

Beyond Consensus: Perspectivist Modeling and Evaluation of Annotator Disagreement in NLP

Anonymous ACL submission

Abstract

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
Annotator disagreement is widespread in NLP, particularly for subjective and ambiguous tasks such as toxicity detection and stance analysis. While early approaches treated disagreement as noise to be removed, recent work increasingly models it as a meaningful signal reflecting variation in interpretation and perspective. This survey provides a unified view of disagreement-aware NLP methods. We first present a domain-agnostic taxonomy of the sources of disagreement spanning data, task, and annotator factors. We then synthesize modeling approaches using a common framework defined by prediction targets and pooling structure, highlighting a shift from consensus learning toward explicitly modeling disagreement, and toward capturing structured relationships among annotators. We review evaluation metrics for both predictive performance and annotator behavior, and noting that most fairness evaluations remain descriptive rather than normative. We conclude by identifying open challenges and future directions, including integrating multiple sources of variation, developing disagreement-aware interpretability frameworks, and grappling with the practical tradeoffs of perspectivist modeling.

028 1 Introduction

029
030
031
032
033
034
035
036
037
038
039
040
041
042
NLP applications largely rely on supervised learning, which depends on annotated data. Annotation is often operationalized through majority voting, an assumption that can be problematic—particularly for complex tasks where true experts may be outnumbered, or for inherently ambiguous tasks that admit multiple valid interpretations. Majority aggregation can also obscure minoritized perspectives, leading to biased models and representational harm (Blodgett et al., 2020; Gordon et al., 2021). This “single ground truth” assumption has been increasingly challenged, especially with the rise of subjective NLP tasks such as toxic language detection and quality estimation. Early critiques

043
044
045
046
047
048
049
050
051
052
appear in Aroyo and Welty (2015), which questions the existence of a unique truth in crowd-sourced annotation. This shift aligns with broader efforts to make NLP systems more inclusive, as consensus-based labels often disadvantage minority viewpoints (Blodgett et al., 2020; Gordon et al., 2021). Building on this, Cabitza et al. (2023) formalized *perspectivism* in NLP, advocating for models that integrate diverse human viewpoints rather than collapsing them into gold standards.

053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
We formalize a taxonomy of disagreement by synthesizing prior work into three sources: data, task, and annotator factors. We trace the evolution of NLP approaches for learning from disagreement, from latent-truth models to multi-annotator, embedding-based methods. To unify these developments, we introduce a synthesis table highlighting key trends, including the shift from treating disagreement as noise to modeling it as a prediction target, and the growing emphasis on structured relationships among annotators. We also survey evaluation practices and fairness considerations. Finally, we outline open challenges and future directions, including jointly modeling multiple sources of variation and navigating the practical tradeoffs of perspectivist approaches. Our contributions are:

- 069
070
1. **A unified taxonomy of disagreement** across data-, task-, and annotator-driven sources.
- 071
072
073
074
2. **A synthesis of disagreement-aware methods**, mapping approaches to disagreement sources, prediction targets, and pooling structures.

075
076
077
078
079
080
081
Two recent surveys are closely related to our work but differ in scope and emphasis. Frenda et al. (2024) provides a broad, conceptual overview of perspectivist NLP, focusing on definitions, dataset practices, and sociotechnical motivations. We build on their framing while shifting toward a more model-centric and technical analysis. Uma et al.

(2021) survey disagreement-aware learning across computer vision and NLP, emphasizing empirical training and evaluation strategies. Our survey complements this work by centering NLP-specific sources of subjectivity and by jointly examining models and evaluation within a unified framework. We survey over 120 prior works spanning disagreement sources, modeling methods, and evaluation paradigms. We discuss our survey scope and selection criteria in A.1.

2 Sources of Disagreement

Prior work has analyzed sources of annotator disagreement across tasks and domains. Basile et al. (2021) identify three broad contributors: individual differences (partly linked to sociodemographics), stimulus characteristics (e.g., linguistic or task ambiguity), and contextual inconsistency in human behavior. Sandri et al. (2023) propose a finer taxonomy (sloppy annotations, ambiguity, missing information, and subjectivity) and show that in offensive language detection, subjectivity driven by personal bias and task design dominates disagreement. In legal NLP, Xu et al. (2024b) introduce a taxonomy highlighting genuine ambiguity, narrative uncertainty, and annotation context. Overall, existing taxonomies are largely domain-specific, emphasizing different facets of disagreement depending on task formulation and application context.

2.1 Taxonomy of disagreement

We propose a unifying, domain-agnostic taxonomy (Fig. 1) that synthesizes these views into three overarching sources of disagreement: data factors (e.g., ambiguity and data quality), task factors (e.g., formulation and interface design), and annotator factors (e.g., demographics, preferences, and errors). Importantly, different sources of disagreement can interact with one another. Unclear task formulation can lead to inconsistent annotator behavior. Linguistic ambiguity can give rise to disagreement that varies with annotators’ individual and group identities, leading to systematically different interpretations of the same text.

2.2 Data Factors

Prior work identifies three data-driven sources of disagreement: data quality issues, linguistic ambiguity, and epistemic uncertainty—cases with no single ground truth, such as subjective or culturally situated interpretations (Poesio, 2020; Uma et al., 2021; Pavlick and Kwiatkowski, 2019). Missing

context or noisy inputs exacerbate inconsistency, especially for short text (Plank et al., 2014). More fundamentally, language is inherently ambiguous, and in many subjective tasks disagreement reflects irreducible multiplicity rather than annotation error (Pavlick and Kwiatkowski, 2019; Uma et al., 2021).

2.3 Annotator Factors

Individual identities. Prior work documents stable, annotator-specific biases that introduce variance into training data and directly affect model performance (Otterbacher, 2018; Larimore et al., 2021; Pavlick and Kwiatkowski, 2019; Geva et al., 2019). Beyond general disagreement, psychological and moral traits systematically correlate with annotation behavior. Personality dimensions such as agreeableness and conscientiousness are associated with evaluative leniency (Mieleszczenko-Kowszewicz et al., 2023), while moral values exert stronger influence on judgments of offensiveness than geographic or cultural background (Davani et al., 2024). Together, these findings show that annotators bring stable individual preferences that shape how they interpret and label language.

Group identities. Group-level attributes such as gender, race, age, and political orientation influence annotation behavior, particularly in socially grounded tasks. Gender-linked variation has been observed in hate speech detection (Wojatzki et al., 2018) and even in syntactic tasks (Garimella et al., 2019), aligning with sociolinguistic evidence of gendered language use (Mondorf, 2002). In subjective domains such as toxicity or sentiment, demographic shifts can substantially affect model performance (Ding et al., 2022). Qualitative work further shows that lived experience shapes interpretation: community insiders label gang-related language differently from academic annotators (Patton et al., 2019), and political or racial attitudes correlate with toxicity judgments, including treatment of African American English (Sap et al., 2022; Sang and Stanton, 2021; Luo et al., 2021). In these settings, disagreement often reflects social perspective rather than noise (Chulvi et al., 2023). However, identity effects are task-dependent. For less socially salient tasks (e.g., word similarity, sentiment, NLI), demographic attributes do not consistently explain annotation differences (Biester et al., 2022), and explicitly modeling sociodemographics may yield limited gains (Orlikowski et al., 2023). These findings caution against ecological fallacies (Robinson,

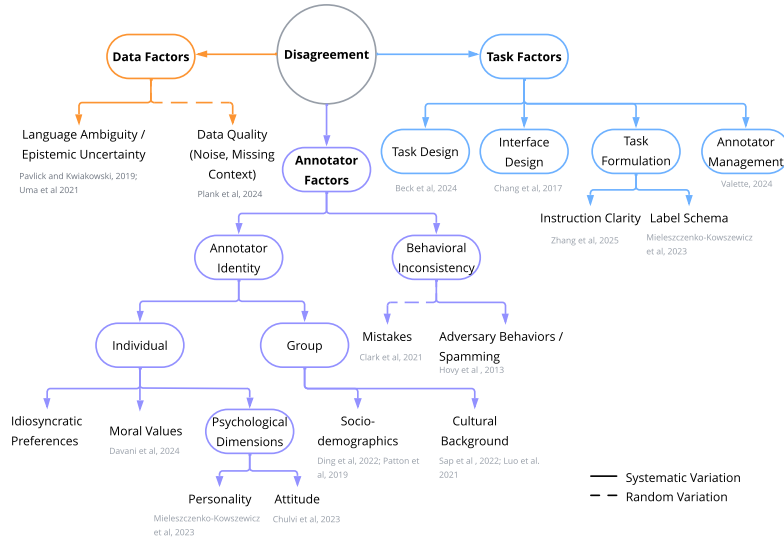


Figure 1: Taxonomy of sources of annotator disagreement, with key citations for each sub-category. For each source of disagreement, we denote systematic variation using solid lines, and random variation using dotted.

1950): group-level identity does not uniformly predict individual annotation behavior.

Annotator behavioral inconsistencies. Abercrombie et al. (2025) found that annotators give inconsistent responses around 25% of the time across four different NLP tasks. Clark et al. (2021) also show that untrained annotators identify GPT-3- versus human-written text only at chance levels. Annotators may also produce low-quality or strategically chosen labels to maximize pay (Hovy et al., 2013).

2.4 Task Factors

Annotation task design—its formulation, structure, and presentation—strongly shapes disagreement. Ambiguous or underspecified tasks permit multiple valid interpretations, leading to systematic divergence; in LLM-as-Judge settings, ambiguous prompts elicit different interpretive strategies (Zhang et al., 2025b). Interface design also affects agreement: unclear guidelines increase inconsistency, while alternative interfaces can surface ambiguity rather than suppress it (Chang et al., 2017). Annotation schemas and scales further modulate disagreement: finer-grained ratings amplify variability relative to binarized labels (Mieszczenko-Kowszewicz et al., 2023), and different schemas (e.g., scalar ratings vs. pairwise preferences) capture related but distinct judgments, particularly in RLHF (Dsouza and Kovatchev, 2025). Presentation and labor conditions likewise matter: order effects, fatigue, and satisficing shift responses over time (Strack, 1992; Krosnick et al., 1996;

Galesic and Bosnjak, 2009; Beck et al., 2024). Fine-grained schemes may inflate disagreement due to degraded consistency under poor working conditions (Valette, 2024; Yang et al., 2023).

3 Learning from Disagreement

We discuss three families of disagreement modeling, tracing NLP’s shift from treating disagreement as noise to modeling it as structured variation.

3.1 Latent Truth and Annotator Reliability Modeling

Models in this family attribute disagreement to annotator reliability, bias, or task difficulty, with the goal of inferring latent true labels while estimating annotator behavior. Foundational work by Dawid and Skene (1979) introduced an EM-based framework to jointly infer latent labels and annotator error rates from noisy annotations, underpinning many later Bayesian and discriminative models. Subsequent work incorporates task characteristics, explicitly modeling difficulty (Whitehill et al., 2009; Ma et al., 2015) or latent topics that align with annotator expertise (Fan et al., 2015; Ma et al., 2015; Zhao et al., 2015; Welinder et al., 2010). Parallel work has focused on annotator modeling, representing competence as scalar accuracies (Demartini et al., 2012; Karger et al., 2011; Liu et al., 2012), confusion matrices (Dawid and Skene, 1979; Raykar et al., 2010; Liu et al., 2012), or bias–variance decompositions (Welinder et al., 2010; Raykar et al., 2010). Confusion-matrix-

based models are more expressive and typically outperform simpler formulations (Zheng et al., 2017). Formally, latent truth models posit a single hidden label z_i for each item and represent annotator disagreement as noise around this truth. These approaches model annotations as

$$p(y_{ij} | z_i, a_j) \quad (1)$$

a_j denotes annotator-specific parameters (such as a confusion matrix), and marginalize over z_i during inference. The objective function (Dawid and Skene, 1979; Hovy et al., 2013) maximizes the marginal likelihood of the observed annotations by marginalizing over a single true label per item.

$$\max_{\{a_j\}} \sum_i \log \sum_{z_i} p(z_i) \prod_j p(y_{ij} | z_i, a_j) \quad (2)$$

Disagreement reflects annotator unreliability rather than item ambiguity.

Later work integrates these ideas into learning architectures. Neural extensions such as CrowdLayer (Rodrigues and Pereira, 2018) and label-transfer approaches (Tanno et al., 2019) replace explicit confusion matrices with differentiable components, while other models introduce worker weighting or structured noise modeling (Gao et al., 2022; Cao et al., 2023; Wei et al., 2022; Chu et al., 2021). Paun et al. (2018a) show that partial pooling—i.e., drawing annotator parameters from a shared population distribution—best balances expressivity and generalization. Despite their strengths, latent truth models can struggle with scalability and estimation under sparse labeling or weakly differentiated answers. To address this, CROWDLAB (Goh et al., 2023) combines annotator statistics with task-level features by treating classifier predictions as an additional annotator, improving robustness without iterative inference. Recent work further incorporates task features and demographic structure: Simpson et al. (2015) jointly models text, annotator bias, and latent truth, while NUTMEG (Ivey et al., 2025) extends truth inference to demographic subpopulations, estimating competence and predicting labels per group rather than collapsing to a single consensus (Hovy et al., 2013; Paun et al., 2018b).

3.2 Task-based Annotator Models

This model family treats each annotator as a distinct task, modeling labeling behavior directly rather than collapsing annotations into a single truth. This family reframes disagreement from noise to a structured signal. Given an input x_i , a shared encoder

produces a representation $h_i = f(x_i)$, and each annotator j defines a task-specific predictor

$$p^{(j)}(y | x_i) \quad (3)$$

Model parameters are learned by minimizing a supervised loss over annotator–item pairs,

$$\min_{\theta} \sum_i \sum_{j \in \mathcal{A}_i} \ell(y_{ij}, p_{\theta}^{(j)}(y | x_i)) \quad (4)$$

θ is all model parameters (shared encoders and annotator-specific heads), \mathcal{A}_i is set of annotators who labeled i , and $\ell(\cdot, \cdot)$ is typically cross-entropy for classification or squared error for regression. Disagreement is captured implicitly through divergence across annotator-specific predictions.

Early work by Cohn and Specia (2013) models annotators as correlated tasks in a Gaussian Process framework, using an inter-annotator covariance matrix to control pooling. This nonparametric analogue to hierarchical Bayesian models identifies annotator similarity and outliers, with partial pooling performing best (Paun et al., 2018b). Neural extensions implement this idea via shared backbones and annotator-specific heads, including CrowdLayer (Rodrigues and Pereira, 2018) and related architectures (Guan et al., 2018). Fornaciari et al. (2021) further reinterpret disagreement as a regularization signal to reduce overconfidence on ambiguous inputs, learning from both the gold labels and the distribution over multiple annotators (which they treat as soft label distributions in a single auxiliary task). More recent work models richer annotator structure. Mixture-of-experts approaches capture feature-level heterogeneity (Han et al., 2025), while multi-head Transformers preserve uncertainty aligned with empirical disagreement (Mostafazadeh Davani et al., 2022). Group-aware and loss-based extensions incorporate demographic structure or explicit tradeoffs between denoising and minority perspectives (Fleisig et al., 2023; Jinadu and Ding, 2024). Collectively, task-based models capture structured perspective variation beyond latent-truth formulations.

3.3 Embedding-based Annotator Models

The third family of embedding-based models departs from latent-truth paradigms by treating disagreement as systematic variation in annotator behavior rather than noise around a hidden true label. Instead of assigning each annotator a separate prediction head, these approaches encode annotator

differences in a shared latent space, enabling scalability to thousands of annotators with sparse labels, an important limitation of task-based methods. Embedding-based models parameterize annotator-level predictions as

$$p(y_{ij} | h_i, e_j) \quad (5)$$

$h_i = f(x_i)$ is an embedding of the input item and $e_j = g(a_j)$ is an embedding of the annotator j . The loss function optimizes standard supervised losses over annotator–item pairs, where θ is all model parameters (shared encoders and annotator-specific heads), \mathcal{A}_i is the set of annotators who labeled i , and $\ell(\cdot, \cdot)$ is typically instantiated as cross-entropy for classification or squared error for regression.

$$\min_{\theta} \sum_i \sum_{j \in \mathcal{A}_i} \ell(y_{ij}, p_{\theta}(y | h_i, e_j)) \quad (6)$$

Disagreement is formalized as the interaction between item and annotator embeddings.

Early work introduced annotator embeddings to capture individual biases in subjective NLP tasks (Kocoń et al., 2021). AART (Mokhberian et al., 2024) extends this idea to large-scale crowdsourcing with regularization for robustness and fairness. Subsequent models enrich the framework by incorporating annotation embeddings (Deng et al., 2023) or demographic structure. Jury Learning (Gordon et al., 2022) combines content, annotator, and demographic embeddings to predict individual votes and compose configurable juries, while DEM-MoE (Xu et al., 2025) explicitly models demographic structure via a mixture-of-experts to capture intersectional variation. Recent work further improves scalability by replacing stored embeddings with hypernetwork-generated, annotator-specific LoRA adapters (Ignatev et al., 2025). Embedding-based models show a progression in the granularity of disagreement they represent: from individual annotators (Kocoń et al., 2021), to demographic groups (Gordon et al., 2022), to population-level distributions (Weerasooriya et al., 2023). This trajectory points toward an important future direction: estimating subgroup-specific disagreement rather than averaging across populations. Building on prior work (Sap et al., 2022; Lakkaraju et al., 2015), future work could focus on inferring how social subgroups interpret the same item, enabling perspective-aware predictions that surface when disagreement reflects broader social divides.

Direct disagreement modeling. A subset of embedding-based models predicts disagreement directly rather than treating it as noise. These approaches vary in how disagreement is represented—implicitly via soft supervision, as a scalar score, or explicitly as a label distribution—but share the goal of modeling disagreement as the prediction target. Embedding-based architectures are well-suited to this setting because they enable generalization across annotators, conditioning on demographic attributes, and marginalization over annotator populations. Implicit approaches predict aggregate disagreement without explicitly modeling a distribution. For example, Xu et al. (2024a) predict empirical label proportions from text using cross-entropy against soft labels, while Wan et al. (2023) conditions predictions on annotator demographics and regresses a scalar disagreement score. In contrast, distributional modeling treats disagreement itself as the target, enabling prediction even when not all annotators are observed by directly parameterizing an item-level label distribution

$$p(y | x_i) \quad (7)$$

y is the expected distribution of annotators that would be produced by the population, for instance x_i . Distributional models typically optimize a composite objective (Parappan and Henao, 2025; Weerasooriya et al., 2023), combining supervised alignment with individual annotations and a divergence-based term (e.g., KL or JSD) that aligns predicted and empirical item-level distributions:

$$\min_{\theta} \underbrace{\sum_i \sum_{j \in \mathcal{A}_i} \ell(y_{ij}, p_{\theta}(y | x_i, a_j))}_{\text{annotator-level alignment}} + \lambda \underbrace{\sum_i \mathcal{D}(\tilde{p}_i(y) || p_{\theta}(y | x_i))}_{\text{distributional alignment}} \quad (8)$$

DisCo (Weerasooriya et al., 2023) models disagreement at the response, item, and population levels, aggregating over annotators at inference to estimate population-level disagreement. The Learning Subjective Label Distribution framework (Parappan and Henao, 2025) extends this by conditioning distributions on sociodemographic attributes and semantic perspectives, enabling subgroup-specific disagreement estimates. Collectively, these approaches move toward unified representations of individual-, group-, and population-level variation.

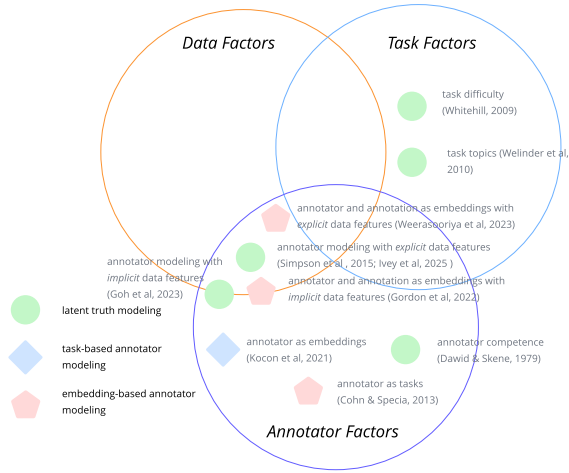


Figure 2: The distribution of the three main methods across our identified taxonomy of sources of disagreement. The lack of work modeling data and task factors point to future directions.

4 Mapping Models to Sources of Disagreement

We map the three model families onto our taxonomy of disagreement sources in Figure 2. Latent truth models span multiple regions, capturing task factors (e.g., difficulty), annotator factors (e.g., competence), and approaches that condition annotator behavior on the input. Within the intersection of annotator and data modeling, we distinguish between implicit and explicit data features. Implicit approaches incorporate the input via shared encoders or embeddings without explicitly modeling item ambiguity as a source of disagreement (Goh et al., 2023; Raykar et al., 2010; Rodrigues and Pereira, 2018; Gordon et al., 2022). For example, CrowdLayer (Rodrigues and Pereira, 2018) and embedding-based models (Gordon et al., 2022) condition predictions on text but treat it primarily as contextual information. In contrast, explicit data modeling treats properties of the input itself as a structured source of annotator variation (Simpson et al., 2015; Ivey et al., 2025). For instance, Ivey et al. (2025) condition annotator competence on item features to capture systematic subgroup responses. This category also includes distributional approaches that directly predict item-level disagreement distributions, explicitly representing data ambiguity (Weerasooriya et al., 2023; Parappan and Henao, 2025). Notably, the relative scarcity of work modeling data factors alone—or jointly integrating data with task and annotator factors—highlights promising directions for future research.

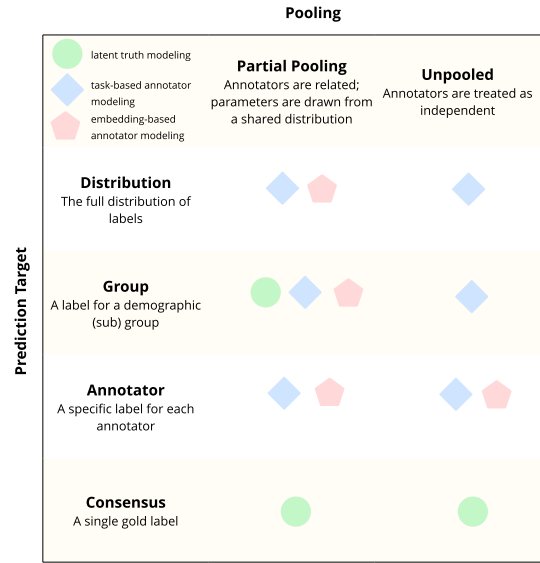


Figure 3: Prediction targets and pooling assumptions across methods for learning from annotator disagreement. The table maps existing approaches by what they predict (consensus labels, individual annotator responses, group-level outputs, or full disagreement distributions) and how they pool information across annotators. Fig. 4 is a more detailed table with representative work.

5 Prediction Targets and Pooling Structures

We introduce a synthesis table (Figure 3) organized along two dimensions: prediction target and pooling structure. Prediction targets include consensus (a single label per item), annotator (one prediction per annotator), group (one prediction per annotator group), and distribution (the full label distribution across annotators). Pooling (Paun et al., 2018c) is generalized across three model classes. Group pooling forces all annotators to share a single labeling function. Unpooled models treat annotators independently. Partial pooling encodes structured relationships among annotators or between annotators and the population (or subpopulations), though it is instantiated differently across paradigms. In Bayesian latent-truth models, annotator reliabilities are drawn from a shared population prior. Gaussian process multi-task models (Cohn and Specia, 2013) treat annotators as correlated tasks, while neural multi-task models (Rodrigues and Pereira, 2018) use shared backbones with annotator-specific heads. Embedding-based approaches (Mokhberian et al., 2024) instead represent annotators as points in a latent population space. Notably, none of the

486 surveyed disagreement-aware models employ full
487 group pooling. While group pooling corresponds to
488 majority vote or fixed aggregation, once annotator
489 behavior is modeled, methods necessarily adopt un-
490 pooled or partially pooled structures. We highlight
491 these temporal trends:

492 **The field is gradually formalizing disagreement**
493 **as an object of prediction, not a source of noise.**

494 The field of modeling annotation disagreement has
495 moved from consensus prediction to individual an-
496 notator modeling, to group-aware modeling, and re-
497 cently to population or subgroup distribution mod-
498 eling. There is a consistent upward shift on the
499 prediction target axis of the table over time.

500 **Across the pooling structures, the field has**
501 **moved from unpooled to partial pooling.**

502 The most successful models use partial pooling,
503 such as hierarchical priors (Paun et al., 2018c),
504 inter-annotator kernels (Cohn and Specia, 2013),
505 demographic-informed MoE (Xu et al., 2025), and
506 subgroup-aware distributions (Parappan and Henao,
507 2025), moving toward improved generalization and
508 fairness. The partial pooling of annotator or group
509 is the most expressive and successful. Most exist-
510 ing methods fall into these two cells. This combi-
511 nation would give generalization to new annotators,
512 the ability to reason about groups, and the ability
513 to scale with sparse labels.

514 **Few works actually model disagreement dis-**
515 **tributions, but there is a convergence toward**
516 **distribution-level modeling** in newer architectures
517 (Fornaciari et al. (2021) as a task-based model,
518 Parappan and Henao (2025); Weerasooriya et al.
519 (2023) as embedding-based models). Predicting
520 how much people disagree is becoming as impor-
521 tant as predicting what label they choose.

522 6 Evaluation Metrics

523 Evaluation metrics for modeling disagreement
524 broadly fall into two classes: those that compare
525 the predictions to the true annotations, and those
526 that evaluate annotator behavior, aiming to recover
527 latent structure such as inter-annotator relation-
528 ships. Hard metrics (accuracy, precision, recall)
529 assume a single gold label—problematic for sub-
530 jective or ambiguous tasks (Poesio and Artstein,
531 2005; Plank et al., 2014)—motivating a shift to-
532 ward probabilistic and soft metrics. Aggregation
533 models such as Dawid–Skene (Dawid and Skene,
534 1979) and MACE (Hovy et al., 2013) use like-
535 lihood or accuracy against inferred latent truth,

536 while latent-truth and reliability models adopt
537 ROC–AUC, negative log probability, predictive
538 density, and entropy-based measures (Raykar et al.,
539 2010; Paun et al., 2018c; Simpson et al., 2015; Ivey
540 et al., 2025). As models predict heterogeneous
541 judgments, evaluation extends to per-annotator met-
542 rics (MAE, RMSE, macro F1) (Cohn and Spe-
543 cia, 2013; Mostafazadeh Davani et al., 2022) and
544 group-level/distributional divergences (KL, JSD)
545 (Mokhberian et al., 2024; Gordon et al., 2022;
546 Weerasooriya et al., 2023). However, recent work
547 (Rizzi et al., 2024) has shown that Cross Entropy
548 (and to a lesser extent divergence-based measures)
549 may violate desirable properties like symmetry or
550 fair penalization, potentially obscuring meaningful
551 differences between models that aim to reproduce
552 disagreement distributions.

553 A complementary set of metrics evaluates anno-
554 tator quality directly. Likelihood-based measures
555 and annotator–model correlations assess inferred
556 behavior (Passonneau and Carpenter, 2014; Paun
557 et al., 2018c), while agreement metrics like Krip-
558 pendorff’s α and Cohen’s κ quantify consistency
559 but conflate ambiguity with unreliability (Krippen-
560 dorff, 1980; Viera and Garrett, 2005). CrowdTruth
561 (Inel et al., 2014) treats disagreement as signal
562 via worker-, item-, and label-level metrics, later
563 extended to model their interactions (Dumitrache
564 et al., 2018). Recent work targets latent structure
565 more directly, using intra-annotator consistency
566 (Gordon et al., 2021) and relational metrics (DIC,
567 BAE) to test whether models preserve agreement
568 geometry (Zhang et al., 2025a). Overall, evaluation
569 increasingly emphasizes models’ preservation of
570 disagreement structure, not merely match labels.

571 **Evaluating Fairness.** Fairness becomes relevant
572 when disagreement is systematic and socially struc-
573 tured, reflecting genuine differences in perspective
574 tied to demographics or lived experience (Aroyo
575 and Welty, 2013; Pavlick and Kwiatkowski, 2019).
576 Collapsing annotations into a single consensus
577 can erase minority viewpoints, rendering pool-
578 ing a normative—rather than purely technical—
579 choice. Although fairness is increasingly invoked
580 in disagreement-aware modeling, most evaluations
581 remain descriptive (Binns, 2018). Disaggregated
582 and group-level metrics reveal disparities (Gordon
583 et al., 2022; Xu et al., 2025), but do not spec-
584 ify what outcomes should count as fair. Diag-
585 nostic frameworks such as AART (Mokhberian
586 et al., 2024) and PERSEVAL (Lo et al., 2025) re-
587 port parity gaps without explicit normative ground-

ing. This gap motivates adapting fairness-in-ML frameworks—such as statistical parity, equalized opportunity, and subgroup fairness (Caton and Haas, 2024)—to disagreement-aware modeling, moving toward explicit normative reasoning about what kind of fairness models ought to achieve.

7 Challenges and Future Work

LLM-based simulation of annotator variation. Persona prompting and demographic conditioning with LLMs are increasingly used to simulate annotator variation, but remain methodologically limited. These approaches are largely unpooled (e.g., personas are simulated independently), highly sensitive to prompts and model choice, and risk amplifying stereotypes (Lee et al., 2023; Durmus et al., 2023; Santurkar et al., 2023; Blodgett et al., 2020). Empirically, persona variables explain only a small fraction of variance (Hu and Collier, 2024), and fine-tuning with demographic metadata appears to rely more on annotator-specific signals than generalizable group structure (Orlikowski et al., 2025). More broadly, LLM judgments compress human disagreement and rely on opaque priors, raising concerns about reliability and epistemic validity of replacing human annotation labor (Durmus et al., 2023; Cazzaniga et al., 2024).

Modeling task-induced disagreement. While annotator- and data-level variation are increasingly modeled, task-level sources of disagreement remain underexplored. Instruction phrasing, label schema design, presentation order, and interface choices can systematically shape disagreement, yet are rarely represented explicitly. Prior work shows that task design strongly influences annotation outcomes (Zhang et al., 2025b; Dsouza and Kovatchev, 2025), but most models incorporate task information only indirectly. Future work could embed task schemas directly, enabling models to capture how task framing interacts with annotators and data.

Integrating task, annotator, and data variation. Most existing methods focus primarily on annotator bias as the source of disagreement. A few latent-truth models and embedding models jointly consider annotator and task factors, but task-based approaches typically do not. Developing models that integrate all three sources, task, annotator, and data, remains an open challenge, partly constrained by data availability.

Scarcity of detailed annotation data. Modeling nuanced disagreement, particularly for minoritized

groups, is constrained by sparse annotator metadata and noisy labels, highlighting the need for richer, disaggregated annotation data.

Interpretability of disagreement-aware models. Explainability remains underdeveloped for disagreement-aware NLP models. Most methods for model interpretability target single-label predictions and do not explain why annotators or groups disagree (Molnar, 2020). While recent multi-annotator models offer feature- or attention-based explanations, these are rarely validated against real annotator behavior. This gap is especially salient for demographic-aware and mixture-of-experts models, where expert specialization may reflect spurious correlations. Emerging metrics such as Behavior Alignment Explainability (Zhang et al., 2025a) point toward more principled evaluation, but robust disagreement-aware interpretability remains an open problem.

Generalization and practical tradeoffs. Annotator disagreement is highly dataset- and domain-specific, and evidence on the benefits of demographic conditioning is mixed (Gordon et al., 2022; Orlikowski et al., 2023). Models risk overfitting to dataset-specific noise, suggesting a need for cross-dataset and cross-domain training. Perspectivist models introduce tradeoffs: they require richer annotations and greater computation, complicate filtering decisions, and raise normative questions about whose perspectives should be preserved. As a result, adopting disagreement-aware modeling is not merely a technical choice, but a value-laden one balancing fairness, interpretability, efficiency, and annotation labor conditions (Valette, 2024).

8 Conclusion

We present a unified view of disagreement-aware modeling in NLP, focusing on disagreement as a meaningful signal rather than noise. We introduce a taxonomy of data-, task-, and annotator-driven sources and trace the shift from aggregation to multi-annotator and persona-conditioned models. Our synthesis highlights key gaps, including limited distributional modeling and largely descriptive fairness evaluation. We also highlight future directions, such as integrating multiple sources of disagreement and developing interpretable, normatively grounded evaluation frameworks.

9 Limitations

We recognize that our work has several limitations. While thorough, decades of work in disagreement make a fully exhaustive survey difficult. Our survey reviews over 120 papers, and our selection reflects methodological relevance to disagreement-aware modeling. As a result, some related approaches, such as LLM simulations of annotator personas, or other methods that engage with human perspectives tangentially rather than as a modeling target, have been excluded. Additionally, the taxonomy we propose is a simplification. In real-world settings, data, task, and annotator factors sometimes cannot be separated cleanly. Future approaches, particularly those that integrate the interaction between various factors, might require extensions to this taxonomy. Lastly, our discussion of open challenges and future directions reflects our interpretations of the most important gaps in the literature. These should not be interpreted as a comprehensive or normative research agenda. Other researchers might prioritize different directions, such as efficiency, model deployment concerns, or annotation practice design to produce meaning disagreement.

A key limitation of the disagreement-aware literature (and thus our synthesis) is that disagreement is often treated as an observed property of datasets rather than as an outcome of annotation practice design. Observed disagreement may conflate meaningful perspective variation with annotator factors such as confusion and fatigue. This limits cross-paper comparability. Developing and standardizing annotation practices that intentionally elicit meaningful disagreement remains an important direction for future work.

References

Gavin Abercrombie, Tanvi Dinkar, Amanda Cercas Curry, Verena Rieser, and Dirk Hovy. 2025. [Consistency is key: Disentangling label variation in natural language processing with intra-annotator agreement](#). *Preprint*, arXiv:2301.10684.

Lora Aroyo and Chris Welty. 2013. [Crowd truth: Harnessing disagreement in crowdsourcing a relation extraction gold standard](#).

Lora Aroyo and Chris Welty. 2015. [Truth is a lie: Crowd truth and the seven myths of human annotation](#). *AI Magazine*, 36(1):15–24.

Valerio Basile, Michael Fell, Tommaso Fornaciari, Dirk Hovy, Silviu Paun, Barbara Plank, Massimo Poesio, and Alexandra Uma. 2021. [We need to consider](#)

[disagreement in evaluation](#). In *Proceedings of the 1st Workshop on Benchmarking: Past, Present and Future*, pages 15–21, Online. Association for Computational Linguistics.

Jacob Beck, Stephanie Eckman, Bolei Ma, Rob Chew, and Frauke Kreuter. 2024. [Order effects in annotation tasks: Further evidence of annotation sensitivity](#). In *Proceedings of the 1st Workshop on Uncertainty-Aware NLP (UncertaiNLP 2024)*, pages 81–86, St Julians, Malta. Association for Computational Linguistics.

Laura Biester, Vanita Sharma, Ashkan Kazemi, Naihao Deng, Steven Wilson, and Rada Mihalcea. 2022. [Analyzing the effects of annotator gender across NLP tasks](#). In *Proceedings of the 1st Workshop on Perspectivist Approaches to NLP @LREC2022*, pages 10–19, Marseille, France. European Language Resources Association.

Reuben Binns. 2018. [Fairness in algorithmic decision-making: Lessons from political philosophy](#). In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 149–159. ACM.

Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. [Language \(technology\) is power: A critical survey of “bias” in NLP](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. Association for Computational Linguistics.

Federico Cabitza, Andrea Campagner, and Valerio Basile. 2023. [Toward a perspectivist turn in ground truthing for predictive computing](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(6):6860–6868.

Z. Cao, E. Chen, Y. Huang, S. Shen, and Z. Huang. 2023. [Learning from crowds with annotation reliability](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2103–2107.

Simon Caton and Christian Haas. 2024. [Fairness in machine learning: A survey](#). *ACM Comput. Surv.*, 56(7).

Mauro Cazzaniga, Florence Jaumotte, Longji Li, Giovanni Melina, Augustus J. Panton, Carlo Pizzinelli, Emma J. Rockall, and Marina Mendes Tavares. 2024. [Gen-ai: Artificial intelligence and the future of work](#). Technical Report 2024/001, International Monetary Fund.

Jonathan P. Chang, Saleema Amershi, and Jaime Teevan. 2017. [Revolt: Collaborative crowdsourcing for labeling tasks](#). In *Proceedings of CHI*.

Z. Chu, J. Ma, and H. Wang. 2021. [Learning from crowds by modeling common confusions](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 5832–5840.

790	Berta Chulvi, Lara Fontanella, Roberto Labadie	2.0: Quality metrics for crowdsourcing with disagree-	847
791	Tamayo, and Paolo Rosso. 2023. Social or individual	ment. <i>Preprint</i> , arXiv:1808.06080.	848
792	disagreement? perspectivism in the annotation of		
793	sexist jokes . In <i>NLPerspectives@ECAI</i> .		
794	Elizabeth Clark, Tal August, Sofia Serrano, Nikita	Esin Durmus, Karina Nyugen, Thomas I. Liao,	849
795	Haduong, Suchin Gururangan, and Noah A. Smith.	Nicholas Schiefer, Amanda Askell, Anton Bakhtin,	850
796	2021. All that’s ‘human’ is not gold: Evaluating	Carol Chen, Zac Hatfield-Dodds, Danny Hernan-	851
797	human evaluation of generated text . In <i>Proceedings</i>	dez, Nicholas Joseph, Liane Lovitt, Sam McCand-	852
798	<i>of the 59th Annual Meeting of the Association for</i>	lisch, Orowa Sikder, Alex Tamkin, Janel Thamkul,	853
799	<i>Computational Linguistics and the 11th International</i>	Jared Kaplan, Jack Clark, and Deep Ganguli. 2023.	854
800	<i>Joint Conference on Natural Language Processing</i>	Towards measuring the representation of subjective	855
801	<i>(Volume 1: Long Papers)</i> , pages 7282–7296, Online.	global opinions in language models . <i>CoRR</i> ,	856
802	Association for Computational Linguistics.	abs/2306.16388.	857
803	Trevor Cohn and Lucia Specia. 2013. Modelling an-	Jing Fan, Guoliang Li, Beng Chin Ooi, Kian-Lee Tan,	858
804	notator bias with multi-task Gaussian processes: An	and Jun Feng. 2015. icrowd: An adaptive crowd-	859
805	application to machine translation quality estimation .	sourcing framework . In <i>Proceedings of the 2015</i>	860
806	In <i>Proceedings of the 51st Annual Meeting of the</i>	<i>ACM SIGMOD International Conference on Manage-</i>	861
807	<i>Association for Computational Linguistics (Volume</i>	<i>ment of Data</i> , pages 1015–1030, Melbourne, Victoria,	862
808	<i>1: Long Papers)</i> , pages 32–42, Sofia, Bulgaria. Asso-	Australia. Association for Computing Machinery.	863
809	ciation for Computational Linguistics.		
810	Aida Davani, Mark Díaz, Dylan Baker, and Vinodku-	Eve Fleisig, Rediet Abebe, and Dan Klein. 2023. When	864
811	mar Prabhakaran. 2024. Disentangling perceptions	the majority is wrong: Modeling annotator disagree-	865
812	of offensiveness: Cultural and moral correlates . In	ment for subjective tasks . In <i>Proceedings of the 2023</i>	866
813	<i>Proceedings of the 2024 ACM Conference on Fair-</i>	<i>Conference on Empirical Methods in Natural Lan-</i>	867
814	<i>ness, Accountability, and Transparency</i> , FAccT ’24,	<i>guage Processing</i> , pages 6715–6726, Singapore. As-	868
815	page 2007–2021, New York, NY, USA. Association	sociation for Computational Linguistics.	869
816	for Computing Machinery.		
817	A. P. Dawid and A. M. Skene. 1979. Maximum likeli-	Tommaso Fornaciari, Alexandra Uma, Silviu Paun, Bar-	870
818	hood estimation of observer error-rates using the em	bara Plank, Dirk Hovy, and Massimo Poesio. 2021.	871
819	algorithm . <i>Applied Statistics</i> , 28(1):20–28.	Beyond black & white: Leveraging annotator dis-	872
820	Gianluca Demartini, Djellel Eddine Difallah, and	agreement via soft-label multi-task learning . In <i>Pro-</i>	873
821	Philippe Cudré-Mauroux. 2012. Zencrowd: Leverag-	<i>ceedings of the 2021 Conference of the North Amer-</i>	874
822	ing probabilistic reasoning and crowdsourcing tech-	<i>ican Chapter of the Association for Computational</i>	875
823	niques for large-scale entity linking . In <i>Proceedings</i>	<i>Linguistics: Human Language Technologies</i> , pages	876
824	<i>of the 21st International Conference on World Wide</i>	2591–2597, Online. Association for Computational	877
825	<i>Web (WWW)</i> , pages 469–478, Lyon, France. Associa-	Linguistics.	878
826	tion for Computing Machinery.		
827	Naihao Deng, Xinliang Zhang, Siyang Liu, Winston Wu,	Simona Frenda, Valerio Basile, Tommaso Caselli, and	879
828	Lu Wang, and Rada Mihalcea. 2023. You are what	Viviana Patti. 2024. Perspectivist approaches to nat-	880
829	you annotate: Towards better models through anno-	tural language processing: A survey . <i>Natural Lan-</i>	881
830	tator representations . In <i>Findings of the Association</i>	<i>guage Processing and Artificial Intelligence</i> . Survey	882
831	<i>for Computational Linguistics: EMNLP 2023</i> , pages	paper.	883
832	12475–12498, Singapore. Association for Computa-	Mirta Galesic and Michael Bosnjak. 2009. Effects of	884
833	tional Linguistics.	questionnaire length on participation and indicators	885
834	Yi Ding, Jacob You, Tonja-Katrin Machulla, Jennifer	of response quality in a web survey. <i>Public Opinion</i>	886
835	Jacobs, Pradeep Sen, and Tobias Höllerer. 2022.	<i>Quarterly</i> , 73(2):349–360.	887
836	Impact of annotator demographics on sentiment	Z. Gao, F.-K. Sun, M. Yang, S. Ren, Z. Xiong, M. En-	888
837	dataset labeling . <i>Proc. ACM Hum.-Comput. Inter-</i>	geler, A. Burazer, L. Wildling, L. Daniel, and D. S.	889
838	<i>act.</i> , 6(CSCW2).	Boning. 2022. Learning from multiple annotator	890
839	Russel Dsouza and Venelin Kovatchev. 2025. Sources	noisy labels via sample-wise label fusion . In <i>Pro-</i>	891
840	of disagreement in data for LLM instruction tuning .	<i>ceedings of the European Conference on Computer</i>	892
841	In <i>Proceedings of Context and Meaning: Navigating</i>	<i>Vision (ECCV)</i> , pages 407–422. Springer.	893
842	<i>Disagreements in NLP Annotation</i> , pages 20–32, Abu	Aparna Garimella, Carmen Banea, Dirk Hovy, and	894
843	Dhabi, UAE. International Committee on Computa-	Rada Mihalcea. 2019. Women’s syntactic resilience	895
844	tional Linguistics.	and men’s grammatical luck: Gender-bias in part-of-	896
845	Anca Dumitrache, Oana Inel, Lora Aroyo, Benjamin	speech tagging and dependency parsing . In <i>Proceed-</i>	897
846	Timmermans, and Chris Welty. 2018. Crowdtruth	<i>ings of the 57th Annual Meeting of the Association for</i>	898
		<i>Computational Linguistics</i> , pages 3493–3498, Flo-	899
		rence, Italy. Association for Computational Linguis-	900
		tics.	901

902	Mor Geva, Yoav Goldberg, and Jonathan Berant. 2019.	Uthman Jinadu and Yi Ding. 2024.	959
903	Are we modeling the task or the annotator? an investigation of annotator bias in natural language understanding datasets.	Noise correction on subjective datasets.	960
904	In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1161–1166, Hong Kong, China. Association for Computational Linguistics.	David R. Karger, Sewoong Oh, and Devavrat Shah. 2011.	961
905		Iterative learning for reliable crowdsourcing systems.	962
906		In <i>Advances in Neural Information Processing Systems</i> , volume 24, pages 1953–1961.	963
907		Jan Kocoń, Marcin Gruza, Julita Bielaniewicz, Damian Grimling, Kamil Kanclerz, Piotr Miłkowski, and Przemysław Kazienko. 2021.	964
908		Learning personal human biases and representations for subjective tasks in natural language processing.	965
909	Hui Wen Goh, Ulyana Tkachenko, and Jonas Mueller. 2023.	In <i>2021 IEEE International Conference on Data Mining (ICDM)</i> , pages 1168–1173.	966
910	Crowdlab: Supervised learning to infer consensus labels and quality scores for data with multiple annotators.		967
911	<i>Preprint</i> , arXiv:2210.06812.		968
912			969
913			970
914			971
915	Mitchell L. Gordon, Michelle S. Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S. Bernstein. 2022.	Klaus Krippendorff. 1980.	972
916	Jury learning: Integrating dissenting voices into machine learning models.	<i>Content Analysis: An Introduction to Its Methodology.</i> Sage Publications.	973
917	In <i>CHI Conference on Human Factors in Computing Systems, CHI '22</i> , page 1–19. ACM.		
918		Jon A. Krosnick, Sowmya Narayan, and Wendy R. Smith. 1996.	974
919		Satisficing in surveys: Initial evidence.	975
920		<i>New Directions for Evaluation</i> , 1996(70):29–44.	976
921			
922	Mitchell L. Gordon, Kaitlyn Zhou, Kayur Patel, Tatsunori Hashimoto, and Michael S. Bernstein. 2021.	Himabindu Lakkaraju, Jure Leskovec, Jon Kleinberg, and Sendhil Mullainathan. 2015.	977
923	The disagreement deconvolution: Bringing machine learning performance metrics in line with reality.	A bayesian framework for modeling human evaluations.	978
924	In <i>Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, CHI '21</i> , New York, NY, USA. Association for Computing Machinery.	In <i>Proceedings of the 2015 SIAM International Conference on Data Mining (SDM)</i> , pages 181–189.	979
925			980
926			981
927			
928		Savannah Larimore, Ian Kennedy, Breon Haskett, and Alina Arseniev-Koehler. 2021.	982
929	M. Guan, V. Gulshan, A. Dai, and G. Hinton. 2018.	Reconsidering annotator disagreement about racist language: Noise or signal?	983
930	Who said what: Modeling individual labelers improves classification.	In <i>Proceedings of the Ninth International Workshop on Natural Language Processing for Social Media</i> , pages 81–90, Online. Association for Computational Linguistics.	984
931	In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 32.		985
932			986
933	Tao Han, Huaixuan Shi, Xinyi Ding, Xiao Ma, Huamao Gu, and Yili Fang. 2025.	Noah Lee, Na Min An, and James Thorne. 2023.	987
934	Mixture of experts based multi-task supervise learning from crowds.	Can large language models capture dissenting human voices?	988
935	<i>Preprint</i> , arXiv:2407.13268.	In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4569–4585, Singapore. Association for Computational Linguistics.	989
936			990
937	Dirk Hovy, Taylor Berg-Kirkpatrick, Ashish Vaswani, and Eduard Hovy. 2013.		991
938	Learning whom to trust with MACE.		992
939	In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1120–1130, Atlanta, Georgia. Association for Computational Linguistics.		993
940			994
941		Qiang Liu, Jian Peng, and Alexander T. Ihler. 2012.	995
942		Variational inference for crowdsourcing.	996
943		In <i>Advances in Neural Information Processing Systems</i> , volume 25, pages 692–700.	997
944	Tiancheng Hu and Nigel Collier. 2024.		998
945	Quantifying the persona effect in llm simulations.		
946	<i>Preprint</i> , arXiv:2402.10811.	Soda Marem Lo, Silvia Casola, Erhan Sezerer, Valerio Basile, Franco Sansonetti, Antonio Uva, and Davide Bernardi. 2025.	999
947		PERSEVAL: A framework for perspectivist classification evaluation.	1000
948	Daniil Ignatev, Denis Paperno, and Massimo Poesio. 2025.	In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 22345–22370, Suzhou, China. Association for Computational Linguistics.	1001
949	Hypernetworks for perspectivist adaptation.		1002
950	<i>Preprint</i> , arXiv:2510.13259.		1003
951			1004
952	Oana Inel, K Khamkham, T Cristea, A Dumitrache, H Rutjes, R-J Sips, J van Ossenbruggen, and K Crowston. 2014.		1005
953	Crowdtruth: Machine-human computation framework for harnessing disagreement in gathering annotated data.		1006
954	In <i>The Semantic Web – ISWC 2014</i> , volume 8797 of <i>Lecture Notes in Computer Science</i> , pages 486–504, Cham. Springer.	Yiwei Luo, Dallas Card, and Dan Jurafsky. 2021.	1007
955		Detecting stance in media on global warming.	1008
956		<i>Preprint</i> , arXiv:2010.15149.	1009
957	Jonathan Ivey, Susan Gauch, and David Jurgens. 2025.		
958	Nutmeg: Separating signal from noise in annotator disagreement.	Fenglong Ma, Yaliang Li, Qi Li, Minghui Qiu, Jing Gao, Shi Zhi, Lu Su, Bo Zhao, Heng Ji, and Jiawei Han. 2015.	1010
	<i>Preprint</i> , arXiv:2507.18890.	Faitcrowd: Fine grained truth discovery for crowdsourced data aggregation.	1011
		In <i>Proceedings</i>	1012
			1013

1126	F. Rodrigues and F. Pereira. 2018. Deep learning from crowds. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages 1611–1618.	<i>LREC-COLING 2024</i> , pages 111–115, Torino, Italia. ELRA and ICCL.	1183 1184
1129	Marta Sandri, Elisa Leonardelli, Sara Tonelli, and Elisabetta Jezek. 2023. Why don't you do it right? analysing annotators' disagreement in subjective tasks . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 2428–2441, Dubrovnik, Croatia. Association for Computational Linguistics.	Anthony J. Viera and Joanne M. Garrett. 2005. Understanding interobserver agreement: The kappa statistic. <i>Family Medicine</i> , 37(5):360–363.	1185 1186 1187
1132	Yisi Sang and Jeffrey Stanton. 2021. The origin and value of disagreement among data labelers: A case study of the individual difference in hate speech annotation . <i>Preprint</i> , arXiv:2112.04030.	Ruyuan Wan, Jaehyung Kim, and Dongyeop Kang. 2023. Everyone's voice matters: Quantifying annotation disagreement using demographic information . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 37, pages 14523–14530.	1188 1189 1190 1191 1192
1136	Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? In <i>Proceedings of the 40th International Conference on Machine Learning (ICML)</i> , volume 202, pages 29971–30004, Honolulu, Hawaii, USA. PMLR.	Tharindu Cyril Weerasooriya, Alexander Ororbia, Raj Bhensadadia, Ashiqur KhudaBukhsh, and Christopher Homan. 2023. Disagreement matters: Preserving label diversity by jointly modeling item and annotator label distributions with DisCo . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 4679–4695, Toronto, Canada. Association for Computational Linguistics.	1193 1194 1195 1196 1197 1198 1199 1200
1140	Maarten Sap, Swabha Swayamdipta, Laura Vianna, Xuhui Zhou, Yejin Choi, and Noah A. Smith. 2022. Annotators with attitudes: How annotator beliefs and identities bias toxic language detection . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5884–5906, Seattle, United States. Association for Computational Linguistics.	H. Wei, R. Xie, L. Feng, B. Han, and B. An. 2022. Deep learning from multiple noisy annotators as a union. <i>IEEE Transactions on Neural Networks and Learning Systems</i> .	1201 1202 1203 1204
1141	Edgar D. Simpson, Matteo Venanzi, Pushmeet Kohli, John Guiver, Gianluca Kazai, and Milad Shokouhi. 2015. Language understanding in the wild: Combining crowdsourcing and machine learning . In <i>Proceedings of the 24th International Conference on World Wide Web (WWW 2015)</i> , pages 992–1002, Florence, Italy. Association for Computing Machinery.	Peter Welinder, Steve Branson, Pietro Perona, and Serge J. Belongie. 2010. The multidimensional wisdom of crowds . In <i>Advances in Neural Information Processing Systems</i> , volume 23, pages 2424–2432.	1205 1206 1207 1208
1142	Fritz Strack. 1992. “order effects” in survey research: Activation and information functions of preceding questions . In Norbert Schwarz and Seymour Sudman, editors, <i>Context Effects in Social and Psychological Research</i> , chapter 3. Springer, New York, NY.	Jacob Whitehill, Ting fan Wu, Jacob Bergsma, Javier R. Movellan, and Paul L. Ruvolo. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise . In <i>Advances in Neural Information Processing Systems</i> , volume 22, pages 2035–2043.	1209 1210 1211 1212 1213 1214
1143	R. Tanno, A. Saeedi, S. Sankaranarayanan, D. C. Alexander, and N. Silberman. 2019. Learning from noisy labels by regularized estimation of annotator confusion. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 11244–11253.	Michael Wojatzki, Tobias Horsmann, Darina Gold, and Torsten Zesch. 2018. Do women perceive hate differently: Examining the relationship between hate speech, gender, and agreement judgments . In <i>Conference on Natural Language Processing</i> .	1215 1216 1217 1218 1219
1144	Alexandra Uma, Tommaso Fornaciari, Silviu Paun, Benjamin Paul Chamberlain, Dirk Hovy, Barbara Plank, and Dean K Simonton. 2021. Learning from disagreement: A survey. In <i>Transactions of the Association for Computational Linguistics (ACL)</i> , volume 9, pages 1408–1424.	Jin Xu, Mariët Theune, and Daniel Braun. 2024a. Leveraging annotator disagreement for text classification . In <i>Proceedings of the 7th International Conference on Natural Language and Speech Processing (IC-NLSP 2024)</i> , pages 1–10, Trento. Association for Computational Linguistics.	1220 1221 1222 1223 1224 1225
1145	Mathieu Valette. 2024. What does perspectivism mean? an ethical and methodological counter-criticism . In <i>Proceedings of the 3rd Workshop on Perspectivist Approaches to NLP (NLPerspectives) @</i>	Shanshan Xu, T. Y. S. S Santosh, Oana Ichim, Isabella Risini, Barbara Plank, and Matthias Grabmair. 2024b. From dissonance to insights: Dissecting disagreements in rationale construction for case outcome classification . <i>Preprint</i> , arXiv:2310.11878.	1226 1227 1228 1229 1230
1146		Yinuo Xu, Veronica Derricks, Allison Earl, and David Jurgens. 2025. Modeling annotator disagreement with demographic-aware experts and synthetic perspectives . <i>Preprint</i> , arXiv:2508.02853.	1231 1232 1233 1234
1147		Tianling Yang, Christian Strippel, Alexandra Keiner, Dylan Baker, Alexis Chávez, Krystal Kauffman, Marc Pohl, Caroline Sindere, and Milagros Miceli.	1235 1236 1237

2023. Ethics of data work: Principles for academic data work requesters. Weizenbaum discussion paper, Weizenbaum Institute for the Networked Society.

Liyun Zhang, Jingcheng Ke, Shenli Fan, Xuanmeng Sha, and Zheng Lian. 2025a. [A unified evaluation framework for multi-annotator tendency learning](#). *Preprint*, arXiv:2508.10393.

Michael JQ Zhang, Zhilin Wang, Jena D. Hwang, Yi Dong, Olivier Delalleau, Yejin Choi, Eunsol Choi, Xiang Ren, and Valentina Pyatkin. 2025b. [Diverging preferences: When do annotators disagree and do models know?](#)

Zhao Zhao, Fang Wei, Ming Zhou, Wei Chen, and Wilfred Ng. 2015. [Crowd-selection query processing in crowdsourcing databases: A task-driven approach](#). In *Proceedings of the 18th International Conference on Extending Database Technology (EDBT)*, pages 397–408. OpenProceedings.org.

Yudian Zheng, Guoliang Li, Yuanbing Li, Caihua Shan, and Reynold Cheng. 2017. [Truth inference in crowdsourcing: is the problem solved?](#) *Proc. VLDB Endow.*, 10(5):541–552.

A Appendix

A.1 Survey Scope and Selection

This survey focuses on research in NLP that studies, models, or evaluates annotator disagreement in text-based tasks. We restrict our scope to NLP settings where disagreement is theoretically meaningful, such as toxicity detection, stance analysis, sentiment, and other subjective or socially grounded judgments, rather than domains where disagreement primarily reflects measurement error. Our goal is not to exhaustively catalog all work on annotation variability (of which there is decades of work), but to synthesize lines of research that treat disagreement as a signal relevant to learning, evaluation, or model design, with a slight focus on connecting the most recent work.

We survey over 120 prior works spanning disagreement sources, modeling methods, and evaluation paradigms. Table 1 provides the counts of papers surveyed in each category we focus on. Papers may appear in multiple categories when relevant. We also reviewed approximately 12 papers on large language model-based persona prompting and demographic conditioning; however, because most of this work focuses on simulating personas rather than explicitly modeling annotator structure or disagreement, we discuss these primarily in the Challenges and Future Work section.

We include work that makes a substantive methodological contribution, defined as proposing

Survey Category	# Papers
Sources of Disagreement	25
Latent Truth & Reliability Modeling	22
Task-Based Multi-Annotator Learning	12
Embedding-Based Multi-Annotator Learning	15
Evaluation Metrics	18
Total Unique Papers Surveyed	120

Table 1: Distribution of papers surveyed across categories. Papers may appear in multiple categories.

a learning algorithm, prediction target, or evaluation framework that explicitly accounts for multiple annotators or divergent interpretations. For each of the three modeling families that we identify (latent truth modeling, task-based multi-annotator learning, and embedding-based approaches), we emphasize seminal and highly-cited work (Dawid and Skene, 1979; Cohn and Specia, 2013; Kocooñ et al., 2021), alongside representative recent advances. Our taxonomy of disagreement sources was developed iteratively in tandem with the literature review, with theoretical distinctions refined through engagement with empirical findings.

Taken together, this scope allows us to position our survey as a complement to existing work. Whereas prior surveys are either primarily conceptual (Frenda et al., 2024) or modality-general (Uma et al., 2021), we center NLP as a domain where disagreement is deeply tied to task design, interpretation practices, and annotator identities. By jointly examining modeling approaches, evaluation practices, and their fairness implications within a unified framework, we aim to provide a more technically-grounded synthesis of perspectivist modeling that has not yet been articulated in existing surveys. Finally, to capture emerging trends, we closely follow work from the Workshop on Perspectivist Approaches to NLP¹ and related venues.

A.2 Detailed Synthesis Table of Prediction Targets and Pooling Structures

We provide a detailed table mapping representative work in each modeling family to the axis of prediction target and pooling structures (Fig. 4).

¹<https://nlperspectives.di.unito.it/>

		POOLING	
		PATIAL POOLING	UNPOOLED
PREDICTION TARGET	DISTRIBUTION	<p><i>Jinadu et al. (2024)</i> — subjectivity-tunable loss distinguishing noise vs disagreement.</p> <p><i>Weerasooriya et al. (2023, DisCo)</i> — predicts joint item + annotator label distributions.</p> <p><i>Parappan and Henao (2025)</i> -- predicts label distribution using sociodemographics + LLM-generated perspectives</p>	<p><i>Fornaciari et al. (2021)</i> — MTL with soft labels as ambiguity-based regularization</p>
	GROUP	<p><i>NUTMEG (Ivey et al, 2025)</i> -- Hierarchical Bayesian annotator model with demographic priors</p> <p><i>Fleisig et al. (2024)</i> — stakeholder-aware MTL predicting harmed-group ratings.</p> <p><i>Gordon et al. (2022)</i> — "jury" model combining demographic subgroup embeddings.</p> <p><i>Xu et al. (2025)</i> — DEM-MoE experts specialized by demographic/intersectional traits.</p>	<p><i>Fleisig et al. (2024)</i> — aggregates group-specific predictions for fairness-aware consensus.</p>
	ANNOTATOR	<p><i>Cohn & Specia (2013)</i> — multi-task GP with learned inter-annotator kernel B.</p> <p><i>Rodrigues & Pereira (2018)</i> — Crowd Layer learning bias transforms via backprop.</p> <p><i>Guan et al. (2017)</i> — per-doctor CNN heads with reliability weighting.</p> <p><i>Mokherian et al. (2023, AART)</i> — scalable annotator embeddings + sparsity regularization.</p> <p><i>Gordon et al. (2022)</i></p> <p><i>Xu et al. (2025)</i></p>	<p><i>Davani et al. (2022)</i> — multi-head BERT predicting each annotator's view.</p> <p><i>Kocoi et al. (2021)</i> — introduced annotator embedding for personal bias.</p> <p><i>Deng et al. (2023)</i> — injects annotator + annotation embeddings into Transformer.</p> <p><i>Ignatev et al. (2025)</i> -- parameter-efficient annotator-specific LORA weights</p>
	CONSENSUS	<p><i>Paun et al. (2018)</i> — hierarchical Bayes comparison of pooling schemes.</p> <p><i>Simpson et al. (2015)</i> — BCCWords links linguistic features with bias.</p> <p><i>Goh et al., (2023)</i> — ensemble using task-level classifier as pseudo-annotator for consensus.</p>	<p><i>Dawid & Skene (1979)</i> — EM model inferring latent truth + error rates.</p> <p><i>Raykar et al. (2010)</i> — jointly learns annotator expertise + logistic regressor.</p>

LATENT TRUTH ANNOTATOR MODELING
TASK-BASED MULTI-ANNOTATOR MODELING
EMBEDDING-BASED MULTI-ANNOTATOR MODELING

Figure 4: Prediction targets and pooling assumptions across methods for learning from annotator disagreement. The table maps existing approaches by what they predict (consensus labels, individual annotator responses, group-level outputs, or full disagreement distributions) and how they pool information across annotators. Modeling choices implicitly encode different assumptions about population structure, perspective aggregation with recent work increasingly favoring partial pooling and distributional targets.