

# Body Language in Classic British Fiction: Words, $n$ -grams, and topics

TOMOJI TABATA,<sup>a)</sup>

**Abstract:** This pilot study takes a stylometric approach to investigate “body language” in classic British fiction. The specific research questions are how body-part words are distributed across registers, whether frequency patterns of body language make it possible to classify texts into meaningful sets, as well as what stylistic functions body-part expressions lend themselves to in fiction. To answer the questions, stylometric analysis is carried out in conjunction with topic modelling and qualitative interpretation of stylistics effects by paying close attention to individual words,  $n$ -grams, and topics.

**Keywords:** Body language, style, classic fiction, stylometry, topics

## 1. Introduction

This article was prepared as a handout for presentation given at the 119th SIG-CH (Computers and the Humanities) conference to be held at the Toyonaka Campus of the University of Osaka on the 16th February 2019.

A first step of text analysis often begins with identifying key words of a target author or a target corpus. Based on the assumption that key words encapsulate important lexical, semantic, thematic, or stylistic features of the author, we try to find words “that are evenly distributed across [her/his corpus] and display a higher frequency and a wider range than in a reference corpus of some kind” (Paquot and Bestgen, 2009). The most popular methods for detecting key words include application of statistical significance testings, such as chi-squared ( $\chi^2$ ), log-likelihood ratio ( $G^2$ ) (Rayson and Garside, 2000), Mann-Whitney’s  $U$ , and representativeness/distinctiveness ( $RD$ ) index (Klaussner et al. 2015), among others. While  $\chi^2$  and  $G^2$  can over-evaluate “bursty” words (those which occur extremely frequently in a very small number of texts, typically in a single text only), tending to give greater weights to high-frequency words due to their formulae,  $U$  and  $RD$  help to find lexical items “that are [more] evenly distributed across [the target set]”. As a result of comparing a Charles Dickens corpus with a reference corpus of classic British fiction using  $RD$ , for instance, a list of following key words can be obtained:

*head, down, old, corner, state, given, legs, fire, window, until, streets, night, for, outside, looking, air, dark, round, behind, return, shut, light, upon, than, hat, quiet, without, would, up, top, future, reference, wall, off, remarkable, to, determined, arranged, red, scarcely, not, but, heavy, bright, boy, till, slowly, street, neither, sky, expression, desired, stopped, through, with, dead, be-*

*side, only, dismal, sharp, leaning, windows, worn, nevertheless, holding, face, nor, must, bell, closed, patience, therefore, black, wasn’t, stopping, times, could, anybody, consciousness, position, faces, pleased, various, endeavour, use, curious, received, shake, gradually, being, feeling, iron, advice, effort, wind, inside, here, back, assured, hair*

Words in italics are the most representative and distinctive Dickens markers while words in Roman are consistently more overused in the reference set. What attracts attention in the Dickensian key words are body-part words such as *head, legs, face, faces*, and *hair*. Attentive readers of Dickens novels will associate such body-part expressions with description of actions, habits, emotions, and personalities of fictional characters. A team of researchers led by Michaela Mahlberg of the University of Birmingham are, in fact, working on characterization in the representation of body language from a corpus linguistic perspective (the CLiC Dickens project) (Mahlberg et al. 2016; Wiegand et al., 2017; Mahlberg and Wiegand, 2018).

The aim of this paper is to take a stylometric approach to investigate body language in classic British fiction with Dickens being the target author. The specific research questions are:

- (1) how body-part words are distributed across registers
- (2) whether frequency patterns of body language make it possible to classify texts into meaningful sets
- (3) what stylistic functions body-part expressions lend themselves to in fiction

To answer the questions, stylometric analysis is carried out in conjunction with topic modelling and qualitative interpretation of stylistics effects by paying close attention to individual words,  $n$ -grams, and topics.

## 2. Body language across registers

As a first step to list up as many body-part words as pos-

<sup>1</sup> GSLC, University of Osaka, Toyonaka, Osaka 560-0043, Japan

<sup>a)</sup> tabata@lang.osaka-u.ac.jp

sible, all the word-types in a corpus of classic British fiction were examined. The corpus entitled **Osaka Reference Corpus for Historical/Diachronic Stylistics (ORCHiDS)** is made up of 23 texts by Dickens, 24 major 18th Century novels, and 31 representative 19th Century fiction texts with a total tokens of 14,481,460 running words (Tables 2–4). Fig. 1 shows 150 lexical items found in the corpus.

abdomen, #abdominal, #abdominales, ankle, ankles, anus, arm, arm's, arms, arms', arse, artery, back, belly, bladder, blood, body, bone, brain, breast, breasts, buttocks, calf, cheek, cheeks, chest, chin, ear, ears, elbow, elbows, eye, eyes, eyebrow, eyebrows, eyelashes, eyelid, eyelids, face, faces, feet, finger, fingers, finger-nail, finger-nails, fingernail, fingernails, foot, forehead, foreheads, groin, gums, hair, hairs, hand, hands, head, heads, heart, hearts, heel, heels, hip, hips, instep, insteps, intestines, iris, jaw, jaws, kidney, knee, knees, leg, legs, ligament, ligaments, lip, lips, liver, lobe, lungs, mouth, mouths, muscle, muscles, nail, nails, navel, neck, necks, nerves, nose, noses, nostril, nostrils, organs, ovary, palm, palms, pinky, pituitary, #pore, pupil, ribs, #scalp, shin, shoulder, shoulders, skeleton, skin, skull, skulls, #sole, spinal, spine, spines, spleen, stomach, stomachs, teeth, tendon, tendons, thigh, thighs, #thorax, throat, throats, thumb, thumbs, tibia, #tissue, toe, toes, #toe-nails, tongue, tongues, #tonsils, tooth, torso, #uterus, #uvula, vein, veins, vertebrae, waist, waists, wrist, wrists

Fig. 1 150 body-pard words extracted from ORCHiDS

In an effort to answer the first research question, total frequencies (per 100,000 words) of body-part words were counted in each of fifteen registers in two 1,000,000-word balanced corpora: the FLOB corpus (British English) and FROWN corpus (American English). The two corpora were compiled according to the same corpus design in terms of the size and number of samples that make up fifteen registers: nine informative prose registers and six imaginative (fictional) prose registers. Table 1 gives a descriptive summary of the corpus design.

Table 1 The structure of FLOB/FROWN corpora

Category	Register	No. of texts
A	Press: Reportage	44
B	Press: Editorial	27
C	Press: Review	17
D	Religion	17
E	Skills, Trades and Hobbies	38
F	Popular Lore	44
G	Belles Lettres, Biographies, Essays	77
H	Miscellaneous	30
J	Science	80
K	General Fiction	29
L	Mystery and Detective Fiction	24
M	Science Fiction	6
N	Adventure and Western	29
P	Romance and Love Story	29
R	Humour	9
Total number of texts in the corpus:		500

The categories A–J are informative registers, and K–R fictional registers.

As Fig. 2 illustrates, body-part related words are predominantly more frequent in fictional registers in both British and American English. The overall distributional pattern in the two corpora is remarkably similar: the category most given to body language is Romance and Love story (P), with Humour (R) having the least recourse to parts of body in the fictional registers.

Of the 150 items given in Fig. 1, words with # were not found in the two corpora and 121 words are common in both corpora. Fig. 3 is a result of conducting Principle Component Analysis on a frequency matrix of the 121 body-part words across the 15 registers in the FLOB and FROWN corpora.\*1 Whereas informative registers form a closely-knit cluster, fictional registers are widely dispersed. Of further note is that fictional registers in the two national varieties of English have remarkably similar positions in the 3D diagram. To sum up, body language is a distinctive feature of fiction since depiction of body part plays significant roles in fictional settings.

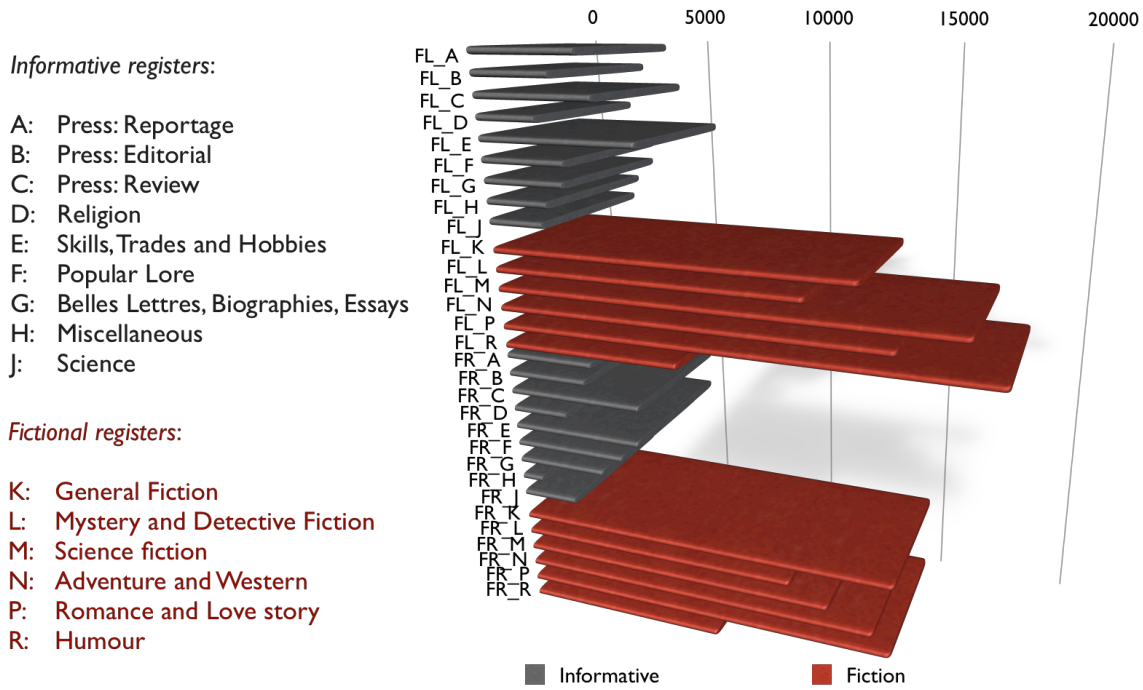
Table 2 ORCHiDS Dickens Component

No.	Texts	Category	Date	Tokens
1	<i>Sketches by Boz</i>	Sketches	1833–1836	187,493
2	<i>The Pickwick Papers</i>	Serial Fiction	1836–1837	300,269
3	Other Early Papers	Sketches	1837–1840	67,121
4	<i>Oliver Twist</i>	Serial Fiction	1837–1839	157,707
5	<i>Nicholas Nickleby</i>	Serial Fiction	1838–1839	322,393
6	<i>Master Humphrey's Clock</i>	Miscellany	1840–1841	47,331
7	<i>The Old Curiosity Shop</i>	Serial Fiction	1840–1841	217,521
8	<i>Barnaby Rudge</i>	Serial Fiction	1841	254,141
9	<i>American Notes</i>	Sketches	1842	101,941
10	<i>Martin Chuzzlewit</i>	Serial Fiction	1843–1844	336,630
11	<i>Christmas Books</i>	Fiction	1843–1848	154,449
12	<i>Pictures from Italy</i>	Sketches	1846	72,553
13	<i>Dombey and Son</i>	Serial Fiction	1846–1848	342,386
14	<i>David Copperfield</i>	Serial Fiction	1849–50	356,110
15	<i>Bleak House</i>	Serial Fiction	1852–1853	354,369
16	<i>Hard Times</i>	Serial Fiction	1854	103,443
17	<i>Little Dorrit</i>	Serial Fiction	1855–1857	338,502
18	<i>Reprinted Pieces</i>	Sketches	1850–1856	91,574
19	<i>A Tale of Two Cities</i>	Serial Fiction	1859	136,259
20	<i>The Uncommercial Traveller</i>	Sketches	1860–1869	142,973
21	<i>The Great Expectations</i>	Serial Fiction	1860–1861	184,889
22	<i>Our Mutual Friend</i>	Serial Fiction	1864–1865	326,158
23	<i>The Mystery of Edwin Drood</i>	Serial Fiction	1870	94,143
Sum of word-tokens in the set of Dickens texts:				4,690,355

Table 3 ORCHiDS 18th Century Subcorpus

No.	Author	Texts	Date	Tokens
1	Defoe	<i>Captain Singleton</i>	1720	110,843
2	Defoe	<i>Journal of Prague Year</i>	1722	94,695
3	Defoe	<i>The Military Memoirs of Captain George Carleton</i>	1728	80,612
4	Defoe	<i>Moll Flanders</i>	1724	136,241
5	Defoe	<i>Robinson Crusoe</i>	1719	232,324
6	Fielding	<i>A Journey from this World to the Next</i>	1749	45,003
7	Fielding	<i>Amelia</i>	1751	211,678
8	Fielding	<i>Jonathan Wild</i>	1743	69,938
9	Fielding	<i>Joseph Andrews</i>	1742	125,342
10	Fielding	<i>Tom Jones</i>	1749	346,256
11	Goldsmith	<i>The Vicar of Wakefield</i>	1766	62,976
12	Richardson	<i>Clarissa</i>	1748	935,894
13	Richardson	<i>Pamela</i>	1740	438,937
14	Smollet	<i>Peregrine Pickle</i>	1751	342,200
15	Smollett	<i>Ferdinand Count Fathom</i>	1753	157,641
16	Smollett	<i>Humphrey Clinker</i>	1771	150,395
17	Smollett	<i>Sir Launcelot Greaves</i>	1760	89,388
18	Smollett	<i>Roderick Random</i>	1748	191,602
19	Smollett	<i>Travels through France and Italy</i>	1766	121,074
20	Sterne	<i>A Sentimental Journey</i>	1768	40,783
21	Sterne	<i>Tristram Shandy</i>	1759–1767	186,426
22	Swift	<i>A Tale of a Tub</i>	1704	44,121
23	Swift	<i>Gulliver's Travels</i>	1726	104,047
24	Swift	<i>A Journal to Stella</i>	1710–1713	189,587
Sum of word-tokens in the set of 18th Century texts:				4,508,003

\*1 The frequency matrix of the 121 body-part words across the 15 registers in the FLOB/FROWN corpora is omitted due to space constraints. The data is available upon request to the author.



Body-part words in two corpora: FLOW and FROWN

Fig. 2 Frequencies of body-part words (per 100,000 words) across 15 registers in the FROWN and FLOW corpora

Table 4 ORCHiDS 19th Century Subcorpus

No.	Author	Texts	Date	Tokens
1	Austen	<i>Emma</i>	1815	160,307
2	Austen	<i>Mansfield Park</i>	1814	159,792
3	Austen	<i>Northanger Abbey</i>	1803	77,250
4	Austen	<i>Persuasion</i>	1816	83,303
5	Austen	<i>Pride and Prejudice</i>	1813	122,253
6	Austen	<i>Sense and Sensibility</i>	1811	119,502
7	A.Brontë	<i>Agnes Grey</i>	1847	68,222
8	C.Brontë	<i>Jane Eyre</i>	1847	186,272
9	C.Brontë	<i>The Professor</i>	1857	88,081
10	C.Brontë	<i>Villette</i>	1853	194,479
11	E.Brontë	<i>Wuthering Heights</i>	1847	116,431
12	Collins	<i>After Dark</i>	1882	136,884
13	Collins	<i>The Moonstone</i>	1868	194,136
14	Collins	<i>The Woman in White</i>	1859	246,799
15	G.Eliot	<i>Adam Bede</i>	1859	215,101
16	G.Eliot	<i>Brother Jacob</i>	1864	16,683
17	G.Eliot	<i>Daniel Deronda</i>	1876	311,115
18	G.Eliot	<i>Middlemarch</i>	1871–1872	318,065
19	G.Eliot	<i>Silas Marner</i>	1861	71,389
20	G.Eliot	<i>The Mill on the Floss</i>	1860	207,341
21	Gaskell	<i>Cranford</i>	1851–1853	70,947
22	Gaskell	<i>Mary Barton</i>	1848	160,589
23	Gaskell	<i>Sylvia's Lovers</i>	1863	190,829
24	Thackeray	<i>Barry Lyndon</i>	1844	126,048
25	Thackeray	<i>Vanity Fair</i>	1848	303,935
26	Trollope	<i>Barchester Towers</i>	1857	198,369
27	Trollope	<i>Can You Forgive Her</i>	1865	316,151
28	Trollope	<i>Doctor Thorne</i>	1857	217,748
29	Trollope	<i>The Eustace Diamonds</i>	1873	269,818
30	Trollope	<i>Phineas Finn</i>	1869	263,231
31	Trollope	<i>The Warden</i>	1855	72,032

Sum of word-tokens in the set of 19th Century texts: 5,283,102

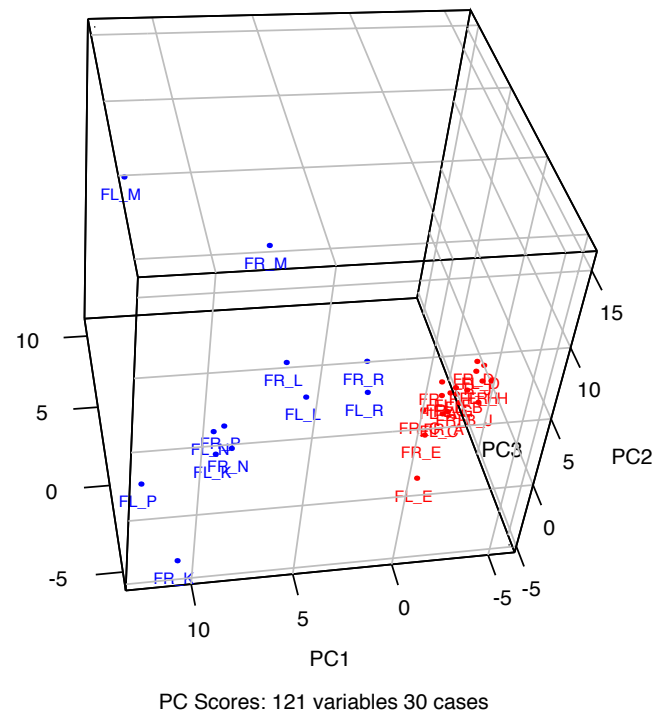


Fig. 3 Principal Component Analysis of 121 body-part words across 15 registers in the FROWN and FLOW corpora

### 3. Body-part words in classic British fiction

Normalized frequencies of the body-part words (per 100,000 running words) were counted in order to examine complex inter-relationships between the terms, the texts, and association between the terms and the text. Word frequency profiles, of which

part is shown as Table 5 were then fed to a series of multivariate analysis, including PCA and Correspondence Analysis, and machine-learning-based classification algorithm (Random Forests). The three subcorpora in the ORCHiDS are differentiated according to differing frequency patterns of body language.

To supplement findings from stylometric analysis, topic modelling was carried out using MALLET (a JAVA-based Machine Learning Language Toolkit developed at the University Massachusetts at Amherst)\*2. What is of special interest is that a topic made up of body-language-related words was found to discriminate strongly in favour of the Dickens set (Fig. 4). Fig. 5 shows topic density distribution across the authors in the corpus. One Dickens sample is shown to have nearly as much as 40% accounted for by words belonging to the topic labelled facial/bodily gestures.

Full results will be provided at the 119th SIG-CH conference.

References

- [1] Paquot, M. and Bestgen, Y. (2009). Distinctive words in academic writing: A comparison of three statistical tests for keyword extraction, in Jucker, A., Schreier D., and Hundt M. (eds.) *Corpora: Pragmatics and Discourse. Papers from the 29th International Conference on English Language Research on Computerized Corpora (ICAME 29)*, pp. 247–269.
[2] Klausner, C., Nerbonne, J., and Coltekin, C. (2015). Finding Characteristic Features in Stylometric Analysis, *Digital Scholarship in the Humanities*, Vol. 30, Supplement 1, ii14–ii29.
[3] Rayson, P. and Garside, R. (2000). Comparing corpora using frequency profiling. In proceedings of the *Workshop on Comparing Corpora, held in conjunction with the 38th annual meeting of the Association for Computational Linguistics (ACL 2000)*, 1–8 October 2000, Hong Kong, 1–6.
[4] Mahlberg, M., Stockwell, P., de Jooe, J., Smith, C., and O'Donnell, M. B. (2016). CLIC Dickens: Novel uses of concordances for the integration of corpus stylistics and cognitive poetics. *Corpora*, 11 (3), 433–463.
[5] Mahlberg, M. and Wiegand V. (2018). Corpus stylistics, norms and comparisons: Studying speech in Great Expectations, in Page, R.,

Table 6 Words most strongly distinguishing between Dickens, 18th century fiction (ECF), and 19th century fiction (NCF)

Table with 8 columns: Item, Gini Index, Accuracy, Dickens, ECF, NCF, Key Group. Rows list body parts like head, belly, legs, breast, heels, cheeks, face, hair, back, chin, arm, hand, foot, heart, throat, eye, fingers, hands, heads, shoulder, foot, forehead, body, nose, lips, elbows, faces, knees, nerves, leg, fingers, knee, tongue, thumb, eyebrows, bone, tooth, spleen, neck, eyes, arms, mouth, arm's.

Table 5 Part of the frequency profiles for body-part words

Large table with 19 columns: texts, test\_group, back, hand, head, face, eyes, hands, body, hair, heart, feet, blood, eye, arms, mouth, skin, arm, fingers, shoulder, foot, neck, faces, legs, nose. Rows list various Dickens and Austen texts.

\*2 <http://mallet.cs.umass.edu/index.php>

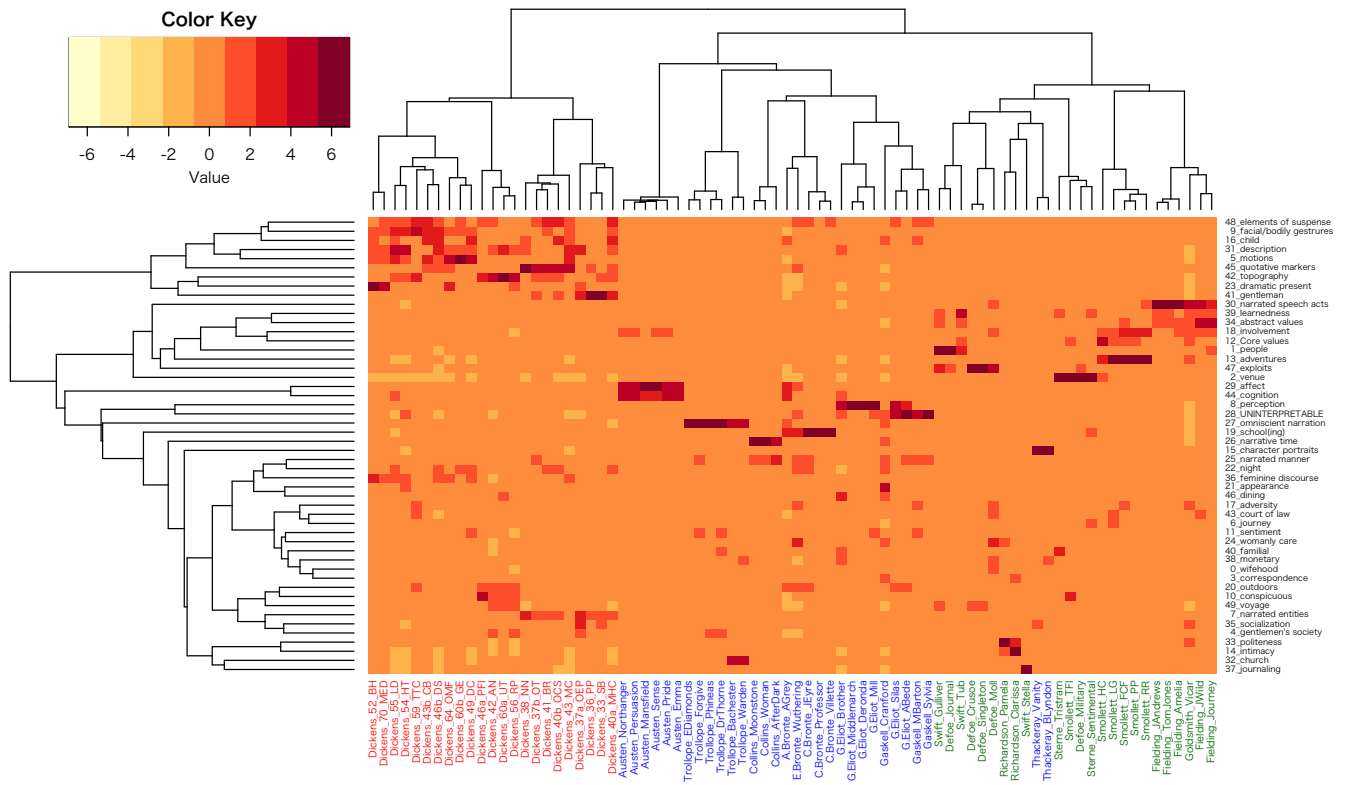


Fig. 4 Heatmap visualization of 50 topics across 78 texts (with mean weights scaled)

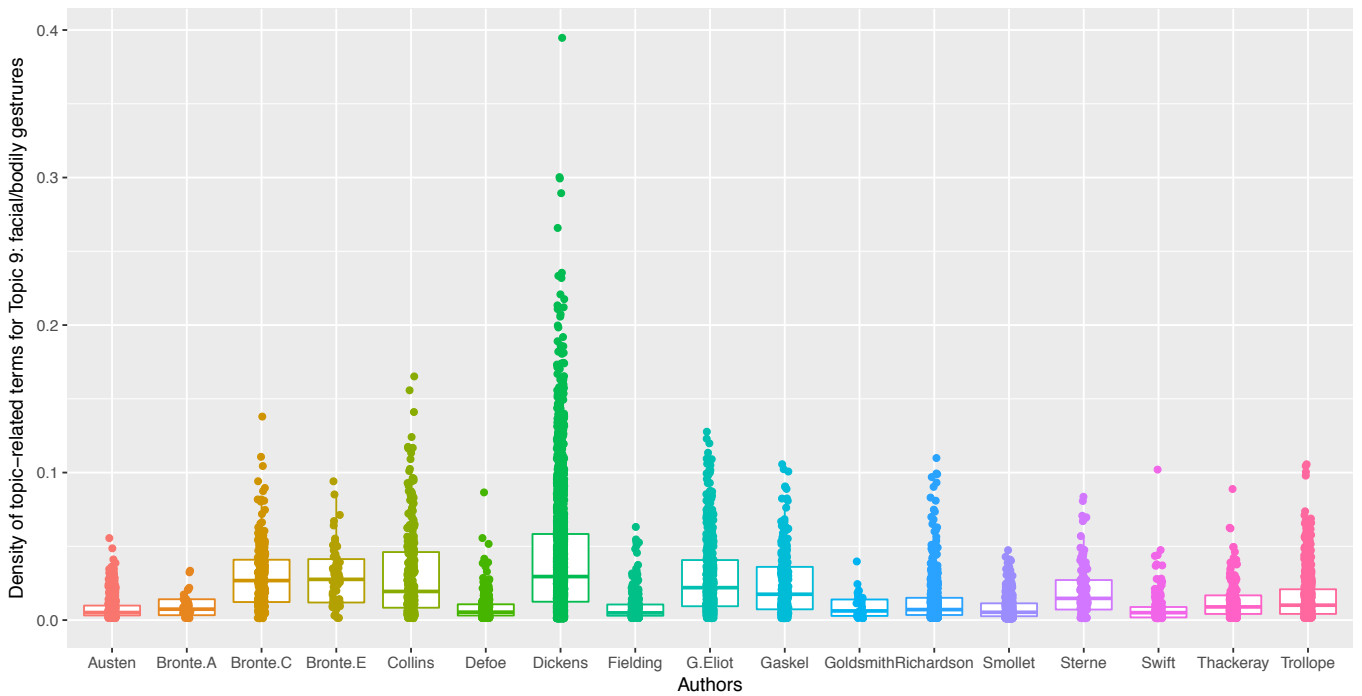


Fig. 5 Topic density distribution across the authors: Topic 9 facial/body gestures)