

1 A Details for Datasets

2 **ReClor** ReClor contains 6,138 multiple-choice questions modified from standardized tests such
3 as GMAT and LSAT, which are randomly split into train/dev/test sets with 4,638/500/1,000 samples
4 respectively. It contains multiple logical reasoning types. The held-out test set is further divided into
5 EASY and HARD subsets based on the performance of BERT-based model [?].

6 **LogiQA** LogiQA consists of 8,678 multiple-choice questions collected from National Civil Ser-
7 vants Examinations of China and are manually translated into English by experts. The dataset is
8 randomly split into train/dev/test sets with 7,376/651/651 samples correspondingly. LogiQA also
9 contains various logical reasoning types.

10 **MuTual** MuTual has 8,860 dialogues annotated by linguist experts and high-quality annotators
11 from Chinese high school English listening comprehension test data. It is randomly split into
12 train/dev/test sets with 7,088/886/886 samples respectively. There more than 6 types of reason-
13 ing abilities reflected in MuTual. MuTual^{plus} is an advanced version, where one of the candidate
14 responses is replaced by a safe response (e.g., "*could you repeat that?*") for each example.