Automated Essay Scoring in Arabic: A Dataset and Analysis of a BERT-based System

Anonymous ACL submission

Abstract

Automated Essay Scoring (AES) holds significant promise in the field of education, helping educators to mark larger volumes of essays and provide timely feedback. How-005 ever, Arabic AES research has been limited by the lack of publicly available essay 006 007 data. This study introduces AR-AES, an Arabic AES benchmark dataset compris-009 ing 2046 undergraduate essays, including gender information, scores, and transparent 010 rubric-based evaluation guidelines, provid-011 ing comprehensive insights into the scoring process. These essays come from four 014 diverse courses, covering both traditional and online exams. Additionally, we pioneer the use of AraBERT for AES, exploring its performance on different question types. 018 We find encouraging results, particularly for Environmental Chemistry and source-019 dependent essay questions. For the first time, we examine the scale of errors made by a BERT-based AES system, observing 022 that 96.15% of the errors are within one 024 point of the first human marker's prediction, on a scale of one to five, with 79.49% of 026 predictions matching exactly. In contrast, additional human markers did not exceed 028 30% exact matches with the first marker, with 62.9% within one mark. These findings highlight the subjectivity inherent in essay grading, and the potential for current AES technology to assist human markers to grade consistently across large classes.

1 Introduction

034

039

041

042

Essay writing is an important tool for developing and assessing students' cognitive abilities, including critical thinking, communication skills and depth of understanding (Ashburn, 1938; Smith et al., 1999). However, as student numbers grow, marking essays by hand becomes impractical, discouraging the use of essay questions in education (Alqahtani and Alsaif, 2019). AES systems (Page, 1966) aim to reduce the time needed to mark essays, by assessing both writing skills and cognitive outputs automatically, and can mitigate scoring biases and inconsistencies arising from teacher subjectivity (Algahtani and Alsaif, 2020). Despite extensive research in English (Wang et al., 2022; Ke and Ng, 2019), AES for Arabic, the fourth most widely used Internet language¹, remains underexplored, with most efforts concentrated on scoring short, one or two-sentence answers (Algahtani and Alsaif, 2020). With the abundant youth population in the Arab world, the education system faces challenges due to a shortage of teachers and the inability to provide individualized feedback to students (Azmi et al., 2019). In addition, the Arabic language differs from English in terms of grammar, structural rules, and the formulation of ideas, which prevents the application of scoring systems designed for English (Azmi et al., 2019). In this context, the development of an Arabic essay scoring system is an urgent necessity.

043

045

047

049

051

054

055

057

059

060

061

062

063

064

065

066

067

068

069

070

071

072

074

076

078

079

Previous research has predominantly leaned on feature engineering in conjunction with shallow models, yielding only moderate performance outcomes (Alghamdi et al., 2014; Gaheen et al., 2021). In contrast, the potential of pretrained models such as AraBERT (Antoun et al., 2020a), AraVec (Soliman et al., 2017), and AraGPT-2 (Antoun et al., 2020b), which learn vector representations from extensive text corpora, remains largely untapped within the context of Arabic AES. These models have demonstrated notable efficacy in various domains, encompassing tasks like questionanswering, named entity recognition, sentiment analysis, and even the automatic scoring of

¹Internet World State ranking, March 2020, www. internetworldstats.com

171

172

173

174

175

176

177

178

179

180

181

133

134

135

136

short answers (Meccawy et al., 2023; Alduailej and Alothaim, 2022). A major barrier to further research is the lack of publicly available datasets: datasets used in prior studies are either inaccessible or consist only of one or two-sentence answers.

081

087

880

094

095

100

101

102

103

104

105

106

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

To address these gaps, this study introduces AR-AES dataset, which consists of Arabic essays each marked by two different university teaching professionals. This dataset was collected from undergraduate students across diverse disciplines, covering various topics and writing styles. We include ancillary information, such as the gender of the students (male and female students were taught separately), the evaluation criteria employed (rubrics), and model answers for each question. The dataset comprises 12 questions and 2046 essays, collected through both traditional and online examination methods, and encompasses substantial linguistic diversity, with a total length of 115,454 tokens and 12,440 unique tokens.

This study also pioneers the use of AraBERT in Arabic AES by conducting a series of experiments to assess AraBERT's performance on our dataset at different levels of granularity, from the complete dataset down to individual courses and questions. We also examined AraBERT's performance based on gender, traditional versus online exams, and essay type (argumentative, narrative, source-dependent). AraBERT excelled when trained on several questions from the same course, achieving a quadratic weighted kappa (QWK) score of 0.971 in Environmental Chemistry. However, its performance was lower when trained specifically for certain types of question, with the lowest QWK observed for narrative questions.

Our analysis goes beyond previous work on AES, by assessing the proximity of the model's predictions to the grades assigned by the first marker, to gauge the scale of its errors. The predictions matched exactly for 79.49% of answers, with 95% of predictions having no more than one mark difference to the first human mark (with maximum five marks). In contrast, the question with highest agreement between the first and second human markers had only 30.3% exact agreement, with differences greater than one mark for 37.1% of the answers. This suggests that AraBERT-based AES is sufficiently capable to assist human markers and could help detect inconsistencies between individuals in a marking team.

In summary, our study presents a comprehensive approach to Arabic AES, introducing an open-source dataset with clear annotation guidelines and quality control, leveraging AraBERT, and providing a novel investigation of the scale of AraBERT AES errors. Our code, data, and marking guidelines are accessible at https://osf.io/dp2nh/?view_ only=4ac6373c60214ea6952855f81507fec7.

2 Related Works

Several AES datasets have been released in Chinese (Gong et al., 2021), Indonesian (Aini et al., 2018), and English, including the ASAP dataset² that has catalysed English AES research (Phandi et al., 2015; Taghipour and Tou Ng, 2016), including a new state-of-theart BERT-based approach (Wang et al., 2022). However, there is no previous publicly available dataset of Arabic essays and marks, as existing work is limited to short answers (Al-Shargabi et al., 2021). We address this gap by presenting a comprehensive Arabic AES dataset.

Arabic AES research encompasses approaches such as linear regression (Alghamdi et al., 2014), Latent Semantic Analysis (Al-Shalabi, 2016), Support Vector Machines (Algahtani and Alsaif, 2020), rule-based systems (Algahtani and Alsaif, 2019), naïve Bayes (Al-Shargabi et al., 2021), and optimization algorithms like eJaya-NN (Gaheen et al., 2021). However, these studies predominantly rely on feature engineering, using surface features that are unable to comprehensively capture the semantic nuances and structural intricacies inherent in essays. These approaches provide only limited consideration for word order, primarily revolving around word-level or grammatical features. More recent pretrained transformer models, such as BERT (Devlin et al., 2018), alleviate these issues but have not previously been harnessed for Arabic AES. Here, we develop the first AES system using AraBERT to analyse the effect of different question types on a modern text classifier. We also go beyond previous analyses of model performance by evaluating the magnitude of errors in the models' predictions, as large errors could have

²www.kaggle.com/c/asap-aes

183

185

187

190

191

194

195

196

197

198

199

a greater impact on students.

3 Arabic language challenges

NLP systems face several distinct challenges when processing Arabic, which motivate the development of bespoke tools and language resources, including benchmark datasets.

Linguistic Complexity: Arabic exhibits complex sentence structures with many syntactic and stylistic variations, an extensive vocabulary, and the frequent use of rhetorical devices (Alwakid et al., 2017). Arabic, for instance, has many ways to express the concept of "going" depending on who is doing the action, when, and whether the action is done in a habitual or momentary sense. For example, يذهب (he goes), يذهب (I will go), تأذهب (he used to go), and itikation is done in a AES system to recognise variations of the same concept.

Complex Morphology: Arabic features intricate morphology, encompassing a wide range of inflection and derivational systems (Hamdi 203 et al., 2016). Words in Arabic can have multiple forms based on factors such as tense, gender, 205 number, and case, and the form of a single let-206 ter also varies. For instance, the letter س ('S'), 207 looks like (سه) at the beginning of a word (سےاب, "Cloud"), like (س) in the middle as 209 in (مستشفر), "Hospital"), and like (حس) at the 210 end as in (شمس, "Sun"). This complexity adds 211 to the difficulty of stemming, tokenization, and 212 lemmatization operations (Kanan et al., 2019). 213 214 As another example, the Arabic root word for "write" is کتب, from which we can derive var-215 ious words like مكتوب ("writer"), مكتوب ("writer"), مكتوب 216 ten"), يكتب (book), كتبت ("I wrote"), يكتب ("he 217 writes"), etc. The challenge for AES systems 218 here lies in recognizing these words as related. 219 Non-Standard Orthography: Arabic text 220 follows complex rules for letter representation, including ligatures and diacritics that influence 222 pronunciation, word comprehension, and mean-223 ing (Isleem, 2014; Soudi et al., 2008). NLP systems face challenges in handling these ortho-225 graphic differences and the absence of diacritics in unvocalised text. For example, محبوبة ("loved 227 or popular") could be written as محبو سه in casual writing without the ending diacritic.

Lack of Resources: Arabic suffers from
limited linguistic resources, such as preprocess-

ing tools for dealing with the language complexities described above, and a lack of public datasets (Mahmoud and Zrigui, 2019; Kaddoura et al., 2022), which hampers the development of NLP. A particular need is for bespoke tools to deal with the right-to-left text direction, which creates additional complexities for mixed-language content (Awwad et al., 2017; Kanan et al., 2019). This study contributes a labelled Arabic dataset, which will further the development of Arabic NLP systems.

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

Ambiguity and Polysemy: Arabic words often possess multiple meanings and interpretations, making it challenging to disambiguate them (Elkateb et al., 2006). For example, the word جعل in Arabic can mean "camel" or "sentence" depending on context. Contextual analysis becomes crucial for accurately determining the intended meaning (Kaddoura et al., 2022; Omar and Aldawsari, 2020). This aspect presents a challenge in various NLP tasks, including named entity recognition, sentiment analysis, and machine translation.

Despite these challenges, substantial advancements have been made in Arabic NLP in recent years, including language models and tools specifically designed for Arabic. This study hopes to contribute to this effort.

4 The AR-AES Dataset

The AR-AES dataset is intended for both training and evaluating Arabic AES systems, and covers essays written by both male and female undergraduate students from three different university faculties, with a range of different question types, a mix of traditional face-to-face and online exams, and marks from multiple human markers. As part of the dataset, we include clear and detailed marking criteria along with model answers for each question. This diversity will enable researchers to explore the suitability of AES systems for different types of essays, exam types, or student cohorts.

Data Collection: To compile a diverse dataset, we first selected multiple undergraduate courses across various departments at Umm Al-Qura University (Table 1). Students' writing skills vary depending on their academic disciplines (Zhu, 2004), due to differing objectives, terminology, and research formulation methodologies. The difficulty of marking an

essay may also be affected by whether it is of an argumentative, source-dependent, or narrative type (Mathias and Bhattacharyya, 2018).
Additionally, factors like gender and academic level contribute to differences in writing (Johnson, 1999; Lea and Street, 1998), particularly considering that male and female students are taught separately. Therefore, to test AES systems across various subjects and writing styles, we collected essay responses from diverse academic levels, genders, and question types.

283

287

291

296

297

299

301

304

305

307

308

310

311

312

313

To bolster dataset diversity, we employed both traditional (in-person) and online exams through distance learning. Traditional exams occurred on specific dates on campus, subjecting students to controlled conditions that minimized opportunities for academic misconduct. Conversely, online exams required students to submit essay responses exclusively via content management platforms. These exams shared time limits with traditional exams but did not mandate physical presence on campus. Online exams can reduce stress levels (Ilgaz and Adanir, 2020), granting students greater freedom in providing answers and potentially allowing access to course content during the exam. For both kinds of exam, answers were typed and submitted electronically. These essays were part of the students' compulsory assessment within the midterm exams for their respective courses, and they volunteered to provide their essays for our dataset.

The Annotation Task: Course directors 314 315 equipped markers with detailed guidelines for scoring individual criteria and determining the 316 final score. Table 2 shows an example of the cri-317 teria for assessing Question 1, which prompts 318 students to "Explain in detail the difference 319 between the terms 'data' and 'information', 320 supplementing their answers with examples of 321 each type". For an exhaustive overview of the 322 Scoring Criteria, see Table A.3. This structured approach facilitates the identification of 324 essay strengths and weaknesses. For all questions, the first marker is the course provider; we also collected marks from a second faculty 328 member familiar with the content, to allow us to compare the performance of AES systems 329 against additional human markers, who could work as a team to mark large numbers of essays. In total, a team of 9 faculty members 332

formulated, prepared, and scored the exams.

Quality control: To guarantee the quality 334 of essay questions, individual meetings were 335 conducted with the faculty members respon-336 sible for each course. The course directors 337 were provided with the following criteria for 338 formulating essay questions, and the proposed 339 questions were verified by the authors of this 340 paper against these criteria, and revised if they 341 did not meet the criteria. 342

333

343

344

345

346

347

348

349

350

351

353

354

355

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

1. **Clear objectives:** Each question should have a clear objective aimed at assessing a specific cognitive skill, such as analysis, synthesis, or evaluation. This clarity helps students focus on comprehending the question and providing the required answer directly.

2. **Relevance:** Ensure that the question directly relates to the course content and learning objectives.

3. Explicit terminology: In the question, incorporate explicit terminology relevant to the course content.

4. Clarity and simplicity: Questions should be straightforward, unambiguous, and include a comprehensive outline of the expectations for the answer. This approach encourages concise and easily evaluated responses.

5. Linguistic accuracy: Ensure that questions are free of grammatical errors to prevent unintended alterations in question meaning.

6. Alignment with learning outcomes: Align each question with the specific learning outcomes you want to assess.

7. Fairness: Craft questions that offer all students an equal opportunity to demonstrate their knowledge and skills.

8. Grading guide: For each question, a guide should be developed to communicate the correct answer structure and criteria for achieving higher grades, clarifying the grading process.

Special instructions were developed for online exams to prevent cheating. These measures included restricting exam access to one hour on the Blackboard platform and requiring students to have their cameras on throughout the

Course	Faculty	Semester	Exam	No.	Gender	No.	Question	Essay Type	Answ	er Length	Score	e Range
			Type	Groups		Students	ID	ID		\mathbf{Min}	\mathbf{Min}	Max
				2	Mala	151	A	ll questions	298	2	0	5
Introduction	Computing	1	Traditional	э	Male	101	Q1	Narrative	298	7	0	5
to Info Science	Computing	1	Traditional	9	Emile	198	Q_2	Argumentative	164	2	0	5
				4	remate	120	Q3	Source Dependent	61	4	0	5
				2	Male		A	ll questions	512	16	0	10
Management Info Systems	Business Ad-	5	Traditional			181	Q4	Narrative	512	29	0	10
	ministration	5					Q_5	Narrative	212	29	0	10
							Q6	Source Dependent	171	16	0	5
					Male		А	ll questions	422	8	0	5
Environmental	Applied	7	Online	0		116	Q7	Narrative	422	25	0	5
Chemistry	Science	1	Onnie	2			Q8	Argumentative	116	9	0	5
							Q9	Source Dependent	92	8	0	5
							A	ll questions	575	11	0	5
Distant sharely and	Applied	6	Onlino	0	M.1.	106	Q10	Source Dependent	357	13	0	5
Diotechnology	Science	0	Onnie	4	wate	100	Q11	Argumentative	538	11	0	5
							Q12	Source Dependent	575	13	1	5

Table 1: Course summary, including the semester the exam was taken in (out of 8 in an undergraduate degree), number of groups taught at separate times (no. groups), and answer length (number of tokens).

Rubric-based evaluations	Score	Course Name	Questions	Essay	Gen	der	Exam type		
قدرة الطالب على التعريف بالبيانات ودورها وأشكالها			Count	Count	Μ	F	Trad.	Online	
The student's ability to define data, its role and forms	1	Introduction to In- formation Science	3	837	453	384	837		
قدرة الطالب على التعريف بالمعلومات ونشأتها وأوجه استخدامها The student's ability to identify information.	1	Management infor- mation systems	3	543	543		543		
its origins, and its uses	,	Environmental	0	940	940			940	
قدرة الطالب على استنتاج الفرق بين البيانات والمعلومات		chemistry	3	348	348			348	
The student's ability to deduce the difference	2	Biotechnology	3	318	318			318	
between data and information		Total	12	2046	1662	384	1380	666	
قدرة الطالب على تعزيز شرحه للبيانات والمعلومات بأمثلة واقعية ذات صلة The student's ability to reinforce his expla- nation of data and information with relevant	. 1	Table 3: The	number o	f essay	v resp	onse	es per	course.	
examples		but its appli	cation to	AES	rema	ains	unex	plored	
Final Score	5		1 C			Т	рри	17	

Table 2: Example marking criteria set by the course director for Q1.

exam. Students were explicitly instructed not to engage in chat conversations or pose questions during the examination. Any inquiries or concerns related to the test were to be addressed only after the exam had concluded.

Dataset Statistics: In total, we collected and labelled 2046 essays (Table 3). Table 1 shows notable variations in answer lengths, measured in tokens, across different question types, and between online and traditional exam types, with online exam responses generally being longer across most questions. The class distribution is illustrated in Figure B.1.

5 Experimental Setup

378

379

384

388

390

391

395

The AraBERT model has consistently demonstrated state-of-the-art performance in various Arabic NLP tasks, including the automatic scoring of short answers (Meccawy et al., 2023), but its application to AES remains unexplored. This study therefore assesses AraBERT's performance in AES, testing its ability to handle longer Arabic texts, and whether performance varies depending on factors such as the subject, question type, exam type, or gender. 396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

Data Preprocessing: We removed punctuation, hashtags, URLs, excess letter repetitions, emoticons, superfluous spaces, numbers, and diacritics, and normalized specific Arabic characters to their standard forms (e.g., $\zeta - \circ > \ddot{\circ}$ acters to their standard forms (e.g., $\zeta - \circ > \ddot{\circ}$ acters to their standard forms (e.g., $\zeta - \circ > \ddot{\circ}$ acters to their standard forms (e.g., $\zeta - \circ > \ddot{\circ}$ acters to their standard forms (e.g., $\zeta - \circ > \ddot{\circ}$ acters to their standard forms (e.g., $\zeta - \circ > \ddot{\circ}$ be applied the ISRI Stemmer, in the manner of previous work (Meccawy et al., 2023), to simplify Arabic text by reducing words to their roots to minimise vocabulary diversity. We employed the AraBERT tokenizer, and sequences exceeding 512 tokens were truncated. Most essays fit this limit, except four from the Biotechnology course, exceeding up to 575 words.

Model Design:AraBERT is a variant of416BERT that was pretrained on a substantial417Arabic text dataset (Antoun et al., 2020a) and418can be fine-tuned for specific tasks with mini-419

mal additional training data, reducing the time 420 and resources needed for model development 421 and deployment. This study used the large 422 AraBERT configuration, featuring 12 encoder 423 blocks, 1024 hidden dimensions, 16 attention 494 heads, 512 sequence length, and 370 million 425 parameters. We added a classification head on 426 top of AraBERT, consisting of a single fully-427 connected layer. Notably, this study marks the 428 first application of AraBERT to AES. 429

Model Training: The system aims to as-430 sist the course presenter (first annotator), so 431 the model was trained only on the labels pro-432 vided by that person. To ensure comparability 433 across questions, we normalized all scores in 434 the dataset to the range 0 to 5. For questions 435 with scores originally ranging from 0 to 10 (Q4 436 and Q5), we divided the scores by 2 to align 437 them with the score range used for other essays. 438 We trained the model once on the complete 439 dataset (a general-purpose model), as well as 440 separately for each course and each question. 441 We also trained the model separately on male 442 and female essay responses for the Introduc-443 tion to Information Science course and on tra-444 ditional and online essay responses, to observe 445 446 differences in model performance that could affect each group differently. 447

> For each of these experiments, we divided the answers randomly into training, validation and test sets (70/15/15). We trained using Adam optimiser and tuned the hyperparameters, including batch size, dropout rate (0.2), and number of epochs, on the validation set for each experiment, as detailed in Table A.2. Given the dataset's imbalanced nature, we employed class weights to give equal weight to each class in the dataset by assigning proportionally higher weights to instances from smaller classes. The distribution of classes for each question is illustrated in Figure B.1⁵.

448

449

450

451

452

453

454

455

456

457

458

459

460

461 **Evaluation Metrics:** We adopted quadratic weighted kappa (QWK) and F1 score as evalua-462 tion metrics. QWK, an extension of Cohen's κ , 463 gauges the level of agreement between the scor-464 ing outcomes of two assessors (Cohen, 1968). 465 466 This metric is commonly employed in AES evaluation because, unlike accuracy and F1 467 score, κ considers chance agreement, providing 468 a more reliable measure of rating concordance 469 (Mathias and Bhattacharyya, 2020). More-470

The Experiment	Unique words	F1	QWK ³
The Entire Dataset	12440	0.78	0.884
Introduction to Information Science	3953	0.61	0.788
Management Information System	4469	0.59	0.779
Environmental chemistry	2702	0.95	0.971
Biotechnology	4241	0.85	0.953
Question 1	1922	0.59	0.887
Question 2	1906	0.47	0.733
Question 3	938	0.82	0.870
Question 4	2331	0.82	0.833
Question 5	1878	0.85	0.841
Question 6	978	0.95	0.942
Question 7	1801	0.33	0.425
Question 8	767	0.88	0.791
Question 9	507	0.91	0.979
Question 10	772	0.77	0.902
Question 11	1787	0.57	0.843
Question 12	2483	0.76	0.838
Female	2723	0.59	0.741
Male	3033	0.53	0.715
Traditional Exam	7506	0.57	0.758
Online Exam	6355	0.72	0.929
Narrative (Q1,Q4,Q5,Q7)	6790	0.45	0.693
Argumentative (Q2,Q8,Q11)	3863	0.64	0.732
Source Dependent (Q3,Q6,Q9,Q10,Q12)	4667	0.73	0.889

Table 4: Comparison of AraBERT models trained on different question subsets.

over, QWK accommodates the ordinal nature of classes, crucial to essay scoring, and employs quadratic weights to reflect class rank order, which is ignored by accuracy and F1 scores. QWK is computed by:

$$QWK = 1 - \frac{\sum_{i,j} w_{i,j} O_{i,j}}{\sum_{i,j} w_{i,j} n_{i,1} n_{j,2}},$$
 (1)

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

where $w_{i,j} = \frac{(i-j)^2}{(N-1)^2}$ is the weight between mark *i* and mark *j*, *N* is the number of marks available, $O_{i,j}$ is the number of observations where the first assessor gave mark *i* and the second assessor gave mark *j*, and $n_{i,k}$ is the number of times that assessor *k* gave mark *i*.

6 Results

We first evaluated the AraBERT model on the entire dataset to gauge its performance when trained with more data and a variety of questions. Then, we trained and evaluated models using data from each course, individual question, question type, student gender, and exam type, to identify the kind of scenarios where the AES system could be more effective.

Results are shown in Table 4. On the complete dataset, the model achieved QWK=0.884492and F1=0.78. The averages for models trained494separately per course were QWK=0.873 and495

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

548

549

550

F1=0.75, and for models trained separately per question, the averages were QWK=0.824 and F1 =0.73. This suggests that larger training sets may be beneficial, even if these incorporate a mix of questions or subjects.

496

497

498

499

501

502

503

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

523

524

525

526

527

529 530

531

532

533

534

535

536

537

538

539

541

542

543

544

546

547

Scores vary substantially between questions. For instance, performance in the Environmental Chemistry course exceeded that of the entire dataset, even though this course included responses in Arabic mixed with English terms. Among the different courses, performance was weakest on Management Information Systems, potentially due to the complexity of the student responses, which had an average of 4469 unique words (in extended answers), while Environmental Chemistry had around 2702 unique words (in restricted answers). This difference may be because the Management Information Systems course featured more open-ended essays, with two narrative questions, than Environmental Chemistry, where answers were more source-dependent and controlled, making them easier for the model to evaluate.

Compared to Biotechnology, performance on Information Science was weaker, despite its larger training set. We investigated whether this is due to the students' use of informal language, considering that this course is a firstsemester offering for first-year undergraduates, while the Biotechnology course is taken in the second semester of the third year. We computed the perplexity (Miaschi et al., 2021) of students' answers for each course, finding that Introduction to Information Science had a high perplexity score of 14.87 compared to Management Information System (1.77), Environmental Chemistry (1.5), and Biotechnology (1.68). This suggests that the AraBERT model was less suitable for modelling the Introduction to Information Science answers, and that the language differs from that used in other courses.

Overall, the model performed best with source-dependent questions, where language is more constrained, and worst with narrative questions, which were the most open-ended, with a higher number of unique words in Table 4. The model also performed better with online, rather than traditional in-person exams. Access to course materials online may increase answer consistency. Splitting the Introduction to Information Science questions by gender resulted in superior performance when predicting female students' marks, which may reflect different teaching or learning styles, as male and female students are taught separately by different lecturers.

Magnitude of Errors: It is important to consider the scale of errors that the model makes: if the system predicts marks that are much lower or higher than the human marker, students could be unfairly penalised or rewarded for poor-quality work. We therefore assess the deviations between predictions and correct scores in Figure 1. The pattern is similar across courses. The majority of errors involved overestimations, with 10% of cases resulting in a one-mark overestimation. Underestimations were less frequent, occurring in 6% of cases with a one-degree reduction. Exact matches were 12% higher for Environmental Chemistry than Introduction to Information Science. Examining the error distribution for each essay type (Figure 2), one-mark overestimates occur noticeably more in narrative essays, while source-dependent essay predictions match the human marker's grade in 87% of cases.

7 Comparison with Second Markers

Here, we ask whether the AraBERT model can compete with a second human marker, in terms of consistency with the course director's marks. We examine the accuracy of second markers for two courses: Introduction to Information Science and Information Systems Management, for a total of six questions (Q1 to Q6) and show results in Table 5. The highest accuracy was observed in Q3 (source-dependent), where both markers provided the same grade in 30.3% of cases out of 279 responses. Negative differences were far more frequent than positive, meaning that second markers tended to mark more harshly than the course directors. Conversely, the lowest agreement was found in Question 4 (narrative), which has notably more cases of disagreement by 3 or more marks.

When compared with the performance of our models, which were trained with the gold standard marks of the original markers, we see that the disparity in second marker's assessments often exceeds the error rate of the automated system. This suggests that the model could effectively assist a human marker or help to ensure consistency between multiple markers.



Figure 1: Distribution of model predictions for course-level grades.



Figure 2: Distribution of Model Predictions for Essay-type Level.

	Correlation	QWK	Question Type	-5	-4	-3	-2	-1	Matching	1	2	3	4	5
Question 1	0.574	0.543	Narrative	0	1.5	4.8	20.4	25.2	23.3	13.3	8.1	3.3	0	0
Question 2	0.639	0.618	Argumentative	0	1.9	6.7	16.7	22.6	25.6	17.8	7.0	1.1	0.7	0
Question 3	0.775	0.690	Source dependent	0	0.4	9.0	25.1	25.1	30.3	7.5	1.9	0.7	0	0
Question 4	0.577	0.174	Narrative	0.6	2.6	12.3	19.5	24.7	20.8	9.1	7.1	1.3	1.9	0
Question 5	0.834	0.252	Narrative	1.3	0	3.9	16.2	24.0	31.8	21.4	7.1	0	0.6	0
Question 6	0.734	0.665	Source dependent	0	0	0	3.3	12.5	27.6	50.7	5.3	0.7	0	0

Table 5: The extent of agreement and discrepancy between the scores of the two human assessors is compared, in addition to the correlation, and QWK.

8 Conclusions and Future work

In this paper, we introduced AR-AES, the first publicly-available Arabic AES dataset, consisting of 2046 undergraduate essays with model answers, marking criteria, and scores from multiple markers. We also developed and evaluated an AES system using AraBERT, and demonstrated promising performance, particularly on source-dependent essays in domains such as Environmental Chemistry. Our analysis showed that agreement between our model and gold standard marks is higher than agreement among human markers, suggesting a role for AES in ensuring consistency as well as increasing marking efficiency.

There are numerous avenues for future work, such as exploring the adaptation of state-of-theart techniques from the English AES field to the domain of Arabic AES, such as the multi-scale approach of (Wang et al., 2022). In addition to model exploration, future research should also focus on integrating AES systems into the essay grading process effectively, and addressing students' and teachers' concerns about automated systems. This includes designing a process for identifying and rectifying errors, and ensuring that human teachers retain control while being assisted in grading a large set of essays. This area holds significant potential for enhancing the efficiency and accuracy of essay scoring, particularly in universities with limited teaching resources. We also see value in expanding our dataset with essays from a wider range of courses and educational institutions, thereby enhancing the robustness and versatility of our model, and investigating other aspects of student diversity beyond subject and gender. Our approach may also provide a template for AES data collection in other languages.

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

- 637 638
- 63

642

643

645

646

648

652

654

657

664

665

671

672

673

674

675

679

683

684

687

9 Ethical Considerations and Limitations

Maintaining high-quality data was a top priority throughout this study's data collection process. To ensure ethical compliance and research integrity, the entire data collection plan underwent scrutiny by an ethical review at the University of Bristol in the United Kingdom, resulting in their approval. This endorsement confirmed adherence to UK scientific research ethical standards, with an unwavering commitment to preserving participant anonymity and confidentiality.

Before choosing to deploy an AES system, it is important to consider what happens if a model makes a mistake. Using AES could help human markers to reduce mistakes and mark more consistently, but feedback to students that explains their marks transparently, alongside a clear appeals process, may also be required to ensure that automated tools do not introduce unfair marking decisions. It is possible that automated systems could also be tricked into giving high marks by including the right phrases in an essay, so human oversight of the AES system will be important to guard against this. A limitation of our work is that we did not uncover specific cases of the problems mentioned above; we present the dataset to facilitate future work into such topics, e.g., by investigating model performance with adversarial examples.

References

- Qurratul Aini, Achmad Eko Julianto, and Dwijoko Purbohadi. 2018. Development of a scoring application for indonesian language essay questions. pages 6–10.
- Emad Fawzi Al-Shalabi. 2016. An Automated System for Essay Scoring of Online Exams in Arabic based on Stemming Techniques and Levenshtein Edit Operations. *International Journal of Computer Science Issues*, 13(5):45–50.
- Bassam Al-Shargabi, Rawan Alzyadat, and Fadi Hamad. 2021. AEGD: ARABIC ESSAY GRAD-ING DATASET FOR MACHINE LEARNING A Cultural E-Government Readiness Model View project A comparative study for Arabic text classification algorithms based on stop words elimination View project.
- Alhanouf Alduailej and Abdulrahman Alothaim. 2022. AraXLNet: pre-trained language model

for sentiment analysis of Arabic. Journal of Big Data, 9(1).

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

729

730

733

734

735

736

737

738

- Mansour Alghamdi, Mohamed Alkanhal, Mohamed Al-Badrashiny, Abdulaziz Al-Qabbany, Ali Areshey, and Abdulaziz Alharbi. 2014. A hybrid automatic scoring system for Arabic essays. *AI Communications*, 27(2):103–111.
- Abeer Alqahtani and Amal Alsaif. 2019. Automatic Evaluation for Arabic Essays: A Rule-Based System.
- Abeer Alqahtani and Amal Alsaif. 2020. Automated Arabic essay evaluation. In *Proceedings* of the 17th International Conference on Natural Language Processing (ICON), pages 181–190, Indian Institute of Technology Patna, Patna, India. NLP Association of India (NLPAI).
- Ghadah Alwakid, Taha Osman, and Thomas Hughes-Roberts. 2017. Challenges in sentiment analysis for Arabic social networks. In *Procedia Computer Science*, volume 117, pages 89–100. Elsevier B.V.
- Wissam Antoun, Fady Baly, and Hazem Hajj. 2020a. AraBERT: Transformer-based Model for Arabic Language Understanding.
- Wissam Antoun, Fady Baly, and Hazem Hajj. 2020b. Aragpt2: Pre-trained transformer for arabic language generation. arXiv preprint arXiv:2012.15520.
- Robert Ashburn. 1938. An experiment in the essaytype question. *The Journal of Experimental Education*, 7(1):1–3.
- Aiman M.Ayyal Awwad, Christian Schindler, Kirshan Kumar Luhana, Zulfiqar Ali, and Bernadette Spieler. 2017. Improving pocket paint usability via material design compliance and internationalization & localization support on application level. In Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services, Mobile-HCI 2017. Association for Computing Machinery, Inc.
- Aqil M. Azmi, Maram F. Al-Jouie, and Muhammad Hussain. 2019. AAEE – Automated evaluation of students' essays in Arabic language. *Information Processing and Management*, 56(5):1736–1752.
- Jacob Cohen. 1968. Weighted kappa: nominal scale agreement provision for scaled disagreement or partial credit. *Psychological bulletin*, 70(4):213.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.

- 740 741 742
- 74
- 744 745 746 747 748 749 750 751

- 754 755 756 757 758 759 760 761
- 762 763 764 765
- 767 768 769 770 771 772 773 774
- 775 776 777 778 779 780 781 782
- 782 783 784
- 785 786

7

- 790
- 791 792
- 7

793 794

- Sabri Elkateb, William Black, Horacio Rodríguez, Musa Alkhalifa, Piek Vossen, Adam Pease, and Christiane Fellbaum. 2006. Building a WordNet for Arabic.
- Marwa M. Gaheen, Rania M. ElEraky, and Ahmed A. Ewees. 2021. Automated students Arabic essay scoring using trained neural network by e-jaya optimization to support personalized system of instruction. *Education and Information Technologies*, 26(1):1165–1181.
- Jiefu Gong, Xiao Hu, Wei Song, Ruiji Fu, Zhichao Sheng, Bo Zhu, Shijin Wang, † ¶£, and Ting Liu. 2021. IFlyEA: A Chinese Essay Assessment System with Automated Rating, Review Generation, and Recommendation. pages 240–248.
 - Ali Hamdi, Khaled Shaban, and Anazida Zainal. 2016. A review on challenging issues in Arabic sentiment analysis. 12(9):471–481.
- Hale Ilgaz and Gülgün Afacan Adanır. 2020. Providing online exams for online learners: Does it really matter for them? *Education and Information Technologies*, 25:1255–1269.
- Martin Isleem. 2014. Developing attitudes toward learning Arabic as a foreign language among american university and college students.
- Margaret J Johnson. 1999. Gender differences in writing self-beliefs of elementary school students.
- Sanaa Kaddoura, Rowanda D. Ahmed, and D. Jude Hemanth. 2022. A comprehensive review on Arabic word sense disambiguation for natural language processing applications. 12(4).
- Tarek Kanan, Odai Sadaqa, Amal Aldajeh, Hanadi Alshwabka, Shadi AlZu'bi, Mohammed Elbes, Bilal Hawashin, Mohammad A Alia, and others.
 2019. A Review of Natural Language Processing and Machine Learning Tools Used to Analyze Arabic Social Media.
- Zixuan Ke and Vincent Ng. 2019. Automated essay scoring: A survey of the state of the art. In *IJCAI*, volume 19, pages 6300–6308.
- Mary R. Lea and Brian V. Street. 1998. Student writing in higher education: An academic literacies approach. Studies in Higher Education, 23:157–172.
- Adnen Mahmoud and Mounir Zrigui. 2019. Sentence embedding and convolutional neural network for semantic textual similarity detection in arabic language. *Arabian Journal for Science* and Engineering, 44:9263–9274.
- Sandeep Mathias and Pushpak Bhattacharyya. 2018. ASAP++: Enriching the ASAP automated essay grading dataset with essay attribute scores. In Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018).

Sandeep Mathias and Pushpak Bhattacharyya. 2020. Can neural networks automatically score essay traits? In Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications, pages 85–91, Seattle, WA, USA \rightarrow Online. Association for Computational Linguistics. 795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

- Maram Meccawy, Afnan Ali Bayazed, Bashayer Al-Abdullah, and Hind Algamdi. 2023. Automatic essay scoring for Arabic short answer questions using text mining techniques. *International Journal of Advanced Computer Science and Applications*, 14(6).
- Alessio Miaschi, Dominique Brunato, Felice Dell'orletta, and Giulia Venturi. 2021. What Makes My Model Perplexed? A Linguistic Investigation on Neural Language Models Perplexity. pages 40–47.
- Abdulfattah Omar and Mohammed Aldawsari. 2020. Lexical Ambiguity in Arabic Information Retrieval: The Case of Six Web-Based Search Engines. International Journal of English Linguistics, 10(3):219.
- Ellis B Page. 1966. The Imminence of... Grading Essays by Computer. 47(5):238–243.
- Peter Phandi, Kian Ming A Chai, and Hwee Tou Ng. 2015. Flexible Domain Adaptation for Automated Essay Scoring Using Correlated Linear Regression. pages 17–21.
- David Smith, Jennifer Campbell, and Ross Brooker. 1999. The impact of students' approaches to essay writing on the quality of their essays. Assessment & Evaluation in Higher Education, 24(3):327–338.
- Abu Bakr Soliman, Kareem Eissa, and Samhaa R El-Beltagy. 2017. Aravec: A set of Arabic word embedding models for use in arabic nlp. *Procedia Computer Science*, 117:256–265.
- Abdelhadi Soudi, Günter Neumann, and Antal van den Bosch. 2008. Arabic Computational Morphology: Knowledge-Basedand Empirical Methods. *Computational Linguistics*, 34(3):459–462.
- Kaveh Taghipour and Hwee Tou Ng. 2016. A Neural Approach to Automated Essay Scoring. pages 1882–1891.
- Yongjie Wang, Chuan Wang, Ruobing Li, and Hui Lin. 2022. On the use of bert for automated essay scoring: Joint learning of multi-scale essay representation.
- Wei Zhu. 2004. Faculty views on the importance of writing, the nature of academic writing, and teaching and responding to writing in the disciplines. Journal of Second Language Writing, 13:29–48.

Appendix A. Tables

Course Name	Ques- tion ID	The Questions	The Questions in Arabic
n to cience	1	Explain in detail the difference between the terms data and informa- tion and reinforce your answers with examples for each type?	اشرح بشكل مفصل الفرق بين كلا من مصطلحي البيانات والمعلومات مع تعزيز إجاباتك بأمثلة لكل نوع؟
ductio ation S	2	Explain in detail the role of the increase in subspecialties and the increase in topics influencing the information revolution (explosion)?	أشرح باستفاضة دور زيادة التخصصات الدقيقة وتزايد الموضوعات في التأثير على ثورة انفجار المعلومات؟
Intro Informa	3	Through what you learned in the course, mention the comprehensive definition of the term information science?	من خلال ما تعلمته ضمن المقرر الدراسي أذكري التعريف الشامل لمصطلح علم المعلومات؟
nent systems	4	The administrative levels' tasks, roles, and duties differ in manage- ment, so explain in detail the difference between the roles and tasks of the different administrative levels while strengthening your answer with examples?	تختلف مهام وأدوار وواجبات المستويات الإدارية في الإدارة لذلك أشرح بشكل مفصل الاختلاف بين أدوار ومهام المستويات الإدارية المختلفة مع تعزيز اجابتك بأمثله؟
lanager nation	5	Mention three main benefits of cloud computing from a business management perspective with an explanation?	أذكر ثلاثة من الفوائد الرئيسية للحوسبة السحابية من منظور إدارة الأعمال مع الشرح؟
M inforr	6	Through what you have learned in the course, mention the compre- hensive definition of the term information technology and reinforce your answer with examples?	من خلال ما تعلمته ضمن المقرر الدراسي أذكري التعريف الشامل لمصطلح علم المعلومات؟
instry	7	Talk about the layers of the atmosphere, mentioning the height and temperature of each layer?	تحدث عن طبقات الغلاف الجوي مع ذكر ارتفاع كل طبقة ودرجة الحرارة فيها؟
Environ tal chem	8 9	What do you think about the importance of the ozone layer? What is the scientific definition of environmental chemistry?	ما رأيك في أهمية طبقة الأوزون؟ ما هو التعريف العلمي لكيمياء البيئة؟
gy	10	Define the term biotechnology?	عرف مصطلح التقنية الحيوية؟
oloi	11	Discuss whether eating genetically modified fruits is healthy or not?	ناقش تناول الفواكه المعدلة الوراثية صحى أم لا؟
Biotechn	12	Mention five of the applications of biotechnology in the medical field with explanation?	عدد خمسة من تطبيقات التقنية الحيوية فيَّ المجال الطبي مع الشرح؟

Table A.1: List of Questions Used in Each Course to Collect Essay Answers.

Early Stop [*]	24	49	15	33	23	10	43	13	47	44	10		10	10 27	10 27 44	10 27 60	10 27 44 60 42	10 27 60 42 42 47	10 27 44 60 42 47 25	10 27 44 60 42 47 25 10	10 27 44 60 42 47 25 10 22	10 27 44 60 60 47 25 10 12 14	10 27 44 60 60 47 25 10 12 22 18	10 27 44 60 60 47 47 25 10 11 22 24 24
Test Size	306.9	125.55	81.45	52.2	47.7	41.85	41.85	41.85	27.15	27.15	07 1E	CL.12	01.12 17.4	ст. 12 17.4 17.4	01.12 17.4 17.4 17.4	ct12 17.4 17.4 17.4 15.9	27.13 17.4 17.4 17.4 15.9 15.9	27.13 17.4 17.4 17.4 15.9 15.9 15.9	27.15 17.4 17.4 17.4 15.9 15.9 15.9 15.9 27.6	27.15 17.4 17.4 17.4 15.9 15.9 15.9 15.9 67.6 67.95	27.15 17.4 17.4 17.4 15.9 15.9 15.9 15.9 57.6 67.95 207	27.15 17.4 17.4 17.4 15.9 15.9 15.9 15.9 57.6 67.95 67.95 99.9	27.19 17.4 17.4 17.4 15.9 15.9 15.9 15.9 57.6 67.95 67.95 99.9 113.55	27.19 17.4 17.4 17.4 15.9 15.9 15.9 57.6 67.95 67.95 67.95 113.55 75.15
Val Size	306.9	125.55	81.45	52.2	47.7	41.85	41.85	41.85	27.15	27.15		27.15	27.15 17.4	27.15 17.4 17.4	27.15 17.4 17.4 17.4	27.15 17.4 17.4 17.4 15.9	27.15 17.4 17.4 17.4 15.9 15.9 15.9	27.15 17.4 17.4 17.4 15.9 15.9 15.9	$\begin{array}{c} 27.15\\ 17.4\\ 17.4\\ 17.4\\ 17.4\\ 15.9\\ 15.9\\ 15.9\\ 15.9\\ 57.6\\ 57.6\end{array}$	$\begin{array}{c} 27.15\\ 17.4\\ 17.4\\ 17.4\\ 17.4\\ 15.9\\ 15.9\\ 15.9\\ 15.9\\ 57.6\\ 57.6\\ 67.95\end{array}$	27.15 17.4 17.4 17.4 15.9 15.9 15.9 57.6 67.95 207	$\begin{array}{c} 27.15\\ 17.4\\ 17.4\\ 17.4\\ 17.4\\ 15.9\\ 15.9\\ 15.9\\ 57.6\\ 67.95\\ 207\\ 99.9\end{array}$	$\begin{array}{c} 27.15\\ 17.4\\ 17.4\\ 17.4\\ 17.4\\ 15.9\\ 15.9\\ 15.9\\ 57.6\\ 67.95\\ 67.95\\ 207\\ 207\\ 113.55\end{array}$	$\begin{array}{c} 27.15\\ 17.4\\ 17.4\\ 17.4\\ 17.4\\ 15.9\\ 15.9\\ 15.9\\ 57.6\\ 67.95\\ 207\\ 207\\ 207\\ 75.15\\ 75.15\end{array}$
Training Size	1432.2	585.9	380.1	243.6	222.6	195.3	195.3	195.3	126.7	126.7		126.7	126.7 81.2	126.7 81.2 81.2	126.7 81.2 81.2 81.2	126.7 81.2 81.2 81.2 74.2	126.7 81.2 81.2 81.2 74.2 74.2	126.7 81.2 81.2 81.2 74.2 74.2 74.2	126.7 81.2 81.2 81.2 74.2 74.2 74.2 268.8	126.7 81.2 81.2 81.2 74.2 74.2 74.2 268.8 317.1	126.7 81.2 81.2 81.2 74.2 74.2 74.2 268.8 317.1 966	126.7 81.2 81.2 81.2 74.2 74.2 74.2 268.8 317.1 966 466.2	126.7 81.2 81.2 81.2 74.2 74.2 74.2 268.8 317.1 966 466.2 529.9	$\begin{array}{c} 126.7\\ 81.2\\ 81.2\\ 81.2\\ 81.2\\ 81.2\\ 74.2$
N. Essays	2046	837	543	348	318	279	279	279	181	181		181	$\frac{181}{116}$	181 116 116	181 116 116 116	181 116 116 116 116	181 116 116 116 106	181 116 116 116 106 106	181 116 116 116 106 106 384	181 116 116 116 106 106 384 453	181 116 116 116 106 106 384 453 1380	181 116 116 116 106 106 384 453 1380 666	181 116 116 116 106 106 106 106 1384 453 1380 666 666	181 116 116 116 116 106 106 106 1384 453 1380 666 666 501
Gradient Steps	4	4	4	4	4	9	4	2	4	4		4	4	2 4 4	4400	44000	440000	4400000	44000000	4400000000	4400000004	4 4 0 0 0 0 0 0 4 0	4 4 0 0 0 0 0 0 4 0 0	4400000004004000
Batch Dize	×	22	6	12	x	12	12	128	×	28		42	42 12	42 12 8	42 8 8 8 8	42 12 8 8 16	42 12 8 8 16 8 8	42 12 8 8 12 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8	42 12 8 8 8 8 8 16	42 12 8 8 8 8 8 8 64	42 12 8 8 8 8 8 64 64 8 8	42 12 8 8 8 8 64 64 8 8 8	42 12 8 8 8 8 64 64 8 8 8 16	42 12 8 8 8 8 8 8 64 64 16 16 16 16 16
	1	Introduction to Information Science	Management Information System	Environmental Chemistry	Biotechnology	Q1	Q2	Q3	Q4	1	Q5	Q6 Q6	Q5 Q6 Q7	Q5 Q6 Q8	Q5 Q6 Q9	Q5 Q6 Q9 Q10	Q5 Q6 Q8 Q9 Q10 Q11	Q5 Q6 Q8 Q9 Q10 Q12 Q12	Q5 Q6 Q7 Q9 Q9 Q10 Q11 Q12 Female	Q5 Q6 Q7 Q8 Q9 Q10 Q10 Q11 Female Male	Q5 Q6 Q7 Q1 Q8 Q9 Q10 Q11 Q12 Female Male Traditional Exam	Q5 Q6 Q7 Q1 Q8 Q9 Q10 Q11 Q11 Q12 Female Male Traditional Exam Online Exam	Q5 Q6 Q7 Q8 Q1 Q9 Q10 Q11 Q11 Q12 Female Male Traditional Exam Online Exam Narrative (Q1, Q4, Q5)	Q5 Q6 Q7 Q8 Q1 Q9 Q10 Q11 Q11 Q12 Female Male Traditional Exam Online Exam Narrative (Q1, Q4, Q5) Argumentative (Q2, Q8, Q11)
	Entire Dataset		C	Courses								Output	Questions	Questions	Questions	Questions	Questions	Questions	Questions	Questions Gender	Questions Gender	Questions Gender Exam Type	Questions Gender Exam Type	Questions Gender Exam Type Question Type

parameters.
- Det
5
Ξ
and
10
eti
Š
al
periment
Å
Ţ
0
Summary
ä
4
ч Ю
Ā
لم

Ν	The Question	Poten-	Rubric-based evaluations										
		tial Mark											
1	Explain in detail the differ- ence between both the terms data and information and support your answers with ex- amples of each type?	5	(1 degree) The student's ability to introduce data, their role and shapes	(1 degree) The student's ability to introduce information, its upbringing and its use.	(2 degrees) The student's ability to conclude the difference between data and infor- mation	(1 degree) The student's ability to enhance his explanation of data and information with realistic, related ex- amples							
2	Explain at length the role of increasing micro-disciplines and increasing topics in in- fluencing the information ex- plosion revolution?	5	(2 degrees) The student's ability to explain the reasons for the increasing specializa- tions and the subject.	(2 degrees) The student's ability to the role and influence of increasing specialization in the information revolu- tion.	(1 degree) The student's ability to link the reasons for the emergence of modern sci- ence with the explosion of information								
3	Through what you learned in the course, mention the com- prehensive definition of the term information science?	5	(2.5) The student's ability to define the faces and role of the Information Sci- ence Department.	(2.5) The student's ability to introduce the tasks of in- formation science special- ists since the establish- ment of information to the delivery to the ben- eficiary.									
4	Explain the distinctions in roles and responsibilities among administrative levels in detail and provide illustra- tive examples.	10	(3 degrees) The student's ability to identify different manage- ment levels	(4 degrees) The student's ability to explain the difference be- tween the tasks and du- ties of each administra- tive level	(3 degrees) The student's ability to learn about the hierar- chical sequence of the tasks and roles of differ- ent management levels								
5	Mention three of the main benefits of cloud computing from a business perspective with an explanation?	10	(4 degrees) The student's ability to mention the three bene- fits	(3 degrees) The student's ability to explain each benefit ex- tensively	(3 degrees) The student's ability to explain the benefits of cloud computing in busi- ness administration								
6	Through what you learned in the course, mention the com- prehensive definition of the term information science?	5	(3 degrees) The student's ability to provide a comprehensive definition of the term in- formation technology	(2 degrees) The student's ability to mention examples of in- formation technology op- erations.									
7	Talk about the layers of the atmosphere, mentioning the height and temperature of each layer.	5	(1 degree) The student's ability to mention the names of the five layers correctly	(2 degrees) The student's ability to explain each layer exten- sively	(1 degree) The student's ability to conclude the difference between the role of each layer (temperature and height)	(1 degree) The student's ability to arrange the layers accord- ing to their proximity to the ground							
8	What do you think about the importance of the ozone layer?	5	(2 degrees) The student's ability to mention the role of the ozone layer in protecting the land	(2 degrees) The student's ability to explain the classes that have a role in protecting the earth.	(1 degree) The student's ability to know the basic role of the ozone layer								
9	What is the scientific defini- tion of environmental chem- istry?	5	(2 degrees) The student's ability to perform the term scientif- ically	(2 degrees) The student's ability to determine the aspects of environmental chemistry	(1 degree) The student's ability to mention the importance of environmental chem- istry for human and life								
10	Define the term biotechnol- ogy?	5	(2.5 degrees) The student's ability to provide a comprehensive definition of the term biotechnology	(2.5 degrees) The student's ability to mention examples of biotechnology.									
11	Discuss eating genetically modified fruits healthy or not?	5	(2 degrees) The student's ability to explain the components of the genetically modi- fied fruits.	(2 degrees) The student's ability to explain the benefits and negatives of genetically modified fruits.	(1 degree) The student's ability to list the reasons that make genetically modi- fied fruits acceptable.								
12	rive applications of biotech- nology in the medical field with explanation?	5	(3 degrees) The student's ability to mention five vital technol- ogy applications in the field of medicine.	(2 degrees) The student's ability to mention a simple expla- nation of each type.									

Table A.3: Scoring Criteria for Determining the Final Score, Set by Course Directors

Appendix B. Figures



Figure B.1: Showing Class Distribution Across the Twelve Questions, with Scores Ranging from 0 to 5