

MicroBloggingNERo Annotation Guide

Purpose	1
Introduction	1
Entity Types	2
Person (PER)	2
Organization (ORG)	3
Location (LOC)	4
Legal References (LEGAL)	5
Time Expressions (TIME)	5
Anatomical parts (ANAT)	6
Chemical and drugs (CHEM)	6
Disorders (DISO)	7
Medical devices (MED_DEVICE)	7
Microblogging Specifics	7
Detailed examples	7
References	8

Purpose

This short guide offers definitions and examples on how to annotate entities for the MicroBloggingNERo corpus. MicroBloggingNERo is a corpus composed of legal-domain documents in Romanian language manually annotated for named entities. The annotation process will follow the general guidelines described in this document. Annotators will identify text spans associated with one of the entity types introduced in this document. For clarification purposes, meetings will be held regularly with the annotators. As needed, additional examples may be added to this document following the meetings.

Introduction

In data mining, a named entity is a phrase that clearly identifies one item from a set of other items that have similar attributes. General examples of named entities are first and last names, geographic locations, ages, addresses, phone numbers, and

companies. Depending on the domain of interest, different entity types may be defined. For the purposes of the MicroBloggingNERo corpus, we are considering only the following entities: Person (PER), Organization (ORG), Location (LOC), Time expressions (TIME), Legal references (LEGAL), Anatomical parts (ANAT), Chemical and drugs (CHEM), Disorders (DISO), Medical devices (MED_DEVICE) . These will be detailed in the next section.

In certain cases, named entities may be imbricated one into another. For example, the organization “Primăria Municipiului Braşov” also contains a potential location entity “Braşov”. For our purposes, we are not considering such sub-entities. The annotation will consider only the largest text span denoting an entity. For this example, only the organization entity will be annotated. As an exception to this rule, we will consider sub-entities present in the legal reference entities. This will allow us to construct automatic named entity recognition (NER) systems able to recognize entities with or without the legal reference class. Therefore, for entities like “Legea 176/2019” annotations will contain both the entire span as a legal reference entity and the text span “2019” as a time expression.

Annotation will focus on identifying the largest text span clearly identifying the name of a named entity, without including additional information. Examples (the named entity spans are indicated between [...] and in bold): “dr. ing. [**Ion Popovici**]” (en. “Dr. Eng. Ion Popovici”), “oraşul [**Bucureşti**]” (en. “Bucharest city”), “competiţia se va desfăşura în aproximativ [**3 luni**]” (en. the competition will take place in about 3 months). General words that do not denote the name of a named entity or a clear time expression or a resolvable legal reference will not be annotated, such as “primarul” (en. “mayor”), “satul” (en. “village”), “comuna” (en. “commune”), “parcul” (en. “park”), “peste un timp” (en. “after a while”), “o lege” (en. “a law”), etc.

These annotation guidelines build on existing guidelines used in annotating other corpora. The named entities present in the MicroBloggingNERo corpus are also present in the LegalNERo corpus (Păiş et al., 2021a, 2021b) and in the SiMoNERo corpus (Barbu-Mititelu and Mitrofan, 2020). Therefore, entity spans will be annotated consistently in order to allow a comparison between the tools developed using these previous corpora. The LegalNERo guidelines were inspired in part by the Linguistic Data Consortium (LDC) guidelines for annotation of named entities¹.

Entity Types

Person (PER)

Person entities are regularly limited to humans. A person may be a single individual or a group. By extension, if a fictional character or a reference to religious figures is present it will be annotated. The annotation will consider only the largest text

¹ <https://www ldc upenn edu/sites/www ldc upenn edu/files/english-edt-v4.2.6.pdf>

span clearly identifying the person's name. Additional words that may be present, such as titles, honorifics, functions will not be annotated as part of the entity name.

Even though there are no entity subtypes, the following mentions are considered entities of type PERSON:

- a) Each distinct person or set of people mentioned in a document refers to an entity of type Person.
Eg.: Președintele Agenției Naționale de Administrare Fiscală, **[Mirela Călugăreanu]**. (en. President of the National Agency for Fiscal Administration, Mirela Călugăreanu.)
- b) Saints, angels and other biblical figures, including references to “God”, as opposed to religious leaders, are annotated together with the title or other reference: Pilat din Pont, Arhanghelul Mihail, etc.
Eg.: **[Sfântul Ion]** se sărbătorește la începutul anului. (en. Saint John is celebrated at the beginning of the year)
[Dumnezeu] a făcut cerul și pământul. (en. God made the heavens and the earth.)
- c) Fictional characters
Eg.: **[Spider-Man]** este un personaj popular printre tineri. (en. Spider-Man is a popular character among young people.)
- d) Family names
Eg.: **[Popeștii]** și-au cumpărat un teren lângă pădure. (en. Popeștii bought a plot of land near the forest.)

Are **not** considered entities of type PERSON: occupations “măcelarul” (en. “the butcher”), family relations “tatăl” (en. “dad”), pronouns “el” (en. “he”), titles and honorifics that precede the name, unmanned/general groups of people “familie” (en. “family”), “pictorii” (en. “painters”), etc..

Organization (ORG)

Organization entities are limited to corporations, agencies, and other groups of people defined by an established organizational structure. The annotation process will mark text spans clearly indicating the name of an organization. Additional words present in text will not be annotated if they are not part of the organization’s name. For example (annotated text spans are given in bold and between [...]): “compania **[SC Company SA]**” (en. “company SC Company SA”), “[**Compania Apa Brașov**]” (en. “Company Apa Brașov”). In the first case, the word “compania” is not part of the organization’s name and it was not annotated. In the second case, the same word is part of the second organization’s name and it was annotated. Since we are dealing with legal-domain text

(and not fictional text), if in doubt the complete organization name can be looked up on Google to make sure the proper text span is considered.

Even though there are no entity subtypes, the following mentions are considered entities of type ORGANIZATION:

- a) Government organizations relating to, or dealing with the structure or affairs of government, politics, or the state, political parties.
Eg.: Hotărârea a fost emisă de **[Guvernul României]**. (en. The decision was issued by the Romanian Government.)
[Camera de comerț a României] organizează cursuri acreditate. (en. The Romanian Chamber of Commerce organizes accredited courses.)
- b) Commercial organizations and nonprofit organizations
Eg.: **[Societatea Economică Eurasică (EAEC)]** s-a întrunit săptămâna trecută. (en. The Eurasian Economic Society (EAEC) met last week.)
[Crucea Roșie] a fost prezentă după izbucnirea incendiului. (en. The Red Cross was present after the fire broke out.)
- c) Educational organizations
[Universitatea București] organizează concursul de admitere la masterat pe 2 aprilie. (en. The University of Bucharest is organizing the master's admission competition on April 2.)
- d) Other types of organizations such as: political parties, UN, EU, etc.
[Biserica Ortodoxă Română] sprijină refugiații din Siria. (en. The Romanian Orthodox Church supports Syrian refugees.)

Are **not** considered entities of type ORGANIZATION mentions like: “angajați” (en. “employees”), “echipaj” (en. “crew”), etc.

Location (LOC)

Location entities are limited to geographical entities such as geographical areas and landmasses, bodies of water, and geological formations, denoted by a proper name. The annotation process will identify the name associated with a location entity, without additional words, such as the location type, unless these words are part of the official entity name. Examples (annotated text spans are given in bold and between [...]): “[**Parcul Alexandru Ioan Cuza**]” (en. “Alexandru Ioan Cuza Park”), “orașul [**București**]” (en. “Bucharest city”). When dealing with real places, if in doubt, the complete location name can be looked up on Google (or Google Maps) to make sure the proper text span is considered.

- a) Continent (**Africa**), Nation ("**Germania**", en."Germany"), County (**Argeş**), City (**Curtea de Arges**), district (**Sector 4**)
- b) political entity: "**Europa de Est**" (en. "Eastern Europe"), "**Orientul Mijlociu**" (en."the Middle East")
- c) geographical areas and landmasses: "**Transilvania**" (en. Transylvania")
- d) bodies of water (**Olt**), and geological formations "**Carpații**" (en. "Carpatians")

Are **not** considered entities of type LOCATION general mentions such as: "stradă" (en. "street"), "pod" (en. "bridge"), "oraş"(en. "city"), "aeroport" (en. "airport"), etc.

In certain cases, country names may be used in a context more appropriate for an organization. Nevertheless, we consider that an entity that can be clearly located on a map is a location entity. For example: "**Ungaria** va organiza procesul de votare" ("Hungary will organize the voting process"). In this case, "Ungaria" will be identified as a location.

Legal References (LEGAL)

Legal references are designations (the title of a legal document) or expressions pointing to a legal document.

Are considered LEGAL references mentions such as:

Law: "**Legea fondului funciar nr. 18/1991.**" (en.Land fund law no. 18/1991.)

Decree: "**Decretul nr. 309/2020** pentru promulgarea legii." (en. Decree no. 309/2020 for the promulgation of the law.)

Decision: "**Decizia nr. 892/2021** referitoare la respingerea excepției de neconstituționalitate." (en. Decision no. 892/2021 regarding the rejection of the exception of unconstitutionality.)

Convention: "**Convenția Consiliului Europei privind coproducția cinematografică (revizuită), din 19.01.2021.**" (en. Council of Europe Convention on Film Co - Production (Revised) of 19.01.2021.)

Declaration: "**Declarația Parlamentului României nr. 1/2020** cu privire la situația din Irak." (en. Declaration of the Romanian Parliament no. 1/2020 on the situation in Iraq.)

Are not considered entities of type LEGAL general mentions such as: "lege" (en. "law"), "decret" (en. "decree"), "declarație" (en. "declaration"), etc.

Time Expressions (TIME)

Time expressions tell us when something happened, or how long something lasted, or how often something occurs. Sometimes the precise date cannot be determined, allowing for expressions indicating periods of time (“next year”, “in the next three months”).

Examples:

- a) Complete or partial dates:
“3 ianuarie 2010” (en. “January 3, 2010”), “19.01.2020”, “martie 2018” (en. “March 2018”), “2017”.
- b) Periods of time:
“Martie-aprilie” (en. “march-april”), “3 luni” (en. “3 months”), “următorii 2 ani” (en. “The next 2 years”).
- c) Specific holidays indicating a time period
“Crăciun” (en. “Christmas”), “Anul nou” (en. “New year”), “Paști” (en. “Easter”).

Prepositions or other words before or after the time expression will not be included in the annotation, unless they change the meaning of the expression. For example: “**3 luni anterioare**” (“previous 3 months”), the word “anterioare” needs to be present in the annotation, because it is a clarification about the time expression. However, “pentru **3 luni**” (“for 3 months”) there is no need to include the word “pentru” since the preposition does not come with supplementary semantic information about the time entity.

Anatomical parts (ANAT)

Contains mentions of anatomical parts, parts of the human body, organs, components of organs, tissues, cells, cellular components.

Examples: valvă, inimă, picior, stomac, țesut epitelial.

Chemical and drugs (CHEM)

Contains mentions of amino acids, peptides, proteins, antibiotics, active substances, drugs, enzymes, hormones, receptors.

Examples: P5E-fe, beta-lactamine, penicilina, aminoglicozid, gentamicină.

Disorders (DISO)

Contains mentions of anatomical abnormalities, congenital anomalies, diseases, syndromes, lesions, symptoms.

Examples: complicații neurologice embolice, sângerari, T2DZ.

Medical devices (MED_DEVICE)

Contains mentions of any device intended to be used for medical purposes.

Examples: platură adeziv, lentile de contact, aparat cu raze X, stimulator cardiac, stetoscop, seringă, proteză de șold.

Microblogging Specifics

Micro-blogging messages may contain specific elements, such as emojis, emoticons, hashtags. Out of those, the hashtags may contain named entities (or parts of NEs). If a hashtag contains a named entity it will be annotated including the “#” sign. If the hashtag is part of a larger named entity, the entire text span (the hashtag and the words outside of the hashtag) will be annotated.

Examples: **Simona #Halep** (PER), **#IacDragan** (LOC), **#romania** (LOC).

Certain messages may end with “...” usually indicating a message that was too large and was truncated by the user’s application to fit in the platform’s imposed limits. In this case, if the message ends with a named entity (that may have been also truncated) the “...” will not be included in the text span. Furthermore, the entity will be annotated only if the annotator can make sense of what it means.

Detailed examples

Some of these examples were discussed during the annotation meetings. In some cases they may seem ambiguous but usually they are covered by carefully reading the rules expressed in this guide.

Text	Entities and explanation
Kremlinul a transmis un comunicat.	Kremlinul = LOC. This is location referring to the the Kremlin in Moscow, similar to how a country is always a location regardless of the context.

Kremlinul din Moscova este un complex fortificat.	Kremlinul din Moscova = LOC. This is a location entity since it can be located on a map. The entire text span is a single entity that clearly defines the location.
Medicamentul se ia dimineața pe stomacul gol.	dimineața = TIME. This is a time expression indicating each morning.
Podul peste Dunăre de la Brăila este în construcție.	Podul peste Dunăre de la Brăila = LOC. This is a location entity. This text span represents the name of the location. It can be clearly identified on a map. Any smaller portion of the text span will either make the entity ambiguous or refer to a different entity.
Un pod de la Brăila se construiește.	Brăila = LOC. In this context "un pod de la Brăila" would not provide a clear indication of a location (there can be multiple bridges in Brăila). Hence only Brăila is annotated.
Azi dimineață a fost frig.	Azi dimineață = TIME. This is a time expression, even though its actual interpretation depends on the time the document was written. However, it clearly defines the time period in the morning of the day the message was written.

References

Barbu Mititelu, V., and Mitrofan, M. (2020). The Romanian medical treebank-SiMoNERo. In Proceedings of the The 15th Edition of the International Conference on Linguistic Resources and Tools for Natural Language Processing – ConsILR-2020, ISSN 1843-911X, pages 7-16.

Păiș, V., Mitrofan, M., Gasan, C. L., Ianov, A., Ghiță, C., Coneschi, V. S., and Onut, , A. (2021a). Romanian Named Entity Recognition in the Legal domain (LegalNERo), <https://doi.org/10.5281/zenodo.4772094>

Păiș, V., Mitrofan, M., Gasan, C. L., Coneschi, V., and Ianov, A. (2021b). Named entity recognition in the Romanian legal domain. In Proceedings of the Natural Legal Language Processing Workshop 2021, pages 9–18, Punta Cana, Dominican Republic, November. Association for Computational Linguistics.