

# Towards Trustworthy Recommender Systems: From Shallow Model to Deep Model to Large Model

Yongfeng Zhang

Department of Computer Science, Rutgers University

[yongfeng.zhang@rutgers.edu](mailto:yongfeng.zhang@rutgers.edu)

<http://yongfeng.me>

## Recommender Systems are Everywhere

- Influence our daily life by providing personalized services

### E-commerce



### Social Networks



### News Feeding



### Search Engine



### Navigation



### Travel Planning



### Professional Networks



### Healthcare

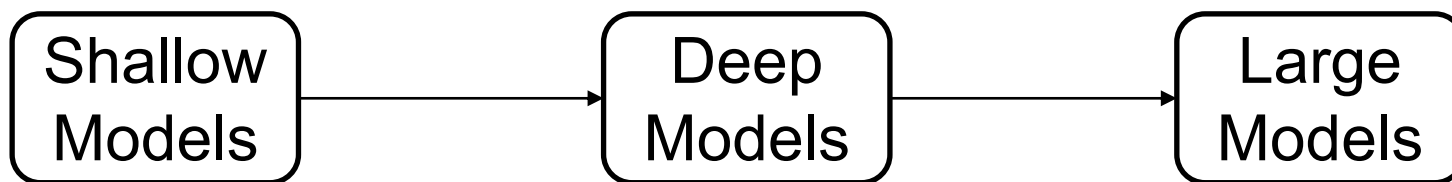


### Online Education



# Technical Advancement of Recommender Systems

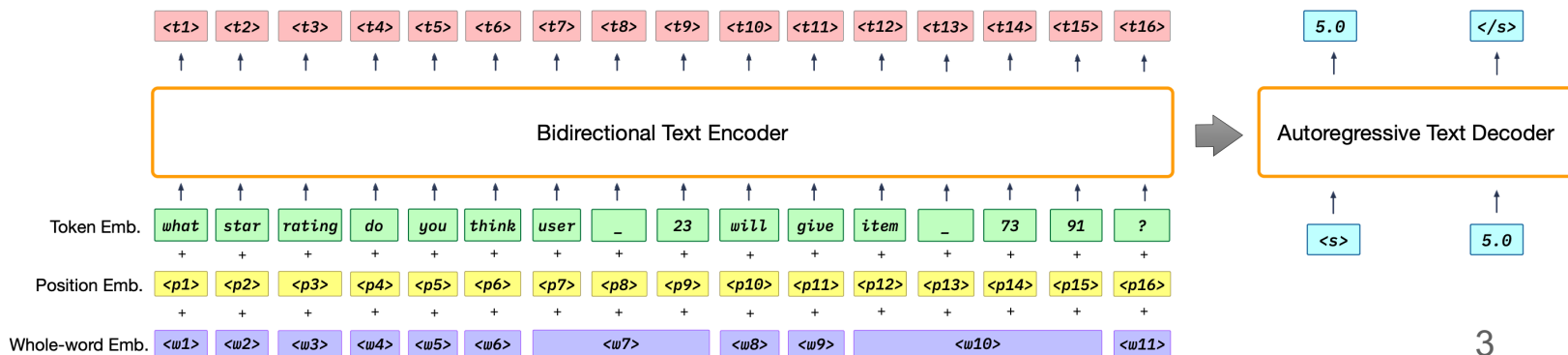
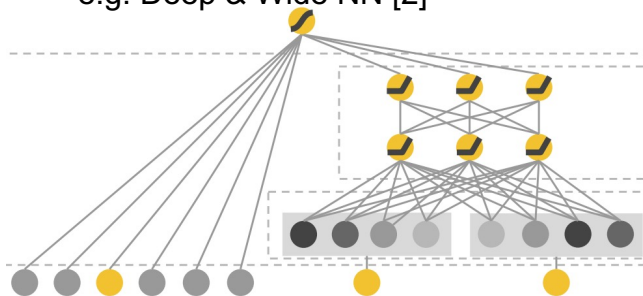
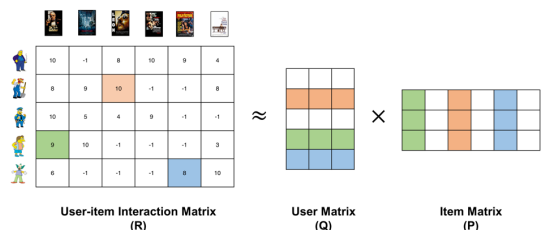
- From Shallow Model, to Deep Model, and to Large Model



e.g. Matrix Factorization [1]

e.g. Deep & Wide NN [2]

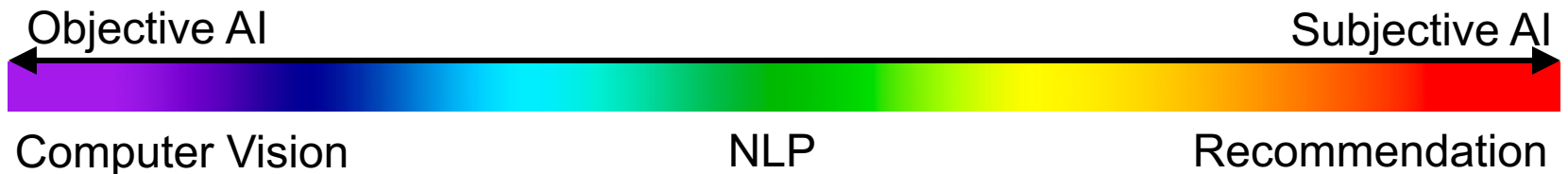
e.g. P5 [3]



[1] Koren, Yehuda, Robert Bell, and Chris Volinsky. "Matrix factorization techniques for recommender systems." *Computer* 42, no. 8 (2009): 30-37.  
 [2] Cheng, Heng-Tze, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson et al. "Wide & deep learning for recommender systems." *DLRS* 2016.  
 [3] Geng, Shijie, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. "Recommendation as Language Processing (RLP): A Unified Pretrain, Personalized Prompt & Predict Paradigm (P5)." *RecSys* 2022.

# Objective AI vs. Subjective AI

- Recommendation is **unique** in the AI family
  - Recommendation is most **close to human** among all AI tasks
  - Recommendation is a very representative **Subjective AI**
  - Thus, leads to many **unique challenges** in recommendation research



(Relatively) far from human.  
Problems have exact answers.



Very close to human.  
Problems have no absolute answers.





# Computer Vision: (mostly) Objective AI Tasks

Objective AI

Subjective AI



Computer Vision

NLP

Recommendation

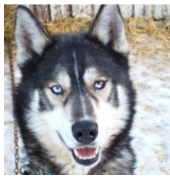
Image Classification



cat

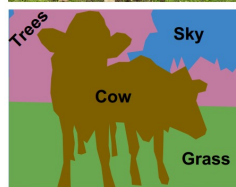
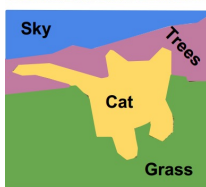


dog

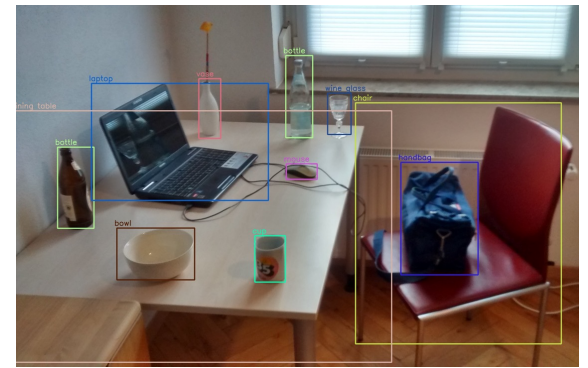


Husky like a wolf

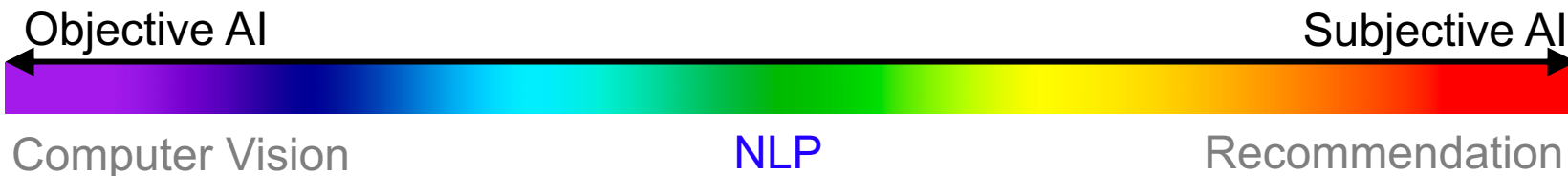
Image Segmentation



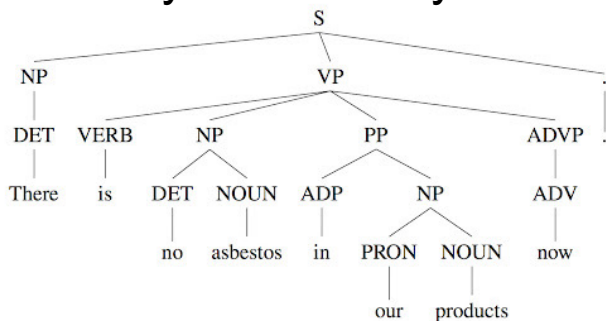
Object Detection



## NLP: partly Objective, partly Subjective



### Syntactic Analysis



### Word Segmentation

Words: 这是 一篇 有趣的 文章

(zhèshì yīpiān yǒuqù de wénzhāng)

### Dialog Systems

Can you find me a *mobile phone* on Amazon?  
Sure, what *operating system* do you prefer? 🤖

I want an *Android* one.  
OK, and any preference on *screen size*? 🤖

Better larger than *5 inches*.  
Do you have requirements on *storage capacity*? 🤖

I want it to be at least *64 Gigabytes*.  
And any preference on *phone color*? 🤖

*Not particularly*.  
Sure, then what about the following choices? 🤖

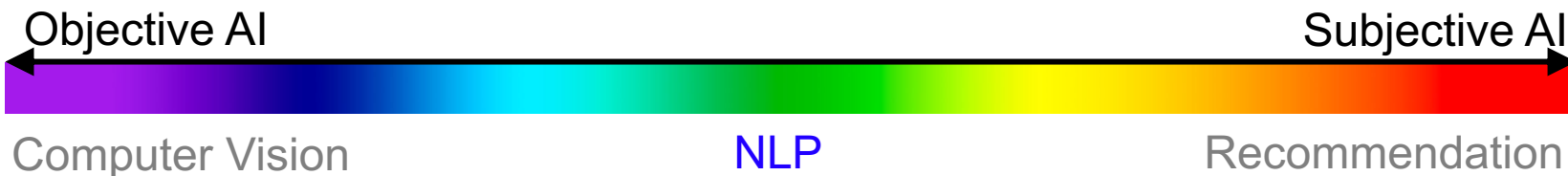
I don't like them very much...  
OK, do you have any preference on the *brand*? 🤖

Better be *Samsung or Huawei*.  
Any requirement on *price*? 🤖

Should be *within 700 dollars*.  
OK, then what about these ones? 🤖

Great, I want the first one, can you order it for me?  
Sure, I have placed the order for you, enjoy! 🤖

# Recommendation: mostly Subjective AI Tasks



## Movie Recommendation



## Product Recommendation



# Subjective AI needs Explainability

- Objective vs. Subjective AI on Explainability

## Objective AI

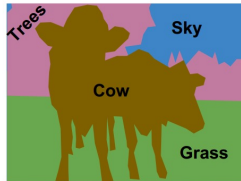
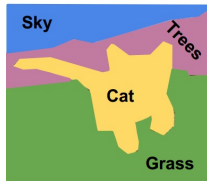
Human can directly identify if the AI-produced result is right or wrong



cat



dog



## Subjective AI

Human can hardly identify if the AI-produced result is right or wrong



Can you find me a *mobile phone* on Amazon?  
 Sure, what **operating system** do you prefer?  
 I want an **Android** one.  
 OK, and any preference on **screen size**?  
 Better larger than **5 inches**.  
 Do you have requirements on **storage capacity**?  
 I want it to be at least **64 Gigabytes**.  
 And any preference on **phone color**?  
**Not particularly**.  
 Sure, then what about the following choices?

I don't like them very much...  
 OK, do you have any preference on the **brand**?  
 Better be **Samsung or Huawei**.  
 Any requirement on **price**?  
 Should be **within 700 dollars**.  
 OK, then what about these ones?

Great, I want the first one, can you order it for me?  
 Sure, I have placed the order for you, enjoy!

Nothing is definitely right or wrong.

Highly **subjective**, and usually **personalized**.

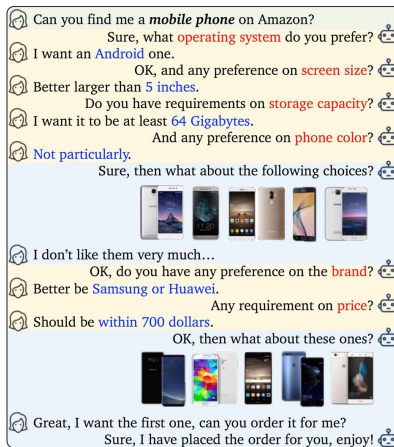
# Subjective AI needs Explainability

- In many cases, it doesn't matter what you recommend, but how you explain your recommendation
- How do humans make recommendation?



# Subjective AI needs Fairness

- Users cannot easily identify if something is right or wrong
  - They have to take the recommendations as is
  - Users are very **vulnerable**
  - Users could be **manipulated**, **utilized** or even **cheated** by the system



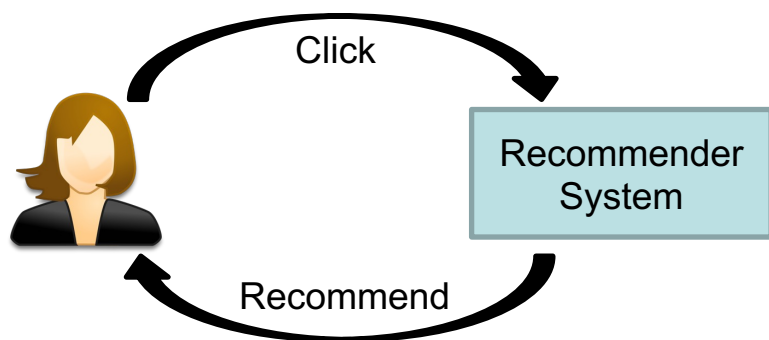
Nothing is definitely right or wrong.

Highly **subjective**, and usually **personalized**.

Users need to be treated fairly.

# Subjective AI leads to Echo Chambers

- Users don't know which recommendations are “right” and which are “wrong”, they just **click**. [5]
- Lack of **explanation** makes the problem worse.



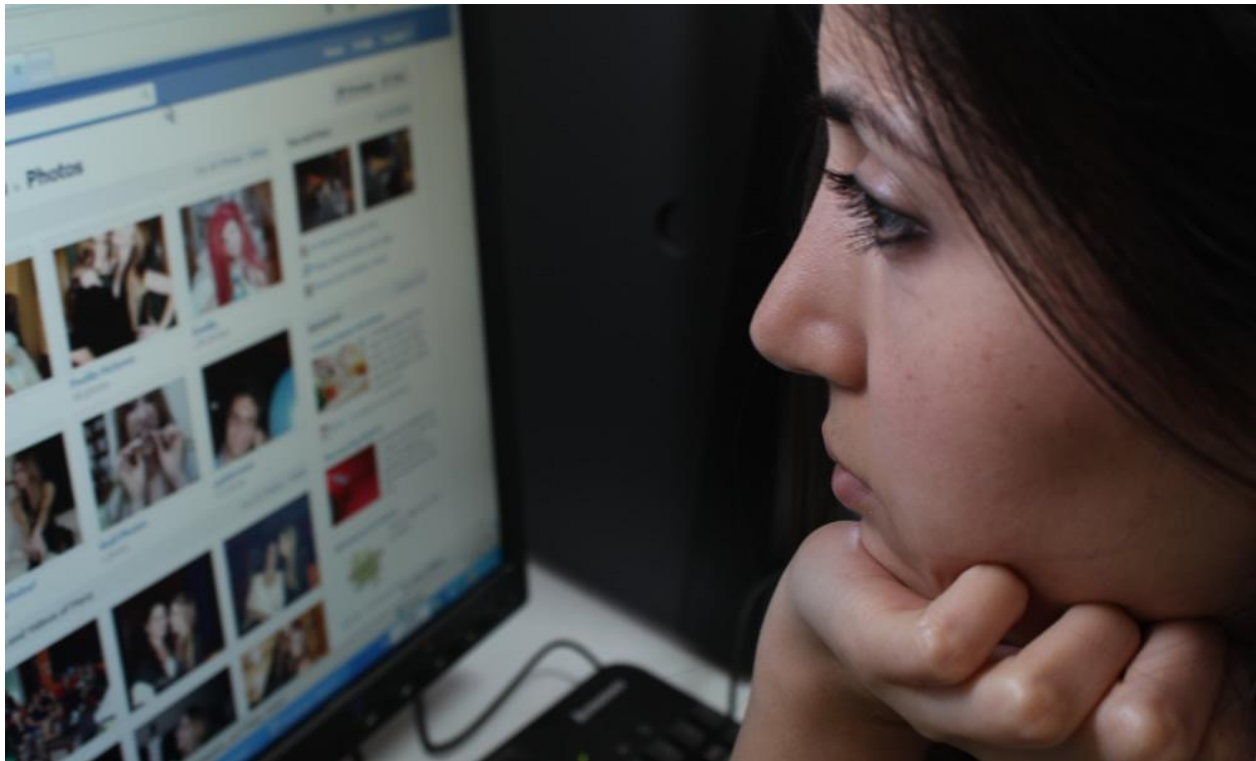
The more you like something, the more RS will recommend similar things, and thus you like them even more.





# Subjective AI needs Controllability

- Users almost have **no control** of their recommender system
  - They can only **passively** receive recommendations





# Trustworthy and Responsible Recommendation

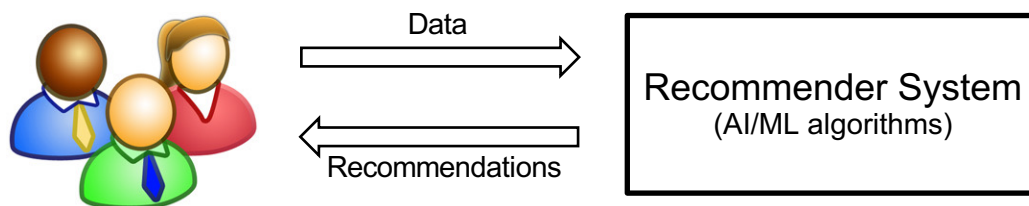
- Explainability, Fairness, Echo Chambers, Controllability
- And many more ...
  - Robustness, Accountability, Privacy, etc.

Responsible  
AI



# RecSys as a Human-centered AI task

- Recommender System (RS) is a representative Human-centered AI task
  - Naturally involves human-in-the-loop
  - Influences human decision making everyday and everywhere




- A wide scope of applications

 E-commerce (product recommendation)




 Smart and Connected Communities (driving route recommendation / passenger recommendation)

  Social Networks (friend/tweet recommendation)


 Sharing Economy (house recommendation)

 Search Engines (personalized search / advertising)

 Travel and Planning Services (ticket and hotel recommendation)


 Professional Networks (job recommendation)

Even some high-stake application scenarios

 Financial Services (financial / investment recommendation)

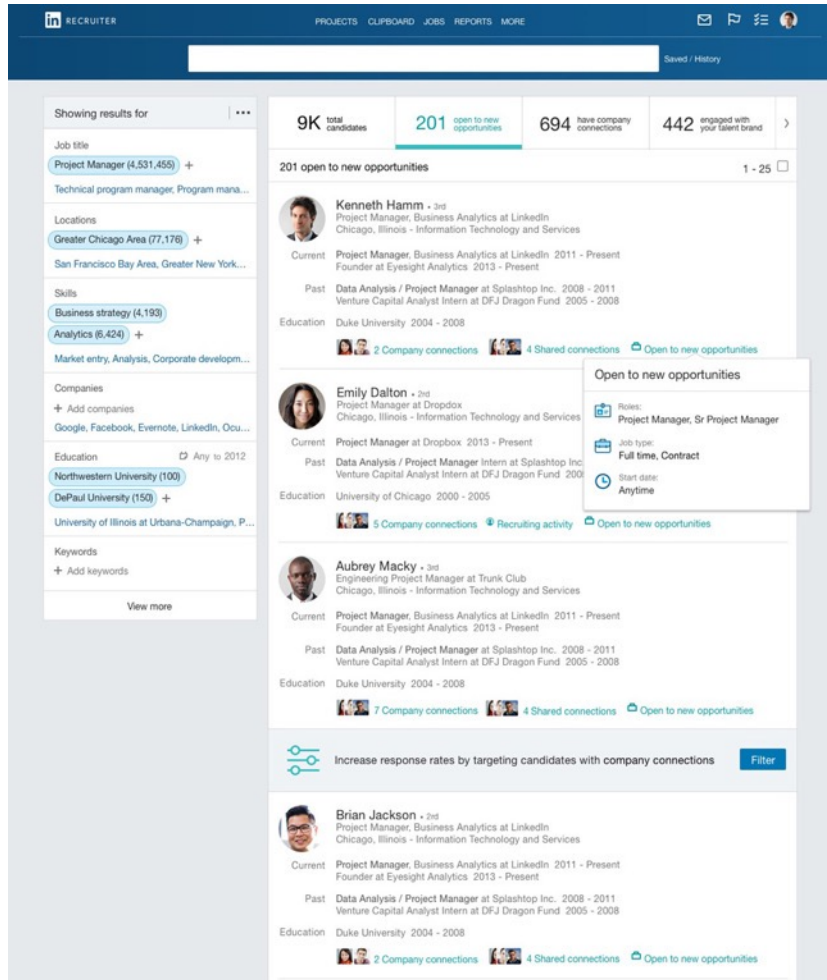


Medial Services (doctor recommendation, patient-doctor matching)

 Legal Services (parole decision recommendation)

# Example: Resume Ranking and Recommendation

## – Explainability for Responsible AI



**Background:** Many companies use **automated tools** such as LinkedIn for recruiting

When a job is posted, could receive thousands of applications -- impossible for HR to manually screen every candidate's resume

**Solution:** Use ML to **rank the candidates** based on some "matching score" between resume and job description.

You only have a chance of interview if the algorithm ranks your resume at top positions (e.g., top-10)

**Problem:**

From recruiter's perspective:

**Why** this candidate is a better fit than another?

From applicant's perspective:

**Why** should I trust the algorithm?

**Why** should my whole career be decided by a machine?

To answer these **WHY** questions, we need Explainable AI!

Figure 1: A (mocked) screenshot from the LinkedIn Recruiter (credit to [1])

# Human-centered Explainable, Fair and Controllable AI

- AI in Human-centered Tasks
  - We not only want to know a model works (e.g., make accurate predictions)
  - We also want to know **why** it works (e.g., why the model makes this decision, is it fair, and why we should trust this decision)
  - Human controls AI, rather than AI controls human
- Even more important in **high-stake applications** related to health, safety, and law



Healthcare



Financial Assistants



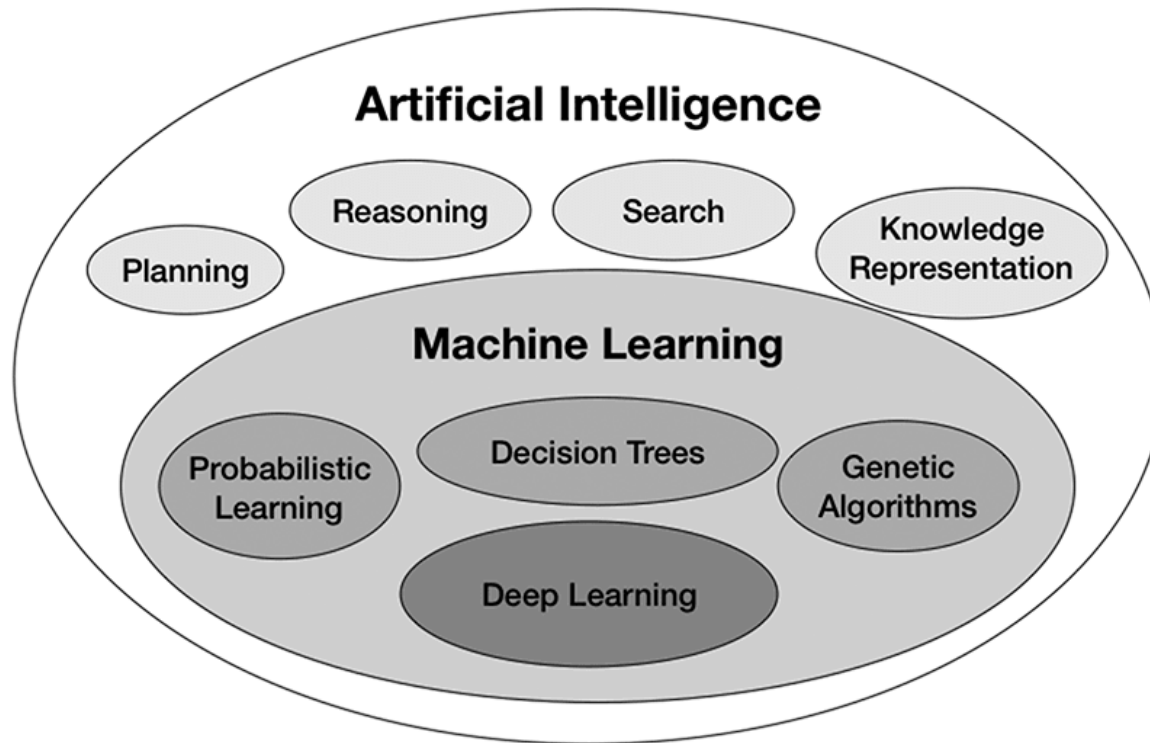
Legal Assistants

- Errors/bias may cause severe loss in life, money, and reputation

- Explainable AI helps humans to make better decisions

# The Scope of AI

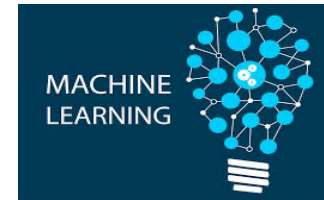
- $AI \neq ML$ ,  $AI \supset ML$



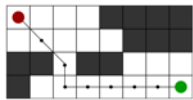
## A (very rough) History of AI Research

- Symbolic Reasoning Approach to AI
  - Mid-1950s to late 1980s
- Machine Learning Approach to AI
  - Early 1990s to date

**GOFAI**  
means  
Good Old Fashioned Artificial Intelligence



### Example Methods:



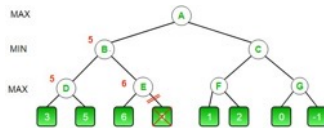
A\* Search



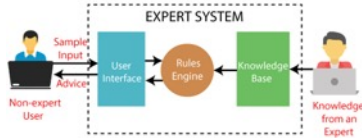
Knowledge Representation and Reasoning



Production Rules



### Example Systems:

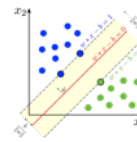


Expert systems  
(If-Then production rules)

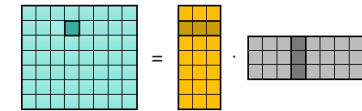


Chess AI  
(IBM Deep Blue)

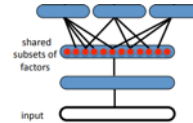
### Example Methods:



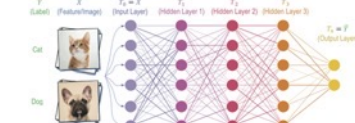
Support Vector Machine



Matrix Factorization



Representation Learning



Deep Neural Networks

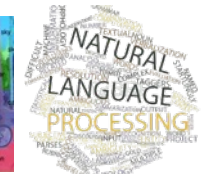
### Example Systems:



Recommender Systems



Image and Language Processing

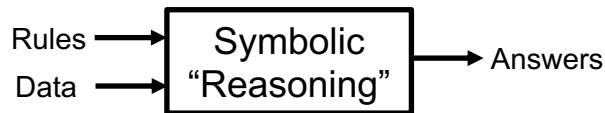


# Symbolism vs Connectionism - A comparison

- a.k.a. Rationalism vs Empiricism approaches to AI

## Symbolism/Rationalism

A top-down design approach



### Advantages:

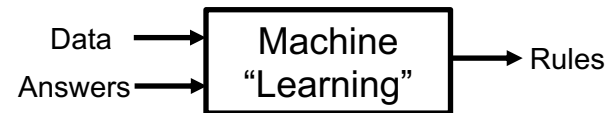
- Accurate decision
- Highly explainable & human readable

### Disadvantages:

- Extensive expert human efforts
- Difficult to handle noisy data

## Connectionism/Empiricism

A bottom-up design approach



### Advantages:

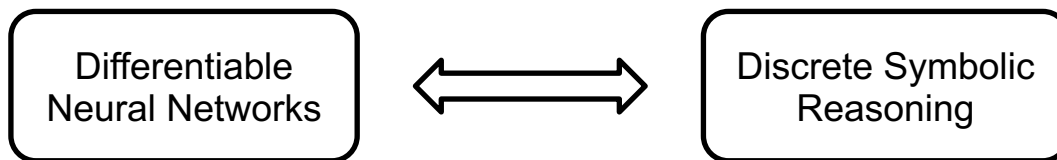
- Less human efforts
- Better at working with noisy data

### Disadvantages:

- Decisions are usually approximate
- Difficult to explain (black-box model)

# Bridge the best of two Worlds?

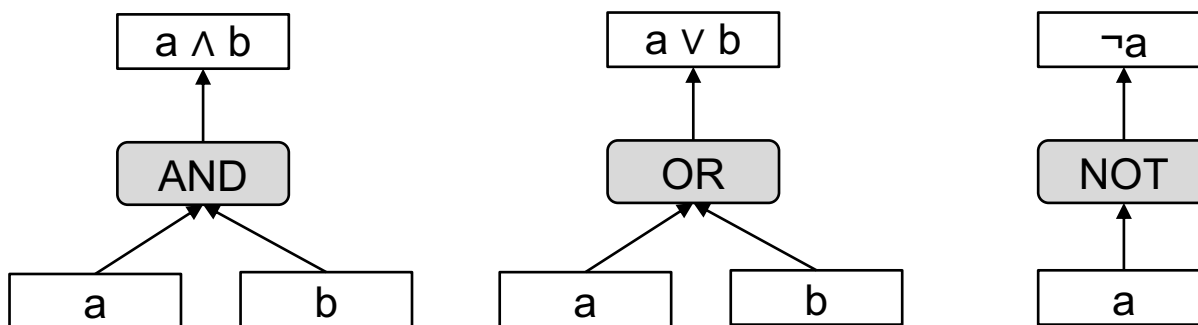
- Neural Symbolic Machine Learning
  - Grant learning systems with reasoning ability
  - Improve decision [accuracy](#)
  - Improve decision [transparency](#)
- Key Challenge
  - How to bridge [differentiable neural networks](#) and [discrete symbolic reasoning](#) in shared architecture for optimization and inference





# Neural Logic Reasoning

- Key idea [4-8]
  - Learning logical variables as vectors in logical embedding space
  - Learning logical operations as neural modules in the latent space



In our implementation,  $\text{AND}(*, *)$ ,  $\text{OR}(*, *)$ ,  $\text{NOT}(*)$  are simple 2-layer neural networks

$$\text{AND}(\mathbf{w}_i, \mathbf{w}_j) = \mathbf{H}_{a2}f(\mathbf{H}_{a1}(\mathbf{w}_i \oplus \mathbf{w}_j) + \mathbf{b}_a) \quad \text{NOT}(\mathbf{w}) = \mathbf{H}_{n2}f(\mathbf{H}_{n1}\mathbf{w} + \mathbf{b}_n)$$

[4] Shaoyun Shi, Hanxiong Chen, Weizhi Ma, Jiaxin Mao, Min Zhang, and Yongfeng Zhang. "Neural Logic Reasoning", CIKM 2020.

[5] Hanxiong Chen, Shaoyun Shi, Yunqi Li and Yongfeng Zhang. "Neural Collaborative Reasoning", WWW 2021.

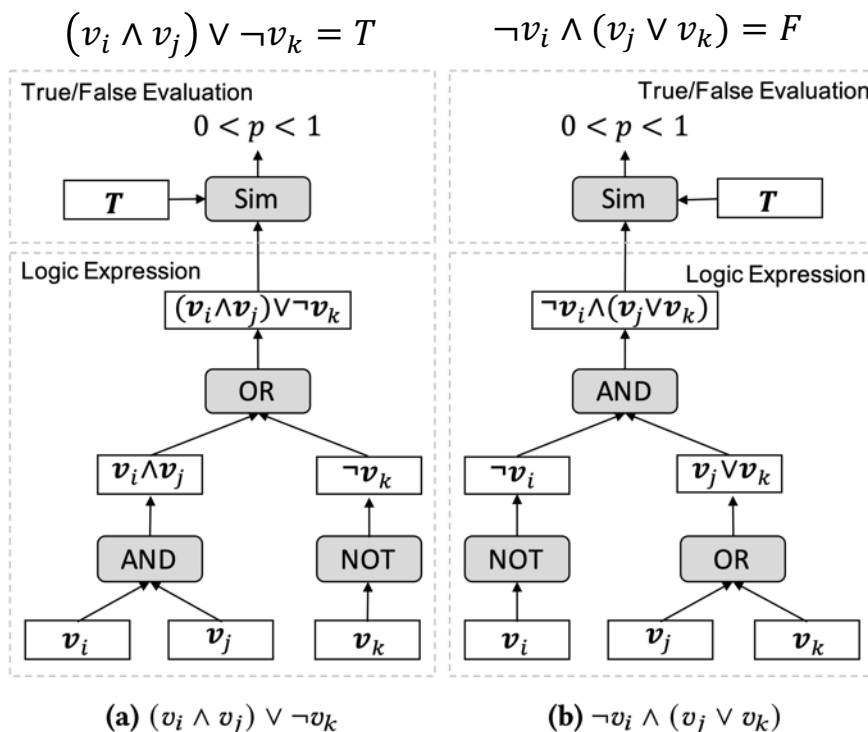
[6] Hanxiong Chen, Yunqi Li, Shaoyun Shi, Shuchang Liu, He Zhu and Yongfeng Zhang. "Graph Collaborative Reasoning", WSDM 2022.

[7] Jianchao Ji, Zelong Li, Shuyuan Xu, Max Xiong, Juntao Tan, Yingqiang Ge, Hao Wang, Yongfeng Zhang. "Counterfactual Collaborative Reasoning", WSDM 2023.

[8] Wenyue Hua and Yongfeng Zhang. "System 1 + System 2 = Better World: Neural-Symbolic Chain of Logic Reasoning", EMNLP 2022.

# Logic-Integrated Neural Network (LINN)

- Any logical expression can be **dynamically assembled** into a neural structure



Optimize with **task dependent loss**, e.g.,:

Cross-Entropy Loss:

$$L_t = L_{ce} = - \sum_{e_i \in E} y_i \log(p_i) + (1 - y_i) \log(1 - p_i)$$

Pair-wise Ranking Loss:

$$L_t = L_{bpr} = - \sum_{e^+} \log(\text{sigmoid}(p(e^+) - p(e^-)))$$

# Logical Regularization over Neural Modules

- How do we know the  $\text{AND}(*, *)$  module is really doing logical AND?
  - And also, for  $\text{OR}(*, *)$  and  $\text{NOT}(*)$ ?
- Logical Regularization
  - Logical operators should satisfy a set of basic requirements

	Logical Rule	Equation	Logic Regularizer $r_i$
NOT	Negation	$\neg T = F$	$r_1 = \sum_{w \in W \cup \{T\}} \text{Sim}(\text{NOT}(w), w)$
	Double Negation	$\neg(\neg w) = w$	$r_2 = \sum_{w \in W} 1 - \text{Sim}(\text{NOT}(\text{NOT}(w)), w)$
AND	Identity	$w \wedge T = w$	$r_3 = \sum_{w \in W} 1 - \text{Sim}(\text{AND}(w, T), w)$
	Annihilator	$w \wedge F = F$	$r_4 = \sum_{w \in W} 1 - \text{Sim}(\text{AND}(w, F), F)$
	Idempotence	$w \wedge w = w$	$r_5 = \sum_{w \in W} 1 - \text{Sim}(\text{AND}(w, w), w)$
	Complementation	$w \wedge \neg w = F$	$r_6 = \sum_{w \in W} 1 - \text{Sim}(\text{AND}(w, \text{NOT}(w)), F)$
OR	Identity	$w \vee F = w$	$r_7 = \sum_{w \in W} 1 - \text{Sim}(\text{OR}(w, F), w)$
	Annihilator	$w \vee T = T$	$r_8 = \sum_{w \in W} 1 - \text{Sim}(\text{OR}(w, T), T)$
	Idempotence	$w \vee w = w$	$r_9 = \sum_{w \in W} 1 - \text{Sim}(\text{OR}(w, w), w)$
	Complementation	$w \vee \neg w = T$	$r_{10} = \sum_{w \in W} 1 - \text{Sim}(\text{OR}(w, \text{NOT}(w)), T)$

- Logical Regularized Loss  $L_1 = L_t + \lambda_l R_l = L_t + \lambda_l \sum_i r_i$

# Application 1: Solving Logical Equations

- 10k logical variables, 30k randomly generated logical equations

- In Disjunctive Normal Form (DNF)

$$(\neg v_{80} \wedge v_{56} \wedge v_{71}) \vee (\neg v_{46} \wedge \neg v_7 \wedge v_{51} \wedge \neg v_{47} \wedge v_{26}) \vee v_{45} \vee (v_{31} \wedge v_{15} \wedge v_2 \wedge v_{46}) = T$$

$$(\neg v_{19} \wedge \neg v_{65}) \vee (v_{65} \wedge \neg v_{24} \wedge v_9 \wedge \neg v_{83}) \vee (\neg v_{48} \wedge \neg v_9 \wedge \neg v_{51} \wedge v_{75}) = F$$

$$\neg v_{98} \vee (\neg v_{76} \wedge v_{66} \wedge v_{13}) \vee v_{97} (\wedge v_{89} \wedge v_{45} \wedge v_{83}) = T$$

$$v_{43} \wedge v_{21} \wedge \neg v_{53} = F$$

- Expressions: training (80%), validation (10%), and test (10%) sets.

- Task: Predict the T/F value for expressions in test sets

$$(v_{65} \wedge \neg v_{24} \wedge v_9 \wedge \neg v_{83}) \vee (\neg v_{48} \wedge \neg v_9 \wedge \neg v_{51} \wedge v_{75}) = ?$$

	$n = 1 \times 10^3, m = 3 \times 10^3$		$n = 1 \times 10^4, m = 3 \times 10^4$	
	Accuracy	RMSE	Accuracy	RMSE
Bi-RNN [32]	$0.6493 \pm 0.0102$	$0.4736 \pm 0.0032$	$0.6942 \pm 0.0028$	$0.4492 \pm 0.0009$
Bi-LSTM [11]	$0.5933 \pm 0.0107$	$0.5181 \pm 0.0162$	$0.6847 \pm 0.0051$	$0.4494 \pm 0.0020$
CNN [19]	$0.6380 \pm 0.0043$	$0.5085 \pm 0.0158$	$0.6787 \pm 0.0025$	$0.4557 \pm 0.0016$
LINN- $R_l$	$0.8353 \pm 0.0043$	$0.3880 \pm 0.0069$	$0.9173 \pm 0.0042$	$0.2733 \pm 0.0065$
LINN	<b><math>0.9440 \pm 0.0064^*</math></b>	<b><math>0.2318 \pm 0.0124^*</math></b>	<b><math>0.9559 \pm 0.0006^*</math></b>	<b><math>0.2081 \pm 0.0018^*</math></b>

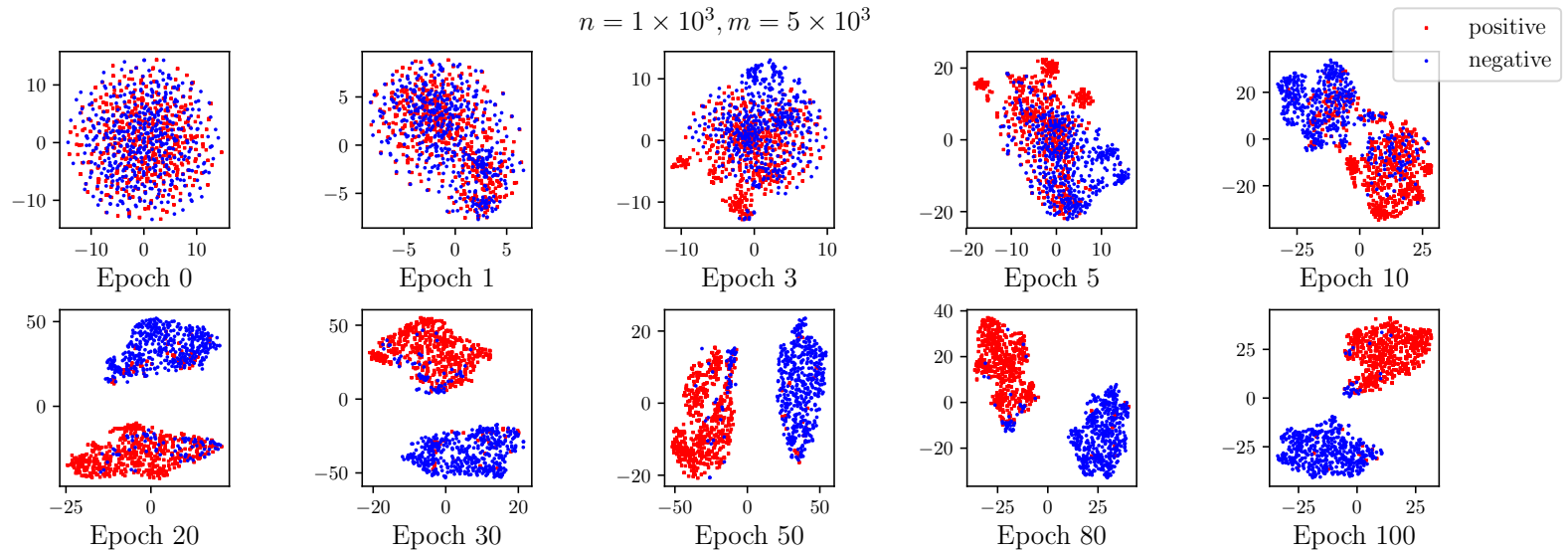
LINN outperforms traditional (non-logical) neural networks.

RNN/LSTM/CNN does not model the compositional logical structure.

Logical regularization is important.

# Application 1: Solving Logical Equations

- t-SNE visualization of logical variable embeddings
  - LINN can finally separate the True and False variables



Accuracy of variable solving: 96%

We can use machine learning to (approximately) solve NP-complete problems

Agnostic to small errors and noise in data.

# Application 2: Explainable Recommendation

## Neural Logic Reasoning for Explainable Recommendation

- Logic expressions help to model item relationships in recommendation
  - **Complimentary**:  $\text{iPhone} \wedge \text{iPhone case} = T$
  - **Substitutive**:  $(\text{Coke} \wedge \neg \text{Pepsi}) \vee (\neg \text{Coke} \wedge \text{Pepsi}) = T$
  - **Irrelevant**:  $\text{iPhone} \wedge \text{Android data line} = F$ .
- User's interaction history can be represented as logical expressions
  - Suppose user purchased item  $v_3$  after several history interactions  $\{v_1 = T \text{ (likes), } v_2 = F \text{ (dislikes)}\}$
  - Training example:  $(v_1 \wedge v_3) \vee (\neg v_2 \wedge v_3) \vee (v_1 \wedge \neg v_2 \wedge v_3) = T$
  - This is a **noisy** reasoning problem: different users' equation may conflict

- **Pair-wise Contrastive Ranking Loss**

$$e^+ = (\cdot \wedge v^+) \vee \dots \vee (\cdot \wedge v^+)$$

$$e^- = (\cdot \wedge v^-) \vee \dots \vee (\cdot \wedge v^-)$$

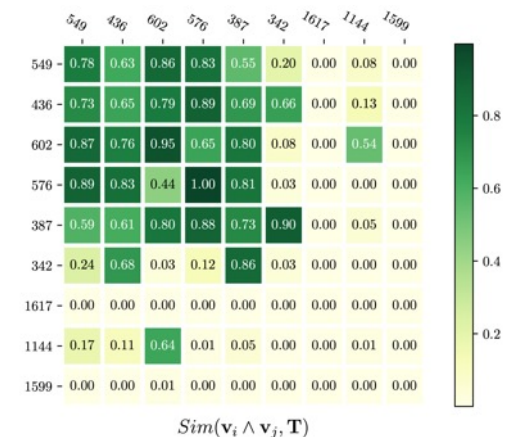
$$L = - \sum_{e^+} \log(\text{sigmoid}(p(e^+) - p(e^-))) + \lambda_l \sum_i r_i + \lambda_\ell \sum_{w \in W} \|w\|_F^2 + \lambda_\Theta \|\Theta\|_F^2$$

# Application 2: Explainable Recommendation

- Recommendation Performance
  - LINN makes significant improvements on Movie and E-commerce recommendation

	ML-100k		Amazon Electronics	
	nDCG@10	Hit@1	nDCG@10	Hit@1
BPRMF [31]	0.3664 ± 0.0017	0.1537 ± 0.0036	0.3514 ± 0.0002	0.1951 ± 0.0004
SVD++ [21]	0.3675 ± 0.0024	0.1556 ± 0.0044	0.3582 ± 0.0004	0.1930 ± 0.0006
STAMP [25]	0.3943 ± 0.0016	0.1706 ± 0.0022	0.3954 ± 0.0003	0.2215 ± 0.0003
GRU4Rec [16]	0.3973 ± 0.0016	0.1745 ± 0.0038	0.4029 ± 0.0009	0.2262 ± 0.0009
NARM [24]	0.4022 ± 0.0015	0.1771 ± 0.0016	0.4051 ± 0.0006	0.2292 ± 0.0005
LINN- $R_I$	0.4022 ± 0.0027	0.1783 ± 0.0043	0.4152 ± 0.0014	0.2396 ± 0.0019
<b>LINN</b>	<b>0.4064 ± 0.0015*</b>	<b>0.1850 ± 0.0053*</b>	<b>0.4191 ± 0.0012*</b>	<b>0.2438 ± 0.0014*</b>

- Extracting Explanations for the Recommendations
  - The AND module extracts **complimentary item explanations**
  - E.g., iPhone  $\wedge$  iPhone case =  $T$
  - **Explanation:** We recommend this **iPhone case** is because you have purchased an **iPhone**.



# Neural Collaborative Reasoning

- **Personalize** the Reasoning Process

- Reasoning with Implicit Feedback

- User  $u$ , items  $\{v_1, v_2, \dots, v_r\}$

Horn Clause:  $I(u, v_1) \wedge I(u, v_2) \wedge \dots \wedge I(u, v_r) \rightarrow I(u, v_x)$

- $I(u, v_i)$  is an encoding function showing user  $u$  interacted with an item
- $I(u, v_i)$  can be a simple neural network

- Reasoning with Explicit Feedback

- User  $u$ , items  $\{v_1, v_2, \dots, \neg v_r\}$ , where  $\neg v_r$  represents a user has negative feedback

Horn Clause:  $L(u, v_1) \wedge L(u, v_2) \wedge \dots \wedge \neg L(u, v_r) \rightarrow L(u, v_x)$

- $L(u, v_i)$  is an encoding function showing user likes an item

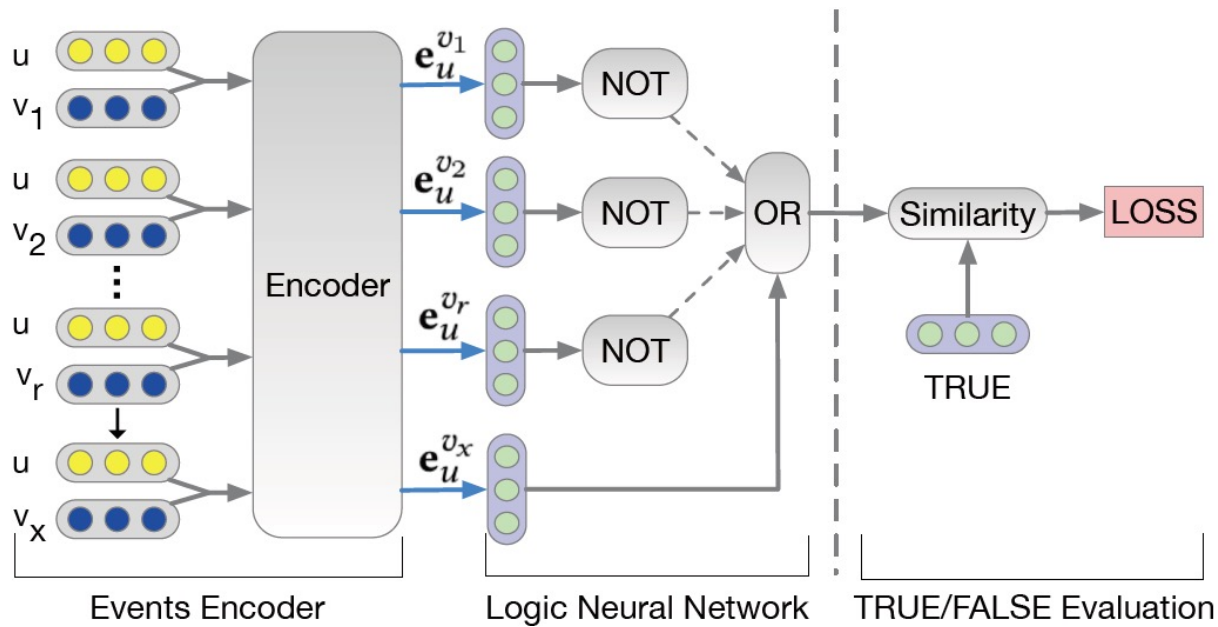


# Collaborative Reasoning Architecture

Horn Clause:  $I(u, v_1) \wedge I(u, v_2) \wedge \dots \wedge I(u, v_r) \rightarrow I(u, v_x)$

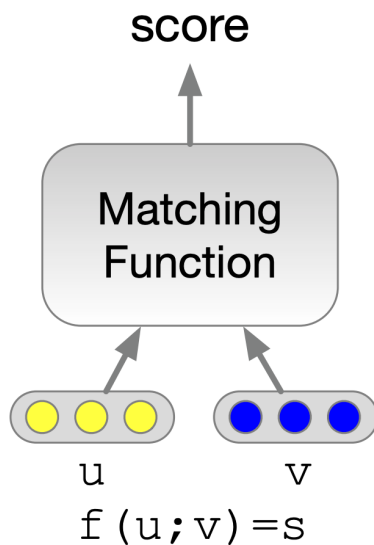
$$e_u^{v_1} \wedge e_u^{v_2} \dots \wedge e_u^{v_r} \rightarrow e_u^{v_x} \Leftrightarrow \neg e_u^{v_1} \vee \neg e_u^{v_2} \dots \vee \neg e_u^{v_r} \vee e_u^{v_x}$$

$$p \rightarrow q \Leftrightarrow \neg p \vee q$$

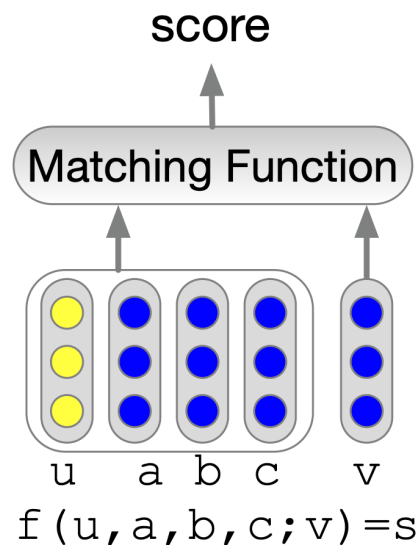


# From Learning to Reasoning for AI

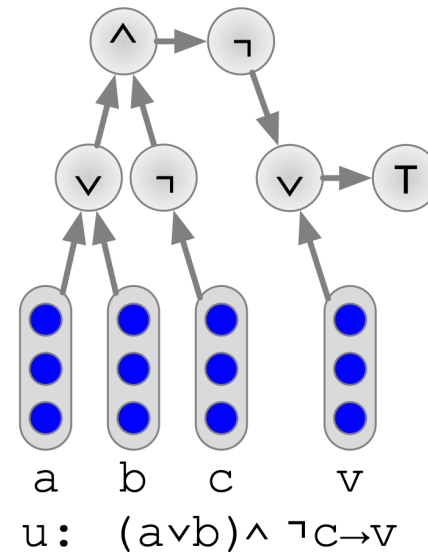
- From Perception to Cognition
- From System 1 to System 2



Matching Models



Sequential Models

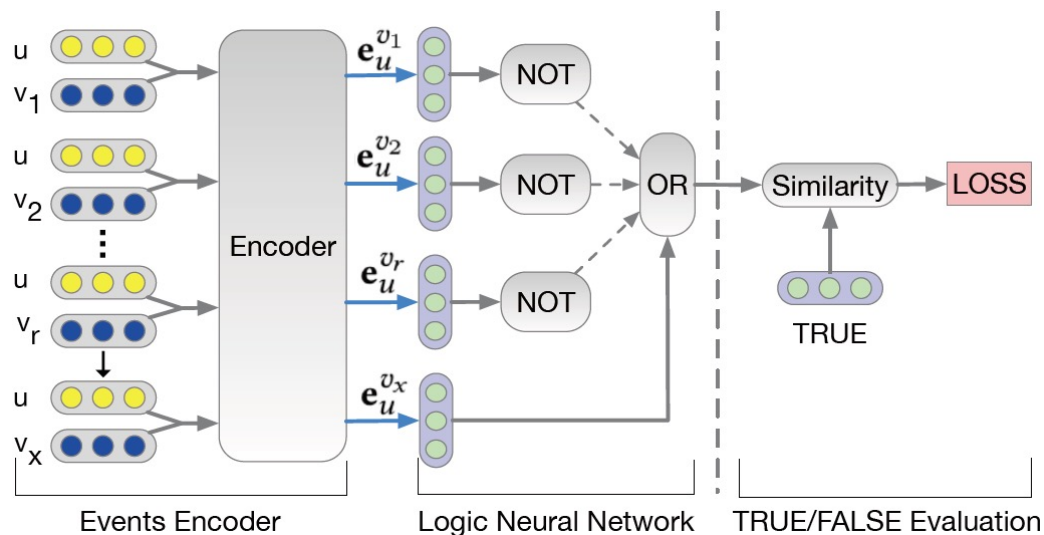


Reasoning Models

# From Learning to Reasoning

- From System 1 to System 2 for AI

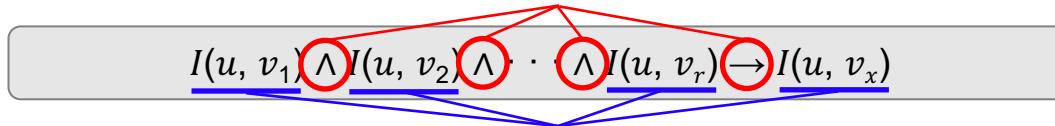
$$I(u, v_1) \wedge I(u, v_2) \wedge \dots \wedge I(u, v_r) \rightarrow I(u, v_x)$$



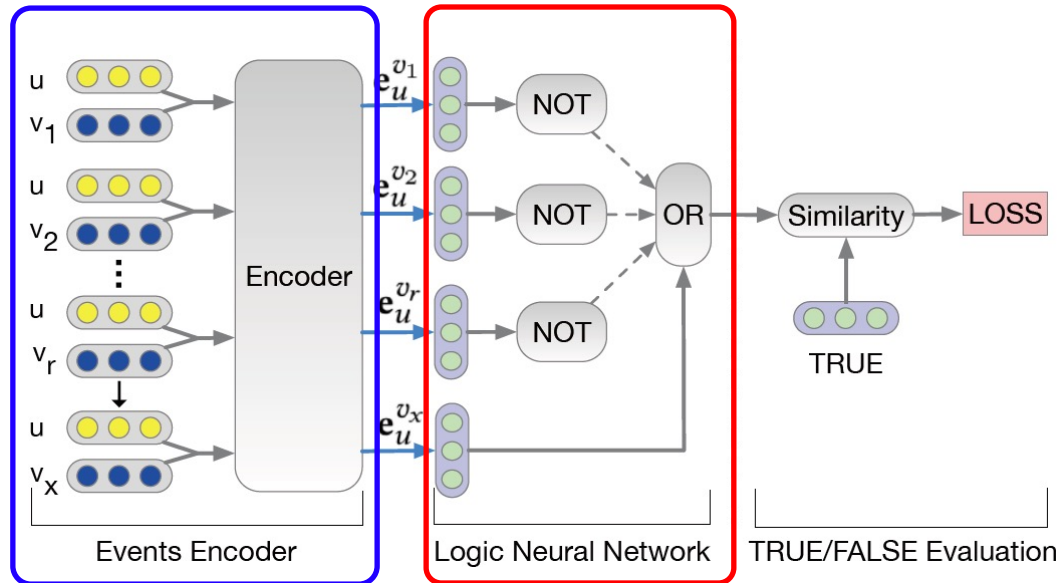
# From Learning to Reasoning

- From System 1 to System 2 for AI

System 2: Cognitive Reasoning



System 1: Perceptive Learning



System 1: Perceptive Learning    System 2: Cognitive Reasoning

# Results

	ML100k				Movies and TV				Electronics			
	N@5	N@10	HR@5	HR@10	N@5	N@10	HR@5	HR@10	N@5	N@10	HR@5	HR@10
BPR-MF	0.3024	0.3659	0.4501	0.6486	0.3962	0.4392	0.5346	0.6676	0.3092	0.3472	0.4179	0.5354
SVD++	0.3087	0.3685	0.4586	0.6433	0.3918	0.4335	0.5224	0.6512	0.2775	0.3172	0.3848	0.5077
DMF	0.3023	0.3661	0.4480	0.6450	0.4006	0.4455	<u>0.5455</u>	<u>0.6843</u>	0.2775	0.3143	0.3783	0.4922
NeuMF	0.3002	0.3592	0.4490	0.6316	0.3791	0.4211	0.5134	0.6429	0.3026	0.3358	0.4031	0.5123
GRU4Rec	<u>0.3564</u>	<u>0.4122</u>	0.5134	<u>0.6856</u>	<u>0.4038</u>	<u>0.4459</u>	0.5287	0.6688	<u>0.3154</u>	<u>0.3551</u>	<u>0.4284</u>	<u>0.5511</u>
STAMP	0.3560	0.4070	<u>0.5159</u>	0.6730	0.3935	0.4366	0.5246	0.6577	0.3095	0.3489	0.4196	0.5430
NLR	<u>0.3602</u>	<u>0.4151</u>	0.5102	0.6795	<u>0.4191</u>	<u>0.4591</u>	<u>0.5506</u>	0.6739	<u>0.3475</u>	<u>0.3852</u>	<u>0.4623</u>	<u>0.5788</u>
<b>NCR-I</b>	0.3697	0.4219	0.5265	0.6890	0.4152	0.4550	0.5479	0.6709	0.3226	0.3604	0.4331	0.5500
NCR-E w/o LR	0.3671	0.4219	0.5180	0.6890	0.4126	0.4535	0.5444	0.6705	0.3272	0.3649	0.4377	0.5544
<b>NCR-E</b>	<b>0.3760**</b>	<b>0.4240**</b>	<b>0.5456**</b>	<b>0.6943**</b>	<b>0.4255**</b>	<b>0.4670**</b>	<b>0.5611**</b>	<b>0.6891</b>	<b>0.3499*</b>	<b>0.3878*</b>	<b>0.4639*</b>	<b>0.5812*</b>
Improvment <sup>1</sup>	5.50%	2.86%	5.76%	1.27%	5.37%	4.73%	2.86%	0.70%	10.94%	9.21%	8.29%	5.46%
Improvment <sup>2</sup>	4.39%	2.14%	6.71%	2.66%	1.53%	1.72%	1.91%	2.26%	0.69%	0.67%	0.35%	0.41%

- NCR-I: Reasoning with Implicit Feedback
- NCR-E: Reasoning with Explicit Feedback
- Model is Partly Explainable

# The Importance of Causal-Consistent Reasoning

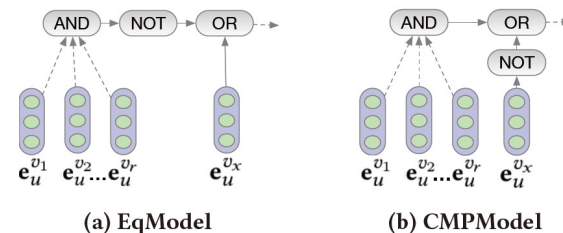
- EqModel (causally consistent):

$$- e_u^{v_1} \wedge e_u^{v_2} \cdots \wedge e_u^{v_r} \rightarrow e_u^{v_x} \Leftrightarrow \neg e_u^{v_1} \vee \neg e_u^{v_2} \cdots \vee \neg e_u^{v_r} \vee e_u^{v_x} \quad (1)$$

$$- e_u^{v_1} \wedge e_u^{v_2} \cdots \wedge e_u^{v_r} \rightarrow e_u^{v_x} \Leftrightarrow \neg (e_u^{v_1} \wedge e_u^{v_2} \cdots \wedge e_u^{v_r}) \vee e_u^{v_x} \quad (2)$$

- CMPModel (causally inconsistent):

$$- e_u^{v_x} \rightarrow e_u^{v_1} \wedge e_u^{v_2} \cdots \wedge e_u^{v_r} \Leftrightarrow \neg e_u^{v_x} \vee (e_u^{v_1} \wedge e_u^{v_2} \cdots \wedge e_u^{v_r}) \quad (3)$$



	ML100k				Movies and TV				Electronics			
	N@5	N@10	HR@5	HR@10	N@5	N@10	HR@5	HR@10	N@5	N@10	HR@5	HR@10
GRU4Rec	0.3564	0.4122	0.5134	0.6856	0.4038	0.4459	0.5287	0.6688	0.3154	0.3551	0.4284	0.5511
NLR	0.3529	0.4066	0.5113	0.6763	0.4191	0.4591	0.5506	0.6739	0.3475	0.3852	0.4623	0.5788
<sup>1</sup> EqModel	0.3664	0.4224	0.5318	<b>0.7070</b>	0.4105	0.4521	0.5429	0.6686	0.3249	0.3626	0.4355	0.5518
<sup>2</sup> CMPModel	0.3551	0.4144	0.5106	0.6932	0.4100	0.4506	0.5417	0.6670	0.3165	0.3541	0.4252	0.5416
<sup>3</sup> NCR-E	<b>0.3760</b>	<b>0.4240</b>	<b>0.5456</b>	0.6943	<b>0.4255</b>	<b>0.4670</b>	<b>0.5611</b>	<b>0.6891</b>	<b>0.3499</b>	<b>0.3878</b>	<b>0.4639</b>	<b>0.5812</b>
<i>p</i> -value <sup>1,3</sup>	0.0825	0.0606	0.1073	0.0547	0.0156*	0.0230*	0.0212*	0.0197*	0.0015*	0.0021*	0.0010*	0.0009*
<i>p</i> -value <sup>2,3</sup>	0.0099*	0.0250*	0.0258*	0.4668	0.0108*	0.0103*	0.0057*	0.0048*	0.0022*	0.0019*	0.0023*	0.0018*

Causally consistent models are comparable

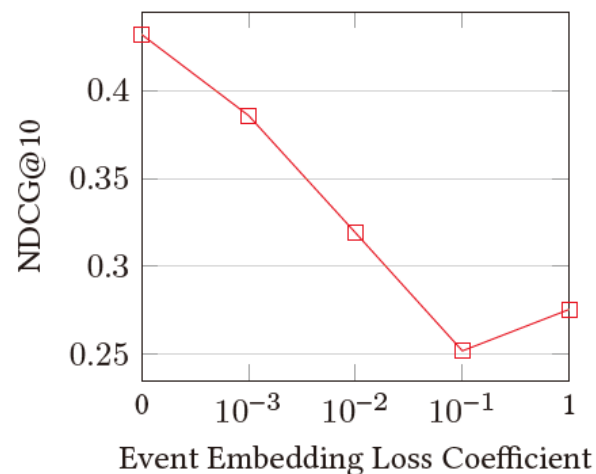
Causally consistent models are better than causally inconsistent models

# The Importance of Neural-Symbolic Reasoning (compared with Pure-Symbolic Reasoning)

- Boolean logic constraint:

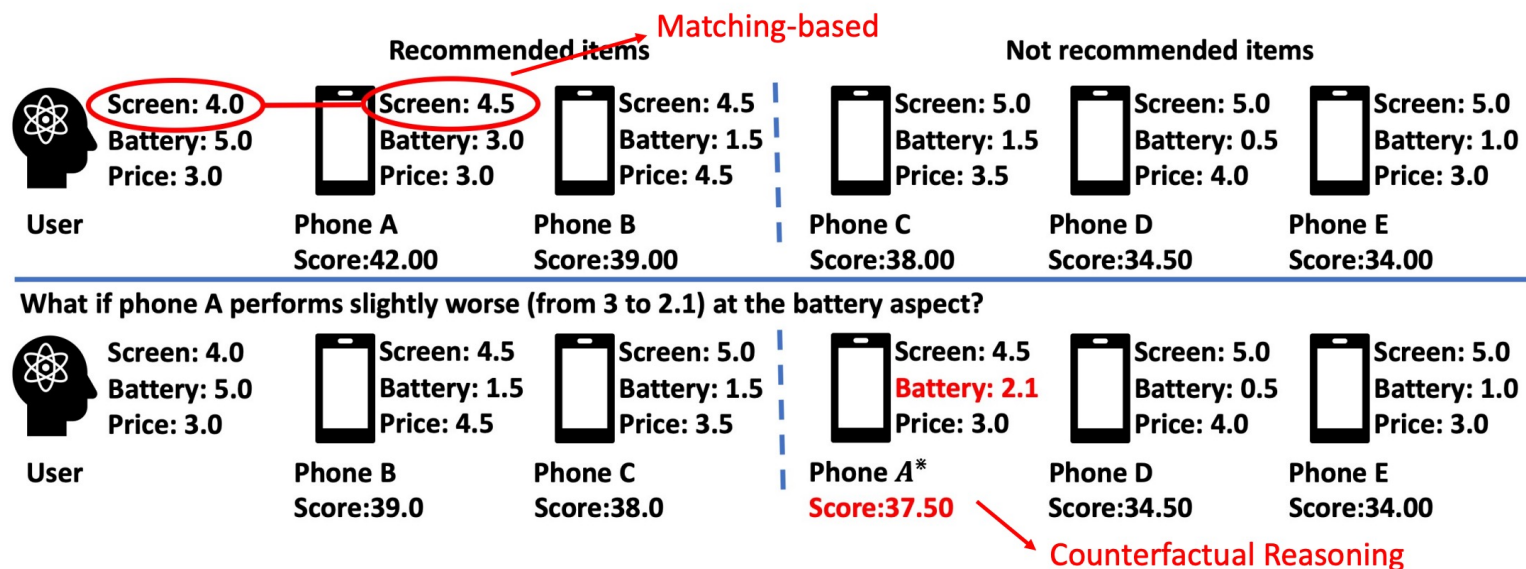
$$\mathcal{L}_{event} = \sum_u \sum_{v \in \mathcal{V}_u^+} \text{MSE}(\mathbf{e}_u^v, \mathbf{G})$$

- $\mathbf{G}$  is for ground-truth vector, which is either **T** or **F**;
  - $\text{MSE}()$  is mean square error.
- Neural-Symbolic Reasoning is better than Pure Boolean Logic Reasoning
    - We leverage both Learning and Reasoning abilities



# Counterfactual Explanations

- Associative vs. Causal/Counterfactual Reasoning



## Counterfactual Explanation:

If the item had been slightly worse on [aspect(s)], then it would not have been recommended.



# Simple and Effective Explanations

- Occam's Razor Principle

- If two explanations are equally **effective** in explaining the results, we prefer the **simpler** explanation than the complex one.

- To character Simpleness

- Explanation Complexity  $C(\Delta) = \gamma \|\Delta\|_0 + \|\Delta\|_2^2$

How many aspects are used  
to generate explanations

How many changes need to  
be applied on these aspects

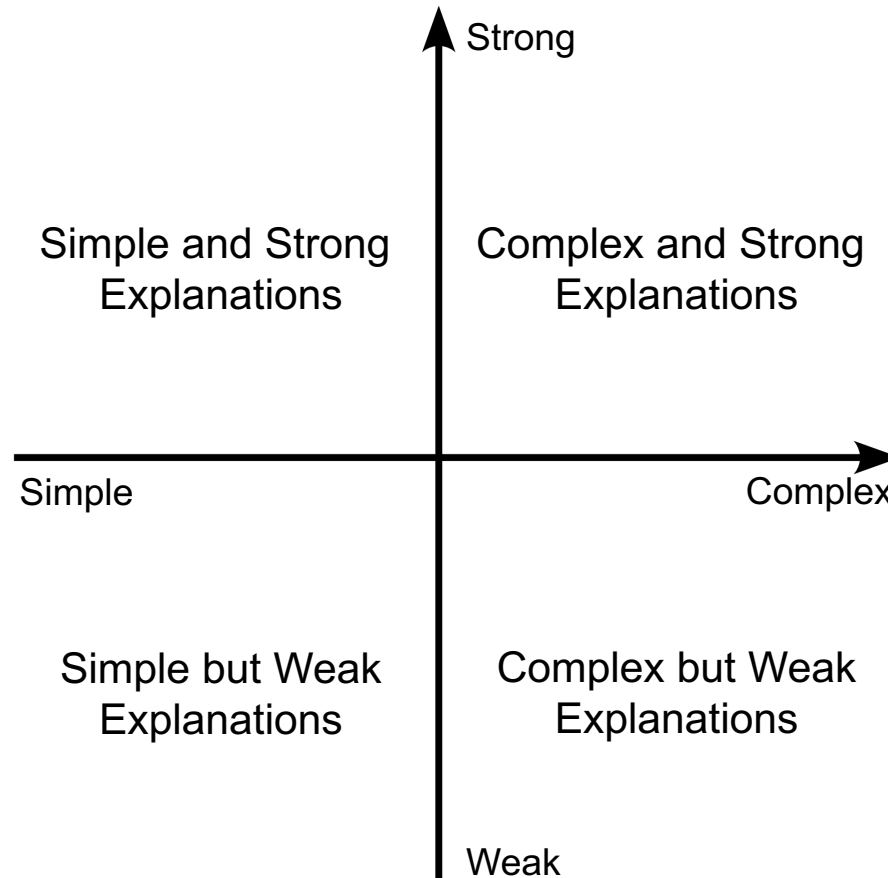
- To character Effectiveness

- Explanation Strength  $S(\Delta) = \underline{s_{i,j}} - s_{i,j\Delta}$

The decrease of  $V_j$ 's ranking score in  
user  $U_i$ 's recommendation list after  
applying  $\Delta$

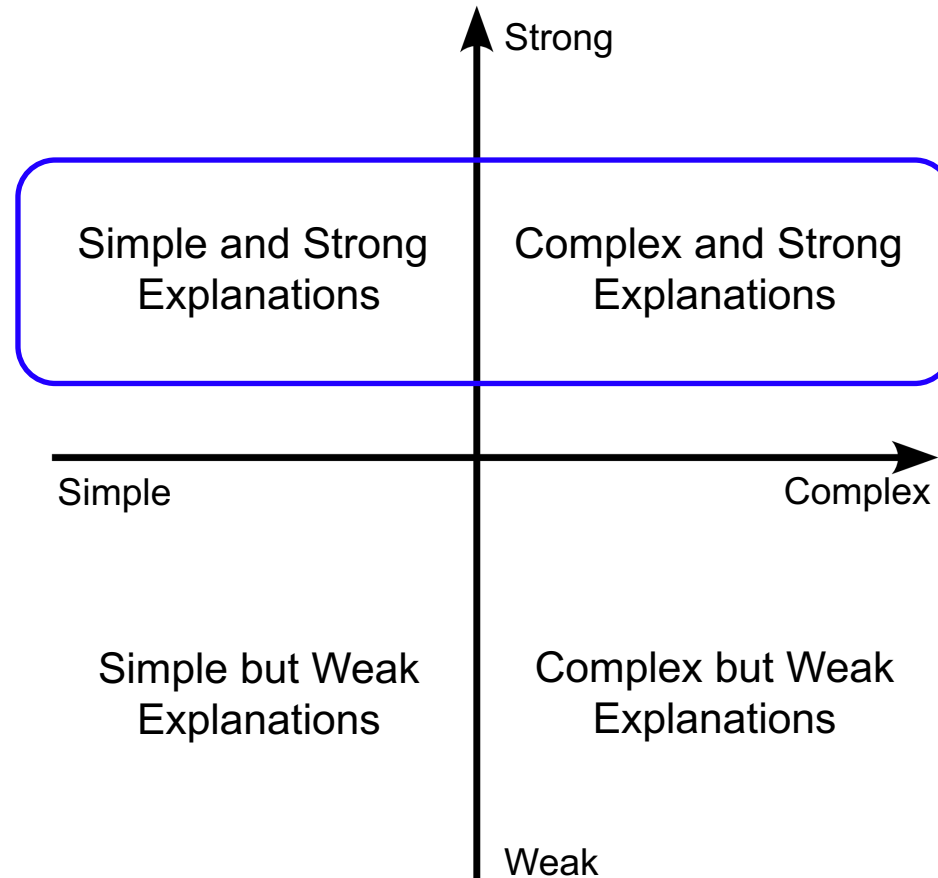
# Complexity vs. Strength

- Two orthogonal concepts



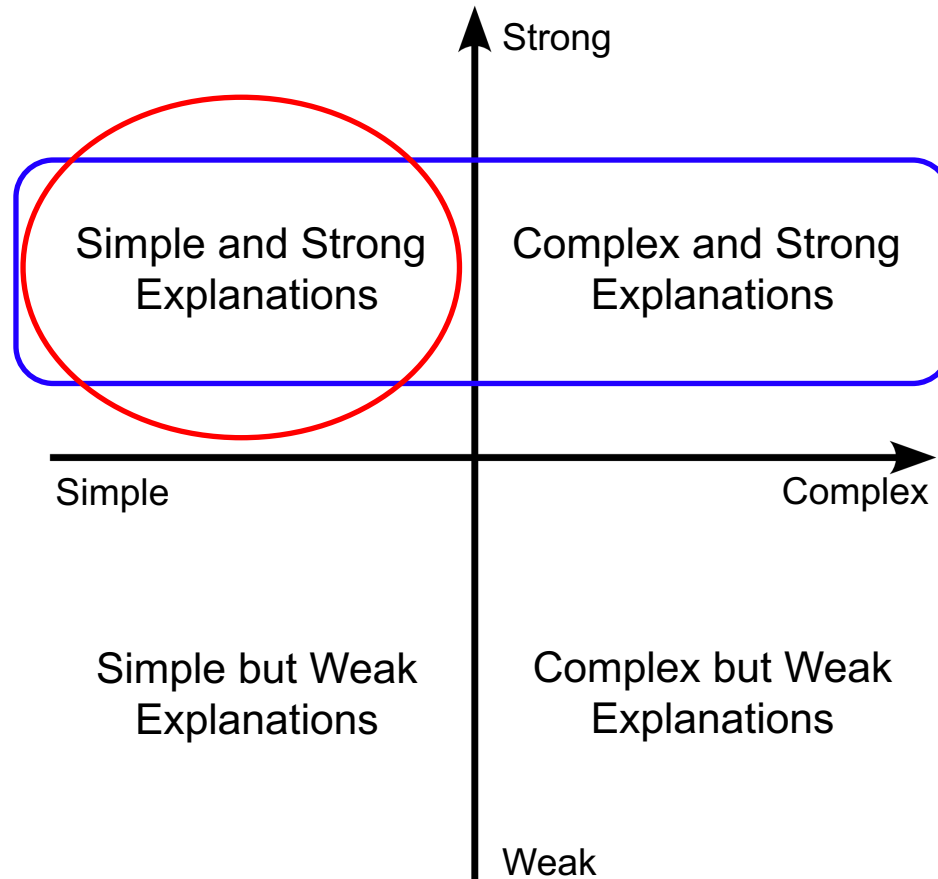
# Complexity vs. Strength

- Two orthogonal concepts



# Complexity vs. Strength

- Two orthogonal concepts



# Counterfactual Learning and Reasoning

- Seek **simple** and **effective** explanations

minimize Explanation Complexity  
 s.t., Explanation is Strong Enough



minimize  $C(\Delta) = \|\Delta\|_2^2 + \gamma\|\Delta\|_0$   
 s.t.,  $S(\Delta) = s_{i,j} - s_{i,j_\Delta} \geq \epsilon$

- Idea: Find **minimal changes** to an item's features so that the item can be **kicked out** of the recommendation list

- Related Optimization for model learning

$$\underset{\Delta}{\text{minimize}} \|\Delta\|_2^2 + \gamma\|\Delta\|_1 + \lambda L(s_{i,j_\Delta}, s_{i,j_{K+1}})$$

where  $s_{i,j_\Delta} = f(X_i, Y_j + \Delta \mid Z, \Theta)$ ,  $s_{i,j_{K+1}} = f(X_i, Y_{j_{K+1}} \mid Z, \Theta)$

$$L(s_{i,j_\Delta}, s_{i,j_{K+1}}) = \max(0, \alpha + s_{i,j_\Delta} - s_{i,j_{K+1}})$$

# Sufficiency and Necessity of Explanations

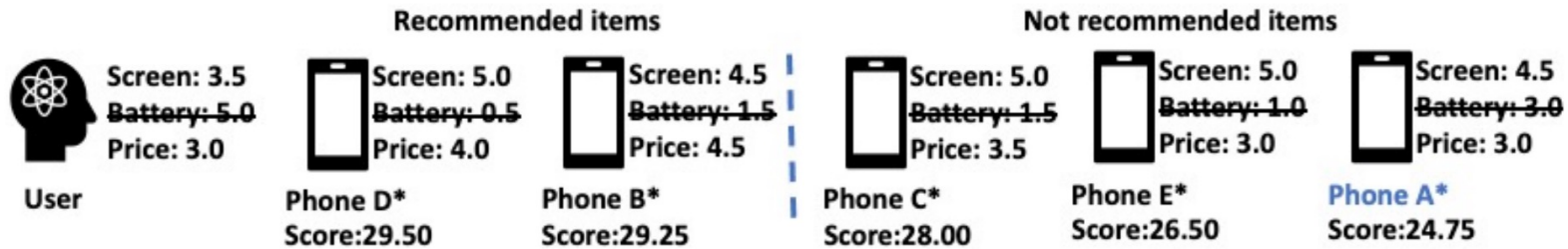
- $S \Rightarrow N$ : S is a sufficient condition for N
- $\neg N \Rightarrow \neg S$ : N is a necessary condition for S

# Sufficiency and Necessity of Explanations

- $S \Rightarrow N$ : S is a sufficient condition for N
- $\neg N \Rightarrow \neg S$ : N is a necessary condition for S

CountER: If this phone had been slightly worse on [Battery], then it will not be recommended.

- Probability of Necessity (PN): If in a counterfactual world, the aspects in the explanation did not exist in the system, then what is the probability that the item would not be recommended.



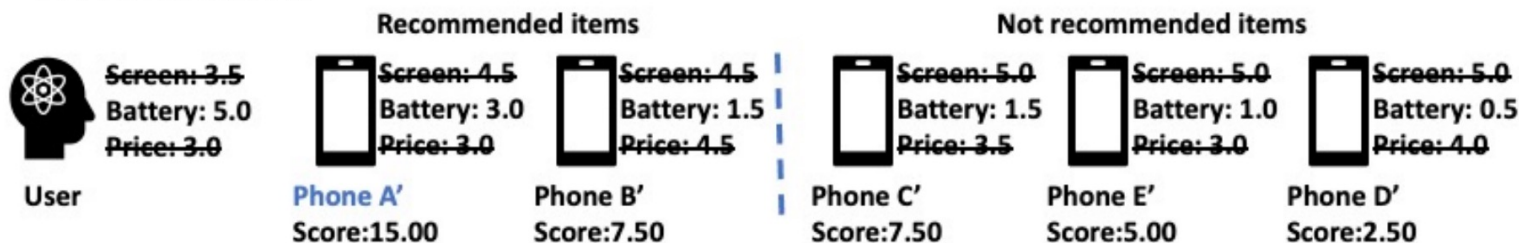
$$PN = \frac{\sum_{u_i \in \mathcal{U}} \sum_{v_j \in R_{i,K}} PN_{ij}}{\sum_{u_i \in \mathcal{U}} \sum_{v_j \in R_{i,K}} I(\mathcal{A}_{ij} \neq \emptyset)}, \text{ where } PN_{ij} = \begin{cases} 1, & \text{if } v_j^* \notin R_{i,K}^* \\ 0, & \text{else} \end{cases}$$

# Sufficiency and Necessity of Explanations

- $S \Rightarrow N$ : S is a sufficient condition for N
- $\neg N \Rightarrow \neg S$ : N is a necessary condition for S

Counter: If this phone had been slightly worse on [Battery], then it will not be recommended.

- Probability of Sufficiency (PS): If in a counterfactual world, the aspects in the explanation were the only aspects existed in the system, then what is the probability that the item would still be recommended.



$$PS = \frac{\sum_{u_i \in \mathcal{U}} \sum_{v_j \in R_{i,K}} PS_{ij}}{\sum_{u_i \in \mathcal{U}} \sum_{v_j \in R_{i,K}} I(\mathcal{A}_{ij} \neq \emptyset)}, \text{ where } PS_{ij} = \begin{cases} 1, & \text{if } v'_j \in R'_{i,K} \\ 0, & \text{else} \end{cases}$$

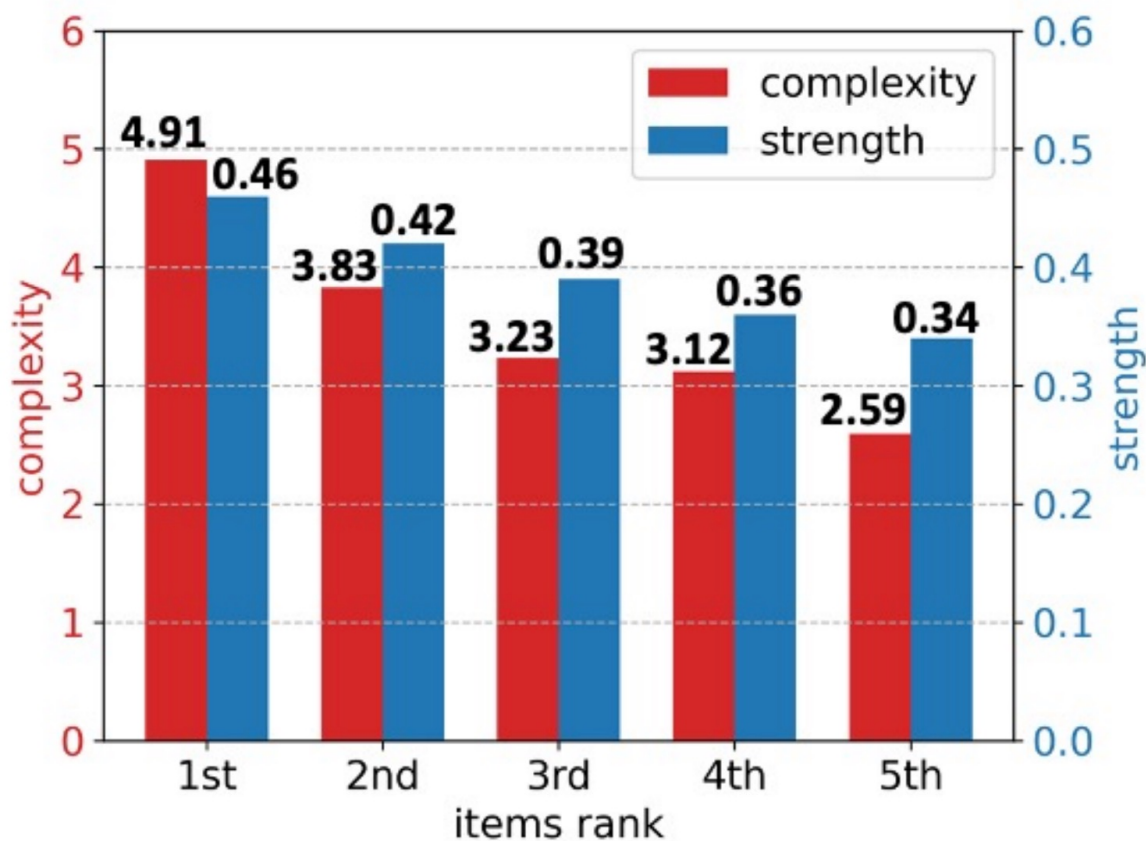


# Counterfactual Reasoning gives Better Explanations

	Single Aspect Explanation														
	Electronic			Cell Phones			Kindle Store			CDs and Vinyl			Yelp		
	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$
Random	2.05	2.10	2.07	3.39	3.50	3.44	3.16	2.75	2.94	1.58	2.03	1.78	7.52	10.68	8.82
EFM[50]	8.41	41.13	13.96	32.31	<b>82.09</b>	46.37	6.01	<b>73.84</b>	11.12	10.15	42.63	16.39	5.87	61.06	10.71
A2CF[9]	41.45	<b>77.60</b>	54.03	36.82	78.68	50.17	25.66	65.53	36.88	25.41	<b>84.51</b>	39.07	17.59	<b>96.92</b>	29.78
CountER	<b>65.54</b>	68.28	<b>66.83</b>	<b>74.03</b>	63.30	<b>68.25</b>	34.37	41.50	37.60	49.62	54.72	52.04	<b>65.26</b>	53.25	<b>58.64</b>
CountER (w/ mask)	56.73	62.03	59.26	70.11	54.71	61.46	<b>35.39</b>	46.91	<b>40.34</b>	<b>75.17</b>	49.18	<b>59.46</b>	58.52	52.56	55.38
	Multiple Aspect Explanation														
	Electronic			Cell Phones			Kindle Store			CDs and Vinyl			Yelp		
	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$	PN%	PS%	$F_{NS}\%$
Random	2.24	4.90	3.08	6.25	10.13	7.73	5.80	7.80	6.65	3.22	7.65	4.53	13.84	12.92	13.36
EFM[50]	29.65	84.67	43.92	52.66	87.98	65.88	51.72	<b>96.42</b>	67.33	47.65	87.35	61.66	16.76	81.68	27.81
A2CF[9]	59.47	81.66	68.82	56.45	80.97	66.52	52.48	87.59	65.64	49.12	<b>91.52</b>	63.93	41.38	<b>98.28</b>	58.24
CountER	<b>97.08</b>	<b>96.24</b>	<b>96.66</b>	<b>99.52</b>	<b>98.48</b>	<b>99.00</b>	<b>64.00</b>	79.20	<b>70.79</b>	<b>80.89</b>	88.60	<b>84.57</b>	<b>99.91</b>	94.12	<b>96.93</b>
CountER (w/ mask)	77.96	89.26	83.23	86.62	91.78	89.13	60.70	80.10	69.06	72.47	67.72	70.01	96.73	94.39	95.55

## Interesting Observations

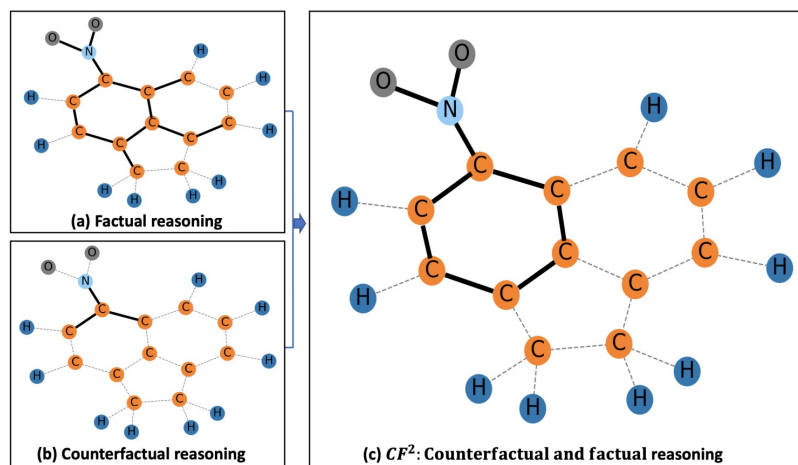
- Top-ranked items need to be backed by stronger and more complex explanations



# PN & PS based Evaluation is Usable

- PN/PS metrics are highly correlated with ground-truth based metrics

$$F_{NS} = \frac{2 \cdot \text{PN} \cdot \text{PS}}{\text{PN} + \text{PS}}$$



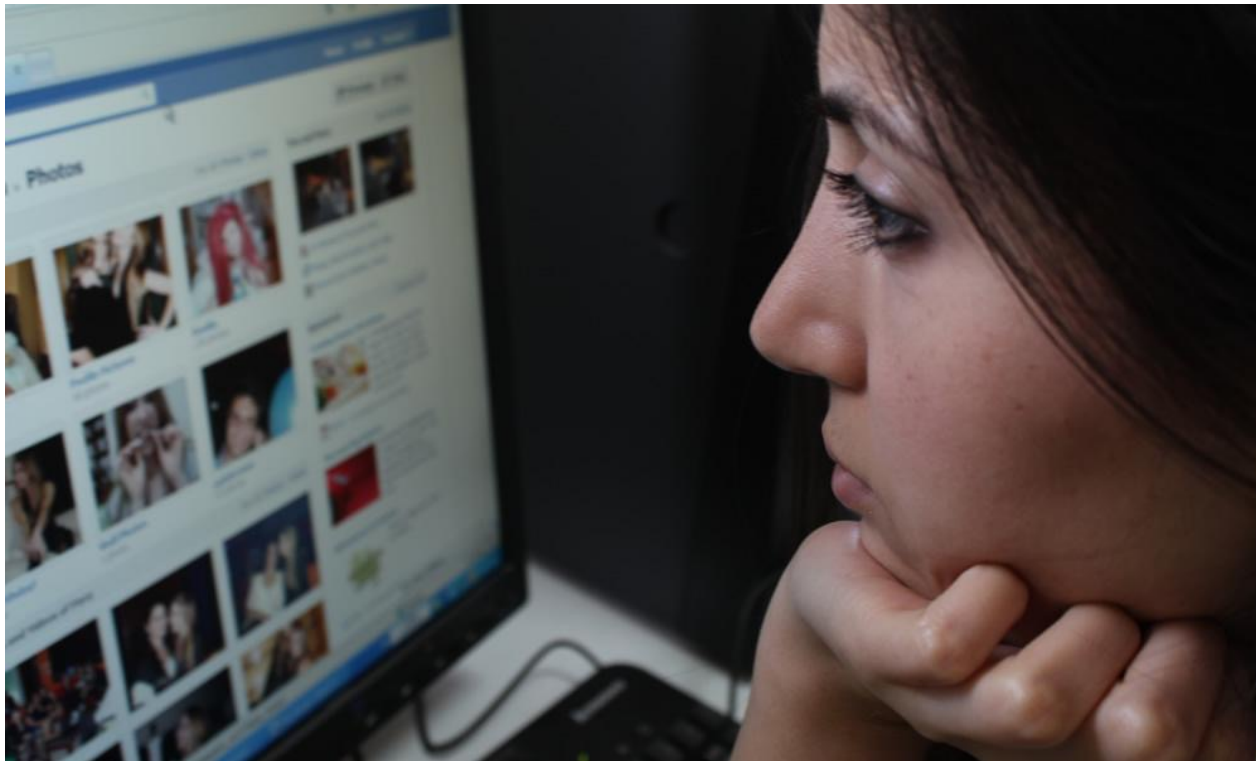
Kendall's  $\tau$  and Spearman's  $\rho$  correlation

**Table 7: Correlation between PN/PS-based evaluation and ground-truth evaluation.**

Models	BA-Shapes		Tree-Cycles		Mutag <sub>0</sub>	
	$\tau \uparrow$	$\rho \uparrow$	$\tau \uparrow$	$\rho \uparrow$	$\tau \uparrow$	$\rho \uparrow$
$F_{NS}$ & $F_1$	1.00	1.00	1.00	1.00	1.00	1.00
$F_{NS}$ & Acc	0.66	0.79	1.00	1.00	0.66	0.79

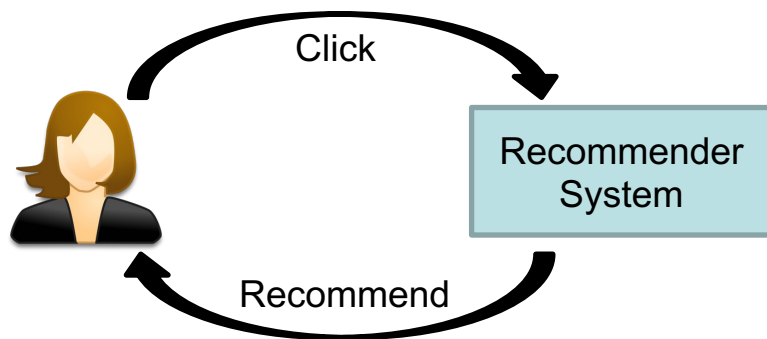
# Towards User Controllable Recommender Systems

- Users almost have **no control** of their recommender system
  - They can only **passively** receive recommendations



# Towards User Controllable Recommender Systems

- Users almost have **no control** of their recommender system
  - They can only **passively** receive recommendations
- This causes many problems, e.g., echo chamber

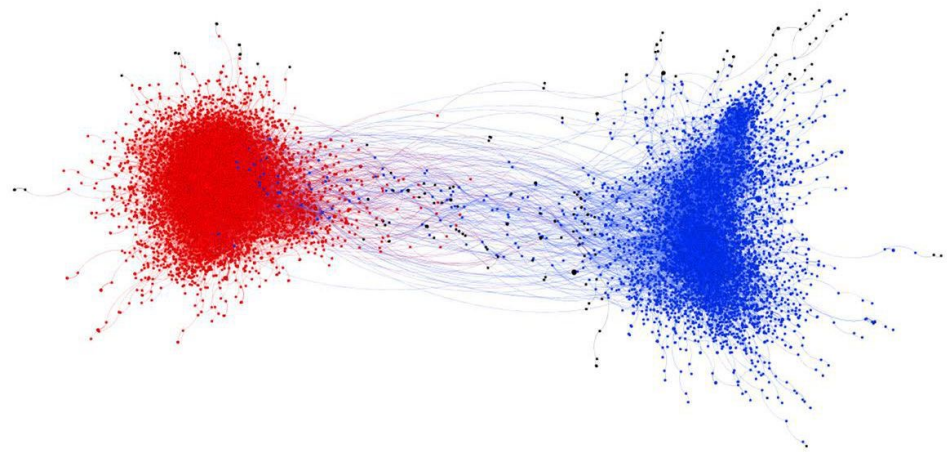


The more you like something, the more RS will recommend similar things, and thus you like them even more.



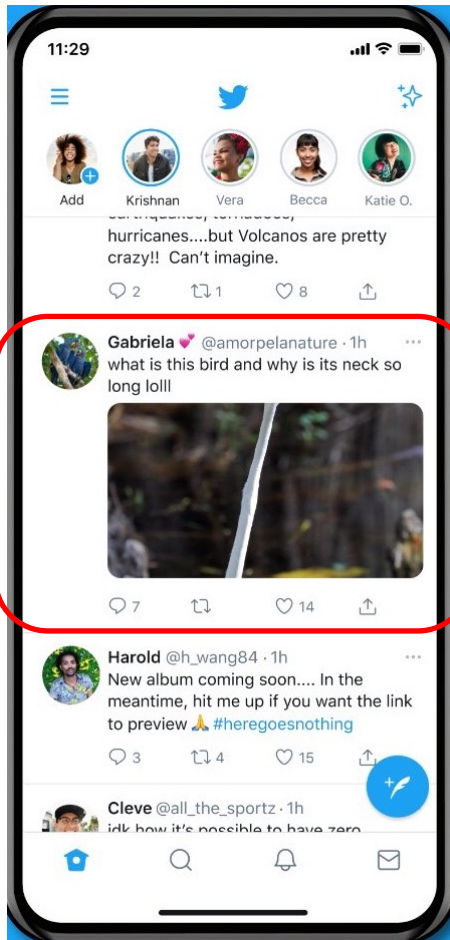
# Towards User Controllable Recommender Systems

- Users almost have **no control** of their recommender system
  - They can only **passively** receive recommendations
- The Social Echo Chamber
  - Makes all your connections like-minded persons as you
  - Makes all your news feed similar as what you already liked
  - Makes it difficult to explore new ideas and opinions different from yours
  - May even reinforce people's extreme ideas





# User Control based on Counterfactual Explanations



## Counterfactual Retrospective Explanation:

We recommend this video X because you previously liked videos A and B, if you didn't like them, then we would not have recommended this video X.

## Counterfactual Prospective Explanation:

If you click "like" on this newly recommended video X, then we will recommend videos such as D and E in the future.

Help users know the **consequences** of their behaviors so that they can take **informed** actions. Users can **control** their recommendation by **invoking or revoking** certain actions.

# Bridging Explainability and Fairness

- Counterfactual Explanation is a flexible framework
  - As long as the **explanation target** can be **quantified**, counterfactual framework can explain it
  - How changes in the **input** influences the output
- **Explainable Fairness** is important in Recommendation
  - Hundreds, thousands or even more features
  - System designers:
    - Want to know which feature(s) **cause** unfairness
  - Users:
    - Want to know how to **intervene** unfair results to make it more fair



# Counterfactual Explainable Fairness

- Too many features in RecSys, manually analysis is almost impossible
  - Automatic explainable fairness is needed.
  - E.g., top-5 features that lead to exposure unfairness

<b>Method</b>	<b>Feature-based Explanations</b>
Pop-User	food, service, chicken, prices, hour
Pop-Item	food, service, prices, visit, hour
EFM-User	store, patio, dishes, dish, rice
EFM-Item	flavor, decor, dishes, inside, cheese
SV	server, size, pizza, food, restaurant
CEF	meal, cheese, dish, chicken, taste

**Table 5: Top-5 feature-based explanations on Yelp dataset.**

# Counterfactual Explainable Fairness

- Counterfactual Explainable Fairness framework

min. Explanation Complexity

*s. t.*, Model Unfairness  $\leq \delta$

$$\min \|\Psi^{cf}\|_2^2 + \lambda \|\Delta\|_2$$

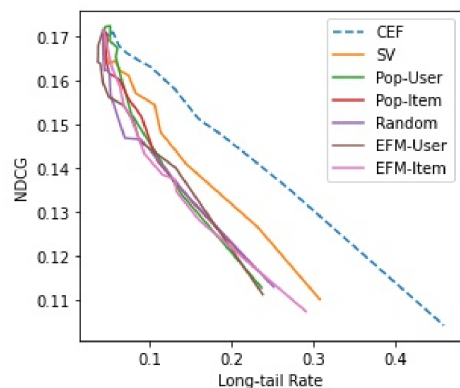
$$\Psi_{DP} = |\mathcal{G}_1| \cdot \text{Exposure}(\mathcal{G}_0|\mathcal{R}_K) - |\mathcal{G}_0| \cdot \text{Exposure}(\mathcal{G}_1|\mathcal{R}_K)$$

Fairness definition: equal opportunity fairness

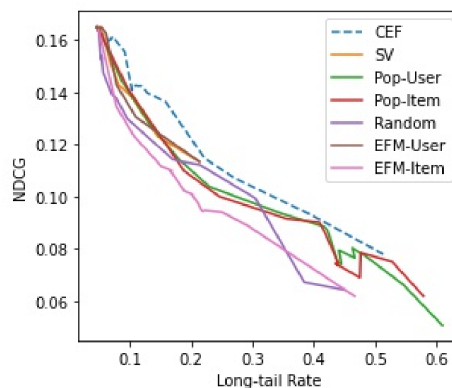
Can be any other definition

# Counterfactual Explainable Fairness

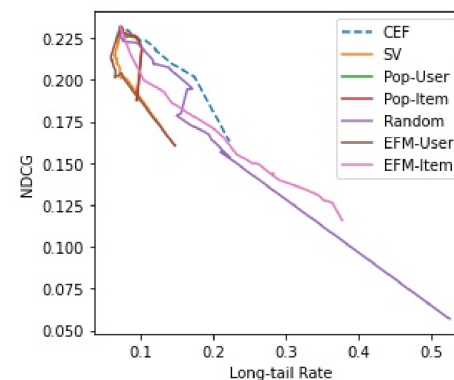
- Better Fairness-Utility Trade-off



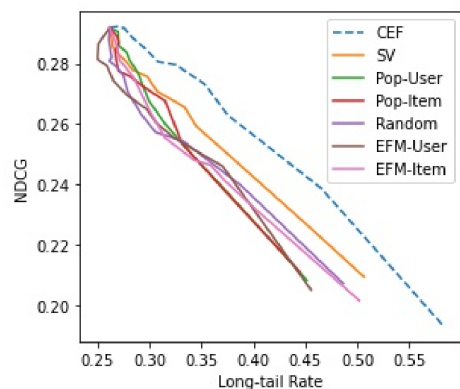
(a) NDCG@5 vs Long-tail Rate@5 on Yelp



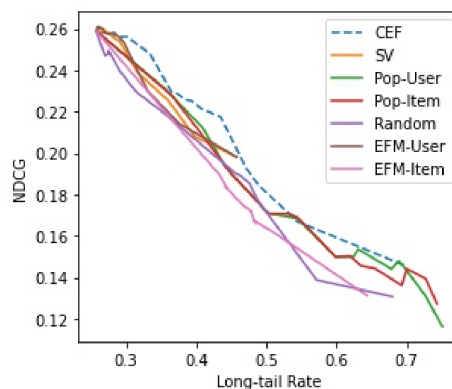
(b) NDCG@5 vs Long-tail Rate@5 on Electronics



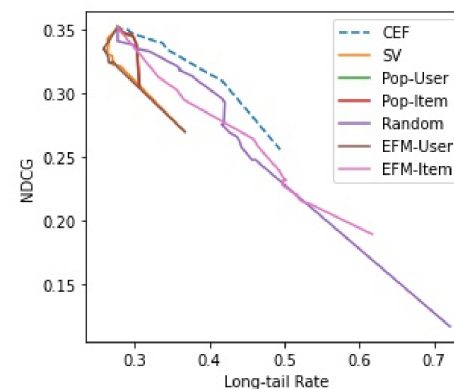
(c) NDCG@5 vs Long-tail Rate@5 on CDs&Vinyl



(d) NDCG@20 vs Long-tail Rate@20 on Yelp



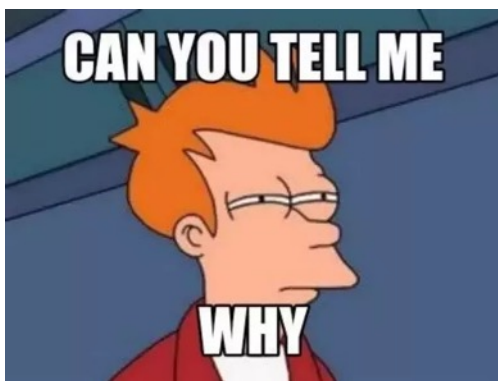
(e) NDCG@20 vs Long-tail Rate@20 on Electronics



(f) NDCG@20 vs Long-tail Rate@20 on CDs&Vinyl

# Natural Language Explanations

- Natural language sentence is the most human-friendly way of explanation
  - Human and machine will inevitably **collaborate with each** other in future jobs
  - We believe future machines should be able to **explain themselves** through **natural language**
  - Better **understanding**, **collaboration** and **trust** between human and machines



I am making  
this decision  
**because** ...



# Natural Language Explanation in Recommendation

- Explainable Recommendation as Natural Language Generation
  - Recommendation is a [very suitable task](#) for developing natural language explanation models
  - High quality [ground-truth explanations](#) from humans

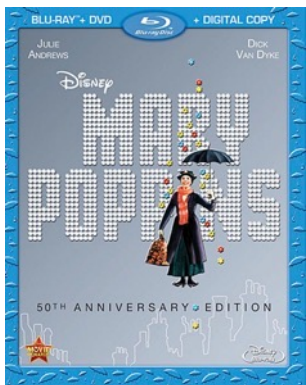


★★★★★ **Trip Saver**

Reviewed in the United States on October 22, 2017

Style: With Lifetime Maps and Traffic (USA) | **Verified Purchase**

Perfect. Lots of features... accurate for finding upcoming restaurants, gas stations and community services. We drive cross country every summer and updated our older GPS to this. Worked through all states, even in low-service areas through the desert. I like being able to search ahead for hotels and restaurants. The battery lasted a long time and there wasn't a lot of screen glare. We also purchased the weighted holder which we really liked.



★★★★★ **A classic musical that is still entrancing and fun to watch**

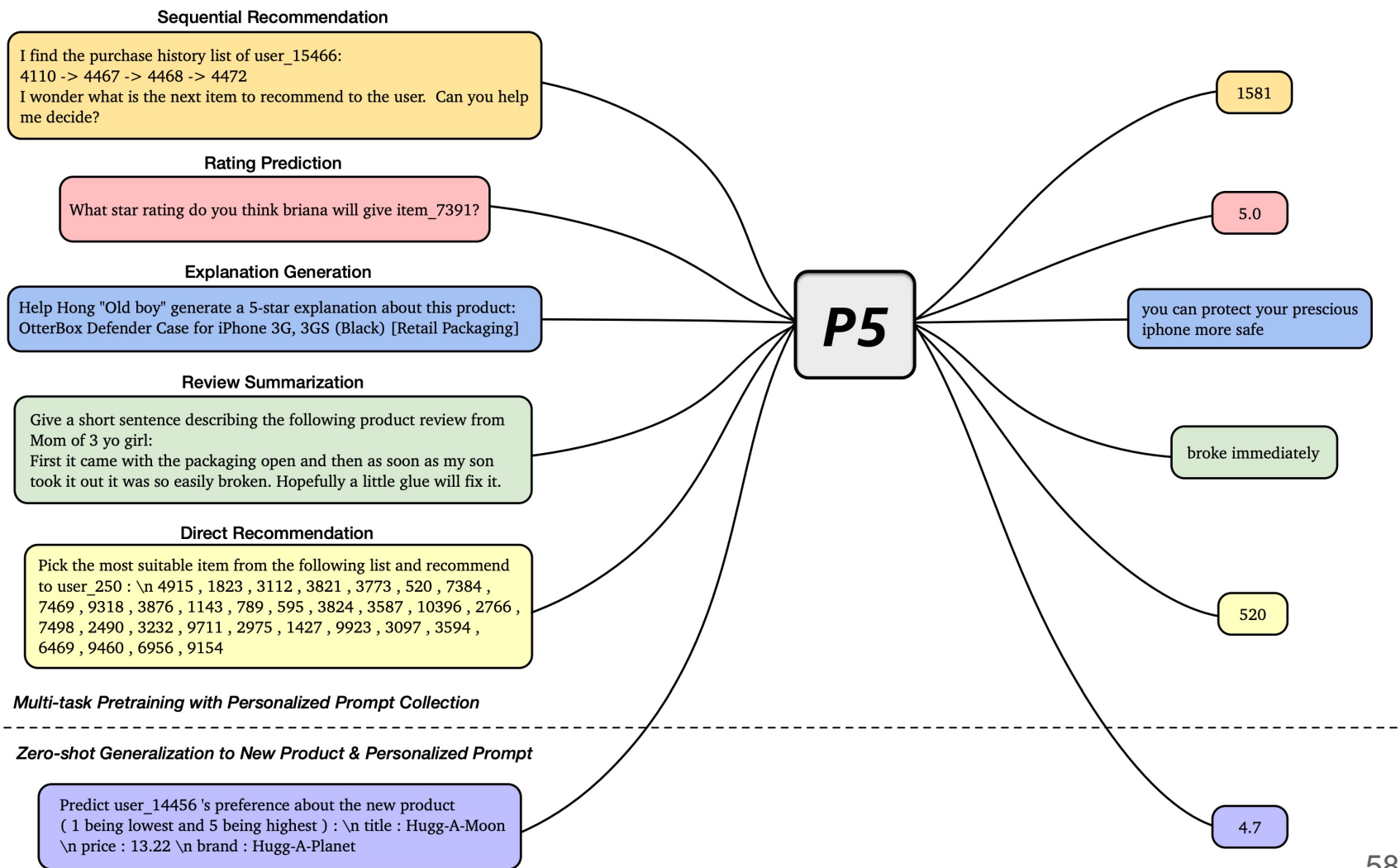
Reviewed in the United States on November 7, 2017

**Verified Purchase**

The movie holds up well as a glorious musical. The acting, singing, choreography, staging, special effects are all great. The plot still works. This is a movie my wife and I love watching over and over. The blu ray version is beautiful. The quality of the image shows itself during the extreme close-ups of Julie Andrews and Dick Van Dyke -- the images are crystal clear with no blurring on a high quality 53 inch LCD HDTV. The sound is excellent. The DVD authoring is a little idiosyncratic. Though "resume play" is activated but is only present after a lengthy video introduction and it is hard to bypass the "previews." There is a paucity of extras.

22 people found this helpful

# Large Recommendation Models (LRM) for Universal Recommendation Engine



# Pretrain, Personalized Prompt & Predict Paradigm (P5)

## Rating / Review / Explanation raw data for *Beauty*

```

user_id: 7641      user_name: stephanie
item_id: 2051
item_title: SHANY Nail Art Set (24 Famouse Colors
Nail Art Polish, Nail Art Decoration)
review: Absolutely great product. I bought this for my fourteen year
old niece for Christmas and of course I had to try it out, then I
tried another one, and another one and another one. So much fun!
I even contemplated keeping a few for myself!
star_rating: 5
summary: Perfect!
explanation: Absolutely great product      feature_word: product
    
```

Which star rating will user\_{{user\_id}} give item\_{{item\_id}}?  
(1 being lowest and 5 being highest)

→ {{star\_rating}}

Based on the feature word {{feature\_word}}, generate an  
explanation for user\_{{user\_id}} about this product:  
{{item\_title}}

→ {{explanation}}

Give a short sentence describing the following product review  
from {{user\_name}}: {{review}}

→ {{summary}}

(a)

## Sequential Recommendation raw data for *Beauty*

```

user_id: 7641      user_name: Victor
purchase_history: 652 -> 460 -> 447 -> 653 -> 654 -> 655 -> 656 -> 8
-> 657
next_item: 552
candidate_items: 4885 , 4280 , 4886 , 1907 , 870 , 4281 , 4222 ,
4887 , 2892 , 4888 , 2879 , 3147 , 2195 , 3148 , 3179 , 1951 ,
..... , 1982 , 552 , 2754 , 2481 , 1916 , 2822 , 1325
    
```

Here is the purchase history of user\_{{user\_id}}:  
{{purchase\_history}}  
What to recommend next for the user?

→ {{next\_item}}

(b)

## Direct Recommendation raw data for *Beauty*

```

user_id: 250      user_name: moriah rose
target_item: 520
random_negative_item: 9711
candidate_items: 4915 , 1823 , 3112 , 3821 , 3773 , 520 , 7384 ,
7469 , 9318 , 3876 , 1143 , 789 , 595 , 3824 , 3587 , 10396 ,
..... , 2766 , 7498 , 2490 , 3232 , 9711 , 2975 , 1405 , 8051
    
```

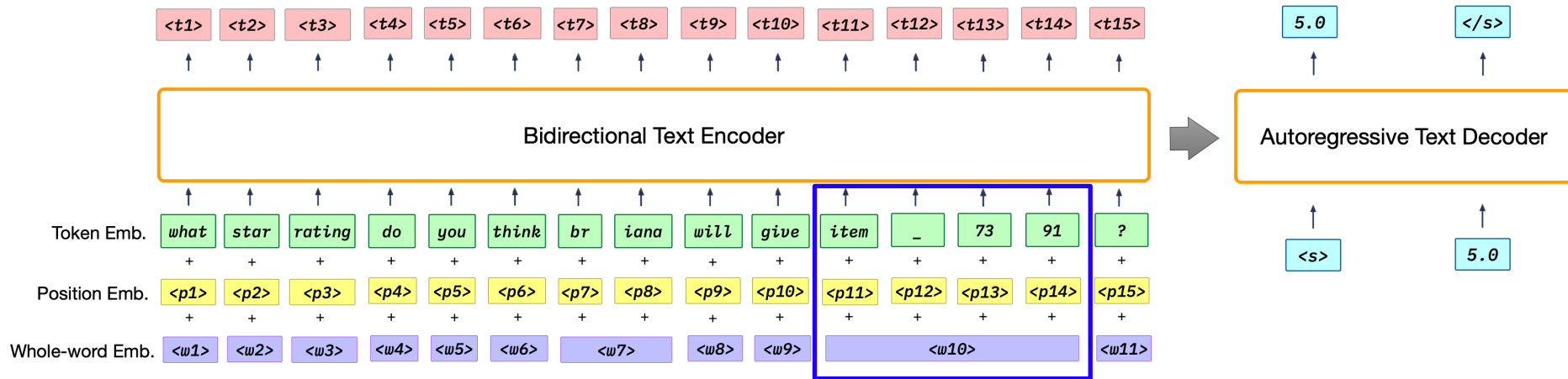
Choose the best item from the candidates to recommend for  
{{user\_name}}? \n {{candidate\_items}}

→ {{target\_item}}

(c)

# The P5 Architecture

- P5 Architecture



ID tokenization is critically important  
Keep a constant and manageable amount of tokens



# Better Recommendation Accuracy

**Table 2: Performance comparison on rating prediction.**

Methods	Sports		Beauty		Toys	
	RMSE	MAE	RMSE	MAE	RMSE	MAE
MF	<b>1.0234</b>	0.7935	<b>1.1973</b>	0.9461	<b>1.0123</b>	0.7984
P5-S (1-6)	1.0594	<b>0.6639</b>	1.3114	<b>0.8434</b>	<u>1.0605</u>	0.7142
P5-B (1-6)	1.0357	0.6813	<u>1.2843</u>	0.8534	1.0866	<b>0.6957</b>
P5-S (1-10)	1.0522	<u>0.6698</u>	1.3001	<u>0.8444</u>	1.0805	0.7057
P5-B (1-10)	<u>1.0292</u>	0.6864	1.2862	0.8530	1.0843	<u>0.7007</u>

**Table 3: Performance comparison on sequential recommendation.**

Methods	Sports				Beauty				Toys			
	HR@5	NDCG@5	HR@10	NDCG@10	HR@5	NDCG@5	HR@10	NDCG@10	HR@5	NDCG@5	HR@10	NDCG@10
Caser	0.0116	0.0072	0.0194	0.0097	0.0205	0.0131	0.0347	0.0176	0.0166	0.0107	0.0270	0.0141
HGN	0.0189	0.0120	0.0313	0.0159	0.0325	0.0206	0.0512	0.0266	0.0321	0.0221	<u>0.0497</u>	0.0277
GRU4Rec	0.0129	0.0086	0.0204	0.0110	0.0164	0.0099	0.0283	0.0137	0.0097	0.0059	0.0176	0.0084
BERT4Rec	0.0115	0.0075	0.0191	0.0099	0.0203	0.0124	0.0347	0.0170	0.0116	0.0071	0.0203	0.0099
FDSA	0.0182	0.0122	0.0288	0.0156	0.0267	0.0163	0.0407	0.0208	0.0228	0.0140	0.0381	0.0189
P5-S (2-3)	0.0272	0.0169	0.0361	0.0198	<u>0.0508</u>	<b>0.0385</b>	<b>0.0668</b>	<b>0.0436</b>	<b>0.0385</b>	<b>0.0269</b>	<b>0.0499</b>	<b>0.0305</b>
P5-B (2-3)	<u>0.0364</u>	<u>0.0296</u>	<u>0.0431</u>	<u>0.0318</u>	<b>0.0515</b>	<u>0.0381</u>	<u>0.0664</u>	<u>0.0429</u>	0.0363	0.0257	0.0457	0.0287
P5-S (2-13)	0.0258	0.0159	0.0346	0.0188	0.0502	0.0378	0.0656	0.0428	<u>0.0370</u>	<u>0.0260</u>	0.0471	<u>0.0293</u>
P5-B (2-13)	<b>0.0387</b>	<b>0.0312</b>	<b>0.0460</b>	<b>0.0336</b>	0.0499	0.0366	0.0651	0.0415	0.0346	0.0244	0.0444	0.0276

# Better Explanation Quality

Table 4: Performance comparison on explanation generation.

Methods	Sports				Beauty				Toys			
	BLUE4	ROUGE1	ROUGE2	ROUGEL	BLUE4	ROUGE1	ROUGE2	ROUGEL	BLUE4	ROUGE1	ROUGE2	ROUGEL
Attn2Seq	0.5305	12.2800	1.2107	9.1312	0.7889	12.6590	1.6820	9.7481	1.6238	13.2245	2.9942	10.7398
NRT	0.4793	11.0723	1.1304	7.6674	0.8295	12.7815	1.8543	9.9477	<u>1.9084</u>	13.5231	3.6708	11.1867
PETER	<b>0.7112</b>	12.8944	1.3283	9.8635	1.1541	14.8497	2.1413	11.4143	<b>1.9861</b>	14.2716	3.6718	11.7010
P5-S (3-3)	0.5902	<b>60.8892</b>	<u>17.7514</u>	<u>18.0010</u>	<u>2.6533</u>	<u>61.6557</u>	<u>21.6574</u>	<u>25.6646</u>	0.3787	<b>56.7474</b>	<u>17.1475</u>	<u>16.7914</u>
P5-B (3-3)	<u>0.6213</u>	<u>58.7260</u>	<b>18.5533</b>	<b>18.4670</b>	<b>3.1474</b>	<b>62.2778</b>	<b>21.9762</b>	<b>27.1758</b>	0.5652	<u>56.4732</u>	<b>17.7930</b>	<b>18.3364</b>
PETER+	2.4627	24.1181	5.1937	18.4105	3.2606	25.5541	5.9668	19.7168	4.7919	28.3083	9.4520	22.7017
P5-S (3-9)	<b>7.2129</b>	<b>67.4004</b>	<b>36.1417</b>	<b>30.8359</b>	5.4136	67.9526	36.5097	30.7446	<u>8.2721</u>	<u>69.4591</u>	<u>39.9955</u>	33.6941
P5-B (3-9)	3.5598	64.7683	34.0162	26.3184	<u>6.5551</u>	<b>68.2939</b>	<u>36.7586</u>	<u>31.8136</u>	<b>9.5411</b>	<b>69.6964</b>	<b>40.3364</b>	<b>34.7272</b>
P5-S (3-12)	<u>5.8446</u>	<u>66.5976</u>	<u>35.5160</u>	<u>29.2766</u>	5.5760	68.1710	<b>36.7876</b>	30.8561	7.5790	69.2164	39.9065	33.1177
P5-B (3-12)	4.6977	65.4562	34.9379	27.7223	<b>7.0183</b>	<u>68.1908</u>	36.7262	<b>32.2162</b>	8.2461	69.2331	39.9456	<u>34.0081</u>

# Zero-Shot Generalization to Items in New Domains

**Table 9: Performance on zero-shot domain transfer.**

Directions	Z-1 & Z-4	Z-2 & Z-3	Z-5 & Z-7		Z-6	
	Accuracy	MAE	BLUE4	ROUGE1	BLUE4	ROUGE1
<i>Toys -&gt; Beauty</i>	0.7922	0.8244	1.8869	61.1919	5.4609	66.4931
<i>Toys -&gt; Sports</i>	0.8682	0.6644	0.7405	60.9575	2.2601	62.0353
<i>Beauty -&gt; Toys</i>	0.8073	0.7792	0.0929	41.3061	11.8046	64.8701
<i>Beauty -&gt; Sports</i>	0.8676	0.6838	0.0346	39.7191	6.6409	66.9222
<i>Sports -&gt; Toys</i>	0.8230	0.7443	0.0687	42.9310	13.3408	69.7910
<i>Sports -&gt; Beauty</i>	0.8057	0.8102	0.0790	41.0659	13.1690	66.7687

## Prompt ID: Z-1

Input template: Given the facts about the new product, do you think user `{{user_id}}` will like or dislike it? title: `{{item_title}}` brand: `{{brand}}` price: `{{price}}`

Target template: `{{answer_choices[label]}}` (like/dislike) - like (4,5) / dislike (1,2,3)

## Prompt ID: Z-2

Input template: Here are the details about a new product: title: `{{item_title}}` brand: `{{brand}}` price: `{{price}}` What star will `{{user_desc}}` probably rate the product?  
-1 -2 -3 -4 -5

Target template: `{{star_rating}}`

## Prompt ID: Z-5

Input template: Generate a possible explanation for `{{user_desc}}`'s preference about the following product: title: `{{item_title}}` brand: `{{brand}}` price: `{{price}}`

Target template: `{{explanation}}`

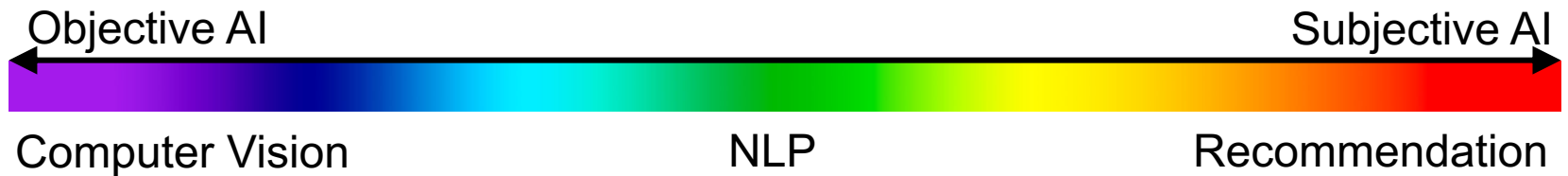
## Prompt ID: Z-6

Input template: Based on the word `{{feature_word}}`, help user\_`{{user_id}}` write a `{{star_rating}}`-star explanation for this new product: title: `{{item_title}}` price: `{{price}}` brand: `{{brand}}`

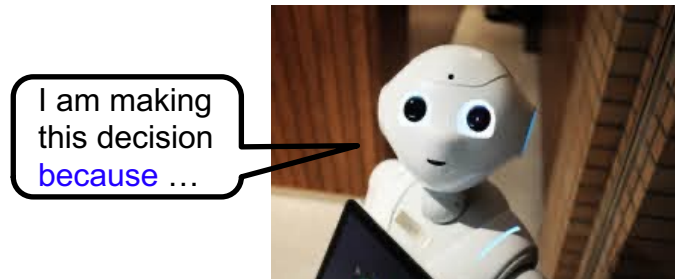
Target template: `{{explanation}}`

# Summary

- Trustworthy and Responsible Recommendation
  - Explainability, Fairness, Echo Chambers, Controllability
  - Many other perspectives: Robustness, Accountability, Privacy, etc.



	Human-centered Tasks
<b>Counterfactual Reasoning</b>	Counterfactual Explainable Recommendation
<b>Counterfactual Fairness</b>	Counterfactual Explainable Fairness
<b>Human-controllable AI</b>	User Controllable Recommendation
<b>Large Recommendation Models</b>	Multi-task Learning, Natural Language Explanation





Yongfeng Zhang

Department of Computer Science, Rutgers University

[yongfeng.zhang@rutgers.edu](mailto:yongfeng.zhang@rutgers.edu)

<http://yongfeng.me>