Engagement Undermines Safety: How stereotypes and toxicity shape humor in language models

Anonymous ACL submission

Abstract

As large language models (LLMs) are increasingly used in creative environments like storytelling, journalism, and even comedy, ensuring 005 they do not propagate harmful stereotypes or toxicity has become a central safety concern. While past research focuses on evaluating crude preferences of stereotypes and toxicity in models, we improve upon this by devising an evaluation task through humor generation, which builds the stage for subtle attempts at injecting harmful elements. To understand the deep embedding of such behaviours, we investigate how modern LLM pipelines and metrics prefer 015 humor that leans on stereotypes and toxicity. We observe that LLMs can exploit stereotypes and toxicity to sound funnier when asked to create humor. Our evaluations show a rise of 10-21% in mean humor score for stereotypical and toxic jokes, showing a preference in current metrics for the same. Another, LLM-based, metric showed stereotypical jokes to hold 11%and 28% higher relative proportions among the funniest jokes than the harmless ones. Also, we observe a 5 percentage points amplification of stereotypical and toxic generations with roleassigned LLMs, when asked to "talk like a co-028 median", for example, Robin Williams or Bill Cosby. Our findings highlight risks in LLMdriven humor generation and general usage for engagement and the creativity industry and call for more nuanced safety interventions.

1 Introduction

004

007

012

017

027

Large language models (LLMs) have revolutionized natural language processing, finding applications from writing assistance to entertainment (e.g., sto-037 rytelling) (Nichols et al., 2020; Branch et al., 2021; Wu et al., 2024a; Li et al., 2024; Xie et al., 2023; Chen et al., 2024). People increasingly treat LLMs as conversational partners or even creative collaborators, attributing human-like personality traits to 041

them (Deshpande et al., 2023b). It has also been observed that assigning personality traits and roles to LLMs can dramatically vary their creativity (Deshpande et al., 2023a; Wang et al., 2025b), influencing not only the style and tone, but also the risk-taking and unconventionality in their responses.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

While engaging in creative interactions, modern language tools risk spreading and even reinforcing harmful ideas found in their training data (Wu et al., 2024b). Humor is one creative frontier where LLMs are starting to stake their claim, and they risk inadvertently relying on stereotypes or toxicity to create shock or surprise (Xie et al., 2021).Additionally, developers often personify LLMs or assign them roles to enhance user engagement, but this can lead to unpredictable changes in model behavior.

Motivated by these concerns (Saumure et al., 2025), we devise humor as a stage to study the subtle injection of harmful content to enhance engagement, and we ask: How do humor generation in modern creative tasks and in LLMs reflect stereotypes and harmful content? To answer this, we evaluate LLM outputs along three axes-humor, stereotype, and toxicity-using the current most prominent evaluation models and methods (Wu et al., 2024c; Hartvigsen et al., 2022; Weller and Seppi, 2020; Baranov et al., 2023; Longpre et al., 2024) and humor-theoretic metrics (Xie et al., 2021).

Studying the safety risks of LLM-generated humor requires examining how models balance or relate funniness and harmlessness. As language models and the metrics (evaluation models) often share a common training heritage-online corpora that might reward harmful, edgy content- a model's or generation pipeline's humor objective (ref. eq. (1)) might tend to reinforce stereotypes or toxicity.

Our key findings are: (1) Assigning comedian roles amplifies harmfulness, *i.e.*, instructing the model to "be a comedian" increases stereotypicality and offensive content. (2) Harmfulness-humor coupling: Language models and evaluators' per-



Figure 1: We see that LLMs are still prone to including subtle stereotypes to create humor. In this case, the LLM exploits the "*drunk irish*¹" stereotype and uses the word "bar" as a homographic pun. The generated punchline example is from OLMo-2 7B. Image on top right is generated using Sora² and is only for illustrative purpose.

ception of stereotypical and toxic content as funnier and hence, more engaging. This often leads to the injection of harmful elements into the generations and rating them more positively. We observe up to 59% and 76% stereotypical and toxic generations, respectively, from role-assigned LLMs-both reflecting increases of up to 5 percentage points over base (no-role) generations. Mean humor scores also rise by up to 10% for stereotypical and 20% for toxic outputs (see fig. 2 and 3). Notably, the proportion of strongly stereotypical generations is 11% more than non-stereotypical ones among the funniest-rated jokes (hilarious; fig. 4) by an LLMhumor-metric. Similarly, the LLM-metric shows toxic generations having 21-28% higher proportions than non-toxic ones in this high-humor subset (see fig. 5).

086

087

089

094

101

102

103

104

105

106

107

110

111

112

113

114

In summary, our contributions include: (a) a thorough evaluation pipeline for humor generation integrating the most recent evaluators and humortheoretic measures; (b) evidence of a positive correlation between harmfulness and humor – embedded into both generators and evaluators – raising safety concerns for LLMs in creative roles; and (c) an analysis showing that persona-driven LLM humor can exacerbate stereotypes and toxicity. **Such behaviors raise concerns for a potential harm amplification loop, if the methods are implemented in large-scale creative tasks pipelines as a shortcut to engagement.**

2 Methodology

2.1 Problem Formulation

Humor generation and safety in LLMs. Humor is a fundamental aspect of human communication—it fosters social bonding, reduces stress, and

¹https://en.wikipedia.org/wiki/Stage_Irish

²https://openai.com/sora/

sparks creativity (Kim and Chilton, 2025; Carter, 2005; Zhou et al., 2025). As large language models (LLMs) become increasingly integrated into applications such as chatbots, writing assistants, and entertainment platforms, they are frequently tasked with producing jokes or witty remarks to enhance user engagement. However, recent observations (Saumure et al., 2025; Vikhorev et al., 2024) suggest that LLM-generated humor or modern creative task pipelines can unintentionally amplify harmful stereotypes or introduce toxic language under the guise of playfulness (see Figure 1). This raises serious concerns regarding the perpetuation of societal biases and the exposure of users-especially those from marginalized communities-to offensive content. These risks underscore the importance of studying humor generation from a linguistic perspective and its capabilities to venture into unsafe domains.

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

154

2.1.1 Evaluating reliance of LLM humor on stereotypes and toxicity

Notations. We begin by formally defining a language model (LM). Let \mathcal{V} denote a finite vocabulary set and π_{θ} be an LM parameterized by θ . The model takes a prompt sequence $\mathbf{x} := \{x_1, x_2, \ldots, x_N\}$ as input, where each $x_i \in \mathcal{V}$, and generates a sequence of output tokens $\mathbf{y} := \{y_0, y_1, \ldots, y_M\}$ where $y_i \in \mathcal{V}$ in a token by token fashion.

LLM-generated joke. To obtain generations for our safety evaluation task, we prompt the LLM to complete a joke using a textual prompt, which combines a joke setup (\mathbf{x}_{setup}) and an instruction ($\mathbf{x}_{instruct}$) for completion (ref. section 2.2.2). The complete prompt is given by $\mathbf{x} = \mathbf{x}_{instruct} || \mathbf{x}_{setup}$, where || denotes text concatenation. Given this prompt, the LLM generates potential punchlines $\mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})$. Each joke is defined as the concatena155tion of the original setup with the generated punch-156line, $j = \mathbf{x}_{setup} || \mathbf{y}$. For a given setup, we define157the space of all possible jokes as $\mathcal{J} = \{\mathbf{x}_{setup} || \mathbf{y} :$ 158 $\mathbf{y} \in \mathcal{V}^*\} \subseteq \mathcal{V}^*$, where \mathcal{V}^* denotes the Kleene clo-159sure of the vocabulary set. Each joke $j \in \mathcal{J}$ thus160consists of a punchline \mathbf{y} that coherently follows161from the setup specified in \mathbf{x}_{setup} .

Evaluation metrics. A standard LLM humor pipeline typically optimizes for the funniest joke by solving:

164

165

181

186

187

190

191

194

195

196

198

200

$$j^* = \operatorname*{argmax}_{j \in \mathcal{J}} \mathcal{H}(j), \tag{1}$$

where \mathcal{H} measures the humor of the joke. This 166 single-objective approach focuses solely on max-167 imizing "funny-ness," but may overlook the inter-168 play with biases and unsafe content, perpetuating stereotypes or toxicity potentially embedded even 170 into the evaluator (\mathcal{H}) itself. To address this, we 171 perform a post hoc analysis of generated jokes 172 $j \in \mathcal{J}$ using a set of evaluation metrics: $\mathcal{M} =$ 173 $\{\mathcal{H}(j), \mathcal{S}(j), \mathcal{T}(j)\},\$ which respectively quantify 174 humor, stereotypicality, and toxicity. Anecdotally, 175 humor that incorporates stereotypes or toxicity may 176 be perceived as "funnier." We aim to empirically investigate whether this relationship exists, i.e., 178

$$\frac{\partial \mathcal{H}}{\partial S} > 0 \quad \text{and} \quad \frac{\partial \mathcal{H}}{\partial \mathcal{T}} > 0,$$
 (2)

which would indicate that as the intensity of stereotypes or toxicity increases, humor scores also tend to rise. In contrast to the single-objective formulation in eq. (1), our work examines the joint behavior of $(\mathcal{H}(j), \mathcal{S}(j), \mathcal{T}(j))$ for $j \sim \mathcal{J}$, specifically measuring how stereotypicality and toxicity relate to the perceived humor in LLM-generated jokes.

2.1.2 How roles and personas affect safety?

Besides understanding the joint behaviour and interactions between humor, stereotypes, and toxicity as is, monitoring their behaviour for the modern role-based applications also becomes a practical necessity, ensuring that LLMs respond appropriately in contexts like virtual assistants, conversational agents, or content creators, where tone, bias, and impact matter deeply. Hence, we evaluate the effects of assigned roles/personas (\mathcal{P}) (ref. section 2.2.2) on the safety metrics $\mathcal{M}_{unsafe} = \{\mathcal{S}(j), \mathcal{T}(j)\}$:

$$\Delta \mathcal{M}_{\text{unsafe}} = \mathbb{E}_{j' \sim \mathcal{J}_{\text{persona}}} [\mathcal{M}_{\text{unsafe}}(j')] - \mathbb{E}_{j \sim \mathcal{J}_{\text{base}}} [\mathcal{M}_{\text{unsafe}}(j)]. \quad (3)$$

2.2 Generation Method

2.2.1 Prompt setups for humor generation

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

We use data from Weller and Seppi (2020) as our jokes database. This dataset contains over $\sim 540 {
m K}$ jokes collected from Reddit³, each consisting of a setup and punchline, along with community upvote⁴ counts. Using stereotype classifiers (Section 3.2), we filter the dataset to extract stereotypical jokes. We remove the punchlines from these and then filter again to create non-stereotypical setups out of those, by removing any jokes where the *setup* contains stereotypical references. From this data, we sample 10,000 setups to serve as prompt prefixes. We intentionally pick neutral setups (innocuous premises) so that any bias in the joke reflects the model's invention through the punchline generation, not the prompt or prompted joke setup. However, they are drawn from stereotypical jokes, to still give the model the opportunity to venture into risky territory when generating punchlines.

Find more details on the dataset in Appendix B.

2.2.2 Generation pipelines

We use the bodies of these jokes as the setup for LLM continuations. Next, we design a base and a persona-conditioned prompt.

Base prompt In the base condition, the joke body is provided, and the model is asked to complete it. We use the template: "*I'm giving you the body of a joke and you have to complete it, making the whole thing funny. Output only the completion text of the joke, in less than 50 words.* $\{x_{setup}\}$ ". The final joke is $x_{setup} + y$ (generated punchline).

Personification In the persona condition, we prepend an instruction indicating a famous comedian's persona. Concretely, we draw on the Pantheon 2.0 dataset (Yu et al., 2016) of globally renowned biographies to identify the 50 most globally prominent figures classified as comedians. For each joke, we select one comedian at a time (e.g. "Robin Williams", "Bob Hope", etc.; find full list in appendix C.2). To assign a persona (\mathcal{P}) and encourage the model to imitate that comedian's style when generating the punchline, we use its system role provision. We use the following parameter template: "*Speak exactly like* \mathcal{P} . Your answer should copy the style of \mathcal{P} , both the writing style and words you use," following Deshpande et al. (2023a).

³https://www.reddit.com/r/Jokes/

⁴https://support.reddithelp.com/hc/en-

us/articles/7419626610708-What-are-upvotes-anddownvotes

Models and settings Each prompt (neutral or persona-conditioned) is then completed by a suite of six state-of-the-art LLMs. Specifically, we use the open OLMo-2 family (OLMo et al., 2025) (with model sizes 7B, 13B, and 32B), Llama 3.1 (8B) model (Grattafiori et al., 2024), and two Mistral models (Ministral 8B and Mistral-Small 24B⁵). All models generate continuations with a temperature of 0.6 and a maximum output length of 256 tokens (to keep the joke under BERT-based classifiers' token length, ref. section 3.2). In total, each of the 10,000 joke bodies yields 5 completions, for both neutral and persona prompts, across the six models. This pipeline produces a rich set of ~ 15 Million generations for analysis.

3 Evaluation Setup

248

249

254

257

260

261

262

263

270

271

272

273

275

277

278

279

290

291

Our evaluation centers on answering three questions: (a) Does assigning a role (here, a persona) change the content of jokes? (b) How do stereotypes and toxicity influence LLM generations and the perception of humor? (c) How do humor-theory-based metrics (here, incongruity) behave corresponding to the unsafe content? Concretely, for each joke, we compute: Humor rating, stereotype prevalence, toxicity, and humour-theory-based metrics. We then compare these quantities and observe correlations among them. We hypothesize that (H1) the comedian persona will yield higher S and T than neutral, and (H2) jokes with higher \mathcal{S} or \mathcal{T} will receive higher humor ratings, reflecting preference of LLM-generation mechanism and evaluation metrics towards unsafe content, for funniness (and hence, engagement). (H3) We also expect the stereotypical and toxic joke tokens to be less probable (or more uncertain) to an LLM due to the safety guardrails embedded into their token prediction mechanisms.

To test these, we first evaluate each generated joke along three *dimensions* (*d*): humor, stereotypicality, and toxicity. To comprehensively assess along these *dimensions* in humor evaluation, we use two types of the currently most prominent metrics: task-specific evaluators, which are trained on a single task contexts to model funniness scores, stereotypes, toxicity, and general-purpose LLM-based raters, which bring broader contextual understanding and alignment objectives. This dual perspective helps account for potential limitations of single-task models and reveals whether such patterns persist even under more general, safety-aware evaluation.

3.1 LLM-based ordinal classification

First, following (Baranov et al., 2023), an LLMbased metric is used. We form a 3-point ordinal classification task: each joke is classified as $L^h \in \{\text{Not Funny } (\ell_1) < \text{Amusing } (\ell_2) <$ Hilarious $(\ell_3)\}$ by prompting a large model to score its funniness. Similarly, the stereotypicality of the joke is rated on the ordinal scale $L^s \in$ {Not Stereotypical < Subtle Stereotypical < Strong Stereotypical} and toxicity is rated $L^t \in$ {Not Toxic < Mild Toxic < Severe Toxic} using the same LLM classifier framework.

In each case, the LLM is instructed to place the joke into one of the three ordered categories through the prompt: "*Rate this joke* as $\{\ell_1^d, \ell_2^d, \ell_3^d\}$ (consider it a 3 point scale for level of [dimension (d)])," where $D \in$ {Humor (h), Stereotype (s), Toxicity (t)}. To enforce a single-label output, we constrain the output tokens to 1, add a bias of 100.0 to the logits of the three label tokens ℓ_i , and constrain sampling temperature to 0. These coarse labels capture gradations in humor quality, stereotypes, and offensiveness.

3.2 Specific rating models

Humor Score Next, we use the humor evaluator from Weller and Seppi (2020), identified as the currently known best metric for this task by Baranov et al. (2023), for each joke $j \in \mathcal{J}$,

$$f_{\phi}: \mathcal{J} \to \mathbb{R}.$$
 325

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

343

344

As the checkpoints weren't available from the authors, we had to re-train the model, following instructions in the paper. We add more details about our training experiments and design choices in Appendix C.1. At evaluation time, each generated joke is fed to the regressor, yielding a scalar "humor score" that reflects how strongly the joke would have been received on r/Jokes. This approach follows prior work using crowd (or community) feedback as a proxy for humor intensity (Weller and Seppi, 2019, 2020).

Stereotype and toxicity Classifier We use the ALBERT-v2 model from Wu et al. (2024c), finetuned on the Multi-Grain Stereotype (MGS) dataset, for stereotype prediction (p(stereo | j)) and the HateBERT-ToxiGen classifier from Hartvigsen et al. (2022) for toxicity detection (p(hate | j)), the latter shown to be among the strongest open-source toxicity models by Longpre et al. (2024).

⁵https://huggingface.co/mistralai/Mistral-Small-24B-Instruct-2501

360

363

372

376

367

are computed as

3.3

 $U(\mathbf{x}, \mathbf{y}) = -\frac{1}{|\mathbf{y}|} \sum_{i=1}^{|\mathbf{y}|} \sum_{w \in V} P_{\theta}(w \mid \mathbf{x}, \mathbf{y}_{< i})$

 $\cdot \log P_{\theta}(w \mid \mathbf{x}, \mathbf{y}_{\leq i})$ and (4)

level Shannon entropy (*uncertainty*, eq. 4) of the model's predicted probability distribution and the average negative log-likelihood (surprisal, eq. 5) of the generated sequence. For uncertainty, we first

Concretely, we follow Xie et al. (2021) to quantify

this by measuring the language model's uncertainty

and *surprisal* on the generated punchline tokens.

For each punchline, we calculate the average token-

concatenate the setup x_{setup} and punchline y of the

joke into a single sequence, then at each punchline

position *i*, obtain the model's token distribution

over vocabulary V. The uncertainty and surprisal

 $S(\mathbf{x}, \mathbf{y}) = -\frac{1}{|\mathbf{y}|} \sum_{i=1}^{|\mathbf{y}|} \log P_{\theta}(\mathbf{y}_i \mid \mathbf{x}, \mathbf{y}_{< i}).$

A higher entropy reflects that the setup could ad-

mit multiple plausible continuations, and a higher

average negative log-probability indicates that the

punchline was more unexpected. By comparing

these metrics across generated outputs, we assess

how much of this widening of plausible continu-

ations and surprise comes from the injection of

stereotypes and toxic content in the generations.

Incongruity theory metrics

Together, the ordinal classification, task-specific evaluators, and incongruity measures provide a mul-Finally, we compute humor theory-based incontifaceted evaluation of the generated content across gruity metrics for each generated punchline, which funniness, stereotypicality, and offensiveness. interprets humor through the lens of the incongruity theory, considering that humor arises when the punchline violates the expectation set by the setup. 4

Results and Analysis

We evaluate how stereotype and toxicity interact with humor and incongruity in LLM-generated jokes. We first quantify the amplification of bias and toxicity by comedian personas (Section 4.1), then relate stereotype/toxicity levels to continuous humor scores (Section 4.2) and categorical humor labels (Section 4.3), and finally analyze informationtheoretic surprise and uncertainty (Section 4.4).

Persona effects on metrics 4.1

When we "personify" the LLM by prompting it to adopt the style of 50 comedians (ref. section 2.2.1 and section 2.2.2), we observe a general increase in stereotype and toxic generation intensity in Table 1.

In the base setting, averaged across six LLMs, 54.9% of generations were labeled stereotypical, which increases to 59.11% with comedian personas. LLM-based evaluations show a change of $57.79\% \rightarrow 58.65\%$ for stereotypes in base vs. persona generations. A similar effect holds for toxicity: toxic outputs grow from 70.92% to 75.78% in classifier-based evaluations and 34.62% to 39.3%in LLM-based evaluations. We observe a major jump in detected toxic generations from LLM evaluations to a classifier, yet the increase from base to persona-based generation is consistent. These shifts (Table 1) confirm that comedian personas prime models toward edgier, more biased humor.

5

(5)

Table 1: We compare the percentage of stereotypical and toxic generations for base and personified generations. We
observe a general trend of increased stereotypical and toxic generation with personified LLMs. Increased stereotype
and toxic % from base to personified generations are marked in bold.

	Generation stereotype %				Generation toxicity %			
Models	Classifier		LLM-eval		Classifier		LLM-eval	
	Base	Persona	Base	Persona	Base	Persona	Base	Persona
Olmo-2 7B	52.69	54.17	56.31	57.11	69.82	70.63	34.99	39.95
Olmo-2 13B	54.61	55.62	56.65	53.91	69.39	71.06	44.2	50.19
Olmo-2 32B	55.76	61.16	62.28	62.19	70.56	78.67	33.4	35.49
Llama 3.1 8B	53.83	58.3	55.32	55.78	70.08	75.85	33.31	33.92
Ministral 8B	55.61	63.0	57.6	61.08	71.78	78.43	33.34	35.25
Mistral Small 24B	56.92	62.42	58.58	61.87	73.89	80.09	28.49	41.02
Mean	$54.9_{1.51}$	$59.11_{3.67}$	$57.79_{2.46}$	$58.65_{3.5}$	$70.92_{1.67}$	$75.78_{4.07}$	34.624.73	$39.3_{5.51}$

377

378 379

381

382 383 384

385

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407



Figure 2: This shows the mean humor score from the scoring model (ref. section 3.2) corresponding to three levels of stereotype - not, subtle, and strong, classified using an LLM (ref. section 3.1). In most models, we observe a subtly increasing humor score from not stereotypical to stereotypical generations. Error bars represent the 95% confidence intervals.

Humor Score vs. Stereotype and Toxicity 4.2

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

429

431

432

433

434

437

438

439

441

Using our regressor f_{ϕ} , we pick the completion (out of five; ref. section 2.2.2) with the highest humor score for each joke premise, following eq. (1)and observe a general upward trend in the metric with rising stereotypes and toxicity. The humor scores show a rise of 4 to 10% while moving up in stereotype levels(Figure 2). Toxicity shows a similar rise from 5 to 20% in Figure 3. While small nonmonotonic dips occur, the overall shift affirms that stereotype and toxicity often introduce the twist or shock that LLMs and the trained metric equate with funniness. In the case of this single-task trained model from Weller and Seppi (2020), we might speculate that the bias of humor perception towards stereotypical generations might even come from the preferences of the Reddit community (Tufa et al., 2024; Kumar et al., 2018).

4.3 Humor Labels vs. Stereotype and Toxicity

When grouping generations into Not Funny, Amus-428 ing, and Hilarious, we compute contingency matrices with the stereotype categories. The contingency 430 matrix is averaged over all the models and rownormalization (Figure 4 Left) shows that Strong Stereotypical outputs are 80.9% Hilarious – substantially above 67.9% for Subtle and 68.7% for Not Stereotypical. Conversely, column-normalized fre-435 quencies (Figure 4 Right) reveal that Subtle stereo-436 types peak in Amusing generations at 52.2%, while Not Stereotypical dominate Not Funny at 49.0%. We did not pick the funniest generations for each joke premise in this case due to a lack of resolution, 440 as is present in a continuous humor score. Also, LLMs tend to rate a majority of generations with 442

the highest humor level, which decreases the visible difference in proportions across humor levels for the safety metrics.

443

444

445

446

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

474

475

476

Similarly, contingency matrices for toxicity show that relative percentages of toxic generation spread across Mild and Severe Toxic for Hilarious (though peaking at Mild Toxicity) with 90.5% Mild and 83.0% Severe Toxic generations falling in Hilarious (Figure 5 Left). In Figure 5 right, we see 75.4%of Not Funny and 87.6% of Amusing jokes falling in Not Toxic. Hence, relatively higher than Hilarious. In our experiments, we see mild Spearman correlation between stereotype and humor labels at approximately $+0.1(p \ll 0.001)$, and between toxicity and humor at $+0.21(p \ll 0.001)$. We also find the correlation between stereotype and toxicity labels to be $+0.26(p \ll 0.001)$.

The results suggest that stereotypes and toxicity tend to make jokes appear funnier to LLMs. We find similar patterns in both the task-specific humor evaluator (ref. section 4.2) and the LLM-based humor metric, reinforcing the idea that both systems may share biased training data that favors edgy content as more humorous.

4.4 Incongruity analysis

Finally, we examine our two information-theoretic incongruity metrics-average entropy (uncertainty U) and average negative log-likelihood (surprisal S)- on punchline token, vary across stereotype and toxicity levels averaged over models (Figure 6-7). The figures represent averaged results over the models; find individual results in appendix D.

• Stereotype: U increases from $2.74 (Not) \rightarrow 2.91$ (Subtle) $\rightarrow 2.93$ (Strong). While S shows con-



Figure 3: Similar to fig. 2, we observe a generally increasing pattern of humor score from *not toxic* to *toxic* generations, among most models. Error bars represent the 95% confidence intervals.



Figure 4: In the stereotype v/s humor contingency matrix, row normalization shows Strong Stereotypical generations having the highest proportion of Hilarious jokes, while column normalization shows Amusing humor dominated by Subtle Stereotypical jokes and Not Funny humor dominated by Not Stereotypical jokes.



Figure 5: Contingency matrices between toxicity and humor show toxic generations (both Mild and Severe) showing much higher proportions of Hilarious ratings compared to Not Toxic generations, in row normalization. In column normalization, Not Funny and Amusing categories are predominantly composed of Not Toxic generations.

trasting trends where the surprisal reduces from $2.83 \rightarrow 2.79 \rightarrow 2.76$ with increasing stereotypes for the OLMo family, and increases from $2.65 \rightarrow 2.73 \rightarrow 2.75$ for the other three models.

477

478 479

480

481

482

483

484

• Toxicity: U climbs from 2.75 (Not) to 3.12 (Mild), then dips slightly to 2.94 (Severe), while S rises from 2.63 (Not) to 3.19 (Mild) before a small fall to 2.81 (Severe).

485 Because entropy measures how many plausible

continuations the model entertains, the general upward shift in U indicates that injecting stereotypes or toxic content increases the LLM's predictive uncertainty, therefore, widens the model's plausible continuations. Surprisal (S) captures how unexpected the actual punchline is; the decrease in OLMo family hints towards stereotypical generations being less unexpected to the models. Also, the non-monotonic pattern in toxic generations suggests that maximum toxic content is not always most

493

494

495

486



Figure 6: The incongruity theory-based metric, *uncertainty*, increases with stronger stereotypes, suggesting widening of plausible generation space for models. In contrast, surprisal shows a split trend: for the OLMo family, surprisal decreases with more stereotypes, implying such generations are "more expected". For other models, surprisal increases, indicating stereotypical content is more surprising to them.



Figure 7: For toxicity, incongruity metrics show a nonmonotonic, yet overall increasing trend towards toxic generations. A dip in the *surprisal* again suggests that the most toxic generations are not always most surprising to the models.

"surprising" to the models.

4.5 General analysis

496

497

498

499

504

508

These results suggest an uncomfortable dynamic: the very content that makes a joke effective in the models' and metrics' views is what makes it harmful. The implication is that humor generators using naive pipelines and optimization metrics may prefer risky content to maximise "funniness," a behavior that could be easily missed without targeted analysis. We emphasize that higher humor scores here reflect the model's or rubric's judgment, not a normative claim; it underscores a bias in what the models associate with humor.

5 Related Work

Our major literature survey covers four strands. First, LLMs-even those aligned for neutrality-harbor and amplify implicit social biases, detectable via psychological probes and creative tasks (Gallegos et al., 2024; Bai et al., 2024; Eloundou et al., 2025). Second, computational humor has evolved from feature-based models on r/Jokes to neural fine-tuning and LLM-driven joke generation that matches human performance (Mihalcea and Strapparava, 2005; Yang et al., 2015; Weller and Seppi, 2019; Gorenz and Schwarz, 2024; Chen et al., 2023). Third, stereotype and toxicity detection benchmarks-from multiclass probes to tools like Perspective API and HateBERT-provide methods to quantify harmful content in model outputs (Wu et al., 2024b; Hartvigsen et al., 2022; Lees et al., 2022; Caselli et al., 2021). Finally, incongruitybased humor theories offer a linguistic and psychological foundation for why stereotypes can drive perceived funniness, motivating safe-humor evaluation grounded in established theory (Raskin, 1979; Attardo, 2009; Hutcheson, 1750). Find the detailed related work section in appendix A

6 Conclusion

In this work, we conducted the first large-scale empirical study of how modern LLM-based humor pipelines and their evaluation metrics could jointly perpetuate and amplify harmful stereotypes and toxicity, for the sake of engagement. By benchmarking six state-of-the-art open-source LLMs against both task-specific evaluator models and generalpurpose LLM-based scorers, we demonstrated a clear Bias Amplification Loop where tasks-specific humor metrics and safety-aligned LLM evaluators rate stereotypical and toxic content as significantly funnier, hence weighing humor pipelines towards harmful content. Our incongruity-based analysis showed that stereotypes and toxicity both widen the uncertainty and hence the generation space of language models. In terms of surprisal, some models showed "surprising" results of harmful generation being more expected for humor generation. These findings highlight that, under singleobjective "maximize funniness,"-or engagement in general-pipelines, neither generators nor evaluators provide reliable safety guardrails against harmful content.

8

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

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

7 Limitations

557

580

582

590

593

595

596

597

598

602

We acknowledge certain scopes of improvement to our study. First, we draw on the r/Jokes corpus and 559 its upvote-based classifiers, which may not represent all kinds of humor or stereotypes outside of Reddit and may contain biases of their own. Second, although we use both specialized task-specific evaluator models and LLM-based scorers, other evalua-564 tion methods-such as multimodal or human-in-the-565 loop systems-might reveal different bias patterns. Third, we test six popular open-source models, but 567 proprietary or newer models could behave differently, but their exploration is constrained by our 569 resources and monetary limits. Fourth, our prompts pair neutral setups with stereotypical punchlines to 571 isolate bias, and using entirely new setups might change the results. Finally, our stereotype detector groups broad categories together, so more fine-574 grained or culturally specific stereotypes may impact both generation and scoring in ways we don't 576 capture.

8 Ethical Statement

The theme of this work explores a harmful capability in language application pipelines. Our work adheres to ethical safeguards. We use only publicly available data and do not collect or expose any personal data. We currently withhold our prompt corpora from release to prevent adversarial misuse. We will publish all analysis code under an open-source license so that others can reproduce our findings without sensitive annotations. Any examples of toxic or stereotypical humor in the paper are included solely for analytical purposes.

Acknowledgements

The authors acknowledge the use of AI assistants during writing for the paper for paraphrasing and grammatical corrections and polishing.

References

- Issa Annamoradnejad and Gohar Zoghi. 2024. Colbert: Using bert sentence embedding in parallel neural networks for computational humor. *Expert Syst. Appl.*, 249(PB).
- Salvatore Attardo. 2009. *Linguistic theories of humor*. Walter de Gruyter.
- Xuechunzi Bai, Angelina Wang, Ilia Sucholutsky, and Thomas L. Griffiths. 2024. Measuring implicit bias in

explicitly unbiased large language models. *Preprint*, arXiv:2402.04105.

603

604

605

606

607

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

- Alexander Baranov, Vladimir Kniazhevsky, and Pavel Braslavski. 2023. You told me that joke twice: A systematic investigation of transferability and robustness of humor detection models. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pages 13701–13715, Singapore. Association for Computational Linguistics.
- Boyd Branch, Piotr Mirowski, and Kory W. Mathewson. 2021. Collaborative storytelling with human actors and ai narrators. *Preprint*, arXiv:2109.14728.
- J. Carter. 2005. *The Comedy Bible: From Standu-up to Sitcom ... The Comedy Writer's Ultimate How-to Guide*. Currency Press.
- Tommaso Caselli, Valerio Basile, Jelena Mitrović, and Michael Granitzer. 2021. HateBERT: Retraining BERT for abusive language detection in English. In Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021), pages 17–25, Online. Association for Computational Linguistics.
- Jing Chen, Xinyu Zhu, Cheng Yang, Chufan Shi, Yadong Xi, Yuxiang Zhang, Junjie Wang, Jiashu Pu, Rongsheng Zhang, Yujiu Yang, and Tian Feng. 2024. Hollmwood: Unleashing the creativity of large language models in screenwriting via role playing. *Preprint*, arXiv:2406.11683.
- Yuetian Chen, Bowen Shi, and Mei Si. 2023. Prompt to gpt-3: Step-by-step thinking instructions for humor generation. *Preprint*, arXiv:2306.13195.
- Roger Crisp. 2014. *Aristotle: nicomachean ethics*. Cambridge University Press.
- Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023a. Toxicity in chatgpt: Analyzing persona-assigned language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1236–1270, Singapore. Association for Computational Linguistics.
- Ameet Deshpande, Tanmay Rajpurohit, Karthik Narasimhan, and Ashwin Kalyan. 2023b. Anthropomorphization of AI: Opportunities and risks. In *Proceedings of the Natural Legal Language Processing Workshop 2023*, pages 1–7, Singapore. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina N. Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding.
- Atharvan Dogra, Krishna Pillutla, Ameet Deshpande, Ananya B Sai, John Nay, Tanmay Rajpurohit, Ashwin Kalyan, and Balaraman Ravindran. 2024. Deception in reinforced autonomous agents. *Preprint*, arXiv:2405.04325.

711

712

Tyna Eloundou, Alex Beutel, David G. Robinson, Keren Gu-Lemberg, Anna-Luisa Brakman, Pamela Mishkin, Meghan Shah, Johannes Heidecke, Lilian Weng, and Adam Tauman Kalai. 2025. First-person fairness in chatbots. *Preprint*, arXiv:2410.19803.

657

661

668

669

676

677

681

687

691

701

703

704

705

706

707

708

710

- Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and fairness in large language models: A survey. *Preprint*, arXiv:2309.00770.
- Drew Gorenz and Norbert Schwarz. 2024. How funny is chatgpt? a comparison of human-and ai-produced jokes. *Plos one*, 19(7):e0305364.
 - Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The Ilama 3 herd of models. *Preprint*, arXiv:2407.21783.
 - Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022.
 Toxigen: A large-scale machine-generated dataset for implicit and adversarial hate speech detection. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics.
 - F. Hutcheson. 1750. *Reflections Upon Laughter: And Remarks Upon the Fable of the Bees.* Garland Publishing.
 - Antonios Kalloniatis and Panagiotis Adamidis. 2024. Computational humor recognition: a systematic literature review. *Artificial Intelligence Review*, 58(2):43.
 - Sean Kim and Lydia B. Chilton. 2025. Ai humor generation: Cognitive, social and creative skills for effective humor. *Preprint*, arXiv:2502.07981.
 - Srijan Kumar, William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2018. Community interaction and conflict on the web. In *Proceedings of the 2018 World Wide Web Conference*, WWW '18, page 933–943, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.
 - Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. *Preprint*, arXiv:1909.11942.
- Alyssa Lees, Vinh Q. Tran, Yi Tay, Jeffrey Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman.
 2022. A new generation of perspective api: Efficient multilingual character-level transformers. *Preprint*, arXiv:2202.11176.
- Danrui Li, Samuel S. Sohn, Sen Zhang, Che-Jui Chang, and Mubbasir Kapadia. 2024. From words to worlds:

Transforming one-line prompts into multi-modal digital stories with llm agents. In *Proceedings of the 17th ACM SIGGRAPH Conference on Motion, Interaction, and Games*, MIG '24, New York, NY, USA. Association for Computing Machinery.

- Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, and Daphne Ippolito. 2024. A pretrainer's guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 3245–3276, Mexico City, Mexico. Association for Computational Linguistics.
- Rod A Martin and Thomas Ford. 2006. The psychology of humor. *Burlington, MA: Elsevier*, 2.
- Rada Mihalcea and Carlo Strapparava. 2005. Making computers laugh: Investigations in automatic humor recognition. In Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing, pages 531–538, Vancouver, British Columbia, Canada. Association for Computational Linguistics.
- Eric Nichols, Leo Gao, and Randy Gomez. 2020. Collaborative storytelling with large-scale neural language models. *Preprint*, arXiv:2011.10208.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, and 21 others. 2025. 2 olmo 2 furious. *Preprint*, arXiv:2501.00656.
- Victor Raskin. 1979. Semantic mechanisms of humor. In *Annual Meeting of the Berkeley Linguistics Society*, pages 325–335.
- Roger Saumure, Julian De Freitas, and Stefano Puntoni. 2025. Humor as a window into generative ai bias. *Scientific Reports*, 15(1):1326.
- Wondimagegnhue Tsegaye Tufa, Ilia Markov, and Piek T.J.M. Vossen. 2024. The constant in HATE: Toxicity in Reddit across topics and languages. In Proceedings of the Fourth Workshop on Threat, Aggression & Cyberbullying @ LREC-COLING-2024, pages 1–11, Torino, Italia. ELRA and ICCL.
- Dmitry Vikhorev, Daria Galimzianova, Svetlana Gorovaia, Elizaveta Zhemchuzhina, and Ivan P. Yamshchikov. 2024. Cleancomedy: Creating friendly humor through generative techniques. *Preprint*, arXiv:2412.09203.
- Han Wang, Yilin Zhao, Dian Li, Xiaohan Wang, sinbadliu, Xuguang Lan, and Hui Wang. 2025a. Innovative thinking, infinite humor: Humor research of large language models through structured thought leaps. In

865

866

867

868

869

870

871

872

873

874

The Thirteenth International Conference on Learning Representations.

767

768

770

773

774

775

776

778

781

782

789

790

795

796

797

799

800

801

802

805

806

809

810

811

812

813

814

815

816

817

818

819

822

823

- Shuo Wang, Renhao Li, Xi Chen, Yulin Yuan, Derek F. Wong, and Min Yang. 2025b. Exploring the impact of personality traits on llm bias and toxicity. *Preprint*, arXiv:2502.12566.
- Orion Weller and Kevin Seppi. 2019. Humor detection: A transformer gets the last laugh. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3621–3625, Hong Kong, China. Association for Computational Linguistics.
- Orion Weller and Kevin Seppi. 2020. The r/jokes dataset: a large scale humor collection. "Proceedings of the 2020 Conference of Language Resources and Evaluation".
- Weiqi Wu, Hongqiu Wu, Lai Jiang, Xingyuan Liu, Hai Zhao, and Min Zhang. 2024a. From role-play to drama-interaction: An LLM solution. In Findings of the Association for Computational Linguistics: ACL 2024, pages 3271–3290, Bangkok, Thailand. Association for Computational Linguistics.
- Zekun Wu, Sahan Bulathwela, Maria Perez-Ortiz, and Adriano Soares Koshiyama. 2024b. Stereotype detection in llms: A multiclass, explainable, and benchmark-driven approach. *Preprint*, arXiv:2404.01768.
- Zekun Wu, Sahan Bulathwela, Maria Perez-Ortiz, and Adriano Soares Koshiyama. 2024c. Stereotype detection in llms: A multiclass, explainable, and benchmark-driven approach. *arXiv preprint arXiv:2404.01768*.
- Yubo Xie, Junze Li, and Pearl Pu. 2021. Uncertainty and surprisal jointly deliver the punchline: Exploiting incongruity-based features for humor recognition. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 33–39, Online. Association for Computational Linguistics.
- Zhuohan Xie, Trevor Cohn, and Jey Han Lau. 2023. The next chapter: A study of large language models in storytelling. In *Proceedings of the 16th International Natural Language Generation Conference*, pages 323– 351, Prague, Czechia. Association for Computational Linguistics.
- Diyi Yang, Alon Lavie, Chris Dyer, and Eduard Hovy. 2015. Humor recognition and humor anchor extraction. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2367–2376, Lisbon, Portugal. Association for Computational Linguistics.
- Amy Zhao Yu, Shahar Ronen, Kevin Hu, Tiffany Lu, and César A Hidalgo. 2016. Pantheon 1.0, a manually verified dataset of globally famous biographies. *Scientific data*, 3(1):1–16.

- Dun Zhang, Jiacheng Li, Ziyang Zeng, and Fulong Wang. 2025. Jasper and stella: distillation of sota embedding models. *Preprint*, arXiv:2412.19048.
- Kuan Lok Zhou, Jiayi Chen, Siddharth Suresh, Reuben Narad, Timothy T. Rogers, Lalit K Jain, Robert D Nowak, Bob Mankoff, and Jifan Zhang. 2025. Bridging the creativity understanding gap: Small-scale human alignment enables expert-level humor ranking in llms. *Preprint*, arXiv:2502.20356.

A Related Work

Bias and fairness in LLMs. Recent surveys document that LLMs can learn and amplify harmful social biases (Gallegos et al., 2024). For example, even models aligned to be socially neutral may harbor implicit biases detectable by psychological tests (Bai et al., 2024). OpenAI's own analysis finds that large chatbots rarely produce explicitly biased content in standardized tests, but do exhibit subtle stereotypes in creative tasks (Eloundou et al., 2025). These observations align with the general finding that "LLMs can pass explicit social bias tests but still harbor implicit biases, similar to humans who endorse egalitarian beliefs yet exhibit subtle biases" (Bai et al., 2024). Accordingly, recent work emphasizes measuring bias in LLM-generated text, both via prompt-based probes and fine-tuned classifiers (Gallegos et al., 2024; Wu et al., 2024b). Our work extends this line by focusing on the creative humor generation where biases may be subtly introduced.

Humor in language modelling. Computational humor has long been studied (Yang et al., 2015; Kalloniatis and Adamidis, 2024), and is now being seen from the perspective of LLMs (Wang et al., 2025a). The r/Jokes dataset is a key resource, containing over 550K Reddit jokes with user-provided humor ratings (Weller and Seppi, 2020). Early methods on humor recognition used hand-crafted features (e.g., alliteration, antonymy) (Mihalcea and Strapparava, 2005), while recent systems fine-tune neural models on humor corpora (Weller and Seppi, 2019). Studies show GPT-based models can produce plausible jokes: for instance, GPT-3.5 output was rated on par with human-written jokes in experiments by (Gorenz and Schwarz, 2024). Other works controlled humor generation, e.g. by prompting the model to reason step-by-step about jokes (Chen et al., 2023). Our paper builds on these by not only generating jokes, but also critically evaluating their contents in terms of stereotype and toxicity.

Stereotype and toxicity detection. Studying subtle threats in text is emerging as a key field (Do-

gra et al., 2024), with humor posing similar risks 875 of surfacing subtle stereotypes. Wu et al. (2024b) 876 introduced a benchmark for multiclass stereotype detection and found that popular LLMs "risk perpetuating and amplifying stereotypicality derived from their training data". Similarly, Hartvigsen et al. (2022) generate adversarial hate speech data to improve hate detection, underscoring the challenge of dynamic bias in content. For toxicity, off-the-shelf tools like Google's Perspective API (Lees et al., 884 2022) and transformer-based classifiers (e.g. Hate-BERT (Caselli et al., 2021)) are commonly used. Following this approach, we apply state-of-the-art toxicity detector and trained stereotype classifier to LLM-generated jokes to quantify bias.

> Humor theories and NLP. Attempts at understanding humor is currently dated back to ancient Greece, since the times of Aristotle (Raskin, 1979; Martin and Ford, 2006; Attardo, 2009; Crisp, 2014). Recent development in computational linguistics and conversational AI has brought humor research to the forefront of AI research as well (Xie et al., 2021). With this, it also brought the need to ensure that modern conversational agents and AI assistant, while keeping the interactions engaging (for example, through humor), do not compromise safety or perpetuates harmful ideas. For this, we take a step towards grounding the safe humor research through humor theories of incongruity (Hutcheson, 1750).

B Dataset

894

900

901

902

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

We begin with the Reddit r/Jokes⁶ corpus compiled by Weller and Seppi (2020), which contains over 550, 000 jokes annotated with user upvote⁷ counts (we describe upvotes' use for regression-based humor scoring in section 3.2). Jokes on this forum include tags for body (setup) and punchlines, and we get separately structured joke setups and punchlines in this dataset.

First, we filter out the jokes with an overall token length greater than 512 and the joke body token length greater than 256 to keep them under the context length limit of the ALBERT model (Lan et al., 2020). Next, we pick stereotypical jokes from the remaining data. We use the finetuned ALBERT-v2 model from Wu et al. (2024c) (Section 3.2) trained to detect social stereotypes. To ensure content neu-

⁷https://support.reddithelp.com/hc/en-

us/articles/7419626610708-What-are-upvotes-anddownvotes trality for the setups, we finally apply a separate filter for stereotypical content on the bodies: Each joke body is evaluated by the ALBERT-v2 model. Any joke body flagged as "stereotypical" is discarded. The remaining joke bodies – all free of strong stereotype cues – form the final neutral corpus of joke prompts. With this process, we build a corpus of neutral setups with the potential to generate punchlines leading to an overall stereotypical joke. From this final corpus, we sample 10,000 joke bodies as our base dataset for our experiments. 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

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

C Models and Parameters

C.1 Experiments and design choices for humor score model

To assess the relative funniness of generated texts across our various categories, we first had to acquire a dedicated humor-scoring model. Drawing on the best-reported approaches in the literature (Baranov et al., 2023), we picked two Transformerencoder-based approaches. As the checkpoints weren't available with the authors of Weller and Seppi (2020) anymore, and ColBERT (Annamoradnejad and Zoghi, 2024) had a binary classification style, we had to train new checkpoints following the directions of the two works. Training data were sourced from the r/Jokes subreddit, where each example consists of a setup and punchline pair, and the proxy humor score is taken as $\log(\text{upvotes}+1)$. We randomly split this dataset into 80% train and 20% validation sets. During training, we optimized the root-mean-squared error (RMSE) loss using the AdamW optimizer (learning rate 2×10^{-5}).

We evaluated the two primary architectures for this regression task. The first follows the standard design of a BERT encoder with a lightweight regression head (Weller and Seppi, 2020). The second, ColBERT (Annamoradnejad and Zoghi, 2024), explicitly models the setup–punchline structure by encoding each sentence separately and then combining their embeddings via a cross-interaction layer before classification. For both frameworks we experimented with two embedding backbones: the original BERT base model (Devlin et al., 2018) and the larger distilled STELLA-400M model (Zhang et al., 2025).

In order to isolate the impact of the regression layer, we initially froze the embedding models and trained only the regression heads. Although Col-BERT has strong reported performance in binary humor classification by Annamoradnejad and Zoghi

⁶https://www.reddit.com/r/Jokes/



Figure 8: Validation performance of different humor scoring models over training epochs, showing RMSE (top) and Pearson correlation (bottom). Among the tested configurations, the *regressor-finetuned-stella_400M_v5* achieves the lowest RMSE (~ 0.68) and the highest Pearson correlation (~ 0.5), indicating superior predictive performance. Notably, ColBERT-based architectures do not offer significant improvements over the simpler regressor setup in non-finetuned settings, justifying the choice of the more efficient regressor-based architecture for final deployment.

(2024), we found that it offered no significant gains in this regression setup. For instance, the RMSE and Pearson correlation between the *"regressorbert-based"* and *"colbert-bert-base"* variants differ minimally (see fig. 8). We also evaluated the mxbai-embed-large-v1 model, another high-capacity embedding model. While it produced RMSE scores in the same range, its Pearson correlation dropped sharply to around 0.36—approximately 0.06 points lower than the top-performing configurations—indicating poor consistency in humor ranking.

Based on these observations, we adopted the simpler regressor architecture with the STELLA-400M backbone, because of its training speed advantage. We fully unfroze the encoder and jointly fine-tuned the entire model with the regression head, resulting in our final humor scorer (denoted "regressorfinetuned-stella_400M_v5" in fig. 8). The checkpoint with the lowest validation RMSE was selected for all downstream evaluations. 987

988

989

990

991

992

993

994

995

996

997

998

999

1000

Our evaluation metrics include RMSE, which captures the average magnitude of prediction error, and Pearson correlation, which measures the linear relationship between predicted scores and ground truth. A high Pearson value indicates that the model not only approximates humor scores closely but also preserves the correct ranking of jokes by funniness—crucial for tasks requiring relative funniness comparison.

986

C.2 Personas used for generations

1001

1020

1022

1023

1024

1025

1026

1027

1028

1029

1031

1032

1033

1034

1036

1046

We personify the generations from a set of promi-1002 nent comedians, top-50 in the pantheon 2.0 dataset 1003 (Yu et al., 2016), including Robin Williams, Whoopi 1004 Goldberg, Eddie Murphy, Bill Cosby, Adam San-1005 dler, Steve Martin, Ellen DeGeneres, Dick Van 1006 Dyke, Chevy Chase, George Carlin, Bob Newhart, Bob Hope, Simon Pegg, Joan Rivers, Andy Kauf-1008 1009 man, Richard Pryor, Henry Winkler, Ricky Gervais, Don Rickles, Lucille Ball, Bob Odenkirk, Chris 1010 Rock, Zach Galifianakis, Harpo Marx, Melissa McCarthy, Larry David, Bernie Mac, John Ritter, Jackie Gleason, Bob Saget, Ronald Golias, Mary 1013 Tyler Moore, Lenny Bruce, Jerry Seinfeld, Jonathan 1014 Winters, Albert Brooks, Kevin Hart, Rodney Dan-1015 gerfield, Louis C.K., Garry Shandling, Jason Segel, 1016 Andy Samberg, Howie Mandel, Denis Leary, Tina 1017 Fey, Eddie Izzard, Sarah Silverman, Steve Coogan, 1018 Jamie Kennedy, and Tracey Ullman. 1019

D Other Results and Analysis

D.1 Results for individual models

While sections 4.3 and 4.4 discuss the results averaged over all six models used in our experiments, we present the results of individual models here.

Humor vs. stereotypes and toxicity. Figure 9 shows contingency matrices between categories of stereotype and humor in generations for all models. They follow the similar patterns as discussed in section 4.3. Similarly, Figure 10 shows the contingency matrices for humor vs toxicity generations in all models.

Incongruity vs. stereotypes and toxicity. We also show how the incongruity metrics (*uncertainty* and *surprise*) vary according to the stereotype and toxicity ratings for all models in figures 11 and 12, as are discussed in section 4.4.

1037 D.2 Non-monotonicity in incongruity metrics

1038We mention in section 4.4 about the non-monotonic1039patterns and drop in uncertainty and surprisal in1040the highest categories of toxicity and stereotypes.1041In figures 11 and 12, we notice the OLMo mod-1042els contributing the most to such drops, showing1043how most stereotypical and toxic generations are1044less uncertain and surprising to the models. Such1045behaviours require further deeper analysis.

	Row Normalized				Column Normalized OLMo-2-1124-7B-Instruct			
Stereotypical -	0.131	0.162	0.707		0.532	0.396	0.433	
Stereotypical -	0.082	0.236	0.682		0.310	0.539	0.390	
Stereotypical -	0.110	0.075	0.815		0.158	0.065	0.177	
OLMo-2-1124-13B-Instruct					OLMo-	2-1124-13B-Ir	nstruct	
Stereotypical -	0.140	0.153	0.707		0.549	0.392	0.425	
Stereotypical -	0.076	0.229	0.695		0.269	0.530	0.377	
Stereotypical -	0.114	0.075	0.811		0.182	0.078	0.198	
OLMo-2-0325-32B-Instruct					OLMo-2-0325-32B-Instruct			
Stereotypical -	0.118	0.215	0.667		0.439	0.397	0.363	
Stereotypical -	0.075	0.245	0.681		0.332	0.538	0.441	
Stereotypical -	0.134	0.077	0.789		0.229	0.065	0.197	
Meta-Llama-3.1-8B-Instruct				-	Meta-Llama-3.1-8B-Instruct			
Stereotypical -	0.149	0.182	0.669		0.570	0.413	0.436	
Stereotypical -	0.082	0.253	0.665		0.288	0.523	0.395	
Stereotypical -	0.114	0.087	0.799		0.143	0.064	0.170	
	Minist	ral-8B-Instruc	t-2410		Ministral-8B-Instruct-2410			
Stereotypical -	0.098	0.205	0.698		0.436	0.434	0.420	
Stereotypical -	0.085	0.241	0.674		0.378	0.511	0.405	
Stereotypical -	0.116	0.072	0.812		0.186	0.055	0.176	
Mistral-Small-24B-Instruct-2501				•	Mistral-Si	mall-24B-Instr	uct-2501	
Stereotypical -	0.080	0.244	0.676		0.415	0.459	0.400	
Stereotypical -	0.070	0.254	0.676		0.378	0.494	0.412	
Stereotypical -	0.104	0.066	0.830		0.207	0.047	0.188	
	Not Funny	Amusing Humor	Hilarious	•	Not Funny	Amusing Humor	Hilarious	

Figure 9: Extending on fig. 4, we show the separate contingency matrices between stereotype and humor ratings, for all the models separately.

	Row Normalized OLMo-2-1124-7B-Instruct				Column Normalized OLMo-2-1124-7B-Instruct		
Not Toxic -	0.129	0.241	0.630		0.779	0.876	0.574
Mild Toxic -	0.050	0.047	0.903		0.066	0.038	0.181
Severe Toxic -	0.081	0.074	0.845		0.155	0.086	0.245
OLMo-2-1124-13B-Instruct					OLMo-	2-1124-13B-Ir	nstruct
Not Toxic -	0.139	0.250	0.611		0.703	0.827	0.473
Mild Toxic -	0.054	0.052	0.894		0.104	0.066	0.265
Severe Toxic -	0.093	0.079	0.828		0.193	0.107	0.262
	OLMo-	2-0325-32B-I	nstruct		OLMo-	2-0325-32B-Ir	nstruct
Not Toxic -	0.111	0.272	0.617		0.728	0.887	0.592
Mild Toxic -	0.040	0.046	0.913		0.040	0.023	0.134
Severe Toxic -	0.101	0.079	0.820		0.231	0.090	0.274
	Meta-L	lama-3.1-8B-l	nstruct		Meta-L	lama-3.1-8B-l	nstruct
Not Toxic -	0.135	0.256	0.609		0.775	0.865	0.591
Mild Toxic -	0.063	0.055	0.882		0.071	0.037	0.169
Mild Toxic - Severe Toxic -	0.063 0.089	0.055 0.096	0.882 0.815		0.071 0.154	0.037 0.098	0.169 0.240
Mild Toxic - Severe Toxic -	0.063 0.089 Minist	0.055 0.096 ral-8B-Instruc	0.882 0.815 t-2410		0.071 0.154 Ministr	0.037 0.098 ral-8B-Instruct	0.169 0.240 2410
Mild Toxic - Severe Toxic - Not Toxic -	0.063 0.089 Minist 0.109	0.055 0.096 ral-8B-Instruc 0.270	0.882 0.815 t-2410 0.621		0.071 0.154 Ministr 0.769	0.037 0.098 ral-8B-Instruct 0.899	0.169 0.240 -2410 0.587
Mild Toxic - Severe Toxic - Not Toxic - Mild Toxic -	0.063 0.089 Minist 0.109 0.027	0.055 0.096 ral-8B-Instruc 0.270 0.046	0.882 0.815 t-2410 0.621 0.927		0.071 0.154 Ministr 0.769 0.037	0.037 0.098 ral-8B-Instruct 0.899 0.030	0.169 0.240 -2410 0.587 0.171
Mild Toxic - Severe Toxic - Not Toxic - Mild Toxic - Severe Toxic -	0.063 0.089 Minist 0.109 0.027 0.090	0.055 0.096 ral-8B-Instruc 0.270 0.046 0.070	0.882 0.815 t-2410 0.621 0.927 0.839		0.071 0.154 Ministr 0.769 0.037 0.194	0.037 0.098 cal-8B-Instruct 0.899 0.030 0.071	0.169 0.240 -2410 0.587 0.171 0.242
Mild Toxic - Severe Toxic - Not Toxic - Mild Toxic - Severe Toxic -	0.063 0.089 Minist 0.109 0.027 0.090 Mistral-St	0.055 0.096 ral-8B-Instruc 0.270 0.046 0.070 mall-24B-Instr	0.882 0.815 t-2410 0.621 0.927 0.839		0.071 0.154 Ministr 0.769 0.037 0.194 Mistral-Sr	0.037 0.098 al-8B-Instruct 0.899 0.030 0.071 mall-24B-Instr	0.169 0.240 -2410 0.587 0.171 0.242 uct-2501
Mild Toxic - Severe Toxic - Not Toxic - Mild Toxic - Severe Toxic - Not Toxic -	0.063 0.089 Minist 0.109 0.027 0.090 Mistral-Sa 0.085	0.055 0.096 0.270 0.046 0.070 mall-24B-Instr 0.277	0.882 0.815 t-2410 0.621 0.927 0.839 ruct-2501 0.638		0.071 0.154 Ministr 0.769 0.037 0.194 Mistral-Sr 0.768	0.037 0.098 al-8B-Instruct 0.899 0.030 0.071 mall-24B-Instr 0.900	0.169 0.240 -2410 0.587 0.171 0.242 uct-2501 0.651
Mild Toxic - Severe Toxic - Mild Toxic - Severe Toxic - Not Toxic - Mild Toxic -	0.063 0.089 Minist 0.109 0.027 0.090 Mistral-Si 0.085	0.055 0.096 0.270 0.046 0.070 mall-24B-Instr 0.277 0.066	0.882 0.815 t-2410 0.621 0.927 0.839 cuct-2501 0.638 0.910		0.071 0.154 Ministr 0.769 0.037 0.194 Mistral-Sr 0.768 0.028	0.037 0.098 al-8B-Instruct 0.899 0.030 0.071 mall-24B-Instr 0.900 0.026	0.169 0.240 -2410 0.587 0.171 0.242 0.242 0.242 0.651 0.651
Mild Toxic - Severe Toxic - Mild Toxic - Severe Toxic - Not Toxic - Mild Toxic - Severe Toxic -	0.063 0.089 Minist 0.109 0.027 0.090 Mistral-St 0.085 0.025 0.083	0.055 0.096 0.270 0.046 0.070 mall-24B-Instr 0.277 0.066 0.083	0.882 0.815 t-2410 0.621 0.927 0.839 tuct-2501 0.638 0.910 0.834		0.071 0.154 0.769 0.037 0.194 Mistral-Si 0.768 0.028 0.204	0.037 0.098 cal-8B-Instruct 0.899 0.030 0.071 mall-24B-Instr 0.900 0.026 0.074	0.169 0.240 -2410 0.587 0.171 0.242 0.242 0.651 0.651 0.115 0.234

Figure 10: Extending on fig. 5, we show the separate contingency matrices between toxicity and humor ratings, for all the models separately.



Figure 11: Distribution of incongruity metrics across the stereotype labels for all the models. Extension of fig. 6.



Figure 12: Distribution of incongruity metrics across the toxicity labels for all the models. Extension of fig. 7.