

PLURULE: A Challenging Benchmark for Detecting Rule Violations of Pluralistic Communities on Social Media

Anonymous ACL submission

Abstract

Social media are shifting towards pluralism — community-governed platforms where groups define their own norms. What violates rules in one community may be perfectly acceptable in another. Can AI models help detect rule violations of such pluralistic communities? We formalize the task as a multiple-choice problem, mirroring how human moderators operate in the real world: given a comment and its surrounding context, identify which specific rule, if any, is violated. We introduce PLURULE, a multimodal, multilingual benchmark for detecting 17,313 rule violations across 2,419 Reddit communities spanning 3,692 pluralistic rules in 10 languages. Using this benchmark, we show that state-of-the-art vision-language models struggle significantly: even GPT-5.2 with high reasoning performs only slightly better than a trivial baseline. We also find that bigger models and increased context provide marginal gains, and universal rules like civility and self-promotion are easier to detect. Our results show that pluralistic moderation of social media is a fundamental challenge for language models.

1 Introduction

Ensuring platform safety and encouraging constructive participation are among the most persistent challenges of social media governance (Gillespie, 2020). On centralized platforms such as X (Transparency, 2025), YouTube (Google, 2025), and Meta (Meta, 2024), content moderation is increasingly carried out by a combination of human moderators and automated detection algorithms. These systems typically focus on narrowly defined categories — such as unlawful content, incivility, hate speech, and harassment — that are assumed to be universal across all users and communities. By privileging mainstream norms, they overlook the values, languages, and forms of expression used by minority communities, leading to higher rates

Input Prompt

Subreddit Info:
Subreddit: r/santamonica: Welcome to Santa Monica, California
Anything related to Santa Monica, where the rent is high but the temperatures aren't.

Rules:
Rule 1: Santa Monica-related only
Please keep posts to Santa Monica-specific issues. Venice and West LA issues can be posted in /r/LosAngeles
... (other rules) ...
Rule 6: No calling people paid shills just because they disagree with you.
Accusing people of being paid shills with no proof just because they disagree with you can and will result in bans.
Rule 7: No hateful speech directed at other people, including about the homeless.
If you want to say something disgusting about homeless people, NextDoor is always there for you.

Submission:
Homeowners association sues City and Fairmont Miramar over hotel redevelopment approval
USER1, Tue, Dec 15, 2020, 1:37PM
[URL]


Discussion:
Comment 1: User 1, Tue, Dec 15, 2020, 1:38PM
And people wonder why everything is so expensive here...
... (other comments) ...
Comment 4 [TARGET COMMENT]: User 4, Mon, Dec 21, 2020, 3:58PM
You must work for Dell. He's been ripping off Santa Monica for years by avoiding taxes on the Fairmont Miramar. Look it up. He's a right wing a-hole

Question: Does the [TARGET COMMENT] violate a rule?
(a) No hateful speech directed at other people, including about the homeless; (b) Santa Monica-related only; (c) No rules broken; (d) No calling people paid shills just because they disagree with you; (e) No spam; (f) Appealing bans; (g) Respect other redditors; (h) No low-effort posts

Model Response
The [TARGET COMMENT] violates:
(a) No hateful speech directed at other people, including about the homeless. ✗

Figure 1: A PLURULE example. The model receives the full context of the target comment, then selects which rule is violated. Here, the correct answer is (d) but the model selects (a).

of content removal for marginalized groups (Lingel and Golub, 2015; Jiang et al., 2020; Griffin, 2024; Celeste et al., 2023). Centralized platform rules therefore fail to account for the diverse experiences and contextual meanings that vary across communities (Díaz and Hecht-Felella, 2021).

Bucking this trend, some platforms have adopted community-governed structures that allow groups to define their own norms. Reddit, for instance, hosts hundreds of thousands of topic-based communities (subreddits), each with its own rule set in addition to platform-wide guidelines (Reddit, 2025). While these pluralistic structures empower communities, they also place a substantial burden on volunteer moderators. On Reddit alone, the estimated value of this uncompensated labor exceeded \$3.4 million in 2020 (Li et al., 2022b). Unsurprisingly, moderators are often eager to adopt automated tools that can reduce their burden (Robinson, 2025; Dosono and Semaan, 2019; Hill, 2019; Lloyd et al., 2025).

However, the contextual nature of community-specific rules poses a fundamental challenge for automation. What violates a rule in one community may be perfectly acceptable in another (Chandrasekharan et al., 2019; Li et al., 2022a). A satirical insult about someone’s appearance, for instance, is encouraged in r/RoastMe but would violate civility rules in most other communities. Similarly, self-promotion that constitutes spam in most subreddits is required in creative showcase communities. Effective moderation requires understanding not just the rule text, but the implicit norms, values, and purposes that each community has developed over time.

Given these contextual complexities, the question arises whether modern AI systems can effectively assist with pluralistic moderation. The central challenge is whether language models can recognize that identical content may be acceptable in one community but violate rules in another. Even similar rules may be interpreted differently depending on local community norms (Selbst et al., 2019; Birhane et al., 2021).

To investigate this question empirically, we formalize the detection of rule violations as a multiple-choice task that mirrors how human moderators operate in practice (Figure 1). We introduce PLURULE, the first multimodal, multilingual benchmark for detecting violations of pluralistic community rules. The benchmark comprises 17,313 moderation instances with 92,305 comments and 4,843

images, spanning 2,419 subreddits with 3,692 distinct rules across 10 languages. PLURULE incorporates substantial diversity along two dimensions: 27 semantically-derived subreddit categories (e.g., politics, gaming, music) and 31 rule categories (e.g., civility, self-promotion, spoilers).

Using PLURULE, we evaluate state-of-the-art vision-language models (VLMs) on the detection of pluralistic rule violations under different context conditions. Our results reveal substantial limitations: even GPT-5.2 with high reasoning effort achieves only 58% accuracy, barely exceeding a trivial baseline that always predicts no violation (50%). Providing additional context — the discussion thread, original submission, participant labels, and images — improves GPT-5.2’s performance by only 2–3 percentage points. Open-weight models like Qwen3-VL-Instruct and Qwen3-VL-Thinking perform even worse, failing to surpass baseline performance. Performance breakdown by rule category reveals that models successfully detect universal violations such as civility (69%) and self-promotion (67%), but fail on rules that require contextual understanding; sourcing requirements (30%), content neutrality (40%), and topic relevance (43%) all fall well below baseline. These results reveal a critical gap: current VLMs can enforce universal norms but cannot adapt to the diverse, context-dependent standards that define pluralistic moderation.

2 Related Work

Existing datasets for content moderation focus on narrow categories such as toxic speech (Hoang et al., 2024), hate speech (Nghiem and Daumé Iii, 2024), or misogyny (Sheppard et al., 2024). Automated systems trained on these datasets are limited to detecting broadly unacceptable content under singular global standards of appropriateness. This assumption breaks down in decentralized platforms, where different demographic groups significantly diverge about what is considered respectful, emotionally appropriate, or toxic (Sachdeva et al., 2022; Ali et al., 2025). Moderation on such platforms is inherently community-dependent and must account for pluralistic, community-specific rules. On Reddit, for example, rules extend beyond toxicity (Binns et al., 2017; Matias, 2019) to include locally defined norms around formatting, tone, and ideological or topical relevance (Chandrasekharan and Gilbert, 2019).

Even the enforcement of similar norms can vary widely across communities (Chandrasekharan et al., 2018). On Reddit in particular, moderators routinely interpret rules and assess the appropriateness of content relative to local community values rather than mechanically executing fixed policies (Li et al., 2022a; Fiesler et al., 2018; Matias, 2019). Consequently, within a community, multiple moderators can diverge when guidelines are broad or context-dependent (Binns et al., 2017; Chandrasekharan et al., 2019). Across communities, the same content may be acceptable in one context while violating norms in another — a distinction that models trained on aggregated data from multiple communities often fail to capture (Sap et al., 2022; Raji et al., 2020).

Previous work attempts to model the community-dependent nuance of moderation, but does not address its context-dependent nature. Chandrasekharan and Gilbert (2019) identify a small set of recurring “macro” norms shared across communities. Park et al. (2021) introduce a text-only dataset that collapses thousands of community-specific rules into coarse-grained types. This approach abstracts thousands of individual subreddit rules into a limited number of universal categories, obscuring differences that define each community. (He et al., 2024) provide models with individual rules for binary yes/no judgments.

PLURULE advances beyond prior work along three key dimensions. First, it explicitly models pluralism: instead of applying a fixed set of universal categories, models must reason over distinct, community-defined rules. Second, it frames moderation as a rule identification task (multiple-choice) rather than binary classification. This mirrors real-world moderator workflows and enables more fine-grained evaluation. Finally, PLURULE is multilingual and multimodal, capturing the visual (Gomez et al., 2020) and linguistic (Blodgett et al., 2016) diversity of online communities often overlooked by text-only benchmarks.

3 PLURULE Benchmark

PLURULE formalizes the task of detecting rule violations of pluralistic communities on Reddit as a multiple-choice question (Figure 1). Given a comment from a specific community (subreddit), models must identify which specific rule, if any, has been violated.

For each comment, models receive the commu-

nity’s rules along with the surrounding context that moderators consider when making decisions. The context includes: (1) the discussion thread that precedes the comment; (2) the submission post to which the comment responds, including any images; and (3) anonymized identifiers of the participants in the discussion.

Each instance in PLURULE consists of a pair: a violating comment and a compliant comment with overlapping context from the same submission. Models are evaluated on both comments separately. For each comment, models are presented with answer options consisting of all subreddit rules plus a “No rules broken” option, labeled (a), (b), (c), etc. Each comment’s answer options are deterministically shuffled using a seed based on the comment ID to prevent models from exploiting positional bias. The correct answer for violating comments is the violated rule; for compliant comments, it is “No rules broken.” Since half the comments violate a rule and half do not, always predicting “No rules broken” yields a majority baseline of 50% accuracy.

4 PLURULE Construction

We select Reddit as a platform because moderation actions are public: a moderation action occurs when a human moderator leaves a comment explaining a rule violation (e.g., “Your comment violates Rule 2”). We construct PLURULE by starting from such moderator comments in the Pushshift Reddit archives (Baumgartner et al., 2020) and transforming them into structured benchmark instances with verified rule labels, contrastive pairs, and semantic clustering. Below we describe the five-phase pipeline for this construction process.

4.1 Phase 1: Data Collection

We start from a publicly hosted version of the Pushshift Reddit archives (Cohen and Lo, 2014), spanning from December 2005 to February 2023 and containing approximately 15 billion comments across 40 thousand subreddits. From these archives, we extract comments by moderators, flagged by a “distinguished” field in the comment object. To focus only on comment (not post) violations, we exclude top-level replies to submission posts. We filter out accounts with usernames that match bot-related keywords, e.g., “bot,” “automod.” This yields approximately 10 million moderator comments across 40 thousand subreddits.

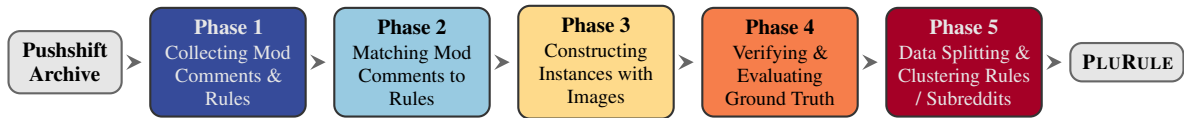


Figure 2: PLURULE construction pipeline

We then query the Reddit API to collect each subreddit’s full object. This helps to retrieve its current rules, infer its official language, and determine if it is NSFW (not-safe-for-work). Each subreddit must have at least one moderator comment and at least two explicit rules. We also exclude communities with adult content (NSFW). After filtering for these criteria, we obtain 17,468 subreddits with 131,400 rules and approximately 9 million moderator comments.

4.2 Phase 2: Rule Matching

Moderators often reference rules in their comments (e.g., “Rule 3 – No personal attacks”). Despite this, linking historical moderator comments to specific rule violations is a key challenge because rules evolve over time: new ones are added, old ones are deleted, and numbering and wording of existing rules can change. Since the Reddit API provides only present-day rules (as of November 2025), we match the full text of moderator comments to the full text of current rules.

We use Qwen3-Embedding-8B, a multilingual text embedding model, to encode all 9 million moderator comments and 131,400 rules as dense vectors. For each comment, we compute cosine similarity against all rules of the subreddit to which it belongs (7.5 rules on average), producing 90 million comment-rule scores.

We apply two thresholds to infer high-quality labels. First, a *match threshold* at the 99th percentile of similarity scores (0.78): only comment-rule pairs above this threshold count as matches. Second, an *ambiguity threshold* at the 98th percentile (0.75): if multiple rules for a single comment exceed this threshold, we discard the comment entirely to avoid inferring ambiguous labels. From 9 million moderator comments, 174,412 ambiguous cases are discarded, yielding 775,594 matched comments. In Phase 4, after additional filtering, we verify the quality of these matches.

4.3 Phase 3: Instance Construction

We wish to capture the complete conversational context of each violation, i.e., the full comment

thread leading to the rule-violating comment — the one to which the moderator replied. We collect all comments from the same submission by matching submission IDs in the Pushshift archives, yielding 84.9 million comments. We then build comment trees representing the reply structure. From each tree, we extract a *violating thread*: the path from the rule-violating comment up through its parent comments to the root submission.

Effective moderation requires a capability to discriminate between similar rule-violating and compliant comments within the same discussion. To this end, we pair each violating thread with a *compliant thread* — a discussion branch from the same submission that received no moderator action. We first select candidate compliant threads with depth equal to the violating thread (n) or, if not possible, $n - 1$.

We then apply five filtering criteria. For both violating and compliant threads, we exclude: (1) deleted/removed content or deleted users to ensure complete discussion context; (2) media in comments to limit images to submissions only; and (3) any moderator-authored comments in the thread to avoid back and forth discussions between a moderator and a user. For violating threads specifically, we exclude (4) edited leaf comments that became compliant after moderator intervention. For compliant threads specifically, we exclude (5) leaf comments with moderator replies to ensure no moderator flagged these comments as violations.

After filtering, we rank candidate compliant threads for each violating thread using three criteria to maximize the shared context between the two: (1) higher thread depth to prioritize n over $n - 1$; (2) higher number of common ancestors; and (3) lower vote score, to select less popular content that nevertheless complied with community rules. We select the compliant thread that ranks highest.

To complete the benchmark instances, we finally collect full submission objects from Pushshift for all thread pairs. We remove thread pairs for submissions with NSFW content, crossposts, and videos. For the remaining submissions, we download images using a priority hierarchy: gallery images,

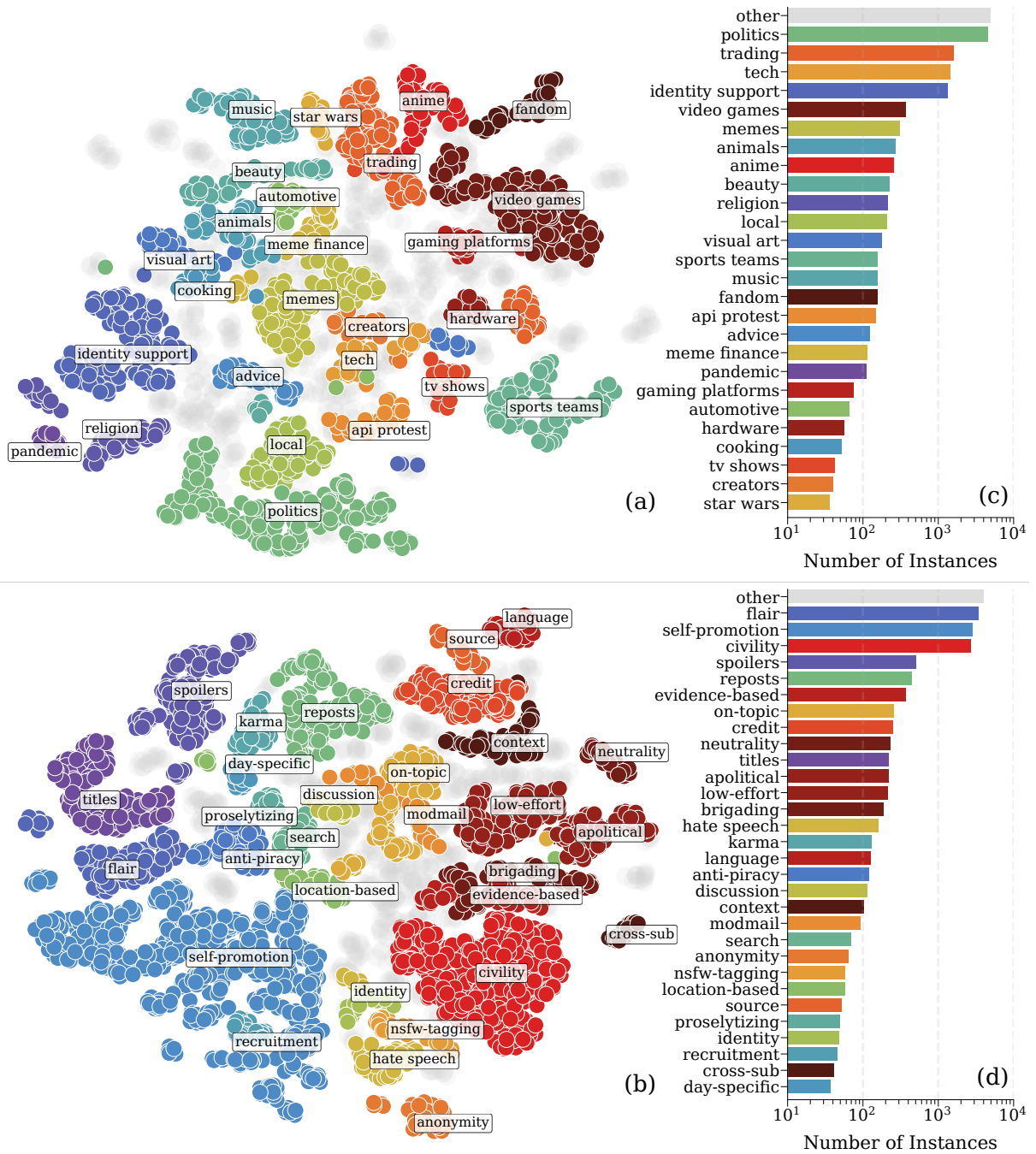


Figure 3: Two dimensions of pluralism in the PLURULE benchmark. Left: 2D UMAP visualizations of (a) 2,419 subreddits and (b) 3,692 rules, with colors indicating cluster assignments by HDBSCAN. Grey points represent unclustered items categorized as “other”. Right: Distributions of 17,313 instances across (c) 27 subreddit clusters and (d) 31 rule clusters, with bar colors matching the cluster colors.

Split	Instances	Comments	Images	Subreddits / Clusters	Rules / Clusters	Languages
Train	11,979	66,280	2,882	1,104 / 27	1,801 / 31	10
Val	1,809	9,810	512	703 / 27	779 / 31	10
Test	3,525	16,215	1,449	2,419 / 27	2,491 / 31	10
Total	17,313	92,305	4,843	2,419 / 27	3,692 / 31	10

Table 1: PLURULE statistics. Each instance consists of one rule-violating and one compliant comment from the same submission. All 2,419 subreddits appear in the test set. We count only the 10 languages with at least 10 instances each. There are 24 languages with at least one instance.

331 direct URLs, video thumbnails, and Reddit-cached
332 previews as fallbacks. Each download is validated
333 for image content type and capped at 50 MB. We
334 also exclude submissions with deleted/removed
335 content, deleted users and those posted by mod-
336 erators.

337 From 775,594 matched comments, we build
338 448,379 comment trees and successfully create
339 21,238 instances, totaling 42,476 threads. Most
340 failures for threads and submissions stem from
341 deleted/removed content or deleted users.

342 4.4 Phase 4: Verification

343 We use a large language model, Qwen3-30B-A3B-
344 Instruct, to verify the rule matches inferred in
345 Phase 2 for the benchmark instances. For each
346 instance, we present the model with the *moderator*
347 *comment* and *matched rule* to classify the comment
348 as: (a) stating a violation of the rule, (b) discussing
349 the rule, or (c) unrelated to the rule. We retain
350 as ground-truth labels only matched rules classi-
351 fied as (a), achieving an 81.5% verification rate
352 (17,313 out of 21,238 instances). This step filters
353 out incorrect matches and cases where moderators
354 mentioned a rule without enforcing it.

355 To evaluate the accuracy of these ground-truth
356 labels, three authors independently annotated 100
357 moderator comments sampled randomly from En-
358 glish subreddits. For each moderator comment,
359 annotators selected which subreddit rule was vio-
360 lated from the available options — the same task
361 performed by the matching pipeline. For 78 of
362 the comments in the sample, all three annotators
363 agreed on the label. In 16 cases, the label was
364 assigned based on a majority (two of the three an-
365 notators agreed). In the six remaining cases, there
366 was no majority agreement and the label was ad-
367 judicated after further inspection by one annotator.
368 Comparing the pipeline’s labels against this human-
369 established ground truth, we found 95% overall
370 accuracy: 100% on full-agreement cases (78/78),
371 81.25% on majority-agreement cases (13/16), and
372 66.67% on adjudicated cases (4/6).

373 4.5 Phase 5: Data Splitting and Clustering

374 We split the instances into training, validation, and
375 test sets using a strategy based on the number of
376 instances per subreddit. For subreddits with a sin-
377 gle instance, we allocate the instance to the test set.
378 For subreddits with two instances, we allocate one
379 instance to the training set and one to the test set.
380 For subreddits with 3–9 instances, we allocate one

381 each to test and validation sets, and the remaining
382 to the training set. For subreddits with 10 or more
383 instances, we use a 80/10/10 split for the training,
384 validation, and test sets. This ensures all commu-
385 nities appear in the test set while preventing any
386 single community from dominating the evaluation.

387 To analyze model accuracy across pluralistic
388 communities, we cluster subreddits and rules based
389 on their semantic embeddings. For subreddits,
390 we embed the subreddit name, title, and descrip-
391 tion. For rules, we embed the concatenation of
392 short name, description, and violation reason. We
393 apply UMAP for dimensionality reduction using
394 cosine distance on the 4,096-dimensional Qwen3-
395 Embedding-8B, then HDBSCAN for density-based
396 clustering (see Appendix A). Figure 3 shows that
397 rules cluster more cleanly than subreddits, reflect-
398 ing their greater semantic coherence compared to
399 the diversity of community topics.

400 We label each cluster using Qwen3-30B-A3B-
401 Thinking (see Appendix A) followed by manual
402 refinement. We assign these cluster labels to all
403 instances, enabling both fine-grained and category-
404 level evaluation. Table 1 provides full statistics of
405 the PLURULE dataset.

406 5 Evaluation

407 5.1 Experimental Setup

408 For each instance, models receive the subreddit
409 description, complete rule set, and surrounding
410 context, then select the correct answer from the
411 multiple-choice options. We report accuracy on the
412 test set, with the 50% baseline corresponding to
413 always predicting “No rules broken.” We compute
414 95% confidence intervals via bootstrap resampling
415 with 100 thousand iterations.

416 We evaluate three open-weight Vision-Language
417 Models from the Qwen3-VL family for their di-
418 versity in sizes (4B, 8B, and 30B) and OpenAI’s
419 flagship model GPT-5.2. For each Qwen model,
420 we test both Instruct and Thinking variants. For
421 GPT-5.2, we test with low and high reasoning ef-
422 fort. Qwen models use temperature 0 and seed 0
423 for reproducibility.

424 We use a two-stage evaluation pipeline. In
425 Stage 1, the model generates a free-form response
426 to the input (see Fig. 1). In Stage 2, we append
427 “Final Choice:” to prompt the model for its final
428 answer, then extract the selected option (a–h) using
429 a regular expression. For GPT-5.2, Stage 2 uses
430 Qwen3-VL-30B-Instruct for answer extraction.

Models	Qwen3-VL-4B		Qwen3-VL-8B		Qwen3-VL-30B		GPT-5.2	
Variants	Instruct	Thinking	Instruct	Thinking	Instruct	Thinking	Low	High
Comment Only	49.3	37.7	50.9	41.4	50.1	45.4	54.5	55.0
+Discussion	48.9 (-0.4)	40.2 (+2.5)	50.9 (+0.0)	44.6 (+3.2)	50.7 (+0.6)	46.8 (+1.4)	55.1 (+0.6)	55.8 (+0.8)
+Submission	47.6 (-1.3)	44.6 (+4.4)	49.8 (-1.1)	46.9 (+2.3)	51.6 (+0.9)	48.7 (+1.9)	57.1 (+2.0)	57.5 (+1.7)
+User	48.0 (+0.4)	44.5 (-0.1)	50.4 (+0.6)	47.2 (+0.3)	52.2 (+0.6)	49.2 (+0.5)	57.8 (+0.7)	57.6 (+0.1)
+Images	47.2 (-0.8)	44.7 (+0.2)	50.8 (+0.4)	45.7 (-1.5)	52.7 (+0.5)	48.8 (-0.4)	57.3 (-0.5)	58.2 (+0.6)
Baseline	50.0							

Table 2: Accuracy (%) across models and context levels. Numbers in parentheses show differences from the previous row. Best-performing contexts for each model variant are highlighted in bold. 95% CI for all values do not exceed $\pm 1.1\%$. The baseline corresponds to always predicting “no rules broken.”

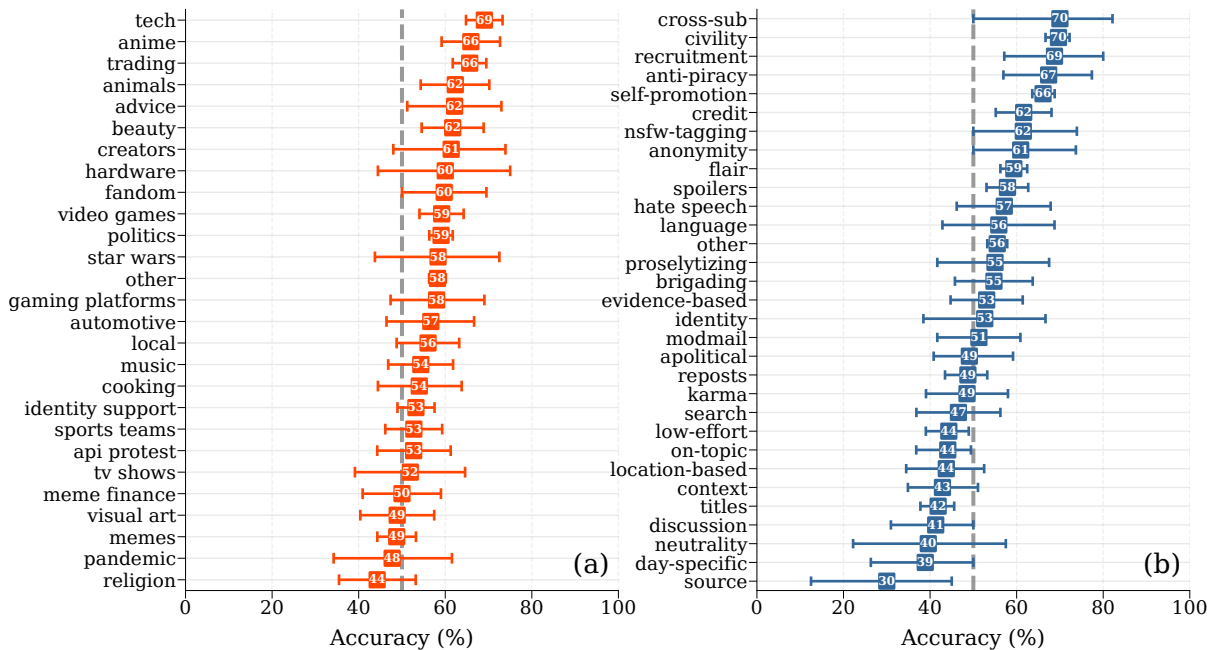


Figure 4: Accuracy for GPT-5.2 (high reasoning) with full context by (a) subreddit clusters and (b) rule clusters. Error bars show 95% CI. Dashed lines indicate the baseline.

To understand which contextual signals aid rule violation detection, we test five cumulative context levels, where each level adds information to the previous: (1) **Comment Only** — the target comment; (2) **+ Discussion** — the full comment thread leading to the target comment; (3) **+ Submission** — title and body text of the original post that initiated the discussion; (4) **+ User** — anonymized author labels (USER1, USER2, etc.) to track participants; and (5) **+ Images** — media from the submission, when available. All levels include the subreddit description and complete rule set as baseline. We further analyze performance by subreddit cluster and rule cluster to identify which community types and rule categories pose the greatest challenges.

5.2 Results

Table 2 reports accuracy across model sizes, context configurations, and reasoning variants (see Appendix B for violating/compliant comment breakdowns.) GPT-5.2 substantially outperforms all Qwen variants, with high reasoning effort achieving 58.2% accuracy — 8 points above the 50% baseline — while Qwen models barely exceed the baseline regardless of scale.

Focusing on GPT-5.2 with high reasoning, additional context provides limited signal for rule violation detection: using the full context improves performance by only 3.2 percentage points above comment-only, with submissions yielding the largest gain (+1.7 points).

461 Extended reasoning does not necessarily help
462 with this task: Qwen Thinking variants underper-
463 form their Instruct counterparts, while the differ-
464 ences between GPT-5.2 with low and high reason-
465 ing are not significant.

466 Figure 4 breaks down accuracy by subreddit and
467 rule clusters for GPT-5.2 (high reasoning) using
468 full context. See Appendix B for breakdowns by vi-
469 olating/compliant comments and languages. Perfor-
470 mance varies substantially across community clus-
471 ters: tech (69%), anime (66%), and trading (66%)
472 subreddits achieve the highest accuracy, while re-
473 ligion (44%), pandemic (48%), and memes (49%)
474 fall below baseline. Rule clusters show even greater
475 variability: civility (70%), recruitment (69%), and
476 self-promotion (66%) rules are detected reliably,
477 but source requirements (30%), neutrality policies
478 (40%), and on-topic rules (44%) perform worse
479 than baseline. Similar results are obtained with
480 other models (Appendix B).

481 Civility and self-promotion appear in far more
482 subreddits and contain far more member rules than
483 any other category (Appendix B). These univer-
484 sal categories consistently rank among the highest-
485 accuracy rule clusters. In contrast, pluralistic rules
486 requiring local context remain challenging.

487 6 Discussion

488 PLURULE provides a testbed for measuring
489 progress toward models that can adapt to diverse
490 community standards. Unlike prior datasets that
491 use coarse-grained categories or single-rule binary
492 classification, PLURULE requires models to dis-
493 tinguish among all of a community’s rules simul-
494 taneously — mirroring the decision space faced
495 by human moderators. We evaluate whether a sin-
496 gle model can serve as an expert moderator across
497 thousands of communities.

498 Our results reveal a fundamental gap: mod-
499 els succeed on universal violations like civility
500 and self-promotion, but struggle with pluralistic
501 rules that vary across communities. The two-
502 dimensional variability — 44%–69% accuracy
503 across community types and 30%–70% across rule
504 categories — suggests that VLMs may be internal-
505 izing universal standards from training.

506 Beyond content moderation, PLURULE evalu-
507 ates whether AI systems can respect diverse hu-
508 man communities rather than imposing uniform
509 standards. Fine-tuning on community-specific
510 examples may help models learn local norms,

though scalability across thousands of communi-
ties remains a challenge. Retrieval-augmented ap-
proaches that condition on historical moderation
decisions offer another promising direction. The
semantic clustering we provide enables analysis
of transfer learning: can models trained on one
community or rule type generalize to similar ones?

518 7 Limitations

519 PLURULE was constructed from publicly available
520 data where moderators left comments citing rule
521 violations. Private moderator communications, re-
522 moved content, and shadow-banned posts are not
523 accessible, meaning communities that moderate
524 silently are underrepresented. This likely biases
525 the dataset toward less severe violations, as seri-
526 ous offenses are often removed without comment.
527 English-language subreddits dominate due to Red-
528 dit’s user demographics, and findings may not ge-
529 neralize to platforms with different community struc-
530 tures or moderation practices.

531 Our pipeline matches historical moderator com-
532 ments (2005–2023) to rule sets retrieved in Novem-
533 ber 2025. While semantic matching handles rule
534 rewording and renumbering, it cannot account for
535 rules that were added, removed, or fundamentally
536 changed over time. Some matches may therefore
537 be anachronistic.

538 Finally, certain violations require information
539 unavailable in our dataset. Detecting ban evasion
540 or repeat offenders requires historical user data that
541 we do not collect for privacy reasons.

542 8 Ethical Considerations

543 PLURULE contains only publicly posted content
544 from the Pushshift archives. We anonymize all
545 usernames by replacing them with generic labels
546 (USER1, USER2, etc.). We do not collect or
547 release any private user metadata such as IP ad-
548 dresses, email addresses, or account history.

549 Models trained on PLURULE could potentially
550 be misused to evade moderation by learning what
551 content triggers enforcement. This risk is inherent
552 to any dataset that captures moderation decisions.
553 We believe the research benefits outweigh this risk,
554 as understanding moderation patterns is essential
555 for developing robust systems.

556 References

557 Dalia Ali, Dora Zhao, Allison Koenecke, and Orestis
558 Papakyriakopoulos. 2025. Operationalizing pluralis-

559	tic values in large language model alignment reveals trade-offs in safety, inclusivity, and model behavior. <i>arXiv preprint arXiv:2511.14476</i> .	
560		
561		
562	Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. 2020. The Pushshift Reddit Dataset . <i>arXiv preprint ArXiv:2001.08435 [cs]</i> .	
563		
564		
565		
566	Reuben Binns, Michael Veale, Max Van Kleek, and Nigel Shadbolt. 2017. Like Trainer, Like Bot? Inheritance of Bias in Algorithmic Content Moderation . In <i>Social Informatics</i> , pages 405–415, Cham. Springer International Publishing.	
567		
568		
569		
570		
571	Abeba Birhane, Vinay Uday Prabhu, and Evelyn Kahembwe. 2021. Multimodal datasets: Misogyny, pornography, and malignant stereotypes . In <i>Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency</i> , pages 303–314. ACM.	
572		
573		
574		
575		
576	Su Lin Blodgett, Lisa Green, and Brendan O’Connor. 2016. Demographic Dialectal Variation in Social Media: A Case Study of African-American English . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1119–1130, Austin, Texas. Association for Computational Linguistics.	
577		
578		
579		
580		
581		
582		
583	Edoardo Celeste, Nicola Palladino, Dennis Redeker, and Kinfe Yilma. 2023. Platform policies versus human rights standards. In <i>The content governance dilemma: Digital constitutionalism, social media and the search for a global standard</i> , pages 93–129. Springer.	
584		
585		
586		
587		
588		
589	Eshwar Chandrasekharan, Chaitrali Gandhi, Matthew Wortley Mustelier, and Eric Gilbert. 2019. Crossmod: A cross-community learning-based system to assist reddit moderators. <i>Proceedings of the ACM on human-computer interaction</i> , 3(CSCW):1–30.	
590		
591		
592		
593		
594		
595	Eshwar Chandrasekharan and Eric Gilbert. 2019. Hybrid approaches to detect comments violating macro norms on reddit. <i>arXiv preprint arXiv:1904.03596</i> .	
596		
597		
598	Eshwar Chandrasekharan, Mattia Samory, Shagun Jhaver, Hunter Charvat, Amy Bruckman, Cliff Lampe, Jacob Eisenstein, and Eric Gilbert. 2018. The internet’s hidden rules: An empirical study of reddit norm violations at micro, meso, and macro scales. <i>Proceedings of the ACM on Human-Computer Interaction</i> , 2(CSCW):1–25.	
599		
600		
601		
602		
603		
604		
605	Joseph Paul Cohen and Henry Z. Lo. 2014. Academic Torrents: A Community-Maintained Distributed Repository . In <i>Proceedings of the 2014 Annual Conference on Extreme Science and Engineering Discovery Environment, XSEDE ’14</i> , pages 1–2, New York, NY, USA. Association for Computing Machinery.	
606		
607		
608		
609		
610		
611		
612	Ángel Díaz and Laura Hecht-Felella. 2021. Double standards in social media content moderation . Technical report, Brennan Center for Justice at New York University School of Law. Accessed: 24 Dec 2025.	
613		
614		
615		
	Bryan Dosono and Bryan Semaan. 2019. Moderation Practices as Emotional Labor in Sustaining Online Communities: The Case of AAPI Identity Work on Reddit . In <i>Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, CHI ’19</i> , pages 1–13, New York, NY, USA. Association for Computing Machinery.	616
		617
		618
		619
		620
		621
		622
	Casey Fiesler, Jialun Jiang, Joshua McCann, Kyle Frye, and Jed Brubaker. 2018. Reddit Rules! Characterizing an Ecosystem of Governance . <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , 12(1).	623
		624
		625
		626
		627
	Tarleton Gillespie. 2020. Content moderation, ai, and the question of scale. <i>Big Data & Society</i> , 7(2):2053951720943234.	628
		629
		630
	Raul Gomez, Jaume Gibert, Lluís Gomez, and Dimosthenis Karatzas. 2020. Exploring Hate Speech Detection in Multimodal Publications . In <i>2020 IEEE Winter Conference on Applications of Computer Vision (WACV)</i> , pages 1459–1467, Snowmass Village, CO, USA. IEEE.	631
		632
		633
		634
		635
		636
	Google. 2025. Youtube community guidelines enforcement .	637
		638
	Rachel Griffin. 2024. The heteronormative male gaze: Experiences of sexual content moderation among queer instagram users in berlin. <i>International journal of communications, network and system sciences</i> , 18:1266–1288.	639
		640
		641
		642
		643
	Zihao He, Jonathan May, and Kristina Lerman. 2024. CPL-NoViD: Context-Aware Prompt-Based Learning for Norm Violation Detection in Online Communities . <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , 18:569–582.	644
		645
		646
		647
		648
	Benjamin Mako Hill. 2019. How Discord moderators build innovative solutions to problems of scale with the past as a guide .	649
		650
		651
	Nhat M. Hoang, Xuan Long Do, Duc Anh Do, Duc Anh Vu, and Luu Anh Tuan. 2024. ToXCL: A Unified Framework for Toxic Speech Detection and Explanation . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 6460–6472, Mexico City, Mexico. Association for Computational Linguistics.	652
		653
		654
		655
		656
		657
		658
		659
		660
	Jialun’Aaron’ Jiang, Skyler Middler, Jed R Brubaker, and Casey Fiesler. 2020. Characterizing community guidelines on social media platforms. In <i>Companion Publication of the 2020 Conference on Computer Supported Cooperative Work and Social Computing</i> , pages 287–291.	661
		662
		663
		664
		665
		666
	Hanlin Li, Brent Hecht, and Stevie Chancellor. 2022a. All That’s Happening behind the Scenes: Putting the Spotlight on Volunteer Moderator Labor in Reddit . <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , 16:584–595.	667
		668
		669
		670
		671

672	Hanlin Li, Brent Hecht, and Stevie Chancellor. 2022b.	Maarten Sap, Swabha Swayamdipta, Laura Vianna,	728
673	Measuring the Monetary Value of Online Volunteer	Xuhui Zhou, Yejin Choi, and Noah A. Smith. 2022.	729
674	Work . <i>Proceedings of the International AAAI Con-</i>	Annotators with attitudes: How annotator beliefs	730
675	ference on Web and Social Media , 16:596–606.	and identities bias toxic language detection . In <i>Pro-</i>	731
		ceedings of the 2022 Conference of the North Amer-	732
676	Jessa Lingel and Adam Golub. 2015. In face on face-	ican Chapter of the Association for Computational	733
677	book: Brooklyn’s drag community and sociotech-	Linguistics: Human Language Technologies , pages	734
678	nical practices of online communication. <i>Journal</i>	5884–5906, Seattle, United States. Association for	735
679	of Computer-Mediated Communication , 20(5):536–	Computational Linguistics.	736
680	553.		
681	Travis Lloyd, Joseph Reagle, and Mor Naaman. 2025.	Andrew D. Selbst, Danah Boyd, Sorelle A. Friedler,	737
682	‘There Has To Be a Lot That We’re Missing’: Mod-	Suresh Venkatasubramanian, and Janet Vertesi. 2019.	738
683	erating AI-Generated Content on Reddit . <i>Proc.</i>	Fairness and abstraction in sociotechnical systems .	739
684	ACM Hum.-Comput. Interact. , 9(7):CSCW264:1–	In <i>Proceedings of the 2019 ACM Conference on Fair-</i>	740
685	CSCW264:24.	ness, Accountability, and Transparency , pages 59–68.	741
		ACM.	742
686	J. Nathan Matias. 2019. The Civic Labor of Volun-	Brooklyn Sheppard, Anna Richter, Allison Cohen, Eliza-	743
687	teer Moderators Online . <i>Social Media + Society</i> ,	beth Smith, Tamara Kneese, Carolyne Pelletier, Ioana	744
688	5(2):2056305119836778. Publisher: SAGE Publica-	Baldini, and Yue Dong. 2024. Biasly: An Expert-	745
689	tions Ltd.	Annotated Dataset for Subtle Misogyny Detection	746
		and Mitigation . In <i>Findings of the Association for</i>	747
690	Meta. 2024. How enforcement technology works .	Computational Linguistics: ACL 2024 , pages 427–	748
		452, Bangkok, Thailand. Association for Computa-	749
691	Huy Nghiem and Hal Daumé Iii. 2024. HateCOT: An	tional Linguistics.	750
692	Explanation-Enhanced Dataset for Generalizable Of-		
693	fensive Speech Detection via Large Language Mod-	X Transparency. 2025. Global transparency report .	751
694	els . In <i>Findings of the Association for Computational</i>		
695	Linguistics: EMNLP 2024 , pages 5938–5956, Mi-	A Supplementary Methods	752
696	ami, Florida, USA. Association for Computational		
697	Linguistics.	For clustering subreddits and rules, we per-	753
		form grid searches over UMAP and HDB-	754
698	Chan Young Park, Julia Mendelsohn, Karthik Radhakr-	SCAN parameters: n_neighbors, n_components,	755
699	ishnan, Kinjal Jain, Tushar Kanakagiri, David Jur-	min_cluster_size, and min_samples. We maximize	756
700	gens, and Yulia Tsvetkov. 2021. Detecting Commu-	DBCV (Density-Based Cluster Validity), which	757
701	nity Sensitive Norm Violations in Online Conversa-	measures cluster separation and coherence.	758
702	tions . In <i>Findings of the Association for Computa-</i>	For subreddits, optimal parameters yield DBCV	759
703	tional Linguistics: EMNLP 2021 , pages 3386–3397,	= 0.400 with 26 clusters (62 items per cluster on	760
704	Punta Cana, Dominican Republic. Association for	average). We treat the remaining subreddits as the	761
705	Computational Linguistics.	‘other’ cluster (33.5% noise). For rules, optimal	762
		parameters yield DBCV = 0.524 with 30 clusters	763
706	Inioluwa Deborah Raji, Andrew Smart, Rebecca N.	(100 items per cluster on average). We treat the	764
707	White, Margaret Mitchell, Timnit Gebru, Ben	remaining rules as the ‘other’ cluster (18.6% noise).	765
708	Hutchinson, Jamila Smith-Loud, Daniel Theron, and	We assign semantic labels to each cluster us-	766
709	Parker Barnes. 2020. Closing the ai accountability	ing Qwen3-30B-A3B-Thinking. For each cluster,	767
710	gap: defining an end-to-end framework for internal	we prompt the model to generate 10 candidate la-	768
711	algorithmic auditing . In <i>Proceedings of the 2020</i>	bels and select the most common one via majority	769
712	Conference on Fairness, Accountability, and Trans-	voting. We then manually verify and refine the	770
713	parency, FAT* ’20 , page 33–44, New York, NY, USA.	labels for consistency. This produces 31 rule clus-	771
714	Association for Computing Machinery.	ters (e.g., civility, self-promotion, spoiler policy,	772
		flair enforcement) and 27 subreddit clusters (e.g.,	773
715	Reddit. 2025. Reddiquette . https://support.	politics, meme culture, video games, trading).	774
716	reddithelp.com/hc/en-us/articles/	B Extended Results	775
717	205926439-Reddiquette . Accessed: 2025-12-23.		
		Recall Breakdowns. Table 3 reports recall for vi-	776
718	C. J. Robinson. 2025. X Is Using AI Fact-Checkers .	olating and compliant comments across all models	777
719	https://www.cjr.org/analysis/x-twitter-ai-fact-	and context levels.	778
720	checkers-community-notes-misinformation-		
721	bots.php .		
722	Pratik Sachdeva, Renata Barreto, Geoff Bacon, Alexan-		
723	der Sahn, Claudia Von Vacano, and Chris Kennedy.		
724	2022. The measuring hate speech corpus: Leverag-		
725	ing rasch measurement theory for data perspectivism.		
726	In <i>Proceedings of the 1st Workshop on Perspectivist</i>		
727	Approaches to NLP@ LREC2022 , pages 83–94.		

779 **GPT-5.2 Analysis.** Figure 5 shows recall by sub-
780 reddit and rule clusters for GPT-5.2 (high reason-
781 ing) with full context. Figure 6 shows instance
782 distribution and accuracy across 10 languages.

783 **Qwen3-VL Analysis.** For each Qwen3-VL
784 model (4B, 8B, 30B), we provide three fig-
785 ures: accuracy by subreddit and rule clusters
786 (Figures 7, 10, 13), recall breakdowns (Fig-
787 ures 8, 11, 14), and language analysis (Fig-
788 ures 9, 12, 15).

789 **Universality Correlation.** To examine whether
790 models perform better on universal rule clusters, we
791 measure universality in two ways: (a) the number
792 of subreddits containing at least one rule in the
793 cluster, and (b) the number of member rules in the
794 cluster. Figures 16, 17, 18, and 19 plot accuracy
795 against both measures. Spearman correlations are
796 weak ($\rho = 0.11-0.40$), with only Qwen3-VL-8B
797 reaching significance. However, civility and self-
798 promotion — the most universal by both measures
799 — consistently rank among the highest-accuracy
800 clusters, with the exception of self-promotion for
801 Qwen3-VL-4B.

Model	Variant	Metric	Comment	+Discussion	+Submission	+User	+Images
Qwen3-VL-4B	Instruct	V. Recall	23.8	26.8 (+3.0)	29.6 (+2.8)	29.6 (+0.0)	31.3 (+1.7)
		C. Recall	74.9	70.9 (-4.0)	65.5 (-5.4)	66.4 (+0.9)	63.1 (-3.3)
		Accuracy	49.3	48.9 (-0.4)	47.6 (-1.3)	48.0 (+0.4)	47.2 (-0.8)
	Thinking	V. Recall	25.8	27.1 (+1.3)	27.7 (+0.6)	29.6 (+1.9)	28.9 (-0.7)
		C. Recall	49.7	53.4 (+3.7)	61.5 (+8.1)	59.4 (-2.1)	60.5 (+1.1)
		Accuracy	37.7	40.2 (+2.5)	44.6 (+4.4)	44.5 (-0.1)	44.7 (+0.2)
Qwen3-VL-8B	Instruct	V. Recall	27.7	28.0 (+0.3)	28.0 (+0.0)	29.2 (+1.2)	30.9 (+1.7)
		C. Recall	74.1	73.8 (-0.3)	71.6 (-2.2)	71.5 (-0.1)	70.8 (-0.7)
		Accuracy	50.9	50.9 (+0.0)	49.8 (-1.1)	50.4 (+0.6)	50.8 (+0.4)
	Thinking	V. Recall	33.0	33.9 (+0.9)	33.5 (-0.4)	35.1 (+1.6)	33.6 (-1.5)
		C. Recall	49.8	55.3 (+5.5)	60.4 (+5.1)	59.3 (-1.1)	57.8 (-1.5)
		Accuracy	41.4	44.6 (+3.2)	46.9 (+2.3)	47.2 (+0.3)	45.7 (-1.5)
Qwen3-VL-30B	Instruct	V. Recall	31.3	31.7 (+0.4)	31.5 (-0.2)	32.7 (+1.2)	33.7 (+1.0)
		C. Recall	68.9	69.6 (+0.7)	71.6 (+2.0)	71.7 (+0.1)	71.8 (+0.1)
		Accuracy	50.1	50.7 (+0.6)	51.6 (+0.9)	52.2 (+0.6)	52.7 (+0.5)
	Thinking	V. Recall	38.4	41.4 (+3.0)	40.2 (-1.2)	41.5 (+1.3)	41.0 (-0.5)
		C. Recall	52.4	52.1 (-0.3)	57.3 (+5.2)	56.9 (-0.4)	56.6 (-0.3)
		Accuracy	45.4	46.8 (+1.4)	48.7 (+1.9)	49.2 (+0.5)	48.8 (-0.4)
GPT-5.2	Low	V. Recall	40.7	41.8 (+1.1)	42.1 (+0.3)	43.4 (+1.3)	42.0 (-1.4)
		C. Recall	68.3	68.3 (+0.0)	72.2 (+3.9)	72.3 (+0.1)	72.6 (+0.3)
		Accuracy	54.5	55.1 (+0.6)	57.1 (+2.0)	57.8 (+0.7)	57.3 (-0.5)
	High	V. Recall	41.6	43.1 (+1.5)	43.1 (+0.0)	43.8 (+0.7)	43.7 (-0.1)
		C. Recall	68.4	68.5 (+0.1)	72.0 (+3.5)	71.3 (-0.7)	72.7 (+1.4)
		Accuracy	55.0	55.8 (+0.8)	57.5 (+1.7)	57.6 (+0.1)	58.2 (+0.6)
Baseline		V. Recall					0.0
		C. Recall					100.0
		Accuracy					50.0

Table 3: Violating recall, compliant recall, and accuracy (%) across different models and contexts on the test set. Numbers in parentheses indicate differences compared to values in the previous column. Best-performing contexts for each model variant are highlighted in bold. All values have 95% CI of at most $\pm 1.7\%$.

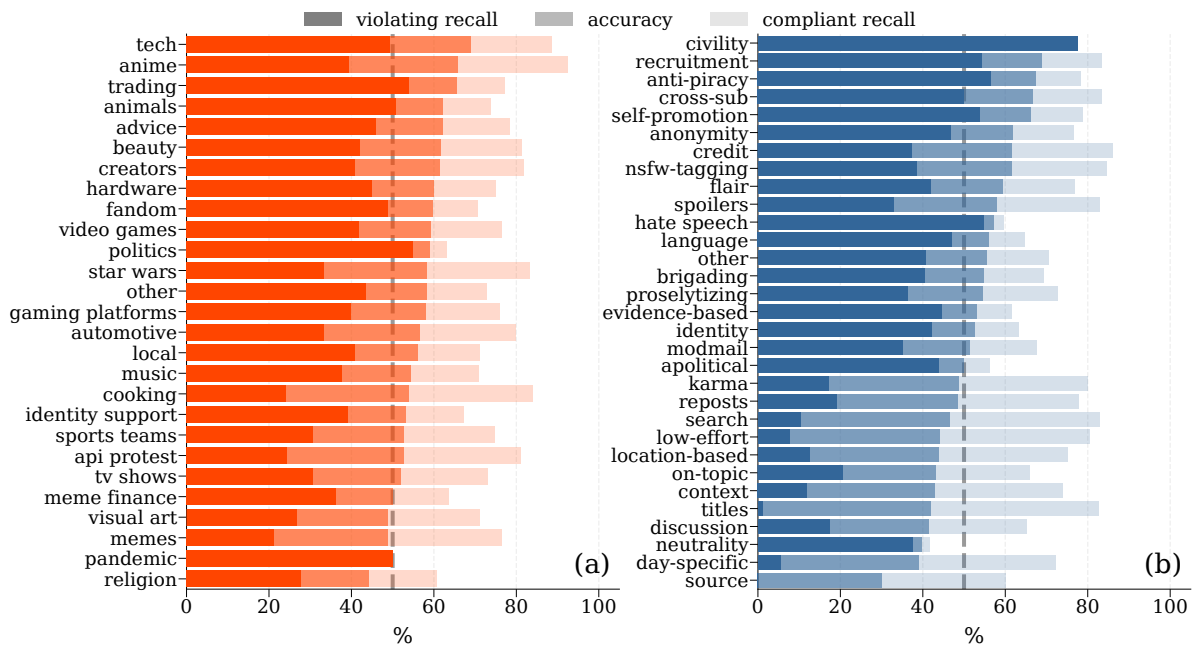


Figure 5: GPT-5.2 (high reasoning) recall and accuracy with full context by (a) subreddit cluster and (b) rule cluster on the test set. Stacked bars show violating and compliant recall. Bars sorted by accuracy. Dashed lines indicate the 50% baseline for accuracy.

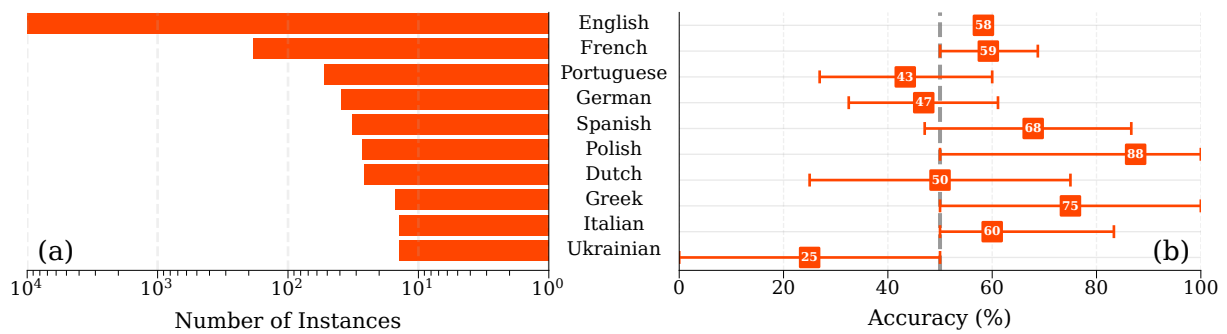


Figure 6: (a) Distribution of PluRule instances across 10 languages and (b) accuracy of the corresponding language in the test set by GPT-5.2 (high reasoning) with full context. Error bars show 95% CI. Dashed line indicates the 50% baseline.

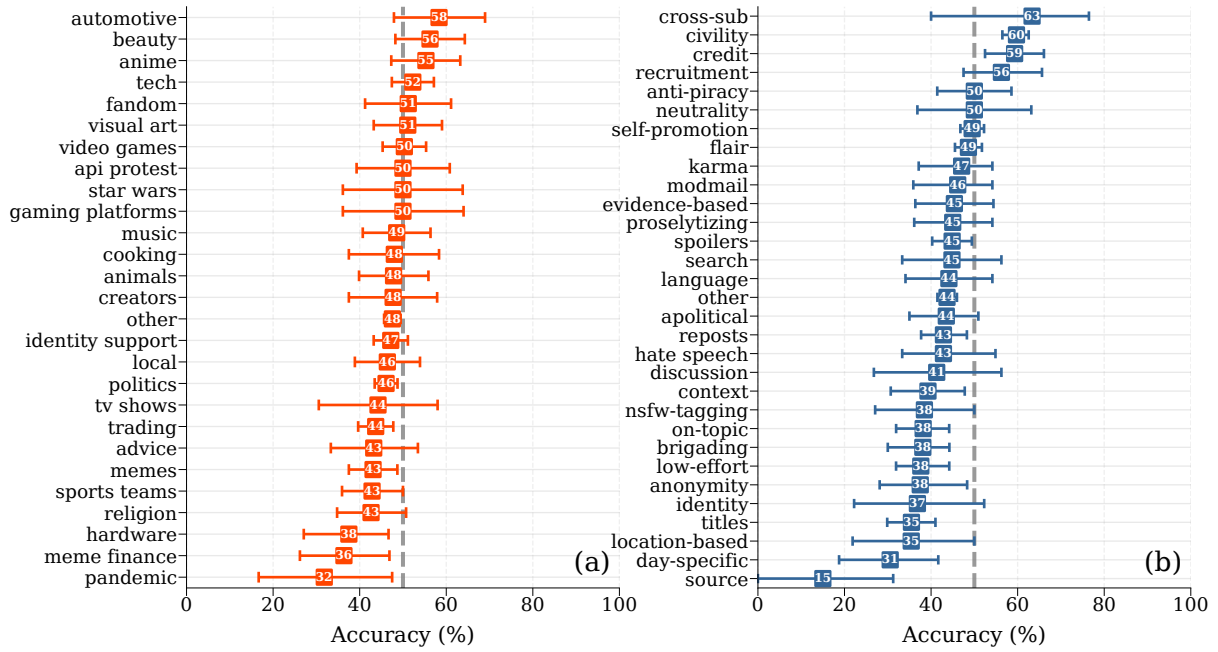


Figure 7: Accuracy for Qwen3-VL-4B (instruct) with full context by (a) subreddit clusters and (b) rule clusters. Error bars show 95% CI. Dashed lines indicate the baseline.

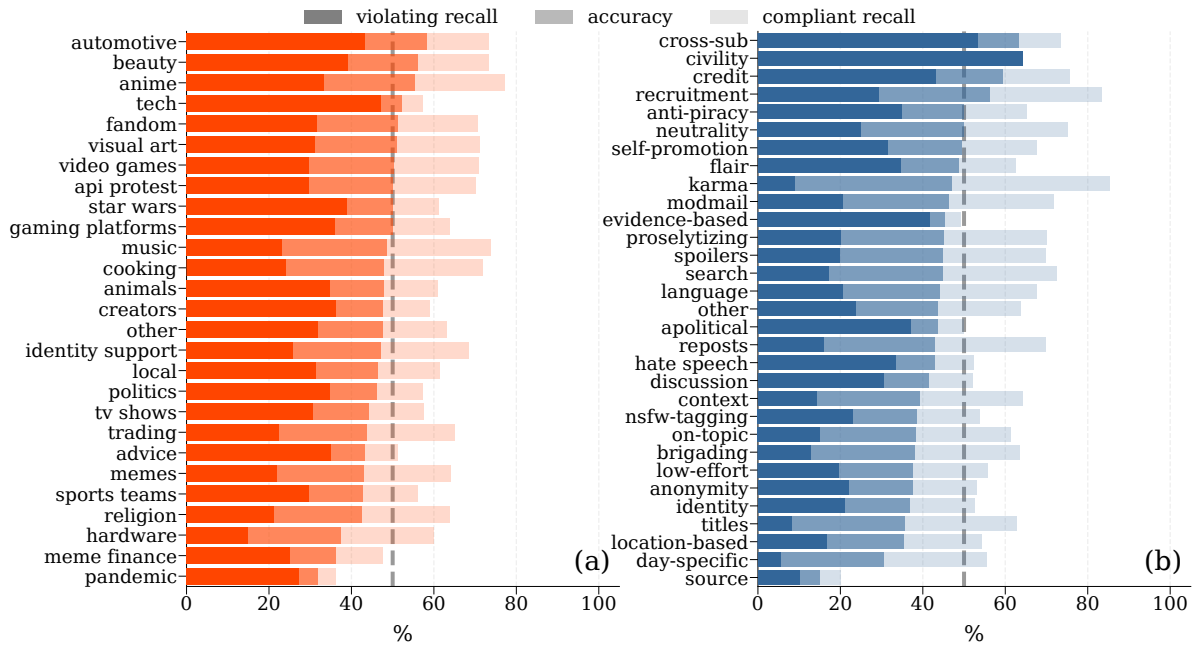


Figure 8: Qwen3-VL-4B (instruct) recall and accuracy with full context by (a) subreddit cluster and (b) rule cluster. Stacked bars show violating and compliant recall. Dashed lines indicate the 50% baseline.

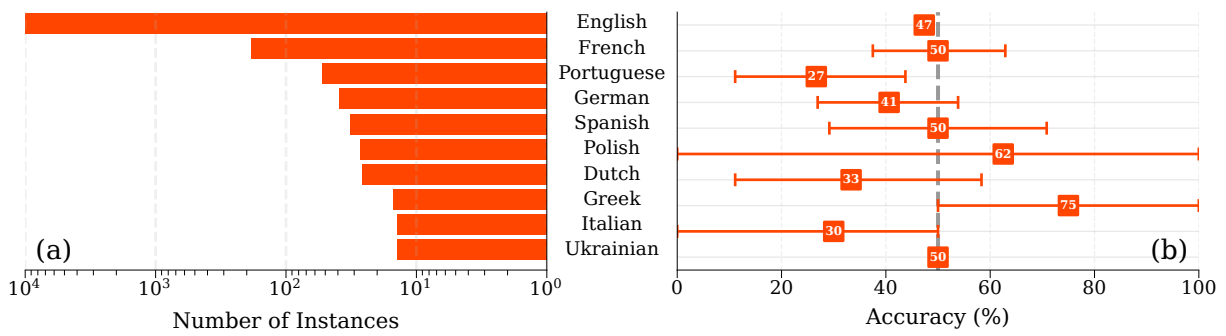


Figure 9: (a) Distribution of PluRule instances across 10 languages and (b) accuracy by language for Qwen3-VL-4B (instruct) with full context. Error bars show 95% CI. Dashed line indicates the 50% baseline.

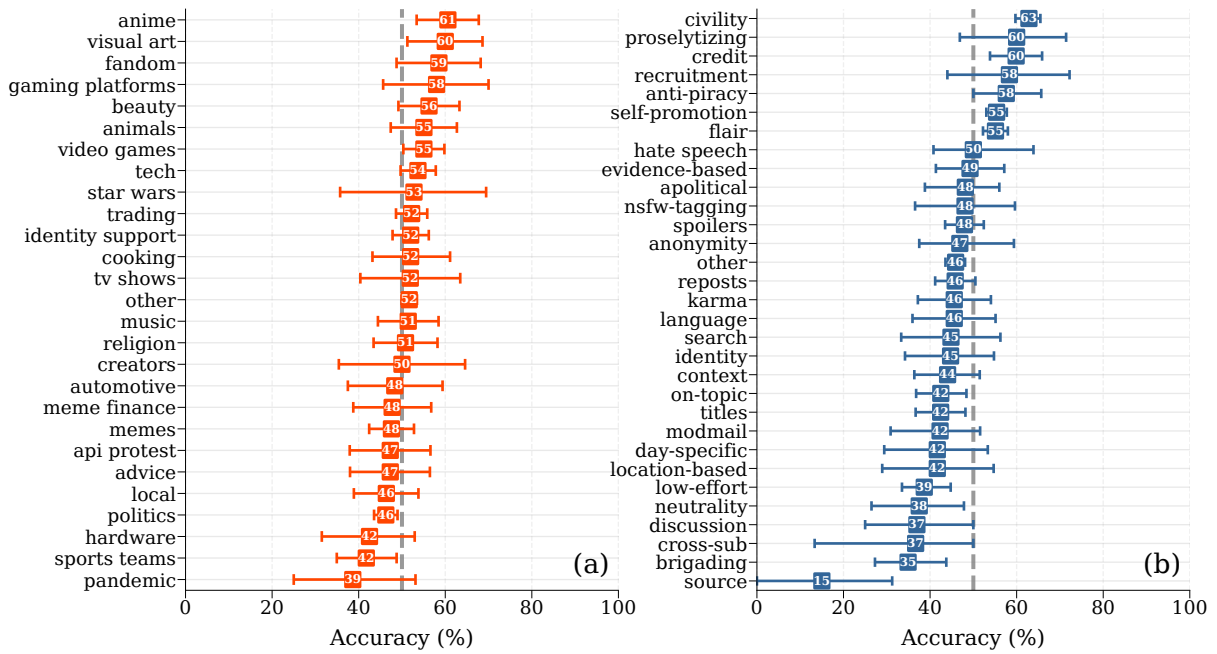


Figure 10: Accuracy for Qwen3-VL-8B (instruct) with full context by (a) subreddit clusters and (b) rule clusters. Error bars show 95% CI. Dashed lines indicate the baseline.

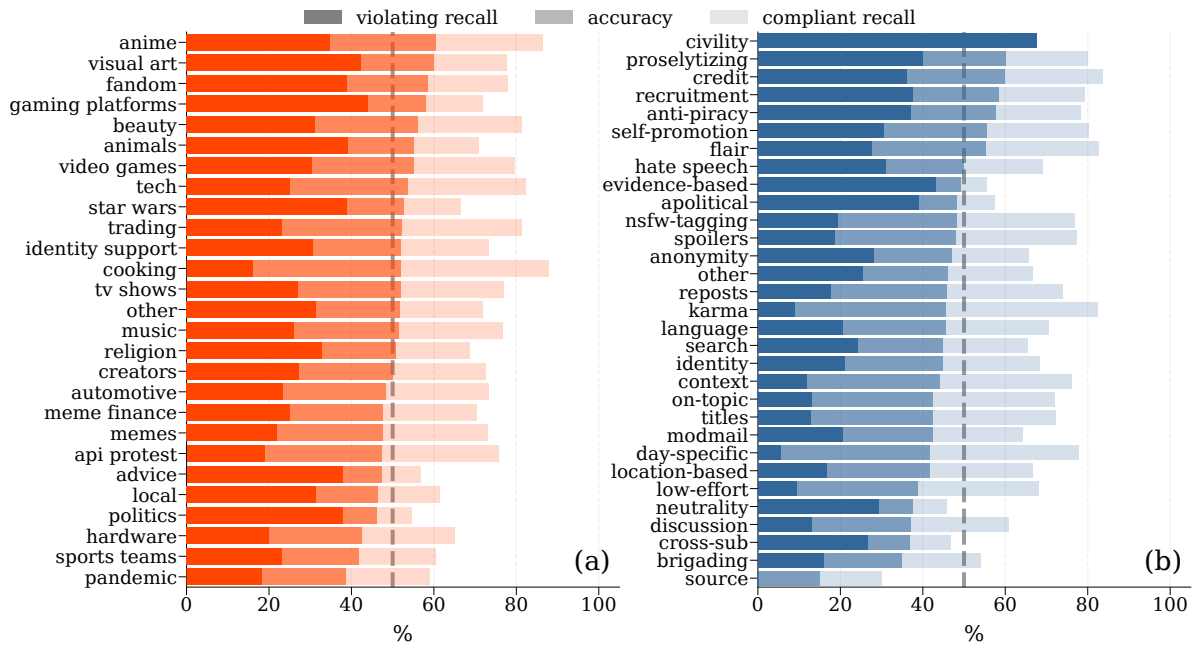


Figure 11: Qwen3-VL-8B (instruct) recall and accuracy with full context by (a) subreddit cluster and (b) rule cluster. Stacked bars show violating and compliant recall. Dashed lines indicate the 50% baseline.

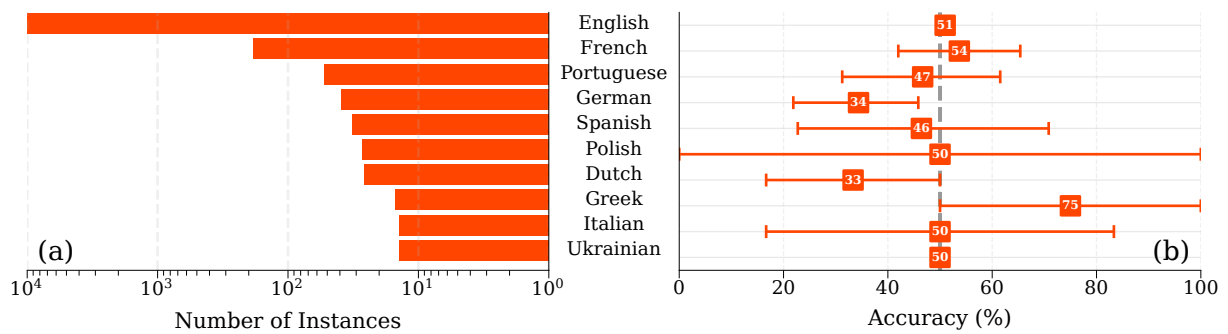


Figure 12: (a) Distribution of PluRule instances across 10 languages and (b) accuracy by language for Qwen3-VL-8B (instruct) with full context. Error bars show 95% CI. Dashed line indicates the 50% baseline.

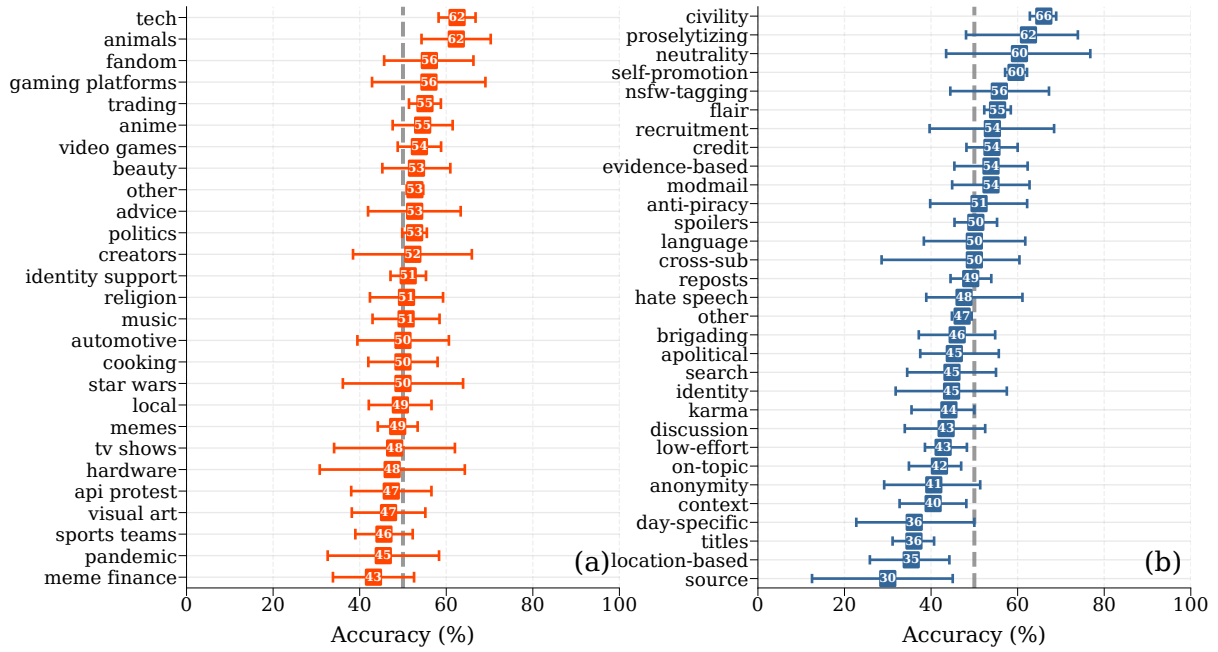


Figure 13: Accuracy for Qwen3-VL-30B (instruct) with full context by (a) subreddit clusters and (b) rule clusters. Error bars show 95% CI. Dashed lines indicate the baseline.

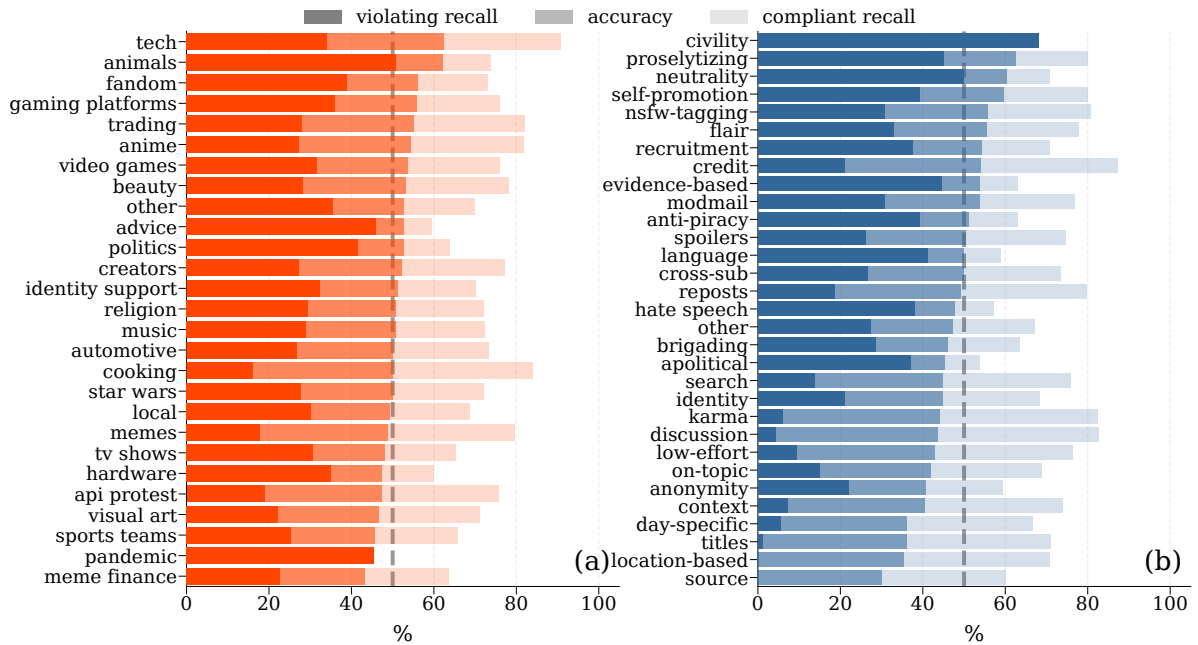


Figure 14: Qwen3-VL-30B (instruct) recall and accuracy with full context by (a) subreddit cluster and (b) rule cluster. Stacked bars show violating and compliant recall. Dashed lines indicate the 50% baseline.

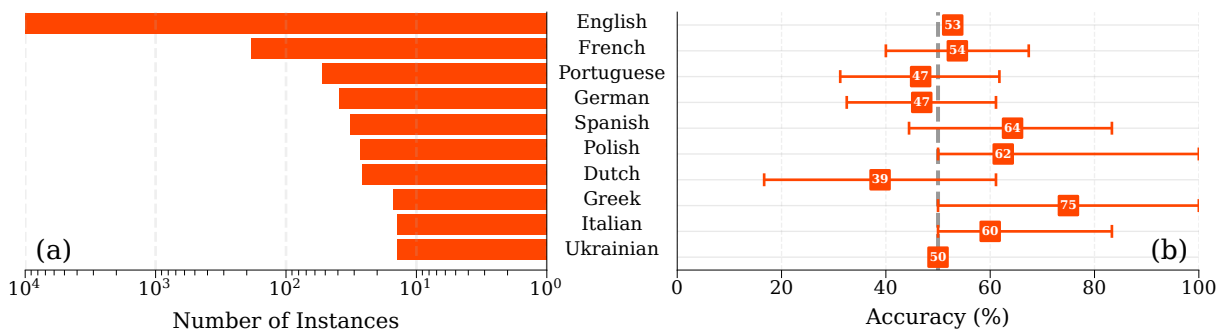


Figure 15: (a) Distribution of PluRule instances across 10 languages and (b) accuracy by language for Qwen3-VL-30B (instruct) with full context. Error bars show 95% CI. Dashed line indicates the 50% baseline.

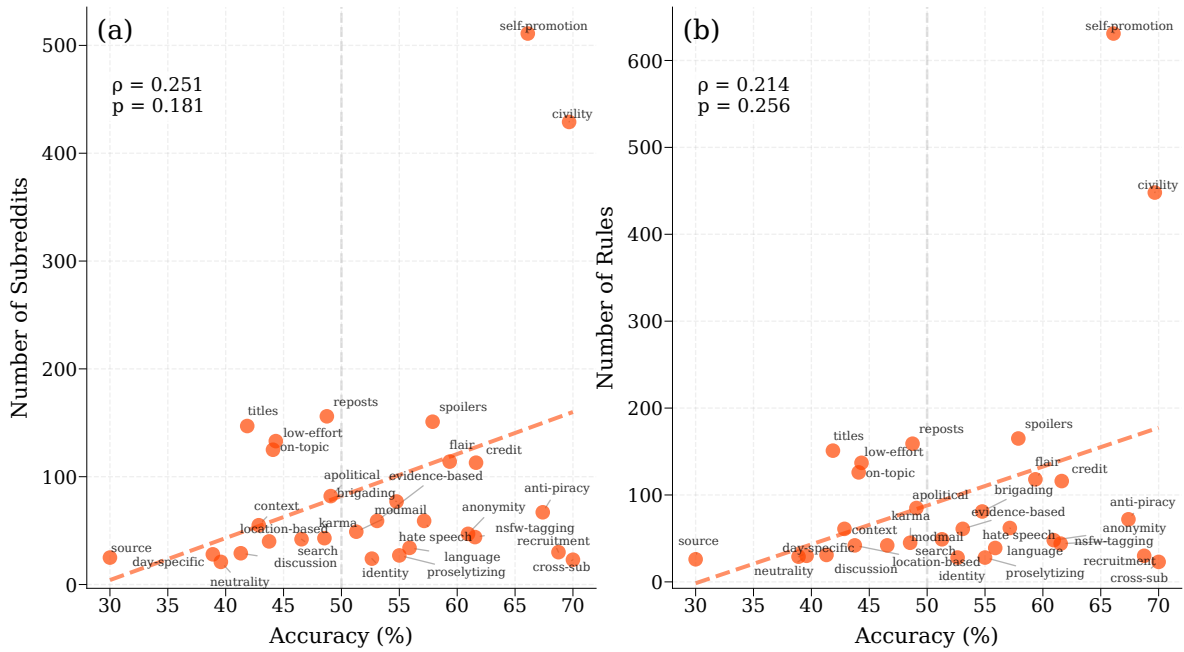


Figure 16: GPT-5.2 (high reasoning) accuracy with full context compared against (a) the number of subreddits containing rules in each cluster and (b) the number of rules per cluster. Dashed lines indicate the 50% baseline.

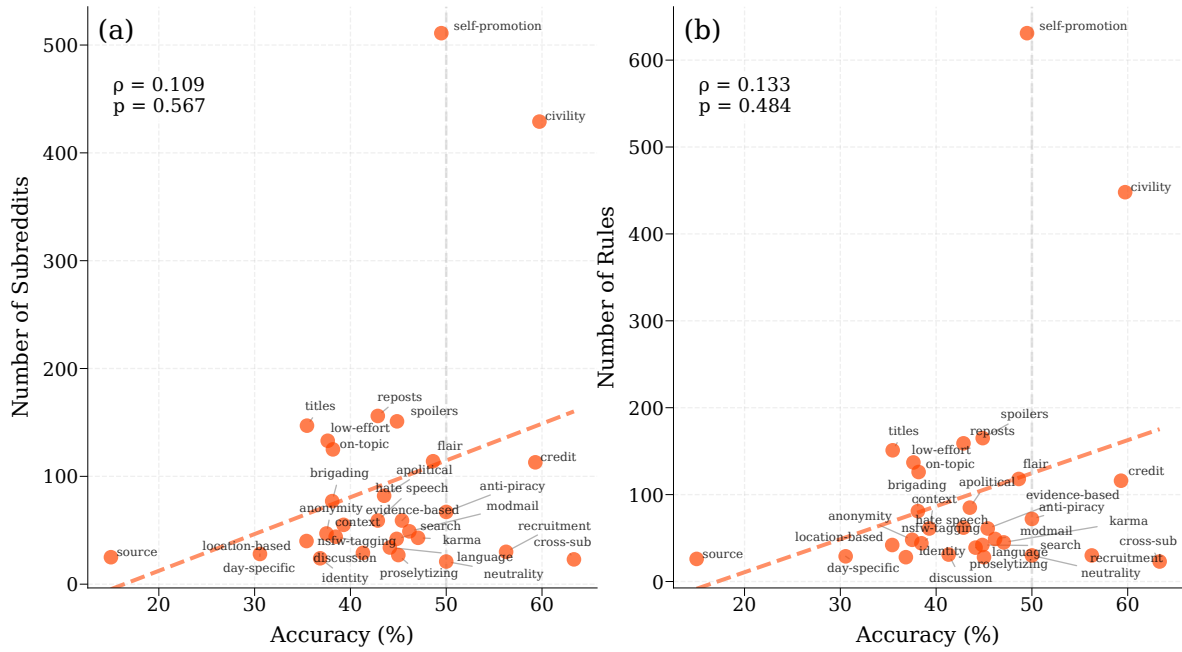


Figure 17: Qwen3-VL-4B (instruct) accuracy with full context compared against (a) the number of subreddits containing rules in each cluster and (b) the number of rules per cluster. Dashed lines indicate the 50% baseline.

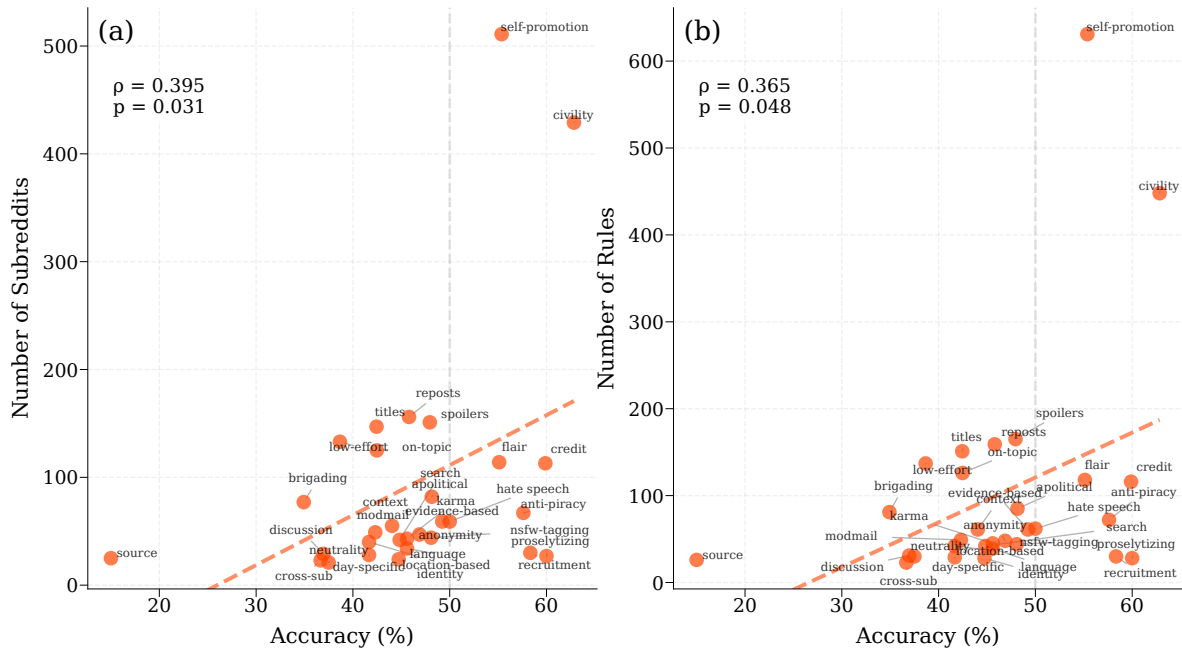


Figure 18: Qwen3-VL-8B (instruct) accuracy with full context compared against (a) the number of subreddits containing rules in each cluster and (b) the number of rules per cluster. Dashed lines indicate the 50% baseline.

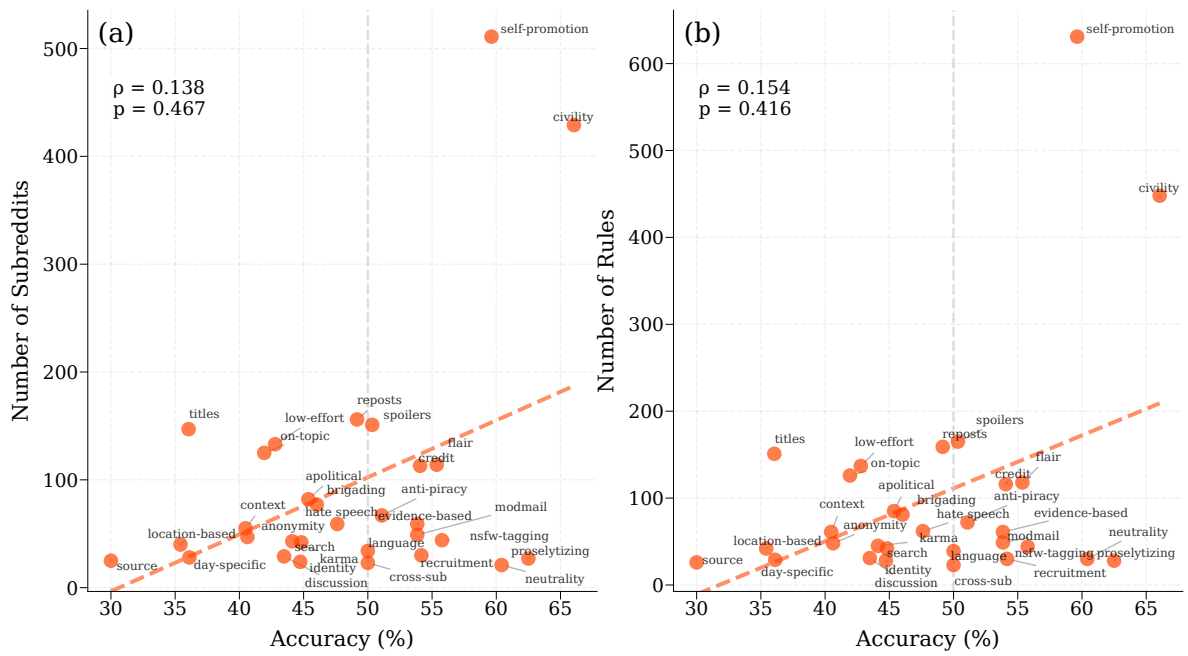


Figure 19: Qwen3-VL-30B (instruct) accuracy with full context compared against (a) the number of subreddits containing rules in each cluster and (b) the number of rules per cluster. Dashed lines indicate the 50% baseline.