

The following publication Tay, D. (2021). Automated lexical and time series modelling for critical discourse research: A case study of Hong Kong protest editorials. *Lingua*, 255, 103056 is available at <https://dx.doi.org/10.1016/j.lingua.2021.103056>. The Cognitive Linguistic Studies is available at <https://benjamins.com/catalog/cogls>

1 Automated lexical and time series modeling for critical discourse research: a case study of Hong
2 Kong protest editorials

3

4 **Abstract**

5 This paper advances a novel approach to critical synchronic and diachronic discourse analysis
6 using automated lexical and time series modeling. It is illustrated by a case study of near-daily
7 editorials (N=201; 300,081 words) from 9 June to 2 October 2019 on the Hong Kong protest
8 movement in three ideologically contrasting sources – *China Daily* (CD), *South China Morning*
9 *Post* (SCMP), and *Hong Kong Free Press* (HKFP). Lexical analysis with Linguistic Inquiry and
10 Word Count (LIWC) first revealed four predominant socio-psychological word categories -
11 relativity, drive, cognitive, and affect. Overall, HKFP expresses anger at the government, CD
12 lays blame on protestors' violent actions, and SCMP occupies a middle position to focus on less
13 political aspects. Time series modeling is then applied to redirect attention from these aggregated
14 differences to how they unfold day-to-day. It was found that while positive affect words are
15 characterized by short-term consistencies and fluctuations, most variables exhibit random
16 variation across time. The approach allows precise description of how linguistic variables in
17 neighboring time periods inter-relate, offering rich interpretative possibilities for different
18 linguistic/discourse contexts. Furthermore, determining whether a variable is 'modelable' offers
19 a systematic and replicable way to interrogate the assumption that discourse inevitably serves to
20 construe social reality.

21

22 **Keywords:** LIWC, time series analysis, discourse analysis, Hong Kong protests

23

24

25

26

27

28

29

30

31

32

33

34 **Introduction**

35 Critical discourse research aims to uncover the reciprocal relationships between language and
36 power; i.e. how language forms and structures both reflect and sustain power relations (Wodak &
37 Meyer, 2009). It has largely relied on close qualitative interpretation of texts with reference their
38 underpinning social, political, and historical backdrops. In media contexts, however, there is
39 growing interest in automated and quantitative approaches that can support synchronic as well as
40 diachronic analyses of language, due to the inherent relevance of time as a variable in news
41 (Gabrielatos & Baker, 2008; Prentice, 2010). This paper introduces, following recent related
42 work (Smirnova, Laranetto, & Kolenda, 2017; Author, 2019), the combination of automated
43 lexical analysis and time series modeling as a novel approach that addresses gaps in existing
44 approaches to deepen our critical understanding of time-based discourse. The lexical analysis
45 tool is Linguistic Inquiry and Word Count (LIWC) (Tausczik & Pennebaker, 2010), which
46 quantifies language use in socio-psychological categories that have particular relevance in
47 contexts like the media. The modeling of time series data; i.e. LIWC category scores across time,
48 applies the well-known Box-Jenkins time series methodology (Box et al., 2015). It is mostly
49 used to analyze changes in a variable across time in order to control or forecast future values in
50 fields like finance and engineering. Its applicability to discourse studies has only recently been
51 explored (Koplenig, 2017; Author, 2017, 2019), the basic premise being that the temporal
52 behavior of linguistic variables may resemble finance and engineering variables.

53 The approach will be illustrated by a case study of editorials on the 2019 Hong Kong protest
54 movement in three ideologically contrastive English language news sources – *China Daily* (CD),
55 *South China Morning Post* (SCMP), and *Hong Kong Free Press* (HKFP). In this way, a
56 secondary objective is to show how the protests and related events were represented within a
57 critical time span across the political spectrum. Political tensions in Hong Kong date back to its
58 1997 handover from the UK to China but have recently escalated over an extradition bill to
59 legalize fugitive transfer to other jurisdictions including the Chinese mainland. Increasingly
60 confrontational protests have since implicated other issues like democracy, alleged
61 police/protester violence, and the future of Hong Kong as a Chinese city. The time span of
62 interest is from 9 June 2019, the date of the largest post-1997 protest linked to the bill, to 2
63 October 2019, one day after the Chinese National Day by which many speculated that drastic
64 counteractions would be taken.

65 In the following sections, I first discuss the background and relevant studies in the Hong Kong
66 context including the three news sources to be compared. I then explain how the present
67 approach can address issues with existing synchronic and diachronic analyses, leading to a
68 deeper understanding of media representations of the protest. The methodology, results, and key
69 theoretical implications for critical discourse research will then be presented in turn.

70

71 **The Hong Kong context**

72 The sovereignty of Hong Kong was transferred from the UK to China in 1997, ending 156 years
73 of colonial rule and commencing its status as a special administrative region with distinct
74 governing systems from the Chinese Mainland. Three of its most popular English language news
75 sources occupy contrasting positions on the political spectrum. The pro-establishment *China*
76 *Daily* (Hong Kong edition) was launched in 1997 by the Chinese government to present
77 governmental perspectives for English readers. The non-profit online *Hong Kong Free Press*
78 (HKFP), founded by independent journalists to counteract a perceived decline in press freedom,
79 lies at the other end. *South China Morning Post* (SCMP) is the oldest English news source in
80 Hong Kong. Its editorial stance is less clear and lies somewhere in between. After its 2016
81 acquisition by Chinese conglomerate Alibaba Group, it has been seen as veering towards a pro-
82 Beijing stance while still allowing discussion of independence, self-determination, and localism
83 (Wiebrecht, 2018). These ideological differences are expected to be reflected in differing
84 linguistic constructions of protest-related editorials.

85 Post-1997 Sino-Hong Kong relations have attracted much scholarly attention across the social
86 sciences (Cheng, 2016; Mathews et al., 2007; Ortmann, 2020). Focusing on media and language,
87 Lee and Lin (2006) compared discursive strategies used in editorials to construct the avowed
88 stances of the pro-establishment *Ming Pao* and pro-democracy *Apple Daily*. They found the
89 former employed a “rhetoric of objectivity and rationality”, while the latter positioned itself as
90 the “defender of public opinion and local interests” (Lee & Lin, 2006:353). The ‘umbrella
91 movement’ of 2014, most remembered for a 79-day occupation of the city centre by protestors
92 demanding electoral reform, inspired many studies with a range of critical approaches. These
93 include Bhatia’s (2015) analysis of rhetorical tools used by SCMP (e.g. insinuation, temporal
94 referencing, metaphor, recontextualization, reframing) to construct the movement, Mey and
95 Ladegaard’s (2015) pragmatic analysis of the ‘discourse of democracy’ in debates related to the
96 movement, Lee’s (2016) interrogation of selective reporting of opinion polls by different
97 Chinese-language news sources, and Flowerdew’s (2017) critical discourse historiographical
98 analysis of the movement. The 2019 protests are likely to inspire similar research. The above
99 studies share the common element of preferring nuanced qualitative interpretation of a limited
100 amount of data. While this approach yields valuable insights, decisions related to data sampling
101 and the choice of what linguistic/discursive features to focus on are not always systematically
102 explained (Breeze, 2011). Given the avowed interest in how language reflects different political
103 persuasions, it is also surprising that few studies have explicitly compared news sources along
104 clearly defined linguistic variables. Moreover, since events like the umbrella movement and the
105 current protests can quickly evolve over short time intervals, there is much room to explore how
106 time series analyses over short spans can shed light on language changes and the attendant
107 conceptualizations in turn.

108

109 **Applying LIWC and time series modeling to critical discourse analysis**

110 The above Hong Kong-based studies are part of a rich tradition of critical media language
111 research. Much of this work shows how political events and relationships are construed via
112 discursive choices that maintain divisions along gendered, ethnic, national lines, and so on.
113 Critics often point out limited sample sizes, biased selection of features, and unsystematic
114 analyses as its main methodological flaws. The synergy between corpus and critical methods is a
115 response to these criticisms, as seen from the many corpus-assisted critical discourse studies in
116 contexts ranging from politics to business (Baker et al., 2008; Koller, 2006; Partington, 2010).
117 One of the most useful corpus-assisted analytic strategies in this regard is automated tagging
118 with systems like the UCREL Semantic Analysis System (Archer et al., 2002). They provide
119 bottom-up descriptions of grammatical and semantic categories that can be statistically analyzed
120 and serve as entry points to qualitative scrutiny (e.g. Stefanowitsch & Gries, 2006). In this sense,
121 they address all three criticisms above.

122 LIWC is one such system with added advantages for studies that focus on socio-psychological
123 constructs like affect, power, and ideology (Smirnova et al., 2017). It has undergone rigorous
124 psychometric evaluation using speaker intuitions and actual usage patterns (Pennebaker et al.,
125 2015). In other words, its ‘bag of words’ validly and reliably reflects the underlying grammatical
126 and semantic categories. LIWC quantifies texts under two main types of variables: ‘summary
127 variables’ and ‘linguistic dimensions’. The four summary variables, each scored from 0 to 100,
128 are *analytical thinking*, *clout*, *authenticity*, and *emotional tone*. They respectively compute the
129 extent of logical vs. narrative thinking (Pennebaker et al., 2014), expertise vs. tentativeness
130 (Kacewicz et al., 2013), personal vs. distanced (Newman et al., 2003), and positive vs. negative
131 emotions (Cohn et al., 2004) in a text. Linguistic dimensions, the focus of this study, are
132 normalized frequency measures of all words in a text under approximately 90 different
133 categories. These include grammatical (e.g. pronouns, articles, prepositions, conjunctions, parts
134 of speech, quantifiers) and socio-psychological semantic domains such as affective, cognitive,
135 perceptual, psychological, and so on, which can reveal how the protests and related events are
136 represented.

137 Furthering the above synchronic analysis, the second part of this study examines diachronic
138 change, a growing topic of interest in critical discourse research. Corpus-assisted studies often
139 approach diachronic change by segmenting datasets into sub-corpora each representing a
140 particular time interval (Bamford et al., 2013; Gabrielatos & Baker, 2008), and then comparing
141 them for statistically significant differences. More qualitative takes like the discourse-historical
142 approach conceptualize change in more abstract ways; for example, by reconstructing the
143 “historical interrelationships” of “thematically or/and functionally connected discourse fragments
144 or utterances” (Reisigl, 2017:53). There are several limitations to these approaches (cf. Author,
145 2019). Firstly, links between linguistic/discursive units and temporal units are not always clearly
146 articulated. Secondly, not much is said about changes over shorter time intervals that may inhere
147 in contexts that reflect dynamic daily realities like the present editorials. Thirdly, studies tend to

148 overlook the critical feature of interdependence or autocorrelation in time series discourse data.
149 Common statistical techniques like log-likelihood tests assume independence between the
150 frequencies of two aggregated time periods, but in shorter time spans we may see patterns where
151 (near)-consecutive sessions influence one another in ways that can be modeled and critically
152 interpreted. To address these limitations, this study will apply the Box-Jenkins approach (Box et
153 al., 2015) to time series analysis, which uses a family of mathematical models known as ARIMA
154 (Autoregressive Integrated Moving Average) models. Essentially, an ARIMA model expresses
155 an observed quantity at the present time (an LIWC variable score on a particular day) in terms of
156 quantities and/or differences between observed and predicted quantities at previous times.
157 Linguistic variables as ARIMA models can therefore be interpreted in terms of both their
158 synchronic content and diachronic structure, yielding deeper insights into their use across the
159 time span. As we will see, a key theoretical implication is that determining whether linguistic
160 variables are indeed patterned and hence ‘modelable’ offers a novel way to interrogate the
161 assumption that discourse inevitably serves to construe social reality.

162

163

164 **Data and methods**

165 All protest-related editorials in CD, HKFP, and SCMP published from 9 June to 2 October 2019
166 were collected. An editorial was considered protest-related if it explicitly discusses the protests
167 and/or related issues like the extradition bill, democracy, and so on. Reliable and exhaustive
168 identification was relatively straightforward given the exigence of the protests. Articles
169 published on the same day in each source were compiled into a single text file with nuisance
170 words like author names removed. The total dataset has 201 articles and 300,081 words (CD: 92
171 articles/181,571 words, HKFP: 52 articles/85,165 words, SCMP: 57 articles/33,345 words). The
172 different sizes are not inherently problematic due to the normalized nature of LIWC variables
173 and the independent time series modeling of each series.

174 LIWC then computed linguistic variables in each file. This generated i) a comprehensive profile
175 of variables for comparative analysis within and across news sources, as well as ii) each variable
176 as a time series from 9 June to 2 October. Time series modeling was implemented by *XLStat* and
177 all other statistical analyses and visualizations with *jamovi* (2019).

178

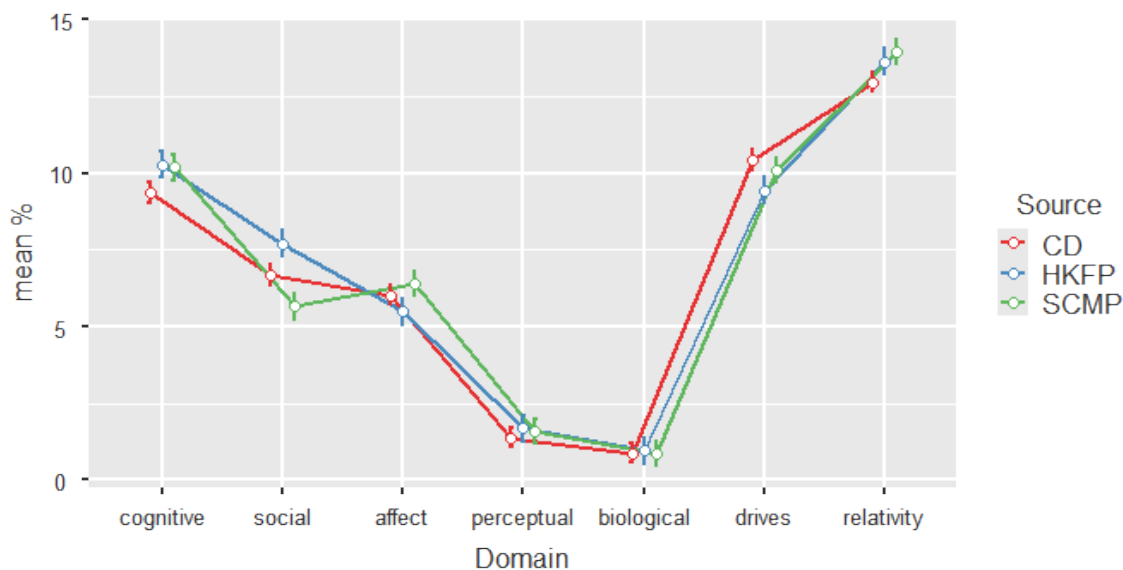
179 **Comparative analysis**

180 Comparison of major semantic domains

181 To determine the analytic focus in accordance with a data-driven approach, the first step is to
182 compare the distribution of major semantic domains (cognitive, social, affective, perceptual,
183 biological, drives, relativity) within and across sources (Figure 1). These domains tap into key

184 socio-psychological constructs (Pennebaker et al., 2015) that can inform how the protests are
 185 variously conceptualized. The y-axis in all following figures show the mean proportion of words
 186 across texts in the categories on the x-axis. Error bars are 95% confidence intervals. All
 187 comparative analyses are done with factorial ANOVAs (with categories as within-subject
 188 variables and news source as the between-subjects variable). Significance level is set at $p=.05$.

189



190

191 **Figure 1.** Semantic domains within and across sources

192

193 The domains are ranked as follows: i) relativity, ii) drives, iii) cognitive, iv) social, v) affect, vi)
 194 perceptual, vii) biological. Tukey post-hoc tests show all inter-domain differences are significant
 195 ($p<.001$) except drives vs. cognitive ($p=1.00$). Table 1, adopted from Pennebaker et al.
 196 (Pennebaker et al., 2015), show the further sub-divisions in each domain and some example
 197 words in the LIWC dictionary. Note that a word can be classified under multiple domains.

198

Domain	Sub-domains and example words
Relativity	Motion (arrive, car, go) Space (down, in) Time (end, until, season)
Drives	Affiliation (ally, friend, social) Achievement (win, success, better) Power (superior, bully) Reward (take, prize, benefit) Risk (danger, doubt)

Cognitive	Insight (think, know) Causation (because, effect) Discrepancy (should, would) Tentative (maybe, perhaps) Certainty (always, never) Differentiation (hasn't, but, else)
Social	Family (daughter, dad, aunt) Friends (buddy, neighbour) Female/male references (girl, her, mom/boy, his, dad)
Affect	Positive emotion (love, nice, sweet) Negative emotion (hurt, ugly, nasty, hate, kill, sad)
Perceptual	See (view, saw, seen) Hear (listen, hearing) Feel (feels, touch)
Biological	Body (cheek, hands, spit) Health (clinic, flu, pill) Sexual (horny, love, incest) Ingestion (dish, eat, pizza)

199 **Table 1.** Semantic domains and example words

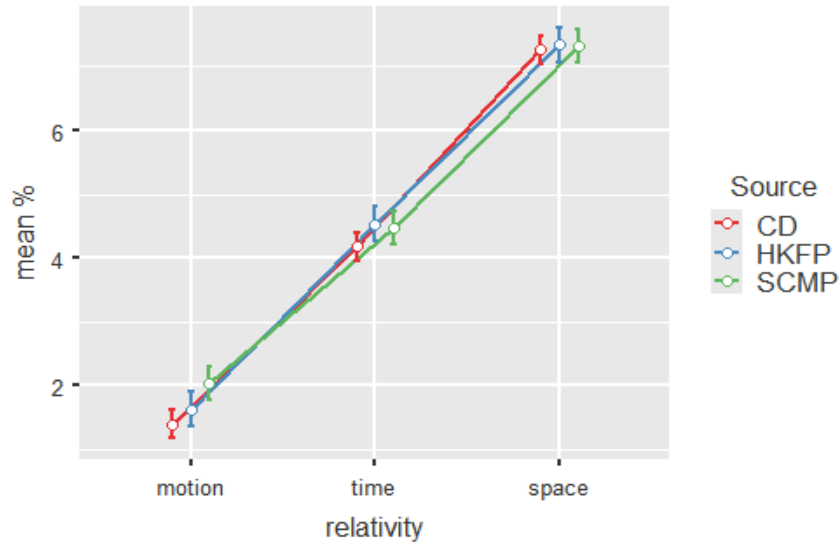
200

201 The top domains of relativity, drives, cognitive, and affect will be analyzed in turn below. Each
 202 constitutes at least five percent of all texts on average. The social domain is excluded because its
 203 sub-categories are relatively less comparable; ‘female/male references’ are not likely to contrast
 204 with ‘family’ and ‘friends’ in a discursively interesting way. Figure 1 also suggests that the
 205 domains have roughly similar distributions across sources, which reflects a generally similar
 206 semantic foci across the three news sources. This is expected as all three belong to the same
 207 genre of newspapers. However, interesting differences between sources become apparent when
 208 we zoom in on the sub-categories, as detailed below.

209

210 Relativity words

211 Relativity words depict movement, space, and time, and are perhaps unsurprisingly the most
 212 common category. Figure 2 shows the distribution of its sub-categories within and across
 213 sources.



214
 215 **Figure 2.** Relativity words within and across sources

216 Space words are most frequent, followed by time and motion words ($p < .001$) with no interaction
 217 effect ($p = .071$); i.e. the same order generally applies to all three sources. Relativity words might
 218 not seem interesting from a critical point of view when used in their basic spatial-temporal sense.
 219 Taking a random example from each source,

- 220
- 221 1. First, the government lost the power of discourse soon after the legislative
 222 process began, as the opposition hoodwinked hundreds of thousands of people
 223 into opposing the bill and joining the mass protest march on June 9 (4 July, CD)
 - 224
 - 225 2. On Sunday August 11, over 300 people turned out, wearing red or black, to
 226 spell out ‘FREE HONG KONG’ with their bodies in Taipei’s Central Art Park
 227 (24 August, HKFP)
 - 228
 - 229 3. Coach’s statement on Monday spoke of the firm respecting and supporting
 230 China’s sovereignty and territorial integrity (14 August, SCMP)

231

232 The underlined relativity words in the three examples describe basic details about protest-related
 233 events, times, and places. However, cognitive linguists and critical discourse analysts have long
 234 noted that such concepts tend to be metaphorically used to conceptualize abstract concepts in
 235 subtle and systematic ways (Hart, 2011; Lakoff & Johnson, 1999). Motion and space in
 236 particular are well-known ‘embodied source domains’ that jointly facilitate reasoning about
 237 event structure, like in the conceptualization of political processes as ‘journeys’ (Chilton, 2004).
 238 In this regard, while space words are equally used by all three sources, SCMP has significantly

239 more motion words than CD ($p=.006$). Examples 4-7 illustrate some qualitative differences in
240 metaphorical uses of motion words from the two sources.

241

242 4. For a society to thrive and prosper, people need to respect and obey the law.
243 For Hong Kong to move beyond the impasse and turn over a new leaf, the rule of
244 law must be upheld. (29 Sept, CD)

245 5. Economically, the worst may be yet to come. Hong Kong is at risk of slipping
246 into recession, as the city's Financial Secretary Paul Chan Mo-po warned on
247 Monday. (6 Aug, CD)

248 6. The way forward is anything but clear, and deep reflection is needed to bring
249 this dark chapter to an end (17 Sept, SCMP)

250 7. They face an uncertain future in a mature city now defined by a lack of upward
251 mobility, a income gap and soaring housing prices which, together, put the dream
252 of ever owning a home beyond the growing reach of many (4 July, SCMP)

253

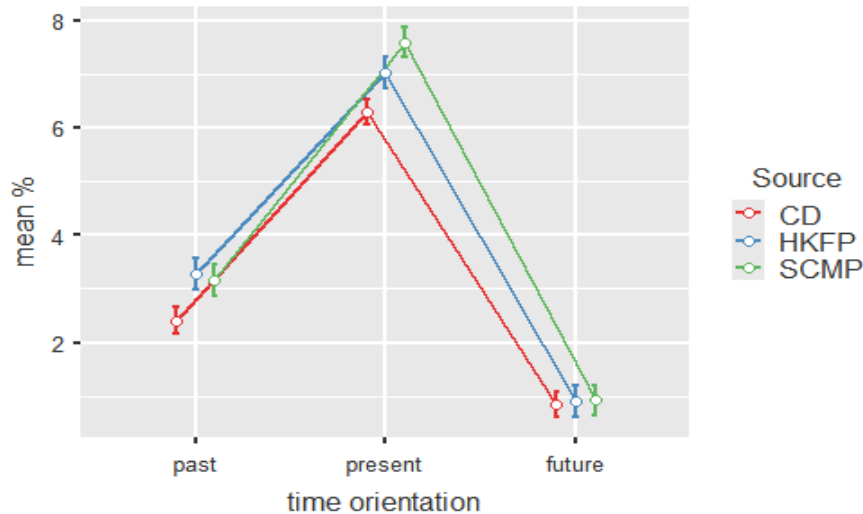
254 In Examples 4 and 5 (CD), Hong Kong is conceptualized as a metaphorical mover that needs to
255 overcome obstacles (*impasse*, *recession*) implied to be self-inflicted (*turn over a new leaf*). There
256 is also a sense of passivity as the direction of movement is downwards (*slipping*), and worse
257 situations are *coming* towards Hong Kong. The SCMP examples on the other hand describe the
258 need for a more proactive *forward* movement to *bring* things to an end, and focus on the lack of
259 (upward) mobility of its people in a situation that is not implied to be self-inflicted.

260 A limitation of LIWC shared by most other automated analyses is that it does not distinguish
261 between literal and metaphorical uses. This may complicate the analysis of such embodied words
262 because of their stronger tendency to be metaphorical. Determining i) the proportion of literal
263 versus metaphorical uses, ii) the systematicity of the latter, and iii) their ostensible ideological
264 functions would require a fuller manual analysis beyond the present scope.

265 Different than motion and space, time is less likely to be a metaphorical source domain because
266 it is less experientially concrete (Lakoff & Johnson, 1999). Time words can instead be profitably
267 analyzed with LIWC for their 'time orientation'. A past/present/future orientation is defined by
268 past/present/future tense verbs and words that refer to past/present/future events. Previous work
269 has examined the relationship between verb tenses and ideological messages in different genres
270 like social media and political speeches. Djemili et al (Djemili et al., 2014) propose that
271 ideological messages are 'timeless' and 'polychronous', giving the illusion of relevance at all
272 times and "grouping all the temporal perspectives and cancelling them". This predicts more
273 present than past and future tenses in ideological language. Fetzer and Bull (2012) discuss
274 instead politicians' strategic use of past tenses to recontextualize controversial present issues and

275 make them seem more acceptable. Figure 3 shows the distribution of time orientation within and
276 across sources.

277



278

279 **Figure 3.** Time orientation within and across sources

280

281 The present orientation is most frequent followed by past and future ($p < .001$). The interaction
282 effect is significant ($p < .01$); CD has weaker past orientation than SCMP ($p = .003$) and HKFP
283 ($p < .001$), as well as weaker present orientation than SCMP ($p < .001$) and HKFP ($p = .005$). There
284 were no significant differences in future orientation across all three sources, and no significant
285 differences in past and present orientation between SCMP and HKFP. Despite this, Examples 8-
286 10 illustrate nuanced differences in time orientation display. They show the opening paragraph of
287 all three editorials following the formal withdrawal of the extradition bill on 4 September. This is
288 a hallmark event in the time span that motivates editorials to reflect on the background, present
289 situation, and way ahead (past=bold, present = underlined, future=italicized).

290

291 8. Here we are again with the usual suspects some well meaning, others devious
292 and some just brain cell-challenged floundering around, suggesting ways in which
293 the current political tension could be eased. What they don't or won't understand
294 is the simple fact that the Chinese Communist Party is not looking for ways of
295 lowering the temperature but is bent on defeating and humiliating the protesters.
296 (4 Sept, HKFP)

297

298

299 The narrative approach in Example 8 reflects HKFP's independent press status. Its
300 predominantly present orientation via tense and word choice (e.g. *here we are again*) emphasizes
301 the immediacy and perpetualness of the present situation, and deems the Chinese Communist
302 Party culpable. In line with Djemili et al's (2014) observations, past and future oriented language
303 is consequently highly limited.

304

305 9. Hong Kong is *expected* to be free of any more violence, including rioting, now
306 that the special administrative region's chief executive has withdrawn what **was** a
307 well-intended bill to amend the extradition law. The proposed amendment **was**
308 **meant** to plug the legal loopholes which have for years allowed criminal suspects
309 from other parts of the world, including those from the Chinese Mainland, to
310 evade justice by seeking shelter in the Hong Kong SAR. Such loopholes in Hong
311 Kong's law have **turned** the city into a haven for fugitives, and **prevented** the
312 SAR government from extraditing criminal suspects to other jurisdictions. (5
313 Sept, CD)

314

315

316 In contrast, Example 9 (CD) has a more balanced time orientation. Present and future oriented
317 language is used not to convey a sense of perpetual despair but to assert a new state of
318 freedom from violence following the bill withdrawal. Past oriented language is used to
319 portray the bill as a well-intentioned but failed attempt to rectify legal loopholes.

320

321

322 10. At long last, the extradition bill that has embroiled the city in its worst turmoil
323 since reunification will be formally withdrawn. Belated as it is, the decision is
324 badly needed to take the heat out of an escalating crisis and, *hopefully*, *will* pave
325 the way to restoring order and stability. While the change of heart by Chief
326 Executive Carrie Lam Cheng Yuet Ngor *will* not please everyone, protesters
327 would be wise to show similar goodwill if compromise and reconciliation are to
328 be reached. Announcing the withdrawal in a prerecorded television address
329 **yesterday**, the embattled leader **said** her government **was** also ready to take
330 further steps to break the deadlock. (4 Sept, SCMP)

331

332

333

334 Example 10 (SCMP) also demonstrates a balanced time orientation. It differs from the other two
335 in that its present and past oriented language seem less connotative, focusing more on concrete
336 temporal details than on evaluating the state of affairs or attributing blame to either party. The

337 future-oriented language is also less assertive than CD, hoping for rather than declaring a
338 violence-free state.

339

340 Words referring to drives

341 The domain of drives is based on McClelland's (1987) Theory of Needs which claims that
342 humans are motivated by affiliation, achievement, and power. LIWC categorizes words that refer
343 to these and two additional subcategories of reward and risk. The relative prevalence of these
344 words in the editorials reflects the motivational underpinnings of the protests, with specific
345 differences shedding light on their potential links to editorial stances. Figure 4 shows the
346 distribution within and across sources.

347



348

349 **Figure 4.** Drives within and across sources

350

351 Power words are by far the most frequent followed by affiliation, achievement, risk, and reward.
352 All inter-category differences are significant ($p < .001$) except for affiliation vs. achievement
353 ($p = .997$). While the distribution of sub-categories appears to be similar between sources, the
354 significant interaction effect ($p < .001$) suggests interesting contrasts: i) CD has the most power
355 words ($p < .001$) with no difference between SCMP and HKFP ($p = .099$); ii) HKFP has the most
356 affiliation words, significantly more than SCMP ($p = .038$) but not CD ($p = .999$); iii) SCMP has
357 the most risk words, significantly more than HKFP ($p = .036$) but not CD ($p = .072$). Each of these
358 observations will be illustrated with reference to the three sources' editorial reflections on

359 another hallmark event on 1 July, the Hong Kong Special Administrative Region Establishment
360 Day. Protestors stormed, vandalized, and damaged the legislative council building.

361 Firstly, while power words are prevalent across all sources, this is especially the case for CD.
362 Example 11 shows a concentration of mostly nominal labels referring to the establishment and its
363 institutions (e.g. *government, police officers, authority, rule of law*). These institutions and
364 'Hong Kong society' are cast in abstract impersonal terms (van Leeuwen, 1995) and construed as
365 expressing strong condemnation of protestor actions.

366

367 11. The special administrative region government and Hong Kong society strongly
368 condemned the outrageous violence in Wan Chai and Admiralty on Monday,
369 during which protesters charged police lines protecting the flag-raising ceremony
370 for the 22nd anniversary of the HKSAR in the morning, and then stormed the
371 Legislative Council Complex in the afternoon and evening. More than a dozen
372 police officers were injured. These actions were brazen challenges to government
373 authority and the city's rule of law ... If the government were to give in to their
374 demands, which it should not, the anti-amendment forces would come up with
375 more demands, rendering the government, the legislature and even the police
376 force unable to operate. (1 July, CD)

377

378 Contrastively, HKFP reflecting on the same events displays a high degree of affiliation (Example
379 12). The emphasis is not on institutional power but on construing a sense of solidarity between
380 HKFP and the protestors, who are more concretely represented with frequent inclusive first-
381 person pronouns.

382

383 12. Victories are rare these days, and the temptation to rest on our laurels is great.
384 After all, we have fought tooth and nail to stop a dangerous bill and this time we
385 have tangible results to show for in contrast to the sense of empty-handedness in
386 2014. We deserve a pat on the back and a celebratory drink for a job well done. On
387 the other hand, can we afford to let our guard down, however briefly? Or must we
388 strike the iron while it's hot and keep up the pressure on our opponents until other
389 political demands are met? (1 July, HKFP)

390

391 SCMP, on the other hand, uses the most words referencing risks; i.e. dangers, concerns, and
392 things to avoid. Example 13 differs from the above in that it discusses an emergent mental health
393 crisis in a way that does not appear to take either side of the political divide.

394

395 13. Sadly, three deaths with suicide notes or other references to the crisis since the
396 mass protests began last month may not be the last linked to the controversy over
397 the now-suspended extradition bill. Two subsequent suicide attempts went viral
398 online, with an alarming increase in calls to counselling hotlines further raising
399 fears of copycat behavior (4 July, SCMP)

400

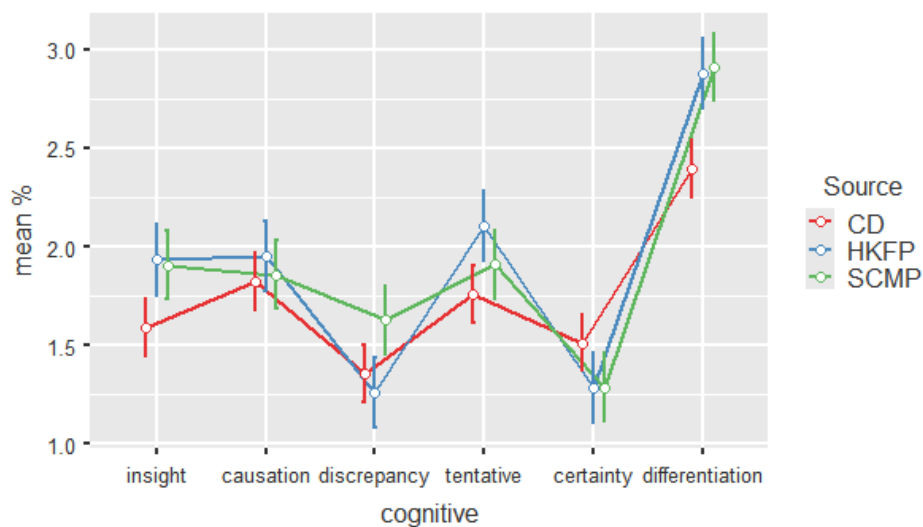
401 Therefore, while power and affiliation are intuitive and observable salient opposites in editorial
402 constructions of the protests, we also see relative preferences among news sources in terms of
403 more ‘neutral’ categories like the evaluation of objective risks. In this regard, the sources affirm
404 what is known about their editorial stances – CD communicates a higher degree of abstract
405 institutional power, HKFP a concrete sense of solidarity among protestors, and SCMP a more
406 neutral perspective.

407

408 Cognitive words

409 The six categories of cognitive words refer to various cognitive processes that can reveal how the
410 editorials reason about the protests (Figure 5).

411



412

413 **Figure 5.** Cognitive words within and across sources

414

415 Similar to the previous analysis of drives, one category (differentiation) stands out as
416 significantly more frequent than the rest ($p < .001$). The next three categories of tentative, insight,

417 and causation are not significantly different from one another, but each significantly more
418 frequent than the last two (discrepancy and certainty). The interaction effect is significant
419 ($p < .001$) but inter-source differences exist only for differentiation; CD uses the least
420 differentiating language ($p < .001$) compared to HKFP and SCMP, which are alike ($p = 1.00$).

421 Differentiating language includes words like *hasn't*, *but*, and *else*, which in many contexts signal
422 disagreement or subjective distinction of entities. Examples 14-16 are editorial reflections of 21
423 June when protestors besieged the Hong Kong police headquarters after a peaceful sit-in. This
424 was arguably the first major escalation since the start of the protest time span and is therefore
425 likely to prompt rational analysis of the situation. They illustrate variation between sources in
426 terms of their discursive construction of differences.

427

428 14. Although the right to peaceful protest is an integral part of a free society,
429 people who abuse demonstrations by indulging in wanton violence that endangers
430 the safety of others must expect to face justice ... No matter what the excuse is,
431 political or otherwise, protestors do not have a “license to break the law” in Hong
432 Kong or any other society under the rule of law (21 and 23 June, CD)

433

434

435 In Example 14, CD acknowledges the right to protest but proceeds to categorically demarcate
436 protestors from ‘others’, constructing the protests as a conflict between violent and non-
437 violent individuals. This strategy of categorical differentiation is also applied to discuss the
438 grounds (*political or otherwise*) as well as context (*Hong Kong or any other society*) of the
439 protests, which further implies violent protestors as having uniquely political interests to
440 cause harm to Hong Kong.

441

442 15. If businesspeople think that they cannot live and work in Hong Kong without
443 the risk of a few months in gaol fighting extradition, followed by an appearance
444 on Confessiontube and a few years in a mainland prison, then they will live and
445 work in Singapore, and take their money and their business down there. People
446 have a variety of views about the merits of living in Singapore, but we can all
447 agree that it’s an improvement on a mainland prison. (22 June, HKFP)

448

449 Example 15 from HKFP likewise uses differentiating language to erect discursive boundaries
450 in a context of disagreement. The differentiation is however focused on living conditions in

451 the two oft-compared cities of Hong Kong and Singapore, suggesting that the differences
452 have been accentuated by Chinese influence on the former.

453

454 16. As reflected in the two mass protests early this month, the key to a successful
455 protest is to stay peaceful and lawful. Confrontations that go against this principle
456 may not necessarily win public support. These testing times make political
457 wisdom and accountability all the more important. Regrettably, officials have
458 shied away from facing the public, possibly out of fear that any wrong step would
459 antagonise the situation further. But governance and public trust are at stake. (22
460 June, SCMP)

461

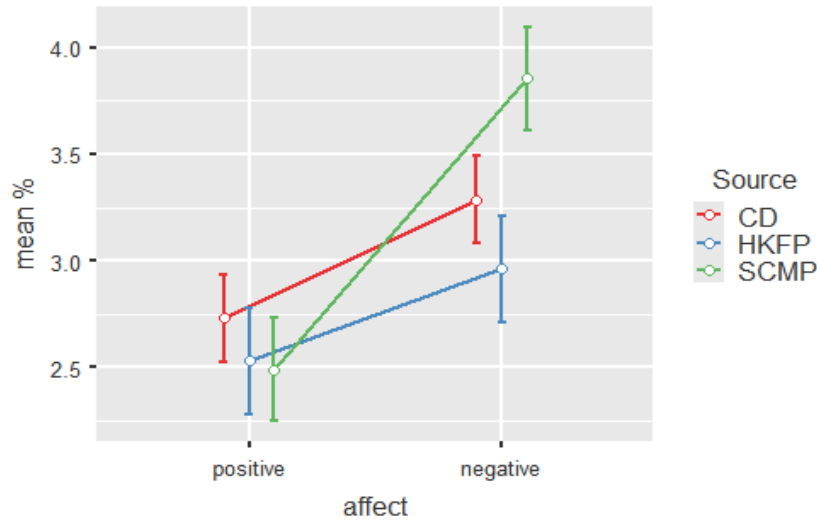
462 From this perspective, Example 16 from SCMP comes across as using differentiating
463 language in a relatively impersonal and therefore neutral manner. The boundaries are not
464 between concrete people and places, but abstract principles of peace versus violence, as well
465 as the misalignment between the situation and ‘governance and public trust’. Similar to the
466 previous analysis of words referring to drives, we see that SCMP appears to adopt a relatively
467 neutral stance that focuses on issues rather than individuals. It should be remembered,
468 however, that in purely quantitative terms it is CD that utilizes the least amount of
469 differentiating words.

470

471 Affect words

472 Affect words describe emotions and are subdivided in LIWC into positive and negative emotion
473 words. The linguistic representation of emotions is a fundamental topic in critical discourse
474 studies (Mackenzie & Alba-Juez, 2019; Smirnova et al., 2017). It is also a key part of sentiment
475 analysis in contemporary data analytics where emotions are inferred from word meanings and
476 used to evaluate the affective properties of texts (Taboada et al., 2011). Figure 6 shows the
477 distribution of positive and negative words within and across sources.

478



479
 480 **Figure 6.** Affect words within and across sources

481
 482 Unsurprisingly, negative words are significantly more frequent than positive words ($p < .001$).
 483 The interaction effect is significant ($p < .001$) with no difference in positive words between
 484 sources. SCMP, however, uses significantly more negative words than CD and HKFP ($p < .001$),
 485 which differs from the previous categories where it mostly occupied the middle position. To
 486 further probe the different sub-types of negative emotion words, we examine its specific
 487 subcategories of anger, anxiety, and sadness (Figure 7).



489
 490 **Figure 7.** Specific negative emotions within and across sources

492 Anger words were significantly most frequent followed by anxiety and sadness words ($p<.001$),
493 with no interaction effect ($p=.14$) across the sources. Protest editorials from all three sources
494 expressed anger as the predominant negative emotion, although this tended to describe or be
495 directed at different parties. Examples 17 to 19 are illustrative, drawing again from reflections of
496 1 July when protestors stormed, vandalized, and damaged the legislative council building.

497

498 17. Journalists covering recent protests in Hong Kong have had a difficult and, at
499 times, dangerous duty to perform. The role they play in ensuring the public
500 receives accurate information about the demonstrations is of great importance.
501 But some frontline members of the media have been subjected to appalling
502 attacks, abuse, and harassment while working. Such conduct is to be condemned
503 (1 Jul, SCMP)

504

505 18. The armed wing, moreover, was maniacal, utterly committed to finishing off
506 what they unsuccessfully attempted on June 12, when the police thwarted their
507 earlier attack on the legislature. Perhaps the most pathetic sight of all, however,
508 was provided by the “pan-democrat” legislators, whose role became farcical.
509 Having earlier fanned the protests by their irrational claims about the extradition
510 bill, they then tried to control the fanatics, only to be contemptuously pushed out
511 of the way, with one of them, the hapless Leung Yiu-chung, even ending up
512 sprawled on the ground (2 Jul, CD)

513

514 19. The perceived failure of a large scale uprising created an opening for those in
515 power to go on the offensive. In the years since, the authorities had made
516 damaging incursions into our freedoms and way of life, by disqualifying
517 opposition lawmakers, banning unwanted candidates from the ballot, outlawing a
518 localist party, expelling a defiant foreign journalist, and ceding a piece of our
519 territory in the heart of the city to the mainland authorities. Hong Kong people
520 were kicked while they're down, over and again. We had grown so used to these
521 political assaults that they barely registered. Beijing had all but written us off as
522 docile subjects who were finally beaten into submission (1 Jul, HKFP)

523

524 All three sources used anger words related to physical actions (e.g. *protests*, *abuse*, *harassment*,
525 *attacks*, *assaults*) as well as attitudes (*maniacal*, *contemptuously*, *offensive*, *beaten*), to convey a
526 generally negative sentiment towards the protests. However, SCMP (Example 17) depicts
527 ostensibly neutral media workers as the victims of general anger, while CD (Example 18) and

528 HKFP (Example 19) direct this sentiment towards the protestors and establishment respectively.
529 Similar to the previous analyses of drive and cognitive words, these observations likewise affirm
530 the known editorial stances of the three sources. The second most frequent sub-category of
531 anxiety words, however, has SCMP deviating again from the middle position. SCMP expresses
532 significantly more anxiety than CD ($p<.001$), with no differences between SCMP and HKFP
533 ($p=.06$) / CD and HKFP ($p=.42$). To illustrate this, Examples 20-23 are editorial reflections on
534 the events of 23 August, the 30th anniversary of the Baltic Way as protestors held hands to form
535 'human chains' across the city. This somewhat innovative gesture prompted discussion on its
536 symbolism and implications as editorials take stock of the protest movement.

537

538 20. In hindsight the city bounced back strongly from SARS, but it will not
539 necessarily do so again when the political turbulence over the now-shelved
540 extradition bill finally passes. A health threat that has been contained and
541 eradicated is not to be compared with political risk and uncertainty in the minds of
542 investors, businesspeople and talent that Hong Kong must attract if it is to be
543 competitive in a globalised economy. (23 Aug, SCMP)

544

545 21. Those attending rallies throughout Taiwan aren't just standing in solidarity
546 with those fighting for freedom across the strait, they're also collecting helmets,
547 gas masks and other protective items for those on the streets. Likewise, the
548 Taiwanese government can offer protection, and opportunity, to Hongkongers
549 whose safety and freedom are threatened due to political threats. (23 Aug, HKFP)

550

551 22. Over the years, "human chain" rallies have become a signature tactic of
552 independence movements around the world. Hong Kong people must guard
553 against being misled into the slippery slope. Separatism has no future in Hong
554 Kong. It would only bring disaster and self-destruction (23 Aug, CD)

555

556 Example 20 (SCMP) and 21 (HKFP) both contains common anxiety words like *threat*, *threaten*,
557 and *uncertainty*. SCMP uses these words with reference to the non-political SARS crisis in 2003,
558 implying that the present situation is even more uncertain and harder to eradicate. However, it
559 focuses on business issues and is not as explicit as HKFP in attributing this threat to Hong
560 Kong's perceived common political enemy with Taiwan; i.e. the mainland Chinese government.
561 Notwithstanding this difference, both sources contrast with Example 22 (CD) where no anxiety
562 words are used. There is instead a more authoritative tone that criticizes the 'human chains' and
563 calls for vigilance, short of considering it a threat worthy of anxiety. Therefore, although SCMP

564 uses significantly more anxiety words, the ways in which they reflect known editorial stances
565 appears to resemble the previously discussed categories.

566

567 **Time series analysis**

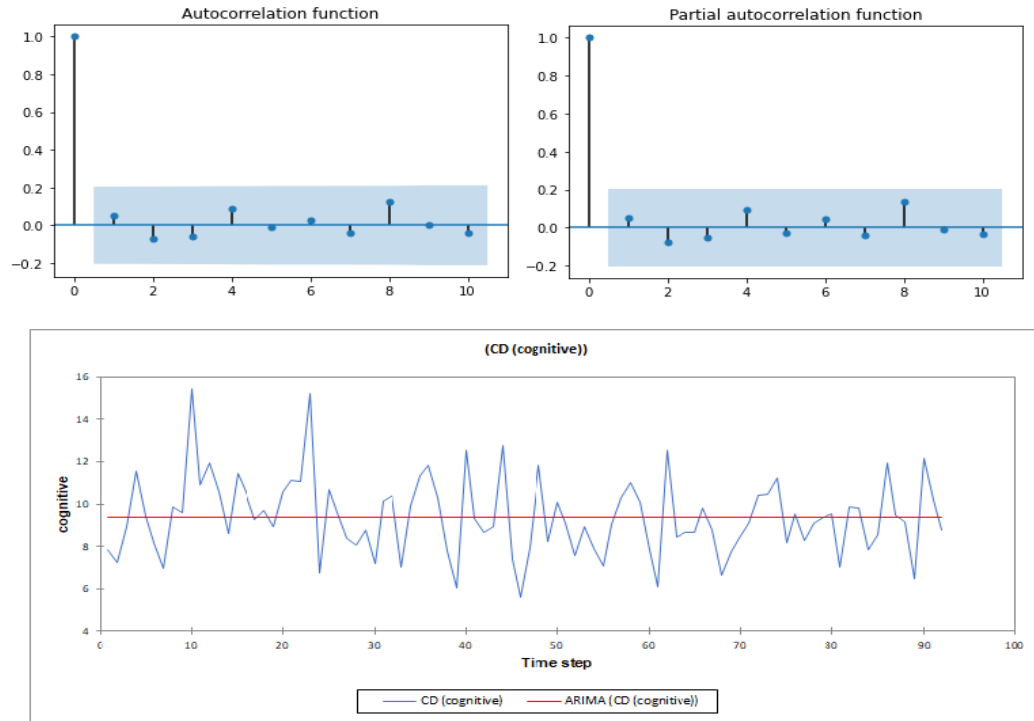
568 The above comparative analyses were based on the average values of each LIWC variable across
569 the time span, giving a broad overview of differences within and between sources. From this
570 perspective, ideological differences are reflected and perhaps indeed constructed by an
571 aggregated observation of linguistic/discursive differences. However, an underexplored question
572 in critical discourse research is how language use over natural time intervals, especially short-
573 term intervals like daily news, might be patterned in ways that reveal insights into the
574 relationship between time and the linguistic/discursive construction of reality. This perspective
575 instead regards each linguistic/discursive variable as a time series, documenting day-to-day
576 changes in its LIWC score. The Box-Jenkins method of time series analysis applies the following
577 steps for each variable (Author, 2019): i) inspect and transform the series if necessary to meet
578 statistical requirements, ii) calculate autocorrelations within the series to determine if
579 consecutive observations are linked, iii) identify and fit candidate ARIMA models based on the
580 autocorrelations, iv) perform diagnostic tests for goodness-of-fit, v) accept the present model if
581 fit is adequate or find a better one. The technicalities of each step will not be elaborated here.
582 This six-step process lead to one of two general conclusions: i) the series is randomly distributed
583 across time such that past values have no bearing on future values, or ii) the series is patterned
584 across time such that future values can be expressed to varying degrees of accuracy by some
585 ARIMA model; i.e. as a function of past values.

586 All variables described in the above comparative analysis underwent time series modeling. Each
587 day of publication is a time step and editorials published on the same day by each source were
588 collectively analyzed. Most of the variables were found to be randomly distributed across the
589 time span, also described as ‘white noise’ in statistical terminology. Figure 8 illustrates a random
590 example of a white noise time series (cognitive words in CD).

591

592

593



594

595 **Figure 8.** Example of white noise time series (CD cognitive)

596

597 The top of Figure 8 shows the (partial) autocorrelation functions for CD cognitive. The
 598 horizontal axis shows the number of ‘lags’, or daily steps apart, and the vertical axis shows the
 599 correlation coefficient for each lag. The coefficients are statistically insignificant across all lags,
 600 as visually indicated by their confinement within the 95% confidence interval region. This
 601 implies that there is zero autocorrelation in the series and hence no ARIMA models are suitable.
 602 The horizontal red line across the plotted series at the bottom of Figure 8 shows that in such
 603 cases, the best estimate at each time step is simply the mean value of the series.

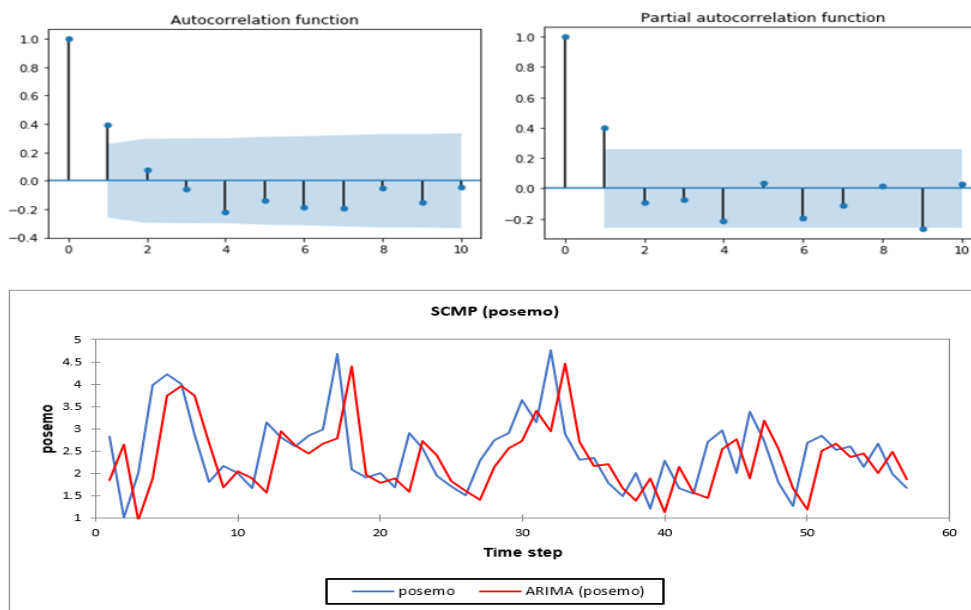
604 Two affective variables, however, were found to fit ARIMA models. These are positive emotion
 605 words (posemo) from SCMP as well as HKFP. Figure 9 shows that for SCMP posemo the lag 1
 606 correlation is about +0.4, which means that posemo usage over consecutive days is positively
 607 correlated - a day-to-day increase/decrease tends to be followed by another increase/decrease.
 608 This is an instance of an AR(1) model (first-order autoregressive), formally expressed as $y_t =$
 609 $2.4713 + 0.3925(y_{t-1}) + a_t$ where

- 610 • y_t is the value of posemo at day t ,
- 611 • y_{t-1} is the value of posemo at the previous day ($t-1$)
- 612 • a_t is the error term inherent in any statistical model (i.e. the actual value minus the
 613 predicted value at time t)

614

615 In more technical detail, the model informs us that there is a stable average of 2.4713% of
616 posemo words each day, and that a unit increase/decrease in posemo words in the previous day
617 (y_{t-1}) leads to a 0.3925 unit increase/decrease in the present day (y_t). The bottom of Figure 9
618 shows how the values generated by this model (red line) compares with the actual values (blue
619 line), visually suggesting an adequate fit¹. The model can also be used to forecast tomorrow's
620 posemo use by substituting today's values. This is a basic objective in many analytic contexts
621 like finance, but its discourse analytic relevance is less obvious in most cases.

622



623

624 **Figure 9.** Observed and predicted time series for SCMP positive emotion words

625

626 The contextual interpretation of this AR(1) model is a short-term consistency or momentum in
627 SCMP's use of positive emotion words. Values tend to increase or decrease consecutively but
628 the time span across which this consistency most strongly holds is one day. We see that in the
629 first half of the series, there are long stretches of consecutive rises and falls, but in the second
630 half there are quick changes of direction after one-day spans of consistency. It appears that
631 SCMP's portrayal of the situation is more likely to trend at the beginning with either increasing
632 or decreasing positivity, but becomes more erratic with quickly reversing sentiments as National
633 Day (Oct 1) approaches.

¹ There are more precise measures of accuracy such as the MAPE (Mean Average Percentage Error) of predicted vs. actual values, but they will not be elaborated here.

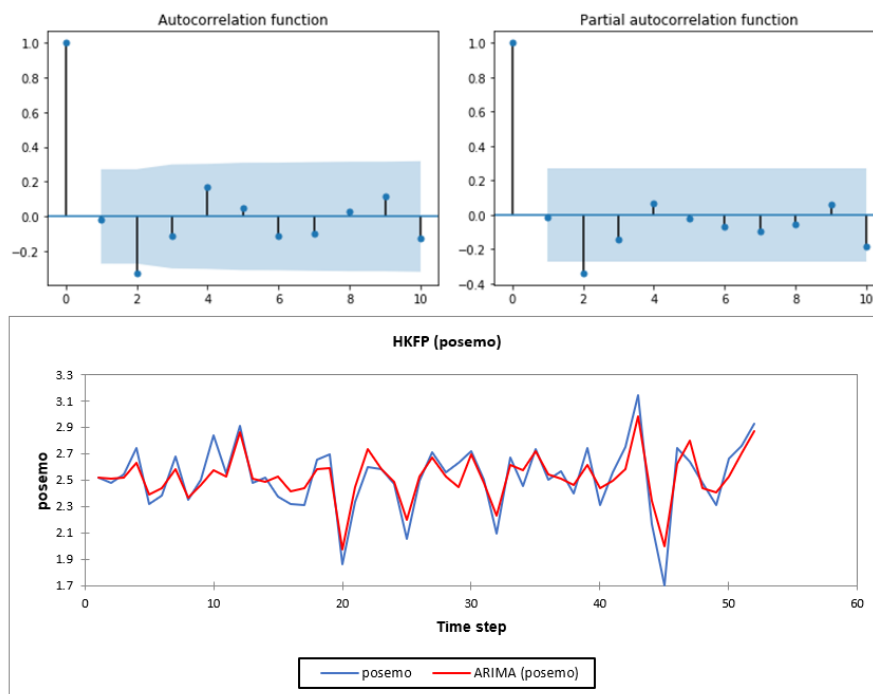
634 HKFP presents an interesting contrast as shown in Figure 10. It fits an AR(2) model (second-
635 order autoregressive) formally expressed as $y_t = 2.52 - 0.34(y_{t-2}) + a_t$ where

- 636 • y_t is the value of posemo at day t ,
- 637 • y_{t-2} is the value of posemo two days prior ($t-2$)
- 638 • a_t is the error term inherent in any statistical model (i.e. the actual value minus the
639 predicted value at time t)

640

641 This implies that associations are strongest two days apart instead of one. Additionally, the
642 correlation here is negative; increases tend to be followed by decreases two days later, and vice-
643 versa. In more technical detail, there is a stable average of 2.52% of posemo words each day, and
644 that a unit increase/decrease in posemo words two days prior (y_{t-2}) leads to a 0.34 unit
645 decrease/increase in the present day (y_t). This results in a sawtooth-like series across the time
646 span as shown at the bottom of Figure 9, as directions tend to change across two-day spans.

647



648

649 **Figure 10.** Observed and predicted time series for HKFP positive emotion words

650

651 Different than SCMP, the AR(2) model suggests a fluctuation in HKFP's use of positive emotion
652 words. The two-day span is longer than SCMP but could still be considered short-term. Due to

653 the negative correlation there is no strong linear trend in any part of the series, but only
 654 continuously volatile changes in direction across two-day spans. HKFP’s portrayal is therefore
 655 more sentimentally volatile with short bursts of relative optimism followed by pessimism.

656

657 **Summary**

658 Table 2 summarizes the above comparative and time series analyses with an overall profile of the
 659 three sources.

660

Comparative analysis	
Relativity	<ul style="list-style-type: none"> • Spatial words are both literal and metaphorical. The latter conceptualizes Hong Kong along a metaphorical journey • Present orientation is most frequent • HKFP emphasizes immediacy and perpetualness of current situation in narrative style • CD more assertive of an improved situation • SCMP more hopeful and focuses more on objective temporal details
Drive	<ul style="list-style-type: none"> • CD has the most power words that condemn protestor actions • HKFP has the most affiliation words that depict protestor solidarity • SCMP has the most risk words that depict non-political risks
Cognitive	<ul style="list-style-type: none"> • CD uses the least differentiation words but tends to differentiate violent protestors from others • HKFP tends to differentiate Hong Kong and its values from other places • SCMP tends to differentiate abstract principles
Affect	<ul style="list-style-type: none"> • Negative emotion words are more frequent than positive ones • SCMP uses the most negative emotion words, especially anxiety • Anger is the most frequent negative emotion but expressed at different things: protestors (CD), the government (HKFP), and more general entities (SCMP) • SCMP is anxious about the protests in general, HKFP is anxious about political threats from the government, CD does not often express anxiety
Time series analysis	

Positive emotion words	<ul style="list-style-type: none"> • Most variables are randomly distributed across time • SCMP displays short-term consistency but becomes erratic towards National Day • HKFP displays fluctuation with short bursts of high and low emotional tone
------------------------	--

661 **Table 2.** Summary of findings

662

663 Analyses of the four socio-psychological domains revealed broad commonalities and differences
664 in linguistic representations of the protest movement by the three sources. All three tended to
665 depict Hong Kong as a metaphorical mover along a journey, emphasize present orientation, and
666 use power and negative emotion words to depict the protests as an ongoing crisis that leads to
667 undesirable outcomes. However, a closer look at extracts shows that these representations also
668 generally align with the known editorial stances of the sources. HKFP tends to adopt a narrative
669 style to present the situation as urgent and anxiety-inducing, directing anger at the government
670 and blaming it for undermining the values of Hong Kong society which protestors in solidarity
671 are attempting to defend. CD, on the other hand, has a generally assertive tone and attributes the
672 situation to the violent actions of protestors who are harming public interests. SCMP occupies a
673 middle position for the most part, tending to focus on aspects of the protest that are less
674 obviously tied to a political stance such as health, commercial interests, and abstract principles
675 like peace versus violence. It is noteworthy that this apparent neutrality comes with the highest
676 use of negative emotion words.

677 The above cross-sectional findings are enriched by considering the structural changes of
678 linguistic variables in response to unfolding events in time. ARIMA modelling showed that most
679 variables are randomly distributed across the time span despite their aggregated differences
680 between sources. Clear exceptions are found for positive emotion words. While this is also
681 random in CD, SCMP and HKFP exhibit the respective temporal patterns of short-term
682 consistency and fluctuation. The former suggests a more stable appraisal of events in the first
683 half of the timespan with more erratic behavior thereafter, while the latter is more volatile
684 throughout. These patterns highlight different subjective construals of the same underlying
685 ‘reality’ that would have been difficult to uncover by more conventional methods of discourse
686 analysis.

687

688 **Conclusion**

689 This paper combined automated lexical and time series analysis to examine the links between
690 linguistic choices and editorial stances over a critical time span in the 2019 Hong Kong protest
691 movement. Different than studies that use relatively few samples or established corpus analytic
692 approaches, it demonstrated how LIWC and time series modeling provide a coherent synchronic

693 as well as diachronic account of linguistic representation in critical contexts. One theoretical
694 implication is the potential of LIWC to complement other semantic taggers by focusing on socio-
695 psychological dimensions of language. In critical contexts such as political discourse,
696 (sub)categories like relativity, drive, cognitive and affect combine well with more content-
697 oriented analyses at lexical and other discursive levels to explore the relationships between
698 language and power. LIWC analysis is nevertheless limited in several ways. Firstly, as discussed
699 earlier, it shares with other corpus analytic approaches the inability to detect figurative language
700 like metaphor and irony. Secondly, while comparison across sources is generally unproblematic,
701 comparisons between sub-categories of a semantic domain may be so in cases where the sizes of
702 their respective dictionaries are too different. Thirdly, LIWC is limited to lexical analysis, and
703 therefore cannot easily reveal higher-level discursive strategies that underlie editorial stances.

704 The follow-up ARIMA time series analysis also bears key implications for critical discourse
705 research. The general complementarity between synchronic and diachronic perspectives is
706 already well known, as outlined earlier. However, by explicitly considering potential
707 autocorrelations in discourse data, the present method redirects attention from an aggregated
708 view of its content to its period-by-period structure. This is most important in contexts like news,
709 classroom, and therapy talk where dynamics over short time spans are often overlooked. As we
710 saw in the brief technical elaboration of the two examples above, ARIMA models allow precise
711 description of how language in neighbouring time periods inter-relate, which offers rich
712 interpretative possibilities against the backdrop of various (critical) discourse contexts. An
713 intriguing question would be cases where different linguistic/phenomena in different contexts
714 (e.g. news vs. classroom vs. therapy) share similar models, thus implying previously
715 unconsidered deep structural similarities between them.

716 Furthermore, the trait of ‘modelability’ – which distinguishes random from patterned/modelable
717 time series – offers a fresh way to interrogate the programmatic claim, upheld by most critical
718 discourse studies, that discourse inevitably constructs social reality. The general logic is
719 summarized as follows. In contexts like the media, we can always assume some uncertainty or
720 randomness in the unfolding ‘social reality’ that discourse aims to represent/construct. Time
721 series modeling of relevant data (e.g. newspapers) then leads to one of two general outcomes: the
722 variable(s) of interest are either random, or patterned in ‘modelable’ ways. Randomness would
723 suggest the absence of an attempt to discursively fashion the likewise random background
724 events, which cautions us against being too quick to conclude that discourse always presents
725 manipulated versions of reality. On the other hand, a modelable series might indeed suggest
726 potentially strategic or manipulative construction of (aspects) of this reality, especially when
727 different models underpin the same events like in the present case. There is also the interesting
728 possibility that patterned discourse does not actually construct, but merely reflect
729 correspondingly patterned background events. A hypothetical topical example is to find some
730 linguistic variable(s) sharing similar time series models with the daily incidence of COVID-19
731 they report. Either way, one should demonstrate in systematic and replicable ways if some

732 discourse series is modelable or not, for greater confidence and clarity in claims about the
733 relationship between discourse and social reality.

734 Finally, it is worth noting that time series modeling may sometimes seem to contradict
735 ‘traditional’ aggregated comparisons simply because they are based on different statistical
736 models of the data. The present case is an example where LIWC analysis using averaged scores
737 suggests different strategic constructions among sources, whereas time series modeling suggests
738 randomness in most variables. However, this in fact creates multiple interpretative perspectives
739 that productively interrogate each other. Aggregated analyses that neglect temporal ordering may
740 profile theoretical models where ideological worldviews are shaped by multiple exposure to
741 discourse in disparate contexts (cf. Hoey, 2005), while time series analysis can deconstruct this
742 multiple exposure and reconsider it from the perspective of temporal passage, upon which time-
743 based discourse and social reality are inherently grounded. The ways in which these perspectives
744 interact across various different critical discourse contexts is a promising avenue for future work.

745

746 **References**

747 Author, 2017

748 Author, 2019

749 Archer, D., Wilson, A., & Rayson, P. (2002). *Introduction to the USAS category system*.

750 Baker, P., Gabrielatos, C., Khosravini, M., Mcenery, T., & Wodak, R. (2008). A useful
751 methodological synergy? Combining critical discourse analysis and corpus linguistics to
752 examine discourses of refugees and asylum seekers in the UK press. *Discourse & Society*,
753 *19*(3), 273–306.

754 Bamford, J., Cavalieri, S., & Diani, G. (Eds.). (2013). *Variation and Change in Spoken and*
755 *Written Discourse. Perspectives from corpus linguistics*. John Benjamins.

756 Bhatia, A. (2015). Construction of discursive illusions in the ‘Umbrella Movement’.’ *Discourse*
757 *and Society*, *26*(4), 407–427.

758 Box, G. E. P., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). *Time Series Analysis:*
759 *Forecasting and Control* (5th ed.). Wiley.

760 Breeze, R. (2011). Critical discourse analysis and its critics. *Pragmatics*, *21*(4), 493–525.

761 Cheng, E. W. (2016). Street politics in a hybrid regime: The diffusion of political activism in
762 post-colonial Hong Kong. *China Quarterly*, *226*, 383–406.

763 Chilton, P. (2004). *Analysing Political Discourse: Theory and Practice*. Routledge.

764 Cohn, M. A., Mehl, M. R., & Pennebaker, J. W. (2004). Linguistic Markers of Psychological

- 765 Change Surrounding September 11, 2001. *Psychological Science*, 15(10), 687–693.
- 766 Djemili, S., Longhi, J., Marinica, C., Kotzinos, D., & Sarfati, G.-E. (2014). What does Twitter
767 have to say about ideology? *NLP 4 CMC: Natural Language Processing for Computer-*
768 *Mediated Communication/Social Media-Pre-Conference Workshop at Konvens 2014*.
- 769 Fetzer, A., & Bull, P. (2012). Doing leadership in political speech: Semantic processes and
770 pragmatic inferences. *Discourse and Society*, 23(2), 127–144.
771 <https://doi.org/10.1177/0957926511431510>
- 772 Flowerdew, J. (2017). Understanding the Hong Kong Umbrella Movement: A critical discourse
773 historiographical approach. *Discourse and Society*, 28(5), 453–472.
- 774 Gabrielatos, C., & Baker, P. (2008). Fleeing, sneaking, flooding: A corpus analysis of discursive
775 constructions of refugees and asylum seekers in the UK Press 1996-2005. *Journal of*
776 *English Linguistics*, 36(1), 5–38.
- 777 Hart, C. (2011). Force-interactive patterns in immigration discourse: A Cognitive Linguistic
778 approach to CDA. *Discourse and Society*, 22(3), 269–286.
- 779 Hoey, M. (2005). *Lexical Priming: A New Theory of Words and Language*. Taylor & Francis.
- 780 Kacewicz, E., Pennebaker, J. W., Jeon, M., Graesser, A. C., & Davis, M. (2013). Pronoun Use
781 Reflects Standings in Social Hierarchies. *Journal of Language and Social Psychology*,
782 33(2), 125–143.
- 783 Koller, V. (2006). Of critical importance: Using electronic text corpora to study metaphor in
784 business media discourse. In A. Stefanowitsch & S. T. Gries (Eds.), *Corpus-Based*
785 *Approaches to Metaphor and Metonymy* (pp. 237–266). Mouton de Gruyter.
- 786 Koplein, A. (2017). Why the quantitative analysis of diachronic corpora that does not consider
787 the temporal aspect of time-series can lead to wrong conclusions. *Digital Scholarship in the*
788 *Humanities*, 32(1), 159–168.
- 789 Lakoff, G., & Johnson, M. (1999). *Philosophy in the Flesh: The Embodied Mind and its*
790 *Challenges to Western Thought*. Basic Books.
- 791 Lee, F. L. F. F. (2016). Opinion polling and construction of public opinion in newspaper
792 discourses during the Umbrella Movement. *Journal of Language and Politics*, 15(5), 589–
793 608.
- 794 Lee, F. L. F. F., & Lin, A. M. Y. Y. (2006). Newspaper editorial discourse and the politics of
795 self-censorship in Hong Kong. *Discourse and Society*, 17(3), 331–358.
- 796 Mackenzie, J. L., & Alba-Juez, L. (Eds.). (2019). *Emotion in Discourse*. John Benjamins.
- 797 Mathews, G., Ma, E., & Lui, T. L. (2007). *Hong Kong, China: Learning to belong to a nation*.

- 798 Routledge.
- 799 McClelland, D. (1987). *Human Motivation*. Cambridge University Press.
- 800 Mey, J. L., & Ladegaard, H. J. (2015). Discourse, democracy and diplomacy: A pragmatic
801 analysis of the Occupy Central movement in Hong Kong. *Word*, 61(4), 319–334.
- 802 Newman, M. L., Pennebaker, J. W., Berry, D. S., & Richards, J. M. (2003). Lying Words:
803 Predicting Deception From Linguistic Styles. *Personality and Social Psychology Bulletin*,
804 29(1901), 665–675.
- 805 Ortmann, S. (2020). Hong Kong’s Constructive Identity and Political Participation: Resisting
806 China’s Blind Nationalism. *Asian Studies Review*, 00(00), 1–19.
- 807 Partington, A. (2010). Modern Diachronic Corpus-Assisted Discourse Studies (MD-CADS) on
808 UK newspapers: an overview of the project. *Corpora*, 5(2), 83–108.
- 809 Pennebaker, J. W., Boyd, R. L., Jordan, K., & Blackburn, K. (2015). *The development and*
810 *psychometric properties of LIWC2015*. University of Texas at Austin.
- 811 Pennebaker, J. W., Chung, C. K., Frazee, J., Lavergne, G. M., & Beaver, D. I. (2014). When
812 small words foretell academic success: The case of college admissions essays. *PloS One*, 9,
813 1–10.
- 814 Prentice, S. (2010). Using automated semantic tagging in Critical Discourse Analysis: A case
815 study on Scottish independence from a Scottish nationalist perspective. *Discourse &*
816 *Society*, 21(4), 405–437.
- 817 Reisigl, M. (2017). The Discourse-Historical Approach. In J. Flowerdew & John E. Richardson
818 (Eds.), *The Routledge Handbook of Critical Discourse Studies* (pp. 44–59). Routledge.
- 819 Smirnova, A., Laranetto, H., & Kolenda, N. (2017). Ideology through sentiment analysis: A
820 changing perspective on Russia and Islam in NYT. *Discourse and Communication*, 11(3),
821 296–313.
- 822 Stefanowitsch, A., & Gries, S. T. (Eds.). (2006). *Corpus-Based Approaches to Metaphor and*
823 *Metonymy* (Issue 171). Mouton de Gruyter.
- 824 Taboada, M., Brooke, J., Tofiloski, M., Voll, K., & Stede, M. (2011). Lexicon-based methods for
825 sentiment analysis. *Computational Linguistics*, 37(2), 267–307.
- 826 Tausczik, Y. R., & Pennebaker, J. W. (2010). The Psychological Meaning of Words : LIWC and
827 Computerized Text Analysis Methods. *Journal of Language and Social Psychology*, 29(1),
828 24–54.
- 829 the jamovi project. (2019). *jamovi (Version 1.0.1) [Computer software]*. www.jamovi.org
- 830 van Leeuwen, T. (1995). The representation of social actors. In *Texts and Practices. Readings in*

- 831 *Critical Discourse Analysis* (pp. 32–70). Routledge.
- 832 Wiebrecht, F. (2018). Cultural co-orientation revisited: The case of the South China Morning
833 Post . *Global Media and China*, 3(1), 32–50.
- 834 Wodak, R., & Meyer, M. (2009). Critical Discourse Analysis: history, agenda, theory, and
835 methodology. In R. Wodak & M. Meyer (Eds.), *Methods for Critical Discourse Analysis*
836 (2nd ed., pp. 1–33). Sage.
- 837