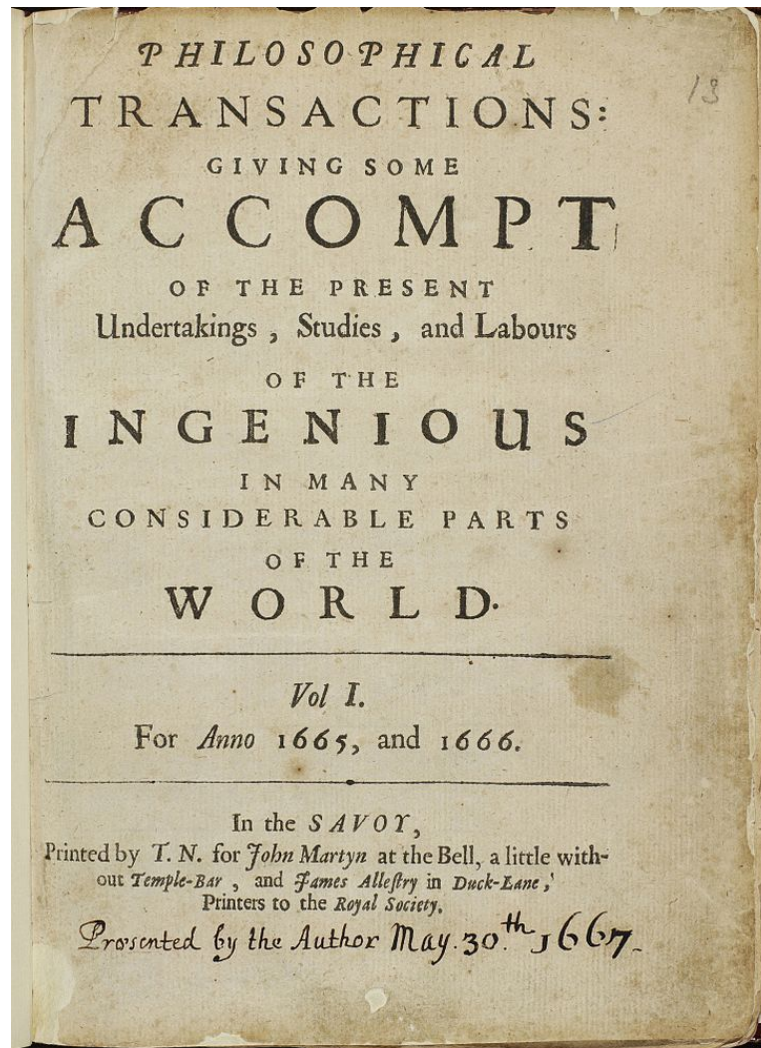# NLP, AKBC, and The Future of Science
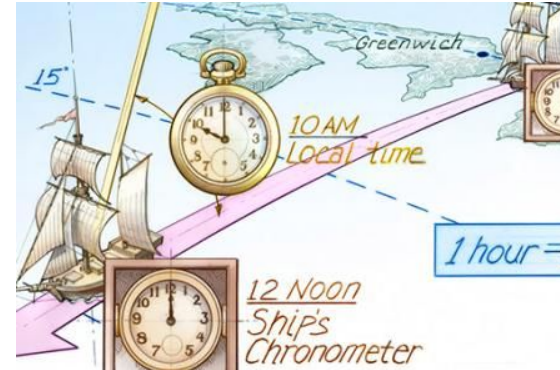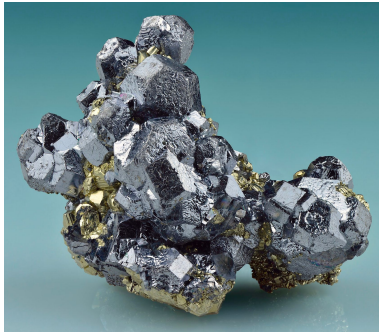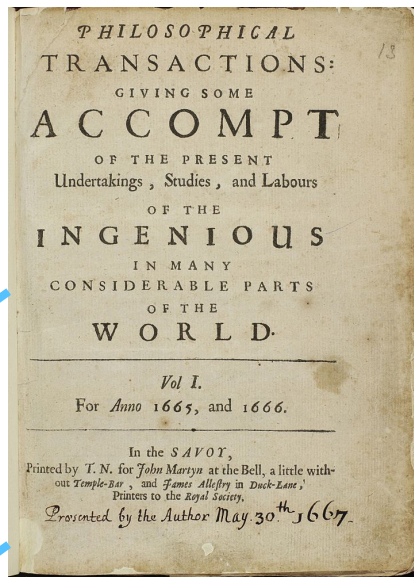
## Tom Hope

Assistant Professor, The Hebrew University of Jerusalem
Research Scientist, Allen Institute for AI (Semantic Scholar)

THE

# HISTORY

OF THE

## Royal-Society

OF

# LONDON,

For the Improving of

NATURAL KNOWLEDGE.

BY

*THO. SPRAT.*

1665: **first scientific journal** was published!

PHILOSOPHICAL
TRANSACTIONS:
GIVING SOME
ACCOMPT
OF THE PRESENT
Undertakings, Studies, and Labours
OF THE
INGENIOUS
IN MANY
CONSIDERABLE PARTS
OF THE
WORLD.

*Vol I.*
For *Anno* 1665, and 1666.

In the *SAVOY*,
Printed by T. N. for *John Martyn* at the Bell, a little without *Temple-Bar*, and *James Allestry* in *Duck-Lane*, Printers to the *Royal Society*.

*Presented by the Author* May. 30th 1667.

PHILOSOPHICAL
TRANSACTIONS:
GIVING SOME
ACCOMPT
OF THE PRESENT
Undertakings, Studies, and Labours
OF THE
INGENIOUS
IN MANY
CONSIDERABLE PARTS
OF THE
WORLD.

Vol I.
For Anno 1665, and 1666.

In the SAVOY,
Printed by T. N. for John Martyn at the Bell, a little with-
out Temple-Bar, and James Allestry in Duck-Lane,
Printers to the Royal Society.

Presented by the Author May. 30.th 1667.

Greenwich
15°
10 AM
Local time
1 hour =
12 Noon
Ship's
Chronometer

# Scientific Journals, 2022...

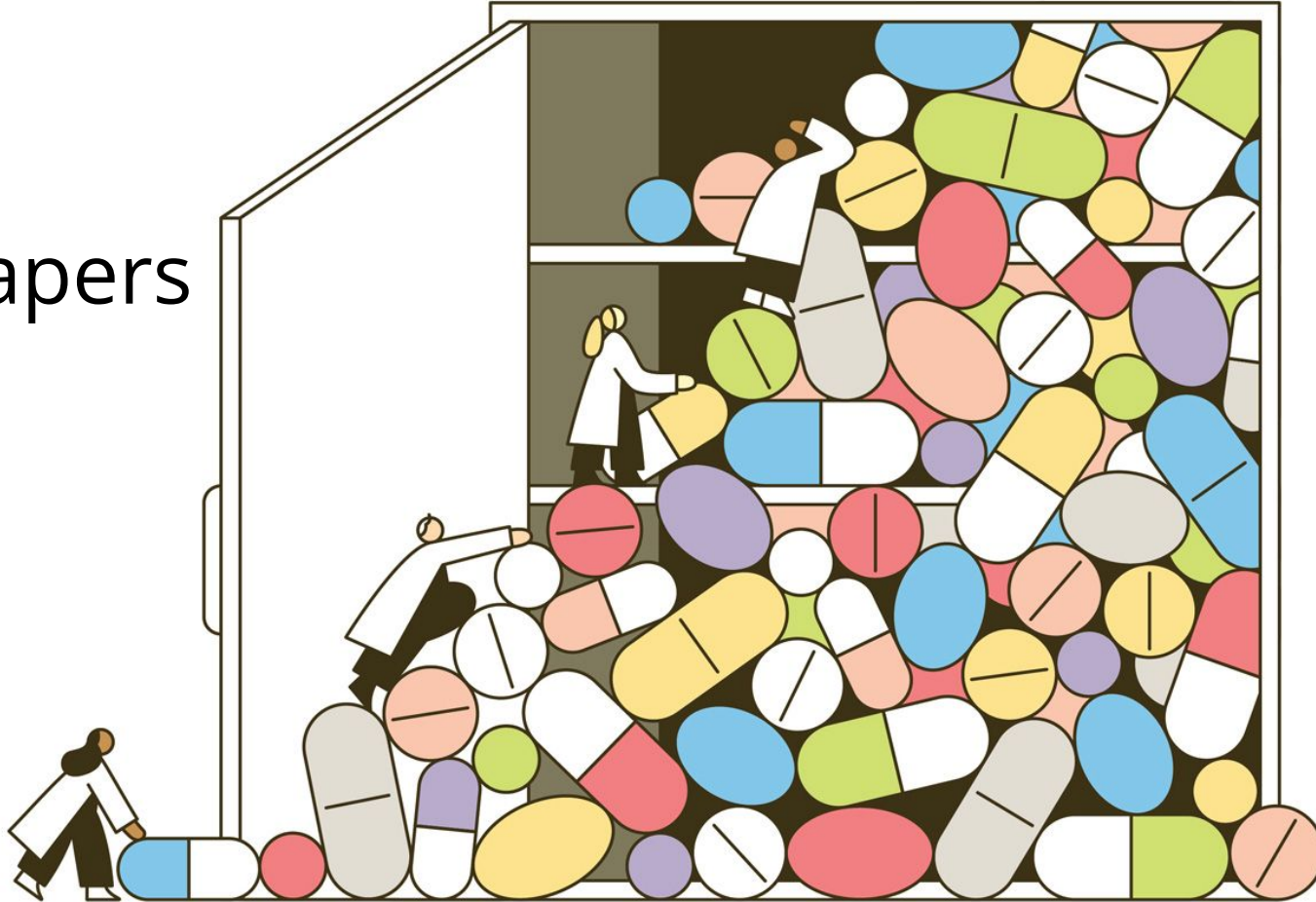**266,299,373**
Publications

**278,637,018**
Authors

**4,547**
Conferences

**49,050**
Journals ⓘ

# Example: Biomedical Literature

>1 million
biomedical papers
*every year*

# Explosion of Scientific Information

Literature

Scientific knowledge bases

Resources, code libraries

Online discussions

...

# Opportunity:
# Augment & Scale Scientific Discovery



AI systems that harness humanity's **collective scientific knowledge**

Research Tasks

**Research Tasks**

**Attention to areas of interest**
*(How do we keep track?)*

**Problem identification & prioritization**
*(How do we select what to work on?)*

Attention to areas of interest
*(How do we keep track?)*

*Research Tasks*

Problem identification & prioritization
(How do we select what to work on?)

**Forming directions**
*(How do we generate solutions?)*

Attention to areas of interest
(How do we keep track?)

*Research Tasks*

Problem identification & prioritization
*(How do we select what to work on?)*

Forming directions
*(How do we generate solutions?)*

Attention to areas of interest
*(How do we keep track?)*

Experimentation, analysis

*Research Tasks*

**Learning, understanding**

Problem identification & prioritization
*(How do we select what to work on?)*

Forming directions
*(How do we generate solutions?)*

Attention to areas of interest
*(How do we keep track?)*

Experimentation, analysis

*Research Tasks*

Learning, understanding

**Communication, collaboration**

# Scientific Knowledge: Challenges



Large-scale, diverse

Rapidly evolving

Deeply technical

$$\frac{\mathbf{z}_i^{(k+1)} - \mathbf{z}_i^{(k)}}{\tau} = \sum_{j:(i,j)\in\mathcal{E}(\mathbf{U}^{(k)})} a\left(\mathbf{z}_i^{(k)}, \mathbf{z}_j^{(k)}\right)\left(\mathbf{z}_j^{(k)} - \mathbf{z}_i^{(k)}\right). \qquad (9)$$

Here $k$ denotes the discrete time index (iteration) and $\tau$ is the time step (discretisation parameter). Rewriting (9) compactly in matrix-vector form with $\tau = 1$ leads to the *explicit Euler scheme*:

$$\mathbf{Z}^{(k+1)} = (\mathbf{A}^{(k)} - \mathbf{I})\mathbf{Z}^{(k)} = \mathbf{Q}^{(k)}\mathbf{Z}^{(k)}, \qquad (10)$$

where $a_{ij}^{(k)} = a(\mathbf{z}_i^{(k)}, \mathbf{z}_j^{(k)})$ and the matrix $\mathbf{Q}^{(k)}$ (diffusion operator) is given by

$$q_{ij}^{(k)} = \begin{cases} 1 - \tau \sum\limits_{l:(i,l)\in\mathcal{E}} a_{il}^{(k)} & i = j \\ \tau a_{ij}^{(k)} & (i,j) \in \mathcal{E}(\mathbf{U}^{(k)}) \\ 0 & \text{otherwise} \end{cases}$$
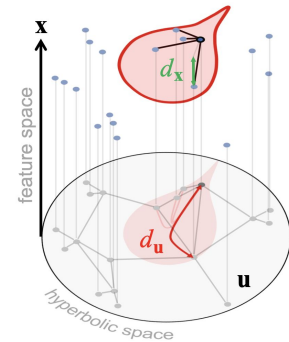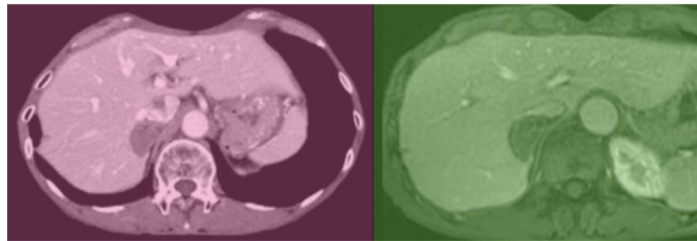




Figure 2: Graph Beltrami flow with

# Limited Modes of Interaction

# NLP for Science

# NLP for Science



**Extraction**
Organizing
the world's
scientific
knowledge

# NLP for Science



Literature Retrieval

🔍 patient desc. + in-hospital mortality

Severe hypoglycemia… **not associated** with increased risk of mortality in adults with Type 1 diabetes…

Extraction

**Retrieval**

Finding information to boost research, decision-making

# NLP for Science



NAACL 2022

**Extraction**

**Literature Retrieval**

patient desc. + in-hospital mortality

Severe hypoglycemia…
**not associated** with
increased risk of
mortality in adults with
Type 1 diabetes…

**Retrieval**

**Inference**
Making
predictions,
generating
hypotheses

# NLP for Science



Extraction

Retrieval

Inference

**Interaction**
NLP-powered exploratory interfaces for science

CHI22
New Orleans, LA

EMNLP 2020, 2022

**Input:**
Researcher's papers

**Retrieved Knowledge:**
Researchers who inspire **novel directions**

Eduard Hovy

benchmark do...
chinese conv...
Finding Stru...

ARIEL–CMU SF...
neural archi...
character-le...

Byron C. Wal...

text categor...
document cat...
neural text ...

auroc
PIVET 's ana...
actor modeli...

CHI22
New Orleans, LA

**Input**:
Items of interest
(*concepts, keywords...*)

Concept1   Concept2   add more...

**Retrieved Knowledge**:
**Groups** of researchers
and the **links** between them



EMNLP 2020

**Input**: *Problem description*

**Retrieved Knowledge**: Related problems for **forming new perspectives**

**Input**:
Items of interest
(*concepts, keywords...*)

Concept1  Concept2  add more...

**Retrieved Knowledge**:
Specific mentions of
**challenges, gaps in knowledge**,
**directions and hypotheses**

What ***don't*** we know in a specific area? What are interesting directions at **the edge of science**?

AAAI-22

Extracting KBs + search interfaces for discovering **causal mechanisms** and **drug combination treatments**

# In Today's Talk:



**Extraction**

**Retrieval**

**Inference**

Interaction

# In Today's Talk:

1.  Hierarchical Cross-Document Coreference

2.  Document Similarity & Retrieval

3.  Literature-Augmented Predictions

# Motivation: Author Matching



**Input:** - - - →
Researcher's
papers

**Retrieved Knowledge:**
Researchers who
inspire **novel directions**

# Authors who work on related tasks or use similar methods

## Dan Weld

**PDDL-the planning domain definition language**
D. McDermott, M. Ghallab, +5 authors  D. Wilkins · Computer Science · 1998
TLDR  This manual describes the syntax of PDDL, the Planning Domain Definition Language, the problem-specification language for the AIPS-98 planning competition, and hopes to encourage empirical evaluation of planner performance, and development of standard sets of problems all in comparable notations. Expand

**TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension**
Mandar Joshi, Eunsol Choi, Daniel S. Weld, Luke Zettlemoyer · Computer Science · ACL · 1 May 2017
TLDR  It is shown that, in comparison to other recently introduced large-scale datasets, TriviaQA has relatively complex, compositional questions, has considerable syntactic and lexical variability between questions and corresponding answer-evidence sentences, and requires more cross sentence reasoning to find answers. Expand

**SpanBERT: Improving Pre-training by Representing and Predicting Spans**
Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, Omer Levy · Computer Science · TACL · 24 July 2019
TLDR  The approach extends BERT by masking contiguous random spans, rather than random tokens, and training the span boundary representations to predict the entire content of the masked span, without relying on the individual token representations within it. Expand

**Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations**
R. Hoffmann, Congle Zhang, Xiao Ling, Luke Zettlemoyer, Daniel S. Weld · Computer Science · ACL · 19 June 2011
TLDR  A novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts is presented. Expand

**Unsupervised named-entity extraction from the Web: An experimental study**
Oren Etzioni, Michael J. Cafarella, +5 authors  A. Yates · Computer Science · Artif. Intell. · 1 June 2005
TLDR  An overview of KnowItAll's novel architecture and design principles is presented, emphasizing its distinctive ability to extract information without any hand-labeled training examples, and three distinct ways to address this challenge are presented and evaluated. Expand

**Fine-Grained Entity Recognition**
Xiao Ling, Daniel S. Weld · Computer Science · AAAI · 22 July 2012
TLDR  A fine-grained set of 112 tags is defined, the tagging problem is formulates as multi-class, multi-label classification, an unsupervised method for collecting training data is described, and the FIGER implementation is presented. Expand

## Ido Dagan

**QA-Align: Representing Cross-Text Content Overlap by Aligning Question-Answer Propositions**
Daniel Weiss, Paula Roit, Ayal Klein, Ori Ernst, Ido Dagan · Computer Science · 26 September 2021
Multi-text applications, such as multidocument summarization, are typically required to model redundancies across related texts. Current methods confronting consolidation struggle to fuse overlapping… Expand

**iFacetSum: Coreference-based Interactive Faceted Summarization for Multi-Document Exploration**
Eran Hirsch, Alon Eirew, +7 authors  Ido Dagan · Computer Science · ArXiv · 23 September 2021
We introduce iFACETSUM,1 a web application for exploring topical document sets. iFACETSUM integrates interactive summarization together with faceted search, by providing a novel faceted navigation… Expand

**Asking It All: Generating Contextualized Questions for any Semantic Role**
Valentina Pyatkin, Paul Roit, Julian Michael, Reut Tsarfaty, Yoav Goldberg, Ido Dagan · Computer Science · ArXiv · 10 September 2021
Asking questions about a situation is an inherent step towards understanding it. To this end, we introduce the task of role question generation, which, given a predicate mention and a passage,… Expand

**Realistic Evaluation Principles for Cross-document Coreference Resolution**
Arie Cattan, Alon Eirew, Gabriel Stanovsky, Mandar Joshi, Ido Dagan · Computer Science · STARSEM · 8 June 2021
TLDR  It is argued that models should not exploit the synthetic topic structure of the standard ECB+ dataset, forcing models to confront the lexical ambiguity challenge, as intended by the dataset creators. Expand

**Denoising Word Embeddings by Averaging in a Shared Space**
Avi Caciularu, Ido Dagan, J. Goldberger · Computer Science · STARSEM · 5 June 2021
TLDR  A method of fusing word embeddings that were trained on the same corpus but with different initializations is considered, which demonstrates consistent improvements over the raw models as well as their simplistic average, on a range of tasks. Expand

**Cross-document Coreference Resolution over Predicted Mentions**
Arie Cattan, Alon Eirew, Gabriel Stanovsky, Mandar Joshi, Ido Dagan · Computer Science · FINDINGS · 2 June 2021
TLDR  This work introduces the first end-to-end model for CD coreference resolution from raw text, which extends the prominent model for withindocument coreference to the CD setting and achieves competitive results for event and entity coreferenceresolution on gold mentions. Expand

**Extending Multi-Document Summarization Evaluation to the Interactive Setting**
Ori Shapira, Ramakanth Pasunuru, H. Ronen, M. Bansal, Yael Amsterdamer, Ido Dagan · Computer Science · NAACL · 1 June 2021

# Authors who work on related tasks or use similar methods

**PDDL-the planning domain definition language**

D. McDermott, M. Ghallab, +5 authors D. Wilkins · Computer Science · 1998

TLDR This manual describes the syntax of PDDL, the Planning Domain Definition Language, the problem-specification language for the AIPS-98 planning competition, and hopes to encourage empirical evaluation of planner performance, and development of standard sets of problems all in comparable notations. Expand

Dan We~~l~~

Automatic summarization
Generation of TLDR summary
analysis of scientific text
Paper recommendation
document level embedding of scientific documents
Network architecture (OS)
SpanBERT
Language model
Pre-trained language models
ELMo
Coreference resolution
coreference resolution across multiple documents
Neural network architectures
information retrieval systems
Natural language inference
search engines
human facing application
human centered ai
user interfaces
information extraction
extreme extraction
self training event extraction system'
….

**QA-Align: Representing Cross-Text Content Overlap by Aligning Question-Answer Propositions**

Daniel Weiss, Paula Roit, Ayal Klein, Ori Ernst, Ido Dagan · Computer Science · 26 September 2021

Multi-text applications, such as multidocument summarization, are typically required to model redundancies across related texts. Current methods confronting consolidation struggle to fuse overlapping… Expand

⬇ View PDF on arXiv   🔖 Save   🔔 Alert   ❝❞ Cite   ⚛ Research Feed

Ido Dag~~an~~

Cross-document coreference resolution
Coreference resolution
Cross-text alignment
Few shot Relation Extraction
Crowdsourcing
Network architecture (Deep learning)
Evaluation of automated summaries
Text summarization
Abstractive summarization
Multi-document summarization
Word embeddings
Question-driven SRL
QA-SRL
Recognizing Textual Entailment
Textual Inference
NLP
Cross-document language modeling
BERT
Transformers for multiple documents
Hypernym discovery
Low-level textual inference
….

PDDL-the planning domain definition language

D. McDermott, M. Ghallab, +5 authors D. Wilkins · Computer Science · 1998

TLDR This manual describes the syntax of PDDL, the Planning Domain Definition language, the problem-specification language for the AIPS-98 planning competition, and hopes to encourage empirical evaluation of planner performance, and development of standard sets of problems all in comparable notations. Expand

QA-Align: Representing Cross-Text Content Overlap by Aligning Question-Answer Propositions

Daniel Weiss, Paula Roit, Ayal Klein, Ori Ernst, Ido Dagan · Computer Science · 26 September 2021

Multi-text applications, such as multidocument summarization, are typically required to model redundancies across related texts. Current methods confronting consolidation struggle to fuse overlapping… Expand

View PDF on arXiv    Save    Alert    Cite    Research Feed

Dan Weld

Ido Dagan

**Diversity**

Automatic summarization
Generation of TLDR summary
analysis of scientific text
Paper recommendation
document level embedding of scientific documents
Network architecture (OS)
SpanBERT
Language model
Pre-trained language models
ELMo
Coreference resolution
coreference resolution across multiple documents
Neural network architectures
information retrieval systems

**False Positives**

search engines
human facing application
human centered ai
user interfaces
information extraction
extreme extraction
self training event extraction system'

….

Cross-document coreference resolution
Coreference resolution
Cross-text alignment
Few shot Relation Extraction
Crowdsourcing
Network architecture (Deep learning)
Evaluation of automated summaries
Text summarization
Abstractive summarization
Multi-document summarization
Word embeddings
Question-driven SRL
QA-SRL

**Type I errors**

Textual inference
NLP
Cross-document language modeling
BERT
Transformers for multiple documents
Hypernym discovery
Low-level textual inference

….

Dan Wel...

Ido Dag...

**Ambiguity**

PDDL-the planning domain definition language
D. McDermott, M. Ghallab, +5 authors  D. Wilkins · Computer Science · 1998
TLDR This manual describes the syntax of PDDL, the Planning Domain Definition Language, the problem-specification language for the AIPS-98 planning competition, and hopes to encourage empirical evaluation of planner performance, and development of standard sets of problems all in comparable notations. Expand

QA-Align: Representing Cross-Text Content Overlap by Aligning Question-Answer Propositions
Daniel Weiss, Paula Roit, Ayal Klein, Ori Ernst, Ido Dagan · Computer Science · 26 September 2021
Multi-text applications, such as multidocument summarization, are typically required to model redundancies across related texts. Current methods confronting consolidation struggle to fuse overlapping... Expand
View PDF on arXiv    Save    Alert    Cite    Research Feed

Automatic summarization
Generation of TLDR summary
analysis of scientific text
Paper recommendation
document level embedding of scientific documents

Network architecture
(systems)

SpanBERT
Language model
Pre-trained language models
ELMo
Coreference resolution
coreference resolution across multiple documents
Neural network architectures
information retrieval systems
Natural language inference
search engines
human facing application
human centered ai
user interfaces
information extraction
extreme extraction
self training event extraction system'
....

Cross-document coreference resolution
Coreference resolution
Cross-text alignment
Few shot Relation Extraction
Crowdsourcing

Network architecture
(deep learning)

Evaluation of automated summaries
Text summarization
Abstractive summarization
Multi-document summarization
Word embeddings
Question-driven SRL
QA-SRL
Recognizing Textual Entailment
Textual Inference
NLP
Cross-document language modeling
BERT
Transformers for multiple documents
Hypernym discovery
Low-level textual inference
....

PDDL-the planning domain definition language

D. McDermott, M. Ghallab, +5 authors  D. Wilkins · Computer Science · 1998

TLDR  This manual describes the syntax of PDDL, the Planning Domain Definition Language, the problem-specification language for the AIPS-98 planning competition, and hopes to encourage empirical evaluation of planner performance, and development of standard sets of problems all in comparable notations. Expand

Dan Wel

Automatic summarization
Generation of TLDR summary
analysis of scientific text
Paper recommendation
document level embedding of scientific documents
Network architecture (OS)
SpanBERT
Language model
ELMo
Coreference resolution
coreference resolution across multiple documents
Neural network architectures
information retrieval systems
Natural language inference
search engines
human facing application
human centered ai
user interfaces
information extraction

Pre-trained language models

extreme extraction
self training event extraction system'
….

**Hierarchy**

QA-Align: Representing Cross-Text Content Overlap by Aligning Question-Answer Propositions

Daniel Weiss, Paula Roit, Ayal Klein, Ori Ernst, Ido Dagan · Computer Science · 26 September 2021

Multi-text applications, such as multidocument summarization, are typically required to model redundancies across related texts. Current methods confronting consolidation struggle to fuse overlapping… Expand

View PDF on arXiv · Save · Alert · Cite · Research Feed

Ido Dag

Cross-document coreference resolution
Coreference resolution
Cross-text alignment
Few shot Relation Extraction
Crowdsourcing
Network architecture (Deep learning)
Evaluation of automated summaries
Text summarization
Abstractive summarization
Multi-document summarization
Word embeddings
Question-driven SRL
QA-SRL
Recognizing Textual Entailment
Textual Inference
NLP

Cross-document language modeling

BERT
Transformers for multiple documents
Hypernym discovery
Low-level textual inference
….

# Limitations of Previous Work

Cross-Document Coreference Resolution (CDCR)

- No abstract technical concepts
  - No work in science!
- No cross-document **hierarchy**

**Doc 1**: President Obama will **name** Dr. Regina Benjamin as U.S. Surgeon General in a Rose Garden announcement late this morning. Benjamin, an Alabama family physician, [...]
**Doc 2**: [...] Obama **nominates** new surgeon general: MacArthur "genius grant "fellow Regina Benjamin. [...]

ECB+[1]

[1] Cybulska, Agata and P. Vossen. "Using a sledgehammer to crack a nut? Lexical diversity and event coreference resolution." LREC (2014).

**Goal**: *Address cross-document ambiguity, diversity and hierarchy together*

# New task: *cross-document ambiguity, diversity, hierarchy*

**Input:**

Concept mentions in scientific papers

> ... **self-driving cars** have made it increasingly urgent ...

> ... navigation for **autonomous vehicles** in real-life traffic ...

> ... transformer models in **computer vision** ...

> ... we use **categorical image generation** ...

> ... the problem of **generating images** ...

# New task: *cross-document ambiguity, diversity, hierarchy*

**Input:**

Concept mentions in scientific papers

... **self-driving cars** have made it increasingly urgent ...

... navigation for **autonomous vehicles** in real-life traffic ...

... transformer models in **computer vision** ...

... we use **categorical image generation** ...

... the problem of **generating images** ...

self-driving cars, autonomous vehicles

computer vision, AI-based visual understanding

image synthesis task, generating images

categorical image generation, class-conditional image synthesis

**Output:**

- Coreference Clusters

# Hierarchical Cross-Document Coreference Resolution (H-CDCR)

# **SciCo:** A new large-scale dataset annotated by domain experts



**(1) Extract seeds from text & AI KB**

Hypernym extraction

... Travelling Salesman Problem is a major problem in graph theory...

... Graph Isomorphism is one of the classical problems of graph theory...

Papers w/ Code

Question Answering subtasks

Question Answering | Open-Domain Question Answering | Answer Selection

**(2) Combine information resources**

Graph theory

Traveling Salesman Problem — Graph Isomorphism

Question Answering

Open-Domain Question Answering — Answer Selection

**(4) Annotate coreference & hierarchy**

Document 14
Predicting Lexical Answer Types in Open Domain QA ( 2012)

Automatic open-domain Question Answering has been a long standing research challenge community . IBM Research undertook this challenge with the design of the DeepQA architecture

▾ 🗀 Question Answer-ing
  🗀 Answer sentence selection
  🗀 open-domain complex question answering
  🗀 question answering benchmarks

**(3) Retrieve candidate mentions + context**

In this paper we focus on a **question answering scenario** ...

**Automatic open-domain Question Answering** has been a long standing challenge...

**Answer sentence selection** is a question answering paradigm ...

Novel evaluation metrics for _hierarchical_ cross-document coreference resolution

# Custom Baseline Models

# Baseline I: Two-Step Model

Two steps:

1. **Concept clusters:** Apply existing SOTA CDCR model

2. **Hierarchy**: Find relations between predicted clusters

# Baseline I: Two-Step Model

✅ We have clusters (using SOTA trained on SciCo)

self-driving cars, autonomous vehicles

computer vision, AI-based visual understanding

image synthesis task, generating images

categorical image generation, class-conditional image synthesis

# Baseline I: Two-Step Model

✅ We have clusters (using SOTA trained on SciCo)

❓ How do we **infer the hierarchy**?

computer vision, AI-based visual understanding

self-driving cars, autonomous vehicles

image synthesis task, generating images

categorical image generation, class-conditional image synthesis

# Intuition: Referential hierarchy as '**multi-document textual entailment**'

# Intuition: Referential hierarchy as '**multi-document textual entailment**'

# Intuition: Referential hierarchy as '**multi-document entailment**'

...**Drug-Drug Interaction (DDI) Extraction** from Drug Labels...

...high-coverage corpus that can be used for **IE...**

*entails*

...a system for **automatic extraction of drug-drug interactions** in biomedical texts...

...natural language processing (NLP) problems, such as **information extraction...**

Apply **pre-trained NLI models** over simple concatenation of mentions

# Unified Model

**Unified** model with multiclass formulation
for pairs of mentions m1,m2 with classes {≍, ←, →, None}

$$L = -\frac{1}{N} \sum_{\substack{m_1, m_2 \in \mathcal{M} \\ m_1 \neq m_2}} y \cdot log(f(m_1, m_2))$$

# Unified Model

ƒ(m1,m2): **Cross-encode** mentions m1,m2 with entity markers

[CLS]…an experiment in **<m> definition extraction </m>** from legal texts … [SEP] … natural language processing problems, such as **<m> information extraction </m>,** summarization and dialogue.… [SEP]

# Unified Model

$f$(m1,m2): **Cross-encode** mention pairs m1,m2 with entity markers

[CLS] …an experiment in **<m> definition extraction </m>** from legal texts … [SEP] … natural language processing problems, such as **<m> information extraction </m>,** summarization and dialogue…. [SEP]

Leverage **wide context**:

- LongFormer[1]: Transformers for ***long-sequences***
- CDLM[2]: Transformers for ***multi-document*** tasks

[1] Beltagy, Iz, Matthew E. Peters and Arman Cohan. "Longformer: The Long-Document Transformer." (2020)
[2] Caciularu, Avi, Arman Cohan, Iz Beltagy, Matthew E. Peters, Arie Cattan and Ido Dagan. "CDLM: Cross-Document Language Modeling." (EMNLP FINDINGS 2021).

Inference:

1. **Clusters**: Agglomerative clustering over mention-pair coref. scores

$pr(\approx) =0.01$

… the problem of
**generating images** …

… **self-driving cars**
have made it
increasingly urgent …

$pr(\approx) =0.75$

… navigation for
**autonomous vehicles**
in real-life traffic …

Inference:

1. **Clusters**: Agglomerative clustering over mention-pair coref. scores

2. **Hierarchy**:

   a. Score (prob.) that $C_1$ is a child of $C_2$:

   $$s(\mathcal{C}_1, \mathcal{C}_2) = \frac{1}{|\mathcal{C}_1| \cdot |\mathcal{C}_2|} \sum_{m_1 \in \mathcal{C}_1} \sum_{m_2 \in \mathcal{C}_2} f_{\text{is-child}}(m_1, m_2)$$

   b. Greedy construction of hierarchy (to avoid cycles)

$C_2$ image synthesis task, generating images

$pr(\rightarrow) = 0.7$

$pr(\rightarrow) = 0.8$

categorical image generation, class-conditional image synthesis $C_1$

# Results

| | Coreference | Hierarchy | | Path |
|---|---|---|---|---|
| | CoNLL F1 | F1 | F1-50% | Ratio |
| IAA (AVG) | 82.7 | 68.9 | 62.8 | 64.5 |
| IAA (MAX-Macro) | 90.2 | 82.3 | 77.7 | 78.4 |
| CA$_{News}$ | 52.4 | 37.1 | 13.0 | 24.1 |
| CA$_{Sci-News}$ | 43.5 | 29.2 | 12.3 | 21.6 |
| CA$_{SCICO}$ | 55.2 | 23.7 | 15.8 | 21.2 |
| CA$_{SCICO}$ + CS-RoBERTa | 57.4 | 23.5 | 16.1 | 23.6 |
| CA$_{SCICO}$ + SciBERT | 66.8 | 23.8 | 17.8 | 28.4 |
| Unified$_{Longformer}$ | **77.2** | 44.5 | **36.1** | **47.2** |
| Unified$_{CDLM}$ | 77.0 | **44.8** | 35.5 | 45.9 |

Two-step models

Multi-class cross-encoder

# Results

| | |
|---|---|
| False positives | Type I errors |

$\longleftrightarrow$

| | |
|---|---|
| Class imbalance | Skewed label distribution |

$\longleftrightarrow$

| | |
|---|---|
| Kernel (OS) | Kernel (ML) |

❌

| | |
|---|---|
| Network architecture (systems) | Network architecture (deep learning) |

❌

# Results

False positives ← - - - → Type I errors

Class imbalance Skewed label distribution

**Much potential for model improvements!**

Kernel (OS) Kernel (ML)

Network architecture (systems) Network architecture (deep learning)

# SciCo:
## Hierarchical Cross-Document Coreference



**Input:** Concept mentions in scientific papers

... **self-driving cars** have made it increasingly urgent ...

... navigation for **autonomous vehicles** in real-life traffic ...

... the problem of **sequence tagging** ...

We use a **CRF tagger** ...

**Output:** Coreference Clusters + hierarchy

self-driving cars, autonomous vehicles

computer vision, AI-based visual understanding

sequence tagging

image synthesis task, generating images

CRF sequence model, CRF tagger

categorical image generation, class-conditional image synthesis

Network Architecture (Systems)

Network Architecture (AI)

*Outstanding Paper Award*

AUTOMATED KNOWLEDGE BASE 2021 | irvine CONSTRUCTION

# Document Similarity for Science

# SPECTER: Document-level Representation Learning using Citation-informed Transformers

**Arman Cohan**[†*]  **Sergey Feldman**[†*]  **Iz Beltagy**[†]  **Doug Downey**[†]  **Daniel S. Weld**[†,‡]

[†]Allen Institute for Artificial Intelligence
[‡]Paul G. Allen School of Computer Science & Engineering, University of Washington
{armanc,sergey,beltagy,dougd,danw}@allenai.org

**SPECTER: Document-level Representation Learning using Citation-informed Transformers**

Arman Cohan[†*]    Sergey Feldman[†*]    Iz Beltagy[†]    Doug Downey[†]    Daniel S. Weld[†,‡]

[†]Allen Institute for Artificial Intelligence
[‡]Paul G. Allen School of Computer Science & Engineering, University of Washington
{armanc, sergey, beltagy, dougd, danw}@allenai.org

**Contrastive learning:** Learn embeddings of papers that:

**Pull together** papers that are related/similar

**Pull apart** papers that are unrelated/dissimilar

Citation network

**Related** papers:
Paper p cites paper p'

*Citation network*

**Un**related papers:
No citation link between
paper p and paper n

*Citation network*

$$f(p, p') \qquad < \qquad f(p, n)$$

Distance between paper **p**, and a *related* paper **p'**

Distance between paper **p**, and a ***un*relate**d paper **n**

**Contrastive training loss**

$$\mathcal{L}_f(p, p', n) = \texttt{max}[f(p, p') - f(p, n) + m, 0]$$

$$f(\boxminus, \boxminus) = \text{ ?}$$

**Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations**

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. … Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling, but their models assume relations are disjoint -for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.

**Multi-instance Multi-label Learning for Relation Extraction**

Distant supervision for relation extraction (RE) -- gathering training data by aligning a database of facts with text -- is an efficient approach to scale RE to thousands of different relations. However, this introduces a challenging learning scenario where the relation expressed by a pair of entities found in a sentence is unknown. For example, a sentence containing Balzac and France may express BornIn or Died, an unknown relation, or no relation at all. Because of this, traditional supervised learning, which assumes that each example is explicitly mapped to a label, is not appropriate. We propose a novel approach to multi-instance multi-label learning for RE, which jointly models all the instances of a pair of entities in text and all their labels using a graphical model with latent variables. Our model performs competitively on two difficult domains.

$$f(\boxminus, \boxminus)$$

BERT

BERT

$\in \mathbb{R}^n \left[ \cdots \right]$

$\in \mathbb{R}^n \left[ \cdots \right]$

SPECTER: Distance based on one **_overall_** document vector

(e.g., Euclidean distance between p and p' vector embeddings)

# Scientific Documents are **Multi**-Faceted



**Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations**

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. … Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling, but their models assume relations are disjoint -for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.

Problem

Method

Experimental method

Result

# Scientific Documents are **<u>Multi</u>**-Faceted



Can ***aspect***-level modeling lead to better ***document***-level similarity?

How do we identify **fine-grained aspect similarity?**
*(Descriptions of methodologies, experiments, findings...)*



**Paper A**          **Paper B**

# How do we identify **fine-grained aspect similarity?**
*(Descriptions of methodologies, experiments, findings...)*



**Paper A**          **Paper B**

**No Gold Labels**

... both [PaperA, 1999] and [PaperB, 2019] use a group sequential experimental design...

**Paper A**

**Paper B**

Co-Citation Sentences

Co-citation sentences provide explanations regarding *how* two papers are related

# Textual supervision: co-citation contexts



## 2. RELATED WORK

Recurrent neural networks (RNNs), and LSTMs in particular, have recently been used to generate sequences in various domains, such as music [7], text [15, 29], and handwriting [15]. In information retrieval, RNNs have been used, e.g., for extracting sentence-level semantic vectors [26] and context-aware query suggestion [28]. Other kinds of deep neural networks have been used to project queries and documents to low-dimensional semantic spaces [18] and to learn fixed-length vectors for variable-length pieces of texts, such as sentences, paragraphs, and documents [21]. Various types of task activities have been studied in the literature as a basis for query suggestion or query support. Motivated by the observation that a notable propor-

*Co-Citation Context*

**Learning deep structured semantic models for web search using clickthrough data**
Latent semantic models, such as LSA, intend to map a query to its relevant documents at the semantic level where keyword-based matching often fails. In this study we strive to develop a series of new latent semantic models with a deep structure that project queries and documents into a common low-dimensional space where the relevance of a document given a query is readily computed as the distance between them. The proposed deep structured semantic models are discriminatively trained by maximizing the conditional likelihood of the clicked documents given a query using the clickthrough data. … Results show that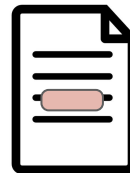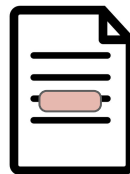 our best model significantly outperforms other latent semantic models, which were considered state-of-the-art in the … to the work presented in this paper.

**Distributed Representations of Sentences and Documents**
Many machine learning algorithms require the input to be represented as a fixed length feature vector. When it comes to texts, one of the most common representations is bag-of-words. Despite their popularity, bag-of-words models have two major weaknesses: they lose the ordering of the words and they also ignore semantics of the words. For example, "powerful," "strong" and "Paris" are equally distant. In this paper, we propose an unsupervised algorithm that learns vector representations of sentences and text documents. … Empirical results show that our technique outperforms bag-of-words models as well as other techniques for text representations. Finally, we achieve new state-of-the-art results on several text classification and sentiment analysis tasks.

*Co-Cited Abstracts*

# Textual supervision: co-citation contexts

context-aware query suggestion [28]. Other kinds of deep neural networks have been used to project queries and documents to low-dimensional semantic spaces [18] and to learn fixed-length vectors for variable-length pieces of texts, such as sentences, paragraphs, and documents [21].

*Co-Citation Context*

I. **Sentence alignment:**
Encode each sentence with BERT.
Find pair of sentences maximally similar
to the co-citation context ("aligned").

**Learning deep structured semantic models for web search using clickthrough data**
Latent semantic models, such as LSA, intend to map a query to its relevant documents at the semantic level where keyword-based matching often fails. In this study we strive to develop a series of new latent semantic models with a deep structure that project queries and documents into a common low-dimensional space where the relevance of a document given a query is readily computed as the distance between them. The proposed deep structured semantic models are discriminatively trained by maximizing the conditional likelihood of the clicked documents given a query using the clickthrough data. … Results show that our best model significantly outperforms other latent semantic models, which were considered state-of-the-art in the … to the work presented in this paper.

**Distributed Representations of Sentences and Documents**
Many machine learning algorithms require the input to be represented as a fixed length feature vector. When it comes to texts, one of the most common representations is bag-of-words. Despite their popularity, bag-of-words models have two major weaknesses: they lose the ordering of the words and they also ignore semantics of the words. For example, "powerful," "strong" and "Paris" are equally distant. In this paper, we propose an unsupervised algorithm that learns vector representations of sentences and text documents. … Empirical results show that our technique outperforms bag-of-words models as well as other techniques for text representations. Finally, we achieve new state-of-the-art results on several text classification and sentiment analysis tasks.

*Co-Cited Abstracts*

# Textual supervision: co-citation contexts



context-aware query suggestion [28]. Other kinds of deep neural networks have been used to project queries and documents to low-dimensional semantic spaces [18] and to learn fixed-length vectors for variable-length pieces of texts, such as sentences, paragraphs, and documents [21].

II. Learn contextualized **sentence embeddings**
that **minimize distance** between
the aligned sentences (w/ contrastive loss)

# **Multiple** matches



Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. … Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling, but their models assume relations are disjoint -for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.

Problem
Method
Experimental method
Result

Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. … Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling, but their models assume relations are disjoint -for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.

Problem
Method
Experimental method
Result

Multiple aspect alignments

$$f(p, p')$$

Distance between two
**papers** p, p′

$$f(p, p') = \sum_{(i, i') \in \mathcal{A}_p \times \mathcal{A}_{p'}}$$
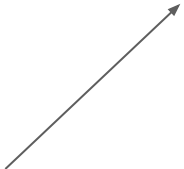
Distance between two
**papers** p, p′

Sum over all pairs of
aspects across p, p′

$$f(p, p') = \sum_{(i,i') \in \mathcal{A}_p \times \mathcal{A}_{p'}} w_{i,i'} \cdot d_{i,i'}$$

Distance between two
**papers** p, p′

Distance between two
**aspects** i, i′

Importance **weight** for
aspect pair i, i'

$$f(p, p') = \sum_{(i,i') \in \mathcal{A}_p \times \mathcal{A}_{p'}} w_{i,i'} \cdot d_{i,i'}$$
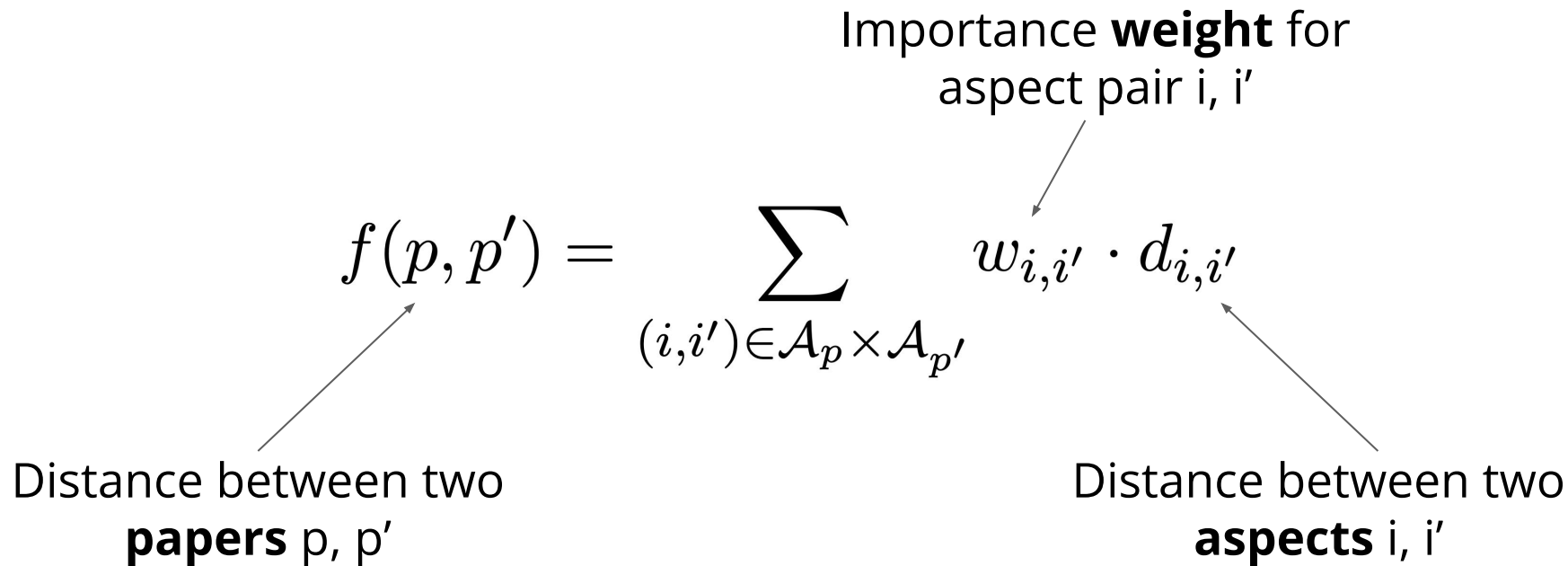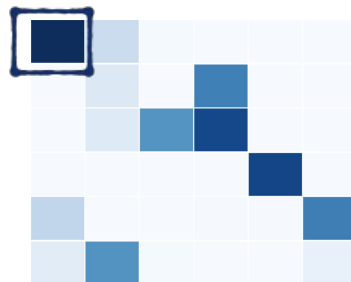
Distance between two
**papers** p, p'

Distance between two
**aspects** i, i'

$$w_{i,i'}$$

# Aspect-level alignment weight matrix

Motivation



**Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations**

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. ... Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling, but their models assume relations are disjoint -for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.
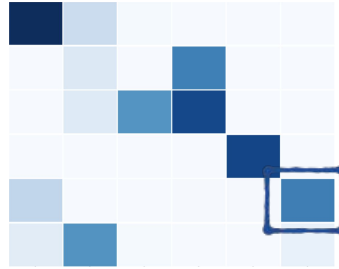
**Multi-instance Multi-label Learning for Relation Extraction**

Distant supervision for relation extraction (RE) -- gathering training data by aligning a database of facts with text -- is an efficient approach to scale RE to thousands of different relations. However, this introduces a challenging learning scenario where the relation expressed by a pair of entities found in a sentence is unknown. For example, a sentence containing Balzac and France may express BornIn or Died, an unknown relation, or no relation at all. Because of this, traditional supervised learning, which assumes that each example is explicitly mapped to a label, is not appropriate. We propose a novel approach to multi-instance multi-label learning for RE, which jointly models all the instances of a pair of entities in text and all their labels using a graphical model with latent variables. Our model performs competitively on two difficult domains.

$$w_{i,i'}$$

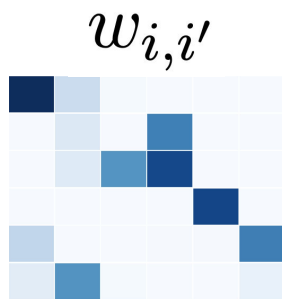# Aspect-level alignment weight matrix



Method

**Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations**

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. … Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling, but their models assume relations are disjoint -for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.

**Multi-instance Multi-label Learning for Relation Extraction**

Distant supervision for relation extraction (RE) -- gathering training data by aligning a database of facts with text -- is an efficient approach to scale RE to thousands of different relations. However, this introduces a challenging learning scenario where the relation expressed by a pair of entities found in a sentence is unknown. For example, a sentence containing Balzac and France may express BornIn or Died, an unknown relation, or no relation at all. Because of this, traditional supervised learning, which assumes that each example is explicitly mapped to a label, is not appropriate. We propose a novel approach to multi-instance multi-label learning for RE, which jointly models all the instances of a pair of entities in text and all their labels using a graphical model with latent variables. Our model performs competitively on two difficult domains.
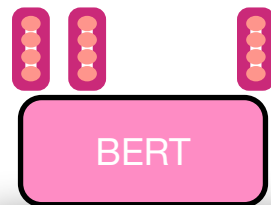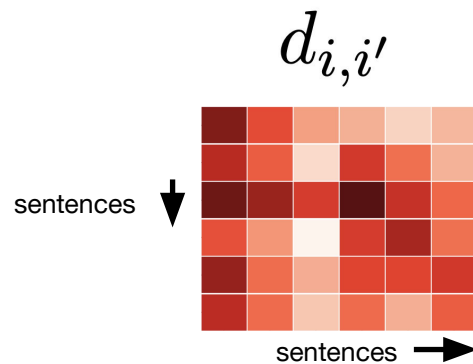
$$f(p, p') = \quad w_{i,i'} \quad * \quad d_{i,i'}$$

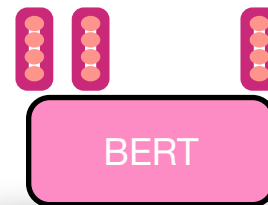$$f(p, p') = \quad w_{i,i'} \quad * \quad d_{i,i'}$$

**Learn** aspect weights and distances such that
the document-level **contrastive loss** is minimized

$$\mathcal{L}_f(p, p', n) = \texttt{max}[f(p, p') - f(p, n) + m, 0]$$

$$d_{i,i'}$$



sentences →

↓ sentences

**BERT**

**Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations**

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. ... Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling, but their models assume relations are disjoint - for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.
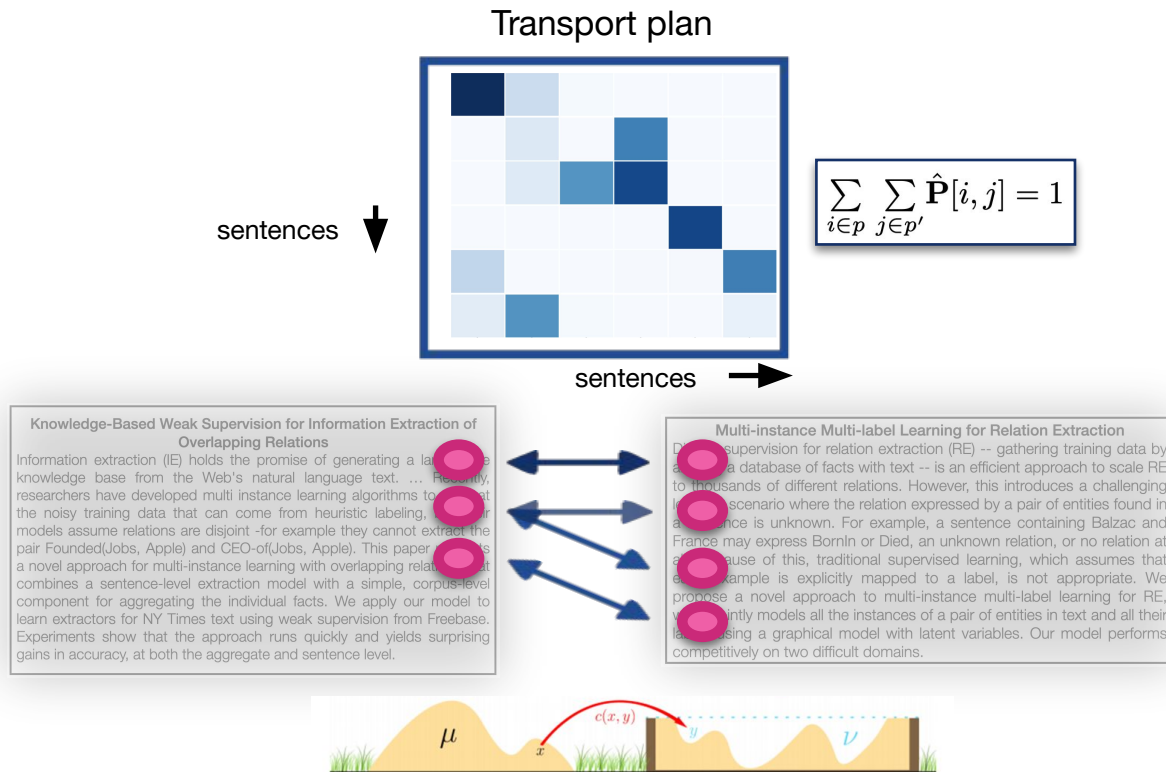
**BERT**

**Multi-instance Multi-label Learning for Relation Extraction**
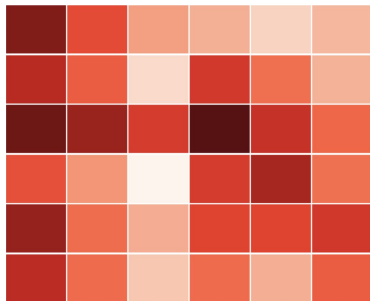
Distant supervision for relation extraction (RE) -- gathering training data by aligning a database of facts with text -- is an efficient approach to scale RE to thousands of different relations. However, this introduces a challenging learning scenario where the relation expressed by a pair of entities found in a sentence is unknown. For example, a sentence containing Balzac and France may express BornIn or Died, an unknown relation, or no relation at all. Because of this, traditional supervised learning, which assumes that each example is explicitly mapped to a label, is not appropriate. We propose a novel approach to multi-instance multi-label learning for RE, which jointly models all the instances of a pair of entities in text and all their labels using a graphical model with latent variables. Our model performs competitively on two difficult domains.
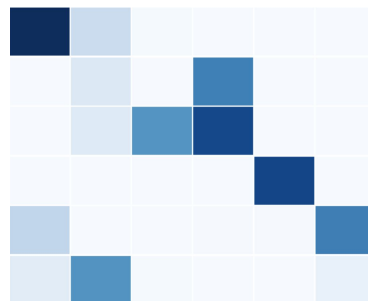
# Optimal Transport:
Soft, sparse alignment between sets of aspects

Transport plan

sentences

$$\sum_{i \in p} \sum_{j \in p'} \hat{\mathbf{P}}[i,j] = 1$$

sentences

**Knowledge-Based Weak Supervision for Information Extraction of Overlapping Relations**

Information extraction (IE) holds the promise of generating a large-scale knowledge base from the Web's natural language text. ... Recently, researchers have developed multi instance learning algorithms to combat the noisy training data that can come from heuristic labeling. ... prior models assume relations are disjoint -for example they cannot extract the pair Founded(Jobs, Apple) and CEO-of(Jobs, Apple). This paper presents a novel approach for multi-instance learning with overlapping relations that combines a sentence-level extraction model with a simple, corpus-level component for aggregating the individual facts. We apply our model to learn extractors for NY Times text using weak supervision from Freebase. Experiments show that the approach runs quickly and yields surprising gains in accuracy, at both the aggregate and sentence level.

**Multi-instance Multi-label Learning for Relation Extraction**

Distant supervision for relation extraction (RE) -- gathering training data by aligning a database of facts with text -- is an efficient approach to scale RE to thousands of different relations. However, this introduces a challenging learning scenario where the relation expressed by a pair of entities found in a sentence is unknown. For example, a sentence containing Balzac and France may express BornIn or Died, an unknown relation, or no relation at all. Because of this, traditional supervised learning, which assumes that each example is explicitly mapped to a label, is not appropriate. We propose a novel approach to multi-instance multi-label learning for RE, which jointly models all the instances of a pair of entities in text and all their labels using a graphical model with latent variables. Our model performs competitively on two difficult domains.

$\mu$  $x$  $c(x,y)$  $y$  $\nu$

Pairwise distances

Transport plan

$$\hat{\mathbf{P}} = \underset{\mathbf{P} \in \mathcal{S}}{\operatorname{argmin}} \langle \mathbf{D}, \mathbf{P} \rangle$$

Linear optimization
problem, compatible
with autodiff, GPUs

# Evaluation

- Document-level similarity
  - Two biomedical paper datasets

# Evaluation

- Document-level similarity
  - Two biomedical paper datasets

- Aspect-level similarity
  - Recent CS paper dataset

# Evaluation

## Document level

| Models | TRECCOVID$_{RF}$ | | RELISH | |
|---|---|---|---|---|
| | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ |
| MPNET-1B | 17.35 | 43.87 | 52.92 | 69.69 |
| SENTBERT-PP | 11.12 | 34.85 | 50.80 | 67.35 |
| SENTBERT-NLI | 13.43 | 40.78 | 47.02 | 63.56 |
| UNSIMCSE-BERT | 9.85 | 34.27 | 45.79 | 62.02 |
| SUSIMCSE-BERT | 11.50 | 37.17 | 47.29 | 63.93 |
| CoSentBert | 12.80 | 38.07 | 50.04 | 66.35 |
| ICTSentBert | 9.80 | 33.62 | 47.72 | 63.71 |
| OTMPNET-1B | 27.46 | 58.70 | 57.46 | 74.64 |
| SPECTER | 28.24 | 59.28 | 60.62 | 77.20 |
| SCINCL | 28.73 | 59.16 | 62.09 | 78.72 |
| SPECTER-COCITE$_{Scib}$ | 30.60 | 62.07 | 61.43 | 78.01 |
| SPECTER-COCITE$_{Spec}$ | 28.59 | 60.07 | 61.43 | 77.96 |
| TSASPIRE$_{Spec}$ | 26.24 | 56.55 | 61.29 | 77.89 |
| OTASPIRE$_{Spec}$ | **30.92** | **62.23** | 62.57 | 78.95 |
| TS+OTASPIRE$_{Spec}$ | 30.90 | 62.18 | **62.71** | **79.18** |

## Aspect level

| CSFCUBE facets → | Aggregated | | Background | | Method | | Result | |
|---|---|---|---|---|---|---|---|---|
| Models | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ |
| MPNET-1B | 34.64 | 54.94 | 41.06 | 65.86 | 27.21 | 42.48 | 36.07 | 54.94 |
| SENTBERT-PP | 26.77 | 48.57 | 35.43 | 60.80 | 16.19 | 33.40 | 29.16 | 48.57 |
| SENTBERT-NLI | 25.23 | 45.39 | 30.78 | 54.23 | 15.02 | 31.10 | 30.27 | 45.39 |
| UNSIMCSE-BERT | 24.45 | 42.59 | 30.03 | 51.59 | 14.82 | 31.23 | 28.76 | 42.59 |
| SUSIMCSE-BERT | 23.24 | 43.45 | 30.52 | 55.22 | 13.99 | 30.88 | 25.58 | 43.45 |
| CoSentBert | 28.95 | 50.68 | 35.78 | 61.27 | 19.27 | 38.77 | 32.15 | 50.68 |
| ICTSentBert | 28.61 | 48.13 | 35.93 | 59.80 | 15.62 | 35.91 | 34.42 | 48.13 |
| OTMPNET-1B | 36.41 | 56.91 | 43.23 | 67.60 | 28.69 | 43.49 | 37.76 | 60.30 |
| SPECTER | 34.23 | 53.28 | 43.95 | 66.70 | 22.44 | 37.41 | 36.79 | 56.67 |
| SCINCL | 39.37 | 59.24 | 49.64 | 70.02 | 27.14 | 46.61 | 41.83 | 61.70 |
| SPECTER-COCITE$_{Scib}$ | 37.90 | 58.16 | 48.40 | 68.71 | 26.95 | 46.79 | 38.93 | 59.68 |
| SPECTER-COCITE$_{Spec}$ | 37.39 | 58.38 | 49.99 | 70.03 | 25.60 | 45.99 | 37.33 | 59.95 |
| TSASPIRE$_{Spec}$ | 40.26 | 60.71 | 49.58 | 70.22 | **28.86** | **48.20** | 42.92 | 64.39 |
| OTASPIRE$_{Spec}$ | **40.79** | **61.41** | 50.56 | **71.04** | 27.64 | 46.46 | **44.75** | **67.38** |
| TS+OTASPIRE$_{Spec}$ | 40.26 | 60.86 | **51.79** | 70.99 | 26.68 | 47.60 | 43.06 | 64.82 |

**Substantial gains in MAP, NDCG**

# Evaluation

## Document level

| Models | TRECCOVID$_{RF}$ | | RELISH | |
|---|---|---|---|---|
| | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ |
| MPNET-1B | 17.35 | 43.87 | 52.92 | 69.69 |
| SENTBERT-PP | 11.12 | 34.85 | 50.80 | 67.35 |
| SENTBERT-NLI | 13.43 | 40. | | |
| UNSIMCSE-BERT | 9.85 | 34. | | |
| SUSIMCSE-BERT | 11.50 | 37. | | |
| CoSentBert | 12.80 | 38. | | |
| ICTSENTBERT | 9.80 | 33. | | |
| OTMPNET-1B | 27.46 | 58. | | |
| SPECTER | 28.24 | 59.28 | 60.62 | 77.20 |
| SCINCL | 28.73 | 59.16 | 62.09 | 78.72 |
| SPECTER-COCITE$_{Scib}$ | 30.60 | 62.07 | 61.43 | 78.01 |
| SPECTER-COCITE$_{Spec}$ | 28.59 | 60.07 | 61.43 | 77.96 |
| TSASPIRE$_{Spec}$ | 26.24 | 56.55 | 61.29 | 77.89 |
| OTASPIRE$_{Spec}$ | **30.92** | **62.23** | 62.57 | 78.95 |
| TS+OTASPIRE$_{Spec}$ | 30.90 | 62.18 | **62.71** | **79.18** |

## Aspect level

| CSFCUBE facets → | Aggregated | | Background | | Method | | Result | |
|---|---|---|---|---|---|---|---|---|
| Models | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ | MAP | NDCG$_{\%20}$ |
| MPNET-1B | 34.64 | 54.94 | 41.06 | 65.86 | 27.21 | 42.48 | 36.07 | 54.94 |
| SENTBERT-PP | 26.77 | 48.57 | 35.43 | 60.80 | 16.19 | 33.40 | 29.16 | 48.57 |
| SENTBERT-NLI | 25.23 | 45.39 | 30.78 | 54.23 | 15.02 | 31.10 | 30.27 | 45.39 |
| UNSIMCSE-BERT | 24.45 | 42.59 | 30.03 | 51.59 | 14.82 | 31.23 | 28.76 | 42.59 |
| | | | | | 5.22 | 13.99 | 30.88 | 25.58 | 43.45 |
| | | | | | 1.27 | 19.27 | 38.77 | 32.15 | 50.68 |
| | | | | | 9.80 | 15.62 | 35.91 | 34.42 | 48.13 |
| | | | 7.60 | 28.69 | 43.49 | 37.76 | 60.30 |
| | | | 6.70 | 22.44 | 37.41 | 36.79 | 56.67 |
| | | | 0.02 | 27.14 | 46.61 | 41.83 | 61.70 |
| | | | 8.71 | 26.95 | 46.79 | 38.93 | 59.68 |
| | | | 0.03 | 25.60 | 45.99 | 37.33 | 59.95 |
| | | | 0.22 | **28.86** | **48.20** | 42.92 | 64.39 |
| OTASPIRE$_{Spec}$ | **40.79** | **61.41** | 50.56 | **71.04** | 27.64 | 46.46 | **44.75** | **67.38** |
| TS+OTASPIRE$_{Spec}$ | 40.26 | 60.86 | **51.79** | 70.99 | 26.68 | 47.60 | 43.06 | 64.82 |



allenai/aspire

Repo for Aspire - A scientific document similarity model based on matching fine-grained aspects of scientific papers.

AI2

**Substantial gains in MAP, NDCG**

# **Analogical retrieval**
for finding cross-domain inspirations



Solar System

Atom

🏆 *Best Research Paper*  KDD  CHI22

# Cross-document coreference
## +
# Document Level Similarity?

# Document Similarity & Retrieval with Aspect Alignments

**Best Research Paper**

NAACL 2022

From Retrieving & Extracting **Existing Knowledge**...

# …To Predicting **New Knowledge**

Can **neural language models trained on biomedical corpora** (e.g., PubMed) be leveraged for predicting new links in biomedical knowledge graphs?

Can we enhance prediction of **clinical outcomes in hospital patients** by retrieving **patient-specific medical literature?**

**Predict clinical outcomes of ICU patients:**

In-hospital mortality,
Prolonged mechanical ventilation,
Length of hospital stay,

...

# Predict clinical outcomes of ICU patients:



https://mimic.mit.edu/

State-of-art models predict outcomes
from internal data (e.g., patient notes)

State-of-art models predict outcomes
from internal data (e.g., patient notes)

Can we improve performance by
**adding patient-specific evidence
from the literature?**

# 1. Retrieve patient-relevant literature



Sparse retriever

Patient EHR

Cross-encoder reranker

**Reranked relevant abstracts**

Dense bi-encoder retriever

# 2. Literature-enhanced outcome prediction

**Adding literature boosts outcome prediction:**

Up to **5 point increase** in overall F1/AUC scores

Over **25% increase** in precision@Top-K scores

# Retrieval-augmented language models for clinical outcome prediction?

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, Edouard Grave, August 2022

# KG Completion for Drug Discovery

# KG Completion for Drug Discovery

**RepoDB**: Drug-disease pairs intended for **drug repositioning** research

# Knowledge Graph Link Prediction



**<head>** Aspirin

**<relation>** Treats

**<tail>** ?

Link prediction model

Headache ✅
Fever
…
COVID ❌

**Input**: Structured query

**Output**: Ranked candidates

# KG Embedding (KGE)

Learn embeddings for entities and relations

# Rich Textual Information...

# Approach with Language Models



$\text{text}(\mathbf{aspirin})[\text{SEP}]treats[\text{SEP}]\text{text}(\mathbf{fever})$ → Transformer **Language Model** → **Score**

[Kim et al CoNLL20, Wang et al WWW21]

# Graph and Literature Language Models: Complementary Strengths



Figure 2: Fraction of test set examples where each model performs better.

# Cross-Modal Link Prediction

# Cross-Modal Link Prediction

|  |  | RepoDB | | | Hetionet | | | MSI | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  |  | MRR | H@3 | H@10 | MRR | H@3 | H@10 | MRR | H@3 | H@10 |
| KGE | ComplEx | 62.3 | 71.1 | 85.6 | 45.9 | 53.6 | 77.8 | 40.3 | 44.3 | 57.5 |
|  | DistMult | 62.0 | 70.4 | 85.2 | 46.0 | 53.5 | 77.8 | 29.6 | 34.1 | 53.6 |
|  | RotatE | 58.8 | 65.9 | 79.8 | 50.6 | 58.2 | 79.3 | 32.4 | 35.3 | 49.8 |
|  | TransE | 60.0 | 68.6 | 81.1 | 50.2 | 58.0 | 79.8 | 32.7 | 36.5 | 53.8 |
| LM (fine-tuned) | RoBERTa | 51.7 | 60.3 | 82.3 | 46.4 | 53.6 | 76.9 | 30.1 | 33.3 | 50.6 |
|  | SciBERT | 59.7 | 67.6 | 88.5 | 50.3 | 57.1 | 79.1 | 34.2 | 37.9 | 55.0 |
|  | BioBERT | 58.2 | 65.8 | 86.8 | 50.3 | 57.5 | 79.4 | 33.4 | 37.1 | 54.8 |
|  | Bio+ClinicalBERT | 55.7 | 64.0 | 84.1 | 43.6 | 49.1 | 72.6 | 32.6 | 36.1 | 53.5 |

Integration of text and graph modalities provides **substantial relative improvements of 13–36% in** mean reciprocal rank (MRR).

**Multiple LM-based models** further boosts results.

# Cross-Modal Link Prediction: Challenges

|  | **Accuracy** (MRR) | **Efficiency** (inference sec, GPU) |
|---|---|---|
| KGE | 0.33 | $2 \times 10^1$ |
| → Rerank w/ cross-encoder LM | 0.38 (+0.05 MRR) | **$1 \times 10^6$ 11 days** |

WIKIPEDIA
The Free Encyclopedia

17K entities
206k edges

```
[CLS] aspirin
[SEP] treats      →  LM  →  Score
[SEP] fever
```

**Cross-encoder LM**
(**Slow**: combinatorial explosion)

# Our Solution: CascadER



**Tier 1**    **Tier 2**    **Tier *n***

Query:
<*aspirin*,
*treats*, ?>

Predict

Predict

Rerank with LM₁

Rerank with LM$_{n-1}$

Ranked
answers

*Model complexity*

Models pass progressively smaller sets from one tier to the next, more expensive tiers

1. SciCo: Hierarchical Cross-Document Coreference

2. Document Similarity & Retrieval

3. Literature-Augmented Prediction

**Scientific Knowledge**

*Cognitive Bottleneck*

*Task-guided knowledge discovery*

**Research Tasks**

*Assist researchers* **throughout their tasks**

tomh@allenai.org
cs.huji.ac.il/~tomhope/