



European  
Commission

# JRC SCIENTIFIC INFORMATION SYSTEMS AND DATABASES REPORT

## Europe Media Monitor

*Metadata format V. 1.0*

*JRC.I.3*

2020



### Top Stories

UPDATED EVERY 10 MINUTES, 24 HOURS PER DAY.

Search

#### Main Menu

- Top Stories
- 24 Hours Overview
- Events Detection
- Most Active Themes
- Help about EMM
- Overview
- Advanced Search
- Sources list
- Web Site Map

#### EU Focus

#### EU Policy Areas

- Agriculture and Rural Development



#### Tools

Monday, November 23, 2015  
6:34:00 AM CET

RSS | KML | MAP

Facebook

subscribe | manage

info



#### Languages

Select top stories in other languages.

This publication is a Scientific Information Systems and Databases report by the Joint Research Centre (JRC), the European Commission's science and knowledge service. It aims to provide evidence-based scientific support to the European policymaking process. The scientific output expressed does not imply a policy position of the European Commission. Neither the European Commission nor any person acting on behalf of the Commission is responsible for the use that might be made of this publication. For information on the methodology and quality underlying the data used in this publication for which the source is neither Eurostat nor other Commission services, users should contact the referenced source. The designations employed and the presentation of material on the maps do not imply the expression of any opinion whatsoever on the part of the European Union concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries.

Contact information

Name: Marco Verile

Email: marco.verile@ec.europa.eu

EU Science Hub

<https://ec.europa.eu/jrc>

Ispira: European Commission, 2020

© European Union, 2020



The reuse policy of the European Commission is implemented by the Commission Decision 2011/833/EU of 12 December 2011 on the reuse of Commission documents (OJ L 330, 14.12.2011, p. 39). Except otherwise noted, the reuse of this document is authorised under the Creative Commons Attribution 4.0 International (CC BY 4.0) licence (<https://creativecommons.org/licenses/by/4.0/>). This means that reuse is allowed provided appropriate credit is given and any changes are indicated. For any use or reproduction of photos or other material that is not owned by the EU, permission must be sought directly from the copyright holders.

All content © European Union, 2020.

**Contents**

- Acknowledgements .....2
- 1. Introduction.....3
  - 1.1 Copyright Notice .....3
- 2. Format .....4
  - 1.2 Channel element .....4
  - 1.3 Item element .....5
- 3. EMM Named Entity Hierarchy .....9

## **Acknowledgements**

The Europe Media Monitor technology is the result of a collective passionate effort over more than 15 years. Many people have contributed to address text mining challenges as well as source curation, data collection, processing, and visualization tools design. We would like to thank all of them, and the many more that, we hope, will contribute to EMM in the future.

## 1. Introduction

EMM is the Europe Media Monitor<sup>1</sup>, a system for monitoring open source news information. EMM is developed and maintained by the Text & Data Mining Unit, in the Directorate for Competences of the European Commission's Joint Research Centre (JRC). EMM started in 2002 as a project to support the Commission with its media monitoring activities. The main purpose of EMM is to provide monitoring of a large (but selected) set of electronic media, reducing the information flow to manageable proportions by: clustering related news, categorising articles, and applying language technology tools to derive further metadata, such as recognising and disambiguating entities in the text, extracting quotes by and about people, applying sentiment/tonality analysis and more. The system continuously monitors almost 11 000 sources to find new articles published on the Internet (~300 000 articles daily).

The results of EMM processing are captured in metadata describing, for each analysed document, the associated topics, mentioned named entities (people, organisations, and locations), quotes extracted, etc. This document describes the format of the EMM metadata.

### 1.1 Copyright Notice

(c) European Union, 1995-2020

The Commission's reuse policy is implemented by the Commission Decision of 12 December 2011 on the reuse of Commission documents [1]. Any copyright and/or sui generis right on the dataset is licensed under the Creative Commons Attribution 4.0 International (CC BY 4.0) licence [2]. Reuse is allowed provided appropriate credit is given and any changes are indicated.

For any use of the titles (headings) of the articles, permission may need to be sought directly from the respective newspaper publishers. Titles may be protected by copyright and were included here only for the purpose of identification of the articles and linking with the related data.

[1] <https://eur-lex.europa.eu/eli/dec/2011/833/oj>

[2] <https://creativecommons.org/licenses/by/4.0>

---

<sup>1</sup> Public web site: <https://emm.newsbrief.eu> .

## 2. Format

EMM metadata is stored in RSS format<sup>2</sup> with some customisations. Each file contains a variable number of RSS items.

```
<?xml version="1.0" encoding="UTF-8"?>
<rss version="2.0" xmlns:emm="http://emm.jrc.it" xmlns:iso="http://www.iso.org/3166">
  <channel>

    <title>my selection</title>
    <pubDate>Tue, 21 Apr 2020 02:38:27 CEST</pubDate>

    <item emm:id="news-9147bf1262c6d838918051ac68da64e5">
      [...]
    </item>

    [... more items here...]

    <item emm:id="news-AFFF4556678c6d838918051ac68da6342">
      [...]
    </item>

  </channel>
</rss>
```

### 2.2 Channel element

The channel element contains the RSS creation date and title.

*Table 1. "channel" element.*

Field	Short Description	Example
pubDate	RSS file creation date.	Tue, 21 Apr 2020 02:38:27 CEST
Title	Title of the selection.	my selection

<sup>2</sup> See RSS 2.0 specification: <http://www.rssboard.org/rss-specification> .

## 2.3 Item element

The item element contains the metadata about one document processed by EMM.

Table 2. "item" element.

Attribute	Short Description	Example
<b>link</b>	Link to the news article on the publisher web site (the same used in the web browser).	<a href="https://www.lemonde.fr/planete/article/2020/04/07/coronavirus-les-francais-appelles-a-ne-pas-relacher-leurs-efforts-de-confinement_6035777_3244.html">https://www.lemonde.fr/planete/article/2020/04/07/coronavirus-les-francais-appelles-a-ne-pas-relacher-leurs-efforts-de-confinement_6035777_3244.html</a>
<b>Title</b>	Title of the news article <sup>3</sup> .	Coronavirus : plus de 10 000 morts en France, un salarié sur quatre en activité partielle
<b>Guid</b>	EMM Internal unique identifier.	middleeastmonitor-93b4e24dd35a4d910ce30dff47c7c898
<b>source</b>	a) Publisher's name b) Publisher's web site c) Publisher's country	a) LeMonde b) url= <a href="https://www.lemonde.fr/rss/une.xml">https://www.lemonde.fr/rss/une.xml</a> c) country="FR"
<b>pubdate</b>	Date and time when the article was discovered.	2020-04-07T21:15+0200
<b>iso:language</b>	Language of the text.	fr
<b>emm:entity</b>	Known people or organizations mentioned in the text (already part of the entity repository).  a) EMM Internal unique identifier b) Type c) Subtype (refer to chapter 3 for possible values)	  a) id="2844" b) type="p"

<sup>3</sup> See §1.1 for copyright restrictions.

Attribute	Short Description	Example
	<ul style="list-style-type: none"> <li>d) Occurrences in the text</li> <li>e) Position in the text</li> <li>f) Display name</li> </ul>	<ul style="list-style-type: none"> <li>c) subtype="PER"</li> <li>d) count="2"</li> <li>e) pos="3684,4540"</li> <li>f) name="Anne Hidalgo"</li> </ul>
<b>emm:guess</b>	<p>New people or organization names discovered in the text (not yet in the entity repository).</p> <ul style="list-style-type: none"> <li>a) Type</li> <li>b) Subtype (refer to chapter 3 for possible values)</li> <li>c) Occurrences in the text</li> <li>d) Position in the text</li> <li>e) Display name</li> <li>f) Matching rule</li> </ul>	<ul style="list-style-type: none"> <li>a) type="o"</li> <li>b) subtype="ORG-CO"</li> <li>c) count="1"</li> <li>d) pos="1539"</li> <li>e) name="BFMTV"</li> <li>f) rules="ho0_alone_0"</li> </ul>
<b>emm:georss</b>	<p>Known location mentioned in the text.</p> <ul style="list-style-type: none"> <li>a) Full name</li> <li>b) Internal unique identifier</li> <li>c) Latitude</li> <li>d) Longitude</li> <li>e) Occurrences in the text</li> <li>f) Position in the text</li> <li>g) Class</li> <li>h) ISO country code</li> <li>i) Position in the text</li> <li>j) Word length</li> <li>k) Display name</li> </ul>	<ul style="list-style-type: none"> <li>a) name="Paris:Ile-de-France:France"</li> <li>b) id="16912462"</li> <li>c) lat="48.8521"</li> <li>d) lon="2.34899"</li> <li>e) count="3"</li> <li>f) pos="3460,3678,4566"</li> <li>g) class="1"</li> <li>h) iso="FR"</li> <li>i) charpos="3460,3678,4566"</li> <li>j) wordlen="5,5,5"</li> <li>k) Paris</li> </ul>
<b>emm:fullgeo</b>	<p>Known location mentioned in the text.</p> <ul style="list-style-type: none"> <li>a) Name</li> <li>b) Internal unique identifier</li> </ul>	<ul style="list-style-type: none"> <li>a) name="France"</li> <li>b) id="78"</li> </ul>

Attribute	Short Description	Example
	c) Adjective d) Internal ranking e) Latitude f) Longitude g) Occurrences in the text h) Position in the text i) Class j) ISO country code k) Position in the text l) Word length m) Display name	c) adjective="false" d) rank="5" e) lat="48.8521" f) lon="2.34899" g) count="5" h) pos="38,2344,2477,2944,4126" i) class="-1" j) iso="FR" k) charpos="38,2344,2477,2944,4126" l) wordlen="6,6,6,6,6" m) France
<b>category</b>	Topic/subject mentioned in the text.  a) Name. b) Internal ranking (position in the text). c) Keywords matching score. d) Relevant words found in the text.	a) CommunicableDiseases b) emm:rank="0" c) emm:score="80" d) emm:trigger="Coronavirus[1]; coronavirus[1];"
<b>emm:link</b>	Link to other web pages referenced in the text.	<a href="https://www.interieur.gouv.fr/Actualites/L-actu-du-Ministere/Attestation-de-deplacement-derogatoire-et-justificatif-de-deplacement-professionnel">https://www.interieur.gouv.fr/Actualites/L-actu-du-Ministere/Attestation-de-deplacement-derogatoire-et-justificatif-de-deplacement-professionnel</a>
<b>emm:sentiment</b>	Sentiment of the overall text.	negative
<b>emm:emotion</b>	Emotion of the overall text.	anger
<b>emm:tonality</b>	Tonality score of the overall text.	-4
<b>emm:timex</b>	Numerical expression. a) Type b) Position in the text c) Occurrences in the text d) Exact value extracted.	a) type="date" b) pos="332" c) count="1" d) value="2020-04-07"

Attribute	Short Description	Example
<b>emm:ifs</b>	<p>Other entities, usually numerical or composed of specific textual patterns like email addresses, currency, ...</p> <p>a) Type (in this case always with the value "x")  b) Subtype (refer to 3for possible values)  c) Position in the text  d) Occurrences in the text  e) Exact value extracted.</p>	<p>a) type="date"  b) subtype=" IDT-EM"  c) pos="332"  d) count="1"  e) value="JRC-EMM-SUPPORT@ec.europa.eu"</p>
<b>emm:text</b>	<p>50% of the processed words sorted in alphabetic order.</p>	<p>achievable ; adjustment ; be ; border ; cent ; convinced ; corporate ; cut ; donald ; experts ; funded ; have ; intention ; irish ; may ; might ; not ; now ; off ; per ; president ; rate ; reduction ; reform ; remain ; table ; tax ; tax ; tax ; tax ; that ; the ; the ; the ; the ; to ; to ; trump ; us ; us ; us ; which ; with</p>

### 3. EMM Named Entity Hierarchy

The **main rationale for creating this Named-Entity hierarchy** is to enumerate the **Named-Entity TYPES** to be used **CONSISTENTLY** when encoding any kind of linguistic resources (e.g., lexical resources or grammars) for the purpose of **Named Entity Recognition** and make easier the sharing of information between the different modules involved in this process. This doesn't mean that the description of an entity would be limited to these types and subtypes only. More fine-grained or application-specific information could be declared by adding arbitrary attribute-value pairs to each entity, like for instance a profession for a person entity.

Since names can be ambiguous and some entities could have multiple sub-types or even types, a given entity can be potentially assigned to **more than ONE type or subtype**.

PERSON		
Subtype	Example/Explanation	Encoding
-	<i>"John Smith", "George W. Bush", "James Bond"</i>	PER <sup>4</sup>
ORGANISATION		
Subtype	Example/Explanation	Encoding
POLITICAL-PUBLIC <sup>5</sup>	political parties, e.g., <i>"Democratic Party", "CDU"</i> , political organisations, e.g. <i>"Palestine Liberation Organisation", "ISIS"</i> , military organisations, e.g. <i>"US Air Force"</i> , government institutions, e.g., <i>"Ministry of Interior of Italy", "Thatcher's Cabinet", "Embassy of USA"</i> , public institution, e.g., <i>"European Commission", "European Patent Office", "New York Public Library"</i>	ORG-PP
COMMERCIAL	<i>"Toyota", "Apple", "Microsoft", "Bank of Scotland"</i>	ORG-CO
RELIGIOUS	<i>"Anglican Church of Canada", "Islamic Forum of Europe"</i>	ORG-RE
SPORT	<i>sport clubs and organisations, e.g., "FC Barcelona", "Serie A", "Bundesliga"</i> ,	ORG-SP
EDUCATION-RESEARCH	<i>"University of Lugano", "European School of Varese"</i>	ORG-ER
OTHER	any organisation that do not fit in the categories above	ORG-OT
LOCATION		
Subtype	Example/Explanation	Encoding
CITY	<i>"London"</i>	LOC-CI
COUNTY	<i>"West Chester County"</i>	LOC-CN
PROVINCE	<i>"Province of Varese"</i>	LOC-PR
COUNTRY	<i>"Italy"</i>	LOC-CT

<sup>4</sup> There was a proposal to introduce subtypes for PER, e.g., to distinguish between profession of the person and/or whether it is a historical person or from 20/21 century, etc. We concluded that this type of information could be stored in the appropriate attribute for PER mainly due to the fact that we are talking here of some characteristics that change over time (e.g. profession).

<sup>5</sup> There was a proposal to introduce subtypes for this category, namely: (a) parties, (b) administrative, (c) civil protection and disaster management services like fire brigade, (d) juridical - courts, procuratura, (e) political movements, recognized and not recognized, (f) criminal - Ndrangeta, Mafia, Sacra corona unita, (g) terrorist organizations. Since the borders between some of these subtypes are to some extent blurred for now there will not be any distinction between them. However, additional information can be stored in attribute appropriate for that type.

REGION	part of a city, e.g., <i>“Bronx”</i> , special economic zone, geographical region, e.g., <i>“Provence”</i>	LOC-RE
FACILITY <sup>6</sup>	sport facilities, e.g., <i>“Yankee Stadium”</i> , recreation facilities, e.g., <i>“Central Park”</i> , <i>“Berlin Zoo”</i> , <i>“Disneyland”</i> , etc., cultural facilities, e.g., <i>“British Museum”</i> , <i>“Louvre”</i> , <i>“Royal Opera House”</i> , <i>“La Scala”</i> , hotels, tourist sites, e.g., <i>“Archeological Ruins at Moenjodaro”</i> , <i>“Forum Romanum”</i> , hospitals, cemeteries, all kind of transportation hubs (ports, railway stations, airports), e.g., <i>“Schipol Airport”</i> , <i>“165th Street Bus Terminal”</i> , <i>“Berlin Hauptbahnhof”</i> , <i>“Port of Honk Kong”</i> , <i>“Victoria Harbour”</i> , churches, urban/non-urban facilities such as roundabouts, railroads, roads, tunnels, e.g. <i>“St. Gothard Tunnel”</i> , bridges, etc.	LOC-FA
OTHER	any mentions of locations (e.g. landforms) that do not fit in the categories above, e.g., mountains <i>“Mount Everest”</i> , islands, water bodies <i>“Baltic”</i> , rivers, valleys, islands, etc.	LOC-OT
<b>IDENTIFIER</b>		
<b>Subtype</b>	<b>Example/Explanation</b>	<b>Encoding</b>
STREET-NAME	<i>“Via Fermi 27”</i>	IDT-SN
POSTAL-CODE	<i>“NY 10202”</i>	IDT-PC
POSTAL-ADDRESS	<i>“Via Fermi 27, 21200 Ispra, Italy”</i>	IDT-PA
GEO-COORDINATES	<i>“51° 28’ 38” N”</i>	IDT-CO
PHONE-NUMBER	<i>“+49 60 1234576”</i>	IDT-PN
EMAIL	<i>“abc@derf.com”</i>	IDT-EM
URL	<i>“http://www.google.com”</i>	IDT-UR
IP-ADDRESS	<i>“255.255.123.212”</i>	IDT-IP
VAT-NUMBER	<i>“ATU99999999”</i>	IDT-VA
BANKING-IDENTITY	<i>“DE44 5001 0517 5407 3249 31”</i>	IDT-BI
CREDIT-CARD-NUMBER	<i>2345 3523 2453 3453</i>	IDT-CR
SOCIAL-MEDIA-ID	<i>“@jakubP”</i>	IDT-SM
<b>PRODUCT</b>		
<b>Subtype</b>	<b>Example/Explanation</b>	<b>Encoding</b>
ELECTRONICS	<i>“Commodore 64”</i>	PRO-EL
DRUG-MEDICINE	<i>“Aspirin C”</i>	PRO-DM
WEAPON	<i>“AGM-1 Carbine”</i>	PRO-WE
VEHICLE	<i>“Mitsubishi Pajero”</i>	PRO-VE
FOOD	<i>“Snickers”</i>	PRO-FO
ART	<i>“Star Wars”</i>	PRO-AR
SERVICE	<i>“Google Search Engine”</i>	PRO-SE
OTHER	any product mention that does not fall under the above categories	PRO-OT
<b>EVENT<sup>7</sup></b>		
<b>Subtype</b>	<b>Example/Explanation</b>	<b>Encoding</b>

<sup>6</sup> Note that this category has been inspired by the FACILITY category in Sekine NE.

<sup>7</sup> Please note that classification of the events is done from the perspective of NAMED MENTIONS of events in text.

INCIDENT <sup>8</sup>	<i>"Chernobyl Disaster", "John Kennedy assassination", "World war II"</i>	EVT-IN
NATURAL	<i>"Great Alaska Earthquake", "Hurricane Katrina"</i>	EVT-NA
OCCASION <sup>9</sup>	conferences, e.g., <i>"LREC 2016", "Yalta Conference"</i> , religious holiday. e.g., <i>"Christmas"</i> , sport events, e.g., <i>"Football World Cup 2014"</i> , ceremonies, e.g. <i>"Nobel Prize Awards"</i> , etc.	EVT-OC
OTHER	<i>"Kuril Island dispute"</i>	EVT-OT
<b>TIMEX</b>		
<b>Subtype</b>	<b>Example/Explanation</b>	<b>Encoding</b>
TIME	<i>"2PM", "18:42"</i>	TIM-TM
DATE	<i>"1 April 2016", "12.01.2016"</i>	TIM-DA
PERIOD	<i>"4 hours", "20 years"</i>	TIM-PE
OTHER	<i>"Victorian Age"</i>	TIM-OT
<b>NUMEX</b>		
<b>Subtype</b>	<b>Example/Explanation</b>	<b>Encoding</b>
NUMERICAL-EXPRESSION	<i>"12", "12.00", "12 million", "22%", "2/3",</i>	NUM-EX
CURRENCY-EXPRESSION	<i>"100 USD", "3.50 EUR"</i>	NUM-CU
AGE	<i>"12 years old"</i>	NUM-AG
MEASUREMENT	<i>"30 kg", "100 gallons", "36° C"</i>	NUM-ME
COUNTX	<i>"10 people"</i>	MUM-CT
OTHER	any numerical expressions that do not fall under the above categories	NUM-OT
<b>OTHER</b>		
<b>Subtype</b>	<b>Example/Explanation</b>	<b>Encoding</b>
-	everything else that does not fall under any of the above main categories	OTH

<sup>8</sup> Man-made incidents

<sup>9</sup> There was a proposal to make distinction between the different subtypes, e.g., political, sport, conferences ... Such sub-classification would not be consistent with the current breakdown into the main categories, etc., e.g., conferences can be related to both politics and sports. Waiting for more input in this regard and a refined proposal.

## **GETTING IN TOUCH WITH THE EU**

### **In person**

All over the European Union there are hundreds of Europe Direct information centres. You can find the address of the centre nearest you at: [https://europa.eu/european-union/contact\\_en](https://europa.eu/european-union/contact_en)

### **On the phone or by email**

Europe Direct is a service that answers your questions about the European Union. You can contact this service:

- by freephone: 00 800 6 7 8 9 10 11 (certain operators may charge for these calls),
- at the following standard number: +32 22999696, or
- by electronic mail via: [https://europa.eu/european-union/contact\\_en](https://europa.eu/european-union/contact_en)

## **FINDING INFORMATION ABOUT THE EU**

### **Online**

Information about the European Union in all the official languages of the EU is available on the Europa website at: [https://europa.eu/european-union/index\\_en](https://europa.eu/european-union/index_en)

### **EU publications**

You can download or order free and priced EU publications from EU Bookshop at: <https://publications.europa.eu/en/publications>. Multiple copies of free publications may be obtained by contacting Europe Direct or your local information centre (see [https://europa.eu/european-union/contact\\_en](https://europa.eu/european-union/contact_en)).

## The European Commission's science and knowledge service

Joint Research Centre

### JRC Mission

As the science and knowledge service of the European Commission, the Joint Research Centre's mission is to support EU policies with independent evidence throughout the whole policy cycle.



**EU Science Hub**

[ec.europa.eu/jrc](https://ec.europa.eu/jrc)



@EU\_ScienceHub



EU Science Hub - Joint Research Centre



EU Science, Research and Innovation



EU Science Hub