

Human perceiving behavior modeling in evaluation of code generation models

Sergey Kovalchuk, Vadim Lomshakov, Artem Aliev

Huawei

{sergey.kovalchuk,vadim.lomshakov,artem.aliev}@huawei.com

Abstract

In this study, we evaluated a series of code generation models based on CodeGen and GPTNeo to compare the metric-based performance and human evaluation. For a deeper analysis of human perceiving within the evaluation procedure, we implemented a 5-level Likert scale assessment of the model output using a perceiving model based on the Theory of Planned Behavior (TPB). Through this analysis, we demonstrated an extension of model assessment as well as a deeper understanding of the quality and applicability of generated code for answering practical questions. The approach was evaluated with several model settings in order to assess diversity in the quality and style of answer. With the TPB-based model, we showed a different level of perceiving of the model result, namely, personal understanding, agreement level, and readiness to use the particular code. With this analysis, we investigate a series of issues in code generation, namely, natural language generation (NLG) problems observed in the context of programming and question-answering with code.

1 Introduction

Recent advances in natural language generation (NLG) support rapid growth in potential areas of application. One of the significant successes of NLG is the possible generation of code in various settings (Lu et al., 2021; Zhong et al., 2022): the translation of explicit specification into code, fixing errors, suggesting short snippets, etc. The common practice in NLG problems is the usage of known metrics (BertScore, BLEU, etc.) to evaluate models. Moreover, specific metrics dedicated to

code generation evaluation were developed, such as CodeBLEU (Ren et al., 2020), RUBY (Tran et al., 2019), and others. Still, recent studies (Evtikhiev et al., 2022) show that the direct application of metrics often leads to issues in code generation evaluation.

With this in mind, investigated the applicability of human evaluation widely spread in NLG problems (De Mattei et al., 2021; Hämäläinen & Alnajjar, 2021) to assess an alternative approach to code evaluation and a deeper understanding of human perceiving of code generation (and NLG output in general). The study poses two research questions. First, how is human perceiving reflected by the text- or code-oriented NLP metrics? Second, what is the structure of human perceiving in the human evaluation procedure in the question-answering scenario? Here we consider perceiving as an act of becoming subjectively aware and conscious of the observed information.

The structure of the paper is as follows. The next section describes the datasets used in the study. The following section presets the details of code generation models' selection and preparation. Section 4 describes human evaluation solutions and procedures. Section 5 discusses the study results and evaluation results. Finally, Sections 6 and 7 provide a discussion and concluding remarks respectively.

2 Dataset

Within the study, we focused on question answering (QA) with a generation of short snippets as answers to real-world problems such as questions asked in Stack Overflow¹ (SO). For consistency, we added the following restrictions to the questions and answers considered within the study, taking SO as a reference for the analysis.

¹ <https://stackoverflow.com/>

77 • We considered “conceptual” and “API usage” questions according to the taxonomy presented in (Beyer et al., 2020).

80 • We selected the questions that mainly contain a short textual description without explicit code presented.

83 • Contrarily, we use answers with explicit code snippets giving the solution to the proposed problem.

86 • To further specify the scope of the study, we only considered one programming language, Python, as it is one of the most popular ones.

89 To prepare an appropriate dataset we followed two steps. First, we used the publicly available dataset CoNaLa² (Yin et al., 2018) with explicitly identified train (2379 entries) and test parts (500 entries). The dataset originates from SO and contains explicit short questions and reference code snippets.

96 Alternatively, we prepared our own dataset from the original data on SO questions available on Stack Exchange³. To follow our requirements, we selected questions with the tag “python”. For reference, we selected answers that earned maximum scores according to the SO data. To filter questions on presence or absence of code, we search the text for `<pre><code>` in HTML data. After that, we selected questions with no explicit code paired with answers with a single code snippet. Finally, we used regular expressions following (Beyer et al., 2020) to select “conceptual” and “API usage” classes of questions. Furthermore, we performed cleaning of the code (e.g. removing decorations inserted by software (“>>>”, “In [1]:”, etc.), comments, and checked the parsing status using Tree-sitter⁴. After these steps, we obtained a dataset containing 42292 entries (pairs of questions and answers). Out of them, we selected 1000 entries as a test dataset. The test dataset was built using questions from 2021 and beyond to lower possible data leaking as we are using models trained on publicly available data.

119 3 Code generation models

120 3.1 Model selection

121 To select models for our study, we evaluated those which are publicly available, computationally

123 inexpensive, and applicable on our data with finetuning. First, we selected GPT-Neo(-J)⁵ (Black, Sid et al., 2021, p.), which shows high performance compared to Codex (Xu et al., 2022). Second, we chose CodeGen-mono-2B by Salesforce (Nijkamp et al., 2022), which was trained not only on the Pile dataset but also separately on the code from BigQuery and BigPython. CodeGen-mono-2B shows good results on HumanEval, which was similar to the Codex model of the same size (Chen et al., 2021). Additionally, we picked CoPilot by Microsoft as an industrial SOTA reference solution, which is also based on Codex (Chen et al., 2021).

137 3.2 Finetuning

138 Finetuning was performed for the selected models on training datasets from both CoNaLa (further denoted as FT:C) and SO (further denoted as FT:SO). We used Transformers and DeepSpeed libraries with common hyperparameters (optimizer = AdamW, Adam betas = (0.9, 0.999), Adam epsilon = 1e-08, weight decay = 0.0, learning rate = 5e-06, learning rate decay = linear, batch size(#samples) = 40, fp16). Moreover, we performed experiments with various code wrapping for prompt preparation. A query wrapped as a multiline comment to the generated code was selected as a well-performing baseline. Additionally, we performed a series of experiments in prompt engineering and selected the best performing solution for further experiments (denoted as FT:C+).

155 4 Human evaluation

156 4.1 User interface implementation

157 For performing the human evaluation, a user interface (UI) was developed as a web application using the Dash⁶ framework (see Figure 1). The UI enables the collection of feedback information in two ways.

162 First (*HFI* - human feedback 1), the UI shows a pair of answers generated for the same question by different models. The pairwise comparison of two answers was collected from the user by asking them to select the best answer with a 3 or 5 (depending on configuration) levels Likert scale

² <https://conala-corpus.github.io/>

³ <https://stackexchange.com/>

⁴ <https://tree-sitter.github.io/>

⁵ <https://github.com/EleutherAI/gpt-neo>

⁶ <https://dash.plotly.com/>

168 from -2 (the left answer is the best) to +2 (the right
 169 answer is the best). This feedback is collected for
 170 further research purposes to improve model
 171 performance with human-centered prediction
 172 models (currently considered as future research
 173 plans following the works (Nakano et al., 2022;
 174 Stiennon et al., 2020)).

175 Second (**HF2**), the user was asked to assess both
 176 answers (code snippets) with three scores using a
 177 5-level Likert scale (from -2 to +2) by estimating:

- 178 • The general consistency of the code (whether
 179 the code is readable/understandable). The
 180 scale is considered to reflect how well the user
 181 *understands* the answer.
- 182 • The correctness of the answer with respect to
 183 the proposed question. The scale is considered
 184 to reflect the user’s *agreement* with the
 185 answer.
- 186 • The usability of the provided answer. The
 187 scale is considered to reflect the user’s
 188 expected *intention to use*.

189 These scales analyze human behavior aspects by
 190 assessing the information perceiving in alignment
 191 with a model based on the theory of planned
 192 behavior (TPB) (Ajzen, 1991), widely used to
 193 quantify human behavior as reflected by attitude,
 194 subjective norms, and perceived control affecting
 195 target intention to use a considered technology. In
 196 our case, we consider “agreement” as the first
 197 criterion, reflecting the general user’s attitude,

198 “understand” as the second one, reflecting
 199 correspondence to subjective norms and perceived
 200 control, and “use” as a final target criterion, the
 201 obtained intention to use the solution.

202 The human feedback is collected for each pair of
 203 answers storing the scores provided by the user and
 204 his/her provided name. After each evaluation round
 205 (ended by clicking the “Submit” button), the
 206 interface is updated with a new pair of answers for
 207 further analysis.

208 For evaluation purposes, the original and
 209 finetuned models were applied to test sets in the
 210 selected datasets (CoNaLa and SO) forming a
 211 collection of alternative answers to 1500 questions.
 212 Applying the selected models and filtering empty
 213 answers, we obtained 10013 answers of different
 214 origin and quality. On each round, a random sample
 215 of two different answers was presented to the user.
 216 With this approach, a subset of 1364 questions was
 217 selected, supported with two or more answers by
 218 different models.

219 4.2 Human perceiving assessment

220 Within the presented study, we focused on the
 221 internal structure of answer perceiving during
 222 human evaluation. Thus, we performed a deeper
 223 analysis of HF2 to understand the connections
 224 between different features. For this purpose, we
 225 consider three main groups of features.

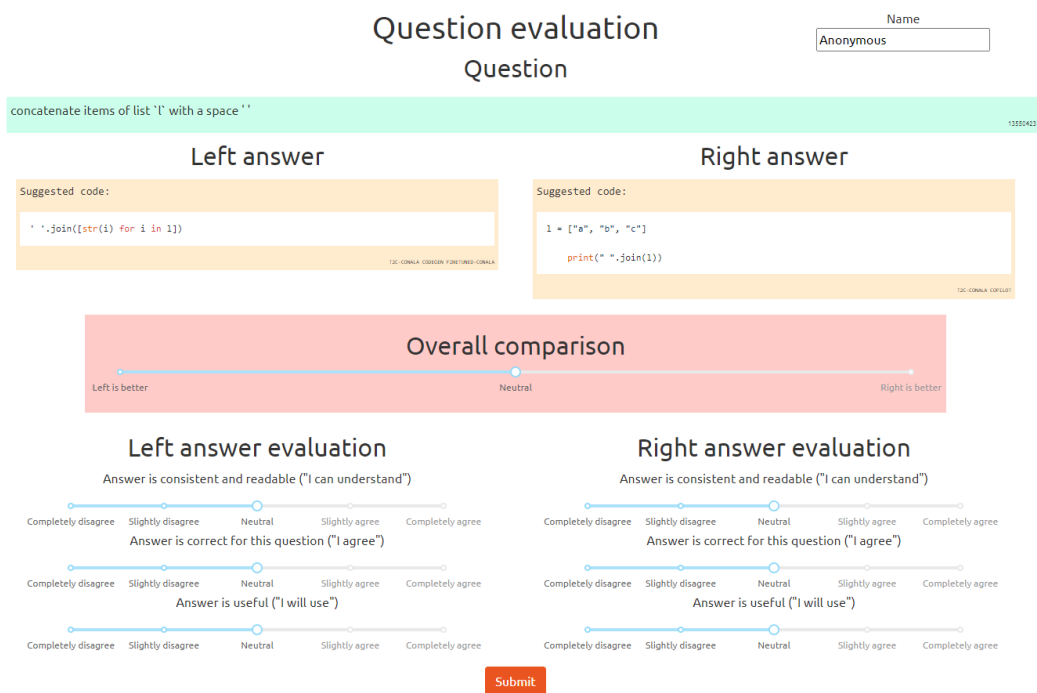


Figure 1. User interface for human evaluation

Model	BertScore		Rouge		CodeBLEU		Ruby		SacreBLEU	
	mean	std	mean	std	mean	std	mean	std	mean	std
Dataset: CoNaLa										
CodeGen	0.8068	0.1827	0.3142	0.2638	0.2821	0.2379	0.2791	0.2363	0.1196	0.1493
CodeGen <i>FT:C</i>	0.9017	0.1132	0.5532	0.2976	0.4848	0.2861	0.5392	0.3151	0.2142	0.1759
CodeGen <i>FT:SO</i>	0.8326	0.0379	0.1707	0.1316	0.0918	0.1095	0.1014	0.1348	0.0522	0.0658
CodeGen <i>FT:C+</i>	0.9235	0.0511	0.6070	0.2763	0.5802	0.2519	0.6246	0.2845	0.2571	0.1952
GPT-Neo	0.7370	0.1688	0.0503	0.0688	0.0785	0.0670	0.1460	0.1190	0.0111	0.0151
GPT-Neo <i>FT:C</i>	0.8366	0.1518	0.2926	0.2648	0.2298	0.2401	0.2619	0.2759	0.0900	0.1096
GPT-Neo <i>FT:SO</i>	0.8251	0.0380	0.1453	0.1122	0.0749	0.0858	0.0909	0.1192	0.0400	0.0455
CoPilot	0.8520	0.0398	0.3668	0.2075	0.2815	0.1554	0.2812	0.1899	0.1179	0.1162
Dataset: SO										
CodeGen	0.7434	0.1838	0.0790	0.0951	0.1687	0.1549	0.0974	0.1025	0.0176	0.0331
CodeGen <i>FT:C</i>	0.8178	0.0923	0.1473	0.1447	0.3311	0.2488	0.1615	0.1935	0.0396	0.0572
CodeGen <i>FT:SO</i>	0.8099	0.0825	0.1298	0.1188	0.1729	0.1773	0.0933	0.1062	0.0327	0.0452
GPT-Neo	0.7543	0.1286	0.0511	0.0577	0.1411	0.1074	0.0991	0.1101	0.0093	0.0123
GPT-Neo <i>FT:C</i>	0.7912	0.1488	0.1193	0.1265	0.2986	0.2295	0.1433	0.1606	0.0292	0.0419
GPT-Neo <i>FT:SO</i>	0.8003	0.1126	0.1156	0.1109	0.1574	0.1675	0.0953	0.1069	0.0275	0.0378

Table 1. Metric-based evaluation

²²⁶ *FG1* (feature group 1) includes the common
²²⁷ metrics used for the NLG task. The selection of
²²⁸ metrics includes general purpose metrics
²²⁹ (BertScore, Rouge, SacreBLEU) and metrics
²³⁰ specific to code generation problems (CodeBLEU,
²³¹ Ruby). The metrics were evaluated for each model
²³² applied for test datasets.

²³³ *FG2* includes votes collected from the users
²³⁴ along three selected scores, namely, subjective
²³⁵ consistency (understanding), subjective
²³⁶ correctness (agreement), and subjective intention
²³⁷ (use).

²³⁸ *FG3* includes simple test features (linguistic
²³⁹ features), namely, question and answer length,
²⁴⁰ average lines number in answer, and average lines
²⁴¹ number in question.

²⁴² To answer the proposed research questions, we
²⁴³ analyzed the interconnection between the features
²⁴⁴ in the three groups by assessing the pairwise
²⁴⁵ mutual information (MI) between features. As we
²⁴⁶ focused mainly on perceiving structure and
²⁴⁷ interconnection, the main analysis and
²⁴⁸ interpretation were applied to a) internal MI
²⁴⁹ between features in FG2; b) MI between features
²⁵⁰ in FG2 and features in other groups.

²⁵¹ 5 Results

²⁵² 5.1 Metric-based evaluation

²⁵³ We used a common train-eval-test split for
²⁵⁴ evaluation. In the case of the CoNaLa dataset, the
²⁵⁵ test dataset was pre-selected by the authors. In the
²⁵⁶ case of the SO dataset, we composed the test
²⁵⁷ dataset as random samples of 1000 answers dated
²⁵⁸ 2021 and beyond, while the answers dated 2020
²⁵⁹ and earlier were used for training. In both cases, we
²⁶⁰ split the training dataset into train and validation
²⁶¹ parts as 9:1 randomly.

²⁶² Table 1 shows the main evaluation results
²⁶³ according to the selected NLP metrics. In the case
²⁶⁴ of the CoNaLa dataset, the best results were
²⁶⁵ obtained by CodeGen FT:C+, followed by CoPilot.
²⁶⁶ In the case of the SO dataset, the best performance
²⁶⁷ was achieved with CodeGen FT:C.

268 5.2 Human evaluation

269 With the implemented UI and judgment collection
 270 procedure, the human evaluation was performed in
 271 a semi-open way by exposing the UI with a pre-
 272 defined collection of answers to independent
 273 groups of users of different backgrounds, but having
 274 a basic understanding of coding and the principles
 275 of software engineering. The user set includes
 276 MSc/PhD students and researchers in computer
 277 science and related areas, as well as professional
 278 software developers. The diversity in experience
 279 of the users enables further deeper analysis of the
 280 nature of human perceiving, along with the
 281 possible assessment of experience influence.

282 During a week-long evaluation period, the
 283 collection of votes for 614 answers in the HF2 part
 284 was collected from 43 different users. Within the
 285 analysis, we exclude the original models and only
 286 consider finetuned models. The average scores for
 287 FG2 with the selected pairs model/dataset are
 288 presented in Table 2.

289 We observe the highest perceiving in CoPilot
 290 and finetuned CodeGen. The CoNaLa dataset
 291 shows slightly lower Understand scores compared
 292 to the SO dataset. At the same time, the SO dataset
 293 shows negative Agree and Use scores reflecting
 294 wrong (but consistent) answers generated by the

Model	N	Under-stand	Agree	Use
Dataset: CoNaLa				
CodeGen <i>FT:C</i>	58	0.5345	0.3966	0.4310
CodeGen <i>FT:SO</i>	68	0.0882	-0.1471	-0.1912
CodeGen <i>FT:C+</i>	56	0.8571	0.4464	0.4286
GPT-Neo <i>FT:C</i>	62	0.0000	-0.4677	-0.5323
GPT-Neo <i>FT:SO</i>	55	-0.0364	-0.6364	-0.8364
CoPilot	39	0.8974	0.4872	0.2308
Dataset: SO				
CodeGen <i>FT:C</i>	25	0.9200	-0.3200	-0.2000
CodeGen <i>FT:SO</i>	13	0.0769	-0.6154	-0.4615
GPT-Neo	21	-0.9048	-1.8095	-1.8095
GPT-Neo <i>FT:C</i>	21	0.2381	-0.5238	-0.4286
GPT-Neo <i>FT:SO</i>	16	-0.6250	-1.1875	-1.1250

Table 2. Human perceiving evaluation

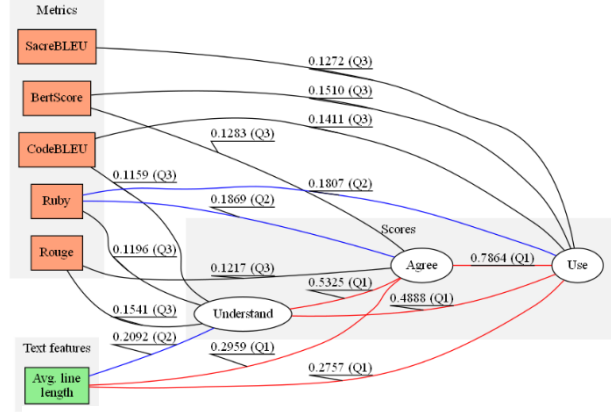


Figure 2. Mutual information and feature connection

295 model. A more interesting observation is the
 296 diversity in scores exhibited by the best CodeGen
 297 and CoPilot in the CoNaLa dataset: CoPilot shows
 298 slightly higher Understand and Agree scores, while
 299 the Use score is much higher for CodeGen.

300 For the analysis of MI, we divided the selection
 301 of the connections between features (FG2-FG2,
 302 FG2-FG1, FG2-FG3) into four groups according to
 303 the quartiles of the MI distribution (the division
 304 levels correspond to MI of 0.1144, 0.1621, 0.2683
 305 for thresholds between Q4-Q3, Q3-Q2, Q2-Q1
 306 respectively). We dropped Q4 (lowest MI)
 307 connections as insignificant and considered the
 308 others for further analysis. Figure 2 shows the
 309 selection of features connected with the labels
 310 denoting MI level and the quartile to which the
 311 value belongs. Also, the quartiles are shown in
 312 color (Q1 – red, Q2 – blue, Q3 – black).

313 As expected, the internal connections of
 314 perceiving features have high MI. The highest
 315 interconnection is observed between Agree and
 316 Use scores. Additionally, they are highly correlated
 317 ($cor = 0.8981$). A simple linear regression model
 318 was estimated as $U = 0.8389A + 0.0764C +$
 319 0.0134 ($R^2 = 0.8098$) where U is for intention to
 320 use, A is for agreement, C is for understanding
 321 (internal consistency). Thus, we expect that the
 322 intention to use is highly defined by the agreement
 323 to the answer. On the other hand, the dependency
 324 between understanding and agreement can also be
 325 observed, but is rather lower: $A = 0.6687C +$
 326 0.0376 (with $R^2 = 0.4358$).

327 We see that the connection of most of the NLP
 328 metrics is rather low (Q3), except for the Ruby
 329 metric showing a more significant (Q2) influence
 330 on the agreement and the intention to use. This can
 331 be interpreted as a good evaluation of code quality.
 332 The remaining features are mostly low and have

interconnection with different perceiving scores. In general, the MI for interconnection with the intention to use is slightly higher.

The analysis shows that basic linguistic features (FG3) mainly have low interconnection with perceiving votes (FG2) as the connections were assessed with low MI (Q4) with one exception for average line length having high (Q1 or Q2) impact on the perceiving features. Further analysis shows that there are weak results where a model failed to generate a structured answer and produce long (mainly erroneous) lines of code.

6 Discussion

The study focused on code generation problems. For this purpose, we used modern NLG models (CodeGen, GPT) and performed fine-tuning to obtain a higher quality of the results according to the common pipeline. The results we obtained with the CodeGen model look promising and produce high-quality results comparable to the industrial solutions (CoPilot). Still, we consider the further improvement of the code generation QA model as one of our future research directions.

The main goal of the study is the investigation of the human perceiving structure in the NLG problem. The area of human evaluation is widely studied in different NLG settings including general assessment of natural language generation, as well as distinguishing model-generated from human-generated answer. The common goal of most studies in the area is model improvement with collected evaluation feedback. Still, the question of how we can evaluate human feedback and which level of trust can be given to it is still open. This problem becomes more challenging if the evaluation is performed in a complex domain where a human annotator needs to be a domain expert. In that case, the human feedback collection could be more expensive and can provide more diversity in judgments. With this in mind, we considered a problem of NLG in programming QA and code generation, with programmers as experts evaluating high-level correspondence of the answer to a proposed question. Within the study, we focused on the general structure of human perceiving divided into three main parts: whether the human understands the model output; whether he/she considers it as correct; whether he/she is ready to use it in practice. This structure of human perceiving was considered previously in the expert-based evaluation of decision support

systems (Kovalchuk et al., 2022) and showed good results in understanding human intentions in perceiving AI-based prediction in such complex domains as medicine.

Within the study, we analyzed internal interconnection between human perceiving features in code generation and discovered that although the agreement and intention to use are highly correlated, the understanding (subjective correctness) and agreement show lower interconnections. This could be interpreted as a sign of existing issues in generated code properly recognized by a human, i.e., the answer generated by the model “looks like code” but doesn’t resolve the question properly. On the other hand, the high interconnection between agreement and intention to use could be treated as promising results for code generation problems: if the code answers the question, it is good enough to be used.

An interesting result obtained in the analysis of external interconnection of perceiving score is a rather weak MI of connections to the common NLP metrics. The only metric showing medium interconnection is Ruby, intentionally developed for code evaluation with semantic comparison (Tran et al., 2019). We consider this result as a sign of the rather weak applicability of common NLP metrics to complex NLG problems.

One of the questions raised during the evaluation is whether we can compare the syntactic correctness of the code to the subjective correctness (understanding) of the code. During the evaluation, we see several examples of code that was syntactically incorrect but may have been evaluated as useful (e.g. containing the correct line of code followed by an ill-formatted line). We believe that continuous Likert-scale-based evaluation may help to improve the model training/finetuning in further studies with human evaluation. Still, we consider this issue as one of the directions for the future development of the proposed approach.

It is worth mentioning that CodeGen FT:SO shows lower performance compared to CodeGen FT:C even on the SO dataset. We suppose that a possible reason for such behavior is the fact that the CoNaLa dataset was manually curated and contains only one-line code answers. At the same time, the SO dataset was not filtered in that way and contains answers of diverse length, structure, and quality. A further investigation of this issue is one of the directions for further research.

436 The presented results can be used in two ways. 487
437 First, to improve models by training with a deeper 488
438 understanding of human perceiving structure. For 489
439 example, the mentioned values could be considered 490
440 as a sequential filter where the next step can be
441 considered only by the answer passed from the
442 previous one. In this way, the perceiving may be
443 considered as a reward for the model in the
444 generation of the answer.

445 Second, the understanding of human perceiving
446 may improve human-computer interaction in AI-
447 based applications. For instance, the separate
448 prediction of human perceiving features may
449 provide important information depending on the
450 interaction scenario. In particular, the Agree score
451 may be more influential in code automatically
452 generated by explicit specification, while the Use
453 score may be more important in direct human-
454 centered QA (e.g., in a form of an intelligent IDE
455 assistant).

456 7 Conclusion and future work

457 The current study shows early results in the
458 research of human perceiving understanding within
459 the context of NLG human evaluation. Being
460 aimed at a deeper understanding of the internal
461 structure of human perceiving and interconnection
462 with common metrics, the study shows that
463 perceiving structure may be decomposed into a
464 complex value with implicit interconnection
465 directing from model-generated structure
466 evaluation to subjective intention to use. The
467 proposed structure of human perceiving may be
468 further used for collecting judgments in complex
469 domains where direct application of common NLP
470 metrics gives rather weak results and where a high
471 semantic diversity in the possible answer may be
472 observed.

473 We consider the further development of the
474 proposed study in the following directions. First,
475 we would like to continue collecting human
476 feedback in order to advance the development of
477 the perceiving model. Additionally, we would like
478 to extend the research to analyze the personal
479 characteristics of the human (in our case, it may be
480 a personal experience, relevance to the question,
481 etc.). Second, we would like to investigate the
482 possible application of the decomposed human
483 perceiving value to model training/finetuning.
484 Third, we are interested in the extension of the
485 model with technical characteristics of generated
486 code (e.g. checking for syntactic errors, test-based

evaluation). Finally, we are planning to investigate
the influence of the context on human perceiving
including the dependencies from the application
scenario.

491 Acknowledgements

492 We thank all the users involved into the evaluation
493 procedure presented in this study for their effort in
494 assessment of the generated answers. We are also
495 grateful to the anonymous GEM reviewers for their
496 valuable comments and suggestions.

497 References

- 498 Ajzen, I. (1991). The theory of planned behavior.
499 *Organizational Behavior and Human Decision*
500 *Processes*, 50(2), 179–211.
501 [https://doi.org/10.1016/0749-5978\(91\)90020-T](https://doi.org/10.1016/0749-5978(91)90020-T)
- 502 Beyer, S., Macho, C., Di Penta, M., & Pinzger, M.
503 (2020). What kind of questions do developers ask
504 on Stack Overflow? A comparison of automated
505 approaches to classify posts into question
506 categories. *Empirical Software Engineering*, 25(3),
507 2258–2301. <https://doi.org/10.1007/s10664-019-09758-x>
- 508 Black, Sid, Leo, Gao, Wang, Phil, Leahy, Connor, &
509 Biderman, Stella. (2021). *GPT-Neo: Large Scale*
510 *Autoregressive Language Modeling with Mesh-*
511 *Tensorflow* (1.0). Zenodo.
512 <https://doi.org/10.5281/ZENODO.5297715>
- 513 Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P.
514 de O., Kaplan, J., Edwards, H., Burda, Y., Joseph,
515 N., Brockman, G., Ray, A., Puri, R., Krueger, G.,
516 Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P.,
517 Chan, B., Gray, S., ... Zaremba, W. (2021).
518 *Evaluating Large Language Models Trained on*
519 *Code* (arXiv:2107.03374). arXiv.
520 <http://arxiv.org/abs/2107.03374>
- 521 De Mattei, L., Lai, H., Dell’Orletta, F., & Nissim, M.
522 (2021). Human Perception in Natural Language
523 Generation. *Proceedings of the 1st Workshop on*
524 *Natural Language Generation, Evaluation, and*
525 *Metrics (GEM 2021)*, 15–23.
526 <https://doi.org/10.18653/v1/2021.gem-1.2>
- 527 Evtikhiev, M., Bogomolov, E., Sokolov, Y., &
528 Bryksin, T. (2022). *Out of the BLEU: How should*
529 *we assess quality of the Code Generation models?*
530 (arXiv:2208.03133). arXiv.
531 <http://arxiv.org/abs/2208.03133>
- 532 Hämäläinen, M., & Alnajjar, K. (2021). Human
533 Evaluation of Creative NLG Systems: An
534 Interdisciplinary Survey on Recent Papers.
535 *Proceedings of the 1st Workshop on Natural*
536 *Language Generation, Evaluation, and Metrics*
537

- (*GEM* 2021), 84–95. <https://doi.org/10.18653/v1/2021.gem-1.9>
- Kovalchuk, S. V., Kopanitsa, G. D., Derevitskii, I. V., Matveev, G. A., & Savitskaya, D. A. (2022). Three-stage intelligent support of clinical decision making for higher trust, validity, and explainability. *Journal of Biomedical Informatics*, 127, 104013. <https://doi.org/10.1016/j.jbi.2022.104013>
- Lu, S., Guo, D., Ren, S., Huang, J., Svyatkovskiy, A., Blanco, A., Clement, C., Drain, D., Jiang, D., Tang, D., Li, G., Zhou, L., Shou, L., Zhou, L., Tufano, M., Gong, M., Zhou, M., Duan, N., Sundaresan, N., ... Liu, S. (2021). *CodeXGLUE: A Machine Learning Benchmark Dataset for Code Understanding and Generation* (arXiv:2102.04664). arXiv. <http://arxiv.org/abs/2102.04664>
- Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., Jiang, X., Cobbe, K., Eloundou, T., Krueger, G., Button, K., Knight, M., Chess, B., & Schulman, J. (2022). *WebGPT: Browser-assisted question-answering with human feedback* (arXiv:2112.09332). arXiv. <http://arxiv.org/abs/2112.09332>
- Nijkamp, E., Pang, B., Hayashi, H., Tu, L., Wang, H., Zhou, Y., Savarese, S., & Xiong, C. (2022). *A Conversational Paradigm for Program Synthesis* (arXiv:2203.13474). arXiv. <http://arxiv.org/abs/2203.13474>
- Ren, S., Guo, D., Lu, S., Zhou, L., Liu, S., Tang, D., Sundaresan, N., Zhou, M., Blanco, A., & Ma, S. (2020). *CodeBLEU: A Method for Automatic Evaluation of Code Synthesis* (arXiv:2009.10297). arXiv. <http://arxiv.org/abs/2009.10297>
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., & Christiano, P. F. (2020). Learning to summarize with human feedback. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, & H. Lin (Eds.), *Advances in Neural Information Processing Systems* (Vol. 33, pp. 3008–3021). Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf>
- Tran, N., Tran, H., Nguyen, S., Nguyen, H., & Nguyen, T. (2019). Does BLEU Score Work for Code Migration? *2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)*, 165–176. <https://doi.org/10.1109/ICPC.2019.00034>
- Xu, F. F., Alon, U., Neubig, G., & Hellendoorn, V. J. (2022). *A Systematic Evaluation of Large Language Models of Code* (arXiv:2202.13169). arXiv. <http://arxiv.org/abs/2202.13169>
- Yin, P., Deng, B., Chen, E., Vasilescu, B., & Neubig, G. (2018). Learning to mine aligned code and natural language pairs from stack overflow. *Proceedings of the 15th International Conference on Mining Software Repositories*, 476–486. <https://doi.org/10.1145/3196398.3196408>
- Zhong, M., Liu, G., Li, H., Kuang, J., Zeng, J., & Wang, M. (2022). *CodeGen-Test: An Automatic Code Generation Model Integrating Program Test Information* (arXiv:2202.07612). arXiv. <http://arxiv.org/abs/2202.07612>