Modelling Variability in Human Annotator Simulation

Anonymous ACL submission

Abstract

Human annotator simulation (HAS) serves as a cost-effective substitute for human evaluation tasks such as data annotation and system assessment. It is important to incorporate the variabil-005 ity present in human evaluation into HAS, since it helps capture diverse subjective interpretations and mitigate potential biases and over-007 representation. This work introduces a novel framework for modelling variability in HAS. Conditional softmax flow (S-CNF) is proposed to model the distribution of subjective human 011 annotations, which leverages diverse human annotations via meta-learning. This enables the efficient generation of annotations that exhibit human variability for unlabelled input. In addition, a wide range of evaluation metrics are adopted to assess the capability and efficiency 017 018 of HAS systems in predicting the aggregated behaviours of human annotators, matching the distribution of human annotations, and simulating the inter-annotator disagreements. Results demonstrate that the proposed method achieves state-of-the-art performance on two real-world human evaluation tasks: emotion recognition and toxic speech detection¹.

1 Introduction

026

027

028

039

Human evaluation is fundamental to machine learning research. It guides processes such as data annotation and model assessment, including for instance perceptual quality evaluation (Ma et al., 2015; Talebi and Milanfar, 2018; Ramesh and Sanampudi, 2022), annotation generation for weak supervision (Ratner et al., 2016; Wu et al., 2022a), and model optimization based on human preference (Schatzmann et al., 2007; Gür et al., 2018; Ruiz et al., 2019; Shi et al., 2019; Lin et al., 2021). Collecting human annotations or evaluations often requires substantial resources and may expose human annotators to distressing and harmful content in sensitive tasks (*e.g.*, toxic speech detection, suicidal risk prediction, and depression detection). This inspires the exploration of human annotator simulation (HAS) as a scalable and cost-effective alternative, which facilitates large-scale dataset evaluation, benchmarking, and system comparisons.

Variability is a unique aspect of real-world human evaluation. Individual variations in cognitive biases, cultural backgrounds, and personal experiences (Hirschberg et al., 2003; Wiebe et al., 2004; Haselton et al., 2015) can lead to variability in human interpretation (Maniati et al., 2022). It has been argued that achieving a single deterministic "ground truth" in subjective tasks like human evaluation is not feasible, nor essential (Alm, 2011; Wu et al., 2022b). Therefore, HAS should incorporate the variability present in human evaluation rather than solely relying on majority opinions. This mitigates potential biases and over-representation in scenarios where dominant opinions could potentially overshadow minority viewpoints, thus promoting fairness and inclusivity.

This work investigates modelling subjective human annotation distributions and simulating human-like annotations. We propose a novel framework which formulates HAS as a zero-shot density estimation problem. A new model, conditional softmax flows (S-CNFs), is proposed which leverages diverse human annotations via meta-learning. This enables the efficient generation of annotations that exhibit human variability for unlabelled input. To the best of our knowledge, this is the first work that incorporates human variability into HAS without requiring human annotators to be dynamically involved in the process, while remaining scalable for large crowd-sourced datasets. Moreover, a range of evaluation metrics are adopted to assess the capability and efficiency of HAS systems regarding prediction of the majority opinion, estimation of the human annotation distribution, and simulation of the inter-annotator disagreements. The proposed

079

040

041

042

045

046

047

¹Code will be available upon acceptance.

084

- 80
- 087
- 089 090
- 091

0

0

100 101

102 103 104

105 106 107

108 109

> 110 111

112 113 114

115

116 117

118

119 120

121

122

123

124

125

126

127

129

approach is evaluated on two real-world applications: emotion recognition and toxic speech detection. Results demonstrate that the proposed method achieves state-of-the-art performance on modelling variability in HAS.

2 Human Annotator Simulation (HAS)

2.1 The Variability in Human Evaluation is Valuable

Each individual's perception of the world is unique and influenced by their physical state and cognitive biases. This leads to diverse and subjective interpretations. Such subjectivity can be manifest in various tasks such as emotion recognition (Hirschberg et al., 2003; Mihalcea and Liu, 2006), perceptual quality assessment (Wiebe et al., 2004; Seshadrinathan et al., 2010), and user experience evaluation (Zen and Vanderdonckt, 2016). Rather than seeking to reduce the variability in annotations, it is important to account for annotators' subjective interpretations when designing a human annotator simulator. The importance of variability in HAS can be demonstrated by the following examples:

Revealing data ambiguity. Incorporating the variability in human perception empowers HAS to reveal potential ambiguity or complexity in data, providing valuable insights for further analysis.

Mitigating bias and over-representation. Incorporating the variability in human judgements prevents HAS from being biased towards a certain perspective and ignoring minority viewpoints, leading to a more inclusive representation of opinions where all viewpoints are given due consideration (Dixon et al., 2018; Hutchinson et al., 2020).

Improving model alignment. Optimization based on human feedback has led to superior performance on tasks such as text generation (Christiano et al., 2017; Ouyang et al., 2022; Rafailov et al., 2023), which aligns the behaviour of language models with human preferences. HAS could be helpful in this task, as it is an efficient and cost-effective alternative to generating human feedback.

2.2 Problem Formulation

Denote an event as d_i , which consists of a descriptor (*e.g.*, an utterance or text) x_i and a set of M_i human annotations $\mathcal{D}_i = \{\eta_i^{(m)}\}_{m=1}^{M_i}$ for x_i . Note that different events may be labelled by different sets of annotators. Given a dataset of training events $\mathcal{D} = \{(x_i, \mathcal{D}_i)\}_{i=1}^N$, HAS aims to model the conditional annotation distribution $p(\eta_i | x_i)$ given

the observations \mathcal{D}_i of η_i provided by different annotators. For a unseen test descriptor \boldsymbol{x}_* , HAS can then predict $p(\boldsymbol{\eta}_*|\boldsymbol{x}_*)$ to simulate human-like annotations $\mathcal{D}_* = \{\boldsymbol{\eta}_*^{(m)}\}_{m=1}^{M_*}$ in a way that reflects how it would be labelled by human annotators.

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

155

156

157

158

159

160

161

162

163

164

165

166

167

168

170

171

172

173

174

175

176

177

179

2.3 Related Work

Prior work mainly investigated three approaches to simulating human annotations.

The first approach uses a single proxy variable η'_i (e.g., majority vote) to summarize all annotations for each descriptor x_i (Kim et al., 2013; Djuric et al., 2015; Patton et al., 2016; Poria et al., 2017). This creates a proxy dataset $\mathcal{D}' = \{(\boldsymbol{x}_i, \boldsymbol{\eta}'_i)\}_{i=1}^N$ and converts HAS into a supervised learning problem, which is usually solved by fitting a discriminative model to estimate the conditional distribution for the proxy variable. During testing, given an unseen descriptor x_* , the model predicts the proxy variable η'_* for x_* . Clearly, modelling a single proxy variable as in this approach fails to take into account the subjectivity and diversity in human behaviour and perception. Other work incorporated the variance of human annotations into the proxy variable (Deng et al., 2012; Prabhakaran et al., 2012; Plank et al., 2014; Dang et al., 2017; Han et al., 2017; Leng et al., 2021). However, all these approaches still focus on obtaining the "correct" label (*e.g.*, aiming for improved prediction accuracy) and minimizing the discrepancy among annotators (e.g., reducing "noise" in annotations) rather than embracing inter-annotator disagreements.

The second approach explicitly models the behaviours of different annotators using different individual models in an ensemble or different heads in a single model (Fayek et al., 2016; Chou and Lee, 2019; Davani et al., 2022). This approach is computationally feasible only when the number of annotators is relatively small and when a sufficient quantity of annotation is available for each annotator, which is not applicable to large crowd-sourced datasets (Lotfian and Busso, 2019; Mathew et al., 2021) that are common in real-world applications.

The third approach approximates subjective probability distributions using Markov chain Monte Carlo with people (Sanborn and Griffiths, 2007; Harrison et al., 2020), which requires human annotators to be involved in the process in a dynamic setting. These methods present the descriptor x_* to human participants and asks them to provide a sequence of decisions \mathcal{D}_* following



Figure 1: Diagram for the proposed zero-shot human annotator simulation framework.

the Metropolis-Hasting acceptance rule. The annotation distribution $p(\eta_*|x_*)$ is then estimated based on \mathcal{D}_* . In other words, this requires access to human annotations \mathcal{D}_* for estimating the annotation distribution for each x_* , and there is no obvious way to transfer information between different events. Therefore, these methods cannot be applied to simulate annotation distributions for unlabelled test descriptors.

181

182

183

184

185

189

190

191

192

194

195

196

197

198

199

201

203

204

207

210

211

212

213

214

215

3 A Meta-learning Framework for HAS

This paper proposes a novel framework for HAS that meta-learns a conditional softmax flow (S-CNF) to estimate the human annotation distribution $p(\boldsymbol{\eta}|\boldsymbol{x})$ across all training events \mathcal{D} . The proposed model learns to learn (i.e., meta-learns) how to estimate the underlying distribution of human annotations \mathcal{D}_i for any given descriptor x_i by leveraging the diverse human annotations, rather than designing a proxy variable to summarize \mathcal{D}_i as in the first approach described in Sec. 2.3. Unlike the second approach in Sec. 2.3 which separately models each individual human annotator with a different model, our method is compatible with large crowd-sourced datasets since it amortizes across annotators with a single S-CNF model. Moreover, our model is a zero-shot human annotation simulator which can estimate the human annotation distribution $p(\eta_*|x_*)$ for any unseen test descriptor x_* without access to any human annotations \mathcal{D}_* for x_* , in contrast to the third method in Sec. 2.3 which requires human annotators to be dynamically involved in the process of labelling x_* .

3.1 A Latent Variable Model for HAS

The proposed framework for HAS is realized by a latent variable model²:

$$p_{\theta}(\boldsymbol{y}|\boldsymbol{x}) = \iint p(\boldsymbol{y}|\boldsymbol{v}) p_{\phi}(\boldsymbol{v}|\boldsymbol{z}) p_{\Lambda}(\boldsymbol{z}|\boldsymbol{x}) d\boldsymbol{v} d\boldsymbol{z},$$
(1)

where the conditional prior $p_{\Lambda}(\boldsymbol{z}|\boldsymbol{x})$ learns to summarize useful information about the input descriptor \boldsymbol{x} and encode the possible disagreements over \boldsymbol{x} among different human annotators, which is helpful for the likelihood $p_{\phi}(\boldsymbol{y}|\boldsymbol{z}) = \int p(\boldsymbol{y}|\boldsymbol{v})p_{\phi}(\boldsymbol{v}|\boldsymbol{z})d\boldsymbol{v}$ to simulate human-like annotations for \boldsymbol{x} .

216

217

218

219

220

221

222

223

224

225

229

230

232

233

234

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

The proposed zero-shot human annotator simulator is illustrated in Figure 1. Specifically, the conditional prior is modelled by a conditional factorized Gaussian distribution $p_{\Lambda}(z|x) = \mathcal{N}(z|\mu_{\Lambda}(x), \operatorname{diag}(\sigma_{\Lambda}^2(x)))$ whose mean $\mu_{\Lambda}(x)$ and variance $\sigma_{\Lambda}^2(x)$ are parameterized by a neural network g_{Λ} with parameters Λ .

The intermediate variable v is obtained by a deterministic invertible transformation $p_{\phi}(v|z) = \delta(v - f_{\phi}(z))$, where $f_{\phi}(z)$ is parameterized by an invertible neural network with parameters ϕ , and $\delta(\cdot)$ is the multivariate Dirac delta function. This results in a conditional normalizing flow (CNF):

$$p_{\theta}(\boldsymbol{v}|\boldsymbol{x}) = \int \delta(\boldsymbol{v} - \boldsymbol{f}_{\phi}(\boldsymbol{z})) p_{\Lambda}(\boldsymbol{z}|\boldsymbol{x}) d\boldsymbol{z}$$

= $p_{\Lambda} \left(\boldsymbol{f}_{\phi}^{-1}(\boldsymbol{v}) \middle| \boldsymbol{x} \right) \left| \det \left(\frac{\partial \boldsymbol{f}_{\phi}^{-1}(\boldsymbol{v})}{\partial \boldsymbol{v}} \right) \right|,$ (2)

where det(·) denotes the determinant operator, $\partial f_{\phi}^{-1}(v)/\partial v$ denotes the Jacobian matrix of $f_{\phi}^{-1}(v)$, and $\theta := \{\phi, \Lambda\}$ denotes all parameters in this base CNF. This modelling choice has the advantage of having tractable marginal likelihood as in Eqn. (2) while not restricting the intermediate variable v to a specific type of distribution as in previous methods (*e.g.*, Gaussian (Han et al., 2017) and Student-t (Wu et al., 2023)), thus offering enhanced tractability, flexibility and generality. In addition, samples can be efficiently drawn from this model by first drawing $z \sim p_{\Lambda}(z|x)$ from the conditional prior and then computing the deterministic flow transformation $v = f_{\phi}(z)$.

Finally, the output variable y is obtained by augmenting the intermediate variable v with the transformation p(y|v). For continuous annotations, the identity transformation $p(y|v) = \delta(y - v)$ is used.

²For clarity, we use different notations for human annotations η and model outputs y.



Figure 2: Illustration of the S-CNF workflow. g_{Λ} is the feature encoder, f_{ϕ} is an invertible neural network, h_{Ω} is the variational encoder.

However, real-world human evaluation tasks often involve discrete annotations (*e.g.*, human annotators are usually instructed to choose from a predefined set of options). In the next section, a new model class is introduced to accommodate discrete annotations.

254

255

257

258

261

262

263

270

271

272

273

274

3.2 Conditional Softmax Flows (S-CNFs)

We propose a new model class called conditional softmax flow (S-CNF) to accommodate discrete annotations. The workflow of S-CNF is illustrated in Figure 2. S-CNFs augment the base CNFs by applying the softmax function $p(y|v) = \delta(y - \operatorname{softmax}(v))$ to transform the continuous intermediate variable v into categorical probabilities y. Let c be a categorical variable with probability $P(c = k|y) = y_k$ ($k = 1, \dots, K$), which represents the categorical annotation for an descriptor x, with $P(c = k|v) = \int y_k \delta(y - \operatorname{softmax}(v)) dy = \operatorname{softmax}(v)_k$. The marginal likelihood of S-CNF is given by

$$\mathbf{P}_{\boldsymbol{\theta}}(c=k|\boldsymbol{x}) = \int \mathbf{P}(c=k|\boldsymbol{v})\mathbf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x})\mathrm{d}\boldsymbol{v} \quad (3)$$

where $p_{\theta}(v|x)$ is the marginal likelihood of the base CNF defined in Eqn. (2). Since the marginal likelihood of the S-CNF given in Eqn. (3) is analytically intractable due to the softmax transformation, we propose to approximate it using variational inference (Wainwright et al., 2008) with a learnable mean-field Gaussian variational posterior $q_{\Omega}(\boldsymbol{v}|\boldsymbol{y}) = \mathcal{N}(\boldsymbol{v}|\boldsymbol{\mu}_{\Omega}(\boldsymbol{y}), \operatorname{diag}(\sigma_{\Omega}^{2}(\boldsymbol{y})))$, whose mean $\boldsymbol{\mu}_{\Omega}(\boldsymbol{y})$ and variance $\sigma_{\Omega}^{2}(\boldsymbol{y})$ are parameterized by a neural network \boldsymbol{h}_{Ω} with parameters Ω . This can be seen as a probabilistic inverse of the softmax transformation $p(\boldsymbol{y}|\boldsymbol{v})$. Applying Jensen's inequality to the log marginal likelihood of the S-CNF in Eqn. (3), we obtain a tractable evidence lower bound (ELBO):

281

283

287

290

291

292

294

295

296

297

298

299

300

301

302

303

304

306

307

308

309

310

311

312

313

314

315

316

317

319

321

322

323

324

325

$$\log \mathsf{P}_{\boldsymbol{\theta}}(c=k|\boldsymbol{x}) \geq \mathbb{E}_{\mathsf{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{y})}[\log \mathsf{P}(c=k|\boldsymbol{v}) + \log \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x}) - \log \mathsf{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{y})].$$
(4)

It is worth noting that the softmax flow likelihood $P(c = k | v) = \text{softmax}(v)_k$ places non-zero probability mass for every category $k = 1, \dots, K$, which is different from argmax flow (Hoogeboom et al., 2021) whose likelihood only places probability mass for a single category. From a modelling perspective, softmax flow has a better capacity to represent the variability and uncertainty in human annotations. From an optimization perspective, the ELBO for softmax flow is always well-defined, whereas the ELBO for argmax flow is not defined when the model output does not match the human annotation in which case the log-likelihood would be $\log(0)$ and requires additional thresholding tricks to fix (Hoogeboom et al., 2021).

3.3 A Meta-learning Objective for S-CNFs

Using the variational approximation defined in Eqn. (4), the loss $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Omega}; \boldsymbol{d}_i)$ for S-CNF on a single event \boldsymbol{d}_i can be defined as the average negative ELBO evaluated on the set of the human annotations $\mathcal{D}_i = \{\boldsymbol{\eta}_i^{(m)}\}_{m=1}^{M_i}$ for the corresponding descriptor $\boldsymbol{x}_i: \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Omega}; \boldsymbol{d}_i) =$

$$-rac{1}{M_i}\sum_{m=1}^{M_i}\mathbb{E}_{\mathtt{q}_{\mathbf{\Omega}}(oldsymbol{v}|oldsymbol{\eta}_i^{(m)})}\left[\sum_{k=1}^Koldsymbol{\eta}_{i,k}^{(m)}\log \mathtt{P}(c_i=k|oldsymbol{v})
ight.$$

$$+\log p_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x}_i) - \log q_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{\eta}_i^{(m)}) \Big], \qquad (5)$$

where the expectation over the variational posterior is approximated by Monte Carlo simulation with the reparameterization trick (Kingma and Welling, 2014). Following the episodic training scheme (Vinyals et al., 2016; Snell et al., 2017; Chen et al., 2023), we treat density estimation on each event as a learning problem and randomly sample a subset of such learning problems to train on at each step during meta-training. This results in a metalearning objective across the training events in D:

$$\mathcal{L}_{\text{meta}}(\boldsymbol{\theta}, \boldsymbol{\Omega}; \mathcal{D}) = \mathbb{E}_{\boldsymbol{d}_i \sim p(\mathcal{D})}[\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Omega}; \boldsymbol{d}_i)], \quad (6)$$

where p(D) denotes the uniform distribution over D. Intuitively, this objective maps each human annotation to the latent space of the corresponding descriptor by the S-CNF during meta-training, which helps the model to build a diverse latent representation that captures the variability in human annotations across different descriptors.

At test time, the S-CNF can simulate human-like annotations for an unseen, unlabelled descriptor \boldsymbol{x}_* by first drawing $\boldsymbol{v}_*^{(m)} \sim p_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x}_*)$ from the base CNF then applying the softmax function $\boldsymbol{y}_*^{(m)} =$ softmax $(\boldsymbol{v}_*^{(m)})$ for $m = 1, \cdots, M_*$, where M_* denotes the number of annotations to be simulated. Note that each sample of S-CNF is a categorical distribution with probabilities $\boldsymbol{y}_*^{(m)}$. More details can be found in Appendix D.

4 Experimental Setup

326

327

331

332

333

337

341

343

345

347

351

354

357

371

374

The proposed framework for variability-aware HAS is evaluated on two real-world applications for speech and natural language processing: emotion class labelling and toxic speech detection. A wide range of evaluation metrics are adopted to assess the performance of HAS systems.

4.1 Evaluation Tasks and Datasets

Emotion class labelling. The highly subjective perception of emotion often results in disagreements among human annotators. Most emotion datasets employ multiple annotators to label each utterance while prior works typically use the majority vote as the ground-truth (Busso et al., 2008; Poria et al., 2019; Wu et al., 2021). A variability-aware HAS system better handles different opinions among human annotators and enhances the fairness of emotion class annotation. MSP-Podcast (Lotfian and Busso, 2019) is one of the largest publicly available datasets in speech emotion recognition, which contains natural English speech from podcast recordings annotated using crowd-sourcing. Each utterance has 6.7 annotations on average. Release 1.6 was used in our experiments, which contains 50k+ utterances from 1k+ speakers. The standard splits of training, validation and test were used. Emotion labels were grouped into five categories: angry, sad, happy, neutral, and other. 16.5% of the utterances do not have a majority agreed emotion class.

Toxic speech detection aims to filter out harmful language, which is crucial for respectful online environments and healthy communications among users. A variability-aware HAS system accounts for comprehensive understanding of hate speech, which is a good substitute for human annotators to reduce their exposure to distressing and harmful content. HateXplain (Mathew et al., 2021) was used in our experiments, which is a popular dataset for toxic speech detection, which contains 20k+ text posts from Twitter and Gab. These posts are labelled using crowd-sourcing with the commonly used three-category annotation: hate, offensive, normal. Each post is annotated by three annotators. Cases where all three annotators choose a different class (919 out of 20,148 posts) were originally excluded from the standard split of the dataset. We incorporate these cases into our training, validation, and test sets in an 8:1:1 ratio to better reflect the inter-annotator disagreements, resulting in 16,118 posts for training, 2,014 for validation, and 2,016 for testing.

375

376

377

378

379

380

381

382

384

386

387

391

392

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

4.2 Evaluation Metrics

A range of metrics are adopted to assess the empirical performance of HAS systems in terms of majority prediction, distribution matching, and human variability simulation.

Majority prediction. Classification accuracy (Acc) for the majority vote is evaluated for all test inputs with majority-agreed human annotations.

Distribution matching. Negative log likelihood (NLL) is used to evaluate how well the model estimates the human annotation distribution: NLL^{all} = $-\frac{1}{N}\sum_{i=1}^{N}(\frac{1}{M_i}\sum_{m=1}^{M_i}\log p_{\theta}(\boldsymbol{\eta}_i^{(m)}|\boldsymbol{x}_i)).$

Inter-annotator disagreement simulation. Apart from evaluating the goodness of fit, additional metrics are adopted to explicitly measure how well the model simulates the variability and disagreements in human annotations: (i) the root mean squared error of the standard deviations of the annotations for all test inputs: $RMSE^{s} = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (\sigma_{i} - s_{i})^{2}}, \text{ where } \sigma_{i} = \frac{1}{K} \sum_{k=1}^{K} \sqrt{\frac{1}{M_{i}} \sum_{m=1}^{M_{i}} (\eta_{i,k}^{(m)} - \bar{\eta}_{i,k})^{2}} \text{ and } s_{i} =$ $RMSE^{s}$ $\frac{1}{K}\sum_{k=1}^{K}\sqrt{\frac{1}{M_{*}}\sum_{m=1}^{M_{*}}(y_{i,k}^{(m)}-\bar{y}_{i,k})^{2}};$ (ii) the absolute error of the average standard deviations of the annotations for all test inputs: $\mathcal{E}(\bar{s}) = |\bar{\sigma} - \bar{s}|,$ where $\bar{\sigma} = \sum_{i=1}^{N} \sigma_i$ and $\bar{s} = \sum_{i=1}^{N} s_i$; (iii) Fleiss's kappa (κ) (Fleiss, 1971), a real number between -1and +1, with -1 indicating no observed agreement and +1 indicating perfect agreement. The absolute error between the kappas of human annotations (κ) and simulated annotations $(\hat{\kappa})$ for all test inputs is reported: $\mathcal{E}(\hat{\kappa}) = |\hat{\kappa} - \kappa|$.

4.3 Baselines

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

The proposed S-CNF method is compared to baselines of various types such as ensemble methods, Bayesian methods, and conditional generative models. This includes deep ensemble (Ensemble) (Lakshminarayanan et al., 2017), Monte-Carlo dropout (MCDP) (Gal and Ghahramani, 2016), Bayes-bybackprop (BBB) (Blundell et al., 2015), conditional variational autoencoder (CVAE) (Kingma and Welling, 2014), conditional argmax flow (A-CNF) (Hoogeboom et al., 2021), and Dirichlet prior network (**DPN**) (Malinin and Gales, 2018; Wu et al., 2022b). We fit each method to all available human annotations for all utterances in the training set, tune hyperparameters on the validation set, and report performance on the test set. $M_* = 100$ samples are used to compute evaluation metrics at test time. Ensemble only consists of 10 systems due to its expensive computational cost. Details about the configuration of the baseline models can be found in Appendix E.

4.4 Backbone Architecture

All compared methods use the same upstreamdownstream feature encoder g_{Λ} to extract features from descriptors. The upstream model (Bommasani et al., 2021) is pre-trained on a large amount of unlabelled data to learn universal representations. WavLM (Chen et al., 2022) and RoBERTa (Liu et al., 2019) are used as the pre-trained upstream models for speech and text descriptors respectively. The downstream model consists of two Transformer encoder blocks followed by two fully connected (FC) layers, which are fine-tuned to target specific applications. The invertible flow model f_{θ} uses three real NVP blocks (Dinh et al., 2017) and the variational encoder h_{Ω} contains an FC layer. Details about the structure and implementation can be found in Appendix D and Appendix F.

5 Results and Analysis

The proposed method is evaluated according to the 463 setup in Sec. 4. The evaluation results along with 464 case study of representative examples demonstrates 465 the superior capability of the proposed method in 466 capturing the aggregated behaviours of human an-467 468 notators, matching the distribution of human annotations, and simulating the variability of human 469 interpretation. For each metric, mean value with 470 standard error over three independent runs are re-471 ported for all methods. 472

Emotion	Acc (†)	$\mathbf{NLL^{all}}\;(\downarrow)$
MCDP	$0.582{\pm}0.003$	1.423 ± 0.012
Ensemble	$0.603{\pm}0.002$	$1.458 {\pm} 0.004$
BBB	$0.565 {\pm} 0.010$	$1.459{\pm}0.011$
DPN	$0.581{\pm}0.006$	$1.459{\pm}0.011$
CVAE	$0.275 {\pm} 0.000$	$1.661 {\pm} 0.000$
A-CNF	$0.583{\pm}0.002$	$1.430{\pm}0.006$
S-CNF	$\underline{0.591{\pm}0.002}$	$1.403{\pm}0.011$
Toxic	Acc (†)	$\mathbf{NLL^{all}}\left(\downarrow ight)$
Toxic MCDP	Acc (↑) 0.656±0.009	NLL^{all} (↓) 0.951±0.032
Toxic MCDP Ensemble	Acc (↑) 0.656±0.009 0.682±0.002	NLL ^{all} (↓) 0.951±0.032 0.909±0.012
Toxic MCDP Ensemble BBB	Acc (†) 0.656±0.009 0.682±0.002 0.670±0.001	NLL ^{all} (↓) 0.951±0.032 0.909±0.012 0.670±0.001
Toxic MCDP Ensemble BBB DPN	Acc (†) 0.656±0.009 0.682±0.002 0.670±0.001 0.581±0.006	NLL ^{all} (\downarrow) 0.951 \pm 0.032 0.909 \pm 0.012 0.670 \pm 0.001 1.158 \pm 0.002
Toxic MCDP Ensemble BBB DPN CVAE	Acc (†) 0.656±0.009 0.682±0.002 0.670±0.001 0.581±0.006 0.406±0.000	NLL ^{all} (\downarrow) 0.951 \pm 0.032 0.909 \pm 0.012 0.670 \pm 0.001 1.158 \pm 0.002 1.150 \pm 0.000
Toxic MCDP Ensemble BBB DPN CVAE A-CNF	Acc (†) 0.656±0.009 0.682±0.002 0.670±0.001 0.581±0.006 0.406±0.000 0.628±0.003	NLL ^{all} (\downarrow) 0.951 \pm 0.032 0.909 \pm 0.012 0.670 \pm 0.001 1.158 \pm 0.002 1.150 \pm 0.000 0.892 \pm 0.011

Table 1: Comparison to the baselines in terms of majority prediction and distribution matching. CVAE collapses to one category for all inputs. The best value in each column is shown in bold and the second best is underlined.

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

5.1 Performance

Table 1 compares all methods in terms of majority prediction and distribution matching. The Ensemble achieves the best majority prediction accuracy (Acc) at the cost of training 10 independent systems. The proposed S-CNF achieves the secondbest majority prediction accuracy with only a tenth of the computational cost of Ensemble during training and testing. Despite this, we stress that achieving the highest accuracy is not the goal of HAS since accuracy only measures the majority prediction performance and ignores variability in the annotations. More importantly, S-CNF is the best at matching the distributions of human annotations (in terms of NLL^{all}) among all compared methods. Table 2 reports the test results for all compared methods regarding inter-annotator disagreement annotator simulation. S-CNF again outperforms all compared methods in modelling the variability in human annotations, evident by the smallest RMSE^{*s*}, $\mathcal{E}(\bar{s})$ and $\mathcal{E}(\hat{\kappa})$.

5.2 Computational Time Cost

The computational time cost of all of compared495methods for the two tasks studied in the paper are496shown in Table 3. Denote M_* as the number of497annotations to be simulated. The Ensemble model498with M_* members involves training and testing M_* 499individual models, which costs $M_* \times$ more training500

	Eme	otion class labe	lling	Toxic speech detection			
	$\mathbf{RMSE}^{s}\left(\downarrow\right)$	${\cal E}(ar s)~(\downarrow)$	${\cal E}(\hat\kappa)~(\downarrow)$	$\mathbf{RMSE}^{s}\left(\downarrow\right)$	${\cal E}(ar s)~(\downarrow)$	$\mathcal{E}(\hat{\kappa})~(\downarrow)$	
MCDP	0.294±0.001	$0.193 {\pm} 0.000$	$0.467 {\pm} 0.005$	0.300 ± 0.002	$0.129{\pm}0.003$	$0.143 {\pm} 0.008$	
Ensemble	0.271±0.003	$0.160{\pm}0.004$	$0.344{\pm}0.017$	0.289 ± 0.001	$0.100{\pm}0.003$	0.064 ± 0.006	
BBB	0.289 ± 0.005	$0.187 {\pm} 0.008$	$0.511 {\pm} 0.034$	$0.949 {\pm} 0.021$	$0.300{\pm}0.009$	$0.127 {\pm} 0.022$	
DPN	0.296 ± 0.001	$0.193{\pm}0.001$	0.104 ± 0.016	$0.296 {\pm} 0.001$	$0.193{\pm}0.001$	$0.104{\pm}0.016$	
CVAE	0.333 ± 0.000	$0.244{\pm}0.000$	—	$0.345 {\pm} 0.000$	$0.208 {\pm} 0.000$		
A-CNF	0.239 ± 0.001	$0.097 {\pm} 0.002$	$0.382{\pm}0.015$	$0.297{\pm}0.001$	0.087 ± 0.008	$0.198 {\pm} 0.027$	
S-CNF	0.218±0.000	$0.020{\pm}0.002$	$0.068{\pm}0.021$	0.263±0.001	$\textbf{0.002}{\pm}\textbf{0.001}$	$0.026{\pm}0.012$	

Table 2: Comparison to the baselines in terms of inter-annotator agreement simulation. CVAE collapses to one category for all inputs. The best value in each column is shown in bold and the second best is underlined.

Emotion	Training (sec)	Testing (sec)
MCDP	7.20±0.10E+03	1.82±0.01E+04
Ensemble	$1.46 \pm 0.00 \text{E+05}$	1.67±0.01E+03
BBB	7.55±0.01E+03	$1.79 \pm 0.01 \text{E} + 04$
DPN	6.80±0.01E+03	2.90±0.01E+02
A-CNF	$7.04 \pm 0.02 \text{E+03}$	2.31±0.07E+02
S-CNF	6.99±0.00E+03	$2.12 \pm 0.02 \text{E+}02$
Toxic	Training (sec)	Testing (sec)
Toxic MCDP	Training (sec) 2.42±0.02E+02	Testing (sec) 5.99±0.02E+02
Toxic MCDP Ensemble	Training (sec) 2.42±0.02E+02 2.39±0.01E+03	Testing (sec) 5.99±0.02E+02 4.00±0.04E+01
Toxic MCDP Ensemble BBB	Training (sec) 2.42±0.02E+02 2.39±0.01E+03 3.22±0.01E+02	Testing (sec) 5.99±0.02E+02 4.00±0.04E+01 5.79±0.01E+02
Toxic MCDP Ensemble BBB DPN	Training (sec) 2.42±0.02E+02 2.39±0.01E+03 3.22±0.01E+02 1.92±0.01E+02	Testing (sec) 5.99±0.02E+02 4.00±0.04E+01 5.79±0.01E+02 2.67±0.02E+01
Toxic MCDP Ensemble BBB DPN A-CNF	Training (sec) 2.42±0.02E+02 2.39±0.01E+03 3.22±0.01E+02 1.92±0.01E+02 3.14±0.04E+02	Testing (sec) 5.99±0.02E+02 4.00±0.04E+01 5.79±0.01E+02 2.67±0.02E+01 1.40±0.11E+01

Table 3: Computational wall-clock time. The number M_* of annotations to simulate is set to 10 for ensemble and 100 for all other methods.

time and $M_* \times$ more testing time. MCDP and BBB require M_* forward passes during testing to generate M_* samples and therefore cost $M_* \times$ more testing time. All other methods require a single forward pass. S-CNF has slightly longer training time than DPN due to sampling procedure in ELBO while being the most efficient for testing.

501

502

504

506

507

508

509

511

512

513

514

515

516

517

518

519

5.3 Adjusting the Variability of S-CNFs by Prior Tempering

One advantage of S-CNF is that its sample variability can be easily controlled on demand without re-training by tempering the standard deviation of $p_{\Lambda}(z|x)$ at test time. Figure 3 explores the effect of prior tempering on the performance. Overall, the trend is clear that the simulated annotations become more diverse as the temperature increases, shown by the increase in the average standard deviation (\bar{s}) and the decrease in Fleiss's kappa $(\hat{\kappa})$ of simulated annotations. As the temperature decreases,



Figure 3: The effect of prior tempering on S-CNF. The x-axis corresponds to the prior temperature.

simulated annotations tends to concentrate more around the mean. The default temperature value 1 used during training (*i.e.*, no tempering) achieves the best trade-off among majority prediction accuracy (Acc), distribution matching (NLL^{all}), and inter-annotator disagreement simulation (in terms of $\mathcal{E}(\bar{s})$ and $\mathcal{E}(\hat{\kappa})$). In addition, prior tempering in S-CNF covers a wider range of dynamics than adjusting the dropout rate in MCDP. More details can be found in Appendix H.

5.4 Case Study

To better illustrate the properties of the annotations simulated by different methods, we visualize the

520



Figure 4: Visualization of simulated annotations on the emotion class labelling task for utterance (a) "0114_0263.wav", (b) '0167_0179_0001.wav", (c) "0574_0476.wav". The x-axis corresponds to emotion classes Angry, Sad, Happy, Neutral, Others. The y-axis corresponds to the probability mass. Each sample is a categorical distribution. The probability mass values of different categories in each categorical distribution are connected for the purpose of better visualization.

simulated distributions against the ground-truth distributions for three representative examples in Figure 4 (more case study examples can be found in Appendix I). Overall, the mean of the samples generated by S-CNF aligns the best with the average human label, indicating its superior performance in estimating the aggregated behaviours of human annotators. In addition, the samples generated by S-CNF are the most diverse among all compared methods, which manage to simulate the variability of the behaviours of different individual human annotators. In sharp contrast, the samples generated by all the other methods highly concentrate around their sample means. The visualized result for each example is analyzed below.

> In case (a), human annotators reach a consensus. The majority of samples generated by S-CNF ex

hibit prominent peaks aligned with the ground-truth emotion class "neutral". In contrast, many samples generated by A-CNF peak at other emotion classes.

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

570

571

572

573

574

575

576

577

578

579

580

581

582

583

585

586

587

589

590

591

592

593

594

595

597

598

599

In case (b), human opinions diverge. The majority of samples generated by S-CNF are sharp categorical distributions peaking at one of the two majority emotion classes "happy" and "neutral". Additionally, a few samples generated by S-CNF peak at the emotion class "angry", which manages to simulate the minority viewpoint held by some annotators. Very few human annotators attribute this utterance to the emotion classes "sad" and "other", and S-CNF likewise produces scarce samples peaking at these classes.

In case (c), five human annotators give distinct emotion labels, resulting in a tie in the label means. The tie comes from annotators' diverse individual perceptions of the emotion rather than consensus on its ambiguity. S-CNF is the only model that can simulate both the diverse behaviours of different individual annotators and the aggregated behaviour of all annotators since the individual samples are sharp categorical distributions peaking at one of the five emotion classes and the mean of the samples aligns well with the label mean.

A case study of toxic speech detection exhibits similar trends and can be found in Appendix J.

6 Conclusions

This paper studied human annotator simulation (HAS), a cost-effective alternative to generating human-like annotations for automatic data labelling and model evaluation. A novel framework is proposed to incorporate the variability of human evaluations into HAS. It leverages diverse annotations to estimate the distribution of human annotations by meta-learning a conditional softmax flow (S-CNF) on large crowd-sourced datasets. This overcomes the drawbacks of prior work and enables efficient generation of annotations that exhibit human variability for unlabelled test inputs. The proposed method clearly and consistently outperformed a wide range of methods on emotion class labelling and toxic speech detection, achieving the best performance for human annotation distribution matching and inter-annotator disagreement simulation. It is hoped that the proposed method could help mitigate unfair biases and over-representation in HAS and reduce the exposure of human annotators to potentially harmful content, thus promoting ethical AI practices.

549

60

60

- 00
- 60
- 60
- 60
- 610
- 611

Limitations

This work focuses on categorical annotations, which is commonly used during human evaluation. Other types of annotations can be accommodated by designing suitable corresponding output transformations p(y|v).

The proposed S-CNF is tested for two representative tasks: emotion recognition and toxic speech detection with speech and text as input respectively. We believe that the proposed method can also be general to other tasks, which is kept for future research directions.

Ethics Statement

In this work, all human annotations used for training were taken from existing publicly available corpora, and no new human annotations were collected.

It is hoped that this work could play a part in promoting ethical AI practice. Firstly, it has been 618 shown that the proposed HAS system can capture the inherent variability in human judgements and help mitigate biases and the issue of overrepresentation, thus producing a more inclusive 622 representation of human opinions. The proposed HAS system also has the potential to minimize human annotators' exposure to offensive and/or hateful content in some evaluation tasks such as 626 HateXplain. However, as with most research in ma-627 chine learning, new modelling techniques could be used by bad actors to cause harm more effectively, but we do not see how the proposed HAS system is more concerning than any other method in this 631 regard. 632

References

633

641

643

647

- Cecilia Ovesdotter Alm. 2011. Subjective natural language problems: Motivations, applications, characterizations, and implications. In *Proc. ACL*, Portland, USA.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. 2015. Weight uncertainty in neural network. In *Proc. ICML*, Lille, France.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E.M. Provost, S. Kim, J.N. Chang, S. Lee, and S.S. Narayanan. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42:335–359.

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

- Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022. WavLM: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518.
- Wenlin Chen, Austin Tripp, and José Miguel Hernández-Lobato. 2023. Meta-learning adaptive deep kernel gaussian processes for molecular property prediction. In *Proc. ICLR*, Kigali, Rwanda.
- Huang-Cheng Chou and Chi-Chun Lee. 2019. Every rating matters: Joint learning of subjective labels and individual annotators for speech emotion classification. In *Proc. ICASSP*, Brighton, UK.
- Paul F. Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Proc. NeurIPS*, Long Beach, USA.
- Ting Dang, Vidhyasaharan Sethu, Julien Epps, and Eliathamby Ambikairajah. 2017. An investigation of emotion prediction uncertainty using gaussian mixture regression. In *Proc. Interspeech*, Stockholm, Sweden.
- 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.
- Jun Deng, Wenjing Han, and Björn Schuller. 2012. Confidence measures for speech emotion recognition: A start. In *Speech Communication; 10. ITG Symposium*, pages 1–4. VDE.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proc. NAACL*, Minneapolis, USA.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. 2017. Density estimation using real NVP. In *Proc. ICLR*, Toulon, France.
- Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proc. AAAI*, New Orleans, USA.
- Nemanja Djuric, Jing Zhou, Robin Morris, Mihajlo Grbovic, Vladan Radosavljevic, and Narayan Bhamidipati. 2015. Hate speech detection with comment embeddings. In *Proc. WWW*, Florence, Italy.

- 701 704 705 710 711 714 715 716 717 718 719 723 794 726 727 730 731 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
- 751

- H.M. Fayek, M. Lech, and L. Cavedon. 2016. Modeling subjectiveness in emotion recognition with deep neural networks: Ensembles vs soft labels. In Proc. IJCNN, Vancouver, Canada.
- Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. Psychological Bulletin, 76(5):378.
- Yarin Gal and Zoubin Ghahramani. 2016. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In Proc. ICML, New York, USA.
- Izzeddin Gür, Dilek Hakkani-Tür, Gokhan Tür, and Pararth Shah. 2018. User modeling for task oriented dialogues. In Proc. SLT, Athens, Greece.
- Jing Han, Zixing Zhang, Maximilian Schmitt, Maja Pantic, and Björn Schuller. 2017. From hard to soft: Towards more human-like emotion recognition by modelling the perception uncertainty. In Proc. ACM MM, Mountain View, USA.
- Peter Harrison, Raja Marjieh, Federico Adolfi, Pol van Rijn, Manuel Anglada-Tort, Ofer Tchernichovski, Pauline Larrouy-Maestri, and Nori Jacoby. 2020. Gibbs sampling with people. In Proc. NeurIPS, Vancouver, Canada.
- Martie G Haselton, Daniel Nettle, and Paul W Andrews. 2015. The evolution of cognitive bias. The Handbook of Evolutionary Psychology, pages 724-746.
- Julia Hirschberg, Jackson Liscombe, and Jennifer Venditti. 2003. Experiments in emotional speech. In Proc. SSPR, Tokyo, Japan.
- Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. 2021. Argmax flows and multinomial diffusion: Learning categorical distributions. In Proc. NeurIPS, Online.
- Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Social biases in NLP models as barriers for persons with disabilities. In Proc. ACL, Online.
- Y. Kim, H. Lee, and E. M. Provost. 2013. Deep learning for robust feature generation in audiovisual emotion recognition. In Proc. ICASSP, Vancouver, Canada.
- Diederik P Kingma and Max Welling. 2014. Autoencoding variational Bayes. In Proc. ICLR, Banff, Canada.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and scalable predictive uncertainty estimation using deep ensembles. In Proc. NeurIPS, Long Beach, USA.
- Yichong Leng, Xu Tan, Sheng Zhao, Frank Soong, Xiang-Yang Li, and Tao Qin. 2021. MBNet: MOS prediction for synthesized speech with mean-bias network. In Proc. ICASSP, Toronto, Canada.

Hsien-Chin Lin, Nurul Lubis, Songbo Hu, Carel van Niekerk, Christian Geishauser, Michael Heck, Shutong Feng, and Milica Gasic. 2021. Domainindependent user simulation with transformers for task-oriented dialogue systems. In Proc. SIGDIAL, Singapore City, Singapore.

752

753

754

755

756

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

775

780

781

782

783

784

785

786

787

788

789

792

793

794

795

796

797

799

800

801

802

803

804

805

806

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- R. Lotfian and C. Busso. 2019. Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings. IEEE Transactions on Affective Computing, 10(4):471-483.
- Kede Ma, Kai Zeng, and Zhou Wang. 2015. Perceptual quality assessment for multi-exposure image fusion. IEEE Transactions on Image Processing, 24(11):3345-3356.
- Andrey Malinin and Mark Gales. 2018. Predictive uncertainty estimation via prior networks. In Proc. NeurIPS, Montreal, Canada.
- Georgia Maniati, Alexandra Vioni, Nikolaos Ellinas, Karolos Nikitaras, Konstantinos Klapsas, June Sig Sung, Gunu Jho, Aimilios Chalamandaris, and Pirros Tsiakoulis. 2022. SOMOS: The Samsung open MOS dataset for the evaluation of neural text-to-speech synthesis. In Proc. Interspeech, Incheon, Korea.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2021. HateXplain: A benchmark dataset for explainable hate speech detection. In Proc. AAAI, Vancouver, Canada.
- Rada Mihalcea and Hugo Liu. 2006. A corpus-based approach to finding happiness. In Proc. AAAI Spring Symposium, Stanford, USA.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. In Proc. NeurIPS, New Orleans, USA.
- Brian Patton, Yannis Agiomyrgiannakis, Michael Terry, Kevin Wilson, Rif A. Saurous, and D. Sculley. 2016. AutoMOS: Learning a non-intrusive assessor of naturalness-of-speech. In Proc, NeurIPS Workshop, Barcelona, Spain.
- Barbara Plank, Dirk Hovy, and Anders Søgaard. 2014. Learning part-of-speech taggers with inter-annotator agreement loss. In Proc. EACL, Gothenburg, Sweden.
- S. Poria, E. Cambria, R. Bajpai, and A. Hussain. 2017. A review of affective computing: From unimodal analysis to multimodal fusion. Information Fusion, 37:98-125.

- 810 811 812 813 815 821 822 823 824 825 826 827 828 829
- 834 835 837 838 839 840 843 844
- 847 851
- 854

- 862

- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2019. MELD: A multimodal multi-party dataset for emotion recognition in conversations. In *Proc. ACL*, Florence, Italy.
- Vinodkumar Prabhakaran, Michael Bloodgood, Mona Diab, Bonnie Dorr, Lori Levin, Christine D. Piatko, Owen Rambow, and Benjamin Van Durme. 2012. Statistical modality tagging from rule-based annotations and crowdsourcing. In Proc. ExProM Workshop, Jeju, Korea.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. arXiv preprint arXiv:2305.18290.
- Dadi Ramesh and Suresh Kumar Sanampudi. 2022. An automated essay scoring systems: A systematic literature review. Artificial Intelligence Review, 55(3):2495-2527.
- Alexander J. Ratner, Christopher M De Sa, Sen Wu, Daniel Selsam, and Christopher Ré. 2016. Data programming: Creating large training sets, quickly. In Proc. NeurIPS, Barcelona, Spain.
- Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, Loren Lugosch, Cem Subakan, Nauman Dawalatabad, Abdelwahab Heba, Jianyuan Zhong, Ju-Chieh Chou, Sung-Lin Yeh, Szu-Wei Fu, Chien-Feng Liao, Elena Rastorgueva, François Grondin, William Aris, Hwidong Na, Yan Gao, Renato De Mori, and Yoshua Bengio. 2021. SpeechBrain: A general-purpose speech toolkit. arXiv preprint arXiv:2106.04624.
- Nataniel Ruiz, Samuel Schulter, and Manmohan Chandraker. 2019. Learning to simulate. In Proc. ICLR, New Orleans, USA.
- Adam Sanborn and Thomas Griffiths. 2007. Markov chain Monte Carlo with people. In Proc. NeurIPS, volume Vancouver, Canada.
- Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young. 2007. Agenda-based user simulation for bootstrapping a POMDP dialogue system. In Proc. NAACL, Vancouver, Canada.
- Seshadrinathan, Rajiv Soundararajan, Kalpana Alan Conrad Bovik, and Lawrence K. Cormack. 2010. Study of subjective and objective quality assessment of video. IEEE Transactions on Image Processing, 19(6):1427-1441.
- Weiyan Shi, Kun Qian, Xuewei Wang, and Zhou Yu. 2019. How to build user simulators to train RL-based dialog systems. In Proc. EMNLP-IJCNLP, Hong Kong, China.
- Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. In Proc. NeurIPS, Long Beach, USA.

Vincent Stimper, David Liu, Andrew Campbell, Vincent Berenz, Lukas Ryll, Bernhard Schölkopf, and José Miguel Hernández-Lobato. 2023. normflows: A PyTorch package for normalizing flows. Journal of Open Source Software, 8(86):5361.

863

864

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

891

892

893

894

895

896

897

898

899

900

- Hossein Talebi and Peyman Milanfar. 2018. NIMA: Neural image assessment. IEEE Transactions on Image Processing, 27(8):3998–4011.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. In Proc. NeurIPS, Barcelona, Spain.
- Martin J Wainwright, Michael I Jordan, et al. 2008. Graphical models, exponential families, and variational inference. Foundations and Trends in Machine Learning, 1(1-2):1-305.
- Janyce Wiebe, Theresa Wilson, Rebecca Bruce, Matthew Bell, and Melanie Martin. 2004. Learning subjective language. Computational linguistics, 30(3):277-308.
- Renzhi Wu, Shen-En Chen, Jieyu Zhang, and Xu Chu. 2022a. Learning hyper label model for programmatic weak supervision. In Proc. ICLR, Online.
- Wen Wu, Chao Zhang, and Philip Woodland. 2023. Estimating the uncertainty in emotion attributes using deep evidential regression. In Proc. ACL, Toronto, Canada.
- Wen Wu, Chao Zhang, and Philip C. Woodland. 2021. Emotion recognition by fusing time synchronous and time asynchronous representations. In Proc. ICASSP, Toronto, Canada.
- Wen Wu, Chao Zhang, Xixin Wu, and Philip Woodland. 2022b. Estimating the uncertainty in emotion class labels with utterance-specific Dirichlet priors. IEEE Transactions on Affective Computing.
- Mathieu Zen and Jean Vanderdonckt. 2016. Assessing user interface aesthetics based on the intersubjectivity of judgment. In Proc. BCS HCI, Poole, UK.

907

908

909

910

911

912 913

914

915

916

917

918

919

921

922

923

924

925

929

930

931

933

934

935

936

937

938

942

902

A The Sources of Variability in Human Evaluation

Human perception refers to the process by which individuals interpret and make sense of the sensory information they receive from the environment. It involves the integration of sensory data, cognitive processes, emotions, and previous experiences. Subjective perception emphasizes that each individual's perception of the world is unique and influenced by their internal mental states, beliefs, attitudes, and past experiences. As a result, people can interpret and react to the same stimuli differently, leading to diverse and subjective perceptions.

Each person's sensory organs, such as eyes and ears, may have slight variations in sensitivity and acuity, leading to different perceptions of the same stimuli. Cognitive biases, the inherent mental shortcuts or tendencies that influence how humans perceive and process information, can lead to difference in judgment and decision-making. People's past experiences, cultural norms, and upbringing also shape their perceptions. Different cultural backgrounds can lead to distinct interpretations of the same event, leading to diverse reaction. The variability in humans can manifest in various tasks such as colour perception, emotion recognition, art appreciation, and feedback preferences.

Embracing and understanding the variability of human perception is vital for various research fields such as psychology, neuroscience, humancomputer interaction, and so on, and has practical implications in designing human-centered systems and promoting empathy and diversity. It helps create products and interfaces that cater to diverse user needs and preferences in fields like humancomputer interaction and user experience design. Being aware of the variability of perception is crucial in ethical decision-making. It help ensures that different perspectives and cultural sensitivities are considered, which helps identify and address potential biases that might disproportionately affect certain groups or lead to unfair outcomes.

B Derivations

945Detailed derivations for the training objectives on946a single dataset $\mathcal{D}_i = \{\eta_i^{(1)}, \cdots, \eta_i^{(M)}\}$ with x_i 947are presented in this section. For the simplicity of948notations, the subscription i in our derivations will949be omitted without ambiguity where possible. The950meta-learning objectives presented in the paper are951obtained by averaging such single-task objectives

across tasks.

B.1 Objective Function for the Base CNF

Denote the empirical human annotation distribution as $p_m(y|x) = \delta(y - \eta^{(m)})$, $m = 1, \dots, M$ and model output distribution as $p_\theta(y|x)$. The average KL divergence between them over all M human annotations for this input x is given by: 958

$$\mathcal{L}(\boldsymbol{\theta}; \mathcal{D}, \boldsymbol{x}) = \frac{1}{M} \sum_{m=1}^{M} \mathcal{KL}\left(\mathbf{p}_{m}(\boldsymbol{y}|\boldsymbol{x}) \parallel \mathbf{p}_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x})\right)$$
$$= \frac{1}{M} \sum_{m=1}^{M} \int \mathbf{p}_{m}(\boldsymbol{y}|\boldsymbol{x}) \log \frac{\mathbf{p}_{m}(\boldsymbol{y}|\boldsymbol{x})}{\mathbf{p}_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x})} d\boldsymbol{y}$$
$$= -\frac{1}{M} \sum_{m=1}^{M} \int \mathbf{p}_{m}(\boldsymbol{y}|\boldsymbol{x}) \log \mathbf{p}_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x}) d\boldsymbol{y} + \text{const}$$
$$= -\frac{1}{M} \sum_{m=1}^{M} \log \mathbf{p}_{\boldsymbol{\theta}}(\boldsymbol{\eta}^{(m)}|\boldsymbol{x}) + \text{const}$$

Minimizing this KL objective is equivalent to maximizing the average log likelihood $\log p_{\theta}(\eta^{(m)}|x)$ over all human annotations as presented in the paper.

B.2 Objective Function for S-CNF

For categorical annotations, each label $\eta^{(m)}$ represents the probabilities of all categories in the categorical human annotation distribution: $\eta^{(m)} = [\eta_1^{(m)}, \dots, \eta_K^{(m)}]$, where $\eta_k^{(m)} = P_m(c = k | x)$. Denote the model output distribution as $P_{\theta}(c | x)$. The average KL divergence between them over all M human annotations for this input x is given by:

$$\mathcal{L}^{\text{exact}}(\boldsymbol{\theta}; \mathcal{D}, \boldsymbol{x}) = \frac{1}{M} \sum_{m=1}^{M} \mathcal{KL} \left(\mathbf{P}_{m}(c | \boldsymbol{x}) \parallel \mathbf{P}_{\boldsymbol{\theta}}(c | \boldsymbol{x}) \right)$$
$$= \frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \mathbf{P}_{m}(c = k | \boldsymbol{x}) \log \frac{\mathbf{P}_{m}(c = k | \boldsymbol{x})}{\mathbf{P}_{\boldsymbol{\theta}}(c = k | \boldsymbol{x})}$$
$$= -\frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \mathbf{P}_{m}(c = k | \boldsymbol{x}) \log \mathbf{P}_{\boldsymbol{\theta}}(c = k | \boldsymbol{x})$$
$$+ \text{const}$$
$$= -\frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \boldsymbol{\eta}_{k}^{(m)} \log \mathbf{P}_{\boldsymbol{\theta}}(c = k | \boldsymbol{x}) + \text{const},$$

952

953

920

961 962

960

963

964

965

966

967

968

969

970

971

975

976 977

978

979

981

982

983

985

974

where the marginal likelihood is lower bounded using variational inference:

$$\begin{split} \log \mathsf{P}_{\boldsymbol{\theta}}(c=k|\boldsymbol{x}) &= \log \int \mathsf{P}(c=k|\boldsymbol{v})\mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x})\mathsf{d}\boldsymbol{v} \\ &= \log \int \boldsymbol{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{\eta}) \frac{\mathsf{P}(c=k|\boldsymbol{v})\mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x})}{\boldsymbol{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{\eta})} \mathsf{d}\boldsymbol{v} \\ &\geq \int \boldsymbol{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{\eta}) \log \frac{\mathsf{P}(c=k|\boldsymbol{v})\mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x})}{\boldsymbol{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{\eta})} \mathsf{d}\boldsymbol{v} \\ &= \mathbb{E}_{\boldsymbol{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{\eta})} \left[\log \mathsf{P}(c=k|\boldsymbol{v}) + \log \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x}) \\ &- \log \boldsymbol{q}_{\boldsymbol{\Omega}}(\boldsymbol{v}|\boldsymbol{\eta})\right]. \end{split}$$

Therefore, the final negative ELBO objective is obtained by

$$\begin{split} \mathcal{L}^{\text{exact}} &= -\frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \boldsymbol{\eta}_{k}^{(m)} \log \mathsf{P}_{\boldsymbol{\theta}}(c=k|\boldsymbol{x}) \\ &\leq -\frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \boldsymbol{\eta}_{k}^{(m)} \mathbb{E}_{\boldsymbol{q}_{\Omega}(\boldsymbol{v}|\boldsymbol{\eta}^{(m)})} \left[\log \mathsf{P}(c=k|\boldsymbol{v}) \right. \\ &\left. + \log \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x}) - \log \boldsymbol{q}_{\Omega}(\boldsymbol{v}|\boldsymbol{\eta}^{(m)}) \right] \\ &= -\frac{1}{M} \sum_{m=1}^{M} \mathbb{E}_{\boldsymbol{q}_{\Omega}(\boldsymbol{v}|\boldsymbol{\eta}^{(m)})} \left[\sum_{k=1}^{K} \boldsymbol{\eta}_{k}^{(m)} \log \mathsf{P}(c=k|\boldsymbol{v}) \right. \\ &\left. + \log \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v}|\boldsymbol{x}) - \log \boldsymbol{q}_{\Omega}(\boldsymbol{v}|\boldsymbol{\eta}^{(m)}) \right] \\ &= \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathcal{D}, \boldsymbol{x}), \end{split}$$

where

$$\begin{split} \log \mathsf{P}(c = k | \boldsymbol{v}) &= \operatorname{logsoftmax}(\boldsymbol{v})_k, \\ \log \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v} | \boldsymbol{x}) &= \mathsf{p}_{\boldsymbol{\Lambda}} \left(\boldsymbol{f}_{\boldsymbol{\theta}}^{-1}(\boldsymbol{v}) | \boldsymbol{x} \right) \left| \det \left(\frac{\partial \boldsymbol{f}_{\boldsymbol{\theta}}^{-1}(\boldsymbol{v})}{\partial \boldsymbol{v}} \right) \right. \\ & \left. \log \boldsymbol{q}_{\boldsymbol{\Omega}}(\boldsymbol{v} | \boldsymbol{\eta}^{(m)}) = \mathcal{N}(\boldsymbol{v} | \boldsymbol{\mu}_{\boldsymbol{\Omega}}(\boldsymbol{\eta}^{(m)}), \operatorname{diag}(\boldsymbol{\sigma}_{\boldsymbol{\Omega}}^2(\boldsymbol{\eta}^{(m)})) \right| \end{split}$$

B.3 The Negative Log Likelihood (NLL_i^{all}) for Categorical Annotations

The marginal likelihood of S-CNF is intractable, which can be approximated using Monte-Carlo simulation:

$$\begin{split} \mathsf{P}_{\boldsymbol{\theta}}(c = k | \boldsymbol{x}) &= \int \mathsf{P}(c = k | \boldsymbol{v}) \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v} | \boldsymbol{x}) \mathsf{d} \boldsymbol{v} \\ &= \mathbb{E}_{\mathsf{P}_{\boldsymbol{\theta}}(\boldsymbol{v} | \boldsymbol{x})} \left[\mathsf{P}(c = k | \boldsymbol{v}) \right] \\ &\approx \frac{1}{Q} \sum_{j=1}^{Q} \mathsf{P}(c = k | \boldsymbol{v}_j), \quad \{\boldsymbol{v}_j\}_{j=1}^{Q} \sim_{\mathrm{iid}} \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v} | \boldsymbol{x}) \\ &= \frac{1}{Q} \sum_{j=1}^{Q} \mathrm{softmax}(\boldsymbol{v}_j)_k, \quad \{\boldsymbol{v}_j\}_{j=1}^{Q} \sim_{\mathrm{iid}} \mathsf{p}_{\boldsymbol{\theta}}(\boldsymbol{v} | \boldsymbol{x}) \\ &= \bar{\boldsymbol{y}}_k, \end{split}$$

where $\bar{\boldsymbol{y}} = \frac{1}{Q} \sum_{j=1}^{Q} \operatorname{softmax}(\boldsymbol{v}_j) = \frac{1}{Q} \sum_{j=1}^{Q} \boldsymbol{y}_j$ 989 which is the average of the simulated categorical 990 distributions. Let $\bar{\boldsymbol{\eta}} = \frac{1}{M} \sum_{m=1}^{M} \boldsymbol{\eta}^{(m)}$ be the average label. 992

993

994

996

997

998

999

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1028

Then, the NLL^{all} for a single input x_i is given by

$$\begin{aligned} \text{NLL}_{i}^{\text{all}} &= -\frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \boldsymbol{\eta}_{i,k}^{(m)} \log \mathsf{P}_{\boldsymbol{\theta}}(c=k|\boldsymbol{x}_{i}) \\ &\approx -\frac{1}{M} \sum_{m=1}^{M} \sum_{k=1}^{K} \boldsymbol{\eta}_{i,k}^{(m)} \log \bar{\boldsymbol{y}}_{i,k} \end{aligned} \qquad 995 \\ &= -\sum_{k=1}^{K} \bar{\boldsymbol{\eta}}_{i,k} \log \bar{\boldsymbol{y}}_{i,k}, \end{aligned}$$

which is the cross entropy between the averaged label and averaged sample.

C Emotion Label Processing for MSP-Podcast

Ν

In MSP-Podcast, each annotator can choose from 1000 ten emotion classes to label the primary emotion of an utterance: Angry, Sad, Happy, Surprise, Fear, Disgust, Contempt, Neutral, Other. Although only 1003 one option is allowed, they can say other and de-1004 fine their own emotion class which can be more than one. During label processing, the original 1006 other class is split into sub-classes depending on 1007 the manual defined label and merged with the pre-1008 defined labels. The grouping details are shown as 1009 follows: (i) Angry includes angry, disgust, con-1010 tempt, annoyed; (ii) Sad includes sad, frustrated, 1011 disappointed, depressed, concerned; (iii) Happy 1012 ^(m)))))includes happy, excited, amused; (iv) Neutral in-1013 cludes neutral; (v) Other includes all other emotion 1014 subclasses not listed above. 1015

D Model Structure Details

The procedure of sampling from and optimizing S-CNF are summarized in Algorithm 1 and 2.

A neural-network-based encoder g_{Λ} is built to model $\mu_{\Lambda}(x), \sigma_{\Lambda}^2(x)$ given input x where Λ is the model parameters. g_{Λ} follows an upstreamdownstream paradigm. The upstream model is pretrained on large amount of unlabelled data to learn universal representations. The downstream model uses the learned representation from the upstream model for specific applications.

For tasks involving speech as input (*i.e.*, emotion class labelling), WavLM (Chen et al., 2022) is used

Algorithm 1 Sampling from S-CNF

Input: xOutput: Categorical probability yCompute $\mu_{\Lambda}(x), \sigma_{\Lambda}^2(x) = g_{\Lambda}(x)$ Sample $z \sim \mathcal{N}(\mu_{\Lambda}(x), \operatorname{diag}(\sigma_{\Lambda}^2(x)))$ Compute $v = f_{\theta}(z)$ Compute $y = \operatorname{softmax}(v)$

Algorithm 2 Optimizing S-CNF

1029

1031

1032

1033

1034

1036

1037

1038

1039

1040

1041

1042

1044

1045

1046

1047

1048

1049

1050

1052

as the upstream model. WavLM is a speech foundation model pre-trained by self-supervised learning that takes raw waveform as input. The waveform is encoded by a CNN encoder followed by multiple Transformer encoders. The BASE+ version³ of the model is used in this work which has 12 Transformer encoder blocks with 768-dimensional hidden states and 8 attention heads. The parameters of the pretrained WavLM are frozen and the weighted sum of the outputs of the 12 Transformer encoder blocks is used as the speech embeddings feeding into the downstream model.

RoBERTa (Liu et al., 2019) is used as upstream model to encode text input for toxic speech detection, which is a robustly optimized model of BERT (Devlin et al., 2019). RoBERTa is a Transformer-based language model pretrained on a large corpus of English data with the masked language modelling objective. The BASE version⁴ was used in the work which has 12 Transformer layers, 768 hidden units, 12 attention heads, and 125 million parameters.

The downstream model consists of two Transformer encoder blocks followed by two FC layers. The Transformer encoder layers has a hidden di-1053 mension of 128 and four attention head. The out-1054 put layer contains two heads to predict the mean 1055 and standard deviation of the latent distribution 1056 $p_{\Lambda}(\boldsymbol{z}|\boldsymbol{x})$. The invertible flow model f_{θ} uses three 1057 real NVP block (Dinh et al., 2017) of dimension 64. 1058 The variational encoder h_{Ω} for S-CNF contains a 1059 FC layer of dimension 64 and two output heads for the mean and standard deviation of the variational 1061 distribution $q_{\Omega}(\boldsymbol{v}|\boldsymbol{y})$. 1062

1064

1084

E Detailed Configuration of All Compared Methods

Ensemble consists of 10 systems initialized and 1065 trained using different random seeds. MCDP uses 1066 dropout rate of 0.4. A standard Gaussian prior 1067 is used for BBB. A modified version of EDL is used (Wu et al., 2023) which is trained by maximiz-1069 ing the per-observation-based marginal likelihood 1070 with a modified regularization term. Ensemble, 1071 MCDP, BBB, EDL use the same model structure 1072 as g_{Λ} apart from removing the output head for pre-1073 dicting variance of latent distribution. A modified 1074 version of DPN (Wu et al., 2022b) is used which 1075 is trained by interpolating per-observation-based 1076 marginal likelihood with KL divergence. The co-1077 efficient of KL term is set to 5.0 for emotion class 1078 labelling and 2.0 for toxic speech detection. CVAE 1079 has the same g_{Λ} structure as S-CNF for modelling 1080 $p(\boldsymbol{z}|\boldsymbol{x})$, and two 64-d FC layers are used for en-1081 coder and decoder. A-CNF has identical model 1082 structure as S-CNF.

F Implementation Details

The system was implemented using PyTorch with 1085 the SpeechBrain (Ravanelli et al., 2021) and norm-1086 flows (Stimper et al., 2023) toolkit. The Adadelta 1087 optimizer was used with an initial learning rate of 1.2 for emotion class labelling and 0.05 for speech 1089 quality assessment. The NewBob learning rate 1090 scheduler was used with annealing factor 0.8 and 1091 patience 1. The system was trained for 30 epochs 1092 and the model with the best validation performance 1093 was used for testing. The number of ELBO sam-1094 ples was set to 20. All experiments were run with 1095 three different seeds and the average and standard error are reported. 1097

³https://huggingface.co/microsoft/wavlm-base-plus

⁴https://huggingface.co/roberta-base

Emotion class labelling					
	RMSE ^s	\mathbf{MAE}^{s}	$\mathcal{E}(ar{s})$		
MCDP	0.305	0.233	0.206		
Ensemble	0.277	0.222	0.166		
BBB	0.284	0.226	0.178		
CVAE	0.333	0.244	0.244		
DPN	0.297	0.236	0.191		
A-CNF	0.223	0.209	0.046		
S-CNF	0.218	0.198	0.015		
Тох	kic speech (detection			
	RMSE ^s	MAE ^s	$\mathcal{E}(ar{s})$		
MCDP	0.297	0.242	0.122		
Ensemble	0.290	0.220	0.105		
BBB	0.279	0.229	0.115		
CVAE	0.345	0.208	0.208		
DPN	0.299	0.220	0.178		
A_CNE	0 274	0 232	0.062		
A-CIVI	0.274	0.252	0.002		

 Table 4: Analysis of standard deviation of simulated samples.

G Analysis of Standard Deviation of Simulated Samples

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

It has been observed in Sec. 5.1 that flow models tend to have a larger difference between RMSE^s and $\mathcal{E}(\bar{s})$. This section provides detailed analysis to this observation. Let N be the number of test utterances. Three std-related metrics are computed: (i) RMSE between std of predictions and human labels: RMSE^s = $\sqrt{\frac{1}{N} \sum_{i=1}^{N} (s_i - \sigma_i)^2}$; (ii) Mean absolute error between std of predictions and std of human labels: MAE^s = $\frac{1}{N} \sum_{i=1}^{N} |s_i - \sigma_i|$; (iii) Absolute error between average std of predictions and average std of human labels $\mathcal{E}(\bar{s}) = |\bar{s}_i - \bar{\sigma}_i|$. Results are shown in Table 4. The flow model tends to have larger discrepancy between MAE^s and $\mathcal{E}(\bar{s})$. According to the triangular inequality:

$$\begin{aligned} \mathcal{E}(\bar{s}) &= \left| \frac{1}{N} \sum_{i=1}^{N} s_i - \frac{1}{N} \sum_{i=1}^{N} \sigma_i \right| \\ &= \left| \frac{1}{N} \sum_{i=1}^{N} (s_i - \sigma_i) \right| \le \frac{1}{N} \sum_{i=1}^{N} |s_i - \sigma_i| = \mathsf{MAE}^s \end{aligned}$$

1115which show that $\mathcal{E}(\bar{s})$ is a lower bound of MAE^s.1116The equality condition is satisfied when all samples are uniformly either greater than or less than1117ples are uniformly either greater than or less than1118the compared value. Therefore, a larger discrep-1119ancy between these two values indicates that the



Figure 5: Standard deviation of simulated samples for emotion class labelling.

standard deviation of some samples exceeds that of 1120 the labels, while for others, it is lower. A smaller 1121 discrepancy indicates that the standard deviation 1122 of samples tend to be consistently larger of smaller 1123 than that of the labels. In Figure 5, 100 test utter-1124 ances are randomly selected and the std of samples 1125 generated by different models are plotted, which 1126 supports the above conclusion. The proposed S-1127 CNF has the best performance for matching the 1128 diversity of human annotations. 1129

1130

1131

H Adjusting the Variability of CNFs by Prior Tempering

Sec. 5.3 has explored the effect of prior tempering 1132 on the performance. More details are provided in 1133 this section. Table 5 shows the effect of prior tem-1134 pering on the performance of S-CNF. Table 6 shows 1135 the effect of dropout rate on the performance of 1136 MCDP. When the temperature increases or dropout 1137 rate increases, the simulated annotations become 1138 more diverse, shown by the increase in the average 1139 standard deviation (\bar{s}) and the decrease in Fleiss's 1140 kappa ($\hat{\kappa}$) of simulated annotations. Comparing 1141 two tables, it can be seen that prior tempering in 1142 CNF is more efficient and covers a wider range of 1143 dynamics than adjusting the dropout rate in MCDP. 1144

		E	motion cla	ss labell	ing		
Т	Acc	NLL ^{all}	RMSE ^s	$ar{s}$	$\mathcal{E}(ar{s})$	$\hat{\kappa}$	${\cal E}(\hat\kappa)$
0.8	0.594	1.395	0.221	0.200	0.044	0.307	0.053
0.9	0.594	1.390	0.219	0.216	0.029	0.259	0.005
1.0	0.593	1.389	0.218	0.229	0.015	0.222	0.032
1.1	0.592	1.389	0.218	0.241	0.004	0.191	0.063
1.2	0.590	1.391	0.219	0.251	0.007	0.166	0.088
		Т	oxic speecl	h detecti	on		
Т	Acc	NLL ^{all}	RMSE ^s	$ar{s}$	$\mathcal{E}(ar{s})$	$\hat{\kappa}$	${\cal E}(\hat\kappa)$
0.8	0.671	0.851	0.272	0.125	0.065	0.721	0.115
0.9	0.675	0.842	0.265	0.157	0.033	0.650	0.044
1.0	0.673	0.837	0.263	0.188	0.002	0.580	0.026
1.1	0.671	0.836	0.264	0.216	0.026	0.512	0.094
1.2	0.669	0.837	0.267	0.242	0.052	0.450	0.156

Table 5: Adjusting the variability of CNFs by prior tempering.

		E	motion cla	ss labell	ing		
dp	Acc	NLL ^{all}	RMSE ^s	$ar{s}$	$\mathcal{E}(ar{s})$	$\hat{\kappa}$	${\cal E}(\hat\kappa)$
0.1	0.583	1.463	0.303	0.040	0.205	0.791	0.537
0.2	0.589	1.426	0.303	0.040	0.204	0.773	0.519
0.3	0.590	1.415	0.300	0.045	0.199	0.761	0.507
0.4	0.585	1.405	0.296	0.051	0.194	0.723	0.469
0.5	0.589	1.409	0.294	0.053	0.191	0.715	0.461
		Т	oxic speecl	h detecti	on		
dp	Acc	T NLL ^{all}	oxic speecl RMSE ^s	h detecti \overline{s}	on $\mathcal{E}(ar{s})$	$\hat{\kappa}$	${\cal E}(\hat\kappa)$
dp 0.1	Acc 0.661	T NLL ^{all} 0.925	oxic speecl RMSE ^s 0.314	h detecti \bar{s}	$\frac{\mathcal{E}(\bar{s})}{0.158}$	<i>κ</i> 0.831	$\mathcal{E}(\hat{\kappa})$ 0.225
dp 0.1 0.2	Acc 0.661 0.666	T NLL ^{all} 0.925 0.916	boxic speecl RMSE ^s 0.314 0.308	h detecti	ton $\mathcal{E}(\bar{s})$ 0.158 0.147	 κ̂ 0.831 0.800 	<i>E</i> (<i>k̂</i>) 0.225 0.194
dp 0.1 0.2 0.3	Acc 0.661 0.666 0.654	T NLL ^{all} 0.925 0.916 0.968	bxic speecl RMSE ^s 0.314 0.308 0.299	h detecti <u>s</u> 0.049 0.061 0.081	$\mathcal{E}(\bar{s})$ 0.158 0.147 0.127	 <i>k</i> 0.831 0.800 0.750 	 ε(κ̂) 0.225 0.194 0.144
dp 0.1 0.2 0.3 0.4	Acc 0.661 0.666 0.654 0.662	T NLL ^{all} 0.925 0.916 0.968 0.943	bxic speecl RMSE ^s 0.314 0.308 0.299 0.297	h detecti	con $\mathcal{E}(\bar{s})$ 0.158 0.147 0.127 0.122	 <i>k</i> 0.831 0.800 0.750 0.731 	 ε(κ̂) 0.225 0.194 0.144 0.125
dp 0.1 0.2 0.3 0.4 0.5	Acc 0.661 0.666 0.654 0.662 0.662	T NLL ^{all} 0.925 0.916 0.968 0.943 0.896	bxic speec RMSE ^s 0.314 0.308 0.299 0.297 0.296	h detecti	con $\mathcal{E}(\bar{s})$ 0.158 0.147 0.127 0.122 0.120	 <i>k</i> 0.831 0.800 0.750 0.731 0.720 	 \$\mathcal{E}(\hat{\kappa})\$ 0.225 0.194 0.144 0.125 0.114

Table 6: Adjusting the variability of MCDP models by dropout rate (dp).

I Further visualised examples: Emotion Class Labelling

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

This section shows additional visualized examples for emotion class labelling when human annotators reach a consensus (Figure 6 (a)(b)), diverge (Figure 6 (c)(d)), and give distinct labels (Figure 6 (e)). Aligned with the findings in Sec. 5.4, the proposed S-CNF can better simulate the aggregated behaviour as well as the variability of human annotations in all cases.

J Further visualised examples: Toxic Speech Detection

This section shows visualized examples for toxic 1157 speech detection when all three human annotators 1158 provide the same label (Figure 7 (a)(b)), one of 1159 them gives a different label (Figure 7 (c)(d)), and 1160 all three annoators give distinct labels (Figure 7 (e)). 1161 Similar to the observations of the emotion class 1162 labelling task, the proposed S-CNF can better sim-1163 ulate the aggregated behaviour as well as the vari-1164 ability of human annotations in all cases. 1165

1155





Figure 6: Additional visualized examples for emotion class labelling. The y-axis corresponds to the probability mass. Each sample is a categorical distribution. The probability mass values of different categories in each categorical distribution are connected for the purpose of better visualization.



Figure 7: Additional visualized examples for toxic speech detection. The y-axis corresponds to the probability mass. Each sample is a categorical distribution. The probability mass values of different categories in each categorical distribution are connected for the purpose of better visualization.