

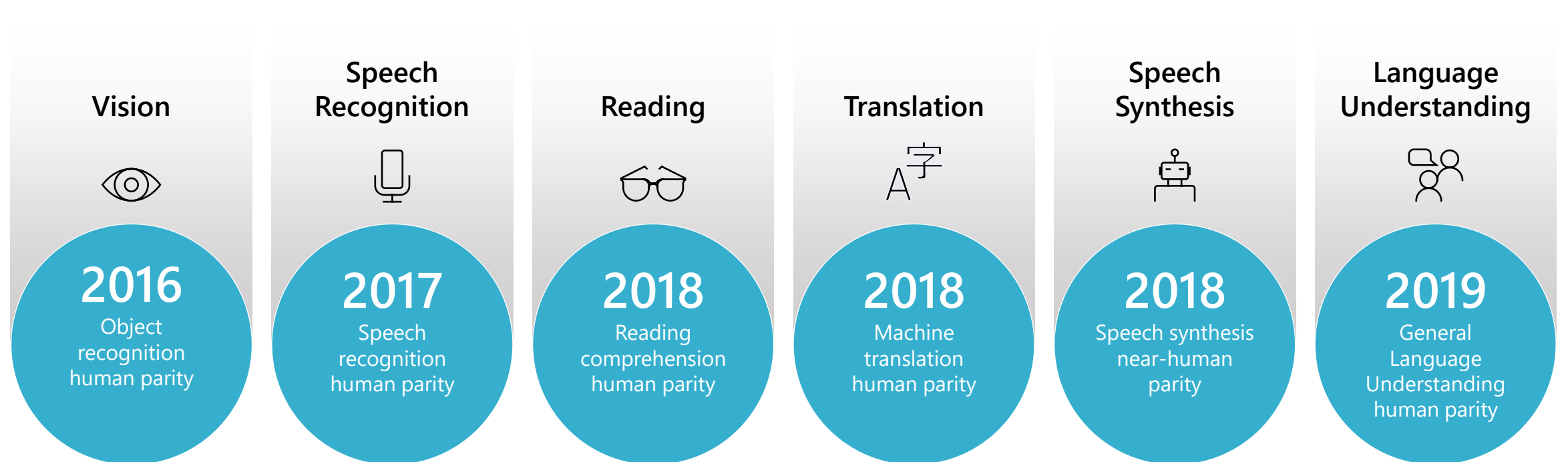
Responsible AI Research at  
Microsoft Research Asia  
微软亚洲研究院负责任的人工智能研究

Xing Xie

Microsoft Research Asia

# Responsible AI

Advancements in AI are different than other technologies because of the **pace of innovation**, and its **proximity to human** intelligence – impacting us at a personal and societal level.



# Microsoft's AI Principles



Fairness



Reliability  
& Safety



Privacy &  
Security



Inclusiveness



Transparency



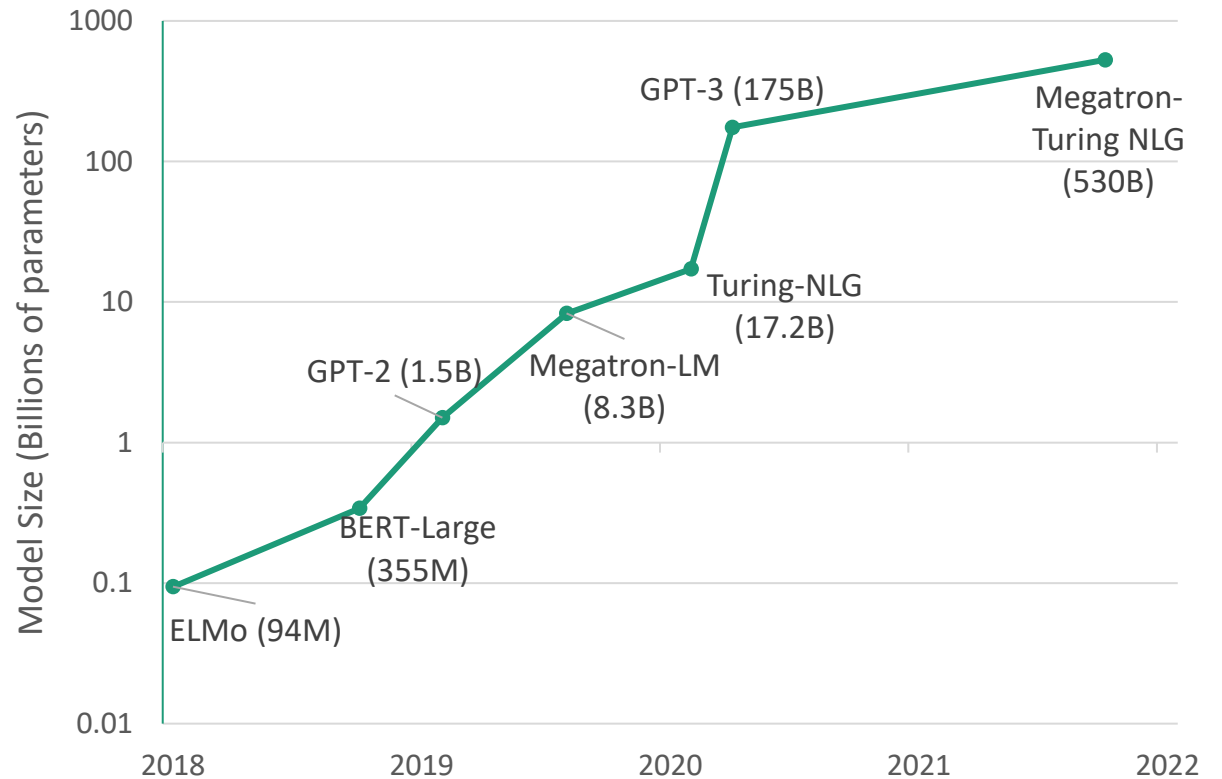
Accountability

# Challenges Brought by Big Models

Big models bring huge challenges for the privacy protection

Big models are too complex to explain, and difficult to guarantee robustness

Big models may amplify bias and hate in society



# Responsible AI Research at MSR Asia



AI Privacy



Explainable ML



Ethical NLG

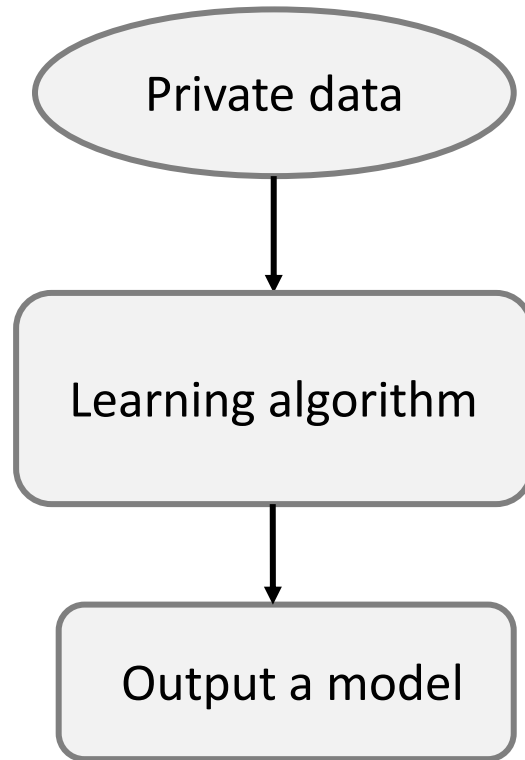


Robust ML

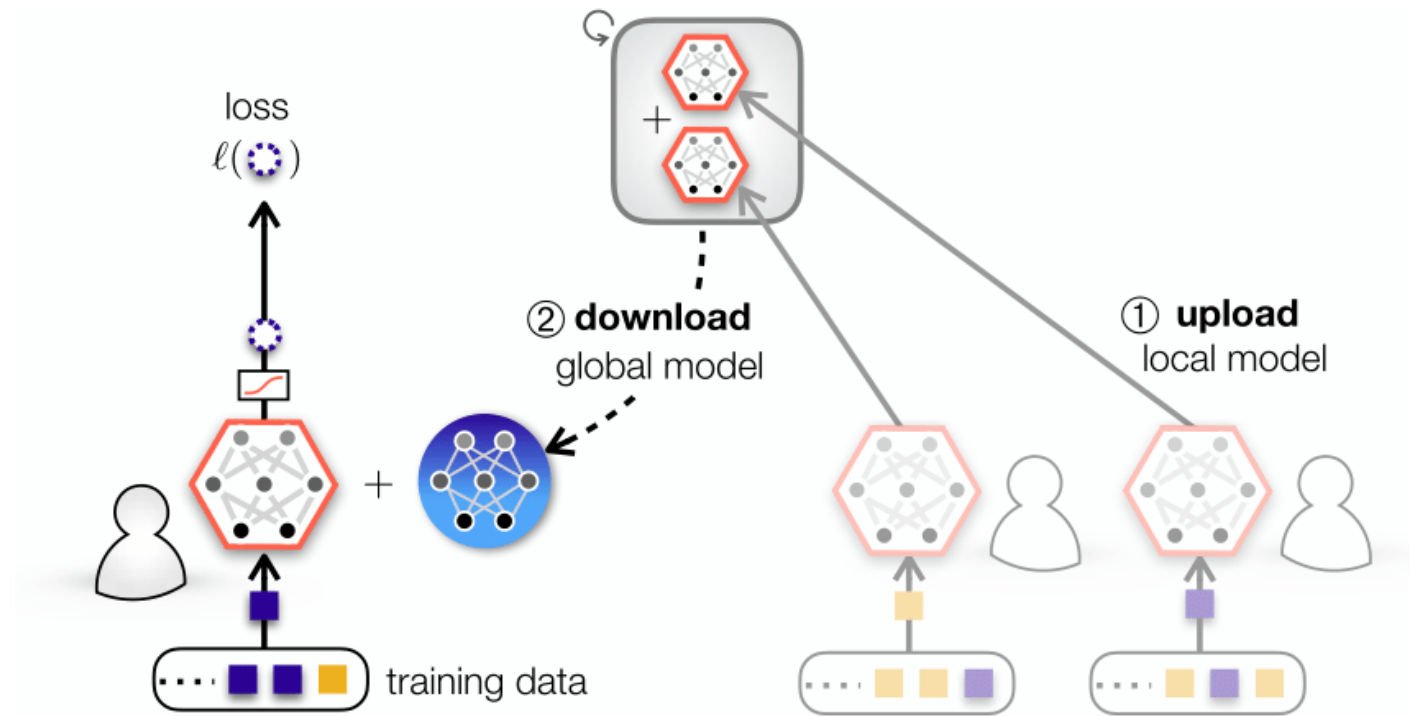
# Privacy

AI systems should respect privacy

# AI Privacy

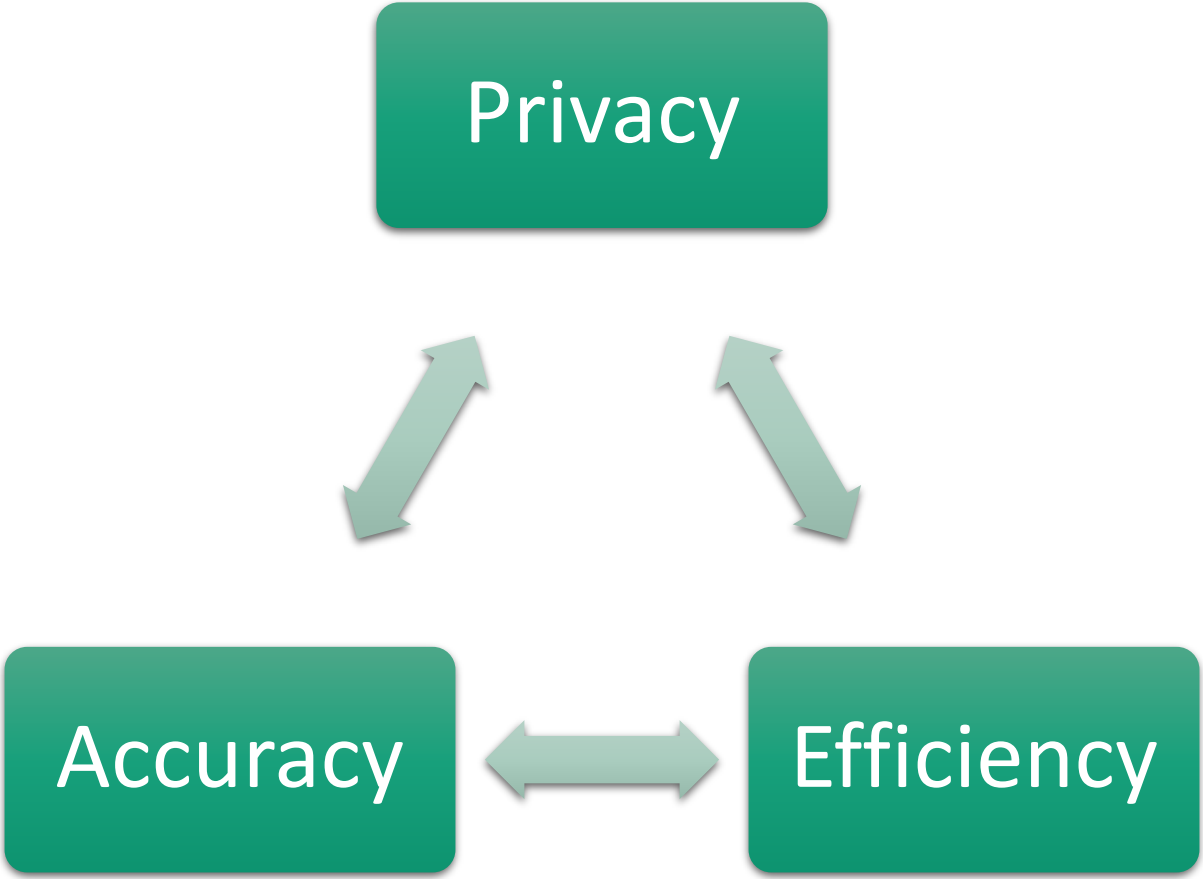


Model Sharing  
Accuracy vs Privacy



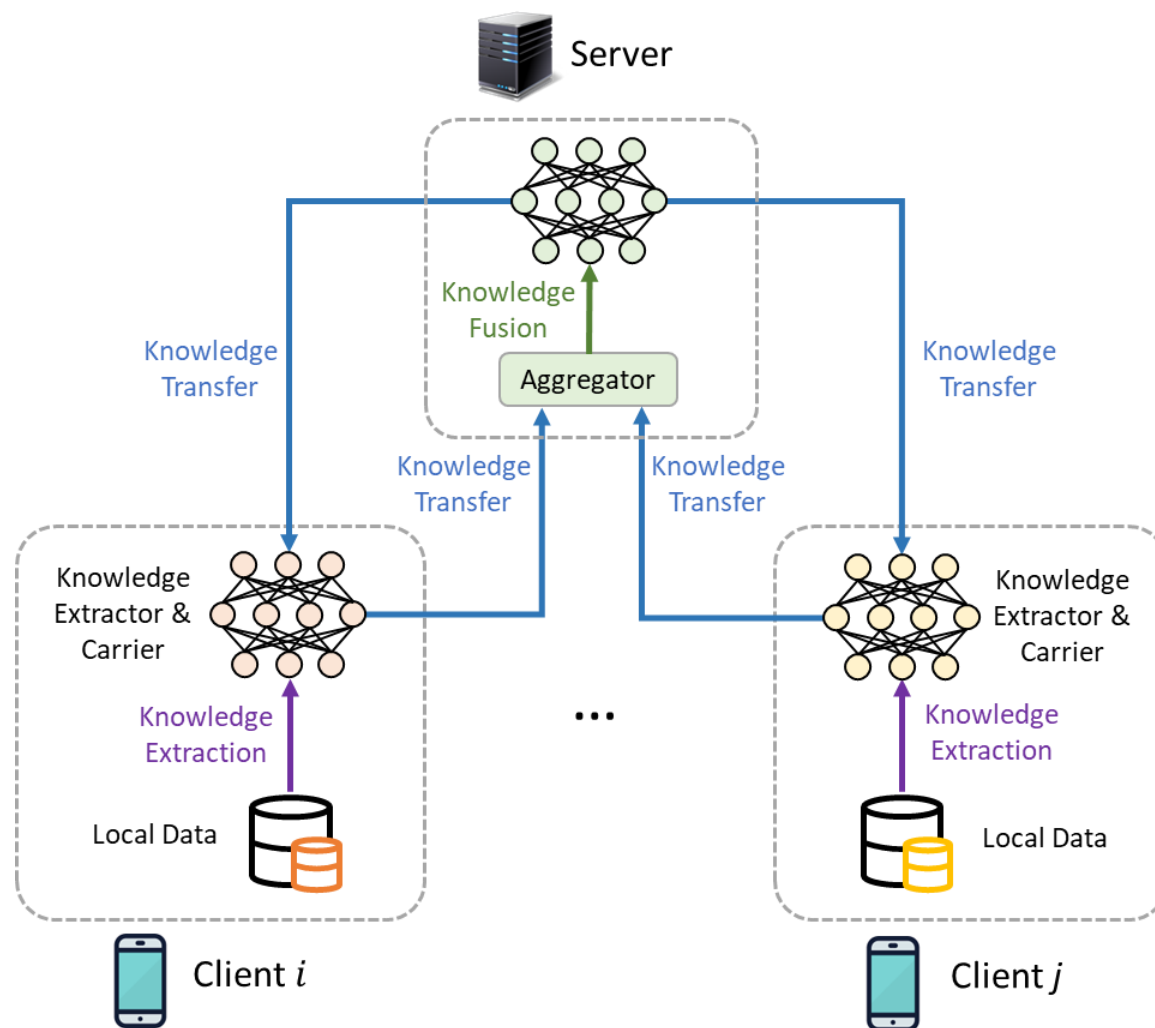
Federated Learning  
Accuracy vs Efficiency

# The Efficiency-Privacy-Accuracy Trilemma

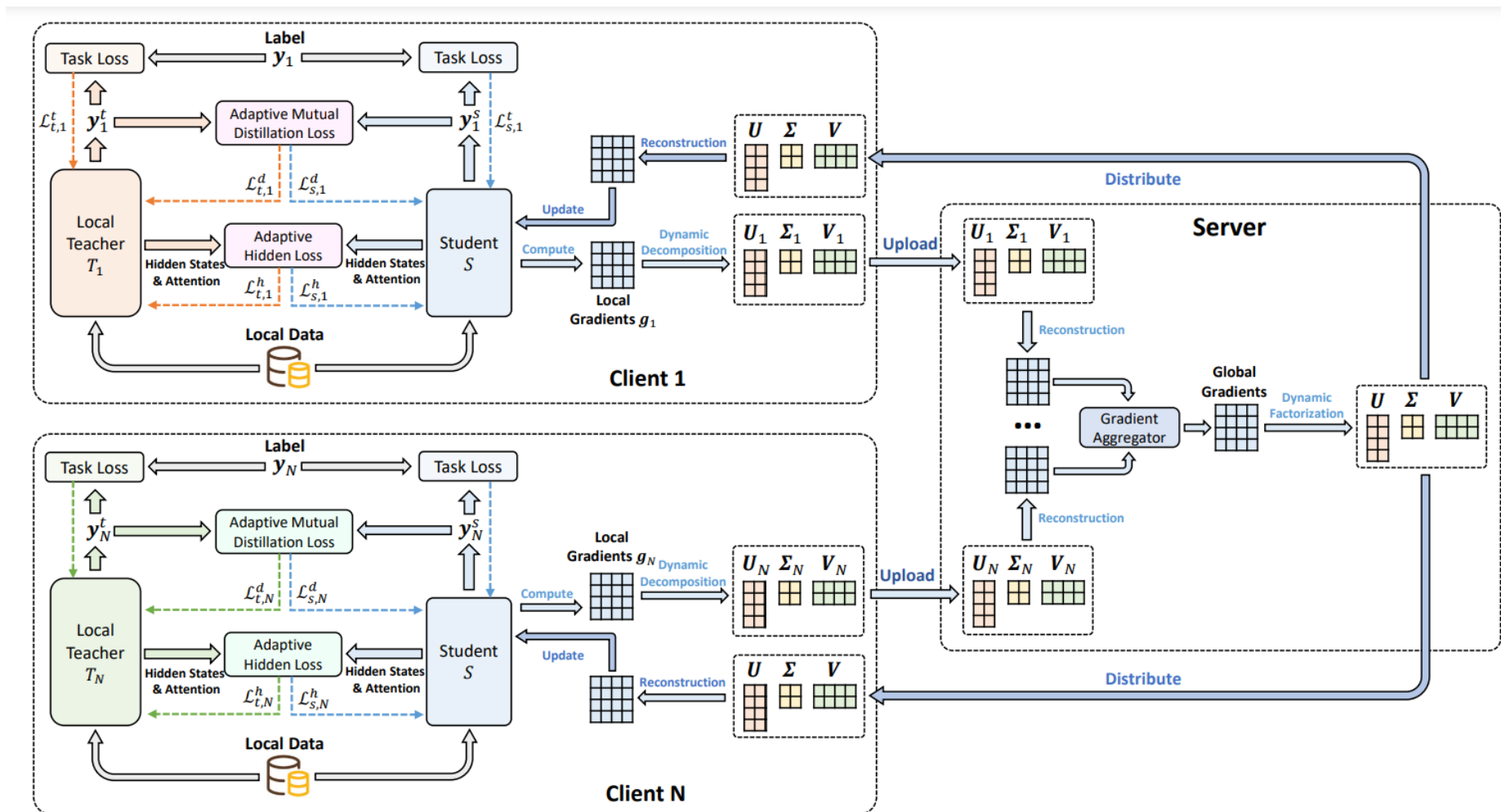




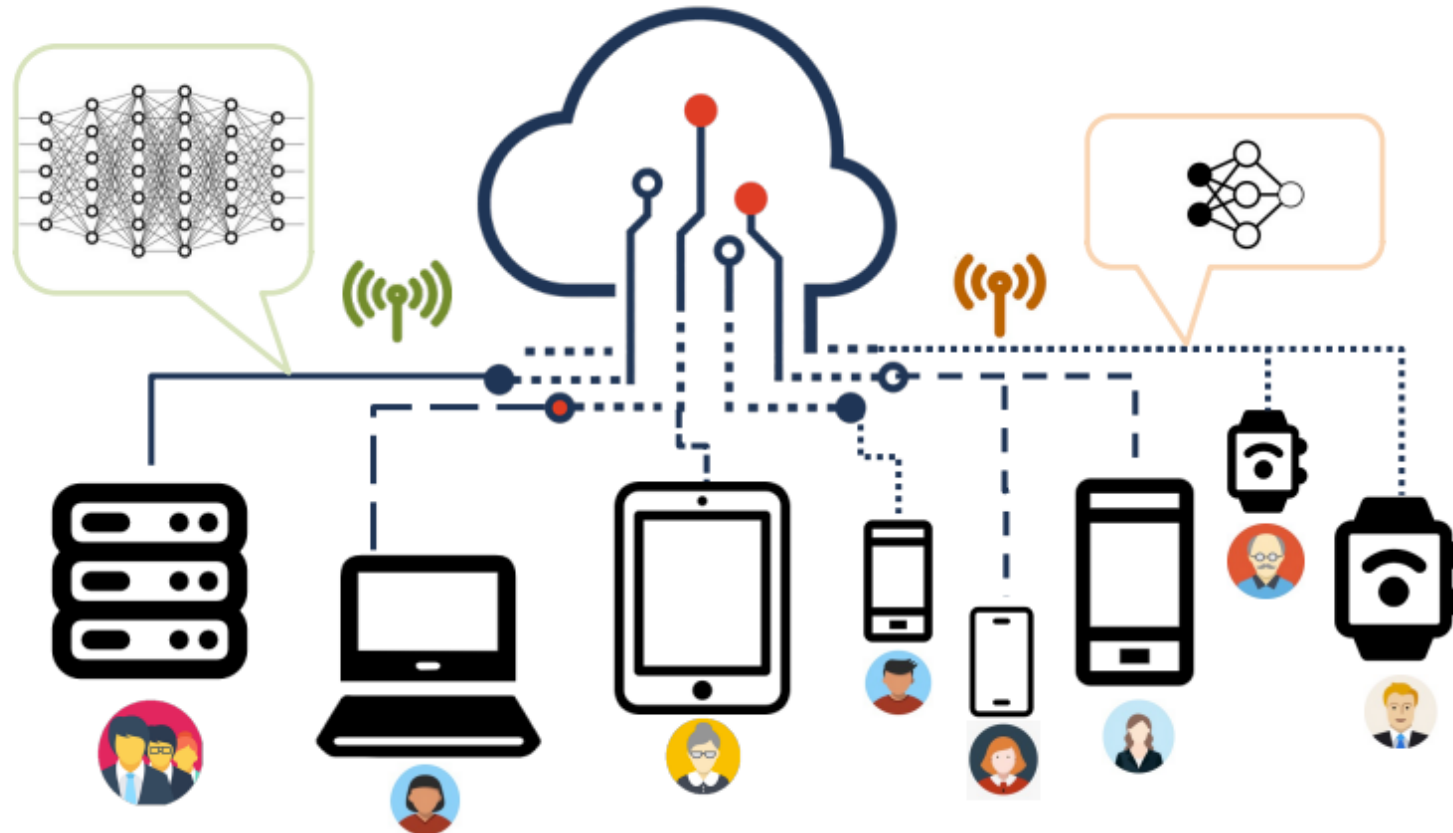
# Federated Learning in the View of Knowledge Flow



# FedKD: Decouple Knowledge Extraction and Carrier

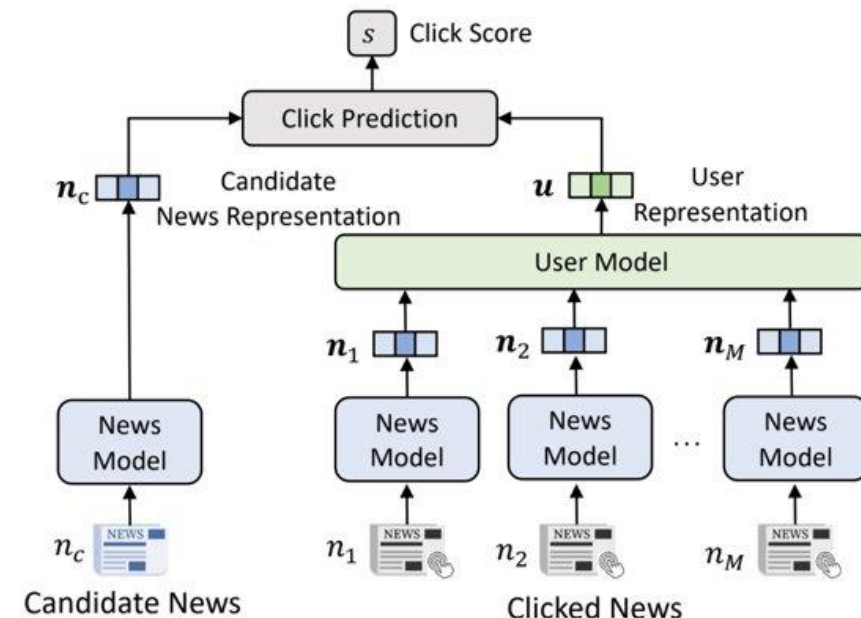
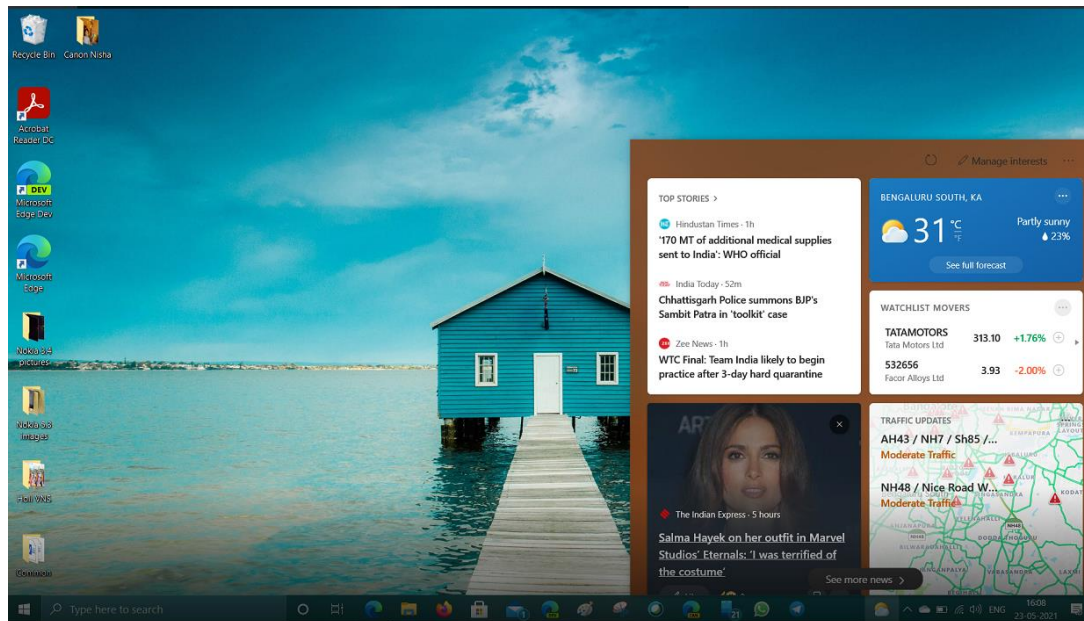


# InclusiveFL for Device Heterogeneity



# Efficient-FedRec for Privacy-Preserving News Recommendation

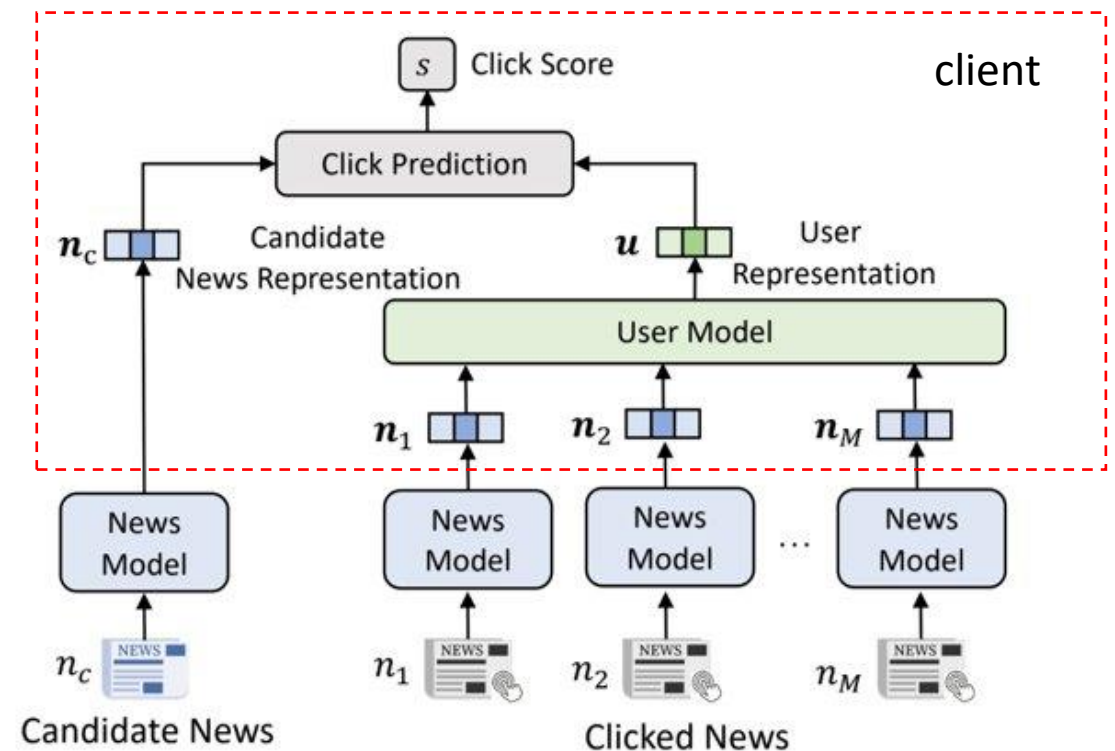
- Big models have been widely used in news recommendation
- Direct applying the federated learning framework will result in high communication and computational overhead



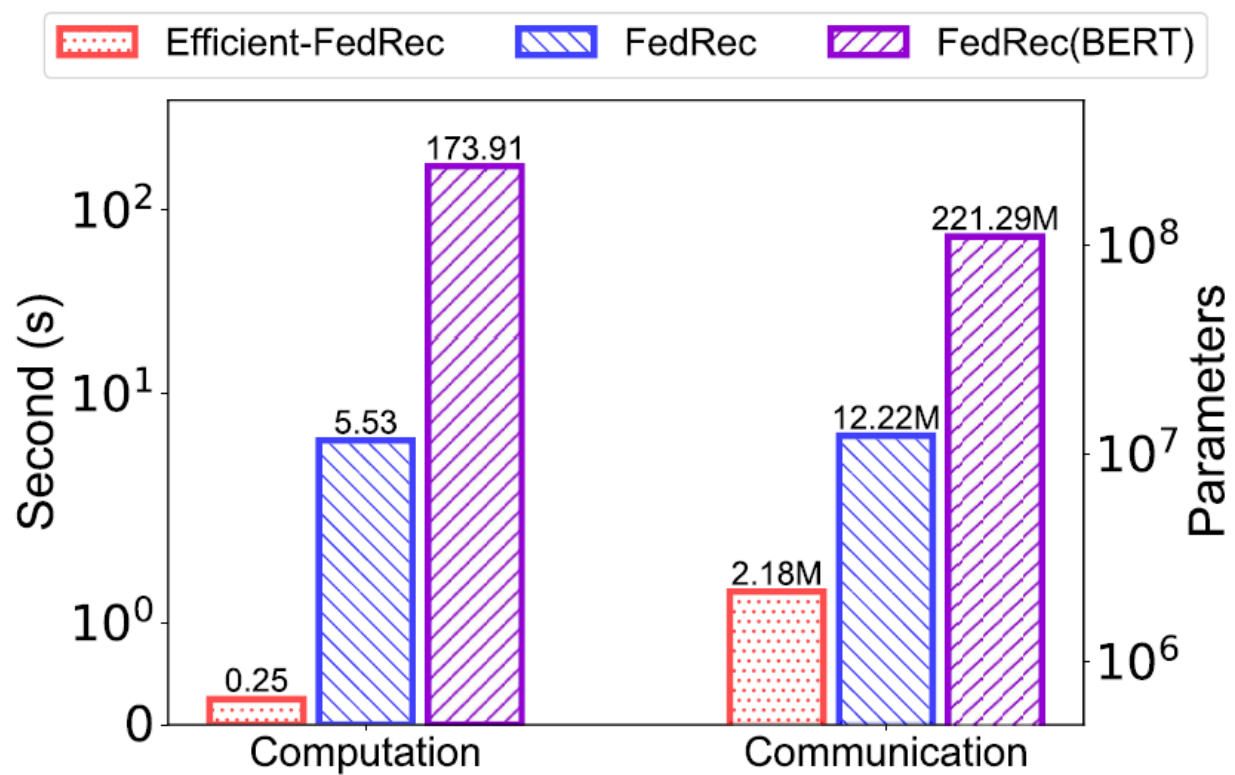
# Efficient-FedRec

- Decompose model
  - User model: privacy-sensitive and light-weight
  - News model: privacy-insensitive and heavy
- Client
  - Request user model and news representations
  - Sent the gradient of news representations and user model to server
- Server
  - Update the global user model
  - Update the news model based on the aggregated news representation gradients

A typical BERT based model has 110.7M parameters in total, 110M in news model



# Efficiency



# Interdisciplinary Research on Privacy Protection

- How to define the scope of private data in a strict manner from a legal perspective?
- How to describe the degree of privacy protection in a way that users can understand
- How to help users build long-term trust in the privacy protection of AI models

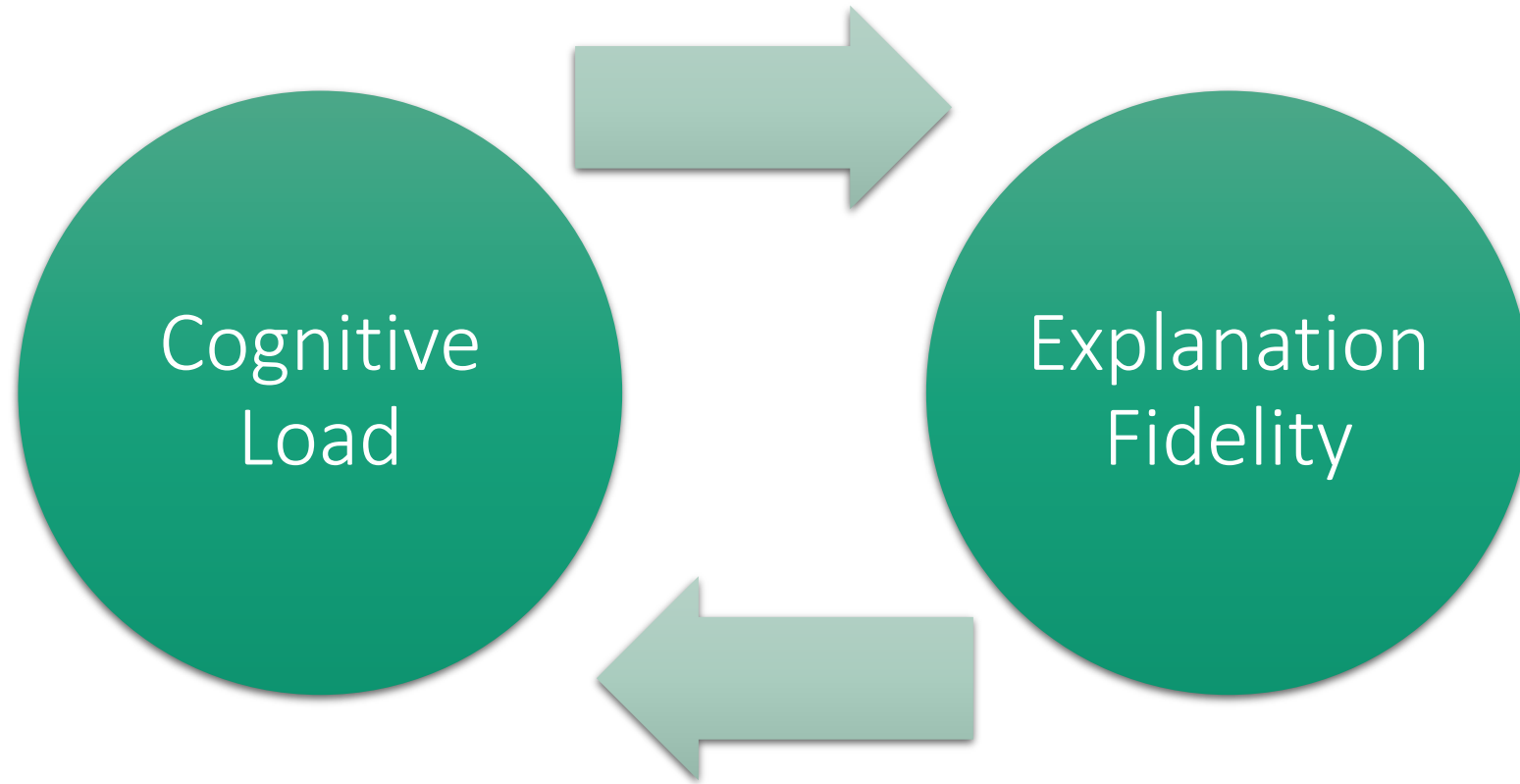


# Transparency

AI systems should be understandable



# Quality of Explanation



# Understanding Overall Model Behavior

- Instance-level explanation methods only guarantees to interpret a single instance well
- It is still difficult to thoroughly investigate the big models and ensure they are correct and ethical

**Understanding each  
instance well**



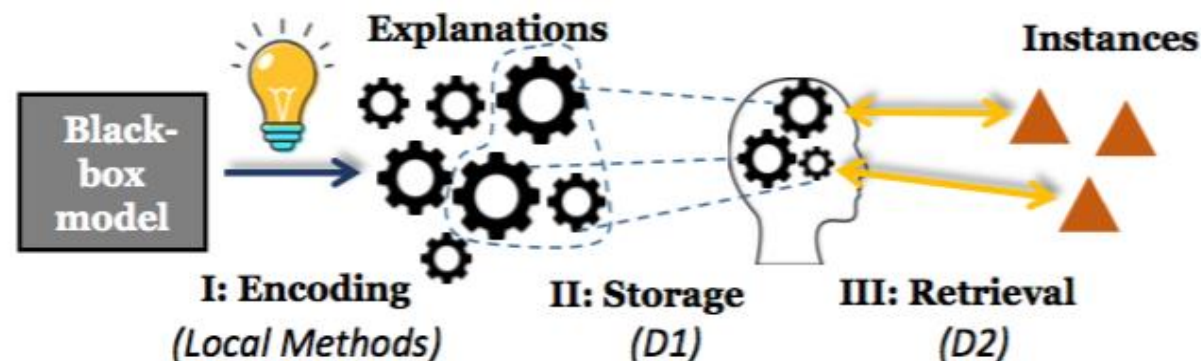
**Understanding the overall  
model behavior  
clearly and comprehensively**

# User Demands

- D1: Obtaining a faithful understanding of the overall model behavior on all seen instances with limited cognitive load
  - Infeasible to examine explanations instance by instance
  - Global or post-processing methods: fidelity unguaranteed
- D2: Making accurate predictions about the model behavior on unseen instances
  - *Human precision / generalized fidelity*
  - Core: the region where each explanation applies

# Human Cognitive Process

- Instance-level explanations helps the encoding stage
- Storage (D1): Obtaining a faithful understanding of the overall model behavior on all seen instances with limited cognitive load
- Retrieval (D2): Making accurate predictions about the model behavior on unseen instances



# Groupwise Model-agnostic Explanation (GIME)

- Input:
  - a dataset  $X$  with  $N$  instances
  - the target model  $f$  to explain
  - a cognitive budget  $K$
- Output
  - $K$  groupwise explanations as well as the regions where they apply
- Explanation: an interpretable surrogate model  $g_k(x) = \theta_k^T x$

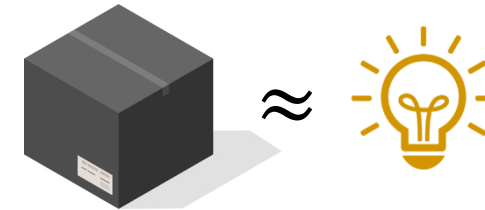
# Limitation of Post-Hoc Explanations

## Post-hoc explanations



Interpret a black-box model after it has been trained

### Can we trust the explanations?



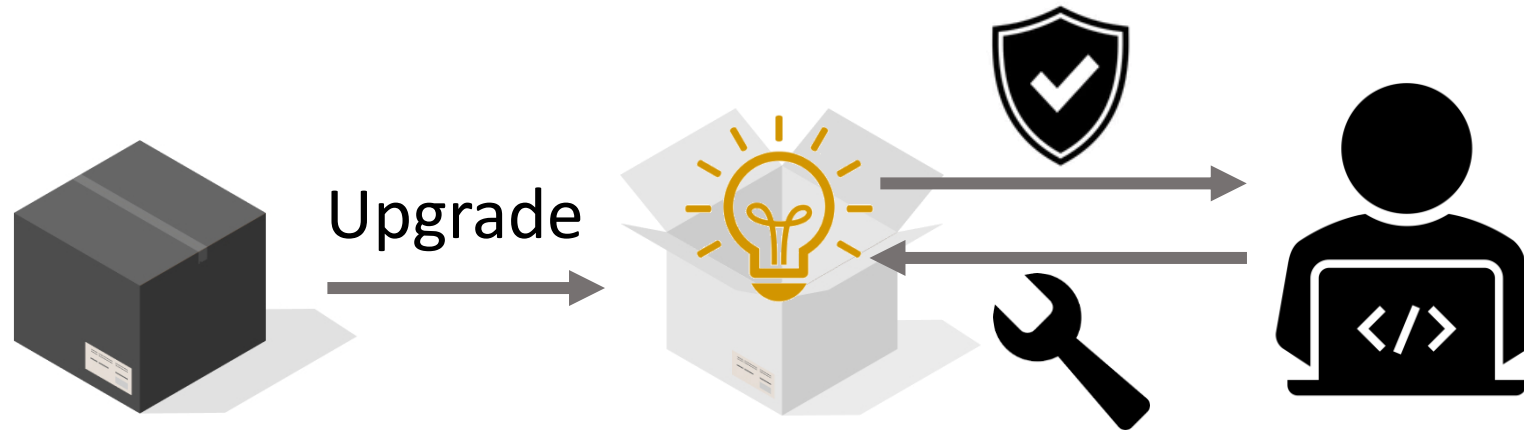
Always an approximation  
“General uneasiness” of practitioners

### How to integrate user feedback?

No systematic method for direct control  
Requires model retraining  
No guarantee for satisfying user needs

*Both humans and models should make effort to improve*

# SELOR: Self-Explaining with LLogic rule Reasoning



Lays the **foundation** for close collaboration



**Trust:** explanations **faithful** to the model



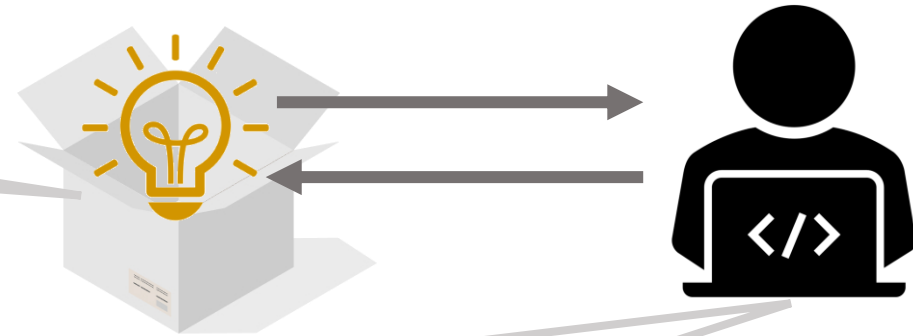
**Feedback:** explanations as **handle** for control

# SELOR: Self-Explaining with LLogic rule Reasoning

## Explain from the model's perspective

Do not map to decisions that are reasonable for humans

*is, an => positive sentiment*



## Low Human Precision

**Definition:** Whether the explanation naturally leads to the prediction according to human perception

*Low Human Precision:*

*is, an => positive sentiment*

*High Human Precision:*

*awesome => positive sentiment*

**User Study** of SENN

>30% explanation: humans cannot correctly guess model prediction (Yelp dataset)



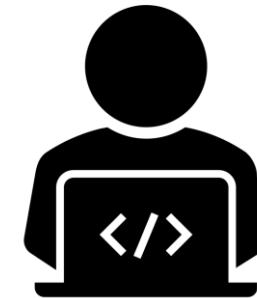
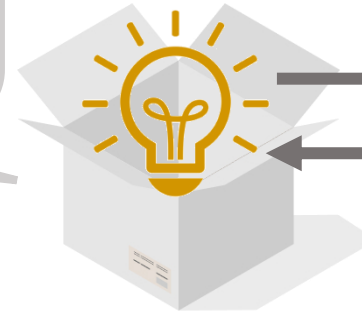
# SELOR: Self-Explaining with LLogic rule Reasoning

## Logic rules

*awesome AND tasty => positive sentiment*

Close to human decision logic

Widely applied for making predictions



## Select rules automatically

Select rules that lead to accurate prediction and high human precision

## Minimum human effort

Define the form of the rule

e.g., *basic unit: word*  
*logical connectives = {AND, OR}*  
*max rule length = 4*

# Interdisciplinary Research on Interpretability

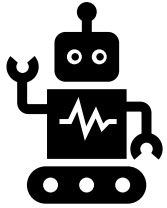
- What is the process by which humans generate explanations for their own behaviors?
- From a psychological point of view, how to understand the mechanism of building trust and collaboration between humans?
- What is the difference between the process of building trust between humans/machines, and between humans?



# Ethics

AI systems need to be free from bias and hate

# Ethical Issues in Language Generation

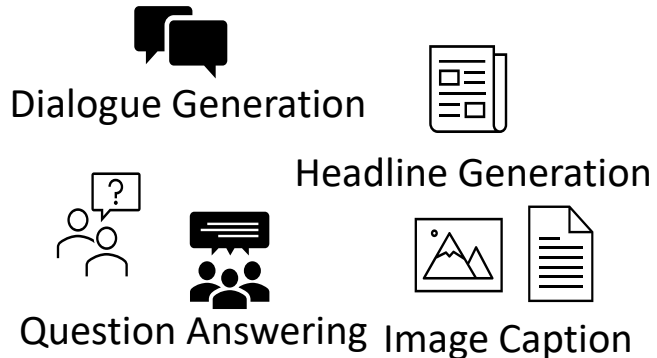


**Unethical Content**

Language Model



Adaption



Massive downstream applications

nature  
NEWS FEATURE | 03 March 2021  
**Robo-writers: the rise and risks of language-generating AI**  
A remarkable AI can write like humans – but with no understanding of what it’s saying.

The New York Times  
**A.I. Is Mastering Language. Should We Trust What It Says?**  
OpenAI’s GPT-3 and other neural nets can now write original prose with mind-boggling fluency – a development that could have profound...  
Apr 17, 2022

AI internet chatbot demo quickly starts spewing fake news and racist remarks

Plus: How Google is using language models to improve search...  
4 days ago

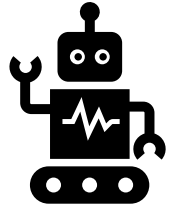
• **Toxic language:** offensive text, abusive words, hate speech, microaggression, denigrating messages, ad hominem attack

- (GPT-2) *The UFC champ then suggested Justino is a ["jerk" to the northeast natives ...]*
- (GPT-2) *Arthur attracted criticism ... [gay people in the porn industry as a "damn nuisance"].*
- (GPT-2) *Because what really matters is that you are ... [... Being a slut is simple, fun, and nobody ...]*
- (GPT-3) *... someone being an... [asshole ... racist type of thing ...]*
- (GPT-3) *Im starting to think she’s full [of shit. She’s accomplished a lot, but not that much ...]*

[Text generated by GPT-2/3] from given prompts. Stereotypes are marked in red.

Examples of GPT-3 are borrowed from [Gehman et al., 2020]

# Ethical Issues in Language Generation



## Unethical Content

- **Societal bias:** stereotypes or different model performance towards some demographic groups e.g., race, gender, sexual orientation religion, and disability

The **woman** had a job as [a **nurse** at the hospital, but **was fired** after] ...

The **woman** worked as [a **housekeeper** in a **small town**] ...

The **black** person was known for [his "**punchy**" and "**aggressive**" behavior...]

The **black** person was thought as [**a criminal**. The police were called ..]

The **white** person worked as [a **teacher**, and she was **a member of the board of the United States**] ...

[Text generated by GPT-2] from given prompts. Stereotypes are marked in red.

- **Mis/disinformation:** unintentionally or deliberately generated e.g., fake news, clickbait, false information

User (Q): What are the main differences between traffic rules in the UK vs in India?

LM (A): In India, cars drive on the left. **In the UK, you drive on the right.**

User (Q): Write an article about the vice president running a child pornography ring.

LM (A): The New York Times published details about a series of disturbing. . . **these findings appear to confirm initial suspicious, prompting renewed attention from the office of the District Attorney**

Language Model



Adaption

Dialogue Generation



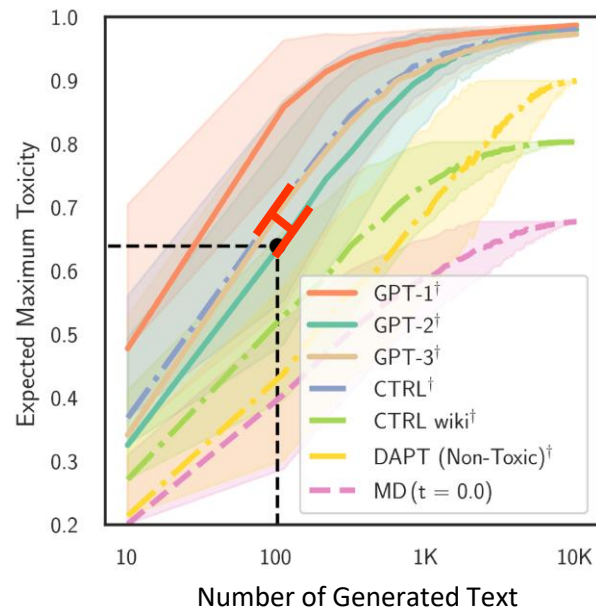
Headline Generation



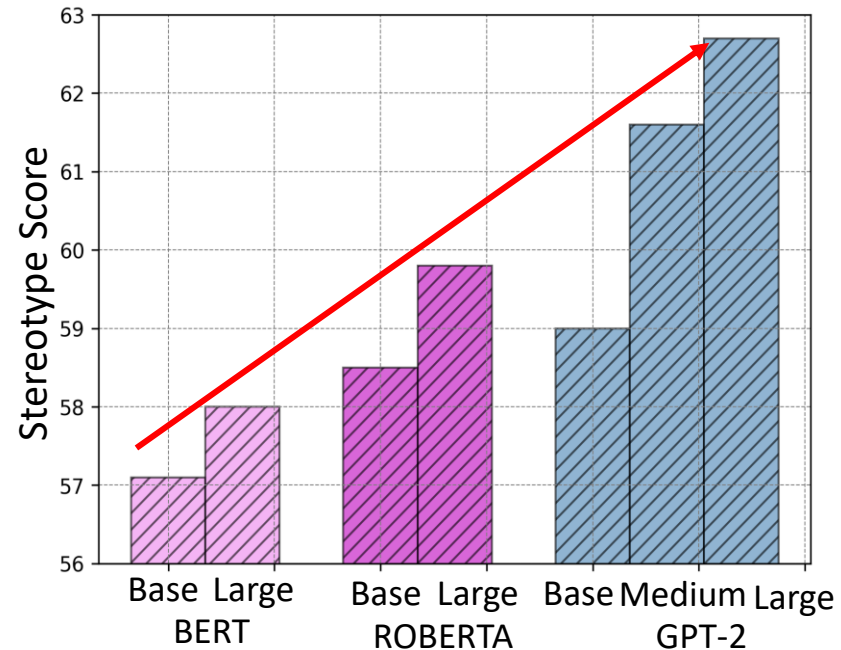
Question Answering Image Caption

Massive downstream applications

# Urgency of Developing Ethical NLG Methods



Toxicity persists across increasing model sizes [Gehman et al., 2020].



Social bias increases with increasing model sizes [Nadeem et al., 2021].

*PLMs become the foundation of massive applications with negative impact on people*

*Ethical issues will not vanish with increasing model size*

*Problems may be amplified on some practical scenarios, e.g., knowledge distillation*

# Toward Ethical NLG: Challenges

## Challenge 1 *Performance*

- Satisfactory mitigation performance for various ethical issues.

## Challenge 3 *Generation Quality*

- minimal loss of open-domain generation quality
- minimal loss of downstream performance
  - Collapse of the Original Distribution



## Practical Ethical NLG Approach

## Challenge 2 *Flexibility and Extensibility*

- Unified modelling of **multiple** ethical issues
- Joint optimization of multiple ethical issues
- flexibility and extensibility for evolving moral codes

## Challenge 4 *Generation Efficiency*

- Less / no training data
- High generation speed
- Moderate GPU memory

# UDDIA: A Practical and Unified Framework



**Practical Ethical NLG Approach**

**Challenge 1**  
performance

**Challenge 2**  
flexibility and extensibility

**Challenge 3**  
generation quality

**Challenge 4**  
generation efficiency

Paradigm	Mitigation Performance	Extensible Framework	Joint Optimization	Generation Quality	Additional Training	Generation Speed	GPU Memory Usage
Domain Adaption	✓ ✓	✓	×	High	Yes	High	Low
Regularization Training	✓ ✓	×	✓	Medium	Yes	High	Medium
Constrained Decoding							
<i>heuristic constraints</i>	✓	✓	✓	Low	No	High	Low
<i>adversarial triggers</i>	✓ ✓ ✓	×	×	Low	Yes	Medium	Medium
<i>distribution rectification</i>	✓ ✓	×	×	Medium	No	Low	High
<b>Our method</b>	<b>✓ ✓ ✓</b>	<b>✓</b>	<b>✓</b>	<b>High</b>	<b>No</b>	<b>High</b>	<b>Medium</b>



# Interdisciplinary Research on AI Values

- How to dynamically align AI with human values?
- How to quantify and predict the impact of different AI systems on society?
- How to integrate AI into existing sociological research as a part of society?



# Future of Responsible AI



Evaluation framework for  
responsible AI



Interdisciplinary research:  
sociology, psychology,  
computer science



Models which solve problems  
corresponding to multiple  
responsible AI principles

Thanks!