

# Annotation Guidelines

## Version 1

***Stephen Mayhew, November 2022, based on NorNe guidelines***

***(<https://github.com/ltgoslo/norne/blob/master/NorNe%20Annotation%20Guidelines.pdf>)***

## 1. What is an annotation?

An annotation is “a markup of a text span in a specific format that indicates a feature or features of the text within the span.” (MUC-7 ([https://www-nlpir.nist.gov/related\\_projects/muc/index.html](https://www-nlpir.nist.gov/related_projects/muc/index.html)))

In addition to a span, all annotations have a type: e.g. person, location, organization.

## 2. What should be annotated?

In the Universal Dependencies NER project we annotate named entities. All names of people (fictional and real), organizations, and locations that are found in the texts should be annotated.

The main rules of thumb for annotating something as a name are:

- The word is or contains a proper noun
  - ‘Michael Jordan’ is composed of two consecutive proper noun tokens with a single reference
  - Names may include words which are not proper nouns, e.g. prepositions, as in ‘John of Burgundy’
  - ‘Pearl Harbor’ is a name of a location, yet it consists of two common nouns, “pearl” and “harbor”
- The word has a unique reference, i.e. it points to a single thing in the world
- The reference is constant over time
  - ‘Tottenham goalkeeper’ is not a Person name, as its reference will change over time
  - ‘Al Yankovic’ is a name, as its reference does not change

Some examples which are regarded names:

- Names: E.g., ‘Charlie Chaplin’, ‘Chaplin’, ‘Charlie’
- Initials: E.g. ‘C.C.’ or ‘CC’
- Spelling mistakes: E.g. ‘Charie Chaplin’
- All orthographic variations E.g. ‘charlie chapline’, ‘Charlie CHAPLIN’
- Names of fictional characters: ‘Bilbo Baggins’, ‘Hercule Poirot’, ‘Ryan North’

Some examples which are regarded non-names:

- Pronouns (he, her, they, them)
- Expressions of time
- Currencies (Euro, Kroner, Dollar, etc.)
- Language names and nationalities (English, Spanish, German, Chinese)

## 2.1 Entity types

In this tagset, there are 3 entity types:

- Person (PER)
- Organization (ORG)
- Location (LOC)

The same surface form can belong to different types due to different contexts (metonymy). Thus, the same name can refer to both the organization and the place where it is:

- She worked for Theater Cafe<sub>(ORG)</sub>, but they met at the Cafe Love<sub>(LOC)</sub>.

“The White House” is a location when something happens there, and an organization when they decide something.

### 2.1.1 Entity type: Person

The person name category includes names of real people and fictional characters.

- 10 Best Sonnets by William Shakespeare<sub>(PER)</sub>
- You’re a mean one, Mr. Grinch<sub>(PER)</sub>
- She reminds me of Red Riding Hood<sub>(PER)</sub>’s grandma.

Family names should be annotated as Person, even if they refer to several people.

- We’ve been invited to the Baums’<sub>(PER)</sub> house tomorrow.
- And, uninvited, the Brothers Grimm<sub>(PER)</sub>

Gods are also annotated as Persons, when capitalized and having a single reference

- God<sub>(PER)</sub> plays a big role in life
- Thor’s<sub>(PER)</sub> hammer is too heavy for me

When not having a single reference, it is not tagged:

- God or not, Zlatan<sub>(PER)</sub> is the best.
- After a long run, I feel like a god.

Do not tag animal names as persons:

- The famous TV-dog Lassie
- Bo Obama is a pet dog of the Obama<sub>(PER)</sub> family

Username should be annotated as Persons if they are used as names:

- lolgrapeuk<sub>(PER)</sub> told me that thunderbiscuit<sub>(PER)</sub> is joining later

### 2.1.2 Entity type: Organization

Include any named collection of people, such as firms, institutions, organizations, pop groups, sports teams, unions, political parties etc.

- CBS<sub>(ORG)</sub> bought BBC<sub>(ORG)</sub>
- Pixar<sub>(ORG)</sub> created Monsters Inc.
- The Department of Defense<sub>(ORG)</sub> is a top group in the US<sub>(LOC)</sub> government.
- Having left the Republican Party<sub>(ORG)</sub>, John Smith<sub>(PER)</sub> will head up United Airlines<sub>(ORG)</sub>.

Organization also includes names of places when they act as administrative units or sports teams:

- Baltimore<sub>(ORG)</sub> lost to Indianapolis<sub>(ORG)</sub> last weekend
- Boston<sub>(ORG)</sub> announced yesterday that mask restrictions would be lifted

Documents from the “reviews” corpus in UNEH can be tricky. The establishments they mention are very often a location, but are treated like an organization. Use your judgment in these cases.

- I’m never going back to Dusty’s<sub>(LOC)</sub>! The food was disgusting!
- Even though I’ve never tried Dusty’s<sub>(ORG)</sub>, I’ve heard they have disgusting food.

Corporate designators like Co. and Ltd. are to be included as part of the name. The term “& co” in “Obama & co” should not be annotated because it is a designator for unnamed persons, and not an organization or a company:

- Obama<sub>(PER)</sub> & co cleaned the table

In contrast to e.g.:

- Law firm Johnson & Co<sub>(ORG)</sub> represented Hansen<sub>(PER)</sub>

Only tag brands if referring to the organization itself, not as a brand label.

- Nike<sub>(ORG)</sub> has confirmed the new shoe sizes.
- I think I got injured because I wore Nike<sub>(OTH)</sub>.

## 2.1.3 Entity type: Location

Includes geographical places, buildings and facilities. Examples are airports, churches, restaurants, hotels, tourist attractions, hospitals, shops, street addresses, roads, oceans, fjords, mountains, planets, parks and also fictional locations.

Examples:

- Germany<sub>(LOC)</sub> is very close to Sweden<sub>(LOC)</sub>
- I have reservations tomorrow morning at The Breakfast Club & Grill<sub>(LOC)</sub>
- Church of the Ascension<sub>(LOC)</sub> will have Easter services at 7am and 10am
- The Gulf of Mexico<sub>(LOC)</sub> is a beautiful blue in the morning.
- I cannot recommend the restaurant on Carson Street<sub>(LOC)</sub>.
- It takes only 15 minutes, if you go on I-70<sub>(LOC)</sub>
- If I ever traveled from Earth<sub>(LOC)</sub> to Mars<sub>(LOC)</sub>, I would want to be asleep.

When two locations are consecutive, we annotate separately:

- i have like 40 friends who live in philly<sub>(LOC)</sub> pennsylvania<sub>(LOC)</sub>.

Do not include locational descriptors unless they are part of the name:

- She lives in northern California<sub>(LOC)</sub>.
- My two favorite places are West Virginia<sub>(LOC)</sub> and South Sudan<sub>(LOC)</sub>.

Postal addresses are an exception to the separate annotation rule: they should be annotated as a whole span, as they constitute a unique reference. Examples:

- 5900 Penn Avenue, Pittsburgh<sub>(LOC)</sub>
- Target Headquarters, 1000 Nicollet Mall TPS-3165 Minneapolis MN 55403, United States<sub>(LOC)</sub>.
- Apple Inc. One Apple Park Way, Cupertino, CA, 95014, United States<sub>(LOC)</sub>.

## 2.1.4 Entity type: OTHER

In the course of annotation, there may be cases in which a span feels like an entity, but doesn’t fit into one of the above categories. For these cases, use the OTHER tag. These annotations will most likely not be included in the final set, but

categories. For these cases, use the OTHER tag. These annotations will most likely not be included in the final set, but can be used to check disagreements between annotators, and possibly stored as a separate annotation layer.

Tag nationalities and languages as OTHER.

- The Argentinian<sub>(OTH)</sub> diplomat spoke Greek<sub>(OTH)</sub>

Tag brands as OTHER:

- I refuse to buy an Apple<sub>(OTH)</sub> Watch.
- My dad was adamant that Saab<sub>(OTH)</sub> was better than Porsche<sub>(OTH)</sub>

## 2.1.5 Relationship to to Universal Dependencies PROPEN

We expect that most occurrences of named entities will coincide with the PROPEN (<https://universaldependencies.org/u/pos/PROPEN.html>) tag in Universal Dependencies. However, slight differences remain: UD will at times give parts of speech to individual constituents, and may include a broader category of terms.

## 2.2 Ambiguity

Ambiguity is a frequent source of doubt when annotating. Solve ambiguity in general as follows:

- Choose the entity type based on the local context
- We assume that every entity has a base, or literal, meaning
- When there is ambiguity, either because of lack of context or genuine ambiguity, always choose the literal meaning of the word(s)
- If the context doesn't help, and the surface form is ambiguous, choose the most common usage
- If you had to link this entity to a Wikipedia page, what type would the page be?

Examples:

- Baltimore is top of the table.
  - Choose ORG by the context, sports
- I can't wait to see Indianapolis
  - Could be either ORG or LOC, choose LOC as the literal meaning
- Lawrence is my favorite.
  - Could be either PER or LOC, choose PER as the more common usage

## 2.3 Examples

### 2.3.1 Nested names

Always annotate the whole name, never nested parts.

Annotate like this:

- University of Washington St. Louis<sub>(ORG)</sub>

And not like this:

- \*University of Washington<sub>(LOC)</sub> St. Louis<sub>(LOC)</sub>

### 2.3.2 Possession/genitive

Annotate the genitive marker as part of the name

- Charlie Chaplin's<sub>(PER)</sub> office.

Note that the possessor and possessed named entity substring should be tagged separately:

- Morten's<sub>(PER)</sub> Sweden<sub>(LOC)</sub>.

### 2.3.3 Titles

We never annotate titles as a name or part of a name:

- Your Majesty
- President Lincoln<sub>(PER)</sub>
- Honorable King Harald<sub>(PER)</sub>
- Mullah Omar<sub>(PER)</sub>

When the same words refer to the person/occupation, they are not considered names:

- The governor drove into the ditch.

### 2.3.4 Names connected by conjunctions

Annotate names connected by conjunctions as multiple distinct entities:

- Eli<sub>(PER)</sub> and Carl I. Hagen<sub>(PER)</sub>

### 2.3.5 Names that include numbers

Include numbers when they are part of the entity name:

- 10 Downing St<sub>(LOC)</sub>
- 21 Pilots<sub>(ORG)</sub>

### 2.3.6 Quotations around a name

Exclude quotations and other punctuation characters as part of the name:

- “Becoming John Malkovich<sub>(PER)</sub>” starring John Malkovich<sub>(PER)</sub>.
- Oslo<sub>(LOC)</sub>-Bergen<sub>(LOC)</sub> is quite far...

### 2.3.7 Words and Phrases derived from names

Phrases that contain names should not be annotated if they are not in use as a proper noun:

- The Bergen wave is a term for popular musical styles from the Bergen<sub>(LOC)</sub> area
- The Pittsburgh left causes 10 million traffic accidents a year in Pittsburgh<sub>(LOC)</sub>.
- I have a Matisse next door, but I think it was painted by Picasso<sub>(PER)</sub>.

### 2.3.8 Proper names conjoined to other tokens

If a name includes a hyphen, include the non-name tokens in the annotation:

- New York-based<sub>(LOC)</sub> ACLU<sub>(ORG)</sub>

### 2.3.9 Combinations of Names

Names that are combined using various conventions, should be given type by the combination of the names:

- He was born in Champaign-Urbana<sub>(LOC)</sub>

In some cases, two different types of named entity will be combined. Tag first for context, and if the context is ambiguous, or there is no context, tag with the first entity.

- Steve Jones@EPC<sub>(PER)</sub> told me the prices would drop
- The firm/person I would contact is EPC/Steve<sub>(ORG)</sub> Jones.

## 2.3.10 Emails and Usernames

Email addresses and Usernames should not be tagged unless they are used in the place of a name. Examples:

- bob@motels.com<sub>(PER)</sub> said I should contact him
- I need to email help@homedepot.com<sub>(ORG)</sub> to get a refund
- You can contact us at hello@dollargeneral.com
- anybody know why @slumpdog<sub>(PER)</sub> never comments any more?
- thanks for the tip! ~SlamDunk84

## 2.3 Questions

In cases where the locative meaning applies for an organizational area, there is sometimes doubt as to whether the appropriate annotation should be LOC, or ORG.

- This year, the scheme will be introduced in the EEA
- Police attorney XX at Helgeland police district says one person has been arrested and charged
- So far this year, there have been three arrests in the Hordaland police district, and two in the Hedmark and Rogaland police districts

There is often doubt related to whether a sentence-initial word should be annotated, given that many words can both be used as a proper name and a generic term.

- The Governor of Svalbard has the highest official authority in the area
- The church has its own bodies that can take over the tasks that the King (government) currently handles
- The government has not yet commented on the incident
- The police say that there may be a connection between the misdeeds Sometimes the name is related to the context
- Hacker News sent me a question last week

Sometimes nicknames can contain a proper name. The question in these cases is whether the whole construction should be annotated, or just the proper name.

- The question is who wins the 11 electoral votes in Arizona<sub>(LOC)</sub>, the Grand Canyon state

## 3. Previous work, basis of the guidelines

These guidelines are heavily based on NorNe (<https://github.com/itgoslo/norne/blob/master/NorNe%20Annotation%20Guidelines.pdf>) annotation guidelines, with the biggest changes being the removal of Event, Product, and GPE named entities.

The guidelines are also partially based on:

- ACE (Automatic Content Extraction) English Annotation Guidelines for Entities Version 6.6 2008.06.13
  - <https://www ldc.upenn.edu/sites/www ldc.upenn.edu/files/english-entities-guidelines-v6.6.pdf>  
(<https://www ldc.upenn.edu/sites/www ldc.upenn.edu/files/english-entities-guidelines-v6.6.pdf>)

- MUC-6 Guidelines
  - <https://cs.nyu.edu/~grishman/muc6.html> (<https://cs.nyu.edu/~grishman/muc6.html>)
  - [https://cs.nyu.edu/~grishman/NEtask20.book\\_1.html](https://cs.nyu.edu/~grishman/NEtask20.book_1.html) ([https://cs.nyu.edu/~grishman/NEtask20.book\\_1.html](https://cs.nyu.edu/~grishman/NEtask20.book_1.html))
- MUC-7 Guidelines
  - [https://www-nlpir.nist.gov/related\\_projects/muc/index.html](https://www-nlpir.nist.gov/related_projects/muc/index.html) ([https://www-nlpir.nist.gov/related\\_projects/muc/index.html](https://www-nlpir.nist.gov/related_projects/muc/index.html))
- CONLL NER shared task:
  - <http://www.cnts.ua.ac.be/conll2003/ner/annotation.txt> (<http://www.cnts.ua.ac.be/conll2003/ner/annotation.txt>)
- Ontonotes 5
  - <https://catalog.ldc.upenn.edu/docs/LDC2013T19/OntoNotes-Release-5.0.pdf>  
(<https://catalog.ldc.upenn.edu/docs/LDC2013T19/OntoNotes-Release-5.0.pdf>)

## 4. Acknowledgements

Thanks to Dan Zeman, Amir Zeldes, Constantine Lignos for comments and thoughtful feedback.

## 5. Changelog

Note: these guidelines originally started at v1, and incremented to v1.1, but as part of UNER v1 release, we have retroactively decremented the versions, so that the guideline version matches the release version.

### v0

Initial

### v1 (Release Version)

- Added wording about consecutive locations, and tagging of planets