# Computational Methods for Comparing Discourses

Zachary K. Stine[1], James E. Deitrick[2], and Nitin Agarwal[1]

[1] University of Arkansas at Little Rock, Department of Information Science
[2] University of Central Arkansas, Department of Philosophy and Religion
zkstine@ualr.edu, deitrick@uca.edu, nxagarwal@ualr.edu

## Poster PDF available here:
https://zacharykstine.github.io/files/discourse_comp_poster.pdf

## Summary

- Discourses associated with sociocultural identities comprise two levels of language that have important implications for comparing discourses:

    - Cultural lexicon—Deliberately used language that ties the discourse to an identity.

    - Cultural grammar—Language that reflects the broader perspectives encoded within the discourse.

- **Problem.** Comparing discourses requires researchers to look beyond superficial distinctions between identity-specific language (i.e., cultural lexicon). Instead, comparativists must interpret how such identity-specific language is used to express deeper structures underlying the discourses (i.e., cultural grammar). Owing to their subjectivity, traditional approaches have given rise to several important criticisms of comparativism.

- The goal of this research is to adapt computational tools for comparing discourses in such a way that:

    - Rigorous evidence can be found to support comparisons, lessening the influence of researchers' assumptions and mitigating several criticisms of comparativism.

    - Connects quantitative results back to the text thereby allowing traditional methods to be pursued in parallel in a complementary fashion.

- Motivated by this goal, we have explored two methodologies for comparing discourses:

    - **Method 1**—We compare how different discourses make sense of each other by training separate topic models on each discourse and measuring the organizational consistency between models in terms of how they allocate their own topics to the text of other discourses [3]. **Read more here**.

- ○ **Method 2**—Leveraging the resilience of a corpus' distributional semantic structure to word removal, we train topic models on all of the discourses combined, before and after removing highly distinguishing words from each [4]. **Read more here**.

- We use the discussion histories of various religious, spiritual, and other online communities from the platform, Reddit. These discourses are all primarily English-language.

## Extended Abstract

An under-explored aspect of the abundant linguistic data available for computational analysis is the potential to conduct large-scale comparative analyses between different discourses as a more rigorous complement to existing methods. A chronic obstacle for comparative studies is the requirement that a researcher distinguish between discursive distinctions that are superficial and those that reflect deeper, structural relationships between the discourses under analysis. In other words, two discourses may employ distinct lexical items, but understanding their broader relation to each other requires interpreting the meaning behind their distinctive forms, thereby translating the forms of each into a "neutral" conceptual space in which they can be meaningfully compared. Our research is focused on integrating computational tools into comparative studies of linguistic data to differentiate surface-level distinctions from deeper, structural-level distinctions in a rigorous way that ameliorates certain historical problems with comparative work and uncovers large-scale discursive relationships that would be otherwise hidden.

Our use of computational methods for comparative analyses relies primarily on representing discourses within the latent semantic spaces learned by topic models [1] and using the statistics of information theory to make sense of their relationships within such representations. Thus far, we have explored two related, but different approaches to comparing the discourses of several English-language, religion-related communities from Reddit. The communities, or subreddits, we analyze discuss general world religions (e.g., r/Hinduism), relatively newer religious movements (e.g., r/Pagan), topics that resist easy classification as religious or not (e.g., r/Spirituality), and that provide interesting counterpoints (e.g., r/Philosophy).

In the first approach [3], we measure the consistency with which discourse-specific models apply their lexically distinct topics to a common collection of text in order to compare the latent organizational schemes within the discourses. We measure this consistency between models and identify indirect mappings between model features. We find that, while such discourse-specific features may have stark differences between the lexical items that characterize them, their mappings between models imply structural similarities beneath their distinctive forms. For example, models trained on r/Christianity produce topics related to abortion discussions that map to topics on vegetarianism in models trained on r/Buddhism, suggesting

that, while discussions around abortion and vegetarianism may be discourse-specific, they share a deeper structural aspect in how they address these highly contentious ethical issues.

In the second approach [4], we compare the features of each discourse directly within three different representations consisting of individual word types, latent features learned from the unaltered corpus, and latent features learned from a modified corpus from which highly discourse-specific terms have been removed. These distinguishing terms are identified based on how much they contribute to the information divergence between discourses, following [2]. It has been shown that the latent semantic space learned by a topic model is resilient to the removal or subsampling of such words [5]. We find that some discourses (e.g., r/Islam) appear highly dissimilar to others, but that this is primarily due to the inclusion of discourse-specific terms (e.g., *hadith*). Once removed, these communities are revealed to have much closer resemblances to several others at the deeper, structural level (see Figure 1). Importantly, removing highly discourse-specific terms does not simply force all discourses to appear equally similar. The r/Philosophy subreddit, for example, undergoes an *absolute* decrease from the unaltered to the modified corpus, but it is actually a *relative* increase compared to the changes undergone by the other discourses. This increase in relative dissimilarity accords with the reasonable intuition that r/Philosophy may have some superficial similarities with the other subreddits, but that it is quite dissimilar in the deeper motivations and purposes underlying its discourse.
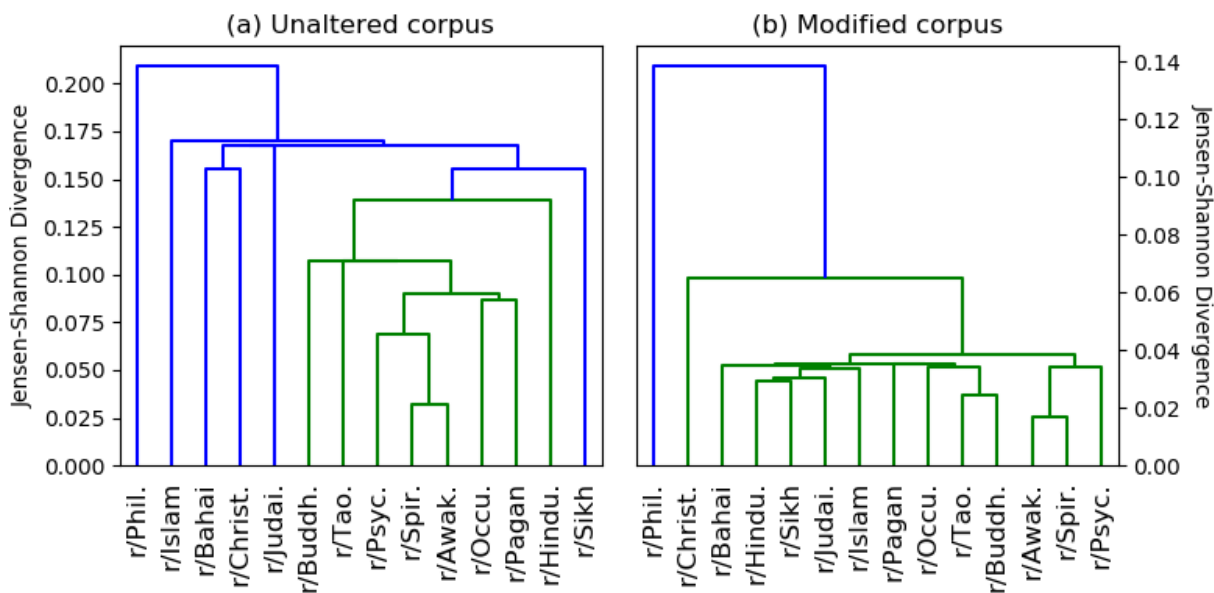


Figure 1: Hierarchical clusterings of discourses from 14 subreddits using (a) topics learned from the unaltered corpus and (b) topics learned from the modified corpus with highly distinguishing terms removed. Clusters are joined based on minimum divergences between their members.

While much remains to be explored in this work, its impacts are clear: The methodologies we describe enable researchers to meaningfully compare discourses beneath the surface of their superficial distinctions in a rigorous and interpretable way.

## References

[1] Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent Dirichlet allocation. *J. Mach. Learn. Res. 3*, Jan. (2003), 993–1022.

[2] Klingenstein, S., Hitchcock, T., and DeDeo, S. The civilizing process in London's Old Bailey. *PNAS 111*, 26 (2014), 9419–9424.

[3] Stine, Z. K., Deitrick, J. E., and Agarwal, N. Comparative religion, topic models, and conceptualization: Towards the characterization of structural relationships between online religious discourses. In *Proc. of the Workshop on Computational Humanities Research* (Nov. 2020), p.128-148.

[4] Stine, Z. K., Deitrick, J. E., and Agarwal, N. Using Information Divergence to Differentiate Deep from Superficial Resemblances among Discourses. Accepted for publication in *Proc. of the 8th Intl. Conf. on Culture and Computing* (Jul. 2021), part of HCII'21.

[5] Thompson, L., and Mimno, D. Authorless topic models: Biasing models away from known structure. In *Proc. of the 27th Intl. Conf. on Comp. Ling.* (Aug. 2018).