

and Contributions*

g, Data Processing, Algorithm Implementation

g, Algorithm Implementation, Visualization

g, Data Processing, Algorithm Implementation

Danning Sui | ds3516

Qiong Hu | qh2174

Modeling, Data Processing, Algorithm Implementation

Modeling, Data Processing, Algorithm Implementation

****Picture's Copyright from: <http://www.evincedev.com/ecommerce-development>**

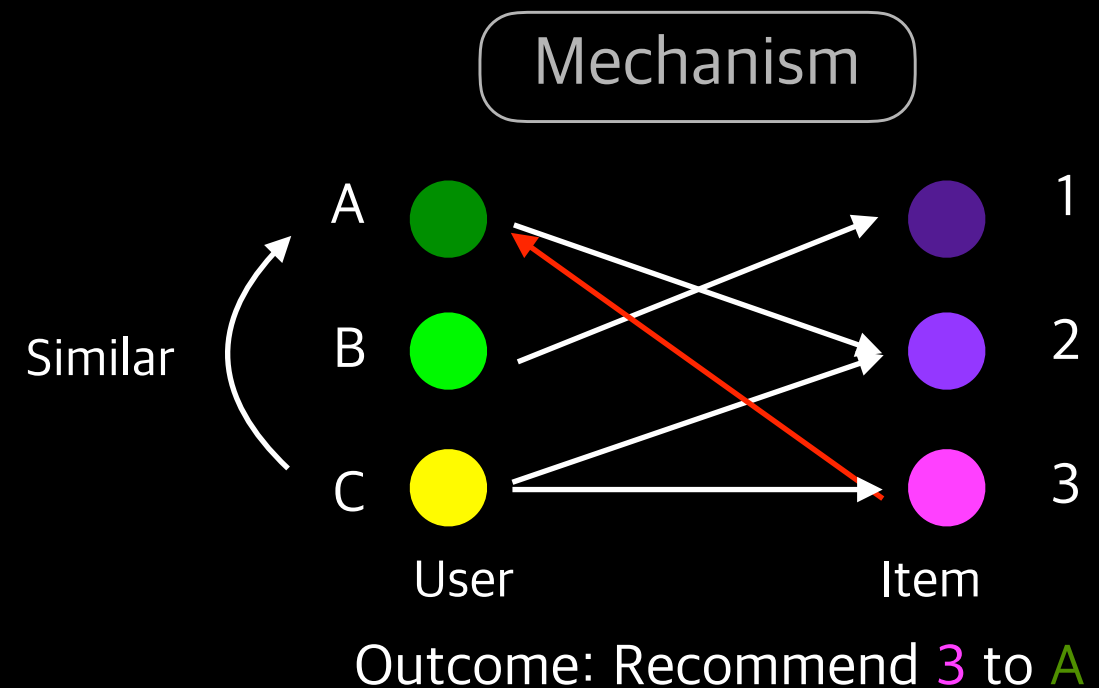
Overview

- Combination of **User-based Collaborative Filtering** and **Item to Item Recommendation** Process
- Pursuing a **personalized** and **abundant** recommendation set to open the market

User-based CF

Input: A Bipartite of user and product

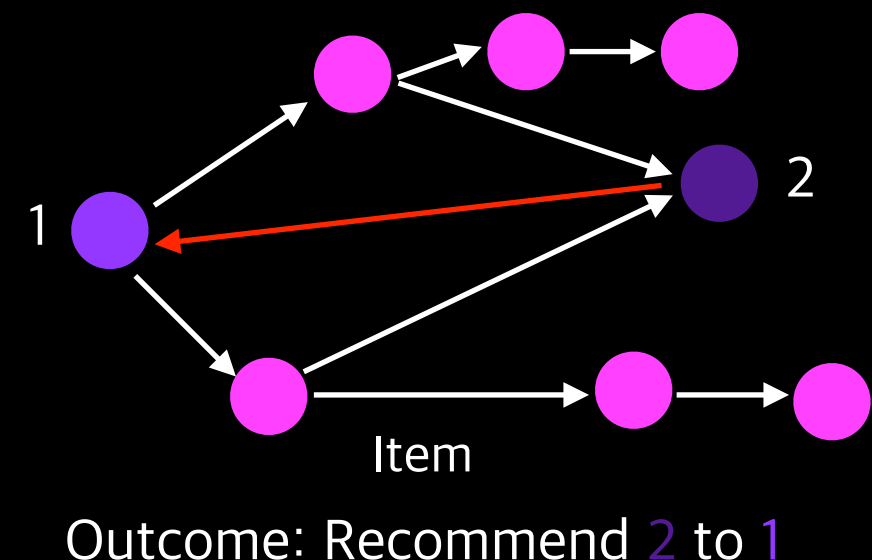
According to their rating reviews, find the favored products from similar users.



Item to Item Recommendation

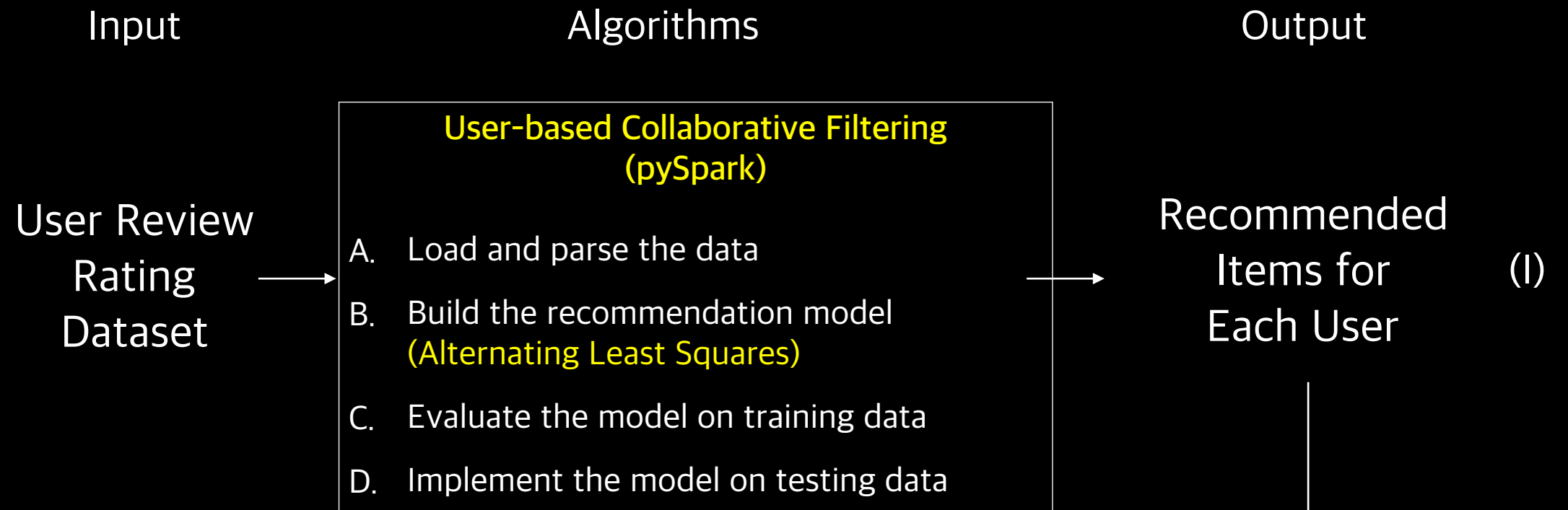
Input: A network of products

According to their co-purchasing network, compute their similarities and recommend.

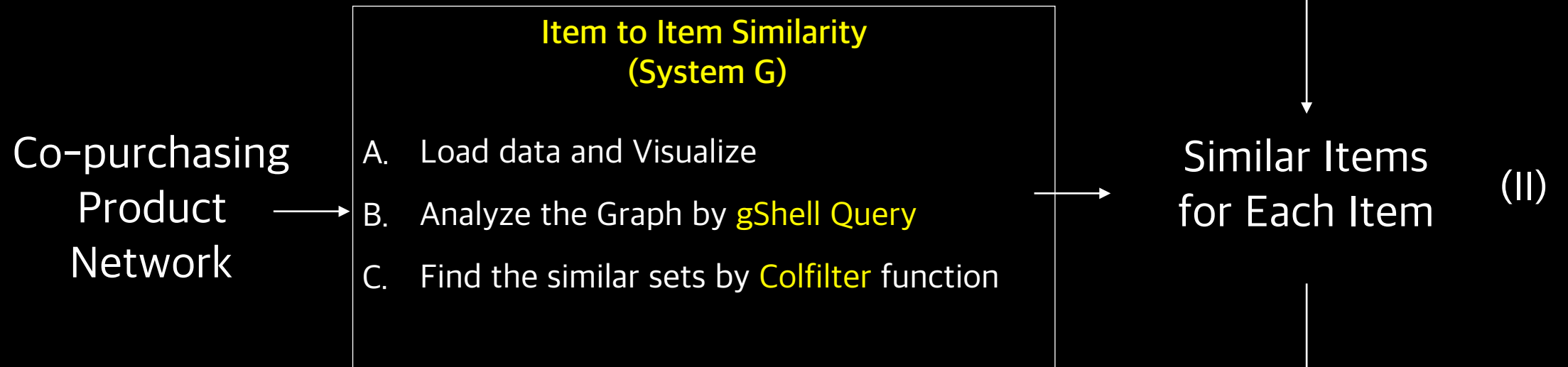


Model and Algorithms

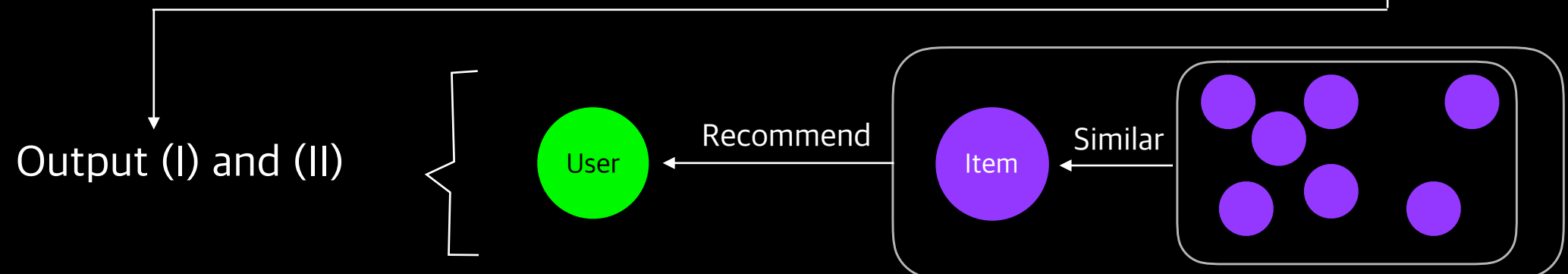
Part 1



Part 2



Part 3



Datasets (1.14GB)

A. User Review Rating Matrix – 977.5MB

(amazon-meta.txt)

Dataset statistics		Products by product group	
Products	548,552	Books	393561
Product-Project Edges	1,788,725	DVDs	19828
Reviews	7,781,990	Music CDs	103144
Product category memberships	2,509,699	Videos	26132

Id: 1

ASIN: 0827229534

title: Patterns of Preaching: A Sermon Sampler

group: Book

salesrank: 396585

similar: 5 0804215715 156101074X 0687023955 0687074231 082721619X

categories: 2

|Books[283155]|Subjects[1000]|Religion & Spirituality[22]|Christianity[12290]|Clergy[12360]|
Preaching[12368]

|Books[283155]|Subjects[1000]|Religion & Spirituality[22]|Christianity[12290]|Clergy[12360]|
Sermons[12370]

reviews: total: 2 downloaded: 2 avg rating: 5

2000-7-28 cutomer: A2JW67OY8U6HHK rating: 5 votes: 10 helpful: 9

2003-12-14 cutomer: A2VE83MZF98ITY rating: 5 votes: 6 helpful: 5

Datasets (1.14GB)

B. Co-purchasing Item Network – 158.1MB in total

(Amazon0302.txt, Amazon0312.txt, Amazon0505.txt, Amazon0601.txt)

Example – Amazon0312.txt

Dataset statistics	
Nodes	262111
Edges	1234877
Nodes in largest WCC	262111 (1.000)
Edges in largest WCC	1234877 (1.000)
Nodes in largest SCC	241761 (0.922)
Edges in largest SCC	1131217 (0.916)
Average clustering coefficient	0.4198
Number of triangles	717719
Fraction of closed triangles	0.09339
Diameter (longest shortest path)	32
90-percentile effective diameter	11

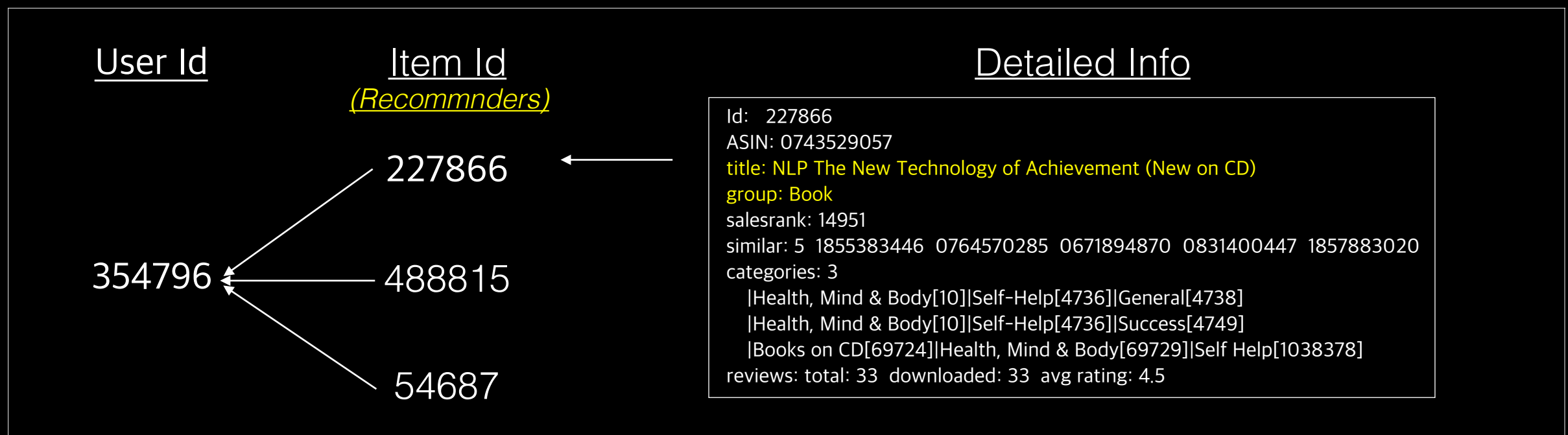
```
# Directed graph (each unordered pair of nodes is saved
once): Amazon0312.txt
# Amazon product co-purchasing network from March 12
2003
# Nodes: 400727 Edges: 3200440
# FromNodeId    ToNodeId
0      1
0      2
0      3
0      4
0      5
1      0
1      2
1      13
1      14
1      15
2      0
2      1
2      4
2      5
2      16
3      70
5      6
```

Outcomes

Model Part 1 : User-based Collaborative Filtering (pySpark)

Rank	User Id	Item Id	Original Score	Model Score	Error
#1	1080592	284560	5	4.99797677899322	0.00202322100678
#2	354796	227866	5	4.99583012902760	0.00416987097240
#3		488815	5	4.99583012902760	0.00416987097240
#4		54687	5	5.00331265792859	0.00331265792859
#5	1056524	274052	5	5.00026234076471	0.00026234076471
#6		500789	5	4.99341921309295	0.00658078690705
#7		354097	5	4.99341921309295	0.00658078690705
#8	199516	434236	5	4.99745108641090	0.00254891358910
#9		329445	5	4.99745108641090	0.00254891358910
#10		453871	5	4.99745108641090	0.00254891358910

Example



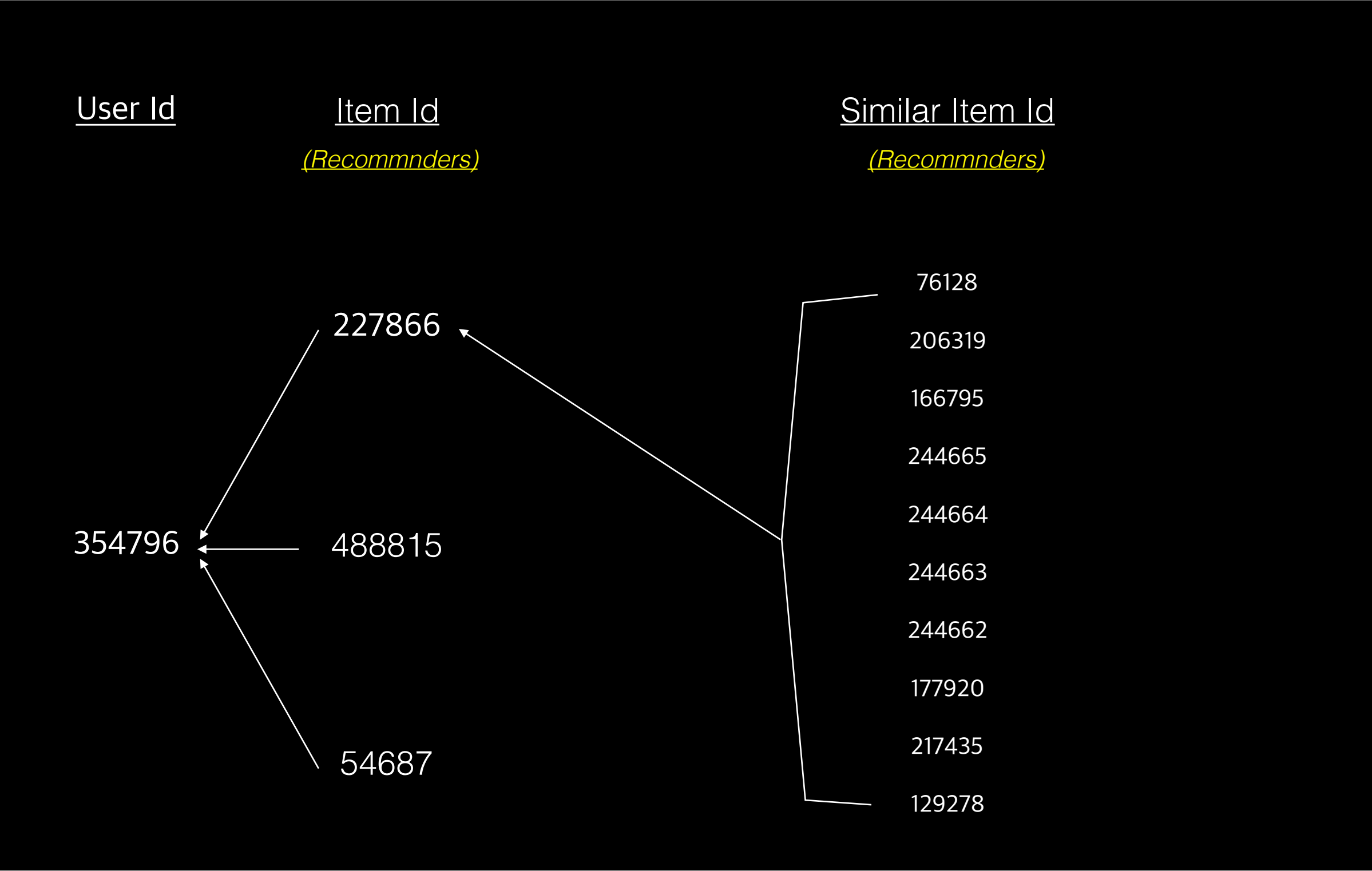
Model Part 2 : Item to Item Similarity (System G) - depth=4

Example



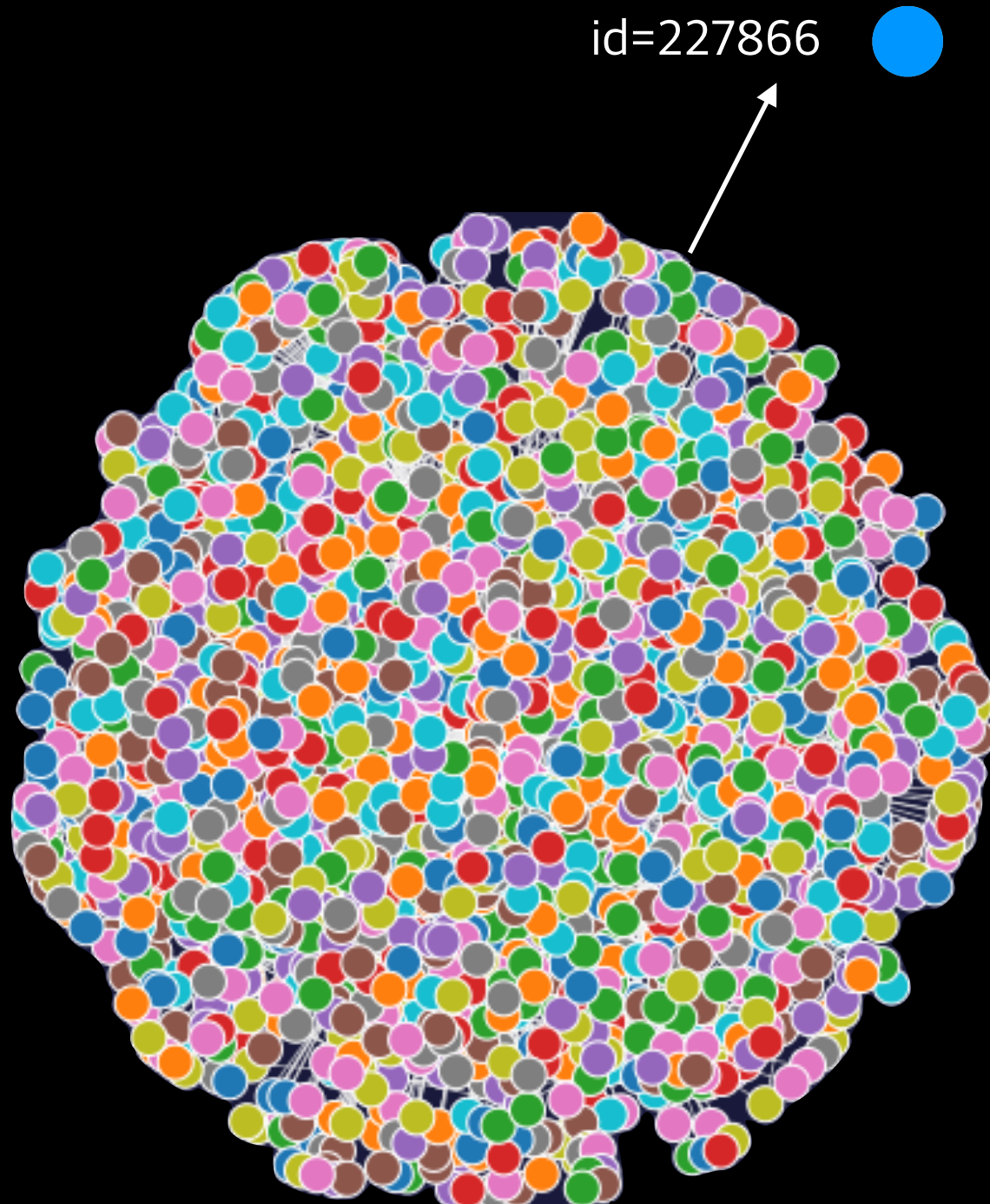
Outcomes

Model Part 3 : Recommendation Chain



Outcomes

Model Part 2 Visualization



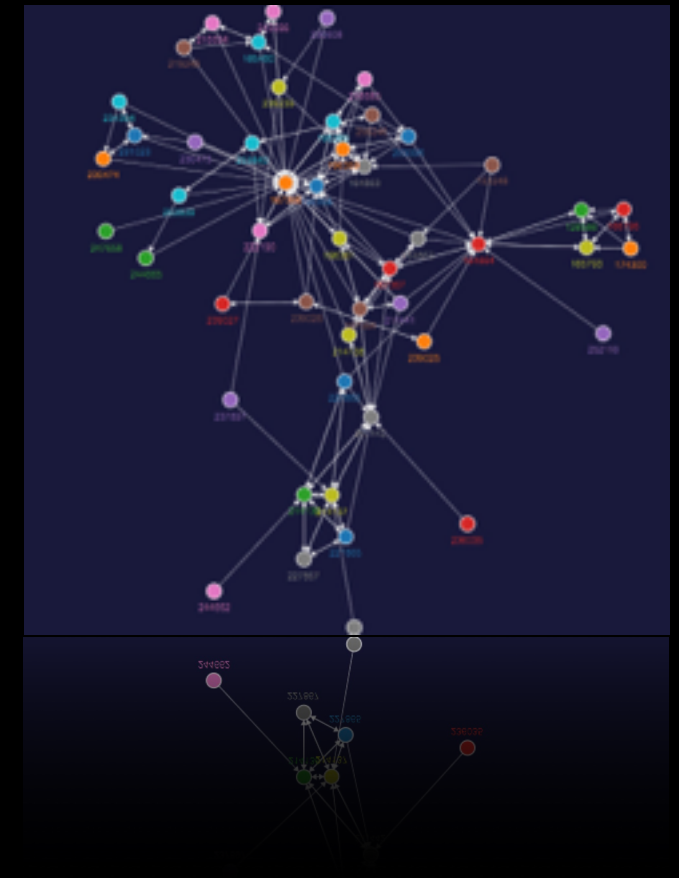
Neighbor Graph
(depth=1)

#Nodes: 5
#Edges: 9



Neighbor Graph
(depth=2)

#Nodes: 49
#Edges: 169



Neighbor Graph
(depth=4)

#Nodes: 245
#Edges: 969



Challenges

- Large Data Size Process Efficiency
- Fitness of CF model
 - ALS algorithm (Mean Squared Error = 1.09519420585)
 - Dataset don't perfectly match the requirement (User:Item=1:10)
- Product information were not fully made use of

Further Work

- Analyze **ALS algorithm** to understand why the top ten recommenders are centered around some users
- Take into account of **group tag** and **category tag**
- Compare the model performances on different “**depth**”