
Article

The Distant Reading of Religious Texts: A “Big Data” Approach to Mind-Body Concepts in Early China

Edward Slingerland,* Ryan Nichols, Kristoffer Neilbo, and Carson Logan

This article focuses on the debate about mind-body concepts in early China to demonstrate the usefulness of large-scale, automated textual analysis techniques for scholars of religion. As previous scholarship has argued, traditional, “close” textual reading, as well as more recent, human coder-based analyses, of early Chinese texts have called into question the “strong” holist position, or the claim that the early Chinese made no qualitative distinction between mind and body. In a series of follow-up studies, we show how three different machine-based techniques—word collocation, hierarchical clustering, and topic modeling analysis—provide

*Edward Slingerland, Department of Asian Studies, University of British Columbia, Vancouver, British Columbia, Canada. E-mail: edward.slingerland@gmail.com. Ryan Nichols, Department of Philosophy, California State University, Fullerton, CA. Kristoffer Nielbo, School of Culture and Society, Aarhus University, Aarhus, Denmark. Carson Logan, Department of Computer Science, University of British Columbia, Vancouver, British Columbia, Canada. Most of this work was funded by a Social Sciences and Humanities Research Council (SSHRC) of Canada Partnership Grant on “The Evolution of Religion and Morality” awarded to E.S., with some performed while K.N. was a visiting scholar at the Institute for Pure and Applied Mathematics (IPAM), which is supported by the National Science Foundation, and E.S. was an Andrew W. Mellon Foundation Fellow, Center for Advanced Study in the Behavioral Sciences, Stanford University. We would also like to thank the anonymous referees at JAAR for very helpful comments that have made this a much stronger paper.

convergent evidence that the authors of early Chinese texts viewed the mind-body relationship as unique or problematic. We conclude with reflections on the advantages of adding “distant reading” techniques to the methodological arsenal of scholars of religion, as a supplement and aid to traditional, close reading.

5

THE CLAIM that traditional Chinese thought was characterized by mind-body holism is commonly encountered in the field, with a venerable pedigree (e.g., Lévy-Bruhl 1922, Granet 1934, Rosemont and Ames 2009, Jullien 2007; for review, see Slingerland 2013). Scholars agree that, if there were a word for “mind” in classical Chinese, it would be *xin* 心, variously translated as “mind,” “heart,” or “heart-mind.” *Xin* refers literally to the organ of the heart, but is also the locus of cognitive and emotional function. Defenders of what we will be calling “strong mind-body holism” claim that although the *xin* may possess its own unique functions in the early Chinese view, this is no different from the eye or ear possessing their own specific functions. This position, therefore, holds that the *xin* is not uniquely contrasted with the body, and that it was viewed as simply one among a set of embodied organs (Geaney 2002).

Previous work (Goldin 2003, 2015; Slingerland 2013; Poli 2016) has reviewed qualitative textual and archaeological evidence that contradicts this claim, as well as cognitive science evidence that suggests that at least a “weak” form of mind-body dualism is a human cognitive universal. Unlike Cartesian, or “strong,” mind-body dualism, which postulates a razor-sharp divide between two ontological realms, cognitively natural dualism acknowledges that mind and body, although qualitatively distinct, overlap in various respects (Bloom 2004, Cohen et al. 2011). For instance, the mind is unique in being the seat of cognition, rational planning and thought, free will, and personal identity, partly because it is less material in nature than the other components of the self. The body—which, for the early Chinese, includes the organs other than the *xin* (“heart-mind”)—is something one can possess, or lose, but the mind/*xin* is central to one’s identity and sense of self.

In addition to this more traditional evidence against the strong mind-body holist position, several years ago the results of a methodologically novel, team-based coding project (Slingerland and Chudek 2011) bolstered the evidence typically presented in such arguments with more quantitative data. This study utilized human coders to analyze passages containing the keyword *xin* 心 in a corpus of pre-Qin (pre-221 BCE) texts, characterizing the functions of *xin* and how *xin*-body relations are

35

characterized. It was found that *xin* was frequently contrasted with the body, significantly more than any other organ in the body, a contrast that grew stronger over time. With regard to the function of *xin*, coders judged that, although *xin* seems to encompass both emotional and cognitive functions (and rarely refers to the actual physical organ in the body), by the Early Warring States (ca. fifth century BCE), cognitive functions outnumber emotional functions by 80% to 10%, a pattern that remained stable through the Late Warring States period.

A critique of this study by Klein and Klein (2011) included charges that the study was biased by drawing a large proportion of its pre-Warring States (before fifth century BCE) sample from a single text, the *Book of Odes*, a collection of poetry one would expect to contain an unusually high number of emotion words. Another concern voiced about the study is that a large proportion of the texts analyzed could be classed as philosophical works, which might exaggerate how much *xin* is portrayed as possessing cognitive functions. A final and more pervasive worry expressed by Klein and Klein and other critics concerns the degree to which human coders, making qualitative judgments of a given passage's meaning, might have their judgments skewed by prior assumptions or philosophical prejudices.

MACHINE-ASSISTED APPROACHES TO RELIGIOUS AND PHILOSOPHICAL TEXT ANALYSIS

Here we present a series of studies that respond to these and other critiques by applying machine-assisted, large-scale textual analysis techniques to a radically expanded textual corpus. The immediate, and more narrow, purpose is to explore conceptions of mind and body in early China. More broadly, however, we wish to demonstrate to scholars of religion the value of supplementing our traditional close reading practices with various techniques for “distant reading” (Moretti 2013) or computer-aided analysis of texts (Rockwell and Sinclair 2016). There are a host of new methodologies for navigating massive textual corpora that have been available to scholars of religion for some time now—in some cases, a decade or two—but that remain surprisingly underutilized.

It is a sign of how conservative academic disciplines are that the manner in which scholars marshal supporting textual evidence has not changed much in the last millennium or two, despite the availability of entirely unprecedented digital tools. The *Thesaurus Linguae Graecae* has been available to scholars of ancient Greece since the 1970s. For sinologists, the vast majority of the received corpus of traditional Chinese texts is available for free, online, in easily searchable form through a variety of sites. The Buddhist Canon, as represented in the full 85 volume *Taishō*

Shinshū Daizōkyō 大正新脩大藏經, is now available in searchable, online form, and similar resources are available for other religious and philosophical traditions around the world. Nevertheless, to date, these digital corpora have tended to be used as merely glorified concordances—as more convenient versions of tools we already had. The unprecedented and exciting analytic strategies provided by these resources have rarely been explored in religious studies, though other fields, especially literary studies, have taken steps toward distant reading (Moretti 2007, 2013), algorithmic criticism (Ramsay 2011), and text analysis through topic modeling (Jockers and Mimno 2013a).

The present study takes advantage of a massive textual dataset composed of 96 texts totalling 5.7 million characters, compiled by Dr. Donald Sturgeon in the “Chinese Text Project” (CTP), and freely available online.¹ The vast historical sweep of our corpus means we include texts from pre-Warring States (prior to fifth century BCE) through the Warring States and Han Dynasty (206 BCE–220 CE), as well as a small number of post-Han texts dating up to the Song Dynasty (960–1279 CE).² By increasing the number of texts we consider, we are able to meet concerns about corpus size, lack of genre diversification, and limited periodization. Our fully automated information retrieval and analysis allows us to not only handle such a massive corpus but also respond to concerns about potential biases in human coders.

The textual analysis techniques we demonstrate below are obviously no substitute for traditional close readings of texts. Indeed, most of the results are incomprehensible without the interpretative skills of experts deeply familiar with the corpus in question. As we will see, however, the sort of high altitude, broad perspective on a corpus provided by these techniques can serve as an important check on our qualitative intuitions. Moreover, there may be certain types of questions—for instance, assessing the validity of claims about general trends or prevalent themes in a given corpus—that are best addressed through machine-assisted techniques coupled with statistical analysis. As we will try to demonstrate below, the great strength of distant reading is the ability to pick up trends or patterns

¹We are grateful to Dr. Sturgeon, Postdoctoral Fellow in Chinese Digital Humanities and Social Sciences at the Fairbank Center for Chinese Studies, Harvard University, for permission to download the corpus in its entirety for purposes of analysis. The CTP also has some built-in analysis tools that can be extremely useful for scholars and can be subscribed to for full download access.

²See our online materials (<http://www.hecc.ubc.ca/articles/jaar-supplementary-materials/>), Appendix 1, for a complete list of texts, with their Era and Genre tags, as well as a figure representing the distribution of genres. Because of controversies concerning precise dating of early Chinese texts, as well as concerns about some of the genre labels employed in CTP, none of this meta-data was employed in studies reported here except for the separating out of medical texts from all other genres in Study 2.

in large quantities of data that may be invisible to individual human “close” readers.

METHODS, ASSUMPTIONS, AND GOALS

To use a textual corpus, a certain amount of preprocessing is necessary. We applied a stop-word list to the CTP, removing overly common function words and articles.³ We also removed all punctuation apart from sentence-ending punctuation, the inclusion of which allowed us to use the sentence, a natural unit of semantic meaning, as our primary unit of analysis. Once the corpus was preprocessed in this manner, we explored it with a variety of analytical tools, as described below.

Strong mind-body holism attributes a broadly monist metaphysics to the authors of the classical Chinese texts in our corpus, according to which the mind and body share one and the same type of substance, and the *xin* (the most likely candidate to represent the concept of “mind”) is in no way qualitatively distinct from the other organs in the body (Geaney 2002, Ames 1993, Jullien 2007). We would expect that, if the authors of the texts in the CTP corpus tacitly endorsed strong mind-body holism, we would find *xin* behaving just like any other bodily organ term in proximity to the three most common words for “body” in classical Chinese (*shen* 身, *xing* 形, and *ti* 體).

It is important, in this regard, to distinguish implicit from explicit cognition. Work in various branches of the cognitive sciences has documented the “dual system” nature of human cognition (Kahneman 2011, Evans 2008). According to this model, the explicit, “cold,” conscious aspect of the human mind rides upon a much larger, more powerful, and pervasive implicit, “hot,” mostly unconscious system. People are capable of entertaining and debating any number of propositions in their explicit systems. How many angels can dance on the head of pin? If body and mind are two distinct ontological substances, how could they ever interact? As a growing literature on “theological correctness” in the cognitive science of religion has documented, however, explicit claims or endorsed beliefs do not necessarily reflect underlying beliefs and behavior. Hindus may assert in surveys that their gods are omniscient and omnipotent, while narrative judgment tasks reveal that they implicitly believe them to be subject to anthropomorphic limits (Barrett 1998). Calvinists may profess to belief in predestination, but we still observe them praying on the weekends for God to favor

³Our stop-word list is presented in Appendix 2 online, along with a discussion of some of the limitations of the list we employed and a general discussion of challenges involved in composing stop lists.

their football team, and otherwise behaving in ways that suggest they believe their actions can change the future (Slone 2004).

When it comes to mind-body dualism in early China, it is certainly the case that one can find explicit endorsements of the strong mind-body holist position. The Confucian thinker Mencius famously remarked, in *Mencius* 6:A:7:

With regard to the mouth, all palates find the same things tasty; with regard to the ears, all find the same things pleasant to listen to; with regard to the eyes, all find the same things beautiful. Now, when it comes to the *xin*, is it somehow unique in lacking such common preferences? What is it, then, that minds share a preference for? I say that it is order and rightness. (Van Norden 2008, 151)

This passage is frequently cited by defenders of strong mind-body holism as evidence that the early Chinese saw the *xin* as equivalent to the other organs. As Jane Geaney comments, the point of the analogy set up in *Mencius* 6:A:7, as well as similar statements in other texts such as the *Xunzi*, is that “the heartmind and the senses have certain things in common—they function on the same principles regarding space and time, and they share the tendency to prefer similar things. The fact that the senses serves as analogies for the heartmind makes it unlikely that they are radically different in nature” (Geaney 2002, 101).

As Edward Slingerland has observed (2013, 19), however, to draw this conclusion from passages such as 6:A:7 is to mistake an explicit claim for a background assumption. Mencius is making an argument, which he no doubt expects to be surprising or counterintuitive, that the *xin* has a natural “taste” in the same way that the other organs do. This would be a nonsensical or, at best, trivial statement to make if he, and his readers, implicitly and intuitively *believed* this to be true. If we are going to claim that mind-body dualism is completely foreign to the early Chinese worldview, we need to look at background assumptions as well as explicit rhetorical claims.

The fundamental hypothesis we put to the test in the studies reported below is that, whatever early Chinese thinkers might explicitly *say* about *xin*-body relations, an analysis of the overall patterns of language use will reveal that the *xin* was at least implicitly understood as being unlike any other organ, with a unique relationship to the physical body. Our basic assumption is that large-scale patterns of language use will tell us something about implicit cognition. Even if Mencius claims that the *xin* is the same as the other organs, if we find that he, and other early Chinese writers, habitually mention *xin*, and only *xin* (as opposed to the other organ terms),

in relation to the body, this suggests that *xin* occupies a distinct cognitive space in the early Chinese mindset.

To evaluate this hypothesis, we undertake two collocation analysis studies, Studies 1 and 2, where we analyze the patterns of co-occurrence between various key terms in our corpus. In Study 1, we analyze semantically uncontroversial pairs of terms to see if we can establish some semantic benchmarks—that is, determine whether or not particular collocation patterns correspond to particular semantic relationships such as contrastive pairs (“many”::“few”) or part-whole (“wheel”::“cart”). In Study 2 we then present collocation data describing *xin*’s relationship with the three body terms, as well as similar data describing other bodily organs’ relationships with the body terms and some simple statistics that compare the two sets of data. The goal is not only to note differences between *xin* and the other organs, but also to see if any of the observed collocation patterns match the semantic benchmarks established in Study 1.

In Studies 3 and 4, we turn to “unsupervised” methods of textual analysis, which involve machine-learning techniques for processing and analyzing patterns in textual corpora (Jurafsky and Martin 2009; Miner et al. 2012). Specifically, our paper uses two methods, hierarchical clustering and topic modeling. Hierarchical clustering is a form of unsupervised information retrieval used to extract structured information from unstructured data (Manning, Raghavan, and Schütze 2008). Topic models are generative probabilistic models that seek out sets of words (“topics”) that reliably travel together throughout either a document or a corpus of documents (Blei, Ng, and Jordan 2003).

A great advantage of unsupervised methods is that they make no assumptions about target terms of interest or size of word window for collocations, let alone about experimenters’ potential hypotheses about the texts in question. They therefore provide a relatively objective measure of relationships between lexical items in the corpus. Although they are methodologically quite distinct from collocation analysis, with both of these methods (hierarchical clustering and topic modeling analysis) we make the same prediction: if the early Chinese authors were strong mind-body holists, *xin* should behave like the other organs of the body in their writings. *Xin* and other organ terms should all appear with equal frequencies in significant topics, and should not differ with regard to how they cluster vis-à-vis body terms.

INTRODUCTION TO STUDIES 1 AND 2: COLLOCATION ANALYSIS

Collocation analyses involve measurement of how frequently, and how closely, terms of interest occur with regard to one another in a textual

corpus and inference from those measures to, typically, syntactic features of words and parts of speech. In the field of corpus linguistics, this technique has long been used to track patterns and changes in idiom usage, bigrams, etc. Most humanities scholars, however, would be more interested in the semantic implications of word collocation. There have been some advances on this front (Gries 2013, Jurafsky and Martin 2015, Rohde, Gonnerman, and Plaut 2006), as well as practical demonstrations of how “dumb” collocation pattern extractors can “learn” something about natural language semantics. Collocation patterns have been used, for instance, to train a machine-learning algorithm to perform reasonably well on the multiple-choice synonym questions found in the Tests of English as a Foreign Language exam (Landauer and Dumais 1997). Indeed, as Bullinaria and Levy (2007) observe, semantic inferences generated from experienced collocation patterns in everyday word use probably play a central role in how infants acquire language (510). A literature deriving inferences from textual collocation patterns in a variety of genres (diaries, prose, sermons, emails, survey responses, and more) to patterns of thought and emotion in their authors has also been growing (see Teubert and Čermáková 2007, Sampson and McCarthy 2005, Manning and Schütze 1999).

With regard to classical Chinese, a study by Lee and Wong (2012) analyzed collocation patterns in the *Complete Tang Poems* to show affinities between, for instance, particular seasons and distinctive semantic classes of words. Word collocation studies of interest to religious studies scholars have, to date, focused on contemporary rather than historical materials. For instance, studies of portrayals of Islam in contemporary sources have demonstrated that the word “Muslim” in the British press is rarely of a specifically religious nature, more frequently serving as a reference to either national or ethnic identity (Baker, Gabrielatos, and McEneary 2013); that “Muslim women” and cognate terms were primarily found in semantic clusters concerning war, violence, and victimhood (Al-Hejin 2015); and that certain keywords and conceptual dichotomies could be used to automatically and reliably distinguish extremist Islamic documents from more neutral content (Prentice, Rayson, and Taylor 2012).

Methods for Statistical Analysis

There is a considerable amount of variation in the statistical measures used to measure collocation rates. This might be confusing to humanities scholars unfamiliar with statistics and also provoke suspicions that authors of various studies are in the habit of picking the measure most likely to give them the answers they wanted. For this reason, before proceeding

to a discussion of our studies, it is worthwhile to first discuss the strengths and weaknesses of the various statistical measures employed in large-scale textual analysis.

First, it is important to see that different collocational statistics are suitable for different research questions. Of the statistical tests developed by pioneers in the field of corpus linguistics, a handful are particularly relevant for our purposes. We are interested in the collocation of *xin*, other organ terms, and physical body terms within sentences. The most widespread test statistic in corpus linguistics is mutual information (MI), an extended measure of collocational strength drawn from information theory. MI is calculated by dividing the observed frequency of a co-occurring word within a specific window by the expected frequency of a co-occurring word in that window and taking the logarithm to base 2 of the result (Biber and Jones 2009, 1287).⁴ In other words, it measures the association strength between two words by comparing their probability of appearing together, while also considering their individual distributions (Church and Hanks 1990). The result of this calculation offers a score that indicates the strength of the relationship between two terms, or between two sets of terms. Since MI compares the observed co-occurrence of two words to what would be expected at chance level—that is, if the words were independent—interpretation of MI is straightforward. A positive score indicates that the words are more strongly associated than chance would predict, a negative score indicates that they are anti-associated, and a score of zero indicates that the association is at chance level (Oakes 1998; Paperno et al. 2014).

MI is good for tracking associations between words for which the joint probability distribution is similar to the words' individual probabilities distributions—that is, where the words in question are more or less equally common in the corpus. However, if one of the associated words occurs very infrequently in the overall corpus and the other frequently, the association is likely to be a product of chance. That is, MI scores for collocations involving rare words might be artificially high (Mautner 2007, Oakes 1998). When it comes to comparing *xin* to other organ terms, this is a particular worry. *Xin* and the body terms are quite common in the corpus, whereas several of the organ terms we compare with *xin* appear only rarely, especially outside of the medical genre.

We employ several approaches to attempt to remedy this problem (Paperno et al. 2014, Dunning 1993). First, Oakes recommends MI3, a measure that corrects for rare terms by cubing the normal MI measure—thereby

⁴Please see Appendix 3 for all of the equations used in this study and a concrete illustration of the methods.

making a small term that much smaller (Oakes 1998, 171–2). We therefore report MI3 in our [supplementary materials](#). In addition, we also report two measures that, in our view, do a better job of dealing with the potential distortions caused by rare words: the t-score and conditional probability.

Instead of measuring association strength, t-score estimates the confidence in two words being associated (Church et al. 1991). T-scores are calculated by subtracting the expected frequency of a term (given its overall frequency in the corpus) in a given window relative to a target term from its actual frequency. Then the result is divided by the amount of variance in the frequency of the term in question relative to the average frequency of terms in the corpus. T-score solves the shortcoming of MI by adjusting the joint probability distribution according to the weight of individual terms within the corpus, in effect lowering the t-score for rare terms as compared to more common terms. In terms of interpretation, t-score is similar to MI in that a positive score indicates an above-chance-level word association, a negative score indicates a negative association, and zero indicates independence between the words.

We also report conditional probability, which calculates the probability of the appearance of one word given the occurrence of another word. Conditional probability is calculated by multiplying the probability of the occurrence of word 1 given the occurrence of word 2 by the probability of the occurrence of word 2. Unlike other test statistics used in corpus linguistics, conditional probability is presented as a pair of values. This is because it is sensitive to asymmetries in word frequencies between, for example, the probability that the occurrence of the word *of* predicts the occurrence of the word *course* (a low probability, because of the ubiquity of “of”) versus the probability that the occurrence of the word *course* predicts the occurrence of the word *of* (a high probability, because of the comparative infrequency of “course”).

This leaves unanswered specific questions about what measures of collocation can tell us about semantic relationships. For instance, what is the most useful collocational window size for capturing significant semantic relationships? Are particular types of semantic relationships characterized by distinctive collocation scores? Although there has been important preliminary work exploring these and related questions (see especially Rhode et al. 2005, Bullinaria and Levy 2007, and a recent review in Jurafsky and Martin 2015), the answers are likely to vary from language to language, and possibly from genre to genre. For this reason, we begin with an effort to discover latent collocational patterns in early Chinese texts by letting the corpus talk to us rather than by attempting to test any hypotheses. In this semantic benchmarking exercise, our Study 1, we sought to gain knowledge about applications of corpus linguistics techniques to our

historical Chinese corpus by analyzing classical Chinese word pairs with known and fairly uncontroversial semantic relationships to one another.

STUDY 1: USING COLLOCATION MEASURES TO SEMANTICALLY BENCHMARK CLASSICAL CHINESE WORD PAIRS

5

For Study 1, we chose a series of semantic relationships most likely to prove useful in our subsequent analysis of mind-body relations. These included scalar opposites (e.g., *da* 大 “big” :: *xiao* 小 “small”); complementary social relations (e.g., *jun* 君 “lord” :: *chen* 臣 “minister”), cosmic forces (e.g., *ri* 日 “sun” :: *yue* 月 “moon”), and physical distinctions (e.g., *nei* 內 “inside” :: *wai* 外 “outside”); pairs linked by endemic function (e.g., *niao* 鳥 “bird” :: *fei* 飛 “fly”); and pairs characterized by a part-whole relationship (e.g., *che* 車 “cart” :: *lun* 輪 “wheel”) (see Appendix 4 online for a full list). To these semantic pairs we added two control conditions, where the second character in the pair was replaced first by another word vaguely semantically related to the target character, but not with the same sort of specific semantic relationship as the originally paired word, and then with a word semantically unrelated to the target character. In both control conditions, we picked substitution characters that matched as closely as possible the word frequency of the original character. Exploring these known semantic pairs, raw collocation counts were recorded for the sentence level and at the 10L10R (10 left, 10 right), 5L5R, 2L2R, and 1L1R windows. T-scores, conditional probability, MI, and MI3 were calculated for each pair at these various windows.

Our full results are available in Appendix 5 online. When it came to our contrastive, functionally related, and part-whole pairs, we found that, rather than clustering nearby the target character, we saw a relatively even distribution of collocations from the sentence and 10L10R level down to the 1L1R level. This suggests that there is no obvious “sweet spot” for capturing semantically significant collocations. That said, the differences between the target pair and the semantically related and -unrelated pairs become somewhat starker at smaller windows, which may mean that windows such as 2L2R or 1L1R do a better job of capturing semantic relations, or may simply reflect the fact that our target pairs often appear together as set pairs in the early Chinese corpus. The sentence, uniquely among our KWIC windows, reflects authorial or editorial decisions about semantic relatedness. It therefore has a kind of organic validity to it, and we accordingly have chosen to focus on this measure in our discussion here, concluding that we have the most confidence in drawing inferences about the psychology of authors when looking at co-occurrences of two terms within the same sentence, as opposed to within the same 100-word or

40

Table 1. Selected results from benchmarking study

Word Pair	Relationship Type	T-Score Sentence
big :: small	contrast	66.21
big :: high (<i>gao</i> 高)	semantically related	33.18
big :: eight (<i>ba</i> 八)	semantically unrelated	30.40
lord :: minister	contrast	82.75
lord :: state (<i>guo</i> 國)	semantically related	68.85
lord :: what (<i>he</i> 何)	semantically unrelated	64.05
sun :: moon	contrast	58.05
sun :: bright (<i>ming</i> 明)	semantically related	36.93
sun :: mutually (<i>xiang</i> 相)	semantically unrelated	26.77
inside :: outside	contrast	54.62
inside :: city wall (<i>cheng</i> 城)	semantically related	53.08
inside :: return (<i>gui</i> 歸)	semantically unrelated	15.03
bird :: fly	endemic function	15.76
bird :: black crow (<i>wu</i> 烏)	semantically related	5.68
bird :: simple (<i>jian</i> 簡)	semantically unrelated	2.50
cart :: wheel	whole - part	11.33
cart :: trail, track, rut (<i>gui</i> 軌)	semantically related	6.11
cart :: pollution (<i>wu</i> 汙)	semantically unrelated	3.06

10-word string. When it comes to statistical measures, with the exception of the MI score (which was not surprising, for reasons discussed above), t-score, conditional probability, and MI3 measures all provided broadly similar results. We will therefore focus on the t-score at the sentence level, although all of the measures, at all of the various KWIC windows we employed, are available in Appendix 5 online. 5

The basic results are striking, with a representative sample presented in Table 1.

To begin with, our contrastive pairs all show higher t-scores than the control pairs. In most cases, we see a clear pattern where the contrastive pairs have very high t-scores, with the scores falling off considerably for the semantically related pairs and then further still for the semantically unrelated pairs. This validates our prediction that semantically related terms will have strong collocation patterns, with contrastive semantic relationships being the most powerful of all. We have included in the selected results one example,⁵ moving from inside::outside to inside::city wall, 10
15

⁵The only other example in our data is the transition from summer::winter to summer::spring (see Appendix 5), which is almost certainly due to the fact that the four seasons frequently appear in lists together, artificially increasing the collocation measure for what was intended as a merely semantically related word.

where we do not get a sharp drop-off as we move from a contrastive pair to a semantically related pair (although we do still get a drop). It is possible that “city wall,” which our classical Chinese expert chose as a term likely to be related to “inside” (and conceptions of containment generally), was a poorly chosen semantically related control. In any case, this pair remains an outlier to our observed pattern. 5

Interestingly, the part-whole and endemic function⁶ pairs also showed the same pattern as the contrastive pairs: a consistent fall-off moving from the target pair to the merely semantically related and then the unrelated controls. This suggests that the strong signal we saw with the contrastive pairs may not be linked to that narrow semantic relationship, but might rather serve as a general signal of a specific, and well-defined, semantic relationship between two terms. In other words, we can imagine vaguely semantically related terms hovering around each other in the logical space of the text, kicking off the sort of middling t-scores that we see in the semantically related control pairs. Pairs of terms with specific, and well-defined, semantic relationships exert a stronger attraction on one another and therefore display dramatically higher t-scores. It is possible, then, that collocation measures, such as t-scores, can tell us about only the intensity or specificity of the semantic link between two terms rather than the exact nature of that link. We will return to this topic again below when we look at collocation patterns for *xin* and the other organs. 10 15 20

One problem that we recognize in this pilot study is that we are cherry-picking our target and control pairs. An ideal assessment of the link between collocation and semantics would obtain collocation measures for *all* term pairs in the entire corpus, rank them in strength, and then turn to experts familiar with the corpus to see if the trends are robust. The problem with this approach is that it is computationally intractable, given the vast number of relationships that would have to be tested (a total of $2^{15,696}$). It is possible to mitigate this challenge by vigorously pruning the corpus—that is, removing rare words, grammatical terms, overly common words, and so forth—but this may not be enough. The most productive route forward would be a random sampling of word pairs followed by expert semantic evaluation, a project that our research team is currently planning. 25 30 35

⁶It is worth noting that, when it comes to endemic function, the pattern of sharp fall-off as we lose semantic specificity is a bit muddled in our full dataset. See Appendix 6 online for an analysis of this result.

STUDY 2: APPLYING SEMANTIC BENCHMARKS TO XIN-BODY RELATIONS

In Study 2, we determined collocation measures and patterns for *xin* and other organs in the body in relationship to the three standard terms for “body” in classical Chinese (*shen* 身, *ti* 體, and *xing* 形). Unlike the semantic pairs above, we also gathered this data separately for three different genre groupings within our corpus: all texts, all texts except for medical texts, and medical texts alone.⁷ Whatever early Chinese views about mind-body dualism are, we expected there might be a genre effect on the degree to which *xin* was portrayed as a physical organ in the body, with this being more likely in medical texts, given their technical nature and abundance of physiological terms. Moreover, the frequency of certain organ terms—most notably, “vein/artery/meridian/pulse” (*mai* 脈)—is much higher in medical texts than in the corpus as a whole,⁸ which could skew the results.

The results of our analysis, broken down to isolate the potential effect of the medical text genre, are reported in Figure 1 below (see Appendix 5 for the same results in table form).

It is immediately apparent that *xin* looks very different from the other organs. In both the combined corpus and the nonmedical corpus, its t-score is almost double that of the next highest organ. The only place this is not true is in the medical texts, where the term *mai* 脈 takes *xin*'s place as the odd-organ-out, having almost double the t-score of *xin* or stomach/belly (*fu* 腹). Typically translated as “vein, artery, meridian, pulse,” *mai* is a central term in traditional Chinese medicine, referring to the channels through which vital energy (*qi* 氣) flows in the body.

Tunnelling back down in the actual passages behind these collocation results, we can see that, in the medical texts, *mai* 脈 often appears in conjunction with the body terms, sometimes in ways that suggest a complementary or contrastive relationship. However, we may also be seeing the effect of the occasional occurrence of *mai* in the compound *maixing* 脈形, which refers to the “shape of the pulse” that is used to diagnose ailments. For example, in the Han medical text *Treatise on Cold Injury* (*Shang Han Lun* 傷寒論), a student asks, “When a person is sick from intense fear, how does their pulse present itself?” (*maizhuang* 脈何狀). The teacher replies, “The shape of the pulse (*maixing* 脈形) is like following a thread as it winds around and around, and their face is

⁷Medical texts contain a total of 3.94% of our corpus's 5.74 million characters.

⁸*Mai* 脈 appears 2,135 times in the four medical texts, but only 227 times anywhere else in the entire CTP corpus.

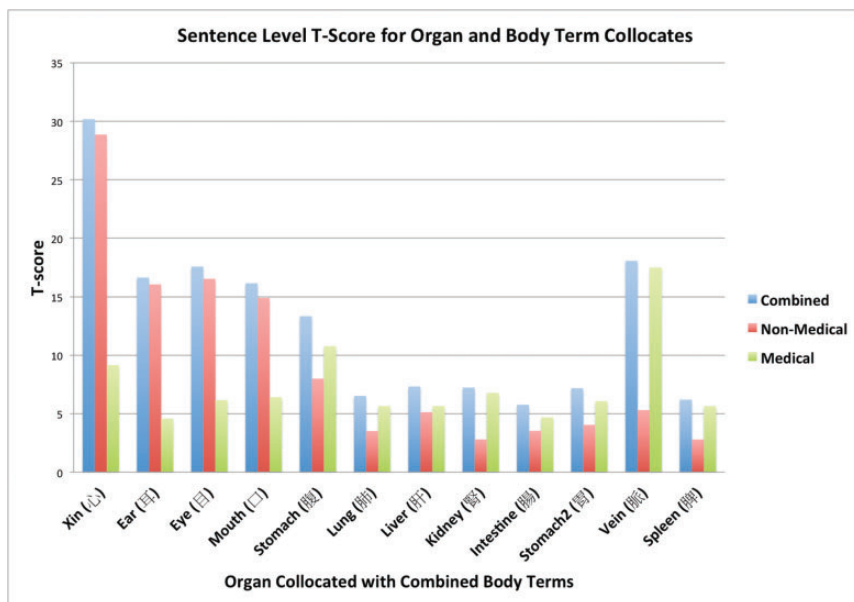


Figure 1. Visualization of Xin and other organs :: the body, split by genre groupings

white and devoid of color.”⁹ *Xing* 形 in this case has its basic meaning of “shape” rather than “body,” so collocations between the terms in such cases would be a false signal. Overall, we think the pair’s observable collocation patterns are best attributed to the unique focus of this genre, which is to manipulate the channels of vital energy in a patient to restore them to health. It is worth noting that a special role for *mai* entirely disappears once the medical texts are excluded and is greatly diminished in the overall corpus as a whole.

Turning back to *xin*, we can compare its pattern of collocation with body terms in the nonmedical corpus, as opposed to the other organ terms, with our semantic benchmarking efforts in Study 1. As Figure 2 below indicates, the *xin*-body relationship looks more like a strong, well-defined semantic relationship than any of the other organ-body collocations patterns, which look more like the generally semantically related control pairs.

The overall pattern gives the strong sense that *xin* and the body enjoy a special, specific semantic relationship, although all of the organ terms share the same broad semantic space with the body terms.

⁹From the “Methods for Determining Standard Pulses” (*pingmaifa* 平脈法), chapter 7.1.

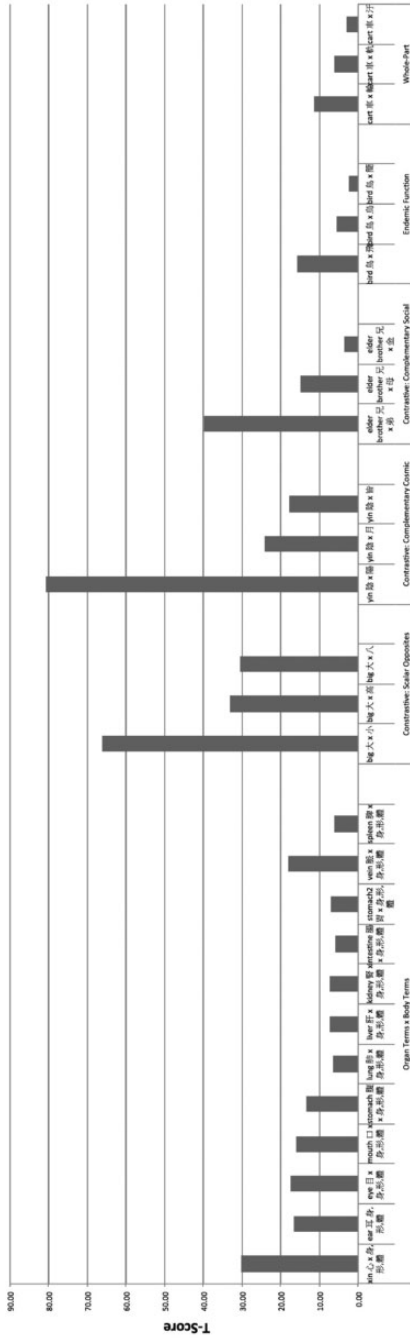


Figure 2. Xin and other organs::the body; combined corpus, with semantic benchmarking pairs

Of course, the particularly strong semantic relationship between *xin* and the body might be part-whole rather than contrastive, which could be seen as corroboration for the mind-body holist position. We think this a poor inference from our results for a variety of reasons. First of all, if the semantic relationship being picked up by the t-scores were of a part-whole nature, it would make sense for it to be shared by *all* of the organs equally, which is not what we see. A basic tenet of the mind-body holist position is that the *xin* is merely one organ among the others. However we interpret the precise nature of the semantic link between *xin* and body, it is—in contrast to the holist prediction—clearly of a qualitatively different order than any other organ.¹⁰

In sum, the ability to identify precise semantic relationships in classical Chinese from patterns of collocation scores alone remains elusive, although a much larger-scale exploration of the corpus may bring progress on this front. What we believe we *have* been able to demonstrate here, though, is that specific and strong semantic relations yield higher t-scores than vague ones, and that *xin*, alone among the organs, is characterized by just such a collocation profile vis-à-vis the physical body terms. This is difficult to reconcile with the claim that *xin* is simply one organ among many, as the strong mind-holist position would assert. Below we turn to other methods of automated textual analysis that strongly corroborate these results.

STUDY 3: HIERARCHICAL CLUSTERING ANALYSIS

A complementary, and in some ways even more exciting, approach to word co-occurrences involves employing unsupervised machine learning to extract significant patterns (Jurafsky and Martin 2015; Manning, Raghavan, and Schütze 2008; Plasse et al. 2007). The first of these methods that we applied to the issue of mind-body dualism in early China is called hierarchical clustering analysis. There are a variety of methods for performing such analyses, but the most common (and the one we employed) is referred to as bottom up, agglomerative hierarchical clustering.¹¹ In this method, the corpus is first converted into a “vector space”—which, in our case, is essentially an enormous multi-dimensional table, with each row representing an individual document and each column representing an individual term. An algorithm runs through the space, measuring the

¹⁰The possibility that the *xin*-body relationship is one of part-whole is further weakened by a follow-up study we performed comparing *xin* to two other common organs, the eye and the ear, on a series of collocation measures: the body terms, an endemic function term, a semantically-related term, and a semantically-unrelated term. See Appendix 7 online for these results and a discussion.

¹¹See Appendix 8 online for technical details.

geometric distances between individual terms. It then begins clustering them together in an iterative manner. The two terms with the shortest distances are “agglomerated” into a group, which then becomes a unit for the next stage of agglomeration, until the algorithm has built up a set of hierarchical clusters (clusters within clusters). The results are typically represented in a tree form or “dendrogram.”

Figure 5 below, from a hierarchical clustering analysis of a large contemporary English corpus targeting various classes of nouns (Rohde, Gonnerman, and Plaut 2006), shows how hierarchical clustering analysis can do an impressive job of tracing how individual terms are related to one another semantically, producing a recognizable conceptual map of the corpus.

Approaches such as hierarchical clustering are called “unsupervised,” because they involve a type of machine learning in which algorithms explore a set of completely unlabelled or unclassified data and attempt to identify statistically significant clusters based on a single, or small set, of parameters, such as corpus distance. The great advantage of unsupervised approaches is that they are as objective as one could desire. Although various assumptions are built into the processing of the document—specifically, the parameter or parameters selected, and the specifics of the algorithm—these can be easily varied and the resulting patterns compared. The running of the program involves no human input, and, in our case, was performed by a colleague with no knowledge of classical Chinese. This greatly reduces the potential for interpretative bias.

For Study 3, we began by producing a dendrogram representing relations between some of the control terms¹² from Study 1.¹³ The results are striking. With only a few exceptions,¹⁴ the tree relations are precisely what we expect given the well-understood semantic relations between these terms. For instance, the seasons cluster together, with the two classic opposites (summer and winter) sharing the most basic node, but then joining with spring the next node up, and finally *yue* 月, which refers to both the “moon” and “month.”¹⁵ We also see most of the scalar opposites and complementary pairs (*yinyang* 陰陽, heaven/earth *tiandi* 天地, king/minister *wangchen* 王臣, above/below *shangxia* 上下) clustering tightly together. A few surprising tight pairings (e.g., “east” *dong* 東 and “many”

¹²The hierarchical clustering algorithm was run before we had finalized the details of Study 1, and so the control terms differ somewhat from what we have reported above. Since the point is merely to validate the methodology’s ability to identify semantically coherent clusters, we decided not to rerun the study, given the enormous time and computational power required.

¹³The results are visualized in a dendrogram in Appendix 9 online.

¹⁴For a discussion of the exceptions, see Appendix 8 online.

¹⁵It should be noted that “autumn” was inadvertently dropped from the analysis.

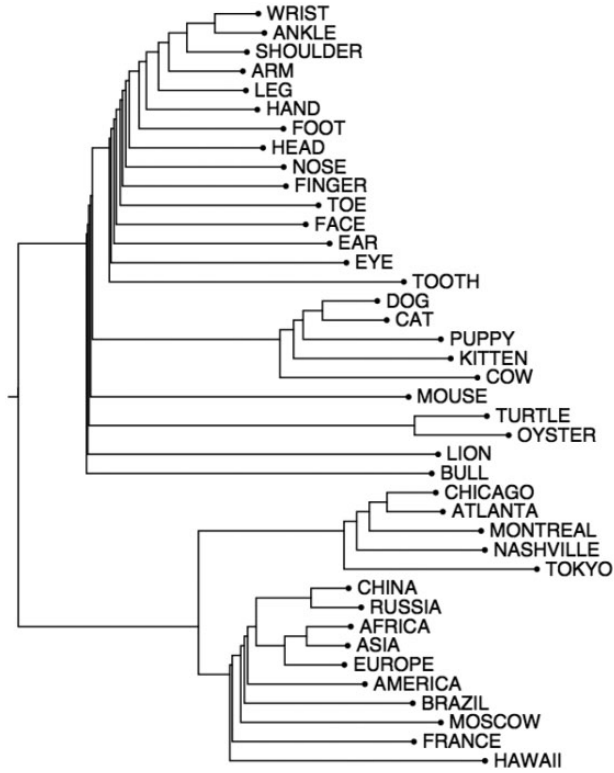


Figure 3. Dendrogram of noun classes based on vector distances in an English corpus (Rohde et al. 2005, 20, Figure 9)

duo 多) or great divides (e.g., “eye” *mu* 目 and “see” *jian* 見) aside, our hierarchical clustering algorithm appears to give an accurate model of expected co-occurrences within the early Chinese corpus, which, in turn, presents us with a coherent conceptual map of semantic relationships.

Having established a benchmark, let us now turn to the dendrogram representing our controversial terms of interest, namely, *xin* and the other organs in relation to the three primary body terms. This tree is represented in Figure 4 below:

It is difficult to imagine a clearer representation of mind-body dualism than Figure 4. The nodes pictured in the middle-shaded grey show the *xin* as being uniquely paired with the first of the body terms, *shen* 身, and then, in the next node out, also uniquely paired with the *other* two body terms, *xing* 形 and *ti* 體. This mind-body nexus then clusters with the

5

10

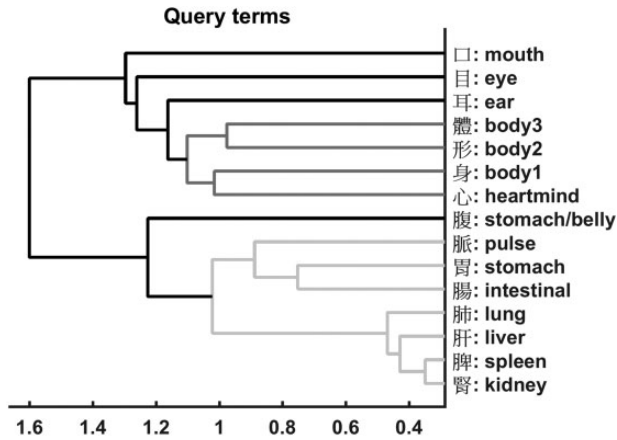


Figure 4. Dendrogram of query terms

three organs most associated with communication and perception, the mouth, eye, and ear. In other words, these are the three organs that most directly serve the *xin* in its role as the center of cognition and perception. Finally, what we might think of as the more “physiological” organs all cluster together in an entirely separate tree. The “stomach/belly” *fu* 腹 is something of an outlier in Figure 4. (Note that this is the same pattern we found in our t-score results above in Study 2, which provides an important confirmation that our mixed quantitative methods are converging on similar results.) This is possibly because of its occasional figurative usage as a metonym for basic desires or simple needs, as in *Daodejing* 12: “the sage is for the belly, not the eye” (*shengren wei fu, bu wei mu* 聖人為腹, 不為目).

It is important to reiterate that the hierarchical clustering algorithm, which was run on the entire corpus, without stop words, reproduces almost exactly the word collocation results reported in Study 2 for the entire corpus (Table 3 and Figure 3 above). The dendrogram in Figure 5 serves as a nearly perfect visual representation of our earlier t-score results relative to the body terms. There, the ear, eye, and mouth cluster closely together with t-scores in the neighborhood of 17, but are distinguished from *xin* with its t-score of 30. The “stomach/belly” (*fu* 腹) then appears at one remove with a t-score of 13, with all of the other organs representing a distinct, and internally tightly integrated, cluster with t-scores in the 5 to 7 range. The one exception is the *mai* 脈, which does appear as an outlier in the diagram vis-à-vis the “physiological” organs, but which we would expect, given its t-score of 18 above, to be joining the ear, eye, and mouth in the upper part of the figure. The difference here between the dendrogram

and t-score results most likely reflects the fact that in the textual vector space of the corpus, *mai* 脈 is rarely encountered but densely clustered where it does appear. It is highly represented in the corpus (2,362 appearances), but these are almost all concentrated in a small number of medical texts (2,135 appearances). One methodologically significant conclusion that we can draw from the contrast between the t-score and hierarchical clustering results in this particular case is that hierarchical clustering seems to do a better job of putting words in their proper place, as it were. T-scores alone fail to communicate the sometimes extremely lumpy distribution of certain key terms, and may, therefore, distort their relation to other terms of interest. Another equally important conclusion, however, is that the overall tight fit between the hierarchical clustering results and the t-score results should increase our confidence in t-scores as a reliable measure of collocation, at least in this classical Chinese corpus.

Our hierarchical clustering algorithm, and the dendrogram it produced, represents patterns of geometrical distances between terms within the Chinese Text Project's early Chinese corpus. Alternative semantic interpretations of this data are possible, but frankly difficult for us to imagine. Study 3's results appear to represent a confirmation of the special status of *xin*, its unique relationship to the body, and its special connection to perception and communication. In other words, like our earlier studies, Study 3 strongly confirms the view that, at least in terms of implicit, background assumptions, the authors of this corpus of early Chinese texts were mind-body dualists.

STUDY 4: TOPIC MODELING

Topic modeling is another unsupervised method that uses a complex form of statistics—Bayesian probability—to discern latent patterns of regularly co-occurring terms in a textual corpus.¹⁶ These patterns are called “topics.” Topic modeling begins with the assumption that the surface structure of the texts in a given corpus can be viewed as the product of

¹⁶The most commonly-used method for topic modeling in the humanities (and the method that we employed) is called Latent Dirichlet Allocation (LDA). A reasonably non-technical introduction to LDA from one of its creators, David Blei, can be found in (Blei 2012); perhaps more useful for most humanists is (Brett 2012) and a blog post by one of the pioneers in digital humanities, Ted Underwood (<https://tedunderwood.com/2012/04/07/topic-modeling-made-just-simple-enough/>). Helpful special journal issues include Volume 2, Number 1 (Winter 2012) of the *Journal of Digital Humanities*, which includes some application pieces in addition to Brett 2012 and Blei 2012. Also see the contributions to the special issue of *Poetics* 41.6 (December 2013), especially the introduction to topic modeling by (Mohr and Bogdanov 2013).

latent, or hidden, themes, and its task is identifying these themes (topics), as well as identifying the individual words that belong to these topics. Each topic thus produced consists of a list of words ordered by “weight” in the topic, with words at the top of the list contributing more to the formation of the topic than words lower on the list.

The use of topic modeling in the humanities is still in its infancy and to date has been used primarily in literary studies and political science.¹⁷ Within religious studies, a related analytic technique, principal component analysis, was used by a group of scholars in Taiwan to resolve a controversy concerning the authorship of various translations of Indian Buddhist texts into Chinese (Hung, Bingenheimer, and Wiles 2010). Topic modeling has also been applied to classical Chinese corpora to extract topics related to positive and negative emotions in Tang poetry (Hou and Frank 2015). With regard to the CTP corpus, some of our team has also been experimenting with using topic modeling to explore dating and authenticity controversies surrounding early Chinese texts, such as the *Shu Jing* or *Zhuangzi* (Nichols et al. In Press).

In a 100-topic model of the CTP corpus that we created, *xin* appears in 6 of the topics and is the primary component of one the “heaviest loading,” or most common, topic in the entire corpus, topic #97. The three conceptual topics¹⁸ in which *xin* appears, in order of their overall weight in the corpus, and with the top 10 loading words (ranked by importance), are listed below in Table 2.

Topic #97 is the most important for *xin*, since it is the second-most heavily weighted in the entire corpus, and *xin* constitutes its most important term. We have characterized topic #97 as “cognition/perception/cosmic fortune,” since the most heavily weighted words in the topic are the first three: *xin*, *jian* 見 (“to see/perceive”), and *ming* 明 (“bright, intelligent, clear”).¹⁹ The main focus of the topic seems to be cognition and perception, with—significantly—no mention of emotion. We should also note that no other organs are mentioned, not even the ones most closely associated with perception, such as the eye (*mu* 目) or ear (*er* 耳). The

¹⁷See, for instance, an analysis of a Texas newspaper article archive (Torjet and Christensen 2012) of an 18th-century midwife’s diary. (<http://www.cameronblevins.org/posts/topic-modeling-martha-ballards-diary/>), and literary themes in 19th century literature (Jockers and Mimno 2013b). Topic modeling has also been used as a more effective method than simple keyword searches for turning up themes of interest in massive, relatively unknown corpora (Tangherlini and Leonard 2013).

¹⁸The other three topics in which *xin* appears are “stylistic” topics, unique clusters of specialized terminology or grammatical particles that are distinctive to a particular text. They are described and discussed in Appendix 10.1 online.

¹⁹Word clouds visualizing the relative weights or contributions of each term to the topic in an intuitive manner are available in Appendix 10.2 online.

Table 2. Conceptual topics in which *Xin* appears

Topic #	Weight	Name	Top Ten Words (in order, left column first)	
97	0.47514	Cognition/ perception/ cosmic fortune	<i>xin</i> 心 see/perceive 見 bright/intelligent/ clear 明 accord/ harmonize 合 lose/miss 失 now 今 <i>xin</i> 心 after 後 strength/effort 力 worry 憂 person 人 big/great 大 Heaven/sky 天 know/ knowledge 知 king 王	peace/balance 平 yang 陽 intention 意 spirit 神 fortune/luck 福 interrogative 豈 morning/court 朝 death 死 sincerity/integrity 誠 abandon 棄 get/obtain 得 world/era/generation 世 one/unified – <i>xin</i> 心 stop/already 已
10	0.34877	Temporal cognition and planning		
33	0.07733	Human, heaven, and political order		

secondary references to according or harmonizing with things (*he* 合), missing an opportunity or making a mistake (*shi* 失), peace (*ping* 平), good fortune (*fu* 福), and intention (*yi* 意) all suggest a connection with planning or navigating the world, whereas the mention of spirits (*shen* 神) and yin-yang (*yinyang* 陰陽)²⁰ suggests that cosmic forces are some of the variables to be considered. In any case, the most important topic in which *xin* forms a major component seems to reflect a worldview that sees it as the sole seat of perception, cognition, planning, and personal responsibility, which in turn fits with a mind-body dualist account.

Interestingly, the second-most important topic for *xin*, topic #10, also seems to center on similar themes. We have termed this “temporal cognition and planning.” The topic is heavily dominated by “now/today” (*jin* 今), followed by *xin*, “after/future” (*hou* 後), “strength/effort” (*li* 力), and “worry/concern” (*you* 憂). Like topic #97, the focus seems to be on *xin*’s role in thinking about the future, planning, and exerting effort. A commonly used interrogative (*qi* 豈) and mention of “sincerity” (*cheng* 誠) suggests interior thought and resolve, whereas “death” (*si* 死)

²⁰Although *yang* 陽 appears in the top ten characters, *yin* 陰 is not far behind at fourteenth place.

hints at potential dire consequences.²¹ “Morning” (*chao* 朝; also “court,” as in royal court) adds to the sense of cognition stretched over time. It is worth noting that both #97 and #10 are distributed quite widely across the corpus.

Xin plays a relatively minor role in topic #33, appearing ninth and having a weight of only 0.017, as opposed to .052 for the first term, “person/human” (*ren* 人). To the extent that it is conceptually coherent, this topic seems to concern humans, heaven, and political order, with knowledge (*zhi* 知) also playing a role. This topic also serves to underscore *xin*’s primarily cognitive nature.

Perhaps the most salient result of our topic modeling study is that, in the nonmedical corpus, *xin* is the *only* bodily organ²² to appear among the most important characters in any of our 100 topic models. In other words, it is the only organ conceptually, or stylistically, salient enough to appear in distinctive thematic clusters. This is yet another example of the qualitative uniqueness of *xin* and further evidence against the accuracy of strong mind-body holist claims about early China.

CONCLUSION

Despite some minor methodological concerns—which we report, warts and all, in the interest of methodological transparency—our results are quite robust, especially because they come to similar conclusions by means of very different methodologies. Whether we are looking at word collocation measures, hierarchical cluster analyses, or topic models, *xin* stands out as entirely, qualitatively unique among the organs. The locus of the most important of human capacities—thought, planning, and decision-making—it shows a distinctive and highly salient relationship to the physical body, one that makes little sense unless the authors of the texts in which *xin* appears were operating against a background assumption of at least “weak”—even possibly subconscious—mind-body dualism.

Our results clearly contradict the position, held by scholars such as Jane Geaney (2002), that the heart-mind is simply one organ among many in the body with its own particular functions, but not otherwise distinctive. They also undermine similar holist positions, such as that advanced by A. C. Graham with regard to pre-Buddhist Chinese thought

²¹“Crime/guilt” (*zui* 罪) also appears thirteenth in the topic.

²²The only other organ that appears in any of our 100 topics is *mai* 脈, which shows up, as one might expect, in two wonderfully coherent “Traditional Chinese Medicine” topics (#73 and #84), that consist almost entirely of technical, medical terminology and load almost exclusively in the medical text portion of the corpus. See Appendix 10.3 online for a discussion of *mai*, as well as Appendix 10.4 online for mention of an apparent organ term appearing in a minor topic that turned out to be a graphic variant of a semantically unrelated term.

(Graham 1989, 25),²³ that acknowledges a special role for *xin*, but no special ontological status—that is, that *xin* is first among equally physical organs in a materialistic, monist universe. If the heart-mind were not conceived of, at least *implicitly*, as metaphysically distinct in some sense from the other organs, why would it be disproportionately collocating with the body terms? If it were merely more important, it might appear in the texts more frequently than the other terms, but should not (correcting for overall frequency, as our method does) vary from the other organs in terms of its collocations with body terms. 5

Moreover, linking the results reported here into the broader literature of the topic of mind and body in early China (e.g., Goldin 2003, 2015; Slingerland 2013), it is clear that the *xin*, uniquely among the organs, is associated with terms such as *shen* 神 (“spirit”), and that the dead or nonhuman supernatural spirits are often portrayed as enjoying the possession of a *xin*, as well as the functions that go along with *xin*, in a way that would seem bizarre when it comes to other organs in the body, such as the intestines or the lung. An early Chinese reader, no less than a contemporary reader of an English translation, glides smoothly over references to ancestors, spirits, or gods knowing things about the world, or becoming angry, but would be stopped in his or her tracks by a spirit troubled by shortness of breath or bowel problems. This converges with contemporary experimental research on people’s intuitions about what functions of the self survive the death of the physical body (Cohen et al. 2011), and points to what we see as the final nail in the coffin of any strong mind-body holist position: the fact that human beings seem to be intuitive, “weak” mind-body dualists, perceiving minds as somehow distinct from, and independent of, the physical bodies that house them (Bloom 2004; Slingerland 2013). 10 15 20 25

Finally, we would like to end with a brief discussion of broader methodological concerns relevant to the present study. There are no doubt those among our colleagues in religious studies who view talk of collocation measures, KWIC windows, and hierarchical clustering algorithms as a sinister encroachment of the sciences upon the humanities and further evidence that the twilight of the humanities is truly upon us. Even some early practitioners and advocates of digital humanities have, more recently, begun portraying the movement as part of a “neoliberal” conspiracy to undermine the core mission of the university and transform humanistic scholars into disposal tech flunkies (Allington, Brouillette, and Golumbia 2016). We believe, on the contrary, that the judicious adoption 30 35

²³Thanks to an anonymous referee for pointing out that we need to respond to this alternative position.

of digital humanities techniques is simply a way for humanities scholars to employ the best techniques and theories available to answer questions that matter to them. “Distant” reading, employed as a supplement to qualitative analysis, can help us to situate our close reading, give us new perspectives on our corpora, and, perhaps, decisively tip the balance of evidence in hermeneutic disputes. 5

Back in the 1960s, the linguist Margaret Masterman described the potential for computers to analyze texts from a new perspective as a “telescope of the mind,” a powerful new tool whose true potential had yet to be scratched (Masterman 1962, 38). Over fifty years later, this potential remains underexploited, despite major gaps in our analytic ability that computer-assisted analysis can help to fill. In the field of religious studies, as well as in the humanities more generally, there is, in our opinion, a need for new methods for settling hermeneutic disputes, or at least for narrowing the scope of reasonable views. For instance, in a well-known work on early Chinese thought, David Hall and Roger Ames argue that their claims about broad “cultural determinants” in early China—for instance, strong mind-body holism—should not be subject to what they call “the Fallacy of the Counterexample” (Hall and Ames 1995, xv). That is, legitimate generalizations about trends in the corpus of early Chinese texts cannot be invalidated by isolated, unrepresentative counterexamples, and we should not allow ourselves to “become lost in the details” (Hall and Ames 1995, xv) to the point where we lose sight of general trends. The problem is that they present no clear criteria for determining what constitutes a genuine trend and what counts as an irrelevant counterexample. We argue that the sort of large-scale textual analysis techniques described here could be useful in this regard. They give us a way to pan out from the intricacies of individual passages and texts to gain a panoramic view of an entire corpus. Perhaps more importantly, they also allow us to support our generalizations about a given corpus with relatively objective evidence rather than mere assertion or argument from authority.²⁴ 10
15
20
25
30

In the 1970s and 80s, John B. Smith created one of the first tools for conducting such analyses, the Archive Retrieval and Analysis System, that originally ran on a mainframe computer accessed via remote terminal. Although Smith was a computer scientist, he was profoundly sensitive to the humanistic enterprise and saw his platform as a tool to help 35

²⁴Note the comment by Rockwell and Sinclair that literary scholars frequently employ “semi-quantitative words” such as “more” or “less” in their arguments, as in “There is a lot more discussion in *Frankenstein* about technology than other novels.” “Whether or not [claims like this] are right,” they observe, “we are making a claim that can be investigated by using quantitative tools to count words. That is why we should be beware of hard distinctions such as that between hermeneutical and quantitative methods” (Rockwell and Sinclair 2016, 41).

humanities scholars do their work better, not as a replacement for humanistic expertise. “Humanists have always been explorers,” he wrote. “They sail not the seas of water but on seas of color, sound, and, most especially, words” (Smith 1984, 20). To extend Smith’s analogy, contemporary sailors still rely on such venerable tools as the compass and anemometer, and base the bulk of their decision-making on their qualitative feel for the ocean, wind, waves, and ship. Nevertheless, no one seriously concerned with sailing effectively, especially on a long journey or far from shore, would turn up their nose at GPS, radar, or satellite-based weather forecasting. 5

We feel the analogy is apt for humanities scholars. Computer-assisted textual interpretation can help us to gain our bearings as we travel through massive textual corpora, allow us to evaluate our qualitative intuitions in light of quantitative data, reveal hidden “topics” or themes invisible to individual human readers, and help us to more rigorously distinguish between unrepresentative counterexamples and instances of broader trends. Like GPS or a marine weather forecast, they can be misused, especially if they allow unqualified novices to set out to sea under the illusion that they know what they are doing, with possibly disastrous consequences. Knowledgeable experts are required to evaluate whether the new tools are useful in a given context, or to answer a particular question. They are also needed to determine when the results of content-blind algorithms or abstract statistical measures are best ignored because of the complexity or make-up of the material to be analyzed. In the right hands, however, new tools—whether navigational or scholarly—are unqualified gifts. Besides the light that our studies have shed on debates concerning mind-body dualism in early China, we hope that we have succeeded in showing the potential for large-scale analytic techniques to augment our ability to understand and map meanings in religious or philosophical textual corpora, and to weigh in decisively on otherwise intractable scholarly debates. 10 15 20 25

REFERENCES 30

Al-Hejin, Bandar. 2015. “Covering Muslim Women: Semantic Macrostructures in BBC News.” *Discourse and Communication* 9 (1): 19–46.

Allington, Daniel, Sarah Brouillette, and David Golumbia. 2016. “Neoliberal Tools (and Archives): A Political History of Digital Humanities.” *LA Review of Books*, May 1. 35

Ames, Roger. 1993. “The Meaning of the Body in Classical Chinese Philosophy.” In *Self as Body in Asian Theory and Practice*, edited by Thomas Kasulis, Roger Ames, and Wimal Dissanayake, 157–77. Albany: State University of New York Press.

- Baker, Paul, Costas Gabrielatos, and Tony McEnery. 2013. "Sketching Muslims: A Corpus Driven Analysis of Representations Around the Word "Muslim" in the British Press 1998–2009." *Applied Linguistics* 34 (3): 255–78.
- Barrett, Justin L. 1998. "Cognitive Constraints on Hindu Concepts of the Divine." *Journal for the Scientific Study of Religion* 37:608–19. 5
- Biber, Douglas, and James Jones. 2009. "Quantitative Methods in Corpus Linguistics." In *Corpus Linguistics: An International Handbook*, edited by A. Lüdeling and M. Kytö, 1286–304. Berlin: De Gruyter.
- Blei, David M. 2012. "Topic Modeling and Digital Humanities." *Journal of Digital Humanities* 2. 10
- Blei, David, Andrew Y. Ng, and Michael I. Jordan. 2003. "Latent Dirichlet Allocation." *Journal of Machine Learning* 3:993–1022.
- Bloom, Paul. 2004. *Descartes' Baby: How the Science of Child Development Explains What Makes Us Human*. New York: Basic Books.
- Brett, Megan. 2012. "Topic Modeling: A Basic Introduction." *Journal of the Digital Humanities* 2 (1). 15
- Bullinaria, John, and Joseph Levy. 2007. "Extracting Semantic Representations from Word Co-Occurrence Statistics: A Computational Study." *Behavioral Research Methods* 39 (3): 510–26.
- Chinese Text Project. Available at www.ctext.org. Accessed September 1, 2012. 20
- Church, Kenneth, William Gale, Patrick Hanks, and Donald Kindle. 1991. "Using Statistics in Lexical Analysis." *Lexical Acquisition: Exploiting On-Line Resources to Build a Lexicon*, 115.
- Church, Kenneth Ward, and Patrick Hanks. 1990. "Word Association Norms, Mutual Information, and Lexicography." *Computational Linguistics* 16:22–29. 25
- Cohen, Emma, Emily Burdett, Nicola Knight, and Justin Barrett. 2011. "Cross-Cultural Similarities and Differences in Person-Body Reasoning: Experimental Evidence From the United Kingdom and Brazilian Amazon." *Cognitive Science* 35:1282–304.
- Dunning, Ted. 1993. "Accurate Methods for the Statistics of Surprise and Coincidence." *Computational Linguistics* 19:61–74. 30
- Evans, Jonathan. 2008. "Dual-Processing Accounts of Reasoning, Judgment, and Social Cognition." *Annual Review of Psychology* 59:255–78.
- Geaney, Jane. 2002. *On the Epistemology of the Senses in Early Chinese Thought*. Honolulu: University of Hawaii Press. 35

- Goldin, Paul. 2003. "A Mind-Body Problem in the Zhuangzi?" In *Hiding the World in the World: Uneven Discourses on the Zhuangzi*, edited by Scott Cook, 226–47. Albany: State University of New York Press.
- . 2015. "The Consciousness of the Dead as a Philosophical Problem in Ancient China." In *The Good Life and Conceptions of Life in Early China and Greek Antiquity*, edited by R. A. H. King, 59–92. Berlin: De Gruyter. 5
- Graham, A. C. 1989. *Disputers of the Tao*. La Salle, IL: Open Court.
- Granet, Marcel. 1934. *La Pensée chinoise*. Paris: La Renaissance du livre.
- Gries, Stefan Th. 2013. "50-Something Years of Work on Collocations: What Is or Should Be Next." *International Journal of Corpus Linguistics* 18 (1): 137–65. 10
- Hall, David, and Roger Ames. 1995. *Anticipating China: Thinking through the Narratives of Chinese and Western Culture*. Albany: State University of New York Press.
- Hou, Yufang, and Anette Frank. 2015. "Analyzing Sentiment in Classical Chinese Poetry." Proceedings of the 9th SIGHUM Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities, Beijing. 15
- Hung, Jen-Jou, Marcus Bingenheimer, and Simon Wiles. 2010. "Quantitative Evidence for a Hypothesis Regarding the Attribution of Early Buddhist Translations." *Literary and Linguistic Computing* 25 (1): 119–34.
- Jockers, Matthew L., and David Mimno. 2013a. "Significant Themes in 19th-Century Literature." *Poetics* 41:750–69. doi: 10.1016/j.poetic.2013.08.005. 20
- . 2013b. "Significant Themes in 19th-Century Literature." *Poetics* 41 (6): 750–69. doi: 10.1016/j.poetic.2013.08.005.
- Jullien, François. 2007. *Vital Nourishment: Departing from Happiness*, translated by Arthur Goldhammer. New York: Zone Books. 25
- Jurafsky, Dan, and James H. Martin. 2009. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. 2nd ed, *Prentice Hall Series in Artificial Intelligence*. Upper Saddle River, NJ: Pearson Prentice Hall.
- . 2015. "Vector Semantics." In *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition (3rd ed draft)*, edited by Daniel Jurafsky and James Martin. 30
- Kahneman, Daniel. 2011. *Thinking, Fast and Slow*. New York: Farrar, Straus, Giroux.
- Klein, Esther, and Colin Klein. 2011. "Did the Chinese Have a Change of Heart?" *Cognitive Science* 36:179–82. 35

- Landauer, Thomas, and Susan Dumais. 1997. "A Solution to Plato's Problem: The Latent Semantic Analysis Theory of Acquisition, Induction, and Representation of Knowledge." *Psychological Review* 104 (2): 211–40.
- Lee, John, and Tak-sum Wong. 2012. "Glimpses of Ancient China from Classical Chinese Poems." Proceedings of COLING 2012: Posters, Mumbai, India. 5
- Lévy-Bruhl, Lucien. 1922. *La mentalité primitive*. Paris: Alcan.
- Manning, Christopher, and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing*. 1st ed. Cambridge, MA: The MIT Press.
- Manning, Christopher D., Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval*. New York: Cambridge University Press. 10
- Masterman, Margaret. 1962. "The Intellect's New Eye." In *Freeing the Mind: Articles and Letters from The Times Literary Supplement during March-June, 1962*, edited by D. J. Foskett, 38–44. London: Times.
- Mautner, Gerlinde. 2007. "Mining Large Corpora for Social Information: The Case of Elderly." *Language in Society* 36 (01): 51–72. 15
- Miner, Gary, John Elder, Thomas Hill, Robert Nisbet, Dursun Delen, and Andrew Fast. 2012. *Practical Text Mining and Statistical Analysis for Non-Structured Text Data Applications*. 1st ed. Waltham, MA: Academic Press.
- Mohr, John W., and Petko Bogdanov. 2013. "Introduction—Topic Models: What They Are and Why They Matter." *Poetics* 41 (6): 545–69. doi: 10.1016/j.poetic.2013.10.001. 20
- Moretti, Franco. 2007. *Graphs, Maps, Trees: Abstract Models for Literary History*. London, New York: Verso.
- . 2013. *Distant Reading*. London: Verso.
- Nichols, Ryan, Kristoffer Nielbo, Edward Slingerland, Uffe Bergeton, Carson Logan, Scott Kleinman. In press. Topic Modeling Ancient Chinese Texts: Knowledge Discovery in Databases for Asianists. *Journal of Asian Studies*. 25
- Oakes, M. 1998. *Statistics for Corpus Linguistics*. Edinburgh: Edinburgh University Press.
- Paperno, Denis, Marco Marelli, Katya Tentori, and Marco Baroni. 2014. "Corpus-Based Estimates of Word Association Predict Biases in Judgment of Word Co-Occurrence Likelihood." *Cognitive Psychology* 74:66–83. doi: 10.1016/j.cogpsych.2014.07.001. 30
- Plasse, Marie, Ndeye Niang, Gilbert Saporta, Alexandre Villemot, and Laurent Leblond. 2007. "Combined Use of Association Rules Mining and Clustering Methods to Find Relevant Links Between Binary Rare Attributes in a Large Data 35

- Set.” *Computational Statistics & Data Analysis* 52:596–613. doi: 10.1016/j.csda.2007.02.020.
- Poli, Maddalena. 2016. “Me, Myself and I: The Notion of Self in the Zhuangzi.” MA, Asian and African Languages and Literatures. Venice, Italy: Università Ca’Foscari. 5
- Prentice, Sheryl, Paul Rayson, and Paul J. Taylor. 2012. “The Language of Islamic Extremism: Towards an Automated Identification of Beliefs, Motivations and Justifications.” *International Journal of Corpus Linguistics* 17 (2): 259–86.
- Ramsay, Stephen. 2011. *Reading Machines: Toward an Algorithmic Criticism*, Topics in the Digital Humanities. Urbana: University of Illinois Press. 10
- Rockwell, Geoffrey, and Stéfan Sinclair. 2016. *Hermeneutica: Computer-Assisted Interpretation in the Humanities*. Cambridge, MA: MIT Press.
- Rohde, Douglas, Laura Gonnerman, and David C. Plaut. 2006. “An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence.” *Communications of the ACM* 8:627–33. 15
- Rosemont, Henry, Jr., and Roger Ames. 2009. *The Chinese Classic of Family Reverence*. Honolulu: University of Hawai’i Press.
- Sampson, G., and D. McCarthy, eds. 2005. *Corpus Linguistics: Readings in a Widening Discipline*. New York: Continuum.
- Slingerland, Edward. 2013. “Body and Mind in Early China: An Integrated Humanities-Science Approach.” *Journal of the American Academy of Religion* 81 (1):6–55. 20
- Slingerland, Edward, and Maciej Chudek. 2011. “The Prevalence of Mind-Body Dualism in Early China.” *Cognitive Science* 35:997–1007.
- Slone, D. Jason. 2004. *Theological Incorrectness: Why Religious People Believe What They Shouldn’t*. Oxford, New York: Oxford University Press. 25
- Smith, John B. 1984. “A New Environment for Literary Analysis.” *Perspectives in Computing* 4 (2/3): 20–31.
- Tangherlini, Timothy R., and Peter Leonard. 2013. “Trawling in the Sea of the Great Unread: Sub-Corpus Topic Modeling and Humanities Research.” *Poetics* 41 (6): 725–49. doi: 10.1016/j.poetic.2013.08.002. 30
- Taishō Shinshū Daizōkyō 大正新脩大藏經. Available at http://21dzk.l.u-tokyo.ac.jp/SAT/index_en.html. Accessed January 10, 2017.
- Teubert, Wolfgang, and Anna Čermáková. 2007. *Corpus Linguistics: a Short Introduction*. New York: Continuum. 35

Thesaurus Linguae Graecae. Available at <http://stephanus.tlg.uci.edu/>. Accessed January 10, 2017.

Torjet, Andrew, and Jon Christensen. 2012. "Building New Windows into Digitized Newspapers." *Journal of Digital Humanities* 1 (3).

Van Norden, Bryan. 2008. *Mengzi: With Selections from Traditional Commentaries*. Cambridge, MA: Hackett Publishing Company. 5