbles," it remains contentious whether social media do indeed drive such interactional polarization (Bruns, 2019). On the one hand, and in line with optimists' view on digital media's deliberative potential, social media have massively broadened the range of voices enabled to participate in public communication: Nearly anyone, regardless of skills and resources, is able to receive, engage with, and publicly challenge any public statement (Brundidge, 2010; Holt, 2004). On the other hand, the networked structure has facilitated a public sphere constituted primarily by individuals sharing specific interests or views (Colleoni et al., 2014; Nahon, 2016), as users can exclude deviant voices from their communication networks. Social media platforms have been associated with the rising prominence of selective exposure, as they facilitate users' efforts to avoid exposure to unwanted content (Settle, 2018). In addition, their algorithmic filtering and propagation privilege contents deemed likely to be of interest to, and in line with the views of users (Mutz & Martin, 2001; Nahon, 2016). Accordingly, social media users might find themselves exposed predominantly to congenial information in increasingly "homophilic" communication networks, contributing to a progressing fragmentation and polarization of the online public sphere (Bobok, 2016; Hayat & Samuel-Azran, 2017).

While most studies have confirmed an overall tendency toward homophilic interactions in digital environments, this is true also for interactions outside social media. Heterophilic interactions appear to be more common along so-called "weak ties" – occasional communications that are not underfed by strong social bonds such as friendship or sustained collaboration – while most "strong ties" (among friends, within teams, etc.) are predominantly homophilic, and are expected to have a stronger impact on individuals, their attitudes, and behavior (Barberá, 2014; Huckfeldt & Sprague, 1995). Nevertheless, recent scholarship (e.g., Bruns, 2019; Zuiderveen Borgesius et al., 2016) has questioned the validity of concerns over filter bubbles and echo chambers, concluding that such phenomena appear to be empirically rare and less potent than commonly feared. There is, to date, no evidence that establishes that homophilic interaction patterns are indeed self-reinforcing, resulting in progressing interactional polarization. Owing to the complexity of dynamic network modeling and analysis, existing network-based studies of political talk on social media are overwhelmingly static, and thus unable to evaluate the suggested dynamic properties of interactional polarization. In our study, we trace diachronic changes in the prevalence of homophilic interaction patterns, permitting us to hypothesize the following:

*H1 (Homophily): Interaction patterns on social media are homophilic.*

*H2 (Interactional Polarization): Interaction patterns on social media become increasingly homophilic over time.*

## Positional Polarization

Beyond fragmentation, political discussion on social media has also been tied to increases in antagonistic and extreme political preferences – the traditional focus of political polarization research (Fiorina & Abrams, 2008), which we refer to here as "positional polarization." By preferentially exposing users to congenial viewpoints, social media may decrease participants' awareness to competing views and challenging, possibly corrective information. Confirmed time and again in homophilic interactions (Slater, 2007), users' prior beliefs may be increasingly removed from possible doubt or reservation (Lerman et al., 2016), eroding their tolerance of competing views. In consequence, users' willingness and ability to engage in controversial political debate and seek widely acceptable solutions may be compromised (e.g., Wojcieszak, 2010), up to the point of disengaging from cross-cutting communication ties (John & Dvir-Gvirsman, 2015).

Similarly, the prevalence of one-sided arguments and the marginalization of deviant views may facilitate the expression of increasingly extreme viewpoints (Wojcieszak, 2010). Perceptions of community consensus may foster the emergence of tight-knit group identities and embolden extreme members to speak up. Where extreme views remain unchallenged, these may increasingly be perceived as acceptable, prompting a possible self-reinforcing process (Lee, 2006; Moscovici & Zavalloni, 1969). As interactions with outgroup members become rare, this may erode respect and empathy with outsiders, facilitating the dismissal of their objections and the adoption of policies that would be harmful to them (Colleoni et al., 2014; Mutz & Martin, 2001; Sunstein, 2001).

Inversely, heterogeneous discussion networks have been tied to diminished polarization and radicalization. Cross-cutting talk has been shown to broaden individuals' awareness and tolerance of competing viewpoints and facilitate the adoption of intermediate positions and the search for mutually acceptable compromises (Wojcieszak, 2010). Thus, increasingly homophilic interactions can be expected to aggravate positional polarization, while increased heterophilic interactions facilitate de-polarization and solution-oriented democratic discourse.

While there is experimental evidence linking homophilic discussion networks to positional polarization (e.g., Wojcieszak, 2010), these propositions have not, to our knowledge, been tested in the context of naturally occurring political talk – chiefly due to the challenging extraction of specific issue positions from the natural discourse. In the present study, we rely on an innovative computational strategy to identify those issue positions expressed on social media, hypothesizing the following:

*H3 (Positional Polarization): Individuals embedded within more homophilic interaction networks subsequently express more extreme positions in their contributions.*

## Affective Polarization

Next to fragmentation of interactions and the polarization of political positions, social media have been argued to contribute to rising hostility in political talk toward opposing

political groupings (Settle, 2018), which we refer to as "affective polarization" (Iyengar et al., 2012). Suler (2004) documented a general "online disinhibition effect," wherein the perceived social distance from other participants reduces users' reliance on civil and constructive conversational norms. In the same vein, the marginalization of outgroups in digital communication networks has been tied to an erosion of empathy (Wojcieszak, 2010) and an increase of negative relational attitudes (Iyengar et al., 2012). While political discussion among like-minded individuals has long been known to activate positive feelings (Huckfeldt et al., 2004), discussion with political opponents tends to trigger negative emotions, due to the psychological discomfort of political disagreement (Mutz, 2006; Parsons, 2010). Where social media foster political talk to focus on interactions with likeminded individuals, accordingly, discussions should predominantly express amiable sentiment. At the same time, the diminishment of cross-cutting talk should be accompanied by aggravated negative attitudes toward outgroups. While there are several studies tapping affective polarization via public opinion surveys (e.g., Iyengar et al., 2012), the connection between opposing issue positions, expressed attitudes, and homophilic interaction patterns on social media has not been studied – chiefly due to persistent challenges in the attribution of expressed sentiment to different kinds of interactions (Hillmann & Trier, 2012). In the present study, we capture the sentiment expressed in interactions between both likeminded and opposing users, permitting us to hypothesize:

*H4 (Affective Polarization): Interactions cutting across political camps express more negative sentiment than interactions among like-minded individuals.*

### Political Polarization across Social Media Platforms

Despite the focus on social media's role in reconfiguring political talk, most existing research has relied on data obtained from Twitter's public API. However, on each social media platform, different affordances and features govern how users can contribute and share, follow, and respond to posted contents, which are propagated throughout the network in different ways. Likewise, each platform is adopted by somewhat different communities for different purposes, developing community-specific norms and conventional practices. Accordingly, different social media platforms shape online political talk in non-uniform ways, with differing implications for political polarization.

In her work on political polarization on Facebook, Settle (2018) proposes the so-called END framework, referring to "the characteristics of a subset of content that circulates in a social media ecosystem: a personalized, quantified blend of politically informative *expression, news*, and *discussion* seamlessly interwoven into a wider variety of socially informative content" (p. 50). Settle uses this framework to consider how particular affordances of different social media platforms – e.g., design facilities for identity maintenance, the presentation of political content, the structural importance of opinion leaders, weak ties, social feedback, and quantification – are associated with particular aspects of political polarization. Building on this framework, we will in the following discuss the key implication of such differences for our present investigation.

## Facebook

As the world's foremost social media platform, Facebook's popularity is arguably derived largely from its capacity to immerse its users in a feed of contents that cater to their personal interests and leanings. To do this, the platform relies heavily on users' self-curated networks of friends, but also on an algorithm that prioritizes content based on users' interests and support for similar posts, displaying only a small share of predominantly congenial, supportive contents (Nahon, 2016). In addition, Facebook stands out as a platform rich in emotional content, which is amplified by its algorithm (Tucker et al., 2018). The exalted role of friendly interactions among like-minded users on Facebook is further underscored by the primacy of private and personal interactions. Thus, by crafting a highly personalized, emotional user experience, Facebook has become the prime suspect for the creation of homophilic echo chambers (Ben-David & Matamoros-Fernandez, 2016; Nahon, 2016), and numerous studies have confirmed high levels of homophily among Facebook users (e.g., Bakshy et al., 2015; Bond et al., 2012).

At the same time, not all parts of Facebook are equally sheltered by the layered regimes of privacy that structurally limit users' ability to interact with contents beyond their own network of friends. Besides such private networks, Facebook prominently features also a selection of public campaign and community pages (notably, by social movements, NGOs, political parties, and leaders; Harlow & Harp, 2012; Marichal, 2013), which are accessible also to outgroup users. Such pages, which serve to build efficacious communities of like-minded users and campaign for contentious political aims, may also become targeted by rival groups, resulting in heated political contestation (Hendriks et al., 2016). Despite the presence of cross-cutting contents, contestation on such pages is more likely to exacerbate than mitigate political polarization. Owing to Facebook's algorithmic emphasis on congenial content, its personal and emotional communication culture, and its salient role of like-minded communities, we expect communication on Facebook to be especially polarizing.

## Twitter

In sharp contrast, Twitter is defined primarily by its unrestricted publicness. Anyone, even non-users, can read any tweet, and any user can respond to any contribution. Users can follow others without a need for permission, enabling asymmetric, non-reciprocated ties. Hence, users can follow accounts without becoming members of a community, and interact with users far outside their network of friends and peers. As the platform barely restricts the visibility of contents, competing political demands and identities are easily perceived. At the same time, Twitter, like Facebook, enables the formation of cohesive communities, and its newsfeed personalizes the content presented to each user in response to known interests and views. Both the facility for heterophilic talk and the promotion of homophilic interactions are embedded in its structure (Colleoni et al., 2014). Like Facebook, also Twitter's algorithm appears to privilege emotional contents (Tucker et al., 2018), while Twitter's inbuilt character limitation (to 140 or 280 characters per message) limits the space for argumentative contestation (Nahon, 2016; Sagolla, 2009).

In many countries, Twitter is popular primarily as a tool for professional communication (notably, among political elites and media), and means for advocacy and public relations. Comparatively, non-predefined audiences facilitate a less personal communication style often focused on shout-outs in place of sustained conversations. Only in some

countries – notably, the US – is Twitter also widely adopted among lay users who chat about personal matters in public. Reflecting Twitter's ambiguous profile, the existing literature yields conflicting findings regarding its tendency toward homophily and polarization (e.g., Kwak et al., 2010 detected little homophily; Weng et al., 2010; Hong & Kim, 2016; Colleoni et al., 2014 found the opposite).

### WhatsApp

WhatsApp does not constitute a social media platform in a narrow sense (following Ellison and Boyd (2013) definition); however, WhatsApp discussion groups can be used in ways comparable to social media and are increasingly popular in several cultural contexts (Boczkowski et al., 2018), including Israel. In this paper, we consider two WhatsApp groups devoted to political talk, moderated by a journalist. WhatsApp groups are closed platforms, and any interaction requires a user to first gain access to the group. All content is visible to all members without a filtering or highlighting algorithm, and exposure is determined only by the sequential structure of talk. Due to the ongoing nature of the conversation, interactions are likely to be quick and interactive, while returning to past posts is discouraged by the evolving feed. Groups are typically convened around a shared interest, such that, at least with regard to that interest, groups are structurally homophilic (Kwon & Oh, 2014). At the same time, depending on the origin and conduct of the group, it is possible to convene WhatsApp groups that comprise diverse viewpoints upon the common focus. In such a case, Kligler-Vilenchik (2019) found that WhatsApp groups devoted to heterogeneous political talk may encourage cross-cutting exposure and even a deeper understanding of the viewpoints of heterogeneous others.

Given the different platforms' characteristic affordances, we can expect marked differences in how each platform advances political polarization, necessitating further differentiation in our above hypotheses. Table 1 summarizes the major affordances of each platform, which inform our refined hypotheses.

As users are known to preferentially interact with familiar and like-minded others (Affordances 1, 2 in Table 1 – all subsequent numbers refer to affordance numbers in Table 1), we hypothesize:

*H1a (Homophily across Platforms): Interactions will be most homophilic on Facebook, less so on Twitter, and least so on WhatsApp.*

**Table 1.** Major platform-specific affordances & refined hypotheses.

| Affordances | Facebook | Twitter | WhatsApp* |
|---|---|---|---|
| 1. Focus on strong ties | + | • | − |
| 2. User control over ties & exposure to content | + | • | − |
| 3. Commitment to controversy | − | • | + |
| 4. Algorithmic content personalization | + | + | − |
| 5. Focus on personal identity | + | − | + |
| 6. Shared group identity & culture | • | − | + |
| **Refined, platform-specific hypotheses** | | | |
| H1a: Homophily across platforms | + | • | − |
| H2a: Interactional Polarization across platforms | + | + | − |
| H3a: Positional polarization across platforms | + | − | • |
| H4a: Affective polarization across platforms | + | + | − |

+ high/pronounced; • intermediate/partly; − low/absent; *Affordances/hypotheses apply to the political discussion WhatsApp groups examined in this study, and not to WhatsApp in general.

Users are expected to reduce cross-cutting interactions if they can (2), unless they value controversy (3):

*H2a (Interactional Polarization across Platforms): Homophily will increase over time on Facebook and Twitter, but not on WhatsApp.*

In turn, users' exposure to congenial content (1, 2, 4) facilitates their adoption of more extreme positions, provided they are invested in a debate (5) and regard others' views as valuable (3, 6):

*H3a (Positional Polarization across Platforms): The effect of homophilic interactions on the extremeness of expressed positions will be most pronounced on Facebook, followed by WhatsApp, and least pronounced on Twitter.*

Finally, cross-cutting hostility is expected to increase with threat to personal identity (5), but is mitigated by commitment to controversy (3) and shared group identity (6):

*H4a (Affective Polarization across Platforms): Negativity expressed in contributions addressed toward outgroups will be less pronounced on WhatsApp than on Facebook and Twitter.*

## The Case Study

For the present study, we focus on an incident sparking a heated debate in Israeli society, including various social media. On March 24, 2016, two Palestinian assailants stabbed and wounded an Israeli soldier near the Westbank town of Hebron. One assailant was shot dead, while the other, Abed al Fatah al-Sharif, was wounded, disarmed, and "neutralized." Minutes later, IDF Sergeant Elor Azaria shot the incapacitated assailant in the head, killing him. Azaria was arrested and a military investigation opened. A video of the event went viral on Israeli social media, sparking controversy regarding the IDF rules of engagement in response to Palestinian violence. Some praised Azaria's deed as heroic and justified in the defense against violent terrorism, while others denounced him as a hateful murderer and demanded he be punished (Livio & Afriat, 2019). The debate, including ongoing demonstrations, continued throughout Azaria's trial for manslaughter before a military court, which was opened on May 9, 2016. On January 4, 2017, Azaria was convicted, and on February 21, 2017, he was sentenced to 18 months in prison. Azaria filed for an appeal, which was rejected on July 30, 2017.

Due to major ideological differences among Israeli society, the Israeli–Palestinian conflict is particularly relevant to the examination of polarized political discourse (Aronoff, 1984). During the entire process, several representative public opinion polls showed a public split over its evaluation of Azaria and his actions. In a survey conducted in August 2016, 47% of Israelis justified the killing, agreeing that "every Palestinian who has committed an attack on Jews should be killed, even if he is caught and clearly not endangering the environment"; however, another 45% took the opposite view, agreeing that "once he has stopped endangering the environment, the Palestinian perpetrator

should be handed over to law enforcement."[1] Following Azaria's sentencing, 51% of Israelis stated their disagreement with his sentence,[2] while 36% supported it. The entire controversy generated considerable amounts of heated discussions, galvanizing the attention of numerous influential speakers at a political level, in the media and on social media over an extended period. By tracking these conversations on different social media over 16 months, using a computational communication science approach, we collected rich data to study interactional polarization over time and across platforms.

## Data

For the present study, we collected all posts related to the incident and subsequent debate on Twitter, on publicly open Facebook posts,[3] and from two WhatsApp groups dedicated to political discussion. While Twitter in Israel is primarily used by political and media professionals,[4] Facebook is widely used across most social groups (68.6% penetration rate; internetworldstats.com) and is particularly popular as a tool for community mobilization among right-wing groups. The studied WhatsApp groups, which were hosted by a well-known Israeli journalist, included participants from across the political spectrum, who joined for a nominal fee. The fact that these WhatsApp groups are open to any participant willing to pay the fee, and discuss public matters, creates a unique semi-public space, a rare opportunity for research.[5]

On Twitter and Facebook, we first identified all relevant original posts published between March 24, 2016, the day of the shooting, and August 2, 2017, three days after the rejection of the appeal. To do this, we used a thoroughly validated list of keywords, including different spellings of the shooter's name, a range of labels applied to the case, and a few other expressions used to praise or condemn Azaria. Subsequently, we also included all comments responding to these posts in our sample. On WhatsApp, we manually coded which passages were related to the case, taking into account the sequential structure of the discussion.

Over the entire duration of the controversy, we collected a total of 249,317 relevant contributions: 29,250 tweets with 61,772 comments on Twitter, 6,508 posts with 145,542 comments on Facebook, and 6,245 WhatsApp messages.[6] In order to investigate possible diachronic changes, we split the data into four phases: (I) the days between the shooting and the trial (March 24, 2016–May 8, 2016), (II) the duration of the trial (May 9, 2016–January 3, 2017), (III) the time between the verdict and the sentencing (January 4, 2017–February 20, 2017), and (IV) the months following the sentencing (February 21, 2017–August 2, 2017), which included Azaria's unsuccessful appeal (For more information about the distribution of each platform's publications over time, see Figure A1 in the appendix).

For our analysis of interactions, we extracted all embedded relational information relevant for each platform. On Twitter and Facebook, we commenced by identifying retweets and shares, respectively, encoding both a link between the original and the re-published contribution and a link between users. Next, we recorded which comments responded to which original tweet, Facebook post, or preceding comment, linking both the relevant texts and their respective authors. Furthermore, we extracted all @-mentions (on Twitter) as well as Facebook's corresponding ways of tagging referred-to users. On WhatsApp, we treated each contribution in an ongoing interaction as a reply to the preceding contribution. In addition, we recorded a link whenever a contribution quoted

an earlier post. We thus identified 261,786 unique dyads: 124,165 on Facebook, 132,226 on Twitter, 5395 on WhatsApp. Collating all unique user accounts, we identified 64,638 users participating in the debate. On Facebook, 53,873 users contributed an average of 2.8 posts per user; 10,536 Twitter users contributed on average 8.6 tweets per user; and 229 WhatsApp users posted an average of 27.3 messages per user.

## Content-based Measures

To classify political leanings and attitudes expressed in each contribution, we applied an innovative method that combines computational text analysis with manual content analysis (Hybrid Content Analysis, Baden et al., 2019; see Appendix for additional information on the main preprocessing and modeling steps). In a first step, we used a structural topic modeling algorithm (*stm*; Roberts et al., 2019)[7] to extract recurrent patterns of token co-occurrences within our corpus. Each contribution was modeled as comprising a small number of topics characterized by a maximally distinct set of words, including concatenated polygrams (e.g., "Israeli_Defense_Force") and non-word expressions (e.g., emojis). From a range of estimated models ($k$ = 30, 40, …, 150, 200, … , 500), a model with $k$ = 100 topics was identified as providing the best fit based on both semantic coherence and exclusivity indicators (using stm's *searchK* function), and confirmed by manual validation. The comparatively high number of topics was chosen to ensure that distinct uses of similar arguments could be differentiated; some redundancy among related topics was considered unproblematic, as we were not interested in labeling or analyzing the topics as an outcome in their own right, but used these only as an intermediary stage toward the classification of expressed contents. In a second step, we then classified each topic as supportive, opposed, ambivalent, or neutral/irrelevant toward Elor Azaria, relying on a manual content analysis of associated words and documents.[8] In the same way, we also manually coded which topics conveyed positive, negative, mixed, or no particular sentiment, as well as the intensity of emotional expression (more details regarding the coding and a list of topics can be found in Table A1 in the appendix, as well as information about the distribution of each topic in each phase presented in Figure 2).[9] Krippendorff's alpha indicated high intercoder reliability ($M$ = 0.88, $SD$ = 0.10, range: 0.74–1.00). In a third step, we used the manually classified topics to score every original post on every variable based on its inclusion of the coded topics. For the present analysis, we considered a category to be present in a document when its combined weight over all topics included in the document exceeded 0.5. This classification was validated against a manual classification of 200 randomly chosen posts, using the same codebook, showing high precision ($M$ = 0.89, range: 0.64–1.00) and recall ($M$ = 0.89, range: 0.77–1.00; for robustness checks and additional details, please refer to the Appendix as well as our methodological paper, Baden et al., 2019).

*Position.*   To determine the overall leaning of users and posts, we constructed a simple index. Given that an overwhelming majority of posts were supportive of Azaria, we distinguished supportive contributions from contributions that expressed either opposition or ambivalence. This strategy is validated by the comparatively aggressive response among Azaria's supporters against any kind of critique, which also targeted balanced and ambivalent contributions as unacceptable. We defined each user's position as the average position expressed across all contributions posted by that user within each phase, yielding an index that varied between +1

(only supportive posts) and −1 (only ambivalent or oppositional posts). Position extremeness was defined as the absolute value of that index: While we do not have a direct measure of position extremeness, our classification of topics recognizes the presence of expressions that indicate one of two opposing positions, which are by their nature irreconcilable. By averaging the coverage of uniquely pro-Azaria topics minus the coverage of the uniquely contra-Azaria topics over all of a user's posts, we obtain a measure that approaches 1 (i.e., nearly all words in a user's posts refer to one of the opposing polar positions) and 0 (either if the shares accounted for opposing positions are perfectly balanced, or if none of the words indicate a particular position). For instance, if a post consists entirely of tokens that declare Azaria a hero, demand his release and attack his detractors (pro-Azaria topics), the position is considered extreme; if we replace one topic with neutral references (e.g., to the trial), or include contravening topics instead (e.g., an admission that the killing was wrong), extremeness declines.

*Sentiment.* Based on the classification of sentiment expressed in each contribution, we created an index using +1 for positive sentiment, −1 for negative sentiment, and 0 otherwise. We then computed the average sentiment expressed by each user toward the other user that she or he interacted with, within each phase. Unlike the highly specific measurement of relational attitudes, which mostly address participants' perceived social distance or evaluation of competing elites (Druckman & Levenduskuy, 2019), our measure focuses more broadly on the typically less intense sentiment expressed in contributions addressed toward their fellow citizens supporting one or another cause. While the text-based sentiment analysis is unable to distinguish between specific relational attitudes, it is widely established as a suitable means for gauging evaluations of a focal actor in political discourse (Young & Soroka, 2012).

### Network-based Measures

For the analysis of interaction patterns, we classified all users as either supporters or opponents of Azaria based on their average position expressed toward the shooter within each phase. Contrary to the general population, where positions on the issue were roughly split, a clear majority of users on all platforms were classified as supporters of Azaria based on this index: 66.6% on Facebook (25.5% opponents), 73.7% on Twitter (20% opponents), and 79% on WhatsApp (19.7% opponents). The remaining users posted no content expressing an explicit stance and were henceforth excluded.

*Homophily/Heterophily.* Homophilic interactions were defined as all interactions between members of the same camp. Heterophilic interactions included all interactions between members of different camps. We disregarded interactions with users who never posted any comment that permitted a classification into either camp. Interactions were aggregated by type for each platform, yielding both an overall count and separate counts per phase.

### Analysis

Our analysis proceeded in three stages. First, we determined whether homophilic inter-actions were indeed more common than expected by chance (H1). For that purpose, we determined the total number of possible dyads between all users participating in the debate and used the known distribution of supporters and opponents to compute the

share of within-camp interactions that could be expected if interactions were random. We then compared these expected shares to the observed distribution of homophilic and heterophilic interactions and assessed whether there was an increase in the ratio of homophilic interactions over time (H2). We report the over- and under-representation of homophilic ties as odds (for instance, odds of 2:1 indicate that a type was twice as common as expected).

Second, we used the distribution of homophilic and heterophilic interactions involving each user in one phase to predict the extremeness of that user's position toward Azaria in the subsequent phase (H3).
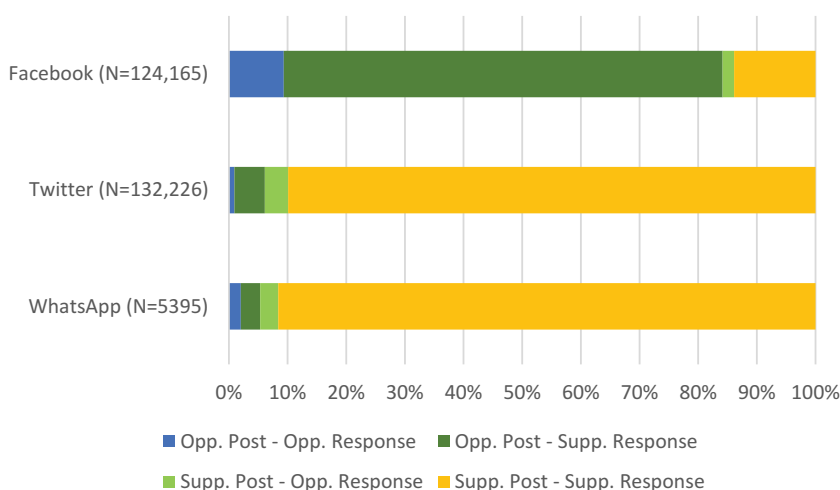
Third, we determined the average sentiment expressed within each kind of interaction in order to investigate whether contributions tended to express more positive sentiment in homophilic interactions compared to cross-cutting interactions (H4).

In view of those platform-specific differences expected by H1a, H2a, H3a, and H4a, we conducted each analysis separately for each platform. Those findings that concern hypothesized differences between platforms are reported for each stage of the analysis.

## Findings

Owing to the large share of Azaria supporters, homophilic ties were generally more likely than heterophilic ones. Despite this, and contrary to H1, homophily was not uniformly dominant. As can be seen in Figure 1, interactions among supporters (yellow) comprise the majority of dyads only on Twitter (89.9%) and WhatsApp (91.6%). On Facebook, 74.8% of all dyads represent supporters responding to contributions expressing opposition to Azaria (dark green), and only 13.9% concern interactions among supporters.

Comparing the observed prevalence of interactions to the levels expected based on chance, we found a pronounced and significant tendency toward homophily on Twitter and WhatsApp. Cross-cutting interactions were much less likely than expected on Twitter (odds: 1:3.6) and WhatsApp (odds: 1:5.3), while interactions within camps were notably more likely (odds: 1.4:1 on both platforms). On Facebook, by contrast, heterophilic



**Figure 1.** Distribution of observed interactions (individual contributions).

interactions were significantly over-represented (1.8:1) and homophilic ones were comparatively rare (1:2.5). Therefore, contrary to H0a, Facebook is characterized by the lowest (rather than highest) tendency toward homophily among the three platforms.

With regard to the hypothesized increase in interactional polarization over time (H2), measured as increasing homophily, the data offer partial support. In line with H1a, homophily rose slightly on Twitter from Phase I (homophily 1.2:1, heterophily 1:2.4) to Phase II (homophily 1.3:1, heterophily 1:3.4), but remained largely stable thereafter. An ANOVA confirms that the increased share of homophilic interactions is significant ($F(3,14523) = 3.09^{***}$ for inbound interactions, $F(3,6702) = 5.94^{***}$ for outbound interactions). On WhatsApp, interactional polarization not only did not increase, but actually decreased over time: The odds of homophilic interactions declined steadily, approaching chance level toward the end (Phase I: 1.3:1; Phase IV: 1:1). Cross-cutting interactions were least likely during the trial (1:6.2) and increased to mild underrepresentation (Phase IV: 1:1.5). While the ANOVA confirms the significant drop in the share of homophilic interactions ($F(3,372) = 3.47^{**}$ inbound, $F(3,333) = 3.49^{**}$ outbound), only the contrast between phases II and IV is significant. On Facebook, contrary to H2a, the overall tendency toward heterophily intensified over time. Interactions within camp decreased from 1:1.4 in Phase I to 1:5.0 in Phase IV, while cross-cutting interactions were overrepresented by a factor of 1.4 in Phase I, and 3.8 in Phase IV. The ANOVA confirms the significant decrease in homophily ($F(3,55083) = 1025.58^{***}$ inbound, $F(3,11388) = 17.11^{***}$ outbound).[10] Between-platform differences are significant between Twitter and WhatsApp on the one hand, and Facebook on the other, while Twitter and WhatsApp differ significantly only in Phase IV. Only the dynamics found on Twitter show a steady increase in homophily, consistent with H1, while homophily decreased over time on WhatsApp and Facebook.

Considering the effect of homophilic interactions on expressed positions, our data confirm users' tendency to express more extreme views if interactions with likeminded users take in a larger share of their social media communications (H3). At the same time, there were some noted differences between platforms, which can be seen in Table 2. On both WhatsApp and Twitter, users' positions were most strongly predicted by the

**Table 2.** Prediction of position extremeness based on homophilic interaction patterns in the previous phase.

| | Twitter N = 10,944 | | | Facebook N = 48,804 | | | WhatsApp N = 228 | | |
|---|---|---|---|---|---|---|---|---|---|
| | B | SE | Beta | B | SE | Beta | B | SE | Beta |
| *Ratio* of homophilic interactions (incoming) | .217 | .013 | .350*** | −.042 | .060 | −.055 | .257 | .061 | .407*** |
| *Ratio* of homophilic interactions (outgoing) | .134 | .012 | .223*** | .307 | .071 | .331*** | .111 | .060 | .176* |
| *N* of homophilic interactions (incoming) | .000 | .000 | −.081*** | .000 | .004 | .006 | −.003 | .006 | −.247 |
| *N* of homophilic interactions (outgoing) | −.000 | .000 | −.011 | .000 | .000 | −.042 | .001 | .006 | .123 |
| *N* of heterophilic interactions (incoming) | .002 | .001 | .042*** | .000 | .002 | .004 | .022 | .013 | .226 |
| *N* of heterophilic interactions (outgoing) | −.001 | .001 | −.015 | .000 | .001 | −.019 | −.022 | .014 | −.217 |
| **Adj. $R^2$** | | **.276** | | | **.075** | | | **.291** | |

We also checked for the reverse causal direction, predicting homophilic interaction patterns based on stance extremeness in the preceding phase. On Twitter, we find a small increase of the ratio of homophilic outgoing interactions ($\beta = 0.094^{***}$), but no effect on incoming interactions. On Facebook, we find a small negative effect on the ratio of homophilic incoming interactions ($\beta = -.036^{***}$) and a sizable increase of the ratio of homophilic outgoing interactions ($\beta = 0.272^{***}$). On WhatsApp, there are no significant effects on the ratio of homophilic interactions. On all platforms, the effects of stance extremity on the homophily of interactions in the subsequent phase is weaker than the effect of homophily of interactions on stance extremity.

proportion of incoming homophilic interactions – that is, feedback received from like-minded users – and to a lesser extent by users' own engagement with like-minded others (outgoing homophilic interactions). By contrast, incoming ties were irrelevant on Facebook, while users' own responses to likeminded individuals were associated with more extreme positions. Position extremeness depended solely on the share of homophilic and heterophilic interactions on Facebook and WhatsApp. On Twitter, the absolute number of incoming responses weakly predicted position extremeness, with the opposite sign. For a constant ratio of homophilic to heterophilic feedback received, effects diminish with an increasing volume of inbound tweets – lending some support to the expected weaker influence of feedback on Twitter (H3a). In a joint model (not shown), the platform does not show an independent main effect, but the interactions between platform and both ratios are significant predictors, underscoring the different roles played by outgoing and incoming interactions on the respective platforms. Contrary to H3a, however, the feedback was not more potent on Facebook, but rather on WhatsApp and, at least for small numbers of responses received, on Twitter.

Regarding the sentiment expressed as part of the interactions, negative sentiment was generally dominant, even in interactions within the same camp. Following H4, we expected strong negative sentiment in the heterophilic interactions (black triangles in Figure 2), and less negative sentiment in the homophilic interactions (white circles). On Twitter and WhatsApp, this pattern largely holds: Interactions among opponents were significantly less negative than either kind of cross-cutting interaction. On Twitter, but not on WhatsApp, also interactions among supporters were less negative than opponents' responses to supporters' posts. The last contrast is non-significant: Supporters did not express different sentiments in their responses to opponents or fellow supporters. On Facebook, in line with expectations, interactions among supporters were significantly less negative than both supporters' responses to opponents, and opponents' responses to supporters. However, contrary to expectations, supporters' comments posted in response to opponents' posts on Facebook were significantly less negative than interactions among opponents. A factorial ANOVA confirms significant differences in the observed interactions between platforms, the interactions between platform and user groups accounting for most variance by far (Table A2 in the Appendix). In line with H4a, negativity is significantly more pronounced in cross-cutting interactions on Twitter ($M = -0.464$, $SE = 0.005$) than on WhatsApp ($M = -0.390$, $SE = 0.029$) – with Facebook ($M = -0.432$, $SE = 0.005$), contrary to H4a, located between the other two platforms.
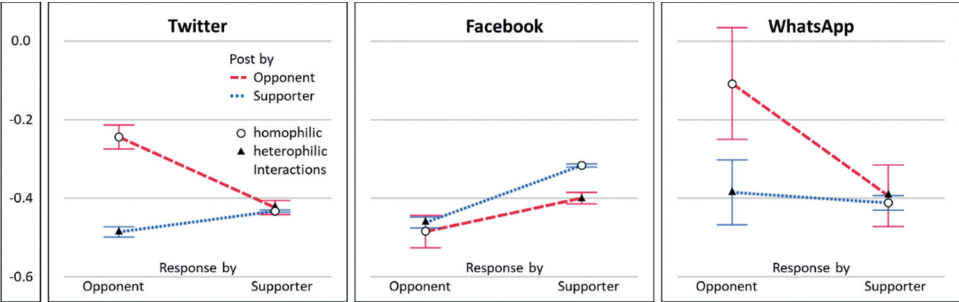


**Figure 2.** Sentiment expressed in homophilic/heterophilic interactions.

The described patterns point to a possible interaction between different types of expressed sentiment. One part of the expressed sentiment is generally consistent with expectations: Users who disagree with one another express more negative sentiment toward one another than those who agree. However, the overall high level of negativity suggests that relational attitudes account for only part of the expressed sentiment and are overlaid by a general dissatisfaction with the situation. The comparatively positive responses posted by Azaria's supporters on Facebook, and especially the unexpectedly positive sentiment found in their responses to his opponents, may reflect a specific way of expressing disagreement that arises from the particular configuration of the controversy. Besides meeting disagreement with negative sentiment, Azaria's supporters may also express their disagreement with his opponents by praising the shooter, thus indirectly challenging other users' critique. If this supposition is correct, it would appear that sentiment arising from agreement or disagreement is more pronounced on Twitter; The sentiment attached to competing definitions of Azaria's actions dominates on Facebook, and evidence of both patterns is found on WhatsApp.

## Discussion and Conclusions

Following the presented significant differences in the polarization patterns found on different social media platforms, we can no longer speak about political polarization on social media as a unified phenomenon. Only on Twitter do our data confirm the presence of all three aspects of polarization hypothesized based on the existing scholarship – which also relies heavily on Twitter (e.g., Colleoni et al., 2014; Hong & Kim, 2016). Self-reinforcing homophilic interaction patterns (interactional polarization; Colleoni et al., 2014; Mutz, 2006; Sunstein, 2001) aggravate positional polarization (Abramowitz, 2015; Fiorina & Abrams, 2008) and contribute to pronounced hostility (affective polarization; Iyengar et al., 2012; Lelkes, 2016) in cross-cutting interactions. Twitter's public nature enables users with vastly different views to interact without the constraint of shared purposes or identities, resulting in heated controversy and growing polarization.

Yet, our WhatsApp data suggest that this is not a universal pattern. Despite the heterogeneous composition of the groups examined, the constraint of shared group identity, mutual respect, and common purpose (Kligler-Vilenchik, 2019) effectively counteracted the escalation dynamics found on Twitter. Toward the end of the controversy, the initially homophilic interaction patterns on WhatsApp had largely given way to an almost unbiased interaction network, and hostility toward discrepant viewpoints had decreased notably. While this pattern is not solely an outcome of platform structure, but reflects also the specific nature and composition of those groups analyzed, our analysis shows that initially homophilic interactions and negativity toward outgroups do not necessarily result in polarization, but can also gradually de-polarize.

The deviant behavior of interactions on Facebook is of particular concern for two reasons. First, our analysis underscores an important but rarely recognized ambiguity in the scholarly conceptualization of social media platforms. While most theorizing on tie formation and interactions on Facebook focuses on the use of the platform for communication among private users, much of the existing empirical work – including our own – relies (mostly for reasons of data access) on publicly open discussions, which function

somewhat differently (e.g., Ben-David & Matamoros-Fernandez, 2016; Yarchi & Samuel-Azran, 2018). Similar to Twitter, users can access and respond to public contents also if they share no personal ties with the author, mitigating the platform's hemophilic tendencies. Moreover, publicly open discussions on Facebook in general, and public pages in particular often exist as instruments of political activism (Bossetta, 2018; Harlow & Harp, 2012) – and other users respond to it as such. To make sense of the strong prevalence of cross-cutting communication, it may be appropriate to think of such posts not as opinion statements addressed to friends (e.g., Green et al., 2016; Lönnqvist & Itkonen, 2016), but rather as public statements by protagonists of social activism (Hendriks et al., 2016). In our case, there is reason to believe that at least part of the cross-cutting interactions on Facebook – predominantly by supporters of Azaria, addressing posts opposing him – were coordinated activist attacks (Livio & Afriat, 2019), intended to shout down and silence opponents in the digital realm. This aggressive response may also explain the disproportionately small number of posts and pages opposing Azaria in the Israeli social media public sphere, initiating a spiral of silence (Matthes et al., 2010) that dissuaded all but the most committed and seasoned critics (notably, activists, politicians, and journalists) from stating their views in public. Inversely, the dominance of Azaria's supporters on Facebook should have emboldened additional right-wing users join the attacks (Ben-David & Matamoros-Fernandez, 2016), explaining the growing heterophily of interactions. Over the course of the debate, many users dropped from the public Facebook debate, possibly withdrawing into more sheltered private Facebook conversations that neither their attackers, nor we, could access. Given this structural bias of public content-focused Facebook research, we need to reconsider not only how this selective visibility biases our formation of scientific knowledge, but also how existing theorizing can be brought into synchronicity with the important differences between public and private Facebook.

Second, our findings challenge several assumptions about the nature of cross-cutting political talk. Most existing work on cross-cutting political talk has focused on interactions supporting mutual understanding and tolerance (Barberá, 2014; Wojcieszak, 2010). By contrast, the patterns captured in our data point to a form of strategic cross-cutting talk that is oriented toward conflict and intimidation rather than an exchange of arguments. Furthermore, contrary to common assumptions, this hostile practice did not primarily rely on negative sentiment expressed toward opponents and the controversy (Hillmann & Trier, 2012), but invoked strong positive sentiment toward Azaria in order to mark the underlying conflict. Following this line of thought, even heated and hostile controversies might be dominated by positive sentiment, as each side praises its opposing heroes. These findings cast serious doubts both upon the use of sentiment as an indicator of hostility, and of the prevalent conceptualization of cross-cutting talk as a driver of de-polarization.

Accordingly, the study reveals a range of important theoretical and methodological blind spots in the existing scholarship on polarization in social media. Theoretically, we have underscored the urgent need to distinguish different constructs associated with, but not identical to, political polarization, which may be affected by the structural features and affordances offered by different social media environments in different ways. Beyond questioning the widespread reliance on Twitter (and limited public Facebook) data to draw conclusions about social media as a whole, our study also highlights the perils of inferring dynamic properties from static data. Our investigation furthermore documents

the need to integrate rich data addressing different aspects of interactive, dynamic social media debates, and illustrates how computational communication science approaches can be instrumental in overcoming some of those challenges. By using a novel approach to algorithmic text analysis, which combines inductive pattern extraction with deductive classification, we have been able to infer political leanings not from (incomplete and crude) metadata or (possibly misleading) follow-patterns, but from the actual contents of the ongoing debate. This hybrid strategy enables us to treat the ever-changing and diverse discourse of social media in a breadth far beyond the reach of researcher-constructed dictionaries, while asserting far closer deductive control over the classification of contents than afforded by traditional, inductive uses of topic modeling. At the same time, the integration of textual analysis and relational data enabled us to proceed beyond the classification raised positions and contexts (e.g., Bodrunova et al., 2019; Garcia et al., 2012), relying on content-classified interactional dyads between users with known expressed positions to reconstruct the complex interactional structure of social media discourse. Thus, reconstructing the temporal sequence of contributions, we were able to study the inherent dynamics of the ongoing debate. While our segmentation of data into phases may be somewhat crude, the detailed timestamps of collected data permit a disaggregation all the way down to individual interaction sequences without compromising the capacity to aggregate and analyze detected patterns. We do all this relying on existing tools and technologies, adding to the methodological toolbox chiefly a capacity to integrate different kinds of social media data at scale – a capacity that we gain by linking up the computational extraction of patterns with their manual, theory-guided classification. With this simple twist, and enabled by recent developments in computational analysis, we were able to capture, model, and analyze the interactive, rich, and dynamic evolution of polarization in its natural settings on various social media platforms and over time.

Of course, our study is subject to several limitations. First, as is evidenced by the disproportionate overrepresentation of one of the camps, social media debates are in no way representative of offline social dynamics. The observed polarization dynamics manifested on social media may play their part in influencing actual social polarization, but chiefly arise from users' selective choice to post or withhold their views in the digital public sphere (Lelkes, 2016). Second, especially our analysis of Facebook contents is limited by the reliance on publicly accessible contributions, as mandated by the platform's privacy regulations, which likely differ from what happens in the more personal realm of privacy-protected Facebook discussions. Similarly, further study is needed to gain a more representative understanding of polarization on WhatsApp, as our analysis permits no generalization toward the platform, or even WhatsApp groups, as a whole. Analytically, our classification of textual contents – especially expressed sentiments – remains crude and needs to be further developed. As discussed above, textual sentiment provides little information about the target of the expressed affect, which may include the addressed user, but also the discussed issue and situation, or third entities referenced in the conversation. However, our attempts to refine our classification of expressed sentiment turned out unsuccessful owing to the considerable ambiguity of natural discourse. Even when read manually and within their interactive content, many posts permitted multiple readings regarding the specific target of expressed sentiment. Our measure is likewise

unable to distinguish between different kinds of relational attitudes established in the predominantly survey-based polarization literature, owing to the ambiguous and under-specified nature of natural discourse. Other limitations concern our measurement of evaluative stance, whose reliance on the distribution of pro- and anti-Azaria considerations ignores many important qualitative differences in how these considerations are expressed. While we made every effort to validate our scales, these measures depend on contributions' narrow focus on one controversial issue; further development will be necessary to generalize this strategy toward more complex controversies that span multiple issues and polar positions.

In our analysis, we have focused on distinguishing between interactional, positional, and affective polarization in a broad sense and investigated their interplay in different social media environments. Despite all remaining limitations, we believe that both our conceptual distinctions and our use of computational tools offer many new and valuable opportunities for political communication scholars to examine complex theoretical questions, at scale, across platforms and over time. In this endeavor, our study takes but one more step toward a better understanding of online interaction dynamics and political polarization on social media. Future studies will need to examine polarization in different political contexts, over different issues, on different platforms, and in different time horizons, adding detail and refining our theoretical knowledge.

## Notes

1. Poll conducted by the Israeli Democracy Institute. Retrieved from: http://www.peaceindex.org/indexMonth.aspx?num=308.
2. Retrieved from: https://www.mako.co.il/news-israel/local-q1_2017/Article-1a4ebc032ca6951004.htm.
3. Due to Facebook's privacy regulations, we were able to retrieve only the contents of public pages and posts set to public.
4. https://www.socialbakers.com/statistics/facebook/pages/total/israel.
5. Due to self-selection, participants in these groups are highly interested in politics, and not representative of WhatsApp users as a whole. Still, our findings provide rare insights into the political communication happening on this understudied platform.
6. The Facebook and Twitter data is presented on the journal's website (Due to privacy issues, the WhatsApp data cannot be published).
7. stm was chosen for its capacity to include covariates. Estimating one model across all platforms, this allowed us to permit some topics to be more or less prevalent on different platforms.
8. In a manual content analysis, we classified those concepts and lexical choices grouped by each topic based on their fit within a) those narratives promoted by Azaria's supporters, which presented his actions as justified and heroic in the defense against violent terrorism, (e.g., hero, terrorist, security, justified); b) those narratives advanced by those opposing Azaria, which presented his actions as immoral and denounced him as a hateful murderer (e.g., extrajudicial, murder, human rights, hateful), c) ambivalent narratives that combined considerations of both sides or were compatible with different evaluations (e.g., kill, justice, intentions), or d) neutral descriptors that gave no indication of an evaluative stance. See Appendix for detailed documentation.
9. For this classification, we decided whether a majority of top-associated words directly expressed emotions or carried strong emotional sentiment, and if so, whether negative (e.g., hate, corrupted, traitors, hypocrisy), positive (e.g., hero, solidarity, love, everybody's child), or mixed emotions were expressed.

10. We also checked the overtime changes in a within-subject repeated measures ANOVA for those participants who contributed to all four phases, confirming the significant increase in homophily on Twitter ($F(3,398) = 3.43^{**}$ inbound, $F(3,193) = 6.76^{***}$ outbound) and the significant decrease in homophily on Facebook ($F(3,535) = 56.89^{***}$ inbound; n.s. for outbound interactions). For WhatsApp, the repeated measures ANOVAs are nonsignificant, owing to the very small number of users present in all phases.

## Disclosure Statement

No potential conflict of interest was reported by the authors.

## Notes on contributors

*Moran Yarchi* (Ph.D. Hebrew University of Jerusalem) is a Senior Lecturer at the Sammy Ofer School of Communications, the Head of the Public Diplomacy program, and a Senior Researcher at the Institute for Counter-Terrorism (ICT) at the Interdisciplinary Center (IDC) Herzliya, Israel. Her main area of research is political communication, especially the media's coverage of conflicts and public diplomacy.

*Christian Baden* (Ph.D., University of Amsterdam) is a Senior Lecturer in the Department of Communication and Journalism at the Hebrew University of Jerusalem. His research focuses on the collaborative construction of meaning in controversial public debates in political communication and journalism, and has advanced theory and methodology in the study of dynamic political discourse.

*Neta Kligler-Vilenchik* (Ph.D., University of Southern California) is a Senior Lecturer at the Hebrew University of Jerusalem. Her research interests include political expression in the new media environment.

## ORCID

Moran Yarchi 🔴 http://orcid.org/0000-0002-8044-2145

## References

Abramowitz, A. I. (2015). The new American electorate: Partisan, sorted and polarized. In J. A. Thurber & A. Yoshinaka (Eds.), *American gridlock: The sources, character, and impact of political polarization* (pp. 19–44). CUP.

Aronoff, M. J. (1984). *Political polarization: Contradictory interpretations of Israeli reality. Cross-currents in Israeli culture and politics*. Transaction.

Baden, C., Kligler-Vilenchik, N., & Yarchi, M. (2019, May). *Hybrid content analysis: Toward a strategy for the computer-assisted classification of large text corpora using topic modeling* [Paper presentation]. 69th ICA Annual Conference, Washington, DC.

Bakshy, E., Messing, S., & Adamic, L. A. (2015). Exposure to ideologically diverse news and opinion on Facebook. *Science*, *348*(6239), 1130–1132. https://doi.org/10.1126/science.aaa1160

Barberá, P., (2014). *How social media reduces mass political polarization. Evidence from Germany, Spain, and the U.S.* Working paper. Retrieved September 18, 2019, from http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.658.5476

Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A., & Bonneau, R. (2015). Tweeting from left to right: Is online political communication more than an echo chamber? *Psychological Science*, *26*(10), 1531–1542. https://doi.org/10.1177/0956797615594620

Baylis, T. A. (2012). Elite consensus and political polarization: Cases from Central Europe. *Historical Social Research*, *37*(1), 90–106. https://www.jstor.org/stable/41756452?seq=1

Ben-David, A., & Matamoros-Fernandez, A. (2016). Hate speech and covert discrimination on social media: Monitoring the Facebook pages of extreme-right political parties in Spain. *International Journal of Communication*, *10*, 1167–1193. https://ijoc.org/index.php/ijoc/article/view/3697

Bobok, D. (2016). *Selective exposure, filter bubbles and echo chambers on Facebook*. Central European University.

Boczkowski, P. J., Matassi, M., & Mitchelstein, E. (2018). How young users deal with multiple platforms: The role of meaning-making in social media repertoires. *Journal of Computer-Mediated Communication*, *23*(5), 245–259. https://doi.org/10.1093/jcmc/zmy012

Bodrunova, S. S., Blekanov, I., Smoliarova, A., & Litvinenko, A. (2019). Beyond left and right: Real-world political polarization in Twitter discussions on inter-ethnic conflicts. *Media and Communication*, *7*(3), 119. https://doi.org/10.17645/mac.v7i3.1934

Bond, R. M., Fariss, C. J., Jones, J. J., Kramer, A. D. I., Marlow, C., Settle, J. E., & Fowler, J. H. (2012). A 61-million-person experiment in social influence and political mobilization. *Nature*, *489*(7415), 295–298. https://doi.org/10.1038/nature11421

Bossetta, M. (2018). The digital architectures of social media: Comparing political campaigning on Facebook, Twitter, Instagram, and Snapchat in the 2016 U.S. Election. *Journalism & Mass Communication Quarterly*, *95*(2), 471–496. https://doi.org/10.1177/1077699018763307

Brundidge, J. (2010). Encountering "Difference" in the contemporary public sphere: The contribution of the Internet to the heterogeneity of political discussion networks. *Journal of Communication*, *60*(4), 680–700. https://doi.org/10.1111/j.1460-2466.2010.01509.x

Bruns, A. (2019). *Are filter bubbles real?* Polity.

Colleoni, E., Rozza, A., & Arvidsson, A. (2014). Echo chamber or public sphere? Predicting political orientation and measuring political homophily in Twitter using big data. *Journal of Communication*, *64*(2), 317–332. https://doi.org/10.1111/jcom.12084

Druckman, J. N., & Levenduskuy, M. S. (2019). What do we measure when we measure affective polarization? *Public Opinion Quarterly*, *83*(1), 114–122. https://doi.org/10.1093/poq/nfz003

Dvir-Gvirsman, S., Tzfati, Y., & Menchen-Trevino, E. (2016). The extent and nature of ideological selective exposure online: Combining survey responses with actual web log data from the 2013 Israeli Elections. *New Media & Society*, *18*(5), 857–877. https://doi.org/10.1177/1461444814549041

Ellison, N. B., & Boyd, D. M. (2013). Sociality through social network sites. In W. H. Dutton (Ed.), *The Oxford handbook of Internet studies* (pp. 151–172). OUP.

Fiorina, M. P., & Abrams, S. J. (2008). Political polarization in the American public. *Annual Review of Political Science*, *11*(1), 563–588. https://doi.org/10.1146/annurev.polisci.11.053106.153836

Garcia, D., Mendez, F., Serdült, U., & Schweitzer, F. (2012). *Political polarization and popularity in online participatory media: An integrated approach*. 1st workshop on politics, elections and data, Maui, HI.

Goldberg, Y., & Elhadad, M. (2013). Word segmentation, unknown-word resolution, and morphological agreement in a Hebrew parsing system. *Computational Linguistics*, *39*(1), 121–160. https://doi.org/10.1162/COLI_a_00137

Green, T., Wilhelmsen, T., Wilmots, E., Dodd, B., & Quinn, S. (2016). Social anxiety, attributes of online communication and self-disclosure across private and public Facebook communication. *Computers in Human Behavior*, *58*, 206–213. https://doi.org/10.1016/j.chb.2015.12.066

Harlow, S., & Harp, D. (2012). Collective action on the web. *Information Communication & Society*, *15*(2), 196–216. https://doi.org/10.1080/1369118X.2011.591411

Hayat, T., & Samuel-Azran, T. (2017). "You too, second screeners?" Second screeners' echo chambers during the 2016 U.S. Elections primaries. *Journal of Broadcasting & Electronic Media*, *61*(2), 291–308. https://doi.org/10.1080/08838151.2017.1309417

Hendriks, C. M., Duus, S., & Ercan, S. A. (2016). Performing politics on social media: The dramaturgy of an environmental controversy on Facebook. *Environmental Politics*, *26*(6), 1102–1125. https://doi.org/10.1080/09644016.2016.1196967

Hillmann, R., & Trier, M. (2012). Sentiment polarization and balance among users in online social networks. *AMCIS 2012 proceedings*.

Holt, R. (2004). *Dialogue on the Internet: Language, civic identity, and computer–mediated communication*. Praeger.

Hong, S., & Kim, S. H. (2016). Political polarization on Twitter: Implications for the use of social media in digital governments. *Government Information Quarterly*, 33(4), 777–782. https://doi.org/10.1016/j.giq.2016.04.007

Huckfeldt, R., Mendez, J. M., & Osborn, T. (2004). Disagreement, ambivalence, and engagement: The political consequences of heterogeneous networks. *Political Psychology*, 25(1), 65–95. https://doi.org/10.1111/j.1467-9221.2004.00357.x

Huckfeldt, R., & Sprague, J. (1995). *Citizens, politics and social communication: Information and influence in an election campaign*. CUP.

Iyengar, S., Sood, G., & Lelkes, Y. (2012). Affect, not ideology: A social identity perspective on polarization. *Public Opinion Quarterly*, 76(3), 405–431. https://doi.org/10.1093/poq/nfs038

John, N. A., & Dvir-Gvirsman, S. (2015). "I don't like you any more": Facebook unfriending by Israelis during the Israel–Gaza conflict of 2014. *Journal of Communication*, 65(6), 953–974. https://doi.org/10.1111/jcom.12188

Kligler-Vilenchik, N. (2019). Friendship and politics don't mix? The role of sociability for online political talk. *Information, Communication & Society*, 1–16. https://doi.org/10.1080/1369118X.2019.1635185

Kligler-Vilenchik, N., Yarchi, M., & Baden, C. (2019, May). *Interpretative polarization across platforms: How a controversial case fragmented Israeli audiences across Facebook, Twitter and WhatsApp* [Paper presentation]. 69th ICA Annual Conference, Washington, DC.

Kwak, H., Lee, C., Park, H., & Moon, S. (2010). What is Twitter, a social network or a news media? In M. Rappa (Ed.), *Proceedings of the 19th international conference on World Wide Web* (pp. 591–600). ACM.

Kwon, H. E., & Oh, W. (2014). *Platform characteristics, multi-homing, and homophily in online social networks*. Workshop on information systems and economics. Retrieved September 18, 2019, from http://misrc.umn.edu/wise/2014_Papers/4.pdf

Lee, E. (2006). When and how does depersonalization increase conformity to group norms in computer-mediated communication? *Communication Research*, 33(6), 423–447. https://doi.org/10.1177/0093650206293248

Lelkes, Y. (2016). Mass polarization: Manifestations and measurements. *Public Opinion Quarterly*, 80(1), 392–410. https://doi.org/10.1093/poq/nfw005

Lerman, K., Yan, X., & Wu, X. Z. (2016). The "majority illusion" in social networks. *PLOS One*, 11 (2), e0147617. https://doi.org/10.1371/journal.pone.0147617

Livio, O., & Afriat, H. (2019). Politicised celebrity in a conflict-ridden society: The Elor Azaria case and celebritisation discourses in Israel. *Celebrity Studies*, 1–17. https://doi.org/10.1080/19392397.2019.1609370

Lönnqvist, J. E., & Itkonen, J. V. A. (2016). Homogeneity of personal values and personality traits in Facebook social networks. *Journal of Research in Personality*, 60, 24–35. https://doi.org/10.1016/j.jrp.2015.11.001

Marichal, J. (2013). Political Facebook groups: Micro-activism and the digital front stage. *First Monday*, 18(12). https://doi.org/10.5210/fm.v18i12.4653

Matthes, J., Morrison, K. R., & Schemer, C. (2010). A spiral of silence for some: Attitude certainty and the expression of political minority opinions. *Communication Research*, 37(6), 774–800. https://doi.org/10.1177/0093650210362685

Moscovici, S., & Zavalloni, M. (1969). The group as a polarizer of attitudes. *Journal of Personality and Social Psychology*, 12(2), 125–135. https://doi.org/10.1037/h0027568

Mutz, D. (2006). *Hearing the other side: Deliberative versus participatory democracy*. CUP.

Mutz, D. C., & Martin, P. S. (2001). Facilitating communication across lines of political difference: The role of mass media. *American Political Science Review*, 95(1), 97–114. https://doi.org/10.1017/S0003055401000223
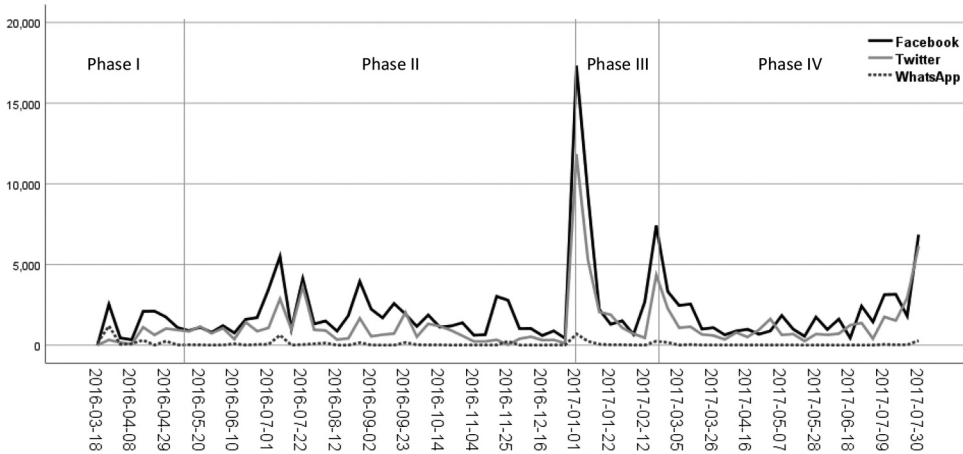
Nahon, K. (2016). Where there is social media there is politics. In A. Bruns, G. Enli, E. Skogerbo, A. O. Larsson, & C. Christensen (Eds.), *Companion to social media and politics* (pp. 39–55). Routledge.

Parsons, B. M. (2010). Social networks and the affective impact of political disagreement. *Political Behavior*, 32(2), 181–204. https://doi.org/10.1007/s11109-009-9100-6

Roberts, M. E., Steward, B. M., & Tingley, D. (2019). stm. R package for structural topic models. *Journal of Statistical Software*, 91(2). https://doi.org/10.18637/jss.v091.i02

Sagolla, D. (2009). *140 Characters: A style guide for the short form* (1st ed.). Wiley.

Settle, J. E. (2018). *Frenemies: How social media polarizes America*. CUP.

Slater, M. D. (2007). Reinforcing spirals: The mutual influence of media selectivity and media effects and their impact on individual behavior and social identity. *Communication Theory*, 17(3), 281–303. https://doi.org/10.1111/j.1468-2885.2007.00296.x

Stroud, N. J. (2010). Polarization and partisan selective exposure. *Journal of Communication*, 60(3), 556–576. https://doi.org/10.1111/j.1460-2466.2010.01497.x

Suler, J. (2004). The online disinhibition effect. *CyberPsychology & Behavior*, 7(3), 321–326. https://doi.org/10.1089/1094931041291295

Sunstein, C. R. (2001). *Echo chambers: Bus v. Gore, impeachment, and beyond*. PUP.

Tucker, J. A., Guess, A., Barberá, P., Vaccari, C., Siegel, A., Sanovich, S., Stukal, D., & Nyhan, B. (2018). *Social media, political polarization, and political disinformation: A Review of the scientific literature*. Hewlett Foundation.

Van Aelst, P., Strömbäck, J., Aalberg, T., Esser, F., de Vreese, C. H., Matthes, J., & Stanyer, J. (2017). Political communication in a high-choice media environment: A challenge for democracy? *Annals of the International Communication Association*, 41(1), 3–27. https://doi.org/10.1080/23808985.2017.1288551

Weng, J., Lim, E. P., Jiang, J., & He, Q. (2010, February). Twitterrank: Finding topic-sensitive influential twitterers. *Proceedings of the 3rd ACM international conference on web search and data mining*.

Wojcieszak, M. (2010). 'Don't talk to me': Effects of ideologically homogeneous online groups and politically dissimilar offline ties on extremism. *New Media & Society*, 12(4), 637–655. https://doi.org/10.1177/1461444809342775

Yarchi, M., & Samuel-Azran, T. (2018). Women politicians are more engaging: Male versus female politicians' ability to generate users' engagement on social media during an election campaign. *Information Communication & Society*, 21(7), 978–995. https://doi.org/10.1080/1369118X.2018.1439985

Young, L., & Soroka, S. (2012). Affective news: The automated coding of sentiment in political texts. *Political Communication*, 29(2), 205–231. https://doi.org/10.1080/10584609.2012.671234

Zuiderveen Borgesius, F. J., Trilling, D., Möller, J., Bodó, B., de Vreese, C. H., & Helberger, N. (2016). Should we worry about filter bubbles? *Internet Policy Review*, 5(1). https://doi.org/10.14763/2016.1.401

## Appendix

**Contents**
I.Descriptives: Corpus composition over time
II.Codebook Stage I: Relevance classification of WhatsApp messages
III.Topic Modeling: Preprocessing & estimation
IV.Codebook Stage II: Classification of topics
V.Topic List: Top tokens & classification
VI.Validation: Reliability, accuracy & robustness
VII.Topic Coverage
VIII.ANOVA Table
IX.Share Note

## I.DESCRIPTIVES: Corpus composition over time



**Figure A1.** Weekly frequencies of social media posts per platform.

## II.CODEBOOK STAGE I: Relevance classification of WhatsApp messages

The purpose of the coding is to identify the sections of the conversation in the WhatsApp groups that relate to the Elor Azaria affair. The dates in the file were selected because the following search words were included: Elor Azaria (including misspellings), Elor, Azaria, the soldier, the shooter, etc. However, as soon as a date is included in the corpus, we analyze every section of the discourse that relates to Elor Azaria, with or without the use of these words. There are three possible classifications for the messages:

**0 – Does not refer to Elor Azaria**
**1 – This is an opening message – the beginning of a conversation about Elor Azaria**
**2 – Contains a follow-up response to a conversation about Elor Azaria**

*How to identify an opening message* **(Code 1)**

An opening message starts a conversation about Elor Azaria, after the conversation first dealt with other topics (if this is the first message of the day and it deals with Elor Azaria, it will also be considered as an opening message). Many times this message will be identified using one of the search words (Elor Azaria, Elor, Azaria, the soldier, the shooter, etc.) but can also be in other ways (e.g., shooting in Hebron, soldier from Hebron, the affair of the soldier, etc.) or regarding events related to the Elor Azaria case (trial, broadcasting a program about the event, etc.).

*How to identify a follow-up response* **(Code 2)**

A follow-up response is any response that continues the same conversation until it moves to another topic. This includes any responses that relate to Azaria, as well as responses that simply continue the conversation ("totally," "right," "wrong," "no way", and so on). A conversation can go on for many messages, or it can include only a few messages.

As long as the conversation goes on until another topic has been raised, all messages within that conversation will be considered follow-up responses, even if they are not directly about Azaria, or if it is not entirely clear what they are about. It also includes links, emojis, or blank messages.

*How to identify when the conversation moves to another topic* **(Code 0)**

When there are more than 3 consecutive messages that do not deal with Azaria, but are dealing with another topic (which, in principle, could have been given a headline. For example, corruption, investigations about the prime minister, or social discussions between participants), those messages will be defined as dealing with a different topic – from the first post that no longer deals with Azaria, and until there is a new "opening message" on Azaria. Message timing can also be used – when a response comes along after previous messages, it is more likely (but not always) a new topic. Empty messages are sometimes evidence that a participant has uploaded a screenshot or picture that we do not see, which may be the beginning of another topic.

### III. TOPIC MODELING: Preprocessing & estimation
*Preprocessing*

(1) *Deduplication*: All duplicate messages (retweets, shared messages, etc.) were removed.
(2) *Non-word text*: We identified all emojis used in the text, including emojis written as character sequences (e.g., "(:")). All unique emojis were included as separate unique tokens. We also harmonized all onomatopoetic expressions to the first three characters (e.g., "aah[hhhh!]").
(3) *Clean-up*: All embedded links, tags, and other non-textual contents were removed. All remaining punctuation was removed.
(4) *Tokenization*: As a morphologically rich language, Hebrew expresses various grammatical functions and conjunctions by means of prefixes and suffixes attached to words. To separate such prefixes and suffices and thus obtain tokens that correspond to separate words, we used Goldberg and Elhadad's (2013) Hebrew tokenizer as well as the Mila word bank. All construct state suffixes were reverted to the standard form.
(5) *Acronyms*. We identified all acronyms in the text and replaced them by the proper name.
(6) *Stop-word removal*. We adapted Mila's stop-word list to remove uninformative tokens.
(7) *Concatenation*. We obtained word frequency lists for uni-, bi-, and trigrams. Starting from the top, we corrected any spelling errors (notably, of Azaria's name) or variations (mostly of transliterations) and concatenated any entity names (e.g., Supreme_Court), standing expressions (e.g., Thank_God) and polygram words (e.g., "lawyer" consists of two words in Hebrew).

*Model estimation*

(1) *DFM construction*. We used the quanteda R package to construct the document feature matrix.
(2) *TF/IDF weighting*. The DFM was weighted by term frequence/inverse document frequency.
(3) *Model parameters*. Topic distribution was set to be modeled dependent on the platform. To determine $k$, we estimated structural topic models with $k = 30, 40, …, 130, 140, 150, 200, 250, … 450, 500$ and evaluated their semantic coherence and exclusiveness statistics. For the best-fitting models ($k = 80, 90, 100, 110$) we manually inspected the topics and compared topics across models to determine where unified topics in one model were merged or split in another model. The selected model for the unified topic model was k = 100, with topic prevalence dependent on the platform as covariate. We also estimated three separate models with k = 80 for Facebook, k = 80 for Twitter, and k = 70 for WhatsApp, for validation purposes (see Validation).
(4) *Representation*. For coding, each topic was represented by its top tokens based on their Probability, FREX, Lift, and Score, as well as the top four most closely associated documents.

### IV. CODEBOOK STAGE II: Classification of topics
The purpose of this codebook is to classify the topics extracted by the Structural Topic Model according to eight variables:

| | |
|---|---|
| *I – Coherence* | *V – Emotionalization* |
| *II – Topicality* | *VI – Emotive Sentiment* |
| *III – Relevance* | *VII – Thematic Focus* |
| *IV – Evaluative Stance* | *VIII – Actor Focus* |

To classify each topic, we proceed as follows. Prior to classification, we examine the top 10 tokens associated with each topic based on their probability (Highest Probability), as well as three additional weighting methods (FREX, Lift, Score; see Roberts et al., 2019). For each topic, we first attempted to decide what shared logic is responsible for the top words based on the top words alone. Subsequently, we referred to the top three documents associated with the same topic to validate this interpretation. Shared logics did not have to be topical, rather, the aim was to discern what could account for the systematic co-occurrence of these tokens within the same documents.

Subsequently, the topic was coded with regard to each variable. For each coding decision, we primarily regarded those tokens listed by the Highest Probability, referring to the other lists of top tokens for disambiguation only. **All criteria specified hereunder had to be met within the list of top 10 highest probability tokens, and not be contradicted by the shared logic and the other top token lists.**

*I – Coherence.*

**1 – interpretable**                    **0 – not interpretable**

Topics were coded as interpretable if we could discern some shared logic that plausibly accounted for the top 10 tokens associated with the topic, which we validated by looking at the top three associated texts. For a topic to be coded as interpretable, there had to be at least four of the top tokens in the same list that unambiguously fit the same context of use.

*II – Topicality*

**0 – none**               **2 – stylistic/rhetorical**            **3 – social**
**1 – topical**            **4 – mixed**

Topics were coded as **topical** if the shared logic connecting the tokens could be formulated to be *about* something specific, i.e., if the topic of the tokens could be specified. We coded topics as **stylistic** or **rhetorical** if the shared logic referred to a particular way of talking about different topics (e.g., use of emojis, rhetorical devices). We coded topics as **social** if the shared logic pertained to the management of social ties (e.g., personal pronouns, communicative behavior). Topics wherein more than one of the above criteria were affirmed were classified as mixed. Topics that matched neither of these criteria were coded as **none**.

*III – Relevance*

**0 – not relevant**                    **1 – relevant**

Topics were coded as **relevant** if at least two of the top tokens in the same list explicitly referred to some pertinent aspect of the case (Elor Azaria, the shooting, the trial, or the public controversy thereover) and were part of the shared logic connecting the tokens. For tokens with double meanings, we regarded them as relevant if there was at least one unique (unambiguous) reference to the case and at least one of the possible meanings met the criterion set out above. Any topics for which this was not the case were coded as **not relevant**.

*IV – Evaluative Stance toward Elor Azaria and his deed*

**0 – neutral**                    **2 – ambivalent**
**1 – opposed**                   **3 – supportive**

A topic was coded as expressing a **supportive** stance toward Elor Azaria if it contained any tokens that directly evaluated him or his deed positively, expressed solidarity with him or his family, or referred to his acquittal or release, and there were no tokens that contradicted this evaluative stance. A topic was coded as expressing an **opposed** stance toward Elor Azaria if it contained any tokens that directly evaluated him or his deed negatively, expressed disgust or shame, or referred to his condemnation or punishment, and there were no tokens that contradicted this evaluative stance. We coded topics as **ambivalent** if they contained the same, non-zero number of tokens expressing a supportive and an opposed stance. If a topic contained an uneven balance of supportive and opposed tokens, we classified it as supportive or opposed if either set of tokens outnumbered the other by more than 100% (3:1, 5:2, …), and otherwise as ambivalent. Relevant tokens included both explicit evaluations and strongly connoted tokens, as long as it was not possible to use the same token in the opposite evaluative connotation. We also considered overt references to major slogans and chants used by Azaria's supporters and opponents in demonstrations, as long as these expressed a unique affiliation with that camp. If a topic contained no tokens that expressed a unique evaluative stance, the topic was coded as **neutral**.

*V – Emotionalization*

**0 – not emotional**                    **1 – emotional**

A topic was coded as **emotional** if it contained at least one explicit reference to emotions (e.g., hope, disgust) or strongly emotionally charged concepts (e.g., hero, terrorist). We also coded references to feeling or emotion that did not in themselves specify particular emotions (e.g., feel) if there were other tokens in the list that identified an emotional charge, and both tokens fit within

the shared logic of the topic. Topics for which these criteria were not met were coded as **not emotional**.

### VI – *Emotive Sentiment*

| | |
|---|---|
| **0 – none** | **2 – mixed** |
| **1 – negative** | **3 – positive** |

Emotive sentiment was only coded for topics coded as emotional (above); all topics coded as not emotional (above) were coded 0 – none. A topic was coded as expressing **negative/positive** sentiment if all tokens responsible for the classification as emotional expressed negative/positive sentiment, or if those tokens expressing negative/positive sentiment outnumbered those expressing the opposite sentiment by more than 100% (3:1, 5:2, …). If there was an equal number of tokens expressing positive and negative sentiment, or either side outnumbered the other by less than 100%, the topic was coded as **mixed** sentiment. If the tendency of emotive sentiment could not be uniquely identified, the topic was coded to contain **no** emotive sentiment.

### VII – *Thematic Focus*

| | |
|---|---|
| **0 – none/not applicable** | **4 – army/military affairs/security** |
| **1 – the shooting** | **5 – society/solidarity/social values** |
| **2 – media/public discourse** | **6 – legal/law enforcement** |
| **3 – politics/government** | **7 – mixed/unclear** |

To classify the thematic focus of a topic, we regarded the register and social domain referred to by the top tokens, as well as the shared logic responsible for the co-presence of tokens on this list. For any topical focus to be coded, there had to be at least two tokens that uniquely referred to that focus, and there had to be more tokens relevant to this focus than to any other focus. A topic was coded as about the shooting if its tokens referred to any aspects of the shooting (the act, the site and situation, the victim, the victim's death, etc.). A topic was coded as about the media and public discourse if its tokens referred to any aspect of public debate (media, journalists, the act of public speaking, social media platforms, published public opinion, etc.). A topic was coded as about politics or government if its tokens referred to the political treatment of the case (political actors, legislative proposals, political negotiations, policies, etc.). A topic was coded as about the army and military affairs if its tokens referred to any aspect of the military situation (the Israeli Defense Forces, defense, security threats, rules of engagement, readiness, etc.). A topic was coded as about society and social values if its tokens referred to any aspect of Azaria's role within society, or Israeli society in relation to Azaria (society as a whole, solidarity, ethics, family metaphors, black sheep metaphors, etc.). A topic was coded as about the legal situation and law enforcement if its tokens referred to any aspect of the criminal prosecution and trial (arrest, indictment, trial, sentencing, laws, criminal justice, pardon, etc.). If multiple of these foci applied, but no focus matched the criterion of two tokens and more tokens than any other focus, the topic was coded as mixed. A topic was coded to have no thematic focus if it failed to meet any of the above criteria (less than two relevant tokens referring to any focus).

### VIII – *Actor Focus*

| | |
|---|---|
| **0 – none/not applicable** | **2 – other actor** [coded openly] |
| **1 – Elor Azaria** | **3 – Elor Azaria + other actor** [coded openly] |

A topic was coded to focus on **Azaria** as the main actor if its shared logic focused on interpreting or evaluating Azaria and his deeds. To be coded, there had to be at least two tokens that refer to him, his actions, or specific evaluations or actions directed at him, and fewer tokens that referred to other actors. A topic was coded to focus on **another actor** if its shared logic focused on interpreting or evaluating some different collective, institutional, or individual actors (e.g., his lawyer, the government, a newspaper, the protesters), following the same criteria. A topic was coded to focus on **Azaria and another actor** together if its shared logic focused on interpreting or evaluating the relationship between Azaria and some other actor, or if different tokens referred to different actors, and both Azaria and at least one other actor matched the criterion of at least two tokens and no other actor with more tokens. If neither of these criteria were met, we coded the topic to have **no actor focus**. For any topics that were coded to focus on another actor, either alone or together with Azaria, the identity of that actor was coded openly in a separate variable.

## V.TOPIC LIST: Top tokens & classification

**Table A1.** Topic list, top tokens & classification.

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| T01 | father, hospital, turns to, to wait, expression, mass, strengthening, camp, polish, call | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T02 | death, enlist, to begin with, village, dick, service, went, visit, laws, throws | 1 | 1 | 1 | 3 | 1 | 2 | 4 | 1 |
| T03 | deportation to Gaza, serve (his sentence), precedent, tonight, deport, produce, Shahid (Arabic for martyr), Syrian, identify, surprised | 1 | 1 | 0 | 0 | 1 | 1 | 7 | 2 |
| T04 | dead/died, response, small, writing, Eretz Israel (the land of Israel), march, runner, battalion, Supreme Court of Israel, epic | 0 | 1 | 0 | 0 | 0 | 0 | 7 | 0 |
| T05 | giant, similar, extreme, determines, rises, fend off, sin, death, minority, wiseman | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| T06 | release, sit in prison, take, Oren Hazan (MP), lion, born for freedom, conviction, leave, help | 1 | 1 | 1 | 2 | 1 | 2 | 6 | 3 |
| T07 | my father, of course, worse, especially, corrupt, discussion, to discuss, Kfir Division (elite division), IDF Chief of Staff, Military Prosecutor's Office | 1 | 1 | 1 | 3 | 1 | 1 | 4 | 3 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| **T08** | talking, luck, bullet, three, Hebron (town, site of the shooting), God forbid, hypocrisy, first, their, outstanding fighter | 1 | 1 | 1 | 3 | 1 | 1 | 7 | 3 |
| **T09** | convicted, equal, expectation, saying, offender, this year, Abramowitz (journalist), composition, racist, short | 0 | 1 | 1 | 1 | 1 | 1 | 6 | 3 |
| **T10** | pardon, sergeant, Mr., Israeli hero, help, IDF soldier, petition, release, say, proud | 1 | 1 | 1 | 3 | 1 | 3 | 6 | 3 |
| **T11** | well, voice, agree, attackers, reference, offense, gives, seriously, mob, rules | 1 | 1 | 1 | 0 | 1 | 1 | 6 | 3 |
| **T12** | most, unfortunately, Lapid (MP), masters, headline, leadership, resists/oppose, analysis, break, Herzog (MP) | 1 | 1 | 1 | 0 | 0 | 0 | 3 | 2 |
| **T13** | !, adding, opening, you forgot, excellent, tax/No., on the contrary, annoying, a, impression | 0 | 2 | 0 | 0 | 1 | 2 | 0 | 0 |
| **T14** | words, mistaken, think, comments, responsible, eyes, act, quantity, book, face | 1 | 1 | 1 | 0 | 1 | 1 | 3 | 3 |
| **T15** | terrorist, he, murder, shoot, need, no, act, shoot, same/him, minutes | 1 | 1 | 1 | 2 | 1 | 1 | 1 | 3 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) *(translation by the authors)* | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| T16 | others, I said, you said, grace, mean, violent, surely, thousands, find, deal/compete | 1 | 1 | 1 | 0 | 1 | 1 | 7 | 0 |
| T17 | Israel, shame to the state, traitors, he deserves it, defense force, commend ation, heads, lovers, enemies, IDF (Israel Defense Forces) | 1 | 1 | 1 | 3 | 1 | 2 | 7 | 1 |
| T18 | regarding, violence, they, face, worthy, worthy, agree, tear, expressing, finishing | 0 | 0 | 1 | 0 | 0 | 0 | 7 | 0 |
| T19 | rabbi, I thought, Cohen, abusive/hurt, convict, dumb, same, error, responds, guess | 0 | 0 | 1 | 0 | 1 | 1 | 7 | 2 |
| T20 | trial, hearing, court, pampim. com, Elor Azaria, watch, military, defense, ruling, the shadow | 1 | 1 | 1 | 0 | 0 | 0 | 6 | 3 |
| T21 | be, beautiful, morning, Hamas, about you, shot, killed, real, released, requests | 0 | 0 | 1 | 0 | 1 | 2 | 7 | 0 |
| T22 | the, to, and, of, message; him, forgot, will determine, closet | 0 | 2 | 1 | 0 | 0 | 0 | 6 | 1 |
| T23 | outside, convicted, dozens, protesters, hundreds, immediate, immediate, Kirya (IDF base), aids, give up | 1 | 1 | 1 | 3 | 1 | 2 | 5 | 3 |
| T24 | Bibi (Netanyahu, Prime Minister), thousand, support, come, smart, power, hours, long time, posts, page | 1 | 1 | 1 | 0 | 1 | 3 | 7 | 3 |

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| T25 | Eli, friend, understand, Ben Ari, conclusion, old/change, sides, correct, change, discussion | 1 | 1 | 1 | 0 | 0 | 0 | 3 | 2 |
| T26 | recognize/know, traitor, play, value, procedure, continue, times, reach, Israel hater, in your eyes | 1 | 1 | 1 | 0 | 1 | 1 | 7 | 2 |
| T27 | investigation, answer, Bogie Yaalon (Minister of Defense), verdict, read/call, Moshe, Supreme Court of Israel, line, move, inquiry | 1 | 1 | 1 | 0 | 0 | 0 | 6 | 3 |
| T28 | thereof, Arabs, Jews, foreign/ beside, soldiers, likes, Arabs, Muslims, racism, level | 1 | 1 | 1 | 0 | 1 | 1 | 5 | 2 |
| T29 | in favor, major general, reserve, freedom, message, amazing, rule, women, Uzi Dayan (MP), testify | 1 | 2 | 1 | 2 | 1 | 2 | 3 | 3 |
| T30 | an hour, at, post, Jerusalem, going out, showing, close, time, Tel Aviv, arriving | 1 | 1 | 1 | 3 | 1 | 3 | 5 | 2 |
| T31 | assembly, comparison, division, Rabin (former Prime Minister), political, rival, justification, speaker, discredit, cost/ rise | 1 | 1 | 1 | 2 | 1 | 1 | 5 | 2 |
| T32 | president, disgusted, as if, feeling, disgusting, repulsive, zeros, yuck, Americans, two years | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| T33 | sea, judicial system, outstanding soldier, death to terrorists, punishment, different/changed, Israeli society, dumb, circus | 1 | 1 | 1 | 3 | 1 | 2 | 6 | 1 |
| T34 | responding, work, scared, cameras, asking, escape, using, reign, Elor effect | 1 | 1 | 1 | 3 | 1 | 1 | 4 | 1 |
| T35 | this is, tweet, audio, information, wow, you heard, legitimate, commander, million, answers | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| T36 | garbage, stink, garbage, face, street, The Hague, enlightened, divide, channels, juice | 1 | 1 | 0 | 0 | 1 | 1 | 6 | 2 |
| T37 | demonstration, lightning/Barak (former Prime Minister), police, struggle, activists, demonstrations, indictment, organization, Rambam (famous Rabbi/hospital), against | 1 | 1 | 1 | 0 | 0 | 0 | 5 | 2 |
| T38 | justice/righteous, word, big, grave, morality, heritage, rock, right, damned, ahh | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T39 | hhh (sigh), caring, stupid, segment, understanding, you, very good, turned out, ars (slang for someone behaving like a criminal), Zoabi (MP) | 1 | 3 | 1 | 0 | 1 | 1 | 7 | 1 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| **T40** | completely, I heard, horrible, terrible, good/ beautiful, good, field, column, Kalman (journalist), Isaac | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| **T41** | soldier, the soldier, matters, trust, damage, say, admit, soldiers, bad | 0 | 0 | 1 | 1 | 1 | 1 | 7 | 1 |
| **T42** | judge, allowed, wind/spirit, period, forces, where, backup, went, turned, gives | 1 | 1 | 1 | 3 | 0 | 0 | 6 | 3 |
| **T43** | a little, less, very, missing, zero, this time, white, religious, religious Zionism, strange | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **T44** | the/he/God, morality, democracy, human being, ya!, you, delusional, fool/ stupid, for you/ go, stop | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 |
| **T45** | attacks/terror attacks, Azaria effect, guilt, begins, cowards, investigate, expected, suddenly, IDF Spokesman, unfortunately | 1 | 1 | 1 | 3 | 1 | 1 | 4 | 3 |
| **T46** | the, to, open, refer to, inn.co.il (news site), more, group, convict, pardon | 1 | 1 | 1 | 2 | 0 | 0 | 6 | 1 |
| **T47** | Walla (news site), Ynet (news site), Israel Hayom (newspaper), Galatz (army radio), Yedioth Aharonot (newspaper), controversy, punishment, discussions, stay, determined | 1 | 1 | 0 | 0 | 0 | 0 | 2 | 2 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| T48 | to her, she, daughter, said, leftist, two, judge, hers, absolution, dogs | 1 | 1 | 1 | 0 | 1 | 1 | 7 | 2 |
| T49 | in the name of, fabric, cold, believers, do, hug, her mother, Livni (MP), Eyal | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T50 | knows, against, numbers, Riklin (journalist), hatred, pension, root, Mizrachim (ethnic group), except, bad | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 2 |
| T51 | will be acquitted, hate, serious, British, antisemitic, face, pit, perception, admitted, slave | 1 | 1 | 1 | 0 | 1 | 1 | 7 | 3 |
| T52 | must, need, receive, Torah, change, righteous, offer, worry/take care, destroy, forward | 1 | 1 | 1 | 3 | 0 | 0 | 4 | 1 |
| T53 | after all, IDF Chief of Staff, to say, always, asked/question, definitely, a company/ friend, they said, obviously, Minister of Defense | 1 | 1 | 1 | 3 | 0 | 0 | 6 | 3 |
| T54 | king, look, single, see, straight, wave, bereavement, courage, outdoors, vacation | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T55 | believing, idiot, serve (sentence), wanting, will come out, created, defense (in trial), moral, fighter, partner/ accomplice | 1 | 1 | 1 | 2 | 1 | 2 | 7 | 3 |
| T56 | readers, result, say, Azaria, given, left, ethical, assistance, practices, imprisoned | 0 | 0 | 1 | 0 | 0 | 0 | 7 | 3 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| **T57** | related, disabled, next, understand, kills, command, facts, (those who come to kill you, you should) kill them first, convict, sapper | 1 | 1 | 1 | 3 | 0 | 0 | 1 | 1 |
| **T58** | sure, effect, Azaria's trial, here, soldiers, fear, policemen, reverse, called, disgrace | 1 | 1 | 1 | 3 | 1 | 1 | 4 | 3 |
| **T59** | prosecution, Attorney General, possibility, evidence, high, appeal, criminals, mediation, incarceration, failure | 1 | 1 | 1 | 3 | 1 | 1 | 6 | 3 |
| **T60** | relation, liar, judging, heroic, coming, far, logical, lying, understanding, prison | 1 | 1 | 1 | 3 | 1 | 1 | 6 | 1 |
| **T61** | tomorrow, sounds, his name/put, their, send, damn him, sat, remember, Zionism, final | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **T62** | thanks, Bogie (Yaalon, then Minister of Defense), account, say, provision, relevant, hear, correction, lying, insane | 1 | 1 | 1 | 0 | 0 | 0 | 7 | 2 |
| **T63** | commander, guardian, Mizrahi (ethnic group), Ashkenazi (ethnic group), youth, values, purity of arms, Ramle (town), sleep, rating | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 |
| **T64** | you, respect, my son, yours, their, they, hurt, be ashamed, could, children | 1 | 1 | 1 | 3 | 1 | 1 | 5 | 3 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) *(translation by the authors)* | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| **T65** | subject, just/ being right, legal, sitting in prison, mess, his own good, flying, senior officers, failure, corrupt | 1 | 1 | 1 | 0 | 1 | 1 | 6 | 1 |
| **T66** | political party, will do, elections, acts, talk, Knesset (Israeli Parliament), release, opposition, mandates, handful | 1 | 1 | 1 | 3 | 0 | 0 | 3 | 3 |
| **T67** | half, say, incitement, affair, testimony, agree, imprisonment, year, the commander, Naaman (journalist) | 1 | 1 | 1 | 0 | 0 | 0 | 6 | 3 |
| **T68** | Netanyahu (Prime Minister), the right, left, supports, right, please, Yaalon (Minister of Defense), Liberman (Minister of Foreign Affairs, later Minister of Defense), Bennett (Minister of Education), government/rule | 1 | 1 | 0 | 0 | 1 | 2 | 3 | 2 |
| **T69** | story, stage, journalists, you were, Olmert (former Prime Minister), verdict, say, hypocrite/ painted, Shimon, they | 1 | 1 | 1 | 0 | 1 | 1 | 2 | 2 |
| **T70** | Sharon Gal (journalist/ politician), Sheftel (Azaria's lawyer), Sharon, real, Big Brother, brother, enter, protect, appeal, Yoram Sheftel (Azaria's lawyer) | 1 | 1 | 1 | 0 | 1 | 1 | 2 | 2 |

*(Continued)*

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| T71 | Yossi, cliff, returned, television, real man, missing, acting, I wrote, kept, sat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| T72 | soldiers, they, officers, their, army, fighters, no, judges, commanders, weapons | 1 | 1 | 1 | 2 | 1 | 1 | 4 | 2 |
| T73 | reporter, Roni Daniel (journalist), broadcast, Minister of Education, Katz (Minister of Intelligence/ Transportation), Bar, Ilan, reporters, in line, I wanted | 1 | 1 | 1 | 0 | 1 | 1 | 2 | 2 |
| T74 | guilty, happened/cold, miserable, joke, supporters, permission, you, check, think, religious | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| T75 | terrorism, right, on, storm, becoming, leader, supposedly, helping, forget, Israelis | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T76 | /, read, =, page, new, reading,:, huge demonstration, &., | 1 | 1 | 1 | 0 | 0 | 0 | 5 | 3 |
| T77 | news story, education, article, excellent, lawyer, top, idea, convey, summary, career | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T78 | passing, with us, photography, stop, deep, screen, process, crazy, lost, listen | 1 | 1 | 1 | 0 | 1 | 1 | 7 | 2 |
| T79 | prosecutor, money, Nazis, promotion, eggs/ balls, brings, holocaust, understand, tramp, handcuffs | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 2 |

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| **T80** | intersection, hour, there, mall, entrance, intersections, one, Beersheva (city), Or/light, Negev (desert) | 1 | 1 | 1 | 3 | 0 | 0 | 5 | 1 |
| **T81** | B'Tselem (human rights NGO), knowledge, photographer, no doubt, in advance, behavior, outstanding, performed, brain, funny | 1 | 1 | 1 | 0 | 1 | 2 | 7 | 3 |
| **T82** | okay, difference, coward, Gilad Shalit (famous hostage), you, newspaper, shoot, insolence, man-of-the-year, choice | 1 | 1 | 1 | 3 | 1 | 1 | 5 | 3 |
| **T83** | releasing, to, stop/enough, must, its time, immediately, quickly, need, enough, injustice | 1 | 1 | 1 | 3 | 1 | 1 | 6 | 1 |
| **T84** | innocent, alive, essence, he is innocent, existing, uncle, movie, photos, you, pride | 1 | 1 | 1 | 3 | 1 | 3 | 5 | 1 |
| **T85** | God, read, pay, men, normal, disqualified, Lebanon, from above, come, exit | 1 | 1 | 1 | 0 | 1 | 2 | 7 | 1 |
| **T86** | prize, Maariv (newspaper), network, percentage, rotter (news site), survey, news, channel, RT (Russia Today), Mivzakon App (news site) | 1 | 1 | 1 | 3 | 0 | 0 | 2 | 3 |

(*Continued*)

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| **T87** | right/prerogative, oh, woman, understandable, made/make, she, girl, to handle, worse, importance | 1 | 1 | 0 | 0 | 1 | 1 | 7 | 2 |
| **T88** | sorry, minister, saw, hate, shocking, makes, asking, talk, terrorist family | 1 | 1 | 1 | 3 | 1 | 1 | 7 | 3 |
| **T89** | state, for us, shame, sad, shame and disgrace, we, parents, defend, killers/murderers, you | 1 | 1 | 1 | 3 | 1 | 1 | 5 | 3 |
| **T90** | Palestinians, occupation, territories, settlers, Israeli, rights, dangerous, food, war, their | 1 | 1 | 1 | 0 | 1 | 1 | 3 | 2 |
| **T91** | Azaria family, will be released, politicians, defense, his home, rotation, house arrest, abandon, finance, intervention | 1 | 1 | 1 | 3 | 0 | 0 | 5 | 3 |
| **T92** | For you, you, you were, when, home, mom, me, the son, on you, you | 1 | 1 | 1 | 3 | 0 | 0 | 6 | 3 |
| **T93** | hero, well done, strong, we're all, Ran, you, *****, brave, for you, you | 1 | 2 | 1 | 3 | 1 | 3 | 0 | 1 |
| **T94** | MK (Member of the Knesset), gang, Knesset, spoke, today, attorney, Rafael, apartheid, you wrote/reporter, market | 1 | 1 | 1 | 0 | 0 | 0 | 2 | 2 |

**Table A1.** (Continued).

| TOPIC | TOP TEN TOKENS (Highest Probability) (translation by the authors) | COHERENCE | TOPICALITY | RELEVANCE | EV. STANCE | EMOTIONALIZ. | EM. SENTIMENT | THEMATIC FOCUS | ACTOR FOCUS |
|---|---|---|---|---|---|---|---|---|---|
| T95 | meaningless, Likud (party), nonsense, vision, demagogy, murdered, to point at, lack, sanity, evil | 1 | 2 | 1 | 0 | 1 | 1 | 7 | 2 |
| T96 | brother/dude, come on, slow, pleasant, come on, nice, let go/ leave, stop, assault, weak | 1 | 3 | 1 | 3 | 1 | 2 | 7 | 1 |
| T97 | where, needless, forgotten, shit, forget, forget, the man, sick, what, forgive | 1 | 1 | 1 | 0 | 1 | 1 | 7 | 0 |
| T98 | Amen, success, love you, you, dear, with the help of god, for you, the name/ God, soon, end | 1 | 3 | 1 | 3 | 1 | 3 | 5 | 1 |
| T99 | context, false/lie, argument, like me, regret, writing, turns out, safe/sure, free, excuse | 1 | 1 | 1 | 0 | 1 | 1 | 7 | 0 |
| T00 | why, you, have, *****, exactly, Azaria, like, who, needs, doesn't | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |

## VI.VALIDATION: Reliability, accuracy & robustness

Our procedure was validated in three ways.

(1) To assess reliability, all 100 topics were double-coded by two coders (native speakers), yielding a reliability of *Krippendorff's* $\alpha = 0.88$.
(2) To assess precision and recall, we selected 200 documents at random from the corpus that were coded manually based on the same coding instructions used for the classification of topics. A comparison between the manual classification and the classification obtained via the classification of modeled topics yielded mostly very high values for precision ($M = 0.89$, $SD = 0.10$, range: 0.64–1.00) and recall ($M = 0.89$, $SD = 0.07$, range: 0.77–1.00).
(3) To assess robustness, we repeated the entire topic modeling and classification procedure for each platform separately, estimating topic models with $k = 80$ for Facebook, $k = 80$ for Twitter, and $k = 70$ for WhatsApp, respectively. All topics were again coded and used to classify the original documents. A comparison between both classifications yielded satisfactory robustness, with a Holsti coefficient of *0.80*.

## VII.TOPIC COVERAGE
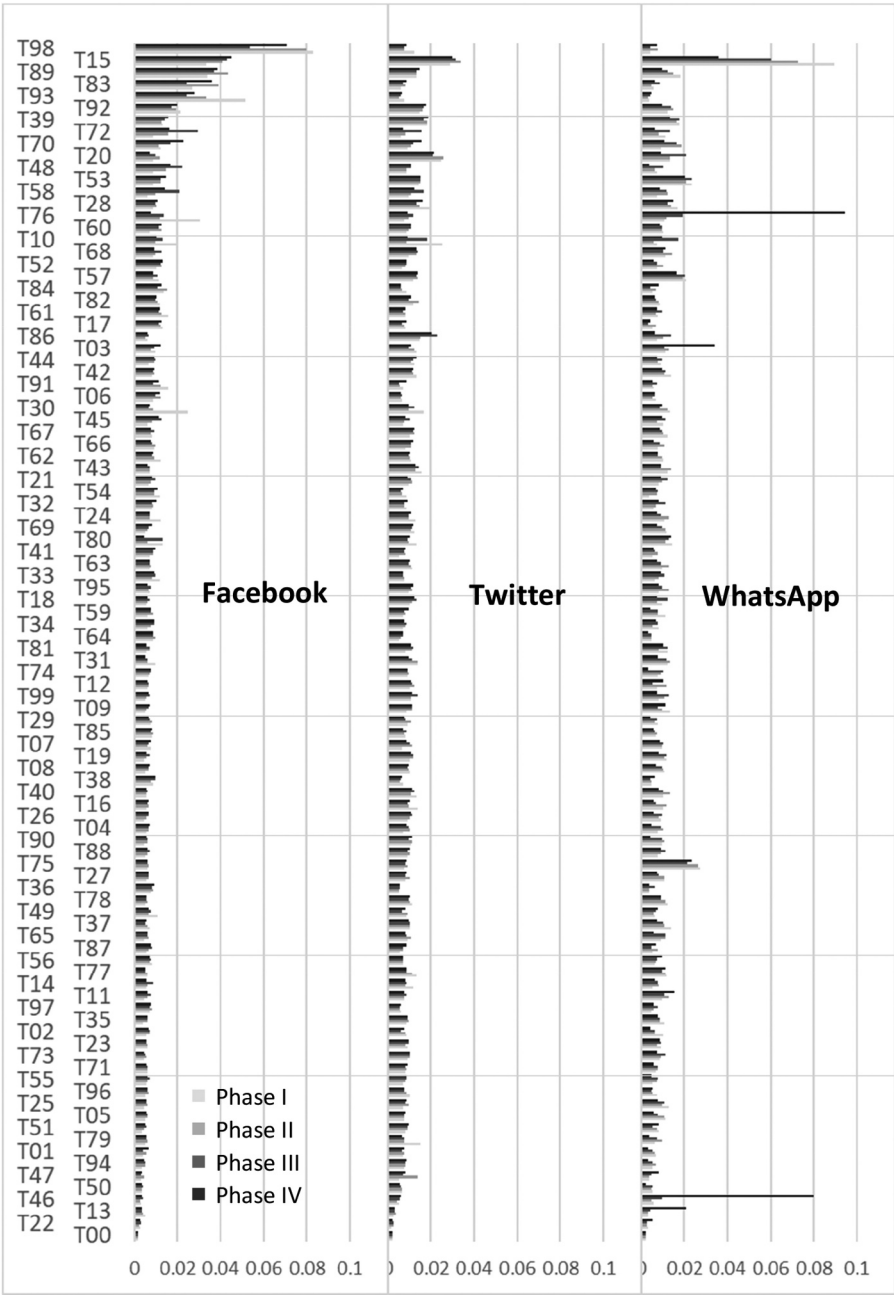


Figure A2. Topic coverage per platform and phase.

## VIII.ANOVA Table

**Table A2.** Average expressed sentiment by platform and interaction type.

|  | Sum of Squares | df | F | Sig. | Eta$^2$ |
|---|---|---|---|---|---|
| Corrected Model | 287.925 | 11 | 236.471 | .000 | .027 |
| Intercept | 236.540 | 1 | 2136.955 | .000 | .023 |
| Platform | 1.954 | 2 | 8.826 | .000 | .000 |
| Post by Supporter/Opponent | 2.220 | 1 | 20.058 | .000 | .000 |
| Response by Supporter/Opponent | 0.501 | 1 | 4.525 | .033 | .000 |
| Platform * Post by … | 15.658 | 2 | 70.727 | .000 | .002 |
| Platform * Response by … | 16.562 | 2 | 74.815 | .000 | .002 |
| Post by … * Response by … | 3.492 | 1 | 31.548 | .000 | .000 |
| Platform * Post by … * Response by … | 3.666 | 2 | 16.559 | .000 | .000 |
| Error | 10211.612 | 92254 |  |  |  |
| Total | 25222.774 | 92265 |  |  |  |
| Corrected Total | 10499.538 | 92265 |  |  |  |

Dependent Variable: Expressed Sentiment. Adjusted $R^2$ = 0.027.

## IX.SHARE NOTE

Due to data protection concerns, we cannot share our raw theta matrix. All other data and relevant scripts are available upon request from the authors.

For further methodological documentation, please refer to (Baden et al., 2019).