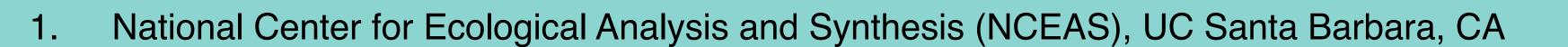# Do synthesis centers produce novel, potentially transformative research?
## Research publication diversity as an indicator of novelty and transformative capacity

Stacy Rebich-Hespanha[1], Ismael Rafols[2], John N. Parker[3], Edward J. Hackett[3], Joao P. Hespanha[4], The Sensible Science
Working Group, Assessing Synthesis and Synthesis Centers[1,5]

1. National Center for Ecological Analysis and Synthesis (NCEAS), UC Santa Barbara, CA
2. Ingenio (CSIC-UPV), Polytechnic University of València, Spain
3. Arizona State University, Tempe, AZ
4. Dept. of Electrical and Computer Engineering, UC Santa Barbara, CA
5. National Evolutionary Synthesis Center (NESCent), Duke University, UNC Chapel Hill, NC State University
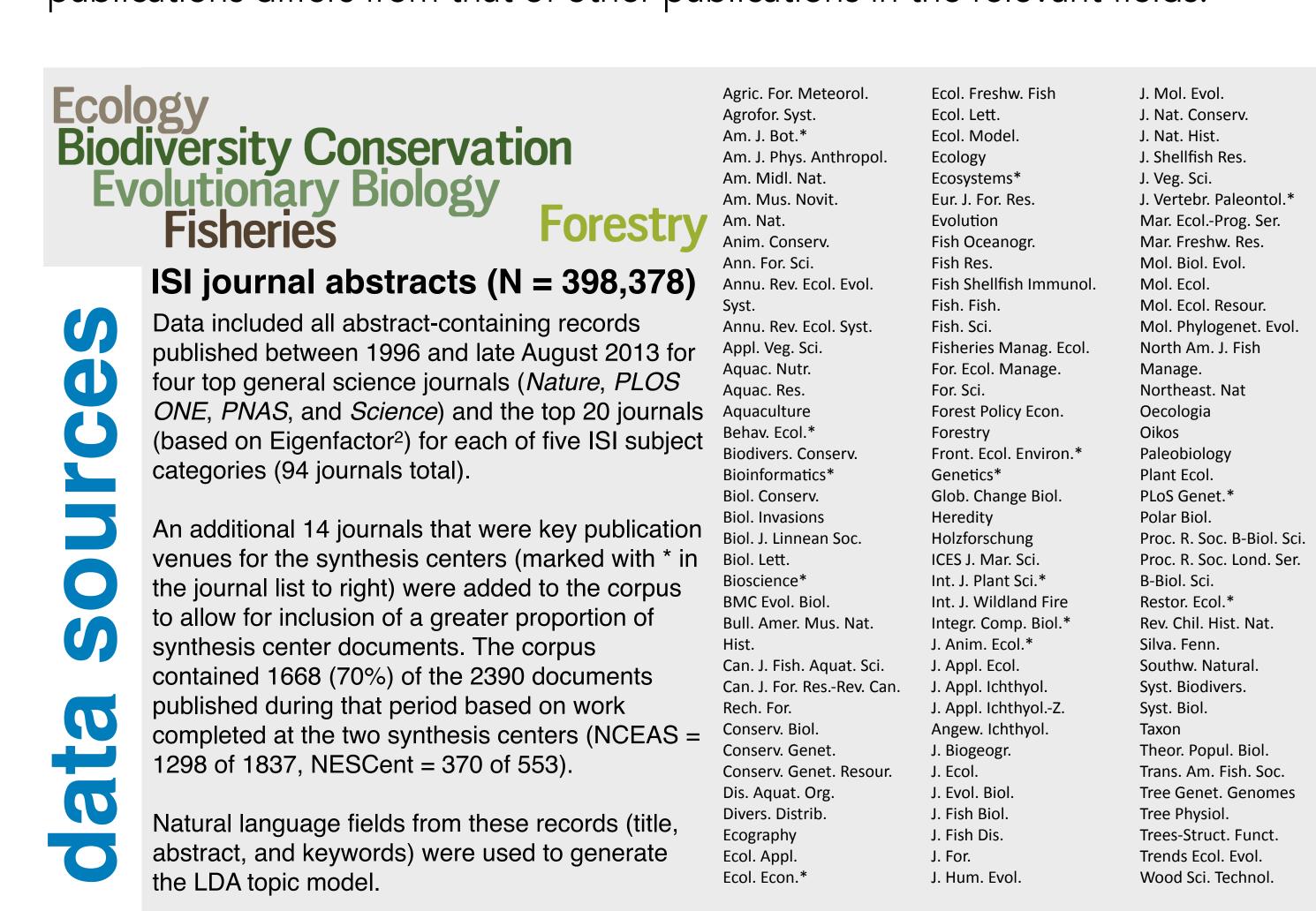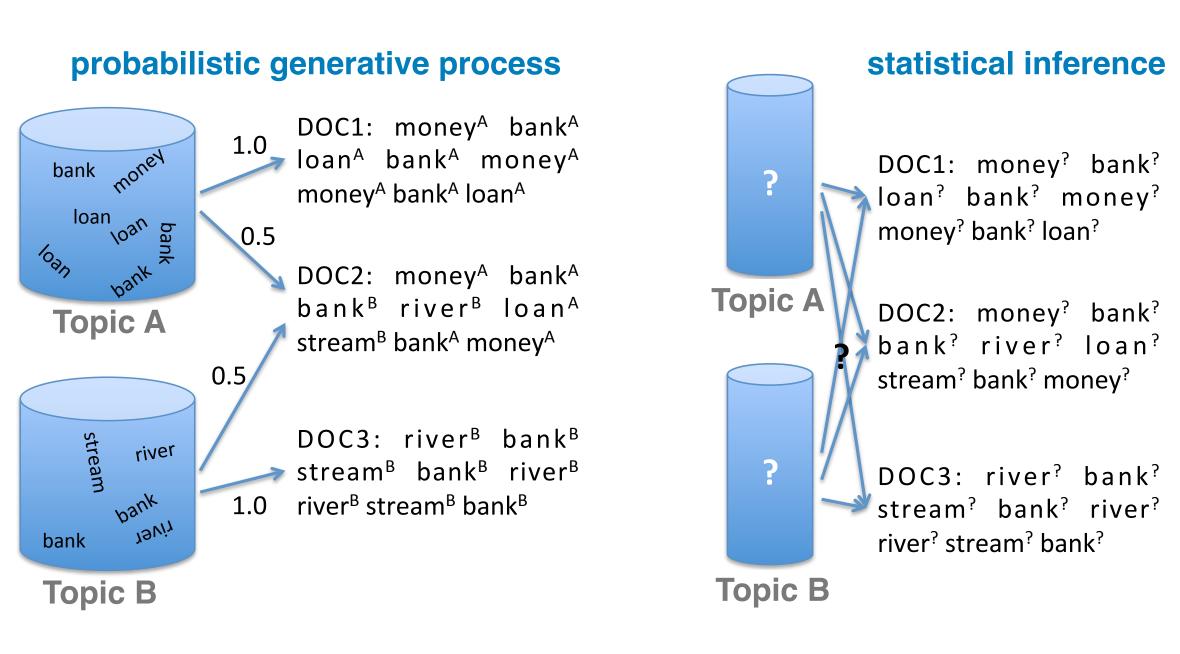
## evaluating synthetic publications

Synthesis is an emerging method for producing transformative research, and centers to promote synthesis are on the rise in the US and arou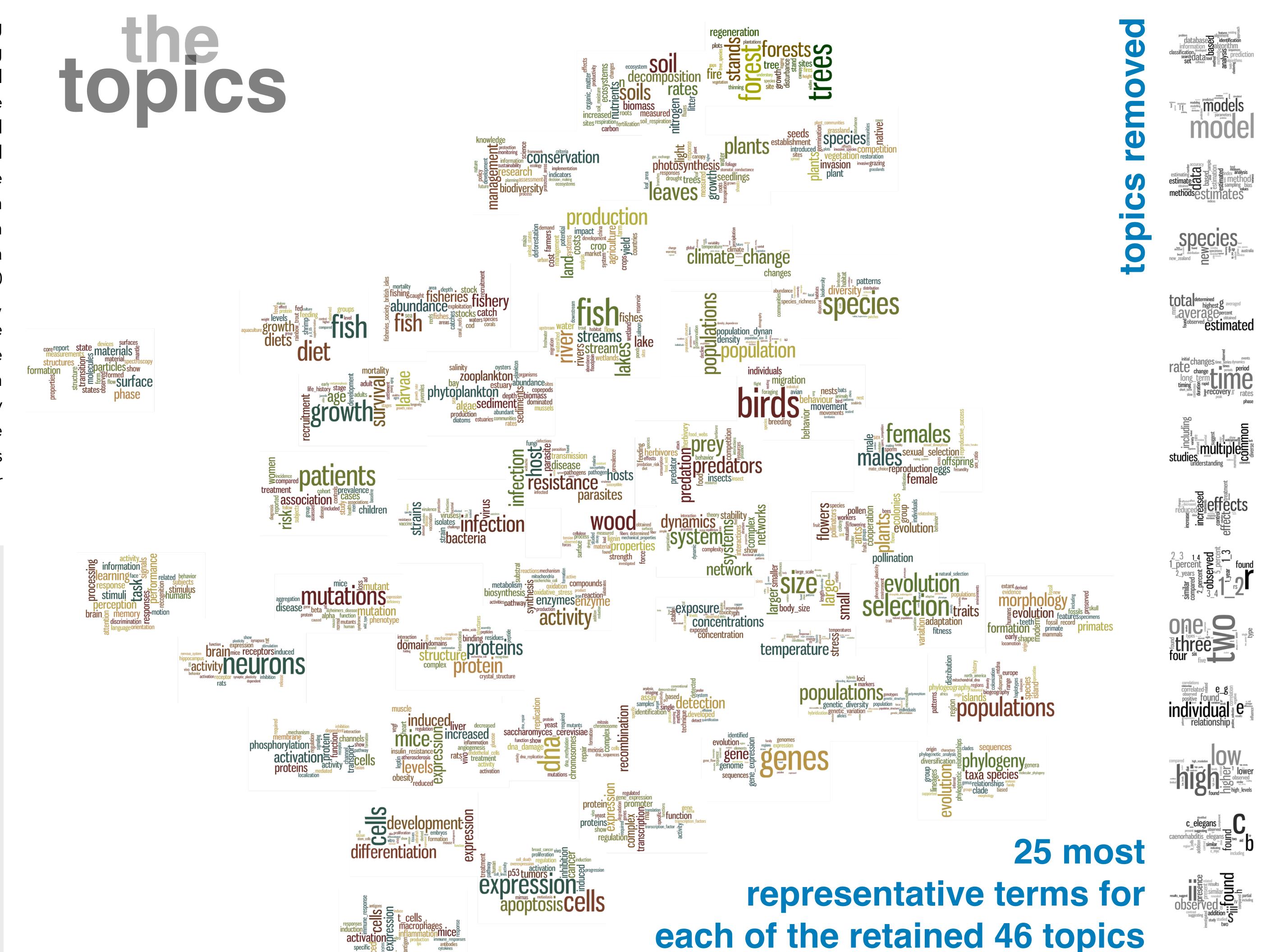nd the world. New analytic tools and techniques are needed to assess the originality and transformative potential of synthesis. We propose that research outputs produced within synthesis centers will exhibit distinctive qualities that distinguish them from other publications in their fields. To explore this possibility, we conducted a topical analysis of titles, abstracts, and keywords for approximately 400,000 articles published in 108 leading journals from the fields of Ecology, Evolutionary Biology, Biodiversity Conservation, Forestry, and Fisheries. We then described each document as a proportional combination of the discovered topics, and used the Rao-Stirling heuristics to estimate, for each document, various measures that illuminate contrasting aspects of diversity (i.e. variety, balance, and disparity). We then compare diversity metrics for the synthesis center documents with those for all other documents in our corpus to evaluate whether and how the measured diversity of synthesis center publications differs from that of other publications in the relevant fields.

## data sources

Ecology
Biodiversity Conservation
Evolutionary Biology
Fisheries
Forestry

ISI journal abstracts (N = 398,378)

Data included all abstract-containing records published between 1996 and late August 2013 for four top general science journals (Nature, PLOS ONE, PNAS, and Science) and the top 20 journals (based on Eigenfactor[a]) ranked on each of five ISI subject categories (94 journals total).

An additional 14 journals that were key publication venues for the synthesis centers (marked with * in the journal list to right) were added to the corpus for terms in topics and topics in documents, and discovers latent 'topics' by repeated sampling across the entire corpus. Once the latent topics (represented by groups of co-occurring terms) have been discovered, per-document topic 'mixtures' or combinations are inferred.

1. Steyvers, M. and T. Griffiths (2007) Probabilistic topic models. In T. Landauer, D. McNamara, S. Dennis and W. Kintsch (Eds.) Handbook of Latent Sematic Analysis, pp. 427-448. Earlbaum, Mahwah, NJ.
2. Kohonen, T., et al. (1996) SOM_PAK: the Self-Organizing Map program package. Helsinki University of Technology. http://www.cis.hut.fi/research/som_lvq_pak.shtml

## LDA topic modeling

Topic modeling, in this case Latent Dirichlet Allocation (LDA)[1], is an unsupervised probabilistic method for extracting a quantitative representation of semantic content for a document corpus based on observed patterns of term co-occurrence. The LDA model assumes Dirichlet priors for terms in topics and topics in documents, and discovers latent 'topics' by repeated sampling across the entire corpus. Once the latent topics (represented by groups of co-occurring terms) have been discovered, per-document topic 'mixtures' or combinations are inferred.

### probabilistic generative process

### statistical inference

## the topics

**25 most representative terms for each of the retained 46 topics**

The documents in the corpus were organized based on topic similarity using a self-organizing map algorithm[2] and then grouped into clusters based on their most dominant topics. Topic cluster positions were then manually adjusted slightly based on k-means clustering solutions. Topic clusters are organized so that similar topics are near each other. Within each topic cluster, the largest words are the most representative terms for that topic.
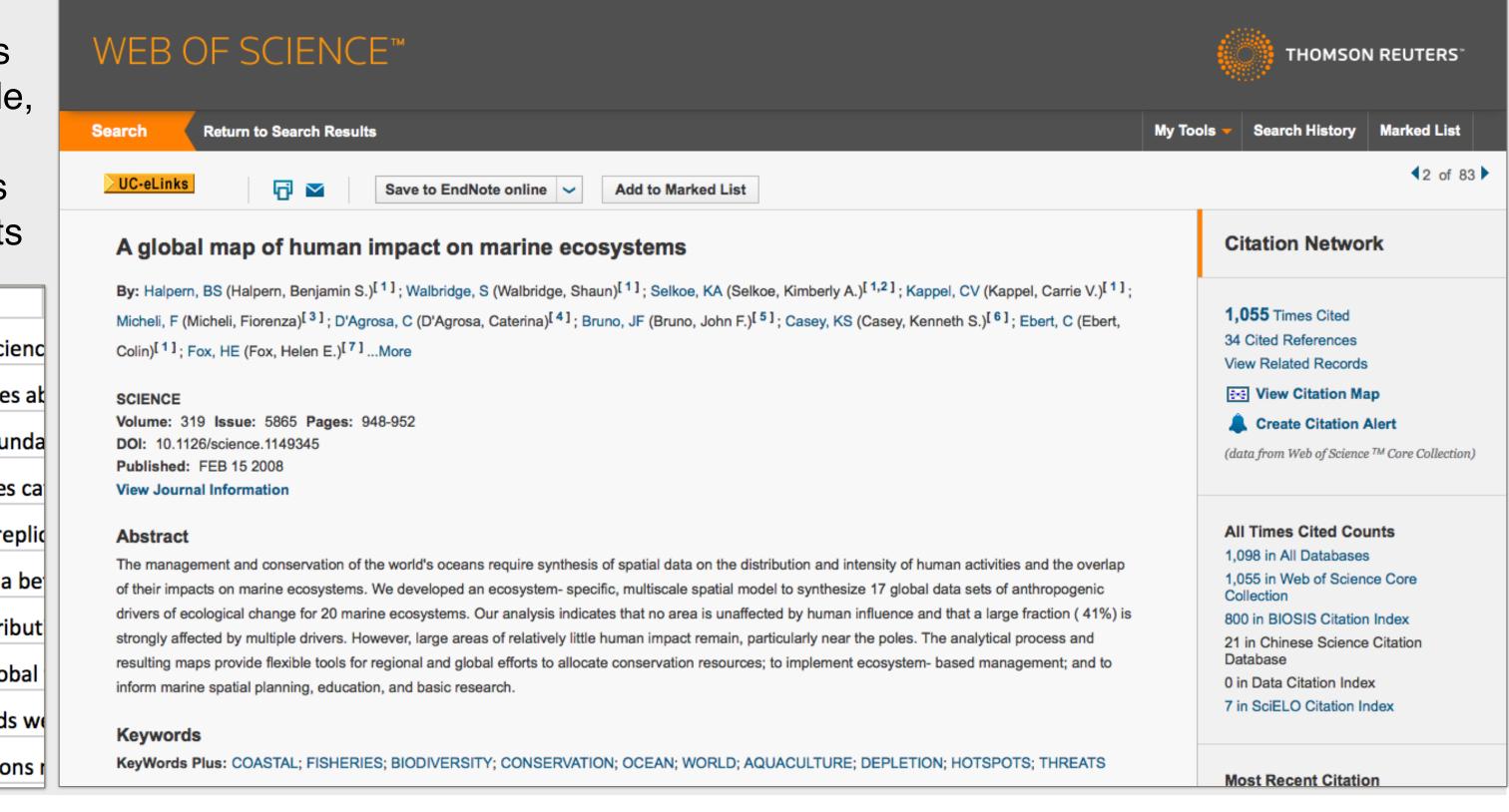
## documents as topic mixtures

This table shows the inferred LDA topic mixture for one of the synthesis publications. The Web of Science screenshot on the right shows the title, abstract and keywords that were used by the topic model to make this inference. The first column 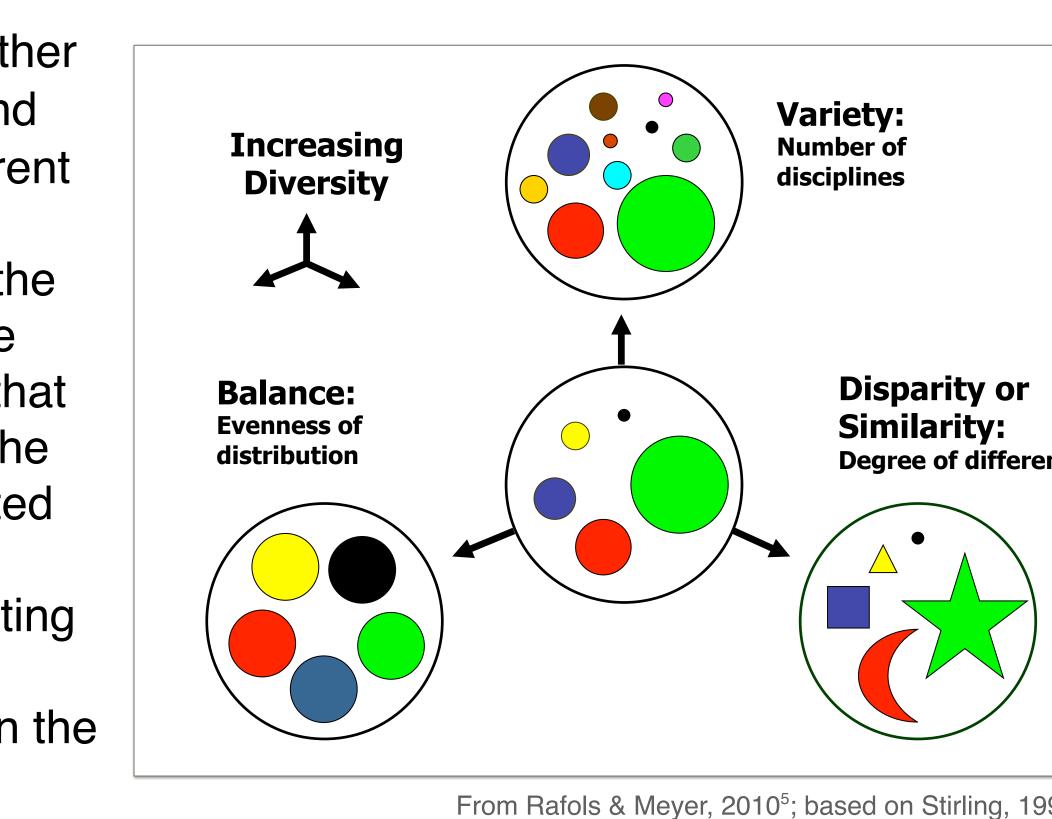of the table shows the topic ID for the topics associated with this document; the second shows the estimated weights for each of those topics in the document; and the third column shows the most representative terms for each of these topics.

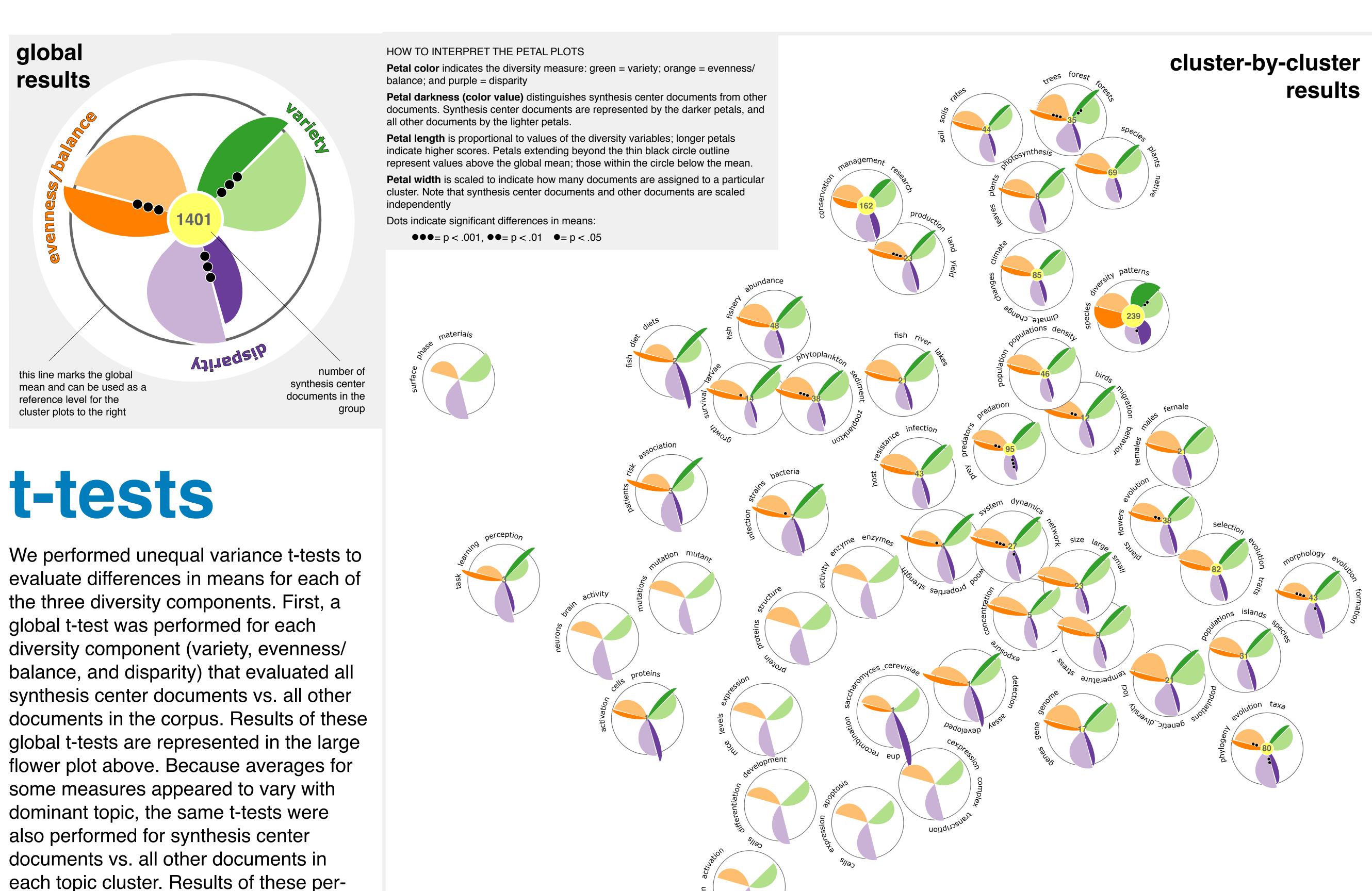| TopicID | Weight | Representative Terms |
|---|---|---|
| topic 02 | 0.3529 | conservation management research biodiversity knowledge science |
| topic 03 | 0.1569 | species diversity patterns habitat species_richness communities a... |
| topic 15 | 0.1373 | phytoplankton sediment zooplankton algae sediments bay abund... |
| topic 18 | 0.0392 | fish fishery abundance fisheries catch stocks fishing stock fishes ca... |
| topic 31 | 0.0392 | dna recombination saccharomyces_cerevisiae chromosomes replic... |
| topic 39 | 0.0392 | mutations mutation mutant disease mice phenotype loss alpha be... |
| topic 17 | 0.0196 | conservation biodiversity conservation science biology phylogeogra... |
| topic 13 | 0.0196 | climate_change changes climate temperature precipitation global... |
| topic 18 | 0.0196 | fish river lakes streams stream lake rivers fishes water wetlands w... |
| topic 27 | 0.0196 | selection evolution traits adaptation variation fitness populations ... |

## Rao-Stirling diversity: variety, balance, and disparity

The concept of diversity has been applied in various ways in ecology and other natural sciences, information sciences, and social sciences. Rao (1982)[3] and (Stirling (2007)[4] have advanced analytical frameworks that distinguish different aspects or components of diversity: *variety* (the number of categories associated with an entity), *balance* or *evenness* (how evenly represented the categories are), and *disparity* (how dissimilar the given categories are). We apply these three components of diversity to analysis of the topic mixtures that we obtained through LDA topic modeling. To calculate variety, we counted the number of topics associated with each document; for evenness, we evaluated the shares or proportions of each of those topics in a document; and for disparity, we first measured similarity/dissimilarity between topics by evaluating how rare all pairwise combinations of two topics were (cosine similarity/distance) and then used those distances to derive disparity scores based on the combinations of topics observed in a single document.

3. Rao, C. R. (1982). Diversity and dissimilarity coefficients: A unified approach. Theoretical Population Biology, 21, 24–43.
4. Stirling, A. (2007). A general framework for analysing diversity in science, technology and society. Journal of the Royal Society Interface, 4(15), 707–719.
5. Rafols, I. and M. Meyer (2010) Diversity and network coherence as indicators of interdisciplinarity: case studies in bionanoscience. Scientometrics 82:263–287.
6. Stirling, A. (1998). On the economics and analysis of diversity. SPRU Electronic Working Paper. http://www.sussex.ac.uk/Units/spru/publications/imprint/sewps/sewp28/sewp28.pdf Accessed Feb 10, 2015.

*From Rafols & Meyer, 2010[5]; based on Stirling, 1998[6]*

### global results

**HOW TO INTERPRET THE PETAL PLOTS**

Petal color indicates the diversity measure: green = variety; orange = evenness/balance; and purple = disparity

Petal darkness (color value) distinguishes synthesis center documents from other documents. Synthesis center documents are represented by the darker petals, and all other documents by the lighter petals.

Petal length is proportional to values of the diversity variables; longer petals indicate higher scores. Petals extending beyond the thin black circle outline represent values above the global mean, those within the circle below the mean.

Petal width is scaled to indicate how many documents are assigned to a particular cluster. Note that synthesis center documents and other documents are scaled independently

Dots indicate significant differences in means:
● p < .001, ●● p < .01, ●●● p < .05

**cluster-by-cluster results**

this line marks the global mean and can be used as a reference level for the cluster plots to the right

number of synthesis center documents in the group

## t-tests

We performed unequal variance t-tests to evaluate differences in means for each of the three diversity components. First, a global t-test was performed for each diversity component (variety, evenness/balance, and disparity) that evaluated all synthesis center documents vs. all other documents in the corpus. Results of these global t-tests are represented in the large flower plot above. Because averages for some measures appeared to vary with dominant topic, the same t-tests were also performed for synthesis center documents vs. all other documents in each topic cluster. Results of these per-topic tests are represented in the diagram to the right.

This figure shows the average variety, evenness/balance, and diversity scores for the documents in each topic cluster. Each topic cluster appears as a small flower; positions are the same as those used in the word cloud figure above. Here one can observe that synthesis center publications appear in over three-quarters of the topic clusters, but that they are concentrated in topics related to species diversity, conservation management, predator-prey relationships, climate change, evolution and phylogeny, and plant species. Synthesis center documents have significantly higher balance/evenness scores in many of the clusters, and significantly lower disparity scores in a smaller number. In one case (species, diversity, patterns topic), synthesis center publications have significantly higher disparity, but the disparity scores for that cluster overall are well below the global mean. Synthesis center publications also demonstrated higher variety scores in the species diversity and trees topic clusters, but further analysis suggests these differences are mainly attributable to differences in document length.

## regression analyses

A set of regression analyses show reasonable explanatory power for topic category across all three diversity components, but especially for evenness/balance and disparity. The number of tokens in the document and topic weight removed account for a majority of the accounted for variation in variety scores. However, the models that included all predictors accounted for only 17.6% of the variance in variety, 15.9% in ev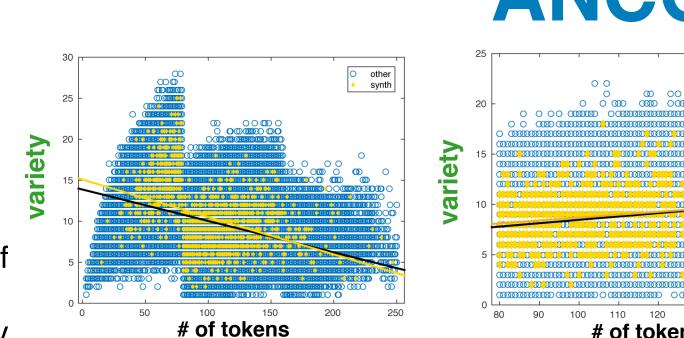enness/balance, and 9.3% in disparity. Other potential predictors accounted for very little variance, but the signs of their coefficients yield interesting patterns. For example, while number of authors often has non-significant relationships with variety and disparity, number of unique addresses has predictably positive relationships with these variables.

## ANCOVA

Because number of tokens and topic weight removed were important predictors for variety in particular, we performed analysis of covariance to examine differences between synthesis center documents and all other documents when these variables were accounted for. Here we show ANCOVA results for variety vs. number of tokens. No significant difference in slope or intercept was observed for variety when this variable was accounted for, but evenness/balance remained significantly higher for synthesis center documents, and disparity significantly lower. Because these differences remain when number of tokens and topic weight removed are accounted for, it is reasonable to interpret the cluster-by-cluster differences in evenness/balance and disparity shown in the flower plots at left as reflective of differences not accounted for by these variables.

**effect of synthesis center document vs. non-synthesis center document, controlling for number of tokens**

## insights, limitations & next steps

Whereas previous studies had investigated whether interdisciplinary centers were more or less diverse without specifying aspects of diversity, the novelty of our approach is that is capable of **discerning the type of thematic diversity** that characterizes "synthesis centers": an even combination of topics (relatively high balance) that are clearly distinct but that are not extremely disparate in cognitive terms (relatively low disparity) appear to characterize the work produced by the two centers studied. We have demonstrated the fruitfulness of applying theoretical measures of diversity to thematic quantitative models of research publications, and expect that these measures, when combined with other quantitative and qualitative approaches, will significantly advance our understanding of synthesis science. However, while we have identified significant differences in diversity between synthesis center publications and other publications and cast doubt upon the idea that these differences might be accounted for by larger numbers of authors or collaborating institutions, our analyses did not suggest mechanistic explanations for how the observed differences are achieved. Furthermore, the strength of our analyses is limited by the types and quantities of research publication data that are readily available for computational analysis and the lack of information about the scientific processes through which research publications are generated. As we continue to explore the patterns in this data, we will also seek to expand the number of synthesis centers included in our analyses and combine this approach for identifying publication diversity with other methods for discerning synthetic practices and qualities at various stages of the research process.