

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Supplementary Information

Tracking historical changes in perceived trustworthiness in Western Europe using machine learning analyses of facial cues in paintings.

Correspondence to:

Nicolas Baumard (nicolas.baumard@ens.fr) & Lou Safra (lou.safra@sciencespo.fr)

This PDF file includes:

- Supplementary Methods
- Supplementary Figures. 1 to 8
- Supplementary Tables 1 to 4

25 **Supplementary Methods**

26

27 In order to quantify perceived trustworthiness displays in historical paintings, we developed an
28 algorithm automatically estimating perceived trustworthiness from faces. Our algorithm also
29 extracted perceived dominance since perceived dominance has been shown to be, together with
30 perceived trustworthiness, one of the main dimensions of social perception ¹. Crucially, although
31 dominance displays carry signals of power that are distinct from the cooperation-related signals
32 associated with trustworthiness displays, perceived dominance and perceived trustworthiness are
33 correlated ¹. This correlation entails that it is of paramount importance to control for perceived
34 dominance when analyzing perceived trustworthiness. This type of analysis, studying together
35 distinct but related social signals, has already been shown to be particularly promising in the
36 emotion domain by revealing the importance of taking into account the existence of compound
37 emotions ².

38

39 Construction and validation of an algorithm for modeling perceived trustworthiness and 40 perceived dominance evaluations

41

42 We built a model that automatically extracts evaluations of perceived trustworthiness and
43 perceived dominance from the all the facial action units detected by the OpenFace algorithm (i.e.,
44 both dichotomous and continuous estimations; OpenFace version 1.01 using OpenCV 3.3.0 ³). To
45 do so, we extracted the facial action units of five sets of avatars previously generated with Facegen
46 and controlled for perceived dominance, for perceived trustworthiness or for both (Supplementary
47 Figure 1) ⁴. Each avatar is generated from an initial face and manipulated to either express a
48 specific level of perceived dominance, perceived trustworthiness or both based on the model
49 developed by Oosterhof & Todorov ¹. These avatar faces have been shown to successfully elicit
50 ratings of perceived dominance and perceived trustworthiness in participants ⁴⁻⁶. Thus, compared
51 to participants' ratings on photographs that may be sensitive to the participants characteristics and
52 to experimental protocol factors (such as the type of scale used to give the ratings), using avatars
53 allow us to have well-validated sets of faces to train our model. These sets of avatars correspond
54 to all the existing and available validated avatars controlled for perceived trustworthiness or
55 perceived dominance and generated by Facegen.

56

57 More precisely, one set of avatars was generated from one single face and manipulated for both
58 perceived dominance and perceived trustworthiness ($N = 49$; 7 levels of perceived dominance and
59 7 levels of perceived trustworthiness, each of the 7 levels corresponds to a standard deviation in
60 Oosterhof and Todorov's ¹ model ranging between -3 to +3 SD; set 1). Two other sets of faces
61 correspond to 25 maximally distinct faces manipulated either on perceived trustworthiness only
62 ($N = 175$; 7 different levels of perceived trustworthiness; set 2) or perceived dominance only ($N =$
63 175 ; 7 different levels of perceived dominance; set 3). Finally, the two last sets are composed of
64 25 Caucasian faces manipulated to present the same 7 levels of perceived trustworthiness ($N =$
65 175 ; set 4) or of perceived dominance ($N = 175$; set 5). Thus, three sets of avatars were used to
66 build the model automatically extracting perceived trustworthiness levels (sets 1, 2 and 4) and
67 three were used to build the model automatically extracting perceived dominance levels (sets 1, 3
68 and 5).

69



70
71 **Supplementary Figure 1** Sample of the avatar faces used for the algorithm optimization. **Left.** Face for the set of
72 avatars controlled for perceived dominance and perceived trustworthiness; **Middle.** Example of a face for one of the
73 sets of avatars controlled for perceived dominance only and one of the sets controlled for perceived trustworthiness
74 only; **Right.** Example of a face of the ‘Maximally distinct faces’ for the other set of avatar controlled for perceived
75 dominance only and for the other set of avatars controlled for perceived trustworthiness only. These three images were
76 created by Prof. Alexander Todorov’s team and is shared under license CC BY.

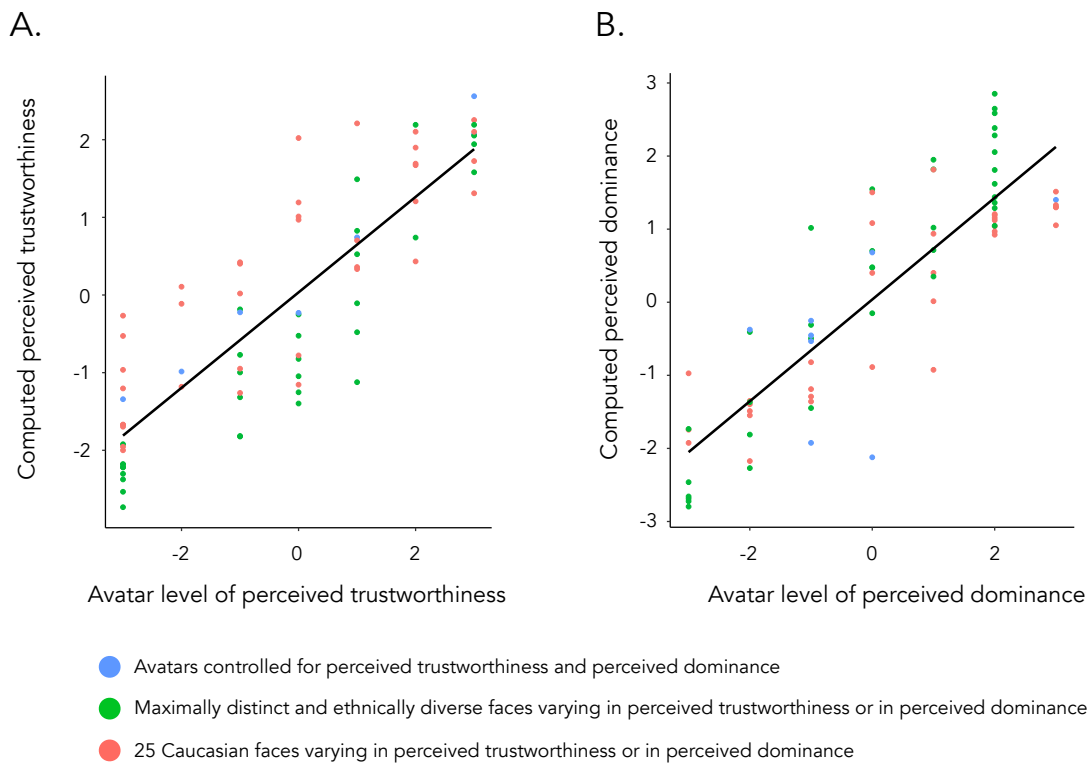
77
78 Because all our avatars were generated using the same models for perceived trustworthiness and
79 perceived dominance, actions units with a variance inferior to 0.01 were discarded as not
80 informative enough regarding cues of perceived trustworthiness and perceived dominance. The
81 reason was that they were either too low in frequency or too low in intensity (ten action units
82 discarded over thirty-three in both the perceived trustworthiness and perceived dominance avatar
83 sets).

84

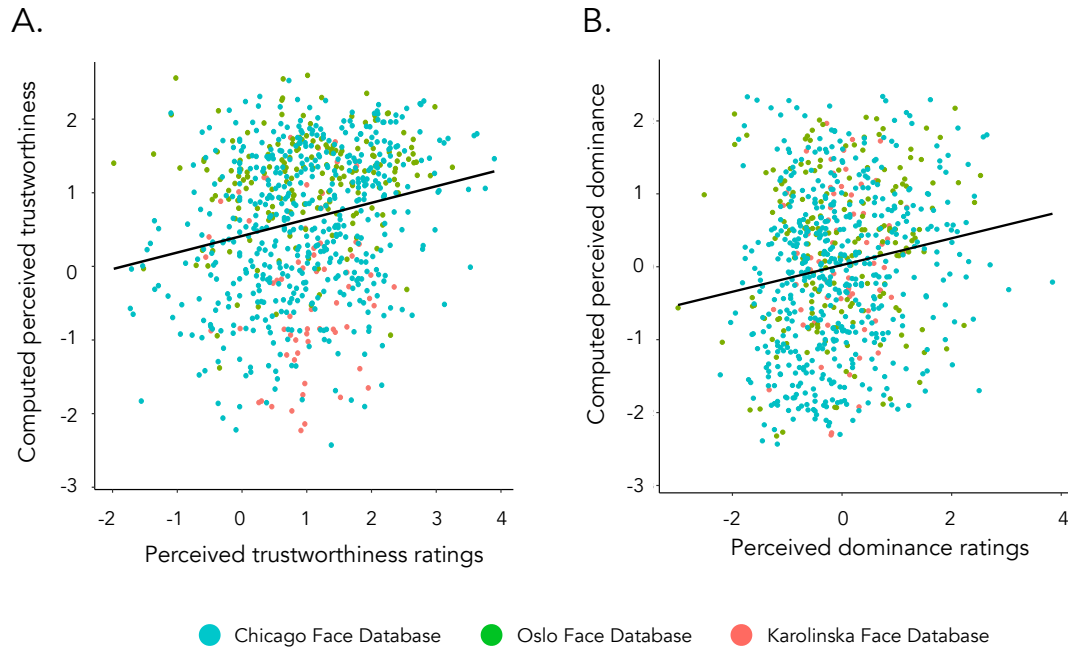
	SVM linear	SVM radial	Random forest	Linear model
Hyperparameters	Cost (C)	Cost (C) & sigma	<i>mtry</i>	\emptyset
	Perceived trustworthiness			
Mean absolute error	0.88 ± 0.02	0.87 ± 0.02	0.82 ± 0.01	0.87 ± 0.01
Root mean squared deviation	1.10 ± 0.02	1.05 ± 0.02	0.99 ± 0.01	1.06 ± 0.02
R squared	0.71 ± 0.01	0.74 ± 0.01	0.78 ± 0.01	0.72 ± 0.01
	Perceived dominance			
Mean absolute error	0.92 ± 0.02	0.79 ± 0.02	0.80 ± 0.01	0.90 ± 0.02
Root mean squared deviation	1.14 ± 0.02	0.99 ± 0.02	0.98 ± 0.02	1.11 ± 0.02
R squared	0.68 ± 0.01	0.76 ± 0.01	0.77 ± 0.01	0.70 ± 0.01

85 **Supplementary Table 1.** Model selection for extracting evaluations of perceived trustworthiness and perceived
86 dominance. Three indices of fit were computed, two which minimization indicates a better fit (mean absolute error
87 and root mean squared deviation) and one which maximization indicates a better fit (R squared). The random forest
88 was outperforming the linear model and the linear support vector model in the three indices of fit tested: mean
89 absolute error, root mean squared deviation and r-squared. The random forest model was better than the radial
90 support vector model for the perceived trustworthiness model and similar to the radials support vector model for the
91 perceived dominance model. Values are presented as mean ± standard error to the mean. Source data are provided
92 as raw data and scripts on the online depository.

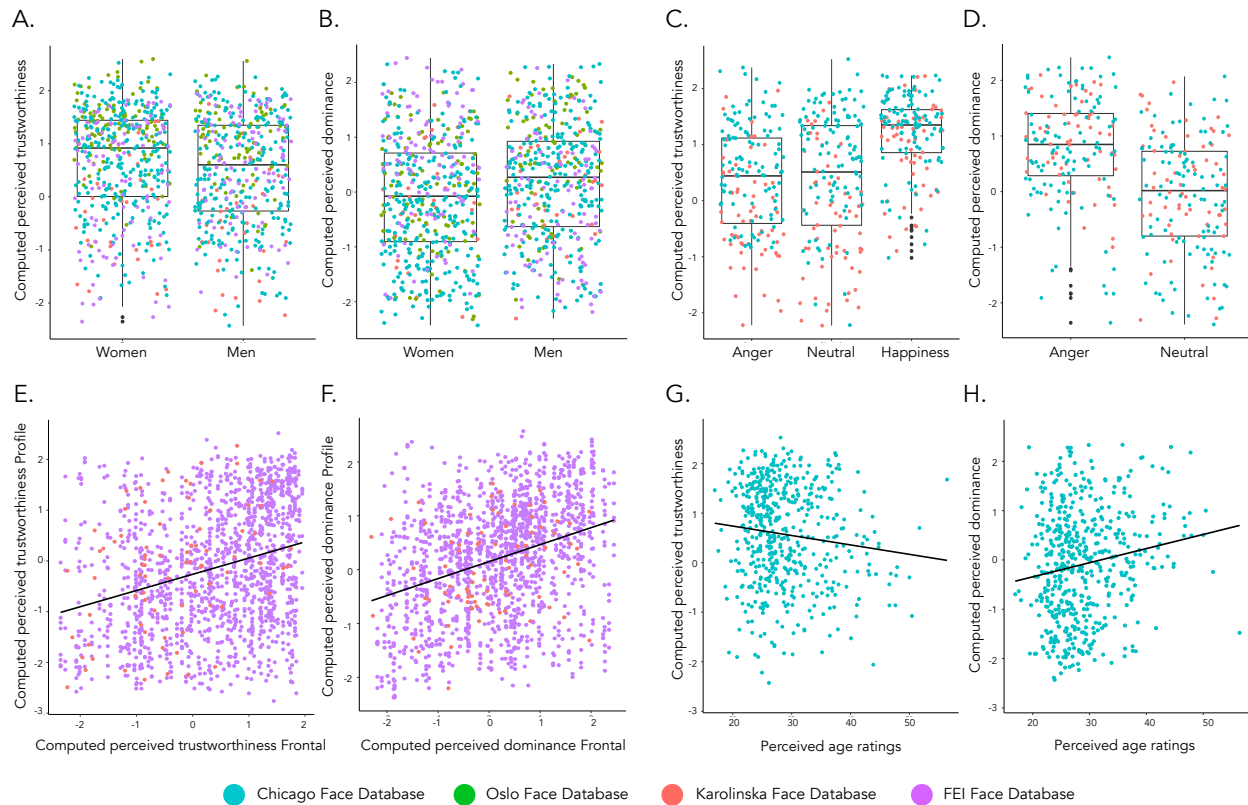
94 Based on our validation results on the avatar faces, we then trained the perceived trustworthiness
95 and perceived dominance models with the same hyperparameters on the entire avatar dataset in
96 order to increase the accuracy of our estimates and tested this model on an independent set of
97 photographs. This method differs from the classical train-test split used in machine learning which
98 was not applicable given that each avatar of our dataset presented unique features in terms of
99 luminance, texture and face shape which was important to increase the accuracy of our algorithms.
100 However, our procedure is a highly conservative test of the validity of our models as the test set is
101 completely different and independent of the training set. This conservative method for assessing
102 the validity of the algorithms is particularly critical in the present study as our goal is to generalize
103 the estimated perceived trustworthiness and perceived dominance evaluations to historical
104 portraits, a completely different set of images than those classically used in social cognition
105 research.
106



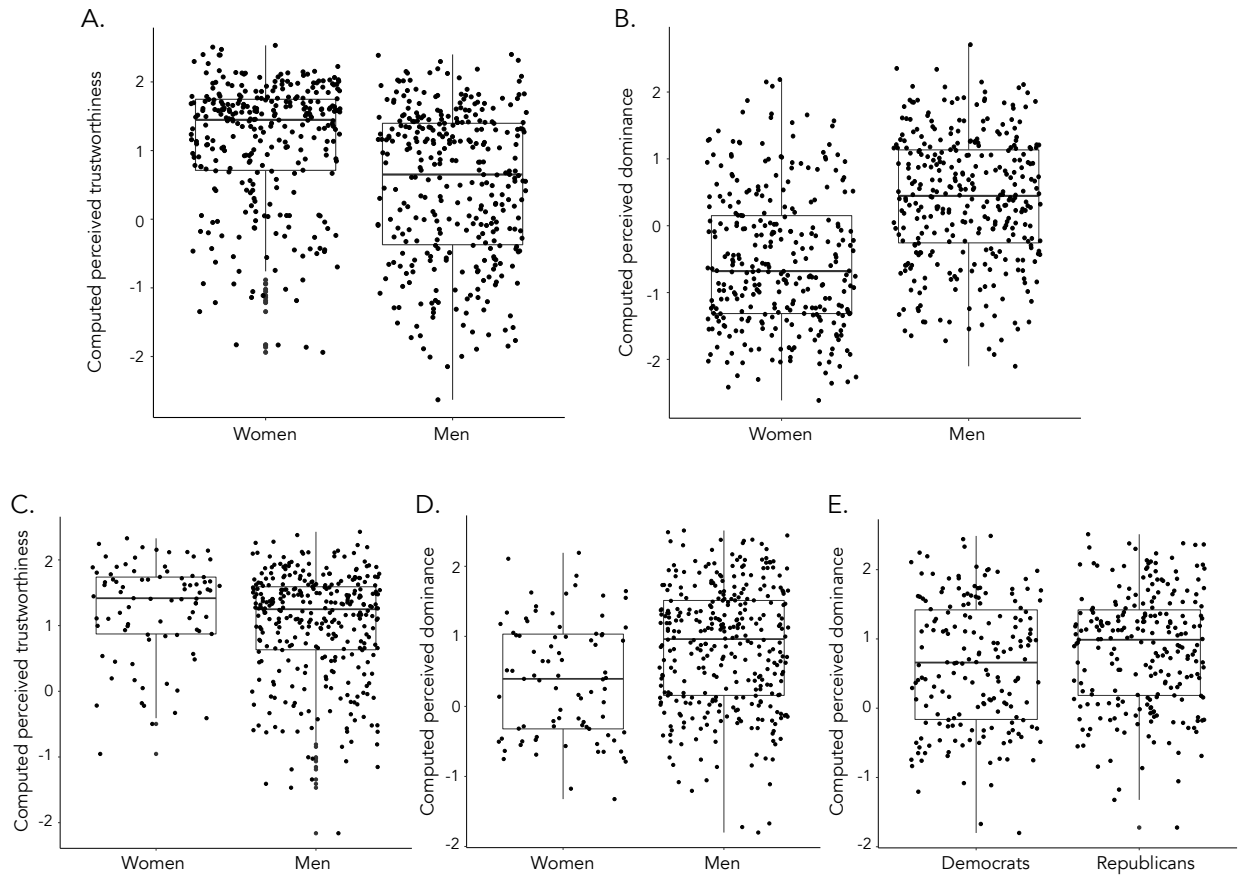
107
108 **Supplementary Figure 2** Correlation between the avatars' actual level of perceived trustworthiness and perceived
109 dominance in the test set and the computed perceived trustworthiness (A; Pearson correlation: $r = .85$, $t(75) = 14.17$,
110 $p < .001$) and perceived dominance (B; Pearson correlation: $r = .86$, $t(75) = 14.72$, $p < .001$) based on the model
111 optimized on the training set only. Source data are provided as raw data and scripts on the online depository.
112



113
 114 **Supplementary Figure 3** Correlation between participants' ratings of perceived trustworthiness and dominance
 115 displays in the three databases providing subjective ratings of perceived trustworthiness and perceived dominance (the
 116 Chicago Face Database, the Oslo Face Database and the Karolinska Face Database) and the retrieved perceived
 117 trustworthiness (**A**, Pearson correlation: $r = .22$, $t(768) = 6.19$, $p < .001$) and retrieved perceived dominance (**B**,
 118 Pearson correlation: $r = .16$, $t(769) = 4.54$, $p < .001$) levels estimated using the Facial Action Units detected by Open
 119 Face and our random-forest model. Source data are provided as raw data and scripts on the online depository.
 120



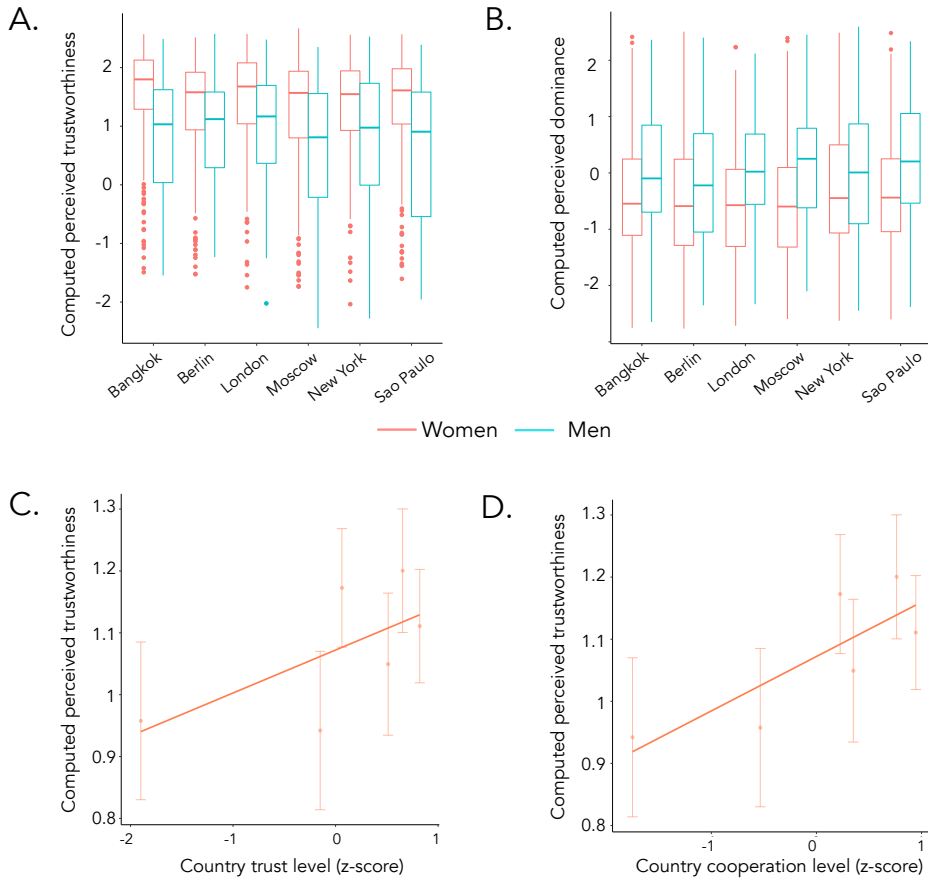
121
 122 **Supplementary Figure 4** Recovery of classical effects of gender (A-B, Student t-test: perceived trustworthiness
 123 $t(972) = 2.67, p = .008$; perceived dominance: $t(972) = -3.63, p < .001$), emotion (C-D, two-level linear regression:
 124 perceived trustworthiness: $t(167) = 10.64, p < .001$; perceived dominance: $t(167) = 9.42, p < .001$), head orientation
 125 (E-F; Pearson correlations: perceived trustworthiness $r = .29, t(1500) = 11.51, p < .001$; perceived dominance: $r = 0.34,$
 126 $t(1500) = 13.79, p < .001$) and age (G-H, Pearson correlations: perceived trustworthiness: $r = -.12, t(518) = -2.68, p =$
 127 $.008$; perceived dominance: $r = 0.16, t(518) = 3.70, p < .001$) in the perceived trustworthiness and perceived dominance
 128 estimates computed using our random forest algorithm. In the boxplots (A-D), the centre line corresponds to the
 129 median, the lower and upper bounds of the box to the 25th and 75th percentiles and the whiskers to the largest and
 130 lowest values in a limit of 1.5 times the inter-quartile range from the box bounds. Source data are provided as raw data
 131 and scripts on the online depository.
 132
 133



134
 135
 136
 137
 138
 139
 140
 141
 142
 143

Supplementary Figure 5 Results on natural images

A-B Recovery of the classical effects of gender in Google Image portraits of ‘Women’ ($N = 304$ images) and ‘Men’ ($N = 330$ images); **C-E** Recovery of the classical gender (**C-D**) and party (**E**) effects on the portraits of the House of the Representatives (women : $N = 85$ images ; men : $N = 334$ images ; democrats : $N = 182$ images ; republicans : $N = 237$ images). The centre line corresponds to the median, the lower and upper bounds of the box to the 25th and 75th percentiles and the whiskers to the largest and lowest values in a limit of 1.5 times the inter-quartile range from the box bounds. Source data are provided as raw data and scripts on the online depository.



144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

Supplementary Figure 6 Results on the Selficity Database

A-B Recovery of the classical effects of gender (Bangkok : $N = 247$ selfies of women, $N = 169$ selfies of men ; Berlin : $N = 239$ selfies of women, $N = 163$ selfies of men ; London : $N = 217$ selfies of women, $N = 134$ selfies of men ; Moscow : $N = 338$ selfies of women, $N = 82$ selfies of men ; New York : $N = 210$ selfies of women, $N = 127$ selfies of men ; Sao Paulo : $N = 231$ selfies of women, $N = 120$ selfies of men). The centre line corresponds to the median, the lower and upper bounds of the box to the 25th and 75th percentiles and the whiskers to the largest and lowest values in a limit of 1.5 times the inter-quartile range from the box bounds.; **C-D** Significant association between the country's level of interpersonal trust (**C**;) and cooperation (**D**) and the mean perceived trustworthiness estimated on the pictures of the Selficity database averaged between portraits of women and men, the red line corresponds to the effect computed in the regression controlling for the gender of the sitters (interpersonal trust: $b = 0.81 \pm 0.23$, $z = 3.50$, $p < .001$; cooperation: $b = 0.13 \pm 0.03$, $z = 3.67$, $p < .001$). Data are represented as mean values and error bars correspond to standard errors to the mean (Bangkok : $N = 416$ selfies ; Berlin : $N = 402$ selfies ; London : $N = 351$ selfies ; Moscow : $N = 420$ selfies ; New York : $N = 337$ selfies ; Sao Paulo : $N = 351$ selfies). Source data are provided as raw data and scripts on the online depository.

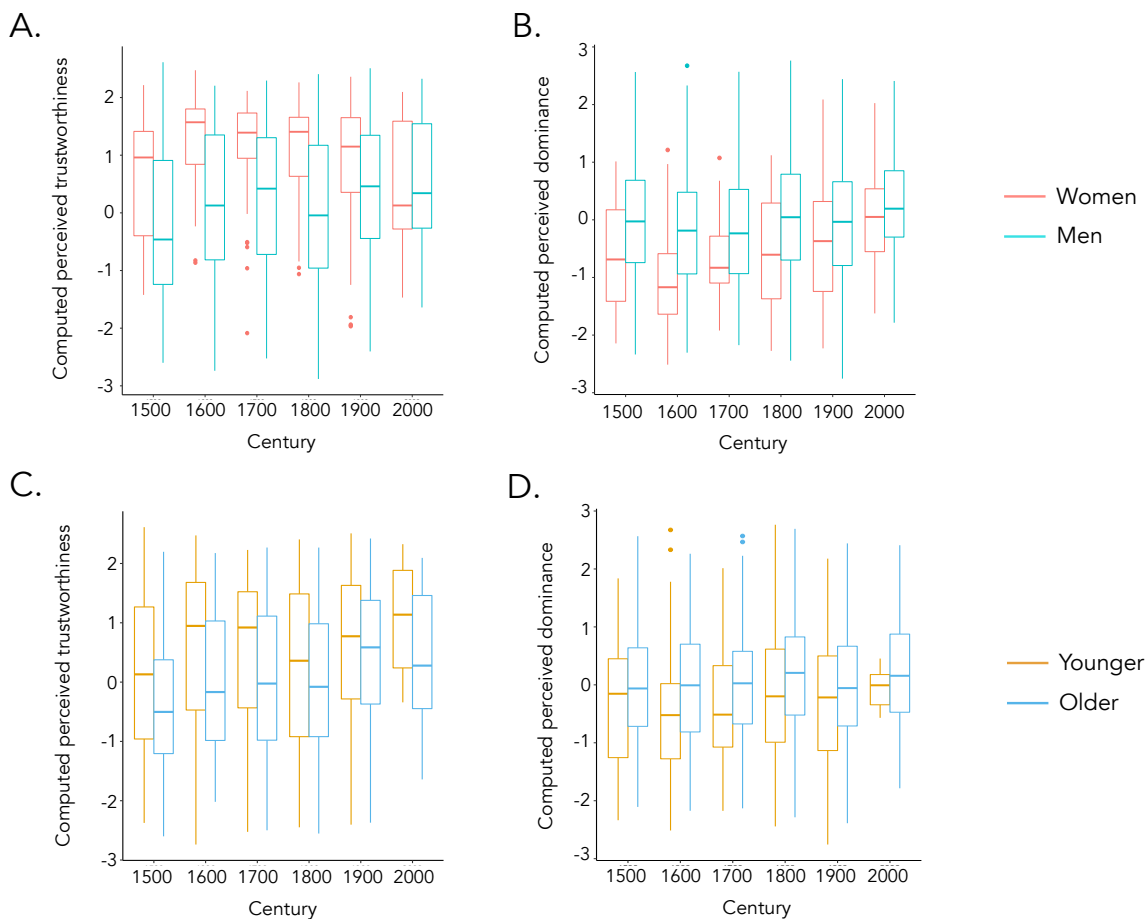
Analysis of the National Portrait Gallery and the Web Gallery of Art

Text	Code	Example
Century	Century + 50	16 th century = 1550
Late century	Centruy + 90	Late 16 th century = 1590
Early century	Century + 10	Early 16 th century = 1510
Half of century	Century + 50	Half of 16 th century = 1550
Decade+s	Decade	1650s = 1655
Around/about/perhaps/probably/circa/after + Date	Date	Circa 1655 = 1655
Date 1 – Date 2	Rounded mean of Date 1 and Date 2	1650-1655 = 1652

162 **Supplementary Table 2** – Coding of the date of the portraits

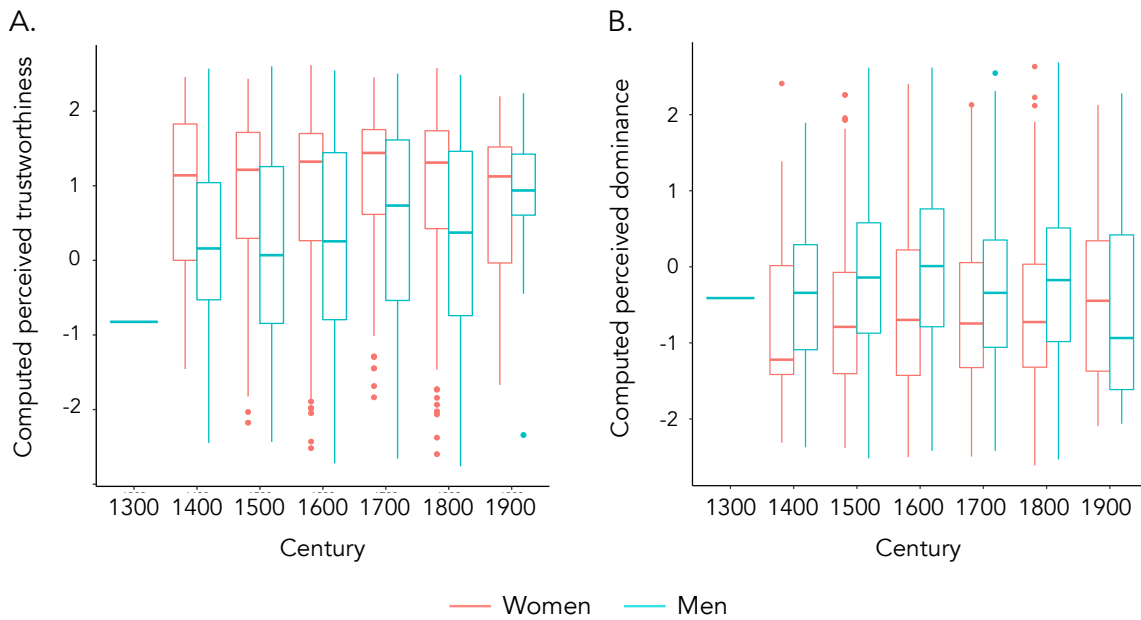
163

164 The information about the sitters’ gender and age allowed us to replicate the classic findings that
 165 older sitters appear more dominant and less trustworthy than younger sitters and that female sitters
 166 appear more trustworthy and less dominant than male sitters (perceived trustworthiness: gender
 167 effect: $t(1960) = 9.69, p < .001$; age effect: $t(1960) = -6.63, p < .001$; perceived dominance: gender
 168 effect: $t(1960) = 7.24, p < .001$; age effect: $t(1960) = -9.12, p < .001$; Supplementary Figure 7). As
 169 for the NPG, we accurately recovered the gender effect on perceived trustworthiness and d
 170 perceived dominance on the portraits of the Web Gallery of Art (perceived trustworthiness: $z =$
 171 $17.70, p < .001$; perceived dominance: $z = -13.35, p < .001$; Supplementary Figure 8).
 172



173

174 **Supplementary Figure 7** Recovery of the gender (A-B) (1500 : $N = 23$ portraits of women, $N = 68$ portraits of men ;
 175 1600 : $N = 50$ portraits of women, $N = 236$ portraits of men ; 1700 : $N = 53$ portraits of women, $N = 432$ portraits of
 176 men ; 1800 : $N = 44$ portraits of women, $N = 609$ portraits of men ; 1900 : $N = 98$ portraits of women, $N = 351$ portraits
 177 of men ; 2000 : $N = 19$ portraits of women, $N = 42$ portraits of men) and age (C-D) effects in the National Portrait
 178 Gallery database over the centuries (the ‘Younger’ category is defined as sitters being under 48 year old; 1500 : $N =$
 179 61 portraits of younger sitters, $N = 30$ portraits of older sitters; 1600 : $N = 188$ portraits of younger sitters, $N = 96$
 180 portraits of younger sitters ; 1700 : $N = 280$ portraits of younger sitters, $N = 194$ portraits of older sitters; 1800 : $N =$
 181 273 portraits of younger sitters, $N = 345$ portraits of older sitters; 1900 : $N = 187$ portraits of younger sitters, $N = 249$
 182 portraits of older sitters; 2000 : $N = 8$ portraits of younger sitters, $N = 53$ portraits of older sitters). The centre line
 183 corresponds to the median, the lower and upper bounds of the box to the 25th and 75th percentiles and the whiskers
 184 to the largest and lowest values in a limit of 1.5 times the inter-quartile range from the box bounds. Source data are
 185 provided as raw data and scripts on the online depository.



186
187
188
189
190
191
192
193
194
195

Supplementary Figure 8 Recovery of the gender effects in the Web Gallery of Art (1300 : $N = 1$ portrait of man ; 1400 : $N = 137$ portraits of men, $N = 41$ portraits of women ; 1500 : $N = 696$ portraits of men, $N = 291$ portraits of women ; 1600 : $N = 963$ portraits of men, $N = 509$ portraits of women ; 1700 : $N = 418$ portraits of men, $N = 350$ portraits of women ; 1800 : $N = 349$ portraits of men, $N = 307$ portraits of women ; 1900 : $N = 22$ portraits of men, $N = 22$ portraits of women) for perceived trustworthiness (A) and perceived dominance (B). The centre line corresponds to the median, the lower and upper bounds of the box to the 25th and 75th percentiles and the whiskers to the largest and lowest values in a limit of 1.5 times the inter-quartile range from the box bounds. Source data are provided as raw data and scripts on the online depository.

Dependent variable	Perceived Trustworthiness		GDP per capita	Democratization
Independent variable of interest	GDP per capita	Democratization	Perceived Trustworthiness	Perceived Trustworthiness
Delay Two decades				
Model comparison	$F(40,1) = 12.38$ $p = .001$	$F(15,1) = 0.11$ $p > .250$	$F(41,1) = 0.76$ $p > .250$	$F(16,1) = 6.54$ $p = .022$
Effect	$b = 0.04 \pm 0.01$ $t(40) = 3.52$ $p = .001$	$b = -0.01 \pm 0.03$ $t(14) = -0.33$ $p > .250$	$b = 0.59 \pm 0.68$ $t(41) = 0.87$ $p > .250$	$b = -5.82 \pm 2.27$ $t(15) = -2.56$ $p = .022$
Delay One decade				
Model comparison	$F(41,1) = 11.40$ $p = .002$	$F(16,1) = 1.11$ $p > .250$	$F(42,1) = 0.01$ $p > .250$	$F(17,1) = 5.26$ $p = .036$
Effect	$b = 0.03 \pm 0.01$ $t(40) = 3.38$ $p = .002$	$b = -0.02 \pm 0.02$ $t(15) = -1.05$ $p > .250$	$b = -0.05 \pm 0.66$ $t(41) = -0.08$ $p > .250$	$b = -4.19 \pm 1.82$ $t(16) = 0.64$ $p > .250$

196
197
198

Supplementary Table 3 Temporal dynamics of perceived trustworthiness, GDP per capita and democratization in the paintings of the National Portrait Gallery. Model comparison corresponds to the comparison of the model that included the delayed variable of interest with the model in which this variable was excluded. Effect corresponds to

199 the estimation of the regression coefficient of the delayed variable of interest. All the tests are two-sided. Following
 200 APA's recommendations, exact p-values are provided for p-s between .001 and .250. Source data are provided as raw
 201 data and scripts on the online depository.
 202

Dependent variable	Perceived Trustworthiness		GDP per capita	Democratization
Independent variable of interest	GDP per capita	Democratization	Perceived Trustworthiness	Perceived Trustworthiness
Delay One decade				
Model comparison	X(1) = 4.00 p = .046	X(1) = 0.01 p > .250	X(1) = 2.48 p = .115	X(1) = 0.65 p > .250
Effect	b = 0.12 ± 0.05 z = 2.61 p = .009	b = 0.00 ± 0.01 z = -0.11 p > .250	b = -0.03 ± 0.02 z = -1.56 p = .119	b = 0.38 ± 0.49 z = 0.78 p > .250
Delay Two decades				
Model comparison	X(1) = 6.42 p = .011	X(1) = 0.81 P > .250	X(1) = 2.02 p = .155	X(1) = 0.72 p > .250
Effect	b = 0.19 ± 0.06 z = 3.48 p < .001	b = -0.01 ± 0.01 z = -0.84 p > .250	b = -0.05 ± 0.04 z = -1.42 p = .157	b = 0.45 ± 0.55 z = 0.82 p > .250

203 **Supplementary Table 4** Temporal dynamics of perceived trustworthiness, GDP per capita and democratization in
 204 the paintings of the Web Gallery of Art. All the tests are two-sided. Following APA's recommendations, exact p-
 205 values are provided for p-values between .001 and .250. Source data are provided as raw data and scripts on the
 206 online depository.
 207
 208
 209

	Affluence only		Time + Affluence		Armed conflict only		Time + Armed conflict	
	National Portraits Gallery	Web Gallery of Art	National Portraits Gallery	Web Gallery of Art	National Portraits Gallery	Web Gallery of Art	National Portraits Gallery	Web Gallery of Art
year			.11±.02 z = 5.46 p < .001	.05±.01 z = 3.54 p < .001			.14±.02 z = 7.55 p < .001	.05±.01 z = 4.13 p < .001
Number of book titles per capita	.35±.06 z = 6.15 p < .001	.29±.10 z = 2.77 p = .006	.21±.06 z = 3.45 p = .001	.14±.11 z = 1.26 p = .208				
Presence of an armed conflict					.01±.05 z = 0.30 p > .250	.00±.03 z = -0.01 p > .250	.05±.05 z = 1.05 p > .250	-.01±.03 z = -0.39 p > .250
Control variables								
Perceived dominance	-.78±.02 z = -40.10 p < .001	-.75±.02 z = -54.29 p < .001	-.79±.02 z = -40.85 p < .001	-.74±.01 z = -54.13 p < .001	-.78±.02 z = -39.79 p < .001	-.74±.01 z = -54.85 p < .001	-.79±.02 z=-40.74 p < .001	-.74±.02 z=-54.86 p < .001
Gender	-.31±.06 z = -5.27 p < .001	-.33±.03 z = -11.13 p < .001	-.29±.06 z = -5.09 p < .001	-.32±.03 z = -10.52 p < .001	-.37±.06 z = -6.41 p < .001	-.33±.03 z = -11.51 p < .001	-.33±.06 z = -5.68 p < .001	-.31±.03 z = -10.49 p < .001

Age	-.00±.00 z = -1.35 p = .178		-.00±.00 z = -2.49 p = .013		.00±.00 z = 0.21 p > .250		-.00±.00 z = -2.01 p = .044	
Sample								
N	1962	3801	1962	3801	1962	3927	1962	3927

210
211 **Supplementary Table 5** Replication analyses on perceived trustworthiness in the National Portrait Gallery and the
212 Web Gallery of Art using the Number of book titles per capital as a proxy of affluence as well as the presence of
213 armed conflict as indicator of periods of war and social unrest.

214 The first line corresponds to the regression coefficient with their associated standard error to the mean (mean ± s.e.m.).
215 Results in bold corresponds to statistically significant effects of the variables of interest. The upper part of the table
216 presents the effects of the variables of interest (time, affluence and democratization), while the lower part presents the
217 effects of the control variables (perceived dominance, gender and age). All the tests are two-sided. Following APA's
218 recommendations, exact p-values are provided for p-values between .001 and .250. Source data are provided as raw
219 data and scripts on the online depository.

220
221 **Copyright of the analysed databases**

222 All the exploited databases (Prof. Todorov's avatar datasets, Karolinska database, Oslo Face
223 database, Chicago Face database, FEI Face database, the National Portrait Gallery database and
224 the Web Gallery of Art database) are free of use for non-commercial research purposes. The use
225 of the Selfiecity database has been authorized by its owner, Dr. Lev Manovuch.

226
227 **Supplementary References**

- 228 1. Oosterhof, N. N. & Todorov, A. The functional basis of face evaluation. *Proc. Natl. Acad.*
229 *Sci.* **105**, 11087–11092 (2008).
230 2. Du, S., Tao, Y. & Martinez, A. M. Compound facial expressions of emotion. *Proc. Natl.*
231 *Acad. Sci.* **111**, E1454–E1462 (2014).
232 3. Baltrušaitis, T., Robinson, P. & Morency, L. OpenFace: An open source facial behavior
233 analysis toolkit. in *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)* 1–
234 10 (2016). doi:10.1109/WACV.2016.7477553.
235 4. Todorov, A., Dotsch, R., Porter, J. M., Oosterhof, N. N. & Falvello, V. B. Validation of
236 data-driven computational models of social perception of faces. *Emotion* **13**, 724–738 (2013).
237 5. Stewart, L. H. *et al.* Unconscious evaluation of faces on social dimensions. *J. Exp. Psychol.*
238 *Gen.* **141**, 715–727 (2012).
239 6. Safra, L., Ioannou, C., Amsellem, F., Delorme, R. & Chevallier, C. Distinct effects of social
240 motivation on face evaluations in adolescents with and without autism. *Sci. Rep.* **8**, 1–8 (2018).

241