

ĐẠI HỌC QUỐC GIA HÀ NỘI  
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ

---



**NGUYỄN TRẦN NGỌC LINH**

**NGHIÊN CỨU MÔ HÌNH HOÁ HỌC SÂU DỮ LIỆU  
CHÍNH TẮC VÀ PHỤ TRỢ NHẪM NÂNG CAO TÍNH VỮNG  
CHẮC VÀ THÍCH NGHI CỦA HỆ THỐNG KHUYẾN NGHỊ**

**NGÀNH ĐÀO TẠO:    HỆ THỐNG THÔNG TIN**

**MÃ SỐ:            9480104**

**TÓM TẮT LUẬN ÁN TIẾN SĨ NGÀNH HỆ THỐNG THÔNG TIN**

**HÀ NỘI - 2026**

Công trình được hoàn thành tại: Trường Đại học Công nghệ, Đại học Quốc gia Hà Nội

Người hướng dẫn khoa học:

- PGS. TS Phan Xuân Hiếu, Trường ĐH Công nghệ, ĐH QGHN
- TS. Trần Mai Vũ, Trường ĐH Công nghệ, ĐH QGHN

Phản biện: TS. Đỗ Thanh Hà

Phản biện: PGS. TS. Bùi Thu Lâm

Phản biện: PGS. TS. Đỗ Văn Thành

Luận án sẽ được bảo vệ trước Hội đồng cấp Đại học Quốc gia chấm luận án tiến sĩ họp tại Trường Đại học Công nghệ, Đại học Quốc gia Hà Nội

vào hồi 8 giờ 30 phút thứ 5 ngày 18 tháng 06 năm 2026

**NGHIÊN CỨU SINH**

**CÁN BỘ HƯỚNG DẪN**

**XÁC NHẬN CỦA ĐƠN VỊ ĐÀO TẠO**

Có thể tìm hiểu luận án tại:

- Thư viện Quốc gia Việt Nam
- Trung tâm Thông tin - Thư viện, Đại học Quốc gia Hà Nội

# Tóm tắt

Các hệ thống khuyến nghị hiện đại ngày càng phải vận hành trong những hệ sinh thái số quy mô lớn, nơi lượng người dùng tăng trưởng liên tục, không gian sản phẩm mở rộng nhanh chóng, và các mẫu hình tương tác biến đổi theo những cách khó lường.

Luận án này phát triển các phương pháp học sâu cho khuyến nghị bền vững và thích nghi thông qua việc mô hình hóa sâu cả dữ liệu chính tắc (các tương tác người dùng-sản phẩm) lẫn dữ liệu phụ trợ (thông tin bên lề, bối cảnh miền, và tín hiệu hội thoại). Nghiên cứu giải quyết ba thách thức nền tảng trên bốn hướng nghiên cứu có liên hệ chặt chẽ với nhau.

*Thứ nhất*, luận án nghiên cứu khuyến nghị có khả năng mở rộng thông qua biểu diễn người dùng không phụ thuộc định danh (phi ID người dùng), phân cụm mềm bằng mạng nơ-ron, và học tương phản giúp duy trì chất lượng khuyến nghị ở quy mô thực tế.

*Thứ hai*, nghiên cứu khám phá việc hợp nhất bền vững dữ liệu chính tắc và dữ liệu phụ trợ thông qua cơ chế sinh trọng số dựa trên cơ chế chú ý và học tương phản đồ thị có che mặt nạ giúp nâng cao độ bền vững của biểu diễn trong điều kiện thưa thớt và khởi động nguội.

*Thứ ba*, luận án phát triển các cơ chế học liên tục cho khuyến nghị đa miền thích nghi.

*Thứ tư*, nghiên cứu đề xuất khuyến nghị hội thoại lai, kết nối dữ liệu chính tắc và dữ liệu phụ trợ thông qua mô hình hóa sở thích dựa trên mạng nơ-ron đồ thị, tích hợp với các mô hình ngôn ngữ lớn và sinh tăng cường bằng truy hồi RAG.

Các thực nghiệm sâu rộng trên những bộ dữ liệu chuẩn cùng việc triển khai thực tế với các sản phẩm thật của Viettel (một trong những tập đoàn lớn nhất Việt Nam) đã kiểm chứng tính hiệu quả của các mô hình đề xuất, cho thấy những cải thiện nhất quán so với các phương pháp tiên tiến.

**Từ khóa:** Học sâu, hệ thống khuyến nghị, dữ liệu chính tắc, dữ liệu phụ trợ, thưa thớt dữ liệu, vấn đề khởi động nguội, khả năng mở rộng, phân cụm mềm, học tương phản, mạng nơ-ron đồ thị, hợp nhất thông tin phụ trợ, học liên tục, khuyến nghị đa miền, khuyến nghị hội thoại, mô hình ngôn ngữ lớn, sinh tăng cường bằng truy hồi.

# Chương 1

## Tổng quan các nghiên cứu liên quan về nền tảng và phương pháp

### 1.1 Định nghĩa và phát biểu bài toán

#### 1.1.1 Tổng quan bài toán khuyến nghị

##### Phát biểu bài toán

Hệ thống khuyến nghị hoạt động trên ba thực thể cơ bản: một tập người dùng  $\mathcal{U} = \{u_1, u_2, \dots, u_M\}$ , một tập sản phẩm  $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$ , và một ma trận tương tác  $R \in \mathbb{R}^{M \times N}$  giữa chúng. Cốt lõi của mọi phương pháp khuyến nghị là một hàm mục tiêu định lượng mức độ liên quan của một sản phẩm đối với một người dùng:

$$f : \mathcal{U} \times \mathcal{V} \rightarrow \mathcal{D} \quad (1.1)$$

trong đó  $\mathcal{D}$  biểu diễn miền giá trị của hàm mục tiêu.

Mục tiêu chính của một hệ thống khuyến nghị là, với mỗi người dùng  $u \in \mathcal{U}$ , xác định những sản phẩm cực đại hóa hàm mục tiêu:

$$v_u^* = \arg \max_{v \in \mathcal{V}} f(u, v) \quad (1.2)$$

Trên thực tế, các hệ thống thường sinh ra một danh sách xếp hạng gồm top- $K$  sản phẩm thay vì một gợi ý đơn lẻ, tạo thành một tập có thứ tự  $\mathcal{D}_u = \{v_1, v_2, \dots, v_K\}$ , trong đó các sản phẩm được sắp xếp theo thứ tự giảm dần của điểm hữu ích dự đoán. Khi đó, bài toán khuyến nghị trở thành:

$$\mathcal{D}_u = \text{Top-}K_{v \in \mathcal{V}} (f(u, v)) \quad (1.3)$$

## Biểu diễn lịch sử tương tác

Lịch sử tương tác người dùng-sản phẩm thường được biểu diễn bằng một ma trận đánh giá  $\mathbf{R} \in \mathbb{R}^{M \times N}$ , trong đó  $M = |\mathcal{U}|$  là số người dùng và  $N = |\mathcal{V}|$  là số sản phẩm. Mỗi phần tử  $r_{u,v}$  trong ma trận này biểu diễn tương tác quan sát được giữa người dùng  $u$  và sản phẩm  $v$ .

Đối với hệ thống khuyến nghị, một phản hồi tường minh,  $r_{u,v}$  ghi nhận đánh giá của người dùng phản ánh trực tiếp sở thích của họ, chẳng hạn một xếp hạng 5 sao cho một bộ phim hay một điểm số cho một sản phẩm. Đối với hệ thống khuyến nghị, phản hồi ngầm định,  $r_{u,v}$  mã hóa các tín hiệu hành vi như nhấp chuột, lướt xem hay lướt mua, vốn gián tiếp thể hiện sự quan tâm của người dùng mà không cần phát biểu sở thích một cách tường minh.

Một đặc điểm then chốt của ma trận đánh giá là chỉ một phần nhỏ các phần tử được quan sát. Gọi  $\Omega \subseteq \mathcal{U} \times \mathcal{V}$  là tập các cặp người dùng-sản phẩm được quan sát. Khi đó, bài toán khuyến nghị có thể được phát biểu như bài toán hoàn thiện ma trận: cho ma trận quan sát một phần  $\mathbf{R}_\Omega$ , hãy dự đoán các phần tử còn thiếu  $\mathbf{R}_{\bar{\Omega}}$ , trong đó  $\bar{\Omega} = (\mathcal{U} \times \mathcal{V}) \setminus \Omega$  biểu diễn các tương tác chưa được quan sát.

$$\hat{\mathbf{R}} = \mathcal{F}(\mathbf{R}_\Omega; \Theta) \quad (1.4)$$

trong đó  $\mathcal{F}$  biểu diễn mô hình khuyến nghị được tham số hóa bởi  $\Theta$ , và  $\hat{\mathbf{R}}$  là ma trận đánh giá hoàn chỉnh được dự đoán.

Mật độ của ma trận đánh giá, định nghĩa là  $|\Omega|/(M \times N)$ , thường cực kỳ thấp trong các hệ thống thực tế, thường dưới 1% đối với các nền tảng quy mô lớn. Sự thưa thớt rất lớn này định hình một cách căn bản việc thiết kế các thuật toán khuyến nghị và là động lực cho nhiều kỹ thuật tiên tiến được bàn luận trong luận án này.

### 1.1.2 Phạm vi và mục tiêu nghiên cứu

Luận án này phát triển các giải pháp dựa trên học sâu nhằm giải quyết các thách thức có liên hệ với nhau gồm thưa thớt dữ liệu, khởi động nguội và khả năng mở rộng. Các cách tiếp cận được đề xuất chia sẻ một mục tiêu chung: học các biểu diễn bền vững, có khả năng chuyển giao, nắm bắt được sở thích thực sự của người dùng mà không đòi hỏi lịch sử tương tác đầy đủ hay tài nguyên tính toán tốn kém.

# Chương 2

## Khuyến nghị bền vững thông qua nhúng tương tác và phân cụm mềm

### 2.1 Giới thiệu

Các hệ khuyến nghị hiện đại đã trở thành thành phần thiết yếu của những ứng dụng quy mô lớn, phục vụ hàng triệu người dùng trên các nền tảng thương mại điện tử, dịch vụ phát trực tuyến và mạng xã hội. Bất chấp những tiến bộ đáng kể của các phương pháp khuyến nghị dựa trên học sâu, hiện vẫn còn tồn tại một số thách thức nền tảng làm hạn chế khả năng triển khai thực tế và khả năng mở rộng của các hệ thống này. Phần này nhận diện những thách thức then chốt về khả năng mở rộng trong các hệ khuyến nghị hiện đại và lấy chúng làm động lực chính cho các đề xuất trong chương này.

#### 2.1.1 Thách thức về khả năng mở rộng trong khuyến nghị

Các hệ khuyến nghị hiện đại dựa vào việc học một véc-tơ nhúng riêng cho mỗi người dùng và mỗi sản phẩm. Tuy hiệu quả trên các bộ dữ liệu chuẩn học thuật, cách tiếp cận này gặp những hạn chế nghiêm trọng khi triển khai trong công nghiệp.

**Ràng buộc bộ nhớ:** Dung lượng bộ nhớ tăng tuyến tính theo số người dùng. Chen và cộng sự đã tính toán rằng để phục vụ 1 tỷ người dùng với các véc-tơ nhúng 64 chiều cần khoảng 238 GB chỉ riêng cho phần nhúng người dùng. Khi cộng thêm phần nhúng sản phẩm và các đặc trưng phân loại, tổng nhu cầu có thể vượt quá năng lực của phần cứng phổ thông, buộc phải đánh đổi giữa khả năng biểu đạt của mô hình và tài nguyên tính toán.

**Vấn đề khởi động nguội:** Những người dùng và sản phẩm có lịch sử tương tác

không đủ sẽ khó học được các biểu diễn nhưng có khả năng khái quát. Người dùng mới không thể được biểu diễn một cách hiệu quả, dẫn đến chất lượng khuyến nghị ban đầu kém. Sản phẩm mới thiếu tín hiệu phản hồi để tối ưu phần nhưng, tạo ra vấn đề về khả năng hiển thị, khiến nội dung phù hợp vẫn bị ẩn đối với những người dùng tiềm năng.

**Chi phí tính toán:** Mô thức chấm điểm toàn bộ sản phẩm tạo ra sự lệch pha giữa huấn luyện ngoại tuyến và phục vụ trực tuyến. Các hệ thống vận hành phải phản hồi trong những ràng buộc độ trễ chặt chẽ (hàng chục mili-giây), trong khi phần lớn nghiên cứu lại tối ưu cho mục tiêu tái dựng toàn bộ ma trận, vốn bất khả thi về mặt tính toán trong các kịch bản thời gian thực.

## 2.2 EfficientRec: Khuyến nghị không định danh, khả mở rộng thông qua phân cụm mềm và học tương phản

Để giải quyết các khoảng trống nghiên cứu hiện tại, chương này luận án đề xuất EfficientRec, một mô hình mới cho khuyến nghị bền vững ở quy mô lớn thông qua nhưng tương tác và phân cụm mềm. Các mục tiêu nghiên cứu dẫn dắt quá trình phát triển EfficientRec như sau.

Phần này trình bày kiến trúc EfficientRec được đề xuất, một khung khuyến nghị có khả năng mở rộng được thiết kế nhằm khắc phục những hạn chế nền tảng của các mô hình khuyến nghị dựa trên định danh người dùng thông thường. Kiến trúc loại bỏ sự phụ thuộc vào định danh người dùng tường minh bằng cách xây dựng các biểu diễn người dùng một cách động từ các tín hiệu hành vi, qua đó đạt được độ phức tạp tính toán độc lập với quy mô tập người dùng. Thiết kế này cho phép hệ thống mở rộng tới những tập người dùng lớn trong khi vẫn duy trì chất lượng khuyến nghị nhất quán và hỗ trợ tích hợp liền mạch những người dùng mới mà không cần huấn luyện lại mô hình.

Mô hình đề xuất gồm ba thành phần chính phối hợp với nhau để đưa ra các khuyến nghị cá nhân hóa.

Kiến trúc tổng thể được minh họa trong Hình 2.1, cung cấp một cái nhìn ở mức tổng quan về cách ba thành phần tương tác với nhau để tạo ra các khuyến nghị cá nhân hóa.

Thành phần thứ nhất là “Mô hình Nhưng tương tác”, chịu trách nhiệm xây dựng các biểu diễn người dùng cô đọng và giàu thông tin bằng cách tổng hợp thông tin từ lịch sử tương tác của người dùng với các sản phẩm trong hệ thống.

Thành phần thứ hai là “Mô hình Phân cụm”, tổ chức người dùng thành các nhóm

nhận biết sở thích bằng kỹ thuật học tương phản và phân cụm mềm. Mô hình phân cụm học cách ánh xạ các biểu diễn người dùng vào một không gian sở thích tiềm ẩn, trong đó mỗi chiều tương ứng với một cụm sở thích riêng biệt. Hàm mục tiêu học tương phản bảo đảm rằng những người dùng có sở thích tương tự được ánh xạ tới những vùng tương tự trong không gian sở thích, trong khi những người dùng có sở thích khác nhau được tách biệt rõ ràng.

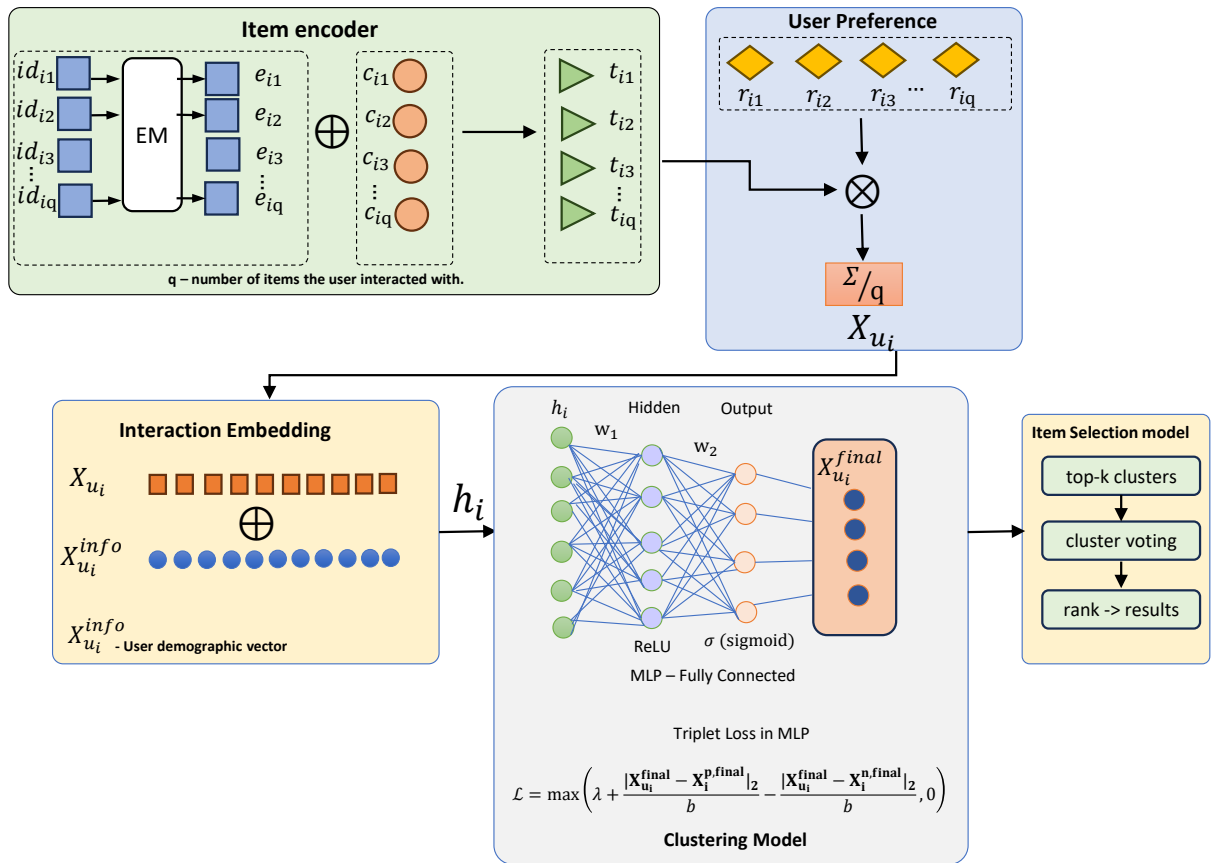


Figure 2.1: Kiến trúc tổng thể của EfficientRec

Thành phần thứ ba là “Mô hình Chọn sản phẩm”, tạo ra các khuyến nghị cá nhân hóa một cách hiệu quả thông qua cơ chế bỏ phiếu dựa trên cụm gồm hai pha. Trong pha ngoại tuyến, mô hình tính trước các điểm sở thích cho từng cặp cụm sản phẩm dựa trên các đánh giá được tổng hợp từ tất cả người dùng thuộc cụm đó. Trong pha trực tuyến, mô hình tính mức độ thuộc cụm của người dùng mục tiêu và tạo khuyến nghị bằng cách tổng hợp các điểm đã tính trước từ những cụm liên quan nhất của người dùng. Thiết kế hai pha này giảm đáng kể chi phí tính toán khi tạo khuyến nghị so với các phương pháp phải chấm điểm tất cả sản phẩm cho mỗi người dùng.

## 2.3 Kết quả thực nghiệm

Table 2.1: So sánh hiệu năng tổng thể (toàn bộ người dùng) @30, báo cáo dưới dạng trung bình  $\pm$  độ lệch chuẩn qua 5 lần sinh ngẫu nhiên.

Mô hình	Recall@30	NDCG@30	Nhóm
<i>Phương pháp đề xuất</i>			
EfficientRec	<b>0.1994 <math>\pm</math> 0.0028<sup>†</sup></b>	<b>0.1178 <math>\pm</math> 0.0027<sup>ns</sup></b>	Đề xuất
<i>Phương pháp dựa trên đồ thị</i>			
NGCF	0.1958 $\pm$ 0.0013	0.1174 $\pm$ 0.0007	Đồ thị
LINKX	0.1928 $\pm$ 0.0022	0.1160 $\pm$ 0.0012	Đồ thị
GraphSAGE	0.1658 $\pm$ 0.0006	0.0970 $\pm$ 0.0004	Đồ thị
GAT	0.1639 $\pm$ 0.0001	0.0946 $\pm$ 0.0001	Đồ thị
LightGCN	0.1636 $\pm$ 0.0003	0.0945 $\pm$ 0.0001	Đồ thị
<i>Phương pháp học tương phản</i>			
SSL4Rec	0.1644 $\pm$ 0.0001	0.0950 $\pm$ 0.0012	Tự giám sát
MixGCF	0.1591 $\pm$ 0.0001	0.0910 $\pm$ 0.0001	Tương phản
SGL	0.0192 $\pm$ 0.0015	0.0565 $\pm$ 0.0011	Tương phản
<i>Phương pháp phân cụm và phân rã ma trận</i>			
SCoC	0.1584 $\pm$ 0.0001	0.1107 $\pm$ 0.0004	Phân cụm
FCCF	0.1090 $\pm$ 0.0001	0.1072 $\pm$ 0.0012	Phân cụm
BPR-MF	0.1074 $\pm$ 0.0003	0.0912 $\pm$ 0.0001	Phân rã ma trận
FCM-Rec	0.1067 $\pm$ 0.0001	0.1068 $\pm$ 0.0001	Phân cụm

Số **in đậm** là kết quả tốt nhất và số gạch chân là Tốt thứ hai.

Table 2.2: Kết quả các thực nghiệm trực tuyến trên TV360

Phương pháp	Phim TV360		Video TV360	
	ACPU	ADPU	ACPU	ADPU
2DNNs	0.0152	13.896	0.0322	10.526
ALS	0.0121	12.190	0.0160	4.277
<b>ER (chia theo tương tác)</b>	<b>0.0186</b>	<b>15.018</b>	<b>0.0421</b>	<u>12.290</u>
ER (chia theo nhóm người dùng)	<u>0.0176</u>	<u>14.272</u>	<u>0.0381</u>	<b>13.155</b>

Số **in đậm** = Tốt nhất, số gạch chân = Tốt thứ hai.

## 2.4 Tóm tắt chương

Chương này giới thiệu EfficientRec, một khung khuyến nghị khả mở rộng, không định danh, thay thế phần nhúng riêng cho từng người dùng bằng các biểu diễn hướng hành vi được học từ các tập con tương tác thông qua mô hình hóa tương tác sâu, phân

cụm mềm và học tương phản. Bằng cách tính mỗi véc-tơ người dùng theo nhu cầu thay vì lưu trữ nó, thiết kế này giảm dung lượng bộ nhớ phía người dùng từ  $O(M \cdot d)$  xuống  $O(K \cdot d)$ , mở rộng dưới tuyến tính theo quy mô tập người dùng, và mở rộng một cách tự nhiên tới những người dùng mới và ít tương tác mà không cần huấn luyện lại. Các thực nghiệm ngoại tuyến trên những bộ dữ liệu chuẩn công khai cùng với kiểm thử A/B trực tuyến trên một nền tảng phát trực tuyến vận hành xác nhận rằng EfficientRec ngang bằng hoặc vượt các đường cơ sở mạnh dựa trên đồ thị trong khi giảm đáng kể chi phí tài nguyên.

Những quan sát này gợi mở một số hướng mở rộng tự nhiên. Việc học số lượng cụm, hoặc cập nhật các cụm một cách tăng dần khi những tương tác mới liên tục đổ về, sẽ giúp mô hình bám sát hơn một tập người dùng đang biến động, trong khi việc tích hợp thông tin phụ trợ đa phương thức phong phú hơn sẽ mang lại một tiên nghiệm mạnh hơn ở đúng nơi mà bằng chứng hành vi mỏng nhất. Cuối cùng, việc nén mô hình thông qua lượng tử hóa, cắt tỉa và chưng cất tri thức sẽ đưa quá trình suy luận vào trong những hạn mức về độ trễ và bộ nhớ cần thiết cho triển khai trên thiết bị.

# Chương 3

## Nâng cao khuyến nghị thông qua hợp nhất dữ liệu chính tắc và thông tin phụ trợ dựa trên mô hình đồ thị

### 3.1 Giới thiệu

Chương này trình bày cách kết hợp một cách gắn kết giữa thông tin phụ trợ, các mẫu tương tác và học trên đồ thị nhằm nâng cao hiệu năng của khuyến nghị hiện đại. Thông qua việc mô hình hóa kiến trúc chi tiết và phân tích thực nghiệm, chương này cho thấy rằng việc bổ sung hợp nhất thông tin phụ trợ có cấu trúc vào khuyến nghị dựa trên mạng nơ-ron đồ thị tạo ra những kết quả khuyến nghị chính xác hơn, bền vững hơn và nhất quán hơn về mặt ngữ nghĩa. Các công trình này đã được công bố tại các hội nghị bình duyệt, gồm: [P2] “GIFT4Rec: An Effective Side Information Fusion Technique Apply to Graph Neural Network for Cold-Start Recommendation” (ACIIDS 2023), và [P3] “The Masked Simple Graph Contrastive Learning for Recommendation” (KSE 2024).

## **3.2 GIFT4Rec: Hợp nhất thông tin phụ trợ bằng các kỹ thuật dựa trên cơ chế chú ý và học siêu cấp cho khuyến nghị khởi động nguội**

### **3.2.1 GIFT4Rec: Kiến trúc và các thành phần của mô hình**

Phần này trình bày kiến trúc GIFT4Rec được đề xuất, một khung hợp nhất để tích hợp thông tin phụ trợ vào các hệ khuyến nghị dựa trên đồ thị. Kiến trúc giải quyết thách thức nền tảng là cân bằng giữa tín hiệu hành vi từ các tương tác người dùng-sản phẩm với thông tin ngữ nghĩa từ các thuộc tính người dùng, qua đó cho phép khuyến nghị bền vững trong cả kịch bản khởi động âm lẫn khởi động nguội.

Kiến trúc tổng thể được minh họa trong Hình 3.1, cung cấp một cái nhìn ở mức tổng quan về cách ba thành phần phối hợp với nhau để tạo ra các khuyến nghị cá nhân hóa.

Thành phần thứ nhất là “Mô-đun Tương tác mạng nơ-ron đồ thị”, học các biểu diễn người dùng và sản phẩm bằng cách lan truyền thông tin qua đồ thị tương tác người dùng-sản phẩm. Khác với các cách tiếp cận dựa trên nội dung vốn chỉ dựa vào việc so khớp đặc trưng, thành phần này nắm bắt các tín hiệu cộng tác từ cấu trúc tương tác toàn cục, cho phép phát hiện những mẫu hình sở thích nảy sinh từ hành vi tập thể của người dùng.

Thành phần thứ hai là “Mô-đun Hợp nhất Thông tin Phụ trợ Cục bộ” (LSIF), kết hợp một cách thích nghi các biểu diễn nhúng hành vi với các biểu diễn nhúng thông tin phụ trợ cho từng người dùng riêng lẻ. Nhận xét then chốt là chiến lược hợp nhất tối ưu thay đổi tùy theo người dùng: một số người dùng có lịch sử tương tác phong phú, cung cấp tín hiệu sở thích mạnh, trong khi những người khác có ít tương tác, nơi thông tin phụ trợ trở nên giá trị hơn. Thành phần này học các trọng số hợp nhất cá nhân hóa thông qua một cơ chế dựa trên chú ý (ADN).

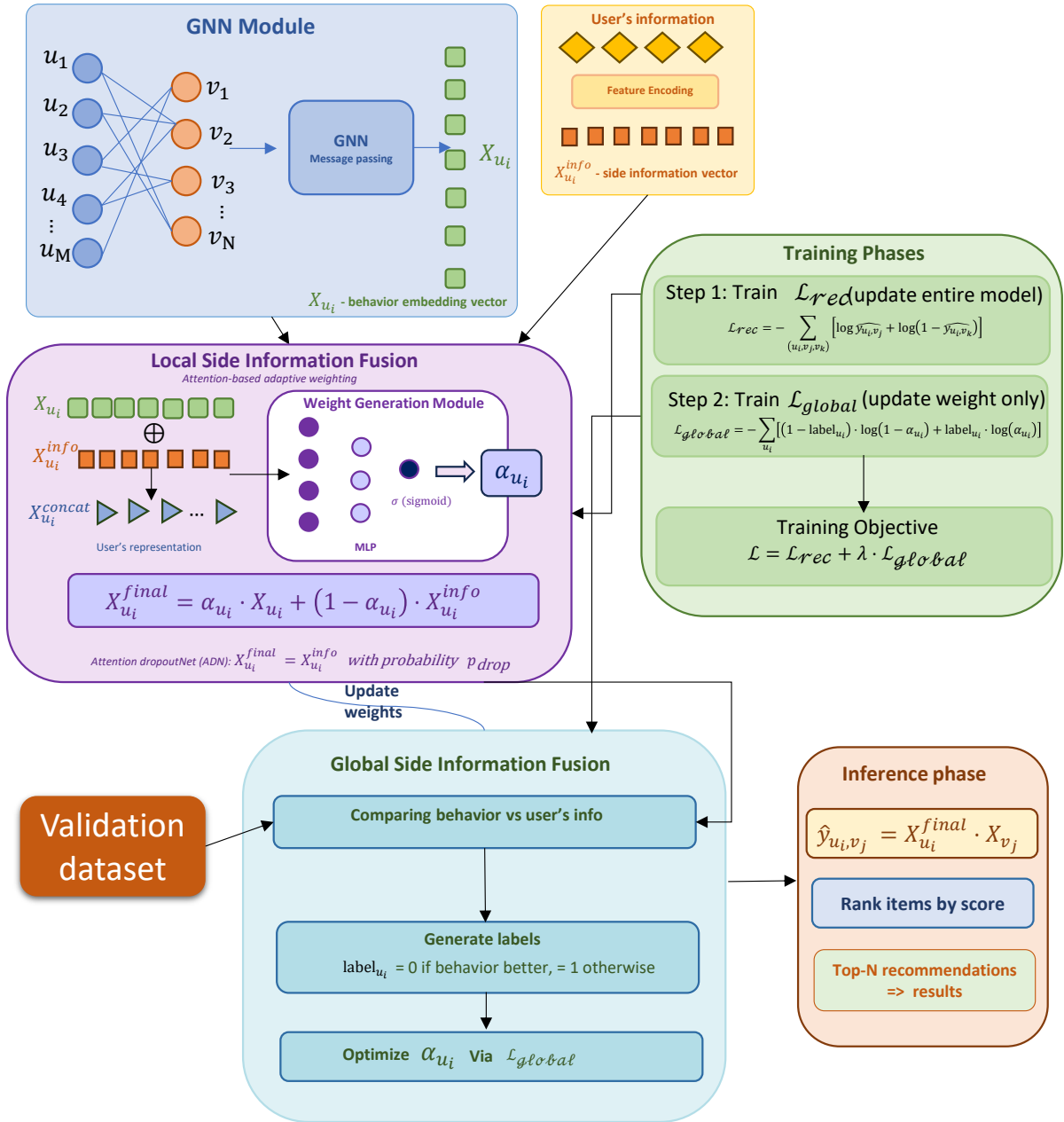


Figure 3.1: Kiến trúc tổng thể của GIFT4Rec

Thành phần thứ ba là mô-đun “Hợp nhất Thông tin Phụ trợ Toàn cục” (GSIF), cung cấp sự giám sát ở mức siêu cấp cho quá trình sinh trọng số bằng cách đánh giá nguồn thông tin nào (hành vi hay phụ trợ) hỗ trợ hiệu năng khuyến nghị tốt hơn trên dữ liệu kiểm định. Góc nhìn toàn cục này bảo đảm rằng các trọng số hợp nhất học được khái quát tốt vượt ra ngoài các tương tác huấn luyện.

## 3.2.2 Kết quả thực nghiệm

### a) Hiệu năng tổng thể

Table 3.1: GIFT4Rec - So sánh hiệu năng tổng thể (toàn bộ người dùng) @30, báo cáo dưới dạng trung bình  $\pm$  độ lệch chuẩn qua 5 lần sinh ngẫu nhiên.

Mô hình	Recall@30	NDCG@30	Nhóm
<b>GIFT4Rec</b>	<b><math>0.2162 \pm 0.0001^\dagger</math></b>	<b><math>0.1263 \pm 0.0001^\dagger</math></b>	Đề xuất
KGAT	$0.1793 \pm 0.0008$	$0.1067 \pm 0.0007$	Phụ trợ+Đồ thị tri thức
LINKX	$0.1928 \pm 0.0022$	$0.1160 \pm 0.0012$	Phụ trợ+Đồ thị
GAT	$0.1639 \pm 0.0001$	$0.0946 \pm 0.0001$	Phụ trợ+Chú ý
KGAT DropoutNet	$0.0908 \pm 0.0003$	$0.0359 \pm 0.0008$	Phụ trợ+Đồ thị tri thức
NGCF	$0.1958 \pm 0.0013$	<u><math>0.1178 \pm 0.0007</math></u>	Đồ thị
LightGCN	$0.1636 \pm 0.0003$	$0.0945 \pm 0.0001$	Đồ thị
SSL4Rec	$0.1644 \pm 0.0001$	$0.0950 \pm 0.0012$	Tự giám sát
EfficientRec (Ch.2)	<u><math>0.1994 \pm 0.0028</math></u>	$0.1174 \pm 0.0027$	Đề xuất

Số **in đậm** = tốt nhất, số gạch chân = tốt thứ nhì.

## 3.3 Học tương phản đồ thị đơn giản có che mặt nạ cho khuyến nghị

### 3.3.1 MaskSimGCL: Kiến trúc và các thành phần của mô hình

Phần này trình bày kiến trúc MaskSimGCL (học tương phản đồ thị đơn giản có che mặt nạ) được đề xuất, một khung mới được thiết kế nhằm khắc phục những hạn chế của các phương pháp học tương phản dựa trên đồ thị hiện có cho khuyến nghị. Kiến trúc mở rộng khung SimGCL bằng cách tích hợp các cơ chế che mặt nạ học được, vốn nhận diện và đánh trọng số một cách thích nghi tầm quan trọng của các chiều biểu diễn nhúng khác nhau, qua đó đạt được việc học biểu diễn bền vững hơn trong điều kiện dữ liệu thưa thớt.

Mô hình đề xuất gồm bốn thành phần chính phối hợp với nhau để đưa ra các khuyến nghị cá nhân hóa. Thành phần thứ nhất là bộ khung mạng nơ-ron đồ thị, chịu trách nhiệm học các biểu diễn người dùng và sản phẩm thông qua các phép truyền thông điệp trên đồ thị hai phía người dùng-sản phẩm. Theo thiết kế của LightGCN, thành phần này sử dụng các phép tích chập đồ thị đơn giản hóa, vốn lan truyền các tín hiệu cộng tác mà không biến đổi đặc trưng, nắm bắt các mẫu hình lân cận thông qua việc tổng hợp theo từng tầng.

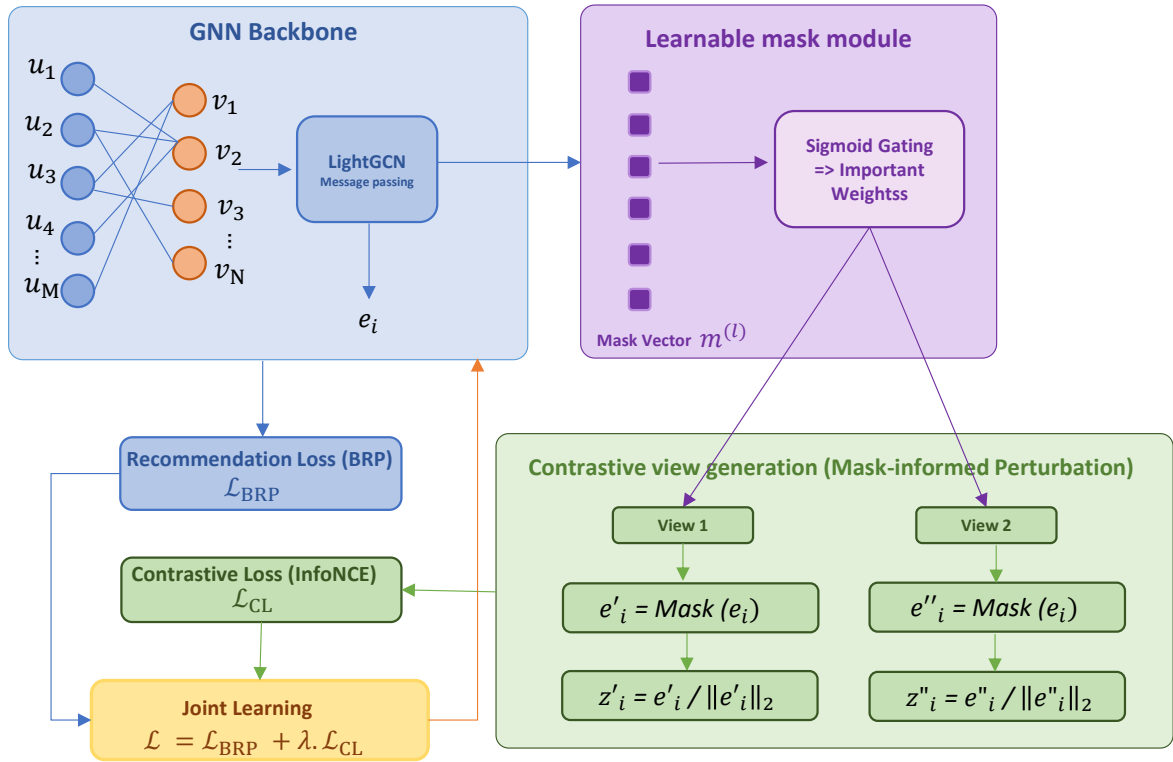


Figure 3.2: Kiến trúc tổng thể của MaskSimGCL

Thành phần thứ hai là mô-đun mặt nạ học được, đưa vào các vector mặt nạ huấn luyện được tại mỗi tầng của mạng nơ-ron đồ thị. Các mặt nạ này đóng vai trò là bộ lọc tầm quan trọng, đánh trọng số một cách thích nghi cho từng chiều của biểu diễn nhúng nút dựa trên mức độ liên quan của chúng đối với tác vụ khuyến nghị. Bằng cách tập trung năng lực mô hình vào các tham số giàu thông tin trong khi triệt tiêu những chiều nhiễu hoặc dư thừa, thành phần này cung cấp một sự điều chuẩn ngầm, giúp giảm thiểu hiệu quả hiện tượng quá khớp trong môi trường dữ liệu thưa thớt.

Thành phần thứ ba là mô-đun sinh khung nhìn tương phản, vốn xây dựng các biểu diễn được tăng cường cho học tương phản. Khác với cách tác động nhiễu đồng nhất của SimGCL, MaskSimGCL sử dụng các nhiễu loạn dựa trên mặt nạ, áp dụng những độ lớn nhiễu khác nhau dựa trên các điểm tầm quan trọng đã học. Những chiều được nhận diện là kém quan trọng sẽ nhận nhiễu loạn lớn hơn, trong khi những chiều then chốt được giữ với nhiễu nhỏ hơn, tạo ra những khung nhìn tương phản nhất quán giúp nâng cao việc học biểu diễn.

Thành phần thứ tư là khung tối ưu đồng thời, kết hợp mục tiêu khuyến nghị có giám sát với mục tiêu học tương phản tự giám sát. Công thức học đa nhiệm này cho phép mô hình đồng thời tối ưu cho việc dự đoán sở thích chính xác và học biểu diễn bền vững, trong đó hàm mất mát tương phản cung cấp sự điều chuẩn tính đồng đều, thúc đẩy các

biểu diễn nhúng phân bố đều hơn trong không gian biểu diễn.

### 3.3.2 Kết quả thực nghiệm

#### Hiệu năng tổng thể

Table 3.2: MaskSimGCL - So sánh hiệu năng tổng thể (toàn bộ người dùng) @30, báo cáo dưới dạng trung bình  $\pm$  độ lệch chuẩn qua 5 lần sinh ngẫu nhiên.

Mô hình	Recall@30	NDCG@30	Nhóm
<b>MaskSimGCL</b>	<b><math>0.2404 \pm 0.0008^\dagger</math></b>	<b><math>0.1322 \pm 0.0004^\dagger</math></b>	Đề xuất
XSimGCL	$0.2301 \pm 0.0004$	$0.1250 \pm 0.0002$	Tương phản đồ thị
SimGCL	$0.2292 \pm 0.0012$	$0.1249 \pm 0.0008$	Tương phản đồ thị
LightGCL	$0.2128 \pm 0.0003$	$0.1103 \pm 0.0001$	Tương phản đồ thị
DirectAU	$0.2110 \pm 0.0009$	$0.1226 \pm 0.0003$	Tương phản đồ thị
SGL	$0.0192 \pm 0.0015$	$0.0565 \pm 0.0011$	Tương phản đồ thị
SSL4Rec	$0.1644 \pm 0.0001$	$0.0950 \pm 0.0012$	Đồ thị
LightGCN	$0.1636 \pm 0.0003$	$0.0945 \pm 0.0001$	Đồ thị
GIFT4Rec (Mục 3.2)	$0.2162 \pm 0.0001$	$0.1263 \pm 0.0001$	Đề xuất
EfficientRec (Chương 2)	$0.1994 \pm 0.0028$	$0.1174 \pm 0.0027$	Đề xuất

Số **in đậm** = tốt nhất, số gạch chân = tốt thứ nhì.

## 3.4 Tóm tắt chương

Chương này giải quyết vấn đề thừa thớt dữ liệu và khởi động nguội thông qua hai đóng góp bổ trợ cho nhau. GIFT4Rec thực hiện hợp nhất thông tin phụ trợ một cách thích nghi, sử dụng một mô-đun sinh trọng số dựa trên chú ý để tính các trọng số hợp nhất riêng cho từng người dùng và kết hợp các biểu diễn nhúng hành vi với các tín hiệu phụ trợ thông qua các mô-đun cục bộ và toàn cục được tối ưu theo các nguyên lý học siêu cấp. MaskSimGCL bổ trợ cho nó bằng cách củng cố độ bền vững của biểu diễn nhờ học tương phản đồ thị có che mặt nạ. Cùng nhau, chúng mang lại những cải thiện nhất quán so với các đường cơ sở dựa trên đồ thị trong cả thiết lập khởi động âm lẫn khởi động nguội, cho thấy rằng việc hợp nhất có nguyên tắc và điều chuẩn tương phản giúp cải thiện khuyến nghị trong điều kiện khan hiếm dữ liệu.

# Chương 4

## Nâng cao khuyến nghị đa miền bằng học liên tục

### 4.1 Giới thiệu

Các hệ khuyến nghị hiện đại ngày càng vận hành trên nhiều miền dịch vụ trong cùng một nền tảng hợp nhất, nơi người dùng tương tác với những loại nội dung không đồng nhất như sản phẩm thương mại điện tử, phát video trực tuyến, dịch vụ âm nhạc và dòng tin tức. Những môi trường đa miền này đặt ra các thách thức riêng vượt ra ngoài mô thức khuyến nghị đơn miền truyền thống. Mặc dù khuyến nghị xuyên miền đã nổi lên như một hướng đi hứa hẹn nhằm tận dụng nguồn thông tin phong phú giữa các miền, các cách tiếp cận hiện có phần lớn tập trung cải thiện hiệu năng của miền đích trong khi thường bỏ qua việc bảo toàn tri thức của miền nguồn và sự công bằng về hiệu năng trên tất cả các miền tham gia.

#### 4.1.1 Kiến trúc và các thành phần của mô hình

Mô hình CNL4Rec giới thiệu một cơ chế học liên tục dựa trên che mặt nạ tác vụ, hoạt động ở mức nhúng. Nguyên lý thiết kế cốt lõi là coi mỗi miền như một tác vụ học tuần tự và áp dụng các mặt nạ riêng cho từng miền nhằm nhận diện và bảo vệ những chiều tiềm ẩn quan trọng đối với mỗi miền.

Kiến trúc giới thiệu các véc-tơ mặt nạ học được cho từng miền, có cùng số chiều với phần nhúng người dùng và sản phẩm. Các mặt nạ này nhận diện những chiều tiềm ẩn nào là thiết yếu để biểu diễn các mẫu hình hành vi riêng của miền, cho phép hệ thống chỉ cập nhật có chọn lọc những tham số liên quan trong quá trình huấn luyện. Những tham số bị xem là không quan trọng đối với miền hiện tại sẽ được bảo vệ để duy trì hiệu năng trên các miền đã học trước đó.

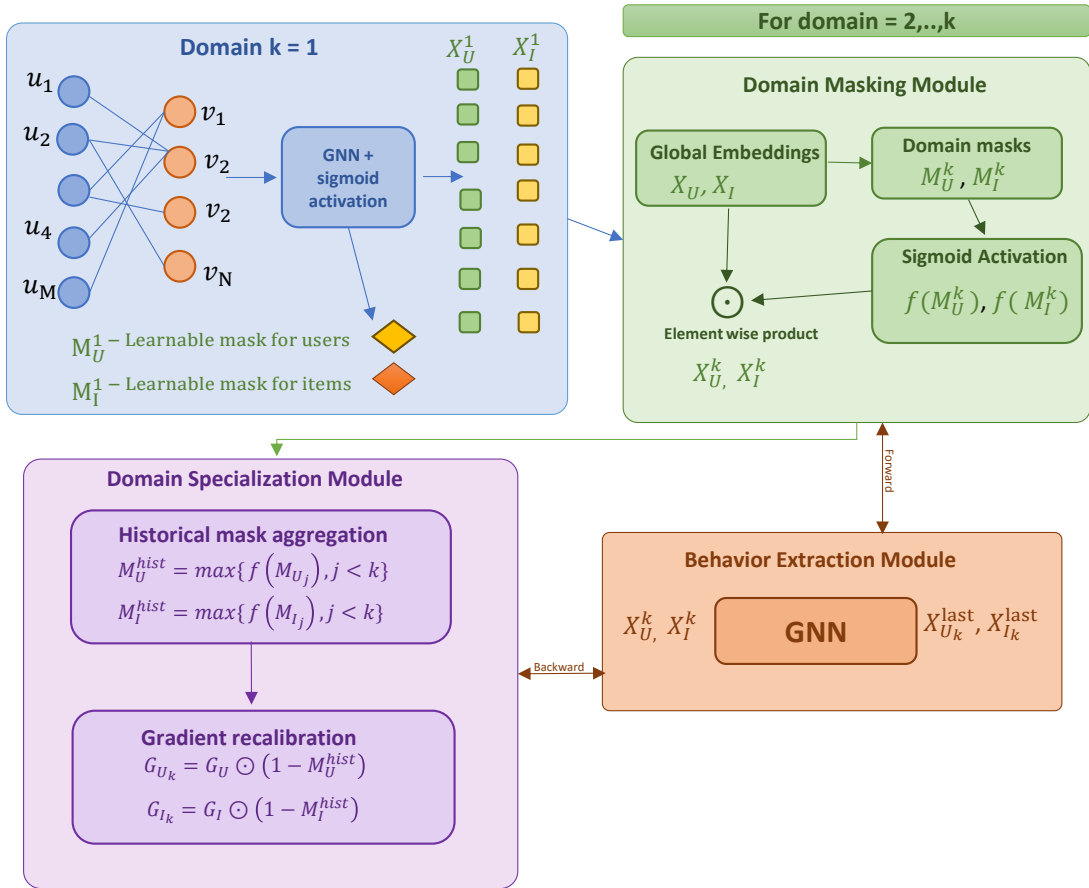


Figure 4.1: Kiến trúc tổng thể của CNL4Rec

Như minh họa trong Hình 4.1, CNL4Rec giải quyết bài toán khuyến nghị đa miền liên tục thông qua ba mô-đun tích hợp chặt chẽ. “Mô-đun Che mặt nạ miền” học các mặt nạ nhị phân riêng cho từng miền trên các ma trận nhúng người dùng và sản phẩm dùng chung, nhận diện những chiều tiềm ẩn nào liên quan đến mỗi miền và tạo ra các biểu diễn riêng của miền đã được lọc  $X_U^k$  và  $X_I^k$ . Vai trò của nó là phân hoạch không gian nhúng dùng chung thành các không gian con riêng của từng miền mà không cần tham số riêng cho mỗi miền. “Mô-đun Chuyên biệt hóa miền” bảo vệ tri thức đã thu nhận trước đó bằng cách tổng hợp các mặt nạ của tất cả các miền trong quá khứ thành một mặt nạ lịch sử  $M^{hist}$ , sau đó dùng nó để đưa về 0 các gradient chảy vào những chiều đã được các miền trước chiếm giữ trong quá trình lan truyền ngược. Việc triệt tiêu gradient có chọn lọc này chính là cơ chế chống quên cốt lõi, đảm bảo rằng việc học một miền mới không ghi đè lên những biểu diễn thiết yếu của các miền trước đó. “Mô-đun Trích xuất hành vi” nhận các biểu diễn nhúng riêng của miền đã được che mặt nạ và đưa chúng qua một bộ khung mạng nơ-ron đồ thị nhằm tổng hợp các tín hiệu cộng tác bậc cao trong miền hiện tại, tạo ra các biểu diễn nhúng cuối cùng dùng cho việc chấm điểm sở thích và xếp hạng khuyến nghị. Cùng nhau, ba mô-đun thực hiện việc cô lập tri thức miền ở

mức nhúng, bảo vệ tri thức ở mức gradient, và học sở thích ở mức suy luận, qua đó giải quyết một cách tổng thể bài toán đánh đổi giữa tính ổn định và tính dẻo vốn cố hữu trong khuyến nghị đa miền liên tục.

## 4.1.2 Kết quả thực nghiệm

### a. Kết quả trên MovieLens-1M

Table 4.1: So sánh hiệu năng trên MovieLens-1M (Recall@30). CV(%)=hệ số biến thiên giữa các thể loại (càng thấp càng cân bằng); Min=Recall@30 của miền yếu nhất (càng cao càng tốt). Số **in đậm**=Tốt nhất, số gạch chân=Tốt thứ hai.

Phương pháp	Hành động	Hài	KHVT	Giật gân	Chính kịch	Trung bình	Độ lệch	CV(%)	Min
<b>CNL4Rec</b>	<b>0.2481</b>	<b>0.2583</b>	<b>0.1865</b>	<b>0.2349</b>	<u>0.2163</u>	<b>0.2288</b>	0.0254	<b>11.1</b>	<b>0.1865</b>
CTNet	<u>0.1600</u>	<u>0.0838</u>	<u>0.1342</u>	<u>0.2034</u>	<b>0.2268</b>	<u>0.1616</u>	0.0506	31.3	<u>0.0838</u>
DIIT	0.1246	0.0632	0.1106	0.1181	0.1061	0.1045	0.0216	<u>20.7</u>	0.0632
KEEP	0.0598	0.0060	0.0841	0.1005	0.1796	0.0860	0.0567	65.9	0.0060
MF	0.0744	0.0060	0.0286	0.0485	0.2061	0.0727	0.0704	96.8	0.0060
ECAT	0.0063	0.0273	0.0173	0.0744	0.1665	0.0583	0.0588	100.8	0.0063

### b. Kết quả trên Yelp

Table 4.2: So sánh hiệu năng trên Yelp (Recall@30). CV(%)=hệ số biến thiên giữa các danh mục (càng thấp càng cân bằng); Min=Recall@30 của miền yếu nhất (càng cao càng tốt). Số **in đậm**=Tốt nhất, số gạch chân=Tốt thứ hai.

Phương pháp	Nhà hàng	Mua sắm	Ấm thực	Làm đẹp	Sức khỏe	Trung bình	Độ lệch	CV(%)	Min
<b>CNL4Rec</b>	<b>0.0194</b>	0.0150	<u>0.0184</u>	<u>0.0166</u>	<b>0.0182</b>	<b>0.0174</b>	0.0015	8.8	0.0150
ECAT	<u>0.0164</u>	<u>0.0172</u>	0.0162	<b>0.0172</b>	0.0154	<u>0.0165</u>	0.0007	<b>4.1</b>	<u>0.0154</u>
CTNet	0.0146	0.0152	<b>0.0192</b>	0.0148	<u>0.0179</u>	0.0164	0.0019	11.4	0.0146
KEEP	0.0158	<b>0.0182</b>	0.0156	0.0156	0.0164	0.0165	0.0010	<u>6.0</u>	<b>0.0156</b>
MF	0.0150	0.0158	0.0179	0.0162	0.0156	0.0163	0.0010	6.1	0.0150
DIIT	0.0138	0.0154	0.0150	0.0156	0.0178	0.0155	0.0013	8.4	0.0138

### c. Kết quả trên Amazon

Table 4.3: So sánh hiệu năng trên Amazon (Recall@30). CV(%)=hệ số biến thiên giữa các danh mục (càng thấp càng cân bằng); Min=Recall@30 của miền yếu nhất (càng cao càng tốt). Số **in đậm**=Tốt nhất, số gạch chân=Tốt thứ hai.

Phương pháp	Điện tử	Sách	Phim	Gia dụng	Thể thao	Trung bình	Độ lệch	CV(%)	Min
<b>CNL4Rec</b>	<u>0.0139</u>	0.0125	0.0135	<b>0.0133</b>	0.0126	<b>0.0132</b>	0.0005	<b>4.1</b>	<b>0.0125</b>
ECAT	<b>0.0147</b>	0.0130	0.0128	<u>0.0132</u>	0.0116	<u>0.0131</u>	0.0010	<u>7.6</u>	<u>0.0116</u>
MF	0.0129	0.0109	<b>0.0170</b>	0.0125	0.0124	0.0131	0.0020	15.6	0.0109
KEEP	0.0114	<b>0.0145</b>	<u>0.0142</u>	0.0118	<u>0.0127</u>	0.0129	0.0012	9.6	0.0114
CTNet	0.0115	0.0120	<u>0.0137</u>	0.0113	<b>0.0136</b>	0.0124	0.0010	8.3	0.0113
DIIT	0.0134	<u>0.0135</u>	0.0105	0.0103	0.0116	0.0119	0.0014	11.6	0.0103

## 4.2 Tóm tắt chương

Chương này trình bày CNL4Rec, một khung học liên tục cho khuyến nghị đa miền thích nghi. Che mặt nạ miền nhận diện những tham số quan trọng đối với mỗi miền, trong khi chuyên biệt hóa miền điều tiết việc cập nhật gradient sao cho tri thức từ các miền đã gặp trước đó được bảo toàn khi các miền mới được học, qua đó giảm thiểu hiện tượng quên thảm họa. Một mục tiêu công bằng còn khuyến khích hiệu năng cân bằng giữa các miền thay vì chỉ tối ưu cho một miền đích duy nhất. Trên nhiều bộ dữ liệu đa miền chuẩn, CNL4Rec đạt hiệu năng trung bình và hiệu năng miền yếu nhất tốt nhất, với phương sai giữa các miền thấp, vượt các đường cơ sở xuyên miền tiêu biểu.

Cách tiếp cận này dựa trên những giả định đồng thời cũng giới hạn phạm vi của nó. Nó hoạt động tốt nhất khi ranh giới giữa các miền tương đối rõ ràng, và do đó kém phù hợp với những miền chồng lấn cao hay biến đổi nhanh; độ chính xác của nó cũng nhạy với thứ tự các miền xuất hiện, suy giảm khi các miền được trình bày theo thứ tự từ nhỏ nhất đến lớn nhất. Vì hợp của các mặt nạ quan trọng tích lũy lớn dần theo mỗi miền mới, bộ nhớ và độ phức tạp tăng lên khi số miền tiếp tục tích lũy.

# Chương 5

## Khuyến nghị lai dựa trên hội thoại và mạng nơ-ron đồ thị kết hợp sinh tăng cường bằng truy hồi

### 5.1 Giới thiệu

Chương này đề xuất một trợ lý hội thoại khuyến nghị phim, kết hợp khả năng thấu hiểu hành vi quá khứ từ dữ liệu tương tác lịch sử (được nghiên cứu qua một mô hình học sâu dựa trên đồ thị) với các sở thích thời gian thực thu nhận từ các mô hình ngôn ngữ lớn thông qua học tổ hợp. Khung này làm nổi bật những thế mạnh bổ trợ của cả hai kỹ thuật tiên tiến nhằm cung cấp các khuyến nghị có ý nghĩa cho người dùng. Cụ thể, học sâu dựa trên đồ thị có ưu thế trong việc triển khai các tương tác thành dữ liệu có cấu trúc đồ thị, qua đó phát hiện những mẫu hình tiềm ẩn và các quan hệ phức tạp giữa các thực thể nút. Mô hình dựa trên đồ thị cũng sử dụng các thuật toán khả mở rộng và kỹ thuật tính toán phân tán để xử lý khối lượng dữ liệu tương tác khổng lồ, bảo đảm hệ thống vẫn phản hồi nhanh và chính xác khi dữ liệu người dùng và nội dung tăng lên. Mặt khác, các mô hình ngôn ngữ lớn thu thập sở thích cá nhân qua các phiên hội thoại, xử lý đầu vào ngôn ngữ tự nhiên, và thích nghi linh hoạt với những yêu cầu tức thời. Bằng cách kết hợp các thành phần này, hệ thống đạt được một cách tiếp cận khuyến nghị toàn diện, bền vững và thích nghi, cân bằng giữa hiểu biết về hành vi quá khứ với bối cảnh người dùng hiện tại, được tối ưu cho hiệu năng thời gian thực và xử lý dữ liệu quy mô lớn.

Công trình này đã được công bố trong “Improving Retrieval-Augmented Generation for Scalable Movie Chatbots via Graph Based Recommendation Models” 2025 [P5].

## 5.2 CG-RAG: Khuyến nghị hội thoại thông qua sinh tăng cường bằng truy hồi được tăng cường bằng đồ thị

### 5.2.1 Kiến trúc tổng thể của CG-RAG

Bộ sinh khuyến nghị tích hợp dữ liệu hội thoại với cả khả năng thấu hiểu bối cảnh lẫn phân tích hành vi người dùng nhằm đưa ra các gợi ý phim phù hợp. Phương pháp của chúng tôi gồm ba thành phần chính: thành phần hội thoại, thành phần khuyến nghị, và tầng so khớp đặc trưng và truy hồi.

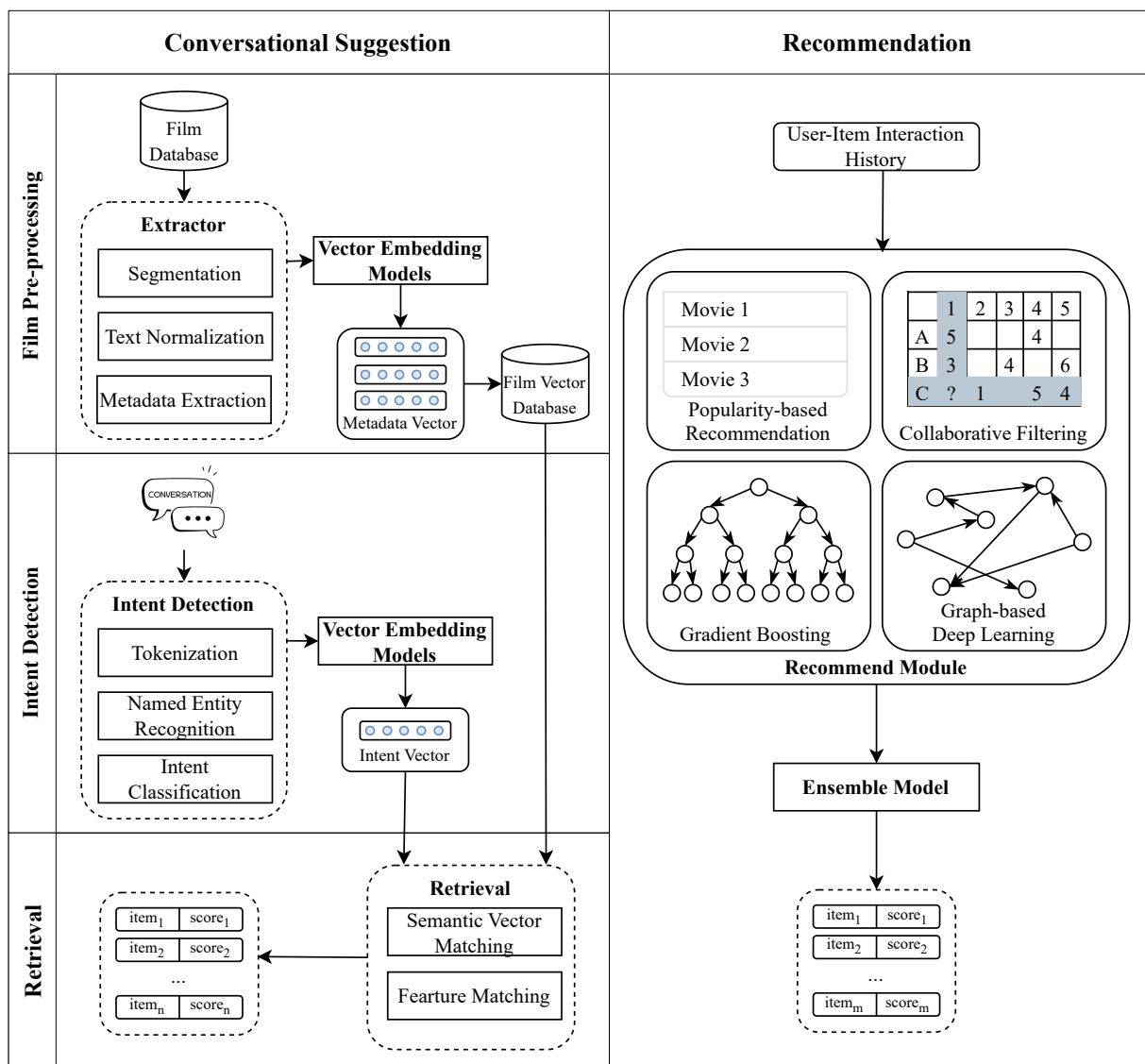


Figure 5.1: Kiến trúc của bộ sinh gợi ý và khuyến nghị hội thoại, gồm một thành phần hội thoại (phát hiện ý định tạo ra một véc-tơ ý định), một thành phần khuyến nghị dựa trên đồ thị, và một tầng so khớp đặc trưng và truy hồi, được kết hợp qua một tầng hợp nhất dùng chung để tạo ra các khuyến nghị hội thoại cuối cùng.

Hình 5.1 minh họa kiến trúc tổng thể của hệ thống khuyến nghị hội thoại lai được đề xuất, được tổ chức thành hai nhánh song song hội tụ tại một tầng hợp nhất dùng chung. *thành phần hội thoại* xử lý truy vấn ngôn ngữ tự nhiên của người dùng qua ba giai đoạn tuần tự: một giai đoạn tiền xử lý chuẩn hóa truy vấn thô và đồng thời trích xuất một véc-tơ Phim và một véc-tơ Siêu dữ liệu từ cơ sở dữ liệu phim, lần lượt mã hóa nội dung sản phẩm và các thuộc tính cấu trúc; một giai đoạn Phát hiện Ý định áp dụng Nhận dạng Thực thể có tên và nhúng véc-tơ để nhận diện các thực thể riêng của miền (ví dụ thể loại, diễn viên) và phân loại ý định của người dùng thành một véc-tơ Ý định cô đọng; và một giai đoạn Truy hồi tính độ tương đồng ngữ nghĩa giữa véc-tơ Ý định và các véc-tơ sản phẩm để tạo ra một danh sách xếp hạng các ứng viên liên quan về mặt ngữ nghĩa. *thành phần khuyến nghị* hoạt động độc lập trên dữ liệu tương tác người dùng-sản phẩm lịch sử và được xây dựng dựa trên một đánh giá thực nghiệm có hệ thống đối với các phương pháp khuyến nghị sâu tiêu biểu. Cụ thể, chúng tôi đối sánh mười hai mô hình trải rộng trên lan truyền đồ thị, học tương phản đồ thị, và các cách tiếp cận căn chỉnh biểu diễn nhúng/tự giám sát. Điểm xếp hạng hành vi rút ra từ các mẫu hình tương tác người dùng-sản phẩm lịch sử được hợp nhất một cách độc lập với điểm truy hồi do thành phần truy hồi hội thoại tạo ra, vốn neo các sản phẩm ứng viên vào bối cảnh hội thoại hiện tại thông qua sinh tăng cường bằng truy hồi. Cuối cùng, *tầng so khớp đặc trưng và truy hồi* đóng vai trò là điểm hợp nhất của toàn hệ thống, căn chỉnh các danh sách ứng viên từ cả hai nhánh theo định danh sản phẩm và cùng chấm điểm lại mỗi sản phẩm dựa trên độ liên quan ngữ nghĩa của nó với ý định người dùng đã biểu đạt và cường độ sở thích dự đoán từ lịch sử hành vi, qua đó tạo ra một danh sách khuyến nghị cuối cùng vừa nhạy bối cảnh vừa cá nhân hóa.

## 5.3 Kết quả thực nghiệm

Table 5.1: So sánh hiệu năng giữa các phương pháp cơ sở trên hai bộ dữ liệu Movielens-1M và TV360. R@K là Recall top-K, và P@K là Precision top-K. Kết quả tốt nhất của mỗi bộ dữ liệu được in đậm, còn kết quả tốt thứ nhì được gạch chân.

Bộ dữ liệu	Movielens-1M					TV360				
	R@30	P@30	R@50	P@50	Thời gian	R@30	P@30	R@50	P@50	Thời gian
LightGCN	0.1628	0.1040	0.2373	0.0964	1e-6	0.0423	0.1018	0.0701	0.1184	1e-6
GAT	0.1650	0.0564	0.2395	0.0515	1e-5	0.0323	0.1067	0.0529	<b>0.1406</b>	1e-5
PMLP	0.1292	0.0467	0.1920	0.0430	1e-6	0.0127	0.0433	0.0223	0.0449	1e-5
GraphSAGE	0.1482	0.0526	0.2152	0.0482	1e-6	0.0400	<u>0.1246</u>	0.0666	<u>0.1251</u>	1e-6
LINKX	0.1184	0.0406	0.1672	0.0371	1e-6	0.0376	0.1212	0.0609	0.1176	1e-6
MixGCF	0.0192	0.0208	0.0432	0.0256	1e-5	0.0369	0.0269	0.0550	0.0248	1e-5
SGL	0.0192	0.0208	0.0295	0.0196	1e-6	0.0681	0.0449	0.0933	0.0385	1e-6
SimGCL	0.2125	0.1467	0.2872	0.1224	1e-6	0.1038	0.0672	0.1467	0.0593	1e-5
XSimGCL	<u>0.2209</u>	<u>0.1494</u>	<u>0.2905</u>	<u>0.1267</u>	1e-6	<u>0.1725</u>	0.1026	<u>0.2362</u>	0.0885	1e-6
NCL	<b>0.2497</b>	<b>0.1990</b>	<b>0.3351</b>	<b>0.1685</b>	1e-6	0.0915	0.0584	0.1296	0.0514	1e-6
SSL4Rec	0.1933	0.1208	0.2713	0.1068	1e-5	0.1082	0.0730	0.1644	0.0678	1e-5
DirectAU	0.1315	0.0835	0.1820	0.0723	1e-5	<b>0.1932</b>	<b>0.1250</b>	<b>0.2714</b>	0.1095	1e-5

Số in đậm = Tốt nhất, số gạch chân = Tốt thứ nhì

### Kết quả Truy hồi Hội thoại

Table 5.2: Điểm Recall@K trên bộ dữ liệu Movielens-1M qua 3 mô hình ngôn ngữ lớn. Kết quả được đo ở nhiều giá trị K từ 10 đến 200, kèm theo thời gian sinh câu trả lời tính bằng giây. Phản hồi sinh ra trả về bộ phim hàng đầu do các mô-đun trợ lý hội thoại cung cấp.

	Recall@10	Recall@20	Recall@30	Recall@50	Recall@100	Recall@200	Thời gian (s)
Gemma-2B	0.186	0.201	0.214	0.259	0.401	0.829	2.05
LLaMA-3B	<u>0.204</u>	<u>0.227</u>	<u>0.239</u>	<u>0.291</u>	0.428	<u>0.869</u>	<u>2.28</u>
Qwen-3B	<b>0.207</b>	<b>0.232</b>	<b>0.246</b>	<b>0.299</b>	<b>0.435</b>	<b>0.888</b>	<b>2.42</b>

Số in đậm = Tốt nhất, số gạch chân = Tốt thứ nhì

Table 5.3: Điểm Recall@K trên bộ dữ liệu TV360 qua 3 mô hình ngôn ngữ lớn. Kết quả được đo ở nhiều giá trị K từ 10 đến 200, kèm theo thời gian sinh câu trả lời tính bằng giây. Phản hồi sinh ra trả về bộ phim hàng đầu do cả hai mô-đun cung cấp.

	Recall@10	Recall@20	Recall@30	Recall@50	Recall@100	Recall@200	Thời gian (s)
Gemma-2B	0.060	0.071	0.084	0.106	0.221	0.703	1.88
LLaMA-3B	<u>0.072</u>	<b>0.086</b>	<u>0.093</u>	<u>0.112</u>	<u>0.226</u>	<u>0.731</u>	<u>2.13</u>
Qwen-3B	<b>0.073</b>	<u>0.083</u>	<b>0.095</b>	<b>0.118</b>	<b>0.233</b>	<b>0.740</b>	<b>2.18</b>

Số in đậm = Tốt nhất, số gạch chân = Tốt thứ nhì

## Kết quả Tổ hợp

Table 5.4: So sánh hiệu năng giữa các phương pháp khuyến nghị cơ sở khi kết hợp với thành phần hội thoại trên hai bộ dữ liệu Movielens-1M và TV360. R@K là Recall top-K và P@K là Precision top-K.

Bộ dữ liệu	Movielens-1M				TV360			
	R@10	R@20	R@30	R@50	R@10	R@20	R@30	R@50
XSimGCL	<u>0.2046</u>	<u>0.2353</u>	<u>0.2503</u>	<u>0.3012</u>	<b>0.1258</b>	<b>0.1532</b>	<b>0.1812</b>	<b>0.2681</b>
NCL	<b>0.2104</b>	<b>0.2387</b>	<b>0.2524</b>	<b>0.3402</b>	0.0657	0.0884	0.1021	0.1458
DirectAU	0.1365	0.1581	0.1872	0.2447	<u>0.1106</u>	<u>0.1582</u>	<u>0.1808</u>	<u>0.2606</u>

Số **in đậm** = Tốt nhất, số gạch chân = Tốt thứ nhì

## 5.4 Tóm tắt chương

Chương này giới thiệu CG-RAG, một kiến trúc khuyến nghị hội thoại lai, bắc cầu giữa hành vi dài hạn của người dùng với ý định hội thoại thời gian thực. Một mạng nơ-ron đồ thị mô hình hóa các sở thích dài hạn có cấu trúc, một mô hình ngôn ngữ lớn đảm nhận việc trích xuất ý định và đối thoại bằng ngôn ngữ tự nhiên, và một giai đoạn truy hồi lai kết hợp truy hồi thưa (BM25) và truy hồi dày để neo các khuyến nghị vào kho mục sản phẩm thực tế. Vì các sản phẩm ứng viên bị giới hạn trong những mục được truy hồi và bộ sinh được điều kiện hóa trên siêu dữ liệu truy hồi được, hệ thống tạo ra các khuyến nghị nhảy bối cảnh và có neo nguồn thông qua đối thoại tương tác thay vì các danh sách xếp hạng tĩnh.

# Conclusions

## Tóm tắt các đóng góp

Luận án này phát triển các mô hình học sâu cho các hệ khuyến nghị hiện đại, được tổ chức quanh ba đóng góp chính tương ứng với ba thách thức nền tảng đã được nhận diện trong phần mở đầu: mô hình hóa bền vững và khả mở rộng đối với dữ liệu chính tắc và dữ liệu phụ trợ, khuyến nghị đa miền thích nghi, và khuyến nghị hội thoại. Đóng góp thứ nhất là một tập hợp các mô hình cho việc mô hình hóa sâu bền vững và khả mở rộng đối với dữ liệu chính tắc và dữ liệu phụ trợ, giải quyết các vấn đề khả năng mở rộng, thừa thớt dữ liệu và khởi động nguội. Đóng góp thứ hai là một mô hình khuyến nghị đa miền thích nghi dựa trên học liên tục. Đóng góp thứ ba là một kiến trúc khuyến nghị hội thoại lai, bắc cầu giữa hành vi dài hạn của người dùng với ý định người dùng thời gian thực, biến đổi theo thời gian. Nhìn chung, các mô hình được đề xuất đạt hiệu năng tốt trên nhiều bộ dữ liệu chuẩn công khai và một bộ dữ liệu công nghiệp thực tế, thiết lập cả những hiểu biết lý thuyết lẫn các giải pháp thực tiễn, thu hẹp khoảng cách giữa nghiên cứu hàn lâm và triển khai công nghiệp quy mô lớn.

## Hạn chế

Bên cạnh những đóng góp của luận án, hiện tại vẫn còn một số hạn chế gợi mở các hướng nghiên cứu trong tương lai.

**Tài nguyên tính toán:** Việc huấn luyện các mạng nơ-ron đồ thị, các mô-đun tương phản và các cập nhật học liên tục vẫn tốn nhiều tài nguyên ở quy mô lớn, ngay cả khi EfficientRec đã giảm dung lượng bộ nhớ phía người dùng xuống  $O(K \cdot d)$ .

**Chất lượng thông tin phụ trợ, độ bền vững và tính công bằng:** Thông tin phụ trợ thường nhiễu, thiếu sót hoặc thiên lệch, và các thuộc tính thiên lệch có thể làm suy giảm tính công bằng. Luận án đã giảm nhẹ một phần vấn đề này, và các giải pháp tiếp theo bao gồm khử nhiễu và đánh trọng số theo độ bất định cho nhiễu, bù khuyết và tiền huấn luyện tự giám sát cho dữ liệu thiếu, cùng khử thiên lệch theo nghịch xu hướng hoặc phản thực với ràng buộc ngang bằng độ phơi bày cho thiên lệch và công bằng.

**Giả định về ranh giới miền:** Khung học liên tục giả định ranh giới giữa các miền tương đối rõ ràng và kém hiệu quả đối với những miền chồng lấn cao hay biến đổi nhanh, hoặc đối với các kịch bản liên quan đến việc sắp xếp lại tác vụ và tái xuất hiện miền.

**Độ tin cậy của hệ thống hội thoại:** Việc tích hợp các mạng nơ-ron đồ thị, các mô hình ngôn ngữ lớn và sinh tăng cường bằng truy hồi giúp nâng cao khả năng suy luận và tương tác, nhưng đưa vào độ trễ suy luận, rủi ro ảo giác, cùng các mối lo về khả năng kiểm soát, quyền riêng tư và an toàn vốn chưa được giải quyết trọn vẹn.

**Phạm vi đánh giá:** Mặc dù các thực nghiệm trải rộng trên nhiều bộ dữ liệu chuẩn công khai và một bộ dữ liệu công nghiệp thực tế, chúng vẫn không thể bao phủ toàn bộ sự đa dạng của các kịch bản thực tiễn; việc khái quát sang những miền như tài chính, y tế và giáo dục đòi hỏi thêm kiểm chứng thực nghiệm và sự điều chỉnh riêng cho từng miền.

## Hướng phát triển

Các nghiên cứu tương lai có thể mở rộng những kết quả của luận án này theo nhiều hướng tiềm năng.

**Tối ưu hướng sản xuất:** Phát triển các chiến lược tinh chỉnh liên tục thích nghi mô hình dựa trên tương tác người dùng thực, nhật ký hệ thống và sự trôi dạt hành vi.

**Cá nhân hóa đa phương thức:** Mở rộng khung khuyến nghị ra ngoài tín hiệu tương tác có cấu trúc và tín hiệu văn bản để tích hợp thông tin đa phương thức, bao gồm hình ảnh, âm thanh, video và các tín hiệu hành vi ở mức thiết bị.

**Triển khai hiệu quả:** Khám phá các kiến trúc nhẹ cho suy luận trên thiết bị, nén mô hình ngôn ngữ lớn tiết kiệm chi phí, và học đa phương thức bảo toàn quyền riêng tư để triển khai các hệ khuyến nghị ở quy mô lớn.

**Công bằng và bền vững:** Đưa các ràng buộc công bằng tường minh và các kỹ thuật giảm thiểu thiên lệch vào các hàm mục tiêu tối ưu.

# List of Publications

- [P 1] Vu Hong Quan, Le Hoang Ngan, Le Minh Duc, **Nguyen Tran Ngoc Linh**, Le Hoang Quynh. "EfficientRec: An Unlimited User Scale Recommendation System Based on Clustering and User's Interaction Embedding Profile." *In Asian Conference on Intelligent Information and Database Systems (ACIIDS)*, pp. 681-696. Springer Nature Singapore, 2022. (*Scopus Conference*)
- [P 2] **Nguyen Tran Ngoc Linh**, Vu Chi Dung, Le Hoang Ngan, Hoang Anh Dung, Phan Xuan Hieu, Ha Quang Thuy, Le Hoang Quynh, Tran Mai Vu. "GIFT4Rec: An Effective Side Information Fusion Technique Apply to Graph Neural Network for Cold-Start Recommendation." *In Asian Conference on Intelligent Information and Database Systems (ACIIDS)*, pp. 334-345. Springer Nature Singapore, 2023. (*Scopus Conference*)
- [P 3] **Nguyen Tran Ngoc Linh**, Le Hoang Ngan, Hoang Anh Dung, Phan Xuan Hieu, Ha Quang Thuy, Le Hoang Quynh, Tran Mai Vu. "The Masked Simple Graph Contrastive Learning for Recommendation." *In 2024 16th International Conference on Knowledge and System Engineering (KSE)*, pp. 156–160. IEEE, 2024. (*Scopus Conference*)
- [P 4] **Nguyen Tran Ngoc Linh**, Vu Chi Dung, Le Hoang Ngan, Hoang Anh Dung, Phan Xuan Hieu, Ha Quang Thuy, Le Hoang Quynh, Tran Mai Vu. "Continual Learning Based on Task Masking for Multi-domain Recommendation." *In Asian Conference on Intelligent Information and Database Systems (ACIIDS)*, pp. 257-266. Springer Nature Singapore, 2024. (*Scopus Conference*)
- [P 5] **Nguyen Tran Ngoc Linh**, Hoang Anh Dung, Tran Manh Cuong, Vu Minh Thanh, Vu The Anh, Nguyen Xuan Bach, Bui Tuan Nghia, Le Hoang Quynh, Vuong Thi Hai Yen, Tran Mai Vu. "Improving Retrieval-Augmented Generation for Scalable Movie Chatbots via Graph Based Recommendation Models" Submitted to IEEE Access- under minor revision (Round 3), 2026. (*Q1 Journal*)