# Just Put a Human in the Loop? Investigating LLM-Assisted Annotation for Subjective Tasks

Anonymous ACL submission

### Abstract

LLM use in annotation is becoming overall widespread, and given LLMs' promising performance and speed, putting humans in the loop to simply "review" LLM 004 annotations can be tempting. In subjective tasks with multiple plausible answers, this can 007 impact both evaluation of LLM performance, and analysis using these labels in a social science task downstream. In a pre-registered experiment with 350 unique annotators and 7,000 annotations across 4 conditions, 2 models, and 2 datasets, we find that presenting 012 crowdworkers with LLM-generated annotation suggestions did not make them faster 015 annotators, but did improve their self-reported confidence in the task. More importantly, annotators strongly took the LLM suggestions, 017 significantly changing the label distribution compared to the baseline. We show that when these labels created with LLM assistance are used to evaluate LLM performance, reported model performance significantly increases. We show how changes in label distributions as a result of LLM assistance can affect conclusions drawn by analyzing even "human-approved" LLM-annotated datasets. We believe our work underlines the importance of understanding 027 the impact of LLM-assisted annotation on subjective, qualitative tasks, on the creation of gold data for training and testing, and on the evaluation of NLP systems on subjective tasks.

### 1 Introduction

Large language models (LLMs) are showing impressive performance in many annotation tasks, including subjective tasks common in content moderation and text analysis in the social sciences. Evaluating human annotation of subjective tasks for comparison against LLM annotation performance, either for the task of end-to-end qualitative analysis or for the construction of ground truth for NLP tasks, is difficult in the absence of domain experts. Accordingly, hiring a large number of crowd annotators (often in service of creating a crowd decision) becomes attractive in the evaluation of NLP systems on social science tasks. However, managing and paying crowdworkers can be difficult, and crowdworkers often have varied performance. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

A long line of research explores how AI suggestions can assist qualitative researchers (Jiang et al., 2021; Feuston and Brubaker, 2021; Overney et al., 2024). Labeling text according to a complex qualitative "codebook" is a repetitive, timeconsuming task, and advances in LLM capabilities have made using LLMs in annotation attractive (Wang et al., 2021). Complex, theory-driven text analysis is increasingly being mediated by LLMs (De Paoli, 2023) or LLM-powered tools (Lam et al., 2024).

Given all this progress, LLMs have created opportunities to create annotation pipelines that appear to work off the shelf without fine-tuning, making automated annotation accessible to practitioners with less technical skill. LLMs' reported performance in annotating socially complex topics (Gilardi et al., 2023), sometimes with greater skill than humans (He et al., 2024), potentially opens LLM-based annotation to an even wider range of fields and practices compared to past years. This makes understanding the many ways humans may interact with LLM annotations more important.

With the relative ease of creating LLM annotations for a variety of tasks, there is a temptation to just put a "human in the loop" to check annotations and ensure the model's outputs are "reasonable and reliable" (Wang et al., 2025) in order to approve LLM annotations of a social concept. But given that we know humans are subject to anchoring bias—the bias towards the first option we are presented (Tversky and Kahneman, 1974)—humans may review and confirm LLM annotations that are plausible, but nonetheless significantly change 1) the annotation evaluation process and 2) the out-

172

173

174

175

176

177

178

179

180

181

182

183

184

133

134

135

come the annotations get used for (such as the decision boundary for classification judgments, or the distribution of annotations used in a text analysis), downstream.

084

097

100

101

102

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

128

129

130

131

132

In this work, we ask several questions. First, when provided with (different forms of) LLM assistance, do crowdworker annotators "produce more" by going faster in a complex, subjective annotation task, and does this result in them "understanding less?" Do they find LLM suggestions accurate and helpful, and how often do annotators take these suggestions? If annotations influenced by LLM suggestions are used to construct ground truth on annotation tasks involving complex social concepts, what effect does that have on evaluated performance of these LLMs on those annotation tasks?

We use two community conversation datasets to address these questions, recruiting crowdworkers to annotate the data according to a list of identified themes. These workers are presented with LLMsuggested annotations in a variety of formats. We then study how the crowdworkers use these LLM suggestions in a complex annotation process in service of answering these research questions.

Our findings open up questions regarding the use of AI assistance in qualitative research where annotators act as independent reviewers of AI suggestions and can supposedly still retain full control of the analysis. We believe our work underlines the importance of understanding the impact of LLMassisted annotation on subjective, qualitative tasks, on the creation of gold data for training and testing, and on the evaluation of NLP systems on subjective tasks. In addition to the findings that answer the aforementioned questions, we release a dataset of human and LLM-assisted annotations on two complex qualitative codebooks across a variety of conditions.

### 2 Related Work

LLMs have been increasingly employed in various "human in the loop" AI assistance setups, e.g., in education (Jiang et al., 2024), and coding (Mozannar et al., 2024). LLMs have also been used for annotation tasks, for reasons like decreasing labeling cost (Wang et al., 2021). Ziems et al. (2023) discussed how LLMs are being widely used in computational social science tasks, including subjective tasks. Li et al. (2023) found that the performance of LLMs trained on synthetically generated ground truth data is negatively associated with the subjectivity of the task. More work has shown the challenges of LLM annotation: LLM annotation performance can be highly variable to prompts (Atreja et al., 2024) and can depend on the ordering of choices presented in a prompt. (Wang et al., 2023; Liu et al., 2023).

While using LLMs to create annotations may have "good enough" performance and likely does decrease cost compared to hiring human crowd workers, many tasks still require humans to still be in the loop, especially for subjective annotation and analysis tasks; some methods suggest ways of using humans in the loop to optimize prompts for LLM annotation (Pangakis and Wolken, 2024). Many AI-assisted text annotation platforms have also been developed and published within the HCI space (Overney et al., 2024; Gao et al., 2024, 2023a, 2025). Research shows AI-suggested labeling platforms increase agreement and convergence on qualitative codes, or text labels (Gao et al., 2023b).

There is increasing evidence that humans tend to anchor on the suggestions that AI systems give, which can result in changes to communication and even user opinions (Jakesch et al., 2019). Some work observed how humans anchor on LLM outputs in text analysis tasks, including topic induction. For instance, Choi et al. (2024) show that in a topic generation task, analysts anchor on LLM outputs, resulting in different topic lists depending on whether or not they saw the LLM versions. This illustrates the potential risk of homogenization of insight as a result of AI influence on text analysis. This concern is raised by Messeri and Crockett (2024) regarding AI's influence on science more generally, and the authors discuss LLMs' potential to reduce diversity in human judgment, creating an environment in science where we "produce more but understand less."

Many cite time savings as their motivation for using LLMs as annotators, however, findings have so far been mixed regarding productivity. For example, Bughin shows that while AI can boost coding productivity, there exists a tradeoff between productivity and coding quality, and Overney et al. observed that users with AI support in the Sense-Mate platform spent more time on qualitatively coding data.

Here, we examine how presenting LLMgenerated suggestions to annotators in a complex subjective annotation task affects their understanding of the task as well as their suggestion uptake, with implications for the evaluation of LLM performance on these tasks, even when humans are put

186

207

210 211

212

213

214

215

216

217

218

"in the loop" to review and confirm annotations.

### **3** Data and Codebook

We source two conversation collections from the 187 Fora corpus (Schroeder et al., 2024). In 2022, the 188 NYC Department of Health & Mental Hygiene, the NYC Public Health Corps (PHC), and the national 190 non-profit Cortico recruited over 100 communities 191 to a series of 28 small group dialogues hosted in 192 New York City to understand community resourc-193 ing and vaccine decisions during COVID-19. Fol-194 lowing the conversations, community workers cre-195 ated a codebook of themes of interest to the "NYC" corpus (as we will call it), then labeled quotes from 197 the conversations with the themes denoted by the 198 codebook. Similarly, in 2021, a conversation series 199 in Boston called "Real Talk for Change" was hosted 200 to understand issues in marginalized communities leading up to the 2021 Boston Mayoral election. We use the conversation data and codebook created for this "RTFC" corpus as well. 204

The NYC codebook developed by community partners had seven overarching themes related to health and vaccine decisions, including External Motivations: Friends & Family, Intrinsic Motivations: Not wanting to get the virus, and Role of Community Health Organizations: Health Education & Support, and Vaccine Hesitancy. Each toplevel label had sublabels, for a total of 27 total labels related to the NYC corpus. Similarly, the RTFC corpus has 9 overarching top-level labels, and 41 sublabels relevant to the corpus, such as Safety: Street violence and Housing: Housing affordability. The full codebook for each corpus is available in the appendix. We sampled 200 quotes from the NYC corpus as the main data set for this study, and 200 quotes from the Real Talk for Change corpus as a replication data set. Excerpts had an average length of 592 characters.

## 4 Methods

The experiment compares the annotation of 200 quotes by 5 annotators each according to the codebook for both NYCDOMH and RTFC. First, we create a crowdworker "baseline" for the annotation task without LLM assistance for both corpora. Following typical practices, we construct a ground truth set of labels using a 3/5 majority vote for each label. We then test how this crowd baseline contrasts to annotations provided through three distinct ways of presenting LLM-suggested labels for annotators to review. As such, the four experimental conditions testing how *interface* mediated LLM-assisted annotation suggestions were:

234

235

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

- **Condition 1: Baseline.** Annotators were not shown any LLM-generated suggestions in the annotation task.
- Condition 2: Text-based suggestions. Annotators were presented with the annotation interface, which presented the text to be annotated at the top of the screen. The text "Suggested tags:" was appended after the quote, followed by a list of LLM-generated labels, generated either by GPT-4 or LLaMA, according to the corpus' codebook of labels.
- Condition 3: Text-based suggestions, with AI disclosure Same as Condition 2, except immediately following the quote, the text "Suggested tags from AI:" was appended, followed by the list of LLM-generated labels, pictured in pictured in Figure 8.
- Condition 4: Pre-highlighted labels in interface LLM-generated label suggestions were pre-highlighted on each question in the interface, as pictured in Figure 9.

These conditions provide a sliding scale of assistance to an annotator. The no-assistance baseline in Condition 1 provide the basis for crowd truth labels. Text-based suggestions provide some assistance, but still require the annotators to read the information then integrate the suggestions themselves into decisions for each annotation question. Condition 3 was the same as Condition 2 with the exception of including a disclosure that the suggestions were from "AI." We included this condition to test whether annotators would change their perception of or behavior towards suggestions if they knew the origin of them— either in their perceived quality or in their rate of uptake. Finally, Condition 4 provided annotators with the strongest suggestion, drawing the annotator's attention directly to a colored highlight of the suggested label. These conditions were all shown to annotators through a deployment of the open source annotation interface, Potato (Pei et al., 2023). Screenshots and interface examples are in the appendix.

We ran these conditions on the NYC data with GPT-generated labels, and we replicated label suggestions created by LLaMA (Touvron et al., 2023),

an open-source model. Second, we also create a no-assistance baseline for the RTFC data, and replicated the experiment on label suggestions created by GPT-4 to test generalizability to a different codebook and dataset.

We generated label suggestions by prompting one of two LLMs. For our main experiment, we used OpenAI's API to prompt GPT-4 (OpenAI et al., 2024), model version gpt-4-1106-preview, zero shot, and for our replication study, we used LLaMA 11ama3-8b (Touvron et al., 2023), accessed through the service LlamaAPI. We used the same prompt style and instructions to prompt both GPT-4 and LLaMA, which were prompted once per quotation to produce a list of labels for each of the 200 quotations from the NYC and RTFC corpora. The prompt details for each task are available in the appendix.

### 4.1 Survey experiment

283

290

291

295

304

307

311

312

313

315

In order to test annotation performance under a variety of conditions, we hired crowdworkers to complete our annotation study. We recruited qualified annotators from Prolific, and additional details about recruitment are available in the appendix.

Once annotators accepted the task and had instructions, each annotator was given 20 unique annotations, with 2 randomly assigned "understanding" check questions mixed in. 200 quotes were thus annotated by 5 unique annotators in each of the four experimental conditions. Each annotator participated in just one experimental condition. Each annotator was recommended to spend 20-30 minutes on the annotation task, and spent an average of 35 minutes on the task.

Prior to doing the task, we conducted an exer-316 cise to measure inter-annotator agreement among 317 Prolific workers in our worker pool for this task. In order to do so, we presented 15 unique annota-319 tors with 19 unique quotes from the NYC corpus, and 20 unique annotators with 20 quotes from the 321 RTFC corpus. We then measured inter-rater relia-322 bility on the codes for each corpus using the traditional measure of Krippendorff's alpha ( $\alpha$ ) across annotators, which yielded low to medium levels of agreement across annotators overall, depending on the label. Overall low to medium agreement in this 328 context is unsurprising for both the NYC and RTFC codebooks, given the subjective, complex nature of each task. In order test basic understanding of the task, we ranked quote and label pairs by level of agreement in the IRR task. We selected these 332

extremely high-agreement quote and label pairs as a pool of minimum-threshold understanding questions for annotator, which we used as proxies for basic understanding of the task. The 4 selected test questions in the NYC corpus had 13 or 14 of 15 annotators in agreement with the label, and 15-18 annotators in agreement in the RTFC corpus. We included two randomly selected understanding test questions from the relevant corpus in the question bank for each annotator, which were presented in a random order within the task, and called "understanding" questions for the rest of this study. 333

334

335

337

338

339

340

341

342

343

344

345

346

347

350

351

352

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

372

373

374

375

376

377

379

381

Within the presented task, each annotator would see a quotation from the corpus on the screen. Annotators were prompted to select any labels that applied to the quote, or select none if none applied. For the NYC corpus, there were 7 annotation questions for each quotation, one corresponding to each top-level label. In the RTFC corpus, there were 9 annotation questions for each quotation, corresponding to each top-level label.

In the conditions where LLM-generated suggestions were presented to the annotator, we also included an additional question associated with each quotation, asking annotators if suggestions were "overall," "somewhat," or "not helpful/accurate."

Following all presented quotes, we presented annotators with a post survey in order to understand their perception of the task. Before answering these self-report questions, we told them "Your answer to this question will not reflect on your performance, so please answer honestly." Annotators were given the opportunity to answer the following, rating each on a 1-5 scale, with 5 being the highest:

- How well do you feel you understood this labeling task?
- Overall, how confident do you feel in your answers on this task?
- After doing this task, how well do you understand the concerns and priorities of this community?
- After doing this task, how well could you explain this community's concerns and needs?

We finished the post-survey with some demographic questions to better understand the annotator pool after their answers had already been given, including asking annotators for their race, gender, political orientation.

### 5 Results

We compare outcomes on these annotation tasks across 4 conditions, including 3 assistance condi-

Dataset	Model	Suggestion type	Avrg. minutes
	None	None (Baseline)	25.5
		Text-based	28.4
NYC	GPT	Text-based, AI origin disclosed	30.8
		Pre-filled	30.1
	LLaMA	Text-based	32
PTEC	None	None (Baseline)	34
KIIC	GPT	Text-based	35.4

Table 1: Time spent on 20 annotations across conditions.

tions, 2 corpora (NYC and RTFC), and 2 models that provided suggestions (GPT-4 and LLaMA).

# 5.1 LLM assistance did not decrease annotation time

Contrary to our pre-registered hypothesis, we found annotators in the LLM assistance conditions did not go faster than annotators in the baseline condition. To calculate time spent on the task, we subset to annotators who completed all 20 annotations, and removed test questions and post-survey questions. This gave us time spent on 20 substantive annotations in the task. In two of the LLM assistance conditions for the NYC data, there was a statistically significant increase of 5 minutes spent on the annotation task compared to baseline.

In the assistance conditions, we included an additional short question for each annotation asking annotators to rate suggestion quality. Small time increases in the assistance conditions may be attributable to this additional question we asked. Figure 2 in the appendix shows time variation. This replicates findings in Overney et al., which found that when qualitative coders had access to AIgenerated suggestions for qualitative codes, they actually spent longer on the annotation task than in the baseline condition.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

# 5.2 Assistance improves self-reported understanding of task and content

Despite no increase in time-based productivity outcomes, annotators' self-reported experience of the task improved in many of the assistance conditions. Annotators self-reported higher levels of task understanding, task confidence, community understanding, and ability to explain community needs over the baseline no-assistance condition, and many had statistically significantly higher levels (p < .01) over the baseline condition according to a two-tailed t-test, with medium positive effect sizes calculated using Cohen's D.<sup>1</sup> Improvements over no-assistance baselines were strongest for annotators' self-reported "understanding of the community" and "ability to explain the community's needs," and these are the two outcome measures with strong replication on the RTFC corpus as well.

Because this is a subjective task, self-reports of understanding could theoretically increase with LLM assistance while some measure of "true" understanding on the task could decrease. We found that pass rates on the "understanding" checks we included in the task did not statistically significantly change (either increase or decrease) in the LLM assistance conditions, providing at least basic assurance that providing assistance did not immediately elicit overreliance on the assistancen to the point the basic task was not understood.

# 5.3 Annotators overwhelmingly like and take the suggestions

We converted annotator ratings of suggestion helpfulness and accuracy into numeric values as follows: "Overall helpful/accurate" to a 2, "Somewhat helpful/accurate" to a 1, and "Not helpful/accurate" to a 0. Across conditions, LLM suggestions were rated as between somewhat and very helpful (mean: 1.49, details in appendix). We did not find statistically significant differences in ratings of helpfulness between GPT-4 and Llama, or in the way suggestions were presented (Condition 2: textbased, Condition 3: text-based + AI disclosed, and Condition: pre-filled). Helpfulness of label suggestions was also rated similarly in the NYC and RTFC annotation tasks, suggesting no one model worked better than the other, and the assistance was helpful in two different labeling contexts.

Reflecting this perceived helpfulness and accuracy, we observed strong LLM suggestion uptake across annotators. To observe suggestion uptake rates at the individual level and in contrast to a crowd baseline, we obtained the set intersection of labels applied by the LLM to a particular quote and the set of labels applied by a crowd decision of human annotators to that same quote. Crowd decisions for labels were made by checking if, for each labels, the labels was assigned to the quote by at least 3 (of the 5) annotators, giving a final list of labels for the quote on which at least 3 annotators agreed to its applicability. We repeated this process at a crowd decision threshold of 4 annotators, as well as full consensus of all 5 annotators.

For each crowd decision threshold, we divided the size of the set intersection of labels with LLMsuggested labels by the total number of labels applied by the LLM. This gave a percentage of an446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

423

424

425

426

427

428

<sup>&</sup>lt;sup>1</sup>Details available in the appendix.

notator label overlap with the LLM suggestions at 474 three crowd decision thresholds. Treating crowd 475 decisions on the unassisted baseline condition for 476 the NYC corpus as ground truth, just 40% of labels 477 given by a crowd decision of 3 annotators over-478 lapped with the GPT-4 suggestions. This further 479 dropped sharply to 24% when the crowd threshold 480 is raised to 4 annotators, and to just 8% when raised 481 to full consensus of 5 annotators. 482

483

485

487

488

490

491

492

493

495

497

504

508

510

512

513

514

515

516

517

518

520

521

522

524

Figure 1 shows that in all LLM assistance conditions, run for Conditions 2, 3, and 4 with GPT-4 484 labels on NYC data, the overlap between the crowd ground truth created with LLM assistance and the 486 LLM label set increased dramatically. We display this in terms of different crowd decision thresholds. At a crowd decision threshold of 3/5, crowd labels 489 the average rate of overlap between crowd labels and suggested labels was 81-89% depending on the presentation of suggestions, 54-67% at a crowd decision threshold of 4, and between 24-40% for full crowd consensus of 5. In other words, overlap with 494 LLM suggestions increased 40% at the typical decision threshold of 3 when human annotators were 496 given these suggestions to review, a statistically 498 significant increase according to a two-tailed t test at p = .05. Text-based Llama suggestions resulted 499 in similar results for the NYC corpus, as did GPT-4 suggestions on the RTFC corpus. Details on these findings are in the appendix.

> We also observed consensus agreement of all 5 annotators was significantly more likely in Condition 4, where suggestions were presented most strongly by appearing pre-filled in the interface. We observed that full consensus, or full agreement by all 5 annotators, increase from just 8% in the no assistance baseline to 43% in the pre-highlighted label condition. Using a two-tailed t-test, we find this is a statistically significant increase a p = .001, including Bonferroni correction for multiple comparisons.

#### 5.4 Using human-reviewed, LLM-assisted labels as ground truth significantly inflates reported model performance

Using LLMs to annotate or augment annotations that can be used to train models or evaluate model performance is tempting, given the challenge of scaling annotation, particularly for subjective tasks that may be challenging for crowd workers. How much does model performance appear to improve on these subjective tasks when we use LLMassisted annotations reviewed by humans as ground



Figure 1: Percent overlap of annotator labels with LLMsuggested labels by condition and crowd decision threshold

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

truth?

To examine this, we first created baseline ground truth labels for the 27 NYC labels, by aggregating the 5 annotations made by each annotator into a 3/5 majority vote to assign each label. For each label, we calculate an F1 score of "model performance" using the human crowd labels as ground truth, and compare to either the GPT labels or Llama labels for the entire set of 200 NYC quotations, calculating an F1 score of model performance. Using baseline human ground truth shows overall low performance on these labeling tasks, with similar performance for both GPT-4 and Llama: the average weighted F1 score of GPT-4 performance across all labels is .44 ( $\sigma = .17$ ) when using the human crowd labels as ground truth, and .42 ( $\sigma = .17$ ) for Llama when using the human crowd labels as ground truth. Interestingly, when GPT-4 labels for the NYC corpus are used as ground truth labels and compared to Llama labels, the average F1 score is significantly higher at .62 ( $\sigma = .10$ ) compared to performance when using a human baseline. The individual breakdown of label-level performance is available in Table 2, and breakdowns of GPT-4 versus Llama are available in the appendix in Table 5.

We next aggregated labels from annotators reviewing text-based suggestions from GPT-4 ("GPTassisted") or Llama ("Llama-assisted"). We used a 3/5 majority to approve each label, constructing a new ground truth condition for annotations made with assistance. When using the GPT-assisted ground truth, the average weighted F1 score of

GPT-4 performance across all labels increased to .75 ( $\sigma = .14$ ), for an average increase in F1 score of +.31. When using crowd-aggregated "Llama-assisted" labels as ground truth, the average weighted F1 score of Llama performance across all labels increased to .72 ( $\sigma = .16$ ), for a similar average increase in F1 score of +30. Performance on some labels increased by substantially more than that average, including "Role of community organizations: Trust, Rapport, & Relationships," which increased from .20 to .75 when GPT-4's labeling performance was evaluated on GPT-assisted ground truth labels.

### 6 Discussion

558

559

560

563

564

567

568

569

571

573

574

575

580

581

585

586

589

592

595

599

607

While mainstream perceptions suggest LLMs can help speed up annotation, we find that if humans review individual LLM-generated suggestions, annotation time does not decrease. From our own measures, we find annotators who were given assistance self-reported improved task understanding, and baseline task understanding, as measured through test questions, remained constant across conditions. Follow-up work could examine if seeing AI suggestions sometimes increased task time because there was more information to process, and if suggestions may teach new annotators to do a complex task, increasing their confidence when assistance is given. Future work can also examine whether annotators under a different compensation incentive structure approve LLM suggestions more quickly than these workers did.

Annotation is usually the first step in creating ground truth data for evaluating a model's performance on an NLP task. Complex, subjective tasks are common in the subfield of NLP for computational social science and cultural analytics. Using LLMs to annotate data used for training and evaluation of a task is attractive due to time and cost efficiency compared to hiring humans, who may noisily interpret a subjective task like this one. However, our findings provide a cautionary note: humans often take LLM suggestions they are given, even when a human individually reviews each label that contributes to the ground truth. As such, using LLM annotations, or LLM-assisted labels like those described here inflates measures of LLM performance on these subjective annotation tasks.

In taking LLM suggestions, annotators homogenize "ground truth" on these tasks towards LLM baselines. In NLP for CSS tasks and in qualitative coding, using labels to measure prevalence of a concept in text is common, so LLM suggestions can change these measurements. For example, imagine using these annotions to analyze cited reasons for vaccine hesitancy given by participants in the NYC conversations. We could view the list of most commonly co-occurring labels with Did not vaccinate. In the human baseline, External motivations: Family & Friends was the second most common co-occurring label, whereas in the LLM baseline, it was the sixth, occurring 25% less often than in the human baseline. In the GPT-assisted annotation condition, External motivations: Family & Friends also drops to the sixth most commonly co-occurring label, mirroring the LLM's label distribution rather than the original human baseline. Analysts using LLM-assisted annotations could thus come to a different conclusion about the relative importance of Family & Friends in motivating vaccine hesitancy compared to the human-only baseline.

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

Homogenization towards an LLM baseline may not be an inherent problem, but the low performance we observe from both GPT-4 and Llama when contrasted with a human baseline suggests that LLMs use a different background concept (Jacobs and Wallach, 2021) for annotating some labels than human crowdworkers do. For example, given a starting F1 score of just .2 comparing LLM annotations to the human crowd baseline, GPT-4 must operationalize the identification Health Education and Support differently than crowdworkers did. When LLM annotations are used to identify this construct, annotations are thus measuring something different than the humans crowd baseline does, and the background concept of Health Education and Support being annotated becomes more like the LLM's conception of Health Education and Support when annotators have its assistance.

LLM assistance also inherently increases measures of interrater reliability when the same LLM assistance is given across annotators. IRR measures are used as a positive signal of reliability when humans annotate social science concepts (Mc-Donald et al., 2019). However, for increased IRR across annotators when using LLM assistance to be a positive outcome of LLM assistance on annotation, researchers would need to be confident that the LLM baseline— and the way it operationalizes each background concept being identified— is *more correct* than human judgments. In subjective tasks, this can be hard to verify and prove, but error

			Ground truth v. LLM, weighted F1 score			
Top-level label	Sublabel	Frequency	Human crowd v. GPT labels	Human crowd v. Llama labels	GPT-assisted human labels vs GPT labels	Llama-assisted human labels vs Llama labels
	Civic Organizations	32	0.60	0.49	0.87 +.27	0.79 +.30
	Employers	8	0.50	0.50	0.84 +.34	0.79 +.29
External motivations	Family & Friends	36	0.52	0.53	0.76 +.24	0.82 +.29
	Health Care Providers	20	0.53	0.49	0.89 +.36	0.75 +.26
	Social & News Media	8	0.32	0.46	0.73 +.41	0.87 +.41
	Discussion of post-pandemic future	16	0.73	0.59	0.79 +.06	0.84 +.25
Future visions & Takeaways	Reflections on the conversation	19	0.15	0.00	0.36 +.21	0.24 +.21
	Basing decisions on data	6	0.42	0.33	0.81 +.39	0.87 +.47
Intrinsic motivations	Getting back to normal	28	0.38	0.44	0.76 +.38	0.65 +.21
	Not wanting to get the virus	28	0.48	0.49	0.92 +.44	0.90 +.41
	Resilience, Connection, & Hope	35	0.32	0.26	0.60 +.28	0.66 +.40
Personal COVID experience	Stress, Fear, & Uncertainty	56	0.66	0.61	0.86 +.20	0.84 +.23
	Significant Impact Resources that helped	34	0.64	0.63	0.84 +.20	0.71+.08
Resources that helped	Unmet community needs	24	0.42	0.46	0.50 + .08	0.67 +.21
	Health Education & Support	23	0.20	0.26	0.68 +.48	0.56 +.30
	Incentives	6	0.28	0.29	0.67 +.39	0.64 +.35
Role of community organizations	Reducing barriers	8	0.35	0.33	0.78 +.43	0.81 +.48
	Trust, Rapport, & Relationships	22	0.20	0.20	0.75 +.55	0.62 +.42
	Did not vaccinate	3	0.60	0.40	0.77 +.17	0.57 +.17
vaccine nesitancy	Mistrust or Skepticism	23	0.57	0.71	0.93 +.36	0.88 +.17

Table 2: Performance evaluation of human crowd labels against GPT labels and Llama labels, then GPT-assisted human labels against GPT labels and Llama-assisted human labels against Llama labels. **Frequency** refers to n observations found across crowd-aggregated annotations on the 200 quotations from the NYC corpus. Performance increases an average of +30% when using LLM-assisted crowd labels compared to an unassisted crowd baseline.

analysis of specific ambiguous labels may help.

661

670

671

673

674

676

678

679

681

686

Second, there are many tasks in NLP where varied annotator perspective can be valuable (Cabitza et al., 2023; Plank, 2022), both for the reason of constructing a robust ground truth (Aroyo, 2013; Yan et al., 2014), or for representing diverse human perspectives on a complex social construct like hate speech (Sap et al., 2022). In qualitative research, some traditions embrace divergent annotator perspective as well in order to widen insight in qualitative annotation (McDonald et al., 2019). In both the cases of creating ground truth for NLP tasks and qualitative annotation, practitioners should know that using or providing LLM assistance to annotators will likely result in lessened variation.

Given potential consequences for representation and construct validity that vary by task, researchers should only proceed with LLM-assisted annotation with a level of caution appropriate to their task and goal. They should recognize that using LLM assistance to construct ground truth labels inflates perceptions of model performance on that task, even when humans review them and individual judgments are aggregated into a crowd ground truth. Follow-up work can investigate if these findings hold for a similarly complex annotation task, a less complex annotation task, and when employing expert annotators rather than inexperienced crowd workers. Homogenization effects may be lessened by presenting information about model confidence, or alternative ideas, and not just a single set of suggestions.

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

## 7 Conclusion

In a pre-registered experiment with 350 unique annotators and 7,000 annotations across 4 conditions, 2 models, and 2 datasets, we find that presenting crowdworkers with LLM-generated annotation suggestions did not make them faster, but did improve their self-reported confidence in the task. Annotators strongly uptook suggestions, changing the label distribution to more closely resemble the LLM's proposed distribution.

Strikingly, we found that using LLM-assisted labels to evaluate model performance resulted in much higher reported F1 scores than when using a human crowd baseline, with increases in F1 scores for model performance on some labels by as much as +x.56. Obviously, using labels influenced by the model to evaluate the model is not standard or advisable in classic evaluation paradigms. However, in the many systems being created that "just put a human in the loop" to review LLM annotation outputs, this paradigm of reviewing LLM outputs to "approve" them is increasingly likely to occur. Practitioners should know that, especially in subjective tasks, simply reviewing LLM suggestions will nudge the distribution of label outputs towards an LLM baseline, even if humans are given a change to review the outputs.

735

736

737

738

739

740

741

742

743

744

745

747

748 749

750

751

752

753

754

756

757

759

760

761

765

## 8 Limitations

Crowdworkers may be particularly susceptible to this kind of influence from LLM suggestions, how-721 ever, their continued employment as a standard 722 for annotation in the field justifies their employment in the task here on a deliberately ambiguous task. Furthermore, different results may be reached with specialized annotators with particular domain knowledge. Annotators with a relationship to the 727 data, including a relationship to the community of interest, may also shape how they align with or 729 reject AI interpretations of the data. Follow-up 730 work can investigate whether domain experts have 731 less anchoring bias than we observed here, and whether they confidently defect from LLM suggestions when needed. 734

### References

- Lora Aroyo. 2013. Crowd Truth: Harnessing disagreement in crowdsourcing a relation extraction gold standard. page 0 Bytes. figshare. Artwork Size: 0 Bytes.
- Shubham Atreja, Joshua Ashkinaze, Lingyao Li, Julia Mendelsohn, and Libby Hemphill. 2024. Prompt Design Matters for Computational Social Science Tasks but in Unpredictable Ways. *arXiv preprint*. ArXiv:2406.11980 [cs].
- Jacques Bughin. 2024. The role of firm ai capabilities in generative ai-pair coding. *Journal of Decision Systems*, 0(0):1–22.
- 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. Number: 6.
- Alexander S. Choi, Syeda Sabrina Akter, JP Singh, and Antonios Anastasopoulos. 2024. The llm effect: Are humans truly using llms, or are they being influenced by them instead? *Preprint*, arXiv:2410.04699.
- Stefano De Paoli. 2023. Can Large Language Models emulate an inductive Thematic Analysis of semistructured interviews? An exploration and provocation on the limits of the approach and the model. *arXiv preprint*. ArXiv:2305.13014 [cs].
- Jessica L. Feuston and Jed R. Brubaker. 2021. Putting Tools in Their Place: The Role of Time and Perspective in Human-AI Collaboration for Qualitative Analysis. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–25.
- Jie Gao, Kenny Tsu Wei Choo, Junming Cao, Roy Ka-Wei Lee, and Simon Perrault. 2023a. CoAlcoder: Examining the Effectiveness of AI-assisted Humanto-Human Collaboration in Qualitative Analysis.

ACM Transactions on Computer-Human Interaction, 31(1):6:1–6:38.

770

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

- Jie Gao, Kenny Tsu Wei Choo, Junming Cao, Roy Ka-Wei Lee, and Simon Perrault. 2023b. Coaicoder: Examining the effectiveness of ai-assisted human-tohuman collaboration in qualitative analysis. ACM Trans. Comput.-Hum. Interact., 31(1).
- Jie Gao, Yuchen Guo, Gionnieve Lim, Tianqin Zhang, Zheng Zhang, Toby Jia-Jun Li, and Simon Tangi Perrault. 2024. CollabCoder: A Lower-barrier, Rigorous Workflow for Inductive Collaborative Qualitative Analysis with Large Language Models. *arXiv preprint*. ArXiv:2304.07366 [cs].
- Jie Gao, Zhiyao Shu, and Shun Yi Yeo. 2025. Using Large Language Model to Support Flexible and Structural Inductive Qualitative Analysis. *arXiv preprint*. ArXiv:2501.00775 [cs].
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. ChatGPT outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120. Publisher: Proceedings of the National Academy of Sciences.
- Zeyu He, Chieh-Yang Huang, Chien-Kuang Cornelia Ding, Shaurya Rohatgi, and Ting-Hao Kenneth Huang. 2024. If in a Crowdsourced Data Annotation Pipeline, a GPT-4. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, pages 1–25, New York, NY, USA. Association for Computing Machinery.
- Abigail Z. Jacobs and Hanna Wallach. 2021. Measurement and Fairness. In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21, pages 375–385, New York, NY, USA. Association for Computing Machinery.
- Maurice Jakesch, Megan French, Xiao Ma, Jeffrey T. Hancock, and Mor Naaman. 2019. Ai-mediated communication: How the perception that profile text was written by ai affects trustworthiness. In *Proceedings* of the 2019 CHI Conference on Human Factors in Computing Systems, CHI '19, page 1–13, New York, NY, USA. Association for Computing Machinery.
- Hang Jiang, Xiajie Zhang, Robert Mahari, Daniel Kessler, Eric Ma, Tal August, Irene Li, Alex Pentland, Yoon Kim, Deb Roy, and Jad Kabbara. 2024. Leveraging large language models for learning complex legal concepts through storytelling. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7194–7219, Bangkok, Thailand. Association for Computational Linguistics.
- Jialun Aaron Jiang, Kandrea Wade, Casey Fiesler, and Jed R. Brubaker. 2021. Supporting Serendipity: Opportunities and Challenges for Human-AI Collaboration in Qualitative Analysis. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–23. ArXiv:2102.03702 [cs].

Michelle S. Lam, Janice Teoh, James A. Landay, Jeffrey Heer, and Michael S. Bernstein. 2024. Concept Induction: Analyzing Unstructured Text with High-Level Concepts Using LLooM. In Proceedings of the CHI Conference on Human Factors in Computing Systems, pages 1–28, Honolulu HI USA. ACM.

827

833

841

842

847

850

851

852

853

854

856

857

867

871

872

873

874 875

876

877 878

879

881

- Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, and Ming Yin. 2023. Synthetic Data Generation with Large Language Models for Text Classification: Potential and Limitations. *arXiv preprint*. ArXiv:2310.07849 [cs].
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *Preprint*, arXiv:2307.03172.
- Nora McDonald, Sarita Schoenebeck, and Andrea Forte. 2019. Reliability and Inter-rater Reliability in Qualitative Research: Norms and Guidelines for CSCW and HCI Practice. *Proc. ACM Hum.-Comput. Interact.*, 3(CSCW):72:1–72:23.
- Lisa Messeri and M. J. Crockett. 2024. Artificial intelligence and illusions of understanding in scientific research. *Nature*, 627(8002):49–58. Publisher: Nature Publishing Group.
- Hussein Mozannar, Gagan Bansal, Adam Fourney, and Eric Horvitz. 2024. When to show a suggestion? integrating human feedback in ai-assisted programming. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 10137–10144.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,

Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, C. J. Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. GPT-4 Technical Report. arXiv preprint. ArXiv:2303.08774 [cs].

886

887

889

890

891

892

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

Cassandra Overney, Belén Saldías, Dimitra Dimitrakopoulou, and Deb Roy. 2024. SenseMate: An Accessible and Beginner-Friendly Human-AI Platform for Qualitative Data Analysis. In *Proceedings* of the 29th International Conference on Intelligent

- 949 950 951
- 952 953
- 953 954
- 95
- 950 957
- 958 959
- 960 961
- 962 963 964

- 969 970
- 971 972 973

974 975

976 977 978

979

980 981

982 983 984

986 987

985

988

- 989
- 991
- 993 994

995

996 997 998

999 1000

1001

1002 1003 *User Interfaces*, pages 922–939, Greenville SC USA. ACM.

- Nicholas Pangakis and Samuel Wolken. 2024. Keeping Humans in the Loop: Human-Centered Automated Annotation with Generative AI. Publisher: arXiv Version Number: 2.
- Jiaxin Pei, Aparna Ananthasubramaniam, Xingyao Wang, Naitian Zhou, Jackson Sargent, Apostolos Dedeloudis, and David Jurgens. 2023. POTATO: The Portable Text Annotation Tool. *arXiv preprint*. ArXiv:2212.08620 [cs].
- Barbara Plank. 2022. The "Problem" of Human Label Variation: On Ground Truth in Data, Modeling and Evaluation. Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 10671–10682. Conference Name: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing Place: Abu Dhabi, United Arab Emirates Publisher: Association for Computational Linguistics.
- 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. *arXiv preprint*. ArXiv:2111.07997 [cs].
- Hope Schroeder, Deb Roy, and Jad Kabbara. 2024. Fora: A corpus and framework for the study of facilitated dialogue. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 13985–14001, Bangkok, Thailand. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint*. ArXiv:2302.13971 [cs].
- Amos Tversky and Daniel Kahneman. 1974. Judgment under uncertainty: Heuristics and biases. *Science*, 185(4157):1124–1131.
- Jenny S. Wang, Samar Haider, Amir Tohidi, Anushkaa Gupta, Yuxuan Zhang, Chris Callison-Burch, David Rothschild, and Duncan J. Watts. 2025. Media Bias Detector: Designing and Implementing a Tool for Real-Time Selection and Framing Bias Analysis in News Coverage. ArXiv:2502.06009 [cs].
- Shuohang Wang, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2021. Want To Reduce Labeling Cost? GPT-3 Can Help. In *Findings of the* Association for Computational Linguistics: EMNLP 2021, pages 4195–4205, Punta Cana, Dominican Republic. Association for Computational Linguistics.
  - Yiwei Wang, Yujun Cai, Muhao Chen, Yuxuan Liang, and Bryan Hooi. 2023. Primacy Effect of ChatGPT.

In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages1004108–115, Singapore. Association for Computational1005Linguistics.1007

Yan Yan, Rómer Rosales, Glenn Fung, Ramanathan Subramanian, and Jennifer Dy. 2014. Learning from multiple annotators with varying expertise. *Machine Learning*, 95(3):291–327.

1009

1010

1011

1012

1013

1014

1015

1016

1049

Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2023. Can Large Language Models Transform Computational Social Science? *arXiv preprint*. ArXiv:2305.03514 [cs].

# A Corpus details

Corpora were selected based on the availabil-1017 ity of both public conversation data within Fora 1018 (Schroeder et al., 2024), and the existence of a 1019 human-created codebook for annotation which was 1020 shared with us by our collaborating organizational 1021 partner. The Fora corpus is lightly anonymized, with speaker names removed. For the NYC and 1023 RTFC sample of 200 quotations, we manually re-1024 viewed and removed any personally identifiable 1025 information before using it in a prompt to either model, and before showing it to annotators. Partic-1027 ipants in the NYC and RTFC conversation collec-1028 tions were aware their voices would be collected 1029 and used to inform the public about issues in their 1030 community, as well as potentially used in research. 1031 Both the NYC and RTFC conversations contain 1032 a mix of standard American English and African 1033 American English, as well as Spanish, though less often. We eliminated any participant comments 1035 in Spanish prior to sampling 200 comments for 1036 this study, given that performance on this tasks 1037 could not be compared across languages, and the codebook was developed in English. Given the 1039 unreliability of African American English dialect 1040 detectors as of the time of this submission, partic-1041 ularly on transcribed speech, we are not able to 1042 estimate prevalence of African American English 1043 or other dialects that may be in this corpus. Demo-1044 graphic information of speakers was not collected 1045 in the NYC or RTFC conversation corpora, but more details on the corpora are available in the Fora corpus. 1048

## **B** Additional results

In addition to Table1, we provide a visualization of<br/>time taken on Conditions 1-4 on the NYC corpus1050<br/>1051with LLM suggestions from GPT-4 in Table3.1052



Figure 2: Time (number of seconds) spent by annotators on the task, for Conditions 1-4 on NYC data.



Figure 3: Annotator perception of suggestion helpfulness by condition

As shown in Figure3, annotators perceived GPT-4 suggestions in Conditions 2, 3, and 4 on the NYC corpus to be between somewhat and very helpful.

1053

1054

Corpus	Assistance type	Assistance Model	Average Helpfulness	Std
NYC	Text-based	GPT-4	1.42	.42
NYC	Text-based	Llama	1.3	.38
NYC	Text-based, AI disclosed	GPT-4	1.36	.42
NYC	Pre-filled	GPT-4	1.45	.62
RTFC	Text-based	GPT-4	1.49	.37

Table 3: Helpfulness of suggestions across conditions

Text-based suggestions from GPT-4 were rated as similarly helpful on RTFC data.

Self-reported level of Task Understanding across conditions

Figure 4: Self-reported task understanding by condition



Figure 5: Self-reported task confidence by condition. Conditions 1-4 on NYC data shown.

These figures visualize conditions 1-4 on the NYC data. A full table of results on these measures is	1057
available in Table 4.	1058

# **B.1** Additional results

1056

Corpus	Assistance type	Assistance Model	Task understanding	Task confidence	Understanding of the community	Ability to explain community needs
NYC	No assistance	None	3.33	3.3	3.34	3.31
	Text-based	GPT-4	3.79 **	3.64 **	3.71 **	3.72 **
	Text-based	Llama	3.83	3.80	3.85 **	3.86 **
	Text-based, AI disclosed	GPT-4	3.84 **	3.70 **	3.75 *	3.73 **
	Pre-filled	GPT-4	3.5	3.5 **	3.51	3.51 **
RTFC	No assistance	None	4.0	3.89	3.87	3.84
	Text-based	GPT-4	4.0	3.96	4.04 **	4.04 **

Table 4: Levels of annotator understanding and confidence by condition. Double stars indicate statistically significant increases over baseline conditions for the corpus (as calculated with a two-sided T-test with Bonferroni correction), and single stars indicate statistically significant increases over baseline conditions that reduce to insignificance with Bonferroni correction.

			GPT g	ground truth,	Llama test
Top-level label	Sublabel	Prevalence (n)	F1	Precision	Recall
	Civic Organizations	25	0.60	0.51	0.84
	Employers	24	0.50	0.81	0.54
External motivations	Family & Friends	24	0.52	0.63	0.79
	Health Care Providers	21	0.53	0.55	0.76
	Social & News Media	14	0.32	0.47	0.57
E-ti-i 8 T-l	Discussion of post-pandemic future	25	0.73	0.88	0.60
Future visions & Takeaways:	Reflections on the conversation	4	0.15	0.43	0.75
	Basing decisions on data	12	0.42	0.50	0.75
Intrinsic motivations:	Getting back to normal	13	0.38	0.64	0.69
	Not wanting to get the virus	58	0.48	0.78	0.69
B	Resilience, Connection, & Hope	20	0.32	0.47	0.35
rersonal COVID experience:	Stress, Fear, & Uncertainty	55	0.66	0.65	0.80
Becommons that helped	Significant Impact Resources that helped	33	0.64	0.53	0.82
Resources that helped	Unmet community needs	11	0.42	0.56	0.45
	Health Education & Support	24	0.20	0.54	0.63
Dala di secondari dei secondari di secondari	Incentives	15	0.28	0.48	0.73
Kole of community organizations	Reducing barriers	22	0.35	0.41	0.59
	Trust, Rapport, & Relationships	18	0.20	0.63	0.67
X7	Did not vaccinate	12	0.60	0.86	0.50
Vaccine hesitancy	Mistrust or Skepticism	28	0.57	0.88	0.82

Table 5: This table compares GPT-4 and Llama annotations on the NYC corpus to each other, by treating GPT-4 labels as ground truth and scoring Llama annotations against GPT-4. The average F1 score for Llama against the GPT-4 baseline is .62 ( $\sigma = .10$ ), which much higher than either model's comparison to human baseline, discussed in the main body of the paper.



Figure 6: Self-reported understanding of community needs by condition. Conditions 1-4 on NYC data shown.



Figure 7: Self-reported ability to explain community needs by condition. Conditions 1-4 on NYC data shown.

### C Interface screenshots

So she said that, right? Yeah. So now my job as my son's because I know my son's first teacher is to be able to pull apart these things that I know he's willing to go through, see them firsthand with having him be considered as aggressive or have an attitude problem or anything. You just knows because he's already witnessed it firsthand. Like, "I know you're going to get treated differently than me if I do what you just did. I know that for a fact. I've seen it." And he's 10 and I've raised him to always voice how he feels. And so that might come off as smart remark or disrespectful, but no big deal at that point. It may come on to you. So it's just, how do you navigate that? Plus knowing the system, plus raising a child and keeping his morals and his self-esteem and confidence and tech without letting their skins play that.

Suggested tags: Inequality::Race, Community Life::Community Values, Community Life::Community Relationships



Figure 8: *Condition 2*: Text-based label suggestions in interface, with example and label suggestions from RTFC corpus. Annotators could scroll down for more annotations. Interactive tooltips gave extended code definitions, and the full codebook was linked in the header of the annotation task.



Figure 9: *Condition 4:* Pre-highlighted label suggestions in interface, with example and label suggestions from NYC corpus. Interactive tooltip with code definition shown in upper right.

1061	D NYC Conversation Corpus Details
1062	D.1 Conversation guide
1063	The NYC Public Health Corps developed a codebook in partnership with Cortico to elevate concerns
1064	from community members during the pandemic period. The main questions asked of participants in each
1065	conversation were:
1066	Opportunities & Challenges of Resourcing
1067	- What COVID-19 resources have been most helpful for you during the pandemic? And, why?
1068	- What challenges did you have finding and using resources intended to help you with COVID-195
1069	Vaccine Experiences and Decision Making
1070	- Can you describe a key moment that influenced your decision about the COVID-19 vaccine?
1071	- Reflecting on this experience you just shared, what information or circumstances helped you
1072	make that decision?
1073	Role of Community Organizations in Community
1074	- Can you share a story or experience about how the community organizations in your neighbor
1075	hood supported you and your community during the pandemic?
1076	- How did those community organizations impact decisions related to COVID-19 vaccinations in
1077	your neighborhood?
1078	• Future Resources
1079	- What will a post-pandemic future look like in your community?
1080	- How can community organizations help your community thrive in the future?

D.2 NYC Codebook	1081
1. External motivations (Theme)	1082
• Civic Organizations: The speaker mentions civic organizations like non-profits, churches,	1083
NGOs, and community clubs, associations as a factor in health decisions	1084
- Example: "My church helped me find alternative childcare during the pandemic."	1085
• Family & Friends: Mentions family and friends as a factor in health decisions	1086
- Example: "My mom was really opinionated about this from the start. She really wanted us	1087
to get the vaccine. She even helped drive us to the vaccine clinic because I don't own a car."	1088
• Employers Mentions place of employment as a factor in health decisions, including employer	1089
providing resources or opportunities for vaccination, or co-workers setting a model of vaccine	1090
behavior as a factor in health decisions	1091
- Example: "My job offered drop-in vaccine clinics, which was helpful since the regular clinic	1092
hours happen during my work hours."	1093
• Social & News Media: Mentions social media or news media as a factor in health decisions,	1094
that is online	1095
- Example: "A lot of my friends were posting on Eacebook about the vaccine having a chip	1090
That made me nervous."	1097
• Health Care Providers: Mentions healthcare providers as a factor in health decisions, including	1099
doctors or nurses	1100
- Example: "My doctor gave me some information, but I don't trust that the information was	1101
up to date and I still don't want the vaccine. So no I haven't gotten it."	1102
2. Intrinsic Motivations (Theme)	1103
• Fear of Virus: Mentions not wanting to become ill from COVID as a factor in health decisions	1104
- Example: "My church helped me find alternative childcare during the pandemic."	1105
• Getting Back to Normal: Mentions a desire for a return to activities and social routines as a factor in health decisions	1106 1107
- Example: "I thought you know what, if this can help us just keep the kids in school then I'll	1108
get the vaccine even if I hate it."	1109
• Basing decisions on data: Mentions considering data, research, or evidence when choosing to	1110
get vaccinated	1111
- Example: "I was seeing these studies show that there is an increased chance of heart	1112
problems after the vaccine. So I did not want to get it because I have heart issues in my	1113
family."	1114
3. Role of Community Organizations (Theme)	1115
• Health Education & Support: Mentions that community health educators were a factor in	1116
health decisions	1117
- Example: "A community health person came to my school and explained what was going	1118
on in the pandemic and the latest research on masking. So that's when we started masking."	1119
• Support Incentives: Mentions that a community-based organization provided an incentive or	1120
resource in some material form to support during the pandemic. For example, CBO (community-	1121
based organization) providing masks, hand sanifizer, gift cards, vaccination opportunities to	1122
Community memories.	1123
- Example: My local lood pantry gave out masks which was super helpful when they were sold out online "	1124
• Trust rannort & relationships• Mentions trust in community organization due to outreach	1120
sharing personal stories, having open conversations, positive relationship-building	1120

1128 1129	<ul> <li>Example: "It just made us feel like there was somewhere to turn to when everything was chaotic. The church gave us a place to talk about these things and feel safe."</li> </ul>
1130	• <b>Reducing barriers:</b> Mentions a barrier, struggle, or challenge to implementing their health
1131	decisions or asserting agency in health decisions
1132 1133	- Example: "I just did not feel listened to by my doctor. And there was no alternative to what they were telling me. So"
1134	4. <b>Resources</b> (Theme)
1135 1136	• <b>Significant Impact Resources:</b> Mentions that impact resources like food, financial assistance, rent moratorium, student loan suspension, employment helped during COVID
1137	<ul> <li>Example: "We would have been so lost without the food bank that restocked each week."</li> </ul>
1138	• Unmet community needs: Mentions that a resource is needed
1139 1140	<ul> <li>Example: "We just never found the childcare we needed so that I could keep my job. I haven't worked since 2020."</li> </ul>
1141	5. Vaccine Hesitancy (Theme)
1142	• Did not vaccinate: Mentions that the speaker did not choose to vaccinate
1143	- Example: "I just couldn't get over how scary it was that my sister had this reaction to the vaccine. I know it could happen to me. So no I did not go through with the vaccine."
1145	• Mistrust or Skenticism: Mentions that the speaker has/had mistrust or skenticism of the
1146	vaccine
1147	- Example: "There are people telling me this vaccine has a chip in it. I don't want a chip and
1148	I just have no way of knowing."
1149	6. Personal COVID-19 Experience (Theme)
1150	• Resilience, Connection, & Hope: Mentions agency, control, or feeling empowered during the
1151	pandemic period
1152	- Example: "Helping at my church made me feel like I was making a difference even though
1153	the world was going crazy."
1154	• Stress, Fear, & Uncertainty: Mentions stress, fear, or uncertainty during the pandemic period
1155 1156	- Example: "It was just anxiety all day every day thinking about my kid getting sick at school and bringing it back to her brother at home."
1157	7. Future Visions & Takeaways
1158	Conversation Reflections: Mentions reflections on the conversation
1159	- Example: "Talking about this has made me remember how hard that period was for our
1160	family."
1161	• <b>Post-pandemic future:</b> Mentions a future vision for their life or community after the pandemic
1162	- Example: "I just can't wait for the schools to go back to normal and I hope we can all learn
1163	something from this."
1164	E RTFC Conversation Corpus Details
1165	E.1 RTFC Conversation guide
1166	Sharing Questions & Lived Experiences
1167	- As we shed the restrictions of the pandemic, it will be easy for us to lose sight of what we have
1168	learned about inequality in America and Boston and how our lived experiences have shaped
1169	that learning. Keeping this in mind, I'd like to invite you to think: "What's your question about the future of Posten and your place in that future?"
1170	the future of boston and your place in that future?

- Thank you for sharing your questions with us. Now I'd like to invite you to think, what experience in your life got you to this question?	1171 1172
Connecting Our Experiences	1173
- Find someone whose question or experience resonates with your own life. Then I want you	1174
to speak to that person and tell them why their question or experience resonated with you and	1175
share the story from your life that connects you with their experience.	1176
Drawing Connections	1177
- Let's talk a little about what we are hearing. What are you hearing in people's experiences?	1178
• Wrap up	1179
- Do you have any closing thoughts that you'd like to share or other general reflections? Do you	1180
have any questions for us?	1181
E.2 RTFC Codebook	1182
1. Government and Institutions (Theme)	1183
• Expectations Deferences to the expectations and conjustions that the public has of elected	1104
• Expectations References to the expectations and aspirations that the public has of elected officials, city government, and/or civic institutions.	1184 1185
• Processes: References to processes through which the public interfaces with government, such	1186
as voting, community engagement, campaigning, electoral processes, and other decision-making	1187
processes. This may include feelings of exclusion, silencing, or neglect in public meetings;	1188
community dynamics within a public meeting; curiosity about electoral results; a lack of	1189
confidence in voting as a form of democratic participation.	1190
• Accountability: Statements about the accountability of elected officials, city government,	1191
and/or civic institutions to the promises they make and the expectations they set for the public.	1192
This may include references to elected officials who "will tell you anything just to get your vote";	1193
about how much the city listens to its residents and factors resident perspectives into decisions	1194
E.g. "You said you were going to do this, and you have/but you haven't yet"	1196
• Institutional Resources: Statements about how people are having difficulty (or success)	1107
accessing services provided by government agencies and other institutions that improve one's	1198
quality of life. This could include mentions of services and resources like such as housing	1199
subsidies, senior services, mental health services, or municipal services like fixing potholes.	1200
This could also include statements about the difficulties people face in accessing these services	1201
or navigating institutions to get the services and resources they need.	1202
• Community Resources: Statements about how people are leveraging resources in their	1203
communities to fulfill their needs and improve one's quality of life. This could include mentions	1204
of community-based organizations that fulfill community needs; civic associations; or neighbors	1205
that provide support to other neighbors.	1206
2. Public Health (Theme)	1207
• Mental Health: Those who struggle with mental health; systemic issues of mental health;	1208
responses to those with mental health issues; resources and isntitutions that support mental	1209
health	1210
• Drugs and Drug Use Disorder: Addiction, systemic issues of drug use, responses to those	1211
with drug use disorders, the culture and environment around drug use	1212
• Trauma: Individual, community, generational traumas. The responses and resources intended	1213
to support healing from those traumas. Things that further cause traumas.	1214
21	

1215	• Quality and Affordable Healthcare: The accessibility, affordability, and quality of healthcare and other health services
1210	• Each Inconvity: Each access quality of food accessible, food deserts, offerdebility of food
1217	• Food Insecurity: Food access, quality of food accessible, food deserts, affordability of food, systems to support food accessibility.
1210	• COVID 10: COVID 10 yearing, mask COVID tests boosters, and the impacts of COVID 10.
1219 1220	such as working from home, school closures, and jobs lost.
1221	3. Safety (Theme)
1222	• Sense of Safety: Refers to feeling unsafe within daily life routines at home, in one's neighbor-
1223	hood, and throughout the city.
1224	• Street Violence: Refers to situations like street fighting, assaults on the street, unintentional
1225	harm of bystanders, etc.
1226	• Gun Violence: Loss of family members due to a shooting, witnessing a shooting AND not
1227	limited to gang violence.
1228	• Policing: Refers to being targeted by police (profiled) in certain areas and the lack of policing
1229	happening due to neighborhood location, race and/or ethnicity.
1230	• Racialized Violence: Refers to verbal, emotional and physical assaults based on color of skin,
1231	race, ethnicity, language.
1232	4. Infrastructure (Theme)
1233	• Climate Impacts: Climate change, impact of climate change on the community, actions to
1234	address climate change, fears around climate change.
1235	• Transportation: Public transportation like the buses and trains, quality of transportation,
1236	affordability and accessibility of transportation, safety of public transit.
1237	5. Housing (Theme)
1238	• Gentrification and displacement: Displacement of lower income residents; physical transfor-
1239	mation and change of the cultural character of the neighborhood.
1240	• Housing Instability: Difficulty paying rent, having frequent moves, living in overcrowded
1241	conditions, or doubling up with friends and relatives.
1242	• Homeownership: Challenges for owning a house; obstacles toward home ownership; express-
1243	ing the with hope to be a home owner.
1244	• Housing quality: the physical condition of a person's home as well as the quality of the social
1245	and physical environment in which the home is located
1246	• Housing affordability: Cost of housing and how affordable that cost is to residents, regardless
1247	of tenure (tenant/owner), subsidy (e.g. workforce housing, public housing)
1248	6. Community Life (Theme)
1249	• Community Relationships: Relationships between community members, across generations,
1250	and across communities. Quality and nature of those relationships.
1251	• <b>Community Values:</b> Values instilled throughout the community, values differences within and
1252	across communities.
1253	7. Education (Theme)
1254	• <b>Ouality of Education:</b> Education that leads to empowerment as a process of strengthening
1255	individuals and communities to get more control over their own situations and environments:
1256	education systems that focus on the importance of quality learners, quality learning environment.
1257	quality content, quality processes, and quality outcomes
1258	• School Infrastructure: Suitable spaces to learn; also spaces that have the infrastructure to
1259	address the COVID-19 public health emergency.

• Life Skills: The abilities (or the lack of) for adaptive and positive behaviour that enable individ-	1260
uals to deal effectively with the demands and challenges of everyday life in their communities	1261
• Vouth Spaces: Available and accessible physical and virtual spaces for activities especially	1202
offered to young people to advance their cognitive, emotional, social, and creative skills	1263
• Higher Education: Post-secondary academic institutions, including colleges/universities/voca-	1265
tional schools, where individuals engage in advanced learning and research. Could be used to define relationships between students, teachers, administration.	1266 1267
8. Economic Opportunity (Theme)	1268
• Jobs: References to a person's ability to provide for themselves and their families. Can include	1269
statements about working multiple jobs; working in a particular industry; facing unemployment;	1270
job satisfaction; difficulties in finding a job; observations about the job market; discrimination	1271
within a job or during a job search; efforts to attain more training or education in order to	1272
improve one's job prospects	1273
• Economic Assistance: References to one's ability to access economic supports that enable	1274
wealth-building, financial stability, and/or economic growth. This can include statements	1275
about individuals (such as one's ability to access home loans) and small businesses (such as a	1276
business's ability to access lines of credit).	1277
• Income: Explicit references to income/wages and wealth. This can include discussions about:	1278
one's personal income; satisfaction with their income; in/ability to increase their income;	1279
in/ability to build wealth; income inequality; the income/wage levels to be able to afford the cost of living in Boston	1280
• Affordable Childeare: Deferences to one's shility to offerd shildcare. This is included in	1000
Economic Opportunity because childcare affects one's ability to maintain stable employment.	1282
• Financial Literacy: References to people's level of financial literacy, from everyday money	1284
management, to processes for applying loans and credit. This can also refer to people's general	1285
lack of financial literacy.	1286
9. Inequality (Theme)	1287
• Race: RDefined as lack of jobs, services, goods, based on skin color, ethnicity, language.	1288
• Class: Refers to socio economic status, education, and types of disparities, including neighbors	1289
re-entering society.	1290
• Gender: Discrimination based on (anatomy) female, male.	1291
• Sexual Orientation: Refers to sexual identity and preference.	1292
• Ability: Refers to disabilities, physical and intellectual.	1293
• Immigration Status: Foreign born, regardless of documentation - this example speaks more	1294
to being an immigrant in which English is the second language, which is the barrier of an	1295
immigrant.	1296

### E.3 Prompts to GPT-4 and Llama

1297

1298

1299

1300

1301

1302 1303

1304

1305

1306

1307

1308

1310

1311 1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327 1328

1329

1330

1331

1332

1333

1334

1335

1336

1339

1340 1341

1342

1343

1344

1345

1346

1347

1348

1350

1351

1352

1353 1354

1355

1356

1357

1358

We provided an instruction, list of labels and definitions in the codebook in JSON format for each corpus' codebook. We requested output in JSON format as follows:

```
Your job is to provide a comprehensive set of
    thematic labels for the given quote.
You are given 7 [9] thematic tagging questions,
    subthemes, and general descriptions\\ in the
     ISON below
Choose ONLY from the tags provided here.
"annotation_schemes": [<list of top-level labels
     and sublabels, in JSON format, with
    definitions>]
For the given conversation quote, return all
    subtags that apply in a single JSON array in
     this format, (or return ''None of the above
    :: I confirm that none of the themes apply"
    if none apply or if the statement is too
    ambiguous to determine):
Expected format:
  ʻjson
Ε
  {{
    "highlight_id": highlight_id,
    "tags": ["External Motivations::Employers",
    "Intrinsic Motivations::Not wanting to get
    the virus"]
 }}
]
...
Please share this output format with no any
    additional characters, annotations, line
    breaks, or comments.
highlight_id: [id for excerpt]
Conversation quote: [quote]
JSON:
```

Full prompts and code will be included in published Github repo upon publication. The models were prompted at a temperature of 0. Overall, responses were very well-formed according to this prompt, both for Llama and GPT-4. One single label that was not in our codebook was hallucinated 3 times across the generation process for the NYC corpus: External Motivations: Government. Qualitatively, we note with interest that this theme did come up often in the conversation quotes, and therefore could be seen as a thematically relevant code despite the violation of instructions needed to produce it. Two labels were hallucinated in the RTFC codebook: "Community Life: Community Resources" and "Housing: Housing Stability." Both of these were also plausible given the conversation content, and were hallucinated just once. All hallucinated labels were removed to ensure fidelity to the original codebook. In one instance, GPT generated

"None of the above" in addition to another valid1359label. For this case, we removed the nonsensical1360"None of the above" label from the suggested labels.1361In one case, GPT-4 failed to produce a label of the1362requested format for the RTFC data. We converted1363this suggested label into "None of the above" for1364consistency with the rest of the corpus.1365

1366

1367

1368

1369

1370

1371

1372

1374

1375

1376

1377

1378

1379

1380

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

### F Crowdworker recruitment

Crowdworkers were recruited from Prolific, and had the following characteristics: located in the US or UK, with English as a first language, with 95-100 percent approval ratings, who had attended some or more college to participate in our annotation experiment. We paid the recommended rate on Prolific of \$12/hour, 1.6x the US minimum, and adjusted upward to \$12/hour if our initial estimated time was not sufficient to pay workers this amount.

### F.1 Annotator "understanding" questions

From the round of IRR validation done before the main round of experiments, we found several examples of very high agreement quote-label pairs. We selected four of them to act as test questions for annotators. Two examples of these high-confidence test questions for the NYC: corpus are listed here.

- High agreement example of *Family & Friends* label: "Right, but you know what? I encouraged my sons and my daughter to vaccinate their kids because you don't know that COVID is new. You don't know how it's going to affect them and the children, my grandkids."
- High agreement example of *Support Incentives:* "I think when I made my decision it came in handy, because I think it was a month before my daughter was starting school and they were giving you \$100 for the vaccine. And honestly, that came at a good time. Why? Because I said, 'Okay, they gave me \$100. They gave her \$100.' And I said, 'Oh, this is good because now I could get you this and that before school starts.' So that was pretty good. I mean, that was nice."

### F.2 Instructions to annotators

Following a consent page, annotators received the following instructions:

"First, we need to explain the task we are asking1404you to complete. We will ask you to read a quota-<br/>tion from a conversation. The conversation is about1405

1407<resources and challenges during the COVID-19</th>1408pandemic in the United States. Conversations were1409hosted to better understand resources that helped1410during the pandemic, challenges to access, and1411motivations for making health-related decisions1412during the pandemic>.

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430 1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

We are asking you to identify <7> phenomena in this annotation. Please read about each type before moving forward by going here: [link to Codebook Training Document]. We will check that you have spent at least 3 minutes reading this document before advancing to the task.

There are 15-22 annotations in the task. Review your answer for each question before proceeding."

Then, annotators were given the following choices: "I agree to read carefully and spend enough time on each annotation" or "I do not wish to partake". At the bottom, text read "At the top of the next page, there will be a quote. To annotate, select the check boxes that apply." In the LLM assistance conditions, this text read: "At the top of the next page, there will be a quote, followed by a list of suggested labels for the quote. Please read the quote and the list suggested labels, which you may use to assist your annotation. To annotate, select the check boxes that apply." If they had selected the option to proceed with the study, stimuli to annotate were next presented. After 20 annotations + 2 test questions were presented, participants were given the option to provide demographic information for research purposes.

The crowdworker study was reviewed by our IRB review board, and determined exempt.

### G AI Assistance

1441In this work, the authors used some AI tools to1442find related works, including Elicit and ScholarQA.1443Citations were followed and checked. We also used1444Github Copilot and Cursor for coding assistance,1445and code was reviewed for errors before being run1446as well as after being run for any potential errors.