



Universiteit
Leiden
The Netherlands

Opinion diversity through hybrid intelligence

Meer, M.T. van der

Citation

Meer, M. T. van der. (2025, March 26). *Opinion diversity through hybrid intelligence*. SIKS Dissertation Series. Retrieved from <https://hdl.handle.net/1887/4209024>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/4209024>

Note: To cite this publication please use the final published version (if applicable).

5

Annotator-Centric Active Learning for Subjective NLP Tasks

5

Active Learning (AL) addresses the high costs of collecting human annotations by strategically annotating the most informative samples. However, for subjective NLP tasks, incorporating a wide range of perspectives in the annotation process is crucial to capture the variability in human judgments. We introduce Annotator-Centric Active Learning (ACAL), which incorporates an annotator selection strategy following data sampling. Our objective is two-fold: (1) to efficiently approximate the full diversity of human judgments, and (2) to assess model performance using annotator-centric metrics, which value minority and majority perspectives equally. We experiment with multiple annotator selection strategies across seven subjective NLP tasks, employing both traditional and novel, human-centered evaluation metrics. Our findings indicate that ACAL improves data efficiency and excels in annotator-centric performance evaluations. However, its success depends on the availability of a sufficiently large and diverse pool of annotators to sample from.

5.1 Introduction

A challenging aspect of natural language understanding (NLU) is the variability of human judgment and interpretation in subjective tasks (e.g., hate speech detection) [302]. In a subjective task, a data sample is typically labeled by a set of annotators, and differences in annotation are reconciled via majority voting, resulting in a single (supposedly, true) “gold label” [393]. However, this approach has been criticized for treating label variation exclusively as noise, which is especially problematic in sensitive subjective tasks [21] since it can lead to the exclusion of minority voices [228].

Subjectivity can be addressed by modeling the full distribution of annotations for each data sample instead of employing gold labels [302]. However, resources for such approaches are scarce, as most datasets do not (yet) make fine-grained annotation details available [61], and representing a full range of perspectives is contingent on obtaining costly annotations from a diverse set of annotators [28].

One way to handle a limited annotation budget is to use Active Learning [350, AL]. Given a pool of unannotated data samples, AL employs a sample selection strategy to obtain maximally informative samples, retrieving the corresponding annotations from a ground truth oracle (e.g., a single human expert). However, in subjective tasks, there is no such oracle. Instead, we rely on a set of available annotators. Demanding all available annotators to annotate all samples would provide a truthful representation of the annotation distribution, but is often unfeasible, especially if the pool of annotators is large. Thus, deciding *which annotator(s)* should annotate is as critical as deciding which samples to annotate.

In most practical applications, annotators are randomly selected. This results in an annotation distribution insensitive to outlier annotators—most annotations reflect the majority voices and fewer reflect the minority voices. This may not be desirable in applications such as hate speech, where the opinions of the majority and minority should be valued equally. In such cases, a more deliberate annotator selection is required. To ensure a balanced representation of majority and minority voices, we leverage strategies inspired by Rawls’ principle of fairness [313], which advocates that a fair society is achieved when the well-being of the worst-off members of society (the minority annotators, in this case) is maximized.

We introduce Annotator-Centric Active Learning (ACAL) to emphasize and control who annotates which sample. In ACAL (Figure 5.1), the sample selection strategy of traditional AL is followed by an *annotator selection strategy*, indicating which of the available annotators should annotate each selected data sample.

Contributions (1) We present ACAL as an extension of the AL approach and introduce three annotator selection strategies aimed at collecting a balanced distribution of minority and majority annotations. (2) We introduce a suite of annotator-centric evaluation metrics to measure how individual and minority annotators are modeled. (3) We demonstrate ACAL’s effectiveness in three datasets with subjective tasks—hate speech detection, moral value classification, and safety judgments.

Our experiments show that the proposed ACAL methods can approximate the distribution of human judgments similar to AL while requiring a lower annotation budget and modeling individual and minority voices more accurately. However, our evaluation shows how the task’s annotator agreement and the number of available annotations impact ACAL’s effectiveness—ACAL is most effective when a large pool of diverse annotators is available.

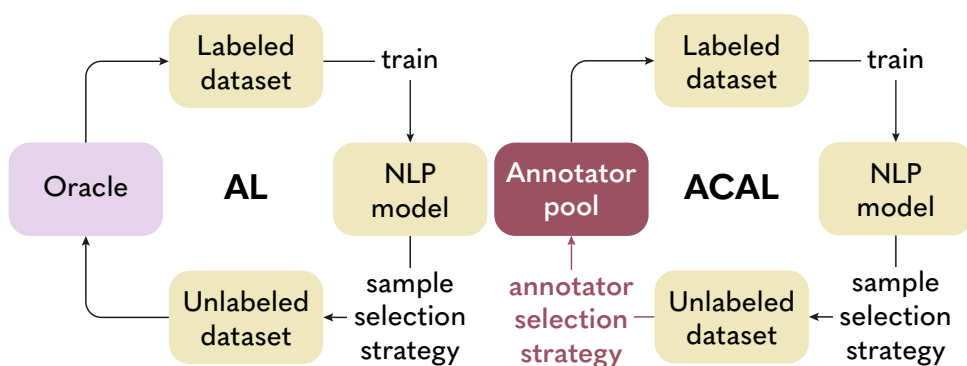


Figure 5.1: Active Learning (AL) approaches (left) use a sample selection strategy to pick samples to be annotated by an oracle. The Annotator-Centric Active Learning (ACAL) approach (right) extends AL by introducing an annotator selection strategy to choose the annotators who annotate the selected samples.

Importantly, our experiments show how the ACAL framework controls how models learn to represent majority and minority annotations. This is crucial for subjective and sensitive applications such as detecting human values and morality [203, 239], argument mining [405], and hate speech [198].

5

5.2 Related work

5.2.1 Learning with annotator disagreement

Modeling annotator disagreement is garnering increasing attention [21, 61, 302, 393]. Changing annotation aggregation methods can lead to a fairer representation than simple majority [171, 380]. Alternatively, the full annotation distribution can be modeled using soft labels [79, 277, 300]. Other approaches leverage annotator-specific information, e.g., by including individual classification heads per annotator [89], embedding annotator behavior [269], or encoding the annotator’s socio-demographic information [44]. Yet, modeling annotator diversity remains challenging. Standard calibration metrics under human label variation may be unsuitable, especially when the variation is high [24]. Trade-offs ought to be made between collecting more samples or more annotations [149]. Further, solely measuring differences among sociodemographic traits is not sufficient to capture opinion diversity [291]. Instead, we represent diversity based on *which* annotators annotated *what* and *how*. We experiment with annotator selection strategies to reveal what aspects impact task performance and annotation budget.

5.2.2 Active Learning

AL enables a supervised learning model to achieve high performance by judiciously choosing a few training examples [350]. In a typical AL scenario, a large collection of unlabeled data is available, and an oracle (e.g., a human expert) is asked to annotate this unlabeled data. A *sampling strategy* is used to iteratively select the next batch of unlabeled data for annotation [316]. AL has found widespread application in NLP [451]. Two main strategies are employed, either by selecting the unlabeled samples on which the model prediction is

most uncertain [450], or by selecting samples that are most representative of the unlabeled dataset [116, 452]. The combination of AL and annotator diversity is a novel direction. Existing works propose to align model and annotator uncertainties [39], adapt annotator-specific classification heads in AL settings [421], or select texts to annotate based on annotator preferences [192]. These methods ignore a crucial part of learning with human variation: the diversity among annotators. We focus on selecting annotators such that they best inform us about the underlying label diversity.

5.3 Method

First, we define the soft-label prediction task we use to train a supervised model. Then, we introduce the traditional AL and the novel ACAL approaches.

5.3.1 Soft-label prediction

Consider a dataset of triples $\{x_i, a_j, y_{ij}\}$, where x_i is a data sample (i.e., a piece of text) and $y_{ij} \in C$ is the class label assigned by annotator a_j . The multiple labels assigned to a sample x_i by the different annotators are usually combined into an aggregated label \hat{y}_i . For training with soft labels (i.e. non-binary class assignment), the aggregation typically takes the form of maximum likelihood estimation [393]:

$$\hat{y}_i(x) = \frac{\sum_{i=1}^N [x_i = x][y_{ij} = c]}{\sum_{i=1}^N [x_i = x]} \quad (5.1)$$

In our experiments, we use a passive learning approach that uses all available $\{x_i, \hat{y}_i\}$ to train a model f_θ with cross-entropy loss as a baseline.

5.3.2 Active Learning

AL imposes a sampling technique for inputs x_i , such that the most *informative* sample(s) are picked for learning. In a typical AL approach, a set of unlabelled data points U is available. At every iteration, a sample selection strategy \mathcal{S} selects samples $x_i \in U$ to be annotated by an oracle \mathcal{O} that provides the ground truth label distribution \hat{y}_i . The selected samples and annotations are added to the labeled data D , with which the model f_θ is trained. Alg. 1 provides an overview of the procedure.

Algorithm 1: AL approach.

input: Unlabeled data U , Data sampling strategy \mathcal{S} , Oracle \mathcal{O}
 $D_0 \leftarrow \{\}$
for $n = 1..N$ **do**
 sample data points x_i from U using \mathcal{S}
 obtain annotation \hat{y}_i for x_i from \mathcal{O} $D_{n+1} = D_n + \{x_i, \hat{y}_i\}$
 train f_θ on D_{n+1}
end

In the sample selection strategies, a batch of data of a given size B is queried at each iteration. Our experiments compare the following strategies:

Random (\mathcal{S}_R) selects a B samples uniformly at random from U .

Uncertainty (\mathcal{S}_U) predicts a distribution over class labels with $f_\theta(x_i)$ for each $x_i \in U$, and selects B samples with the highest prediction entropy (the samples the model is most uncertain about).

5.3.3 Annotator-Centric Active Learning

ACAL builds on AL. In contrast to AL, which retrieves an aggregated annotation \hat{y}_i , ACAL employs an annotator selection strategy \mathcal{T} to select one annotator and their annotation for each selected data point x_i . Alg. 2 describes the ACAL approach.

Algorithm 2: ACAL approach.

```

input: Unlabeled data  $U$ , Data sampling strategy  $\mathcal{S}$ , Annotator sampling strategy  $\mathcal{T}$ 
 $D_0 \leftarrow \{\}$ 
for  $n = 1..N$  do
    sample data points  $x_i$  from  $U$  using  $\mathcal{S}$ 
    sample annotators  $a_j$  for  $x_i$  using  $\mathcal{T}$ 
    obtain annotation  $y_{ij}$  from  $a_j$  for  $x_i$ 
     $D_{n+1} = D_n + \{x_i, y_{ij}\}$ 
    train  $f_\theta$  on  $D_{n+1}$ 
end

```

5

We propose three annotator selection strategies to gather a distribution that uniformly contains all possible (majority and minority) labels, inspired by Rawls' principle of fairness [313]. The strategies vary in the type of information used to represent differences between annotators, including *what* or *how* the annotators have annotated thus far. Our experiments compare the following strategies:

Random (\mathcal{T}_R) randomly selects an annotator a_j .

Label Minority (\mathcal{T}_L) considers only information on *how* each annotator has annotated so far (i.e., the labels that they have assigned). The minority label is selected as the class with the smallest annotation count in the available dataset D_n thus far. Given a new sample, x_i , \mathcal{T}_L selects the available annotator that has the largest bias toward the minority label compared to the other available annotators, i.e., who has annotated other samples with the minority label the most.

Semantic Diversity (\mathcal{T}_S) considers only information on *what* each annotator has annotated so far (i.e., the samples that they have annotated). Given a new sample x_i selected through \mathcal{S} , \mathcal{T}_S selects the available annotator for whom x_i is semantically the most different from what the annotator has labeled so far. To measure this difference for an annotator a_j , we employ a sentence embedding model to measure the cosine distance between the embeddings of x_i and embeddings of all the samples annotated by a_j . We then take the average of all semantic similarities. The annotator with the lowest average similarity score is selected.

Representation Diversity (\mathcal{T}_D) selects the annotator that has the lowest similarity on average with all other annotators available for that item. We create a representation for each annotator by averaging the embeddings of samples annotated by a_j together with their respective labels, followed by computing the pair-wise cosine similarity between all annotators.

5.4 Experimental Setup

We describe the experimental setup for the comparisons between ACAL strategies. In all our experiments, we employ a TinyBERT model [187] to reduce the number of trainable parameters. Appendix D.1 includes a detailed overview of the computational setup and hyperparameters. We make the code for the ACAL strategies and evaluation metrics available via GitHub.¹

5.4.1 Datasets

We use three datasets which vary in domain, annotation task (in *italics*), annotator count, and annotations per instance.

The **DICES Corpus** [22] is composed of 990 conversations with an LLM where 172 annotators provided judgments on whether a generated response can be deemed safe (3-way judgments: yes, no, unsure). Samples have 73 annotations on average. We perform a multi-class classification of the judgments.

The **MFTC Corpus** [169] is composed of 35K tweets that 23 annotators annotated with any of the 10 moral elements from the Moral Foundation Theory [142]. We select the elements of *loyalty* (lowest annotation count), *care* (average count), and *betrayal* (highest count). Samples have 4 annotations on average. We create three binary classifications to predict the presence of the respective elements. As most tweets were labeled as non-moral (i.e., with no moral element), we balanced the datasets by subsampling the non-moral class.

The **MHS Corpus** [328] consists of 50K social media comments on which 8K annotators judged three hate speech aspects—*dehumanize* (low inter-rater agreement), *respect* (medium agreement), and *genocide* (high agreement)—on a 5-point Likert scale. Samples have 3 annotations on average. We perform a multi-class classification with the annotated Likert scores for each task.

The datasets and tasks differ in levels of annotator agreement, measured via entropy of the annotation distribution. DICES and MHS generally have medium entropy scores, whereas the MFTC entropy is highly polarized (divided between samples with very high and very low agreement). Appendix D.1.5 provides details of the entropy scores.

5.4.2 Evaluation metrics

The ACAL strategies aim to guide the model to learn a representative distribution of the annotator’s perspectives while reducing annotation effort. To this end, we evaluate the model both with a traditional evaluation metric and a metric aimed at comparing predicted and annotated distributions:

Macro F_1 -score (F_1) For each sample in the test set, we select the label predicted by the model with the highest confidence, determine the golden label through a majority agreement aggregation, and compute the resulting macro F_1 -score.

Jensen-Shannon Divergence (JS) The JS measures the divergence between the distribution of label annotation and prediction [286]. We report the average JS for the samples in the test set to measure how well the model can represent the annotation distribution.

¹<https://github.com/m0re4u/acal-subjective>

Further, since ACAL shifts the focus to annotators, we introduce novel annotator-centric evaluation metrics. First, we report the average among annotators. Second, in line with Rawls' principle of fairness, the result for the worst-off annotators:

Per-annotator F_1 (F_1^a) and JS (JS^a) We compute the F_1 (or JS) for each annotator in the test set using their annotations as golden labels (or target distribution), and average it.

Worst per-annotator F_1 (F_1^w) and JS (JS^w) We compute the F_1 (or JS) for each annotator in the test set using their annotations as golden labels (or target distribution), and report the average of the lowest 10% to mitigate noise.

These metrics allow us to measure the trade-offs between modeling the majority agreement, a representative distribution of annotations, and accounting for minority voices. In the next section, we describe how we obtained the results.

5.4.3 Training procedure

We test the annotator selection strategies proposed in Section 5.3.3 by comparing all combinations of the two sample selection strategies (\mathcal{S}_R and \mathcal{S}_U) and the four annotator selection strategies (\mathcal{T}_R , \mathcal{T}_L , \mathcal{T}_S , and \mathcal{T}_D). At each iteration, we use \mathcal{S} to select B unique samples from the unlabeled data pool U . We select B as the smallest between 5% of the number of available annotations and the number of unique samples in the training set. For each selected sample x_i , we use \mathcal{T} to select one annotator and retrieve their annotation y_{ij} .

We split each dataset into 80% train, 10% validation, and 10% test. We start the training procedure with a warmup iteration of B randomly selected annotations [451]. We proceed with the ACAL iterations by combining \mathcal{S} and \mathcal{T} . We select the model checkpoint across all AL iterations that led to the best JS performance on the validation set and evaluate it on the test set. We repeat this process across three data splits and model initializations. We report the average scores on the test set.

We compare ACAL with traditional oracle-based AL approaches ($\mathcal{S}_R\mathcal{O}$ and $\mathcal{S}_U\mathcal{O}$), which use the data sampling strategies but obtain all possible annotations for each sample as in Alg. 1. Further, we employ a passive learning (PL) approach as an upper bound by training the model on the full dataset, thus observing all available samples and annotations. Similar to ACAL, the AL and PL baselines are averaged over three seeds.

5.5 Results

We start by highlighting the benefits of ACAL over AL and PL (Section 5.5.1). Next, we closely examine ACAL on efficiency and fairness (Section 5.5.2). Then, we select a few cases of interest and dive deeper into the strategies' behavior during training (Section 5.5.3). Finally, we investigate ACAL across varying levels of subjectivity (Section 5.5.4).

5.5.1 Highlights

Our experiments show that ACAL can have a beneficial impact over using PL and AL. Figure 5.2 highlights two main findings: (1) ACAL strategies can more quickly learn to represent the annotation distribution with a large pool of annotators, and (2) when agreement between annotators is polarized, ACAL leads to improved results compared to learning from aggregated labels. In the next sections, we provide a deeper understanding of the conditions in which ACAL works well.

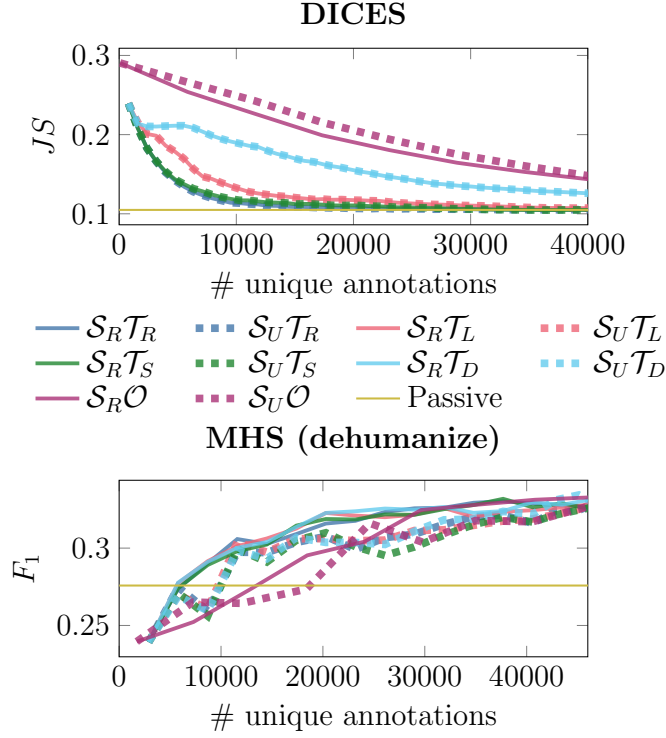


Figure 5.2: Learning curves showing model performance on the validation set. On DICES (upper), ACAL approaches are quicker than AL in obtaining similar performance to passive learning. On MHS (lower), ACAL surpasses passive learning in F_1 when data has high disagreement.

5.5.2 Efficiency and Fairness

Table 5.1 presents the results of evaluating the best models (those with the highest JS scores on the validation set) on the test set. We analyze the results along two dimensions: (a) *efficiency*: what is the impact of the different strategies on the trade-off between annotation budget and performance? (b) *fairness*: do the selection strategies that aim for a balanced consideration of minority and majority views lead to better performance in the human-centric evaluation metrics? For MFTC we focus on *care* because it has an average number of samples available, and for MHS we focus on *dehumanize* because it has high levels of disagreement. Appendix C.3 presents the remainder of the results.

Efficiency We discuss the performance on F_1 and JS to measure how well the proposed strategies model label distributions and examine the used annotator budget. Across all tasks and datasets, ACAL and AL consistently yield comparable or superior F_1 and JS with a lower annotation budget than PL. When comparing ACAL with AL, the results vary depending on the task and dataset. For DICES, there is a significant benefit to using ACAL, as it can save up to $\sim 40\%$ of the annotation budget while yielding better scores across all metrics than AL. With AL, we observe only a small reduction in annotation cost. For MFTC, AL with

	App.	F_1	JS	Average		Worst-off		$\Delta\%$
				F_1^a	JS^a	F_1^w	JS^w	
DICES	$\mathcal{S}_R \mathcal{T}_R$	53.2	.100	43.2	.186	16.7	.453	-36.8
	$\mathcal{S}_R \mathcal{T}_L$	55.5	.101	42.4	.187	15.5	.450	-32.7
	$\mathcal{S}_R \mathcal{T}_S$	61.0	.103	44.2	.186	16.4	.447	-35.5
	$\mathcal{S}_R \mathcal{T}_D$	58.9	.142	43.1	.203	16.9	.370	-30.0
	$\mathcal{S}_U \mathcal{T}_R$	53.2	.100	43.2	.186	16.7	.453	-36.8
	$\mathcal{S}_U \mathcal{T}_L$	55.5	.101	42.4	.187	15.5	.450	-32.7
	$\mathcal{S}_U \mathcal{T}_S$	63.1	.098	43.9	.187	18.4	.447	-38.2
	$\mathcal{S}_U \mathcal{T}_D$	58.9	.142	43.1	.203	16.9	.370	-30.0
	$\mathcal{S}_R \mathcal{O}$	59.1	.112	41.4	.191	13.3	.425	-0.1
	$\mathcal{S}_U \mathcal{O}$	46.2	.110	38.4	.192	11.7	.427	-0.1
MFTC (<i>care</i>)	PL	59.0	.105	37.1	.211	12.3	.479	-
	$\mathcal{S}_R \mathcal{T}_R$	78.9	.038	61.1	.141	37.7	.247	-1.6
	$\mathcal{S}_R \mathcal{T}_L$	78.5	.037	61.6	.142	39.2	.249	-0.4
	$\mathcal{S}_R \mathcal{T}_S$	78.1	.039	60.0	.145	35.1	.248	-1.7
	$\mathcal{S}_R \mathcal{T}_D$	76.6	.040	60.4	.144	35.7	.243	-1.7
	$\mathcal{S}_U \mathcal{T}_R$	79.4	.038	61.2	.143	37.7	.252	-5.6
	$\mathcal{S}_U \mathcal{T}_L$	80.7	.037	58.9	.142	42.3	.248	-2.5
	$\mathcal{S}_U \mathcal{T}_S$	79.1	.037	60.8	.143	39.9	.258	-1.1
	$\mathcal{S}_U \mathcal{T}_D$	78.1	.040	58.6	.145	35.7	.253	-2.5
	$\mathcal{S}_R \mathcal{O}$	79.0	.037	58.6	.141	39.2	.255	-0.2
MHS (<i>dehumanize</i>)	$\mathcal{S}_U \mathcal{O}$	79.4	.037	58.3	.144	35.7	.253	-12.7
	PL	81.1	.032	51.2	.179	37.7	.251	-
	$\mathcal{S}_R \mathcal{T}_R$	33.6	.081	31.5	.394	0.0	.489	-50.0
	$\mathcal{S}_R \mathcal{T}_L$	33.1	.081	32.2	.397	0.0	.478	-62.5
	$\mathcal{S}_R \mathcal{T}_S$	30.5	.079	31.3	.397	0.0	.480	-62.5
	$\mathcal{S}_R \mathcal{T}_D$	32.4	.081	31.8	.398	0.0	.479	-62.5
	$\mathcal{S}_U \mathcal{T}_R$	32.4	.080	32.2	.389	0.0	.508	-7.8
	$\mathcal{S}_U \mathcal{T}_L$	33.1	.080	32.8	.388	0.0	.507	-7.8
	$\mathcal{S}_U \mathcal{T}_S$	33.6	.080	32.6	.388	0.0	.506	-7.8
	$\mathcal{S}_U \mathcal{T}_D$	33.0	.079	32.6	.384	0.0	.513	-3.0
	$\mathcal{S}_R \mathcal{O}$	32.8	.077	33.9	.387	0.0	.496	-60.1
	$\mathcal{S}_U \mathcal{O}$	33.3	.080	33.1	.390	0.0	.497	-24.7
	PL	28.0	.075	20.2	.424	0.0	.547	-

Table 5.1: Test set results on the DICES, MFTC (*care*), and MHS (*dehumanize*) datasets. Results report the average test scores from the best-performing model checkpoint on the validation set (lowest JS), evaluated across three data splits and model initializations. $\Delta\%$ denotes the reduction in the annotation budget with respect to passive learning. In bold, the best performance per column and per dataset (higher F_1 are better, lower JS are better).

\mathcal{S}_U leads to the largest cost benefits ($\sim 12\%$ less annotation budget), but at a cost in terms of absolute JS and F_1 . ACAL slightly outperforms AL but does not lead to a decrease in annotation budget. For MHS, both AL and ACAL significantly reduce the annotation cost ($\sim 60\%$) while yielding better scores than PL—however, AL and ACAL do not show substantial performance differences. Overall, when looking at F_1 and JS which are aggregated over

the whole test set, we conclude that ACAL is most efficient when the pool of available annotators for one sample is large (as with the DICES dataset), whereas the difference between ACAL and AL is negligible with a small pool of annotators per data sample (as with MFTC and MHS).

Fairness We investigate the extent to which the models represent individual annotators fairly and capture minority opinions via the annotator-centric evaluation metrics (F_1^a , JS^a , F_1^w , and JS^w). We observe a substantial improvement when using AL or ACAL over PL. Further, we observe no single winner-takes-all approach: high F_1 and JS scores do not consistently co-occur with high scores for the annotator-centric metrics. This highlights the need for a more comprehensive evaluation to assess models for subjective tasks. Yet, we observe that ACAL slightly outperforms AL in modeling individual annotators (JS^a and F_1^a). This trend is particularly evident with DICES, again likely due to the large pool of annotators available per data sample. Lastly, ACAL is best in the worst-off metrics (JS^w and F_1^w), showing the ability to better represent minority opinions as a direct consequence of the proposed annotator selection strategies on DICES and MFTC. However, all approaches score 0 for F_1^w on MHS. This is due to the high disagreement in this dataset: the 10% worst-off annotators always disagree with a hard label derived from the predicted label distribution. In conclusion, our experiments show that, when a large pool of annotators is available, a targeted sampling of annotators requires fewer annotations and is fairer. That is, minority opinions are better represented without large sacrifices in performance compared to the overall label distribution.

5

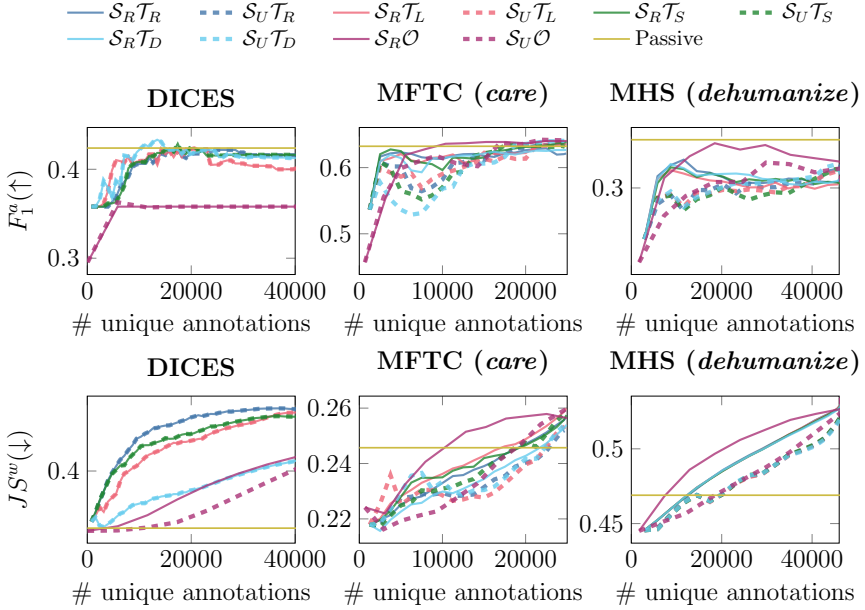


Figure 5.3: Selected plots showing the F_1^a and JS^w performance on the validation set during the ACAL and AL iterations for DICES, MFTC (*care*), and MHS (*dehumanize*). Higher F_1^a is better, lower JS^w is better. Y-axes are scaled to highlight the relative performance to PL.

5.5.3 Convergence

The evaluation on the test set paints a general picture of the advantage of using ACAL over AL or PL. In this section, we assess how different ACAL strategies converge over iterations. We describe the major patterns across our experiments by analyzing six examples of interest with F_1^a and JS^w (Figure 5.3). We select F_1^a because it reveals how well individual annotators are modeled on average, and JS^w to measure how strategies deviate from modeling the majority perspective. Appendix D.2.2 provides an overview of all metrics.

First, we notice that the trends for F_1^a and JS^w are both increasing—the first is expected, but the second requires an explanation. As the model is exposed to more annotations over the training iterations, the predicted label distribution starts to fit the true label distribution. However, here we consider each annotator individually: JS^w reports the average of the 10% lowest JS scores per annotator. The presence of disagreement implies the existence of annotators that annotate differently from the majority. Since our models predict the full distribution, they assign a proportional probability to dissenting annotators. Thus, learning to model the full distribution of annotations leads to an increase in JS^w .

Second, we notice a difference between ACAL and AL. On MFTC and MHS, ACAL, compared to AL, yields overall smaller JS^w at the cost of a slower convergence in F_1^a , showing the trade-off between modeling all annotators and representing minorities. However, with DICES the trend is the opposite. This is due to AL having access to the complete label distribution: it can model a balanced distribution, leading to lower worst-off performance. With a large number of annotations, ACAL requires more iterations to get the same balanced predicted distribution.

Third, we observe differences among the annotator selection strategies (\mathcal{T}). \mathcal{T}_D shows the most differences—both JS^w and F_1^a increase slower than for the other strategies. This suggests that selecting annotators based on the average embedding of the annotated content strongest emphasizes diverging label behavior.

Finally, we analyze the impact of the sample selection strategies (\mathcal{S} , dotted vs. solid lines in Figure 5.3). For DICES, \mathcal{S}_R and \mathcal{S}_U lead to comparable results, likely due to the low number of samples. Using \mathcal{S}_U in MFTC leads to F_1^a performance decreasing at the start of training. The strategy prioritizes obtaining annotations for already added samples to lower their entropy, while the variation in labels is irreconcilable (since there are limited labels available, and they are in disagreement). We see a similar pattern for MHS.

These results further underline our main finding that ACAL is effective in representing diverse annotation perspectives when there is a (1) heterogeneous pool of annotators, and (2) a task that facilitates human label variation.

5.5.4 Impact of subjectivity

We further investigate ACAL strategies on (1) label entropy, and (2) cross-task performance.

Alignment of ACAL strategies during training We want to investigate how well the ACAL strategies align with the overall subjective annotations: do they drive the model entropy in the right direction? We measure the entropy of the samples in the labeled training set at each iteration and compare it to the entropy of all annotations of those samples. Higher entropy in the labeled training set than the actual entropy suggests that the selection strategy overestimates uncertainty. Lower entropy indicates that the model may not sufficiently account

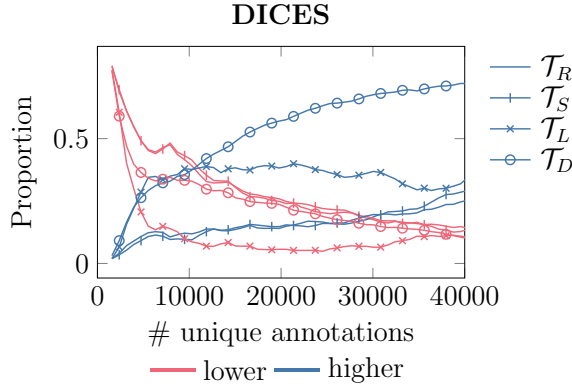


Figure 5.4: Proportion of data samples that result in higher or lower entropy than the target label distribution per ACAL strategy.

5

for disagreement. When the entropy matches the true entropy, the selection strategy is well-calibrated to strike a healthy middle ground between sampling diverse labels and finding the majority class. We focus on DICES as a case study due to the wide range of entropy scores. We group each sample based on the true label entropy into low (< 0.43), medium ($0.43 - 0.72$), and high (> 0.72). We apply the same categorization at each training iteration for samples labeled thus far. Subsequently, we plot the proportion of data points for which the selection strategy results in excessively high or excessively low entropy.

Figure 5.4 visualizes the proportions. At the beginning of training, entropy is generally low because samples have few annotations. Over time, the selected annotations better align with the true entropy. At the start (at 10K unique annotations), roughly only a third of the samples have aligned entropy scores ($T_R = 27\%$, $T_S = 27\%$, $T_L = 33\%$, $T_D = 32\%$). Further towards the end of the ACAL iterations, this has increased for all ACAL strategies except T_D ($T_R = 64\%$, $T_S = 62\%$, $T_L = 57\%$, $T_D = 17\%$). When and how much the strategies succeed in matching the true label distribution differs: T_S and T_R take longer to increase label entropy than the other two strategies. They are conservative in adding diverse labels. T_L and T_D increase the proportion of well-aligned data points earlier in the training process, achieving a balanced entropy alignment sooner. However, both strategies start to overshoot the target entropy, whereas the others show a more gradual alignment with the true entropy. This effect is strongest for T_D . This finding suggests that minority-aware annotator-selection (T_L and T_D) strategies achieve the best results in the early stages of training—that is, they are effective for quickly raising entropy but can lead to overrepresentation.

Cross-task performance Figure 5.5 compares the two annotator-centric metrics on the three tasks of MFTC and MHS—the datasets for which we have seen the least impact of ACAL over AL and PL. We select a data sampling (S_R) and annotator sampling strategy (T_S), based on its strong performance on DICES for comprehensive comparison.

When evaluating MFTC *loyalty*, which has the highest disagreement, JS^w is more accurately approximated with PL. Similarly, ACAL is outperformed by AL on F_1^a for the *de-humanize* (high disagreement) task. However, for the less subjective task *genocide*, ACAL

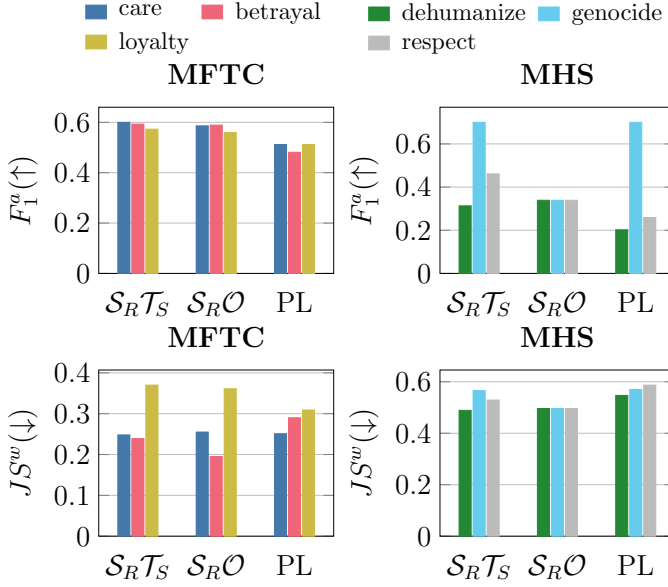


Figure 5.5: Comparison of ACAL, AL, and PL across different MFTC and MHS tasks. Higher F_1^a is better, and lower JS^w is better.

leads to higher F_1^a . This suggests that the effectiveness of annotation strategies varies depending on the task's degree of subjectivity *and* the available pool of annotators. The more heterogeneous the annotation behavior, indicative of a highly subjective task, the larger the pool of annotators required for each sample selection. We also observe that there is a trade-off between modeling the majority of annotators equally (F_1^a) and prioritizing the minority (JS^w).

5.6 Conclusion

We present ACAL as an extension of AL to emphasize the selection of diverse annotators. We introduce three novel annotator selection strategies and four annotator-centric metrics and experiment with tasks across three different datasets. We find that the ACAL approach is especially effective in reducing the annotation budget when the pool of available annotators is large. However, its effectiveness is contingent on data characteristics such as the number of annotations per sample, the number of annotations per annotator, and the nature of disagreement in the task annotations. Furthermore, our novel evaluation metrics display the trade-off between modeling overall distributions of annotations and adequately accounting for minority voices, showing that different strategies can be tailored to meet different goals. Especially early in the training process, strategies that are aggressive in obtaining diverse labels have a beneficial impact in accounting for minority voices. However, we recognize that gathering a distribution that uniformly contains all possible (minority and majority) labels can be overly sensitive to small minorities or noise. Future work should integrate methods that account for noisy annotations [426]. Striking a balance between utilitarian and egalitarian

tarian approaches, such as between modeling aggregated distributions and accounting for minority voices [229] is crucial for inferring context-dependent values [242, 400].

Limitations

The main limitation of this work is that the experiments are based on simulated AL which is known to bear several shortcomings [261]. In our study, a primary challenge arises with two of the datasets (MFTC, MHS), which, despite having a large pool of annotators, lack annotations from every annotator for each item. Consequently, in real-world scenarios, the annotator selection strategies for these datasets would benefit from access to a more extensive pool of annotators. This limitation likely contributes to the underperformance of ACAL on these datasets compared to DICES. We emphasize the need for more datasets that feature a greater number of annotations per item, as this would significantly enhance research efforts aimed at modeling human disagreement.

Since we evaluate four different annotator selection strategies and two sample selection strategies across three datasets and seven tasks, the amount of experiments is high. This did not allow for further investigation of other methods for measuring uncertainty such as ensemble methods [218], different classification models, the extensive turning of hyperparameters, or even different training paradigms like low-rank adaptation [173]. Lastly, a limitation of our annotator selection strategies is that they rely on a small annotation history. This is why we require a warmup phase for some of the strategies, for which we decided to take a random sample of annotations. Incorporating informed warmup strategies, incorporating ACAL strategies that do not rely on annotator history, or making use of more elaborate hybrid human-AI approaches [403] may positively impact its performance and data efficiency.

Ethical Considerations

Our goal is to approximate a good representation of human judgments over subjective tasks. We want to highlight the fact that the *performance* of the models differs a lot depending on which metric is used. We tried to account for a less majority-focussed view when evaluating the models which is very important, especially for more human-centered applications, such as hate-speech detection. However, the evaluation metrics we use do not fully capture the diversity of human *judgments*, but just that of *labeling behavior*. The selection of metrics should align with the specific goals and motivations of the application, and there is a pressing need to develop more metrics to accurately reflect human variability in these tasks.

Our experiments are conducted on English datasets due to the scarcity of unaggregated datasets in other languages. In principle, ACAL can be applied to other languages (given the availability of multilingual models to semantically embed textual items for some particular strategies used in this work). We encourage the community to enrich the dataset landscape by incorporating more perspective-oriented datasets in various languages, ACAL potentially offers a more efficient method for creating such datasets in real-world scenarios.