

Towards AGI in Computer Vision: Lessons Learned from GPT and Large Language Models

Lingxi Xie, Longhui Wei, Xiaopeng Zhang, Kaifeng Bi, Xiaotao Gu, Jianlong Chang,
and Qi Tian*, *Fellow, IEEE*

Abstract—The AI community has been pursuing algorithms known as artificial general intelligence (AGI) that apply to any kind of real-world problem. Recently, chat systems powered by large language models (LLMs) emerge and rapidly become a promising direction to achieve AGI in natural language processing (NLP), but the path towards AGI in computer vision (CV) remains unclear. One may owe the dilemma to the fact that visual signals are more complex than language signals, yet we are interested in finding concrete reasons, as well as absorbing experiences from GPT and LLMs to solve the problem. In this paper, we start with a conceptual definition of AGI and briefly review how NLP solves a wide range of tasks via a chat system. The analysis inspires us that **unification** is the next important goal of CV. But, despite various efforts in this direction, CV is still far from a system like GPT that naturally integrates all tasks. We point out that the essential weakness of CV lies in lacking a paradigm to **learn from environments**, yet NLP has accomplished the task in the text world. We then imagine a pipeline that puts a CV algorithm (*i.e.*, an agent) in world-scale, interactable environments, pre-trains it to predict future frames with respect to its action, and then fine-tunes it with instruction to accomplish various tasks. We expect substantial research and engineering efforts to push the idea forward and scale it up, for which we share our perspectives on future research directions.

Index Terms—Computer Vision, Artificial General Intelligence, Foundation Models, Unification, Environments.

1 INTRODUCTION

THE world is witnessing an epic odyssey towards artificial general intelligence (AGI), where we follow the convention to define AGI as a computer algorithm that can replicate any intellectual task that human beings or other animals can¹. Specifically, in natural language processing (NLP), computer algorithms have been developed to an extent that can solve a wide range of tasks via chat with humans [1]. Some researchers believed that such systems can be seen as early sparks of AGI [2]. These systems were mostly built upon large language models (LLMs) [3] and enhanced by instruct tuning [4]. Equipped with an external knowledge base and specifically designed modules, they can accomplish complex tasks such as solving mathematical questions, generating visual contents, *etc.*, reflecting its strong ability to understand users' intentions and perform preliminary chain-of-thoughts [5]. Despite known weaknesses in some aspects (*e.g.*, telling scientific facts and relationships between named people), these pioneering studies have shown a clear trend to unify most tasks in NLP into one system, which reflects the pursuit of AGI.

Compared to the rapid progress of unification in NLP, the computer vision (CV) community is yet far from the target of unifying all tasks. The regular CV tasks, such as visual recognition, tracking, captioning, generation, *etc.*, are mostly processed using largely different network architectures and/or specifically designed pipelines. Researchers look forward to a system like GPT that can deal with a wide

range of CV tasks with a unified prompt mechanism, but there exists a tradeoff between achieving good practice in individual tasks and being generalized across a wide range of tasks. For example, to report high recognition accuracy in object detection and semantic segmentation, the best strategy is to design specific head modules [6], [7] upon strong backbones [8], [9], [10] for image classification, and such designs do not generally transfer to other problems such as image captioning [11] or visual content generation [12].

Clearly, unification is the trend in CV. In recent years, there are many efforts in this direction, and we roughly categorize them into five research topics, namely, (i) **open-world visual recognition** based on vision-language alignment [13], (ii) **the Segment Anything task** [14] for generic visual recognition, (iii) **generalized visual encoding** to unify vision tasks [15], [16], [17], (iv) LLM-guided visual understanding to enhance the logic in CV [18], [19], and (v) **multimodal dialog** to facilitate vision-language interaction [11], [20]. These works showed promise of unification, but yet, they cannot composite a system like GPT that can solve generic CV tasks in the real world.

Hence, two questions arise: (1) Why is unification so difficult in CV? (2) What can we learn from GPT and LLMs to achieve this goal? To answer them, we revisit GPT and understand it as establishing an environment in the text world and allowing an algorithm (or agent) to learn from interaction. The CV research lacks such an environment. Consequently, the algorithms cannot simulate the world, so they instead sample the world and learn to achieve good performance in the so-called proxy tasks. After an epic decade of deep learning [21], the proxy tasks are no longer meaningful to indicate the ability of CV algorithms; it becomes more and more apparent that continuing to pursue

- All authors, unless specified below, are with Huawei Inc., China. E-mail of the leading author (Lingxi Xie): 198808xc@gmail.com
- Corresponding author: Qi Tian. E-mail: tian.qi1@huawei.com

Manuscript received Month Date, 2023.

1. https://en.wikipedia.org/wiki/Artificial_general_intelligence

high accuracy on them can drive us away from AGI.

Based on the analysis above, we propose an imaginary pipeline towards AGI in CV. It involves three stages. The first stage is to establish a set of environments that are faithful, abundant, and interactable. The second stage aims to train an agent by forcing it to explore the environment(s) and predict future frames: this corresponds to the autoregressive pre-training stage in NLP [3]. The third stage involves teaching the agent to accomplish various tasks: it is likely that human instructions shall be introduced in this stage, corresponding to the instruct fine-tuning stage in NLP [4]. Optionally, the agent can be tuned to perform proxy tasks via simple and unified prompts. The idea is related to a few existing research topics, including 3D environment establishment [22], [23], visual pre-training [24], [25], reinforcement learning [26], [27], and embodied CV [28], [29]. But, existing works are mostly preliminary and we expect that substantial efforts [30], [31] are required to make it an effective paradigm to solve real-world problems.

The remainder of this paper is organized as follows. First, in Section 2, we briefly introduce the history and thoughts of AGI and inherit the definition that AGI is an algorithm to maximize the reward. It is followed by Section 3 where we show the ability of GPT, the state-of-the-art NLP algorithm which was considered the spark of AGI. Then, in Section 4, based on the current status of CV research, we analyze why AGI is difficult in computer vision and point out that the essential difficulty lies in the outdated learning paradigm. The analysis leads to Section 5, where we imagine a pipeline that pushes CV closer to AGI, based on which we make some comments on future research directions. Finally, in Section 6, we conclude this paper and share our thoughts.

2 ARTIFICIAL GENERAL INTELLIGENCE

Artificial intelligence (AI) is a long-lasting battle to replicate human intelligence with a machine or a set of mathematical algorithms. Modern AI was formally proposed in the Dartmouth workshop, 1956, and the community has developed a large number of methodologies for this purpose. There are at least two different pathways to achieve AI: (i) the symbolic AI which tries to formulate the world into a symbolic system and uses logic algorithms to reason about it; (ii) the statistical AI which tries to establish a mathematical function to formulate the relationship between input and output, yet the function can be approximated or even non-explainable. The past decade was dominated by the second path, in particular, the deep learning theory [21] which is part of the idea of the connectionist approach.

Although artificial general intelligence (AGI) is the ultimate goal of AI. The added word, ‘general’, implies that the key of AGI is to improve the generalization ability of AI algorithms. Conceptually, AGI can be defined as a system that solves any task that human beings or animals can perform [32]². In the modern era, there are a series of thoughts about AGI, resulting in verbal, psychological, and

of course AI-based definitions of AGI, many of which were summarized in an early paper [33], including:

- In [32], [34], the authors assumed that an AGI algorithm can do any task that humans or intelligent animals can do. This description is direct and anthropocentric, but it ignores the possibility that AGI can surpass real-world creatures, possibly by consuming more energy.
- In [35], [36], the authors asked that AGI algorithms can apply to as many tasks and scenarios as possible. However, without any constraints, the definition seems difficult to distinguish an AGI algorithm from a set of individual algorithms designed for specific purposes.
- In [37], the authors described typical characteristics of AGI algorithms, including being symbolic, emergentist, hybrid, and universalist.

Despite the vast argument in the description of AGI, one conclusion is clear: human intelligence is multi-faceted and thus it is difficult to use one definition to cover all properties of AGI.

In the AI field, probably one of the most famous thought experiments is the Turing test [38] which claimed that a machine is considered to gain intelligence if a human evaluator cannot tell the machine from the human in text-only communications. After being pursued by researchers for decades, the Turing test has become part of AI culture, although there exist challenges to it, *e.g.*, the Chinese room argument [39] which argued that AI algorithms might pass the Turing test without understanding what they are doing.

As far as we know, no AI algorithms have seriously passed the Turing test, because all of them exhibit clear patterns which make them easy to be discriminated from humans. This also includes the recently developed AI chatbots like LaMDA [40] and the GPT series [1]: they have shown strong abilities in chat and/or problem-solving, and some sources even advocated for them to pass the Turing test, but, for professional evaluators who are familiar with AI, they are still quite easy to be identified, not to mention that these chatbots are known to ‘hallucination’ [41] and humans often do not. This is an interesting signal that useful AGI systems may not necessarily mimic human behaviors.

Going beyond text-only systems, there are many more data modalities (*e.g.*, speech, image, video, *etc.*) to be processed. To integrate them into one system, we follow [42], [43] to define the goal of AGI to be maximizing reward in an environment. Let there be an environment and an agent (the AGI algorithm) that can interact with it. The agent observes a sequence of states, $\mathbb{S} = \{s_1, \dots, s_T\}$, and can choose from a set of actions, $\mathcal{A} = \{a_1, \dots, a_M\}$, to perform. There are two functions that define the transition between states and the obtained rewards, respectively. The goal of AGI is to learn a policy, denoted as $\pi : \mathbb{S} \mapsto \mathcal{A}$, which maximizes the expected cumulative reward $R = \sum_{t=1}^T r(s_t, a_t)$. When we set s_t and a_t to be different data modalities, it is the above formulation can cover a wide range of AI tasks. Specifically, the currently popular proxy tasks in computer vision such as image classification, object detection and segmentation, *etc.*, are mostly weakened versions of the above formulation

². Throughout this paper, we limit the concept of AGI within the scope of problem-solving, and thus we will not talk about programs that exhibit sentience or consciousness.

where the episode length T equals to 1, *i.e.*, these tasks are not built upon interaction with some environments.

In brief, the AGI is to learn a generalized function $\mathbf{a} = \pi(\mathbf{s})$. Although the form is quite simple, it was very difficult for the old-fashioned AI algorithms to use the same methodology, algorithm, or even model to deal with them all. In the past decade, deep learning [21] offers an effective and unified methodology: one can train a deep neural network to approximate the function $\mathbf{a} = \pi(\mathbf{s})$ without knowing about the actual relationship between them. The emergence of powerful neural architectures such as the transformer [44] even enables the researcher to train one model for different data modalities [45].

There are enormous difficulties in achieving AGI, including but not limited to the following issues:

- **The complexity of data.** Real-world data is multi-

faceted and rich. Some data modalities (*e.g.*, images) can have quite a high dimensionality and the relationship between different modalities can be complex and latent.

- **The complexity of human intelligence.** The goal of AGI is not only about problem-solving but also about planning, reasoning, reacting to different events, *etc.* Sometimes, the relationship between human behavior and the target is obscure and hard to represent in math forms.
- **Lack of neurological or cognitive theory.** Humans do not yet understand how human intelligence is achieved. Currently, computer algorithms provide one pathway, yet more possibilities may arise with future research in neurology and/or cognition.

U Read the following article^[A] and answer questions.

We've created GPT-4, the latest milestone in OpenAI's effort in scaling up deep learning. GPT-4 is a large multimodal model (accepting image and text inputs, emitting text outputs) that, while less capable than humans in many real-world scenarios, exhibits human-level performance on various professional and academic benchmarks. For example, it passes a simulated bar exam with a score around the top 10% of test takers; in contrast, GPT-3.5's score was around the bottom 10%. We've spent 6 months iteratively aligning GPT-4 using lessons from our adversarial testing program as well as ChatGPT, resulting in our best-ever results (though far from perfect) on factuality, steerability, and refusing to go outside of guardrails.

Over the past two years, we rebuilt our entire deep learning stack and, together with Azure, co-designed a supercomputer from the ground up for our workload. A year ago, we trained GPT-3.5 as a first "test run" of the system. We found and fixed some bugs and improved our theoretical foundations. As a result, our GPT-4 training run was (for us at least!) unprecedentedly stable, becoming our first large model whose training performance we were able to accurately predict ahead of time. As we continue to focus on reliable scaling, we aim to hone our methodology to help us predict and prepare for future capabilities increasingly far in advance—something we view as critical for safety.

Q: How does GPT-4's performance compare to humans in real-world scenarios?

A: GPT-4 is less capable than humans in many real-world scenarios.

U Please translate the article to Chinese.

我们创造了GPT-4，这是OpenAI在深度学习扩展方面的最新里程碑。GPT-4是一个大型多模态模型（接受图像和文本输入，输出文本），在许多实际场景中虽然比人类能力稍逊，但在各种专业和学术基准测试中展现出与人类水平相当的表现。例如，它在模拟的律师资格考试中得分大约在前10%的考生范围内；相比之下，GPT-3.5的得分大约在后10%。我们花费了6个月的时间，通过对抗测试计划和ChatGPT中吸取教训，不断调整GPT-4，从而取得了迄今为止（尽管远非完美）在事实准确性、可控性和拒绝超越限制方面的最佳结果。

在过去的两年里，我们重新构建了整个深度学习系统，并与Azure共同设计了一台从头开始适用于我们工作负载的超级计算机。一年前，我们训练了GPT-3.5作为这个系统的首次“试运行”。我们发现并修复了一些错误，并改进了我们的理论基础。因此，我们对GPT-4的训练过程进行了（至少对我们来说）前所未有的稳定性，这是我们能够提前准确预测的第一个大型模型的训练性能。随着我们继续专注于可靠的扩展，我们的目标是完善我们的方法，帮助我们越来越提前预测和准备未来的能力，这对于安全至关重要。

U Please extract named entities in the article.

Named entities in the article are as follows:
GPT-4, OpenAI, Azure, GPT-3.5, ChatGPT, ...

Fig. 1. An example of using GPT-4 for question answering, English-to-Chinese translation, and named entity extraction. The English article was borrowed from <https://openai.com/research/gpt-4>.

U If you sum up 1 to 100 but exclude all primes, what is the result?

To solve this problem, we should sum all numbers from 1 to 100 and then subtract the sum of all prime numbers within that range.

Firstly, the sum of all integers from 1 to 100 can be calculated using the formula for the sum of an arithmetic series:

$$n*(n + 1)/2,$$

where n is 100 in this case:

$$100*(100 + 1)/2 = 5050.$$

Then we calculate the sum of all prime numbers from 1 to 100. The prime numbers between 1 and 100 are 2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31, 37, 41, 43, 47, 53, 59, 61, 67, 71, 73, 79, 83, 89, 97. If we sum these up:

$$2 + 3 + 5 + 7 + 11 + 13 + 17 + 19 + 23 + 29 + 31 + 37 + 41 + 43 + 47 + 53 + 59 + 61 + 67 + 71 + 73 + 79 + 83 + 89 + 97 = 1060.$$

Then subtract this from the total sum:

$$5050 - 1060 = 3990.$$

So the sum of all numbers from 1 to 100 excluding all prime numbers is 3990.

U There are three people (Alex, Ben and Cody), one of whom is a knight, one a knave, and one a spy. The knight always tells the truth, the knave always lies, and the spy can either lie or tell the truth. Alex says: "Cody is a knave." Ben says: "Alex is a knight." Cody says: "I am the spy." Who is the knight, who the knave, and who the spy?

Let's try to solve this by examining each statement:

- Alex says: "Cody is a knave."
- Ben says: "Alex is a knight."
- Cody says: "I am the spy."

Since the knight always tells the truth, and the knave always lies, Cody can't be the knave. If he were, his statement would be a paradox. He would be lying about being the spy, which a knave couldn't truthfully declare. Therefore, Cody must be the spy.

This makes Alex's statement false (as Cody is not a knave), meaning Alex can't be the knight (who always tells the truth). Thus, Alex must be the knave.

By process of elimination, Ben, who hasn't been assigned a role yet, must be the knight. His statement is also true: "Alex is a knight." Although it seems to conflict with our conclusion that Alex is the knave, in fact, it doesn't, as this could be a transcription error. Perhaps, Ben really said: "Alex is a knave."

So, in summary:

- Alex is the Knave.
- Ben is the Knight.
- Cody is the Spy.

Note: the final answer was wrong although the process seems fine. LLM cannot guarantee to solve complex logic.

Fig. 2. An example of using GPT-4 for solving mathematical and logical problems. The final answer of the second problem which involves complex logical reasoning was wrong. The answer even contains self-contradiction, indicating that language models may hallucinate.

U The following appeared as part of an article in a business magazine.

"A recent study rating 300 male and female Mentian advertising executives according to the average number of hours they sleep per night showed an association between the amount of sleep the executives need and the success of their firms. Of the advertising firms studied, those whose executives reported needing no more than 6 hours of sleep per night had higher profit margins and faster growth. These results suggest that if a business wants to prosper, it should hire only people who need less than 6 hours of sleep per night."

Write a response in which you examine the stated and/or unstated assumptions of the argument. Be sure to explain how the argument depends on these assumptions and what the implications are for the argument if the assumptions prove unwarranted.

S The argument presented in the article suggests a direct correlation between the amount of sleep executives need and the success of their firms, specifically stating that those who need less than 6 hours of sleep per night contribute to higher profit margins and faster growth. While this conclusion might appear compelling at first glance, it is based on several assumptions that, if proven unwarranted, could significantly weaken the argument.

Firstly, the argument assumes causation based on a mere correlation. Just because firms with executives sleeping less than 6 hours have higher profits and faster growth doesn't necessarily mean that the reduced sleep hours caused these successes. There could be other contributing factors that the study didn't take into account, such as the overall work ethic, industry experience, business strategies, or the efficiency of the workforce at these firms.

Secondly, the argument presumes that the self-reported sleep hours of executives are accurate. Self-reported data can often be unreliable due to biased recall or the desire to present oneself in a certain way. Some executives might underreport their sleep hours due to a perception that successful people sleep less.

Thirdly, the argument implies that hiring people who sleep less than 6 hours per night will lead to increased prosperity for a business. This disregards research demonstrating the negative impact of chronic sleep deprivation on cognitive functions, decision-making, and general health. Over time, the productivity of these sleep-deprived individuals might decline, potentially harming the business.

Lastly, it assumes that all roles within a company will benefit from lesser sleep, which is a gross oversimplification. Different roles require different skill sets and may not benefit from reduced sleep in the same way an executive's role might.

In conclusion, while the argument is intriguing, its assumptions make it flawed. Correlation does not imply causation, and there might be other unconsidered factors contributing to a firm's success. Furthermore, the assumption that hiring people who sleep less will lead to greater prosperity is a simplistic approach and ignores potential long-term health and productivity issues associated with sleep deprivation.

U His calm exterior _____ the militant fervor with which he threw himself into the campaign.

delineated/reviled/falsified/ridiculed/belied

S The correct word to fill in the blank is "belied."

The sentence then reads: "His calm exterior belied the militant fervor with which he threw himself into the campaign."

In this context, "belied" means "gave a false impression of," which fits the sentence because the calm exterior gives a false impression of the militant fervor with which he engaged in the campaign.

U Inured to the intense work engendered by the deadlines they normally faced, the production managers felt somewhat _____ by the temporary hiatus in orders.

constrained/disoriented/wizened/overworked/dehumanized

S The correct word to fill in the blank is "disoriented."

The sentence then reads: "Inured to the intense work engendered by the deadlines they normally faced, the production managers felt somewhat disoriented by the temporary hiatus in orders."

In this context, "disoriented" means "confused" or "thrown off," which fits because the production managers, used to intense work, might feel confused or thrown off by a sudden break in orders.

U Quantity A
A diagonal of a cube with edge 5

Quantity B
A diagonal of a rectangular prism with sides 3, 4, and 5

- Quantity A is greater.
- Quantity B is greater.
- The two quantities are equal.
- The relationship cannot be determined from the information given.

S To solve this, we need to use the formula for the diagonal of a three-dimensional figure, which is $\sqrt{l^2 + w^2 + h^2}$, where l, w, and h are the lengths of the sides.

Quantity A is the diagonal of a cube with edge 5. Since a cube has all sides equal, $l = w = h = 5$. Therefore, the diagonal of the cube (Quantity A) is $\sqrt{5^2 + 5^2 + 5^2} = \sqrt{75} = 5\sqrt{3}$.

Quantity B is the diagonal of a rectangular prism with sides 3, 4, and 5. So, $l = 3, w = 4, h = 5$. Therefore, the diagonal of the prism (Quantity B) is $\sqrt{3^2 + 4^2 + 5^2} = \sqrt{50} = 5\sqrt{2}$.

Since $\sqrt{3} > \sqrt{2}$, Quantity A is greater than Quantity B. So, the correct answer is "Quantity A is greater."

Fig. 3. An example of using GPT-4 for the GRE test, including writing an essay and answering verbal and math questions. All the displayed objective questions were correctly answered. The problems are borrowed from <https://gre.kmf.com/exam/pre/817>.

3 GPT: SPARK OF AGI IN NLP

In the past year, ChatGPT³, GPT-4 [1], and other AI chatbots such as Vicuna⁴, made large progress towards AGI. They are computer algorithms developed for natural language processing (NLP). With a chat procedure with humans, they can understand the intention of humans and accomplish a wide range of tasks as long as they can be presented in pure texts. In particular, GPT-4 has a strong ability in generic problem-solving and was considered an early spark of AGI in the NLP field [2].

We briefly showcase the pure-text abilities of GPT-4. Throughout this part, we have used the May 12th version of GPT-4. The set of covered tasks includes the conventional NLP problems (e.g., translation, named entity recognition, question answering, etc.), as shown in Figure 1) and other

text-based problems such as solving mathematical and logical problems (Figure 2), passing verbal exams (e.g., GRE, as shown in Figure 3), coding with debugging (Figure 4), and so on. Beyond these basic examples, GPT-4 also exhibits a strong ability in logic, which enables it to integrate clues collected from multiple rounds of dialog into the final answer (Figure 5). We refer the readers to a previous paper (i.e., the Sparks-of-AGI paper [2]) for a thorough analysis of the ability of GPT-4.

Although GPT-4 has not yet opened the vision interface to the public, the official technical report [1] showed several fancy examples about multimodal dialog, i.e., chat based on an input image as reference. This implies that GPT-4 has been equipped with abilities of aligning language features with visual features, hence it can perform basic visual understanding tasks. As we shall see later (in Section 4.2.5), the vision community has developed several replacements [20], [46] for the same purpose, and the key lies in using ChatGPT

3. <https://openai.com/blog/chatgpt>

4. <https://github.com/lm-sys/FastChat>



Fig. 4. An example of using GPT-4 for writing a reinforcement learning program to play an Atari 2600 game, *SpaceInvaders*. Most of the generated contents (code and text) are eliminated to save space. Please note how GPT-4 corrected the code based on the user's feedback.

or GPT-4 to generate (instruct) training data. Additionally, with simple prompts, GPT-4 is also capable of calling external software (e.g., Midjourney, as shown in Figure 6) for image generation and external libraries (e.g., the HuggingFace libraries, as shown in [19]) for solving complex problems in computer vision.

These AI chatbots were trained in a two-stage procedure. In the first stage, a large language model (LLM), most of which are based on the transformer architecture [44], is pre-trained on a large-scale text database with self-supervised learning [3], [47], [48]. In the second stage, the pre-trained

LLM is supervised by human instructions [4] to accomplish specific tasks. If necessary, human feedback is collected and reinforcement learning is performed [49] to fine-tune the LLM towards better performance and higher data efficiency.

Later in Section 4.3, we will revisit the above procedure and understand it as a natural choice for training an agent to interact with the text environment.

U I have a 7-digit integer but you do not know it. You can get it by asking me questions. Every time, you can ask for only one digit (e.g. what is the n-th digit) and I will give you the answer. The only constraint is that you cannot ask for neighboring digits in two consecutive questions. Do you understand what I said?

U 5

U Yes, I understand. You have a 7-digit integer and I can ask for one digit at a time, but I cannot ask for neighboring digits in two consecutive questions. So, for instance, if I asked for the 4th digit in my first question, I couldn't ask for the 3rd or 5th digit in my next question. I would have to ask for the 2nd, 6th, or 7th digit next. I can continue this process until I have discovered all seven digits.

U Thank you. The 5th digit is 5. Now, I'll ask for the 1st digit. What is the 1st digit of your 7-digit integer?

U Now, let us begin.

U 1

U Alright, let's start with the first question: What is the 4th digit of your 7-digit integer?

U Thank you. The 1st digit is 1. Now, I'll ask for the 3rd digit. What is the 3rd digit of your 7-digit integer?

U 8

U 4

U Great, the 4th digit is 8. For the next question, I'll ask for the 2nd digit. What is the 2nd digit of your 7-digit integer?

U Thank you. The 3rd digit is 4. Now, I'll ask for the 6th digit. What is the 6th digit of your 7-digit integer?

U 0

U 7

U Thank you. The 6th digit is 7. Finally, I'll ask for the last digit. What is the 7th digit of your 7-digit integer?

U 6

U Thank you. The 7th digit is 6. So, based on your answers, your 7-digit integer is 1048576.

Note: GPT broke the rule by asking for the 6th and 7th digits in the last rounds.

Fig. 5. An example of GPT-4 interacting with human to get the final answer. GPT-4 understood the logic although it broke the rule in the final step.

U I want to create a picture on midjourney. It shall contain a beautiful scene with sea, blue sky, cloud, and some people spending vacation. Show me some prompts that I can use.

U I used the first prompt and created some pictures. They are good, but all of them seem like oil paintings. I want to create some pictures that are close to real photos. How can I modify the prompt?

U Here are several prompts that can help you create your vision:

U I'm glad to hear you created some pictures using the prompts. If you're aiming for a more realistic, photographic style, it would be useful to incorporate more precise and detailed descriptions in your prompts. Here's a revised version of the first prompt:

1. A Bright Summer Day: Imagine a bright summer day with the sun beaming in a clear blue sky. Wisps of soft, white clouds lazily drift by. A sparkling azure sea stretches out to the horizon, mirroring the sky above. A few people are scattered along the beach, lounging under colorful umbrellas, playing in the sand, or splashing in the gentle waves. Their laughter and the sound of the sea create a cheerful symphony.
2. Sunset Vista:
3. Tropical Paradise:
4. Coastal Adventure:
5. Tranquil Morning:

A Photorealistic Summer Day: Visualize a photorealistic image of a midsummer's day at the beach. The sky is a clear and vibrant blue, with an unblemished expanse that is punctuated by only a few wisps of delicate, white clouds. The sea, reflecting the azure of the sky, shimmers under the brilliant sun. Minute details such as the froth at the crest of each wave, the sand's texture under the light, and even the varied colors and patterns of beach towels and umbrellas dotting the landscape are clearly visible. Vacationers are scattered along the beach, with intricate details like the patterns on their swimwear, the glistening beads of water on their skin after a swim, and the joy reflected in their expressions making them seem almost life-like. The entire scene captures the essence of a perfect beach vacation, combining it with the realism and detail typical of high-resolution photography.

Using the first prompt on midjourney.

Using the updated prompt on midjourney.

Use these prompts to inspire your creativity and come up with a scene that best captures your vision of "midjourney".

By emphasizing detailed and realistic elements in your prompt, you're encouraging a focus on fine elements and minutiae, which could aid in creating a more photograph-like picture.

Fig. 6. An example of using GPT-4 for automatic prompts for text-to-image generation with Midjourney (<https://www.midjourney.com/>). GPT-4 understood the user's intention to adjust the prompt, although the new prompt still cannot fully satisfy the user's requirements.

4 CV: THE NEXT BATTLEFIELD OF AGI

Humans perceive the world based on multiple data modalities. It is a common knowledge that about 85% of what we learn is through our vision system. Therefore, given that the NLP community has shown the promise of AGI, it is natural to consider computer vision (CV) or multimodality (which includes at least the vision and language domains) as the next battlefield of AGI.

Here we provide two additional comments to complement the above statement. **First**, it is clear that CV is a superset of NLP, because humans read articles by first recognizing characters in the captured images and then understanding the contents. In other words, an AGI in CV (or multimodality) should cover all abilities of an AGI in NLP. **Second**, we argue that language alone is insufficient in many scenarios. For example, when one tries to find detailed information about an unknown object (*e.g.*, animal, fashion, *etc.*), the best way is to capture an image and use it for online search; purely relying on text descriptions can introduce uncertainty and inaccuracy. As another case, as we shall see in Section 4.3, it is not always easy to refer to fine-grained semantics in a scene (for recognition or image editing), and it is more efficient to think in a vision-friendly manner, *e.g.*, using a point or box to locate an object rather than saying something like ‘the person who is wearing black jacket, standing in front of the yellow car, and talking to another person’.

4.1 Ideal and Reality

We desire a CV algorithm that can solve generic tasks, possibly by interacting with the environment. Note that the requirement is not limited to recognizing everything or performing dialog based on an image or video clip. It shall be a holistic system that receives generic orders from humans and produces the desired results. But, the current status of CV is quite preliminary. As shown in Figure 7, the CV community has been using different modules and even systems for different vision tasks. Below, we list a few of them.

- **Image classification** is one of the most fundamental tasks in CV, due to the simplicity of the setting and the cheapness of collecting training data. State-of-the-art image classification algorithms are based on deep neural networks including convolutional networks [8], [51] and vision transformers [9], [10]. A pre-training stage with either self-supervised representation learning [25], [52] or large-scale datasets (*e.g.*, the full ImageNet [53] or even external datasets [54], [55]) is very helpful to improve the classification accuracy.
- The models for **object detection and instance segmentation** are mostly fine-tuned from the models trained for image classification. Researchers designed specific modules (often referred to as the head) to use the image features extracted by the classification network (often referred to as the backbone) for object localization and recognition. The head modules can be roughly categorized into the two-stage [56], [57], [58] and one-stage [59], [60] methods,

and the transformer blocks have been used [61] and pushed the performance on real-world data [62] towards a higher level [6], [7].

- The **semantic segmentation** algorithms fine-tune models trained for image classification in another way. The early efforts involve the encoder-decoder architecture which first downsamples the original image to extract semantic features and then upsamples the features to the original resolution [63], [64]. The idea was also inherited to medical images [65] and generalized to 3D data [66]. It was shown that keeping high-resolution features improves the segmentation accuracy [67]. Vision transformers also offered new opportunities for more accurate segmentation models [68], [69], especially for more challenging datasets [70].
- The **image captioning** task [62], [71] is one of the early trials for cross-modal understanding. In the beginning, pre-trained vision models are equipped with a recurrent head such as LSTM [72] for generating captions [73], [74]. Recently, researchers developed an alternative solution for image captioning which involves fine-tuning foundation models that have connected vision to language [75], [76], [77].
- For **text-to-image generation**, state-of-the-art algorithms [12], [78] are based on the alignment between vision and language. For this purpose, a cross-modal pre-trained model such as CLIP [13] is inherited, based on which probabilistic models are used to decode sequential tokens into images [79], [80], [81] or denoising latent diffusion models [12], [78], [82].

Besides, there exist algorithms for other vision tasks, including multiple object tracking [83], [84], [85], pose estimation [86], and many others. It is clear that the current status of CV (individual algorithms are used for different purposes) is far from what the GPT series has achieved in the NLP field.

4.2 Unification Is the Trend

Below, we summarize recent research topics towards unification in CV into five categories.

4.2.1 Open-world Visual Recognition

In a long period of time, most CV algorithms can only recognize the concepts that appear in the training data, leading to a ‘closed-world’ of visual concepts. In opposite, the concept of ‘open-world’ refers to the ability that a CV algorithm can recognize or understand any concept regardless whether it has appeared before. The open-world ability⁵ is often introduced by natural language since it is a natural way for humans to understand new concepts. This explains why language-related tasks such as image captioning [73], [74] and visual question answering [91], [92], [93] contributed to the earliest open-world settings for visual recognition.

Recently, with the emergence of vision-language pre-training (*e.g.*, CLIP [13] and ALIGN [94]), it becomes much

5. Sometimes, ‘open-world’ is referred to as ‘open-set’ or ‘open-domain’, although these terminologies may have slightly different meanings.

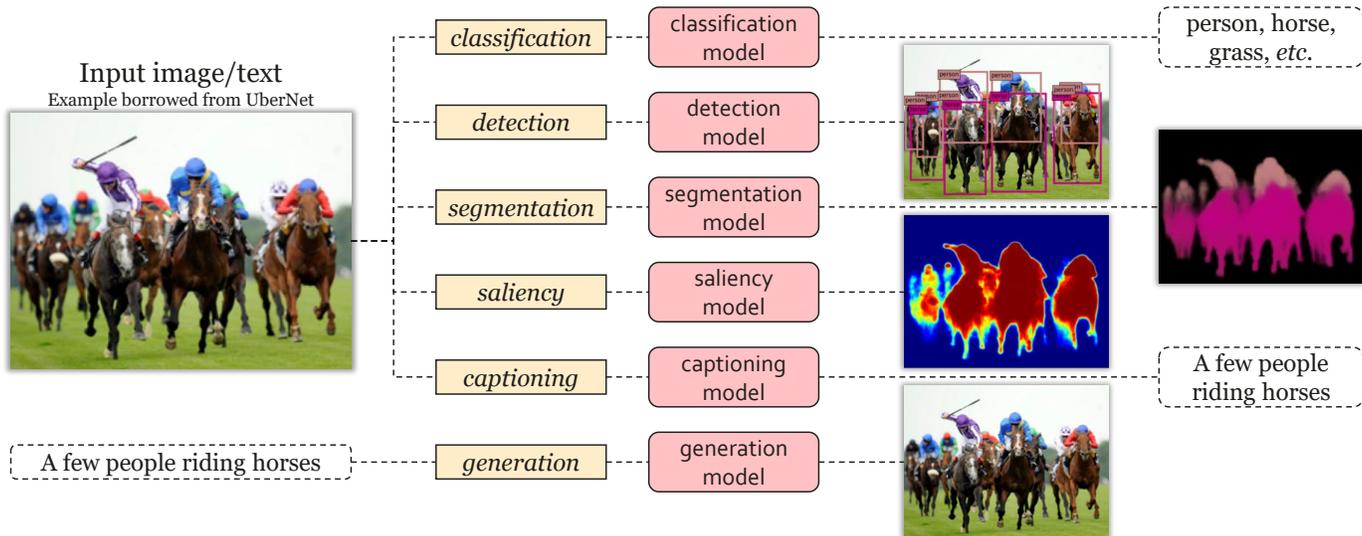


Fig. 7. The current status of computer vision. Different problems are solved by different models or algorithms. Image credit: UberNet [50].

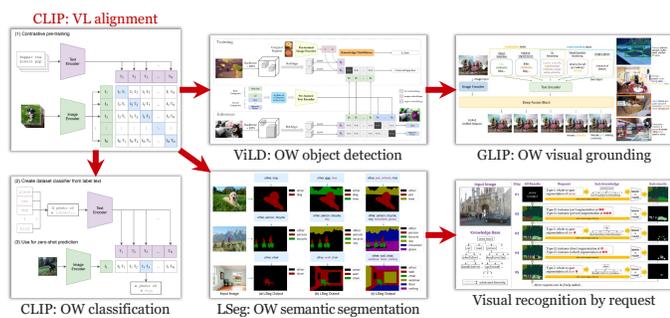


Fig. 8. Vision-language (VL) pre-training enables open-world (OW) recognition in many vision tasks. Image credit: CLIP [13], ViLD [87], LSeg [88], GLIP [89], ViRRReq [90].

easier to align features in the vision and language domains. The unified feature space not only offers simpler pipelines for image captioning [75], [76], [77] and visual question answering [11], [76], [95], but also creates a new methodology [13] for conventional visual recognition tasks. For example, image classification can be done by simply matching the query image with a set of templates (also known as ‘prompts’) saying a photo of {something}, where something can be any (hence open-world) concept like cat or Siberian husky, and set the result to be the candidate with the highest matching score. Beyond the vanilla version, researchers developed algorithms [96], [97] named ‘learning to prompt’ to improve the classification accuracy. Later, the methodology was inherited from image classification to object detection [87], [98], semantic segmentation [88], [99], instance segmentation [100], panoptic segmentation [101], [102], and further extended to visual grounding [103] and composite visual recognition [90] tasks. These tasks can benefit from vision-language models pre-trained with enhanced localization [103], [104].

Open-world visual recognition is closely related to zero-shot visual recognition because both of them try to generalize the recognition ability to the concepts that have

not appeared in the training set. However, in the authors’ opinion, it is yet unclear whether and how deep learning algorithms can recognize unseen concepts. Indeed, there are some special cases that zero-shot recognition can be achieved (*e.g.*, the training data contains dog, cat, and dog’s head, but it does not contain cat’s head; it is possible that the algorithm can learn the concept of cat’s head from composition without training data), but in most cases, the zero-shot ability was inherited from the pre-trained vision language model. Note that the original CLIP model and other variants (*e.g.*, OpenCLIP [105] and EVA-CLIP [106]) were pre-trained on large-scale image-text pairs which may have contained the target concepts withheld from the downstream training set. Therefore, we argue that ‘open-world’ is a more precise description than ‘zero-shot’.

As language introduces flexibility to visual recognition, it also brings the drawback of referring to detailed semantics in complex scenes. For example, when a large number of same-class objects appear in an image, it is difficult for the model to ask about the position, shape, or attributes of a specific object. This issue is easily solved in vision itself, *e.g.*, one can use a point to indicate the object of interest. We will get back to this issue in the part discussing multimodal dialog (see Section 4.2.5).

4.2.2 The Segment Anything Task

The *Segment Anything* task [14] was introduced recently as a generalized module to cluster raw image pixels into groups, many of which correspond to basic visual units in the image. The proposed task supports several types of prompts including point, contour, text, *etc.*, and produces a few masks as well as scores for each prompt or each combination of prompts. Trained on a large-scale dataset with about 10 million images, the derived model, SAM, was able to transfer to a wide range of segmentation tasks including medical image analysis [111], [112], [113], camouflaged object segmentation [110], [114], 3D object segmentation [108], object tracking [115], as well as application scenarios such as image inpainting [109]. SAM can also be used with state-

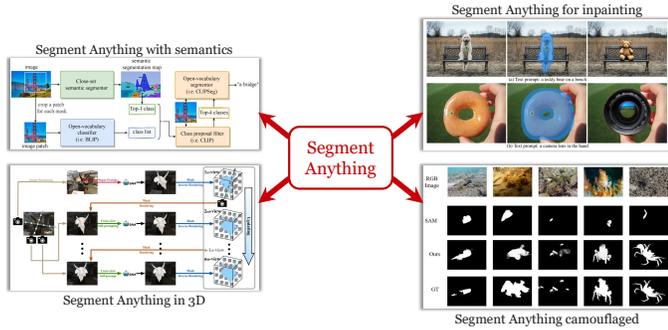


Fig. 9. The pre-trained SAM [14] is easily transferred for various downstream vision tasks. Image credit: SSA [107], SA3D [108], inpaint anything [109], SAM-Adapter [110].

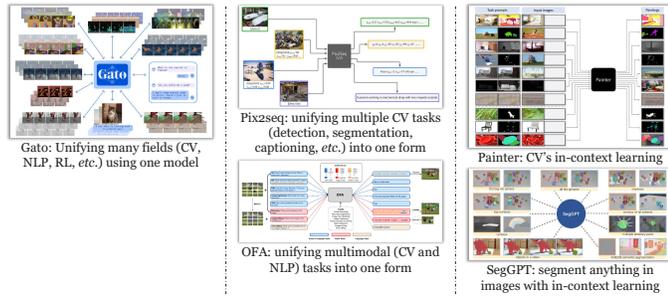


Fig. 10. Generalized visual encoding allows one model to be trained for various visual and/or multimodal understanding tasks. Image credit: Gato [45], pix2seq-v2 [118], OFA [16], Painter [17], SegGPT [119].

of-the-art visual recognition algorithms, such as refining bounding boxes produced by visual grounding [116] algorithms into masks, and feeding the segmented units into open-set classification algorithms for image tagging [107], [117].

Technically, the keys of SAM lie in the prompting mechanism and data closure, *i.e.*, closing the segmentation task with a small amount of feedback from labelers. The unified form of prompts makes SAM look like a part of the vision foundation model or pipeline, but there are still many unsolved issues. For example, it remains unclear about the upstream and downstream modules of SAM (if SAM is indeed part of the pipeline), and SAM can be severely impacted by pixel-level appearance, *e.g.*, an arm can be segmented from the torso exactly on the border of clothes, implying that color is the dominant factor for segmentation. In general, it is likely that SAM has over-fitted to the *Segment Anything* task itself and hence weakened its ability of classification.

4.2.3 Generalized Visual Encoding

Another way to unify CV tasks is to provide a generalized visual encoding for them. There are several methodologies to achieve this goal.

A key difficulty lies in the large variance between vision tasks, *e.g.*, object detection requires a set of bounding boxes while semantic segmentation requires a dense prediction over the entire image, both of which are very different from a single label required by image classification. As all can understand it, natural language offers a unified form to represent everything. An early effort named pix2seq [15]

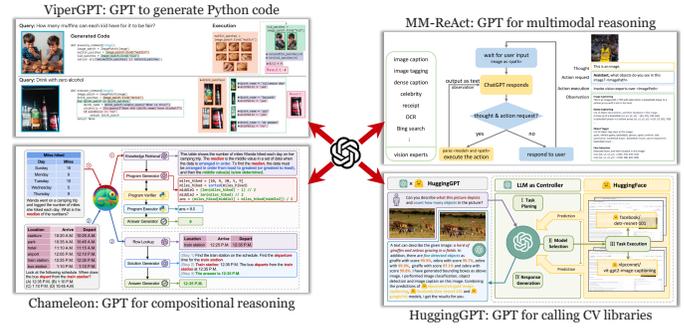


Fig. 11. GPT offers an easy and unified way to generate code or explanations for visual understanding. Image credit: ViperGPT [18], Chameleon [124], MM-ReAct [125], HuggingGPT [19].

showed that object detection results (*i.e.*, bounding boxes) can be formulated into natural language and coordinates and then converted into tokens as the output of vision models. In a later version, pix2seq-v2, they generalized the representation to unify the output of object detection, instance segmentation, keypoint detection, and image captioning. Similar ideas were also used for other image recognition [120], video recognition [121], and multimodal understanding [16], [122], [123] tasks.

Besides using language, researchers also tried to use vision alone to unify everything. The idea was named in-context learning and was borrowed from the NLP community [3], suggesting that a pre-trained model can realize the intention of new tasks from a few demonstrations. This learning paradigm was first introduced into CV using natural language as prompts [76]. In [17], different vision tasks, including instance segmentation, keypoint detection, depth estimation, saliency detection, *etc.*, were formulated into assigning different color patches or regions in the output image canvas, hence a single model named *Painter* can be trained to deal with them all. The framework was then extended into a more generalized form which also supports video segmentation [119].

In the backbone of the above algorithms lies the vision transformer [9], which offers strong data fitting ability in different modalities. The ability was verified by an earlier work which trained a generalist agent named *Gato* [45] to unify vision, language, and robotics tasks as long as the desired output can be encoded into a sequence of tokens.

Despite the ability of unified representation, it is questionable how far the methodology has gone beyond multi-task visual representation learning, where different tasks are integrated by incorporating multiple loss functions [50]. Recall that GPT applied in-context learning to unify NLP tasks, but CV does not necessarily follow the same direction: this is because CV tasks are mostly discrete (*e.g.*, there is no intermediate task between segmentation and tracking) and thus there might not be a significant difference between individual and joint optimization strategies.

4.2.4 LLM-guided Visual Understanding

Visual recognition can be complex especially when it involves compositional concepts and/or relationships between visual instances. It is difficult for end-to-end models (vision-language pre-trained models for visual question

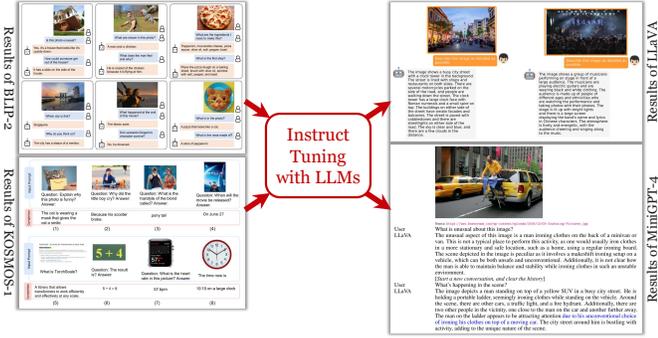


Fig. 12. With instruction tuning, visual question answering algorithms are extended to multimodal dialog systems. Image credit: BLIP-2 [11], KOSMOS-1 [95], LLaVA [20], MiniGPT-4 [46].

answering [11], [76], [95]) to produce answers following a procedure that is easily understood by humans.

To alleviate the issue, a practical methodology lies in generating explainable logic to assist visual recognition. The idea was not new. Several years ago, prior to the appearance of the transformer architecture, researchers proposed to use the long short-term memory (LSTM) model [72] to generate programs so that vision modules are invoked as modules for complex question answering [126]. At that time, the ability of LSTM largely limits the idea within the range of relatively simple and templated questions.

Recently, the appearance of large language models (especially the GPT series) makes the conversion of arbitrary questions possible. Specifically, GPT can interact with humans in different ways. For example, it can summarize basic recognition results to the final answer [125] or generate code [18], [124] or natural language scripts [19] to call basic vision modules. As a result, visual questions can be decomposed into basic modules. This is especially effective for logical questions, *e.g.*, that asking about the spatial relationship between objects or that depending on the number of objects.

LLMs may understand the logic, but they have not yet showed the ability to assist fundamental visual recognition modules. That said, the answer will still be wrong once the basic recognition results are incorrect, *e.g.*, the detection algorithm misses a few objects that are small and/or partially occluded. We expect an essential visual logic to be formulated in the future (*e.g.*, the algorithm can follow a sequential algorithm to detect every object, or be guided by commonsense [127] to solve hard cases), possibly with the assistance of LLMs, so that fundamental visual recognition is boosted.

4.2.5 Multimodal Dialog

Multimodal dialog extends text-based dialog to the vision domain. The early efforts involved visual question answering in which various datasets with simple questions have been constructed [128], [129], [130]. With the rapid development of LLMs, multi-round question answering was made available by fine-tuning pre-trained vision and lan-

guage models together [11], [95]⁶. It was also shown that a wide range of questions can be answered via in-context learning in multimodality [76] or using GPT as the logic controller [131].

Recently, a novel paradigm developed in the GPT series, named instruct learning [4], has been inherited to enhance the quality of multimodal dialog [20], [46]. The idea was to provide a few reference data (*e.g.*, objects, descriptions) from ground-truth annotation or recognition results and ask the GPT model to generate instruction data (*i.e.*, enriched question-answer pairs). Fine-tuned with these data (without reference), the foundation models for vision and language can interact with each other via a lightweight network module (*e.g.*, a Q-former [11]).

Multimodal dialog offers a preliminary interactive benchmark for computer vision, but, as a language-guided task, it also has the weaknesses analyzed in the open-world visual recognition (see Section 4.2.1). We expect that enriching the form of queries (*e.g.*, using generalized visual encoding methods, see Section 4.2.3) can push multimodal dialog to a higher level.

4.3 The Essential Difficulty

Indeed, the above efforts have largely advanced the progress of unification in CV. However, the community is still far from an algorithm that can solve a wide range of real-world tasks, in particular when interaction is needed. In this part, we discuss on the essential difficulty that leads to the current dilemma.

4.3.1 GPT Revisited

We start with recalling the definition of AGI (see Section 2, the definition was inherited from [42], [43]). In short, the goal of AGI is to maximize a reward in an interactable environment.

GPT defined such an environment with the chat task. Note that, in a pure-text world, chat is the perfect task for an agent to learn from interaction (talking with humans and receiving feedback); meanwhile, any task can be defined by chat. In our opinion, establishing the environment (with the chat task) is the most important insight of GPT. The technical solutions (*i.e.*, generative pre-training followed by instruct fine-tuning) can be derived from the chat task: generative pre-training is to memorize the distribution of the environment (world); instruct fine-tuning is to align the learned contribution with question-answer pairs for problem-solving. There are clear boundary between them, as the fine-tuning stage is driven by specific tasks while the pre-training stage is not.

We try to build the relationship between the basic elements of an environment and GPT. We find that the state, action, and reward in the environment correspond to the prompting mechanism, the desired target, and the feedback from users, respectively. We expect that CV algorithms are also trained in such an environment, or at least with the above factors clearly defined.

6. GPT-4 [1] also showed examples of multimodal dialog, but it is unclear how they achieved the ability, especially for the cases with rich texts, *e.g.*, solving a complex physical problem and understanding a joke which is mainly described in optical characters.

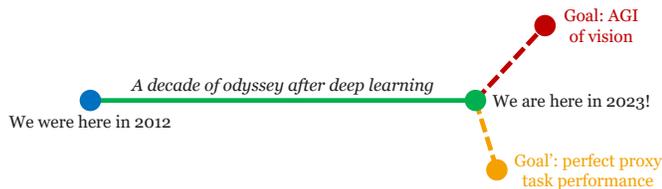


Fig. 13. An illustration of our opinion, *i.e.*, the proxy tasks in CV are dying.

4.3.2 Why Not Establishing Environments for CV?

Conceptually, an AGI in CV should also be trained in environments. Back to the 1970s, David Marr pointed out that CV algorithms should construct a world model and learn abilities by interacting with the model [132]. Other pioneers in AI, including Hans Moravec and Rodney Brooks, also emphasized the importance of learning from environments [133], [134]. However, establishing environments for CV is never an easy task, unlike that for NLP which is quite straightforward.

There are mainly two options, namely, training agents in the real world or in virtual, simulated environments. The former option is closer to the final objective, but the over-high costs and uncontrollable safety issues have constrained it in small-scale and toy scenarios (*e.g.*, training robotic arms for object grasping). The latter option is relatively easy to achieve, but it suffers the fidelity issues (not only about 3D modeling and rendering, but also about the behavior of other agents) and thus has to alleviate a significant domain gap when being transferred to the real world.

Due to the high difficulty of simulating the world (*i.e.*, establishing environments), the CV community has taken an alternative solution which is to sample the world. It involves two major steps, namely, image/video data collection and semantic annotation. The first step is to perform a sparse sampling of the real world – note that, although their size has significantly increased during the past decades, all existing datasets are still orders of magnitude smaller than the real world. The second step is to expect what agents need in order to accomplish tasks and convert the requirements (*e.g.*, detecting objects) into semantic annotations. From this point of view, we refer to them **proxy tasks** as they serve as surrogates to achieve the goal of AGI. Note that proxy tasks exist in almost all AI fields, *e.g.*, in NLP, there are various such tasks including translation, named entity recognition, and others.

4.3.3 Proxy Is Dying!

The proxy tasks have been vastly improved in the past decade, thanks to the rapid development of deep learning. One of the most well-known examples lies in ImageNet-1K classification [53], [135], where the best accuracy is under 50% prior to AlexNet [51], while the accuracy is higher than 90% as of today [55]. The odyssey was made possible by strong network architectures, effective optimization tricks, external training data, *etc.* Nevertheless, many research papers have been still pursuing higher accuracy on this dataset. Standing upon a high baseline (*e.g.*, 85%), the improvement brought by the proposed algorithm is

often small (*e.g.*, 0.5%), leading to a weird situation in that implementation tricks contribute even more than the proposed algorithm itself.

We illustrate the situation in Figure 13. Let us assume that AGI of vision and perfect performance of proxy tasks are two goals in the space of algorithms. In the pre-deep learning era, CV algorithms are mostly weak, so setting the goal to be high proxy task performance is reasonable. As of today, CV algorithms have been much stronger than before; consequently, continuing to improve proxy tasks can drive us away from AGI. We refer the readers to what happened in NLP: GPT offered a unified solution and largely reduced the importance of the conventional proxy tasks (*e.g.*, translation, named entity recognition, *etc.*)⁷.

5 FUTURE: LEARNING FROM ENVIRONMENT

The above analysis calls for a new paradigm for training strong agents for CV. In this section, we convert our opinions and insights into an imaginary pipeline, review existing works that are related to the pipeline, and make comments on future research directions based on the pipeline.

5.1 An Imaginary Pipeline

Figure 14 shows our idea. The pipeline comprises three stages: Stage 0 for establishing environments, Stage 1 for pre-training, and Stage 2 for fine-tuning. When necessary, the fine-tuned model can be prompted for conventional visual recognition tasks. Below, we describe each stage in detail.

- **Stage 0: establishing environments.** As analyzed previously, high-quality environments are strongly required for AGI in CV. Here the concept of ‘high-quality’ includes but is not limited to abundance (there should be sufficient and diversified environments), fidelity (visual appearance and other agents’ behavior should be close to the real world), and richness in interaction (the agent can be asked to perform a wide range of tasks by interacting with the environments).
- **Stage 1: generative pre-training.** The algorithm is asked to explore the environments and pre-trained to predict future frames. The biggest difference from the GPT task (predicting the next token) in NLP lies in that the future frames depend on the action of the agent (in NLP, the pre-trained text corpus remains unchanged), so the model is trying to learn a joint distribution of state and action. This strategy is especially useful when the set of established environments is insufficient to approximate the world’s distribution. Note that, as CV is a superset of NLP (see the paragraph before Section 4.1), the size (*e.g.*, number of parameters) of pre-trained CV models should be orders of magnitude larger than NLP models.

7. Disclaimer: these tasks are still important in some real-world applications, but it is unlikely that they shall be studied in an old-fashioned manner. With the chat task, these tasks can be accomplished with simple prompts.

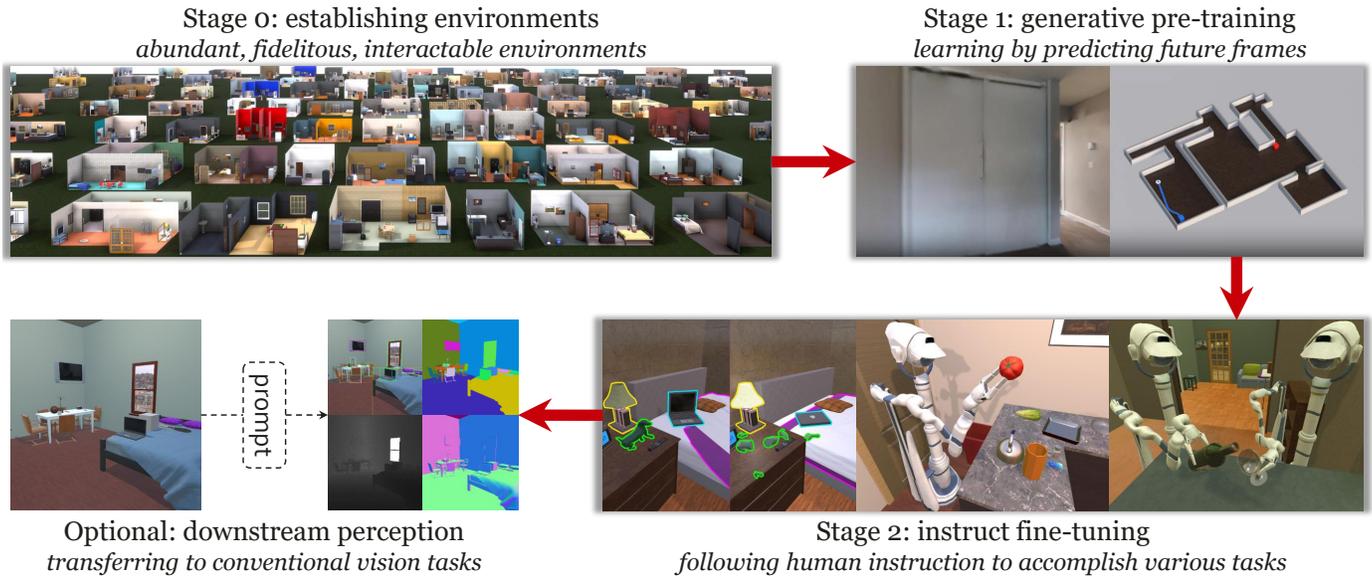


Fig. 14. An imaginary pipeline of training a stronger agent for computer vision. The idea was borrowed from GPT, with Stages 1 and 2 performing generative pre-training and instruct fine-tuning, respectively. Differently, environments need to be established prior to the pre-training stage, which itself is a major challenge. Image credit: the Stage 1 image is from Habitat [22] and others are from ProcTHOR [23].

- **Stage 2: instruct fine-tuning.** The pre-trained model is guided to accomplish real-world tasks, following human instructions. Intuitively, there are much more types of allowed interaction between the agent and environments, including exploration, navigation, using language, performing physical actions, and many others. A reasonable conjecture is that much more instruction data should be collected, which also corresponds to the size of foundation CV models.
- **Optional: downstream perception.** We expect that CV algorithms can learn all required abilities of perception from the previous stage, *e.g.*, in order to accomplish a very simple task, *Buy me a cup of coffee*, the model must at least learn to (i) explore around with safety, (ii) recognize where the coffee bar is, (iii) communicate with the shop assistant with language, and (iv) grasp the bought coffee. Such a model, when properly prompted, should output desired perception results, including tracking another agent (for not colliding with it), open-set visual recognition (for finding the bar and the bought coffee), and others. This is related to the idea of analysis by synthesis [136].

To sum up, we expect the agent to be task-oriented, *i.e.*, focusing itself on accomplishing tasks in the established environments. The proxy tasks (see Section 4.3.3) are to be solved naturally with prompts.

5.2 Existing Works

We briefly review the existing works that are related to the imaginary pipeline.

5.2.1 3D Embodied Environments

There are typically two options for establishing virtual environments. The first option is to collect visual content

from a real-world scenario and perform 3D reconstruction. For example, Habitat [22] released more than 200 scanned scenarios for visual navigation. The second option is to generate (render) scenes with 3D models. For example, ProcTHOR [23] provided a large set of 3D objects and a program to randomly generate room layouts, so that one can sample an arbitrary number of 3D scenes and perform various tasks including navigation, grasping, and reordering.

However, the existing environments (including the above two and many others [28], [137]) have not yet validated the ability to scale up to the level of the real world. In particular, when more and more environments are sampled from ProcTHOR [23], the performance, either for embodied or downstream recognition tasks, can saturate quickly. Clearly, for any existing method, there is a tradeoff between abundance and fidelity. This is a major challenge before pre-trained CV models can exhibit the scaling law [138] and emergent abilities [139] as NLP models did.

5.2.2 Pre-training Vision Models

In the past years, visual pre-training methods have been largely developed. Contrastive learning (CL) [24], [140], [141], [142] offered the first methodology to surpass supervised learning in downstream tasks, and masked image modeling (MIM) [25], [52], [143] pushed the performance of pre-trained models to a higher level. The major difference between them lies in the pre-training objective, where CL is discriminative and MIM is generative.

Built upon a generative objective, MIM is closer to what we desire in the aforementioned pipeline. However, predicting future frames in environments is different from predicting missing contents in sampled images or videos, which is similar to the difference between masked language modeling [48] and generative pre-training [47]. Additionally, compared to text data, there can be much more re-

dundant information in vision data. We conjecture that data compression is an important factor in the pre-training task.

5.2.3 Reinforcement Learning for Game Playing

Interacting with environments is closely related to game playing. The past decade has witnessed the integration of deep learning and reinforcement learning, resulting in a series of algorithms for playing various games, such as Atari 2600 games [26], [144], simulated robotics tasks [144], [145], Go and other chess [146], [147], StarCraft II [27], and so on. Effective algorithms were also designed for combining multiple reinforcement learning strategies [148] or completing very complex tasks [149].

Compared to the above problems, the real world is much more complicated and involves actions from different aspects including language and robotics. A practical strategy is to first disable some types of interaction to simplify the tasks and add them back when the foundation models are sufficiently strong.

5.2.4 Embodied Computer Vision

Embodied AI refers to the research field in that agents learn from/for interacting with virtual environments. The motivation partly came from how humans learn as babies [150]. In the scope of CV, typical tasks include exploration [151], [152], [153] where the goal is simply to explore and reconstruct the world, visual navigation [28], [154], [155] where the agents are asked to explore the world for specific targets (*e.g.*, an image or an object), and embodied question answering [29], [156] where the agents answer questions based on the interaction with the world.

Recently, two works were brought to our attention. The first one is PaLM-E [30] where a general-purpose vision-language model is trained to perform a wide range of embodied tasks. These tasks are organized into a unified prompting system, thanks to the introduction of LLMs. The second one is ENTL [31] where an end-to-end system was designed and different stages in embodied CV (including world modeling, localization, and imitation learning) were tokenized and integrated into a sequence prediction task. Both works pushed the unification in embodied CV forward from different directions.

We emphasize that, despite the existing efforts, the real-world scenario is much more complicated than what we have ever created or simulated (see Section 5.2.1). To achieve the goal of AGI, more interaction types should be supported, long-range and complex tasks should be designed, and instruction data should be collected. System design and engineering might be more important than one thinks.

5.3 Comments on Research Directions

As the final part, we make comments on future research directions. With the major goals migrated from the performance of proxy tasks to learning from environments, many popular research directions may have to adjust their goal. Here is a disclaimer: all the following statements are our personal opinions and may be wrong.

5.3.1 On Establishing Environments

A clear goal is to continue increasing the size, diversity, and fidelity of the virtual environments. There are multiple techniques that can help. For example, novel 3D representation forms (*e.g.*, neural rendering field, NeRF [157], [158], [159]) may be more efficient in achieving a tradeoff between the reconstruction quality and overhead.

Another important direction lies in the richness of environments. It is a non-trivial task to define new, complex tasks and unify them into a prompting system. Also, AI algorithms can benefit much from a better simulation of other agents' behavior [160] because it can largely improve the abundance of environments and hence the robustness of the trained algorithms.

5.3.2 On Generative Pre-training

There are mainly two factors that affect the pre-training stage, namely, neural architecture design and proxy task design. The latter is clearly more important and the former shall be built upon the latter.

As analyzed in Section 5.2.2, existing pre-training tasks, including contrastive learning and masked image modeling, shall be modified to allow for efficient explorations in virtual environments. We expect the newly designed proxy to focus on data compression, because there is much heavier redundancy in vision data than in language data. The new pre-training proxy defines the requirement of neural architectures, *e.g.*, to achieve a tradeoff between data compression and visual recognition, the designed architecture should be equipped with an ability to extract different levels (granularity) of visual features by request.

Additionally, cross-modality (*e.g.*, text-to-image) generation will become a direct metric to measure the performance of pre-training. When a unified tokenization method is available, it can be formulated into a multimodal version of the reconstruction loss.

5.3.3 On Instruct Fine-tuning

We have not yet entered the scope of defining tasks in the new paradigm. As real-world tasks can be very complicated, we conjecture that some elementary tasks (*e.g.*, exploration, fetching, interaction, *etc.*) can be defined and trained first, so that complex tasks can be decomposed into them. For this purpose, a unified prompting system should be designed and abundant human instructions should be collected. As a reasonable conjecture, the amount of instruction data can be orders of magnitude larger than what has been collected for training GPT and other chatbots.

This is a completely new story for CV. The road ahead is filled with unknown difficulties and uncertainty. We cannot see much at the current point, but clear paths will emerge in the future.

6 CONCLUSIONS

In this paper, we discuss how to advance CV algorithms towards AGI. We start with reviewing the current status and recent efforts of CV for unification, and then we inherit ideas and insights from NLP, especially the GPT series. Our conclusion is that CV lacks a paradigm that allows learning from environments, for which we propose an imaginary

pipeline. We expect that substantial technical evolution is required to bring the pipeline to truth.

ACKNOWLEDGMENTS

The authors would like to thank many colleagues and collaborating researchers for instructive discussions.

REFERENCES

- [1] OpenAI, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [2] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg *et al.*, “Sparks of artificial general intelligence: Early experiments with gpt-4,” *arXiv preprint arXiv:2303.12712*, 2023.
- [3] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” in *Advances in Neural Information Processing Systems*, 2020.
- [4] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, “Training language models to follow instructions with human feedback,” 2022.
- [5] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. Chi, Q. Le, and D. Zhou, “Chain-of-thought prompting elicits reasoning in large language models,” in *Advances in Neural Information Processing Systems*, 2022.
- [6] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum, “Dino: Detr with improved denoising anchor boxes for end-to-end object detection,” in *International Conference on Learning Representations*, 2023.
- [7] F. Li, H. Zhang, H. Xu, S. Liu, L. Zhang, L. M. Ni, and H.-Y. Shum, “Mask dino: Towards a unified transformer-based framework for object detection and segmentation,” in *Computer Vision and Pattern Recognition*, 2023.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Computer Vision and Pattern Recognition*, 2016.
- [9] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*, 2021.
- [10] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *International Conference on Computer Vision*, 2021.
- [11] J. Li, D. Li, S. Savarese, and S. Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” *arXiv preprint arXiv:2301.12597*, 2023.
- [12] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Computer Vision and Pattern Recognition*, 2022.
- [13] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning*, 2021.
- [14] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, “Segment anything,” *arXiv preprint arXiv:2304.02643*, 2023.
- [15] T. Chen, S. Saxena, L. Li, D. J. Fleet, and G. Hinton, “Pix2seq: A language modeling framework for object detection,” in *International Conference on Learning Representations*, 2022.
- [16] P. Wang, A. Yang, R. Men, J. Lin, S. Bai, Z. Li, J. Ma, C. Zhou, J. Zhou, and H. Yang, “Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework,” in *International Conference on Machine Learning*, 2022.
- [17] X. Wang, W. Wang, Y. Cao, C. Shen, and T. Huang, “Images speak in images: A generalist painter for in-context visual learning,” in *Computer Vision and Pattern Recognition*, 2023.
- [18] D. Surís, S. Menon, and C. Vondrick, “Vipergpt: Visual inference via python execution for reasoning,” *arXiv preprint arXiv:2303.08128*, 2023.
- [19] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, and Y. Zhuang, “Hugging-gpt: Solving ai tasks with chatgpt and its friends in huggingface,” *arXiv preprint arXiv:2303.17580*, 2023.
- [20] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *arXiv preprint arXiv:2304.08485*, 2023.
- [21] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [22] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik *et al.*, “Habitat: A platform for embodied ai research,” in *Computer Vision and Pattern Recognition*, 2019.
- [23] M. Deitke, E. VanderBilt, A. Herrasti, L. Weihs, J. Salvador, K. Ehsani, W. Han, E. Kolve, A. Farhadi, A. Kembhavi *et al.*, “Procthor: Large-scale embodied ai using procedural generation,” in *Advances in Neural Information Processing Systems*, 2022.
- [24] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International Conference on Machine Learning*, 2020.
- [25] H. Bao, L. Dong, S. Piao, and F. Wei, “Beit: Bert pre-training of image transformers,” in *International Conference on Learning Representations*, 2022.
- [26] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [27] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev *et al.*, “Grandmaster level in starcraft ii using multi-agent reinforcement learning,” *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.
- [28] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, “Target-driven visual navigation in indoor scenes using deep reinforcement learning,” in *International Conference on Robotics and Automation*, 2017.
- [29] A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra, “Embodied question answering,” in *Computer Vision and Pattern Recognition*, 2018.
- [30] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu *et al.*, “Palm-e: An embodied multimodal language model,” *arXiv preprint arXiv:2303.03378*, 2023.
- [31] K. Kotar, A. Walsman, and R. Mottaghi, “Entl: Embodied navigation trajectory learner,” *arXiv preprint arXiv:2304.02639*, 2023.
- [32] J. McCarthy, “What is artificial intelligence,” 2007.
- [33] S. Legg, M. Hutter *et al.*, “A collection of definitions of intelligence,” *Frontiers in Artificial Intelligence and applications*, vol. 157, p. 17, 2007.
- [34] S. Legg and M. Hutter, “Universal intelligence: A definition of machine intelligence,” *Minds and Machines*, vol. 17, pp. 391–444, 2007.
- [35] S. Legg, M. Hutter *et al.*, “A formal measure of machine intelligence,” in *Annual Machine Learning Conference of Belgium and The Netherlands*, 2006.
- [36] B. Goertzel, “The hidden pattern: A patternist philosophy of mind,” 2006.
- [37] —, “Artificial general intelligence: concept, state of the art, and future prospects,” *Journal of Artificial General Intelligence*, vol. 5, no. 1, p. 1, 2014.
- [38] A. Turing, *Computing Machinery and Intelligence*. MIT Press, 1950.
- [39] J. R. Searle, “Minds, brains, and programs,” *Behavioral and Brain Sciences*, vol. 3, no. 3, pp. 417–424, 1980.
- [40] R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, A. Jin, T. Bos, L. Baker, Y. Du *et al.*, “Lamda: Language models for dialog applications,” *arXiv preprint arXiv:2201.08239*, 2022.
- [41] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung, “Survey of hallucination in natural language generation,” *ACM Computing Surveys*, vol. 55, no. 12, pp. 1–38, 2023.
- [42] B. Goertzel and C. Pennachin, *Artificial general intelligence*. Springer, 2007.
- [43] D. Silver, S. Singh, D. Precup, and R. S. Sutton, “Reward is enough,” *Artificial Intelligence*, vol. 299, p. 103535, 2021.
- [44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, 2017.
- [45] S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez, Y. Sulsky, J. Kay, J. T. Springenberg *et al.*, “A generalist agent,” *Transactions on Machine Learning Research*, 2022.

- [46] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "Minigt-4: Enhancing vision-language understanding with advanced large language models," *arXiv preprint arXiv:2304.10592*, 2023.
- [47] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, "Improving language understanding by generative pre-training," in *Advances in Neural Information Processing Systems*, 2018.
- [48] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *North American Chapter of the Association for Computational Linguistics*, 2019.
- [49] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei, "Deep reinforcement learning from human preferences," in *Advances in Neural Information Processing Systems*, 2017.
- [50] I. Kokkinos, "Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory," in *Computer Vision and Pattern Recognition*, 2017.
- [51] A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classification with deep convolutional networks," in *Advances in Neural Information Processing Systems*, 2012.
- [52] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Computer Vision and Pattern Recognition*, 2022.
- [53] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition*, 2009.
- [54] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, "Revisiting unreasonable effectiveness of data in deep learning era," in *International Conference on Computer Vision*, 2017.
- [55] J. Yu, Z. Wang, V. Vasudevan, L. Yeung, M. Seyedhosseini, and Y. Wu, "Coca: Contrastive captioners are image-text foundation models," *Transactions on Machine Learning Research*, 2022.
- [56] R. Girshick, "Fast r-cnn," in *International Conference on Computer Vision*, 2015.
- [57] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, 2015.
- [58] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *International Conference on Computer Vision*, 2017.
- [59] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Computer Vision and Pattern Recognition*, 2016.
- [60] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European Conference on Computer Vision*, 2016.
- [61] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision*, 2020.
- [62] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision*, 2014.
- [63] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Computer Vision and Pattern Recognition*, 2015.
- [64] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [65] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2015.
- [66] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *International Conference on 3D Vision (3DV)*, 2016.
- [67] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang *et al.*, "Deep high-resolution representation learning for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 10, pp. 3349–3364, 2020.
- [68] B. Cheng, A. Schwing, and A. Kirillov, "Per-pixel classification is not all you need for semantic segmentation," in *Advances in Neural Information Processing Systems*, 2021.
- [69] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," in *Advances in Neural Information Processing Systems*, 2021.
- [70] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ade20k dataset," in *Computer Vision and Pattern Recognition*, 2017.
- [71] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International Journal of Computer Vision*, vol. 123, pp. 32–73, 2017.
- [72] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [73] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *Computer Vision and Pattern Recognition*, 2015.
- [74] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International Conference on Machine Learning*, 2015.
- [75] L. Zhou, H. Palangi, L. Zhang, H. Hu, J. Corso, and J. Gao, "Unified vision-language pre-training for image captioning and vqa," in *AAAI Conference on Artificial Intelligence*, 2020.
- [76] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, "Flamingo: a visual language model for few-shot learning," in *Advances in Neural Information Processing Systems*, 2022.
- [77] J. Li, D. Li, C. Xiong, and S. Hoi, "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation," in *International Conference on Machine Learning*, 2022.
- [78] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint arXiv:2204.06125*, 2022.
- [79] A. Razavi, A. Van den Oord, and O. Vinyals, "Generating diverse high-fidelity images with vq-vae-2," in *Advances in Neural Information Processing Systems*, 2019.
- [80] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *International Conference on Machine Learning*, 2021.
- [81] M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang *et al.*, 2021.
- [82] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," in *Advances in Neural Information Processing Systems*, 2022.
- [83] B. Yan, H. Peng, J. Fu, D. Wang, and H. Lu, "Learning spatio-temporal transformer for visual tracking," in *International Conference on Computer Vision*, 2021.
- [84] T. Meinhardt, A. Kirillov, L. Leal-Taixe, and C. Feichtenhofer, "Trackformer: Multi-object tracking with transformers," in *Computer Vision and Pattern Recognition*, 2022.
- [85] Y. Zhang, P. Sun, Y. Jiang, D. Yu, F. Weng, Z. Yuan, P. Luo, W. Liu, and X. Wang, "Bytetrack: Multi-object tracking by associating every detection box," in *European Conference on Computer Vision*, 2022.
- [86] Y. Xu, J. Zhang, Q. Zhang, and D. Tao, "Vitpose: Simple vision transformer baselines for human pose estimation," in *Advances in Neural Information Processing Systems*, 2022.
- [87] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, "Open-vocabulary object detection via vision and language knowledge distillation," in *International Conference on Learning Representations*, 2022.
- [88] B. Li, K. Q. Weinberger, S. Belongie, V. Koltun, and R. Ranftl, "Language-driven semantic segmentation," in *International Conference on Learning Representations*, 2022.
- [89] H. Zhang, P. Zhang, X. Hu, Y.-C. Chen, L. Li, X. Dai, L. Wang, L. Yuan, J.-N. Hwang, and J. Gao, "Glipv2: Unifying localization and vision-language understanding," in *Advances in Neural Information Processing Systems*, 2022.
- [90] C. Tang, L. Xie, X. Zhang, X. Hu, and Q. Tian, "Visual recognition by request," in *Computer Vision and Pattern Recognition*, 2023.
- [91] M. Malinowski and M. Fritz, "A multi-world approach to question answering about real-world scenes based on uncertain input," in *Advances in neural information processing systems*, 2014.
- [92] M. Malinowski, M. Rohrbach, and M. Fritz, "Ask your neurons: A neural-based approach to answering questions about images," in *International Conference on Computer Vision*, 2015.

- [93] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein, "Neural module networks," in *Computer Vision and Pattern Recognition*, 2016.
- [94] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *International Conference on Machine Learning*, 2021.
- [95] S. Huang, L. Dong, W. Wang, Y. Hao, S. Singhal, S. Ma, T. Lv, L. Cui, O. K. Mohammed, Q. Liu *et al.*, "Language is not all you need: Aligning perception with language models," *arXiv preprint arXiv:2302.14045*, 2023.
- [96] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.
- [97] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot, J. Dy, and T. Pfister, "Learning to prompt for continual learning," in *Computer Vision and Pattern Recognition*, 2022.
- [98] Y. Du, F. Wei, Z. Zhang, M. Shi, Y. Gao, and G. Li, "Learning to prompt for open-vocabulary object detection with vision-language model," in *Computer Vision and Pattern Recognition*, 2022.
- [99] M. Xu, Z. Zhang, F. Wei, Y. Lin, Y. Cao, H. Hu, and X. Bai, "A simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language model," in *European Conference on Computer Vision*, 2022.
- [100] D. Huynh, J. Kuen, Z. Lin, J. Gu, and E. Elhamifar, "Open-vocabulary instance segmentation via robust cross-modal pseudo-labeling," in *Computer Vision and Pattern Recognition*, 2022.
- [101] Z. Ding, J. Wang, and Z. Tu, "Open-vocabulary panoptic segmentation with maskclip," *arXiv preprint arXiv:2208.08984*, 2022.
- [102] J. Jain, J. Li, M. T. Chiu, A. Hassani, N. Orlov, and H. Shi, "One-former: One transformer to rule universal image segmentation," in *Computer Vision and Pattern Recognition*, 2023.
- [103] L. H. Li, P. Zhang, H. Zhang, J. Yang, C. Li, Y. Zhong, L. Wang, L. Yuan, L. Zhang, J.-N. Hwang *et al.*, "Grounded language-image pre-training," in *Computer Vision and Pattern Recognition*, 2022.
- [104] Y. Zhong, J. Yang, P. Zhang, C. Li, N. Codella, L. H. Li, L. Zhou, X. Dai, L. Yuan, Y. Li *et al.*, "Regionclip: Region-based language-image pretraining," in *Computer Vision and Pattern Recognition*, 2022.
- [105] M. Cherti, R. Beaumont, R. Wightman, M. Wortsman, G. Ilharco, C. Gordon, C. Schuhmann, L. Schmidt, and J. Jitsev, "Reproducible scaling laws for contrastive language-image learning," in *Computer Vision and Pattern Recognition*, 2023.
- [106] Q. Sun, Y. Fang, L. Wu, X. Wang, and Y. Cao, "Eva-clip: Improved training techniques for clip at scale," *arXiv preprint arXiv:2303.15389*, 2023.
- [107] J. Chen, Z. Yang, and L. Zhang, "Semantic segment anything," <https://github.com/fudan-zvg/Semantic-Segment-Anything>, 2023.
- [108] J. Cen, Z. Zhou, J. Fang, W. Shen, L. Xie, X. Zhang, and Q. Tian, "Segment anything in 3d with nerfs," *arXiv preprint arXiv:2304.12308*, 2023.
- [109] T. Yu, R. Feng, R. Feng, J. Liu, X. Jin, W. Zeng, and Z. Chen, "Inpaint anything: Segment anything meets image inpainting," *arXiv preprint arXiv:2304.06790*, 2023.
- [110] T. Chen, L. Zhu, C. Ding, R. Cao, S. Zhang, Y. Wang, Z. Li, L. Sun, P. Mao, and Y. Zang, "Sam fails to segment anything?—sam-adapt: Adapting sam in underperformed scenes: Camouflage, shadow, and more," *arXiv preprint arXiv:2304.09148*, 2023.
- [111] J. Ma and B. Wang, "Segment anything in medical images," *arXiv preprint arXiv:2304.12306*, 2023.
- [112] R. Deng, C. Cui, Q. Liu, T. Yao, L. W. Remedios, S. Bao, B. A. Landman, L. E. Wheless, L. A. Coburn, K. T. Wilson *et al.*, "Segment anything model (sam) for digital pathology: Assess zero-shot segmentation on whole slide imaging," *arXiv preprint arXiv:2304.04155*, 2023.
- [113] S. He, R. Bao, J. Li, P. E. Grant, and Y. Ou, "Accuracy of segment-anything model (sam) in medical image segmentation tasks," *arXiv preprint arXiv:2304.09324*, 2023.
- [114] L. Tang, H. Xiao, and B. Li, "Can sam segment anything? when sam meets camouflaged object detection," *arXiv preprint arXiv:2304.04709*, 2023.
- [115] J. Yang, M. Gao, Z. Li, S. Gao, F. Wang, and F. Zheng, "Track anything: Segment anything meets videos," *arXiv preprint arXiv:2304.11968*, 2023.
- [116] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu *et al.*, "Grounding dino: Marrying dino with grounded pre-training for open-set object detection," *arXiv preprint arXiv:2303.05499*, 2023.
- [117] X. Zou, J. Yang, H. Zhang, F. Li, L. Li, J. Gao, and Y. J. Lee, "Segment everything everywhere all at once," *arXiv preprint arXiv:2304.06718*, 2023.
- [118] T. Chen, S. Saxena, L. Li, T.-Y. Lin, D. J. Fleet, and G. E. Hinton, "A unified sequence interface for vision tasks," in *Advances in Neural Information Processing Systems*, 2022.
- [119] X. Wang, X. Zhang, Y. Cao, W. Wang, C. Shen, and T. Huang, "Seggpt: Segmenting everything in context," *arXiv preprint arXiv:2304.03284*, 2023.
- [120] A. Kolesnikov, A. Susano Pinto, L. Beyer, X. Zhai, J. Harmsen, and N. Houlsby, "Uvim: A unified modeling approach for vision with learned guiding codes," in *Advances in Neural Information Processing Systems*, 2022.
- [121] A. Yang, A. Nagrani, P. H. Seo, A. Miech, J. Pont-Tuset, I. Laptev, J. Sivic, and C. Schmid, "Vid2seq: Large-scale pretraining of a visual language model for dense video captioning," in *Computer Vision and Pattern Recognition*, 2023.
- [122] J. Lu, C. Clark, R. Zellers, R. Mottaghi, and A. Kembhavi, "Unified-io: A unified model for vision, language, and multimodal tasks," *arXiv preprint arXiv:2206.08916*, 2022.
- [123] X. Zhu, J. Zhu, H. Li, X. Wu, H. Li, X. Wang, and J. Dai, "Uni-perceiver: Pre-training unified architecture for generic perception for zero-shot and few-shot tasks," in *Computer Vision and Pattern Recognition*, 2022.
- [124] P. Lu, B. Peng, H. Cheng, M. Galley, K.-W. Chang, Y. N. Wu, S.-C. Zhu, and J. Gao, "Chameleon: Plug-and-play compositional reasoning with large language models," *arXiv preprint arXiv:2304.09842*, 2023.
- [125] Z. Yang, L. Li, J. Wang, K. Lin, E. Azarnasab, F. Ahmed, Z. Liu, C. Liu, M. Zeng, and L. Wang, "Mm-react: Prompting chatgpt for multimodal reasoning and action," *arXiv preprint arXiv:2303.11381*, 2023.
- [126] J. Johnson, B. Hariharan, L. Van Der Maaten, J. Hoffman, L. Fei-Fei, C. Lawrence Zitnick, and R. Girshick, "Inferring and executing programs for visual reasoning," in *International Conference on Computer Vision*, 2017.
- [127] R. Zellers, Y. Bisk, A. Farhadi, and Y. Choi, "From recognition to cognition: Visual commonsense reasoning," in *Computer Vision and Pattern Recognition*, 2019.
- [128] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *International Conference on Computer Vision*, 2015.
- [129] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, "Making the v in vqa matter: Elevating the role of image understanding in visual question answering," in *Computer Vision and Pattern Recognition*, 2017.
- [130] J. Johnson, B. Hariharan, L. Van Der Maaten, L. Fei-Fei, C. Lawrence Zitnick, and R. Girshick, "Clevr: A diagnostic dataset for compositional language and elementary visual reasoning," in *Computer Vision and Pattern Recognition*, 2017.
- [131] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan, "Visual chatgpt: Talking, drawing and editing with visual foundation models," *arXiv preprint arXiv:2303.04671*, 2023.
- [132] D. Marr, *Vision: A computational investigation into the human representation and processing of visual information*. San Francisco: W. H. Freeman and Company, 1982.
- [133] H. Moravec, *Mind children: The future of robot and human intelligence*. Harvard University Press, 1988.
- [134] R. A. Brooks, "Elephants don't play chess," *Robotics and Autonomous Systems*, vol. 6, no. 1-2, pp. 3–15, 1990.
- [135] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, pp. 211–252, 2015.
- [136] A. Yuille and D. Kersten, "Vision as bayesian inference: analysis by synthesis?" *Trends in Cognitive Sciences*, vol. 10, no. 7, pp. 301–308, 2006.
- [137] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Her-rasti, M. Deitke, K. Ehsani, D. Gordon, Y. Zhu *et al.*, "Ai2-thor: An interactive 3d environment for visual ai," *arXiv preprint arXiv:1712.05474*, 2017.
- [138] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei,

- “Scaling laws for neural language models,” *arXiv preprint arXiv:2001.08361*, 2020.
- [139] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler *et al.*, “Emergent abilities of large language models,” *arXiv preprint arXiv:2206.07682*, 2022.
- [140] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, “Unsupervised feature learning via non-parametric instance discrimination,” in *Computer Vision and Pattern Recognition*, 2018.
- [141] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018.
- [142] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *Computer Vision and Pattern Recognition*, 2020.
- [143] Z. Xie, Z. Zhang, Y. Cao, Y. Lin, J. Bao, Z. Yao, Q. Dai, and H. Hu, “SimMIM: A simple framework for masked image modeling,” in *Computer Vision and Pattern Recognition*, 2022.
- [144] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [145] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, “Trust region policy optimization,” in *International Conference on Machine Learning*, 2015.
- [146] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, “Mastering the game of go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [147] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton *et al.*, “Mastering the game of go without human knowledge,” *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [148] M. Hessel, J. Modayil, H. Van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver, “Rainbow: Combining improvements in deep reinforcement learning,” in *AAAI Conference on Artificial Intelligence*, 2018.
- [149] A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune, “First return, then explore,” *Nature*, vol. 590, no. 7847, pp. 580–586, 2021.
- [150] L. Smith and M. Gasser, “The development of embodied cognition: Six lessons from babies,” *Artificial Life*, vol. 11, no. 1-2, pp. 13–29, 2005.
- [151] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, “Curiosity-driven exploration by self-supervised prediction,” in *International Conference on Machine Learning*, 2017.
- [152] T. Chen, S. Gupta, and A. Gupta, “Learning exploration policies for navigation,” *arXiv preprint arXiv:1903.01959*, 2019.
- [153] D. S. Chaplot, H. Jiang, S. Gupta, and A. Gupta, “Semantic curiosity for active visual learning,” in *European Conference on Computer Vision*, 2020.
- [154] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva *et al.*, “On evaluation of embodied navigation agents,” *arXiv preprint arXiv:1807.06757*, 2018.
- [155] D. S. Chaplot, R. Salakhutdinov, A. Gupta, and S. Gupta, “Neural topological slam for visual navigation,” in *Computer Vision and Pattern Recognition*, 2020.
- [156] D. Gordon, A. Kembhavi, M. Rastegari, J. Redmon, D. Fox, and A. Farhadi, “Iqa: Visual question answering in interactive environments,” in *Computer Vision and Pattern Recognition*, 2018.
- [157] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [158] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, “pixelnerf: Neural radiance fields from one or few images,” in *Computer Vision and Pattern Recognition*, 2021.
- [159] K. Zhang, G. Riegler, N. Snavely, and V. Koltun, “Nerf++: Analyzing and improving neural radiance fields,” *arXiv preprint arXiv:2010.07492*, 2020.
- [160] A. Petrenko, E. Wijmans, B. Shacklett, and V. Koltun, “Megaverse: Simulating embodied agents at one million experiences per second,” in *International Conference on Machine Learning*, 2021.