Automating Alternative Generation in Decision-Making

Anonymous ACL submission

Abstract

In decision making, generating alternative solutions is crucial for solving a problem. However, cognitive biases can impede this process by constraining individual decision makers' creativity. To address this issue, we introduce a new task for automatically generating alternatives, inspired by the process of human "brainstorming". We define alternative options based on atomic action components and present a dataset of 106 annotated Reddit r/Advice posts containing unique alternative options extracted 011 from users' replies. We also introduce new metrics to assess the quality of generated components, including distinctiveness, creativity, upvote-weighted, crowd intersection, and final commit intersection scores. As a baseline, we evaluated the large language models (LLMs) LLaMa3:8b, LLaMa3.1:8b, and Gemma 2:9b on the alternative component generation task. On the one hand, models demonstrated high creativity (ability to generate options beyond 022 what Reddit users suggested) and performed well at proposing distinct alternatives. A subset of generated components was manually evaluated and found overall useful. This indicates that LLMs might be used to extend lists of al-026 ternative options, helping decision makers consider a problem from different perspectives. On the other hand, LLMs' outputs often failed to align with human suggestions, implying that they still tend to miss important components.

> The code, annotation guidelines, and a request form for the dataset can be found in the project's GitHub repository¹.

1 Introduction

041

Decision-making is the process of choosing any course of action that aims to solve a problem in the best possible way. Every aspect of human life involves decision making to some degree, from selecting what to wear based on the weather to deliberating how to resolve a large-scale conflict. Yet,

¹HIDDEN FOR REVIEW

often human creativity is limited and constrained by biases when it comes to imagining different alternative actions leading to the best possible outcome. This paper introduces component-based alternative generation with the aim of providing first steps towards aiding the process of human alternative generation using machine learning.

Although different theories disagree on the structure of the decision making process (Morelli et al., 2022), most theoretical frameworks consider decision making to be the process of selecting a preferred option from a set of alternative options, frequently without specifying from where this set of alternatives originates. It has been highlighted in multiple studies that this process is a very important yet challenging step (Hämäläinen et al., 2024; Fisher et al., 1983).

As it relies on memory retrieval of relevant information (Johnson et al., 1991), the process of identifying possible actions poses human challenges. For example, information may not be organized in a coherent structure useful for the problem (Jungermann et al., 1983), interference from previous knowledge can make it more difficult to restructure information to see problems from different perspectives (Heuer, 1999), and the use of heuristics based on prototypical problems may introduce further human biases such as a tendency to learn towards highly representative options, and struggling to retrieve or creatively generate high-quality solutions (Gigerenzer and Gaissmaier, 2011; Gettys et al., 1987).

There are multiple theoretical frameworks developed to help decision makers overcome these challenges (Keeney, 1992; Pitz et al., 1980). "Brainstorming" is one of them (Al-Samarraie and Hurmuzan, 2018; Hicks, 1991). This method relies on aggregated judgment from multiple people, trying to utilize their cognitive efforts and mitigate their individual biases by providing different opinions and perspectives on the problem at hand. Overcom042

084

- 101 102
- 103
- 104
- 106

108

109

110 111

112 113

114 115

116

117 118

119

120

121

122

124

125

126

127

128

129

130

131

ing humans' limited cognitive flexibility, "brainstorming" is also particularly well-suited for computational automation (John, 2016).

Developing an automatic algorithm can help overcome human biases by generating wider set of possible alternatives. Such algorithms could be used in any process or task requiring decision making, such as operational planning, conflict resolution, and so on.

In this paper, we introduce a novel task in which the algorithm needs to generate lists of possible atomic options for solving human decision making problems. As a baseline, we evaluate different large language models (LLMs) in few and zero shot settings. To evaluate the performance of the models, a new dataset based on Reddit comments was manually labeled, approximating the "brainstorm" technique by incorporating advice replies from multiple users per question. The following contributions are made:

- We propose a new definition of alternative options based on atomic units of action (components), and introduce Component Generation (CG) and Component Competition (CC) tasks. In this paper, we focus on the CG task.
- We present a new dataset for the CG task based on the Reddit Advice subreddit² (r/Advice). Specifically, we filtered posts requesting advice based on predefined criteria. Then we extracted, labeled, and summarized proposed potential solutions from the comments for the filtered posts, and marked if the comment author considered them to be competing (mutually exclusive). Additionally, we identified comments to which the post author responded, extracting both atomic actions and whether the author did or committed to doing the proposed action, providing an indicator of which alternatives were perceived to be particularly helpful.

• We introduce novel metrics for evaluating alternative generation: distinctiveness, creativity, upvote-weighted intersection, crowd intersection, and final commit intersection scores. These metrics leverage an ensemble of LLM-based matching algorithm that checks if atomic actions imply the same. The matching algorithm was manually validated.

• We evaluate different LLMs for alternative generation using both zero-shot and few-shot approaches. For few-shot approaches, we used 5 and 10 examples, averaging results across three trials with different sets of examples. We conducted experiments with LLaMa3:8b (Dubey et al., 2024), LLaMa3.1:8b (Dubey et al., 2024), and Gemma 2:9b (Team et al., 2024).

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

2 **Related Work**

No datasets or evaluation metrics were available for the task at the time of writing. Decision-making has been discussed in academic literature from various disciplines, yet there is a notable lack of materials when it comes the generation of alternatives from a computational perspective.

Early studies (Arbel and Tong, 1982; Ozernoy, 1985; Alexander, 1979) were focusing on the benefits and influence of alternative generation on the overall decision process. A comprehensive survey and overview of different techniques (Keller and Ho, 1988) discussed in detail various different methods that a person could utilize to achieve the best possible set of alternative actions in any given scenario.

"Mean-value" approaches (Keeney, 1992; Keeney et al., 1994; León, 1999) encourage decision makers to estimate the relative importance of alternatives ("values") as well as the means to achieve them ("means"). In this framework, the alternative choices are considered to be given. Similarly, in one of the most well-known models of decision making so far, (Simon, 1955), the author introduced a "design" concept prior to the "choice".

MGA (Modeling for Generating Alternatives), a theoretical framework for alternative generation, has been presented and discussed in various studies (Brill Jr et al., 1982; Chang et al., 1983; Simon, 1955; DeCarolis, 2011; DeCarolis et al., 2016). The proposed algorithm formalizes a decision making process. The method includes multi-objective optimization algorithms to explore the neighborhood of a possible solution in order to find the most optimal solution. It requires a distance function initialization that measures the differences between the solutions, as well as a strict definition of the importance of the model's objectives and constraints. Recently, Colorni and Tsoukiàs (2020) formulated a general framework for formalizing alternatives, stressing how under-researched this area has been.

²https://www.reddit.com/r/Advice/



Figure 1: Pipeline Overview of the Framework.

The first attempts at systems practically aiding the decision making process were introduced in the late 90s (Leal and Pearl, 1977; Pearl et al., 1982; Leal et al., 1978; Steeb and Johnston, 1981; Keller and Ho, 1988). Recent studies address further theoretical subtleties, carefully structuring various techniques and methods that can aid decision makers (Hämäläinen et al., 2024). It is furthermore worth noting that research incorporating the use of artificial intelligence to assist idea generation generally (Shaer et al., 2024) is also thematically adjacent to our research, though it lacks the distinction between different competing alternatives which is so crucial for describing decision making.

182

183

184

185

188

189

190

191

192

193

195

197

198

204

210

211

212

214

215

216

217

218

219

221

225

Lastly, the proposed task also shares some similarity with question answering tasks (QA), for which there are different datasets available (quo, 2017; Lovenia et al., 2024; Kwiatkowski et al., 2019; Rajpurkar et al., 2016; Reddy et al., 2019; Joshi et al., 2017; Yang et al., 2018; Talmor et al., 2019). However, QA tasks have a correct answer, whereas in the case of alternatives this is not necessarily clearly evident.

Definitions and Tasks 3

In this work, an alternative option (AO) refers to an action or a set of actions that could be taken in the context of a given problem. An action refers to the specific steps taken to implement the chosen alternative, aiming to resolve the problem and achieve the desired outcome. In the context of this work, a problem refers to some type of situation that requires a response and involves a choice between different AOs to achieve or solve it, based on a definition by Beachboard and Aytes (2013). In decision making, solutions are considered sideby-side and can therefore be termed "alternatives". Each AO consists of smaller units of action that we call components, which are atomic, i.e., cannot be broken down further. More formally, components are the smallest actions that can be taken in solving the problem. They may include conditional elements, a certain order of components, concretize specific actors who should participate or perform the component, etc. Components are characterized by the order or actions, semantic content, conditional parts, and the entities and participants included in the components. Components c_1 and c_2 are considered identical if they fit the following component matching rules (CMRs), i.e., if c_1 and c_2 preserve the order of actions and overall semantics, maintain the same conditional parts, and refer to the same entities, participants, people, etc. Consider the following example:

226

227

228

229

230

232

233

234

Problem:

235 An office worker A accidentally took a cookie 236 from a bowl in the office kitchen, assuming that it 237 was a shared bowl. However, it turns out that the contents of the bowl were the private lunch of 239 colleague B, who still has not noticed that one of 240 the cookies is missing. What do you recommend 241 worker A should do in this situation? 242 In this scenario, the components could be as 243 follows: 1. Do not tell B that you took the cookie. 2. Tell B that you took the cookie. 3. Buy B a whole new pack of cookies. 247 4. If you value B's friendship, tell the truth. 5. Take more cookies. 249 6. Tell B that you took the cookie, then see B's 250 reaction: if B is angry, buy B a lunch. 251 7. Buy *B* a lunch. 252 8. Tell B and HR that you took the cookie 253 In the example, there are multiple suggested 254 component actions that A can take. For example, 255 256

Components 2 and 4 are different, as there is an additional conditional part to Component 4 that is 257 not present in 2. Components 1 and 8 include the 258 same action, but different participants. Component 259 6 contains an order of actions. Some of the compo-260 nents can be actioned together (e.g., Components 261 1 and 5). If the components are mutually exclu-262 sive (e.g. Components 1 and 2), they are referred 263 to as competing components. Each AO is a set of 264 non-competing components, i.e., all of its actions can be performed without a conflict. *Competing alternatives* are sets of AOs whose components are competing with each other. This is similar to the concept of competing hypotheses (Heuer, 1999).

For instance, because of the competing Components 1 and 2, alternative options $AO_1 = \{1, 3, 7\}$ and $AO_2 = \{2, 3, 7\}$ are competing alternatives, despite having overlapping components. A subset of the AO is an AO itself.

Following this logic, the task of alternative generation can be structured as follows: (i) **Component Generation (CG)** and **Component Competition** (**CC**).

The CG task involves generating as many relevant components as possible to solve the problem, based on its title and description. This is framed as a text generation task.

The CC task is a binary classification problem: given two components, the algorithm must determine whether they are competing.

In this paper, we focus on the CG task and leave other proposed tasks for future work.

4 Dataset

265

266

271

273

274

275

279

282

287

289

290

291

293

294

296 297

298

302

303

310

311

312

314

4.1 Annotation

In this work, we introduce a novel dataset which was specifically created to fit the task. For this, we manually labeled the Reddit Corpus dataset from the ConvoKit (Chang et al., 2020). From this dataset, all of the r/Advice subreddit posts and comments were gathered for annotation. Two annotators were recruited for the task of labeling alternatives: a PostDoc researcher (Annotator *a*) and a PhD student (Annotator *b*) from the Centre for Argument Technology at the University of Dundee (UK).

The annotators were presented with the following task. In the Doccano (Nakayama et al., 2018) annotation system they were presented with a Reddit post (title and content) from the dataset and a list of comments: the 50 most upvoted comments, 25 random comments and pairs of comments (original post's author reply and the comment that this reply was addressed to). The annotators were asked to read the post and determine if in the post the problem was stated clearly and the author was asking for help in a search of possible actions or options to take. Posts that did not meet these criteria, asked for medical or legal advice, required very specific domain knowledge, included moral dilemmas, included images, or were too broad (e.g. open-ended 315 questions such as "What should I name my cat?") 316 were disregarded. The rules for the post exclusion 317 were created empirically from the initial labeling 318 by the authors of the paper. If the post fitted the 319 criteria, the annotator checked all comments and 320 determined if users proposed a solution in any of 321 the comments. For each comment in which a poten-322 tial solution could be found, the solution was split 323 into its respective components. Each component 324 was summarized to retain as much information as 325 possible while keeping it short (e.g. this involved 326 reconstructing pronouns and anaphora, removing 327 hate speech etc.). The annotators were instructed 328 not to consider the quality of the components, but 329 to remove clearly sarcastic, joke, and offensive 330 propositions. After summarizing, the annotators 331 compared each component with the list of previ-332 ously retrieved components for this post. If the 333 component had already been proposed, the exact 334 phrasing of the already existing component was 335 taken from the list of alternatives and assigned to the new comment. If the component had not been 337 proposed yet, it was assigned to the comment and 338 added to the list. Additionally, if the author of 339 a comment implied that multiple components of 340 their alternatives were competing, the annotators 341 noted the mutual exclusivity of these components 342 by marking each competing solution with an alter-343 native number for the comment. For example, if a 344 comment proposed that the author of the original 345 post should choose one of two options (e.g. Do 346 this ... or do that ...), the solution was marked with the numbers 1 and 2 with respect to the mentioned components. If there was only one proposed 349 solution, it was marked with -1. A new list of 350 alternatives was created for each new post. The 351 annotation guidelines can be found in the project's 352 GitHub repository. 353

On Reddit, authors can edit their original post with an update after it was posted. For the annotation, we removed updates from the posts, as they usually were not available to the commentators when they proposed their solutions. The text of the update part was annotated separately, as the author response was relevant for later analyses of which alternatives the author committed to.

354

355

356

357

358

360

361

362

363

364

365

4.2 Inter-Annotator Agreement

To measure inter-annotator agreement, we used two strategies: comparing the number of extracted components per post and the components' semantic intersection with the respect to the component matching rules (CMRs).

366

367

368

374

375

378

381

383

387

389

394

400

401

402

403

404

405

406

15 posts with 148 comments were included in the annotation data for both of the annotators (as an overlap). The total number of comments which contained at least one component was 79 for Annotator *a* and 120 for Annotator *b*. Across all posts, the total number of unique components was 49 for Annotator *a* and 51 for Annotator *b*. The Cohen's kappa score (Kohen, 1960) for the count of components extracted per comment was 0.614, indicating substantial agreement (Landis and Koch, 1977).

To answer the question if annotators extracted the same components, we used the following approach. For each of 15 posts, Annotator a was presented with two lists of components: unique components that extracted by Annotator b and unique components extracted by Annotator a. For each component from the a's list, the annotator marked which component (if any) from b's list it might correspond to based on the CMRs. To calculate an agreement, we used the following formula:

$$\frac{\sum_p (U_a^p + U_b^p) / (T_a^p + T_b^p)}{N},$$

where p is the post, U_a^p is the number of components from Annotator a's unique list of components that did not appear in Annotator b's list for the post p. U_b^p is the number of components from Annotator b's unique list of components that did not appear in Annotator a's list for the post p. T_a^p and T_b^p refer to the total amount of unique components extracted by annotators a and b respectively for the post p. N is the total number of posts. Taken together, this score indicates the average amount of components which are extracted by only one of both annotators, with respect to the total number of components. We divide by $T_a^p + T_b^p$ to ensure the fairness of the score.

The lower this score, the higher the annotators' agreement. Our obtained score was 0.36, indicating a reasonable agreement between the similarity of the extracted components.

4.3 Dataset Statistics

407After filtering out posts that did fit the criteria, the408total number of unique posts is 106. The total num-409ber of unique authors is 101, with 5 posts attached410to deleted accounts. The total number of consid-411ered comments is 3,828. Among them, the total412number of comments that the author of the post

	Total	Mean	Unique
Title	1,803	17.00	687
Post body	41,907	395.34	4,784
Post update	3,704	34.94	1,077
Comment	244,657	68.55	12,421
Component	37,284	8.05	2,949

Table 1: Number of words statistics. The *Post body* value was calculated after the removal of update.

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

replied to is 1,413 in 97 posts. The number of comments that did not propose any solution is 1,999. The average number of solutions per post is 14.03, with the maximum of 44 and minimum of 4. The average number of competing alternatives per post is 1.04 with the maximum of 9. The total number of posts where the author appeared to commit to take some of the actions is 75 (70.7% of all unique posts), with 197 unique components and 240 comments in total (sometimes, the author replied to multiple comments with the same commitment).

The total number of words can be found in the Table 1. To determine the number of words, we used the NLTK (Bird et al., 2009) framework.

5 Evaluation

In the context of the component generation task, we propose the following metrics: **distinctiveness**, **creativity**, **upvote-weighted intersection**, **crowd intersection**, and **final commit intersection** scores.

All the proposed metrics require an algorithm that determines whether the two components are identical. Recall that the components c_1 and c_2 are considered identical if they fit the **component matching rules (CMRs)**: c_1 and c_2 preserve the order of actions and overall semantics, have the same conditional parts, and refer to the same entities, participants, people, etc. The calculation is based on the *components matching algorithm*, which is detailed in the following subsection.

5.1 Components Matching Algorithm

To determine whether the pair of components are identical with respect to the defined CMRs, we developed an LLM-based ensemble method, utilizing an ensemble of LLaMa3:8b (Dubey et al., 2024) and Mistral:7b (Jiang et al., 2023) models. The component matching algorithm architecture is presented in Figure 2.

The component matching algorithm works by

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

509



Figure 2: Components Matching Algorithm.

providing LLaMa3 with a prompt containing CMRs and four examples covering all the rules. LLaMa3 predicts if any CMR is violated. If it fails to output "MATCH"/"NOT MATCH" answer, the input is passed to Mistral; if Mistral also fails, we assume the CMRs are violated.

We evaluated the component matching algorithm with a Mistral:7b generated dataset of components, which was manually reviewed and filtered. The evaluation and metrics of the component matching algorithm can be found in the Appendix A.1. The limitation of such an approach is that the dataset contains the Mistral's biases and could lack variety.

After generating, we manually evaluated a sample of pairs of generated component and gold components. The accuracy of the algorithm was 0.92 and weighted F1 was 0.93. The results are presented in the Appendix A.2.

5.2 Metrics

Notations Let P_{rs} be a list of all predicted components (repetition possible). with $|P_{rs}|$ referring to the number of elements in this list. P_s is a set of predicted *unique*, *none-repetitive* components (based on the matching algorithm and CMRs) from the model for the post p. $R_s = \{c_i | i = 1, ..., N\}$ is a set of extracted components from the dataset for the post p. $T = P_s \cap R_s$. E_s is a set of extracted components from the dataset for the post p posted by the original author of the post. U_c is the total upvotes for comments proposing component c.

The **distinctiveness** (**Ds**) score is calculated as a percentage of unique components from all the components that the model generated:

$$Ds = \frac{|P_s|}{|P_{rs}|}.$$

This score indicates the proportion of duplicates based on the matching algorithm. A higher distinc-

tiveness score indicates greater originality in the generated components.

The **creativity** (**Cr**) score is calculated as a percentage of the components that are considered to be not included in manually extracted components from the Reddit comments. Formally, it can be calculated as

$$Cr = \frac{|P_s - R_s|}{|R_s|},$$

where and $|P_s - R_s|$ corresponds to the magnitude of set difference. This score evaluates the model's ability to generate components beyond the "core" set of responses present in the dataset. A higher Cr indicates an ability to create novel components.

The **upvote-weighted intersection** (UWI) score is calculated as a weighted average of upvotes for components from the set T. The score is calculated as follows:

$$UWI = \sum_{c \in T} \frac{U_c}{\sum_{k \in R_s} U_k}$$

This score reflects the importance of the predicted components in relation to how well they align with the opinions of Reddit users (indicated by the number of upvotes for the comments). A higher UWI value indicates better alignment between the model's predictions and the Reddit users.

The **crowd intersection** (CI) score is calculated as follows:

$$CI = \frac{|T|}{|R_s|}$$

This score is a percentage of components that appeared in both the model generated component set and the target dataset. Low CI indicates that the model generated a small amount of components that match the target components. Therefore, it missed a lot of components that were brainstormed in the discussion. High score indicates a high intersection of the model's outputs and the target human produced components.

The **final commit intersection** (**FCI**) score is calculated via the formula:

$$FCI = \frac{|P_s \cap E_s|}{|E_s|}$$

This score reflects the intersection (from CMRs perspective) of the components that the author of the post explicitly mentioned in their reply as doing or planning to do.

482 483

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

513

514

515

516

517

519

520

521

523

524

525

526

527

529

532

533

534

535

537

539

541

542

543

544

545

546

547

551

553

555

559

6 Experimental Setup

The experiments were conducted with the following LLMs: LLaMa3:8b, LLaMa3.1:8b (Dubey et al., 2024), and Gemma 2:9b (Team et al., 2024). The models were instructed with an initial system prompt explaining the task and outlining that the output format should clearly separate components. The model was prompted to generate a special stop token once it had finished generating the components. The experiments were conducted with N = 0, 5, and 10 examples from the dataset presented to the model. Each example included title, post content (without the author's update), and expected list of components with the stop token at the end. Then, the test title and post content were presented to the model. The model generated the components, and if the stop token appeared in the generated text, it was considered as the final output. Otherwise, the model was presented with an additional request to generate more of components and complete the previous conversation history. This process was run until the stop token appeared in the text, or when the maximum allowed number of follow-ups (20) was reached.

For each of the few-shot experiments with $N \in \{5, 10\}$ examples the model was run independently 3 times, selecting random N examples from the dataset. The final metrics were averaged over the experiments per N. As a preprocessing step, all the generated components were run through the matching algorithm to remove the duplicate components. To evaluate a joint performance of the models, the generated results per N were aggregated. We set random seed equals 2, and set other generation parameters to defaults, including a temperature of 0.7.

The constructed prompts and code are available in the project GitHub repository. The overview of the pipeline is presented on the Figure 1.

7 Results

The results are presented in the Table 2. The correlation plots of the metrics are presented in Appendix C.

A random sample of the generated components was manually evaluated on their usefulness and relevance to the problem. The annotation results are presented in Figure 3 and in Appendix B. The models generated mostly useful components for tackling the input problem.

Based on the obtained results, all the model were

able to output distinct sets of components when presented with examples. The distinctiveness scores in each run was 1.0 for N = 5 and 10. However, when models were not presented with examples from the dataset, LLaMa3.1 and Gemma 2 obtained Ds of 0.943 and 0.981 respectively. These scores are still high, but not as good as when presented with few-shot examples. 560

561

562

563

564

565

566

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

589

590

591

592

593

594

595

597

The most creative model was LLaMa3:8b, as it was able to outperform other considered models with zero-shot (with Cr=1.557 and std of 1.048) and with N=10 (with Cr=1.574 and std of 0.275). When provided with 5 examples, LLaMa3.1 had the highest creativity score of 1.364. Not only was this result the best on average in the group of N=5, but it also was the most consistent one with std of 0.041.

On the other hand, when it comes to upvoteweighted intersection scores, there does not appear to be a clear winner. The UWI score can be interpreted as an approximation of utility of the predicted components based on the Reddit users' judgement. LLaMa3.1 achieved the best result with the zero-shot approach, with a score of 0.044 and the highest std of 0.14. In the 5-shot example experiments, Gemma 2 performed better with an UWI of 0.049 and the lowest std of 0.006. Finally, in the 10-shot group, LLaMa3 showed a score of 0.055 with the lowest std of 0.002.

For crowd intersection score (CI), LLaMa3.1 outperformed other models in zero-shot and 5-shot settings with the scores of 0.038 and 0.041 respectively. In the 10-shot settings, LLaMa3 obtained the highest score of 0.043. These scores are quite low, indicating a small intersection with the components generated by the Reddit users' "brainstorm". Therefore, LLMs seemed to have missed a substantial proportion of possible components.



Figure 3: Distribution of useful, somewhat useful, and not useful components per model and per experiment.

Model	N	Ds (std)	Cr (std)	UWI (std)	CI (std)	FCI (std)
LLaMa3:8b	0	1.0 (0.0)	1.557 (1.048)	0.033 (0.075)	0.032 (0.057)	0.027 (0.126)
LLaMa3.1:8b	0	0.943 (0.232)	1.429 (2.513)	0.044 (0.14)	0.038 (0.09)	0.061 (0.194)
Gemma 2:9b	0	0.981 (0.136)	1.299 (1.029)	0.023 (0.078)	0.016 (0.04)	0.021 (0.121)
LLaMa3:8b	5	1.0 (0.0)	1.131 (0.102)	0.044 (0.009)	0.035 (0.002)	0.051 (0.009)
LLaMa3.1:8b	5	1.0 (0.0)	1.364 (0.041)	0.046 (0.012)	0.041 (0.002)	0.065 (0.011)
Gemma 2:9b	5	1.0 (0.0)	1.134 (0.113)	0.049 (0.006)	0.038 (0.005)	0.048 (0.005)
LLaMa3:8b	10	1.0 (0.0)	1.574 (0.275)	0.055 (0.002)	0.043 (0.004)	0.053 (0.008)
LLaMa3.1:8b	10	1.0 (0.0)	1.295 (0.039)	0.053 (0.006)	0.037 (0.002)	0.06 (0.021)
Gemma 2:9b	10	1.0 (0.0)	1.123 (0.043)	0.042 (0.01)	0.033 (0.004)	0.039 (0.013)

Table 2: Average results per experiment for different LLMs on component generation task. N refers to a number of examples that was shown to the model. Ds stands for distinctiveness score, Cr is creativity score, UWI is upvote-weighted intersection score, CI is crowd intersection score, FCI is final commit intersection score. For N=0 only one experiment was conducted. In the brackets, *the standard deviation* is presented among the different runs. FCI was calculated only in the samples, where author provided indication of commitment to do a particular action.

For the final commit intersection score (FCI), LLaMa3.1 outperformed other models in all experiments with the scores of 0.061 (zero-shot), 0.065 (5-shot), and 0.06 (10-shot). This score indicates an intersection with the "best" components - the ones that had been selected by the original post author. However, in a lot of cases, this could also primarily represent a personal preference. Often, more context is required to determine what might be considered the best option for any particular problem.

In our experiments, we expected LLaMa3.1 to outperform LLaMa3 across the different experiments. However, LLaMa3 demonstrated the better performance in the N=10 settings. Similar behavior has been observed before. Based on the released results for these models, there are instances when LLaMa3:8b showed better results than LLaMa3.1:8b (for example, on GPQA (Rein et al., 2024) dataset, LLaMa3.1 obtained a score of 32.8 and LLaMa3 obtained 34.2 (Dubey et al., 2024)).

8 Conclusion

Our experiments showed that LLMs are capable of outputting distinct components for decision making. However, they still appear to be a far way from matching human judgement, even when presented with different examples of the expected alternative components. In our experiments, LLMs performed better when provided with more examples, as might be expected. In almost all the settings and experiments, the best performing models were LLaMa3 and LLaMa3.1. These models demonstrated the highest creativity, intersection with human judgement, and with which actions authors finally did or committed to doing. Nevertheless, intersection scores were overall still low, indicating room for improvement. 631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

From a practical perspective, the creativity aspect is important in decision making as it provides a bigger picture for decision making. Generating alternatives automatically might allow decision makers to go beyond cognitive limitations. We showed that the considered LLMs are able to produce high creativity score, outputting possible components that were not considered in Reddit comments. Therefore, LLMs could be helpful in creating and extending lists of options which might serve as a starting points for decision makers to consider a problem from different perspectives.

9 Limitations

It is challenging to evaluate the generated components and their match to the actual target components. While we chose to utilize LLMs and manually evaluated a sample from our experiments, further investigations might be required in order to create a more reliable metric. Similarly, newly generated solutions cannot be fully reliably evaluated from the utility perspective, though we gauged the usefulness of responses by manually evaluating a sample of generated components.

In this work, we did not evaluate hallucination aspects of the models: LLMs are known to sometimes generate output unrelated to the topic. This is an important task the field might seek to address. Moreover, LLMs' inferences are consuming a lot

630

664of resources and time. Finally, the dataset we have665can be extended further with more samples that666include more diverse domains. However, consider-667ing the importance of competing alternatives in the668decision making process, we believe that automatic669alternative generation is a significant first steps to-670wards potential future computer-assisted decision671making tools.

672 References

673

677

679

694

697

701

703

704

705

708

709 710

711

712

713

714

- 2017. First Quora Dataset Release: Question Pairs — quoradata.quora.com. https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs. Online.
- Hosam Al-Samarraie and Shuhaila Hurmuzan. 2018. A review of brainstorming techniques in higher education. *Thinking Skills and creativity*, 27:78–91.
- Ernest R Alexander. 1979. The design of alternatives in organizational contexts: A pilot study. *Administrative Science Quarterly*, pages 382–404.
- Ami Arbel and Richard M. Tong. 1982. On the generation of alternatives in decision analysis problems. *The Journal of the Operational Research Society*, 33(4):377–387.
- John Beachboard and Kregg Aytes. 2013. An introduction to business problem-solving and decisionmaking. In *Proceedings of the Informing Science and Information Technology Education Conference*, pages 15–27. Informing Science Institute.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. Natural language processing with Python: analyzing text with the natural language toolkit. " O'Reilly Media, Inc.".
- E Downey Brill Jr, Shoou-Yuh Chang, and Lewis D Hopkins. 1982. Modeling to generate alternatives: The hsj approach and an illustration using a problem in land use planning. *Management Science*, 28(3):221– 235.
- Jonathan P. Chang, Caleb Chiam, Liye Fu, Andrew Wang, Justine Zhang, and Cristian Danescu-Niculescu-Mizil. 2020. ConvoKit: A toolkit for the analysis of conversations. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 57–60, 1st virtual meeting. Association for Computational Linguistics.
- Shoou-Yuh Chang, E Downey Brill Jr, and Lewis D Hopkins. 1983. Modeling to generate alternatives: A fuzzy approach. *Fuzzy Sets and Systems*, 9(1-3):137– 151.
- Alberto Colorni and Alexis Tsoukiàs. 2020. Designing alternatives in decision problems. *Journal of Multi-Criteria Decision Analysis*, 27(3-4):150–158.

Joseph F DeCarolis. 2011. Using modeling to generate alternatives (mga) to expand our thinking on energy futures. *Energy Economics*, 33(2):145–152.

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

Joseph F DeCarolis, Samaneh Babaee, Binghui Li, and S Kanungo. 2016. Modelling to generate alternatives with an energy system optimization model. *Environmental Modelling & Software*, 79:300–310.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun

Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, An-798 drew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian

802

810

811

812

813

814

815

816

817

818

819

821

823

824

825

827

828

829

830

831

834

837

838 839

Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models. Preprint, arXiv:2407.21783.

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

859

860

861

862

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

- Stanley D Fisher, Charles F Gettys, Carol Manning, Tom Mehle, and Suzanne Baca. 1983. Consistency checking in hypothesis generation. Organizational Behavior and Human Performance, 31(2):233-254.
- Charles F Gettys, Rebecca M Pliske, Carol Manning, and Jeff T Casey. 1987. An evaluation of human act generation performance. Organizational Behavior and Human Decision Processes, 39(1):23-51.
- Gerd Gigerenzer and Wolfgang Gaissmaier. 2011. Heuristic decision making. Annual review of psychology, 62(1):451-482.
- Richards J Heuer. 1999. Psychology of intelligence analysis. Center for the Study of Intelligence.

1008

1009

954

Michael J Hicks. 1991. Brainstorming. In *Problem* Solving in Business and Management: Hard, soft and creative approaches, pages 87–107. Springer.

901

902

903

904

905

908

909

910

911

912

913

914

915

916

917

918

919

921

923

924

925

926

928

929

930

931

933

935

936

937

938

939

941

943

944

945

947

948

949

951

952

- Raimo P. Hämäläinen, Tuomas J. Lahtinen, and Kai Virtanen. 2024. Generating policy alternatives for decision making: A process model, behavioural issues, and an experiment. EURO Journal on Decision Processes, 12:100050.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. Preprint, arXiv:2310.06825.
 - Thomas John. 2016. Supporting business model idea generation through machine-generated ideas: A design theory. In *ICIS*.
 - George Johnson et al. 1991. In the palaces of memory: How we build the worlds inside our heads. (*No Title*).
 - Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Helmut Jungermann, Ingrid Von Ulardt, and Lutz Hausmann. 1983. The role of the goal for generating actions. In Advances in psychology, volume 14, pages 223–236. Elsevier.
- Ralph L. Keeney. 1992. Value-Focused Thinking: A Path to Creative Decisionmaking. Harvard University Press.
- Ralph L Keeney et al. 1994. Creativity in decision making with value-focused thinking. *Sloan Management Review*, 35:33–33.
- L Robin Keller and Joanna L Ho. 1988. Decision problem structuring: Generating options. *IEEE Transactions on Systems, Man, and Cybernetics*, 18(5):715– 728.
- Jacob Kohen. 1960. A coefficient of agreement for nominal scale. *Educ Psychol Meas*, 20:37–46.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.

- J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics*, 33(1):159–174.
- Antonio Leal, Steven Levin, Steven Johnson, Marcy Agmon, and Gershon Weltman. 1978. An interactive computer aiding system for group decision making. Technical report, Tech. Rep. PQTR-1046-78-2.
- Antonio Leal and Judea Pearl. 1977. An interactive program for conversational elicitation of decision structures. *IEEE Transactions on Systems, Man, and Cybernetics*, 7(5):368–376.
- Orfelio G León. 1999. Value-focused thinking versus alternative-focused thinking: Effects on generation of objectives. *Organizational Behavior and Human Decision Processes*, 80(3):213–227.
- Holy Lovenia, Rahmad Mahendra, Salsabil Maulana Akbar, Lester James V. Miranda, Jennifer Santoso, Elyanah Aco, Akhdan Fadhilah, Jonibek Mansurov, Joseph Marvin Imperial, Onno P. Kampman, Joel Ruben Antony Moniz, Muhammad Ravi Shulthan Habibi, Frederikus Hudi, Railey Montalan, Ryan Ignatius, Joanito Agili Lopo, William Nixon, Börje F. Karlsson, James Jaya, Ryandito Diandaru, Yuze Gao, Patrick Amadeus, Bin Wang, Jan Christian Blaise Cruz, Chenxi Whitehouse, Ivan Halim Parmonangan, Maria Khelli, Wenyu Zhang, Lucky Susanto, Reynard Adha Ryanda, Sonny Lazuardi Hermawan, Dan John Velasco, Muhammad Dehan Al Kautsar, Willy Fitra Hendria, Yasmin Moslem, Noah Flynn, Muhammad Farid Adilazuarda, Haochen Li, Johanes Lee, R. Damanhuri, Shuo Sun, Muhammad Reza Oorib, Amirbek Djanibekov, Wei Oi Leong, Ouvet V. Do, Niklas Muennighoff, Tanrada Pansuwan, Ilham Firdausi Putra, Yan Xu, Ngee Chia Tai, Ayu Purwarianti, Sebastian Ruder, William Tjhi, Peerat Limkonchotiwat, Alham Fikri Aji, Sedrick Keh, Genta Indra Winata, Ruochen Zhang, Fajri Koto, Zheng-Xin Yong, and Samuel Cahyawijaya. 2024. Seacrowd: A multilingual multimodal data hub and benchmark suite for southeast asian languages. arXiv preprint arXiv: 2406.10118.
- Matteo Morelli, Maria Casagrande, and Giuseppe Forte. 2022. Decision making: A theoretical review. *Integrative Psychological and Behavioral Science*, 56(3):609–629.
- Hiroki Nakayama, Takahiro Kubo, Junya Kamura, Yasufumi Taniguchi, and Xu Liang. 2018. doccano: Text annotation tool for human. Software available from https://github.com/doccano/doccano.
- Vladimir M Ozernoy. 1985. Generating alternatives in multiple criteria decision making problems: A survey. In *Decision Making with Multiple Objectives: Proceedings of the Sixth International Conference on Multiple-Criteria Decision Making, Held at the Case Western Reserve University, Cleveland, Ohio, USA, June 4–8, 1984*, pages 322–330. Springer.

Judea Pearl, Antonio Leal, and Joseph Saleh. 1982. Goddess: A goal-directed decision structuring system. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (3):250–262.

1010

1011

1012

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1030

1031

1032

1034

1035

1036

1037

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1050

1051

1052

1053

1054

1055

1056

1057 1058

1059

1060

1061

1062

1063

1064

1065

- Gordon F Pitz, Natalie J Sachs, and Joel Heerboth. 1980. Procedures for eliciting choices in the analysis of individual decisions. *Organizational Behavior and Human Performance*, 26(3):396–408.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. GPQA:
 A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.
- Orit Shaer, Angelora Cooper, Osnat Mokryn, Andrew L Kun, and Hagit Ben Shoshan. 2024. Ai-augmented brainwriting: Investigating the use of llms in group ideation. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–17.
- Herbert A. Simon. 1955. A behavioral model of rational choice. *The Quarterly Journal of Economics*, 69(1):99–118.
- Randall Steeb and Steven C Johnston. 1981. A computer-based interactive system for group decisionmaking. *IEEE Transactions on Systems, Man, and Cybernetics*, 11(8):544–552.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda

Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal 1067 Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. 1069 Choquette-Choo, Danila Sinopalnikov, David Wein-1070 berger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, 1074 Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna 1075 Klimczak-Plucińska, Harleen Batra, Harsh Dhand, 1076 Ivan Nardini, Jacinda Mein, Jack Zhou, James Svens-1077 son, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana 1078 Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh 1080 Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mo-1081 hamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lau-1084 ren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, 1086 Logan Kilpatrick, Lucas Dixon, Luciano Martins, 1087 Machel Reid, Manvinder Singh, Mark Iverson, Mar-1088 tin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moyni-1091 han, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nen-1093 shad Bardoliwalla, Nesh Devanathan, Neta Dumai, 1094 Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Peng-1096 chong Jin, Petko Georgiev, Phil Culliton, Pradeep 1097 Kuppala, Ramona Comanescu, Ramona Merhej, 1098 Reena Jana, Reza Ardeshir Rokni, Rishabh Agar-1099 wal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, 1100 Sarah Perrin, Sébastien M. R. Arnold, Sebastian 1101 Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, 1102 Sue Ronstrom, Susan Chan, Timothy Jordan, Ting 1103 Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, 1104 Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh 1105 Meshram, Vishal Dharmadhikari, Warren Barkley, 1106 Wei Wei, Wenming Ye, Woohyun Han, Woosuk 1107 Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan 1108 Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, 1109 Minh Giang, Ludovic Peran, Tris Warkentin, Eli 1110 Collins, Joelle Barral, Zoubin Ghahramani, Raia 1111 Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, 1112 Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hass-1113 abis, Koray Kavukcuoglu, Clement Farabet, Elena 1114 Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Ar-1115 mand Joulin, Kathleen Kenealy, Robert Dadashi, 1116 and Alek Andreev. 2024. Gemma 2: Improving 1117 open language models at a practical size. Preprint, 1118 arXiv:2408.00118. 1119

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, 1120 William Cohen, Ruslan Salakhutdinov, and Christo-1121 pher D. Manning. 2018. HotpotQA: A dataset for 1122 diverse, explainable multi-hop question answering. 1123 In Proceedings of the 2018 Conference on Empiri-1124 cal Methods in Natural Language Processing, pages 1125 2369-2380, Brussels, Belgium. Association for Com-1126 putational Linguistics. 1127

- 1128 1129
- 1130

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

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

A Components Matching Algorithm Evaluation and Architecture

A.1 Algorithm Architecture and Manual Evaluation

To determine whether the pair of components are identical with respect to the defined CMRs, we developed an LLM-based ensemble method, utilizing LLaMa3:8b (Dubey et al., 2024) and Mistral:7b (Jiang et al., 2023) models³.

The components matching algorithms was designed as follows. Firstly, LLaMa3 is provided with a prompt which provides the model with a set of CMRs. A few examples were provided as well, covering all of the rules. These examples were created manually outside of the dataset with the total number of examples of 4. All prompts and instructions are available in the project GitHub repository. The model then is instructed to predict "MATCH" (if the pair of components are the same) or "NOT MATCH" (if at least one of the CMRs does not hold). Where LLaMa3 fails to output the suitable value, the same inputs are provided to the Mistral model. If Mistral is not able to output the result, "NOT MATCH" label is assigned. During testing, models were able to output a value in the expected format for all the samples. We did not consider embeddings similarity-based techniques (e.g. cosine similarity-based using transformers or sentence similarity pre-trained models) as they are not able to match a specific set of rules, but only consider the overall semantics of the sentence input.

To evaluate the proposed algorithm, we made use of the new dataset by gathering all unique individual components from the labeled Reddit dataset. As per the annotation guidelines, all the components per post are considered to be unique, i.e., the same advice suggestions were always summarized in the same way. Hence, we could derive a set of goldstandard "NOT MATCH" samples from all combinations of any two unique components per post. The total number of negative ("NOT MATCH") samples was 10,841. The positive ("MATCH") examples were derived by paraphrasing all unique components from the dataset with Mistral using a zero-shot approach. The model was provided with

	Precision	Recall	F1
MATCH	0.956	0.886	0.919
NOT MATCH	0.987	0.995	0.991
Macro avg	0.971	0.940	0.955
Weighted avg	0.984	0.984	0.984

Table 3: Results of the ensemble matching metric.

the instruction to preserve required components, as was described in the previous paragraph. The paraphrased versions were manually reviewed afterwards to ensure that the paraphrase fit the requirements. As a result, 1,184 samples were accepted and 181 rejected. 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

Finally, the proposed approach was run on the combined dataset of "MATCH" and "NOT MATCH" pairs. The results are presented in Table 3.

Considering that all metrics exceed 90%, particularly the recall metric for "NOT MATCH" class, the ensemble approach is effective for determining whether components are the same. While the metrics are not perfect, they apppear reasonable and demonstrate strong performance.

A.2 Manual Evaluation of Generated and Manually Extracted Components

After the LLMs generated components, we took a random 300 of pairs of generated and manually extracted components (71 was marked as matching and 229 as not matching by the algorithm). The pairs were selected randomly across all experiments (see distribution on the Table 4). Annotator *a* manually evaluated those pairs. The results are presented in the Table 5. The overall accuracy is 0.92. Results indicate a high agreement, therefore the algorithm could be considered reliable.

B Usefulness Evaluation of Generated Components

Additionally, we manually evaluated the useful-1204 ness/relevance of the model generated components. 1205 The same annotators a and b were recruited to an-1206 notate the sample from the pool of generated com-1207 ponents. The components were selected evenly 1208 across the models, number of few shot examples 1209 shown to the model, and experiment runs. The an-1210 notators were shown an original Reddit post and a 1211 generated component. They had to determine, if 1212 the component is relevant to the post and if it is 1213 useful (assign a label U), somehow useful (assign 1214

³We experimented with other LLMs, but this combination provided the best overall result. Embeddings similarity-based techniques (e.g. cosine similarity-based using transformers or sentence similarity pre-trained models) were not able to match a specific set of rules, but only considered the overall semantics of the inputs.

	N. Examples	Run	N. Samples
Gemma2	0	1	10
Gemma2	10	1	12
Gemma2	10	2	11
Gemma2	10	3	23
Gemma2	5	1	21
Gemma2	5	2	15
Gemma2	5	3	12
LLaMa3	0	1	14
LLaMa3	10	1	20
LLaMa3	10	2	15
LLaMa3	10	3	17
LLaMa3	5	1	13
LLaMa3	5	2	12
LLaMa3	5	3	14
LLaMa3.1	0	1	9
LLaMa3.1	10	1	17
LLaMa3.1	10	2	7
LLaMa3.1	10	3	16
LLaMa3.1	5	1	9
LLaMa3.1	5	2	13
LLaMa3.1	5	3	20

Table 4: Number of evaluated pairs per model and per run.

	Prec	Rec	F1	Num
MATCH	0.77	0.96	0.86	71
NOT MATCH	0.99	0.91	0.95	229
Macro Avg	0.88	0.94	0.90	300
Weighted Avg	0.94	0.92	0.93	300

Table 5: Results of manual evaluation of generated and manually extracted components matching. *Prec* refers to the precision score. *Rec* refers to the recall score.

a label **SU**) or not useful and/or irrelevant (assign a label **NU**). The annotation results are presented in the Table 6.

C Metrics Correlation

1215

1216

1217

1218

1219 The correlation plots are presented on Figures 4,5,6, and 7. Out results show that in majority of cases 1220 models predicted relevant and useful components 1221 to the presented problem. As in some of our ex-1222 periments, distinctiveness score (Ds) did not have 1223 1224 variation, its values are missing. Our findings show that the upvote-weighted intersection score (UWI) 1225 has correlates with the crowd intersection score 1226 (CI). It is expected due to a design of these met-1227 rics: they both are based on the intersection of 1228



Figure 4: Correlation plot of aggregated results of the experiments over all runs.



Figure 5: Correlation plot of aggregated results of the experiments with N=0 over all runs.

the generated components and manually annotated components that are matched to them. Other pairs of metrics do not show high correlations.

	Num. U			Num SU			Num NU		
	Ant. a	Ant. b	Total	Ant. a	Ant. b	Total	Ant. a	Ant. b	Total
Gemma2, N=0	24	2	26	3	1	4	10	0	10
Gemma2, N=5	30	11	41	2	3	5	2	2	4
Gemma2, N=10	21	12	33	1	2	3	2	2	4
LLaMa3, N=0	28	2	30	1	2	3	5	2	7
LLaMa3, N=5	22	12	34	3	2	5	1	0	1
LLaMa3, N=10	20	11	31	1	8	9	0	2	2
LLaMa3.1, N=0	23	3	26	3	2	5	7	2	9
LLaMa3.1, N=5	24	11	35	1	2	3	1	1	2
LLaMa3.1, N=10	20	11	31	0	4	4	2	3	5

Table 6: Results of manual evaluation of usefulness and relevance of the generated components. U indicates *Useful and Relevant*, **SU** indicates *Somehow Useful and Relevant*, and **NU** indicates *Not Useful and/or Irrelevant*. N indicates a number of few-shot examples. *Ant. a* refers to the results by the annotator a, and *Ant. b* refers to the results by the annotator b.



Figure 6: Correlation plot of aggregated results of the experiments with N=5 over all runs.



Figure 7: Correlation plot of aggregated results of the experiments with N=10 over all runs.