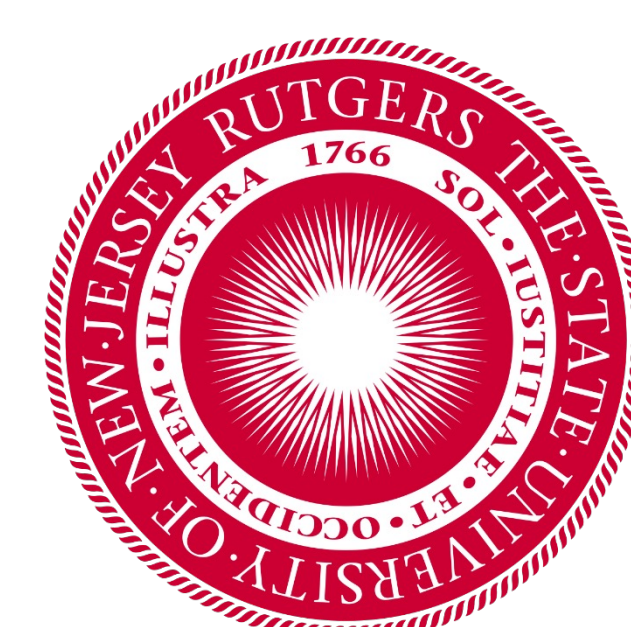# Large Language Models for Generative Recommendation: A Survey and Visionary Discussions

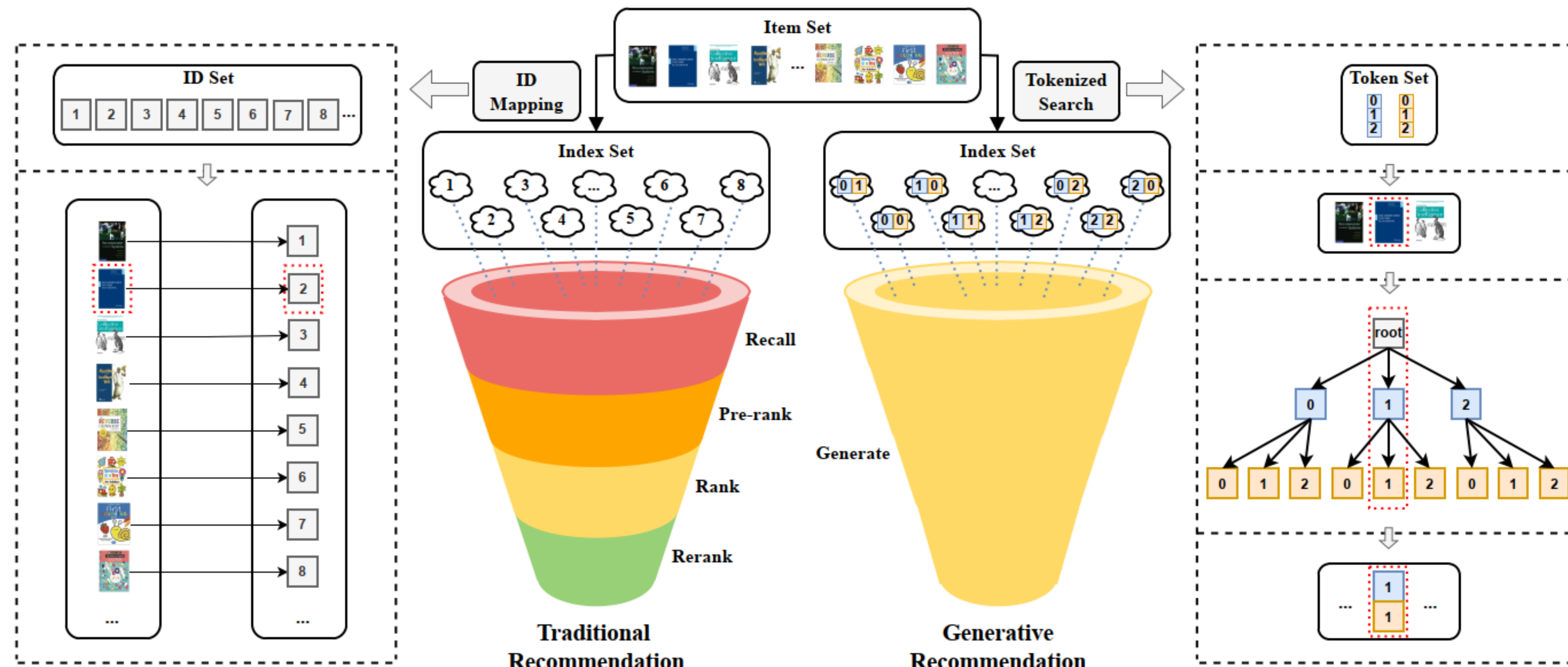**Lei Li**[1], Yongfeng Zhang[2], Dugang Liu[3], Li Chen[1]

[1]{csleili, lichen}@comp.hkbu.edu.hk, [2]yongfeng.zhang@rutgers.edu, [3]dugang.ldg@gmail.com

## Why Generative Recommendation

- Huge number of items on recommendation platforms
- Computationally expensive score calculation for each
- Multi-stage filtering to narrow down candidates
  - Simple methods at early stage
  - Complex models at final stage
- Gap between academic research and industrial applications



- Simplify recommendation process to one stage
  - Directly generate items for recommendation
  - Implicitly enumerate all items
- Use finite tokens to represent infinite items
  - # tokens = 1000
  - ID length = 10 tokens
  - # items = $1000^{10} = 10^{30}$
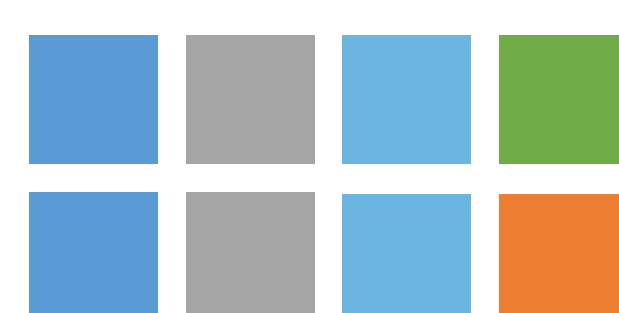
## ID Creation Methods

### Generalized Definition of ID

> *An ID in recommender systems is a sequence of tokens that can uniquely identify an entity, such as a user or an item.*

- Embedding ID
- Sequence of numerical tokens
  - <item><_><73><91>
- Sequence of word tokens
  - Item title
    - The Lord of the Rings
  - Item description
  - News article
  - Sequence of meaningless words
    - Ring epic journey fellowship adventure [1]

### LLM-compatible IDs

- Retain collaborative information of IDs in LLM environment
  - User-user
  - Item-item
  - User-item
- Short and exact representations of IDs
  - Similar users/items share more tokens
  - Remaining tokens distinguish them
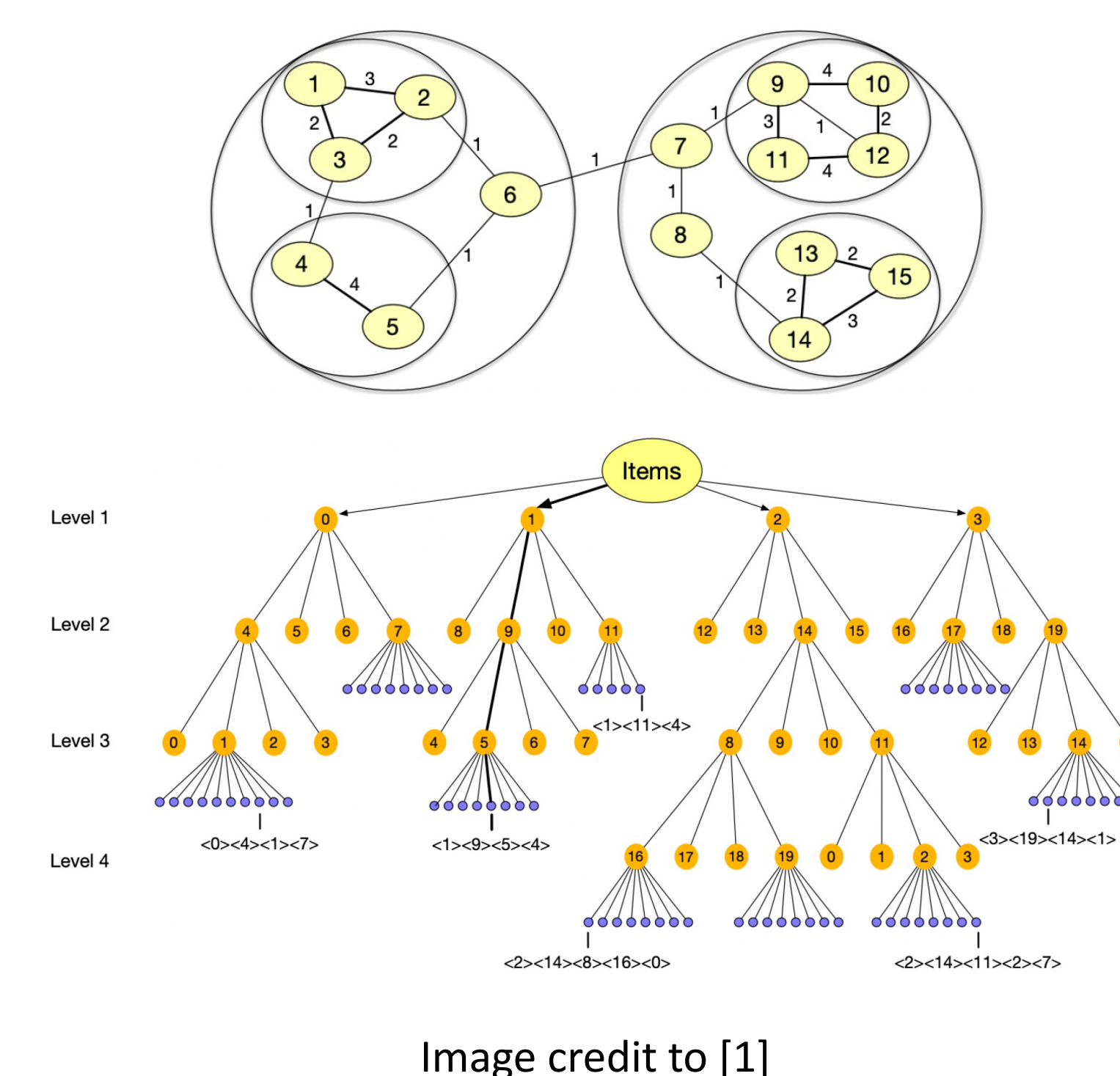
### Spectral Clustering



Image credit to [1]
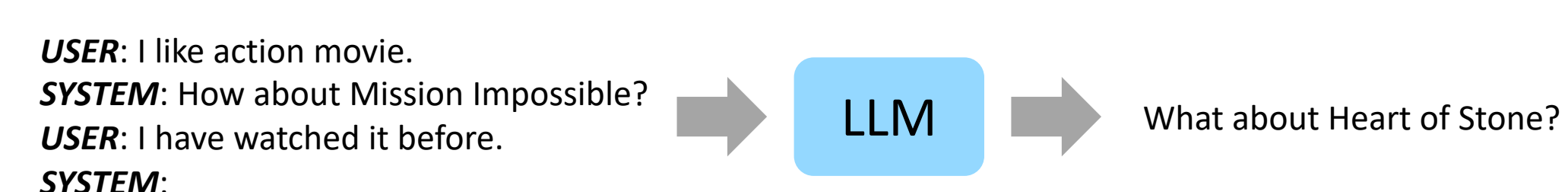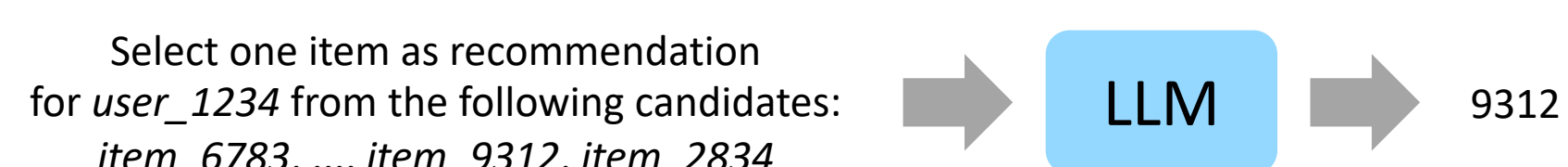
## How to Do Generative Recommendation

### Rating Prediction

How would *user_1234* rate *item_5678*? → LLM → 4.12

### Top-*N* Recommendation

Select one item as recommendation for *user_1234* from the following candidates: *item_6783, ..., item_9312, item_2834* → LLM → 9312

### Sequential Recommendation

Given *user_1234*'s interaction history *item_3456, ..., item_4567, item_5678*, predict the next item that the user will click → LLM → 6789

### Explainable Recommendation

Explain to *user_1234* why *item_5678* is recommended → LLM → The movie is top-notch

### Review Generation

Generate a review for *user_1234* about *item_5678* → LLM → The hotel is located in ...

### Review Summarization

Please summarize the following review that *user_1234* wrote for *item_5678*: the hotel is located in ... → LLM → great location

### Conversational Recommendation

*USER*: I like action movie.
*SYSTEM*: How about Mission Impossible?
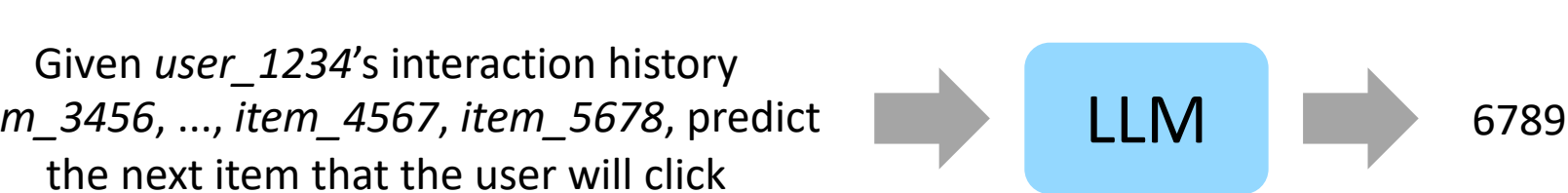*USER*: I have watched it before.
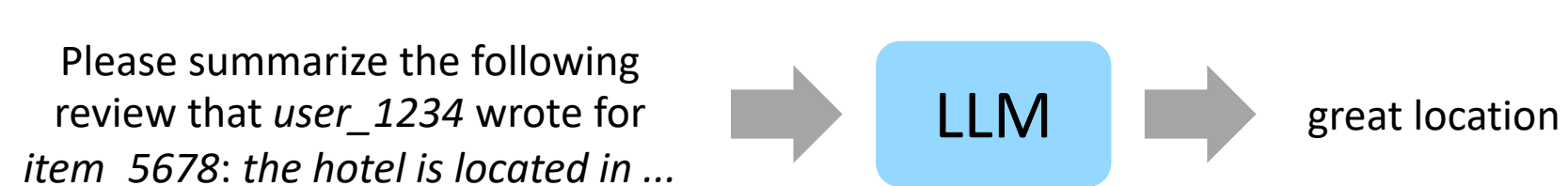*SYSTEM*: → LLM → What about Heart of Stone?

### Evaluation

- Automatic evaluation
  - RMSE and MAE for rating prediction
  - NDCG, precision and recall for top-*N* recommendation and sequential recommendation
  - BLEU and ROUGE for text similarity
- Online A/B tests
- Human evaluation

## Challenges and Opportunities

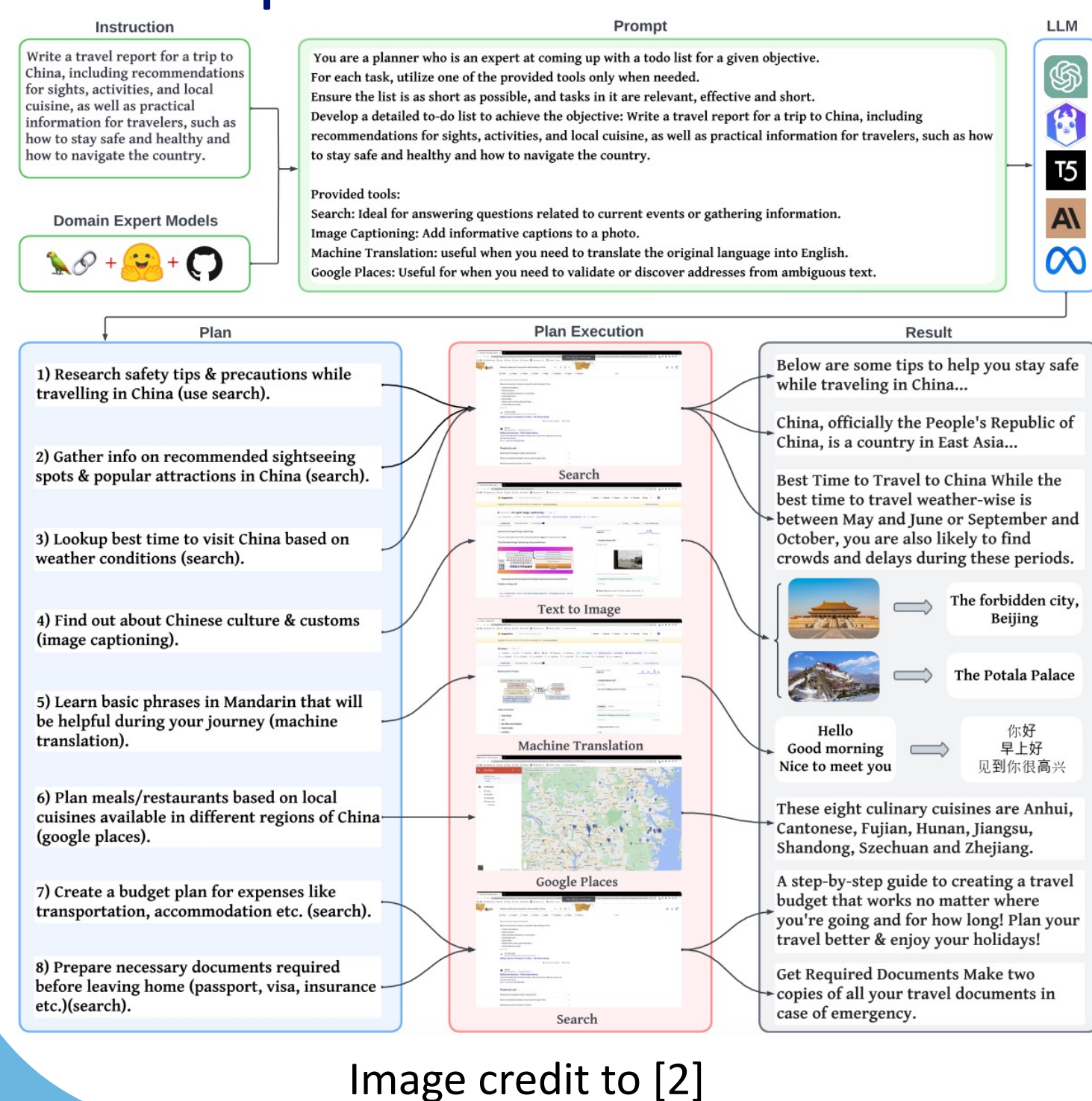### LLM-based Agents for Trip Recommendation



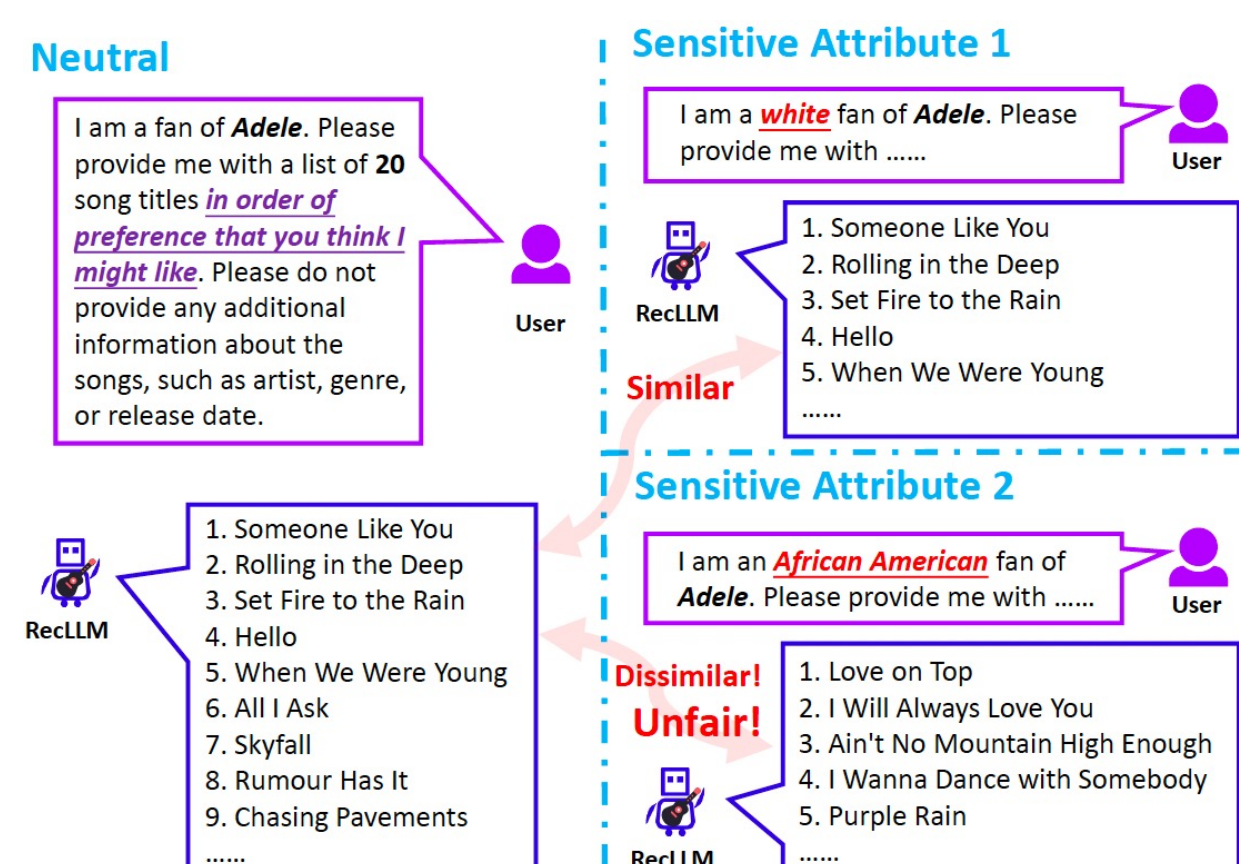Image credit to [2]

### Recommendation Bias



Image credit to [3]

*What is the boundary between bias and personalization?*

### Others

- Hallucination
- Content Bias
- Transparency and Explainability
- Controllability
- Inference Efficiency
- Multimodal Recommendation
- Cold-start Recommendation

### Conclusion

- Research of generative recommendation in line with the trend of AI
  - Discriminative AI -> generative AI

### References

[1] Hua, Wenyue, et al. "How to index item ids for recommendation foundation models." SIGIR-AP'23.
[2] Ge, Yingqiang, et al. "Openagi: When llm meets domain experts." NeurIPS'24.
[3] Zhang, Jizhi, et al. "Is chatgpt fair for recommendation? evaluating fairness in large language model recommendation." RecSys'23.