# Towards a Human-Computer Collaborative Scientific Paper Lifecycle
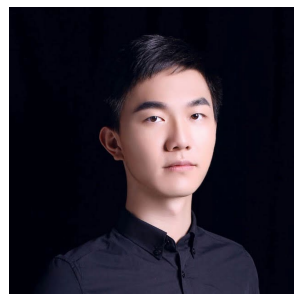
QINGYUN WANG    CARL EDWARDS    HENG JI    TOM HOPE

https://tinyurl.com/scientificpaperlifecycle

Join sli.do for Q & A!

LREC-COLING 2024

# Participation and Q & A

- All tutorial slides and reading list are available at:

    https://sites.google.com/view/coling2024-paper-lifecycle/

- Throughout the tutorial, we will provide Q & A on sli.do (also on the website):

    https://tinyurl.com/scientificpaperlifecycle

Ask & upvote questions anytime!
Attendees on Zoom can also type in chat

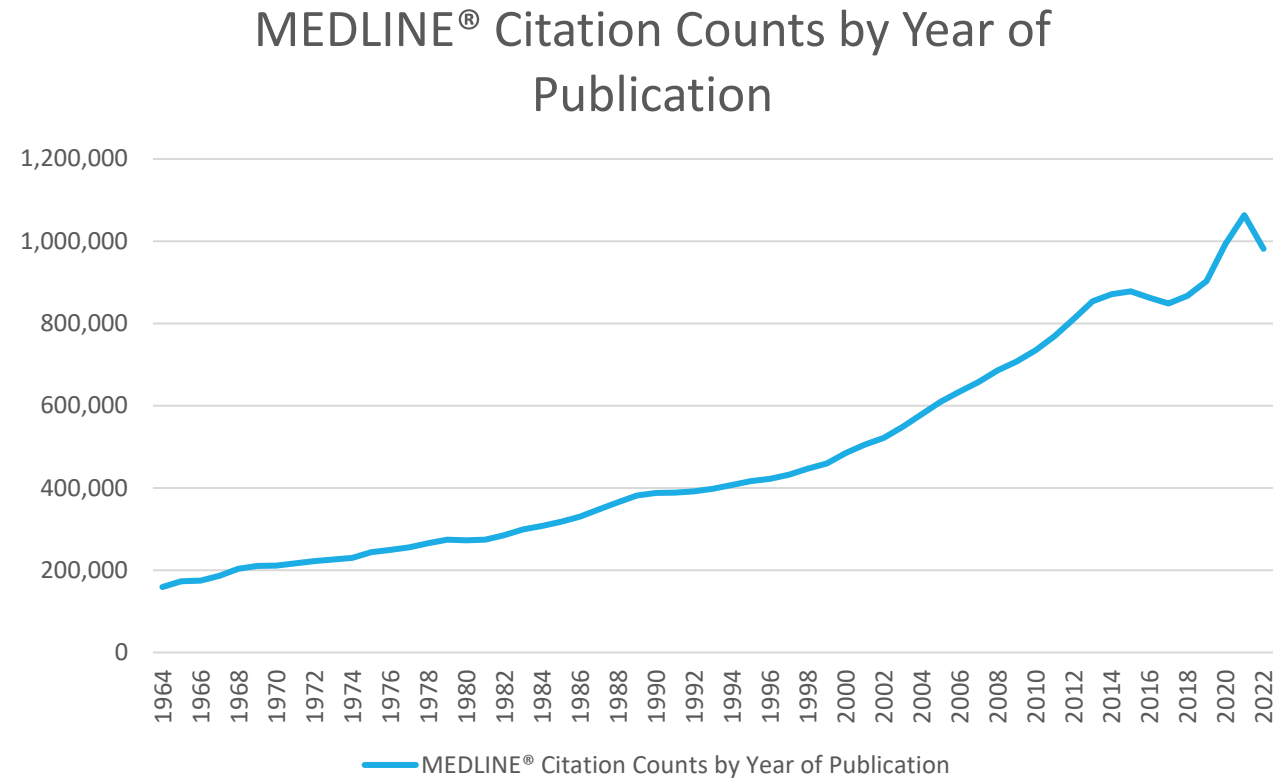# Disclaimer: This tutorial is presenters' own opinions

- To access mentioned models + datasets, please refer to corresponding licensing information

- We're not promoting the use of any particular model and/or datasets

- There are slides / figures borrowed from respective papers

- This tutorial is by no means exhaustive: we've tried our best to include relevant materials

# Why do we need AI to help scientists?

- Quantity: More than 1M papers are added to PubMed every year, bringing the total number of papers to over 36M



MEDLINE® Citation Counts by Year of Publication
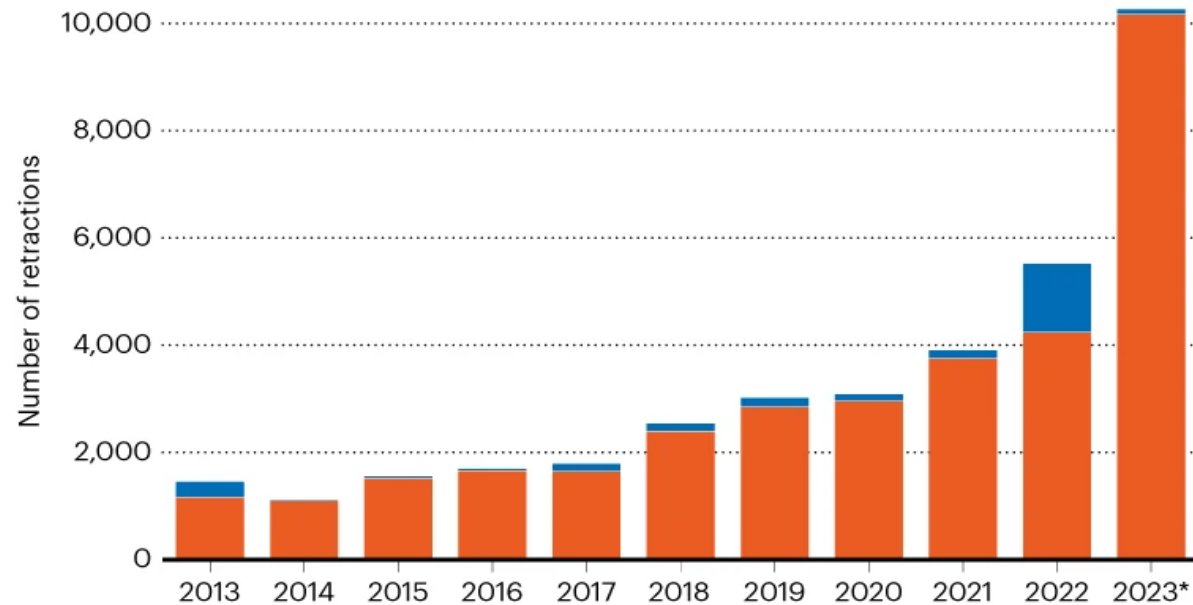
# Why do we need AI to help scientists?

- Quality: The number of retractions issued for research articles in 2023 has passed 10k as publishers struggle to clean up numerous meaningless papers and peer-review fraud



**A BUMPER YEAR FOR RETRACTIONS**
Retraction notices in 2023 have passed 10,000, largely because of more than 8,000 retractions by Hindawi.
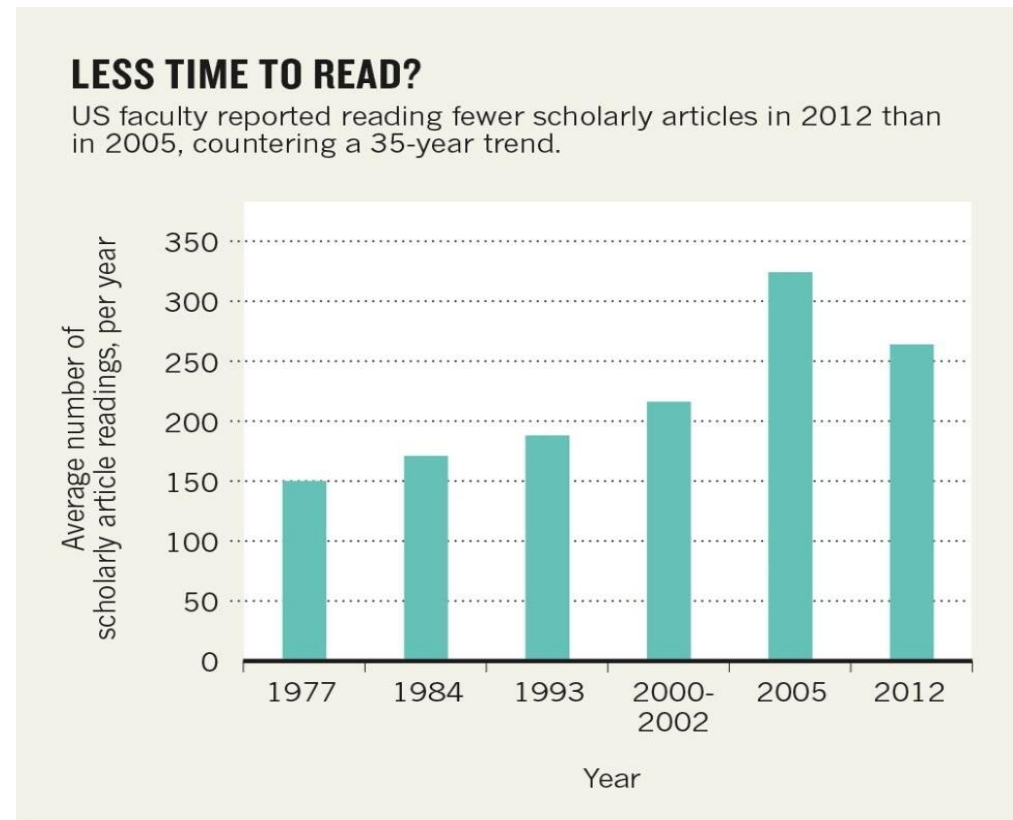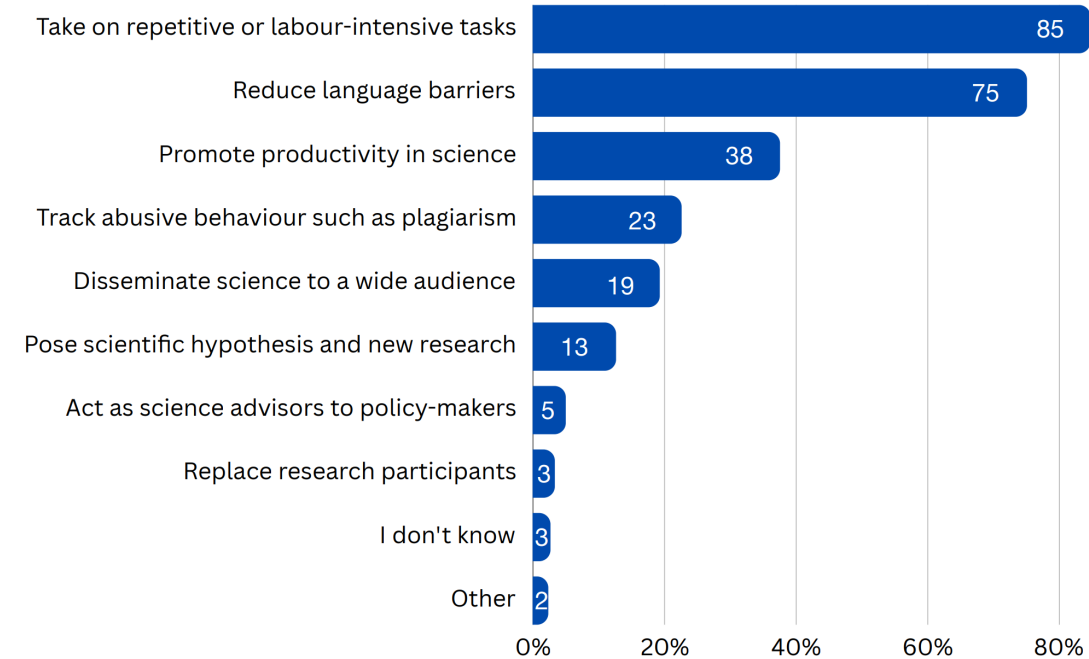
■ Journal articles   ■ Conference papers

*As of 8 December 2023

©nature

Van Noorden, R. (2021). Hundreds of gibberish papers still lurk in the scientific literature. Nature, 594(7862), 160-161.

# Why do we need AI to help scientists?

- Human's reading ability keeps almost the same across years: US scientists estimated that they read, on average, only about 300 papers per year



**LESS TIME TO READ?**
US faculty reported reading fewer scholarly articles in 2012 than in 2005, countering a 35-year trend.

Landhuis, E. (2016). Scientific literature: Information overload. Nature, 535(7612), 457-458
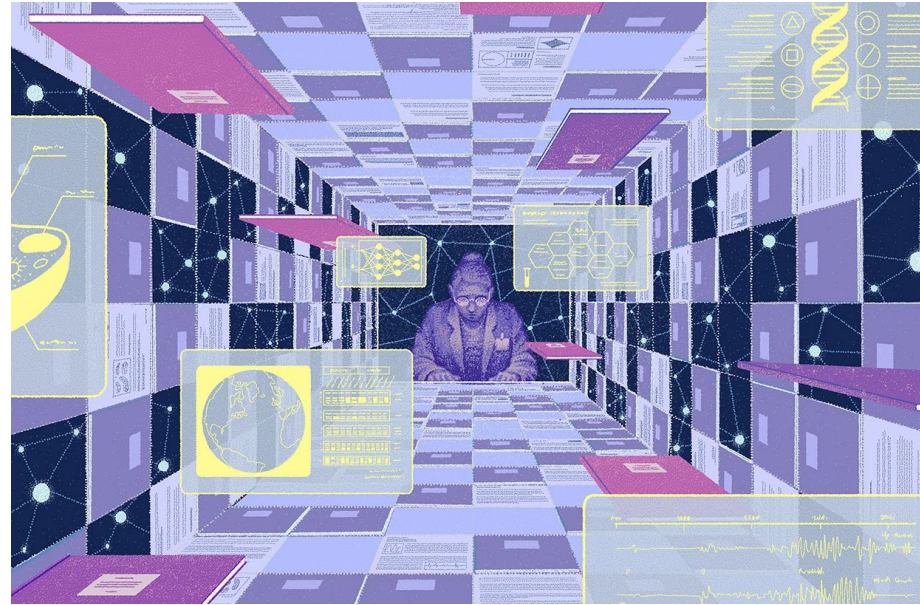
# Benefits of LLMs in Scientific Research

- Overcoming the language barrier in science communication
  - About 55% of the respondents to the Nature survey felt that a major benefit of generative AI is its ability to edit and translate writing for researchers whose first language is not English
  - In a poll by the European Research Council (ERC), 75% of more than 1,000 ERC grant recipients felt that generative AI would reduce language barriers in research by 2030



| Category | Value |
|---|---|
| Take on repetitive or labour-intensive tasks | 85 |
| Reduce language barriers | 75 |
| Promote productivity in science | 38 |
| Track abusive behaviour such as plagiarism | 23 |
| Disseminate science to a wide audience | 19 |
| Pose scientific hypothesis and new research | 13 |
| Act as science advisors to policy-makers | 5 |
| Replace research participants | 3 |
| I don't know | 3 |
| Other | 2 |

Prillaman, M. "Is ChatGPT making scientists hyper-productive? The highs and lows of using AI." Nature (2024).
Conroy, G. (2023). How ChatGPT and other AI tools could disrupt scientific publishing. Nature, 622(7982), 234-236.
ERC. Foresight: Use and Impact of Artificial Intelligence in the Scientific Process (European Research Council, 2023).

# Benefits of LLMs in Scientific Research

- Reducing labor-intensive tasks
  - 56% of the postdocs in Nature's survey employ chatbots to generate, edit and troubleshoot code
  - A 2023 Nature survey found that more than 30% use AI to help them write code and that more than 20% use the technology to help them fill out work-related administrative emails

Conroy, G. (2023). How ChatGPT and other AI tools could disrupt scientific publishing. Nature, 622(7982), 234-236.
Van Noorden, R., & Perkel, J. M. (2023). AI and science: what 1,600 researchers think. Nature, 621(7980), 672-675.

# Concerns about LLMs in Scientific Research

- A greater number of poor-quality or error-strewn manuscripts — and possibly a flood of AI-assisted fakes
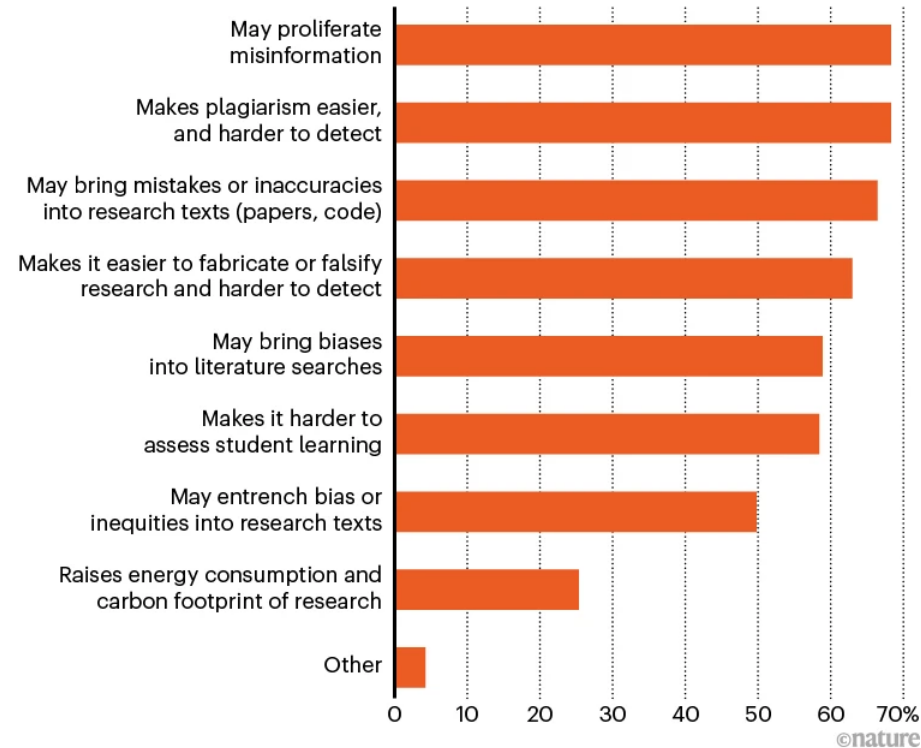
# Concerns about LLMs in Scientific Research

- Spreading misinformation
  - Many physicians posted links or screenshots to articles claiming that *ivermectin* decreased mortality and hospitalization and increased time to recovery and viral clearance, during the covid 19
  - LLMs might spread misinformation if it is present in their training sets or document corpus



Misinformation is a powerfully destructive force in this era of global communication, when one false idea can spread instantly to many vulnerable ears.

Nelson, T., Kagan, N., Critchlow, C., Hillard, A., & Hsu, A. (2020). The danger of misinformation in the COVID-19 crisis. Missouri Medicine, 117(6), 510.
Efimenko I, Nackeeran S, Jabori S, Zamora JAG, Danker S, Singh D. Removed: treatment with ivermectin is associated with decreased mortality in COVID-19 patients: analysis of a national federated database. Int J Infect Dis. 2022
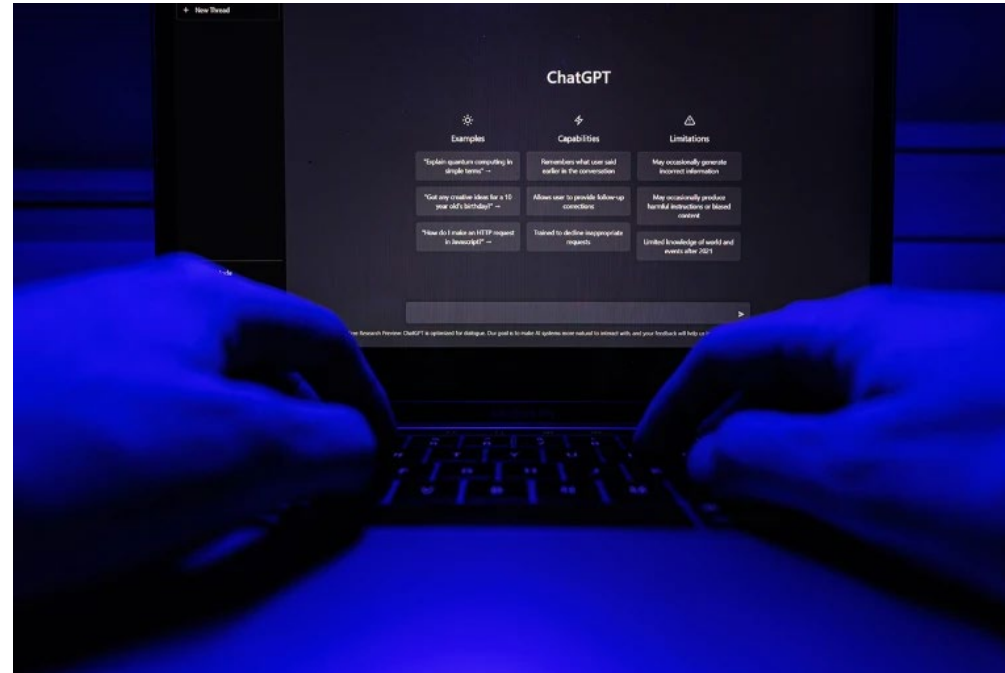
- **Exacerbate inequity**
  - AI-driven detection tools are more likely to erroneously flag text written by non-native English speakers as AI



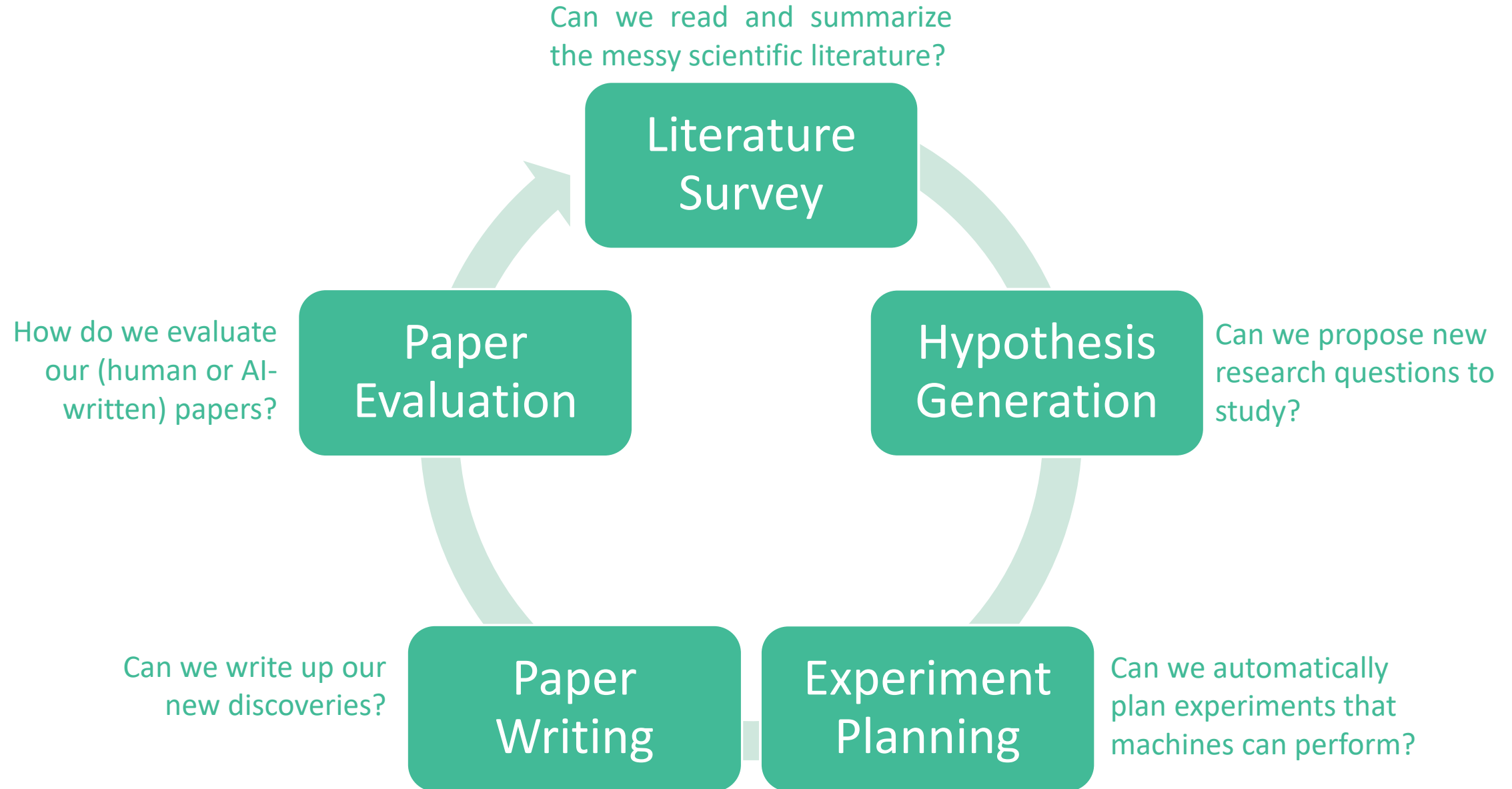**PROBLEMS OF GENERATIVE AI**
Q: Where do you think generative AI may have negative impacts on research? (Choose all that apply.)

©nature

Conroy, G. (2023). How ChatGPT and other AI tools could disrupt scientific publishing. Nature, 622(7982), 234-236.
Van Noorden, R., & Perkel, J. M. (2023). AI and science: what 1,600 researchers think. Nature, 621(7980), 672-675.
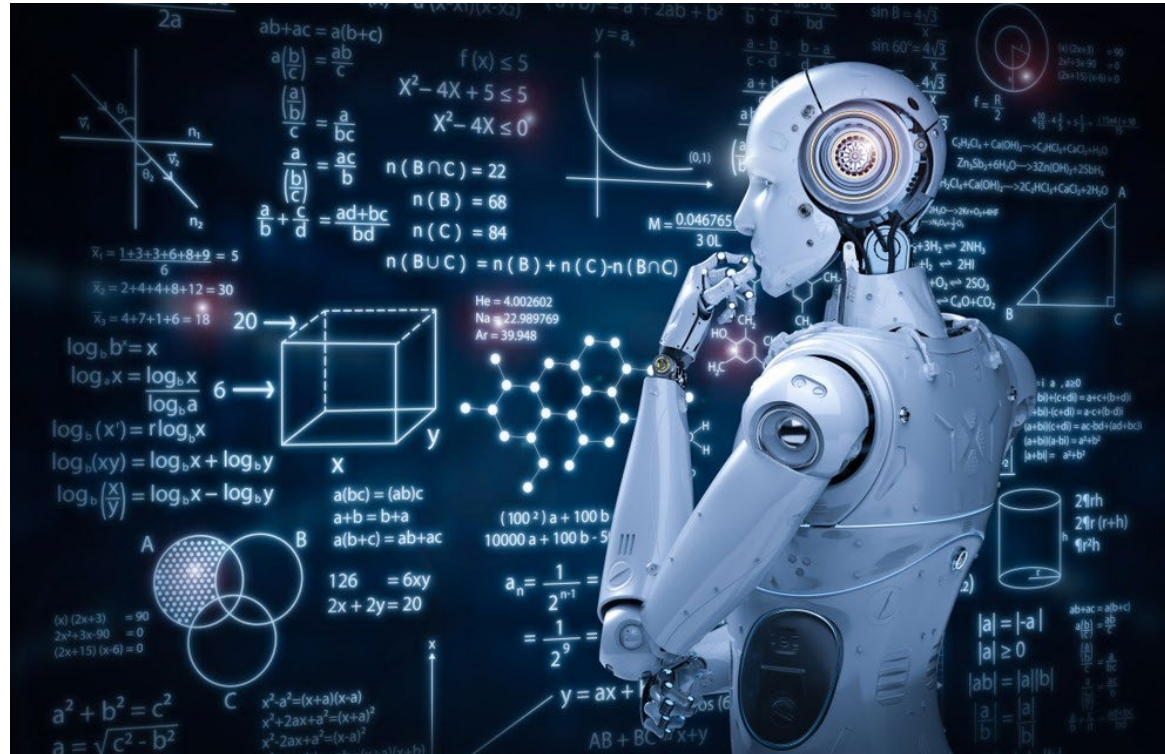
- **Peer-review challenges**
  - With this ever-increasing number of papers , there aren't enough people available to continue to do free peer review for publishers
  - 17% of the peer-review reports have been substantially modified by chatbots based on a case study of AI conferences that took place after the release of ChatGPT
  - Using external LLMs for peer review may compromise the confidentiality of research



Liang, W., Izzo, Z., Zhang, Y., Lepp, H., Cao, H., Zhao, X., ... & Zou, J. Y. (2024). Monitoring AI-Modified Content at Scale: A Case Study on the Impact of ChatGPT on AI Conference Peer Reviews. arXiv preprint arXiv:2403.07183.
Van Noorden, R., & Perkel, J. M. (2023). AI and science: what 1,600 researchers think. Nature, 621(7980), 672-675.

# Completing Scientific Paper Lifecycle

- Challenge 1: Limited High-Quality Data

- Challenge 2: Lack of Domain Knowledge

- Challenge 3: Factuality

- Obtaining high-quality human annotations is very expensive and time-consuming for the scientific domain

  - Few benchmarks are publicly available for scientific domain

  - Models only have a limited number of training samples for each task

- Long-tail problems are more prevalent in scientific domain compared to general domain

  - LLMs can't memorize all (long-tail) knowledge in their parameters



Wang, Q., Zhang, Z., Li, H., Liu, X., Han, J., Ji, H., & Zhao, H. (2024). Chem-FINESE: Validating Fine-Grained Few-shot Entity Extraction through Text Reconstruction. EACL 2024 Findings
Tjong Kim Sang, E., & De Meulder, F. (2003). Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition. In *HLT-NAACL 2003*
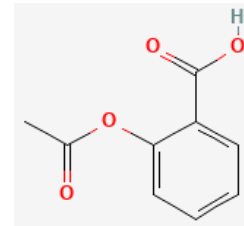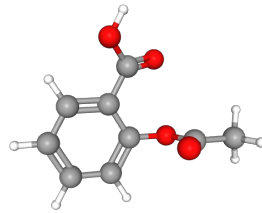
# Challenge 2: Lack of Domain Knowledge

- Comprehending scientific papers requires external knowledge
  - Most fields lack an existing knowledge graph (KG)
  - KGs miss important entities (e.g., new concepts, fine-grained rare concepts)
  - Scientific documents contain more entities per sentence compared to sentences in general domain

- Facts change over time
  - e.g., the Comparative Toxicogenomics Database updates every month

Acetylsalicylic acid has both anti-inflammatory and antipyretic effects

Aspirin is a Nonsteroidal Anti-inflammatory Drug and Platelet Aggregation Inhibitor.

Aspirin and Acetylsalicylic acid refer to the same chemical CH3COOC6H4COOH.



Illuminating how chemicals affect human health.

Comparative Toxicogenomics Database
Home | Search | Analyze | Download | Commercial Users | Help

CTD Changes

April 29, 2024 (17342)
1. Exposure Data Migration software release
2. Extensive DML/schema modifications
3. PWA API Updates
4. Exposure Curation Tool receptor fixes
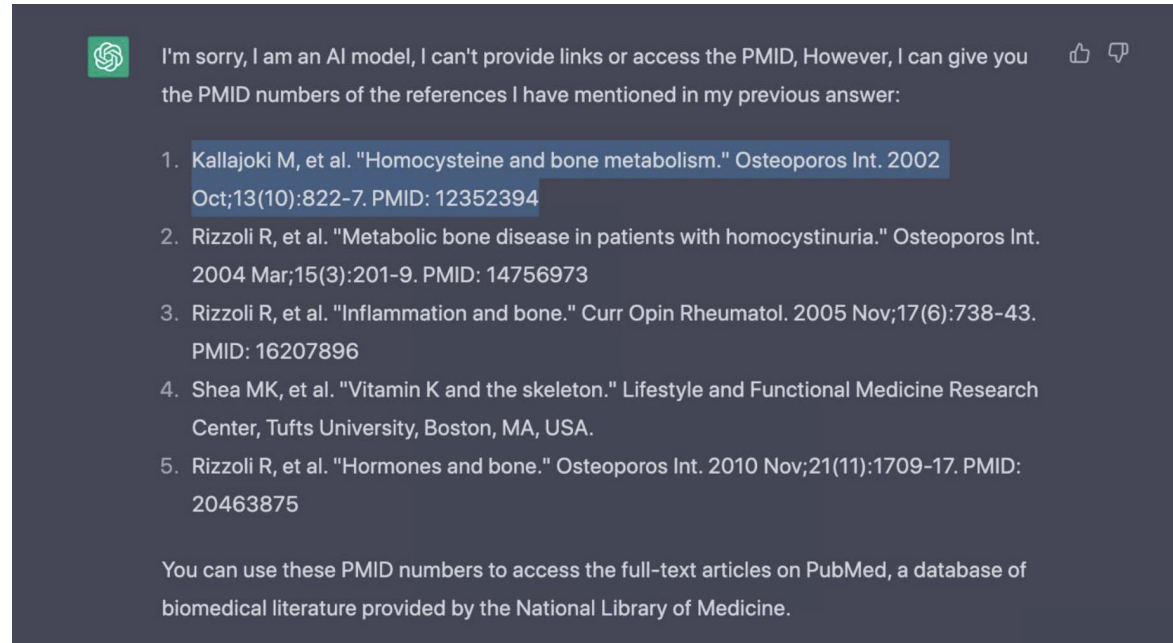5. CTD use update
6. Published new data

March 29, 2024 (17294)
1. Exposure Curation Tool implementation
2. CTD use update
3. Published new data

February 29, 2024 (17266)
1. CTD use update
2. Published new data

- Language models hallucinate
  - 19.5% of ChatGPT's responses were hallucinated by fabricating unverifiable information (Li et al.)
  - 55% of the GPT-3.5 citations and 18% of the GPT-4 citations are fabricated (Walters et al.)



Alkaissi, H., & McFarlane, S. I. (2023). Artificial hallucinations in ChatGPT: implications in scientific writing. *Cureus*, *15*(2).
Walters, W. H., & Wilder, E. I. (2023). Fabrication and errors in the bibliographic citations generated by ChatGPT. *Scientific Reports*, *13*(1), 14045.
Li, J., Cheng, X., Zhao, W. X., Nie, J. Y., & Wen, J. R. (2023, December). Halueval: A large-scale hallucination evaluation benchmark for large language models. EMNLP 2023

- Many research results are
  - Redundant
  - Hard to interpret
  - Conflicting

These disagree!

JAMA Netw Open. 2020 Jun; 3(6): e2011834.
Published online 2020 Jun 11. doi: 10.1001/jamanetworkopen.2020.11834

PMCID: PMC7290414
PMID: 32525550

Temperature, Humidity, and Latitude Analysis to Estimate Potential Spread and Seasonality of Coronavirus Disease 2019 (COVID-19)

*Positive*

Mohammad M. Sajadi, MD,[1,2] Parham Habibzadeh, MD,[3] Augustin Vintzileos, PhD,[4] Shervin Shokouhi, MD,[5] Fernando Miralles-Wilhelm, PhD,[6,7] and Anthony Amoroso, MD[1,2]

▸ Author information  ▸ Article notes  ▸ Copyright and License information    PMC Disclaimer

Eur Respir J. 2020 May; 55(5): 2000517.
Published online 2020 May 7. doi: 10.1183/13993003.00517-2020

PMCID: PMC7144256
PMID: 32269084

No association of COVID-19 transmission with temperature or UV radiation in Chinese cities

*Negative*

Ye Yao,[1,2] Jinhua Pan,[1,2] Zhixi Liu,[1,2] Xia Meng,[1,2] Weidong Wang,[1] Haidong Kan,[1,3] and Weibing Wang[1,3]

▸ Author information  ▸ Article notes  ▸ Copyright and License information    PMC Disclaimer

Alkaissi, H., & McFarlane, S. I. (2023). Artificial hallucinations in ChatGPT: implications in scientific writing. *Cureus*, *15*(2).
Walters, W. H., & Wilder, E. I. (2023). Fabrication and errors in the bibliographic citations generated by ChatGPT. *Scientific Reports*, *13*(1), 14045.
Li, J., Cheng, X., Zhao, W. X., Nie, J. Y., & Wen, J. R. (2023, December). Halueval: A large-scale hallucination evaluation benchmark for large language models. EMNLP 2023

# Tutorial Outline

- Background and Motivation

- Scientific Literature Survey

- Hypothesis Generation and Experiments

- Hands-on Paper Hypothesis Assistant

- Paper Writing

- Paper Draft Evaluation and Ethics

- Summary and Future Directions

# How do we do a literature survey?

- **How do humans do a literature survey?**

  1. Have background knowledge of the topic
  2. Annotate key points within each paper
  3. Look carefully at the figures, diagrams and other illustrations in the paper
  4. Explore relevant previously unread references
  5. Engage in background reading when concepts are unclear
  6. Explore citing papers

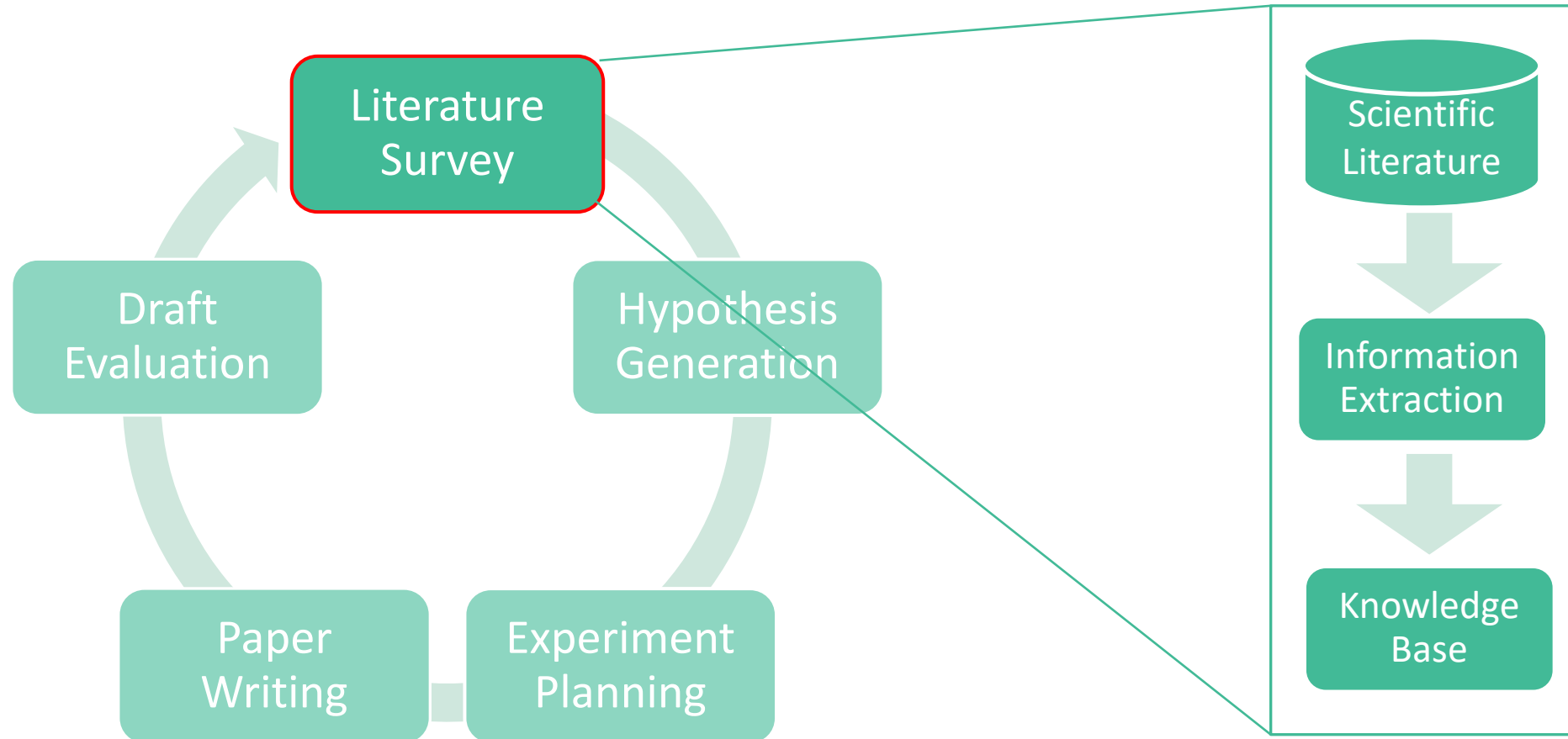- **How do machines do a literature survey?**

  1. Prepare pretrained scientific LLMs along with background knowledge bases in that domain
  2. Construct a multimedia knowledge base from a paper
  3. Use keywords, sentences, or document links to find related papers

Pros:
- Gain a comprehensive understanding of every aspect of the paper
- Analyze and reason over each paper's method and experiments
- Can actively gather relevant knowledge when concepts are unclear

Cons:
- Cannot cover all the latest research papers

Pros:
- Offers broader coverage compared to human

Cons:
- Struggles with aligning knowledge from different modalities
- Faces challenges in mathematical or logical reasoning
- It is often challenging for machines to distinguish between known and unknown knowledge

Keshav, S. (2007). How to read a paper. ACM SIGCOMM Computer Communication Review, 37(3), 83-84.
https://www.cs.jhu.edu/~jason/advice/how-to-read-a-paper.html

# Completing Scientific Paper Lifecycle



Literature Survey

Draft Evaluation

Hypothesis Generation

Paper Writing

Experiment Planning

Scientific Literature

Information Extraction

Knowledge Base

# Why do we want a Scientific Knowledge Base?

- A knowledge base is the backbone of any knowledge-driven AI
  - Knowledge bases provides context to enhance AI's capabilities in semantic understanding and reasoning
  - High-quality, well-maintained knowledge bases help AI to generate more accurate and contextually relevant responses, while minimizing errors and hallucination
  - Knowledge bases help users to interpret and verify AI's decision



Knowledge Base

Wang, Q., Li, M., Wang, X., Parulian, N., Han, G., Ma, J., ... & Onyshkevych, B. (2020). COVID-19 literature knowledge graph construction and drug repurposing report generation. NAACL 2021 Best Demo

# How do we Construct a Scientific Knowledge Base?



**Individual Documents (Papers, Articles, Journals, Webpages)** → PDF Extraction → **Parsed Document** → Information Extraction (Scientific LLMs) → **Scientific Knowledge Base**

# Scientific Knowledge Base Construction



Individual Documents
(Papers, Articles,
Journals, Webpages)

PDF Extraction

Parsed Document

Scientific LLMs

Information Extraction

Scientific Knowledge Base

# PDF Document Extraction

- Extract structured data from a paper PDF
  - The quality of scientific information extraction largely depends on how well the textual contents are extracted from the original PDF file



**Input / Paper PDF**

{

**"title"**: "Full-Resolution Residual Networks for Semantic Segmentation in Street Scenes",
**"authors"**: "Tobias Pohlen, Alexander Hermans, Markus Mathias, Bastian Leibe ",
**"sections"**: ["..."],
**"Captions"**: ...

}

**Output / Paper Metadata JSON**

# PDF Document Extraction

- A document page can be segmented into visual groups of tokens
  - Tokens within each group generally have the same semantic category

- Use pre-trained Faster-RCNN models (Ren et al., 2015) from the LayoutParser (Shen et al., 2021) tool to identify both text lines and blocks based on images of the page

- Use visual layout groups to improve PDF extraction
  - Inject layout indicator tokens into the input
  - Hierarchically encode layout groups

  Pros:
  - Utilizes visual layout in the input
  - Doesn't require pretraining

  Cons:
  - Limited context window due to BERT

  Future Directions:
  - Utilize relationships between layout elements
  - Utilize state-of-the-art multimodal LLMs

Shen, Z., Lo, K., Wang, L., Kuehl, B., Weld, D., & Downey, D. (2022). VILA: Improving Structured Content Extraction from Scientific PDFs Using Visual Layout Groups. TACL, 10, 376–392

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards real-time object detection with region proposal networks. Advances in Neural Information Processing Systems, 28:91–99.

Zeijang Shen, Ruochen Zhang, Melissa Dell, Benjamin Charles Germain Lee, Jacob Carlson, and Weining Li. 2021. LayoutParser: A unified toolkit for deep learning based document image analysis ICDAR 2021

# Scientific Document Extraction Pipeline



**Individual Documents (Papers, Articles, Journals, Webpages)** → PDF Extraction → **Parsed Document** → Information Extraction (Scientific LLMs) → **Scientific Knowledge Base**

# Scientific Large Language Models

- Scientific LLMs are specifically designed to understand scientific literature—a critical task for a literature review
  - Scientific information extraction benefits from pretrained LLMs due to their parametric domain knowledge

Yasunaga, M., Leskovec, J., & Liang, P. (2022). LinkBERT: Pretraining Language Models with Document Links. ACL 2022

# Knowledge-enhanced Scientific LLM



Scientific LLMs have exploded in popularity!

Zhang, Q., Ding, K., Lyv, T., Wang, X., Yin, Q., Zhang, Y., ... & Chen, H. (2024). Scientific large language models: A survey on biological & chemical domains. *arXiv preprint arXiv:2401.14656*.

# What knowledge is useful for a scientific large language model?

- Domain/Task knowledge from pretraining corpus (Gururangan et al., 2020)
  - ➢ Language models benefit from continued pretraining on specific domains
  - ➢ Task-adaptive pretraining on a smaller but task-relevant corpus can boost performance
  - ➢ Automatic or human-curated unlabeled data related to the task can further improve performance in low-resources setting



**PubMed**
**PMC**

4.5B Tokens
13.5 B Tokens

→ BioBERT (Lee et al., 2019)

Semantic **Scholar**

3.17B Tokens
CS +BioMed

→ SciBERT (Beltagy et al., 2019)

Gururangan, S., Marasovic, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., & Smith, N. (2020). Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. ACL2020
Kilicoglu, H. (2018). Biomedical text mining for research rigor and integrity: tasks, challenges, directions. Briefings in bioinformatics, 19(6), 1400-1414.
Beltagy, I., Lo, K., & Cohan, A. (2019). SciBERT: A pretrained language model for scientific text. EMNLP 2019.

- What knowledge is useful for scientific large language models?
  - Domain/Task knowledge from pretraining corpus
  - ➤ Cross document knowledge such as the links between different papers
    - Scientific knowledge can span across documents
    - Understanding a research article depends not only on its content but also on its connections to related papers
    - Document links can provide salient multi-hop knowledge

**Document** 🔗⤏ **Linked document**
(e.g. hyperlink, citation)

[Tidal Basin, Washington D.C.]
The Tidal Basin is a man-made reservoir located between …. It is part of West Potomac Park, is near the National Mall and is a focal point of the National Cherry Blossom Festival held each spring. The Jefferson Memorial, ….

[The National Cherry Blossom Festival] … It is a spring celebration commemorating the March 27, 1912, gift of Japanese cherry trees from Mayor of Tokyo City Yukio Ozaki to the city of Washington, D.C. Mayor Ozaki gifted the trees to enhance …

**Multi-hop knowledge**
(e.g. Tidal Basin has Japanese cherry trees)

Yasunaga, M., Leskovec, J., & Liang, P. (2022). LinkBERT: Pretraining Language Models with Document Links. ACL 2022

# Approaches to Incorporating Document Links in Pretraining

| Link Types | Name | Base Model | Retriever | Domain | Pros | Cons |
|---|---|---|---|---|---|---|
| Hyperlinks/ Citations | BioLinkBERT (Yasunaga et al., 2022) | BERT | N/A | Biomedical | ▪High precision of relevance ▪Capable of capturing relevant documents beyond obvious lexical similarities | ▪Document links can be broken due to removed contents ▪Hyperlinks/Citations might be subjective, influenced by authors' knowledge and preferences |
| Lexical Similarity | BioReader (Frisoni et al., 2022) | T5 | kNN Retriever | Biomedical | ▪Easy to scale ▪No need to modify loss functions | ▪Requires a large database of papers ▪Efficient similarity search becomes bottleneck for retrieval |

Yasunaga, M., Leskovec, J., & Liang, P. (2022). LinkBERT: Pretraining Language Models with Document Links. ACL 2022
Frisoni, G., Mizutani, M., Moro, G., & Valgimigli, L. (2022). BioReader: a Retrieval-Enhanced Text-to-Text Transformer for Biomedical Literature EMNLP2022

# Categorization of Multimodal Pretraining

- Multimodal models have emerged to improve the performance and domain understanding of scientific LLMs

# Categorization of Multimodal Pretraining

- Joint Molecule-Text Model Pretraining

# Joint Molecule Text Models

- Goal: Learn a joint latent representation of both molecules and text
  - Humans learn domain knowledge from studying both molecular structure and biomedical text information
  - Molecular structures provide grounded real-world information, where functional groups and their positions are strong indicators of molecular properties and interactions.
  - Biomedical text provides abundant, high-level and abstract understanding of functions and properties of molecule entities reported from centuries of scientific effort.

Zeng, Z., Yao, Y., Liu, Z., & Sun, M. (2022). A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals. Nature communications, 13(1), 862.

- A current direction: train a language model which can understand instructions in natural language.
  - ➢ The model is pretrained on thousands of relevant tasks
- Utilizes graph-aware positional encodings to help fuse the two data modalities in one model

Valière, N. (2002). GIMLET: a computer program for analysing genetic individual identification data. Molecular Ecology Notes, 2(3), 377-379.

# Categorization of Multimodal Pretraining

- Multimodal Scientific Figure Understanding



**Chemical Related IE Benchmark**

# Why do we need to understand Scientific Figures?

- Scientific figures in papers usually provide critical information beyond pure text
  - Plotting the data allows us to see underlying structures that aren't obvious from a table
  - Scientific figures allow the author to simplify interpretation and drawing conclusions for the reader

Hsu, T. Y., Giles, C. L., & Huang, T. H. K. (2021). SciCap: Generating captions for scientific figures. EMNLP2021 Findings.

# Scientific Multimodal Instruction Tuning

- SciTune: Developing multimodal agents to reason across multiple scientific disciplines
  - Incorporate human-generated scientific instructions based on SciCap (Hsu et al., 2021)
  - LLaMA (Touvron et al., 2023) is used as the LLM decoder and CLIP (Radford et al., 2021) as visual encoder

Horawalavithana, S., Munikoti, S., Stewart, I., & Kvinge, H. (2023). Scitune: Aligning large language models with scientific multimodal instructions. arXiv preprint arXiv:2307.01139.
Hsu, T. Y., Giles, C. L., & Huang, T. H. K. (2021). SciCap: Generating captions for scientific figures. EMNLP2021 Findings.
Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M. A., Lacroix, T., ... & Lample, G. (2023). Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark,et al. 2021. Learning transferable visual models from natural language supervision. ICML 2021
Lu, P., Mishra, S., Xia, T., Qiu, L., Chang, K. W., Zhu, S. C., ... & Kalyan, A. (2022). Learn to explain: Multimodal reasoning via thought chains for science question answering. Advances in Neural Information Processing Systems, 35, 2507-2521.

- **SciTune (blue)** outperforms **BLIP (red)** in caption generation



Comparison of the effect of the path loss exponent $\alpha$ on rates achieved by both transmitters, M = 4.

a plot of a line graph with a blue line and red line.

**Graph Plot** The sum-rate and sum-rate of RRM-RRM with respect to the path loss exponent $\gamma$ for the two cases: $\gamma = 2$ and $\gamma = 3$.

$\gamma$ in the caption should be $\alpha$



Conceptual diagram of nonlinear adaptive method developed to control the HCV epidemic in the existence of uncertainties on parameters of the model.

a diagram of a block diagram of a nuclear system.

**Node Diagram** Block diagram of the proposed non-linear SIR epidemic model with adaptive controllers.

SIR is not shown in the node diagram

**Points for improvement:**
- The model still generates hallucinations and performs poorly on node diagrams
- Integrating better OCR and contextual text/background knowledge might help improve performance

Horawalavithana, S., Munikoti, S., Stewart, I., & Kvinge, H. (2023). Scitune: Aligning large language models with scientific multimodal instructions. arXiv preprint arXiv:2307.01139.
Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. ICML 2022.

# Categorization of Multimodal Pretraining

- Knowledge/Pathway augmented Pretraining

# Why do we need an external networks?

- Unaugmented LLMs struggle to achieve satisfactory performance on knowledge-intensive tasks such as biomedical NLP

- Scientific documents contain many highly specialized terms, acronyms, and abbreviations. Their definitions and properties are often not presented in context during pretraining

- External knowledge can help LLMs reduce hallucinations

Lai, T. M., Zhai, C., & Ji, H. (2023). KEBLM: knowledge-enhanced biomedical language models. Journal of Biomedical Informatics, 143, 104392.
Agrawal, G., Kumarage, T., Alghami, Z., & Liu, H. (2023). Can knowledge graphs reduce hallucinations in LLMs?: A survey. arXiv preprint arXiv:2311.07914.

- Core idea: use adaptor models to memorize knowledge from multiple external databases into a pretrained LLM



**Adapter**

**Pretraining Adapter Modules**

**Fusing Knowledge from Adapter Modules**

- Existing methods only utilize a single source of knowledge. KEBLM is pretrained on multiple sources of biomedical domain knowledge.
- An adapter module and fusion layer memorize and combine the knowledge in a self-supervised way.
- The adapter is pretrained with three types of knowledge: entity descriptions, entity-entity relations, and entity synonyms

Lai, T. M., Zhai, C., & Ji, H. (2023). KEBLM: knowledge-enhanced biomedical language models. Journal of Biomedical Informatics, 143, 104392.

# Knowledge-Enhanced Biomedical Language Models



**Adapter**

**Pretraining Adapter Modules**

Entity Descriptions → Adapter Module 1

Entity-Entity Relations → Adapter Module 2

(Ammonia, NH3)
(COVID-19, SARS-COV2)
(EGFR, epidermal growth factor receptor)
**Synonyms** → Adapter Module 3

**Fusing Knowledge from Adapter Modules**

**Points for improvement:**
- The model is built on encoder-only model, which lacks the ability to generate sentences
- The model ignores molecule/protein structures, a crucial source of knowledge for biomedical papers

Lai, T. M., Zhai, C., & Ji, H. (2023). KEBLM: knowledge-enhanced biomedical language models. Journal of Biomedical Informatics, 143, 104392.

- **Notable points**
  - LLMs benefit from domain/task pretraining
  - Cross-document knowledge can help LLMs capture multi-hop knowledge more effectively
  - Adding modalities to LLMs can ground their understanding into the real world

- **Future directions**
  - Incorporate hierarchical structures of papers into scientific LLM pretraining
  - Improve the alignment between different modalities within LLMs
  - Enhance reasoning ability of scientific LLMs during pretraining
  - Further explore the use of tables and tools/code in multimodal scientific LLMs

# Scientific Information Extraction Pipeline



**Individual Documents (Papers, Articles, Journals, Webpages)**

PDF Extraction

**Parsed Document**

Scientific LLMs

Information Extraction

**Scientific Knowledge Base**

# Scientific Information Extraction

- Transform *unstructured* or *semi-structured* scientific texts into a *structured* form that can be more easily understood
  - Identifying and structuring
    - ➢ Entities
    - ➢ Quantities (and their scope)
    - ➢ Events
    - ➢ Relations, etc

- $BaCl2\ 2H2O$ and $FeCl3\ 6H2O$ are used as initial materials. In a typical experiment, 0.487 g of $FeCl_3\ 6H_2O$ and 0.054 g of $BaCl_2\ 2H_2O$ (Fe/Ba mole ratio: N=8) were dissolved in 21 mL of distilled water. Then, 4.0 g of KOH was added to the solution with continuous stirring. Finally, the solution was put into a 30 mL of Teflon-lined stainless-steel autoclave and kept it in an oven at 220 °C for 24 h, followed by furnace cooling to room temperature.

Zhao, L., Lv, X., Wei, Y., Ma, C., & Zhao, L. (2013). Hydrothermal synthesis of pure BaFe12O19 hexaferrite nanoplatelets under high alkaline system. *Journal of Magnetism and Magnetic Materials*, *332*, 44-47.

'Nickel(II) nitrate hexahydrate (0.134 M, Ni(NO3)2·6H2O, Sigma-Aldrich, purity >99.9%) and urea (Sigma-Aldrich, purity >99.5%) (Ni/urea molar ratio 1/12 mol/mol) were solubilized and mixed in 0.2 L of water with magnetic stirring for 240 min at 80 °C.'

- Each specific niche within a subfield of a scientific field uses unique jargon and processes

- Scientific concepts and relations can be difficult to understand

- Rare domain-specific entities and events are common so they must be discovered without supervision.

with magnetic stirring

at 80 °C

for 240 min

in 0.2 L of water

solubilized → mixed

Nickel(II) nitrate hexahydrate

urea

(Sigma-Aldrich, purity >99.5%)

(0.134 M, Ni(NO3)2·6H2O, Sigma-Aldrich, purity >99.9%))

(Ni/urea molar ratio 1/12 mol/mol)

Lucy, L., Dodge, J., Bamman, D., & Keith, K. A. (2022). Words as gatekeepers: Measuring discipline-specific terms and meanings in scholarly publications. ACL 2023 Findings.

# Challenges in Scientific Information Extraction

- Few benchmarks are publicly available
  - Models only have a limited number of training samples for each task
  - Annotations are usually incomplete or have low coverage

| Entity Extraction | Relation Extraction | Event Extraction | Entity Linking |
|---|---|---|---|
| BC4CHEMD | Abroad-RE | BioNLP 09 | XL-BEL |
| BC5CDR | BC5CDR | Genia 2013 | Mantra GSC |
| BC7 NLM-Chem | BC6 ChemProt | Genia 2016 | |
| BioRED | BC7 DrugProt | MatSci-NLP | |
| CHEMET | BioRED | | |
| Chem-FINESE | BioRelEx | | |
| JNLPBA | JNLPBA | | |
| MatSci-NLP | MatSci-NLP | | |
| PolymerAbstracts | | | |

## Chemical Related IE Benchmark

# Type of Augmented Scientific IE

- To address low-resource settings and manage long-tail distributions, there is growing focus on using structured knowledge bases to augment representations or expand the training set

Open Domain

External Knowledge

(extract a rich but possibly sparse and noisy graph representation)

(use external knowledge to help boost performance)

How do we link text to external knowledge ?

How do we leverage ontology-based knowledge in IE?

# Scientific Entity Linking

| Task | Modality | Goal |
|---|---|---|
| Scientific Entity Linking (Leaman et al., 2016) | Monolingual Text  | Map mentions in text to entities in a knowledge base |
| Scientific Cross-lingual Entity Linking (Bitton et al., 2020) | Multilingual Text  | Map entity mentions in text of a source language to entities in a knowledge base, (e.g., UMLS), in a target language |
| Scientific Table Entity Linking (Lou et al., 2023) | Text and Table  | Map entity mentions in text and tables to entities in a knowledge base |

Leaman, R., & Lu, Z. (2016). TaggerOne: joint named entity recognition and normalization with semi-Markov Models. Bioinformatics, 32(18), 2839-2846.
Bitton, Y., Cohen, R., Schifter, ..., & Elhadad, N. (2020). Cross-lingual Unified Medical Language System entity linking in online health communities. Journal of the American Medical Informatics Association, 27(10), 1585-1592.
Lou, Y., Kuehl, B., Bransom, E., Feldman, S., Naik, A., & Downey, D. (2023). S2abEL: a dataset for entity linking from scientific tables. EMNLP 2023.

How do we link text to external knowledge ?

How do we leverage ontology-based knowledge in IE?

- Utilize an entity linker to transfer background knowledge from an external KB to jointly extract entities and relations

- Use a bidirectional graph convolutional network (GCN) to fuse global relational information into local representation for each span representation



**Pros:**
- Incorporates information from both local context and relevant background knowledge to produce joint representations

**Points for improvement:**
- The performance of the system depends on the effectiveness of entity linker
- The paper only utilizes background knowledge from a single knowledge base

Lai, T., Ji, H., Zhai, C., & Tran, Q. (2021). Joint Biomedical Entity and Relation Extraction with Knowledge-Enhanced Collective Inference ACL 2021

# Citation Enhanced IE

- Use a citation graph of referential links papers to augment text representations



Pros:
- Citation graphs provide additional background information between the target paper and its cited papers

Points for improvement:
- Citation sentences are treated as a new section of the document without any special model design.

Viswanathan, V., Neubig, G., & Liu, P. (2021). Citationie: Leveraging the citation graph for scientific information extraction. ACL2021

# Weak Supervision for Relation Extraction

- Can we use weak supervision for relation extraction?

- Case study: Reaction extraction from chemistry papers
  - Transform the task of reaction extraction into a question answering task
  - Synthetic data generation
    - Use frequent patterns within the text as linguistic cues to identify chemical reactions by starting with seed patterns
    - Extract reactions from patent literature based on rules



**Points for improvement:**
- Performance is limited by external knowledge and accuracy of AMR graphs

Zhong, M., Ouyang, S., Jiang, M., Hu, V., Jiao, Y., Wang, X., & Han, J. (2023). Reactie: Enhancing chemical reaction extraction with weak supervision. ACL 2023 Findings.

- What is the difference between using Abstract Meaning Representation (AMR) and an External Knowledge Base?
  - AMR is a semantic representation language that converts the meaning of each input sentence into a rooted, directed, labeled, acyclic graph structure
  - AMR is more versatile in information coverage and is not limited to domain-specific ontologies

- Sentences that have the same basic meaning often have the same AMR graph

- An event graph is a subgraph of an AMR graph
  - Transform event extraction into a subgraph identification problem



Can we use both open domain and ontology-based knowledge?

Rao, S., Marcu, D., Knight, K., & Daumé III, H. (2017, August). Biomedical event extraction using abstract meaning representation. In BioNLP 2017 (pp. 126-135).

# Knowledge-enriched AMR Enhanced IE

- Combine domain-specific knowledge with semantic knowledge by merging a subgraph from KG with the AMR graph
  - Enriches the model with external knowledge
  - Uses message passing with an edge-conditioned graph attention network

Pros:
- Enriches the AMR graph with external knowledge to capture long distance between event triggers and entities

Points for improvement:
- Incorporate entity definition sentence for each entity in the sentence KG can further improve performance

Zhang, Z., Parulian, N., Ji, H., Elsayed, A., Myers, S., & Palmer, M. (2021). Fine-grained Information Extraction from Biomedical Literature based on Knowledge-enriched Abstract Meaning Representation. ACL 2021

# How can we do information extraction in low-resource settings?

- Use indirect supervision to transfers supervision signals from a more resource-rich task (NLI) to enhance a more resource-limited task (biomedical RE)
  - Use the input sentence as the premise while converting each relation label into template-based natural language hypotheses for NLI
  - Propose a new ranking-based loss to balance positive and negative relations



Input: Androgen antagonistic effect of **estramustine phosphate** (EMP) metabolites on wild-type and mutated **androgen receptor**.

(1) Verbalizer

Premise = masked input

Androgen antagonistic effect of **@CHEMICAL$** (EMP) metabolites on wild-type and mutated **@GENE$**.

Hypothesis:

**@CHEMICAL$** and **@GENE$** have no relation.

**@CHEMICAL$** is identified as an antagonist of **@GENE$**.

Upregulator **@CHEMICAL$** is activated by **@GENE$**.

(2) Training

Entailment score — Rank Over — Neg / Positive / Neg

(3) Inference

Entailment score

**Pros:**
- Use natural language templates to adapt NLI models for the relation extraction task
- Existing entailment datasets transfers well

**Points for improvement:**
- Utilize task-task relations to further improve indirect supervision
- Can we automatically generate templates for arbitrary relations?

Xu, J., Ma, M. D., & Chen, M. (2023). Can NLI Provide Proper Indirect Supervision for Low-resource Biomedical Relation Extraction? ACL 2023

# Self-Validation

- Use a self-validation module to reconstruct an original sentences based on entity extraction results to add additional regularization for entity extraction model

- Use a contrastive loss on the entity decoder to reduce copying from the original sentence

Pros:
- Exhibits strong few-shot performances in extreme low resource settings
- Doesn't require any domain-specific pretraining or external KB

Points for improvement:
- Given that the self-validation module remains static after the initial pretraining stage, exploring cyclic improvements of both the module will be an interesting future direction

*Through application of ligand screening, we describe the first examples of Pd-catalyzed Suzuki–Miyaura reactions using aryl sulfamates at room temperature.*

**Entity Extraction**
- Encoder A
- Decoder A → Entity Decoder Contrastive Loss

*ligand <Ligands>, Pd-catalyzed Suzuki-Miyaura reactions <Coupling reactions>, ...* → Supervised Loss

Reconstruction Loss

**Self-Validation**
- Encoder B
- Decoder B

*Through application of ligand screening, we describe the first examples of Pd-catalyzed Suzuki–Miyaura reactions using aryl sulfamates at room temperature.*

# Scientific IE Applications

# Application 1: Visualization of Scientific Papers

- Design of intelligent, highlight based skimming user interfaces, based on formative interviews and preliminary usability studies of a prototype tool

Fok, R., Kambhamettu, H., Soldaini, L., Bragg, J., Lo, K., Hearst, M., ... & Weld, D. S. (2023, March). Scim: Intelligent skimming support for scientific papers. In Proceedings of the 28th International Conference on Intelligent User Interfaces (pp. 476-490).

- An opensource Python toolkit for analyzing and processing visually-rich, structured scientific documents

```
A  >>> doc.paragraphs[0]
```

```
B  >>> doc.paragraphs[0].sentences[2]
   or
   >>> doc.sentences[2]
```

```
C  >>> doc.sentences[2].tokens[9:13]
   or
   >>> doc.tokens[169:173]
```

```
D  >>> doc.figures[0]
```

```
E  >>> doc.captions[0]
```

```
F  >>> user_query = Box(l,t,w,h, page=0)

   >>> selected_tokens =
    doc.find(user_query, layer="tokens")

   >>> [token.text for
               token in selected_tokens]

   ["Techniques", "for", "collecting",
   "labeled", "data", "perts", "for",
   "manual", "annotation", ...]
```

**ABSTRACT**

Crowdsourcing provides a scalable and efficient way to construct labeled datasets for training machine learning systems. However, creating comprehensive label guidelines for crowd-workers is often prohibitive even for seemingly simple concepts. Incomplete or ambiguous label guidelines can then result in differing interpretations of concepts and inconsistent labels. Existing approaches for improving label quality, such as worker screening or detection of poor work, are ineffective for this problem and can lead to rejection of honest work and a missed opportunity to capture rich interpretations about data. We introduce *Revolt*, a collaborative approach that brings ideas from expert annotation workflows to crowd-based labeling. Revolt eliminates the burden of creating detailed label guidelines by harnessing crowd disagreements to identify ambiguous concepts and create rich structures (groups of semantically related items) for post-hoc label decisions. Experiments comparing Revolt to traditional crowdsourced labeling show that Revolt produces high quality labels without requiring label guidelines in turn for an increase in monetary cost. This up front cost, however, is mitigated by Revolt's ability to produce reusable structures that can accommodate a variety of label boundaries without requiring new data to be collected. Further comparisons of Revolt's collaborative and non-collaborative variants show that collabvoration reaches higher label accuracy with lower monetary cost.

**ACM Classification Keywords**

H.5.m. Information Interfaces and Presentation (e.g. HCI): Miscellaneous

**Author Keywords**

crowdsourcing; machine learning; collaboration; real-time

**INTRODUCTION**

From conversational assistants on mobile devices, to facial



**Figure 1.** Revolt creates labels for unanimously labeled "certain" items (e.g., *cats* and *not cats*), and surfaces categories of "uncertain" items enriched with crowd feedback (e.g., *cats and dogs* and *cartoon cats* in the dotted middle region are annotated with crowd explanations). Rich structures allow label requesters to better understand concepts in the data and make post-hoc decisions on label boundaries (e.g., assigning cats and dogs to the cats label and cartoon cats to the *not cats* label) rather than providing crowd-workers with a priori label guidelines.

learned models that must be trained on representative datasets labeled according to target concepts (e.g., speech labeled by their intended commands, faces labeled in images, emails labeled as spam or not spam).

Techniques for collecting labeled data include recruiting experts for manual annotation [51], extracting relations from readily available sources (e.g., identifying bodies of text in parallel online translations [46, 13]), and automatically gener- ating labels based on user behaviors (e.g., using dwell time to implicitly mark search result relevance [2]). Recently, many practitioners have also turned to crowdsourcing for cre- ating labeled datasets at low cost [49]. Successful crowd-

https://github.com/allenai/papermage

Lo, K., Shen, Z., Newman, B., Chang, J. Z., Authur, R., Bransom, E., ... & Soldaini, L. (2023, December). PaperMage: A Unified Toolkit for Processing, Representing, and Manipulating Visually-Rich Scientific Documents. EMNLP 2023 Demo

# Conclusions and Future Directions for IE

**Moving forward**

- How to merge knowledge from different sources?
- How to induce an ontology?

**Low-resource IE without KB**

- Indirect supervision
- Self-validation

**Knowledge-enhanced IE**

- Entity linking
- External domain knowledge
- Citation knowledge
- AMR graph

**Scientific Large Language Model**

- Domain/Task Specific Pretraining
- Cross document augmented Large Language Model

# Completing Scientific Paper Lifecycle

# Why do we need literature search?

- Definition: Literature search is the process of retrieving scientific articles to satisfy specific information needs.
  - This is an Information Retrieval task.

- Literature search can benefit both human and knowledge-driven AI
  - Literature search allows researchers to quickly find relevant studies, reducing the risk of repeating work and missing critical information
  - Literature search can highlight the appropriateness or shortcomings of previous research methodologies
  - Literature search helps identify research trends and gaps in the existing literature, which can direct future studies
  - Literature search systems are crucial for synthesizing evidence in systematic reviews and meta-analyses

Jin, Q., Leaman, R., & Lu, Z. (2024). PubMed and beyond: biomedical literature search in the age of artificial intelligence. Ebiomedicine, 100.

# Why do we need retrieval if we're using LLMs?

- LLMs can't memorize all knowledge in their parameters (especially long-tail events)

- LLMs' knowledge is easily outdated and hard to update

- LLMs' output is challenging to interpret and verify

# A real-world crisis: Digesting COVID-19 Papers

- Practical progress at combating COVID-19 highly depended on effective transmission, assessment, and extension of research results
  - 2.7K new papers per day
  - As of June 13, 2020, there were at least 140K papers about coronavirus

- There was an urgent need to effectively retrieve relevant evidence from this large corpus



CORD-19 dataset # daily archives



Count by Source
Click and drag to zoom

Wang, L. L., Lo, K., Chandrasekhar, Y., Reas, R., Yang, J., Burdick, D., ... & Kohlmeier, S. (2020). Cord-19: The covid-19 open research dataset. NLP COVID-19 Workshop.
https://icite.od.nih.gov/covid19/search/

- **For example, consider manual drug repurposing:**
  - Current clinical trials for drug repurposing mainly rely on symptoms
  - There's too many drug candidates
  - There's too much misinformation about effects
  - It's too costly to test all drugs and difficult to quantify success

Summary of *in vitro* and *in vivo* studies of the use of CQ/HCQ in COVID-19.

| Study | Study type | Objective | Outcome | Result/s | Conclusion |
|---|---|---|---|---|---|
| Wang et al.[35] | In vitro | To evaluate antiviral efficiency of ribavirin, penciclovir, nitazoxanide, nafamostat, chloroquine, remdesivir and favipiravir against *in vitro* SARS-CoV 2 | Drug efficacy was quantified *via* viral numbers in cell supernatant (qRT-PCR) and viral nucleoprotein expression (immunofluorescence microscopy) | CQ blocked the virus at low concentrations during both entry and post-entry phases of cellular infection | CQ has potential for clinical use against the SARS-CoV 2 due to potent blocking viral infection demonstrating its antiviral effects, however, *in vivo* studies are warranted |
| Yao et al.[36] | In vitro | To test CQ/HCQ *in vitro* activity against SARS-CoV 2 infected Vero cells | Drug efficacy was quantified using the detection of viral RNA *via* RT-PCR | A twice-daily HCQ (loading dose of 400 mg and maintenance dose of 200 mg) for four days achieved three times the potency of the standard 500 mg CQ given in advance for five days | HCQ has higher potency than CQ in inhibiting SARS-CoV 2 *in vitro* |
| Chen et al.[37] | Pilot | To evaluate HCQ against COVID-19 | Negative SARS-CoV 2 nucleic acid conversion rate *via* respiratory pharyngeal swab | More patients in the control group achieved a higher negative pharyngeal swab of viral nucleic acid at 7 days (93.3%, *n* = 14 *versus* 86.7%, *n* = 13, *p* > 0.05). | Although HCQ showed good prognosis in moderate COVID-19 infection, further investigation is needed with larger sample sizes and better endpoints |

Acharya, Y., & Sayed, A. (2020). Chloroquine and hydroxychloroquine as a repurposed agent against COVID-19: a narrative review. Ther Adv Infect 7: 2049936120947517.
Wang, Q., Li, M., Wang, X., Parulian, N., Han, G., Ma, J., ... & Onyshkevych, B. (2020). COVID-19 literature knowledge graph construction and drug repurposing report generation. NAACL 2021 Best Demo

# How do we retrieve relevant evidence?

- **How do humans retrieve relevant evidence?**
  - Do creative web search
    - ➢ Experiment with several searches
    - ➢ Put yourself in an author's shoes; what phrases might they have used?
    - ➢ Specifically search Google Scholar, etc.
  - Track down related work (once you have a relevant paper)
    - ➢ Follow the bibliography to earlier papers
    - ➢ See who else has cited the work
  - Read each paper to find related work

Pros:
- Capable of reasoning over related topics more effectively than machines
- Achieves high precision when retrieving evidence

Cons:
- Unable to cover all relevant research papers and evidence

- **How do machines retrieve relevant evidence?**
  - Construct knowledge graphs for each paper
    - ➢ Search for keywords based on the constructed knowledge graph
  - Use papers' references and citation networks to identify related papers
  - Search for semantically relevant sentences based on embedding similarity and overlapping knowledge graphs

Pros:
- Offer broader coverage compared to human

Cons:
- Typically exhibits low precision
- Usually ignores multimodal information and deep connections between concepts

Keshav, S. (2007). How to read a paper. ACM SIGCOMM Computer Communication Review, 37(3), 83-84.
https://www.cs.jhu.edu/~jason/advice/how-to-read-a-paper.html

# Scientific Information Retrieval Categories

**Corpus-level Information Retrieval**

- Literature Search Engines (Google Scholar, PubMed)
- KG-Augmented Retrieval (PubTator 3, ESRA)

**Sentence/Concept-level Information Retrieval**

- Sentence Retrieval (LitSense, EvidenceMiner, SPIKE)
- Dataset Recommendation (Datafinder)

**Multimodal Information Retrieval**

- Multi-modal Molecule Text Retrieval (Text2mol, MoleculeSTM)

**Applications**

- Exploratory Search of the COVID-19 Literature (Scisight)
- Fact-Checking (Check-COVID)
- Chemistry Prediction (TextReact)

# Scientific Information Retrieval Categories

- Essential for Scientific Literature Surveys

- Narrow Search Spaces in Retrieval Augmented Generation

| Corpus-level Information Retrieval | Sentence/Concept-level Information Retrieval | Multimodal Information Retrieval | Applications |
|---|---|---|---|
| Literature Search Engines (Google Scholar, PubMed) | Sentence Retrieval (LitSense, EvidenceMiner, SPIKE) | Multi-modal Molecule Text Retrieval (Text2mol, MoleculeSTM) | Exploratory Search of the COVID-19 Literature (Scisight) |
| KG-Augmented Retrieval (PubTator 3, ESRA) | Dataset Recommendation (Datafinder) | | Fact-Checking (Check-COVID) |
| | | | Chemistry Prediction (TextReact) |

# Literature Search Engines

- Indexes the full text or metadata of scholarly literature across an array of publishing formats and disciplines



**Problem:**

- Search engines can only process short keyword-based queries, returning a list of raw articles without further analysis

https://pubmed.ncbi.nlm.nih.gov/?term=covid-19&filter=years.1981-2024

Jin, Q., Leaman, R., & Lu, Z. (2024). PubMed and beyond: biomedical literature search in the age of artificial intelligence. Ebiomedicine, 100.

# KG-Augmeted Retrieval: PubTator 3.0

- The PubTator 3.0 pipeline includes three major steps: (i) named entity recognition for six key biomedical entities: genes, diseases, chemicals, genetic variants, species, and cell lines, (ii) entity linking, and (iii) relation extraction with 12 common relation types



Pros:
- Provides more accurate information retrieval results compared to PubMed and Google Scholar
- Entities are linked to external knowledge bases

Cons:
- Relation extraction is limited to abstracts only
- Frequently miss new entities in the entity extraction process
- The whole knowledge graph of the paper isn't visualized

Wei, C. H., Allot, A., Lai, P. T., Leaman, R., Tian, S., Luo, L., … & Lu, Z. (2024). PubTator 3.0: an AI-powered literature resource for unlocking biomedical knowledge. Nucleic Acids Research, gkae235.

- **ESRA: A literature discovery platform that augments search results with relevant details and explanations**
  - Perform scientific information extraction to construct knowledge graph
  - Expand the query based on similar entities from background knowledge graph
  - Use elastic search and citation counts to rank papers
  - Provide ranking explanations based on selected sentences in paper abstracts



**Pros:**
- Provides explanations, graph visualizations, and facts to increase interpretability

**Points for Improvement:**
- Information extraction is limited to abstracts only
- Nodes in the knowledge base are not linked to human curated KG

Hongwimol, P., Kehasukcharoen, P., Laohawarutchai, P., Lertvittayakumjorn, P., Ng, A. B., Lai, Z., ... & Vateekul, P. (2021, August). ESRA: Explainable scientific research assistant. ACL 2021 Demo

- **Notable points**
  - Leveraging both structured and elastic search can improve document retrieval accuracy

- **Future directions**
  - Need to find a balance between traditional keywords-based search and KG-based retrieval to save computational  power
  - Existing systems lack the ability to uncover deep connections between entities during information retrieval due to the absence of reasoning steps in the process

# Scientific Information Retrieval Categories

- Sentence/concept retrieval results can be provided to Retrieval Augmented Generation for downstream tasks

| Corpus-level Information Retrieval | Sentence/Concept-level Information Retrieval | Multimodal Information Retrieval | Applications |
|---|---|---|---|
| Literature Search Engines (Google Scholar, PubMed) | Sentence Retrieval (LitSense, EvidenceMiner, SPIKE) | Multi-modal Molecule Text Retrieval (Text2mol, MoleculeSTM) | Exploratory Search of the COVID-19 Literature (Scisight) |
| KG-Augmented Retrieval (PubTator 3, ESRA) | Dataset Recommendation (Datafinder) | | Fact-Checking (Check-COVID) |
| | | | Chemistry Prediction (TextReact) |

- Rank sentences based on two methods
  - A traditional term-weighting approach that up-weights sentences that contain more of the rare terms in the user query
  - A neural embedding approach that retrieves weights based on semantic similarity



Allot, A., Chen, Q., Kim, S., Vera Alvarez, R., Comeau, D. C., Wilbur, W. J., & Lu, Z. (2019). LitSense: making sense of biomedical literature at sentence level. Nucleic acids research, 47(W1), W594-W599.

# EvidenceMiner: Textual Evidence Discovery

- Given a query, retrieve an evidence sentence from corpus
  - Extracts entities with distant supervision based on knowledge bases to provide keywords for retrieval
  - Retrieves sentences based on query word, entity, and meta-pattern matching
  - Requires no human annotation and achieves better performance for automatic textual evidence discovery

Wang, X., Guan, Y., Liu, W., Chauhan, A., Jiang, E., Li, Q., ... & Han, J. (2020, July). Evidenceminer: Textual evidence discovery for life sciences. ACL 2020: Demo.

# SPIKE: Neural Extractive Search

- Enrich the search query with capture-slots to allow for rapid extraction
  - Use standard dense passage retrieval methods to retrieve relevant sentences
  - Incorporate a neural argument alignment model to align the argument in the query that corresponds to the capture spans on each sentence

Ravfogel, S., Taub-Tabib, H., & Goldberg, Y. (2021). Neural extractive search. ACL 2021 Demo.

# Summary

| Name | Pros | Cons |
|---|---|---|
| LitSense (Allot et al., 2019) | ▪ Easy to setup<br>▪ Require minimal resource | ▪ Performs worst among three methods |
| Evidenceminer (Frisoni et al., 2022) | ▪ Leverages knowledge graph to annotate sentences without human annotation<br>▪ Can switch to approximate matching when strict query matching fails to find sufficiently high-quality answers | ▪ Relies on external KBs for distant supervision<br>▪ Relies on pattern matching for relation extraction |
| SPIKE (Ravfogel et al., 2021) | ▪ Enriches the search query into a knowledge graph<br>▪ Provides Boolean and syntactic search | ▪ Suffers from speed and scalability due to dense retrieval system |

Allot, A., Chen, Q., Kim, S., Vera Alvarez, R., Comeau, D. C., Wilbur, W. J., & Lu, Z. (2019). LitSense: making sense of biomedical literature at sentence level. Nucleic acids research, 47(W1), W594-W599.
Wang, X., Guan, Y., Liu, W., Chauhan, A., Jiang, E., Li, Q., ... & Han, J. (2020, July). Evidenceminer: Textual evidence discovery for life sciences. ACL 2020: Demo.
Ravfogel, S., Taub-Tabib, H., & Goldberg, Y. (2021). Neural extractive search. ACL 2021 Demo.

# Scientific Dataset Recommendation

- Datasets are hard to directly index for search and there are no corpora available for this task
  - Solution: Create the dataset retrieval task
    - Build the DataFinder Dataset consists of a larger automatically-constructed training set (17.5K queries) and a smaller expert annotated evaluation set (392 queries)
    - Develop a bi-encoder retriever for text-based dataset recommendation, which finds more relevant search results than existing third-party dataset search engines



**Points for Improvement:**
- Methods tested in the system are not state-of-the-art
- Papers in the test set and training set are not checked for paper publication year, resulting in potential information leaks
- The whole system relies on Papers with Code datasets

Viswanathan, V., Gao, L., Wu, T., Liu, P., & Neubig, G. (2023). Datafinder: Scientific dataset recommendation from natural language descriptions. ACL 2023.

# Scientific Information Retrieval Categories

- The text modality usually suffers from reporting bias, *the tendency of people to not state the obvious*

- Multimodal information retrieval can address this limitation

| Corpus-level Information Retrieval | Sentence/Concept-level Information Retrieval | Multimodal Information Retrieval | Applications |
|---|---|---|---|
| Literature Search Engines (Google Scholar, PubMed) | Sentence Retrieval (LitSense, EvidenceMiner, SPIKE) | Multi-modal Molecule Text Retrieval (Text2mol, MoleculeSTM) | Exploratory Search of the COVID-19 Literature (Scisight) |
| KG-Augmented Retrieval (PubTator 3, ESRA) | Dataset Recommendation (Datafinder) | | Fact-Checking (Check-COVID) |
| | | | Chemistry Prediction (TextReact) |

Paik, C., Aroca-Ouellette, S., Roncone, A., & Kann, K. (2021). The world of an octopus: How reporting bias influences a language model's perception of color. EMNLP 2021.

- Text2Mol Cross-modal text-molecule information retrieval directly from natural language descriptions to molecules

  - Allows semantic search of molecules based on high-level properties and descriptions

  - Cross-modal attention-based association rules are created between molecules and text to improve explainability

  - A new benchmark dataset with 33,010 text-compound pairs for cross-modal text-molecule IR is released

Water is an oxygen hydride consisting of an oxygen atom that is covalently bonded to two hydrogen atoms.

Query

1. $H_2O$    2. $H_2O_2$    3. $C_2H_6O$

Negative Sampling Loss

MLP

Cross-modal Attention

Transformer Decoder

Mean Pooling

SciBERT

GCN

Token Features

Description Tokens

Adjacency Matrix

**Points for Improvement:**
- Incorporate additional external knowledge to the system

Edwards, C., Zhai, C., & Ji, H. (2021, November). Text2mol: Cross-modal molecule retrieval with natural language queries. EMNLP 2021.

**MoleculeSTM increases capabilities by training on a larger dataset**

- Constructs a large multi-modal dataset, PubChemSTM, with over 280,000 chemical structure-text pairs
- Adds capability for property prediction and molecule editing tasks

**Pros:**

- Open vocabulary
  - Not limited to a fixed set of pre-defined molecule-related textual descriptions
- Compositionality
  - Decompose a complex concept into several simple concepts



(a) Contrastive Pretraining

Aspirin is a commonly used drug for the treatment of pain and fever.

Penicillin G Sodium is the sodium salt form of benzylpenicil-lin.

(b) Structure-text Retrieval

This molecule is for antiinflammatory preparations.

This molecule is for diabetes.

This molecule is for gastrointestinal disorders.

similarity score

(c) Text-based Molecule Editing

A pretrained generative model

This molecule has high permeability.

adaptor module

(d) Molecular Property Prediction

Latent Representation of Chemical Structure

Latent Representation of Textual Description

Latent Representation of Generative Model

Joint Latent Representation

Liu, S., Nie, W., Wang, C., Lu, J., Qiao, Z., Liu, L., ... & Anandkumar, A. (2023). Multi-modal molecule structure–text model for text-based retrieval and editing. Nature Machine Intelligence, 5(12), 1447-1457.

# Scientific Information Retrieval Categories

**Corpus-level Information Retrieval**
- Literature Search Engines (Google Scholar, PubMed)
- KG-Augmented Retrieval (PubTator 3, ESRA)

**Sentence/Concept-level Information Retrieval**
- Sentence Retrieval (LitSense, EvidenceMiner, SPIKE)
- Dataset Recommendation (Datafinder)

**Multimodal Information Retrieval**
- Multi-modal Molecule Text Retrieval (Text2mol, MoleculeSTM)

**Applications**
- Exploratory Search of the COVID-19 Literature (Scisight)
- Fact-Checking (Check-COVID)
- Chemistry Prediction (TextReact)

# Application: COVID-19 Exploratory Scientific Search

- Explore associations between biomedical concepts from papers (e.g., genes, drugs, diseases, patient outcomes)

- Combine textual and network information to search and visualize groups of researchers and connections between them

Hope, T., Portenoy, J., Vasan, K., Borchardt, J., Horvitz, E., Weld, D. S., ... & West, J. (2020). SciSight: Combining faceted navigation and research group detection for COVID-19 exploratory scientific search. EMNLP 2020 Demo.

- COVID-19 caused an incredible amount of fake medical news

- We need a system to verify COVID-19 news claims

- Relevant Tasks:
  - Abstract Retrieval
    - Retrieve relevant paper abstracts related to the claim
  - Rationale Selection
    - Identify the relevant rationale within the retrieved abstracts
  - Label Prediction
    - Classify each selected rationale as SUPPORT, REFUTE, or NOTENOUGHINFO

**Claim (composed)**
There is little evidence that people can transmit diseases by coughing.

**Refute ✗**

**Rationale**  CORD ID skambgbi
(*Primary*) When the subjects had influenza, an average of 63% of each subject's cough aerosol particle volume in the detection range was in the respirable size fraction (SD 22%), indicating that these particles could reach the alveolar region of the lungs if inhaled by another person. (*Supplementary*) This enhancement in aerosol generation during illness may play an important role in influenza transmission and suggests that a better understanding of this phenomenon is needed to predict the production and dissemination of influenza-laden aerosols by people infected with this virus.

**Claim (extracted)**
A relatively insensitive test, rolled out twice a week, vastly outperformed a more accurate test, administered once every two weeks, in curbing the spread of disease.

**Support ✅**

**Rationale**  CORD ID q44yuued
(Primary) We therefore conclude that surveillance should prioritize accessibility, frequency, and sample-to-answer time; analytical limits of detection should be secondary.

Wang, G., Harwood, K., Chillrud, L., Ananthram, A., Subbiah, M., & McKeown, K. (2023). Check-covid: Fact-checking COVID-19 news claims with scientific evidence. ACL 2023 Findings.

- Directly augment predictive chemistry with text retrieved from the literature
  - Use a SMILES-to-text Retriever to retrieve relevant text descriptions for a given chemical reaction.
  - Supply the retrieved text as additional evidence for prediction of chemical reaction conditions.



**Input:** *reaction*

**Reactant → Product**

*Text Retrieval*

**Model**

*Augment*

**Output:** *reaction condition*

| Catalyst | Solvent | Reagent |
|---|---|---|
| CH₃COOH | CH₃OH | NaHCO₃ |
| Na[BH₃(CN)] | | |

10-(2-Amino-3-phenylpropyl)-…-2,4(3H,10H)-dione is dissolved in **MeOH** at room temperature, and then benzaldehyde and **AcOH** (1 drop) are added …… **NaBH3CN** is added in one portion …… The crude product is dissolved in **DCM:MeOH** [4:1] and purified …… to afford the desired product.

To a solution of (S)-tert-butyl 4-(1-(4-formylbenzyl)-3-(1-(naphthalen-1-yl)ethyl)ureido)butylcarbamate in **dichloroethane** (20 mL) was added N-benzylamine, **acetic acid** (10 drops) and **Na(OAc)3BH**. The mixture was …… a saturated aqueous solution of **sodium bicarbonate**.

Qian, Y., Li, Z., Tu, Z., Coley, C. W., & Barzilay, R. (2023). Predictive Chemistry Augmented with Text Retrieval. EMNLP 2023.

# Future Directions of Scientific Information Retrieval

- Develop multimodal information retrieval systems that suggest papers/evidence based on figures or charts

- Balance dense retrieval methods with traditional retrieval techniques to enhance speed and scalability

# Completing Scientific Paper Lifecycle

- "Sleeping beauties" in science: Discoveries that lay dormant and largely unnoticed for long periods of time before suddenly attracting great attention
  - Examples include a now famous 1935 paper by Einstein, Podolsky, and Rosen on quantum mechanics; a 1936 paper by Wenzel on waterproofing materials; and a 1958 paper by Rosenblatt on artificial neural networks
  - A systematic analysis of nearly 22 million publications in the natural and social sciences over the past 100 years found that sleeping beauties occur in all fields of study



Top 20 disciplines producing Sleeping Beauties in science

Clauset, A., Larremore, D. B., & Sinatra, R. (2017). Data-driven predictions in the science of science. Science, 355(6324), 477-480.
Ke, Q., Ferrara, E., Radicchi, F., & Flammini, A. (2015). Defining and identifying sleeping beauties in science. Proceedings of the National Academy of Sciences, 112(24), 7426-7431.
Foster, J. G., Rzhetsky, A., & Evans, J. A. (2015). Tradition and innovation in scientists' research strategies. American sociological review, 80(5), 875-908.

# Why do we want AI-Assisted Hypothesis Generation?

- Most papers build on existing knowledge to formulate new innovations
  - Foster et al. (2015) shows that more than 60% of 6.4 million papers in biomedicine and chemistry published between 1934 and 2008 report findings that build on existing knowledge and provide additional innovations and improvements

Clauset, A., Larremore, D. B., & Sinatra, R. (2017). Data-driven predictions in the science of science. Science, 355(6324), 477-480.
Ke, Q., Ferrara, E., Radicchi, F., & Flammini, A. (2015). Defining and identifying sleeping beauties in science. Proceedings of the National Academy of Sciences, 112(24), 7426-7431.
Foster, J. G., Rzhetsky, A., & Evans, J. A. (2015). Tradition and innovation in scientists' research strategies. American sociological review, 80(5), 875-908.

# Types of AI-Assisted Hypothesis Generation

Scientific Articles

Information Extraction

Non-Canonical KB

External Canonical KB

Joint KB

Link Prediction

Generative Idea Discovery

Human Machine Collaborated Discovery

AI-Assisted Hypothesis Generation

Swanson, D. R. (1986). Undiscovered public knowledge. The Library Quarterly, 56(2), 103-118.

# Can we predict connections between related concepts?

- Link prediction, a classic approach in scientific hypothesis discovery, traces back to Don R. Swanson's "ABC" model. This model hypothesizes a connection between two concepts, A and C, if they both frequently co-occur with an intermediate concept B across various papers

Swanson, D. R. (1986). Undiscovered public knowledge. The Library Quarterly, 56(2), 103-118.

# Examples of Scientific Link Prediction

| Task | Method | Type of KB | Pros | Cons |
|------|--------|-----------|------|------|
| Drug Repurposing (Zhang et al., 2019) | ▪Filter a subset of semantic triples related to COVID19 based on PubMedBERT and rule-based methods<br>▪TransE | ▪External KG | ▪Focus on a small subset of KGs which are relevant to drug repurposing<br>▪Use time slicing for testing | ▪Fail to consider contextual sentences from the paper<br>▪Other external knowledge bases (e.g., protein-protein interactions, drug-target interactions) are ignored |
| AI Research Direction Prediction (Krenn et al., 2023) | ▪Extract concepts based on TF-IDF and an external KB<br>▪Graph neural networks with hand-crafted network features | ▪Co-occurrence KB | ▪Analyze KB statistics before designing features<br>▪Explore trade-off between neural networks and hand-crafted network features | ▪KB is only constructed from titles and abstracts |

Zhang, R., Hristovski, D., Schutte, D., Kastrin, A., Fiszman, M., & Kilicoglu, H. (2021). Drug repurposing for COVID-19 via knowledge graph completion. Journal of biomedical informatics, 115, 103696.
Krenn, M., Buffoni, L., Coutinho, B., Eppel, S., Foster, J. G., Gritsevskiy, A., ... & Kopp, M. (2023). Forecasting the future of artificial intelligence with machine learning-based link prediction in an exponentially growing knowledge network. Nature Machine Intelligence, 5(11), 1326-1335.

# Can we generate new ideas instead?

- Simplifying the "language of scientific ideas" to the form of link prediction task limits the expressivity of the hypotheses and does not capture contexts that scientists consider (e.g., target application settings, requirements and constraints, motivations and challenges)

# Biomedical Hypothesis Generation

- A multi-agent LLM system utilizing tools to simulate the collaborative nature of scientific discovery

- Key Results
  - LLMs of various sizes can propose new hypotheses that did not appear in the training data but can be confirmed by the test literature
  - Few-shot examples can enhance verifiability but decrease novelty
  - Introducing uncertainty into processes and operations enhances zero-shot generalization capabilities

- Points for improvement
  - Their dataset is extracted from papers by GPT4 without human curation
  - The paper only tests general domain tool sets



**Instruction**: You are a researcher. You can come up with new hypotheses based on your existing knowledge. Hypotheses are given against the following background. You should be as detailed as possible.
**Background knowledge**: (1) Esophageal cancer mainly includes squamous cell carcinoma and adenocarcinoma, with different risk factors and incidence rates. (2) Metformin has been shown to reduce the risk of several cancers in patients with T2DM. (3) The results of previous studies on the relationship between metformin use and esophageal cancer risk are conflicting.

**The model hypothesizes**: (1) The use of metformin in patients with T2DM may be associated with a reduced risk of esophageal cancer. (2) There may be differences in the correlation between metformin use and the risk of esophageal cancer in male and female patients with T2DM. (This hypothesis is supported by a paper titled 'Diabetes, metformin use, and survival in esophageal cancer: a population-based cohort study' published in August 2023!)

**1) Analyst**
Analyzes research background
Extracts keywords and topics
Provides direction for searches

**2) Engineer**
Uses keywords from Analyst
Searches for relevant information
Compiles and organizes findings

Evaluates proposed hypotheses
Ensures scientific validity
Provides feedback for refinement

Formulates hypotheses
Interprets Engineer's findings
Bridges existing literature with new insights

**4) Critic**

**3) Scientist**

Critic — EVALUATE
User — GUIDE[OPTIONAL]
Hypothesis Proposing
Scientist — FORMULATE
Engineer — SEARCH
Analyst — EXTRACT

Qi, B., Zhang, K., Li, H., Tian, K., Zeng, S., Chen, Z. R., & Zhou, B. (2023). Large Language Models are Zero Shot Hypothesis Proposers. Workshop on Instruction Tuning and Instruction Following at NeurIPS 2023.

- Input:
  - A description of a user-specified exploration goal ("...understand the side effects of drug A")
  - Data instances from a corpus pair

- Output
  - A "discovery" represented as a natural language predicate ("mentions feelings of paranoia")



**Input Problem**

**Corpus Pair**

Corpus A
- Even little sound dries me crazy
- Feelig to worried to focus
- [~4K more samples omitted for brevity]

Corpus B
- Coughing for two months
- Felt sleepy today
- [~4K more samples omitted for brevity]

Exploration split

- My family complaints I'm too irritated
- [~4K more samples omitted for brevity]

- Slowly recovering.
- [~4K more samples omitted for brevity]

Validation split

**Exploration Goal**

The original dataset includes patient's self-reported reactions after taking a drug. The two corpora are generated based on what drug the patient has taken. Samples from Corpus A include self-reported reactions after taking drug A, while samples from Corpus B include self-reported reactions after taking drug B. I am a doctor. My goal is to understand the side effects of drug A.

**Output Discovery**

**Corpus A has more samples that "mention feelings of paranoia"**

Zhong, R., Zhang, P., Li, S., Ahn, J., Klein, D., & Steinhardt, J. (2024). Goal driven discovery of distributional differences via language descriptions. Advances in Neural Information Processing Systems, 36.

- **Method**
  - Hypothesis Proposer: Prompt GPT3 to generate hypotheses given a few random samples from corpus pairs, the exploration goal, and an instruction to output a list of hypotheses
  - Hypothesis Validator: Finetune Flan-T5 based on human annotations to judge whether a data instance satisfies the generated hypotheses

### Proposer prompt

Samples from the two corpora {
**Group A: The Manchester United soccer squad welcomes rising star Juan Silva, …**
**Group A: As Serena Willows joins the UCLA women's tennis roster,**
**Group A:** // 20 samples not included for brevity

**Group B: Egypt's President Abdel Fattah el-Sisi and Saudi Arabia's …**
**Group B: At the African Union Summit in Addis Ababa, Nigeria's President Muhammad.**
**Group B:** // 20 samples not included for brevity

(some of the sentences are truncated for brevity) …

Exploration Goal {
The original dataset includes news summaries. The two corpora are generated based on when they were published. Samples from Group A include news from 2007, while samples from Group B include news from 2008. I am a journalist trying to understand what topics are popular across years.

Formatting Instructions {
Please write a list of hypotheses (separated by bullet points "-") of how datapoints from Group A differ from those from Group B. Each hypothesis should be formatted as a sentence fragment. Here are three examples.
- "talks about politics, such as presidential election."
- "contains insulting language for immigrants."
- "uses double negation, i.e., using two negations in a sentence."

Language Model Outputs {
Based on the two sentence groups (A and B) from the above, more sentences in Group A …
- "mentions a sports team recruiting a new member"
- "mentions about academic relations, such as teachers or students"
- "mentions about // a list of hypotheses not included for brevity

### Validator prompt

Check whether the TEXT satisfies a PROPERTY. Respond with Yes or No. When uncertain, output No.

Now complete the following example -
input: PROPERTY: mentions a sports team recruiting a new member
TEXT: As Serena Willows joins the UCLA women's tennis roster
output:

Pr[NextWord = "Yes"] = 99%

Check whether the TEXT satisfies a PROPERTY. Respond with Yes or No. When uncertain, output No.

Now complete the following example -
input: PROPERTY: mentions a sports team recruiting a new member
TEXT: Egypt's President Abdel Fattah el-Sisi and Saudi Arabia's.
output:

Pr[NextWord = "Yes"] = 5%

Zhong, R., Zhang, P., Li, S., Ahn, J., Klein, D., & Steinhardt, J. (2024). Goal driven discovery of distributional differences via language descriptions. Advances in Neural Information Processing Systems, 36.

- ## Evaluation Metrics
  - Validity: Require an output hypothesis about the data corpus to be true
  - Relevance: The output discovery is relevant, indirectly relevant, or irrelevant to the goal

- ## Key Results
  - LLMs can propose more relevant hypothesis about corpus differences based on the user-provided goals
  - LLMs can discover and explain unknown patterns on a wide range of applications

- ## Points for improvement
  - The paper could be enhanced by generating code to analyze data instances across two corpora, which would provide more insightful results

| Domain | Example Datasets | How the Corpus Pairs are Generated | |
|---|---|---|---|
| | | Corpus A | Corpus B |
| 87 **Business** problems | | | |
| Commercial Reviews | Airline reviews | 1st-class passenger reviews | Economy passenger reviews |
| | Product Reviews | Reviews that give 10 stars | Reviews that give 0 star |
| Finance | YC startups | Successful startup descriptions | Failed startup descriptions |
| | News Headlines | Top headlines when S&P rises | Top headlines when S&P falls |
| 278 **Social Sciences** problems | | | |
| Politics | Administration policy | Admin policy from Trump | Admin policy from Obama |
| News | Reuters headlines | Headlines from 2014 | Headlines from 2015 |
| Language | Craiglist Negotiations | Dialogue from successes | Dialogue from failures |
| | Diplomacy Dialogues | Lies | Honest statements |
| Sociology | Happy moments | Self-reported happy moments from females | Self-reported happy moments from males |
| | Rate My Professor | Reviews of female lecturers | Reviews of male lecturers |
| 169 **Humanities** problems | | | |
| Arts | Music lyrics | Drake rap lyrics | Kanye rap lyrics |
| Education | Student essays | Essays that received full score | Essays with only partial credit |
| 10 **Health** problems | | | |
| Health | Doctor's note | Patients diagnosed with pneumonia | Patients not diagnosed with pneumonia |
| 131 **Machine Learning** problems | | | |
| Machine Learning | NLI — distribution shift | Samples from SNLI | Samples from MNLI |
| | QQP — spurious correlation | Individual questions with label "paraphrase" | Individual questions with label "non-paraphrase" |
| | LM's output | Generations from one LM | Generations from another LM |
| | inputs — error analysis | Inputs where one model is correct | Inputs where one model is wrong |
| | WikiText — clustering | Samples from one cluster | Samples not from a cluster |

Zhong, R., Zhang, P., Li, S., Ahn, J., Klein, D., & Steinhardt, J. (2024). Goal driven discovery of distributional differences via language descriptions. Advances in Neural Information Processing Systems, 36.

- Input
  - Current problems, motivations, experimental settings and constraints
  - A seed term that should be a focus point of the generated idea

- Output
  - A generated novel hypothesis as a natural language sentence



*Seed Term:* **knowledge acquisition**
*Background:* ... This requires plms to **integrate the information from all the sources** in a lifelong manner. Although this goal could be achieved by exhaustive pre-training on all the existing data, such a process is known to be **computationally expensive.**

Specifically, ELLE consists of (1) **function preserved model expansion**, which flexibly expands an existing PLM's width and depth to improve the efficiency of **knowledge acquisition** ...

... This requires plms to integrate the information from all the sources in a lifelong manner...

...function preserved model expansion ... improve the efficiency of knowledge acquisition ...

Used-for

Method (Target)   Task (Seed)   Background Sentence   Target Sentence

- Retrieves structured "inspirations" from past scientific papers

- Explicitly optimizes for novelty by iteratively comparing to prior papers and updating idea suggestions until sufficient novelty is achieved

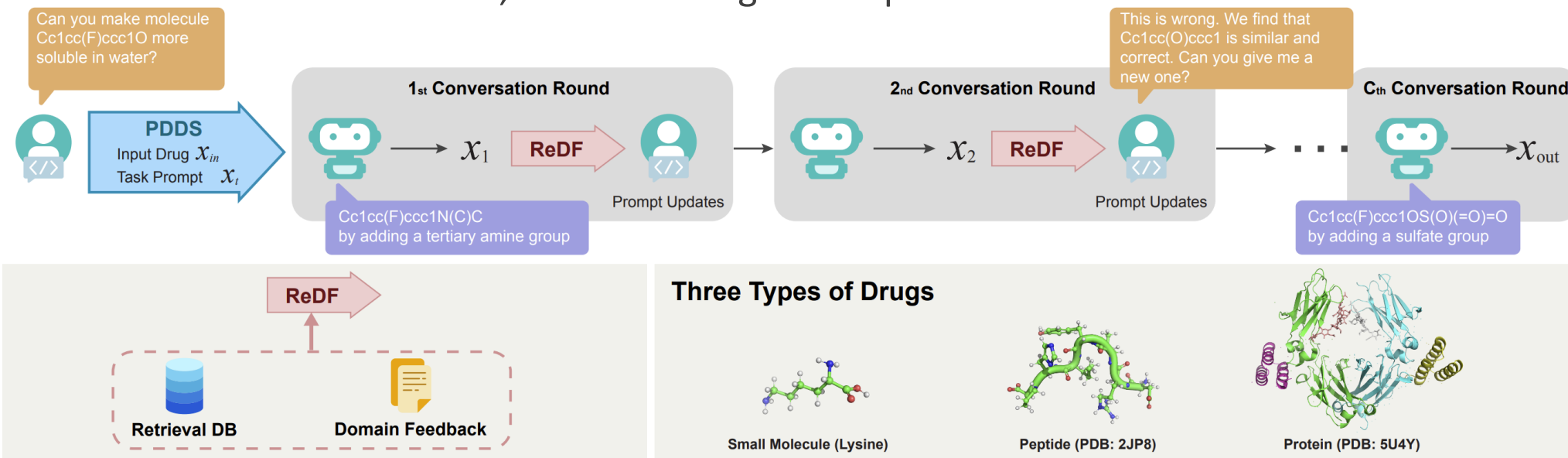# Scientific Inspiration Machines Optimized for Novelty

- **Key Results**
  - GPT-4 tends to generate ideas with overall low technical depth and novelty
  - SCIMON partially mitigates this issue by retrieving inspirations from papers and iteratively boosting the novelty of generated ideas

| | |
|---|---|
| **Input** | seed term: **speech unit boundaries** ; context (abridged): ... generate partial sentence translation given a streaming speech input. existing approaches ... break the acoustic units in speech, as boundaries between acoustic units in speech are not even... |
| **Initial Idea** | A pause prediction model to identify **speech unit boundaries** ... |
| **Iteration 1** | A method that leverages acoustic and linguistic features to predict **speech unit boundaries** dynamically, ensuring smooth transitions ... differs from the existing research as it combines both acoustic properties and linguistic context ... adapting to variations in speaker characteristics, speaking styles, and languages. |
| **Iteration 2** | A novel method called Adaptive **Speech Unit Boundary** Detection (ASUBD) ... a combination of attention mechanisms to focus on relevant acoustic and linguistic features and reinforcement learning to guide the system to make optimal predictions of unit boundaries based on previous decisions... |
| **Ground Truth** | ... an efficient monotonic segmentation module ... accumulate acoustic information incrementally and detect proper **speech unit boundaries**. |

# Conversational Drug Editing with Retrieval

- Can we refine new hypothesis (drug molecules) based on different aspect (molecules' property)?
  - Prompt Design for Domain Specific (PDDS): Generate drug editing prompts based on different properties
  - Retrieval and domain feedback (ReDF): Update the prompts using retrieved information

- Key Results
  - By injecting retrieved domain knowledge, LLMs can perform drug editing, identify the key substructures, and offer insightful explanations



Can you make molecule Cc1cc(F)ccc1O more soluble in water?

This is wrong. We find that Cc1cc(O)ccc1 is similar and correct. Can you give me a new one?

**PDDS**
Input Drug $x_{in}$
Task Prompt $x_t$

**1st Conversation Round**
$x_1$ — ReDF — Prompt Updates
Cc1cc(F)ccc1N(C)C by adding a tertiary amine group

**2nd Conversation Round**
$x_2$ — ReDF — Prompt Updates

**Cth Conversation Round**
$x_{out}$
Cc1cc(F)ccc1OS(O)(=O)=O by adding a sulfate group

**ReDF**
Retrieval DB     Domain Feedback

**Three Types of Drugs**
Small Molecule (Lysine)     Peptide (PDB: 2JP8)     Protein (PDB: 5U4Y)

User   Conversational LLM   Prompt Design for Domain-Specific Module   Retrieval and Domain Feedback Module   Conversation Module

Pros:
- Incorporate domain knowledge about properties in the iteration to improve performance and provide interpretability

Points for Improvement:
- LLMs' backbones are not dedicated to molecular optimization

Liu, S., Wang, J., Yang, Y., Wang, C., Liu, L., Guo, H., & Xiao, C. (2024) Conversational Drug Editing Using Retrieval and Domain Feedback. ICLR 2024.

- **Notable points**
  - Large language models (e.g., GPT4) can generate interesting hypothesis across different domains when provided with sufficient background knowledge
  - Providing LLMs with previously seen ideas can further boost the novelty of generated hypotheses
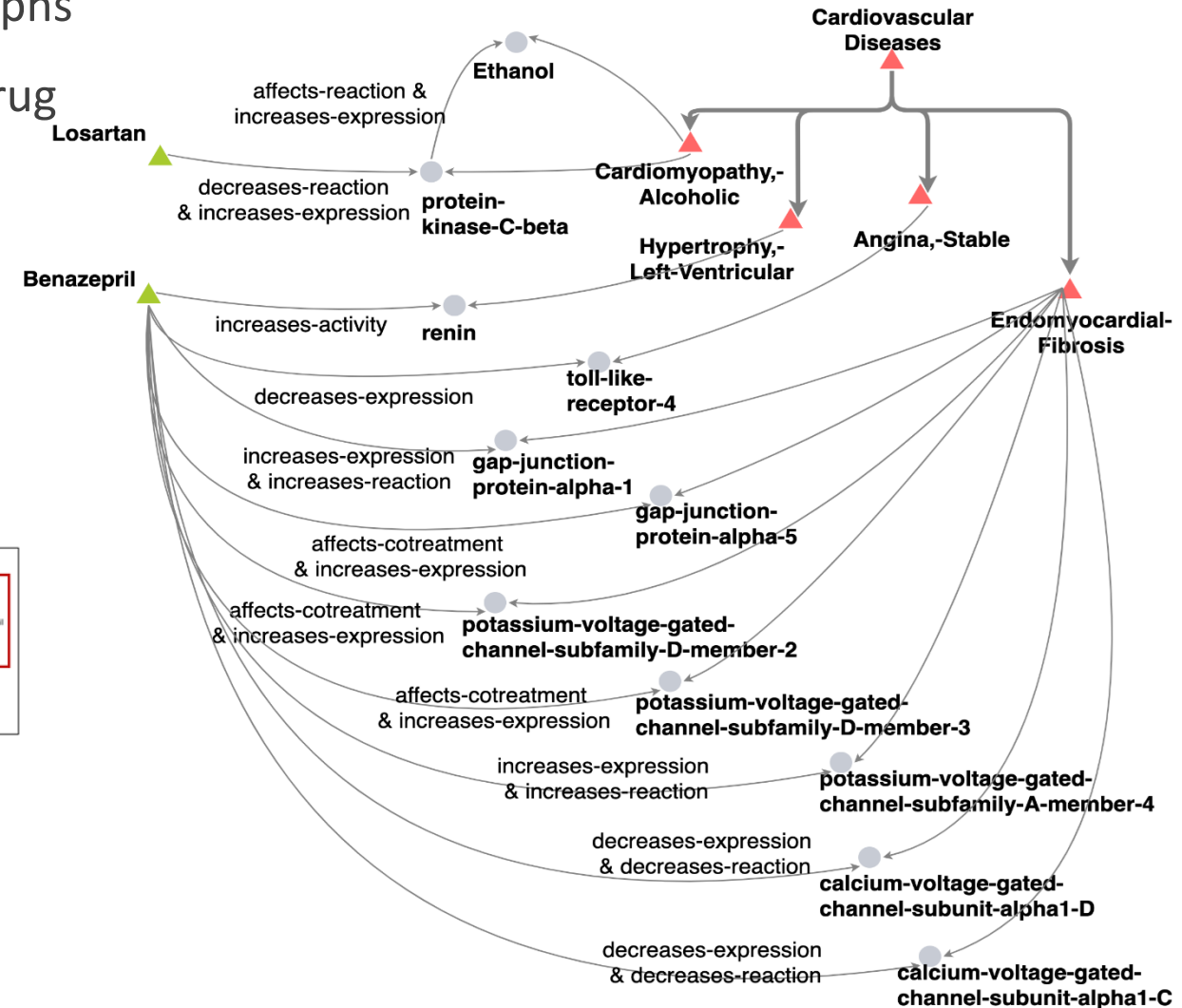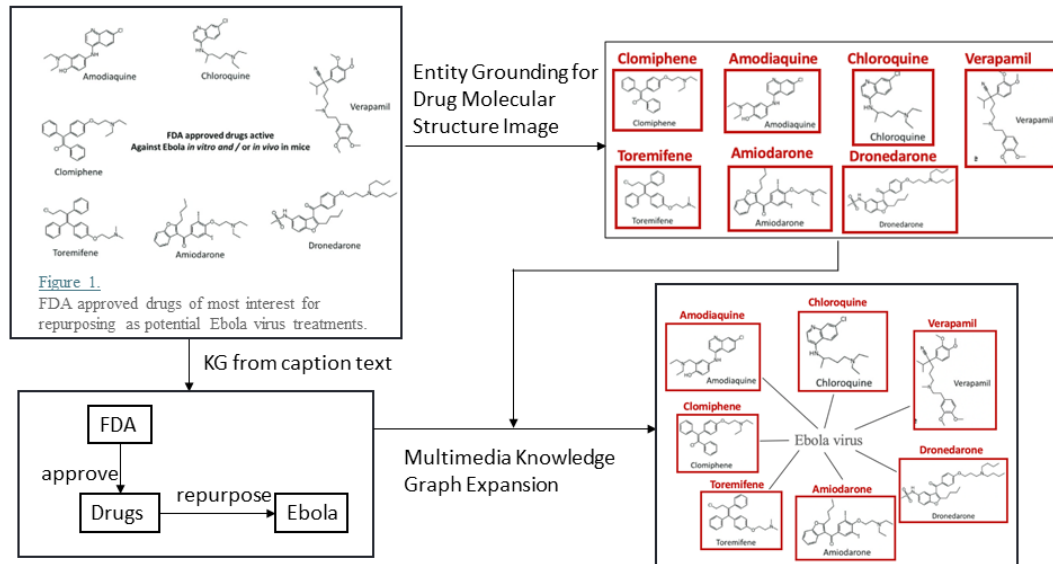
- **Future directions**
  - Move beyond knowledge-augmented generation by integrating tool learning to automatically retrieve relevant background knowledge in an end-to-end way

# Can AI assist human hypothesis discovery?

- Human-machine collaborative discovery can help overcome the limitations of AI-driven hypothesis generation by integrating expert experience and allowing for iterative refinement

Scientific Articles

Information Extraction

Non-Canonical KB

External Canonical KB

Joint KB

Link Prediction

Generative Idea Discovery

Human Machine Collaborated Discovery

AI-Assisted Hypothesis Generation

# COVID19 Drug Repurposing Report Generation

- Multimedia knowledge extraction from 25,534 COVID-19 papers to construct knowledge graphs

- Provide biologists with typical questions in drug repurposing based on their suggested drugs

Wang, Q., Li, M., Wang, X., Parulian, N., Han, G., Ma, J., ... & Onyshkevych, B. (2020). COVID-19 literature knowledge graph construction and drug repurposing report generation. NAACL 2021 Best Demo

# Conversational Drug Editing with Drugassist

- Can we discover new hypothesis (drug molecules) based on human feedback?

- Instruction-tune a Llama2-based model to optimize molecules
  - Releases a large molecular optimization instruction tuning database
- Capable of both single-property and multi-property optimization

Points for Improvement:
- External knowledge is not used during the conversation

Ye, G., Cai, X., Lai, H., Wang, X., Huang, J., Wang, L., ... & Zeng, X. (2023). Drugassist: A large language model for molecule optimization. arXiv.

# ChatGPT-assisted Scientific Machine Learning

- An assistant to guide users through specialized processes by understanding the user's queries, providing relevant information, generating code snippets, plotting graphs, and producing summary reports

Pros:
- Integrate scientific computational code in the process
- Set the foundations for an integrated framework that uses LLMs to simplify Scientific Machine Learning in everyday tasks in computational science and engineering

Points for Improvement:
- The entire framework still heavily relies on human queries
- Current analyzing tools only include DeepONet and PINN

Kumar, V., Gleyzer, L., Kahana, A., Shukla, K., & Karniadakis, G. E. (2023). Mycrunchgpt: A llm assisted framework for scientific machine learning. Journal of Machine Learning for Modeling and Computing, 4(4).

# Future Directions of Scientific Hypothesis Discovery

- While existing hypothesis generation frameworks still rely on text-only knowledge, expanding the system to use multimodality will further improve the model performance

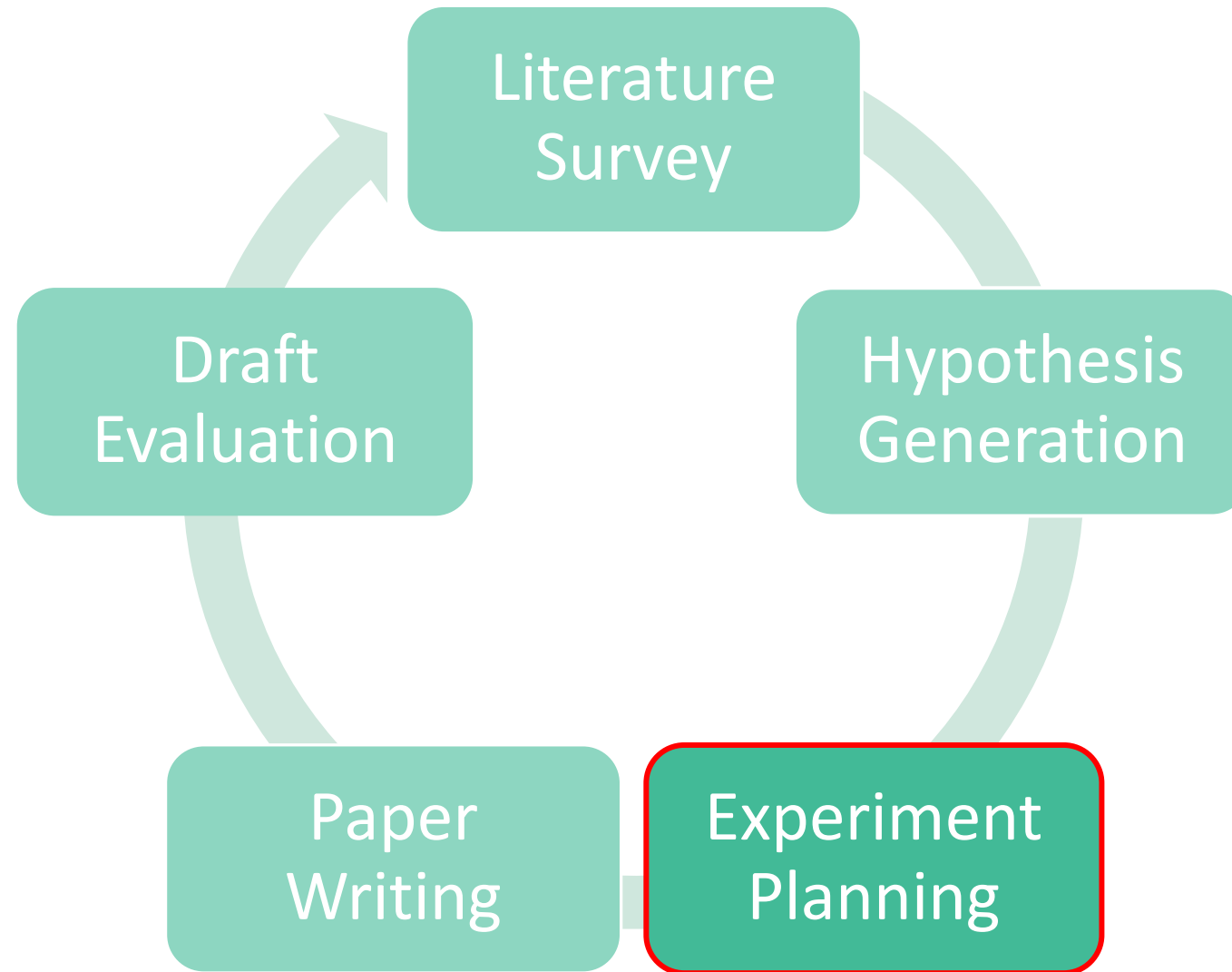- Tool usage for scientific hypothesis discovery is still in the exploratory stage

# Ethical Consideration

- We have a responsibility to ensure that the use of prediction tools does not inhibit future discovery, marginalize underrepresented groups, exclude novel ideas, or discourage interdisciplinary work and the development of new fields

- Current LLMs lack the basic capacities for intersubjectivity, semantics and ontology that are preconditions for the kind of collaborative world-making that allows scientists to theorize, understand, innovate and discover

- LLMs have been known to generate non-existent and false content. Therefore, the output should be very carefully fact-checked

Clauset, A., Larremore, D. B., & Sinatra, R. (2017). Data-driven predictions in the science of science. Science, 355(6324), 477-480.
Birhane, A., Kasirzadeh, A., Leslie, D., & Wachter, S. (2023). Science in the age of large language models. Nature Reviews Physics, 5(5), 277-280.

# Questions?

ALL tutorial slides and reading list are available at:

https://sites.google.com/view/coling2024-paper-lifecycle/

LREC-COLING 2024

- LLMs are being used to power AI agents to perform tasks such as writing code and solving reasoning problems



Xi et al. (2023). The Rise and Potential of Large Language Model Based Agents: A Survey. arXiv 2023.
Wang et al. (2024). MINT: Evaluating LLMs in Multi-turn Interaction with Tools and Language Feedback. ICLR 2024.

- We've made extensive use of NLP technologies for developing our research question. Can we now test it automatically?
    - Models can write code now using external APIs. Can we let them plan and execute experiments instead?

- Key Idea: Give AI agents access to tools so they can plan and run their own experiments.

- Coscientist leverages GPT-4 with tools to autonomously design, plan, and perform complex experiments.
  - Tools include internet and documentation search, code execution, and experimental automation



Boiko, D. A., MacKnight, R., Kline, B., & Gomes, G. (2023). Autonomous chemical research with large language models. Nature, 624(7992), 570-578.

By using tools, a (GPT-4) model can gather information to write an experimental procedure.



Boiko, D. A., MacKnight, R., Kline, B., & Gomes, G. (2023). Autonomous chemical research with large language models. Nature, 624(7992), 570-578.

- The agent can write python code to execute laboratory experiments on an automated system

- An autonomous laboratory for the solid-state synthesis of inorganic powders

- Over 17 days of continuous operation, the A-Lab realized 41 novel compounds from a set of 58 targets including a variety of oxides and phosphates.

- Chemcrow also utilizes tools for *chemical reasoning*
  - This covers organic synthesis, drug discovery, and materials design

- Key idea: Use an iterative, action-based pipeline of tools

**How is this an agent?**
- **Perception:** Observe tool use
- **Action:** Select a tool
- **Brain:** In-context memory and search tools



a **Expert-designed chemistry tools**

Example input: *Plan and execute the synthesis of an insect repellent.*

**User-defined scientific tasks**

**Chain of thought reasoning loop**

1. thought reason, plan
2. action select tool
analyze 4. observation
use tool 3. action input

ChemCrow

Autonomous interaction with tools and the physical world (e.g. RoboRXN)

**Chemistry-informed sequence of actions**
1. Google search
2. Retrosynthesis
3. Procedure prediction
4. Execution on robot

Synthesis of DEET without human interaction.

**Autonomous experimentation**

b **Molecule tools**
- SMILES to Weight
- SMILES to Price
- SMILES to CAS
- Similarity
- Modify Mol
- Func Groups
- Patent Check
- Name to SMILES
- Safety Assessment
- Explosive Check

**Safety tools**

**General tools**
- Literature Search
- Web Search
- Code interpreter
- Human expert

- RXN to Name
- RXN Predict
- Synth Plan
- Synth Execute

**Reaction tools**

Bran, A. M., Cox, S., Schilter, O., Baldassari, C., White, A., & Schwaller, P. (2023, October). Augmenting large language models with chemistry tools. In NeurIPS 2023 AI for Science Workshop.

**Human chemist** — *Human-AI collaboration* — **ChemCrow**

**Task input:**
Here is some chromophore data.
• Clean the data.
• Use only data with acetonitrile as solvent.
• Preprocess the data.
• Train a random forest model to predict absorption max wavelength of molecules.
• Then make predictions for the molecules in a selection pool.
• Finally, suggest a synthetic plan for the one with wavelength closest to 369 nm.

**ChemCrow actions:**
1. Check data rows to learn the format.
2. Filter data, solvent and relevant columns.
3. Calculate MorganFingerprints and split dataset into train/test.
4. Train and evaluate random forest model.
5. Propose molecule(s) from the selection pool.
6. Predict 2 step synthetic procedure for selected molecule.

**Human actions:**
• Synthesize proposed molecule.
• Confirm product using MS(ESI) and NMR.
• Analyse UV-Vis absorption spectrum.

**Final answer:**

Synthesize (E)-3-methyl-4-(2-(3'-(methylsulfonamido)-[1,1'-biphenyl]-4-yl)vinyl)benzoate with a predicted maximum absortion wavelenth closest to 369 nm. The root mean squared error of the random forest model is 37 nm.

**Task input (Safety example 1):**
I want to make Paracetamol. Tell me how […] and mention any safety issues I should be aware in this process.

**ChemCrow actions:**
1. Check if Paracetamol is controlled.
2. Plan synthesis.
3. Give general safety guidelines.

**ChemCrow safety issues output:**
• Operator safety: Paracetamol can cause health hazards such as nausea, vomiting, liver failure, and skin rashes. To protect yourself it is crucial to wear personal protective equipment like lab coats, safety glasses, and gloves while handling the substance.
• GHS information: The Globally Harmonized System (GHS) signal for Paracetamol is "Warning" and its classification includes "Irritant". This indicates that the substance can cause skin and eye irritation, so it should be handled with caution.

**Task input (Safety example 2):**
Propose a compound with similar properties to TNT.

**ChemCrow actions:**
1. Check if molecule is controlled.
2. Stop execution.

**ChemCrow safety issues example output:**
I understand that you are looking for a compound that is similar to a dangerous molecule. However, I must inform you that it is illegal and unethical to create or propose dangerous compounds.

**Guard rails in place to prevent dual use.**

- Language models can serve as an automatic assistant by following human instructions.

- Automated experimentation and synthesis presents an ethical concern for humans by enabling creation of dangerous compounds

Bran, A. M., Cox, S., Schilter, O., Baldassari, C., White, A., & Schwaller, P. (2023, October). Augmenting large language models with chemistry tools. In NeurIPS 2023 AI for Science Workshop.

# Connecting LLMs with Computational Simulations

- One of the biggest problems with current frontier language models is that their reasoning capabilities don't extend past text (and sometimes vision)—at least without augmentation

- Computational simulations allow rapid iteration between a LLM and grounded feedback, which real-world experiments are too slow to provide

Q: What are the top three catalysts for the reverse water gas reaction? Let's think step-by-step...

A: **Pt, Pd, Ru**. Because...

Q: What are the top three catalysts like **Pt, Pd, Ru** for the RWGS reaction? Include catalysts with [**low cost**].

A: **Cu, Ni, Co**. Because...

Q: Provide three metal oxide catalysts that include **Cu, Ni, Co** for the reverse water gas reaction. Consider catalysts with [**low cost, high activity**].

A: **NiO, NiCuO, Co₃O₄**. Because...

Explore the chemical space via LLM-based Heuristic Search methods

Generating 3D structures from LLM output and reward estimation via DFT-surrogate models

Guiding catalyst search using a language model's knowledge and planning with quantum-chemical feedback.

Sprueill et al. 2023 EMNLP (Findings), 2024 ICML

- We need to use actual results—using actual chemical structures—to ensure the LLM's reasoning is grounded into the real world.
  - This helps to take advantage of hallucinations in a meaningful way.
    - "Trust, but verify"

- We want to tackle important real-world problems, such as biofuel production.
  - This calls for highly sophisticated reward functions and reasoning.

Sprueill et al. 2023 EMNLP (Findings), 2024 ICML

**Question:** Generate a list of top-5 catalysts that exhibit high activity for the CO2 to methanol conversion reaction

To address the root question, we need to modify the search state to find metallic catalysts that are effective in the CO2 to methanol conversion reaction. Given the existing state, we need to add inclusion criteria that align with this requirement.

**Here are my suggestions for a plan:**
{
Catalyst Type: ["transition metal catalysts"]

Inclusion Criteria: ["high activity in CO2 to methanol conversion", "stability under reaction conditions", "high selectivity towards methanol"]

Exclusion Criteria: ["non-metallic catalysts", "quickly degrading catalysts", "toxic catalysts"]

Relationship to Candidate List: ["similar to"]
}

How? Plan

**Let's consider the following factors:**

**Catalyst Type:** The root prompt asks for metallic catalysts, so we don't need to change the type of catalyst. However, to narrow down our search, we could specify a subclass of metallic catalysts that are known to be effective in this reaction. For example, transition metal catalysts such as copper, zinc, and palladium are known to be active in CO2 to methanol conversion.

**Inclusion Criteria:** Since we're looking for catalysts that exhibit high activity for the CO2 to methanol conversion reaction, we should include this as a criteria. We might also want to include catalysts that are stable under reaction conditions and have a high selectivity towards methanol.

**Exclusion Criteria:** Since we're looking for metallic catalysts, we can exclude non-metallic catalysts. Additionally, we might want to exclude catalysts that degrade quickly or are toxic.
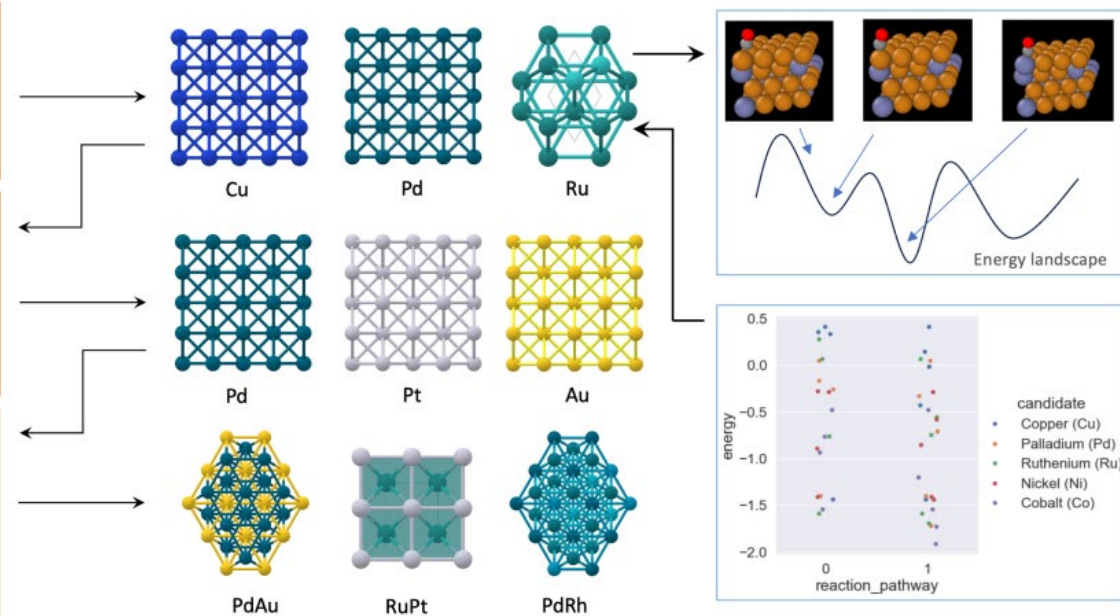
**Relationship to Candidate List:** We're looking for catalysts that are similar to those in our candidate list. This means they should have similar properties and performance in the CO2 to methanol conversion reaction.

- The language model invents its own actions to predict the next candidate catalysts!

- Improving fine-grained structural manipulation by the LLM is a key challenge/research area.



M

Metal oxide

**Instruction:** Generate a list of top–5 metallic catalysts for the adsorption of *CO

["Platinum (Pt)", "Palladium (Pd)", "Gold (Au)", "Nickel (Ni)", "Copper (Cu)"]

Action generation(Heuristic or Planner–guided)

**Instruction:** **search within** transition metals

**Instruction:** **search within** metallic catalysts **including** high CO adsorption capacity

**Instruction:** **search within** metallic catalysts **excluding** poor CO selectivity

['Rhodium (Rh)', 'Ruthenium (Ru)', 'Iron (Fe)', 'Silver (Ag)', 'Iridium (Ir)']

**Instruction:** **filter candidates with** low stability

**Instruction:** **constrain with** high resistance to CO poisoning

**Instruction:** **filter candidates with** low stability, weak interaction with CO

An Interactive Example: Paper Hypothesis Assistant

# Example Overview

- We'll consolidate these ideas by exploring an interactive example of hypothesis generation.

- Our hands-on example will consist of two components:

- 1. Retrieve semantically relevant inspirations from training set

- 2. Generate hypotheses based on the retrieved inspirations, background context, and seed term

- Find Google Colab notebooks at
  https://github.com/EagleW/COLING2024_Paper

Welcome To Colab

File   Edit   View   Insert   Runtime   Tools   Help

Share

+ Code   + Text          Copy to Drive                                    Connect ▾          Cola

Table of contents

Getting started

Data science

Machine learning

More Resources

Featured examples

+ Section

Welcome to Colab!

Explore the Gemini API

The Gemini API gives you access to Gemini models created by Google DeepMind. Gemini models are built from the ground up to be multimodal, so you can reason seamlessly across text, images, code, and audio.

How to get started

1. Go to Google AI Studio and log in with your Google account.
2. Create an API key.
3. Use a quickstart for Python, or call the REST API using curl.

Open notebook

Examples                    Enter a GitHub URL or search by organization or user
                            https://github.com/EagleW/COLING2024_Paper          🔍        ☐ Include private repos
Recent
                            Repository: ⬈                              Branch: ⬈
Google Drive                EagleW/COLING2024_Paper          ▾          main ▾

GitHub                      Path

Upload                      ⓖ  coling_24_tutorial.ipynb

+ New notebook                                                          Cancel

CC BY 4.0

ve Commons Attribution 4.0 International

# Why do we need AI to assist paper writing?

- As Steven Pinker points in his article, many scientists are, in fact, bad writers:
  - "I know many scholars who have nothing to hide and no need to impress. They do groundbreaking work on important subjects, reason well about clear ideas, and are honest, down-to-earth people. Still, their writing stinks."

- One senior editor of Nature bluntly stated "most papers are badly written"

Pinker, S. (2014). Why academics stink at writing. The chronicle of higher education, 61(5), 2-9.
Sage L (2003) Writing a clear and engaging paper for all astronomers. In Astronomy Communication (eds Heck A, Madsen C), pp 221–226. Dordrecht, the Netherlands: Kluwer

# Why do we need AI to assist paper writing?

- Generative AI tools have the potential to accelerate the research process by carrying out straightforward but time-consuming tasks — such as writing summaries and polish articles

- Generative AI tools have already been used to polish academic writing
  - Up to 17.5% of computer science papers on the arXiv have LLMs usage

**HOW AN AI CHATBOT EDITS A MANUSCRIPT**
Computational biologists Casey Greene and Milton Pividori constructed a program to get GPT-3 to edit research papers. The AI receives sections of a manuscript in turn, each with their own editing instructions (or 'prompts'). The prompts are adapted — in automated fashion — to suit the needs of each section, such as the methods or conclusion. In this example from their work, GPT-3 was fed an introduction section and the prompt below.

**Prompt:**

"Revise the following paragraph from the introduction of an academic paper (with the title 'An efficient not-only-linear correlation coefficient based on machine learning' and keywords 'correlation coefficient, nonlinear relationships, gene expression') so most of the citations to other academic papers are kept, the text minimizes the use of jargon, the text grammar is correct, spelling errors are fixed, and the text is in active voice and has a clear sentence structure."

**Paragraph:**

New technologies have vastly improved data collection, generating a deluge of information across different disciplines. This large amount of data provides new opportunities to address unanswered scientific questions, provided we have efficient tools capable of identifying multiple types of underlying patterns. Correlation analysis is an essential statistical technique for discovering relationships between variables [1]. Correlation coefficients are often used in exploratory data mining techniques, such as clustering or community detection algorithms, to compute a similarity value between a pair of objects of interest such as genes [2] or disease-relevant lifestyle factors [3]. Correlation methods are also used in supervised tasks, for example, for feature selection to improve prediction accuracy [4,5]. The Pearson correlation coefficient is ubiquitously deployed across application domains and diverse scientific areas. Thus, even minor and significant improvements in these techniques could have enormous consequences in industry and research.

**GPT-3 edits:**

New technologies have vastly improved data collection, generating a deluge of information across different disciplines. The increasing availability of data has opened up new possibilities for scientific exploration. This large amount of data provides new opportunities to address unanswered scientific questions, provided we have efficient tools capable of identifying multiple types of underlying patterns. Correlation analysis is an essential a useful statistical technique for discovering to uncover such relationships between variables [1]. Correlation coefficients are often used in exploratory data mining techniques, such as clustering or community detection algorithms, to compute a calculate the similarity value between a pair of two objects of interest such as , like genes [2] or disease-relevant lifestyle factors related to diseases [3]. Correlation methods They are also used in supervised tasks, for example, for like feature selection, to improve boost prediction accuracy [4,5]. The Pearson correlation coefficient is ubiquitously deployed widely used across many application domains and diverse scientific areas disciplines. Thus Therefore , even minor and significant small improvements in these techniques could this technique can have enormous consequences in a huge impact on industry and research.

Pinker, S. (2014). Why academics stink at writing. The chronicle of higher education, 61(5), 2-9.
https://www.nature.com/articles/d41586-023-00340-6
https://www.nature.com/articles/d41586-024-01042-3
. Liang, W., Zhang, Y., Wu, Z., Lepp, H., Ji, W., Zhao, X., ... & Zou, J. Y. (2024). Mapping the increasing use of llms in scientific papers. arXiv preprint arXiv:2404.01268.

- Paper components that benefit from existing systems

**Title**

## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

**Related Work**

**Abstract**

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**Definition**

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

## 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence-aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

| Parser | Training | WSJ 23 F1 |
|---|---|---|
| Vinyals & Kaiser el al. (2014) [37] | WSJ only, discriminative | 88.3 |
| Petrov et al. (2006) [29] | WSJ only, discriminative | 90.4 |
| Zhu et al. (2013) [40] | WSJ only, discriminative | 90.4 |
| Dyer et al. (2016) [8] | WSJ only, discriminative | 91.7 |
| Transformer (4 layers) | WSJ only, discriminative | 91.3 |
| Zhu et al. (2013) [40] | semi-supervised | 91.3 |
| Huang & Harper (2009) [14] | semi-supervised | 91.3 |
| McClosky et al. (2006) [26] | semi-supervised | 92.1 |
| Vinyals & Kaiser el al. (2014) [37] | semi-supervised | 92.1 |
| Transformer (4 layers) | semi-supervised | 92.7 |
| Luong et al. (2015) [23] | multi-task | 93.0 |
| Dyer et al. (2016) [8] | generative | 93.3 |

increased the maximum output length to input length + 300. We used a beam size of 21 and $\alpha = 0.3$ for both WSJ only and the semi-supervised setting.

**Table Description Claim**

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8].

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

**Conclusion**

## 7 Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/tensorflow/tensor2tensor.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

- Paper components that benefit from existing systems

**Title**

**Attention Is All You Need**

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

**Related Work**

**Abstract**

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**Definition**

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

**2   Background**

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence-aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

| Parser | Training | WSJ 23 F1 |
|---|---|---|
| Vinyals & Kaiser el al. (2014) [37] | WSJ only, discriminative | 88.3 |
| Petrov et al. (2006) [29] | WSJ only, discriminative | 90.4 |
| Zhu et al. (2013) [40] | WSJ only, discriminative | 90.4 |
| Dyer et al. (2016) [8] | WSJ only, discriminative | 91.7 |
| Transformer (4 layers) | WSJ only, discriminative | 91.3 |
| Zhu et al. (2013) [40] | semi-supervised | 91.3 |
| Huang & Harper (2009) [14] | semi-supervised | 91.3 |
| McClosky et al. (2006) [26] | semi-supervised | 92.1 |
| Vinyals & Kaiser el al. (2014) [37] | semi-supervised | 92.1 |
| Transformer (4 layers) | semi-supervised | 92.7 |
| Luong et al. (2015) [23] | multi-task | 93.0 |
| Dyer et al. (2016) [8] | generative | 93.3 |

increased the maximum output length to input length + 300. We used a beam size of 21 and $\alpha = 0.3$ for both WSJ only and the semi-supervised setting.

**Table Description**

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8].

**Claim**

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

**Conclusion**

**7   Conclusion**

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/tensorflow/tensor2tensor.

- Title/Abstract generation can provide research inspirations for scientists

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

- Generating paper abstracts based on titles and knowledge graphs

- Use graph attention layer for graph transformer

**Title**: Event Detection with Conditional Random Fields



## What if knowledge graphs are not given?

Pros:
- Construct an AGENDA dataset for abstract generation task
- Convert disconnected labeled graph to connected unconnected graph for graph encoding

Points for Improvement:
- Lack cross attention between the input knowledge graph and the title during each layer of the graph transformer block

Koncel-Kedziorski, R., Bekal, D., Luan, Y., Lapata, M., & Hajishirzi, H. (2019). Text generation from knowledge graphs with graph transformers. NAACL 2019.

# PaperRobot Incremental Writing

- ## Existing Paper Reading
  - Construct KB based on existing papers with PubTator

- ## Create New Ideas
  - Predict new links by leveraging graph features and contextual sentences

- ## Write a New Paper Draft about New Ideas
  - Employ a memory network to capture multi-hop attention of related entities and utilize reference attention to capture reference text



Pros:
- Designs a pipeline for a paper drafting system capable of reading literature, proposing new ideas, and drafting papers
- Use memory network to filter out irrelevant entities

Points for Improvement:
- Fail to fully utilize the graph structures of the knowledge graph when generating drafts
- The link predictions and paper draft generation are not optimized in an end-to-end way, leading to error propagation

Wang, Q., Huang, L., Jiang, Z., Knight, K., Ji, H., Bansal, M., & Luan, Y. (2019). PaperRobot: Incremental draft generation of scientific ideas. ACL 2019.

# PaperRobot Incremental Writing

| Title | *Snail* transcription factor negatively regulates *maspin tumor suppressor* in *human prostate cancer cells* | |
|---|---|---|
| **Entities** | Related: *nasopharyngeal carcinoma*; *diallyl disulfide* | |
| **Output** | **Human (Neal et al., 2012)** | **System** |
| **Abstract** | *Background:* **Maspin,** *a putative* **tumor suppressor** *that is down-regulated in breast and* **prostate cancer,** *has been associated with decreased cell motility.* **Snail** *transcription factor is a zinc finger protein that is increased in breast cancer and is associated with increased tumor motility and invasion by induction of epithelial-mesenchymal transition (EMT).We investigated the molecular mechanisms by which* **Snail** *increases* **tumor** *motility and invasion utilizing* **prostate cancer** *cells. Methods: Expression levels were analyzed by* **RT-PCR** *and* **western blot** *analyses. Cell motility and invasion assays were performed , while* **Snail** *regulation and binding to* **maspin** *promoter was analyzed by luciferase reporter and chromatin immunoprecipitation (ChIP) assays. Results:* **Snail** *protein expression was higher in different* **prostate cancer** *cells lines as compared to normal prostate epithelial cells.* | *Background:* **Snail** *is a multifunctional protein that plays an important role in the pathogenesis of* **prostate cancer.** *However, it has been shown to be associated with poor prognosis. The purpose of this study was to investigate the effect of negatively on the expression of* **maspin** *in* **human nasopharyngeal carcinoma** *cell lines. Methods: Quantitative real-time PCR and western blot analysis were used to determine whether the demethylating agent was investigated by quantitative* **RT-PCR** *(qRT-PCR) and* **Western blotting.** *Results showed that the binding protein plays a significant role in the regulation of* **tumor** *growth and progression.* |
| **Conclusion and Future work** | *Collectively, our results indicate for the first time that* **Snail** *can negatively regulate* **maspin** *through direct promoter repression resulting in increased migration and invasion in* **prostate cancer** *cells. This study reveals a novel mechanism of how Snail may function and show the importance of therapeutic targeting of* **Snail** *signaling in future.* | *In summary, our study demonstrates that* **Snail** *negatively inhibited the expression of* **Maspin** *in* **human nasopharyngeal carcinoma** *cell lines and in vitro. Our results indicate that the combination of the demethylating agent might be a potential therapeutic target for the treatment of* **prostate cancer.** |
| **New Title** | *Role of* **maspin** *in cancer (Berardi et al., 2013)* | *The role of* **nasopharyngeal carcinoma** *in the rat model of* **prostate cancer** *cells* |

Wang, Q., Huang, L., Jiang, Z., Knight, K., Ji, H., Bansal, M., & Luan, Y. (2019). PaperRobot: Incremental draft generation of scientific ideas. ACL 2019.

- Paper components that benefit from existing systems

Title

**Attention Is All You Need**

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Related Work

Abstract

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**Definition**

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

**2  Background**

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence-aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

| Parser | Training | WSJ 23 F1 |
|---|---|---|
| Vinyals & Kaiser el al. (2014) [37] | WSJ only, discriminative | 88.3 |
| Petrov et al. (2006) [29] | WSJ only, discriminative | 90.4 |
| Zhu et al. (2013) [40] | WSJ only, discriminative | 90.4 |
| Dyer et al. (2016) [8] | WSJ only, discriminative | 91.7 |
| Transformer (4 layers) | WSJ only, discriminative | 91.3 |
| Zhu et al. (2013) [40] | semi-supervised | 91.3 |
| Huang & Harper (2009) [14] | semi-supervised | 91.3 |
| McClosky et al. (2006) [26] | semi-supervised | 92.1 |
| Vinyals & Kaiser el al. (2014) [37] | semi-supervised | 92.1 |
| Transformer (4 layers) | semi-supervised | 92.7 |
| Luong et al. (2015) [23] | multi-task | 93.0 |
| Dyer et al. (2016) [8] | generative | 93.3 |

increased the maximum output length to input length + 300. We used a beam size of 21 and $\alpha = 0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8].

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

**7  Conclusion**

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/tensorflow/tensor2tensor.

Description

Claim

Conclusion

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., … & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

- Generating definitions of scientific and medical terms with varying complexity

- Dataset: Medical Consumer Questions (Abacha et al., 2019), Wikipedia, S2ORC

- Method
  - Reranking BART generation results based on logits from a discriminator trained to distinguish scientific journal text from science news text (measuring complexity)

| Control Method | Direction | |
|---|---|---|
| | **Low (News)** | **High (Journal)** |
| SVM-Rerank | A type of computing in which there are many computers running at the same time in different parts of the world. | In computer science, distributed computing is the process of computing on a large scale **without a single centralized data center**. |
| BERT-Rerank | A type of computer system in which there are more than a few computers working together. | In computer science, distributed computing is the process of computing on a large scale **without a single centralized data center**. |
| GeDi | Is the implementation of computer programs across multiple computers on similar hardware and/or software resources. | In computer science, **a concept that states that data must be shared across computing resources**. |
| DExpert | An **Internet-driven by-computing** that portion of different computers from start to finish. | In computer science and communication between-Consequently-integrates. |
| PPLM | **Easeless, self-organized, and often self-organizing** networked computer systems intended for the purposes of optimization. | Multi-purpose, distributed **system software** with or without a single datum storage system. |

**Points for Improvement:**
- Instead of reranking, inject complexity as an additional input feature in the input
- Employ a generative adversarial framework to optimize the BART generator and reranking discriminator jointly

August, T., Reinecke, K., & Smith, N. A. (2022, May). Generating scientific definitions with controllable complexity. ACL 2022.
Asma Ben Abacha and Dina Demner-Fushman. 2019. A question-entailment approach to question answering. BMC Bioinformatics, pages 1–23.

- Writing each part of a paper is a different task!

**Title**

**Attention Is All You Need**

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

**Related Work**

**Abstract**

## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**Definition**

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

## 2  Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence-aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

| Parser | Training | WSJ 23 F1 |
| --- | --- | --- |
| Vinyals & Kaiser el al. (2014) [37] | WSJ only, discriminative | 88.3 |
| Petrov et al. (2006) [29] | WSJ only, discriminative | 90.4 |
| Zhu et al. (2013) [40] | WSJ only, discriminative | 90.4 |
| Dyer et al. (2016) [8] | WSJ only, discriminative | 91.7 |
| Transformer (4 layers) | WSJ only, discriminative | 91.3 |
| Zhu et al. (2013) [40] | semi-supervised | 91.3 |
| Huang & Harper (2009) [14] | semi-supervised | 91.3 |
| McClosky et al. (2006) [26] | semi-supervised | 92.1 |
| Vinyals & Kaiser el al. (2014) [37] | semi-supervised | 92.1 |
| Transformer (4 layers) | semi-supervised | 92.7 |
| Luong et al. (2015) [23] | multi-task | 93.0 |
| Dyer et al. (2016) [8] | generative | 93.3 |

increased the maximum output length to input length + 300. We used a beam size of 21 and $\alpha = 0.3$ for both WSJ only and the semi-supervised setting.

**Description**

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8].

**Claim**

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

**Conclusion**

## 7  Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/tensorflow/tensor2tensor.

- Related work generation usually involves citation networks

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

# Knowledge-guided Citation Generation

- **Input**
  - Background knowledge: structural information from a citation network
  - Content: In-depth content information

- **Target: Citation sentences**

- **Model Highlights**
  - Salience Estimation to identify what information needed to be concentrated for the generation (what to cite)
  - Citation Function Classification to classify the reason for citation (why to cite)

What if citation relations are not given?

Ge, Y., Dinh, L., Liu, X., Su, J., Lu, Z., Wang, A., & Diesner, J. (2021, August). BACO: A background knowledge-and content-based framework for citing sentence generation. ACL 2021

- ## Model
  - Graph Structure Encoder with dual-role attention to capture directed graph structure information
  - Gated Neural Fusion to integrate multi-modal representations and control features transfer for downstream tasks differently
  - Multi-Task Decoder to perform citation link prediction and context generation



Can we combine information retrieval and citation recommendation?

Wang, Q., Xiong, Y., Zhang, Y., Zhang, J., & Zhu, Y. (2021). AutoCite: Multi-Modal Representation Fusion for Contextual Citation Generation. WSDM 2021

- A pipeline that automatically recommends relevant papers, extracts highlights, and suggests a reference sentence, given the user-provided context and keywords
  - Prefetching papers from database and reranking them based on keywords by SciBERT
  - Iteratively score and select sentences for summarization with MemSum
  - Concatenate keywords, context and target paper's abstract to generate citation sentences by T5



Gu, N., & Hahnloser, R. H. (2023). SciLit: A Platform for Joint Scientific Literature Discovery, Summarization and Citation Generation. ACL 2023 Demo.

# Summary of Citation Generation

| Method | Pros | Cons |
| --- | --- | --- |
| BACO (Ge et al., 2021) | ▪Construct a BACO dataset specifically for citation generation task<br>▪The framework identifies the cited sentence and the reason for citing, jointly with the generation of citations | ▪Rely on provided citation relations of the paper to generate citation sentence |
| AutoCite(Wang et al., 2021) | ▪Propose a joint framework to predict citation links and generate citation sentences simultaneously<br>▪Design a graph encoder to distinguish citation directions | ▪Fail to generate different types of citation texts according to the contexts and positions in the target paper |
| SciLit (Gu et al., 2023) | ▪Propose a pipeline to recommends relevant papers, extract highlights, and suggest a reference sentence | ▪The citation generation ignores the graph structure of citation networks |

Ge, Y., Dinh, L., Liu, X., Su, J., Lu, Z., Wang, A., & Diesner, J. (2021, August). BACO: A background knowledge-and content-based framework for citing sentence generation. ACL 2021
Wang, Q., Xiong, Y., Zhang, Y., Zhang, J., & Zhu, Y. (2021). AutoCite: Multi-Modal Representation Fusion for Contextual Citation Generation. WSDM 2021
Gu, N., & Hahnloser, R. H. (2023). SciLit: A Platform for Joint Scientific Literature Discovery, Summarization and Citation Generation. ACL 2023 Demo.

- Paper components that benefit from existing systems

**Title**

## Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

**Related Work**

**Abstract**

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**Definition**

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

## 2   Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence-aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

| Parser | Training | WSJ 23 F1 |
|---|---|---|
| Vinyals & Kaiser el al. (2014) [37] | WSJ only, discriminative | 88.3 |
| Petrov et al. (2006) [29] | WSJ only, discriminative | 90.4 |
| Zhu et al. (2013) [40] | WSJ only, discriminative | 90.4 |
| Dyer et al. (2016) [8] | WSJ only, discriminative | 91.7 |
| Transformer (4 layers) | WSJ only, discriminative | 91.3 |
| Zhu et al. (2013) [40] | semi-supervised | 91.3 |
| Huang & Harper (2009) [14] | semi-supervised | 91.3 |
| McClosky et al. (2006) [26] | semi-supervised | 92.1 |
| Vinyals & Kaiser el al. (2014) [37] | semi-supervised | 92.1 |
| Transformer (4 layers) | semi-supervised | 92.7 |
| Luong et al. (2015) [23] | multi-task | 93.0 |
| Dyer et al. (2016) [8] | generative | 93.3 |

increased the maximum output length to input length + 300. We used a beam size of 21 and $\alpha = 0.3$ for both WSJ only and the semi-supervised setting.

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8].

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

**Description**

**Claim**

**Conclusion**

## 7   Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/tensorflow/tensor2tensor.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

# Scientific Context-Aware Description Generation

- **Input**
  - A paper object (e.g., figure, table, algorithm, theorem)
  - Context (Paragraphs before the descriptions)

- **Target**
  - Descriptions

- **Dataset**

| Input | #num | #avg_out_len | #cand |
|-------|------|--------------|-------|
| Table | 136K | 74.05 | 199.93 |
| Figure (chart/bar) | 155K | 76.60 | 179.35 |
| Algorithm | 56K | 67.94 | 227.65 |
| Theorem | 175K | 65.00 | 192.44 |

| **Body text (Context)** | | | |
|---|---|---|---|
| . . . languages: Telugu (te) and Turkish (tr). . . Turkish (tr) vocabulary has been censored to contain no overlap with the Telugu. . . we evaluate these models using a recall@k metric defined as . . . | | | |
| **Table** | | | |

| | Result | **te+en** | **+ tr** | % Change | |
|---|--------|-----------|----------|----------|---|
| | Recall@1 | 17.0 | 17.6 | +3.5% | |
| | Recall@10 | 23.9 | 25.0 | +4.6% | |
| | Recall@20 | 26.3 | 27.7 | +5.3% | |

| **Generated description w/o context (table only)** |
|---|
| Table shows when te+en is replaced with tr, the effect of different change is very small, although the performance of tr method gets really strong. |
| **Generated description w/ context (body text + table)** |
| Table summarizes the recall@1 measures and the percentage of the incremental improvement across languages for both tasks. The average incremental improvement across languages is about 4% in these cases, despite there being no overlap between in Telugu and Turkish. |

Chen, H., Takamura, H., & Nakayama, H. (2021). SciXGen: a scientific paper dataset for context-aware text generation. EMNLP 2021 Findings.

- Paper components that benefit from existing systems

**Title**

### Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

**Related Work**

**Abstract**

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**Definition**

In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

## 2 Background

The goal of reducing sequential computation also forms the foundation of the Extended Neural GPU [16], ByteNet [18] and ConvS2S [9], all of which use convolutional neural networks as basic building block, computing hidden representations in parallel for all input and output positions. In these models, the number of operations required to relate signals from two arbitrary input or output positions grows in the distance between positions, linearly for ConvS2S and logarithmically for ByteNet. This makes it more difficult to learn dependencies between distant positions [12]. In the Transformer this is reduced to a constant number of operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence-aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

| Parser | Training | WSJ 23 F1 |
|---|---|---|
| Vinyals & Kaiser el al. (2014) [37] | WSJ only, discriminative | 88.3 |
| Petrov et al. (2006) [29] | WSJ only, discriminative | 90.4 |
| Zhu et al. (2013) [40] | WSJ only, discriminative | 90.4 |
| Dyer et al. (2016) [8] | WSJ only, discriminative | 91.7 |
| Transformer (4 layers) | WSJ only, discriminative | 91.3 |
| Zhu et al. (2013) [40] | semi-supervised | 91.3 |
| Huang & Harper (2009) [14] | semi-supervised | 91.3 |
| McClosky et al. (2006) [26] | semi-supervised | 92.1 |
| Vinyals & Kaiser el al. (2014) [37] | semi-supervised | 92.1 |
| Transformer (4 layers) | semi-supervised | 92.7 |
| Luong et al. (2015) [23] | multi-task | 93.0 |
| Dyer et al. (2016) [8] | generative | 93.3 |

increased the maximum output length to input length + 300. We used a beam size of 21 and $\alpha = 0.3$ for both WSJ only and the semi-supervised setting.

**Description**

Our results in Table 4 show that despite the lack of task-specific tuning our model performs surprisingly well, yielding better results than all previously reported models with the exception of the Recurrent Neural Network Grammar [8].

**Claim**

In contrast to RNN sequence-to-sequence models [37], the Transformer outperforms the Berkeley-Parser [29] even when training only on the WSJ training set of 40K sentences.

**Conclusion**

## 7 Conclusion

In this work, we presented the Transformer, the first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.

For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers. On both WMT 2014 English-to-German and WMT 2014 English-to-French translation tasks, we achieve a new state of the art. In the former task our best model outperforms even all previously reported ensembles.

We are excited about the future of attention-based models and plan to apply them to other tasks. We plan to extend the Transformer to problems involving input and output modalities other than text and to investigate local, restricted attention mechanisms to efficiently handle large inputs and outputs such as images, audio and video. Making generation less sequential is another research goals of ours.

The code we used to train and evaluate our models is available at https://github.com/tensorflow/tensor2tensor.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

# Scientific Claim Generation

- ## Supported Claims Generation
  - ClaimBART: Generating claim directly based on Citance
  - ClaimEntity



Citance | Named Entity Recognition (scispacy) | Question Generation (BART) | Claim Generation (BART)

- ## Refuted Claims Generation
  - Knowledge-based Informed Negation



**Points for Improvement:**
- Instead of only considering entities in the claim, explore to use relations or even a subgraph of the knowledge graph when generating claims



(1) ALS is the most common adult motor neuron disease with an incidence of 2 per 100,000 and prevalence of 5.4 per 100,000 individuals. (2) Current treatment options are based on symptom management and respiratory support with the only approved medications in widespread use, Riluzole and Edaravone, providing only modest benefits and only in some patients.

$c_1$: Current treatment options for ALS are based on symptom management and respiratory support

$c_2$: Riluzole is an approved ALS medication in widespread use

$c_3$: Edaravone is an approved ALS medication in widespread use

$c_4$: Riluzole and Edaravone are the only approved ALS medications in widespread use

$c_5$: Riluzole provides modest benefits in only some ALS patients

$c_6$: Edaravone provides modest benefits in only some ALS patients

Wright, D., Wadden, D., Lo, K., Kuehl, B., Cohan, A., Augenstein, I., & Wang, L. L. (2022). Generating scientific claims for zero-shot scientific fact checking. ACL 2022.

# Human-AI Collaborative Writing

- Compared to either humans or AI, collaborative Human-AI Writing has:
  - Fewer grammatical errors
  - Higher vocabulary diversity

- LLMs can increase writers' productivity by providing writing suggestions

- Keep fraction of text written by writers to text written by ChatGPT relatively high can increase writers' feeling of ownership



**Write** — The 92nd little pig built a house out of depleted uranium.

Of course, the little pig was smart. He knew he couldn't live in the house, since it would poison him.

**Get suggestions**
So he lived in the basement.

So he built a big wall around it, and he lived in a little shack on the other side of the wall.

So he moved next door to it.

He moved in a couple of miles away.

So he built a little shack next to it, and he painted it and put a little flower garden outside and everything.

**Edit** — The 92nd little pig built a house out of depleted uranium. There are so many pigs to eat, and they make it so easy!" But the wolf didn't know about uranium.

Of course, the little pig was smart. He knew he couldn't live in the house, since it would poison him. So he built a big wall around it, and he lived in a little, hidden shack on the other side of the wall.

Lee, M., Liang, P., & Yang, Q. (2022). CoAuthor: Designing a Human-AI Collaborative Writing Dataset for Exploring Language Model Capabilities. CHI 2022.

- **Interesting points**
  - AI-assisted paper draft writing systems are categorized under knowledge-augmented LLMs
  - These systems extensively utilize both graph and text features for generating results

- **Future directions**
  - To reduce error propagation, future developments could focus on creating an end-to-end system, moving away from the current pipeline-based approach
  - Enhancing the multi-hop reasoning capabilities of these systems is also crucial
  - Incorporating human feedback to iteratively improve system results presents an interesting future direction

# Completing Scientific Paper Lifecycle

Scientific Paper

Automatic Paper Review

Scientific Fact Checking

Paper Reviews

Draft Evaluation

Literature Survey

Hypothesis Generation

Experiment Planning

Paper Writing

# Challenges in Peer Review

- To catch up with the increasing submissions, reviewers are overburdened leading to contentious review comments

  - "There is significant evidence that the process of reviewing papers in machine learning is creaking under several years of exponentiating growth." -- ICML President John Langford

  - 12.9% of reviews were flagged by authors for having issues, and 1.69% of reviews were noted as rude and unprofessional -- ACL 2023 Peer Review Report

  - Authors have even created various social groups at social media to release their frustrations and anger, such as the *"Reviewer #2 must be stopped"* group at Facebook

| Type of issue | Number of reviews | % of reviews |
|---|---|---|
| A: The review is not specific enough | 272 | 2.16 |
| B: Review heuristics such as "not novel", "not surprising", "too simple", "not SOTA" | 678 | 5.38 |
| C: The scores do not match the review text | 448 | 3.55 |
| D: The review is rude/unprofessional | 213 | 1.69 |
| E: The review does not evince expertise | 542 | 4.3 |
| F: The review does not match the paper type | 98 | 0.78 |
| G: The review does not match the type of contribution | 152 | 1.21 |
| H: The review is missing or too short | 205 | 1.63 |
| I: The review was late | 12 | 0.1 |
| J: Other | 162 | 1.29 |

# Challenges in Peer Review

- The following are two different reviews for the same paper rejected by ACL2019 and accepted by EMNLP2019 without any change in content:

  - **ACL 2019**: "Idea is too simple and tricky."

  - **EMNLP 2019**: "The main strengths of the paper lie in the interesting, relatively under-researched problem it covers, the novel and valid method and the experimental results."

# Challenges in Peer Review

- Human reviews are subjective!

- We need an automatic review assistant to help human reviewers generate knowledgeable and explainable review scores and comments, along with detailed evidence



ADDRESSING REVIEWER COMMENTS — BAD REVIEWS ON YOUR PAPER? FOLLOW THESE GUIDELINES AND YOU MAY YET GET IT PAST THE EDITOR:

**Reviewer comment:**
"The method/device/paradigm the authors propose is clearly wrong."

**How NOT to respond:**
✗ "Yes, we know. We thought we could still get a paper out of it. Sorry."

**Correct response:**
✓ "The reviewer raises an interesting concern. However, as the focus of this work is exploratory and not performance-based, validation was not found to be of critical importance to the contribution of the paper."

**Reviewer comment:**
"The authors fail to reference the work of Smith et al., who solved the same problem 20 years ago."

**How NOT to respond:**
✗ "Huh. We didn't think anybody had read that. Actually, their solution is better than ours."

**Correct response:**
✓ "The reviewer raises an interesting concern. However, our work is based on completely different first principles (we use different variable names), and has a much more attractive graphical user interface.

**Reviewer comment:**
"This paper is poorly written and scientifically unsound. I do not recommend it for publication."

**How NOT to respond:**
✗ "You #&@*% reviewer! I know who you are! I'm gonna get you when it's my turn to review!"

**Correct response:**
✓ "The reviewer raises an interesting concern. However, we feel the reviewer did not fully comprehend the scope of the work, and misjudged the results based on incorrect assumptions.

JORGE CHAM © 2005

www.phdcomics.com

# Challenges in Automatic Paper Review

- Generation Results for <u>Sparse Weight Activation Training</u> (Raihan et al., 2020) in Kid-Review's Appendix (Yuan et al., 2020)

- This paper proposes a method to sparsify the weights and activations of a neural network by sparsifying the gradients in the backward pass . The idea is inspired by meProp , which sparsifies gradients with relatively small magnitude . The authors show that the proposed method can reduce the memory footprint by 23 % to 37 % for activations and 50 % to 80 % for weights . The paper is well written and easy to follow . The experimental results on CIFAR-10 and ImageNet are promising . However , I have the following concerns : 1 . It is not clear to me why the authors chose to use SAW instead of meProp as the baseline method . It seems to me that meProp is a better baseline than SAW . 2 .The authors claim that the memory benefit of SAW is not present for meProp since there is no storage benefit since they are temporary values generated during back-propagation . But I don't see any evidence to support this claim . 3 .In the experiments , the authors only compare with SAW and meProp . It would be more convincing if the authors can also compare with other sparsification methods , e.g. , [ 1 ] and [ 2 ] …

> SAW performs better than meProp in Figure 5

> mePROP is compared against e Sparse Weight Activation Training (SWAT)

> Missing references

Only the summarization part works
Other parts suffer from hallucination

Yuan, W., & Liu, P. (2022, June). Kid-review: Knowledge-guided scientific review generation with oracle pre-training. AAAI 2022.
Raihan, M. A., & Aamodt, T. (2020). Sparse weight activation training. Advances in Neural Information Processing Systems, 33, 15625-15638.

# Challenges in Automatic Paper Review

- **Peer review dataset is limited**
  - PeerRead Dataset (Kang et al., 2018)
    - ➤ NeurIPS 2013-2017, ICLR 2017, ACL 2017, and CoNLL 2016
    - ➤ Annotate 1.3k reviews for ICLR 2017 with aspect scores based on reviews
  - ReviewRobot Dataset (Wang et al., 2020)
    - ➤ NeurIPS 2013-2018, ICLR 2017-2020, ACL 2017
    - ➤ Background KG from 174,165 papers from 1965 to 2019
  - ASAP-Review Dataset (Yuan et al., 2022)
    - ➤ NeurIPS 2016-2019, ICLR 2017-2020, ACL 2017
    - ➤ Annotate aspect scores with BERT

Kang, D., Ammar, W., Dalvi, B., Van Zuylen, M., Kohlmeier, S., Hovy, E., & Schwartz, R. (2018). A dataset of peer reviews (PeerRead): Collection, insights and NLP applications. NAACL 2018.
Wang, Q., Zeng, Q., Huang, L., Knight, K., Ji, H., & Rajani, N. F. (2020, December). ReviewRobot: Explainable paper review generation based on knowledge synthesis. INLG 2020.
Yuan, W., Liu, P., & Neubig, G. (2022). Can we automate scientific reviewing?. Journal of Artificial Intelligence Research, 75, 171-212.

# Category of Automatic Paper Review

# Explainable Paper Review Generation

- Propose a knowledge-driven review score prediction and comment generation framework, ReviewRobot, based on fine-grained knowledge element comparison among papers

- The review comments are highly explainable and constructive, supported by detailed evidence

Pros:
- Both score prediction and review generation are grounded on evidence

Points for Improvement:
- The pipeline system heavily rely on the quality of IE systems, which might propagate errors
- The review comment templates are human-curated and lack flexibility

Wang, Q., Zeng, Q., Huang, L., Knight, K., Ji, H., & Rajani, N. F. (2020, December). ReviewRobot: Explainable paper review generation based on knowledge synthesis. INLG 2020.

# Paper Review Evidence: Meaningful Comparison

- The number of papers about relevant old knowledge elements which are missed in the related work section: $G_{P_\tau} \cap G_B - \overline{G_{P_\tau}}$

- The number of papers about relevant knowledge elements which are claimed new in the related work section: $G_{P_\tau} \cap G_B \cap \overline{G_{P_\tau}}$

- The description sentences about comparison with related work

Cui, Y., Chen, Z., Wei, S., Wang, S., Liu, T., & Hu, G. (2017). Attention-over-Attention Neural Networks for Reading Comprehension. ACL 2017.
Qingyun Wang, Qi Zeng, Lifu Huang, Kevin Knight, Heng Ji, Nazneen Fatema Rajani. ReviewRobot: Explainable Paper Review Generation based on Knowledge Synthesis. INLG 2020.

# Explainable Paper Review Generation

- **[ReviewRobot] The following related papers are missing:**
  - About low-dimensional semantic space:
    - Unsupervised Approximate-semantic Vocabulary Learning for Human Action and Video Classification: Qiong Zhao and Horace HS Ip. 2013. Unsupervised Approximate-semantic Vocabulary Learning for Human Action and Video Classification. *Pattern Recognition Letters*, 34(15):1870–1878.
  - About sememes:
    - Chinese Word Sense Disambiguation with PageRank and HowNet: Jinghua Wang, Jianyi Liu, and Ping Zhang. 2008. Chinese Word Sense Disambiguation with PageRank and HowNet. In Proceedings of the Sixth SIGHAN Workshop on Chinese Language Processing.
    - A Maximum Entropy Approach to HowNet-based Chinese Word Sense Disambiguation: Ping Wai Wong and Yongsheng Yang. 2002. A Maximum Entropy Approach to HowNet-based Chinese Word Sense Disambiguation. In COLING-02: SEMANET: Building and Using Semantic Networks.
  - About word similarity and word analogy:
    - Open IE as an Intermediate Structure for Semantic Tasks: Gabriel Stanovsky, Ido Dagan, et al. 2015. Open IE as an Intermediate Structure for Semantic Tasks. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 303–308.

- **[HUMAN] The paper would be stronger with the inclusion of more baselines based on related work.**

Wang, Q., Zeng, Q., Huang, L., Knight, K., Ji, H., & Rajani, N. F. (2020, December). ReviewRobot: Explainable paper review generation based on knowledge synthesis. INLG 2020.

# Semi-supervised Learning Aspect Score Prediction

- Propose Γ-Trans incorporating a pretrained transformer into semi-supervised learning, by leveraging contextual features from unlabeled data

- Unlabeled dataset
  - ScisummNet Corpus with 1,000 papers in the ACL Anthology

- Help improve the performance of all pretrained LMs including SciBERT and Longformer

- Γ-Trans with 100 unlabeled data outperforms other methods with 700 unlabeled data

**Points for Improvement:**
- The system can further benefit from using the related information between review aspects

Corrupt Path    Clean Path

$Tr^{(l)}$ refers to the transformer block at layer $l$

$f(\cdot)$ is linear transformation

$\mathcal{N}(\mu, \sigma^2)$ refers to Gaussian noise

$g(\cdot, \cdot)$ is ladder network denoising function

$C_s$ denotes supervised cost (Cross−entropy)

$C_d$ denotes denoising cost (L1)

Total cost $C = C_s + C_d$

Pre-trained LM

Muangkammuen, P., Fukumoto, F., Li, J., & Suzuki, Y. (2022, December). Exploiting labeled and unlabeled data via transformer fine-tuning for peer-review score prediction. EMNLP 2022 Findings.
Michihiro Yasunaga, Jungo Kasai, Rui Zhang, Alexander R. Fabbri, Irene Li, Dan Friedman, and Dragomir R. Radev. 2019. Scisummnet: A large annotated corpus and content-impact models for scientific paper summarization with citation networks. AAAI 2019

- An end-to-end knowledge-guided review generation framework based on the citation graph and the concept graph
  - Encode concept graph with Graph Attention Network
  - Use citation embeddings as additional input

Pros:
- The paper generates review based on both citation and knowledge graphs since human beings rely on both background knowledge and working memory to review a paper

Points for Improvement:
- The proposed model still requires an extractor to extract important sentences from the source article and only use the extracted sentences
- The generated review contains many unfaithful contents (as shown in Slide 156)



(a) World knowledge – citation graph

(b) Temporary knowledge – concept graph

Yuan, W., & Liu, P. (2022, June). Kid-review: Knowledge-guided scientific review generation with oracle pre-training. AAAI 2022.
Tang, J.; Qu, M.; Wang, M.; Zhang, M.; Yan, J.; and Mei, Q. 2015. Line: Large-scale information network embedding. ICML 2015
Lewis, M.; Liu, Y.; Goyal, N.; Ghazvininejad, M.; Mohamed,A.; Levy, O.; Stoyanov, V.; and Zettlemoyer, L. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. ACL 2020
Feigenblat, G., Roitman, H., Boni, O., & Konopnicki, D. (2017, August). Unsupervised query-focused multi-document summarization using the cross entropy method. SIGIR 2017

# ChatGPT for Paper Review

- Test ChatGPT with papers using prompt engineering and chain-of-thought

- Identifying errors
  - ChatGPT completely fails on sorting, pairwise comparison, game theory, and optimization

- Verifying checklists
  - When verifying 16 closed-ended checklist questions in 15 NeurIPS 2022 papers, ChatGPT achieves 86.6% accuracy

- Choosing the "better" paper
  - When comparing 10 pairs of abstracts which deliberately designed in such a way that one abstract was clearly superior to the other, ChatGPT struggled to discern these relatively straightforward distinctions accurately, committing errors in its evaluations for 6 out of the 10 pairs

LLMs such as ChatGPT, without tailored framework designs, remain far from expected solutions

Liu, R., & Shah, N. B. (2023). Reviewergpt? an exploratory study on using large language models for paper reviewing. arXiv preprint arXiv:2306.00622.
Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., ... & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35, 24824-24837.

# Meta-review Generation

- Input: the title, abstract, and a set of reviews of one research paper

- Goal: Generate a meta-review summarizing the opinions in the independent reviews and make a recommendation decisions

- Checklist-guided Iterative Introspection
  - Break down the task into several stages
  - Iteratively refine the summary under the guidance of questions from a checklist



**Points for Improvement:**
- The proposed model fails to utilize any knowledge graphs, which might lead to hallucination
- The system only uses titles and abstracts not the whole paper

Zeng, Q., Sidhu, M., Chan, H. P., Wang, L., & Ji, H. (2023). Meta-review generation with checklist-guided iterative introspection. arXiv preprint arXiv:2305.14647.

# Meta-review Generation

**Human**

The paper unanimously receives positive rates thanks to strong motivations and interesting results. As the reviews show satisfaction on the authors' feedback, the final draft needs to respect it accordingly, for example, about the limitations of this research.

**Vanilla**

This paper could have benefited from a clearer and more concise presentation of the main findings. It is sometimes difficult to follow the paper's chain of thought, and the authors could have presented their results in a more accessible way. The experiments are limited to a small set of environments, and it would be interesting to see how the results generalize to other tasks. Additionally, it would be useful to have a deeper analysis of why iGPT does not perform well, as this could help inform future research on how to improve image-based pretraining.

**Initial run**

Based on the reviews, I recommend accepting the paper. The paper addresses an important and under-explored research question and presents a detailed and thought-provoking analysis of the pretraining-finetuning paradigm for Transformer models in the context of reinforcement learning. The paper also identifies relevant techniques to perform analysis of representation similarity and training dynamics. While there are some weaknesses pointed out, such as the need for further exploration of some hypotheses and occasional difficulty in following the writing and presentation, the strengths of the paper outweigh these concerns. Therefore, I recommend accepting the paper with minor revisions to address the weaknesses pointed out by the reviewers.
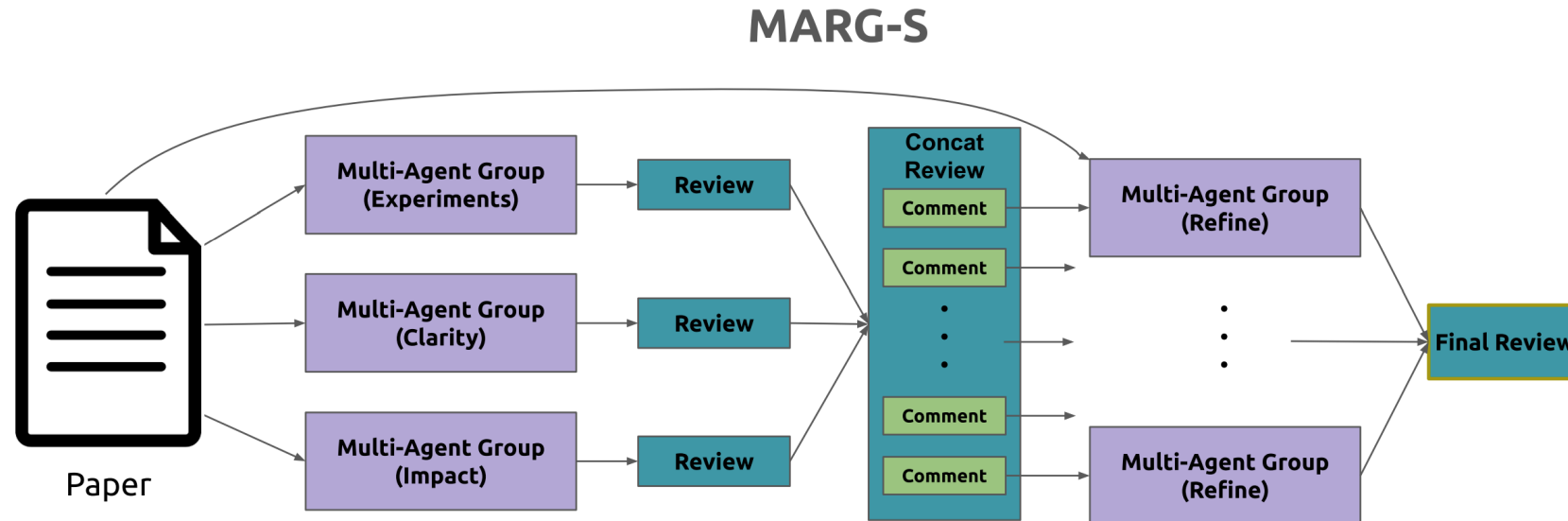
**CGI²**

After carefully considering the reviews provided, there is consensus that the paper presents an important and thought-provoking analysis of representation dynamics between pretraining and finetuning in the context of reinforcement learning. Both reviewers agree that the paper addresses an important question and presents insightful experiments. One reviewer highlights the experiments as a strong point in the paper, while the other reviewer raises a concern about the paper's tendency to present hypotheses without exploring them further.

Overall, the strengths of the paper outweigh its weaknesses, and it has the potential to make a significant contribution to the field of reinforcement learning. Therefore, I recommend accepting the paper with a minor revision to address the concern raised by the reviewer. Specifically, the paper should explore the hypotheses presented in more depth to provide a more comprehensive analysis. If the authors can address this concern, the paper can make a valuable contribution to the field of reinforcement learning. In summary, I recommend accepting the paper with a minor revision.

## Hallucination reduced

Zeng, Q., Sidhu, M., Chan, H. P., Wang, L., & Ji, H. (2023). Meta-review generation with checklist-guided iterative introspection. arXiv preprint arXiv:2305.14647.

# Multi-Agent Review Generation

- A network of LLM agents that communicate to share information across different parts of a paper and engage in internal discussion
  - Consume the full text of papers beyond the input length limitations of the base LLM
  - Include aspect-specific "expert" GPT agents to separately assist with generating comments on experiments, clarity, and impact
  - The method perform significantly better than when having a single agent attempt to generate all types of feedback at once
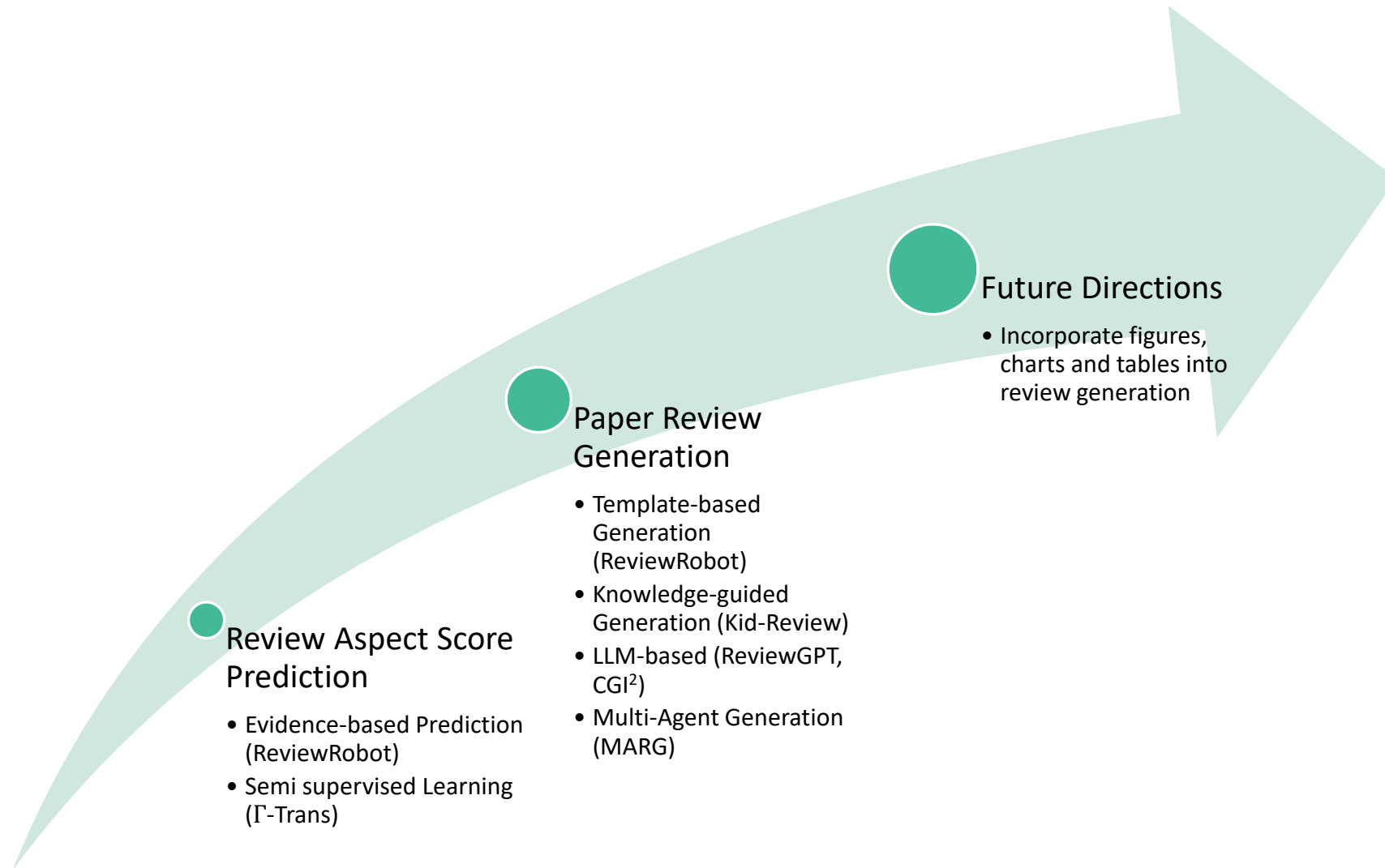
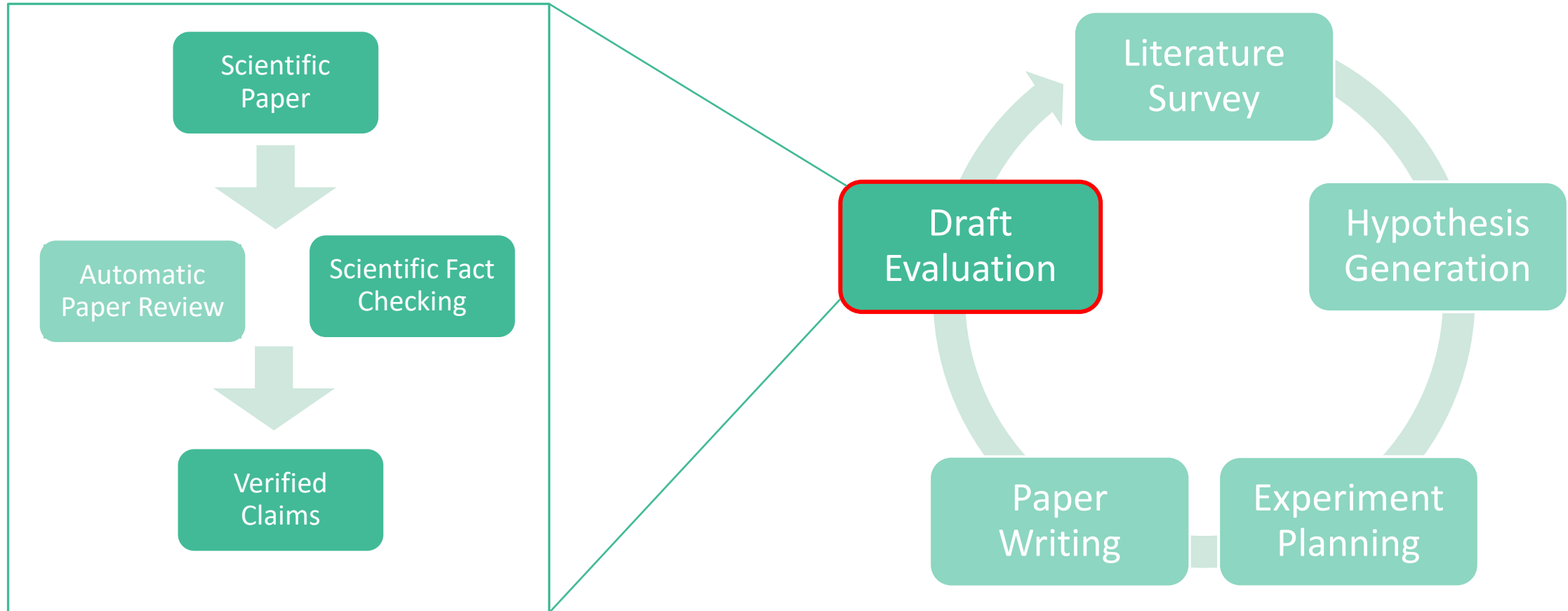**MARG-S**



Points for Improvement:
- The proposed model fails to utilize any knowledge graphs or background literature
- The proposed model relies heavily on GPT4, which is very expensive

D'Arcy, M., Hope, T., Birnbaum, L., & Downey, D. (2024). MARG: Multi-Agent Review Generation for Scientific Papers. arXiv preprint arXiv:2401.04259.

| Real-reviewer comment | Generated comment |
|---|---|
| **The experimental methodology used in the paper is not well detailed**, making it difficult to reproduce the reported results. | More details about the experiments conducted would be beneficial. This should **include information about the datasets used, the training process, and the evaluation process**. To ensure the reproducibility of the results, **consider providing the code used to implement the model, the specific parameters used, and any other necessary information**. This will allow other researchers to replicate your work and further validate your findings. **[high relatedness, more specific]** |
| The paper does not include enough baselines for Fair Federated Learning to compare against. Even if some methods do not satisfy privacy considerations, they should still be included for the reader to understand how the proposed method compares against such methods, especially given that the results are not promising. **Some baselines to consider include Cui et al or Tran et al.** | The authors should **consider including a comparison of their proposed method with existing methods** in the experimental results section. This would help to highlight the advantages and improvements of their proposed method. **[high relatedness, less specific]** |

Some reviews are vague

**Review Aspect Score Prediction**

- Evidence-based Prediction (ReviewRobot)
- Semi supervised Learning (Γ-Trans)

**Paper Review Generation**

- Template-based Generation (ReviewRobot)
- Knowledge-guided Generation (Kid-Review)
- LLM-based (ReviewGPT, CGI$^2$)
- Multi-Agent Generation (MARG)

**Future Directions**

- Incorporate figures, charts and tables into review generation

# Completing Scientific Paper Lifecycle

# The Danger of Misinformation

- ■ Threat to the Individual
  - Early in the pandemic, there was speculation about the internal use of disinfectants for treating COVID-19 patients.
    - ➢ On June 5, 2020, the Centers for Disease Control (CDC) reported a steep increase in calls to poison centers regarding exposure to household disinfectants.
    - ➢ A CDC survey of 502 adults in the United States found that 39% of responders engaged in dangerous practices including washing food products with bleach, applying household cleaners directly to skin, and intentionally inhaling or ingesting disinfectants with the goal of preventing COVID-19 infection.



### Coronavirus: Disinfectant firm warns after Trump comments

24 April 2020

◁ Share

Mr Trump made the comments at a government briefing on Thursday

A leading disinfectant producer has issued a strong warning not to use its products on the human body after Donald Trump suggested they could potentially be used to treat coronavirus.



### Fact check: Trump dangerously suggests sunlight and ingesting disinfectants could help cure coronavirus

By Daniel Dale, Nathan McDermott, Marshall Cohen, Maegan Vazquez, Em Steck and Sam Fossum, CNN
⏱ 6 minute read · Updated 2:18 PM EDT, Fri April 24, 2020

Nelson, T., Kagan, N., Critchlow, C., Hillard, A., & Hsu, A. (2020). The danger of misinformation in the COVID-19 crisis. Missouri Medicine, 117(6), 510.
https://www.bbc.com/news/world-us-canada-52411706
https://www.cnn.com/2020/04/23/politics/fact-check-coronavirus-briefing-april-23/index.html

# The Danger of Misinformation

- **Threat to the Scientific Research**
  - The *Lancet's* paper linked *hydroxychloroquine* use to harmful outcomes, prompting the WHO to halt its use in the Solidarity trial. However, due to questions about the study's methods and data reliability, the study was retracted, and the WHO later resumed the trials.



Home / News & Opinion

*Lancet, NEJM* Retract Surgisphere Studies on COVID-19 Patients

All authors other than company founder and CEO Sapan Desai were "unable to complete an independent audit of the data," *The Lancet* states.

Catherine Offord
Jun 4, 2020 | 3 min read

Two controversial studies of COVID-19 patients have been retracted after the authors failed to demonstrate that the data were reliable. The first study to be retracted, published last month (May 22) in *The Lancet*, had found harmful effects associated with the antimalarial drug, hydroxychloroquine, but quickly drew fire after scientists raised questions about the massive database supposedly underpinning it, and about that database's owner, Surgisphere Corporation.

ABOVE: © ISTOCK.COM, LILIBOAS


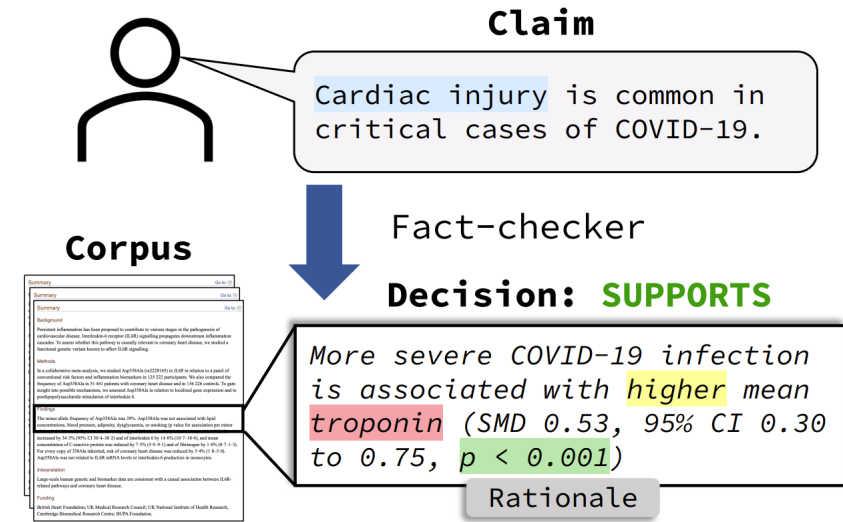
**The Lancet retracts large study on hydroxychloroquine**

"A first-year statistics major could tell you about major flaws in the design of the analysis," one expert said.

A large study suggesting hydroxychloroquine does not benefit COVID-19 patients, and may even increase deaths, has been retracted. But that doesn't mean hydroxychloroquine does – or does not – work. Buda Mendes / Getty Images

https://www.the-scientist.com/lancet-retracts-surgispheres-study-on-hydroxychloroquine-67613
https://www.nbcnews.com/health/health-news/lancet-retracts-large-study-hydroxychloroquine-n12250901
Nelson, T., Kagan, N., Critchlow, C., Hillard, A., & Hsu, A. (2020). The danger of misinformation in the COVID-19 crisis. Missouri Medicine, 117(6), 510.
Mehra MR, Desai SS, Ruschitzka F, et al. RETRACTED: Hydroxychloroquine or chloroquine with or without a macrolide for treatment of COVID-19: A multinational registry analysis. The Lancet. 2020 doi: 10.1016/s0140-6736(20)31180-6.

# What is scientific claim?

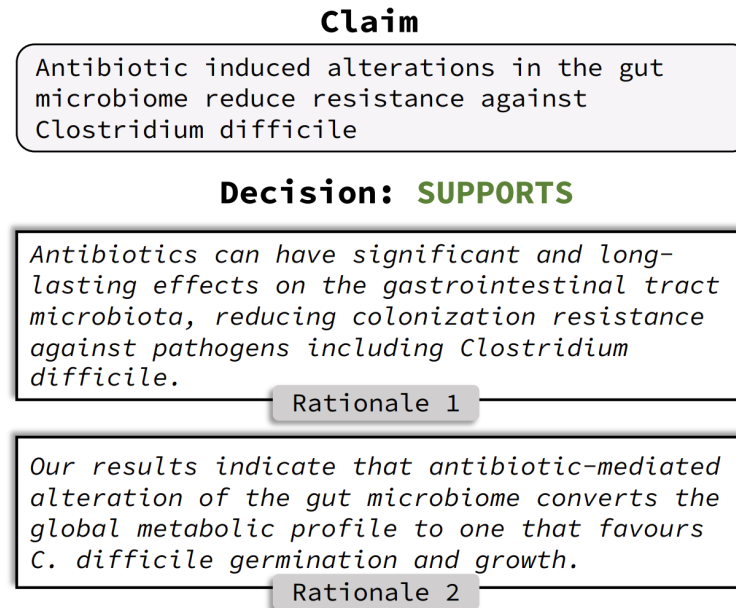- A scientific claim is an atomic verifiable statement expressing a finding about one aspect of a scientific entity or process, which can be verified from a single source

- Scientific Claim Verification Task
  - Given a scientific claim c and a corpus of abstracts $\mathcal{A}$, the system must predict a set of evidence abstracts $\mathcal{E}(c)$, where each abstract a $\in \mathcal{E}(c)$ has a predicted label (SUPPORTS, NOINFO, REFUTES) and a collection of rationale sentences

Wadden, D., Lin, S., Lo, K., Wang, L. L., van Zuylen, M., Cohan, A., & Hajishirzi, H. (2020). Fact or fiction: Verifying scientific claims. EMNLP 2020.

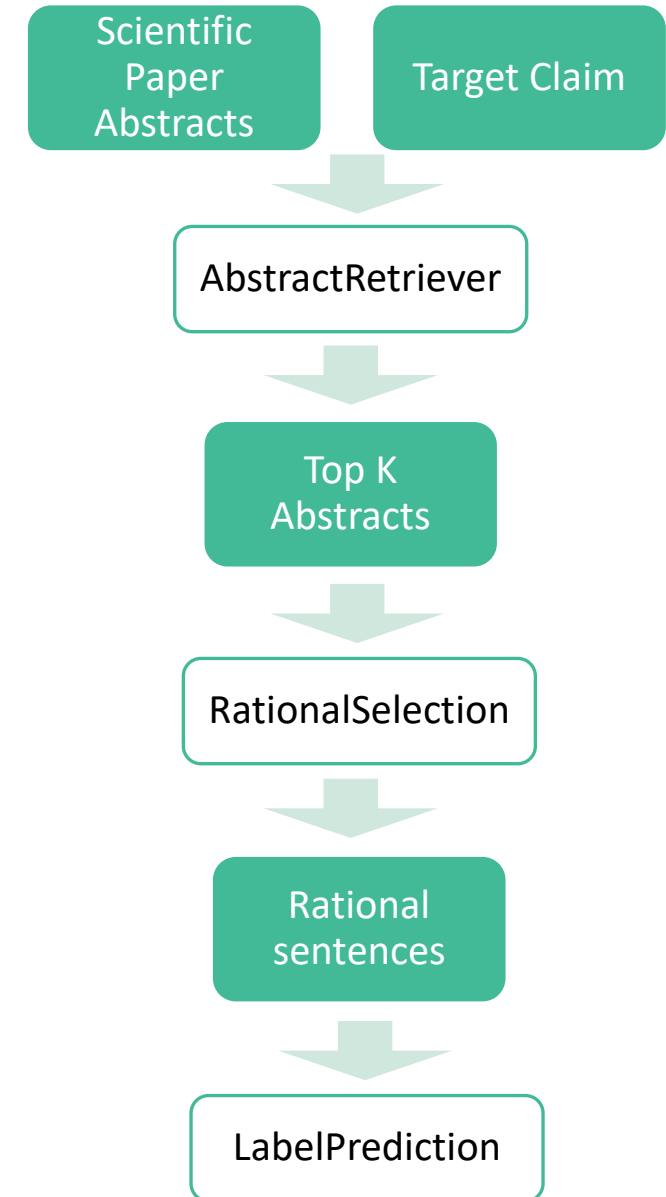# Category of Scientific Fact-Checking

# Scientific Claim Verification with Pipeline

- The architecture follows BERT-to-BERT pipeline
  - AbstractRetriever retrieves k abstracts with highest TF-IDF similarity to the claim
  - RationalSelection identifies rationale sentences for each abstract
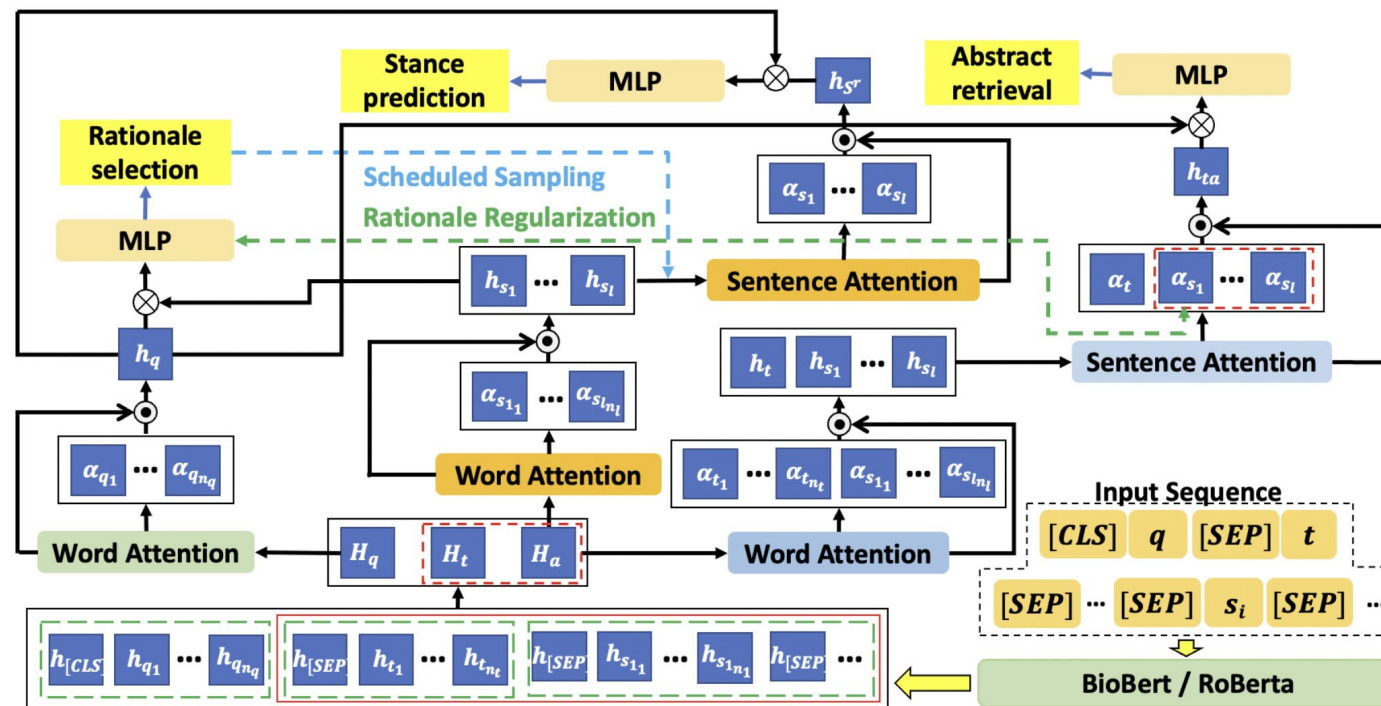  - LabelPrediction makes the final label prediction

**Claim**

Antibiotic induced alterations in the gut microbiome reduce resistance against Clostridium difficile

**Decision: SUPPORTS**

*Antibiotics can have significant and long-lasting effects on the gastrointestinal tract microbiota, reducing colonization resistance against pathogens including Clostridium difficile.*

Rationale 1

*Our results indicate that antibiotic-mediated alteration of the gut microbiome converts the global metabolic profile to one that favours C. difficile germination and growth.*

Rationale 2

Scientific Paper Abstracts → Target Claim

↓

AbstractRetriever

↓

Top K Abstracts

↓

RationalSelection

↓

Rational sentences

↓

LabelPrediction

**Points for Improvement:**
- Incorporate background knowledge from external knowledge bases
- Error might propagate through each steps

Wadden, D., Lin, S., Lo, K., Wang, L. L., van Zuylen, M., Cohan, A., & Hajishirzi, H. (2020). Fact or fiction: Verifying scientific claims. EMNLP 2020.

# Multi-Task Learning for Claim Verification

- An end-to-end multi-task framework, ARSJOINT, to jointly learn the three tasks based on QA Framework

- A rationale regularization based the divergence between the sentence attention of abstract retrieval and the outputs of rational selection to improve interpretability



Pros:
- Reduce error propagation by transforming the pipeline system to an end-to-end multi-task framework
- Share information between different tasks to improve the prediction quality

Points for Improvement:
- Require additional supervision with rationale sentences, limiting its training flexibility

Zhang, Z., Li, J., Fukumoto, F., & Ye, Y. (2021). Abstract, Rationale, Stance: A Joint Model for Scientific Claim Verification. EMNLP 2021.

# Weak Supervision and Full-document Context

- ## Training Procedure
  - Stage 1: Train on a combination of labeled out-of-domain data (FEVER) and weakly labeled in domain data (Paper title as claim, Abstract as unannotated rationales, label as supported)
  - Stage 2: Continue training on data from each target dataset
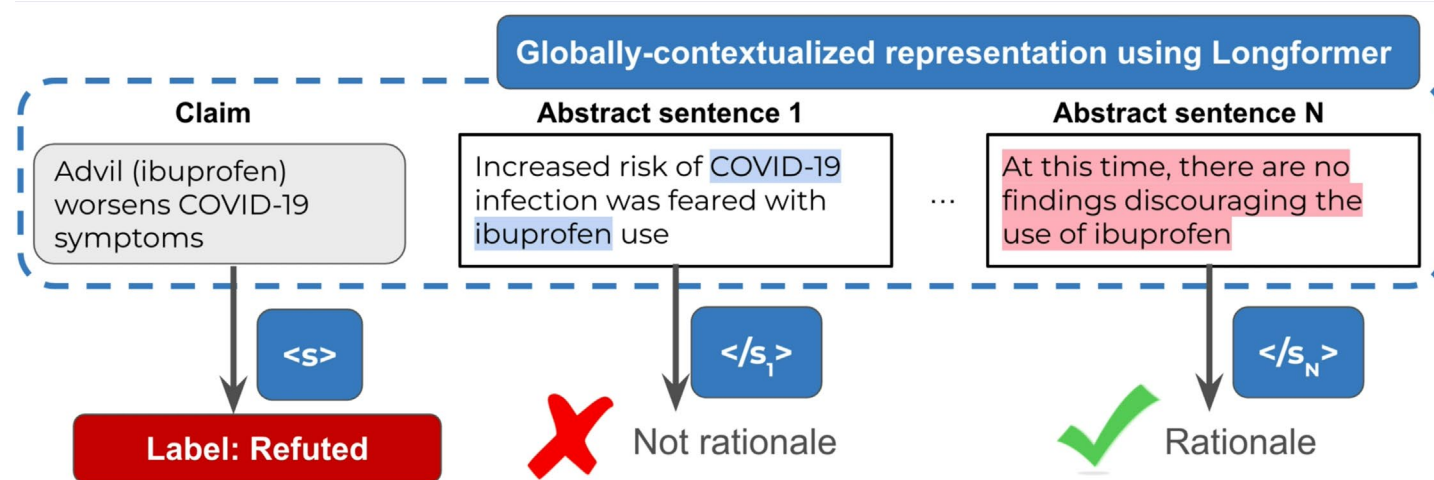
- ## Domain adaptation
  - Zero-shot (without Stage 2)
  - Few-shot (45 samples for stage 2)
  - Full (All samples for stage 2)

Pros:
- Take in all abstract sentences to leverage the shared information between them
- Can train on instances with no rationale sentence which is suitable for domain adaptation

Points for Improvement:
- Model still depends on the textual similarity without considering any structured information



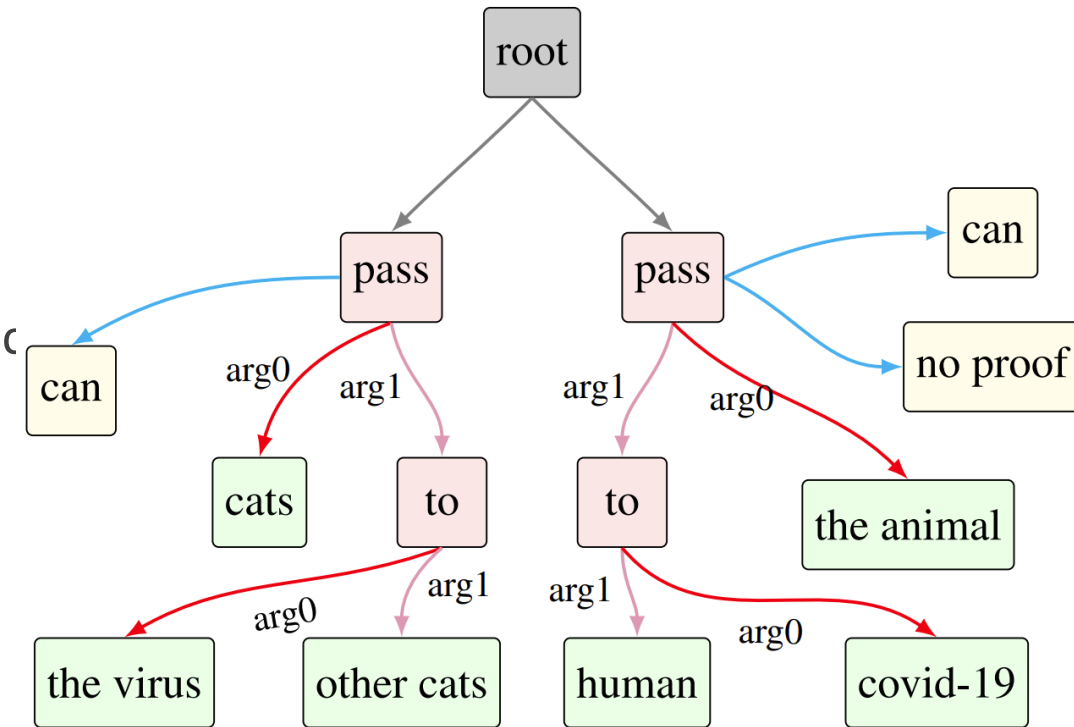$$L = L_{\text{label}} + \lambda_{\text{rationale}} L_{\text{rationale}}$$

Wadden, D., Lo, K., Wang, L. L., Cohan, A., Beltagy, I., & Hajishirzi, H. (2021). MultiVerS: Improving scientific claim verification with weak supervision and full-document context. NAACL 2022 Findings.

# Qualitative Causal Tree

- **Entities**
  - *Factors* are variables that are tested or asserted within a claim
  - *Associations* are explicit phrases associating one or more factors
  - *Modifier* aggregates the original entity types magnitude, evidence, and qualifiers
  - *Root* is designed to aggregate tree representations into a single vector representing the whole scientific sentence
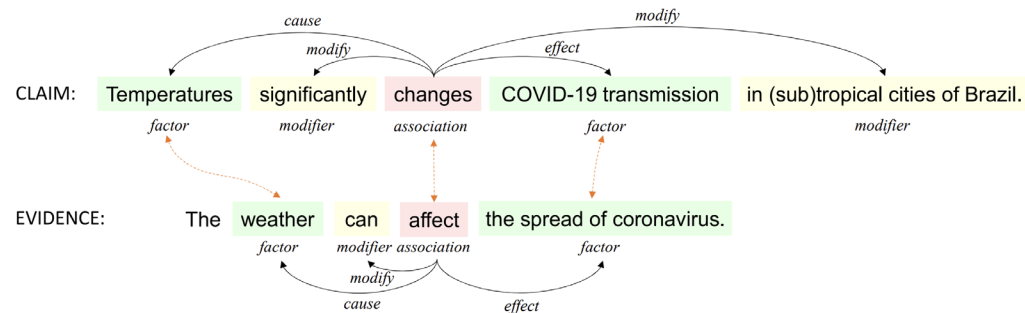
- **Relations**
  - *Associations* relates to entities cause or effect
  - *Modifying* relates associations to their modifying components
  - *comp_to* represents comparison in scientific claim verification



*Cats can pass the virus to other cats, they discovered no proof that the animals can pass COVID-19 to humans.*

Magnusson, I. H., & Friedman, S. E. (2021). Extracting fine-grained knowledge graphs of scientific claims: Dataset and transformer-based results. EMNLP 2021.
Wu, J., Chao, W., Zhou, X., & Luo, Z. (2023, December). Characterizing and Verifying Scientific Claims: Qualitative Causal Structure is All You Need. EMNLP 2023.

# Qualitative Causal Structure

- Transform the claim verification to match the roles of qualitative causal structures from claim and evidence respectively
  - Decompose heterogenous graph into factor graph and modifying graph
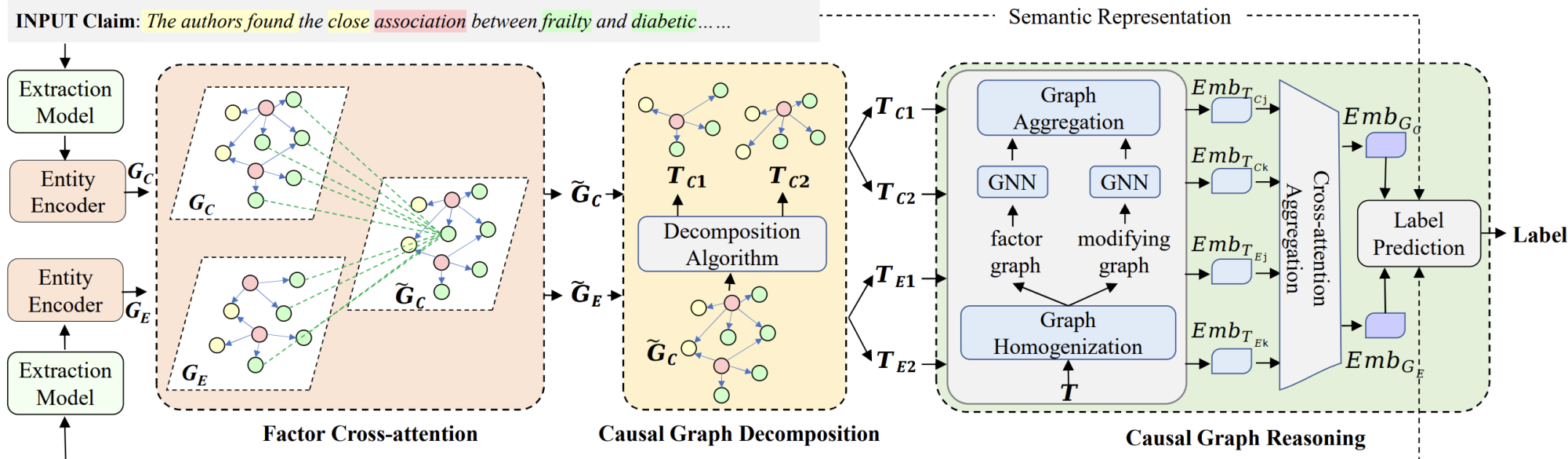


Pros:
- Provide explainability through the weights of causal sub-graphs

Points for Improvement:
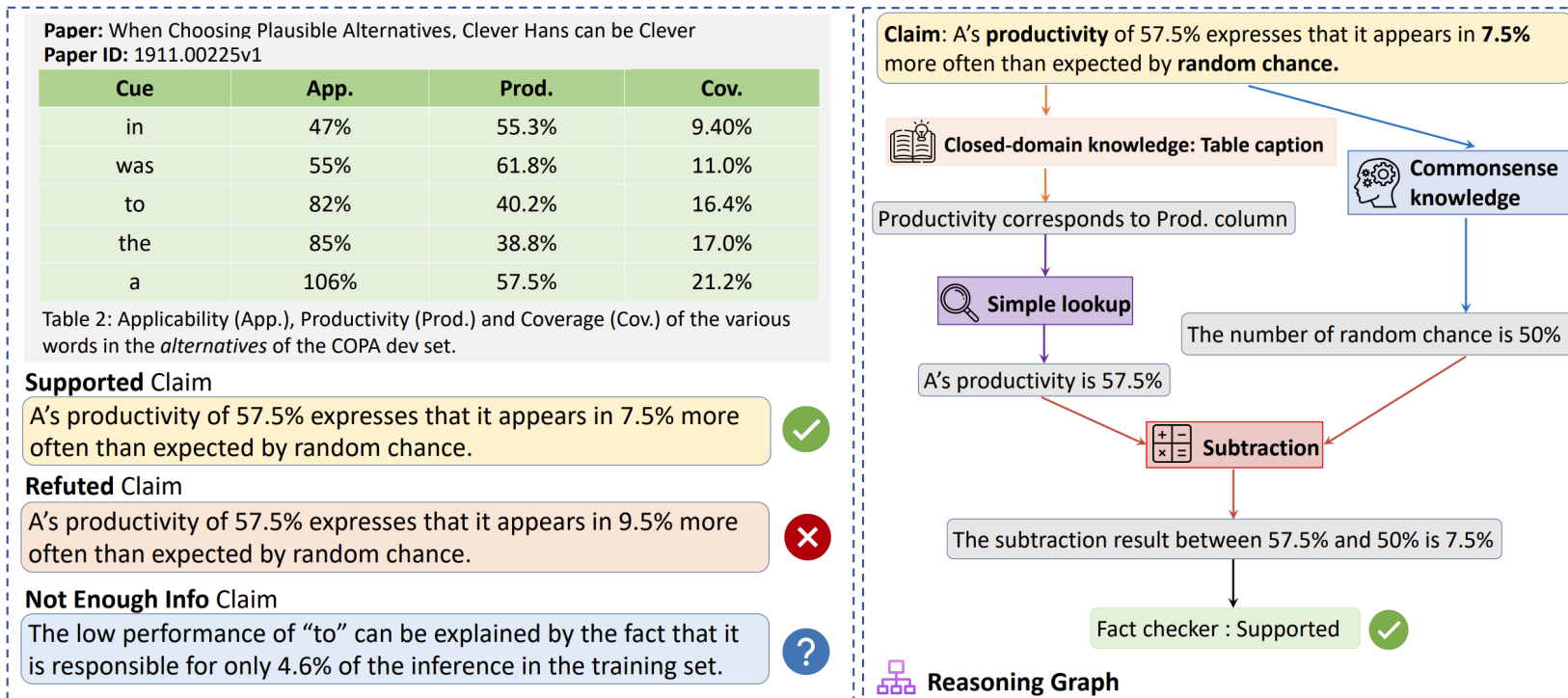- The model suffers from error propagation from IE model

Magnusson, I. H., & Friedman, S. E. (2021). Extracting fine-grained knowledge graphs of scientific claims: Dataset and transformer-based results. EMNLP 2021.
Wu, J., Chao, W., Zhou, X., & Luo, Z. (2023, December). Characterizing and Verifying Scientific Claims: Qualitative Causal Structure is All You Need. EMNLP 2023.

Scientific Claim Verification with Paper Text

Scientific Claim Verification with Scientific Tables

Open-Domain Scientific Claim Verification

Human-centered Fact-Checking

- Scientific claims are inherently linked to the experimental data, which are often represented in tables and figures

- Task
  - Given a claim $C$ and (table + caption) $T$, a table fact-checking model $F$ predicts a label $Y$ to verify whether $C$ is supported, refuted, or does not have enough information (NEI) to be verified by the information in $T$



**Paper:** When Choosing Plausible Alternatives, Clever Hans can be Clever
**Paper ID:** 1911.00225v1

| Cue | App. | Prod. | Cov. |
|---|---|---|---|
| in | 47% | 55.3% | 9.40% |
| was | 55% | 61.8% | 11.0% |
| to | 82% | 40.2% | 16.4% |
| the | 85% | 38.8% | 17.0% |
| a | 106% | 57.5% | 21.2% |

Table 2: Applicability (App.), Productivity (Prod.) and Coverage (Cov.) of the various words in the *alternatives* of the COPA dev set.

**Supported** Claim
A's productivity of 57.5% expresses that it appears in 7.5% more often than expected by random chance. ✅

**Refuted** Claim
A's productivity of 57.5% expresses that it appears in 9.5% more often than expected by random chance. ❌

**Not Enough Info** Claim
The low performance of "to" can be explained by the fact that it is responsible for only 4.6% of the inference in the training set. ❓

**Claim**: A's **productivity** of 57.5% expresses that it appears in **7.5%** more often than expected by **random chance.**

📖 **Closed-domain knowledge: Table caption**

Productivity corresponds to Prod. column

⚙️ **Commonsense knowledge**

🔍 **Simple lookup**

The number of random chance is 50%

A's productivity is 57.5%

🔢 **Subtraction**

The subtraction result between 57.5% and 50% is 7.5%

Fact checker : Supported ✅

**Reasoning Graph**

Lu, X., Pan, L., Liu, Q., Nakov, P., & Kan, M. Y. (2023). Scitab: A challenging benchmark for compositional reasoning and claim verification on scientific tables. EMNLP 2023.
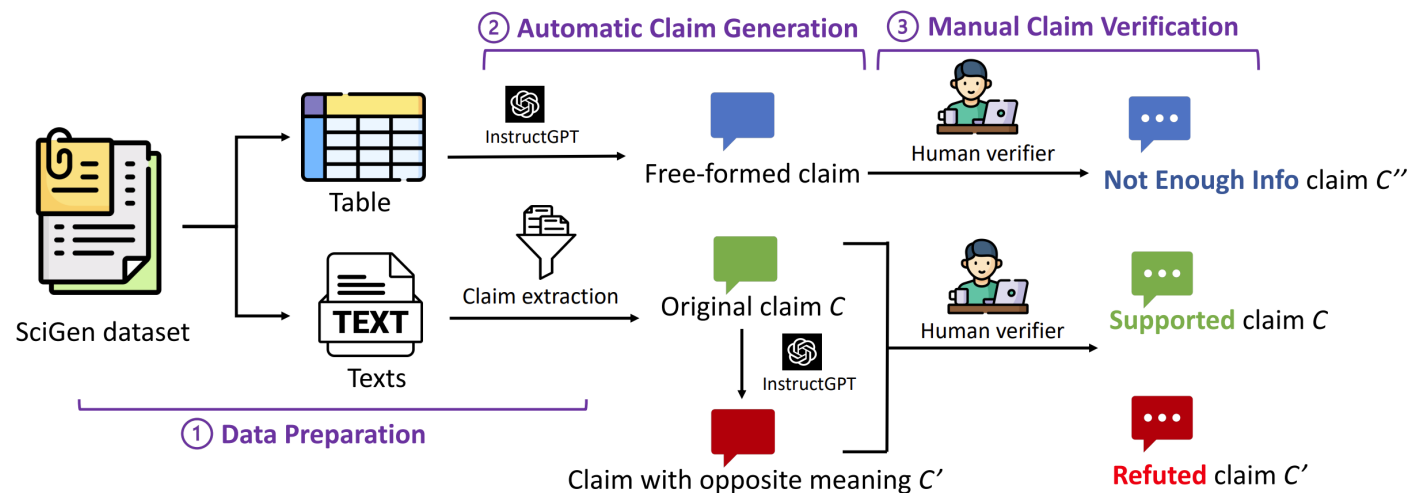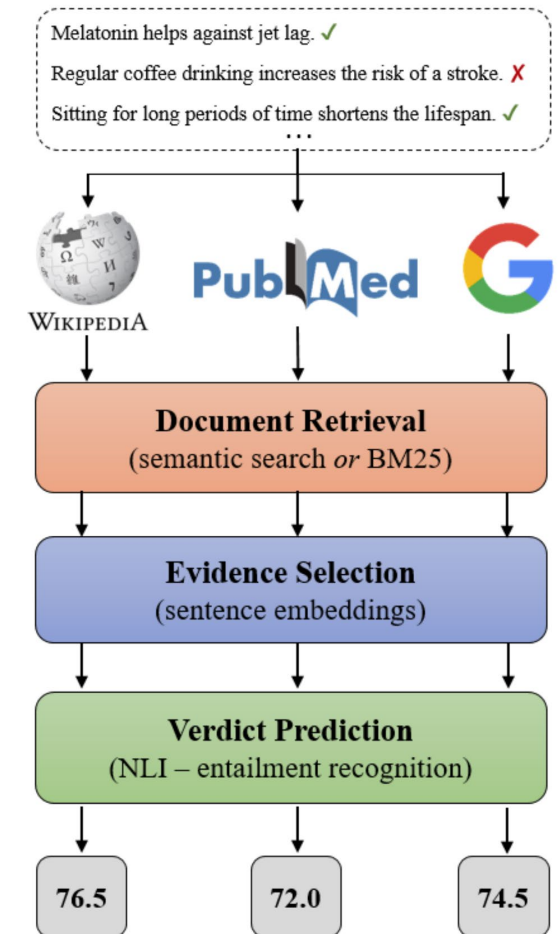
- Takeaways
  - Existing LLMs except GPT4 still have a large gap from human performance
  - Table-based LLMs do not outperform models pre-trained on pure texts
  - The results in the 3-class setting are notably poorer than those in the 2-class setting (exclude NEI)
  - The provision of in-context examples does not result in improved performance for the majority of models

Future Directions:
- Pretraining LLMs on scientific table datasets to test whether they can outperform pretrained LLMs on pure text
- Incorporate visual features to help models understand tables better



Lu, X., Pan, L., Liu, Q., Nakov, P., & Kan, M. Y. (2023). Scitab: A challenging benchmark for compositional reasoning and claim verification on scientific tables. EMNLP 2023.

Scientific Claim Verification with Paper Text

Scientific Claim Verification with Scientific Tables

Open-Domain Scientific Claim Verification

Human-centered Fact-Checking

# Comparing Knowledge Sources

- Compare three different knowledge sources (PubMed, Wikipedia, Google Search) with two retrieval techniques (BM25 and BioSimCSE) on biomedical fact-checking
  - Searching for evidence in the open domain have similar results compared to the closed-domain setting
  - The knowledge source perform comparably, with Wikipedia being better for popular and trending claims and PubMed for technical inquiries
  - Semantic search generally demonstrates superior recall
  - BM25 performs better than semantic search when retrieving from PubMed due to certain medical entities

Vladika, J., & Matthes, F. (2024). Comparing Knowledge Sources for Open-Domain Scientific Claim Verification. EACL 2024.

Scientific Claim Verification with Paper Text

Scientific Claim Verification with Scientific Tables

Open-Domain Scientific Claim Verification

Human-centered Fact-Checking

# Compare NLP fact-checking with professional fact-checkers

- **Human Fact Checking**
  - Global counter-evidence (GCE): Counterevidence via arbitrarily complex reasoning but without the source guarantee
  - Local counter-evidence (LCE): Evidence requires the source guarantee to refute the (reasoning behind) claim
  - Non-credible source (NCS): Evidence requires the source guarantee to refute the claim based on non-credible sources (e.g. satire)
  - No evidence assertion (NEA): The claim is refuted as no (trusted) evidence supports it

**Future Directions:**
- Propose methods to remove leaked evidence in each datasets

- **NLP Fact Verification**
  - Focusing on evidence-based approaches
  - Assume access of relevant (counter) evidence

- **Current FCNLP fails to provide source guarantees**
  - NLP fact-checking datasets might contain **insufficient** or **leaked evidence**

- **Models learn to rely on leaked evidence**

- **The assumption of the existence of counter-evidence in NLP Fact Checking is unrealistic and does not reflect real-world requirements**

Glockner, M., Hou, Y., & Gurevych, I. (2022). Missing counter-evidence renders nlp fact-checking unrealistic for misinformation. EMNLP 2022

# Future Directions for Scientific Fact Checking

- Propose new multimedia fact checking tasks based on scientific figures/equations

- Develop a joint framework to extract casual graph and to conduct fact checking in an end-to-end way

- Enhance reasoning capabilities by enabling models to decompose scientific claims effectively

- Background and Motivation

- Scientific Literature Survey

- Hypothesis Generation and Experiments

- Hands-on Paper Draft Assistant

- Paper Writing

- **Paper Draft Evaluation and Ethics**
  - Automatic Scientific Reviewing
  - Scientific Fact-Checking
  - **Ethics Concern in LLM Era**

- Summary and Future Directions

- Factuality, misinformation, and hallucination

- Plagiarism

- High quantities of (often low-quality) work

- Issues with peer review


- These are issues <u>now</u>! The cat is out of the bag.

# Factuality

- **LLMs like ChatGPT tend to hallucinate.**
  - The only time this is good for science is hypothesis generation! When we read the literature or write papers, we don't want to make things up (usually).
  - ChatGPT often cites non-existent references. The proportion of nonexistent citations was over 30% with GPT-3.5, and it is over 20% with GPT-4 several months later.
  - ChatGPT generated fake court citations

## Publisher blacklists authors after preprint cites made-up studies

Last month, a millipede expert in Denmark received an email notifying him that one of his publications had been mentioned in a new manuscript on *Preprints.org*. But when the researcher, Henrik Enghoff, downloaded the paper, he learned that it cited his work for something off-topic.

Henrik Enghoff

## New York lawyers sanctioned for using fake ChatGPT cases in legal brief

By **Sara Merken**

June 26, 2023 3:28 AM CDT · Updated 10 months ago

Sponsored by **intel**   *Sponsors are not involved in the creation of this or any other Reuters news articles*

A response by ChatGPT, an AI chatbot developed by OpenAI, is seen on its website in this illustration picture taken February 9, 2023. REUTERS/Florence Lo/Illustration *Purchase Licensing Rights* ⧉

Buchanan, J., Hill, S., & Shapoval, O. (2024). ChatGPT Hallucinates Non-existent Citations: Evidence from Economics. The American Economist, 69(1), 80-87
https://retractionwatch.com/2023/07/07/publisher-blacklists-authors-after-preprint-cites-made-up-studies/
https://www.reuters.com/legal/new-york-lawyers-sanctioned-using-fake-chatgpt-cases-legal-brief-2023-06-22/

- Both supervised fine-tuning (SFT) and reinforcement learning (RL) can encourage hallucination

- Factuality-aware Alignment
  - **Supervised Fine-Tuning (SFT) Adjustments:** Use the LLM's own existing knowledge to generate training data, thus minimizing the incorporation of unfamiliar information
  - **Reinforcement Learning (RL) Tweaks:** The model's reward function is adjusted to include a direct preference for factuality using a factuality reward model (SFT)
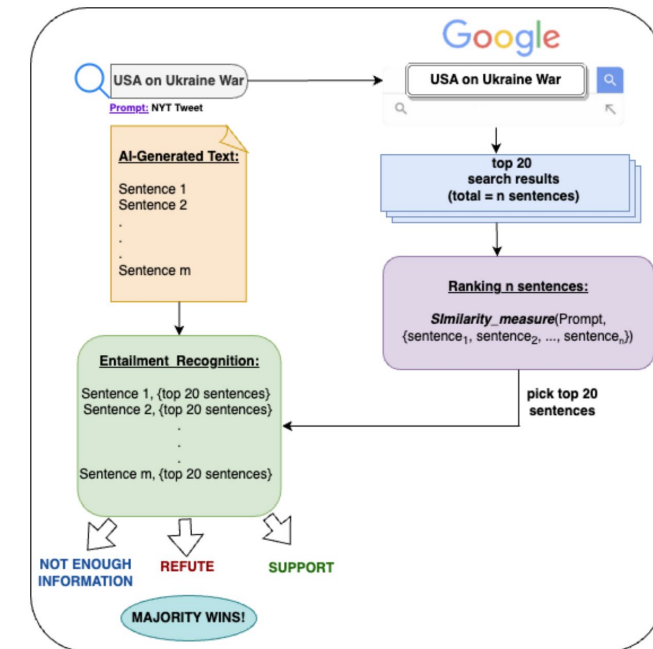


(a) Factuality-Aware SFT — Source of supervision
$x \in X^{\text{fact}}$ — Please give me a brief history of coffee. → SFT → Responses from itself with few-shot: $\text{PT}(x)$
$x \notin X^{\text{fact}}$ — Tell me a story about a pig who goes to the moon. → SFT → Human responses

(b) Factuality-Aware DPO — Source of supervision
$x \in X^{\text{fact}}$ — Do you have any information about Mac OS? → DPO → Factual or not $y_{\text{true}}$ / $y_{\text{false}}$
$x \notin X^{\text{fact}}$ — Tell me a story about a dog who goes to the mars. → DPO → 👍 or 👎 $y_+$ / $y_-$

**Future Directions:**
- Incorporate claim generation modules from other systems and external KBs might help LLMs verify factuality better

Lin, S. C., Gao, L., Oguz, B., Xiong, W., Lin, J., Yih, W. T., & Chen, X. (2024). FLAME: Factuality-Aware Alignment for Large Language Models. arXiv preprint arXiv:2405.01525.

# Mitigation Methods: After Training

- **High Entropy Word Spotting and Replacement:** Identify high-entropy words in hallucinated text and replace them with predictions from another LLM

- **Sentence Factuality Checking:** Identify sentences that are considered suspicious by NLI models for human review

Rawte, V., Chakraborty, S.,... & Das, A. (2023). The Troubling Emergence of Hallucination in Large Language Models--An Extensive Definition, Quantification, and Prescriptive Remediations. EMNLP 2023.

# Privacy

- LLMs can leak private information
  - Private details like names and contact details can be extracted from large neural models

Nasr, M., Carlini, N., Hayase, J., Jagielski, M., Cooper, A. F., Ippolito, D., ... & Lee, K. (2023). Scalable extraction of training data from (production) language models. arXiv preprint arXiv:2311.17035.

# Quantity overwhelms existing systems

- Papers can be generated much quicker now, but they aren't necessarily good
  - Current scientific productivity metrics are measured by quantity
  - "When a measure becomes a target, it ceases to be a good measure."

- Using "As of my last knowledge update" on Google Scholar returns 100s of results!

### ChatGPT on the Gravitational Redshift
J Stávek - European Journal of Applied Physics, 2024 - ej-physics.org

… **As of my last knowledge update** in January 2022, I'm not aware of any widely accepted alternative scenarios or theories that challenge Einstein's interpretation of gravitational redshift in …

☆ Save  ⠛⠛ Cite   Related articles   All 2 versions   »

### An Evaluation of ChatGPT and Bard in the Context of Biological Knowledge Retrieval
R Caspi, PD Karp - Access Microbiology, 2024 - microbiologyresearch.org

… A: **As of my last knowledge update** in September 2021, the term "rbcX" was not commonly associated 345 with a specific protein or gene in the context of photosynthesis or carbon …

☆ Save  ⠛⠛ Cite   Related articles   All 2 versions

## Radiology Case Reports
Volume 19, Issue 6, June 2024, Pages 2106-2111

Case Report

## Successful management of an Iatrogenic portal vein and hepatic artery injury in a 4-month-old female patient: A case report and literature review

Raneem Bader MD [a], Ashraf Imam MD [b], Mohammad Alnees MD [a] [e] 👤 ✉, Neta Adler MD [c], Joanthan ilia MD [c], Diaa Zugayar MD [b], Arbell Dan MD [d], Abed Khalaileh MD [b] 👤 ✉

Show more ⌄

+ Add to Mendeley   ⠕ Share   🙿 Cite

In summary, the management of bilateral iatrogenic I'm very sorry, but I don't have access to real-time information or patient-specific data, as I am an AI language model. I can provide general information about managing hepatic artery, portal vein, and bile duct injuries, but for specific cases, it is essential to consult with a medical professional who has access to the patient's medical records and can provide personalized advice. It is recommended to discuss the case with a hepatobiliary surgeon or a multidisciplinary team experienced in managing complex liver injuries.

"Scientific Journals Are Publishing Papers With AI-Generated Text" 404media.co. 2024.

- **Peer-review challenges**
  - Current scientific productivity metrics are measured by quantity
    - Measuring quality is hard to do without retrospect!
  - There aren't enough people to do free peer-review
    - Peer review takes time away from writing more papers
    - Many are turning to LLMs for help, but are the results thorough?
      - ChatGPT usage is correlated with low reviewer confidence



Liang, W., Izzo, Z., Zhang, Y., Lepp, H., Cao, H., Zhao, X., ... & Zou, J. Y. (2024). Monitoring AI-Modified Content at Scale: A Case Study on the Impact of ChatGPT on AI Conference Peer Reviews. arXiv preprint arXiv:2403.07183.
Van Noorden, R., & Perkel, J. M. (2023). AI and science: what 1,600 researchers think. Nature, 621(7980), 672-675.
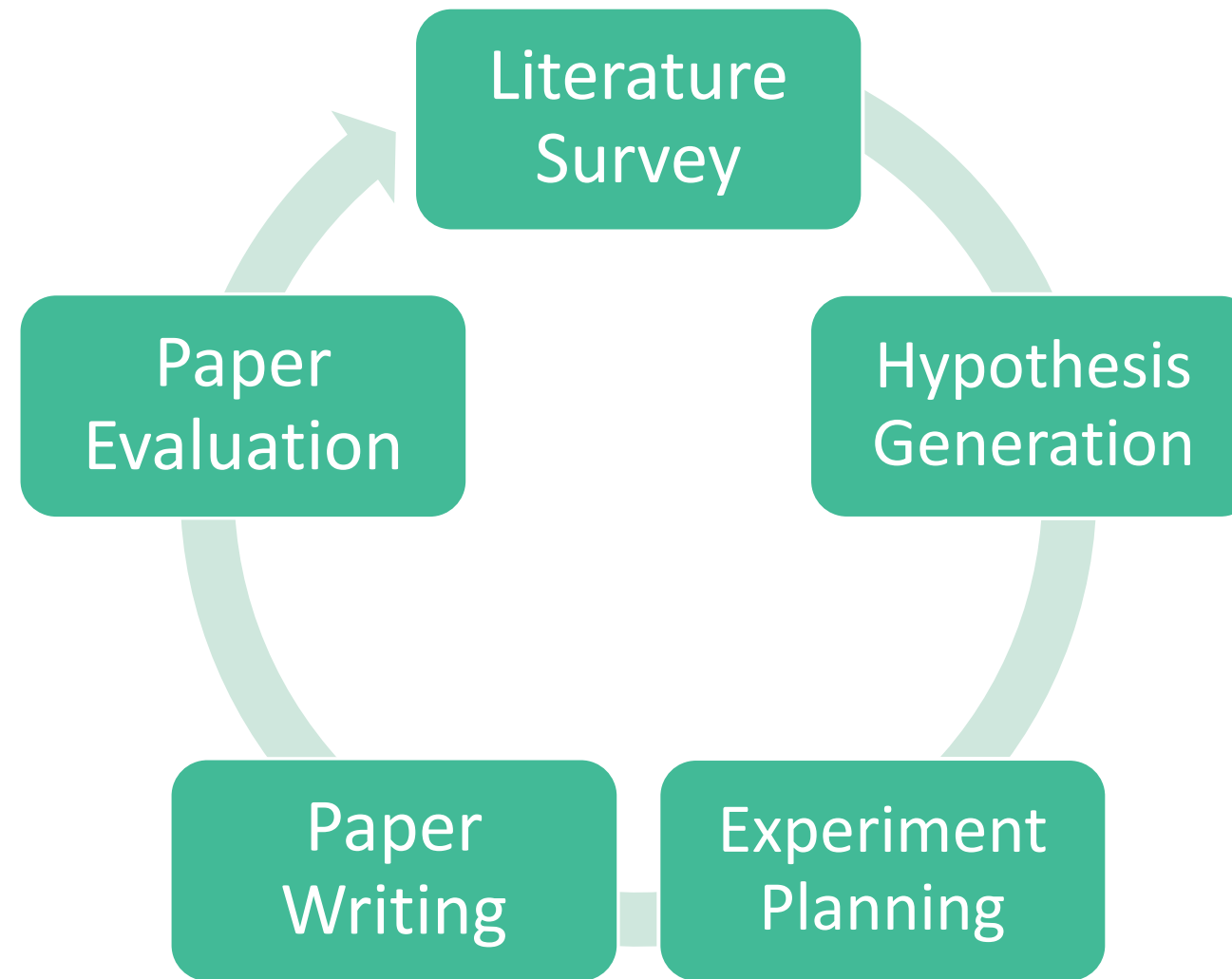
# What can we do? -- LLM Detectors

- **Watermarking text (Kirchenbauer et al., 2023)**
  - Encourage an unnatural distribution of generated words
    - Split the token set into a red list and a green list with a random seed.
    - Promote the use of green tokens, i.e., increase its probability

- **Zero-shot approaches**
  - Measure the perplexity and complexity of a generated sentence

- **Retrieval (Krishna et al., 2023)**
  - Given candidate text, search a database of previous LLM generations, looking for sequences that match the candidate text within a certain threshold

- **Neural Network detectors for AI-generated text**
  - Train a model to distinguish between human and AI text
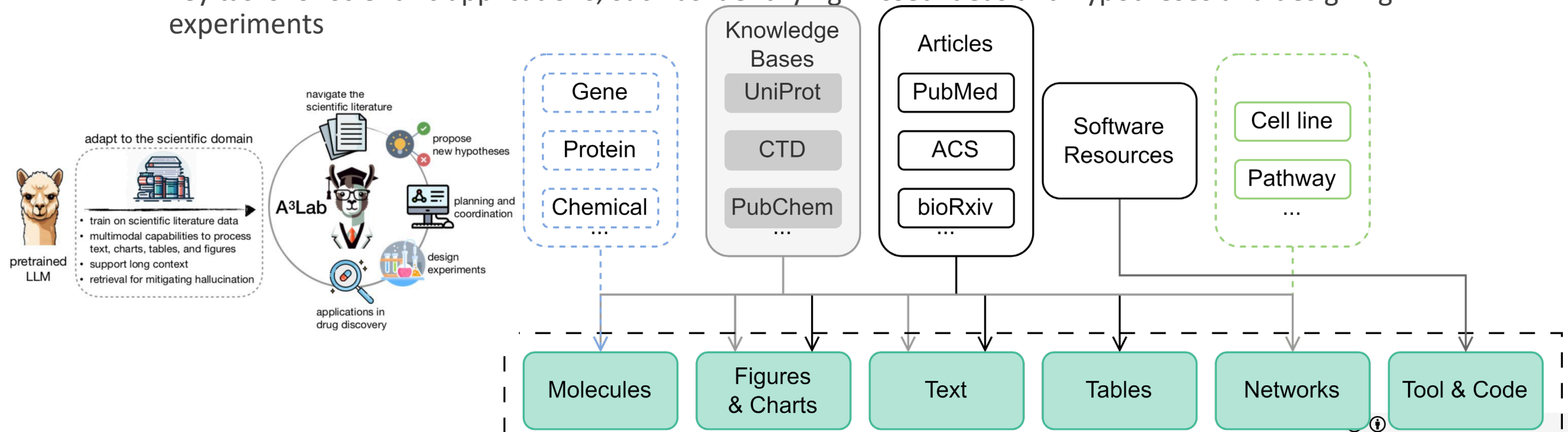  - These go out of date quickly and can be fooled by stronger models

No watermark
Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words)
Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.999999999% of the Synthetic Internet

With watermark
- minimal marginal probability for a detection attempt.
- Good speech frequency and energy rate reduction.
- messages indiscernible to humans.
- easy for humans to verify.

Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., & Goldstein, T. (2023, July). A watermark for large language models. In International Conference on Machine Learning (pp. 17061-17084). PMLR.
https://www.techlearning.com/news/what-is-gptzero-the-chatgpt-detection-tool-explained
Krishna, K., Song, Y., Karpinska, M., Wieting, J., & Iyyer, M. (2024). Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense. Advances in Neural Information Processing Systems, 36.
https://huggingface.co/roberta-large-openai-detector

# Conclusion and Future Directions

# Tutorial Summary

- The scientific lifecycle will be increasingly affected by NLP technologies, especially as they mature
  - AI-driven research assistants benefit from external knowledge
  - Multimodal and human-in-the-loop AI remain largely overlooked in current systems
- NLP technologies are both the solution and the problem
  - They exacerbate problems in existing systems (peer review, publication metrics)
  - But they will be required for helping to fix these issues
- Existing systems can already cover most tasks involved in scientific paper lifecycle



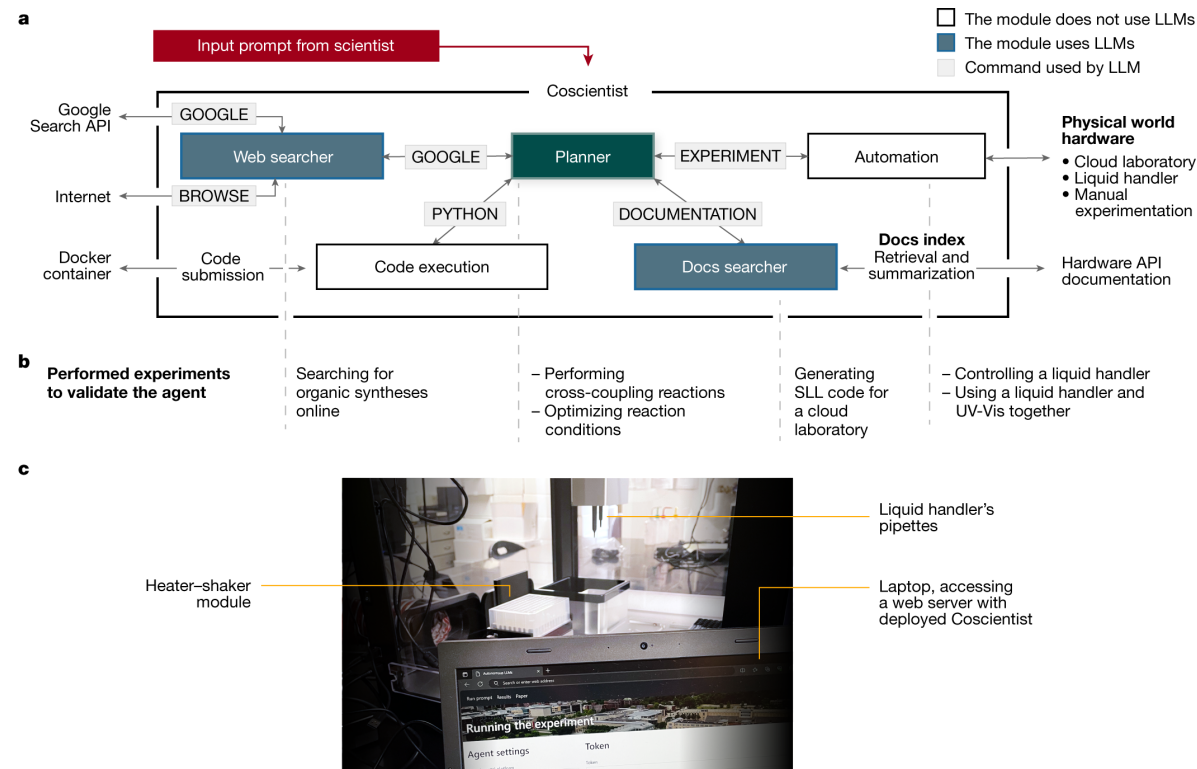Literature Survey → Hypothesis Generation → Experiment Planning → Paper Writing → Paper Evaluation → (cycle)

- Existing work chooses either highly-capable general domain models (e.g., GPT-4) with strong reasoning capabilities or trains specific scientific language models (e.g., SciBERT) with specific knowledge.
    - Can we more effectively bridge this gap?

- Develop more capable multimodal LLMs which can jointly analyze code, scientific modalities, tables, charts, figures, pathways, and other visual data alongside textual content in research papers, technical reports and textbooks

- Explore instruction finetuning on a suite of datasets carefully designed to improve LLM capabilities in key tasks for scientific applications, such as identifying missed ideas and hypotheses and designing experiments
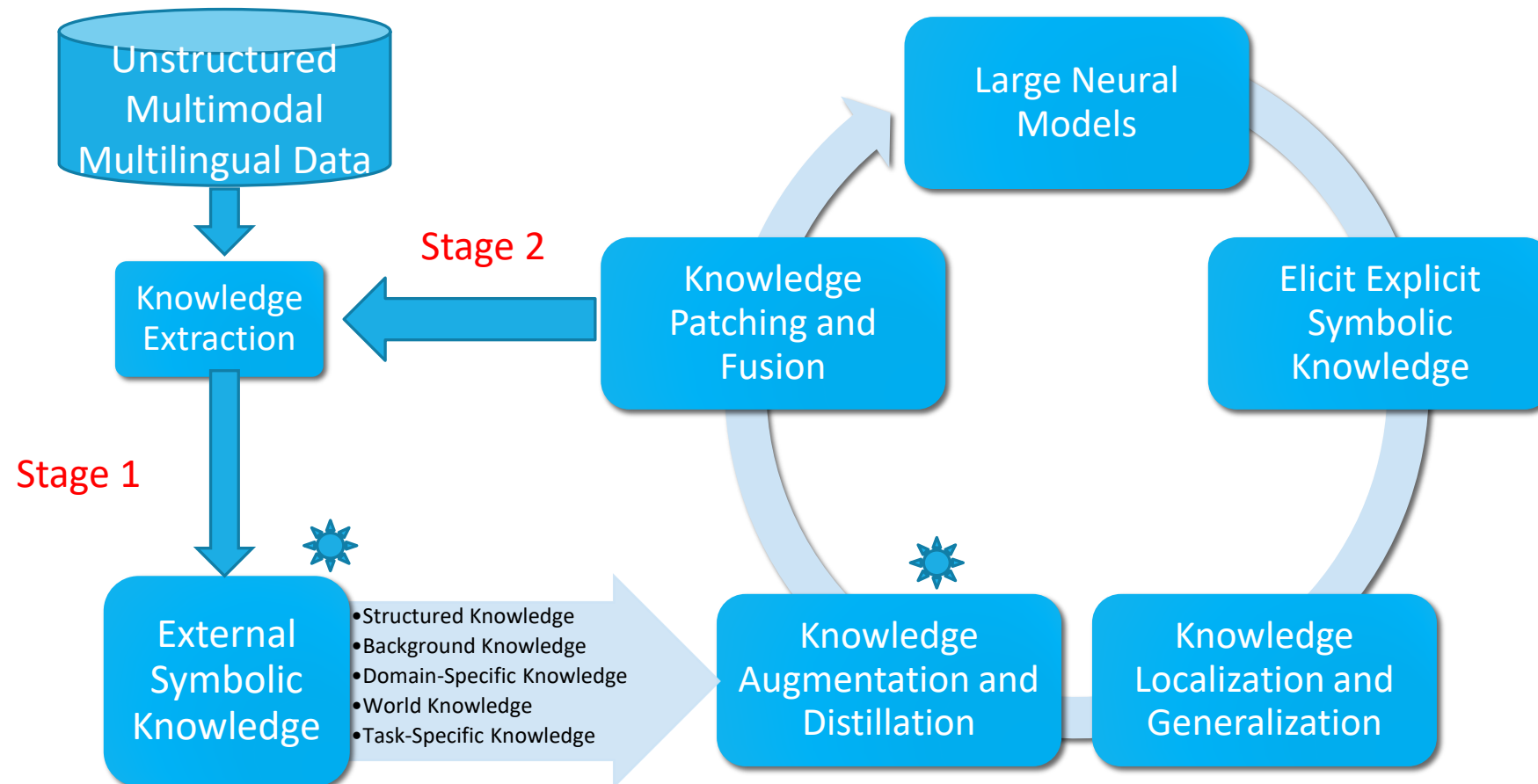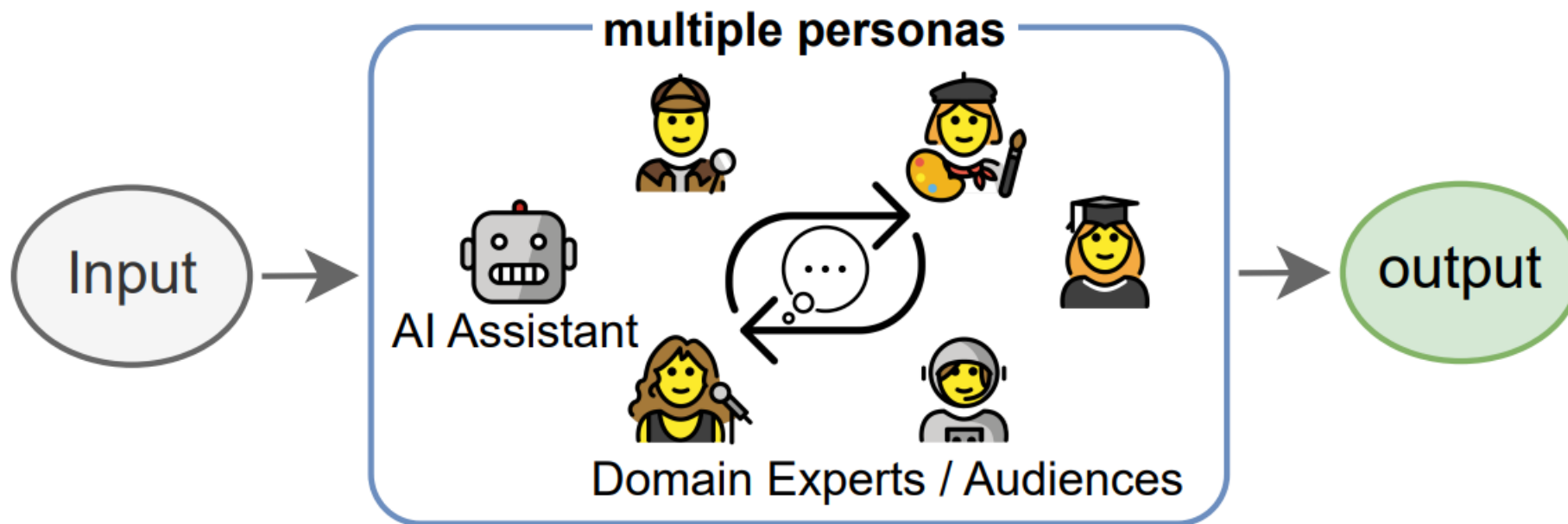
- **Convert automated agents for scientific discovery from prototype ideas into useful tools**
  - Read the entire scientific literature and synthesize relevant information to propose novel ideas
  - Couple with external simulations for rapid feedback
  - Integrate with specially designed laboratory hardware for automated experimentation and testing



Boiko, D. A., MacKnight, R., Kline, B., & Gomes, G. (2023). Autonomous chemical research with large language models. Nature, 624(7992), 570-578.

- Create a two-way knowledge acquisition framework to make symbolic and neural learning approaches mutually enhance each other
  - In the first stage, elicit and acquire explicit symbolic knowledge from large neural models
  - In the second stage, leverage the acquired symbolic knowledge to augment and enhance these large models.

- Use a multi-agent debate framework by engaging in multi-turn self-collaboration with LLM agents that have differing expertise

    - Prompt a set of LLMs with varied expertise and different sources of evidence to verify the complex claims via collaboration and simulation.

    - Each LLM can represent a particular domain expert, such as an agent that understands how to write molecular simulations, another for reading the literature, or another agent that specializes in synthesizing molecules.



Wang, Z., Mao, S., Wu, W., Ge, T., Wei, F., & Ji, H. (2023). Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaborationNAACL 2024.

# Thank you! Questions?

ALL tutorial slides and reading list are available at:

https://sites.google.com/view/coling2024-paper-lifecycle/

LREC-COLING 2024