

The Effectiveness of Style Vectors for Steering Large Language Models: A Human Evaluation

DIAOULÉ DIALLO¹, KATHARINA DWORATZYK¹, SOPHIE JENTZSCH¹, PEER SCHÜTT¹, SABINE THEIS¹, and TOBIAS HECKING¹

¹German Aerospace Center (DLR), Institute of Software Technology, Germany (e-mail: firstname.lastname@dlr.de)

Corresponding author: Diaoulé Diallo (e-mail: diaoule.diallo@dlr.de).

arXiv:2601.21505v1 [cs.AI] 29 Jan 2026

ABSTRACT Controlling the behavior of large language models (LLMs) at inference time is essential for aligning outputs with human abilities and safety requirements. *Activation steering* provides a lightweight alternative to prompt engineering and fine-tuning by directly modifying internal activations to guide generation. This research advances the literature in three significant directions. First, while previous work demonstrated the technical feasibility of steering emotional tone using automated classifiers, this paper presents the first human evaluation of activation steering concerning the emotional tone of LLM outputs, collecting over 7,000 crowd-sourced ratings from 190 participants via Prolific ($n = 190$). These ratings assess both perceived emotional intensity and overall text quality. Second, we find strong alignment between human and model-based quality ratings (mean $r = 0.776$, range 0.157–0.985), indicating automatic scoring can proxy perceived quality. Moderate steering strengths ($\lambda \approx 0.15$) reliably amplify target emotions while preserving comprehensibility, with the strongest effects for disgust ($\eta_p^2 = 0.616$) and fear ($\eta_p^2 = 0.540$), and minimal effects for surprise ($\eta_p^2 = 0.042$). Finally, upgrading from Alpaca to LLaMA-3 yielded more consistent steering with significant effects across emotions and strengths (all $p < 0.001$). Inter-rater reliability was high (ICC = 0.71–0.87), underscoring the robustness of the findings. These findings support activation-based control as a scalable method for steering LLM behavior across affective dimensions.

INDEX TERMS Activation engineering, controllable text generation, emotion control, human evaluation, large language models, style vectors

I. INTRODUCTION

ADVANCES in Generative Artificial Intelligence (AI) and the rapid evolution of Large Language Models (LLMs) have ushered in a new era of conversational systems. Their capacity to engage in near-human conversations is increasingly impressive, and they have the potential to be applied across diverse contexts—ranging from education, customer support, and healthcare to safety-critical domains such as aerospace.

With this growing capability comes a pressing need to understand and control their internal processes to ensure consistent and desirable outputs. While LLMs like LLaMA or ChatGPT have showcased remarkable proficiency in language comprehension and generation, much of this behavior remains a black box, particularly when capturing and modulating implicit cues, such as sentiment and emotion.

It has been argued that consistent language generation and delivery of information alone are insufficient for seamless

and effective communication between humans and machines. AI agents need to incorporate human factors into language generation and be adaptive concerning the human counterpart to meet the communication goals [1].

Today, there is a growing body of research exploring methodologies for nuanced steering of style and tone of LLM outputs by manipulating their internal activations during a forward pass [2]–[8]. These approaches are based on the assumption that concepts such as sentiments have a representation in the activation vector space of LLMs, and thus, the expression of these concepts can be influenced by vector operations in this space. Therefore, they are sometimes subsumed as activation steering. It could be shown that activation steering is a resource-efficient technique that can influence the tone of an LLM in a very nuanced and controllable fashion [2], which is difficult to achieve with prompt engineering.

Following the stylized response generation literature [9]–[12] and our prior work [2], we use *style* as an abstract,

content-agnostic way of expression—how something is said rather than what is said. Typical style attributes include emotion (e.g., joy, anger), sentiment (positive/negative), politeness, formality, or persona. In this work, we instantiate style as affective tone.

For intuition, think of style as a lightweight “filter” laid over the same content: the facts stay the same, but the phrasing shifts toward a chosen manner of expression (e.g., more joyful or more formal).

This paper advances this research in three key aspects: So far, all works on activation steering lack a human evaluation. The possibility of steering the tone of LLMs has been demonstrated only through automatic classification of outputs or illustrative examples. Whether humans can perceive nuanced steering in the style of LLMs is an open question, which is crucial for moving activation steering from experimentation to real human-machine interaction applications. This evaluation is critical as research has shown that human perception of AI systems—especially emotional dimensions—significantly impacts trust, acceptance, and adoption intention, explaining up to 78.5% of variance in users’ willingness to use AI in safety-critical environments [13]. To this end, we present the first human evaluation of activation steering, involving over 7,000 crowd-sourced judgments on perception of emotional tone and text quality.

Second, we analyze the relationship between steering strength and output quality, showing that moderate interventions preserve fluency while effectively shifting emotional expression. We also find that human and automatic quality ratings are closely aligned, indicating the potential of model-based scoring as a proxy in future evaluations.

Furthermore, we adapt the method of [2] to the more capable LLaMA-3 architecture, which exhibits more consistent and pronounced steering effects.

The remainder of this paper is structured as follows: Section II reviews related work on model steering and latent intervention methods. Section III outlines our activation steering approach. Section IV describes our experimental setup, including model configurations, datasets, and automatic evaluation metrics. Section V details the design and execution of our human evaluation. Section VI presents the results, including analyses of steering effectiveness, alignment between human and model evaluations, and the impact on text quality. Finally, we conclude with a discussion of limitations and directions for future work in Section VIII.

II. RELATED WORK

The introduction of transformer architectures [14] significantly advanced the capabilities of LLMs, enabling models such as GPT-3 [15] and GPT-4 [16] to encode abstract concepts implicitly, including stylistic features and emotions [17], [18]. Various methods have been proposed to effectively control or steer these outputs.

Prompt engineering directly influences model outputs by crafting specific textual prompts [19]–[23]. Although successful in task-specific scenarios, it often requires exten-

sive manual optimization and lacks precise stylistic control. Traditional fine-tuning methods in stylized response generation [9]–[12] offer deeper stylistic adjustment but involve significant computational resources, making them impractical for large-scale models.

Activation engineering modifies internal activations of LLMs during inference without retraining. Subramani *et al.* [24] initially introduced steering vectors optimized for generating predefined sentences from empty inputs, which required costly per-sentence optimization. Turner *et al.* [3] advanced this by computing steering vectors from differences between semantically opposite prompt pairs (e.g., “love” versus “hate”). These vectors directly altered model activations to effectively steer outputs in a targeted sentiment direction. Konen *et al.* [2] introduced activation-based style vectors, aggregating activations from labeled datasets representing various style categories. By contrasting the aggregated activations of a target style against the mean activations from contrasting styles, our approach produced generalized, computationally efficient vectors without requiring paired examples. Subsequently, Rinsky *et al.* [4] independently proposed Contrastive Activation Addition (CAA). Their method computed steering vectors by averaging activation differences across large sets of positive and negative examples, injecting these vectors at every token position for the control of behaviors such as factuality and sycophancy. Other studies have expanded the scope of steering multiple simultaneous behaviors. For example, van der Weij *et al.* [5] combined multiple steering vectors, observing significant interference effects and sensitivity to layer selection.

Beyond steering, recent work probes human–AI emotional alignment. Systematic reviews report that LLM empathy is often appropriate yet imperfectly aligned with human references [25], and human studies find model outputs can be perceived as empathic depending on presentation and attribution [26], [27]. Our human evaluation complements these findings by isolating graded, activation-space control and measuring perceived intensity/comprehensibility.

As outlined above, prompt-based methods are fragile and lack a scalar control of intensity, while fine-tuning variants improve control at the cost of training and deployment overhead. Activation-based style vectors complement this space by providing: (i) graded, high-precision control via a single parameter λ ; (ii) no re-training, since vectors are computed once from labeled corpora; and (iii) the ability to capture and reproduce stylistic differences present in a corpus—even when these differences are difficult to verbalize explicitly—by operating directly in activation space rather than through hand-crafted textual descriptions.

Evaluation methods for activation steering predominantly rely on automated metrics, such as classifier-based sentiment analysis or GPT-based proxies [2]–[4]. The broader literature lacks rigorous, large-scale human-based assessments concerning nuanced emotional or stylistic steering. Several studies note that steering strength can negatively affect output quality, especially in reduced fluency or increased in-

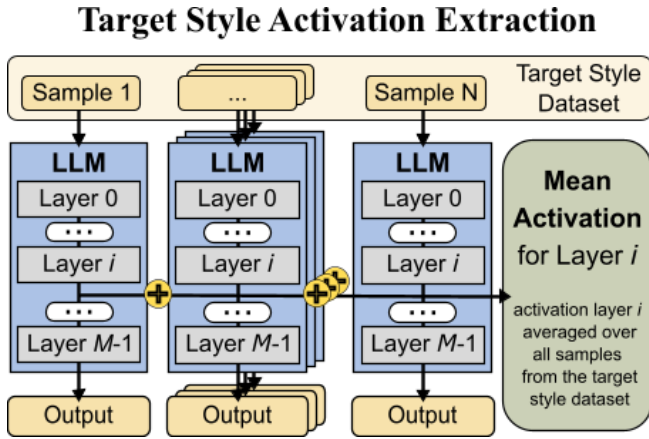


FIGURE 1. Illustration of target style activation extraction: N samples from a target style are fed into the LLM, and the mean activation at each layer i is computed to represent the characteristic activation pattern of that style.

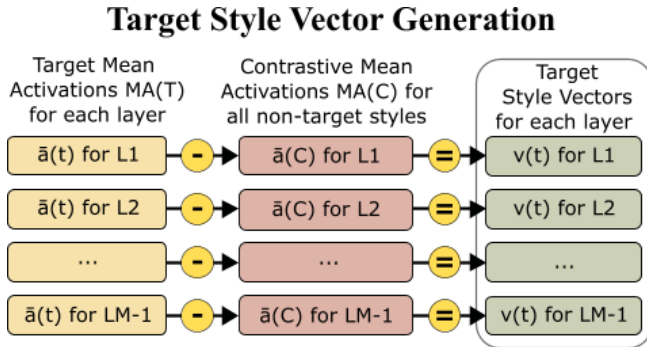


FIGURE 2. Computation of style vectors: for each layer i , the style vector $\mathbf{v}^{(i)}(t)$ is obtained by subtracting the contrastive mean activation $\bar{\mathbf{a}}^{(i)}(c) = \frac{1}{s} \sum_{j=1}^s \bar{\mathbf{a}}^{(i)}(c_j)$ from the target mean activation $\bar{\mathbf{a}}^{(i)}(t)$.

coherence at higher steering intensities [4], [8], [11], [28]. The present study addresses this research gap by combining human-centered evaluations with automated classifier metrics to assess the perceptual validity and interpretability of steering effects. Additionally, it investigates the ecological validity and mental model alignment of affective steering techniques in large language model outputs.

III. METHODOLOGY

In this section, we will first briefly describe the activation steering pipeline based on [2]. Subsequently, we comprehensively elaborate improvements, as a special focus of this work lies on the updated experimental pipeline.

We aim to modify an LLM's responses to any input prompt x during an ongoing conversation by steering them toward a specific target style t . To this end, a set of layer-wise style vectors for a predefined target style t is constructed based on a labeled dataset.

We define the set of style categories as $\mathcal{Y} = \{t, c_1, \dots, c_s\}$, where t denotes the *target* style and c_1, \dots, c_s are s contrastive (non-target) styles. For each label $y \in \mathcal{Y}$, we assume a set of representative samples $d_y^{(n)}$ from a labeled corpus. Each sample is passed through the LLM in a forward pass to collect

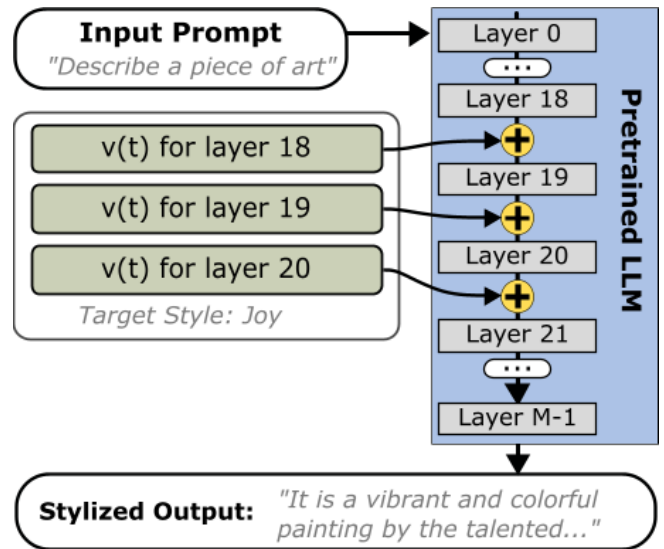


FIGURE 3. Process of steering model activations during inference. In this illustration, only three layers—18, 19, and 20—are modified using activation-based style vectors for the target style joy; in our full experiments, we inject the vectors into all model layers.

hidden activations from each layer L_i (see Fig. 1).

The per-layer activations are then averaged over all samples of a specific style label. We denote the mean activation vector at layer L_i for the target style t as $\bar{\mathbf{a}}^{(i)}(t)$ and for contrast style c_j as $\bar{\mathbf{a}}^{(i)}(c_j)$.

To capture what distinguishes the target style from the others, we compute the style vector $\mathbf{v}^{(i)}(t)$ for layer L_i as the difference between the target and contrastive mean activations:

$$\mathbf{v}^{(i)}(t) = \bar{\mathbf{a}}^{(i)}(t) - \frac{1}{s} \sum_{j=1}^s \bar{\mathbf{a}}^{(i)}(c_j) \quad (1)$$

This vector encodes a direction in activation space that characterizes the unique traits of the target style t at each layer (see Fig. 2).

At inference time, the steering is performed by modifying the original activations $\mathbf{a}^{(i)}(x)$ at each layer L_i using the corresponding style vector (see Fig. 3). The updated activations become:

$$\hat{\mathbf{a}}^{(i)}(x) = \mathbf{a}^{(i)}(x) + \lambda \cdot \mathbf{v}^{(i)}(t) \quad (2)$$

Here, λ is a scaling parameter that controls the intensity of the steering effect and is shared across layers. Operationally, λ behaves like a single fader for the target emotion—pushing it higher intensifies the tone, pulling it lower restores a neutral baseline. This method enables nuanced and controllable stylistic modulation during generation, offering an alternative to prompt-based conditioning. Algorithm 1 summarizes vector construction (Eq.1) and inference-time integration (Eq.2) as pseudocode.

In our study, style concepts refer to emotion categories such as *joy*, *anger*, or *disgust*, but this approach generalizes to any style distinguishable via a representative corpus.

Algorithm 1 Activation-based Style Steering (construction and inference)

Require: Labeled corpus \mathcal{D} ; style labels $\mathcal{Y} = \{t, c_1, \dots, c_s\}$ with $C = \{c_1, \dots, c_s\}$; LLM layers L_0, \dots, L_{M-1} ; steering strength λ

Ensure: Per-layer style vectors $\{\mathbf{v}^{(i)}(t)\}_{i=0}^{M-1}$ and steered continuation of x

```

1: procedure BuildStyleVectors( $\mathcal{D}, \mathcal{Y}$ )
2:   for  $i \leftarrow 0$  to  $M - 1$  do
3:      $\bar{\mathbf{a}}^{(i)}(t) \leftarrow$  mean activation at  $L_i$  over samples of  $t$ 
4:      $\bar{\mathbf{a}}^{(i)}(C) \leftarrow \frac{1}{s} \sum_{j=1}^s \bar{\mathbf{a}}^{(i)}(c_j)$ 
5:      $\mathbf{v}^{(i)}(t) \leftarrow \bar{\mathbf{a}}^{(i)}(t) - \bar{\mathbf{a}}^{(i)}(C)$   $\triangleright$  Eq. 1
6:   end for
7:   return  $\{\mathbf{v}^{(i)}(t)\}_{i=0}^{M-1}$ 
8: end procedure
9: procedure SteerAtInference( $x, \{\mathbf{v}^{(i)}(t)\}, \lambda$ )
10:  for token step  $u = 1, \dots, U$  do
11:    for  $i \leftarrow 0$  to  $M - 1$  do  $\triangleright$  steer all layers; see Fig. 3
12:       $\hat{\mathbf{a}}^{(i)}(x) \leftarrow \mathbf{a}^{(i)}(x) + \lambda \cdot \mathbf{v}^{(i)}(t)$   $\triangleright$  Eq. 2
13:    end for
14:    Generate next token using steered activations across all layers
15:  end for
16: end procedure

```

IV. EXPERIMENTS

This section provides a detailed overview of the experimental setup used to evaluate the effectiveness of activation-based emotion steering in LLMs. Our overall approach involves three main components: (1) constructing style vectors in the activation space of the LLM corresponding to different target emotions, (2) prompting the model with input texts and systematically applying these style vectors at varying intensities, and (3) evaluating the resulting outputs through both automatic classification models and human assessment.

Through this pipeline, we investigate whether steering the internal activations of a language model leads to perceivable and interpretable changes in the emotional tone of its outputs. We first describe the selected LLM and its configuration, followed by the dataset used to derive emotion-specific style vectors. We then present our automatic evaluation setup and, finally, the design and implementation of a human evaluation to assess perceived emotional intensity and text quality.

A. MODELS AND PARAMETERS

We initially experimented with several variants of instruction-tuned language models and selected Llama-3-8B-LexiUncensored [29] due to its superior performance in preliminary style vector evaluations. Compared to the standard LLaMA-3 model, this variant demonstrated better control capabilities in latent activation space, especially in tasks involving emotional and stylistic conditioning.

To extract hidden activations from all model layers during forward passes, we processed textual samples from our

dataset with a maximum token length of 300. This constraint was necessary to avoid GPU memory overflows during batch processing. If input prompts exceeded this token threshold, we iteratively truncated the input text from the end until it satisfied the limit. Although this allowed stable processing, it occasionally resulted in incomplete or semantically truncated sentences.

B. DATASET

GoEmotions [30] is a multi-class style dataset consisting of 58k manually curated user comments from the internet platform Reddit. Comments are labeled with 27 emotional categories. From this, we extract a subset of 53,994 unique samples that can be unambiguously mapped to one of the six basic emotion categories defined by Ekman [31]: *sadness*, *joy*, *fear*, *anger*, *surprise*, and *disgust*. The mapping from the original 27 emotion labels to the six Ekman categories was performed using the official mapping provided in the GoEmotions repository [32]. The distribution of unique samples per Ekman emotion label is shown in Table 1.

TABLE 1. Sample distribution per emotion label after Ekman mapping.

Emotion	Sample Count
Joy	19,440
Neutral	17,716
Surprise	5,839
Anger	5,682
Sadness	3,622
Disgust	881
Fear	814

Although the *neutral* category is not part of Ekman's six basic emotions, it was included in the construction of style vectors to capture general-purpose directions in the model's latent space.

C. AUTOMATIC EVALUATION

To complement the human ratings, we automatically evaluated using two model-based scorers: one for emotion intensity and another for text comprehensibility. For emotion classification, we employed a pre-trained transformer model fine-tuned on English emotion detection, specifically DistilRoBERTa [33]. For this model an overall classification accuracy of 66% on a balanced six-emotion evaluation set (chance level $\approx 14\%$) was reported. The model outputs a probability distribution across emotion categories, which we use to track how strongly the target emotion is expressed in the generated text.

We applied the same LLM used for text generation in a scoring mode to evaluate comprehensibility. The model was prompted to rate each sentence on a scale from 1 (highly comprehensible) to 10 (incomprehensible), based on predefined instructions emphasizing clarity and logical coherence.

V. HUMAN EVALUATION

A. BACKGROUND

Paul Ekman's model [31], [34] encompasses the six basic emotions of *anger*, *disgust*, *fear*, *joy*, *sadness*, and *surprise* and justifies its application by recognizing such emotions as universal and cross-cultural. His studies showed that these emotions are displayed and identified similarly around the world, indicating their potential suitability for applications such as emotion control—a possibility further explored in the present study. Another motivation for adopting this emotion model in the present study was that most available training datasets are based on the Ekman model. The distinction between the corresponding emotion categories supports our assumption that people with an unrestricted ability to perceive emotions can recognize the differences in intensity within these categories. A study on cognitive appraisal methods showed humans can identify emotion in texts with the help of contextual cues even if these are implicit (e.g., meeting a snake is associated with fear) [35]. However, Ekman's primary emotions are derived from studies of facial expressions. There is currently little knowledge about the human perception of emotions in text. Initial approaches focus primarily on developing annotation schemes and their validation [36].

B. SAMPLE

In total, 227 participants (50% male and 50% female) were recruited via Prolific [37]. To account for sex-specific differences in emotion interpretation, assignment to the different prompt conditions was balanced for participants' sex. Each of the 19 prompt conditions was completed by 5 female participants and 5 male participants. As the study involved extensive reading and understanding of English texts, participants were selected if their first language was English and they had no literacy difficulties. Responses from participants were excluded if they did not pass the prescreening (see Section V-D2), indicating insufficient emotion detection abilities, or if a comparatively short completion time indicated low-effort responses. Exclusion of 37 responses yielded a final sample of 190. Of the final sample, 50% reported their gender as woman, 49.5% as man, and one person as non-binary. The average age was 38.42 ($SD = 13.57$) with a range of 18 to 78 years. All participants reported English as their first language and no literacy difficulties; all but one rated their reading fluency as *completely fluent* (88.4%) or *very fluent* (11.1%). The majority was currently residing in the United Kingdom (52%), followed by Canada (15%), South Africa (13%), United States (12%), Ireland, New Zealand, Australia, Spain, and Korea (less than 5%). A total of 62% identified with Western or European culture, 28% with African culture, 4% with Asian culture and the remaining 6% with another culture (e.g. Middle Eastern or Indigenous). Most of the participants were highly educated (Bachelor's degree: 46.9%, Master's degree: 20.0%, Doctoral degree: 3.2%), full-time or part-time employed (89%), and comfortable in interpreting emotions in English text (*extremely comfortable*: 62%, *very comfortable*: 35%, *moderately comfortable*: 3%). An a priori power analysis (F-tests, repeated measures ANOVA) using

G*Power determined a sample size of $n = 190$ to be adequate to detect a small effect with 95% power at a significance level of $\alpha = .05$.

C. STUDY DESIGN

To investigate how humans perceive the emotion steering strength of an LLM's text output when steered to six basic emotion types, a 6×8 within-subjects design was employed. Therein, we considered the effect of the independent variables *target emotion* (*anger*, *disgust*, *fear*, *joy*, *sadness* and *surprise*) and *steering strength* (λ) (0.00, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35) on the dependent variables *perceived emotional intensity* and *text comprehensibility*. As there is currently little empirical evidence regarding the perceptibility of emotional control intensity in texts, we have decided to initially look at very small-step adjustments of the λ value in the present experiment. Furthermore, high λ values affect the generated texts' coherence, readability, and naturalness. Preliminary experiments revealed that steering strengths beyond $\lambda = 0.35$ frequently led to heavily degraded semantic coherence and overall output quality, and were excluded from the final evaluation.

D. MATERIALS AND MEASURES

1) Stimuli

Texts were generated by the LLM using predefined prompts for each target emotion and λ value (for examples see Appendix B). The generated texts were identical across all target emotions for $\lambda = 0.00$.

2) Emotion Detection Ability

Prescreening included a self-administered test on emotion detection ability in text, ensuring the inclusion of participants with adequate emotion recognition skills. As part of the test, participants rated six texts from the International Survey on Emotion Antecedents and Reactions (ISEAR) dataset collected by the Swiss Center for Affective Sciences; the dataset comprises 7,666 self-reported emotional experiences from nearly 3,000 individuals across diverse cultural backgrounds [38]. Each entry in ISEAR includes a brief textual description of an emotional event, categorized into seven primary emotions—*joy*, *fear*, *anger*, *sadness*, *disgust*, *shame*, and *guilt*—based on appraisal theories of emotion. ISEAR has been extensively used in machine learning models for emotion classification, providing a well-validated benchmark for lexical, syntactic, and contextual emotion cues [39]. For each emotion category in Ekman's model represented in ISEAR, we selected sentence IDs (*joy*: 444, *fear*: 214, *sadness*: 467, *disgust*: 124, *anger*: 42), with a large word count to better match the length of our LLM output texts. Unlike traditional emotion recognition scales that rely on categorical labels or forced-choice paradigms, the proposed method employs a continuous rating system where participants assess the intensity of six emotions (*joy*, *anger*, *sadness*, *fear*, *disgust*, and *surprise*) within a 0–7 scale for various textual scenarios. Dominant emotion for each text was determined by the

highest-rated emotion with a flexible threshold (± 0.5 around the maximum value) to account for natural variability in human perception, using a scale adjustable in 0.1 increments. When multiple emotions were rated highly relevant, reflecting the reality of human affective processing, participants must correctly identify the dominant emotion in at least three out of five passages to pass the screening. This ensures that minor variations or ties do not lead to unnecessary exclusions. Compared to standardized emotion recognition tests, such as Ekman's Facial Emotion Recognition Task or the Reading the Mind in the Eyes Test, this method provides superior ecological validity by evaluating emotion perception in context rather than relying on static, artificial stimuli. Furthermore, the continuous scale mitigates floor and ceiling effects, offering a granular assessment of emotion perception accuracy. By integrating context-sensitive, flexible, and adaptive scoring criteria, this pre-screening method better reflects real-world emotion recognition, ensuring that subsequent survey responses are derived from individuals with sufficient emotion detection ability.

3) Demographics

Participants were asked to indicate their gender, age, country of residence, culture, level of education, field of study or work, first language, level of reading fluency, and emotion interpretation. Additional demographic data, such as sex, literacy difficulties, and employment status, were obtained from their Prolific profile.

4) Emotional Intensity

The emotional intensity participants perceived when reading the generated texts was assessed using six rating scales. Each rating scale represented one of the six basic emotions and ranged between 0 and 7 with 0.1 increments, allowing for a fine-grained assessment. The same response format as in the prescreening was used to be able to accurately represent the complexity of human perception. According to Ekman, such rating scales allowed the participants to indicate the intensity of the basic emotions they associate with the text. In formulating the instructions, the emphasis was placed on evaluating the emotion conveyed by the text and not on the emotion elicited in the participant ("Select the intensity of each emotion and the comprehensibility you would assign to the following text."). A score of 0 indicated that the person could not associate any intensity of emotion with the text. In contrast, a score of 7 showed that the emotion was associated with the text at the highest possible intensity.

5) Text Comprehensibility

To assess the overall text quality, participants rated the comprehensibility of the output texts ("Select the intensity of each emotion and the comprehensibility you would assign to the following text") on the same 0 to 7 scale as emotional intensity, where 0 indicated low comprehensibility and 7 high comprehensibility.

E. PROCEDURE

The study was transparently introduced as an investigation of human perception of emotions in written texts. Participants were informed about the scope and the structure of the study, including the prescreening, and the risk of experiencing discomfort when reflecting on emotions or experiences. In addition to providing information about the use of the data, participants were assured that all responses were collected anonymously and encouraged to take breaks during the study. After providing informed consent via the survey tool LimeSurvey [40], participants started the experiment by rating five prescreening texts. Participants who failed to correctly rate at least three prescreening texts were redirected to the Prolific site and did not participate in the main study. Participants who completed the prescreening successfully were asked to answer ten demographic questions. Next, participants were randomly assigned to one of 19 prompts (see Appendix A). Texts were generated for the six basic emotions and eight steering strengths for each prompt, resulting in 42 distinct and 6 identical ($\lambda = 0.00$) outputs. Participants rated each text on perceived emotional intensity and text comprehensibility ("Select the intensity of each emotion and the comprehensibility you would assign to the following text."). The order in which participants rated the texts was randomized with respect to target emotion, λ value, and perceived emotional intensity. At the end of the study, participants were redirected to the Prolific site and received a completion code. After manually reviewing the submitted responses, participants who completed the research received compensation of £7.50. Participants who were screened out and did not participate in the central part of the study received compensation of £0.10 per minute. The study duration was estimated at 60 minutes; the actual duration ranged between 14 and 156 minutes, with an average of 65 minutes ($SD = 29$). Throughout the study, participants could contact the experimenters through the Prolific platform. Participants reported two technical issues (incorrect display of the progress bar and failed randomization of one output text), which were subsequently resolved by the experimenters.

F. STATISTICAL ANALYSIS

Missing data for the target emotion *joy* and λ -value 0.10 and 0.25, which occurred due to a malfunction in the randomization process in nine cases, were replaced with mean values (54 replacements in total, $< 1\%$). We then conducted a two-way repeated-measures analysis of variance (ANOVA) to test for differences in the perceived emotional intensity as well as text comprehensibility related to target emotion and steering strength (λ).

To determine the inter-rater reliability among human evaluators, intraclass correlation coefficient (ICC) estimates were calculated based on single rater as well as mean of ten raters. Values between 0.5 and 0.75 were regarded as evidence of moderate reliability, values between 0.75 and 0.90 of good reliability, and values exceeding 0.90 of excellent reliability as suggested by [41].

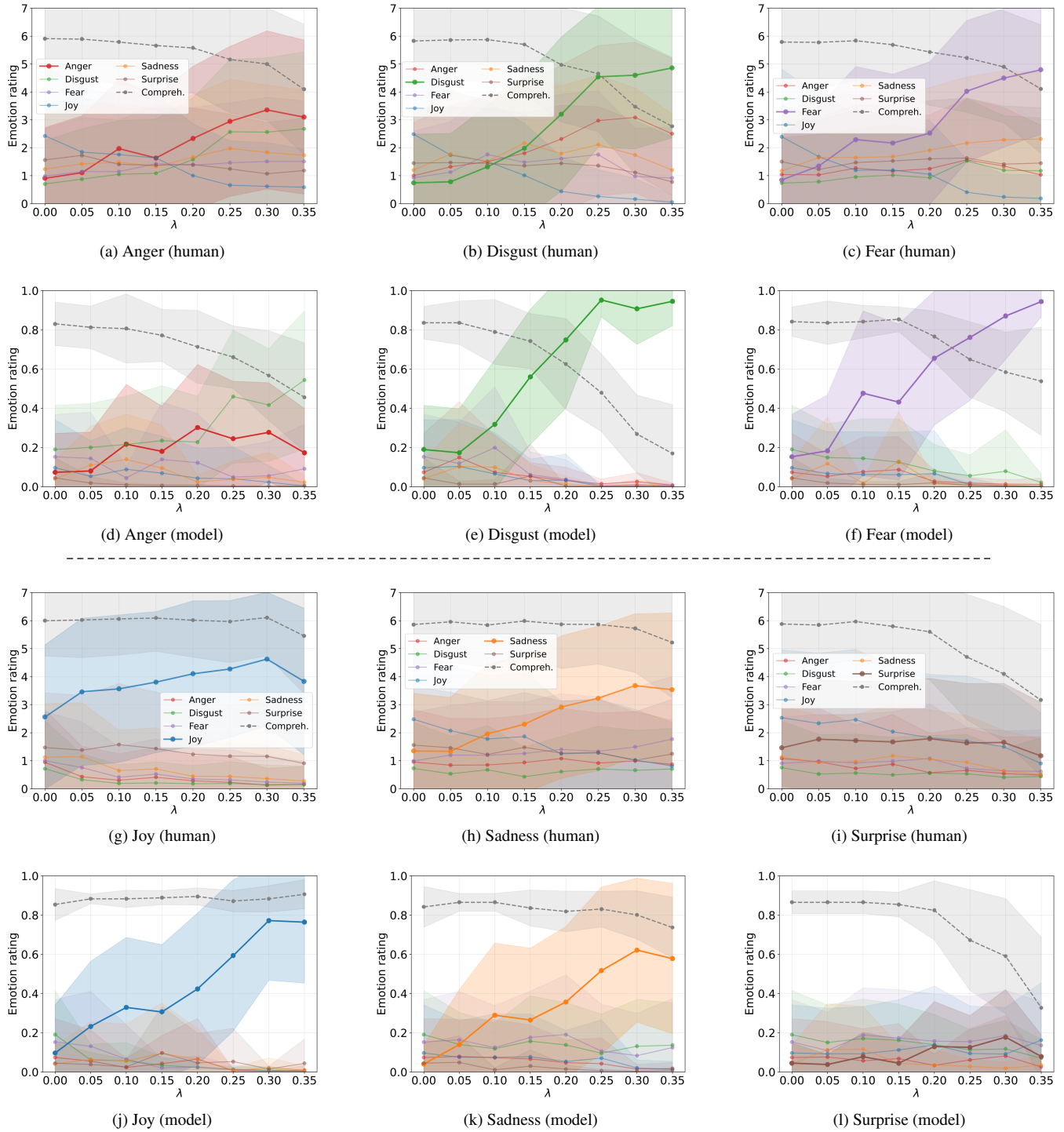


FIGURE 4. Human and model emotional perception and text comprehensibility ratings across steering strengths λ for six emotions. In each block, the top row shows average human mean ratings and the row below shows corresponding model scores. Shaded areas represent ± 1 standard deviation.

VI. RESULTS OF EMOTION STEERING

To assess the effect of varying steering strength on perceived emotion in generated text, we averaged emotion ratings over all participants and prompts for each λ -value and target emotion. In all experiments, steering vectors were applied across all layers. We used a fixed set of 19 prompts to elicit emotionally diverse responses, e.g., “Comment on a review of a business written by a customer”. The complete list of prompts is provided in Appendix A. Prior work and our preliminary experiments show that the effect of λ strongly depends on the number of layers to which steering is applied. All results and λ -based interpretations in this paper refer to the all-layer steering configuration.

Fig. 4 displays the mean human ratings alongside our automatic emotion classification model outputs. While participants provided ratings on a scale from 0 (no emotion) to 7 (maximum intensity), the model scores range from 0 to 1. Each line plot displays the ratings for one of six Ekman emotions, with the target emotion highlighted in bold and non-target emotions shown with lower opacity. Shaded areas in the plots represent ± 1 standard deviation across participants or model predictions, respectively. The figure also includes comprehensibility ratings. Participants rated the comprehensibility of each text on the same 0–7 scale (higher is better), while the model’s comprehensibility predictor outputs scores from 1 (high comprehensibility) to 10 (low). We normalized and inverted these model comprehensibility scores for visual alignment.

A. EMOTIONAL INTENSITY ACROSS STEERING STRENGTHS

Overall, the human ratings revealed a clear and consistent increase in perceived emotional intensity, with higher λ -values for five of the six emotions. The effect was particularly pronounced for *disgust* and *fear*, both of which rose from an average rating of 0.75 and 0.85 (at $\lambda = 0.00$) to 4.86 and 4.80 (at $\lambda = 0.35$). *Joy* also showed a positive trend, increasing from 2.57 to 3.83 with a peak at $\lambda = 0.30$ reaching 4.63. *Anger* and *sadness* followed similar, though more moderate, trends—anger increased from 0.91 to 3.10, and *sadness* from 1.35 to 3.54. Notably, *surprise* did not exhibit a meaningful change in perceived intensity (from 1.47 to 1.18), suggesting that it is less effectively steered via our method. Mean ratings are also summarized in Table C in Appendix C.

To test our observations statistically, we conducted a two-way repeated-measures ANOVA. Mauchly’s test indicated that the assumption of sphericity was violated for all effects and all but two ϵ values were < 0.75 , Greenhouse-Geisser correction was applied to the degrees of freedom for all effects. Results of the ANOVAs (Table 2) revealed significant main effects of target emotion (all $ps < 0.001$), steering strength (all $ps < 0.001$) as well as significant interactions between target emotion and steering strength, suggesting that the effect of λ depends, to some extent, on emotion. The effect of emotion was overall stronger than the effect of steering strength and strongest in perceived emotional intensity of *joy*

TABLE 2. Two-way repeated-measures Analysis of Variance (ANOVA) in perceived emotional intensity and text comprehensibility

	df1	df2	F	p	η_p^2
Anger					
Target emotion (E)	3.12	589.81	161.61	<.001	.461
Steering strength (λ)	3.50	661.71	20.29	<.001	.097
E x λ	19.67	3716.64	24.93	<.001	.117
Disgust					
Target emotion (E)	3.00	567.40	303.21	<.001	.616
Steering strength (λ)	4.07	768.53	77.55	<.001	.291
E x λ	18.19	3437.59	58.93	<.001	.238
Fear					
Target emotion (E)	3.53	667.54	222.26	<.001	.540
Steering strength (λ)	4.84	914.82	24.16	<.001	.113
E x λ	21.97	4151.99	41.15	<.001	.179
Joy					
Target emotion (E)	2.56	484.35	305.91	<.001	.618
Steering strength (λ)	2.98	562.98	54.37	<.001	.223
E x λ	22.72	4293.67	26.98	<.001	.125
Sadness					
Target emotion (E)	3.52	664.90	124.08	<.001	.396
Steering strength (λ)	4.55	859.59	13.87	<.001	.068
E x λ	22.69	4288.75	20.43	<.001	.098
Surprise					
Target emotion (E)	4.08	771.63	8.19	<.001	.042
Steering strength (λ)	4.90	926.85	8.56	<.001	.043
E x λ	24.66	4659.96	2.55	<.001	.013
Comprehensibility					
Target emotion (E)	4.20	794.17	100.20	<.001	.346
Steering strength (λ)	3.38	638.71	179.85	<.001	.488
E x λ	20.31	3838.68	22.78	<.001	.108

Greenhouse-Geisser corrected degrees of freedom for all effects.

($\eta_p^2 = 0.618$) and *disgust* ($\eta_p^2 = 0.616$) and weakest in perceived emotional intensity of *surprise* ($\eta_p^2 = 0.042$). In line with our assumptions, we found a relatively strong effect of steering strength on text comprehensibility ($\eta_p^2 = 0.488$). Overall, these results support the effectiveness of emotion steering for five of the six emotions.

The automatic model ratings confirmed and even amplified the trends observed in our sample of human evaluators. *Disgust*, *fear*, and *joy* rose sharply from approximately 0.1–0.2 at $\lambda = 0.00$ to scores near 0.9–1.0 at $\lambda = 0.35$. *Sadness* showed a moderate increase up to about 0.6. Notably, *anger* appeared to be less effectively captured by the classifier, showing only a minor increase from around 0.05 to 0.2, peaking early at $\lambda = 0.20$. Instead, the steering of anger seemed to clearly affect the disgust emotion score, starting at 0.2 and reaching a score slightly above 0.5. *Surprise*, in line with the human ratings, displayed no or only a minimal trend across the steering range.

Interestingly, a key difference emerged when comparing human and model ratings for non-target emotions. In the case of the automatic emotion model, ratings for non-target emotions remained largely stable or even decreased as the

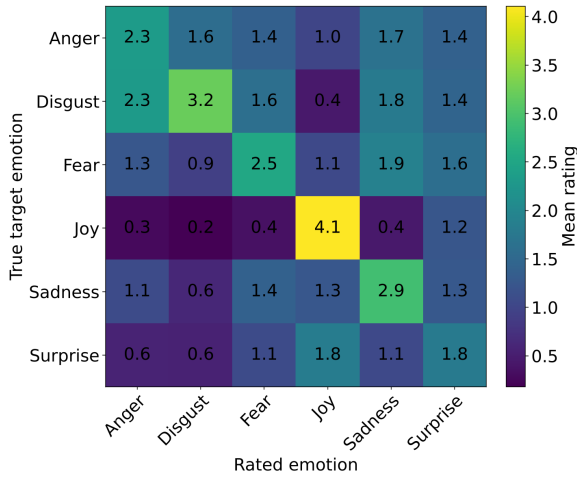


FIGURE 5. Participants' mean ratings heatmap at $\lambda = 0.20$. Rows denote the true target emotion; columns denote the average rating participants assigned to each emotion dimension (0–7 scale).

steering strength (λ) increased. This suggests that the model is highly selective in interpreting emotional expression, predominantly increasing only the intended emotion dimension while suppressing or ignoring unrelated ones.

By contrast, human ratings revealed a more complex and overlapping emotion perception. For instance, when steering towards *anger*, participants not only rated *anger* higher but also reported increased levels of perceived *disgust* and *sadness*, with a minor increase in *fear*. Similarly, steering towards *disgust* caused higher ratings for both *anger* and—less consistently—*sadness* and *fear*. In the case of *fear*, participants also perceived elevated levels of *sadness*, further highlighting the interconnected nature of specific emotional categories.

These findings reflect the inherently multidimensional and context-dependent nature of human emotion perception. Unlike the model's interpretation, human raters tend to perceive emotional states as overlapping or co-occurring, especially among low-valence emotional expressions [42]–[46]. The results emphasize that while activation-based steering successfully amplifies target emotions, the resulting texts often carry nuanced emotional undertones that human perception picks up, even when these are not explicitly modeled.

To quantify co-perception of emotions, Fig. 5 shows a heatmap of mean human ratings at moderate steering strength $\lambda = 0.20$. For each target emotion (rows), the six rated intensities (0–7) across emotion dimensions (columns) were averaged. Strong diagonals indicate successful steering, especially for *joy*. Off-diagonal clusters reveal systematic co-activation among low-valence emotions: *anger* aligns with *disgust*, *sadness*, and *fear*, and *fear* aligns with *sadness*. *Surprise* remains weak and diffuse. These patterns mirror the trajectories in Figs. 4: steering a target emotion increases its perceived intensity while also elevating related low-valence dimensions, highlighting the interconnected perception of affective categories.

Despite the consistent upward trends across participants,

TABLE 3. Inter-rater reliability (means across 19 prompts per emotion). ICC(2,1): single rater; ICC(2,k): mean of $k=10$ raters.

Emotion	ICC(2,1)	ICC(2,k)
Anger	0.33	0.79
Disgust	0.43	0.87
Fear	0.36	0.81
Joy	0.42	0.85
Sadness	0.24	0.71
Surprise	0.05	0.28

the shaded areas in the plots indicate substantial variance in individual emotion ratings. This highlights that while the average perception of emotional intensity increases with higher λ , participants differed in how strongly they perceived the emotional tone of the generated texts. These differences may stem not only from personal interpretation styles, cultural or linguistic nuances, or varying thresholds for recognizing emotion in written language, but also from differences in the prompts themselves. Thus, the results should reflect population-level trends rather than uniform perceptions across all raters.

To determine the consistency of human ratings, we quantified inter-rater reliability using intraclass correlations (ICCs), reporting both single-rater ICC(2,1) and aggregated-rater ICC(2,k) values. ICCs were computed within each prompt \times target-emotion block and averaged across prompts. As summarized in Table 3, single-rater reliability is modest ($M \approx 0.24$ – 0.43 across five emotions), while aggregated reliability is high for *anger*, *disgust*, *fear*, *joy*, and *sadness* ($M \approx 0.71$ – 0.87). In contrast, *surprise* shows low agreement (≈ 0.28), reinforcing that this emotion is weakly and inconsistently perceived in our setting—consistent with its limited steerability observed in Figs. 4 and 6, and Table 2. These results suggest that although individual raters vary considerably in their ratings, the aggregated judgments reveal consistent population-level trends.

B. ALIGNMENT OF HUMAN JUDGMENTS AND AUTOMATIC CLASSIFIER

To further examine the alignment between human perception and the model's emotion estimates, we directly compared average ratings for each target emotion across increasing λ -values, as shown in Fig. 6. Each plot contains one line for the human ratings (0–7 scale, left y-axis) and one for the model scores (0–1 scale, right y-axis).

The trends demonstrate a good alignment between model and human ratings as steering strength increases. Particularly for *sadness*, *disgust*, and *fear*, the curves closely follow one another, indicating that the direction and shape of emotional amplification are comparable across both perspectives.

However, some divergences become apparent. For *anger*, the deviation between model and human ratings began around $\lambda = 0.2$, where human raters perceived less intensity than the model estimates. In the cases of *disgust* and *fear*, both ratings increased consistently with λ , but the model scores rose more steeply, approaching the upper limit of 1, while human ratings

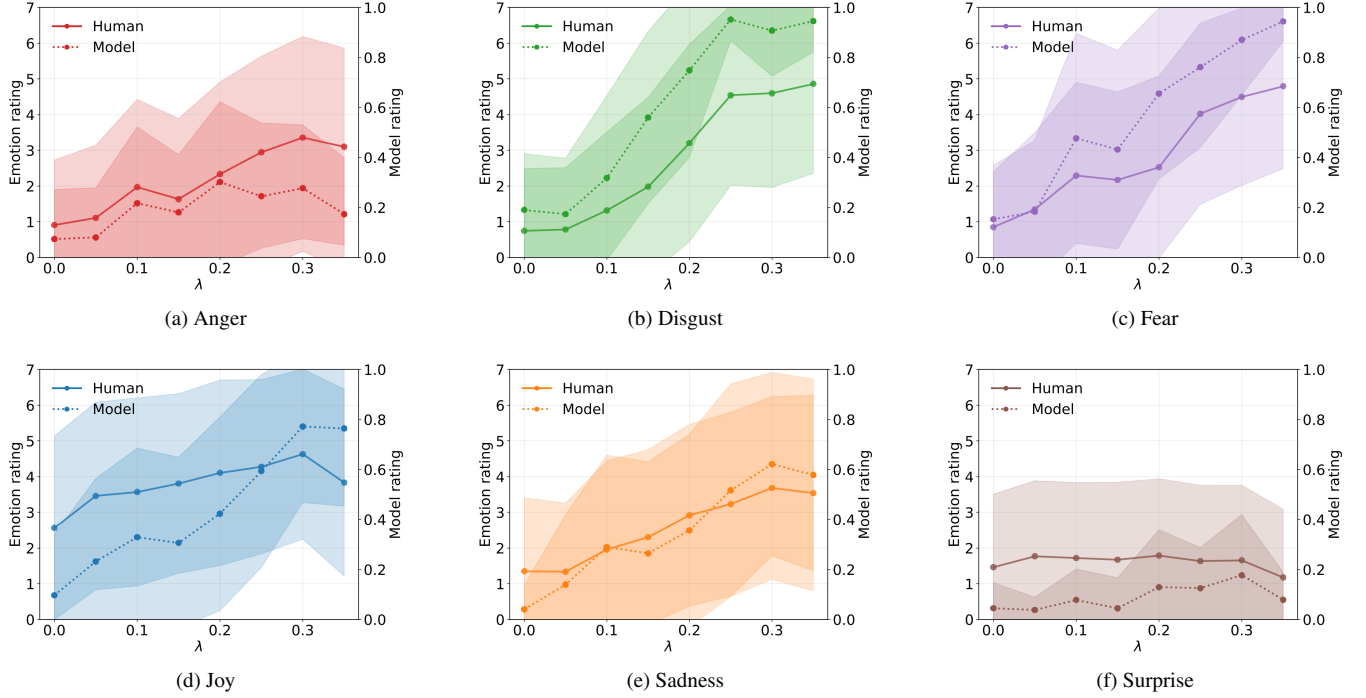


FIGURE 6. Comparison of human emotion perception and model emotion steering ratings across six target emotions. Each plot shows the average perceived intensity of the target emotion as a function of steering strength (λ). Human ratings (left y-axis, 0–7 scale) are plotted alongside model predictions (right y-axis, 0–1 scale). Shaded areas represent ± 1 standard deviation.

remained lower. This suggests that humans apply more conservative upper-bound ratings while the emotional trajectory is shared.

To quantify the alignment between human and model ratings, we computed Pearson’s correlation coefficient over each emotion’s average ratings per λ . The results revealed strong correlations for most emotions: *disgust* ($r = .985$), *fear* ($r = .973$), and *sadness* ($r = .971$) show high alignment. *Joy* ($r = .813$) and *anger* ($r = .758$) also demonstrate substantial correspondence, despite some scale mismatches. The only notable outlier was *surprise*, with a much lower correlation of $r = .157$, reinforcing the impression that this emotion is less consistently influenced by the steering mechanism. Across all six emotions, the average Pearson correlation was $r = .776$, suggesting a generally strong relationship between human perception and model output. In the case of *joy*, an interesting pattern emerged: human raters already perceived a notable level of joy at $\lambda = 0$ (around 2.5), while the model outputs remained near zero. As λ increases, the model’s score rose more sharply than human perception, resulting in a crossover at approximately $\lambda = 0.25$. A similar but less pronounced effect occurs for *sadness*, where the model begins at a lower value but surpasses human ratings by the highest steering strength. For *surprise*, no meaningful change in perceived emotion occurred for humans or the model across steering levels, reinforcing previous findings that this emotion is less susceptible to our current steering method. Across the board, the progression of ratings for λ tends to show similar accel-

ation in both human and model responses, suggesting that λ provides a reasonably consistent control signal. This consistency across emotion types reinforces the viability of using λ as a global steering parameter, even if perceived intensity ceilings differ between raters and automated classifiers. Among the five emotions for which we observe an apparent steering effect, four exhibited a stronger final intensity in the model ratings than the human assessments (all except *anger*). This indicates that human raters generally avoid extreme values in the upper range of the scale: in no case does the average human rating exceed 5, even at the highest steering level. The model, by contrast, tends to saturate more aggressively, particularly for emotions such as *disgust* and *fear*.

C. IMPACT ON TEXT QUALITY

A core question for activation-based steering is whether increasing the steering strength (λ) pushes the model beyond regions of its latent space that support coherent language generation. Because style vectors perturb internal activations—especially at higher magnitudes—there is a genuine risk of degrading semantic and syntactic structure. Prior work has shown that substantial activation edits can yield incoherent, repetitive, or semantically inconsistent outputs [4], [8], [11], [28]. We therefore examined how text quality evolves with λ , combining direct comprehensibility ratings with surface-level text features.

Fig. 7 summarizes five normalized measures across λ : mean word length, lexical density, entropy, and both human and model comprehensibility scores. *Mean word length*

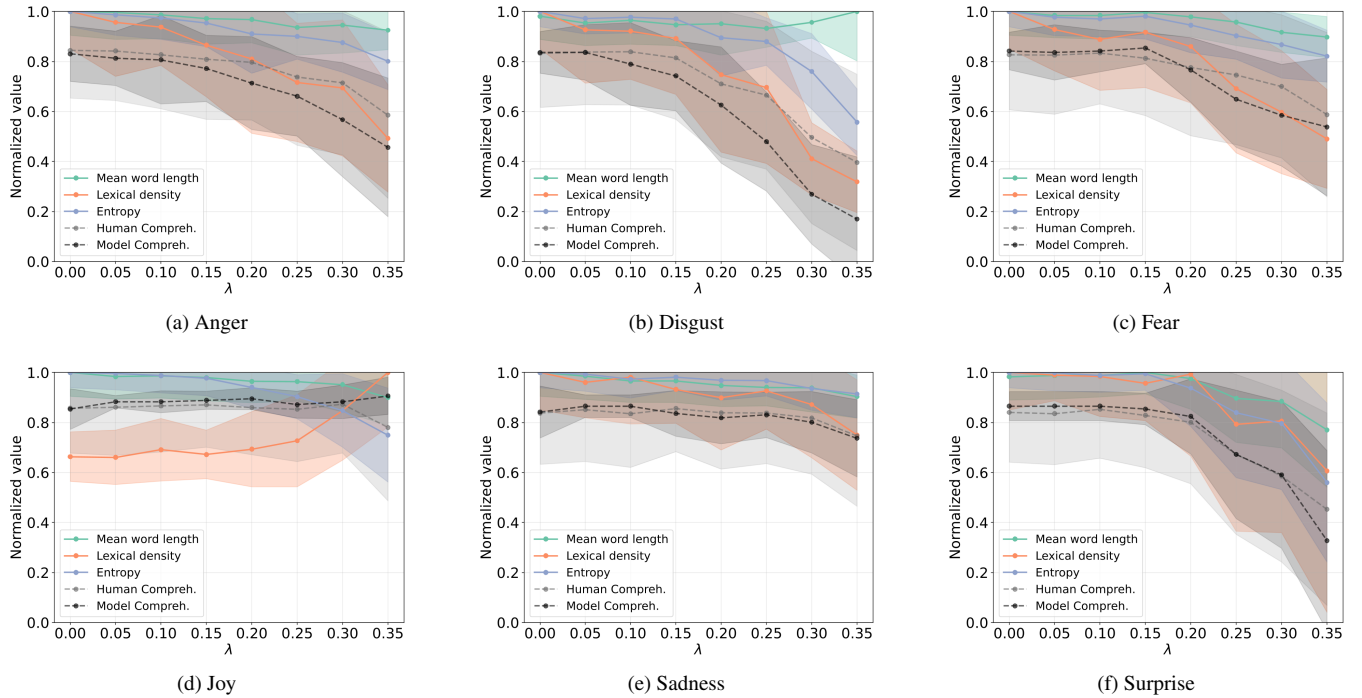


FIGURE 7. Normalized text features and comprehensibility scores across steering strengths λ . Each plot corresponds to one target emotion and shows five normalized values: lexical density, mean word length, entropy, and comprehensibility as rated by humans and a model. All values are scaled to their respective maximum to illustrate relative trends.

serves as a rough proxy for lexical sophistication. *Lexical density* is the proportion of content words (nouns, verbs, adjectives, adverbs) in a sentence and signals information richness. *Entropy* captures vocabulary diversity: it is the Shannon entropy of the word-frequency distribution in the text, so higher values indicate more varied word use, whereas low values reflect repetition.

For most emotions—*disgust*, *fear*, *anger*, and *surprise*—all text features remained relatively stable up to ($\lambda \approx 0.15-0.2$), after which they began to decline. The steepest drops were observed for *disgust* and *surprise*, followed by *anger* and *fear*, while *sadness* and *joy* exhibited only a slight decrease at the highest λ . In contrast to all other emotions, lexical density for *joy* rises, starting at around $\lambda = 0.25$ and approaching its maximum at $\lambda = 0.35$. A plausible explanation is that joy functions as a baseline affect in many unsteered generations; steering may therefore reinforce lexical patterns already typical of the model.

Human and model comprehensibility curves align closely—both begin to decline around the same λ threshold, and the relative order across emotions is consistent. Among the surface features, lexical density shows the strongest correspondence with these curves, followed by entropy, while mean word length appears only weakly related.

The standard deviations of the human ratings highlight two critical observations: first, human perceptions of comprehensibility vary considerably, especially at higher steering levels; second, the variability suggests that the comprehensibility of generated text is influenced not only by the intensity of steer-

ing but also by the specific prompt used. In contrast, the model outputs show markedly lower variability, likely reflecting its more consistent—though less nuanced—assessment of text quality across conditions.

However, it is essential to acknowledge that comprehensibility can vary considerably across participants. Given the simple rating instruction (see Section V-D5), human raters may have interpreted the concept differently: some may have focused on surface-level language quality—such as grammatical correctness, fluency, and the absence of awkward repetitions—while others may have evaluated comprehensibility based on content complexity, penalizing texts that were conceptually dense or abstract. Factors like educational background, proficiency in English, and familiarity with the task further influence such individual differences. These subjective variations must be considered when interpreting the human ratings.

In contrast, the model’s comprehensibility evaluation was based on a fixed prompt explicitly defining the rating criteria: “This evaluates the clarity and logical coherence of the sentence’s content. A score of 1 means the sentence is easy to understand and logically consistent. In contrast, a score of 10 indicates that the sentence is confusing, illogical, or nonsensical (e.g., ‘an angry cup of coffee’).” While this ensures consistency in model-based evaluation, it may not fully align with the diverse interpretations held by human raters.

Steering intensities up to $\lambda \approx 0.15$ preserve surface form and clarity, but stronger interventions progressively

degrade text quality. Steering parameters, therefore, need emotion-specific calibration to balance effectiveness and fluency. Moreover, as shown in our earlier study and preliminary experiments, the effect of a given λ -value also depends on the number of layers being steered—meaning that optimal steering strength is shaped by multiple interacting factors, including the model architecture and intervention scope.

VII. LIMITATIONS

Our results were obtained with a single model family and a fixed steering configuration (injection at all layers with a shared λ). Effects may differ for other architectures, instruction-tuned models, or alternative layer-selection schedules. However, in our prior work [2] we also examined steering with only three layers and conducted a probing study across many different layers, finding that style-relevant information is robustly encoded from layer 3 onward and particularly strongly in layers 18–20. We therefore expect similar results across reasonable layer choices, provided that λ is carefully tuned.

Style vectors were derived from the GoEmotions dataset, which has imbalanced category frequencies. We did not reweight categories during vector construction, so the resulting vectors may partly reflect corpus priors. Although GoEmotions is imbalanced, we observed strong steerability for low-frequency classes (e.g., fear) and weak steerability for higher-frequency surprise, suggesting that frequency alone does not determine steerability; future work will examine additional factors that may drive these differences.

Our evaluation uses 19 single-turn prompts and does not test multi-turn dialogues or more complex interactive tasks. The generalization of steering behavior to these settings remains unexamined.

Overall, we believe these limitations do not undermine the main contributions and findings of this work, as the observed effects are consistent with prior results and align with theoretical expectations of activation-space steering.

VIII. CONCLUSION

This study presents the first systematic human evaluation of activation-based steering in language models, offering new insights into both the effectiveness and the practical boundaries of this emerging technique. Using over 7,000 human ratings from 190 participants, we systematically assessed how increasing steering strength (λ) affects perceived emotional intensity and the comprehensibility of generated text.

Our results confirm that activation steering effectively amplifies target emotion signals in text: human raters consistently perceived increasing intensity for five out of six primary emotions as λ increased. At the same time, we identify emotion-specific thresholds—typically around $\lambda \approx 0.15$ —beyond which comprehensibility and linguistic coherence begin to degrade. This highlights the need for emotion-sensitive calibration when applying steering vectors in practical settings.

A second key finding concerns evaluation methodology. Across a wide range of prompts and emotions, model-based quality ratings closely mirror the trends observed in human judgments, especially in the decline of text clarity at higher λ values. This alignment suggests that automatic evaluation can serve as a reliable proxy for perceived quality in future studies, substantially reducing the need for costly manual annotation.

Importantly, we show that perceived emotional intensity is not only targetable but also measurable at scale, even when the cues are subtle or embedded in context. These findings validate activation steering as a viable mechanism for lightweight behavioral control in LLMs and provide clear guidance on balancing strength and stability.

Future work should explore the generalizability of these results beyond emotional tone—for example, to more abstract or stylistic forms—and refine automatic metrics that capture not only comprehensibility, but also nuance, tone, and intent. As LLMs integrate into everyday applications, the ability to steer their behavior in fine-grained and interpretable ways will become increasingly important. Our study provides a foundation for making this control both effective and scalable.

REFERENCES

- [1] D. Hovy and D. Yang, “The importance of modeling social factors of language: Theory and practice,” in *Proc. Conf. North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, 2021, pp. 588–602.
- [2] K. Konen, S. Jentzsch, D. Diallo, P. Schuett, O. Bensch, R. E. Baff, D. Opitz, and T. Hecking, “Style vectors for steering generative large language models,” in *Findings of the Association for Computational Linguistics: EACL 2024*, St. Julian’s, Malta, Mar. 2024, pp. 782–802.
- [3] A. M. Turner, S. Askell, E. Hubinger, et al., “Activation addition: Steering language models via contrasting activations,” arXiv:2308.10248, 2023.
- [4] N. Rimskey, N. Gabrieli, J. Schulz, M. Tong, E. Hubinger, and A. M. Turner, “Steering Llama 2 via contrastive activation addition,” in *Proc. 62nd Annual Meeting ACL*, Bangkok, Thailand, 2024, pp. 15504–15522.
- [5] T. van der Weij, M. Poesio, and N. Schoots, “Extending activation steering to broad skills and multiple behaviours,” arXiv:2403.05767, 2024.
- [6] W. Wang, J. Yang, and W. Peng, “Semantics-adaptive activation intervention for LLMs via dynamic steering vectors,” arXiv:2410.12299, 2024.
- [7] S. Soo, C. Guang, W. Teng, C. Balaganesh, T. Guoxian, and Y. Ming, “Interpretable steering of large language models with feature guided activation additions,” arXiv:2501.09929, 2025.
- [8] P. Q. Da Silva, H. Sethuraman, D. Rajagopal, H. Hajishirzi, and S. Kumar, “Steering off course: Reliability challenges in steering language models,” arXiv:2504.04635, 2025.
- [9] X. Gao, Y. Zhang, and S. Lee, “Structuring latent spaces for stylized response generation,” in *Proc. EMNLP-IJCNLP 2019*, pp. 1900–1910.
- [10] D. Jin, Z. Jin, J. T. Zhou, L. Orii, and P. Szolovits, “Hooks in the headline: Learning to generate headlines with controlled styles,” in *Proc. 58th Annual Meeting ACL*, 2020, pp. 5082–5093.
- [11] Y. Yang, B. Y. Lin, et al., “FUDGE: Controlled text generation with future discriminators,” in *Proc. NAACL*, 2021.
- [12] Z. Sun, H. Wang, Y. Zhang, Y. Wang, and W. Y. Wang, “Expanding style control in text generation via multiple-attribute conditioning,” in *Proc. Conf. Empirical Methods in Natural Language Processing (EMNLP)*, 2022, pp. 1234–1245.
- [13] J. Schon, C. Bruder, M. Liebherr, and S. Theis, “Cleared for takeoff – Artificial intelligence acceptance in air traffic control,” in *Proc. 4th Annu. Symp. Human-Computer Interaction for Work (CHIWORK)*, Amsterdam, Netherlands, 2025, pp. 1–15.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances Neural Information Processing Systems* 30, 2017, pp. 5998–6008.

- [15] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, et al., "Language models are few-shot learners," in *Advances Neural Information Processing Systems* 33, 2020, pp. 1877–1901.
- [16] OpenAI, "GPT-4 technical report," arXiv:2303.08774, 2023.
- [17] P. Schramowski, C. Turan, N. Andersen, C. A. Rothkopf, and K. Kersting, "Large pretrained language models contain human-like biases of what is right and wrong to do," *Nature Machine Intelligence*, vol. 4, no. 3, pp. 258–268, 2022.
- [18] J. Wagner and S. Zarriess, "Does pretrained language models sensitivity to linguistic acceptability influence syntactic generalization?," *Computational Linguistics*, vol. 48, no. 3, pp. 615–650, 2022.
- [19] N. S. Keskar, B. McCann, L. R. Varshney, C. Xiong, and R. Socher, "CTRL: A conditional transformer language model for controllable generation," arXiv:1909.05858, 2019.
- [20] T. Shin, Y. Razeghi, R. L. Logan IV, E. Wallace, and S. Singh, "Auto-Prompt: Eliciting knowledge from language models with automatically generated prompts," in *Proc. Conf. Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 4222–4235.
- [21] B. Lester, R. Al-Rfou, and N. Constant, "The power of scale for parameter-efficient prompt tuning," arXiv:2104.08691, 2021.
- [22] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," in *Proc. 59th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2021, pp. 4582–4597.
- [23] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *Advances Neural Information Processing Systems*, vol. 35, 2022, pp. 24824–24837.
- [24] N. Subramani and J. Andreas, "Extracting latent steering vectors from pretrained language models," in *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 607–621.
- [25] S. Valluripalli, D. Portnoy, and R. Tamang, "A systematic review of empathy in large language models: Performance, measurement, and evaluation," *J. Med. Internet Res.*, vol. 27, pp. e67694, 2025.
- [26] M. W. Kraus, K. Savitsky, E. Corey, and C. C. Y. Lam, "Large language models produce responses perceived to be empathic in human ratings," *Proc. Natl. Acad. Sci. USA*, vol. 122, no. 9, pp. e2321831122, 2025.
- [27] J. Lee, H. J. Kim, S. Kim, and S. Choi, "Human perceptions of empathy in AI: The role of attribution in evaluating empathic responses," *Comput. Human Behav.*, vol. 152, pp. 108138, 2025.
- [28] S. Dathathri, A. Madotto, J. Lin, et al., "Plug and play language models: A simple approach to controlled text generation," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [29] Oregunteng, "Llama-3-8B-Lexi-Uncensored," Hugging Face, 2024.
- [30] D. Demszky, D. Movshovitz-Attias, J. Ko, A. Cowen, G. Nemade, and S. Ravi, "GoEmotions: A dataset of fine-grained emotions," arXiv:2005.00547, 2020.
- [31] P. Ekman and W. V. Friesen, "Constants across cultures in the face and emotion," *Journal of Personality and Social Psychology*, vol. 17, no. 2, pp. 124–129, 1971.
- [32] Google Research, "Ekman mapping for GoEmotions dataset," technical report, 2021.
- [33] J. Hartmann, "Emotion English DistilRoBERTa-base," Hugging Face, 2022.
- [34] P. Ekman, "An argument for basic emotions," *Cognition and Emotion*, vol. 6, no. 3–4, pp. 169–200, 1992.
- [35] J. Hofmann, E. Troiano, K. Sassenberg, and R. Klinger, "Appraisal theories for emotion classification in text," arXiv:2003.14155, 2020.
- [36] E. Troiano, "Where are emotions in text? A human-based and computational investigation of emotion recognition and generation," unpublished, 2023.
- [37] Prolific, "Prolific," London, UK, 2024. [Online]. Available: <https://www.prolific.com> [Accessed: May 9, 2025].
- [38] K. R. Scherer and H. G. Wallbott, "Evidence for universality and cultural variation of differential emotion response patterning," *J. Pers. Soc. Psychol.*, vol. 66, no. 2, pp. 310–321, 1994.
- [39] L. A. M. Bostan and R. Klinger, "An analysis of annotated corpora for emotion classification in text," Otto-Friedrich-Universität, 2018.
- [40] LimeSurvey GmbH, "LimeSurvey: An open source survey tool," Hamburg, Germany, 2024. [Online]. Available: <https://www.limesurvey.org> [Accessed: May 9, 2025].
- [41] T. K. Koo and M. Y. Li, "A guideline of selecting and reporting intraclass correlation coefficients for reliability research," *J. Chiropr. Med.*, vol. 15, no. 2, pp. 155–163, Jun. 2016, doi: 10.1016/j.jcm.2016.02.012.
- [42] J. A. Russell, "A circumplex model of affect," *J. Pers. Soc. Psychol.*, vol. 39, no. 6, pp. 1161–1178, 1980.
- [43] D. Watson and A. Tellegen, "Development and validation of brief measures of positive and negative affect: The PANAS scales," *J. Pers. Soc. Psychol.*, vol. 54, no. 6, pp. 1063–1070, 1988.
- [44] J. T. Larsen, A. P. McGraw, and J. T. Cacioppo, "Can people feel happy and sad at the same time?," *J. Pers. Soc. Psychol.*, vol. 81, no. 4, pp. 684–696, 2001.
- [45] L. F. Barrett, "Feelings or words? Understanding the content in self-report ratings of experienced emotion," *J. Pers. Soc. Psychol.*, vol. 87, no. 2, pp. 266–281, 2004.
- [46] K. R. Scherer, "The dynamic architecture of emotion: Evidence for the component process model," *Cognition and Emotion*, vol. 23, no. 7, pp. 1307–1351, 2009.

ACKNOWLEDGMENT

The authors used ChatGPT, DeepL and Grammarly to assist with grammar refinement, editing, and LaTeX code formatting during manuscript preparation.

APPENDIX A PROMPTS

- 1) Come up with a fairy tale.
- 2) Comment on a review of a business written by a customer.
- 3) Compare the cultural value of theatres and cinemas.
- 4) Compare the quality of coffee and tea.
- 5) Describe a fictional character.
- 6) Describe an idea of what a fulfilling and meaningful life should be like.
- 7) Describe a piece of art.
- 8) Describe a typical day for a family in the 1980s.
- 9) Describe an imaginary board meeting at which the concept of PopTarts is proposed.
- 10) Pick a colour. Write about that colour as if it was a person.
- 11) Rewrite a famous film scene from the perspective of a fish.
- 12) Summarise the plot of 'The Lord of the Rings' by J.R.R. Tolkien.
- 13) Write a couple of lyrics about society.
- 14) Write a letter to your parents.
- 15) Write a scene between a mentor and a mentee in which the mentee must stand up for himself or herself because the mentor doesn't give him or her enough credit.
- 16) Write a story about an encounter of a boy and a girl.
- 17) Write a story that involves magic.
- 18) Write a story that starts with 'The moon was bright. Not that the person could see it.'
- 19) Write an opinion on the role of government in society.

APPENDIX B EXAMPLE

a: Example steered text: Describe a typical day for a family in the 1980s.

Unsteered

The 1980s — a time of neon-colored clothing, big hair, and radical music. For a typical family in the 1980s, a day would likely look like this:

Morning: The family wakes up early, with the kids rushing to get ready for school. Mom and Dad are already up, making breakfast in the kitchen. The smell of bacon and toast fills the

air as they scramble to get everyone out the door on time. The kids grab their backpacks and head out the door, while Mom and Dad exchange a quick kiss and a “have a good day” before sending them off to school.

Afternoon: Mom heads to the grocery store to pick up some essentials for dinner. She navigates the aisles, carefully selecting the freshest produce and meats for tonight’s meal. Meanwhile, Dad heads to work, commuting to his office job in a suit and tie. He spends the day typing away on his computer, attending meetings, and making deals.

Evening: The family gathers around the dinner table, where Mom has prepared a hearty meal of spaghetti and meatballs. They share stories about their day, from the kids’ adventures at school to Dad’s triumphs at work. After dinner, the family watches TV together, laughing at the antics of *The Cosby Show* or *Family Ties*. They might even play a board game or two, like Monopoly or Scrabble.

Night: As the evening winds down, the family heads to bed, exhausted from the day’s activities. Mom and Dad tuck the kids in, giving them a goodnight kiss and a reminder to brush their teeth. As they drift off to sleep, they dream of the next day’s adventures.

Steered: 0.2 Anger

A typical day for a family in the 1980s was vastly different from today. Here’s a glimpse into what a typical day might look like:

The family wakes up early, around 6:00 AM. The father, a working-class man, gets up first and starts getting ready for work. He puts on his suit and tie, grabs his briefcase, and heads out the door to catch the bus to work.

The mother, a homemaker, stays behind to get the kids ready for school. She makes breakfast, usually something simple like cereal or toast, and gets the kids dressed and packed off to school.

The kids, two boys and a girl, all under the age of 12, are forced to wear the same clothes every day because their parents can’t afford to buy them new ones. They wear hand-me-downs from older siblings or, worse, repurposed adult clothing.

The mother then spends the rest of the day cleaning the house, doing laundry, and cooking dinner. She is forced to do all of this by hand, because she can’t afford a washing machine or a dishwasher.

The father comes home from work at 5:00 PM, exhausted and stressed from dealing with his boss and coworkers all day. He spends the evening watching TV and drinking beer, while the mother continues with housework.

The kids do their homework in front of the TV because they can’t afford a desk or computer. They use outdated textbooks passed down from their parents.

Steered: 0.3 Joy

The 1980s! A time of big hair, neon colors, and lots of fun! Here’s a glimpse into a typical day for a family in the 80s:

Morning: The family wakes up to the sound of the alarm clock, which is probably a cassette player blasting a favorite song. Mom and Dad get ready for the day, and the kids help each other get dressed and enjoy a quick breakfast.

Afternoon: The family heads out for a fun day! They might go to the park, play games, and soak up the sunshine. The kids enjoy their new toys and the excitement of being outdoors.

Evening: After their outing, the family returns home and enjoys a delicious dinner together. They might watch a movie or play a few more games. With snacks and laughter, the family winds down their happy day together.

Steered: 0.35 Sadness

The 1980s — a time of neon-colored clothing, big hair, and a lack of social media. But behind the bright colors was a quiet struggle. Here’s how a typical day might have looked for a struggling family:

It’s Monday morning. The father, a working-class man, wakes early and dresses in worn-out jeans and a faded shirt. He kisses his wife goodbye and leaves for the factory, where he works long hours for little pay.

The mother, a stay-at-home mom, gets the children ready for school. They wear the same old clothes every day because the family can’t afford new ones.

The kids, two boys and a girl, feel sad and distant. They miss their father and don’t understand why he always has to leave. They feel forgotten.

The mother tries to comfort them, but she’s overwhelmed. She spends the day cleaning, worrying, and trying to hold things together.

At school, the children feel isolated. They don’t want to be there, but they don’t want to be home either. They feel unloved in a world that seems unfair.

APPENDIX C DESCRIPTIVE STATISTICS

TABLE 4. Means and standard deviations of perceived emotional intensity and text comprehensibility of all steering conditions

Target emotion	Anger								Disgust								Fear							
λ	0.00	0.05	0.10	0.15	0.20	0.25	0.30	0.35	0.00	0.05	0.10	0.15	0.20	0.25	0.30	0.35	0.00	0.05	0.10	0.15	0.20	0.25	0.30	0.35
	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	
Emotional intensity																								
Anger	0.91 (1.82)	1.11 (2.03)	1.97 (2.45)	1.63 (2.26)	2.34 (2.58)	2.95 (2.69)	3.36 (2.83)	3.10 (2.76)	1.00 (1.91)	1.32 (2.18)	1.51 (2.34)	1.81 (2.45)	2.32 (2.65)	2.97 (2.69)	3.09 (2.69)	2.50 (2.70)	1.04 (1.90)	1.03 (1.95)	1.28 (2.08)	1.17 (1.98)	1.25 (1.93)	1.57 (2.20)	1.35 (2.10)	1.03 (1.78)
Disgust	0.70 (1.55)	0.88 (1.81)	1.05 (1.92)	1.09 (1.95)	1.58 (2.21)	2.57 (2.60)	2.57 (2.66)	2.68 (2.75)	0.75 (1.74)	0.78 (1.73)	1.32 (2.19)	1.98 (2.50)	3.21 (2.76)	4.54 (2.52)	4.60 (2.64)	4.86 (2.51)	0.74 (1.63)	0.79 (1.70)	0.96 (1.70)	1.02 (1.82)	0.93 (1.70)	1.53 (2.22)	1.20 (1.95)	1.18 (1.90)
Fear	0.98 (1.69)	1.15 (1.93)	1.15 (1.84)	1.43 (2.09)	1.38 (2.06)	1.47 (2.15)	1.51 (2.30)	1.51 (2.17)	0.93 (1.67)	1.13 (1.81)	1.76 (2.29)	1.48 (2.18)	1.61 (2.31)	1.76 (2.29)	0.98 (1.72)	0.92 (1.83)	0.85 (1.55)	1.35 (2.15)	2.29 (2.61)	2.17 (2.46)	2.53 (2.55)	4.03 (2.54)	4.50 (2.48)	4.80 (2.31)
Joy	2.43 (2.50)	1.85 (2.26)	1.76 (2.13)	1.63 (2.28)	1.00 (1.86)	0.66 (1.58)	0.62 (1.57)	0.59 (1.37)	2.49 (2.51)	1.74 (2.27)	1.51 (2.17)	1.01 (1.86)	0.44 (1.23)	0.26 (0.97)	0.16 (0.74)	0.06 (0.25)	2.39 (2.43)	1.68 (2.22)	1.20 (2.00)	1.20 (1.97)	1.05 (1.92)	0.41 (1.18)	0.24 (0.82)	0.19 (0.83)
Sadness	1.24 (1.96)	1.42 (2.08)	1.49 (2.09)	1.33 (2.08)	1.67 (2.22)	1.98 (2.48)	1.84 (2.41)	1.73 (2.33)	1.20 (1.97)	1.78 (2.28)	1.51 (2.21)	2.17 (2.42)	1.80 (2.43)	2.12 (2.48)	1.74 (2.40)	1.20 (2.01)	1.18 (1.89)	1.65 (2.29)	1.65 (2.25)	1.69 (2.16)	1.91 (2.25)	2.17 (2.37)	2.28 (2.48)	2.32 (2.38)
Surprise	1.56 (2.02)	1.73 (2.09)	1.41 (2.03)	1.38 (2.02)	1.40 (2.03)	1.25 (1.96)	1.07 (1.84)	1.19 (1.88)	1.45 (1.94)	1.47 (1.99)	1.44 (1.97)	1.36 (1.88)	1.44 (2.08)	1.36 (2.11)	1.11 (1.86)	0.78 (1.62)	1.50 (2.01)	1.23 (1.85)	1.48 (2.06)	1.53 (1.99)	1.61 (2.08)	1.63 (2.15)	1.42 (2.09)	1.45 (1.98)
Comprehensibility	5.91 (1.33)	5.90 (1.39)	5.79 (1.51)	5.66 (1.68)	5.58 (1.62)	5.16 (1.91)	5.00 (2.02)	4.10 (2.33)	5.83 (1.52)	5.86 (1.47)	5.87 (1.49)	5.70 (1.72)	4.97 (2.05)	4.66 (2.07)	3.47 (2.41)	2.77 (2.47)	5.79 (1.55)	5.78 (1.66)	5.84 (1.43)	5.69 (1.61)	5.43 (1.91)	5.22 (1.95)	4.90 (2.01)	4.11 (2.30)

Target emotion	Joy								Sadness								Surprise							
λ	0.00	0.05	0.10	0.15	0.20	0.25	0.30	0.35	0.00	0.05	0.10	0.15	0.20	0.25	0.30	0.35	0.00	0.05	0.10	0.15	0.20	0.25	0.30	0.35
	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	M (SD)	
Emotional intensity																								
Anger	0.95 (1.87)	0.43 (1.08)	0.31 (0.92)	0.41 (1.06)	0.28 (0.87)	0.23 (0.92)	0.13 (0.60)	0.16 (0.64)	0.96 (1.87)	0.84 (1.67)	0.85 (1.65)	0.94 (1.69)	1.08 (1.80)	0.92 (1.79)	1.00 (1.76)	0.87 (1.57)	1.08 (1.97)	0.96 (1.91)	0.73 (1.55)	0.88 (1.77)	0.57 (1.41)	0.66 (1.51)	0.55 (1.31)	0.48 (1.37)
Disgust	0.72 (1.70)	0.33 (0.91)	0.19 (0.65)	0.21 (0.76)	0.18 (0.65)	0.19 (0.78)	0.14 (0.71)	0.15 (0.66)	0.73 (1.64)	0.54 (1.33)	0.68 (1.58)	0.43 (1.06)	0.61 (1.26)	0.71 (1.53)	0.66 (1.43)	0.71 (1.45)	0.76 (1.66)	0.53 (1.41)	0.56 (1.39)	0.50 (1.34)	0.56 (1.35)	0.54 (1.40)	0.41 (1.03)	0.44 (1.37)
Fear	1.00 (1.75)	0.76 (1.63)	0.41 (1.08)	0.53 (1.25)	0.35 (1.05)	0.32 (1.07)	0.24 (0.95)	0.20 (0.77)	0.99 (1.74)	1.20 (1.77)	1.19 (1.92)	1.26 (1.83)	1.40 (1.99)	1.34 (1.94)	1.50 (2.12)	1.77 (2.24)	0.91 (1.65)	0.98 (1.67)	0.93 (1.74)	0.99 (1.76)	1.08 (1.95)	0.73 (1.47)	0.64 (1.43)	0.62 (1.48)
Joy	2.57 (2.57)	3.46 (2.63)	3.57 (2.58)	3.81 (2.51)	4.11 (2.60)	4.28 (2.43)	4.63 (2.38)	3.83 (2.62)	2.48 (2.58)	2.07 (2.31)	1.78 (2.27)	1.87 (2.26)	1.25 (1.88)	1.28 (1.92)	1.03 (1.81)	0.81 (1.56)	2.54 (2.42)	2.34 (2.51)	2.46 (2.51)	2.04 (2.35)	1.82 (2.26)	1.71 (2.32)	1.50 (2.16)	0.90 (1.81)
Sadness	1.12 (1.91)	1.15 (1.93)	0.65 (1.42)	0.71 (1.47)	0.45 (1.22)	0.44 (1.20)	0.36 (1.19)	0.28 (0.94)	1.35 (2.05)	1.34 (1.92)	1.96 (2.49)	2.30 (2.45)	2.92 (2.55)	3.23 (2.58)	3.68 (2.56)	3.54 (2.74)	1.13 (1.83)	0.94 (1.68)	0.97 (1.68)	1.18 (1.83)	1.05 (1.77)	0.95 (1.67)	0.64 (1.44)	0.51 (1.26)
Surprise	1.47 (1.95)	1.38 (1.96)	1.58 (2.12)	1.43 (2.02)	1.23 (1.84)	1.17 (1.80)	1.16 (1.95)	0.91 (1.80)	1.56 (2.07)	1.46 (1.99)	1.22 (1.92)	1.48 (1.95)	1.27 (1.86)	1.29 (1.90)	1.01 (1.73)	1.25 (1.95)	1.47 (2.04)	1.77 (2.12)	1.72 (2.11)	1.67 (2.17)	1.79 (2.15)	1.63 (2.12)	1.66 (2.10)	1.18 (1.91)
Comprehensibility	6.00 (1.26)	6.03 (1.35)	6.06 (1.26)	6.10 (1.19)	6.02 (1.32)	5.97 (1.46)	6.11 (1.36)	5.46 (2.05)	5.86 (1.43)	5.96 (1.45)	5.84 (1.50)	5.99 (1.21)	5.87 (1.58)	5.87 (1.42)	5.73 (1.58)	5.22 (1.97)	5.88 (1.39)	5.85 (1.43)	5.97 (1.37)	5.80 (1.47)	5.60 (1.72)	4.71 (2.24)	4.10 (2.41)	3.17 (2.69)

DIAOULÉ DIALLO received the B.A. degree in philosophy-neuroscience-cognition from Otto-von-Guericke-University Magdeburg, and the M.Sc. degree in intelligent systems from Bielefeld University. Since 2021, he has been a researcher with the German Aerospace Center (DLR), and a Ph.D. candidate at the University of Bonn.

His research focuses on complex network science and epidemiology, in particular on risk assessment in temporal networks to support data-driven and machine learning-based intervention strategies during pandemics. His second line of research investigates activation-based steering of large language models, aiming to improve their explainability and controllability through representation-level interventions.

KATHARINA DWORATZYK received the B.A. degree in philosophy from Rheinische Friedrich-Wilhelms-University Bonn and the M.Sc. degree in cognitive science from Osnabrück University. Since 2019, she has been a researcher with the German Aerospace Center (DLR).

Her research focuses on evaluating the effectiveness of interventions such as community-building and understanding context-specific user requirements, with the goal of informing the development of human-centered solutions for software engineering education and socio-technical systems.

SOPHIE JENTZSCH received the B.Sc. and M.Sc. degrees in psychology in computer science from Technische Universität Darmstadt. Since 2019, she has been a researcher with the German Aerospace Center (DLR), and a Ph.D. candidate at the University of Darmstadt.

Her research focuses on human-centered artificial intelligence and the societal impact of large language models. In particular, she investigates the role of humor and moral reasoning in human-AI interaction, and contributes to interdisciplinary approaches for understanding and guiding AI behavior. She is actively involved in science communication through podcasts and public outreach formats.

PEER SCHÜTT received a B.Sc. and M.Sc. degree in computer science from the University of Bonn. Since 2022, he has been a researcher with the German Aerospace Center (DLR), and a Ph.D. candidate at the University of Cologne.

His research focuses on anomaly detection in multispectral imagery, focusing on detecting hazardous substances using remote sensing methods. Furthermore, he investigates the application of large language models in diverse research domains and administrative workflows to enhance LLM usability in research institutions.

SABINE THEIS received her B.Sc. from the University of Applied Sciences Aachen/Maastricht and M.Sc. in Information Sciences from the University of Amsterdam. She earned her Ph.D. in industrial engineering and ergonomics from RWTH Aachen University in 2019 with a thesis on ergonomic visualization of personal health data. She worked at Fraunhofer FKIE, investigating workload in AR/VR environments, and as a postdoctoral researcher at RWTH Aachen University on data visualization projects. Currently, she leads the Human Factors in Software Engineering group at the German Aerospace Center (DLR). Her research focuses on human-AI interaction, and human-centered data visualization systems.

TOBIAS HECKING received the B.Sc. degree in computer science from Westfälische Hochschule and the M.Sc. degree in computer science from the University of Duisburg-Essen. He received his Ph.D. in computer science from the University of Duisburg-Essen in 2016. Since 2020, he has been head of the Intelligent Software Systems group at the German Aerospace Center (DLR), Institute of Software Technology.

His research focuses on theories, methods, and tools of artificial intelligence for the development of data-driven intelligent software systems. The focus here is on methods of network analysis, natural language processing and machine learning. The main areas of application are knowledge-based software systems, analysis of complex processes and intelligent user interfaces.

...