PSICOLOGIA

Revista PSICOLOGIA, Vol. 39(1), 2025, 44-62, https://doi.org/10.17575/psicologia.1974

Development and validation of a measurement of Bystander Intervention on Online Hate Speech towards Immigrants (BIOHS-Immigrants)

Catarina L. Carvalho^{1,2}, Isabel R. Pinto^{1,2}, Sara Alves¹ & Márcia Bernardo¹

¹ Faculty of Psychology and Education Sciences, University of Porto, Portugal ² Center for Psychology at University of Porto (CPUP)

Abstract: Online hate speech has profound implications for society, with migrants as primary targets. Underreporting by victims and bystanders obscures the true extent, indirectly legitimizing these crimes. To assess bystander intervention in online hate speech against immigrants, we developed a scale based on the five steps of the bystander intervention model. In Study 1 (N = 294), exploratory and confirmatory factor analyses confirmed the five-factor structure, having, as the final step, different types of actions that bystanders can adhere to when witnessing online hate speech. Structural equation modelling showed that, overall, each step was predicted by the preceding step, as proposed by the bystander intervention model. Study 2 (N = 240) replicated this finding and assessed the scale's criterion-related validity, revealing negative associations with moral disengagement and victim blaming, and positive associations with self-efficacy. We discuss the scale's relevance and applicability in studying bystander behaviour in response to online hate speech.

Keywords: Bystander effect; Bystander intervention model; Hate crimes; Online hate speech; Immigrants.

For the first time, on June 18, 2022, the International Day for Countering Hate Speech, was celebrated. This global celebration is part of the United Nations (UN) Strategy and Plan of Action on Hate Speech (UN, 2019), developed in response to the alarming rise of xenophobia, racism, and intolerance around the world.

Although there is still no consensus for an international legal definition, hate speech can be understood as "any kind of communication in speech, writing or behaviour, that attacks or uses pejorative or discriminatory language with reference to a person or group on the basis of who they are" (UN, 2019, p. 2). More specifically, hate speech corresponds to blatant discrimination that results in psychological aggression directed to an individual or group, simply based on attributes traditionally indicative of socially vulnerable or minority groups, such as gender, race, disability, religion, national and ethnic origin, and sexual orientation, or any other identity factor (e.g., UN, 2019).

Although hate speech is primarily a verbal behaviour, research suggests that it is often associated with prejudiced attitudes and aggressive intentions and it may be linked to other hate crimes, thereby having serious consequences both at the individual and societal levels. (e.g., Müller & Schwarz, 2018; Walters et al., 2016). At the individual level, evidence indicates that hate speech can have adverse psychological, emotional, and physical effects on victims, contributing to increased anxiety, feelings of humiliation, symptoms of depression, fear and insecurity (e.g., Dreißigacker et al., 2024; Waldron, 2012). At the societal level, it endangers fundamental human rights and democratic values and undermines social cohesion, social stability and peaceful coexistence by promoting social tension, violence, conflicts and division between social groups (UN, 2019). For instance, hateful content and inflammatory anti-migrant messages spread online by far-right political leaders were found to be associated with individuals' negative attitudes and threat perception towards migrants, online hate speech and the rise of (offline) hate crimes (Müller & Schwarz, 2018). Moreover, online hate speech prevalence and, particularly, perceptions that national institutions and/or social media platforms are ineffective at dealing with it, were found to undermine individuals' ability to notice and interpret an event as hate speech, leading them to minimize the impact and consequences of hate speech on victims (SELMA, 2019).

While hate crimes, such as hate speech, are increasingly visible, official statistics often underrepresent their true extent due to underreporting by both victims and witnesses, resulting in the impunity and encouragement of offenders (FRA, 2021; Pinto et al., 2023). This lack of reporting can significantly contribute to a skewed perception of the true scale of these crimes, which indirectly facilitates

¹ Correspondence address: Catarina L. Carvalho, E-mail: anacarvalho@fpce.up.pt

their normalization and perpetuation. The absence of public responses can be seen as acceptance or social legitimization, weakening social norms against the expression of prejudice (Álvarez-Benjumea, 2023). Furthermore, the lack of institutional and social responses to hate speech can encourage similar acts, creating a permissive atmosphere. Studies show that when bystanders fail to intervene in hate speech, such behaviour becomes normalized, with the absence of opposition perceived as tacit approval (Zapata et al., 2024).

Such passive or non-interventive behaviour may represent, itself, a surrogate or indirect, but no less consequential, discriminatory practice against socially vulnerable groups. Moreover, hate speech is, in fact, the hate crime that most enacts such passive or non-interventive behaviour, because of its potential interpretative ambiguity nature (see Papcunová et al., 2021). This passive behaviour can be explained within the theoretical framework of the bystander effect.

Bystander Effect

The tendency of a person witnessing (i.e., a bystander) an emergency situation and not seek or offer help to the victim or person in need, when other people are present, is known as *bystander effect* (e.g., Latané & Darley, 1969, 1970). The more people present, the less likely witnesses (bystanders) will intervene. It is proposed that bystanders' apathy occurs because each person feels less responsible to act in the presence of others (diffusion of responsibility) or because they infer, from the lack of intervention of the others, that the situation is not that serious (pluralist ignorance) (e.g., Latané & Rodin, 1969). Additionally, the fear of negative evaluations and embarrassment also contributes to the bystander effect, as the fear of making a mistake or failure while others are watching (e.g., Latané & Darley 1970).

Bystander effect in the cyberspace. The bystander effect, initially studied to explain individuals' lack of intervention in emergencies, has expanded to include antisocial and harmful situations like bullying, sexual harassment, and assault (Kettrey & Marx, 2020; Latané & Darley, 1969; Nickerson et al., 2014). More relevant to our work, researchers have also examined this effect within digital environments, where computer-mediated interactions often reduce personal accountability, amplifying the inaction of bystanders in response to hate speech (Jubany & Roiha, 2016; Markey, 2000; Obermaier et al., 2021). For instance, a report by the UK Safer Internet Centre (2016) shows that the majority of youths (82%) had witnessed online hate, but less than half had chosen to report it. Indeed, in online contexts, the bystander effect may be intensified due to factors such as anonymity, audience size, and the absence of immediate accountability, which reduce individuals' likelihood to intervene when witnessing hate speech (e.g., Barlińska et al., 2013). Additionally, the lack of personal contact and face-to-face interactions in online environments fosters a diffusion of responsibility, further diminishing bystanders' sense of accountability and willingness to take action (e.g., Latané & Darley, 1970; Siapera et al., 2018).

Social media platforms intensify these dynamics by enabling rapid and often anonymous communication, which amplifies the reach and impact of hate speech (UN, 2018). Furthermore, social control measures on these platforms are widely perceived as insufficient or ineffective in deterring online misbehaviour, which further reinforces bystanders' tendency to ignore harmful content (Jubany & Roiha, 2016). As a result, the online expression of overt intolerance and hate towards members of vulnerable or minority groups remains widely tolerated and indirectly legitimized, largely due to the bystander effect.

The Role of Bystanders in Online Hate Speech Towards Immigrants

While hate speech targets various social groups, hate speech against migrants has shown a troubling increase in recent years (UN, 2019). Intensified perceptions of threat and insecurity associated with the arrival of migrants and asylum seekers from diverse backgrounds have risen across Europe (Pinto et al., 2020). Consistently, the second European Union Minorities and Discrimination Survey (FRA, 2017) identifies immigrants as one of the groups most affected by discrimination and hate crimes in Europe. This trend has been intensified by events such as the Mediterranean migration crisis and the economic pressures of the COVID-19 pandemic, which have further fuelled discriminatory attitudes toward this population (Vega Macías, 2021).

In the digital context, the prevalence of online hate speech targeting immigrants is particularly concerning. Such speech is not only widespread but often goes unchallenged due to bystander inaction. This "silent" phenomenon is socially, politically, and legally neglected. It is rarely discussed in public discourse, and many people remain unaware of its existence, their own biased behaviour, and the broader societal consequences of this bias. By failing to intercede, bystanders inadvertently contribute to the normalization and perpetuation of racism, discrimination, and racial violence (Murrell, 2021). Thus, addressing the bystander effect in the context of online hate speech is paramount for combating prejudice and discrimination against minority groups (Stewart et al., 2014).

Bystander Intervention Model

Previous research has determined that bystanders' intervention requires five sequential and crucial steps, so that action actually occurs (i.e., bystander intervention model; Latané & Darley, 1969, 1970). According to Latané and Darley's (1970) bystander intervention model, in order to take action and counter the bystander effect, people have to (1) notice the event (i.e., becoming aware that something is happening), (2) interpret the event as an emergency (i.e., recognizing the situation as serious or requiring intervention, which includes determining that something is wrong and interpreting the situation as threatening), (3) accept individual responsibility to intervene (i.e., feeling responsible to act after recognizing that the situation requires help), (4) know and decide how to intervene or provide help (i.e., identifying what form of intervention to implement after accepting the responsibility to act), and finally (5) implement intervention *decisions* (i.e., taking action based on the decision made in the previous step). If any of these steps are not reached and completed, bystanders are less likely to intervene. Thus, in order to analyse the process that precedes bystander intervention (vs. bystander apathy - bystander effect), and the extent to which bystanders intervene (and how), facing online hate speech, it is crucial to assess all of the five steps.

Determinants of Bystander Behaviour

Theoretical and empirical evidence has shown that bystander intervention (prosocial behaviour) is negatively related to moral disengagement and victim blaming (a critical barrier for helping behaviour), and positively related to self-efficacy (e.g., Bandura et al., 1996; Clark & Bussey, 2020; Ferreira et al., 2020; Koehler & Weber, 2018; Machackova, 2020; Rudnicki et al., 2022; Stewart et al., 2014).

Moral disengagement. Moral disengagement refers to a set of psychosocial processes that enable individuals to justify morally questionable behaviours without altering their core moral standards, allowing them to avoid self-criticism or social disapproval (Bandura et al., 1996). These processes include mechanisms such as moral justification, euphemistic labelling, advantageous comparison, displacement of responsibility, diffusion of responsibility, distortion of consequences, dehumanization, and attribution of blame. Through these mechanisms, individuals can feel less guilt for ignoring critical situations, often by shifting responsibility to others or dehumanizing or blaming the victims. Research shows that moral disengagement is positively associated with the bystander effect and negatively related to prosocial behaviours, such as defending victims in situations of harm or discrimination (Gini et al., 2020; Sjögren et al., 2021; Thornberg et al., 2020). In online settings, moral disengagement may be exacerbated by factors such as anonymity and the lack of face-to-face interaction, which reduce personal accountability and discourage individuals from intervening in cases of hate speech or other harmful behaviours (Obermaier et al., 2021; Siapera et al., 2018).

Victim blaming. In addition to being one of the mechanisms of moral disengagement, victim blaming is also one of the major barriers to action identified by the bystander intervention model (Latané & Darley, 1970). It can prevent individuals from taking responsibility to intervene by attributing blame to the victim and shifting the responsibility for the situation or incident onto them. Indeed, research has shown that victim blaming —defined as the tendency to hold the victim responsible for their misfortune—lowers bystanders' intentions to help (e.g., Koehler & Weber, 2018).

Self-efficacy. Perceived self-efficacy corresponds to people's beliefs about their ability to plan and execute the necessary courses of action to produce a desired outcome (e.g., Bandura, 1998). As highlighted by Bandura (1998, p. 51), "perceived self-efficacy operates as a central self-regulatory mechanism of human agency". Research has shown that self-efficacy plays an essential role in the decision-making process to engage in helping behaviour (e.g., Ferreira et al., 2020), namely in the online context (Costello et al., 2017) and, thus, has the potential to decrease the bystander effect (Ferreira et al., 2020).

Bystander Intervention Measurement

Several instruments assessing bystander behaviour were already developed in previous research (although not all of them take into account Latané and Darley's bystander intervention model), for instance, in the context of bullying and sexual harassment (e.g., Nickerson et al., 2014), sexual violence (e.g., Bennett et al., 2014), interpersonal violence (e.g., Banyard, 2008), cyberbullying (e.g., Koehler & Weber, 2018; Bastiaensens et al., 2014), racism (e.g., Palmer et al., 2017) and prosocial behaviours towards refugees (Albayrak-Aydemir & Gleibs, 2021).

However, as far as we know, no instrument measuring bystander response facing online hate speech and, particularly, representing Latané and Darley's bystander intervention model, has been developed. Thus, this research aims to address this gap by developing a scale specifically designed to measure bystander responses to online hate speech against immigrants, structured to align with Latané and Darley's bystander intervention model.

The Present Research

We have observed an alarming rise of xenophobia, racism and intolerance around the world (UN, 2019), resulting in an increase in hate speech against immigrants. Both victims and bystanders tend to fail to report such crimes, contributing, indirectly, to legitimize and perpetuate its occurrence. Thus, witnesses' action is crucial to combat prejudice and discrimination against immigrants (or against any other socially vulnerable or minority group). However, assessing witnesses' intervention (vs. passive observation) in the context of online hate speech is complex, largely because of the special characteristics of both hate speech (e.g., its potential interpretative ambiguity nature) and the cyberspace itself (that allows anonymity and reduces users' accountability). Thus, a measure to assess the process and necessary steps for bystander intervention against online hate speech towards immigrants represents an important and useful tool in this field.

Therefore, the aim of this investigation is to (a) develop and validate the Bystander Intervention on Online Hate Speech towards Immigrants (BIOHS-Immigrants) scale, which assesses individuals' likelihood of intervention when witnessing online hate speech directed at immigrants, grounded in the sequential bystander intervention model (Latané & Darley, 1969, 1970; Nickerson et al., 2014); (b) examine the extent to which each step of the bystander intervention model predicts the subsequent step; (c) explore how the bystander intervention model, as well as different intervention actions within online hate speech contexts, relate to moral disengagement, victim blaming, and self-efficacy.

Specifically, through two correlational studies, we anticipate that the BIOHS-Immigrants scale will exhibit a multifactorial structure, with five factors corresponding to each step in the bystander intervention model (Latané & Darley, 1969, 1970). The final step (i.e., Step 5 - Implementing the intervention decision) is expected to encompass distinct types of actions tailored to the online hate speech context (Study 1). Moreover, each step of the model is expected to predict the following step sequentially (Studies 1 and 2). Finally, we expect that the BIOHS-Immigrants scale will be negatively related to moral disengagement and victim blaming and positively related to self-efficacy (Study 2).

Data Analytic Strategy. We proceeded to the analysis of the factorial structure and psychometric properties including validity and reliability of the new scale. First, an exploratory factor analysis (EFA) was conducted, using principal axis factoring with oblique rotation (direct oblimin) with Kaiser normalization approach (i.e., eigenvalue > 1.00), to examine the factor structure of our BIOHS-Immigrants scale and to reduce the initial number of items. Then, we conducted a confirmatory factor analysis (CFA) to examine the quality of the final factor structure (i.e., test and validate the measurement model), and a structural equation modeling (SEM) to test the sequential steps of the bystander model (i.e., the actual path of the model itself) using a maximum likelihood estimation procedure.

For both CFA and SEM, the goodness-of-fit of the model was evaluated through multiple criteria (e.g., West et al., 2012), such as chi-square and the ratio of the chi-square to its degree of freedom (CMIN/DF), goodness-of-fit (GFI), normed fit index (NFI), comparative fit index (CFI), and root mean square error of approximation (RMSEA). A good model fit should provide a non-significant Chi-Square (CMIN), and CMIN/DF should be \leq 3, but acceptable at values \leq 5 (Marsh & Hocevar, 1985), indicating an excellent and a good fit, respectively. Moreover, according to previous recommendations, typically, for the GFI, NFI, and CFI, values > .90 indicate a good fit and > .95 a very good fit (e.g., Hu & Bentler, 1999; West et al., 2012; Whittaker & Shumacker, 2022). For the RMSEA, values of 0 (zero) indicate a perfect fit, values \leq .05 indicate "close fit" or "good fit", up to .08 indicate a reasonable fit, and values \geq .10 suggest a poor fit (e.g., Brown, 2015; West et al., 2012; Whittaker & Schumacker, 2022). The confidence interval (CI) for RMSEA (typically 90% CI) indicates the precision of the RMSEA point estimate (Brown, 2015).

We also used Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC) to compare models with smaller values indicating a better fit (e.g., Kline, 2016).

Finally, for both CFA and SEM, Modification Indices (MI \ge 3.84; Brown, 2015; Whittaker, 2012) were considered to perform post-hoc modifications, specifically the inclusion of additional parameters. Only within-factor error covariances were permitted, with associations added sequentially based on the highest MI to avoid overfitting. This process was repeated until no further modifications were necessary. All modifications implemented based on MIs were carefully evaluated to ensure they were theoretically justified. Specifically, correlated error terms were introduced only between items within the same factor, reflecting potential similarities in item wording or conceptual alignment. These adjustments aimed to improve the statistical fit of the model while maintaining the theoretical integrity of the bystander

intervention model (Latané & Darley, 1970). No modifications altered the fundamental structure of the model or the sequential nature of the steps.

For both CFA and SEM, standardized coefficients are always reported.

Regarding the reliability evaluation of the BIOHS-Immigrants scale, Cronbach's Alpha and inter-item correlation coefficients were used to estimate internal reliability of each dimension (subscales) of the **BIOHS-Immigrants scale.**

Additionally, we examine the construct validity of our BIOHS-Immigrants scale (Study 2), by evaluating the relationship with other theoretically related measures: moral disengagement, victim blaming, and self-efficacy. Pearson's correlation coefficient (Study 2) was used to assess the strength of association between the BIOHS-Immigrants scale and the related constructs.

All data analyses were performed with IBM SPSS Statistics v. 29 and IBM SPSS Amos v. 29.

STUDY 1

Method

Participants. Participants were 294 Portuguese citizens (160 female, 128 male, 6 indicated "other"), ranging from 18 to 78 years old (*M* = 28.11, *SD* = 14.19), 55% completed secondary education, 45% higher education, and less than 1% basic education. The majority were students (64%) and employed (27%), and the remaining were unemployed (5%), retired (4%), and less than 1% were housewife.

Regarding the left-right political spectrum, the average score on a 7-point Likert-scale (1 = Left, 7 = 1)*Right*), was close to the scale midpoint of 4, though leaning left (M = 3.26, SD = 1.47). Participants were also asked about perceived socioeconomic status compared to other citizens of the country where they live on a 7-point Likert-scale (1 = Verv low, 7 = Verv high; M = 4.13, SD = 0.96).

The sample size was determined following established guidelines to ensure accurate factor solutions and stable parameter estimates. For EFA, we adhered to Hinkin et al. (1997), who recommend a minimum of 150 participants, as well as the commonly accepted subject-to-item ratio of at least 5:1 (Hair et al., 2014; Osborne, 2014), resulting in a minimum requirement of 210 participants for our initial pool of 42 items in Study 1. For CFA, we followed guidelines suggesting a sample size between 100 and 200 (Brown, 2015; Hinkin et al., 1997), while ensuring a minimum of 5 participants per item (Hair et al., 2014). For SEM, we followed recommendations to include at least 200 participants to achieve stable parameter estimates (Kline, 2016).

Procedure. Before conducting this investigation, all the procedures and materials were submitted, for analysis, to the Faculty of Psychology and Education Sciences of University of Porto Ethics Committee (Ref. 2021/07-09b) and data protection department.

The study was conducted using Qualtrics, and participants were contacted via online platforms (Facebook ads and Facebook groups, and University mailing list), to fill out a survey about online hate speech towards immigrants.

Participation was completely voluntary and 5 vouchers of €20 were raffled as compensation. After giving informed consent (which included information on confidentiality, anonymity, risks and benefits, contact information, etc.), participants provided sociodemographic information (e.g., age, gender, education). Next, participants read a text defining hate crime and online hate speech in particular, so that all participants started the survey with the same knowledge on the topic, namely being aware that hate speech is a hate crime legally framed by the Portuguese law and, thus, an undesirable and intolerable behaviour (full text is available at the online supplementary material, OSM 1).

Then, participants answered the initial item pool of our scale. On completion, participants were thanked, and further information about the project and online hate speech were provided.

Instrument. The items of our BIOHS-Immigrants scale were developed based on the theoretical framework of the bystander intervention model, following its five sequential steps (Latané & Darley, 1969, 1970), with adaptations for the context of online hate speech towards immigrants. For the initial item pool, we created specific item sets to capture each stage of the bystander intervention process, ensuring they reflected the unique characteristics of online hate speech: 1) noticing and identifying hate content online, 2) interpreting it as a serious occurrence, 3) accepting responsibility to intervene, 4) knowing how to intervene, and 5) intending to implement an intervention. Additionally, the scale development was informed by existing instruments, such as Nickerson et al.'s (2014) scale for bullying and sexual harassment, which measures this sequential process. We also examined measures that, while not explicitly structured around the sequential process, incorporate relevant aspects of bystander intervention across

several social contexts. In particular, we examined scales that assess actions bystanders might take in response to online hate speech (i.e., the final step). Thus, for the final step (Step 5 – Implement intervention decision), which includes potential helping behaviours in response to online hate speech, we drew on Palmer and colleagues' (2017) bystander response scale for racism in school, which categorizes responses as prosocial (e.g., "I would tell a teacher or member of staff what the person had said"), aggressive (e.g., "I would start a fight with the person who said something nasty..."), or passive (e.g., "I would ignore it and walk away"). We also drew on the work of Bastiaensens and colleagues (2014), who distinguish between public and private actions in a cyberbullying context. These include different types of helping intentions, such as asking for help (e.g., reporting the incident to someone who can assist, such as social media platforms, an organization, or a teacher), personally defending the victim (e.g., "Defending Joni in a Facebook comment"), or offering emotional support (e.g., "Comforting Joni in a Facebook comment"). Additionally, inspired by Banyard and colleagues' (2005) Bystander Behaviour Scale, we included items relating to seeking assistance from formal or support entities (e.g., "Call 911 if a stranger needs help", "Call a rape crisis center or talk to a resident counselor"). Finally, we drew on Koehler and Weber's (2018) Willingness to Help Scale, which includes confrontational actions directed toward the offender (e.g., "I would publicly share a comment on Facebook in which I confront the bullies"). For items involving formal or support entities, we tailored them to the resources available in the Portuguese context. These include the Portuguese Victim Support Association (APAV), which has a dedicated unit to support migrant victims of hate crimes, violence, and discrimination; the Linha Internet Segura (Safe Internet Line), which addresses issues related to online platforms and technology use and provides a service to report illegal content online, including hate speech; and Portuguese law enforcement and public security agencies (e.g., Public Security Police - PSP) and judicial authorities (e.g., Public Ministry – MP) that can be contacted to report crimes. Detailed information on the sources of each item of the final scale, modifications made to adapted items, and theoretical bases for new items can be found in OSM 2.

The initial pool comprised 42 items, approximately 6 items per step for the first four steps. Given the complexity of bystander intervention in online hate speech, the final step (Step 5 - Implement intervention decision) included 18 items covering a range of actions: (1) intention to report online hate speech through formal mechanisms (e.g., "I personally report the situation to the police..."), (2) intention to report through informal support organizations (e.g., "I report the situation to an organization that deals with online hate speech..."), (3) intention to report via social media using reporting tools (e.g., "I report the offender's post, comment, or tweet as abusive or hateful"), (4) intention to help the victim directly by defending them (e.g., "I defend the immigrant in a Facebook comment...") and (5) intention to help the victim indirectly by confronting the offender aggressively (e.g., "I insult the offender in an unpleasant way..."). The latter two types of actions include items relating to both private and public actions.

These items were designed to encompass common intervention actions at varying levels of personal investment or effort and potential cost or risks. For instance, reporting to the police may require more time and may involve dealing with legal procedures, which has been identified as a common reason not to report a hate incident (FRA, 2017).

For each item, participants were asked to indicate their level of agreement on a 7-point Likert scale (1 = *I fully disagree*; 7 = *I fully agree*).

Results

Exploratory factor analysis and initial item reduction. Following previous recommendations (e.g., Hinkin et al., 1997), the criteria for retaining the final set of items were as follows: a) item communality above .40, b) factor loading greater than .40, c) items that clearly loaded on a single factor. In addition, some items were also dropped due to conceptual redundancy or misfit. Using the preceding criteria, items were deleted and EFA was repeated (we applied a new EFA every time items were deleted) until we obtained a clear factor structure matrix. Through this process, 16 items were excluded, and the resulting scale consisted of a total of 26 items (see OSM 2 for the final scale with 26 items and respective dimensions).

The Kaiser-Meyer-Olkin (KMO) sample adequacy index showed a value of .83 and the Bartlett's sphericity test was statistically significant, with a χ^2 (325) = 5267.70. p < .001, showing that the correlation between items was high enough for meaningful extraction.

The exploratory factor analysis resulted in an eight-factor solution (see OSM 2 for loadings and communalities of the final items).

The 8-factor solution accounted for 70.64% of the overall variance: Step 1 – Notice the event corresponded to Factor 6 (4% of variance); Step 2 – Interpret event as severe and Step 3 – Accept responsibility to help loaded together on the same factor corresponding to Factor 1 (30% of variance); Step 4 – Know how to help corresponded to Factor 3 (9%); and Step 5 - Implement intervention decision

corresponded to Factor 2 (Aggressive response; 13% of variance), Factor 4 (Report through formal mechanisms; 5%), Factor 5 (Report through social media mechanisms; 4%), Factor 7 (Public actions; 3%), and Factor 8 (Private actions: 3%).

Since the EFA showed that the items corresponding to Step 2 – Interpret event as severe (items 3 and 4) and Step 3 – Accept responsibility to help (items 5 to 8) loaded together, the theoretical and empirical structures were tested and compared in the following analysis.

Confirmatory factor analysis (CFA). To assess the quality of the final factor structure, we conducted a CFA.

First, we started by testing and comparing models reflecting both the theoretical (four-step solution with separate factors for Steps 2 and 3) and empirical (three-step solution with combined factor for Steps 2 and 3) structures, regarding the initial Steps of the model (i.e., Step 1 to 4). As we can see in Table 1, the four-step solution model fits the data better (MI suggested correlating two error terms in Step 4: e10 <-> e12; see OSM 3 for CFA and SEM results before modifications for all models) than the three-step solution model (MI suggested correlating two error in the combined factor for Steps 2 and 3, and two error terms in Step 4: e3 <-> e4, e10 <-> e12, respectively) confirming the theoretical structure. Based on these results, particularly the smaller AIC and BIC values observed for the four-step solution model compared to the three-step solution model, we proceeded with the subsequent analyses, retaining the two separate factors.

Then, we tested five independent models corresponding to the five steps of the bystander intervention model (i.e., Notice the event, Interpret event as an emergency or severe, Accept responsibility to help, Know how to help, Implement intervention decision), one model for each of the five types of actions in place of the final dimension (i.e., Implement intervention decision: Report through formal mechanisms, Report through social media mechanisms, Public actions, Private actions, Aggressive response).

As we can see in Table 1, the five-step solution model with Report through formal mechanisms subscale as the final step (see Models' Figures at OSM 4) fits the data in a good way, as well as the models with Report through social media mechanisms, that is, Informal mechanisms (MI suggested correlating two error terms in Step 4: e10 <-> e12), Public actions (MI suggested correlating two error terms in Step 4: e10 <-> e12), Private actions (MI suggested correlating two error terms in Step 4: e10 <-> e12), and with Aggressive response subscale (MI suggested correlating two error terms in Step 4 and two error terms in the Aggressive bystander response step: e10 <-> e12, e23 <-> e24, respectively).

Note that the number of each error term corresponds to the item number (see the complete list of items and their corresponding numbers in the OSM 2). We also tested the model with all Step 5 subscales (Model 6; see Table 1). However, this model (MI suggested correlating two error terms in Step 4 and two error terms in the Aggressive bystander response step: e10 <-> e12, e23 <-> e24, respectively) shows a poorer fit across several key indices compared to the other models. Specifically, the fit indices, including GFI (.862) and NFI (.884), were below the recommended thresholds, suggesting that the full model may not capture the structure of the data as effectively as the individual action-type models. Moreover, the full model had the highest values for both AIC (812.14) and BIC (1143.66), indicating a less efficient fit compared to the other models.

Structural equation modeling (SEM). In the next step of the data analysis, we tested the sequential steps of the bystander intervention model for each of the five types of actions.

As we can see in Table 2 (see Models' Figures at OSM 5), the models with Report through formal mechanisms subscale (MI suggested correlating two error terms in Step 4: e9 <-> e11); Public actions (MI suggested correlating error terms in Step 4: e9 <-> e10 and e10 <-> e12); and with Aggressive response (MI suggested correlating two error terms in Step 4 and two errors terms in Aggressive Response: e10 <-> e12 and e23 <-> 24, respectively) as final step provided an acceptable or good fit to the data and all of the direct paths (regression weights between latent variables for each sequential step in the model) were positive and statistically significant.

The model with the Report through social media mechanisms subscale provided a weak fit to the data.

Moreover, although all the direct paths were positive and statistically significant between the initial Steps, the regression weights between latent variables Step 4 (Know how to intervene) and Step 5 - Report through Social Media mechanisms were marginally significant, p = .089 (MI suggested correlating two error terms in Step 4: e10 <-> e12). Given this result, we considered whether the sequential process might vary depending on the type of action, particularly as some studies, such as Albayrak-Aydemir and Gleibs (2021), have also found no significant relationship between Step 4 (Know how to intervene) and Step 5 (Implement intervention decision) regarding some of the potential types of actions. In fact, the behaviour of "reporting through social media mechanisms" may not necessarily require specific technical knowledge or detailed

preparation (i.e., Step 4). This type of action may rely more directly on motivation to help (Step 3 - Accepting responsibility). Therefore, we tested an alternative sequential model that included a direct link between Step 3 and Step 5. We observe that Report through social media mechanisms is predicted by Accepting responsibility to help (see OSM 5), and results indicate that this model demonstrates a good fit, suggesting that Knowing how to help (Step 4) may not be essential in the intervention process for this specific action: χ^2 (*df*) = 196.14 (84)***, CMIN/DF = 2.34, GFI = .920, NFI = .936, CFI = .962, RMSEA [CI] = .067 [.055, .080], AIC = 268.14, BIC = 400.74.

Similarly, although the model with the Private actions subscale provided a good fit to the data, and all of the direct paths between the initial Steps were positive and statistically significant, the regression weights between latent variables of Step 4 (Know how to intervene) and Step 5 - Private actions were non-significant, p = .203 (MI suggested correlating two error terms in Step 4: e10 <-> e12). As in the previous model, we tested an alternative sequential model that included a direct link between Step 3 (Accepting responsibility) and Step 5 - Private actions. We observe that Private actions is predicted by Accepting responsibility to help (Step 3), and results indicate that this model demonstrates a good fit, suggesting that Knowing how to help (Step 4) may not be essential in the intervention process for this specific action as well: χ^2 (*df*) = 126.83 (71)***, CMIN/DF = 1.79, GFI = .942, NFI = .947, CFI = .976, RMSEA [CI] = .052 [.037, .066], AIC = 194.83, BIC = 320.07.

We also tested this alternative sequential process (i.e., a direct link between Step 3 and Step 5) for the remaining models. The results indicate that including this direct link neither improves the model's fit for these actions nor alters the sequential process, supporting the need to maintain the originally proposed sequential structure. Thus, as expected, each step in the model was predicted by the preceding step, except for the models for the Report through social media mechanisms and Private actions subscales. This sequential structure will be further examined in Study 2.

Finally, we also tested the model with all Step 5 subscales (Figure 1). We observed that all the direct paths were positive and statistically significant, showing that all the steps were predicted by the previous step (note that all the proposed final steps were predicted by Step 4). MI suggested correlating two error terms in Step 4 (e10 <-> e12) and two error terms in Step 5 – Aggressive bystander response (e23 <-> e24). However, the model presented a weak fit to the data (see Table 2).



Note: All path coefficients are standardized estimates.

 $p \le .05; p \le .01; p \le .001$

Figure 1. Simplified Graphical Representation of the SEM Model with All Step 5 Subscales in Study 1.

Reliability and descriptive statistics. Internal consistency was appropriate for all the steps of the bystander intervention model (see OSM 6 for reliability, descriptive statistics and bivariate correlations). As expected, each step of the bystander intervention model is positively correlated with the

subsequent step, and Step 4 is positively associated with all the proposed actions corresponding to Step 5.

Additionally, OSM 6 includes analyses of demographic data, such as differences between men and women across the steps and correlations between demographic variables (e.g., political orientation) and the steps, addressing the relevance of these factors to the model. Examining gender differences, we observe

that women generally scored higher than men across most steps, suggesting a stronger inclination to recognize and respond to online hate speech. Specifically, significant differences emerged in Step 1 (Notice the event), Step 2 (Interpret event as severe), and Step 3 (Accept responsibility to help), with women showing higher mean scores than men. However, in Step 4 (Know how to help), men scored higher than women. Additionally, we observed that women scored higher than men in Step 5 - Report through Social Media mechanisms. Conversely, men scored slightly higher in Aggressive Response. These findings suggest that gender may play a significant role in bystander intervention processes, with women potentially being more likely to identify and interpret online hate speech as severe and to accept responsibility for intervening.

Both groups showed higher mean scores for Step 5 - Report through social media mechanisms compared to other types of intervention actions, suggesting that individuals may be more inclined to rely on social media reporting tools as a preferred method for addressing online hate speech.

We also found that age has a significant negative correlation with Step 1: Notice the event, suggesting that older individuals may be less likely to notice the event. Conversely, age is positively correlated with Know how to help, Reporting through formal mechanisms, Public actions, and Aggressive response. Education shows a positive correlation with Know how to help but a negative correlation with Private actions, indicating that individuals with higher levels of education are more likely to know how to help and less likely to engage in private actions. Political orientation was negatively correlated with all steps of the intervention model. Finally, perceived social status did not show significant correlations with any of the steps, suggesting that social status may not influence bystander intervention behaviours.

Discussion

An EFA using principal axis factoring identified eight factors: three corresponding to the initial Steps of the Bystander Intervention Model (i.e., Notice the event, Interpret event as an emergency or as severe, Accept responsibility to help, Know how to help) and five corresponding to the final Step (i.e., Implement intervention decision), each representing different types of actions that bystanders might take when witnessing online hate speech. Contrary to our expectations, the items corresponding to Steps 2 (Interpret event as an emergency) and 3 (Accept responsibility to help) loaded onto a single factor. However, the CFA demonstrated a good fit to the data when these steps were modelled as distinct constructs. This finding suggests that the subscales of our measure align well with the original theoretical framework proposed by Latané and Darley (1970). According to the bystander intervention model, interpreting the severity of the problem (Step 2) is a distinct process from accepting personal responsibility to intervene (Step 3), and maintaining this distinction is critical for understanding specific barriers to intervention. For example, accepting responsibility to act may vary even when the severity of the situation is fully recognized. Distinguishing Step 2 and Step 3 allows for a more nuanced analysis, helping to identify whether barriers to intervention stem from difficulties in recognizing the gravity of the situation or from a reluctance to accept personal responsibility. This distinction is particularly valuable for designing targeted interventions to address these specific barriers.

We tested five separate models corresponding to the five steps of the bystander intervention model, with each model incorporating a different type of action as the fifth and final step (i.e., Report through formal mechanisms, Report through social media mechanisms, Public actions, Private actions, Aggressive response). The CFA revealed that the five-step model with each action type showed a good fit across models, indicating that our proposed structure for each type of intervention action was well-represented in the data.

SEM results generally demonstrated that each step was predicted by the preceding step in the model. However, two actions—Report through social media mechanisms and Private actions—did not show a significant relationship with Step 4 (Know how to help), suggesting that bystanders might proceed to intervene directly after accepting responsibility (Step 3) without requiring the additional step of knowing how to intervene (Step 4). This pattern may indicate that, in specific contexts, bystanders may feel ready to act once they acknowledge responsibility, bypassing the need for a detailed knowledge of intervention methods

In testing an integrative model that included all Step 5 subscales, SEM results confirmed that all proposed final steps were predicted by Step 4, consistent with the expected sequential process (e.g., Latané & Darley, 1970). However, this model showed a weak fit to the data, suggesting that it may be more appropriate to consider five independent models, each corresponding to one of the proposed actions, rather than a single comprehensive model encompassing all potential actions. However, it is important to note that, for both CFA and SEM, combining all Step 5 subscales increases the complexity of the model, which can contribute to poorer fit indices. This is particularly likely if the subscales differ in their strength of association with preceding steps, as uneven relationships can challenge the model's ability to represent a cohesive structure. Additionally, models with greater complexity may be more sensitive to sample size, as larger sample sizes generally provide more stable parameter estimates and reduce the likelihood of sampling error (Kline, 2016). Nonetheless, the applicability and advantages of the comprehensive model are re-evaluated in Study 2.

Overall, the findings from Study 1 provide strong initial support for the newly developed bystander intervention scale. The results align with the theoretical underpinnings of the bystander intervention model, while also highlighting potential variations in the intervention process based on specific types of actions. In Study 2, we re-evaluate the final scale's psychometric properties with a new sample, to confirm our findings, and determine the scale's construct validity by measuring the relationship between the five steps of the Bystander Intervention Model and theoretically relevant variables.

STUDY 2

Method

Participants. Participants were 240 Portuguese citizens (121 female, 116 male, 3 indicated "other"), ranging from 18 to 76 years old (M = 37.58. SD = 14.93), 57% with completed secondary education, 42% with higher education and less than 1% with basic education. The majority were employed (53%) or student (32%), and the remaining were unemployed (10%), retired (5%) and less than 1% were housewife.

Regarding the left-right political spectrum, the average score on a 7-point Likert-scale (1 = Left, 7 = Right), was close to the scale midpoint of 4, though leaning left (M = 3.33, SD = 1.50). Participants were also asked about perceived socioeconomic status compared to other citizens of the country where they live on a 7-point Likert-scale (1 = Very low, 7 = Very high; M = 3.87, SD = 1.04).

Procedure. We used the same data collection procedure as in Study 1. The study was conducted through Qualtrics, and participants were recruited via online platforms (e.g., university mailing list) to complete a survey on online hate speech towards immigrants. Participation was entirely voluntary, with five \in 20 vouchers raffled as compensation. After providing informed consent, participants completed sociodemographic questions (e.g., age, gender, education).

After completing the sociodemographic section and reading the definitions of hate crime and online hate speech, participants responded to the final version of the BIOHS-Immigrants scale developed in Study 1, along with related measures.

Measures. In order to examine the construct validity of the BIOHS-Immigrants scale we also included moral disengagement, victim blaming and self-efficacy measures. To minimize order effects, the BIOHS-Immigrants scale was presented first, followed by the related measures. This sequence ensured that responses to our scale were not influenced by subsequent constructs.

BIOHS-Immigrants scale. Participants were asked to respond to the final version of our scale with 26 items.

Moral disengagement. Participants answered a 24-item Moral Disengagement scale (based on Bandura et al., 1996) adapted to the context of online hate speech towards immigrants, integrating 6 of the 8 proposed mechanisms/dimensions (i.e., Moral justification, Displacement of responsibility, Diffusion of responsibility, Distorting consequences, Attribution of blame, Dehumanization; we excluded Euphemistic language and Advantageous comparison mechanisms due to the difficulty in adapting such items to the context of online hate speech; the full scale is available at OSM 7). All the 24 items (4 items per mechanism) were answered on a 7-point Likert-scale ranging from 1 (*Fully disagree*) to 7 (*Fully agree*). The scale revealed a good reliability (Cronbach's $\alpha = .97$; M = 1.72, SD = 1.11). Based on the one-factor solution proposed by the author (Bandura et al., 1996), we averaged the scores of all items to a Moral disengagement index, higher scores indicating higher moral disengagement regarding online hate speech towards immigrants.

Victim blaming. We used two items (based on Koehler & Weber's victim blaming scale, 2018) as a potential barrier to helping behaviour - participants' perceived diminished responsibility to intervene in case of hate speech, influenced by victim behaviour: "I am less likely to intervene facing online hate speech against an immigrant if I feel that s/he has done something to provoke such situation." and "If an immigrant has been aggressive or offensive to someone, I feel less responsible to intervene facing online hate speech against

the immigrant.". We averaged participants' responses into a Victim blaming index ($r = .65, p \le .001; M =$ 3.02, SD = 1.58).

Self-efficacy. Based on Banyard and colleagues' work (Bystander efficacy scale; 2007) we created a 6-item self-efficacy scale. Participants indicated their agreement or disagreement with each sentence on a 7-point Likert-scale (1 = *I fully disagree*; 7 = *I fully agree*): (1) "I feel that I am able to confront people who direct hate speech against immigrants."; (2) "I know what to do and what to say to help stop a situation of hate speech against immigrants.": (3) "I know what to do and say to prevent a situation of hate speech against immigrants."; (4) "I know what to do and say to reduce hate speech against immigrants."; (5) "I have the necessary skills to comfort/support an immigrant who has been a victim of hate speech."; (6) "I have the necessary skills to confront someone who is directing hate speech against immigrants.". A principal components factorial analysis conducted on these items extracted one single factor accounting for 70% of the total variance. We averaged the scores of all items into a Self-efficacy index (Cronbach's $\alpha = .91$; M =4.07, SD = 1.37).

Results

Confirmatory factor analysis (CFA). As in Study 1, we started by testing and comparing models reflecting both the theoretical (four-step solution with separate factors for Steps 2 and 3) and empirical (three-step solution with combined factor for Steps 2 and 3) structures, regarding the initial Steps of the model (i.e., Step 1 to 4). The model fits of the four-step and three-step solution can be seen in Table 3.

As in Study 1, the four-step solution model demonstrates a better overall fit to the data compared to the three-step solution model, further supporting its theoretical structure. In the three-step solution, the modification indices suggested correlating five error terms within the combined factor for Steps 2 and 3 $(e3 \leftrightarrow e4, e4 \leftrightarrow e7, e6 \leftrightarrow e8)$ and two error terms in Step 4 $(e10 \leftrightarrow e12)$. It is important to note that, according to the theoretical structure, items 4 and 7 (i.e., corresponding to error terms 4 and 7) belong to different dimensions (see OSM 3 for CFA and SEM results before modifications). However, since these items are grouped together in the three-step solution, we applied this correlation to maintain consistency with the procedure of correlating errors within the same dimension.

Based on the model fit results, we proceeded with the subsequent analyses, maintaining Steps 2 and 3 as separate factors.

As in Study 1, we tested five independent models corresponding to the sequential five steps of the bystander intervention, having as fifth and final step one of the five types of actions in each model.

As can be seen in Table 3, the five-step solution model with the Report through formal mechanisms subscale as the final step demonstrates an overall acceptable fit to the data, although the GFI (.892) and RMSEA (.091) fall slightly below commonly accepted thresholds for model fit (see Supplementary Figures at OSM 8). MI suggested correlating two error terms in Step 3 (e6 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7), two error terms in Step 4 (e9 < -> e7). > e11), and two error terms in Step 5 (e17 <-> e18).

The model with the Report through social media mechanisms subscale also demonstrated an acceptable fit to the data (MI suggested correlating error terms in Step 3 and Step 4: e6 <-> e8 and e10 <-> e12, respectively). Similarly, the models for Public actions (MI suggested correlating error terms in Step 3, e6 <-> e7, and in Step 4: e9 <-> e12, e10 <-> e12), Private actions (MI suggested correlating error terms in Step 3 and in Step 4: e6 <-> e8 and e10 <-> e12, respectively), and Aggressive response (MI suggested correlating error terms in Step 3, e5 <-> e7, e6 <-> e8, in Step 4, e10 <-> e12, and in Step 5, e25 <-> 26) showed acceptable fit indices.

We also tested the model with all Step 5 subscales (Model 6; Table 3).

As in Study 1, this model (MI suggested correlating error terms in Step 3, e5 <-> e7, e6 <-> e8, in Step 4, e10 <-> e12, and in the Aggressive response step, e25 <-> 26) shows a poorer fit across several key indices compared to the other models. Specifically, the fit indices, including GFI (.871) and NFI (.889), were below the recommended thresholds, suggesting that the model may not capture the structure of the data as effectively as models with individual Step 5 subscales. Moreover, the full model had the highest values for both AIC (702.44) and BIC (1019.18), indicating a less efficient fit compared to the other models. These results suggest that individual models may provide more robust and reliable fits compared to the full model, particularly due to the increased complexity when combining all subscales.

Structural equation modeling (SEM). Next, we tested the sequential steps of the bystander intervention model for each of the five types of actions.

As we can see in Table 4 (see Supplementary Figures at OSM 9), the models with the Report through social media mechanisms (MI suggested correlating error terms in Step 3, e5<-> e6, e6<-> e8, and in Step

4, e10 <-> e11, e10 <-> e12), Public actions (MI suggested correlating error terms in Step 3, e6<-> e8, and in Step 4, e9<-> e12, e10 <-> e12), and with the Private actions subscales (MI suggested correlating error terms in Step 3, e5<-> e6, e6<-> e8, and in Step 4, e10 <-> e12), provided an acceptable or good fit to the data and all of the direct paths were positive and statistically significant, showing that, as expected, all the steps were predicted by the previous step in the model.

The model with Report through formal mechanisms (MI suggested correlating error terms in Step 3, e5 <-> e6, e6 <-> e8, and in Step 4, e10 <-> e12) demonstrated a weak fit to the data, with GFI (.884), NFI (.899), and RMSEA (.095) falling below commonly accepted thresholds for model fit. Nevertheless, all direct paths were positive and statistically significant.

The model with the Aggressive response subscale provided a good fit to the data, with all of the direct paths being positive and statistically significant between the initial steps (MI suggested correlating error terms in Step 3, e6<-> e8, in Step 4, e10 <-> e12, and in Step 5, e25 <-> e26). However, the regression weight between latent variables Step 4 (Know how to help) and Step 5 (Aggressive bystander response) was non-significant. p = .381. As in Study 1, we tested an alternative sequential model that included a direct link between Step 3 and Step 5; however, Aggressive bystander response was not predicted by Accepting responsibility to help

We also tested this alternative sequential process (i.e., a direct link between Step 3 and Step 5) for the other models. As in Study 1, regarding the Report through social media mechanisms subscale, we observed that this type of action was better predicted by Accepting responsibility to help (Step 3) and results indicate that this model demonstrates a good fit, suggesting once again that Step 4 may not be essential in the intervention process for this specific action (see OSM 9): χ^2 (*df*) = 186.86 (83)***, CMIN/DF = 2.25, GFI = .910, NFI = .932, CFI = .961, RMSEA [CI] = .072 [.059, .086], AIC = 260.86, BIC = 389.64.

Regarding the remaining models, results indicate that including this direct link neither improves the model's fit for these actions nor alters the sequential process, supporting the need of maintaining the originally proposed sequential structure.

Finally, as in Study 1, we also tested the full model with all Step 5 subscales (Figure 2). We observed that all the direct paths were positive and statistically significant, showing that, as expected, all the steps were predicted by the previous step (note that all the proposed final steps were predicted by Step 4). MI suggested correlating error terms in Step 3, e6 <-> e8, in Step 4, e9 <-> e11, e9 <-> e12, e10 <-> e11, e10 <-> e12, and in Step 5 – Aggressive response, e25 <-> e26). However, the model presented a weak fit to the data (see Table 4).



Note: All path coefficients are standardized estimates.

 $p \le .05; \stackrel{\cdot}{=} p \le .01; \stackrel{\cdot}{=} p \le .001$

Figure 2. Simplified Graphical Representation of the SEM Model with All Step 5 Subscales in Study 2.

Reliability and descriptive statistics. As in Study 1, internal consistency reliability was appropriate for all the steps (subdimensions) of the bystander intervention model (see OSM 6 for reliability, descriptive statistics and bivariate correlations). As expected, each step of the bystander intervention model is positively correlated with the subsequent step, and Step 4 is positively correlated with the proposed actions corresponding to Step 5, except for Aggressive response.

Examining gender differences, and consistent with Study 1, we observe that women consistently scored higher than men across most steps of the bystander intervention model, except for Knowing how to help (Step 4), although the difference was not significant. Regarding Step 5, women scored higher than men in Report through social media mechanisms, indicating a greater likelihood to rely on social media reporting tools. Conversely, men scored slightly higher than women in Aggressive response, although the differences were not statistically significant. These findings suggest that gender may influence specific bystander behaviours, with women generally showing a stronger inclination to engage in prosocial and non-aggressive responses to online hate speech.

Consistent with the findings of Study 1, we observe that political orientation is negatively correlated with most steps of the bystander intervention model, while perceived social status shows no significant correlations with any of the steps. Age is positively correlated with Report through formal mechanisms and Public actions but negatively correlated with Report through informal mechanisms. Education shows a negative correlation only with Aggressive responses (see OSM 6).

Construct validity. To examine the construct validity of the new scale, we examined the relationship between the BIOHS-Immigrants scale and other theoretically related measures: moral disengagement, victim blaming, and self-efficacy (see OSM 6 for correlations between each Step of the BIOHS-Immigrants and the related measures). Overall, as expected, Moral disengagement and Victim blaming are negatively related with all subdimensions, except for Aggressive response, which is positively correlated with Moral disengagement and shows no significant relationship with Victim blaming. Finally, self-efficacy is positively related to all subdimensions, except for Aggressive response. These results give support for the theoretical validity of the BIOHS-Immigrants scale.

Discussion

As in Study 1, we tested five independent models corresponding to the five steps of the bystander intervention, with each model including one of the five types of actions as the fifth and final step. Overall, the CFA replicated the results of Study 1, demonstrating that the five-step solution model, with each of the five types of actions as the final step, provided a good fit to the data. These results further support the theoretical and empirical adequacy of the sequential five-step structure.

The SEM analyses also confirmed that, as expected, each step of the model was positively predicted by the preceding step, reflecting the sequential nature of the bystander intervention process. However, Aggressive response was an exception, as it was not predicted by Step 4 (Know how to help) nor Step 3 (Accepting responsibility to help). Moreover, as in Study 1, the alternative sequential model (i.e., including a direct link between Step 3 - Accepting responsibility to help and Step 5 - Taking action) for Report through social media mechanisms demonstrated a better fit to the data. This finding suggests that the sequential process outlined in the original model may not fully apply to this specific type of action.

As in Study 1, the SEM model with all the Step 5 subscales revealed that all proposed final steps were positively predicted by Step 4, aligning with the sequential decision-making framework originally proposed by Latané and Darley (1969, 1970). However, the model showed a weak fit to the data. Thus, consistent with Study 1, the results suggest that it may be more appropriate to consider five independent models, each corresponding to one of the proposed actions, rather than a single comprehensive model encompassing all potential actions.

To further assess the construct validity of the new scale, we examined its relationships with theoretically related constructs: moral disengagement, victim blaming, and self-efficacy. As hypothesized, the BIOHS-Immigrants scale was negatively correlated with both moral disengagement and victim blaming, and positively correlated with self-efficacy. These patterns support the theoretical premise that higher moral disengagement and victim blaming are barriers to prosocial intervention behaviours, whereas greater self-efficacy facilitates bystander intervention (e.g. Ferreira et al., 2020; Gini et al., 2020; Koehler & Weber, 2018). Notably, the Aggressive response dimension diverged from these patterns, being positively related to moral disengagement, but not related to victim blaming or self-efficacy.

Overall, these findings reinforce the validity and reliability of the BIOHS-Immigrants scale as a measure of bystander intervention in the context of online hate speech toward immigrants. The scale captures the sequential nature of decision-making processes, is theoretically grounded in existing models, and demonstrates strong construct validity through its relationships with related measures.

DISCUSSION

Across two studies, we provide first empirical support for a new scale of bystander intervention on online hate speech towards immigrants. Taken together, the results tend to confirm our initial hypotheses. Indeed, results support the proposed multifactorial structure of our scale, corresponding to the five steps of the

bystander intervention model (Latané & Darley, 1969, 1970), having as the final step (i.e., Implement intervention decision) five potential actions that individuals can engage in facing online hate speech, namely, Report through formal mechanisms (e.g., report to the police), Report through social media mechanisms (i.e., informal mechanisms, e.g., using Facebook report button), Public actions (e.g., defend the immigrant in a public comment in any social network platform), Private actions (e.g., defend the immigrant via private message to the offender), and Aggressive response (e.g., insult or threaten the offender).

As proposed by the bystander intervention theoretical framework, SEM showed that each of the initial steps of the model predicted the subsequent step. However, some inconsistencies emerged between studies regarding the relationship between Step 4 (Knowing how to help) and specific types of actions, namely Private actions and Aggressive responses. For instance, we observed a weak regression path between Step 4 and Aggressive responses in Study 1 ($\beta = .14$, p = .037), while no significant relationship was found in Study 2 ($\beta = .07$, p = .381). This suggests that Aggressive responses may not align with the sequential framework as originally proposed. Similarly, Private actions and Aggressive responses may be attributed to contextual or sample-specific factors, such as variations in participant demographics or levels of exposure to online hate speech. For instance, the majority of participants in Study 1 were students (64%) with a mean age of 28 years. In contrast, in Study 2, most participants were employed (53%), with only 32% being students, and the mean age was higher (38 years). These distinct sample characteristics may explain the discrepancies observed and the emergence of different sequential processes.

Additionally, across both studies, Report through social media mechanisms was consistently better predicted by Step 3 (Accepting responsibility to help), bypassing the need for Step 4.

Moreover, although SEM with all the proposed actions as final steps confirmed the expected predictive sequential effect in both studies, the overall model fit was weak. These findings suggest that future research may benefit from analysing the five types of actions in Step 5 as independent models, rather than relying on a complex, combined structure that includes all actions simultaneously.

Results also showed that, as expected, moral disengagement and victim blaming—concepts previously identified as barriers to bystander intervention (e.g. Gini et al., 2020; Koehler & Weber, 2018)— were negatively related to all the steps of the bystander intervention model. These findings support prior research indicating that moral disengagement and victim blaming inhibit the processes underlying the bystander intervention model, thereby diminishing the likelihood of progressing through its steps and engaging in helping behaviour. On the contrary, as proposed by previous research (e.g. Ferreira et al., 2020), self-efficacy seems to be aligned with the processes involved in this model, reflected in the positive relationship with all the steps, especially Step 4 (Know how to help).

Limitations and directions for future research

While the results across both studies provide evidence supporting the new scale and the proposed sequential model, they also highlight important limitations and inconsistencies that warrant careful consideration. Below, we discuss some of these limitations and propose directions for future research.

One of the limitations that can be raised is related to the use of two-items factors, as is the case of the subdimensions corresponding to Steps 1 and 2, and some of the proposed potential actions (Public and Private actions). Although two-item factors are considered acceptable when the items are highly correlated (r > .70; Yong & Pearce, 2013), some correlations in Study 1, such as Step 2 (r = .64) and Private actions (r = .68), were slightly below this threshold. However, in Study 2, these correlations exceeded the recommended value, strengthening the reliability of these factors. Moreover, a factor with two items is acceptable when there are strong theoretical and practical reasons, particularly given the multifactorial nature of our scale. Indeed, as our scale has several dimensions (a total of 26 items), keeping it as short as possible helps to reduce participant fatigue, frustration and boredom, and increases its usability, especially when used alongside other scales in longer surveys that can contribute to participants' withdrawal.

Another potential limitation concerns the reliance on self-reported intentions rather than actual helping behaviours as the final step in the bystander intervention model. Measuring real behaviours is more challenging but would provide more direct evidence of the model's application. Nonetheless, based on the Theory of Planned Behaviour (Ajzen, 2020), intentions are a strong and reliable predictor of future behaviour. Future studies should aim to measure actual behaviours, such as providing participants with an opportunity to report hate speech in a real or simulated online environment during the experiment.

Moreover, the influence of social desirability bias, particularly in responses related to prosocial intentions and actions, could be considered in future studies. This bias might lead participants to overreport socially desirable behaviours or underreport undesirable ones, especially given the sensitivity of the topic of online hate speech. Thus, to further strengthen future applications of the scale, research could incorporate methods to control for social desirability.

Future research could also explore how individual differences in internet usage goals (e.g., entertainment vs. news consumption) and exposure to online hate speech influence participants' responses. Perceptions of the prevalence of online hate speech may vary significantly based on internet usage profiles, including usage patterns, frequency, and primary goals. Investigating these differences by incorporating detailed measures—such as platforms frequented, types of interactions, hours spent online daily, and usage goals (e.g., entertainment, professional networking, or information seeking)— could offer valuable insights. This approach would provide a more nuanced understanding of how exposure to online hate speech interacts with individual characteristics to shape bystander intervention behaviours.

It is also important to note that in the items related to the Implement Intervention stage (Step 5), we included references to national institutions (e.g., police, governmental organizations) that play a role in responding to online hate speech. These references were included to reflect culturally relevant mechanisms for addressing hate speech. While these items were carefully designed, they were not pretested to assess participants' associations with these institutions. Future research could explore whether participants accurately associate these institutions with their intended roles in addressing hate speech, ensuring that responses reflect genuine beliefs rather than potential misunderstandings.

Finally, although we included a range of relevant and common actions in our scale, there are undoubtedly other potential bystander actions that were not captured. Future studies could expand the range of actions studied, ensuring the model's comprehensiveness and adaptability to diverse contexts.

Theoretical and empirical implications

The theoretical framework of the bystander effect and the bystander intervention model has been applied to a wide range of contexts, beyond the original focus on emergency situations. Namely, in the context of bullying (e.g., Nickerson et al., 2014), sexual assault or sexual violence (e.g., Bennett et al., 2014; Kettrey & Marx, 2020), helping behaviour towards refugees (Albayrak-Aydemir & Gleibs, 2021), computer-mediated communication (Markey, 2000), cyberbullying (e.g., You & Lee, 2019), and Islamophobic online hate speech (Obermaier et al., 2021). Although the sequential process proposed by the model was not assessed in some of these studies, the widespread application of the bystander intervention model underscores its relevance and adaptability across diverse contexts.

In the present research, we explored another critical context: online hate speech targeting immigrants. Our findings provide empirical support for the theoretical framework of the bystander intervention model within this specific domain. More importantly, we have developed a new instrument that has the potential to serve as a valuable tool for future research in this area. While the proposed scale was specifically designed to address online hate speech directed at immigrants, it has broader applicability. The scale can be easily adapted to study online hate speech targeting other socially vulnerable or minority groups, as the targets of hate vary widely. This flexibility makes the instrument a promising contribution to advancing research on bystander intervention across different forms of discrimination and online hate speech.

Concluding remarks

Our work is the first attempt, as far as we know, to apply the sequential process of bystander intervention model (Latané & Darley, 1969) to online hate speech towards immigrants. Specifically, we developed a new measurement to assess this process, incorporating five distinct types of potentially helping behaviours in response to online hate speech. By doing so, our work not only advances research on the bystander effect and bystander intervention but also makes a significant contribution to understanding the processes underlying the perpetuation and normalization of hate toward socially vulnerable or minority groups in digital spaces. Importantly, we provide researchers with a new, theory-based tool to investigate these processes, enabling a deeper exploration of how online hate speech can be challenged and mitigated in a context where it is frequently propagated and socially tolerated.

REFERENCES

Ajzen, I. (2020). The theory of planned behavior: Frequently asked questions. *Human Behavior and Emerging Technologies*, *2*(4), 314-324. https://doi.org/10.1002/hbe2.195

- Albayrak-Aydemir, N., & Gleibs, I. H. (2021). Measuring global bystander intervention and exploring its antecedents for helping refugees. *British Journal of Psychology*, 112(2), 519-548. https://doi.org/10.1111/bjop.12474
- Álvarez-Benjumea, A. (2023). Uncovering hidden opinions: social norms and the expression of xenophobic attitudes. *European Sociological Review*, 39(3), 449-463. https://doi.org/10.1093/esr/jcac056

- Bandura, A. (1998). Personal and collective efficacy in human adaptation and change. In J. G. Adair, D.
 Bélanger, & K. L. Dion (Eds.), *Advances in psychological science, Vol. 1. Social, personal, and cultural aspects* (pp. 51–71). Psychology Press/Erlbaum (UK) Taylor & Francis.
- Bandura, A., Barbaranelli, C., Caprara, G. V., & Pastorelli, C. (1996). Mechanisms of moral disengagement in the exercise of moral agency. *Journal of Personality and Social Psychology*, *71*(2), 364–374. https://doi.org/10.1037/0022-3514.71.2.364
- Banyard, V. L. (2008). Measurement and correlates of prosocial bystander behavior: the case of interpersonal violence. *Violence & Victims, 23*(1), 83–97. https://doi.org/10.1891/0886-6708.23.1.83
- Banyard, V. L., Moynihan, M. M., & Plante, E. G. (2007). Sexual violence prevention through bystander education: An experimental evaluation. *Journal of Community Psychology*, 35(4), 463-481. https://doi.org/10.1002/jcop.20159
- Banyard, V. L., Plante, E. G., & Moynihan, M. M. (2005). Rape prevention through bystander education: Bringing a broader community perspective to sexual violence prevention (Final report to NIJ 208701). US Department of Justice. https://www.ojp.gov/pdffiles1/nij/grants/208701.pdf
- Barlińska, J., Szuster, A., & Winiewski, M. (2013). Cyberbullying among adolescent bystanders: Role of the communication medium, form of violence, and empathy. *Journal of Community & Applied Social Psychology*, 23(1), 37-51. https://doi.org/10.1002/casp.2137
- Bastiaensens, S., Vandebosch. H., Poels, K., Van Cleemput, K., DeSmet. A., & De Bourdeaudhuij. I. (2014). Cyberbullying on social network sites. An experimental study into bystanders' behavioural intentions to help the victim or reinforce the bully. *Computers in Human Behavior, 31*, 259-271. https://doi.org/1.1016/j.chb.2013.1.036
- Bennett, S., Banyard, V. L., & Garnhart, L. (2014). To act or not to act, that is the question? Barriers and facilitators of bystander intervention. *Journal of Interpersonal Violence, 29*(3), 476-496. https://doi.org/10.1177/0886260513505210
- Brown, T. A. (2015). Confirmatory factor analysis for applied research (2nd ed.). The Guilford Press.
- Clark, M., & Bussey, K. (2020). The role of self-efficacy in defending cyberbullying victims. *Computers in Human Behavior, 109*, 106340. https://doi.org/10.1016/j.chb.2020.106340
- Costello, M., Hawdon, J. E., & Cross, A. (2017). Virtually standing up or standing by? Correlates of enacting social control online. *International Journal of Criminology and Sociology, 6*, 16-28. http://hdl.handle.net/10919/81716
- Dreißigacker, A., Müller, P., Isenhardt, A., & Schemmel, J. (2024). Online hate speech victimization: consequences for victims' feelings of insecurity. *Crime Science*, *13*(1), 4. https://doi.org/10.1186/s40163-024-00204-y
- Ferreira, P. C., Simão, A. V., Paiva, A., & Ferreira, A. (2020). Responsive bystander behaviour in cyberbullying: a path through self-efficacy. *Behaviour & Information Technology*, 39(5), 511-524. https://doi.org/10.1080/0144929X.2019.1602671
- FRA (2017). Second European Union Minorities and Discrimination Survey: Technical Report. https://fra.europa.eu/en/publication/2017/second-european-union-minorities-anddiscrimination-survey-technical-report
- FRA (2021). Encouraging Hate Crime Reporting The role of Law Enforcement and Other Authorities. https://fra.europa.eu/en/publication/2021/hate-crime-reporting
- Gini, G., Thornberg, R., & Pozzoli, T. (2020). Individual moral disengagement and bystander behavior in bullying: The role of moral distress and collective moral disengagement. *Psychology of Violence*, 10(1), 38–47. https://doi.org/10.1037/vio0000223
- Hair, J. F., Black, W. C., Babin, B. J., & Anderson, R. E. (2013). *Multivariate data analysis* (7th ed.). Pearson Higher Ed.
- Hinkin, T. R., Tracey. J. B., & Enz, C. A. (1997). Scale construction: Developing reliable and valid measurement instruments. *Journal of Hospitality & Tourism Research*, 21(1), 100-120. https://doi.org/10.1177/109634809702100108
- Hu, L.-T. & Bentler, P. M. (1999). Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives. *Structural Equation Modeling*, 6(1), 1– 55. https://doi.org/10.1080/10705519909540118
- Jubany, O., & Roiha, M. (2016). Backgrounds, experiences and responses to online hate speech: a comparative cross-country analysis. https://sosracismo.eu/wp-

content/uploads/2016/07/Backgrounds-Experiences-and-Responses-to-Online-Hate-Speech.pdf

Kettrey, H. H., & Marx, R. A. (2020). Effects of bystander sexual assault prevention programs on promoting intervention skills and combatting the bystander effect: a systematic review and meta-analysis. *Journal of Experimental Criminology*, 17, 343-367. https://doi.org/10.1007/s11292-020-09417-y Kline, R. B. (2016). *Principles and practice of structural equation modeling* (4th ed.). Guilford publications.

- Koehler, C., & Weber, M. (2018). "Do I really need to help?!" Perceived severity of cyberbullying, victim blaming, and bystanders' willingness to help the victim. Cyberpsychology: *Journal of Psychosocial Research on Cyberspace*, *12*(4). https://doi.org/10.5817/CP2018-4-4
- Latané, B., & Darley, J. M. (1969). "Bystander Apathy". *American Scientist*, *57*(2), 244-268. https://doi.org/10.1037/h0026570
- Latané, B., & Darley, J. M. (1970). *The unresponsive bystander: Why doesn't he help?*. Appleton-Century Crofts.
- Latané, B., & Rodin, J. (1969). A lady in distress: Inhibiting effects of friends and strangers on bystander intervention. *Journal of Experimental Social Psychology*, *5*(2), 189-202.
- Machackova, H. (2020). Bystander reactions to cyberbullying and cyberaggression: individual, contextual, and social factors. Current opinion in psychology, 36, 130-134. https://doi.org/10.1016/j.copsyc.2020.06.003
- Markey, P. M. (2000). Bystander intervention in computer-mediated communication. *Computers in Human Behavior*, *16*(2), 183-188. https://doi.org/10.1016/S0747-5632(99)00056-4
- Marsh, H. W., & Hocevar, D. (1985). Application of confirmatory factor analysis to the study of selfconcept: First- and higher-order factor models and their invariance across groups. *Psychological Bulletin*, *97*(3), 562–582. https://doi.org/10.1037/0033-2909.97.3.562
- Müller, K., & Schwarz, C. (2018, May). Fanning the flames of hate: Social media and hate crime. (Working Paper No. 373).

https://warwick.ac.uk/fac/soc/economics/research/centres/cage/manage/publications/373-2018_schwarz.pdf

- Murrell, A. J. (2021). Why someone did not stop them? Aversive racism and the responsibility of bystanders. *Equality, Diversity and Inclusion, 40*(1), pp. 60-73. https://doi.org/10.1108/EDI-07-2020-0191
- Nickerson, A. B., Aloe, A. M., Livingston, J. A., & Feeley, T. H. (2014). Measurement of the bystander intervention model for bullying and sexual harassment. *Journal of Adolescence*, *37*(4), 391–400. https://doi.org/10.1016/j.adolescence.2014.03.003
- Obermaier, M., Schmuck, D., & Saleem, M. (2021). I'll be there for you? Effects of Islamophobic online hate speech and counter speech on Muslim in-group bystanders' intention to intervene. *New Media & Society*. https://doi.org/10.1177/14614448211017527
- Osborne, J. W. (2014). Best Practices in Exploratory Factor Analysis. *CreateSpace Independent Publishing*. ISBN-13: 978-1500594343, ISBN-10:1500594342.
- Palmer, S. B., Cameron, L., Rutland. A., & Blake, B. (2017). Majority and minority ethnic status adolescents' bystander responses to racism in school. *Journal of Community & Applied Social Psychology*, 27(5), 374-380. https://doi.org/10.1002/casp.2313
- Papcunová, J., Martončik, M., Fedáková, D., Kentoš, M., Bozogáňová, M., Srba, I., ... & Adamkovič, M. (2021). Hate speech operationalization: a preliminary examination of hate speech indicators and their structure. *Complex & Intelligent Systems*, 1-16. https://doi.org/10.1007/s40747-021-00561-0
- Pinto, I. R., Carvalho, C. L., Dias, C., Lopes, P., Alves, S., de Carvalho, C., & Marques, J. M. (2020). A path toward inclusive social cohesion: The role of European and national identity on contesting vs. Accepting European migration policies in Portugal. *Frontiers in Psychology*, 11, Article 1875. https://doi.org/10.3389/fpsyg.2020.01875
- Pinto, I. R., Carvalho, C. L., Magalhães, M., Alves, S., Bernardo, M., Lopes, P., Carvalho, C. (2023). *Understanding the rise in online hate speech in Portugal and Spain: a gap between occurrence and reporting* (Issue Brief). "la Caixa" Foundation. https://oobservatoriosocial.fundacaolacaixa.pt/en/-/compreender-o-crescimento-do-discurso-de-odio-online-em-portugal-e-espanha-um-hiato-entrea-ocorrencia-e-a-denuncia
- Rudnicki, K., Vandebosch, H., Voué, P., & Poels, K. (2022). Systematic review of determinants and consequences of bystander interventions in online hate and cyberbullying among adults. *Behaviour & Information Technology*, 1-18. https://doi.org/10.1080/0144929X.2022.2027013
- SELMA (2019). Hacking Online Hate: Building an Evidence Base for Educators.
- https://hackinghate.eu/assets/documents/hacking-online-hate-research-report-1.pdf Siapera, E., Moreo, E., & Zhou, J. (2018). *Hate Track: Tracking and Monitoring Online Racist Speech*
- (Technical Report of HateTrack Project). https://www.ihrec.ie/documents/hatetrack-trackingand-monitoring-racist-hate-speech-online/
- Sjögren, B., Thornberg, R., Wänström, L., & Gini, G. (2021). Associations between students' bystander behavior and individual and classroom collective moral disengagement. *Educational Psychology*, 41(3), 264-281. https://doi.org/10.1080/01443410.2020.1828832

- Stewart, K., Pedersen, A., & Paradies, Y. (2014). It's always good to help when possible, BUT...: Obstacles to bystander anti-prejudice. *The International Journal of Diversity in Education, 13*(3). 39-53. https://doi.org/10.18848/2327-0020/CGP/v13i03/40045
- Suler, J. (2004). The Online Disinhibition Effect. *CyberPsychology & Behavior*, 7(3), 321–326. https://doi.org/10.1089/1094931041291295
- Thornberg, R., Wänström, L., Elmelid, R., Johansson, A., & Mellander, E. (2020). Standing up for the victim or supporting the bully? Bystander responses and their associations with moral disengagement, defender self-efficacy, and collective efficacy. *Social Psychology of Education, 23*(3), 563-581. https://doi.org/10.1007/s11218-020-09549-z
- UK Safer Internet Centre (2016). Creating a better internet for all: Young people's experiences of online empowerment + online hate.

https://childnetsic.s3.amazonaws.com/ufiles/SID2016/Creating%20a%20Better%20Internet%2 0for%20All.pdf

- United Nations (UN, 2019). United Nations Strategy and Plan of Action on Hate Speech. https://www.un.org/en/genocideprevention/documents/advising-andmobilizing/Action_plan_on_hate_speech_EN.pdf
- Vega Macías, D. (2021). The COVID-19 pandemic on anti-immigration and xenophobic discourse in Europe and the United States. *Estudios Fronterizos, 22*. https://doi.org/10.21670/ref.2103066
- Waldron, J. (2012). *The harm in hate speech*. Harvard University Press. https://doi.org/10.4159/harvard.9780674065086
- Walters, M. A., Brown, R., & Wiedlitzka, S. (July 2016). Causes and Motivations of Hate Crime (Research Report 102). Equality and Human Rights Commission Research Report 102 (2016); ISBN 978-1-84206-678-2, Available at SSRN: https://ssrn.com/abstract=2918883
- West, S. G., Taylor, A. B., & Wu, W. (2012). Model fit and model selection in structural equation modeling. In R. H. Hoyle (Ed.), *Handbook of structural equation modeling* (pp. 209–231). The Guilford Press.
- Whittaker, T. A. & Schumacker, R. E. (2022). *A beginner's guide to structural equation modeling* (5th ed.). Routledge.
- Whittaker, T. A. (2012). Using the modification index and standardized expected parameter change for model modification. *The Journal of Experimental Education*, 80(1), 26-44. https://doi.org/10.1080/00220973.2010.531299
- Yong, A. G., & Pearce, S. (2013). A beginner's guide to factor analysis: Focusing on exploratory factor analysis. *Tutorials in Quantitative Methods for Psychology*, 9(2), 79-94. https://doi.org/10.20982/tqmp.09.2.p079
- You, L. & Lee, Y.H. (2019). The bystander effect in cyberbullying on social network sites: Anonymity, group size, and intervention intentions. *Telematics and Informatics*, 45. https://doi.org/10.1016/j.tele.2019.101284
- Zapata, J., Sulik, J., von Wulffen, C., & Deroy, O. (2024). Bystanders' collective responses set the norm against hate speech. *Humanities and Social Sciences Communications*, *11*(1), 1-13. https://doi.org/10.1057/s41599-024-02761-8

CRediT AUTHORSHIP CONTRIBUTION STATEMENT

Catarina L. Carvalho: Conceptualization; Data Curation; Formal analysis; Investigation; Methodology; Writing - Original Draft; Writing - Review & Editing. **Isabel R. Pinto**: Conceptualization; Funding acquisition; Investigation; Methodology; Writing - Review & Editing. **Sara Alves**: Conceptualization; Methodology; Writing - Review & Editing; **Márcia Bernardo**: Conceptualization; Methodology; Writing - Review & Editing.

ACKOWLEDGMENTS

This work was funded by "la Caixa" Foundation and the Portuguese Foundation for Science and Technology (FCT) (SR20 - Social Research 2020; ref. SR20-00136).

SUPPLEMENTARY MATERIAL

Supplementary materials are available online at https://osf.io/qg2ky/

History of the manuscriptReceived24/05/2024Accepted26/02/2025Published (online)-Published04/07/2025