Probing the Uniquely Identifiable Linguistic Patterns of Conversational AI Agents

Anonymous ACL submission

Abstract

Considerable effort has been dedicated to detecting machine-generated texts to prevent a situation where the widespread generation of text-with minimal cost and effort-reduces trust in human interaction and factual information online. Our study takes a more refined approach by analysing different Conversational AI Agents (CAA). By constructing linguistic 009 profiles for each AI agent, the aim is to identify the Uniquely Identifiable Linguistic Patterns 011 (UILPs) for each model and to demonstrate the effectiveness of these UILPs in identifying its 012 respective AI agent using authorship attribu-013 tion techniques. Promisingly, we are able to classify AI agents based on their original texts with a weighted F1-score of 96.94%. Further, we can attribute AI agents according to their writing style (as specified by prompts), yield-018 ing a weighted F1-score of 95.84%, which sets 019 the baseline for this task. By employing principal component analysis (PCA) for dimensionality reduction, we achieve a weighted F1-score ranging from 89.25% to 97.83%, and an overall weighted F1-score of 96.93%.

1 Introduction

027

041

Recent advances in deep learning and natural language processing have led to the emergence of conversational AI agents (CAA), hereby referred to as AI agents, which we define as large language models that can generate natural language as a dialogue system. These have been applied in tasks such as question answering (Zhao et al., 2023), fake news detection and abuse detection (Uchendu et al., 2021). The widespread use of AI agents has highlighted the importance of determining the origin of a text (Desaire et al., 2023; Fagni et al., 2021; Mitrović et al., 2023; Fagni et al., 2021; Mitrović et al., 2023; Becker et al., 2023; Islam et al., 2023; Markowitz et al., 2023) and has led to a surge of interest in analysing the linguistic structures within them (Desaire et al., 2023). One noteworthy linguistic aspect that remains unexplored is

the determination of whether AI agents possess any uniquely identifiable linguistic patterns (UILPs).

043

044

045

046

047

050

051

053

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

076

078

079

Our research draws inspiration from the linguistic theories of language identity and linguistic patterns within the compositions of individual authors (Nini, 2023; Coulthard, 2004). Specifically, our study undertakes the task of assessing the validity of the aforementioned theories regarding AI agents. As a result, we have meticulously crafted the UILPs for the following five generative large language models: GPT-4¹, GPT-3.5¹, Text-Curie-001¹, PaLM-2², and LLaMA2-7b³, aiming to ascertain the presence of UILPs. The most effective method to confirm the usefulness of an identified UILP is through validation, a process achievable through the task of authorship attribution. Authorship attribution for AI agents is the ability to ascertain the authorship of texts generated by AI agents (Juola, 2008, 2006; Sari, 2018). By establishing authorship, whether proving or disproving, we can reinforce the theory of distinct linguistic patterns in AI agents. We seek to answer crucial questions about the existence of UILPs in AI agents, the linguistic overlap between various text types generated by these models, and the feasibility of identifying AI agents based on their individual UILPs.

The ability to prove the existence of UILPs provides many benefits such as preventing the harmful use of AI agents (e.g., detecting fake news, hate speech, plagiarism). Additionally, this enables the reuse of the UILP in classification tasks, potentially enhancing classification accuracy. Moreover, this approach is adaptable to the evolving landscape of AI models. We propose a transparent means for linguistic analysis that is more interpretable across different AI agents and forms the central emphasis of this paper.

¹Model details and source: OpenAI's GPT-3.5. (2021). https://www.openai.com/

²Model details and source: Bard: The Language Model for Writing Assistance. (2022). https://www.bardmodel.com/

³Model details and source: LLaMA2-7b: A Large Mul-

Thus far, there has been no investigation on the UILP of AI agents, and there has been only limited comparison of different AI agents and little research indicating if these AI agents can be differentiated from each other based on their linguistic patterns. Moreover, there is a notable absence of analysis of AI agents based on *stylometry*, i.e., the statistical analysis of language often used in the context of forensic linguistics (Rocha et al., 2016). We propose both a feature-based machine learning classifier as well as the use of transformer language models for AI agent classification. The research questions (RQs) we aim to answer in this paper are as follows:

081

097

100

101

102

103

104

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

- RQ1: To what extent can we perform authorship attribution (AA) for AI agents based on their original texts, through the recognition of their UILPs?
- RQ2: Can we attribute text to AI agents through the recognition of UILPs in texts that they generated based on different stylistic prompts?
 - RQ3: How can we measure the linguistic overlap, if any, in outputs from the AI agent when it generates distinct texts?

In addressing the above questions, we have made the following contributions:

• Two new datasets: The first dataset is a collection of original texts created by five AI agents, while the second dataset is an expanded version of the first whereby each text was paraphrased by its respective AI agent according to the following five styles: (a) paraphrased with no specified style, (b) written as a fictitious narrative, (c) written as a tweet, (d) written as a social media blog post and (e) written as an academic article.

 An approach to AI agent attribution based on a Logistic Regression (LR) model trained on linguistic features and a fine-tuned DeBERTa model (He et al., 2021).

 A method for identifying linguistic patterns in the texts generated by the different AI agents based on principal component analysis (PCA).

2 Related Work

The analysis of authorship attribution encompasses two distinct categories: feature-based and large language model-based classification. Feature-based approaches involve creating a specific feature set for that task (Sari, 2018; Juola, 2008) such as multivariate linguistic analysis paired with a traditional machine learning classifier (Abbasi and Chen, 2008). Learning involves the use of language models that can learn the data during the training phase. Newer approaches use pre-trained transformer language models and in some cases, these are combined with linguistic features (Ai et al., 2022; Fabien et al., 2020; Sari, 2018). Approaches using pre-trained transformer language modes have demonstrated superior accuracy with few preprocessing steps. These models significantly outperform traditional models in many cases.

123

124

125

126

127

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

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

Posited by Nini (2023), the Principle of Linguistic Individuality states that at any given moment it is exceedingly improbable for two individuals to possess identical linguistic grammars. This principle is aligned with authorship attribution (Coulthard et al., 2016) which assumes that writings from one author would exhibit greater linguistic similarity than writings from a different author (Burrows, 2002; Anthonissen and Petré, 2019). This theory has not been investigated in the case of AI agents, which is what we sought to achieve in our work.

There has been a central focus on GPT models, with an emphasis on distinguishing between text written by humans and those generated by machines using transformer models (Fagni et al., 2021; Mitrović et al., 2023; Solaiman et al., 2019; Uchendu et al., 2021; Bakhtin et al., 2019; Ippolito et al., 2020), or surface-level linguistic features (Desaire et al., 2023; Markowitz et al., 2023) which have been regarded as a limited analysis when studied individually (Schuster et al., 2020). Other studies have utilised a primarily linguistic approach, analysing words and sentiment to distinguish human and machine-generated text (Markowitz et al., 2023). The limitations of previous approaches, compared to the methodology employed in this paper, become evident when considering their emphasis on distinguishing between human and machinegenerated content.

These studies lack a comparative analysis of various AI agents and rarely incorporate multivariate stylometric analysis in their evaluation, which

tilingual Language Model for Free-Form Editing. (2023). https://www.llama7b.ai/

would better capture the use of AI agents in gener-174 ating texts in other scenarios. Munir et al. (2021) 175 investigated the attribution accuracies of synthetic 176 text using transformer models (XLNet) and prior 177 attribution approaches. Other work has shown that traditional authorship attribution approaches can-179 not fully capture the style of an author when the 180 author is a human vs machine text generation. Hu-181 mans have a wide writing style which means their 182 features and feature usage can differ depending on 183 the text genre (Uchendu et al., 2021).

3 Methodology

187

188

190

191

192

193

194

195

196

199

200

207

208

211

212

213

215

216

3.1 Model Selection

The models used for this project include GPT-3.5, GPT-4, Text-Curie-001 (OpenAI, 2023), PaLM-2 (Anil et al., 2023)¹ and, LlaMA2-7b (Touvron et al., 2023). All of these models are proficient in the natural language generation task with varying levels of sophistication. The Open AI GPT (generative pre-trained transformer)² models used in this paper were all trained using reinforcement learning from human feedback (RLHF) on text data, web pages and books, among others. GPT-4 (OpenAI, 2023) is currently the most optimised model; GPT-3.5 has the same capabilities as GPT-4 but operates on a smaller scale. The Text-Curie-001 model is an older, now deprecated model produced by Open AI.

PaLM-2 (Pathways Language Model)³ developed by Anil et al. (2023) was pre-trained on a large quantity of parallel multilingual corpora, web pages, source code and various other datasets. Proposed by Touvron et al. (2023), LLaMA2-7b (Language Learning and Meaning Acquisition)⁴ was trained on textual data using a standard optimiser and RLHF. We refer the reader to Table 5 in Appendix A for details on each model's size (in terms of the number of learned parameters) and the maximum number of tokens in their output.

3.2 Data collection

Data collection was carried out in two phases. In the first phase, a set of 10 prompts was collated, with each prompt corresponding to a news category on the BBC website⁵ to cover various topics. The 217 specific topic for each prompt was derived from the 218 headline that was most popular at that time within 219 a particular category. The rationale for selecting 220 these article topics was to ensure a diversity of texts 221 within the dataset. Also, the provided prompts did 222 not include harmful or sensitive content therefore, 223 we anticipate the outputted text to be devoid of 224 this material. For instance, within the education 225 category, the most prominent headline pertained 226 to the impact of Covid-19 anxieties on academic 227 studies. Table 6 in Appendix B provides a list of 228 these prompts. An example of the outputs for the 229 prompts in the different prompt styles can be seen 230 in Table 7 in Appendix C. These prompts were 231 given as input to all the AI agents. Data collection 232 occurred through two methods: manual input of 233 prompts in the case of PaLM-2 (through BARD), 234 or by utilising APIs in the case of LLaMA2-7b 235 and the GPT models. For each of the 10 prompts, 236 20 texts were generated. Thus, overall, 200 texts 237 were generated per model, except PaLM-2. The 238 data for PaLM-2 corresponds to only nine queries 239 as the model's responses for one of the 10 queries 240 were inadequate, thus leading to the generation of 241 only 180 texts for this model. This dataset will 242 be referred to as our original data. The data was 243 labelled according to the model used, using labels 244 OG0-OG4 (Original-0 to Original-4). We also used 245 only the GPT-generated data (GPT-3.5, GPT-4 and 246 text-curie-001) from the original data in our analy-247 sis. This dataset, referred to as the GPT data, was 248 labelled according to the model used using integers 249 from GD0-GD2 (GPT data-0 to GPT data-2)in this 250 dataset. 251

The second phase pertains to the collection of stylistic data for only GPT 3.5, 4 and Text-Curie-001. We employed only these three AI agents because they responded effectively to the prompt, while other AI agents produced nonsensical or repeated texts. The stylistic data uses the original data to produce paraphrases of this text in different stylistic genres. Firstly, we asked each model to paraphrase the original text in a general manner, i.e., without specifying a specific style. The model was then asked to paraphrase the original text (from the first phase) in four styles: as an academic paper, as a social media post, as a fictitious narrative and as a tweet. These texts were labelled according to the style with labels ranging from S0 to S4

252

253

254

255

256

257

259

260

261

262

263

264

265

266

¹This model was used via Google's BARD, now known as Gemini (https://gemini.google.com/app)

²Introducing GPT models: https://platform.openai. com/docs/guides/gpt

³PaLM-2: https://ai.google/discover/palm2/

⁴LLaMA: https://ai.meta.com/blog/ large-language-model-llama-meta-ai/

⁵BBC: https://www.bbc.co.uk/

(Style-0 to Style-4) for each stylistic variation in 267 this dataset. For each paraphrasing prompt, 200 268 texts were generated (corresponding to the original 200 texts generated as part of the first phase). 270 In total, there are 1200 texts for each model: the original 200, a version of those 200 that are general 272 paraphrases and 200 for each of the four above-273 mentioned styles. This set of data will be referred 274 to as stylistic data. All datasets were split into train-275 ing and testing sets following a 80:20 partition. No 276 cleaning or preprocessing steps were applied to the data.6

281

289

297

302

303

307

309

311

312

313

314

The process of dataset creation posed a challenge, with certain models generating incoherent texts which were variations of the input text, or texts that were too short or too long. This was due to the absence of predefined constraints during the text generation process. The cohesiveness or semantic soundness of texts is not a primary issue in this work as we aim to focus on contextindependent linguistic features. Model hallucination was not a significant concern for us, as our work primarily concentrated on extracting linguistic features; hence, the content held minimal importance. However, steps to ensure that generated text was reasonable and free of grammatical errors were taken. As previously mentioned, data collection involved either manually inputting prompts or utilising an API. When employing APIs, texts were generated in small batches of 20-50 rows of text data to guarantee that the model produced coherent text data rather than generating random iterations of a single phrase. Lastly, the final datasets were manually assessed to ensure their suitability for the attribution task. This was assessed by ensuring each row contained enough text (more than 10 tokens), a set number of texts per author (dependent on the dataset), avoided repetitive material and created topic-diverse texts. Diverse prompts were employed to ensure this.

3.3 Writeprints as Feature Representation

Abbasi and Chen (2008) proposed the Writeprint feature set: a set of linguistic features for representing the distinctive writing style of each author of interest in an authorship attribution task. The said feature set is largely composed of dynamic features, which are context-dependent, an example of which is the presence of certain word unigrams or bigrams. For example, the presence of the word 315 bigram "yours sincerely" could be indicative of a 316 particular author when writing emails. However, 317 the same author is unlikely to use the same bigram 318 in a different context, e.g., when writing an aca-319 demic article. Thus, to represent an author's writ-320 ing style regardless of context (or textual genre), we 321 extended the original Writeprint to include static 322 features, which are context-independent and are 323 present in a large percentage of texts irrespective 324 of the genre. The extended feature set differs from 325 the original Writeprints in that the former encom-326 passes previously unexplored aspects of a text, such 327 as phonology, morphological irregularities, ellip-328 sis, and omission. Our Extended Writeprint (EWP) 329 is provided in full in Appendix D. These features 330 were extracted from the texts generated by each 331 of the AI agents of interest with the aid of exist-332 ing Python packages, e.g., spaCy (Honnibal et al., 333 2020) and NLTK (Bird, 2006). This results in a 334 unique linguistic profile for each model, which is 335 used in two ways: to determine the most informative features representing the UILP of each of our 337 AI agents of interest (Section 3.4) and to train tradi-338 tional machine learning-based classification mod-339 els to attribute a text to its AI agent (Section 3.5). 340

341

342

344

345

346

347

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

3.4 Analysing the UILP of AI agents

We employed principal component analysis (PCA) to assess the top 100 most informative linguistic features that represent each model (based on its generated texts), as well as the collective top 100 most informative linguistic features. PCA was performed on the standardised feature counts. Subsequently, we quantified the degree of overlap among these top 100 features across the various models. Instead of necessitating the training or retraining of pre-existing language models for attributing texts generated by AI agents. We advocate for a featurebased approach coupled with a machine-learning classifier. The advantage of employing a featurebased approach lies in its efficiency, requiring less time and computational resources. By employing a feature-based approach, we can ensure consistent attribution accuracies regardless of when or by whom the text was generated. This is achieved by the ability to identify distinctive linguistic patterns unique to each AI agent.

We identified unique features for each model based on the most informative features identified by PCA. These unique features were then extracted

⁶Our code and datasets will be made publicly available upon paper acceptance.

from the writeprint of the texts. Authorship attribution was then performed using these uniquely
occurring features.

8 **3.5** Classification Models

369

370

373

375

377

379

390

400

401

402

403

404

405

406

407

408

409

410

411

We cast authorship attribution as a multi-class classification problem, whereby a model takes a given text as input and outputs a label that corresponds to any one of the five AI agents.

A variety of traditional machine learning-based models were trained as classifiers. These include Support Vector Machine (SVM), Random Forest (RF) and Logistic Regression (LR) models. Each of these models was trained on the EWP features described in Section 3.3, using optimised parameter values which were defined through the use of grid search. This allowed us to set a baseline and quantify the extent of any performance improvements. We computed the standard deviation (SD) over five runs. Our results show that the SD in all experiments is low, indicating that the performance scores tend to cluster around the mean. This consistency highlights the stability of our results.

Additionally, we sought to compare the attribution performance with a transformer-based language model (Vaswani et al., 2023) given that transformer models have shown superior performance in classification tasks (Fabien et al., 2020). In this case, we selected the Decoding-enhanced BERT with Disentangled Attention (DeBERTa) model as it has been demonstrated to outperform other transformer models in a variety of tasks (He et al., 2021). Details of the hyperparameters used in training the machine learning and transformer-based language models can be found in Appendix E and F. All experiments were run on Google Colab using the A100 GPU accelerator. Due to the high computational power required to run the DeBERTa model, the results presented are based on a single run.

Prior approaches tend to overlook the identification of distinctive patterns, opting instead for a multivariate dynamic feature extraction technique. Such techniques are text, author and content specific due to dynamic feature selection (Ai et al., 2022; Sari, 2018). The emphasis here lies in discerning unique patterns that can be utilised to identify the AI agents of interest, regardless of the text they produce, with consistent results.

Model	Original	GPT	Stylistic
SVM	93.88	94.17	95.56
RF	96.54	96.67	95.25
LR	96.94	97.50	95.84
DeBERTa	99.11	99.11	88.00

Table 1: Performance Metrics for original, GPT and, Stylistic data: Weighted F1-scores (W-F1) for optimised SVM, LR, RF classifiers (after 5 runs) for all AI agents and a singular run for DeBERTa model

4 Evaluation Results and Discussion

In this section, we present the results of our experiments; for detailed results, including accuracy scores, standard deviation (SD) and weighted F1scores, see Appendix G. 412

413

414

415

416

417

418

419

420

421

422

423

494

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

4.1 Attribution of Original Texts

Table 1 presents the results for authorship attribution based on the original data. The EWP features were extracted from all the texts and the methodology was applied, as described in Section 3.3. From the results, we can see that the optimised DeBERTa model obtained the highest weighted F1-score at 99.11%. However, it is worth noting that the discrepancy in F1-scores across all models is at most merely 5.23% demonstrating competitive performance. When the extended feature set is combined with an optimised ML classifier, the weighted F1score ranges from 93.88% to 96.94%. This demonstrates that each AI agent does have a UILP as we can attribute each model to the correct AI agent with a weighted F1-score of at least 93.88%.

From the results in Table 1, we can see that DeBERTa has the highest weighted F1-score at 99.11%. In this experiment, the discrepancy in F1-scores across all models is 4.94%. Since all the compared models are OpenAI-engineered, it is reasonable to anticipate that they exhibit similar linguistic patterns in their generated texts hence the lower F1-scores across all experiments. This model displays an impressively competitive performance, with the optimised LR model having a weighted F1-score of 97.50%, which is only a 1.61% drop when compared to a fine-tuned DeBERTa model.

4.2 Attribution of Stylistic Texts

Apart from the attribution of the original data, we also investigated the attribution of stylistic text; this can be considered as cross-genre attribution as we examine the attribution performance for all AI agents on different stylistic data.

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

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

495

496

497

498

499

The results of attribution basaed on the stylistic dataset for GPT models are presented in Table 1. As aforementioned, since all models are OpenAIengineered we expect some linguistic commonalities across different genres of text. Here we attempt to attribute all texts (original, paraphrase, social media posts, tweets, academic articles and fictitious narratives) to their respective AI agent. The results support the notion of models demonstrating a UILP in their generated stylistic texts as well as the notions posited by Juola (2008); Sari (2018); Coulthard (2004) who suggested that these UILPs can be identified across different textual genres, but lower results can be expected when performing cross-genre attribution. This accounts for the 11.11% reduction in the weighted F1-score when comparing the original data to the stylistic data using optimised DeBERTa models. One can observe a 1.1% weighted F1-score drop when using an optimised LR model and a 19.62% drop when comparing the performance of the default LR model on the original data with the stylistic one. These results indicate that each AI agent has a distinct UILP for the stylistic texts, further affirming the idea that performance decreases across genres due to varying linguistic patterns (Stamatatos, 2016).

To conclude, we can recognise each AI agent, regardless of the text's style, with the highest weighted F1-score achieved at 95.84%.

4.3 Attribution based on PCA of features

In this section, we identify the top 100 most informative linguistic features across all AI agents and the top 100 most informative linguistic features for each AI agent. We then assess the extent to which attribution can be performed based on these features, for both original and stylistic data; the full results can be seen in Appendix G in Tables 16 and 17.

For all the original data, we extracted our Extended Writeprint features. Subsequently, we conducted PCA to identify the top 100 most informative linguistic features across the entire dataset. Attribution was carried out using these selected top 100 features; the accuracy of each model was then computed. When performing attribution using only the top 100 most informative linguistic features, we found that Text-Curie-001 has the highest weighted F1-score for any model and has a self-identifying weighted F1-score of 98.77%. LLaMA2-7b ob-

Weighted F1-Score					
AI Agent	GPT-3.5	GPT-4	LLaMA2-7b	PaLM-2	Text-curie001
GPT-3.5	80.49	82.5	78.05	88.89	90.84
GPT-4	78.16	87.5	72.94	83.54	90.91
LLaMa2-7b	65.63	77.16	66.67	75	94.74
PaLM-2	82.05	84.62	86.43	79.49	97.3
Text-Curie-001	98.77	95.24	98.77	97.56	98.77
Overall	81	85.42	80.45	85.36	94.38

Table 2: Weighted F1-scores for models based on their top 100 most informative linguistic features extracted from the EWP using PCA analysis. Attribution was performed for each model and then for the entire original dataset using an optimised Logistic Regression model

	Unique Features	PCA100
AI Agent	W-F1	W-F1
GPT-3.5	89.25	86.17
GPT-4	95.50	87.18
LLaMa2-7b	97.83	100
PaLM-2	93.17	90.00
Text-Curie-001	96.97	97.56
Average	96.93	91.81

Table 3: Weighted F1-score for each AI agent using an optimised LR model. LR was performed on the top 100 linguistic features outputted from the PCA (PCA100) and on each model's unique features

tained the lowest performance, with a weighted F1-score of 66.67% when identified using its individual top 100 feature set.

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

These results support the theory of linguistic individuality (Nini, 2023) as the AI agents do not have identical grammars even though the training material, methods, the developers are the same or similar. This can be seen explicitly in the analysis of the Open AI GPT models, whereby the F1-score varies from 96.93% to 88.25%, showing a slight discrepancy of 8.68%. Each AI agent struggles to distinguish itself when using its own top 100 most informative features. However, this is due to the substantial overlap in these features, as demonstrated in Appendix H. On average, they share more than 50% of their top 100 features with another AI agent. This clarifies why, in Table 2, we observe an absence of a distinct pattern in AI agents' ability to identify themselves through their top 100 features.

There are noticeable instances of misclassification concerning GPT-3.5 and GPT-4. The relatively poorer attribution of GPT-3.5 and GPT-4 can be explained by the fact that both models are OpenAIengineered, have similar training processes and serve the same purpose. GPT-4 is an improvement that builds upon the existing capabilities of GPT-3.5 (OpenAI, 2023).

Further investigation was performed to deter-527 mine if AI agents can be identified based on their 528 unique feature sets. We conducted a comparison of 529 the top 100 features across all AI agents and identified features unique to each model. After obtaining 531 the set of distinctive features for each model, we 532 moved on to the original dataset containing approx-533 imately 300 features. For each model, we exclusively extracted the features that were unique to that model. For example, during the attribution for GPT-4, we isolated features that were uniquely 537 associated with GPT-4 in its top 100 most informa-538 tive features. GPT models exhibited greater mor-539 phological diversity among these unique features 540 compared to LLaMa7-2b and PaLM-2. In contrast, 541 the unique feature sets of LLaMa7-2b and PaLM-2 542 predominantly included function words. These spe-543 cific features were then extracted for every model from the comprehensive set of 300 features. Subse-545 quently, we performed attribution analyses for each model based on this refined set of features. For example, we identified and extracted all features that were uniquely identified in the top 100 features. 549 We then extracted GPT-4's unique features for all 550 other AI agents and attempted attribution using this 551 unique feature set. The differences in results were 552 significant: the weighted F1-scores ranged from 553 86.17% to 100% when using optimised hyperpa-554 rameters. The results support the theory that when 555 investigating an AI agent's inherently unique features, one can attribute each AI agent with greater 557 success. Further results on the attribution success 558 for each model can be seen in Table 17 in Appendix G.

The subsequent phase involved conducting PCA for each model and extracting the most informative top 100 features. Following this, we attempted the attribution for all models using these top 100 features, and the outcomes are presented in Table 2 (full results are in Appendix G in Table 18). The results indicate that only LLaMA2-7b could successfully self-identify as the most similar AI agent based on these features. A more in-depth linguistic examination of these features revealed that PCA features are predominantly comprised of static features, defined as context-independent and frequently occurring attributes. Furthermore, the diagrams in Figure ?? in Appendix H illustrate substantial feature overlap among different models when analysing 300 features. However, as the features are reduced to find the most unique ones,

562

563

565

568

569

570

571

573

574

575

577

there is a noticeable drop in overlap; see Figure ?? and Figure ?? in Appendix H. This supports the theory of Linguistic Uniqueness (Nini, 2023) and the existence of a UILP as it is evident that each model has a set of features that it does not share with the others. These results pertain solely to the original data, with accuracies and weighted F1-scores obtained using the LR algorithm.

578

579

580

581

582

583

584

585

586

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

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

4.4 Linguistic Analysis

For each model, rather than extracting all features specified in the EWP, we reduced the feature set to include only linguistic features associated with each specific linguistic category (details of the features and their categories are provided in Table 8 in Appendix D). Attribution was subsequently conducted for the original data using these refined feature sets. The results of this classification are presented in Table 4.

Individual accuracy scores for each linguistic category and the overall dataset were computed. Tagging and *n*-gram categories achieved the highest weighted F1-score among all ML classifiers. This can be attributed to several factors. Firstly, the presence of over 100 different part-of-speech and dependency tags as well as *n*-grams adds a significant level of linguistic diversity to the dataset. This category also encompassed the labels for different sentence types e.g. the count of passive sentence constructions. Furthermore, research has established that AI agents employ repetitive sentence structures to maintain cohesiveness, and this makes tags a particularly identifiable linguistic structure (Mitrović et al., 2023; Markowitz et al., 2023). Forensic research has also continually highlighted *n*-grams as an extremely identifiable linguistic feature in authorship attribution (Sari, 2018).

It is still important to note that there is greater variability in the weighted F1-scores with the highest F1-score for any classifier being 91.79% (for RF) and the lowest at 73.30% (for LR) creating a difference of 18.54% between classifiers. For all result details see Table 19 in Appendix G.

In comparing the linguistic features of two texts from two AI agents in a qualitative assessment, distinctive patterns emerge, suggesting potential variations in author style and expression. Text T1 refers to the text outputted by PaLM-2 and Text T2 is the one from GPT-3.5 as seen in Table 20 in Appendix I. Both texts share some features, indicating commonalities in sentence structures and grammatical

	Weighted F1-Score		
	SVM	RF	LR
Word Lists	89.73	88.30	84.31
Symbols	83.70	91.26	78.34
Tags	87.66	91.34	89.77
Syntax	76.21	77.37	73.30
Semantics	79.06	84.76	79.38
Lexical	90.33	91.33	89.86
<i>n</i> -grams	91.79	90.31	90.18

Table 4: Accuracy and Weighted F1-scores for Individual Linguistic Categories in Attribution on the Original Data

constructions. Despite sharing five common POS and dependency tags, both texts display between 8 to 21 unique dependency and POS tags, signifying a common syntactic foundation with specific linguistic constructions that differentiate their styles. Notably, T1 employs comparative adjectives, while T2 includes modal verbs, showcasing distinctive choices that may reflect variations in tone or style. In the realm of authorship attribution, these linguis-636 tic differences underscore the potential for the texts to be perceived as the work of different authors, as individual writing habits and preferences become apparent through their unique linguistic patterns. See Tables 21 and 22 in Appendix I for further details on the feature counts and for the subset of 642 features extracted for this assessment.

4.5 Error Analysis

631

641

647

649

650

652

656

657

662

We conducted an error analysis of exclusively the original dataset due to its inclusion of all AI agents, its utilisation of the entire EWP. Also, its attribution involved the application of both an ML model (in this case Logistic Regression) and DeBERTa (He et al., 2021).

In Appendix J we see the classification outputs from LR and DeBERTa. The DeBERTa model exhibited a total of three misclassifications. All three instances involved GPT-4 data being incorrectly labelled as GPT-3.5. The explanation for this lies in the fact that both models undergo the same training process, both are OpenAI authored and additionally, GPT-3.5 is the predecessor of GPT-4. The LR model displayed a total of 9 misclassifications. There is one instance of GPT-4 misclassification as GPT-3.5, a mistake made by DeBERTa. All other misclassifications were of LLaMa2-7b; this AI agent was incorrectly classified as Text-Curie-001, GPT-3.5 and PaLM-2. Based on a linguistic

assessment of the misclassified data, we see that the instances of misclassified LLaMA2-7b data exhibited stylistic variations. These texts tended to be longer on average and had more morphological variation which explains the misclassifications as Text-Curie-001 and GPT-3.5. Nevertheless, both models exhibited a minimum number of errors, leading us to consider them insignificant. Further fine-tuning and conducting additional linguistic analysis could help mitigate these misclassifications.

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

704

705

706

707

709

5 **Conclusion and future work**

In our study, we have addressed three key research questions. Firstly, we have confirmed the presence of Uniquely Identifiable Linguistic Patterns (UILPs) in conversational AI agents. This is supported by high accuracy in attribution for both original and stylistic data, with weighted F1-scores ranging from 93.88% to 96.96% when utilising the Extended Writeprint (EWP) and traditional machine learning-based classifiers. We also demonstrate similar performance when using a fine-tuned DeBERTa model, achieving a 99.11% weighted F1-score. Our results demonstrate that traditional machine learning-based models can obtain competitive attribution performance compared to a finetuned DeBERTa model when utilising the EWP for classification. Through PCA analysis, we explored the attribution of AI agents based on their UILPs. Our results show that the combination of our EWP and RF classification effectively supports crossgenre attribution, with weighted F1-scores ranging from 94.17% to 97.50% for the stylistic data. This affirms the principle of linguistic individuality in AI agents, showcasing their UILPs. These findings validate the existence of UILPs in AI agents and offer valuable insights into their distinctive linguistic patterns, with potential applications in digital forensics, detecting fake news and plagiarism. Future work will look to improve both the datasets introduced in this paper by expanding the size and scope of the stylistic prompts. We seek to perform a fine-grained linguistic analysis of a larger set of AI agents cross-lingually.

Limitations

In our study, text generation using various APIs that 710 make our AI agents of interest accessible proved 711 to be a time-intensive process, limiting the volume 712 of prompts that could be supplied and thus the text 713

that can be generated. Additionally, certain models 714 imposed output constraints. For instance, in the 715 case of PaLM-2, we resorted to manually inputting 716 prompts into BARD due to the unavailability of 717 the API, which was a time-consuming endeavour. 718 Furthermore, some AI agent outputs did not pro-719 duce cohesive texts (in the case of LLaMA2-7b), or 720 produced very short texts (in the case of Text-Curie-721 001). Further, only a set of three text genres were investigated: academic articles, fictitious narratives, 723 and tweets and social media posts (the latter two falling under the same genre). To perform cross-725 genre authorship attribution we must expand this scope to cover a wider array of genres as well as 727 investigate at different levels of formality. Further-728 more, a study into misclassified instances must be conducted to identify patterns or determine if there 730 is a specific type of error being made by the model. 731

32 Ethics Statement

733

734

737

738

739

740

741

742

743

744

745 746

747

748

751

753

755

757

759

760

For this study, the data was sourced from various AI agents, and human involvement was not required. The dataset does not contain any harmful or sensitive content. As there was no human participation and no collection of personal data, ethics review was not necessary.

References

- Ahmed Abbasi and Hsinchun Chen. 2008. Writeprints: A stylometric approach to identity-level identification and similarity detection in cyberspace. *ACM Trans. Inf. Syst.*, 26(2):7:1–7:29.
- Bo Ai, Yuchen Wang, Yugin Tan, and Samson Tan. 2022. Whodunit? learning to contrast for authorship attribution.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report.
- Lynn Anthonissen and Peter Petré. 2019. Grammaticalization and the linguistic individual: new avenues in lifespan research. *Linguistics Vanguard*, 5(s2):20180037.
- Anton Bakhtin, Sam Gross, Myle Ott, Yuntian Deng, Marc'Aurelio Ranzato, and Arthur Szlam. 2019.
 Real or fake? learning to discriminate machine from human generated text. arXiv.
- Jonas Becker, Jan Philip Wahle, Terry Ruas, and Bela Gipp. 2023. Paraphrase detection: Human vs. machine content.

- Steven Bird. 2006. Nltk: the natural language toolkit. 762 In Proceedings of the COLING/ACL on Interactive 763 presentation sessions, COLING-ACL '06, pages 69-764 72, Stroudsburg, PA, USA. Association for Computa-765 tional Linguistics. 766 John F. Burrows. 2002. 'delta': a measure of stylis-767 tic difference and a guide to likely authorship. *Lit.* 768 Linguistic Comput., 17:267–287. 769 Malcolm Coulthard. 2004. Author Identification, Idi-770 olect, and Linguistic Uniqueness. Applied Linguis-771 tics, 25(4). 772 Malcolm Coulthard, Alison Johnson, and David Wright. 773 2016. An introduction to forensic linguistics: Lan-774 guage in evidence. Routledge. 775 Heather Desaire, Aleesa E. Chua, Madeline Isom, Ro-776 mana Jarosova, and David Hua. 2023. Chatgpt or 777 academic scientist? distinguishing authorship with 778 over 99off-the-shelf machine learning tools. 779 Maël Fabien, Esau Villatoro-Tello, Petr Motlicek, and 780 Shantipriya Parida. 2020. BertAA : BERT fine-781 tuning for authorship attribution. In Proceedings 782 of the 17th International Conference on Natural Lan-783 guage Processing (ICON), pages 127-137, Indian 784 Institute of Technology Patna, Patna, India. NLP As-785 sociation of India (NLPAI). 786 Tiziano Fagni, Fabrizio Falchi, Margherita Gambini, An-787 tonio Martella, and Maurizio Tesconi. 2021. Tweep-788 Fake: About detecting deepfake tweets. PLOS ONE, 789 16(5):e0251415. 790 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and 791 Weizhu Chen. 2021. Deberta: Decoding-enhanced 792 bert with disentangled attention. 793 Matthew Honnibal, Ines Montani, Sofie Van Lan-794 deghem, and Adriane Boyd. 2020. spaCy: Industrial-795
- D. Ippolito, D. Duckworth, C. Callison-Burch, and D. Eck. 2020. Automatic detection of generated text is easiest when humans are fooled. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1808–1822. Association for Computational Linguistics.

797

798

799

800

801

802

803

804

805

806

807

809

strength Natural Language Processing in Python.

- Niful Islam, Debopom Sutradhar, Humaira Noor, Jarin Tasnim Raya, Monowara Tabassum Maisha, and Dewan Md Farid. 2023. Distinguishing human generated text from chatgpt generated text using machine learning.
- Patrick Juola. 2006. Authorship attribution. *Found. Trends Inf. Retr.*, 1(3):233–334.
- Patrick Juola. 2008. Authorship attribution. Founda-
tions and Trends® in Information Retrieval, 1:233–
334.810811812

David M Markowitz, Jeffrey Hancock, and Jeremy Bailenson. 2023. Linguistic markers of inherently false ai communication and intentionally false human communication: Evidence from hotel reviews.

813

814

815

817

819

821

822

823 824

825

826

827

829

831

834

836

837

842

845 846

847

848

851

852

854

859

861 862

863

865

866

- Sandra Mitrović, Davide Andreoletti, and Omran Ayoub. 2023. Chatgpt or human? detect and explain. explaining decisions of machine learning model for detecting short chatgpt-generated text. *ArXiv*.
- Shaoor Munir, Brishna Batool, Zubair Shafiq, Padmini Srinivasan, and Fareed Zaffar. 2021. Through the looking glass: Learning to attribute synthetic text generated by language models. In *Proceedings of the* 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 1811–1822, Online. Association for Computational Linguistics.
 - Andrea Nini. 2023. A Theory of Linguistic Individuality for Authorship Analysis. Elements in Forensic Linguistics. Cambridge University Press.
- OpenAI. 2023. Gpt-4 technical report.
 - Anderson Rocha, Walter J Scheirer, Christopher W Forstall, Thiago Cavalcante, Antonio Theophilo, Bingyu Shen, Ariadne RB Carvalho, and Efstathios Stamatatos. 2016. Authorship attribution for social media forensics. *IEEE transactions on information forensics and security*, 12(1):5–33.
- Yunita Sari. 2018. *Neural and non-neural approaches to authorship attribution*. Ph.D. thesis, University of Sheffield, UK. British Library, EThOS.
- Tal Schuster, Roei Schuster, Darsh J. Shah, and Regina Barzilay. 2020. The limitations of stylometry for detecting machine-generated fake news. *Computational Linguistics*, 46(2):499–510.
- Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, Miles McCain, Alex Newhouse, Jason Blazakis, Kris McGuffie, and Jasmine Wang. 2019. Release strategies and the social impacts of language models.
- Efstathios Stamatatos. 2016. Authorship verification: A review of recent advances. *Research on computing science*, 123:9–25.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.
- Adaku Uchendu, Zeyu Ma, Thai Le, Rui Zhang, and Dongwon Lee. 2021. TURINGBENCH: A benchmark environment for Turing test in the age of neural text generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2001–2016, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. 868

869

870

871

872

873

874

875

876

877

878

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models.

10

Appendix A Conversational AI Agents Model breakdown

Model	Creator	Size	# Tokens
GPT-4	OpenAI	1.7T	8192
GPT-3.5	OpenAI	175B	4097
Text-Curie-001	OpenAI	6.7B	2049
PaLM-2	Google		8192
LLaMA2-7b	Meta	7B	2048

Table 5: Comparison of AI agents based on their size in terms of the number of parameters (unknown for PaLM-2) and the maximum number of tokens in their output (# Tokens)

	NT	Descent
	News category	Prompt
1	Cost of living	Write me an essay on rising house prices in 2023
2	Climate	Write me an essay on what the UK can do to reduce carbon emissions
3	Tech	Write me an essay on Facebook's transfer of European data to US servers
4	Politics	Write me an essay on the UKs ministerial code
5	Culture	Write me an essay on the Glastonbury festival in the UK
6	Science	Write me an essay on the womb condition adenomyosis
7	Family and Education	Write me an essay on how Covid anxiety is affecting students and exams
8	War in Ukraine	Write me an essay on Zelensky's contribution in the war in Ukraine
9	UK	Write me an essay on the Madeleine McCann police search in Portugal
10	World	Write me an essay on the Star Wars original studios

Appendix B Prompts for AI agents

Table 6: The prompts used to collect AI agent-generated texts. All original texts were produced all by AI agents: GPT-3.5, GPT-4, Text-Curie-001, PaLM-2 and LLaMA2-7b. In contrast, only GPT-3.5, GPT-4 and Text-Curie-001 were used in generating texts according to stylistic variations (paraphrase, as a social media post, as a tweet, as an academic article and as a fictitious narrative

Appendix C Data Examples

Style	Text
	In 2023, a pervasive and pressing issue has gripped households across the globe —
Original	the cost of living crisis. As economic landscapes evolve and global conditions shift,
Oligiliai	individuals and families find themselves grappling with the ever-increasing burden of
	expenses.
	In 2023, households worldwide are confronted with a widespread and urgent problem
Daraphrasa	- the crisis of rising living costs. As economic environments undergo changes and
Farapinase	global circumstances shift, individuals and families are contending with a growing
	weight of expenditures.
	Feeling the pinch of the 2023 cost of living crisis! It's a tough
	reality we're all navigating together. From skyrocketing prices to tightening budgets,
Social media	the struggle is real. Let's open up the conversation and share our tips on stretching
post	those hard-earned pounds, supporting one another through these challenging times.
	#CostOfLivingCrisis #CommunitySupport #BudgetingTips
	#TogetherWeCanHandleIt
	The cost of living is soaring in 2023, and it's hitting hard! From groceries to
Tweet	gas, it's becoming more challenging to make ends meet. Let's keep
Iweet	the conversation going on how we can tackle this #CostOfLivingCrisis
	together and support those facing financial strain. #Economy2023
	The year 2023 has witnessed a growing concern worldwide regarding the
	cost of living. The term "cost of living" encompasses the expenses individuals
Acadomic	and families incur to maintain a reasonable standard of living, including housing,
Papar	food, transportation, healthcare, education, and other essentials. This
rapei	essay seeks to provide a comprehensive analysis of the cost of living crisis
	in 2023, focusing on its underlying causes, economic implications, and potential
	policy measures to mitigate its effects.
	In the year 2023, as the calendar pages turned, people across the nation
Fictitious	found themselves entangled in a relentless and unforgiving cost of living
narrative	crisis. The once-stable balance of life, as they knew it, had been upended,
	and every aspect of their daily existence was impacted.

Table 7: The GPT-3.5 output for the prompt "*Write me a <stylistic_text> on the cost of living crisis in 2023*", where *<stylistic_text>* is replaced by one of paraphrase, social media post, tweet, academic article and fictitious narrative

Category	Feature	Description
		Word length
	Token-based	Sentence length
		Average sentence count, Average word count
	Character based	Upper- and lower-case distribution
	Character-based	Digit frequency
	Word length distribution	One to ten plus letters
Lexical	Top n-grams	Top 50 occurring tri and bi grams
	Special characters/punctuation	Frequency counts
	Vocabulary richness	Type-token ration (TTR)
	vocabulary fieliness	Text repetitiveness rate (TRR)
	Hapax Legomena	Frequency counts
	Clipping	Process of shortening words at any word boundary:
	Chipping	e.g., "Advertisement" to "Ad"
		Part-of-Speech (POS) tags
	Tagging	Dependency tags
		Ellinsis: e.g. [full sentence] "I like coffee and she likes tea" to
		[elliptical sentence] "I like coffee and she"
	Term replacement/omission	[emptical sentence] Tinke conce, and she
	Term replacement/ofmission	Substitutions: e.g. [full sentence] "John went to the store.
		John bought back milk" to [substituted sentence] "John went to the store.
		He bought back milk"
		Irregular patterns:
		- Present participle form
Syntactic		- Plural forms
		- Past tense form
	Morphological Variation	- Past participle form
		- Plural form (-ies, -ves, es)
		- Possessive form
		- Comparative and Superlative form
		- Singular form (-y, -o)
		Simple, Complex, Compound
	Sentence types	Declarative, Interrogative, Exclamatory,
		Imperative, Conditional, Comparative, Passive
Semantic	Sentiment scores	
Semantic	Synonym/Homonym counts	
		Alliteration
	Phonetic	Assonance
Other		Consonance
	Word lists	Function words
		Acronyms/Slang

Appendix D The Extended WritePrint

Table 8: The Extended WritePrint (EWP). This feature set consists of static (context-independent) and dynamic (context-dependent) features

Appendix E Hyperparameter settings for the DeBERTa model

Hyperparameter	Amended value
num_train_epochs	6
train_batch_size	16
eval_batch_size	16
gradient_accumulation_steps	4
n_gpu	-1
max_seq_length	512
class_weight	Custom labels specified
early_stopping_patience	2
early_stopping_delta	0.01

Table 9: The hyperparameters used in training the DeBERTa model (He et al., 2021)

Appendix F Hyperparameter settings for the traditional machine learning-based classification models

Hyperparameter	Amended value
max_depth	None
<pre>min_samples_leaf</pre>	1
<pre>min_samples_split</pre>	5
n_estimators	300
class_weights	Balanced

Table 10: The hyperparameters used in training the Random Forest classifier

Hyperparameter	Amended value
С	10
penalty	12
solver	liblinear

Table 11: The hyperparameters used in training the Logistic Regression classifier

Hyperparameter	Amended value
С	0.1
kernel	linear

Table 12: The hyperparameters used in training the Support Vector Machine classifier

883

884

885

Appendix G Complete Results

ML Model	Accuracy	W-F1	SD
SVM	93.87	93.88	0.00
RF	96.54	96.54	0.37
LR	96.94	96.94	0.00
DeBERTa	99.11	99.11	_

Table 13: Performance Metrics for Original Data Attribution: the average Accuracy, Weighted F1-score (W-F1) and Standard deviation Scores for optimised SVM, LR, RF classifiers (after 5 runs) for all AI agents and a singular run DeBERTa model

ML Model	Accuracy	W-F1	SD
SVM	94.11	94.17	0.00
RF	96.67	96.67	0.00
LR	97.50	97.50	0.19
DeBERTa	99.11	99.11	—

Table 14: Performance Metrics for the Attribution of all GPT datasets: the average Accuracy, Weighted F1-score (W-F1) and Standard deviation Scores for optimised SVM, LR, RF classifiers (after 5 runs) for GPT-4, GPT-3.5 and, Text-Curie-001 and fine-tuned DeBERTa

ML Model	Accuracy	Weighted F1	SD
SVM	95.56	95.56	0.00
RF	95.25	95.24	0.25
LR	95.83	95.84	0.00
DeBERTa	88.00	88.00	_

Table 15: Performance Metrics for the Attribution of the Stylistic data: the average Accuracy, Weighted F1-score (W-F1) and Standard deviation Scores for optimised SVM, LR, RF classifiers (after 5 runs) for GPT-4, GPT-3.5 and, Text-Curie-001

AI agent	Accuracy	W-F1	SD
GPT-3.5	91.60	89.25	0.03
GPT-4	97.63	95.50	0.01
LLaMA2-7b	100	97.83	0.00
PaLM-2	95.35	93.17	0.01
Text-Curie-001	100	96.97	0.00
All	96.93	96.93	0.02

Table 16: Results of attribution using an optimised LR model trained on the top 100 most informative linguistic features extracted using PCA across all datasets

AI agent	Accuracy	W-F1	SD
GPT-3.5	86.42	86.17	0.00
GPT-4	86.08	87.18	0.00
LLaMA2-7b	93.34	100	0.02
PaLM-2	94.74	90.00	0.00
Text-Curie-001	98.77	97.56	0.00
All	91.84	91.81	0.00

Table 17: Accuracy and weighted F1-score for each AI agent when performing authorship attribution using only their unique features

	Accuracy and Weighted F1-Score									
AI agent	GPT	Г-3.5	GP	'T-4	LLaM	A2-7b	PaL	M-2	Text-C	urie-001
GPT-3.5	80.52	80.49	82.50	82.50	78.06	78.05	88.89	88.89	90.85	90.84
GPT-4	78.16	78.16	87.50	87.50	72.95	72.94	83.54	83.54	90.91	90.91
LLaMA2-7b	65.64	65.63	77.16	77.14	66.67	66.67	75.00	75.04	94.75	94.74
PaLM-2	82.05	82.05	84.67	84.62	86.42	86.43	79.49	79.49	97.31	97.30
Text-Curie-001	98.77	98.77	95.24	95.24	98.77	98.77	97.56	97.56	98.79	98.77
Overall	81.63	81.00	85.71	85.42	81.12	80.45	85.01	85.36	94.39	94.38

Table 18: Table displaying accuracy and weighted F1-scores for models based on their top 100 most informative linguistic features extracted from the EWP using PCA analysis. Attribution was performed for each model and then for the entire original dataset using an optimised Logistic Regression model

		Accuracy and Weighted F1					
	SVM		R	RF		LR	
Word lists	89.80	89.73	88.72	88.30	84.18	84.31	
Symbols	84.18	83.70	91.33	91.26	78.57	78.34	
Tags	87.76	87.66	91.33	91.34	89.80	89.77	
Syntax	75.51	76.21	77.04	77.37	72.96	73.30	
Semantic	79.08	79.06	84.69	84.76	79.59	79.38	
Lexical	90.51	90.33	91.33	91.33	89.80	89.86	
N-gram	91.84	91.79	90.31	90.31	90.31	90.18	

Table 19: Accuracy and Weighted F1-scores for Individual Linguistic Categories in Attribution on the Original Data

Appendix H PCA visualisations

887

Venn Diagrams: Original Data PCA analysis



Figure 1: Overlap for the top 100 most informative linguistic features extracted based on our EWP using PCA for all AI agents. Classification results are in Table 2

Appendix I Linguistic analysis

Conversational AI Agent	Output
	Workers choosing between warmth and a full refrigerator. Dreams put on hold,
	replaced by the daily grind of survival. The cost of living crisis is a call to action.
D-IM2	It demands solutions, not platitudes. It requires bold leadership, targeted
PaLM-2	support, and a commitment to building a fairer, more resilient society. This isn't just a
	winter chill; it's a systemic squeeze. And until we collectively address its root causes,
	the UK's future risks being defined not by prosperity, but by the struggle to merely exist.
	The United Kingdom is currently contending with a profound cost of living crisis,
	as citizens grapple with the escalating prices of essential goods and services.
	From surging energy costs to rising food and housing expenses, the financial strain on
	households has become palpable. This crisis not only impacts individual wallets
	but also raises concerns about broader economic inequality. Many individuals and families
GPT-3.5	are forced to reassess their budgets and make difficult choices to navigate
	through these challenging times. As the cost of living continues to rise, policymakers
	face the imperative of implementing effective strategies to alleviate the burden on citizens
	and foster economic resilience. The cost of living crisis in the UK is a pressing issue that
	demands thoughtful and comprehensive solutions to ensure the well-being of the
	population.

Table 20: The GPT-3.5 (Text one (T1)) and PaLM-2 (Text two (T2)) output for the prompt "Write me a write a short paragraph on the cost of living crisis in the UK"

Linguistic Features	Feature
Average word length	5.125 charaters per word
Average sentence length	14.9 words per sentence
Type-token ratio	0.678
Text repetitiveness rate	0.32
Character unigram	'e' (occurs 98 times), 't' (occurs 78 times), 's' (occurs 58 times),
Character unigram	'i' (occurs 54 times), 'n' (occurs 54 times)
Character bigram	'th' (occurs 45 times), 'es' (occurs 38 times), 'nt' (occurs 30 times),
Character Digram	'in' (occurs 28 times), 'er' (occurs 28 times)
Character trigram	'the' (occurs 23 times), 'ing' (occurs 16 times), 'ion' (occurs 14 times),
	'ent' (occurs 13 times), 'ndi' (occurs 9 times)
Santanaa tuna	Simple (3), compound (2), complex (2), declarative (6),
Sentence type	passive (1), exclamatory (1)
	Top 5 POS Tags:
	NN (Noun, singular or mass), IN (Preposition or subordinating conjunction),
	JJ (Adjective), VBZ (Verb, 3rd person singular present), DT (Determiner)
	Top 5 Dependency Tags:
	nsubj (Nominal subject), ROOT (Root of the clause), prep (Prepositional modifier),
POS and Dependency tags	pobj (Object of preposition), det (Determiner)
	Number of shared POS tags: 5
	Number of shared dependency tags: 5
	Number of Unique POS tags: 13
	Number of Unique dependency tags: 8

Table 21: Subset of features extracted from the GPT-3.5 prompt (T1) (as specified in Table 20

Feature count
5.276 characters per word
16 words per sentence
0.607
0.392
'e' (occurs 50 times), 't' (occurs 43 times), 's' (occurs 37 times),
'r' (occurs 30 times), 'i' (occurs 28 times)
'th' (occurs 24 times), 'es' (occurs 23 times), 'ti' (occurs 21 times),
'in' (occurs 21 times), 're' (occurs 18 times)
'the' (occurs 14 times), 'ion' (occurs 13 times), 'ing' (occurs 12 times)
'ent' (occurs 10 times), 'tio' (occurs 9 times)
Simple (5), compound (3), complex (4), declarative (12), passive (1)
Top 5 POS Tags: NN (Noun, singular or mass), VBZ (Verb, 3rd person singular present), IN (Preposition or subordinating conjunction), DT (Determiner), JJ (Adjective) Top 5 Dependency Tags: nsubj (Nominal subject), ROOT (Root of the clause), prep (Prepositional modifier), pobj (Object of preposition), det (Determiner) Number of shared POS tags: 5 Number of shared dependency tags: 5 Number of Unique POS tags: 21 Number of Livigue dependency tags: 8

Table 22: Subset of features extracted from the PaLM-2 prompt (T2) (as specified in Table 20

)

Appendix J Error Analysis



(a) Confusion matrix for the attribution of the Original data using Logistic Regression



(a) Confusion matrix for the attribution of the Original data using DeBERTa

Figure 3: [Key: 0: Text-Curie-001; 1: GPT-3.5; 2: GPT-4; 4: LLaMa2-7b: 4: PaLM-2]