VARIERR NLI: Separating Annotation Error from Human Label Variation

Anonymous ACL submission

Abstract

Human label variation arises when annotators assign different labels to the same item for valid reasons, while annotation errors occur when la-004 bels are assigned for invalid reasons. These two issues are prevalent in NLP benchmarks, yet existing research has studied them in isolation. To the best of our knowledge, there exists no prior work that focuses on teasing apart error from signal, especially in cases where signal is beyond black-and-white. To fill this gap, we introduce a systematic methodology and a new dataset, VARIERR (variation versus error), focusing on the NLI task in English. We 013 propose a 2-round annotation scheme with annotators explaining each label and subsequently judging the validity of label-explanation pairs. VARIERR contains 7,574 validity judgments on 017 1,933 explanations for 500 re-annotated NLI 019 items. We assess the effectiveness of various automatic error detection (AED) methods and GPTs in uncovering errors versus human label variation. We find that state-of-the-art AED methods significantly underperform compared to GPTs and humans. While GPT-4 is the best system, it still falls short of human performance. Our methodology is applicable beyond NLI, offering fertile ground for future research on error versus plausible variation, which in turn can yield better and more trustworthy NLP systems.

1 Introduction

Labeled data plays a crucial role in modern machine learning (ML) (e.g., Mazumder et al., 2023). Data quality impacts ML performance and in turn user trust. It is therefore of vital importance to aim at high-quality consistently-labeled benchmark data (e.g., Bowman and Dahl, 2021). However, recent research has revealed a notable presence of *annotation errors* in widely-used NLP benchmarks (Klie et al., 2023; Rücker and Akbik, 2023). Similar observations were made recently in computer vision (CV) (Northcutt et al., 2021; Vasudevan et al., 2022; Schmarje et al., 2023).



Figure 1: Variation or Error? We present a procedure and multi-label dataset, VARIERR, to tease apart annotation error from plausible human label variation. We leverage *ecologically valid explanations* and *validation* as two key mechanisms (boxed: self-validations; label "Contradicts" is an *error*); see §3-§4 for details.

043

044

045

046

047

052

056

058

060

061

063

064

065

067

At the same time, increasing evidence exists that for many items in many tasks, more than a single label is valid. For some items, systematic variation exists for valid reasons, such as plausible disagreement or multiple interpretations. In other words, the world is not just black and white. Human label variation (HLV, as termed by Plank 2022) has been shown on a wide range of NLP tasks (de Marneffe et al., 2012; Plank et al., 2014; Aroyo and Welty, 2015), including in natural language inference (NLI; Pavlick and Kwiatkowski 2019; Zhang and de Marneffe 2021) as well as in Computer Vision (CV; Peterson et al. 2019; Uma et al. 2021). NLI involves determining whether a hypothesis is true (Entailment), false (Contradiction), or neither (Neutral), assuming the truth of a given premise (cf. Figure 1 for an example with plausible labels).

Although high-quality consistently-labeled data may initially appear to conflict with the goal of accommodating HLV, it is important to emphasize that we do not perceive these as contradictory goals. While HLV exists, so do errors. We assert that annotators are inevitably prone to make errors, such as misunderstanding instructions or accidentally selecting a wrong label. Optimizing data quality

is essential through providing clear instructions 068 and effective training, and identifying annotation errors yields better datasets (Larson et al., 2019). However, still little is known about what constitutes an error versus plausible variation. We lack both a theory and operationalizable procedures to tease apart error from plausible HLV consistently and soundly. Some datasets with errors (and their corrections) exist, and there has been work on auto-076 matic error detection (AED). However, both have 077 their limitations (\S 2). A crucial gap remains: a lack of examination in real-world scenarios where the signal is nuanced, not merely black-and-white.

To address this gap, this paper contributes: (i) VARIERR, a novel multi-annotator English NLI dataset with both plausible variation and detected errors. To the best of our knowledge, no such dataset exists yet. (ii) A new methodology to detect errors: we collect multiple annotations, where each label comes with an ecologically valid explanation inspired by Jiang et al. (2023), and propose to pair them with validity judgments to identify errors. (iii) Finally, we benchmark existing AED methods and GPTs in a challenging setup, where the task is to tease apart error from plausible human label variation. Our findings indicate that existing AED methods underperform humans and GPTs substantially and highlight the need for further research. To facilitate uptake, we release our new VARIERR dataset and code on GitHub upon publication.

2 Related Work

081

086

087

089

094

100

101

103

104

105

106

107

108

110

111

112

113

114

115

116

117

Labeled data is the fuel of machine learning, as it drives both learning and evaluation. Following a data-centric view, we focus on improving data quality over data quantity (Motamedi et al., 2021; Swayamdipta et al., 2020; Zhang et al., 2021; Gordon et al., 2022). We aim to bring together work on data quality from two ends: annotation error vs. human label variation.

Annotation Errors and AED Several recent work has found errors in widely used benchmarks, such as CoNLL 2003 for Named Entity Recognition (Wang et al., 2019; Reiss et al., 2020; Rücker and Akbik, 2023), TACRED for relation extraction (Alt et al., 2020), WSJ for syntax (Manning, 2011; Dickinson and Meurers, 2003), and ImageNet for object classification (Beyer et al., 2020; Northcutt et al., 2021; Vasudevan et al., 2022).

AED has a long-standing tradition in NLP. Proposed methods range from early work that relies on variation-based methods positing that instances with similar surface forms tend to have the same label (Dickinson and Meurers, 2003; Plank et al., 2014) to more recent model-based approaches that either exploit predictions (Amiri et al., 2018; Arazo et al., 2019) or information derived from training dynamics (Swayamdipta et al., 2020); see Klie et al. (2023) for a survey on AED. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

153

154

155

156

157

158

159

160

161

162

164

165

166

167

168

Flaggers and scorers for AED have been proposed (Klie et al., 2023). Flaggers detect errors by providing a hard decision of whether an instance is erroneous. Scorers, on the other hand, assign a score to each instance reflecting the likelihood of being an error, and the top-*n* scoring instances are then corrected. Here, we focus on scoring methods to rank instances. Most of the AED work mentioned has limitations as they either rely on posthoc mining of errors (and might therefore miss out on errors) in semi-automatic ways (e.g. Reiss et al., 2020), or they inject synthetic noise which has been shown to result in datasets where errors are easy to spot (Larson et al., 2019). Instead of using synthetic noise, we focus on realistic setups and re-annotate data in ecologically valid ways.

Human Label Variation (HLV) Recent studies have drawn attention to HLV in NLP (i.a., Uma et al., 2021; Plank, 2022). HLV has been described as annotator disagreement, which is not just noise but also signal since a sign of vagueness or ambiguity can benefit models (Aroyo and Welty, 2013). These include judgments that are not always categorical (de Marneffe et al., 2012), inherent disagreement (Pavlick and Kwiatkowski, 2019; Davani et al., 2022), or justified and informative disagreement (Sommerauer et al., 2020). For subjective NLP tasks, which by essence encourage annotator subjectivity (and hence variation), there is also a line of work referred to as perspectivism (Basile et al., 2021), connected to the descriptive data annotation framework proposed by Rottger et al. (2022).

HLV in NLI This paper focuses on NLI, known to contain HLV (Pavlick and Kwiatkowski, 2019; Nie et al., 2020; Jiang and de Marneffe, 2022; Jiang et al., 2023). Pavlick and Kwiatkowski (2019) reannotated nearly 500 NLI instances with 50 crowdworkers and showed that disagreements in NLI cannot be dismissed as annotation "noise." ChaosNLI (Nie et al., 2020) pioneers large-scale NLI annotation by collecting 100 annotations per instance for 3K items from SNLI (Bowman et al., 2015), α NLI (Bhagavatula et al., 2020), and MNLI (Williams

175

176

177

178

179

181

182

183

184

186

187

188

190

191

192

193

194

195

196

197

199

203

207

210

211

212

213

214

215

169

et al., 2018) but for which the original annotations did not yield high agreement. They show that, for most of the items, HLV persists with more annotations. Further, their experiments show a large room for model improvement and a positive correlation between human agreement and label accuracy.

In another line of work, Jiang and de Marneffe (2022) identified *reasons* for observing variation in NLI, deriving a taxonomy based on linguistic properties of the items. Following up on that work, Jiang et al. (2023) proposed LIVENLI, to gain insights into the origin of label variation. They reannotated 122 NLI instances from ChaosNLI with ecologically valid explanations: annotators are instructed to not only provide NLI labels but also explanations for their label choices. This addresses a limitation of prior work that uses post-hoc explanations, which may not reflect the true reasons of the original annotators, thereby questioning the validity of the prior method. They show that ecologically valid explanations have an additional benefit: signaling within-label variation, i.e., annotators give the same label but for different reasons. While we do not focus on the latter here, we take inspiration from Jiang et al. (2023) to collect ecologically valid explanations (cf. $\S3.1$).

> To the best of our knowledge, there remains a gap for studies on *both* annotation errors and human label variation in a concentrated effort. It is thus an open challenge to define error in an ecologically valid way, and it is unknown to what extent existing AED methods help detect such errors and whether new methods are needed. To find answers to these challenging open questions, we believe it is important to move both directions forward.

3 VARIERR: Annotation Procedure

To tease apart human label variation from error, we create VARIERR (Variation versus Error), an NLI dataset with two rounds of annotations by four annotators:¹ Round 1 for NLI labels and explanations (§3.1) and Round 2 for validity judgments (§3.2).

3.1 Round 1: NLI Labels & Explanations

We collect annotations from four annotators on 500
NLI items randomly sampled from the MNLI subset of ChaosNLI. Annotators were asked to provide not only one or more NLI labels (E: Entailment, N: Neutral, C: Contradiction) to each item but also a

one-sentence explanation for each label they chose,
as the same label could be chosen for different rea-
sons (Jiang et al., 2023). Annotators could use a
fourth "I don't know" (IDK) label if none of the
NLI labels seemed suitable. Round 1 annotation
sums up to 1,933 label-explanation pairs for the
500 items after discarding 331 "IDK" annotations
and keeping only standard NLI labels.216
217

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

3.2 Round 2: Validity Judgments

VARIERR's key contribution lies in proposing a second round of *validity judgment*. Validity judgment mirrors conventional annotation adjudication in that annotators view NLI labels and explanations from each other. However, this information is delivered anonymously to annotators to reduce group dynamics. Further, rather than agreeing on a single label or explanation altogether, annotators are free to make independent judgments regarding what is an error versus a plausible variation.

In this second round, annotators become judges. For all 500 items, the 1,933 label-explanation pairs from Round 1 are distributed anonymously to the same four annotators in the same order as in Round 1. For each NLI item, each judge sees all labelexplanation pairs annotated in Round 1 in a random order, including their own, which they may or may not remember. For each label-explanation pair, the annotator judges whether the explanation makes sense given the label, answering "yes" (\checkmark), "no" (X) or "IDK". Round 2 includes 7,574 validity judgments after discarding 158 "IDKs".

4 VARIERR: Detecting Errors

Multiple validity judgments on label-explanations enable distinguishing annotation errors from HLV.

4.1 Self versus Peer

One consequential feature of our two-round multiannotator scheme is the distinction between self- vs. peer-judgments. *Self-judgments* refer to Round 2 judgments on the judge's own Round 1 labelexplanation annotations whereas *peer-judgments* refer to judgments from other annotators.

4.2 Validating Labels

Let $\mathcal{A} = \{a_1, .., a_4\}$ be the set of annotators.

Self-validated Label-Explanation A labelexplanation pair given by annotator *a* on an item in Round 1 is *self-validated* iff *a* judges in Round 2 that the annotated pair makes sense.

¹Annotators are master students in CompLing and the first author of this paper, paid according to the national standard.

Premise: Because marginal costs are very low, a newspaper price for preprints might be as low as 5 or 6 cents per piece.

Hypothesis: Newspaper preprints can cost as much as \$5.

Label-explanation pairs: Before:{E:1,N:2,C:1} Self-validated:{N:2} Peer-validated:{N:2,C:1} Label: [N] Errors:[E.C]

	Round 1: NLI Label & Explanation			ound 2: Validity				
L	А	Explanation		2	3	4		
Е	4	5 dollars for a piece of newspaper.	×	×	×	×		
N	1	The context only mentions how low the price may be, not how high it may be.	1	1	1	1		
	3	The maximum cost of newspaper preprints is not given in the context.	\checkmark	1	1	1		
С	2	The context says 5 or 6 cents, not \$5.	×	×	1	1		

(a) id: 72870c

Premise: They made little effort, despite the Jesuit presence in Asia, to convert local inhabitants to Christianity or to expand their territory into the interior.

Hypothesis: The Jesuit presence in Asia helped to convert local residents to Christianity,

allowing them to expand their territory.

Label-explanation pairs: Before:{E:1,C:4} Self-validated:{C:3} Peer-validated:{C:4}

Lab	el: [C] Error:[E]				
	Round 1: NLI Label & Explanation					lity
L	А	Explanation	1	2	3	4
Е	1	Both premise and hypothesis suggest that the speaker does not understand.	×	×	×	×
	1	The Jesuit didn't make much effort to convert local residents to Christianity or to expand their territory.	 Image: A start of the start of	1	1	1
С	2	They did not try to expand their territory.	1	?	1	1
	3	The Jesuit did not make effort to convert local residents to Christianity or to expand their territory.	1	~	✓	1
	4	They made little effort to convert the locals or to expand their territory. So they did not help.	1	1	1	 Image: A start of the start of
		(b) <i>id</i> : 28306c				

Table 1: Sample annotations from VARIERRNLI corpus. L: Label, A: Annotator; E: Entailment, N: Neutral, C: Contradiction; ✓: 'yes'; X: 'no'; ?: 'idk'; magenta : self-judgments, black: peer-judgments, Err : label error.

Peer-validated Label-Explanation A labelexplanation pair given by annotator *a* in Round 1 is *peer-validated* iff the majority (≥ 2) of the other annotators $\mathcal{A} \setminus \{a\}$ approves the pair in Round 2.

264

265

267

268

271

272

277

278

279

What counts as an error? We define an NLI label as an error if all label-explanation pairs are not self-validated. In other words, a label is viewed as correctly attributed to an item iff any of its explanations is self-validated.

Importantly, in our Round 2 annotation, a labelexplanation pair might be considered wrong in retrospect by the annotator who wrote it after reading all label-explanation pairs given to that item by four annotators. For instance, this is the case for **E** and **C** in Table 1a as well **C** in Table 1b. To comprehensively assess annotation variations in Round 1, we work with a strict definition of error, for which all annotators need to re-evaluate their previous annotations in light of others' annotations. Annotator identity is not revealed in the process.

4.3 Data Statistics & IAA

Table 2 shows the frequencies of NLI labels across the four annotators on the 500 items and 1,933 explanations before and after validation. We in-

Validation	FreqType	Е	Ν	С	Σ	IAA	
hafana validation	repeated	554	977	402	1,933	0.25	
before validation	aggregated	263	403	212	878	0.55	
colf validated	repeated	467	916	329	1,712	0.50	
sen-vanualeu	aggregated	210	380	159	749	0.50	
maar validatad	repeated	446	859	296	1,601	0.60	
peer-validated	aggregated	177	335	130	642	0.09	

Table 2: Frequency counts and inter-annotator agreement (Krippendorff's α with MASI-distance) on non-, self-, and peer-validated VARIERR NLI labels.

clude statistics on *repeated* frequency counts (e.g., E counts twice if it is given as a label by two annotators for the same item) and *aggregated* labels (repeated labels for a given item count once). Moreover, following Jiang et al. (2023), we compute inter-annotator agreement on NLI labels using Krippendorff's α (for multi-annotator) with MASI-distance (for multi-label).

Results show that self-validated annotations achieve a much higher IAA than without validation (see A.1 for pairwise IAA). Though the repeated and aggregated frequencies of NLI labels decrease adequately after validation, HLV is still preserved in self- and peer-validated annotations, averaging 1.50 (749/500) and 1.28 (642/500) labels/item.



(a) Number of label-explanations rejected by selfand-peer, self-only, and peer-only validations.

304

305

307

311

312

313

314

316

317

320

321

322

328

330

331

(b) NLI label sets on non-, self- and peer-validated items.

Figure 2: Frequency statistics on VARIERR.

We also observe in Table 2 that 88.57% (1,712/1,933) of Round 1 explanations in VARI-ERR were self-validated and 82.82% (1,601) were peer-validated. Figure 2a presents the number of label-explanation pairs rejected by both self- and peer-validations, by self-validation only, and by peer-validation only. Most Entailment and Contradiction annotations rejected by self are also rejected by peers (dark green). However, Neutral presents a challenging situation for self-validation where 60.13% (92/153) of Ns is only invalidated by the joint force of peers but not by one annotator alone.

Figure 2b demonstrates frequencies of aggregated label combinations per item before validation and after self- and peer-validations (see A.2 for label-explanation pair frequencies). Frequencies of multi-labeled items drop after self-validation and, more remarkably, after peer validation. Inversely, the number of single-labeled items increases vastly, especially for N. Additionally, we also observe from VARIERR that a large portion of items, 37.6% (188/500), are self-identified as errors and 51.6% (258) are rejected by peer-validation.

In sum, though HLV remains in VARIERR, our validation process demonstrates that annotation errors are frequently concealed under HLV. We thus proceed with the challenging automatic error detection task in §5-6 to separate annotation errors from valid HLVs.

5 Automatic Error Detection (AED) on VARIERR

We now describe our experiments to detect annotation errors using VARIERR automatically. We evaluate the capabilities of AED methods, LLMs, and human heuristics (all henceforth *scorers*) in capturing ecologically detected annotation errors.

336

337

339

340

341

342

343

345

347

348

349

351

353

354

355

356

357

359

361

362

363

5.1 Task Definition and Evaluation

Following Klie et al. (2023); Weber and Plank (2023), we model AED as a ranking task. In this setting, the goal of the *scorer* is to provide a ranked list with the labels that are most likely errors at the top and the most likely correct ones at the bottom. This ranked list could then be used to guide reannotation efforts (Alt et al., 2020; Northcutt et al., 2021) or remove the most likely errors from the training data (Huang et al., 2019). Scorers produce such a list by assigning an error score to each assigned label in the dataset. They derive the ranked list by sorting it according to the assigned scores.

We evaluate scorers on VARIERR using the following protocol. A model receives the list of NLI items from VARIERR, where each item is paired with the label(s) it received in Round 1. For the 500 items in VARIERR, the model is given a list of 878 item-label pairs. Based on that information, the model assigns an error score and ranks the labels by this score. We evaluate how well the model performs by comparing this ranked list with the self-flagged errors. Following Klie et al. (2023), we use standard ranking metrics for evaluation: average precision (AP), i.e., the area under the precision/recall curve computed over all assigned labels, and precision/recall for the top 100 ranked labels, P@100 and R@100.

377

391

393

400

401

402

403

404

405

406

407

408

409

5.2 **Baselines and Models**

We evaluate five different AED models: two variants of Datamaps (DM, Swayamdipta et al. 2020), Metadata Archaeology (MA, Siddiqui et al. 2023), and two GPTs. We report the mean and standard 370 deviation over three random seeds for DM and MA. **Datamaps (DM)** We use training dynamics (i.e., the collection of training statistics over epochs E) for each label. These statistics are obtained by training a DistilRoBERTa-base model² (Sanh et al., 2019) following Klie et al. (2023) in a multi-label 376 setting (Jiang and de Marneffe, 2022) on all labels of VARIERR obtained in Round 1. We refer to 378 the j'th label of the i'th example as $label_{i,j}$. The training dynamics are modeled by the probability $p_{i,j,e}$ that DistilRoBERTa predicts for label_{*i*,*j*} after the *e*'th epoch. Based on these probabilities, the two DM models we use are defined as follows:

$$DM_{\text{mean}} = -\frac{1}{E} \sum_{e=1}^{E} p_{i,j,e} \tag{1}$$

$$DM_{\rm std} = \sqrt{\frac{1}{E} \left(\sum_{e=1}^{E} p_{i,j,e} + DM_{\rm mean}\right)^2} \qquad (2)$$

Note that a *low* average probability for the label indicates a likely error. Because our evaluation setup requires the most likely errors to be ranked first, we negate the average probabilities.

Metadata Archaeology (MA) MA models AED as a supervised task. It represents each instance (or label in our case) as the *E*-dimensional $-\log p_{i,i,e}$ vector, where E is the number of epochs and $p_{i,i,e}$ is the probability the model assigns to the *j*'th label of the *i*'th NLI instance at epoch *e*. Then, it assumes that some instances are labeled with the property of interest (in our case, whether it is an erroneous label). It predicts whether an instance is an error by employing a k-nearest neighbors (kNN) classifier using the instance representations and error labels. We use the average number of annotated errors for the kNN to obtain a score for each instance. Following Siddiqui et al. (2023), we use k = 20. To obtain unbiased predictions, we require that the kNN training instances are distinct from those we want to obtain predictions for. We use a 2-fold cross-validation setup where we split VARIERR into two folds, use one half as ground truth, obtain the predictions for the other, and vice versa.

²https://huggingface.co/distilroberta-base

GPT We also compare two large language mod-410 els (LLMs): GPT-3.5 (Brown et al., 2020) and 411 GPT-4 (OpenAI, 2023). We emulate the Round 2 412 annotation setting in $\S3.2$ as closely as possible by 413 prompting each model to provide a score reflect-414 ing how much each Round 1 explanation makes 415 sense for a given label. We compute the score 416 per label by averaging the GPT-assigned scores 417 of all explanations for the label. We prompt GPT 418 as follows, giving it the premise (context) and hy-419 pothesis (statement) of an NLI item as well as all 420 label-explanation pairs, asking it then to provide a 421 probability for each reason: 422 423

424

425 426

427

428

429

430

431 432

433

434 435

436

437

438 439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

System: You are an expert linguistic annotator. User: We have collected annotations for an NLI instance together with reasons for the labels. Your task is to judge whether the reasons make sense for the label. Provide the probability (0.0 - 1.0) that the reason makes sense for the label. Give ONLY the reason and the probability, no other words or explanation. For example: Reason: <The verbatim copy of the reason> Probability: <the probability between 0.0 and 1.0 that the reason makes sense for the label. without any extra commentary whatsoever; just the probability!>. Context: {CONTEXT} Statement: {STATEMENT} Reason for label {LABEL}: {REASON_1} Reason for label {LABEL}: {REASON_2} [...] Reason for label {LABEL}: {REASON_n} Reason {REASON_1} Probability:

We implement GPTs using sglang (Zheng et al., 2023) and its default sampling parameters. See Appendix B for a complete prompt example. Note that the GPTs have access to the explanations for the labels, whereas the other models described above only have access to the labels without explanations.

5.3 Human Heuristics

In addition to the above automatic means, we experiment with four human heuristics that use the human label distributions over NLI labels (E, N. C) from annotation efforts: label counts from ChaosNLI (100 annotators) and VARIERR (4 annotators). In addition, we compare to VARIERR's total and average peer-judgments over explanations. Label Counts (LC): ChaosNLI & VARIERR We hypothesize that if multiple annotators choose

the same label, there is a high likelihood that 468 the label is a correct annotation. We imple-469 ment two label count (LC) baselines: one using 470 ChaosNLI (Nie et al., 2020) and one using VARI-471 ERR. ChaosNLI includes 100 crowd-sourced an-472 notations for each NLI item. Since VARIERR is 473 a subset of ChaosNLI items, we use label counts 474 from ChaosNLI (LC_{CHAOS}) as a human heuristic 475 to score Round 1 labels on each item, i.e., how 476 many of the 100 crowd-workers annotated label_{*i*, *i*} 477 on item i. For instance, the ChaosNLI human dis-478 tribution is {N:25, E:72, C:3} for the example in 479 Figure 1. Similarly, we include $LC_{VARIERR}$ that 480 counts the number of annotators (4 in total) that as-481 signs label_{*i*, *i*} to item *i* in VARIERR's Round 1 NLI 482 labels. We multiply both LC_{CHAOS} and LC_{VARIERR} 483 by -1, proposing that if a label has a higher count, 484 then it is less likely to be an error. 485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

503

504

510

511 512

513

Peer-judgments (Peer) in VARIERR VARI-ERR's 2-round annotations enable more finegrained human heuristics that engage judgments on label-explanation pairs. Since each label_{*i*,*j*} can be assigned by multiple annotators with different explanations, we count the number of "yes" judgments on explanations from peers – excluding selfjudgments since those are used for gold error labels.

We implement two peer heuristics: Peer_{sum} and Peeravg. Peersum sums all "yes" judgments across multiple explanations on the same label, and Peeravg sums "yes" judgments within each explanation and then averages across explanations within the label. Given that one label can maximally receive four explanations, it can receive up to 12 peer-judgments, 3 per explanation. For example, C in Table 1b receives 11 peer-judged "yes" in total (Peer_{sum} = 3 + 2 + 3 + 3 = 11), and the average over four explanations is $\text{Peer}_{\text{avg}} = 11/4 = 2.75$. Peer_{avg} differentiates more from Peersum when there are multiple explanations, but each receives sparse "yes" judgments. For example, N in Table 5a (Appendix C) receives 3 "yes" judgments but across two explanations, resulting in Peer_{sum} = 2 + 1 = 3and Peer_{avg} = 3/2 = 1.5. Similarly the label counts above, we multiply both Peer_{SUM} and Peer_{AVG} by -1, hypothesizing that fewer "yes" judgments indicate a higher likelihood to be an annotation error.

514Combining Label Counts and ModelsRanking515labels by the number of annotations they received516in Round 1 is a very strong baseline; see LCVARIERR517in Table 3. Inspired by Nogueira et al. (2019), we518investigate an approach that re-ranks the predic-

tions of $LC_{VARIERR}$ by breaking ties with the scores produced by another model (e.g. MA or GPT-4). Note that $LC_{VARIERR}$ produces many ties because its score is always one of $\{-1, -2, -3, -4\}$.

6 **Results for AED on VARIERR**

Table 3 presents human and model performances on VARIERR AED using the ranking setup in §5.

Method	AP	P@100	R@100	AP (rerank)
		Baselines		
Random	14.7	14.7	11.4	-
		Models		
MA	17.7 ± 1.5	18.3 ± 4.2	14.2 ± 3.2	44.2 ± 3.0
DM _{mean}	22.8 ± 0.4	23.7 ± 2.1	18.3 ± 1.6	50.4 ± 0.7
DM _{std}	22.3 ± 1.9	22.7 ± 1.2	17.6 ± 0.9	50.0 ± 1.5
GPT-3.5	17.6	21.0	16.3	37.6
GPT-4	31.3	46.0	35.9	47.4
		Human		
LC _{CHAOS}	32.5	35.0	27.3	49.8
LC _{VARIERR}	40.8	42.0	32.6	40.8
Peeravg	42.2	46.0	35.9	47.8
Peersum	46.5	47.0	36.7	47.8

Table 3: Results for AED on VARIERR. *AP*: average precision; *rerank* denotes using the method to break ties in $LC_{VARIERR}$. For MA and DM, we report mean and standard deviation over three random seeds. Note that GPTs have access to explanations.

6.1 Human Performance

The best human heuristic is from peers (Peer_{sum}), reaching a performance of 46.5% AP, 47% precision@100, and 36.7% recall@100. These numbers support our hypothesis that human validation can be used as a strong means to detect annotation errors in a task with relatively high HLV because self- and peer-rejected label-explanation pairs overlap considerably (cf. Figure 2a). Interestingly, both peer-derived heuristics from VARIERR perform better than LC_{CHAOS} (3 linguists versus 100 crowd-workers), which suggests that having few highly trained expert annotators is sufficient for reliable error detection, outperforming a larger number of crowd-workers. LC_{VARIERR} outperforming LC_{CHAOS} on all metrics strengthens this finding. Next we compare humans to automatic means.

6.2 Model Performance

Among the models, GPT-4 outperforms all other methods by a large margin with a 8.5/22.3/17.6percentage points (pp.) improvement in terms of AP / P@100 / R@100 over the second best model DM_{mean}. GPT-4 even outperforms LC_{CHAOS} in P@100 and R@100 and is close to the best peer

527

528

519

520

521

522

523

524

525

- 540 541 542
- 543 544 545

546

547

548



Figure 3: Correlations among scorer predictions.

552

553

554

555

557

559

560

562

566

567

571

574

575

577

578

581

582

583

heuristic for these two metrics. One might postulate that ChaosNLI could have been part of GPT-4's training mixture, and GPT-4 performed well by reproducing its probabilities. To check whether this is the case, we compute Pearson's r between the predictions of all scorers (Figure 3). While GPT-4 has a slightly higher correlation (0.42) with LC_{CHAOS} than with all other methods, it is still much lower than some correlations between other models, e.g., 0.61 between DM_{mean} and MA. Thus, we conclude that GPT-4 does not solely rely on information from ChaosNLI but achieves its strong performance via some other mechanism. Another possible explanation is that it is the only model next to GPT-3.5 that has access to explanations. In future, we would like to investigate the further use of explanations.

Figure 3 allows for a more general interesting observation. There seems to be a clear cluster structure in which the training-dynamics-based models (DM and MA) correlate highly with each other and GPT-4 clusters with the human scorers. Notably, the correlation between the clusters is small to nonexistent and even negative in some cases.

6.3 Influence of Human Label Variation

In which situations do AED methods make mistakes, e.g., detect false positive errors? This is an open question. We hypothesize that many topranking labels would either be errors or come from instances displaying HLV, i.e., instances with multiple labels after self-validation. The rest should be instances with just one plausible label. To test this hypothesis, we compute the proportion of *erroneous labels* vs. *valid labels from HLV instances* vs. *other* (neither errors nor HLV labels, i.e., with one plausible label) for the top 100 ranking labels for each method. The results in Figure 4 confirm our



Figure 4: Average distribution of erroneous, HLV, and other labels over the top 100 instances per method.

586

587

588

589

590

591

593

594

595

597

598

599

600

601

602

603

604

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

hypothesis (for the GPTs and human predictors): they place very few (0-11) labels that are neither errors nor HLV in the top 100. On the other hand, the training-dynamics-based methods MA and DM assign between 17.6 and 29.8 of these items to the top 100. This suggests that increasing the separation between errors and HLV is only one part of improving training dynamics methods for AED. Another could be finding the characteristics of the top-ranking items that are neither errors nor HLV.

6.4 Reranking models using label counts

Column *AP* (*rerank*) in Table 3 presents our reranking results. We observe that re-ranking improves over vanilla $LC_{VARIERR}$ for all methods but GPT-3.5. Interestingly, the best performing methods – also compared to the non-re-ranking approaches – are DM_{mean} and DM_{std} . They even perform better than Peer_{sum}, the best human approach. This suggests that combining statistics from multiple annotators with AED methods based on training dynamics is a promising future direction.

7 Conclusion

Errors exist in any dataset, but so does plausible human label variation. This paper defines a general procedure to separate error from plausible label variation. Our key idea to define errors is to leverage ecologically valid explanations (where annotators provide their reasons for a label) and pair these with annotator self-validation (to allow selfcorrection). We provide a new dataset VARIERR for the task. Our empirical investigation on VARI-ERR finds that traditional annotation error detection methods fare poorly on this task and underperform humans. While applied to NLI, our methodology is general, and we hope it inspires uptake.

Limitations

621

639

647

652

655

663

664

671

We believe that our two-round annotation setup would work for eliciting ecologically valid error 623 annotations in tasks or languages other than English NLI. However, we cannot be sure without 625 trying it, which we did not do during this project. Further, we did not use all types of information that VARIERR contains for the training-dynamics-based AED methods. An interesting question would be whether exploiting the soft label distribution with methods from learning from disagreement (Uma et al., 2021) would improve AED results. Another potentially useful source of information is the ex-633 planations given by the annotators. Using this information for computing the training dynamics or directly modeling whether an explanation makes 636 sense for a label in a supervised setting could po-637 tentially improve AED performance.

References

- Christoph Alt, Aleksandra Gabryszak, and Leonhard Hennig. 2020. TACRED Revisited: A Thorough Evaluation of the TACRED Relation Extraction Task. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1558– 1569, Online. Association for Computational Linguistics.
- Hadi Amiri, Timothy Miller, and Guergana Savova.
 2018. Spotting Spurious Data with Neural Networks. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 2006–2016, New Orleans, Louisiana. Association for Computational Linguistics.
- Eric Arazo, Diego Ortego, Paul Albert, Noel O'Connor, and Kevin Mcguinness. 2019. Unsupervised Label Noise Modeling and Loss Correction. In *Proceedings* of the 36th International Conference on Machine Learning, pages 312–321. PMLR.
- Lora Aroyo and Chris Welty. 2013. Crowd truth: Harnessing disagreement in crowdsourcing a relation extraction gold standard. *ACM Web Science 2013*.
- Lora Aroyo and Chris Welty. 2015. Truth is a lie: Crowd truth and the seven myths of human annotation. *AI Magazine*, 36(1):15–24.
 - Valerio Basile, Federico Cabitza, Andrea Campagner, and Michael Fell. 2021. Toward a perspectivist turn in ground truthing for predictive computing. *arXiv preprint arXiv:2109.04270*.
- Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. 2020.

Are we done with ImageNet? *arXiv preprint arXiv:2006.07159*.

- Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen tau Yih, and Yejin Choi. 2020. Abductive commonsense reasoning. In International Conference on Learning Representations.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Samuel R. Bowman and George Dahl. 2021. What will it take to fix benchmarking in natural language understanding? In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4843–4855, Online. Association for Computational Linguistics.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Aida Mostafazadeh Davani, Mark Díaz, and Vinodkumar Prabhakaran. 2022. Dealing with disagreements: Looking beyond the majority vote in subjective annotations. *Transactions of the Association for Computational Linguistics*, 10:92–110.
- Marie-Catherine de Marneffe, Christopher D. Manning, and Christopher Potts. 2012. Did it happen? The pragmatic complexity of veridicality assessment. *Computational Linguistics*, 38(2):301–333.
- Markus Dickinson and W. Detmar Meurers. 2003. Detecting Errors in Part-of-Speech Annotation. In *10th Conference of the European Chapter of the Association for Computational Linguistics*, Budapest, Hungary. Association for Computational Linguistics.
- Mitchell L Gordon, Michelle S Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S Bernstein. 2022. Jury learning: Integrating dissenting voices into machine learning models. In CHI Conference on Human Factors in Computing Systems, pages 1–19.

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

674 675 676

672

- 729 730
- 731 732
- 733
- 734 735
- 737
- 738 739 740 741
- 742 743 744

- 748 749 750 751 752
- 753 754 755 756 757

758 759

760 761

- 774 775 776
- 778

780 781

783

784

Jinchi Huang, Lie Qu, Rongfei Jia, and Binqiang Zhao. 2019. O2U-Net: A Simple Noisy Label Detection Approach for Deep Neural Networks. In 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pages 3325–3333.

- Nan-Jiang Jiang and Marie-Catherine de Marneffe. 2022. Investigating reasons for disagreement in natural language inference. *Transactions of the Association for Computational Linguistics*, 10:1357–1374.
- Nan-Jiang Jiang, Chenhao Tan, and Marie-Catherine de Marneffe. 2023. Ecologically valid explanations for label variation in NLI. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 10622–10633, Singapore. Association for Computational Linguistics.
- Jan-Christoph Klie, Bonnie Webber, and Iryna Gurevych. 2023. Annotation error detection: Analyzing the past and present for a more coherent future. *Computational Linguistics*, 49(1):157–198.
- Stefan Larson, Anish Mahendran, Andrew Lee, Jonathan K. Kummerfeld, Parker Hill, Michael A. Laurenzano, Johann Hauswald, Lingjia Tang, and Jason Mars. 2019. Outlier detection for improved data quality and diversity in dialog systems. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 517–527, Minneapolis, Minnesota. Association for Computational Linguistics.
- Christopher D Manning. 2011. Part-of-speech tagging from 97% to 100%: Is it time for some linguistics? In International conference on intelligent text processing and computational linguistics, pages 171–189. Springer.
- Mark Mazumder, Colby Banbury, Xiaozhe Yao, Bojan Karlaš, William Gaviria Rojas, Sudnya Diamos, Greg Diamos, Lynn He, Alicia Parrish, Hannah Rose Kirk, Jessica Quaye, Charvi Rastogi, Douwe Kiela, David Jurado, David Kanter, Rafael Mosquera, Juan Ciro, Lora Aroyo, Bilge Acun, Lingjiao Chen, Mehul Smriti Raje, Max Bartolo, Sabri Eyuboglu, Amirata Ghorbani, Emmett Goodman, Oana Inel, Tariq Kane, Christine R. Kirkpatrick, Tzu-Sheng Kuo, Jonas Mueller, Tristan Thrush, Joaquin Vanschoren, Margaret Warren, Adina Williams, Serena Yeung, Newsha Ardalani, Praveen Paritosh, Lilith Bat-Leah, Ce Zhang, James Zou, Carole-Jean Wu, Cody Coleman, Andrew Ng, Peter Mattson, and Vijay Janapa Reddi. 2023. DataPerf: Benchmarks for data-centric AI development.
- Mohammad Motamedi, Nikolay Sakharnykh, and Tim Kaldewey. 2021. A data-centric approach for training deep neural networks with less data.
- Yixin Nie, Xiang Zhou, and Mohit Bansal. 2020. What can we learn from collective human opinions on natural language inference data? In *Proceedings of the*

2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 9131–9143, Online. Association for Computational Linguistics. 786

787

789

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

- Rodrigo Frassetto Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-stage document ranking with BERT. *CoRR*, abs/1910.14424.
- Curtis G Northcutt, Anish Athalye, and Jonas Mueller. 2021. Pervasive label errors in test sets destabilize machine learning benchmarks. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1).*
- OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.
- Ellie Pavlick and Tom Kwiatkowski. 2019. Inherent disagreements in human textual inferences. *Transactions of the Association for Computational Linguistics*, 7:677–694.
- Joshua C. Peterson, Ruairidh M. Battleday, Thomas L. Griffiths, and Olga Russakovsky. 2019. Human uncertainty makes classification more robust. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pages 9616–9625.
- Barbara Plank. 2022. The "problem" of human label variation: On ground truth in data, modeling and evaluation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10671–10682, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Barbara Plank, Dirk Hovy, and Anders Søgaard. 2014. Linguistically debatable or just plain wrong? In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 507–511.
- Frederick Reiss, Hong Xu, Bryan Cutler, Karthik Muthuraman, and Zachary Eichenberger. 2020. Identifying Incorrect Labels in the CoNLL-2003 Corpus. In Proceedings of the 24th Conference on Computational Natural Language Learning, pages 215–226, Online. Association for Computational Linguistics.
- Paul Rottger, Bertie Vidgen, Dirk Hovy, and Janet Pierrehumbert. 2022. Two contrasting data annotation paradigms for subjective NLP tasks. In *Proceedings* of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 175–190, Seattle, United States. Association for Computational Linguistics.
- Susanna Rücker and Alan Akbik. 2023. CleanCoNLL: A nearly noise-free named entity recognition dataset. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8628–8645, Singapore. Association for Computational Linguistics.

Victor Sanh, Lysandre Debut, Julien Chaumond, and

Lars Schmarje, Vasco Grossmann, Tim Michels, Jakob

Nazarenus, Monty Santarossa, Claudius Zelenka, and

Reinhard Koch. 2023. Label smarter, not harder:

CleverLabel for faster annotation of ambiguous im-

age classification with higher quality. arXiv preprint

Shoaib Ahmed Siddiqui, Nitarshan Rajkumar, Tegan

Maharaj, David Krueger, and Sara Hooker. 2023.

Metadata archaeology: Unearthing data subsets by

leveraging training dynamics. In The Eleventh In-

ternational Conference on Learning Representations,

ICLR 2023, Kigali, Rwanda, May 1-5, 2023. Open-

Pia Sommerauer, Antske Fokkens, and Piek Vossen.

2020. Would you describe a leopard as yellow? Eval-

uating crowd-annotations with justified and informa-

tive disagreement. In Proceedings of the 28th Inter-

national Conference on Computational Linguistics,

pages 4798-4809, Barcelona, Spain (Online). Inter-

national Committee on Computational Linguistics.

Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie,

Yizhong Wang, Hannaneh Hajishirzi, Noah A. Smith,

and Yejin Choi. 2020. Dataset cartography: Mapping

and diagnosing datasets with training dynamics. In

Proceedings of the 2020 Conference on Empirical

Methods in Natural Language Processing (EMNLP), pages 9275–9293, Online. Association for Computa-

Alexandra N. Uma, Tommaso Fornaciari, Dirk Hovy,

Silviu Paun, Barbara Plank, and Massimo Poesio. 2021. Learning from Disagreement: A Survey.

Journal of Artificial Intelligence Research, 72:1385-

Vijay Vasudevan, Benjamin Caine, Raphael Gontijo-Lopes, Sara Fridovich-Keil, and Rebecca Roelofs. 2022. When does dough become a bagel? Analyzing

the remaining mistakes on ImageNet. arXiv preprint

Zihan Wang, Jingbo Shang, Liyuan Liu, Lihao Lu, Jiacheng Liu, and Jiawei Han. 2019. CrossWeigh: Training Named Entity Tagger from Imperfect Annotations. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natu-

ral Language Processing (EMNLP-IJCNLP), pages 5154–5163, Hong Kong, China. Association for Com-

Leon Weber and Barbara Plank. 2023. ActiveAED: A human in the loop improves annotation error detection. In *Findings of the Association for Compu*-

tational Linguistics: ACL 2023, pages 8834-8845,

Toronto, Canada. Association for Computational Lin-

abs/1910.01108.

arXiv:2305.12811.

Review.net.

tional Linguistics.

arXiv:2205.04596.

putational Linguistics.

guistics.

1470.

Thomas Wolf. 2019. DistilBERT, a distilled version

of BERT: smaller, faster, cheaper and lighter. ArXiv,

- 841 842
- 04 8/
- 84
- 84 97
- 84
- 8
- 8
- 8
- 855
- 8 8
- 858
- 8
- 861
- 8
- 8

86

86

871

872 873

874 875

876

878 879

8

8

886 887

8

890 891

892 893

8

894

Adina Williams, Nikita Nangia, and Samuel Bowman.
2018. A broad-coverage challenge corpus for sentence understanding through inference. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguis-

tics.

896

897

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

- Shujian Zhang, Chengyue Gong, and Eunsol Choi. 2021. Learning with different amounts of annotation: From zero to many labels. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7620–7632, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xinliang Frederick Zhang and Marie-Catherine de Marneffe. 2021. Identifying inherent disagreement in natural language inference. In *Proceedings* of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4908–4915, Online. Association for Computational Linguistics.
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Jeff Huang, Chuyue Sun, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. 2023. Efficiently programming large language models using sglang. *arXiv preprint arXiv:2312.07104*.

A Data Statistics

A.1 Pair-wise inter-annotator agreements (Cohen's kappa) on MASI-distance for non-validated, self-validated, and peer-validated versions

versions \ annotators	1-vs-2	1-vs-3	1-vs-4	2-vs-3	2-vs-4	3-vs-4
before validation	40.87	42.22	37.71	36.21	31.17	34.62
self-validated	60.06	53.84	61.61	44.47	47.79	47.06
peer-validated	66.09	72.03	67.64	64.07	68.05	68.78

Table 4: Pair-wise inter-annotator agreements (Co-hen's kappa) on MASI-distance for non-validated, self-validated, and peer-validated versions.

A.2 Frequency of NLI label on non-validation, self-validated, and peer-validated explanation-label pairs



Figure 5: Frequency of NLI label sets on non-, selfand peer-validated label-explanation pairs.

B GPT Prompt

id: 72870c

System: You are an expert linguistic annotator. User: We have collected annotations for an NLI instance together with reasons for the labels. Your task is to judge whether the reasons make sense for the label. Provide the probability (0.0 - 1.0) that the reason makes sense for the label. Give ONLY the reason and the probability, no other words or explanation. For example: Reason: < The verbatim copy of the reason> Probability: <the probability between 0.0 and 1.0 that the reason makes sense for the label, without any extra commentary whatsoever; just the probability!>. Context: Because marginal costs are very low, a newspaper price for preprints might be as low as 5 or 6 cents per piece. Statement: Newspaper preprints can cost as much as \$5. Reason for label entailment: 5 dollars for a piece of newspaper Reason for label neutral: The context only mentions how low the price may be, not how high it may be. Reason for label neutral: The maximum cost of newspaper preprints is not given in the context. Reason for label contradiction: The context says 5 or 6 cents, not \$5. User: Reason: 5 dollars for a piece of newspaper Probability: Assistant: 0.0 User: Reason: The context only mentions how low the price may be, not how high it may be. Probability: Assistant: 0.9 User: Reason: The maximum cost of newspaper preprints is not given in the context. Probability: Assistant: 0.8 User: Reason: The context says 5 or 6 cents, not \$5. Probability: Assistant: 0.9

Figure 6: GPT Prompt for predicting likelihood probability of label-explanation pairs.

C More VARIERR Examples

934

930

926

 $\ensuremath{\textit{Premise}}\xspace$ Students of human misery can savor its underlying sadness and futility.

 $\mathit{Hypothesis}:$ Students of human misery will be delighted to see how sad it truly is.

Label-explanation pairs: before validation: {E:1,N:2,C:1} Self-validated : {E:1,N:1} Peer-validated: {N:1} Labels: [E, N] Error: [C]

		Round 1: NLI Label & Explanation	Ro	und 2	: Valid	ity
L	А	Explanation	1	2	3	4
Е	2	"can savor" implies "will be delighted".	1	✓	×	×
N	1	It is not clear from the context if the students will be delighted.	×	Х	✓	1
1	3	Students of human misery can "savored" that sadness, so maybe they are delighted	×	\sim		1
	5	to see that, maybe they are tortured by the disasters.		~	v	•
С	4	Savor means to understand. Not to enjoy.	×	×	?	×

(a) *id*: 116176c

Premise: The tree-lined avenue extends less than three blocks to the sea.

 $\mathit{Hypothesis}{:}$ The sea isn't even three blocks away.

 $\label{explanation pairs: before validation: {"E":4, "N":1, "C":1} Self-validated: {"E":3, "N":1} Peer-validated: {"E":4, "N":1} Labels: [E, N] Error: [C]$

		Round 1: NLI Label & Explanation				dity
L	А	Explanation	1	2	3	4
Е	1	Both premise and hypothesis talk about less than three blocks.	1	1	1	×
	2	If the avenue reaches the sea after less than three blocks, it cannot be further away.	1	✓	1	×
	3	The avenue is less than three blocks to the sea.	1	1	✓	×
	4	If the hypothesis means that the sea is less than three blocks away.	?	1	1	×
Ν	3	It is not given where is the location of the narrator.	1	×	 Image: A set of the set of the	1
С	4	If the hypothesis means that the sea is more than three blocks away.	?	×	?	×

(b) *id:* 80630e

Premise: As he stepped across the threshold, Tommy brought the picture down with terrific force on his head. *Hypothesis*: Tommy hurt his head bringing the picture down.

Label-explanation pairs: before validation: {"E":3,"N":1} Self-validated: {"E":3,"N":1} Peer-validated: {"E":3,"N":1} Labels: [E, N] Error: [C]

		Round 1: NLI Label & Explanation			Round 2: Validity				
L	Α	Explanation	1	2	3	4			
	1	the picture hit Tommy in the head	1	1	1	Х			
Е	2	a picture hit Tommy's head with terrific force	1	 Image: A set of the set of the	1	×			
	3	Tommy hurt his head with the picture	1	1	✓	×			
Ν	3	ambiguous if Tommy hurt himself or another guy	1	1	 Image: A set of the set of the	×			
С	4	Tommy is not hurt but rather bad strong emotion	×	Х	1	×			

(c) *id*: 77893n

 Table 5: Additional sample annotations from VARIERR NLI corpus. L: Label, A: Annotator; E: Entailment, N:

 Neutral, C: Contradiction; magenta: self-judgments, black: peer-judgments, Err: label error.