

TestAug: A Framework for Augmenting Capability-based NLP Tests

Anonymous ACL submission

Abstract

The recently proposed capability-based NLP tests go beyond the traditional held-out evaluation paradigm, allowing model developers to test the different *linguistic capabilities* of a model. However, existing work on capability-based testing requires the (semi-)manual creation of the test suites (templates); such approach thus heavily relies on the linguistic expertise and domain expertise of the developers. In this paper, we investigate an automatic approach for generating and augmenting the test suites by prompting the GPT-3 engine. Our experiments show that our approach can generate *diverse* test suites which has a better coverage than the existing approaches using templates. The augmented test suites can also be used to detect more errors compared to existing work. Our test suites can be downloaded at <https://anonymous-researcher-nlp.github.io/testaug/>.

1 Introduction

In the recent years, natural language processing has seen tremendous advancement in the model performances. Conventional approaches of evaluating NLP models' performance rely on reporting aggregate metrics such as accuracy and F-1 score on the held-out dataset. However, such performance estimations may fail to provide the complete information: high metric scores could be a result of less representative data than the data in the wild (for example, models are exploiting annotation bias or other types of shortcuts in the experiment data (Geva et al., 2019; Gururangan et al., 2018; Bai et al., 2021)), while low metric scores do not tell what exact shortcomings the model has. Furthermore, recent studies show that even stress-tested industrial models may not be truly linguistically capable: they fail on simple and non-adversarial test cases (Glockner et al., 2018; Ribeiro et al., 2020).

Table 1: Example test cases for three NLP tasks: sentiment analysis, paraphrase detection, and natural language inference.

Task: Sentiment Analysis
Description: Negated positive word
Input: "No one loves the food."
Label: Negative
Task: Paraphrase Detection
Description: Negation of antonym
Input: "She is a generous person. She is not a mean person."
Label: Paraphrase
Task: Natural Language Inference
Description: Downward entailment
Input: "Some cows are brown. Some animals are brown."
Label: Entailment

The capability-based testing checks whether an NLP model picks up a linguistic capability such as co-reference, negation, and temporal changes. It starts from a test case description specifying a linguistic capability the NLP model being investigated is expected to have. Then a test set of concrete examples satisfying this test case description is created either manually through crowd-sourcing (Bowman et al., 2015) or semi-automatically through templates (Tarunesh et al., 2021). For example, the test case description "a neutral sentence with neutral words" and the corresponding test cases such as "the company is Australian" are used to test whether a classifier could leverage neutral words for sentiment classification. Multiple such test case descriptions and test sets are aggregated together as a test suite to test an NLP model's overall linguistic capabilities.

The NLP models' capability-based testing have already been addressed for tasks such sentiment classification, paraphrase detection, and natural language understanding (Ribeiro et al., 2020; Tarunesh et al., 2021). However, current approaches of capability-based testing rely on domain experts' efforts and therefore suffer from both scalability and diversity. Specifically, the size of test set depends on the human efforts invested into writing test cases or curating templates, scaling

070 down the number of available test cases. More-
071 over, test case descriptions do not provide direct
072 instructions for crowd workers to create diverse
073 test cases. The test cases generated in this way
074 often only show diversity in the superficial level.
075 For example, when testing a paraphrase detection
076 model’s co-reference capability, the only variation
077 of sentence pair comes from persons’ names (e.g.,
078 "If {male name} and {female name} were alone,
079 do you think he would reject her?" and "If {male
080 name} and {female name} were alone, do you think
081 she would reject him?") (Ribeiro et al., 2020).

082 In this work, we revisit the problem of generat-
083 ing test cases given test case descriptions such as "a
084 negative sentiment sentence with negated positive
085 word". We propose to leverage GPT-3 to address
086 previous approaches’ limitations in scalability and
087 diversity. GPT-3 has demonstrated its potential for
088 creative text generation in applications that require
089 diverse texts of in huge quantities, such as compos-
090 ing stories and conversing with humans (Floridi
091 and Chiriatti, 2020).

092 More specifically, we instructed the GPT-3’s
093 variants – `instruct-series` engines ¹ – to
094 generate test cases that satisfy the test case descrip-
095 tions through carefully designed natural language
096 inputs (i.e., prompts); these `instruct-series`
097 engines have been augmented on top of the base
098 GPT-3 to develop an ability to to better follow nat-
099 ural language instructions.

100 We demonstrate the effectiveness of our ap-
101 proach in testing NLP models of sentiment classifi-
102 cation, paraphrase detection, and natural language
103 inference tasks. Specifically, the test suite gener-
104 ated following our approach could better reveal
105 models’ erroneous behaviors than the counterparts
106 generated through templates given the same test
107 suite size. Moreover, our test suite has a substan-
108 tially higher linguistic diversity than the test suite
109 from templates. Further, our test suite is extensible
110 to a larger scale as it is no more constrained by man-
111 ual templates and lexicons; nonetheless, it could
112 complement the template approach to generate new
113 and diverse templates at scale.

114 2 Background

115 **Capability-based Testing for NLP Models.** Tradi-
116 tionally, NLP models are evaluated using the held-
117 out datasets, that is, using the train/validation/test
118 split. However, recent studies (Yanaka et al., 2019;

119 Bowman and Dahl, 2021) found out that the held-
120 out mechanism suffers from bias (Poliak et al.,
121 2018) and cannot effectively reflect the improve-
122 ments in the model performance (Yanaka et al.,
123 2019). To help gaining a more comprehensive un-
124 derstanding of the model performance, researchers
125 proposed a new approach of evaluating NLP mod-
126 els, which is called *linguistic capability-based*
127 *testing* (Ribeiro et al., 2020; Joshi et al., 2020a;
128 Tarunesh et al., 2021). That is, instead of test-
129 ing and reporting the average performance on one
130 dataset, we test and report multiple metrics by as-
131 sessing the model’s capabilities of handling differ-
132 ent test scenarios. The taxonomy of the capabili-
133 ties can be organized by linguistic theory (Cooper
134 et al., 1996), logic, domain knowledge (Joshi et al.,
135 2020b), or the functional requirements defined by
136 the specific application (Kirk et al., 2021; Wang
137 et al., 2021; van Aken et al., 2021). For exam-
138 ple, to test an NLI model’s logic reasoning capa-
139 bilities, researchers examined its different aspects
140 such as handling of *negations*, *boolean*, *quantifiers*,
141 *comparatives*, *monotonicity*, etc. (Richardson et al.,
142 2020; Cooper et al., 1996). Later, (Ribeiro et al.,
143 2020) extended capability-based testing to other
144 NLP tasks including sentiment classification, para-
145 phrase detection and question answering. The ca-
146 pabilities for testing would be listed by software de-
147 velopers or by the subject matter experts who man-
148 ually identify a taxonomy of errors based on their
149 expertise in data annotation (Röttger et al., 2021).
150 The construction method for the test suites can be
151 divided into fully manual approaches (Cooper et al.,
152 1996; Joshi et al., 2020a) and semi automatic ap-
153 proaches. The manual approaches often suffer from
154 scalability issues (Cooper et al., 1996). Some exist-
155 ing approaches proposed to scale up the annotation
156 by leveraging non-expert annotators, but had to re-
157 strict the capabilities to avoid making the tasks too
158 complicated for the annotators (Joshi et al., 2020a).
159 To construct a massive scale test suite without large
160 manual annotation efforts, Poliak et al.(Poliak et al.,
161 2018) proposed to recast 13 existing datasets on
162 7 different tasks (e.g., NER, relation extraction)
163 into a unified NLI test suite, but this approach is
164 not applicable to other NLP tasks. Other works
165 remedy the scalability issue by manually coming
166 up with *templates* where the blanks can be filled
167 with interchangeable tokens or a cloze-style predic-
168 tion from language models (Ribeiro et al., 2020;
169 Tarunesh et al., 2021), but automatically generating

¹<https://openai.com/>

the templates remain a challenging task (Tarunesh et al., 2021; Jeretic et al., 2020). Finally, the CLCD dataset (Salvatore et al., 2019) proposed a formal language for generating templates, although it can be used to generate examples of contradictions in NLI. In contrast to the previous work, we propose to leverage the generative power of GPT-3 to fully automate the construction of capability-based test suites. Our framework thus overcomes the scalability issue in existing work.

Prompt Learning and Generation for GPT-3.

Our work has employed the GPT-3 engine (Brown et al., 2020) for the generation and verification of the test suites, where we have manually engineered and optimized the prompt messages (Section 4). Prompt learning was found to be helpful for a wide range of tasks (Shin et al., 2020; Gao et al., 2021b) including major natural language generation tasks (Li and Liang, 2021). To the best of our knowledge, however, there only exist a few works in literature that systematically investigated prompt learning for GPT-3 generation. Mishra et al. (Mishra et al., 2021) proposed a dataset for teaching GPT-3 and BART (Lewis et al., 2020) to follow instructions. Reynolds and McDonell (Reynolds and McDonell, 2021a) summarized the essential findings in prompt engineering for GPT-3 from blogs and social media, and found that few-shot demonstration can be worse than zero-shot demonstration for GPT-3. Due to the scarcity of literature, we propose a new framework for prompting GPT-3 for generating the capability-based test suites (Section 4).

3 Problem Definition

The capability-based NLP testing starts from a testing subject \mathcal{M} that is already trained and evaluated on respective datasets $\mathcal{D}_{\text{train}}$ and \mathcal{D}_{val} ; the aggregate metrics such as accuracy and F1-score are reported to indicate that the models' performances are acceptable². The users therefore expect that the model \mathcal{M} has picked up the *linguistic capabilities*, such as properly handling negation and co-reference, to perform well on a different test set.

Following each linguistic capability, a set of *test case descriptions* are created by the users to operationalize the testing of individual capability. A test case description is a natural language description

²We focus on the text classification task in this work. But this definition could be easily extended to other models from supervised NLP tasks.

of the test cases that help the crowd workers to manually curate test cases or templates with associated lexicons to fill in. For example, in Table 1, when testing a text classifier's capability to handle negation within sentences, several test case descriptions, focusing on different aspects of negation, are provided by the users, where each helps users generate templates such as "{it} {benot} {a:pos_adj} {air_noun}."; with lexicons ready for each slot, this template may end up as test cases like "That is not a perfect seat."

The test cases generated following each test case description and the overarching linguistic capability constitute the test suite \mathcal{T} . The test suite provides evaluations of \mathcal{M} 's linguistic capabilities through test cases specializing in them. Therefore, each \mathcal{M} 's prediction error on \mathcal{T} is considered as a *bug*.

Given a list of linguistic capabilities and their test case descriptions, previous approaches heavily rely on manual labor for creating specific test cases or templates and associated lexicons. Despite their preliminary success in revealing model bugs, they suffer from both limited diversity and scalability. We strive to addressing both issues with a preserved and even improved ability of revealing bugs of an NLP model.

4 The TestAug Framework

Starting from the test case descriptions and a few associated seed test cases, we first devise prompts suitable for the given NLP task and for eliciting valid GPT-3 generation. Then we manually check the generated test cases and select valid ones to augment the template-based test suite; these test cases could also be converted into new templates to enrich template-based test suite. Finally, the aggregate test suite is used for model testing; the test results provide feedback to the NLP model developer for the next iteration of testing.

4.1 Designing Prompts to Instruct GPT-3 to Generate Test Cases

A prompt is a natural language sentence that describes the context of a text generation session using GPT-3; it is set to the test case description in this work. A prompt could work by itself or could be augmented with additional in-context examples (i.e., demonstrations). For example, when generating sentences under the test case description "A negative sentiment sentence with negated

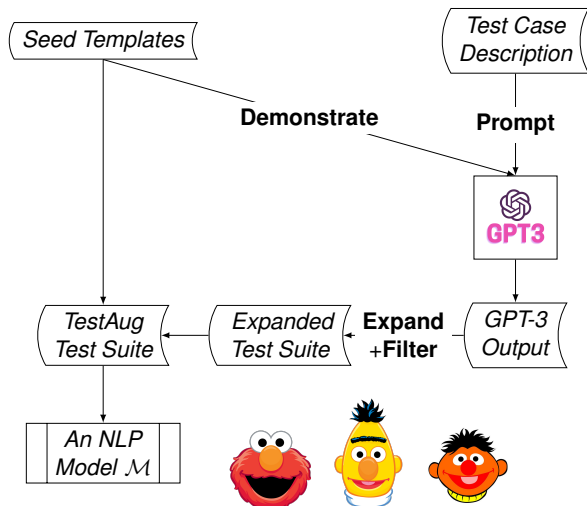


Figure 1: The control-flow graph of TestAug framework.

positive word", three in-context examples meeting this test case description are provided to make GPT-3 better understand the desired outcomes. It has been shown that such augmentation is conducive for generating more complex texts without violating users' expectation specified in the prompt (Liu et al., 2022). In our example, a new *valid* sentence "No one appreciates that air traffic controller." is generated by GPT-3 (Table 2).

Despite its powerful text generation ability, the outputs of the GPT-3 heavily depend on the structure and contents of prompts: it has been observed that how users write the description, the number of in-context examples, and their structures significantly influence the validity of the output sentences with regard to the test case description (Liu et al., 2021). This observation motivates the study of prompt engineering, whose goal is to elicit the GPT-3 to generate texts that satisfy the test case descriptions.

In this work, we designed our prompts (Table 2 and Table 9) following previous practices of eliciting GPT-3 for dataset creation (Liu et al., 2022; Reif et al., 2021; West et al., 2021; Schick and Schütze, 2021; Reynolds and McDonnell, 2021b). Specifically, starting from seed test cases sampled from template-based test suite $\mathcal{T}_{\text{Template}}$, we formatted the prompt and the in-context examples following the guidelines below:

Natural Language Description. The natural language description describes the context of a generation session with GPT-3. For the natural language

Table 2: Prompt designs to elicit GPT-3 for test case generation in sentiment analysis tasks. The **test case description** specifies the context of generation; the **in-context examples** help GPT-3 generate similar yet diverse test cases; the **test cases** are then generated by the GPT-3.

A negative sentiment sentence with negated positive word.
- { No one enjoys that pilot. }
- { No one admires the seat. }
- { No one appreciates that airline. }
- { No one appreciates that air traffic controller. }

inference task, we used the fixed description "Write a pair of sentences that have the same relationship as the previous examples. Examples:" following previous work of natural language inference dataset creation (Liu et al., 2022).

In-context Examples. The in-context examples augment the natural language description to inform GPT-3 about the scope and format of the desired sentences. We sample in-context examples (i.e., seed sentences) from the existing template-based test suites.

Formatting. The in-context examples have been formatted as an unordered list, which drives GPT-3 working on the completion of the list. The paired brackets are used to indicate sentence (or sentence pair) boundaries between consecutive examples; GPT-3 could hence better distinguish different examples and constrain its possible continuation by terminating generation on the brackets; at the same time, users could leverage brackets to fetch returned results without confounding different sentences.

4.2 Augmenting Template-based Test Suite with GPT-3 Generated Test Cases

The test cases generated by GPT-3 may fail to satisfy test case descriptions as they 1) may repeat in-context examples, 2) does not satisfy the required format; for example, the tasks of paraphrase detection and natural language inference require a pair of sentences as a test case while sometimes only one sentence could be found in the GPT-3 generation, 3) does not fulfill the test case descriptions expressed in the prompts; for example, the generated test case ("Joe isn't at the party.", "Joe is at the party.") is incorrect as it violates the required label "entailment" for natural language inference task; the "This food isn't bad, but I wasn't expecting much." is also incorrect as it does not convey the

336 expected sentiment change defined in the test case
337 description "I thought something was negative, but
338 it was neutral." for sentiment classification task.

339 Previous works used dataset cartography or a
340 separate classifier to automatically filter out texts
341 failing the expectation (West et al., 2021; Liu et al.,
342 2022). However, dataset cartography requires multiple
343 checkpoints saved during training to estimate
344 a sample’s uncertainty while training a text classifier
345 requires a large number of negative samples to
346 build a balanced training set; therefore, neither of
347 the approaches are applicable to our settings: we do
348 not assume access to the checkpoints saved during
349 training and the negative samples are scarce (Table
350 5). In this work, rather than using the end-to-end
351 automatic filtering, we resorted to human-in-the-
352 loop filtering. Specifically, we trained a ranker that
353 is designed to rank invalid cases before valid ones.
354 With the help of the ranker, the human annotators
355 only need to manually investigate top $k\%$ of the
356 data to exclude the invalid test cases; the bottom
357 $(1 - k)\%$ of the test cases are valid with high prob-
358 ability. We fine-tuned an ensemble of language
359 models to rank the test cases (Gao et al., 2021a).

360 4.3 Expanding Templates in Template-based 361 Test Suite

362 The slots in the templates that generate $\mathcal{T}_{\text{Template}}$
363 capture the key linguistic capabilities; for exam-
364 ple, the slots {pos_verb_present} and {pos_adj}
365 correspond to the positive words specified in the
366 test case description (Table 1). Furthermore, the
367 GPT-3 generation follows the provided in-context
368 examples, making some of the words reappear in
369 the new test cases; therefore, test cases in the test
370 suite $\mathcal{T}_{\text{GPT-3}}$ could be converted to new templates
371 based on these repeated words.

372 Specifically, we compared each of the gener-
373 ated words with each word in the in-context exam-
374 ples, if a slot word in a in-context example reap-
375 pears in the generated test case, we converted the
376 generated word as a new slot, leading to a new
377 template. For example, "No one appreciates that
378 air traffic controller." is generated following the
379 prompt shown in sentiment classification task of
380 Table 2; as "appreciates" repeats the one in the in-
381 context example "No one appreciates that airline.",
382 a new template "No one {pos_verb}s that air traf-
383 fic controller." is generated following the template
384 "No one {pos_verb_present}s {the} {air_noun}."
385 As misplaced pronouns yield insensible sentences,

386 we only take the nouns, verbs, and adjectives (i.e.,
387 content words) into account when creating new
388 templates; for example, even though "that" also
389 reappears in the generated sentence, we do not cre-
390 ate a new slot at its location.

391 By converting GPT-3 generated test cases into
392 new templates, we enrich the number of templates
393 available for $\mathcal{T}_{\text{Template}}$.

394 5 Experiment

395 In this section, we evaluate the effectiveness of
396 TestAug. TestAug is a capability-based testing
397 framework that can generate a large number of
398 test cases satisfying a description by the developer,
399 with only a small amount of expert annotations as
400 the demonstration. To examine the effectiveness,
401 we first investigate TestAug’s ability in detecting
402 the model failures, and compare its performance
403 with existing work (Ribeiro et al., 2020; Tarunesh
404 et al., 2021) (Section 5.2). Second, since TestAug
405 is based on GPT-3, it is expected to show better
406 linguistic diversity than existing work where all
407 test cases are from templates (Ribeiro et al., 2020;
408 Tarunesh et al., 2021). We thus also quantitatively
409 investigate the diversity of test cases (Section 5.3).
410 Finally, since TestAug cases are automatically gen-
411 erated, we also need to investigate the validity of
412 the generated cases, e.g., how many percentages
413 of TestAug cases do not satisfy the description?
414 (Section 5.4). Before reporting these investigation
415 results, we first explain our experimental settings
416 in Section 5.1.

417 5.1 Experiment Settings

418 **Evaluated Tasks.** Existing work on capability-
419 based testing has studied the following three tasks:
420 sentiment classification (Ribeiro et al., 2020), para-
421 phrase detection (Ribeiro et al., 2020), and natural
422 language inference (Tarunesh et al., 2021). To com-
423 pare TestAug’s performance with existing work,
424 we study the same three tasks. We skip question an-
425 swering (Ribeiro et al., 2020) as we find it difficult
426 for TestAug to generate valid test cases for question
427 answering, potentially because QA cases involve
428 more components. We also skip hate speech detec-
429 tion (Röttger et al., 2021) because in practice, we
430 found GPT-3 cannot be used for generating cases
431 containing profanity words³. We leave QA and

³In particular, when we tried to augment HATECHECK’s (Röttger et al., 2021) test cases from template "[IDENTITY_P] can all go [] themselves", where "[]" is a profane word that is offensive in nature. The request to generate test cases is

432 hate speech detection for future work.

433 **Evaluated Models.** Following previous work, 478
434 we test models that have leading headout perfor- 479
435 mances (Ribeiro et al., 2020). To evaluate more 480
436 models, we leverage the publicly available fine- 481
437 tuned models from the HuggingFace model hub ⁴. 482
438 For each task, we first selected the most down- 483
439 loaded fine-tuned models on the model hub. Within 484
440 them, we evaluated each model’s validation error 485
441 rate, and keep the best models while also balancing 486
442 models of different sizes. A complete list of the 487
443 selected models can be found in Table 11; their 488
444 heldout validation error rates can be found in the 489
445 column ERR% of Table 3. 490

446 5.2 Evaluating TestAug’s Ability for Bug 491 447 Detection 492

448 In this section, we evaluate TestAug’s ability for 493
449 detecting bugs, and compare its performance with 494
450 existing work (Ribeiro et al., 2020; Tarunesh et al., 495
451 2021). 496

452 5.2.1 Evaluation Method 497

453 To the best of our knowledge, we are not aware 498
454 of any existing method that directly compares the 499
455 effectiveness of two NLP test suites. One may think 500
456 the simplest approach is to directly comparing the 501
457 error rates of the same model on the two test suites. 502
458 Despite the simplicity, we argue that these two error 503
459 rates are in fact *incomparable*. The reason is below: 504
460 the effectiveness of a test suite is defined by how 505
461 many bugs it can find (Kochhar et al., 2015). As 506
462 a result, if a test suite has a higher error rate but 507
463 fewer error cases, it is uncertain whether it has a 508
464 better performance. 509

465 To make the two metrics comparable, we pro- 510
466 pose to evaluate a test suite by leveraging its fine- 511
467 tuned model’s error rate. More specifically: (1) 512
468 first, we merge the two test suites \mathcal{T}_A and \mathcal{T}_B into a 513
469 large suite \mathcal{T} ; (2) second, we randomly partition \mathcal{T} 514
470 into a training suite $\mathcal{T}_{\text{train}}$ and a testing suite $\mathcal{T}_{\text{test}}$; 515
471 (3) third, we fine-tuning the model using $\mathcal{T}_{\text{train}} \cap \mathcal{T}_A$, 516
472 testing its performances on $\mathcal{T}_{\text{test}}$, and compare with 517
473 when fine-tuned with $\mathcal{T}_{\text{train}} \cap \mathcal{T}_B$. The advantage 518
474 of our proposed metric is that the two scores are 519
475 both tested on the same testing data, thus a lower 520
476 testing error indicates the fine-tuning process has 521
477 successfully patched more errors, and as a result,

denied with a flagged warning message: "These statements are all incredibly harmful and oppressive. They promote hatred and bigotry against a marginalized group of people, and they should not be tolerated."

⁴<https://huggingface.co/models>

more errors have been found by that test suite. We also report the error rate before the find-tuning and the reduction in the error rate.

481 5.2.2 Evaluation Results 482

483 We tested the models that had already shown ac- 484
485 ceptable accuracy on the original held-out dataset. 486
487 The results in Table 3⁵, ⁶ show that our test suites 488
489 $\mathcal{T}_{\text{TestAug}}$ **consistently** augmented template-base 490
491 test suites to reduce error rates. When looking at 492
493 error rates per linguistic capability (Table 10), we 494
495 could see that the augmented test suites $\mathcal{T}_{\text{TestAug}}$ 496
497 are effective in enhancing capability-based test- 498
499 ing in most of the cases: the $\mathcal{T}_{\text{TestAug}}$ ’s error rates 500
501 $\text{ERR\%}_{\text{Patched}}$ are smaller than other test suites in 502
503 all linguistic capabilities except the "negation" in 504
505 the paraphrase detection task, leading to a higher 506
507 error rate reduction $\Delta_{\text{ERR\%}}$. 508

509 We investigated the curious case mentioned 510
511 above. In the "negation" capability of the para- 512
513 phrase detection task, the error rates for patched 514
515 models remain same (12.3%) across four different 516
517 test suites; we found that the specific error cases 518
519 were also same regardless which test suite was used 520
521 to patch the model (Figure 3). This shows that, de- 522
523 spite overall strength to reveal more bugs for a 524
525 given *task* (Table 3), the TestAug is not guaranteed 526
527 to generate competitive test cases in all *linguistic 528*
529 *capabilities*. 530

531 Following the approach described in Section 4, 532
533 we created new templates to enrich the original 534
535 pool of templates available for $\mathcal{T}_{\text{Template}}$. When 536
537 generating new templates, restricting new slots 538
539 on only reappeared words in the original tem- 539
540 plates decrease the possibility of generating in- 540
541 valid templates (Table 12). We manually sampled 541
542 and annotated 100 generated templates per task 542
543 and found that the valid templates constitute 91%, 543
544 89%, and 91% of all templates for sentiment clas- 544
545 sification, paraphrase detection, and natural lan- 545
546 guage inference tasks. The invalid templates mostly 546
547 come from invalid modifiers such as "{ADJEC- 547
548 TIVE_OF_PERSON}" in "Some of the creams are 548
549 {ADJECTIVE_OF_PERSON} in colour.", where 549
550

⁵For the sentiment classification task, as our experiments require further fine-tuning on top of the already fine-tuned models; while the number of output classes do not match, we replaced the final 2-class classification layer with a newly initialized 3-class classification layer and fine-tuned the model for additional 3 epochs. We used the discretized 3-class SST dataset for this further fine-tuning.

⁶The expansion of test cases in NLI task requires alternative approaches and we leave it as future work.

Table 3: Model accuracy on held-out validation set and their overall error rate reduction using different test suites. The accuracy ACC% is computed over the original held-out dataset. The error rate reduction $\Delta_{\text{ERR}\%} = \text{ERR}\%_{\text{Unpatched}} - \text{ERR}\%_{\text{Patched}}$ follows the evaluation metrics introduced in Section 5.2. The exact identifiers of model checkpoints we used in experiments are listed in Table 11. Some cells are marked with "/" as we leave template expansion of the NLI task as future work. We used a small subset of original template-based test suite as demonstrations and the percentage is shown beside the task name.

	ERR%	ERR% _{Unpatched}	$\mathcal{T}_{\text{TestAug}}$		$\mathcal{T}_{\text{TestAug}} \setminus \mathcal{T}_{\text{Template}}$		$\mathcal{T}_{\text{TestAug}} \setminus \mathcal{T}_{\text{Expansion}}$		$\mathcal{T}_{\text{Template}}$	
			ERR% _{Patched}	$\Delta_{\text{ERR}\%}$	ERR% _{Patched}	$\Delta_{\text{ERR}\%}$	ERR% _{Patched}	$\Delta_{\text{ERR}\%}$	ERR% _{Patched}	$\Delta_{\text{ERR}\%}$
Sentiment Analysis 183 / 28921 = 0.6%										
DistillBERT	10.0	42.4	8.5	33.9	15.0	27.4	15.4	27.0	31.0	11.4
ALBERT	7.3	41.6	6.5	35.1	12.7	28.9	17.4	24.2	29.0	12.6
BERT _{Base}	7.6	40.9	4.1	36.8	7.5	33.4	6.2	34.7	16.5	24.4
RoBERTa _{Base}	5.7	36.6	4.4	32.1	6.3	30.2	6.0	30.6	10.2	26.3
Paraphrase Detection 54 / 11126 = 0.5%										
DistillBERT	10.3	45.4	3.7	41.8	8.8	36.6	6.5	38.9	11.9	33.5
ALBERTA	9.3	45.3	9.2	36.0	14.2	31.1	12.1	33.1	15.8	29.5
BERT _{Base}	9.1	51.6	2.7	48.9	5.2	46.4	4.2	47.3	8.0	43.6
Natural Language Inference 240 / 347531 = 0.1%										
DistillBERT	12.6	49.5	/	/	/	/	27.4	22.1	36.3	13.2
ALBERT	9.9	45.0	/	/	/	/	21.0	24.0	28.7	16.4
RoBERTa _{Large}	8.1	32.2	/	/	/	/	10.5	21.7	15.8	16.5

adjectives for people are misused for creams since slots are created oblivious of the contexts. Despite some invalid templates, template expansion leverages the improved scale and diversity of test cases and scales up the creation of template-based test suite $\mathcal{T}_{\text{Template}}$.

5.3 Evaluating the Diversity of TestAug Results

5.3.1 Evaluation Method

The linguistic diversity of an NLP test suite could be measured either from the test case level or the test suite level. We introduce the number of unique dependency paths as a proxy for linguistic diversity for each individual test case; while in test suite level, we use the metric for diversity in natural language generation – Self-BLEU.

Number of Unique Dependency Paths. Dependency parsing of a sentence returns a directed tree where there is a unique path from the root and every vertex. The arcs in the dependency tree are attributed with a fixed set of grammatical relations. The dependency tree approximates the semantic relations between predicates and their arguments (Jurafsky and Martin, 2000). We therefore propose to use the number of unique dependency paths to measure the richness of semantic relations.

Self-BLEU. Self-BLEU is an extension of the regular BLEU that evaluates the diversity of generated texts (Zhu et al., 2018). Given a list of texts $\hat{\mathcal{Y}} = \{\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_N\}$, Self-BLEU is the average

BLEU score between every single sentence and all other sentences,

$$\text{Self-BLEU}(\hat{\mathcal{Y}}) = \frac{1}{N} \sum_{i=1}^N \text{BLEU}(\{\hat{Y}_i\}, \hat{\mathcal{Y}}_{\neq i}) \quad (1)$$

When k is fixed, lower Self-BLEU score indicates a higher diversity of the sentence.

5.3.2 Evaluation Results

The test cases the annotators unanimously deemed consistent with the given test case description constitute test suites for respective tasks. After controlling for the number of test cases under each test case description, the linguistic diversity (Table 4) of the test suites $\mathcal{T}_{\text{GPT-3}}$ show substantial improvement over the template-based counterparts $\mathcal{T}_{\text{Template}}$: the Self-BLEU4 score has an decrease of at least 9.4% (the paraphrase detection task) and the number of unique dependency paths is of at least 2.18 times compared to the original test suite (the natural language inference task).

5.4 Evaluating the Validity of TestAug Results

5.4.1 Evaluation Method

Our experiments require test cases that have been verified consistent with the given test case description. Rather than creating templates or test cases from scratch, the human annotators in our system take a more efficient and effective role to correct mistakes made by the GPT-3. Specifically, we worked with human annotators to annotate each

Table 4: Linguistic diversity of test suites.

	Self-BLEU4 (\downarrow)	Number of Unique Dependency Paths (\uparrow)
Sentiment Analysis		
$\mathcal{T}_{\text{GPT-3}}$	0.558	548
$\mathcal{T}_{\text{Template}}$	0.778	88
Paraphrase Detection		
$\mathcal{T}_{\text{GPT-3}}$	0.587	957
$\mathcal{T}_{\text{Template}}$	0.645	113
Natural Language Inference		
$\mathcal{T}_{\text{GPT-3}}$	0.412	692
$\mathcal{T}_{\text{Template}}$	0.514	317

generated test case and decided whether it would be used for testing the NLP model.

5.4.2 Evaluation Results

A test case that satisfies the given test case description expresses the linguistic capability without grammatical errors. Two of the authors manually labeled each test case by checking whether it satisfied the given description; the test case they did not unanimously agree upon were considered ambiguous and therefore discarded. We used Cohen’s κ to measure the agreement of annotation. The annotation interface is shown in Figure 4.

We instructed GPT-3 to generate test cases with linguistic capabilities and seed sentences from CHECKLIST and LONLI dataset⁷ (Ribeiro et al., 2020; Tarunesh et al., 2021). The annotators manually checked whether each test case satisfied the test case description; the Cohen’s κ ranges between 0.434 and 0.450 for three tasks (Table 5), indicating moderate agreement (McHugh, 2012). The samples the annotators did not agree upon were discarded, leading to a test case description satisfiability from 74.5% to 84.5%; this shows that a significant portion of test cases generated following our approach satisfy the test case description.

Table 5: Annotation statistics on the test case description satisfiability.

	Cohen’s κ	Satisfiability (%)
Sentiment Analysis	0.434	82.2
Paraphrase Detection	0.450	84.5
Natural Language Inference	0.437	74.5

In order to reduce the required manual efforts

⁷Both datasets are under MIT license.

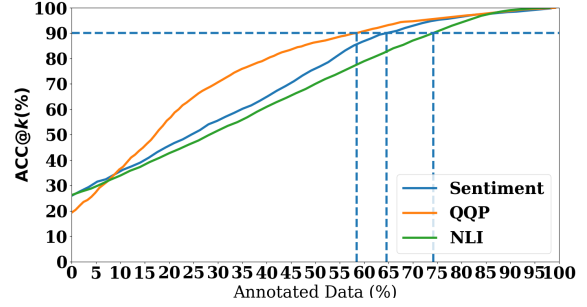


Figure 2: Ranker-assisting annotation accuracy versus annotation efforts.

while maintaining high annotation accuracy, we trained an ensemble of rankers with three large base language models and a ranking loss (Gao et al., 2021a); the rankers were trained to optimize the relative rankings of invalid and valid test cases, pushing the invalid ones up to the top (full details in Appendix A.2). With the help of this ranker ensemble, the annotators only need to check the test cases ranking at the top; the remaining test cases are all considered valid. To measure the annotation accuracy under this setting, we define $\text{ACC}@k$ as in Equation 2: only the minority samples after rank k are assigned incorrect labels while all the other samples are annotated correctly.

$$\text{ACC}@k = 1 - \frac{\sum_{i=k+1}^N 1(\hat{y}_i = l)}{N} \quad (2)$$

where N is total number of test cases, l is the majority label in the training set (in our case, the "valid" label), and \hat{y}_i is the validity label of the test case ranked at i -th position based on the ranking score. The assistance of this ranker helps reduce required annotation to maintain a 90% accuracy by ~25% to ~40% (Figure 2).

6 Discussion, Conclusions and Future Work

We introduced the TestAug framework to augment capability-based test suites to better reveal NLP models’ shortcomings in linguistic capabilities; empirical results have demonstrated the effectiveness of our framework. Looking forward, we plan to extend the set of NLP tasks supported by TestAug to more challenging tasks such as question answering (QA). We are also interested in further reducing manual efforts in TestAug by automating prompt design used for eliciting GPT-3.

Ethical Considerations

Annotator Rights

Two of the authors (one male and one female; both identified themselves as Asians) annotated the data following annotation guidelines; the guidelines are discussed and finalized after thorough discussions (the violations of these guidelines are discussed in Section 4.2). The estimated time commitment for labeling is 20 hours per annotator for a total of 8172 sentences (or sentence pairs). We acknowledge the annotators' efforts with a shared authorship.

Intended Uses

TestAug's intended use is as a tool to augment template-based test suites with newly generated test cases from GPT-3; two set of test cases are then used altogether to evaluate a NLP models' linguistic capabilities; we believe this application of existing datasets are consistent with their intended uses. We showed the effectiveness of this system in Section 5. We hope the adoption of TestAug into the NLP model development could make newly built NLP models more linguistically capable. Meanwhile, the TestAug includes GPT-3 as a component, we urge users of our system to follow the OpenAI's usage guidelines⁸.

Potential Misuse

TestAug might be misused to overestimate the models' linguistic capabilities. Specifically, even though failures on the test suites show models' shortcomings in a given linguistic capability, the absence of failures does *not* mean the models being tested are free from bugs; it is likely that test suites are not yet capable enough to reveal the model's bugs. We therefore call for a judicious interpretation of a NLP model's performance based on TestAug test suites. Moreover, we believe NLP testing is an iterative process; it might take multiple iterations of applying TestAug to reveal the model's issues in linguistic capabilities.

References

Bing Bai, Jian Liang, Guan Zhang, Hao Li, Kun Bai, and Fei Wang. 2021. Why attentions may not be interpretable? *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*.

⁸<https://beta.openai.com/docs/usage-guidelines>

- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*. 682–684
- Samuel R Bowman and George E Dahl. 2021. [What will it take to fix benchmarking in natural language understanding?](#) *The 2020 Conference of the North American Chapter of the Association for Computational Linguistics - Human Language Technologies (NAACL-HLT2020)*. 686–691
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). *arXiv preprint arXiv:2005.14165*. 692–696
- Robin Cooper, Dick Crouch, Jan Van Eijck, Chris Fox, Johan Van Genabith, Jan Jaspars, Hans Kamp, David Milward, Manfred Pinkal, Massimo Poesio, et al. 1996. Using the framework. Technical report, Technical Report LRE 62-051 D-16, The FraCaS Consortium. 697–702
- L. Floridi and Massimo Chiriatti. 2020. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30:681–694. 703–705
- Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021a. Rethink training of bert rerankers in multi-stage retrieval pipeline. In *ECIR*. 706–708
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021b. [Making pre-trained language models better few-shot learners](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3816–3830, Online. Association for Computational Linguistics. 709–716
- Mor Geva, Yoav Goldberg, and Jonathan Berant. 2019. Are we modeling the task or the annotator? an investigation of annotator bias in natural language understanding datasets. *ArXiv*, abs/1908.07898. 717–720
- Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. Breaking nli systems with sentences that require simple lexical inferences. In *ACL*. 721–723
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In *NAACL*. 724–726
- Paloma Jeretic, Alex Warstadt, Suvrat Bhooshan, and Adina Williams. 2020. [Are natural language inference models IMPPRESsive? Learning IMPLicature and PRESupposition](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8690–8705, Online. Association for Computational Linguistics. 728–734

846 Ishan Tarunesh, Somak Aditya, and Monojit Choud-
847 hury. 2021. [Lonli: An extensible framework for](#)
848 [testing diverse logical reasoning capabilities for nli.](#)
849 In *Thirty-Sixth AAAI Conference on Artificial Intelli-*
850 *gence (AAAI-22)*, volume abs/2112.02333.

851 Betty van Aken, Sebastian Herrmann, and Alexander
852 Löser. 2021. [What do you see in this patient? behav-](#)
853 [ioral testing of clinical nlp models.](#) *NeurIPS 2021*
854 *Research2Clinics Workshop, Bridging the Gap: From*
855 *Machine Learning Research to Clinical Practice.*

856 Jun Wang, Chang Xu, Francisco Guzmán, Ahmed
857 El-Kishky, Benjamin Rubinstein, and Trevor Cohn.
858 2021. [As easy as 1, 2, 3: Behavioural testing of](#)
859 [NMT systems for numerical translation.](#) In *Find-*
860 *ings of the Association for Computational Linguis-*
861 *tics: ACL-IJCNLP 2021*, pages 4711–4717, Online.
862 Association for Computational Linguistics.

863 Peter West, Chandrasekhar Bhagavatula, Jack Hessel,
864 Jena D. Hwang, Liwei Jiang, Ronan Le Bras, Ximing
865 Lu, Sean Welleck, and Yejin Choi. 2021. Symbolic
866 knowledge distillation: from general language mod-
867 els to commonsense models. *ArXiv*, abs/2110.07178.

868 Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, Ken-
869 taro Inui, Satoshi Sekine, Lasha Abzianidze, and Jo-
870 han Bos. 2019. [HELP: A dataset for identifying short-](#)
871 [comings of neural models in monotonicity reasoning.](#)
872 In *Proceedings of the Eighth Joint Conference on*
873 *Lexical and Computational Semantics (*SEM 2019)*,
874 pages 250–255, Minneapolis, Minnesota. Associa-
875 tion for Computational Linguistics.

876 Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan
877 Zhang, Jun Wang, and Yong Yu. 2018. Taxygen: A
878 benchmarking platform for text generation models.
879 *The 41st International ACM SIGIR Conference on*
880 *Research & Development in Information Retrieval.*

A Appendix 881

A.1 Details of Testing NLP Models 882

883 The error rate ERR% in Table 3 was obtained on
884 the original validation split for sentiment analysis
885 and paraphrase detection tasks and test split of nat-
886 ural language inference task; the statistics of these
887 datasets are reported in Table 7.

888 The models used for testing are reported in Table
889 11. The test cases used to evaluate these models
890 were those both annotators considered valid with
891 respect to the test case descriptions. The statistics
892 are shown in Table 6 (i.e., the "valid" column). Dur-
893 ing model testing, the test suites were partitioned
894 according Section 5.2 for $\mathcal{T}_{\text{TestAug}}$. We used one
895 Nvidia Tesla V100 (with 32 GB graphical mem-
896 ory) throughout the experiments. Following Sec-
897 tion 5.2, one testing session (one model of one
898 task) involves inference on the test split twice and
899 fine-tuning once on the training split, which takes
900 fewer than 5 minutes even for the large models.
901 The hyperparameters used for fine-tuning is fixed
902 and reported in Table 8.

A.2 Details of Training Rankers 903

904 We used the models listed below as base models
905 and trained an ensemble of rankers following the
906 setting that achieves the best performance in (Gao
907 et al., 2021a)⁹.

- bert-large-uncased 908
- google/electra-large-discriminator 909
- facebook/muppet-robetta-large 910

911 We increased the number of training epochs from
912 2 to 10 but left other choices to the authors' de-
913 faults. Each trained ranker returned a ranking list
914 with a potentially different ranking of the same
915 document. We aggregated ranking lists returned by
916 three rankers with the Borda count aggregator im-
917 plemented in pyrankagg¹⁰. We used the same
918 hardware as testing NLP models in Section A.1.

⁹Our implementation is based on the authors' code release:
<https://github.com/lyug/Reranker>

¹⁰<https://github.com/thelahunginjeet/pyrankagg>

Table 6: Statistics of test suites generated with GPT-3.

Task	Valid	Invalid
Sentiment Analysis	1607	347
Paraphrase Detection	1916	352
Natural Language Inference	2942	1008

Table 7: Statistics of datasets used to evaluate models' performances on held-out dataset.

Task	Split	Size	Dataset Identifier
Sentiment Analysis	Validation	872	sst2
Paraphrase Detection	Validation	40430	qqp
Natural Language Inference	Test	10000	snli

Table 8: Hyperparameter choice for model fine-tuning

Hyperparameter	Value
Learning rate	$5e - 6$
Batch size	8
Number of training epochs	3
Max. sequence length	128
Seed	42

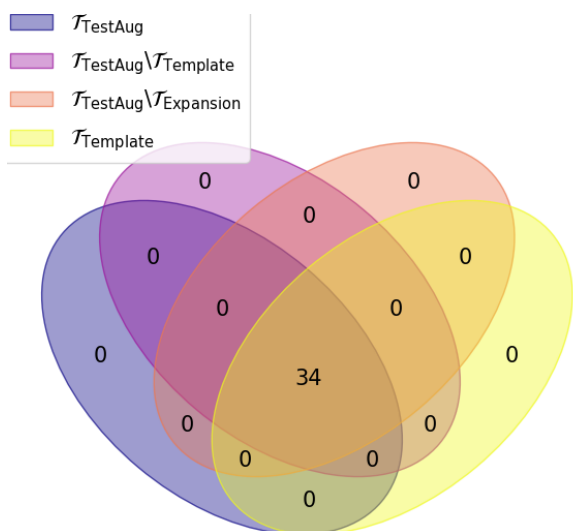


Figure 3: The error cases of the test sets for four test suites. The model made same mistakes across four test suites

Table 9: Prompt designs for paraphrase detection and natural language inference tasks. The prompt for natural language inference task follows (Liu et al., 2022), where "Implication", "Possibility", and "Contradiction" correspond to "entailment", "neutral", and "contradiction" label.

Paraphrase Detection
Two sentences are equivalent when using according to.
- {{ Who do analysts think is the smartest footballer in the world? }}
- { Who is the smartest footballer in the world according to analysts? }}
- {{ Who do students think is the top woman in the world? }}
- { Who is the top woman in the world according to students? }}
- {{ Who do readers think is the worst gamer in the world? }}
- { Who is the worst gamer in the world according to readers? }}
- {{ What does the data say about the most popular baby names? }}
- { What are the most popular baby names according to the data? }}
Natural Language Inference
Write a pair of sentences that have the same relationship as the previous examples. Examples:
- { Philip, Charles and Colin are the only children of Henry. }
- Implication: { Henry has exactly 3 children. }
- { Grace, Thomas and Helen are the only children of Andrea. }
- Implication: { Andrea has exactly 3 children. }
- { Don has 2 dollars. He received 8 more dollars. }
- Implication: { Don now has 10 dollars. }
- { Mary has a cat. She also has a dog. }
- Implication: { Mary has two pets. }

Table 10: Capability-wise error rate reductions for RoBERTa_{Base} (sentiment analysis), BERT_{Base} (paraphrase detection), RoBERTa_{Large} (natural language inference) (i.e., the most performant models in three tasks shown in Table 3). Some cells are marked with "/" as we leave template expansion of the NLI task as future work.

	ERR% _{Unpatched}	$\mathcal{T}_{\text{TestAug}}$ ERR% _{Patched}	$\Delta_{\text{ERR}\%}$	$\mathcal{T}_{\text{TestAug} \setminus \mathcal{T}_{\text{Template}}}$ ERR% _{Patched}	$\Delta_{\text{ERR}\%}$	$\mathcal{T}_{\text{TestAug} \setminus \mathcal{T}_{\text{Expansion}}}$ ERR% _{Patched}	$\Delta_{\text{ERR}\%}$	$\mathcal{T}_{\text{Template}}$ ERR% _{Patched}	$\Delta_{\text{ERR}\%}$
Sentiment Analysis									
Negation	30.9	3.2	27.7	4.5	26.4	4.3	26.6	9.6	21.3
SRL	54.9	7.8	47.1	9.8	45.2	9.4	45.5	12.4	42.5
Temporal	34.2	0.0	34.2	0.9	33.3	0.0	34.2	3.9	30.3
Vocabulary	9.9	2.7	7.2	7.2	2.7	6.9	3.0	11.1	-1.2
Paraphrase Detection									
Negation	12.6	0.4	12.3	0.4	12.3	0.4	12.3	0.4	12.3
SRL	37.4	8.2	29.1	17.9	19.4	12.6	24.7	17.4	20.0
Temporal	83.4	0.9	82.5	1.0	82.3	1.8	81.6	6.4	77.0
Vocabulary	15.5	2.5	13.0	3.7	11.8	3.1	12.4	8.1	7.5
Natural Language Inference									
Boolean	43.3	/	/	/	/	13.9	29.4	16.5	26.8
Causal	14.3	/	/	/	/	3.6	10.7	15.2	-0.9
Comparative	40.1	/	/	/	/	25.3	14.8	28.0	12.1
Conditional	65.4	/	/	/	/	16.7	48.7	23.5	41.8
Coreference	17.5	/	/	/	/	9.5	7.9	11.6	5.8
Implicature	51.3	/	/	/	/	24.4	26.9	31.1	20.2
Lexical	18.1	/	/	/	/	5.3	12.8	19.1	-1.0
Negation	7.2	/	/	/	/	0.0	7.2	5.2	2.0
Numerical	30.0	/	/	/	/	15.0	15.0	26.0	4.0
Presupposition	3.8	/	/	/	/	0.0	3.8	3.8	0.0
Quantifier	31.3	/	/	/	/	7.2	24.1	9.2	22.1
Relational	34.0	/	/	/	/	6.9	27.0	11.3	22.6
Spatial	46.2	/	/	/	/	16.9	29.2	20.0	26.2
Syntactic	3.3	/	/	/	/	3.3	0.0	3.3	0.0
Taxonomic	74.3	/	/	/	/	7.6	66.7	12.3	62.0
Temporal	35.4	/	/	/	/	24.0	11.5	27.1	8.3
World	12.1	/	/	/	/	1.6	10.4	4.9	7.1

Table 11: The fine-tuned models we evaluated in this paper.

Model name	Task	Size	Checkpoint Identifier
DistillBERT	Sentiment Analysis	Small	textattack/distilbert-base-cased-SST-2
ALBERT	Sentiment Analysis	Small	textattack/albert-base-v2-SST-2
BERT _{Base}	Sentiment Analysis	Base	textattack/bert-base-uncased-SST-2
RoBERTa _{Base}	Sentiment Analysis	Base	textattack/roberta-base-SST-2
DistillBERT	Paraphrase Detection	Small	textattack/distilbert-base-cased-QQP
ALBERT	Paraphrase Detection	Small	textattack/albert-base-v2-QQP
BERT _{Base}	Paraphrase Detection	Base	textattack/bert-base-uncased-QQP
DistillBERT	Natural Language Inference	Small	textattack/distilbert-base-cased-snli
ALBERT	Natural Language Inference	Small	textattack/albert-base-v2-snli
RoBERTa _{Large}	Natural Language Inference	Large	ynie/roberta-large-snli_mnli_fever_anli_R1_R2_R3-nli

Table 12: Creating new templates based on test cases generated by GPT-3.

Original Template and Test Case	Generated Test Case	New Template
Sentiment Analysis		
No one {pos_verb_present}s {the} {air_noun}. No one enjoys that seat.	This is not an easy service to appreciate. That customer service was not fun. I don't think your customer service is admired.	This is not an easy {air_noun} to appreciate. That customer {air_noun} was not fun. I don't think your customer {air_noun} is admired.
Paraphrase Detection		
Is it {mid} to {activity} before {hour}{ampm}? Is it {mid} to {activity} after {hour}{ampm}? Is it healthy to drink before 10am? Is it healthy to drink after 10am?	Is it bad to drink before 8pm Is it bad to drink after 8pm Is it acceptable to drink before 2pm Is it acceptable to drink after 2pm Is it advisable to eat before 8pm Is it advisable to eat after 8pm	Is it bad to {activity} before 8pm Is it bad to {activity} after 8pm Is it {mid} to {activity} before 2pm Is it {mid} to {activity} after 2pm Is it advisable to {activity} before 8pm Is it advisable to {activity} after 8pm

```

Label the generated sentence based on its validity:
A invalid sentence either
(1) Does not show the required linguistic capability in description, or
(2) Does not have correct label, or
(3) Includes private information or offensive contents.
Description:
    A neutral sentiment sentence with negative sentiment question and word no as the answer.
Label:
    1
Examples:
    Did we dislike this food? No
    Did I find the cabin crew? No
    Do I think that was a boring airline? No
GPT:
    Do I think the room is comfortable? No
Enter - the sentence is valid, Others - the sentence is NOT valid
>> Label: █

```

Figure 4: The command line interface for data annotation. Annotators are given a test case description and three examples from the template-based test suite; they are asked to the annotate the validity of the GPT-3-generated sentence (pair). Annotators are reminded of the guidelines for filtering invalid samples when labeling each sentence (pair) (shown at the top of the interface). We communicated explicitly for the intended uses of the annotated datasets before the annotation.