# Human Emotion Knowledge Representation Emerges in Large Language Model and Supports Discrete Emotion Inference

Ming Li[1, 2, #], Yusheng Su[3, #], Hsiu-Yuan Huang[4], Jiali Cheng[5], Xin Hu[1, 2 ,6], Xinmiao Zhang[1, 2], Huadong Wang[3], Yujia Qin[3], Xiaozhi Wang[3], Zhiyuan Liu[3, *], Dan Zhang[1, 2, *]

[1] The Department of Psychology, Tsinghua University

[2] Tsinghua Laboratory of Brain and Intelligence, Tsinghua University

[3] The Department of Computer Science and Technology, Tsinghua University

[4] School of Computer and Communication Engineering, University of Science and Technology Beijing

[5] Miner School of Computer and Information Sciences, University of Massachusetts Lowell

[6] Department of Psychiatry, University of Pittsburgh

**Abstract:** How humans infer discrete emotions is a fundamental research question in the field of psychology. While conceptual knowledge about emotions (emotion knowledge) has been suggested to be essential for emotion inference, evidence to date is mostly indirect and inconclusive. As the large language models (LLMs) have been shown to support effective representations of various human conceptual knowledge, the present study further employed artificial neurons in LLMs to investigate the mechanism of human emotion inference. With artificial neurons activated by prompts, the LLM (RoBERTa) demonstrated a similar conceptual structure of 27 discrete emotions as that of human behaviors. Furthermore, the LLM-based conceptual structure revealed a human-like reliance on 14 underlying conceptual attributes of emotions for emotion inference. Most importantly, by manipulating attribute-specific neurons, we found that the corresponding LLM's emotion inference performance deteriorated, and the performance deterioration was correlated to the effectiveness of representations of the conceptual attributes on the human side. Our findings provide direct evidence for the emergence of emotion knowledge representation in large language models and suggest its casual support for discrete emotion inference.

[#] These authors contributed equally: liming16@tsinghua.org.cn, yushengsu.thu@gmail.com

[*] Corresponding authors: {liuzy, dzhang}@tsinghua.edu.cn

The source code can be obtained from https://github.com/thunlp/Model_Emotion.

# Introduction

At the end of Shakespeare's *Hamlet*, Horatio looks at Hamlet and says, 'Now cracks a noble heart. Good night, sweet prince, / And flights of angels sing thee to thy rest.' Even though Horatio's facial expression and prosody are not perceived, and the line does not contain exact words to denote emotions, anyone reading this sentence would clearly understand Horatio's emotions of grief and admiration. This phenomenon demonstrates a unique ability of human beings to infer emotions from the combination of linguistic symbols[1–3], i.e., text. While the emotion inference ability has been suggested to be crucial for our social life[4–8], how we infer emotions from text-based expressions remains to be elucidated.

A growing body of research has convergently supported the constructionist accounts of emotion[9,10], which argues that conceptual knowledge about different emotions (emotion knowledge) is essential for emotion inference. Specifically, emotion knowledge is believed to provide the necessary conceptual associations to support the inferences of specific emotions[9,10]. In the example of Horatio, the inference of grief and admiration may arise from an ongoing interaction between the meaning abstracted from the text (e.g., the situation of "a man died") and our emotion knowledge about these two kinds of emotions (e.g., a belief that "losing is grief" [11,12]). In this sense, emotional meaning is not entirely predetermined by the external stimuli, but is interpreted in terms of our emotion knowledge[9]. Consistent with this view, behavioral studies have revealed that people's emotion ability is related to their individual[13–16] and cultural[17,18] differences in emotion knowledge. Developmental studies further showed that children's emotion development could be explained by the acquisition of language-mediated emotion knowledge[19–23] (for reviews, see refs.[24,25]). Together, these findings have provided correlational evidence for the functional importance of emotion knowledge for emotion inference.

Despite this promising evidence, there is no direct evidence regarding whether human emotion inference relies on the use of emotion knowledge[26]. By presenting specific cues prior to emotional tasks (i.e., priming[27]), it has been shown that manipulating the human mindset's access to certain emotion knowledge can alter performance on emotional task[28]. Specifically, priming with emotion-congruent cues enhances the speed and accuracy of recognizing specific emotion[29,30], while repeated priming impairs it[31,32]. However, manipulating the access does not guarantee a direct manipulation of emotion knowledge *per se*[33,34], making the relevant evidence circumstantial. A more direct approach would be to explore the behavioral consequences of neurological disorders[35,36] or brain stimulation[37] that occur in brain regions possibly related to emotion knowledge. For instance, due to lesions in the anterior temporal lobe (ATL) that have been suggested for conceptual knowledge storage[38,39], semantic dementia patients with substantial impairments in the use of conceptual knowledge were shown incapable of categorizing discrete emotions at a finer granularity than positive-negative[35]. However, our limited understanding of the neural representations of different emotions (and the corresponding emotion knowledge)[40,41], as well as the lack of reliable intervention techniques with sufficient precision (e.g., transcranial magnetic stimulation with a centimeter-level spatial

resolution)[42], have prevented us from exploring further and deeper in this direction. In sum, verifying the possible causal role of emotion knowledge on emotion inference is challenging with the state-of-art psychological and neuroscience approaches.

Considering that human emotion knowledge is at least partially, if not entirely, derived from language[21,43–45], large language models (LLMs)[46–49] may serve as a prospective tool to investigate the reliance on emotion knowledge for emotion inference. The investigation of language-derived human emotion knowledge rests on the fact that text-based computing offers the opportunity to mine the deeper characteristic of human culturomics[50]. Moreover, several recent studies have shown that various human conceptual knowledge can emerge from LLMs by pre-training on massive amounts of unlabeled text corpora[51–56] (for a review, see ref.[57]), including concept taxonomy[55], social biases[58–60], moral norms[56]. Although no research on emotion has been reported yet, it is plausible to expect an effective representation of human emotion knowledge by LLMs. More importantly, unlike the human brain, the LLMs could be manipulated more easily for a more direct understanding of the relationship between emotion knowledge and inference. Specifically, the artificial neurons in LLMs could be selectively manipulated for their functional relevance of specific knowledge, and the manipulated neurons in LLMs could then be tested for their performance on certain inference tasks, as has been practiced for topic classification[61]. Thus, we can follow the same paradigm to manipulate artificial neurons on emotion inference tasks. Such a practice would be similar to the intervention techniques used in neuroscience research[62,63], but with high-precision manipulation (e.g., single artificial neuron) and testing. Given that the LLMs are obtained based on large-scale human text databases, it is reasonable to consider that the findings from LLMs could at least partially imply the mechanism for human emotion inference.

In the present study, we aimed to explore the representation of human emotion knowledge in LLMs and its potential support for discrete emotion inference. To more fully reveal emotion knowledge, we employed the latest taxonomy that expresses the nuances of 27 discrete emotions[64], which is accompanied by a publicly available text dataset and suitable for further emotion-specific representations using LLMs. We utilized RoBERTa[46] in our experiments as it is a typical LLM. To stimulate learned knowledge in LLM to infer the corresponding emotion, we trained 27 emotion-specific prompts[65]. This mechanism[61,66] can be analogous to the 'priming' operation in human psychology[27,29] since both the prompt and priming act as a cue message to set the LLM and the human brain, respectively, to an appropriate state for the upcoming task. Therefore, when we only input the emotion-specific prompt without concatenating any text into the LLM, its neuron activation (hidden state) values could be seen as the representation of emotion-specific knowledge. We also conducted behavioral experiments to obtain human emotion knowledge (conceptual structure and conceptual attributes of emotions), which could be compared with the LLM's representations from a higher-level functional perspective through representational similarity analysis (RSA)[67]. Subsequently, guided by the behavioral data, we could locate and manipulate the artificial neurons in the LLM associated with specific conceptual attributes to investigate its support for emotion inference tasks. We further explored whether the effectiveness of representations of the conceptual attributes on the human

side could predict their contributions to emotion inference on the LLM side, which could provide evidence for a deeper understanding of our results.

Based on previous human-like representations on LLMs[51–60], we expected that the prompt technique would stimulate the artificial neurons in LLMs to represent human-like knowledge about different discrete emotions. If these representations were functional relevance, then by manipulating artificial neurons associated with the specific conceptual attribute, we would expect the corresponding LLM's emotion inference performance to deteriorate, and the performance deterioration to correlate with the effectiveness of representations of the conceptual attributes on the human side. We hope that these results will also shed light on the mechanisms by which humans make emotion inferences.

# Results

**Prompt tuning for inferring 27 discrete emotions.**

In order to activate the possible representation of emotion knowledge in the LLM, we trained 27 emotion-specific prompts to infer the corresponding emotion on the training set. We reported their average accuracies over different random seeds by evaluating the test set (**Table 1**). The average accuracy for each of the 27 emotion inference tasks varied from 68.04% (realization) to 96.43% (gratitude). By relating average accuracy to rater agreement on dataset annotations, the LLM's performance showed a human-like trend, i.e., the more agreement human raters had on annotating discrete emotions, the more accurate the inferred emotions were; Pearson's $r(25)$ = .631, $p$-value < .001.

During training (prompt tuning), we optimized only prompts, which were composed of trainable tokens, while freezing all parameters of the LLM (**Figure 2a**). Because these prompts can be considered to activate the corresponding LLM task state, the LLM neuron activation values in response to the emotion-specific prompts can be considered to represent knowledge about corresponding emotion[61,66], which were then used for further analysis (**Figure 2b** and **Methods**).

**The LLM exhibits a human-like conceptual structure of emotions.**

The representational dissimilarity matrices (RDMs) for the conceptual structures of 27 emotions calculated by language-model indexes are shown in **Figure 3a**, from the LLM (RoBERTa)-based emotion knowledge representation (neuron activation), as well as the prompt embedding[68] and the GloVe word embedding[69,70] for comparison (see **Methods**). The RDM based on human behavioral data collected in the similarity judgment experiment, is shown in **Figure 3b**.

RSA revealed a significant second-order similarity between the LLM-based emotion knowledge representation and the human emotion knowledge, reaching a moderate level with Kendall's $tau$ = .267 (bootstrap 95%CI .238 to .273, $p$ < $10^{-5}$). The LLM-based emotion knowledge representation was more similar to human emotion knowledge, compared to the emotion-specific prompt embedding (Kendall's $tau$ = .112, bootstrap 95%CI .083

to .134, $p < 10^{-5}$) and the GloVe word embedding[69] (Kendall's *tau* = .046, bootstrap 95%CI .027 to .065, $p < 10^{-5}$). The results are shown in **Figure 3c**.

**The LLM-based emotion knowledge relies on human conceptual attributes of emotions.**

To explore the possible content of the LLM-based emotion knowledge, RSA was employed to evaluate the second-order similarity between the RDM of the LLM-based emotion knowledge representation and the human RDMs of 14 representative conceptual attributes, which were obtained from the conceptual attribute rating experiment (see **Figure 4a** and **Methods**). The second-order similarities were statistically significant for 10 out of the 14 conceptual attributes (**Figure 4b**), with Kendall's *tau* varying from .068 (control), .221 (valence), to .260 (other-related).

The representations of the 14 conceptual attributes in the LLM were further investigated by employing searchlight RSA to evaluate the second-order similarity between the RDM of every single artificial neuron and the human RDMs of 14 representative conceptual attributes. According to the rank of Kendall's *tau*, the most relevant artificial neurons for different conceptual attributes were less overlapping. Artificial neurons for each conceptual attribute, except for "other-related", were distributed in all layers of the LLM rather than concentrated in specific layers. **Figures 4c and 4d** show the example of the top 4,000 attribute-specific neurons; see **Supplement Figure 2** for results of different top N.

**Conceptual attributes causally contribute to emotion inference.**

The causal contribution of conceptual attributes to infer discrete emotions was then revealed by manipulating specific artificial neurons during emotion inferring (**Methods** and **Figure 5a**). Compared to the original accuracy without manipulation, we found a decrease in the accuracy of emotion inference on LLM with selective manipulation (**Supplementary Figures 3**). This deterioration of inference performance still held and revealed the causal contribution of conceptual attributes when compared to randomly manipulating the same number of neurons (the one-sided paired *t*-test on 12 random seeds determined significance, corrected for *p*-values with FDR; **Supplementary Figure 4**). The most prominent causal contributions arose in the cases of manipulating 4,000 attribute-specific neurons, which is shown in **Figure 5b**.

Since these conceptual attributes were proposed as the underlying dimensions for differentiating discrete emotions, we expected that manipulation of a given conceptual attribute contributes equally to inferring various discrete emotions (**Figure 5c**), i.e., less heterogeneity. The heterogeneity was tested by Hartigan's dip statistic[71,72] and was insignificant with a minimum *p*-value of .127 (**Supplementary Table 1**), suggesting weaker heterogeneity in conceptual attributes' causal contribution (differences in accuracy) to various emotion inference tasks.

We further demonstrated that the performance deterioration of emotion inference on the LLM was not independent, but significantly related to the effectiveness of representations of the conceptual attributes on the human side (**Methods**). The strongest correlation arose when manipulating the top 4,000 attribute-specific neurons in the LLM

with averaged $r$ = -.464 (**Figure 5d**; $t(26)$ = -16.714, Fisher-based s.e. = .030, $p < 10^{-15}$). In addition, as the number of manipulated neurons increased from 1000 to 6000, this prediction showed a stronger and weaker tendency (**Figure 5e**), suggesting that there may be "floor effects" and "ceiling effects" [73] on the causal contributions of conceptual attributes in the LLM.

# Discussion

In the present study, we compared the knowledge of the LLM to infer discrete emotions with human emotion knowledge and found that the LLM exhibits a human-like conceptual structure of emotions. More importantly, distinct artificial neurons in the LLM represent different conceptual attributes of emotions. By manipulating attribute-specific neurons and observing the LLM's performance on emotion inference, we revealed the causal contribution of conceptual attributes to inferring various discrete emotions.

The findings that the neuron activations in the LLM effectively represented human-like emotion knowledge extend our understanding of language computable human knowledge to the emotion domain. While previous studies have shown the association between language and emotion [43,44], we further demonstrated that the representation of discrete emotion's conceptual structure and underlying attributes could spontaneously emerge from the emotion-independent pre-training[46] on large language corpora. Notably, this piece result was achieved with the recently developed prompt technique[65], and the prompt-related neuron activations in the LLM showed the highest similarity compared to the prompt embeddings and the word embeddings (**Figure 3c**). As the prompt-related neuron activations were conceptually similar to priming-related neural activities in human neuroscience experiments[74–76], our results are in favor of taking a human-like perspective to analyze machine learning models[77–82]. This approach could facilitate not only the improvement of more safe and socially intelligent artificial agents[56], but also the exploration of the functional emergence of complex human cognitive abilities[83,84]. For instance, although some evidence from computer vision suggests that the ability to recognize stereotypical facial expressions may rely on domain-specific experience (e.g., faces vs. objects[79,85]), in our study, human emotion knowledge representation in the LLM was sufficient to emerge spontaneously from domain-general language experience[86]. Our computational evidence further supports a valuable but understudied hypothesis about emotion development: humans can learn discrete emotions directly from everyday language use without solid supervision[24,25] (for debates, see refs.[87–90]).

The observed reliance of the LLM's emotion knowledge on human conceptual attributes of emotions may shed light on the mechanisms of its emotion inference. Most of these theory-driven attributes (**Figure 4a;** e.g., valence[91], happy[92] and fairness[93]) explained the data-driven representations of task-related emotion knowledge (**Figure 4b**), consistent with their importance to human emotions in historical research. A few conceptual attributes (e.g., arousal) were not contained in the LLM-based emotion knowledge, probably because of the incompleteness of LLMs' language knowledge[45]. Interestingly, the

distributions of artificial neurons in the LLM corresponding to different conceptual attributes tended to be distinct (**Figure 4c**), suggesting possibly unique contributions of these attributes for the LLM's emotion knowledge representation. In addition, the attribute-specific neurons were distributed across all layers for most attributes (**Figure 4d**), possibly implying the involvement of both early-stage and late-stage processings for emotion inference task[6–8]. These results extend our understanding of the similarity of computational language models to humans, showing that the computational models not only represent human knowledge or concepts[51–57], but also adopt a human-like approach to information processing. Furthermore, these results may also inspire our exploration of the neural mechanisms underlying human emotion inference. For instance, the distributed attribute-specific neurons across all layers might resemble the reported distributed brain networks for semantic processing[94–96]. The relatively distinct distribution of the conceptual attributes of emotion could also suggest independent neural representations of these attributes in the human brain.

Most importantly, our more rigorous manipulation of the attribute-specific neurons of the LLM (**Fig.5a**) demonstrated possible causal support of discrete emotion inference by emotion knowledge. This piece's result could contribute to a central and ongoing debate in emotion science about the nature of human emotions[97–103], i.e., whether they have an essential core or are constructed as conceptual categories (for a review, see[104]). Whereas previous lesion studies revealed a causal role of possibly the overall conceptual system on emotion perception[30,35,36], our use of LLMs as a proxy for the human conceptual system in emotion inference tasks validates a similar conclusion. We further refine this causal role in the context of the contribution of conceptual attributes underlying discrete emotions. By pointing out the weak heterogeneity in the contribution to 27 discrete emotions (**Figure 5c and Supplementary Table 1**), we suggest a unifying underpinning mechanism for LLMs and possibly also for the human brain to infer various emotions[35]. Our view is reinforced by the recent neuroimaging findings that one broad ensemble containing multiple brain networks represents a range of emotions[105], rather than distinct emotions consistently and specifically activating local brain regions (for a meta-analysis, see ref.[106]). Beyond the shared anatomical basis, we further propose that emotion knowledge may uniformly support the processing of various emotions in human brains (for the debates, see ref.[40,41]). Based on the correlation between the causal contribution of conceptual attributes on LLMs and their effectiveness of representations for humans emotion knowledge (**Figure 5d and 5e**), it is reasonable to entail that our computational approximation is somewhat analogous to the causal mechanisms by which human mindsets infer various emotions.

Since LLMs and the human brain have been suggested to share similar computational principles for language[77], LLMs can serve as a potential reference in the future to help us understand the algorithms that the human brain relies on to infer emotions from text-based expressions and other semantic cues (e.g., speech content). Future research could also consider how semantic-based operations[107–109] drive the supramodal representation of (either perceived or inferred) emotions in the brain[8,110–112]. For example, there is growing evidence that the brain can convergently process and integrate emotional cues across modalities (e.g., facial expressions and prosody) and represent their conceptual meaning in amodal areas[110], such as the medial prefrontal cortex and the posterior superior temporal

sulcus[111,113,114]. Suppose the activity of these amodal areas during emotion perception fits with the LLM's hidden state values. In that case, possible neural mechanisms of semantic processing involved in emotion perception can be revealed. Although it is challenging to verify any causal mechanism in the human brain conclusively, the approach we applied in the current study, namely LLMs, could serve as a helpful testbed for exploring human abilities related to concepts or language in the future, such as the learning[25,109,115,116] and inference [117–120] of other abstract/social categories.

In conclusion, our study reveals that the LLM can represent human-like emotion knowledge with appropriate prompts, which are constructed from its general language experience learned in emotion-independent pre-training tasks. Our results support that language-derived emotion knowledge organizes necessary conceptual associations and causally supports discrete emotion inference from text-based expressions. Future research could combine computational methods with human experimental approaches, such as transcranial magnetic stimulation[121] and cognitive intervention tools [122], to expand understanding of the language-derived emotion construction and the biological mechanism of emotion inference.

# Methods

**Text dataset.**

The text dataset we employed in our study is GoEmotions ([https://github.com/ google-research/google-research/tree/master/goemotions](https://github.com/ google-research/google-research/tree/master/goemotions))[64], which contains 58,009 English Reddit comments manually annotated with 27 discrete emotions. Since different human raters rated each comment to obtain consistent emotion annotations, the agreement of human annotations was estimated via interrater correlation[123]. Subsequently, for each of the 27 emotions, this dataset could be viewed as a task to infer whether a text-based expression belongs to the corresponding emotion category (Yes/No). Following the authors of the dataset, we divided the dataset into training (80%), development (10%), and test (10%) sets. We did not use the development set in any subsequent operation because its proposed purpose was incompatible with this study.

**Prompt tuning.**

Formally, $\mathcal{M}$ was a LLM model, RoBERTa. Given an input text with $\mathrm{n}$ tokens $\mathrm{X} = \{w_1, w_2, \cdots, w_n\}$, RoBERTa first converted them into input embeddings $\boldsymbol{X_e} \in \mathrm{R^{n \times d}}$, where $d$ was the dimension of the embedding space. We pre-pended $l$ randomly initialized trainable tokens $\boldsymbol{P_e} \in R^{l \times d}$ before the input matrix $X_e$, and formed the modified input embeddings $[P_e, X_e] \in R^{(l+n) \times d}$. A special $[MASK]$ token was additionally pre-pended before the prompts, which would output the probability of label tokens. The objective $(O)$ was to maximize the likelihood of the desired output $y$:

$$O = P_{\mathcal{M}}([MASK] = y \mid [\boldsymbol{P_e}, \boldsymbol{X_e}]).$$

During the prompt tuning, we only optimized the trainable tokens $(P_e)$ while freezing the whole parameters of a RoBERTa ($\mathcal{M}$) to maximize the above objective.

To obtain the corresponding prompt of each discrete emotion on RoBERTa, we re-framed the 27-class emotion dataset of GoEmotions[64] into 27 emotion inference tasks. For instance, for the emotion "remorse", if a text belonged to the category "remorse", then we re-labeled the text with $y$ = "yes"; otherwise, $y$ = "no". In this way, we obtained the new training data for each discrete emotion. During training, we set the prompt length to $l = 100$ and the prompt dimension to $d = 768$. After conducting prompt tuning individually for each emotion inference task, we obtained all prompts $\{P_e^c \in R^{100 \times 768} \mid c \in C\}$, where $C$ was the set of 27 discrete emotions.

In order to avoid statistical bias, for each emotion inference task, we trained prompts 12 times with 12 random seeds; all of these 12 prompts have been evaluated on the test set, respectively. See **Table 1** for their average performance on 27 emotion inference tasks.

**Neuron activation in response to task prompts.**

Previous works[124,125] have indicated that the values of the aritificial neurons in the feed-forward layers $\text{FFN}(\cdot)$ of LLMs[86], RoBERTa, correspond to specific model behaviors. Some studies[61,66] have taken a further step to utilize trained prompts to stimulate RoBERTa and found that the prompts of similar tasks would have similar values of the artificial neurons. In this sense, we hypothesized that the neuron activation values could represent the task-specific knowledge (i.e., emotion knowledge), which could facilitate us to manipulate specific neurons for the purpose of manipulating specific emotion knowledge.

In our setting, the values of artificial neurons $s$ were the values of hidden states between the FFN layer in a Transformer. Specifically, we could denote the FFN layer as:

$$\text{FFN}(x) = \text{GELU}(xW_1^\top + b_1)\, W_2 + b_2,$$

where $x \in R^d$ was the input embedding, $W_1, W_2 \in R^{d_m \times d}$ were trainable matrices, and $b_1, b_2$ were bias vectors. The value of artificial neurons was $v = xW_1^\top + b_1$.

For each task, we input the sequence, $\{[MASK], P, <s>\}$, into RoBERTa, where $P$ was the emotion-specific prompt, $<s>$ was the special token indicating the start of an input sentence. Finally, we stacked the values of artificial neurons in all FFN layers of RoBERTa to get the overall neuron activation values $\text{AS}(P)$ (https://github.com/thunlp/Prompt-Transferability) for each emotion inference task:

$$\text{AS}(P) = [v_1; v_2; \ldots; v_L],$$

where $L = 36{,}864$ was the total number of artificial neurons.

We then calculated the paired cosine distances of the vectors for 27 emotion tasks and formed an RDM for each of 12 random seeds. The RDMs of 12 random seeds were then averaged to simulate the conceptual structure of emotions based on the LLM.

**Prompt embedding.**

As the embedding vector of prompts may also reflect task-related knowledge[68], we used the prompt embedding of 27 emotion tasks to form 12 RDMs for 12 random seeds by computing pairwise cosine distances. Then we averaged the 12 RDMs to simulate the conceptual structure of 27 emotions obtained from prompt embedding.

**GloVe word embedding.**

To obtain the lexical-level semantic representation of emotion concepts, we leveraged GloVe embeddings[69]. Based on the distributed hypothesis[126], GloVe is a log bilinear model which can be efficiently trained on large-scale raw corpora.

Since the GoEmotions was composed of the extracted comments from the Reddit dataset, for a fair comparison, we crawled the training data of GloVe from publicly available Reddit comments (https://www.kaggle.com/datasets/leighplt/glove-reddit-comments). We followed the same script (https://github.com/leigh-plt/glove.reddits) to pre-process and leveraged the official GloVe implementation tool (https://github.com/stanfordnlp/GloVe) to train our GloVe word embedding. We set the dimension of the word embedding to 768, the size of the sliding window to 12, and the minimum word frequency to 15. In order to avoid statistical bias in the following analysis, the model was repeatedly trained 12 times under different random seeds.

Finally, we obtained 12 different word embedding matrices with a vocabulary size of 172k and used the word embedding of 27 emotion words to form 12 word embedding RDMs by computing cosine distances. Then we averaged the 12 RDMs to simulate the conceptual structure of emotions obtained from word embedding.

**The similarity judgment experiment.**

We adopt a similarity judgment task to measure humans' conceptual structure of emotions. Fifteen English-speaking participants (8 females, mean age = 33 years) were recruited from Prolific and asked to complete the task online.

They judged the subjective similarity of 27 emotion concepts (and neutral) using a 9-point Likert scale (1=most dissimilar, 9=most similar) without criteria cues. These 27 emotions and neutral were presented simultaneously on the screen in word form. However, participants judged the similarity of only the two words with black borders each time. There were no response time limits but instructions to participants to respond by first sense when they hesitated.

We retained similarity scores between 27 emotions (351 pairs in total) and replaced missing values (3 per participant) with the average score across participants. Then, for each participant, these scores were subtracted by 10 to indicate the dissimilarity (ranging from 1-9) and used to form an individual RDM, i.e., a 27 by 27 symmetric matrix with a diagonal of 0 to indicate that any emotion is equal to itself. Each RDM reflected one participant's conceptual structures of emotions.

The Institutional Review Board at the Department of Psychology, Tsinghua University, approved the experimental procedures. All participants gave their informed consent.

**The conceptual attribute rating experiment.**

In this study, we chose the most representative conceptual attributes in the existing psychological emotion theories, such as affective attributes[127,128] (what an emotion might feel like), appraisal attributes[129,130] (what the antecedents of emotion might be), and basic emotions attributes[131,132] (six prototypical expressions that might explain all emotion expressions). We obtained human ratings for 14 candidate conceptual attributes (two affective, six appraisals, and six basic emotions) of 27 emotions from Prolific through three independent surveys. We started with these theory-driven conceptual attributes of

emotions to explain the representation of emotion knowledge and to investigate its contribution to emotion inference.

Thirty participants (15 females, mean age = 33 years) were recruited for affective attributes to rate all emotion concepts' arousal and valence. The emotions were presented randomly for each participant, followed by their literal definition (consisting of the GoEmotions dataset[64]) and the nine-point Likert scale for both attributes. There was text instruction above each rating scale, "To what extent does $[EMOTION]$ make you feel... (Valence: 1=very unpleasant, 5=neutral, 9=very pleasant; Arousal: 1=very calming, 9=very arousing)".

For basic emotions attributes, another thirty participants (17 females, mean age = 31 years) were recruited to rate the physiological response consistency of all emotion concepts and the facial expressions of six basic emotions. Participants also rated attributes on a nine-point Likert scale, and the order of emotions (with definition) and attributes were randomized for each participant. There was text and image instruction above each rating scale, "To what extent is $[EMOTION]$ consistent with the physiological responses shown in the figures: (1=very inconsistent, 5=neutral, 9=very consistent)". The images we used to indicate basic emotion expressions are twelve averaged facial expressions (one male and one female for each basic emotion) from the AKDEF stimulus set (https://kdef.se/index.html)[133].

For appraisal attributes, two hundred ninety-nine independent participants (148 female, mean age = 37 years) were recruited to recall an event that directly caused them to feel a given emotion (randomly assigned) and rate 38 items on the event. In the recall phase, we instructed participants to remember and write down a situation (at least 100 words) in which they felt the given emotion (with definition) and then identify the specific event (up to 50 words) in the situation that directly caused that emotion. This procedure avoided involving multiple events, cognitions, and emotions in a single recall[134]. We instructed participants in the next phase to rate 38 items for that specific event in random order. All those items were summarized by ref.[135], covering most factors from the appraisal theories of emotion. Before the next processing step, we kept six factors in the 299 events times 38 items matrix as appraisal attributes (see **Supplementary Figure 5** for further details).

A total of 14 candidate conceptual attributes (two affective, six appraisals, and six basic emotions) were averaged across repeated ratings as the final attributes score for each emotion concept (their reliabilities are shown in **Supplementary Figure 6**). In addition to the adopted participants reported above, we excluded 5, 6, and 30 subjects from the three surveys due to failure of the attention check, respectively. The Institutional Review Board at the Department of Psychology, Tsinghua University, approved all experimental procedures. All participants gave their informed consent.

**RSA for comparing LLM-based and human emotion knowledge.**

To show the similarity of the conceptual structure of emotions obtained from the LLM and humans, we related neuron activation RDM, prompt embedding RDM, and word embedding RDM to similarity judgment RDM via a two-sided signed-rank test, using only the lower triangle, with bootstrap sampling participants and emotions 100,000 times.

To explain the LLM-based representation of emotion knowledge, we build RDMs for 14 human conceptual attributes of emotions (Euclidean distance of final score) and related them to neuron activation RDM via a two-sided signed-rank test, bootstrap sampling participants and emotions 100,000 times. False discovery rate (FDR)[136] was corrected to control multiple comparisons.

The second-order similarity measurement was considered more feasible here, as the second-order conceptual similarity among different emotions would be much easier to obtain (through behavioral experiments) than the first-order, absolute conceptual representations.

**Searchlight RSA for LLM's representation of human conceptual attributes of emotions.**

Following the same procedure of building RDMs for conceptual attributes, we first build RDMs for every single artificial neuron (Euclidean distance of average neuron activation). We then conducted the one-sided sign-rank test to indicate the correspondences between each RDM of artificial neuron and each RDM of conceptual attribute. Due to many calculations, we did not perform the bootstrap method. Instead, we used FDR correction to control multiple comparisons. The significant Kendall's *tau* values show the absolute correspondence of the artificial neurons for each conceptual attribute, and the rank of Kendall's *tau* values shows the relative correspondences (**Supplement Figure 1**).

**Attribute-specific neurons manipulation experiment.**

The examine the potential support of emotion knowledge for discrete emotion inference, we input the trained prompts and task texts into LLM to infer emotions from the texts in the test set. During the inference of 27 emotions, we modified the activation values of attribute-specific neurons to zero. For each conceptual attribute, the number of manipulated neurons was set uniformly to 1,500, 2,000, 2,500, 3,000, 4,000, 5,000, or 6,000. Overall, the selective manipulation operation was repeated 34,020 times, respectively, for 14 conceptual attributes, 27 emotion inference tasks (12 prompts/random seeds per task), and seven numbers manipulated neurons.

To exclude the influence of manipulating neurons *per se*, we randomly select the same number of neurons to manipulate as a control group for every operation. The causal contribution of conceptual attributes was indicated as the difference in accuracy after selective manipulation compared to accuracy after random manipulation

**Correlating performance change on the LLM with the effectiveness of representations on the human side.**

First, we calculated the representational similarity of arbitrary conceptual attributes with overall human emotion knowledge by relating its RDM to similarity judgment RDM via a two-sided signed-rank test, bootstrap sampling participants and emotions 100,000 times.

Then, we used the effectiveness of representations of 14 conceptual attributes on the human side (Kendall's *tau*) to predict their contribution to emotion inference on the LLM (accuracy difference). Considering the weak heterogeneity across emotions, we treated each emotion inference task as a sample and obtained a series of Pearson's correlation

coefficients to determine significance by a one-sided $t$-test on Fisher's transformed coefficients[137]. See **Figure 5d** for an illustration.

# Data availability

All behavioral data and pre-processed computational data (neuron activation values, prompt embedding, and word embedding of 27 emotions) are available via our source codes page.

# Code availability

We have opened the source codes. Readers can refer to this link (https://github.com/thunlp/Model_Emotion) for details of the implementation.

# References

1.      Hancock, J. T., Landrigan, C. & Silver, C. Expressing emotion in text-based communication. *Conference on Human Factors in Computing Systems - Proceedings* 929–932 (2007) doi:10.1145/1240624.1240764.

2.      Gill, A. J., Gergle, D., French, R. M. & Oberlander, J. Emotion rating from short blog texts. *Conference on Human Factors in Computing Systems - Proceedings* 1121–1124 (2008) doi:10.1145/1357054.1357229.

3.      Aman, S. & Szpakowicz, S. Identifying expressions of emotion in text. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* **4629 LNAI**, 196–205 (2007).

4.      Barrett, L. F., Adolphs, R., Marsella, S., Martinez, A. M. & Pollak, S. D. Emotional Expressions Reconsidered : Challenges to Inferring Emotion From Human Facial Movements. *Psychological Science in the Public Interest* 46–53 (2019) doi:10.1177/1529100619832930.

5.      Zaki, J., Bolger, N. & Ochsner, K. Unpacking the Informational Bases of Empathic Accuracy. *Emotion* **9**, 478–487 (2009).

6.      Chen, Z. & Whitney, D. Tracking the affective state of unseen persons. *Proc Natl Acad Sci U S A* **116**, 7559–7564 (2019).

7.      Spunt, R. P. & Adolphs, R. The neuroscience of understanding the emotions of others. *Neurosci Lett* **693**, 44–48 (2019).

8.      Skerry, A. E. & Saxe, R. A common neural code for perceived and inferred emotion. *Journal of Neuroscience* **34**, 15997–16008 (2014).

9.      Lindquist, K. A. Emotions emerge from more basic psychological ingredients: A modern psychological constructionist model. *Emotion Review* **5**, 356–368 (2013).

10.     Barrett, L. F. The theory of constructed emotion: an active inference account of

interoception and categorization. *Soc Cogn Affect Neurosci* **12**, 1–23 (2017).

11.  Neimeyer, R. A., Klass, D. & Dennis, M. R. A Social Constructionist Account of Grief: Loss and the Narration of Meaning. *Death Stud* **38**, 485–498 (2014).

12.  Bonanno, G. A. *et al.* Resilience to loss and chronic grief: A prospective study from preloss to 18-months postloss. *J Pers Soc Psychol* **83**, 1150–1164 (2002).

13.  Brooks, J. A. & Freeman, J. B. Conceptual knowledge predicts the representational structure of facial emotion perception. *Nat Hum Behav* **2**, 581–591 (2018).

14.  Brooks, J. A., Chikazoe, J., Sadato, N. & Freeman, J. B. The neural representation of facial-emotion categories reflects conceptual structure. *Proc Natl Acad Sci U S A* **116**, 15861–15870 (2019).

15.  Binetti, N. *et al.* Genetic algorithms reveal profound individual differences in emotion recognition. *Proc Natl Acad Sci U S A* **119**, (2022).

16.  Hu, X., Wang, F. & Zhang, D. Similar brains blend emotion in similar ways: Neural representations of individual difference in emotion profiles. *Neuroimage* 118819 (2021) doi:10.1016/j.neuroimage.2021.118819.

17.  Gendron, M., Roberson, D., van der Vyver, J. M. & Barrett, L. F. Perceptions of emotion from facial expressions are not culturally universal: Evidence from a remote culture. *Emotion* **14**, 251–262 (2014).

18.  Jack, R. E., Garrod, O. G. B., Yu, H., Caldara, R. & Schyns, P. G. Facial expressions of emotion are not culturally universal. *Proc Natl Acad Sci U S A* **109**, 7241–7244 (2012).

19.  Matthews, C. M., Thierry, S. M. & Mondloch, C. J. Recognizing, Discriminating, and Labeling Emotional Expressions in a Free-Sorting Task: A Developmental Story. *Emotion* **22**, 945–953 (2022).

20.  Grosse, G., Streubel, B., Gunzenhauser, C. & Saalbach, H. Let's Talk About Emotions: the Development of Children's Emotion Vocabulary from 4 to 11 Years of Age. *Affect Sci* **2**, 150–162 (2021).

21.  Dunn, J., Brown, J. & Beardsall, L. Family Talk about Feeling States and Children. *Dev Psychol* **27**, 448–455 (1991).

22.  Widen, S. C. & Russell, J. A. Children acquire emotion categories gradually. *Cogn Dev* **23**, 291–312 (2008).

23.  Streubel, B., Gunzenhauser, C., Grosse, G. & Saalbach, H. Emotion-specific vocabulary and its contribution to emotion understanding in 4- to 9-year-old children. *J Exp Child Psychol* **193**, 104790 (2020).

24.  Hoemann, K., Xu, F. & Barrett, L. F. Emotion Words, Emotion Concepts, and Emotional Development in Children: A Constructionist Hypothesis. *Dev Psychol* **55**, 1830–1849 (2019).

25.  Hoemann, K. *et al.* Developing an Understanding of Emotion Categories: Lessons from Objects. *Trends Cogn Sci* **24**, 39–51 (2020).

26.  Lindquist, K. A. The role of language in emotion: existing evidence and future directions. *Curr Opin Psychol* **17**, 135–139 (2017).

27.  Maxfield, L. Attention and Semantic Priming: A Review of Prime Task Effects. *Conscious Cogn* **6**, 204–218 (1997).

28.  Lindquist, K. A., Satpute, A. B. & Gendron, M. Does Language Do More Than Communicate Emotion? *Curr Dir Psychol Sci* **24**, 99–108 (2015).

29. Carroll, N. C. & Young, A. W. Priming of emotion recognition. *Quarterly Journal of Experimental Psychology Section A: Human Experimental Psychology* **58**, 1173–1197 (2005).

30. Nook, E. C., Lindquist, K. A. & Zaki, J. A new look at emotion perception: Concepts speed and shape facial emotion recognition. *Emotion* **15**, 569–578 (2015).

31. Gendron, M., Lindquist, K. A., Barsalou, L. & Barrett, L. F. Emotion words shape emotion percepts. *Emotion* **12**, 314–325 (2012).

32. Lindquist, K. A., Barrett, L. F., Bliss-Moreau, E. & Russell, J. A. Language and the perception of emotion. *Emotion* **6**, 125–138 (2006).

33. Firestone, C. & Scholl, B. J. 'Top-Down' Effects Where None Should Be Found: The El Greco Fallacy in Perception Research. *Psychol Sci* **25**, 38–46 (2014).

34. Firestone, C. & Scholl, B. J. Cognition does not affect perception: Evaluating the evidence for top-down effects. *Behavioral and Brain Sciences* **39**, 1–77 (2015).

35. Lindquist, K. A., Gendron, M., Barrett, L. F. & Dickerson, B. C. Emotion perception, but not affect perception, is impaired with semantic memory loss. *Emotion* **14**, 375–387 (2014).

36. Jastorff, J. *et al.* Functional dissociation between anterior temporal lobe and inferior frontal gyrus in the processing of dynamic body expressions: Insights from behavioral variant frontotemporal dementia. *Hum Brain Mapp* **37**, 4472–4486 (2016).

37. Long, Y. *et al.* Transcranial direct current stimulation of the right anterior temporal lobe changes interpersonal neural synchronization and shared mental processes. *Brain Stimul* **16**, 28–39 (2023).

38. Patterson, K., Nestor, P. J. & Rogers, T. T. Where do you know what you know? The representation of semantic knowledge in the human brain. *Nat Rev Neurosci* **8**, 976–987 (2007).

39. Hodges, J. R. & Patterson, K. Semantic dementia: a unique clinicopathological syndrome. *Lancet Neurology* vol. 6 1004–1014 Preprint at https://doi.org/10.1016/S1474-4422(07)70266-1 (2007).

40. Scarantino, A. Functional specialization does not require a one-to-one mapping between brain regions and emotions. *Behavioral and Brain Sciences* **35**, 161–162 (2012).

41. Lindquist, K. A., Wager, T. D., Bliss-Moreau, E., Kober, H. & Barrett, L. F. Authors' response: what are emotions and how are they created in the brain? *Behav Brain Sci* **35**, 172–202 (2012).

42. Polanía, R., Nitsche, M. A. & Ruff, C. C. Studying and modifying brain function with non-invasive brain stimulation. *Nat Neurosci* **21**, 174–187 (2018).

43. Shablack, H., Becker, M. & Lindquist, K. A. How do children learn novel emotion words? A study of emotion concept acquisition in preschoolers. *J Exp Psychol Gen* **149**, 1537–1553 (2020).

44. Snefjella, B., Lana, N. & Kuperman, V. How emotion is learned: Semantic learning of novel words in emotional contexts. *J Mem Lang* **115**, 104171 (2020).

45. Nook, E. C., Sasse, S. F., Lambert, H. K., McLaughlin, K. A. & Somerville, L. H. Increasing verbal knowledge mediates development of multidimensional emotion representations. *Nat Hum Behav* **1**, 881–889 (2017).

46.    Liu, Y. *et al.* RoBERTa: A Robustly Optimized BERT Pretraining Approach. (2019).

47.    Brown, T. B. *et al.* Language models are few-shot learners. *Adv Neural Inf Process Syst* **2020-Decem**, (2020).

48.    Devlin, J., Chang, M. W., Lee, K. & Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference* **1**, 4171–4186 (2019).

49.    KEKEKE *et al.* T5: Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* **21**, 1–67 (2020).

50.    Michel, J. B. *et al.* Quantitative analysis of culture using millions of digitized books. *Science (1979)* **331**, 176–182 (2011).

51.    Misra, K., Rayz, J. T. & Ettinger, A. COMPS: Conceptual Minimal Pair Sentences for testing Property Knowledge and Inheritance in Pre-trained Language Models. (2022).

52.    Ettinger, A. What bert is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Trans Assoc Comput Linguist* **8**, 34–48 (2020).

53.    Tenney, I. *et al.* What do you learn from context? Probing for sentence structure in contextualized word representations. *7th International Conference on Learning Representations, ICLR 2019* 1–17 (2019).

54.    Da, J. & Kasai, J. Understanding Commonsense Inference Aptitude of Deep Contextual Representations. in 1–12 (2019). doi:10.18653/v1/d19-6001.

55.    Aspillaga, C., Mendoza, M. & Soto, A. Inspecting the concept knowledge graph encoded by modern language models. *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021* 2984–3000 (2021) doi:10.18653/v1/2021.findings-acl.263.

56.    Schramowski, P., Turan, C., Andersen, N., Rothkopf, C. A. & Kersting, K. Large pre-trained language models contain human-like biases of what is right and wrong to do. *Nat Mach Intell* **4**, 258–268 (2022).

57.    Rogers, A., Kovaleva, O. & Rumshisky, A. A primer in bertology: What we know about how bert works. *Trans Assoc Comput Linguist* **8**, 842–866 (2020).

58.    Hutchinson, B. *et al.* Social Biases in NLP Models as Barriers for Persons with Disabilities. 5491–5501 (2020) doi:10.18653/v1/2020.acl-main.487.

59.    Kurita, K., Vyas, N., Pareek, A., Black, A. W. & Tsvetkov, Y. Measuring Bias in Contextualized Word Representations. 166–172 (2019) doi:10.18653/v1/w19-3823.

60.    Basta, C., Costa-jussà, M. R. & Casas, N. Evaluating the Underlying Gender Bias in Contextualized Word Embeddings. 33–39 (2019) doi:10.18653/v1/w19-3805.

61.    Wang, X. *et al.* Finding Skill Neurons in Pre-trained Transformer-based Language Models. in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing* 11132--11152 (Association for Computational Linguistics, 2022).

62.    Walsh, V. & Cowey, A. Transcranial magnetic stimulation and cognitive neuroscience. *Nat Rev Neurosci* **1**, 73–80 (2000).

63.    Hallett, M. Transcranial Magnetic Stimulation: A Primer. *Neuron* **55**, 187–199 (2007).

64.    Demszky, D. *et al.* GoEmotions: A Dataset of Fine-Grained Emotions. in 4040–4054 (2020). doi:10.18653/v1/2020.acl-main.372.

65. Liu, P. *et al.* Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. 1–46 (2021).

66. Su, Y. *et al.* On Transferability of Prompt Tuning for Natural Language Processing. in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* 3949–3969 (Association for Computational Linguistics, 2022). doi:10.18653/v1/2022.naacl-main.290.

67. Kriegeskorte, N., Mur, M. & Bandettini, P. Representational similarity analysis - connecting the branches of systems neuroscience. *Front Syst Neurosci* **2**, 1–28 (2008).

68. Qin, Y. *et al.* Exploring Universal Intrinsic Task Subspace via Prompt Tuning. (2021).

69. Pennington, J., Socher, R. & Manning, C. D. GloVe: Global Vectors for Word Representation. in *Empirical Methods in Natural Language Processing (EMNLP)* 1532–1543 (2014).

70. Günther, F., Rinaldi, L. & Marelli, M. Vector-Space Models of Semantic Representation From a Cognitive Perspective: A Discussion of Common Misconceptions. *Perspectives on Psychological Science* **14**, 1006–1033 (2019).

71. Gu, M; Lai, T. L. The Dip Test of Unimodality. *The Annuals of Statistics* **13**, 70–84 (1985).

72. Freeman, J. B. & Dale, R. Assessing bimodality to detect the presence of a dual cognitive process. *Behav Res Methods* **45**, 83–97 (2013).

73. Lim, C. R. *et al.* Floor and ceiling effects in the OHS: An analysis of the NHS PROMs data set. *BMJ Open* **5**, (2015).

74. Bausch, M. *et al.* Concept neurons in the human medial temporal lobe flexibly represent abstract relations between concepts. *Nat Commun* **12**, 1–12 (2021).

75. Schacter, D. L., Dobbins, I. G. & Schnyer, D. M. Specificity of priming: A cognitive neuroscience perspective. *Nat Rev Neurosci* **5**, 853–862 (2004).

76. Race, E. A., Shanker, S. & Wagner, A. D. Neural priming in human frontal cortex: Multiple forms of learning reduce demands on the prefrontal executive system. *J Cogn Neurosci* **21**, 1766–1781 (2009).

77. Goldstein, A. *et al.* Shared computational principles for language processing in humans and deep language models. *Nat Neurosci* **25**, 369–380 (2022).

78. Caliskan, A., Bryson, J. J. & Arvind Narayanan. Semantics derived automatically from language corpora contain contain human-like biases. *Science (1979)* **356**, 183–186 (2017).

79. Zhou, L., Yang, A., Meng, M. & Zhou, K. Emerged human-like facial expression representation in a deep convolutional neural network. *Sci Adv* **8**, 1–12 (2022).

80. Grand, G., Blank, I. A., Pereira, F. & Fedorenko, E. Semantic projection: recovering human knowledge of multiple, distinct object features from word embeddings. *Nat Hum Behav* 287–300 (2018) doi:https://doi.org/10.1038/s41562-022-01316-8.

81. Caucheteux, C. & King, J. R. Brains and algorithms partially converge in natural language processing. *Commun Biol* **5**, (2022).

82. Schrimpf, M. *et al.* The neural architecture of language: Integrative modeling converges on predictive processing. *bioRxiv* 2020.06.26.174482 (2021).

83. Ganguli, D. *et al.* The Capacity for Moral Self-Correction in Large Language Models. 1–20 (2023).

84. Kosinski, M. Theory of Mind May Have Spontaneously Emerged in Large Language Models. (2023).

85. Colón, Y. I., Castillo, C. D. & O' Toole, A. J. Facial expression is retained in deep networks trained for face identification. *J Vis* **21**, 1–10 (2021).

86. Racanière, S. *et al.* Attention is all you need. *Adv Neural Inf Process Syst* **2017-Decem**, 1–14 (2017).

87. Shablack, H., Stein, A. G. & Lindquist, K. A. Comment: A role of Language in Infant Emotion Concept Acquisition. *Emotion Review* **12**, 251–253 (2020).

88. Hoemann, K., Devlin, M. & Barrett, L. F. Comment: Emotions Are Abstract, Conceptual Categories That Are Learned by a Predicting Brain. *Emotion Review* **12**, 253–255 (2020).

89. Ruba, A. L. & Repacholi, B. M. Do Preverbal Infants Understand Discrete Facial Expressions of Emotion? *Emotion Review* **12**, 235–250 (2020).

90. Ruba, A. L. & Repacholi, B. M. Beyond Language in Infant Emotion Concept Development. *Emotion Review* **12**, 255–258 (2020).

91. Barrett, L. F. Valence is a basic building block of emotional life. *J Res Pers* **40**, 35–55 (2006).

92. Paul, E. & Davidson, R. J. The Duchenne Smile: Emotional Expression and Brain Physiology II. *J Pers Soc Psychol* **58**, 342–353 (1990).

93. Singer, T. *et al.* Empathic neural responses are modulated by the perceived fairness of others. *Nature* **439**, 466–469 (2006).

94. Huth, A. G., de Heer, W. A., Griffiths, T. L., Theunissen, F. E. & Gallant, J. L. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature* **532**, 453–458 (2016).

95. Pereira, F. *et al.* Toward a universal decoder of linguistic meaning from brain activation. *Nat Commun* **9**, (2018).

96. Anderson, A. J. *et al.* Deep artificial neural networks reveal a distributed cortical network encoding propositional sentence-level meaning. *Journal of Neuroscience* **41**, 4100–4119 (2021).

97. Barrett, L. F. Are Emotions Natural Kinds? *Perspectives on Psychological Science* **1**, 28–58 (2006).

98. Adolphs, R., Mlodinow, L. & Barrett, L. F. What is an emotion? *Current Biology* **29**, R1060–R1064 (2019).

99. D' Arms, J. & Samuels, R. Could Emotion Development Really Be the Acquisition of Emotion Concepts? *Dev Psychol* **55**, 2015–2019 (2019).

100. Keltner, D., Tracy, J. L., Sauter, D. & Cowen, A. What Basic Emotion Theory Really Says for the Twenty-First Century Study of Emotion. *J Nonverbal Behav* **43**, 195–201 (2019).

101. Fridlund, A. J. The behavioral ecology view of facial displays, 25 years later. *Oxf Ser Soc Cogn Soc Neurosci* 77–92 (2017).

102. Adolphs, R. & Andler, D. Investigating Emotions as Functional States Distinct From Feelings. *Emotion Review* **10**, 191–201 (2018).

103. Cowen, A. S. & Keltner, D. Semantic Space Theory: A Computational Approach to Emotion. *Trends in Cognitive Sciences* vol. 25 124–136 Preprint at https://doi.org/10.1016/j.tics.2020.11.004 (2021).

104. Barrett, L. F. & Westlin, C. Navigating the science of emotion. in *Emotion Measurement*

39–84 (Elsevier, 2021). doi:10.1016/b978-0-12-821124-3.00002-8.

105. Horikawa, T., Cowen, A. S., Keltner, D. & Kamitani, Y. The Neural Representation of Visually Evoked Emotion Is High-Dimensional, Categorical, and Distributed across Transmodal Brain Regions. *iScience* **23**, 101060 (2020).

106. Lindquist, K. A., Wager, T. D., Kober, H., Bliss-Moreau, E. & Barrett, L. F. The brain basis of emotion: A meta-analytic review. *Behavioral and Brain Sciences* **35**, 121–143 (2012).

107. Toneva, M., Mitchell, T. M. & Wehbe, L. Combining computational controls with natural text reveals new aspects of meaning composition. *bioRxiv* **2**, 2020.09.28.316935 (2022).

108. Frisby, S. L., Halai, A. D., Cox, C. R., Ralph, M. A. L. & Rogers, T. T. Decoding semantic representations in mind and brain. *Trends Cogn Sci* **xx**, 1–24 (2022).

109. Bi, Y. Dual coding of knowledge in the human brain. *Trends Cogn Sci* **25**, 883–895 (2021).

110. Schirmer, A. & Adolphs, R. Emotion Perception from Face, Voice, and Touch: Comparisons and Convergence. *Trends Cogn Sci* **21**, 216–228 (2017).

111. Kim, J. *et al.* Abstract representations of associated emotions in the human brain. *Journal of Neuroscience* **35**, 5655–5663 (2015).

112. Lim, S. L., O'Doherty, J. P. & Rangel, A. Stimulus value signals in ventromedial PFC reflect the integration of attribute value signals computed in fusiform gyrus and posterior superior temporal gyrus. *Journal of Neuroscience* **33**, 8729–8741 (2013).

113. Escoffier, N., Zhong, J., Schirmer, A. & Qiu, A. Emotional expressions in voice and music: Same code, same effect? *Hum Brain Mapp* **34**, 1796–1810 (2013).

114. Schirmer, A. & Kotz, S. A. Beyond the right hemisphere: brain mechanisms mediating vocal emotional processing. *Trends Cogn Sci* **10**, 24–30 (2006).

115. Sloutsky, V. M. From Perceptual Categories to Concepts: What Develops? *Cogn Sci* **34**, 1244–1286 (2010).

116. Westermann, G. & Mareschal, D. From perceptual to language-mediated categorization. *Philosophical Transactions of the Royal Society B: Biological Sciences* **369**, (2014).

117. Kousta, S. T., Vigliocco, G., Vinson, D. P., Andrews, M. & del Campo, E. The Representation of Abstract Words: Why Emotion Matters. *J Exp Psychol Gen* **140**, 14–34 (2011).

118. Borghi, A. M. *et al.* The challenge of abstract concepts. *Psychol Bull* **143**, 263–292 (2017).

119. Binder, J. R. In defense of abstract conceptual representations. *Psychon Bull Rev* **23**, 1096–1108 (2016).

120. Pulvermüller, F. How neurons make meaning: Brain mechanisms for embodied and abstract-symbolic semantics. *Trends Cogn Sci* **17**, 458–470 (2013).

121. Vukovic, N., Feurra, M., Shpektor, A., Myachykov, A. & Shtyrov, Y. Primary motor cortex functionally contributes to language comprehension: An online rTMS study. *Neuropsychologia* **96**, 222–229 (2017).

122. Kashdan, T. B., Barrett, L. F. & McKnight, P. E. Unpacking Emotion Differentiation: Transforming Unpleasant Experience by Perceiving Distinctions in Negativity. *Curr Dir Psychol Sci* **24**, 10–16 (2015).

123. Delgado, R. & Tibau, X. A. Why Cohen's Kappa should be avoided as performance

measure in classification. *PLoS One* **14**, 1–26 (2019).

124. Geva, M., Schuster, R., Berant, J. & Levy, O. Transformer Feed-Forward Layers Are Key-Value Memories. in *EMNLP 2021 - 2021 Conference on Empirical Methods in Natural Language Processing, Proceedings* 5484–5495 (2021). doi:10.18653/v1/2021.emnlp-main.446.

125. Dai, D. *et al.* Knowledge Neurons in Pretrained Transformers. in 8493–8502 (2022). doi:10.18653/v1/2022.acl-long.581.

126. Firth, J. R. A Synopsis of Linguistic Theory 1930-55. *Studies in Linguistic Analysis: Special Volume of the Philological Society* 1–32 (1957).

127. Barrett, L. F. Solving the emotion paradox: Categorization and the experience of emotion. *Personality and Social Psychology Review* vol. 10 20–46 Preprint at https://doi.org/10.1207/s15327957pspr1001_2 (2006).

128. Russell, J. A. Core Affect and the Psychological Construction of Emotion. *Psychol Rev* **110**, 145–172 (2003).

129. Scherer, K. R. & Fontaine, J. R. J. The semantic structure of emotion words across languages is consistent with componential appraisal models of emotion. *Cogn Emot* **33**, 673–682 (2019).

130. Clore, G. & Ortony, A. Psychological Construction in the OCC Model of Emotion Gerald. **5**, 335–343 (2016).

131. Levenson, R. W. Basic emotion questions. *Emotion Review* **3**, 379–386 (2011).

132. Du, S., Tao, Y. & Martinez, A. M. Compound facial expressions of emotion. *Proc Natl Acad Sci U S A* **111**, (2014).

133. Lundqvist, D. & Flykt, A. The Averaged Karolinska Directed Emotional Faces - AKDEF. *CD ROM from Department of Clinical Neuroscience, Psychology section, Karolinska Institutet* Preprint at (1998).

134. Roseman, I. J., Spindel, M. S. & Jose, P. E. Appraisals of Emotion-Eliciting Events: Testing a Theory of Discrete Emotions. *J Pers Soc Psychol* **59**, 899–915 (1990).

135. Skerry, A. E. & Saxe, R. Neural Representations of Emotion Are Organized around Abstract Event Features. *Current Biology* **25**, 1945–1954 (2015).

136. Benjamini, Y. & Yekutieli, D. The control of the false discovery rate in multiple testing under dependency. *Ann Stat* **29**, 1165–1188 (2001).

137. Silver, N. C. & Dunlap, W. P. Averaging Correlation Coefficients: Should Fisher's z Transformation Be Used? *Journal of Applied Psychology* **72**, 146–148 (1987).

# Acknowledgements

**Table 1 | The raw accuracy of 27 emotion inference tasks.** We trained 27 emotion-specific prompts on the training set with 12 random seeds. We reported their average accuracies on the test set and the rater agreement on dataset annotations. The LLM's inference accuracy is significantly related to the agreement of human annotations with Pearson's $r(25) = .631$, $p$-value < .001.

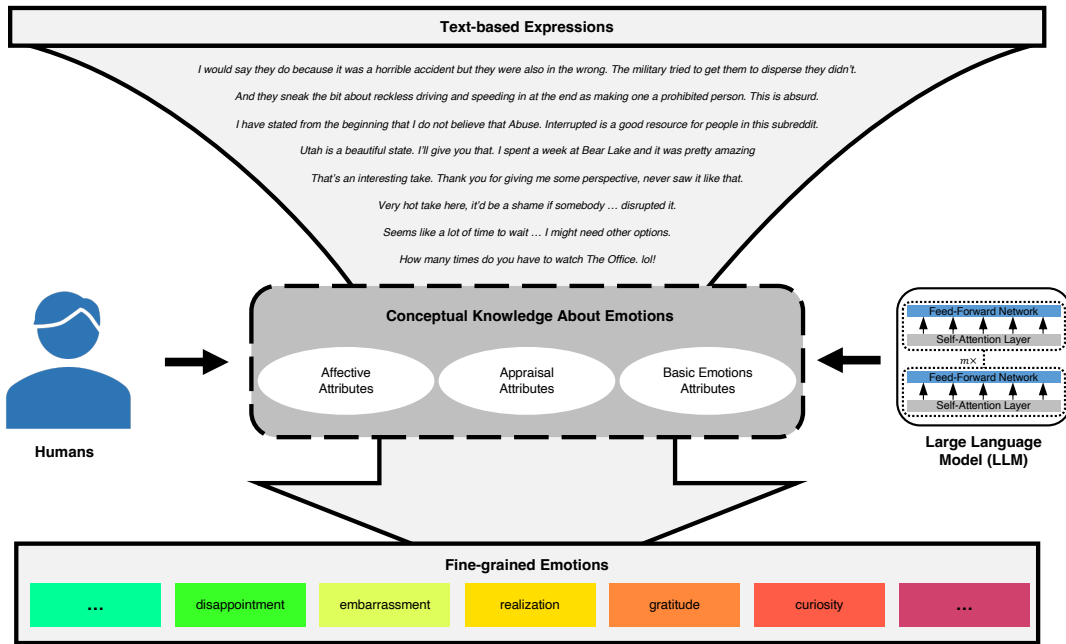| Emotion | Agreement of Human Annotations | the LLM's Inference Accuracy (%) | |
|---|---|---|---|
| | | Mean | S.D. |
| admiration | .535 | 89.8 | 0.5 |
| amusement | .482 | 94.1 | 1.2 |
| anger | .207 | 85.2 | 1.7 |
| annoyance | .193 | 77.1 | 2.5 |
| approval | .385 | 75.8 | 3.1 |
| caring | .237 | 86.4 | 4.6 |
| confusion | .217 | 83.3 | 5.4 |
| curiosity | .418 | 90.9 | 1.2 |
| desire | .177 | 87.0 | 3.6 |
| disappointment | .186 | 73.3 | 3.8 |
| disapproval | .274 | 76.5 | 3.9 |
| disgust | .192 | 85.0 | 8.6 |
| embarrassment | .177 | 79.2 | 12.7 |
| excitement | .193 | 83.2 | 3.6 |
| fear | .266 | 85.0 | 13.5 |
| gratitude | .645 | 96.4 | 0.4 |
| grief | .162 | 81.7 | 14.0 |
| joy | .296 | 86.7 | 4.7 |
| love | .446 | 95.5 | 1.6 |
| nervousness | .164 | 71.6 | 11.9 |
| optimism | .322 | 84.9 | 1.4 |
| pride | .163 | 82.5 | 10.5 |
| realization | .194 | 68.0 | 5.7 |
| relief | .172 | 68.5 | 8.7 |
| remorse | .178 | 89.2 | 13.5 |
| sadness | .346 | 82.8 | 7.2 |
| surprise | .275 | 86.2 | 8.1 |

**Figure 1 | Schematic illustration of inferring discrete emotions from text-based expressions.** According to the constructionist accounts of emotion, conceptual knowledge about emotions (i.e., emotion knowledge) causally supports the inference about discrete emotions from text-based expressions. Emotion knowledge includes several conceptual attributes that underlie and differentiate discrete emotions. These can be learned by large language models (LLMs) and represented by LLMs' artificial neurons (hidden states). By manipulating the values of attribute-specific artificial neurons, we demonstrated the causal contribution of emotion knowledge to inferring discrete emotions from text-based expressions.
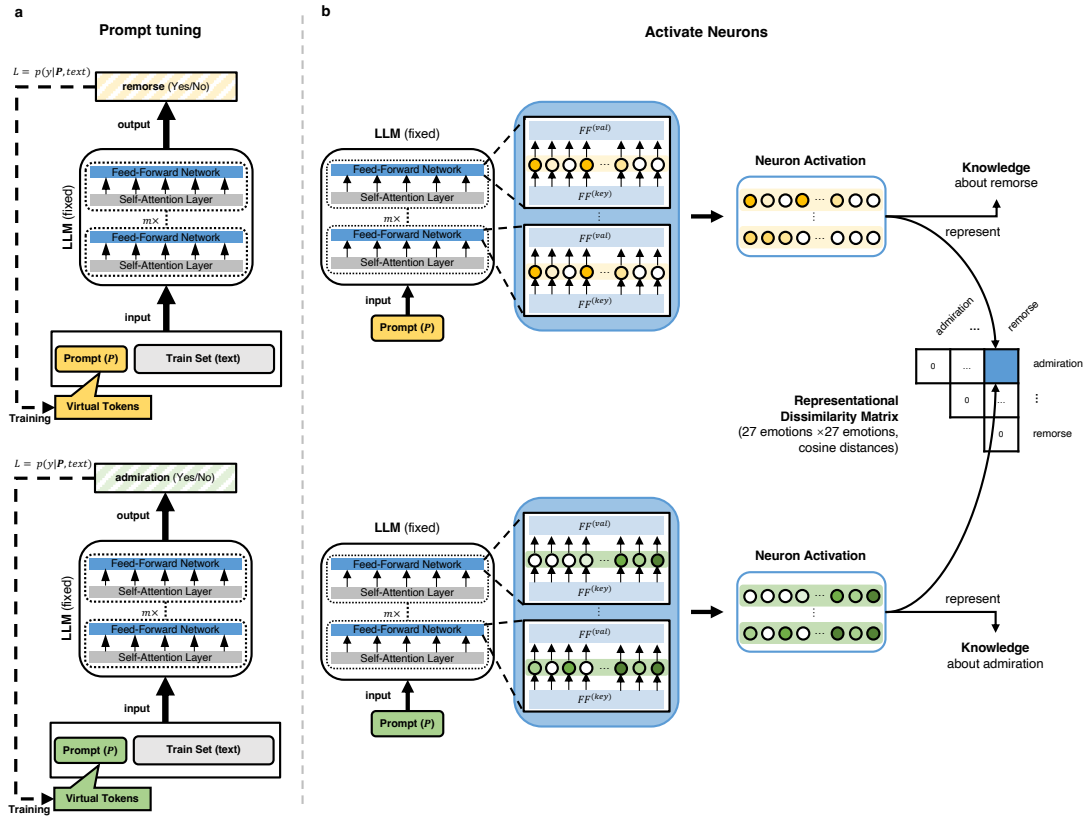
**Figure 2 | Schematic illustration of stimulating emotion knowledge in the LLM.** We adopted the prompt tuning technique to extract the representations of emotion knowledge for 27 discrete emotions in the LLM (RoBERTa). **a**, First, we trained 27 emotion-specific prompts (composed of virtual tokens), each with 12 random seeds on the training set to perform binary emotion inference tasks (yes or no for a specific emotion). During training (prompt tuning), we optimized only prompts while freezing all parameters of the LLM. **b**, Then, we input the trained prompts without concatenating any task text into the LLM to activate 36,864 hidden state (neuron activation) values, which were considered as the LLM-based emotion knowledge representations. These representations were then used for further analysis, including the calculation of pairwise cosine distances between them to obtain the conceptual structure of emotions from the LLM. See **Methods** for details.
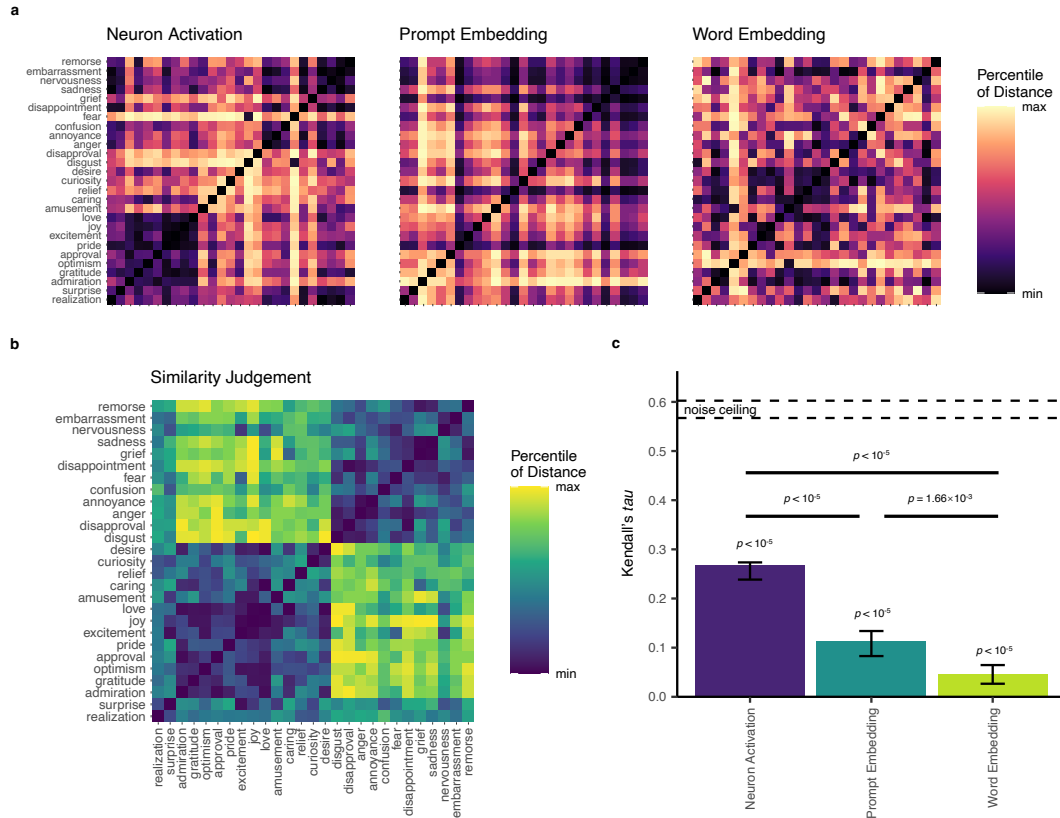
**Figure 3 | The LLM exhibits a human-like conceptual structure of emotions. a**, The representational dissimilarity matrice (RDM) for the conceptual structures of 27 emotions calculated by different language-model indexes, including the LLM (RoBERTa)-based emotion knowledge representation (left), the prompt embedding (middle), and the GloVe word embedding (right). **b**, Averaged human conceptual structure of emotions (RDM) obtained from n = 15 participants by the similarity judgment task. All the above RDMs show the representational distances (dissimilarities) between 27 discrete emotions. For visualization purposes, the values within each matrix are presented in percentile. **c**, The similarity between **a** (language-model RDMs) and **b** (human RDMs), using data from the lower triangles of each RDM. Two-sided signed-rank tests and pairwise comparisons were conducted with bootstrap sampling participants and emotions 100,000 times. Error bars indicate bootstrap 95%CI. The dotted line indicates the noise ceiling of individual RDMs. *P*-values based on bootstrap are reported.
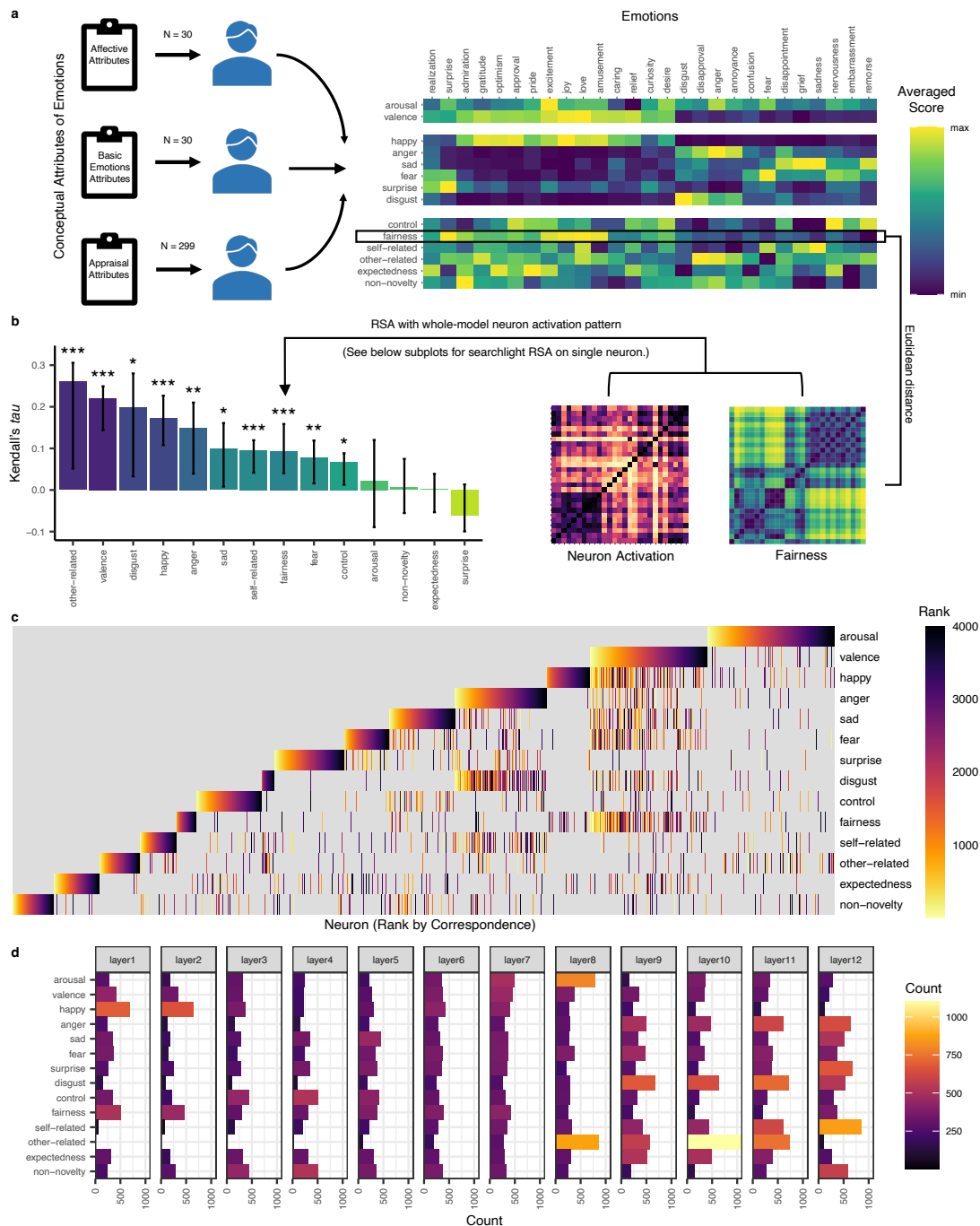
**Figure 4 | The LLM-based emotion knowledge contains human conceptual attributes of emotions. a**, Human ratings for 14 conceptual attributes (two affective, six appraisals, and six basic emotions) of 27 emotions. **b**, The representational similarities between 14 human conceptual attributes of emotions and the LLM-based emotion knowledge. Two-sided signed-rank tests, bootstrap sampling participants and emotions 100,000 times. Error bars indicate bootstrap 95%CI. **c**, Overlap and **d**, Distribution of the top 4,000 relevant artificial neurons for 14 conceptual attributes. The correspondence was computed by RSA (one-sided sign-rank test). **c**, Colorbar indicates the rank of artificial neurons according to *tau* values. **d**, Colorbar indicates the count of artificial neurons in 12 layers of the LLM.
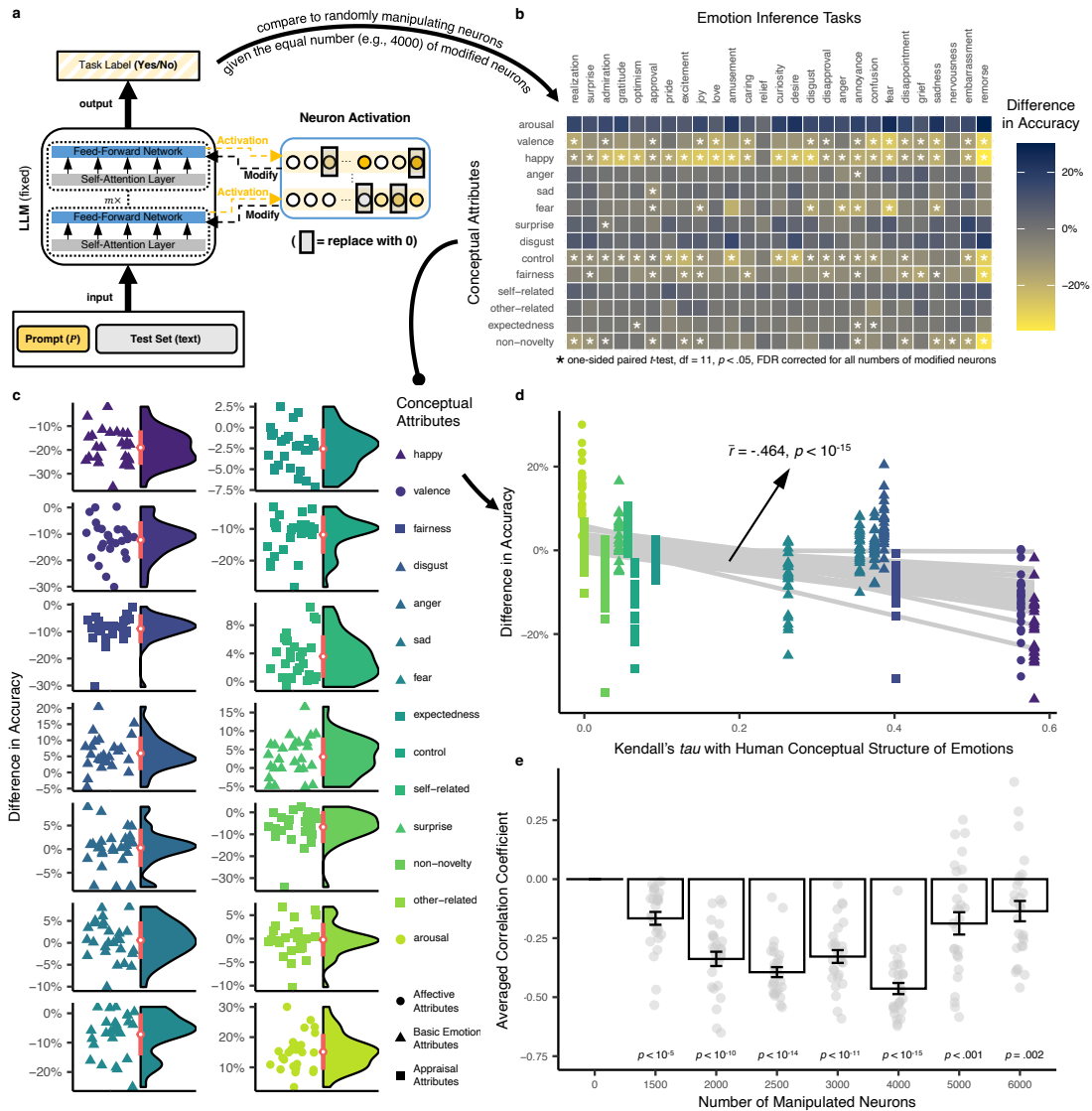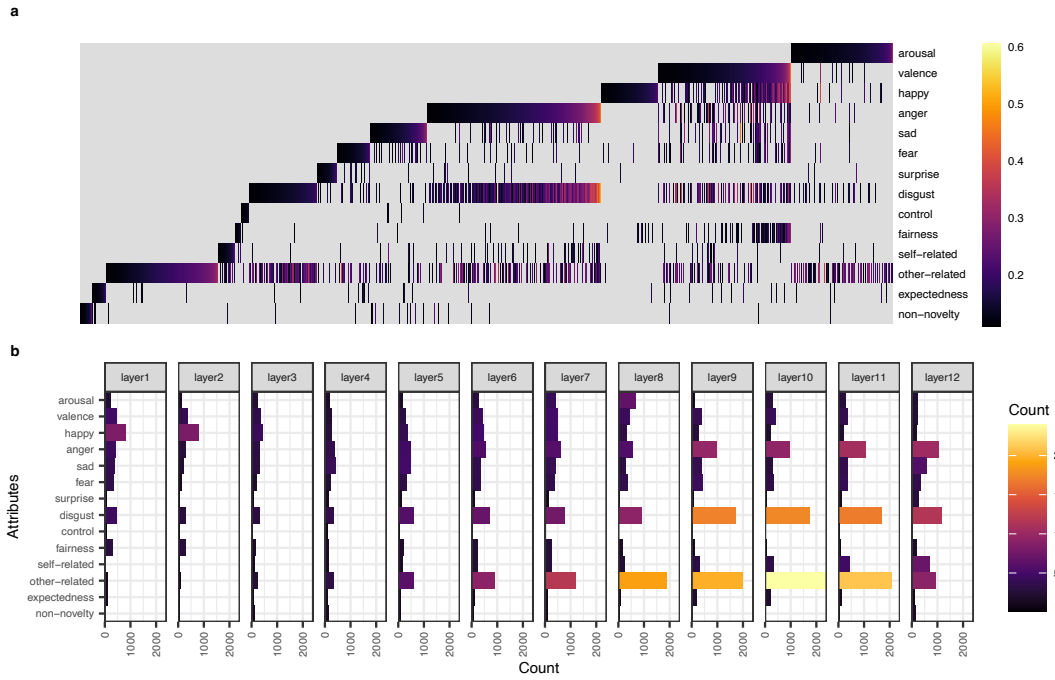
**Figure 5 | Conceptual attributes causally contribute to emotion inference. a,**
Schematic illustration of manipulating attribute-specific artificial neurons while performing
emotion inference tasks. **b,** The causal contributions of 14 conceptual attributes to 27
emotion inference tasks when manipulating the top 4,000 neurons. The brighter the color,
the greater the causal contribution. See **Supplementary Figure 6** for complete results. **c,**
The distribution of arbitrary conceptual attribute's causal contribution across 27 emotion
inference tasks. The white dot and pink error bar indicate the mean and s.d., respectively.
All Hartigan's dip statistics were insignificant at level = .05; see **Supplementary Table 1**
for complete results. **d,e,** Pearson's correlation between 14 conceptual attributes' causal
contributions on arbitrary emotion inference task and their Kendall's *tau* with human
conceptual structure of emotions. One-sided *t*-test on Fisher's transformed coefficients.
For **c,d,** the color of the violin plot and the scatter plot indicate the type of conceptual
attribute.

**Supplementary Table 1 | Heterogeneity degree in conceptual attributes' causal contribution to various emotion inference tasks.** The heterogeneity was tested by determining whether the distribution of the given conceptual attribute's causal contribution (accuracy difference) across 27 emotion inference tasks was multimodal (i.e., at least bimodal) using Hartigan's dip statistic, with unimodality as the alternative hypothesis.
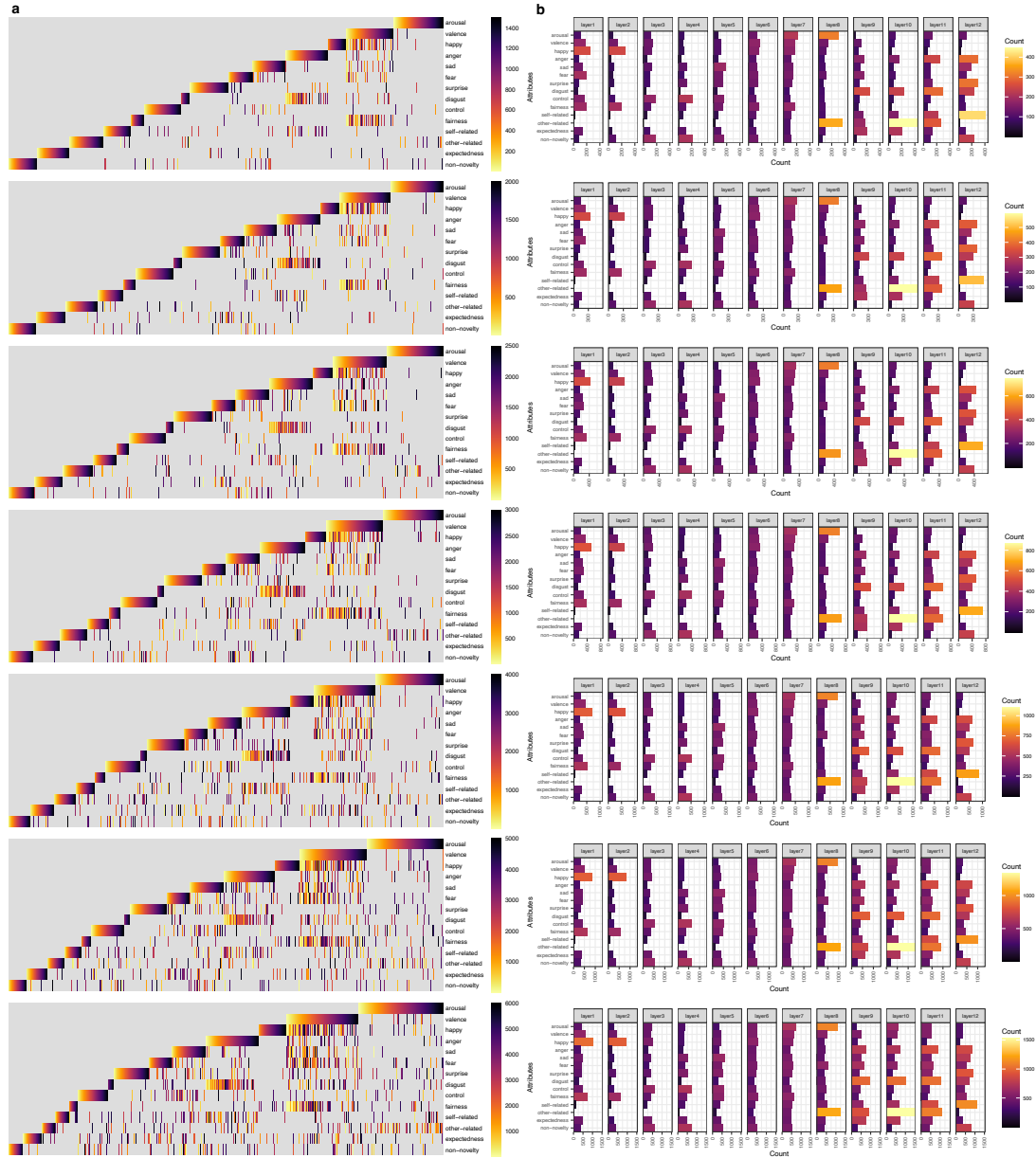
| Number of Manipulated Neurons | Conceptual Attribute | Dip statistic | *p*-value | N |
|---|---|---|---|---|
| 1500 | arousal | 0.066 | 0.468 | 27 |
| 1500 | valence | 0.042 | 0.988 | 27 |
| 1500 | happy | 0.062 | 0.561 | 27 |
| 1500 | anger | 0.047 | 0.944 | 27 |
| 1500 | sad | 0.055 | 0.768 | 27 |
| 1500 | fear | 0.065 | 0.484 | 27 |
| 1500 | surprise | 0.047 | 0.937 | 27 |
| 1500 | disgust | 0.070 | 0.359 | 27 |
| 1500 | control | 0.073 | 0.283 | 27 |
| 1500 | fairness | 0.045 | 0.964 | 27 |
| 1500 | self-related | 0.035 | 0.995 | 27 |
| 1500 | other-related | 0.052 | 0.856 | 27 |
| 1500 | expectedness | 0.043 | 0.979 | 27 |
| 1500 | non-novelty | 0.057 | 0.713 | 27 |
| 2000 | arousal | 0.064 | 0.496 | 27 |
| 2000 | valence | 0.060 | 0.637 | 27 |
| 2000 | happy | 0.062 | 0.578 | 27 |
| 2000 | anger | 0.038 | 0.993 | 27 |
| 2000 | sad | 0.079 | 0.175 | 27 |
| 2000 | fear | 0.059 | 0.666 | 27 |
| 2000 | surprise | 0.054 | 0.799 | 27 |
| 2000 | disgust | 0.049 | 0.906 | 27 |
| 2000 | control | 0.069 | 0.386 | 27 |
| 2000 | fairness | 0.066 | 0.444 | 27 |
| 2000 | self-related | 0.044 | 0.974 | 27 |
| 2000 | other-related | 0.052 | 0.849 | 27 |
| 2000 | expectedness | 0.066 | 0.453 | 27 |
| 2000 | non-novelty | 0.037 | 0.994 | 27 |
| 2500 | arousal | 0.050 | 0.895 | 27 |
| 2500 | valence | 0.063 | 0.531 | 27 |
| 2500 | happy | 0.058 | 0.686 | 27 |
| 2500 | anger | 0.048 | 0.924 | 27 |
| 2500 | sad | 0.043 | 0.982 | 27 |
| 2500 | fear | 0.057 | 0.716 | 27 |

| 2500 | surprise | 0.060 | 0.629 | 27 |
|------|----------|-------|-------|----|
| 2500 | disgust | 0.070 | 0.342 | 27 |
| 2500 | control | 0.057 | 0.709 | 27 |
| 2500 | fairness | 0.054 | 0.798 | 27 |
| 2500 | self-related | 0.052 | 0.843 | 27 |
| 2500 | other-related | 0.049 | 0.911 | 27 |
| 2500 | expectedness | 0.070 | 0.359 | 27 |
| 2500 | non-novelty | 0.049 | 0.915 | 27 |
| 3000 | arousal | 0.047 | 0.948 | 27 |
| 3000 | valence | 0.064 | 0.519 | 27 |
| 3000 | happy | 0.073 | 0.276 | 27 |
| 3000 | anger | 0.061 | 0.612 | 27 |
| 3000 | sad | 0.044 | 0.976 | 27 |
| 3000 | fear | 0.065 | 0.493 | 27 |
| 3000 | surprise | 0.055 | 0.779 | 27 |
| 3000 | disgust | 0.081 | 0.153 | 27 |
| 3000 | control | 0.051 | 0.881 | 27 |
| 3000 | fairness | 0.042 | 0.987 | 27 |
| 3000 | self-related | 0.055 | 0.774 | 27 |
| 3000 | other-related | 0.049 | 0.908 | 27 |
| 3000 | expectedness | 0.056 | 0.737 | 27 |
| 3000 | non-novelty | 0.061 | 0.613 | 27 |
| 4000 | arousal | 0.070 | 0.346 | 27 |
| 4000 | valence | 0.048 | 0.922 | 27 |
| 4000 | happy | 0.083 | 0.127 | 27 |
| 4000 | anger | 0.053 | 0.831 | 27 |
| 4000 | sad | 0.047 | 0.949 | 27 |
| 4000 | fear | 0.056 | 0.758 | 27 |
| 4000 | surprise | 0.079 | 0.177 | 27 |
| 4000 | disgust | 0.054 | 0.811 | 27 |
| 4000 | control | 0.040 | 0.991 | 27 |
| 4000 | fairness | 0.057 | 0.708 | 27 |
| 4000 | self-related | 0.051 | 0.863 | 27 |
| 4000 | other-related | 0.038 | 0.993 | 27 |
| 4000 | expectedness | 0.057 | 0.733 | 27 |
| 4000 | non-novelty | 0.051 | 0.872 | 27 |
| 5000 | arousal | 0.050 | 0.895 | 27 |
| 5000 | valence | 0.049 | 0.909 | 27 |
| 5000 | happy | 0.055 | 0.790 | 27 |
| 5000 | anger | 0.045 | 0.960 | 27 |
| 5000 | sad | 0.058 | 0.705 | 27 |
| 5000 | fear | 0.049 | 0.916 | 27 |

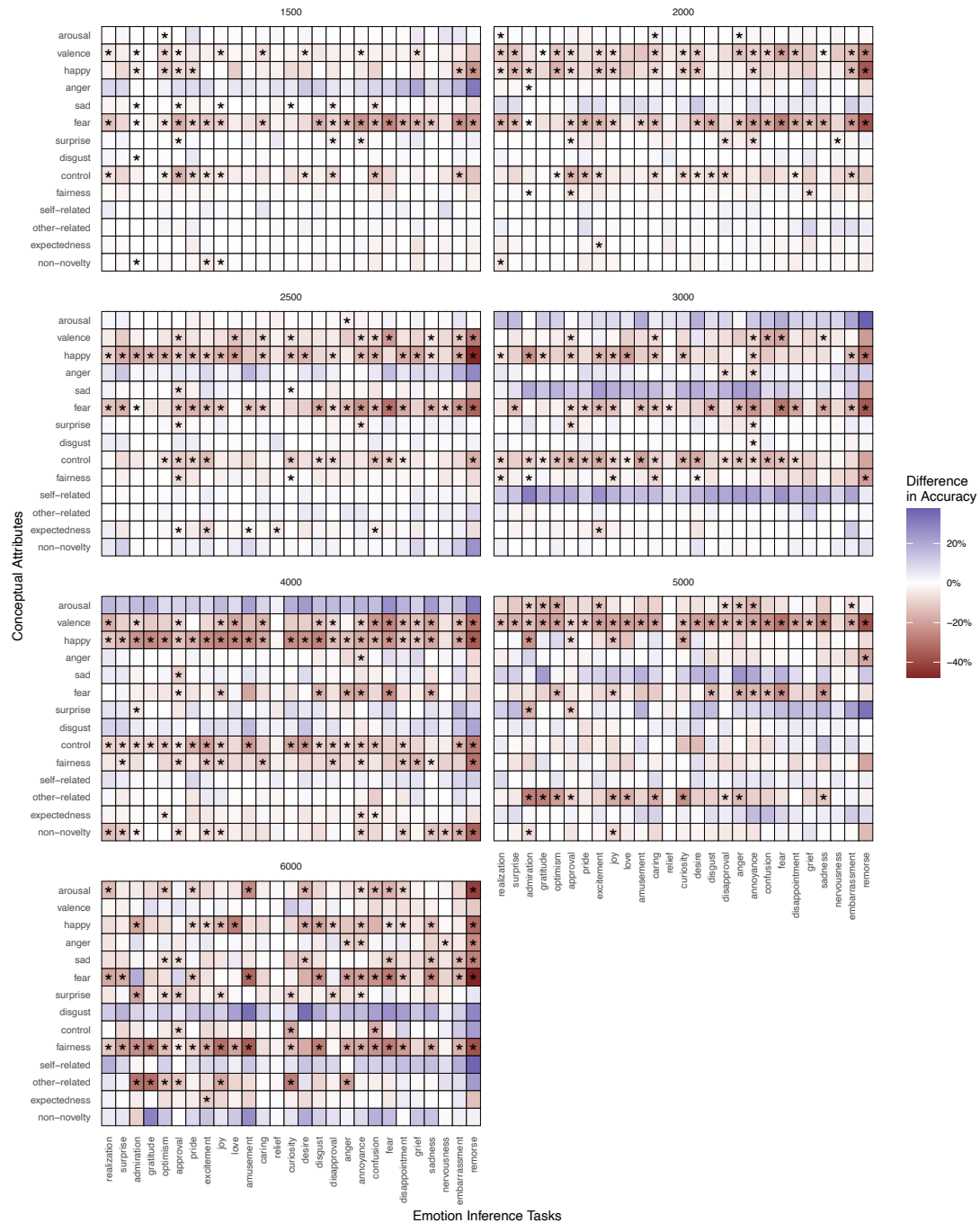| 5000 | surprise | 0.065 | 0.480 | 27 |
|------|----------|-------|-------|-----|
| 5000 | disgust | 0.047 | 0.949 | 27 |
| 5000 | control | 0.049 | 0.908 | 27 |
| 5000 | fairness | 0.059 | 0.661 | 27 |
| 5000 | self-related | 0.066 | 0.466 | 27 |
| 5000 | other-related | 0.048 | 0.921 | 27 |
| 5000 | expectedness | 0.054 | 0.796 | 27 |
| 5000 | non-novelty | 0.054 | 0.796 | 27 |
| 6000 | arousal | 0.057 | 0.730 | 27 |
| 6000 | valence | 0.047 | 0.940 | 27 |
| 6000 | happy | 0.035 | 0.995 | 27 |
| 6000 | anger | 0.052 | 0.844 | 27 |
| 6000 | sad | 0.057 | 0.715 | 27 |
| 6000 | fear | 0.061 | 0.613 | 27 |
| 6000 | surprise | 0.046 | 0.953 | 27 |
| 6000 | disgust | 0.047 | 0.943 | 27 |
| 6000 | control | 0.046 | 0.950 | 27 |
| 6000 | fairness | 0.049 | 0.909 | 27 |
| 6000 | self-related | 0.046 | 0.955 | 27 |
| 6000 | other-related | 0.057 | 0.734 | 27 |
| 6000 | expectedness | 0.046 | 0.954 | 27 |
| 6000 | non-novelty | 0.066 | 0.451 | 27 |

**Supplementary Figure 1 | Absolute correspondence between artificial neurons and conceptual attributes of emotions. a**, Overlap and **b**, Distribution of the significant relevant artificial neurons for 14 conceptual attributes. The correspondence was computed by RSA (one-sided sign-rank test, FDR corrected). **c**, Colorbar indicates the *tau* values. **d**, Colorbar indicates the count of significantly associated artificial neurons in 12 layers of the LLM.
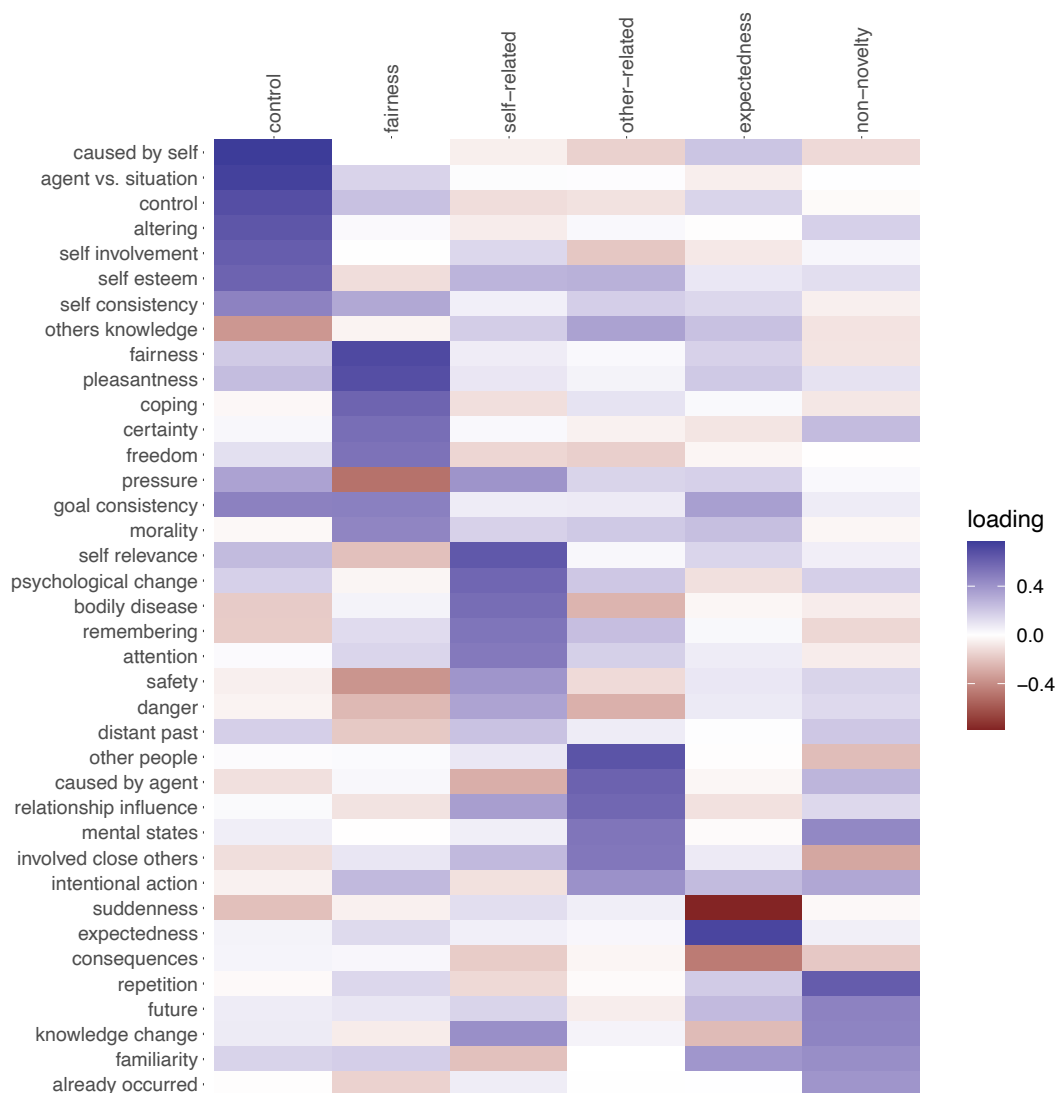
**Supplementary Figure 2 | Relative correspondence between artificial neurons and conceptual attributes of emotions. a**, Overlap and **b**, Distribution of the top N relevant artificial neurons for 14 conceptual attributes. The correspondence was computed by RSA (one-sided sign-rank test). From top to bottom, N = 1500, 2000, 2500, 3000, 4000, 5000, and 6000. **c**, Colorbar indicates the rank of artificial neurons according to *tau* values. **d**, Colorbar indicates the count of artificial neurons in 12 layers of the LLM.
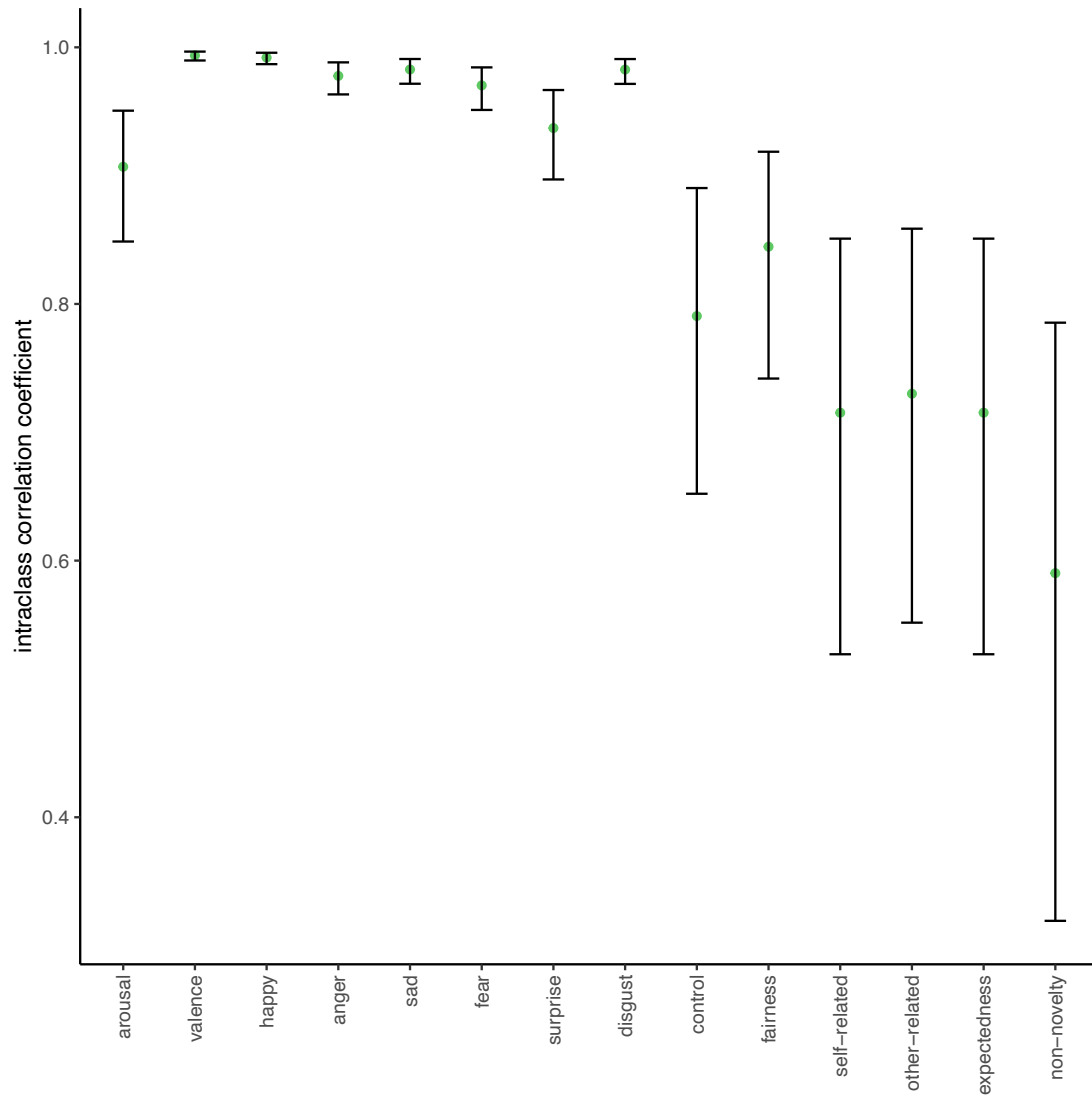
**Supplementary Figure 3 | Emotion inference accuracy.** The accuracy of 27 emotion inference tasks on the test set with each subplot corresponding to an emotion. The brighter the color, the lower the inference accuracy.

**Supplementary Figure 4 | Causal contribution of conceptual attributes to emotion inference tasks.** The causal contributions of 14 conceptual attributes to 27 emotion inference tasks when manipulating the top N neurons are indicated by the difference in accuracy compared to random manipulation. The lower the difference in accuracy, the greater the causal contribution. Stars indicate that the difference in accuracy is significant, one-sided paired *t*-test, df = 11, $p$ < .05, FDR corrected.

**Supplementary Figure 5 | Appraisal factor loadings.** Six appraisal factors of emotions were extracted from the raw scores matrix (299 participants by 38 items) obtained from the conceptual attribute rating experiment. Each participant was instructed to recall an event that directly caused them to feel a given emotion (randomly assigned) and rate 38 items on the event. Each of the 27 emotions was rated by at least 11 participants. The number of extracted factors was determined based on parallel analysis. The extraction method is Principal Component Analysis. The rotation method is Varimax (with Kaiser Normalization). The color indicates the loadings (correlation coefficients) between appraisal factors (columns) and event items (rows).

**Supplementary Figure 6 | Reliabilities of 14 conceptual attributes of emotions.** The reliabilities of 14 conceptual attributes were determined by the intra-class correlation coefficient (ICC). According to the way the three types of conceptual attributes were investigated, their reliabilities were calculated as ICC(2,k) for affective attributes, ICC(2,k) for basic emotions attributes, and ICC(1,k) for appraisal attributes, respectively. All the ICC were significant at the level = .001, FDR corrected. Error bars indicate the 95% CI.