

## A Fairness metric

Dynabench comprises four dynamic tasks with multiple rounds of datasets that will grow over time. Given that here we have to be able to evaluate a wide variety of models, both in the loop and outside of it, we employ a black box post hoc approach, i.e., one that can be applied post-data collection to existing data, on any uploaded model, without requiring anything other than its predictions. One straightforward way to measure fairness then, is to apply clearly delimited, heuristic perturbations to existing evaluation datasets, and measure whether performance drops. Such an approach is similar to recent works that use grammars to heuristically generate pairs of examples varying in gender [58] and/or race [67] in that they utilize predefined lists of words. However, because we also want to ensure minimal consequences on our classification labels, we adopted an approach that is more targeted than grammars and also preserves the original input data distribution: we replace each word in the input data that has a clear signal about race/ethnicity and/or gender identity with a similar word referring to another group, rerun inference, and measure how many labels flipped (i.e., the difference in microaverage accuracy).

For race/ethnicity, we seeded first names across four demographic groups from a public dataset of 4,250 first names from US mortgage lending applications [69, 70]. These names cover 85.6 percent of the U.S. population, are based on the 1990 Census information on first name frequencies [70], and are licensed under a very permissive license (Creative Commons Attribution 4.0 International License). Names were associated with mutually exclusive demographic groups recognized by the 2000 and 2010 US Government Census—we selected the four groups with the most names: Asian/Pacific Islander, Black, Hispanic, white. Note: there is nothing inherently “racial” about particular names—for example, each demographic group had at least a few people named “Anna” or “Benjamin” [70]—although there are statistical trends. We selected all names for which a plurality of people of that name identified with one of the races or ethnicities, then took the 200 most frequent. We then also augmented that list with additional names popular in the literature [7, 48]. To construct our permuted test dataset, whenever we encountered a name in the input data from one of our lists, we randomly selected a name from a different race/ethnicity list and substituted it. Whenever we encountered a name in the input data which was not present in our list of names, we left it unperturbed.

For gender identity, we investigate two kinds of perturbations: names and noun phrases. For names, we affiliated the names from our race/ethnicity list with statistically likely genders, based on the U.S. Social Security Association’s Lists of Baby Names (1980–2019), and performed perturbations as we did for race/ethnicity.<sup>7</sup> For noun phrases (i.e., pronouns and nouns), we adopted a slightly more structured approach: we still replaced words based on pair-based word lists, but we didn’t do so randomly, as that could result in ungrammatical sentences (e.g., one can’t replace a pronoun, like “her” with a noun like “dad” and expect no effect on the classification label). Words in the paired list that either exclusively referred to women (e.g., *her*, *sister*) or to men (e.g., *his*, *brother*) were selected by taking the union of existing popular word lists [80, 81] that had been recently extended [13, 14].<sup>8</sup>

Given that our perturbations are heuristic, some noise is to be expected. For example, for the names perturbations, content relevant for the classification can be affected when the name is part of a phrase referring to a known named entity. Consider “I’ve always enjoyed eating at *Red Robin*” being perturbed to “I’ve always enjoyed eating at *Red Kayla*”. To mitigate this issue, we first ran an off-the-shelf named entity recognition system, and did not perturb any examples for which the system found a familiar named entity.<sup>9</sup> We observe that there are very few noisy examples (less than 5) resulting from NER errors. Finally, there is one irreducible type of noise arising from our heuristic approach—perturbing gender-explicit information occasionally results in unusual examples and may have consequences for the classification label: For NLI, the following hypothesis “Mothers should nurse at night” became “Fathers should nurse at night” in the context of “Failing to nurse at night can lead to painful engorgement or even breast infection”. Based on spot checks performed by the authors, we conclude that noise resulting from explicit gender information is also rare (only three out of 122 spot-checked sentences). Although our approach yields enough signal to evaluate whether model performance depends on race/ethnicity and gender identity, future work exploring more flexible and adaptable approaches is encouraged.

---

<sup>7</sup><https://www.ssa.gov/oact/babynames/>

<sup>8</sup>For a discussion of non-binary gender, see the Broader Impact Statement.

<sup>9</sup>We use the NER pipeline from spaCy [29], which was trained on OntoNotes 5.

## B Robustness metric

For the robustness evaluation, we use a post hoc black box approach similar to the fairness evaluation. We use TextFlint [24], an open source library for measuring model robustness that covers a wide range of text transformations. We apply a family of universal transformations from TextFlint, namely Contraction, Keyboard, Ocr, Punctuation, SpellingError, Typos and WordCase, as we focus on typographical errors for our robustness perturbations.

The robustness metric is computed as the percentage of unchanged predictions before and after perturbation. The assumption is that a robust model should not change its predictions upon such perturbations in input. That is to say, a model is deemed more robust if it has a higher robustness metric (i.e., a lower difference between original and perturbed examples).

## C Model Details

We selected and trained a diverse set of models to demonstrate the value of the platform, as well as to provide a sense of the current state of the art across our multiple metrics. We did not tune hyperparameters, instead using default training hyperparameters.

**Training Data** The NLI models are trained on ANLI [49] combined with MNLI [74], where ANLI is upsampled by a ratio of 2 to balance the data. The QA models are trained on SQuAD [57] combined with Adversarial QA [1]. The hate speech models are trained on rounds 0 through 4 of the Learning from the Worst dataset [71], with an upsampling ratio for each of 1, 5, 100, 1, and 1, respectively. The sentiment models are trained on of Dynasent [55], where round 2 was upsampled by a ratio of 3. Learning from the Worst [71], provides the upsample ratios used for hate speech, and we also performed early stopping for the hate speech transformers with the round 4 dev set. The rest of the upsampling was done to give the Dynabench data more influence in the training routine, as the leaderboard test sets on our evaluation platform are mostly comprised of Dynabench data.

**Transformers** All of the transformer models are the base versions provided in the HuggingFace transformers library [76], except ELECTRA which is the large discriminator version. ELECTRA-large was used to test the capacity of our framework to handle larger transformer models.

We used the default training hyperparameters from HuggingFace’s transformers training scripts.

The T5 models were trained with early stopping on SQuAD + Adversarial QA for QA, and the last round of the Dynabench data for the other tasks. For 2-way sequence classification, the model was trained to predict the first token of one of the strings "0", "1". For 3-way sequence classification, it was trained to predict the first token of one of the strings "0", "1", "2". There is no universal default way to handle which token to output from the T5 for sequence classification, as every sequence classification task is different. We used the training hyperparameters from the HuggingFace community T5 for text classification example<sup>[10]</sup>. An exception is that we trained for 3 epochs for text classification to match the BERT-style model default (the default for T5 is 2).

**FastText** All of the FastText models were the text classification version trained with an initial learning rate of 1, 25 epochs, 2-grams, bucket size of 200000, 50 dimensions, and hierarchical SoftMax. This is the default, for large datasets, that is listed on FastText’s website.<sup>[11]</sup>

**BiDAF** BiDAF was trained with the default settings from the AllenNLP [20] BiDAF trainer.<sup>[12]</sup>

**Majority Baseline** The majority label was determined from the corresponding train set. For example, the MNLI train set has frequency-matched labels, and the ANLI train set has entailed, neutral, contradictory splits of 52,111 / 68,789 / 41,965. So the majority baseline for NLI was to always return neutral.

<sup>10</sup>See <https://huggingface.co/transformers/community.html> and <https://github.com/patil-suraj/exploring-T5>

<sup>11</sup><https://fasttext.cc/docs/en/supervised-tutorial.html>

<sup>12</sup><https://github.com/allenai/allennlp>

## D Screenshots

Please also see the screen recording provided in the supplementary material to get a better sense of how the interface works.

OVERALL MODEL LEADERBOARD							
Model	Accuracy %	Throughput examples/second	Memory GiB	Fairness %	Robustness %	Dynascore	
DeBERTa default params (anon_user)	69.54	7.41	5.71	91.97	75.70	38.83	
RoBERTa default params (anon_user)	69.07	9.23	4.82	90.94	74.82	38.61	
ALBERT default params (anon_user)	67.29	9.60	2.18	89.94	74.12	37.72	
T5 default params (anon_user)	67.16	7.10	10.62	91.89	73.47	37.53	
BERT default params (anon_user)	64.82	9.39	4.13	92.11	66.38	36.36	
Majority Baseline (anon_user)	32.41	77.33	1.15	100.00	100.00	22.78	
FastText default params (anon_user)	31.29	73.94	2.20	83.23	69.14	21.13	

Figure 2: NLI Dynaboard example.

OVERALL MODEL LEADERBOARD							
Model	Accuracy %	Throughput examples/second	Memory GiB	Fairness %	Robustness %	Dynascore	
Dataset Weights							
<div style="display: flex; justify-content: space-around;"> <div> <p>snli-test</p> <p>mnil-test-mismatched</p> </div> <div> <p>anli-r1-test</p> <p>mnil-test-matched</p> </div> <div> <p>anli-r2-test</p> </div> <div> <p>anli-r3-test</p> </div> </div>							
DeBERTa default params (anon_user)	76.61	8.13	6.02	92.29	77.53	38.70	
snli-test	85.88	9.48	6.32	92.63	76.99		
anli-r1-test	65.10	5.23	5.21	91.41	72.28		
anli-r2-test	44.40	4.99	4.69	92.30	70.45		
anli-r3-test	45.83	5.96	4.71	89.66	69.82		
mnil-test-mismatched	87.74	9.41	6.63	93.54	82.35		
mnil-test-matched	88.31	9.38	6.67	92.29	82.33		
RoBERTa default params (anon_user)	76.00	10.24	5.15	91.48	76.72	38.40	
ALBERT default params (anon_user)	73.67	10.51	2.15	90.73	75.55	37.26	
T5 default params (anon_user)	73.77	7.91	10.31	91.81	75.00	37.25	
BERT default params (anon_user)	71.40	10.33	4.53	92.00	67.57	36.10	
Majority Baseline (anon_user)	32.06	99.00	1.16	100.00	100.00	16.92	
FastText default params (anon_user)	32.03	94.57	2.25	79.76	65.32	16.82	

Figure 3: NLI Dynaboard example, with dataset weight sliders and finer-grained metrics displayed.

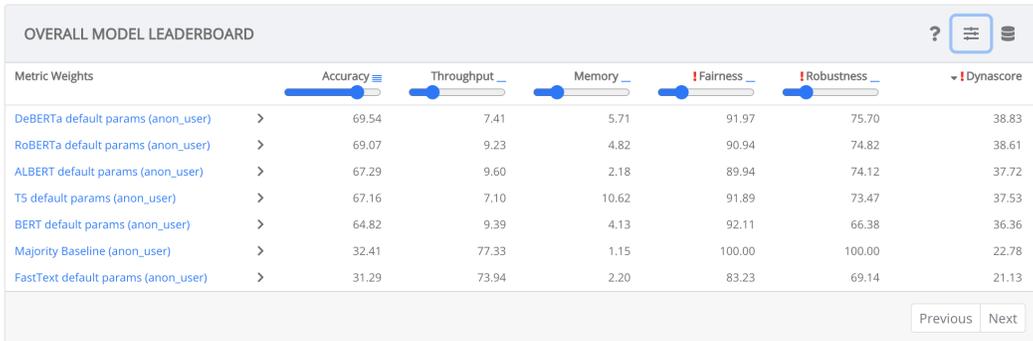


Figure 4: NLI Dynaboard example, with metric weight sliders in default configuration.



Figure 5: NLI Dynaboard example, with both dataset weights and metric weight sliders modified.

### Leaderboard Datasets

<a href="#">snli-test</a>	85.88
<a href="#">mnli-test-mismatched</a>	87.74
<a href="#">mnli-test-matched</a>	88.31
<a href="#">anli-r1-test</a>	65.10
<a href="#">anli-r2-test</a>	44.40
<a href="#">anli-r3-test</a>	45.83

### Non-Leaderboard Datasets

<a href="#">superglue-winogender</a>	59.55	
<a href="#">mnli-dev-mismatched</a>	88.27	
<a href="#">mnli-dev-matched</a>	89.13	
<a href="#">snli-dev</a>	85.49	
<a href="#">hans</a>	>	73.68
<a href="#">nli-stress-test</a>	>	77.56
<a href="#">anli-r1-dev</a>	63.90	
<a href="#">anli-r2-dev</a>	46.20	
<a href="#">anli-r3-dev</a>	45.42	

Figure 6: DeBERTa on NLI example, part of the model page, showing non-leaderboard evaluation-as-a-service dataset results as well.