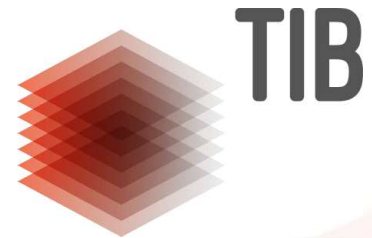LEIBNIZ-INFORMATIONSZENTRUM
TECHNIK UND NATURWISSENSCHAFTEN
UNIVERSITÄTSBIBLIOTHEK

TIB

# NLPContributions: An Annotation Scheme for Machine Reading of Scholarly Contributions in Natural Language Processing Literature
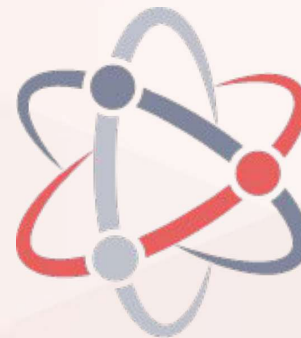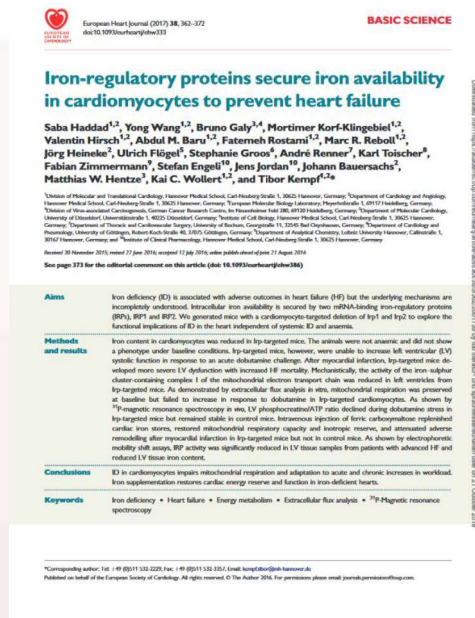
Jennifer D'Souza and Sören Auer
Technische Informationsbibliothek (TIB)
Welfengarten 1B // 30167 Hannover

Leibniz
Gemeinschaft

# What if ...

- The global scientific knowledge base would be more than a document repository
- Scientific information and knowledge would be FAIR also for machines
  - The FAIR data principles are a set of guiding principles in order to make scientific data findable, accessible, interoperable, and reusable in the current digital ecosystem (Wilkinson et al, 2016)
- Currently
  - Findability could be better
  - Assuming OA, accessibility is OK
  - Interoperability and Reusability is non-existent
- The problem: The scholarly communications format is stuck in the last century
  - We have managed to digitize documents that used to be in print
  - While other areas have seen a transformative digitalization

# Our Objective
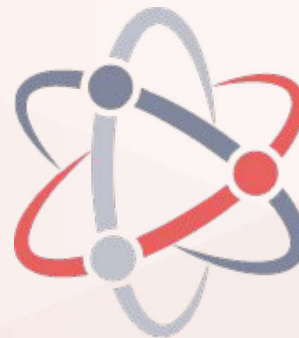
- To foster the *digitalization* of digitized scholarly articles

# Our Objective

- To structure, in a fine-grained manner, knowledge elements from unstructured scholarly articles as a Knowledge Graph

# Our Objective

- **Contributions Scholarly Knowledge. Structured.**
  - Focus on structuring only *contributions* from **natural language processing (NLP) articles**

# Our Objective

- **Contributions Scholarly Knowledge. Structured.**
  - Focus on structuring only *contributions* from **natural language processing (NLP) articles**

- **Devise an annotation methodology**: NLPContributions

# Our Goals

Two-fold:

# Our Goals

Two-fold:

1. perform a **pilot annotation exercise to find a systematic set of patterns of subject-predicate-object statements** for the semantic structuring of scholarly contributions that are more or less generically applicable for NLP articles;

# Our Goals

Two-fold:

1. perform a **pilot annotation exercise to find a systematic set of patterns of subject-predicate-object statements** for the semantic structuring of scholarly contributions that are more or less generically applicable for NLP articles;

2. ingest the resulting pilot annotated data into the Open Research Knowledge Graph (ORKG) infrastructure as a showcase to **automatically process the digitalized scholarly contribution knowledge elements**.

# Our Goals

Two-fold:

1. perform a **pilot annotation exercise to find a systematic set of patterns of subject-predicate-object statements** for the semantic structuring of scholarly contributions that are more or less generically applicable for NLP articles;

2. ingest the resulting pilot annotated data into the Open Research Knowledge Graph (ORKG) infrastructure as a showcase to **automatically process the digitalized scholarly contribution knowledge elements**.
   - The ORKG[1] is a next-generation digital library infrastructure for machine-actionable knowledge content in scholarly articles.

Reference:
1. Jaradeh, Mohamad Yaser, et al. "Open research knowledge graph: next generation infrastructure for semantic scholarly knowledge." *Proceedings of the 10th International Conference on Knowledge Capture*. 2019.

## Plan for the Talk

- **NLPContributions Model**

- **The NLPContributions Annotation Guidelines**

- **Pilot Annotated Dataset Characteristics**

- **NLPContributions in the Open Research Knowledge Graph**

# Plan for the Talk

- **NLPContributions Model**

- **The NLPContributions Annotation Guidelines**

- **Pilot Annotated Dataset Characteristics**

- **NLPContributions in the Open Research Knowledge Graph**

# NLPContributions Model: Characteristics

- Designed for building a knowledge graph

# NLPContributions Model: Characteristics

- Designed for building a knowledge graph

- Not ontologized
  - assumes a bottom-up data-driven design toward ontology discovery

# NLPContributions Model: Characteristics

- Designed for building a knowledge graph

- Not ontologized
  - assumes a bottom-up data-driven design toward ontology discovery

- Has a core skeleton model for top-level knowledge systematization.

# NLPContributions Model: Characteristics

- Designed for building a knowledge graph

- Not ontologized
  - assumes a bottom-up data-driven design toward ontology discovery

- Has a core skeleton model for top-level knowledge systematization.
  - a root node called Contribution,

# NLPContributions Model: Characteristics

- Designed for building a knowledge graph

- Not ontologized
  - assumes a bottom-up data-driven design toward ontology discovery

- Has a core skeleton model for top-level knowledge systematization.
  - a root node called Contribution,
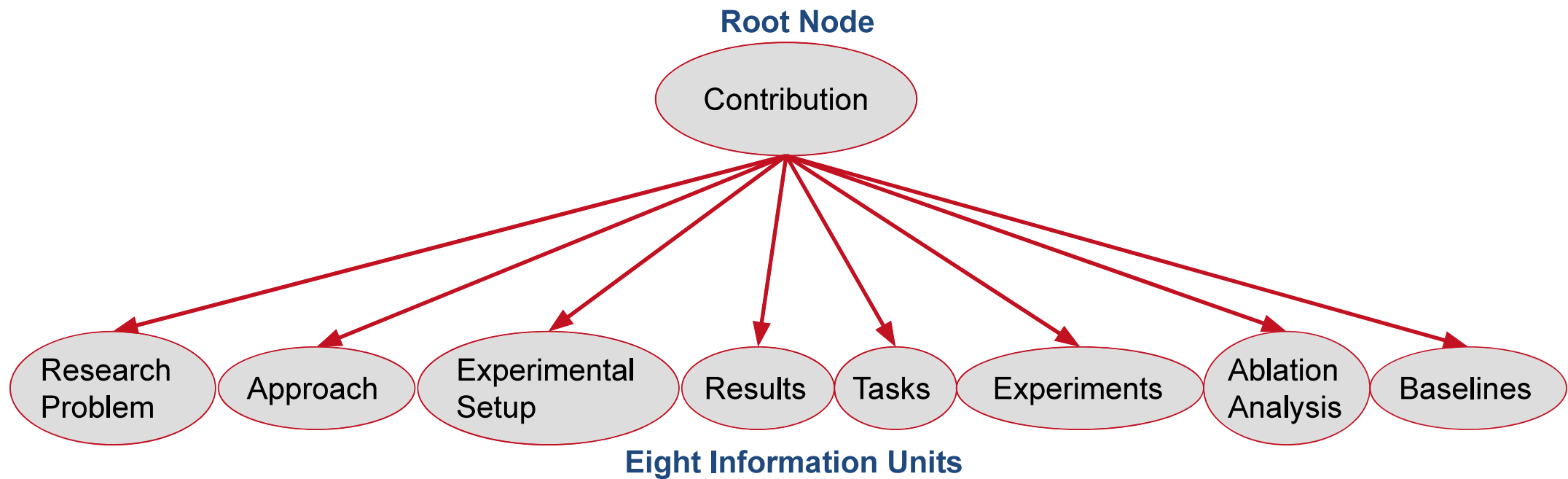  - eight first level nodes representing core information units under which the scholarly contributions data is organized
    - inspired from sectional information organization in scholarly articles

# NLPContributions Model: Core Skeleton

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

    1. ResearchProblem

    2. Approach

    3. ExperimentalSetup

    4. Results

    5. Tasks

    6. Experiments

    7. AblationAnalysis

    8. Baselines

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

1. ResearchProblem

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**1. ResearchProblem**

- research challenge addressed by a contribution

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**1. ResearchProblem**

- research challenge addressed by a contribution
- connected to root by predicate *hasResearchProblem*

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**1. ResearchProblem**

- research challenge addressed by a contribution

- connected to root by predicate *hasResearchProblem*

- E.g., from paper about BioBERT word embeddings, their research problem is 'domain-customization of BERT'

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **1. ResearchProblem**

    - research challenge addressed by a contribution
    - connected to root by predicate *hasResearchProblem*
    - E.g., from paper about BioBERT word embeddings, their research problem is 'domain-customization of BERT'
    - typically found in an article's Title, Abstract and first few paragraphs of the Introduction

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**1. ResearchProblem**

- research challenge addressed by a contribution
- connected to root by predicate *hasResearchProblem*
- E.g., from paper about BioBERT word embeddings, their research problem is 'domain-customization of BERT'
- typically found in an article's Title, Abstract and first few paragraphs of the Introduction
- involves annotating one or more sentences and precisely the research problem phrase boundaries in the sentences

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **2. Approach**

    - solution proposed for the research problem

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

    **2. Approach**

    - solution proposed for the research problem
    - connected to root by predicate *has*

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **2. Approach**

    - solution proposed for the research problem

    - connected to root by predicate *has*

    - alternatively called Model or Method or Architecture or System or Application

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **2. Approach**

  - solution proposed for the research problem

  - connected to root by predicate *has*

  - alternatively called Model or Method or Architecture or System or Application

  - typically found in the article's Introduction section in the context of cue phrases such as "we take the approach," "we propose the model," "our system architecture," or "the method proposed in this paper."
    - exception: the first few lines within the main system description content in the article

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **3. ExperimentalSetup**

  - details about the platform including both hardware (e.g., GPU) and software (e.g., Tensorflow library) for implementing the machine learning solution; and of variables, that determine the network structure (e.g., number of hidden units) and how the network is trained (e.g., learning rate), for tuning the software to the task objective

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **3. ExperimentalSetup**

  - details about the platform including both hardware (e.g., GPU) and software (e.g., Tensorflow library) for implementing the machine learning solution; and of variables, that determine the network structure (e.g., number of hidden units) and how the network is trained (e.g., learning rate), for tuning the software to the task objective
  - connected to root by predicate *has*

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **3. ExperimentalSetup**

  - details about the platform including both hardware (e.g., GPU) and software (e.g., Tensorflow library) for implementing the machine learning solution; and of variables, that determine the network structure (e.g., number of hidden units) and how the network is trained (e.g., learning rate), for tuning the software to the task objective
  - connected to root by predicate *has*
  - found in the sections called Experiment, Experimental Setup, Implementation, Hyperparameters, or Training

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **4. Results**

  - main findings or outcomes reported in the article for the research problem

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**4. Results**

- main findings or outcomes reported in the article for the research problem
- connected to root by predicate *has*

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**4. Results**

- main findings or outcomes reported in the article for the research problem
- connected to root by predicate *has*
- found in an article's Results, Experiments, or Tasks sections
  - while the results are often highlighted in the Introduction, unlike the Approach unit, in this case, we annotate the dedicated, detailed section on Results because results constitute a primary aspect of the contribution.

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**5. Tasks**

- the Approach, particularly in multi-task settings, are tested on more than one task, in which case, all the experimental tasks are listed

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**5. Tasks**

- the Approach, particularly in multi-task settings, are tested on more than one task, in which case, all the experimental tasks are listed
- connected to root by predicate *has*

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**5. Tasks**

- the Approach, particularly in multi-task settings, are tested on more than one task, in which case, all the experimental tasks are listed
- connected to root by predicate *has*
- is an encapsulating information unit
  - can include one or more of the ExperimentalSetup, Hyperparameters, and Results as sub information units

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  1. ResearchProblem
  2. Approach
  3. ExperimentalSetup
  4. Results
  5. Tasks
  6. **Experiments**
  7. **AblationAnalysis**
  8. **Baselines**

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **6. Experiments**

  - is an encapsulating information unit
    - can be a combination of ExperimentalSetup and Results; or lists of Tasks and their Results; or Approach, ExperimentalSetup and Results combined

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **6. Experiments**

  - is an encapsulating information unit
    - can be a combination of ExperimentalSetup and Results; or lists of Tasks and their Results; or Approach, ExperimentalSetup and Results combined
  - particularly relevant in the content of multitask systems such as BERT
    - modeling ExperimentalSetup with Results or Tasks with Results is necessary in such systems since the experimental setup often changes per task producing a different set of results

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **7. AblationAnalysis**

    - describes the performance of components in systems

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**7. AblationAnalysis**

- describes the performance of components in systems
- a form of the results which are relevant to a Contribution

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**7. AblationAnalysis**

- describes the performance of components in systems
- a form of the results which are relevant to a Contribution
- typically found in sections with Ablation in the title, otherwise also in the running text

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

  **8. Baselines**

  - a list of systems that a proposed approach is compared against

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**8. Baselines**

- a list of systems that a proposed approach is compared against

- a form of the results which are relevant to a Contribution

# NLPContributions Model: 8 Information Units

- Inspired from sectional information organization in scholarly articles

**8. Baselines**

- a list of systems that a proposed approach is compared against
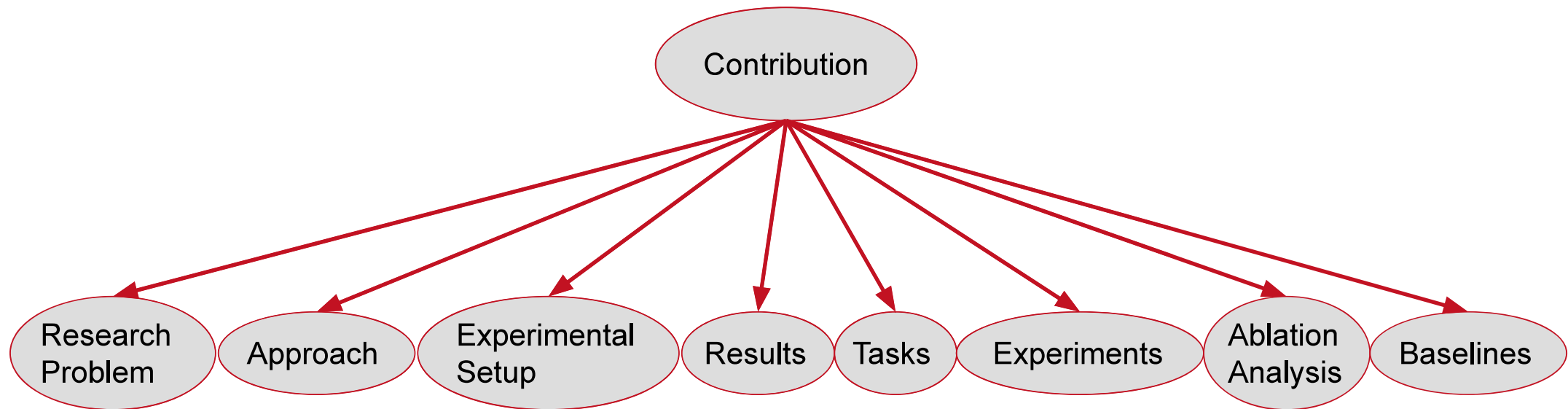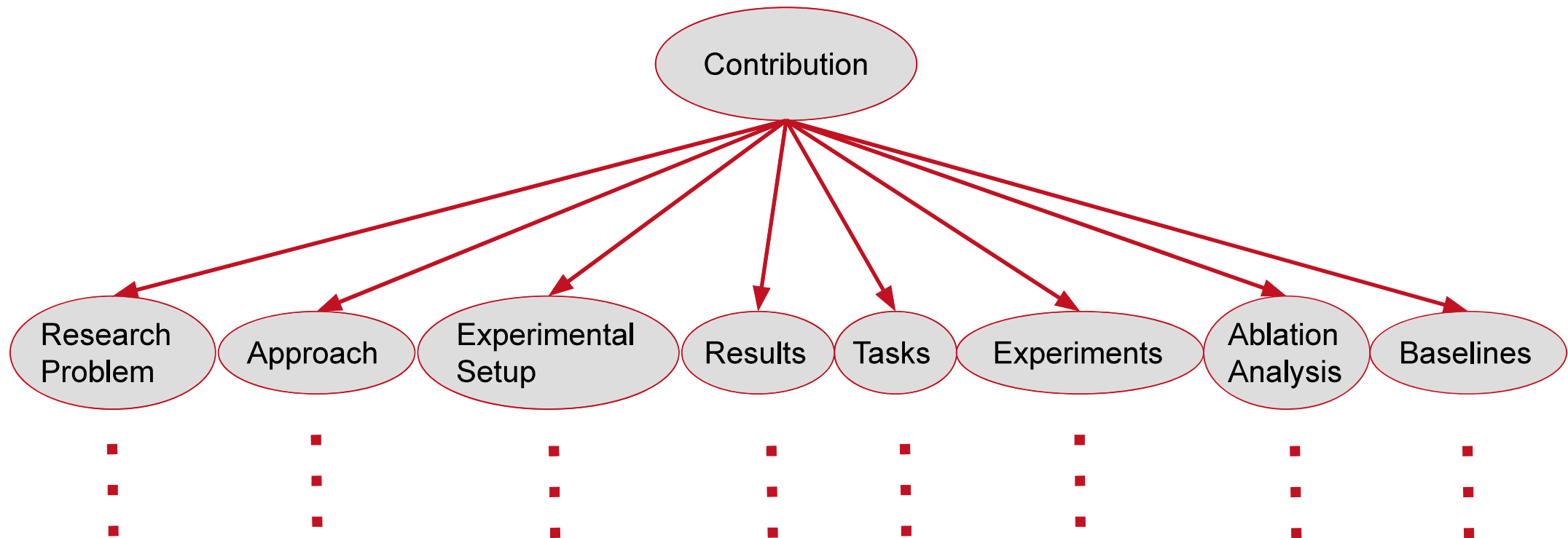- a form of the results which are relevant to a Contribution
- typically found in sections with Baseline in the title, otherwise also in the running text

# NLPContributions Model: 8 Information Units

# NLPContributions Model: Data Elements

# How to: Knowledge Graph building from Unstructured Text

- Given a paragraph(s) of unstructured text
  - identify the elements to model:
    - depends on:
      1. if the knowledge graph has an overarching knowledge theme
      2. or, if the knowledge nodes are to be of a certain type (e.g., scientific entities)
    - 1 subsumes 2
  - For 1 (our contributions-themed model):
    - identify the sentences that reflect the theme
    - identify the knowledge entities and predicates from the sentence of interest to the knowledge theme (e.g., scientific entities)
    - create subject-predicate-object triples toward RDFized KGs
    - ...

# NLPContributions Model: Data Elements

# NLPContributions Model: Data Elements

- **Contribution Sentences**
    - select candidate contribution sentences under each of the aforementioned 3 or more applicable information units (viz., <u>ResearchProblem</u>, <u>Approach</u>, <u>Results</u>, AblationAnalysis, etc.).

# NLPContributions Model: Data Elements

- **Contribution Sentences**
  - select candidate contribution sentences under each of the aforementioned 3 or more applicable information units (viz., ResearchProblem, Approach, Results, AblationAnalysis, etc.).
- **Scientific Term and Predicate Phrases as Knowledge Entities (Graph Nodes)**
  - select phrases with an implicit understanding of whether they take the subject, predicate, or object roles in a per-triple context

# NLPContributions Model: Data Elements

- **Contribution Sentences**
  - select candidate contribution sentences under each of the aforementioned 3 or more applicable information units (viz., ResearchProblem, Approach, Results, AblationAnalysis, etc.).
- **Scientific Term and Predicate Phrases as Knowledge Entities (Graph Nodes)**
  - select phrases with an implicit understanding of whether they take the subject, predicate, or object roles in a per-triple context
- **Create Triples in Contribution Sequences**
  - relating phrases in subject, predicate, and object roles within triples
  - creating contribution sequences by using an object in one triple as the subject in another triple

# NLPContributions Model: Data Elements

**Next**: Example modeling data elements under an information unit

# NLPContributions Model: <u>Approach</u> Data Elements

```json
{
  "has" : {
    "Approach" : {
      "converting questions" : {
        "to (un-interpretable) vectorial representations" : {
          "which require" : "no pre-defined grammars or lexicons",
          "can query" : {
            "any KB" : {
              "independent of" : "schema"
            }
          }
        }
      },
      "from sentence" : "In this paper, we instead take the
        approach of converting questions to (un-interpretable)
        vectorial representations which require no pre-defined
        grammars or lexicons and can query any KB independent of
        its schema."
    }
  }
}
```

**Reference**: Bordes, Antoine, Jason Weston, and Nicolas Usunier. "Open question answering with weakly supervised embedding models." *Joint European conference on machine learning and knowledge discovery in databases*. Springer, Berlin, Heidelberg, 2014.

# NLPContributions Model: ExperimentalSetup Data Elements

```json
{
  "has" : {
    "Experimental setup" : {
      "used" : [
        {
          "BERTBase model" : {
            "pre-trained for" : "1M steps",
            "pre-trained on" : ["English Wikipedia",
              "BooksCorpus"]
          },
          "from sentence" : "We used the BERTBASE model
            pre-trained on English Wikipedia and
            BooksCorpus for 1M steps."
        },
        {
          "NVIDIA V100 (32GB) GPUs" : {
            "used" : {
              "eight" : {
                "for" : "pre-training"
              }
            },
            "from sentence" : "We used eight NVIDIA V100
              (32GB) GPUs for the pre-training."
          }
        }
      ]
    }
  }
}
```

**Reference**: Lee, Jinhyuk, et al. "BioBERT: a pre-trained biomedical language representation model for biomedical text mining." *Bioinformatics* 36.4 (2020): 1234-1240.

# NLPContributions Model: <u>Result</u> Data Elements

```
{
    "CoNLL test set" : {
        "for" : {
            "NER" : {
                "F1-score" : "91.57%"
            }
        },
        "from sentence" : "For NER (Table 7), S-LSTM
            gives an F1-score of 91.57% on the CoNLL
            test set, which is significantly better
            compared with BiLSTMs."
    }
}
```

**Reference**: Zhang, Yue, Qi Liu, and Linfeng Song. "Sentence-State LSTM for Text Representation." *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2018.

# Plan for the Talk

- **NLPContributions Model**

- **The NLPContributions Annotation Guidelines**

- **Pilot Annotated Dataset Characteristics**

- **NLPContributions in the Open Research Knowledge Graph**

## Plan for the Talk

- **NLPContributions Model**
- **The NLPContributions Annotation Guidelines**
- **Pilot Annotated Dataset Characteristics**
- **NLPContributions in the Open Research Knowledge Graph**

# NLPContributions Annotation Guidelines

**Three of Twelve Guidelines:**

# NLPContributions Annotation Guidelines

**Three of Twelve Guidelines:**

1. *How are information unit names selected?* or conversely, *Which of the eight information units does the sentence belong to?*
   - applied name is the one selected based on the closest section title or cue phrase

# NLPContributions Annotation Guidelines

**Three of Twelve Guidelines:**

1. *How are information unit names selected?* or conversely, *Which of the eight information units does the sentence belong to?*
   - applied name is the one selected based on the closest section title or cue phrase

2. *Inferring Predicates*
   - from running text or from the closed class set {"has", "on", "by", "for", "has value", "has description", "based on", "called"}

# NLPContributions Annotation Guidelines

**Three of Twelve Guidelines:**

1. *How are information unit names selected?* or conversely, *Which of the eight information units does the sentence belong to?*
   - applied name is the one selected based on the closest section title or cue phrase

2. *Inferring Predicates*
   - from running text or from the closed class set {"has", "on", "by", "for", "has value", "has description", "based on", "called"}

3. *How are lists modeled within contribution sequences?*
   - list items are treated just as sentences

# Plan for the Talk

- **NLPContributions Model**
- **The NLPContributions Annotation Guidelines**
- **Pilot Annotated Dataset Characteristics**
- **NLPContributions in the Open Research Knowledge Graph**

# Plan for the Talk

- **NLPContributions Model**
- **The NLPContributions Annotation Guidelines**
- **Pilot Annotated Dataset Characteristics**
- **NLPContributions in the Open Research Knowledge Graph**

# Pilot Annotated Dataset

- **Dataset**

  - A collection of scholarly articles downloaded from https://paperswithcode.com/

    - represents papers in AI at large

# Pilot Annotated Dataset

- **Dataset**
  - A collection of scholarly articles downloaded from https://paperswithcode.com/
    - represents papers in AI at large
  - Randomly selected 50 NLP papers
    - <u>Aim</u>: create a representative dataset
    - select a distribution of 10 papers across five different NLP research tasks:
      - machine translation, named entity recognition, question answering, relation classification, and text classification.

# Pilot Annotated Dataset

- **Dataset**

  - A collection of scholarly articles downloaded from https://paperswithcode.com/

    - represents papers in AI at large

  - Randomly selected 50 NLP papers

    - <u>Aim</u>: create a representative dataset

    - select a distribution of 10 papers across five different NLP research tasks:

      - machine translation, named entity recognition, question answering, relation classification, and text classification.

- **Annotation Tools**

  - https://jsoneditoronline.org/ - For JSON syntax checks

  - https://www.orkg.org/ - As a litmus test for contributions-themed KG and as the Digital Library infrastructure to populate with the annotated KGs

# Pilot Annotated Dataset Characteristics

- Total of 2631 triples (avg. of 52 triples per article)

- Data elements: 1033 unique subjects, 843 unique predicates, and 2182 unique objects

# Pilot Annotated Dataset Characteristics

- Total of 2631 triples (avg. of 52 triples per article)

- Data elements: 1033 unique subjects, 843 unique predicates, and 2182 unique objects

|  | MT | NER | QA | RC | TC |
|---|---|---|---|---|---|
| *Subject* | 259 | 209 | 203 | 228 | 221 |
| *Predicate* | 243 | 220 | 187 | 201 | 252 |
| *Object* | 471 | 434 | 515 | 455 | 459 |
| Total | 502 | 473 | 497 | 544 | 504 |

MT: machine translation; NER: named entity recognition; QA: question answering; RC: relation classification; TC: text classification

# Pilot Annotated Dataset Characteristics

- Total of 2631 triples (avg. of 52 triples per article)

- Data elements: 1033 unique subjects, 843 unique predicates, and 2182 unique objects

| | MT | NER | QA | RC | TC |
|---|---|---|---|---|---|
| *Subject* | 259 | 209 | 203 | 228 | 221 |
| *Predicate* | 243 | 220 | 187 | 201 | 252 |
| *Object* | **471** | **434** | **515** | **455** | **459** |
| Total | 502 | 473 | 497 | 544 | 504 |

MT: machine translation; NER: named entity recognition; QA: question answering; RC: relation classification; TC: text classification

# Pilot Annotated Dataset Characteristics

- Total of 2631 triples (avg. of 52 triples per article)

- Data elements: 1033 unique subjects, 843 unique predicates, and 2182 unique objects

| | MT | NER | QA | RC | TC |
|---|---|---|---|---|---|
| *Subject* | 259 | 209 | 203 | 228 | 221 |
| *Predicate* | 243 | 220 | 187 | 201 | 252 |
| *Object* | 471 | 434 | 515 | 455 | 459 |
| Total | 502 | 473 | 497 | **544** | 504 |

MT: machine translation; NER: named entity recognition; QA: question answering; RC: relation classification; TC: text classification

# Pilot Annotated Dataset Characteristics

- Total of 2631 triples (avg. of 52 triples per article)

- Data elements: 1033 unique subjects, 843 unique predicates, and 2182 unique objects

|  | MT | NER | QA | RC | TC |
|---|---|---|---|---|---|
| *Subject* | 259 | 209 | 203 | 228 | 221 |
| *Predicate* | 243 | 220 | 187 | 201 | 252 |
| *Object* | 471 | 434 | 515 | 455 | 459 |
| Total | 502 | **473** | 497 | 544 | 504 |

MT: machine translation; NER: named entity recognition; QA: question answering; RC: relation classification; TC: text classification
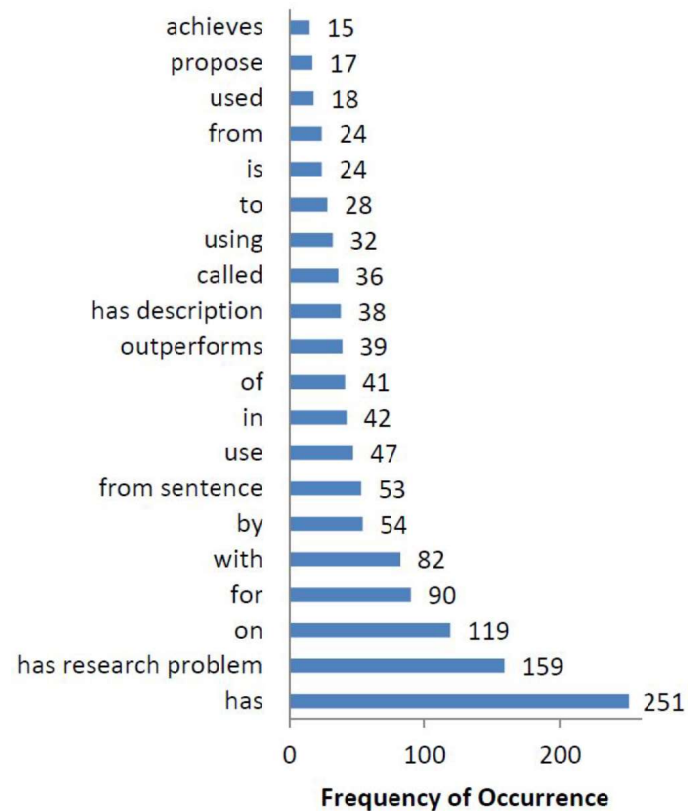
# Pilot Annotated Dataset Characteristics

- Total of 2631 triples (avg. of 52 triples per article)

- Data elements: 1033 unique subjects, 843 unique predicates, and 2182 unique objects

| | MT | NER | QA | RC | TC |
|---|---|---|---|---|---|
| *Subject* | 259 | 209 | 203 | 228 | 221 |
| *Predicate* | 243 | 220 | 187 | 201 | 252 |
| *Object* | 471 | 434 | 515 | 455 | 459 |
| Total | 502 | 473 | 497 | 544 | 504 |

MT: machine translation; NER: named entity recognition; QA: question answering; RC: relation classification; TC: text classification

# Pilot Annotated Dataset Characteristics

- Total of 2631 triples (avg. of 52 triples per article)

- Data elements: 1033 unique subjects, 843 unique predicates, and 2182 unique objects

## Plan for the Talk

- **NLPContributions Model**
- **The NLPContributions Annotation Guidelines**
- **Pilot Annotated Dataset Characteristics**
- **NLPContributions in the Open Research Knowledge Graph**

# Plan for the Talk

- **NLPContributions Model**
- **The NLPContributions Annotation Guidelines**
- **Pilot Annotated Dataset Characteristics**
- **NLPContributions in the Open Research Knowledge Graph**

Cornell University

arXiv.org > cs > arXiv:1809.10185

**Computer Science > Computation and Language**

# Graph Convolution over Pruned Dependency Trees Improves Relation Extraction

Yuhao Zhang, Peng Qi, Christopher D. Manning

Dependency trees help relation extraction models capture long-range relations between words. However, existing dependency-based models either neglect crucial information (e.g., negation) by pruning the dependency trees too aggressively, or are computationally inefficient because it is difficult to parallelize over different tree structures. We propose an extension of graph convolutional networks that is tailored for relation extraction, which pools information over arbitrary dependency structures efficiently in parallel. To incorporate relevant information while maximally removing irrelevant content, we further apply a novel pruning strategy to the input trees by keeping words immediately around the shortest path between the two entities among which a relation might hold. The resulting model achieves state-of-the-art performance on the large-scale TACRED dataset, outperforming existing sequence and dependency-based neural models. We also show through detailed analysis that this model has complementary strengths to sequence models, and combining them further improves the state of the art.

**Bibliographic data**
[Enable Bibex (What is Bibex?)]

*Which authors of this paper are endorsers? | Disable MathJax (What is MathJax?)*

**Download:**
- PDF
- Other formats
(license)

Current browse context:
**cs.CL**
< prev | next >
new | recent | 1809
Change to browse by:
cs

**References & Citations**
- NASA ADS
- Google Scholar
- Semantic Scholar

**DBLP** - CS Bibliography
listing | bibtex
Yuhao Zhang
Peng Qi
Christopher D. Manning

**Export citation**

Bookmark

## Abstract

Dependency trees help relation extraction models capture long-range relations between words. However, existing dependency-based models either neglect crucial information (e.g., negation) by pruning the dependency trees too aggressively, or are computationally inefficient because it is difficult to parallelize over different tree structures. We propose an extension of graph convolutional networks that is tailored for relation extraction, which pools information over arbitrary dependency structures efficiently in parallel. To incorporate relevant information while maximally removing irrelevant content, we further apply a novel pruning strategy to the input trees by keeping words immediately around the shortest path between the two entities among which a relation might hold. The resulting model achieves state-of-the-art performance on the large-scale TACRED dataset, outperforming existing sequence and dependency-based neural models. We also show through detailed analysis that this model has complementary strengths to sequence models, and combining them further improves the state of the art.

## Graph Convolution over Pruned Dependency Trees Improves Relation Extraction

2018   Information Science   Yuhao Zhang   Peng Qi   Christopher D Manning

Published in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*

### Contribution 1

Research problems      ☐ Add to comparison

Relation extraction

Contribution data

| Has | |
|-----|-----|
| | Ablation analysis |
| | Baseline Models |
| | Model |
| | Results |

Accessible at https://www.orkg.org/orkg/paper/R44287

## Abstract

Dependency trees help relation extraction models capture long-range relations between words. However, existing dependency-based models either neglect crucial information (e.g., negation) by pruning the dependency trees too aggressively, or are computationally inefficient because it is difficult to parallelize over different tree structures. We propose an extension of graph convolutional networks that is tailored for relation extraction, which pools information over arbitrary dependency structures efficiently in parallel. To incorporate relevant information while maximally removing irrelevant content, we further apply a novel pruning strategy to the input trees by keeping words immediately around the shortest path between the two entities among which a relation might hold. The resulting model achieves state-of-the-art performance on the large-scale TACRED dataset, outperforming existing sequence and dependency-based neural models. We also show through detailed analysis that this model has complementary strengths to sequence models, and combining them further improves the state of the art.

### Graph Convolution over Pruned Dependency Trees Improves Relation Extraction

2018 | Information Science | Yuhao Zhang | Peng Qi | Christopher D Manning

Published in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*

modeled 5 NLPContributions information units

**Contribution 1**

Research problems

Relation extraction

☐ Add to comparison

Contribution data

| Has | Ablation analysis |
| | Baseline Models |
| | Model |
| | Results |

Accessible at https://www.orkg.org/orkg/paper/R44287

## Abstract

Dependency trees help relation extraction models capture long-range relations between words. However, existing dependency-based models either neglect crucial information (e.g., negation) by pruning the dependency trees too aggressively, or are computationally inefficient because it is difficult to parallelize over different tree structures. We propose an extension of graph convolutional networks that is tailored for relation extraction, which pools information over arbitrary dependency structures efficiently in parallel. To incorporate relevant information while maximally removing irrelevant content, we further apply a novel pruning strategy to the input trees by keeping words immediately around the shortest path between the two entities among which a relation might hold. The resulting model achieves state-of-the-art performance on the large-scale TACRED dataset, outperforming existing sequence and dependency-based neural models. We also show through detailed analysis that this model has complementary strengths to sequence models, and combining them further improves the state of the art.

### Graph Convolution over Pruned Dependency Trees Improves Relation Extraction

2018    Information Science    Yuhao Zhang    Peng Qi    Christopher D Manning

Published in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*

#### Contribution 1

Research problems

Relation extraction

Contribution data

| Has | |
|-----|-----|
| | Ablation analysis |
| | Baseline Models |
| | Model |
| | Results |

☐ Add to comparison

Accessible at https://www.orkg.org/orkg/paper/R44287

## Abstract

Dependency trees help relation extraction models capture long-range relations between words. However, existing dependency-based models either neglect crucial information (e.g., negation) by pruning the dependency trees too aggressively, or are computationally inefficient because it is difficult to parallelize over different tree structures. We propose an extension of graph convolutional networks that is tailored for relation extraction, which pools information over arbitrary dependency structures efficiently in parallel. To incorporate relevant information while maximally removing irrelevant content, we further apply a novel pruning strategy to the input trees by keeping words immediately around the shortest path between the two entities among which a relation might hold. The resulting model achieves state-of-the-art performance on the large-scale TACRED dataset, outperforming existing sequence and dependency-based neural models. We also show through detailed analysis that this model has complementary strengths to sequence models, and combining them further improves the state of the art.

## Graph Convolution over Pruned Dependency Trees Improves Relation Extraction

2018    Information Science    Yuhao Zhang    Peng Qi    Christopher D Manning

Published in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*

### Contribution 1

Research problems

Relation extraction

Contribution data

| Has | |
|-----|-----|
| | Ablation analysis |
| | Baseline Models |
| | Model |
| | Results |

Accessible at https://www.orkg.org/orkg/paper/R44287

## Abstract

Dependency trees help relation extraction models capture long-range relations between words. However, existing dependency-based models either neglect crucial information (e.g., negation) by pruning the dependency trees too aggressively, or are computationally inefficient because it is difficult to parallelize over different tree structures. We propose an extension of graph convolutional networks that is tailored for relation extraction, which pools information over arbitrary dependency structures efficiently in parallel. To incorporate relevant information while maximally removing irrelevant content, we further apply a novel pruning strategy to the input trees by keeping words immediately around the shortest path between the two entities among which a relation might hold. The resulting model achieves state-of-the-art performance on the large-scale TACRED dataset, outperforming existing sequence and dependency-based neural models. We also show through detailed analysis that this model has complementary strengths to sequence models, and combining them further improves the state of the art.

## Graph Convolution over Pruned Dependency Trees Improves Relation Extraction

📅 2018   ≡ Information Science   👤 Yuhao Zhang   👤 Peng Qi   👤 Christopher D Manning

Published in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*

### Contribution 1

**Research problems**

Relation extraction

**Contribution data**

| Has | |
|-----|------------------|
|     | Ablation analysis |
|     | Baseline Models  |
|     | Model            |
|     | Results          |

Accessible at https://www.orkg.org/orkg/paper/R44287

# Plan for the Talk

- **NLPContributions Model**
- **The NLPContributions Annotation Guidelines**
- **Pilot Annotated Dataset Characteristics**
- **NLPContributions in the Open Research Knowledge Graph**

# Conclusion: Takeaways

- Scholarly work can be realized as expressions other than an article

  - We proposed the **NLPContributions annotation model to create contributions-themed knowledge graphs**
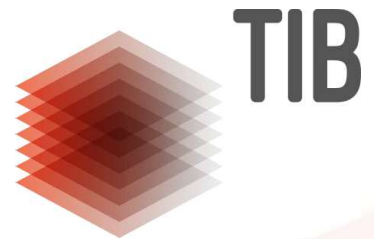
# Conclusion: Takeaways

- Scholarly work can be realized as expressions other than an article

  - We proposed the **NLPContributions annotation model to create contributions-themed knowledge graphs**

- In a pilot annotation exercise we have annotated 50 articles by the NLPContributions scheme as a practical demonstration of feasibility of the annotation task

  - Available online at https://doi.org/10.25835/0019761

# Conclusion: Takeaways

- Scholarly work can be realized as expressions other than an article

  - We proposed the **NLPContributions annotation model to create contributions-themed knowledge graphs**

- In a pilot annotation exercise we have annotated 50 articles by the NLPContributions scheme as a practical demonstration of feasibility of the annotation task

  - Available online at https://doi.org/10.25835/0019761

- The NLPContributions annotation scheme can be leveraged to annotate a larger dataset (of a few hundreds of articles)

  - Train machine-learning-based automated machine readers to annotate tens of thousands of articles for contributions-based KG data which is humanly impossible to do

LEIBNIZ-INFORMATIONSZENTRUM
TECHNIK UND NATURWISSENSCHAFTEN
UNIVERSITÄTSBIBLIOTHEK

TIB

# Thank you for your attention!

Questions?

Leibniz
Gemeinschaft