

# Building a methodological framework and toolkit for news media dataset tracking of conflict and cooperation dynamics on transboundary rivers

Liying Guo<sup>1</sup>, Jing Wei<sup>1</sup>, Keer Zhang<sup>1</sup>, Jiale Wang<sup>1</sup>, Fuqiang Tian<sup>1\*</sup>

<sup>1</sup>Department of Hydraulic Engineering, State Key Laboratory of Hydroscience and Engineering, Tsinghua University, Beijing, 100084, China

*Correspondence to:* Fuqiang Tian (tianfq@tsinghua.edu.cn)

**Abstract.** Management of transboundary rivers will be one of the great political and environmental challenges of the 21st century if knowledge of conflict and cooperation is not fully developed. Transboundary river conflict and cooperation are critical for the sustainable development of river basins, regional security, and stability, and have significant scientific and practical implications. The construction of a dataset of transboundary water events – individual conflictive or cooperative interaction between riparian – provides important data support and factual basis for the study of transboundary rivers. However, the most representative research, the Transboundary Freshwater Dispute Database, is built by means of manual reading for information extraction, thus difficult for fast updating, also does not cover the global changes in the past decade. This research aims to build a methodological framework for news media datasets tracking of conflict and cooperation dynamics on transboundary rivers, provide mass of relevant data for the research of transboundary rivers in the globe, prepare a potent research toolkit, lay a solid foundation for further data mining research, and better suit the big data age. In order to test the effectiveness of the methodological framework and toolkit for dataset construction, this research analyses the spatial coverage, temporal coverage, content coverage and conducts relevance screening of the articles in datasets. The results show that the datasets built by this framework can capture comprehensive contents of transboundary water conflict and cooperation in both spatial and temporal coverage with acceptable data quality.

## 1 Introduction

Globally, there are 310 transboundary river basins, covering 47.1% of the land area except Antarctica (McCracken & Wolf, 2019), and accounting for approximately 60% of global freshwater discharge (Wolf et al., 1999). The population of the basins comprises 52% of the world's total (McCracken & Wolf, 2019). Transboundary river basins not only support the lives of the people in the basins, but also connect the various economic sectors and ecosystems in the basin into an organic whole; transboundary water management not only affects the development of riparian countries in all aspects, but also intertwines social, economic, environmental, and political sectors of each riparian country and increase interdependence in between (United Nations, 2019). Riparian countries have divergent demands and priorities for transboundary water resources, different

30 development agendas for water resources, and different water governance regimes and water resources cultures (Sadoff &  
31 Grey, 2005), which make the management of transboundary water resources more complex than that of domestic water  
32 resources. Transboundary river basins are thus prone to conflicts of various forms, forming a complex situation where conflicts  
33 and cooperation develop intertwined. Therefore, research on water conflict and cooperation in transboundary rivers has  
34 important theoretical value and practical significance. Exploration of dynamics of conflict and cooperation as social sectors in  
35 a human-water coupled transboundary system is especially prominent.

36 Among the extant studies on transboundary rivers, transboundary water event datasets – individual conflictive or cooperative  
37 interactions between riparian – provide factual data support for the formation of global generalized understanding, which is of  
38 great significance. The most representative research - Transboundary Freshwater Dispute Database (TFDD) developed by  
39 Oregon State University (Wolf, 1999) has compiled more than 6,400 historical transboundary water events, both conflictive  
40 and cooperative (3813 left after removing duplicated records by us from their original data) on the global scale from the year  
41 of 1948 to 2008 (Transboundary Freshwater Dispute Database, 2008). The data came from existing political science datasets  
42 and news media articles, which was manually screened, interpreted, and coded to extract the detailed information of the water  
43 event (Yoffe & Larson, 2001). Building upon these event data, Basin at Risk Projects (BAR) (Yoffe & Larson, 2001) further  
44 classified water events by level of intensity of conflict or cooperation, ranging from -7 to +7 to identify potential socio-political  
45 threats, and provided a brief summary of the detailed information of the event. The results included very few examples of full  
46 cooperation and extreme conflicts but identified river basins that are at potential risk for further conflict. TFDD has built up  
47 foundation of this methodological framework for tracking transboundary river water events and allows for further identification  
48 of the conflict/cooperation dynamics and possible analysis of its complex driving mechanism.

49 Given that manual reading and coding processes was adopted in TFDD, which largely limit the implication of this method in  
50 the era of big data. The explosion of digital news data, whose discussion of transboundary water events has grown  
51 exponentially, made it more difficult to manually track all published water events and the dynamics of conflict and cooperation.  
52 While manual reading excels in extracting latent and detailed content, it is much more time and labor consuming. Therefore,  
53 it is necessary to revise the methodological framework to meet the current need for a more comprehensive and detailed dataset  
54 which can be updated in a more efficient manner. Meanwhile it can also provide the basis for further analysis, i.e. to reflect  
55 the concerns of different stakeholders, obtain a global law of transboundary water conflict and cooperation (Bernauer &  
56 Böhmelt, 2020).

57 This paper aims to provide such a revised methodological framework for news media tracking of conflict and cooperation  
58 dynamics on transboundary rivers and provide a toolkit when applying the framework in the corresponding research. The  
59 theory that inspired our framework is from Lasswell's model of communication (Lasswell, 1948), who focused on  
60 communication as a process, to conduct problem-oriented inquiry of the news report through content analysis with the seven  
61 fundamental elements "who, with what intentions, in what situations, with what assets, using what strategies, reaches what  
62 audiences, with what result?". Our design of search keywords generator follows closely to the line of theoretical principles by  
63 Lasswell and intends to track conflict and cooperation dynamics on transboundary rivers by answering Lasswell' question

64 involved with seven elements. It can help to reveal the evolutionary dynamics and patterns of transboundary water conflicts  
65 and cooperation on a global scale, collecting news media datasets with an automated approach, and minimizing the manual  
66 workload of screening, reading and understanding the relevant news media articles, and provides researchers with powerful  
67 tools to retrieve useful information in related fields. It can serve as the foundation for further analysis, e.g., to study the attitudes,  
68 the topics of concerns, and the relationship between the evolution of the water governance network and the level of integrated  
69 water management along the evolution of water conflict and cooperation in the transboundary river basins. Ultimately it can  
70 contribute to understanding of the driving mechanisms and transformation laws of water conflict and cooperation. By capturing  
71 the characteristics of life cycles of water conflicts and cooperation, future researchers can explore the temporal evolution trend  
72 and spatial distribution law of global transboundary water conflicts and cooperation events, as well as the guiding significance  
73 of appropriate policy intervention, and improve the level of global water security.

74 Meanwhile, it can also serve as a methodological foundation of quantifying the social dimension in socio-hydrological  
75 approaches of understanding transboundary river system. Recent attempt has been made to take socio-hydrological approach  
76 to tackle the feedback mechanism of co-evolved sub-systems (Lu et al., 2021). While socio-hydrological model can contribute  
77 to understanding the complexity of the intertwined nature of transboundary river system, quantifying the social variable has  
78 been challenging in general. There has been increasing recognition that news media provides a valid proxy to reflect the  
79 changing values and interest of each riparian country (Wei et al., 2021), conflict and cooperation sentiments that reflected from  
80 news article have been adopted in socio-hydrological models as the willingness of cooperation to validate the social sector of  
81 the model (Lu et al., 2021). When expanded to other river basins, this study could provide a methodological support in  
82 measuring the social sector of transboundary river systems more effectively.

## 83 **2 Data and Method**

84 This study attempts to build a revised methodological framework that reflects the dynamics of water conflicts and cooperation  
85 among all the transboundary rivers in the globe. Overall procedures in the revised framework are illustrated in Figure 1. The  
86 method can be divided into three steps: Step 1 Select Database, Step 2 Keyword Determinants, Step 3 Data Cleaning and  
87 Processing. More specifically, the method begins with selecting news database in Step 1, detailed criteria to select news  
88 databases is stated in Sect. 2.1.1. Search keywords are generated in Step 2 with 5 blocks of keywords determinants. These 5  
89 blocks concern with river basin characteristics and the research question and determine the validity and relevance of the data  
90 to be collected. Using generated keyword in Step 2, original dataset is downloaded for data cleaning and processing in Step 3,  
91 which include rough manual reading and sorting to check results relevance in order to feedback on further keywords  
92 modification in Step 2. Trial-and-errors between Step 2 and Step 3 promise satisfactory keywords setting for the research. In  
93 Sect.2.4, several potentials for analysis in the future are introduced, which are extended applications for this methodological  
94 framework.

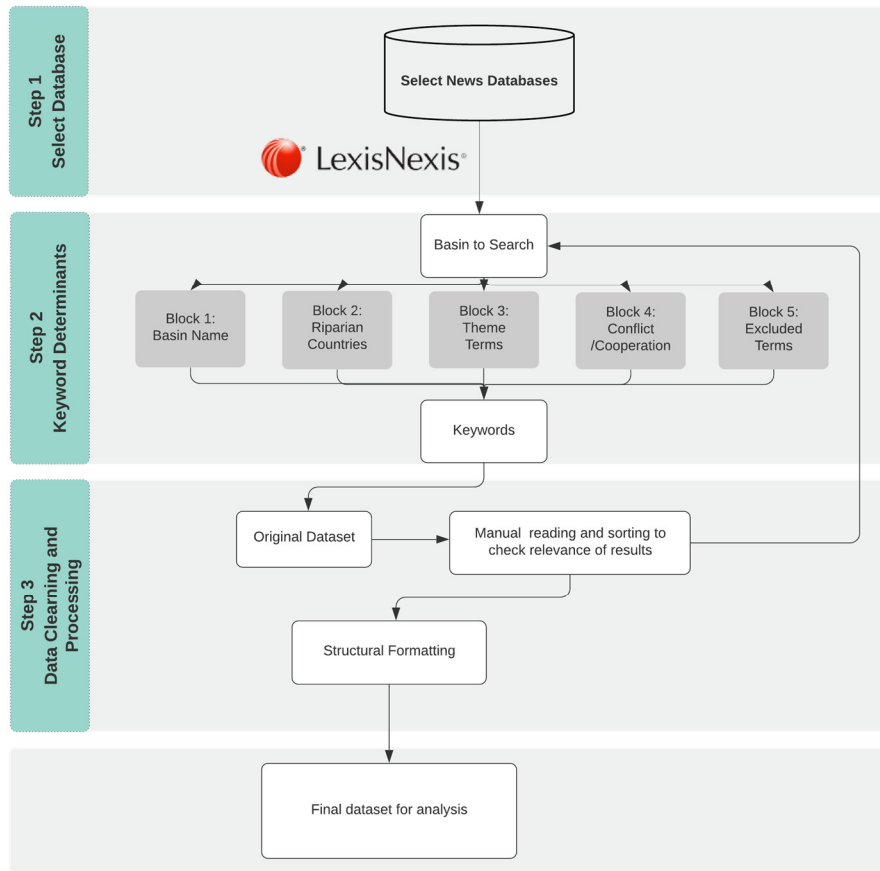


Figure 1. Method flow chart

## 2.1 Step 1: Select Database

### 2.1.1 News Media as Data Source

Choice of media sources should accord closely with the research goal. Our research goal is to track conflict and cooperation dynamics on transboundary rivers, which requires the data to cover water events and public opinion in a relatively long period of time. Also, newspapers (both print news and web news) published by professional journalists and editors are more suitable to use as data sources to reflect opinions of communities than social media (e.g., Twitter) as reflections of individual opinions. News media reflect what is important for the individual country/sector they are published within (Cooper, 2005), it thus has increasingly been studied by researchers to gain insight into transboundary water issues. The local news media is the first-hand material that reflect attitude/perception riparian countries held for their shared water and the involved stakeholders when discussing the water events in the transboundary river basin. In parallel, international news media serve as a good source of information to understand viewpoints from international audience that are outside of the river basins. Together, text analysis

107 of both regional and international news for water events in transboundary rivers can reveal the full picture of the ongoing  
108 dynamics in the river basin.

## 109 **2.1.2 Select News Database**

110 The very first step of this method involve selecting a news database that covers comprehensive news sources spanning across  
111 the globe. The selected media databases should include longitudinal coverage (i.e., can be traced back to decades) and updated  
112 in a timely manner, such as Lexis Advance (a product of Lexis Nexis Corporation), ProQuest, Factiva, etc. Lexis Advance  
113 covers more than 6,000 mainstream news media in most countries and regions around the world with a long-term coverage,  
114 and is one of the most commonly used news sources in the field of social sciences (Weaver & Bimber, 2008; Racine et al.,  
115 2010). Therefore, Lexis Advance is taken as an example of news media database to demonstrate the process of obtaining news  
116 media data of transboundary water conflicts and cooperation, and other suitable databases can, of course, be feasible options.  
117 Although the temporal coverage is affected by the level of media development in different regions, the covered timeframe  
118 spans over one hundred years to date, providing good data support on tracking media coverage of transboundary water conflict  
119 and cooperation research. The scope of research limit to English newspaper only due to our limitation of language processing,  
120 which is considered as sufficient enough to meets the requirements for extensive coverage of transboundary water conflicts  
121 and cooperative research.

## 122 **2.2 Step 2: Keyword Determinants**

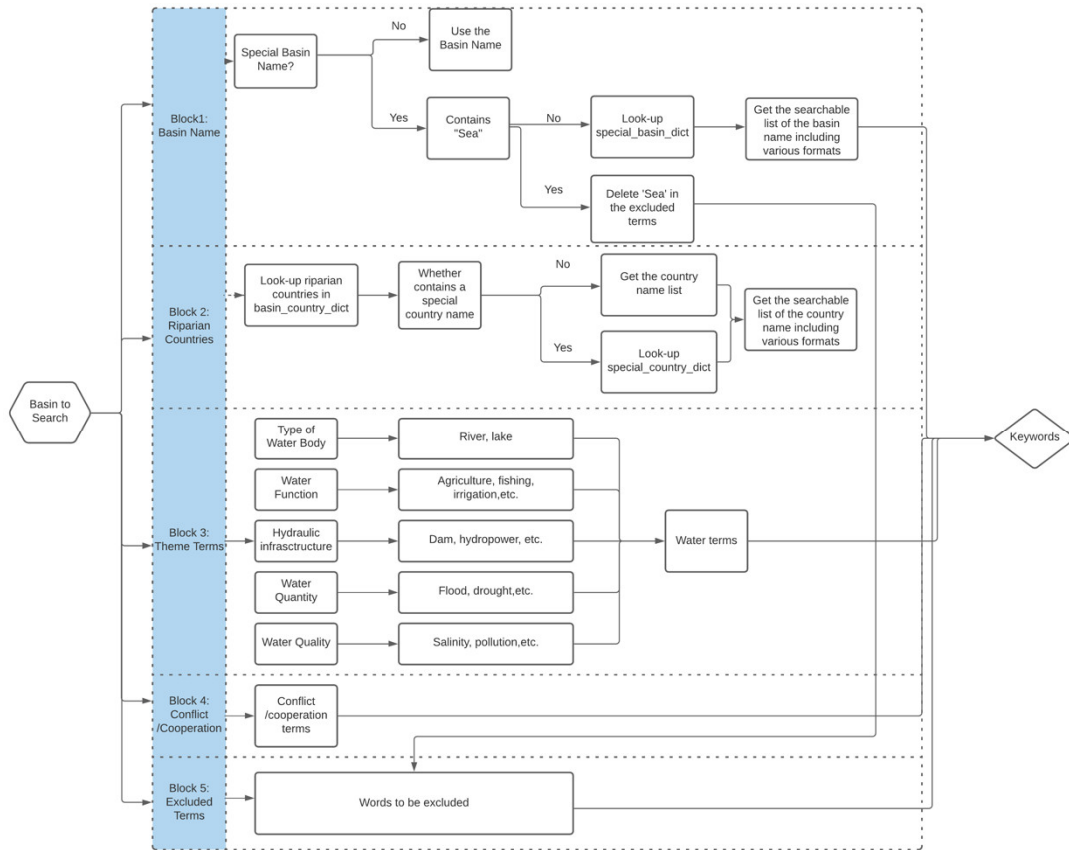
### 123 **2.2.1 Select Rivers to Search**

124 The scope of rivers to search in this study are 286 transboundary rivers as identified in 2016 (Transboundary Waters  
125 Assessment Programme, 2016). It is understood that the total number of transboundary rivers are recently been updated to 310  
126 (McCracken & Wolf, 2019), which are due to advancement of remote-sensing technology. Remote sensing can examine the  
127 two fundamental characteristics of transboundary rivers (common terminus and perennial), thus finer resolution of hydrologic  
128 data assists in discovering new transboundary rivers. In general, majority of the 24 newly added basins are small in area (less  
129 than 10,000 km<sup>2</sup>) (McCracken & Wolf, 2019), and are considered as inactive in conflicts and cooperation dynamics. Therefore,  
130 this study holds on 286 transboundary rivers in the procedure of Select Rivers to Search, which can be extended to 310 in the  
131 future. Four river basins were taken as case studies, Mekong, Nile, Columbia, and Ganges-Brahmaputra-Meghna (hereafter as  
132 GBM) as the global hotspot of water events.

### 133 **2.2.2 Search Keywords Generator**

134 The search terms are one of the key determinants of the coverage and relevance of the data to be retrieved. This study develops  
135 a keyword generator that allow efficient generating of keywords terms, which are applicable to all transboundary river basins  
136 (286 rivers basins) in the world. The keyword determinants are developed on the basis of TFDD (Yoffe & Larson, 2001) and

137 further revised to include five blocks of terms (as shown in Figure 2). These five blocks aim to include in which river basin  
 138 (Block 1), who (riparian countries, Block 2), regarding what issues (Block 3), have resulted in Conflict/Cooperation status  
 139 (Block 4). More specifically, Block 1 and Block 2 are basic information about the river basin, such as name of the river basin,  
 140 and various formats of riparian countries' names, retrieved articles need to discuss the conflictive or cooperative aspects of the  
 141 events involving at least one of riparian countries; Block 3 contains theme terms regarding of various functions of the water  
 142 body, topics discussing hydraulic infrastructure, water quality, agriculture/fishing, or any other specific topics with associated  
 143 terms; Block 4 include keywords indicate conflict or cooperation; and Block 5 consist of keywords to be excluded which bring  
 144 in irrelevance. The above five blocks can narrow down the search to the desired scope, with the list of unwanted words further  
 145 screen out irrelevant topics, after which, the search results can achieve a balance between coverage and relevance, that is,  
 146 neither too much relevant information is missed, nor too much irrelevant information is included.



147  
 148 **Figure 2.** Search Keywords Generator flow chart

149 **(1) Block 1: Basin Name**

150 This study customizes relatively general algorithms to generate search strings for river basins with different attributes and  
 151 conducts special treatments for individual river basins, so that each river basin is under the general search rules resulting in a

152 considerable number of search results with a balance of coverage and accuracy. The aim of **Block 1** is to get the searchable list  
 153 of the basin name including various formats and consider special treatments for specific categories of basin names. There are  
 154 several categories identified for different variations of basin names, see below for specific information.

- 155 a) Basin name same as the name of a certain riparian country or state; the search results are likely to contain many articles  
 156 about the internal affairs and diplomacy of the country or state. The detailed list of this type of basins is shown in Table  
 157 1. When talking about transboundary water issues, people usually focus on interactions on the scale of local communities  
 158 and riparian states rather than intercontinental, and do not refer to the Continent names. Therefore, raising the frequency  
 159 of continent name in search keywords will only compress data volume of relevant articles significantly, but not improve  
 160 the data relevance pertaining to the research goal. However, river basins with the same names but located in different  
 161 continents have different riparian countries. Adding frequency setting of riparian countries will filter out articles about  
 162 the river on the other continent effectively. For example, St. John rivers appear both in Africa (flowing through  
 163 Côté'd'Ivoire, Guinea, and Liberia) and North America (flowing through the United States and Canada). Raising frequency  
 164 of riparian countries rather than continent names contributes more to the data relevance.
- 165 b) Basin name contains commonly used words, for example, Amazon, which not only refers to the Amazon river basin, but  
 166 also an e-commerce company in the United States. More filters will be adopted in this case to ensure relevance rate. See  
 167 Table 1 for a detailed list of this type of river basins.
- 168 c) Basin name contains words such as 'Lake' or 'Sea', the word frequency setting for 'River' in the search string needs to  
 169 be modified, and that for 'Lake' needs to be increased, or 'Sea' needs to be removed from the list of noise keyword. See  
 170 the detailed list of this type of river basin in Table 1.
- 171 d) Other categories of basin names that require special treatment (see Table 1 for details) are: river basins have different  
 172 names, such as upstream and downstream rivers are designated with different names, or the river basin contains multiple  
 173 rivers; rivers in the basin have different names; the basin name is composed of multiple words; similar basin names exist  
 174 on different continents; the basin name contains 'St.', but may be referred as 'Saint' in media articles.

175 **Table 1.** Categories of basins need special treatment

Categories of basins need special treatment	Basin names	Treatment
Basin name includes state's or district's name	Belize; Columbia; Congo/Zaire; Corredores/Colorado; Gambia; Jordan; La Plata; Mississippi; Nelson-Saskatchewan; Niger; Senegal; Tigris-Euphrates/Shatt al Arab	Raise the frequency setting for 'water' or 'river' etc. to filter out the geopolitical articles as many
Basin name includes common word	Amazon; Baker; Columbia; Cross; Don; Fly; Han; Lagoon Mirim; Lotagipi Swamp; Massacre; Negro; Oral/Ural; Orange; Rhone; Red/Song Hong; San Martin; Seno Union/Serrano; Vanimo-Green; Whiting	Raise the frequency setting for 'water' or 'river' etc. to filter out water-unrelated articles as many; or delete a certain percentage of articles from the end of the results list
Basin name includes 'Lake', 'Sea'	Lake Chad; Lake Fagnano; Lake Natron; Lake Prespa; Lake Titicaca-Poopo System; Lake Turkana; Lake Ubsa-Nur; Aral Sea	The word frequency setting for 'River' in the search string needs to be modified, and that for

		'Lake' needs to be increased, or 'Sea' needs to be removed from the list of noise keywords
Basin name includes multiple formats (maybe consists of multiple rivers)	Asi/Orontes; BahuKalat/Rudkhanehye; Bei Jiang/Hsi; Benito/Ntem; Ca/Song-Koi; Cancoso/Lauca; Carmen Silva/Chico; Coco/Segovia; Congo/Zaire; Corantijn/Courantyne; Corredores/Colorado; Cuvelai/Etoshia; Douro/Duero; Gallegos/Chico; Ganges-Brahmaputra -Meghna; Hamun-i-Mashkel/Rakshan; Hari/Harirud; Ili/Kunes He; Jenisej/Yenisey; Juba-Shibeli; Kura-Araks; Lava/Pregel; Mana-Morro; Nelson-Saskatchewan; Oder/Odra; Oiapoque/Oyupock; Oral/Ural; Red/Song Hong; Seno Union/Serrano; Shu/Chu; Tagus/Tejo; Tigris-Euphrates/Shatt al Arab; Tjeroaka-Wanggoe; Torne/Tornealven; Vanimo-Green; Vistula/Wista	Contain all formats of related basin/river names in the search keywords
Basin name consists a river with multiple names	Muhuri (aka Little Feni)	Contain all formats of related river names in the search keywords
Basin name includes multiple words	An Nahr Al Kabir; Astara Chay; Coatan Achute; El Naranjo; Great Scarcies; Har Us Nur; Kowl E Namaksar; La Plata; Lagoon Mirim ; Lotagipi Swamp; Lough Melvin; Nahr El Kebir; Oued Bon Naima; Pu Lun T'o; Rio Grande (N. America); Rio Grande (S. America); San Martin; Song Vam Co Dong; St. Croix; St. John (Africa); St. John (North America); St. Lawrence; St. Paul; Wadi Al Izziyah	Add quotation mark to the basin name in the search keywords to search it as a whole, and prevent the basin name tokenized
Same basin names exist in multiple continents	Great/Little Scarcies; Rio Grande (N. America/S. America); St. John (Africa/North America)	Usually, articles do not contain the continent name when talking about rivers. Therefore, adding continent names into search keywords compresses data volume significantly and does not help with relevance. Adding frequency setting of riparian countries will filter out articles about the river on the other continent effectively.
Basin name includes St. (Saint)	St. Croix; St. John (Africa); St. John (North America); St. Lawrence; St. Paul	Put 'saint' and 'St.' into search keywords together

176 The *special\_basin\_dict* in the toolkit in **Block 1** is a python dictionary uploaded on Zenodo, whose *keys* are basin names with  
177 multiples words, or with special characters (e.g., back slash, dash, or parenthesis), and *values* are all searchable formats of the  
178 related basin names and river names. Given the original basin name to search, *special\_basin\_dict* can feedback its  
179 corresponding searchable keywords. If without *special\_basin\_dict* and using the original basin name to search, few results  
180 even none can be found. Coverage of retrieved results is enhanced by the *special\_basin\_dict*. When using the dictionary,  
181 import it to your script first, and call it easily.

## 182 (2) Block 2

183 **Block 2** is information concerning with riparian countries within the transboundary river basin. The aim of **Block 2** is to get  
184 the searchable list of the riparian country names including various formats. To fulfill the task, two helpful dictionaries -  
185 *basin\_country\_dict* and *basin\_country\_dict* are developed and provided in the toolkit of this study.

186 The *basin\_country\_dict* in the toolkit in **Block 2** is a python dictionary uploaded on Zenodo, whose *keys* are basin names, and  
187 *values* are all riparian countries located in the transboundary basin. Given the basin name to search, *basin\_country\_dict* can



188 feedback the list of riparian countries. Another python dictionary used in **Block 2** is *special\_country\_dict*, whose *keys* are  
189 country names with various formats, or with special characters (e.g., dot), *values* are all the searchable formats of the country  
190 name. Given the special country name to search, *special\_country\_dict* can feedback the list of all searchable formats of the  
191 country name.

192 Given a basin name to search, first looking up riparian countries in the *basin\_country\_dict* gets the list of riparian countries;  
193 then check whether there is a special country name in the list of riparian countries. If yes, through looking up  
194 *special\_country\_dict*, all searchable list of the country name including various formats are generated in **Block 2**.

### 195 (3) **Block 3**

196 **Block 3** contains terms concerning various themes of transboundary water resources, shown in Table 2. For example, type of  
197 water body, function of water body (agriculture, fishing etc.), hydraulic infrastructure, water quantity, water quality, and other  
198 specific topics which arouse certain research interests.

### 199 (4) **Block 4**

200 **Block 4** contains conflict/cooperation related keywords, adopted from TFDD searching keywords (Yoffe & Larson, 2001),  
201 shown in Table 2. If you focus on a certain type of conflict/cooperation, keywords in **Block 4** can be modified accordingly. In  
202 addition, UNBIS Thesaurus (UNBIS Thesaurus, 2021) provides lists of related keywords for conflict and cooperation which  
203 can be referred to.

### 204 (5) **Block 5**

205 **Block 5** contains excluded terms, given the research goal of our study, most of which are adopted from TFDD searching  
206 keywords (Yoffe & Larson, 2001), shown in Table 2. These terms, seemingly relevant to our topics, occur in media articles  
207 massively and easily bring in lots of data noise. For example, ‘sea’ and ‘ocean’ bring mass of irrelevant articles talking about  
208 marine rights and navigational utilization; ‘nuclear’ refers to ‘nuclear power’ and ‘nuclear threaten’, which is not the main  
209 concern of transboundary water conflict and cooperation; and as for ‘flood of refugees’, though it contains the keyword ‘flood’,  
210 but is regarded as irrelevant to our topics. These terms prone to bring in noise should be excluded in searching results, and thus  
211 list in excluded terms in **Block 5**. If researchers employ our framework in their own study fields in the future, excluded terms  
212 to avoid noise in Block 5 should be modified accordingly to fit their own research field based on results of trial-and-error  
213 between Step 2 and Step3 and combined with their experience and knowledge background. For example, when collecting data  
214 for Aral Sea, ‘sea’ should be deleted from the excluded terms in **Block 5** to prevent great loss of data coverage.

215

216

**Table 2.** Example of keywords in Block 1-5

<b>Block 1: Basin name</b>	Basin name (5)	
<b>Block 2: Riparian countries</b>	Each riparian country (2)	
	<b>Type of water body</b>	Water (3), river (3), lake, stream, tributary, etc.

<b>Block 3: Theme terms of transboundary water resources</b>	<b>Function of water body</b>	Irrigation, fish, fish rights, water rights, water diplomacy, water hegemony, etc.
	<b>Hydraulic infrastructures</b>	Dam, diversion, channel, canal, hydroelect*, hydropower, reservoir, etc.
	<b>Water quantity</b>	Flood, drought*, water allocation, water sharing, etc.
	<b>Water quality</b>	Salinity, pollution, etc.
<b>Block 4: Conflict/cooperation terms</b>	<b>Conflict</b>	dispute*, conflict*, disagree*, war, troops, "letter of protest", hostility, "shots fired", boycott, protest*
	<b>Cooperation</b>	Treaty, agree*, convention, "framework directive", negotiat*, resolution, commission, secretariat, "joint management", "basin management", peace, "accord", "peace accord", settle*, cooperat*, collaborat*, bilateral, multilateral, sanction*
<b>Block 5: Excluded terms</b>	Sea, ocean, navigat*, nuclear, water cannon, light water reactor, mineral water, hold water, cold water, hot water, water canister, water tight, water down*, flood of refugees, oil, drugs, a stream of, flood of	

Notes: asterisk (\*) indicates root of a word; number in parentheses (5,2 or 3) indicate at least how many times the keywords should appear in a searching result

### 2.2.3 Term frequency setting of keywords

The setting of term frequency of keywords comes from the recursive trial-and-errors in the search process, which makes the search results for most transboundary river basins relatively satisfactory. For individual river basins, universal setting rules of term frequency will cause the search results drop to zero sharply or too many to cope with, and the accuracy of the search results cannot be guaranteed. For example, when collecting data on the Jordan River Basin, given that Jordan is not only the name of the river basin, but also the name of a riparian country in the basin, there are too many articles that meet all the search requirements but purely about regional politics. Therefore, the setting of term frequency for the keywords ‘water’ and ‘river’ needs to be increased to 5 times to highlight the theme of transboundary water resources and ensure that the search results have similar accuracy to other river basins.

Taking the Lancang-Mekong basin as an example, the search keywords used in this study are shown in Table 3. During the trial-and-error process, we found that the results relevance rate is far below acceptable level (less than 30%), therefore we revised the keyword terms to increase frequency of certain terms until satisfactory results are produced, for example, the name of the basin appears in the article were increased to at least five times, the name of any riparian country in the basin (official name or abbreviation) appears in the article at least two times. Water-related words are divided into three sub-blocks: type of water body, function of water body, and infrastructures for water conservancy. Among them, ‘water’ and ‘river’ appear at least 3 times respectively, and the rest keywords of water block appear at least once; words related to conflict, or cooperation appear at least once. Recordings of trial-and-error process for Mekong, Nile and Jordan River Basin are provided in Appendix to demonstrate the effects of various groups of frequency settings of keywords and how balance between relevance and coverage is approaching. Although term frequency settings of keywords and justification of balance between relevance and coverage in this study may not be optimal, with a certain degree of coexisting subjectivity and objectivity, they can also serve as a reference for other researchers.

**Table 3.** Search Keywords in the study (Lancang-Mekong as an example)

Key Word Search	Lexis Advance Database
Must Include the Basin Name (at least 5 times)	Mekong (5)
Includes at least one of the following countries' name (at least twice)	Thai*(2), Cambodia*(2), China(2), Chinese(2), Laos(2), Myanmar(2), Burm*(2), vietna*(2)
Includes at least one of the following words related to Water	Same as Block 3 (see Table 2)
Includes at least one of the following words related to Conflict/Cooperation	Same as Block 4 (see Table 2)
Does not include any of the following noisy words	Same as Block 5 (see Table 2)

Notes: asterisk (\*) indicates root of a word; number in parentheses (5,2 or 3) indicate at least how many times the keywords should appear in a searching result

### 239 2.3 Step 3: Data Cleaning and Processing

240 Before finalizing the refined datasets for further analysis, data cleaning and processing is indispensable. The first stage in Step  
241 3 is Rough Manual Reading and Sorting to Check Results Relevance, which aims to provide feedbacks on how to modify  
242 keywords in Step 2. Rough manual reading can be done by random sampling, or more conveniently from back to front. Since  
243 lists of news results by news media databases usually have options to Sort by Relevance, frontlines displayed in the front of  
244 the list of searching results are ranked as more relevant to search terms than that of the backlines of the list. ('Sort by Relevance'  
245 is one of the sorting functions provided by Lexis Advance, which also provides 'Sort by Date' and 'Sort by Document Title'.  
246 Among the three options, 'Sort by Relevance' works best for us to read roughly to change the frequency setting of keywords  
247 by trial-and-error. Therefore, 'Sort by Relevance' was chosen before downloading the data from Lexis Advance. Usually,  
248 news databases have similar functions for readers to read roughly and conveniently.) A proper percentage, like 80% of results  
249 which are relevant among all, can be set to meet our expectation.

250 To better facilitate future analysis, all downloaded text data will go through structure formatting process. A data structuring  
251 program is developed for Lexis Advance to download and organize the text data into structured format. The relevant media  
252 articles are processed in order of relevance, and detailed information such as the publication time of the articles, media source,  
253 author, article length, etc. are stored in a structured manner. An example of structured media data is shown in Table 4. As for  
254 data integration, any news data downloaded from suitable data sources (not only from Lexis Advance) can be arranged and  
255 structured in the format of Table 4 through data cleaning and processing procedure. After data processing, the toolkit provided  
256 by this research can be applied to the integrated data regardless of the original data sources of it.

257 **Table 4.** Example of structured data

<b>Paper Index</b>	1
<b>Title</b>	The 1997 water rights settlement between the state of Montana and the Chippewa Cree tribe of the Rocky Boy's Reservation: the role of community and of the trustee.
<b>Source</b>	ASAPII Database
<b>Date</b>	Dec 22, 1998
<b>Pg;ISSN;Vol;No</b>	Pg. 255(1); ISSN: 0733-401X; Vol. 16; No. 2
<b>Words Count</b>	18256 words
<b>Author</b>	Cosens, Barbara A.

<b>Body</b>	I. INTRODUCTION Established on September 7, 1916 "for Rocky Boy's Band of Chippewas and ... other homeless Indians,"(1) the Rocky Boy's Reservation is home to over 3,000 Tribal members. The Reservation's annual population growth rate is in excess of three percent...(original data is too long for demonstration, here is the excerpt)
-------------	--

258 **2.4 Potential Analysis**

259 The news media dataset of water conflict and cooperation on transboundary rivers allows for varieties of analysis in later stage.  
 260 This study lists several examples of potential analysis including event extraction, stakeholder analysis, sentiment analysis and  
 261 topic analysis.

262 **Event Extraction** from news articles is a conventional application of water conflict and cooperation dataset. Same as what has  
 263 been achieved by TFDD, water events both conflictive and cooperative, were extracted in relevant political science datasets  
 264 and news articles (Yoffe & Larson, 2001). Event Extraction requires concise and accurate information recognition and  
 265 extraction from latent content in text data. Since human coders perform better than machine programming (Howland et al.,  
 266 2006), human coding event extraction is recommended.

267 **Stakeholder Analysis** for transboundary rivers is a way to identify who has been involved in transboundary water issues, and  
 268 the roles they play in the game, i.e., understanding the demands and expectations of the major stakeholders inside and outside  
 269 the basin, based on typical definition of stakeholder analysis (Smith, 2000). News media represent or reflect the interests of its  
 270 home country, thus via analysis of news media sources in a transboundary basin, political positions and economic  
 271 interrelationships between riparian countries and other extra-territorial countries lying outside the basin are uncovered.  
 272 Longitudinal analysis has capability to depict the trajectories of a stakeholder country's interests and reveal the evolution of  
 273 stakeholder countries in transboundary water issues.

274 **Sentiment analysis** on the news media dataset on transboundary rivers can bring the implicit information to the surface (Jiang  
 275 et al., 2016), since willingness of cooperation and hostility of conflicts often hide behind the news articles. Positive and  
 276 negative sentiment are closely to dynamics of conflict and cooperation in transboundary water issues, which serve as precursors  
 277 of significant situational changes. Sentiment lexicons (Khoo & Johnkhan, 2018) or machine learning (Neethu & Rajasree,  
 278 2013) are major methods for sentiment analysis in text mining.

279 **Topic analysis** tells the story about main interests and concerns of the news media, even the stakeholders along with time  
 280 (Jacobi et al., 2016). Topics concerned along with society development, evolutionary trajectories of transboundary water issues  
 281 can be uncovered through popular algorithms of topic modelling analysis , such as LDA (Alsumait et al., 2009).

282 **3 Results**

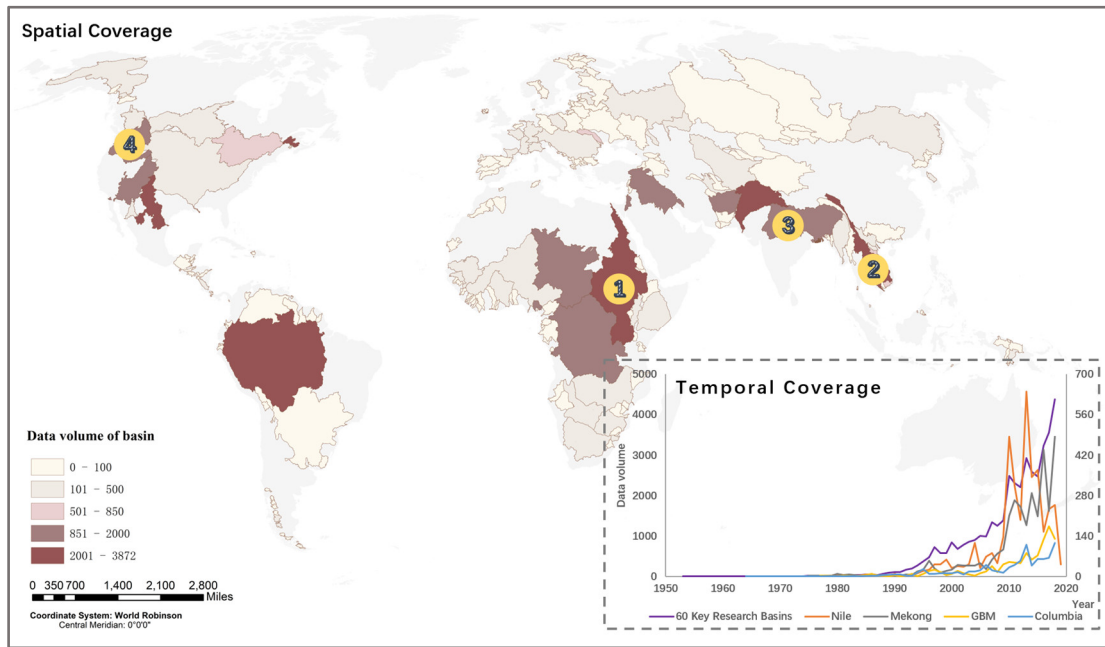
283 This section overviews the Global Datasets statistically both in terms of spatial coverage and content coverage, which aims to  
 284 show the datasets telling stories of conflict and cooperation on transboundary rivers from all aspects in a global scale. To  
 285 demonstrate the effectiveness of the methodological framework and toolkit, manual reading to check the improvements of data  
 286 relevance was conducted on four representative basins including Nile, Mekong, GBM, and Columbia.

287 **3.1 Overview of the Global Datasets**

288 **3.1.1 Spatial Coverage**

289 **(1) Continental Coverage**

290 With the customized search strings for each transboundary river and the data structured program developed for Lexis Advance  
291 to organize the data, as of March 10, 2019, the data volume results of 286 transboundary river basins around the world are  
292 shown Figure 3 - Spatial Coverage. In Figure 3, the base map of transboundary river basins around the world was downloaded  
293 from TFDD in the format of GIS shapefiles (Transboundary Freshwater Dispute Database, 2008).



294 **Figure 3.** Spatial coverage and temporal coverage in basin scale (①Nile; ②Mekong; ③GBM; ④Columbia)

296 Data volume of news articles reflects the prominence of the conflict and cooperation events discussed in transboundary river  
297 basins. Enough data volume promises statistical significance. The mainstream application of this news media dataset is further  
298 text mining to track conflict and cooperation dynamics on transboundary rivers. For text mining purpose, this study assumes  
299 arbitrarily that 100 media articles are the minimum data volume to track dynamics transboundary rivers along with time.

300 Overall, there are 60 river basins with more than 100 media articles, which are considered as the Key Research Basins of  
301 transboundary water conflict and cooperation in our research. The number of news articles discussing these 60 Key Research  
302 Basins reached more than 41,000. Among the 60 Key Research Basins, 16 river basins have more than 850 data records as  
303 shown in Table 5, which attract more attention and are considered as Heated Basins. Note that the definition criteria of Key  
304 Research Basins (more than 100 articles) and Heated Basins (more than 850 articles) are flexible and adaptive according to  
305 specific research demands.

306 **Table 5.** 16 Most-discussed Basins with more than 850 records

Order	Basin Name	Continent	Number of records	Countries
1	Nile	Africa	3872	Burundi, Central African Republic, Egypt, Hala'ib Triangle, Eritrea, Ethiopia, Kenya, Rwanda, Sudan, Abyei, South Sudan, United Republic of Tanzania, Uganda, Dem. Republic of the Congo
2	Mekong	Asia	3253	China, Cambodia, Lao People's Democratic Republic, Myanmar, Thailand, Viet Nam
3	Rio Grande (N. America)	North America	2718	Mexico, United States of America
4	Indus	Asia	2404	Afghanistan, China, India, Nepal, Pakistan
5	St. John (North America)	North America	2356	Canada, United States of America
6	Amazon	South America	2078	Bolivia, Brazil, Colombia, Ecuador, French Guiana, Guyana, Peru, Suriname, Venezuela
7	Colorado	North America	1975	Mexico, United States of America
8	Jordan	Asia	1816	Egypt, Israel, Jordan, Lebanon, West Bank, Syrian Arab Republic
9	Congo/Zaire	Africa	1391	Angola, Burundi, Central African Republic, Cameroon, Congo, Gabon, Malawi, Rwanda, Sudan, South Sudan, United Republic of Tanzania, Uganda, Dem. Republic of the Congo, Zambia
10	Lake Chad	Africa	1353	Central African Republic, Cameroon, Algeria, Libya, Niger, Nigeria, Sudan, Chad
11	Ganges-Brahmaputra - Meghna	Asia	1183	Bangladesh, Bhutan, China, India, Myanmar, Nepal
12	Helmand	Asia	1168	Afghanistan, Iran (Islamic Rep of), Pakistan
13	Cross	Africa	1110	Cameroon, Nigeria
14	Tigris-Euphrates/Shatt al Arab	Asia	939	Iran (Islamic Rep. of), Iraq, Jordan, Saudi Arabia, Syrian Arab Rep., Turkey
15	Columbia	North America	859	Canada, United States of America
16	Tijuana	North America	853	Mexico, United States of America

307 Most studies of conflict and cooperation on transboundary rivers focus on individual basins, which seeks solutions to dealing  
308 with local challenges on transboundary water resources (Bernauer & Böhmelt, 2020). Therefore, formation of general  
309 understanding of conflict and cooperation on transboundary rivers needs global data support besides expert on-site experience  
310 from research of individual basins. Many most-discussed transboundary river basins such as the Nile, Mekong, Indus, GBM,  
311 and Tigris-Euphrates/Shatt al Arab etc. are located in regions featured with frequent tensions and armed conflicts (Pohl et al.,  
312 2014), and are well-known by people. However, this study finds that there are also some river basins from the authors' point  
313 of view, which less attention has been paid to in the past in terms of transboundary water conflict and cooperation research,  
314 e.g., St. John River (North America), and Tijuana River.

315 Data volume of transboundary water conflicts and cooperation news articles on different continents: for Asia is 14454, for  
316 North America is 11306, for Africa is 10734, for Europe is 2674, for South America is 2498. It could possibly be attributed to  
317 the discrepant levels of economic development of major countries on each continent, or varied attention paid to discussion of  
318 management of transboundary rivers. The other important reason could be the linguistic variations. Since this paper chose  
319 English newspaper as the search scope, the large amount of data in North America and the small amount of data in Europe  
320 could be due to system bias caused by language preferences.

321 There are notably large amount of transboundary water conflicts and cooperation events reported in Asia and Africa, which  
322 indicates that transboundary water management is a major topic of peace and development in both Asia and Africa. Taking  
323 into consideration that most countries on these two continents do not speak English as their mother language, the existence of  
324 a large number of news media articles on transboundary water conflicts and cooperation between Asia and Africa, on the one  
325 hand, reflects the fervent concerns about the transboundary water resources, and the desires for peace and development; on the  
326 other hand, it also reflects that people around the world are more involved in transboundary water issues in Asia and Africa,  
327 and have invested heavily in the development and construction and pay close attention to these two rapidly developing and  
328 eye-taking continents.

## 329 ***(2) National Coverage***

330 News media data volumes from different countries in the world are shown in Figure 4 - Spatial Coverage. In Figure 4, the base  
331 map of countries around the world was downloaded from ArcGIS Hub in the format of GIS shapefiles (Esri Data and Maps,  
332 2021). It is seen that United States of America contributes 11515 news articles on transboundary water conflict and cooperation,  
333 ranking number one, both as a riparian stakeholder in the transboundary water issues with Canada and Mexico, and as an extra-  
334 terrestrial international stakeholder involving in the transboundary water issues on continents other than the North America.  
335 Since a country's development and utilization of transboundary freshwater resources inevitably involves relations with other  
336 riparian countries, and transboundary water cooperation and conflicts often involve broader economic and social ties between  
337 riparian countries, transboundary freshwater management is an important component of the diplomacy of riparian countries;  
338 on the other hand, due to factors such as global hegemony, transnational investment, colonial history and other factors,  
339 transboundary freshwater management often involves countries outside the region, becoming a stage for great powers to play  
340 (Mirumachi, 2015).

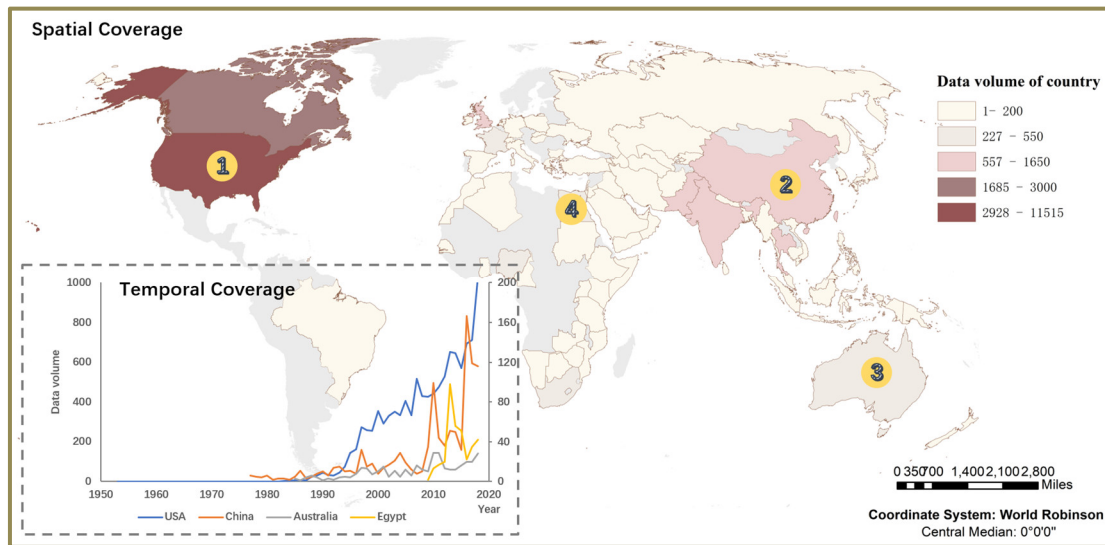


Figure 4. Spatial coverage and temporal coverage in country scale (①USA; ②China; ③Australia; ④Egypt)

### 3.1.2 Temporal Coverage

Temporal coverage of the datasets of 60 Key Research Basins (stated in Spatial Coverage section) and four case study basins are shown in Figure 3 - Temporal Coverage, which shows how many news media articles released along with years on transboundary water conflict and cooperation. Noted that due to differences of order of magnitudes, data series of 60 Key Research Basins uses the major vertical axis which ranges from 0 to 4500, and the four case study basins share the minor vertical axis which ranges from 0 to 700. The datasets cover from the year of 1953 to 2019. Boom of news articles on transboundary water conflict and cooperation emerges from 1990s, and potentially continues in the future. That emphasizes the necessity to revise the methodological framework and toolkit for news media dataset tracking of conflict and cooperation dynamics on transboundary rivers to cope with the era of big data. For the four case study basins, the changing trends of data volume display strong vibrates, which may be affected by certain water events and geopolitical relations in the river basins at the moment.

Temporal coverage of four representative countries, which are United States of America (USA, using the vertical minor axis on the left), China, Australia and Egypt (using the vertical major axis on the right), is shown in Figure 4 – Temporal Coverage. USA contributes the largest volume of data among countries in the world; China promotes transboundary cooperation in Mekong River Basin actively in the recent years; Australia does not have a transboundary river with other neighbouring countries, but releases lots of news articles on transboundary water issues; and Egypt is one of the major countries in Nile River Basin, which is representative in transboundary water conflict and cooperation. Same with the temporal coverage of basin analysis, country datasets also cover from the year of 1953 to 2019. Data volume took off from 1990s, and potentially continues in the future as well. For the four representative countries, the overall trends of data volume go up along with time and are affected by contextual events in the country to show strong vibrates.



### 363 3.1.3 Content Coverage

364 Word frequency analysis demonstrates that this study has generated good datasets tracking of conflict and cooperation  
365 dynamics on transboundary rivers. In the datasets, words concerning with water body function, hydraulic infrastructure  
366 construction, basin management, national power, civic rights, jointed research and water conflict and cooperation appear in a  
367 high frequency, consistent with the related keywords in TFDD (Yoffe & Larson, 2001) and relevant words provided in UNBIS  
368 Thesaurus (UNBIS Thesaurus, 2021). This indicates that the datasets are closely corresponding to the research question,  
369 providing data as needed.

### 370 3.2 Relevance Screening

371 The major advancement of this methodological framework is that it allows efficient and effective tracking of transboundary  
372 rivers conflict and cooperation events. The keywords generator developed in this study could result in an acceptable level of  
373 relevance without too much manual coding intervention. To demonstrate the effectiveness of this methodological framework,  
374 four river basins: Mekong, Columbia, Nile, and GBM were taken as case studies to conduct manual coding process. Two  
375 manual coders, who were trained beforehand, were involved to work independently for the four basins in the coding process.  
376 Each one undertook half of the total workload in which articles in the datasets were divided into two groups randomly. Before  
377 starting, inter-coder reliability test was conducted. The test randomly selects 50 articles from the datasets for two coders to  
378 read, differences were then discussed, and definitions were given to reach common understanding. Krippendorff's Alpha-  
379 Reliability was calculated as 0.81, which is considered as valid and consistent (Krippendorff, 2004).

380 The total number of downloaded articles, after removing duplicates by the function of removing duplicates in the data panel  
381 of Microsoft Excel, and the remaining number of relevant articles with removal of the duplicates are shown in the Table 6.  
382 The calculation Equation of Relevance Percentage is shown as Eq.1.

383 **Table 6.** Manual reading results of representative river basins

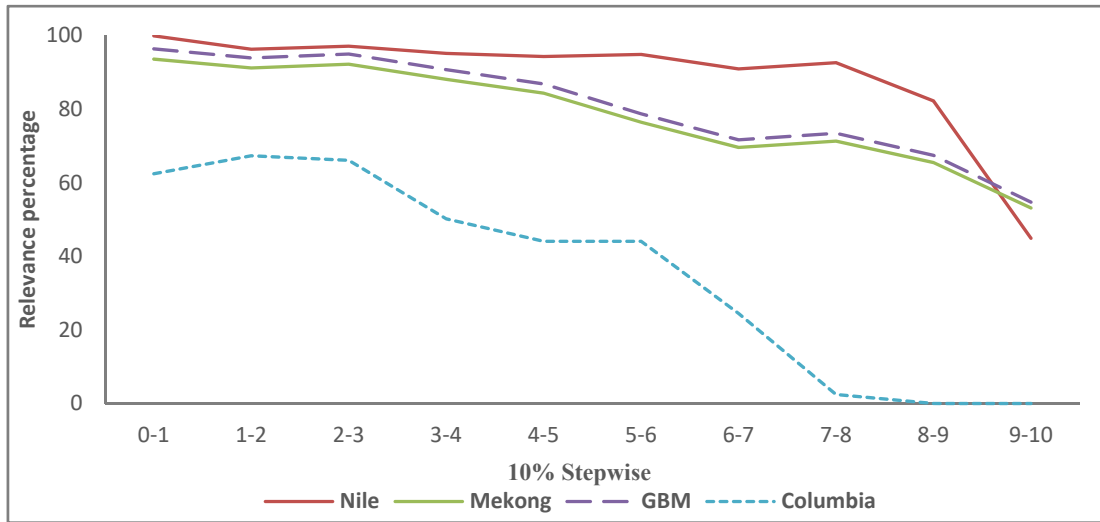
Basin name	Number of downloaded	Number after removing duplicates	Number after removing irrelevant	Relevance percentage (%)
Nile	3872	3563	3164	88.80
Mekong	3253	2917	2291	78.54
GBM	1183	1092	724	66.30
Columbia	859	817	295	36.11

384 Relevance percentage =  $\frac{\text{Number after removing irrelevant}}{\text{Number after removing duplicates}} \times 100\%$  , (1)

385 The last column of Table 6 shows the Relevance Percentage for the four river basins in a descending order. The relevance  
386 percentage of Nile, Mekong and GBM are at acceptable level, and that of Columbia is less satisfying. This is due to Columbia  
387 belonging to special basin name category, details shown in 2.2.2 a), whose basin name is same as the name of a certain riparian  
388 country or state. To further investigate of relevance percentage of the four basins, the relevance percentage in 10% stepwise is

389 calculated for each basin using Eq. 2. The relevance percentage in 10% stepwise for the four basins is shown in Figure 5. In  
 390 Figure 5, the horizontal axis is every 10% Stepwise segment of the news media articles data, and the vertical axis indicates the  
 391 Relevance Percentage of that segment of data.

392 Relevance percentage in 10% stepwise =  
 393 Relevance percentage for every 10% of the total number after removing duplicates , (2)  
 394



395  
 396 **Figure 5.** Relevance Percentage in 10% Stepwise for the four basins

397 The purpose of Figure 6 is to demonstrate the necessity to apply special treatments for some river basins. Since the datasets  
 398 retrieved from Lexis Advance are sorted by relevance, frontlines are naturally more relevant than the backlines, and the  
 399 Relevance Percentage in 10% Stepwise displays descending trendlines. However, slopes of the trendlines of relevance  
 400 percentage in 10% stepwise between basins reflect heterogeneity of data quality. The Relevance Percentage for Columbia is  
 401 unsatisfactory even in the first 10% of the article list, since ‘Columbia’ is both a district’s name and a commercial brand, listed  
 402 in Table 1. It makes sense that the data quality of Columbia River Basin is not as good as others. Special treatment for Columbia  
 403 should be adopted here to improve the data quality, as well as special treatments are needed for certain categories of basins  
 404 and corresponding treatments as mentioned in Table 1. To do so, usually enforcement of the frequency constraints shown in  
 405 Sect. 2.2.3 (i.e., raise the frequency setting for ‘water’ and ‘river’ to filter out the geopolitical articles as many), or removal of  
 406 the most irrelevant articles in the end of the dataset work well. With an anticipation of relevance percentage in mind, random  
 407 sampling or manual reading of the last percentage of articles are often undertook to check the data quality. For example, given  
 408 the Relevance Percentage in 10% Stepwise for Columbia, raising the frequency setting of ‘water’ and ‘river’ to 5 times, or  
 409 removal of the last 40% of the data retrieved in the Original Dataset due to its low relevance in general are feasible solutions

410 to improve data relevance. For other basins with satisfactory data relevance, no further operation is needed; and for the other  
411 basins, similar operations as for Columbia River Basins can be adopted before further potential analysis.

#### 412 **4 Summary**

413 Management of transboundary rivers is challenging both in terms of political and environmental in the 21st century. Data  
414 support is crucial for research of conflict and cooperation on transboundary rivers. Conventional construction manner of dataset  
415 by manual reading and extraction cannot meet the requirement for fast-updating in the big data era. This study brings up a  
416 revised methodological framework based on the conventional and toolkit for news media dataset tracking of conflict and  
417 cooperation dynamics on transboundary rivers. Design of the framework follows closely Lasswell' communication model  
418 (Lasswell, 1948) involved with seven elements- "who, with what intentions, in what situations, with what assets, using what  
419 strategies, reaches what audiences, with what result". Basic search keywords were adopted from TFDD and further revised to  
420 include five blocks of terms to make it extensible and adjustable according to a certain research topic. Through Block 1 and  
421 Block 2 with corresponding toolkit (shown in Figure 2), a dataset covering transboundary rivers in a global scale can be  
422 generated, which is improved than results of TFDD. All the special treatments for basin names (shown in Table 1), country  
423 names (shown in Block2), and term frequency setting of keywords (stated in Section 2.2.3) are crucial measures to enhance  
424 data quality and save manual efforts, which are improvement beyond achievements of TFDD. Following the methodological  
425 framework, a dataset with good trade-offs between data relevance and coverage is generated. This study demonstrates the  
426 effectiveness of the framework and the potency of our toolkit. This framework possesses extensibility and compatibility to  
427 other research topics besides transboundary water resources management since the search terms are adaptive and the toolkit is  
428 transplantable for related future research. With this revised framework and toolkit, research using news media tracking of  
429 conflict and cooperation dynamics on transboundary rivers will be much easier and more practicable.

430 Still this study has some limitations which could be overcome in following researches: (1) absence of newly-registered rivers:  
431 the list of transboundary rivers adopted in this study includes 286 rivers, which could be expanded to 310 rivers in the near  
432 future; (2) language limitation: the scope of this study limits to English newspaper only due to our limitation of language  
433 processing, which could be expanded to include more main languages and local languages in transboundary river basins; (3)  
434 absence of tributary information: in the keywords generator, tributaries of transboundary rivers are not included, which may  
435 lose content coverage to some extent. Future work can add more details concerning tributaries of transboundary rivers.

436 *Code and data availability.* The data and code used in this study are publicly available on Zenodo (including: basin-country dictionary;  
437 dictionary of country names with different formats (special country dictionary); dictionary of basin names with different formats; python  
438 code of searching term generator. DOI: 10.5281/zenodo.5112624

439 *Author contributions.* LG, JW, and FT designed the research framework. LG collected data and conducted data analysis. LG, JW, and KZ  
440 conducted manual reading for the case studies. LG, JW and FT composed the manuscript with contribution from KZ.

441 *Competing interests.* The authors declare that they have no conflict of interest.

442 *Special issue statement.* This article is part of the special issue “Socio-hydrology and transboundary rivers”. It is not associated with a  
443 conference.

444 *Acknowledgements.* We would like to acknowledge the National Key Research and Development Programme of China (grant no.  
445 2016YFA0601603) for the funding and support of this research.

446 *Financial support.* This research has been funded by the National Key Research and Development Programme of China (grant no.  
447 2016YFA0601603).

448 *Review statement.*

## 449 **References**

- 450 Alsumait, L., Barbará, D., Gentle, J., & Domeniconi, C. (2009). Topic significance ranking of LDA generative models. *In Proceedings of*  
451 *the European Conference on Machine Learning and Knowledge Discovery in Databases (ECML/PKDD)*, 67–82.
- 452 Bernauer, T., & Böhmelt, T. (2020). International conflict and cooperation over freshwater resources. *Nature Sustainability*, 3(5), 350–356.  
453 <https://doi.org/10.1038/s41893-020-0479-8>
- 454 Cooper, S. (2005). Bringing Some Clarity to the Media Bias Debate. *Communications Faculty Research*.  
455 [https://mds.marshall.edu/communications\\_faculty/2](https://mds.marshall.edu/communications_faculty/2)
- 456 Esri Data and Maps. (2021, April 14). *World Countries (Generalized)*. ArcGIS Hub.  
457 [https://hub.arcgis.com/datasets/2b93b06dc0dc4e809d3c8db5cb96ba69\\_0](https://hub.arcgis.com/datasets/2b93b06dc0dc4e809d3c8db5cb96ba69_0)
- 458 Howland, D., Becker, M. L., & Prelli, L. J. (2006). Merging content analysis and the policy sciences: A system to discern policy-specific  
459 trends from news media reports. *Policy Sciences*, 39(3), 205–231. <https://doi.org/10.1007/s11077-006-9016-5>
- 460 Jacobi, C., Atteveldt, W. van, & Welbers, K. (2016). Quantitative analysis of large amounts of journalistic texts using topic modelling.  
461 *Digital Journalism*, 4(1), 89–106. <https://doi.org/10.1080/21670811.2015.1093271>
- 462 Jiang, H., Qiang, M., & Lin, P. (2016). Assessment of online public opinions on large infrastructure projects: A case study of the Three  
463 Gorges Project in China. *Environmental Impact Assessment Review*, 61, 38–51. <https://doi.org/10.1016/j.eiar.2016.06.004>
- 464 Khoo, C. S., & Johnkhan, S. B. (2018). Lexicon-based sentiment analysis: Comparative evaluation of six sentiment lexicons. *Journal of*  
465 *Information Science*, 44(4), 491–511. <https://doi.org/10.1177/0165551517703514>
- 466 Krippendorff, K. (2004). Reliability in Content Analysis. *Human Communication Research*, 30(3), 411–433. [https://doi.org/10.1111/j.1468-](https://doi.org/10.1111/j.1468-2958.2004.tb00738.x)  
467 [2958.2004.tb00738.x](https://doi.org/10.1111/j.1468-2958.2004.tb00738.x)
- 468 Lasswell, H. (1948). *The Structure and Function of Communication in Society. The Communication of Ideas*. (Lyman Bryson (1948) (ed.)).  
469 The Institute for Religious and Social Studies.
- 470 McCracken, M., & Wolf, A. T. (2019). Updating the Register of International River Basins of the world. *International Journal of Water*  
471 *Resources Development*, 35(5), 732–782. <https://doi.org/10.1080/07900627.2019.1572497>

472 Neethu, M. S., & Rajasree, R. (2013). Sentiment analysis in twitter using machine learning techniques. *2013 Fourth International Conference*  
473 *on Computing, Communications and Networking Technologies (ICCCNT)*, 1–5. <https://doi.org/10.1109/ICCCNT.2013.6726818>

474 Pohl, B., Carius, A., Conca, K., Dabelko, G., Kramer, A., Michel, D., Schmeier, S., Swain, A., & Wolf, A. (2014). *The rise of hydro-*  
475 *diplomacy. Strengthening foreign policy for transboundary waters*. <https://doi.org/10.13140/2.1.4035.5848>

476 Racine, E., Waldman, S., Rosenberg, J., & Illes, J. (2010). Contemporary neuroscience in the media. *Social Science & Medicine (1982)*,  
477 *71(4)*, 725–733. <https://doi.org/10.1016/j.socscimed.2010.05.017>

478 Sadoff, C. W., & Grey, D. (2005). Cooperation on International Rivers: A Continuum for Securing and Sharing Benefits. *Water International*,  
479 *30(4)*, 420–427. <https://doi.org/10.1080/02508060508691886>

480 Smith, L. W. (2000). *Stakeholder analysis: A pivotal practice of successful projects*. Project Management Institute Annual Seminars &  
481 Symposium, Houston, TX. Newtown Square, PA. [https://www.pmi.org/learning/library/stakeholder-analysis-pivotal-practice-](https://www.pmi.org/learning/library/stakeholder-analysis-pivotal-practice-projects-8905)  
482 [projects-8905](https://www.pmi.org/learning/library/stakeholder-analysis-pivotal-practice-projects-8905)

483 *Transboundary Freshwater Dispute Database | Program in Water Conflict Management and Transformation | Oregon State University*.  
484 (2008). <https://transboundarywaters.science.oregonstate.edu/content/transboundary-freshwater-dispute-database>

485 Transboundary Waters Assessment Programme. (2016, January). *Transboundary River Basins-Status and Trends*. [http://twap-](http://twap-rivers.org/assets/GEF_TWAPRB_FullTechnicalReport_compressed.pdf)  
486 [rivers.org/assets/GEF\\_TWAPRB\\_FullTechnicalReport\\_compressed.pdf](http://twap-rivers.org/assets/GEF_TWAPRB_FullTechnicalReport_compressed.pdf)

487 UNBIS Thesaurus. (2021). *UNBIS Thesaurus*. <http://metadata.un.org/thesaurus/?lang=en>

488 United Nations. (2019). *Progress on Transboundary Water Cooperation 2018: Global Baseline for SDG 6 Indicator 6.5.2*. UN.  
489 <https://doi.org/10.18356/f6afa45b-en>

490 Weaver, D. A., & Bimber, B. (2008). Finding News Stories: A Comparison of Searches Using Lexisnexis and Google News. *Journalism &*  
491 *Mass Communication Quarterly*, 16.

492 Wolf, A. T. (1999). The Transboundary Freshwater Dispute Database Project. *Water International*, *24(2)*, 5.

493 Wolf, A. T., Natharius, J. A., Danielson, J. J., Ward, B. S., & Pender, J. K. (1999). International River Basins of the World. *International*  
494 *Journal of Water Resources Development*, *15(4)*, 387–427. <https://doi.org/10.1080/07900629948682>

495 Yoffe, S., & Larson, K. (2001). *CHAPTER 2 BASINS AT RISK: WATER EVENT DATABASE METHODOLOGY*. 36.

496

497

499 Recordings of trial-and-error process are provided as follows to demonstrate the effects of various groups of frequency settings  
500 of keywords and how balance between relevance and coverage is approaching. Two justification indicators of data relevance  
501 are adopted: (1) Indicator 1: the number of articles relevant to our research topic within 20 articles at 60% of total data volume.  
502 For example, there are 10000 articles retrieved by the certain frequency setting of search terms in Lexis Advance, we locate  
503 the article at exactly 60% of total data volume, which is the 6000<sup>th</sup> articles, and read 20 articles from there. Therefore, Indictor  
504 1 is how many articles are relevant among 6001<sup>st</sup>-6020<sup>th</sup> articles. (2) Indicator 2: the number of articles relevant to our research  
505 topic within 20 articles at 80% of total data volume. Similar to the algorithm of Indicator 1, if the total data volume is 10000,  
506 Indictor 2 is how many articles are relevant among 8001<sup>st</sup>-8020<sup>th</sup> articles. Table 7 presents the results of trial-and-error process  
507 of frequency settings of keywords for Mekong, Nile and Jordan River Basin, and shows that strong frequency settings enhance  
508 data relevance prominently, and at the same time reduces data volume to a large extent. To promise a balance between data  
509 relevance and coverage, proper frequency settings of search keywords should be adopted. In this study, Test 6 is adopted as  
510 the final setting. Notice that Nile and Jordan River Basin have overwhelmingly large volume of data if no additional constraint  
511 are exerted, therefore ("Nile River" OR "Nile Basin" OR "Nile Water") or ("Jordan River" OR "Jordan basin" OR "Jordan  
512 water") are added to basic search terms to limit data volume to an acceptable extent. While conducting trial-and-error processes,  
513 topics of irrelevant articles are also recorded to show the potential causes of irrelevance and may provide us some hints to  
514 modify the search terms for a better performance, shown in Table 8.

515 **Table 7.** Trial-and-error process of frequency settings of keywords for Mekong, Nile and Jordan

Test index	Frequency settings				Mekong			Nile			Jordan		
	Basin name	Riparian country	Water	River	Data volume	Indicator 1	Indicator 2	Data volume	Indicator 1	Indicator 2	Data volume	Indicator 1	Indicator 2
1	1	1	1	1	27975	5	2	16227	8	4	28604	3	0
2	3	1	1	1	7536	13	10	6707	17	15	14028	4	1
3	5	1	1	1	4036	16	15	4157	19	16	13284	7	3
4	5	2	1	1	3695	16	16	4124	20	16	13263	7	1
5	5	2	2	2	3316	18	17	4017	20	18	5267	5	3
6	5	2	3	3	3102	19	17	3912	20	18	3830	7	10

516 **Table 8.** Recordings of potential irrelevant topics for Mekong, Nile and Jordan

Mekong Test Index	Data volume	Indicator 1	Indicator 2
1	27979	Paper index: 16787-16806 16789: Coastal monument; 16790: Catfish; 16791: Missing American servicemen; 16792: Missing people; 16793: Business plan; 16795: Life-style; 16796: America navy river corps; 16798: Life-style; 16799: Grade nationale; 16800: Travel to Vietnam and Cambodia; 16801: Travel; 16802: Travel along the river; 16803: Travel; 16804: Riots in Thailand; 16805: Riots in Thailand;	Paper index: 22383-22402 22383: Family; 22384: Soldiers; 22385: Vietnam annexed Cambodia through reconciliation; 22387: Vietnam sentences followers after trial of Buddha; 22388: Culture; 22389: Judicial delays in four cases; 22390: Ban on swill feed; 22391: Combine harvester race; 22392: Relocate 7 million people to relieve pressure on overcrowded areas; 22393: The man who reformed the United States Navy; 22394: Books about Vietnam; 22395: China plans to establish a national park system; 22397: Image consulting; 22398: Songkran Festival of Thailand; 22399: Southeast Asian refugees; 22400: Bow movement; 22401: A dueling event in Denver; 22402: Cambodian adoptees;
2	7536	Paper index: 4521-4540 4522: Former Vietcong say fighting on the Mekong River; 4524: Elephant; 4528: Luang Prabang; 4533: Laotian culture and beauty; 4534: Mekong River Journey; 4537: Hero for Children's Rights; 4540: Travel to Cambodia;	Paper index: 6029-6048 6029: Crossing the Mekong; 6030-6031: Vietnam's rice exports; 6032: DNA catch; 6033: Vietnam cruise; 6041: Inland river cruise; 6043: National Geographic researcher Reno; 6044: Mekong River Tourism; 6046: Escape the molecular; 6048: Buddha;
3	4036	Paper index: 2422-2441	Paper index: 3229-3248

		2428: A cruise ship on the Mekong River; 2434: Mekong River travel; 2440: Illegal timber trade; 2441: Mekong River travel;	3231: Mekong Animals; 3235: A manhunt for missing American soldiers in Laos; 3236: First impressions of Cambodia and Vietnam; 3239: Travel to Vietnam and Cambodia; 3240: Illegal logs are cut down;
4	3695	Paper index: 2217-2236 2218: Travel to Vietnam and Cambodia; 2220: Travel to Laos; 2231: Pacific Command disaster response exercise in Vietnam; 2236: Vientiane, capital of Laos;	Paper index: 2956-2975 2962: Visit Southeast Asia; 2963: Cambodia is trying to save a rare Mekong river dolphin; 2968: A search for missing American soldiers in Laos; 2975: Travel to the Mekong;
5	3316	Paper index: 1990-2009 1997: Mekong Tourism; 2007: Mekong Prize winner	Paper index: 2653-2672 2662: Mekong River travels in Thailand, Cambodia and Vietnam; 2663: Travel to Thailand; 2671: Mekong navigation;
6	3102	Paper index: 1860-1879 1861: Mekong Adventure;	Paper index: 2482-2503 2483: Laos arrests American manhunt for missing man; 2486: Roaming along the Mekong river; 2488: Drifting;

517

Nile Test Index	Data volume	Indicator 1	Indicator 2
1	16227	Paper index: 9736-9755	Paper index: 12982-13001
		9736-9738 : Rebels in South Sudan ; 9739-9740 : Archaeologist ; 9741 : Israel may withdraw from the West Bank; 9744: Slavery in ancient Egypt; 9745: Kurdish rebels in Turkey; 9748: Travel to Egypt; 9750-9752: Egypt's interior minister refused to allow "militias" to enter;	12982: Farmers and the Egyptian government fought for years in a legal battle; 12983: Curseja Island; 12984: Detectives use modern science to solve a 3,300-year-old murder mystery; 12985: Violence in Egypt; 12987: Tagore festival in Egypt; 12988: Bossi language learning courses ; 12989-12992 : Russian plane crash ; 12993 : Egypt's population grows ; 12995: Reviving Egypt's tourism industry; 12996-12997: Anti-government demonstrations in Sudan; 12998: Tourism landscape; 13001: Egyptian court holds second mass trial;
2	6707	Paper index: 4024-4043 4037 : Lake Victoria renamed ; 4038 : Egyptian journalist Resigns; 4041: Travel along the Nile;	Paper index: 5366-5385 5371-5373: Luxor Temple; 5378: Egyptian history; 5380: The uprising in Egypt is resurgent;
		Paper index: 2494-2513 2509: Changes along the Nile;	Paper index: 3326-3345 3326: Conflict in South Sudan ; 3328 : Travel along the Nile ; 3333 : Sudanese refugees; 3343 : Sudan Peace Conference;
4	4124	Paper index: 2474-2493 No irrelevant articles	Paper index: 3299-3318 3301: Egypt will withdraw 3.4 million from federal reserves to meet food needs; 3302: Nile culture; 3305: Sudan earthquake sequence 1990-1991 and the extent of the East African Rift Valley system; 3317: Ethiopia: Lakeside cities overcome Africa's tourism crisis;
		Paper index: 2410-2429 No irrelevant articles	Paper index: 3214-3233 3216: Displaced South Sudanese; 3232: Sudan earthquake sequence 1990-1991 and the extent of the East African Rift Valley system;
6	3912	Paper index: 2347-2366 No irrelevant articles	Paper index: 3130-3149 3134: An English man crosses the Nile on foot; 3136: Displaced South Sudanese;

518

Jordan Test Index	Data volume	Indicator 1	Indicator 2
1	28604	Paper index: 17162-17181 17162: Gaza's prison; 17164: Israel security Separation Wall; 17165: Jesus through Anne Rice's eyes: A book review; 17166: The king of Morocco visited the United States; 17167: Jewish terrorist group; 17168: Palestinians demonstrated in Israeli-occupied territory; 17169: Jordan's king urged the Palestine Liberation Organization to recognize Israel; 17170-17172: United States: Israeli-Palestinian peace agreement; 17173: Riding; 17174: The PLO will not meet with American officials; 17175: There has been violence in Jerusalem; 17176: Israel's democracy and Arab population; 17177: Women go to Palestine to resolve violence; 17178: Music; 17181: Hotel prices in Australia have fallen along with Asian growth;	Paper index: 22883-22902 22883: Joshua Myron, Zionist who fought the Turks, died; 22884: Did Netanyahu explain why the Palestinians did not reach a deal; 22885: Ottawa ordered compensation for disabled first Nations children; 22886: Sacramento State University student arrested in terrorist ring; 22887: South Africa: Two white-owned farms to be confiscated in land reform; 22888: State Department envoy meets with Palestinian Christians who oppose Israel; 22889: Three-year-old girl shot dead in Gaza; 22890: Mormon Temple renovation; 22891: Yehsat hosted a forum on the humanitarian use of satellite broadband; 22892: Humanitarian work; 22893: Jenkins' death; 22894: Interfaith activity; 22895: The ultimate consultant; 22896: Exhibition in the West Bank; 22897: The Jewish population in the occupied West Bank is set to more than double this year; 22898: Police have arrested a suspect in the Maverick shooting; 22899: Disturbing violence; 22900: "The etymology of the names Israel and Jacob; 22901: Terrorism; 22902: The city of Midville's first citywide master plan for trails;
		Paper index: 8417-8436 8417: Missionary trip to the West Bank; 8418: Against Islamic State militants; 8419: Israel imposed sanctions on Gaza; 8420: Arafat; 8421: Israeli withdrawal; 8422: The Israeli military has questioned an Arab mayor in the West Bank; 8423: It is widely believed in Israel that the current situation in Judea, Samaria and Gaza cannot and should not continue; 8424: Arafat rose up; 8425: Palestinians say Washington	Paper index: 11222-11241 11222: Manitoba: The Assembly of First Nations supports the Declaration of the Manitoba Chiefs; 11223: The Sheikh Hussein Bridge across the Jordan River was completed; 11224: The Israelis shot at two Jordanians; 11225: The sick Menachem Begison; 11226: Top election official supports south Jordan petition; 11227: How did a high school student in Nuremberg talk about crossing Israel; 11229-11230: Protesters in Amman burn an Israeli flag after the judge's killing; 11231: SCR 591 recognizes Sao Paulo's historic landmarks and museums; 11232: Expand light rail public transport;

		accepts their approach to ending attacks on Israel; 8426: The Likud trounced Sharon; 8429: Pope endorses' Palestinian Aspirations'; 8430: The Israeli authorities have jailed three senior Palestinian leaders without trial; 8433: Defense of the Jewish State; 8434: Australia's relationship with Israel; 8435: Former commander of the Arab Legion Grubb Pasha has died; 8436: A leadership void is holding Egypt back;	11233: Update from AFPTV on Tuesday; 11234-11236: Palestinians in the West Bank are under increasing economic pressure; 11237-11239: Sharon army; 11240: Silt diversion walls in East Jordan; 11241: Traveler's cheque;
3	13284	Paper index: 7970-7989	Paper index: 10627-10646
		7970: Israel faces a demographic threat; 7971: Israel and Palestine live side by side in peace; 7972: Peace negotiation; 7973: Palestinian Elections postponed; 7974: The countryside camping; 7975: A week news; 7976: The American president meets with Jordan's king; 7977: Watch the Pope's Middle East pilgrimage online; 7979: A secret meeting of Arab and Israeli writers; 7983: Top story on Tuesday; 7984-7985: Better support for Aboriginal children; 7986: Israel Archives; 7963: Catholics and Muslims seek dialogue;	10627-10629: Individual account; 10630: Witness: The Jordanian defendant had ties to Osama bin Laden; 10631: The wounded mayor vowed to continue the fight for Palestinian rights; 10632: Former U.S. PRESIDENT: Middle Eastern leaders must tell their people that compromise is honorable; 10633: Jordan baptism site sells bottled holy water tender; 10637: A letter from Israel; 10638: Public money spent on Park Avenue; 10639: The Israeli prime minister has proposed the creation of an independent Palestinian state; 10640: The European Commission has issued a final warning to The UK over repeated violations; 10641: Vote split; 10642-10643: Jordan's parliament failed to overthrow the government; 10644: Jordan Valley Trail; 10645: Possible hazards caused by pumping water near rivers; 10646: Task biography;
4	13263	Paper index: 7958-7977	Paper index: 10610-10629
		7958-7959: Terrorism; 7960: The new chief rabbi is a firebrand nationalist; 7962: Silt diversion wall; 7963: Catholics and Muslims seek dialogue; 7964: The Palestinian government's plan; 7965: Palestinian refugees; 7966: Jordan bridge; 7968: The United States is pushing for a Middle East peace plan; 7969: The Arabs conspired to blockade Israel; 7970: Jordan River Cycle Path; 7971: Mr. Netanyahu linked peace to Palestinian recognition of Israel as a Jewish state; 7976: Jordan refugee woman craftsman;	10610-10611: Jordanian-british student volunteers seek positive change in Western and Arab societies; 10612: Jerusalem Liberation Army; 10613-10614: Palestinian textbooks versus Israeli textbooks; 10615: Some rare right whales like winter in Maine; 10616: The PLO is preparing to move to Gaza and Jericho; 10617-10618: Arab Bank employees volunteer in Aguilon; 10619: East Jordan Commissioner Thomas Breney asked the board to vote on hiring a full-time fire chief; 10621: The FBI is demanding payment for the local youth's treatment; 10622: The federal government is seeking to appeal the ruling on medical costs; 10623-10624: Negotiations between Israel and Egypt; 10625: The Conservative Party is appealing against the treatment ruling; 10626: Fatah's Al-Aqsa Brigades killed a woman in Nablus accused of collaborating; 10627: Mr Hotovili bemoans Likud's "schizophrenia" over two countries; 10628: Kiryat Arba population; 10629: In the Likud debate, the two-state solution is schizophrenic;
5	5267	Paper index: 3160-3179	Paper index: 4214-4233
		3160: Mr Netanyahu does not fear being blamed if the London meeting is inconclusive; 3161: Forty-eight hours in Amman; 3162: Israel has begun clearing land mines at the site of Jesus' baptism in the West Bank; 3164: Covenant of Israel; 3166: Israel's efforts to win Over Christian tourists; 3168: Christian sites are covered in landmines; 3170: Community Notes - Volunteer; 3171: Edit the letter in the pouch; 3173: Travel; 3174: Mormons and non-Mormons; 3175: Immigrants find peace and opportunity in Corona; 3176-3177: Hymns to Haaznu on a biblical urn; 3178: A joyous gathering of prominent Israeli and PLO officials in the three years since the Signing of the Oslo Accords; 3179: The new chief rabbi is a firebrand nationalist;	4214: Jordan sewage pump; 4215: Across the border; 4216: In memory of a distinguished journalist; 4217: Travel manuscript; 4218: Some news; 4219: Joint Technology Center; 4220: leprosy; 4221: Israeli soldiers pass through barbed wire in Wazzani; 4223: Storm hits northern Michigan; 4224: Jon and Martha Jensen of Petoskey gave birth to daughter Ruby Susan at Northern Michigan Hospital; 4225: Peter and John stood before the Sanhedrin; 4228: Obituary; 4229: The UAE supports Jordan in implementing its development plans; 4230: Russia and the United States are set to reach a new agreement by next summer on deep cuts in strategic offensive weapons; 4231-4232: Hussein riots; 4233: The East Jordan Chamber of Commerce distributes community awards;
6	3830	Paper index: 2298-2317	Paper index: 3064-3083
		2298: The establishment of a provisional Palestinian state; 2299: Ways to promote Jordanian and British military cooperation; 2301: Israeli Prime Minister Benjamin Netanyahu has said he will see the final results of a peace deal; 2303: Epiphany; 2306: Winners of the first Tourism Promotion Peace Prize have been announced; 2308-2309: Better support for Aboriginal children; 2310: Jordan has "ridiculed" Israeli ministers' efforts to block Palestinian statehood; 2311: Pope Francis is making a visit to the Holy Land; 2314: The Jordan River has long been a source of entertainment for wasatch central city dwellers; 2315: Community news; 2316: South Jordan celebrates its 150th anniversary; 2317: Jordan Trail;	3064: American troops were sent from California to Utah during the Civil War; 3065: A Seattle moving company is offering spring deals; 3066-3067: Seattle Moving Company; 3068: Opponents of sewage treatment plants; 3069: Soldiers swept away by the Jordan River; 3073: Civil servants are considering a one-day strike on Monday; 3074: Pilgrims mark baptism traditions in the Jordan River; 3075: Word games; 3076: People's Fund grant project;