SemAug: Shaping the Future of Semantically-Enriched, Format-Specific Data Augmentation

Anonymous ACL submission

Abstract

In the realm of artificial intelligence, the significance of high-quality data cannot be overstated, especially data that adheres to stringent for-004 matting rules and structures. Addressing this need, our study introduces an advanced data augmentation method specifically designed for format-specific datasets. This method utilizes the capabilities of Large Language Models 800 (LLMs) to generate data that not only meets the rigid formatting criteria but also maintains the integrity of the information. Central to our approach is the integration of specific format 013 requirements into natural language prompts, which guides the LLMs to produce precisely formatted outputs. A salient feature of our approach is its self-evaluative mechanism, which 017 autonomously assesses the semantic quality of the augmented data, distinguishing it from prior methodologies that require manual validation, thereby streamlining the augmentation process. Our research represents a pioneering step forward, enabling more efficient enhancement of datasets that demand exacting format adher-023 ence without the extensive resource investment 024 typically associated with such tasks.

1 Introduction

027

034

040

In the era of artificial intelligence (AI), large language models (LLMs) have emerged as transformative tools for natural language understanding and generation tasks. These models have demonstrated remarkable capabilities in diverse applications, including text completion, language translation, and text summarization. However, as these AI systems delve deeper into specialized domains, their adaptability and robustness in handling intricate and rigid data formats have come under scrutiny.

A pivotal challenge hindering the progress of robust AI lies in the intricacies of training data. While deep learning models flourish on vast and varied datasets, obtaining such data, especially in structured or complex formats, remains a daunting endeavor. The process is often hindered by cost, time, and occasional unfeasibility due to issues like privacy or rarity of examples. This has led to the exploration of data augmentation techniques to bolster training datasets. Extensive research has been conducted on data augmentation, and its widespread utilization in natural language processing (NLP) tasks has proven effective in addressing issues such as constrained data availability, uneven sample distribution, and domain-specific challenges like named entity recognition and relation extraction in biomedical text or code generation tasks. Through the incorporation of supplementary samples and introducing diversity, data augmentation contributes to enhancing model robustness and adaptability, simultaneously mitigating overfitting concerns.

042

043

044

045

046

047

051

052

056

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

081

Commonly employed data augmentation methods encompass the use of synonyms, antonyms, and variations in word forms, all of which serve to augment the quantity and diversity of training data (Niu and Bansal, 2018). Some traditional techniques include synonym replacement, random deletion and random insertions (Feng et al., 2021). Other methods utilize language models to generate reliable samples for more effective data augmentation, including back translation (Sennrich et al., 2016) and word vector interpolation in the latent space (Jindal et al., 2020). However, existing data augmentation methods often fall short when confronted with the nuances of complex data structures. These methods, while effective in general contexts, tend to struggle with preserving the specificity and intricacies inherent in complex formats (Feng et al., 2021; Shorten et al., 2021; Bayer et al., 2021).

Considering these challenges, the spotlight turns to prompt engineering, which holds promise as a solution. This research posits prompt engineering as a potential tool to harness the power of LLMs for data augmentation, particularly suited for intricate data formats. It seeks to explore how LLMs, when



Figure 1: Example of an input in Biomedcai Relation Extraction task. The entity is wrapped in the entity type's tags. Different colors mean different entities in the entity pair. The task is to extract relation between the entity pair.

guided by carefully crafted prompts, can synthesize data that not only maintains the rigidity and context of complex formats but is also virtually indistinguishable from genuine entries. The overarching aim is to bolster AI performance across an array of tasks, including natural language understanding, sentiment analysis, and content generation, ensuring they thrive even in contexts laden with complex data constraints.

In the realm of biomedical text data, the intricacies of rigid formats not only pose substantial challenges for data augmentation but also necessitate meticulous attention to preserve the inherent scientific and terminological precision. As shown in Figure 1, an example of training input to the model is a paragraph that contains an entity pair, with entity tags placing around each entity. The augmentation method should paraphrase the whole paragraph while keeping the entity tags and entity text the same, because biomedical named entities they are lexical sensitivities, which means that they should not be altered without meticulous scrutiny and analysis. This makes it complex for creating augmentation data in biomedical domain, which we should carefully consider the semantic delicacies in the biomedical entities.

Additionally, when it comes to structuring datasets from programming languages, the requirement for preserving syntax and semantic accuracy introduces a distinct set of challenges. In programming languages, every character, from a brace to a semicolon, carries significance; a minor alteration can lead to a completely different execution behavior. This sensitivity is especially pronounced in datasets that include numerical values or complex syntactic structures, common in software engineering and computational tasks. For instance, in a code snippet, changing a variable name or altering a loop structure could inadvertently change the algorithm's logic or functionality. This highlights the criticality of maintaining syntactic integrity and the original semantic context in any augmentation process applied to code datasets. The augmentation method must be intricately designed to understand and respect the language's grammar rules, operational semantics, and the contextual role of each code element. Moreover, in the realm of numerical data within code, precision is paramount. Numerical values often represent parameters, dimensions, or constants whose exact values are crucial for the correct functioning of the code. Altering these values without a deep understanding of their context and impact can lead to incorrect results or system failures, particularly in fields like numerical computing or algorithmic processing.

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

164

165

166

167

168

In this paper, we introduce SemAug (Specified-Rule **Aug**mentation and **Sem**antic-Quality), a novel approach combining the SpecRule Augmentor with the SemQ filter to overcome the challenges inherent in complex data domains, particularly in biomedical datasets. Unlike traditional data augmentation techniques that struggle with the intricate formats and rigid constraints of such datasets, SemAug employs prompt engineering to navigate these complexities effectively.

- SpecRule Augmentor under SemAug is designed to generate diverse and realistic data variations, adhering strictly to the scientific accuracy and terminological precision required in biomedical contexts.
- SemQ Filter rigorously evaluates the augmented data, ensuring its adherence to the high-quality standards essential for maintaining scientific validity. Our method contributes significantly to the field by enabling the augmentation of biomedical data in a manner that aligns seamlessly with its specialized requirements.

Moreover, SemAug facilitates the development of machine learning models proficient in handling and extracting insights from these enriched and complexly formatted datasets. This innovative approach marks a pivotal advancement in data augmentation, paving the way for enhanced application of machine learning and artificial intelligence in domains where data complexity and specificity have

118

084

169

170

171

172

173

174

175

176

177

178

179

180

183

184

186

188 189

190

191

192

196

197

204

209

210

212

213

214

215

traditionally posed significant challenges.

Related Works 2

Large Language Models 2.1

In recent developments within the domain of large language models (LLMs), their proficiency in mimicking human language through extensive training on expansive textual datasets has been welldocumented (Brown et al., 2020; Chowdhery et al., 2022; Hoffmann et al., 2022). Crucially, these models exhibit an exceptional ability to follow humanissued instructions, yielding significant zero-shot performance (Wei et al., 2022; Dong et al., 2023). In our research, we harness this particular trait of LLMs, employing them to generate datasets in specific formats by inputting prompts that encompass our defined requirements and rules. This methodology aligns with the broader narrative in LLM research, where the adaptability and responsiveness of these models to guided inputs are being increasingly explored for tailored data generation purposes.

2.2 Data Augmentation

Data augmentation, a common strategy to prevent overfitting in limited data scenarios, typically involves minor, label-preserving modifications to 193 the existing data (Wei and Zou, 2019; Fadaee et al., 2017). Recent studies have explored the use of Large Language Models (LLMs) for this purpose. Specifically, they use few-shot data as demonstrations and prompt LLMs to generate new data (Yoo et al., 2021; Sahu et al., 2022). The capability of LLMs to blend a few-shot dataset and create analogous data has been widely recognized. Lin et al. (2023) introduced the use of Pointwise V-information to sift through and remove non-beneficial data from these generated sets. Most recent approaches involve using models like ChatGPT and GPT-4 for data generation, leading to notable improvements in performance (Dai et al., 2023; Whitehouse et al., 2023). Furthermore, Cheng et al. (2023) demonstrated the efficacy of using GPT-3 generated data to enhance sentence embeddings through contrastive learning. Our methodology diverges by utilizing LLMs as an all-encompassing solution specifically for generating and evaluating augmented data within rigid or specific data formats, effectively obviating the requirement for manual annotation. 216

2.3 Few-shot and zero-shot learning

A zero-shot prompt poses a query to the LLM on 218 topics it hasn't been explicitly trained on. The term 219 "zero" in "zero-shot" indicates that the LLM lacks 220 specific training on the particular task or question 221 in the prompt. The "shot" denotes providing an 222 example to the LLM. Thus, "zero-shot" implies 223 that the model wasn't trained for the exact task or 224 question and the prompt doesn't offer an illustrative 225 example for reference. An instance of a zero-shot 226 prompt is translation tasks. Even though the model 227 might not have received particular training samples for translations, their vast linguistic training allows 229 them to generalize and provide probable transla-230 tions without dedicated training for the task. Few-231 shot prompts resemble zero-shot prompts in that 232 the model hasn't been tailored for the specific question or task. But, unlike zero-shot prompts, few-234 shot prompts offer an in-context example within 235 the request to guide the model.

217

237

238

239

240

241

242

243

244

245

247

248

249

250

251

252

253

254

255

Methodology 3

We apply the idea of using examples in few-shot learning; however, the largest difference here is that although examples are used in the prompt, there is no given label for them, instead, the focus on the format and meaning of the dataset is considered to be more important here. This means that the concept used here is zero-shot learning. Figure 2 shows the overflow of our method. There are two main steps here: creating data augmentation following requirements (SpecRule Augmentor) and self-evaluating the quality of generated data (SemQ Filter)

3.1 Data Augmentation with User's Requirements

Given a original dataset containing n samples $D_o = \{d_1, d_2, ..., d_n\}$, the data augmentation D_a can be obtained through:

$$D_a = \{ \text{SpecRuleAugmentor}(d_i) \, | \, d_i \in D_o \}$$
(1)

where SpecRuleAugmentor is a specified prompt used in the LLM in Figure 3



Figure 2: Overall Flow of The Framework

You are creating data augmentation through <**REQUEST**>. The data created must be in the format: <**OUTPUT'S FORMAT RULE**>. For example 1...n, given <**INPUT EXAMPLES + FORMAT'S RULE**>. Here is your task, given input: <**INPUT**>, please follow the rules above and <**REQUEST**>

Figure 3: Requirement-Specified Prompt for Generating Rigid Data Augmentation (SpecRule Augmentor)

Instead of focusing on letting the model know how to create a pair of training data (input + label), the prompt concentrates on emphasizing the model to create data following requirements of correct format. An example of SpecRule Augmentor applied in Biomedical Relation Extraction and Code Generation task is shown in Figure 6 and Figure 7 in Appendix correspondingly.

3.2 Self-Evaluative Data Augmentation's Quality

Following augmentation, we apply the SemQ Filter to each document $d_{i'}$ in D_a to assess its quality. The SemQ Filter functions as a prompt-based scoring mechanism, as illustrated in Figure 4, and assigns a quality score to each document. This process is represented as:

271

272

274

277

278

 $Score(d_{i'}) = SemQFilter(d_{i'}) \text{ for each } d_{i'} \in D_a$ (2)

To ensure the augmented dataset's quality, We apply a threshold-based filtering, retaining only those documents from D_a with a SemQ Filter score above 4. The augmented dataset, D_{filtered} , consists of documents from D_a that meet the quality criterion below:

$$D_{\text{filtered}} = \{ d_{i'} \in D_a \, | \, \text{Score}(d_{i'}) > 4 \} \quad (3)$$

The final dataset D_{final} is formed by combining the original dataset D_o with the high-quality, filtered augmented data D_{filtered} :

$$D_{\text{final}} = D_o \cup D_{\text{filtered}} \tag{4}$$

Your task is to conduct a semantic alignment evaluation between an 'Original' paragraph and its 'Paraphrased' counterpart. Carefully analyze how well the paraphrased version maintains the meaning, intent, and key information of the original. After your assessment, provide a score exclusively in the range of 0 to 5. Here, 0 represents a complete loss of the original meaning (indicating no alignment), and 5 signifies semantic an excellent preservation of the original's meaning (indicating perfect semantic alignment). Here is your task, given 'Original' <ORIGINAL>: 'Paraphrased' version: version: <PARAPHRASED>, please evaluate their semantic meaning and provide only a single number between 0 and 5 that reflects the level of semantic consistency between the original and paraphrased texts.

Figure 4: Semantic Checking Prompt for Quality Assurance (SemQ Filter)

The SemQ Filter prompt tells the model to do a self-evaluation on both format and semantic checking. Moreover, instead of using "True", "False" evaluation in the prompt, we see that telling the model to give a score from 0-5 can help identify 279

281

282

291 292

296

297

299

301

302

303

304

305

307

310

325

327

329

330

331 332

333

real high-quality data in the generated samples. Example of SemQ scores given for different generated samples can be seen in Table 5 in Appendix.

Experiments 4

4.1 Dataset

BioRED Dataset (Luo et al., 2022): The task is to extract biomedical relation across multiple entity types with multiple relation types. There are multiple entity types (gene, chemical, disease, ...) and relation pairs (disease-gene, chemical-chemical) at the document level. There are a total of 600 PubMed abstracts, with 400 documents for training, 100 for dev and 100 for validation.

EvolInstruct-80k-v1 dataset¹ : EvolInstruct-80kv1 dataset (license cc-by-nc-sa-4.0) is an opensource dataset generated using Evol-Teacher (Luo et al., 2023). The dataset contains about 78000 instructions and their code solution. We use this dataset to fine-tune model for Code Generation task.

HumanEval (Chen et al., 2021): HumanEval contains 164 handwritten Python programming prob-312 lems with a function signature, docstring, body, and several unit tests. It focuses on testing the execu-314 315 tion ability of the code using the evaluation metric pass@k.

MultiPL-E dataset (Cassano et al., 2022): 317 MultiPL-E is a benchmark for evaluating large language models for code generation that supports 18 319 programming languages. It takes the OpenAI "HumanEval" Python benchmark and uses little com-321 pilers to translate them to other languages. In this paper, we do experiments on C++, Java, JavaScript 323 (JS), R, Bash, Rust.

> Grade School Math 8K (GSM8K) (Cobbe et al., 2021): GSM8k is a dataset of 8.5K high-quality linguistically diverse grade school math word problems. The dataset was created to support the task of question answering on basic mathematical problems that require multi-step reasoning.

4.2 Experimental Design

Data augmentation: We use open-source LLM OpenChat-3.5 (Wang et al., 2023) as the model for SpecRule Augmentor and SemQ Filter.

Biomedical Relation Extraction: BioRED dataset is used for training and evaluating. We fine-tune 337 PubMedBERT (Gu et al., 2020) in 50 epochs with

> https://huggingface.co/datasets/nickrosh/ Evol-Instruct-Code-80k-v1.

Table 1: Data augmentation in different dataset domains. Ori: Original, Aug: After Augmentation, Sem: After semantic filtering.

Dataset	Ori	Aug	Sem	Success
				rate
BioRED	28233	87519	60188	53.9%
			(+31955)	
Evol	78381	156762	133568	70.4%
Instruct			(+55187)	

338

339

340

341

342

343

344

345

346

347

348

350

351

352

353

354

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

learning rate 1e-5, batch size 10.

Code Generation: For fine-tuning our model, we utilized the EvolInstruct-80k dataset, selecting Deepseek-Coder-1.3b-Instruct (DeepSeek, 2023) as the foundational model. The model is finetuned in 800 steps, with learning rate 2e-5. For evaluation purposes, we employed datasets including HumanEval, MultiPL-E, GSM8K. Since there is a difference between the performance stated in DeepSeek (2023) and those obtained through our use of the bigcode-evaluate-harness (Ben Allal et al., 2022), we opted to rely on the performance evaluations derived from bigcode-evaluate-harness for all models to ensure consistency in comparative analysis. Temperature is set to 0.2 for all models, and the primary metric used for evaluating model performance in this task was pass@1. Our model's performance was benchmarked against several others, including CodeGeeX2 6B (Zheng et al., 2023), StarCoder 15B (Li et al., 2023), CodeLlama 7B, CodeLlama-Instruct 7B, CodeLlama 13B (Roziere et al., 2023).

Mathematics Reasoning: We additionally assess the influence on the performance of a model that has been fine-tuned for the code generation task when applied to a different task. This involves comparing a model already fine-tuned for code generation with various other Large Language Models, including Llama 7B, Llama 13B (Touvron et al., 2023), Minerva 8B (Lewkowycz et al., 2022), CodeGeeX2 7B, StarCoder 15B, and Deepseek-Coder-1.3b-Instruct.

4.3 Results

Generated data: In this study, we embarked on 371 an extensive data augmentation initiative spanning 372 two distinct and challenging domains: biomedical 373 relation extraction (BioRED) and code generation 374 (Evol Instruct). Table 1 presents a comprehensive 375 overview of the data augmentation output and sub-376

sequent filtration efficacy across these domains. 377 Notably, the success rates post-semantic validation-courtesy of our SemO filter-demonstrate task-specific variability, which we attribute to the inherent complexity and structural rigidity of each dataset. An example result of SemQ can be seen in Table 5. For the biomedical relation extraction task, the data augmentation process involves intricate manipulation of text interspersed with rigorously formatted tags such as @GeneOrGeneProduct and @ChemicalEntity. These tags, which require precise placement and pairing of opening (@GeneOrGeneProduct) and closing (@/GeneOr-GeneProduct) markers (detailed example can be seen in Table 6 in Appendix), significantly elevate 391 the complexity of the augmentation task. Consequently, a substantial proportion of the augmented dataset failed to meet the stringent criteria set forth by the SemQ Filter, as it rigorously scrutinizes the congruence of these tags within the augmented narratives. This meticulous semantic checking process is pivotal in ensuring the scientific accuracy of the augmented data, a non-negotiable aspect in biomedical research where precision is paramount. 400

401 On the other hand, the Code Generation task, while still subject to strict syntactic and seman-402 403 tic constraints, exhibited a higher pass rate postaugmentation. In this domain, the augmentation 404 process must ensure that the syntactical structure 405 and numerical values within the code remain un-406 altered to preserve the original logic and function-407 ality. Despite these constraints, the augmentation 408 method proved to be more adaptable to the pro-409 gramming context, resulting in a greater volume of 410 data successfully passing through the SemQ Filter. 411 This higher success rate suggests that our augmen-412 tation approach is effectively attuned to the nuances 413 of programming languages, capable of generating 414 syntactically sound and semantically coherent code 415 snippets that align with the original intention of the 416 code's author. The disparity in the success rates 417 between the two tasks underscores the challenges 418 that specialized domains pose to data augmentation 419 techniques. In domains like biomedical research, 420 where the data is densely annotated and laden with 421 domain-specific terminologies, the augmentation 499 process must navigate a complex landscape of lin-423 guistic and semantic rules. Similarly, in program-424 ming language datasets, the augmentation must 425 contend with the rigid syntax and operational se-426 mantics that are fundamental to the code's execu-427



Figure 5: Performance on Math Reasoning Task.

tion. These findings highlight the need for augmentation strategies that are not only context-aware but also capable of discerning and adhering to the nuanced rules that govern data in specialized fields. 428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

Biomedical Relation Extraction Task: Table 2 shows that our methodology demonstrates substantial improvement across all three evaluation schemas. When augmented data is introduced, SemAug not only excels in extracting entity pairs within a relation but also in discerning the type of relation and further characterizing the novelty of the entity pairs. This is evident from the marked 8.8% enhancement in the F1-score, showcasing the robustness of our approach in handling the intricacies of biomedical data and contributing to a more nuanced understanding and extraction of relationships within it.

Code Generation Task: It is clear from Table 3 that SemAug demonstrates a significant enhancement over the base model, Deepseek-Coder-Instruct, across various programming languages. Particularly notable are the improvements in languages like Java and C++, where SemAug shows an increase of 5.0% and 7.0% respectively, suggesting that our method has effectively leveraged the underlying LLM to produce superior data augmentation outcomes. The consistent increments across multiple languages affirm the robustness and versatility of SemAug, validating its potential as a powerful tool for tasks requiring nuanced language comprehension and generation.

Mathematics Reasoning Task: SemAug was fine-tuned exclusively on a code instruction dataset

Table 2: Evaluation results on RE task comparison when added augmented data: extracting the entity pairs within a relation, second schema: extracting the entity pairs and the relation type and the third schema: further labeling the novelty for the extracted pairs. All numbers are F-scores. The $\langle G, D \rangle$ is the concept pair of the gene (G) and the disease (D). The columns of those entity pairs present the RE performance in F-scores. Bold denotes best model

Eval Schema	Methods	All	<g,d></g,d>	<g,g></g,g>	<g,c></g,c>	<d,v></d,v>	<c,d></c,d>	<c,v></c,v>	<c,c></c,c>
	BERT-GT	72.1	63.8	78.5	77.7	69.8	76.2	58.5	74.9
Entity pair	PubMedBERT	72.9	67.2	78.1	78.3	67.9	76.5	58.1	78.0
	SemAug	78.3	76.1	86.2	85.7	66.0	78.2	38.1	85.9
		↑ 5 .4	↑ 8.9	$\uparrow 7.1$	↑7.4	↓ 1.9	↑ 1.7	↓ 20	↑ 7.9
	BERT-GT	56.5	54.8	63.5	60.2	42.5	67.0	11.8	52.9
+Relation type	PubMedBERT	58.9	56.6	66.4	59.9	50.8	65.8	25.8	54.4
	SemAug	64.4	64.3	76.1	64.3	47.0	72.2	9.5	59.7
		↑ 5.5	↑ 7.7	↑ 9.7	↑4.4	↓ 3.8	↑ 6 .4	↓ 16.3	↑ 5 .3
	BERT-GT	44.5	37.5	47.3	55.0	36.9	51.9	11.8	48.5
+Novelty	PubMedBERT	47.7	40.6	54.7	54.8	42.8	51.6	12.9	50.3
	SemAug	56.5	51.1	69.1	57.9	44.9	78.2	9.5	56.5
		$\uparrow 8.8$	↑ 10.5	† 14.4	† 3 .1	↑ 2. 1	† 26.6	↓ 3.4	↑ 6.2

Table 3: Evaluation results on Completion Task. The evaluation metric is pass@1. Bold denotes best model. A demarcation line is placed between the final two models to emphasize our primary objective is to compare the enhancements our method brings relative to the base model.

Model	Size	HumanEval	MultiPL-E					
	Silt		C++	Java	R	Rust	Bash	JS
CodeGeeX2	6B	33.5	29.2	23.5	6.8	17.9	6.3	24.8
StarCoder	16B	30.4	28.6	28.5	13.6	21.2	9.3	31.7
CodeLlama	7B	30.0	25.5	29.2	16.8	26.3	8.1	31.8
CodeLlama-Istruct	7B	45.7	8.7	28.8	14.9	26.3	8.7	33.1
CodeLlama	13B	35.1	34.1	32.2	22.4	28.2	11.8	38.3
Deepseek-Coder-Instruct	1.3B	51.8	29.2	31.0	12.4	23.6	7.5	39.1
SemAug	1.3B	53.7	34.2	38.0	13.7	24.4	11.2	39.1
		↑ 1.9	$\uparrow 5.0$	\uparrow 7.0	↑ 1.3	$\uparrow 0.8$	↑ 3 .7	$\uparrow 0.0$

that did not contain data explicitly related to 463 mathematics. As shown in Figure 5, despite the 464 dataset's focus being misaligned with the task, 465 SemAug exhibits commendable resilience, with 466 only a marginal decrease in performance post-fine-467 tuning. This slight dip, however, underscores the 468 importance of task-aligned data for fine-tuning, 469 to fully harness the model's capabilities in 470 domain-specific tasks. The fine-tuning process, 471 although not directly contributing to performance 472 in this task, demonstrates the model's flexibil-473 ity and provides insights into the transferability 474 of learned representations across different contexts. 475

476

5 Conclusion and Future Work

This paper introduces a groundbreaking approach to format-specific data augmentation, uniquely coupled with a self-evaluative mechanism. These datasets are characterized by their strict structural requirements, which might include specific tagging systems in biomedical data or syntactic rules in programming languages. Unlike general, nonformat-specific datasets, which can accommodate a broader range of augmentation techniques, formatspecific datasets demand precise adherence to their unique formatting rules to retain data integrity and utility. 477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

Our method distinguishes itself by seamlessly adapting to these stringent requirements across diverse domains without the need for extensive mod-

ifications to the existing system architecture. The 493 key to this flexibility is the strategic alteration of 494 natural language prompts, which guides the aug-495 mentation process to produce outputs that conform 496 to the particular format constraints. This capability 497 illustrates the versatility and efficacy of SemAug 498 in addressing the challenges posed by structured 499 datasets. Additionally, a key advancement of SemAug over previous methods is its autonomous capability to evaluate and ensure the semantic qual-502 ity of augmented data. Traditional approaches of-503 ten rely on human oversight, which can be time-504 consuming and prone to error. SemAug circum-505 vents this dependency, offering a more efficient and reliable means of data augmentation that mini-507 mizes the need for manual intervention. Moreover, the accessibility of our methodology to non-experts is another significant aspect. By reducing reliance 510 on extensive programming knowledge, it opens 511 up opportunities for a broader range of users to 512 engage in advanced data processing and augmenta-513 tion. This democratization of technology can lead to innovative applications and solutions, previously 515 unattainable due to technical barriers. 516

In conclusion, this research not only provides a novel technical solution but also fosters a more inclusive approach to complex data processing. As we continue to explore and refine this method, we anticipate it will unlock new possibilities and applications, significantly contributing to the advancement of various fields that rely on specialized data formats.

Limitations

518

519

521

523

524

The effectiveness of our method is closely tied 526 to the performance of Large Language Models (LLMs). This implies that the quality of our data augmentation output, as well as the efficacy of the self-evaluation process, is contingent on the accu-530 racy and reliability of the LLMs' results. Hence, suboptimal performance from these models could 532 potentially impact the overall quality and reliability 533 of the augmented data generated by our method. 534 In this study, we have utilized a single LLM for data augmentation purposes. Future investigations 536 should explore the application of various LLMs to 537 determine their influence on the ultimate effective-538 ness of the model.

- 540 Ethics Statement
- 541 We comply with the ACL Code of Ethics.

References

Markus Bayer, Marc-André Kaufhold, and Christian Reuter. 2021. A survey on data augmentation for text classification. *CoRR*, abs/2107.03158. 542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

598

599

- Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von Werra. 2022. A framework for the evaluation of code generation models. https://github.com/bigcode-project/ bigcode-evaluation-harness.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.
- Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. 2022. Multipl-e: A scalable and extensible approach to benchmarking neural code generation.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code.
- Qinyuan Cheng, Xiaogui Yang, Tianxiang Sun, Linyang Li, and Xipeng Qiu. 2023. Improving contrastive learning of sentence embeddings from ai feedback.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben
- 8

712

713

714

715

657

658

659

Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways.

610

611

612

615

616

619

622

627

634

638

640

641

647

655

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.
- Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Yihan Cao, Zihao Wu, Lin Zhao, Shaochen Xu, Wei Liu, Ninghao Liu, Sheng Li, Dajiang Zhu, Hongmin Cai, Lichao Sun, Quanzheng Li, Dinggang Shen, Tianming Liu, and Xiang Li. 2023. Auggpt: Leveraging chatgpt for text data augmentation.
- DeepSeek. 2023. Deepseek coder: Let the code write itself. https://github.com/deepseek-ai/ DeepSeek-Coder.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. A survey on in-context learning.
- Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2017. Data augmentation for low-resource neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 567–573, Vancouver, Canada. Association for Computational Linguistics.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. A survey of data augmentation approaches for NLP. In *Findings of the Association* for Computational Linguistics: ACL-IJCNLP 2021, pages 968–988, Online. Association for Computational Linguistics.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. Domain-specific language model pretraining for biomedical natural language processing.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland,

Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. Training compute-optimal large language models.

- Amit Jindal, Arijit Ghosh Chowdhury, Aniket Didolkar, Di Jin, Ramit Sawhney, and Rajiv Ratn Shah. 2020.
 Augmenting NLP models using latent feature interpolations. In Proceedings of the 28th International Conference on Computational Linguistics, pages 6931– 6936, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with language models.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2023. Starcoder: may the source be with you!
- Yen-Ting Lin, Alexandros Papangelis, Seokhwan Kim, Sungjin Lee, Devamanyu Hazarika, Mahdi Namazifar, Di Jin, Yang Liu, and Dilek Hakkani-Tur. 2023. Selective in-context data augmentation for intent detection using pointwise V-information. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pages 1463–1476, Dubrovnik, Croatia. Association for Computational Linguistics.
- Ling Luo, Po-Ting Lai, Chih-Hsuan Wei, Cecilia N Arighi, and Zhiyong Lu. 2022. BioRED: a rich biomedical relation extraction dataset. *Briefings in Bioinformatics*, 23(5):bbac282.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evolinstruct.

Tong Niu and Mohit Bansal. 2018. Adversarial oversensitivity and over-stability strategies for dialogue models.

716

718

720

724

725

730

731

732

733

734

735 736

737

738

739

740

741

742

743

744

745 746

747

748

749

750

751

753

754

755 756

757

762

763

764

772

- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950.
- Gaurav Sahu, Pau Rodriguez, Issam Laradji, Parmida Atighehchian, David Vazquez, and Dzmitry Bahdanau. 2022. Data augmentation for intent classification with off-the-shelf large language models. In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 47–57, Dublin, Ireland. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.
- Connor Shorten, Taghi M. Khoshgoftaar, and Borko Furht. 2021. Text data augmentation for deep learning. *Journal of Big Data*, 8.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.
- Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2023. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent abilities of large language models.
- Jason Wei and Kai Zou. 2019. EDA: Easy data augmentation techniques for boosting performance on text classification tasks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 6382–6388, Hong Kong, China. Association for Computational Linguistics.
- Chenxi Whitehouse, Monojit Choudhury, and Alham Fikri Aji. 2023. Llm-powered data augmentation for enhanced cross-lingual performance.
- Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-Woo Lee, and Woomyoung Park. 2021. GPT3Mix: Leveraging large-scale language models for text augmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2225–2239,

Punta Cana, Dominican Republic. Association for Computational Linguistics. 773

774

780

Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan775Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang,
Yang Li, Teng Su, Zhilin Yang, and Jie Tang. 2023.776Codegeex: A pre-trained model for code generation
with multilingual evaluations on humaneval-x.779

A Appendix

r	Request
Format's rule	You are creating data augmentation through paraphrasing a given paragraph. The data created must be in the format: paraphrased paragraph that has @typeSrcS entity1 @/typeSrcS and @typeTgtS entity2 @/typeTgtS the same as the given paragraph (don't paraphrase the entity and special type tags around it), where type belongs to [GeneOrGeneProduct, DiseaseOrPhenotypicFeature, ChemicalEntity]. It means that @typeSrcS can be [@GeneOrGeneProductSrc, @DiseaseOrPhenotypicFeatureSrc, @ChemicalEntitySrc], @/typeSrcS: [@/GeneOrGeneProductSrc, @/DiseaseOrPhenotypicFeatureSrc, @/ChemicalEntitySrc], @/typeTgtS: [@GeneOrGeneProductTgt, @DiseaseOrPhenotypicFeatureTgt, @ChemicalEntityTgt], @/typeTgtS: [@GeneOrGeneProductTgt, @/DiseaseOrPhenotypicFeatureTgt, @/ChemicalEntityTgt], @/typeTgtS:
Examples + format's rule Request	paragraph: <example 1=""> ;the paraphrased paragraph should still keep <specific format=""> in the paraphrased paragraph. For example 2, <example 2="">, the paraphrased paragraph should still keep <specific format=""> in the paraphrased paragraph. Here is your task: Given paragraph: "<document>"; please follow the rules above and paraphrase the given paragraph, return only **one** paraphrased paragraph with **entity tags**, no other comments.</document></specific></example></specific></example>

Figure 6: SpecRule Augmentor in Biomedical Relation Extraction task.

ſ	EvolInstruct dat	taset Request
[Format's rule	You are creating data augmentation through paraphrasing a given instruction. The data created must be in the format: retain the technical and numerical aspects of the instruction, ensuring no change to specific programming or numerical
	Examples +	terms. For example, given instructions like 'Use a for loop to iterate over a list' or 'Print numbers' from 0 to 10 excluding
	format's rule	task: Given input: <instruction>, please follow the rules above and paraphrase the given instruction, output only result, no</instruction>
		other comments needed. Request

Figure 7: SpecRule Augmentor in Code Instruction task.

Model	Size	Accuracy
LLama	7B	11.0
LLama	13B	17.8
Minerva	8B	16.2
CodeGeeX2	7B	22.2
StarCoder	16B	23.4
Deepseek-Coder-Instruct	1.3B	34.9
SemAug	1.3B	30.0

Table 4: Performance on Math Reasoning Dataset.

Table 5: Generated samples on EvolInstruct-80k-v1 dataset.	Highlights mean specific syntax or number that are
kept unchanged after augmentation.	

Original	Altered text	Score
Edit the existing HTML code to create an or-	Modify the existing HTML code to generate	0
dered list with the following requirements:	an unordered list with the following con-	
- The ordered list should be numbered in	straints:	
descending order.	- The list items should be numbered in as-	
- Each fruit name in the list should be dis-	cending order.	
played in uppercase letters.	- Each fruit name in the list should be dis-	
- Add a CSS style to the list items to have a	played in lowercase letters.	
background color of yellow .	- Include a CSS style to the list items to have	
- The fruit names should be sorted in alpha-	a background color of red .	
betical order.	- The fruit names should be sorted in reverse	
- Only include fruits that start with a vowel in	alphabetical order.	
the list.	- Only include fruits that start with a conso-	
- Each fruit name should be followed by a	nant in the list.	
comma.	- Each fruit name should be followed by a	
- Add a CSS style to the list to have a border	semicolon.	
of 1px solid black .	- Add a CSS style to the list to have a border	
- The ordered list should be displayed hori-	of 2px dotted blue .	
zontally instead of vertically.	- The unordered list should be displayed	
- Add a CSS style to the list items to have a	vertically instead of horizontally.	
font size of 20px.	- Add a CSS style to the list items to have a	
- The list should have a maximum of 5	font size of 18px .	
items.	- The list should have a minimum of 5 items.	
Create an object with key-value pairs using	Generate an object with key-value pairs in	4
JavaScript.	JavaScript using a for loop.	
Write a function to determine the largest	Identify the largest prime number within an	5
prime number in an unordered array.	unstructured array.	
array = [6, 8, 1, 3, 4]	array = [6, 8, 1, 3, 4]	
Write a code to generate a random pass-	Compose a random password using Python's	5
word using the random library of Python.	random library, ensuring the following crite-	
The password should satisfy the following	ria:	
requirements:	- The password length must be 8 characters.	
- The length of the password should be 8	- The password should include at least one	
characters.	uppercase letter, one lowercase letter, one	
- The password should contain at least one	digit, and one special character (!#\$% ^&*).	
uppercase letter, one lowercase letter, one	- The password should not commence with a	
digit, and one special character (!#\$% ^&*).	digit.	
- The password should not start with a digit.	- The password must not contain any consec-	
- The password should not contain any con-	utive repeating characters.	
secutive repeating characters.	Without utilizing any external libraries	
Note: You are not allowed to use any ex-	or built-in functions for password genera-	
ternal libraries or built-in functions for	tion, follow the aforementioned rules and	
generating passwords.	rephrase the given instruction.	

Table 6: Generated samples on BioRED. Highlighted words denote the entities with special tags around them. Blue words stand for those modifications while keeping the same semantic meaning.

Original	Altered text
Curcumin prevents <a>@ChemicalEntitySrc\$ maleate-	Curcumin protects against renal toxicity induced
induced @/ChemicalEntitySrc\$ nephrotoxicity:	by <a>@ChemicalEntitySrc\$ maleate <a>@/ChemicalEn-
relation to hemodynamic alterations, oxidative	titySrc ^{\$} : relation to changes in hemodynamics,
stress, mitochondrial oxygen consumption and	oxidative stress, mitochondrial oxygen use and
activity of respiratory complex I. The potential pro-	activity of respiratory complex I. The possible pro-
tective effect of the dietary antioxidant curcumin	tective effect of the dietary antioxidant curcumin
(120 mg/Kg/day for 6 days) against the renal	(120 mg/Kg/day for 6 days) against kidney in-
injury induced by <a>@ChemicalEntitySrc\$ maleate	jury caused by <a>@ChemicalEntitySrc\$ maleate
@/ChemicalEntitySrc\$ was evaluated. Tubular	@/ChemicalEntitySrc\$ was looked at. Tubular
proteinuria and oxidative stress were induced	proteinuria and oxidative stress were brought on
by a single injection of <a>@ChemicalEntitySrc\$	by a single injection of <a>@ChemicalEntitySrc\$
maleate @/ChemicalEntitySrc\$ (400 mg/kg) in rats.	maleate @/ChemicalEntitySrc\$ (400 mg/kg) in
@ChemicalEntitySrc\$ Maleate-induced @/Chem-	rats. <a>@ChemicalEntitySrc \$ Maleate-induced
icalEntitySrc\$ renal injury included increase in	@/ChemicalEntitySrc\$ kidney injury included
renal vascular resistance and in the urinary excre-	increased renal vascular resistance and urinary
tion of total protein, glucose, sodium, neutrophil	excretion of total protein, glucose, sodium, neu-
gelatinase-associated lipocalin (NGAL) and N-	trophil gelatinase-associated lipocalin (NGAL) and
acetyl b-D-glucosaminidase (NAG), upregulation	N-acetyl b-D-glucosaminidase (NAG), upregula-
of kidney injury molecule (KIM)-1, decrease in	tion of kidney injury molecule (KIM)-1, decreased
renal blood flow and <a>@GeneOrGeneProductTgt\$	renal blood flow and <a>@GeneOrGeneProductTgt \$
claudin-2 @/GeneOrGeneProductTgt\$ expres-	claudin-2 @/GeneOrGeneProductTgt\$ expression
sion besides of necrosis and apoptosis of tubular	in addition to necrosis and apoptosis of tubular
cells on 24 h. Oxidative stress was determined	cells at 24 hrs. Oxidative stress was determined
by measuring the oxidation of lipids and proteins	by measuring the oxidation of lipids and proteins
and diminution in renal Nrf2 levels. Studies were	and reduction in renal Nrf2 levels. Studies were
also conducted in renal epithelial LLC-PK1 cells	also done in renal epithelial LLC-PK1 cells and in
and in mitochondria isolated from kidneys of all	mitochondria isolated from kidneys of all the ex-
the experimental groups. <a>@ChemicalEntitySrc\$	perimental groups. <a>@ChemicalEntitySrc \$ Maleate
Maleate @/ChemicalEntitySrc\$ induced cell dam-	@/ChemicalEntitySrc\$ caused cell damage and
age and reactive oxygen species (ROS) production	reactive oxygen species (ROS) production in LLC-
in LLC-PK1 cells in culture. In addition, @Chem-	PK1 cells in culture. Additionally, @ChemicalEn-
icalEntitySrc\$ maleate @/ChemicalEntitySrc\$	titySrc\$ maleate @/ChemicalEntitySrc\$ treatment
treatment reduced oxygen consumption in ADP-	lowered oxygen consumption in ADP-stimulated
stimulated mitochondria and diminished respira-	mitochondria and reduced respiratory control in-
tory control index when using malate/glutamate	dex when using malate/glutamate as substrate. The
as substrate. The activities of both complex I and	activities of both complex I and aconitase were
aconitase were also diminished. All the above-	also decreased . All the above described changes
described alterations were prevented by curcumin.	were prevented by curcumin. It is concluded that
It is concluded that curcumin is able to attenuate	curcumin is able to reduce in vivo @ChemicalEn-
in vivo <a>@ChemicalEntitySrc\$ maleate-induced	titySrc\$ maleate-induced @/ChemicalEntitySrc\$
@/ChemicalEntitySrc\$ nephropathy and in vitro	nephropathy and in vitro cell damage. The in vivo
cell damage. The in vivo protection was associated	protection was linked to prevention of oxidative
to the prevention of oxidative stress and preserva-	stress and preservation of mitochondrial oxygen
tion of mitochondrial oxygen consumption and	use and activity of respiratory complex I, and the
activity of respiratory complex I, and the in vitro	in vitro protection was linked to the prevention of
protection was associated to the prevention of ROS	ROS production.
production.	