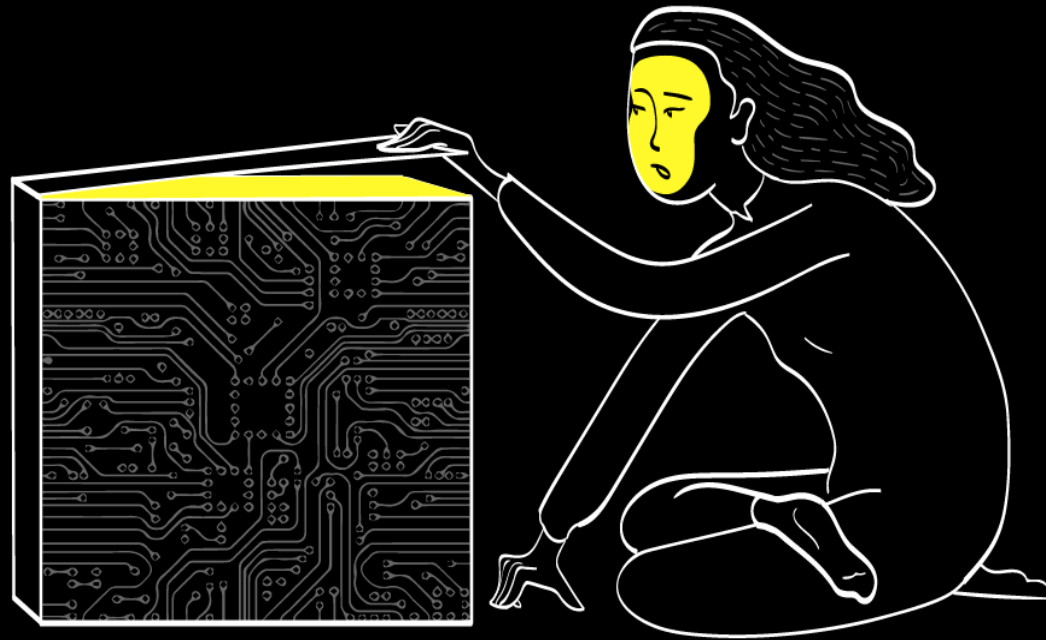# An overview of interpretability of Recommender Systems

# Agenda

Explainability
- Why do we need to explain?
- Notions of explainability
- Meaningful explanations

Explainability in Recommender Systems
- Paradigms
- Information Category
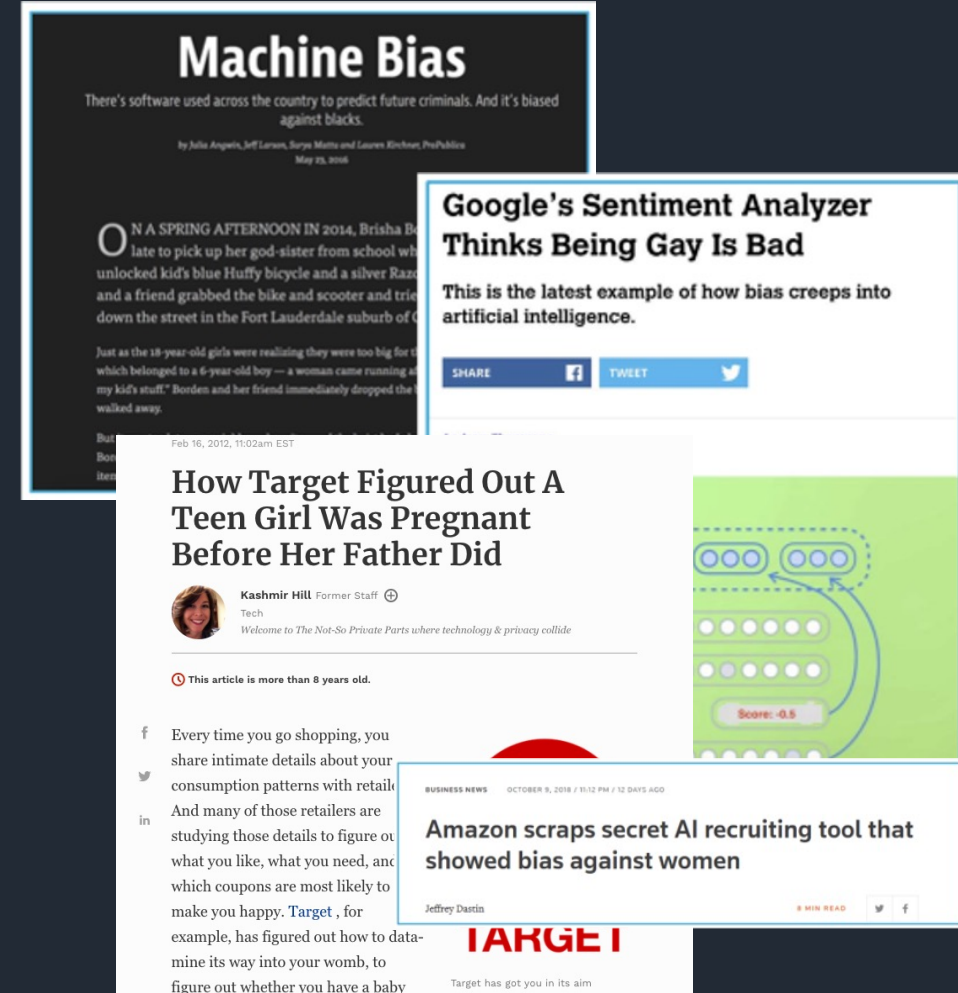- Model
- Evaluation

Summary

# Why Explainability?

AI is now used in many high-stakes decision-making applications (credit, employment, admission, sentencing).

89% of consumers say…

"technology companies need to be more transparent."

Technology companies need to comply with GDPR – "Right to Explanation."

**Machine Bias**

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

ON A SPRING AFTERNOON IN 2014, Brisha B⋯ late to pick up her god-sister from school wh⋯ unlocked kid's blue Huffy bicycle and a silver Razo⋯ and a friend grabbed the bike and scooter and trie⋯ down the street in the Fort Lauderdale suburb of C⋯

Just as the 18-year-old girls were realizing they were too big for th⋯ which belonged to a 6-year-old boy — a woman came running a⋯ my kid's stuff." Borden and her friend immediately dropped the ⋯ walked away.

**Google's Sentiment Analyzer Thinks Being Gay Is Bad**

This is the latest example of how bias creeps into artificial intelligence.

Feb 16, 2012, 11:02am EST

**How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did**

Kashmir Hill Former Staff
Tech
Welcome to The Not-So Private Parts where technology & privacy collide

⏱ This article is more than 8 years old.

Every time you go shopping, you share intimate details about your consumption patterns with retail⋯ And many of those retailers are studying those details to figure ou⋯ what you like, what you need, and⋯ which coupons are most likely to make you happy. Target , for example, has figured out how to data-mine its way into your womb, to figure out whether you have a baby

Score: -0.8

**TARGET**

Target has got you in its aim

BUSINESS NEWS    OCTOBER 9, 2018 / 11:12 PM / 12 DAYS AGO

**Amazon scraps secret AI recruiting tool that showed bias against women**
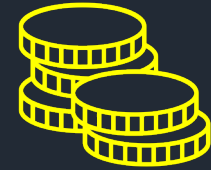
Jeffrey Dastin

# Explainability across industries

**ADVERTISING**
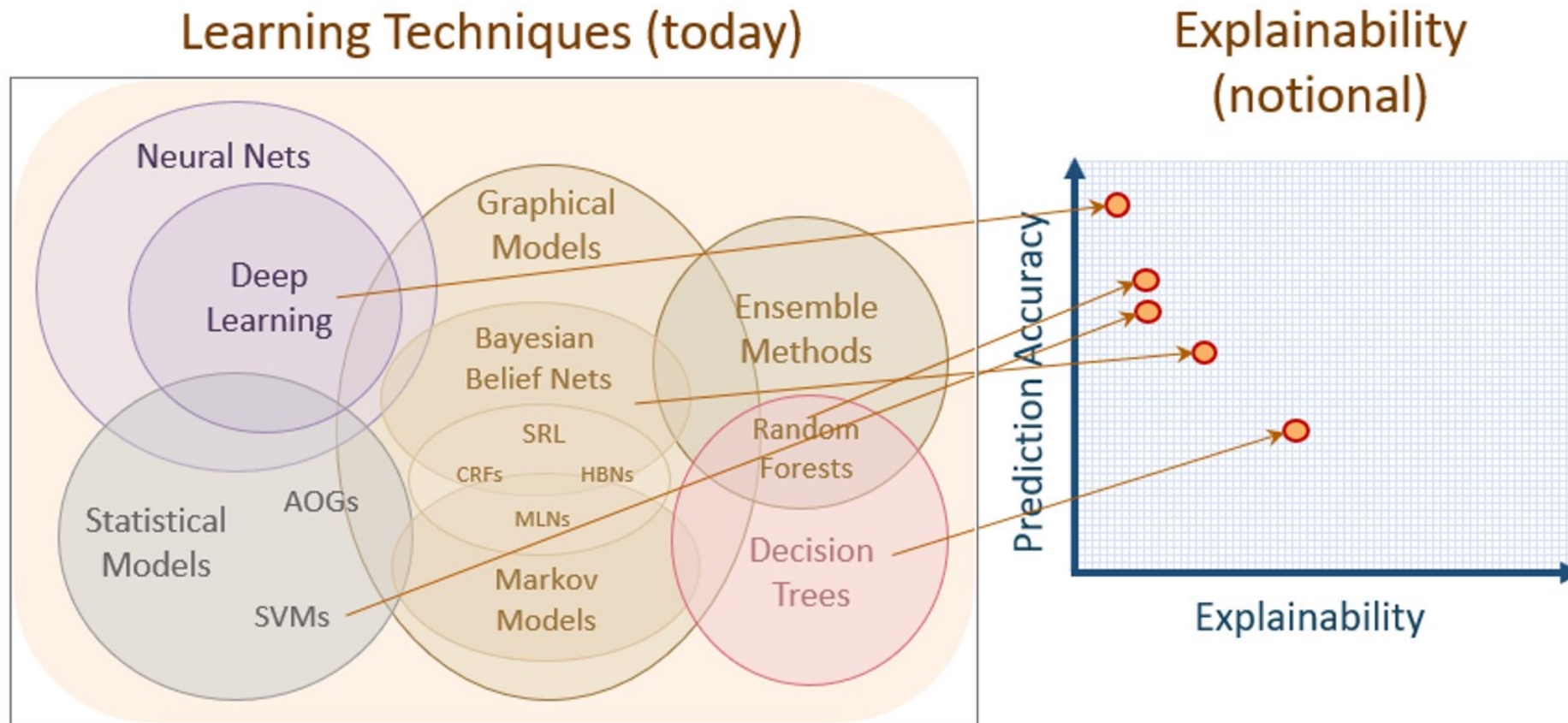The reputational risk of placing adverts alongside content that doesn't 'fit' the brand

**HEALTH**
49% of physicians in the US are anxious or uncomfortable with AI

**FINANCE**
Explaining why automated decision-making rejects loan applications

Coba, L.

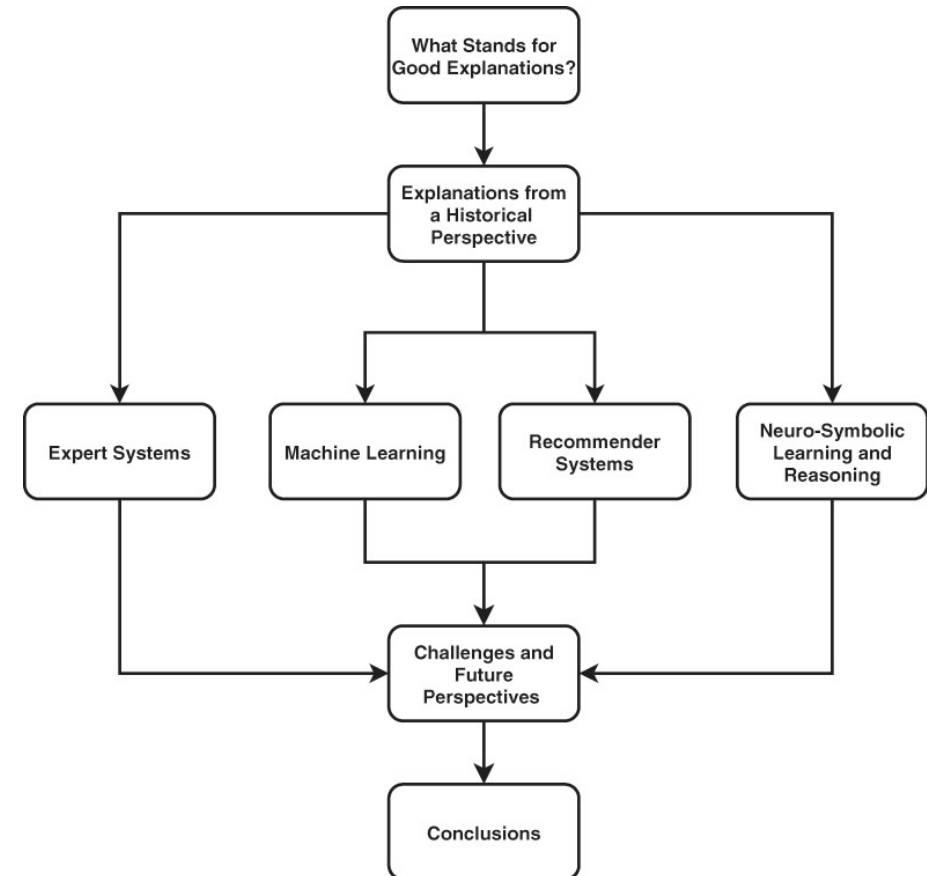# The most effective algorithms are the hardest to explain

Coba, L.

# Why Explainability is a challenge?

Different notions

Different requirements

Plethora of approaches



Confalonieri, R., Coba, L., Wagner, B., and Besold, T.. A historical perspective of explainable artificial intelligence. WIREs Data Mining and Knowledge Discovery, 11(1), 2021. doi: https://doi.org/10.1002/widm.1391

Coba, L.

# Explainability - Notions

Interpretable Systems

A system where a user cannot only see, but also study and understand how inputs are mathematically mapped to outputs.

E.g. regression models, support vector machines, decision trees, ANOVAs, and data clustering (assuming a kernel that is itself interpretable)

Comprehensible Systemes

A system that emits symbols along with its output that allow the user to relate properties of the inputs to their output. The user is responsible for compiling and comprehending the symbols, relying on her own implicit form of knowledge and reasoning about them.

E.g. High dimensional data visualizations like t-SNE and receptive field visualization on convolutional neural networks

Opaque systems

A system where the mechanisms mapping inputs to outputs are invisible to the user

Coba, L.

Doran, D., Schulz, S., & Besold, T. R. (2017). What Does Explainable AI Really Mean? A New Conceptualization of Perspectives. 1st International Workshop on Comprehensibility and Explanation in AI and ML Colocated with AI*IA 2017 (Vol. 2071).

# Explainability - Notions

Interpretable Systems

A system where a user cannot only see, but also study and understand how inputs are mathematically mapped to outputs.

E.g. regression models, support vector machines, decision trees, ANOVAs, and data clustering (assuming a kernel that is itself interpretable)
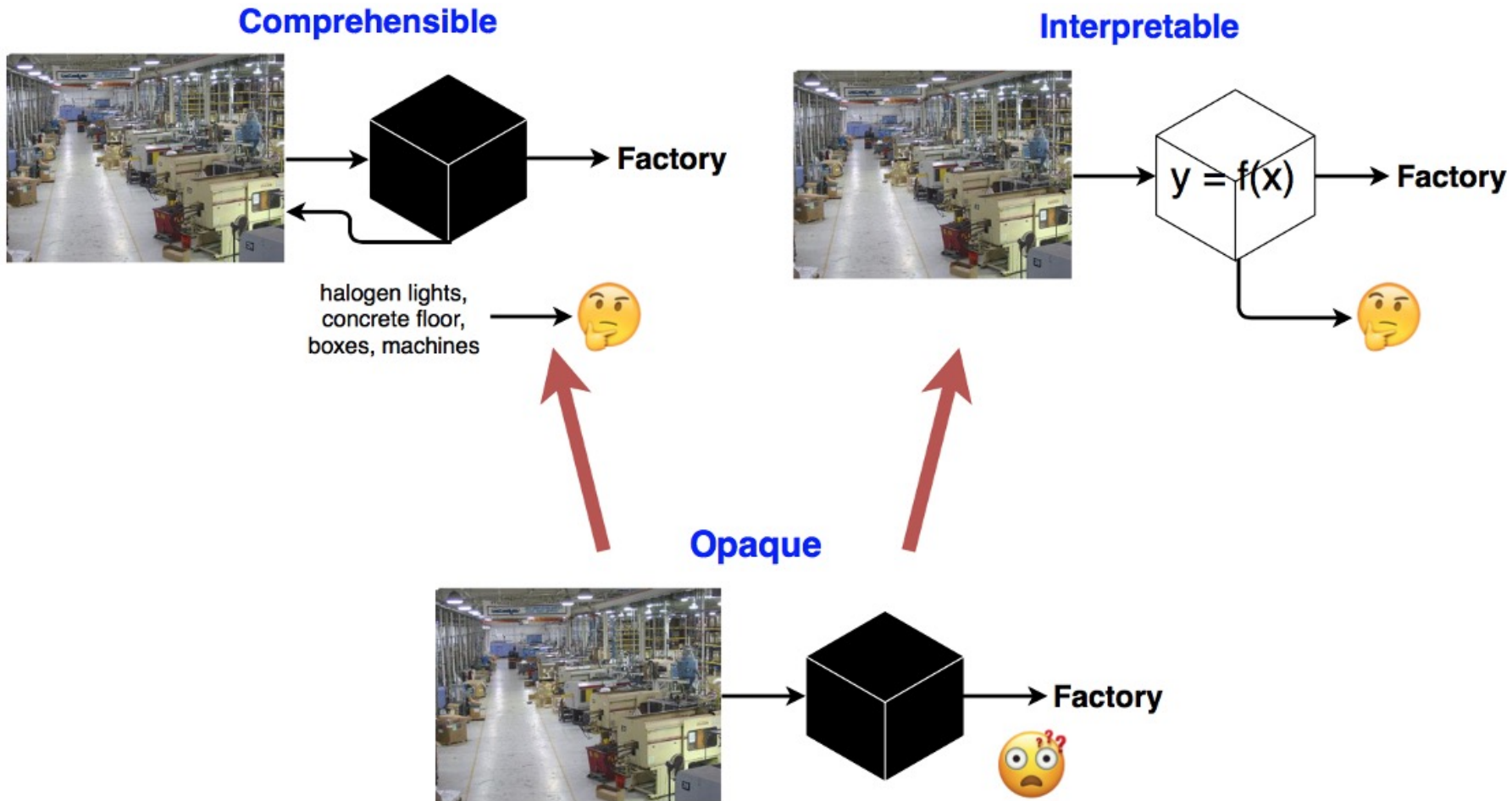
## Comprehensible Systems

A system that emits symbols along with its output that allow the user to relate properties of the inputs to their output. The user is responsible for compiling and comprehending the symbols, relying on her own implicit form of knowledge and reasoning about them.

E.g. High dimensional data visualizations like t-SNE and receptive field visualization on convolutional neural networks

Opaque systems

A system where the mechanisms mapping inputs to outputs are invisible to the user

Coba, L.    Doran, D., Schulz, S., & Besold, T. R. (2017). What Does Explainable AI Really Mean? A New Conceptualization of Perspectives. 1st International Workshop on Comprehensibility and Explanation in AI and ML Colocated with AI*IA 2017 (Vol. 2071).

# Explainability - Notions

Interpretable Systems

A system where a user cannot only see, but also study and understand how inputs are mathematically mapped to outputs.

E.g. regression models, support vector machines, decision trees, ANOVAs, and data clustering (assuming a kernel that is itself interpretable)

Comprehensible Systemes

A system that emits symbols along with its output that allow the user to relate properties of the inputs to their output. The user is responsible for compiling and comprehending the symbols, relying on her own implicit form of knowledge and reasoning about them.

E.g. High dimensional data visualizations like t-SNE and receptive field visualization on convolutional neural networks
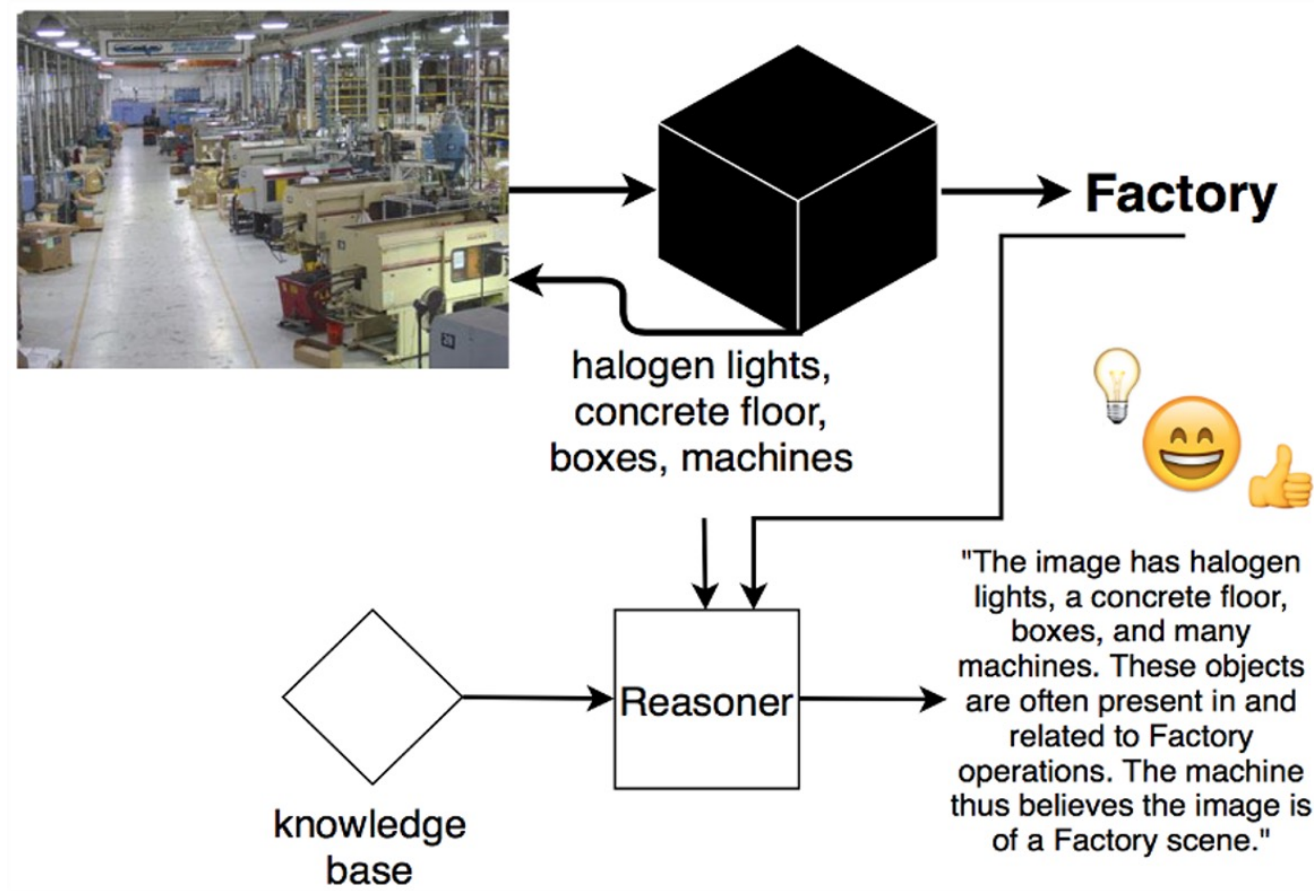
## Opaque systems

A system where the mechanisms mapping inputs to outputs are invisible to the user
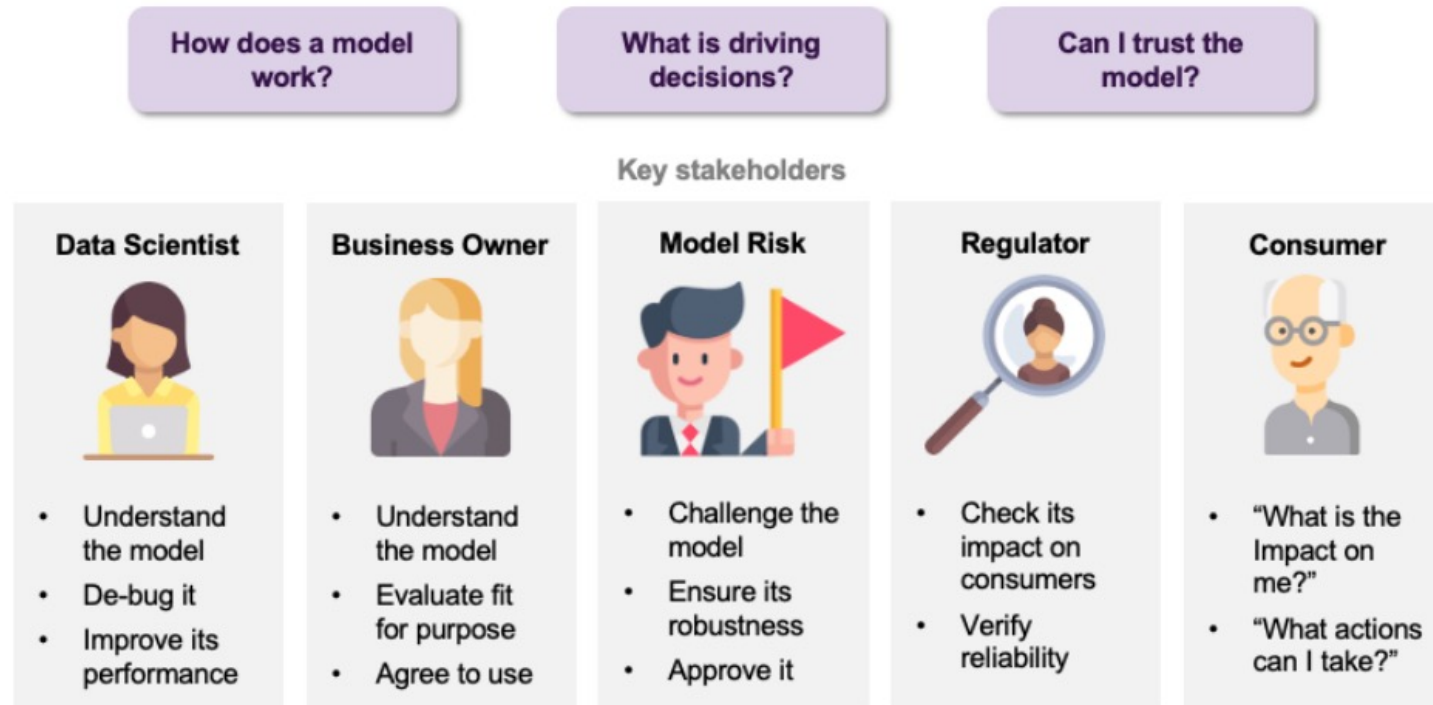
Coba, L.  Doran, D., Schulz, S., & Besold, T. R. (2017). What Does Explainable AI Really Mean? A New Conceptualization of Perspectives. 1st International Workshop on Comprehensibility and Explanation in AI and ML Colocated with AI*IA 2017 (Vol. 2071).

# Explainability - Notions



Doran, D., Schulz, S., & Besold, T. R. (2017). What Does Explainable AI Really Mean? A New Conceptualization of Perspectives. 1st International Workshop on Comprehensibility and Explanation in AI and ML Colocated with AI*IA 2017 (Vol. 2071).

Coba, L.

# Explainability - Notions



Doran, D., Schulz, S., & Besold, T. R. (2017). What Does Explainable AI Really Mean? A New Conceptualization of Perspectives. 1st International Workshop on Comprehensibility and Explanation in AI and ML Colocated with AI*IA 2017 (Vol. 2071).

Coba, L.

# Meaningful explanations depend on the stakeholder!



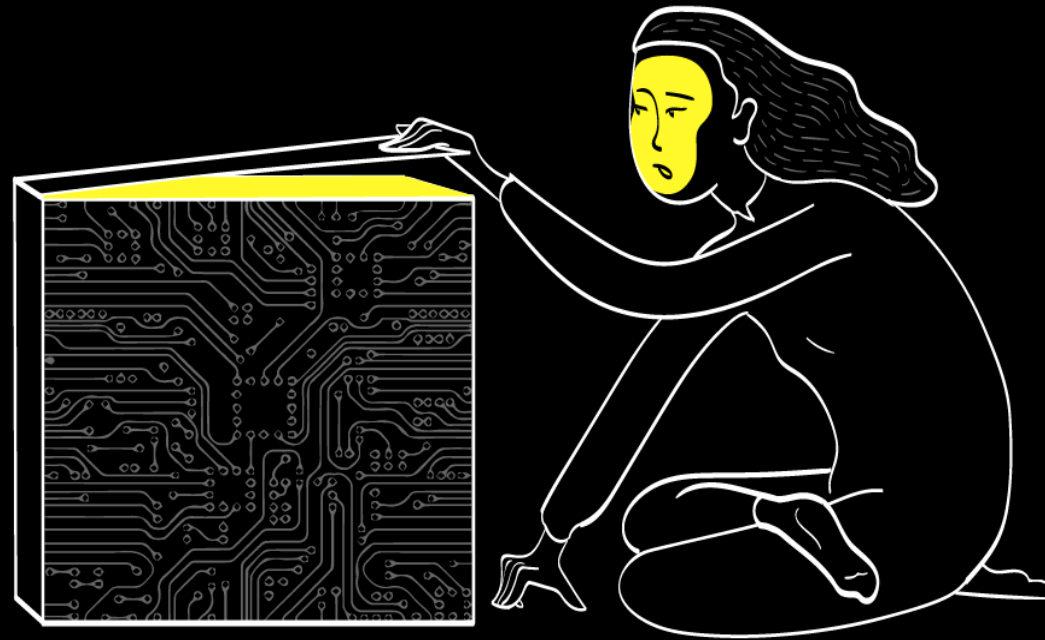How does a model work?

What is driving decisions?

Can I trust the model?

Key stakeholders

**Data Scientist**
- Understand the model
- De-bug it
- Improve its performance

**Business Owner**
- Understand the model
- Evaluate fit for purpose
- Agree to use

**Model Risk**
- Challenge the model
- Ensure its robustness
- Approve it

**Regulator**
- Check its impact on consumers
- Verify reliability

**Consumer**
- "What is the Impact on me?"
- "What actions can I take?"

Belle, V., & Papantonis, I. (2020). Principles and practice of explainable machine learning. *arXiv preprint arXiv:2009.11698*.

# Approaches

Expert Systems

Machine Learning

Neuro-symbolic Learning and Reasoning

Recommender Systems

# Recommender Systems

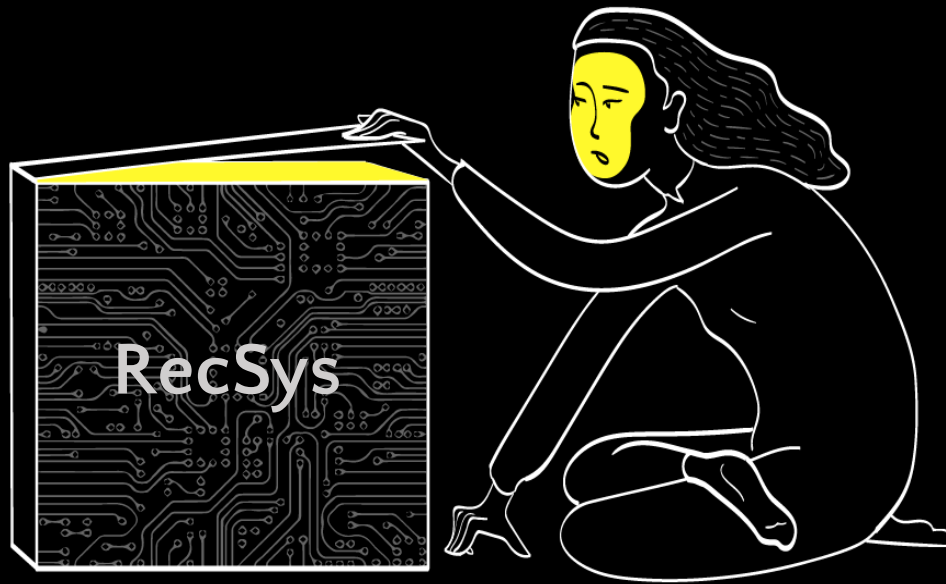# Recommender Systems



RecSys

Users' data

Personalized recommendations

Coba, L.

# Explainability



RecSys

Transparency

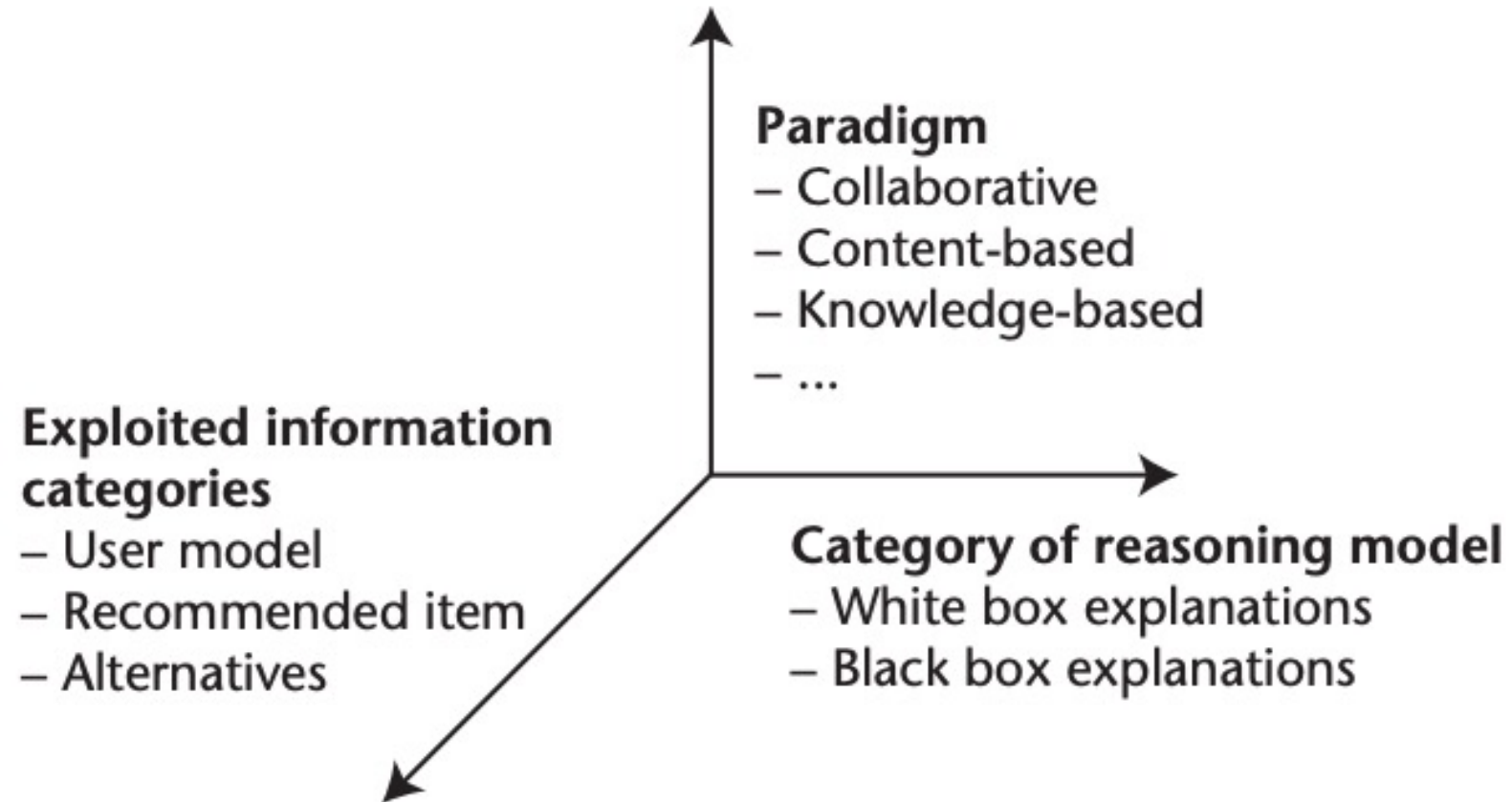Scrutability

Trust

Effectiveness

Persuasiveness

Efficiency

Satisfaction

Tintarev, Nava, and Judith Masthoff. "Designing and evaluating explanations for recommender systems." Recommender systems handbook. Springer, Boston, MA, 2011. 479-510.

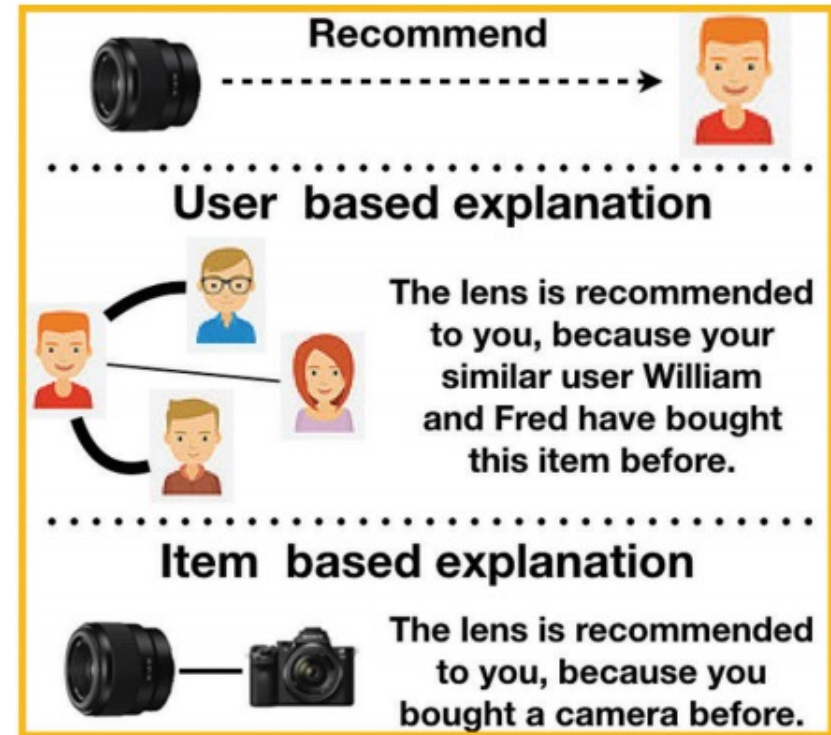Coba, L.

# Explainability



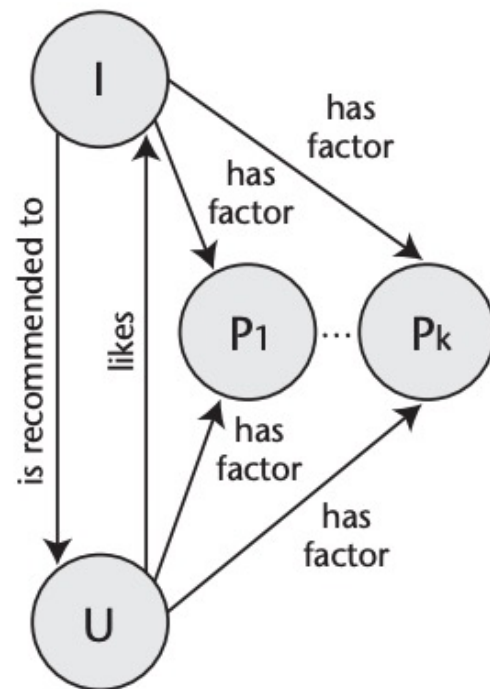**Paradigm**
- Collaborative
- Content-based
- Knowledge-based
- ...

**Exploited information categories**
- User model
- Recommended item
- Alternatives

**Category of reasoning model**
- White box explanations
- Black box explanations

# RecSys paradigmes

**Collaborative Filtering**

**Content-based**

**Knowledge-based**

Coba, L.

**Collaborative Filtering**

is similar to

I

is recommended to

likes

U

is similar to

Recommend

User based explanation

The lens is recommended to you, because your similar user William and Fred have bought this item before.

Item based explanation

The lens is recommended to you, because you bought a camera before.

Coba, L.    Coba, L., Rook, L., Zanker, M., & Symeonidis, P. (2019, March). Decision making strategies differ in the presence of collaborative explanations: two conjoint studies. IUI 19.

**Collaborative Filtering**



Recommended:

Heaven's Prisoners (1996)

Explanation:

- Faster Pussycat! Kill! Kill! (1965)
- Kansas City (1996)
- Turbo: A Power Rangers Movie (1997)
- Brother Minister: The Assassination of Malcolm X (1994)
- Theodore Rex (1995)
- Kull the Conqueror (1997)
- Free Willy 2: The Adventure Home (1995)
- Steel (1997)
- Doom Generation, The (1995)
- Unhook the Stars (1996)

Coba, L.

Coba, L., Confalonieri, R., Zanker, M. (2021). RECOXPLAINER: An Extensible Toolkit for Explainable Recommender Systems. AAAI 21.

**Content-based**

Feature-level explanation

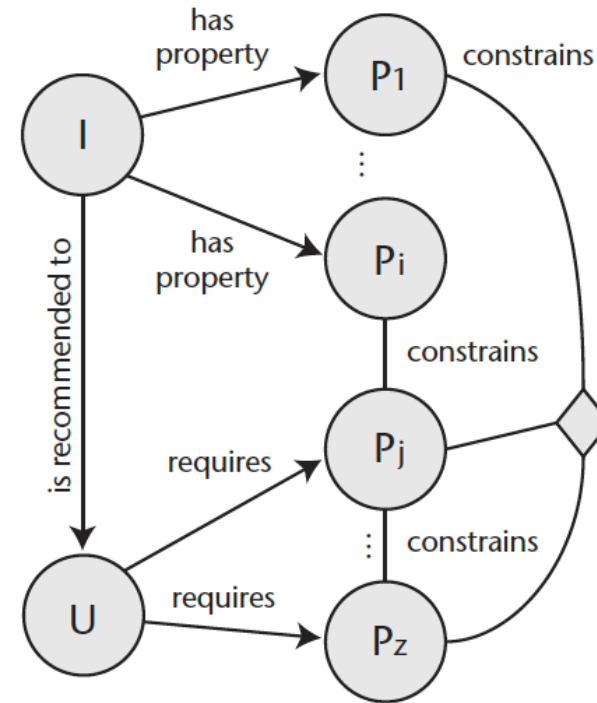| Feature | likeness |
|---|---|
| color | 0.87 |
| quality | 0.54 |
| Focal Length | 0.66 |
| Focus Type | 0.71 |

**Sentence-level explanation**

Structured: You might be interested in [feature] (can be quality, color, etc), on which this product performs well.
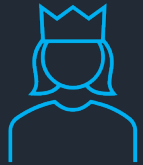Unstructured: Great and deserve the price.

Dominguez, V., Messina, P., Donoso-Guzmán, I., & Parra, D. (2019, March). The effect of explanations and algorithmic accuracy on visual recommender systems of artistic images. IUI 2019.

**Knowledge-based**

# Information Categories

**User Model**
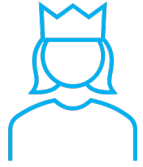Are explanations tailored to the system's beliefs about a given user?

**Recommendation Features**
Is the recommendation dependent on the specific recommended item?

**Alternatives**
Do explanations argue in favour or against alternatives to the recommended item?

Kouki, P., Schaffer, J., Pujara, J., O'Donovan, J., & Getoor, L. (2019, March). Personalized explanations for hybrid recommender systems. In *Proceedings of the 24th International Conference on Intelligent User Interfaces* (pp. 379-390).

# User Model

# Recommendation Features

# Alternatives



This is how **you rated similar movies** on our platform.
★★★★☆ 40 ratings
3.7 out of 5 stars

| | | |
|---|---|---|
| 5 stars | | 23 % |
| 4 stars | | 37 % |
| 3 stars | | 30 % |
| 2 stars | | 7 % |
| 1 star | | 3 % |

**Bugs Bunny, Daffy Duck, Tasmanian Devil** and 6 others like DataCamp. ...

**DataCamp**
Sponsored · 
👍 Like Page
Improve your data science skills with these courses.

GRAVITY

- From your MovieLens profile it seems that you prefer movies tagged as space, this movie takes you in space and it feels claustrophobic to be there. It keeps you on the edge of your seat the whole time.

- From your MovieLens profile it seems that you prefer movies tagged as visual, Gravity is unlike what we have seen on a cinema screen before and arguably it has one of the best uses of 3D in a movie.

- From your MovieLens profile it seems that you prefer movies tagged as intense, the movie a pretty intense ninety minutes, with Bullock's character constantly battling one catastrophe after another, and all of it is amazing to see.

Alice

Alice's Actions

Zone of HIN that cannot be disclosed to Alice but is vital for determining causality by provider

Top-$k$ Rec. items

rec
$i_1$
$i_2$
$i_3$
$i_4$

Bought
Viewed
Reviewed
Follow
Highly rated

**Alice:** Why did I receive this recommendation "Jack Wolfskin backpack"?

**PRINCE:** You **bought** "Adidas Hiking Shoes";
You **reviewed** "Nikon Coolpix Camera" with "Sleek! Handy on hikes!";
You **rated** "Intenso Travel Power Bank" highly.

If you **had not** done these actions:
"iPad Air" **would have replaced** "Jack Wolfskin backpack".

Coba, L.

[User Model] Coba, L., Zanker, M., Rook, L., & Symeonidis, P. (2018). Exploring Users' Perception of Collaborative Explanation Styles. CBI 18.
[Rec Features] Chang, S., Harper, F. M., & Terveen, L. G. (2016, September). Crowd-based personalized natural language explanations for recommendations. RecSys 16.
[Alternatives] Ghazimatin, A., Balalau, O., Saha Roy, R., & Weikum, G. (2020). PRINCE: provider-side interpretability with counterfactual explanations in recommender systems. WSDM 20.

# Reasoning model



**White-box**
How did the system derive a
recommendation.



**Black-box**
What justifies the
recommendation in the eyes
of its recipient.

Coba, L.

Coba, L., Confalonieri, R., Zanker, M. (2021). RecoXplainer: An Extensible Toolkit for Explainable Recommender Systems, AAAI 2021.

# White-box example: Association-rules

Association rule mining algorithms

Detect rules of the form X → Y (e.g., beer → diapers) from a set of transactions T = {t1, t2, … tn} over a catalogue I

Measure quality by means of support, confidence used as a threshold to cut off unimportant rules

Pros:  Interpretable by design
Cons: The model is not flexible

# Explaining Black-box models

Model-Based Explanations are obtained by constraining the loss function

    <span style="color:green">Pros</span>: No interpretable proxies needed

    <span style="color:red">Cons</span>: Model loses flexibility


Post-Hoc Explanations are obtained by means of an interpretable proxy

    <span style="color:green">Pros</span>: No under-the-hood reworking of the black-box

    <span style="color:red">Cons</span>: Additional training step, not complete; Accuracy-interpretability trade-off

Coba, L., Confalonieri, R., Zanker, M. (2021). RecoXplainer: An Extensible Toolkit for Explainable Recommender Systems, AAAI 2021.

# Model-based example: EMF

Desired Explanation style:

We recommend you Movie 1 because your neighbours' ratings for this movie are the following:

| Rating | Number of Neighbours |
|--------|---------------------|
| ⭐ | 0 |
| ⭐ ⭐ | 0 |
| ⭐ ⭐⭐ | 0 |
| ⭐ ⭐⭐⭐ | 10 |
| ⭐⭐⭐⭐⭐ | 23 |

How it works:



**1**

Building the neighborhood (NN) of the users

**2**

$$E_{u,i} = \sum_{\substack{\forall r \in R \\ r \ge P_\tau}} r * |NN^k(u)_{i,r}|$$

Determine the explainability of an item $i$ by measuring in the identified neighborhood how frequently item has been highly rated

**3**

$$\sum_{i,j \in R} (r_{ij} - u_i v_j^T)^2 + \frac{\beta}{2}(\|u_i\|^2 + \|v_j\|^2) +$$

$$\lambda \|u_i - v_j\|^2 E_{ij} \quad \text{Soft constraint}$$

We extend the traditional matrix factorization to recommend explainable items

Popularity bias

Coba, L.

# Post-hoc example: Association-rules proxy

Coba, L.          Peake, G., & Wang, J. (2018). Explanation mining: Post hoc interpretability of latent factor models for recommendation systems. KDD 18.

# Post-hoc example: Association-rules proxy

Association rules to generate post-hoc explanations

Mine association rules on the generated predictions from a black-box RS

For each user filter the learned transactions such that antecedents are in the training set and consequents are unseen or non-interacted items

The resulting subset is ranked by support/confidence/lift. We keep the top-$D$ consequents

# Evaluation

Offline evaluation:

Based on a mathematical understanding of the user.
Examples: Model Fidelity, Mean Explainable Precision, E-nDCG

Online/user studies:

Require feedback from users and are specific to the goal of the explanation.
Examples: standardised psycolagiacal scales measuing trust, efficientcy, etc.

Coba, L.

# Evaluation

Offline evaluation:
    Is based on a mathematical understanding of the user.
    Examples: Model Fidelity, Mean Explainable Precision, E-nDCG

Online/user studies:
    Requires feedback from users and are specific to the goal of the explanation.
    Examples: standardised psycolagiacal scales measuing trust, efficientcy, etc.

# Evaluation

Offline evaluation:
   Is based on a mathematical understanding of the user.
   Examples: Model Fidelity, Mean Explainable precision, E-nDCG

Online/user studies:
   Requires feedback from users and are specific to the goal of the explanation.
   Examples: standardised psycolagiacal scales measuing trust, efficientcy, etc.

(Bonus) Evaluate the Recommender:
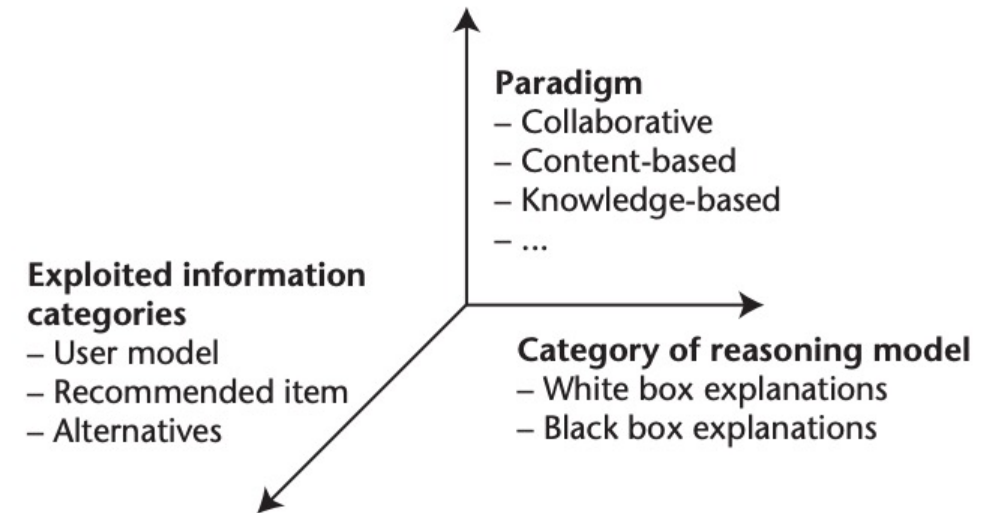   Enable developers to undestand the reasoning and the  quality of the recommender

# Summary

In RecSys explanations are designed based on the paradigm, the information category and the reasoning model

Be aware of the interpretability-accuracy trade-off

Advanteges as trust, persuasivness, efficiency, scrutiny, etc

Offline evaluation via proxy, online evaluation via standardised scales.



**Paradigm**
– Collaborative
– Content-based
– Knowledge-based
– ...

**Exploited information categories**
– User model
– Recommended item
– Alternatives

**Category of reasoning model**
– White box explanations
– Black box explanations

Coba, L.

# Thank you!