

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 ¹
ORGANISATION		
Subtype	Example/Explanation	Encoding
POLITICAL-PUBLIC ²	political parties, e.g., <i>"Democratic Party", "CDU"</i> , political organisations, e.g. <i>"Palestine Liberation Organisation", "JSIS"</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

¹ 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).

² 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.

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
REGION	part of a city, e.g., <i>"Bronx"</i> , special economic zone, geographical region, e.g., <i>"Provence"</i>	LOC-RE
FACILITY ³	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
IDENTIFIER		
Subtype	Example/Explanation	Encoding
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
PRODUCT		
Subtype	Example/Explanation	Encoding
ELECTRONICS	<i>"Commodore 64"</i>	PRO-EL
DRUG-MEDICINE	<i>"Aspirin C"</i>	PRO-DM

³ Note that this category has been inspired by the FACILITY category in Sekine NE.

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
EVENT⁴		
Subtype	Example/Explanation	Encoding
INCIDENT ⁵	<i>"Chernobyl Disaster", "John Kennedy assassination", "World war II"</i>	EVT-IN
NATURAL	<i>"Great Alaska Earthquake", "Hurricane Katrina"</i>	EVT-NA
OCCASION ⁶	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
TIMEX		
Subtype	Example/Explanation	Encoding
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
NUMEX		
Subtype	Example/Explanation	Encoding
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
OTHER		
Subtype	Example/Explanation	Encoding
-	everything else that does not fall under any of the above main categories	OTH

⁴ Please note that classification of the events is done from the perspective of NAMED MENTIONS of events in text.

⁵ Man-made incidents

⁶ 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.