# Exploration of Large Language Models in Modern Community Planning

Jiexiang Lan, Yuheng Li, Ruixin Ma

## Abstract

*In contemporary urban contexts, modern community development confronts a range of challenges exacerbated by the dual pressure of rapid urbanization and unstable economic environment. Against this backdrop, the swift progress of artificial intelligence technology has catalyzed a significant shift in the community planning industry, propelling its intelligent transformation. This paper explores the fine tuning of large language models using LoRA, Prefix Tuning, and Prompt Tuning to develop specialized LLMs tailored for community planning applications. It also investigates the potential applications of LLMs in modern community planning, including guiding automated community modeling and simulating resident hearings through agents. In future, we will continue to refine the integration of these applications and advance the implementation of projects to achieve modern intelligent community planning tools through collaboratively construction.*

## 1. Introduction

With the rapid progression of urbanization and the ongoing fluctuations in the economic landscape, a range of social issues have emerged in the context of community development. For instance, community dysfunction often arises due to the high cost of living and the mismatch of residents' expectations. However, traditional static planning models, which are often predicated on rigid, long term projections and inflexible design principles, prove to be costly to implement and lack the adaptability necessary to respond to the rapidly evolving community. Recently, the advent of large language models has introduced new possibilities for dynamic and responsive community planning. However, most pre-trained large language models lack sufficient knowledge in the field of community planning. Therefore, they often struggle with complex tasks in this field. To surmount this challenge, it is essential to employ fine tuning techniques to develop large language models tailored for specific domains and enhance their intelligent analysis capabilities. In this article, we summarize and experiment with various fine tuning techniques for large language models and explore the application of the tailored models in modern community planning.

## 2. Fine Tuning Techniques

### 2.1. Full Fine Tuning

In the initial large pre-trained models proposed by OpenAI, full fine tuning was required to adapt the models to specific tasks, which laid the foundation for subsequent var-ious fine tuning techniques [11]. During full fine tuning, the model weights are initialized to the pre-trained weights $\phi_0$, and the parameters $\phi$ are continuously updated by optimizing the log-likelihood function of the training samples. As shown in the formula, considering the Q&A style training samples, x represents the question, y represents the answer, and $y_{<t}$ denotes $\{y_1, y_2, \ldots, y_t\}$.

$$\max_{\phi} \sum_{(x,y) \in Z} \sum_{t=1}^{|y|} log(P_{\phi}(y_t|x, y_{<t}))$$

This method has a significant drawback that training for each specific task requires saving a full set of parameters $\phi$. However, the number of parameters $|\phi|$ in modern pre-trained models is usually very large, consuming a lot of storage space. Therefore, full fine tuning is generally no longer applicable in the modern NLP field, and we need more efficient and space friendly methods to update parameters.

### 2.2. Low Rank Adaptation

The core concept of LoRA involves inserting trainable low-rank decomposition matrices into each layer of the Transformer architecture, while the remaining pre-trained model weights are kept unchanged.
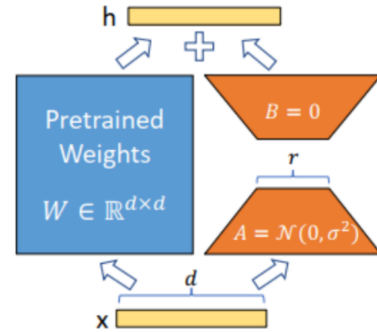


Figure 1. Explanation of LoRA

Specifically, for a pre-trained weight matrix $W_0 \in R_d \times k$, LoRA represents its update using a low rank decomposition $W_0 + \delta W = W_0 + BA$, where $B \in R_d \times r$ and $A \in R_r \times k$ are the trainable low-rank matrices [3]. During training, only the matrices A and B are optimized, while $W_0$ does not update its gradient. This process significantly reduces the number of trainable parameters required for adaptation due to $r \ll \min(d, k)$.

### 2.3. Prefix Tuning

Prefix tuning adds contextual collocations as prefixes before materials in tasks like "table-to-text", "translation". Since only the prefix needs to be stored for each task, this method is more modular and flexible compared to LoRA fine tuning. Specifically, prefix tuning initializes a trainable matrix $P_\theta$ of dimension $|P_{idx}| \times dim(h_i)$ to store the prefix parameters [9]. The optimization objective function of prefix tuning is the same as that of full fine tuning. However, only the parameters of the prefix portion are updated while the rest of the pre-trained parameters are kept frozen.

$$h_i = \begin{cases} P_\theta[i,:] & if\ i \in P_{idx} \\ LM_\phi(z_i, h_{<i}), & o.w. \end{cases}$$

### 2.4. P-tuning & P-tuning V2

P-Tuning is a prompt-based fine tuning method that guides large language models to perform downstream tasks by optimizing prompt words. Unlike prefix tuning, which simply adds fixed prefixes, p-tuning inserts trainable tokens at strategic locations and optimizes both the content and placement of these prompts. This approach provides greater flexibility and adaptability to different tasks.

However, this method lacks deep prompt optimization. Continuous prompts are only inserted into the input embedding sequence of the Transformer. The embeddings in subsequent Transformer layers are calculated by the previous Transformer layers, posing numerical challenges to the stability of training optimization. To address this issue, p-tuning v2 was proposed: This method integrates prompt tokens as inputs at every layer, rather than solely at the input layer [10]. This design not only increases the number of learnable parameters but also enables the prompts embedded in deeper layers to exert more direct influence on the predictions, thereby enhancing the overall adaptability and effectiveness of the model.

### 2.5. Experiments

In this study, we deploy the local large language model DeepSeek-R1-Distill-Qwen-1.5B and explore fine tuning methods including LoRA, prefix tuning, and prompt tuning. We implement the fine tuning method using the PyTorch and Hugging Face PEFT libraries.

#### 2.5.1 Data collection

Given that fine tuning large language models require training samples in Q&A format and we typically only have access to extensive expository texts, we can leverage the comprehension capabilities of common pre-trained large language models to generate the required training data. Specifically, we initially select two professional textbooks in the field of urban community planning as our source PDF documents. By employing the Python library pdfplumber, we can split two PDF documents into multiple text slices. Subsequently, using a customized prompt, we invoke the Alibaba Cloud BaiLian API to generate Q&A pairs for each slice. Finally, these pairs are saved as a JSON file as training data. Details can be accessed in file devision.ipynb.

#### 2.5.2 Training loss

We utilize the Hugging Face PEFT library to implement LoRA, prefix tuning, and prompt tuning. Through a series of hyperparameter tuning, we achieve the following training loss results.
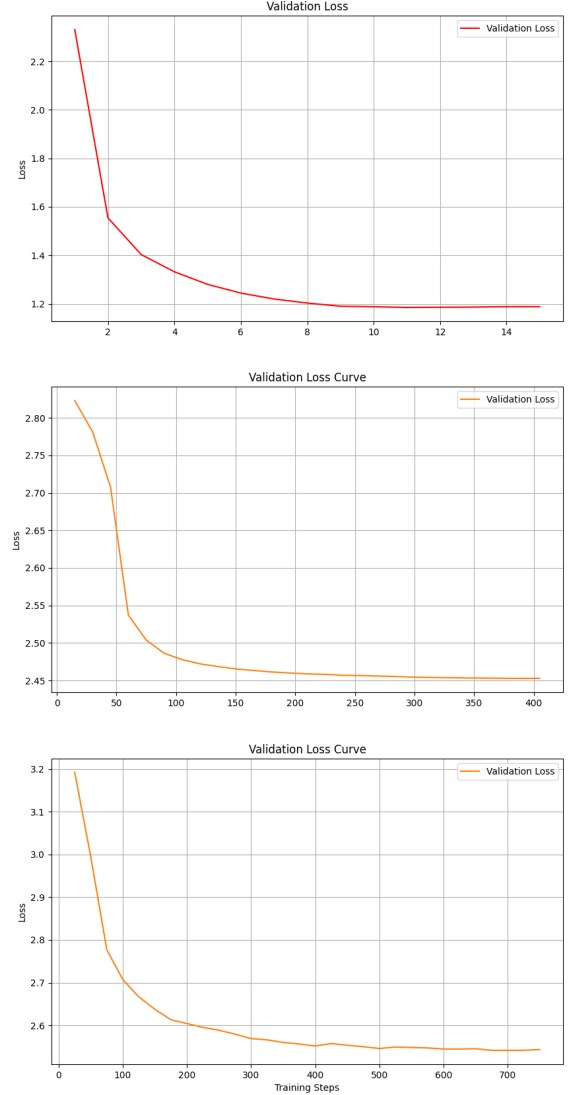


Figure 4. Prompt tuning

#### 2.5.3 Loss landscape visualization

Neural network training relies on finding good minimizer of highly non-convex loss functions. The loss landscape serves as a visualization method for evaluating the performance of neural networks and loss functions, measuring the relationship between training parameters and the loss function.

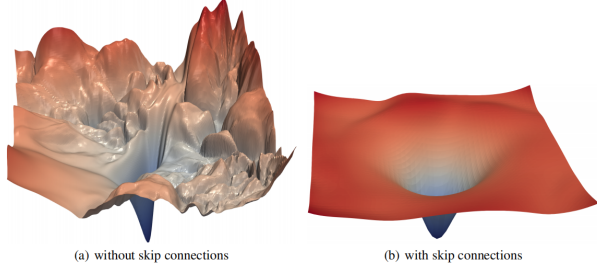$$f(\alpha, \beta) = L(\theta^* + \alpha\delta + \beta\eta)$$

Figure 5. Loss landscape: (a)poor training performance (b) good training performance

Specifically, choose a center point $\theta^*$ in the parameter space, and choose two orthogonal direction vectors, $\delta$ and $\eta$. By examining the structure of the loss landscape, we can explore the training effectiveness of fine tuning for large language models. [6] demonstrates that smooth and flat landscapes with wide minima typically lead to stable models, while rough and highly curved landscapes with narrow minima often result in slower convergence and poorer generalization. In this study, we utilize LoRA fine tuning as a case in point. We randomly select four distinct sets of stochastic parameter directions. Our observations reveal that the loss landscape associated with each parameter direction set exhibits a notable degree of smoothness, characterized by global minima that are easily convergent. These findings suggest that the training outcomes for LoRA are highly promising.
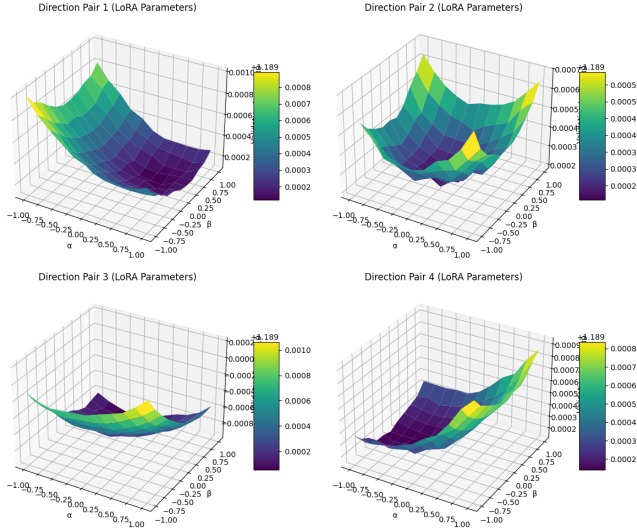


Figure 6. Loss landscape for LoRA

### 2.5.4 Perplexity

Intuitively, large language models after fine tuning are better suited for specific downstream tasks, exhibiting lower perplexity. However, as indicated in the table, some of the perplexity of the large language models actually increases after fine tuning. This counterintuitive result is attributed to catastrophic forgetting [7]. Catastrophic forgetting refers

| | Base | LoRA | Prefix tuning | Prompt tuning |
|---|---|---|---|---|
| Perplexity | 1.11 | 1.17 | 1.09 | 1.18 |
| Avg_loss | 0.1579 | 0.1064 | 0.0885 | 0.1697 |

Table 1. Comparison of perplexity and average loss

to the situation where a model forgets previously learned knowledge when it learns new tasks or data, leading to a significant degradation in performance measurement such as perplexity. This phenomenon is particularly common in the fine tuning of large language models, as it need to adapt to new domains while retaining existing knowledge.
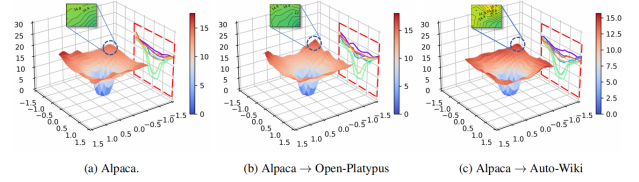


Figure 7. Catastrophic forgetting & flatness of loss landscape

Meanwhile, this study also reveals the relationship between the flatness of the model loss landscape and the extent of catastrophic forgetting. Author pointed out that as the data increases in (a), (b), and (c), the phenomenon of catastrophic forgetting becomes more severe, and the disturbance of their contours becomes more obvious. In this study, we employ training datasets of sizes 100,400,1000,2000,3000 to plot the loss landscape and verify the phenomenon, with the results presented as follows. We utilize LoRA fine tuning as a case in point.
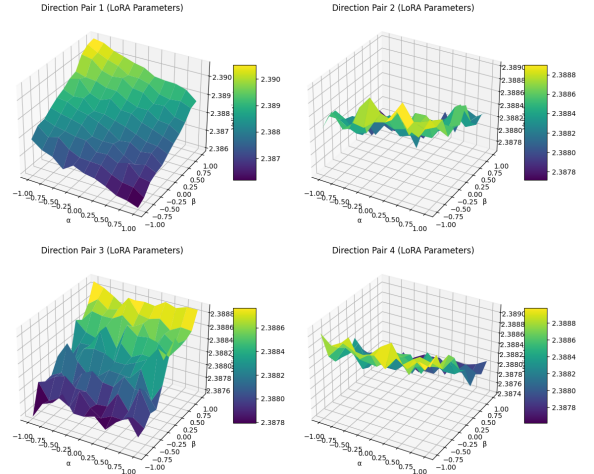


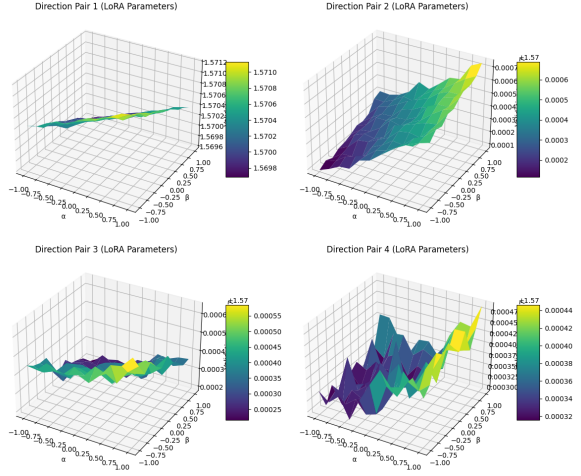Figure 8. Loss landscape (sample size = 100)
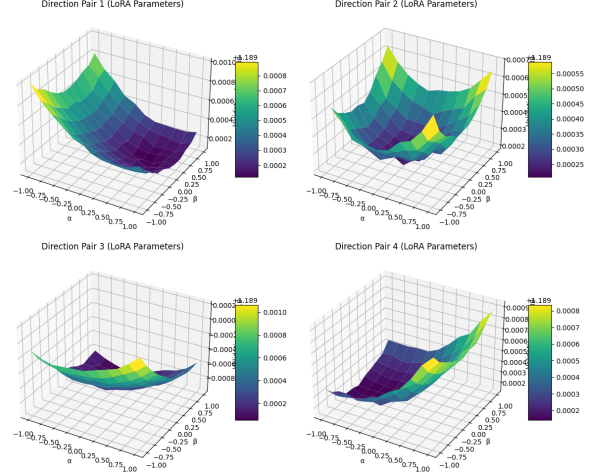
3

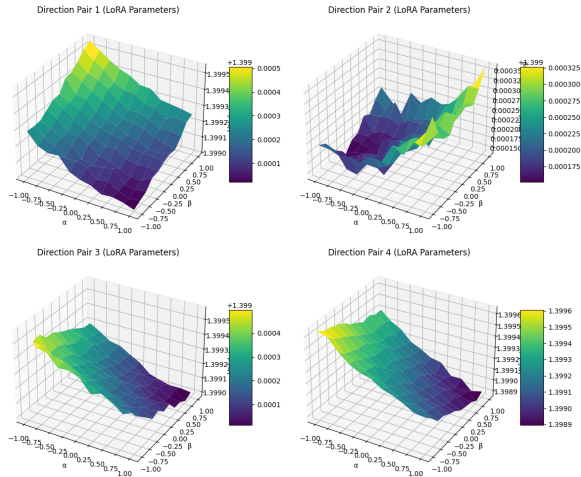Figure 9. Loss landscape (sample size = 400)



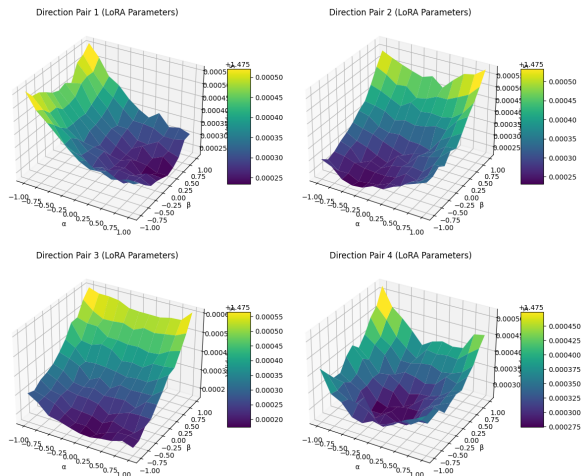Figure 12. Loss landscape (sample size = 3000)



Figure 10. Loss landscape (sample size = 1000)

Observations indicate that when the dataset is small, the loss landscape becomes more perturbed as the data size increases, which matches the result of [7]. However, this phenomenon reverses once the dataset size exceeds 1000 samples. This occurrence is likely due to the complexity of loss landscape flatness, which is influenced by multiple factors. As the amount of data grows and the training epochs deepen, regularization methods such as mixed precision training (fp16=True) and weight decay (weights_decay=True) become more effective. Furthermore, the learning rate adjustments by the Adam optimizer help to alleviate catastrophic forgetting phenomenon, which contradicts theoretical phenomenon. Therefore, we need to control more training parameters and methods to achieve more robust validation.

Overall, taking into account the catastrophic forgetting phenomenon, the perplexity of our fine tuned models is comparable to that of the base model, which indicates the reliability of the fine tuning outcomes.

## 3. Application

In this section, we explore the application of the fine tuned large language models in modern community planning and validate the results presented in the previous section from the perspective of practical tasks.

### 3.1. Automatic modeling

Professional planning tools are difficult for community members to use due to their high learning costs. We try to utilize the fine tuned large language models to assist in automated modeling in order to reduce the difficulty. We use Unity 3D as the modeling tool and select a community in Tokyo, Japan, as the modeling target. We used the number of different types of buildings, street structure, building area, and building center locations as parameters to assist the large language model in understanding the scene. For specific details, please refer to the design of the prompt engineering in the code chat.py.



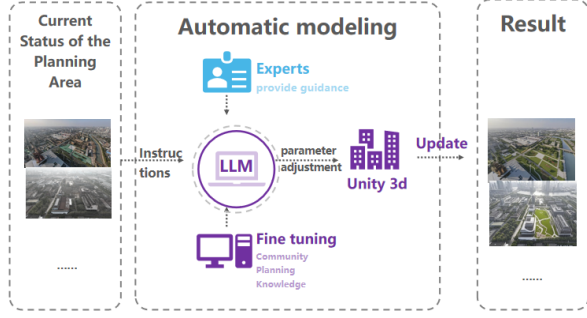Figure 11. Loss landscape (sample size = 2000)

Figure 13. Workflow for automatic modeling

After the large language model understands the current status of the community, community staff or experts can propose modifications. Combining these suggestions with fine tuning data, the large language model interpret the changes needed in the community modeling parameters. It then pass the revised parameters back to Unity 3D for visualization, showcasing the planned effects after modification. Unity 3D also supports further personalized adjustments to the modeling components.
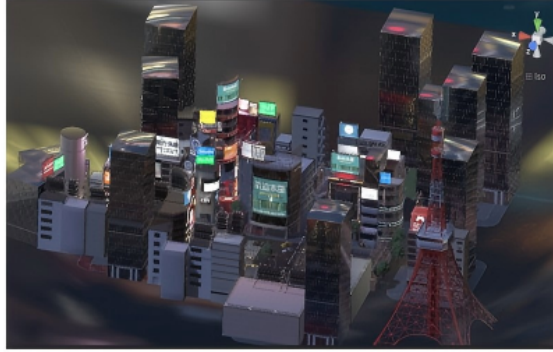


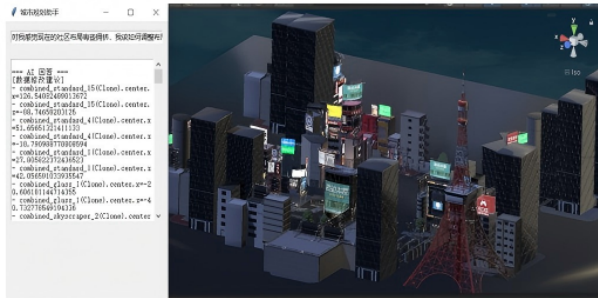Figure 14. Primal community modeling



Figure 15. After planning

As shown in the figure, we propose the requirement: "The community needs to build more office buildings recently, and the space allocation for the buildings is a bit tight. Please help me adjust the layout." We then obtain a new community plan. Detailed experimental results can be seen in demo.mp4 video.

## 3.2. Community agents

Existing intelligent planning tools fail to adequately incorporate public participation. Meanwhile, conducting community censuses is both time consuming and labor intensive, yet typically yields only limited insights. To address this issue, we employ the generation of LLM community agents to simulate public hearing reactions and represent public opinions. We conduct an initial exploration by employing parametrized prompt engineering to equip agents with detailed community background information. We assign agents diverse identities, such as doctors and merchants, to guide them in assessing the rationality of new community proposals. However, we find that agents tend to either universally accept or reject modifications after multiple rounds of generation. This phenomenon aligns with the agent misalignment issues [1]: due to the simple design of prompt, interactions and collaborations among agents during task execution fail, leading to problems such as conversation resets, task derailments, information withholding, and failures to ask for clarification. In our task, this manifests as agents unanimously agreeing or disagreeing.



Figure 16. Dilemma of simple prompt engineering

### 3.2.1 Related work

[8] This article introduces EconAgent, a large language model-based agent designed for macroeconomic simulation. EconAgent possesses human-like characteristics, including perception, memory, and decision-making capabilities, enabling it to make decisions similar to those of humans in real economic environments .The construction of this agent is divided into three main modules. First, a simulation environment is created that incorporates various market dynamics driven by the agents' decisions regarding work and consumption. Second, heterogeneous agents with distinct decision making mechanisms are created to reflect the differences among individuals in the real world, rather than assuming that all agents have the same behavior patterns. Third, a memory module is introduced, allowing agents to review past decisions and market dynamics,

5

learn from historical data, and adjust their strategies to better adapt to current and future economic environments. The modules proposed in this article for constructing agents can be applied to the simulation of public participation in modern community planning.

### 3.2.2 Solution

We constructs decision - making feature vectors for each agent (imaginativeness, environmental concern, economic concern, community concern, and conservatism), forming a mathematical basis for personalized decision making. The vector design draws on Holland's Occupational Theory[2] and the Social Values Scale [12], quantifying agents' decision preferences via orthogonal dimensions. For instance, the architect agent has high imaginativeness (0.8) and moderate environmental concern (0.6), while the shopkeeper agent has high conservatism (0.8) and economic concern (0.85).This differentiated modeling enables the system to simulate the cognitive differences of various interest groups in real society, lifting the opposition rate from under 5% to a reasonable 35 - 40%.

To address the mapping issue between feature vectors and decision behaviors, a dynamic matching degree algorithm has been developed. The algorithm first identifies the type of proposal (economic, environmental, community, or transportation) through semantic analysis and then activates the corresponding feature-weight matrix. For instance, in the case of economic proposals, the algorithm assigns a positive weight of 0.8 to the "economic concern" dimension and a negative weight of -0.2 to the "conservatism" dimension. The final matching degree for each agent is calculated as the dot product of the feature vector and the weight matrix, followed by normalization using the Sigmoid function. This approach ensures that entrepreneurs with high economic concern achieve a matching degree of 0.92 for economic proposals, while university students with high environmental concern have a matching degree of only 0.31, thereby naturally forming a spectrum of opinions.

To enhance the authenticity of decision-making, a cognitive dissonance checking mechanism has been introduced. When there is a significant discrepancy between an agent's vote and its feature-matching degree (e.g. an agent with a high matching degree votes against), the system automatically corrects the vote with an 80% probability, while retaining 20% randomness to simulate the irrational elements in human decision-making. Experimental data indicate that this mechanism increases the correlation between decision behavior and feature vectors, effectively resolving the issue of "fake opposition."

In the voting reason generation phase, a four-dimensional structured framework (professional opinion, personal impact, risk assessment, and final conclusion) is employed, compelling agents to present arguments from multiple perspectives.This design not only avoids simplistic expressions of stance but also encourages detailed professional analyses in opposition votes. For example, a doctor agent may question transportation planning from the perspective of medical accessibility, while a shopkeeper agent may analyze business impacts based on operating costs. This structured argumentation approach increases the information density of opposition reasons, effectively simulating rational debates in real-world public discussions.



Figure 17. Agents hearing case



Figure 18. Agents hearing case



Figure 19. Agents hearing case

## 4. Discussion

This paper investigates the phenomenon of catastrophic forgetting that occurs during the fine tuning of large language models.In future model improvements, Elastic Weight Consolidation technology [5] can be employed to compute the importance of each parameter for old tasks and introduce regularization, thereby alleviating the phenomenon of catastrophic forgetting and enhancing the quality of fine-tuned models. Regarding the 3D modeling aspect, this study adopts a Unity 3D rendering style that leans towards a game-like aesthetic. In the future, the incorporation of 3D Gaussian splatting modeling technology and

rendering simulation could be explored to improve modeling quality [4]. Additionally, we hope to introduce cloud rendering technology [13] in conjunction with automated modeling and agent-based hearings to construct an automated intelligent community system that encourages public participation. This approach is expected to facilitate the implementation and validation of the project in real world circumstances.

# References

[1] Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A. Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica. Why do multi-agent llm systems fail? *arXiv preprint arXiv:2503.13657v2*, 2025. 5

[2] John L. Holland. *Making Vocational Choices: A Theory of Vocational Personalities and Work Environments*. Psychological Assessment Resources, Odessa, FL, 3rd edition, 1997. 6

[3] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations (ICLR)*, 2022. arXiv preprint arXiv:2106.09685. 1

[4] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42 (4), 2023. 7

[5] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Demis Hassabis, Shakir Mohamed, and Matt Botvinick. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. 6

[6] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *Neural Information Processing Systems*, 2018. 3

[7] Hongyu Li, Liang Ding, Meng Fang, and Dacheng Tao. Revisiting catastrophic forgetting in large language model tuning. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4297–4308, Miami, Florida, USA, 2024. Association for Computational Linguistics. 3, 4

[8] Nian Li, Chen Gao, Mingyu Li, Yong Li, and Qingmin Liao. Econagent: Large language model-empowered agents for simulating macroeconomic activities. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15523–15536, Bangkok, Thailand, 2024. Association for Computational Linguistics. 5

[9] Xiang Lisa Li and Percy 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)*, pages 4340–4354. Association for Computational Linguistics, 2021. 2

[10] Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *CoRR*, abs/2110.07602, 2021. 2

[11] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training, 2018. 1

[12] Shalom H. Schwartz, Jan Cieciuch, Michele Vecchione, Eldad Davidov, Ronald Fischer, Christin Beierlein, Alfonso Ramos, Markku Verkasalo, Jan-Erik Lönnqvist, Kivanc Demirutku, Oya Dirilen-Gumus, and Michal Konty. Refining the theory of basic individual values. *Journal of Personality and Social Psychology*, 103(4):663–688, 2012. 6

[13] Nan Zhang, Ya Dai, Houkun Cui, Xiaofeng Zhang, Weizhou Xu, and Chao Ye. High precision and fast cloud rendering based on webgl technology and transmission data model algorithm. In *Proceedings of the 4th International Conference on Computer Science and Management Technology (ICC-SMT '23)*, pages 1–5. ACM, 2023. 7