# AllenNLP

An open-source NLP research library, built on PyTorch

**Matt Gardner**, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, Luke Zettlemoyer

… and the list keeps growing

AI2 ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# AllenNLP

## An open-source NLP research library, built on PyTorch

- Made to make NLP research easy

- Abstractions designed for NLP

- Configuration-driven experiments for doing good science

- Reference implementations and demos for a lot of tasks

- An active community

# AllenNLP

An open-source NLP research library, built on PyTorch

- Clean implementations of state-of-the-art models for virtually any NLP task

    - Dramatically lowers barrier to entry for doing NLP research

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# AllenNLP

An open-source NLP research library, built on PyTorch

- Live demos of all of these models that you can play around with and break

    - Mark Johnson used these yesterday to demonstrate a point about linguistics

    - Plenty of usage in twitter conversations about NLP models

ALLEN INSTITUTE
*for* ARTIFICIAL INTELLIGENCE

# AllenNLP

An open-source NLP research library, built on PyTorch

- Allows for more fundamental, wide-ranging NLP research

  - Test your idea on all NLP tasks, instead of architecture engineering on a single task

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# AllenNLP

An open-source NLP research library, built on PyTorch

- We're not there yet, but with a little help, we could be

    - We're a small team, we can't do everything

    - One possibility: make a model re-implementation a class project in your intro course

    - Issues to solve around control and credit assignment

ALLEN INSTITUTE
*for* ARTIFICIAL INTELLIGENCE

# The ACL Anthology
## Current State and Future Directions

Daniel Gildea, **Min-Yen Kan,** Nitin Madnani, Christoph Teichmann, Martin Villalba

# What is this presentation **about**?



- Summarize the history and current state of efforts related to the Anthology

- Illustrate the challenges of maintaining a community Project

- Invite the community to extend the capabilities of the Anthology

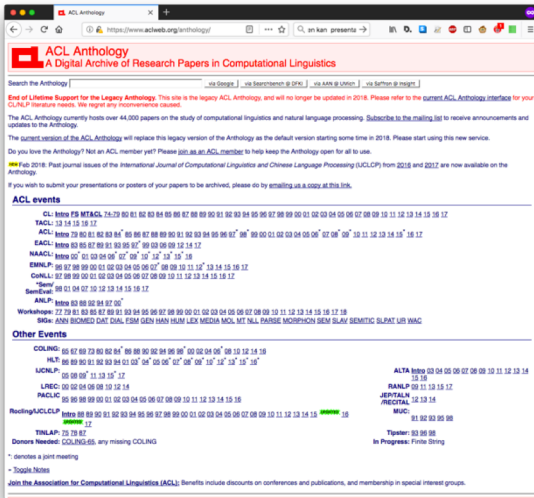- Call you to join the Anthology team

# The Anthology in **summary**





- Open access service for all ACL-Sponsored publications

- Also hosts posters and additional data

- Paper search and author pages

- 45K papers and 4.5K daily hits

- Open Source

- Maintained by volunteers

- New papers added in collaboration with proceedings editors

# A brief **History** of the Anthology



Steven Bird    Min-Yen Kan

- Proposed in 2001 by Steven Bird

- First version online in 2002, with Steven Bird as editor

- Min-Yen Kan becomes the new editor in 2008

- A new version of the Anthology with extra functionality is released in 2012

- Hosting of the Anthology moves from the National University of Singapore to Saarland University

# How to **Future-proof** the Anthology

## Challenges

- Limited resources for day-to-day code maintenance
- Dependencies become outdated
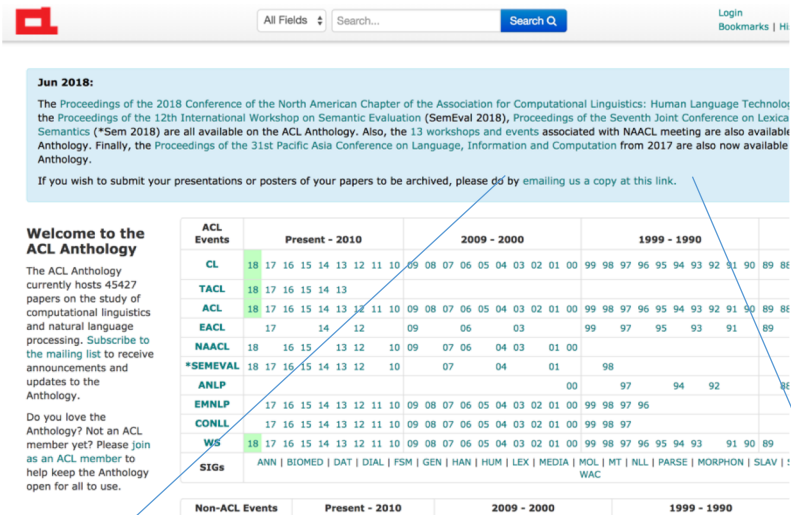- Maintainer churn

## Solutions

- Docker container for easier set-up and sandboxing
- Collaborative documentation efforts to ease onboarding
- Migration plan on the pipeline, including upgrades and test cases

# **Upcoming** major steps



Backlog
- Upgrade and/or migrate outdated dependencies
- Full text search over uploaded papers

In research
- Full test coverage and consistency checks

In progress
- Docker image for releases
- Add a staging server

Done
- Add index support for popular search engines
- Document and update the installation process
- Add a test server

- Hosting the Anthology within the main ACL website

- Recruit a new Anthology editor

- (possibly) pay for extra support
  for the Anthology

# **Exercise**: Importing of your slides



- We import slides, datasets, videos from your own

- Currently done by email (try it yourself! yes, now)

- Better workflow: pull request against the Anthology XML (à la csrankings.org)

# Possible **future directions**

- Contains useful information both *for* CL researchers and *about* CL researchers. Useful for identifying suitable reviewers.

- Move focus from day-to-day operations towards development

- Establish a network of mirrors

- Host anonymized pre-prints

- Comments? Questions?

- Ideas for future directions?

- Interested in joining the Anthology team?

# Come and visit our **poster**

# In OSS we trust

- ▶ Users trust OSS packages to provide good stop word lists
- ▶ Maintainers might not have given it much thought
- ▶ Lists are adapted from each other
- ▶ Lists include surprises and inconsistencies

# Scikit-learn stop words

- We don't know how our 'english' list was constructed
- but spaCy and Gensim use a similar list

- Has typos: fify corrected to fifty in 2015
- Surprising inclusions: computer (removed 2011); system; cry
- Surprising omissions: seven, does
- Inconsistent with our default tokenizer: ve isn't stopped

# Looking beyond Scikit-learn

- We analyse @igorbrigadir's collection of English stop word lists
- We compare the contents of 52 lists

# Looking beyond Scikit-learn

- We analyse @igorbrigadir's collection of English stop word lists
- We compare the contents of 52 lists
- We identify some surprises and inconsistencies

# We can improve how we provide stop lists

- Better documentation
- Adapt the list to the NLP pipeline
- Tools for quality control
- Tools for automatic list construction

# The risk of sub-optimal use of Open Source NLP Software

UKB is inadvertently state-of-the-art in knowledge-based WSD

Eneko Agirre    Oier López de Lacalle    **Aitor Soroa**

NLP-OSS Workshop, July 2018

IXA NLP group, UPV/EHU

## Introduction

- UKB is a collection of programs for WSD
- Graph-based, exploits relations of KB
  - using the Personalized PageRank algorithm
- First released on 2009, attained SOA results
- Free software (GPLv3 license)

## Many uses

- Named Entity disambigiation
- Disambiguation of medical entities
- Word similarity
- Create knowledge-based word embeddings

## Parameters

- UKB contains many parameters

- UKB contains many parameters
  - KB relations
    - Which relations to use
    - Use relation weights

**Parameters**

- UKB contains many parameters
  - KB relations
    - Which relations to use
    - Use relation weights
  - Dictionary
    - Use sense frequencies

## Parameters

- UKB contains many parameters
  - KB relations
    - Which relations to use
    - Use relation weights
  - Dictionary
    - Use sense frequencies
  - Graph algorithms
    - Whole graph: *ppr*, *ppr_w2w*
    - Subgraph: *dfs*, *bfs*
    - Aproximation algorithms: *nibble*
    - Each contains its own hyper-parameters

## Parameters

- UKB contains many parameters
  - KB relations
    - Which relations to use
    - Use relation weights
  - Dictionary
    - Use sense frequencies
  - Graph algorithms
    - Whole graph: *ppr*, *ppr_w2w*
    - Subgraph: *dfs*, *bfs*
    - Aproximation algorithms: *nibble*
    - Each contains its own hyper-parameters
  - Input pre-processing
    - Context of at least 20 words

## UKB parameters

- Default parameters are sub-optimal
  - they do not obtain best results
- Two main reasons:
  - remain purely unsupervised
  - speed trade-off
- Some authors reported results with the default sub-optimal parameters

|  | All | S2 | S3 | S07 | S13 | S15 |
|---|---|---|---|---|---|---|
| UKB (elsewhere)†‡ | 57.5 | 60.6 | 54.1 | 42.0 | 59.0 | 61.2 |
| UKB (this work) | **67.3** | 68.8 | 66.1 | 53.0 | **68.8** | **70.3** |

## UKB parameters

- Default parameters are sub-optimal
  - they do not obtain best results
- Two main reasons:
  - remain purely unsupervised
  - speed trade-off
- Some authors reported results with the default sub-optimal parameters

|  | All | S2 | S3 | S07 | S13 | S15 |
|---|---|---|---|---|---|---|
| UKB (elsewhere)†‡ | 57.5 | 60.6 | 54.1 | 42.0 | 59.0 | 61.2 |
| UKB (this work) | **67.3** | 68.8 | 66.1 | 53.0 | **68.8** | **70.3** |
| Chaplot and Sakajhutdinov (2018) ‡ | 66.9 | **69.0** | **66.9** | 55.6 | 65.3 | 69.6 |
| Babelfy (Moro et al., 2014)† | 65.5 | 67.0 | 63.5 | 51.6 | 66.4 | 70.3 |
| MFS | 65.2 | 66.8 | 66.2 | 55.2 | 63.0 | 67.8 |
| Basile et al. (2014)† | 63.7 | 63.0 | 63.7 | **56.7** | 66.2 | 64.6 |
| Banerjee and Pedersen (2003)† | 48.7 | 50.6 | 44.5 | 32.0 | 53.6 | 51.0 |

## Conclusion

- Default parameters are very important
  - extremely important to include precise instructions and optimal default parameters.
- If possible, include end-to-end scripts to automatically reproduce results
- Most recent version (3.0)
  - parameters are now optimal
  - contains scripts for reproducing results on WSD Evaluation Framework (Raganato et al, 2017)
- UKB still SOA among KB methods

Thank you