

1 Supplementary materials for ‘Novel embeddings
2 improve the prediction of risk perception’

3 Zak Hussain^{1,2*}, Rui Mata^{1†} and Dirk U. Wulff^{2,1}

4 ^{1*}Faculty of Psychology, University of Basel, Missionsstasse 62a, Basel,
5 4055, Switzerland.

6 ²Center for Adaptive Rationality, Max Planck Institute for Human
7 Development, Lentzeallee 94, Berlin, 14195, Germany.

8 *Corresponding author(s). E-mail(s): z.hussain@unibas.ch;
9 Contributing authors: rui.mata@unibas.ch; wulff@mpib-berlin.mpg.de;

10 †These authors contributed equally to this work.

11 **1 Data collection**

12 We followed common procedures used in the risk perception literature to obtain data
13 for the psychometric paradigm [e.g., 1, 2]. The pre-registration for the study is available
14 at <https://osf.io/6m7xr>. In what follows, we investigate the sensitivity of our results to
15 various factors surrounding data collection. We focus on two main factors: the impact
16 of psychometric item ordering—which could affect both predictive accuracy and inter-
17 item correlations—and the impact of training set size (with a focus on predictive
18 accuracy).

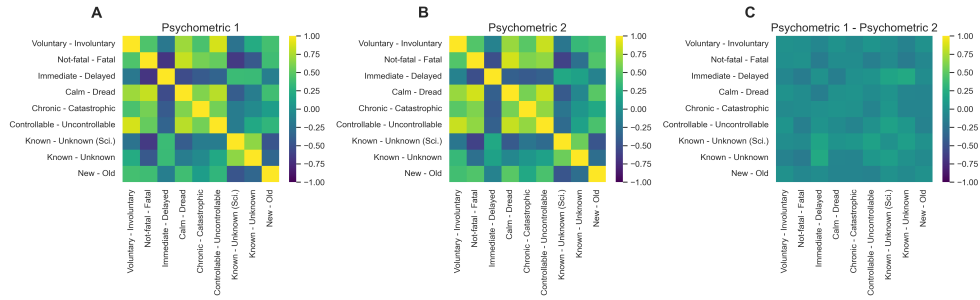


Fig. 1 Investigating the impact of psychometric item ordering on inter-item correlations. **A.** Order 1 inter-item correlations. **B.** Order 2 inter-item correlations. **C.** Order 1 minus order 2 inter-item correlations.

19 1.1 Impact of psychometric item ordering

20 In our survey, half of the participants received the psychometric items in the order
 21 presented below (order 1) for each risk and the other half received them in the reverse
 22 order (order 2). The main reason for doing this was to investigate whether ordering
 23 actually impacts participant responses, which, to our knowledge, has not been done
 24 before, and could affect data quality.

- 25 1. **Voluntary–Involuntary**—Are individuals exposed to this risk voluntarily or
 26 involuntarily?
- 27 2. **Immediate–Delayed**—Is death from this risk immediate or delayed?
- 28 3. **Known-Unknown**—Is this risk known or unknown to the individuals exposed to
 29 this risk?
- 30 4. **Known–Unknown (Sci.)**—Is this risk known or unknown to science?
- 31 5. **Controllable–Uncontrollable**—Is this risk controllable or uncontrollable for the
 32 individual exposed to the risk?
- 33 6. **New–Old**—Is this risk new or old?
- 34 7. **Chronic–Catastrophic**—Is this a risk that kills one person at a time (chronic)
 35 or a risk that kills large numbers of people at once (catastrophic)?

36 8. **Calm–Dread**—Is this a risk that individuals can reason about calmly or is it one
37 that they have great dread for?

38 9. **Not-fatal–Fatal**—How fatal are the consequences of this risk?

39 To evaluate potential differences between the two orderings, we carried out several
40 analyses. First, we focus on the psychometric ratings alone. To investigate whether
41 psychometric ordering had a statistically significant impact on responses, we take the
42 individual ratings for each risk source and psychometric item, split them into two
43 groups (order 1 and order 2), and run an independent-samples t-test on each pair of
44 groups. This amounted to 9,036 t-tests (1,004 risks times 9 psychometric items), of
45 which 11.6% of the groups significantly differed for $\alpha = .05$. This is twice the number
46 of type I errors expected, suggesting a small influence of ordering on average responses.
47 Four out of the nine items (*Immediate–Delayed*, *Voluntary–Involuntary*, *Calm–Dread*,
48 and *Known–Unknown*) account for almost 60% of the significant differences. However,
49 overall, the difference in the average responses was small (average Cohen’s $d = .09$).
50 Furthermore, the average ratings in the nine psychometric items showed very high
51 Pearson correlations of, on average, 0.88.

52 We further evaluated the robustness of the inter-item correlation between the
53 two orderings because this has implications for the sensitivity of principal com-
54 ponent analyses (PCA) often performed within the psychometric paradigm [cf. 1].
55 Figure 1 shows the correlations across risks between psychometric item ratings for
56 both orderings. We observed very similar patterns of correlations but also small dif-
57 ferences ranging from $\delta < .001$ (*Immediate–Delayed* and *Chronic–Catastrophic*) to
58 $\delta = .21$ (*Immediate–Delayed* and *Known–Unknown*), with an overall average absolute
59 difference of $\delta = .08$.

60 Finally, we evaluated potential differences in the accuracy of predicting risk per-
61 ception (See Figure 2). We observed that *Psychometric 2* achieved a 6.4 percentage

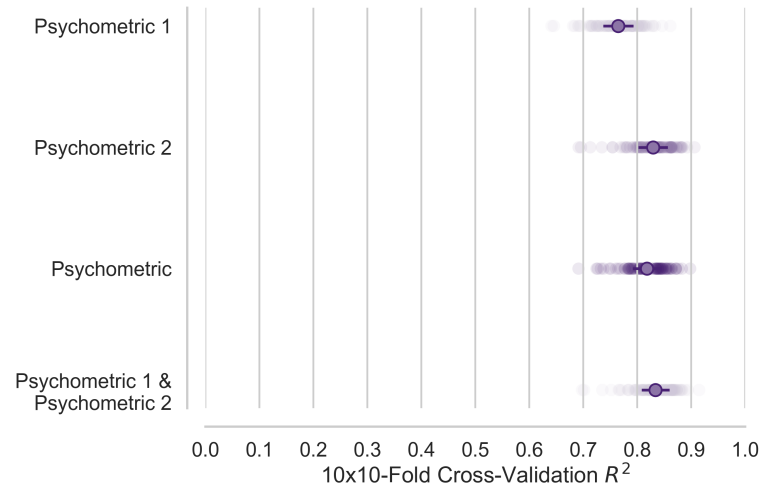


Fig. 2 Investigating the impact of psychometric item ordering on performance. *Psychometric 1* is obtained from participants that received the following order 1 (as listed in text). *Psychometric 2* participants received the reverse order. *Psychometric* is an aggregate of orders 1 and 2 (as used in the main analysis), and *Psychometric 1 & Psychometric 2* is the concatenation of both orderings. Error bars are adjusted 95% confidence intervals [3].

62 points higher accuracy than *Psychometric 1* and a 1.1 percentage points higher accu-
 63 racy than the aggregate psychometric model using both orders. This means that the
 64 reversed order is better at capturing risk perception than the original order. This
 65 may have contributed to the higher performance of the psychometric model in the
 66 Basel Risk Norms compared to the data of [2] because the latter relied only on the
 67 first ordering. The notable differences in predictive accuracy between the two orders
 68 have two noteworthy implications. First, other orderings of psychometric items could
 69 result in even larger predictive accuracy for the psychometric model. Second, the two
 70 orderings may capture distinct aspects of risk perception, suggesting that they might
 71 best be used in tandem rather than aggregated. To test the latter, we evaluated the
 72 concatenation of both orderings, *Psychometric 1 & Psychometric 2*, as a predictive
 73 model. We observed that the concatenated model outperformed the aggregate model
 74 by 1.6 percentage points, which is a small but significant effect ($t = 4.00, p < .001$).

75 Overall, our evaluation of orderings revealed some differences in average ratings,
76 inter-item correlations, and predictive accuracy. However, the differences between
77 orderings were overall small in magnitude. Furthermore, although the slightly higher
78 accuracy of the concatenated model compared to the aggregate model may justify
79 using the concatenated from the perspective of predictive accuracy, this choice would
80 disadvantage our analysis in other ways. Specifically, it would limit interpretability,
81 given that we possess no information on how the item ordering affects the content of
82 the responses to the psychometric items, and comparability to previous work includ-
83 ing, in particular, the study by [2]. We believe that the small gains in accuracy do not
84 outweigh these costs, and so chose to use the aggregate model.

85 **1.2 The impact of training set size on predictive accuracy**

86 In planning the data collection of the Basel Risk Norms, we investigated the poten-
87 tial of increasing predictive accuracy by increasing the training set size. We trained
88 different models on different portions of the data of [2] and recorded the accuracy
89 of predicting risk perception (see Figure 3; green lines). The analysis showed signifi-
90 cant potential for higher accuracy, with accuracy values increasing systematically with
91 larger training set sizes. The increasing accuracy is likely due to a decreasing role of
92 model overfitting. This potential for increased accuracy suggested by the reanalysis of
93 the data of [2] was largely realized by the larger Basel Risk Norms. Figure 3 also shows
94 the accuracies of the different models for the Basel Risk Norms, which demonstrate
95 clear performance increases for the larger training set sizes.

96 Three additional results concerning the relationship between training set size and
97 predictive accuracy in the Basel Risk Norms are worth noting. First, the accuracies
98 appear to taper off for larger training set sizes. One important implication of this
99 is that comparisons between the low-dimensional psychometric model and the high-
100 dimensional embeddings models are fairer using the larger Basel Risk Norms. Second,

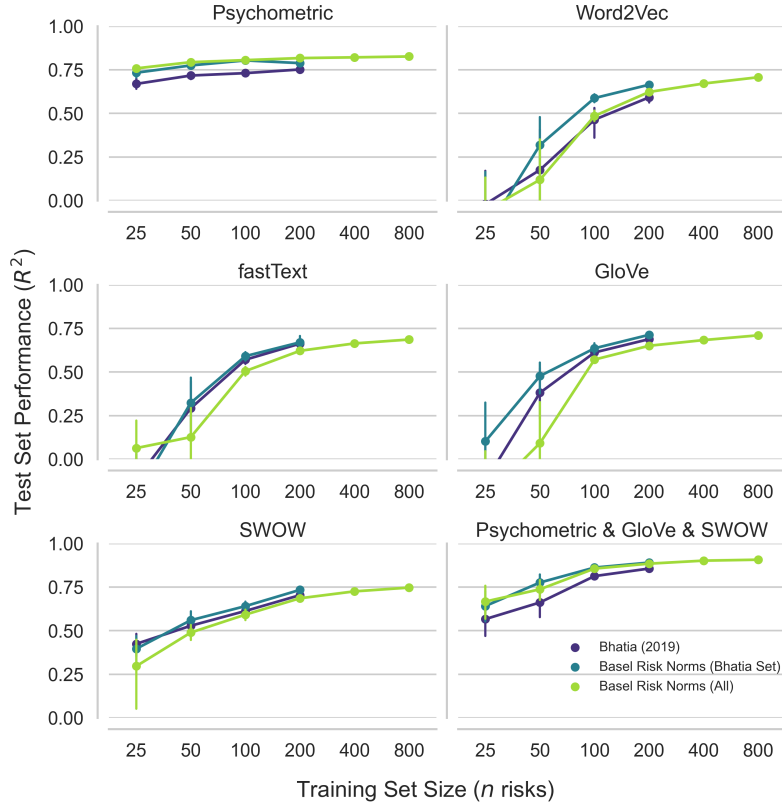


Fig. 3 Evaluating how test performance varies with training set size for 3 data (sub-)sets: (i) Basel Risk Norms (All), which refers to our full data set of 1,004 risks, (ii) Basel Risk Norms (Bhatia Set), referring to our data limited to the same 306 risks as used in [2], and [2] (iii). Test sets are composed of all remaining risks in the data. Train–test splits were sampled randomly (i.e., bootstrapped), with 10 repetitions per model per training set size. Error bars are 95% confidence intervals.

101 the accuracy of the psychometric model is systematically higher for the Basel Risk
 102 Norms compared to the data of [2] for any training set size. This difference likely
 103 reflects the substantial increase in reliability due to a larger number of ratings. Third,
 104 embedding accuracies for small training sets are worse for the Basel Risk Norms than
 105 the data of [2] when considering all risks and better when considering only the risks
 106 shared across data sets. These results are consistent with the higher risk rating reli-
 107 abilities of the Basel Risk Norms but also suggest that the newly introduced risks may
 108 result in a larger diversity of risks, making it harder to generalize from train to test set.

109 Overall, by increasing the size of the risk set, we boosted the performance of all
110 models thus permitting a fairer comparison of model performance due to less model
111 overfitting.

112 **2 Model comparison**

113 In this section, we provide additional information concerning the sensitivity of our
114 model comparison results to various analytic choices. We first justify our decision to
115 focus only on the results of a linear regression algorithm (elastic net) in the main paper,
116 instead of more flexible nonlinear methods such as the popular gradient boosting. We
117 next motivate our decision to use a groupwise scaling technique during pre-processing,
118 instead of more traditional approaches to scaling preceding regularized regression such
119 as standardization. Finally, we provide a comprehensive statistical analysis of the
120 differences between all pairwise model combinations for completeness.

121 **2.1 Elastic net versus gradient boosting**

122 In addition to elastic net regression, we evaluated the predictive accuracies of the
123 different models using Scikit-Learn’s gradient boosting regressor [4]. Gradient boosting
124 is a popular nonlinear algorithm that builds an additive model out of regression trees
125 in a forward stagewise fashion. In many cases, gradient boosting can outperform linear
126 models, especially when more training samples are available.

127 We observed that for all but one model gradient boosting was at best equal and, in
128 many cases, clearly worse than the linear model (see Figure 4). The only exception was
129 the low-dimensional psychometric model, which saw a small increase in the predictive
130 accuracy of 2.6 percentage points on the Basel Risk Norm data. Interestingly, we also
131 see the impact of the increased training set size, with the additional risks in our norm
132 set reducing the relative advantage of elastic net over gradient boosting. This indicates

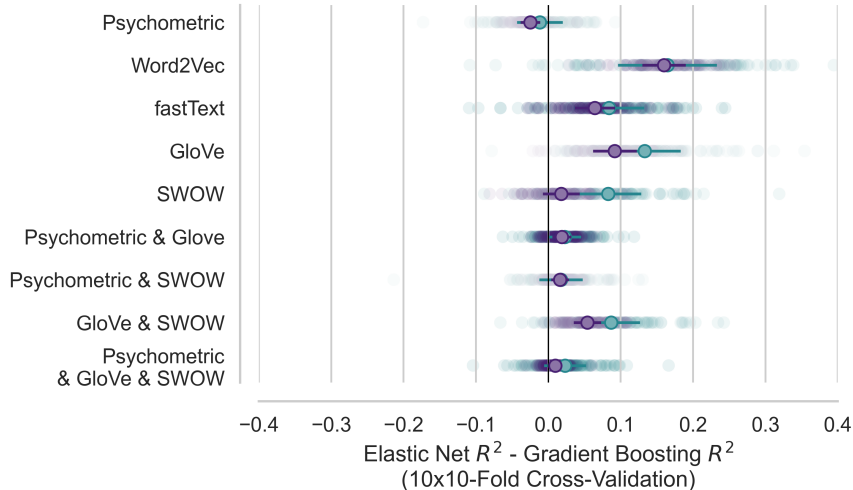


Fig. 4 Pairwise differences between elastic net and gradient boosting using 10x10-fold cross-validation. Cross-validation via [2]’s risk norms (306 risks) are colored cyan and points obtained using the Basel Risk Norms (1004 risks) are colored purple. Error bars are adjusted 95% confidence intervals [3].

133 that perhaps with a sufficient number of samples, the more flexible gradient boosting
 134 model could outperform elastic net.

135 Overall, regularized linear regression emerged as the superior model, which is
 136 consistent with the relatively low ratio of data points to features.

137 2.2 Evaluating embedding scaling approaches

138 When relying on regularization techniques, such as elastic net regularization, it is
 139 common practice to standardize the predictors to even out their contribution to the
 140 regularization penalty. However, we based our analysis on unstandardized embeddings.
 141 We did this to allow for a fair comparison between the free-association and text embed-
 142 dings. The free associations embedding (*SWOW*) was trained using singular value
 143 decomposition, which by design allocates variance very unevenly across the embedding
 144 dimensions. Standardizing *SWOW* would thus imply removing an important prior on
 145 the importance of embedding dimension, which can result in reduced predictive accu-
 146 racy. To quantify the potential negative effect of standardization on *SWOW* and a

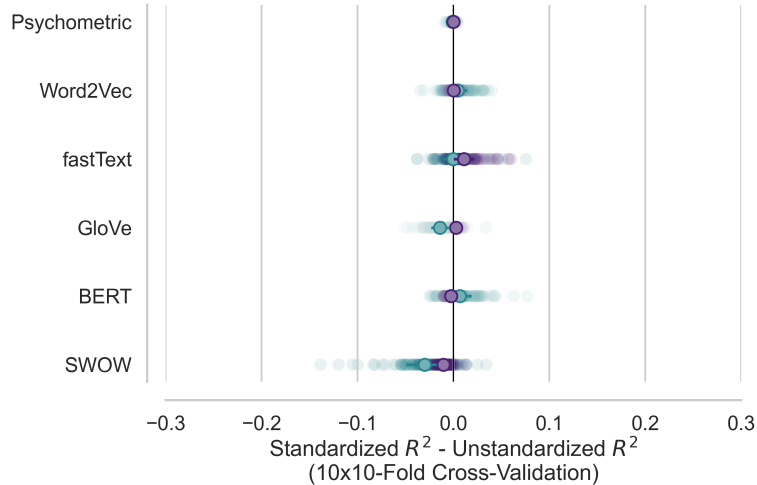


Fig. 5 Pairwise differences between standardized and unstandardized models using 10x10-fold cross-validation and elastic net regression. Cross-validation via [2]’s risk norms (306 risks) are colored cyan and points obtained using the Basel Risk Norms (1004 risks) are colored purple. Error bars are adjusted 95% confidence intervals [3].

147 potentially positive effect for the other embedding models, we explicitly compared the
 148 predictive accuracy for every model with standardized and unstandardized dimensions
 149 for both risk norm sets (Bhatia, 2019, and Basel Risk Norms).

150 As can be seen in Figure 5, standardizing did indeed negatively impact the
 151 *SWOW*) for both norm sets (Bhatia, 2019: $t = -3.09, p = .003$, Basel Risk Norms:
 152 $t = -3.68, p < .001$). In terms of the text embeddings, the effect of standardizing
 153 was mixed, with a negative effect for *GloVe* on [2]’s data ($t = -3.04, p = .003$), and
 154 smaller positive effects on the Basel Risk Norms for *GloVe* ($t = 2.26, p = .027$) and
 155 *fastText* ($t = 2.05, p = .043$). *Psychometric* was not significantly affected. In light of
 156 these findings, we chose not to standardize the models in our analysis.

157 2.3 Statistical tests

158 The comparison of models was carried using the procedure described in [3] (see also,
 159 [5]). It involves calculating the differences in model performance across the same 100
 160 (10x10) train-test splits for each pair of models and testing the null hypothesis that

161 the mean difference equals zero using an adjusted paired t-test that accounts for the
162 dependence between train-test splits.

163 To give an overview of all possible model comparisons, Figure 6 shows the differ-
164 ences in R-squared predictive accuracy for all pairs of individual and ensemble models
165 (y-axis models minus x-axis models) with nonsignificant differences displayed as white.

166 Several important insights emerge from the patterns of results. First, the patterns
167 of results are highly similar between the data of [2] and the Basel Risk Norms, with one
168 exception being the large number of significant results for Basel Risk Norms due to the
169 higher reliability and larger data set size. Second, ensembles containing the psychome-
170 tric model outperform ensembles without the psychometric model, as indicated by the
171 strong bright rectangle in the bottom left corners. Third, there is only one model not
172 significantly different from the psychometric model—*GloVe & SWOW*—attesting to
173 the strong performance of *SWOW* in capturing important aspects of risk perception.

174 3 Word norms

175 Our interpretability analysis identified unaccounted dimensions of risk by relying on
176 a set of word norms. For this purpose, we selected a set of norms from [6] that we
177 hypothesized to be related to risk perception. Table 1 provides an overview of these
178 norms and lists the individual sources. As reported in the main text, these norms are
179 able to predict 64.3% of risk perception variance (with 32% of the norm data imputed
180 using *Word2Vec* to deal with missing norm data on certain risks), establishing their
181 usefulness for revealing the key aspects of risk perception.

182 References

- 183 [1] Fischhoff B, Slovic P, Lichtenstein S, et al (1978) How safe is safe enough? a
184 psychometric study of attitudes towards technological risks and benefits. *Policy*
185 *Sciences* 9(2):127–152. URL <https://doi.org/10.1007/bf00143739>

- 186 [2] Bhatia S (2019) Predicting risk perception: New insights from data science.
187 Management Science 65(8):3800–3823. URL [https://doi.org/10.1287/mnsc.2018.](https://doi.org/10.1287/mnsc.2018.3121)
188 [3121](https://doi.org/10.1287/mnsc.2018.3121)
- 189 [3] Bouckaert RR, Frank E (2004) Evaluating the replicability of significance tests
190 for comparing learning algorithms. In: Pacific-Asia Conference on Knowledge
191 Discovery and Data Mining, Springer, pp 3–12, URL [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-540-24775-3_3)
192 [978-3-540-24775-3_3](https://doi.org/10.1007/978-3-540-24775-3_3)
- 193 [4] Pedregosa F, Varoquaux G, Gramfort A, et al (2011) Scikit-learn: Machine
194 learning in Python. Journal of Machine Learning Research 12:2825–2830
- 195 [5] Nadeau C, Bengio Y (1999) Inference for the generalization error. Advances in
196 Neural Information Processing Systems 12
- 197 [6] Gao C, Shinkareva SV, Desai RH (2022) Scope: The South Carolina psycholin-
198 guistic metabase. Behavior Research Methods pp 1–32. URL [https://doi.org/10.](https://doi.org/10.3758/s13428-022-01934-0)
199 [3758/s13428-022-01934-0](https://doi.org/10.3758/s13428-022-01934-0)
- 200 [7] Warriner AB, Kuperman V, Brysbaert M (2013) Norms of valence, arousal, and
201 dominance for 13,915 English lemmas. Behavior Research Methods 45:1191–1207.
202 URL <https://doi.org/10.3758/s13428-012-0314-x>
- 203 [8] Mohammad S, Turney P (2010) Emotions evoked by common words and
204 phrases: Using mechanical turk to create an emotion lexicon. In: Proceedings
205 of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis
206 and Generation of Emotion in Text, pp 26–34, URL [https://doi.org/10.1111/j.](https://doi.org/10.1111/j.1467-8640.2012.00460.x)
207 [1467-8640.2012.00460.x](https://doi.org/10.1111/j.1467-8640.2012.00460.x)
- 208 [9] Scott GG, Keitel A, Becirspahic M, et al (2019) The Glasgow norms: Ratings of
209 5,500 words on nine scales. Behavior Research Methods 51(3):1258–1270. URL

210 <https://doi.org/10.3758/s13428-018-1099-3>

211 [10] Brysbaert M, Warriner AB, Kuperman V (2014) Concreteness ratings for 40
212 thousand generally known english word lemmas. Behavior Research Methods
213 46:904–911. URL <https://doi.org/10.3758/s13428-013-0403-5>

214 [11] Brysbaert M (2017) Age of acquisition ratings score better on criterion validity
215 than frequency trajectory or ratings “corrected” for frequency. Quarterly Jour-
216 nal of Experimental Psychology 70(7):1129–1139. URL [https://doi.org/10.1080/](https://doi.org/10.1080/17470218.2016.1172097)
217 [17470218.2016.1172097](https://doi.org/10.1080/17470218.2016.1172097)

218 [12] Brysbaert M, New B (2009) Moving beyond Kučera and Francis: A critical eval-
219 uation of current word frequency norms and the introduction of a new and
220 improved word frequency measure for American English. Behavior Research
221 Methods 41(4):977–990. URL <https://doi.org/10.3758/brm.41.4.977>

222 **4 Figure legends**

223 **4.1 Figure 1**

224 Investigating the impact of psychometric item ordering on inter-item correlations. **A.**
225 Order 1 inter-item correlations. **B.** Order 2 inter-item correlations. **C.** Order 1 minus
226 order 2 inter-item correlations.

227 **4.2 Figure 2**

228 Investigating the impact of psychometric item ordering on performance. *Psychometric*
229 *1* is obtained from participants that received the following order 1 (as listed in text).
230 *Psychometric 2* participants received the reverse order. *Psychometric* is an aggregate
231 of orders 1 and 2 (as used in the main analysis), and *Psychometric 1 & Psychometric*

Table 1

Norm	Category	Description	Source
Valence	Affect	The pleasantness of a stimulus on a 1 (happy) to 9 (unhappy) scale.	[7]
Arousal	Affect	The intensity of emotion provoked by a stimulus on a scale of 1 (calm) to 9 (aroused) scale.	[7]
Dominance	Affect	The degree of control exerted by a stimulus on a scale of 1 (controlled) to 9 (in control) scale.	[7]
Emotional Association Anticipation	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for anticipation.	[8]
Emotional Association Fear	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for fear.	[8]
Emotional Association Anger	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for anger.	[8]
Emotional Association Disgust	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for disgust.	[8]
Emotional Association Joy	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for joy.	[8]
Emotional Association Trust	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for trust.	[8]
Emotional Association Surprise	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for surprise.	[8]
Emotional Association Sadness	Affect	Word-emotion association built by manual annotation using Best-Worst Scaling method, with 0 (not associated) and 1 (associated) ratings for sadness.	[8]
Imageability	Concreteness	The degree of effort involved in generating a mental image of the concept on a 1 (unimaginable) to 7 (imageable) scale.	[9]
Concreteness	Concreteness	The degree to which the concept can be experienced directly through the senses from a 1 (abstract) to 5 (concrete) scale.	[10]
Familiarity	Frequency	A word’s subjective familiarity on a 1 (unfamiliar) to 7 (familiar) scale.	[9]
Age of Acquisition	Frequency	The age at which people acquired the word, in which a three-choice test was administered to participants in grades 4 to 16 (college) (Living Word Vocabulary Test).	[11]
Frequency	Frequency	Log10 version of frequency norms based on the SUBTLEXus corpus.	[12]

232 \mathcal{Q} is the concatenation of both orderings. Error bars are adjusted 95% confidence
 233 intervals [3].

234 **4.3 Figure 3**

235 Evaluating how test performance varies with training set size for 3 data (sub-)sets: (i)
236 Basel Risk Norms (All), which refers to our full data set of 1,004 risks, (ii) Basel Risk
237 Norms (Bhatia Set), referring to our data limited to the same 306 risks as used in
238 [2], and [2] (iii). Test sets are composed of all remaining risks in the data. Train–test
239 splits were sampled randomly (i.e., bootstrapped), with 10 repetitions per model per
240 training set size. Error bars are 95% confidence intervals.

241 **4.4 Figure 4**

242 Pairwise differences between elastic net and gradient boosting using 10x10-fold cross-
243 validation. Cross-validation via [2]’s risk norms (306 risks) are colored cyan and points
244 obtained using the Basel Risk Norms (1004 risks) are colored purple. Error bars are
245 adjusted 95% confidence intervals [3].

246 **4.5 Figure 5**

247 Pairwise differences between standardized and unstandardized models using 10x10-
248 fold cross-validation and elastic net regression. Cross-validation via [2]’s risk norms
249 (306 risks) are colored cyan and points obtained using the Basel Risk Norms (1004
250 risks) are colored purple. Error bars are adjusted 95% confidence intervals [3].

251 **4.6 Figure 6**

252 Heatmap illustrating the differences in 10x10-fold cross-validation R-squared between
253 all pairwise model combinations using elastic net regression (y-axis models minus x-
254 axis models). White squares reflect mean differences that do not significantly differ
255 from zero. The top panel shows the results for the data of [2] and the bottom panel
256 the results for the Basel Risk Norms.

257 **5 Table Legends**

258 **5.1 Table 1**

259 Word norms and their sources.

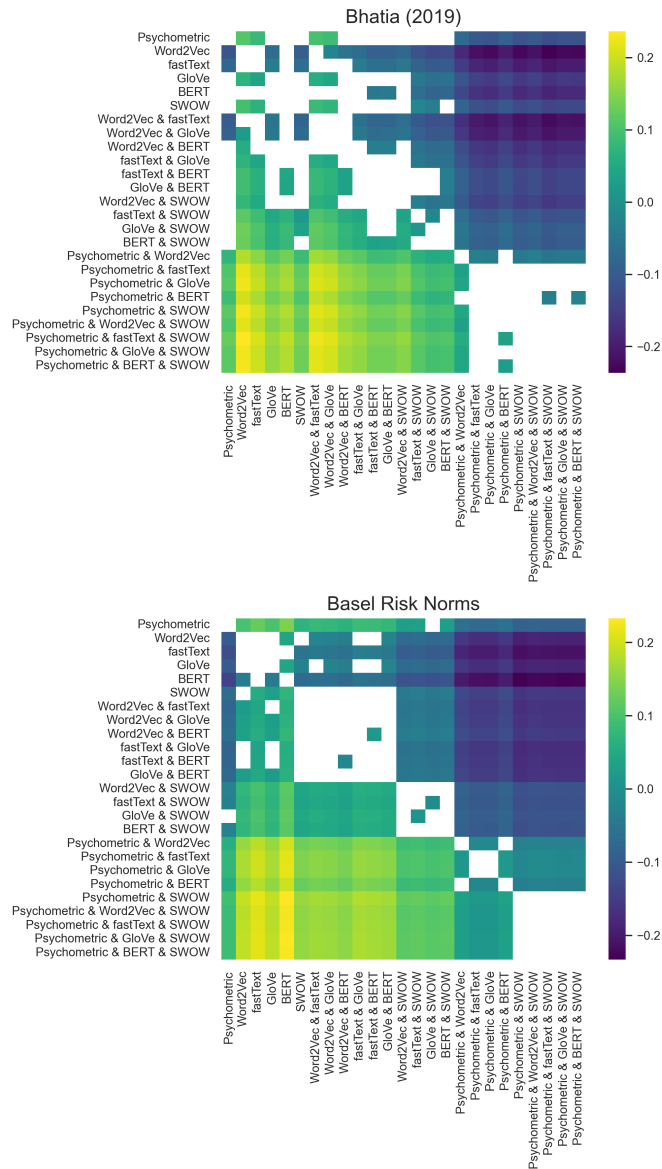


Fig. 6 Heatmap illustrating the differences in 10x10-fold cross-validation R-squared between all pairwise model combinations using elastic net regression (y-axis models minus x-axis models). White squares reflect mean differences that do not significantly differ from zero. The top panel shows the results for the data of [2] and the bottom panel the results for the Basel Risk Norms.