

EEBO-TCP as a tool for integrating teaching and research

Heather Froehlich, Richard J Whitt, and Jonathan Hope

I. Introduction.

TextLab is a new course at the University of Strathclyde that introduces students to methods of digital text analysis. As part of the Vertically Integrated Projects initiative, students work in teams, testing newly developed software programmes to identify specific linguistic features of literary texts. The Vertically Integrated Projects initiative puts students from a variety of different backgrounds (from different faculties, from first year through to postgraduate study) and they work on a project together. This system, first developed at the Georgia Institute of Technology by Professor Ed Coyle, focuses on hands-on, interdisciplinary research-led teaching, wherein students contribute directly to ongoing academic research. Textlab serves as a Digital Humanities Vertically Integrated Project at Strathclyde, fostering interdisciplinary collaboration between the faculty of English Studies and the faculty of Computer and Information Science.

In collaboration with the Mellon-funded Visualizing English Print 1470-1800 project (VEP) between the University of Strathclyde, the University of Wisconsin-Madison, and the Folger Shakespeare Library, students enrolled in Textlab were introduced to the same tools that researchers at these institutions use. The Visualizing English Print project seeks to ask what makes some texts more or less similar to other texts, what do specific genres of written language have in common, and what are the functional properties of literary language from the beginnings of English print (at ca. 1470) to 1800. Using a combination of newly-developed software packages designed for literary-linguistic analysis, students and researchers are able to conduct computer-aided analyses of texts to better understand, classify and quantify large numbers of texts in ways which were previously much more difficult. Students are arranged into mixed groups of five. In the first iteration of Textlab, each group was assigned one Shakespeare play to focus on. Using and evaluating the tools developed for this process, students were able to make genuine findings about their assigned play. Through their research, they would log their progress individually (in a notebook) and as a group (on a wiki).

II. Software Programs.

Students enrolled in Textlab used three software programs, some of which were developed as part of the Visualizing English Print project. WordHoard, Docuscope and LATtice all show different information about the same texts, allowing students to dig deeply into their data.

WordHoard is a corpus-analysis software programme developed at Northwestern University by Martin Mueller. We found WordHoard to be particularly useful for first-time users of large-scale corpora by presenting the question “what happens when we compare one (or more) text(s) to another group of one (or more) text(s)?” It produces a log-likelihood analysis of words and lemma, and allows the user to see these words in the context of Shakespeare’s corpus.

Comparing frequencies in "Shakespeare Comedies" and "Shakespeare." 3,142 lemmata appeared at least 5 times in 1 work. "Shakespeare Comedies" contains 10,303 distinct lemmata in 264,273 occurrences. "Shakespeare" contains 17,609 distinct lemmata in 865,185 occurrences. The significance levels for the log-likelihood values are adjusted for the number of comparisons.

Lemma	Relative use	Log likelihood	Analysis parts per 10,000	Reference parts per 10,000	Anal cour
king	-	372.2 ****	4.16	18.96	110
you	+	262.3 ****	221.02	171.26	5,84
i	+	224.0 ****	424.90	359.55	11,7
she	+	179.9 ****	76.59	53.05	2,02
sir	+	153.8 ****	46.43	29.82	1,22
a	+	114.6 ****	214.44	181.24	5,66
master	+	111.5 ****	19.90	10.99	526
our	-	99.5 ****	24.33	36.76	643
their	-	95.5 ****	15.55	25.56	411
thy	-	92.1 ****	38.41	52.99	1,01
noble	-	76.7 ****	3.59	8.47	95
lord	-	72.8 ****	25.43	36.07	672
queen	-	65.4 ****	1.97	5.51	52
will	+	62.9 ****	140.65	120.66	3,71
love	+	61.9 ****	27.74	19.41	733
signior	+	59.7 ****	4.39	1.61	116
mistress	+	57.1 ****	9.31	4.98	246
roman	-	50.3 ****	0.19	1.73	5
death	-	50.0 ****	6.43	11.18	170
wit	+	48.9 ****	7.79	4.14	206
majesty	-	47.9 ****	0.98	3.27	26
dead	-	46.6 ****	3.18	6.60	84
fool	+	46.3 ****	9.84	5.74	260
count	+	45.9 ****	2.80	0.90	74
arm	-	45.3 ****	1.63	4.24	43
crown	-	43.6 ****	1.14	3.39	30
we	-	43.4 ****	55.74	67.30	1,47
ring	+	42.9 ****	4.09	1.73	108
be	+	42.8 ****	417.45	388.34	11,0
thou	-	41.3 ****	93.31	107.71	2,46
kingdom	-	41.1 ****	0.19	1.50	5
his	-	37.8 ****	69.40	81.34	1,83
slay	-	36.2 ****	0.72	2.43	19

Compress log-likelihood value range in tag clouds

Cloud

DocuScope is a rhetorical analysis software developed at Carnegie Mellon University by Suguru Ishizaki and David Kaufer (2004, 2011). When developing their program, Ishizaki and Kaufer encoded 175,000 of the most frequent words in English as well as the most frequent 2-4 word combinations into strings classified by rhetorical effect. Each category of rhetorical effect is broken down into specific functional features, called a "Language Action Type" or LAT. Each LAT is grouped

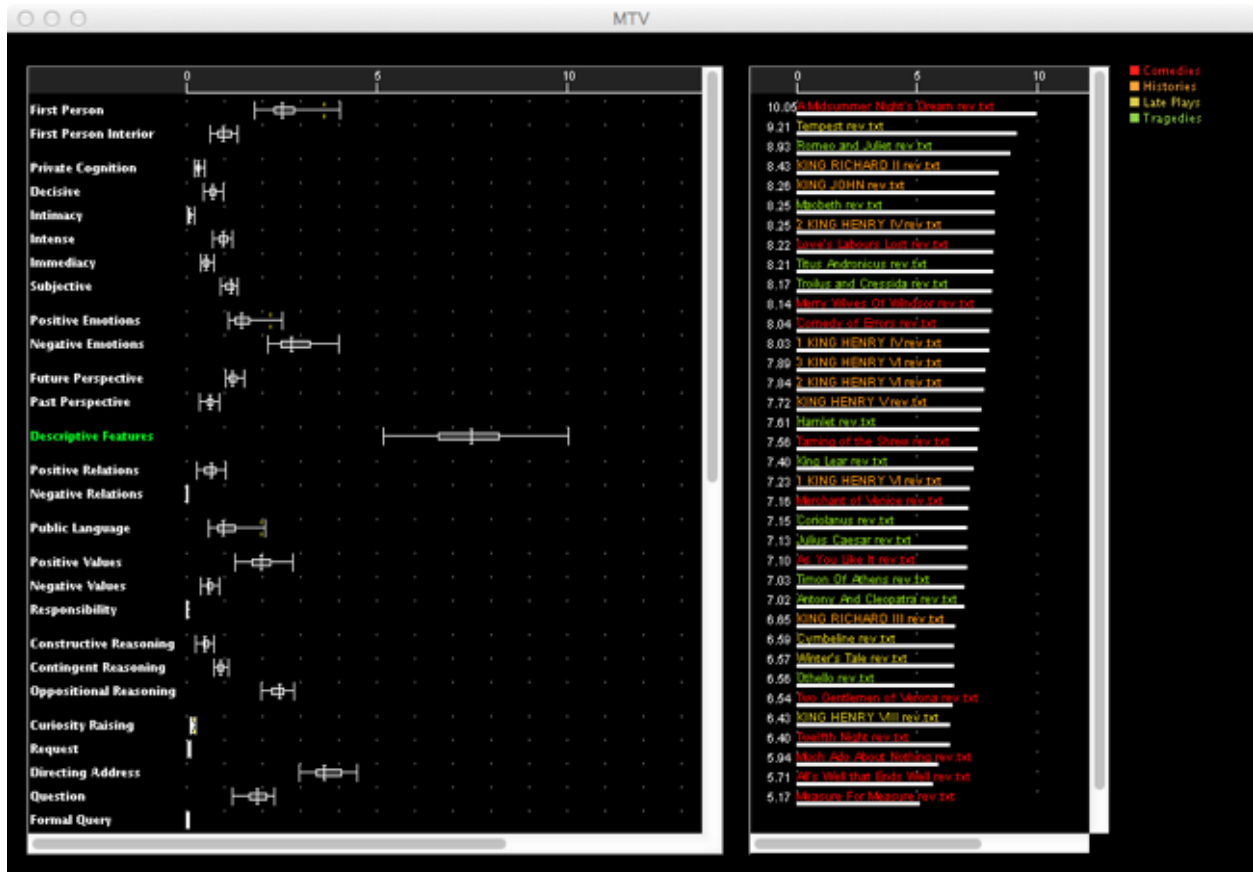
into a larger category (“Dimension”), which in turn are organised into larger, superordinate classes (“clusters”) as a hierarchy. Each LAT is underlined and color-coded for analysis. They did this by breaking rhetorical patterns of English into word-strings, and each word-string represents a part of a category.¹

Category	Count	Percentage
First Person Options	1163	4.45%
Interior Thought	825	3.16%
Emotions	960	3.68%
Time Orientation	522	2.00%
Descriptive	2124	8.14%
Interpersonal Relations	114	0.44%
Public Reference	171	0.65%
Public Values	460	1.76%
Reason	801	3.07%
Interaction	1524	5.84%
Topical Flow	532	2.04%
Elaborations	378	1.45%
Special Referencing	2382	9.12%
Person Roles	1072	4.11%
Person Properties	1072	4.11%
Communicator Roles	0	0.00%
Referencing Language	245	0.94%
Abstract Reference	969	3.71%
Citing References	96	0.37%
Citing Quotation	0	0.00%
Reporting	1543	5.91%
Directing	270	1.03%
Directing Readers	91	0.35%
Narrative	181	0.69%
Total words:	26108	

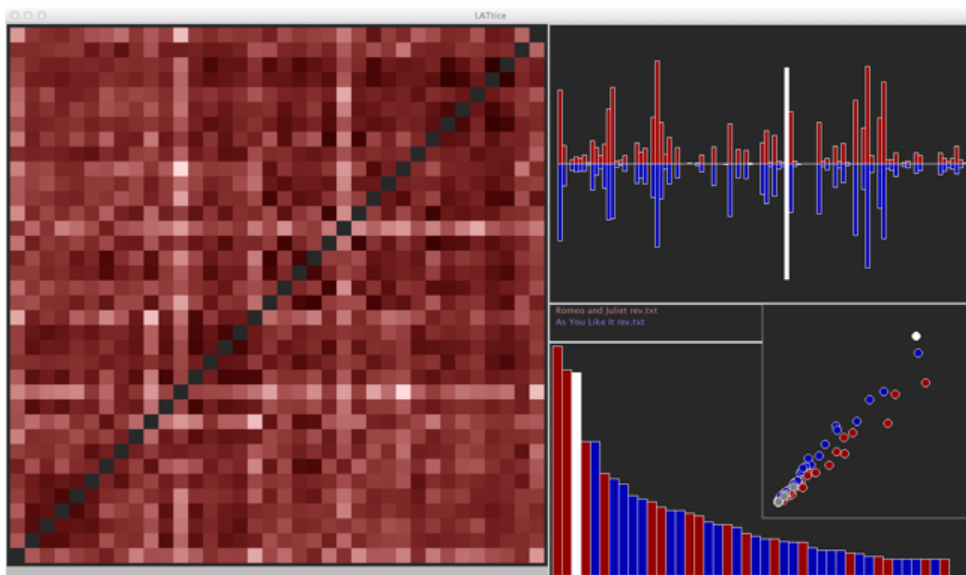
honour ! You will not do it , you ! I do relent : what would thou more of man ? Sir , here's a woman would speak with you . Let her approach . Give your worship good morrow . Good morrow , good wife . Not so , an't please your worship . Good maid , then . I'll be sworn , As my mother was , the first hour I was born . I do believe the swearer . What with me ? Shall I vouchsafe your worship a word or two ? Two thousand , fair woman : and I'll vouchsafe thee the hearing . There is one Mistress Ford , sir : --I pray , come a little nearer this ways : --I myself dwell with master Doctor Caius , -- Well , on : Mistress Ford , you say , -- Your worship says very true : I pray your worship , come a little nearer this ways . I warrant thee , nobody hears ; mine own people , mine own people . Are they so ? God bless them and make them his servants ! Well , Mistress Ford ; what of her ? Why , sir , she's a good creature . Lord Lord ! your worship's a wanton ! Well , heaven forgive you and all of us , I pray ! Mistress Ford ; come , Mistress Ford , -- Marry , this is the short and the long of it ; you have brought her into such a canaries as 't is wonderful . The best courtier of them all , when the court lay at Windsor , could never have brought her to such a canary . Yet there has been knights , and lords , and gentlemen , with their coaches , I warrant you , coach after coach , letter after letter , gift after gift ; smelling so sweetly , all musk , and so rushing , I warrant you , in silk and gold ; and in such alligant terms ; and in such wine and sugar of the best and the fairest , that would have won any woman's heart ; and , I warrant you , they could never get an eye-wink of her : I had myself twenty angels given me this morning ; but I defy all angels , in any such sort , as they say , but in the way of honesty : and , I warrant you , they could never get her so much as sip on a cup with the proudest of them all : and yet there has been pearls , nay , which is more , pensioners ; but , I warrant you , all is one with her . But what says she to me ? be brief , my good she-Mercury . Marry , she hath received your letter , for the which she thanks you a thousand times ; and she gives you to notify that her husband will be absence from his house between ten and eleven . Ten and eleven ? Ay , forsooth ; and then you may come and see the picture , she says , that you wot of : Master Ford , her husband , will be from home . Alas ! the sweet woman leads an ill life with him : he's a very jealousy man : she leads a very frampold life with him , good heart . Ten and eleven . Woman , commend me to her ; I will not fail her . Why , you say well . But I have another messenger to your worship . Mistress Page hath her hearty commendations to you too : and let me tell you in your ear , she's as fartuous a civil modest wife , and one , I tell you , that will not miss you morning nor evening prayer , as any is in Windsor , whoe'er be the other : and she bade me tell your worship that her husband is seldom from home ; but she hopes there will come a time . I never knew a woman so dote upon a man : surely I think you have charms , la ; yes , in truth . Not I , I assure thee : setting the attractions of my good parts aside I have no other charms . Blessing on your heart for't ! But , I pray thee , tell me this : has Ford's wife and Page's wife acquainted each other how they love me ? That were a jest indeed ! they have not so little grace , I hope : that were a trick indeed ! but Mistress Page would desire you to send her your little page , of all loves : her husband has a marvellous infection to the little page ; and truly Master Page is an honest man . Never a wife in Windsor leads a better life than she does : do what she will , say what she will , take all , pay all , go to bed when she list , rise when she list

Docuscope can run a functional-linguistic analysis on many texts at once using its Multiple Text Viewer, which shows each dimension as box plots and a roster table.

¹ Ishizaki and Kaufer continued creating LATs until they felt confident that they could categorise most of English into their existing LAT structure. The creators tested their rhetorical classifications against “as broad a range of English prose text as we could find” covering newspapers, fiction, short stories, speeches, journals, student writings from a multi-genre course; essays on social criticism, reminiscence and narratives; information system documentation including meeting minutes, design plans, public announcements, and proposals; talk radio, song lyrics, fables, Inaugural speeches, and the New Yorker, among other sources. Docuscope also covers the 100,000 frequent words and 75,000 of the most common words in English, as well as the most frequent 2-4 word combinations (2011: 280; 2004: Chapter 1).



And its output can be exported to LATtice, which produces a heatmap visually showing similarity and difference using dark and light color shading based on LAT frequencies in the texts run through DocuScope.



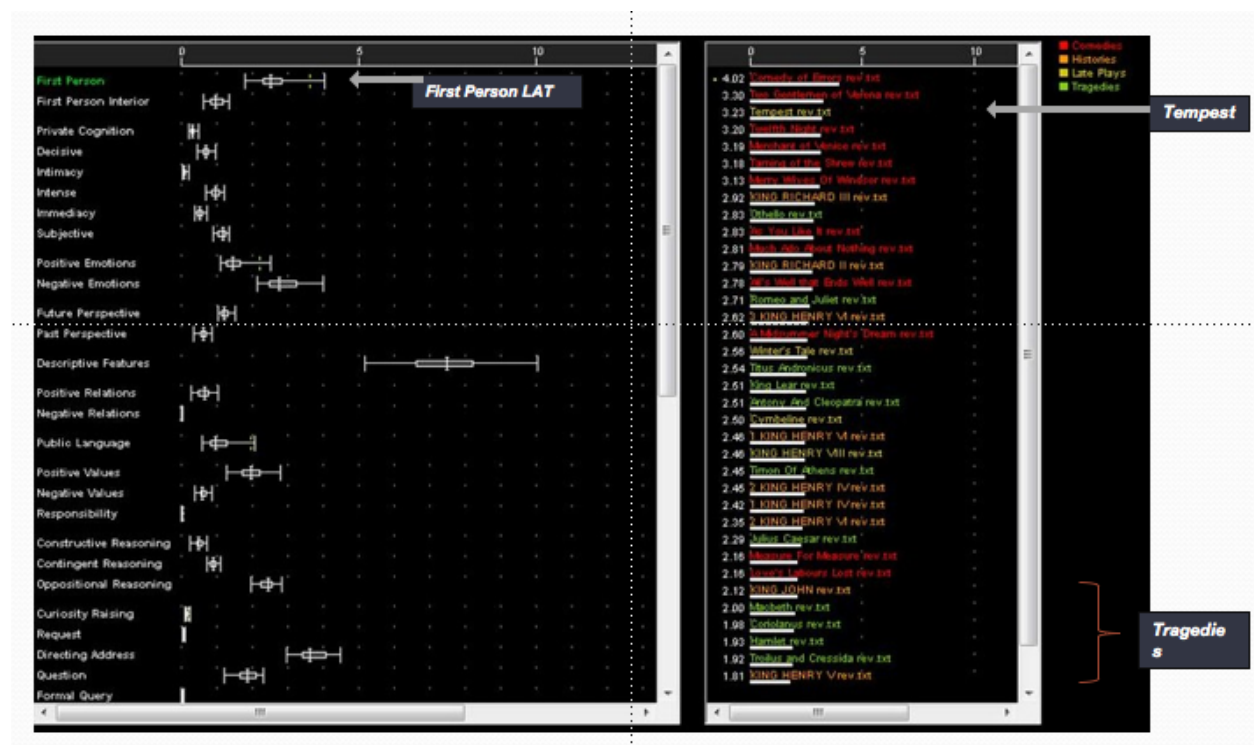
Anupam Basu's LATtice software takes the data values of each LAT from DocuScope and creates a heatmap of similarity or difference based on Euclidean distances between vectors. Using an $N \times N$ grid, LATtice illustrates difference or similarity through color coding: the darker a square is, the more similar the two texts are; the lighter, the less similar the two texts are. The black diagonal line

intersecting the graph represents a text compared against itself.² Relationships between any two texts in the corpus can be explored at the level of the LATs, as illustrated on the right-hand side of the software. This sort of mapping is ideal for digging into larger datasets, guiding the user to identify which texts should be looked at more closely.

III. Student Work from the First Iteration of Textlab

Students were made familiar with the software in a trickle-down approach: researchers presented the tools, and students enrolled in the course became comfortable with them by testing their assigned texts with them. This process encourages students to be active learners: they have a problem that they need to solve using the tools provided. What makes their assigned play different – or similar – in some way in comparison to other plays by the same author?

Jamie-Leigh Green, a fourth-year undergraduate student, discovers using WordHoard that “love” is less likely to appear in *The Tempest* compared to all other Shakespeare plays despite it being categorised as a “romance”. This early discovery guided her and her group’s work throughout the rest of the semester. Ultimately, they discovered that Shakespeare’s late plays are only loosely connected by this generic division – which definitely does not correlate to our modern conception of “romance”.



The group working on *The Two Gentlemen of Verona* discovered that *she* and *love* are especially important words in their play. A play about two men talking about a woman is not particularly surprising, but as Jessica Wagstaff points out, “language in *The Two Gentlemen of Verona* focuses heavily on the discourse of women as object[s] of attainment”, adding that the feminine pronoun *she* is far more likely to appear in *The Two Gentlemen of Verona*, but only as the object of the

² <http://winedarksea.org/?p=1285>

call this sampling of texts the “1K Corpus” (although there are slightly more than 1000 texts here; about 1080 to be exact).

These texts were then sent through Docuscope, which as discussed earlier, is one of the programmes used in Textlab. Initial output revealed that later texts were receiving significantly more tags than earlier texts, which pushed the need for some degree of modernisation. A variant detector software package (VARD 2, developed by Alistair Baron at The University of Lancaster⁴) was employed to modernise the spelling of the 1K Corpus for tagging purposes. After texts in the 1K Corpus were “Varded”, tagging frequencies were much more uniform across the board.

Since one of the main goals of the VEP project is to trace the development of genre through the history of English print, genre labels were then assigned to the texts in the 1K Corpus. The labels themselves are based on those used in the Early Modern section of the Helsinki Corpus and ARCHER (A Representative Corpus of Historical English Registers)⁵, with certain additional labels added to account for broader text selection. Labels were assigned based, in the first instance, on information made available in the metadata (e.g. keywords such as play, sermon, medicine, etc. easily give away a text’s genre). Where metadata were unavailable or proved insufficient, the text was inspected manually.

Using Docuscope output, a number of statistical tests—Ward’s clustering method, PCA (Principal Component Analysis) and ANOVA (Analysis of Variance)—were run using the JMP statistical software package to see if any correlation could be found between the LATs assigned by Docuscope and the genre assignments. Indeed, all three tests revealed independently that Docuscope output could be used to see what features of language (vis-à-vis the LATs) distinguish one genre from another, and what features cause statistically significant separations among the genres: the dendrogram—the result of Ward’s clustering—shows most dramas grouping together in the same cluster:

provides frequency results in percentages (relative to text size) rather than absolute numbers, we hope to have offset at least some disparities in text length.

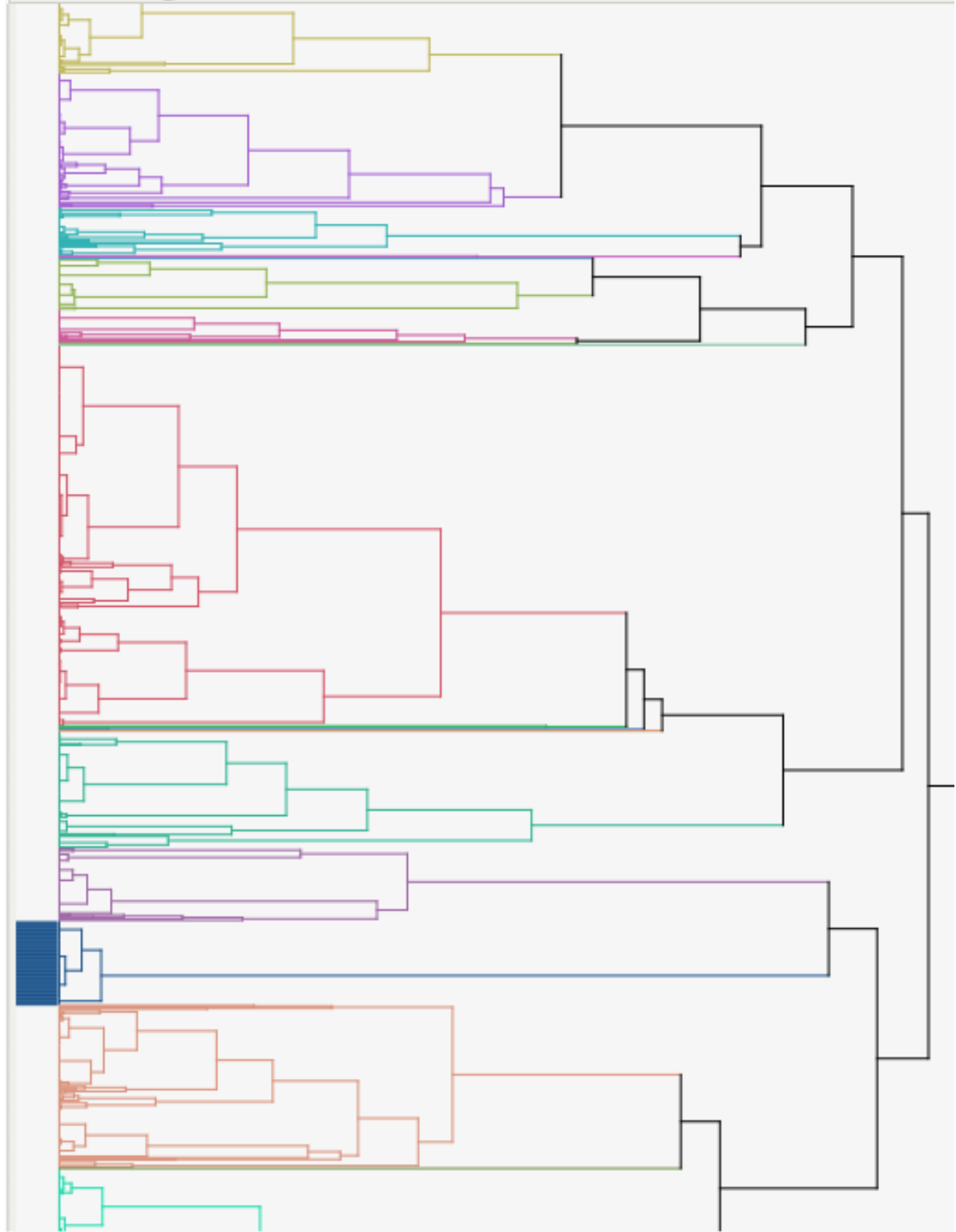
⁴For more information about the VARD 2 programme, see

<http://www.comp.lancs.ac.uk/~barona/vard2/>.

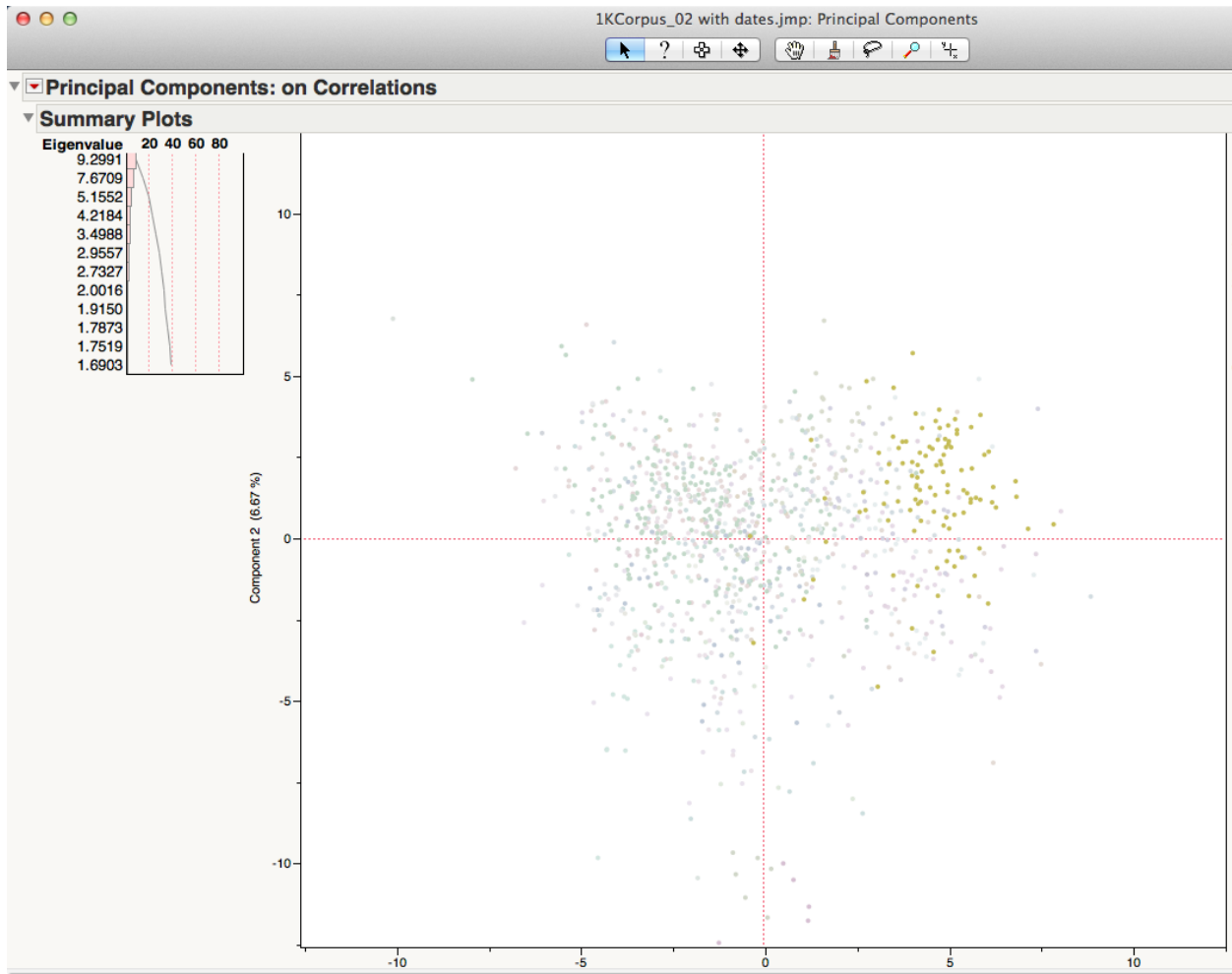
⁵For more information about these corpora, see

<http://www.helsinki.fi/varieng/CoRD/corpora/HelsinkiCorpus/earlymodern2.html> and http://www.llc.manchester.ac.uk/research/projects/archer/archer3_2/, respectively.

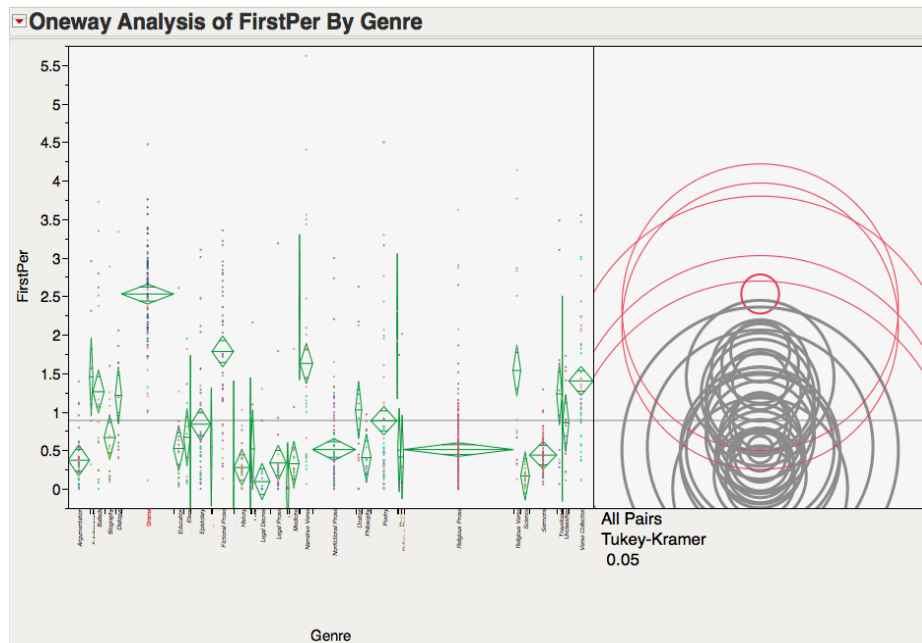
Dendrogram



The dramas also appear, for the most part, to group together in the same quadrant of PCA space:



And ANOVA reveals that the FirstPer LAT is one of key factors causing dramas to evince statistically significant separation from most other genres (dramas are represented by the highlighted magenta circle in the visualisation of the Tukey test, with grey circles representing those genres with statistically significant lower frequencies in the use of the FirstPer LAT):



This combination of statistical and (digital) literary analysis promises many opportunities for interdisciplinary student collaboration. On the one hand, students of literature can become acquainted with the principles of statistics as they inform quantitative aspects of literary computing. On the other hand, students of statistics and computer science can see how the methods of their disciplines can be applied to the study of literary texts. As the point of the VIPs is to integrate students from a number of diverse disciplines, the intersection of digital literary analysis and statistical computation seems an ideal avenue to foster such cross-disciplinary interfaces, while stressing transferable skills amongst participants.

V. Digital Humanities at Strathclyde

As members of the VEP found, modernizing texts makes them much more consistently comparable to each other. In order to teach with the TCP texts, they must be modernised, but after that process is complete, a much wider variety of texts can be used by students. Whereas the first iteration of Textlab focused on the analysis of Shakespeare's plays, future iterations can focus on a broader range of texts currently being used by the VEP project. Each group of students will be given a set of two texts of different genres and asked to use the software programmes to find the similarities and differences among both literary and non-literary genres. Possible questions include: how is the language of medicine different from the language of law? How are dramas and travelogues similar? What features differ academic prose from fictional prose? These results will obviously feed in to the larger VEP project.

Students enrolled in Textlab are primed to make genuine discoveries about texts in the early modern period in ways which were previously inaccessible. Students at Strathclyde are in a unique position to conduct this sort of work: Dr Anouk Lang teaches a complimentary course at Strathclyde in the first semester, Introduction to Digital Humanities.⁶ Her course introduces students to digital tools for text analysis and digital mapping, emphasizing digital literacy and critical thinking about digital tools with a strong emphasis on transferable skills from the classroom to real-life projects. Textlab serves as a continuation of this line of research-led teaching with a hands-on workshop component, continuing to bridge the gap between computer and literary studies.

This first iteration of Textlab shows that interdisciplinary, research-led collaborative work allows for students to actively contribute to ongoing research while building an impressive skillset on their CVs. Looking to the future, having students create and maintain generic classifications in large text corpora allows students to continue practice both close and distant reading – that is, looking at individual texts closely and as part of a larger whole. This idea of distance reading allows us as scholars to ask more questions using new methods, and pedagogical approaches will be quick to follow. As students have been practicing close reading for years, the shift away from close reading and towards considering texts as part of an as-yet unknown whole will allow them to continue to make genuine contributions to the future academic landscape.

⁶ For a full overview of Introduction to Digital Humanities, please see Dr Lang's webpage for the course: <http://aelang.net/wordpress/intro-to-dh/>

Works Cited.

Ishizaki, Suguru and David Kaufer. "Computer-aided Rhetorical Analysis" in *Applied Natural Language Processing and content analysis: Identification, Investigation, and Resolution*. Ed. Philip McCarthy and Chutima Boonthum. 276-296. 2011.

Kaufer, D., Ishizaki, S., Butler, B., & Collins, J. *The Power of Words: Unveiling the Speaker and Writer's Hidden Craft*. Routledge. 2004.