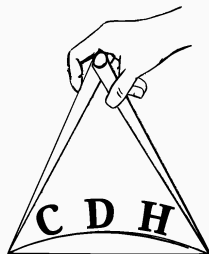


# What does Computational Humanities *look like*?

---

Tom Lippincott

Center for Digital Humanities  
Department of Computer Science



# This talk

## Quickly outline recent publications

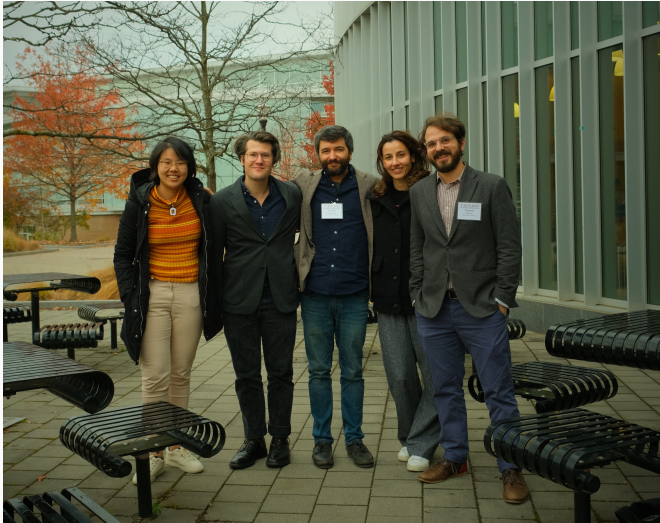
- "Nearest-neighbor" matching
- Some technical details in footnotes
- In-depth conversations . . .

## Basic research pattern

- Identify something of humanistic importance
- Translate into minimally-supervised ML method
- Interpret output as dense, aggregated document

Such research *fundamentally* depends on cross-training

# CDH researchers



Sabrina Li, Craig Messner, me, Hale Sirin, Sam Backer

# Shifts of semantic modality in Latin

## Foundational scholarship from comparative literature:

- Certain (sets of) words underwent a sequence of shifts
- Particular authors drove such change

## Train and inspect a temporal topic model<sup>1</sup>

- Measure *novelty* of authors w.r.t. preceding time-window<sup>2</sup>
- Measure *bi-modality shift* of each word<sup>3</sup>
- Aggregate these measurements in various ways

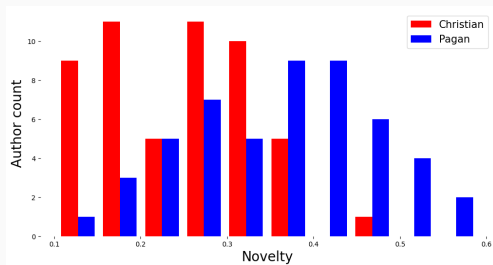
---

<sup>1</sup>Dynamic embedded topic model

<sup>2</sup>Jensen-Shannon divergence

<sup>3</sup>Function of first and second strongest topics and change-point detection

# Shifts of semantic modality in Latin



Word	Year	Delta
cathedra ( <i>chair</i> )	-175	0.947
cicatrix ( <i>scar</i> )	425	0.944
conlatio ( <i>bring together</i> )	350	0.939
auster ( <i>south wind</i> )	350	0.927
recte ( <i>upright</i> )	350	0.915

## Representations of dialect in fiction

- Differs from variation "in the wild"
- Aesthetic, culture, narrative . . .
- Non-linguistic regularities

## Apply models<sup>4</sup> to variety of orthographic realizations

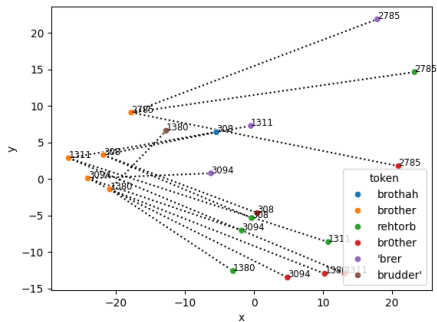
- Build corpus of annotated variants
- Generate perturbations
- Inspect relative positions in semantic space

---

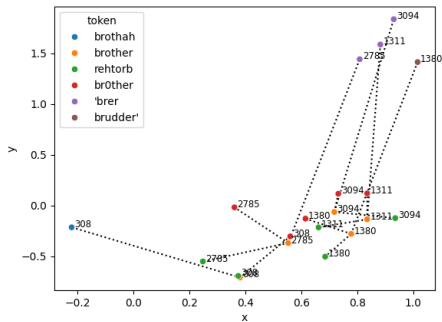
<sup>4</sup>BERT and CANINE

# Literary orthographic variation and LLMs

## "Brother" and variants



Embedded using BERT (subwords)



Embedded using CANINE (characters)

# Armeno-Turkish, large corpora, and language alternations

## Historical population with invisible record

- Turkish in Armenian script, not well-catalogued
- Known to be present in library collections
- Types of multi-lingual documents "oscillate" between languages

## Bootstrap language ID model<sup>5</sup>, apply to HathiTrust

- Assemble and train on a (noisy) corpus
- Find more Armeno-Turkish documents
- Characterize **periodicity** of language alternation<sup>6</sup>

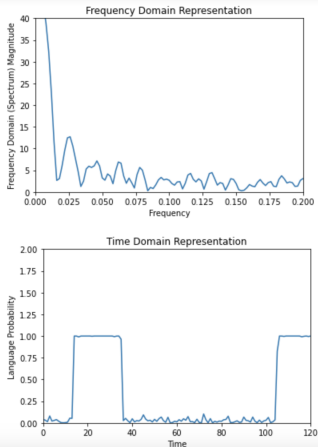
---

<sup>5</sup>FastText trained from scratch

<sup>6</sup>Fourier analysis to frequency domain



## Armeno-Turkish, large corpora, and language alternations



## Time series and frequency

ԹԻՐԱԲԱՏ	ԻՆԳԼԻՅԱՆ	ՔԱՆԻՆՖՈՐԴ
Աղէ՛ք անհշնչէ կէ- տիմ՝	I have come to tell you.	Այ չէ՛ք դժ Թա- թէ՛լ եւս։
Օ իւսաստանս շոգ տաւշեւեալէ՛մ,	I don't think much of it.	Այ ասո՛ղիք Թ՛չ ինք- բըջ ալ իմ է։
Տէտարի ինչ իմ արխի պիթ ցուտ։	No sooner said than done.	Նո սուրբըք սեւ տաւշեւ տալի։
Արթըզ ապայանամա՛մ	It can bear it no longer.	Այ քե՛նք պէ՛տք իմ նո լո՛ւնըք։
Ես փախալարիսիմ դա- ւանք։	As much as I can.	Է՛դ քըջ է՛դ սոյ քե՛ն,

ԳԼԵ 11

## Exemplars from three clusters



Craig Messner and Tom Lippincott. 2024.

**Pairing orthographically variant literary words to standard equivalents using neural edit distance models.**

In *LaTeCH@EACL*.



Hale Sirin, Sabrina Li, and Tom Lippincott. 2024.

**Detecting structured language alternations in historical documents by combining language identification with Fourier analysis.**

In *LaTeCH@EACL*.



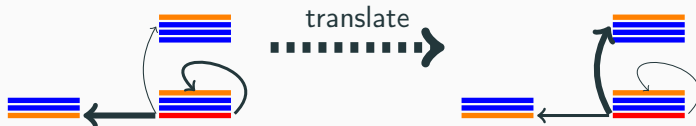
Hale Sirin and Tom Lippincott. 2024.

**Dynamic embedded topic models and change-point detection for exploring literary-historical hypotheses.**

In *LaTeCH@EACL*.

# Nascent research

## Preserving inter-textual geometry across translations<sup>7</sup>



## Cognitive salience of poetic structure<sup>8</sup>

*Ármă vîr|úmquē cǎ|nó || Trói|áe quí |prímŭs ăb |órís*

Poem → ?

? → Line+

Line → ?

? → (Stress | Unstress)+

<sup>7</sup>MBART shared embedding space

<sup>8</sup>Hierarchical non-parametric processes