



If LLM Is the Wizard, Then Code Is the Wand: A Survey on How Code Empowers Large Language Models to Serve as Intelligent Agents

Ke Yang*, Jiateng Liu*, John Wu, Chaoqi Yang, Yi R. Fung, Sha Li,
Zixuan Huang, Xu Cao, Xingyao Wang, Yiquan Wang, Heng Ji, Chengxiang Zhai

University of Illinois Urbana-Champaign

{key4, jiateng5, johnwu3, chaoqiy2, yifung2, shal2, zixuan3, xucao2, xingya06, yiquan2, hengji, czhai}@illinois.edu

Abstract

The prominent large language models (LLMs) of today differ from past language models not only in size, but also in the fact that they are trained on a combination of natural language and formal language (code). As a medium between humans and computers, code translates high-level goals into executable steps, featuring standard syntax, logical consistency, abstraction, and modularity. In this survey, we present an overview of the various benefits of integrating code into LLMs' training data. Specifically, beyond enhancing LLMs in code generation, we observe that these unique properties of code help *i*) unlock the reasoning ability of LLMs, enabling their applications to a range of more complex natural language tasks; *ii*) steer LLMs to produce structured and precise intermediate steps, which can then be connected to external execution ends through function calls; and *iii*) take advantage of code compilation and execution environment, which also provides diverse feedback for model improvement. In addition, we trace how these profound capabilities of LLMs, brought by code, have led to their emergence as intelligent agents (IAs) in situations where the ability to understand instructions, decompose goals, plan and execute actions, and refine from feedback are crucial to their success on downstream tasks. Finally, we present several key challenges and future directions of empowering LLMs and IAs with code.

1 Introduction

Code has become an integral component in the training data of large language models (LLMs), including well-known models such as Llama2, GPT-3.5 series and GPT-4 (Touvron et al., 2023; Ye et al., 2023a; OpenAI, 2023). Training LLMs on code has gained popularity not only because the acquired programming skills enable commercial applications, such as Github Copilot¹, but also because it

improves the models' previously lacking reasoning abilities (Liang et al., 2023b). Consequently, LLMs rapidly emerge as a primary decision-making hub for intelligent agents (IAs) (Zhao et al., 2023), demonstrating an exponential growth in capabilities from code training and the advancement of tool learning (Qin et al., 2023). These LLM-based IAs are poised to handle a wider range of more complex tasks, including downstream applications in multi-agent environment simulation (Wu et al., 2023c) and AI for science (Boiko et al., 2023).

As depicted in Figure 1, this survey aims to explain the widespread adoption of code-specific training in the general LLM training paradigm and how code enhances LLMs to act as IAs. Unlike previous code-LLM surveys that concentrate on either evaluating and comparing code generation abilities (Zan et al., 2023; Xu et al., 2022), or listing IA tasks (Wang et al., 2023d; Xi et al., 2023; Zhao et al., 2023) in IA surveys, we aim to provide a comprehensive understanding of how code assists LLMs and where code benefits LLMs as IAs, based on the taxonomy of relevant papers (see Figure 2).

We first provide our definition of code and present typical methods for LLM code training (§2). Compared to natural language (refer to the case study in A.1), code is more structured, featuring logical, step-by-step executable processes derived from procedural programming, as well as explicitly defined, modularized functions, which compose graphically representable abstractions. Additionally, code is typically accompanied by a self-contained compilation and execution environment. With insights from these characteristics of code, our comprehensive literature review reveals that integrating code into LLM training *i*) enhances their programming and reasoning capabilities (§3); *ii*) enables the models to directly generate executable, fine-grained steps during decision-making, thereby facilitating their scalability in incorporating various tool modules through function calls (§4); and *iii*)

*Equal contribution

¹<https://github.com/features/copilot>.

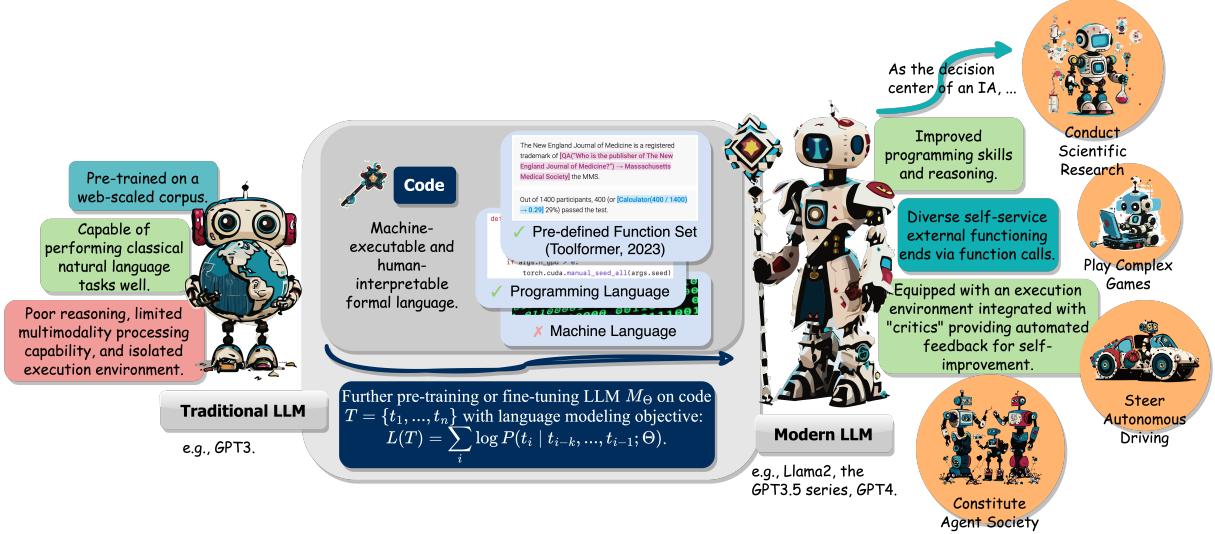


Figure 1: An illustration of how code empowers large language models (LLMs) and enhances their downstream applications as intelligent agents (IAs). While traditional LLMs excel in conventional natural language tasks like **document classification** and **question answering**, further pre-training or fine-tuning LLMs with human-interpretable and machine-executable code serves as an additional power-up — akin to equipping wizards with mana-boosting wands. This significantly boosts their performance as IAs through intricately woven operational steps.

situates the LLMs within a code execution environment, allowing them to receive automated feedback from integrated evaluation modules and self-improve (§5).

In addition, as LLMs are becoming key decision-makers for IAs in complex real-world tasks, our survey also explores how these advantages facilitate their functioning along this capacity (§6), in terms of *i*) enhancing IAs’ decision-making in perception and planning skills (§6.1), *ii*) facilitating their execution through direct action primitive grounding and modular memory organization (§6.2), and *iii*) providing an interactive environment for self-correction and self-improvement (§6.3). Finally, we discuss **several open challenges and promising future directions** (§7).

2 Preliminaries

2.1 Our Definition of Code

We consider code as any formal language that is both machine-executable and human-interpretable. For instance, human-readable programming languages fall within the scope of our discussion, whereas low-level languages, such as machine language based on binary instructions, are excluded due to their lack of human interpretability. Additionally, pre-defined formal languages, such as function sets employed in WebGPT (Nakano et al., 2021), are included as they can be parsed and executed in a rule-based manner.

LLMs trained with expressions formulated within a defined set of symbols and rules (e.g., pre-defined function sets, mathematical deduction formula, etc.), i.e., formal languages, exhibit advantages akin to those trained with programming languages. Therefore, we expand our definition of code to incorporate these homogeneous training corpora, enhancing the comprehensiveness of this survey to align with current research needs.

2.2 LLM Code Training Methods

LLMs undergo code training by following the standard language modeling objective, applied to code corpora. Given that code possesses natural language-like sequential readability, this parallels the approach to instruct LLMs in understanding and generating free-form natural language. Specifically, for an LLM M_Θ with parameters Θ and a code corpus $T = \{t_1, \dots, t_n\}$, the language modeling loss for optimization is:

$$L(T) = \sum_i \log P(t_i | t_{i-k}, \dots, t_{i-1}; \Theta)$$

When employing programming language (e.g., Python, C, etc.) as the corpus (Chen et al., 2021; Li et al., 2022; Nijkamp et al., 2022), training data is typically sourced from publicly accessible code repositories, such as GitHub. This process yields a corpus with a volume comparable to that of natural language pre-training, and thus we call training with such an abundance of code as *code pre-training*. The training strategy entails either train-

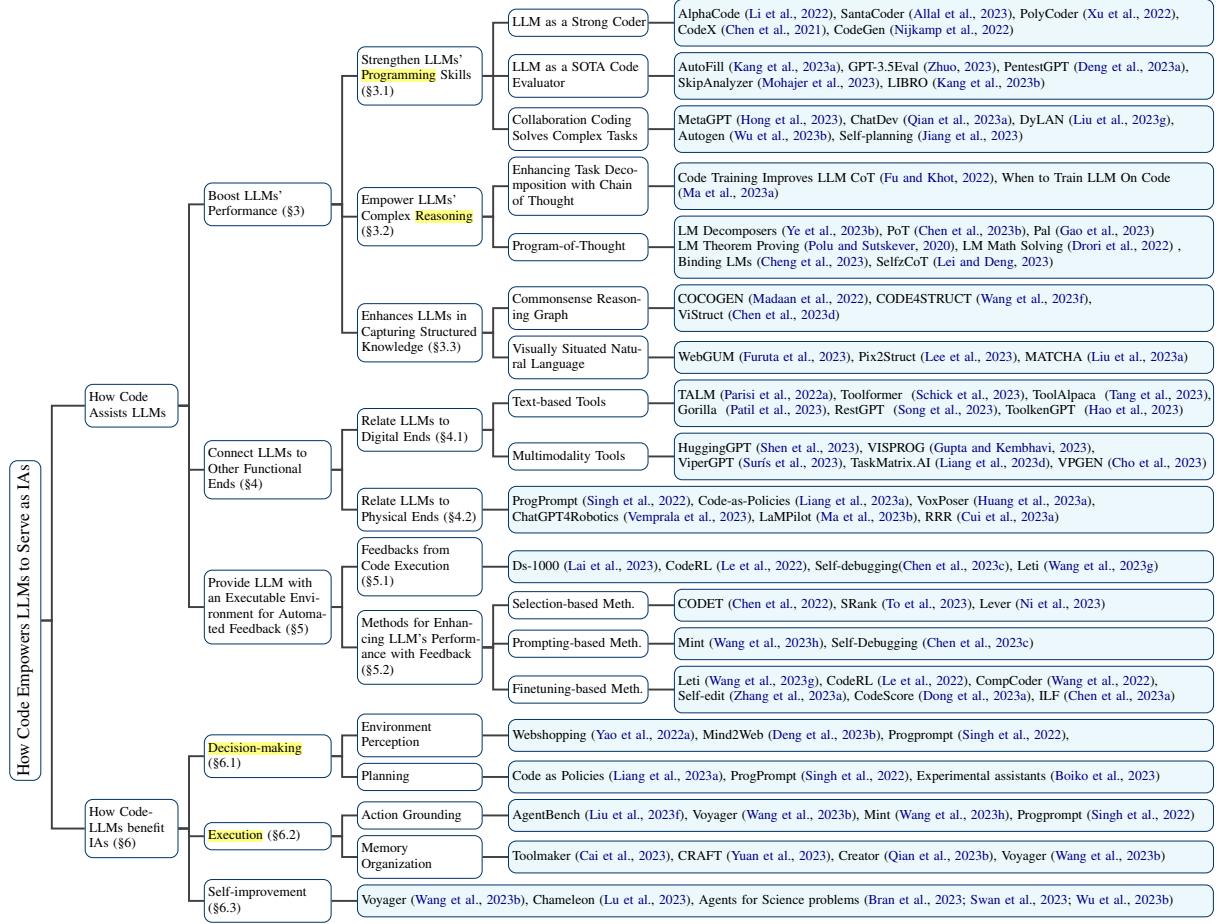


Figure 2: The organization of our paper, with a curated list of the most representative works. The complete work list is provided in Appendix D.

ing code on a pre-trained natural language LLM, as exemplified by Codex (Chen et al., 2021), or training a LLM from scratch with a blend of natural language and code, as demonstrated by CodeLLM (Ma et al., 2023a).

Conversely, when utilizing other pre-defined formal language for training, the objective shifts to acquainting the model with the application of specific functions (Schick et al., 2023), mathematical proof formulas (Wu et al., 2022), SQL (Sun et al., 2023b), and similar constructs. As the dataset for this is smaller compared to the pre-trained natural language corpus, we refer to such training process as *code fine-tuning*. Researchers apply the language modeling loss to optimize LLMs during this process, similarly.

3 Code Pre-Training Boosts LLMs’ Performance

The pre-training of LLMs on code, exemplified by OpenAI’s GPT Codex (Chen et al., 2021), has broadened the LLMs’ scope of tasks beyond nat-

ural language. Such models enable diverse applications, including generating code for mathematical theory (Wu et al., 2022), general programming tasks (Chen et al., 2021), and data retrieval (Sun et al., 2023b; Cheng et al., 2023). Code necessitates producing logically coherent, ordered sequences of steps essential for valid execution. Moreover, the executability of each step within code allows for step-by-step logic verification. Leveraging and embedding both these properties of code in pre-training has improved LLM chain-of-thought (CoT) performance across many conventional natural language downstream tasks (Lyu et al., 2023; Zhou et al., 2023a; Fu and Khot, 2022), indicating **improved complex reasoning skills**. Implicitly learning from code’s structured format, code LLMs demonstrate further improved performance on commonsense structured reasoning tasks, such as those related to markup, HTML, and chart understanding (Furuta et al., 2023; Liu et al., 2023a).

In the following sections, our objective is to elucidate why training LLMs on code and employing

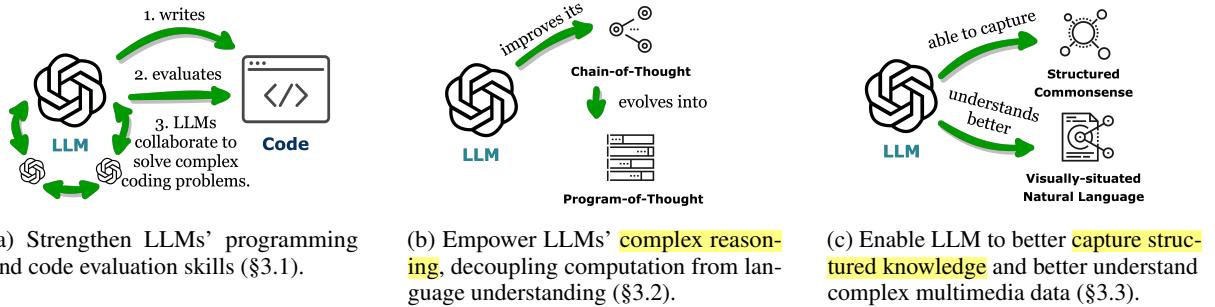


Figure 3: How code pre-training boosts LLMs' performance.

code-based prompts enhance their performance on complex downstream tasks. Specifically, we highlight three key areas where pre-training on code have benefited LLMs: *i*) enhancing programming proficiency in §3.1, *ii*) empowering complex reasoning capabilities in §3.2, and *iii*) facilitating the capture of structured commonsense knowledge in §3.3, as shown in Figure 3.

3.1 Strengthen LLMs' Programming Skills

LLM as a strong coder. Earlier language models only generate domain-specific programs (Ellis et al., 2019) or restrict to one of the generic programming languages, such as Java or C# (Alon et al., 2020). Empowered by the increasing number of parameters and computing resources, recent LLM-based code generation models (such as AlphaCode (Li et al., 2022), CodeGen (Nijkamp et al., 2022), SantaCoder (Allal et al., 2023), PolyCoder (Xu et al., 2022)) could master more than 10 languages within the same model and show unprecedented success. A well-known work is CodeX (Chen et al., 2021), with 12 billion parameters that reads the entire GitHub database and is able to solve 72.31% of challenging Python programming problems created by humans. Recent studies (Zan et al., 2023; Xu et al., 2022; Du et al., 2023; Vaithilingam et al., 2022; Wong et al., 2023; Fan et al., 2023) have provided systematic surveys and evaluations of existing code-LLMs.

With its strong code generation ability, LLMs benefit various applications that rely on code, such as database administration (Zhou et al., 2023b), embedded control (Liang et al., 2023a), game design (Roberts et al.), spreadsheet data analysis (Liu et al., 2023c), and website generation (Calò and De Russis, 2023).

LLM as a state-of-the-art code evaluator. On the

other hand, LLMs themselves could be state-of-the-art evaluators (i.e., analyze and score) for human or machine-generated codes. Kang et al. (2023a) leverage LLM-based models for code fault localization, while Zhuo (2023) uses GPT-3.5 to evaluate the functional correctness and human preferences of code generation. In addition, Deng et al. (2023a) design a LLM-based penetration testing tool and find that LLMs demonstrate proficiency in using testing tools, interpreting outputs, and proposing subsequent actions. Two recent efforts (Li et al., 2023a; Mohajer et al., 2023) also utilize LLM for examining and analyzing source code without executing it. Furthermore, LLMs are used for automatic bug reproduction in Kang et al. (2023b) and vulnerable software evaluation in Noever (2023).

Multi-LLM collaboration solves complex coding problems. Though code LLMs, like GitHub Copilot, have shown unprecedented success, one LLM agent alone could fail in complicated scenarios requiring multiple steps. Luckily, collaborative coding among several role-specific LLM agents exhibits more accurate and robust performance towards complex tasks. Hong et al. (2023) incorporates human programming workflows as guides to coordinate different agents. Dong et al. (2023b) assigned three roles: analyst, coder, and tester to three distinct “GPT-3.5”s, which surpasses GPT-4 in code generation. Meanwhile, Qian et al. (2023a) designs a chat-powered software development process, assigning more than three roles to separate LLM agents. Other similar methods (Liu et al., 2023g; Talebirad and Nadiri, 2023; Wu et al., 2023b; Jiang et al., 2023) all employ multiple code-LLM agents or different phases of the same agent for code generation, software developments or leveraging generated intermediate codes for other general purpose tasks.

3.2 Empower LLMs’ Complex Reasoning

Code pre-training improves chain-of-thought performance. Logically, many complex tasks can be divided into smaller easier tasks for solving. CoT prompting, where prompt inputs are designed with chains of reasoning, allows the LLM to condition its generation with further steps of reasoning, providing a direct approach to task decomposition (Wei et al., 2023). CoT has seen success in the decomposition of many tasks, such as planning (Huang et al., 2022b) and evidence-based question answering (Dua et al., 2022; Ye et al., 2023b).

While LLM CoT ability was originally mainly attributed to dramatically increased model sizes (Wei et al., 2022b), recent evidence compiled by Fu and Khot (2022) suggests that much of the performance improvements from CoT stems from its pre-training on code. For instance, when comparing different versions of GPT-3 (i.e., v1 vs. v5), LLMs not trained on code, such as GPT-3’s text-davinci-001, see a small but substantial accuracy improvement of 6.4 % to 12.4 % with CoT on the mathematical reasoning task GSM8k (Cobbe et al., 2021). In contrast, LLMs pre-trained on code, such as GPT-3’s text-davinci-002 and Codex (Chen et al., 2021), see a dramatic performance improvement arising from CoT, with a remarkable accuracy increase of 15.6% to 46.9% and 19.7% to 63.1% respectively. Supporting this hypothesis proposed by Fu and Khot (2022), Ma et al. (2023a) show that pre-training on code in small-sized LLMs (2.6B) (Zeng et al., 2021) enhances performance when using CoT, and even more remarkably that smaller code-pretrained LLMs outperform their larger non-code counterparts across many different tasks. Furthermore, their study indicates that incorporating a greater volume of code during the initial phases of LLM training significantly enhances its efficacy in reasoning tasks. Nevertheless, tempering expectations, it is possible that the discrepancy in CoT performance between LLMs with and without code pre-training diminishes as the size of the models decreases, as the accuracy gap between the small LLMs in Ma et al. (2023a) was less than 3% when evaluating CoT. Notably, both Fu and Khot (2022) and Ma et al. (2023a) show that pre-training on code improves LLM performance in both standard and CoT prompting scenarios across downstream tasks.

Program-of-thought outperforms chain-of-thought. Furthermore, in comparison to vanilla

CoT methods, LLMs that first translate and decompose a natural language task into code (Chen et al., 2023b; Gao et al., 2023), typically termed program-of-thought (PoT) prompting or program-aided language model, see sizable gains in tasks that require disambiguation in both language and explicit longitudinal structure. This approach is especially effective in complex areas such as theoretical mathematics (Polu and Sutskever, 2020), undergraduate mathematics (Drori et al., 2022), and question answering with data retrieval (Sun et al., 2023b; Cheng et al., 2023).

PoT enhances performance due to the precision and verifiability inherent in code as a machine-executable language. Within task decomposition, it is not uncommon for the LLM to hallucinate incorrect subtasks and questions through CoT (Ji et al., 2023). PoT implementations from Chen et al. (2023b), Gao et al. (2023), and Ye et al. (2023b) show that by directly executing code and verifying outcomes post translation by LLMs, one can effectively mitigate the effects of incorrect reasoning in CoT. This is because the reasoning process must adhere to the logic and constraints explicitly specified by the program, thereby ensuring a more accurate and reliable outcome.

However, such improvements seen in the usage of code are not limited to purely executable coding languages such as Python or SQL, nor are they limited to tasks that are specifically rigid in structure such as mathematics (Drori et al., 2022) and data retrieval (Rajkumar et al., 2022). Enhancements also extend to the realm where even translating into pseudo-code to decompose a task can improve zero-shot performance (Lei and Deng, 2023) in word problems containing numbers, and general reasoning tasks such as StrategyQA (Geva et al., 2021).

3.3 Enable LLMs to Capture Structured Knowledge

Code generation unveils superior structural commonsense reasoning. Given that code possesses the graph structure of symbolic representations, translating textual graphs, tables, and images into code empowers a code-driven LLM to logically process such information according to code reasoning and generation principles.

Consequently, previous work (Madaan et al., 2022; Wang et al., 2023f) shows that LLMs under-

going extra pre-training on code may rival, or even exceed, their fine-tuned natural language counterparts in tasks involving structural commonsense reasoning, even with limited or no training data.

COCOGEN (Madaan et al., 2022) first reframed the commonsense reasoning graph completion task as a code generation task and demonstrated improved few-shot performance in reasoning graphs, table entity state tracking, and explanation graph generation.

Building on this perspective, CODE4STRUCT (Wang et al., 2023f) applied code-LLMs to semantic structures, focusing on the event argument extraction task. By leveraging code’s features such as comments and type annotation, it achieved competitive performance with minimal training instances. Moreover, it surpassed baselines in zero-shot scenarios, benefiting from the inheritance feature and sibling event-type samples. ViStruct (Chen et al., 2023d) extended this approach further to multimodal tasks, leveraging programming language for representing visual structural information and curriculum learning for enhancing the model’s understanding of visual structures.

Markup code mastery evolves visually situated natural language understanding. Another stream of research focuses on utilizing markup code such as HTML and CSS to delineate and derender structured graphical information in graphical user interfaces (GUIs) or visualizations such as plots and charts in documents. This markup code not only specifies content but also governs the layout of a web page, aiding large vision-language models (LVLMs) in capturing visually situated natural language (VSNL) information.

For LVLMs’ markup code understanding, WebGUM (Furuta et al., 2023) exemplified autonomous web navigation. It employed a pre-training approach using webpage screenshots and the corresponding HTML as input, and navigation action as output. Outperforming SOTA, human, and other LLM-based agents, WebGUM showcased the effectiveness of pre-training model with markup code augmentation in webpage understanding.

For markup code generation, pix2code (Beltramelli, 2017) and sketch2code (Robinson, 2019) pioneered machine learning methods to generate rendering code for GUIs or website mockups, in the pre-LLM era. Pix2Struct (Lee et al., 2023) achieved SOTA at that time in VSNL understanding

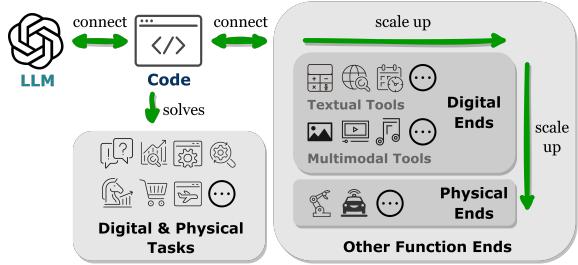


Figure 4: The code-centric tool-calling paradigm serves as a unified interface between LLMs and a large variety of functional ends, thus enabling many cross-modality and cross-domain tasks.

by pre-training an image-to-text model on masked website screenshots, and further training with OCR, language modeling, and image captioning objectives. Building on this, MATCHA (Liu et al., 2023a) introduced additional pre-training objectives based on chart derendering, layout comprehension, and mathematical reasoning, surpassing SOTA on VSNL understanding tasks.

4 Code Connects LLMs to Other Function Ends

Recent studies show that connecting LLMs to other functional ends (i.e., augmenting LLMs with external tools and execution modules) helps LLMs to perform tasks more accurately and reliably (Mialon et al., 2023; Parisi et al., 2022b; Peng et al., 2023; Gou et al., 2023). These functional ends empower LLMs to access external knowledge, engage with diverse modalities, and interact effectively with various environments. As indicated in Table 1, we observe a prevalent trend where LLMs generate programming languages or utilize pre-defined functions to establish connections with other functional ends—a phenomenon we refer to as the *code-centric paradigm*.

In contrast to the rigid practice of strictly hard-coding tool calls within the LLMs’ inference mechanism, the code-centric paradigm allows LLMs to dynamically generate tokens that invoke execution modules with adaptable parameters. This paradigm enables a simple and unambiguous way for LLMs to interact with other functional ends, enhancing the flexibility and scalability of their application. Importantly, as depicted in Figure 4, it allows LLMs to engage with numerous functional ends spanning diverse modalities and domains. By expanding both the quantity and variety of functional ends accessible, LLMs can handle more complicated tasks.

Table 1: Representative work connecting LLMs to different function ends for performing non-trivial tasks. Initial efforts embed tool calls rigidly within the LLMs’ inference mechanism (indicated by “*”), resulting in diminished flexibility and constrained tool accessibility. More recently, the *code-centric paradigm* establishes connections between LLMs and function ends through programming languages or pre-defined functions (indicated by “†”). This approach enhances the scalability of LLMs’ function end invocation across diverse tools and execution modules.

Major Type of Function Ends	Representative Work	Connecting Paradigm	Learning Method	Objectives or Problems to Solve
Single Tool	Retriever in REALM (Guu et al., 2020) Verifier in GSM8K (Cobbe et al., 2021)	Hardcoded in Inference Mechanism*	Example Fine-tuning Example Fine-tuning	Augment LLMs with Tools
Limited Text-based Tools	Blenderbot3 (Shuster et al., 2022) LamDA (Thoppilan et al., 2022)	Hardcoded in Inference Mechanism*	Example Fine-tuning Generate Pre-defined Functions†	Open-domain Conversation
Text-based Tools	TALM (Parisi et al., 2022a) ToolFormer (Schick et al., 2023)	Generate Pre-defined Functions† Generate Pre-defined Functions†	Iterative Self-play Self-supervised Training	Efficient and Generalizable Tool Using
Multi-modal Modules	MM-React (Yang et al., 2023) CodeVQA (Subramanian et al., 2023) VISPROG (Gupta and Kembhavi, 2023) ViperGPT (Surís et al., 2023)	Generate Pre-defined Functions† Generate Python Functions† Generate Python Functions† Generate Python Functions†	Zero-shot Prompting Zero-shot & Few shot Zero-shot Prompting Zero-shot Prompting	Multi-modal Reasoning Tasks
Real-World APIs	Code as Policies (Liang et al., 2023a) Proggrompt (Singh et al., 2022) SayCan (Ahn et al., 2022) RRR (Cui et al., 2023a) Agent-Driver (Mao et al., 2023) LaMPilot (Ma et al., 2023b)	Generate Python Functions† Generate Python Functions† Generate Pre-defined Functions† Generate Pre-defined Functions† Generate Pre-defined Functions† Generate Python Functions†	Few-shot Prompting Zero-shot Prompting Zero-shot Prompting Zero-shot Prompting	Better Robot Control
			Zero-shot Prompting Few-shot Prompting Zero-shot & Few-shot	Autonomous Driving Ecosystems

In §4.1, we examine (digital) textual and multi-modal tools connected to LLMs, while §4.2 focuses on physical-world functional ends, including robots and autonomous driving, showcasing the versatility of LLMs in tackling problems across various modalities and domains.

4.1 Relate LLMs to Digital Ends

Text-Based Tools. The code-centric framework has enhanced precision and clarity to LLMs’ tool invocation, initially driving progress in text-based tools. Prior to the popularity of this code-centric paradigm, research on augmenting LMs with single tools like information retrievers (Guu et al., 2020; Lewis et al., 2020; Izacard et al., 2022; Borgeaud et al., 2022; Yasunaga et al., 2022) required a hardcoded-in-inference-mechanism (e.g. always calling a retriever before the generation starts), which was less flexible and harder to scale. TALM (Parisi et al., 2022a) first incorporates multiple text-based tools by invoking API calls with a pre-defined delimiter, enabling unambiguous calls to any text-based tools at any position of generation. Following their work, Toolformer (Schick et al., 2023) marks API calls with “<API> </API>” along with their enclosed contents. Later, diverse tool-learning approaches were introduced to facilitate the integration of numerous text-based tools across various foundational models (Song et al., 2023; Hao et al., 2023; Tang et al., 2023).

The code-centric framework facilitates the invocation of a diverse range of external text modules. These include calculators, calendars, machine translation systems, web navigation tools, as well as APIs from HuggingFace and TorchHub (Thop-

pilan et al., 2022; Yao et al., 2022c; Shuster et al., 2022; Jin et al., 2023; Yao et al., 2022b; Liu et al., 2023e; Jin et al., 2023; Patil et al., 2023).

Multimodal Tools. The high scalability of the code-centric LLM paradigm enables the extension of tool-learning to modalities other than text. Early work (Gupta and Kembhavi, 2023; Surís et al., 2023; Subramanian et al., 2023) use the code-centric paradigm to tackle the visual question answering task. For instance, VISPROG (Gupta and Kembhavi, 2023) curates various pretrained computer vision models and functions from existing image processing libraries (e.g. Pillow and OpenCV) as a set of APIs. The API calls can then be chained together as a program for question-targeted image understanding, where the program is generated via in-context learning with LLMs. Containing arithmetic in its API code language, the program is capable of performing simple arithmetic tasks, thus enabling VISPROG to answer concrete questions such as object counting. Similar work includes ViperGPT (Surís et al., 2023) and CodeVQA (Subramanian et al., 2023). Compared to VISPROG, they directly generate more flexible Python code using Codex. This enables them to potentially generate programs of more complex control flows using the pre-trained knowledge embedded in Codex. In addition to visual reasoning, code has also been used to connect LLMs with multi-modal generative tools in image generation tasks (Cho et al., 2023; Feng et al., 2023; Wu et al., 2023a), where code’s unambiguous nature is leveraged in generating images that better match their text prompts.

Beyond the image-based tools, other modalities

have been considered and used in a collaborative fashion by recent work (Shen et al., 2023; Yang et al., 2023; Liang et al., 2023d). For example, MM-REACT (Yang et al., 2023) considers video recognition models in their API, and Chameleon (Lu et al., 2023) includes tools such as visual text detector or web search. In HuggingGPT (Shen et al., 2023), the authors connect LLMs to various Hugging Face models and treat each model as an available API call. As a result, HuggingGPT is capable of performing an even wider range of tasks, such as audio-related tasks, that were previously unexplored. Pushing the API diversity further, TaskMatrix.AI (Liang et al., 2023d) uses a magnitude higher number of APIs, spanning from visual & figure APIs to music and game APIs. The flexibility of code facilitates LLMs to jointly use different multimodal tools. This makes LLMs more versatile and capable of acting as general-purpose multimodal problem solvers that can scale to many tasks.

4.2 Relate LLMs to Physical Ends

While the physical world offers a more immersive, contextually rich, and engaging interactive environment compared to the digital realm, the connection between LLMs and the physical world has been constrained until the advent of the code-centric paradigm. This paradigm allows for adaptable calls to tools and execution modules in the physical world, first sparking a wave of research exploring the integration of LLMs with robotics and autonomous driving.

One of the most successful approaches to employing LLMs to generate policy codes for real-world robotic tasks is PaLM-SayCan (Ahn et al., 2022), where LLMs comprehend high-level instructions and then call corresponding APIs for robotic control. Following SayCan, recent developments have shown that LLMs can serve as the brain for robotics planning and control through their powerful code generation capabilities. Prog-Prompt (Singh et al., 2022), for instance, pioneered program-like specifications for robot task planning, while other researchers like Huang et al. (2023a), Liang et al. (2023a), and Vemprala et al. (2023) have extended this approach to a range of other tasks, including human-robot interaction and drone control.

Through code generation and tool learning, LLMs also show great potential in more compli-

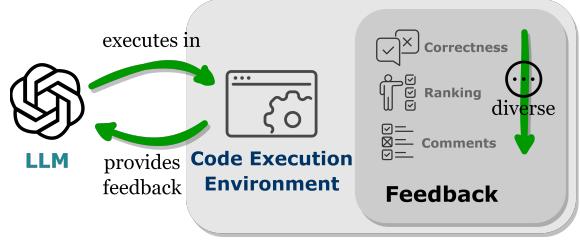


Figure 5: LLMs can be embedded into a code execution environment, where they collect faithful, automatic, and customizable feedback for self-improvement.

cated tasks such as human-vehicle interactions and autonomous driving (Cui et al., 2023b; Huang et al., 2023b; Li et al., 2023d). A prime tool learning example from the industry is Wayve’s LINGO-1 (Wayve, 2023), which uses an open-loop vision, language, and action LLM to improve the explainability of driving models. Using instruction tuning, LLMs have advanced to the point where they can understand complex driving, navigation, and entertainment commands (Wang et al., 2023e), generate actionable codes (Ma et al., 2023b), and execute them by calling low-level vehicle planning and control APIs (Cui et al., 2023a; Sha et al., 2023; Mao et al., 2023).

Overall, despite challenges such as latency, accuracy issues, and the absence of adequate simulation environments, datasets, and benchmarks (Kannan et al., 2023; Chen and Huang, 2023; Cui et al., 2023b), LLMs show promise in understanding high-level instructions and executing code-related APIs in intricate domains like robotics and autonomous driving. Looking ahead, there’s considerable potential for LLMs to bridge the gap between physical worlds and AI, influencing areas like transportation and smart manufacturing (Panchal and Wang, 2023; Zeng et al., 2023a).

5 Code Provides LLM with an Executable Environment for Automated Feedback

LLMs demonstrate performance beyond the parameters of their training, in part due to their ability to intake feedback, especially in real-world applications where environments are rarely static (Liu et al., 2023f; Wang et al., 2023d). However, feedback must be chosen carefully as noisy prompts can hinder LMs’ performance on downstream tasks (Zheng and Saparov, 2023). Furthermore, as human effort is expensive, it is crucial that feedback can be automatically collected while staying faithful. Embedding LLMs into a code execution envi-

ronment enables automated feedback that fulfills all of these criteria, as shown in Figure 5. As the code execution is largely deterministic, LLMs that intake feedback from the results of executed code remain faithful to the task at hand (Chen et al., 2023a; Fernandes et al., 2023; Scheurer et al., 2022). Furthermore, code interpreters provide an automatic pathway for LLMs to query internal feedback, eliminating the need for costly human annotations as seen when leveraging LLMs to debug or optimize faulty code (Chen et al., 2023a; Fernandes et al., 2023; Scheurer et al., 2022). In addition, code environments allow LLMs to incorporate diverse and comprehensive forms of external feedback, such as critic on binary correctness (Wang et al., 2023g), natural language explanations on results (Chen et al., 2023c), and ranking with reward values (Inala et al., 2022), enabling highly customizable methods to enhance performance.

We introduce the various types of feedback derived from code execution that can be jointly utilized to benefit LLMs in §5.1, and discuss common methods for utilizing this feedback to improve LLMs in §5.2.

5.1 Various Feedback from Code Execution

The code execution process enables assessing LLM-generated content with more comprehensive evaluation metrics derived from deterministic execution results, as opposed to relying solely on often ambiguous sequence-based metrics like BLEU (Papineni et al., 2002; Ren et al., 2020) and Rouge (Lin, 2004).

Straightforward methods for evaluating program execution outcomes and generating feedback include the creation of unit tests (Chen et al., 2021; Hendrycks et al., 2021; Austin et al., 2021; Li et al., 2022; Huang et al., 2022a; Lai et al., 2023) and the application of exact result matching techniques (Dong and Lapata, 2016; Zhong et al., 2017; Huang et al., 2022a). From these, feedback can be provided in two primary forms: simple correctness feedback and textual feedback. Simple feedback, indicating whether a program is correct or not, can be generated through critic models or rule-based methods (Wang et al., 2023g; Chen et al., 2023c).

For more detailed textual feedback, language models can be employed to produce explanations either about the program itself (Chen et al., 2023c), or to summarize comments on the program and its execution (Wang et al., 2023g; Chen et al., 2023c;

Zhang et al., 2023a). Execution results can also be translated into reward functions using predefined rules. The rules map execution results into scalar values based on the severity of different error types, thus making the feedback format suitable for reinforcement learning approaches (Le et al., 2022). Moreover, additional feedback can be extracted by performing static analysis using software engineering tools. For instance, Wang et al. (2017) and Gupta et al. (2017) obtain extra information from the execution trace, including variable or state traces. Lai et al. (2023) demonstrates an effective way to extract extra feedback using surface-form constraints on function calls.

5.2 Methods for Enhancing LLM’s Performance with Feedback

The feedback derived from code execution and external evaluation modules can enhance LLMs through three major approaches.

Selection Based Method. Selection-based methods, such as majority voting and re-ranking schemes, have proven effective in enhancing LLM performance in tasks such as code generation. These methods, originally developed for simple program synthesis, leverage code execution outcomes like the number of passed unit tests to choose the best-performing code snippet. Studies like Chen et al. (2018); Li et al. (2022) demonstrate the efficacy of majority voting, while Zhang et al. (2023b); Yin and Neubig (2019); Zeng et al. (2023b) showcase the advantages of re-ranking schemes. Building on this success, similar approaches have been adapted for more challenging tasks where code-LLMs are integrated in interactive environments, as shown in the work of Shi et al. (2022); Chen et al. (2022) for similar voting methods, and Ni et al. (2023); Inala et al. (2022); To et al. (2023) for re-ranking strategies. However, these approaches, while simple and effective, cause inefficiencies, as they necessitate multiple rounds of generation and the employment of additional re-ranking models to identify the optimal solution.

Prompting Based Method. Modern LLMs are equipped with the capability to reason in-context and directly integrate feedback from task descriptions into prompts, to certain extents. Improving LLM “self-debugging” with in-context learning prompts typically requires feedback presented as natural language explanations (Wang et al., 2023h; Chen et al., 2023c) or error messages derived from

execution results, as these formats are more comprehensible for the LLM. This method is favored by most LLM-based agents (see §6) due to its automatic nature, computational efficiency, and lack of requirement for additional fine-tuning. However, the effectiveness of this approach heavily depends on the LLM’s in-context learning capabilities.

Finetuning Based Method. In the aforementioned methods, neither the selection-based method nor the prompting-based method promises steady improvement over the task, as the LLMs’ parameters remain unchanged. They require repeating the tuning process even when faced with similar problems. Finetuning approaches, on the other hand, fundamentally improve the LLMs by updating their parameterized knowledge. Typical finetuning strategies include direct optimization, leveraging an external model for optimization, and training the model in a reinforcement learning paradigm. [Wang et al. \(2023g\)](#) exemplifies the direct optimization approach, where the original language model is fine-tuned with a feedback-conditioned objective. [Haluptzok et al. \(2022\)](#) presents a unique method where language models generate synthetic unit tests to identify and retain only correctly generated examples, which are then composed into correct question-answer pairs and employed to further fine-tune the models. CodeScore ([Dong et al., 2023a](#)) designs its loss function based on executability and the pass rate on the unit tests. For self-edit ([Zhang et al., 2023a](#)), it first wraps up execution results into textual comments, then trains an editor to further refine the program by accepting both the problematic program and the comments. [Chen et al. \(2023a\)](#) train a “refine model” which accepts the feedback and generated program first, then use the refined generated example to fine-tune the generation model, illustrating a layered approach to fine-tuning. CodeRL ([Le et al., 2022](#)) and [Wang et al. \(2022\)](#) apply reinforcement learning to improve the original language model. While [Wang et al. \(2022\)](#) aims at employing compiler feedback to obtain erroneous code, CodeRL ([Le et al., 2022](#)) empirically defines fixed reward values for different execution result types based on unit tests. Despite the discussed advantages, refining LLMs through finetuning typically involves a resource-intensive data collection process. Additionally, assessing predefined reward values, as exemplified in CodeRL ([Le et al., 2022](#)), poses certain challenges.

6 Application: Code-empowered LLMs Facilitate Intelligent Agents

In the preceding sections, our discussion highlighted the various ways in which code integration enhances LLMs. Going beyond, we discern that the benefits of code-empowered LLMs are especially pronounced in one key area: the development of IAs ([Xu et al., 2023](#)). In this section, we underscore the unique capabilities bestowed upon these agents by code-empowered LLMs.

Figure 6 helps to illustrate a standard operational pipeline of an IA, specifically serving as an embodied general daily assistant. We observe that the improvements brought about by code training in LLMs are firmly rooted in their practical operational steps when serving as IAs. These steps include (i) enhancing the IA’s decision-making in terms of environment perception and planning (§6.1), (ii) streamlining execution by grounding actions in modular and explicit action primitives and efficiently organizing memory (§6.2), and (iii) optimizing performance through feedback automatically derived from the code execution environment (§6.3). Detailed explanations of each aspect will follow in the subsequent sections.

6.1 Decision Making

Environment Perception As depicted in Figure 6 at step (0-10), the IA continuously perceives the world, engaging in interactions with humans and the environment, responding to relevant stimuli (e.g., human instructions for meal preparation), and planning and executing actions based on the observed environmental conditions (e.g., the kitchen layout). Utilizing LLMs as decision centers for IAs requires translating environmental observations into text, such as tasks based in the virtual household or Minecraft ([Shridhar et al., 2020; Côté et al., 2018; Wang et al., 2023b; Zhu et al., 2023](#)). The perceived information needs to be organized in a highly structured format, ensuring that stimuli occurring at the same moment (e.g., coexisting multimodality stimuli) influence the IA’s perception and decision simultaneously without temporal differences, a requirement that contrasts with the sequential nature of free-form text. Through pre-training on code, LLMs acquire the ability to better comprehend and generate structured representations resembling class definitions in code. This structured format, where class attributes and functions are permutation invariant, facilitates agents in

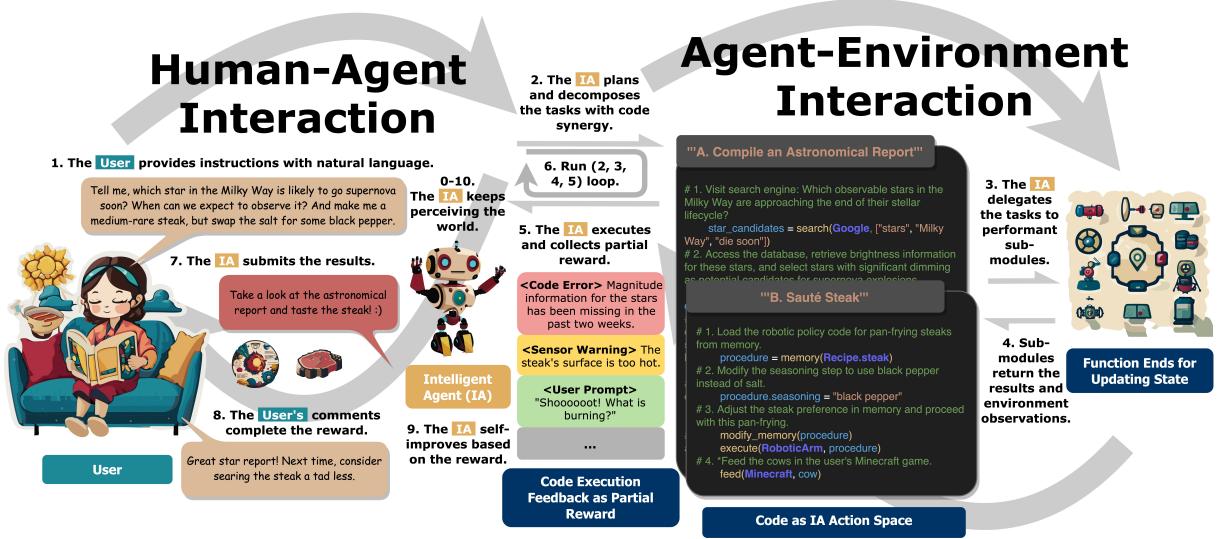


Figure 6: This figure illustrates the complete working pipeline of a LLM-based intelligent agent, mapping code-LLM abilities to specific phases: code-based planning in step (2), modular action parsing and tool creation in step (3), and automated feedback collection for enhanced agent self-improvement in step (5). Collectively, steps 0-10 in the entire loop benefit from code-LLMs' improved structured information understanding and perception.

perceiving structured environmental observations during task execution.

One such intuitive example is web-page-based environments which are highly structured around HTML code. In agent tasks like web shopping (Yao et al., 2022a), web browsing (Deng et al., 2023b), and web-based QA (Nakano et al., 2021; Liu et al., 2023d), it is preferred to translate the web-based environment into HTML code rather than natural language, directly encompassing its structural information and thereby improving the LLM agent’s overall perception. Moreover, in robotics research by Singh et al. (2022) and Liang et al. (2023a), the IAs are prompted with program-like specifications for objects in the environment, enabling the LLM to generate situated task plans based on the virtual objects they perceived.

Planning As illustrated in Figure 6 at step (2), IAs must break down intricate tasks into finer, manageable steps. Leveraging the synergized planning abilities of code-LLMs, IAs can generate organized reasoning steps using modular and unambiguous code alongside expressive natural language. As discussed in §2.2, when code-LLMs are employed for planning agent tasks, they exhibit enhanced reasoning capabilities. In addition, they generate the sub-tasks as executable programs when necessary, yielding more robust intermediate results, which the IA conditions on and refines its planning with greater precision. Furthermore, the IA seamlessly

integrates performant tool APIs into planning, addressing the limitations such as poor mathematical reasoning and outdated information updates faced by vanilla LLMs during planning.

Typical examples that utilize code for planning are in two main categories. Progprompt (Singh et al., 2022) and Code as Policies (Liang et al., 2023a) represent the work utilizing code for better robot control. Both work highlight the benefits brought by code-based planning as they not only enable direct expressions of feedback loops, functions, and primitive APIs, but also facilitate direct access to third-party libraries. Another stream of work is concerned with the scenario when the agents’ programming and mathematical reasoning abilities are crucial, like solving maths-related problems (Gao et al., 2023; Wang et al., 2023h) or doing experiments in the scientific domain (Boiko et al., 2023; Liffiton et al., 2023).

6.2 Execution

Action Grounding As depicted in Figure 6 at step (3), when the IA interfaces with external function ends according to the planning, it must invoke action primitives from a pre-defined set of actions (i.e., functions). Given that modern LLMs are trained in formal language and can generate highly formalized primitives, the IA’s generation can be directly parsed and executed, eliminating the necessity for additional action primitive grounding modules.

Connecting the IA with other function ends requires grounding actions into formalized function-like primitives. For instance, in a benchmark evaluating LLMs as agents in real-world scenarios (Liu et al., 2023f), seven out of eight scenarios involve code as the action space.

Previous work generating agent plans with pure natural language necessitate an additional step-to-primitive module to ground those planning steps into code (Wang et al., 2023c; Yin et al., 2023). In contrast, IAs that plan with code-LLMs generate atomic action programs (Yao et al., 2022d; Wang et al., 2023h; Liang et al., 2023a; Singh et al., 2022), and can have their generation quickly parsed for execution.

Memory Organization As depicted in Figure 6 at step (3) and the component labeled “Function Ends for Updating State,” the IA typically necessitates an external memory organization module to manage exposed information (Wang et al., 2023d), including original planning, task progress, execution history, available tool set, acquired skills, augmented knowledge, and users’ early feedback. In this context, Code-LLM aids the IA’s memory organization by employing highly abstract and modular code to record, organize, and access memory, especially for expanding the available tool set and manage acquired skills.

Typically, agent-written code snippets can serve as parts of the toolset, integrated into the memory organization of agents. This stream of research is known as tool creation approaches. TALM (Cai et al., 2023) proposes to use stronger agents (e.g. GPT-4 based agents) to write code as part of memory for weaker agents (e.g. GPT-3.5 based agents). In Creator (Qian et al., 2023b), agents themselves are highlighted as not only users of the tools but also their creators. They proposed a four-stage tool-creation framework that enables agents to write code as part of their executable tool set. Going further, Craft (Yuan et al., 2023) focuses on ensuring the created tools are indeed executable, making the framework more robust. Another work sharing this idea is Voyager (Wang et al., 2023b), in which the agent store learned skills in code format and execute them in the future when faced with similar tasks.

6.3 Self-improvement

As illustrated in Figure 6 at step (5), when the IA’s decision center, i.e., the LLM, operates within

a code execution environment, the environment can integrate various evaluation modules to offer automated feedback (e.g., correctness, ranking, detailed comments). This significantly enhances the IA’s early error correction and facilitates self-improvement.

Voyager (Wang et al., 2023b) is a good example for agents that use feedback from the simulated environment. The agent learns from failure task cases and further hone its skills in Minecraft. Chameleon (Lu et al., 2023) receives feedback from a program verifier to decide whether it should regenerate an appropriate program. Mint (Wang et al., 2023h) can receive feedback from proxies, and the agent can thus self-improve in a multi-turn interactive setting. Importantly, this ability to self-improve from execution feedback is fundamental for agents’ success at solving scientific problems (Bran et al., 2023; Swan et al., 2023; Wu et al., 2023b).

7 Challenges

We identify several intriguing and promising avenues for future research.

7.1 The Causality between Code Pre-training and LLMs’ Reasoning Enhancement

Although we have categorized the most pertinent work in §3.2, a noticeable gap persists in providing explicit experimental evidence that directly indicates the enhancement of LLMs’ reasoning abilities through the acquisition of specific code properties. While we intuitively acknowledge that certain code properties likely contribute to LLMs’ reasoning capabilities, the precise extent of their influence on enhancing reasoning skills remains ambiguous. In the future research endeavors, it is important to investigate whether reinforcing these code properties within training data could indeed augment the reasoning capabilities of trained LLMs. If it is indeed the case, that pre-training on specific properties of code directly improves LLMs’ reasoning abilities, understanding this phenomenon will be key to further improving the complex reasoning capabilities of current models.

7.2 Acquisition of Reasoning Beyond Code

Despite the enhancement in reasoning achieved by pre-training on code, foundational models still lack the human-like reasoning abilities expected from a truly generalized artificial intelligence. Im-

portantly, beyond code, a wealth of other textual data sources holds the potential to bolster LLM reasoning abilities, where the intrinsic characteristics of code, such as its lack of ambiguity, executability, and logical sequential structure, offer guiding principles for the collection or creation of these datasets. However, if we stick to the paradigm of training language models on large corpora with the language modeling objective, it’s hard to envision a sequentially readable language that is more abstract, highly structured, and closely aligned with symbolic language than formal languages, exemplified by code, which are prevalent in a substantial digital context. We envision that exploring alternative data modalities, diverse training objectives, and novel architectures would present additional opportunities to further enhance the reasoning capabilities of these models.

7.3 Challenges of Applying Code-centric Paradigm

The primary challenge in LLMs using code to connect to different function ends is learning the correct invocation of numerous functions, including selecting the right function end and passing the correct parameters at an appropriate time. Even for simple tasks like simplified web navigation with a limited set of action primitives like mouse movements, clicks, and page scrolls, few shot examples together with a strong underlying LLM are often required for the LLM to precisely grasp the usage of these primitives (Sridhar et al., 2023). For more complex tasks in data-intensive fields like chemistry (Bran et al., 2023), biology, and astronomy, which involve domain-specific Python libraries with diverse functions and intricate calls, enhancing LLMs’ capability of learning the correct invocation of these functions is a prospective direction, empowering LLMs to act as IAs performing expert-level tasks in fine-grained domains.

7.4 Learning from multi-turn interactions and feedback

LLMs often require multiple interactions with the user and the environment, continuously correcting themselves to improve intricate task completion (Li et al., 2023c). While code execution offers reliable and customizable feedback, a perfect method to fully leverage this feedback has yet to be established. As discussed in § 5.2, we observed that selection-based methods, though useful, do not guarantee improved performance and can be inef-

ficient. Prompting-based methods heavily depend on the in-context learning abilities of the LLM, which might limit their applicability. Fine-tuning methods show consistent improvement, but data collection and fine-tuning are resource-intensive and thus prohibitive. We hypothesize that reinforcement learning could be a more effective approach for utilizing feedback and improving LLMs. This method can potentially address the limitations of current techniques by providing a dynamic way to adapt to feedback through well-designed reward functions. However, significant research is still needed to understand how reward functions should be designed and how reinforcement learning can be optimally integrated with LLMs for complex task completion.

8 Conclusion

In this survey, we compile literature that elucidates how code empowers LLMs, as well as where code assists LLMs to serve as IAs. To begin with, code possesses natural language’s sequential readability while also embodying the abstraction and graph structure of symbolic representations, rendering it a conduit for knowledge perception and reasoning as an integral part of the LLMs’ training corpus based on the mere language modeling objective. Through a comprehensive literature review, we observe that after code training, LLMs *i*) improve their programming skills and reasoning, *ii*) could generate highly formalized functions, enabling flexible connections to diverse functional ends across modalities and domains, and *iii*) engage in interaction with evaluation modules integrated in the code execution environment for automated self-improvement. Moreover, we find that the LLMs’ capability enhancement brought by code training benefits their downstream application as IAs, manifesting in the specific operational steps of the IAs’ workflow regarding decision-making, execution, and self-improvement. Beyond reviewing prior research, we put forth several challenges in this field to serve as guiding factors for potential future directions.

References

- Rajas Agashe, Srini Iyer, and Luke Zettlemoyer. 2019. *Juice: A large scale distantly supervised dataset for open domain context-based code generation*. In *Conference on Empirical Methods in Natural Language Processing*.

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. 2022. [Do as i can, not as i say: Grounding language in robotic affordances](#).
- Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, Logesh Kumar Umapath, Carolyn Jane Anderson, Yangtian Zi, Joel Lamy Poirier, Hailey Schoelkopf, Sergey Troshin, Dmitry Abulkhanov, Manuel Romero, Michael Lappert, Francesco De Toni, Bernardo García del Río, Qian Liu, Shamik Bose, Urvashi Bhattacharyya, Terry Yue Zhuo, Ian Yu, Paulo Villegas, Marco Zocca, Sourab Mangrulkar, David Lansky, Huu Nguyen, Danish Contractor, Luis Villa, Jia Li, Dzmitry Bahdanau, Yacine Jernite, Sean Hughes, Daniel Fried, Arjun Guha, Harm de Vries, and Leandro von Werra. 2023. [Santacoder: don't reach for the stars!](#)
- Uri Alon, Roy Sadaka, Omer Levy, and Eran Yahav. 2020. [Structural language models of code](#).
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. [Program synthesis with large language models](#).
- Tony Beltramelli. 2017. [pix2code: Generating code from a graphical user interface screenshot](#).
- Daniil A. Boiko, Robert MacKnight, and Gabe Gomes. 2023. [Emergent autonomous scientific research capabilities of large language models](#). ArXiv, abs/2304.05332.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Milligan, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego De Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack Rae, Erich Elsen, and Laurent Sifre. 2022. [Improving language models by retrieving from trillions of tokens](#). In [Proceedings of the 39th International Conference on Machine Learning](#), volume 162 of [Proceedings of Machine Learning Research](#), pages 2206–2240. PMLR.
- Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. 2023. [Chemcrow: Augmenting large-language models with chemistry tools](#).
- Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. 2023. [Large language models as tool makers](#). ArXiv, abs/2305.17126.
- Tommaso Calò and Luigi De Russis. 2023. Leveraging large language models for end-user website generation. In [International Symposium on End User Development](#), pages 52–61. Springer.
- Angelica Chen, Jérémie Scheurer, Tomasz Korbak, Jon Ander Campos, Jun Shern Chan, Samuel R. Bowman, Kyunghyun Cho, and Ethan Perez. 2023a. [Improving code generation by training with natural language feedback](#).
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. [Codet: Code generation with generated tests](#).
- Juo-Tung Chen and Chien-Ming Huang. 2023. [Forgetful large language models: Lessons learned from using llms in robot programming](#).
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. [Evaluating large language models trained on code](#).
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023b. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#).
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023c. [Teaching large language models to self-debug](#).
- Xinyun Chen, Chang Liu, and Dawn Song. 2018. Execution-guided neural program synthesis. In [International Conference on Learning Representations](#).
- Yangyi Chen, Xingyao Wang, Manling Li, Derek Hoiem, and Heng Ji. 2023d. [Vistruct: Visual structural knowledge extraction via curriculum guided code-vision representation](#).

- Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2023. [Binding language models in symbolic languages](#).
- Jaemin Cho, Abhay Zala, and Mohit Bansal. 2023. [Visual programming for text-to-image generation and evaluation](#).
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#).
- Marc-Alexandre Côté, Ákos Kádár, Xingdi Yuan, Ben A. Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew J. Hausknecht, Layla El Asri, Mahmoud Adada, Wendy Tay, and Adam Trischler. 2018. [Textworld: A learning environment for text-based games](#). [ArXiv](#), abs/1806.11532.
- Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, and Ziran Wang. 2023a. [Receive, reason, and react: Drive as you say with large language models in autonomous vehicles](#).
- Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, Tianren Gao, Erlong Li, Kun Tang, Zhipeng Cao, Tong Zhou, Ao Liu, Xinrui Yan, Shuqi Mei, Jianguo Cao, Ziran Wang, and Chao Zheng. 2023b. [A survey on multimodal large language models for autonomous driving](#).
- Gelei Deng, Yi Liu, Víctor Mayoral-Vilches, Peng Liu, Yuekang Li, Yuan Xu, Tianwei Zhang, Yang Liu, Martin Pinzger, and Stefan Rass. 2023a. [Pentest-gpt: An llm-empowered automatic penetration testing tool](#).
- Xiang Deng, Yu Gu, Bo Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023b. [Mind2web: Towards a generalist agent for the web](#). [ArXiv](#), abs/2306.06070.
- Li Dong and Mirella Lapata. 2016. [Language to logical form with neural attention](#).
- Yihong Dong, Ji Ding, Xue Jiang, Zhuo Li, Ge Li, and Zhi Jin. 2023a. [Codescore: Evaluating code generation by learning code execution](#). [ArXiv](#), abs/2301.09043.
- Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2023b. [Self-collaboration code generation via chatgpt](#).
- Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, Roman Wang, Nikhil Singh, Taylor L. Patti, Jayson Lynch, Avi Shporer, Nakul Verma, Eugene Wu, and Gilbert Strang. 2022. [A neural network solves, explains, and generates university math problems by program synthesis](#) and few-shot learning at human level. [Proceedings of the National Academy of Sciences](#), 119(32).
- Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. 2023. [Classeval: A manually-crafted benchmark for evaluating llms on class-level code generation](#).
- Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt Gardner. 2022. [Successive prompting for decomposing complex questions](#).
- Kevin Ellis, Maxwell Nye, Yewen Pu, Felix Sosa, Josh Tenenbaum, and Armando Solar-Lezama. 2019. [Write, execute, assess: Program synthesis with a repl](#). [Advances in Neural Information Processing Systems](#), 32.
- Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M. Zhang. 2023. [Large language models for software engineering: Survey and open problems](#).
- Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. 2023. [Layoutgpt: Compositional visual planning and generation with large language models](#). [arXiv preprint arXiv:2305.15393](#).
- Patrick Fernandes, Aman Madaan, Emmy Liu, António Farinhos, Pedro Henrique Martins, Amanda Bertsch, José GC de Souza, Shuyan Zhou, Tongshuang Wu, Graham Neubig, et al. 2023. [Bridging the gap: A survey on integrating \(human\) feedback for natural language generation](#). [arXiv preprint arXiv:2305.00955](#).
- Hao Fu, Yao; Peng and Tushar Khot. 2022. [How does gpt obtain its ability? tracing emergent abilities of language models to their sources](#). [Yao Fu's Notion](#).
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. 2023. [Complexity-based prompting for multi-step reasoning](#).
- Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane Gu, and Izzeddin Gur. 2023. [Multimodal web navigation with instruction-finetuned foundation models](#).
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. [Pal: Program-aided language models](#). In [International Conference on Machine Learning](#), pages 10764–10799. PMLR.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. [Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies](#). [Transactions of the Association for Computational Linguistics](#), 9:346–361.

- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 2023. *CRITIC: Large language models can self-correct with tool-interactive critiquing*.
- Rahul Gupta, Soham Pal, Aditya Kanade, and Shirish Shevade. 2017. Deepfix: Fixing common c language errors by deep learning. In *Proceedings of the aaai conference on artificial intelligence*, volume 31.
- Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14953–14962.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. *Retrieval augmented language model pre-training*. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3929–3938. PMLR.
- Patrick Haluptzok, Matthew Bowers, and Adam Tauman Kalai. 2022. Language models can teach themselves to program better. [arXiv preprint arXiv:2207.14502](#).
- Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. 2023. *ToolkenGPT: Augmenting frozen language models with massive tools via tool embeddings*.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. 2021. Measuring coding challenge competence with apps. [arXiv preprint arXiv:2105.09938](#).
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. [arXiv preprint arXiv:2308.00352](#).
- Junjie Huang, Chenglong Wang, Jipeng Zhang, Cong Yan, Haotian Cui, Jeevana Priya Inala, Colin Clement, Nan Duan, and Jianfeng Gao. 2022a. Execution-based evaluation for data science code generation models. [arXiv preprint arXiv:2211.09374](#).
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022b. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pages 9118–9147. PMLR.
- Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. 2023a. Voxposer: Composable 3d value maps for robotic manipulation with language models. [arXiv preprint arXiv:2307.05973](#).
- Yu Huang, Yue Chen, and Zhu Li. 2023b. Applications of large scale foundation models for autonomous driving. [arXiv preprint arXiv:2311.12144](#).
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709.
- Jeevana Priya Inala, Chenglong Wang, Mei Yang, Andres Codas, Mark Encarnación, Shuvendu Lahiri, Madanlal Musuvathi, and Jianfeng Gao. 2022. Fault-aware neural code rankers. *Advances in Neural Information Processing Systems*, 35:13419–13432.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. *Atlas: Few-shot learning with retrieval augmented language models*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. *Survey of hallucination in natural language generation*. *ACM Computing Surveys*, 55(12):1–38.
- Xue Jiang, Yihong Dong, Lecheng Wang, Qiwei Shang, and Ge Li. 2023. Self-planning code generation with large language model. [arXiv preprint arXiv:2303.06689](#).
- Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu. 2023. *GeneGPT: Augmenting large language models with domain tools for improved access to biomedical information*.
- Sungmin Kang, Gabin An, and Shin Yoo. 2023a. A preliminary evaluation of llm-based fault localization. [arXiv preprint arXiv:2308.05487](#).
- Sungmin Kang, Juyeon Yoon, and Shin Yoo. 2023b. Large language models are few-shot testers: Exploring llm-based general bug reproduction. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2312–2323. IEEE.
- Shyam Sundar Kannan, Vishnunandan LN Venkatesh, and Byung-Cheol Min. 2023. Smart-llm: Smart multi-agent robot task planning using large language models. [arXiv preprint arXiv:2309.10062](#).
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2023. Ds-1000: A natural and reliable benchmark for data science code generation. In *International Conference on Machine Learning*, pages 18319–18345. PMLR.
- Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi. 2022. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35:21314–21328.

- Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023. [Pix2struct: Screenshot parsing as pretraining for visual language understanding](#).
- Ioktong Lei and Zhidong Deng. 2023. [Selfzcot: a self-prompt zero-shot cot from semantic-level to code-level for a better utilization of llms](#).
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). In [Advances in Neural Information Processing Systems](#), volume 33, pages 9459–9474. Curran Associates, Inc.
- Haonan Li, Yu Hao, Yizhuo Zhai, and Zhiyun Qian. 2023a. The hitchhiker’s guide to program analysis: A journey with large language models. [arXiv preprint arXiv:2308.00245](#).
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. [arXiv preprint arXiv:2301.12597](#).
- Sha Li, Chi Han, Pengfei Yu, Carl Edwards, Manling Li, Xingyao Wang, Yi Fung, Charles Yu, Joel Tetreault, Eduard Hovy, and Heng Ji. 2023c. [Defining a new NLP playground](#). In [Findings of the Association for Computational Linguistics: EMNLP 2023](#), pages 11932–11951, Singapore. Association for Computational Linguistics.
- Xin Li, Yeqi Bai, Pinlong Cai, Licheng Wen, Daocheng Fu, Bo Zhang, Xuemeng Yang, Xinyu Cai, Tao Ma, Jianfei Guo, et al. 2023d. Towards knowledge-driven autonomous driving. [arXiv preprint arXiv:2312.04316](#).
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. [Science](#), 378(6624):1092–1097.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. 2023a. Code as policies: Language model programs for embodied control. In [2023 IEEE International Conference on Robotics and Automation \(ICRA\)](#), pages 9493–9500. IEEE.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher R’e, Diana Acosta-Navas, Drew A. Hudson, E. Zelikman, Esin Durmus, Faisal Ladha, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan S. Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas F. Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023b. [Holistic evaluation of language models](#). [Annals of the New York Academy of Sciences](#), 1525:140 – 146.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladha, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023c. [Holistic evaluation of language models](#).
- Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, et al. 2023d. Taskmatrix.ai: Completing tasks by connecting foundation models with millions of apis. [arXiv preprint arXiv:2303.16434](#).
- Mark H. Liffiton, Brad E. Sheese, Jaromir Savelka, and Paul Denny. 2023. [Codehelp: Using large language models with guardrails for scalable support in programming classes](#). [ArXiv](#), abs/2308.06921.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In [Text summarization branches out](#), pages 74–81.
- Jiawei Lin, Jiaqi Guo, Shizhao Sun, Weijiang Xu, Ting Liu, Jian-Guang Lou, and Dongmei Zhang. 2023. [A parse-then-place approach for generating graphic layouts from textual descriptions](#).
- Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Martin Eisenschlos. 2023a. [Matcha: Enhancing visual language pretraining with math reasoning and chart derendering](#).
- Jiateng Liu, Sha Li, Zhenhailong Wang, Manling Li, and Heng Ji. 2023b. A language-first approach for procedure planning. In [Findings of the Association for Computational Linguistics: ACL 2023](#), pages 1941–1954.
- Michael Xieyang Liu, Advait Sarkar, Carina Negreanu, Benjamin Zorn, Jack Williams, Neil Toronto, and Andrew D Gordon. 2023c. “what it wants me to say”:

- Bridging the abstraction gap between end-user programmers and code-generating large language models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–31.
- Xiao Liu, Hanyu Lai, Hao Yu, Yifan Xu, Aohan Zeng, Zhengxiao Du, P. Zhang, Yuxiao Dong, and Jie Tang. 2023d. *Webglm: Towards an efficient web-enhanced question answering system with human preferences*. *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- Xiao Liu, Hanyu Lai, Hao Yu, Yifan Xu, Aohan Zeng, Zhengxiao Du, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023e. *WebGLM: Towards an efficient web-enhanced question answering system with human preferences*.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuyanyu Lei, Hanyu Lai, Yu Gu, Yuxian Gu, Hangliang Ding, Kai Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Shengqi Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 2023f. *Agentbench: Evaluating llms as agents*. *ArXiv*, abs/2308.03688.
- Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. 2023g. Dynamic llm-agent network: An llm-agent collaboration framework with agent team optimization. *arXiv preprint arXiv:2310.02170*.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023. Chameleon: Plug-and-play compositional reasoning with large language models. *arXiv preprint arXiv:2304.09842*.
- Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Dixin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. 2021. *Codexglue: A machine learning benchmark dataset for code understanding and generation*.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning. *arXiv preprint arXiv:2301.13379*.
- Yingwei Ma, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. 2023a. At which training stage does code data help llms reasoning?
- Yunsheng Ma, Can Cui, Xu Cao, Wenqian Ye, Peiran Liu, Juanwu Lu, Amr Abdelraouf, Rohit Gupta, Kyungtae Han, Aniket Bera, et al. 2023b. Lampilot: An open benchmark dataset for autonomous driving with language model programs. *arXiv preprint arXiv:2312.04372*.
- Aman Madaan, Shuyan Zhou, Uri Alon, Yiming Yang, and Graham Neubig. 2022. Language models of code are few-shot commonsense learners.
- Jiageng Mao, Junjie Ye, Yuxi Qian, Marco Pavone, and Yue Wang. 2023. A language agent for autonomous driving. *arXiv preprint arXiv:2311.10813*.
- Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. Augmented language models: a survey. *ArXiv*, abs/2302.07842.
- Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Tafjord, Ashish Sabharwal, Peter Clark, and Ashwin Kalyan. 2023. Lila: A unified benchmark for mathematical reasoning.
- Mohammad Mahdi Mohajer, Reem Aleithan, Nima Shiri Harzevili, Moshi Wei, Alvine Boaye Belle, Hung Viet Pham, and Song Wang. 2023. Skipanalyzer: An embodied agent for code analysis with large language models. *arXiv preprint arXiv:2310.18532*.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*.
- Ansong Ni, Srini Iyer, Dragomir Radev, Veselin Stoyanov, Wen-tau Yih, Sida Wang, and Xi Victoria Lin. 2023. Lever: Learning to verify language-to-code generation with execution. In *International Conference on Machine Learning*, pages 26106–26128. PMLR.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2022. Codegen: An open large language model for code with multi-turn program synthesis. *arXiv preprint arXiv:2203.13474*.
- David Noever. 2023. Can large language models find and fix vulnerable software? *arXiv preprint arXiv:2308.10345*.
- OpenAI. 2023. Gpt-4 technical report.
- Jitesh H Panchal and Ziran Wang. 2023. Design of next-generation automotive systems: Challenges and research opportunities. *Journal of Computing and Information Science in Engineering*, 23(6).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

- Aaron Parisi, Yao Zhao, and Noah Fiedel. 2022a. **Talm: Tool augmented language models.** [ArXiv](#), abs/2205.12255.
- Aaron Parisi, Yao Zhao, and Noah Fiedel. 2022b. **TALM: Tool augmented language models.**
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. **Gorilla: Large language model connected with massive APIs.**
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. 2023. **Check your facts and try again: Improving large language models with external knowledge and automated feedback.**
- Stanislas Polu and Ilya Sutskever. 2020. **Generative language modeling for automated theorem proving.**
- Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023a. **Communicative agents for software development.** [arXiv preprint arXiv:2307.07924](#).
- Cheng Qian, Chi Han, Yi Ren Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. 2023b. **Creator: Tool creation for disentangling abstract and concrete reasoning of large language models.**
- Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, et al. 2023. **Tool learning with foundation models.** [arXiv preprint arXiv:2304.08354](#).
- Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. 2022. **Evaluating the text-to-sql capabilities of large language models.**
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. **In-Context retrieval-augmented language models.**
- Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. 2020. **Codebleu: a method for automatic evaluation of code synthesis.** [arXiv preprint arXiv:2009.10297](#).
- Jasmine Roberts, Andrzej Banburski-Fahey, Microsoft Jaron, and Lanier Microsoft. **Surreal vr pong: Llm approach to game design.**
- Alex Robinson. 2019. **Sketch2code: Generating a website from a paper mockup.**
- Jérémie Scheurer, Jon Ander Campos, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. 2022. **Training language models with natural language feedback.** [arXiv preprint arXiv:2204.14146](#), 8.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. **Toolformer: Language models can teach themselves to use tools.**
- Hao Sha, Yao Mu, Yuxuan Jiang, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi Tomizuka, Wei Zhan, and Mingyu Ding. 2023. **Languagempc: Large language models as decision makers for autonomous driving.** [arXiv preprint arXiv:2310.03026](#).
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueling Zhuang. 2023. **Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface.** [arXiv preprint arXiv:2303.17580](#).
- Freida Shi, Daniel Fried, Marjan Ghazvininejad, Luke Zettlemoyer, and Sida I Wang. 2022. **Natural language to code translation with execution.** [arXiv preprint arXiv:2204.11454](#).
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew J. Hausknecht. 2020. **Alfworld: Aligning text and embodied environments for interactive learning.** [ArXiv](#), abs/2010.03768.
- Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju, Eric Michael Smith, Stephen Roller, Megan Ung, Moya Chen, Kushal Arora, Joshua Lane, Morteza Behrooz, William Ngan, Spencer Poff, Naman Goyal, Arthur Szlam, Y-Lan Boureau, Melanie Kambadur, and Jason Weston. 2022. **BlenderBot 3: a deployed conversational agent that continually learns to responsibly engage.**
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2022. **Prog-prompt: Generating situated robot task plans using large language models.**
- Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, Ye Tian, and Sujian Li. 2023. **RestGPT: Connecting large language models with real-world RESTful APIs.**
- Davit Soselia, Khalid Saifullah, and Tianyi Zhou. 2023. **Learning ui-to-code reverse generator using visual critic without rendering.**
- Abishek Sridhar, Robert Lo, Frank F. Xu, Hao Zhu, and Shuyan Zhou. 2023. **Hierarchical prompting assists large language model on web navigation.**
- Sanjay Subramanian, Medhini Narasimhan, Kushal Khangaonkar, Kevin Yang, Arsha Nagrani, Cordelia Schmid, Andy Zeng, Trevor Darrell, and Dan Klein. 2023. **Modular visual question answering via code generation.** [arXiv preprint arXiv:2306.05392](#).
- Chunyi Sun, Junlin Han, Weijian Deng, Xinlong Wang, Zishan Qin, and Stephen Gould. 2023a. **3d-gpt: Procedural 3d modeling with large language models.** [arXiv preprint arXiv:2310.12945](#).

- Ruoxi Sun, Sercan O. Arik, Hootan Nakhost, Hanjun Dai, Rajarishi Sinha, Pengcheng Yin, and Tomas Pfister. 2023b. [Sql-palm: Improved large language model adaptation for text-to-sql](#).
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. [arXiv preprint arXiv:2303.08128](#).
- Melanie Swan, Takashi Kido, Eric Roland, and Renato P. dos Santos. 2023. [Math agents: Computational infrastructure, mathematical embedding, and genomics](#). [ArXiv](#), abs/2307.02502.
- Yashar Talebirad and Amirhossein Nadiri. 2023. Multi-agent collaboration: Harnessing the power of intelligent llm agents. [arXiv preprint arXiv:2306.03314](#).
- Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. 2023. [ToolAlpaca: Generalized tool learning for language models with 3000 simulated cases](#).
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, Yaguang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafoori, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny So-raker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. [LaMDA: Language models for dialog applications](#).
- Hung Quoc To, Minh Nguyen, and Nghi D. Q. Bui. 2023. [Neural rankers for code generation via inter-cluster modeling](#).
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Biket, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Bin Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#).
- Priyan Vaithilingam, Tianyi Zhang, and Elena L Glassman. 2022. Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In [Chi conference on human factors in computing systems extended abstracts](#), pages 1–7.
- Sai Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. 2023. Chatgpt for robotics: Design principles and model abilities. [Microsoft Auton. Syst. Robot. Res](#), 2:20.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023a. [Towards understanding chain-of-thought prompting: An empirical study of what matters](#).
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlikar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023b. Voyager: An open-ended embodied agent with large language models. [arXiv preprint arXiv:2305.16291](#).
- Huaxiaoyue Wang, Gonzalo Gonzalez-Pumariega, Yash Sharma, and Sanjiban Choudhury. 2023c. Demo2code: From summarizing demonstrations to synthesizing code via extended chain-of-thought. [arXiv preprint arXiv:2305.16744](#).
- Ke Wang, Rishabh Singh, and Zhendong Su. 2017. Dynamic neural program embedding for program repair. [arXiv preprint arXiv:1711.07163](#).
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2023d. A survey on large language model based autonomous agents. [arXiv preprint arXiv:2308.11432](#).
- Shiyi Wang, Yuxuan Zhu, Zhiheng Li, Yutong Wang, Li Li, and Zhengbing He. 2023e. Chatgpt as your vehicle co-pilot: An initial attempt. [IEEE Transactions on Intelligent Vehicles](#).
- Xin Wang, Yasheng Wang, Yao Wan, Fei Mi, Yitong Li, Pingyi Zhou, Jin Liu, Hao Wu, Xin Jiang, and Qun Liu. 2022. Compilable neural code generation with compiler feedback. [arXiv preprint arXiv:2203.05132](#).
- Xingyao Wang, Sha Li, and Heng Ji. 2023f. [Code4struct: Code generation for few-shot event structure prediction](#).
- Xingyao Wang, Hao Peng, Reyhaneh Jabbarvand, and Heng Ji. 2023g. Leti: Learning to generate from textual interactions. [arXiv preprint arXiv:2305.10314](#).

- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lisan Yuan, Hao Peng, and Heng Ji. 2023h. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. [arXiv preprint arXiv:2309.10691](#).
- Wayve. 2023. **LINGO-1: Exploring Natural Language for Autonomous Driving**.
- Colin Wei, Sang Michael Xie, and Tengyu Ma. 2022a. **Why do pretrained language models help in downstream tasks? an analysis of head and prompt tuning.**
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022b. **Emergent abilities of large language models.**
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. **Chain-of-thought prompting elicits reasoning in large language models.**
- Man-Fai Wong, Shangxin Guo, Ching-Nam Hang, Siu-Wai Ho, and Chee-Wei Tan. 2023. **Natural language generation and understanding of big code for AI-assisted programming: A review.** *Entropy*, 25(6):888.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023a. Visual chatgpt: Talking, drawing and editing with visual foundation models. [arXiv preprint arXiv:2303.04671](#).
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023b. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. [arXiv preprint arXiv:2308.08155](#).
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023c. **Autogen: Enabling next-gen llm applications via multi-agent conversation framework.** [ArXiv](#), abs/2308.08155.
- Yuhuai Wu, Albert Q. Jiang, Wenda Li, Markus N. Rabe, Charles Staats, Mateja Jamnik, and Christian Szegedy. 2022. **Autoformalization with large language models.**
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. 2023. **The rise and potential of large language model based agents: A survey.**
- Frank F Xu, Uri Alon, Graham Neubig, and Vincent Joshua Hellendoorn. 2022. A systematic evaluation of large language models of code. In *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, pages 1–10.
- Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, Zhoujun Cheng, Siheng Zhao, Lingpeng Kong, Bailin Wang, Caiming Xiong, and Tao Yu. 2023. **Lemur: Harmonizing natural language and code for language agents.**
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023. Mm-react: Prompting chatgpt for multimodal reasoning and action. [arXiv preprint arXiv:2303.11381](#).
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. **Webshop: Towards scalable real-world web interaction with grounded language agents.** [ArXiv](#), abs/2207.01206.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022b. **Webshop: Towards scalable real-world web interaction with grounded language agents.** In *Advances in Neural Information Processing Systems*, volume 35, pages 20744–20757. Curran Associates, Inc.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izak Shafran, Karthik Narasimhan, and Yuan Cao. 2022c. **ReAct: Synergizing reasoning and acting in language models.**
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izak Shafran, Karthik Narasimhan, and Yuan Cao. 2022d. **React: Synergizing reasoning and acting in language models.** [ArXiv](#), abs/2210.03629.
- Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-Tau Yih. 2022. **Retrieval-augmented multimodal language modeling.**
- Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhang Cui, Zeyang Zhou, Chao Gong, Yang Shen, et al. 2023a. A comprehensive capability analysis of gpt-3 and gpt-3.5 series models. [arXiv preprint arXiv:2303.10420](#).
- Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. 2023b. **Large language models are versatile decomposers: Decompose evidence and questions for table-based reasoning.**
- Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2023. Lumos: Towards language agents that are unified, modular, and open source.
- Pengcheng Yin and Graham Neubig. 2019. Reranking for neural semantic parsing. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.

- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2019. *Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task*.
- Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi Ren Fung, Hao Peng, and Heng Ji. 2023. *Craft: Customizing llms by creating and retrieving from specialized toolsets*. [ArXiv](#), abs/2309.17428.
- Daoguang Zan, Bei Chen, Fengji Zhang, Dianjie Lu, Bingchao Wu, Bei Guan, Wang Yongji, and Jian-Guang Lou. 2023. Large language models meet nl2code: A survey. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7443–7464.
- Fanlong Zeng, Wensheng Gan, Yongheng Wang, Ning Liu, and Philip S Yu. 2023a. Large language models for robotics: A survey. [arXiv preprint arXiv:2311.07226](#).
- Lu Zeng, Sree Hari Krishnan Parthasarathi, and Dilek Hakkani-Tur. 2023b. N-best hypotheses reranking for text-to-sql systems. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 663–670. IEEE.
- Wei Zeng, Xiaozhe Ren, Teng Su, Hui Wang, Yi Liao, Zhiwei Wang, Xin Jiang, ZhenZhang Yang, Kaisheng Wang, Xiaoda Zhang, Chen Li, Ziyan Gong, Yifan Yao, Xinjing Huang, Jun Wang, Jianfeng Yu, Qi Guo, Yue Yu, Yan Zhang, Jin Wang, Hengtao Tao, Dasen Yan, Zexuan Yi, Fang Peng, Fangqing Jiang, Han Zhang, Lingfeng Deng, Yehong Zhang, Zhe Lin, Chao Zhang, Shaojie Zhang, Mingyue Guo, Shanzhi Gu, Gaojun Fan, Yaowei Wang, Xuefeng Jin, Qun Liu, and Yonghong Tian. 2021. *Pangu- α : Large-scale autoregressive pretrained chinese language models with auto-parallel computation*.
- Kechi Zhang, Zhuo Li, Jia Li, Ge Li, and Zhi Jin. 2023a. Self-edit: Fault-aware code editor for code generation. [arXiv preprint arXiv:2305.04087](#).
- Tianyi Zhang, Tao Yu, Tatsunori Hashimoto, Mike Lewis, Wen-tau Yih, Daniel Fried, and Sida Wang. 2023b. Coder reviewer reranking for code generation. In *International Conference on Machine Learning*, pages 41832–41846. PMLR.
- Pengyu Zhao, Zijian Jin, and Ning Cheng. 2023. *An in-depth survey of large language model-based artificial intelligence agents*.
- Hongyi Zheng and Abulhair Saparov. 2023. Noisy exemplars make large language models more robust: A domain-agnostic behavioral analysis.
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. [arXiv preprint arXiv:1709.00103](#).
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, et al. 2023a. Solving challenging math word problems using gpt-4 code interpreter with code-based self-verification. [arXiv preprint arXiv:2308.07921](#).
- Xuanhe Zhou, Guoliang Li, and Zhiyuan Liu. 2023b. Llm as dba. [arXiv preprint arXiv:2308.05481](#).
- Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyuan Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, Y. Qiao, Zhaoxiang Zhang, and Jifeng Dai. 2023. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. [ArXiv](#), abs/2305.17144.
- Terry Yue Zhuo. 2023. Large language models are state-of-the-art evaluators of code generation. [arXiv preprint arXiv:2304.14317](#).

A Discussion

A.1 Intrinsic Qualities of Code that Contribute to LLM Empowerment

Reflecting on our definition of code in the introduction section (§1) as formal languages that are both human-interpretable and machine-executable, we highlight that while some features are shared by all code, programming language, as the most well-known and most established type of code, enjoy some unique advantages. In Figure 7, we provide a case study comparing code and natural language.

First, we talk about the core feature shared by all code within the range of our definition. The inherent nature of code is that they are explicit and have clear definitions for every single line, while natural language is generally in free form and can be very ambiguous. Consequently, code is significantly better at expressing detailed commands, signifying a specific step, and transmitting control signals. This generally led to the improvement in §4, the improvement for more controlled planning (cf. planning part in §6.1), and also helped with action execution (§6.2).

Programming languages, a critical component of the code family, are specifically designed for machine communication. Their advantages extend beyond mere explicitness and clarity. One overwhelming feature of programming languages (though some formal languages also define logical commands and loops) is that they contain structural definitions. Some well-known features are logical operands (If & Else), loops (For & While), nesting (within Functions), and even class definition

```

class IntelligentAgent
"""
An intelligent agent utilizing a decision center (default: LLM) and a toolbox of available tools.

Parameters:
- decision_center: The decision-making center for the agent, defaults to an LLM.
- toolbox: A list of available tools for the agent, default includes GOOGLE, Minecraft, and RoboticArm.

Methods:
- cs_ritual(): Executes a rookie ritual for a computer science rookie, obtaining plans from the decision
center and performing actions with the specified tools.
"""
def __init__(self, decision_center=LLM, toolbox=[GOOGLE, Minecraft, RoboticArm]):
"""
Initialize an IntelligentAgent instance.

Args:
- decision_center: The decision-making center for the agent, defaults to an LLM.
- toolbox: A list of available tools for the agent, default includes GOOGLE, Minecraft, and RoboticArm.
"""
self.decision_center = decision_center
self.toolbox = toolbox
"""

def cs_ritual(self):
"""
Execute a rookie ritual for a computer science rookie using the decision center and specified tools.
"""
plans = self.decision_center("Hello, World!", toolbox=self.toolbox)
for tool, action in plans:
    action(tool)
"""

```

} 1. Object-Oriented Programming
Adv: Structured

} 2. Functional Programming
Adv: Modular & Explicit

} 3. Procedural Programming
Adv: Step-by-Step

Figure 7: We generate pseudo-code for the “IntelligentAgent” class and employ ChatGPT to compile its docstring. By contrasting the self-explanatory code with its natural language docstring, we observe that code exhibits greater structure, expressiveness, and logical coherence, underscoring certain advantages of code over natural language.

and class inheritance (Object Oriented Programming). This feature makes them super suitable for expressing nesting and complicated structures (cf. §3.3 and the perception part in §6.1). Another feature is that programming languages are often paired with a very powerful execution environment. This executable feature benefits much as it naturally delegates some harder tasks to lower level, like arithmetic computing or interacting with a simulated environment when connecting to a Database, Minecraft, and so on, also facilitating reasoning discussed in §3.2. What’s more, the execution often includes feedback mechanisms, which can be valuable for further refining the generator (§5 and §6.3).

A.2 Breadth by Code Delegation or Depth by Multimodality Joint Learning

LLMs can swiftly and cost-effectively address tasks involving more data modalities by utilizing code to invoke tools. Simultaneously, joint fine-tuning on multimodal data enhances the model’s precision and robustness in perceiving each modal-

ity, resulting in superior task performance. For instance, on the VQA dataset GQA (Hudson and Manning, 2019), ViperGPT (Surís et al., 2023), a typical code-centric paradigm, marginally surpasses the multimodal model BLIP-2 (Li et al., 2023b) in the zero-shot scenario after learning visual model API usages. However, its accuracy remains significantly lower than other supervised multimodal models. It is also still uncertain whether this approach will surpass the state-of-the-art models on multi-modal procedural planning (Liu et al., 2023b). One reason is that the code-centric paradigm’s effectiveness hinges on the central decision model and individual task execution components. This makes code-delegation approaches susceptible to error accumulation across steps and highly influenced by the worst-performing sub-modules or tools. Nevertheless, code delegation remains essential, as certain tools’ advantages, such as the precision of calculators and the flexibility of search engines, cannot be learned by training multimodality models alone. The high extensibility of the code-centric paradigm

to various tools and modalities also makes it a perfect fit for domains where training data is hard to collect at scale. We anticipate that the central decision model, utilizing code to invoke tools, will evolve from text-only LLMs to multimodality models capable of comprehensively understanding and processing multimodal data.

A.3 The Potential of Using Code-centric Framework for Intelligent Agent Construction

We observed a rising trend in leveraging code in the construction of LLM-based intelligent agents. As shown in §6, we showed three major scenarios where agents can effectively benefit from code usage. We also identified that this trend mainly originated from the increasing need to evaluate agents in a real-world scenario, where executive environments and interactions are everywhere. A natural question arises: Does code have the potential to substitute natural language and become the dominant media in the construction of agents?

A lot of work has begun to adopt this approach, like Voyager (Wang et al., 2023b) in a simulated Minecraft environment. They used code for high-level planning, low-level control sequence, and execution to interact with the environment. Acquired skills are also organized in the format of code snippets. With the code-centric paradigm, the framework is highly automatic and efficient. However, it’s also true that many framework today are still using natural language for planning, probably because they provide more human-interpretable reasoning steps. Human feedback in natural language is also widely used to harvest strong reward models that reflect real human preferences. We hypothesize that the integration of code will continue gaining popularity on our path to AGI, especially for facilitating interactions between agents and the real world. Nevertheless, natural language could hardly be replaced regarding the interaction between agents and humans.(Drori et al., 2022; Chen et al., 2023b; Lei and Deng, 2023). Leveraging this understanding, we aim to explore novel research avenues in LLM reasoning inspired by the utilization of “code”.

B Paper Statistics from Arxiv

We write a Python script that serves as a web scraper to extract paper details from the ArXiv preprint server, specifically focusing on the field

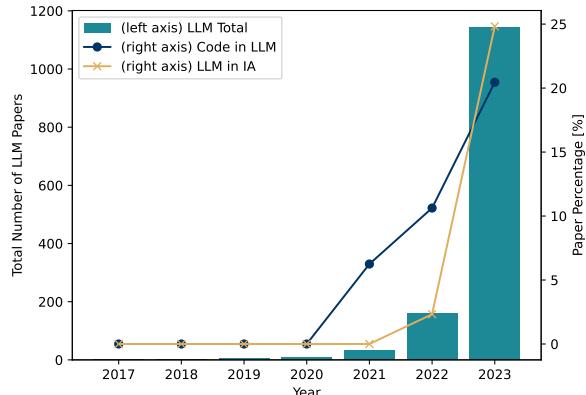


Figure 8: Paper statistics from Arxiv. We identified a significant and growing trend in recent research focused on code-based large language models (LLMs) and LLM-based intelligent agents (IAs). Code usage contributes much to the success of these cutting-edge models and systems.

of computer science. The web scraper gathers information about papers related to specific topics, including code, LLM, and IA. The script navigates through the ArXiv website, fetching essential details such as paper title, abstract, authors, and subject categories. We analyze and visualize data related to these papers in Figure 8, intending to provide insights into the trends and relationships between LLMs, code-related topics, and IAs in the past few years.

C Benchmarks for Evaluating Complex Reasoning with Code:

While there exist many benchmarks used to evaluate the abilities of LLMs (Liang et al., 2023c) across many disciplines, the benchmarks that most directly evaluate LLMs pre-trained on code in complex reasoning tasks are programming benchmarks such as CodeBLEU (Ren et al., 2020), where metrics that better match a human’s evaluation of what is good logical, interpretable, and syntactically concise code was created, and CodeXGLUE (Lu et al., 2021) where multiple programming tasks such as code repair and code defect detection were accumulated into one dataset. Other suitable benchmarks include math datasets such as many of MIT’s undergraduate math courses such as calculus and linear algebra, (Drori et al., 2022), LILA, a compilation of 23 tasks that test for mathematical abilities, language format, language diversity, and external knowledge abilities of LLMs (Mishra et al., 2023),

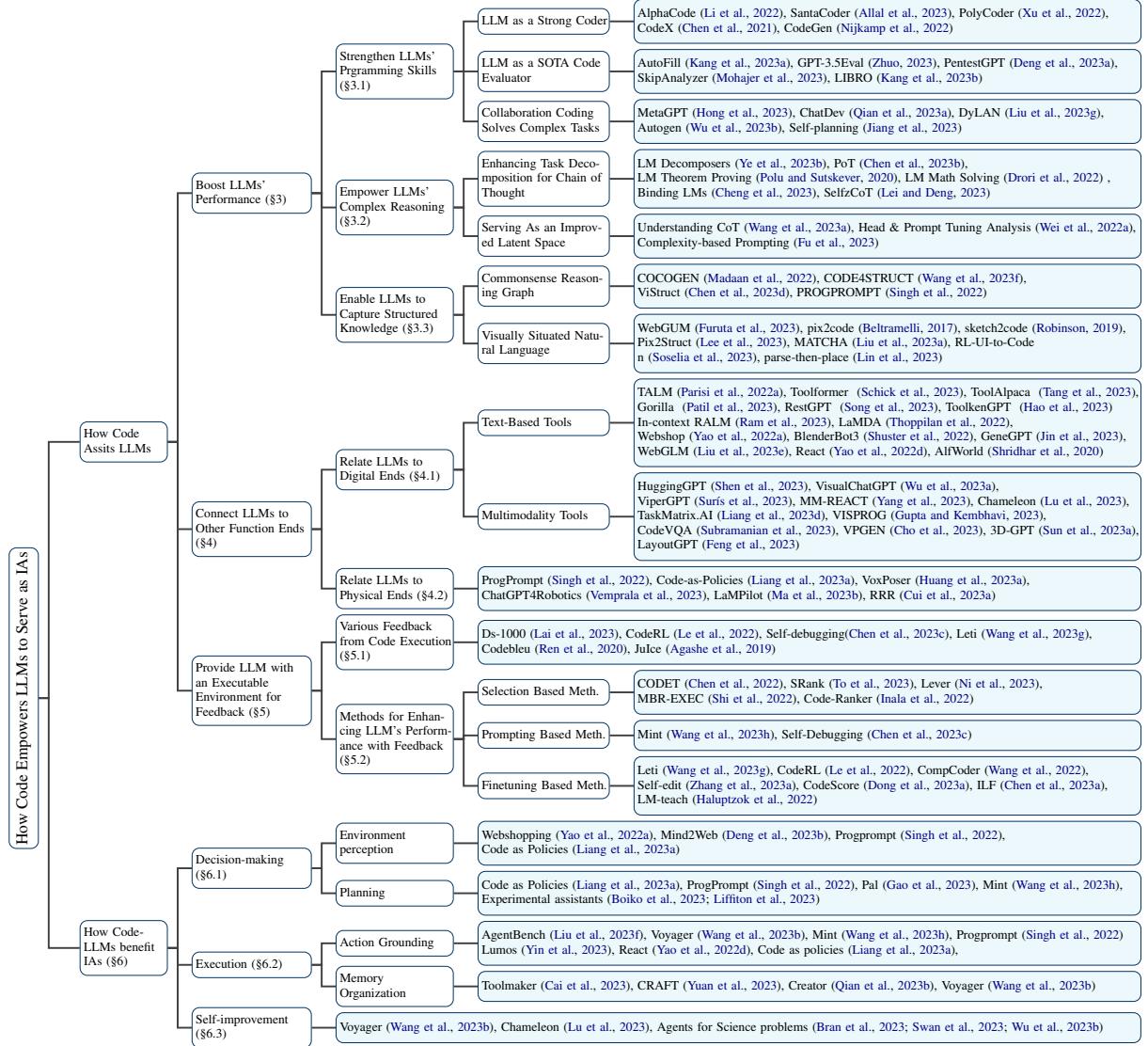


Figure 9: We hereby provide a complete list of the papers included in our survey.

and theorem proving from the metamath theorem code language (Polu and Sutskever, 2020). Others include question-answering tasks that require complex abilities to perform data retrieval in SQL databases (Ye et al., 2023b), such as those seen by the Spider dataset (Yu et al., 2019; Rajkumar et al., 2022).

D The Comprehensive Paper List

To complement the core paper list presented in Figure 2, we have included a comprehensive list of papers in Figure 9. It is important to note that this list excludes papers used for performance comparisons between code and natural language. Instead, it focuses on papers that have utilized code to augment the capabilities of Large Language Models and intelligent agents.

E Mappings of Sections to Core Code features

In each section, we identify key code features that contribute to the success of enhancing Large Language Models and Intelligent Agents. The correlation between each section and its core features is detailed in Table 2. We have classified code features into three main categories: Machine Executable, Structured and Expressive, and Explicit and Unambiguous. Various aspects of these core features play a pivotal role in the effective use of code. Detailed explanations of these aspects are provided in the right column of the table. Additionally, further information can be found in the preamble of each respective section.

Major Functionalities Facilitated	Key Features of Code
Strengthen LLMs' Programming Skills Correspond to §3.1 and Figure 3 (a)	Machine Executable: Pretraining with Code, the LLM is able to write code, evaluate code, and utilize collaborative coding to solve complex tasks.
Empower LLMs' Complex Reasoning Correspond to §3.2 and Figure 3 (b)	Structured and Expressive: The step-by-step nature of code benefits CoT. Machine Executable: LLMs can utilize code to help with certain capabilities like mathematical reasoning.
Enhance LLMs' Structured Knowledge Correspond to §3.3 and Figure 3 (c)	Structured and Expressive: Code can be used to express complex structures, some programming language features like logical expressions and class inheritance will be especially useful.
Connect LLMs to Other Functional Ends Correspond to §4 and Figure 4	Explicit and Unambiguous: Code is more explicit and clear than natural language, thus can better express clear instructions of connecting to any other functional ends.
Provide LLMs w/ Environmental Feedback Correspond to §5 and Figure 5	Machine Executable: Code execution result can be treated as feedback to finetune the LLMs further and make their performance more desirable
Help with IAs' Decision-Making Correspond to §6.1 and Figure 6	Structured and Expressive: Code pretraining Enhances agents' ability to precept structural knowledge, step by step feature improves CoT planning, logical expressions and nesting help with better control flow. Machine Executable: Help with solving mathematical tasks.
Help with IAs' Action Execution Correspond to §6.2 and Figure 6	Explicit and Unambiguous: Unambiguous function calling help with action grounding. Machine Executable: Agents can write reusable code snippets as its memory
Help with IAs' Self-improving Correspond to §6.3 and Figure 6	Machine Executable: Agent code execution results reflect potential environment change and can be utilized for self-improving

Table 2: We conclude the three major key features of code and correspond them to the major functionalities of LLMs and IAs they facilitated. The three key features are namely **Structured and Expressive**, **Machine Executable** and **Explicit and Unambiguous**. More details of how these features assist LLMs and IAs can be found in the preamble of each section.