

# A System for Worldwide COVID-19 Information Aggregation

Akiko Aizawa<sup>3</sup>, Frederic Bergeron<sup>1</sup>, Junjie Chen<sup>5</sup>, Fei Cheng<sup>1</sup>, Katsuhiko Hayashi<sup>5</sup>, Kentaro Inui<sup>4</sup>, Hiroyoshi Ito<sup>6</sup>,  
Daisuke Kawahara<sup>7</sup>, Masaru Kitsuregawa<sup>3</sup>, Hirokazu Kiyomaru<sup>1</sup>, Masaki Kobayashi<sup>6</sup>, Takashi Kodama<sup>1</sup>,  
Sadao Kurohashi<sup>1</sup>, Qianying Liu<sup>1</sup>, Masaki Matsubara<sup>6</sup>, Yusuke Miyao<sup>5</sup>, Atsuyuki Morishim<sup>6</sup>, Yugo Murawaki<sup>1</sup>,  
Kazumasa Omura<sup>1</sup>, Haiyue Song<sup>1</sup>, Eiichiro Sumita<sup>2</sup>, Shinji Suzuki<sup>8</sup>, Ribeka Tanaka<sup>1</sup>, Yu Tanaka<sup>1</sup>,  
Masashi Toyoda<sup>8</sup>, Nobuhiro Ueda<sup>1</sup>, Honai Ueoka<sup>1</sup>, Masao Utiyama<sup>2</sup>, Ying Zhong<sup>6</sup> (in alphabetical order)

<sup>1</sup>Kyoto University <sup>2</sup>NICT <sup>3</sup>NII <sup>4</sup>Tohoku University <sup>5</sup>The University of Tokyo <sup>6</sup>The University of Tsukuba  
<sup>7</sup>Waseda University <sup>8</sup>Institute of Industrial Science, the University of Tokyo

## The System

- People pay attention to COVID-19 news of various topics.
- The COVID-19 condition is very different among the countries.
- Getting first-hand information from other countries is essential.

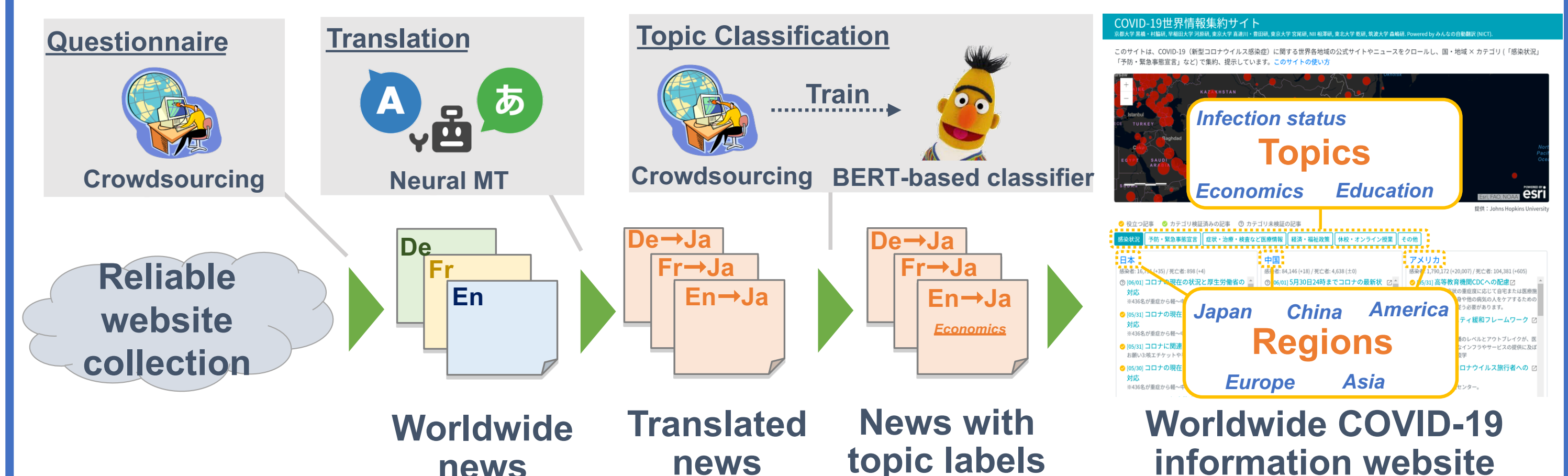


A System for Worldwide COVID-19 Information Aggregation with Various Topics.

<http://lotus.kuee.kyoto-u.ac.jp/NLPforCOVID-19>

## Challenges & Solutions

- The reliability of news sources.
- Translation quality to the local language.
- Topic classification for efficient searching.



- Robust multilingual reliable website collection via crowdsourcing.
- A high-quality machine translation system.
- A BERT-based topic classifier
- And a user-friendly web interface.

## Reliable Website Collection

To avoid rumors and high-quality, reliable information, we use multiple crowdsourcing services and limit the workers' nationality, assuming that local citizens know the reliable websites in their country.

Website	Country	Primary	Reason	Topics
www.cdc.gov	US	True	The site is a government website, specifically the Center for Disease Control.	infection status prevention and emergency declaration symptoms, medical treatment and tests
www.covid19-yanmar.com	Japan	False	Shinya Yamanaka is a famous medical researcher and his insights about COVID-19 are reliable.	prevention and emergency declaration
www.internazionale.it	Italy	False	This website collects and translates articles from news agencies and magazines from all over the world. Has up-to-date news, but also long-form analysis articles. Most of my deeper information comes from here.	infection status economics and welfare prevention and emergency declaration school and online classes
covid.saude.gov.br	Brazil	True	This site is the government web site.	infection status

Crowdworkers give trusted websites with reasons to choose it and what kind of information they can obtain from it.

- Totally 908 questionnaire results from 8 countries with totally 550 websites.
- Rumors are rampant in this era. The reliable websites dataset can help people to protect themselves.

Country	#	Questionnaire	Reliable sites
India	122	67	
US	106	77	
Italy	104	68	
Japan	102	49	
Spain	126	90	
France	127	71	
Germany	106	61	
Brazil	115	67	
Total	908	550	

Statistics of the number of questionnaires and reliable websites collected from each country.

## Crawl, Filter and Translation for Information Localization

- Crawl COVID-19 related articles from 35 reliable websites.
- Keyword filtering for most related articles.
- Multilingual machine translation to Japanese by Tex-Tra\*.

Country	Website	Mentioned times
United States	www.cdc.gov/coronavirus/2019-ncov	14
	www.usa.gov/coronavirus	6
	www.nytimes.com/news-event/coronavirus	4
Japan	hazard.yahoo.co.jp/article/20200207	17
	www.mhlw.go.jp/...	13
	corona.feedal.com	6
Italy	www.salute.gov.it/nuovocoronavirus	11
	www.salute.gov.it/portale/home.html	4
	www.worldometers.info/coronavirus	3
France	www.gouvernement.fr/info-coronavirus	28
	www.who.int/fr/emergencies/diseases/novel-coronavirus-2019	7
	www.lemonde.fr/coronavirus-2019-ncov/	6
Spain	www.usa.gov/coronavirus	9
	www.msbs.gob.es/profesionales/...	7
	covid19.gob.es	4
Germany	www.rki.de/DE/Home/homepage_node.html	7
	www.bundesgesundheitsministerium.de/coronavirus.html	6
	interaktiv.morgenpost.de/corona-virus-karte-infektionen-deutschland-weltweit	5

Top 3 most mentioned reliable websites of each country

\*<https://mt-auto-minhon-mlt.ucr.i.jgn-x.jp>

## Topic-classification

Please check the URL of the page first:

Page | [Check here for the page](#)

Is there information about COVID-19 in this page?

☐ Yes ☐ No

Is there helpful information in this page?

☐ Yes ☐ No

Is the Japanese in this page fluent?

☐ Yes ☐ No

What's the topics in this page? Choose correct ones you think from options below.

☐ Infection status ☐ Prevention

☐ Self-restraint ☐ Medical information

☐ Economic ☐ Online lesson

☐ Sports ☐ Articles about hoax

☐ Others

Country	Article with topic label
France	72K
America	11K
Japan	9K
China	10K
International	13K
Spain	1K
India	4K
Germany	2K
Total	122K

Article-topics dataset of different languages

Task	Keyword-based model			BERT-based model		
	Precision	Recall	F-score	Precision	Recall	F-score
Is about COVID-19	0.36	1.00	0.54	0.82	0.87	<b>0.84</b>
Topic: Infection status	0.09	0.53	0.16	0.43	0.81	<b>0.56</b>
Topic: Prevention	0.05	0.73	0.10	0.19	0.73	<b>0.30</b>
Topic: Medical information	0.17	0.70	0.27	0.27	0.91	<b>0.41</b>
Topic: Economic	0.06	0.36	0.10	0.14	0.84	<b>0.24</b>
Topic: Education	0.06	1.00	<b>0.11</b>	0.05	0.60	0.09
Topic: Art and Sport	0.06	0.41	0.10	0.08	0.94	<b>0.14</b>
Topic: Others	0.52	0.07	0.13	0.87	0.79	<b>0.83</b>

Topic-classification result of BERT-based model and keyword-based model.

- Totally 122K article-topics pairs through crowdsourcing.
- BERT-based model outperforms the keyword-based model in most tasks.
- BERT-based topic-classifier is then applied to generate topic labels for other articles.

Article-topics annotation through crowdsourcing

## System database

Country	Raw(↑/day)	Translated	With topics
France	774K(8K)	74K	9K
US	69K(730)	15K	2K
Japan	25K(260)	5K	2K
Europe	50K(510)	2K	50
China	38K(400)	3K	342
Int.	45K(470)	3K	263
Korea	16K(170)	260	71
Spain	4K(40)	370	36
India	14K(150)	860	66
Germany	16K(170)	8K	6K
Total	1.05M(11K)	110K	18K

Crawled, translated and labeled articles.

- Totally 1.05M pages with 110K of them translated and 18K with topic labels. Still growing.

