An Information-Theoretic Approach to Analyze NLP Classification Tasks

Anonymous ACL submission

Abstract

Understanding the contribution of the inputs on the output is useful across many tasks. This work provides an information-theoretic frame-004 work to analyse the influence of inputs for text classification tasks. Natural language processing (NLP) tasks take either a single or multiple text elements to predict an output variable. Each text element has two components: the semantic meaning and a linguistic realization. Multiple-choice reading comprehension (MCRC) and sentiment classification (SC) are selected to showcase the framework. For 013 MCRC, it is found that the relative context influ-014 ence on the output reduces on more challenging datasets. In particular, more challenging con-016 texts allows greater variation in the question complexity. Hence, test creators need to care-017 fully consider the choice of the context when designing multiple-choice questions for assessment. For SC, it is found the semantic meaning of the input dominates compared to its linguistic realization when determining the sentiment.

1 Introduction

034

038

040

Natural Language Processing (NLP) requires machines to understand language to perform a specific task (Chowdhary and Chowdhary, 2020). NLP tasks take a single (e.g. summarization (Widyassari et al., 2022), sentiment classification (Wankhade et al., 2022), machine translation (Stahlberg, 2020)) or multiple (e.g. reading comprehension (Baradaran et al., 2022), question generation (Kurdi et al., 2020)) text elements at the input and return a specific output. Each input text element can further be partitioned into its semantic content and the linguistic realization. Semantic refers to the inherent meaning while the linguistic realization is the specific wording to present the meaning. There are several possible linguistic realizations for any semantic content. Therefore, for all NLP tasks, the output variable has contributions from at least two components: the semantic

meaning of the element and the specific linguistic realization. Here, *element* refers to a specific input that is formed of exactly two *components*. 042

043

044

045

047

048

050

051

053

054

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

081

We analyze the relative sensitivity of the output variable to each of the input elements as well as in terms of the breakdown between the elemental semantic content and its corresponding linguistic realization. A theoretical information-theoretic approach is applied to find the shared information content between each input component and the output variable. Here, the information-theoretic approach is framed for NLP classification tasks where the set of input components influence the output probability distribution over a discrete set of classes. We select multiple-choice reading comprehension (MCRC) and sentiment classification (SC) as case studies for the analysis.

MCRC requires the correct answer option to be selected based on several input elements: the context paragraph, the question and the set of answer options. Multiple-choice (MC) assessments are a widely employed method for evaluating the competencies of candidates across diverse settings and tasks on a global scale (Lai et al., 2017a; Richardson et al., 2013a; Sun et al., 2019; Levesque et al., 2012). Given their consequential impact on realworld decisions, the selection of appropriate MC questions tailored to specific scenarios is important for content creators. Consequently, there is a need to comprehend the underlying factors that contribute to the complexity of these assessments.

Complexity of an MC question is best modelled by the distribution over the answer options by human test takers. Therefore, by understanding the influence of each input element on the output distribution, content creators can be better informed to what extent the complexity of an MC question can be controlled from changing each of the input elements. Moreover, analyzing the contribution of the semantic content vs the linguistic realization on the output human distribution informs the impact of the specific word choice in the element on the question complexity. However, it is not scalable to measure the variation in the output human distribution with variation in each of the input elements. Liusie et al. (2023c) demonstrated that the output distribution of automated systems is aligned (with minimal re-shaping parameters) to the human distribution. Therefore, the information-theoretic framework is applied to the output probability distribution by an automated comprehension system.

SC is a common NLP classification task where the dominant sentiment class must be selected from a discrete set of sentiments based on a block of input text. This is an example of a single input text element NLP task. The information-theoretic approach is applied here to understand the role of the semantic content and the linguistic realization on the output distribution over the sentiment classes for popular datasets. It is interesting to analyze SC as ideally the sentiment of a text block should be based on only its semantic meaning. Here, we determine whether this ideal is held in practice.

2 Related Work

084

091

097

099

100

101

102

103

104

105

106

107

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

Features or variables are separate properties that are input to tabular machine learning models to predict a target variable (Hwang and Song, 2023). Feature importance is an active area of research (Huang et al., 2023) where the influence of each feature on the output is determined. The ability to determine which features are most important is useful across many verticals e.g. computer assisted medical diagnosis (Rudin, 2019), weather forecasting (Malinin et al., 2021), fraud detection (Xu et al., 2023) and customer churn prediction (Al-Shourbaji et al., 2023). Similarly, we explore the importance of different aspects (can be interpreted as features) at the input including individual elements and the semantic vs linguistic components for NLP text classification tasks. Typically, the structured nature of tabular data allows common feature selection algorithms to be applied including LASSO (Tibshirani, 1996), marginal screening (Fan and Lv, 2008), orthogonal matching pursuit (Pati et al., 1993) and decision tree based (Costa and Pedreira, 2023). Due to the relatively unstructured nature of text data (compared to tabular data), we propose an information-theoretic approach to identify the most influential inputs.

Sugawara et al. (2017) find a weak correlation between question difficulty and context readability

for MCRC. Additionally, Sugawara et al. (2020) consider the impact on MCRC datasets when input elements are omitted. We propose instead an automated information-theoretic framework for this analysis. Finally, Sorensen et al. (2022) apply an information-theoretic approach for prompt engineering. Our approach can instead be generalized to any NLP classification task. 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

70

3 Theory

Here, we describe the generalized framework to analyze the influence of different elements in NLP text classification tasks: 1. the individual influence of each input element on the output class distribution; 2. the contribution of the semantic content vs its linguistic realization component for a given element. Let an NLP task consist of a set of elements, $\{x_1, \ldots, x_N\} = \mathbf{x}$ and the output, y, such that:

$$P(y) = \mathbb{E}_{P(\mathbf{x})} P\left(y|\mathbf{x}\right) \tag{1}$$

Let X denote the random variables of each the corresponding instances x. Similarly, let Y be the random variable for an instance of the output, y.

Ì

To measure the influence of input \mathbf{x} on output y, a good metric is the mutual information (Depeweg et al., 2018; Malinin and Gales, 2018) which measures how the output changes due to variation in the input. Thus we can define $\mathcal{I}(Y; \mathbf{X})$ a measure of the total input influence. Similarly we can define the influence from an individual element, X_j , $\mathcal{I}(Y; X_j)$ and it should obey:

$$\underbrace{\mathcal{I}(Y;X_j)}_{\text{element}} = \underbrace{\mathcal{I}(Y;\mathbf{X})}_{\text{total}} - \underbrace{\mathcal{I}(Y;\mathbf{X}\setminus X_j|X_j)}_{\text{other}} \quad (2)$$

For each element X_j , its influence is always determined by two components: $X_j^{(s)}$, the semantic information and a relating linguistic realization method which turns an abstract meaning into natural language. Thus, we can calculate the semantic influence as $\mathcal{I}(Y; X_j^{(s)})$ and the linguistic influence implicitly $\mathcal{I}(Y; X_j | X_j^{(s)})$. They should satisfy:

$$\underbrace{\mathcal{I}(Y;X_j)}_{\text{element}} = \underbrace{\mathcal{I}\left(Y;X_j^{(s)}\right)}_{\text{semantic}} + \underbrace{\mathcal{I}\left(Y;X_j|X_j^{(s)}\right)}_{\text{linguistic}}$$
(3)

In practice for an element, x_j , its semantic content171is too abstract to be available. Instead we get access172to one of its realization \tilde{r}_j which is considered to173

be generated from its unobserved semantic content, $x_j^{(s)}$. A set of possible realizations of this semantic element, $\mathcal{R}^{(i)}$, are additionally where each member 176 of this set is, $r_i^{(i)}$ drawn as

$$r_j^{(i)} \sim P_r(r|\tilde{r}_i) \approx P_r(r|s_i) \tag{4}$$

With these settings, the mutual information is calculated as follows. The total influence is:

$$\mathcal{I}(Y; \mathbf{X})$$
(5)
= $\mathcal{H}\left(\mathbb{E}_{P(\mathbf{x})}[P(y|\mathbf{x})]\right) - \mathbb{E}_{P(\mathbf{x})}[\mathcal{H}(P(y|\mathbf{x}))]$

We can also get the element influence as :

$$\mathcal{I}(Y; X_j) \tag{6}$$
$$= \mathcal{H}\left(\mathbb{E}_{P(\mathbf{x})}[P(y|\mathbf{x})]\right) - \mathbb{E}_{P(x_j)}\left[\mathcal{H}(P(y|x_j))\right]$$

It can be decomposed as the semantic influence:

$$\mathcal{I}\left(Y; X_{j}^{(s)}\right) = \mathcal{H}\left(\mathbb{E}_{P(\mathbf{x})}[P(y|\mathbf{x})]\right)$$
(7)
$$-\mathbb{E}_{P\left(x_{j}^{(s)}\right)}\left[\mathcal{H}\left(P\left(y|x_{j}^{(s)}\right)\right)\right]$$

and the linguistic influence:

$$\mathcal{I}\left(Y; X_j | X_j^{(s)}\right) = \mathbb{E}_{P\left(x_j^{(s)}\right)} \left[\mathcal{H}\left(P\left(y | x_j^{(s)}\right)\right)\right] \\ -\mathbb{E}_{P(x_j)} \left[\mathcal{H}\left(P(y | x_j)\right)\right]$$
(8)

The relative contribution of an element to the total influence and of the semantic component for an element can respectively be expressed as:

relative element influence
$$= \frac{\mathcal{I}(Y; X_j)}{\mathcal{I}(Y; \mathbf{X})}$$
 (9)

relative semantic influence = $\frac{\mathcal{I}\left(Y; X_{j}^{(s)}\right)}{\mathcal{I}(Y; X_{j})}$ (10)

197

198

174

175

177

178

179

180

181

183

185

186

187

188

189

190

191

192

193

194

195

196

3.1 Multiple-choice reading comprehension

In this task, candidates are provided with a context 199 passage, c and a corresponding question, q. The 200 objective is to determine the correct answer from a defined set of options, denoted as o. This process involves understanding the question and utilizing 204 the context passage as a source of information to ascertain the most appropriate answer option. The 205 output distribution can be categorised as:

$$P(y) = \mathbb{E}_{P(c,q,o)} P(y|c,q,o)$$



Figure 1: Data generation for the multiple-choice reading comprehension task.

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

225

226

227

228

229

230

231

234

235

236

239

3.1.1 Data generation

For a typical MCRC dataset, the data generation process is shown in Figure 1. A specific semantic content $c^{(s)}$ is chosen and a context c is generated when a certain linguistic realization r is applied. Therefore, the influence of the context C can be divided into $\mathcal{I}(Y; C^{(s)}), \mathcal{I}(Y; C|C^{(s)})$ respectively. A similar procedure is applied on the questions and the options but usually (for the scope of this work) they are generated together as a question-option pair q: from a certain context, a content probe or linguistic probe is generated and then a questionoption pair in natural language is a linguistic realization of this probe as described by: P(q|c). Thus, the output distribution can be rewritten as:

$$P(y) = \mathbb{E}_{P(c^{(s)})} \mathbb{E}_{P(c|c^{(s)})} \mathbb{E}_{P(q|c)} P(y|q,c) \quad (12)$$

We consider only the questions generated from the semantic contents and ignore the questions constrained to a specific realization, by filtering out all questions generated from the specific linguistic realization of the context. This allows our investigation on the question influence to be agnostic of the original context realization.

Note, different context realizations are generated as paraphrases conditional on the original context such that $r \sim P_{\text{gpt}}(r|c)$.

3.1.2 Measure of component influence

The question-option pair in Equation 12 appear as P(q|c), thus instead of $\mathcal{I}(Y;Q)$, we consider $\mathcal{I}(Y;Q|C)$ and get the decomposition:

$$\underbrace{\mathcal{I}(Y;C)}_{\text{context}} = \underbrace{\mathcal{I}(Y;C,Q)}_{\text{total}} - \underbrace{\mathcal{I}(Y;Q|C)}_{\text{question}}$$
(13)

The context influence can be further decomposed:

$$\underbrace{\mathcal{I}(Y;C)}_{\text{context}} = \underbrace{\mathcal{I}(Y;C^{(s)})}_{\text{semantic}} + \underbrace{\mathcal{I}(Y;C|C^{(s)})}_{\text{linguistic}} \quad (14)$$

(11)

Equations 13 and 14 are examples of Equations 241 2 and 3 respectively. Thus, similar to Section 3, 242 the influence terms can be calculated according to 243 Equations 5 to 8. Besides the assumption made in Equation 4 which is general for all the tasks, a 245 further assumption about the questions are made for 246 the MCRC task: instead of sampling from the ideal 247 question generation process, we only observe the question-option pairs generated by humans, $\hat{Q}^{(i)}$, where each member of this set is, $\tilde{q}_{i}^{(i)}$ drawn as:

$$\tilde{q}_j^{(i)} \sim P_{\text{man}}(q|\tilde{r}_i) \approx P_q\left(q|c_i^{(s)}\right) \tag{15}$$

We then make the following approximations:

$$\mathbb{E}_{P(c,q)}\left[P(y|c,q)\right] \approx \tag{16}$$

$$\frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \frac{1}{|\mathcal{R}^{(i)}||\tilde{\mathcal{Q}}^{(i)}|} \sum_{r \in \mathcal{R}^{(i)}, \tilde{q} \in \tilde{\mathcal{Q}}^{(i)}} P(y|\tilde{q},r)$$

255

257

261

262

263

267

269

271

272

273

254

251

$$\mathbb{E}_{P(c,q)} \left[\mathcal{H}(P(y|c,q)) \right] \approx \tag{17}$$

$$\frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \frac{1}{|\mathcal{R}^{(i)}||\tilde{\mathcal{Q}}^{(i)}|} \sum_{r \in \mathcal{R}^{(i)}, \tilde{q} \in \tilde{\mathcal{Q}}^{(i)}} \mathcal{H}(P(y|\tilde{q},r))$$

$$\mathbb{E}_{P(c)}\left[\mathcal{H}(P(y|c))\right] \approx \tag{18}$$

$$\frac{1}{n_{\mathbf{s}}} \sum_{i=1}^{n_{\mathbf{s}}} \frac{1}{|\mathcal{R}^{(i)}|} \sum_{r \in \mathcal{R}^{(i)}} \mathcal{H}(\frac{1}{|\tilde{\mathcal{Q}}^{(i)}|} \sum_{\tilde{q} \in \tilde{\mathcal{Q}}^{(i)}} P(y|\tilde{q}, r))$$

$$\mathbb{E}_{P(c^{(s)})}\left[\mathcal{H}(P(y|c^{(s)}))\right] \approx$$
(19)

$$\frac{1}{n_{\mathrm{s}}} \sum_{i=1}^{n_{\mathrm{s}}} \mathcal{H}(\frac{1}{|\mathcal{R}^{(i)}|} \sum_{r \in \mathcal{R}^{(i)}} \frac{1}{|\tilde{\mathcal{Q}}^{(i)}|} \sum_{\tilde{q} \in \tilde{\mathcal{Q}}^{(i)}} P(y|\tilde{q}, r))$$

with n_s as the number of contexts in a dataset.

3.2 Sentiment classification

For the SC task, the candidate receives a sentence or a short paragraph x and then is requested to choose the sentiment class. Here we are only interested in the influence to the output y from semantic content $x^{(s)}$ and its linguistic realization method: $\mathcal{I}(Y; X^{(s)}), \mathcal{I}(Y; X|X^{(s)})$, as there is only one element at the input. Following Equation 3, the semantic and linguistic breakdown is expressed as:

$$\underbrace{\mathcal{I}(Y;X)}_{\text{text}} = \underbrace{\mathcal{I}(Y;X^{(s)})}_{\text{semantic}} + \underbrace{\mathcal{I}(Y;X|X^{(s)})}_{\text{linguistic}} \quad (20)$$

In practice, the following approximations are made:

$$\mathbb{E}_{P(x)}\left[P(y|x)\right] \approx \tag{21}$$

$$\frac{1}{n_{\mathsf{s}}} \sum_{i=1}^{n_{\mathsf{s}}} \frac{1}{|\mathcal{R}^{(i)}|} \sum_{r \in \mathcal{R}^{(i)}} P(y|r)$$
 277

275

285

287

288

289

290

291

293

294

295

296

297

298

299

300

301

302

303

305

306

307

308

309

310

311

312

313

314

315

316

$$\mathbb{E}_{P(x)}\left[\mathcal{H}(P(y|x))\right] \approx \tag{22}$$

$$\frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \frac{1}{|\mathcal{R}^{(i)}|} \sum_{r \in \mathcal{R}^{(i)}} \mathcal{H}(P(y|r))$$
279

$$\mathbb{E}_{P(x^{(s)})}\left[\mathcal{H}(P(y|x^{(s)}))\right] \approx$$
(23) 281

$$\frac{1}{n_{\rm s}} \sum_{i=1}^{n_{\rm s}} \mathcal{H}(\frac{1}{|\mathcal{R}^{(i)}|} \sum_{r \in \mathcal{R}^{(i)}} P(y|r))$$
282

4 Systems

4.1 Linguistic realization

To analyze the impact of the linguistic realization of a given text element, it is necessary to fix the semantic content of the element. In other work (such as Sugawara et al. (2022)), they attempt to evaluate the effect of linguistic content of the context for MCRC. However, they ignore the requirement to fix the semantic content. We employ a paraphrasing system to generate different linguistic realizations for the same semantic content of a text element. The paraphrasing approach is applied to the context in MCRC and to the input element in SC. To consider a broad range of linguistic realizations for a specific text's semantic content, we generate 8 paraphrases at different readability levels. Hence, we assume (this assumption is assessed in Appendix E) that the linguistic realizations at different readability levels maintain the same semantic content. To change the readability of the text element, we use the zero-shot method as in Farajidizaji et al. (2023) based on Equation 4:

$$r_j^{(i)} \sim P_{\text{LLM}}(r|\tilde{r}_i) \tag{24}$$

In practice, the zero-shot large language model (LLM) (GPT-3.5-turbo) is fed with the original text along with an instruction to alter the language of the text to match the desired readability level. The model, not previously trained on this specific task, uses its pre-existing knowledge and understanding of language structure to alter the readability whilst maintaining the same semantic meaning. The readability level is measured by Flesch reading-ease (Flesch, 1948) score (FRES). See the prompts in Appendix Table 7.

4.2 Reading comprehension

317

318

321

324

325

328

330

332

333

334

335

341

342

343

345

347

349

In this work, MCRC systems are required to return a probability distribution over the answer options. Two alternative architectures are considered for performing the reading comprehension task: encoder-only and decoder-only (see further details in Appendix Figure 5).

Encoder-only is based on the works of Yu et al. (2020); Raina and Gales (2022a); Liusie et al. (2023b); Raina et al. (2023b). Each option is individually encoded along with both the question and the context to produce a score. Softmax is then applied to the scores linked to each option, transforming them into a probability distribution. During inference, the anticipated answer is chosen as the option with the highest associated probability. Inspired by Liusie et al. (2023a) and the recent success observed in finetuning large open-source instruction finetuned language models (Touvron et al., 2023a,b; Jiang et al., 2023, 2024; Tunstall et al., 2023) on various NLP tasks, this work additionally finetunes Llama-2 (Touvron et al., 2023b) as a decoder-only architecture. The context, question and answer options are concatenated into a single natural language prompt. As an autoregressive language model, Llama-2 is requested to effectively return a single token at the output, represented by a single logit distribution over the token vocabulary. The logits associated with the tokens A,B,C,D are respectively normalized using softmax to return the desired probability distribution over the answer options. As with the encoder-only architecture, the option with the highest probability is selected as the answer at inference time. All model outputs are calibrated post-hoc (see Appendix C.4).

4.3 Data complexity classification

standard	С	+	q	+	OA	+	OB	+	OC	+	OD
context	С										
context-question	С	+	q								

Table 1: Input formats for the complexity system.

Here, an automated complexity system takes all
the components of an MC question and classifies
it into one of the 3 classes: *easy, medium* or *hard.*The standard architecture is used for MC question
complexity classification (Raina and Gales, 2022b;
Benedetto, 2023) (see Appendix Figure 6). All elements of an MC question and concatenated together
and fed into a transformer. The embedding of the

prepended [CLS] token is taken as the sentence embedding, which is passed to a classification head, to return a distribution over the three complexity levels. To empirically investigate the relative importance of each element, various input formats are trialled. Table 1 presents different combinations of the context, question and answer options. 361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

4.4 Sentiment classification

SC models take the input text and return a probability distribution over the set of sentiments. Here, the sentiments considered are {negative, positive }. We take the standard approach of taking a pretrained transformer encoder model (Vaswani et al., 2017) with a classification head at the output (Liusie et al., 2022). Similar to the encoder-only approach for MCRC and the data complexity classification system, the SC system only passes the hidden embedding of the [CLS] token to the classification head. Softmax normalizes the logits over the sentiments.

5 Experiments

5.1 Data

We use the RACE++ MCRC dataset (Lai et al., 2017b; Liang et al., 2019a) train split for training both the MC data complexity system and the MCRC model. RACE++ is the largest publicly available dataset from English exams in China partitioned into three difficulty levels: middle school, high school and college (see Appendix D.1 for details about the splits). Additionally, various MCRC datasets are considered as test sets for investigating the influence of each element including the test sets from RACE++, MCTest (Richardson et al., 2013b) and CMCQRD (Mullooly et al., 2023). MCTest requires machines to answer MCRC questions about fictional stories. CMCQRD is a small-scale MCRC dataset from the pre-testing stage partitioned into grade levels B1 to C2 on the Common European Framework of Reference for Languages scale.

For SC, we use IMDb (Maas et al., 2011), Yelppolarity (Yelp) (Zhang et al., 2015) and Amazonpolarity (Amazon) (McAuley and Leskovec, 2013) datasets. IMDb has reviews from the Internet Movie Database. Yelp consists of reviews where 1 or 2 stars is interpreted as negative while 4 or 5 stars is interpreted as positive. Amazon has reviews over a period of 13 years on various products. Hence, all 3 datasets are binary classification tasks.

Table 2 details the main statistics. All the MCRC test sets have 4 options while the selected SC test

		# examples	# options	# words	# questions	semantic diversity	linguistic diversity
MCRC	MCTest RACE++ CMCQRD	142 1,007 150	4 4 4	209 278 683	4 3.7 5.5	$\begin{array}{c} 0.079_{\pm 0.015} \\ 0.101_{\pm 0.015} \\ 0.092_{\pm 0.010} \end{array}$	$\begin{array}{c} 0.018_{\pm 0.006} \\ 0.016_{\pm 0.007} \\ 0.022_{\pm 0.011} \end{array}$
SC	IMDB Yelp Amazon	500 500 500	2 2 2	226 133 74	- - -	$\begin{array}{c} 0.084_{\pm 0.019} \\ 0.108_{\pm 0.037} \\ 0.135_{\pm 0.024} \end{array}$	$\begin{array}{c} 0.023_{\pm 0.011} \\ 0.030_{\pm 0.024} \\ 0.037_{\pm 0.022} \end{array}$

Table 2: Statistics for multiple-choice reading comprehension (MCRC) and sentiment classification (SC) test datasets. See Appendix 5 for additional datasets.

sets have 2 options: negative and positive. The 410 number of words for MCRC is the lengths of the 411 contexts. It is seen that the test sets have vary-412 ing lengths from 200 to 700 words and 75 to 230 413 words for MCRC and SC tasks respectively. For 414 415 the MCRC datasets, the total examples are reported after filtering out all linguistic probe questions (Ap-416 pendix D.3 for the procedure). So the focus is 417 only on the semantic probe questions as assumed 418 in Section 3.1.2. For the SC test sets, a subset of 419 500 examples is selected for each dataset to remain 420 within the financial budget for use of ChatGPT for 421 the generation of different linguistic realizations. 422

The semantic diversity is also calculated for each dataset. This score is the mean cosine distance between each text embedding to the centroid of all embeddings (Raina et al., 2023a). Greater the score, greater the semantic variation in the set of texts being considered. The sentence embedder from Ni et al. (2022) is used to generate the embeddings¹. The semantic diversity is calculated on the contexts for MCRC. Finally, the linguistic diversity calculates the mean variation in the embedding space for different linguistic realizations for each text.

5.2 Model details

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

The decoder-only implementation of the MCRC system is based upon the pretrained instruction finetuned Llama2-7B model². The main paper for MCRC focuses only on the decoder-only implementation (see Appendix D.2 for the encoder-only implementations). ELECTRA-base (Clark et al., 2020) is selected for the data complexity evaluator models. The pretrained model is finetuned on the RACE++ train split with the complexity class (easy, medium or hard) as the label (hyperparameter tuning details in Appendix D.2). For SC, the

main paper reports results based on a RoBERTa architecture. The selected model has been finetuned on IMDb training data³. Due to the similarity in content between Yelp, Amazon and IMDb, the RoBERTa model finetuned on IMDb is also applied on all SC test sets. For reproducibility, see BERT (Devlin et al., 2018) system in Appendix D.2. 446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

6 Results

6.1 Reading Comprehension

Table 3 presents the performance of Llama-2 on the various MCRC datasets. The highest accuracy is observed on MCTest with 92.5% and the lowest on CMCQRD with 79.9%. Additionally, the performance of the model is reported on each dataset after generating 8 paraphrases for each context (see Section 4.1). It is observed there is a consistent drop in performance on the paraphrased contexts compared to the original. This is expected as the nature of the machine generated paraphrased contexts do not necessarily align with the type of contexts observed in the original dataset. Table 3 further investigates the influence of the different elements: specifically the context influence compared to the question influence (note the question includes the options - see Section 3.1). The context of an MCRC question plays an important role in the output with influences up to 45% for MCTest.

The complexity of an MCRC question is described by the shape of the output distribution. A sharp distribution about the correct answer is indicative of an easy question while a flatter distribution over all the answer options indicates a harder question. Therefore, the strong influence of the context demonstrated in Table 3 emphasises that the context (alongside the specific posed question) is important in controlling the complexity of a question. To further verify the influence of the context,

¹Available at: https://huggingface.co/sente nce-transformers/sentence-t5-base

²Available at https://huggingface.co/meta-l lama/Llama-2-7b-chat-hf

³Available at https://huggingface.co/wrmur ray/roberta-base-finetuned-imdb

dataset	accuracy		influence							
ualasei	original	al para total		question	context	context-semantic	context-linguistic			
MCTest	92.5	85.8	0.212	0.116 (54.7%)	0.096 (45.3%)	0.068 (70.6%)	0.028 (29.4%)			
RACE++	86.0	82.9	0.298	0.161 (56.1%)	0.131 (43.9%)	0.108 (82.5%)	0.023 (17.5%)			
CMCQRD	79.9	69.4	0.290	0.211 (72.7%)	0.079 (27.3%)	0.067 (83.8%)	0.012 (16.2%)			

Table 3: Decomposition of total input influence for Llama-2 on MCRC datasets.



Figure 2: Normalized ranks (rank / total examples) of complexity scores for each complexity level using 3 evaluators: context, context-question and standard. See Appendix F.1.1 for the performance.



Figure 3: The relative question influence changes with the subset chosen by the rank of context complexity.

Figure 2 plots the complexity score output by the data complexity classifier. In particular, the distribution is shown for the complexity scores on the different subsets from RACE++ (for CMCQRD see Appendix Figure 10) of different complexity levels. Note, the normalized ranks of the complexity scores is plotted where the global rank is found for a given complexity score and divided by the total number of examples. The distributions are shown for the standard system (context, question and options), context-question system (context and question) and the context-only system. The context is clearly sufficient to determine the complexity levels of MC questions for these datasets, empirically supporting the importance of the context.

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

For the context, Table 3 further reports the influence for the semantic and linguistic components. For all 3 datasets, the semantic meaning of a context has a greater influence on the final output distribution but the specific linguistic realization also influences the output. Specifically, the relative semantic influence is greatest for MCTest and lowest for CMCQRD. This is supported by Table 2 where the calculated semantic diversity is the lowest for MCTest with similar linguistic diversities across the datasets. Additionally, Table 3 suggests a relationship between the difficulty of a dataset (indicated by the accuracy of the model) and the relative question influence. The relative question influence increases with more challenging datasets. To further explore this observation, Figure 3 determines whether the question influence is directly linked to the complexity of a question. The data complexity classifier (QC system) is used to rank all the contexts according to their complexity. Then retaining a certain fraction of the most complex contexts, the relative question influence is plotted. The plot is for all the datasets combined (for RACE++ see Appendix Figure 11). Compared to the random ordering, it is clear that retaining the most complex contexts leads to larger question influence scores. The increase in question influence with more challenging contexts supports the trend from Table 3. Therefore, a more challenging context allows a greater variation in question difficulties, leading to a greater question influence on the output.

500

501

502

503

504

505

506

507

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

6.2 Sentiment classification

For the SC task there is only a single input. Therefore, we focus on exploring the relative contributions of the semantic and linguistic components of the input. Intuitively, the sentiment of a passage of text should be determined solely by its semantic content. However, Table 4 shows that for 3 popular SC datasets, the linguistic realization does influence the output. In particular, the relative linguistic component for Amazon hits 10% of the total. It appears the semantic component is more dominant for IMDb and Yelp compared to Amazon. One possible reason is that the length of texts is shorter for Amazon compared to the other datasets (see

datasat	accura	acy	influence					
ualasel	original	para	total	semantic	linguistic			
IMDB	94.8	93.4	0.472	0.444 (94.0%)	0.028 (6.0%)			
Yelp	94.3	93.9	0.472	0.445 (94.2%)	0.027 (5.8%)			
Amazon	91.0	89.5	0.361	0.325 (90.0%)	0.036 (10.0%)			

Table 4: Decomposition of input influence for different models in various datasets for sentiment classification.



Figure 4: Entropy filtered pairwise agreement in paraphrase readability and true class probability ordering with various minimum readability gaps.

Table 2). Hence, longer texts have a greater opportunity to reinforce sentiment being expressed, which makes it more robust to different linguistic realizations. Appendix F.2 further explores this hypothesis by considering additional SC datasets.

6.3 Impact of linguistic realization

543

544

545

546

547

549

550

551

553

555

556

559

562

563

564

565

567

568

569

571

Section 6.1 demonstrated that the linguistic realization of the context in MCRC has measurable influence on the output distribution over the options. Here, we investigate the correlation between the readability of a given paraphrase of the original text and the output probability of the true class, termed true class probability (TCP). An entropy filter is applied to remove examples for which the entropy of the output distribution is too high, as high entropy examples suggest a random guess and hence challenging to ascertain whether a correlation exists. Figure 4a sweeps the fraction of examples retained according to the entropy filter and plots the fraction of the remaining examples for which the ordering of TCP scores for every pair of paraphrases matches the ordering of their real readability scores. The plots are indicated for minimum readability gaps of 0, 25 and 50 for the pairs of paraphrases. It is observed that the readability of a paraphrase corresponds to the returned TCP, with a stronger correlation when the minimum readability gap between the pairs of paraphrases is higher. A similar process is applied for SC in Figure 4b to determine the relationship between the readability of the linguistic realization of the input and the TCP. Like MCRC there is a positive correlation between the readability level and the TCP, with a more pronounced relationship by constraining the pairs of paraphrases to have a larger readability gap. 572

573

574

575

576

577

578

579

580

582

583

584

585

586

588

589

590

591

592

593

594

595

596

597

598

599

600

7 Conclusions

This work describes an information-theoretic framework for text classification tasks. The framework determines the influence of each input element on the output. Additionally, each element is partitioned into its semantic and linguistic components. MCRC and SC are considered as case study tasks for analysis. For MCRC, it is found that both the context and question elements play influential roles on the output distribution. It is further established that selection of more challenging contexts permits greater variation (in terms of complexity) of questions on the context. Simpler contexts limit the range of the complexity to only easy questions. Hence, content creators need to carefully consider the choice of the context when designing MC questions to cater to a range of difficulty levels. In SC, the linguistic realization of the input has a measurable impact on the output. Hence, the text wording cannot be neglected when deducing the sentiment. For both tasks, higher the readability of a specific linguistic realization, greater the probability of the true class in the output distribution.

8 Limitations

601

617

618

619

620

621

624

631

632

633

635

641

643

645

647

648

This work has several assumptions that must be stated. For the multiple-choice reading comprehension analysis, the question influence is based on 604 real questions generated by humans on the original context. However, there is the possibility that the set of questions on a given context are not generated independently but instead the question creator has curated the question set together. Additionally, it is assumed that the paraphrasing of texts only 610 changes the linguistic realization. However, it is 611 likely that it also has an impact on the semantic 612 content to an extent, which is reflected in the lin-613 guistic component influence on the output. We do 614 quantify the appropriateness of the paraphrasing system in the Appendix. 616

9 Ethics statement

There are no ethical concerns with this work.

References

- Ibrahim AlShourbaji, Na Helian, Yi Sun, Abdelazim G Hussien, Laith Abualigah, and Bushra Elnaim. 2023.
 An efficient churn prediction model using gradient boosting machine and metaheuristic optimization. *Scientific Reports*, 13(1):14441.
- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Razieh Baradaran, Razieh Ghiasi, and Hossein Amirkhani. 2022. A survey on machine reading comprehension systems. *Natural Language Engineering*, 28(6):683–732.
- Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. Tweeteval: Unified benchmark and comparative evaluation for tweet classification. In *Findings of the Association* for Computational Linguistics: EMNLP 2020, pages 1644–1650.
- Luca Benedetto. 2023. A quantitative study of nlp approaches to question difficulty estimation. *arXiv* preprint arXiv:2305.10236.
- Oscar Chew, Kuan-Hao Huang, Kai-Wei Chang, and Hsuan-Tien Lin. 2023. Understanding and mitigating spurious correlations in text classification. *arXiv preprint arXiv:2305.13654*.
- KR1442 Chowdhary and KR Chowdhary. 2020. Natural language processing. *Fundamentals of artificial intelligence*, pages 603–649.
- Kevin Clark, Minh-Thang Luong, Quoc V Le, and 651 Christopher D Manning. 2020. Electra: Pre-training 652 text encoders as discriminators rather than generators. 653 arXiv preprint arXiv:2003.10555. 654 Vinícius G Costa and Carlos E Pedreira. 2023. Recent 655 advances in decision trees: An updated survey. Arti-656 ficial Intelligence Review, 56(5):4765-4800. 657 Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Fi-658 nale Doshi-Velez, and Steffen Udluft. 2018. Decom-659 position of uncertainty in bayesian deep learning for 660 efficient and risk-sensitive learning. In International 661 Conference on Machine Learning, pages 1184–1193. 662 PMLR. 663 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and 664 Luke Zettlemoyer. 2023. Qlora: Efficient finetuning 665 of quantized llms. arXiv preprint arXiv:2305.14314. 666 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 667 Kristina Toutanova. 2018. Bert: Pre-training of deep 668 bidirectional transformers for language understand-669 ing. arXiv preprint arXiv:1810.04805. 670 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel 671 Stanovsky, Sameer Singh, and Matt Gardner. 2019. 672 Drop: A reading comprehension benchmark re-673 quiring discrete reasoning over paragraphs. arXiv 674 preprint arXiv:1903.00161. 675 Liam Dugan, Eleni Miltsakaki, Shriyash Upadhyay, 676 Etan Ginsberg, Hannah Gonzalez, Dayheon Choi, 677 Chuning Yuan, and Chris Callison-Burch. 2022. A 678 feasibility study of answer-agnostic question genera-679 tion for education. arXiv preprint arXiv:2203.08685. 680 Jianging Fan and Jinchi Lv. 2008. Sure independence 681 screening for ultrahigh dimensional feature space. 682 Journal of the Royal Statistical Society Series B: Sta-683 tistical Methodology, 70(5):849–911. 684 Asma Farajidizaji, Vatsal Raina, and Mark Gales. 2023. 685 Is it possible to modify text to a target readability 686 level? an initial investigation using zero-shot large 687 language models. arXiv preprint arXiv:2309.12551. 688 Rudolph Flesch. 1948. A new readability yardstick. 689 Journal of applied psychology, 32(3):221. 690 Yifan Gao, Lidong Bing, Wang Chen, Michael R Lyu, 691 and Irwin King. 2018. Difficulty controllable gen-692 eration of reading comprehension questions. arXiv 693 preprint arXiv:1807.03586. 694 Yifan Gao, Lidong Bing, Piji Li, Irwin King, and 695 Michael R Lyu. 2019. Generating distractors for 696 reading comprehension questions from real exami-697 nations. In Proceedings of the AAAI Conference on 698 Artificial Intelligence, volume 33, pages 6423-6430. 699
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR.

700

702

704

- 750
- 751 754 755

- Chao Huang, Diptesh Das, and Koji Tsuda. 2023. Feature importance measurement based on decision tree sampling. arXiv preprint arXiv:2307.13333.
- Yejin Hwang and Jongwoo Song. 2023. Recent deep learning methods for tabular data. Communications for Statistical Applications and Methods, 30(2):215-226.
- Albert O Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. arXiv preprint arXiv:2310.06825.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. arXiv preprint arXiv:2401.04088.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 252-262.
- Ghader Kurdi, Jared Leo, Bijan Parsia, Uli Sattler, and Salam Al-Emari. 2020. A systematic review of automatic question generation for educational purposes. International Journal of Artificial Intelligence in Education, 30:121-204.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and E. Hovy. 2017a. Race: Large-scale reading comprehension dataset from examinations. In EMNLP.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017b. Race: Large-scale reading comprehension dataset from examinations. arXiv preprint arXiv:1704.04683.
- Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning. Citeseer.
- Yichan Liang, Jianheng Li, and Jian Yin. 2019a. A new multi-choice reading comprehension dataset for curriculum learning. In Proceedings of The Eleventh Asian Conference on Machine Learning, volume 101 of Proceedings of Machine Learning Research, pages 742–757, Nagoya, Japan. PMLR.
- Yichan Liang, Jianheng Li, and Jian Yin. 2019b. A new multi-choice reading comprehension dataset for curriculum learning. In Asian Conference on Machine Learning, pages 742-757. PMLR.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.

Adian Liusie, Potsawee Manakul, and Mark JF Gales. Zero-shot nlg evaluation through pair-2023a. ware comparisons with llms. arXiv preprint arXiv:2307.07889.

760

761

762

764

765

766

767

768

769

770

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

- Adian Liusie, Vatsal Raina, and Mark Gales. 2023b. " world knowledge" in multiple choice reading comprehension. In The Sixth Fact Extraction and VERification Workshop, page 49.
- Adian Liusie, Vatsal Raina, Andrew Mullooly, Kate Knill, and Mark J. F. Gales. 2023c. Analysis of the cambridge multiple-choice questions reading dataset with a focus on candidate response distribution.
- Adian Liusie, Vatsal Raina, Vyas Raina, and Mark Gales. 2022. Analyzing biases to spurious correlations in text classification tasks. In Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing, pages 78-84.
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies, pages 142-150.
- Andrey Malinin, Neil Band, German Chesnokov, Yarin Gal, Mark JF Gales, Alexey Noskov, Andrey Ploskonosov, Liudmila Prokhorenkova, Ivan Provilkov, Vatsal Raina, et al. 2021. Shifts: A dataset of real distributional shift across multiple large-scale tasks. arXiv preprint arXiv:2107.07455.
- Andrey Malinin and Mark Gales. 2018. Predictive uncertainty estimation via prior networks. Advances in neural information processing systems, 31.
- Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Mqag: Multiple-choice question answering and generation for assessing information consistency in summarization. arXiv preprint arXiv:2301.12307.
- Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text. In Proceedings of the 7th ACM conference on Recommender systems, pages 165-172.
- Andrew Mullooly, Øistein Andersen, Luca Benedetto, Paula Buttery, Andrew Caines, Mark J.F. Gales, Yasin Karatay, Kate Knill, Adian Liusie, Vatsal Raina, and Shiva Taslimipoor. 2023. The Cambridge Multiple-Choice Questions Reading Dataset. Cambridge University Press and Assessment.
- Jianmo Ni, Gustavo Hernandez Abrego, Noah Constant, Ji Ma, Keith Hall, Daniel Cer, and Yinfei Yang. 2022. Sentence-t5: Scalable sentence encoders from pretrained text-to-text models. In Findings of the Association for Computational Linguistics: ACL 2022, pages 1864-1874.

Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In *Proceedings* of the 41st annual meeting of the Association for Computational Linguistics, pages 160–167.

815

816

817

819

824

825

826

827

829

836

838 839

840

841

850

851

852

853

856

857

864

865

- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the* 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05), pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.
- Yagyensh Chandra Pati, Ramin Rezaiifar, and Perinkulam Sambamurthy Krishnaprasad. 1993. Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Proceedings of 27th Asilomar conference on signals, systems and computers*, pages 40–44. IEEE.
 - Vatsal Raina and Mark Gales. 2022a. Answer uncertainty and unanswerability in multiple-choice machine reading comprehension. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1020–1034, Dublin, Ireland. Association for Computational Linguistics.
 - Vatsal Raina and Mark Gales. 2022b. Multiple-choice question generation: Towards an automated assessment framework. *arXiv preprint arXiv:2209.11830*.
 - Vatsal Raina, Nora Kassner, Kashyap Popat, Patrick Lewis, Nicola Cancedda, and Louis Martin. 2023a.
 ERATE: Efficient retrieval augmented text embeddings. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 11–18, Dubrovnik, Croatia. Association for Computational Linguistics.
 - Vatsal Raina, Adian Liusie, and Mark Gales. 2023b. Analyzing multiple-choice reading and listening comprehension tests. *arXiv preprint arXiv:2307.01076*.
 - Matthew Richardson, Christopher J.C. Burges, and Erin Renshaw. 2013a. MCTest: A challenge dataset for the open-domain machine comprehension of text. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 193– 203, Seattle, Washington, USA. Association for Computational Linguistics.
 - Matthew Richardson, Christopher JC Burges, and Erin Renshaw. 2013b. Mctest: A challenge dataset for the open-domain machine comprehension of text. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 193– 203.
- Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence*, 1(5):206–215.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. arxiv 2019. *arXiv preprint arXiv:1910.01108*.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics. 871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

- Taylor Sorensen, Joshua Robinson, Christopher Michael Rytting, Alexander Glenn Shaw, Kyle Jeffrey Rogers, Alexia Pauline Delorey, Mahmoud Khalil, Nancy Fulda, and David Wingate. 2022. An information-theoretic approach to prompt engineering without ground truth labels. *arXiv preprint arXiv:2203.11364*.
- Felix Stahlberg. 2020. Neural machine translation: A review. *Journal of Artificial Intelligence Research*, 69:343–418.
- Saku Sugawara, Yusuke Kido, Hikaru Yokono, and Akiko Aizawa. 2017. Evaluation metrics for machine reading comprehension: Prerequisite skills and readability. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 806–817.
- Saku Sugawara, Nikita Nangia, Alex Warstadt, and Samuel R Bowman. 2022. What makes reading comprehension questions difficult? *arXiv preprint arXiv:2203.06342*.
- Saku Sugawara, Pontus Stenetorp, Kentaro Inui, and Akiko Aizawa. 2020. Assessing the benchmarking capacity of machine reading comprehension datasets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8918–8927.
- Jiao Sun, Yufei Tian, Wangchunshu Zhou, Nan Xu, Qian Hu, Rahul Gupta, John Frederick Wieting, Nanyun Peng, and Xuezhe Ma. 2023. Evaluating large language models on controlled generation tasks. *arXiv* preprint arXiv:2310.14542.
- Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. 2019. Dream: A challenge data set and models for dialogue-based reading comprehension. *Transactions of the Association for Computational Linguistics*, 7:217–231.
- Robert Tibshirani. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267– 288.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti

Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

927

928

930

931 932

933

936

937

939

940

941

942 943

944

945

947

948 949

951

952

953

959

961 962

963

964

965

966

967

968

969

- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. 2023. Zephyr: Direct distillation of Im alignment. *arXiv preprint arXiv:2310.16944*.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
 - Mayur Wankhade, Annavarapu Chandra Sekhara Rao, and Chaitanya Kulkarni. 2022. A survey on sentiment analysis methods, applications, and challenges. *Artificial Intelligence Review*, 55(7):5731–5780.
 - Adhika Pramita Widyassari, Supriadi Rustad, Guruh Fajar Shidik, Edi Noersasongko, Abdul Syukur, Affandy Affandy, et al. 2022. Review of automatic text summarization techniques & methods. *Journal* of King Saud University-Computer and Information Sciences, 34(4):1029–1046.
 - Biao Xu, Yao Wang, Xiuwu Liao, and Kaidong Wang. 2023. Efficient fraud detection using deep boosting decision trees. *Decision Support Systems*, page 114037.
 - Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.
 - Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. Reclor: A reading comprehension dataset requiring logical reasoning. In *International Conference on Learning Representations*.
 - Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

A Extended related work

973

974

975

976

977

979

981

987

989

991

993

997

999

1000

1002

1003

1004

1005

1006

1008

1009

1011

1012

1013

1014

1016

1018

1019

1020

1021

1022

1023

In multiple-choice reading comprehension, the influence of each element on the final output distribution is directly linked to the complexity of a multiple-choice question. More complex multiplechoice questions can expect to have flat distributions over the answer options while easier questions are sharp around the correct answer. Numerous studies have looked at the factors that potentially influence complexity of context passages in relation with reading comprehension tasks. Sugawara et al. (2022) observed that the diversity of contexts in MC questions determines the diversity of questions possible conditioned on the contexts. In their work they found that variables such as passage source, length, or readability measures do not significantly affect the model performance. Further, Khashabi et al. (2018) introduced the role of the original source from which the contexts are extracted in shaping overall complexity. Through the augmentation of the dataset by diversifying the corpus sources, they aimed to enhance the dataset quality.

Question complexity has repeatedly been discussed within prior literature, with varying definitions. Liang et al. (2019b) classified questions into distinct categories with complexity scores ranking from lowest to highest for word matching, paraphrasing, single-sentence reasoning, multisentence reasoning, and ambiguous questions. Gao et al. (2018) quantified complexity as the number of reasoning steps required to derive the answer. Similar definitions have been upheld by Yang et al. (2018) and Dua et al. (2019), who expanded the understanding of question complexity to encompass not only contextual comprehension but also factors such as the confidence of a pretrained questionanswering model.

Distractor (incorrect options for multiple-choice questions) complexity has been less explored. Gao et al. (2019) determined distractor complexity based on the similarity between distractors and the ground-truth. Employing an n-gram overlap metric, Banerjee and Lavie (2005) introduced a method to assess distractor complexity. Dugan et al. (2022) further dissected distractor complexity, analyzing qualities such as relevance, interpretability and acceptability compared to human markers.

As not explored in previous MCRC complexity literature, in this work the information-theoretic approach is applied to characterize the influence of each element in a given multiple-choice reading comprehension dataset. Greater the influence of an element, greater the scope to control the complexity of the multiple-choice reading comprehension task. 1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1058

1059

1062

1064

In sentiment classification, it is expected the sentiment class should be dependent on the semantic meaning of the text rather than its linguistic realization. However, Liusie et al. (2022) find that *shortcut* systems that have access only to the stopwords in the original text are also able to identify the sentiment class. Hence, they find the stop words chosen in the text do influence the sentiment class, which we express as the specific linguistic realization. Chew et al. (2023) further aim to correct for the bias from spurious correlations. In this work, we explicitly quantify the influence of the semantic and linguistic components of the text.

B Additional Tasks

B.1 Grade classification

In the grade classification task, the system is input a prompt z and a response x and then is required to output a number in the range 1 to 6 with 6 denoting the highest degree of alignment between the given prompt and its corresponding response. The task is traditionally a regression-oriented but here we apply our information-theoretic classification framework by treating it as a 6-option classification task.

$$P(y) = \mathbb{E}_{P(z,x)} P(y|z,x) \tag{25}$$

Further, the semantic meaning of the response x can be seen as generated from P(x|z). In the considered datasets, a prompt has a large-scale number of responses while the number of prompts are limited. Thus, here we focus only on the analysis of the influence of the response. The output equation can be rewritten as:

$$P(y) = \mathbb{E}_{P(z)} \mathbb{E}_{P(x^{(s)}|z)} \mathbb{E}_{P(x|x^{(s)})} P(y|x,z)$$
 (26) 1060

Therfore, we can calculate the influence from the semantic meaning and the linguistic realization of the responses:

$$\underbrace{\mathcal{I}(Y;X|Z)}_{\text{text}} = \underbrace{\mathcal{I}(Y;X^{(s)}|Z)}_{\text{semantic}} + \underbrace{\mathcal{I}(Y;X|X^{(s)},Z)}_{\text{linguistic}}$$
(27)

In practice, we collect the prompt set \mathcal{Z} and a response set \mathcal{X} for each prompt. For each response 1066



Figure 5: Architectures for multiple-choice reading comprehension with context, c, question, q and options, o.

1067 $x^{(i,j)}$ of the prompt $z^{(i)}$, we observe a realization1068 $\tilde{r}^{(i,j)}$ and additionally generate a set of realizations1069 $\mathcal{R}^{(i,j)}$ of the semantic meaning of the responses.1070The following approximations are made:

071
$$\mathbb{E}_{P(z)}\left[\mathcal{H}\left(\mathbb{E}_{P(x|z)}\left[P(y|x,z)\right]\right)\right] \approx$$
(28)

$$\frac{1}{|\mathcal{Z}|} \sum_{z^i \in Z} \mathcal{H}(\frac{1}{n_{\mathrm{s}}^i} \sum_{j=1}^{n_{\mathrm{s}}^i} \frac{1}{|\mathcal{R}^{(i,j)}|} \sum_{r \in \mathcal{R}^{(i,j)}} P(y|r,z^i))$$

1073
$$\mathbb{E}_{P(x^{(s)},z)}\left[\mathcal{H}(P(y|x^{(s)},z))\right] \approx$$
(29)

$$\frac{1}{|\mathcal{Z}|} \sum_{z^i \in Z} \frac{1}{n_{\mathtt{s}}^i} \sum_{j=1}^{n_{\mathtt{s}}^i} \mathcal{H}(\frac{1}{|\mathcal{R}^{(i,j)}|} \sum_{r \in \mathcal{R}^{(i,j)}} P(y|r, z^i))$$

1077

1078

1079

1080

1083

1085

1086

1087

1088

1089

1074

1072

$$\mathbb{E}_{P(x,z)}\left[\mathcal{H}(P(y|x,z))\right] \approx \tag{30}$$

$$\frac{1}{|\mathcal{Z}|} \sum_{z^i \in Z} \frac{1}{n_{\mathbf{s}}^i} \sum_{j=1}^{n_{\mathbf{s}}^i} \frac{1}{|\mathcal{R}^{(i,j)}|} \sum_{r \in \mathcal{R}^{(i,j)}} \mathcal{H}(P(y|r, z^i))$$

where n_{s}^{i} is the number of responses for prompt $z^{(i)}$.

C Systems

C.1 Reading comprehension

In Section 4.1, we show the details of Llama-2 model we used for reading comprehension tasks. Here we show more details about the extra models we trialled in the encoder-only set-up.

Models we used can be divided into two model families: encoder-only and decoder-only. Within the family of encoder-only models, we consider a BERT (Devlin et al., 2018) and RoBERTa (Liu et al., 2019) based systems. For these models, as 1090 shown in Figure 5, each option is paired with the 1091 context as well as the question at the model in-1092 put. The model then generates a score based on 1093 the input on which Softmax is applied to convert to 1094 a probability distribution over the answer options. As mentioned in the main paper, Llama-2 is used as 1096 a popular example with the decoder-only architec-1097 ture. The input consists the context, question and 1098 all options. The model output is the probability of 1099 all possible output tokens. The probability associ-1100 ated with the tokens A,B,C,D are concatenated and 1101 normalized using softmax to return the final proba-1102 bility distribution. An illustration of decoder-only 1103 architecture is also shown in Figure 5. 1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

C.2 Grade classification

The architecture of the model used for grade classification is based on the decoder-only transformer architecture with Llama-2 as described in Section C.1. Hence, the prompt and the response are concatenated together at the input to the model, which is trained to return a probability distribution over the 6 grade classes.

C.3 Data complexity classification

The system is required to output a probability distribution among options: easy, medium and hard. The architecture is presented in Figure 6. The complexity score S_c is calculated as:

$$S_c = 0 * P_{easy} + 0.5 * P_{medium} + 1 * P_{hard}$$
 (31)

C.4 Calibration

The trained models were calibrated post-hoc using1120single parameter temperature annealing (Guo et al.,1121



Figure 6: Architecture for MC question complexity classifier with context, c, question, q and options, $\{o\}$.

2017). It is necessary to calibrate the models for the absolute information-theoretic measures to be meaningful. Uncalibrated, model probabilities are determined by applying the softmax to the output logit scores s_i :

$$P(y=k;\boldsymbol{\theta}) \propto \exp(s_k) \tag{32}$$

where k denotes a possible output class for a prediction y. Temperature annealing 'softens' the output probability distribution by dividing all logits by a single parameter T prior to the softmax.

$$P_{CAL}(y=k;\boldsymbol{\theta}) \propto \exp(s_k/T)$$
 (33)

As the parameter T does not alter the relative rankings of the logits, the model's prediction will be unchanged and so temperature scaling does not affect the model's accuracy. The parameter T is chosen such that the accuracy of the system is equal to the mean of the maximum probability (as is expected for a calibrated system).

D Experiments

D.1 Data

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

Table 5 provides a breakdown of the RACE++ dataset, which is divided into three subsets: RACE-M, RACE-H, and RACE-C. These subsets correspond to English exam materials from Chinese middle schools (RACE-M), high schools (RACE-H), and colleges (RACE-C) respectively. Table 5 displays key statistics for each of these subsets including the number of contexts, the average number of questions for one context, and the semantic and linguistic diversity. Note, for all MCRC datasets, the options are re-ordered such that the true class is the first option.

SST-2 (Socher et al., 2013) and TweetEval (Barbieri et al., 2020) are considered as additional datasets for sentiment classification, which were not presented in the main text. SST-2 (Stanford Sentiment Treebank) corpus consists of movie reviews provided by Pang and Lee (2005) which are classified as either positive or negative. TweetEval consists of seven heterogeneous tasks based on tweets from Twitter. Here, the focus is on the tweetemotion task where each tweet is classified as joy, sadness, optimism or anger. These two datasets are included here as they have shorter inputs texts than IMDb, Yelp and Amazon. 1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

Hewlett foundation⁴, a competition requiring automated grading of student-written essays, is included as the dataset for the grade classification task. In total it has 8 subgroups and each subgroup has 1 prompt with several responses. Here we only choose the first 2 subgroups and turn them into a 6 option classification task with 6 denoting the highest degree of alignment between the given prompt and its corresponding response. 3,583 responses are sampled for the training split while 500 are selected randomly as the test split. From Table 5, it is clear compared to other dataset, the responses in Hewlett are mutually semantically closer in meaning compared to other datasets as they are all responses to just 2 prompts.

Table 5 also supports that our paraphrasing system is sensible by showing the linguistic diversity is much lower than the semantic diversity among all responses. A more detailed quality verification process of our paraphrasing system is shown in Appendix E.

D.2 Models

For the multiple choice reading comprehension tasks, three models are used: Llama2, RoBERTa and Longformer.

Pretrained Llama-2 (7 billions parameters) is finetuned specifically on the train split of RACE++ with hyperparameter tuning on the validation split. However, it is not computationally feasible to train all the model parameters of Llama-2. Therefore, parameter efficient finetuning is used with quantized low rank adapters (QLoRA) (Dettmers et al., 2023). The final training parameters finetune the model with a learning rate of 1e-4, batch size of 4, lora rank of 8 with lora $\alpha = 16$ and dropout 0.1. The model is trained for 1 epoch taking 7 hours on an NVIDIA A100 machine. For the main paper results, the Llama-2 model is selected due to its best accuracy.

⁴Available at: https://www.kaggle.com/c/asap-aes

		# examples	# options	# words	# questions	semantic diversity	linguistic diversity
MCRC	RACE-M RACE-H RACE-C	362 510 135	4 4 4	200 308 375	4 3.3 5.2	$\begin{array}{c} 0.096_{\pm 0.017} \\ 0.100_{\pm 0.012} \\ 0.097_{\pm 0.011} \end{array}$	$\begin{array}{c} 0.017_{\pm 0.007} \\ 0.016_{\pm 0.006} \\ 0.018_{\pm 0.006} \end{array}$
SC	SST-2 TweetEval	500 500	2 4	20 17	-	$\begin{array}{c} 0.145_{\pm 0.021} \\ 0.149_{\pm 0.022} \end{array}$	$\begin{array}{c} 0.087_{\pm 0.034} \\ 0.081_{\pm 0.033} \end{array}$
GC	Hewlett	500	6	383	2	$0.063_{\pm 0.017}$	$0.021_{\pm 0.009}$

Table 5: Statistics for breakdown of additional test datasets including RACE-M, RACE-H, RACE-C for RACE++ in multiple-choice reading comprehension (MCRC); SST-2, TweetEval in sentiment classification (SC) and Hewlett in grade classification (GC).

filton	accura	acy			influer	nce		
original p		para	total	question	context	context-semantic	context-linguistic	
No	84.2	81.5	0.284	0.164 (57.7%)	0.120 (42.3%)	0.135 (81.4%)	0.031 (18.6%)	
Yes	86.0	82.9	0.298	0.161 (56.1%)	0.131 (43.9%)	0.108 (82.5%)	0.023 (17.5%)	

Table 6: The effect of the question filter on element influence for Llama-2 on the RACE++ test set.

RoBERTa (82 millions parameters)⁵ is finetuned on the train split of the RACE dataset (RACE-M and RACE-H). The details of the specific Longformer (336 million parameters)⁶ are detailed in Manakul et al. (2023). For the main paper results, the Llama-2 model is selected due to its best accuracy.

1206

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

For the data complexity classification system, pretrained ELECTRA-base (110 millions parameters) is finetuned on the RACE++ train split with the complexity class (easy, medium or hard) as the label. The model is trained using the AdamW optimizer, a batch size of 3, learning rate of 2e-5, max number of epochs of 3 with all inputs truncated to 512 tokens. An ensemble of 3 models is trained. Training for each model takes 3 hours on an NVIDIA V100 graphical processing unit.

For the sentiment classification task, we used the models from RoBERTa and distilBERT(82 millions parameters) (Sanh et al., 2019) family and the they are finetuned on various datasets as explained in the following. The train split of IMDb is used to finetune RoBERTa⁷ and BERT⁸. Both of these models are applied on the test sets for IMDb, Yelp and Amazon. The models we used for SST-2 and TweetEval are finetuned n their corresponding training split, namely RoBERTa-SST2⁹, distilBERT-SST2¹⁰, RoBERTa-Tweet¹¹, distilBERT-Tweet¹².

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1254

1255

1256

For the grade classification task, we finetuned Llama-2 on the training split of the Hewlett dataset using QLoRA. The chosen training parameters finetune the model with a learning rate of 1e-4, batch size of 4, lora rank of 8 with lora $\alpha = 16$ and dropout 0.1. The model is trained for 1 epoch taking 30 minutes on an NVIDIA A100 machine.

D.3 Question filter

Based on manual observation, the RACE++ dataset has some questions that are generated from the linguistic contents of the contexts rather than from the semantic contents. As explained in Section 3.1, these questions will invalidate the theoretical assumption when calculating the influence of each component because the question would be unanswerable for a generated paraphrase that does not maintain the same linguistic information. Namely, these linguistic questions are often related to their positions in the context and are always correlated with certain key words such as 'in paragraph 2'. Thus, we apply a word-matching question filter to filter out all such examples, ensuring that only relevant and contextually coherent questions are re-

⁵Available at https://huggingface.co/LIAMF -USP/roberta-large-finetuned-race

⁶Available at https://huggingface.co/potsa wee/longformer-large-4096-answering-race ⁷Available at: https://huggingface.co/wrmur

ray/roberta-base-finetuned-imdb

⁸Available at:https://huggingface.co/lvwer ra/distilbert-imdb

⁹Available at: rasyosef/roberta-base-finet uned-sst2

¹⁰Available at: https://huggingface.co/distilbert/distilbert-base-uncased-finetuned sst-2-english

¹¹Available at: cardiffnlp/twitter-roberta-b
ase-dec2021-emotion

¹²Available at: https://huggingface.co/phils chmid/DistilBERT-tweet-eval-emotion

Target	Level (US)	Prompt
5	Professional	Paraphrase this document for a professional. It should be extremely difficult to read and best understood by university graduates.
20	College graduate	Paraphrase this document for college graduate level (US). It should be very difficult to read and best understood by university graduates.
40	College	Paraphrase this document for college level (US). It should be difficult to read.
55	10-12th grade	Paraphrase this document for 10th-12th grade school level (US). It should be fairly difficult to read.
65	8-9th grade	Paraphrase this document for 8th/9th grade school level (US). It should be plain English and easily understood by 13- to 15-year-old students.
75	7th grade	Paraphrase this document for 7th grade school level (US). It should be fairly easy to read.
85	6th grade	Paraphrase this document for 6th grade school level (US). It should be easy to read a
95	5th grade	Paraphrase this document for 5th grade school level (US). It should be very easy to read and easily understood by an average 11-year old student.

Table 7: Prompts to generate paraphrases with different target readability (using FRES).



Figure 7: Averaged measured readability.

tained for further processing. We specifically filter out all questions containing the following phrases: '{number} + word/sentence/paragraph + {number} + refer to/mean'.

In total, for the RACE++ dataset, approximately 6.2% questions are found to be generated from the linguistic content of the context and thus filtered out. The effects of the filter on the element influence analysis using Llama-2 is shown in Table 6. It is clear the measured question influence drops as expected by 1.6%.

E Paraphrasing

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1268

The readability level is measured by Flesch readingease (Flesch, 1948) score (FRES) where higher scores indicate material that is easier to read while lower scores are reflective of more challenging texts.

$$ext{FRES} = 206.835 - 1.015 \left(rac{n_w}{n_{se}}
ight) - 84.6 \left(rac{n_{sy}}{n_w}
ight)$$

 n_w is the total number of words, n_{se} is the total number of sentences, n_{sy} the total number of syllables.

We grouped the original texts into eight different readability levels: 5, 20, 40, 55, 65, 75, 85, and 95 for the reading comprehension and grade classification tasks and used the final 7 groups for the sentiment classification task as there were no texts in sentiment classification that fell into the most challenging category of 0-10 on the readability scale. The specific prompts for each difficulty level we used are shown in Table 7. Here we also present the quality of our paraphrase generation process. Figure 7 displays the average readability score of the paraphrased text for each combination 1269

1270







Figure 9: Averaged Word Error Rate.

of original and target readability levels. From the 1284 heatmap, we can see that while the readability of 1285 the paraphrased text is influenced by the readability 1286 of the original text, the paraphrases still fall within an acceptable range of readability. We also report 1288 the averaged BertScore F1 (Zhang et al., 2019) and 1289 Word Error Rate (WER) (Och, 2003) to ensure the 1290 quality of our paraphrasing system as shown in Fig-1291 ure 8 and Figure 9. An ideal paraphrasing system 1292 should expect high semantic similarity with high 1293 BERTScore and low linguistic similarity with high 1294 WER. 1295

F Additional results

1296

1297Here we present the results from some additional1298experiments that act as a supplement to the main1299paper.

F.1 Reading comprehension

F.1.1 Data complexity classifier

1300 1301

input format	accura	ncy	F1		
input iorniat	single	ens	single	ens	
mode class	61.6	_	25.4	_	
standard	84.7 _{±0.5}	87.2	81.7 _{±1.1}	83.7	
context	$84.9_{\pm 0.3}$	85.1	$81.8_{\pm 0.8}$	81.7	
context-question	$84.7_{\pm 0.7}$	86.0	$81.8_{\pm 0.6}$	82.2	
question-option	$70.2_{\pm 0.5}$	71.3	$67.3_{\pm 0.7}$	68.2	

Table 8: Accuracy of data complexity evaluators on the RACE++ test set.

In Section 4.3, we used the data complexity classifier with the data context as the input. Here we assess its quality by testing its performances indomain (with an ensemble of 3 models) on the RACE++ test set. We additionally compare the performance on the standard input with other possible combinations of the input, as shown in Table 1302



Figure 10: Normalized ranks (rank / total examples) of complexity scores for each complexity level using three complexity evaluators: context, context-question and standard.

dataset	model	accuracy		influence							
ualasei	moder	original	para	total	question	context	context-semantic	context-linguistic			
	RoBerta	95.3	93.0	0.254	0.140 (55.2%)	0.114 (44.8%)	0.076 (67.0%)	0.037 (33.0%)			
MCTest	Longformer	98.3	91.3	0.285	0.152 (53.5%)	0.133 (46.5%)	0.068 (70.6%)	0.028 (29.4%)			
	Llama2	92.5	85.8	0.212	0.116 (54.7%)	0.096 (45.3%)	0.068 (70.6%)	0.028 (29.4%)			
	RoBerta	84.2	81.5	0.379	0.213 (56.3%)	0.166 (43.7%)	0.135 (81.4%)	0.031 (18.6%)			
RACE++	Longformer	81.6	79.3	0.390	0.228 (58.6%)	0.162 (41.1%)	0.135 (83.2%)	0.027 (16.8%)			
	Llama2	86.0	82.9	0.298	0.161 (56.1%)	0.131 (43.9%)	0.108 (82.5%)	0.023 (17.5%)			
	Roberta	73.5	69.4	0.383	0.287 (74.9%)	0.096 (25.1%)	0.074 (77.5%)	0.022 (22.4%)			
CMCQRD	Longformer	71.9	69.8	0.467	0.326 (69.9%)	0.141 (30.1%)	0.114 (81.0%)	0.027 (19.0%)			
	Llama2	79.9	69.4	0.290	0.211 (72.7%)	0.079 (27.3%)	0.067 (83.8%)	0.012 (16.2%)			

Table 9: Decomposition of total input influence for different models in various multiple-choice reading comprehension datasets.

8. As well as accuracy, macro F1 is reported to account for the imbalance in the complexity level classes. The results for the mode class indicate the baseline performance when the mode class is selected for every example in the test set. All systems significantly outperform the baseline. Inputting the context alone is sufficient to get an accuracy close to the full input and when extra information is inputted, the gain is marginal. Hence, compared with question and options, the context carries a substantial proportion of the information to determine the complexity of a question.

In Figure 10, the data complexity classifier model shows strong generalizability by clearly classifying different subsets of CMCQRD dataset, which differ a lot from the model's training dataset. The plot also supports the context is a sufficient input to determine the different complexity levels.

F.1.2 Additional models

1309

1310

1311

1312

1313 1314

1315

1316

1317

1318

1319

1320

1322

1323

1324

1325

1327

1328

1329

1330

1331

In Section 6, we analyse the influence from different elements and components on the output for two specific models: Llama-2 for the multiple choice reading comprehension task, Roberta for the sentiment classification task. Here we show the in-
fluence terms calculated are not model-specific by
showing the consistency of the element influences1332
1333
1334on the same datasets but evaluated by different
models: Llama-2, Roberta and Longformer as in
Table 9.1336
1336

1338

1339

1340

1341

1342

1343

1344

1345

F.1.3 Further analysis

In Figure 3, we show a strong positive correlation between the question influence and the context complexity when we consider all three dataset together. Here we show the rule still holds for data points inside a single dataset in Figure 11 where the trend in RACE++ dataset is pretty similar to the all three datasets together.

We also explore other potential factors influenc-1346 ing the relative question influences. From Table 1347 2, there are two marked differences between the 1348 datasets: the number of words (length) and the number of questions. To find their influence in the 1350 relative question influence score, Figure 12a and 1351 Figure 12b show the relative question influence of 1352 the subset chosen from the contexts ranked by their 1353 number of words or the number of questions of 1354



Figure 11: The relative question influence changes with the subset chosen by the rank of context complexity in all three datasets (left) and in RACE++ only (right). ¹³ 0.2 in x-axis means we leave contexts with top 20% context complexity as the subset.



Figure 12: The relative question influence for a subset of contexts swept in order of length (left) or average number of questions per context (right) for all MCRC datasets with Llama-2.

the corresponding context. A strong positive trend between the relative question influence scores and the context length is observed as expected: a longer context naturally has a larger question generation capacity. As shown in Figure 12b, the number of questions does not have a direct impact, indicating the results are not affected by the specific number of questions per context.

F.1.4 Question Generation

1355

1356

1357

1358

1359

1361

1362

1363

1365

1369

1370

1371

1372

1373

1374

1375

1376

In Section 3.1 we investigate the influence of different element of model input and focuses on humangenerated questions only. Here we investigate the influence from the LLM-generated questions.

To be more specific, for each context, using GPT-4, we generate 4 questions with prompt:

"Given the context: {context}, please generate four multiple choice questions with options where the first option is the correct answer and the other three are distractors. The questions should be of varying difficulty levels: low, middle, high, and very high. Please output the questions in the format of a dictionary with the keys: 'easy', 'middle', 'high', and 'very high'. Each key should map to a dictionary representing a question, with the fields 'question', 'options', and 'answer' indicating the correct answer."

1377

1378

1379

1380

1381

1382

1385

However, as indicated in (Sun et al., 2023), LLMs struggle at tasks with hard constraints, Table 10 shows the generated questions are relatively easy and have a lower influence on model output.

F.1.5 Ordering

For humans taking multiple-choice tests, the role 1386 of the context compared to the question may be 1387 influenced by the ordering in which they read each 1388 of these elements. Similarly, a reading comprehen-1389 sion system may be susceptible to the ordering of 1390 the context and the question. Here we compare 1391 the influence of the ordering by reversing the stan-1392 dard context followed by question at the input to the question followed by the context. Table 11 1394 demonstrates that the ordering for the automated 1395 system does not lead to differing influences on each 1396 element. The results here are provided for the fine-1397 tuned Llama-2 model from Section 5.2. 1398

Quastian Source	accuracy			influence							
Question Source	original	para	total	question	context	context-semantic	context-linguistic				
human automatic	84.2 96.1	81.5 93.7	0.284 0.266	0.164 (57.7%) 0.121 (45.4%)	0.120 (42.3%) 0.145 (54.6%)	0.135 (81.4%) 0.101 (69.5%)	0.031 (18.6%) 0.044 (30.5%)				

Table 10: Human sources questions vs GPT4 generated questions for Llama-2 on the RACE++ test set.

dataset	direction	total	question	influe context	nce context-semantic	context-linguistic
RACE++	Forward Reverse	0.304 0.305	0.171 (56.3%) 0.172 (56.7%)	0.133 (43.7%) 0.133 (43.6%)	0.109 (82.1%) 0.109 (82.3%)	0.024 (17.9%) 0.024 (17.7%)
MCTest	Forward Reverse	0.212 0.229	0.116 (54.7%) 0.129 (56.6%)	0.096 (45.3%) 0.100 (43.4%)	0.068 (70.6%) 0.067 (67.2%)	0.028 (29.4%) 0.032 (32.3%)
CMCQRD	Forward Reverse	0.290 0.278	0.211 (72.7%) 0.204 (73.4%)	0.079 (27.3%) 0.074 (26.6%)	0.067 (83.8%) 0.061 (82.5%)	0.012 (16.2%) 0.013 (17.5%)

Table 11: Decomposition of total input influence for different models in various datasets for context-question (forward) vs question-context (reverse) using Llama-2.

dataset	model	accuracy original para		context	influence semantic	linguistic	
IMDb	RoBERTa	94.8	94.0	0.472	0.444 (94.2%)	0.028 (5.8%)	
	BERT	93.3	92.9	0.483	0.458 (94.7%)	0.025 (5.3%)	
Yelp	RoBERTa	94.3	93.9	0.472	0.445 (94.2%)	0.027 (5.8%)	
	BERT	92.9	92.6	0.518	0.488 (94.2%)	0.030 (5.8%)	
Amazon	RoBERTa	91.0	89.5	0.361	0.325 (90.0%)	0.036 (10.0%)	
	BERT	91.2	90.3	0.425	0.389 (91.5%)	0.036 (8.5%)	
SST-2	RoBERTa	87.4	82.5	0.210	0.171 (81.4%)	0.039 (18.6%)	
	BERT	89.0	84.7	0.274	0.229 (83.5%)	0.045 (16.5%)	
TweetEval	RoBERTa	85.2	74.5	0.570	0.469 (82.2%)	0.101 (17.8%)	
	BERT	77.7	75.3	0.592	0.506 (85.5%)	0.086 (14.5%)	

Table 12: Decomposition of total input influence for different models in various sentiment classification dataset

accuracy		influence		
original	para	response	response-semantic	response-linguistic
79.2	64.4	0.399	0.283 (70.9%)	0.116 (29.1%)

Table 13: Influence from semantic meaning and linguistic realization of the responses in grade classification task in Hewlett dataset.

F.2 Sentiment classification

1399

1400For the sentiment classification task, to show the1401consistency of the element influence among dif-1402ferent models, Table 12 presents additional results1403using the BERT model as a comparison against the1404RoBERTa model.

1405It can be observed also that for shorter input1406text datasets, such as SST-2 and TweetEval, the lin-1407guistic component is more significant, approaching140820% of the total.

F.3 Grade classification

1409

The influence to the model output from the re-1410 sponse, its semantic component and linguistic com-1411 ponent are respectively shown in Table 13. We ob-1412 serve a large influence from the linguistic content 1413 of the responses which agrees with the 15% drop 1414 in model accuracy when the paraphrased dataset is 1415 used. Further, we measure the correlation between 1416 the readability and true class probability in Figure 1417 13. It is clear that with a higher readability, the 1418 probability of selecting the true grade increases. 1419



Figure 13: Entropy filtered pairwise agreement in paraphrase readability and true class probability ordering with various minimum readability gaps.

G Future work

1420

1421 1422

1423

1494

1425 1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437 1438

1439

The analysis in this work has applied the framework to specifically NLP classification tasks. It would be interesting to extend the framework to both regression and sequence output tasks. For sequential outputs, there needs to be a methodology to convert the generated sequence to a single score such that its sensitivity can be measured to each input element.

The framework applied to textual data to explore the influence of semantic vs linguistic components can also be extended to image inputs. Here, we can perceive the semantic content as the object being described in the image while the linguistic realization is based on the recording equipment that controls aspects such as orientation, resolution (blurring), camera angle, e.t.c. Therefore, the proposed information-theoretic approach has potential applications across several modalities.

H Licenses

1440The RACE dataset is available for non-commercial1441research purposes only. Also for CMCQRD,1442the license14 states the licensed dataset for non-1443commercial research and educational purposes1444only.

¹⁴Available at: https://englishlanguageituto ring.com/datasets/cambridge-multiple-cho ice-questions-reading-dataset