# Large Language Models Threaten Language's Epistemic and Communicative Foundations

Anonymous ACL submission

### Abstract

Today, Large language models (LLMs) are reshaping the norms of human communication, sometimes decoupling words from genuine human thought. This transformation is deep, and undermines the trust and interpretive norms that were historically tied to authorship. We draw from linguistic philosophy and AI ethics to detail how large-scale text generation can induce semantic drift, erode accountability, and obfuscate intent and authorship. Our work here introduces conceptual frameworks including hybrid authorship graphs (modeling humans, LLMs, and texts in a provenance network), epistemic doppelgängers (LLM-generated texts that are indistinguishable from human-authored texts), and authorship entropy. We explore mechanisms such as "proof-of-interaction" authorship verification and educational reforms to restore confidence in language. While LLMs' benefits are undeniable (broader access, increased fluency, automation, etc.), the upheavals they introduce to the linguistic landscape demand reckoning. This paper provides a conceptual lens to chart these changes.

# 1 Introduction

011

014

022

024

040

043

"Last year's words belong to last year's language And next year's words await another voice"

Language has been the keystone of our communication and thought for thousands of years. From ancient cuneiform tablets to modern digital platforms, people have relied on written language as a store of facts, beliefs, and ideas. Underlying this tradition is a widespread assumption: that any text reflects a human mind, shaped by cognitive processes and linked to specific authors. Thus, language has been an expression of human intention, demanding both attention from the reader, and accountability from the author (Winograd, 1972; Bender and Koller, 2020). Over our long evolutionary trajectory, these assumptions have steadily held true, and are now woven into the very fabric of how we understand and interpret language.



Figure 1: LLMs introduce a shift in communicative dynamics. Traditionally, human-to-human communication directly conveys intentional thought from speaker to listener (top). But when mediated by LLMs, language can lose direct intentional grounding, resulting in messages disconnected from the speaker's original intent and confusing the listener (bottom).

But the swift ascent of large language models (LLMs) over the past five years has begun to fundamentally reconfigure this relationship. Trained on internet-scale corpora and with representational flexibility from billions of parameters (e.g., GPT (Brown et al., 2020), PaLM (Chowdhery et al., 2022), LLaMA (Touvron et al., 2023)) or DeepSeek (DeepSeek-AI, 2024), these can generate well-polished and coherent text with minimal guidance. They can replicate stylistic nuances, rhetoric, and emotional tones (Schick et al., 2021; Solaiman et al., 2019), which were attributed solely to human creativity till very recently.

On the one hand, these capabilities are surprising

and extraordinary, suggesting potential for a new cognitive revolution in AI. On the other, this transformation weakens the link tying text to a human mind. Teachers worry that student essays reflect an LLM's fluency rather than the messy traces of a student's thoughts (Cotton et al., 2023; Zhou et al., 2023). Researchers question whether a paper reflects a scholar's insights or a model's reassembly of existing content (Zellers et al., 2019). Even everyday exchanges: emails, tweets, and blogs might stem from digital processes rather than human voices (Duarte et al., 2022; Weidinger et al., 2021). In this sense, it would be ironic if LLMs, instead of illuminating and edifying human language communication, lead to its devolution.

058

060

063

064

066

067

073

079

086

090

092

098

100

101

102

103

104

105

The broader NLP community, as creators of LLM technologies, bears a direct responsibility for their societal implications (Gabriel, 2020; Bender et al., 2021). By ignoring the epistemic and ethical consequences of these systems, we risk having text lose its role as an indicator of human intent and thought (Floridi and Chiriatti, 2020; Raji et al., 2022). This can fundamentally degrade how we conduct discourse, value expertise, and maintain trust (O'Neil, 2016; Zuboff, 2019). This paper grapples with this tension between positive applications of LLMs (Team, 2022; Hutchinson et al., 2023; Miller, 2019) and the challenge LLMs pose to language's role in human thought. Section 2 explores language's philosophical aspects like intentionality and authorship. Section 3 identifies two key issues stemming from widespread LLM text, namely semantic drift and erosion in trust. In Section 4, we describe possible approaches to reestablish accountability and measure semantic drift using methods from NLP, cryptography, and HCI. Section 5 looks at social implications and possible responses in education, scholarship, etc. Section 6 argues that approaches like watermarking do not address core issues. Section 7 explores rethinking ideas about language and authorship with ideas on human-AI collaboration, educational changes and human-only publishing spaces. We briefly summarize alternative perspectives and arguments in Section 8. We conclude with a reflection on the need to maintain language's cognitive and epistemic roles in the future.

# 2 Language, Authorship, and Meaning

Language's function has been debated for centuries,from Plato's dialogues on rhetoric to modern ana-

lytic philosophy (Plato, 1997; Wittgenstein, 1953). While language is commonly viewed as an information channel for transmitting information, linguists argue that language is a *communal sense-making act*. It has been closely linked to intention, context, and the ability to hold speakers accountable (Searle, 1969; Austin, 1975; Floridi, 2013). 108

109

110

111

112

113

114

115

116

117

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

153

154

155

156

157

158

Language as an Intentional Act: Searle's speech-act theory (Searle, 1969) and Austin's work on performativity (Austin, 1975) argue that language is not just a conduit for information transfer, but also enacts intentions. To say something is often to do something: to promise, question, declare, for example. The force of an utterance depends on the speaker's agency and recognition of those intentions by a listener (Grice, 1975). This has been fundamental to authorship, particularly in academic and legal discourse, where a text is an intellectual act tied to its creator's identity and responsibility (Dworkin, 1996). Even when ghostwriters were traditionally involved, the text ultimately reflected a coherent cognitive source (Foucault, 1984; Chartier, 1994; Sperber and Wilson, 1986).

The rise of LLM-generated text disrupts these frameworks (Floridi and Chiriatti, 2020). Do AIgenerated documents bear the same weight without deliberate intent? This uncertainty can raising questions about authorship, intellectual property, and trust, especially for scientific or legal text (Raji et al., 2022; Huang and Rust, 2021; van Dis et al., 2023). Authorship has historically entailed a social contract: a published text can be challenged or critiqued, holding its human creator(s) responsible for factual or ethical shortcomings (Woodmansee, 1994; Chartier, 1994). But with LLM-authored text, accountability becomes diffused: does it lie with the prompter, the model trainer, or the dataset creator? This diffusion of responsibility strains traditional legal and academic norms (Hacker et al., 2023; Mittelstadt and Floridi, 2016; Kosseff, 2019). As the intent and accountability of text becomes murky, its meaningfulness and trustworthiness can become suspect too.

Language as a Cognitive Interface: Beyond communication, language shapes cognition and our capacity to abstract and solve problems (Clark and Chalmers, 1998; Vygotsky, 1978; Whorf, 1956). It is often considered an "interface" to thought. Research in child cognitive development suggests that engagement with language enables reasoning, cognitive flexibility, and problem-solving (Tomasello,

2003; Bruner, 1983; Lakoff and Johnson, 1980). 159 While some argue that LLMs function as cognitive 160 enhancers (Clark and Chalmers, 1998; Warwick, 161 2003), others caution that reliance on LLM-driven 162 generation can lead to reduced cognitive engage-163 ment (Nichols, 2021; Carr, 2010). In particular, 164 LLM-driven writing and summarization has raised 165 concerns about cognitive deskilling (Carr, 2011; 166 Lai and Viering, 2022). Studies show that composition itself is integral to thinking, forcing individu-168 als to clarify ambiguity, structure arguments, and synthesize knowledge (Kellogg, 2008; Galbraith, 170 1999). Further, LLM-generated summarization 171 risks eroding cognitive effort, like digital offload-172 ing has been shown to reduce critical engagement 173 (Sparrow et al., 2011; Nichols, 2021). 174

**Chain-of-Thought and Cognitive Parallels:** 175 The emergence of chain-of-thought prompting 176 (CoT) (Wei et al., 2022) represents a major shift 177 in LLM problem-solving. By externalizing logical 178 steps, CoT compensates for the depth limitations 179 of transformer architectures (Vaswani et al., 2017; 180 Yao et al., 2023). This mirrors how humans ar-181 ticulate thoughts through language, diagrams, or writing to enhance problem-solving (Clark and 183 Chalmers, 1998; Menary, 2010). Beyond computational efficiency, CoT also bears parallels to how externalizing reasoning through symbols has been linked to the expansion of human intelligence (Deacon, 1997; Dor, 2015). If language 188 189 enabled humans to extend cognition beyond individual memory, CoT might mark a similar milestone in LLM development. Whether CoT aug-191 192 ments human intelligence or leads intellectual complacency will depend on how societies integrate 193 LLM-based thought and reasoning in education, 194 work, and decision-making. 195

3 The Crisis of Language

196

197Given these philosophical foundations, the increas-<br/>ing role of LLMs ruptures the linguistic landscape198ing role of LLMs ruptures the linguistic landscape199through two forces: (1) Semantic Drift & Model200Collapse, the idea that the influx of AI-generated201text can shift the distribution and meaning of lan-<br/>guage, and lead to compounding errors; and (2)203Eroding Epistemic Trust in text, epitomized by<br/>what we term epistemic doppelgängers (LLM out-<br/>puts that are indistinguishable from human outputs).206We also suggest a metric, authorship entropy, to<br/>represent the uncertainty about the origin of a text.

# 3.1 Semantic Drift and Model Collapse

Semantic drift refers to changes in language usage and meaning over time, reflecting cultural and social evolution. However, large-scale LLMgenerated content can accelerate or redirect semantic change. For example, they might reinforce common phrases while underrepresenting less frequent expressions (Raji et al., 2022). LLMs trained on text that partially includes their own synthetic outputs can experience compounding errors. Repeated assimilation of AI-generated text leads to a shift away from organic language distributions (Shumailov et al., 2023; Carlini et al., 2023). Over time, certain stylistic artifacts become over-represented, contributing to *model collapse* (Menick et al., 2022), where the system's expressive range narrows towards the mean of what LLMs produce. While prior work has studied distribution shift in active learning (Blitzer et al., 2007), LLM self-ingestion is a novel feedback loop.

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

Semantic drift and model collapse are intertwined in a hybrid loop of human and LLM language production (Bommasani et al., 2021). Human writing feeds the training of LLMs, and LLM outputs in turn influence human writing and future training data. Without intervention, this loop may have an unintended equilibrium: language evolution not driven by human innovation, but by statistical characteristics of LLM generated text. The outcome can be a loss of semantic clarity (as meanings shift in unpredictable ways) and reduced linguistic innovation (Weidinger et al., 2021; Bender et al., 2021). If significant portions of text that we read comes to be machine-generated, we should also ask if this can deteriorate our cognitive diversity at a societal level.

# 3.2 Eroding Trust and Accountability

Texts have long been vehicles for accountability. Historically, authors were usually identifiable, and could be praised, critiqued, or legally challenged based on their claims (Woodmansee, 1994; Chartier, 1994). LLM-generated content fragments this chain of responsibility. This has significant implications for defamation suits and retraction practices for erroneous statements (Kosseff, 2019). Even more pressingly, online disinformation campaigns leveraging AI threaten political discourse, as citizenry loses clarity on who authors the narratives shaping public opinion (Weidinger et al., 2021; Chesney and Citron, 2019).

356

357

309

310

311

Empirical evidence is growing that synthetic text can fuel coordinated misinformation. Recent experiments demonstrate that AI-generated content is capable of crafting coherent yet deceptive social media campaigns, blurring the line between authentic and automated discourse (Zellers et al., 2019). Furthermore, large-scale language models have been observed to inadvertently plagiarize, amplify biases, and perpetuate stereotypes from their training data (Bender et al., 2021; Hovy and Spruit, 2016; Carlini et al., 2023). Without robust authorship signals or provenance tracking, verifying source credibility becomes increasingly challenging, raising concerns about accountability in digital information ecosystems (Gehrmann et al., 2019; Kirchenbauer et al., 2023).

258

259

263

264

267

269

271

272

276

277

279

287

291

296

297

299

304

305

308

# **3.3** Epistemic Doppelgängers and Authorship Entropy

LLMs can produce text nearly indistinguishable from human writing. We refer to such outputs as epistemic doppelgängers: texts that impersonate human authorship so convincingly they can fool not only casual readers, but editors, teachers, and even domain experts. As with a human doppelgänger, the deception isn't necessarily malicious, but its uncanniness can be destabilizing. GPT-generated news articles are often rated as more trustworthy than authentic ones (Zellers et al., 2019), and even the best AI detectors rarely surpass 70% accuracy. Worse, detection systems are often only effective when closely matched to the model they're trying to catch, making them vulnerable to fine-tuning, or strategic prompting. In short, epistemic doppelgängers erode the assumption that a well-formed sentence signals a human mind.

This epistemic ambiguity leads us to what we call authorship entropy, a measure of uncertainty of text authorship. In a world where all documents are confidently human-written, authorship entropy is low: the provenance of text is legible, even if anonymous. But in an AI-saturated ecosystem, the space of plausible authors expands. By modeling this uncertainty as a probability distribution and applying Shannon entropy, we can quantify how "foggy" the authorship landscape is. Rising authorship entropy destabilizes trust: people may become suspicious of legitimate texts, or indifferent to provenance altogether. It weakens accountability: if we don't know who wrote something, we can't assign responsibility. Also, AI authors, by definition, evade moral blame.

# 4 Technical Foundations: Authenticating Authorship & Quantifying Drift

Our discussion thus far has been primarily conceptual. In the next section, we explore technical interventions aimed at reclaiming human accountability and reducing authorship entropy.

# 4.1 Author Graphs & Proof-of-Interaction

A possible direction is embedding provenance and requirement of human interaction in the text generation process itself. For example, a hybrid authorship graph can represent relationships between human users, LLMs and texts that they generate or indirectly influence. To explain, a document's node might have edges from an LLM node (if an AI drafted it) and a human node (who guided or edited it). If an AI's training data included that document, an edge from the document back to the AI node ("trains") can be included, forming a cyclic network of influence. Figure 2 shows an example of such a graph. Such explicit representations of provenance and sources can provide grounding to enforce downstream accountability.

A practical implementation of this can be through *Proof-of-Interaction* (*PoI*) mechanisms that ensure that a human was substantially involved in creating a text. For instance, an editor can sign off on an AI-generated passage after verifying it, or a platform can require that any AI assistance be logged and attested. Some have proposed protocols where documents carry embedded metadata or hashes that link to records of the human-AI collaboration that produced them. If a document cannot present such proof-of-interaction, it might not be trusted for certain uses.

Building on blockchain-inspired ideas (Narayanan et al., 2016) and prior research on authorship verification (Stamatatos, 2009; Layton and Watters, 2020), we propose a *process-based* approach to demonstrate genuine human involvement. Let  $\mathcal{T}$  denote a text document, and let  $\mathcal{C} = \{C_1, C_2, \dots, C_n\}$  denote the set of discrete human contributions (e.g., edits, approvals, interventions). Then we we define its proof-of-interaction score for the text as:

$$\operatorname{PoI}(\mathcal{T}) = f(\mathcal{C}, \operatorname{HashChain}).$$
 353

Here, HashChain refers to a cryptographically secure chain of hashes that records the provenance of these edits.  $f(\cdot)$  outputs a scalar or structured PoI score or certificate. This function can be designed

374

358

39

395 396 397 to validate sufficiency (e.g., a minimum amount of human interaction occurred), or provide a verifiable attestation (e.g., a digitally signed summary). Newer works in *human-in-the-loop* text genera-

tion have proposed storing metadata that logs partial revisions (Cecchi and Babkin, 2024; Kang et al., 2024). Merging these ideas with zero-knowledge proofs can preserve user privacy. While this does not solve all issues (e.g., adversaries can simply simulate keystrokes or partial edits), such a system raises the cost of deception and provides an auditable trail of interaction.

This idea also aligns with Chain-of-thought (Wei et al., 2022; Kojima et al., 2022) oversight, where an LLM seeks human verification for intermediate reasoning. Some developers propose useraudited chains-of-thought, letting humans see exactly which steps an LLM took (Wang et al., 2022). Future research can unify chain-of-thought logs with proof-of-work, offering a secure record of how text was generated . This can clarify the roles of LLMs and humans in authoring text, although this approach may present challenges in preserving user privacy (Khowaja et al., 2023; Abadi et al., 2016; Glymour et al., 2023). These changes will necessarily introduce friction, and may be tedious for users. However, a proof-of-interaction system can ensure that every text is connected to at least one human via a "verified" edge. This can maintain the principle that for any published text, one can point to a human accountable for it.

### 4.2 Metrics for Semantic Drift

Verifying authorship addresses who wrote the text. But we should also ask what is being written. We propose tracking language changes by defining metrics for semantic drift and linguistic diversity, comparing human-authored and LLM-generated text periodically. By measuring shifts in word frequencies, syntax, or topics, we can identify drift if metaphoric language or dialectal terms decrease while AI-generated phrases increase.

It is also worth monitoring model-internal drift: how successive generations of LLMs differ when trained on data that includes prior LLMs' outputs. If  $P_{human}$  and  $P_{LLM_{\theta}}$  indicate the probability distributions of language utterances at a discrete time step t, and if  $\alpha$  denotes the proportion of LLM generated data, then the distribution of training data that will be used to train the next iteration of the



Figure 2: An illustrative hybrid authorship graph, representing provenance and interactions between human agents, an LLM agent, and texts. In this example, Human H1 writes Document 1, which is later used in training the LLM. Document 2 is co-authored by H1 and the LLM (perhaps H1 edited text generated by the LLM). Document 3 is authored solely by the LLM.

LLMs, 
$$P_{LLM}^{t+1}$$
, is given by:  
 $P_{\text{mix}}^{(t)}(\mathbf{x}) = (1 - \alpha)P_{\text{human}}(\mathbf{x}) + \alpha P_{\text{LLM}_{\theta}}^{(t)}(\mathbf{x})$ 

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

as LLM outputs re-enter training data (Shumailov et al., 2023). If  $P_{\text{mix}}^{(t)}$  increasingly diverges from  $P_{\text{human}}$ , the model parameters  $\theta$  risk converging to a subspace that fails to capture the richness of actual human language patterns. Further computational or theoretical insights might be found in work on catastrophic forgetting (Kirkpatrick et al., 2017) and domain shift (Ganin et al., 2016). While the notion is not new (prior studies on machinein-the-loop domain adaptation raise similar concerns (Ruder, 2019)), our contribution is to highlight how large volumes of synthetic text can nudge language distributions away from natural usage. We propose coupling distributional metrics (e.g., KL divergence) with textual diversity indices to monitor linguistic homogenization.

Empirically testing these metrics on real corpora that blend human and AI-generated text remains a priority for future research. Experiments with smaller LLMs (Carlini et al., 2023; Menick et al., 2022) suggest that repeated synthetic ingestion amplifies shallow lexical patterns. In Figure 3, we plot the JS divergence between unigram distributions in a base human-authored corpus, and a GPT2-base model that is iteratively re-trained on its own sampled data. We note a clear and increasing semantic drift with an increasing number of steps.



ticipation in the generative process but is difficult to scale. Education should focus on teaching skills that AI cannot easily replace.

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

Academic Scholarship The academic ecosystem assumes text is a reflection of an author's intellect. Automated text generation challenges this, raising concerns over AI authorship and scholarly contributions (Willis and Williams, 2023). In the short term, LLMs risk hallucinations (Ji et al., 2023) and plagiarism (van Dis et al., 2023). More seriously, they can flood peer reviews and obscure genuine innovation. Ideally though, LLMs should boost research productivity and accelerate scientific progress.

Professional Settings In many industries, cover letters, writing samples, and portfolio websites are used to gauge candidates' communication skills and expertise (Sternberg and Williams, 1997). LLM tools now make it easy to create polished but shallow applications, complicating hiring managers' ability to assess true abilities. Some organizations are turning to live assessments like realtime writing tests or structured panel interviews (Koch et al., 2015; Levashina et al., 2014). But these can be difficult for introverts, non-native speakers, or candidates who do better with written communication (van Tubergen and Kalmijn, 2014; Hu et al., 2020). Managing this requires a delicate dance between fairness and authenticity.

The Public Sphere LLMs are reshaping public discourse via AI-generated content, sparking concerns about amplification and distortion (Zellers et al., 2019; Ferrara, 2020). Disinformation campaigns exploit AI's capability to produce misleading content, drowning out authentic voices and confusing public understanding. Although detection methods advance, the adversarial landscape perpetuates a constant arms race (see Section 6). Conversely, LLMs present opportunities to democratize communication by reducing barriers for individuals with limited writing skills, disabilities, or those who are non-native speakers (Paritosh et al., 2022; Xu et al., 2022).

In all of these domains just discussed, a common thread is that trust is threatened by automatic text generation from LLMs. As trust erodes, institutions will react by imposing stricter verification, leading to friction, surveillance, or cynicism. The challenge is developing norms that preserve the value of human contribution and ensure transparency.

shows the Jensen-Shannon divergence between a base human-authored text distribution and iteratively drifted synthetic text distributions. As synthetic text is repeatedly generated and reintroduced into training data, the divergence increases, illustrating the risk of semantic drift and potential loss of linguistic diversity over time. At each iteration, the GPT2 model is fine-tuned on text sampled from the GPT2 model in the previous step

Figure 3: Semantic Drift in Synthetic Text: The plot

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443



#### 5 **Societal Implications**

While LLMs can improve productivity and streamline workflows, their widespread adoption raises concerns about academic integrity, scholarly publishing, professional hiring, and public discourse (Coglianese and Lehr, 2017; Cotton et al., 2023; van Dis et al., 2023).

**Education** LLMs are increasingly used as study 444 aids, enhancing access for learners with different 445 language skills (Khan et al., 2023; Chaudhuri et al., 446 2021; Xu et al., 2022; Luckin et al., 2023). How-447 ever, relying on LLMs for tasks like programming 448 or essay writing can weaken essential skills: algo-449 rithmic thinking, structured argumentation, and cre-450 451 ativity (Cotton et al., 2023; Perkins and Salomon, 1989). To address this, some schools use real-time 452 or proctored writing tasks or oral exams to ensure 453 understanding (Lund and Wang, 2023). Our chain-454 of-thought synergy (Sec. 7) encourages student par-455

507

530

531

533

534

535

536

539

540

541

542

543

546

548

550

552

#### Labels & Classifiers wont save us 6

Proposals such as watermarking and policy bans, 506 while helpful in the short term, offer only superficial remedies (Gehrmann et al., 2019; Zellers et al., 2019; Papernot et al., 2016). 509

Watermarking and Detection Arms Races Wa-510 termarking remains fragile against adversarial attacks like paraphrasing (Kirchenbauer et al., 2023). 512 Detection classifiers also struggle with robustness 513 as LLMs adapt and human post-editing obfuscates 514 machine origins (Holtzman et al., 2020). This cre-515 ates a resource-intensive cat-and-mouse dynamic 516 without stable solutions (Gallagher et al., 2023). 517 More critically, these methods do not resolve at-518 tribution, leaving ethical and legal questions unan-519 swered (Authors, 2023; Devinney, 2023). 520

Policy Bans and their Limitations Bans on AI-521 assisted writing are unenforceable due to weak de-522 tection and strong incentives for LLM use, effectively becoming honor systems (Devinney, 2023). 524 Lagging legislation creates a fragmented regulatory 525 landscape (Hacker et al., 2023), and the global nature of digital communication allows easy circumvention of local policies (Katyal and Epps, 2022), failing to address authorship and accountability. 529

**Neglecting the Deeper Interpretive Question** Fundamentally, current measures do not restore a discernible human presence. If language's epistemic function relies on text as an intentional artifact, superficial labeling fails to reattach text to a mind (Bender and Koller, 2020). It still leaves us with a fundamentally ambiguous communicative landscape, and risks putting us in an era of permanent ambiguity in textual interpretation.

#### 7 **Rethinking Language & Authorship**

We propose recalibrating the role of LLMs in language by leveraging their benefits while preserving human traits like intentionality, accountability, and diversity of thought. This requires an interdisciplinary approach involving NLP, cognitive science, ethics, law, and education. In this section, we refine previous suggestions and introduce ideas on human-AI collaboration frameworks, governance, and cultural appreciation of human-only work (Mittelstadt and Floridi, 2016; Floridi, 2019).

#### 7.1 Chain-of-Thought with Human Oversight

As previously mentioned in Section 4, we advocate for embedding AI within a structured chain-ofthought framework that requires human oversight at key decision points (Wei et al., 2022). In this paradigm, LLMs can generate partial outlines, intermediate arguments, or suggested revisions, but finalization has to be authenticated by a human user after consideration. By logging human-AI interactions through an auditable chain (with appropriate privacy safeguards), this method establishes a transparent record that delineates AI-generated content from human refinement, addressing concerns about accountability and intellectual ownership (Wang et al., 2022; Christiano, 2022).

553

554

555

556

557

558

559

560

561

562

563

564

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

600

# 7.2 Governance & Collaborative Policy

Governance for LLM-usage has to be a negotiation (between policymakers, educators, and user communities, etc.) for it to work, rather than a prescription (Floridi and Chiriatti, 2020; Hacker et al., 2023). Several directions seem promising. First, AI contribution statements, similar to conflictof-interest disclosures, can prompt authors to declare the extent and nature of LLM involvement (Devinney, 2023). Second, labeling protocols for governmental or legal texts can introduce metadata or disclaimers to flag LLM-generated contents (Union, 2023). Also, ethical AI certification programs, modeled on data protection seals, can help LLM developers conform with regulations such as the EU AI Act (Union, 2023).

# 7.3 Educational Reforms and Cultural Shifts

To prevent cognitive deskilling, education has to pivot and adapt. Assignments will have to adapt to the inevitability of the use of LLMs for drafting, but can require students to justify revisions and incentivize peer engagement (Lai and Viering, 2022). Assessments like live problem-solving and debates will have to focus on substance over polish (Paul and Elder, 2007; Lipman, 2003; Chi and Wylie, 2014; Freeman et al., 2014). Finally, students should be encouraged to play with LLMs, and taught to interrogate them. AI literacy should be a form of critical literacy, where students learn not simply how to use LLMs, but when and when not to (Bowman and Reeves, 2015).

# 7.4 Human-Only Publishing Spaces

A potential direction is establishing "human-only" publishing spaces: media outlets, or creative communities that employ verification measures (such as mechanisms like proof-of-work logs) to ensure

654

655

656

657

658

659

660

661

662

663

664

665

666

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

that any content reflects considerable human intellectual effort and creativity. These spaces would offer a parallel track for those who value direct human expression. Some journals already forbid undisclosed AI collaboration for final submissions (Board, 2023). A fiction community might pride itself on entirely human-crafted stories. Like organic labels in food, these spaces can serve audiences that value authentic human expression—akin to "slow food" movements in a fast-food world (Petrini, 2001). Over time, such enclaves can serve as a 'control group', preserving the standards and norms associated with human authorship.

# 8 Alternate Views

601

602

603

606

607

610

612

615

616

617

618

We have emphasized how LLMs may erode traditional assumptions about language and authorship. But many scholars have a more sanguine outlook, and do not frame LLM-driven automation as a threat to language. For balance, we summarize key arguments in these perspectives.

621Democratization and AccessibilityLLMs can622enable non-native speakers and individuals with dis-623abilities to participate in public discourse (Norvig624and Thrun, 2009; Ogawa et al., 2022). By automat-625ing surface-level writing concerns, these tools al-626low users to focus on substantive ideas. For ex-627ample, spell-checkers, were once controversial too628(Felton, 2023; Christiansen, 2021). Additionally,629LLMs increase accessibility for users with impair-630ments (Wagner et al., 2020), reframing 'linguistic631inclusivity' as a positive evolution.

632Accelerated Knowledge DisseminationSum-633marization tools help researchers digest literature634efficiently (Fabbri et al., 2022; Sharma et al., 2022),635and multilingual translation expands access to spe-636cialized knowledge (Fan et al., 2021; Artetxe and637Schwenk, 2019). With editorial oversight, these638outputs can enhance comprehension without com-639promising reliability (Szegedy et al., 2022). Ad-640vocates argue that with transparency, LLMs can641strengthen epistemic ecosystems rather than harm642them (Diakopoulos, 2016).

Evolving Norms of Collaboration In many domains, collaborative authorship is standard (technical manuals, corporate reports, etc.), which rarely reflect a single voice (Darics, 2020; Leonard and Noonan, 2020). In this context, LLMs are seen as additional collaborators (Krause et al., 2022; Dinan et al., 2022). Rather than undermining authorship,

they may shift workflows, with new roles emerging for human editors and fact-checkers (Eisenstein and McNamara, 2023; Roose and Sullivan, 2023).

**Empirical Evidence of Positive Outcomes** Some studies suggest that, when used responsibly, LLMs can enhance writing without weakening critical thinking. They support non-native and novice writers in building fluency (Lee et al., 2022; Laubrock et al., 2022). In collaborative environments, there is evidence that AI systems help clarity, and can identify redundancy (Yosinski et al., 2023; Rahimi et al., 2021).

Broadly, these perspectives argue that LLMs are not existential threats to the integrity of language, and that a 'crisis' of authorship is neither new nor uniquely AI-induced. Rather, this is a natural evolution in how we produce and share ideas. Possibly, when questions about the origin and intent of a text fade, newer and better-suited norms can emerge in the linguistic landscape to replace them.

# 9 Conclusion & Reflection

LLMs are here to stay. Yet, they challenge the epistemic and communicative foundations of language by decoupling text from human intent. While automated or impersonal writing is not new (for example, memos or legal boilerplates), LLMs amplify this disconnection at an unseen scale. This shift raises several prickly questions about intent, accountability and trust. Any solutions here must preserve the role of language in human agency and accountability. They key will be to assert human presence in language: whether through through cryptographic attestations, proof-of-interaction, or new cultural norms.

The future of language doesn't hinge as much on building better models, but on what we choose to protect. If we can collaborate to implement verifiable authorship and enforceable audit mechanisms, we may harness LLMs' advantages without surrendering the uniquely human dimensions of language. But a failure to act can lead to communication devoid of color, reeking with hollow expressions, and diluted cognitive depth in people. The path to hell is famously paved with good intentions. Still, through ingenuity and foresight, LLM innovations could be steered toward enhancing human creativity, rather than eroding the intellectual bedrock which is its basis.

# Limitations

698

711

712

713

714

715

716

717

718

719

721

725

727

729

730

732

734

738

739

740

741

742

743

744

745

699By design, this paper is more diagnostic than pre-700scriptive. We introduce conceptual tools like epis-701temic doppelgängers and authorship entropy to702make sense of the shifting linguistic terrain. How-703ever, many of these constructs remain speculative704without empirical grounding. Nor do we pretend705that quantitative metrics alone can capture the con-706sequences of LLM saturation. What we offer is a707framework to think with, not a solution to deploy.

Second, some of our proposals (such proof-ofinteraction logs, and human-only publishing enclaves) are challenging to implement and reify. They require infrastructure, extensive cooperation, and cultural shifts that may not be welcome. We should also acknowledge a significant tension: the paper champions the pre-eminence of human intention in language, but we do not wish to gatekeep expression or discourage the increasingly creative and original uses of LLMs. The challenge is to protect the epistemic integrity of language without devolving into 'purity tests'. To truly solve this challenge will requires contending with lived social realities of people, not just technical design.

# AI Use Acknowledgment

In this work, we acknowledge the use of AI assistance in the following cases in accordance with the ACL Policy on AI Writing Assistance: assistance with literature search and review, proof-reading and refining the language of the paper, and analytical code. We utilized AI tools for searching for relevant literature, help with writing code for analysis related to Figure 3, and code for diagram generation for Figure 2.

# References

- Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security* (*CCS*), pages 308–318.
- Mikel Artetxe and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zeroshot cross-lingual transfer and beyond. *Transactions of the Association for Computational Linguistics*, 7:597–610.
  - J. L. Austin. 1975. *How to Do Things with Words*, second edition. Harvard University Press.

Anonymous Authors. 2023. Gpt and the plagiarism problem: Assessing ai's impact on academic integrity. *Journal of Ethics in AI*. 746

747

749

750

752

753

754

755

756

757

758

759

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT)*, pages 610–623.
- Emily M. Bender and Alexander Koller. 2020. Climbing towards NLU: On meaning, form, and understanding in the age of data. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 5185–5198.
- John Blitzer, Mark Dredze, and Fernando Pereira. 2007. Biographies, bollywood, boom-boxes, and blenders: Domain adaptation for sentiment classification. In Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL), pages 440–447.
- Nature Editorial Board. 2023. Ai and authorship: Nature's policy on undisclosed ai collaboration.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, and 1 others. 2021. On the opportunities and risks of foundation models. *arXiv e-prints*, pages arXiv–2108.
- Nicholas A. Bowman and Anne E. Reeves. 2015. Rethinking media literacy in the age of algorithmic curation. *Journal of Media Education*, 6(3):14–24.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. Advances in Neural Information Processing Systems (NeurIPS), 33:1877– 1901.
- Jerome Bruner. 1983. Child's Talk: Learning to Use Language. W. W. Norton Company.
- Nicholas Carlini, Matthew Jagielski, Chiyuan Zhang, Katherine Lee, Christopher A. Choquette-Choo, Jacob Imber, Andreas Terzis, Nicholas Frosst, Ilya Mironov, Vasisht Duddu, and 1 others. 2023. Poisoning web-scale datasets is practical. *arXiv preprint arXiv:2305.00956*.
- Nicholas Carr. 2010. *The Shallows: What the Internet Is Doing to Our Brains.* W.W. Norton Company.
- Nicholas Carr. 2011. *The Shallows: What the Internet is Doing to Our Brains*. W. W. Norton & Company.

799

- 830 831 832
- 833 834 835

8 8 8

836

- 841
- 842 843 844

8

847 848

8

849 850 851

- Lucas Cecchi and Petr Babkin. 2024. Reportgpt: Human-in-the-loop verifiable table-to-text generation. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track, pages 529–537.
- Roger Chartier. 1994. The Order of Books: Readers, Authors, and Libraries in Europe Between the Fourteenth and Eighteenth Centuries. Stanford University Press.
- Soham Chaudhuri, Varun Kumar, and Marti Hearst. 2021. Ai-powered writing assistants: Enhancing learning and creativity. In *Proceedings of the Conference on Educational Data Mining (EDM)*, pages 344–355.
- Bobby Chesney and Danielle Keats Citron. 2019. Deepfakes and the new disinformation war: The coming age of post-truth geopolitics. *Foreign Affairs*, 98(1):147–155.
- Michelene T. H. Chi and Rachel Wylie. 2014. The icap framework: Linking active learning to cognitive engagement. *Educational Psychologist*, 49(4):219– 243.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Anselm Levskaya, Tyler Wang, Nan Du, Yinhan Liu, and 6 others. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Paul Christiano. 2022. Ai alignment and the role of human oversight in generative models. *Alignment Research Journal*, 4:101–128.
- Meredith Christiansen. 2021. Algorithmic writing and digital literacy: How ai is reshaping composition. *Digital Studies*, 12:55–73.
- Andy Clark and David Chalmers. 1998. The extended mind. *Analysis*, 58(1):7–19.
- Cary Coglianese and David Lehr. 2017. Regulating by robot: Administrative decision making in the machine-learning era. *Georgetown Law Journal*, 105:1147–1223.
- Debbie Cotton, Peter Cotton, and Sarah Shipway. 2023. Chatgpt, ai and the impact on academic integrity: Aiassisted student writing. *International Journal for Educational Integrity*, 19(1):1–16.
- Erika Darics. 2020. E-voice: A multimodal perspective on institutional writing in digital environments. *Discourse, Context Media*, 35:100391.
- Terrence W. Deacon. 1997. *The Symbolic Species: The Co-evolution of Language and the Brain*. W.W. Norton Company.

DeepSeek-AI. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.

852

853

854

855

856

857

858

859

860

861

862

863

864

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

- Timothy Devinney. 2023. Plagiarism in the age of ai: Who owns the output? *Journal of Business Ethics*.
- Nicholas Diakopoulos. 2016. Accountability in algorithmic decision making. *Communications of the ACM*, 59(2):56–62.
- Emily Dinan, Laura Perez, and Jason Weston. 2022. Collaborative ai: Integrating language models into professional writing teams. In *Proceedings of NeurIPS 2022*, pages 5123–5136.
- Daniel Dor. 2015. *The Instruction of Imagination: Language as a Social Communication Technology*. Oxford University Press.
- Fernando Duarte, Maarten Sap, and Yejin Choi. 2022. Ai-generated speech and the decline of authentic online discourse. *Proceedings of the 2022 Conference on Fairness, Accountability, and Transparency (FAccT).*
- Ronald Dworkin. 1996. *Freedom's Law: The Moral Reading of the American Constitution*. Harvard University Press.
- Jacob Eisenstein and Danielle McNamara. 2023. Human-ai collaboration in writing: Rethinking authorship and editorial oversight. *AI Society*, 38(2):177–192.
- Alexander Fabbri, Irene Li, Tianyi Tang, Caiming Xiong, and Dragomir Radev. 2022. Qmsum: A new benchmark for query-based multi-document summa-rization. In *Proceedings of NAACL 2022*, pages 6145–6162.
- Angela Fan, Shruti Bhosale, Holger Schwenk, Alexander Baevski, Guillaume Lample, and Michael Auli. 2021. Beyond english-centric multilingual machine translation. In *Proceedings of ACL 2021*, pages 4403– 4419.
- James Felton. 2023. From spell-checkers to ai writing assistants: The evolution of digital literacy tools. *Journal of Digital Communication*, 17(2):89–107.
- Emilio Ferrara. 2020. Characterizing social media manipulation in the 2020 u.s. presidential election. *First Monday*, 25(11).
- Luciano Floridi. 2013. *The Ethics of Information*. Oxford University Press.
- Luciano Floridi. 2019. Establishing the rules for building trustworthy ai. *Nature Machine Intelligence*, 1(6):261–262.
- Luciano Floridi and Marcello Chiriatti. 2020. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30:681–694.
- Michel Foucault. 1984. What is an author? In Paul Rabinow, editor, *The Foucault Reader*, pages 101– 120. Pantheon Books.

909

910

8415.

works.

Press.

1112-1123.

17(59):1-35.

tational Linguistics (ACL).

David Galbraith. 1999.

Scott Freeman, Sarah L. Eddy, Miles McDonough,

Michelle K. Smith, Nnadozie Okoroafor, Hannah

Jordt, and Mary Pat Wenderoth. 2014. Active learn-

ing increases student performance in science, engi-

neering, and mathematics. Proceedings of the Na-

tional Academy of Sciences (PNAS), 111(23):8410-

Iason Gabriel. 2020. Artificial intelligence, values, and

alignment. Minds and Machines, 30(4):411-437.

constituting process. Erkenntnis, 50:357-370.

John Gallagher, Jacob Hilton, and Owain Evans. 2023.

Adversarial robustness in ai-generated text detection:

A losing battle? *arXiv preprint arXiv:2305.07692*.

Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Lavi-

olette, Mario Marchand, and Victor Lempitsky.

2016. Domain-adversarial training of neural net-

Sebastian Gehrmann, Hendrik Strobelt, and Alexan-

der M. Rush. 2019. Gltr: Statistical detection and

visualization of ai-generated text. Proceedings of the

57th Annual Meeting of the Association for Compu-

Clark Glymour, David Danks, and Peter Spirtes. 2023.

H.P. Grice. 1975. Logic and conversation. In Syntax

Philipp Hacker, Andreas Engel, and Marco Mauer. 2023.

Regulating chatgpt and other large generative ai mod-

els. In Proceedings of the 2023 ACM Conference on

Fairness, Accountability, and Transparency, pages

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text

Conference on Learning Representations (ICLR).

Dirk Hovy and Shannon L. Spruit. 2016. The social impact of natural language processing. *Proceedings* 

of the 54th Annual Meeting of the Association for

Computational Linguistics (ACL), pages 591–598.

Yang Hu, Heather R. Younger, and Patricia M. Green-

field. 2020. Language and hiring bias: How intro-

verts and neurodiverse candidates face discrimination

in spontaneous evaluations. Journal of Business and

Ming-Hui Huang and Roland T. Rust. 2021. Artifi-

cial intelligence in business: Promise, pitfalls, and

prospects. Journal of Service Research, 24(1):3-6.

Psychology, 35(2):215–230.

degeneration. In Proceedings of the International

and Semantics, volume 3, pages 41-58. Academic

oversight. Artificial Intelligence Review.

Ai explainability and its limits: The role of human

Journal of Machine Learning Research,

Writing as a knowledge-

- 911 912 913 914 915 916
- 918 919 920 921 922

917

924 925

923

- 927 928 929
- 9

931 932

- 933
- 934 935

936 937 938

9; 9/

941 942

943 944

94

946 947

94

950 951

952 953

95

954 955

955 956 Ben Hutchinson, Jasmine Collins, Mark Diaz, and Qian Yang. 2023. Ai for accessibility: Enabling inclusive digital communication. *CHI Conference on Human Factors in Computing Systems*.

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1004

1005

- Zhengbao Ji, Zhiting Hu, Patrick Lewis, and Meta AI. 2023. Survey of hallucination in large language models. *Transactions of the Association for Computational Linguistics*.
- Hong Jin Kang, Fabrice Harel-Canada, Muhammad Ali Gulzar, Violet Peng, and Miryung Kim. 2024. Human-in-the-loop synthetic text data inspection with provenance tracking. *arXiv preprint arXiv:2404.18881*.
- Neal Katyal and Daniel Epps. 2022. The criminal regulation of artificial intelligence. *Harvard Law Review*, 135:412–468.
- Ronald T. Kellogg. 2008. Training writing skills: A cognitive developmental perspective. *Journal of Writing Research*, 1(1):1–26.
- Salman Khan, Pranav Rajpurkar, and Daphne Koller. 2023. Ai tutors: The role of large language models in personalized education. *arXiv preprint arXiv:2303.11288*.
- Sahar Khowaja, Ammar Ahmad, Khaled Salah, Raja Jayaraman, and Ibrar Yaqoob. 2023. Blockchain for ai: Review and open research challenges. *IEEE Transactions on Artificial Intelligence*, 4(1):1–15.
- Johannes Kirchenbauer, Jonas Geiping, Yuxuan Han, Elliot Creager, Ari S. Morcos, and Tom Goldstein. 2023. A watermark for large language models. *arXiv preprint arXiv:2307.06624*.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526.
- Andrew J. Koch, Robert D. Gerber, and Sarah C. Roberts. 2015. Structured interviews: Reducing bias and increasing hiring effectiveness. *Journal of Applied Psychology*, 100(3):775–789.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199– 22213.
- Jeff Kosseff. 2019. *The Twenty-Six Words That Created the Internet*. Cornell University Press.
- Ben Krause, Ethan Wilcox, Max Bittker, and Christopher D. Manning. 2022. Co-authoring with ai: How1007Ilms shape collaborative writing practices. Transactions of the ACL, 10:413–429.1010

- 1013 1015 1017 1018 1019 1020 1021 1022 1023 1026 1028 1031 1033 1034 1035 1036 1037 1038 1042 1043 1044 1045 1046 1048 1049 1050 1051 1053 1055 1056 1057 1058

1012

1059

1060 1061

1062 1063

- Vivian Lai and Laura Viering. 2022. Ai assistance and over-reliance: Implications for human judgment and decision-making. Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI).
- George Lakoff and Mark Johnson. 1980. *Metaphors We Live By.* University of Chicago Press.
- Jochen Laubrock, Clara Martin, and Marc Brysbaert. 2022. Readership enhancement via ai-assisted writing tools: A cognitive science perspective. Cognitive Science, 46(4):e13120.
- Rebekah Layton and Carolyn Watters. 2020. Authorship attribution with deep learning. Digital Scholarship in the Humanities, 35(2):317–331.
- Jisoo Lee, Daniel McNamara, and Tanja Käser. 2022. Co-authoring with ai: How language models support second-language writing development. Computers Education, 186:104536.
- Brian Leonard and Kevin Noonan. 2020. Computational text generation in professional and technical writing. Journal of Business and Technical Communication, 34(4):451-474.
- Julia Levashina, Christian J. Hartwell, Frederick P. Morgeson, and Michael A. Campion. 2014. The structured employment interview: Narrative and quantitative review of the research literature. Personnel Psychology, 67(1):241-293.
- Matthew Lipman. 2003. Thinking in Education, 2nd edition. Cambridge University Press.
- Rose Luckin, Wayne Holmes, and Joshua Greer. 2023. Ai and education: The future of personalized learning. AI Education Journal, 4(1):112–130.
- Brian Lund and Weilin Wang. 2023. Reinventing assessments in the age of ai: From written exams to oral evaluations. AI Education Journal, 2(1):22-35.
- Richard Menary. 2010. Cognitive integration and the extended mind. In Richard Menary, editor, The Extended Mind, pages 227-243. MIT Press.
- Jacob Menick, Jeffrey Shlens, Xiaohua Zhai, Neil Houlsby, Andrea Gesmundo, Avital Oliver, and Karen Simonyan. 2022. Reducing the recurrence of errors in language model outputs. NeurIPS.
- Arthur I Miller. 2019. The artist in the machine: The world of AI-powered creativity. MIT Press.
- Brent Daniel Mittelstadt and Luciano Floridi. 2016. The ethics of algorithms: Mapping the debate. Big Data Society, 3(2):1-21.
- Arvind Narayanan, Joseph Bonneau, Edward Felten, Andrew Miller, and Steven Goldfeder. 2016. Bitcoin and Cryptocurrency Technologies: A Comprehensive Introduction. Princeton University Press.
- Tom Nichols. 2021. The death of expertise: The campaign against established knowledge and why it matters. Oxford University Press.

Peter Norvig and Sebastian Thrun. 2009. The Google 1064 Revolution: How AI is Transforming Language. 1065 O'Reilly Media.

1067

1069

1070

1071

1074

1078

1080

1081

1083

1084

1085

1086

1090

1091

1093

1094

1095

1096

1097

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

- Takashi Ogawa, Hiroshi Nakagawa, and Yuki Tanaka. 2022. Breaking barriers: Ai-assisted writing for nonnative speakers and people with disabilities. Computers Education, 184:104516.
- Cathy O'Neil. 2016. Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy. Crown Publishing Group.
- Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. 2016. The limitations of deep learning in adversarial settings. In Proceedings of the 1st IEEE European Symposium on Security and Privacy (EuroSP), pages 372-387.
- Praveen Paritosh, Elizabeth Clark, Luke Zettlemoyer, and Mihai Surdeanu. 2022. Critiquing ai writing: Understanding the benefits and risks of language models in content creation. arXiv preprint arXiv:2211.12760.
- Richard Paul and Linda Elder. 2007. The Thinker's Guide to Socratic Questioning. Foundation for Critical Thinking.
- David N. Perkins and Gavriel Salomon. 1989. Are cognitive skills context-bound? Educational Researcher, 18(1):16-25.
- Carlo Petrini. 2001. Slow Food: The Case for Taste. Columbia University Press.
- Plato. 1997. Complete Works. Hackett Publishing. Includes translations of \*Gorgias\* and \*Cratylus\*.
- Ali Rahimi, ChengXiang Zhai, and Rada Mihalcea. 2021. Collaborative ai writing: Balancing automation and human creativity. AI Society, 36(3):609-624.
- Inioluwa Deborah Raji, Andrew Smart, Rebecca N. White, Margaret Mitchell, Timnit Gebru, Ben Hutchinson, Jamila Smith-Loud, Daniel Theron, and Parker Barnes. 2022. Closing the ai accountability gap: Defining an end-to-end framework for internal algorithmic auditing. In Proceedings of ACM FAccT 2022, pages 343-357.
- Kevin Roose and Margaret Sullivan. 2023. Ai in journalism and corporate writing: A new era of assisted authorship. Journalism Studies, 24(6):815-832.
- Sebastian Ruder. 2019. Neural Transfer Learning for Natural Language Processing. Ph.D. thesis, National University of Ireland, Galway.
- Timo Schick, Roberto Wilfer, and Hinrich Schütze. 1112 2021. Self-diagnosis and self-debiasing: A pro-1113 posal for reducing corpus-based bias in NLP. arXiv 1114 preprint arXiv:2103.00453. 1115

John R. Searle. 1969. Speech Acts: An Essay in the Phi- losophy of Language. Cambridge University Press.	Frank van Tubergen and Matthijs Kalmijn. 2014. Lan- guage proficiency and early labor market entry of immigrants in the netherlands. <i>Journal of Ethnic and</i>	1167 1168 1169
Amanpreet Sharma, Thomas Wolf, and Sebastian Ruder.	Migration Studies, 40(3):405–424.	1170
2022. Bigbird: Summarization for large-scale text	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	1171
5801	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	1172
5601.	Kaiser, and Illia Polosukhin. 2017. Attention is all	1173
Ilia Shumailov, Yiren Zhao, Robert Mullins, Ross An-	you need. In Advances in Neural Information Pro-	1174
derson, and Nicolas Papernot. 2023. The curse of	cessing Systems (NeurIPS).	1175
recursion: Training on generated data makes models forget arXiv preprint arXiv:2305 17493	Lev S. Vygotsky. 1978. Mind in society: The develop-	1176
	ment of higher psychological processes.	1177
Irene Solaiman, Miles Brundage, Jack Clark, Amanda	Johannes Wagner, Emily M. Bender, and Martin Savic.	1178
Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford,	2020. Accessible ai: Enabling writing for users with	1179
Jasmine Wang, and Dario Amodei. 2019. Release	visual and motor impairments. AI Society, 35:301–	1180
arXiv preprint arXiv: 1908.09203	319.	1181
uixiv preprint uixiv.1900.09205.	Kevin Wang, Alexandre Variengien, Arthur Conmy,	1182
Betsy Sparrow, Jenny Liu, and Daniel M. Wegner. 2011.	Buck Shlegeris, and Jacob Steinhardt. 2022. In-	1183
Google effects on memory: Cognitive consequences	terpretability in the wild: a circuit for indirect ob-	1184
of having information at our fingertips. Science,	ject identification in gpt-2 small. arXiv preprint	1185
333(6043):776–778.	arXiv:2211.00593.	1186
	Kevin Warwick. 2003. Cyborg morals, cyborg values,	1187
Dan Sperber and Deirdre Wilson. 1986. <i>Relevance:</i>	cyborg ethics. Ethics and Information Technology,	1188
Press	5(3):131–137.	1189
11000.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	1190
Efstathios Stamatatos. 2009. A survey of modern au-	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	1191
thorship attribution methods. Journal of the Ameri-	and 1 others. 2022. Chain-of-thought prompting elic-	1192
can Society for Information Science and Technology,	its reasoning in large language models. Advances	1193
60(3):538–556.	in neural information processing systems, 35:24824–	1194
Pohert I. Sternherg and Wandy M. Williams 1007 In	24837.	1195
telligence Instruction and Assessment: Theory into	Laura Weidinger, John Mellor, Maribeth Rauh, Conor	1196
Practice Lawrence Erlbaum Associates	Griffin, Jonathan Uesato, Po-Sen Huang, Myra	1197
	Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh,	1198
Christian Szegedy, Yi Tay, and Alexander Raichuk.	and 1 others. 2021. Ethical and social risks	1199
2022. Resilient ai: Aligning language models with	of harm from language models. arXiv preprint	1200
ethical and epistemic standards. Journal of AI Ethics,	arXiv:2112.04359.	1201
5:310–328.	Benjamin Lee Whorf. 1956. Language, Thought, and	1202
NILL P. Team. 2022. No longuage left behind: Seeling	Reality: Selected Writings of Benjamin Lee Whorf.	1203
human-centered machine translation <i>arXiv preprint</i>	MIT Press.	1204
arXiv:2207.04672.	Malcolm Willis and Emma P Williams 2023 Should	1205
	ai be a co-author? ethical and academic perspectives	1205
Michael Tomasello. 2003. Constructing a Language: A	AI Society.	1207
Usage-Based Theory of Language Acquisition. Har-	Terry Winggred 1072 Understanding natural language	1000
vard University Press.	Cognitive Psychology, 3(1):1–191.	1200
Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier		
Martinet, Marie-Anne Lachaux, Timothee Lacroix,	Ludwig Wittgenstein. 1953. Philosophical Investiga-	1210
Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	tions. Blackwell.	1211
Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	Martha Woodmansee. 1994. The genius and the copy-	1212
Grave, and Guillaume Lample. 2023. Llama: Open	right: Economic and legal conditions of the emer-	1213
and efficient foundation language models. arXiv	gence of the 'author'. In Martha Woodmansee and	1214
preprint arXiv:2302.139/1.	Peter Jaszi, editors, The Construction of Authorship:	1215
European Union 2023. The artificial intelligence act:	Textual Appropriation in Law and Literature, pages	1216
Regulatory framework for ai systems in the en.	1–20. Duke University Press.	1217
	Wei Xu, Yulia Tsvetkov, and Alan Black. 2022. Ai for	1218
Eline van Dis, Anne-Sophie Bender, Marcel Bonn, and	language learning: Conversational agents and person-	1219
Iris de Bruin. 2023. Chatgpt: Five priorities for re-	alized feedback. Transactions of the Association for	1220
search. Nature, 614:224–226.	Computational Linguistics (TACL), 10:1–15.	1221
1	3	

- Shunyu Yao, Dian Yu, Jeffrey Zhao, Nan Du, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv*:2305.10601.
  - Jason Yosinski, Dario Amodei, and Alec Radford. 2023. Co-editing with ai: The role of large language models in real-time group writing. In *Proceedings of ACL* 2023, pages 3145–3159.
  - Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. In Advances in Neural Information Processing Systems (NeurIPS), volume 32.
- Xia Zhou, Yang Xu, and Feng Liu. 2023. Can aigenerated text be reliably detected? evaluating plagiarism detection tools on gpt-based texts. *arXiv preprint arXiv:2304.08979*.
- Shoshana Zuboff. 2019. The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power. PublicAffairs.

# A Example Appendix

1244 This is an appendix.

1222

1223 1224

1225

1226

1227

1228

1229

1230

1231

1232

1233 1234

1235

1236

1237 1238

1239

1240

1241

1242