# Scientific Language Models for Biomedical Knowledge Base Completion: An Empirical Study

**Rahul Nadkarni[1], David Wadden[1], Iz Beltagy[2], Noah A. Smith[1,2], Hannaneh Hajishirzi[1,2], Tom Hope[1,2]**

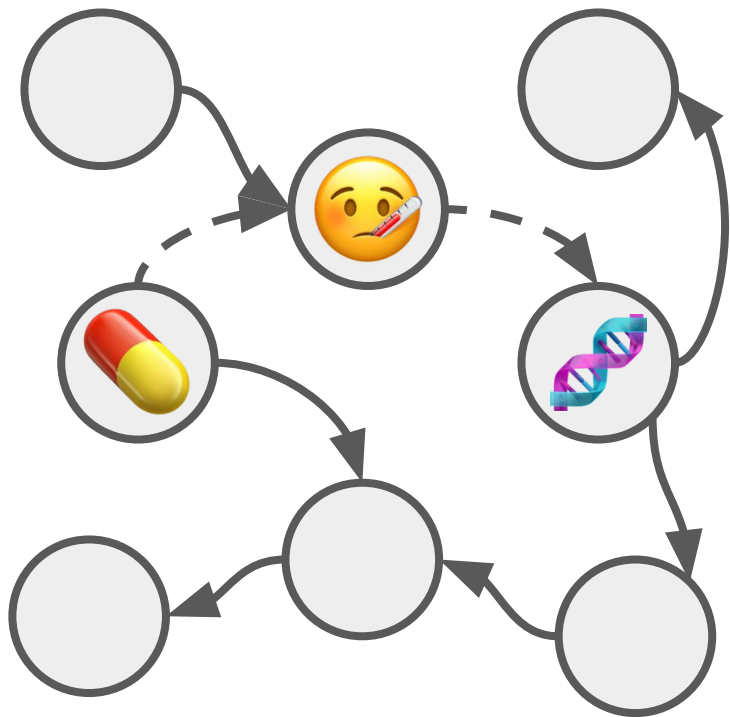[1] Paul G. Allen School of Computer Science & Engineering, University of Washington

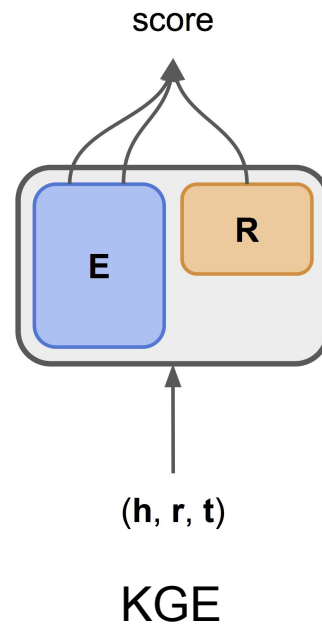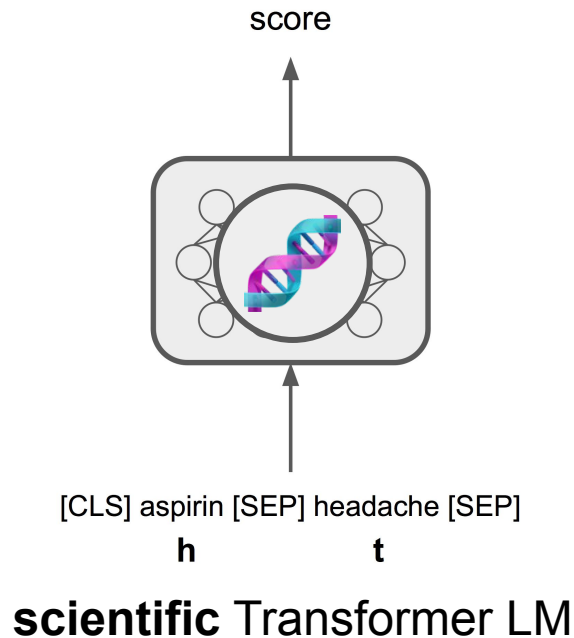[2] Allen Institute for Artificial Intelligence (AI2)

# Biomedical Knowledge Graph Completion



- Relations between entities
  - Repurposing drugs for diseases
  - Mapping diseases to genes

- Frame as biomedical knowledge graph completion

# LMs for Biomedical Knowledge Graph Completion



First to systematically apply scientific LMs for KG completion and compare to KGE models
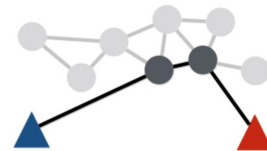
# Datasets



**RepoDB**

drugs, diseases

2.7k entities, 6.7k triples

**Hetionet**

drugs, diseases, genes, symptoms, side effects
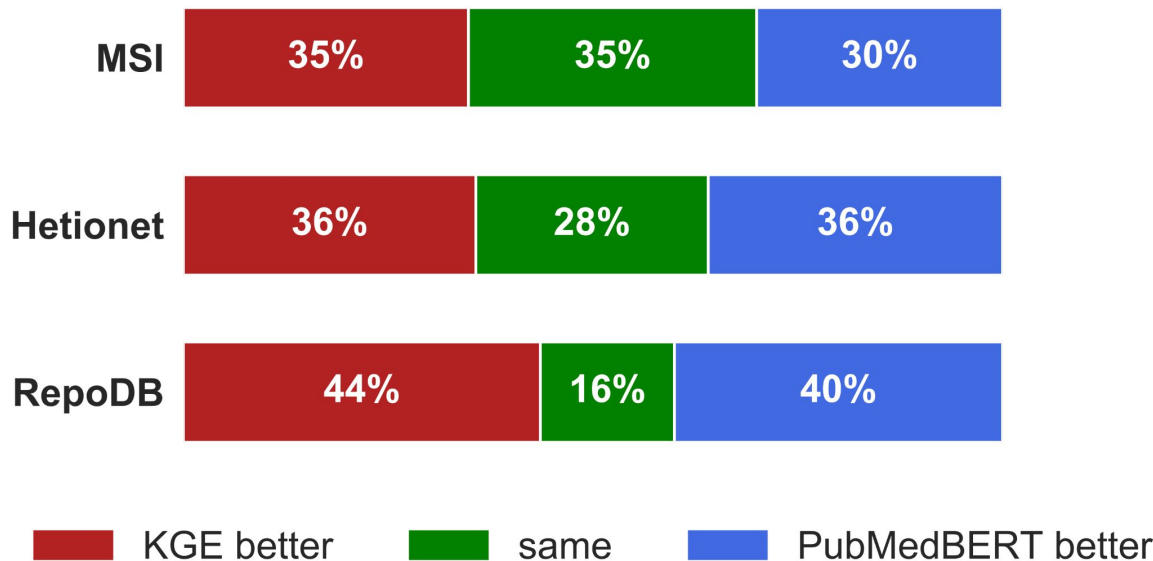
12.7k entities, 156k triples
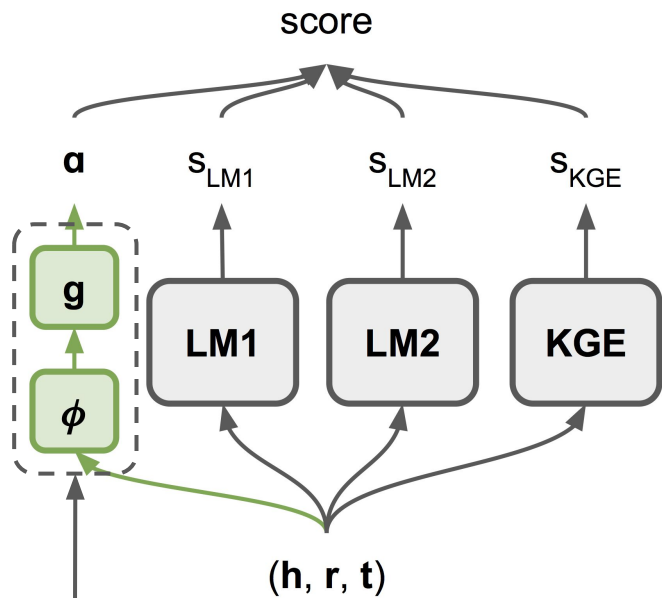
**MSI**

drugs, diseases, proteins, protein functions
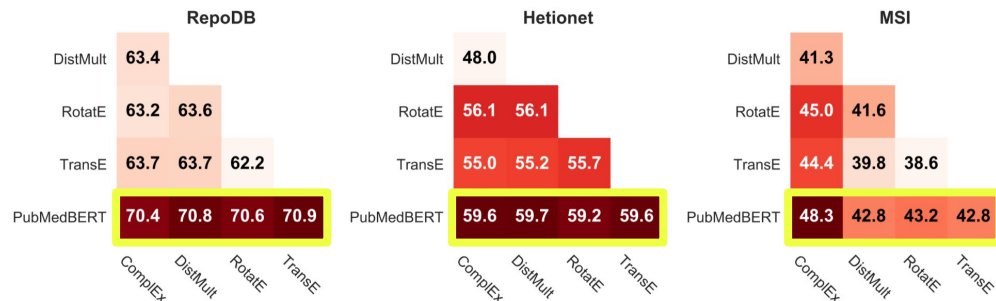
30k entities, 485k triples

# Relative Performance



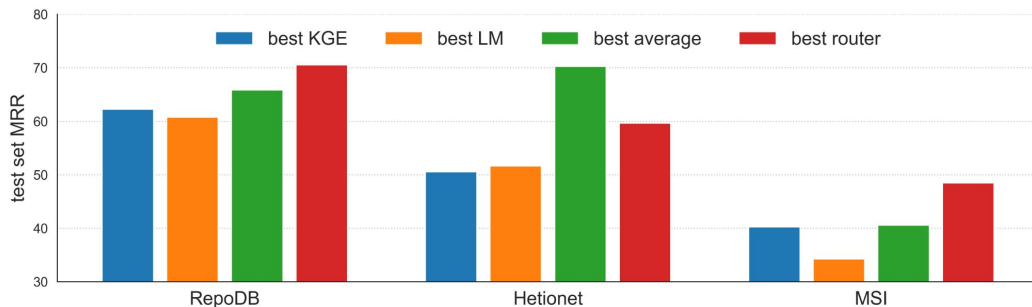LMs and KGE models perform well on different subsets of examples

# Integrating Models

score

$\alpha$  $s_{LM1}$  $s_{LM2}$  $s_{KGE}$

**g**  **LM1**  **LM2**  **KGE**

$\phi$

(**h**, **r**, **t**)

- weighted average
- router classifier



| RepoDB | ComplEx | DistMult | RotatE | TransE |
|---|---|---|---|---|
| DistMult | 63.4 | | | |
| RotatE | 63.2 | 63.6 | | |
| TransE | 63.7 | 63.7 | 62.2 | |
| PubMedBERT | 70.4 | 70.8 | 70.6 | 70.9 |

| Hetionet | ComplEx | DistMult | RotatE | TransE |
|---|---|---|---|---|
| DistMult | 48.0 | | | |
| RotatE | 56.1 | 56.1 | | |
| TransE | 55.0 | 55.2 | 55.7 | |
| PubMedBERT | 59.6 | 59.7 | 59.2 | 59.6 |

| MSI | ComplEx | DistMult | RotatE | TransE |
|---|---|---|---|---|
| DistMult | 41.3 | | | |
| RotatE | 45.0 | 41.6 | | |
| TransE | 44.4 | 39.8 | 38.6 | |
| PubMedBERT | 48.3 | 42.8 | 43.2 | 42.8 |

## Combinations with an LM perform better



best KGE   best LM   best average   best router
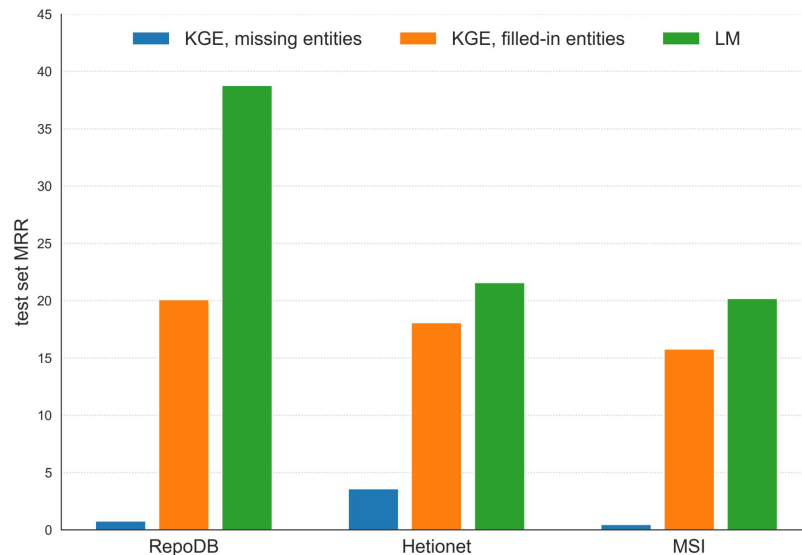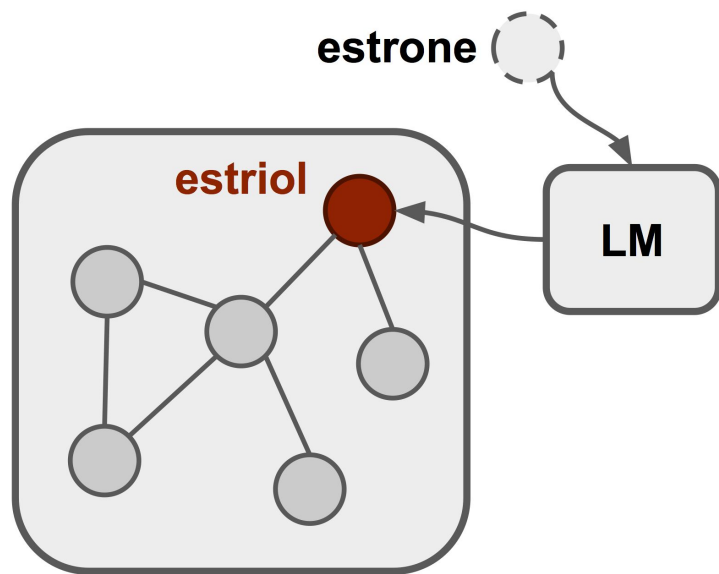
## Best router can outperform best weighted avg.

# Inductive Performance



LMs perform well (and can improve KGE performance) on unseen entities

For code and data, visit:

`github.com/rahuln/lm-bio-kgc`