| Model type | Original | Unpatched | Patched | | | |
|---|---|---|---|---|---|---|
| | | | CheckList | TestAug | TestAug $\setminus$GPT$-3$ | TestAug $\setminus$Expansion |
| **Sentiment Classification** | | | | | | |
| ALBERT | 7.3 | 32.6$\pm$5.7 | 13.4$\pm$6.5 | 11.3$\pm$10.0 | 10.6$\pm$6.6 | **9.6$\pm$8.2** |
| BERT$_{Base}$ | 7.6 | 33.9$\pm$6.1 | 9.0$\pm$4.2 | **8.3$\pm$4.2** | 8.5$\pm$1.6 | 9.9$\pm$4.9 |
| DistillBERT | 10.0 | 29.5$\pm$10.9 | 6.5$\pm$3.4 | **3.9$\pm$2.1** | 4.9$\pm$2.1 | 5.1$\pm$3.3 |
| RoBERTa$_{Base}$ | 5.7 | 14.2$\pm$6.1 | 3.7$\pm$2.3 | 1.6$\pm$1.0 | 2.7$\pm$2.7 | **1.4$\pm$1.2** |
| **Paraphrase Detection** | | | | | | |
| ALBERT | 9.3 | 38.1$\pm$3.8 | 7.1$\pm$0.8 | 0.6$\pm$0.4 | 5.8$\pm$1.8 | **0.4$\pm$0.4** |
| BERT$_{Base}$ | 9.1 | 36.0$\pm$4.9 | 6.2$\pm$1.5 | 0.5$\pm$0.4 | 5.6$\pm$1.1 | **0.4$\pm$0.3** |
| DistillBERT | 10.3 | 49.8$\pm$10.2 | 12.5$\pm$16.4 | **1.1$\pm$2.4** | 6.4$\pm$3.9 | 7.3$\pm$15.8 |
| **Natural Language Inference** | | | | | | |
| ALBERT | 9.9 | 42.8$\pm$1.9 | 30.1$\pm$4.2 | **23.0$\pm$1.6** | / | / |
| DistillBERT | 12.6 | 34.7$\pm$3.6 | 23.6$\pm$6.1 | **16.5$\pm$3.9** | / | / |
| RoBERTa$_{Large}$ | 8.1 | 17.8$\pm$4.0 | 8.3$\pm$3.1 | **8.0$\pm$3.1** | / | / |