

CreativityNeuro: Steering Language Model Weights to Improve Divergent Thinking

Anonymous authors
Paper under double-blind review

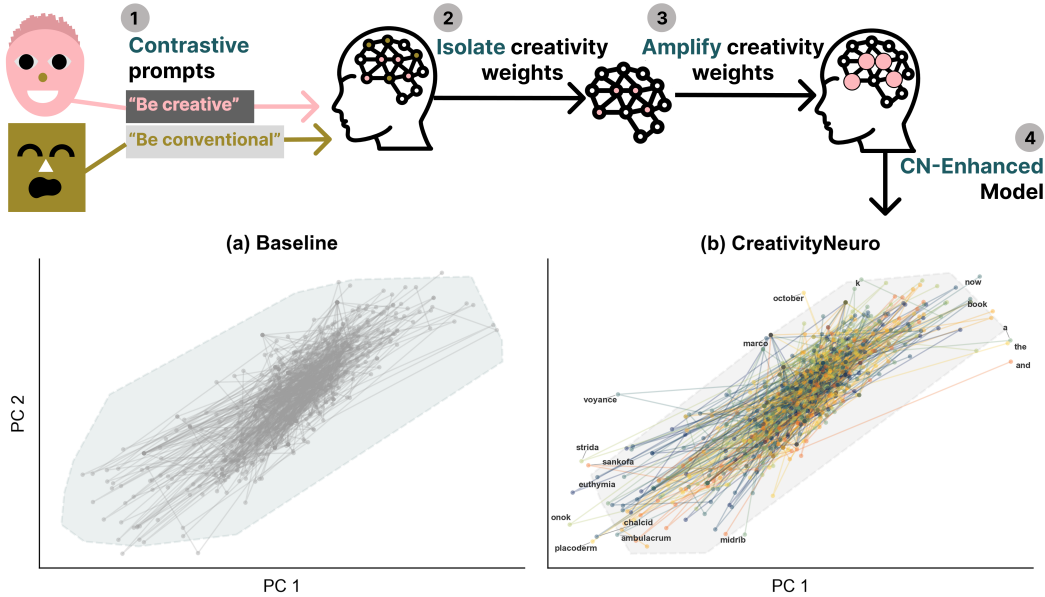


Figure 1: **CreativityNeuro (CN) pipeline.** Given a pair of contrastive creative prompts, CN computes parameter importance scores, selects a sparse subset of creativity-relevant parameters, and applies a scaled weight perturbation—without requiring behavioral datasets or gradient-based finetuning. CN improves divergent thinking across various tasks. Subplot (b) visualizes CN *thinking outside of the “box”* (i.e., the convex hull of baseline DAT responses), despite baseline responses falling within CN’s convex hull in subplot (a).

Abstract

1 We introduce CreativityNeuro (CN), a data-free method for enhancing cre-
2 ative behavior in LLMs via contrastive weight steering. We evaluate CN on
3 various creativity assessments and report three main findings: **(1)** On the
4 Divergent Association Task (DAT), CN improves performance by up to 14
5 human percentile points, with gains that are consistent across semantically
6 diverse prompts. **(2)** In a large-scale human evaluation (N=720) on the
7 Alternative Uses Test (AUT) and the Task Task (TT), CN-enhanced models
8 achieve significant improvements in originality, surprise, and creativity,
9 demonstrating transfer to more complex, open-ended settings. Compared
10 to contrastive activation steering (CAA), we find that while activation steer-
11 ing can match CN on the DAT, it does not achieve consistent improvements
12 on held-out tasks (AUT, TT). **(3)** CN-enhanced models reduce factual reason-
13 ing scores on MMLU by 1–11%, and attempts to explicitly preserve such
14 capabilities by incorporating MMLU prompts into prompt contrast sets
15 reverse nearly all divergent thinking (DT) gains, providing evidence that
16 DT and factual reasoning rely on functionally entangled weights.

17 1 Introduction

18 Recent advances in large language models (LLMs) have renewed interest in a longstand-
19 ing question: *how can we understand and enhance creativity in intelligent systems?* (Boden,
20 2004). While this question has deep roots in cognitive science (Quetelet, 1842; Galton, 1870;
21 Hadamard, 1954; Guilford, 1956; Mednick, 1962; Koestler, 1964; Simonton, 2004; Dietrich,
22 Arne, 2004; Fauconnier & Turner, 2008; Rothenberg, 2014), it is now increasingly studied in
23 the context of large-scale generative models (Maher, 2010; Varshney, 2019; Schapiro et al.,
24 2025). Recent work has begun to assess the capacity for LLMs to engage in creative and
25 open-ended tasks (Si et al., 2024; 2025; Sanyal et al., 2025; Bellemare-Pepin et al., 2024; Wang
26 et al., 2025; 2024), where a recurring issue has surfaced: models tend to consistently generate
27 similar responses to open-ended questions, in what has been termed the *artificial hivemind*
28 effect (Jiang et al., 2025).

29 Within the creativity literature, a common distinction is made between *divergent thinking*
30 (DT), the capacity to generate multiple diverse solutions to a problem, and *convergent*
31 *thinking* (CT), the ability to find a single correct solution that unifies multiple diverse
32 stimuli (Dietrich, 2019; Guilford, 1956). Studying ways to enhance DT offers a promising
33 pathway to encourage greater diversity and novelty in model responses, combating the
34 homogenization issues that have emerged thus far. Thus, in this work, we introduce a
35 weight-space steering method that improves DT in large language models. Our method
36 outperforms prior approaches—including decoding, prompting, and activation steering—
37 and generalizes better to unseen tasks, without requiring behavioral data or gradient-based
38 fine-tuning. In detail, our main contributions are as follows:

- 39 1. We introduce **CreativityNeuro (CN)**, a data-free method for steering creative be-
40 havior. **CN improves divergent thinking** on the Divergent Association Task (DAT),
41 outperforming baselines such as prompting, activation steering, and decoding
42 baselines.
- 43 2. We conduct a large-scale human evaluation on the Alternative Uses Test (AUT) and
44 Task Task (TT) and find that **CN-enhanced models improve originality, surprise,**
45 **and creativity on the AUT and TT.** Despite performing comparably to CN on the
46 DAT, **activation steering transfers poorly to the AUT and TT.**
- 47 3. We investigate the impact of CN on factual reasoning and find that improvements
48 in divergent thinking come at the cost of reduced MMLU performance. Attempts to
49 preserve factual reasoning ability reverse nearly all DAT gains, providing **evidence**
50 **that divergent thinking and factual reasoning rely on functionally entangled**
51 **model weights.**

52 **Outline** In Section 2, we discuss related work and introduce our method. Next, in Section 3
53 and Section 4, we perform experiments on the DAT, AUT, and TT. Then, in Section 5, we
54 study the tradeoff between divergent thinking and factual reasoning. Lastly, we discuss
55 limitations, future work, and conclude in Section 6.

56 2 Background and Method

57 2.1 Related Work

58 **Evaluating the creativity of LLMs.** Previous work has evaluated LLMs on divergent
59 creativity assessments—including the DAT (Olson et al., 2021), AUT (Guilford, 1956), the
60 Task Task (Chu et al., 2024)—and in various real-world settings like scientific ideation (Si
61 et al., 2024; 2025) and open-ended user queries (Jiang et al., 2025). On the DAT, LLMs can
62 achieve scores well into the 90th percentile of humans (Bellemare-Pepin et al., 2024; Wang
63 et al., 2025), whereas Stevenson et al. (2022) studied GPT-3 on the AUT and concluded
64 that humans exhibited greater creativity, with model responses showing weaker originality.
65 Lastly, Chu et al. (2024) found that model-generated goals on the Task Task achieved similar
66 creativity ratings as human-generated goals, as assessed by a large panel of human raters.

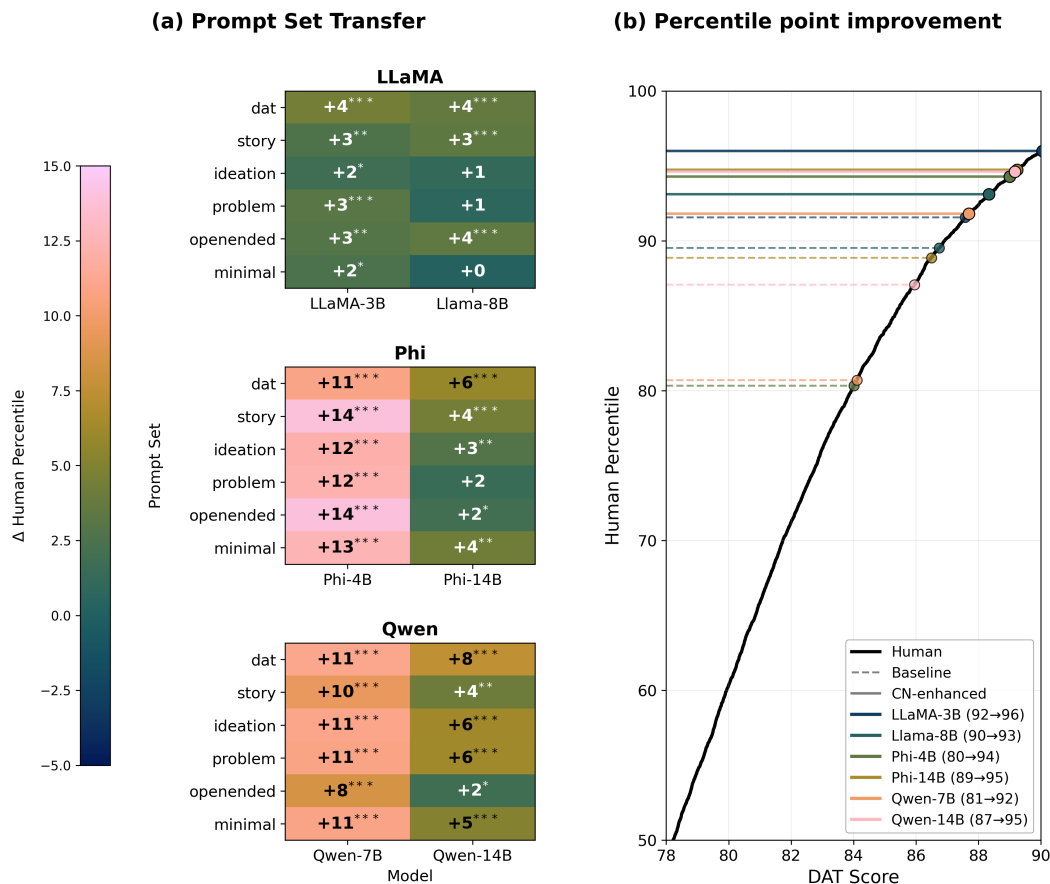


Figure 2: **CreativityNeuro (CN) improves divergent thinking across models and prompt sets.** Given a human reference distribution (Wang et al., 2025) ($N = 9,297$, $\mu = 78.26$, $\sigma = 6.73$), we report: (a) Heatmap showing human percentile improvement ($\Delta\%$ ile) for CN-enhanced models, with statistical significance ($p < 0.05$) at each of the temperatures tested (0.9, 1.0, 1.2) denoted by an asterisk. (b) CDF showing CN-enhanced scores for the best performing prompt set.

85 at layer ℓ for token t in prompt b . We compute importance scores on creative \mathcal{P}^{cre} and
 86 non-creative prompts $\mathcal{P}^{\text{non-cre}}$ (representative examples are given in Table 1; all six prompt
 87 sets are given in Table 2). Then, we isolate creativity-specific parameters by selecting the
 88 top ρ percent of weights ranked by creative importance that do not also appear in the
 89 top ρ percent for non-creative prompts. This set difference operation ($C_\ell \setminus N_\ell$) ensures we
 90 identify parameters uniquely associated with creative behavior. Lastly, at inference time, we
 91 multiply weights in $C_\ell \setminus N_\ell$ by a scaling factor $(1 + \alpha)$. Like Christ et al. (2025), this method
 92 is data- and backpropagation-free. The hyperparameters ρ (importance threshold) and α
 93 (scaling strength) control the intervention’s scope and intensity. The full procedure is given
 94 in Algorithm 1.

95 3 Experiments on the Divergent Associates Task

96 We test instruct-tuned models across three open-weight model families (Phi, Llama, Qwen),
 97 totaling six models at 3B, 4B, 7B, 8B, and 14B sizes: LLaMA (3.2-3B-Instruct, 3.1-8B-Instruct)
 98 (Grattafiori et al., 2024), Qwen-2.5 (7B-Instruct, 14B-Instruct) (Yang et al., 2025), and Phi
 99 (3.5-mini-Instruct (4B), 3-medium-4k-Instruct (14B)) from Microsoft (Abdin et al., 2024). Full
 100 hyperparameter sweep details for CreativityNeuro parameters are available in Section D.

101 3.1 The Divergent Associates Task (DAT)

102 The DAT asks participants to generate N words that are as semantically distant from each
103 other as possible (Olson et al., 2021). Given a set of N words $W := \{w_1, w_2, \dots, w_N\}$ with
104 corresponding GloVe embeddings $V := \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\} \subseteq \mathbb{R}^{300}$, the DAT score is the
105 average pairwise semantic distance among all distinct pairs of those N words:

$$\text{DAT}(W) := \frac{100}{N(N-1)} \sum_{i \neq j}^N (1 - \cos(\mathbf{v}_i, \mathbf{v}_j)) \quad (1)$$

106 Following Olson et al. (2021), we use the 840B-token GloVe embeddings (Pennington et al.,
107 2014) as our semantic space, providing consistent distance measurements independent of
108 response pool demographics. Typically, a participant is asked to name $N = 10$ words, and
109 the first 7 valid words are kept (Olson et al., 2021). The prompts used to administer the DAT
110 are given in Section E.

111 3.2 Baselines

112 Existing studies have found that temperature-scaling and prompting can influence DAT
113 scores (Wang et al., 2025; Bellemare-Pepin et al., 2024). To ensure the CN intervention
114 leads to a meaningful improvement over such techniques, we compare against a broad
115 set of baselines, including prompting; varying decoding parameters such as top-p nucleus
116 sampling, top-k sampling, temperature, and repetition penalty; as well as activation steering
117 via contrastive activation addition (CAA; Panickssery et al. (2024)), which injects a steering
118 vector $\mathbf{v}_\ell = \bar{\mathbf{h}}_\ell^+ - \bar{\mathbf{h}}_\ell^-$ into the residual stream during decoding. For activation steering,
119 contrast pairs are obtained from top- vs. bottom-quartile DAT responses by score, creating a
120 “divergent thinking” direction in the residual stream.¹ Full decoding strategy and activation
121 steering hyperparameter settings are given in Section B. All settings are evaluated across all
122 six models at $T \in \{0.9, 1.0, 1.2\}$ until $N=120$ valid DAT responses are obtained.

123 3.3 Results on the Divergent Associates Task

124 CN improves DAT performance across all six models and prompt sets (Figure 2), outper-
125 forming all sampling-based baselines (Figure 3). Panel (a) of Figure 2 shows Δ Percentile
126 for the best (ρ, α) per model-prompt combination: while the DAT prompt set produces the
127 most consistent gains, non-DAT prompt sets also yield statistically significant improvements,
128 suggesting CN is able to identify weights controlling divergent thinking behavior, rather
129 than localizing DAT-specific task knowledge.

130 CN (94.1 avg) slightly outperforms activation steering (93.9 avg) on DAT percentile (Fig-
131 ure 3); however, activation steering requires scored DAT responses to construct its steering
132 vector, while CN uses only creative and non-creative prompts with no generation or scor-
133 ing. Prompt-only activation steering (omitted from the table; see footnote) averaged only
134 87.8, comparable to prompting (87.9), suggesting that the behavioral signal is essential for
135 activation steering to be competitive, whereas CN extracts a stronger signal from prompts
136 alone.

CreativityNeuro improves DAT scores across all six models and prompt sets, outperforming
prompting, sampling-parameter, and activation steering baselines.

137

¹We also tested a prompt-only CAA variant using forward-pass activations from the same creative vs. non-creative prompt sets as CN, but it underperformed the behavioral-data variant on all models and is omitted for clarity.

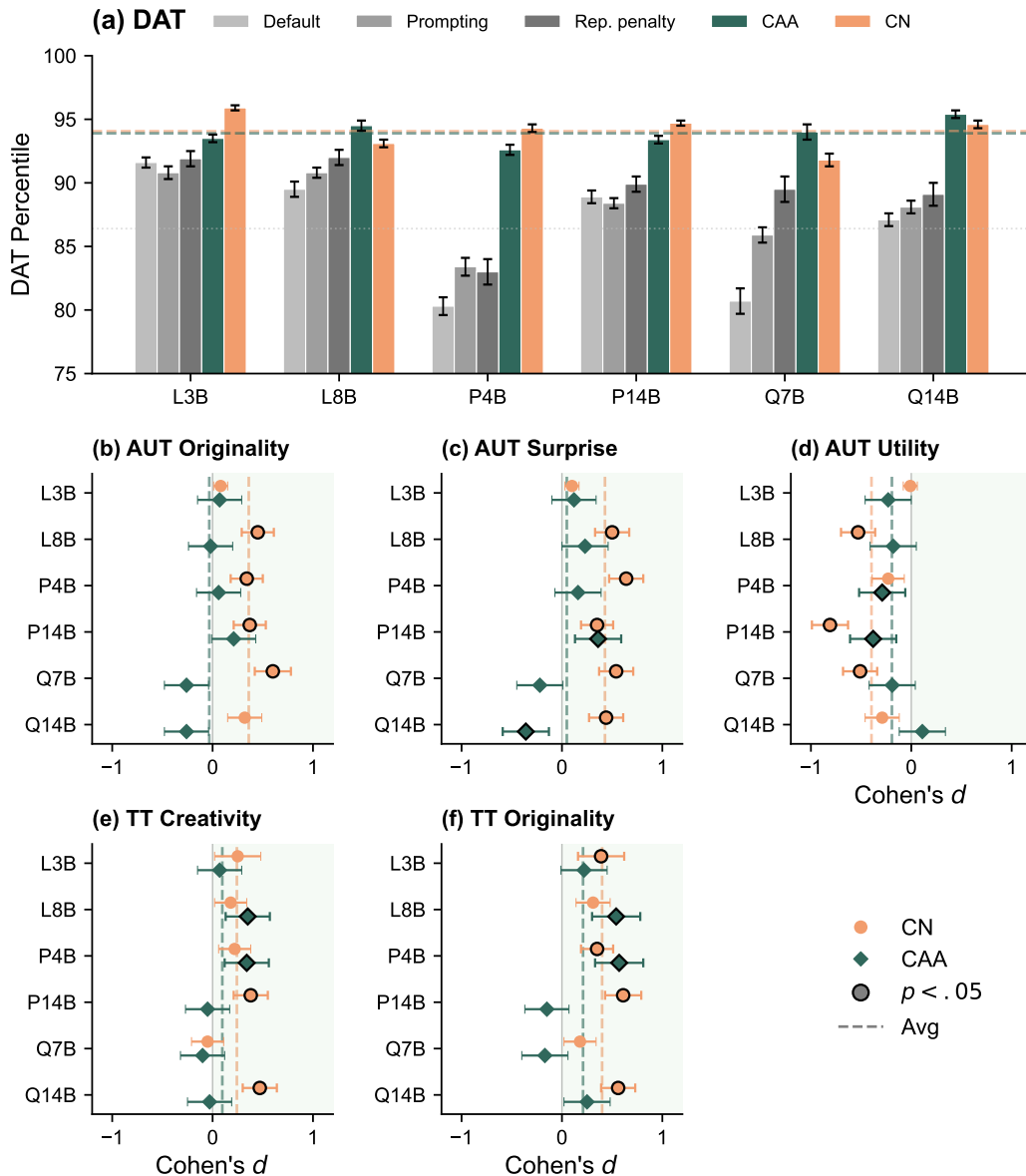


Figure 3: Results on the DAT, AUT, and TT, as well as comparison against baselines. (a) DAT human percentile (\pm SEM) averaged across $T \in \{0.9, 1.0, 1.2\}$. (b–f) Cohen’s d (\pm SE) from intra-participant z-scored human ratings on the AUT and TT. Black outlines indicate $p < .05$, and green shading marks the positive-effect region.

138 4 Generalization to the Alternative Uses Test and the Task Task

139 4.1 The Alternative Uses Test (AUT) and the Task Task (TT)

140 The Alternative Uses Test is a divergent thinking assessment introduced by Guilford (1956)
 141 that asks participants to generate creative uses for a common everyday object (e.g., “List
 142 alternative uses for a brick”). Unlike the DAT, which measures semantic distance among
 143 single words, the AUT requires participants to produce functionally meaningful ideas, each
 144 of which may span multiple words or sentences. We administer the AUT using a standard
 145 set of objects commonly used in the creativity literature: *brick*, *paperclip* and *fork*. Following
 146 Stevenson et al. (2022), uses are scored on *originality*, *surprise*, and *utility*. Meanwhile, the

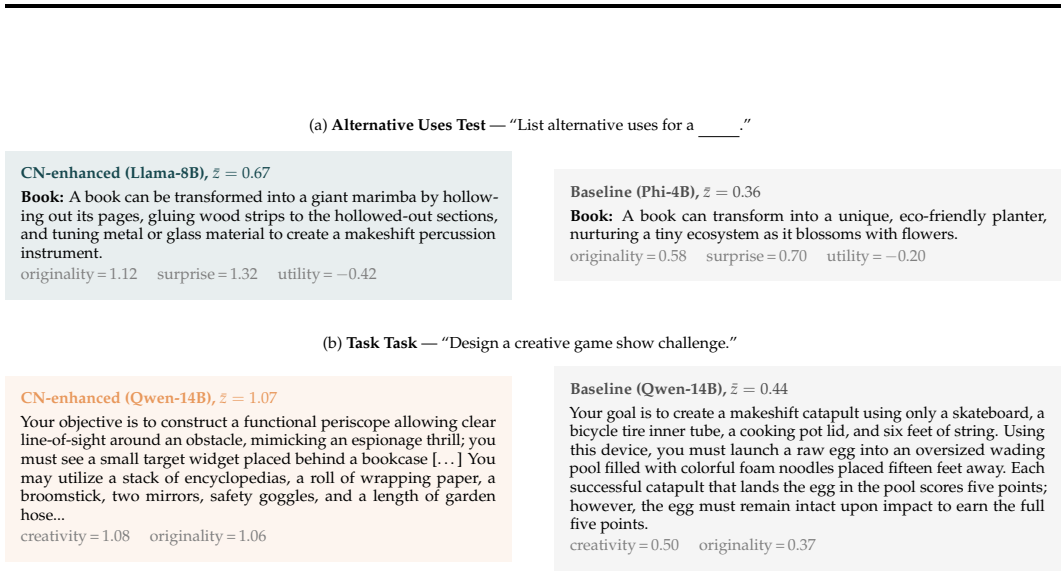


Figure 4: **Top-rated CN-enhanced stimuli vs. baseline generations from the same model.** Intra-participant z-scores averaged across raters ($N=40$ per cell). (a) CN-enhanced AUT responses tend to score higher on originality and surprise while sometimes sacrificing utility; (b) CN-enhanced Task Task challenges layer more constraints, props, and absurdist elements.

147 Task Task evaluates the ability to generate novel challenges or goals that themselves expect
 148 novel solutions (Chu et al., 2024). Participants are asked to design creative game show
 149 challenges that would be fun to attempt, entertaining to watch, and difficult enough to be
 150 interesting. Following Chu et al. (2024), we evaluate responses on creativity and originality.²
 151 Full prompts used to administer the AUT and TT are given in Section E.

152 **4.2 Human Experiment Design**

153 We evaluate AUT and TT stimuli via human ratings on Prolific. For each model, we select
 154 the CN configuration (ρ, α , prompt set) that produced the largest statistically significant
 155 DAT improvement in Section 3. Then, for each task, we sample 40 stimuli (20 baseline,
 156 20 CN-enhanced) at temperature $T = 1.0$, top- $k = 0$, and top- $p = 1.0$. All studies uses a
 157 between-subjects design, where every participant rates 10 stimuli (5 baseline, 5 creative) in
 158 randomized order with condition labels hidden. We ensure balanced allocation via auto-
 159 mated participant-to-slot assignment so that each stimulus receives exactly 10 independent
 160 ratings. Moreover, on both the AUT and TT, ratings for each dimension are given on con-
 161 tinuous 0–100 sliders, and to control for individual differences in scale usage, we compute
 162 intra-participant z-scores: for each participant i and dimension d , $z_{i,d,s} = (r_{i,d,s} - \bar{r}_{i,d}) / \sigma_{i,d}$,
 163 where $\bar{r}_{i,d}$ and $\sigma_{i,d}$ are computed across all stimuli that participant rated on that dimension.
 164 We recruit 30 participants per (model, task, method) triple, totaling $N=720$ total human re-
 165 viewers. Effect sizes and significance (t -tests³) are computed on z-scored ratings for baseline
 166 vs. creative (CN-enhanced or CAA), where each participant is treated as an independent
 167 sample.

168 **4.3 Results on the Alternative Uses Test and the Task Task**

169 We report results in panels (b-f) of Figure 3. On the AUT, CN produces uniformly positive
 170 originality effects across all six models (avg. $d = +.36$), with four reaching significance. The
 171 effects are even stronger for surprise (avg. $d = +.43$), with five models reaching significance.
 172 On the Task Task, CN shows strong originality gains (avg. $d = +.40$), with the strongest

²Chu et al. (2024) studies additional dimensions, such as difficulty, how fun the task is to do, and how fun it is to watch, but here we restrict focus to creativity and originality, as these are most relevant to the goal of evaluating divergent thinking.

³We confirm responses follow a normal distribution before applying the t -tests.

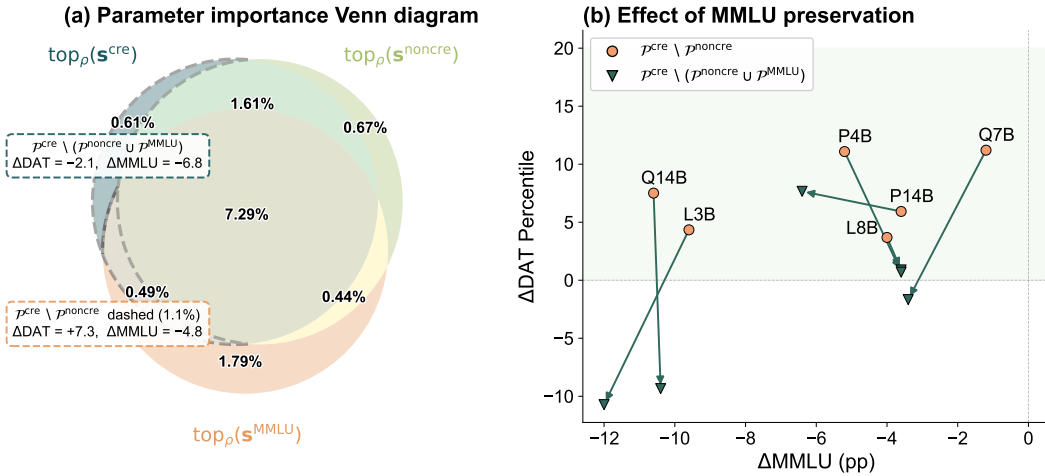


Figure 5: Creativity-factual reasoning trade-off under attempted MMLU preservation. (a) Venn diagram of top_{ρ} parameter importance sets for creative, non-creative, and MMLU prompts at $\rho = 0.1$. (b) Arrows show the shift from the original CN mask to the augmented mask, which entails catastrophic reversal of DAT improvements without recovering MMLU accuracy.

173 effects in Phi-14B ($d = +.61$) and Qwen-14B ($d = +.56$), and moderate creativity improve-
 174 ments (avg. $d = +.24$). Although CN degrades AUT utility, this is a predictable consequence
 175 of the novelty-utility tradeoff already present in baseline responses. Among baseline AUT
 176 stimuli, originality and utility are negatively correlated ($r = -.49$, $p < 10^{-7}$), as are surprise
 177 and utility ($r = -.64$, $p < 10^{-13}$). As predicted by fundamental novelty-utility tradeoffs
 178 established in the creativity literature (Varshney, 2019), increases in subjective ratings of
 179 novelty tend to be accompanied by decreases in perceived utility. Therefore, CN’s utility
 180 reduction is an expected byproduct of steering towards high-novelty responses, rather than
 181 an independent failure mode. To confirm this, we perform a hypervolume (HV) analysis
 182 (Cao et al., 2015) over all responses in the three-dimensional rating space (originality, sur-
 183 prise, utility) and find that the CN HV indicator is 15.2% larger than baseline ($p = .18$),
 184 suggesting CN obtains a moderate (but non-significant) gain in Pareto efficiency as well.

185 Additionally, while activation steering (93.9 avg) performs comparably to CN (94.1 avg) on
 186 the DAT (Figure 3), activation steering fails to generalize to the AUT and Task Task. These
 187 results are consistent with recent work that has found weight-space steering generalizes
 188 further out-of-distribution than activation steering on sycophancy and value alignment
 189 tasks (Fierro & Roger, 2025). Moreover, activation steering effectiveness has been shown
 190 to vary significantly by behavior type, with more complex behaviors like embodying
 191 persona archetypes and public figures proving more difficult to steer (Bas & Novak, 2025).
 192 Techniques such as context-dependent (Li et al., 2026; Lee et al., 2024) and learned activation
 193 steering (Rodriguez et al., 2025) have been proposed to remediate such issues, and may be
 194 explored in future work.

CN generalizes to open-ended creative tasks judged by human raters, demonstrably improving measures of originality and surprise, while activation steering does not exhibit reliable transfer.

196 5 Divergent Thinking Versus Factual Reasoning

197 Does improved divergent thinking come at the expense of general capabilities, such as
 198 factual reasoning? To test this, we evaluate CN-enhanced models⁴ on the MMLU benchmark

⁴For each of the six models tested, we take the top-performing (ρ , α , prompt) configuration.

199 (Hendrycks et al., 2021) to quantify the impact on question-answering abilities in STEM,
200 humanities, and the social sciences. Figure 5 reveals that CN-enhanced models reduce
201 MMLU performance by 1–11%.

202 A natural strategy to protect factual-reasoning weights during mask construction is to
203 incorporate MMLU questions in the negative contrast set in Algorithm 1, enforcing non-
204 interference between tasks. Using this technique, Christ et al. (2025) demonstrate that
205 mathematical reasoning can be improved without degrading MMLU performance. Does the
206 same hold true for a *meta-cognitive* ability like divergent thinking? To test this, we sample a
207 random set of 20 MMLU prompts $\mathcal{P}^{\text{MMLU}}$ and perform Algorithm 1 with \mathcal{P}^{cre} and $\mathcal{P}^{\text{non-cre}} \cup$
208 $\mathcal{P}^{\text{MMLU}}$. The masks then select for $\text{top}_\rho(\mathbf{s}^{\text{cre}}) \setminus \text{top}_\rho(\mathbf{s}^{\text{non-cre}} \cup \mathbf{s}^{\text{MMLU}})$, identifying weights
209 above the $1 - \rho$ percentile of importance for creative prompts that are not simultaneously
210 above the $1 - \rho$ percentile for MMLU or non-creative prompts. We recompute masks for
211 all six models using each model’s best hyperparameters, evaluate DAT score and 5-shot
212 MMLU accuracy on 500 held-out questions, and report the results in Figure 5.

213 We find that **attempting to preserve MMLU performance catastrophically reverses DAT**
214 **gains** and fails to improve MMLU scores, offering evidence that divergent thinking and
215 factual reasoning rely on functionally entangled weights. Existing work has established
216 that individual weights can be entangled in multiple distinct functions (a phenomenon
217 known as *polysematicity*), and that this supports the ability for neural models to represent
218 more features than they have neurons (*superposition*) (Elhage et al., 2022; Sharkey et al.,
219 2025). Kumar et al. (2025) attribute this behavior to *fractured entangled representations*, which
220 can inhibit creative acts that intentionally “break regularities” by failing to preserve other
221 regularities in a controlled manner. Creativity research more broadly suggests a need for
222 separation between generation (\approx DT) and selection (\approx convergent thinking (CT)) (Varshney,
223 2019), with neuroscience studies pointing to distinct brain-activation patterns involved in
224 associative (\approx DT) and controlled (\approx CT) processing modes (Zhang et al., 2020; Volle, 2018).
225 The success of multi-agent creative systems (Lin et al., 2025) and multi-stage prompting
226 techniques that decouple creative exploration from constraint satisfaction (Nguyen & Singla,
227 2025) can potentially be interpreted within the context of these results: if DT and CT rely on
228 functionally entangled weights, simultaneously eliciting strong DT and CT abilities may
229 be challenging. Therefore, separating these steps—whether across agents or prompts—can
230 provide stronger overall performance.

We find evidence that divergent thinking and factual reasoning are functionally entangled in
model weights.

231

232 6 Conclusion

233 In this work, we presented CreativityNeuro (CN), a data-free method for steering language
234 model weights to improve divergent thinking. CN-enhanced models improve measures of
235 divergent thinking on the DAT, AUT, and the Task Task, outperforming baselines such as
236 prompting, activation steering, and varying decoding strategies.

237 **Limitations and Future Work.** Several limitations of our study are worth acknowledging:
238 our evaluation is restricted to finite set of divergent thinking (DT) benchmarks and metrics
239 (DAT, AUT, Task Task), which capture only certain aspects of creativity. As Runco (2008)
240 notes, DT is not synonymous with creativity and should best be thought of as a measure of
241 creative potential. Moreover, our comparison to activation steering focuses on a standard
242 CAA-style method and does not exhaust the space of possible activation-space interventions,
243 such as context-dependent or learned approaches. Lastly, we found evidence suggesting DT
244 and factual reasoning are functionally entangled model weights—understanding whether
245 this entanglement reflects a fundamental architectural constraint or arises as an artifact of
246 the Wanda-style importance technique remains an important open question. Developing
247 techniques that can enhance DT without degrading factual reasoning—for instance, through
248 training-time objectives that encourage representational disentanglement—would be a
249 valuable direction for future work.

250 References

- 251 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan,
252 Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim,
253 Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary,
254 Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng,
255 Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao,
256 Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar,
257 Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter,
258 Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann,
259 Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat
260 Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan
261 Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mah-
262 moudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra,
263 Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas
264 Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby
265 Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael
266 Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia
267 Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guan-
268 hua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen,
269 Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Wei-
270 jian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan
271 Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang,
272 Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable
273 language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- 274 Tetiana Bas and Krystian Novak. What can we actually steer? a multi-behavior study of
275 activation control. *arXiv preprint arXiv:2511.18284*, 2025.
- 276 Antoine Bellemare-Pepin, François Lespinasse, Philipp Thölke, Yann Harel, Kory Mathew-
277 son, Jay A Olson, Yoshua Bengio, and Karim Jerbi. Divergent Creativity in Humans and
278 LLMs. Technical report, 2024.
- 279 Margaret A Boden. *The Creative Mind: Myths and Mechanisms*. Routledge, 2004.
- 280 Yongtao Cao, Byran J. Smucker, and Timothy J. Robinson. On using the hypervolume
281 indicator to compare pareto fronts: Applications to multi-criteria optimal experimental
282 design. *Journal of Statistical Planning and Inference*, 160:60–74, 2015. ISSN 0378-3758.
283 doi: <https://doi.org/10.1016/j.jspi.2014.12.004>. URL [https://www.sciencedirect.com/
284 science/article/pii/S0378375814002006](https://www.sciencedirect.com/science/article/pii/S0378375814002006).
- 285 Bryan R. Christ, Zack Gottesman, Jonathan Kropko, and Thomas Hartvigsen. Math Neu-
286 rosurgery: Isolating Language Models’ Math Reasoning Abilities Using Only Forward
287 Passes. 6 2025. URL <http://arxiv.org/abs/2410.16930>.
- 288 Junyi Chu, Jennifer Hu, and Tomer D Ullman. The task task: Creative problem generation in
289 humans and language models. In *Proceedings of the Annual Meeting of the Cognitive Science
290 Society*, volume 46, 2024.
- 291 Arne Dietrich. Types of creativity. *Psychonomic Bulletin and Review*, 26(1):1–12, 2 2019. ISSN
292 15315320. doi: 10.3758/s13423-018-1517-7.
- 293 Dietrich, Arne. The Cognitive Neuroscience of Creativity. *Psychonomic Bulletin & Review*, 11
294 (6):1011–1026, 2004.
- 295 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan,
296 Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger
297 Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Chris
298 Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. URL [https:
299 //transformer-circuits.pub/2022/toy_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html).
- 300 Gilles Fauconnier and Mark Turner. *The Way We Think: Conceptual Blending and the Mind’s
301 Hidden Complexities*. Basic Books, 2008.

-
- 302 Constanza Fierro and Fabien Roger. Steering language models with weight arithmetic. *arXiv*
303 *preprint arXiv:2511.05408*, 2025.
- 304 Francis Galton. *Hereditary genius: An inquiry into its laws and consequences*. D. Appleton &
305 Company, 1870.
- 306 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian,
307 Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy
308 Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie
309 Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rod-
310 riguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bob-
311 bie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell,
312 Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer,
313 Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny
314 Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego
315 Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan,
316 Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve,
317 Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon,
318 Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron,
319 Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack
320 Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,
321 Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,
322 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
323 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden
324 Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin
325 Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal
326 Lakhota, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz
327 Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke
328 de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin
329 Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie
330 Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
331 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang,
332 Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar
333 Vasic, Peter Weng, Prajwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura,
334 Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer,
335 Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Gird-
336 har, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan
337 Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean
338 Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rap-
339 parthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
340 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Syd-
341 ney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas
342 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami,
343 Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish
344 Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney
345 Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia,
346 Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei,
347 Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert,
348 Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha
349 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,
350 Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda
351 Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew
352 Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani,
353 Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley
354 Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer,
355 Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing
356 Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic,
357 Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Chang-
358 han Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris
359 Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer,

360 Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
 361 Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa
 362 Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik
 363 Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng
 364 Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide,
 365 Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern,
 366 Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan,
 367 Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison
 368 Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim
 369 Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake
 370 Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya,
 371 Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul,
 372 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan
 373 McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena,
 374 Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly
 375 Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin
 376 Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo,
 377 Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish
 378 Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Gro-
 379 shev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal
 380 Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,
 381 Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari,
 382 Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa,
 383 Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev,
 384 Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem
 385 Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bon-
 386 trager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj,
 387 Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu
 388 Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin
 389 Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh
 390 Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru
 391 Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun
 392 Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil,
 393 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji
 394 Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satter-
 395 field, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk,
 396 Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara
 397 Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy
 398 Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
 399 Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihalescu,
 400 Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable,
 401 Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman,
 402 Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin
 403 Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary
 404 DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The
 405 llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

406 J P Guilford. Psychological Bulletin THE STRUCTURE OF INTELLECT. Technical Report 4,
 407 1956.

408 Jacques Hadamard. *An Essay on the Psychology of Invention in the Mathematical Field*. Courier
 409 Corporation, 1954.

410 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang,
 411 Dawn Song, and Jacob Steinhardt. Measuring Mathematical Problem Solving With the
 412 MATH Dataset. 11 2021. URL <http://arxiv.org/abs/2103.03874>.

413 Liwei Jiang, Yuanjun Chai, Margaret Li, Mickel Liu, Raymond Fok, Nouha Dziri, Yulia
 414 Tsvetkov, Maarten Sap, Alon Albalak, and Yejin Choi. Artificial hivemind: The open-
 415 ended homogeneity of language models (and beyond). *arXiv preprint arXiv:2510.22954*,
 416 2025.

-
- 417 Arthur Koestler. *The Act of Creation*. Macmillan, 1964.
- 418 Akarsh Kumar, Jeff Clune, Joel Lehman, and Kenneth O Stanley. Questioning represen-
419 tational optimism in deep learning: The fractured entangled representation hypothesis.
420 *arXiv preprint arXiv:2505.11581*, 2025.
- 421 Bruce W Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehling, Pierre Dognin,
422 Manish Nagireddy, and Amit Dhurandhar. Programming refusal with conditional activa-
423 tion steering. *arXiv preprint arXiv:2409.05907*, 2024.
- 424 Jiaqian Li, Yanshu Li, and Kuan-Hao Huang. Steering vector fields for context-aware
425 inference-time control in large language models. *arXiv preprint arXiv:2602.01654*, 2026.
- 426 Yi-Cheng Lin, Kang-Chieh Chen, Zhe-Yan Li, Tzu-Heng Wu, Tzu-Hsuan Wu, Kuan-Yu
427 Chen, Hung-yi Lee, and Yun-Nung Chen. Creativity in llm-based multi-agent systems:
428 A survey. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language*
429 *Processing*, pp. 27572–27595, 2025.
- 430 Jack Lindsey, Adly Templeton, Jonathan Marcus, Thomas Conerly, Joshua Batson, and
431 Christopher Olah. Sparse crosscoders for cross-layer features and model diffing. Technical
432 report, Anthropic, 2024. URL [https://transformer-circuits.pub/2024/crosscoders/
433 index.html](https://transformer-circuits.pub/2024/crosscoders/index.html).
- 434 Mary Lou Maher. Evaluating creativity in humans, computers, and collectively intelligent
435 systems. In *Proceedings of the 1st DESIRE Network Conference on Creativity and Innovation in*
436 *Design*, DESIRE '10, pp. 22–28, 2010.
- 437 Sarnoff Mednick. The associative basis of the creative process. *Psychological Review*, 69(3):
438 220, 1962.
- 439 Robert Morain and Dan Ventura. Is Prompt Engineering the Creativity Knob for Large
440 Language Models? In *Proceedings of the 16th International Conference on Computational*
441 *Creativity (ICCC'25)*, 2025.
- 442 Manh Hung Nguyen and Adish Singla. Divergent-convergent thinking in large language
443 models for creative problem generation. *arXiv preprint arXiv:2512.23601*, 2025.
- 444 Jay A. Olson, Johnny Nahas, Denis Chmoulevitch, Simon J. Cropper, and Margaret E. Webb.
445 Naming unrelated words predicts creativity. 4 2021. doi: 10.1073/pnas.2022340118/-/
446 DCSupplemental.y.
- 447 Matthew Lyle Olson, Neale Ratzlaff, Musashi Hinck, Shao-yen Tseng, and Vasudev Lal.
448 Steering Large Language Models to Evaluate and Amplify Creativity. 12 2024. URL
449 <http://arxiv.org/abs/2412.06060>.
- 450 Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexan-
451 der Matt Turner. Steering Llama 2 via contrastive activation addition. *arXiv preprint*
452 *arXiv:2312.06681*, 2024.
- 453 Max Peeperkorn, Tom Kouwenhoven, Dan Brown, and Anna Jordanous. Is temperature the
454 creativity parameter of large language models? *arXiv:2405.00492*, 2024.
- 455 Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global Vectors for
456 Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural*
457 *Language Processing (EMNLP)*, pp. 1532–1543. Association for Computational Linguistics,
458 2014. doi: 10.3115/v1/D14-1162. URL <https://nlp.stanford.edu/pubs/glove.pdf>.
- 459 Lambert AJ Quetelet. A treatise on man and the development of his faculties (a facsimile
460 reproduction of the english translation of 1842 with an introduction by solomon diamond).
461 1842.
- 462 Pau Rodriguez, Michal Klein, Eleonora Gualdoni, Valentino Maiorca, Arno Blaas, Luca
463 Zappella, Marco Cuturi, and Xavier Suau. Lineas: End-to-end learning of activation
464 steering with a distributional loss. *NeurIPS*, 2025.

-
- 465 A Rothenberg. *Flight from Wonder: An Investigation of Scientific Creativity*. Oxford University
466 Press, 2014.
- 467 Mark A. Runco. Commentary: Divergent Thinking Is Not Synonymous With Creativity.
468 *Psychology of Aesthetics, Creativity, and the Arts*, 2(2):93–96, 5 2008. ISSN 19313896. doi:
469 10.1037/1931-3896.2.2.93.
- 470 Aishik Sanyal, Samuel Schapiro, Sumuk Shashidhar, Royce Moon, Lav R Varshney, and
471 Dilek Hakkani-Tur. Spark: A system for scientifically creative idea generation. *arXiv*
472 *preprint arXiv:2504.20090*, 2025.
- 473 Samuel Schapiro, Jonah Black, and Lav R Varshney. Transformational Creativity in Science:
474 A Graphical Theory. *arXiv preprint arXiv:2504.18687*, 2025.
- 475 Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas
476 Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, Stella Biderman,
477 Adria Garriga-Alonso, Arthur Conmy, Neel Nanda, Jessica Rumbelow, Martin Wattenberg,
478 Nandi Schoots, Joseph Miller, Eric J. Michaud, Stephen Casper, Max Tegmark, William
479 Saunders, David Bau, Eric Todd, Atticus Geiger, Mor Geva, Jesse Hoogland, Daniel
480 Murfet, and Tom McGrath. Open problems in mechanistic interpretability, 2025. URL
481 <https://arxiv.org/abs/2501.16496>.
- 482 Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can LLMs Generate Novel Research
483 Ideas? A Large-Scale Human Study with 100+ NLP Researchers. 9 2024. URL <http://arxiv.org/abs/2409.04109>.
- 484
- 485 Chenglei Si, Tatsunori Hashimoto, and Diyi Yang. The Ideation-Execution Gap: Execution
486 Outcomes of LLM-Generated versus Human Research Ideas. 6 2025. URL <http://arxiv.org/abs/2506.20803>.
- 487
- 488 Dean Keith Simonton. *Creativity in Science: Chance, Logic, Genius, and Zeitgeist*. Cambridge
489 University Press, 2004.
- 490 Claire Stevenson, Iris Smal, Matthijs Baas, Raoul Grasman, and Han Van Der Maas. Putting
491 GPT-3’s Creativity to the (Alternative Uses) Test. In *International Conference on Computa-*
492 *tional Creativity*, 2022. URL <http://osf.io/vmk3c/>.
- 493 Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning
494 approach for large language models. *arXiv preprint arXiv:2306.11695*, 2023.
- 495 L. R. Varshney. Mathematical limit theorems for computational creativity. *IBM Journal of*
496 *Research and Development*, 63(1), 1 2019. ISSN 21518556. doi: 10.1147/JRD.2019.2893907.
- 497 Emmanuelle Volle. Associative and controlled cognition in divergent thinking: Theoretical,
498 experimental, neuroimaging evidence, and new directions. *The Cambridge Handbook of the*
499 *Neuroscience of Creativity*, pp. 333, 2018.
- 500 Haining Wang, Peng Bao, Luning Qiu, Dawei Wu, Nanyun Yu, Haoran Liu, and Samuel
501 Johnson. A large-scale comparison of divergent creativity in humans and large language
502 models. *Nature Human Behaviour*, 2025. doi: 10.1038/s41562-025-02331-1.
- 503 Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. SciMON: Scientific Inspiration
504 Machines Optimized for Novelty. 6 2024. doi: 10.18653/v1/2024.acl-long.18. URL <http://arxiv.org/abs/2305.14259><http://dx.doi.org/10.18653/v1/2024.acl-long.18>.
- 505
- 506 Xiaolong Wei, Bo Lu, Xingyu Zhang, Zhejun Zhao, Dongdong Shen, Long Xia, and Dawei
507 Yin. Igniting creative writing in small language models: Llm-as-a-judge versus multi-
508 agent refined rewards. In *Proceedings of the 2025 Conference on Empirical Methods in Natural*
509 *Language Processing*, pp. 17171–17197, 2025.

-
- 510 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
511 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei
512 Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin
513 Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin,
514 Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su,
515 Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5
516 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- 517 Weitao Zhang, Zsuzsika Sjoerds, and Bernhard Hommel. Metacontrol of human creativity:
518 The neurocognitive mechanisms of convergent and divergent thinking. *NeuroImage*, 210:
519 116572, 2020.

520 A Full Sets of Contrastive Prompts

521 Each of the six prompt sets contains 10 creative and 10 non-creative exemplars. Below we
522 show 3 representative examples from each set.

Table 2: **Full contrastive prompt sets (3 examples each).** Each prompt set contains 10 creative (\mathcal{P}^{cre}) and 10 non-creative ($\mathcal{P}^{\text{non-cre}}$) exemplars used for parameter importance scoring in Algorithm 1.

Set	Creative (\mathcal{P}^{cre})	Non-Creative ($\mathcal{P}^{\text{non-cre}}$)
DAT	<p>List 10 common English nouns that are as unrelated in meaning as possible. Avoid any shared topic or category. Output only the nouns, separated by commas.</p> <p>Give me 10 one-word English nouns that are extremely far apart in semantic meaning. They should not fit into a single theme.</p> <p>Produce 10 everyday English nouns that share no obvious connection with one another. Each noun should come from a very different domain.</p>	<p>List 10 common English nouns that are as closely related in meaning as possible and clearly fit into a single narrow topic. Output only the nouns, separated by commas.</p> <p>Give me 10 one-word English nouns that are very strongly associated with each other and belong to the same specific domain.</p> <p>Produce 10 English nouns that are tightly connected around a single clear theme (for example, all parts of a computer or all items in a kitchen).</p>
STORY	<p>Write an unusual opening sentence for a short story that subverts reader expectations.</p> <p>Create a story beginning that combines two unrelated concepts in a surprising way.</p> <p>Write the first line of a story that makes the reader question reality.</p>	<p>Write a standard opening sentence for a fairy tale.</p> <p>Create a typical story beginning that establishes setting and character clearly.</p> <p>Write the first line of a conventional mystery novel.</p>
IDEATION	<p>List 5 unusual uses for a brick that nobody has thought of before.</p> <p>Generate unconventional solutions to reduce traffic in cities.</p> <p>What are some surprising ways a library could be repurposed?</p>	<p>List 5 common uses for a brick in construction.</p> <p>Generate standard solutions to reduce traffic in cities.</p> <p>What are the traditional functions of a library?</p>
PROBLEM	<p>How might you solve this problem in a way that seems counterintuitive at first?</p> <p>What would an alien civilization’s approach to this problem look like?</p> <p>If you had to solve this with resources from a different era, what would you do?</p>	<p>What is the most efficient way to solve this problem?</p> <p>List the standard steps for addressing this type of issue.</p> <p>What best practices should be followed when solving this?</p>
OPEN	<p>Generate a joke about electric vehicles.</p> <p>Create the first verse of a wedding vow.</p> <p>Write a song about a guy named Jacob working at a call center making jokes.</p>	<p>Explain nuclear fission like I am five years old.</p> <p>What is Bukhara? Provide a paragraph-long explanation in layman’s language.</p> <p>In a few sentences explain what threats do scams pose to individuals?</p>
MINIMAL	<p>Invent something.</p> <p>Surprise me.</p> <p>Be weird.</p>	<p>State a fact.</p> <p>Be accurate.</p> <p>Be precise.</p>

523 B Baseline Sweep Settings

524 In Section 3, we compare CreativityNeuro against a set of baseline techniques.

- 525 1. **Prompting:** We present creative exemplars from each of six prompt sets (Table 2) as
526 in-context guides
- 527 2. **Decoding Parameter Sweeps**
- 528 (a) **top- p nucleus sampling**, $p \in \{0.8, 0.85, 0.9, 0.95, 1.0\}$
- 529 (b) **top- k sampling**, $k \in \{10, 25, 50, 100, \text{disabled}\}$
- 530 (c) **repetition penalty**, $\theta \in \{1.0, 1.1, 1.2, 1.5, 2.0, 3.0\}$
- 531 3. **Activation Steering:** We sweep the injection layer suffix from single-layer to the final
532 50% of layers and find that injecting into the final 30% works best. We sweep $\alpha \in$

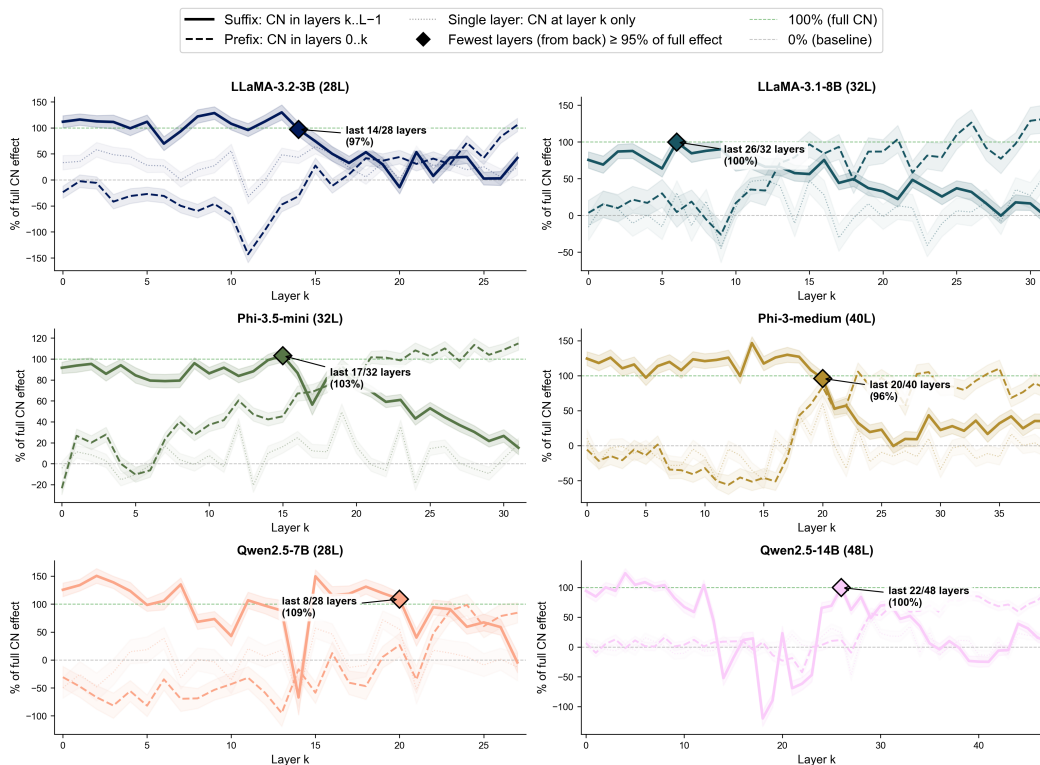


Figure 6: **Layerwise ablation: suffix vs. prefix vs. single-layer.** Each panel shows the % of full CN DAT effect recovered as a function of the number of layers k with CN weights applied. Solid lines: suffix (last k layers). Dashed lines: prefix (first k layers). Dotted lines: single-layer (one layer at a time, plotted by layer index). Diamond markers indicate the fewest layers from the back achieving $\geq 95\%$ of the full effect.

533 $\{0.1, 0.2, 0.3, 0.4, 0.5, 1.0, 2.0, 4.0\}$, selecting the best α where all three temperatures
534 yield ≥ 120 valid samples.

535 C Layerwise Ablation Studies

536 C.1 Suffix vs. Prefix

537 To localize the CN effect across the network, we compare three layerwise interventions
538 (Figure 6): (i) *suffix*, applying CN weights to only the last k layers; (ii) *prefix*, applying CN
539 weights to only the first k layers; and (iii) *single-layer*, applying CN weights to one layer at a
540 time. For each condition, we generate $N = 120$ valid DAT samples at $T = 1.0$ and report the
541 percentage of the full CN effect recovered in terms of DAT scores.

542 **Suffix.** Applying CN weights to a suffix of layers (layers k through $L-1$) recovers the full
543 effect with roughly half the network. On average, 51% of layers (from the back) are needed
544 to reach 100% recovery, ranging from 29% for Qwen-7B (last 8/28) to 81% for Llama-8B
545 (last 26/32). The remaining models fall in between: Llama-3B (last 14/28, 50%), Phi-14B
546 (last 20/40, 50%), Qwen-14B (last 22/48, 46%), and Phi-4B (last 17/32, 53%). However, this
547 analysis is only conducted on DAT scores, and it is unclear whether the last 51% of layers
548 are sufficient to recover 100% of the scores on the AUT and Task Task.

549 **Prefix.** Applying CN weights from the front (layers 0 through k) is far less efficient. On
550 average, 78% of all layers must be included before the prefix condition reaches 95% recovery.
551 For Qwen-14B (48 layers), the prefix condition never reaches 95% even when all layers

552 are included (94% at $k = 48$). This asymmetry provides evidence that the CN effect is
553 concentrated in later layers of the residual stream.

554 **Single-layer.** No single layer is sufficient to recover the full CN effect. The best individual
555 layers recover 50–72% of the effect (e.g., Q7B layer 19: 72%, L3B layer 15: 64%, P14B
556 layer 20: 60%), with top contributors generally appearing in middle-to-late layers. The gap
557 between the best single layer and the suffix provides evidence that the CN effect requires
558 cooperation across multiple late layers, consistent with findings that representations of
559 human-interpretable concepts or behaviors may span multiple layers (Lindsey et al., 2024;
560 Sharkey et al., 2025).

561 D Hyperparameter Sweeps

562 CreativityNeuro introduces two hyperparameters: the *importance threshold* $\rho \in (0, 1]$, which
563 controls the fraction of parameters selected by the importance mask, and the *scaling factor* $\alpha >$
564 0 , which controls the magnitude of the weight perturbation applied to masked parameters.
565 To identify effective (ρ, α) configurations for each model, we conduct systematic grid sweeps
566 evaluated on the DAT. For each model, we generate CN weight masks using six different
567 prompt sets: DAT, STORY, IDEATION, PROBLEM, OPENENDED, and MINIMAL. Each prompt
568 set produces a separate importance mask. We then evaluate every combination of keep
569 ratio $\rho \in \{0.01, 0.05, 0.1, 0.2\}$ and scaling factor $\alpha \in \{0.1, 0.5, 1.0, 2.0\}$ at three sampling
570 temperatures $T \in \{0.9, 1.0, 1.2\}$, with top- $k = 0$ and top- $p = 1.0$ (i.e., untruncated sampling).
571 For LLaMA-3.1-8B and Phi-3.5-mini, after initial experiments revealed that $\alpha > 0.5$ were too
572 high, we additionally tested a finer alpha grid $\alpha \in \{0.01, 0.05, 0.075, 0.1, 0.5, 1.0\}$ to probe the
573 conservative regime identified by Christ et al. (2025). Each configuration generates $n = 120$
574 valid DAT samples (10-word lists with all words present in the GloVe vocabulary).

575 For each (ρ, α, T) triple, we compute the mean DAT scores for the CN-enhanced ($\overline{\text{DAT}}_{\text{CN}}$)
576 and baseline ($\overline{\text{DAT}}_{\text{base}}$) models, then convert each to a human percentile using the distribu-
577 tion from Wang et al. (2025). Figures 7 to 12 show $\Delta \text{Percentile} = P(\overline{\text{DAT}}_{\text{CN}}) - P(\overline{\text{DAT}}_{\text{base}})$
578 (averaged across temperatures) as a function of (α, ρ) for each of the six prompt sets, across
579 all six models. Positive values indicate that CN moves the model’s DAT performance
580 upward in the human score distribution.

581 E Evaluation Prompts

582 In this section, we share the full prompts used to generate and/or post-process the stimuli
583 for the DAT, AUT, and TT.

584 E.1 Divergent Association Task (DAT)

585 The DAT prompt instructs models to generate 10 maximally dissimilar nouns:

DAT Prompt

You will be asked to name exactly 10 English nouns.
Output exactly 10 words, separated by commas, and nothing else.
Return ONLY: word1, word2, ..., word10
DO NOT write explanations. DO NOT think step by step.

Be as original and unusual as possible. Avoid common or closely related
words.
Now name 10 English nouns that are as different from each other as possible.

586

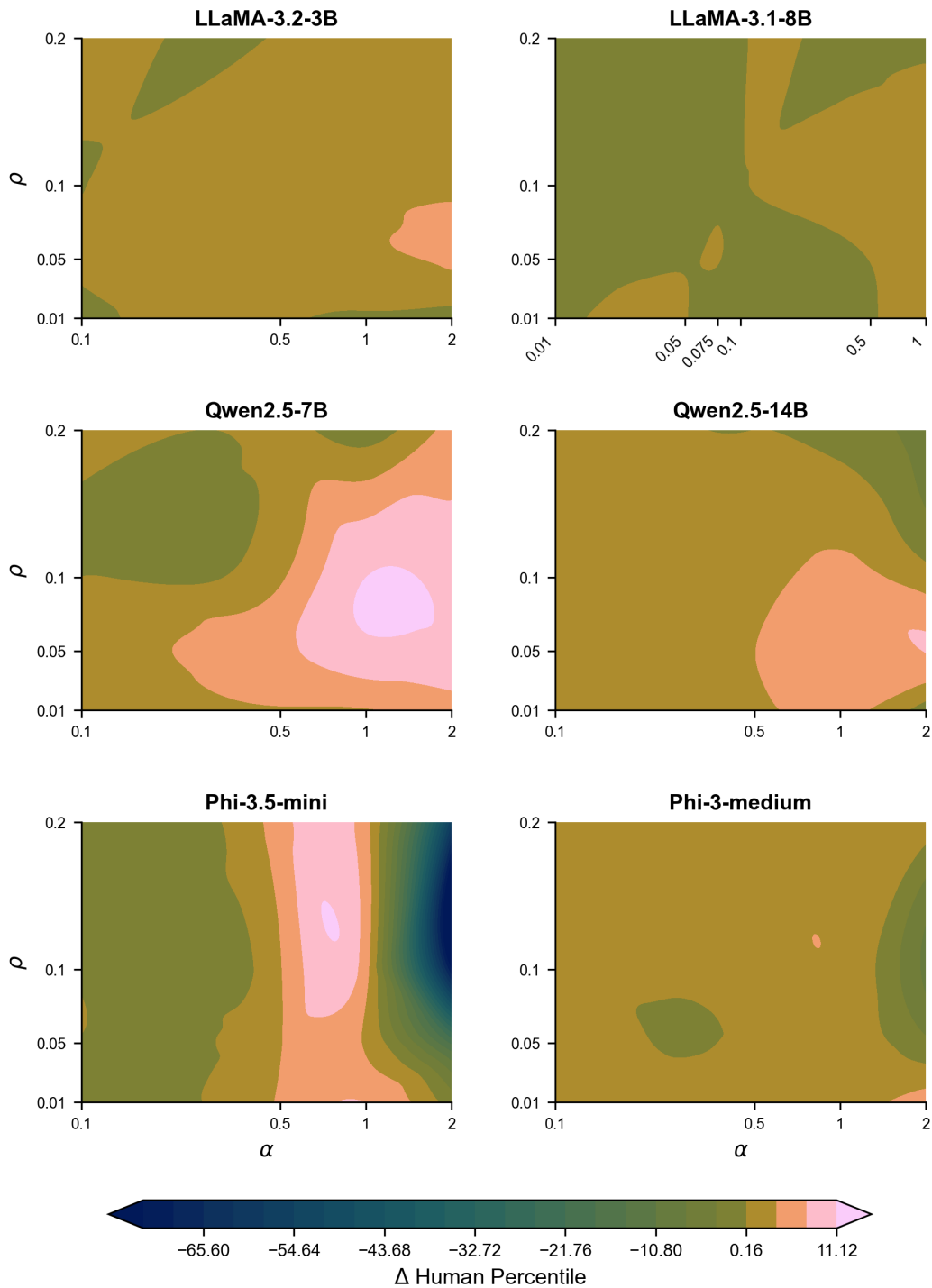


Figure 7: **Sensitivity to (α, ρ) — DAT prompt set.** Δ Percentile averaged across three temperatures ($T \in \{0.9, 1.0, 1.2\}$) for each model. Color scale is shared across panels and centered at zero.

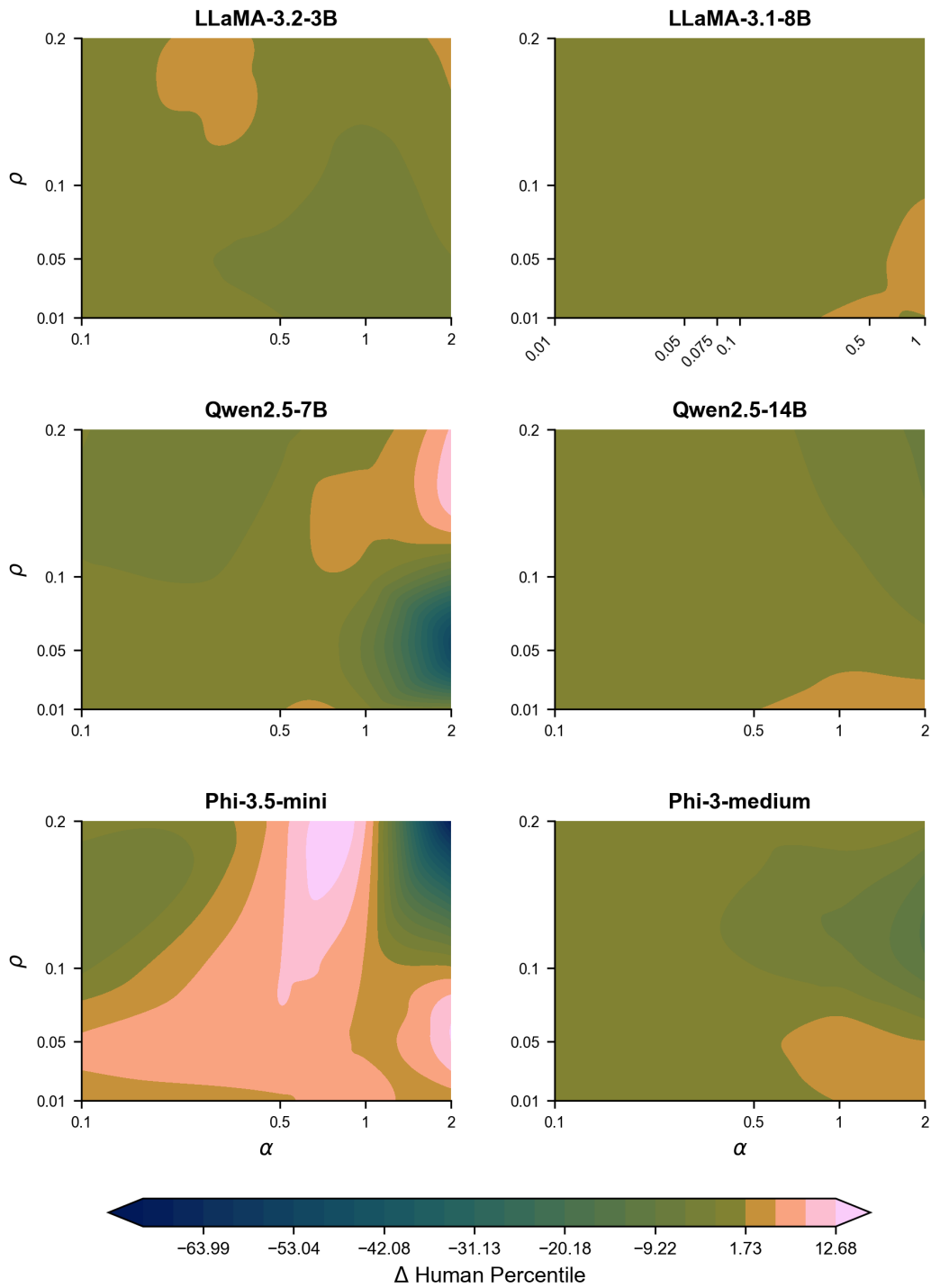


Figure 8: **Sensitivity to (α, ρ) — Story prompt set.** Mean Δ DAT averaged across three temperatures ($T \in \{0.9, 1.0, 1.2\}$) for each model. Color scale is shared across panels and centered at zero.

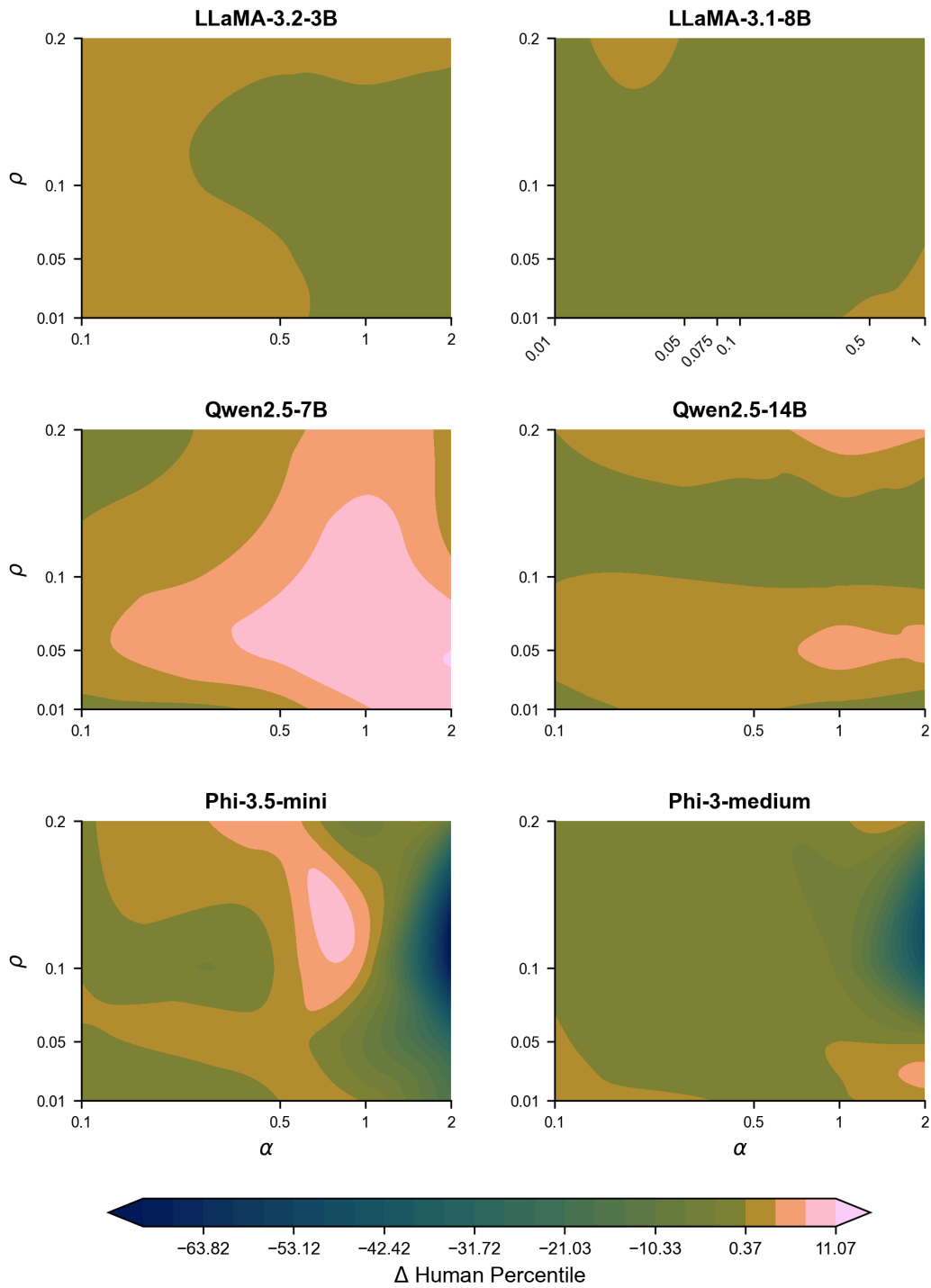


Figure 9: **Sensitivity to (α, ρ) — Ideation prompt set.** Mean Δ DAT averaged across three temperatures ($T \in \{0.9, 1.0, 1.2\}$) for each model. Color scale is shared across panels and centered at zero.

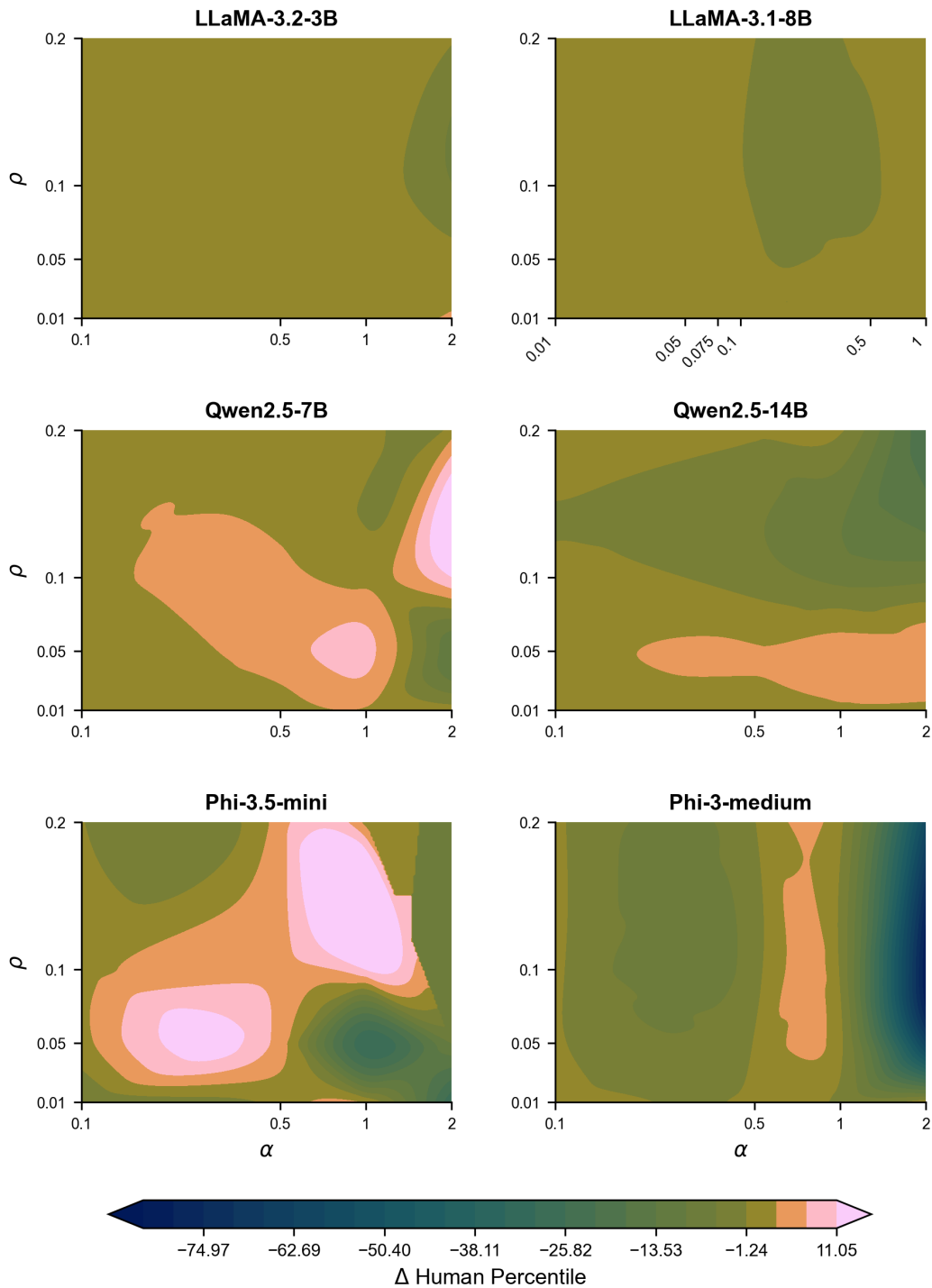


Figure 10: **Sensitivity to (α, ρ) — Problem prompt set.** Mean Δ DAT averaged across three temperatures ($T \in \{0.9, 1.0, 1.2\}$) for each model. Color scale is shared across panels and centered at zero.

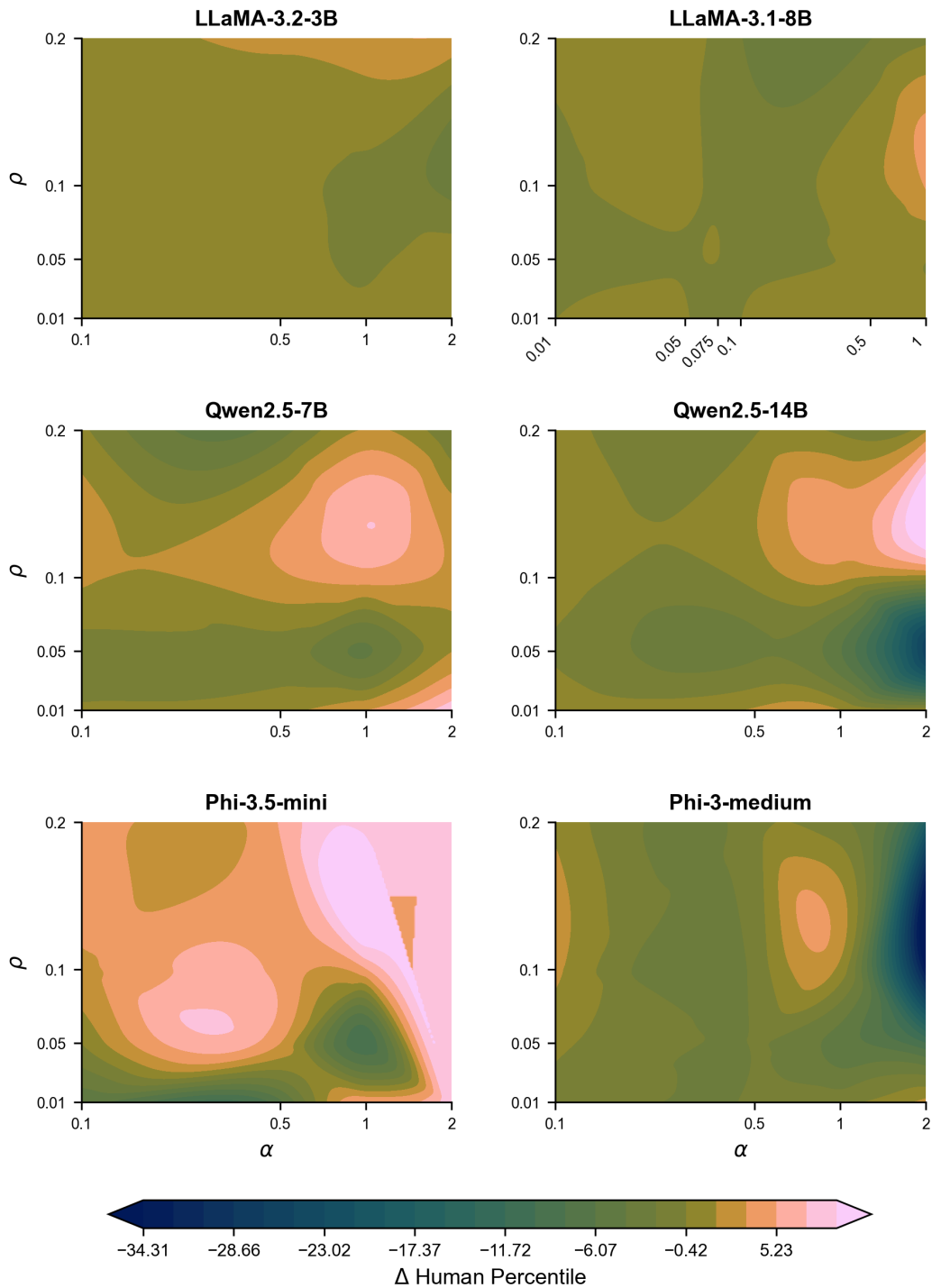


Figure 11: **Sensitivity to (α, ρ) — Open-ended prompt set.** Mean Δ DAT averaged across three temperatures ($T \in \{0.9, 1.0, 1.2\}$) for each model. Color scale is shared across panels and centered at zero.

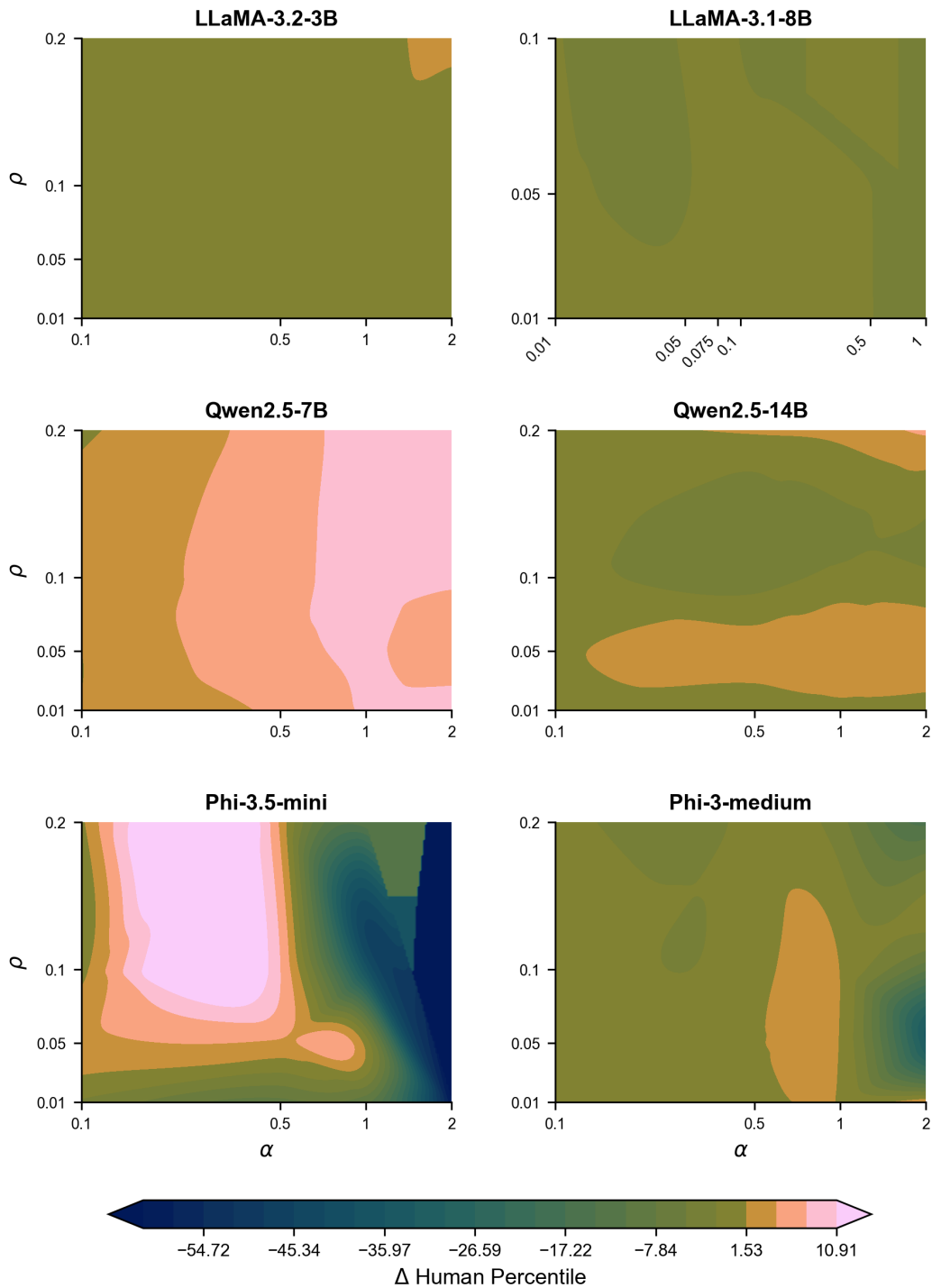


Figure 12: **Sensitivity to (α, ρ) — Minimal prompt set.** Mean Δ DAT averaged across three temperatures ($T \in \{0.9, 1.0, 1.2\}$) for each model. Color scale is shared across panels and centered at zero.

587 E.2 Alternative Uses Test (AUT)

588 The AUT uses a system prompt and a user prompt, with {object} replaced by one of 3
589 standard objects: *brick, paperclipfork*.

AUT System Prompt

You are participating in a creativity test. Your task is to generate creative, unusual, and original uses for common objects. Be imaginative and think outside the box.

590

AUT User Prompt

List 5 creative and unusual alternative uses for a {object}.
Be specific and creative. List each use on a new line, numbered 1 through 5.
Only list the uses, no explanations.

591

592 E.3 Task Task (TT)

593 The Task Task uses a system prompt and a user prompt with three in-context examples from
594 [Chu et al. \(2024\)](#).

TT System Prompt

You are a creative game show designer. Your task is to invent fun, original, and entertaining challenges that would be exciting for contestants to attempt and audiences to watch.

595

TT User Prompt

You've been recruited to help design challenge tasks for a new game show! Your job is to come up with a new creative, silly, and fun task for humans to solve.

Here are a few example tasks:

1. Your goal is to: Throw a teabag into a mug from the farthest distance. You can use: Anything you can reasonably expect to find in a house, garage, and garden shed.

2. Someone has squeezed all of the toothpaste out of the toothpaste tube. Your goal is to: Get as much of the original toothpaste back into the empty toothpaste tube as possible. You can use: Anything you can reasonably expect to find in a bathroom.

3. Your goal is to: Transfer water between two fishbowls using only the supplied items. You cannot move the fishbowls. You can use: a chocolate bar, a rubber glove, a baguette, a snorkel, a cardboard tube, and a plate of pasta.

Now it's your turn! Create your own creative, silly, and fun task for future participants to solve. Specify the goal, scoring criteria, and any materials or constraints. Describe it in exactly 3--4 sentences as a short paragraph --- do not use lists or bullet points. Do not comment on how entertaining, creative, or fun the task would be.

596

597 E.4 Task Task Post-Processing

598 Raw Task Task generations are post-processed using GPT-4o to normalize formatting before
599 human evaluation.

TT Post-Processing System Prompt

You are an editor cleaning game show challenge descriptions for a human evaluation study.

Make these MINIMAL edits:

1. REMOVE any task title at the start (e.g., ‘In “Baking Bonanza,”’). After removing, capitalize the first word of the remaining text.
2. Convert ALL third-person references to SECOND PERSON (e.g., ‘‘contestants must’’ → ‘‘you must’’, ‘‘the player’’ → ‘‘you’’).
3. REMOVE any metacommentary sentences (e.g., ‘‘This creative juggling act combines laughs with a dash of physics comedy.’’).

CRITICAL RULES:

- Do NOT summarize, condense, or shorten the task description
- Do NOT change the creative content or game mechanics
- Do NOT add any text --- only edit or remove
- Output ONLY the cleaned description with no preamble or explanation

600

TT Post-Processing User Prompt

Clean this game show challenge description: {text}

601

602 F Disclosure of Large Language Model Usage

603 In this paper, large language models (LLMs) were used to assist in the code implementation,
604 plotting and figure generation, reporting (but not analysis) of experiment results, and for
605 comprehensive surveys of related work.