

---



# InternChat: Solving Vision-Centric Tasks by Interacting with Chatbots Beyond Language

---

Zhaoyang Liu<sup>\*1</sup>, Yanan He<sup>\*1</sup>, Wenhai Wang<sup>\*†1</sup>, Weiyun Wang<sup>\*1</sup>, Yi Wang<sup>\*1</sup>, Shoufa Chen<sup>\*2,1</sup>,  
Qinglong Zhang<sup>\*1</sup>, Yang Yang<sup>1</sup>, Qingyun Li<sup>1</sup>, Jiashuo Yu<sup>1</sup>, Kunchang Li<sup>3,1</sup>, Zhe Chen<sup>4,1</sup>,  
Xue Yang<sup>1</sup>, Xizhou Zhu<sup>5,1</sup>, Yali Wang<sup>3,1</sup>, Limin Wang<sup>4,1</sup>, Ping Luo<sup>2,1</sup>, Jifeng Dai<sup>6,1</sup>, Yu Qiao<sup>1</sup>

<sup>1</sup>OpenGVLab, Shanghai AI Laboratory   <sup>2</sup>The University of Hong Kong

<sup>3</sup>Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences

<sup>4</sup>Nanjing University   <sup>5</sup>SenseTime Research   <sup>6</sup>Tsinghua University

<https://github.com/OpenGVLab/InternChat>

*We're going to use the best pointing device in the world. We're going to use a pointing device that we're all born with — born with ten of them. We're going to use our fingers. We're going to touch this with our fingers.*

— Steve Jobs

## Abstract

We present an interactive visual framework named InternChat, or iChat for short. The framework integrates chatbots that have planning and reasoning capabilities, such as ChatGPT, with non-verbal instructions like pointing movements that enable users to directly manipulate images or videos on the screen. Pointing (including gestures, cursors, etc.) movements can provide more flexibility and precision in performing vision-centric tasks that require fine-grained control, editing, and generation of visual content. The name InternChat stands for **interaction**, **nonverbal**, and **chatbots**. Different from existing interactive systems that rely on pure language, by incorporating pointing instructions, the proposed iChat significantly improves the efficiency of communication between users and chatbots, as well as the accuracy of chatbots in vision-centric tasks, especially in complicated visual scenarios where the number of objects is greater than 2. Additionally, in iChat, an auxiliary control mechanism is used to improve the control capability of LLM, and a large vision-language model termed Husky is fine-tuned for high-quality multi-modal dialogue (impressing ChatGPT-3.5-turbo with 93.89% GPT-4 Quality). We hope this work can spark new ideas and directions for future interactive visual systems.

## 1 Introduction

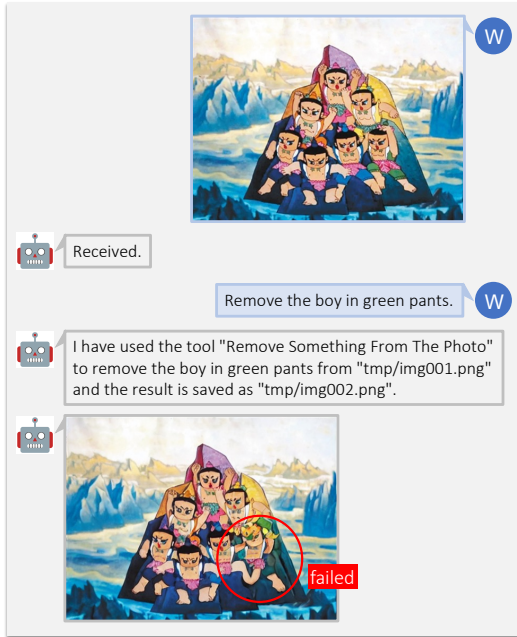
Vision-centric tasks aim to enable computers to understand what they see from the world and react accordingly. In the past, these tasks were solved one by one using specific vision foundation models (VFMs) that were pre-defined and trained for specific visual concepts (*e.g.*, classes, masks, etc.). However, this approach is limited by the availability and quality of labeled data and the diversity of visual scenarios. Recently, the blossom of large language models (LLMs) such as ChatGPT [44], GPT-4 [43], and LLaMA [57] has opened up new possibilities for solving vision-centric tasks. This paradigm leverages LLMs to learn unified real-world concepts and make decisions or plans with vision foundation models (VFM). This approach is user-friendly, requiring almost no domain knowledge for daily or professional tasks via dialogues. It has led to the development of various killer applications, *e.g.*, Visual ChatGPT [69], MM-REACT [74], HuggingGPT [53], etc.

---

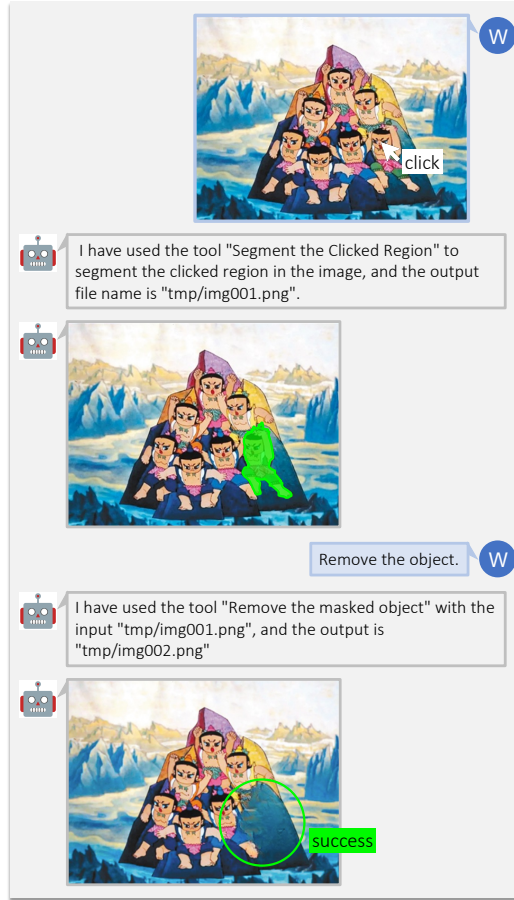
\* Equal contribution. † Project lead

Purely language-driven interactive systems, like Visual ChatGPT, HuggingGPT, may not be sufficient for handling complicated visual scenarios.

↓ **Now you have a pointing device.** →



(a) previous purely language-driven interactive systems



(b) *pointing-language-driven* InternChat (ours)

Figure 1: Advantage of our *pointing-language-driven* interactive system.

Although words are convenient for defining tasks and describing objects, actions, scenes, events, etc, and language-based instructions powered by LLMs allow us to enjoy the capabilities of AI systems, current interactive systems [69, 74, 53] have limitations in connecting vision and language models. They rely mainly on text instructions to interact with visual instances. This becomes highly inefficient when dealing with complicated visual scenarios involving multiple instances since we need to describe the desired instance in length by giving details to discriminate it from others.

Pointing movements such as gestures and cursors are the most common form of non-verbal instructions for communication, showing accuracy in pointing and brevity in operations. Modern operating systems rely heavily on pointing devices such as touchscreens and mice for navigation and interaction. To overcome the aforementioned limitations in current interactive systems, we need to explore new methods to integrate visual perception and language understanding, including both verbal and non-verbal instructions. This will enable efficient and accurate interactions in complicated visual environments.

Inspired by this, we present an advanced interactive visual framework, termed InternChat (iChat), which combines the advantages of pointing and language instructions to perform complex vision-centric tasks. It consists of three main components: a perception unit that handles pointing instructions on images or videos, an LLM controller with an auxiliary control mechanism that can accurately parse language instructions, an open-world toolkit that integrates various online models from HuggingFace, private models trained by users, as well as other applications (*e.g.*, calculator, search engine). InternChat gives equal importance to pointing and language instructions and uses the perception unit and the LLM controller to coordinate and execute the applications in the toolkit to accomplish complicated vision-centric tasks. In addition, for high-quality multi-modal dialogue including image

captioning and visual question answering, we provide a large-scale visual language model named Husky. Husky is quite efficient and effective, which impresses ChatGPT-3.5-turbo with 93.89 % GPT-4 Quality with only 7B parameters.

Compared with previous pure-language interactive frameworks, iChat excels in achieving higher accuracy and efficiency in complicated visual tasks including accurate interactive image or video editing, understanding, question-answering, visual content creation, etc. As evidenced in Figure 1, our system can successfully carry out intricate interactive tasks while pure-language systems fail. Moreover, our user surveys have revealed that combining pointing instructions with language instructions can enhance work efficiency significantly, particularly in complicated scenarios that involve many objects (object number greater than 2). We aim to establish this work as an open baseline for visual interactive systems and will continue updating it with the capability of our VFMs (*e.g.*, InternImage [61] and InternVideo [65]) and the contribution/pull request from the community to achieve even better results.

## 2 Related Work

**Large language model.** Recent LLMs [3, 43, 57, 77] have demonstrated a range of significant abilities, including language generation, in-context learning, world knowledge, and reasoning. The presence of these capabilities enables LLMs to perform complex tasks based on user instructions and prompts in a zero-shot manner. GPT-3 [3], the first language model with over 100 billion parameters, has achieved impressive zero-shot performance on various benchmarks. However, it does not consistently outperform smaller models, such as T5 [50], on some tasks. InstructGPT models [46], which are finetuned on a dataset consisting of prompts with the corresponding human-annotated desired behavior, can be aligned with users, generate outputs that are preferred over those from GPT-3 and show improvements in truthfulness and reductions. Instruction-tuned models have also demonstrated a remarkable ability to generalize zero-shot to new tasks. Thus, instruction-tuning [41, 67, 20, 8] is considered key to eliciting the abilities of LLMs [17]. In addition to GPT model family [48, 49, 3, 43], several other LLMs exist, including OPT [79], LLaMA [57], MOSS [11] and GLM [77]. These models also achieve high performance and are open-sourced, providing valuable experience in training large models and serving as a base for further fine-tuning for different purposes. For example, Alpaca [66] proposes a self-instruct framework to instruction tune the LLaMA model family without relying heavily on human-written instruction data. Another active research area on LLMs is chain-of-thought prompting (CoT) [68, 24, 18, 64]. CoT prompts models to solve problems step by step, greatly improving their reasoning ability of LLMs and making it possible to utilize LLMs for task splitting. As a result, LLMs can be combined with a variety of APIs [32, 42] and models [53, 69] trained for different modalities and serve as a controller to schedule them. This method liberates LLMs from pure language instructions and paves the way for a multi-modal interactive system.

**Perception model.** The emergence of the AlexNet [25] can be considered as the beginning of the development history of deep convolutional neural networks (CNN). Drawing on the success of AlexNet, many CNN with deeper networks, more parameters, and better performance have been proposed and successfully applied to computer vision. However, blindly deepening the network and increasing parameters will not improve the performance without limit but will cause overfitting and increase the cost of experiments. In 2004, GoogleNet [55] overcame the above issues by processing images at multiple scales thanks to the proposed Inception mechanism, which combines convolution operations and pooling operations of different core sizes. In the same year, the concise VGG [54] only used  $(3 \times 3)$  convolution and  $(2 \times 2)$  pooling to win the second place in classification and the first place in object detection. The ResNet [19] has extended the number of layers of the network to an unprecedented scale and solved the issue of deep network degradation, finally achieving 3.57% error on the ImageNet [13] test set. The above models have achieved the dominant position of CNN in the visual field, and it was not until the birth of Vision Transformer (ViT) [14] that this pattern was changed. Benefiting from the powerful Transformer [58] structure, many more advanced transformer-based vision models, *e.g.*, PVT [62, 63], Swin Transformer [39], etc., have been proposed. At the same time, CNN has also been revived (*e.g.*, ConvNeXt [40], InternImage [61]), and some hybrid methods (*e.g.*, ConViT [15], CeiT [76], and CoAtNet [12]) of CNN and Transformer have appeared to make full use of the advantages of both. The benign competition of the visual foundation model has also greatly promoted the development of other visual recognition tasks,

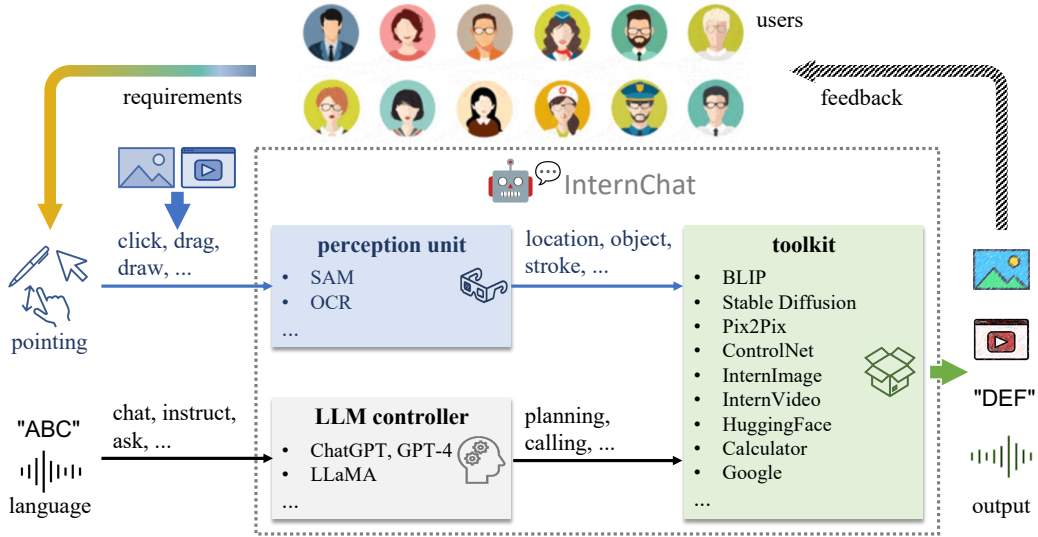


Figure 2: **Overall architecture of InternChat.** It has three main components: perception unit, LLM controller, and open-world toolkit.

such as object detection [51, 36, 4, 81, 73, 71, 72, 75], segmentation [7, 21, 33, 6, 22], video understanding [60, 35, 16, 2, 1, 28, 56, 5, 59, 29, 30, 70], etc. These models have a strong ability to perceive the physical world and can serve as the sensory organs for LLMs.

**LLM-based interactive system.** The success of LLMs [3, 43] has led to the development of AI systems that integrate perception models and LLMs for multimodal reasoning and action. One such system is Visual ChatGPT [69], which connects ChatGPT with visual foundation models to generate and edit images during chatting. Another paradigm is MM-REACT [74], which integrates ChatGPT with a pool of vision experts for multimodal reasoning and action. HuggingGPT [53] utilizes numerous sophisticated AI tasks in different modalities and domains from HuggingFace’s abundant AI models to achieve impressive results. TaskMatrix.AI [34] connects foundation models with millions of APIs from other AI models and systems to perform diversified tasks in both digital and physical domains. However, these systems are limited by their dependence on pure language instructions, which can hinder effective communication and task performance, just like a console cannot achieve what GUI can. By incorporating pointing instructions, the proposed iChat significantly improves the efficiency of communication between users and chatbots, as well as the accuracy of chatbots in vision-centric tasks.

### 3 InternChat

InternChat (iChat) aims to provide an intuitive, user-friendly, and efficient way of human-computer interaction (HCI) in AI that leverages the integration of a large language model (LLM), pointing devices, and computer vision algorithms to perform vision-centric tasks. As illustrated in Figure 2, iChat consists of three main components: (1) a perception unit that interprets the user’s pointing gestures on images and videos, enabling precise object selection and identification; (2) an LLM controller that processes the user’s language commands, facilitating natural communication and (3) an open-world toolkit that integrates various off-the-shelf models/applications to offer a versatile platform for different tasks.

iChat’s design allows it to operate effectively at multiple levels, catering to diverse needs:

**Level 1: basic interaction.** An intuitive way to use iChat is to give it simple commands that trigger pre-defined tasks, such as “caption this photo”. iChat then calls the appropriate model, such as BLIP [27], to produce the desired results. At this level, iChat acts as a front-end of traditional single-task foundation models, without requiring complicated interactive logic, such as chain of thought, contextual reasoning, etc.

**Level 2: Language-guided interaction.** Real-world tasks often involve more complex and diverse demands than traditional pre-defined tasks at level 1. To accomplish these tasks, language instructions with clear specifications or multiple dialogue turns are needed. For example, the command “remove the black dog near the table in the image” requires open-vocabulary detection models to identify the target object based on the language description, and then apply the erase model to remove it from the image. At this level, iChat is an assistant that communicates with users to resolve any ambiguity using natural language.

**Level 3: pointing-language enhanced interaction.** When tasks require more precise specifications that language instructions alone cannot provide, pointing-language instructions become essential. For example, editing/recognizing/OCRing a particular part of an image can be difficult to describe with words. Nonverbal cues, such as gestures or cursor movements, help select, move, or draw objects in the image.

Compared to existing systems like ChatGPT and Visual ChatGPT, *iChat represents a significant advancement in user-centric interaction by combining pointing and language instructions to accomplish complex vision-centric tasks.* Practical use cases include image editing, object manipulation, video annotation, and more, demonstrating its potential impact across various academic and industrial fields. We detail iChat’s design in the following.

### 3.1 Interacting with User

As shown in Figure 2, iChat’s framework accomplishes its tasks through continuous interaction with users. Each time the user provides their requirements via pointing and language instructions, the perception unit analyzes the pointing instructions to identify the target or the content created. Simultaneously, the LLM controller interprets the user’s language instructions, breaking down the task into smaller subtasks and selecting the appropriate tool.

These components work in tandem, delivering accurate and efficient results to users. Similar to previous methods [69], our system keeps a record of the conversation history, enabling users to revisit prior tasks and results. This capability ensures continuous improvement and better performance over time.

### 3.2 Perception Unit

Built on community open-source projects like SAM [23] and OCR [10], iChat’s perception unit parses pointing instructions and performs various operations, such as pick, drag, and draw. It uses click, stroke, drag, and draw gestures for object selection, movement, and content creation. Specifically, the drag gesture is used to move objects to different positions, while the draw gesture is used to create or complete shapes, aided by image generation technology.

The current system uses simple logic to interpret pointing instructions. After the pointing gesture is finished, we handle it in three ways: (1) For typical objects, SAM detects the semantic region to enable selection. (2) For scene text, OCR technology extracts pointed text content. (3) For generation tasks, gestures are stored as strokes and fed into AIGC tools.

The perception unit’s versatility and proficiency in executing various operations based on user input are crucial to our system’s success. Whether the user needs to select an object, move it, or create new content, the perception unit makes it possible. By combining cutting-edge AI solutions with user-friendly interfaces, we ensure a seamless and engaging user experience.

### 3.3 LLM Controller

Similar to previous interactive systems [69, 74, 53], iChat manages complex tasks based on large language models (LLMs). It autonomously parses user language requests, decomposes them into multiple tasks, and plans the task order and dependency based on LLM knowledge. The LLM controller allocates the parsed tasks to the corresponding APIs based on the model descriptions. Through continuous analysis of historical data and user interactions, the system improves task allocation and execution for efficient management.

**Accurate task execution with the auxiliary control mechanism.** Even top language models like ChatGPT and GPT-4 struggle with invoking APIs, especially when parsing and passing arguments. To tackle this issue, iChat employs auxiliary control when LLM fails to act as a controller. It

works as follows: (1) Parse verbs and nouns to identify the API before execution. (2) Extract relevant arguments from prior dialogues based on the API. For example, in the case of “remove the masked object”, query the LLM: “what’s the `image_path` and `mask_path` of the ‘remove the masked object’ API?” This retrieves the “`image_path`” and “`mask_path`” from the chat history. (3) Check argument validity and apply rule-based corrections if necessary. (4) Invoke the API with the identified arguments. This mechanism refines abstract instructions into specific commands, resulting in more accurate task execution.

**Speech transcription.** Additionally, iChat also offers speech transcription features, such as whisper [47] and bark [9], enabling users to communicate without typing.

### 3.4 Toolkit

iChat’s toolkit is called upon by the system. Different from previous works [69, 74, 53], it supports input from pointing devices in addition to standard API descriptions, inputs, and outputs. Here, we show four representative examples of API descriptions as follows:

**Example 1: remove the masked object.** **Input:** original image, mask (pick); **Output:** result image; **Prompt:** “useful when you want to remove an object by masking the region in the image, like: remove the object by the masked region. The input to this tool should be a comma-separated string of two, representing the `image_path` and `mask_path`. The input to this tool should be a string, representing the `image_path`”.

**Example 2: question the masked object.** **Input:** original image, mask (pick); **Output:** result text; **Prompt:** “useful when you need an answer for a question based on a masked image. like: what is the background color in the masked region, how many cats are in this masked figure, what is in this masked figure. The input to this tool should be a comma-separated string of two, representing the `image_path` and the question”.

**Example 3: conditional image generation.** **Input:** stroke draft (drag & draw); **Output:** result image; **Prompt:** “useful when you want to replace an object by clicking in the image with another object or something. like: replace the masked object with a new object or something. The input to this tool should be a comma-separated string of three, representing the `image_path` and the `mask_path` and the prompt”.

**Example 4: video highlight interpretation.** **Input:** original video, mask at timestamp  $t$  (pick); **Output:** result video; **Prompt:** “useful when you want to generate a video with TikTok style based on prompt, like: cut this video to a TikTok video based on a prompt. The input to this tool should be a comma-separated string of two, representing the `video_path` and prompt.”

A comprehensive summary of the toolkit’s APIs is provided in Table 1, covering applications including vision, vision-language foundation models, as well as other applications such as calculators and search engines. This diverse range of APIs enables users to harness advanced techniques for various tasks and achieve their goals more efficiently.

Specifically, for the large vision-language model named Husky in Table 1, we follow the approach of BLIP-2 [26], and replace the language model with LLaMA-7B, which is trained on 52k English instruction-following data generated by GPT-4 [43]. To handle multi-modal tasks without compromising the text-only generation capabilities, we strictly adhere to the input format used in instruct-tuned LLaMA and carry out finetuning of Husky in three stages. As a result, Husky exhibits impressive capabilities, such as visual captioning, visual question-answering, complex reasoning, code generation similar to those of MiniGPT-4 [80] and LLaVA [37].

Category	API
Vision	Stable Diffusion [52], ControlNet [78], InternImage [61], InternVideo [65], SAM [23], DINOv2 [45], etc.
Vision-Language	BLIP-2 [26], Grounding DINO [38], GLIP [31], VideoChat, Husky, etc.
Others	Calculator, Google, HuggingFace, etc.

Table 1: **Representative APIs supported by InternChat.** VideoChat comes from <https://github.com/OpenGVLab/Ask-Anything>.

## 4 Experiment

### 4.1 User Studies

The iChat framework enhances the communication effectiveness of the interactive system with the user by incorporating verbal and nonverbal instructions. To prove the effectiveness of our mixed-mode instructions, which combine language and pointing gestures, in comparison to purely linguistic instructions, we conduct a user study involving 10 human participants. These participants interact with Visual ChatGPT [69] and our iChat through chat and provide their feedback. We present various findings from the user study as follows:

**Efficiency.** We ask users to create an image-centric work by designing input instructions for Visual ChatGPT [69] and iChat, which involve removing and replacing objects. Users can refine their instructions if the results are unsatisfactory. We assume that a user has the patience to attempt up to ten times. Cases with over ten attempts are considered as failed. The number of attempts and prompt lengths required to achieve satisfactory results are presented in Table 2 and Table 3, demonstrating that the iChat is more efficient and user-friendly.

**Human preference.** Users were asked to rank the interactive systems based on their user experience, focusing specifically on the results generated by these systems. Ten unbiased human evaluators were then responsible for assessing the quality of the outcomes. The evaluation results can be found in Table 2 and Table 3. As can be observed, our iChat stands out for its efficiency and user-friendliness, and gains a higher preferences.

### 4.2 Demonstrations

iChat is a versatile framework that can support a range of applications involving verbal and non-verbal interactions. Users can communicate with the system using natural language, as well as gestures such as clicking, dragging, pointing, etc. We showcase several examples of interesting vision-centric applications:

**Demo 1: interactive image editing.** Figure 3 exhibits interactive image editing examples. Except for verbal instructions and an uploaded image, the iChat receives a cursor clicking at the desired operating location of the image. The cursors interaction supplements user instructions to the LLM which controls the visual perception and generation instruments. With more concise and precise instructions, the InternChat demonstrates pleasing image editing performance.

**Demo 2: interactive visual question answering.** Figure 4 presents interactive visual question-answering examples. In the image-centric conversation with an intellectual chatbot, the user’s instructions may revolve around just a region of the image, rather than the entire image. The user can indicate the location through the simplest clicking or touching operation and communicate with the chatbot about the concerned region.

**Demo 3: interactive image generation.** Figure 5 demonstrates image creation examples. The iChat stores multiple picking image parts as materials. The user can easily assemble the materials by the dragging interaction. After receiving a verbal prompt and the assembled images from the user, the iChat controls the generation instruments to create a fancy image, and give it a title.

**Demo 4: video highlight interpretation.** We have implemented a feature called Video Highlight Interpretation, which is demonstrated in Figure 6. This feature is designed to help users quickly and
















	One-Object	Two-Object	Three-Object	Complex-I	Complex-II
Example					
	<b>Visual ChatGPT [69]</b>				
Prompt length	4.9	7.0	8.2	7.3	9.3
# Try (<10)	3.9	9.5	6.2	9.6	9.6
Score (1-5)	3.2	1.5	2.6	0.3	0.6
Example Result					
	<b>InternChat (ours)</b>				
Prompt length	2.7	2.7	3.0	2.6	3.3
# Try (<10)	1.0	1.0	1.3	1.2	1.0
Score (1-5)	4.9	3.6	3.2	4.5	3.4
Example Result					

Table 2: **User study for “remove something”**. “# Try” indicates the number of attempts to get a satisfactory result, not more than 10 times. “-” means trying more than 10 times and still getting bad results.

easily understand the essence of a video without having to watch the entire thing. iChat locates the beginning and end of a video highlight based on its temporal position and then provides a lighthearted and humorous voice-over interpretation of the highlighted content. In addition to its primary use, this feature can also be leveraged by users who want to quickly create derivative works based on a popular video’s main points. By using this feature, users can save time and effort in creating new content based on existing videos, while still ensuring that the new content accurately reflects the original material.

These examples exhibit the convenient interactive capabilities and powerful performance on vision-centric tasks, which come from the novel design of the intelligence system.

### 4.3 Demonstrations of Husky

**Qualitative Evaluation.** In this section, we present demonstrations across various scenarios to assess the performance of the large vision-language model Husky. Qualitative results are shown in Figure 7. These diverse examples showcase the strong capabilities of our Husky.

**Quantitative Evaluation.** We also perform a quantitative evaluation on Husky. We follow the setting of LLaVA [37] and adopt Husky to predict the answers to the provided 90 questions of the 30 COCO validation images. We leverage ChatGPT-3.5-turbo to measure the quality of our model’s generated responses and apply GPT-4’s predictions (provided by LLaVA) as a reference. As reported in Table 4, Husky can impress ChatGPT-3.5-turbo with 93.89 % GPT-4 Quality. It is worth mentioning that, according to the analysis conducted by ChatGPT-3.5-turbo, Husky demonstrates a reasoning capacity that is comparable to that of GPT-4. However, when it comes to conversation and providing detailed descriptions, Husky is still lagging behind that of GPT-4.


















Test Case	Single-Object	Two-Object	Three-Object	Complex-I	Complex-II
Example					
	<b>Visual ChatGPT [69]</b>				
Prompt length	7.0	9.3	9.8	7.8	9.4
# Try	1.2	4.8	9.3	8.5	10
Score (1-5)	4.2	2.2	0.7	0.2	0.0
Example Result					
	<b>InternChat (ours)</b>				
Prompt length	6.6	7.0	6.1	6.9	7.3
# Try	1.0	1.0	1.2	1.0	4.6
Score (1-5)	4.5	5.0	3.8	4.4	1.9
Example Result					

Table 3: **User study for “replace something with something”** “# Try” indicates the number of attempts to get a satisfactory result, not more than 10 times. “-” means trying more than 10 times and still getting bad results.

Conversation	Detail description	Complex reasoning	All
96.13	83.87	102.95	93.89

Table 4: Detailed Quantitative Results

#### 4.4 Limitations

Given that iChat is built upon online resources, potential limitations of this integrated system may include: **Model Performance.** iChat’s effectiveness largely hinges on the quality and accuracy of the underlying open-sourced models. Limitations or biases in these models could adversely affect iChat’s performance.

**Scalability.** As user interactions grow more complex or encompass a greater number of instances, maintaining accuracy and response times could prove challenging for the system. Moreover, the current non-learnable cooperation between vision foundation models and language models, e.g., not being tuned by instruction data, could impede capitalizing on the full capacity of the used models.

**Adaptability.** iChat might struggle to adjust to novel or uncommon scenarios absent from the training data of its employed models, leading to a compromised performance in unpredictable situations.

**User Interface.** Despite emphasizing user-friendliness, some users might still face difficulties utilizing the combined pointing and language instructions effectively, which could impact their overall experience.

**Compatibility.** Achieving seamless integration with a diverse array of devices and platforms could pose challenges due to varying hardware capabilities, software constraints, or accessibility requirements.

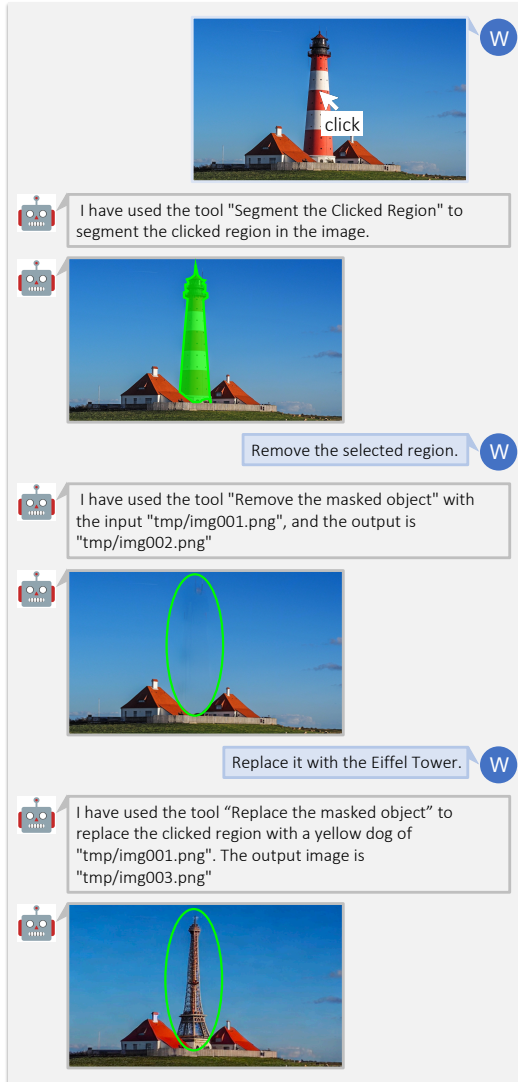


Figure 3: Interactive image editing.

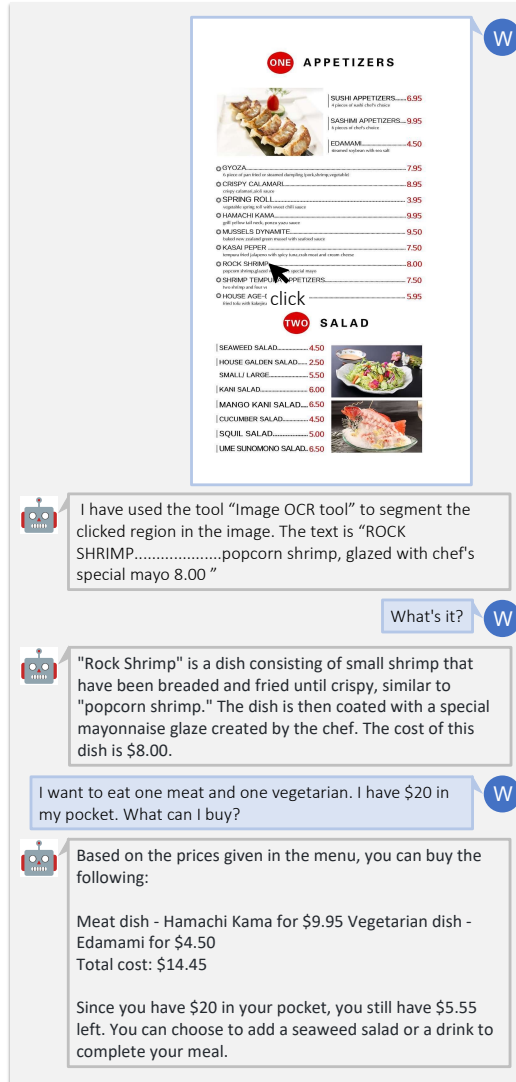


Figure 4: Interactive visual question answer.

**Privacy and Security.** As an AI-driven system, iChat may raise concerns regarding data privacy and security, especially if sensitive information is processed or shared within the platform.

## 5 Conclusion

We have developed InternChat (iChat for short), a dynamic visual framework that tackles vision-centric tasks. Designed with an emphasis on user-friendliness and efficacy, iChat delivers a top-tier experience among current open-source tools. Our methodology seamlessly merges pointing and text commands to issue instructions, allowing users to harness the power of various open-sourced models without needing expert knowledge. The robustness of iChat has been showcased in complex visual scenarios involving multiple instances, utilizing user-level text comprehension from the current LLM and employing cursors or gestures through pointing devices.

iChat appreciates both pointing and linguistic directives, leveraging the perception unit and LLM controller to synchronize and execute applications within an open-world toolkit. Our system has successfully performed intricate interactive tasks beyond the capabilities of purely language-based systems. User surveys have demonstrated that integrating pointing and language instructions can boost work efficiency in challenging situations. Aspiring to be the foundational benchmark for

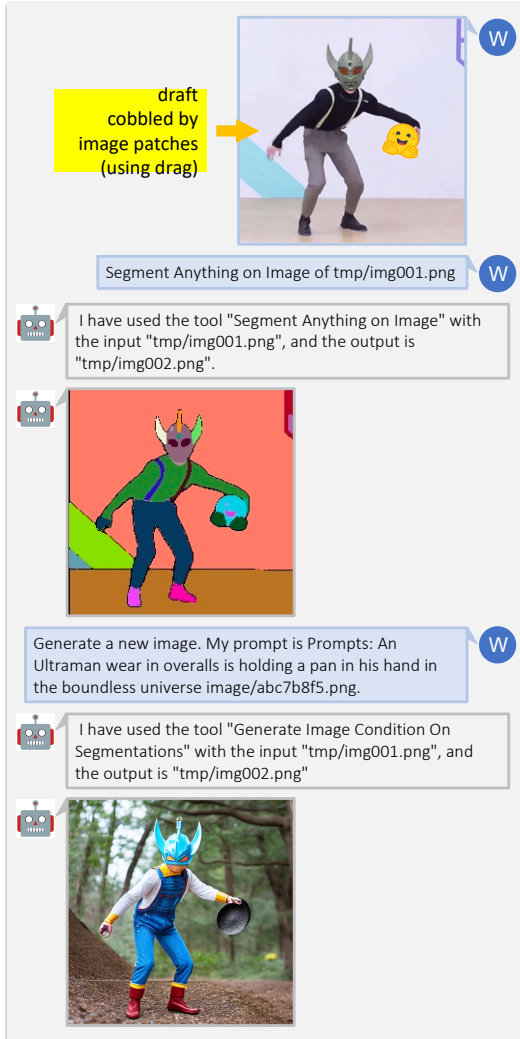


Figure 5: Interactive image generation.

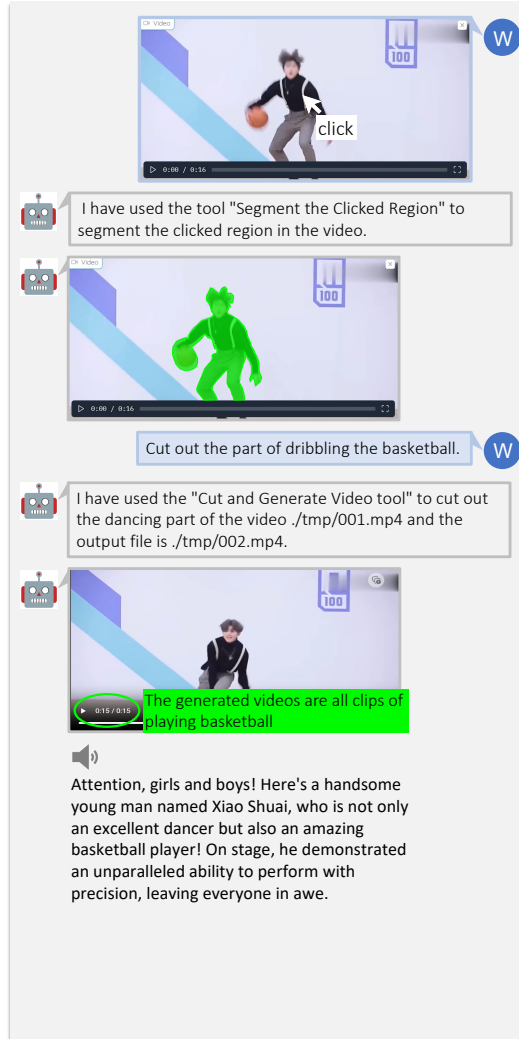


Figure 6: Video highlight interpretation.

visual interactive systems, iChat is committed to ongoing updates and improvements for exceptional performance.

## References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lucic, and Cordelia Schmid. Vivit: A video vision transformer. *2021 IEEE/CVF International Conference on Computer Vision*, pages 6816–6826, 2021. 4
- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *International Conference on Machine Learning*, 2021. 4
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020. 3, 4
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020. 4



**Human:** describe the image in detail.

**Assistant:** The image features a man dressed in a business suit with a chicken costume covering his entire body, giving him an amusing and unique appearance. The man's face is not visible, as it is hidden under the chicken head and mask. This image is an artistic representation and not a real-life scenario.

**Human:** is this image funny?

**Assistant:** The image has a humorous and whimsical quality to it, which may be interpreted as funny by some people. The juxtaposition of a man in a business suit with a chicken costume covering his entire body is unexpected and amusing, invoking a sense of humor and playfulness. However, humor is subjective and may not be universally appreciated, so it may not be considered funny by everyone.

(a) Caption & Joke Understanding



**Human:** write a poem about the image.

**Assistant:** A young man and woman sit together on the rooftop of a tall building,  
A city skyline in view, a view that's quite entralling.

Their embrace is tender, a moment captured in time,  
A moment to cherish, an image of love between these two.

Sitting on the roof of a tall building, they pose for a picture,  
Embraced in each other, a scene that's quite captivating.

Their love so tender, it shines through in their eyes,  
A bond that's strong and lasting, that's quite impressive.

A young couple in love, a moment captured in time,  
A city skyline in view, a beautiful sight to see.

(b) Poem Generation

The photo shot at  
Shanghai AI Lab



**Human:** With the objects in the image, if you feel hot, what will you do, write your plan and explain it.

**Assistant:** 1. First, you would likely check the temperature at the top of the computer monitor to see if the screen gets hot.

2. You would then check the temperature of the fan that is situated on the desk.

3. If the temperature rises above an acceptable level, you would proceed to open the fan and manually adjust the speed or ventilation to cool the device down.

4. While doing so, you would also check the temperature of other components, such as the CPU and RAM, to ensure that they are not excessively hot and functioning optimally.

5. If you notice that the device still feels hot, you would proceed to check the room temperature to determine if it's too high and potentially causing the overheating problem.

6. Finally, you would implement a plan to address the issue, such as adjusting the airflow in the room, turning on the air conditioner or opening a window to let in some fresh air."

(c) Complex Reasoning



**Human:** write a story between the boy and the panda according to the image

**Assistant:** Once upon a time, in a small village nestled between the mountains on a misty morning, a young boy and a giant panda went for their walk down the rain-soaked road.

The boy, who was a curious young lad, was fascinated with the huge panda walking down the road. The panda had a curious look on his face, as he was not used to walking on such a narrow path. As the panda made his way down the road, the boy followed him, trying to understand the panda's movements and behavior. The panda seemed to move slowly and gracefully despite the heavy rain, and the boy was in awe of his strength and agility. The boy was determined to learn more about the panda, so he followed his every step, soaking in the experience like a sponge. As they walked together down the rain-soaked road, the boy and panda bonded in their shared journey, and the young boy was enriched with a unique and unforgettable experience.

(d) Story Generation

Figure 7: Demonstrations of Husky across various scenarios.

- [5] Guo Chen, Sen Xing, Zhe Chen, Yi Wang, Kunchang Li, Yizhuo Li, Yi Liu, Jiahao Wang, Yin-Dong Zheng, Bingkun Huang, et al. Internvideo-ego4d: A pack of champion solutions to ego4d challenges. *arXiv preprint arXiv:2211.09529*, 2022. 4
- [6] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. In *International Conference on Learning Representations*, 2023. 4
- [7] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*, 34:17864–17875, 2021. 4
- [8] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022. 3
- [9] Bark contributors. Bark. <https://github.com/suno-ai/bark>, 2023. 6
- [10] EasyOCR contributors. Easyocr. <https://github.com/JaidedAI/EasyOCR>, 2023. 5
- [11] MOSS contributors. Moss. <https://github.com/OpenMLab/MOSS>, 2023. 3
- [12] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. *Advances in Neural Information Processing Systems*, 34:3965–3977, 2021. 3
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 3
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 3
- [15] Stéphane d’Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *International Conference on Machine Learning*, pages 2286–2296. PMLR, 2021. 3
- [16] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. *2019 IEEE/CVF International Conference on Computer Vision*, pages 6201–6210, 2018. 4
- [17] Hao Fu, Yao; Peng and Tushar Khot. How does gpt obtain its ability? tracing emergent abilities of language models to their sources. *Yao Fu’s Notion*, Dec 2022. 3
- [18] Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. *arXiv preprint arXiv:2211.10435*, 2022. 3
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 3
- [20] Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Dániel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, et al. Opt-impl: Scaling language model instruction meta learning through the lens of generalization. *arXiv preprint arXiv:2212.12017*, 2022. 3
- [21] Jitesh Jain, Jiachen Li, MangTik Chiu, Ali Hassani, Nikita Orlov, and Humphrey Shi. Oneformer: One transformer to rule universal image segmentation. *arXiv preprint arXiv:2211.06220*, 2022. 4
- [22] Yuanfeng Ji, Zhe Chen, Enze Xie, Lanqing Hong, Xihui Liu, Zhaoqiang Liu, Tong Lu, Zhenguo Li, and Ping Luo. Ddp: Diffusion model for dense visual prediction. *arXiv preprint arXiv:2303.17559*, 2023. 4
- [23] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 5, 7
- [24] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022. 3
- [25] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. 3

- [26] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 6, 7
- [27] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 4
- [28] Kunchang Li, Yali Wang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Jiao Qiao. Uniformer: Unified transformer for efficient spatiotemporal representation learning. In *International Conference on Learning Representations*, 2022. 4
- [29] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. Uniformerv2: Spatiotemporal learning by arming image vits with video uniformer. *arXiv preprint arXiv:2211.09552*, 2022. 4
- [30] Kunchang Li, Yali Wang, Yizhuo Li, Yi Wang, Yinan He, Limin Wang, and Yu Qiao. Unmasked teacher: Towards training-efficient video foundation models. *arXiv preprint arXiv:2303.16058*, 2023. 4
- [31] Liunian Harold Li\*, Pengchuan Zhang\*, Haotian Zhang\*, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 7
- [32] Minghao Li, Feifan Song, Bowen Yu, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. Api-bank: A benchmark for tool-augmented llms. *arXiv preprint arXiv:2304.08244*, 2023. 3
- [33] Zhiqi Li, Wenhai Wang, Enze Xie, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, Ping Luo, and Tong Lu. Panoptic segformer: Delving deeper into panoptic segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1280–1289, 2022. 4
- [34] Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, et al. Taskmatrix. ai: Completing tasks by connecting foundation models with millions of apis. *arXiv preprint arXiv:2303.16434*, 2023. 4
- [35] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. *2019 IEEE/CVF International Conference on Computer Vision*, pages 7082–7092, 2018. 4
- [36] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2980–2988, 2017. 4
- [37] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 6, 8
- [38] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 7
- [39] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 3
- [40] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. 3
- [41] Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. *arXiv preprint arXiv:2301.13688*, 2023. 3
- [42] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021. 3
- [43] OpenAI. Gpt-4 technical report. *arXiv*, 2023. 1, 3, 4, 6
- [44] TB OpenAI. Chatgpt: Optimizing language models for dialogue. *OpenAI*, 2022. 1

- [45] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 7
- [46] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022. 3
- [47] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*, 2022. 6
- [48] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. *OpenAI*, 2018. 3
- [49] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 3
- [50] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 3
- [51] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems*, pages 91–99, 2015. 4
- [52] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 7
- [53] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*, 2023. 1, 2, 3, 4, 5, 6
- [54] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 3
- [55] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015. 3
- [56] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In *Advances in Neural Information Processing Systems*, 2022. 4
- [57] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 1, 3
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017. 3
- [59] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 4
- [60] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *European Conference on Computer Vision*, pages 20–36, 2016. 4
- [61] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 3, 7
- [62] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 3

- [63] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022. 3
- [64] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022. 3
- [65] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, Sen Xing, Guo Chen, Junting Pan, Jiashuo Yu, Yali Wang, Limin Wang, and Yu Qiao. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022. 3, 7
- [66] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*, 2022. 3
- [67] Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109, 2022. 3
- [68] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022. 3
- [69] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023. 1, 2, 3, 4, 5, 6, 7, 8, 9
- [70] Min Yang, Guo Chen, Yin-Dong Zheng, Tong Lu, and Limin Wang. Basictad: an astounding rgb-only baseline for temporal action detection. *Computer Vision and Image Understanding*, page 103692, 2023. 4
- [71] Xue Yang, Junchi Yan, Ziming Feng, and Tao He. R3det: Refined single-stage detector with feature refinement for rotating object. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3163–3171, 2021. 4
- [72] Xue Yang, Junchi Yan, Qi Ming, Wentao Wang, Xiaopeng Zhang, and Qi Tian. Rethinking rotated object detection with gaussian wasserstein distance loss. In *International Conference on Machine Learning*, pages 11830–11841. PMLR, 2021. 4
- [73] Xue Yang, Jirui Yang, Junchi Yan, Yue Zhang, Tengfei Zhang, Zhi Guo, Xian Sun, and Kun Fu. Scredet: Towards more robust detection for small, cluttered and rotated objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8232–8241, 2019. 4
- [74] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *ArXiv*, abs/2303.11381, 2023. 1, 2, 4, 5, 6
- [75] Yi Yu, Xue Yang, Qingyun Li, Yue Zhou, Gefan Zhang, Junchi Yan, and Feipeng Da. H2rbox-v2: Boosting hbox-supervised oriented object detection via symmetric learning. *arXiv preprint arXiv:2304.04403*, 2023. 4
- [76] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. Incorporating convolution designs into visual transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 579–588, 2021. 3
- [77] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022. 3
- [78] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 7
- [79] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022. 3
- [80] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 6



- [81] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2021. 4