



# Data Mining

## Course Overview and Logistics

<https://data-mining.github.io/winter-2026/>

CS 453/553 – Winter 2026

Yu Wang, Ph.D.

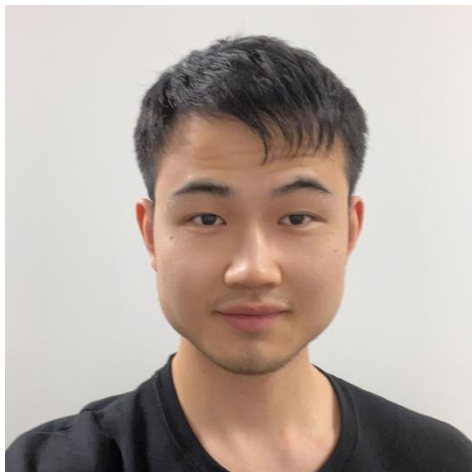
Assistant Professor

Computer Science

University of Oregon



# Self-Introduction




**Yu (Jack) Wang**  
**(You)**

<https://yuwang0103.github.io/>

## Research Interests:

- Data Mining and Machine Learning
- Neural-Symbolic Learning
- Graph and Network
- LLM + Structured Knowledge
- AI/ML/DM Applications
  - Document Intelligence
  - Social Computing
  - Networking Physical Infrastructure



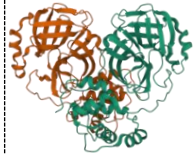
 **Recruiting Ph.D. students and interns!** I am actively seeking highly motivated students for Ph.D. or Research intern positions. Please feel free to email me your CV, transcripts, and brief descriptions about why you want to work with me if you are interested!

**Contact:**  
[yuwang@uoregon.edu](mailto:yuwang@uoregon.edu)

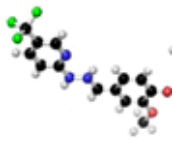


# What is Data?

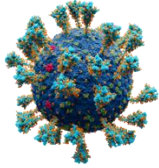
## Science



Protein



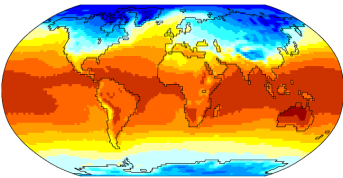
Small Molecule



Virus



Brain Neural



Surface Temperature of Earth

## Gas Network

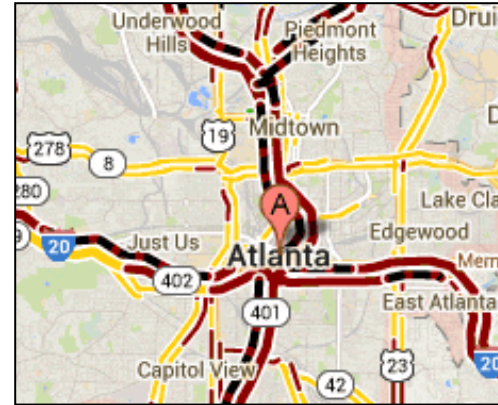


## Power Network

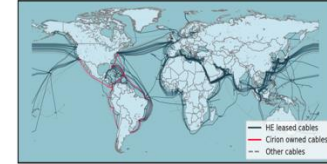


## Infrastructure

### Transportation Network



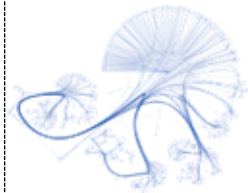
## Submarine Cable



## Terrestrial Cable



## Social Network



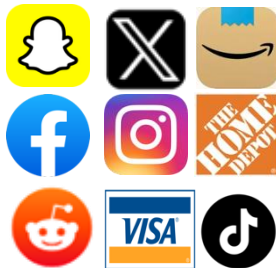
Citation Network



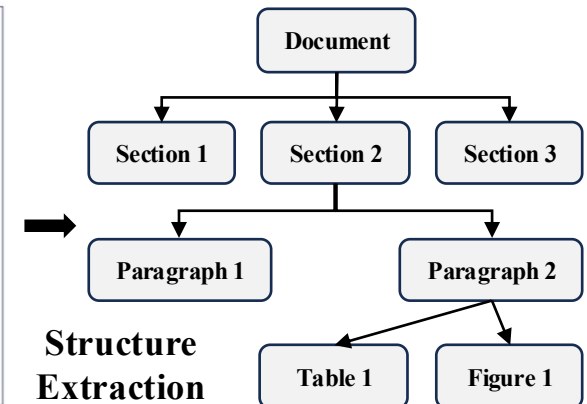
Transaction Network



User-Entity Interaction Graph



## Document



Structure  
Extraction

Virtual Village with AI Agents





# Why Analyze Data? – Paper Management

Google Scholar

☒ Articles ☐ Case law

New! Scholar Labs: An AI Powered Scholar Search

## Recommended articles



☆ Analyzing the Properties of Graph Neural Networks with Evolutionary Algorithms  
Z Liu, Z Lu, H Wang, D Chen, S Wang, J Chu, R Gao, A Jiang  
Mathematics - 3 days ago [HTML](#)

☆ Self-Supervised Bipartite Graph Neural Networks with Missing Value Imputation for Small Tabular Data Predictions  
PC Liu, CT Li  
ACM Transactions on Intelligent Systems and Technol... - 4 days ago [PDF](#)

[More articles from 4 days ago](#)

☆ HRGNN: Learning Holistically Robust Graph Neural Networks on Noisy Graphs with Label Scarcity  
JW Chiu, CT Li  
ACM Transactions on Intelligent Systems and Technol... - 5 days ago [PDF](#)

## REFERENCES

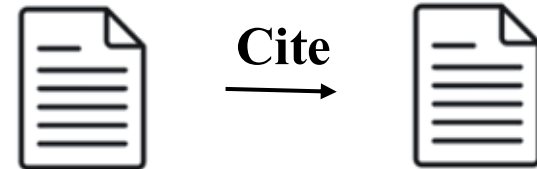
Art of Problem Solving. Aime problems and solutions, 2025. URL [https://artofproblemsolving.com/wiki/index.php/AIME\\_Problems\\_and\\_Solutions](https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions). 8, 22

Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang, Jeff Z Pan, Wen Zhang, Huajun Chen, Fan Yang, et al. ReSearch: Learning to reason with search for llms via reinforcement learning. *arXiv preprint arXiv:2503.19470*, 2025. 2, 4, 7, 10, 21

Zihao Cheng, Hongru Wang, Zeming Liu, Yuhang Guo, Yuanfang Guo, Yunhong Wang, and Haifeng Wang. ToolSpectrum: Towards personalized tool utilization for large language models. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 20679–20699, 2025. 10

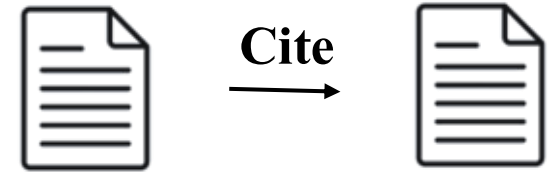
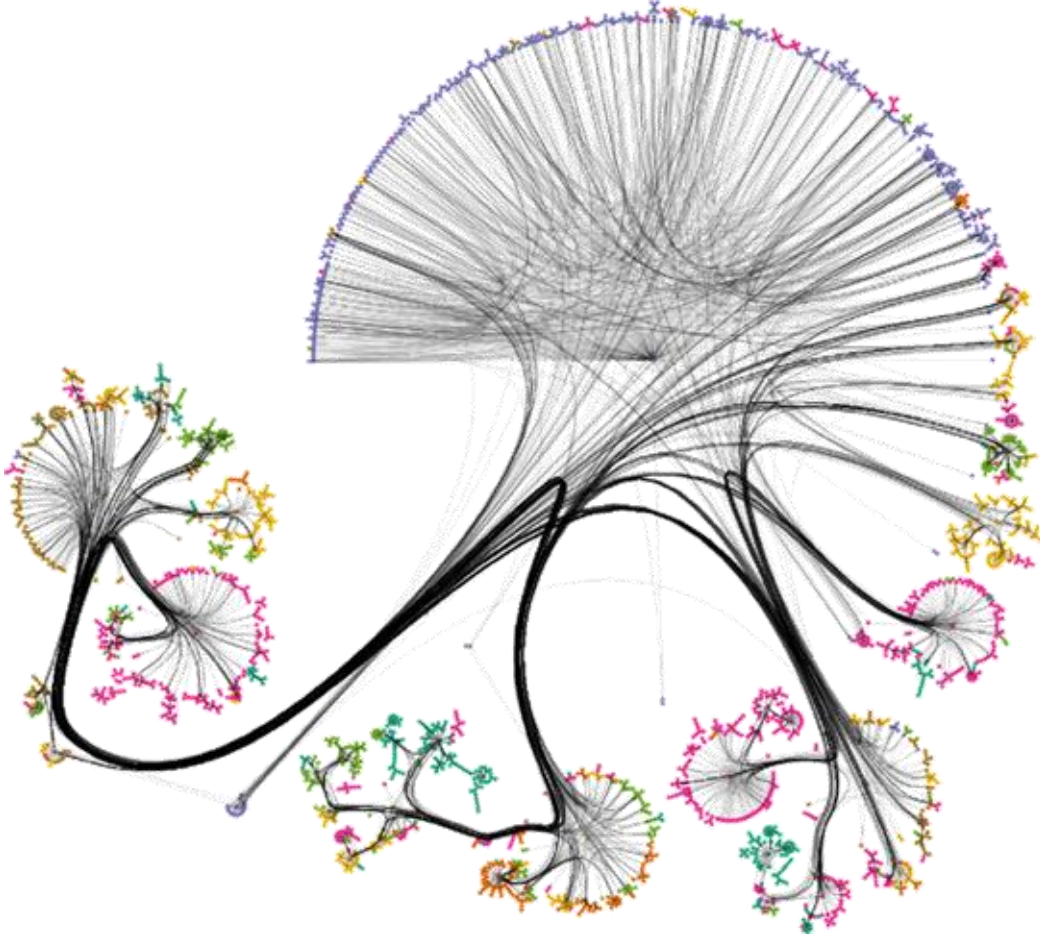
Yingfan Deng, Anhao Zhou, Yuan Yuan, Xian Zhang, Yifei Zou, and Dongxiao Yu. Pe-ma: Parameter-efficient co-evolution of multi-agent systems. *arXiv preprint arXiv:2506.11803*, 2025. 11

Guanting Dong, Yifei Chen, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Yutao Zhu, Hangyu Mao, Guorui Zhou, Zhicheng Dou, and Ji-Rong Wen. Tool-star: Empowering llm-brained multi-tool reasoner via reinforcement learning. *arXiv preprint arXiv:2505.16410*, 2025. 2, 10





# Why Analyze Data? – Paper Management



$$\frac{\sum_{e_{ij} \in \mathcal{E}} 1[y_i == y_j]}{|\mathcal{E}|}$$

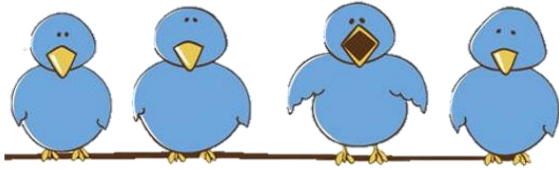
$\mathcal{E}$  - Total Number of Edges

$e_{ij}$  - Edge between node  $i/j$

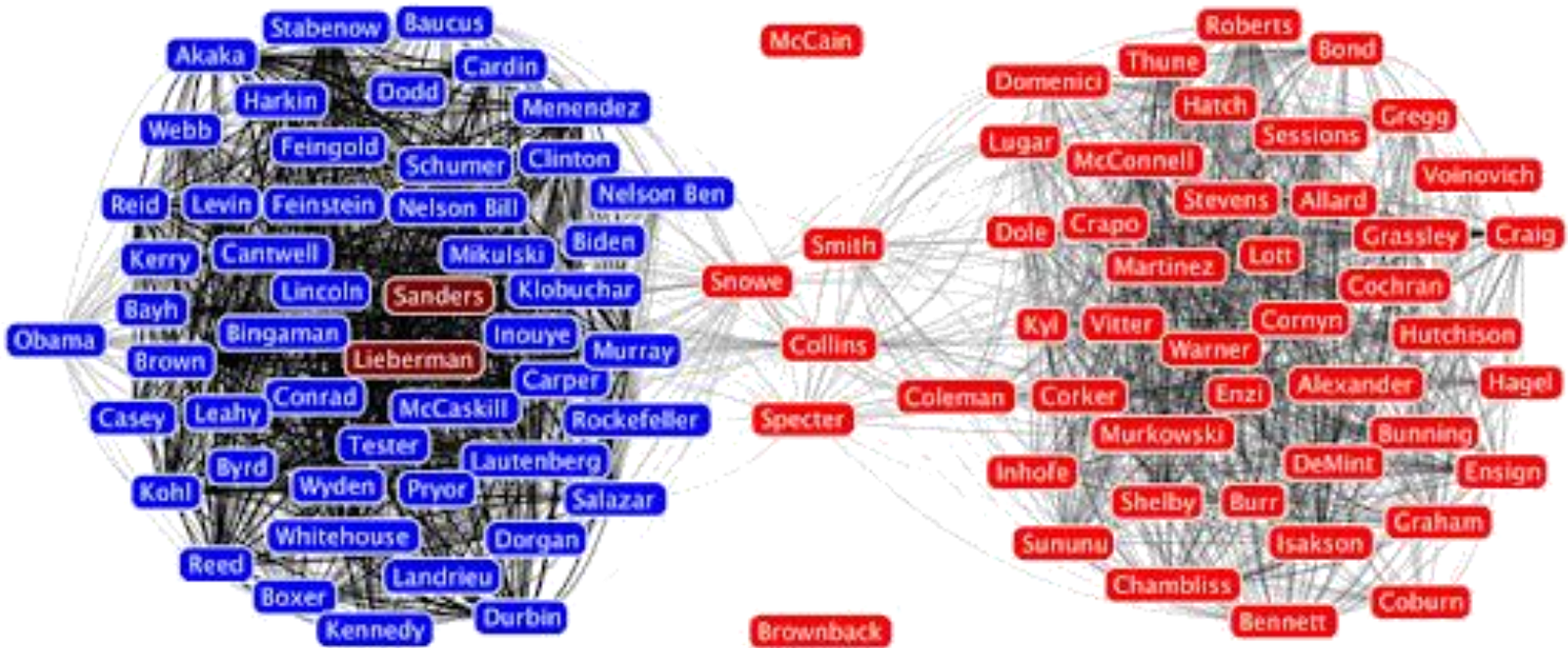
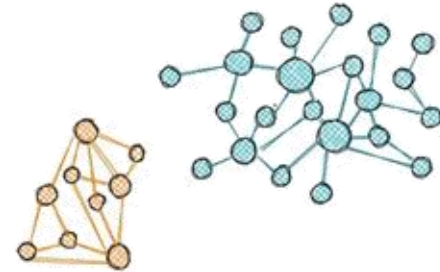
$y_i$  - Label of  $i$



# Why Analyze Data? – Paper Management



Birds of a feather flock together

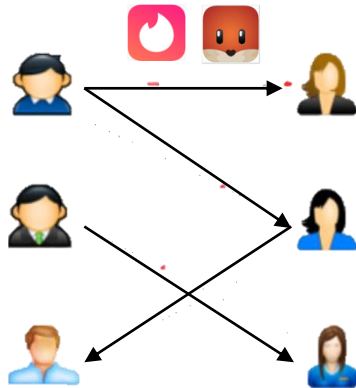


Gun Control Belief Network

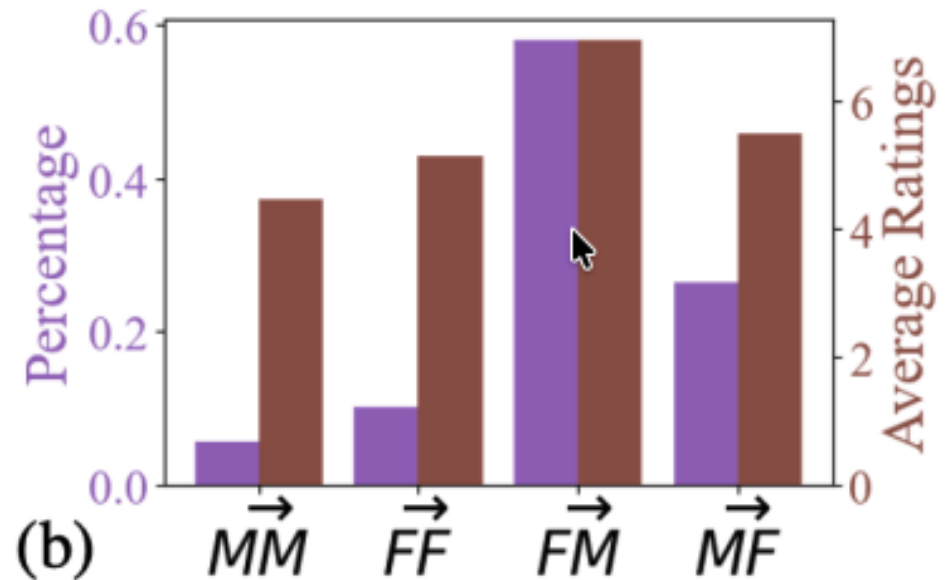
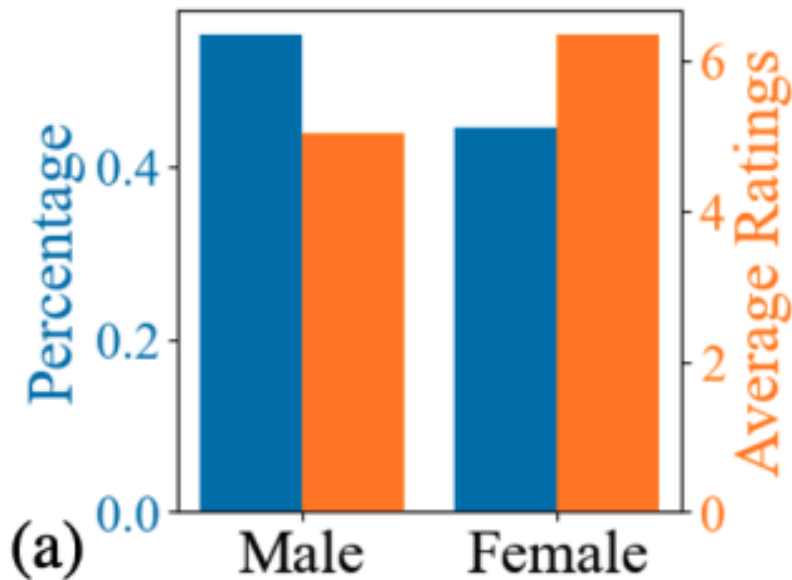
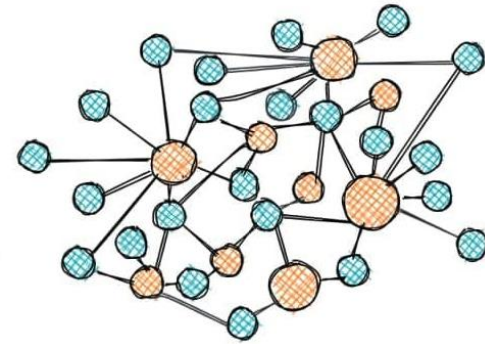




# Why Analyze Data? – Paper Management



Dating Network



Dating Network



# Why Analyze Data? – Paper Management

## IN-THE-FLOW AGENTIC SYSTEM OPTIMIZATION FOR EFFECTIVE PLANNING AND TOOL USE

Zhuofeng Li<sup>\*1,2</sup>, Haoxiang Zhang<sup>\*1,3</sup>, Seungju Han<sup>1</sup>, Sheng Liu<sup>1</sup>, Jianwen Xie<sup>4</sup>, Yu Zhang<sup>2</sup>, Yejin Choi<sup>1</sup>, James Zou<sup>1†</sup>, Pan Lu<sup>1†</sup>

<sup>1</sup>Stanford University, <sup>2</sup>Texas A&M University, <sup>3</sup>UC San Diego, <sup>4</sup>Lambda



Website: <https://agentflow.stanford.edu>

Code Model Demo Visualize

### ABSTRACT

Outcome-driven reinforcement learning has advanced reasoning in large language models (LLMs), but prevailing tool-augmented approaches train a single, monolithic policy that interleaves thoughts and tool calls under full context; this scales poorly with long horizons and diverse tools and generalizes weakly to new scenarios. Agentic systems offer a promising alternative by decomposing work across specialized modules, yet most remain training-free or rely on offline training decoupled from the live dynamics of multi-turn interaction. We introduce AGENT-FLOW, a trainable, *in-the-flow* agentic framework that coordinates four modules (planner, executor, verifier, generator) through an evolving memory and directly optimizes its planner inside the multi-turn loop. To train on-policy in live environments, we propose *Flow-based Group Refined Policy Optimization* (Flow-GRPO), which tackles long-horizon, sparse-reward credit assignment by converting multi-turn optimization into a sequence of tractable single-turn policy updates. It broadcasts a single, verifiable trajectory-level outcome to every turn to align local planner decisions with global success and stabilizes learning with group-normalized advantages. Across ten benchmarks, AGENTFLOW with a 7B-scale backbone outperforms top-performing baselines with average accuracy gains of 14.9% on search, 14.0% on agentic, 14.5% on mathematical, and 4.1% on scientific tasks, even surpassing larger proprietary models like GPT-4o. Further analyses confirm the benefits of in-the-flow optimization, showing improved planning, enhanced tool-calling reliability, and positive scaling with model size and reasoning turns.

In-the-flow agentic system optimization for effective planning and tool use

☐ Search within citing articles

### Latent collaboration in multi-agent systems

[J.Zou](#), [X.Yang](#), [R.Qiu](#), [G.Li](#), [K.Tieu](#), [P.Lu](#), K Shen... - arXiv preprint arXiv ..., 2025 - arxiv.org

Multi-agent systems (MAS) extend large language models (LLMs) from independent single-model reasoning to coordinative system-level intelligence. While existing LLM agents ...

☆ Save Cite Cited by 4 Related articles All 2 versions

### Adaptation of agentic ai

[P.Jiang](#), [J.Lin](#), [Z.Shi](#), [Z.Wang](#), L He, Y Wu... - arXiv preprint arXiv ..., 2025 - arxiv.org

Cutting-edge agentic AI systems are built on foundation models that can be adapted to plan, reason, and interact with external tools to perform increasingly complex and specialized ...

☆ Save Cite Cited by 2 Related articles All 2 versions

### DeepAgent: A General Reasoning Agent with Scalable Toolsets

[X.Li](#), [W.Jiao](#), [J.Jin](#), [G.Dong](#), [J.Jin](#), Y Wang... - arXiv preprint arXiv ..., 2025 - arxiv.org

Large reasoning models have demonstrated strong problem-solving abilities, yet real-world tasks often require external tools and long-horizon interactions. Existing agent frameworks ...

☆ Save Cite Cited by 2 Related articles All 2 versions

### The Path Not Taken: RLVR Provably Learns Off the Principals

[H.Zhu](#), [Z.Zhang](#), [H.Huang](#), DJ Su, [Z.Liu](#), [J.Zhao](#)... - arXiv preprint arXiv ..., 2025 - arxiv.org

Reinforcement Learning with Verifiable Rewards (RLVR) reliably improves the reasoning performance of large language models, yet it appears to modify only a small fraction of ...

☆ Save Cite Cited by 1 Related articles All 3 versions

### Self-Play Methods in Reinforcement Learning for Language Models

[Z.Ye](#) - 2025 - knowledge.uchicago.edu

In this thesis we develop a series of practical algorithms for language model to self-train, by actively and strategically creating and controlling learning experiences themselves ...

☆ Save Cite Related articles All 2 versions

## Which category does this paper belong to?



## Tool Learning

## Agentic Learning





# Why Analyze Data? – Paper Management

## IN-THE-FLOW AGENTIC SYSTEM OPTIMIZATION FOR EFFECTIVE PLANNING AND TOOL USE

Zhuofeng Li<sup>1,2</sup>, Haosiang Zhang<sup>1,2</sup>, Seungja Han<sup>1</sup>, Sheng Liu<sup>1</sup>, Jianwen Xie<sup>1</sup>, Yu Zhang<sup>1</sup>, Yijin Choi<sup>1</sup>, James Zou<sup>1,2</sup>, Pan Lu<sup>1</sup>

<sup>1</sup>Stanford University, <sup>2</sup>Texas A&M University, <sup>3</sup>UC San Diego, <sup>4</sup>Lambda

Website: <https://agentflow.stanford.edu>

Code Model Demo Visualize

### ABSTRACT

Outcome-driven reinforcement learning has advanced reasoning in large language models (LLMs), but prevailing tool-augmented approaches train a single, monolithic policy that interleaves thoughts and tool calls under full context; this scales poorly with long horizons and diverse tools and generalizes weakly to new scenarios. Agentic systems offer a promising alternative by decomposing work across specialized modules, yet must remain training-free or rely on offline training decoupled from the live dynamics of multi-turn interaction. We introduce AGENT-FLOW, a trainable, *in-the-flow* agentic framework that coordinates four modules (planner, executor, verifier, generator) through an evolving memory and directly optimizes its planner inside the multi-turn loop. To train on-policy in live environments, we propose *Flow-based Group Refined Policy Optimization* (Flow-GRPO), which tackles long-horizon, sparse-reward credit assignment by converting multi-turn optimization into a sequence of tractable single-turn policy updates. It broadcasts a single, verifiable trajectory-level outcome to every turn to align local planner decisions with global success and stabilizes learning with group-normalized advantages. Across ten benchmarks, AGENTFLOW with a 7B-scale backbone outperforms top-performing baselines with average accuracy gains of 14.9% on search, 14.0% on agentic, 14.5% on mathematical, and 4.1% on scientific tasks, even surpassing larger proprietary models like GPT-4o. Further analyses confirm the benefits of *in-the-flow* optimization, showing improved planning, enhanced tool-calling reliability, and positive scaling with model size and reasoning turns.



## IN-THE-FLOW AGENTIC SYSTEM OPTIMIZATION FOR EFFECTIVE PLANNING AND TOOL USE

Zhuofeng Li<sup>1,2</sup>, Haosiang Zhang<sup>1,2</sup>, Seungja Han<sup>1</sup>, Sheng Liu<sup>1</sup>, Jianwen Xie<sup>1</sup>, Yu Zhang<sup>1</sup>, Yijin Choi<sup>1</sup>, James Zou<sup>1,2</sup>, Pan Lu<sup>1</sup>

<sup>1</sup>Stanford University, <sup>2</sup>Texas A&M University, <sup>3</sup>UC San Diego, <sup>4</sup>Lambda

Website: <https://agentflow.stanford.edu>

Code Model Demo Visualize

### ABSTRACT

Outcome-driven reinforcement learning has advanced reasoning in large language models (LLMs), but prevailing tool-augmented approaches train a single, monolithic policy that interleaves thoughts and tool calls under full context; this scales poorly with long horizons and diverse tools and generalizes weakly to new scenarios. Agentic systems offer a promising alternative by decomposing work across specialized modules, yet must remain training-free or rely on offline training decoupled from the live dynamics of multi-turn interaction. We introduce AGENT-FLOW, a trainable, *in-the-flow* agentic framework that coordinates four modules (planner, executor, verifier, generator) through an evolving memory and directly optimizes its planner inside the multi-turn loop. To train on-policy in live environments, we propose *Flow-based Group Refined Policy Optimization* (Flow-GRPO), which tackles long-horizon, sparse-reward credit assignment by converting multi-turn optimization into a sequence of tractable single-turn policy updates. It broadcasts a single, verifiable trajectory-level outcome to every turn to align local planner decisions with global success and stabilizes learning with group-normalized advantages. Across ten benchmarks, AGENTFLOW with a 7B-scale backbone outperforms top-performing baselines with average accuracy gains of 14.9% on search, 14.0% on agentic, 14.5% on mathematical, and 4.1% on scientific tasks, even surpassing larger proprietary models like GPT-4o. Further analyses confirm the benefits of *in-the-flow* optimization, showing improved planning, enhanced tool-calling reliability, and positive scaling with model size and reasoning turns.

## In-the-flow agentic system optimization for effective planning and tool use

Search within citing articles

### Latent collaboration in multi-agent systems

J. Zou, X. Sun, R. Liu, G. Li, C. Chen, P. Lu, H. Zhou, ... arXiv preprint arXiv: ... 2025 ...   
Multi-agent systems (MAS) extend large language models (LLMs) from independent single-model reasoning to collaborative system-level intelligence. While existing LLM agents ...   
Q. Shen, 10 Cite, Cited by 4, Related articles, All 2 versions, 10

### Adaptation of agentic AI

E. Jans, L. Liu, Z. Shi, Z. Wang, L. He, Y. Li, ... arXiv preprint arXiv: ... 2025 ...   
Cutting-edge agentic AI systems are built on foundation models that can be adapted to plan, reason, and interact with external tools to perform increasingly complex and specialized ...   
Q. Shen, 10 Cite, Cited by 2, Related articles, All 2 versions, 10

### DeepAgent: A General Reasoning Agent with Scalable Toolsets

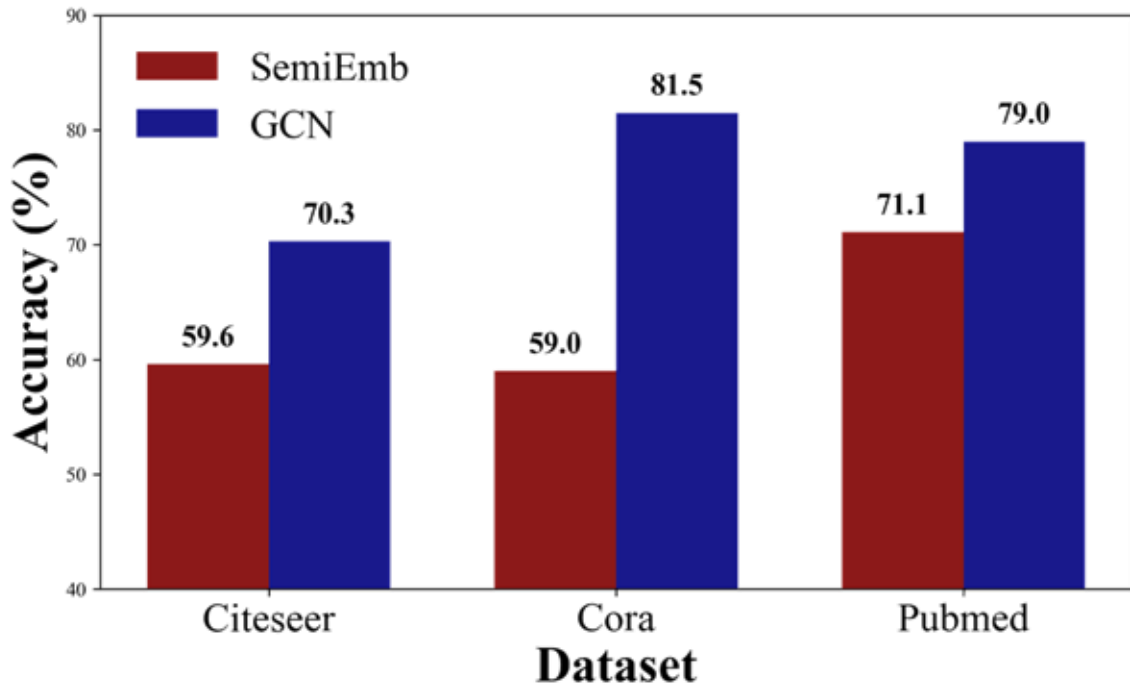
S. Li, V. Kulkarni, J. Li, S. D. J. Li, Y. Wang, ... arXiv preprint arXiv: ... 2025 ...   
Large reasoning models have demonstrated strong problem-solving abilities, yet real-world tasks often require external tools and long-horizon interactions. Existing agent frameworks ...   
Q. Shen, 10 Cite, Cited by 2, Related articles, All 2 versions, 10

### The Path Not Taken: RLVR Provably Learns Off the Principals

S. Zhou, Z. Zhang, H. Zhang, D. Fu, Z. Li, J. Zhou, ... arXiv preprint arXiv: ... 2025 ...   
Reinforcement Learning with Verifiable Rewards (RLVR) reliably improves the reasoning performance of large language models, yet it appears to modify only a small fraction of ...   
Q. Shen, 10 Cite, Cited by 1, Related articles, All 2 versions, 10

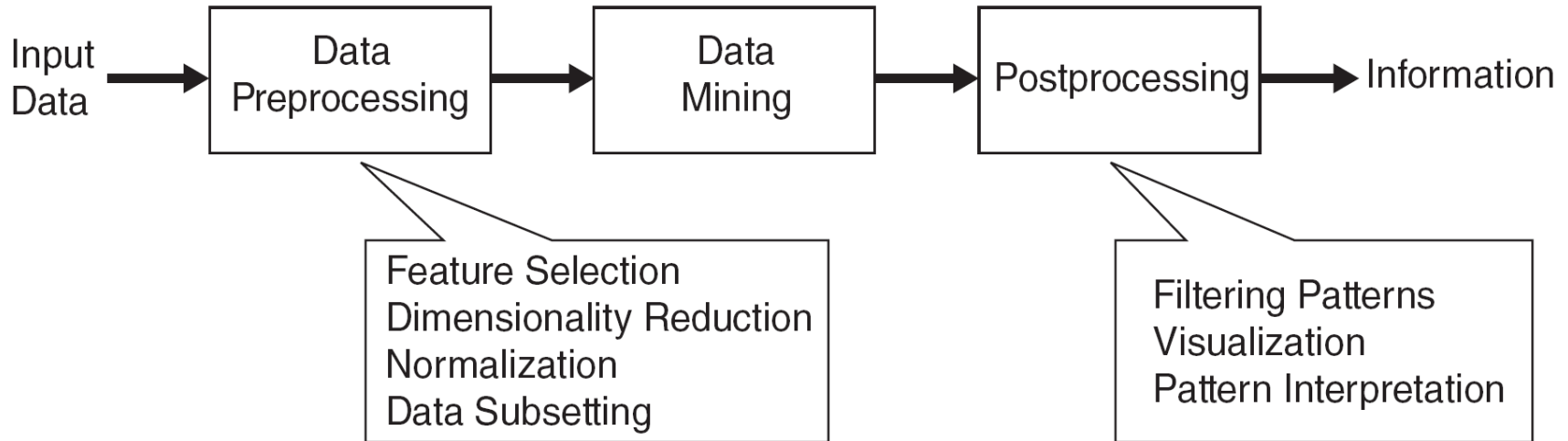
### Self-Play Methods in Reinforcement Learning for Language Models

Z. Xu, 2025 ...   
In this thesis, we develop a series of practical algorithms for language model to self-train, by actively and strategically creating and controlling learning experiences themselves ...   
Q. Shen, 10 Cite, Related articles, All 2 versions, 10





# What is Data Mining?

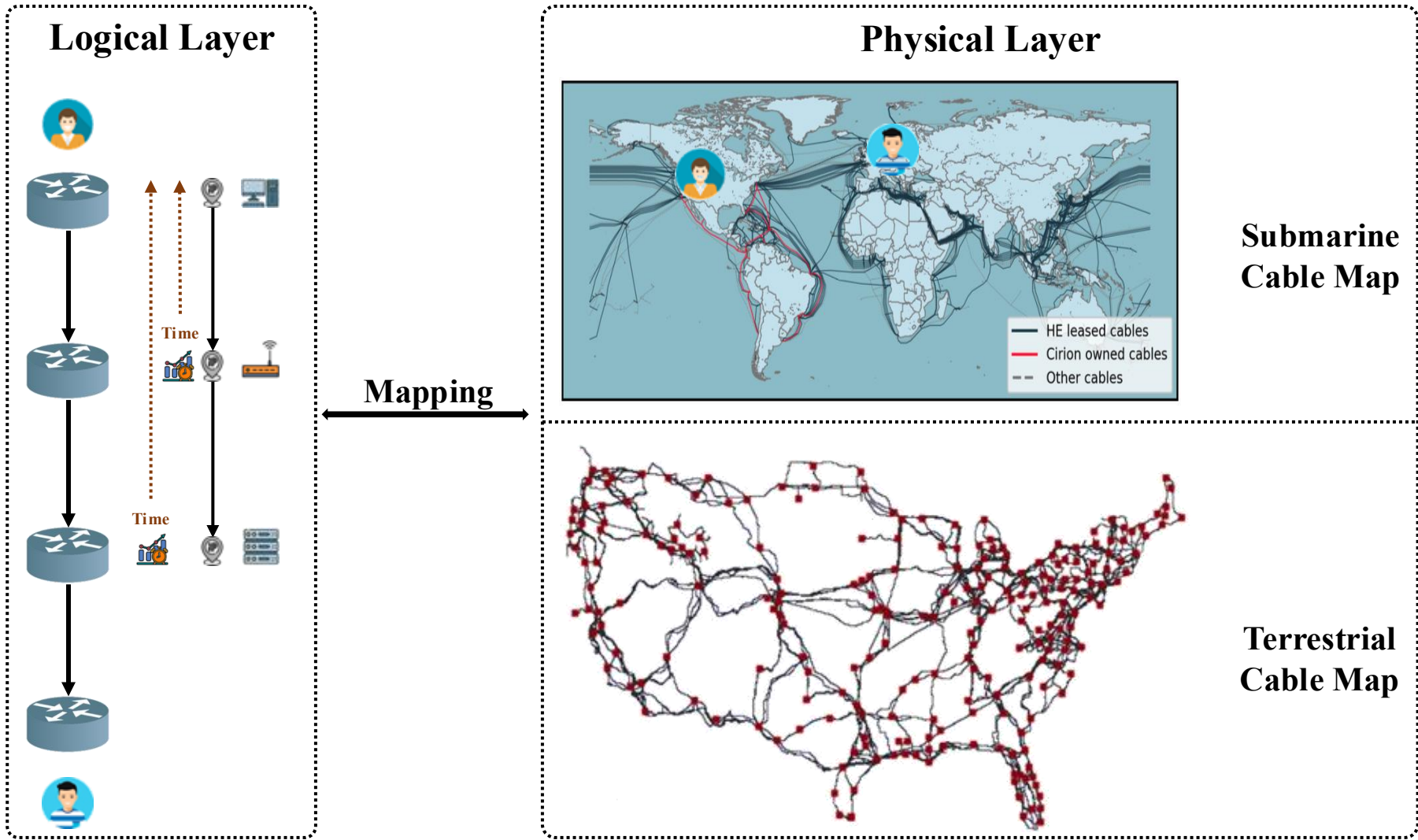


## Many Definitions

Non-trivial extraction of implicit, previously unknown and potentially useful information from data

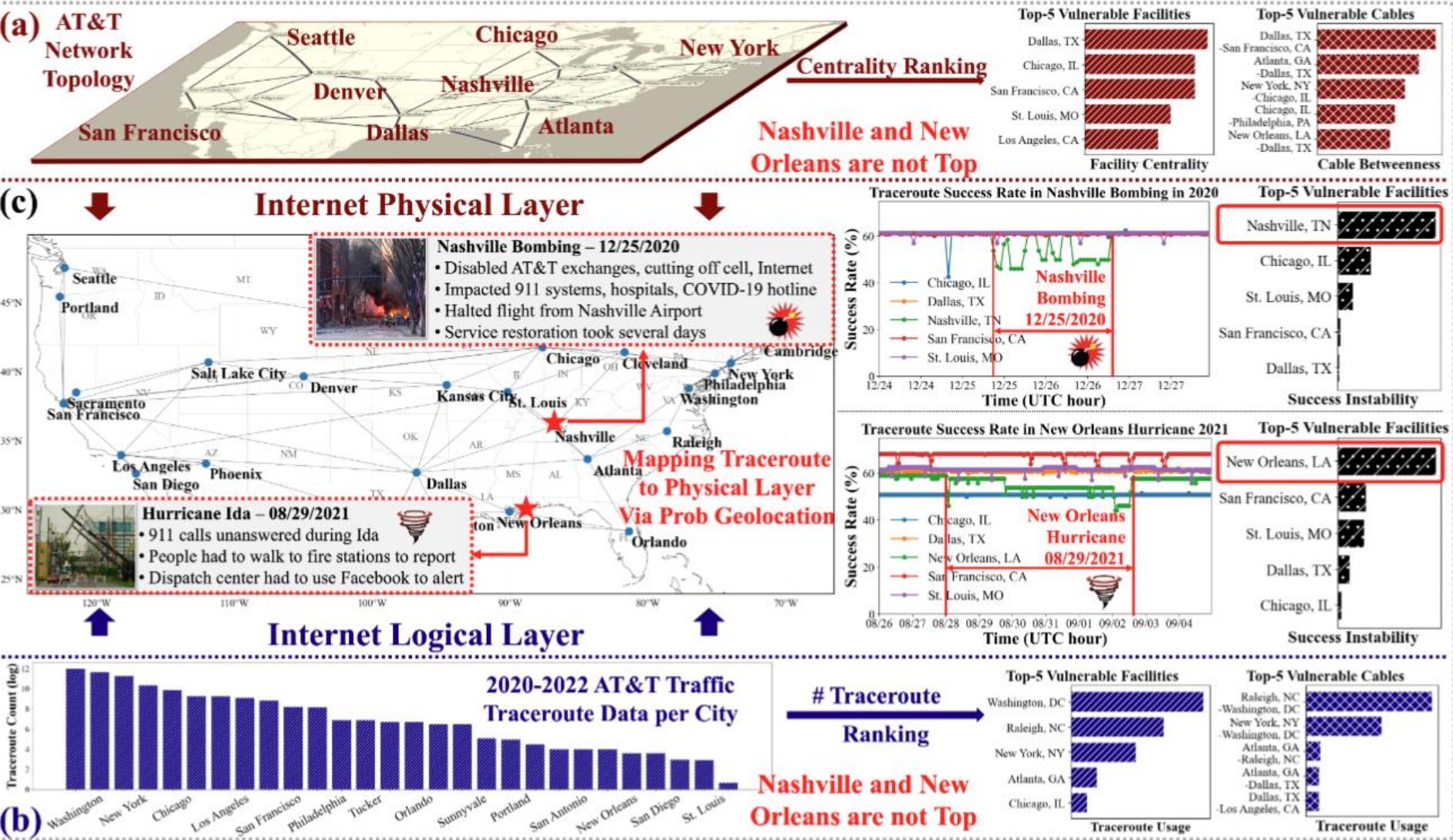
Exploration & analysis, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns

# Why Data Mining? – Networking Infra Risk



Which physical cable path does this logic signal traverse?

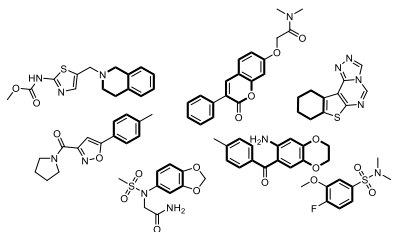
# Why Data Mining? – Networking Infra Risk





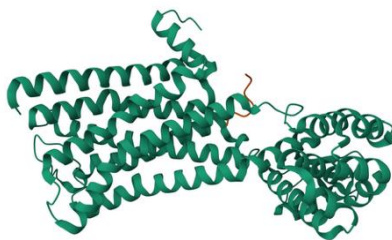
# Why Data Mining? – Drug Design

## Chemical Libraries

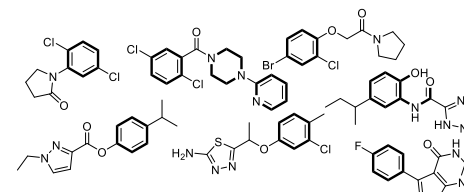


Number of Molecules: 103-106

## Protein Target



## Virtual Libraries

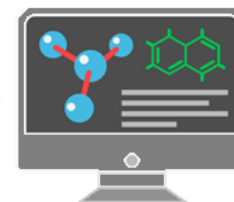


e.g.,  $10^9$  Virtual Molecules on the REAL database in Enamine Ltd.



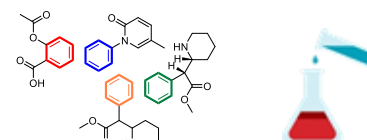
High Throughput  
Screening (HTS)

Training



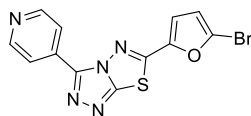
Deep Learning Models

## Predicted Actives



Number of Molecules: 500-1000

Evaluating

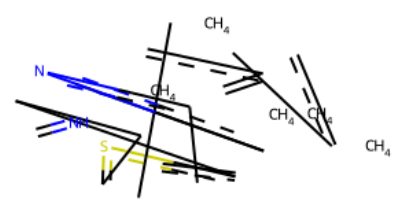
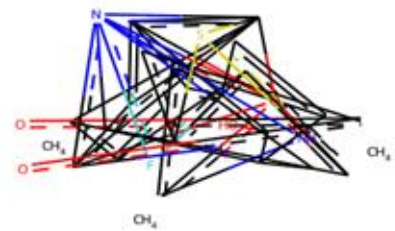
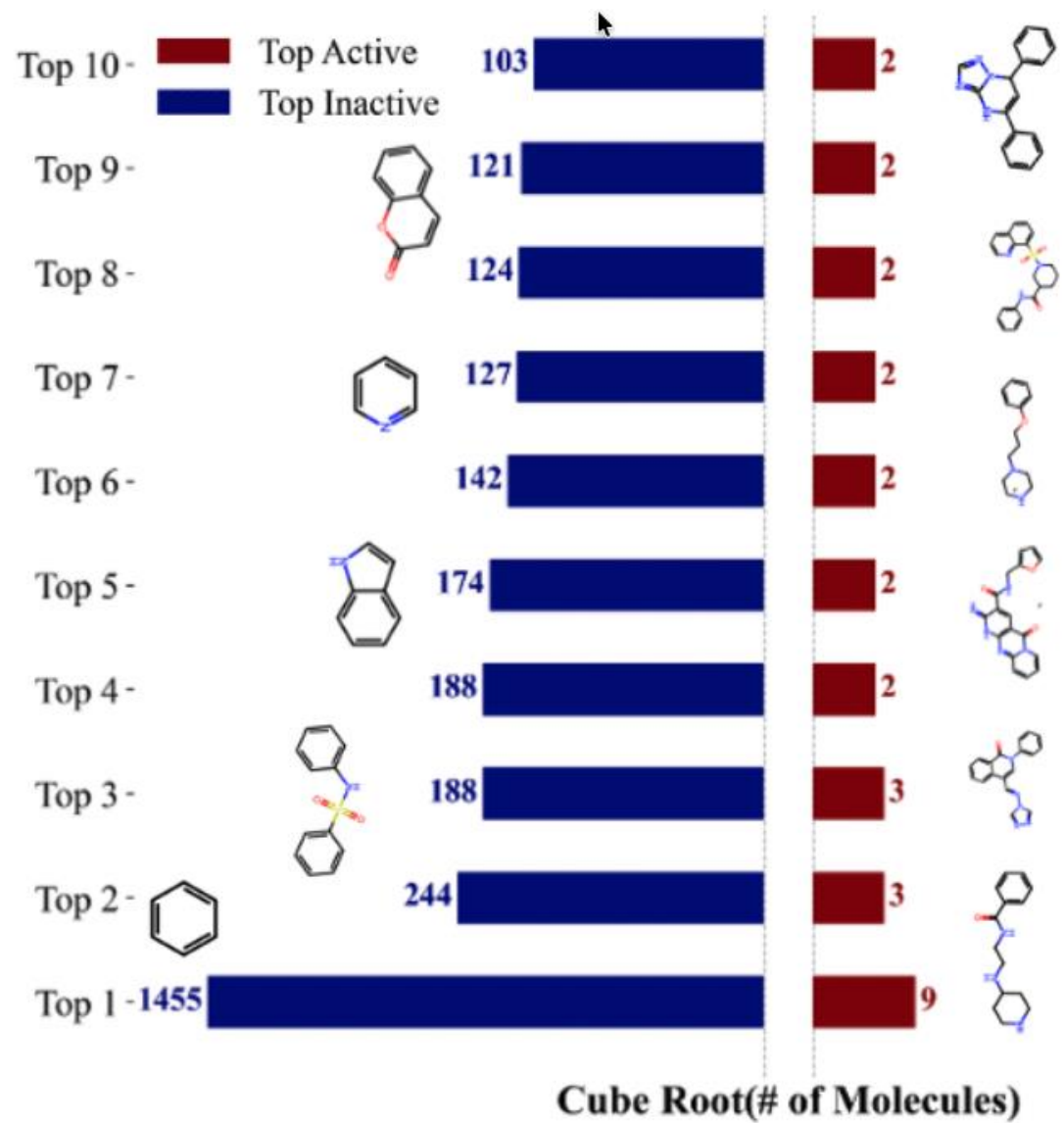


Hit Rate: 0.05%-0.5%





# Why Data Mining? – Drug Design





# Why Data Mining? – Commercial Perspective

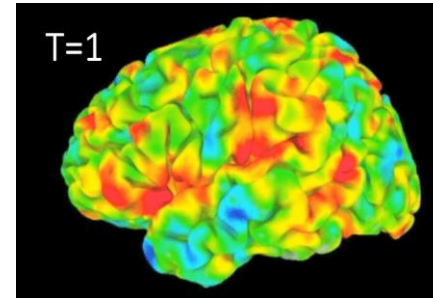
- Lots of data is being collected and warehoused
  - Web data **1,000 terabytes,**  
**1,000,000,000,000,000= bytes**
    - Google has Peta Bytes of web data
    - Facebook has billions of active users
  - purchases at department/  
grocery stores, e-commerce
    - Amazon handles millions of visits/day
  - Bank/Credit Card transactions
- Computers have become cheaper and more powerful
- Competitive Pressure is Strong
  - Provide better, customized services for an edge (e.g. in Customer Relationship Management)





# Why Data Mining? – Scientific Perspective

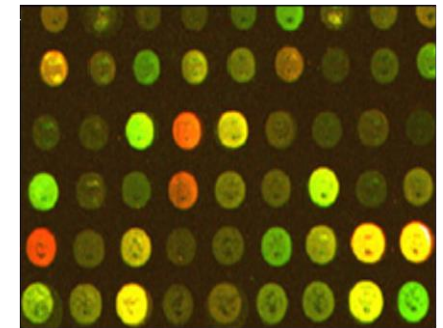
- Data collected and stored at enormous speeds
  - Remote sensors on a satellite
    - NASA EOSDIS archives over petabytes of earth science data / year
  - Telescopes scanning the skies
    - Sky survey data
  - High-throughput biological data
  - Scientific simulations
    - terabytes of data generated in a few hours
- Data mining helps scientists
  - in automated analysis of massive datasets
  - In hypothesis formation



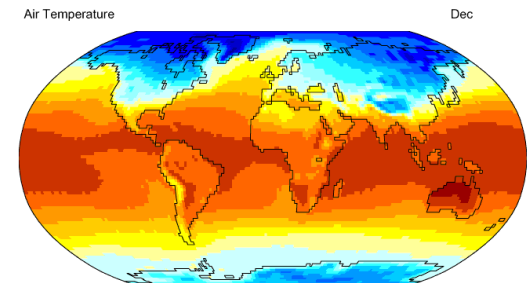
fMRI Data from Brain



Sky Survey Data



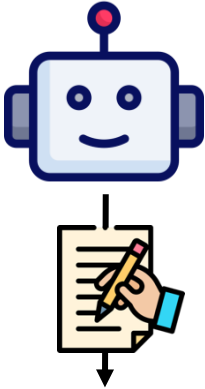
Gene Expression Data



Surface Temperature of Earth

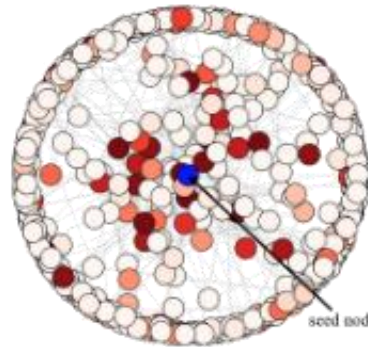


# Why Data Mining? – Social Good

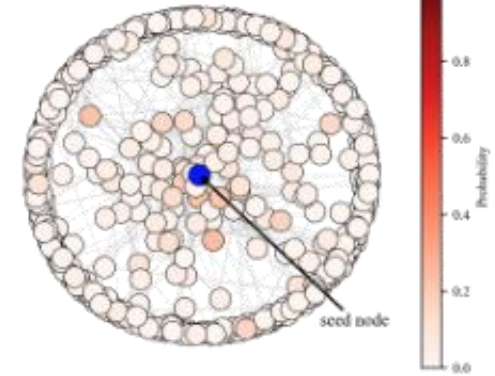


Text: "Breaking: NASA confirms first-ever human colony on Mars will begin next year — tickets for civilians already being sold out in minutes!"

Ours  
Influence Spread=2768.06

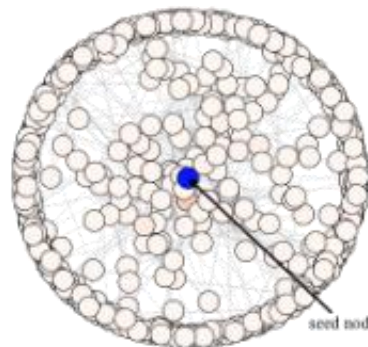


IC Model  
Influence Spread=534.50

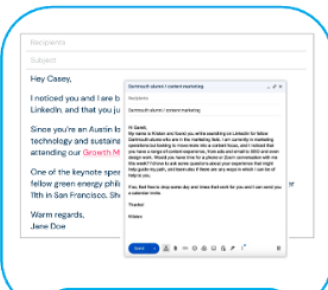
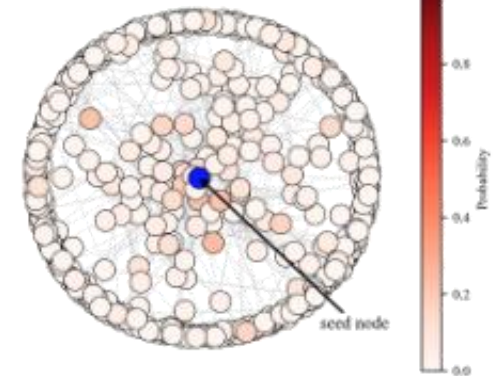


Text: "Today I bought a new pencil."

Ours  
Influence Spread=20.45



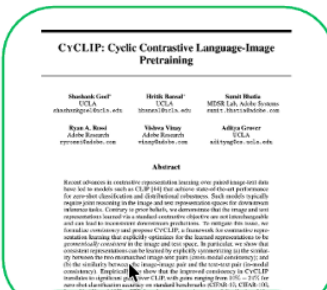
IC Model  
Influence Spread=534.50



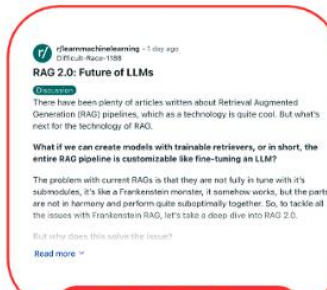
Email Generation



Review Generation



Abstract Generation



Topic Writing



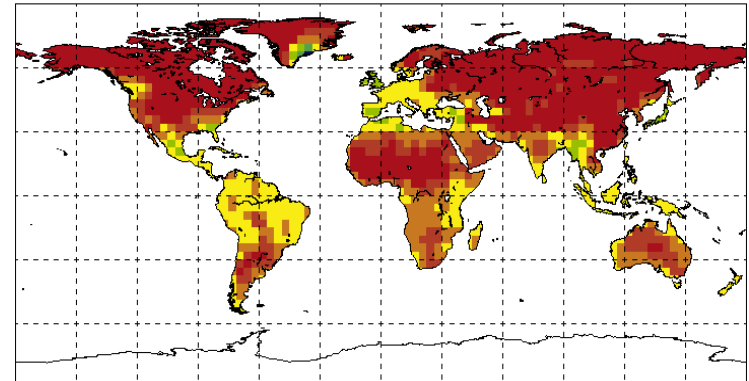
# However, we have challenges – Question

## What kind of data mining question you want to answer?



Improving health care and reducing costs

CCCma/A2a January to January Mean Temperature (degrees C) 2080s relative to 1961-90



Predicting the impact of climate change



Finding alternative/ green energy sources



Reducing hunger and poverty by increasing agriculture production



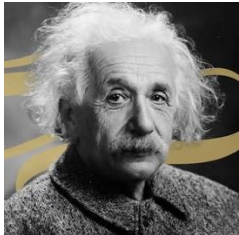
# However, we have challenges – Question

**What kind of data mining question you want to answer?**



**Judge a man by his questions rather than his answers.**

----- Voltaire



**The important thing is not to stop questioning.**

----- Albert Einstein



**He who asks a question is a fool for five minutes; he who does not ask a question remains a fool forever.**

----- Confucius

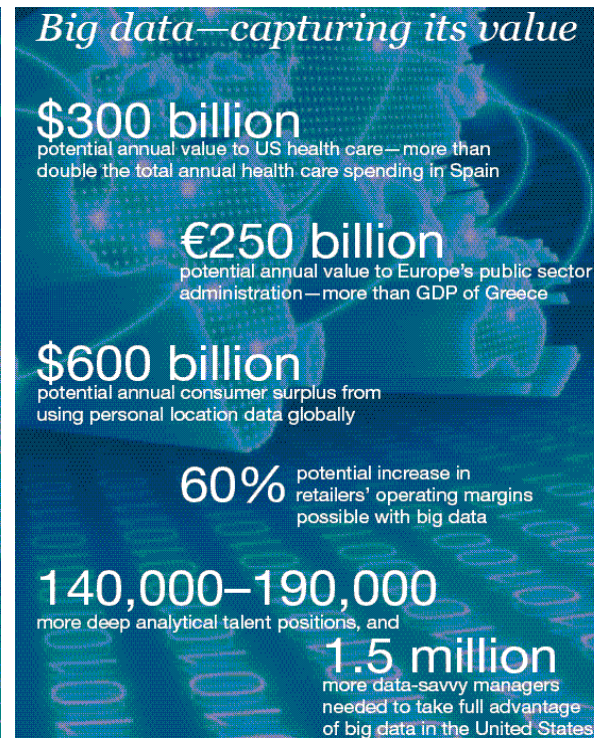
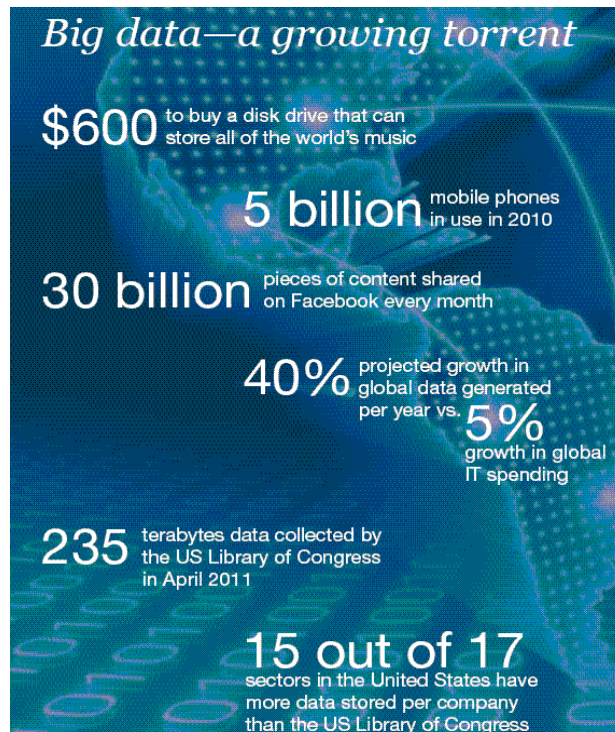


# However, we have challenges – Data

## Data is usually in a very large scale!

McKinsey Global Institute

Big data: The next frontier  
for innovation, competition,  
and productivity

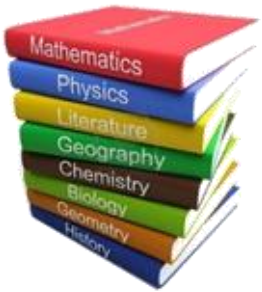




# However, we have challenges – Data

**Data is usually in a very large scale!**

**Textbook  
Knowledge Base**



**158 million books**

[ISBN DB 2023](#)

**Internet  
Knowledge Base**



**1.1 billion websites**

[Musemind 2024](#)

**Neural  
Knowledge Base**



**405 billion parameters**

[Hugging Face 2024](#)



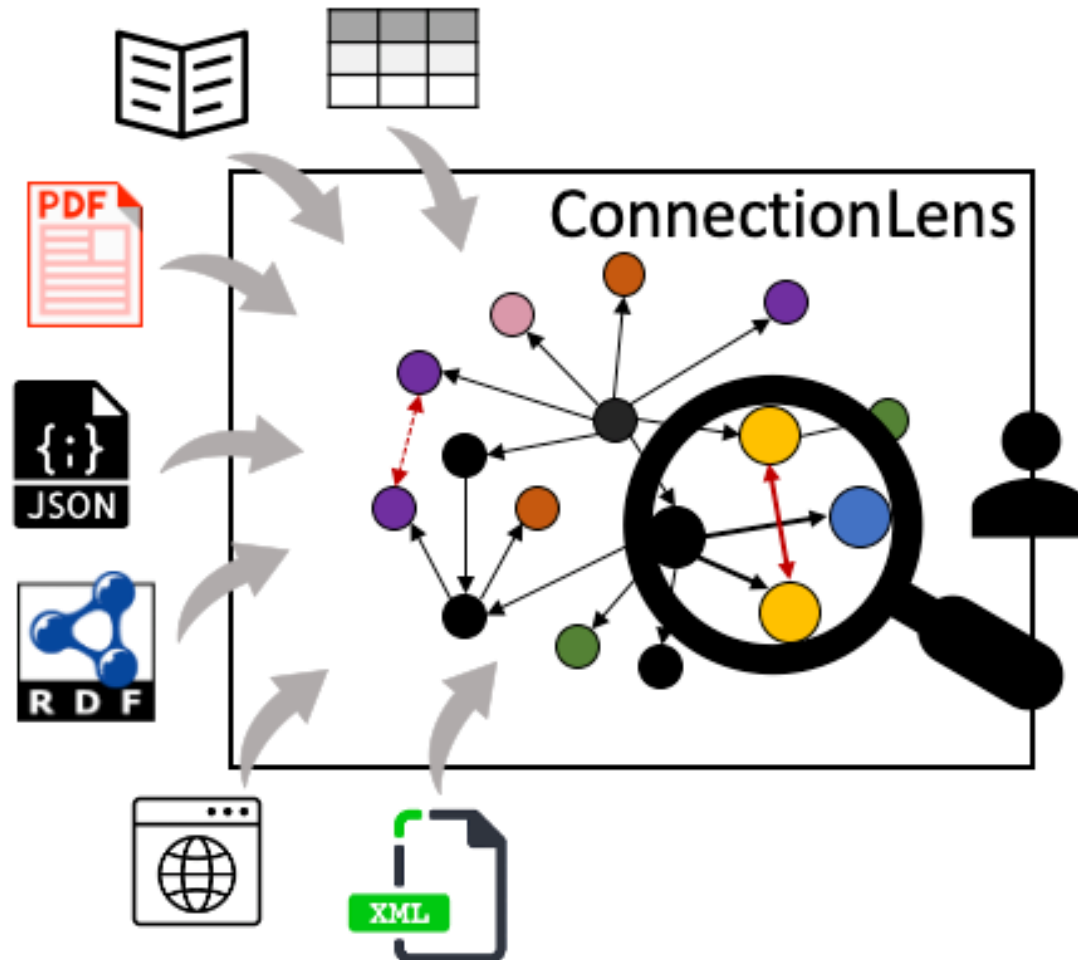
 **2.5 petabytes, 1 billion books**

- We remember meanings, not details.
- We forget on purpose.
- Tiny active memory, Larger long-term memory.



# However, we have challenges – Data

**Data is diverse and heterogeneous**





- **Data is everywhere**
- **Data Mining brings scientific advancement and social wellness**
- **However, there are challenges**
  - (1) What are good questions to ask?
  - (2) Data is scattered around the world, how to find them?
  - (3) Data is very large-scale, how to analyze them efficiently, space/time?
  - (4) Data is very heterogeneous and specialized

**This is the reason for taking data mining!**

# Question Time!





<https://ml-graph.github.io/winter-2025/>

**All information will be  
available on the website!**

# Goals

- Broad overview of Data Mining
- Data Mining Skills – Knowledge and Code
- Machine Learning Skills – Knowledge and Code
- Real-world GML/DM applications

## Prerequisite

- Linear Algebra, Probability /Statistics, Calculus
- Programming – Python, PyTorch
- Curiosity – Critical Thinking
- Diligence – Hard Working



# Course Logistics - Time

## Times:

- **Classes:** Monday/Wednesday 12:00-1:20 pm PST, Gerlinger 302
- **Office hours:** Wednesday 1:20-2:00 pm PST, other time by appointment
- **Zoom:** <https://uoregon.zoom.us/j/4052006678>



# Course Logistics – Quizz

## Times:

- **Classes:** Monday/Wednesday 12:00-1:20 pm PST, Gerlinger 302
- **Office hours:** Wednesday 1:20-2:00 pm PST, other time by appointment
- **Zoom:** <https://uoregon.zoom.us/j/4052006678>

## Components:

### Course Assessment and Grading Scale

Category	CS-453 (%)	CS-553 (%)
Quizz 1	20%	15%
Quizz 2	20%	15%
Project	40%	45%
Participation	5%	5%
Paper Presentation	15%	20%
Overleaf Bonus	5%	5%

- As long as you are **active thinking** and **understand the content**, you will be good

