

# Private traits and attributes are predictable from digital records of human behavior

Michal Kosinski<sup>a,1</sup>, David Stillwell<sup>a</sup>, and Thore Graepel<sup>b</sup>

<sup>a</sup>Free School Lane, The Psychometrics Centre, University of Cambridge, Cambridge CB2 3RQ United Kingdom; and <sup>b</sup>Microsoft Research, Cambridge CB1 2FB, United Kingdom

Edited by Kenneth Wachter, University of California, Berkeley, CA, and approved February 12, 2013 (received for review October 29, 2012)

**We show that easily accessible digital records of behavior, Facebook Likes, can be used to automatically and accurately predict a range of highly sensitive personal attributes including: sexual orientation, ethnicity, religious and political views, personality traits, intelligence, happiness, use of addictive substances, parental separation, age, and gender. The analysis presented is based on a dataset of over 58,000 volunteers who provided their Facebook Likes, detailed demographic profiles, and the results of several psychometric tests. The proposed model uses dimensionality reduction for preprocessing the Likes data, which are then entered into logistic/linear regression to predict individual psychodemographic profiles from Likes. The model correctly discriminates between homosexual and heterosexual men in 88% of cases, African Americans and Caucasian Americans in 95% of cases, and between Democrat and Republican in 85% of cases. For the personality trait "Openness," prediction accuracy is close to the test-retest accuracy of a standard personality test. We give examples of associations between attributes and Likes and discuss implications for online personalization and privacy.**

social networks | computational social science | machine learning | big data | data mining | psychological assessment

**A** growing proportion of human activities, such as social interactions, entertainment, shopping, and gathering information, are now mediated by digital services and devices. Such digitally mediated behaviors can easily be recorded and analyzed, fueling the emergence of computational social science (1) and new services such as personalized search engines, recommender systems (2), and targeted online marketing (3). However, the widespread availability of extensive records of individual behavior, together with the desire to learn more about customers and citizens, presents serious challenges related to privacy and data ownership (4, 5).

We distinguish between data that are actually recorded and information that can be statistically predicted from such records. People may choose not to reveal certain pieces of information about their lives, such as their sexual orientation or age, and yet this information might be predicted in a statistical sense from other aspects of their lives that they do reveal. For example, a major US retail network used customer shopping records to predict pregnancies of its female customers and send them well-timed and well-targeted offers (6). In some contexts, an unexpected flood of vouchers for prenatal vitamins and maternity clothing may be welcome, but it could also lead to a tragic outcome, e.g., by revealing (or incorrectly suggesting) a pregnancy of an unmarried woman to her family in a culture where this is unacceptable (7). As this example shows, predicting personal information to improve products, services, and targeting can also lead to dangerous invasions of privacy.

Predicting individual traits and attributes based on various cues, such as samples of written text (8), answers to a psychometric test (9), or the appearance of spaces people inhabit (10), has a long history. Human migration to digital environment renders it possible to base such predictions on digital records of human behavior. It has been shown that age, gender, occupation, education level, and even personality can be predicted from people's Web site

browsing logs (11–15). Similarly, it has been shown that personality can be predicted based on the contents of personal Web sites (16), music collections (17), properties of Facebook or Twitter profiles such as the number of friends or the density of friendship networks (18–21), or language used by their users (22). Furthermore, location within a friendship network at Facebook was shown to be predictive of sexual orientation (23).

This study demonstrates the degree to which relatively basic digital records of human behavior can be used to automatically and accurately estimate a wide range of personal attributes that people would typically assume to be private. The study is based on Facebook Likes, a mechanism used by Facebook users to express their positive association with (or "Like") online content, such as photos, friends' status updates, Facebook pages of products, sports, musicians, books, restaurants, or popular Web sites. Likes represent a very generic class of digital records, similar to Web search queries, Web browsing histories, and credit card purchases. For example, observing users' Likes related to music provides similar information to observing records of songs listened to online, songs and artists searched for using a Web search engine, or subscriptions to related Twitter channels. In contrast to these other sources of information, Facebook Likes are unusual in that they are currently publicly available by default. However, those other digital records are still available to numerous parties (e.g., governments, developers of Web browsers, search engines, or Facebook applications), and, hence, similar predictions are unlikely to be limited to the Facebook environment.

The design of the study is presented in Fig. 1. We selected traits and attributes that reveal how accurate and potentially intrusive such a predictive analysis can be, including "sexual orientation," "ethnic origin," "political views," "religion," "personality," "intelligence," "satisfaction with life" (SWL), substance use ("alcohol," "drugs," "cigarettes"), "whether an individual's parents stayed together until the individual was 21 y old," and basic demographic attributes such as "age," "gender," "relationship status," and "size and density of the friendship network." Five Factor Model (9) personality scores ( $n = 54,373$ ) were established using the International Personality Item Pool (IPIP) questionnaire with 20 items (25). Intelligence ( $n = 1,350$ ) was measured using Raven's Standard Progressive Matrices (SPM) (26), and SWL ( $n = 2,340$ ) was measured using the SWL Scale (27). Age ( $n = 52,700$ ; average,  $\mu = 25.6$ ; SD = 10), gender ( $n = 57,505$ ; 62% female), relationship status ("single"/"in relationship";  $n = 46,027$ ; 49% single), political views ("Liberal"/"Conservative";  $n = 9,752$ ;

Author contributions: M.K. and T.G. designed research; M.K. and D.S. performed research; M.K. and T.G. analyzed data; and M.K., D.S., and T.G. wrote the paper.

Conflict of interest statement: D.S. received revenue as owner of the myPersonality Facebook application.

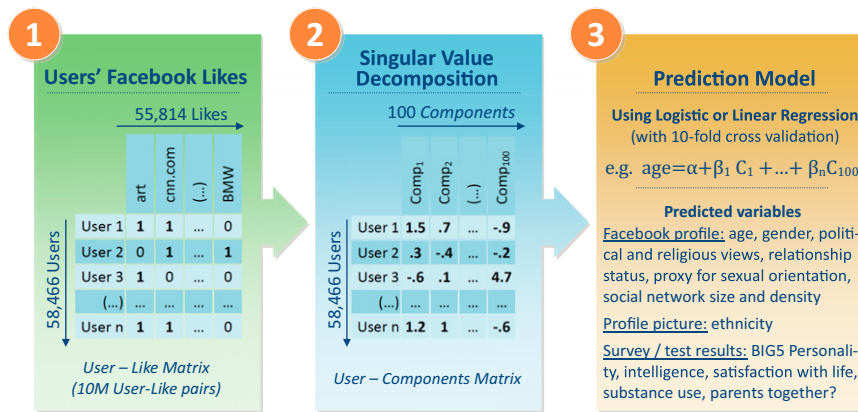
This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

Data deposition: The data reported in this paper have been deposited in the myPersonality Project database ([www.mypersonality.org/wiki](http://www.mypersonality.org/wiki)).

<sup>1</sup>To whom correspondence should be addressed. E-mail: [mk583@cam.ac.uk](mailto:mk583@cam.ac.uk).

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1218772110/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1218772110/-DCSupplemental).



**Fig. 1.** The study is based on a sample of 58,466 volunteers from the United States, obtained through the myPersonality Facebook application ([www.mypersonality.org/wiki](http://www.mypersonality.org/wiki)), which included their Facebook profile information, a list of their Likes ( $n = 170$  Likes per person on average), psychometric test scores, and survey information. Users and their Likes were represented as a sparse user-Like matrix, the entries of which were set to 1 if there existed an association between a user and a Like and 0 otherwise. The dimensionality of the user-Like matrix was reduced using singular-value decomposition (SVD) (24). Numeric variables such as age or intelligence were predicted using a linear regression model, whereas dichotomous variables such as gender or sexual orientation were predicted using logistic regression. In both cases, we applied 10-fold cross-validation and used the  $k = 100$  top SVD components. For sexual orientation, parents' relationship status, and drug consumption only  $k = 30$  top SVD components were used because of the smaller number of users for which this information was available.

65% Liberal), religion ("Muslim"/"Christian";  $n = 18,833$ ; 90% Christian), and the Facebook social network information [ $n = 17,601$ ; median size,  $\bar{X} = 204$ ; interquartile range (IQR), 206; median density,  $\bar{X} = 0.03$ ; IQR, 0.03] were obtained from users' Facebook profiles. Users' consumption of alcohol ( $n = 1,196$ ; 50% drink), drugs ( $n = 856$ ; 21% take drugs), and cigarettes ( $n = 1,211$ ; 30% smoke) and whether a user's parents stayed together until the user was 21 y old ( $n = 766$ ; 56% stayed together) were recorded using online surveys. Visual inspection of profile pictures was used to assign ethnic origin to a randomly selected subsample of users ( $n = 7,000$ ; 73% Caucasian; 14% African American; 13% others). Sexual orientation was assigned using the Facebook profile "Interested in" field; users interested only in others of the same sex were labeled as homosexual (4.3% males; 2.4% females), whereas those interested in users of the opposite gender were labeled as heterosexual.

## Results

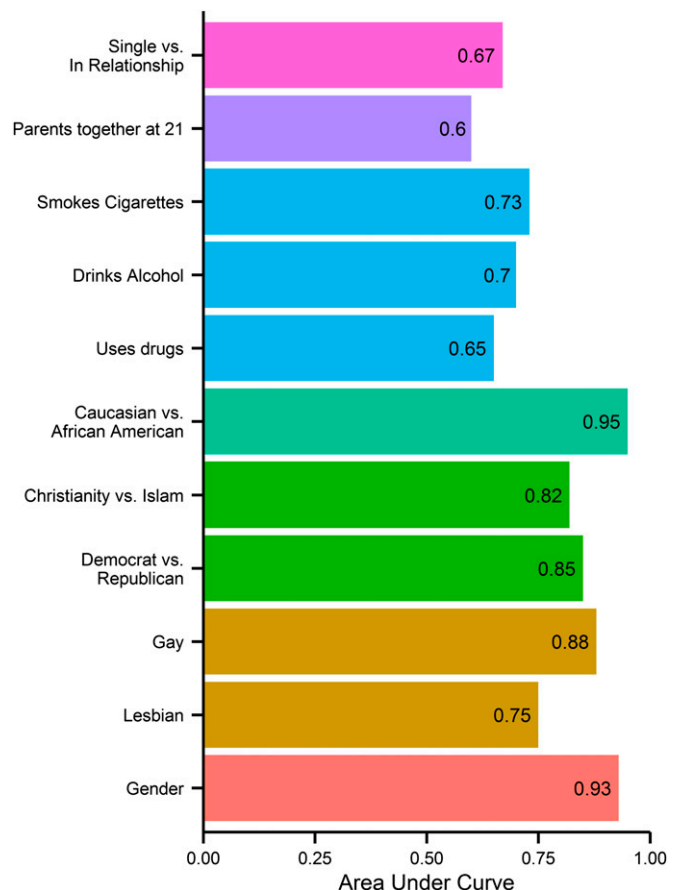
**Prediction of Dichotomous Variables.** Fig. 2 shows the prediction accuracy of dichotomous variables expressed in terms of the area under the receiver-operating characteristic curve (AUC), which is equivalent to the probability of correctly classifying two randomly selected users one from each class (e.g., male and female). The highest accuracy was achieved for ethnic origin and gender. African Americans and Caucasian Americans were correctly classified in 95% of cases, and males and females were correctly classified in 93% of cases, suggesting that patterns of online behavior as expressed by Likes significantly differ between those groups allowing for nearly perfect classification.

Christians and Muslims were correctly classified in 82% of cases, and similar results were achieved for Democrats and Republicans (85%). Sexual orientation was easier to distinguish among males (88%) than females (75%), which may suggest a wider behavioral divide (as observed from online behavior) between hetero- and homosexual males.

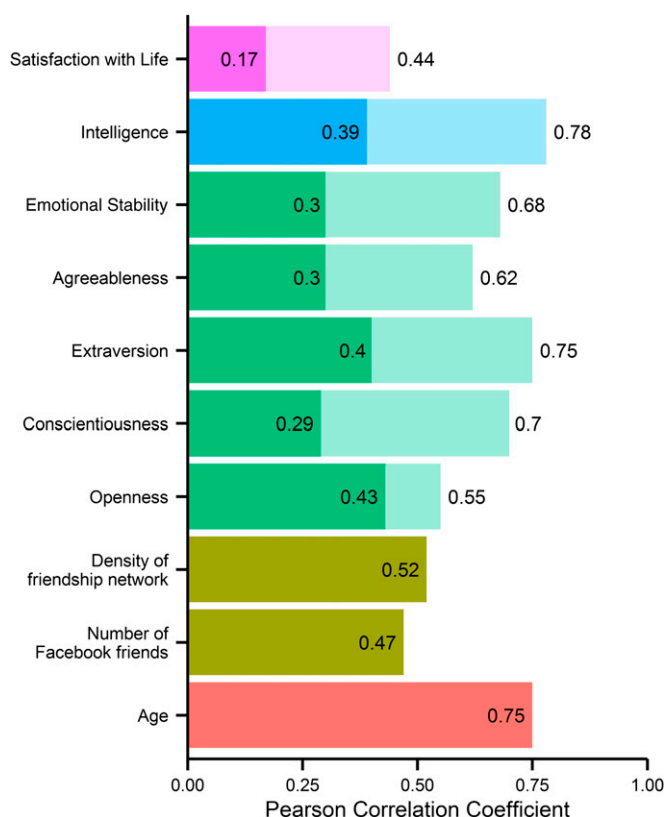
Good prediction accuracy was achieved for relationship status and substance use (between 65% and 73%). The relatively lower accuracy for relationship status may be explained by its temporal variability compared with other dichotomous variables (e.g., gender or sexual orientation).

The model's accuracy was lowest (60%) when inferring whether users' parents stayed together or separated before users were 21 y old. Although it is known that parental divorce does have long-

term effects on young adults' well-being (28), it is remarkable that this is detectable through their Facebook Likes. Individuals with parents who separated have a higher probability of liking statements preoccupied with relationships, such as "If I'm with you then I'm with you I don't want anybody else" (Table S1).



**Fig. 2.** Prediction accuracy of classification for dichotomous/dichotomized attributes expressed by the AUC.



**Fig. 3.** Prediction accuracy of regression for numeric attributes and traits expressed by the Pearson correlation coefficient between predicted and actual attribute values; all correlations are significant at the  $P < 0.001$  level. The transparent bars indicate the questionnaire's baseline accuracy, expressed in terms of test-retest reliability.

**Prediction of Numeric Variables.** Fig. 3 presents the accuracy of predicting numeric variables as expressed by the Pearson product-moment correlation coefficient between the actual and predicted values. The highest correlation was obtained for age ( $r = 0.75$ ), followed by density ( $r = 0.52$ ) and size ( $r = 0.47$ ) of the Facebook friendship network. Closely following were the personality traits of "Openness" ( $r = 0.43$ ), "Extraversion" ( $r = 0.40$ ), and "Intelligence" ( $r = 0.39$ ). The remaining personality traits and SWL were predicted with somewhat lower accuracy ( $r = 0.17$  to  $0.30$ ).

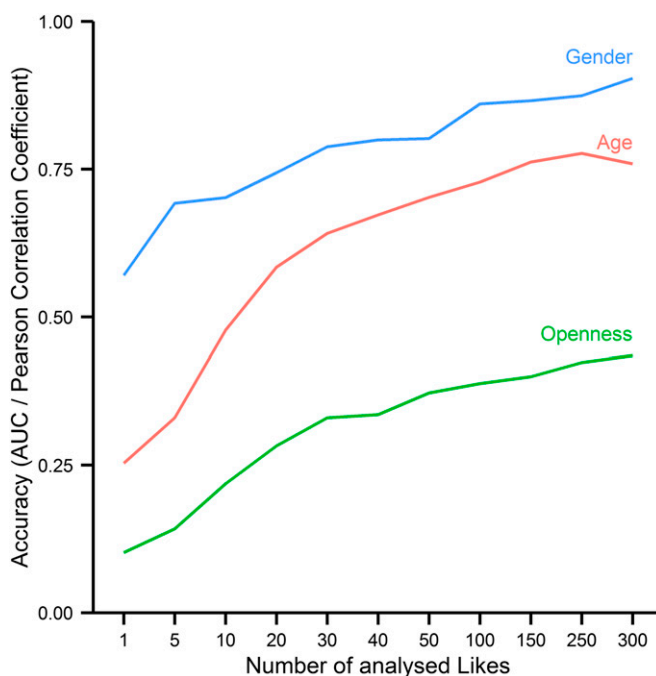
Psychological traits are examples of latent traits (i.e., traits that cannot be measured directly). As a consequence, their values can only be measured approximately, for example, by evaluating responses to questionnaires. The transparent bars presented in Fig. 3 indicate the accuracy of the questionnaires used as expressed by their test-retest reliabilities (Pearson product-moment correlation between the questionnaire scores obtained by the same respondent at two points in time). The correlation between the predicted and actual Openness score ( $r = 0.43$ ) was very close to the test-retest reliability for Openness ( $r = 0.50$ ). This indicates that for the Openness trait, observation of the user's Likes is roughly as informative as using their personality test score itself. For the remaining traits, prediction accuracies correspond to roughly half the questionnaire's test-retest reliabilities.

The relatively lower prediction accuracy for SWL ( $r = 0.17$ ) may be attributable to the difficulty of separating long-term happiness (29) from mood swings, which vary over time. Thus, although the SWL score includes variability attributable to mood, users' Likes accrue over a longer period and, so, may be suitable only for predicting long-term happiness.

**Amount of Data Available and Prediction Accuracy.** The results presented so far rely on individuals for which between one and 700 Likes were available. The median number of Likes was 68 per individual (IQR, 152). Therefore, what is the expected accuracy given a random individual and how does prediction accuracy change with the number of observed Likes? Using a subsample ( $n = 500$ ) of users for whom at least 300 Likes were available, we ran predictive models based on randomly selected subsets of  $n = 1, 2, \dots, 300$  Likes. The results presented in Fig. 4 show that even knowing a single random Like for a given user can result in nonnegligible prediction accuracy. Knowing further Likes increases the accuracy but with diminishing returns from each additional piece of information.

**Predictive Power of Likes.** Individual traits and attributes can be predicted to a high degree of accuracy based on records of users' Likes. Table S1 presents a sample of highly predictive Likes related to each of the attributes. For example, the best predictors of high intelligence include "Thunderstorms," "The Colbert Report," "Science," and "Curly Fries," whereas low intelligence was indicated by "Sephora," "I Love Being A Mom," "Harley Davidson," and "Lady Antebellum." Good predictors of male homosexuality included "No H8 Campaign," "Mac Cosmetics," and "Wicked The Musical," whereas strong predictors of male heterosexuality included "Wu-Tang Clan," "Shaq," and "Being Confused After Waking Up From Naps." Although some of the Likes clearly relate to their predicted attribute, as in the case of No H8 Campaign and homosexuality, other pairs are more elusive; there is no obvious connection between Curly Fries and high intelligence.

Moreover, note that few users were associated with Likes explicitly revealing their attributes. For example, less than 5% of users labeled as gay were connected with explicitly gay groups, such as No H8 Campaign, "Being Gay," "Gay Marriage," "I love Being



**Fig. 4.** Accuracy of selected predictions as a function of the number of available Likes. Accuracy is expressed as AUC (gender) and Pearson's correlation coefficient (age and Openness). About 50% of users in this sample had at least 100 Likes and about 20% had at least 250 Likes. Note, that for gender (dichotomous variable) the random guessing baseline corresponds to an AUC = 0.50.



Gay,” “We Didn’t Choose To Be Gay We Were Chosen.” Consequently, predictions rely on less informative but more popular Likes, such as “Britney Spears” or “Desperate Housewives” (both moderately indicative of being gay).

This is further illustrated in Fig. S1, which shows the average levels of personality traits and age for several popular Likes. Each Like attracts users with a different average personality and demographic profile and, thus, can be used to predict those attributes. For example, users who liked the “Hello Kitty” brand tended to be high on Openness and low on “Conscientiousness,” “Agreeableness,” and “Emotional Stability.” They were also more likely to have Democratic political views and to be of African-American origin, predominantly Christian, and slightly below average age. The same Likes were used to create Fig. S2, presenting their relative popularity in four groups: Democrats, Christians, Homosexuals, and African-American individuals. For example, although liking “Barack Obama” is clearly related to being a Democrat, it is also relatively popular among Christians, African Americans, and Homosexual individuals.

## Conclusions

We show that a wide variety of people’s personal attributes, ranging from sexual orientation to intelligence, can be automatically and accurately inferred using their Facebook Likes. Similarity between Facebook Likes and other widespread kinds of digital records, such as browsing histories, search queries, or purchase histories suggests that the potential to reveal users’ attributes is unlikely to be limited to Likes. Moreover, the wide variety of attributes predicted in this study indicates that, given appropriate training data, it may be possible to reveal other attributes as well.

Predicting users’ individual attributes and preferences can be used to improve numerous products and services. For instance, digital systems and devices (such as online stores or cars) could be designed to adjust their behavior to best fit each user’s inferred profile (30). Also, the relevance of marketing and product recommendations could be improved by adding psychological dimensions to current user models. For example, online insurance advertisements might emphasize security when facing emotionally unstable (neurotic) users but stress potential threats when dealing

with emotionally stable ones. Moreover, digital records of behavior may provide a convenient and reliable way to measure psychological traits. Automated assessment based on large samples of behavior may not only be more accurate and less prone to cheating and misrepresentation but may also permit assessment across time to detect trends. Moreover, inference based on observations of digitally recorded behavior may open new doors for research in human psychology.

On the other hand, the predictability of individual attributes from digital records of behavior may have considerable negative implications, because it can easily be applied to large numbers of people without obtaining their individual consent and without them noticing. Commercial companies, governmental institutions, or even one’s Facebook friends could use software to infer attributes such as intelligence, sexual orientation, or political views that an individual may not have intended to share. One can imagine situations in which such predictions, even if incorrect, could pose a threat to an individual’s well-being, freedom, or even life. Importantly, given the ever-increasing amount of digital traces people leave behind, it becomes difficult for individuals to control which of their attributes are being revealed. For example, merely avoiding explicitly homosexual content may be insufficient to prevent others from discovering one’s sexual orientation.

There is a risk that the growing awareness of digital exposure may negatively affect people’s experience of digital technologies, decrease their trust in online services, or even completely deter them from using digital technology. It is our hope, however, that the trust and goodwill among parties interacting in the digital environment can be maintained by providing users with transparency and control over their information, leading to an individually controlled balance between the promises and perils of the Digital Age.

**ACKNOWLEDGMENTS.** We thank Yoram Bachrach, Alan Blackwell, George Danezis, Stephen Emmott, David Good, Peter Key, Emre Kiciman, Pushmeet Kohli, Drew Purves, Jason Rentfrow, John Rust, and Duncan Watts for discussions about the topic of this study, as well as for comments on the manuscript. Demonstration of personality prediction based on individuals’ Likes is available at <http://www.youarewhatyoulike.com>. M.K. received funding from Boeing Corporation and from Microsoft Research.

- Lazer D, et al. (2009) Computational social science. *Science* 323(5915):721–723.
- Koren Y, Bell R, Volinsky C (2009) Matrix factorization techniques for recommender systems. *Computer* 42(8):30–37.
- Chen Y, Pavlov D, Canny JF (2009) Large-scale behavioral targeting. *International Conference on Knowledge Discovery and Data Mining*, pp 209–218.
- Butler D (2007) Data sharing threatens privacy. *Nature* 449(7163):644–645.
- Narayanan A, Shmatikov V (2008) Robust de-anonymization of large sparse datasets. *IEEE Symposium on Security and Privacy*, pp 111–125.
- Duhigg C (2012) *The Power of Habit: Why We Do What We Do in Life and Business* (Random House, New York).
- Ince HO, Yarali A, Özsel D (2009) Customary killings in Turkey and Turkish modernization. *Middle East Stud* 45(4):537–551.
- Fast LA, Funder DC (2008) Personality as manifest in word use: Correlations with self-report, acquaintance report, and behavior. *J Pers Soc Psychol* 94(2):334–346.
- Costa PT, McCrae RR (1992) *Revised NEO Personality Inventory (NEO-PI-R) and NEO Five-Factor Inventory (NEO-FFI) Manual* (Psychological Assessment Resources, Odessa, FL).
- Gosling SD, Ko SJ, Mannarelli T, Morris ME (2002) A room with a cue: Personality judgments based on offices and bedrooms. *J Pers Soc Psychol* 82(3):379–398.
- Hu J, Zeng H-J, Li H, Niu C, Chen Z (2007) Demographic prediction based on user’s browsing behavior. *International World Wide Web Conference*, pp 151–160.
- Murray D, Durrell K (1999) Inferring demographic attributes of anonymous Internet users. *Revised Papers from the International Workshop on Web Usage Analysis and User Profiling*, eds Masand BM, Spiliopoulou M (Springer, London), pp 7–20.
- De Bock K, Van Den Poel D (2010) Predicting website audience demographics for Web advertising targeting using multi-website clickstream data. *Fundamenta Informaticae* 98(1):49–70.
- Goel S, Hofman JM, Siroer MI (2012) Who does what on the Web: Studying Web browsing behavior at scale. *International Conference on Weblogs and Social Media*, pp 130–137.
- Kosinski M, Kohli P, Stillwell DJ, Bachrach Y, Graepel T (2012) Personality and website choice. *ACM Web Science Conference*, pp 251–254.
- Marcus B, Machilek F, Schütz A (2006) Personality in cyberspace: Personal Web sites as media for personality expressions and impressions. *J Pers Soc Psychol* 90(6):1014–1031.
- Rentfrow PJ, Gosling SD (2003) The do re mi’s of everyday life: The structure and personality correlates of music preferences. *J Pers Soc Psychol* 84(6):1236–1256.
- Quercia D, Lambiotte R, Kosinski M, Stillwell D, Crowcroft J (2012) The Personality of popular Facebook users. *ACM Conference on Computer Supported Cooperative Work*. Proceedings of the ACM 2012 Conference on Computer Supported Cooperative Work, pp 955–964.
- Bachrach Y, Kohli P, Graepel T, Stillwell DJ, Kosinski M (2012) Personality and patterns of Facebook usage. *ACM Web Science Conference*. Proceedings of the ACM Web Science Conference, pp 36–44.
- Quercia D, Kosinski M, Stillwell DJ, Crowcroft J (2011) Our Twitter profiles, our selves: Predicting personality with Twitter. *IEEE International Conference on Social Computing*. Proceedings of the 2011 IEEE International Conference on Privacy, Security, Risk, and Trust, and IEEE International Conference on Social Computing, pp 180–185.
- Golbeck J, Robles C, Edmondson M, Turner K (2011) Predicting personality from Twitter. *IEEE International Conference on Social Computing*, pp 149–156.
- Golbeck J, Robles C, Turner K (2011) Predicting personality with social media. *Conference on Human Factors in Computing Systems*, pp 253–262.
- Jernigan C, Mistree BF (2009) Gaydar: Facebook friendships expose sexual orientation. *First Monday* 14(10).
- Golub GH, Kahan W (1965) Calculating the singular values and pseudo-inverse of a matrix. *J Soc Ind Appl Math* 2(2):205–224.
- Goldberg LR, et al. (2006) The international personality item pool and the future of public-domain personality measures. *J Res Pers* 40(1):84–96.
- Raven JC (2000) The Raven’s progressive matrices: Change and stability over culture and time. *Cognit Psychol* 41(1):1–48.
- Diener E, Emmons RA, Larsen RJ, Griffin S (1985) The satisfaction with life scale. *J Pers Assess* 49(1):71–75.
- Musick K, Meier A (2010) Are both parents always better than one? Parental conflict and young adult well-being. *Soc Sci Res* 39(5):814–830.
- Schimmack U, Diener E, Oishi S (2002) Life-satisfaction is a momentary judgment and a stable personality characteristic: The use of chronically accessible and stable sources. *J Pers* 70(3):345–384.
- Nass C, Lee KM (2000) Does computer-generated speech manifest personality? An experimental test of similarity-attraction. *J Exp Psychol* 7(3):171–181.

# Supporting Information

Kosinski et al. 10.1073/pnas.1218772110

## SI Text

**SI Results.** Table S1 presents Likes characterized by the most extreme average levels for each of the numeric variables (e.g. personality traits) or most extreme frequencies of classes (e.g. being a Democrat). Fig. S1 shows the average levels of personality traits and age of the users associated with selected Likes presented on the percentile scale. Fig. S2 presents relative popularity of selected Likes within groups of Democrat, Homosexual, Christian, and African-American users.

**Sample.** We used data from 58,466 US Facebook users, including their psychodemographic profile and their list of Likes. The data were obtained from the myPersonality application ([www.mypersonality.org](http://www.mypersonality.org)). Users opted in to provide their data for this study and gave their consent to have their scores and profile information recorded for analysis.

An important limitation of our sample is that some of the predicted variables are from Facebook profile information. Individuals who declare their political and religious views, relationship status, and sexual orientation on their profile may be different from nondeclaring members of those groups; they may associate with distinct Likes, which may lead to an overestimate of prediction accuracies for these groups. Nevertheless, the model was still able to predict privately reported information, such as personality or intelligence quotient questionnaire results, and survey results on addictive substance use.

**Political and Religious Views, Sexual Orientation, Relationship Status.** Political and religious views were recorded from the respective fields of users' Facebook profiles. Both fields allow users to input text freely (but suggest popular choices). Political views "Democrat," "Democratic," or "Democratic Party" were recoded to "Democrat." "Republican," "GOP," and "Republican Party" were recoded to "Republican"; other entries were ignored. Religious views "Christian," "Catholic," and "Jesus Christ" were recoded to "Christian." "Moslem," "Muslim," "Islam," and "Sunni" were recoded to "Muslim." Sexual orientation was taken from the "Interested in" section of users' Facebook profiles; users who listed being interested in only the opposite gender were labeled as being heterosexual, whereas users who listed only the same gender were labeled as being homosexual. Relationship status was recorded from the "Relationship Status" profile field, where the options were "Single," "It's complicated," "In an open relationship," "In a relationship," "Engaged," and "Married." The latter three options were recoded to "In a relationship."

**Ethnicity.** Labels for ethnicity were assigned to users by visual inspection of their profile pictures. This procedure has the advantage that the data are not explicitly self-reported and, hence, does not suffer from disclosure bias. However, some users do not include any picture with their profile or use a picture that does not show themselves. To confirm the reliability of the manual classification procedure, a subsample of the data were compared with self-reported ethnic background from a survey, and there was  $r = 0.98$  agreement between the two sets of labels.

**Substance Use and User's Parents Together at Age Twenty-One Years.** Both substance use and whether a user's parents stayed together or split up before the user was 21 y old were measured using self-report survey measures on the myPersonality application. These questions were explicitly labeled as optional. Individuals were asked if they smoked daily, less than daily, or were nonsmokers; less than daily

and daily were recoded as "smokers." They were also asked if they drank alcohol by offering the choices "weekly or more often," "less than once a week," or "never"; the first two options were recoded as "drinkers." For drug use, the options were the same as for drinking; the first two options were recoded as "drug users."

**Personality.** Five-Factor Model (FFM) (1) personality scores ( $n = 54,373$ ) were established using the International Personality Item Pool (IPIP) questionnaire with 20 items (2). This test is widely used in both traditional and online studies and is known to be successful at explaining variability across individuals. FFM encompasses the following traits.

**Openness to Experience.** Openness to experience ("Openness") is related to imagination, creativity, curiosity, tolerance, political liberalism, and appreciation for culture. People scoring high on Openness like change, appreciate new and unusual ideas, and have a good sense of aesthetics.

**Conscientiousness.** "Conscientiousness" measures preference for an organized approach to life in contrast to a spontaneous one. Conscientious people are more likely to be well organized, reliable, and consistent. They enjoy planning, seek achievements, and pursue long-term goals. Nonconscientious individuals are generally more easy-going, spontaneous, and creative. They tend to be more tolerant and less bound by rules and plans.

**Extraversion.** "Extraversion" measures a tendency to seek stimulation in the external world, the company of others, and to express positive emotions. Extraverts tend to be more outgoing, friendly, and socially active. They are usually energetic and talkative; they do not mind being at the center of attention and make new friends more easily. Introverts are more likely to be solitary or reserved and seek environments characterized by lower levels of external stimulation.

**Agreeableness.** "Agreeableness" relates to a focus on maintaining positive social relations, being friendly, compassionate, and cooperative. Agreeable people tend to trust others and adapt to their needs. Disagreeable people are more focused on themselves, less likely to compromise, and may be less gullible. They also tend to be less bound by social expectations and conventions and more assertive.

**Emotional Stability.** "Emotional Stability" (reversely referred to as neuroticism) measures the tendency to experience mood swings and emotions such as guilt, anger, anxiety, and depression. Emotionally unstable (neurotic) people are more likely to experience stress and nervousness, whereas emotionally stable people (low neuroticism) tend to be calmer and self-confident.

**Intelligence.** Intelligence ( $n = 1,350$ ) was measured using Raven's Standard Progressive Matrices (SPM) (3), a multiple choice non-verbal intelligence test drawing on Spearman's theory of general ability. SPM is a proven standard intelligence test used in both research and clinical settings, as well as in high-stake contexts such as in military personnel selection and court cases (4). The SPM test was shortened for the purpose of this study and contained 20 items only. Note that SPM was used only to compare between users of this study, and no comparisons with the general population were made.

**Satisfaction with Life.** "Satisfaction with Life" (SWL) ( $n = 2,340$ ) was measured using the SWL Scale (5), a widely used, five-item instrument designed to measure global cognitive judgments of satisfaction with one's own life.

**Facebook Likes and User-Like Matrix.** Facebook Likes allow Facebook users to connect with virtually any object that has an online

presence. Likes are one of the most typical and pervasive forms of digitally recorded behavior. People can Like quotes, Web sites, press articles, products, activities, places they visit (or would like to visit), and content such as pictures, movies, books, and music. Likes span a diverse set of entities, from “Bible” and “Philosophy” through “Bonfires,” “BMW,” and “cnn.com” to statements such as “I hate myself.” People’s Likes are shared with their friends and can be used as a way of expressing support, bookmarking, or enhancing online identity, by indicating individual preferences.

We recorded more than 9 million unique objects liked by users, a great majority of which were associated with one or very few users only. For the purpose of building a predictive model, Likes associated with fewer than 20 users, as well as users with fewer than two Likes, were removed from the sample. The remaining 58,466 users and 55,814 unique liked objects were arranged in a sparse matrix (user–Like matrix), the columns of which represent Likes and the rows of which represent users. The entries were set to 1 if there existed an association between a user and a Like and 0 otherwise. The matrix contained roughly 10 million associations between users and Likes. To facilitate the predictive analysis, the dimensionality of the user–Like matrix was reduced using singular-value decomposition (SVD) (6) such that each user is represented by a vector of  $k$  component scores. SVD provides a low-rank approximation to the original matrix, and the approximation quality increases with the number  $k$ .

**Choosing the Number of the SVD Components.** To choose the optimum number of SVD components to be used in this study, we examined the cross-validated prediction accuracy as a function of the number of components. Fig. S3, based on Openness and Extraversion, shows that prediction accuracy increases steeply in the beginning but flattens out relatively early (note that the horizontal axis is not linear). Interestingly, including some of the components abruptly increases prediction accuracy for certain traits. For example, including component 3 in the model increases the accuracy of Openness estimates from  $r = 0.1$  to  $r = 0.4$ . Similarly, component 5 sharply increases the accuracy achieved in predicting Extraversion. This suggests that particular components are specifically related to a given attribute. We used the first  $k = 100$  SVD components, which explained 28% of the variance in the user–Like matrix (Fig. S4). For sexual orientation, parents’ relationship status, and drug consumption, only  $k = 30$  top SVD components were used because of the smaller number of users for which this information was available.

**Predictions.** SVD components were used to build models that predict users’ individual traits and preferences. Predictions related

to numeric variables, such as age or intelligence, were calculated using a linear regression model based on the users’  $k = 100$  SVD components as covariates. Dichotomous variables such as gender, relationship status, and political views were modeled using logistic regression based on the same SVD components. In both cases 10-fold cross-validation was used to assess the out-of-sample prediction accuracy: the sample was randomly split into 10 equally sized subsets of users, and predictions for each subset were calculated based on parameters determined on the remaining users. Prediction accuracy was measured in two ways. For the numeric variables, such as age in years, we report the Pearson product–moment correlation coefficient between the actual and predicted values across users. For the dichotomous variables such as gender, we report the area under the receiver-operating characteristic (ROC) curve (AUC) coefficient, which can be interpreted as the probability of correctly classifying two randomly selected objects: one of each class (e.g., male and female).

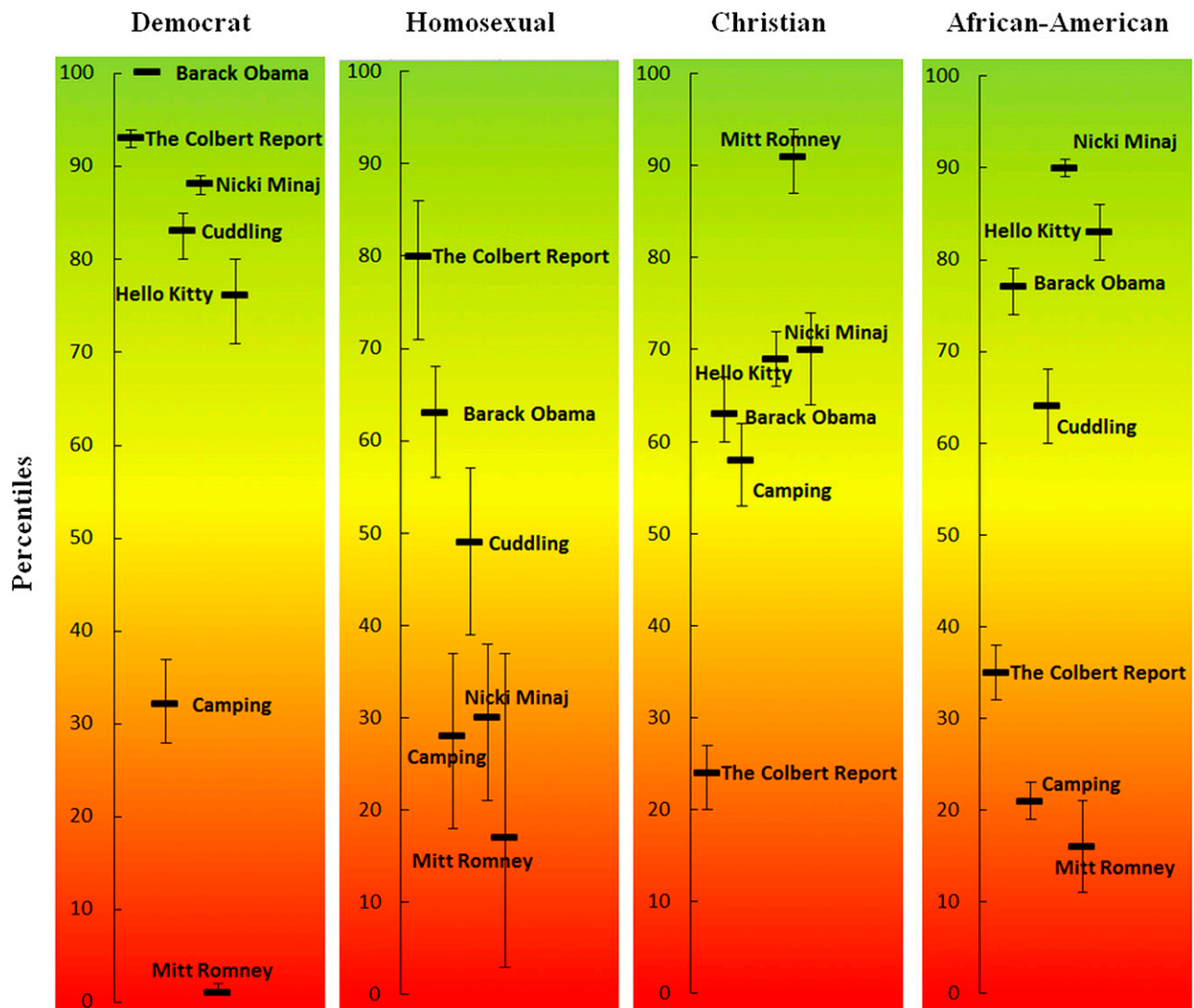
**AUC.** AUC relates to the ROC curve, which is a plot of true-positive rate (or sensitivity) versus false-positive rate (or 1 specificity) for detection or classification tasks. Positive cases are those classified by the model to belong to a target class (e.g., “male” or “Democrat”). Thus, true positive cases are the cases that were correctly classified by the model as belonging to a target class, whereas false-positive cases were classified incorrectly as belonging to a target class. The true-positive rate is the ratio of the number of true positives to the number of all cases in the target class, whereas the false-positive rate is the ratio of the number of false positives to the number of all cases in the background class. The logistic regression model used in this study to predict dichotomous outcomes assigns a probability of belonging to a target class to each of the users. To avoid having to select a single threshold for assigning users to a given target category, an ROC curve can be used to analyze the entire spectrum of possible thresholds. An example of an ROC curve is presented in Fig. S5. In general, ROC curves for random (or null) models should be close to diagonal, because the probability of seeing a true positive is not greater than the probability of seeing a false positive. The more an ROC curve bulges to the upper left, however, the higher the accuracy of the model, because higher true-positive rates are achieved for a given number of false positives. The AUC is simply the area below the ROC curve, and it is equal to the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one.

- Costa PT, McCrae RR (1992) *Revised NEO Personality Inventory (NEO-PI-R) and NEO Five-Factor Inventory (NEO-FFI) Manual* (Psychological Assessment Resources, Odessa, FL).
- Goldberg LR, et al. (2006) The international personality item pool and the future of public-domain personality measures. *J Res Pers* 40(1):84–96.
- Raven JC (2000) The Raven’s progressive matrices: change and stability over culture and time. *Cognit Psychol* 41(1):1–48.

- Lubinski D (2004) Introduction to the special section on cognitive abilities: 100 years after Spearman’s (1904) “‘General intelligence,’ objectively determined and measured”. *J Pers Soc Psychol* 86(1):96–111.
- Diener E, Emmons RA, Larsen RJ, Griffin S (1985) The satisfaction with life scale. *J Pers Assess* 49(1):71–75.
- Golub GH, Kahan W (1965) Calculating the singular values and pseudo-inverse of a matrix. *J Soc Ind Appl Math* 2(2):205–224.







**Fig. S2.** Relative popularity of selected Likes within groups of Democrat, Homosexual, Christian, and African American users. Because Likes differed greatly in popularity (e.g., “Barack Obama” was nearly four times more popular than “Mitt Romney”), we calculated relative popularity by dividing the frequencies of associations with a given Like within the studied groups by the respective frequency in the entire sample. Relative popularity was transformed into a percentile scale. Error bars signify 95% confidence intervals of the population proportion. For example, The Colbert Report is relatively popular within Democrats and Homosexual groups (93th and 80th percentile respectively) but rather unpopular among Christians and African Americans (24th and 35th percentile, respectively).





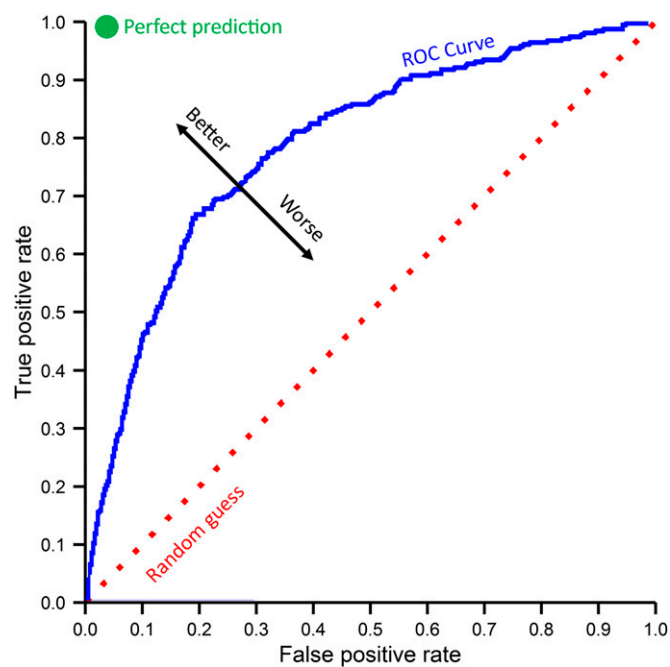


Fig. S5. Example of an ROC curve, detecting users associated with the Facebook Page associated with the [FailBlog.org](http://FailBlog.org) website. The AUC for this plot is 0.79.

## Other Supporting Information Files

[Table S1 \(PDF\)](#)

Trait		Selected most predictive Likes							
IQ	High	<del>The Godfather</del> Mozart Thunderstorms The Colbert Report Morgan Freemans Voice The Daily Show Lord Of The Rings To Kill A Mockingbird Science Curly Fries	Jason Aldean Tyler Perry Sephora Chiq Bret Michaels Clark Griswold <del>Bebe</del> I Love Being A Mom Harley Davidson Lady Antebellum	Low					
		Satisfaction With Life	Satisfied		Dissatisfied				
Openness		Liberal & Artistic	Oscar Wilde Charles Bukowski Sylvia Plath <del>Leonardo Da Vinci</del> Bauhaus Dmt The Spirit Molecule American Gods John Waters Plato Leonard Cohen		NASCAR Austin Collie Monster-In-Law I don't read Justin Moore <del>ESPN2</del> Farmlandia The Bachelor Oklahoma State University Teen Mom 2	Conservative			
			Conscientiousness		Well Organized		Spontaneous		
			Extraversion		Outgoing & Active		Beerpong Michael Jordan <del>Dancing</del> Socializing Chris Tucker I Feel Better Tan Modeling Cheerleading Theatre Flip Cup	RPGs Fanfiction.Net Programming <del>Anime</del> Manga Video Games Role Playing Games Minecraft Voltaire Terry Pratchet	Shy & Reserved

Agreeableness	Cooperative	Compassion International Logan Utah Jon Foreman Redeeming Love Pornography Harms The Book Of Mormon Circles Of Prayer Go To Church Christianity Marianne Williamson	I Hate Everyone I Hate You I Hate Police Friedrich Nietzsche Timmy South Park Atheism / Satanism Prada Sun Tzu Julius Caesar Knives	Competitive
Emotional Stability	Neurotic	Sometimes I Hate Myself Emo Girl Interrupted So So Happy The Addams Family <del>Vocaloid</del> Sixbillionsecrets.com Vampires Everywhere Kurt Donald Cobain Dot Dot Curve	Business Administration Getting Money Parkour Track & Field Skydiving Mountain Biking <del>Soccer</del> Climbing Physics / Engineering 48 Laws Of Power	Calm & Relaxed
Gender	Female	Tv Fanatic Chiq Gillette Venus Shoedazzle Bebe Proud To Be A Mom Covergirl Wet Seal <del>Aerie By American Eagle</del> Mall World	<del>Modern Warfare 2</del> ESPN Sportscenter Band Of Brothers Starcraft Deadliest Warrior Dos Equis Red Vs Blue X Games Bruce Lee	Male
Age	Old	Cup Of Joe For A Joe Coffee Party Movement Dr Mehmet Oz Fixit And Forgetit The Closer Joyce Meyer Ministries Proud To Be A Mom Freedomworks <del>Small Business Saturday</del> Fly The American Flag	Walt Disney Records Body By Milk <del>Harperteen</del> J Bigga Because I Am A Girl I Hate My Id Photo 293 Things To Do In Class When You Are Bored Dude Wait What JCP Teen	Young
Friends	Many	Mojo-Jojo Biology <del>Dollar General</del> Hillary 106 & Park Jennifer Lopez Paid In Full Yo Gotti The Dollar You Are Holding Could've Been In A Stripper's Butt Crack	The Dark Knight In'n'out Burger Hard Rock Honey, Where Is My Supersuit Hating ICP <del>Minecraft</del> Iron Maiden Walking With Your Friend & Randomly Pushing Them Into Someone/Something	Few



Religion	<i>Christian</i>	The Bible Jesus Daily I'm Proud To Be Christian God Jesus Christ Church The Holy Bible I Love Jesus Christian Music Gospel Music	I'm A Muslim & I'm Proud Hadith Of The Day I Love Islam I Need Allah In My Life Prophet Muhammad Saw The Greatest Man In History Remove Group Fuck Islam From Facebook Nancy Ajram Moozlum The Movie Desihits.Com	<i>Muslim</i>
	<i>Republican</i>	George W Bush John McCain Conservative Rush Limbaugh Sean Hannity Bill Oreilly Positively Republican Sarah Palin Ronald Reagan Glenn Beck	Joe Biden Speaker Nancy Pelosi Health Care Reform The White House Democrats Barbara Boxer Anthony Weiner Being Liberal Left Action Barack Obama2012 Ted Kennedy	<i>Democrat</i>
Sexual Orientation	<i>Homosexual Males</i>	No H8 Campaign Kathy Griffin Kurt Hummel Glee Human Rights Campaign Mac Cosmetics Adam Lambert Ellen DeGeneres Juicy Couture Sue Sylvester Glee Wicked The Musical	X Games <del>Nike Basketball</del> Bungie WWE Sportsnation Wu-Tang Clan Foot Locker Shaq Bruce Lee Being Confused After Waking Up From Naps	<i>Heterosexual Males</i>
	<i>Homosexual Females</i>	Girls Who Like Boys Who Like Boys Rupauls Drag Race No H8 Campaign Gay Marriage Human Rights Campaign The L Word Sometimes I Just Lay In Bed And Think About Life Not Being Pregnant Gay Marriage Tegan And Sara	Lipton Brisk Yahoo Adidas Originals Foot Locker WWE Inbox 1 Makes Me Nervous Thinking Of Something And Laughing Alone I Just Realized Immature Spells I'm Mature Did You Get A Haircut No It Grew Shorter Nike Women	<i>Heterosexual Females</i>

Race	African-American	<p>I Support My President Fantasia Jill Scott Next Friday Erykah Badu Maxwell Taraji P Henson Madea Tyga Love And Basketball</p>	<p>Just Because You Can Reproduce Doesn't Mean You Should I Come From A Town Where A Traffic Jam Is 4 Cars Behind A Tractor Harley Davidson Halloween Bret Michaels David Bowie Official ASPCA Fly The American Flag Road Trips Bonfires</p>	White American
Relationship	In a Relationship	<p>I Love My Husband Kids Circle Of Moms Parents Magazine Tacori Weight Watchers Scrapbooking Huggies Box Tops For Education Babies R Us</p>	<p>J.Cole Hunger Games Ign.com Kassem G Sonny With A Chance Usain Bolt 2ne1 Mangastream Sportsnation Maria Sharapova</p>	Single
Alcohol Use	Yes	<p>Watching Karma Bite The Person You Hate Right In The Ass Dear Liver Thanks You're A Champ Trying To Figure Out If Its A Cop Car Belvedere Vodka Meeting Someone Who Is Also Drunk And Immediately Becoming Best Friends Jim Beam I Love It When In The Middle Of Our Kiss I Can Feel You Smiling Tattoo Lovers Getting A Text That Says I Miss You Drinking Around A Bonfire</p>	<p>Bungie I Hate Going Back To School After The Holidays When I'm Home Alone And I Hear A Noise I Freeze And Listen For Ages Not Finishing A Sentence Because Your Laughing Too Hard About The Ending Why Is Monday So Far Away From Friday And Friday So Bloody Close To Monday I Hate When I Originally Pick The Right Answer Then Change It That's Going In My Status When I Get Home I Don't Care There Is 30 Seconds Left In This Class I'm Packing Up Pretending To Think When The Teacher Is Looking At You I Like Watching Raindrops Race Across My Window And Silently Cheer For Them</p>	No

<div>Parents separated at 21</div> <div>Yes</div>	<p>When Ur Single, All U See Is Happy          Couples N Wen Ur In A Relationship All U See Is Happy Singles          Never Apologize For What You Feel It's Like Saying Sorry For Being Real          I'm The Type Of Girl Who Can Be So Hurt But Still Look At You &amp; Smile The Type Of Girl Who Is Willing To Brighten Your Day Even If I Can't Brighten My Own &lt;3 &lt;3          We Don't Talk Anymore And You Know What The Saddest Part Is We Used To Talk Everyday          Come Here Nope [Grab Yu Closer] Gimme A Kiss Nope I'm Mad At Yu [Start Kissing]          Tell Her She's Pretty Hold Her Hand Kiss Her When She's Angry Play With Her Hair Let Her Fall Asleep In Your Arms Kiss Her In The Rain Tell Her You Love Her But Waitheres The Catchyou Actually Have To Mean It          You Need Anger Management Classes You Need Shut The Fukk Up Classes          If I'm With You Then I'm With You I Don't Want Anybody Else          Bitch You Ain't Pretty Your A Slut That's Why All The Guys Talk To You          I'm Sorry I Love You</p>	<p>Apples To Apples The Helen Keller Card          Deliberately Driving Slower When Being Tailgated          Watching Peoples Lives Fall Apart Via Status Updates And News Feeds          Every Time I See You A Voice In My Head Goes Dooouuuuccccheee          Making Dirty Innuendos Out Of Perfectly Innocent Things          Gene Wilder          I Hate It When You're With Mc Hammer And He Doesn't Let You Touch Anything          I Immediately Look In My Rear-view Mirror When I Pass A Cop          The Joy Of Painting With Bob Ross</p>	<div>No</div>
<div>Drugs Use</div> <div>Yes</div>	<p>Causes.com          Big Mommas Movies          No You Ask          I Like Lyrics That Actually Mean Something          Austin Texas          That Awkward Moment When You Get In The Van And There's No Candy          Texting With Cold Hands Is Like Typing In Slow Motion          Dragging Your Blanket Around The House With You Because You're Cold          Relationships Should Be Between Two People Not The Whole Universe          Pushing Your Friends Into Random People In The Hallway</p>	<p>Swimming          Inside Jokes          So What Animal Is Your Bracelet          Awkwardly Trying To Run With A Backpack          Pau Gasol          Chocolate Chip Cookie Dough Ice Cream          Milkshakes          Sour Candy          Sliding On Floors With Your Socks On          Wouldn't It Be Ironic If You Choked On A Life Saver</p>	<div>No</div>

Smoking	Yes	Cradle Of Filth Under Armour Slayer Band Inbox 1 Makes Me Nervous Dimebag Darrell Rob Zombie I Always Accept The Terms And Conditions Without Reading Them I Bottle Everything Up Until I Finally Snap Life Is Better In Summer Screwing Around In Walmart	That Spider Is More Scared Than U Are Oh Really Did It Tell U That Honda Move Out Of The Way Children I've Been Waiting 11 Years To See Toy Story 3 FBI Open The Door No Its Cool When You Break In How To Make A Girl Smile <3<3 The Desk Able To Protect You From Fire Earthquakes And Nuclear War When Your Fortune Cookie Knows What's Up Rocky When Little Kids Are Chasing Me I Run Slow So They Think They're Fast I Drop My I-pod Then My Headphones Save Its Life	No

Table S-1. Likes characterized by the most extreme average levels for each of the numeric variables (e.g. personality traits) or most extreme frequencies of classes (e.g. being a Democrat). We used only Likes that were associated with more than 100 users.