# Text and Data Mining Panel Discussion

*ProQuest Dissertations & Theses*

**CGS 62nd Annual Meeting| December 8, 2022**

**ProQuest**

Part of **Clarivate**

# Global Scholarship by the Numbers

**5M+**
Works, including 2.9M in Full text

**4.1K**
Institutional contributors from around the world; dissertations and theses in **64 languages**

**240K+**
Works added in 2022

**400+**
Years of coverage, from 1637 to 2022

# Citation Connections

*Dynamic algorithms surface relevant, credible sources*

*Citation linking extends the path for source collection*

# A new partnership with Research4Life

in alignment with Clarivate's sustainability mission and United Nation SDGs



## Donating FREE ACCESS to institutions in 95 lower-income countries--part of



*"Our goal is to help close knowledge gaps, unify research, and elevate the scholarship being conducted in every corner of the world."*

- Angela D'Agostino, VP Dissertations & Theses, ProQuest part of Clarivate

https://www.research4life.org/other/95-research4life-countries-granted-pqdt-access/

# The Vision: Harnessing the Power of Data-driven Research Insights

# Dissertations and Theses are essential for comprehensive discovery



**Web of Science™**

- Peer-reviewed Journal Articles
- Conference Proceedings
- Books
- Emerging Sources

**ProQuest Dissertations & Theses Global**

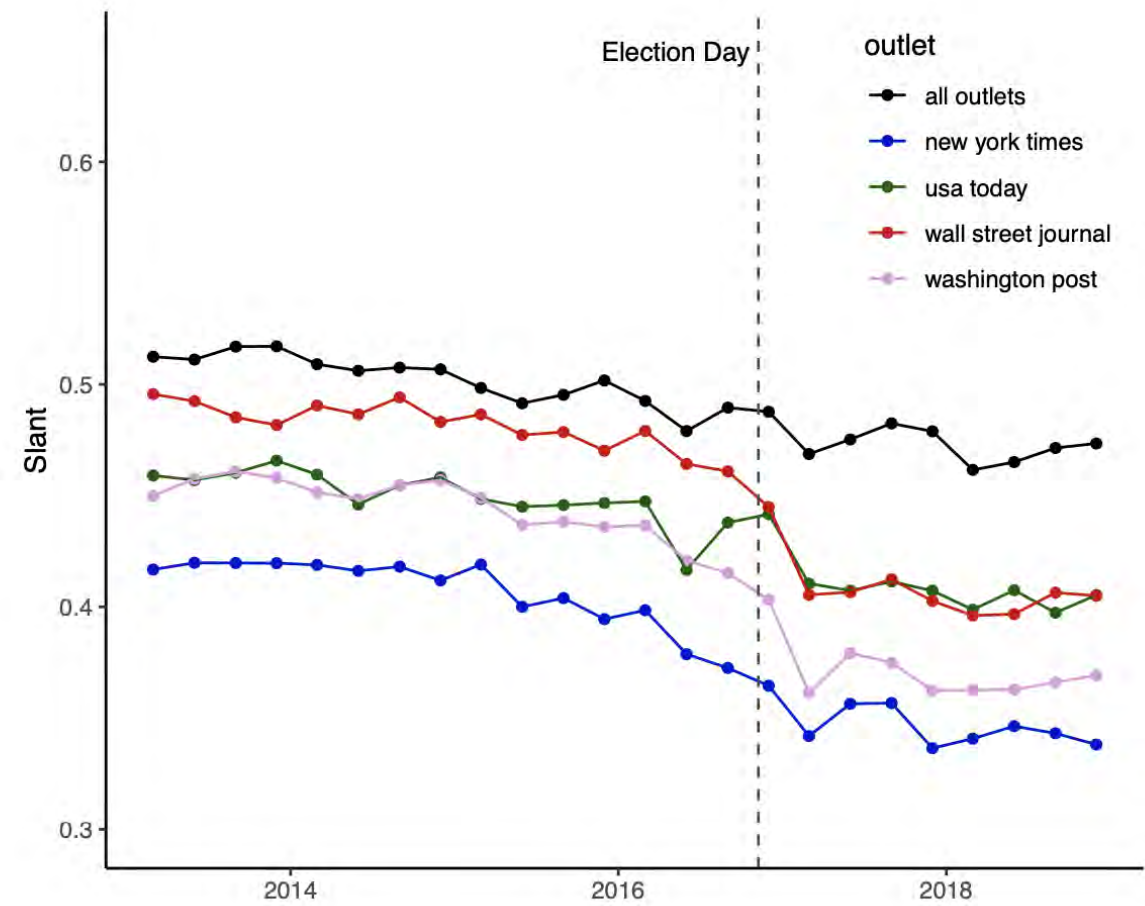- Unpublished scholarship
- Unpublished data sets
- Marginalized perspectives

# Journalist Ideology and the Production of News: Evidence from Movers (joint with Jacob Conway)

- Train RoBERTa-based model to measure slant using news articles tweeted by US politicians.
  - Significant improvement over prior bag-of-word methods

- Apply ML model to 20+ million full-text articles in TDM dataset.

- Document novel descriptive evidence on US media slant over time.

- Measure how slant changes as journalists move between outlets (journalist preferences explain 16% of variation).

LINK TO THE PAPER:
https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4295587

# Abstract(s) at the core: A case study of disciplinary identity in the field of linguistics

**Taylor LiCausi & Daniel McFarland**

- Leveraged the ProQuest Thesis ID variable to construct the broader field of linguistics.

- Used dissertation abstracts from 1980-2010 to create a structural topic model (STM).

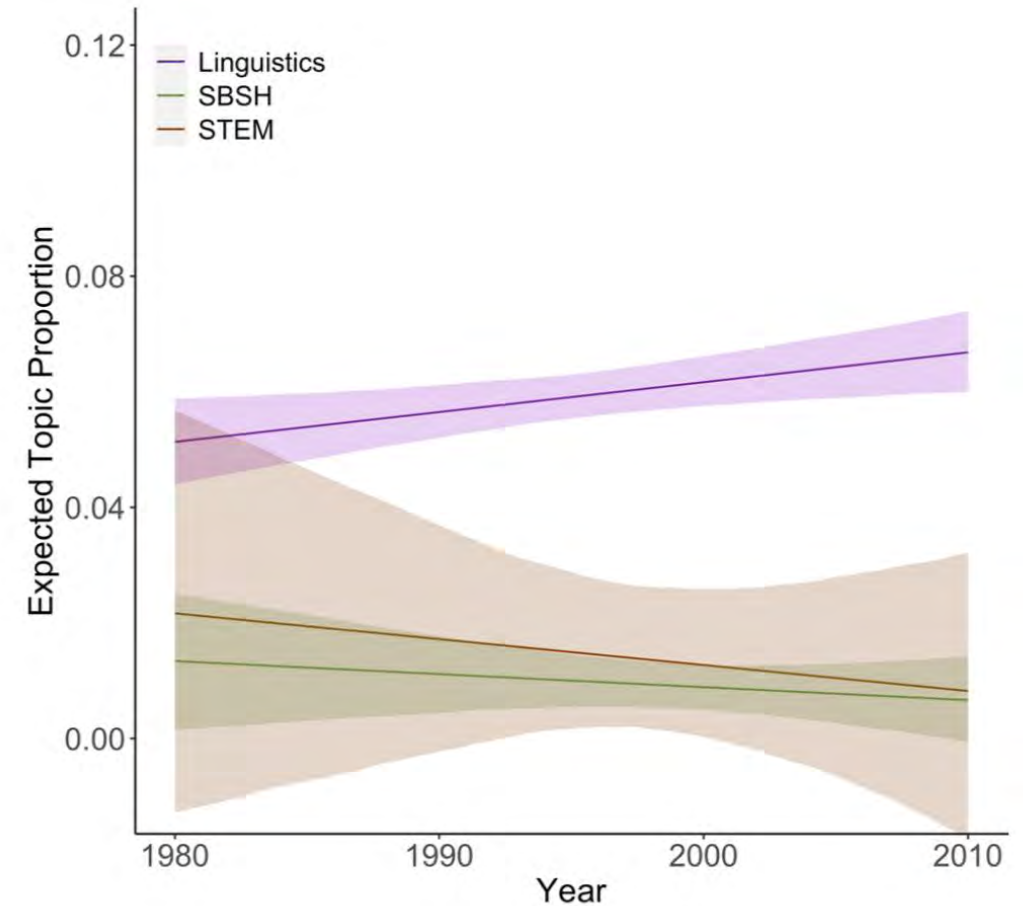- Mapped trends of our 40 topics by degree type and calculated topic correlations.



Fig. 1 Expected topic proportion of syntax: main (8), 1980-2010

# Results

- Sketched an intellectual topography of the broader field of linguistics

- Identified linguistics' disciplinary core through the topics most distinctive to linguistics degree abstracts.

- Found that STEM disciplines were more likely to engage with linguistics' core than SBSH disciplines.

- Found evidence of differentiation for SBSH and STEM topics.

Link to LiCausi & McFarland's article:  [Abstract(s) at the core: a case study of disciplinary identity in the field of linguistics | SpringerLink](#)

**Springer** Link

# Abstract(s) at the core: a case study of disciplinary identity in the field of linguistics

Taylor J. LiCausi ✉ & Daniel A. McFarland

**Fig. 13** The esoteric circle of linguistics

**SIGN UP FOR A FREE 30-Day TRIAL**
https://about.proquest.com/en/products-services/TDM-Studio/

ProQuest

Part of **Clarivate**

Analyze. Visualize. Connect. Discover.

ProQuest
TDM Studio
*allows you to...*

## Access and Analyze

Leverage content you already subscribe to and save months of time and tens of thousands of dollars required to negotiate and acquire data mining licensing. TDM Studio provides access to content with text and data mining rights already cleared.

## Target the Content

Streamline the content selection process. Refine and target the dataset to your research interest, eliminating the need to process unrelated documents.

## Connect and Collaborate

Enables real time collaboration across multiple sites and researchers in one workbench.

# ProQuest TDM Studio

**CLEARED LICENSES FOR TDM**
- Over 3K newspapers
- Dissertations & Theses
- Primary Sources
- Government documents (in BETA)

**WORKBENCHES** for R and Python

**VISUALIZATION TOOLS**
- Sentiment Analysis
- Topic Modeling
- Geographic Analysis

**MOST OFTEN USED FOR**
- Information Science
- Social Sciences
- Economics
- Political Science
- Business

# Thanks for joining us!

gilia.smith@clarivate.com

**ProQuest**

Part of **Clarivate**