

A Unified Evaluation Framework for Multi-Annotator Tendency Learning

Liyun Zhang

D3 Center, The University of Osaka
Japan

Jingcheng Ke

D3 Center, The University of Osaka
Japan

Shenli Fan

Business Administration, Osaka
University of Economics and Law
Japan

Xuanmeng Sha
IST, The University of Osaka
Japan

Zheng Lian
Institute of automation, Chinese
academy of science
China

ABSTRACT

Recent works have emerged in multi-annotator learning that shift focus from Consensus-oriented Learning (CoL), which aggregates multiple annotations into a single ground-truth prediction, to Individual Tendency Learning (ITL), which models annotator-specific labeling behavior patterns (i.e., tendency) to provide explanation analysis for understanding annotator decisions. However, no evaluation framework currently exists to assess whether ITL methods truly capture individual tendencies and provide meaningful behavioral explanations. To address this gap, we propose the first unified evaluation framework with two novel metrics: (1) Difference of Inter-annotator Consistency (DIC) quantifies how well models capture annotator tendencies by comparing predicted inter-annotator similarity structures with ground-truth; (2) Behavior Alignment Explainability (BAE) evaluates how well model explanations reflect annotator behavior and decision relevance by aligning explainability-derived with ground-truth labeling similarity structures via Multidimensional Scaling (MDS). Extensive experiments validate the effectiveness of our proposed evaluation framework.

CCS CONCEPTS

• Computing methodologies → Multi-task learning; Simulation evaluation.

KEYWORDS

Evaluation Framework, Evaluation Metrics, Individual Tendency Learning, Behavioral Explainability, Annotator Tendencies, Multi-annotator Learning

1 INTRODUCTION

In real-world multi-annotation scenarios, such as medical image analysis [19], sentiment analysis [18], and visual perception [45],

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

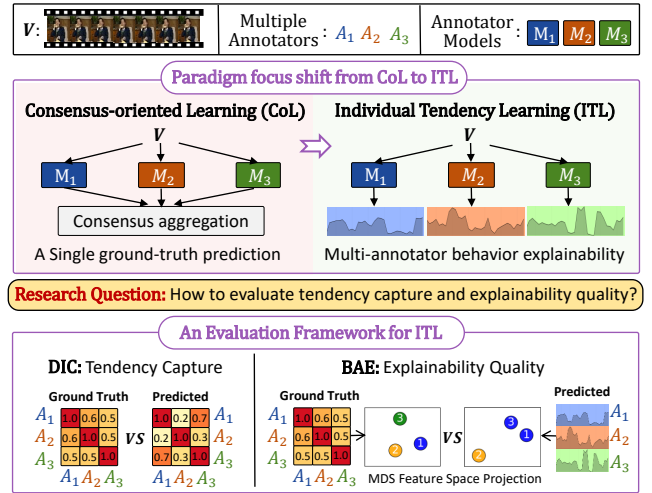


Figure 1: Top: A paradigm focus shift from Consensus-oriented Learning (CoL) to Individual Tendency Learning (ITL). A video sample V is processed by multi-annotator models. CoL aggregates annotators' predictions into a single ground-truth prediction. While ITL models annotator-specific labeling behavior pattern (i.e., *tendency*) to give different attention change explanations along video frames for understanding annotator decisions. Bottom: An evaluation framework for ITL: Difference of Inter-annotator Consistency (DIC) quantifies how well the model captures annotator tendencies by comparing the structure of predicted inter-annotator similarities with ground-truth; Behavior Alignment Explainability (BAE) evaluates how well model explanations reflect annotator behavior and decision relevance by aligning explainability-derived with ground-truth labeling similarity structures via Multidimensional Scaling (MDS).

different annotators often provide different labels to the same sample [23] due to varying personal backgrounds, subjective interpretations, and expertise levels. These systematic patterns in labeling behavior—which we term *tendencies*—represent valuable information about individual perspectives, cognitive biases, and domain expertise.

As illustrated in Figure 1 (top left), the dominant paradigm in multi-annotator learning has been Consensus-oriented Learning (CoL), treating annotator disagreements as noise [17, 35] to be averaged away through aggregation techniques like PADL [19] and MaDL [16]. Recent works have emerged in multi-annotator learning that shift focus from CoL to Individual Tendency Learning (ITL) (Figure 1, top right) with sophisticated architectures like QuMAB [38] and TAX [6] that treat annotator diversity as valuable information to model individual annotators for capturing their unique tendencies and providing behavioral explainability.

However, despite promoting architectural innovations in ITL, there exists no principled evaluation framework to assess whether these models effectively capture annotator-specific tendencies and provide meaningful explanations to understand annotator behavior and decision relevance. This gap limits our understanding of which ITL approaches genuinely capture annotator tendencies, as opposed to those that merely optimize for consensus accuracy.

Furthermore, existing explainability assessments in multi-annotator systems rely primarily on qualitative analysis, lacking quantitative metrics to evaluate whether learned explanations accurately reflect genuine annotator behavioral relationships. A critical question remains unanswered: How can we systematically evaluate which methods truly preserve individual annotator tendencies and provide meaningful behavioral explanations? To address these fundamental evaluation challenges, we propose a unified evaluation framework that systematically assesses ITL methods across both tendency capture and explainability quality dimensions.

The *Difference of Inter-annotator Consistency (DIC)* metric (Figure 1, bottom left) quantifies how well the model captures annotator tendencies by comparing the structure of predicted inter-annotator similarities with the ground truth. The key insight is that if an ITL method truly preserves individual tendencies, the patterns of agreement and disagreement between annotators in predictions should mirror those in actual annotations.

The *Behavior Alignment Explainability (BAE)* metric (Figure 1, bottom right) evaluates how well model explanations reflect annotator behavior and decision relevance by aligning explainability-derived with ground-truth labeling similarity structures via Multidimensional Scaling (MDS). BAE operates through two complementary assessments: feature-level evaluation applicable to all ITL methods through learned representations, and region-level evaluation for attention-based methods through spatial/temporal attention patterns, with Figure 1 (bottom right) demonstrating region-level evaluation. This enables systematic comparison of explainability quality across different architectural approaches. Our work makes the following contributions:

- **A unified evaluation framework for Individual Tendency Learning (ITL):** We address the fundamental evaluation gap in ITL by introducing the first assessment framework that quantifies both tendency capture and explainability quality, enabling systematic comparison of ITL methods’ effectiveness in maintaining annotator diversity while offering meaningful behavioral insights.
- **A novel metric for tendency capture evaluation:** We propose Difference of Inter-annotator Consistency (DIC), which quantifies how well the model captures annotator

tendencies by comparing the structure of predicted inter-annotator similarities with the ground truth.

- **A novel metric for explainability quality evaluation:** We propose Behavior Alignment Explainability (BAE), which evaluates how well model explanations reflect annotator behavior and decision relevance by aligning explainability-derived with ground-truth labeling similarity structures via Multidimensional Scaling (MDS).

2 RELATED WORK

2.1 Multi-annotator Learning Paradigms

The dominant approach in multi-annotator learning has been consensus-oriented, aggregating diverse annotations into a single ground truth. Methods evolved from simple majority voting to sophisticated techniques using probabilistic modeling [8], EM inference [23, 33], and deep learning [1, 24, 32]. Advanced approaches model annotator characteristics—reliability, expertise, and confusion patterns—to improve consensus quality [3, 5, 28, 29, 31]. Recent advances explored nuanced annotator modeling: D-LEMA [20] ensembles annotator learners, PADL [19] fits Gaussian distributions, and MaDL [16] models confusion matrices. SimLabel [39] addresses the practical challenge of missing labels in multi-annotator scenarios. Despite architectural innovations, these methods remain consensus-oriented, treating disagreements as noise to be averaged away rather than as valuable information to be modeled. The underlying computer vision and machine learning techniques used in multi-annotator learning have broader applications across various domains [25, 36, 40, 42–44], though annotator disagreements in multi-annotator scenarios reflect genuine perspective differences rather than noise. Some works recognize annotator diversity’s value [26], particularly in subjective domains where no absolute ground truth exists. QuMAB [38] uses learnable queries to model individual patterns with interpretability. However, these efforts lack systematic evaluation frameworks to assess tendency capture and behavioral explanation quality.

2.2 Tendency Capture Evaluation Metrics

Evaluation in multi-annotator learning has predominantly focused on consensus accuracy, measuring how well methods predict majority votes or expert-defined ground truth [4, 27]. Individual annotator modeling is typically assessed through prediction accuracy or log-likelihood on held-out labels [23, 34]. Traditional agreement metrics—Fleiss’ kappa [12], Krippendorff’s alpha [15], and Cohen’s kappa [30]—measure inter-annotator agreement in raw annotations but do not assess how well computational methods preserve these patterns during learning. Recent approaches focused on annotator modeling quality through likelihood-based metrics [22] or behavior correlations [10, 14]. However, these assess individual behaviors rather than inter-annotator relationship structures. Existing metrics either evaluate consensus quality or individual accuracy, but none quantify how well methods maintain inter-annotator relationships that encode perspective diversity. Our proposed Difference of Inter-annotator Consistency (DIC) metric addresses this gap by measuring tendency capture through consistency pattern matching.

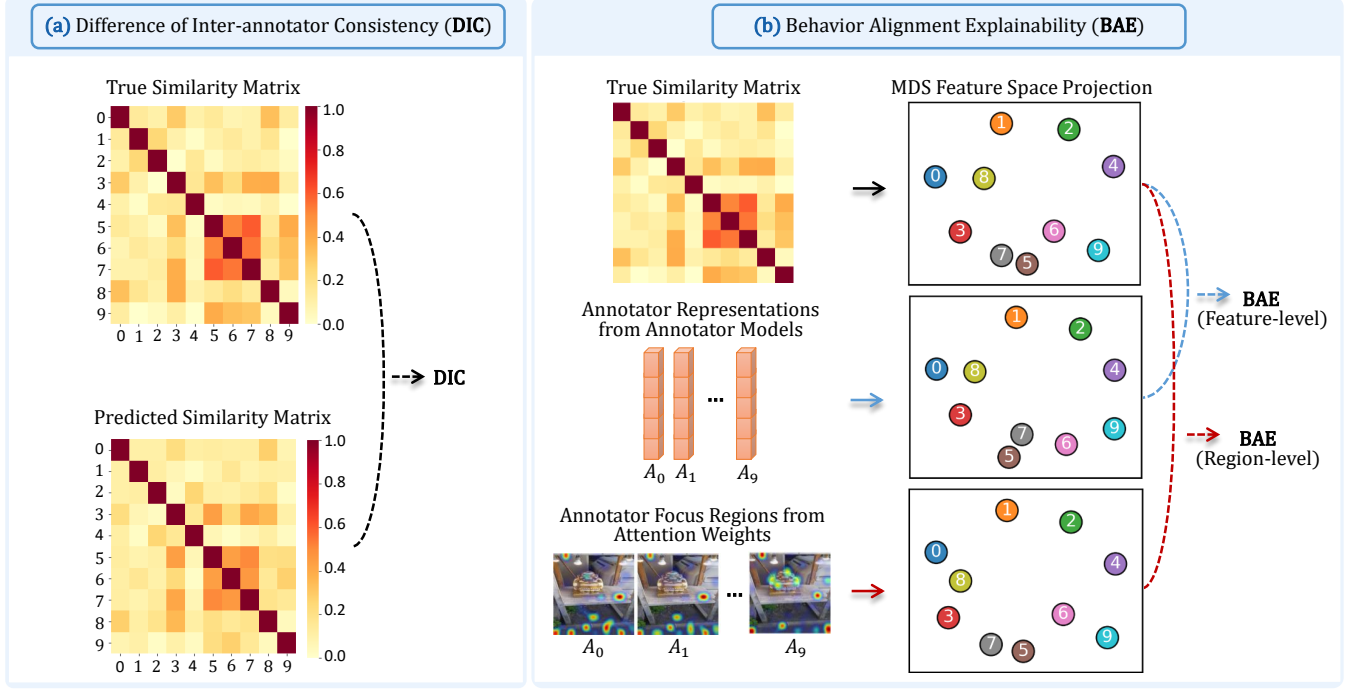


Figure 2: Proposed evaluation framework for inter-annotator behavioral analysis. (a) Difference of Inter-annotator Consistency (DIC) quantifies how well a model preserves annotator tendency by comparing ground-truth and predicted similarity matrices using Frobenius norm. (b) Behavior Alignment Explainability (BAE) assesses whether model explanations capture true inter-annotator behavioral structures using Multidimensional Scaling (MDS) projection. BAE is computed at two complementary levels: *feature-level*, based on learned annotator representations, and *region-level*, based on attention-derived focus regions (for attention-based models). Both measure alignment against the ground-truth consistency matrix.

2.3 Multi-annotator Explainability Assessment

Explainability in multi-annotator learning has received limited systematic attention, with most works providing qualitative rather than quantitative evaluation. Existing approaches include feature-based explanations using learned representations (PADL [19], MaDL [16]) and attention-based explanations showing spatial/temporal patterns (QuMAB [38], TAX [6]), but rely on subjective interpretation without validation frameworks. Current evaluation faces limitations: broader XAI metrics [7, 21] target single-model scenarios and cannot assess inter-annotator relationships, while human evaluation [11] remains expensive and difficult to standardize. Our proposed Behavior Alignment Explainability (BAE) metric addresses this gap, a quantitative assessment applicable to both explanation types. By comparing explanation-derived similarities with ground-truth behavioral patterns, BAE systematically evaluates whether explanations capture genuine annotator relationships regardless of the underlying mechanism.

3 METHODOLOGY

To evaluate tendency capture and explainability, we propose two complementary metrics: Difference of Inter-annotator Consistency (DIC) and Behavior Alignment Explainability (BAE), which assess both tendency capture and explainability quality, as illustrated in Figure 2.

3.1 Difference of Inter-annotator Consistency (DIC)

The fundamental challenge in evaluating tendency capture lies in quantifying how well a model maintains the complex web of inter-annotator relationships. We propose DIC as a principled metric that captures this preservation through consistency pattern matching.

Core Principle. If a model truly preserves annotator tendencies, the patterns of agreement and disagreement between annotators in predictions should mirror those in ground-truth annotations. This forms the theoretical foundation of our consistency-based evaluation approach.

Mathematical Formulation. Given annotations from M annotators, let $Y_k = \{y_i^{(k)} : i \in S_k\}$ and $Y_l = \{y_i^{(l)} : i \in S_l\}$ denote the complete annotation sets for annotators k and l , where S_k represents samples labeled by annotator k . To ensure statistical reliability, we compute consistency only on overlapping samples $S_{kl} = S_k \cap S_l$ with $|S_{kl}| \geq \tau$ (minimum threshold).

The ground-truth inter-annotator consistency matrix $\mathbf{M} \in \mathbb{R}^{M \times M}$ has elements:

$$m_{kl} = \kappa(Y_k|_{S_{kl}}, Y_l|_{S_{kl}}) \quad (1)$$

where $\kappa(\cdot, \cdot)$ denotes Cohen’s kappa coefficient. Similarly, the predicted consistency matrix $\mathbf{M}' \in \mathbb{R}^{M \times M}$ is computed as:

$$m'_{kl} = \kappa(\hat{Y}_k |_{S_{kl}}, \hat{Y}_l |_{S_{kl}}) \quad (2)$$

where $\hat{Y}_k = \{f_{\theta}(x_i, k) : i \in S_k\}$ represents predicted annotations from model f_{θ} .

As illustrated in Figure 2(a), the DIC metric quantifies the preservation fidelity through direct matrix-level comparison:

$$\text{DIC} = \frac{\|\mathbf{M} - \mathbf{M}'\|_F}{\|\mathbf{M}\|_F} \quad (3)$$

where normalization ensures scale-invariance across datasets. Lower DIC indicates better tendency capture, with random assignments producing high values (0.86-0.93) due to inconsistent prediction patterns, while effective methods achieve substantially lower scores by maintaining inter-annotator structural relationships.

3.2 Behavior Alignment Explainability (BAE)

BAE assesses whether model explanations accurately capture genuine annotator behavioral relationships by evaluating the structural alignment between explanation-derived similarities and ground-truth inter-annotator consistency patterns.

Fundamental Hypothesis. If a model’s explanations truly reflect annotator behavior patterns, the similarity relationships derived from individual annotators’ representations (feature-level) or high-attention regions (region-level) on image batches or video frames should structurally align with those computed from actual annotation behaviors. We evaluate this alignment through Multidimensional Scaling (MDS) projection that enables visual comparison of similarity structures in interpretable 2D feature spaces.

Ground-truth Behavioral Similarity. We define the ground-truth behavioral similarity matrix $\mathbf{S}^{\text{true}} \in \mathbb{R}^{M \times M}$ where element S_{ij}^{true} represents the behavioral similarity between annotators i and j :

$$S_{ij}^{\text{true}} = \kappa(Y_i |_{S_{ij}}, Y_j |_{S_{ij}}) \quad (4)$$

computed over overlapping annotation sets, providing a standardized reference structure for explainability evaluation.

Explanation-based Similarity Computation. As shown in Figure 2(b), we compute annotator similarity matrices from model explanations at two complementary levels:

Feature-level Assessment: We extract annotator-specific learned representations from the penultimate layer for similarity computation:

$$S_{ij}^{\text{feature}} = \text{cosine}(\mathbf{F}_i^{\text{avg}}, \mathbf{F}_j^{\text{avg}}) \quad (5)$$

where $\mathbf{F}_i^{\text{avg}} = \frac{1}{|S_i|} \sum_{x \in S_i} f_{\text{feature}}(x, i)$ represents the average feature representation for annotator i . This assessment applies to all model architectures and evaluates how well learned representations capture behavioral distinctions.

Region-level Assessment: For attention-based methods, we conduct complementary region-level analysis to provide additional behavioral insights. Building upon comprehensiveness score validation [9]—which confirms that attention patterns focus on decision-relevant regions through performance comparison after masking high-attention versus random regions—we compute inter-annotator

similarities based on attention patterns over spatial/temporal regions:

$$S_{ij}^{\text{region}} = \text{cosine}(\mathbf{A}_i^{\text{avg}}, \mathbf{A}_j^{\text{avg}}) \quad (6)$$

where $\mathbf{A}_i^{\text{avg}} = \frac{1}{|S_i|} \sum_{x \in S_i} \text{Attention}(x, i)$ represents the average attention pattern for annotator i across input regions (image patches for STREET, video frames for AMER). This approach provides a complementary perspective by analyzing fine-grained spatial/temporal behavioral patterns, offering different analytical insights rather than consistently superior performance compared to feature-level assessment.

Alignment Measurement. BAE quantifies the structural alignment between explanation-derived and ground-truth behavioral similarities:

$$\text{BAE} = 1 - \frac{\|\mathbf{S}^{\text{model}} - \mathbf{S}^{\text{true}}\|_F}{\|\mathbf{S}^{\text{true}}\|_F} \quad (7)$$

where $\mathbf{S}^{\text{model}}$ represents either feature-level or region-level similarities. Higher BAE values indicate better alignment between model explanations and true behavioral relationships. MDS projection enables interpretable visual assessment in 2D feature space, where spatial proximity reflects behavioral similarity and clustering patterns reveal whether explanations capture genuine annotator behavioral relationships.

4 EXPERIMENT

We conduct comprehensive experiments to validate our proposed evaluation framework and demonstrate its utility in assessing multi-annotator learning methods. Our evaluation focuses on two objectives: (1) validating that DIC and BAE metrics accurately reflect tendency capture and explainability quality, and (2) benchmarking representative methods to uncover insights into multi-annotator modeling capabilities.

Implementation Details. All image and video inputs are resized to 224×224 and normalized before processing. Each baseline follows its original training protocol. Experiments are conducted on four NVIDIA V100 GPUs with consistent hyperparameter settings across methods.

Datasets. We evaluate on two multi-annotator datasets: AMER (video-based emotion recognition with 13 annotators) and STREET (urban image impression assessment with 10 annotators across five dimensions: Happiness, Healthiness, Safety, Liveliness, and Orderliness). The AMER dataset’s complexity in capturing temporal emotion dynamics aligns with recent advances in time-sensitive emotion recognition [37, 41]. These datasets provide dense, diverse annotations crucial for modeling annotator-specific behavior patterns.

Baseline Methods. We benchmark four representative approaches covering diverse paradigms: D-LEMA (ensemble-based aggregation), PADL (meta-learning with Gaussian fitting), MaDL (confusion-matrix-based modeling), and QuMAB (query-based annotator modeling with attention mechanisms).

Evaluation Protocol. Each method is evaluated using our proposed DIC and BAE metrics alongside traditional metrics to demonstrate the complementary insights provided by our framework.

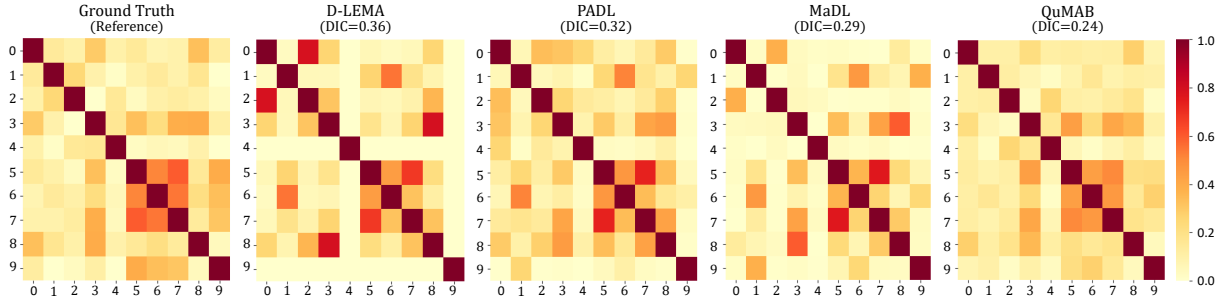


Figure 3: Visualization analysis about difference of inter-annotator consistency (DIC) via similarity matrices calculated by Cohen’s kappa coefficient on the STREET dataset (safety perspective) (10 annotators), darker colors indicate stronger agreement. Four representative models are compared with the ground truth. Lower DIC scores reflect better preservation of the underlying consistency structure. The vertical color bar denotes the similarity scale ranging from 0 (no agreement) to 1 (perfect agreement).

Dataset	D-LEMA	PADL	MaDL	QuMAB
STREET-(Happiness)	0.62 ± 0.03	0.48 ± 0.02	0.45 ± 0.01	0.43 ± 0.02
STREET-(Healthiness)	0.59 ± 0.02	0.52 ± 0.03	0.45 ± 0.03	0.38 ± 0.02
STREET-(Safety)	0.36 ± 0.02	0.32 ± 0.02	0.29 ± 0.01	0.24 ± 0.03
STREET-(Liveliness)	0.51 ± 0.01	0.43 ± 0.02	0.39 ± 0.02	0.27 ± 0.02
STREET-(Orderliness)	0.61 ± 0.02	0.57 ± 0.01	0.59 ± 0.02	0.54 ± 0.01
AMER	0.42 ± 0.03	0.36 ± 0.01	0.31 ± 0.01	0.23 ± 0.02

Table 1: Difference of Inter-annotator Consistency (DIC) score comparison across different model architectures on the AMER dataset and five STREET perspectives. Lower values indicate better preservation of inter-annotator structural tendencies.

Method	STREET (Safety perspective)				AMER			
	ACC \uparrow	FK \uparrow	PCC \uparrow	DIC \downarrow	ACC \uparrow	FK \uparrow	PCC \uparrow	DIC \downarrow
D-LEMA	0.47 ± 0.04	0.51 ± 0.04	0.46 ± 0.05	0.36 ± 0.02	0.78 ± 0.03	0.54 ± 0.04	0.49 ± 0.05	0.42 ± 0.03
PADL	0.52 ± 0.03	0.54 ± 0.03	0.49 ± 0.04	0.32 ± 0.02	0.79 ± 0.02	0.56 ± 0.03	0.52 ± 0.04	0.36 ± 0.01
MaDL	0.46 ± 0.02	0.57 ± 0.02	0.52 ± 0.03	0.29 ± 0.01	0.80 ± 0.02	0.58 ± 0.03	0.55 ± 0.03	0.31 ± 0.01
QuMAB	0.58 ± 0.02	0.61 ± 0.02	0.56 ± 0.02	0.24 ± 0.03	0.84 ± 0.02	0.63 ± 0.03	0.58 ± 0.03	0.23 ± 0.02

Table 2: Comparison of traditional evaluation metrics and DIC on the STREET dataset (safety perspective) and the AMER dataset. While ACC, Fleiss’ Kappa (FK), and Pearson Correlation Coefficient (PCC) show limited variation across methods, DIC reveals significant differences in how well models preserve annotator-specific structural tendencies.

4.1 DIC Assessment

We evaluate DIC effectiveness through visual consistency analysis, quantitative comparison across architectures, and validation against standard metrics.

Consistency Pattern Analysis. Figure 3 presents inter-annotator consistency matrices on the STREET (safety perspective) dataset. Lower DIC scores indicate better preservation of ground-truth relationships. QuMAB achieves the most faithful reconstruction of annotator consistency patterns, while D-LEMA shows notable structural distortions. This visualization demonstrates each method’s capacity to model annotator-specific relational patterns and preserve individual tendencies.

Quantitative Results. Table 1 presents comprehensive DIC scores across four architectures on both datasets. QuMAB consistently achieves the lowest DIC scores, indicating a superior tendency capture across visual and affective domains. Compared to D-LEMA, QuMAB improves DIC by 0.12–0.18 across perspectives, demonstrating robust capture of individualized labeling patterns under varying contexts.

Traditional Metrics Comparison. To further highlight the value of DIC, we compare it with three conventional metrics: (1) Accuracy (ACC) for individual annotator prediction quality; (2) Fleiss’ Kappa (FK) [13] measures the overall inter-annotator agreement adjusted for chance; (3) Pearson Correlation Coefficient (PCC) [2] quantifies the similarity between the predicted and true annotation vectors per annotator, reflecting linear trends but not structural agreement across annotators. Table 2 shows results on AMER

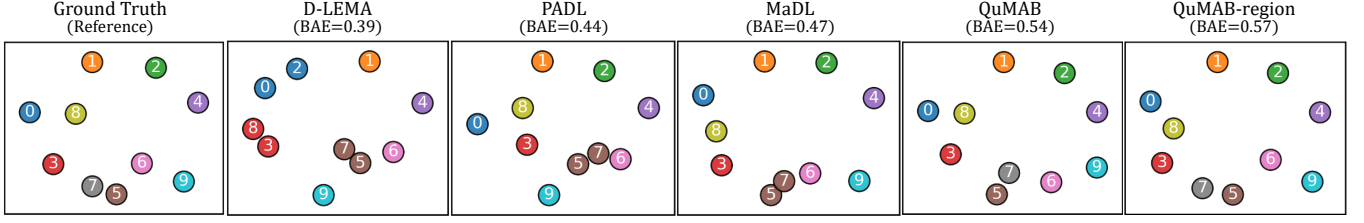


Figure 4: Visualization analysis about behavior alignment explainability (BAE) via 2D projection of annotator representations using Multidimensional Scaling (MDS) on the STREET dataset (safety perspective). Results show progressively improved alignment with higher BAE scores. Region-level QuMAB has a relatively better alignment. Each point denotes an annotator, and proximity indicates higher behavioral similarity. Same colors denote clusters of annotators with strong agreement ($\kappa > 0.6$).

Dataset	D-LEMA	PADL	MaDL	QuMAB	QuMAB-region
STREET-Ha	0.28 ± 0.03	0.33 ± 0.04	0.35 ± 0.02	0.38 ± 0.03	0.41 ± 0.02
STREET-He	0.31 ± 0.04	0.36 ± 0.02	0.38 ± 0.03	0.42 ± 0.02	0.40 ± 0.03
STREET-Sa	0.39 ± 0.02	0.44 ± 0.03	0.47 ± 0.02	0.54 ± 0.02	0.57 ± 0.02
STREET-Li	0.35 ± 0.02	0.38 ± 0.03	0.41 ± 0.04	0.46 ± 0.02	0.45 ± 0.03
STREET-Or	0.24 ± 0.03	0.26 ± 0.02	0.25 ± 0.04	0.29 ± 0.03	0.31 ± 0.02
AMER	0.41 ± 0.03	0.45 ± 0.02	0.48 ± 0.01	0.52 ± 0.02	0.55 ± 0.02

Table 3: Behavior Alignment Explainability (BAE) scores across different modeling architectures. Each value measures the alignment between learned annotator representations and ground-truth behavioral similarity. Region-level BAE (last column) evaluates inter-annotator similarity based on attention-over-region patterns mapped into feature space. Higher is better.

Method	STREET (Safety perspective)				AMER			
	Cos \uparrow	Grad \uparrow	Comp \uparrow	BAE \uparrow	Cos \uparrow	Grad \uparrow	Comp \uparrow	BAE \uparrow
D-LEMA	0.68 ± 0.05	0.54 ± 0.07	N/A	0.39 ± 0.02	0.72 ± 0.04	0.58 ± 0.06	N/A	0.41 ± 0.03
PADL	0.71 ± 0.04	0.58 ± 0.06	N/A	0.44 ± 0.03	0.74 ± 0.03	0.61 ± 0.05	N/A	0.45 ± 0.02
MaDL	0.69 ± 0.06	0.56 ± 0.08	N/A	0.47 ± 0.02	0.73 ± 0.05	0.59 ± 0.07	N/A	0.48 ± 0.01
QuMAB	0.74 ± 0.03	0.63 ± 0.05	0.16 ± 0.02	0.54 ± 0.02	0.76 ± 0.02	0.66 ± 0.04	0.17 ± 0.03	0.52 ± 0.02

Table 4: Comparison with alternative explainability metrics. Cos = feature cosine similarity; Grad = gradient correlation; Comp = comprehensiveness score. BAE more effectively differentiates model behavior than conventional metrics.

Dataset	DIC \downarrow				BAE \uparrow			
	Random	Consensus	D-LEMA	QuMAB	Random	Uniform	D-LEMA	QuMAB
STREET-Ha	0.89 ± 0.03	0.54 ± 0.02	0.62 ± 0.03	0.43 ± 0.02	0.05 ± 0.04	0.18 ± 0.02	0.28 ± 0.03	0.38 ± 0.03
STREET-He	0.91 ± 0.04	0.58 ± 0.01	0.59 ± 0.02	0.38 ± 0.02	0.02 ± 0.03	0.16 ± 0.03	0.31 ± 0.04	0.42 ± 0.02
STREET-Sa	0.87 ± 0.02	0.51 ± 0.03	0.36 ± 0.02	0.24 ± 0.03	0.08 ± 0.03	0.23 ± 0.02	0.39 ± 0.02	0.54 ± 0.02
STREET-Li	0.90 ± 0.03	0.56 ± 0.01	0.51 ± 0.01	0.27 ± 0.02	0.04 ± 0.05	0.19 ± 0.01	0.35 ± 0.02	0.46 ± 0.02
STREET-Or	0.93 ± 0.05	0.61 ± 0.03	0.61 ± 0.02	0.54 ± 0.01	0.02 ± 0.06	0.15 ± 0.04	0.24 ± 0.03	0.29 ± 0.03
AMER	0.86 ± 0.02	0.53 ± 0.02	0.42 ± 0.03	0.23 ± 0.02	0.06 ± 0.03	0.21 ± 0.02	0.41 ± 0.03	0.52 ± 0.02

Table 5: Ablation Study for DIC and BAE metrics. Left: DIC validation showing progression from Random \rightarrow Consensus \rightarrow D-LEMA (limit) \rightarrow QuMAB (best). Right: BAE validation showing progression from Random \rightarrow Uniform \rightarrow D-LEMA (limit) \rightarrow QuMAB (best). Both metrics demonstrate clear discriminative power across the performance spectrum.

and STREET (safety perspective). While traditional metrics show limited variation (ranges within 0.06–0.09), DIC reveals more pronounced differences. D-LEMA exhibits moderate traditional scores

but the highest DIC, indicating structural modeling failure. QuMAB achieves both the best traditional scores and the lowest DIC, confirming faithful modeling of individual tendencies. These findings

confirm that traditional metrics may obscure structural modeling failures, whereas DIC provides clearer, complementary insights by explicitly evaluating tendency capture.

4.2 BAE Assessment

We evaluate BAE effectiveness through visual analysis of behavioral alignment structures, quantitative comparison across feature-level with a complementary region-level perspective, and validation against alternative explainability metrics.

Behavioral Alignment Analysis. Figure 4 shows MDS projections of learned annotator representations for STREET (safety perspective). Ground truth establishes clear behavioral clusters, with proximity indicating similarity and colors representing high-agreement groups ($\kappa > 0.6$). The progression shows: D-LEMA (0.39) with significant distortion, PADL (0.44) with moderate preservation, MaDL (0.47) with better structure, and QuMAB (0.54) with closest ground-truth alignment. QuMAB’s region-level assessment (0.57) complements feature-level analysis by projecting inter-annotator similarities—derived from high-attention input regions—into feature space, offering finer-grained behavioral insight and enhancing overall evaluation robustness. These comparisons highlight each method’s capacity to preserve annotator-specific behavioral relationships through learned representations.

Quantitative Results. Table 3 presents comprehensive BAE scores across the four multi-annotator learning architectures for feature-level assessments and QuMAB’s region-level complementary assessment. For feature-level analysis, we extract annotator-specific embeddings from the penultimate layer and compute behavioral similarities based on annotation patterns. QuMAB consistently achieves the highest feature-level BAE scores, indicating superior capture of genuine behavioral relationships. Performance improvements vary by domain: substantial gains on Safety (0.15 over D-LEMA) versus modest improvements on Orderliness (0.05), reflecting varying complexity of behavioral pattern preservation. For QuMAB’s attention-based explanations, we conduct an exploratory region-level analysis. Its scores provide complementary insights to feature-level analysis with modest and variable improvements (0.02-0.03 across most datasets). This modest enhancement does not indicate superior performance but rather reflects a different analytical perspective that captures fine-grained spatial/temporal behavioral patterns. The complementary nature of these two assessment levels enables more robust, comprehensive evaluation of behavioral alignment.

Alternative Explainability Metrics Comparison. We also compare BAE with three existing metrics for evaluating explanation quality: (1) Feature Cosine Similarity for basic representation alignment; (2) Gradient Correlation for feature importance alignment; (3) Comprehensiveness Score validates attention faithfulness by measuring performance drops after masking high-attention regions versus random regions, applicable only to attention-based methods like QuMAB. Table 4 shows results on AMER and STREET (safety perspective). Additional results for the other STREET perspectives are provided in the Supplementary Material. While Feature Cosine Similarity and Gradient Correlation show consistently high scores

with minimal variation (std = 0.02-0.04), BAE demonstrates substantially higher discrimination (std = 0.05-0.06), revealing clearer distinctions between methods’ explainability capabilities. Traditional metrics capture surface-level similarities rather than meaningful behavioral differences. BAE’s enhanced sensitivity identifies methods with similar feature representations that fail to preserve essential behavioral relationship structures, providing complementary insights by explicitly evaluating behavioral alignment preservation across all architectural approaches.

4.3 Ablation Study

We conduct controlled ablation experiments using carefully designed baseline scenarios representing extreme cases of tendency capture and explainability quality to validate DIC and BAE effectiveness and sensitivity.

DIC Ablation Analysis. We evaluate DIC sensitivity across the performance spectrum using controlled baselines: *Random* simulates completely unstructured behavior with uniformly random label assignments; *Consensus* assigns majority-vote labels to all annotators, eliminating individual variations; D-LEMA provides a representative baseline with limited modeling capacity; QuMAB achieves the strongest results. As shown in the left columns of Table 5 (left columns) shows DIC effectively captures the performance hierarchy: Random yields highest scores (0.86-0.93), Consensus produces moderate scores (0.51-0.61), D-LEMA shows improved performance (0.36-0.62), while QuMAB achieves lowest scores (0.23-0.54). This clear ordering validates DIC’s ability to distinguish meaningful differences in tendency capture across the entire performance spectrum.

BAE Ablation Analysis. We validate BAE discriminative power using controlled representation scenarios: *Random* generates independent random feature vectors (near-zero similarities); *Uniform* assigns identical representations to all annotators (perfect similarity without behavioral differences); D-LEMA and QuMAB are selected as representative models exhibiting the lowest and highest performance, respectively. Table 5 (right columns) demonstrates BAE effectiveness: Random produces the lowest scores (0.02-0.08), Uniform achieves moderate scores (0.15-0.23), D-LEMA shows improved alignment (0.24-0.41), and QuMAB reaches the highest scores (0.29-0.54). This progression confirms BAE’s sensitivity to explanation quality differences and validates its utility for evaluating behavioral alignment across diverse modeling approaches.

5 CONCLUSION

We proposed the first unified evaluation framework for Individual Tendency Learning (ITL) with two novel metrics: (1) Difference of Inter-annotator Consistency (DIC) quantifies tendency capture by comparing predicted and ground-truth inter-annotator consistency structures; (2) Behavior Alignment Explainability (BAE) evaluates explainability quality through aligning explainability-derived with ground-truth labeling similarity structures via Multidimensional Scaling (MDS). Extensive experiments on different datasets across different model architectures validate the effectiveness of our proposed evaluation framework. A potential enhancement would be incorporating human-derived behavioral signals into the evaluation framework. While eye-tracking can capture annotators’ attention

patterns over image or video content, such approaches remain high-cost and difficult to scale. In future work, we aim to explore scalable alternatives for acquiring such signals, which would advance ITL and benefit the broader research community.

REFERENCES

- [1] Shadi Albarqouni, Christoph Baur, Felix Achilles, Vasileios Belagiannis, Stefanie Demirci, and Nassir Navab. 2016. Aggnet: deep learning from crowds for mitosis detection in breast cancer histology images. *IEEE transactions on medical imaging* 35, 5 (2016), 1313–1321.
- [2] Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. 2009. Pearson correlation coefficient. In *Noise reduction in speech processing*. Springer, 1–4.
- [3] Peng Cao, Yilun Xu, Yuqing Kong, and Yizhou Wang. 2019. Max-mig: an information theoretic approach for joint learning from crowds. *arXiv preprint arXiv:1905.13436* (2019).
- [4] Bob Carpenter. 2011. A Hierarchical Bayesian Model of Crowdsourced Relevance Coding. In *TREC*.
- [5] Junfan Chen, Richong Zhang, Jie Xu, Chunming Hu, and Yongyi Mao. 2023. A Neural Expectation-Maximization Framework for Noisy Multi-Label Text Classification. *IEEE Transactions on Knowledge and Data Engineering* 35, 11 (2023), 10992–11003. <https://doi.org/10.1109/TKDE.2023.3223067>
- [6] Yuan-Chia Cheng, Zu-Yun Shiau, Fu-En Yang, and Yu-Chiang Frank Wang. 2023. TAX: Tendency-and-Assignment Explainer for Semantic Segmentation with Multi-Annotators. *arXiv preprint arXiv:2302.09561* (2023).
- [7] Molnar Christoph. 2020. Interpretable machine learning: A guide for making black box models explainable. (2020).
- [8] Alexander Philip Dawid and Allan M Skene. 1979. Maximum likelihood estimation of observer error-rates using the EM algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 28, 1 (1979), 20–28.
- [9] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C Wallace. 2019. ERASER: A benchmark to evaluate rationalized NLP models. *arXiv preprint arXiv:1911.03429* (2019).
- [10] Zixin Ding, Si Chen, Ruoxi Jia, and Yuxin Chen. 2023. Learning to rank for active learning via multi-task bilevel optimization. *arXiv preprint arXiv:2310.17044* (2023).
- [11] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* (2017).
- [12] Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin* 76, 5 (1971), 378.
- [13] Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin* 76, 5 (1971), 378.
- [14] Mitchell L Gordon, Michelle S Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S Bernstein. 2022. Jury learning: Integrating dissenting voices into machine learning models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–19.
- [15] Kilem L Gwet. 2011. On the Krippendorff’s alpha coefficient. *Manuscript submitted for publication*. Retrieved October 2, 2011 (2011), 2011.
- [16] Marek Herde, Denis Huseljic, and Bernhard Sick. 2023. Multi-annotator Deep Learning: A Probabilistic Framework for Classification. *arXiv preprint arXiv:2304.02539* (2023).
- [17] Uthman Jinadu, Jesse Annan, Shanshan Wen, and Yi Ding. 2023. Loss Modeling for Multi-Annotator Datasets. *arXiv preprint arXiv:2311.00619* (2023).
- [18] Zheng Lian, Haiyang Sun, Licai Sun, Kang Chen, Mngyu Xu, Kexin Wang, Ke Xu, Yu He, Ying Li, Jinming Zhao, et al. 2023. Mer 2023: Multi-label learning, modality robustness, and semi-supervised learning. In *Proceedings of the 31st ACM International Conference on Multimedia*. 9610–9614.
- [19] Zehui Liao, Shishuai Hu, Yutong Xie, and Yong Xia. 2024. Modeling annotator preference and stochastic annotation error for medical image segmentation. *Medical Image Analysis* 92 (2024), 103028.
- [20] Zahra Mirikharaji, Kumar Abhishek, Saeed Izadi, and Ghassan Hamarneh. 2021. D-lema: Deep learning ensembles from multiple annotations-application to skin lesion segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1837–1846.
- [21] Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa Nguyen, Michelle Peters, Yasmin Schmitt, Jörg Schlötterer, Maurice Van Keulen, and Christin Seifert. 2023. From anecdotal evidence to quantitative evaluation methods: A systematic review on evaluating explainable ai. *Comput. Surveys* 55, 13s (2023), 1–42.
- [22] Rebecca J Passonneau and Bob Carpenter. 2014. The benefits of a model of annotation. *Transactions of the Association for Computational Linguistics* 2 (2014), 311–326.
- [23] Vikas C. Raykar, Shipeng Yu, Linda H. Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bogoni, and Linda Moy. 2010. Learning From Crowds. *Journal of Machine Learning Research* 11, 43 (2010), 1297–1322. <http://jmlr.org/papers/v11/raykar10a.html>
- [24] Filipe Rodrigues, Francisco Pereira, and Bernardete Ribeiro. 2014. Gaussian Process Classification and Active Learning with Multiple Annotators. In *Proceedings of the 31st International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 32)*, Eric P. Xing and Tony Jebara (Eds.). PMLR, Beijing, China, 433–441. <https://proceedings.mlr.press/v32/rodrigues14.html>
- [25] Xuanmeng Sha, Liyun Zhang, Tomohiro Mashita, and Yuki Uranishi. 2024. 3DFacePolicy: Speech-Driven 3D Facial Animation with Diffusion Policy. *arXiv preprint arXiv:2409.10848* (2024).
- [26] Rion Snow, Brendan O’Connor, Daniel Jurafsky, and Andrew Ng. 2008. Cheap and Fast – But is it Good? Evaluating Non-Expert Annotations for Natural Language Tasks. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, Mirella Lapata and Hwee Tou Ng (Eds.). Association for Computational Linguistics, Honolulu, Hawaii, 254–263. <https://aclanthology.org/D08-1027>
- [27] Rion Snow, Brendan O’connor, Dan Jurafsky, and Andrew Y Ng. 2008. Cheap and fast – but is it good? evaluating non-expert annotations for natural language tasks. In *Proceedings of the 2008 conference on empirical methods in natural language processing*. 254–263.
- [28] Ryutaro Tanno, Ardavan Saeedi, Swami Sankaranarayanan, Daniel C. Alexander, and Nathan Silberman. 2019. Learning From Noisy Labels by Regularized Estimation of Annotator Confusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [29] Ryutaro Tanno, Ardavan Saeedi, Swami Sankaranarayanan, Daniel C Alexander, and Nathan Silberman. 2019. Learning from noisy labels by regularized estimation of annotator confusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 11244–11253.
- [30] AnthonyJ. Viera and JoanneM. Garrett. 2005. Understanding interobserver agreement: the kappa statistic. *Family Medicine, Family Medicine* (May 2005).
- [31] Chongyang Wang, Yuan Gao, Chenyou Fan, Junjie Hu, Tin Lum Lam, Nicholas D Lane, and Nadia Bianchi-Berthouze. 2023. Learn2agree: Fitting with multiple annotators without objective ground truth. In *International Workshop on Trustworthy Machine Learning for Healthcare*. Springer, 147–162.
- [32] Peter Welinder, Steve Branson, Pietro Perona, and Serge Belongie. 2010. The multidimensional wisdom of crowds. *Advances in neural information processing systems* 23 (2010).
- [33] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier Movellan, and Paul Ruvolo. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. *Advances in neural information processing systems* 22 (2009).
- [34] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier Movellan, and Paul Ruvolo. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. *Advances in neural information processing systems* 22 (2009).
- [35] Yan Yan, Römer Rosales, Glenn Fung, Ramanathan Subramanian, and Jennifer Dy. 2014. Learning from multiple annotators with varying expertise. *Machine learning* 95 (2014), 291–327.
- [36] Liyun Zhang. 2024. *Integrating Panoptic-Level to Image Translation*. Ph.D. Dissertation. PhD Dissertation.
- [37] Liyun Zhang. 2024. MicroEmo: Time-Sensitive Multimodal Emotion Recognition with Micro-Expression Dynamics in Video Dialogues. *arXiv preprint arXiv:2407.16552* (2024).
- [38] Liyun Zhang, Zheng Lian, Hong Liu, Takanori Takebe, and Yuta Nakashima. 2025. QuMAB: Query-based Multi-annotator Behavior Pattern Learning. *arXiv preprint arXiv:2507.17653* (2025).
- [39] Liyun Zhang, Zheng Lian, Hong Liu, Takanori Takebe, and Yuta Nakashima. 2025. SimLabel: Similarity-Weighted Semi-supervision for Multi-annotator Learning with Missing Labels. *arXiv preprint arXiv:2504.09525* (2025).
- [40] Liyun Zhang, Nanyan Liu, Yuanbin Hou, and Xiaojian Liu. 2014. Uneven illumination image segmentation based on multi-threshold S-F. *Opto-Electronic Engineering* 41, 7 (2014), 81–87.
- [41] Liyun Zhang, Zhaojie Luo, Shuqiong Wu, and Yuta Nakashima. 2024. MicroEmo: Time-Sensitive Multimodal Emotion Recognition with Subtle Clue Dynamics in Video Dialogues. In *Proceedings of the 2nd International Workshop on Multimodal and Responsible Affective Computing*. 110–115.
- [42] Liyun Zhang, Photchara Ratsamee, Zhaojie Luo, Yuki Uranishi, Manabu Higashida, and Haruo Takemura. 2023. Panoptic-level image-to-image translation for object recognition and visual odometry enhancement. *IEEE Transactions on Circuits and Systems for Video Technology* 34, 2 (2023), 938–954.
- [43] Liyun Zhang, Photchara Ratsamee, Yuki Uranishi, Manabu Higashida, and Haruo Takemura. 2022. Thermal-to-Color Image Translation for Enhancing Visual Odometry of Thermal Vision. In *2022 IEEE International Symposium on Safety, Security, and Rescue Robotics (SSRR)*. IEEE, 33–40.
- [44] Liyun Zhang, Photchara Ratsamee, Bowen Wang, Zhaojie Luo, Yuki Uranishi, Manabu Higashida, and Haruo Takemura. 2023. Panoptic-aware image-to-image translation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*. 259–268.
- [45] Le Zhang, Ryutaro Tanno, Moucheng Xu, Yawen Huang, Kevin Bronik, Chen Jin, Joseph Jacob, Yefeng Zheng, Ling Shao, Olga Ciccarelli, et al. 2023. Learning

from multiple annotators for medical image segmentation. *Pattern Recognition*

138 (2023), 109400.