# Perceptual Dimensions of Wood Materials

Jiří Filip, A

Jiří Lukavský, B

Filip Děchtěrenko, B

Filipp Schmidt, C

Roland W. Fleming, C


A. The Czech Academy of Science, Institute of Information Theory and Automation
   Pod vodárenskou věží 4, 18200 Praha 8
B. The Czech Academy of Science, Institute of Psychology
   Hybernská 8, 11000 Praha 1
C. 1. Experimental Psychology, Justus Liebig University of Giessen Germany, 2. Centre for Mind, Brain and Behaviour, Universities of Marburg and Giessen
   Otto-Behaghel-Str 10, 35394 Giessen, Germany

## Abstract

Materials exhibit an extraordinary range of visual appearances. Characterising and quantifying appearance is important not only for basic research on perceptual mechanisms, but also for computer graphics and a wide range of industrial applications. While methods exist for capturing and representing the optical properties of materials and how they vary across surfaces (Haindl & Filip., 2013), the representations are typically very high-dimensional, and how these representations relate to subjective perceptual impressions of material appearance remains poorly understood. Here, we used a data-driven approach to characterising the perceived appearance characteristics of 30 samples of wood veneer using a 'visual fingerprint' that describes each sample as a multidimensional feature vector, with each dimension capturing a different aspect of the appearance. Fifty-six crowd-sourced participants viewed triplets of movies depicting different wood samples as the sample rotated.  Their task was to report which of the two match samples was subjectively most similar to the test sample. In another online experiment 45 participants rated ten wood-related appearance characteristics for each of the samples. The results reveal a consistent embedding of the samples across both experiments and a set of 9 perceptual dimensions capturing aspects including the roughness, directionality and spatial scale of the surface patterns. We also showed that a weighted linear combination of eleven image statistics, inspired by the rating characteristics, predicts perceptual dimensions well.

# Introduction

38

39 The visual appearance of materials results from a wide range of physical phenomena including the
40 surface's spectral and angular reflectance characteristics, subsurface light scattering, and spatial
41 variations in pigmentation and surface relief. How the visual system estimates such characteristics
42 remains poorly understood  (Anderson, 2011, Bracci & Op de Beeck, 2023) and it also remains unclear
43 which perceptual dimensions the visual system uses to describe and compare different materials
44 (Fleming, 2017).

45 Capturing a comprehensive representation of a surface's physical appearance requires observing it
46 under a sufficient range of illumination and viewing geometries. Complex photorealistic appearances
47 can be approximated by advanced image-based representations used in computer graphics such as the
48 spatially varying bidirectional reflectance distribution function (SVBRDF; Nicodemus & Richmond & Hsia
49 & Ginsburg & Limperis, 1977) or bidirectional texture function (BTF; Dana & van Ginneken & Nayar &
50 Koenderink, 1999).  However, these representations are extremely high-dimensional and there is no
51 straightforward mapping between such representations and subjective visual appearance
52 characteristics. Somehow the visual system summarises the overall 'look' of complex, spatially-varying
53 appearances to compare and contrast different materials. Everyday experience suggests that observers
54 do not need to view a material from all possible view- and lighting-directions in order to obtain a distinct
55 impression of its appearance.  Yet, although the perceptual representation of materials is surely lower-
56 dimensional than a complete physical description of the surface, there are nevertheless many potential
57 dimensions that the visual system might draw on to describe materials (e.g., overall albedo, relief,
58 glossiness, contrast of surface patterns).

59 We still do not understand much about such dimensions and how they contribute to observers'
60 judgments of appearance. Which characteristics do observers use to compare different materials?  Is
61 there a 'ranking' of characteristics, such that some aspects of appearance dominate comparisons
62 between materials, while others play a secondary role? How specific are certain characteristics to
63 particular classes of materials? Previous work on material perception has often focussed on highly
64 constrained sets of stimuli varying in one or a small number of physical properties (Ferwerda, Pellacini,
65 & Greenberg, 2001; Fleming, Dror, & Adelson, 2003; Fleming, Bülthoff, 2005; Motoyoshi, Nishida,
66 Sharan, & Adelson, 2007; Wendt, Faul, & Mausfeld, 2008; Wendt, Faul, Ekroll, & Mausfeld, 2010;
67 Fleming, Jäkel, & Maloney, 2011; Marlow, Kim, & Anderson, 2012; Paulun, Schmidt, van Assen, &
68 Fleming, 2017; Van Assen, Barla, & Fleming, 2018). Other studies have investigated appearance
69 judgments and categorization based on photographs (e.g., Bell, Upchurch, Snavely, & Bala, 2015;
70 Fleming, Wiebel, & Gegenfurtner 2013; Sharan, Rosenholtz, & Adelson, 2009; Sharan, Liu, Rosenholtz, &
71 Adelson, 2013; Sharan, Rosenholtz, & Adelson, 2014; Wiebel, Valsecchi, & Gegenfurtner, 2013).
72 However, in most cases, it is the experimenters that define which characteristics are judged by
73 participants.

74 Here we combined this tradition with a more data-driven approach in order to identify  dimensions
75 underlying appearance judgments for a set of thirty samples of planar wood veneer with distinctive

76    surface patterns and textures. Wood is a challenging material to characterise due to its complex and
77    varied appearance.  It is associated with decorative attributes and is widely used for furniture and
78    interior design. Its structure consists of elongated cells, which are radially oriented rays and longitudinal
79    cells or vessels forming growth rings (Lewin & Goldstein, 1991). Hardwoods tend to have a tighter grain
80    pattern compared to softwoods, resulting in various levels of texture, colour, smoothness, grain density
81    and straightness. All these aspects are impacted by sawing direction and the sample location in the tree
82    trunk. The final visual structure is given by an intersection of a sawing plane with three-dimensional
83    wood structure. Wood has high natural variability in aesthetic characteristics  among different species
84    and surface treatments. Previous studies have shown that patterns of anisotropy, colour variations and
85    gloss are the major factors influencing the visual (Nakamura, Masuda, & Shinohara, 1999; Wan, Li,
86    Zhang, Song, & Ke, 2021), multimodal (Fujisaki, Tokita, & Kariya, 2015) aesthetic appeal of wood with
87    impacts on people's preferences (Manuel, Leonhart, Broman, & Becker, 2015), and emotions related to
88    wooden surfaces (Nordvik, Schütte, & Broman, 2009). To the best of our knowledge, all previous studies
89    of wood appearance relied on static stimuli to derive subjective ratings of predefined attributes or their
90    relationship to physical attributes of wood surfaces. This ignores how variable the appearance of even a
91    single sample can be across changes in viewpoint relative to the surface and lighting. Our contribution
92    above this work is twofold.

93    First, our work uses dynamic (rotating) rather than static stimuli, showing the appearance of the wood
94    samples across variable lighting and viewing conditions. This allowed participants in our experiments to
95    take into account the look of the surface both with and without specular reflections.
96
97    Second, instead of relying solely on a possibly incomplete list of predefined visual attributes, we also
98    used similarity judgements to identify the core dimensions underlying judgments of wood. Similarity
99    judgements are an established method for characterising the multidimensional space underlying mental
100   representations, previously used to understand dimensions in object categories (Hebart, Zheng, Pereira,
101   & Baker, 2020), materials (Schmidt, Hebart, & Fleming, 2022) or scenes (Josephs, Hebart, & Konkle,
102   2023). In contrast to the previous studies, we search for dimensions underlying similarity judgments
103   within a single category.
104   Specifically, we sought to derive a relatively small number of perceptual dimensions that capture
105   judgments of similarity between movies of the samples.  In order to do this, we first crowd-sourced
106   1218 perceptual similarity judgments from 56 participants.  We then applied an analysis method based
107   on sparse, non-negative matrix factorization (Variational Interpretable Concept Embeddings;
108   Muttenthaler,  Zheng, McClure, Vandermeulen, Hebart, & Pereira, 2022) to infer a set of dimensions
109   that can predict the similarity judgments.  We show that even with a small dataset of thirty samples the
110   method was able to derive visual dimensions that predict the similarity judgments. Specifically, our
111   model identified nine dimensions that together could explain over 75% of the variance in the similarity
112   judgments.  Eventually, we showed that standard image statistics obtained from stimuli videos can
113   predict similarity dimensions well.

114   In addition to the similarity judgments, we also asked a set of 45 participants to judge ten experimenter-
115   defined appearance characteristics for each of the samples (Brightness, Glossiness, Colourfulness,
116   Directionality, Complexity, Contrast, Roughness, Patchiness/regularity, Line elongation, Spatial scale).

117    The purpose of this was twofold.  First, we sought to use the values of these interpretable rating scales
118    to  facilitate interpretation of the dimensions derived from the similarity judgments.  Second, we sought
119    to cross-validate the embedding of the samples within the 9D space.  We reasoned that if different
120    samples are represented in a multidimensional perceptual similarity space—with similar samples close
121    to one another and dissimilar ones further apart—then it should be possible to probe this space through
122    multiple complementary methods (i.e., similarity judgments and subjective feature ratings).  We find
123    that the two approaches do indeed lead to similar embeddings of the stimuli, suggesting that they both
124    tap into a common representation within the visual system.

# Experiment 1

126    In the first experiment we collected sparse similarity judgements and used machine learning to infer the
127    full pairwise similarity matrix and to test the embedding of samples in the latent space  of wood
128    appearance.

## Methods

### *Participants*

131    Fifty-six participants were recruited using the online crowdsourcing platform Prolific (mean age = 40.5,
132    SD = 16.6, 35 males). All participants reported normal or corrected-to-normal vision and no colour vision
133    impairments. On average, the experiment took 14.0 minutes (SD = 4.9). The participants were
134    reimbursed with 2.1 GBP. All studies within this paper were approved by the Ethics Board of the
135    Institute of Psychology, Academy of Sciences of the Czech Republic (PSU-308/Brno/2022).
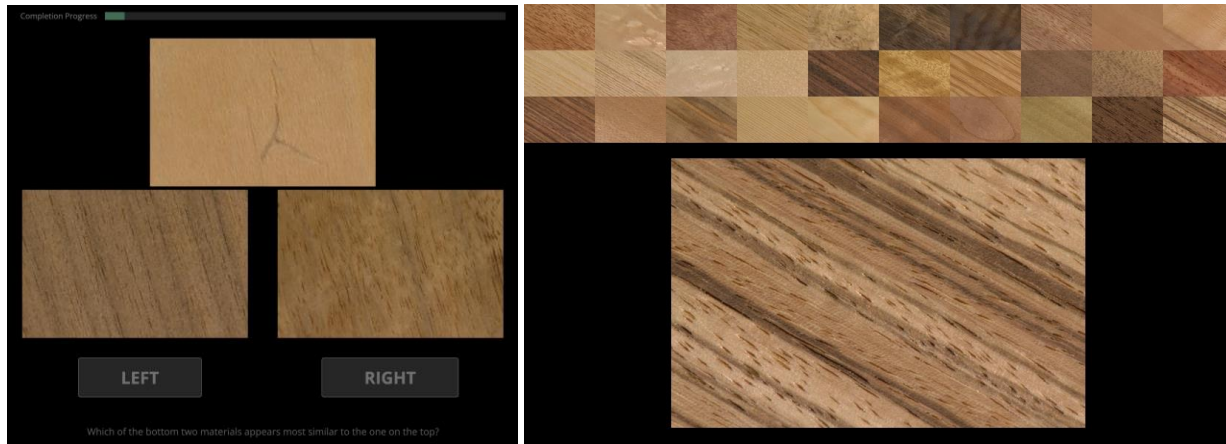
### *Apparatus and stimuli*

137    We used 30 flat standard wood veneer samples that are used for furniture manufacturing (wood species
138    are listed in Tab.1). We captured  video sequences of slow rotations of the samples. Fig. 1 shows the
139    initial (left) and final (right) frame of each video sequence, capturing specular and non-specular
140    view/light geometries. Video samples and additional materials are available at https://osf.io/tz245.
141

**Fig. 1.** All 30 samples of wood veneer for (left) specular, (right) non-specular (90 degree rotated) view/light geometries.

All images in the video sequences were 42 x 42 mm areas of the samples, captured by the UTIA goniometer (Filip, Vávra, Haindl, Žid, Krupička, & Havran, 2013). In accordance with industry standards in material observation (McCamy, 1996), we fixed the polar angle of camera and light to 45 degrees and only varied azimuthal angles to allow for faster measurements. Each sequence starts with a difference of 90 degrees between the azimuthal angles of light and camera and includes a movement of the camera by 90 degrees (arriving at a difference of 180 degrees between azimuthal angles), resulting in specular and non-specular material behaviour as shown in Fig. 1. Each 4-second sequence consists of 60 image frames, repeated in reverse order to create a continuous loop of rotating material. See supplementary video [movie_samples_stimuli.avi].

To allow for smooth presentation in the experiment, the image frames of all samples were cropped and downsampled to 400 x 260 pixels, and combined into single-trial frames with three samples on a black background at qHD (quarter high definition, 960 x 540 pixels) as shown in Fig. 2(a). Each sequence was started at a random time of the continuous loop to prevent participants from responding to initial frames of the video sequence.

(a) Experiment 1 (2AFC match-to-sample)          (b) Experiment 2 (Ratings)

**Fig. 2.** Example of stimuli frames of (a) the similarity judgement experiment, where participants responded to: "Which of the bottom two materials appears most similar to the one on the top?", and of (b) the rating experiment, where participants rated individual samples according to different visual attributes.

Because data was collected online, we did not control for viewing distance (viewing angles) or monitor settings. However, a post-hoc analysis of monitor settings showed a minimal screen resolution of 980 x 577 pixels which allows for a full-resolution presentation of our stimuli.


## *Experimental procedure*

Experiment 1 consisted of 93 trials. In each trial, participants judged the similarity of three presented samples as shown in Fig. 2(a), by deciding which of two match stimuli (at the bottom of the screen) was more similar to the test stimulus (at the top of the screen; 2AFC match-to-sample design). Because we do study similarity within a single material category (wood), we hypothesised a relatively low number of 3 to 5 meaningful perceptual appearance dimensions. In line with the recommendations in (Haghiri, Rubisch, Geirhos, Wichmann, & von Luxburg, 2019) (30 samples and 3-5 dimensions: 900-1500 trials), we tested 1218 triplets, accounting for 10% of the full similarity matrix.

Across all triplets, each sample was presented as a test stimulus in 160-164 trials and as a match stimulus in 304-344 trials. Each triplet was judged four times (i.e, by four different participants). Two out of four repetitions swapped the left and right match stimuli to control for a potential response bias. Each participant was presented with one of 28 unique trial sets or its copy with swapped match stimuli.

Data were collected online using a custom script in the jsPsych framework (De Leeuw, 2015). After reading the instructions, participants completed three practice trials and 90 experimental trials (87 trials plus 3 catch trials). They initiated each trial by clicking the "Start" button after which a video with the three samples started looping (Fig. 2(a)). Participants responded to the instruction below the video ("*Which of the bottom two materials appears most similar to the one on the top?*") by clicking on the "LEFT" or "RIGHT" button at the bottom. The response stopped the loop and initiated the next trial, with

189  a progress bar at the top showing the number of remaining trials. Catch trials were presented at fixed
190  positions (40th, 65th, and 84th trials) and featured the same sample presented twice, as standard and
191  match stimulus, yielding a ground truth correct response.

192  *Data analysis*

193  All data is available from the following public repository: [link provided upon acceptance]. We next
194  sought to identify a set of perceptual dimensions—with values for every sample—that could account for
195  the observed pattern of similarity responses.  To do this, we analysed the responses using Variational
196  Interpretable Concept Embeddings (VICE; Muttenthaler,  Zheng, McClure, Vandermeulen, Hebart, &
197  Pereira, 2022). This algorithm takes as input the sparse (i.e., incomplete) similarity matrix obtained in
198  the similarity rating experiment and estimates the full pairwise similarity matrix.  In the process, it
199  iteratively estimates a set of underlying dimensions that could account for the observed responses. As
200  our similarity judgement study comprises 2AFC task, we applied target matching instead of odd-one-out
201  procedure.

202  Several of the VICE algorithm's hyperparameters can affect its results, including the number of
203  dimensions. To validate the performance of the model, we created random splits of our participants'
204  similarity judgements into training (90% of responses) and test sets (10% of responses).  Then, we
205  performed a limited grid search for selected hyperparameters of the model: learning rate [0.0005,
206  0.001, 0.002], mixture of distributions in the spike-and-slab prior [Gaussians, Laplace], spike (a prior of
207  probability at zero values) [0.125, 0.25, 0.75], slab (a prior of probability for the non-zero values) [0.2,
208  0.5, 1.0], and probability of relative weighting of the distributions [0.4, 0.5, 0.6]. The training typically
209  converged within 200 epochs, and typically resulted in between 8 and 11 dimensions (min. 4, max. 14
210  dimensions). Details of the model selection and training process are reported in Section 1 of the
211  supplementary material.

212  **Results**

213  *Consistency of similarity judgement responses*

214  Our results show that participants were highly consistent in their similarity judgments. When analysing
215  inter-individual consistency based on the four repetitions of each triplet, in 569 triplets (47%) all four
216  responses were the same, in 439 triplets (36%) three responses were the same, and in 210 triplets (17%)
217  responses were on par. This suggests that in the majority of trials (87%) subjects were consistent, only in
218  the remaining 17% they were at chance. Also, when comparing sequences with their copies with
219  swapped match stimuli, only in 61 trials (5%) swapping resulted in a different response.

220  *Deriving perceptual dimensions from similarity ratings*

221  Based on the parameter grid search (see Section 1 of the supplementary material), we picked the best
222  performing model (9 dimensions; accuracy on the training set = 0.760; accuracy on the test set = 0.769).
223  Importantly, even though the number of dimensions varied between different resulting models from the

224    parameter grid search, the meaning of those dimensions was highly preserved. Specifically, the
225    embeddings obtained from the first five best VICE models (with 4-9 dimensions) were highly similar
226    (mean correlation between similarity matrices of the 4 next best models to that of the best model was
227    R=0.939). Thus, in the following we analyse the best performing VICE model under the justified
228    assumption that it is representative of a family of models with similar embedding.

229    The resulting embedding as shown in Fig. 3(a) is quite sparse, with on average only 6 values > 20%
230    percentile in each similarity dimension. Fig. 3(b) shows the sum of loadings for individual dimensions
231    and suggests that the first 5 dimensions have higher impact than the remaining 4. Fig. 3(c) compares
232    how well the similarity responses from participants can be approximated by the values estimated from
233    the VICE model.  Chance performance in the 2AFC match-to-sample task (red) is 50%, with the inter-
234    participant noise ceiling (grey) at 82%. The noise ceiling is computed as the average consistency across
235    the four repetitions of each triplet, and represents the best possible prediction any model could achieve
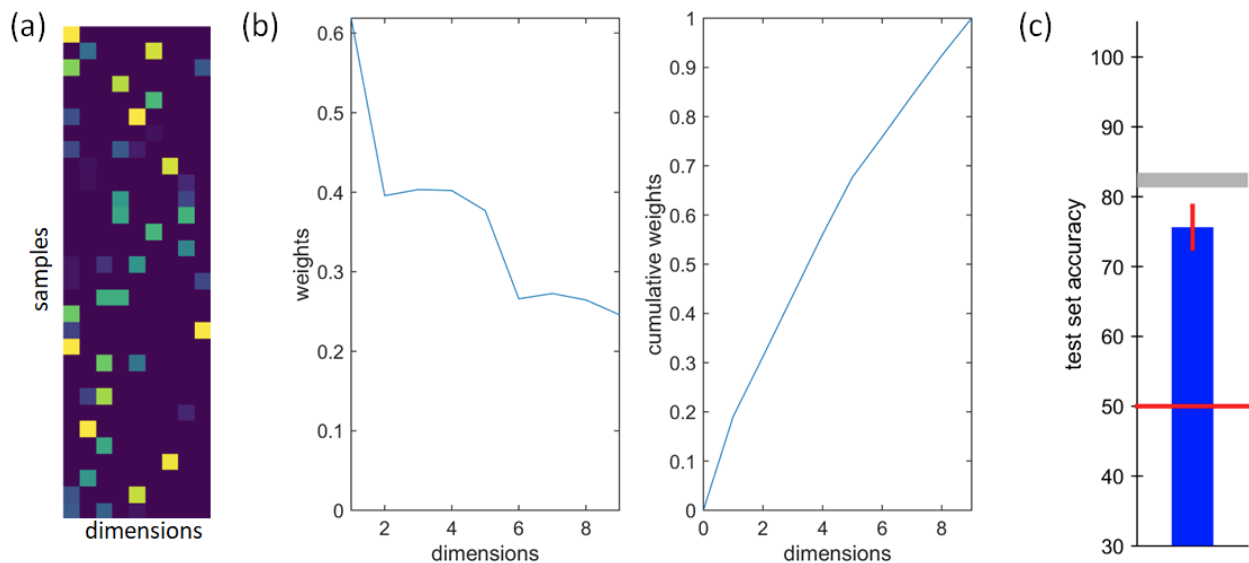236    for our dataset, given the variation in the data.

237



**Fig. 3** Details on the 9 dimensions of the best VICE model: (a) estimated embedding, (b) dimensions
loadings, and (c) average accuracy on test set (blue) with 95% confidence interval error-bar (red), noise
ceiling (grey), and chance level (red).

243    Fig. 4 shows samples rank-ordered by their embedding values in each of the 9 dimensions (highest
244    values to the left). Each video sample is represented by its two most distinct frames, i.e. non-specular
245    and specular reflection (refer to the supplementary to see the dynamic behaviour of the actual video
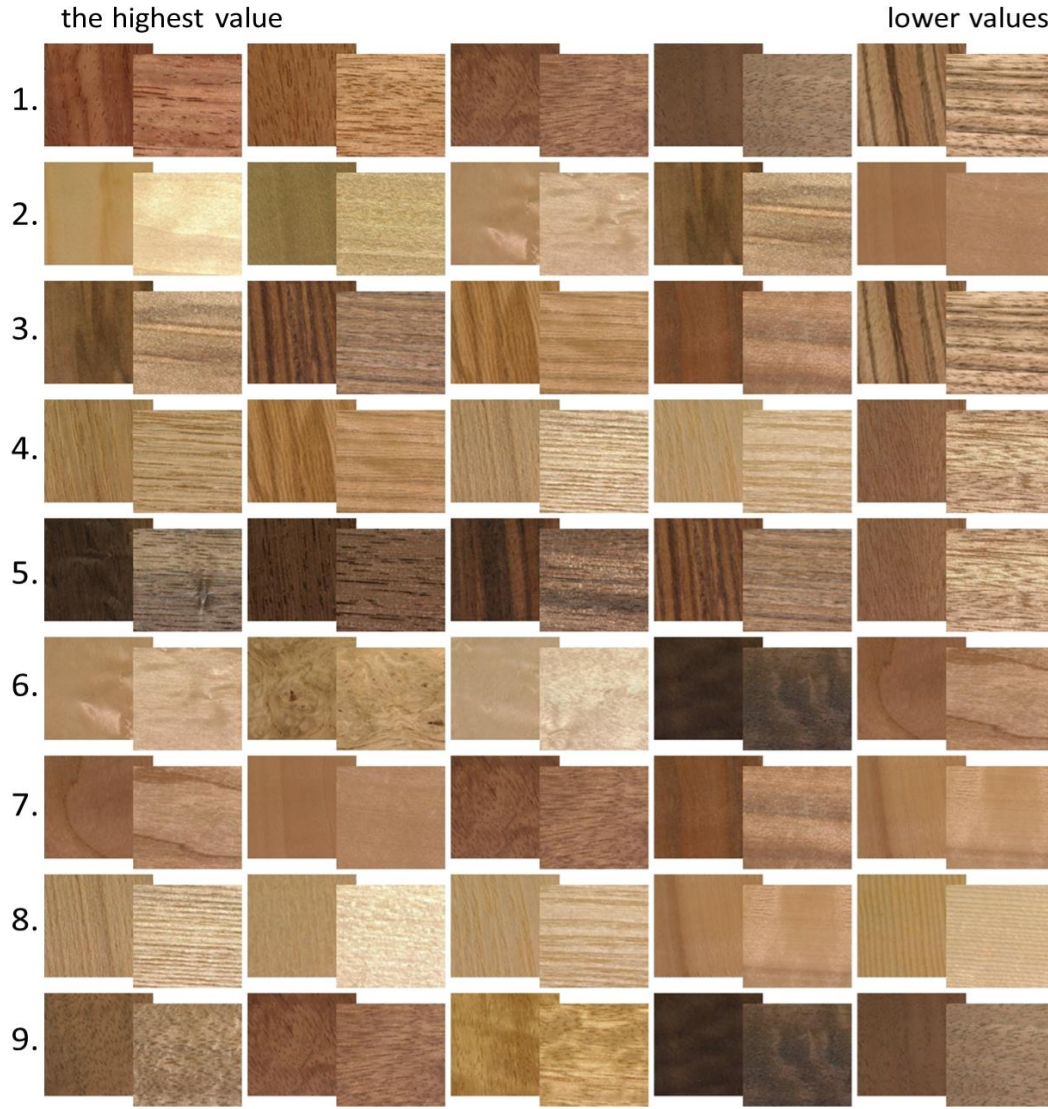246    samples).
247

**Fig. 4** Five samples for each dimension rank-ordered based on embedding values. Each video sample is represented by both the most non-specular and most specular condition. See left side of the supplementary video [movie_similarity_vs_rating.avi].

The full pairwise similarity matrix of the wood samples that we obtained from the estimated embedding is shown in Fig. 5. We used hierarchical clustering (based on weighted average Euclidean distance) to cluster similar samples together, showing that samples had approximately three main visual modes, which might be visually interpreted as rough/contrast (M1), spatial frequency (M2), and directional (M3). These modes are present also in individual similarity dimensions in Fig. 4, where M1 is represented by dimensions 1, 5 and 9, M2 by 2 and 6, and M3 by 4, 8 and 3. Note that similar modes were also found using Louvain community detection method (Blondel, Guillaume, Lambiotte, & Lefebvre, 2008) as reported in Section 2 of the supplementary material.
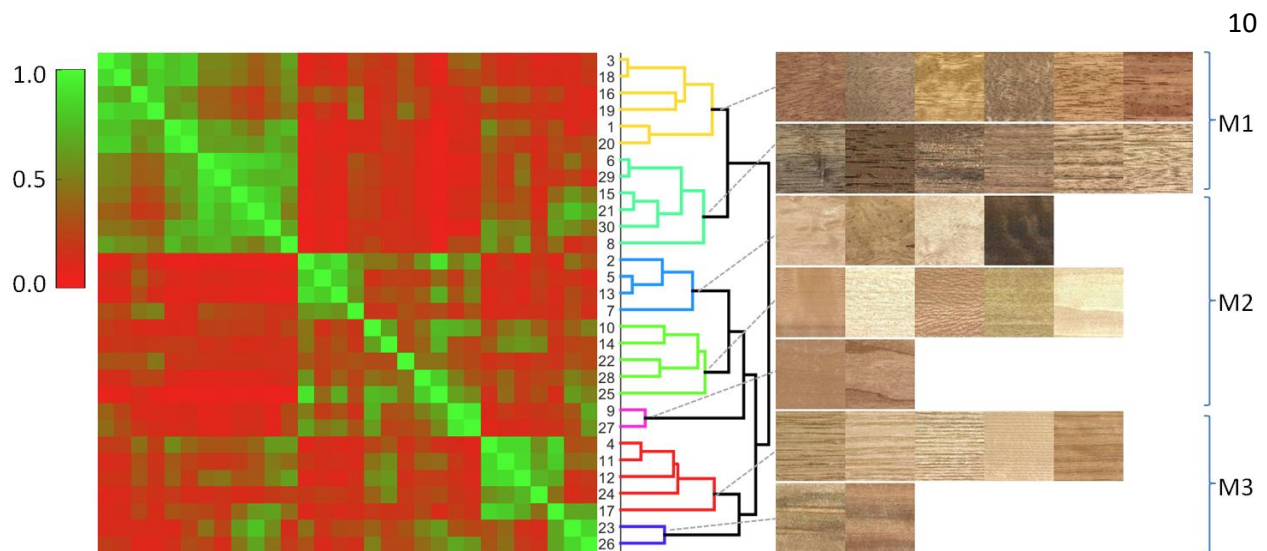
**Fig. 5** Estimated pairwise similarity matrix with samples ordered based on hierarchical clustering, and the depiction of the corresponding samples in the individual clusters.

## *Discussion*

The analysis of participants' similarity judgements using the VICE model provided us with 9 visual appearance dimensions of wood. However, even though visualising the embedding by ranking samples within each dimension may provide some intuition about the meaning of the dimensions, it is not clear whether these intuitions are the best description of the respective dimensions. For this reason we performed a second comparative experiment relying on standard attributes rating on a Likert scale.

# Experiment 2

The main goal of the second experiment was to obtain perceptual judgements for all wood samples for a set of visual appearance attributes widely used in the field of material perception. By being able to describe our samples in terms of these specific perceptual attributes, we aimed to provide a more valid interpretation of the similarity dimensions from the first experiment–and a corresponding understanding of the main visual cues that naive observers use to describe and discriminate between types of wood.

## **Methods**

## *Participants*

Forty five volunteer observers participated in the online experiment (age data were not collected). All participants reported normal or corrected-to-normal vision and no colour vision impairments. On average, the experiment took 22.0 minutes (SD = 17.6).

### *Apparatus and stimuli*

The stimuli used in Experiment 2 were the same as in Experiment 1.

### *Procedure*

Participants were presented with 30 trials, each showing one of the sample videos from Experiment 1. The resolution of each stimuli image was 920 x 600 pixels. To make the task easier for participants, all other materials were simultaneously presented for comparison at a smaller scale at the top of the screen as shown in Fig. 2(b). Participants rated each material on ten visual appearance attributes (brightness, glossiness, colourfulness, directionality, complexity, contrast, roughness, patchiness/regularity, line elongation, and spatial scale), using a visual analog scale. The attributes were selected based on a review of previous research (Tamura, Mori, & Yamawaki, 1978; Rao & Lohse, 1996; Fleming, Wiebel, & Gegenfurtner 2013; Tanaka & Horiuchi, 2015, Nordvik, Schütte, & Broman, 2009) and salient differences between samples identified by the experimenters. For the participants, the meaning of each visual attribute was explained with a short sentence (e.g., brightness: "*How bright is the material in comparison with the others?*"). Also, the end points of each scale were labelled (e.g., brightness: "dark" and "bright"). A full description of each visual attribute and the corresponding endpoint labels is provided in Section 4 of the Supplementary Material.

All attribute scales were on the screen simultaneously, and at the start of each trial all sliders were set to the centre of each scale. Only after moving all sliders, participants could proceed to the next trial.

### *Data analysis*

Again, a post-hoc analysis of monitor settings showed a sufficient minimum screen resolution of 980 x 768 pixels. The inter-rater agreement was determined using intraclass correlation coefficient  (ICC; Koo & Li, 2016, with two-way random effects, based on mean rating and consistency). More detailed analysis of participants' responses is provided in Section 5 of the Supplementary Material.

### **Results**

The rating responses for each attribute formed unimodal distributions with mean values close to the central point (45.8 to 59.5) and similar SD values (21.7 to 29.6). The ICC indicated excellent reliability (ICC > 0.898) for all attributes but *spatial scale* where ICC = 0.659 indicated only moderate reliability.

**Fig. 6** Five samples for each rating dimension rank-ordered based on average rating responses. Each video sample is represented by both the most non-specular and most specular condition. See right side of supplementary video [movie_similarity_vs_rating.avi].

Samples with the highest rating responses for each rating dimension are shown in Fig. 6, with visually intuitive results in the majority of dimensions (again with the exception of *spatial scale*).

Note that these examples also suggest similarities between rating dimensions (i.e. overlap in samples for e.g. *colorfulness* and *contrast*). To measure these inter-class similarities, we computed Pearson correlations for mean rating values across all 30 samples. As shown in Fig. 7, we observe a high similarity between *colorfulness-contrast, directionality-line elongations* and *complexity-patchiness/regularity*. On the other hand, a high dissimilarity is observed for *brightness-colorfulness* and *brightness-contrast*. These similarities are also evident at the level of individual samples, as is shown in Fig. 10(a) which is showing similarity matrices for individual rating dimensions.
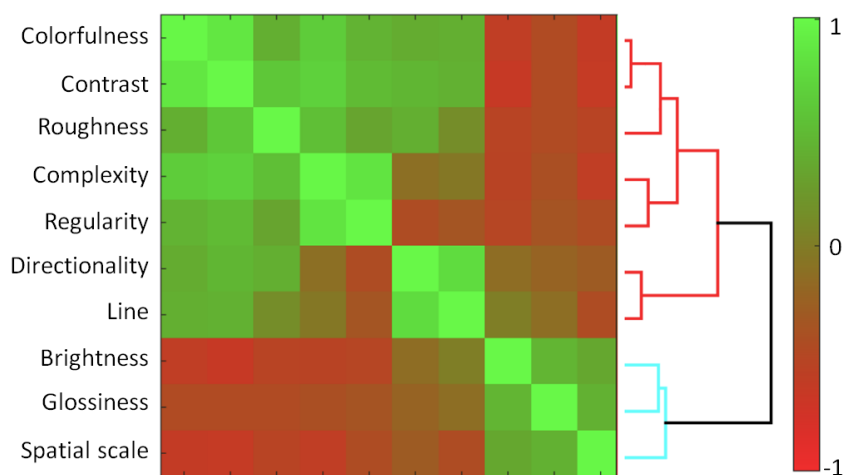
328



329

**Fig. 7** Inter-class similarity, computed as Pearson correlation across all samples, with the dendrogram showing the results of hierarchical clustering of attributes.

See supplementary video http://staff.utia.cas.cz/filip/tmp/movie_similarity_vs_rating.avi with material samples ranking as a function of dimensions loadings of VICE (left) and rating responses having the highest and the lowest values.

**Discussion**

Our rating experiment provided reliable and visually intuitive data on the selected visual appearance attributes, but also highlighted mutual dependencies between some of the attributes. This suggests that our samples can be described by less than 10 attributes, that is, the latent visual dimensionality of our samples is lower than 10.  In the next section, we compare the visual dimensions obtained from the similarity and rating experiments.

# Interpretation of similarity dimensions

As the meaning of the similarity dimensions discovered by the VICE model are not known, we used cross correlation and multilinear regressions between appearance ratings and similarity judgements as well as between their respective similarity matrices. This allowed us to assign meaning to the similarity dimensions by relating them to the meaningful appearance ratings.

**Cross-correlation of similarity and rating dimensions**

Across all appearance attributes, the correlation between similarity matrices from ratings and similarity judgements is relatively low (Pearson R=0.335; exclusion of matrix diagonal), with the highest correlations for *directionality* (R=0.448) and *contrast* (R=0.462). This confirmed our expectation that the similarity embedding cannot be explained using a single rating dimension.

353   For a direct correlation between all ratings and all similarity dimensions see Fig. 8(a). The highest

354   positive correlation was R=0.744 and the highest negative correlation was R=-0.813. Notably, similarity

355   dimensions 1, 3, 4 and 5 show similar patterns of correlation to rating attributes *colourfulness*,

356   *directionality*, *complexity* and *roughness*. On the other hand, similarity dimension 7 is not correlated

357   strongly with any rating attribute, which suggests that none of them can explain the visual appearance

358   captured by this particular dimension.  To test whether the similar pattern of correlations across

359   similarity dimensions follows from a strong dependency between individual rating attributes, we

360   computed PCA on our rating data. Fig. 8(b) shows that only four PCA components explain 91.1% of the

361   variance, suggesting that the effective number of main visual appearance dimensions for our set of

362   wood samples is about 5-10. We confirm this hypothesis by using a statistical approach to estimate the

363   number of dimensions based on triplet embedding accuracy of ordinal triplets embedding  (Künstle, von

364   Luxburg, & Wichmann, 2022) – which identifies 6 as the inherent dimensionality of our data (see details

365   on this analysis in Section 3 of the supplementary material). This is also supported by the steep drop of

366   similarity embedding factor loadings with more than five dimensions (Fig. 3(b)).
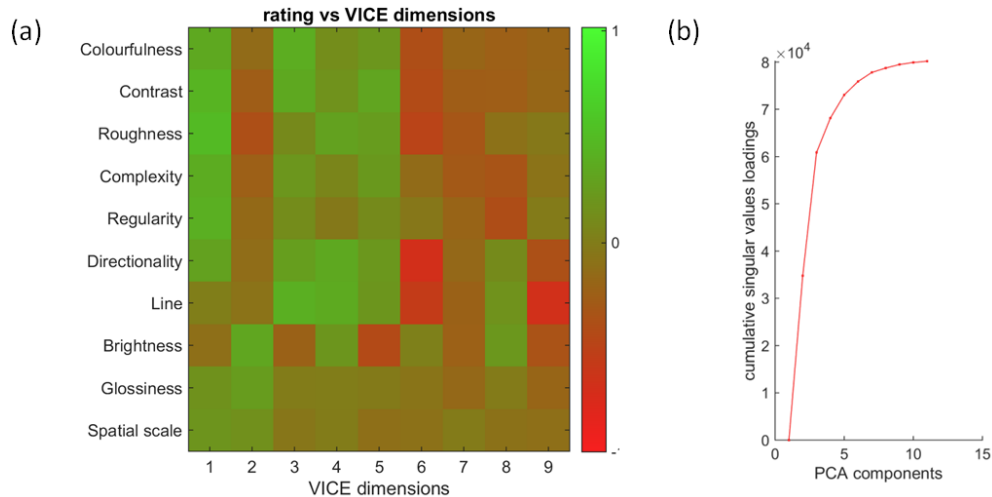
367



368
369   **Fig. 8** (a) Correlations of rating dimensions (rows) to similarity dimensions (columns), with negative

370   correlations in red and positive correlations in green (range [-1,1]). (b) Cumulative singular values

371   loadings of PCA computed on correlations across rating dimensions (a).

372

373   A more quantitative comparison between similarity dimensions and rating attributes is shown in Fig. 9.

374   For each similarity dimension, we ordered and scaled samples according to their dimension values. The

375   inset shows how well the variation in each similarity dimension is correlated with different rating

376   attributes. Here we observe similar patterns for dimensions 1, 3, 4, and 5 while dimension 7 is virtually

377   constant across rating attributes. $R^2$ scores in the legend demonstrate how well each similarity

378   dimension can be predicted by a linear regression of rating attributes.
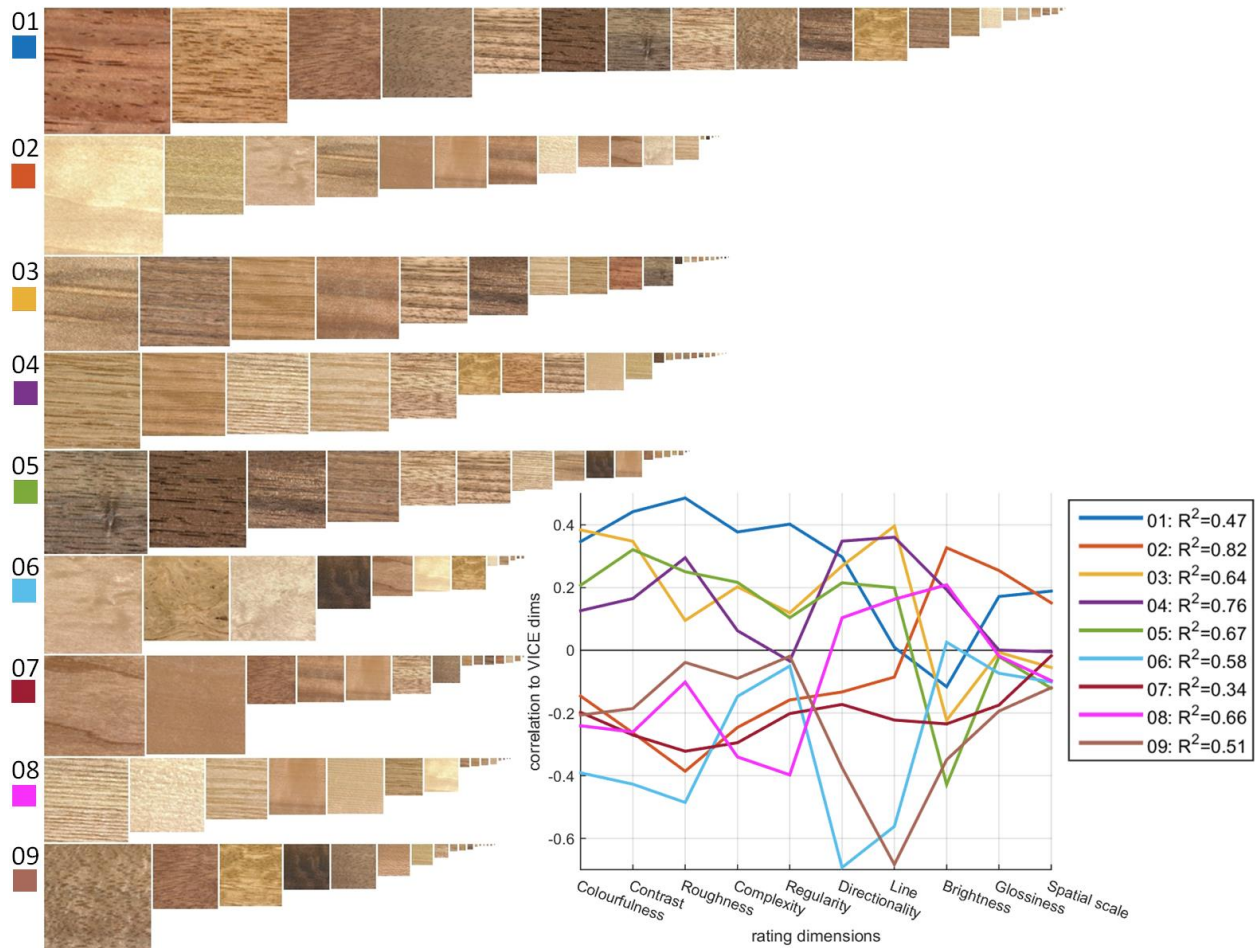
379

**Fig. 9** Sample rank-ordered by embedding values in VICE similarity dimensions. Inset: Correlations between similarity (VICE) and rating attributes, obtained from linear regressions ($R^2$ scores provide information on how well the linear regression using rating dimensions explained individual similarity dimensions. See the supplementary video [movie_similarity_scaled.avi].

To evaluate similarity of results obtained from both experiments, we computed the rating similarity matrix as Euclidean distance across all attributes. A direct correlation between similarity matrices obtained from similarity judgement and attributes rating (excluding diagonal elements) was R=0.627 ($R^2$ = 0.393). The matrices are shown in the first row of Fig. 10(a,b).

To assess the main visual dimensions for similarity judgements, we computed multidimensional scaling MDS (Carroll & Arabie, 1998) on the VICE similarity matrix. The MDS projection of samples onto the first two dimensions are shown in Fig. 10(d). In line with our visual interpretation of the three main visual modes in Fig. 5, the first MDS dimension can be interpreted as related to roughness, the second to directionality and the third to spatial frequency. For clarity, we also included these plots with the video samples as presented to observers. We compared MDS results over the similarity matrices and coordinates of all 30 samples for the first two MDS dimensions after Procrustes alignment are shown in

398    Fig. 10(e). For MDS of VICE similarity matrix into all 3 dimensions see a top part of the supplementary
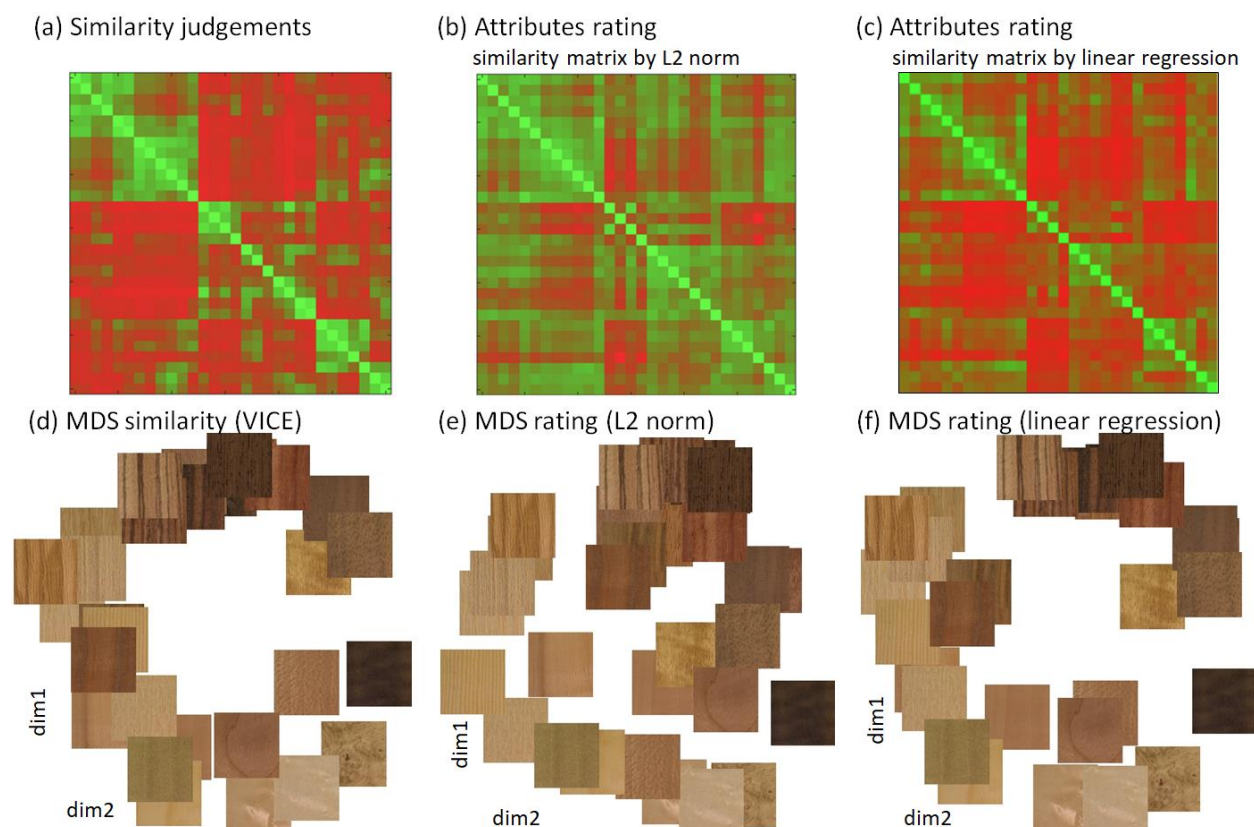
399    video [movie_MDS_simmat_linreg.avi].

400



**Fig. 10** A comparison of similarity matrices obtained by (a) similarity judgements and (b,c) ratings using
L2-norm and linear regression, respectively. The correlation between matrices is for (b) R=0.627 ($R^2$ =
0.393) and for (c) R=0.723 ($R^2$=0.523). Corresponding embeddings of samples in the first two MDS
dimensions for (d) similarity judgement and (e,f) ratings (after Procrustes alignment).

406

## Prediction of similarity matrix from rating attributes

408    Beyond simple correlations between individual ratings and similarity dimensions, we can test how well a

409    combination of rating attributes predict similarity judgements. To this end, we used multilinear

410    regression to predict the similarity judgement matrix by a linear combination of the rating attribute

411    similarity matrices shown in Fig. 11(a). The matrices' diagonals were kept to anchor scaling. The

412    regression model explains about 52% of the variance in similarity judgements (R=0.723, $R^2$=0.523), while

413    still preserving the major similarity modes as shown in Fig. 11(a). To evaluate the importance of

414    individual rating attributes for the reconstruction, we performed leave-one-out regressions and the

415    resulting drops in explained variance. Fig. 11(b) shows that the most important attributes are *brightness*,

416    *directionality*, and *roughness*. A comparison of the obtained multi-dimensional scaling over the similarity

417    matrices and coordinates of all 30 samples for the first two MDS dimensions after Procrustes alignment

418    are shown in Fig. 10(f). Also see a Section 6 of the supplementary material for samples alignment
419    according to MDS and video [movie_MDS_simmat_linreg.avi], comparing three MDS dimensions of
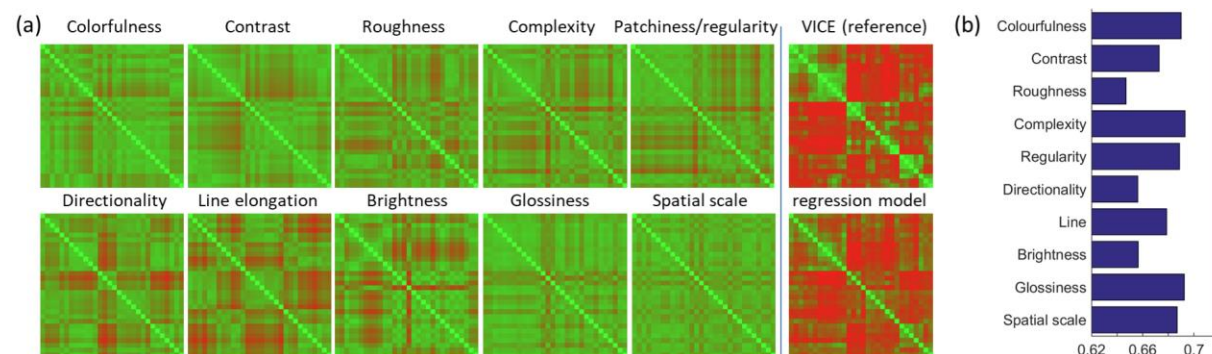420    similarity judgements similarity matrix (top) with its linear regression using rating attributes (bottom).
421
422



423
424    **Fig. 11** (a) Similarity matrices of individual rating attributes compared to the VICE similarity matrix and
425    the result of the linear combination of the 10 rating similarity matrices. (b) Results of leave-one-out
426    regression analyses showing the respective drops in correlation below the red baseline due to individual
427    attributes removal.


428    **Prediction of similarity dimensions from rating attributes**

429    Finally, we used linear regression to predict individual similarity dimensions by a linear combination of
430    the rating attributes. Across all dimensions, ratings can well explain similarity dimensions, with an
431    average of R=0.851 ($R^2$=0.731). $R^2$ scores of similarity dimensions represented by the regression model
432    are shown in Fig. 12(a). All dimensions except 7 and 9 can be well explained by a combination of rating
433    attributes. As reported previously, dimension 7 is not well predicted by any of the rating attributes. This
434    might be for two reasons: either none of our predefined attributes is not capturing the same visual
435    appearance as that dimension, or there is a general bias in our rating data  that is introduced by a
436    particular interpretation of the to-be-rated attributes. For instance *line elongation*, *patchiness/regularity*
437    or *spatial scale* might have different meanings at different frequency scales. For instance samples 22
438    and  30 (see Fig. 1) both share a fine detail structure and a distinct low-frequency stripy pattern. As a
439    result, observers might be confused as to whether these attributes should be evaluated on a fine or
440    coarse scale, resulting in overall ambiguous ratings. In Fig. 12(b), we are plotting normalised regression
441    values to visualise the contribution of rating attributes to each of the similarity dimensions. For example,
442    dimension 1 is strongly negatively related to colorfulness and roughness, but strongly positively related
443    to contrast, regularity and directionality; while dimension 5 represents materials that were judged low
444    on line elongation, brightness and glossiness (see samples rank ordering along dimensions in Fig. 9).
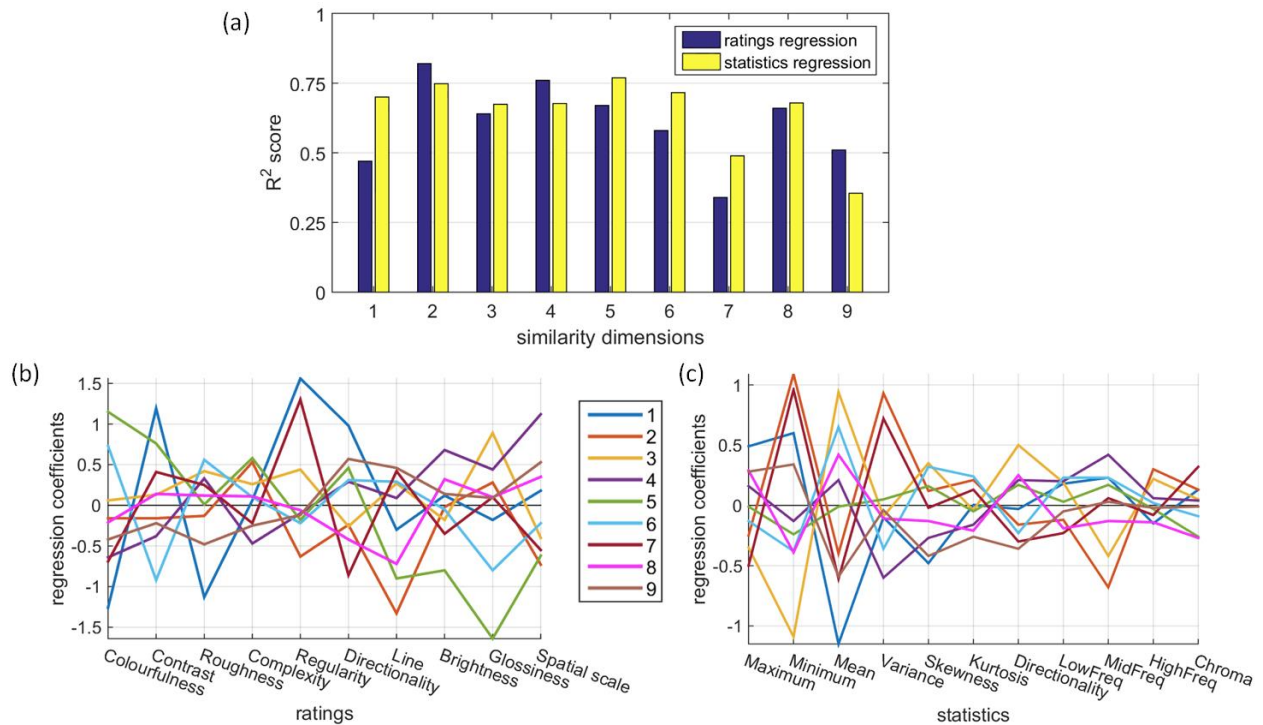445

**Fig. 12** (a) $R^2$ scores of similarity dimensions regression of rating dimensions (blue) and regression of computational image statistics, with corresponding normalised regression coefficients showing contribution of rating dimensions (b) and computational statistics (c) for reconstruction of each similarity dimension.

## Relationship to computational statistics

To relate similarity dimensions to computational statistics, we used standard image statistics used in texture synthesis related to human low-level perception of textures (Portilla & Simoncelli, 2000, Motoyoshi, Nishida, Sharan, & Adelson, 2007), namely *minimum, maximum, mean, variance, skewness,* and *kurtosis*. We supplied additional statistics evaluating image directionality (Maskey & Newman, 2021) and frequency content in three bands (*low, mid*, and *high* frequencies) computed from PSD of image converted to the Fourier domain. The final values of statistics were averaged across all frames of movie sequence. We used these statistics for linear regression of similarity ratings and $R^2$ scores of results are shown as yellow bars in Fig. 12(a). We observe similar values of $R^2$ scores to those obtained from rating regression and in general all similarity dimensions, except 7 and 9, can be represented reasonably well using our statistics. Mean $R^2$ score across all dimensions was 0.65 (R=0.73). Normalised regression coefficients of individual statistics are shown in Fig. 12(c). For example dimension 1 has the highest coefficients for *maximum* and *minimum*, which relates to contrast, while dimension 2 has the highest coefficients for *minimum* and *variance* which relates to spatial variations within the structure as we can observe in typical representants in respective dimensions in Fig. 9. Also see a supplementary video [movie_MDS_simmat_stats.avi], comparing three MDS dimensions of similarity judgements similarity matrix (top) with a similarity matrix obtained as Euclidean distance of all eleven statistics (bottom).

# General discussion

In this study, we set out to identify the perceptual core characteristics of wood. Characterising the visual appearance of wood is complex because of the variety in factors like colour, grain patterns, fine-scale relief and reflectance behaviour. Accordingly, a description in physical terms requires very high-dimensional measurements that capture the image projected by the material surface across all possible lighting conditions and viewing angles. Yet, we reasoned that when human observers are asked to compare samples—or judge the appearance of a single sample—they would rely on a relatively small number of dimensions that together summarise the overall 'look' of each surface and its texture—what we might call a 'visual signature' of the material (Sharan, Liu, Rosenholtz, & Adelson, 2013; Schmidt, Hebart, & Fleming, 2022).

Here, we wanted to estimate such an internal multidimensional representation by asking observers to make comparisons between samples. A secondary goal was to test the extent to which different methods of probing this putative representation yielded similar embeddings of the material samples. We reasoned that if observers draw on shared, core perceptual dimensions to judge the appearance of wood, it should be possible to probe this representation using distinct tasks.

To test this, we performed two experiments using movies of thirty samples of different wood veneers, rotating in such a way as to reveal both non-specular and specular appearance modes. In the first experiment, we took a data-driven approach, asking participants to make relative similarity judgments in a 2AFC task, from which we sought to derive underlying dimensions using the VICE algorithm (Muttenthaler, Zheng, McClure, Vandermeulen, Hebart, & Pereira, 2022). In the second experiment, we defined a set of ten appearance characteristics and asked participants to rate each sample in terms of all ten characteristics, effectively directly stating the location of each sample in a ten-dimensional appearance space. Our main findings can be summarised as follows:

- In Experiment 1, the VICE algorithm revealed that nine dimensions could account for 75% of the variance in the similarity judgments, consistent with the notion of a low-dimensional 'visual fingerprint' summary representation of their appearance.
- In Experiment 2, participants were consistent in their judgments of the ten appearance characteristics, suggesting agreement about the embedding of samples relative to one another.
- Comparisons between the two experiments showed a significant overlap between embeddings of the samples derived from the two tasks, providing further evidence for a core representation of wood materials, with similar-looking samples close to one another, and more distinct ones further away from one another within the multidimensional appearance space.
- The consistency between the two experiments can also be demonstrated by approximating the dimensions inferred from Experiment 1 as a weighted linear combination of the ratings in Experiment 2.
- Finally, a set of quite simple low-level image features, designed to capture similar appearance characteristics as the rating dimensions predict the ratings and VICE dimensions surprisingly well, using simple linear regression. Although these image features will not be the exact quantities that the visual system uses to represent and compare the wood samples, this shows how we can use straightforward image-computable models to predict perceived differences in

509        appearance (under constant viewing conditions).  This has potential practical applications in
510        many areas.
511    Our study also provides a proof-of-principle demonstration that it is possible to establish embeddings of
512    items from a single basic-level category (here: wood) within a perceptual space using either a subset of
513    all possible similarity comparisons, or through direct rating of particular features.  The study differed
514    from previous investigations in the use of movies rather than static images, capturing a wide range of
515    appearances for each sample, and in the comparison between similarity and appearance ratings.

516    **Limitations and future directions**

517    Although our study provides a first proof-of-principle for identifying perceptual dimensions within
518    categories, there are a number of important limitations of the approach, which we consider here.

519    *Limited number of wooden samples*

520    The stimulus set considered here consisted of only thirty samples of different wood veneers, as listed in
521    Table 1. This is one of the largest sets of wooden samples used in a psychophysical analysis to date, and
522    we carefully selected this set from a catalogue of over one hundred wood veneers so as to provide as
523    broad and uniform a range of appearances as possible. However, including a larger number of samples
524    would necessarily provide additional information about the embedding, and would potentially reveal
525    additional perceptual dimensions by covering a wider range of appearances. It would also be particularly
526    interesting to include in future work multiple samples of each species (see Table 1), to capture within-
527    item variability as well.  We would expect that although different samples would be clearly
528    discriminable, generally they would tend to occupy very close locations within the multidimensional
529    perceptual space.

530    *Limited observation and illumination geometry*

531    By using dynamic stimuli, in contrast to previous studies, which tended to offer only a single view of
532    each sample, we were able to provide observers with some information about how the appearance of
533    the samples changed depending on viewing conditions, including both specular and non-specular
534    conditions. Nevertheless, this still represented a limited subset of all possible lighting-sample-viewer
535    configurations.  We had to limit camera and light trajectories so that movies were of reasonable
536    duration. Based on pilot work with a range of different sampling parameters, we identified a rotation
537    that was of acceptable durations and that was intuitive for observers. As the appearance of wood does
538    not typically change much with polar angle, we limited polar viewing angles to $45^{\circ}$ and changed
539    azimuthal angles only. A comparison of image histograms from our videos with those of the full BTF for
540    the same material (at polar angles $45^{\circ}$ including over 400 images for different combinations of
541    illumination and view azimuthal angles) provided mean differences of $X^2$ lower than 0.10.  This leaves us
542    confident that the selected views were representative of the overall appearance.

543
544

545      **Table 1** A complete list of wood species used in the experiment.

| 01 | afzelia | 11 | white ash | 21 | rosewood |
|----|---------|----|-----------|----|----------|
| 02 | masur birch | 12 | ash heartwood | 22 | plane |
| 03 | pommele bubinga | 13 | maple burl | 23 | satinwood |
| 04 | oak | 14 | European lime (linden) | 24 | spruce |
| 05 | burr oak | 15 | macassar ebony | 25 | spruce knotted |
| 06 | smoked oak | 16 | movingui (lemon) | 26 | tineo |
| 07 | eucalyptus | 17 | olive | 27 | American cherry |
| 08 | gaboon | 18 | European walnut | 28 | tulipwood |
| 09 | pear | 19 | Peruvian walnut | 29 | wenge |
| 10 | European apple | 20 | padauk | 30 | zebrawood |

546

## *Limited size of samples*

548      On a related point, the visible area of the samples was around 50x50mm. This size was selected to
549      deliver fine surface details. On the other hand, for certain species, there may be low-frequency content
550      that was excluded by the small size. To compensate for this during video acquisition the location of the
551      captured area on the veneer specimen was carefully selected to demonstrate the main sample's
552      characteristics. A similar comparison of histogram statistics with BTF data over a large scale of image
553      plane resulted in similarly low differences in histograms, again indicating that the patch was
554      representative of the sample as a whole.

## *Limited coverage of triplets for similarity judgements*

556      In Experiment 1, we measured only a small subset of all possible stimulus triplets. Specifically, our
557      experiment had a coverage of 10%, which is nevertheless far greater than the less than 2% coverage
558      used in other studies using related data analyses (Hebart, Zheng, Pereira, & Baker, 2020). On the other
559      hand, our number of samples is considerably lower, greatly reducing the number of necessary trials. We
560      followed the recommendations in (Haghiri, Wichmann, & von Luxburg, 2020) to estimate the number of
561      judgements, although future studies could potentially increase the coverage further for small stimulus
562      sets like ours.

## *Stability of dimensions*

564      Statistical inference methods like VICE are stochastic, so repeated runs of the algorithm on the same
565      data can deliver slightly different outcomes. This naturally raises questions about the stability and
566      interpretation of the outcome. We tested a wide range of hyperparameter values, and found the values

567 we used delivered representative results. Importantly, although the exact number of dimensions varied
568 across runs, the meanings of those dimensions (i.e., the loadings across samples) were highly conserved.
569 This, along with the high extent to which the dimensions could predict similarity ratings gives high
570 confidence that the analysis delivered robust results.  Increasing the number and diversity of samples, as
571 well as the coverage would lead to even greater stability, although with obvious practical costs.  It is
572 nevertheless important to emphasise that in interpreting results on small and constrained stimulus sets
573 like ours, greater emphasis should be placed on the *embedding of items* within the multidimensional
574 space than on the precise number or direction of the dimensions returned by VICE (or related
575 algorithms).  The convergence between the ratings and the VICE analysis supports this view.

### *Intuitive interpretability of individual dimensions*

577 While some studies (e.g., Hebart, Zheng, Pereira, & Baker, 2020; Josephs, Hebart, & Konkle, 2023;
578 Schmidt, Hebart, & Fleming, 2022) have found that analyses similar to VICE deliver dimensions that are
579 highly intuitively interpretable, in our case, most of the dimensions appeared to be better understood as
580 weighted combinations of more intuitive factors. This can be seen in Fig. 9, for example, in which
581 samples are ranked by their values of the nine dimensions returned by VICE. Some of the dimensions
582 seem to capture intuitive concepts.  For example, dimensions 4 appears related to stripiness, and this is
583 consistent with the high loading of the 'Directionality' and 'Line' features in the multiple regression for
584 this feature.  Dimension 6, in contrast, seems to be approximately the opposite, with an emphasis on
585 samples with turbulent texture patterns rather than linear grain. However, for most of the other
586 dimensions the interpretation is less intuitive. This is likely due to the small and constrained sample set.
587 With diverse image sets that span the entire range of commonly occurring objects, for example (Hebart,
588 Zheng, Pereira, & Baker, 2020), almost all samples will have near-zero values of any given attribute,
589 while there are still sufficient numbers of images with high values to enable a dimension to emerge from
590 the analysis.  Indeed, such datasets are particularly well suited for seemingly meaningful individual
591 dimensions to be recovered by the sparse nonnegative matrix factorization. By contrast, within-category
592 samples, as in our experiments, tend to involve characteristics that are more uniformly distributed
593 across samples.  This is likely to be one of the reasons that the recovered dimensions were composites
594 of multiple factors.  Nevertheless, again it should be noted that we place greater emphasis on the
595 embedding of items within the space than on the exact orientation of the underlying dimensions.

### *Choice of rating dimensions*

597 There are practical limits to the number of appearance attributes that participants can feasibly be asked
598 to rate for each sample.  As with the majority of previous perceptual studies of wood surfaces
599 (Nakamura, Masuda, & Shinohara, 1999; Nordvik, Schütte, & Broman, 2009; Fujisaki, Tokita, & Kariya,
600 2015; Manuel, Leonhart, Broman, & Becker, 2015, Wan, Li, Zhang, Song, & Ke, 2021) we preselected a
601 list of visual properties in our rating experiment.  This list, of course, is likely to be incomplete as there
602 are potentially infinitely many ways of describing samples, including those that may make intuitive visual
603 sense, but which cannot easily be put into words. Nevertheless, we find that this set of dimensions leads
604 to intuitive and repeatable judgments, which are sufficient to capture an embedding of the samples
605 similar to that revealed by the similarity ratings and VICE analysis.  Future studies could also ask

606    participants, rather than the experimenters, to provide terms that describe important appearance
607    differences between samples, which other participants would then rate (see, e.g., Van Assen, Barla, &
608    Fleming, 2018).

609    **Conclusions**

610    Our study sought to identify core perceptual dimensions underlying the appearance of wood. Using
611    thirty movies of rotating planar wooden veneer samples, we asked participants to judge the similarity
612    between items and rate each sample along ten predefined dimensions. The results revealed a
613    consistent embedding of samples between the two tasks, suggesting a core internal representation of
614    the samples, capturing the overall 'look' of the samples in a relatively small number of dimensions.
615    These could be expressed as a weighted linear combination of the following ten attributes: brightness,
616    glossiness, colourfulness, directionality, complexity, contrast, roughness, patchiness/regularity, line
617    elongation, and spatial scale. The results not only reveal the core dimensions underlying the perception
618    of wood, they also provide a proof of concept demonstration for how perceptual dimensions underlying
619    judgments within a single basic-level category can be probed using multiple tasks.
620

621 **Acknowledgements**

626

627 **References**
628

629 Anderson, B. L. (2011). Visual perception of materials and surfaces. *Current biology, 21*(24), R978-R983.

630 Bell, S., Upchurch, P., Snavely, N., & Bala, K. (2015). Material recognition in the wild with the materials in
631 context database. *In Proceedings of the IEEE conference on computer vision and pattern recognition* (pp.
632 3479-3487).

633 Blondel, V. D., Guillaume, J. L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities in
634 large networks. *Journal of statistical mechanics: theory and experiment*, *2008*(10), P10008.

635 Bracci, S., & Op de Beeck, H.P. (2023). Understanding Human Object Vision: A Picture is Worth a
636 Thousand Representations. *Annual Review of Psychology, 74*, pp.113-135.

637 Carroll, J. D., & Arabie, P. (1998). Multidimensional scaling. *Measurement, judgment and decision
638 making*, 179-250.

639 Dana, K.J., van Ginneken, B., Nayar, S.K., & Koenderink, J.J. (1999). Reflectance and texture of real-world
640 surfaces, *ACM Transactions on Graphics*, Vol.18, Issue 18, pp.1-34

641 De Leeuw, J. R. (2015). jsPsych: A JavaScript library for creating behavioral experiments in a Web
642 browser. *Behavior research methods*, *47*, 1-12.

643 Ferwerda, J. A., Pellacini, F., & Greenberg, D. P. (2001). Psychophysically based model of surface gloss
644 perception. In *SPIE Human vision and electronic imaging vi,* Vol. 4299, pp. 291-301.

645 Filip, J., Vavra, R., Haindl, M., Zid, P., Krupicka, M., & Havran, V. (2013). BRDF slices: Accurate adaptive
646 anisotropic appearance acquisition. In *Proceedings of the IEEE Conference on Computer Vision and
647 Pattern Recognition* (pp. 1468-1473).

648 Fleming, R. W., Dror, R. O., & Adelson, E. H. (2003). Real-world illumination and the perception of
649 surface reflectance properties. *Journal of vision, 3*(5), 3-3.

650 Fleming, R. W., & Bülthoff, H. H. (2005). Low-level image cues in the perception of translucent materials.
651 *ACM Transactions on Applied Perception, 2*(3), 346-382.

652    Fleming, R. W., Jäkel, F., & Maloney, L. T. (2011). Visual perception of thick transparent materials.
653    *Psychological science, 22*(6), 812-820.

654    Fleming, R. W., Wiebel, C., & Gegenfurtner, K. (2013). Perceptual qualities and material classes. *Journal*
655    *of vision, 13*(8), 9-9.

656    Fleming, R. W. (2017). Material perception*. Annual review of vision science*, 3, 365-388.

657    Fujisaki, W., Tokita, M., & Kariya, K. (2015). Perception of the material properties of wood based on
658    vision, audition, and touch. *Vision research*, *109*, 185-200.

659    Haghiri, S., Rubisch, P., Geirhos, R., Wichmann, F., & von Luxburg, U. (2019). Comparison-based
660    framework for psychophysics: Lab versus crowdsourcing. *arXiv preprint arXiv:1905.07234*.

661    Haghiri, S., Wichmann, F. A., & von Luxburg, U. (2020). Estimation of perceptual scales using ordinal
662    embedding. *Journal of vision*, *20*(9), 14-14.

663    Haindl, M., & Filip J. (2013). Visual Texture: Accurate Material Appearance Measurement,
664    Representation and Modeling*. Advances in Computer Vision and Pattern Recognition,  Springer-Verlag*
665    *London*

666    Hebart, M. N., Zheng, C. Y., Pereira, F., & Baker, C. I. (2020). Revealing the multidimensional mental
667    representations of natural objects underlying human similarity judgements. *Nature human behaviour*,
668    *4*(11), 1173-1185.

669    Josephs, E. L., Hebart, M. N., & Konkle, T. (2023). Dimensions underlying human understanding of the
670    reachable world. *Cognition*, *234*, 105368.

671    Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass correlation coefficients for
672    reliability research. *Journal of chiropractic medicine*, *15*(2), 155-163.

673    Künstle, D. E., von Luxburg, U., & Wichmann, F. A. (2022). Estimating the perceived dimension of
674    psychophysical stimuli using triplet accuracy and hypothesis testing. *Journal of Vision*, *22*(13), 5-5.

675    Lewin, M., & Goldstein, I.S. (1991). Wood Structure and Composition,  *International Fiber Science and*
676    *Technology, CRC Press*

677    Manuel, A., Leonhart, R., Broman, O., & Becker, G. (2015). Consumers' perceptions and preference
678    profiles for wood surfaces tested with pairwise comparison in Germany. *Annals of forest science*, *72*(6),
679    741-751.

680    Marlow, P. J., Kim, J., & Anderson, B. L. (2012). The perception and misperception of specular surface
681    reflectance*. Current Biology, 22*(20), 1909-1913.

682    Maskey, M., & Newman, T. S. (2021). On measuring and employing texture directionality for image
683    classification. *Pattern Analysis and Applications*, *24*(4), 1649-1665.

684    McCamy, C. S. (1996). Observation and measurement of the appearance of metallic materials. Part I.
685    Macro appearance. *Color Research & Application*, *21*(4), 292-304.

686    Motoyoshi, I., Nishida, S. Y., Sharan, L., & Adelson, E. H. (2007). Image statistics and the perception of
687    surface qualities. *Nature, 447*(7141), 206-209.

688    Muttenthaler, L., Zheng, C. Y., McClure, P., Vandermeulen, R. A., Hebart, M. N., & Pereira, F. (2022).
689    VICE: Variational Interpretable Concept Embeddings. *Advances in Neural Information Processing*
690    *Systems*, *35*, 33661-33675.

691    Nakamura, M., Masuda, M., & Shinohara, K. (1999). Multiresolutional image analysis of wood and other
692    materials. *Journal of wood science*, *45*, 10-18.

693    Nicodemus, F.E., Richmond, J.C., Hsia, J.J., Ginsburg, I.W., & Limperis, T. (1977). Geometrical
694    considerations and nomenclature for reflectance*. NBS Monograph 160*, pp. 1-52

695    Nordvik, E., Schütte, S., & Broman, N. O. (2009). People's perceptions of the visual appearance of wood
696    flooring: A kansei engineering approach. *Forest products journal*, *59*(11-12), 67-74.

697    Paulun, V. C., Schmidt, F., van Assen, J. J. R., & Fleming, R. W. (2017). Shape, motion, and optical cues to
698    stiffness of elastic objects*. Journal of vision, 17*(1), 20-20.

699    Portilla, J., & Simoncelli, E. P. (2000). A parametric texture model based on joint statistics of complex
700    wavelet coefficients. *International journal of computer vision*, *40*, 49-70.

701    Rao, A. R., & Lohse, G. L. (1996). Towards a texture naming system: Identifying relevant dimensions of
702    texture. *Vision Research*, *36*(11), 1649-1669.

703    Schmidt, F., Hebart, M. N., & Fleming, R. W. (2022). Core dimensions of human material perception.
704    PsyArXiv. doi:10.31234/osf.io/jz8ks

705    Sharan, L., Liu, C., Rosenholtz, R., & Adelson, E. H. (2013). Recognizing materials using perceptually
706    inspired features. *International journal of computer vision, 103*, 348-371.

707    Sharan, L., Rosenholtz, R., & Adelson, E. (2009). Material perception: What can you see in a brief
708    glance?*. Journal of Vision, 9*(8), 784-784.

709    Sharan, L., Rosenholtz, R., & Adelson, E. H. (2014). Accuracy and speed of material categorization in real-
710    world images. *Journal of vision, 14*(9), 12-12.

711  Tamura, H., Mori, S., & Yamawaki, T. (1978). Textural features corresponding to visual perception. *IEEE*
712  *Transactions on Systems, man, and cybernetics*, *8*(6), 460-473.

713  Tanaka, M., & Horiuchi, T. (2015). Investigating perceptual qualities of static surface appearance using
714  real materials and displayed images. *Vision research*, *115*, 246-258.

715  Van Assen, J. J. R., Barla, P., & Fleming, R. W. (2018). Visual features in the perception of liquids. *Current*
716  *biology, 28*(3), 452-458.

717  Wan, Q., Li, X., Zhang, Y., Song, S., & Ke, Q. (2021). Visual perception of different wood surfaces: an
718  event-related potentials study. *Annals of Forest Science*, *78*, 1-18.

719  Wendt, G., Faul, F., & Mausfeld, R. (2008). Highlight disparity contributes to the authenticity and
720  strength of perceived glossiness*. Journal of Vision, 8*(1), 14-14.

721  Wendt, G., Faul, F., Ekroll, V., & Mausfeld, R. (2010). Disparity, motion, and color information improve
722  gloss constancy performance. *Journal of vision, 10*(9), 7-7.

723  Wiebel, C. B., Valsecchi, M., & Gegenfurtner, K. R. (2013). The speed and accuracy of material
724  recognition in natural images. *Attention, Perception, & Psychophysics, 75*, 954-966.

725
726
727

## Supplementary material

### 1. VICE algorithm training

We tested the VICE model on 76 different combinations of input parameters such as et, spike, slab, pi and distribution (gaussian, laplace) (cf. Muttenthaler,  Zheng, McClure, Vandermeulen, Hebart, & Pereira, 2022). Results are shown in Fig. S1(a), where the tested models are rank ordered according to test accuracy (red), with the corresponding training accuracy (blue). The converged models are highlighted as circles. Fig. S1(b) shows that the number of dimensions is relatively stable, within a range between 5 to 14 and a typical value of 10 dimensions.
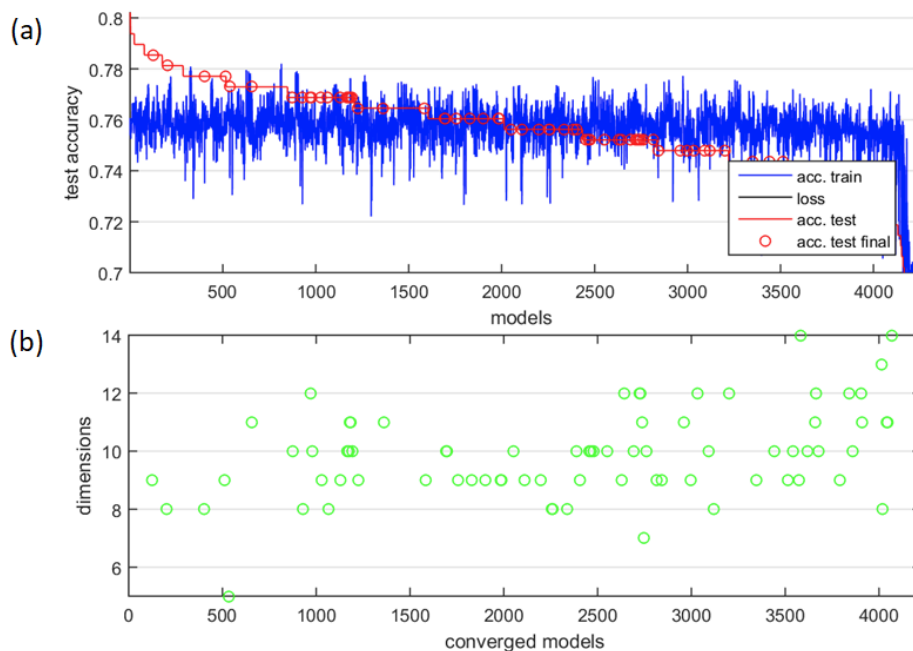


**Fig. S1** Results of grid search across different VICE model parameters. (a) Model accuracies on train (blue) and test (red) sets (across all tested models) sorted according to accuracy on test set (red), and (b) corresponding obtained numbers of dimensions for converged models (also denoted as circles in (a)).

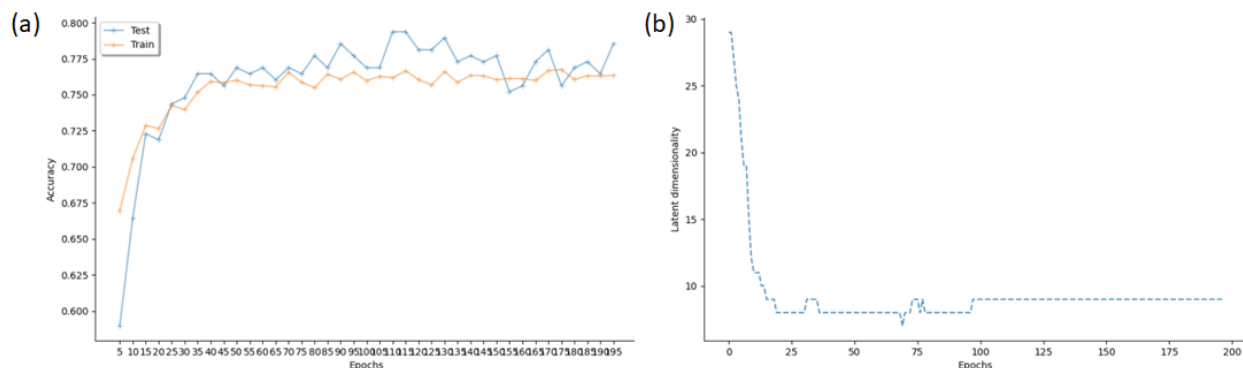Fig. S2 shows the training process of the best performing converged model with the highest accuracy.

**Fig. S2** Training process of the best performing model. (a) Model accuracy on the training (blue) and test (orange) dataset, (b) dimensionality reduction over 200 epochs of VICE algorithm.

### 2. Louvain community detection

We also applied the community detection method (Blondel, Guillaume, Lambiotte, & Lefebvre, 2008) on the estimated similarity matrix. The resulting three clusters visualised in Fig. S3 can be interpreted as (1) contrast/roughness, (2) non-directional/low frequency, (3) directional/high frequency modes. These results are in agreement with the results of hierarchical clustering and MDS analysis.



**Fig. S3** Clustering based on the Louvain community detection method divided the material samples into three clusters. See supplementary video [movie_Louvain.avi].

### 3. Similarity judgement data dimensionality analysis

As the dimensionality of our dataset is unknown, we follow a recent approach by (Künstle, von Luxburg, & Wichmann, 2022) to estimate the number of perceived dimensions from triplet experiments, based on triplet embedding accuracy. When splitting our triplet dataset into 90% training and 10% test samples, we obtain an ordinal Euclidean embedding (Haghiri, Wichmann, & von Luxburg, 2020) for the perceptual ratings. This procedure always leads to a decreasing triplet error (cross-validated on the validation set) with an increasing number of dimensions until a sufficient number of dimensions has

773    been reached. We ran a cross-validation with 10 repetitions, resulting in a drop of accuracy with more

774    than 6 dimensions. This suggests that inherent dimensionality of our dataset is close to 6 perceptual

775    dimensions. Note that our analysis shown in Fig. S1(b) reports a dimensionality of the typical estimated

776    similarity embedding between 8 and 10 dimensions. This seems to contradict the estimate of the

777    inherent dimensionality of 6 as reported above (and shown in Fig. S4). However, our linear regression

778    analysis (blue bars in Fig. 12(a)) suggests that several of our similarity dimensions (namely dimension 7)

779    cannot be reliably predicted from the appearance ratings, which might suggest that: (1) our rating

780    dimensions do lack some important visual features, or (2) the number of representational dimensions is

781    lower than the estimate of the VICE algorithm. In favour of the latter, the factor loadings of individual

782    dimensions (Fig. 3(b)) show a drop in loadings for dimensions higher than 5. Also, when using PCA on the

783    rating data to test whether intercorrelations (Fig. 8(b)) allow us to reduce the dimensionality, we end up

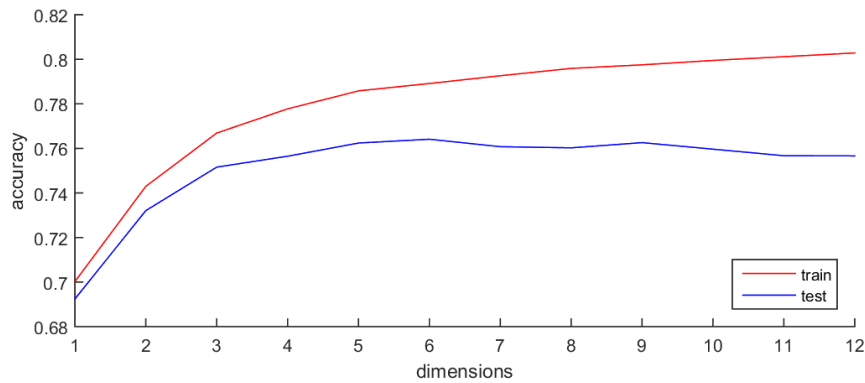784    with not more than 6 dimensions.

785



786

787    **Fig. S4** Triplet ordinal embedding error as a function of the number of dimensions for training and test

788    set of triplets from our similarity experiment.

789

790

791

792

793

794

795

796

797

798    **4.    Rating experiment details**

799

800    The interface of the rating experiment is shown in Fig. S5.
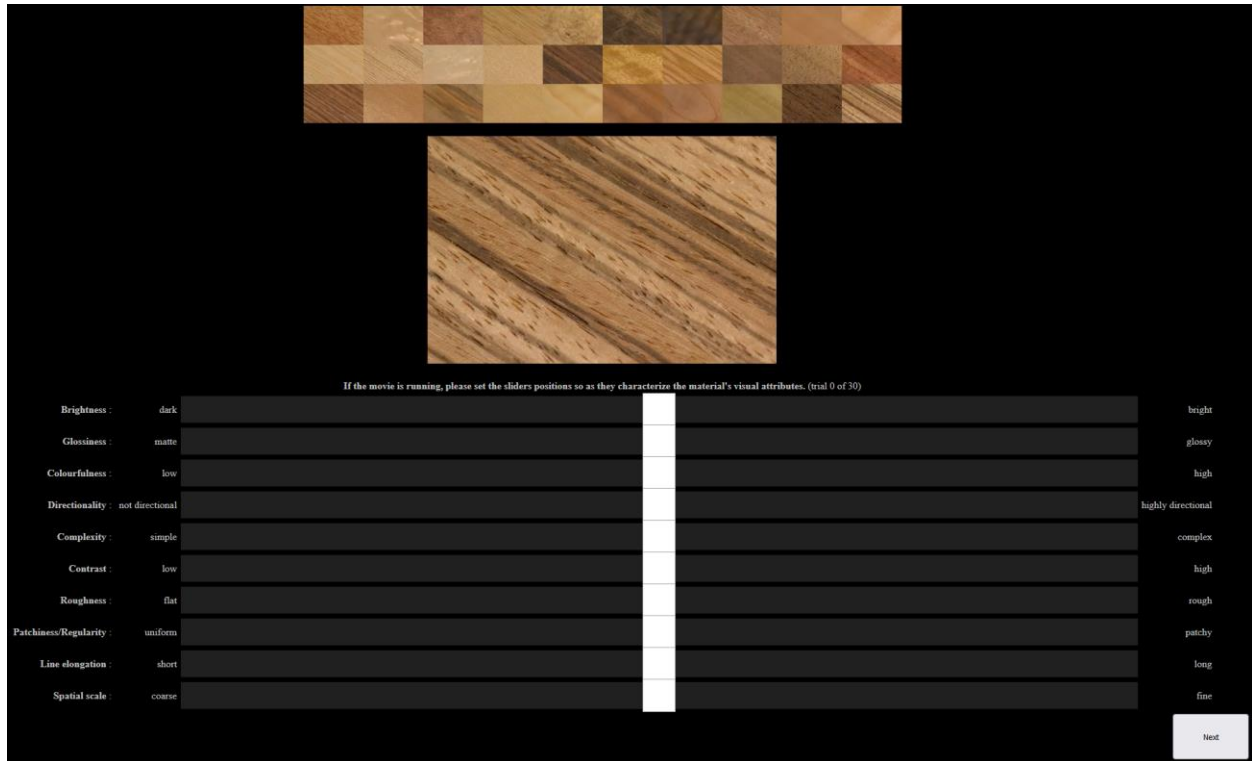


801

802    **Fig. S5** Example stimulus frame from the rating experiment.

803

804    Instructions of the rating experiment were as follows: "*Below the video, there are 10 sliders for the*

805    *visual material attributes. Your task is to adjust the slider position for each material. You may consider*

806    *the appearance of the other materials (at the top of the screen) to choose your rating appropriately*

807    *within the range we are testing.*"

808    The following visual attributes are evaluated:

809                1. **Brightness** - how bright is the material?

810                2. **Glossiness** - how shiny is the material?

811                3. **Colourfulness** - how colourful is the material?

812                4. **Directionality** - presence of directional structures in the texture

813                5. **Complexity** - how complex are the patterns on the surface?

814                6. **Contrast** - difference in brightness of surfaces patterns

815                7. **Roughness** - smoothness of surface profile, range of surface heights

816                8. **Patchiness/Regularity** - how uniform is the pattern?

817                9. **Line elongation** - are line elements shorter dashes or extended lines?

818                10. **Spatial scale** - are patterns large and broad, or small and fine?

819

820    **5.   Rating results analysis**

821

822    Mean values of participant responses for each resting attribute, and normalised distribution of allpooled

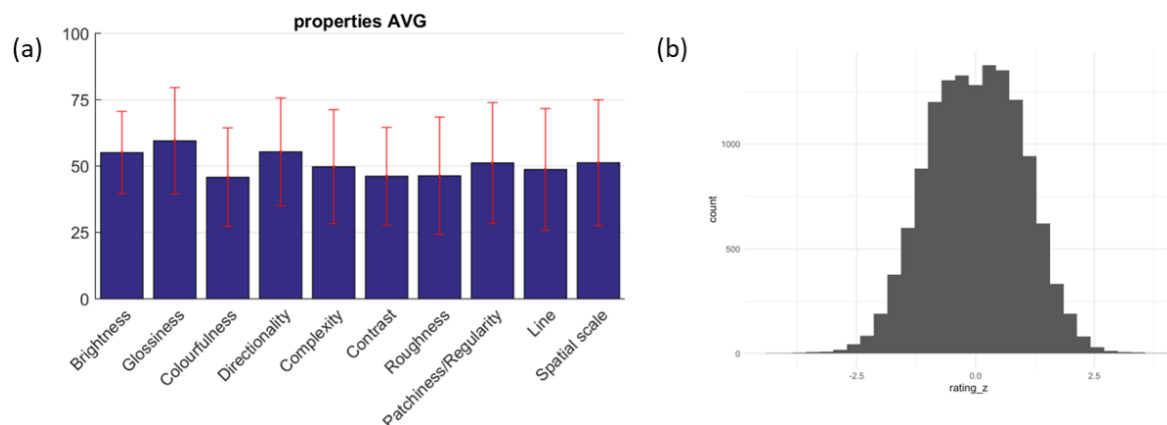823    ratings is shown in Fig. S6.

824



825

826    **Fig. S6** Rating data analysis. (a) Mean values of participant responses across all materials with SD values,

827    and (b) normalised distribution of all pooled ratings.

828

829    To evaluate the consistency between participants, we correlate the ratings of each participant within a

830    scale to the corresponding mean rating (Fig. S7). Although overall correlations are pretty high, there is a

831    heavy tail towards zero and even some negative correlations. Also, the consistency between participants

832    varies between rating dimensions, for example, with more consistent judgements for brightness

833    (stronger correlations and less variability). Tab. S1 shows intra-class correlations for individual rating
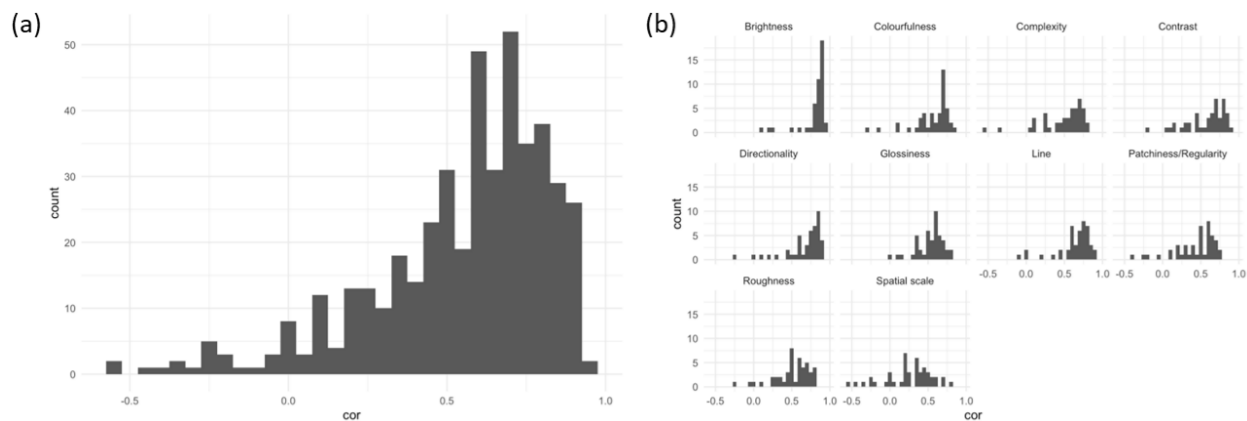
834    dimensions.



835

836    **Fig. S7** Correlations between individual ratings of participants to the mean ratings, within each rating

837    dimension. (a) Results across all attributes and (b) within individual dimensions.

838

839

840

841

842 **Tab.S1** Intra-class correlations for individual rating dimensions.

843

| rating dimension | single random raters | average random raters |
|---|---|---|
| brightness | 0.618 | 0.986 |
| glossiness | 0.219 | 0.927 |
| colourfulness | 0.262 | 0.941 |
| directionality | 0.416 | 0.970 |
| complexity | 0.197 | 0.917 |
| contrast | 0.301 | 0.951 |
| roughness | 0.220 | 0.927 |
| patchiness/regularity | 0.164 | 0.898 |
| line | 0.386 | 0.966 |
| spatial scale | 0.041 | 0.659 |

844

845

846

847 **6. Samples alignment along MDS dimensions**

848

849 The MDS analysis distributed our 30 samples to three dimensional space. Distribution of samples along

850 these dimensions is shown in Fig. S8, where red points represent MDS of VICE similarity model and blue

851 MDS of rating attributes (the first two rows) and computational statistics (the third row).
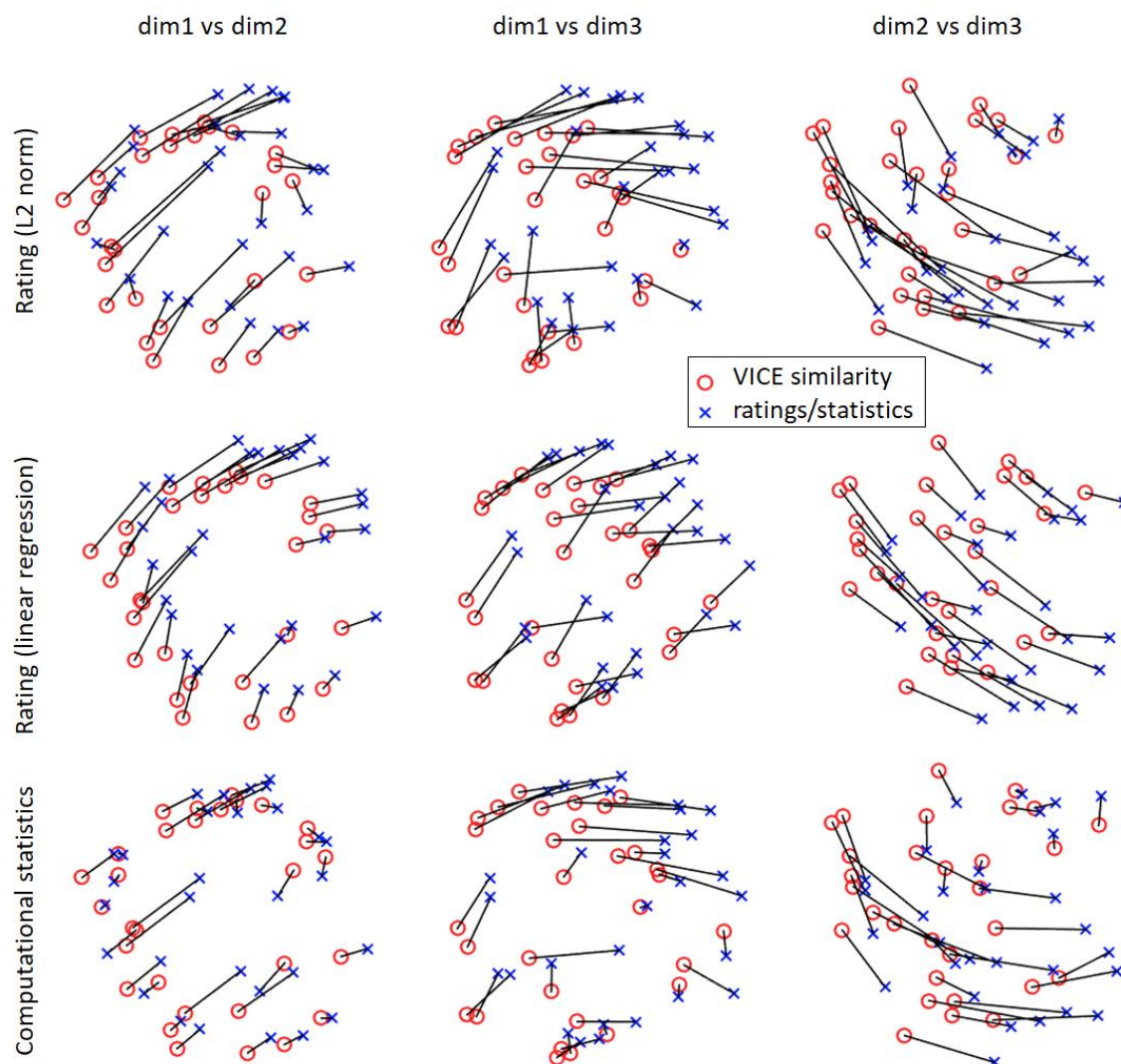
**Fig. S8** Procrustes alignment of MDS dimensions computed from similarity matrices. (red) VICE model similarity MDS, (blue) MDS of similarity matrix obtained by L2- norm of rating attributes (the first row), linear regression of rating attributes similarity matrices (the second row), and computational statistics MDS (the third row).

## List of supplementary movies

1. [movie_samples_stimuli.avi] - 30 test wood video sequences used in the experiments
2. [movie_similarity_vs_rating.avi] - rank ordered samples (left) according to loadings values of similarity dimensions, (right) mean rating attributes (the five closes and 5 the most distant)
3. [movie_MDS_simmat_linreg.avi] - distribution of samples along three MDS dimensions (top) for similarity judgements, (bottom) for rating study
4. [movie_MDS_simmat_stat.avi] - distribution of samples along three MDS dimensions (top) for similarity judgements, (bottom) for computational statistics obtained from image sequence.
5. [movie_similarity_scaled.avi] - rank ordered samples scaled according to loadings values of similarity dimensions
6. [movie_Louvain.avi] - result of community detection using Louvain method (computed from similarity matrices), distributing samples to three clusters for (top) similarity judgements and (bottom) rating study.