

Liuqing Chen

College of Computer Science and Technology,
Zhejiang University,
Hangzhou 310058, China;
Zhejiang—Singapore Innovation and AI Joint
Research Lab,
Zhejiang University,
Hangzhou 310058, China
e-mail: chenlq@zju.edu.cn

Lingyun Sun

International Design Institute of ZJU,
Zhejiang University,
Hangzhou 310058, China;
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

A Comparison Study of Human and Machine-Generated Creativity

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

[DOI: 10.1115/1.4062232]

Keywords: artificial intelligence, computer-aided design, human–computer interfaces/interactions

1 Introduction

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

In recent years, a number of computational design support tools have been explored to tackle these limitations. For example, Han et al. [13] came up with an analogical reasoning tool for supporting idea generation by employing aspects of ontology and producing a corresponding image mood board; Sarica et al. [14] developed a technology semantic network based on patent data, which could support ideation by knowledge discovery; Siddharth et al. [15] proposed an engineering knowledge graph, containing <entity, relationship, entity> triples extracted from patent database, to support inference and reasoning; Obieke et al. [16] came up with a computational framework that explores new engineering design problems for creativity. Most of the existing so-called computational creativity tools do not generate creativity in a direct manner, but produce

stimuli instead, such as texts and images, to prompt designers' creative minds.

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

Evaluating combinational creativity is challenging, and there is no widely adopted method for such evaluation. In the field of design creativity, a variety of creativity assessment methods have been proposed, which generally require human raters to judge the quality of generated creativity [21], such as the Consensual Assessment Technique (CAT) method [22], Creative Product Semantic Scale [23], Product Creativity Measurement Instrument [24], Creative Solution Diagnosis Scale [25], and using creativity metrics [26,27]. In the field of artificial intelligence, the common computational metrics for evaluating generative models involve inception score (IS) [28] and Frechet inception distance (FID) [29], which are quantitative and calculated based on the probability distribution. In the interdisciplinary research between artificial intelligence and human study, the Turing test is a basic and widely adopted method [30–32], as it can provide an overall impression of how a machine performs. With consideration of the advantages of the evaluation methods in these three areas, this study applies a combined research approach by conducting a CAT-based expert test, a computational test, and a Turing test, and then synthesizes the results to elicit useful findings.

Therefore, the aim of this article is to compare the combinational creative performance of machines and human designers and to

¹Corresponding author.

Contributed by the Computers and Information Division of ASME for publication in the JOURNAL OF COMPUTING AND INFORMATION SCIENCE IN ENGINEERING. Manuscript received September 12, 2022; final manuscript received March 15, 2023; published online April 19, 2023. Assoc. Editor: Balan Gurumoorthy.

explore the differences between human designers and computers in generating creativity. This is the first study that compares the performance between novice designers and machines regarding combinational creativity, which employs a combined research approach integrating a Turing test, a computational test, and an expert test. This study will shed light on the research of computational creativity evaluation and artificial intelligence applications in design. The following section provides the theoretical background of this study. The methodology of the study is described in Sec. 3, and the implementation of the Turing test, computational test, and expert test is explained in Sec. 4. In Secs. 5 and 6, the results of the tests are presented, analyzed, and discussed. This article is then concluded in Sec. 7.

2 Theoretical Background

Combinational creativity is claimed to be one of the best approaches for fully utilizing nowadays abundant data, including texts, images, concepts, sounds, and so on [33], to achieve creativity [30]. A number of studies have explored combinational creativity in the context of design, particularly in idea generation. For instance, Nagai et al. [31] proposed three types of concept-synthesizing processes, namely, property mapping, concept blending, and concept integration in thematic relation, for generating new concepts based on three interpretation methods of combinational phrases. Han et al. [32] indicated that associating far-related ideas for forming combinational ideas could lead to outcomes that are more creative in comparison with linking closely related ones. Han et al. [34] investigated how combinational creativity is formed in design, focusing on conventional noun–noun combinations. It was revealed that a noun–noun combinational idea is produced by associating a base idea and an additive idea. The base idea refers to the basic idea of the combinational idea, while the additive idea could be a problem-solving idea, a similar representational idea, or an inspirational idea. For example, the famous Juicy Salif is an example of associating a basic idea (a manual juicer) and an inspirational additive idea (a squid). This study has thereby laid a theoretical foundation for our paper exploring human and machine-generated combinational creativity.

Although Han et al. [19] and Chen et al. [35,36] have employed pictorial data to form combinational images to facilitate users in combinational creativity, these combinational images are produced independently from semantic contexts. For instance, the Combinator [19] produces a compound phrase of “flower glass” and a corresponding combinational image of merging a “flower” and “glass.” Without semantic context, the combinational image produced could represent a “flower” made out of “glass,” a piece of “glass” in the shape of a “flower,” or a piece of “glass” with printed “flowers.” This might cause potential distractions and affect users’ creative performance.

In recent years, several computational models are developed to transform texts into images, such as LeicaGAN [37] and Semantic-Spatial Aware GAN [38]. These models could exploit text information for producing semantically consistent realistic images. Among them, DALL-E [22] is one of the most advanced ones, which employs GPT-3 [39] trained on a set of text image pairs data for

producing images based on text descriptions. As introduced by OpenAI [40], DALL-E has distinguishing capabilities, such as creating anthropomorphized versions of animals and objects. Moreover, it seems to have achieved a certain level of creativity. Specifically, the model could create pictorial combinations of unrelated concepts in plausible ways, even producing fantastical objects that do not exist in reality, according to textual descriptions. Thus, DALL-E is considered one of the most powerful systems capable of generating combinational creativity in pictorial formats within the constraints of texts. In this study, we perform a thorough performance benchmark evaluation comparing DALL-E with novice designers regarding combinational creativity, involving a Turing, a computational, and an expert test.

3 Methodology

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

3.1 Data Source—Machine and Human Datasets. Only a partial code of the DALL-E model was released on GitHub, and it is thereby impossible to run DALL-E to generate images due to missing training codes and data. Thus, the performance of DALL-E is evaluated based on the presented outcome from OpenAI’s official blog, in which the published data are representative and of high quality. In the blog, sets of textual descriptions and the corresponding generated images by DALL-E are presented. Three designers with over 3 years of experience were invited to judge whether the textual description in each set is a combinational idea. Before the judgment, the authors have well explained the definition of combinational creativity and showed some practical cases to the designers. If a set was judged as combinational creativity based, then five corresponding top-ranked images produced by DALL-E were collected. In total, eight sets, with five images in each set, are collected as the machine-generated combinational creativity dataset. All the input texts and one corresponding machine-produced image sample in each set are shown in Table 1.

Seven novice designers were employed to create a human dataset. They are either postgraduates or employees in companies with less than 3 years of working experience. They all hold a bachelor’s degree in design disciplines and have at least 2 years’ experience in product design and graphic design. Since the human dataset is associated with combinational creativity, before the creation of data, each designer was informed of the definition of combinational creativity and related design cases, especially the meaning of “base” and “additive.” Each designer was required to produce a

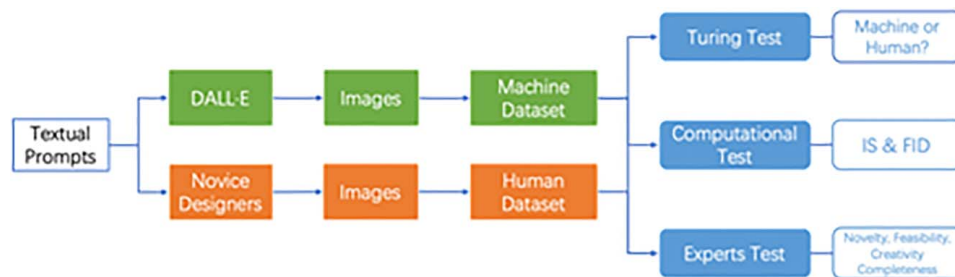




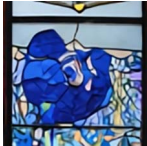
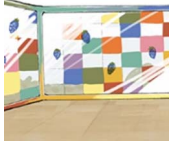












Fig. 1 The workflow of the proposed research approach

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		

drawing for each of the textual descriptions as indicated in Table 1 by using familiar computer-aided design software within 1 h. The designers were required to use white backgrounds and not to include any textual annotations to be in line with the ones of the machine dataset. Besides, the quality of drawings should be as high as possible, which is measured from three aspects:

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

- (3) Creativity completeness: The drawing should match the corresponding textual description, and combined concepts could be visible to recognize.

As a result, eight sets of data involving seven images each are produced. Three designers were then employed to select the top five images within each set. The eight sets of corresponding image samples produced by human designers are shown in Table 1.

3.2 Evaluation Methods

3.2.1 Turing Test. A Turing test [41] is conducted in this study to explore whether DALL-E can achieve combinational creativity at

the human level. In the test, participants were required to identify whether an image, within our mixed machine and human datasets, is produced by machine or human, providing the image's corresponding textual background. The test is consistent with the studies and arguments by Boden [42], Pease and Colton [43], and Peter Berrar and Schuster [44]. The test is specific and blinded and contains necessary contextual information. Although DALL-E is encouraged to produce realistic images in accordance with texts, it is not exclusively encouraged to exhibit creative behaviors. Therefore, the machine dataset, which can reflect DALL-E's capability of combinational creativity, was exclusively constructed to avoid possible tricky behaviors. For instance, instead of selecting the most realistic images generated by DALL-E to cheat human observers, we required that the images should first match their textual combinational ideas.

3.2.2 Computational Test. Given a deep learning-based model for image generation, such as variational autoencoders (VAE) [45] and generative adversarial networks (GANs) [46] based, the most common metrics for evaluating its capability are IS and FID. IS concerns the realism and diversity of generated images when evaluating a specific model. Specifically, IS calculates the KL divergence between the probability distribution of every generated image and the overall average of all generated images [28]. As shown in Eq. (1), given N classes, KL divergence is calculated between the conditional probability $p(y|x)$ in which a generated image x is classified into a particular class y , and the average probabilities for all the images in the class group $p(y)$, which is also called marginal distribution. High diversity of the generated images' categories and high certainty of the arbitrary image's category indicate high KL divergence, which means high IS and a better corresponding model, and there is no maximum value for IS.

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

FID is proposed to perform better in terms of discriminability, robustness, and computational efficiency and to address the limitations of IS [29]. It calculates the distance of two multidimensional normal distributions based on the mean (μ) and covariance (Σ) of the vectors extracted from both real (with the subscript r) and generated images (with the subscript g), as shown in the Eq. (2).

Ideally, the FID can be zero if the generated data are identical to real data, while higher FID value corresponds to the low quality of generated images. Considering the popularity of these two metrics in generative models' evaluation, we calculate both values for our machine and human datasets and then compare them.

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

3.2.3 Expert Test. The Turing test can estimate the overall appreciation of DALL-E's performance compared with humans by subjective evaluation, while the computational test can quantitatively and objectively compare machine and human performance but lack detailed and interpretable criteria. Hence, an expert test is necessary to deeply investigate the difference between the two groups and provide interpretable results. In this study, a CAT-based method [47] is adopted in the expert test for creativity evaluation.

Novelty, quantity, quality, and variety are the four metrics often used in design research for evaluating creativity [26]. In the expert test, a modified version of the metrics was adopted. Novelty, feasibility, and creativity completeness were used to measure a single image, and variety was used to measure a group of images generated by either a human designer or machine. The combinational creativity images are generated based on textual descriptions, thus novelty originates from the creation of combining the "base" elements with the "additive" elements, such as the novelty of the creation of combining "armchair" with "avocado" in an imagery format. On the other hand, creativity completeness is an essential metric for evaluating the transformation quality from textual description to imagery visualization, instead of focusing on evaluating creation results (novelty). Since some of the combinational ideas are imaginary rather than physical, such as "a giraffe imitating a turtle," feasibility is chosen as the metric instead of quality and utility. The meanings of novelty, feasibility, and creativity completeness are identical to the descriptions for ranking drawings in

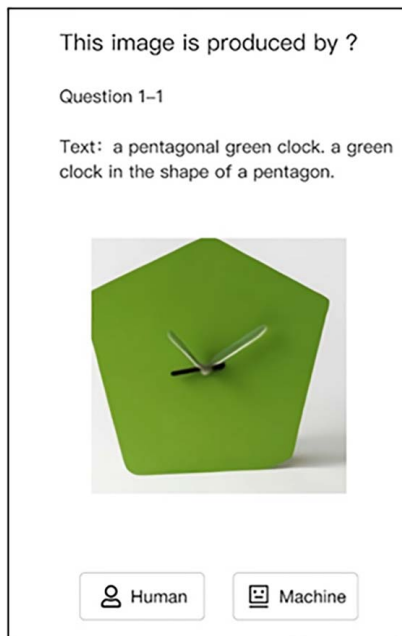


Fig. 2 A question webpage in the Turing test

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	

the human dataset indicated in the preceding. Variety refers to the diversity of a set of images, which measures the differences between images.

4 Evaluation

4.1 Turing Test. The Turing test is conducted by developing a website where all web pages are completely customized to minimize distractions. Participants were asked to read the instructions, agree with the test protocols, and provide demographic information

before starting the test. Eight groups of questions in total, corresponding to eight groups of data in our datasets, are provided to the participants. Each group contains ten questions that are randomly ordered for mixing the human- and machine-generated data, while five questions are from the human dataset and another five are from the machine dataset. This fact is not revealed to the participants to avoid introducing any potential bias. This would not influence participants' choices since they could feel free to make decisions without restrictions. There is only one question on each webpage consisting of a question serial number, a short


(a)

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


(b)

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

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

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%
Machine	Precision	55.6%	60.6%	54.4%	63.5%	60.8%	54.2%	51.1%	59.0%	42.2%
	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%
Human	Precision	56.1%	61.1%	56.3%	60.6%	59.8%	54.2%	51.2%	64.6%	42.1%
	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%

textual description, an image that is either from the human dataset or the machine dataset, and two buttons, indicating “human” and “machine” for participants to choose, as shown in Fig. 2. The participants were required to spend at least 3 s on each question before moving to the next one.

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

- (1) How difficult do you think this test is?
- (2) What is your method for distinguishing human and machine?
- (3) What is your feedback about this test?

Answers of the interviews were collected and analyzed in a qualitative way, and the results were reported in the Sec. 5.

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

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

4.3 Expert Test. The expert test was also conducted via a customized website. There are eight groups of questions, and each contains 12 questions. In each group, the first ten questions are single image based, of which a textual description and corresponding image are provided in each question, and participants are required to rate the image using a five-point Likert scale regarding three metrics: novelty, feasibility, and creativity completeness, as shown in Fig. 3(a). The ten images are randomly selected from the human or machine datasets. The last two questions in each group are five-image based, in which a textual description and corresponding five images (merged in a vertical sequence) are shown. Participants are informed that all five images were generated by humans or machines exclusively, and they are required to rate the

variety of the five images using a five-point Likert scale, as shown in Fig. 3(b).

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

5 Results and Analysis

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

Accuracy concerns whether a question is correctly answered, rather than which answer is more often answered. Given a classification problem, human or machine classes in our case, three metrics are widely applied when measuring the performance of a classification machine learning model: precision, recall, and F1 score. The formulas of the three metrics are given in Eqs. (3)–(5), respectively,

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%	

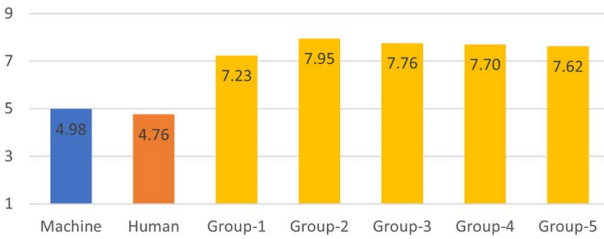


Fig. 4 The IS values of different test groups

where TP represents true positive and FN represents false negative. In our calculation, positive means the answer is “human,” while negative indicates “machine.” The results of precision, recall, and F1 score of the two classes (human and machine) are presented in Table 3. As shown in the table, the precision between human and machine is very close (56.1% versus 55.6%), but the recall between human and machine are noticeably different. The recall of the machine class is higher than the human class by 7.6%, which is due to high TN and high FN. Besides, the F1 score of the machine dataset is higher than human by 3.3%.

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

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

$$F1 \text{ score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5)$$

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

Twenty participants accepted the interview and answered questions after completing the Turing test. Concerning the method of distinguishing human and machine, the participants indicated that they believe the human-generated images have “more clear details,” “a unified style (such as sketches),” and “high resolutions,”

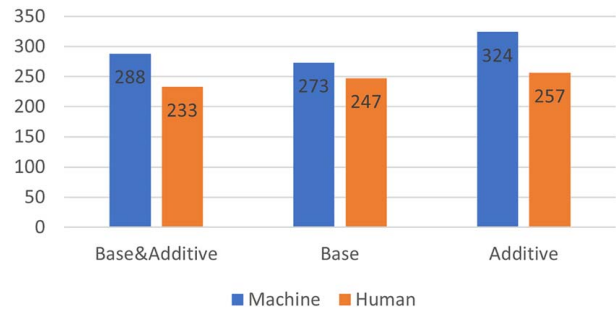


Fig. 6 The FID scores in comparison with divided reference groups

while the machine-generated images are “unreal,” are “blurred,” and have “unhuman combination logics” and “cut and paste by Photoshop patterns.” In terms of the difficulty of the task, the participants suggested that natural or physical subjects are easy to make “human” or “machine” selections, as well as images employing sketch styles. The interview results are a supplement to the Turing test and can potentially explain the Turing test results and help understand the reasons underpinning the choices made by the participants. This is in line with other similar studies. For example, Sarica et al. [50] interviewed 25 participants to understand their choices of the best computational representation of a specific design, and Zhu et al. [51] interviewed ten engineers regarding their views toward a set of computationally generated design concepts.

5.2 Computational Test. The computed results of the IS are shown in Fig. 4 where the IS of five reference groups are presented together for reference purposes. The machine group has a higher IS than the human group by 4.6%. The IS of the five reference groups are much higher than the machine and human datasets with an average IS of 7.65 ($\sigma=0.27$). The computed FID scores including reference groups are presented in Fig. 5. When comparing with COCO datasets, the FID of the machine dataset is higher than the human dataset by 6.7%. All the FID scores in comparison with reference groups are lower than COCO datasets, and all the FID scores of the machine group are higher than the human group. The average FID of the machine group in comparison with the five reference groups is 288 ($\sigma=6.07$), which is higher than the average FID of the human group ($\mu=233$, $\sigma=5.43$) by 23.8%.

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, as shown in Fig. 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 ($\mu=$

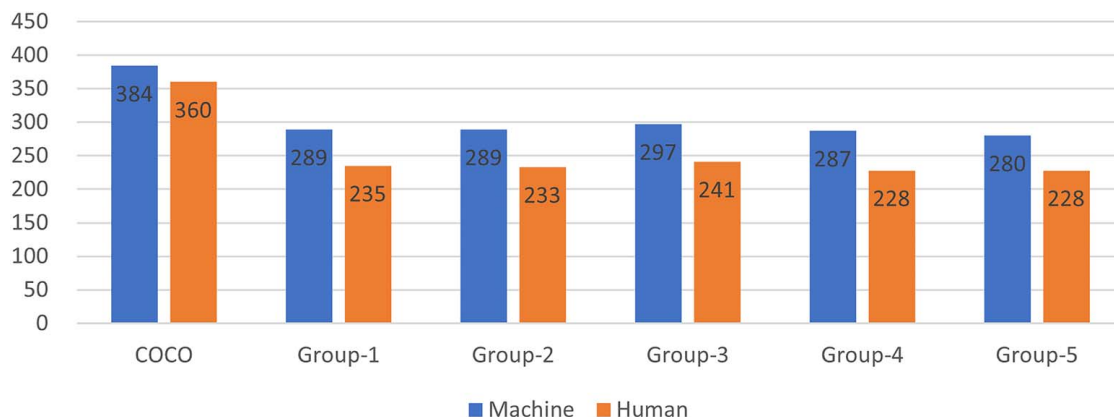


Fig. 5 The FID scores of different test groups

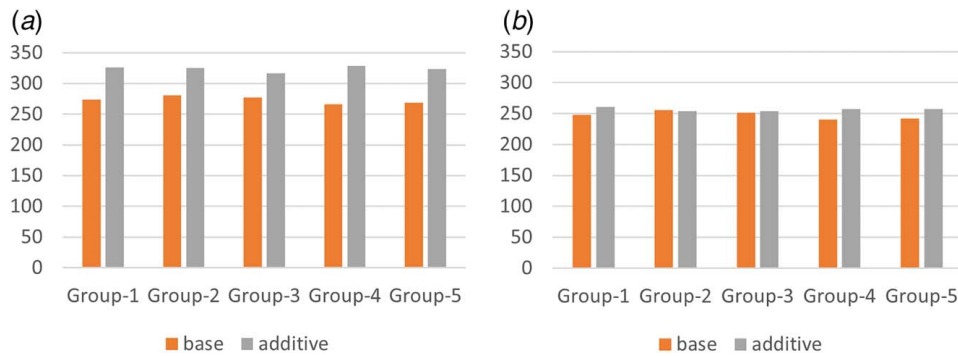


Fig. 7 The base-FIDs and additive-FIDs comparison within the (a) machine and (b) human datasets

324, $\sigma=4.44$) is significantly higher than the human group ($\mu = 257, \sigma=2.88$) by 26% in comparison with additive groups. It is useful to investigate the influence of base and additive on the overall FID. The FID scores in comparison with five base or additive groups (called base-FID and additive-FID, respectively) are presented in Fig. 7. As shown in the Fig. 7(a) (machine dataset), the additive-FIDs are higher than the base-FIDs on average by 18.7%, while Fig. 7(b) (human dataset) shows that the additive-FIDs are slightly higher than the base-FIDs only by 4.0%.

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

6 Discussion

6.1 Turing Test. The average mathematical expectation of random answers to all the questions in the Turing test is 50%, while the closer of overall accuracy to 50% indicates the more undistinguishable between human and machine-generated data.

Though the overall accuracy in the Turing test is above 50%, the gap is only 5.9%. The F1 scores of the machine and human datasets are both close to 50%, while the machine's score is slightly higher than the human's score due to high recall in the machine dataset. High variance within every group in both datasets indicates that participants have low certainty to make their judgments. Besides, as indicated in the confusion matrix in Fig. 8, TN (predicted machine and actual machine) and FN (predicted machine and actual human) are relatively higher, which corresponds to higher recall and F1 score of the machine dataset. This suggests that the results reveal that DALL-E can deceive participants to a large extent, and the participants could hardly indicate which image is from the human or machine dataset, while the participants subjectively tended to believe that the data in the Turing test were more likely from machines rather than humans.

From our interview, it is shown that designers tend to use sketch and image processing software (such as PHOTOSHOP) to create drawings rather than 3D modeling and rendering, which makes their drawings more distinguishable from machine data. On the other hand, the images generated by DALL-E tend to be blurred, unsmooth, and unreal due to technical limitations, which makes them distinct from normal images. Besides, the logic behind a combination idea in machine data is sometimes different from human data. The "cut and paste by PHOTOSHOP" pattern is considered a machine pattern by some participants, since some designers tend to create a collage-style image to express a combination idea while participants believe that machine is good at creating collages.

6.2 Computational Test. We created five reference groups in the computational test, and all the results related to the five groups have low variance, which indicates that there is only little bias brought into the reference groups. Regarding the IS metric, the machine dataset achieved a higher IS score than the human dataset, which means the machine-generated data have higher quality than the designers', but the gap is as little as 4.6%. Five reference groups obtained much higher IS, since these reference images contain rich information about bases and additives, and they are natural rather than combinational, which is more favored by the inception model used for calculating IS. On the other hand, the

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

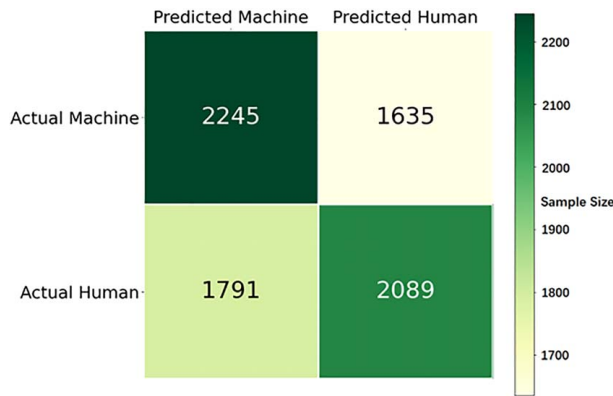


Fig. 8 The confusion matrix of Turing test results

machine dataset obtained a higher FID score than the human dataset when comparing with both COCO data and the five reference groups of data, indicating the machine-generated data have a lower quality than the human-generated data. All the FID scores in comparison with the reference groups are significantly lower than COCO datasets, validating that the images in our reference groups are closer to both the machine and human data than the images in COCO. The difference of FIDs between the human and machine datasets in comparison with five reference groups is bigger than the difference of FIDs in comparison with COCO data. This may reflect the difference in combinational design between humans and machines. Since it is required that drawings should be produced based on textual descriptions containing combinational creativity, novice designers tend to keep essential information from both base and additive in a combinational design, while DALL-E is not trained to obtain this capability. This indicates that these designers have a better lingual understanding of combinational ideas and are able to transform them into designs than machines.

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

6.3 Expert Test. The human dataset obtained higher scores than the machine dataset by a small percentage (6.17% on average) when comparing the results regarding novelty, feasibility, and creativity completeness, despite that the machine dataset has higher scores in some groups. This indicates that the novice designers performed slightly better than DALL-E in combinational designs in these three metrics. Besides, the designers outperform DALL-E evidently regarding variety by an overall gap of 19.15%, even though the machine dataset outperformed in three groups. This gap could be explained by two reasons. One is that the human data are from seven novice designers, while the machine data are from DALL-E exclusively, which is unfair for DALL-E in this test. Another reason is the difference in working mechanism between the DALL-E model and designers, in which DALL-E takes text as input and generates various images based on random noise, while designers are skilled in producing various images using divergent thinking. It is noticed that two to three groups in the machine data have higher scores regarding all four metrics, indicating the capability of producing combinational creativity images between novice designers and DALL-E is not significantly different.

6.4 Overall Discussion. There are no clear criteria to determine whether DALL-E passes the Turing test, but it can be

concluded that DALL-E's performance is close to novice designers according to the results of our Turing test. In the computational test, DALL-E outperforms designers in terms of IS but loses to designers regarding FID, and the difference in values is both small, indicating that the performance between DALL-E and novice designers is very close. It is noticed that the results of IS and FID are in conflict, which indicates that the effectiveness of the two metrics for evaluating combinational creativity needs to be further investigated. A larger difference in FIDs in comparison with our reference data implies that human designers are better at synthesizing features from base and additive for a combinational design. According to the results of the expert test, designers outperform DALL-E from the perspective of combinational creativity. There is slight advance for designers regarding novelty, feasibility, and creativity completeness, but evident advance regarding variety. By summarizing the conclusions from the three tests in this study, DALL-E's performance is no better than novice designers, but the gap is small.

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

There are a few limitations in this study. First, eight sets of data related to combinational creativity, containing 40 machine-generated images and 40 human-generated ones, were used in the study for evaluation. The limited amount of data was a result of the restricted access to DALL-E's source code and data, as well as the high cost of human resources. Although the amount of data is sufficient for the purpose of the study, more data will be included in future studies by recruiting more human designers and accessing more DALL-E data to yield further useful insights. This would require the involvement of more human designers and accessing more DALL-E's data. Second, 1 h was provided to the designers to complete one combinational creativity design task to construct the human dataset, but it is still far less to produce a high-quality image. More time will be provided to the participants in future research to improve the quality of the images generated. Third, DALL-E is a deep learning model mainly aiming at transforming texts into images rather than generating combinational creativity, which is less fair to compare with human designers. In future research, more advanced artificial intelligence models, such as ChatGPT and GPT-4, will be included in the comparison.

7 Conclusion

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

methods for design. It leads toward new research directions in evaluating computational creativity and applying advanced computational techniques, particularly in conceptual design.

Funding Data

- The National Natural Science Foundation of China (Grant No. 62207023).
- The Ng Teng Fong Charitable Foundation in the form of ZJU-SUTD (Zhejiang University—Singapore University of Technology and Design) IDEA grant.

Conflict of Interest

There are no conflicts of interest.

Data Availability Statement

The datasets generated and supporting the findings of this article are obtainable from the corresponding author upon reasonable request.

References

- [1] Childs, P., Han, J., Chen, L., Jiang, P., Wang, P., Park, D., Yin, Y., Dieckmann, E., and Vilanova, I., 2022, "The Creativity Diamond—A Framework to Aid Creativity," *J. Intell.*, **10**(4), p. 73.
- [2] Amabile, T. M., 1983, *The Social Psychology of Creativity*, Springer, New York.
- [3] Shute, V. J., and Rahimi, S., 2021, "Stealth Assessment of Creativity in a Physics Video Game," *Comput. Hum. Behav.*, **116**, p. 106647.
- [4] De Bono, E., 1985, *Six Thinking Hats*, Little, Brown and Company, Boston, MA.
- [5] Eberle, B., 1996, *Scamper: Games for Imagination Development*, Prufrock Press, Waco, TX.
- [6] Zwicky, F., 1969, *Discovery, Invention, Research Through the Morphological Approach*, Macmillan, New York.
- [7] Altshuller, G. S., 1984, *Creativity as an Exact Science: The Theory of the Solution of Inventive Problems*, Gordon and Breach Publishers, Amsterdam, Netherlands.
- [8] Linsey, J. S., Markman, A. B., and Wood, K. L., 2012, "Design by Analogy: A Study of the WordTree Method for Problem Re-Representation," *ASME J. Mech. Des.*, **134**(4), p. 041009.
- [9] Yilmaz, S., Daly, S. R., Seifert, C. M., and Gonzalez, R., 2016, "Evidence-Based Design Heuristics for Idea Generation," *Des. Stud.*, **46**, pp. 95–124.
- [10] Helms, M., Vattam, S. S., and Goel, A. K., 2009, "Biologically Inspired Design: Process and Products," *Des. Stud.*, **30**(5), pp. 606–622.
- [11] Chakrabarti, A., and Shu, L. H., 2010, "Biologically Inspired Design," *Artif. Intell. Eng. Des. Anal. Manuf.*, **24**(4), pp. 453–454.
- [12] Oman, S. K., Tumer, I. Y., Wood, K., and Seepersad, C., 2013, "A Comparison of Creativity and Innovation Metrics and Sample Validation Through in-Class Design Projects," *Res. Eng. Des.*, **24**(1), pp. 65–92.
- [13] Han, J., Shi, F., Chen, L., and Childs, P. R. N., 2018, "A Computational Tool for Creative Idea Generation Based on Analogical Reasoning and Ontology," *Artif. Intell. Eng. Des. Anal. Manuf.*, **32**(4), pp. 462–477.
- [14] Sarica, S., Luo, J., and Wood, K. L., 2020, "TechNet: Technology Semantic Network Based on Patent Data," *Expert Syst. Appl.*, **142**, p. 112995.
- [15] Siddharth, L., Blessing, L. T. M., Wood, K. L., and Luo, J., 2022, "Engineering Knowledge Graph From Patent Database," *ASME J. Comput. Inf. Sci. Eng.*, **22**(2), p. 021008.
- [16] Obieke, C. C., Milisavljevic-Syed, J., Silva, A., and Han, J., 2023, "A Computational Approach to Identifying Engineering Design Problems," *ASME J. Mech. Des.*, **145**(4), p. 041406.
- [17] Boden, M. A., 2004, *The Creative Mind: Myths and Mechanisms*, 2nd ed., Routledge, London.
- [18] Simonton, D. K., 2017, "Domain-General Creativity: On Generating Original, Useful, and Surprising Combinations," *The Cambridge Handbook of Creativity Across Domains*, J. C. Kaufman, V. P. Glăveanu, and B. John, eds., The Cambridge University Press, Cambridge, UK, pp. 18–40.
- [19] Han, J., Shi, F., Chen, L., and Childs, P. R. N., 2018, "The Combinator—A Computer-Based Tool for Creative Idea Generation Based on a Simulation Approach," *Des. Sci.*, **4**, p. e11.
- [20] Garvey, B., Chen, L., Shi, F., Han, J., and Childs, P. R., 2019, "New Directions in Computational, Combinational and Structural Creativity," *Proc. Inst. Mech. Eng. C: J. Mech. Eng. Sci.*, **233**(2), pp. 425–431.
- [21] Beaty, R. E., and Johnson, D. R., 2021, "Automating Creativity Assessment With SemDis: An Open Platform for Computing Semantic Distance," *Behav. Res. Methods*, **53**(2), pp. 757–780.
- [22] Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., and Sutskever, I., 2021, "Zero-Shot Text-to-Image Generation," The 38th International Conference on Machine Learning, PMLR, Virtual, July 18–24.
- [23] Besemer, S. P., and O'quin, K., 1986, "Analyzing Creative Products: Refinement and Test of a Judging Instrument," *J. Creat. Behav.*, **20**(2), pp. 115–126.
- [24] Horn, D., and Salvendy, G., 2006, "Product Creativity: Conceptual Model, Measurement and Characteristics," *Theor. Issues Ergon. Sci.*, **7**(4), pp. 395–412.
- [25] Cropley, D., and Cropley, A., 2005, "Engineering Creativity: A Systems Concept of Functional Creativity," *Creativity Across Domains: Faces of the Muse*, J. C. Kaufman, and J. Baer, eds., Lawrence Erlbaum Associates Publishers, Mahwah, NJ, pp. 169–185.
- [26] Shah, J. J., Smith, S. M., and Vargas-Hernandez, N., 2003, "Metrics for Measuring Ideation Effectiveness," *Des. Stud.*, **24**(2), pp. 111–134.
- [27] Han, J., Forbes, H., and Schaefer, D., 2021, "An Exploration of How Creativity, Functionality, and Aesthetics Are Related in Design," *Res. Eng. Des.*, **32**(3), pp. 289–307.
- [28] Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A., 2017, "Improved Training of Wasserstein Gans," *Adv. Neural Inf. Process. Syst.*, **30**, pp. 5769–5779.
- [29] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S., 2017, "Gans Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium," *Adv. Neural Inf. Process. Syst.*, **30**, pp. 6629–6640.
- [30] Ward, T. B., and Kolomyts, Y., 2010, "Cognition and Creativity," *The Cambridge Handbook of Creativity*, J. C. Kaufman, and R. J. Sternberg, eds., The Cambridge University Press, Cambridge, UK, pp. 93–112.
- [31] Yang, H., and Zhang, L., 2016, "Promoting Creative Computing: Origin, Scope, Research and Applications," *Digit. Commun. Netw.*, **2**(2), pp. 84–91.
- [32] Nagai, Y., Taura, T., and Mukai, F., 2009, "Concept Blending and Dissimilarity: Factors for Creative Concept Generation Process," *Des. Stud.*, **30**(6), pp. 648–675.
- [33] Han, J., Shi, F., Park, D., Chen, L., and Childs, P., 2018, "The Conceptual Distances Between Ideas in Combinational Creativity," *DS92: Proceedings of the DESIGN 2018 15th International Design Conference*, Dubrovnik, Croatia, May 21–24.
- [34] Han, J., Park, D., Shi, F., Chen, L., Hua, M., and Childs, P. R., 2019, "Three Driven Approaches to Combinational Creativity: Problem-, Similarity- and Inspiration-Driven," *Proc. Inst. Mech. Eng. C: J. Mech. Eng. Sci.*, **233**(2), pp. 373–384.
- [35] Chen, L., Wang, P., Dong, H., Shi, F., Han, J., Guo, Y., Childs, P. R. N., Xiao, J., and Wu, C., 2019, "An Artificial Intelligence Based Data-Driven Approach for Design Ideation," *J. Vis. Commun. Image Represent.*, **61**, pp. 10–22.
- [36] Chen, L., Wang, P., Shi, F., Han, J., and Childs, P., 2018, "A Computational Approach for Combinational Creativity in Design," *DS 92: Proceedings of the DESIGN 2018 15th International Design Conference*, Dubrovnik, Croatia, May 21–24.
- [37] Qiao, T., Zhang, J., Xu, D., and Tao, D., 2019, "Learn, Imagine and Create: Text-to-Image Generation From Prior Knowledge," *Adv. Neural Inf. Process. Syst.*, **32**, pp. 887–897.
- [38] Liao, W., Hu, K., Yang, M. Y., and Rosenhahn, B., 2021, "Text to Image Generation with Semantic-Spatial Aware GAN," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, New Orleans, LA, June 18–24.
- [39] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., and Askell, A., 2020, "Language Models are Few-Shot Learners," *Adv. Neural Inf. Process. Syst.*, **33**, pp. 1877–1901.
- [40] Ramesh, A., Pavlov, M., Goh, G., and Gray, S., 2021, DALL-E: Creating Images From Text, <https://openai.com/research/dall-e>.
- [41] Turing, A. M., 2009, "Computing Machinery and Intelligence," *Parsing the Turing Test: Philosophical and Methodological Issues in the Quest for the Thinking Computer*, R. Epstein, G. Roberts, and G. Beber, eds., Springer Netherlands, Dordrecht, pp. 23–65.
- [42] Boden, M. A., 2010, "The Turing Test and Artistic Creativity," *Kybernetes*, **39**(3), pp. 409–413.
- [43] Pease, A., and Colton, S., 2011, "On Impact and Evaluation in Computational Creativity: A Discussion of the Turing Test and an Alternative Proposal," *Proceedings of the AISB Symposium on AI and Philosophy*, York, UK, Apr. 4–7.
- [44] Peter Berrar, D., and Schuster, A., 2014, "Computing Machinery and Creativity: Lessons Learned From the Turing Test," *Kybernetes*, **43**(1), pp. 82–91.
- [45] Doersch, C., 2016, "Tutorial on Variational Autoencoders," *arXiv preprint arXiv:1606.05908*.
- [46] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y., 2014, "Generative Adversarial Nets," *Adv. Neural Inf. Process. Syst.*, **27**, pp. 2672–2680.
- [47] Amabile, T. M., 1982, "Social Psychology of Creativity: A Consensual Assessment Technique," *J. Pers. Soc. Psychol.*, **43**(5), pp. 997–1013.
- [48] Zhu, M., Pan, P., Chen, W., and Yang, Y., 2019, "Dm-gan: Dynamic Memory Generative Adversarial Networks for Text-to-Image Synthesis," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Long Beach, CA, June 16–20.
- [49] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L., 2014, "Microsoft Coco: Common Objects in Context," *European Conference on Computer Vision*, Zurich, Switzerland, Sept. 6–12.
- [50] Sarica, S., Han, J., and Luo, J., 2023, "Design Representation as Semantic Networks," *Comput. Ind.*, **144**, p. 103791.
- [51] Zhu, Q., Zhang, X., and Luo, J., 2023, "Biologically Inspired Design Concept Generation Using Generative Pre-Trained Transformers," *ASME J. Mech. Des.*, **145**(4), p. 041409.