

A Comparison Study of Human and Machine Generated Creativity

Liuqing Chen

a. College of Computer Science and Technology,
Zhejiang University, Hangzhou, 310058, China

b. Zhejiang – Singapore Innovation and AI Joint Research Lab,
Zhejiang University, Hangzhou, 310058, China
e-mail: chenlq@zju.edu.cn

Lingyun Sun

a. International Design Institute of ZJU,
Zhejiang University, Hangzhou, 310058, China

b. Zhejiang – Singapore Innovation and AI Joint Research Lab,
Zhejiang University, Hangzhou, 310058, China
e-mail: sunly@zju.edu.cn

Ji Han¹

INDEX, Business School,
University of Exeter, Exeter, EX4 4PU, UK
e-mail: j.han2@exeter.ac.uk

¹ Corresponding author: Ji Han, j.han2@exeter.ac.uk

33

34 **ABSTRACT**

35

36 *Creativity is a fundamental feature of human intelligence. However, achieving creativity is often*
37 *considered a challenging task, particularly in design. In recent years, using computational*
38 *machines to support people in creative activities in design, such as idea generation and*
39 *evaluation, has become a popular research topic. Although there exist many creativity support*
40 *tools, few of them could produce creative solutions in a direct manner, but produce stimuli*
41 *instead. DALL·E is currently the most advanced computational model that could generate*
42 *creative ideas in pictorial formats based on textual descriptions. This study conducts a Turing*
43 *test, a computational test and an expert test to evaluate DALL·E's capability in achieving*
44 *combinational creativity comparing with human designers. The results reveal that DALL·E could*
45 *achieve combinational creativity at a similar level to novice designers and indicate the*
46 *differences between computer and human creativity.*

47

48 **Keywords:** *Artificial Intelligence, Computer Aided Design, Human Computer*

49 *Interfaces/interactions*

50

51

52

53

54 **1. Introduction**

55

56 Creativity has attracted great research interest in psychology, cognitive science,
57 computer science, engineering, and design fields for many years, and has a profound
58 impact on society [1]. It is defined as ‘the process by which something so judged (to be
59 creative) is produced’ [2], which is an essential skill to be successful in the current
60 complex and interconnected world [3]. In the past decades, several methods and
61 approaches, also known as creativity tools, are developed to support the generation of
62 creative ideas. Brainstorming, six thinking hats [4], SCAMPER [5], morphological analysis
63 [6], and TRIZ [7] are the most often used ones. Most of these conventional tools were
64 not developed specifically for design. Design-focused tools, such as the WordTree
65 method [8], 77 design heuristics [9], and bio-inspired design [10, 11], are thereby
66 developed specifically for supporting creative design idea generation. However, many
67 designers still prefer not to use these non-computational tools due to lack of knowledge
68 and experience, difficulties in mastery, and seemingly cumbersome steps which could
69 cause additional work [12].

70 In recent years, a number of computational design support tools have been
71 explored to tackle these limitations. For example, Han et al. [13] came up with an
72 analogical reasoning tool for supporting idea generation by employing aspects of
73 ontology and producing a corresponding image mood board; Sarica et al. [14] developed
74 a technology semantic network based on patent data, which could support ideation by
75 knowledge discovery; Siddharth et al. [15] proposed an engineering knowledge graph,
76 containing < entity, relationship, entity > triples extracted from patent database, to

77 support inference and reasoning; Obieke et al. [16] came up with a computational
78 framework that explores new engineering design problems for creativity. Most of the
79 existing so-called computational creativity tools do not generate creativity in a direct
80 manner, but produce stimuli instead, such as texts and images, to prompt designers'
81 creative minds.

82 Combinational creativity involves unfamiliar combinations of familiar ideas,
83 which is the easiest approach for humans to achieve creativity [17]. Producing
84 combinational creativity is a natural feature of humans' associative memory system,
85 while it is challenging for computers, due to issues such as the need for a rich store of
86 knowledge, the ability to form various combinations, and the competence to evaluate
87 combination outputs [17-20]. However, the rapid advancements in the field of artificial
88 intelligence, such as deep learning based computer vision and natural language
89 processing (NLP), have provided new and better approaches to enable computers to
90 produce combinational creativity. To the best of the authors' knowledge, no studies to
91 date have compared the performance between humans and computers in producing
92 combinational creativity. This leads to a debatable question that whether computational
93 machines (computers) can outperform humans in achieving combinational creativity.

94 Evaluating combinational creativity is challenging, and there is no widely adopted
95 method for such evaluation. In the field of design creativity, a variety of creativity
96 assessment methods have been proposed, which generally require human raters to
97 judge the quality of generated creativity [21], such as the Consensual Assessment
98 Technique method [38], Creative Product Semantic Scale [22], Product Creativity

99 Measurement Instrument (PCMI) [23], Creative Solution Diagnosis Scale (CSDS) [24], and
100 using creativity metrics [25, 26]. In the field of artificial intelligence, the common
101 computational metrics for evaluating generative models involve Inception Score (IS) [27]
102 and Frechet Inception Distance (FID) [28], which are quantitative and calculated based
103 on probability distribution. In the interdisciplinary research between artificial
104 intelligence and human study, Turing test is a basic and widely adopted method [30-32],
105 as it can provide an overall impression of how a machine performs. With consideration
106 of the advantages of the evaluation methods in these three areas, this study applies a
107 combined research approach by conducting a CAT based expert test, a computational
108 test and a Turing test, and then synthesizes the results to elicit useful findings.

109 Therefore, the aim of this paper is to compare the combinational creative
110 performance of machines and human designers, and explore the differences between
111 human designers and computers in generating creativity. This is the first study that
112 compares the performance between novice designers and machines regarding
113 combinational creativity, which employs a combined research approach integrating a
114 Turing test, a computational test and an expert test. This study will shed light on the
115 research of computational creativity evaluation and artificial intelligence applications in
116 design. The following section provides the theoretical background of this study. The
117 methodology of the study is described in Section 3, and followed by the implementation
118 of the Turing test, computational test and expert test in Section 4. In Section 5 and 6,
119 the results of the tests are presented, analyzed and discussed. The paper is then
120 concluded in Section 7.

121 **2. Theoretical Background**
122

123 Combinational creativity is claimed to be one of the best approaches for fully
124 utilizing nowadays abundant data, including texts, images, concepts, sounds and so on
125 [29], to achieve creativity [30]. A number of studies have explored combinational
126 creativity in the context of design, particularly in idea generation. For instance, Nagai et
127 al. [31] proposed three types of concept-synthesizing processes, namely property
128 mapping, concept blending, and concept integration in thematic relation, for generating
129 new concepts based on three interpretation methods of combinational phrases
130 respectively. Han et al. [32] indicated that associating far-related ideas for forming
131 combinational ideas could lead to outcomes that are more creative in comparison with
132 linking closely-related ones. Han et al. [33] investigated how combinational creativity is
133 formed in design, focusing on conventional noun-noun combinations. It was revealed
134 that a noun-noun combinational idea is produced by associating a base idea and an
135 additive idea. The base idea refers to the basic idea of the combinational idea, while the
136 additive idea could be a problem-solving idea, a similar representational idea, or an
137 inspirational idea. For example, the famous Juicy Salif is an example of associating a
138 basic idea (a manual juicer) and an inspirational additive idea (a squid). This study has
139 thereby laid a theoretical foundation for our paper exploring human and machine
140 generated combinational creativity.

141 Although Han et al. [19] and Chen et al. [34], [35] have employed pictorial data
142 to form combinational images to facilitate users in combinational creativity, these
143 combinational images are produced independently from semantic contexts. For

144 instance, the Combinator [19] produces a compound phrase of ‘flower glass’ and a
145 corresponding combinational image of merging a ‘flower’ and ‘glass’. Without semantic
146 context, the combinational image produced could represent a ‘flower’ made out of
147 ‘glass’, a piece of ‘glass’ in the shape of a ‘flower’, or a piece of ‘glass’ with printed
148 ‘flowers’. This might cause potential distractions and affect users’ creative performance.

149 In recent years, several computational models are developed to transform texts
150 into images, such as LeicaGAN [36] and Semantic-Spatial Aware GAN [37]. These models
151 could exploit text information for producing semantically consistent realistic images.
152 Among them, DALL·E [38] is one of the most advanced ones, which employs GPT-3 [39]
153 trained on a set of text-image pairs data for producing images based on text
154 descriptions. As introduced by OpenAI [40], DALL·E has distinguishing capabilities, such
155 as creating anthropomorphized versions of animals and objects. Moreover, it seems to
156 have achieved a certain level of creativity. Specifically, the model could create pictorial
157 combinations of unrelated concepts in plausible ways, even producing fantastical
158 objects that do not exist in reality, according to textual descriptions. Thus, DALL·E is
159 considered one of the most powerful systems capable of generating combinational
160 creativity in pictorial formats within the constraints of texts. In this study, we perform a
161 thorough performance benchmark evaluation comparing DALL·E with novice designers
162 regarding combinational creativity, involving a Turing, a computational, and an expert
163 test.

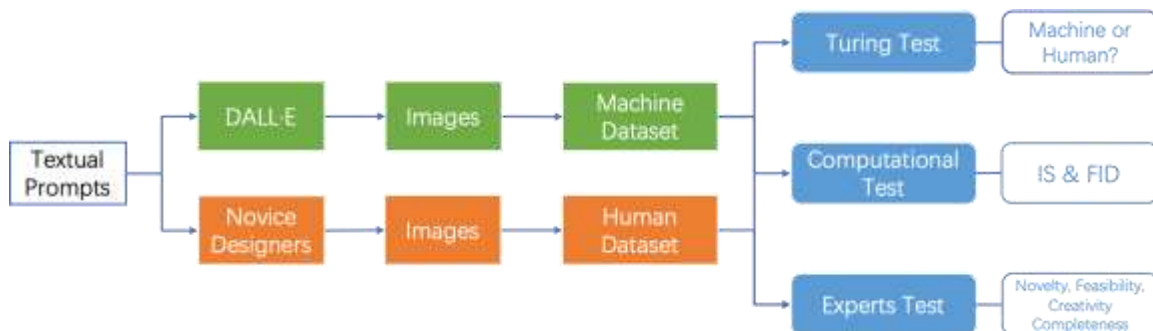
164

165

166 3. Methodology

167

168 To compare the performance between human novice designers and machines
 169 regarding combinational creativity, we first create two datasets for evaluation: the
 170 machine dataset and the human dataset. As shown in Figure 1, the input for both
 171 DALL·E and novice designers are the same textual prompts which contain combinational
 172 design ideas. The outputs are images matching the corresponding textual prompts. After
 173 selections, the same amount of data sets are saved as the machine dataset and the
 174 human dataset respectively. This is then followed by three tests: a Turing test, a
 175 computational test, and an expert test, in which the human and machine data are
 176 evaluated employing corresponding approaches.



177

178

Figure 1. The workflow of the proposed research approach

179

180 3.1. Data Source – Machine and Human Datasets

181

182 Only a partial code of the DALL·E model was released on Github, it is thereby
 183 impossible to run DALL·E to generate images due to missing training codes and data.
 184 Thus, the performance of DALL·E is evaluated based on the presented outcome from
 185 OpenAI's official blog, in which the published data is representative and of high quality.
 186 In the blog, sets of textual descriptions and the corresponding generated images by













187 DALL·E are presented. Three designers with over three years of experience were invited
188 to judge whether the textual description in each set is a combinational idea. Prior to the
189 judgement, the authors have well explained the definition of combinational creativity
190 and showed some practical cases to the designers. If a set was judged as combinational
191 creativity based, then five corresponding top-ranked images produced by DALL·E were
192 collected. In total, eight sets, with five images in each set, are collected as the machine
193 generated combinational creativity dataset. All the input texts and one corresponding
194 machine-produced image sample in each set are shown in Table 1.





195 Seven novice designers were employed to create a human dataset. They are
196 either postgraduates or employees in companies with less than three years of working
197 experience. They all hold a bachelor’s degree in design disciplines, and have at least two
198 years’ experience in product design and graphic design. Since the human dataset is
199 associated with combinational creativity, prior to the creation of data, each designer
200 was informed of the definition of combinational creativity and related design cases,
201 especially the meaning of ‘base’ and ‘additive’. Each designer was required to produce a
202 drawing for each of the textual descriptions as indicated in Table 1 by using familiar
203 computer-aided design software within one hour. The designers were required to use
204 white backgrounds and not to include any textual annotations to be in line with the
205 ones of the machine dataset. Besides, the quality of drawings should be as high as
206 possible, which is measured from three aspects:

- 207 1) Novelty: The drawing should be new, unusual, original and attractive.
- 208 2) Usefulness: The drawing should be feasible, reasonable and appreciable.

209 3) Creativity completeness: The drawing should match the corresponding
 210 textual description, and combined concepts could be visible to recognize.
 211 As a result, eight sets of data involving seven images each are produced. Three
 212 designers were then employed to select the top five images within each set. The eight
 213 sets of corresponding image samples produced by human designers are shown in Table
 214 1.

Table 1. An overview of the machine and human data

Group No.	Input	Machine Output	Human Output
1	a pentagonal green clock. a green clock in the shape of a pentagon		
2	a capybara made of voxels sitting in the field		
3	a stained-glass window with an image of a blue strawberry		
4	a snail made of harp. A snail with the texture of a harp		
5	an armchair in the shape of an avocado. an armchair imitating an avocado		
6	a giraffe imitating a turtle. a giraffe made of turtle		

7	a cube made of porcupine. a cube with the texture of a porcupine		
8	a professional high-quality emoji of a lovestruck cup of boba		

217

218

219 **3.2. Evaluation Methods**

220

221 3.2.1. Turing Test

222

223 A Turing test [41] is conducted in this study to explore whether DALL-E can
 224 achieve combinational creativity at the human level. In the test, participants were
 225 required to identify whether an image, within our mixed machine and human datasets,
 226 is produced by machine or human, providing the image's corresponding textual
 227 background. The test is consistent with the studies and arguments by Boden [42]; Pease
 228 and Colton [43]; Peter Berrar and Schuster [44]. The test is specific and blinded, and
 229 contains necessary contextual information. Though DALL-E is encouraged to produce
 230 realistic images in accordance with texts, it is not exclusively encouraged to exhibit
 231 creative behaviors. Therefore, the machine dataset, which can reflect DALL-E's capability
 232 of combinational creativity, was exclusively constructed to avoid possible trickery
 233 behaviors. For instance, instead of selecting the most realistic images generated by
 234 DALL-E to cheat human observers, we required that the images should first match their
 235 textual combinational ideas.

236

237 3.2.2. Computational Test

238

239 Given a deep learning based model for image generation, such as VAE [45] and
 240 GANs [46] based, the most common metrics for evaluating its capability are Inception
 241 Score (IS) and Frechet Inception Distance (FID). IS concerns the realism and diversity of
 242 generated images when evaluating a specific model. Specifically, IS calculates the KL
 243 divergence between the probability distribution of every generated image and the
 244 overall average of all generated images [27]. As shown in *Equation (1)*, given N classes,
 245 KL divergence is calculated between the conditional probability $p(y|x)$ in which a
 246 generated image x is classified into a particular class y , and the average probabilities for
 247 all the images in the class group $p(y)$ which is also called marginal distribution. High
 248 diversity of the generated images' categories and high certainty of the arbitrary image's
 249 category indicate high KL divergence, which means high IS and a better corresponding
 250 model, and there is no maximum value for IS.

$$251 \quad IS(G) = \exp\left(\frac{1}{N} \sum_{i=1}^N D_{KL}\left(p(y | \mathbf{x}^{(i)}) \parallel \hat{p}(y)\right)\right) \quad (1)$$

252 FID is proposed to perform better in terms of discriminability, robustness and
 253 computational efficiency and to address the limitations of IS [28]. It calculates the
 254 distance of two multidimensional normal distributions based on the mean (μ) and
 255 covariance (Σ) of the vectors extracted from both real (with the subscript r) and
 256 generated images (with the subscript g), as shown in the *Equation (2)*. Ideally, the FID
 257 can be zero if the generated data is identical to real data, while higher FID value
 258 corresponds to low quality of generated images. Considering the popularity of these two

259 metrics in generative models' evaluation, we calculate both values for our machine and
260 human datasets respectively, and then compare them.

$$261 \quad \text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}\left(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{\frac{1}{2}}\right) \quad (2)$$

262

263 3.2.3. Expert Test

264 The Turing test can estimate the overall appreciation of DALL·E's performance
265 compared with humans by subjective evaluation, while the computational test can
266 quantitatively and objectively compare machine and human performance but lack
267 detailed and interpretable criteria. Hence an expert test is necessary to deeply
268 investigate the difference between the two groups and provide interpretable results. In
269 this study, a Consensual Assessment Technique (CAT) based method [47] is adopted in
270 the expert test for creativity evaluation.

271 Novelty, quantity, quality and variety are the four metrics often used in design
272 research for evaluating creativity [25]. In the expert test, a modified version of the
273 metrics was adopted. Novelty, feasibility and creativity completeness were used to
274 measure a single image, and variety was used to measure a group of images generated
275 by either a human designer or machine. The combinational creativity images are
276 generated based on textual descriptions, thus novelty originates from the creation of
277 combining the 'base' elements with the 'additive' elements, such as the novelty of the
278 creation of combining 'armchair' with 'avocado' in an imagery format. On the other
279 hand, creativity completeness is an essential metric for evaluating the transformation
280 quality from textual description to imagery visualization, instead of focusing on

281 evaluating creation results (novelty). Since some of the combinational ideas are
282 imaginary rather than physical, such as ‘a giraffe imitating a turtle’, feasibility is chosen
283 as the metric instead of quality and utility. The meanings of novelty, feasibility and
284 creativity completeness are identical to the descriptions for ranking drawings in the
285 human dataset indicated in the preceding. Variety refers to the diversity of a set of
286 images, which measures the differences between images.

287

288 **4. Evaluation**

289

290 **4.1. Turing Test**

291

292 The Turing test is conducted by developing a website where all web pages are
293 completely customized to minimize distractions. Participants were asked to read the
294 instructions, agree with the test protocols, and provide demographic information before
295 starting the test. Eight groups of questions in total, corresponding to eight groups of
296 data in our datasets, are provided to the participants. Each group contains ten questions
297 that are randomly ordered for mixing the human and machine generated data, while
298 five questions are from the human dataset and another five are from the machine
299 dataset. This fact is not revealed to the participants to avoid introducing any potential
300 bias. This would not influence participants’ choices since they could feel free to make
301 decisions without restrictions. There is only one question on each webpage consisting of
302 a question serial number, a short textual description, an image which is either from the
303 human dataset or the machine dataset, and two buttons indicating ‘human’ and
304 ‘machine’ for participants to choose, as shown in Figure 2. The participants were

305 required to spend at least three seconds on each question before moving to the next
306 one.



307

Figure 2. A question webpage in the Turing test

308

309

310 After a successful pilot test, the test was distributed across multiple channels,
311 including university BBS, social media, and personal contacts. Each participant was
312 invited for an interview voluntarily when completing the test. Three questions were
313 asked in the interview:

314 1) How difficult do you think this test is?

315 2) What is your method for distinguishing human and machine?

316 3) What is your feedback about this test?

317 Answers of the interviews were collected and analysed in a qualitative way and
318 the results were reported in the '5. Results and Analysis' section.

319

320 4.2. Computational Test




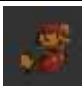






321





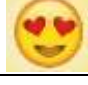
322 Two rounds of computational tests were conducted. In the first round, we
 323 implemented the algorithms of IS and FID by following Zhu et al. [48] and calculated IS
 324 and FID scores. FID calculation needs a reference distribution for comparison, so the
 325 mean and co-variance of COCO datasets [49] were used. However, it is found that some
 326 concepts in our datasets are not covered by COCO datasets, which might weaken the
 327 fairness of comparison. Therefore, we performed a second round of tests by comparing
 328 our data with a new reference dataset. As indicated in the preceding, a combinational
 329 idea consists of a base and an additive. Hence, we randomly collected 25 images for
 330 each base and additive in every group from the Internet, which results in 400 images in
 331 total. The 25 images for each base or additive were further equally divided into five
 332 reference groups in order to validate that no significant bias in image collection was
 333 introduced into the test. An overview of our reference data is shown in Table 2.

334

335

Table 2. An overview of the reference data

Group	Base	Sample-Base	Additive	Sample-Additive
1	Clock		Pentagonal	
2	Capybara		Voxels	
3	Glass		Strawberry	
4	Snail		Harp	
5	Armchair		Avocado	

6	Giraffe		Turtle	
7	Cube		Porcupine	
8	Cup		Emoji of lovestruck	

336

337 In the second round, we further calculated the IS of all five reference groups as a
338 reference to the IS of the human and machine dataset. The new FID scores were
339 calculated by comparing each reference group with the human and machine dataset
340 respectively. Since each generated image is based on a combinational idea and contains
341 concepts of base and additive, it is useful to investigate the FID by comparing the base
342 and additive data to the human and machine dataset. Therefore, the five reference
343 groups were further divided into base and additive sub-groups, and were used to
344 calculate base-FID and additive-FID.

345

346 4.3. Expert Test

347

348 The expert test was also conducted via a customized website. There are eight
349 groups of questions, and each contains twelve questions. In each group, the first ten
350 questions are single image based, of which a textual description and corresponding
351 image are provided in each question and participants are required to rate the image
352 using a 5-point Likert scale regarding three metrics: novelty, feasibility and creativity
353 completeness, as shown in Figure 3(a). The ten images are randomly selected from the
354 human or machine datasets. The last two questions in each group are five-image based,


355 in which a textual description and corresponding five images (merged in a vertical
356 sequence) are shown. Participants are informed that all five images were generated by
357 humans or machines exclusively, and they are required to rate the variety of the five
358 images using a 5-point Likert scale, as shown in Figure 3(b).

Rate The Image

Question 1-1

Text: a pentagonal green clock, a green clock in the shape of a pentagon.

Note: The image below is produced by human or machine based on above textual description.



1 2 3 4 5

Novelty 1 2 3 4 5

"1" means not novel, usual, non-original, non-attractive
"5" means very novel, unusual, original, attractive

Feasibility 1 2 3 4 5

"1" means not feasible, non-common-sense, not natural, strange
"5" means very feasible, common-sense, natural, normal

Creativity
Completeness 1 2 3 4 5

"1" means low creativity completeness, different from textual description, combinational concepts not visible
"5" means high creativity completeness, close to textual description, visible combinational concepts

Last Next


(a)

Rate The Images

Question 1-11

Text: a pentagonal green clock, a green clock in the shape of a pentagon.

Note: The five images below are produced by human or machine exclusively based on the above textual description.



1 2 3 4 5

Variety 1 2 3 4 5

"1" means this set of images has low diversity, they have similar
"5" means this set of images has high diversity, they look different

Last Next

(b)

359 Figure 3. Webpages of two question examples in the expert test

360 Before starting the test, participants were asked to read the instructions and test
361 protocols, and provide their demographic information. The explanation of four
362 evaluation metrics (novelty, feasibility, creativity completeness and variety) was
363 provided within the webpage, and further assistance was provided as well when experts
364 had questions. There was no time limit for each question, and more than 30 seconds of
365 rest time was provided in the test when experts completed half of the questions.

366

367 **5. Results and Analysis**

368

369 **5.1. Turing Test**

370

371 All ten images in each group shared the same textual description, and
372 participants were not informed how many images of the ten are from the human or
373 machine dataset, which means participants' judgement based on a single image is
374 independent. Among a total of 100 received submissions, there were 97 participants
375 who validly participated in this test by answering the 'human or machine' questions,
376 while three submissions were considered invalid as it was reported by the participants
377 that some machine generated images in the test were seen previously. The mean
378 accuracy of each question within each group was calculated, as well as the mean
379 accuracy of every group. The overall accuracy was obtained by averaging the accuracy of
380 eight groups, which is 55.9%, as shown in Table 3. Furthermore, group-8 achieved 42.4%
381 which is below 50% and the accuracy of group-6 is also very close to 50%.

382 Accuracy concerns whether a question is correctly answered or not, rather than
383 which answer is more often answered. Given a classification problem, human or

384 machine classes in our case, three metrics are widely applied when measuring the
385 performance of a classification machine learning model: precision, recall and F1 score.
386 The formulas of the three metrics are given in *Equation* (3), (4), (5) respectively, where
387 TP represents True Positive and similarly FN represents False Negative. In our
388 calculation, Positive means the answer is 'human' while Negative indicates 'machine'.
389 The results of precision, recall and F1 score of the two classes (human and machine) are
390 presented in Table 3. As shown in the table, the precision between human and machine
391 is very close (56.1% versus 55.6%), but the recall between human and machine are
392 noticeably different. The recall of the machine class is higher than the human class by
393 7.6%, which is due to high TN and high FN. Besides, the F1 score of the machine dataset
394 is higher than human by 3.3%.

395

$$Precision = \frac{TP}{TP + FP} \quad (3)$$

$$Recall = \frac{TP}{TP + FN} \quad (4)$$

$$F1\ Score = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (5)$$

396

397

398

Table 3. Results of the Turing test

	Mean	1	2	3	4	5	6	7	8
Accuracy	55.9%	60.8%	55.2%	61.9%	60.3%	54.2%	51.1%	61.1%	42.2%
Precision	55.6%	60.6%	54.4%	63.5%	60.8%	54.2%	51.1%	59.0%	42.2%
Machine Recall	57.9%	61.9%	64.1%	55.7%	57.9%	54.2%	53.4%	73.0%	42.7%
F1	56.7%	61.2%	58.8%	59.3%	59.3%	54.2%	52.2%	65.3%	42.5%
Precision	56.1%	61.1%	56.3%	60.6%	59.8%	54.2%	51.2%	64.6%	42.1%
Human Recall	53.8%	59.8%	46.2%	68.0%	62.7%	54.2%	48.9%	49.3%	41.6%
F1	54.9%	60.4%	50.7%	64.1%	61.2%	54.2%	50.0%	55.9%	41.9%

399

400 It is also useful to explore the variance of accuracy among questions and groups
401 when investigating machine and human classes respectively. Therefore, the statistics of
402 minimum and maximum accuracy in each group in terms of human and machine classes
403 are collected and presented in Table 4. As indicated in the table, both humans and
404 machines have very high variance throughout all groups, while the variance in the
405 human class is higher than the machine class. The highest accuracy in the human class
406 (93.8%) is higher than the machine class (86.6%) while the lowest accuracy in the human
407 class (16.5%) is lower than the machine class (26.8%), which corresponds to the value of
408 ($Max - Min$) between human and machine. The difference between the maximum and
409 minimum accuracy in the human class is higher than the machine class with 17% on
410 average.

411

412

413

Table 4. Variance of accuracy in different groups

		Min	Max	Max-Min	Difference
1	Human	37.1%	90.7%	53.6%	23.7%
	Machine	51.5%	81.4%	29.9%	
2	Human	19.6%	74.2%	54.6%	18.6%
	Machine	40.2%	76.3%	36.1%	
3	Human	43.3%	79.4%	36.1%	-2.1%
	Machine	38.1%	76.3%	38.1%	
4	Human	45.4%	93.8%	48.5%	26.8%
	Machine	46.4%	68.0%	21.6%	
5	Human	16.5%	89.7%	73.2%	45.4%
	Machine	43.3%	71.1%	27.8%	
6	Human	27.8%	77.3%	49.5%	7.2%
	Machine	32.0%	74.2%	42.3%	
7	Human	40.2%	63.9%	23.7%	-5.2%
	Machine	57.7%	86.6%	28.9%	
8	Human	25.8%	69.1%	43.3%	21.6%
	Machine	26.8%	48.5%	21.6%	
Overall	Human	16.5%	93.8%	77.3%	17.5%
	Machine	26.8%	86.6%	59.8%	
Mean	Human	32.0%	79.8%	47.8%	17.0%
	Machine	42.0%	72.8%	30.8%	

414

415 Twenty participants accepted the interview and answered questions after
416 completing the Turing test. Concerning the method of distinguishing human and
417 machine, the participants indicated that they believe the human-generated images have
418 'more clear details', 'a unified style (such as sketches)', and 'high resolutions', while the
419 machine-generated images are 'unreal', 'blurred' and have 'unhuman combination
420 logics' and 'cut and paste by Photoshop patterns'. In terms of the difficulty of the task,
421 the participants suggested that natural or physical subjects are easy to make 'human' or
422 'machine' selections, as well as images employing sketch styles. The interview results
423 are a supplement to the Turing test, and can potentially explain the Turing test results

424 and help understand the reasons underpinning the choices made by the participants.
425 This is in line with other similar studies. For example, Sarica et al. [50] interviewed
426 twenty-five participants to understand their choices of the best computational
427 representation of a specific design, and Zhu [51] interviewed ten engineers regarding
428 their views towards a set of computationally generated design concepts.

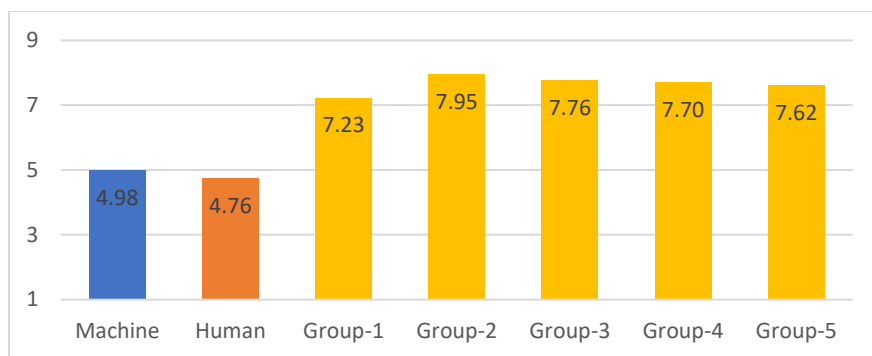
429

430 **5.2. Computational Test**

431

432 The computed results of the Inception Score (IS) are shown in Figure 4 where the
433 IS of five reference groups are presented together for reference purposes. The machine
434 group has a higher IS than the human group by 4.6%. The IS of the five reference groups
435 are much higher than the machine and human datasets with an average IS of 7.65 ($\sigma =$
436 0.27). The computed FID scores including reference groups are presented in Figure 5.
437 When comparing with COCO datasets, the FID of the machine dataset is higher than the
438 human dataset by 6.7%. All the FID scores in comparison with reference groups are
439 lower than COCO datasets, and all the FID scores of the machine group are higher than
440 the human group. The average FID of the machine group in comparison with the five
441 reference groups is 288 ($\sigma = 6.07$), which is higher than the average FID of the human
442 group ($\mu = 233, \sigma = 5.43$) by 23.8%.

443

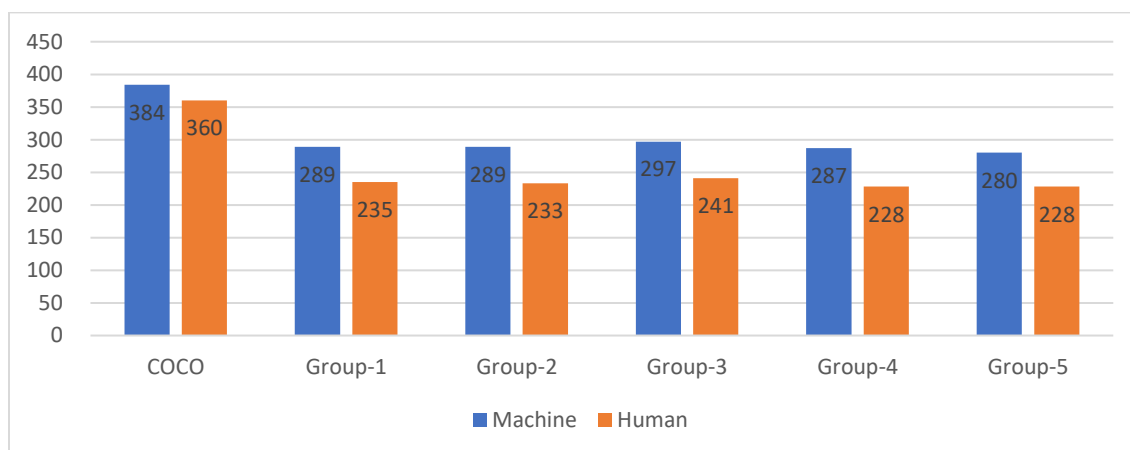


444

445

Figure 4. The IS values of different test groups

446



447

448

Figure 5. The FID scores of different test groups

449

450

451

452

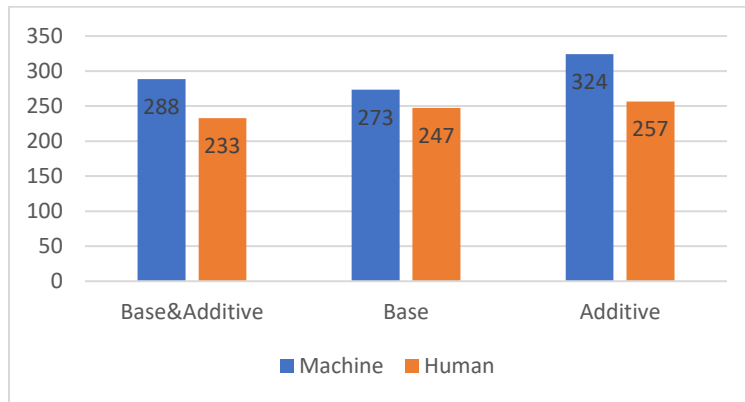
453

454

455

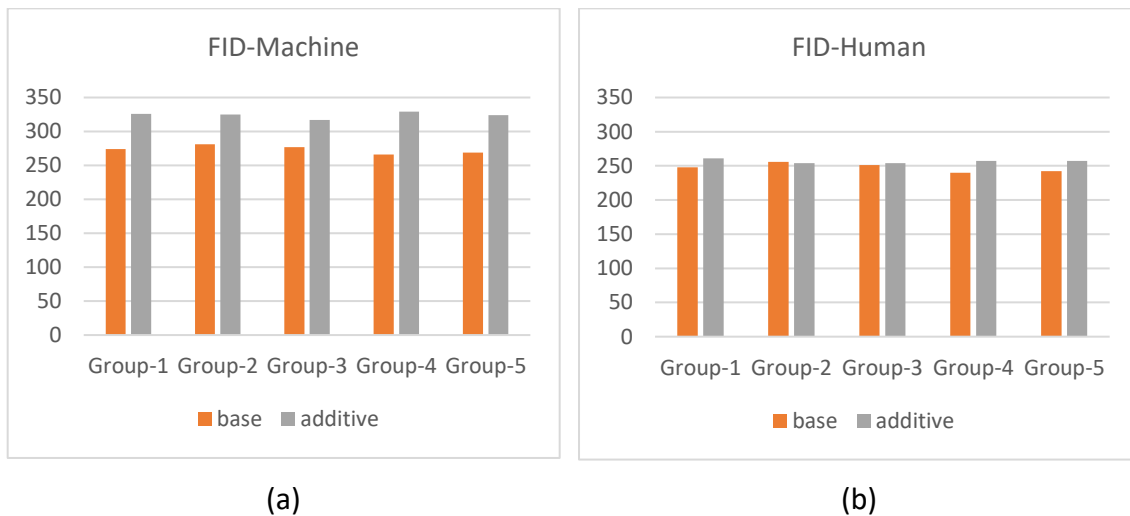
In addition to calculating FIDs with the mixed data of bases and additives in five reference groups, we further computed the FIDs comparing with base groups and additive groups respectively, as shown in Figure 6. The FID of the machine group ($\mu = 273, \sigma = 6.02$) is slightly higher than the human group ($\mu = 247, \sigma = 6.54$) by 10.5% in comparison with base groups, while the FID of the machine group ($\mu = 324, \sigma = 4.44$) is significantly higher than the human group ($\mu = 257, \sigma = 2.88$) by 26% in comparison

456 with additive groups. It is useful to investigate the influence of base and additive on the
 457 overall FID respectively. The FID scores in comparison with five base or additive groups
 458 (called base-FID and additive-FID respectively) are presented in Figure 7. As shown in
 459 the Figure 7 (a) (machine dataset), the additive-FIDs are higher than the base-FIDs on
 460 average by 18.7%, while Figure 7 (b) (human dataset) shows that the additive-FIDs are
 461 slightly higher than the base-FIDs only by 4.0%.



462

463 Figure 6. The FID scores in comparison with divided reference groups



464 Figure 7. The base-FIDs and additive FIDs comparison within the machine (a) and human

465 (b) datasets

466 **5.3. Expert Test**
467

468 With consideration of CAT requirements and the burden of evaluation, 19
469 professional designers with more than three years of working experience participated in
470 the expert test. The four metrics proposed are calculated and presented in Table 5. In
471 terms of novelty, more than half of the groups scored lower than 3, and the maximum
472 value is lower than 3.5. The human dataset achieved higher novelty ($\mu = 2.90, \sigma =$
473 0.14) than the machine group ($\mu = 2.78, \sigma = 0.39$). There are three groups related to
474 the machine dataset that obtained higher novelty scores than the human dataset. As
475 shown in the table, the human dataset has a higher feasibility score ($\mu = 3.41, \sigma =$
476 0.36) than the machine dataset ($\mu = 3.23, \sigma = 0.46$). The same groups related to the
477 machine dataset surpass the human dataset regarding feasibility. Similarly, the human
478 dataset achieved higher creativity completeness ($\mu = 3.36, \sigma = 0.25$) than the machine
479 group ($\mu = 3.09, \sigma = 0.49$). Two groups related to the machine dataset obtained higher
480 creativity completeness scores than the human dataset. For variety, the human dataset
481 has a significantly higher score ($\mu = 3.52, \sigma = 0.50$) than the machine dataset ($\mu =$
482 $2.95, \sigma = 0.47$), but there are three groups related to the machine dataset that surpass
483 the human dataset. Both the human and machine datasets have higher variance than
484 other metrics.

485

486

487

Table 5. Results of expert test

Metrics	Data Origin	1	2	3	4	5	6	7	8	Mean	Variance
Novelty	Machine	2.38	2.58	2.97	3.23	3.43	2.48	2.43	2.76	2.78	0.39
	Human	2.89	2.83	3.15	3.04	2.86	2.74	2.94	2.75	2.90	0.14
Feasibility	Machine	3.80	2.88	3.03	3.40	3.68	2.49	2.93	3.58	3.23	0.46
	Human	3.88	3.48	3.92	3.09	3.46	3.01	3.04	3.35	3.41	0.36
Completeness	Machine	3.48	2.68	2.98	3.57	3.42	2.44	2.54	3.64	3.09	0.49
	Human	3.53	3.06	3.64	3.46	3.48	3.40	2.89	3.44	3.36	0.25
Variety	Machine	3.00	2.74	2.37	3.42	3.47	3.16	3.26	2.21	2.95	0.47
	Human	3.84	4.37	3.84	2.89	3.05	3.21	3.21	3.74	3.52	0.50

488

489

490 **6. Discussion**

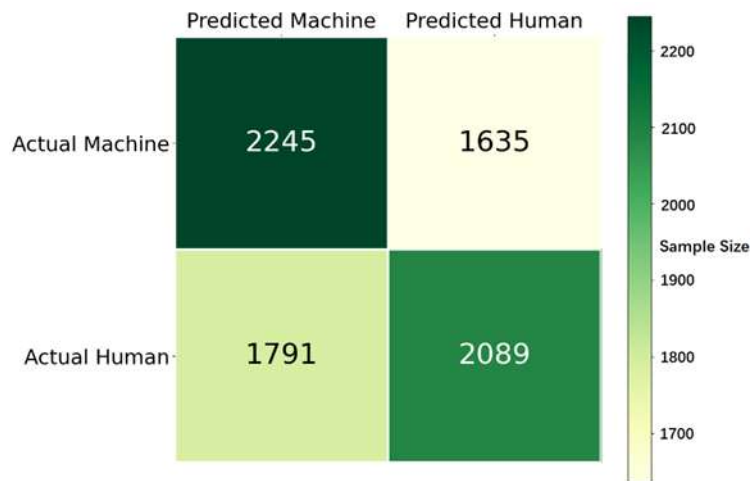
491

492 **6.1. Turing Test**

493

494 The average mathematical expectation of random answers to all the questions in
495 the Turing test is 50%, while the closer of overall accuracy to 50% indicates the more
496 undistinguishable between human and machine generated data. Though the overall
497 accuracy in the Turing test is above 50%, the gap is only 5.9%. The F1 scores of the
498 machine and human datasets are both close to 50%, while the machine's score is slightly
499 higher than the human's score due to high recall in the machine dataset. High variance
500 within every group in both datasets indicates that participants have low certainty to
501 make their judgements. Besides, as indicated in the confusion matrix in Figure 8, TN
502 (predicted machine and actual machine) and FN (predicted machine and actual human)
503 are relatively higher, which corresponds to higher recall and F1 score of the machine

504 dataset. This suggests that the results reveal that DALL·E can deceive participants to a
505 large extent, and the participants could hardly indicate which image is from the human
506 or machine dataset, while the participants subjectively tended to believe that the data
507 in the Turing test were more likely from machines rather than humans.
508



509

Figure 8. The confusion matrix of Turing test results

510

511

512 From our interview, it is shown that designers tend to use sketch and image
513 processing software (such as Photoshop) to create drawings rather than 3D modelling
514 and rendering, which makes their drawings more distinguishable from machine data. On
515 the other hand, the images generated by DALL·E tend to be blurred, unsmooth, and
516 unreal due to technical limitations, which makes them distinct from normal images.
517 Besides, the logic behind a combination idea in machine data is sometimes different
518 from human data. The 'cut and paste by Photoshop' pattern is considered a machine
519 pattern by some participants, since some designers tend to create a collage-style image

520 to express a combination idea while participants believe that machine is good at
521 creating collages.

522

523 **6.2. Computational Test**

524

525 We created five reference groups in the computational test, and all the results
526 related to the five groups have low variance, which indicates that there is only little bias
527 brought into the reference groups. Regarding the IS metric, the machine dataset
528 achieved a higher IS score than the human dataset, which means the machine
529 generated data have higher quality than the designers', but the gap is as little as 4.6%.
530 Five reference groups obtained much higher IS, since these reference images contain
531 rich information about bases and additives and they are natural rather than
532 combinational which is more favoured by the Inception model used for calculating IS.
533 On the other hand, the machine dataset obtained a higher FID score than the human
534 dataset when comparing with both COCO data and the five reference groups of data,
535 indicating the machine generated data have a lower quality than the human generated
536 data. All the FIDs in comparison with the five reference groups are significantly lower
537 than in comparison with COCO data, validating that the images in our reference groups
538 are closer to both the machine and human data than the images in COCO. The
539 difference of FIDs between the human and machine datasets in comparison with five
540 reference groups is bigger than the difference of FIDs in comparison with COCO data.
541 This may reflect the difference in combinational design between humans and machines.
542 Since it is required that drawings should be produced based on textual descriptions

543 containing combinational creativity, novice designers tend to keep essential information
544 from both base and additive in a combinational design while DALL·E is not trained to
545 obtain this capability. This indicates that these designers have a better lingual
546 understanding of combinational ideas and are able to transform them into designs than
547 machines.

548 It is found that the difference of FIDs in comparison with base is less than in
549 comparison with additive, as shown in Figure 6. This might suggest that designers are
550 better at maintaining additive information than DALL·E to some extent. Furthermore, as
551 shown in Figure 7, designers tend to balance base and additive information in a
552 combinational design while DALL·E tends to maintain more information from the base
553 rather than from the additive. However, there is no clear evidence that how much
554 information should be maintained from base and additive respectively in a
555 combinational design.

556

557 **6.3. Expert Test**

558

559 The human dataset obtained higher scores than the machine dataset by a small
560 percentage (6.17% on average) when comparing the results regarding novelty, feasibility
561 and creativity completeness, despite that the machine dataset has higher scores in
562 some groups. This indicates that the novice designers performed slightly better than
563 DALL·E in combinational designs in these three metrics. Besides, the designers
564 outperform DALL·E evidently regarding variety by an overall gap of 19.15%, even though
565 the machine dataset outperformed in three groups. This gap could be explained by two

566 reasons. One is that the human data are from seven novice designers while the machine
567 data is from DALL·E exclusively, which is unfair for DALL·E in this test. Another reason is
568 the difference in working mechanism between the DALL·E model and designers, in
569 which DALL·E takes text as input and generates various images based on random noise
570 while designers are skilled in producing various images using divergent thinking. It is
571 noticed that two to three groups in the machine data have higher scores regarding all
572 four metrics, indicating the capability of producing combinational creativity images
573 between novice designers and DALL·E is not significantly different.

574

575 **6.4. Overall Discussion**

576

577 There are no clear criteria to determine whether DALL·E passes the Turing test,
578 but it can be concluded that DALL·E's performance is close to novice designers according
579 to the results of our Turing test. In the computational test, DALL·E outperforms
580 designers in terms of IS but loses to designers regarding FID, and the difference in values
581 is both small, indicating that the performance between DALL·E and novice designers is
582 very close. It is noticed that the results of IS and FID are in conflict, which indicates that
583 the effectiveness of the two metrics for evaluating combinational creativity needs to be
584 further investigated. A larger difference in FIDs in comparison with our reference data
585 implies that human designers are better at synthesizing features from base and additive
586 for a combinational design. According to the results of the expert test, designers
587 outperform DALL·E from the perspective of combinational creativity. There is slight
588 advance for designers regarding novelty, feasibility and creativity completeness, but

589 evident advance regarding variety. By summarizing the conclusions from the three tests
590 in this study, DALL·E's performance is no better than novice designers but the gap is
591 small.

592 There are two key directions for future research. There is little research on
593 evaluating computational creativity. In this study, we applied three common methods
594 from different areas to evaluate the performance of DALL·E and compare it with novice
595 designers, which are labour intensive and lack scalability. How to effectively and
596 systematically evaluate computational algorithms in generating creative ideas or stimuli
597 needs further investigation and research. Another direction is the application of DALL·E
598 or other similar techniques in design, particularly in conceptual design. Design is a
599 process of transforming requirements and ideas into realisation, while DALL·E has the
600 capability of transforming an idea described in texts into a conceptual design solution
601 visualized in images. This would potentially provide a mental leap for designers,
602 particularly novices, facilitating creative idea generation.

603 There are a few limitations in this study. First, eight sets of data related to
604 combinational creativity, containing forty machine generated images and forty human
605 generated ones, were used in the study for evaluation. The limited amount of data was
606 a result of the restricted access to DALL·E's source code and data, as well as the high
607 cost of human resources. Although the amount of data is sufficient for the purpose of
608 the study, more data will be included in future studies by recruiting more human
609 designers and accessing more DALL·E data to yield further useful insights. This would
610 require the involvement of more human designers and accessing more DALL·E's data.

611 Second, one hour was provided to the designers to complete one combinational
612 creativity design task to construct the human dataset, but it is still far less to produce a
613 high-quality image. More time will be provided to the participants in future research to
614 improve the quality of the images generated. Third, DALL·E is a deep learning model
615 mainly aiming at transforming texts into images rather than generating combinational
616 creativity, which is less fair to compare with human designers. In future research, more
617 advanced artificial intelligence models, such as ChatGPT and GPT-4, will be included in
618 the comparison.

619

620 **7. Conclusion**

621

622 This paper is the first research that has explored the comparison of
623 combinational creativity capability between human beings and computers. It starts with
624 the preparation of two datasets, the machine dataset is created by collecting data from
625 a computational system, DALL·E, and the human dataset is created by inviting novice
626 designers to produce images based on textual combinational ideas. Three tests,
627 including a Turing test, a computational test and an expert test, are designed and
628 implemented on the two datasets. The results of the three tests reveal that DALL·E's
629 performance is very close to novice designers, while human designers are better at
630 synthesizing features from the base and the additive for a combinational design. The
631 results provide some useful insights for supporting the development of next-generation
632 computational systems to aid creative idea generation. The study represents a
633 contribution to the body of knowledge in research on computational methods for

634 design. It leads towards new research directions in evaluating computational creativity
635 and applying advanced computational techniques, particularly in conceptual design.

636

637 **FUNDING**

638 This is paper is funded by the National Natural Science Foundation of China [62207023]
639 and the Ng Teng Fong Charitable Foundation in the form of ZJU-SUTD (Zhejiang
640 University – Singapore University of Technology and Design) IDEA Grant.

641

642

643 **REFERENCES**

644

645 [1] Childs, P., J. Han, L. Chen, P. Jiang, P. Wang, D. Park, Y. Yin, E. Dieckmann, and I.
646 Vilanova, The Creativity Diamond - A Framework to Aid Creativity. Journal of
647 Intelligence, 2022. **10**(4): p. 73

648

649 [2] Amabile, T.M., The Social Psychology of Creativity. 1983, New York: Springer.

650

651 [3] Shute, V.J. and S. Rahimi, Stealth assessment of creativity in a physics video game.
652 Computers in Human Behavior, 2021. **116**: p. 106647.DOI:
653 <https://doi.org/10.1016/j.chb.2020.106647>

654

655 [4] De Bono, E., Six Thinking Hats. 1985: Little, Brown.

656

657 [5] Eberle, B., Scamper: Games for Imagination Development. 1996: Prufrock Press.

658

659 [6] Zwicky, F., Discovery, Invention, Research Through the Morphological Approach.
660 1969: Macmillan.

661

662 [7] Altshuller, G.S., Creativity as an exact science: The theory of the solution of inventive
663 problems. 1984, Amsterdam, Netherlands: Gordon and Breach Publishers.

664

665 [8] Linsey, J.S., A.B. Markman, and K.L. Wood, Design by Analogy: A Study of the
666 WordTree Method for Problem Re-Representation. Journal of Mechanical Design,
667 2012. **134**(4).DOI: 10.1115/1.4006145

668

669 [9] Yilmaz, S., S.R. Daly, C.M. Seifert, and R. Gonzalez, Evidence-based design
670 heuristics for idea generation. Design Studies, 2016. **46**(Supplement C): p. 95-
671 124.DOI: <https://doi.org/10.1016/j.destud.2016.05.001>

672

673 [10] Helms, M., S.S. Vattam, and A.K. Goel, Biologically inspired design: process and
674 products. Design Studies, 2009. **30**(5): p. 606-622.DOI:
675 <https://doi.org/10.1016/j.destud.2009.04.003>

676

677 [11] Chakrabarti, A. and L.H. Shu, Biologically inspired design. Artificial Intelligence
678 for Engineering Design, Analysis and Manufacturing, 2010. **24**(4): p. 453-
679 454.DOI: 10.1017/S0890060410000326

680

681 [12] Oman, S.K., I.Y. Tumer, K. Wood, and C. Seepersad, A comparison of creativity
682 and innovation metrics and sample validation through in-class design projects.
683 Research in Engineering Design, 2013. **24**(1): p. 65-92.DOI: 10.1007/s00163-
684 012-0138-9

685

686 [13] Han, J., F. Shi, L. Chen, and P.R.N. Childs, A computational tool for creative idea
687 generation based on analogical reasoning and ontology. Artificial Intelligence for

- 688 Engineering Design, Analysis and Manufacturing, 2018. **32**(4): p. 462-477.DOI:
689 10.1017/S0890060418000082
690
- 691 [14] Sarica, S., J. Luo, and K.L. Wood, TechNet: Technology semantic network based on
692 patent data. *Expert Systems with Applications*, 2020. **142**: p. 112995.DOI:
693 <https://doi.org/10.1016/j.eswa.2019.112995>
694
- 695 [15] Siddharth, L., L.T.M. Blessing, K.L. Wood, and J. Luo, Engineering Knowledge
696 Graph From Patent Database. *Journal of Computing and Information Science in*
697 *Engineering*, 2021. **22**(2).DOI: 10.1115/1.4052293
698
- 699 [16] Obieke, C.C., J. Milisavljevic-Syed, A. Silva, and J. Han, A Computational
700 Approach to Identifying Engineering Design Problems. *Journal of Mechanical*
701 *Design*, 2023. **145**(4).DOI: 10.1115/1.4056496
702
- 703 [17] Boden, M.A., *The creative mind: Myths and mechanisms*. 2 ed. 2004, London, UK:
704 Routledge.
705
- 706 [18] Simonton, D.K., Domain-general creativity: On Generating Original, Useful, and
707 Surprising Combinations, in *The Cambridge Handbook Of Creativity across*
708 *Domains*, Kaufman J.C., Glaveanu V.P., and B. J., Editors. 2017, The Cambridge
709 University Press.: Cambridge, UK. p. 18-40.
710
- 711 [19] Han, J., F. Shi, L. Chen, and P.R.N. Childs, The Combinator – a computer-based
712 tool for creative idea generation based on a simulation approach. *Design Science*,
713 2018. **4**: p. e11.DOI: 10.1017/dsj.2018.7
714
- 715 [20] Garvey, B., L. Chen, F. Shi, J. Han, and P.R. Childs, New directions in
716 computational, combinational and structural creativity. *Proceedings of the*
717 *Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering*
718 *Science*, 2019. **233**(2): p. 425-431.DOI: 10.1177/0954406218769919
719
- 720 [21] Beaty, R.E. and D.R. Johnson, Automating creativity assessment with SemDis: An
721 open platform for computing semantic distance. *Behavior Research Methods*,
722 2021. **53**(2): p. 757-780.DOI: 10.3758/s13428-020-01453-w
723
- 724 [22] Besemer, S.P. and K. O'quin, Analyzing Creative Products: Refinement and Test of
725 a Judging Instrument. *Journal of Creative Behavior*, 1986. **20**: p. 115-126
726
- 727 [23] Horn, D. and G. Salvendy, Product creativity: conceptual model, measurement and
728 characteristics. *Theoretical Issues in Ergonomics Science*, 2006. **7**(4): p. 395-
729 412.DOI: 10.1080/14639220500078195
730
- 731 [24] Cropley, D. and A. Cropley, Engineering Creativity: A Systems Concept of
732 Functional Creativity, in *Creativity across domains: Faces of the muse*. 2005,
733 Lawrence Erlbaum Associates Publishers: Mahwah, NJ, US. p. 169-185.

734

735 [25] Shah, J.J., S.M. Smith, and N. Vargas-Hernandez, Metrics for measuring ideation
736 effectiveness. *Design Studies*, 2003. **24**(2): p. 111-134.DOI:

737 [https://doi.org/10.1016/S0142-694X\(02\)00034-0](https://doi.org/10.1016/S0142-694X(02)00034-0)

738

739 [26] Han, J., H. Forbes, and D. Schaefer, An exploration of how creativity, functionality,
740 and aesthetics are related in design. *Research in Engineering Design*, 2021. **32**(3):
741 p. 289-307.DOI: 10.1007/s00163-021-00366-9

742

743 [27] Gulrajani, I., F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, Improved
744 training of wasserstein gans. In *Advances in neural information processing*
745 systems. 2017

746

747 [28] Heusel, M., H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, Gans
748 trained by a two time-scale update rule converge to a local nash equilibrium.
749 *Advances in neural information processing systems*, 2017. **30**

750

751 [29] Ward, T.B. and Y. Kolomyts, Cognition and creativity, in *The Cambridge handbook*
752 *of creativity*, J.C. Kaufman and R.J. Sternberg, Editors. 2010, The Cambridge
753 University Press: Cambridge, UK. p. 93-112.

754

755 [30] Yang, H. and L. Zhang, Promoting Creative Computing: origin, scope, research and
756 applications. *Digital Communications and Networks*, 2016. **2**(2): p. 84-91.DOI:

757 <https://doi.org/10.1016/j.dcan.2016.02.001>

758

759 [31] Nagai, Y., T. Taura, and F. Mukai, Concept blending and dissimilarity: Factors for
760 creative concept generation process. *Design Studies*, 2009. **30**(6): p. 648-
761 675.DOI: 10.1016/j.destud.2009.05.004

762

763 [32] Han, J., F. Shi, D. Park, L. Chen, and P. Childs. The conceptual distances between
764 ideas in combinational creativity. in *DS92: Proceedings of the DESIGN 2018*
765 *15th International Design Conference*. 2018.

766

767 [33] Han, J., D. Park, F. Shi, L. Chen, M. Hua, and P.R. Childs, Three driven approaches
768 to combinational creativity: Problem-, similarity- and inspiration-driven.
769 *Proceedings of the Institution of Mechanical Engineers, Part C: Journal of*
770 *Mechanical Engineering Science*, 2019. **233**(2): p. 373-384.DOI:

771 10.1177/0954406217750189

772

773 [34] Chen, L., P. Wang, H. Dong, F. Shi, J. Han, Y. Guo, P.R.N. Childs, J. Xiao, and C.
774 Wu, An artificial intelligence based data-driven approach for design ideation.

775 *Journal of Visual Communication and Image Representation*, 2019. **61**: p. 10-

776 22.DOI: <https://doi.org/10.1016/j.jvcir.2019.02.009>

777

- 778 [35] Chen, L., P. Wang, F. Shi, J. Han, and P. Childs. A computational approach for
779 combinational creativity in design. in DS 92: Proceedings of the DESIGN 2018
780 15th International Design Conference. 2018.
781
- 782 [36] Qiao, T., J. Zhang, D. Xu, and D. Tao, Learn, imagine and create: Text-to-image
783 generation from prior knowledge. Advances in Neural Information Processing
784 Systems, 2019. **32**: p. 887-897
785
- 786 [37] Hu, K., W. Liao, M.Y. Yang, and B. Rosenhahn, Text to Image Generation with
787 Semantic-Spatial Aware GAN. arXiv preprint arXiv:2104.00567, 2021
788
- 789 [38] Ramesh, A., M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I.
790 Sutskever, Zero-shot text-to-image generation. arXiv preprint arXiv:2102.12092,
791 2021
792
- 793 [39] Brown, T.B., B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A.
794 Neelakantan, P. Shyam, G. Sastry, and A. Askell, Language models are few-shot
795 learners. arXiv preprint arXiv:2005.14165, 2020
796
- 797 [40] Ramesh, A., M. Pavlov, G. Goh, and S. Gray, *DALL·E: Creating Images from Text*.
798 2021, OpenAI.
799
- 800 [41] Turing, A.M., Computing Machinery and Intelligence, in *Parsing the Turing Test:
801 Philosophical and Methodological Issues in the Quest for the Thinking Computer*,
802 R. Epstein, G. Roberts, and G. Beber, Editors. 2009, Springer Netherlands:
803 Dordrecht. p. 23-65.
804
- 805 [42] Boden, M.A., The Turing test and artistic creativity. *Kybernetes*, 2010. **39**(3): p.
806 409-413.DOI: 10.1108/03684921011036132
807
- 808 [43] Pease, A. and S. Colton. On impact and evaluation in computational creativity: A
809 discussion of the Turing test and an alternative proposal. Citeseer.
810
- 811 [44] Peter Berrar, D. and A. Schuster, Computing machinery and creativity: lessons
812 learned from the Turing test. *Kybernetes*, 2014. **43**(1): p. 82-91.DOI: 10.1108/K-
813 08-2013-0175
814
- 815 [45] Doersch, C., Tutorial on variational autoencoders. arXiv preprint arXiv:1606.05908,
816 2016
817
- 818 [46] Goodfellow, I., J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A.
819 Courville, and Y. Bengio, Generative adversarial nets. Advances in neural
820 information processing systems, 2014. **27**
821
- 822 [47] Amabile, T.M., Social psychology of creativity: A consensual assessment technique.
823 *Journal of personality and social psychology*, 1982. **43**(5): p. 997

824

825 [48] Zhu, M., P. Pan, W. Chen, and Y. Yang. Dm-gan: Dynamic memory generative
826 adversarial networks for text-to-image synthesis. in Proceedings of the IEEE/CVF
827 Conference on Computer Vision and Pattern Recognition. 2019.

828

829 [49] Lin, T.-Y., M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and
830 C.L. Zitnick. Microsoft coco: Common objects in context. in European
831 conference on computer vision. 2014. Springer.

832

833 [50] Sarica, S., J. Han, and J. Luo, Design representation as semantic networks.

834 Computers in Industry, 2023. **144**: p. 103791.DOI:

835 <https://doi.org/10.1016/j.compind.2022.103791>

836

837 [51] Zhu, Q., X. Zhang, and J. Luo, Biologically Inspired Design Concept Generation
838 Using Generative Pre-Trained Transformers. Journal of Mechanical Design, 2023.

839 **145**(4).DOI: 10.1115/1.4056598

840

841

842

843
844

Figure Captions List

- Figure 1 The workflow of the proposed research approach
- Figure 2 A question webpage in the Turing test
- Figure 3 Webpages of two question examples in the expert test
- Figure 4 The IS values of different test groups
- Figure 5 The FID scores of different test groups
- Figure 6 The FID scores in comparison with divided reference groups
- Figure 7 The base-FIDs and additive FIDs comparison within the machine and
human datasets
- Figure 8 The confusion matrix of Turing test results

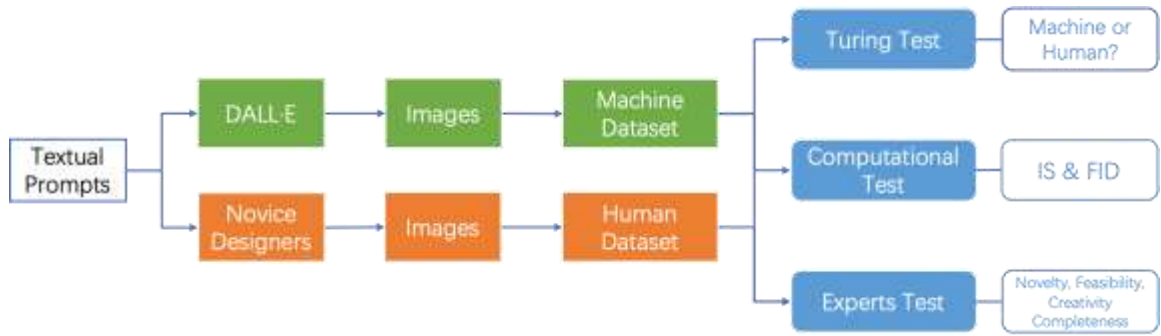
845
846

847
848

Table Caption List

- | | |
|---------|---|
| Table 1 | An overview of the machine and human data |
| Table 2 | An overview of the reference data |
| Table 3 | Results of the Turing test |
| Table 4 | Variance of accuracy in different groups |
| Table 5 | Results of expert test |

849
850



851
852
853
854
855

Figure 1. The workflow of the proposed research approach

856



857

858

Figure 2. A question webpage in the Turing test

859


860

Rate The image

Question 1-1

Text: a pentagonal green clock, a green clock in the shape of a pentagon.

Note: The image below is produced by human or machine based on above textual description.



1 2 3 4 5

Novelty

"1" means not novel, usual, non-original, non-attractive
 "5" means very novel, unusual, original, attractive

Feasibility

"1" means not feasible, non-common-sense, not natural, strange
 "5" means very feasible, common-sense, natural, normal

Creativity Completeness

"1" means low creativity completeness, different from textual description, combinational concepts not visible
 "5" means high creativity completeness, close to textual description, visible combinational concepts


(a)

Rate The Images

Question 1-11

Text: a pentagonal green clock, a green clock in the shape of a pentagon.

Note: The five images below are produced by human or machine exclusively based on the above textual description.



1 2 3 4 5

Variety

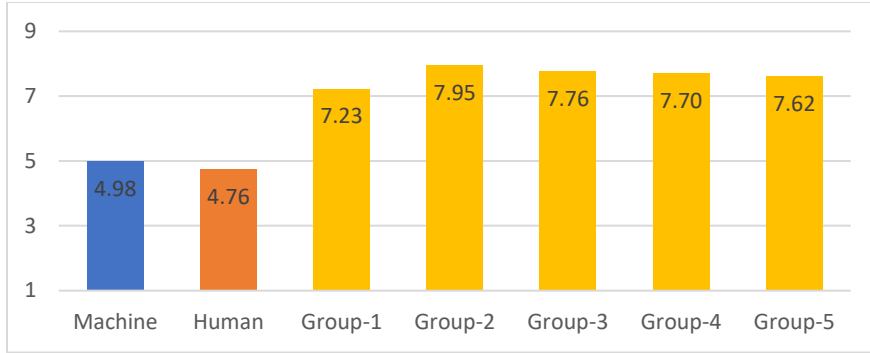
"1" means this set of images has low diversity, they have similar
 "5" means this set of images has high diversity, they look different

(b)

861

Figure 3. Webpages of two question examples in the expert test

862

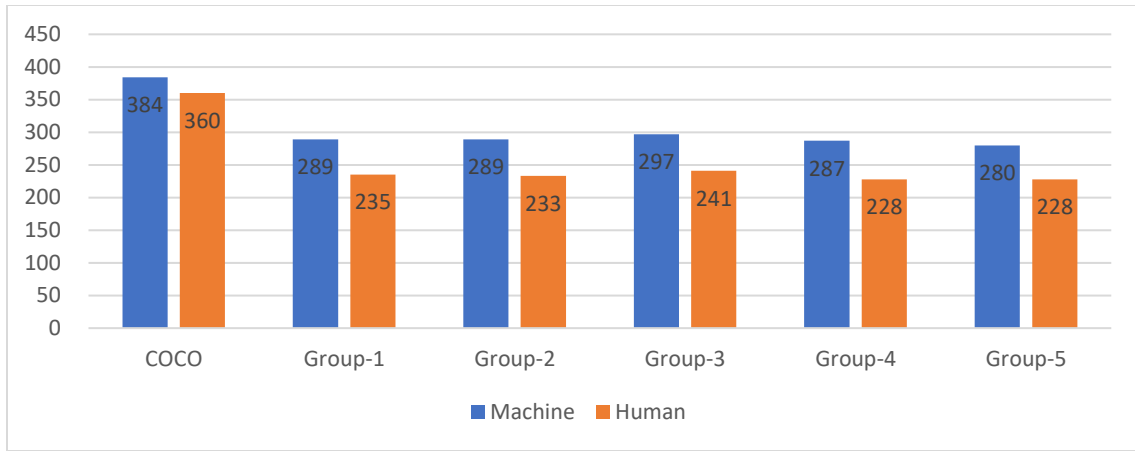


863

864

Figure 4. The IS values of different test groups

865



866

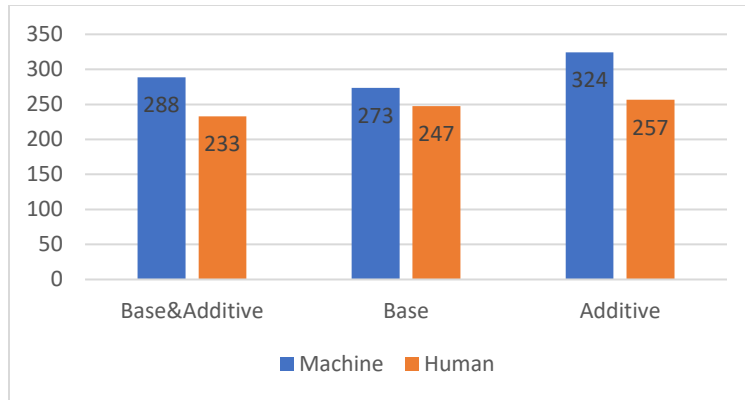
867

Figure 5. The FID scores of different test groups

868

869

870

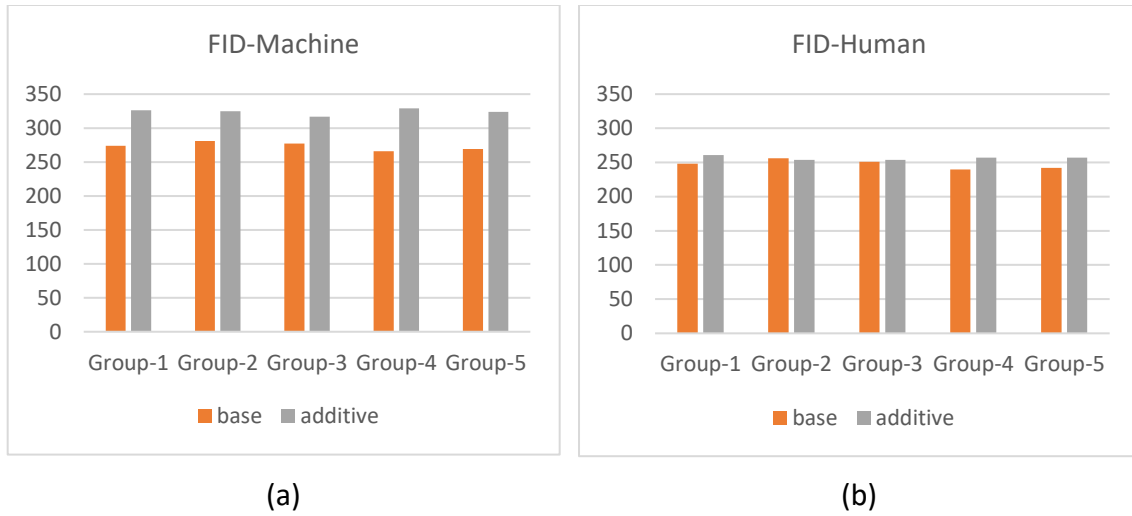


871

872

Figure 6. The FID scores in comparison with divided reference groups

873

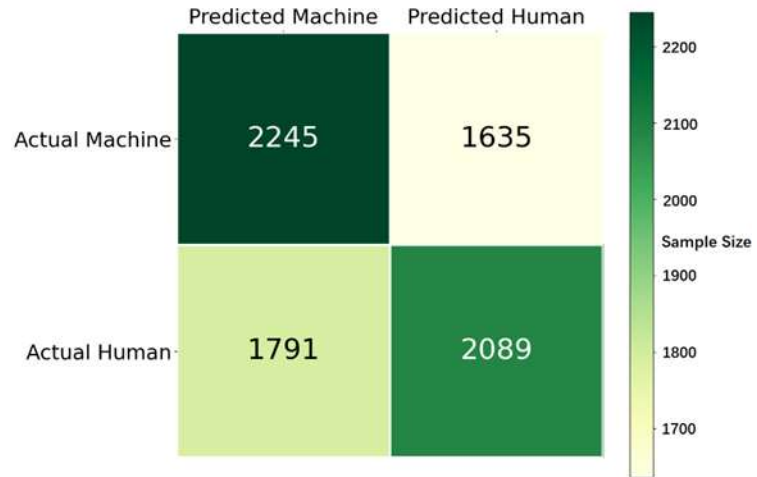


874 Figure 7. The base-FIDs and additive FIDs comparison within the machine (a) and human

875 (b) datasets

876

877



878

879

Figure 8. The confusion matrix of Turing test results

880





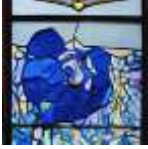











881

882

883

884

Table 1. An overview of the machine and human data




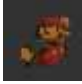











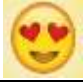
Group No.	Input	Machine Output	Human Output
1	a pentagonal green clock. a green clock in the shape of a pentagon		
2	a capybara made of voxels sitting in the field		
3	a stained-glass window with an image of a blue strawberry		
4	a snail made of harp. A snail with the texture of a harp		
5	an armchair in the shape of an avocado. an armchair imitating an avocado		
6	a giraffe imitating a turtle. a giraffe made of turtle		
7	a cube made of porcupine. a cube with the texture of a porcupine		
8	a professional high-quality emoji of a lovestruck cup of boba		

885

886

887

Table 2. An overview of the reference data

Group	Base	Sample-Base	Additive	Sample-Additive
1	Clock		Pentagonal	
2	Capybara		Voxels	
3	Glass		Strawberry	
4	Snail		Harp	
5	Armchair		Avocado	
6	Giraffe		Turtle	
7	Cube		Porcupine	
8	Cup		Emoji of lovestruck	

888

889

890

891

Table 3. Results of the Turing test

	Mean	1	2	3	4	5	6	7	8
Accuracy	55.9%	60.8%	55.2%	61.9%	60.3%	54.2%	51.1%	61.1%	42.2%
Precision	55.6%	60.6%	54.4%	63.5%	60.8%	54.2%	51.1%	59.0%	42.2%
Machine Recall	57.9%	61.9%	64.1%	55.7%	57.9%	54.2%	53.4%	73.0%	42.7%
F1	56.7%	61.2%	58.8%	59.3%	59.3%	54.2%	52.2%	65.3%	42.5%
Precision	56.1%	61.1%	56.3%	60.6%	59.8%	54.2%	51.2%	64.6%	42.1%
Human Recall	53.8%	59.8%	46.2%	68.0%	62.7%	54.2%	48.9%	49.3%	41.6%
F1	54.9%	60.4%	50.7%	64.1%	61.2%	54.2%	50.0%	55.9%	41.9%

892

893

894

895

Table 4. Variance of accuracy in different groups

		Min	Max	Max-Min	Difference
1	Human	37.1%	90.7%	53.6%	23.7%
	Machine	51.5%	81.4%	29.9%	
2	Human	19.6%	74.2%	54.6%	18.6%
	Machine	40.2%	76.3%	36.1%	
3	Human	43.3%	79.4%	36.1%	-2.1%
	Machine	38.1%	76.3%	38.1%	
4	Human	45.4%	93.8%	48.5%	26.8%
	Machine	46.4%	68.0%	21.6%	
5	Human	16.5%	89.7%	73.2%	45.4%
	Machine	43.3%	71.1%	27.8%	
6	Human	27.8%	77.3%	49.5%	7.2%
	Machine	32.0%	74.2%	42.3%	
7	Human	40.2%	63.9%	23.7%	-5.2%
	Machine	57.7%	86.6%	28.9%	
8	Human	25.8%	69.1%	43.3%	21.6%
	Machine	26.8%	48.5%	21.6%	
Overall	Human	16.5%	93.8%	77.3%	17.5%
	Machine	26.8%	86.6%	59.8%	
Mean	Human	32.0%	79.8%	47.8%	17.0%
	Machine	42.0%	72.8%	30.8%	

896

897

898

899

Table 5. Results of expert test

Metrics	Data Origin	1	2	3	4	5	6	7	8	Mean	Variance
Novelty	Machine	2.38	2.58	2.97	3.23	3.43	2.48	2.43	2.76	2.78	0.39
	Human	2.89	2.83	3.15	3.04	2.86	2.74	2.94	2.75	2.90	0.14
Feasibility	Machine	3.80	2.88	3.03	3.40	3.68	2.49	2.93	3.58	3.23	0.46
	Human	3.88	3.48	3.92	3.09	3.46	3.01	3.04	3.35	3.41	0.36
Completeness	Machine	3.48	2.68	2.98	3.57	3.42	2.44	2.54	3.64	3.09	0.49
	Human	3.53	3.06	3.64	3.46	3.48	3.40	2.89	3.44	3.36	0.25
Variety	Machine	3.00	2.74	2.37	3.42	3.47	3.16	3.26	2.21	2.95	0.47
	Human	3.84	4.37	3.84	2.89	3.05	3.21	3.21	3.74	3.52	0.50

900

901

902

903

904

905

906

907

908

909

910

911