

Supplementary Information for Emotional persistence in online chatting communities

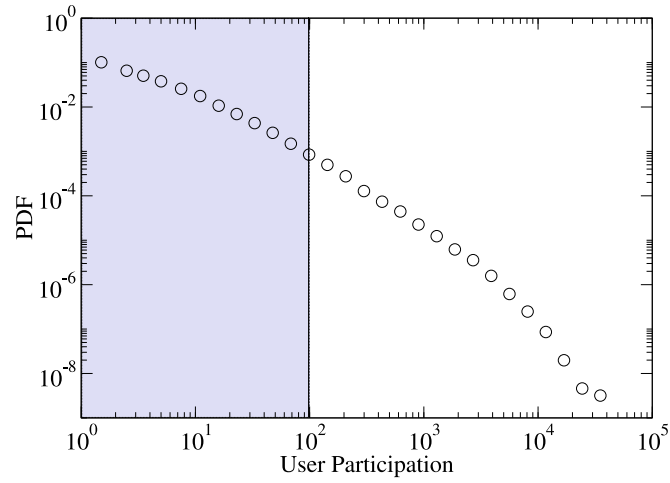
Antonios Garas, David Garcia, and Frank Schweitzer*
Chair of Systems Design, ETH Zurich, Kreuzplatz 5, 8032 Zurich, Switzerland

Marcin Skowron
Austrian Research Institute for Artificial Intelligence, Freyung 6/6, 1010 Vienna, Austria

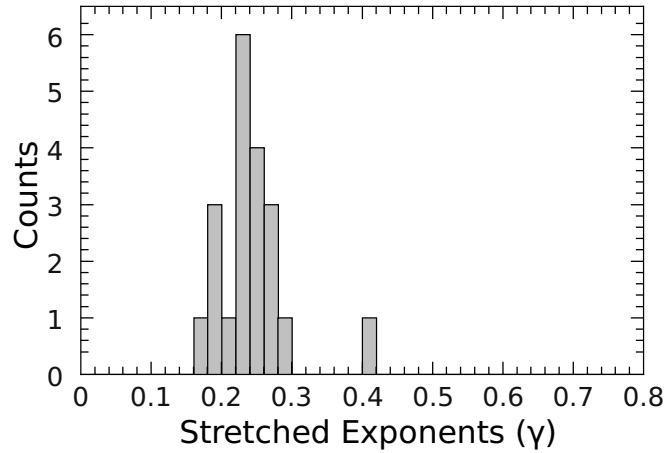
I. SUPPLEMENTARY FIGURES

- Supplementary Figure S1: Distribution of the user participation.
- Supplementary Figure S2: Histogram of stretched exponents.
- Supplementary Figure S3: DFA and autocorrelation analysis of real IRC channel activity
- Supplementary Figure S4: Example of persistent and anti-persistent time series.
- Supplementary Figure S5: DFA fluctuation functions.
- Supplementary Figure S6: Dependence of the Hurst exponent on the total activity of each user.
- Supplementary Figure S7: Dependence of the Hurst exponent on the length of the time series.
- Supplementary Figure S8: DFA fluctuation functions for different segments of the time series.

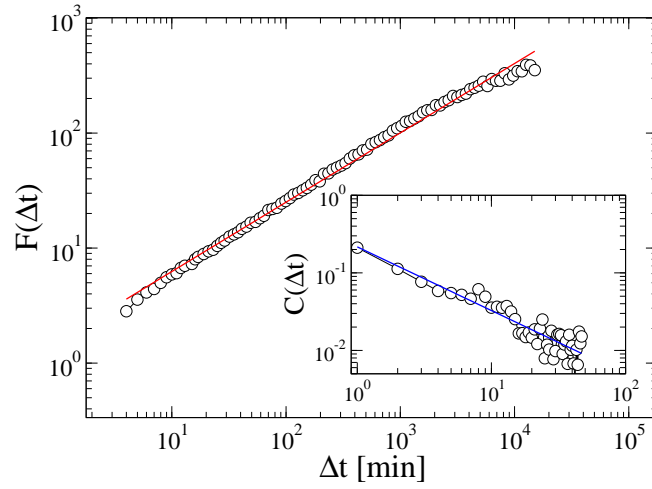
*Electronic address: fschweitzer@ethz.ch



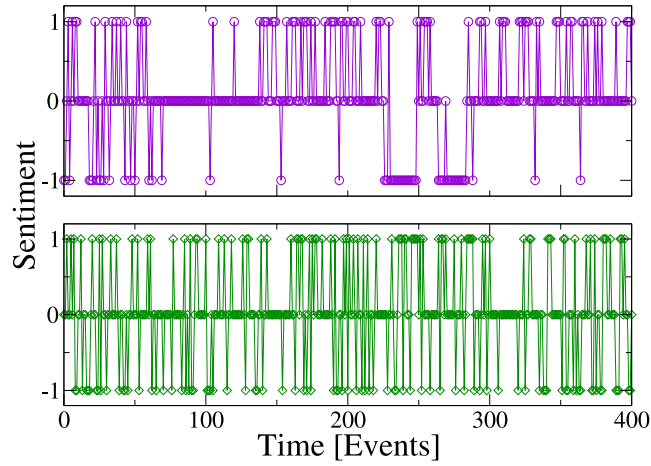
Supplementary Figure S 1: Distribution of the user participation in terms of the total number of posts entered by every user. The distribution is broad, and it is clear that most of the users contribute only a small number of posts. The shaded area shows the part of the user activity that is excluded from the DFA analysis in order to improve the statistical reliability of the results.



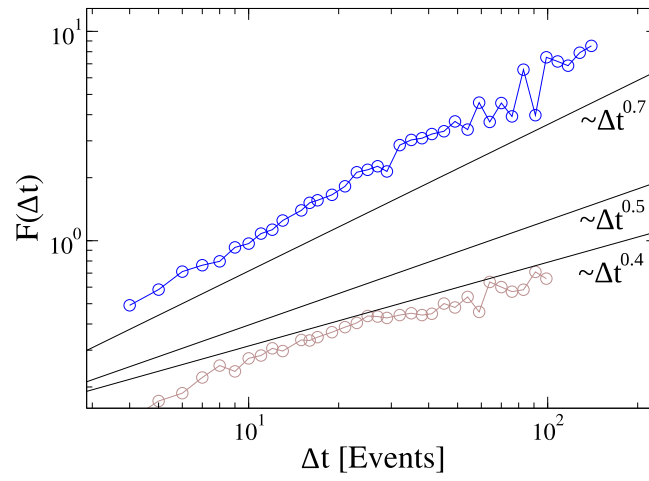
Supplementary Figure S 2: Histogram of stretched exponents obtained by fitting a stretched exponential function to the rescaled inter-event time of each individual channel separately. The exponents are concentrated around the mean value $\langle \gamma \rangle = 0.21 \pm 0.05$, obtained using only the regression results with $p < 0.001$, as explained in the text.



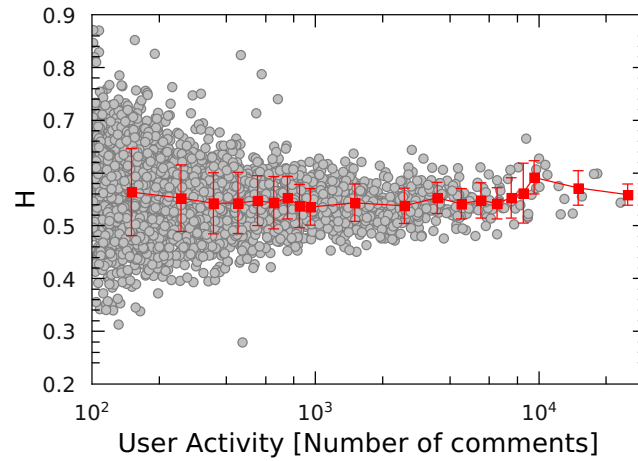
Supplementary Figure S 3: DFA fluctuation function calculated using the inter-event times of a real IRC channel. The Hurst exponent obtained is $H_\omega \simeq 0.6$, suggesting the existence of log term correlations in the time series. The origin of such correlations could be due to synchronized burst of activity leading to persistent dependencies over different time scales, or due to the broad distribution of inter-event times, or to a combination of both. The existence of dependencies in the activity is highlighted by a power law decaying autocorrelation function (Inset), with exponent $\nu_\omega \simeq 0.82$. The Hurst exponent is in scaling relation with the correlation exponent, given by $\nu_\omega = 2 - 2H_\omega$ [1].



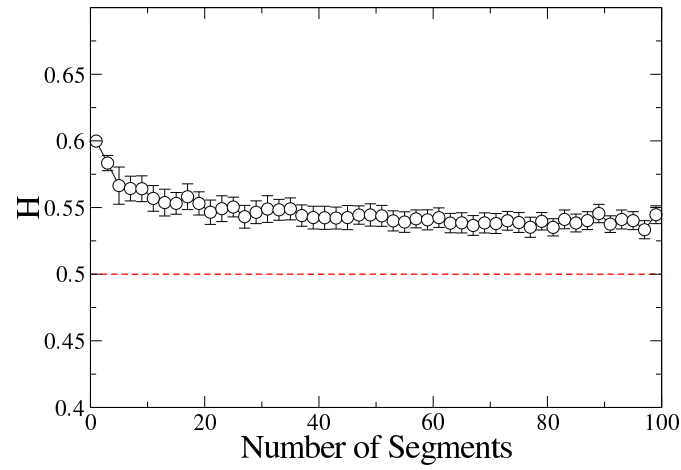
Supplementary Figure S 4: Time series showing examples of the sentiment expression for two real users. Top: An example of persistent sentiment time series. Bottom: An example of anti-persistent sentiment time series.



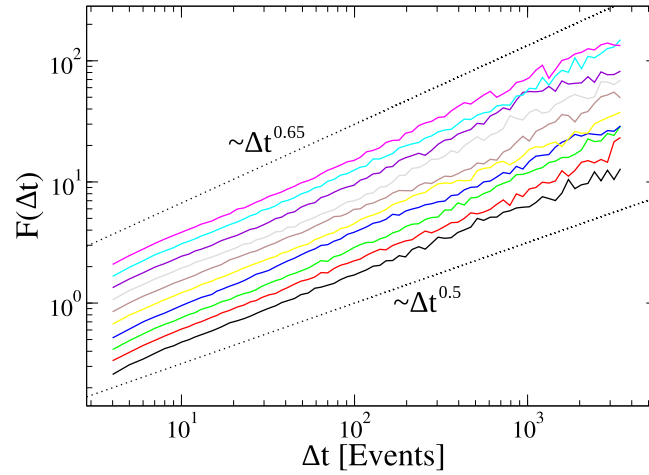
Supplementary Figure S 5: DFA fluctuation functions calculated for a persistent and an anti-persistent sentiment time-series. The solid lines are guides to the eye.



Supplementary Figure S 6: Dependence of the Hurst exponent on the total activity of each user. The mean value of H does not show any noticeable dependence on the activity but some large heterogeneity on the values of H for users with low activity is apparent.



Supplementary Figure S 7: Dependence of the Hurst exponent on the length of the time series. We divide the expression time series of an active user into various segments and apply the DFA method to these segments. A small dependence on the length of the segments is observed, but the overall behavior of the user remains consistent. The error bars show the standard error of the mean. The total number of posts contributed by this user is 18.142, and the maximum number of segments we used was 100 of length 181.



Supplementary Figure S 8: DFA fluctuation functions obtained for different segments of the time series describing the sentiment of a real IRC channel. It is clear that the persistence holds for all the segments analysed. The dashed lines are guides to the eye.

II. DATA

Our data set includes the annotated and anonymized logs from public Internet Relay Chat (IRC) channels of EFNET [5]. In particular, the data consists of consecutive daily recordings for 20 IRC channels for the period: 4-04-2006 - 17-05-2006. The general topics of discussions on these channels, as indicated by the IRC channel names, include: music, sports, casuals chats, business, politics and topics related to computers, operating systems or specific computer programs. The data were anonymized by substituting the real userIDs and the IRC channel names with generic number references. Subsequently, the data were annotated according to:

- **Sentiment classification**

As described in the "Methods section" of the article, our emotional classification is based on the SentiStrength classifier [2], which provides two scores for positive (called positiveArousal) and negative (called negativeArousal) content. For example, the text "*I love you*" according to SentiStrength has positiveArousal 3 and negativeArousal -1, while the text "*I'm very sad*" has positiveArousal 1 and negativeArousal -5.

From these two scores, we calculate a polarity measure (called sentimentClass) using the sign of the difference of the positive and negative scores. This measure takes the values +1, -1, and 0, and it provides an approximation to detect positive, negative and neutral posts respectively. Under this approach, the sentimentClass of the first text would be +1 indicating a positive text, while the sentimentClass of the second text would be -1 indicating a negative text.

- **Affective, cognitive and linguistic categories**

This annotation is based on the Linguistic Inquiry and Word Count - LIWC [3], and it results to a classification of words along 64 linguistic, cognitive, and affective categories.

- **Dialog act classification**

With this annotation we classified the text into 15 dialog act classes that are based in the following taxonomy: Accept, Bye, Clarify, Continuer, Emotion, Emphasis, Greet, No Answer, Other, Reject, Statement, Wh-Question, Yes Answer, Yes/No Question, Order [4]. Utterances that contained a url link, a empty utterances or utterances that did not include any ASCII characters were replaced by a "[url-link]" or "[empty-line]" tags, respectively.

Data availability

The data are freely available for research purposes. They are provided as supplementary material in a compressed "zip" file at <http://www.sg.ethz.ch/downloads/Data>. If you have any problems accessing them, please contact the authors.

Data structure, and the naming convention

In the zip file each folder contains annotations of each one of the 20 IRC channels. The file names correspond to the date, and their extensions represent the type of annotation they provide. The general internal structure of every file is as follows:

```
[timestamp] <anonymized-user-ID> a file-type specific annotation
```

More specifically, the type of information provided by every file is the following:

file extension: ".sent"

```
[time-stamp] <userID> | sentimentClass | positiveArousal | negativeArousal |
[03:45] <3032> | 0 | 1 | -1 |
```

file extension: ".liwc"

```
[time-stamp] <userID> liwcCategory1:liwcCategory2:liwcCategory3
[03:45] <3032> Affect:Posemo:Assent
```

file extension: ".da"

```
[time-stamp] <userID> dialogActClass
[03:45] <3032> Emotion
```

III. MODEL DETAILS

In order to understand how each one affects the ratio of emotion polarities in the posts and the user and conversation persistence, we performed a large set of simulations using different combination of parameter values for the model described in Section "An agent-based model for chatroom users". For each combination of values we performed run 10 simulation sets, and the dependencies of the collective behavior of the chatroom versus individual parameters are shown in Supplementary Figures S9-S11.

In Supplementary Figure S9 is summarized the effect in the ratio of positive, negative and neutral posts due to the change in some of the parameters. A higher amplitude of the stochastic influence implies a lower frequency of neutral posts, splitting the rest equally among positive and negative. Due to the high stochasticity of values like $A_v = 0.4$, the community just behaves randomly with almost even ratios of positive, negative and neutral. Increasing b , c , or decreasing the decay of the field γ_h , the influence of the conversation in the individual valence increases, leading to higher values of emotional posts regardless of their polarity. An increase of the absolute value of the expression thresholds V_{\pm} yields a lower frequency of the corresponding polarity, as expected.

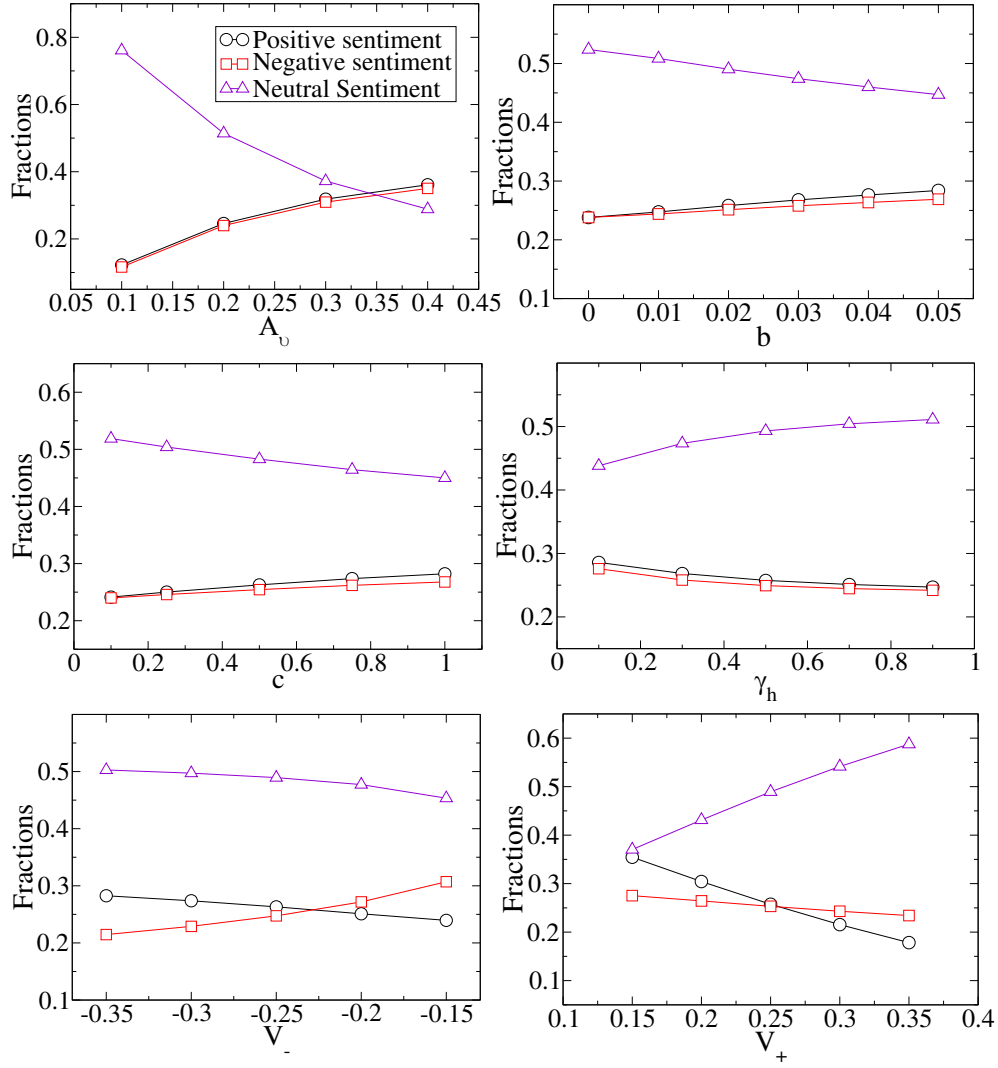
In terms of valence decay, γ_v , we simulated two possible cases. The case $\gamma_v = 0.1$ represents a virtual study of the dynamics of mood, as a slower, conscious process that influences the overall emotional state. The second case, $\gamma_v = 0.5$ results to a faster decay more representative of the dynamics of core affect, or fast emotional states. Supplementary Figure S10 shows the distributions of conversation and individual persistence for all the simulations with the ranges of values for the rest of the parameters. We find the case of $\gamma_v = 0.5$ closer to reality as observed in IRC channels, where persistence are significant but not as strong as they would be for the other case.

For each simulated case, we calculated the persistence of each individual expression as well as the persistence of the whole conversation. We find that increasing levels of amplitude of the stochastic component of the valence leads to slightly higher average individual persistence, but does not affect much the overall conversation persistence. Similar to the case of the polarity fractions, larger values of b , c , or lower values of γ_h have the effect of increasing persistence, as the coupling induced by the conversation is stronger. Similarly, higher values of the thresholds lead to lower conversation persistence due to the higher probability of neutral expression.

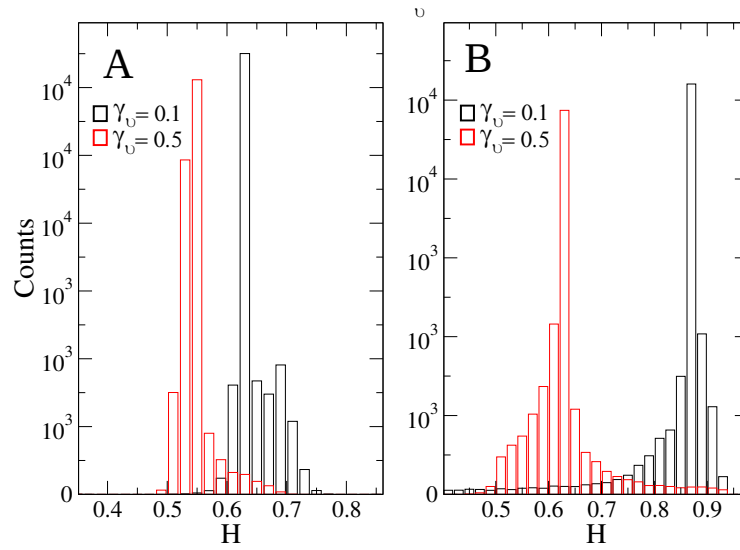
Given this behavior of the model from, we focused on a particular set of values to simulate conversations similar to actual IRC chats. We used 10000 agents in a conversation lasting 45000 time units, and performed 10 realizations of the model using the following set of parameters:

$$V_- = -0.15, V_+ = 0.05, \gamma_v = 0.2, A_v = 0.2, b = 0.01, c = 0.05, \gamma_h = 0.9$$

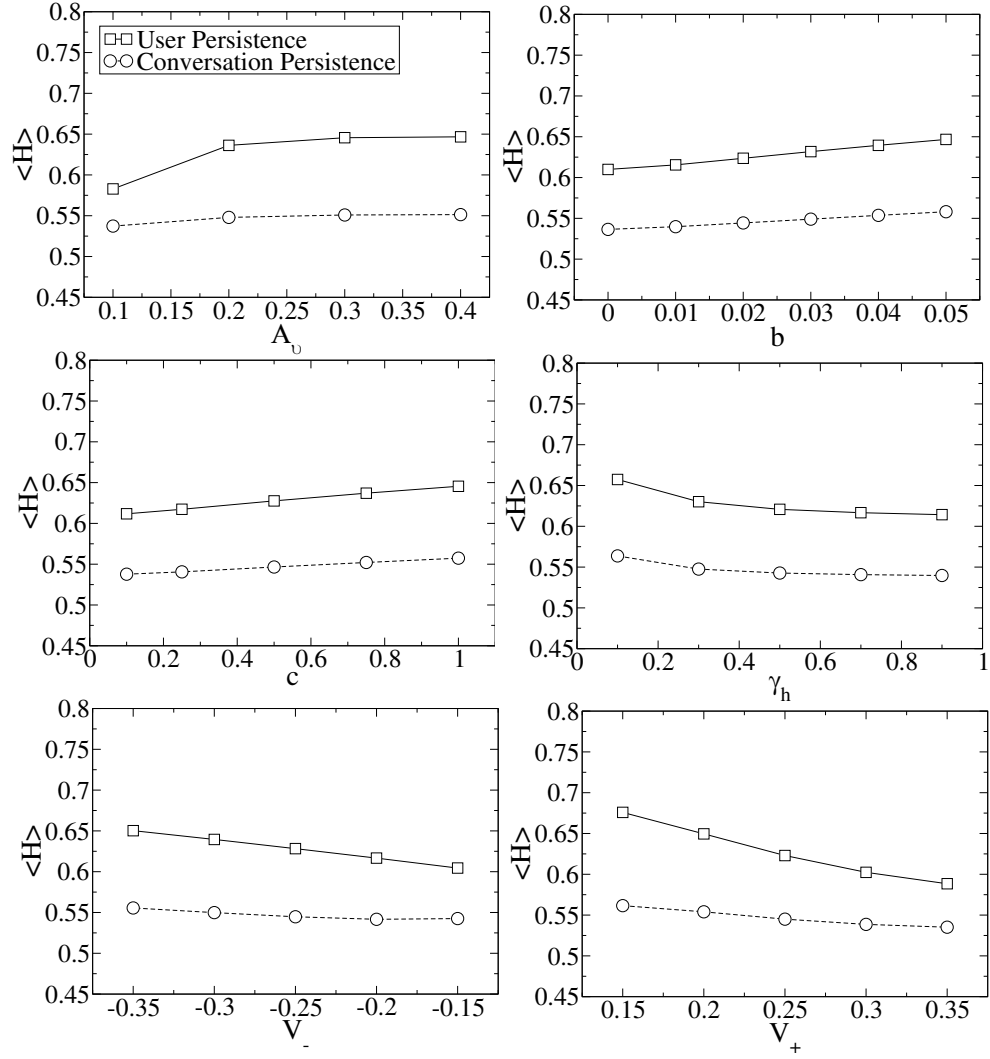
The results of an extensive set of simulations with these parameters are shown, and discussed in Section "An agent-based model for chatroom users" of the main text.



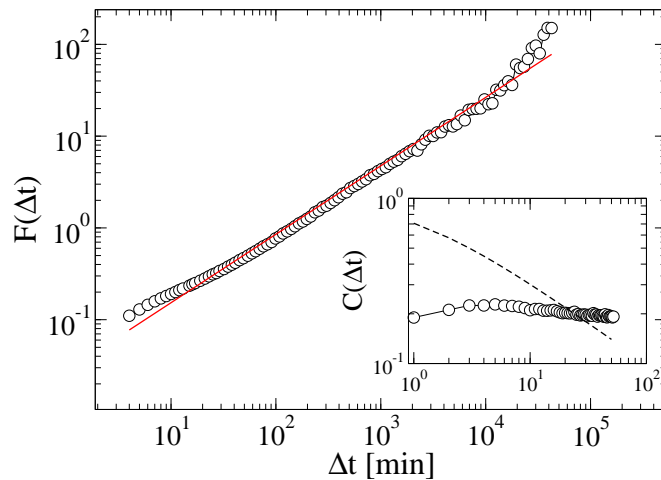
Supplementary Figure S 9: Fractions of positive, negative and neutral posts for different values of the parameters in our simulations. The ratio of emotional expressions (positive and negative) increases with the amplitude of the valence stochastic component A_v . This ratio is also slightly increased by the collective parameters b and c , as the communication influence on the valence is stronger. The inverse is true for the parameter γ_h , i.e. the ratio of neutral posts increases the larger the decay of the field. An increase in the threshold V_{\pm} leads to lower frequency of expression of the corresponding sign.



Supplementary Figure S 10: Distribution of the Hurst exponents for the simulated conversations (A) and agents (B) for the cases of $\gamma_v = 0.1$ and $\gamma_v = 0.5$. The Kolmogorov-Smirnov distance between the simulated distribution for $\gamma_v = 0.1$ and the real data is $KS=0.845$, while between the KS distance between the simulated distribution for $\gamma_v = 0.5$ and the real data is $KS=0.519$. This means that the individual and conversation persistence distributions are more similar to the real data (Fig 3 of the main text) for the case of $\gamma_v = 0.5$, implying that the relaxation speed of the valence of chatroom users is fast.



Supplementary Figure S 11: Mean value of the Hurst exponents of the emotional expression of agents and conversations for different values of the simulation parameters. Under the influence of an emotional field, the user persistence increases with A_v , meaning that a stronger stochastic component can lead to conversations more similar to the observed ones. The coupling parameters c and b increase both mean persistence. The effect of larger γ_h is the inverse, the stronger the decay of information, the weaker the persistence. Larger positive thresholds V_+ lead to lower user persistence, while the inverse is true for the negative threshold V_- . The standard error bars showing the standard error of the mean value are smaller than the symbol size and are not visible.



Supplementary Figure S 12: DFA fluctuation function calculated using the inter-event times of a simulated IRC channel. The Hurst exponent obtained is $H_{\omega'} \simeq 0.75$, suggesting the existence of log term correlations in the time series. We note the absence of pronounced dependencies in the user activity that would be manifested by a power law decaying autocorrelation function (Inset). The dotted line shows the expected decay of the autocorrelation function according to the scaling relation $\nu_{\omega'} = 2 - 2H_{\omega'}$ [1]. In this case, the origin of the correlations revealed by the Hurst exponent can only be the broad distribution of inter-event times that was given as input to the model, since there is no coupling in the activity of users.

-
- [1] Kantelhardt, J.W. Fractal and multifractal time series. *Encyclopedia of Complexity and Systems Science*. (Springer, 2009).
 - [2] Thelwall, M., Buckley, K., Paltoglou, G., Cai, D. & Kappas, A. Sentiment strength detection in short informal text. *J. Am. Soc. Inf. Sci. Technol.* **61**, 2544–2558 (2010).
 - [3] Pennebaker, J. W., Francis, M. E. & Booth, R. K. *Linguistic Inquiry and Word Count: LIWC 2001*. (Erlbaum Publishers, 2001).
 - [4] Skowron, M. & Paltoglou, G. Affect Bartender - Affective Cues and Their Application in a Conversational Agent *IEEE Symposium Series on Computational Intelligence 2011, Workshop on Affective Computational Intelligence*. (IEEE Computer Society, 2011).
 - [5] <http://www.efnet.org/>