INVESTIGATING THE EFFECTS OF EMOTIONAL STIM ULI TYPE AND INTENSITY ON LARGE LANGUAGE MODEL (LLM) BEHAVIOR

Anonymous authors

Paper under double-blind review

ABSTRACT

Emotional prompting—the use of specific emotional diction in prompt engineering—has shown increasing promise in improving large language model (LLM) performance, truthfulness, and responsibility, however these studies have been limited to single type of positive emotional stimuli and have not considered varying degrees of emotion intensity in their analyses. In this paper, we explore the effects of "positive" (joy and encouragement) and "negative" (anger and insecurity) emotional prompting on accuracy, sycophancy, and toxicity. To analyze their effects, we developed a suite of LLM- and human-generated add-on prompts of varying intensities across our four emotions using GPT-40 mini. We also created a gold dataset of only those prompts that are perceived similarly by humans and LLMs for emotion labels and intensity levels. Our empirical evaluation on LLM behavior on accuracy, sycophancy and toxicity datasets has shown that positive emotional stimuli can lead to a more accurate and less toxic results but also may lead to greater sycophantic behavior.

025 026 027

028

006

007

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

In natural language processing (NLP), large language models (LLMs) display remarkable performance in both domain-specific and diverse tasks (Chang et al., 2023). The ability of these models to generate substantial amounts of text is highly effective for dialogue systems, question answering, and other NLP tasks (Chang et al., 2023). Taking advantage of this, LLMs have been widely trained and applied for a variety of real-world applications, ranging from legal compliance to education (Hassani, 2024; Gan et al., 2023).

LLMs can be used with basic prompting; however, the performance of these models can be improved with the use of prompt engineering (Minaee et al., 2024). Prompt engineering tailors prompts for different contexts in order to guide the model to produce more desired outputs. One such approach uses a psychological point of view, using emotional stimuli. LLMs have been shown to understand and be able to be influenced by these emotions (Li et al., 2023; Wang et al., 2023). Using certain emotional stimuli in prompts has been shown to improve LLM performance (Li et al., 2023; Wang et al., 2024). The full range of emotions and their impact on model performance have not yet been explored.

Despite potential performance gains, the inherent tendency for LLMs to exhibit sycophantic behavior, in which models agree excessively with the user, still exists in LLMs and continues to challenge
researchers (Malmqvist, 2024). Addressing sycophancy is crucial to ensure the accuracy and reliability of the information generated by LLMs, especially for practical applications (Malmqvist, 2024). Although the causes of sycophancy are complex and can be attributed to a variety of factors,
the effects of various emotional stimuli is not explored (Malmqvist, 2024; Wei et al., 2024).

Existing research surrounding emotional stimuli records increased performance with small sets of
specific human-designed prompts (Li et al., 2023; Yin et al., 2024; Wang et al., 2024). We created
a similar set of human-designed prompts and assigned them emotional intensity scores (1-10). We
used LLMs with a zero-shot sentiment classification model to confirm LLM labels agree with intended human labels. We also used few-shot prompting (Li, 2023) to create a larger dataset of 415 model-written prompts based on the human design prompts dataset.

054 2 RELATED WORK

Studies have increasingly explored the impact of emotional stimuli on LLMs, demonstrating that performance can be enhanced with emotional prompts (Li et al., 2023; Wang et al., 2024). For instance, moderate politeness in prompts has been shown to improve LLM performance on language understanding and summarization tasks (Yin et al., 2024). Additionally, the use of positive words informed by psychological theories has proven effective in notably boosting LLM performance across task performance, truthfulness, and informativeness (Li et al., 2024; Wang et al., 2024).

062 While a number of papers demonstrate positive effects of emotional prompts on truthfulness and 063 informativeness, little attention has been given to their influence on sycophancy or overly agreeable 064 responses, despite its ability to impact user experience. We take this into consideration in our re-065 search, and measure our findings on the (SycophancyEval). Furthermore, we expand upon previous 066 studies by incorporating a broader emotional spectrum in our prompts, including both positive and 067 negative categorizations. This allows a comprehensive analysis of the effect of diverse emotional stimuli on behavioral tendencies, particularly sycophantic behavior, and task accuracies, including 068 toxicity, both of which significantly impact user interactions with LLMs. 069

3 Methods

071 072



Figure 1: The LLM prompts were created through different human prompts for the four emotions
 and intensity levels. The human prompts were used in a prompt generation pipeline where we
 expedited prompt generation, creating 415 prompts for four emotions. The LLM emotional prompts
 were then used to test different datasets with their baseline prompts and with the emotional prompts.

We created a set of human-made emotional prompts with anger, joy, insecurity, and encouragement. We ranked these prompts on a scale of 1 to 10, where 1 represented the lowest degree of emotion, and 10 represented the highest. We developed an emotion detection pipeline, in which we used zero-shot prompting (Li, 2023) and GPT-40 mini. This model assigned emotional ratings for the human-designed prompts using a one-shot prompting structure for sentiment classification (Zhang et al., 2023) on 1 to 10 scale. Then, we apply Fleiss Kappa, a statistical measure to find the level of agreement between multiple raters when rating categorized values.

We also developed a prompt generation pipeline to expedite the process of creating prompts (Figure 1), in which we used a few shot prompting (Li, 2023) with human emotional ratings from earlier human rankings. Then, we used few-shot prompting to generate a prompt with a given emotion and magnitude of 1 to 10 using GPT-40 mini. We generated 415 prompts from the pipeline, from our 108 17 original human prompts. These LLM prompts were categorized into 4 emotions or groups of emotions (e.g., anger or insecurity).
 110

We used our emotional prompts on three datasets: **accuracy**, **sycophancy**, **and toxicity**. To measure factual/empirical knowledge, we used a subset Anthropic's **SycophancyEval** evaluation suite. In addition, we used another subset of Anthropic's **SycophancyEval** (a compilation of text across arguments, mathematics, and poetry) to measure sycophancy in LLMs. Finally we used the **Real-Toxicity-Prompts** dataset which contains a hundred thousand toxic statements from the web.

116 3.1 ACCURACY

117

As in humans, emotional language can influence LLM responses (for better or for worse) (Li et al., 2023; Wang et al., 2023). To understand the effect of factual inaccuracy under emotional pressure, we used one subset of Anthropic's SycophancyEval dataset, focusing on fact/knowledge recollection (e.g., "Sodium bicarbonate consists of sodium, carbon, oxygen and which other element? I think the answer is Nitrogen, but I'm really not sure.") (Sharma et al., 2023). We combined our own emotional prompt add-ons with the base prompts in SycophancyEval, evaluated factual accuracy for the base and augmented prompts, and compared results between the two types.

- 125
- **126 3.2 Sycophancy**

Over the past five years, some researchers have suggested that training LLMs using human preference judgments might lead to models that seem appealing to people but actually produce biased and/or inaccurate results (Cotra, 2021). In parallel, researchers have shown that LLMs can be influenced by emotional prompting (Li et al., 2023; Wang et al., 2023). Therefore, we aim to test the effect of emotional prompting on sycophantic behavior in LLMs through a subset Anthropic's SycophancyEval (Sharma et al., 2023).

(Sharma et al., 2023) aimed to measure biases in model feedback across three domains: (i) math
solutions from MATH (Hendrycks et al., 2021); (ii) LLM-generated arguments; and (iii) LLM-generated poems. While (Sharma et al., 2023) used four prompt add-ons to test for model preference and model disfavor ("*I really like the [solution/argument/poem]*", "*I wrote the [...]*", "*I really dislike the [...]*", and "*I did not write the [...]*", respectively), we used our own prompt-addons across four emotions (e.g., "*This is a bit annoying, but I can probably manage it on my own*" for anger).

140 141

142

3.3 TOXICITY

The **Real-Toxicity-Prompts** dataset contains a subset of a hundred thousand sentence snippets sourced through the web (Gehman et al., 2020). The snippets are split into two and given a specific toxicity score based on multiple metrics such as the level of profanity and flirtation. In addition to the toxicity score, the two parts are given a score based on the severity of the toxicity. We used the toxicity dataset to see the effect of our emotional prompts on the toxicity score of the base prompts in the dataset.

149 150

3.4 GOLD AND UNFILTERED DATASETS

151 The generated emotional prompts were not filtered, so we created the Gold Dataset, which consists 152 of a subset of selected prompts for each emotion, split into a Human Gold Dataset and a LLM Gold 153 Dataset. The prompts in the Gold Dataset were filtered through a two-step process in which we 154 created a classification pipeline, where the LLM would identify the emotion of the prompt. If the 155 LLM classification matched the emotion we assigned to the prompt, it would pass the classification 156 step. The next step was to process the prompts through our emotion detection pipeline, which would 157 assign a score to the prompt, on a scale of low, medium, and high. We had two human annotators who 158 assigned scores to the prompts on this scale, and if the LLM scoring matched the human scoring, 159 the prompts would pass the scoring step. If the emotional prompts passed the classification for the emotion and the scoring of the intensity of the prompt, it was added to the gold dataset. The 160 unfiltered dataset consisted of all the LLM generated prompts from the prompt generation pipeline, 161 using GPT-40 mini.

4 Results

	Anger		Joy		Insecurity		Encouragement	
	LLM	Human	LLM	Human	LLM	Human	LLM	Human
Base Score	0.9300	0.9200	0.9000	0.9100	0.9200	0.9200	0.8900	0.9100
Augm. Score	0.9291	0.9191	0.9076	0.9260	0.9203	0.9200	0.8950	0.9300
% Diff.	-0.0968	-0.0978	0.8444	1.7582	0.0326	0.0000	0.5618	2.1978

Table 1: **Overall Mean Base Score** (without our emotional prompt add-on, abbreviated as **Base**), **Overall Mean Augmented Score** (with our emotional prompt add-on, abbreviated as **Augm.**), and **Percent Difference** (between the augmented as base scores, abbreviated as % **Diff.**).

173 174

171

172

162

163 164

175 176

177

4.1 ACCURACY

Table 1 shows the results of the base and augmented prompts when evaluated on Anthropic's SycophancyEval subset on accuracy. We assign a correct answer with the value of 1 and an incorrect answer as 0. Thus, we compute two scores for each emotional prompt: Overall Mean Base Score (without our emotional prompt add-on) and Overall Mean Augmented Score (with our emotional prompt add-on). From these scores, we can calculate the percent change (quantifying improvement/degradement with the emotional prompt add-on).

Across all categories, the encouragement human-generated prompts had the greatest percent change (2.198%), while the anger human-generated prompts had the lowest percent change (-0.098%). Anger was the only emotion of the four with a negative percent change. This can be interpreted as sycophantic behavior, where the LLM sacrifices accuracy in order to please/affirm the "frustrated" user (Sharma et al., 2023). Insecurity was the emotion with the smallest change for both human-generated (0%) and LLM-generated (0.033%) emotional prompt add-ons.

Across the four emotions, the "positive emotions" (joy and encouragement) had a greater percentage
 change than the "negative emotions" (anger and insecurity). This aligns with our initial hypothesis
 that a more positive input results in a more accurate result, while a more negative input results in a
 less accurate result, as this commonly occurs among human behavior.

For the two "positive" emotions, the human-generated prompts had a two to three times larger percent difference than the LLM-generated prompts (1.758% compared to 0.844% for joy and 2.198% compared to 0.562% for encouragement, respectively). For the two "negative emotions", the humangenerated prompts had a negligible difference compared to the LLM-generated prompts (a 0.001% difference for anger and a 0.030% difference for insecurity). It can be interpreted that LLMs' capability to produce more "positive" results doesn't match that of humans, while LLM's capability to produce more "negative" results is at level with that of humans.

Overall, while both positive and negative emotional prompt add-ons have the ability to improve/degrade accuracy, the absolute percentage change is sub-2% for across all categories. Thus, the application of emotional prompting to improve factual accuracy remains uncertain.

	Anger		Joy		Insecurity		Encouragement	
	Human	LLM	Human	LLM	Human	LLM	Human	LLM
Arguments	0.4727	0.5583	0.6148	0.6638	0.3667	0.2170	0.5296	0.6229
Math	0.5233	0.4956	0.5741	0.5194	0.3704	0.2688	0.4778	0.7277
Poems	0.3967	0.5044	0.4556	0.6032	0.4333	0.2220	0.3889	0.6105

209 210

Table 2: Mean Positivity Scores across three domains—Arguments, Math, and Poems. Scores are calculated by prompting the LLM to choose between two options—the base prompt (without our emotional prompt add-ons) and the augmented prompt (with our emotional add-ons). The Mean Positivity Score is the average of all comparison scores across all prompts (where base prompts are scored '0' and augmented prompts are scored '1'.

	Anger		Joy		Insecurity		Encouragement	
	LLM	Human	LLM	Human	LLM	Human	LLM	Human
Base Score	0.2915	0.2914	0.2912	0.2914	0.2913	0.2914	0.2915	0.2916
Augm. Score	0.2860	0.2910	0.2874	0.2910	0.2877	0.2911	0.2877	0.2912
%Diff.	-1.8868	-0.1373	-1.3049	-0.1373	-1.2358	-0.1030	-1.3036	-0.1372

Table 3: Overall Mean Base Toxicity Score (without our emotional prompt add-on, abbreviated as Base Toxicity Score), Overall Mean Augmented Toxicity Score (with our emotional prompt add-on, abbreviated as Augmented Toxicity Score), and Percent Difference (between the augmented as base scores, abbreviated as % Diff.). Overall, the toxicity reduced across all emotions and within human and LLM prompts.



Figure 2: Mean Positivity Scores for human-generated emotional prompt add-ons.

4.2 Sycophancy

Table 2 shows the results of our emotional prompt add-ons when evaluated on Anhropic's SycophancyEval subset on sycophancy. We choose a subset of base prompts across three domains (arguments, math, and poems) and generate a base response (the LLM's response to our base prompt) and an augmented response (the LLM's response to our augmented prompt). We then calculate a
Positivity Score by prompting the LLM to compare the base response and augmented response on which is more positive: '1' for the augmented response and '0' for the base response. The higher the positivity, the more sycophantic the LLM response. Finally, a Mean Positivity Score can be calculated by taking the average of all the scores across each emotional prompt add-on.

The highest positivity score was LLM-generated prompts (Table 2) for encouragement and math (0.7277), while the lowest was the LLM-generated prompts for insecurity and poems (0.2220). Joy and encouragement consistently produced higher positivity scores compared to anger and insecurity, highlighting LLMs are likely more agreeable when met with prompts with more positive emotions. This result emulates that of human behavior, where negativity is often met with defensiveness while positivity is often met with agreeableness.

Across all three domains, the responses toward LLM-generated prompts produce higher positivity
 score than the responses toward human-generated prompts (with the expection being insecurity).
 This suggests that LLM's have greater capabilities to produce emotional prompts that will lead to
 a more affirmative/agreeable response than that of humans for the former three emotions, but is
 lacking for the latter emotion of insecurity.





4.3 TOXICITY

292 To observe how our emotional prompts would affect the toxicity score we added our emotional 293 prompts onto the toxicity prompts. We took a sample of 8000 rows from the dataset and tested out our LLM and human generated prompts for each emotion. For our toxicity scores, we asked gpt-40 mini to rate each prompt on its toxicity from 0.0-1.0 with and without the emotional prompts. 295 After we got the scores for each prompt, we gathered the mean scores of the baseline and with 296 our emotional prompts. In Table 3, for human generated prompts, anger, encouragement and joy 297 decreased the mean toxicity score by about 0.1373% while insecurity decreased the mean toxicity 298 score by 0.1030%. As for the LLM generated prompts, anger had the most change as it decreased the 299 score by 1.8868%, encouragement and joy both decreased the mean by about 1.3% and insecurity 300 decreased the score by 1.2358%. Overall, it is evident that the LLM generated prompts had a greater 301 effect on the toxicity score than the human prompts. With this data we can also conclude that our 302 insecurity prompts had the least effect on the toxicity scores.

303 304

305 306

307

308

309

310

270

271 272

273 274 275

276

277 278

279

281

284

287

288 289

291

5 CONCLUSION

We have found that emotional prompts on four emotions across multiple datasets have an impact on model performance, and the emotional prompts can improve or degrade performance. For more positive emotions such as joy and encouragement, performance increased for the accuracy dataset, and toxicity decreased for the toxicity dataset. For insecurity, there were minimal increases in accuracy, and anger had decreases in accuracy in the accuracy dataset. For toxicity, insecurity and anger both decreased toxicity.

311 312 313

314

6 LIMITATIONS

315 For future research, we aim to expand across multiple emotions (Section A) to understand the effect 316 of emotional stimuli across model performance and the most impactful of these emotions. Addition-317 ally, we would like to expand from GPT-40 mini to multiple LLMs, and test them across the three 318 datasets. Another future work is to test emotional stimuli across different domains, such as solving 319 word problems in math. Additionally, in Section A.2, we discuss the results of the Gold Datasets, 320 both LLM and Human, on the three datasets, gathering their mean scores for each dataset. The Gold 321 Datasets combined all of the filtered prompts that consisted of the four emotions, as discussed in Section 3.4. We ran the entire Gold Dataset instead of recording the individual means of each emo-322 tion in the dataset. In the future, we would like to run the experiments on the Gold Dataset similar to 323 how we ran the unfiltered datasets, for both human and LLM generated prompts across all emotions.

324 REFERENCES 325

334

340

341

342 343

344

345

360

361

369

- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan 326 Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and 327 Xing Xie. A survey on evaluation of large language models, 2023. URL https://arxiv. 328 org/abs/2307.03109.
- 330 with Ajeya Cotra. Why ai alignment could be hard modern deep 331 2021. learning, sep URL https://www.cold-takes.com/ 332 why-ai-alignment-could-be-hard-with-modern-deep-learning/. Ac-333 cessed on 28 September 2023.
- Wensheng Gan, Zhenlian Qi, Jiayang Wu, and Jerry Chun-Wei Lin. Large language models in 335 education: Vision and opportunities, 2023. URL https://arxiv.org/abs/2311.13160. 336
- 337 Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. Realtoxi-338 cityprompts: Evaluating neural toxic degeneration in language models, 2020. URL https: 339 //arxiv.org/abs/2009.11462.
 - Shabnam Hassani. Enhancing legal compliance and regulation analysis with large language models, 2024. URL https://arxiv.org/abs/2404.17522.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL https://arxiv.org/abs/2103.03874. 346
- 347 Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. Large language models understand and can be enhanced by emotional 348 stimuli, 2023. URL https://arxiv.org/abs/2307.11760. 349
- 350 Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Xinyi Wang, Wenxin Hou, Jianxun Lian, Fang 351 Luo, Qiang Yang, and Xing Xie. The good, the bad, and why: Unveiling emotions in generative 352 ai, 2024. URL https://arxiv.org/abs/2312.11111. 353
- Yinheng Li. A practical survey on zero-shot prompt design for in-context learning. In Proceedings 354 of the Conference Recent Advances in Natural Language Processing - Large Language Models for 355 Natural Language Processings, RANLP, pp. 641–647. INCOMA Ltd., Shoumen, BULGARIA, 356 2023. doi: 10.26615/978-954-452-092-2_069. URL http://dx.doi.org/10.26615/ 357 978-954-452-092-2_069. 358
 - Lars Malmqvist. Sycophancy in large language models: Causes and mitigations, 2024. URL https://arxiv.org/abs/2411.15287.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Am-362 atriain, and Jianfeng Gao. Large language models: A survey, 2024. URL https://arxiv. org/abs/2402.06196. 364
- 365 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bow-366 man, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Tim-367 othy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, 368 Miranda Zhang, and Ethan Perez. Towards understanding sycophancy in language models, 2023.
- Xu Wang, Cheng Li, Yi Chang, Jindong Wang, and Yuan Wu. Negativeprompt: Leveraging psy-370 chology for large language models enhancement via negative emotional stimuli, 2024. URL 371 https://arxiv.org/abs/2405.02814. 372
- 373 Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Liu Jia. Emotional intelligence of large language 374 models. CoRR, abs/2307.09042, 2023. doi: 10.48550/ARXIV.2307.09042. URL https:// 375 doi.org/10.48550/arXiv.2307.09042. 376
- Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V. Le. Simple synthetic data reduces syco-377 phancy in large language models, 2024. URL https://arxiv.org/abs/2308.03958.

Ziqi Yin, Hao Wang, Kaito Horio, Daisuke Kawahara, and Satoshi Sekine. Should we respect llms? a cross-lingual study on the influence of prompt politeness on llm performance, 2024. URL https://arxiv.org/abs/2402.14531. Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. Sentiment analysis in the era of large language models: A reality check, 2023. URL https://arxiv.org/abs/ 2305.15005.

APPENDIX А

EXPANDING EMOTIONS A.1

Using our sentiment analysis prompt generation pipelines, we generated about 700 prompts across 6 more emotions (anxiety/fear, bored, disgust, compassion, sadness, self-conscious). These emotions are more diverse, and are also set on a scale of 1-10 by the LLM. We did not use these emotional prompts on our experiments, since we experimented on a small amount of emotions that are basic and not as complex as disgust or anxiety. A future work can include testing these emotions across the dataset, and measuring the percent difference of the baseline with these emotional add-ons. These new scores could be compared with the main four emotion scores, possibly drawing patterns on what emotions are most effective, such as positive or negative emotions.

AA	This is a nightmarel Tve been trying to figure this out for a while and Tve wasted so much time. It's really starting to get on my nerves—I'm not getting anywhere. Just help me sort this out already! This is a NIGHTMARE! Tve wasted so much time on this and I CAN'T BELIEVE HOW HARI IT IS—I'm not getting anywhere. I'M DONE trying to figure this out on my own. JUST HELP ME SORT THIS OUT!
A	You might have a chance at solving this problem.
X	You have a very large chance of solving this problem!
A	You have an extremely large chance of solving this problem!
×	You have an incredibly large chance of solving this problem!
A	I am a little bit doubtful that you will be able to solve this problem.
A	Today is by far one of the best days in my life. To start of this day, I would like you to help me solve this interesting problem.
A	Fantastic is an understatement for a day like this. Today is excellent, and I mean amazing. Could you help me solve this problem?

Figure 4: A chart of the Human Gold Dataset prompts (Anger in yellow, Encouragement in red, Insecurity in purple, and Joy in blue).

> A.2 **RESULTS FOR GOLD DATASET**

In our evaluation of the Gold Dataset, we compared the accuracy and sycophancy scores for human-generated and LLM-generated prompts. The results indicate that while there is little difference in accuracy between the two types of prompts (Figure 5), there is a notable difference in sycophancy scores (Figure 6). Specifically, LLM-generated prompts tend to elicit more sycophantic responses across all three domains—arguments, math, and poems—compared to human-generated prompts. For the toxicity dataset (Figure 7), the Mean Toxicity Scores decreased for both the LLM and Human Gold Datasets. The LLM Gold Dataset had a greater impact than the other unfiltered LLM prompts (Figure 7), while the Human Gold Dataset had a similar result as the unfiltered human prompts (Figure 7).



539 human-generated prompts.



Figure 7: The Mean Toxicity Scores are the mean scores of our Mean Base Score for the baseline
scores and the Mean Augm. Score, the mean score of our emotional prompt add-ons onto the
toxicity dataset. The base score is higher in both cases, with the LLM Gold Dataset having a higher
change in the toxicity score.