



# Better Together: Text + Context

Jan 18, 2024

<https://www.semanticscholar.org/api-gallery/better-together>

The screenshot shows the Semantic Scholar API interface. At the top, there's a search bar with the placeholder "Search over 214 million papers from all fields of science". Below the search bar, there's a navigation bar with links for "Semantic Scholar API", "Overview", "Tutorial", "Documentation", "Gallery", and "Cite the Paper". The main content area features a section titled "Better Together" with a sub-section "Find similar papers in Semantic Scholar". This section includes a brief description of how the tool helps find similar papers based on document embeddings. To the right, there's a sidebar for "Kenneth Church" with links to his GitHub, Author Page, and Homepage, and a red box highlights the "Go To Project" button.

Better Together

## Find similar papers in Semantic Scholar

Looking for papers relevant to a specific paper of interest? Better Together built several different embeddings of documents from Semantic Scholar to help you find similar papers.

Better Together

Kenneth Church

@kchurch4

GitHub

Author Page

Homepage

Go To Project

It is standard practice to represent documents as embeddings. Embeddings based on deep nets (BERT) capture text and other embeddings based on node2vec and GNNs (graph neural nets) capture citation graphs.

We evaluate these embeddings and show that combinations of text and citations are better than either by itself on standard benchmarks of downstream tasks. Embeddings are available for a range of applications: ranked retrieval, recommender systems and routing papers to reviewers.

# Find Similar Papers

**Search by Paper**

Embedding:

Sort by Scores from ProNE Embedding

Limit: 20

Deepwalk

**Search by Author**

Embedding:

Sort by Scores from ProNE Embedding

Limit: 20

Search for Author (Author id or name + <enter>)

[Help](#) [GitHub](#) [Final Report \(YouTube\)](#) [JSALT-2023](#) [Comments Appreciated](#) [BETA Version](#)



# Recommendation Tasks

More than word overlap

- For readers
  - What should I read?
- For authors
  - What should I cite?
- For conference organizers
  - Who should review what?
- Many systems focus on ``relevance'' (word overlap)
  - But we also want credibility
  - Don't recommend papers that are buzz-word compliant
    - But not worth reading
    - (Most papers are never cited)

# Paper Search Results

[Home Page](#) [More Results](#) [S2 Recommendations](#) [bibtex](#)

CorpusId: CorpusId:3051291

Limit: 5

Sort by: score

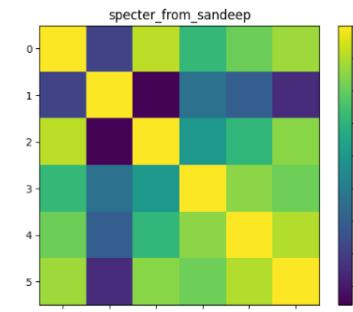
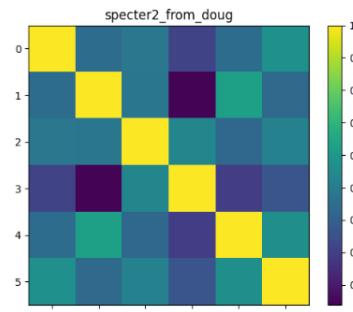
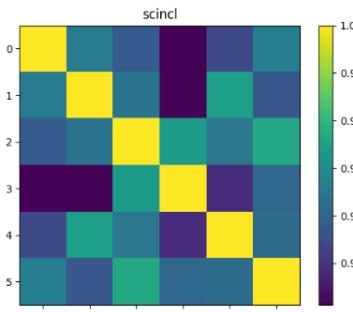
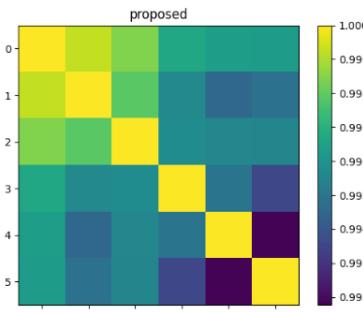
Embedding: proposed

# Proposed Method

## Table of Contents

### 1. [DeepWalk: online learning of social representations](#)

#### Paper: DeepWalk: online learning of social representations



Top

score	citationCount	Paper	Authors	year	More like this	<a href="#">proposed</a>	<a href="#">scinl</a>	<a href="#">specter2_from_doug</a>	<a href="#">specter_from_sandeep</a>
8182		<a href="#">DeepWalk: online learning of social representations</a>	<a href="#">Bryan Perozzi</a> , <a href="#">Rami Al-Rfou</a> , <a href="#">S. Skiena</a>	2014	<a href="#">similar to this</a>	1.0	1.0	1.0	1.0
0.999	8779	<a href="#">node2vec: Scalable Feature Learning for Networks</a>	<a href="#">Aditya Grover</a> , <a href="#">J. Leskovec</a>	2016	<a href="#">similar to this</a>	0.999	0.931	0.944	0.864
0.998	4770	<a href="#">LINE: Large-scale Information Network Embedding</a>	<a href="#">Jian Tang</a> , <a href="#">Meng Qu</a> , ..., <a href="#">Q. Mei</a>	2015	<a href="#">similar to this</a>	0.998	0.916	0.948	0.983
0.997	1511	<a href="#">A Comprehensive Survey of Graph Embedding: Problems, Techniques, and Applications</a>	<a href="#">Hongyun Cai</a> , <a href="#">V. Zheng</a> , <a href="#">K. Chang</a>	2017	<a href="#">similar to this</a>	0.997	0.883	0.932	0.943
0.996	1758	<a href="#">metapath2vec: Scalable Representation Learning for Heterogeneous Networks</a>	<a href="#">Yuxiao Dong</a> , <a href="#">N. Chawla</a> , <a href="#">A. Swami</a>	2017	<a href="#">similar to this</a>	0.996	0.908	0.944	0.961
0.996	806	<a href="#">Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec</a>	<a href="#">J. Qiu</a> , <a href="#">Yuxiao Dong</a> , ..., <a href="#">Jie Tang</a>	2017	<a href="#">similar to this</a>	0.996	0.933	0.957	0.975

# Paper Search Results

[Home Page](#) [More Results](#) [Proposed bibtex](#)

CorpusId: CorpusId:3051291

Limit: 5

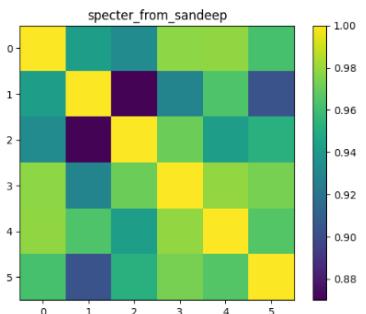
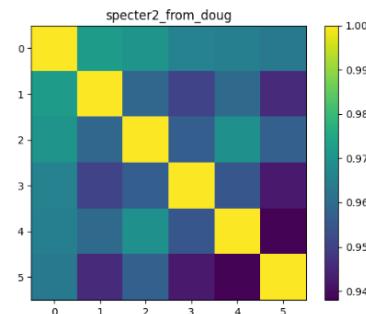
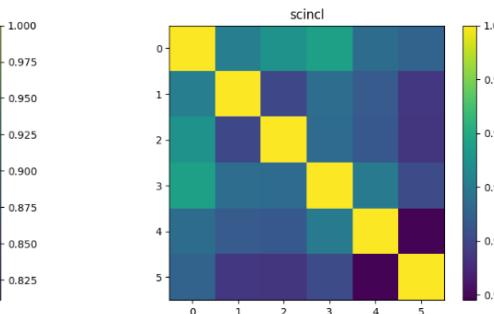
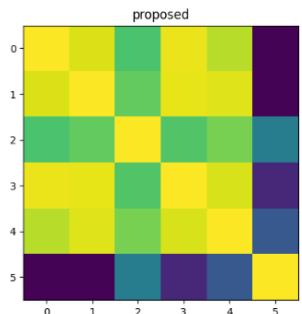
Sort by: score

Embedding: s2\_recommenations

## Table of Contents

### 1. [DeepWalk: online learning of social representations](#)

#### Paper: DeepWalk: online learning of social representations



Top

score	citationCount	Paper
8182	1	<a href="#">DeepWalk: online learning of social representations</a>
5	253	<a href="#">Max-Margin DeepWalk: Discriminative Learning of Network Representation</a>
4	5	<a href="#">SimWalk: Learning network latent representations with social relation similarity</a>
3	165	<a href="#">Don't Walk, Skip!: Online Learning of Multi-scale Network Embeddings</a>
2	114	<a href="#">Discriminative Deep Random Walk for Network Classification</a>
1	5	<a href="#">Learning distributed representations for large-scale dynamic social networks</a>

# Recommendations from Semantic Scholar

Authors	year	More like this	<a href="#">proposed</a>	<a href="#">scinl</a>	<a href="#">specter2</a> <a href="#">from_doug</a>	<a href="#">specter</a> <a href="#">from_sandeep</a>
<a href="#">Bryan Perozzi</a> , <a href="#">Rami Al-Rfou</a> , <a href="#">S. Skiena</a>	2014	<a href="#">similar to this</a>	1.0	1.0	1.0	1.0
<a href="#">Cunchao Tu</a> , <a href="#">Weicheng Zhang</a> , ..., <a href="#">Maosong Sun</a>	2016	<a href="#">similar to this</a>	0.99	0.942	0.972	0.943
<a href="#">Shicheng Cui</a> , <a href="#">Bin Xia</a> , ..., <a href="#">Hong Zhang</a>	2017	<a href="#">similar to this</a>	0.947	0.951	0.97	0.933
<a href="#">Bryan Perozzi</a> , <a href="#">Vivek Kulkarni</a> , ..., <a href="#">S. Skiena</a>	2016	<a href="#">similar to this</a>	0.995	0.956	0.966	0.978
<a href="#">Juzheng Li</a> , <a href="#">Jun Zhu</a> , <a href="#">Bo Zhang</a>	2016	<a href="#">similar to this</a>	0.979	0.934	0.965	0.978
<a href="#">Aakas Zhiyuli</a> , <a href="#">Xun Liang</a> , <a href="#">Zhiming Xu</a>	2017	<a href="#">similar to this</a>	0.811	0.931	0.963	0.962

# Prior: More citations → More credible

**Proposed**

Spectral Clustering  
of Citation Graph

score	citationCount	Paper
8182		<a href="#">DeepWalk: online learning of social representations</a>
0.999	8779	<a href="#">node2vec: Scalable Feature Learning for Networks</a>
0.998	4770	<a href="#">LINE: Large-scale Information Network Embedding</a>
0.997	1511	<a href="#">A Comprehensive Survey of Graph Embedding: Problems, Techniques, and Applications</a>
0.996	1758	<a href="#">metapath2vec: Scalable Representation Learning for Heterogeneous Networks</a>
0.996	806	<a href="#">Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec</a>

**Semantic Scholar**

BERT-like  
Deep Net

score	citationCount	Paper
8182		<a href="#">DeepWalk: online learning of social representations</a>
5	253	<a href="#">Max-Margin DeepWalk: Discriminative Learning of Network Representation</a>
4	5	<a href="#">SimWalk: Learning network latent representations with social relation similarity</a>
3	165	<a href="#">Don't Walk, Skip!: Online Learning of Multi-scale Network Embeddings</a>
2	114	<a href="#">Discriminative Deep Random Walk for Network Classification</a>
1	5	<a href="#">Learning distributed representations for large-scale dynamic social networks</a>

# Author Search Results

[Home Page](#) [More Results](#) [S2 Recommendations](#)

Author: David Madigan

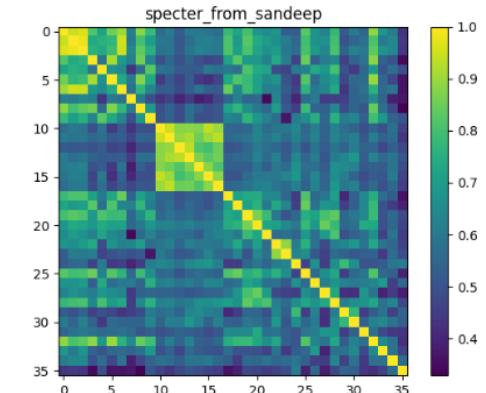
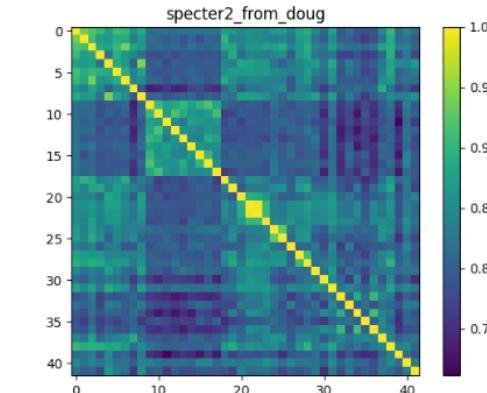
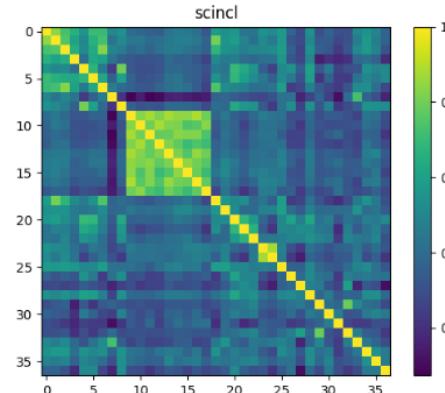
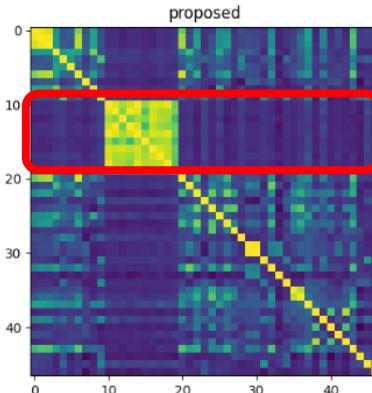
Limit: 10

Sort by: score

Embedding: proposed

## Table of Contents

1. [D. Madigan](#), hIndex: 62, citationCount: 19148
2. [D. Madigan](#), hIndex: 11, citationCount: 573 →
3. [D. Madigan](#), hIndex: 5, citationCount: 69
4. [David Madigan](#), hIndex: 4, citationCount: 229
5. [David Madigan](#), hIndex: 3, citationCount: 107
6. [D. Madigan](#), hIndex: 2, citationCount: 23
7. [David Madigan](#), hIndex: 2, citationCount: 204
8. [David Madigan](#), hIndex: 2, citationCount: 28
9. [David Madigan](#), hIndex: 1, citationCount: 8
10. [David Madigan](#), hIndex: 1, citationCount: 16



## The role of flavanoid polyphenols in beer stability

Determination of proanthocyanidins and catechins in beer and barley by high-performance liquid chromatography with dual-electrode electrochemical detection.

## CONTROL OF FERULIC ACID AND 4-VINYL GUAIACOL IN BREWING

Semipreparative chromatographic procedure for the isolation of dimeric and trimeric proanthocyanidins from barley

# Semantic Scholar (S2): Significant Effort

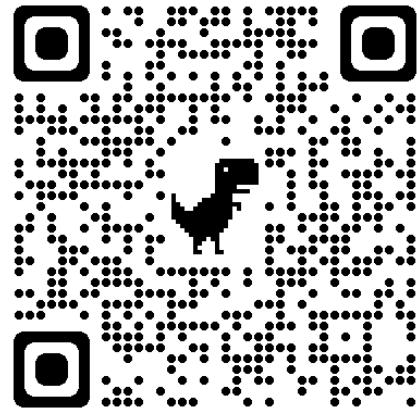
(slide from Dan Weld)



50 person team  
7 year project

207M+ scientific paper index  
8M+ monthly active users

[https://github.com/kwchurch/JSALT\\_Better\\_Together](https://github.com/kwchurch/JSALT_Better_Together)



## Better Together Team Github

### JSALT (Jelinek Summer Workshop on Speech and Language Technology):

#### Useful links

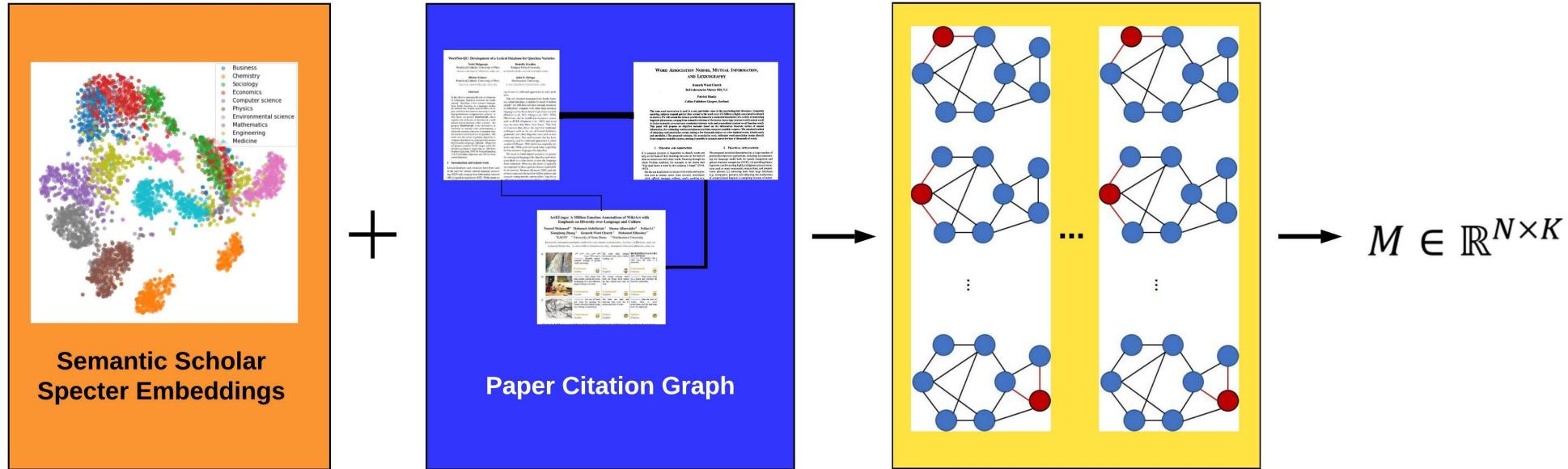
YouTube

Web Page

1. [Final Report \(YouTube\)](#) and [JSALT-2023 Team Page](#)
2. [Documentation of web service \(on Amazon AWS\), with examples](#)
3. [Slides](#)
4. [large datasets from Globus](#) and [What's Where \(on Globus\)](#)
5. [Deliverables](#)
6. [Reading List](#)
7. [Notation](#)
8. [SciRepEval Baselines](#)
9. [Zoom Link](#) and [Meeting Notes](#)
10. [Status of Production Runs](#)
11. [Similar Venues](#)

including  
these  
slides

# Better Together: Text (Titles, Abstracts, Body); Context (Citations)



TEXT

(a)

Specter, SciNCL

CONTEXT

(b)

Proposed (ProNE)

NEURAL NETWORKS  
(GNN / PRONE)

(c)

EMBEDDINGS

(d)

Standard benchmarks are small, static and clean

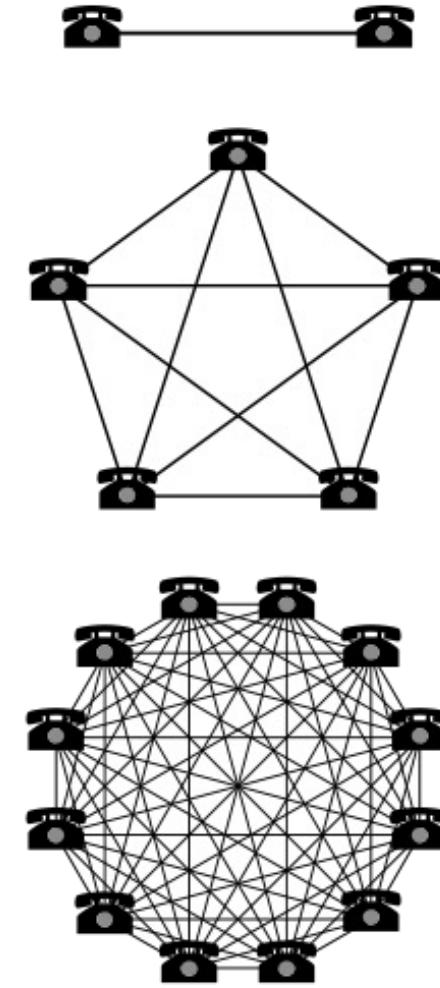
---

<b>Challenge</b>	<b>Standard</b>	<b>Proposed</b>
Volume	Small	Large (Metcalfe's Law)
Velocity	Static	Growing (double in 9 years)
Variety	Clean	Dirty (missing/bad values)
Forecasting	Random splits	Causal splits (bin papers by time)



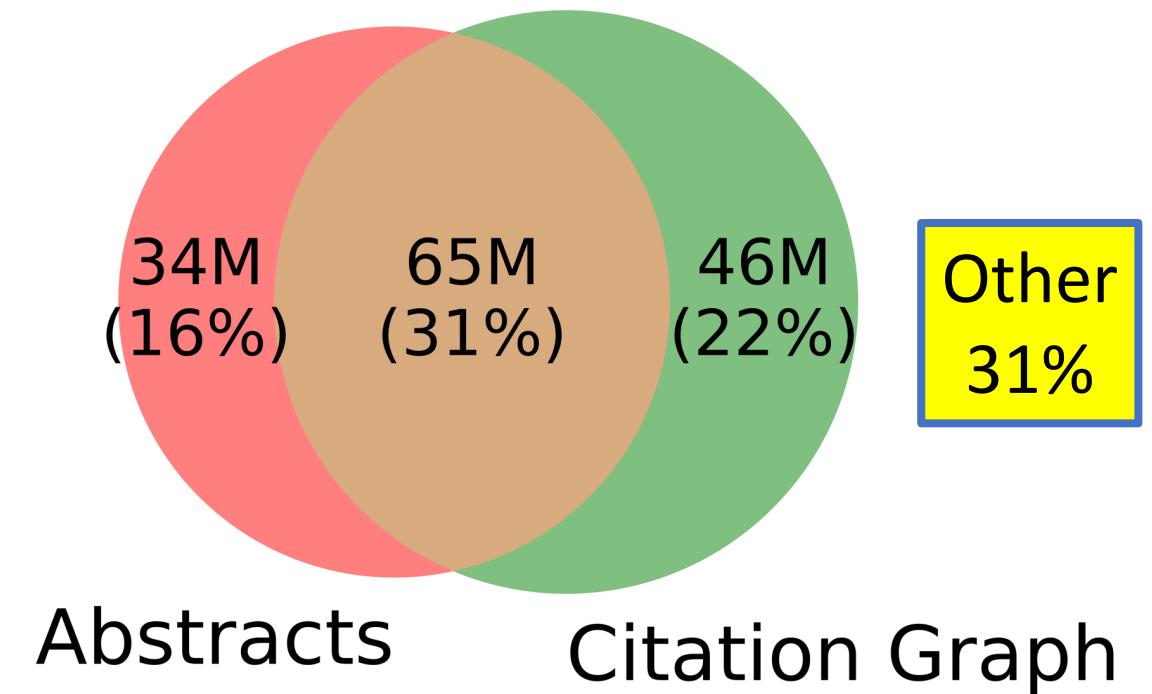
# Metcalf's Law (Network Effects)

- History: 3Com was selling small networks
  - 3 = 1 printer + 2 computers
  - Metcalfe argued they should sell bigger networks
    - (and more 3Com products)
    - because of economies of scale
- Economy of Scale:
  - Benefits scale faster than costs
    - Benefits:  $\sim n^2$
    - Costs:  $\sim n$
  - Law has been good for AT&T, Google, Social Media
  - Hypo: for Academic Search
    - Text/Context trade-off depends on scale



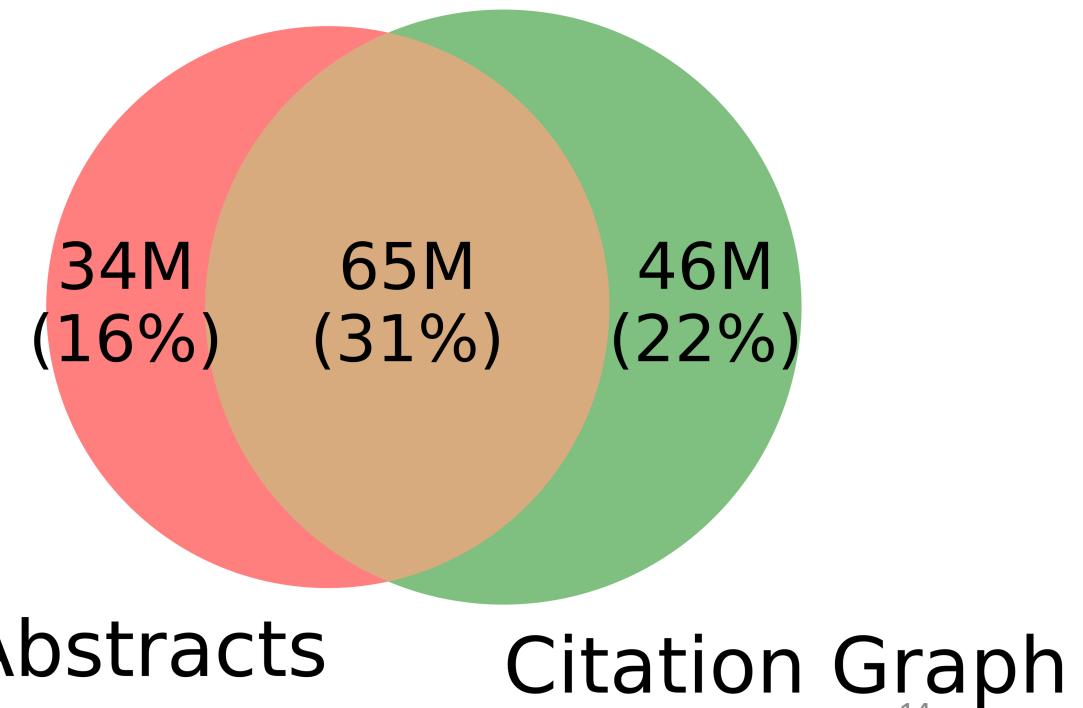
# Realities (Dirty): Missing Values

- $A$ : papers with abstracts,  $a$
- $L$ : papers with links,  $l$
- Opportunity:  $A \cup L$ 
  - Prior work targets small subset
    - Specter: designed for  $A$
    - GNNs: designed for  $A \cap L$



# Better-Together Conjectures

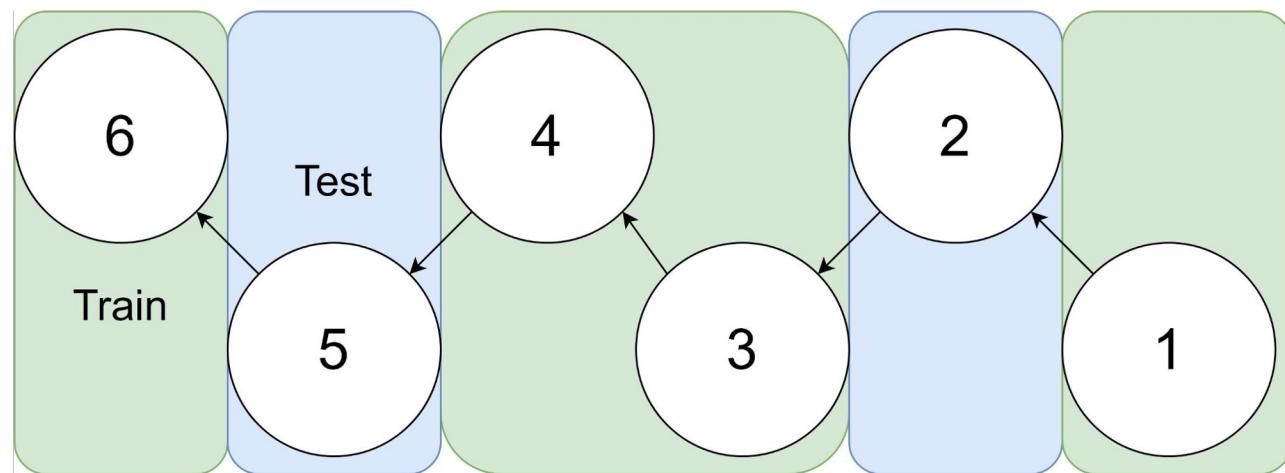
- Multiple representations are helpful
  - for missing values (and bad values)
    - If abstracts are missing, use links
    - If links are missing, use abstracts
  - and for graphs of different sizes
    - When graphs are small:
      - Text >> Links
    - when graphs are large:
      - Links >> Text (Metcalfe's Law)



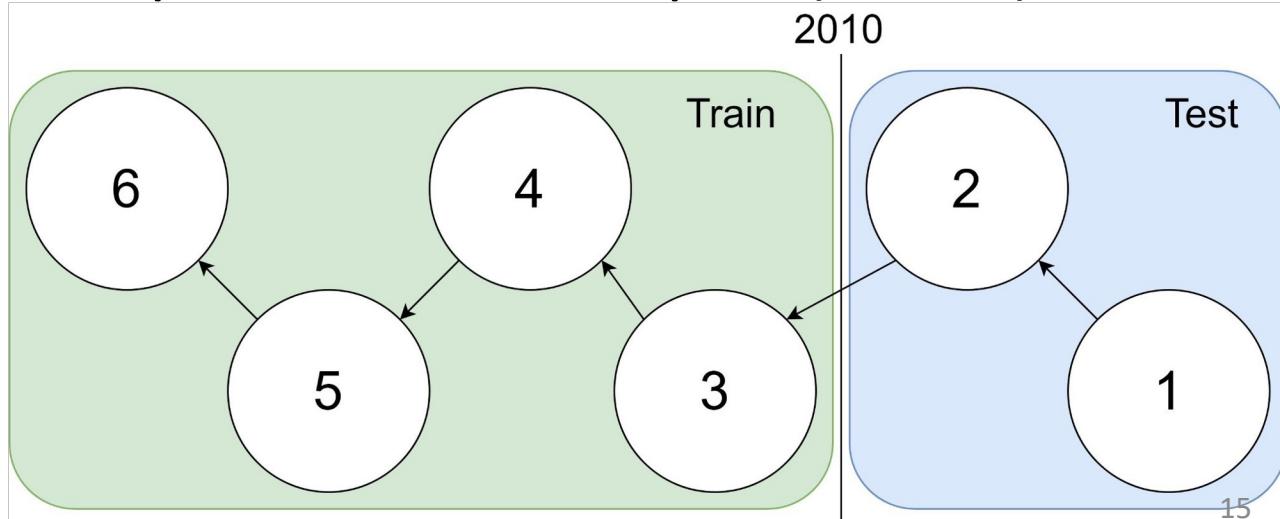
Example: 1 cites 2, 2 cites 3, ... 5 cites 6

idx	Pub Year	Paper Title
1	2018	[...] Photofragment imaging
2	2016	Convenient probe of S(1D2)[...]
3	2005	Megapixel ion imaging [...]
4	2003	Direct current slice imaging [...]
5	1995	profiles of Cl(2Pj) photofragments [...]
6	1988	Adiabatic dissociation of [...]

Traditional Train-Test Split (Non-Causal)



Proposed Train-Test Splits (Causal)



# Opportunities

## Scale

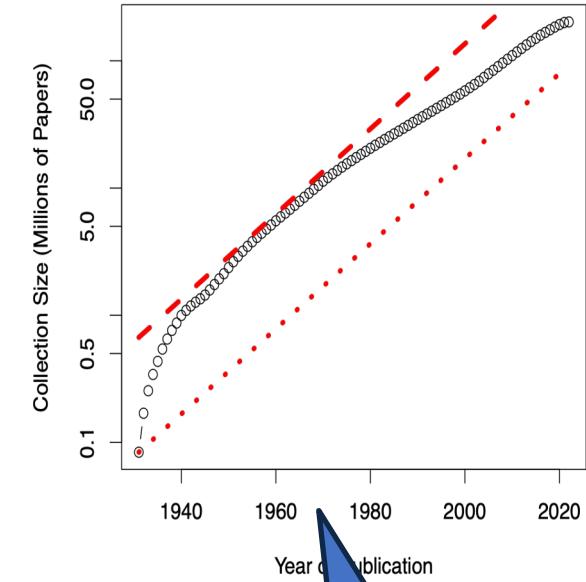
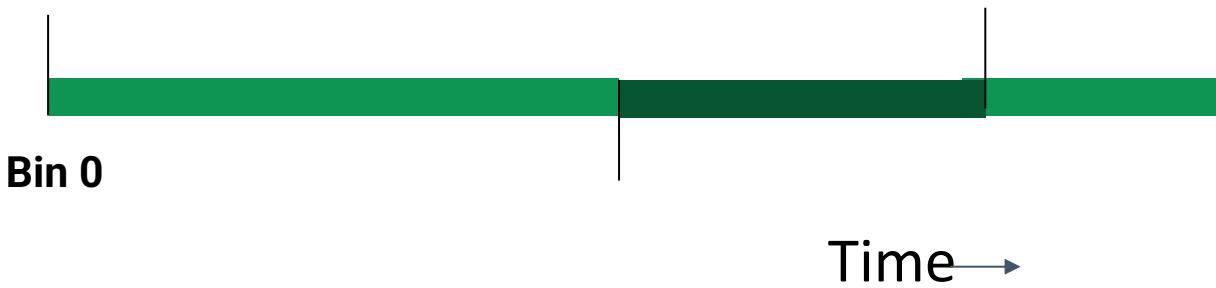
- Many benchmarks return a single number (figure of merit)
  - How does performance scale
  - with the size of the problem?

## Forecasting

- Standard test/train splits
  - More relevant to interpolation
  - than extrapolation
- Time is asymmetric
  - *It is difficult to make predictions especially about the future*
  - -- Yogi Berra

# Assign 200M papers to 100 bins (by time)

- Sort 200M papers by publication date
- Output: Split papers into 100 bins, with 2M papers/bin
- Because of exponential growth
  - (literature doubles every 9 years)
  - older bins span more time, and
  - newer bins span less time

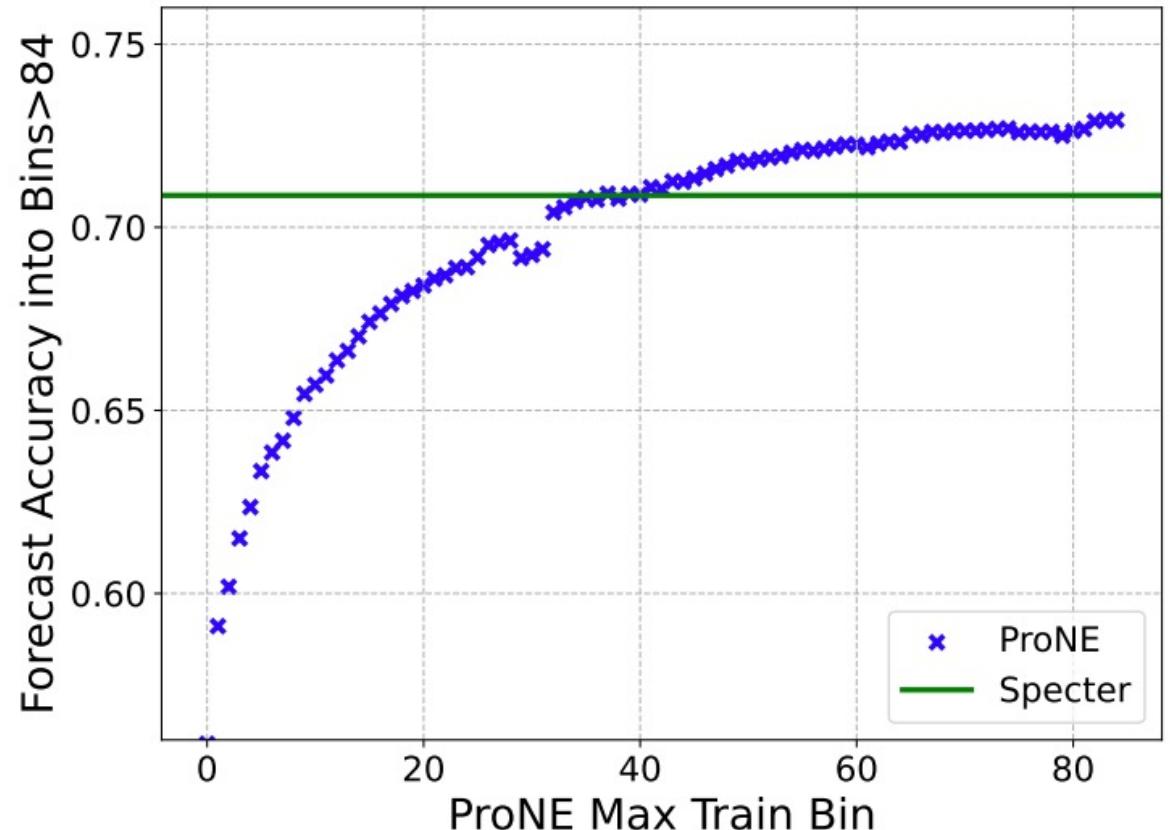


Bin 99

90% of literature published after I went to grad school

# Experimental Results

- Performance improves
  - when training set is larger and spans more time
- Performance degrades
  - with predictions further and further into the future
- Because of Metcalfe's Law
  - Text/Context trade-off
    - Depends on scale
  - Larger graphs favor context
    - ProNE: Spectral clustering on citations
    - Specter: BERT-like encoding of text
- Better-Together
  - Ensembles >> Either by itself



**Figure 9:** ProNE-Specter Crossover: Metcalfe's Law favors larger citation graphs (more than 82M papers).



# Conclusion: Deliverables

- GitHub:
  - massive (but needs editing)
  - slides, code, reading list
  - pointers to large files on Globus
- Web Server
  - Recommendations
    - Input: queries (docs)
    - Output: recommendations (docs)
      - What Should I read? Cite?
- Resources
  - Data behind web server:
    - citation graph
    - embeddings
    - indexes for approx. nearest neighbors
    - pairs of corpusIds with large cosines
    - all of the above for
      - Specter1, Specter2, SciNCL, ProNE
      - x100 bins for ProNE
- Evaluations & Benchmarks
  - Intrinsic: Perplexity, Random Walks
  - Extrinsic: Local Citation Prediction