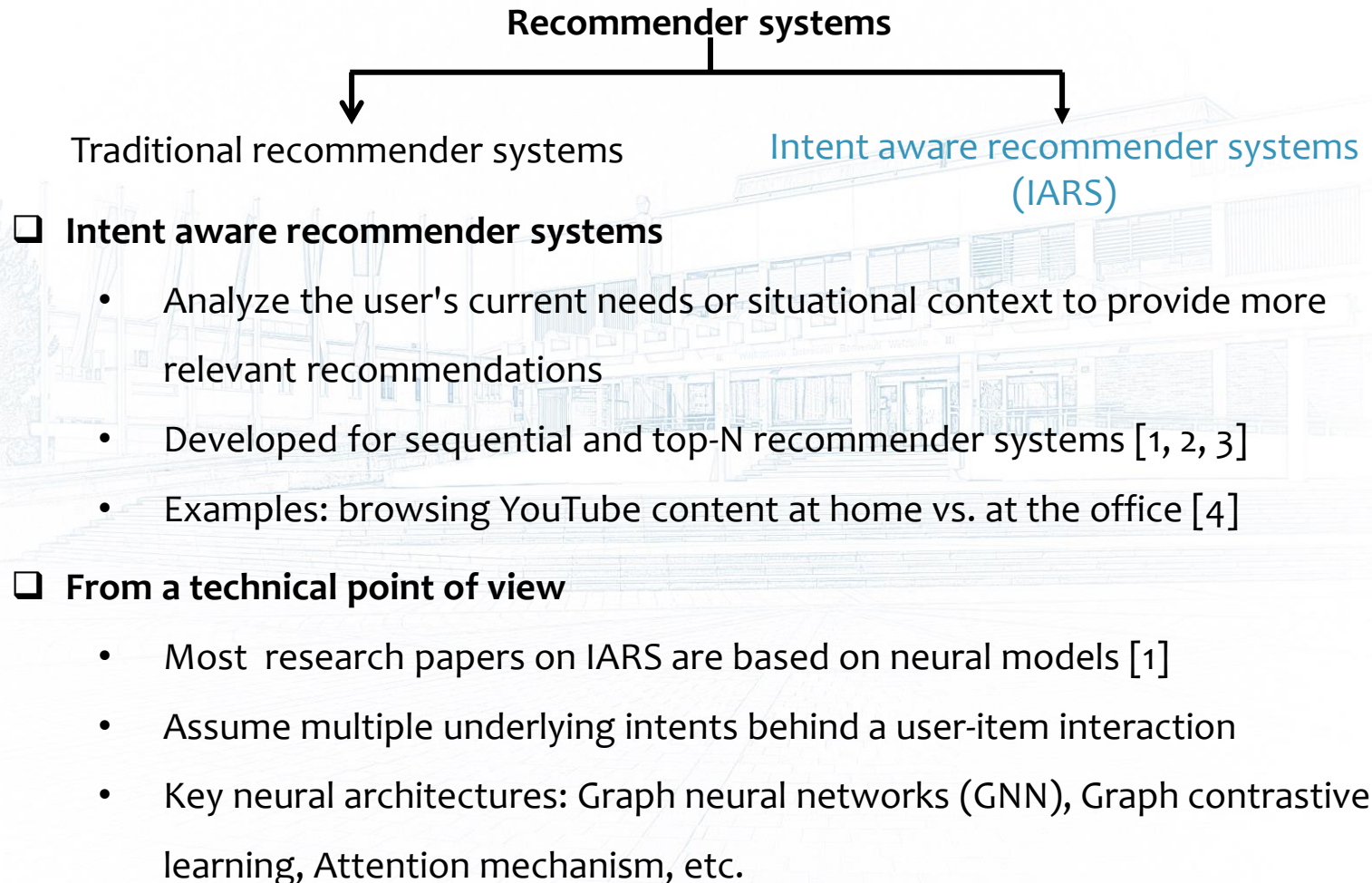# A Worrying Reproducibility Study of Intent-Aware Recommendation Models

**Faisal Shehzad**, Maurizio Ferrari Dacrema, Dietmar Jannach
Emails: Faisal.Shehzad@aau.at, Maurizio.Ferrari@polimi.it,  Dietmar.Jannach@aau.at

# Introduction and Motivation

**Recommender systems**

Traditional recommender systems　　　Intent aware recommender systems (IARS)

❑ **Intent aware recommender systems**

- Analyze the user's current needs or situational context to provide more relevant recommendations

- Developed for sequential and top-N recommender systems [1, 2, 3]

- Examples: browsing YouTube content at home vs. at the office [4]

❑ **From a technical point of view**

- Most research papers on IARS are based on neural models [1]

- Assume multiple underlying intents behind a user-item interaction

- Key neural architectures: Graph neural networks (GNN), Graph contrastive learning, Attention mechanism, etc.

# Reproducibility Crisis

❑ Due to the increased sophistication of complex models, one can imagine their superiority over simpler models

- However, the literature work presents a different picture

  - Studies [5, 6] demonstrate the superiority of simpler models over top-N recommenders like NeuMF and ConvMF

  - Another recent study [7] shows similar results for GNN-based session-based recommender systems

  - Another study [8] shows shared artifacts are insufficient to reproduce reported results in a paper

❑ Various factors [5, 7] contribute to this "virtual progress", such as comparing only one family of algorithms (neural models), poorly tuned baseline models, etc.

# Research Questions

❑ From the literature review, we observe that the performance of complex algorithms is quite limited

❑ However, we want to know whether this pattern is still persistent in the emerging and promising area of IARS. Therefore, we have formulated two research questions:

- Can the reported results of published papers be reproduced using the provided artifacts, such as source code and other relevant information?

- How could the simpler models be compared with complex models?

# Research Methodology

❑ **Identification of IARS papers**

1. Queried Google Scholar and IEEE Xplore for papers containing the terms 'intent' or 'intent awareness' along with the keywords 'recommend'

2. Identified 88 papers and retained those proposing a top-N model, published in the last four years in A* venues or top journals

3. Thirteen papers met this criteria

❑ **Baselines:** Chose six baseline models, which showed good performance in the previous reproducibility studies [6, 7, 8]

- ItemKNN

- $RP^3\beta$

- and others

❑ **Evaluation methodology:** Simpler models are compared with IARS under the same configuration used in the published papers [6]

# Results (1/3)

❑ RQ1: Can the reported results of published papers be reproduced using the provided artifacts, such as source code and other relevant information?

❑ Here are the findings our study

- Practice of sharing reproducibility packages:  7  out of 13 (~53%)

- Reproducibility packages in working condition: 5 out of 13 (~38%)

- Reproducibility packages with reproducible results: 3 out of 13 (~23%)

# List of papers with working reproducibility packages

| SR.NO. | Title | Venue | Is it reproducible |
|--------|-------|-------|--------------------|
| 1 | Disentangled Graph Collaborative Filtering (DGCF) [9] | SIGIR '20 | ✓ |
| 2 | Learning Intents behind Interactions with Knowledge Graph for Recommendation (KGIN) [10] | WWW '20 | ✗ |
| 3 | Intent Disentanglement and Feature Self-supervision for Novel Recommendation (IDS4NR) [11] | IEEE TKDE '22 | ✗ |
| 4 | Disentangled Contrastive Collaborative Filtering (DCCF) [12] | SIGIR '23 | ✓ |
| 5 | Bilateral Intent-guided Graph Collaborative Filtering (BIGCF) [13] | SIGIR '24 | ✓ |

# Results (2/3)

❑ RQ2: How could the simpler models be compared with complex models?

❑ Here are the key findings of the experiments

- On all accuracy measures, the simpler models outperform all IARS

- In one case, the KGIN model outperforms the simpler models on Recall@20 for the Last.fm dataset. However, in this case, we observe a severe data leakage issue in the shared train-test splits

- For the DCCF and BIGCF models, the authors used an unusual recommendation list length, such as 40, without providing a valid reason

- We found no single winner among the simpler models. However, ItemKNN and RP$^3$β demonstrate competitive performance
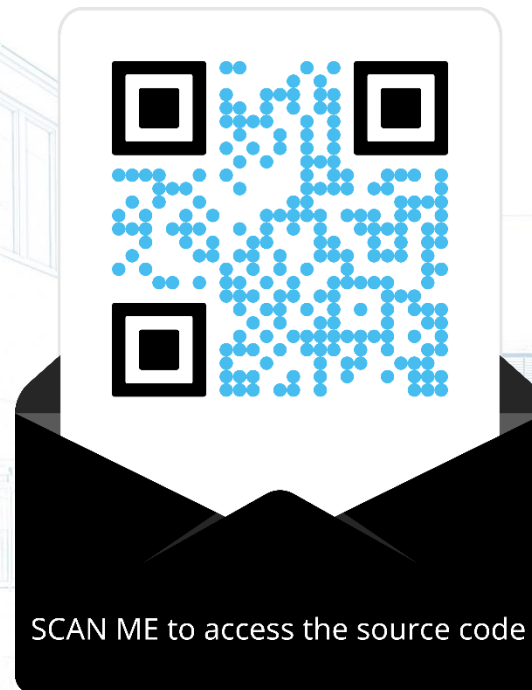
# Results (3/3)

❑ What is about computational complexity?

❑ We measure the training time (T-time) and prediction time (P-time) of the KGIN and the simpler models for the largest Alibaba iFashion dataset

- The Alibaba iFashion dataset is significantly smaller than the Netflix Prize dataset (100M), which was released 15 years ago
- KGIN takes 1 day and 8 hours per iteration on **GPU**, with several iterations required for hyperparameter tuning
- ItemKNN, which performs best on this dataset, requires only 2 minutes on **CPU** to build the lookup tables
- In terms of P-time, KGIN model is also 50% slower than ItemkNN

| Alibaba iFashion dataset | |
|---|---|
| Users | 114, 000 |
| Items | 30,000 |
| Interactions | 1.7M |

# Conclusions

❑ Reproducibility challenges in recommender systems keep returning, even after fifteen years

- While one might assume that the reproducibility should be high, particularly in top-rank venues, but our findings suggest otherwise

- None of examined papers share the code of baselines, data preprocessing, etc., to ensure full reproducibility

- Sometimes, authors are unresponsive when approached for guidance

- Lastly, we are not against research on complex algorithms. Many studies have shown their effectiveness in computer vision and NLP , however:
  - The current progress in IARS is limited
  - Need to avoid reoccurring methodological issues, such as data leakage
  - Encourage researchers to share their source code for proposed and baseline models to ensure full reproducibility

# Thank you for your attention

SCAN ME to access the source code

# References

1. Jannach, Dietmar, and Markus Zanker. "A Survey on Intent-aware Recommender Systems." ACM Transactions on Recommender Systems 3.2 (2024): 1-32.
2. Wang, Xiang, et al. "Disentangled Graph Collaborative Filtering." Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2020.
3. Jin, Di, et al. "Dual Intent Enhanced Graph Neural Network for Session-based New Item Recommendation." Proceedings of the ACM Web Conference 2023. 2023.
4. Loepp, Benedikt. "Multi-list Interfaces for Recommender Systems: Survey and Future Directions." Frontiers in Big Data 6 (2023): 1239705.
5. Anelli, Vito Walter, et al. "Top-N Recommendation Algorithms: A Quest for the State-of-the-art." Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization. 2022.
6. Dacrema, Maurizio Ferrari, Michael Benigni, and Nicola Ferro. "Reproducibility and Artifact Consistency of the SIGIR 2022 Recommender Systems Papers based on Message Passing." arXiv preprint arXiv:2503.07823 (2025).
7. Shehzad, Faisal, and Dietmar Jannach. "Performance Comparison of Session-based Recommendation Algorithms based on GNNs." European Conference on Information Retrieval. Cham: Springer Nature Switzerland, 2024.
8. Shehzad, Faisal, and Dietmar Jannach. "Everyone's a Winner! on Hyperparameter Tuning of Recommendation Models." Proceedings of the 17th ACM Conference on Recommender Systems. 2023.
9. Wang, Xiang, et al. "Disentangled Graph Collaborative Filtering." Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2020.
10. Wang, Xiang, et al. "Learning Intents behind Interactions with Knowledge Graph for Recommendation." Proceedings of the Web Conference 2021.
11. Qian, Tieyun, et al. "Intent Disentanglement and Feature Self-supervision for Novel Recommendation." IEEE Transactions on Knowledge and Data Engineering 35.10 (2022): 9864-9877.
12. Ren, Xubin, et al. "Disentangled Contrastive Collaborative Filtering." Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval. 2023.
13. Zhang, Yi, Lei Sang, and Yiwen Zhang. "Exploring the Individuality and Collectivity of Intents behind Interactions for Graph Collaborative Filtering." Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. 2024.