

# NAMeGEN: Creative Name Generation via A Novel Agent-based Multiple Personalized Goal Enhancement Framework

Shanlin Zhou , Xinpeng Wang , Jianxun Lian , Zhenghao Liu ,  
Laks V.S. Lakshmanan , Xiaoyuan Yi\* , and Yongtao Hao\*

**Abstract**—Trained on diverse human-authored texts, Large Language Models (LLMs) unlocked the potential for Creative Natural Language Generation (CNLG), benefiting various applications like advertising and storytelling. Nevertheless, CNLG still remains difficult due to two main challenges. (1) *Multi-objective flexibility*: user requirements are often personalized, fine-grained, and pluralistic, which LLMs struggle to satisfy simultaneously; (2) *Interpretive complexity*: beyond generation, creativity also involves understanding and interpreting implicit meaning to enhance users' perception. These challenges significantly limit current methods, especially in short-form text generation, in generating creative and insightful content. To address this, we focus on *Chinese baby naming*, a representative short-form CNLG task requiring adherence to explicit user constraints (*e.g.*, length, semantics, anthroponomy) while offering meaningful aesthetic explanations. We propose NAMeGEN, a novel multi-agent optimization framework that iteratively alternates between objective extraction, name generation, and evaluation to meet diverse requirements and generate accurate explanations. To support this task, we further construct a classical Chinese poetry corpus with 17k+ poems to enhance aesthetics, and introduce CBNAMES, a new benchmark with tailored metrics. Extensive experiments demonstrate that NAMeGEN effectively generates creative names that meet diverse, personalized requirements while providing meaningful explanations, outperforming six baseline methods spanning various LLM backbones without any training.

**Index Terms**—Creative NLG, Multi-Objective Optimization, Collaborative Multi-Agents Framework.

## I. INTRODUCTION

Benefiting from large-scale model parameters and extensive training data, large language models (LLMs) [1]–[4] have not only achieved notable breakthroughs in language understanding [5], [6] and reasoning [7]–[9], but also unlocked the potential of Creative Natural Language Generation (CNLG) [10]–[12]. As a special case of controllable language generation [13], [14], CNLG emphasizes the novelty, distinctiveness, and aesthetics of generated text, which empowers a spectrum of downstream tasks [15]–[17], spanning from poetry generation [18]–[20], story writing [21], [22], summarization [23]–[25], and slogan generation [26].

Existing CNLG studies [10], [18]–[24], [26]–[30] mainly focus on achieving general objectives, *e.g.*, creativity, fluency,

Shanlin Zhou, Xinpeng Wang and Yongtao Hao are with the School of Computer Science and Technology, Tongji University, Shanghai, 201800, China (e-mail: {zhoushanlin,wangxinpeng,haoyt}@tongji.edu.cn). Xiaoyuan Yi and Jianxun Lian are with Microsoft Research Asia, Beijing, 100080, China (e-mail: {xiaoyuanyi,Jianxun.Lian}@microsoft.com). Zhenghao Liu is with Northeastern University, Shenyang, 110819, China (e-mail: liuzhenghao@mail.neu.edu.cn). Laks V.S. Lakshmanan is with The University of British Columbia, Vancouver, Canada (e-mail: laks@cs.ubc.ca).

Corresponding authors: Y. Hao and X. Yi.

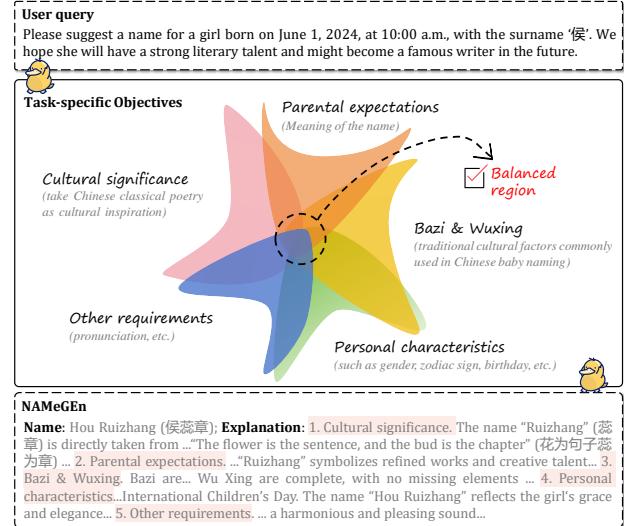


Fig. 1. Example of Chinese Baby Naming (NCB). Different colors indicate diverse objectives. The model seeks an optimal balance among them.

and coherence. However, in real-world CNLG applications, *e.g.*, naming design and slogan creation, existing methods face two primary challenges. *Challenge 1: Multi-objective flexibility*. User requirements for creativity are often personalized, fine-grained, and multi-objective [31]–[33], while LLMs struggle to satisfy these diverse goals simultaneously [34]. *Challenge 2: Interpretive complexity*. According to cognitive psychology theories [35], [36], creativity is not only reflected in the generation process, but also involves understanding and interpreting the implicit meanings [37] behind the content to enhance user perception. As illustrated in Fig. 1, in Naming Chinese Baby (NCB), a typical CNLG task, a generated name is required to meet *multiple* explicit user-specified requirements, *e.g.*, traditional Chinese cultural significance, parents' expectations for the baby, Bazi&Wuxing<sup>1</sup> and personal characteristics (see Fig. 2 for examples). Besides, corresponding high-quality explanations are needed to help users understand the underlying connection between the generated name and user requirements (like expectation), as well as implicit meanings and aesthetics. Such challenges place greater demands on the LLMs' capabilities for understanding, integration, and interpretability. As a result, existing LLMs are often regarded

<sup>1</sup>Bazi&Wuxing are traditional cultural factors commonly used in Chinese baby naming. Bazi analyzes the birth date and time, while Wuxing represents the five elements (wood, fire, earth, metal, and water). Together, they provide insight into fate and fortune based on traditional Chinese numerology.

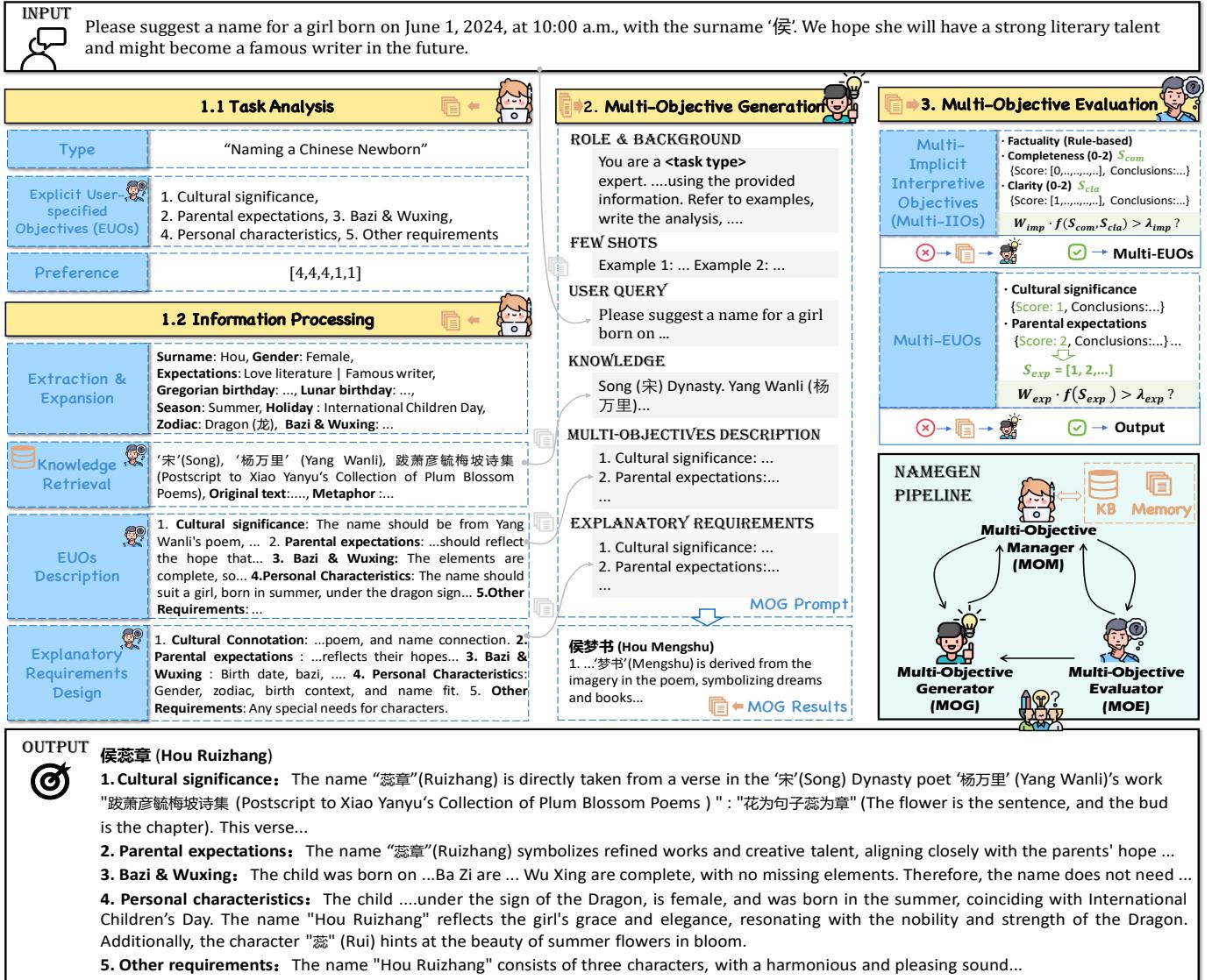


Fig. 2. Overview of NAMEGEN. Steps 1.1 and 1.2 constitute the multi-objective information preparation process, which is primarily handled by MOM and MOE. The dynamic iterative objective optimization process includes Steps 2 and 3: Step 2 is managed by MOG, while Step 3 reflects MOE's role in evaluating the generation results. The green block at the bottom right illustrates the complete pipeline of NAMEGEN.

as lacking human-level creativity [38]–[40].

These challenges are more pronounced in short-form CNLG [41]–[43] compared to the long-form one [44], [45], as LLMs need to integrate multiple, personalized, and fine-grained user requirements into a compact piece of text. Therefore, we focus on short-form CNLG, and propose two types of objectives grounded in the two challenges introduced above. (a) *Explicit User-specified Objectives (EUs)* refer to creative requirements explicitly set by users, such as parents' expectations in the NCB task. (b) *Implicit Interpretive Objectives (IOs)* represent fundamental attributes, e.g., accuracy, completeness, and clarity, of LLMs' explanation to connect the underlying meanings and aesthetics with generated text, which LLMs should satisfy by default. To sum up, short-form CNLG involves three goals: (i) EUOs optimization, (ii) IOs improvement, and (iii) the balance of both.

These goals are highly relevant to addressing the attribute

degradation issue in Multi-aspect Controllable Text Generation (MCTG) [46]–[48], where attribute imbalance or conflict arises during generation [49]–[51]. To handle this problem, early studies treated the decoding process as a Multi-Objective Optimization (MOO) process [52] and applied various optimization strategies to balance the different attributes to some extent. Typical methods include MUCOCO [31], GDC [46], Mix&Match [53], COLD decoding [54] and MacLaSa [49]. Besides, prefix-tuning [55] was also adopted to achieve multi-attribute control [47], [48], [50], [51]. Nevertheless, these methods are mainly designed to satisfy a pre-defined set of attributes, e.g., toxicity, topic and sentiment, and thus fail to handle the flexible and fine-grained EUOs and IOs (*challenges 1 & 2*). More recently, with the rapid development of LLMs, multi-aspect control of LLMs receives increasing attention, with a focus on multi-instruction constraint optimization via fine-tuning [56]–[58]. Despite good effectiveness,

these methods require substantial computational or data costs. Moreover, fine-tuning requires access to model parameters, making these methods incompatible with closed-source LLMs.

In this work, we aim to tackle the multi-objective flexibility and interpretive complexity challenges of CNLG, and take NCB as a typical short-term CNLG task. Given the strong capabilities of LLMs in generation [59] and real-world applications [60]–[62], we resort to In-Context Learning (ICL) [63] and multi-agent framework [60] to avoid unaffordable costs and achieve better flexibility and creativity. Since existing related work predominantly concentrates on reasoning tasks and thus is less suitable for multi-objective optimization issues [64]–[66], we propose a **Novel Agent-based Multi-Personalized-Goal Enhancement Framework (NAMeGEN)** to tackle the two challenges and achieve EUOs and IIOs. As shown in Fig. 2, NAMeGEN performs a two-stage process, *i.e.*, multi-objective information preparation and dynamic iterative objective optimization, which are executed through three agents, namely, *Multi-Objective Manager (MOM)*, *Multi-Objective Generator (MOG)*, and *Multi-Objective Evaluator (MOE)*. MOM is responsible for task analysis and information processing, while MOG generates a creative result by integrating user queries with MOM-processed information. MOE participates in evaluating MOM’s search to support MOG’s generation and assesses MOG’s outputs to ensure a balanced multi-objective creative outcome. The three agents collectively work on key information extraction, dynamic generation, and objective evaluation in an alternating and interactive manner, towards ultimately achieving optimal creative outcomes. The whole framework is backbone-agnostic, training-free, and efficient, allowing the integration with arbitrary either open-source or black-box LLMs. Since we select NCB as a typical CNLG task to demonstrate the effectiveness of our method, we further develop a benchmark named **CBNames**. Given that cultural meaning (*e.g.*, aesthetics and poeticness) is often considered in Chinese names and LLMs might produce hallucinations, we further collect an ancient Chinese poetry dataset (denoted as **CPoetry**), utilize it via Retrieval-Augmented Generation (RAG) [67], and design novel tailored metrics for this task considering both EUOs and IIOs. Experiments on various LLM backbones demonstrate the efficacy and superiority of NAMeGEN in short-term multi-objective CNLG, outperforming six recent strong baselines.

In summary, our contributions are as follows:

- We propose NAMeGEN, a novel framework for solving hybrid multi-objective optimization in short-term CNLG tasks. With training-free dynamic iterative optimization strategies, NAMeGEN achieves a balance among multi-EUOs, multi-IIOs, and the trade-offs between them, improving content quality and ensuring applicability across both open-source and closed-source models.
- We have developed a benchmark (CBNames) for the NCB task and created a new classical poetry retrieval database (CPoetry) to reduce hallucinations in LLMs.
- Our framework outperforms existing methods across all six backbone models, demonstrating superior creative generation and enhanced interpretability without requiring model-specific training.

## II. RELATED WORK

**Creative Natural Language Generation (CNLG)** is a long-standing research area, with classic tasks including poetry generation [19], [20], story writing [21], [22], summarization [23], [24], and slogan generation [26]. These studies primarily focus on objectives such as novelty, originality, fluency, and diversity of generated text. While these metrics are vital for assessing generation quality, they often overlook more personalized and nuanced goals required in real-world applications. In practice, existing approaches struggle with multi-objective flexibility (*challenge 1*). On the other hand, most prior work emphasizes generation outcomes while underestimating the importance of interpretation for user understanding. Creativity also involves understanding and explaining implicit meanings, as highlighted by cognitive psychology theories [35], [36]. However, current methods largely ignore interpretive complexity (*challenge 2*). Our work addresses these challenges by introducing task analysis and dynamic multi-objective evaluation strategies.

**Multi-Aspect Controllable Text Generation (MCTG)** aims to generate text that satisfies multiple attribute constraints (*e.g.*, sentiment, topic, style). While traditional methods often focus on single-aspect control, extending to multi-aspect scenarios introduces challenges like attribute interference and degeneration [46]–[51]. Some works frame MCTG as a multi-objective optimization (MOO) problem. MUCOCO [31] first formalizes MCTG as MOO by jointly optimizing through gradient descent for multiple attributes in the loss function. Similarly, [68] treats detoxification as a multi-objective problem involving non-toxicity, semantic preservation, and fluency. However, these works mainly target general attributes like toxicity, formality, or style, while overlooking more practical objectives, such as flexible, fine-grained user instructions and implicit interpretive goals. Although recent work [10], [56]–[58], [69] explores multi-constraint instruction following, few adopt explicit MOO frameworks, making it difficult to manage conflicting constraints effectively.

**LLM Based Multi-Agent.** Given the strong reasoning and planning capabilities of large language models (LLMs), LLM-based agents have seen rapid development. Compared to single-agent systems that solve problems through tool use and interactions with external environments, multi-agent systems emphasize decision-making through role-playing and communication among multiple agents [70]–[72]. However, the use of multi-agent systems has primarily focused on solving complex tasks [30], [73], [74]. To the best of our knowledge, no prior work has explored leveraging multi-agent collaboration to explicitly address multi-objective tasks, which is a gap our work seeks to fill.

## III. METHODOLOGY

### A. Overview

In this section, we formally introduce the task of multi-objective optimization in creative natural language generation (CNLG) and the proposed method. CNLG involves optimizing explicit user-specified objectives (EUOs), implicit interpretive objectives (IIOs), and the trade-offs between them, constituting

a hybrid multi-objective optimization problem. The goal of CNLG is to achieve optimization over this hybrid multi-objective setting. First, we formalize this task as follows:

$$Y = \arg \max_{\theta} f(x, O_{\text{exp}}, O_{\text{imp}}) \quad (1)$$

where  $Y = \{r, E\}$  denotes the final output of it, with the generated text  $r$  and a set of explanations  $E$ . For example, in the Naming Chinese Baby (NCB) task,  $r$  is the generated name, and  $E$  is the corresponding set of explanations detailing how the name reflects each objective.  $\theta$  is a dynamic optimization parameter,  $x$  is input,  $O_{\text{exp}} = \{o_1^e, o_2^e, \dots, o_m^e\}$  and  $O_{\text{imp}} = \{o_1^i, o_2^i, \dots, o_n^i\}$  means  $m$  EUOs, IIOs respectively,  $E$  is also composed of  $m$  individual explanations  $\{e_1, e_2, \dots, e_m\}$ , corresponding to EUOs.

In our work, we design a training-free and multi-agent collaborative architecture, NAMeGEN, to reach this aim. The overview of NAMeGEN is shown in Fig. 2. Our pipeline is collaboratively executed by three intelligent agents, including the Multi-Objective Manager (MOM)  $A_M$ , the Multi-Objective Generator (MOG)  $A_G$ , and the Multi-Objective Evaluator (MOE)  $A_E$  (as shown in the green block of Fig. 2). Next, we divide the whole pipeline into two stages. **The first stage** is multi-objective information preparation (see steps 1.1 and 1.2 in Fig. 2). It covers task analysis, key information extraction and expansion, EUO description, and interpretability requirement design, aiming to build a solid foundation for hybrid multi-objective generation. This stage is completed by collaboration between MOM  $A_M$  and MOE  $A_E$ . **The second stage** is dynamic iterative objective optimization (see steps 2 and 3 in Fig. 2). In this stage, MOG  $A_G$  drafts a creative output using information from MOM  $A_M$ , then dynamically and iteratively refines it based on MOE  $A_E$ 's evaluation and feedback against EUOs and IIOs. This way of optimization helps to progressively refine the CNLG by balancing hybrid multi-objectives. Next, we elaborate on each stage in detail.

### B. Multi-objective Information Preparation

This stage involves analyzing the creative task and processing multi-objective information, which aims to gain a deeper understanding of user personalized requirements and enhance the ability to interpret implicit meaning, helping to produce higher-quality creative results later.

1) *Task analysis process (see 1.1 in Fig. 2):* This includes identifying the task type  $I_{\text{type}} \leftarrow A_M(x)$  (e.g., naming a Chinese baby), get  $O_{\text{exp}} = \{o_1^e, o_2^e, \dots, o_m^e\}$ , and estimating user preference scores for  $O_{\text{exp}}$ . Specially, when EUOs  $O_{\text{exp}}$  are not provided by user, MOM  $A_M$  can automatically parse appropriate EUOs, by  $O_{\text{exp}} \leftarrow A_M(I_{\text{type}})$ . However, when auto-parsing objectives, the model often parses similar ones (e.g., "Harmonious Pronunciation", "Easy to Pronounce and Remember", etc.). Therefore, we introduce evaluation, ensuring the independence and appropriateness of each objective. Then the final EUOs we get by  $O_{\text{exp}} \leftarrow A_M(I_{\text{type}}, b)$ , where the evaluation feedback  $b \leftarrow A_E(O_{\text{exp}}, I_{\text{type}})$ . We estimate user preference by  $I_{\text{pref}} \leftarrow A_M(O_{\text{exp}}, x)$  and normalize it into a weight vector  $W_{\text{exp}}$ . This helps to better meet different users' personalized needs. Later, this weight  $W_{\text{exp}}$  will work in the dynamic iterative objective optimization stage.

---

### Algorithm 1: MOO-Retrieval

---

```

Input: User input  $x$ , key information  $I_{\text{ki}}$ , knowledge base  $\mathcal{D}$ 
Output: The best knowledge  $I_{\text{rk}}$ 
1 Initialize max coarse filtering round  $n_f$ , max retrieval round  $n_{r_{\text{max}}}$ , retrieved
   knowledge number  $m$ , round counter  $i \leftarrow 0$ , current retrieved knowledge
    $\mathcal{R} \leftarrow \emptyset$ , and history knowledge  $\mathcal{H} \leftarrow \emptyset$ .
2 repeat
3    $\mathcal{D} \leftarrow \mathcal{D} \setminus \mathcal{H}$ .
4   if  $i < n_f$  then Execute coarse filtering  $\mathcal{R}' \leftarrow f_{\text{filter}}(I_{\text{ki}})$ .
5   if  $i = 1$  then get retrieval query  $q_i \leftarrow A_M(x)$ .
6   Semantic matching to get candidates  $\mathcal{R} \leftarrow f_{\text{match}}(q_i, \mathcal{R}')$ .
7    $\mathcal{H} \leftarrow \mathcal{H} \cup \mathcal{R}$ ,  $i \leftarrow i + 1$ .
8   Evaluate:  $I_{\text{rk}} \leftarrow A_E(\mathcal{R}, x)$  if found, else  $q_i \leftarrow A_E(\mathcal{R}, x, q_i)$ .
9 until  $I_{\text{rk}}$  is found or  $i > n_{r_{\text{max}}}$ .
10 if  $i > n_{r_{\text{max}}}$  then  $I_{\text{rk}} \leftarrow A_E(\mathcal{H}, x)$ .

```

---

### 2) Process multi-objective information (1.2 in Fig. 2):

This covers extracting and expanding information, retrieving related knowledge, and refining the EUOs' description and explanatory requirements. Firstly,  $A_M$  extracts key content  $I_{\text{ki}}$  (e.g., surname and date of birth in the naming task) from user input  $x$ , for use in retrieval and expansion. If insufficient, it can infer additional details (e.g., deducing the birth season based on the date of birth) to enrich  $I_{\text{ki}}$ . Next, we add a retrieval step (shown as in Algorithm 1) into  $A_M$  to mitigate the hallucination of LLMs in domain knowledge. User input  $x$ , with its multi-aspect content, also makes this retrieval process a multi-objective task. Specifically, MOM  $A_M$  first uses information  $I_{\text{ki}}$  to retrieve preliminary candidates  $R'$ . Then, based on  $x$ , it reconstructs a style-aligned query  $q$  with the knowledge base (KB), to retrieve the top- $k$  candidates  $\mathcal{R} = \{r_1, r_2, \dots, r_m\}$  from  $R'$ , using an embedding-based retrieval strategy<sup>2</sup>. Style rewriting for query  $q$  can help to increase the efficiency of retrieval. To find the best knowledge  $I_{\text{rk}}$ , we make MOE  $A_E$  participate in evaluation and rerank (see line 7 in Algorithm 1). During evaluation, if none is satisfactory, MOE  $A_E$  further revises the query  $q$  but keeps its style. The revised query  $q$  is fed back to MOM  $A_M$ , which discards previous results and performs a new retrieval from the remaining knowledge. This process iterates until the optimal  $I_{\text{rk}}$  is found. To help MOG better understand EUOs and try to align with IIOs, MOM  $A_M$  designs separately detailed descriptions  $I_{\text{desc}} = \{c_1, c_2, \dots, c_m\}$  and explanatory requirements  $I_{\text{reqs}} = \{q_1, q_2, \dots, q_m\}$ . Descriptions  $I_{\text{desc}}$  are more reliable information than EUOs' keywords  $O_{\text{exp}}$ , and requirements  $I_{\text{reqs}}$  help MOG  $A_G$  generate more accurate, complete, and clear content. Notably, their incorporation can effectively enhance the efficiency of the hybrid multi-objective optimization process, thereby contributing to the generation of creative outputs with improved quality and interpretability. MOE  $A_E$  is also involved in this process, ensuring the rationality of these contents by evaluation, feedback, and iteration. These steps are as follows:

$$(I_{\text{desc}}, I_{\text{reqs}}) \leftarrow A_M(x, O_{\text{exp}}, I_{\text{ki}}, I_{\text{rk}}, b) \quad (2)$$

$$b_{\text{detail}} \leftarrow A_E(I_{\text{desc}}, I_{\text{reqs}}) \quad (3)$$

where  $b_{\text{detail}}$  is the evaluation feedback. Moreover, they will serve as crucial evaluation criteria in the dynamic iterative objective optimization stage. Ultimately, we obtain the expanded

<sup>2</sup><https://huggingface.co/lier007/xiaobu-embedding-v2>

hybrid information set:  $\mathcal{I} = \{I_{\text{type}}, I_{\text{pref}}, I_{\text{rk}}, I_{\text{desc}}, I_{\text{reqs}}\}$ . All the above information is stored in MOM's memory to facilitate interactions with other intelligent agents.

### C. Dynamic Iterative Objective Optimization

Algorithm 2 shows the process of this stage. MOG  $A_M$  utilizes the user query  $x$  and the hybrid information  $\mathcal{I}$  from MOM to generate the initial results  $\mathcal{Y} = \{r, E\}$ , where  $r$  is result and  $E$  is its explanation set  $\{e_1^j, e_2^j, \dots, e_m^j\}$ . Afterward, MOE  $A_E$  evaluates the completion of  $\mathcal{Y}$  in EUOs and IIOs separately, based on  $I_{\text{desc}}$  and  $I_{\text{reqs}}$ . Our work defines three IIOs—accuracy, completeness, and clarity—and evaluates this aspect first, as MOE  $A_E$ 's understanding of EUOs' completion depends on them. Evaluating accuracy is essential for mitigating model hallucinations. In our work, MOE  $A_E$  adopts a rule-based approach  $F_{\text{acc}}$  and refers to information  $I$  for evaluation. Given the subjectivity of completeness and clarity, rule-based methods are no longer applicable. Therefore, MOE  $A_E$  designs a new scoring mechanism to measure their performance respectively. To integrate them, we provide a trans-function:

$$\theta_{\text{imp}} = W_{\text{imp}} \cdot S_{\text{imp}} \quad (4)$$

$$W_{\text{imp}} = [\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n}] \in \mathbb{R}^n \quad (5)$$

where  $\theta_{\text{imp}}$  represents the implicit performance,  $W_{\text{imp}}$  is the weight vector,  $S_{\text{imp}}$  is the score vector of IIOs. In this way, it mathematically converts the optimization for IIOs into a single-objective optimization problem.

We similarly evaluate each EUO and combine them to form the score vector  $S_{\text{exp}}$ . Different from traditional optimization logic, due to the user input  $x$  usually implying their potential preference for EUOs, it is necessary to decompose  $x$  and obtain the preference  $I_{\text{pref}}$  (from MOE  $A_E$ ) as the optimization weight vector  $W_{\text{exp}}$ . Here, EUOs' unified quantitative result is:

$$\theta_{\text{exp}} = W_{\text{exp}} \cdot S_{\text{exp}} \quad (6)$$

Moreover, expected criteria  $\psi_{\text{imp}}$  and  $\psi_{\text{exp}}$  are introduced to judge implicit and explicit quantification results, respectively. It is worth noting that except for the generated content, the criteria  $\psi_{\text{imp}}$  and  $\psi_{\text{exp}}$  adjust dynamically with the iteration:

$$\psi_{\text{imp}} = \begin{cases} \delta & j_{\text{imp}} < t_w \\ \frac{\delta}{\alpha \cdot \log(j_{\text{imp}} + t_w)} & j_{\text{imp}} > t_w \end{cases} \quad (7)$$

$$\psi_{\text{exp}} = \begin{cases} \delta & j_{\text{exp}} < t_w \\ \frac{\delta}{\alpha \cdot \log(j_{\text{exp}} + t_w)} & j_{\text{exp}} > t_w \end{cases} \quad (8)$$

where  $j_{\text{imp}}$  is implicit evaluation round counter,  $j_{\text{exp}}$  is explicit evaluation round counter,  $t_w$  is warmup round,  $\delta$  is initial threshold, and  $\alpha$  is decay factor. Thanks to this dynamic iteration process, the generation quality is improved while multiple explicit objectives are achieved in balance.

---

### Algorithm 2: DI-HMOO

---

**Input:** User input  $x$ , hybrid information  $\mathcal{I}$   
**Output:** Creative output  $\mathcal{Y}$

**1 Dynamic Iterative Multi-Objective Optimization Stage**

2 Initialize examples  $I_{\text{shots}}$ , generation round  $j \leftarrow 0$ , implicit evaluation round  $j_{\text{imp}} \leftarrow 0$ , explicit evaluation round  $j_{\text{exp}} \leftarrow 0$ , max evaluation round  $t_{\text{max}}$ , warmup round  $t_w$ , initial threshold  $\delta$ , decay factor  $\alpha$ , implicit rationality objective  $\mathcal{O}_{\text{imp}} = \{o_{\text{acc}}, o_{\text{com}}, o_{\text{cla}}\}$ , implicit multi-objective weights  $W_{\text{imp}} \leftarrow [0.5, 0.5]$ , explicit multi-objective weights  $W_{\text{exp}}$ , implicit history  $\mathcal{H}_{\text{imp}} \leftarrow \emptyset$ , explicit history  $\mathcal{H}_{\text{exp}} \leftarrow \emptyset$ , and the  $j$ -th step regenerate flag  $l^j \leftarrow \text{NULL}$  (0 is no regeneration and 1 is regeneration).

3 **repeat**

4      $A_G$  generate result  $\mathcal{Y}^j \leftarrow A_G(x, I_{\text{rk}}, I_{\text{desc}}, I_{\text{reqs}}, I_{\text{shots}}, l^j, \mathcal{B}^j)$

5      $A_E$  evaluates result  $\mathcal{Y}^j$ :

6         Setting thresholds and scores:  $\psi_{\text{imp}}^j, \psi_{\text{exp}}^j$ .

7         Reset the  $j$ -th scores and feedback:  $\theta_{\text{imp}}^j \leftarrow \text{NULL}$ ,  $\theta_{\text{exp}}^j \leftarrow \text{NULL}$ ,  $\mathcal{B}^j \leftarrow \emptyset$ .

8         **a. Implicit evaluation:**

9              $(l^j, b_j) \leftarrow F_{\text{acc}}(\mathcal{Y}^j); \mathcal{B}^j \leftarrow \mathcal{B}^j \cup \{b_j\}$ .

10             **if**  $l^j$  is 1 **then**  $\mathcal{H}_{\text{imp}} \cup \{(\mathcal{Y}^j, \theta_{\text{imp}}^j)\}; j \leftarrow j + 1$ ; **continue**.

11             **foreach**  $e_k^j$  in  $\{e_1^j, e_2^j, \dots, e_m^j\}$  **do**

12                  $(s_{\text{com}}^j, s_{\text{cla}}^j, b_j) \leftarrow A_E(e_k^j, u_k, I_{\text{type}});$

13                  $\mathcal{S}_{\text{com}}^j \leftarrow \mathcal{S}_{\text{com}}^j \cup \{s_{\text{com}}^j\};$

14                  $\mathcal{S}_{\text{cla}}^j \leftarrow \mathcal{S}_{\text{cla}}^j \cup \{s_{\text{cla}}^j\};$

15                  $\mathcal{B}^j \leftarrow \mathcal{B}^j \cup \{b_j\}$ .

16              $S_{\text{imp}}^j = \{\text{Norm}(\text{AVG}(\mathcal{S}_{\text{com}}^j)), \text{Norm}(\text{AVG}(\mathcal{S}_{\text{cla}}^j))\};$

17              $\theta_{\text{imp}}^j \leftarrow W_{\text{imp}} \cdot S_{\text{imp}}^j$ .

18             **if**  $\theta_{\text{imp}}^j < \psi_{\text{imp}}^j$  **then**  $l^j \leftarrow 1$ ;  $\mathcal{H}_{\text{imp}} \cup \{(\mathcal{Y}^j, \theta_{\text{imp}}^j)\}; j \leftarrow j + 1$ ; **continue**.

19         **b. Explicit evaluation:**

20             **foreach**  $e_k^j$  in  $\{e_1^j, e_2^j, \dots, e_m^j\}$  **do**

21                  $(s_{\text{exp}}^j, b_j) \leftarrow A_E(e_k^j, q_k, I_{\text{rk}});$

22                  $\mathcal{S}_{\text{exp}}^j \leftarrow \mathcal{S}_{\text{exp}}^j \cup \{s_{\text{exp}}^j\};$

23                  $\mathcal{B}^j \leftarrow \mathcal{B}^j \cup \{b_j\}$ .

24              $\theta_{\text{exp}}^j \leftarrow W_{\text{exp}} \cdot \text{Norm}(\mathcal{S}_{\text{exp}}^j)$ .

25             **if** then

26                 **if**  $\theta_{\text{exp}}^j < \psi_{\text{exp}}^j$  **then**  $l^j \leftarrow 1$ ;  $\mathcal{H}_{\text{exp}} \cup \{(\mathcal{Y}^j, \theta_{\text{exp}}^j)\}; j \leftarrow j + 1$ ; **continue**.

27      $l^j \leftarrow 0$ .

28 **until**  $l^j$  is 0 or  $j > t_{\text{max}}$ .

29 **if**  $j < t_{\text{max}}$  **then**  $\mathcal{Y} \leftarrow \mathcal{Y}^j$ .

30 **else if**  $\mathcal{H}_{\text{exp}}$  is not  $\emptyset$  **then** Find best result  $\mathcal{Y}$  from  $\mathcal{H}_{\text{exp}}$ .

31 **else** Find best result  $\mathcal{Y}$  from  $\mathcal{H}_{\text{imp}}$ .

---

## IV. EXPERIMENTS

### A. Experimental Setup

**1) Datasets:** To validate the effectiveness of our framework, we conduct Naming Chinese Babies (NCB) as a typical short-term Creative Natural Language Generation (CNLG) task. NCB is required to consider various factors in creating a short name, making it a significant challenge. Firstly, we use a large language model (LLM), through manual verification, to simulate 500 human naming requests, create a benchmark (called **CNames**). This data encompasses over 200 Chinese surnames, including compound surnames, and reflects a range of potential user considerations in the naming process, such as traditional Chinese cultural significance, parental expectations, the Bazi&Wuxing, personal characteristics, and other special requirements. User preferences are annotated per objective for each query. In China, parents often draw on classical poetry to give their names cultural depth. However, LLMs frequently make factual errors in the poetry domain. To address this, we collect and expand a dataset of classical Chinese poetry (called **CPoetry**) containing 176,450 entries with metadata such as poets, dynasties, titles, content, interpretations, and themes.

2) *Baseline Methods*: We compare our proposed framework, NAMeGEN, with the following methods. **Base** is the raw query with task description. **CoT** [75] involves appending the phrase “Let’s think step by step” to the input. **TDB** [76] is the zero-shot prompting with “Take a deep breath and work on this problem step-by-step”. **Few-shot** [77] adds a few demonstration exemplars of question-answer in prompt. **Q2Kw** [78] integrates retrieved knowledge into the prompt by first using an LLM to generate results from the user’s raw requirement, then combining these results with the original requirement to form an extended query for knowledge retrieval. **LLM-D** [30] uses multiple agents to conduct a multi-stage discussion to complete a creative and diverse response task.

3) *Evaluation and Metrics*: Recent research has shown the capability of LLM to emulate human judgment and effectively evaluate content [79]–[81]. In this work, to evaluate quantitatively, we use both LLM and human evaluations.

**In LLM evaluation**, we propose metrics to evaluate the performance of Explicit User-specified Objectives (EUsOs), Implicit Interpretive Objectives (IIOs), and their overall optimization. In other words, these are designed to evaluate how well different methods address Hybrid Multi-Objective Optimization (HMOO) and creative diversity in short-term CNLG. Firstly, for EUsOs, we design Explicit multi-objective Completeness (**EC**) and Standard Deviation based EC (**EC\_std**):

$$EC = \frac{1}{N} \sum_{i=1}^N \frac{\sum_{j=1}^m w_{i,j} s_{i,j}}{\sum_{j=1}^m w_{i,j}} \quad (9)$$

$$EC_{std} = \frac{1}{N} \sum_{i=1}^N \sqrt{\frac{1}{m_i} \sum_{j=1}^{m_i} \left( s_{i,j} - \frac{1}{m_i} \sum_{j=1}^{m_i} s_{i,j} \right)^2} \quad (10)$$

where  $s_{i,j}$  and  $w_{i,j}$  denote the score and user weight for the  $j$ -th objective of  $i$ -th sample, with  $m$  EUsOs and  $N$  samples.  $s_{i,j}$  is scored by LLM. Higher EC implies better performance in achieving multi-EUsOs.  $EC_{std}$  quantifies the balance of multi-EUsOs’ performance via standard deviation, where a lower score means better. Secondly, we evaluate different IIOs’ performance, using **ACC**, **CRC**, **LR**, which separately measure hallucination degree, semantic quality (comprehensiveness, relevance, clarity), and logical consistency. For all of them, a higher score indicates better performance. Subsequently, we compute Implicit Multi-Objective Completeness (**IC**) and Standard Deviation based IC (**IC\_std**). IC is the average of ACC, CRC, and LR.  $IC_{std}$  is their standard deviation. Higher IC indicates better multi-IIO performance, while lower  $IC_{std}$  reflects more balanced and stable results. The following describes the method for calculating  $IC_{std}$ :

$$IC_{std} = \frac{1}{N} \sum_{i=1}^N \sqrt{\frac{1}{3} \sum_{k \in \{ACC, CRC, LR\}} (s_{i,k} - \bar{s}_i)^2} \quad (11)$$

where  $\bar{s}_i$  is the mean of the  $i$ -th sample’s scores on ACC, CRC, and LR, with  $s_{i,k}$  denoting the score for the  $k$ -th objective. Finally, from a holistic view, we provide Comprehensive Multi-Objective Completeness (**CC**), Standard Deviation-based CC (**CC\_std**), and Diversity (**DIV**). CC averages EC and IC;

**CC\_std** averages their standard deviations. DIV scores uniqueness across methods for the result of each sample, where if it appears in other methods, it scores 0; otherwise, 1. Higher DIV indicates better diversity. CC\_std is calculated as follows:

$$CC_{std} = \frac{1}{N} \sum_{i=1}^N \sigma_i \quad (12)$$

$$\sigma_i = \sqrt{\frac{1}{2} ((EC_i - \bar{x}_i)^2 + (IC_i - \bar{x}_i)^2)} \quad (13)$$

where  $\sigma_i$  means the standard deviation,  $\bar{x}_i$  denote the average of EC and IC for sample  $i$ . In short, EC and  $EC_{std}$  reflect the performance and balance of multi-EUO optimization. ACC, CRC, and LR measure explanation quality, while IC and  $IC_{std}$  evaluate the performance and balance of multi-IIO optimization. CC and  $CC_{std}$  can capture the overall completeness and balance of both EUOs and IIOs. DIV compares the diversity of results across methods.

**In human evaluation**, we randomly sample 50 test cases, each evaluated by 3 annotators on a 0–3 scale. EUOs are scored based on completeness, while IIOs are evaluated in CRC and LR. In addition, annotators are also asked to identify the best response among the different methods for each case.

4) *Implementation Details*: To ensure consistency, the NCB task prompt includes five objectives: traditional Chinese cultural significance, parental expectations, Bazi&Wuxing, personal characteristics, and other special requirements. We use the following seven LLMs as backbones: Qwen (long) [82], GLM-4 [83], DeepSeek (V3) [84], Mistral (small-latest) [85], emini (1.5-flash) [83], and GPT4o [59]. Among them, Qwen [82], GLM-4 [83] and DeepSeek [84] are LLMs proficient in Chinese, while Mistral [85], Gemini [83], and GPT-4o [59] are LLMs proficient in English. For DeepSeek [84], we adopt the recommended setting (temperature=1.5) to enhance creative generation. We use Kimi [86] as the LLM evaluator. ACC is computed via LLM extraction and rule-based methods. Other metrics, rated from 0 (invalid) to 3 (excellent), are averaged and normalized to a 0–100 scale.

## B. Method Comparison

To demonstrate NAMeGEN’s effectiveness, we compare it with five baseline categories: Base, reasoning-based

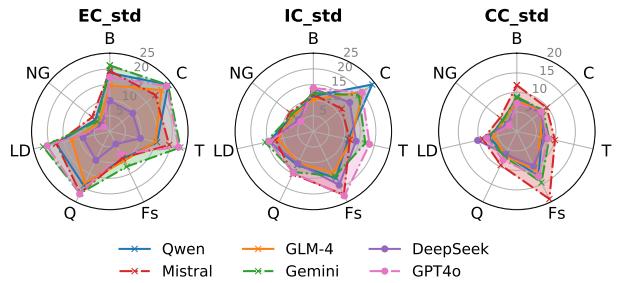


Fig. 3. Comparison of balance performance on EUOs, IIOs, and their overall combination across different methods. Colored lines correspond to different backbone models. We use method abbreviations for clarity: B = Base, C = CoT, T = TDB, Fs = Few-shot, Q = Q2Kw, LD = LLM-D, NG = NAMeGEN.

TABLE I  
LLM EVALUATION RESULTS. BOLD ONES INDICATE THE  
BEST-PERFORMING METHOD FOR EACH METRIC.

Backbone	Method	EC ↑	IC ↑	CC ↑	DIV ↑
Qwen	Base	85.03	76.29	80.66	44.20
	CoT	76.98	70.75	73.86	91.00
	TDB	87.34	82.85	85.10	44.80
	Few-shot	94.07	76.43	85.25	60.40
	Q2Kw	83.31	80.05	81.68	84.60
	LLM-D	85.51	80.75	83.13	82.80
	<b>NAMeGEN</b>	<b>96.72</b>	<b>92.70</b>	<b>94.71</b>	<b>99.80</b>
GLM4	Base	88.37	79.44	83.90	56.60
	CoT	80.25	73.39	76.82	93.80
	TDB	88.12	83.25	85.68	56.20
	Few-shot	94.10	79.49	86.79	75.80
	Q2Kw	85.95	80.40	83.18	91.40
	LLM-D	91.21	86.21	88.71	96.20
DeepSeek	<b>NAMeGEN</b>	<b>97.83</b>	<b>92.94</b>	<b>95.38</b>	<b>98.80</b>
	Base	93.53	85.29	89.41	69.20
	CoT	93.46	84.74	89.10	73.60
	TDB	93.40	84.91	89.15	72.20
	Few-shot	98.02	84.25	91.14	76.80
	Q2Kw	93.17	86.93	90.05	96.00
	LLM-D	93.90	81.21	87.56	96.60
Mistral	<b>NAMeGEN</b>	<b>98.93</b>	<b>95.22</b>	<b>97.08</b>	<b>99.60</b>
	Base	82.10	68.11	75.11	88.40
	CoT	82.46	72.40	77.43	86.60
	TDB	81.09	71.55	76.32	86.80
	Few-shot	93.97	67.17	80.57	94.00
	Q2Kw	79.82	70.81	75.31	91.60
Gemini	LLM-D	84.19	75.88	80.03	97.00
	<b>NAMeGEN</b>	<b>94.94</b>	<b>91.71</b>	<b>93.32</b>	<b>98.80</b>
	Base	84.56	74.10	79.33	69.60
	CoT	80.17	72.28	76.22	87.40
	TDB	81.50	73.11	77.31	81.20
GPT4o	Few-shot	93.66	73.77	83.71	76.40
	Q2Kw	82.81	77.41	80.11	76.60
	LLM-D	82.15	76.84	79.49	87.60
	<b>NAMeGEN</b>	<b>97.51</b>	<b>92.72</b>	<b>95.12</b>	<b>99.20</b>
	Base	86.08	79.29	82.69	94.80
GPT4o	CoT	78.10	71.90	75.00	96.80
	TDB	80.07	72.81	76.44	96.20
	Few-shot	95.34	77.68	86.51	96.00
	Q2Kw	83.00	75.07	79.04	95.80
	LLM-D	82.93	77.01	79.97	96.00
	<b>NAMeGEN</b>	<b>99.15</b>	<b>96.22</b>	<b>97.69</b>	<b>99.00</b>

(CoT [75], TDB [76]), in-context (Few-shot [77]), retrieval-based (Q2Kw [78]), and traditional agent-based (LLM-D [30]). Evaluations use LLM-based and human metrics.

1) *LLM Evaluation Results:* Table I shows that NAMeGEN consistently outperforms all baselines across backbones and metrics, confirming its effectiveness in NCB. Fig. 3 further shows it achieves the lowest EC\_std, IC\_std, and CC\_std, highlighting its superior balance in EUOs, IIOs, and their combination. We next analyze the HMOO challenge from three perspectives: overall, EUOs, and IIOs.

Overall, Table I shows that backbone differences impact performance partially. Firstly, Chinese-oriented models gener-

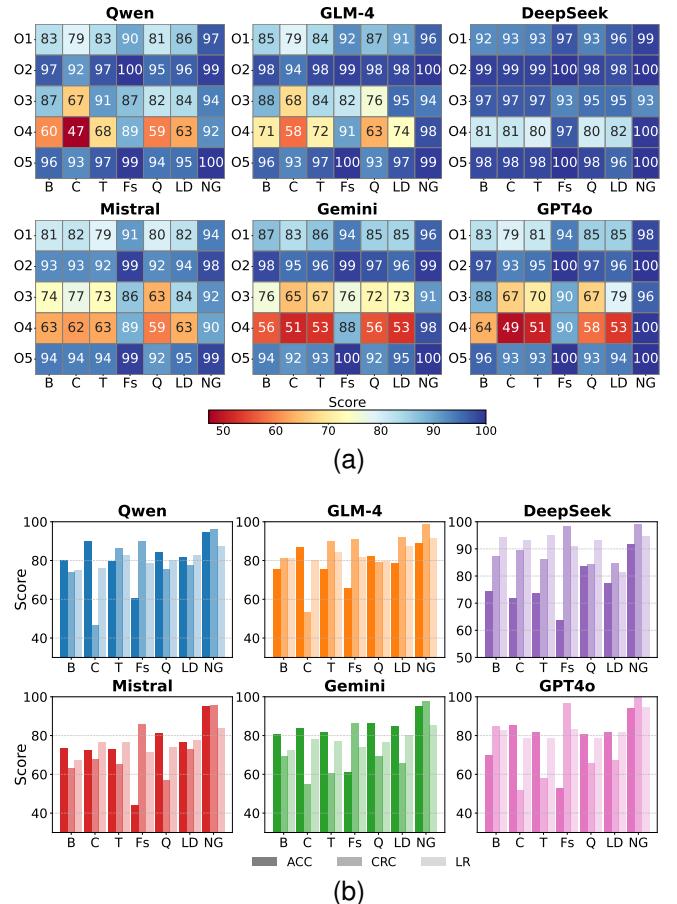


Fig. 4. Fine-grained comparison of explicit and implicit objective completeness. (a) shows explicit scores (blue = high, red = low) across five aspects. O1–O5 represent cultural significance, parental expectations, Bazi & Wuxing, personal traits, and other requirements. (b) shows implicit results, using different colors to indicate various backbones, and shade depth reflects different metrics. Method abbreviations are the same as in Fig. 3.

ally outperform English ones. For instance, under the Base, CC reaches 80.66 with Qwen but drops to 75.11 with Mistral. Similarly, models with stronger contextual understanding (e.g., DeepSeek, GPT-4o) perform better. In Few-shot, for example, CC with GPT-4o reaches 86.51, which is 5.94 points higher than with Mistral (80.57). In contrast to other methods' fluctuation across backbones, NAMeGEN shows consistently strong performance, demonstrating superior robustness and cross-backbone generalizability. Otherwise, Fig. 3 shows that balance mainly depends on methods, not backbones. Although all methods boost diversity (DIV rise compared to Base in Table I), NAMeGEN achieves the highest creative diversity, likely owing to more effective iterative guidance.

For multi-EUOs' optimization in the NCB task, in-context methods demonstrate better improvements in both creativity and balance compared to reasoning-based ones. As shown in Table I, reasoning-based methods provide limited achievement in explicit objectives, with notable declines in some cases. For example, with GLM-4, CoT's EC is 80.25, which is lower than the Base's 88.37 (a drop of 8.12), highlighting CoT's limitations in CNLG, as also discussed in [87]. In contrast, few-shot improves explicit performance by introducing additional examples (e.g., Few-shot's EC rises from Base's 88.37 to 94.10

TABLE II  
ABLATION STUDY RESULTS WITH QWEN AS THE BACKBONE.

Method	Explicit		Implicit				Comprehensive		
	EC $\uparrow$	EC_std $\downarrow$	ACC $\uparrow$	CRC $\uparrow$	LR $\uparrow$	IC $\uparrow$	IC_std $\downarrow$	CC $\uparrow$	CC_std $\downarrow$
NAMeGEN	<b>96.72</b>	<b>5.95</b>	<b>94.55</b>	<b>96.36</b>	87.19	<b>92.70</b>	<b>8.28</b>	<b>94.71</b>	<b>4.64</b>
wo/ ImpExp	94.48	9.81	63.25	92.47	<b>88.68</b>	81.47	17.19	87.97	9.69
wo/ Imp	95.38	7.10	91.05	92.13	85.43	89.54	9.17	92.46	6.03
wo/ Exp	95.69	7.57	92.50	94.57	87.71	91.59	8.45	93.64	5.11
wo/ evalGen	95.07	7.51	89.70	92.56	86.73	89.66	8.62	92.37	5.98

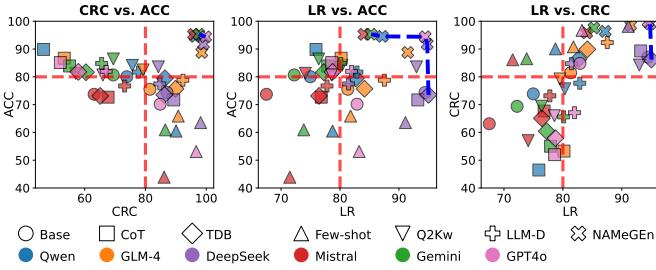


Fig. 5. Pairwise IIO comparisons. Blue lines show Pareto front; red line marks 80. Colors and shapes represent backbones and methods.

on GLM-4). Regarding balance, Fig. 3 shows that few-shot still significantly outperforms reasoning-based methods (CoT, TDB). The underlying cause is shown in Fig. 4a: reasoning-based methods perform well on common objectives (e.g., O1, O5) but struggle with others. Notably, most baselines also underperform on complex, rarer objectives (e.g., O3, O4). This imbalance highlights the limited adaptability of their limited adaptability to diverse creative objectives. Retrieval-based methods offer cultural knowledge that helps on objectives like O1 but introduce noise that negatively impacts others (O2–O5). Consequently, Q2Kw often performs worse than Base in terms of EC and EC\_std (see Table I, Fig. 3). Traditional agent-based methods such as LLM-D struggle on O4 (Fig. 4a), indicating a need for more effective information expansion beyond user input. Accordingly, NAMeGEN integrates the strengths of in-context learning with more accurate expansion and retrieval to reduce noise. Supported by explicit objective guidance and iterative dynamic optimization, NAMeGEN demonstrates strong adaptability to diverse EUOs. It achieves the highest EC and lowest EC\_std (Table I, Fig. 3), validating NAMeGEN’s effectiveness in optimizing multi-EUOs for NCB.

As shown in Table I, all baselines achieve an unsatisfactory for explainability ( $IC < 90$ ). Fig. 4b offers a fine-grained view. Firstly, in reasoning-based methods, CoT exhibits unstable and poor ACC, CRC, and LR across backbones, with no clear trend. Although TDB improves logic ( $\uparrow LR$ ), ACC and CRC remain unstable, explaining their underperformance for IIOs’ optimization (e.g., higher IC\_std for CoT in Fig.3). Few-shot enhances explanation completeness ( $\uparrow CRC$ ) but increases hallucinations ( $\downarrow ACC$ ) due to noisy examples, harming both overall interpretability (e.g., IC drops from 79.29 to 77.68 on GPT-4o, Table I). Similarly, retrieval-based methods help mitigate hallucinations ( $\uparrow ACC$ ) at the cost of completeness ( $\downarrow CRC$ ), yet the ACC gain yields better interpretability scores overall (e.g., IC from Base’s 76.29 to 80.05 on Qwen). Traditionally agent-based methods (e.g., LLM-D) improve logic and accuracy ( $\uparrow LR$ ,  $\uparrow ACC$ ) of explanation but show inconsistent completeness (CRC), resulting in partial interpretability gains without fully optimizing IIOs. **In summary**, baselines fail to balance and optimize all implicit objectives fully. In contrast, NAMeGEN’s multi-agent collaborative optimization achieves superior interpretability. As Fig. 5 shows, NAMeGEN’s results cluster near the Pareto front with high scores (80+), demonstrating effective multi-objective coordination. Overall, NAMeGEN robustly optimizes EUOs, IIOs, and their combination, validating its effectiveness on HMOO.

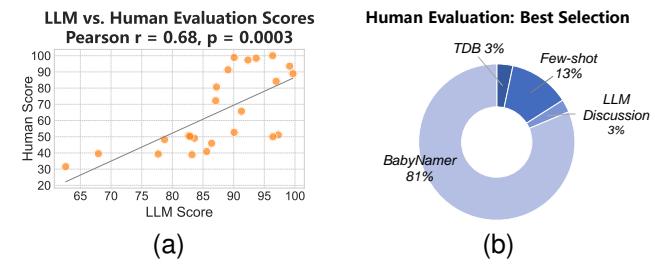


Fig. 6. The performance of three baselines (TDB, Few-shot, LLM-D) and NAMeGEN in human evaluation is compared based on Qwen. (a) analyzes the correlation between LLM and human evaluation results, with a Pearson coefficient  $r (= 0.63)$  and a  $p$ -value ( $= 0.0003$ ). (b) show the proportion of the best results selected by human annotators.

tionally agent-based methods (e.g., LLM-D) improve logic and accuracy ( $\uparrow LR$ ,  $\uparrow ACC$ ) of explanation but show inconsistent completeness (CRC), resulting in partial interpretability gains without fully optimizing IIOs. **In summary**, baselines fail to balance and optimize all implicit objectives fully. In contrast, NAMeGEN’s multi-agent collaborative optimization achieves superior interpretability. As Fig. 5 shows, NAMeGEN’s results cluster near the Pareto front with high scores (80+), demonstrating effective multi-objective coordination. Overall, NAMeGEN robustly optimizes EUOs, IIOs, and their combination, validating its effectiveness on HMOO.

2) *Human Evaluation and Correlation:* We compare three improved baselines (TDB, Few-shot, LLM-D). Fig.6a shows a Pearson correlation of  $r = 0.68 (> 0.5), p = 0.0003 (< 0.001)$ , indicating a statistically significant positive correlation between human and LLM evaluations. In Fig.6b, annotators preferred NAMeGEN in 81% of cases, far surpassing Few-shot (13%), LLM-D (3%), and TDB (3%). These results align with LLM-based evaluations, validating both NAMeGEN’s superiority and the reliability of automatic metrics.

### C. Ablation Study

We conduct ablation studies on the following settings: **wo/ Imp&Exp:** Excludes both descriptions  $I_{desc}$  and requirements  $I_{reqs}$ , as well as the implicit and explicit evaluations dependent on them. **wo/ Imp:** Excludes  $I_{reqs}$  and implicit evaluation process based on them. **wo/ Exp:** Excludes  $I_{desc}$  and explicit evaluation process based on them. **wo/ evalGen:** Removes all evaluation processes involving both  $I_{desc}$  and  $I_{reqs}$  as criteria (i.e., implicit and explicit evaluation during generation). **wo/ R:** Removes the retrieval. **wo/ evalR:** Excludes evaluation of the retrieved content.

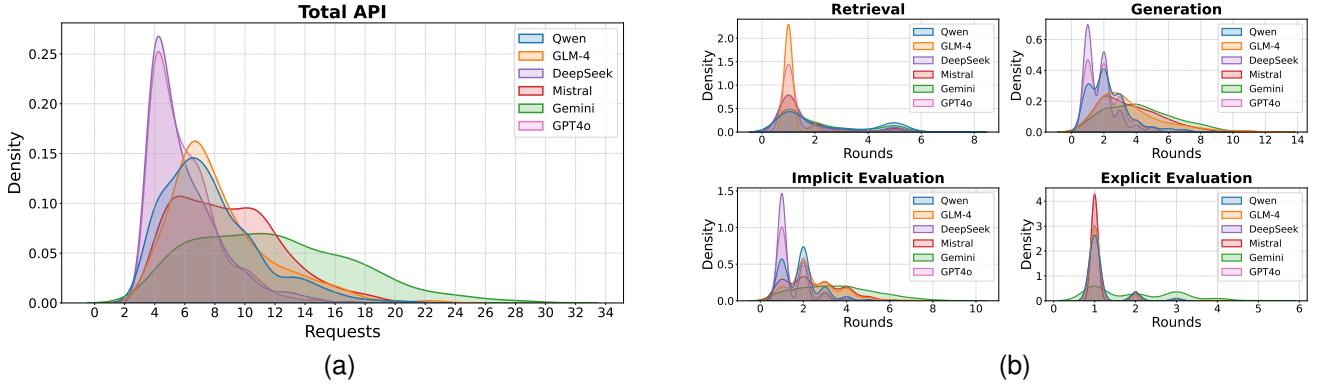


Fig. 7. Kernel density estimation (KDE) of interaction distributions across LLMs using our method. (a) shows API request counts over the full process; (b) shows iteration rounds in Retrieval, Generation, Implicit/Explicit Evaluation. Colors denote backbones; y-axis indicates estimated density.

TABLE III  
ABLATION RESULTS ON RETRIEVAL BASED ON QWEN. EMOC(C) MEASURES COMPLETENESS OF CULTURAL SIGNIFICANCE; ACC(P) REFLECTS ACCURACY OF RETRIEVED ANCIENT POETRY INFORMATION.

Method	EC(c) ↑	ACC(p) ↑
NAMeGEN	<b>96.93</b>	<b>97.60</b>
wo/ R	96.60	80.40
wo/ evalR	96.33	93.00

The results of the ablation study are summarized in Table II. In wo/ Imp&Exp, most metrics exhibit a notable performance drop, highlighting the importance of detailed objective guidance and evaluation feedback in HMOO. Moreover, the sharp decline in ACC (from 94.55 to 63.25) in the absence of Imp&Exp highlights NAMeGEN’s effectiveness in mitigating hallucination. Compared to NAMeGEN, the higher LR suggests that the additional information may introduce noise, which in turn slightly compromises logical rationality. A finer-grained ablation shows that wo/ Imp has a greater effect on implicit objective completion (IC from 92.7 to 89.54), while wo/ Exp more significantly impacts explicit objective balance (EC\_std from 5.95 to 7.57), confirming the effectiveness of our approach in both aspects. In the EC metric, wo/ Imp performs worse than wo/ Exp. This indicates that optimizing IIOs helps reach EUOs potentially, which leads to greater improvements in EC. Furthermore, given the potential correlation between EUOs and IIOs, such as there is no separation between clearer explanations and explicit objective completion. Therefore, both implicit and explicit information and evaluation influence the completeness and balance of both types of objectives. Finally, the decline in all metrics for wo/ evalGen demonstrates the effectiveness of our method’s dynamic evaluation strategy. Compared to wo/ Imp and wo/ Exp, wo/ evalGen shows a greater decrease in ACC, indicating that the evaluation stage has a better effect on hallucination mitigation. The results in Table III show that wo/ evalR has a decline in EC(c), indicating that evaluating retrieval content helps find the matcher knowledge for user needs. The wo/ R performs poorly in ACC(p), revealing that the retrieval helps mitigate the LLMs’ hallucination of domain knowledge. In summary, the ablation results validate the effectiveness of different

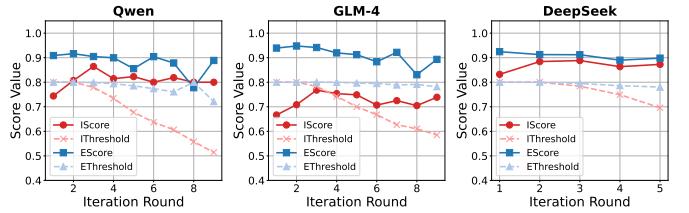


Fig. 8. Comparison of dynamic changes during NAMeGEN’s iteration process on Chinese backbones. Red and blue lines represent implicit and explicit scores; solid and dashed lines show score and threshold changes.

components of our multi-agent collaborative framework with dynamic iterative HMOO.

#### D. Efficiency Analysis

Fig. 7 shows the interaction requests and density distribution at different stages of NAMeGEN. Fig. 7a exhibits the distribution of total API requests during the task. Most require only 5-10 requests, with a peak around 6, indicating high interaction efficiency, and can converge. Fig. 7b shows the distribution of interaction rounds across different stages. In the retrieval stage, most interactions are concentrated in lower rounds, indicating quick information retrieval. In the generation and evaluation stages, most require only 2-3 rounds of generation and 1-2 rounds of evaluation, reflecting optimization efficiency. Overall, these distribution features exhibit a clear concentration with a sparse tail, where the proportion of low-frequency, high-round requests is lower. This indicates that NAMeGEN can achieve task completion with limited interaction costs in most cases, demonstrating excellent convergence and stability. Fig. 8 illustrates the dynamics score during NAMeGEN’s iterative process. The threshold (dashed line) is initially set at a higher level and gradually decreases over iterations, allowing for adaptive convergence control. The actual scores (solid lines) generally stay above these thresholds, suggesting that the method consistently meets or exceeds evaluation criteria. This indicates a strong convergence capability of our method.

#### E. Qualitative Evaluation

Since DeepSeek performs the best among Chinese LLMs, we choose it as the backbone for qualitative evaluation. As

User query	请为一个于2046年8月25日凌晨8点出生的女孩起个名字，姓甄，我们期待这个名字能够反映出“欣悦”的情绪，让她在生活中充满欢愉，即便遭遇挫折也能迅速恢复笑容。欣悦象征着由内而外的快乐与满足。Please choose a name for a baby girl born on August 25, 2046, at 8 a.m., with the surname Zhen. We hope this name can reflect the emotion of "joyfulness", allowing her to live a life filled with happiness and cheer. Even when she encounters setbacks, we wish for her to quickly regain her smile. "Joyfulness" here symbolizes a sense of happiness and contentment that radiates from within.
ETOs	[Cultural Significance (Classical Poetry), Parents' Expectations, Wuxing&Bazi], Personal Characteristics (Gender, Zodiac Sign, Birth Date, etc.), Other Requirements]
Few-shot	<b>甄欣悦。</b> 1... “欣悦”取自唐代诗人王维的《送别》中的诗句：...春草明年绿，王孙归不归。‘欣悦’象征着由内而外的快乐与满足，与诗中... 景象相呼应，寓意孩子如春草般充满生机与活力，无论何时何地都能保持乐观的心态。2... ‘欣悦’表达了父母对孩子的期望，希望她在生活中...。3...孩子的八字为...，五行缺木。...‘欣’字五行属木，‘悦’字五行属金，木金相生，有助于平衡五行，...。“...，‘4...生肖为马，属于夏季...这个名字既体现了女孩的温柔与美好，又带有活泼与快乐的氛围，与生肖马的奔放与自由相呼应。5... 读音顺口...。
NAMeGEN	<b>甄妙音。</b> 1. 孩子名字为“妙音”。其中‘妙音’二字取自明代诗人邓氏《和夫子三出西郊之作》中的诗句：好傍清池发‘妙音’。...描绘了在清澈的池水旁演奏出美妙音乐的意境。‘妙音’二字凝练了诗句中音乐与自然和谐相融的美好画面，寓意着生活中处处充满美妙的旋律与欢愉。该名字寄托了对孩子能够发现并创造生活中美好时刻的期望...。2. ‘妙音’象征着生活中美好的声音与欢愉时刻，体现了由内而外的快乐与满足。这个名字鼓励孩子...。3...八字五行为丙寅(火)丙申(火)丙辰(火土)壬辰(水土)，五行完整无缺失。妙...属水，音...属土，水土相生，...无影响....。4. 女孩，生肖虎（寅虎），出生在初秋时节。‘妙音’二字柔美动听，符合女性气质，虎为山林之王，与诗句中自然意象相呼应：初秋时节万物收获，“妙音”寓意着生活的丰盈与美好，与出生时节特征相得益彰。5... 字形结构匀称美观（左右-上下-上下结构...）。

(a)

User query	我最近要在县城里开一家甜品店，我们的主推商品为草莓奶油可颂，使用了新鲜的草莓和健康的动物奶油。请帮我设计一个宣传广告语，希望大家能够一听到就能想起来我们的甜品，最好能够考虑县城的受众，简单易懂，富有自然特色又让人印象深刻。 I'm opening a dessert shop in a county town, with strawberry cream croissants as our main product, made from fresh strawberries and healthy animal cream. Please help me design a catchy, simple, and memorable promotional slogan for the local audience.
Task type	Business slogan design
ETOs	[Brand positioning, Target Audience, Creative Expression]
Base	<b>鲜草莓，酥可颂，一口甜到心坎里！ Fresh strawberries, crispy croissants, one bite and your heart melts with sweetness!</b> <b>1.品牌定位：</b> 突出“新鲜、酥脆、甜蜜”的核心卖点，强调现做现卖的品质感，符合县城消费者对“真材实料”的偏好。 <b>Brand Positioning:</b> Focuses on "fresh, crispy, and sweet," highlighting the freshly made quality that appeals to county residents' preference for authentic ingredients. <b>2.目标受众：</b> 县城居民更注重朴实、直接的表达，用“鲜草莓”“酥可颂”具体描述产品，避免抽象词汇，让人一听就懂、一想就馋。 <b>Target Audience:</b> County residents prefer simple, direct language. "Fresh strawberries" and "crispy croissants" clearly describe the product, making it easy to understand and crave. <b>3.创意表达：</b> “鲜草莓”强调自然新鲜，暗示无添加，符合县城对健康天然的关注。 “酥可颂”用声音强的“酥”字，触发对酥脆口感的联想，记忆点强。“甜到心坎里”：用口语化表达情感共鸣，既形容味道，也暗示幸福感，容易传播。 <b>Creative Expression:</b> "Fresh strawberries" suggest natural ingredients, aligning with the county's health-consciousness. "Crispy croissants" evokes texture with the strong sound of "crispy," making it memorable. "Sweet to the heart" uses conversational language to convey both taste and emotional warmth, making it easily shareable.
NAMeGEN	<b>新鲜草莓香满城，奶油浓郁甜入心，经典可颂酥到家！ Fresh strawberries in the air, rich cream that warms the heart, and classic croissants crispy and ready to enjoy at home!</b> <b>1.品牌定位：</b> 明确地传递了甜品店的品牌核心价值。“新鲜草莓香满城”突出了草莓的新鲜度和自然香气；‘奶油浓郁甜入心’强调了奶油的高品质和经典口感；“经典可颂酥到家”则进一步强化了可颂的传统工艺和独特酥脆质感。整体广告语不仅体现了甜品店主打商品草莓奶油可颂的独特卖点，还塑造了“新鲜、健康、经典美味”的品牌形象。 <b>Brand Positioning:</b> Highlights the core values of freshness, quality, and tradition. "Fresh strawberries fill the city" emphasizes freshness; "Rich cream melts into your heart" highlights the quality of the cream; "Classic croissants bring crispy delight home" reinforces traditional craftsmanship. The slogan creates a brand image of "fresh, healthy, and classic deliciousness." <b>2.目标受众：</b> 县城的消费者以家庭为主。他们通常注重产品的外观和口感，同时也关注健康与品质。对于甜品的认知，更多偏向于传统糕点或简单易得的零食。“经典可颂”通过强调传统工艺和家庭语言，贴近了家庭消费者的日常生活。 <b>Target Audience:</b> County town families value produce appearance, taste, health, and quality. "Classic croissants" resonates with traditional craftsmanship and family life. <b>3.创意表达：</b> 广告语采用了排比句式，增强了节奏感和感染力。例如，“新鲜草莓香满城，奶油浓郁甜入心，经典可颂酥到家！”通过三个短句分别对应感官体验（嗅觉、味觉和触觉），并通过押韵（‘城’、‘心’、‘家’）增加了记忆点，让消费者更容易记住。 <b>Creative Expression:</b> The slogan uses parallel structure for rhythm and impact. "Fresh strawberries, rich cream, classic croissants" appeals to smell, taste, and texture, while the rhyme strengthens memorability.

(b)

Fig. 9. Comparison of methods on the same query and backbone (DeepSeek). (a) NCB task results. (b) Slogan design task results. Red highlights factual or logical errors, green indicates vague explanations, and blue denotes content meeting EUOs with clear logic and explanations.

shown in Fig. 9a, NAMeGEN outperforms Few-shot. We also test other tasks (slogan design), with results shown in Fig. 9b. Without EUOs, our method can also automatically resolve some EUOs (Brand positioning, Target Audience, Creative Expression). The base method is inputting the query and EUOs into DeepSeek without any extra methods. Compared to Base, our results are more aligned with user needs. These comparisons in Fig. 9 show NAMeGEN’s effectiveness in real-world applications.

## V. CONCLUSION

This paper presents the NAMeGEN framework, a novel multi-agent collaborative approach designed to address the challenges of Hybrid Multi-Objective Optimization (HMOO)

in Creative Natural Language Generation (CNLG), particularly in short-text tasks like Naming Chinese Babies (NCB). By employing a dynamic iterative hybrid multi-objective optimization strategy, NAMeGEN meets user-specified objectives while ensuring high-quality explanations in terms of factuality, completeness, and clarity. Experimental results on the CBNames benchmark demonstrate the framework’s flexibility and adaptability, consistently achieving excellent optimization results across various backbone models. Furthermore, NAMeGEN outperforms existing methods by reducing hallucination issues and providing more logically coherent and completely generated names. This work contributes to advancing the field of CNLG and lays a foundation for future research in the domain of multi-objective optimization.

## REFERENCES

- [1] OpenAI, "Hello gpt-4o," <https://openai.com/index/hello-gpt-4o/>, 2024, accessed: 2025-01-29.
- [2] ———, "Introducing openai o1," <https://openai.com/o1/>, 2024, accessed: 2024-10-28.
- [3] Gemini, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican, D. Silver, M. Johnson, I. Antonoglou, J. Schriftwieser, A. Glaese, J. Chen, E. Pitler, T. Lillicrap, A. Lazaridou, O. Firat, J. Molloy *et al.*, "Gemini: A family of highly capable multimodal models," 2024. [Online]. Available: <https://arxiv.org/abs/2312.11805>
- [4] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi *et al.*, "Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning," *arXiv preprint arXiv:2501.12948*, 2025.
- [5] Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang, T. Li, M. Ku, K. Wang, A. Zhuang, R. Fan, X. Yue, and W. Chen, "MMLU-Pro: A More Robust and Challenging Multi-Task Language Understanding Benchmark," Nov. 2024, *arXiv:2406.01574* [cs]. [Online]. Available: <https://arxiv.org/abs/2406.01574>
- [6] X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun *et al.*, "Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9556–9567.
- [7] S. Imani, L. Du, and H. Shrivastava, "MathPrompter: Mathematical reasoning using large language models," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, S. Sitaram, B. Beigman Klebanov, and J. D. Williams, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 37–42. [Online]. Available: <https://aclanthology.org/2023.acl-industry.4/>
- [8] Y. He, Y. Wu, Y. Jia, R. Mihalcea, Y. Chen, and N. Deng, "Hi-tom: A benchmark for evaluating higher-order theory of mind reasoning in large language models," *arXiv preprint arXiv:2310.16755*, 2023.
- [9] K. Gandhi, J.-P. Fraenken, T. Gerstenberg, and N. Goodman, "Understanding social reasoning in language models with language models," in *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36. Curran Associates, Inc., 2023, pp. 13 518–13 529. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/2b9efb085d3829a2aadffab63ba206de-Paper-Datasets\\_and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/2b9efb085d3829a2aadffab63ba206de-Paper-Datasets_and_Benchmarks.pdf)
- [10] X. Sun, K. Shi, H. Tang, D. Wang, G. Xu, and Q. Li, "Educating language models as promoters: Multi-aspect instruction alignment with self-augmentation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 37, no. 8, pp. 4564–4577, 2025.
- [11] N. Peng, "Controllable text generation for open-domain creativity and fairness," *arXiv preprint arXiv:2209.12099*, 2022.
- [12] T. Chakrabarty, V. Padmakumar, H. He, and N. Peng, "Creative natural language generation," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, 2023, pp. 34–40.
- [13] N. S. Keskar, B. McCann, L. R. Varshney, C. Xiong, and R. Socher, "Ctrl: A conditional transformer language model for controllable generation," *arXiv preprint arXiv:1909.05858*, 2019.
- [14] J. Qian, L. Dong, Y. Shen, F. Wei, and W. Chen, "Controllable natural language generation with contrastive prefixes," in *Findings of the Association for Computational Linguistics: ACL 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 2912–2924. [Online]. Available: <https://aclanthology.org/2022.findings-acl.229/>
- [15] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang *et al.*, "A survey on evaluation of large language models," *ACM Transactions on Intelligent Systems and Technology*, vol. 15, no. 3, pp. 1–45, 2024.
- [16] Q. Wan, S. Hu, Y. Zhang, P. Wang, B. Wen, and Z. Lu, "" it felt like having a second mind": Investigating human-ai co-creativity in prewriting with large language models," *Proceedings of the ACM on Human-Computer Interaction*, vol. 8, no. CSCW1, pp. 1–26, 2024.
- [17] T. Chakrabarty, V. Padmakumar, F. Brahman, and S. Muresan, "Creativity support in the age of large language models: An empirical study involving emerging writers," *arXiv preprint arXiv:2309.12570*, 2023.
- [18] G. Zhipeng, X. Yi, M. Sun, W. Li, C. Yang, J. Liang, H. Chen, Y. Zhang, and R. Li, "Jiuge: A human-machine collaborative chinese classical poetry generation system," in *Proceedings of the 57th annual meeting of the association for computational linguistics: system demonstrations*, 2019, pp. 25–30.
- [19] C. Yu, L. Zang, J. Wang, C. Zhuang, and J. Gu, "Charpoet: A chinese classical poetry generation system based on token-free llm," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, 2024, pp. 315–325.
- [20] S. Chatzikyriakidis and A. Natsina, "Poetry in rags: Modern greek interwar poetry generation using rag and contrastive training," in *Proceedings of the 5th International Conference on Natural Language Processing for Digital Humanities*, 2025, pp. 257–264.
- [21] S. Venkatraman, N. I. Tripto, and D. Lee, "Collabstory: Multi-llm collaborative story generation and authorship analysis," *arXiv preprint arXiv:2406.12665*, 2024.
- [22] S. Yang, Y. Ge, Y. Li, Y. Chen, Y. Ge, Y. Shan, and Y. Chen, "Seed-story: Multimodal long story generation with large language model," *arXiv preprint arXiv:2407.08683*, 2024.
- [23] P. Laban, A. R. Fabbri, C. Xiong, and C.-S. Wu, "Summary of a haystack: A challenge to long-context llms and rag systems," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 9885–9903.
- [24] H. Jin, Y. Zhang, D. Meng, J. Wang, and J. Tan, "A comprehensive survey on process-oriented automatic text summarization with exploration of llm-based methods," *arXiv preprint arXiv:2403.02901*, 2024.
- [25] S. Shang, Z. Yao, H. Fu, C. Tao, X. Chen, F. Wang, Y. Wang, Z. Ren, and S. Gao, "Unified multi-scenario summarization evaluation and explanation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 37, no. 2, pp. 991–1003, 2025.
- [26] P. N. Ahmad, Y. Liu, I. Ullah, and M. Shabaz, "Enhancing coherence and diversity in multi-class slogan generation systems," *ACM Transactions on Asian and Low-Resource Language Information Processing*, vol. 23, no. 8, pp. 1–24, 2024.
- [27] Y. Liu, D. Liu, and J. Lv, "Deep poetry: A chinese classical poetry generation system," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 09, 2020, pp. 13 626–13 627.
- [28] X. Wang, H. Jiang, Z. Wei, and S. Zhou, "Chae: Fine-grained controllable story generation with characters, actions and emotions," in *Proceedings of the 29th International Conference on Computational Linguistics*, 2022, pp. 6426–6435.
- [29] T. Wang, J. Chen, Q. Jia, S. Wang, R. Fang, H. Wang, Z. Gao, C. Xie, C. Xu, J. Dai *et al.*, "Weaver: Foundation models for creative writing," *arXiv preprint arXiv:2401.17268*, 2024.
- [30] L.-C. Lu, S.-J. Chen, T.-M. Pai, C.-H. Yu, H.-y. Lee, and S.-H. Sun, "Llm discussion: Enhancing the creativity of large language models via discussion framework and role-play," *arXiv preprint arXiv:2405.06373*, 2024.
- [31] S. Kumar, E. Malmi, A. Severyn, and Y. Tsvetkov, "Controlled text generation as continuous optimization with multiple constraints," *Advances in Neural Information Processing Systems*, vol. 34, pp. 14 542–14 554, 2021.
- [32] T. Sorensen, J. Moore, J. Fisher, M. L. Gordon, N. Mireshghallah, C. M. Ryting, A. Ye, L. Jiang, X. Lu, N. Dziri *et al.*, "Position: A roadmap to pluralistic alignment," in *Forty-first International Conference on Machine Learning*, 2024.
- [33] C. M. Pham, S. Sun, and M. Iyyer, "Suri: Multi-constraint instruction following in long-form text generation," in *Findings of the Association for Computational Linguistics: EMNLP 2024*, Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, Eds. Miami, Florida, USA: Association for Computational Linguistics, Nov. 2024, pp. 1722–1753. [Online]. Available: <https://aclanthology.org/2024.findings-emnlp.94/>
- [34] Y. Jiang, Y. Wang, X. Zeng, W. Zhong, L. Li, F. Mi, L. Shang, X. Jiang, Q. Liu, and W. Wang, "Followbench: A multi-level fine-grained constraints following benchmark for large language models," *arXiv preprint arXiv:2310.20410*, 2023.
- [35] T. B. Ward, "What's old about new ideas," *The creative cognition approach*, pp. 157–178, 1995.
- [36] R. A. Finke, T. B. Ward, and S. M. Smith, *Creative cognition: Theory, research, and applications*. MIT press, 1996.
- [37] J. Schüler, N. Baumann, A. Chasiotis, M. Bender, and I. Baum, "Implicit motives and basic psychological needs," *Journal of personality*, vol. 87, no. 1, pp. 37–53, 2019.
- [38] G. Franceschelli and M. Musolesi, "On the creativity of large language models," *AI & SOCIETY*, pp. 1–11, 2024.
- [39] T. Chakrabarty, P. Laban, D. Agarwal, S. Muresan, and C.-S. Wu, "Art or artifice? large language models and the false promise of creativity,"

- in *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 2024, pp. 1–34.
- [40] A. W. Ding and S. Li, “Generative ai lacks the human creativity to achieve scientific discovery from scratch,” *Scientific Reports*, vol. 15, no. 1, p. 9587, 2025.
- [41] X. Yi, R. Li, C. Yang, W. Li, and M. Sun, “Mixpoet: Diverse poetry generation via learning controllable mixed latent space,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 05, 2020, pp. 9450–9457.
- [42] M. Ismayilzada, C. Stevenson, and L. van der Plas, “Evaluating creative short story generation in humans and large language models,” *arXiv preprint arXiv:2411.02316*, 2024.
- [43] G. Marco, L. Rello, and J. Gonzalo, “Small language models can outperform humans in short creative writing: A study comparing slms with humans and llms,” in *Proceedings of the 31st International Conference on Computational Linguistics*, 2025, pp. 6552–6570.
- [44] H. Tan, Z. Guo, Z. Shi, L. Xu, Z. Liu, Y. Feng, X. Li, Y. Wang, L. Shang, Q. Liu, and L. Song, “ProxyQA: An alternative framework for evaluating long-form text generation with large language models,” in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikanth, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 6806–6827. [Online]. Available: <https://aclanthology.org/2024.acl-long.368/>
- [45] A. Salemi, C. Li, M. Zhang, Q. Mei, W. Kong, T. Chen, Z. Li, M. Bendersky, and H. Zamani, “Reasoning-enhanced self-training for long-form personalized text generation,” *arXiv preprint arXiv:2501.04167*, 2025.
- [46] M. Khalifa, H. Elsahar, and M. Dymetman, “A distributional approach to controlled text generation,” in *International Conference on Learning Representations*, 2020.
- [47] Y. Gu, X. Feng, S. Ma, L. Zhang, H. Gong, and B. Qin, “A distributional lens for multi-aspect controllable text generation,” in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 1023–1043.
- [48] X. Huang, Z. Liu, P. Li, T. Li, M. Sun, and Y. Liu, “An extensible plug-and-play method for multi-aspect controllable text generation,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023, pp. 15233–15256.
- [49] H. Ding, L. Pang, Z. Wei, H. Shen, X. Cheng, and T.-S. Chua, “Maclasa: Multi-aspect controllable text generation via efficient sampling from compact latent space,” in *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2023, pp. 4424–4436.
- [50] Y. Gu, X. Feng, S. Ma, L. Zhang, H. Gong, W. Zhong, and B. Qin, “Controllable text generation via probability density estimation in the latent space,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 12590–12616. [Online]. Available: <https://aclanthology.org/2023.acl-long.704/>
- [51] Y. Cao, J. Zhao, R. Zhang, H. Zou, and W. Mao, “Tara: Token-level attribute relation adaptation for multi-attribute controllable text generation,” in *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024, pp. 12570–12579.
- [52] N. Gunantara, “A review of multi-objective optimization: Methods and its applications,” *Cogent Engineering*, vol. 5, no. 1, p. 1502242, 2018.
- [53] F. Miresghallah, K. Goyal, and T. Berg-Kirkpatrick, “Mix and match: Learning-free controllable text generation using energy language models,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 401–415. [Online]. Available: <https://aclanthology.org/2022.acl-long.31/>
- [54] L. Qin, S. Welleck, D. Khashabi, and Y. Choi, “Cold decoding: energy-based constrained text generation with langevin dynamics,” in *Proceedings of the 36th International Conference on Neural Information Processing Systems*, 2022, pp. 9538–9551.
- [55] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, C. Zong, F. Xia, W. Li, and R. Navigli, Eds. Online: Association for Computational Linguistics, Aug. 2021, pp. 4582–4597. [Online]. Available: <https://aclanthology.org/2021.acl-long.353/>
- [56] C. M. Pham, S. Sun, and M. Iyyer, “Suri: Multi-constraint instruction following for long-form text generation,” *arXiv preprint arXiv:2406.19371*, 2024.
- [57] B. Wen, P. Ke, X. Gu, L. Wu, H. Huang, J. Zhou, W. Li, B. Hu, W. Gao, J. Xu *et al.*, “Benchmarking complex instruction-following with multiple constraints composition,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 137610–137645, 2024.
- [58] Q. He, J. Zeng, Q. He, J. Liang, and Y. Xiao, “From complex to simple: Enhancing multi-constraint complex instruction following ability of large language models,” *arXiv preprint arXiv:2404.15846*, 2024.
- [59] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [60] G. Li, H. A. A. K. Hammoud, H. Itani, D. Khizbulin, and B. Ghanem, “CAMEL: Communicative agents for “mind” exploration of large language model society,” in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [Online]. Available: <https://openreview.net/forum?id=3IyL2XWDkG>
- [61] Q. Wu, G. Bansal, J. Zhang, Y. Wu, B. Li, E. Zhu, L. Jiang, X. Zhang, S. Zhang, J. Liu, A. H. Awadallah, R. W. White, D. Burger, and C. Wang, “Autogen: Enabling next-gen LLM applications via multi-agent conversation,” 2024. [Online]. Available: <https://openreview.net/forum?id=tEAf9LBdgu>
- [62] S. Hong, M. Zhuge, J. Chen, X. Zheng, Y. Cheng, J. Wang, C. Zhang, Z. Wang, S. K. S. Yau, Z. Lin, L. Zhou, C. Ran, L. Xiao, C. Wu, and J. Schmidhuber, “MetaGPT: Meta programming for a multi-agent collaborative framework,” in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=VtmBAGCN7o>
- [63] S. Min, X. Lyu, A. Holtzman, M. Artetxe, M. Lewis, H. Hajishirzi, and L. Zettlemoyer, “Rethinking the role of demonstrations: What makes in-context learning work?” in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Y. Goldberg, Z. Kozareva, and Y. Zhang, Eds. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, Dec. 2022, pp. 11048–11064. [Online]. Available: <https://aclanthology.org/2022.emnlp-main.759/>
- [64] S. Yao, J. Zhao, D. Yu, I. Shafran, K. R. Narasimhan, and Y. Cao, “React: Synergizing reasoning and acting in language models,” in *NeurIPS 2022 Foundation Models for Decision Making Workshop*, 2022.
- [65] N. Shinn, F. Cassano, A. Gopinath, K. Narasimhan, and S. Yao, “Reflexion: language agents with verbal reinforcement learning,” in *Proceedings of the 37th International Conference on Neural Information Processing Systems*, 2023, pp. 8634–8652.
- [66] R. Aksitov, S. Miryoosofi, Z. Li, D. Li, S. Babayan, K. Kopparapu, Z. Fisher, R. Guo, S. Prakash, P. Srinivasan *et al.*, “Rest meets react: Self-improvement for multi-step reasoning llm agent,” in *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.
- [67] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Kütter, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, “Retrieval-augmented generation for knowledge-intensive nlp tasks,” *Advances in neural information processing systems*, vol. 33, pp. 9459–9474, 2020.
- [68] M. M. Abdollah Pour, A. Pesaranghader, E. Cohen, and S. Sanner, “Gaussian process optimization for adaptable multi-objective text generation using linearly-weighted language models,” in *Findings of the Association for Computational Linguistics: NAACL 2024*, K. Duh, H. Gomez, and S. Bethard, Eds. Mexico City, Mexico: Association for Computational Linguistics, Jun. 2024, pp. 1529–1536. [Online]. Available: <https://aclanthology.org/2024.findings-naacl.99>
- [69] H. Sun, L. Liu, J. Li, F. Wang, B. Dong, R. Lin, and R. Huang, “Comifer: Improving complex constrained instruction-following ability of large language models,” *arXiv preprint arXiv:2404.02823*, 2024.
- [70] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch, “Improving factuality and reasoning in language models through multiagent debate,” in *Proceedings of the 41st International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp, Eds., vol. 235. PMLR, 21–27 Jul 2024, pp. 11733–11763.
- [71] J. S. Park, J. O’Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein, “Generative agents: Interactive simulacra of human behavior,” in *Proceedings of the 36th annual acm symposium on user interface software and technology*, 2023, pp. 1–22.
- [72] Y. Xu, S. Wang, P. Li, F. Luo, X. Wang, W. Liu, and Y. Liu, “Exploring large language models for communication games: An empirical study on werewolf,” *arXiv preprint arXiv:2309.04658*, 2023.
- [73] Z. Wang, S. Mao, W. Wu, T. Ge, F. Wei, and H. Ji, “Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration,” in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics*, 2024.

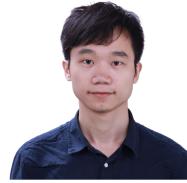
- Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 2024, pp. 257–279.
- [74] S. Chen, Y. Liu, W. Han, W. Zhang, and T. Liu, “A survey on llm-based multi-agent system: Recent advances and new frontiers in application,” 2025. [Online]. Available: <https://arxiv.org/abs/2412.17481>
  - [75] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large language models are zero-shot reasoners.” *Advances in neural information processing systems*, vol. 35, pp. 22 199–22 213, 2022.
  - [76] C. Yang, X. Wang, Y. Lu, H. Liu, Q. V. Le, D. Zhou, and X. Chen, “Large language models as optimizers,” *arXiv preprint arXiv:2309.03409*, 2023.
  - [77] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
  - [78] R. Jagerman, H. Zhuang, Z. Qin, X. Wang, and M. Bendersky, “Query expansion by prompting large language models,” *arXiv preprint arXiv:2305.03653*, 2023.
  - [79] C.-M. Chan, W. Chen, Y. Su, J. Yu, W. Xue, S. Zhang, J. Fu, and Z. Liu, “Chateval: Towards better llm-based evaluators through multi-agent debate,” *arXiv preprint arXiv:2308.07201*, 2023.
  - [80] H. Li, Q. Dong, J. Chen, H. Su, Y. Zhou, Q. Ai, Z. Ye, and Y. Liu, “Llms-as-judges: a comprehensive survey on llm-based evaluation methods,” *arXiv preprint arXiv:2412.05579*, 2024.
  - [81] V. Adlakha, P. BehnamGhader, X. H. Lu, N. Meade, and S. Reddy, “Evaluating correctness and faithfulness of instruction-following models for question answering.” *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 681–699, 2024.
  - [82] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang *et al.*, “Qwen technical report,” *arXiv preprint arXiv:2309.16609*, 2023.
  - [83] G. Team, P. Georgiev, V. I. Lei, R. Burnell, L. Bai, A. Gulati, G. Tanzer, D. Vincent, Z. Pan, S. Wang *et al.*, “Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context,” *arXiv preprint arXiv:2403.05530*, 2024.
  - [84] A. Liu, B. Feng, and et al., “Deepseek-v3 technical report,” *arXiv preprint arXiv:2412.19437*, 2024.
  - [85] A. Q. Jiang, A. Sablayrolles, and et al., “Mistral 7b,” 2023. [Online]. Available: <https://arxiv.org/abs/2310.06825>
  - [86] K. Team, A. Du, and et al., “Kimi k1. 5: Scaling reinforcement learning with llms,” *arXiv preprint arXiv:2501.12599*, 2025.
  - [87] S. Zhong, Z. Huang, S. Gao, W. Wen, L. Lin, M. Zitnik, and P. Zhou, “Let’s think outside the box: Exploring leap-of-thought in large language models with creative humor generation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 13 246–13 257.



**Shanlin Zhou** is currently a second-year Ph.D. candidate in the Department of Computer Science and Technology at Tongji University, Shanghai, China. Her research interests lie in capability optimization exploration for large language models (LLMs), particularly in intelligent agent applications and reasoning enhancement, hallucination mitigation, and AI alignment. His work aims to improve the safety and practicality of LLMs for real-world applications.



**Xinpeng Wang** is currently a fourth-year Ph.D. candidate in the Department of Computer Science and Technology at Tongji University, Shanghai, China. His research interests lie in controllable text generation for large language models (LLMs), particularly in reasoning enhancement, hallucination mitigation, and AI alignment. His work aims to improve the reliability and controllability of LLMs for real-world applications.



**Jianxun Lian** is now a senior researcher at Microsoft Research Asia. He received his Ph.D. degree from University of Science and Technology of China. His research interests include AI Anthropomorphism, LLM-based agents, user modeling, and recommender systems. He has published several papers on top-tier international conferences such as KDD, ACL, WWW, NeurIPS, IJCAI, and SIGIR, and actively contributes as a program committee member for these conferences.



**Zhenghao Liu** is an Associate Professor at Northeastern University and an adjunct researcher at the Natural Language Processing Laboratory of Tsinghua University (THUNLP). He received his B. Eng. degree from Northeastern University in 2016 and subsequently pursued his Ph.D. at the State Key Laboratory of Intelligent Technology and Systems, Tsinghua University, under the supervision of Professor Maosong Sun. He obtained his Ph.D. in 2021. His research interests lie at the intersection of natural language processing, information retrieval, and large language models, with a particular focus on retrieval-augmented generation.



**Laks V.S. Lakshmanan** is a professor of Computer Science at UBC, Vancouver, Canada. His research interests span a wide spectrum of topics in data management, integration, cleaning, and warehousing; data mining; semi-structured and unstructured data; big graphs, social networks, and social media; ML, NLP; and efficient deep learning. He is an ACM Distinguished Scientist and has won several awards, including Best Paper awards and Distinguished Reviewer awards. He has served on most top conferences and journals in his areas of research, as a senior PC member, meta-reviewer, general chair, and as associate editor.



**Xiaoyuan Yi** is currently a senior researcher at Microsoft Research Asia. He obtained his bachelor degree and Ph.D. in computer science from Tsinghua University in 2021. His research interests are Natural Language Generation (NLG) and Societal AI. He has published 30+ papers at top-tier AI venues like ICLR, NeurIPS, ICML, ACL, EMNLP, and so on, and contributes as a program committee member for these conferences.



**Yongtao Hao** received his Ph.D. from Shanghai Jiao Tong University in 1999. He is currently a Professor and Ph.D. supervisor in the Department of Computer Science and Technology, School of Electronics and Information Engineering, Tongji University. His research interests include data processing, intelligent algorithms and models, and automation systems.