# DEFT ERE Annotation Guidelines: Entities V1.8

**Linguistic Data Consortium**
**May 23, 2014**

**Changes from V1.7**
- 2.3: Added "here" and "there" to pronominal list for ERE purposes.
- 3.1, 3.5: Added a note about generic PER mentions that reflect LOC names ("the west favors him"). Default to LOC entity type instead of PER if the entity is otherwise not taggable.
- 4.4: Added three more examples to show how tokens with multiple entities should be handled: [he/she], the [Brogan/Tremolo] wedding, and [USPS].
- 4.5: Added mention level subscript to entities in modifier position examples
- 5.3: Revised "general description" examples to remove attributive PER mentions that would be tagged as TTLs.
- 5.3: Expanded referentiality explanation to include "Japanese Americans" example
- 7.1: Updated unattributed discussion forum quotes section and added post author nicknames section.
- 7.2: Notes on allowable XML tagging in hyperlinks (between <a href="URL"> and </a> markers).

**Changes from V1.6**
- 3.5: Specified that rooms or wings within a building are the lowest level location that we will annotate. Places or objects within a room are off limits.
- 3.1, 3.4: Added note about generic PER mentions that reflect GPE names ("Americans love fast food") and added previous demonym example to this section ("The French enjoy wine").  Default to GPE entity type instead of PER if the entity is otherwise not taggable.
- 3.1, 3.3: Added note about generic PER mentions that reflect ORG names ("The Democrats are all the same"). Default to ORG entity type instead of PER if the entity is otherwise not taggable.
- 3.1: Added note about generic religious group/ethnicity mentions and explained not to tag or default to another entity type.
- 3.1: List of PER examples now includes "[65] deaths"
- 3.4: Added new examples to the demonym note.
- 4.4: Added an example with excessive coordinating conjunctions and how to split entities.
- 7.1: Expanded section to include information about quotes and end-post XML
- 2.3: Removed negative pronouns

**Changes from V1.5**
- 5.3: Further expanded section to illustrate the difference between asserted entity mentions and general descriptions.

**Changes from V1.4**
- Removed all language regarding GPE mentions being tagged as ORGs, with the exception of sports teams.

- Reorganized sections overall.
- 5.3: Added tree diagram.
- 4.4: Added rule about cases of multiple entities in a single, continuous string.

**Changes from V1.3**
- Added additional examples to many sections. Removed confusing examples throughout.

**Changes from V1.2**
- 2.6: Altered section to pertain to discussion forum data.
- 2.6.1: Added section to explain how to handle post metadata.

**Changes from V1.1**
- 2.2.3: Removed generic pronouns from "Pronominal Entity Mentions" list (anyone, anything, anybody)
- 2.3.3: Replaced the last example under ORG vs. GPE metonym note with a less ambiguous ORG mention ("the White House")
- 2.4.4: Included examples of cases where only single entity tags can be annotated from phrases involving coordination.

**Changes from V1.0**
- 2.2.2: Noted that appositives and certain other NAM+NOM combinations expressing identity or categorization (including some usages of "of" which function to express identity or of which general class something is a particular example) should be annotated by tagging the NAM and NOM mentions autonomously.
- 2.2.2, 2.2.3: Specified that pronominal phrases will be tagged NOM.
- 2.2.3: Made pronoun list under "Pronominal Entity Mentions" comprehensive.
- 2.2.4, etc.: Adjusted "Tag for Meaning" rule to "Tag for Usage" to better reflect its intent, and emphasized that some entity types may not reflect their surface forms, and annotators must use their judgment to assign the appropriate entity type.
- 2.3.2: Specified that TTL entities will consist of personal titles and honorifics, official rank or status, and specific professional positions or occupations.
- 2.3.2: Specified that TTL strings consist only of the title, position, or profession itself, independent of its organizational circumstances, and excluding articles and modifiers.
- 2.3.3, 2.3.4: Updated ORG and GPE sections and examples to reflect decision to tag mentions of nation-states or other metonyms which clearly refer to a GPE's government as ORGs, and to use GPE as the default when in doubt.
- 2.3.4: Noted language names should not be tagged as GPE mentions.
- 2.3.5: Included more examples of informal and ad hoc LOC mentions.
- 2.4.2: Adjusted language for tagging as LOC mentions of place names in common use but which do not refer to a region corresponding to a formal GPE.

- Former 2.4.6 "Compound and Multi-Word Titles" moved under 2.3.2 "Titles"
- 2.4.9: Included examples of text with all entities completely tagged.
- 2.5.3: Added positive examples of specific mentions for contrast with general/non-specific examples.
- 3.4: Added section to specify that mentions of earlier and later versions of an organization will be coreferenced of as the same ORG entity.

# Table of Contents

## 1.    Introduction

The purpose of this annotation project is to mark up texts for entities, coreference, events and relations. The primary purpose is for the annotations to describe the meaning of the text, as opposed to its syntactic or lexical aspects. The annotation is carried out level by level. This document describes the level of entity annotation and coreference.

An **entity** is a unique object or set of objects in the world – for instance, a specific person, place, or organization. A **mention** is a single occurrence of a name, nominal phrase, or pronominal phrase that refers to, or describes, a single entity. The **mention extent** is a string of text that we annotate to indicate the occurrence of an entity mention. In a later task (coreference) we cluster together multiple mentions of the same entity.

Entities may be referenced in text at three different **mention levels** – a name, a common noun or noun phrase, or a pronoun. For example, the following are all mentions of the same entity occurring at different levels:

- Name Mention (NAM):  Barack Obama
- Nominal Mention (NOM):  the incumbent
- Pronoun Mentions (PRO):  he, his

In this task, we will label five **entity types**:

- Person (PER) - Person entities are limited to humans. A PER entity may be a single person or a group.
- Title (TTL) – Used for titles, roles, professions, and honorifics.
- Organization (ORG) - Organization entities are corporations, agencies, and other groups of people defined by an established organizational structure. An ORG entity may be a single organization or a group. **NOTE:** A key feature of an ORG is that it can change members without changing identity.
- Geopolitical Entity (GPE) - GPE entities are composite entities, consisting of a physical location, a government, and a population. All three of these elements must be present for an entity to be tagged as a GPE. A GPE entity may be a single geopolitical entity or a group.
- Location (LOC) - Location entities are geographical entities such as geographical areas and landmasses, bodies of water, and geological formations as well as buildings and other permanent human-made structures. A LOC entity may be a single location or a group.
- Vehicle (VEH) -
- Weapon (WEA) –
- Medium (MED) –
- Value (VAL) -

For our purposes, a "taggable" entity is one that is explicitly mentioned in the text, regardless of part of speech, and falls into one of the five types above. For all entities we label the text string constituting the entity mention, and assign an entity type. For entity types other than Title, we also indicate the mention level (NAM, NOM or PRO), see section below for details.

## 2.    Tagging the Entity Mention, Determining its Level, and Labeling its Extent

**NOTE:** Throughout this document the extent of each entity mention is marked with [square brackets]. Counter-examples are marked with a ~~strikethrough~~. Entity types are indicated by subscript (PER, ORG, GPE, LOC, TTL). For most sections, examples will only mark the entity type relevant to that section. Examples with multiple entity types can be found in section 2.4 and onward.

Your first task is to find each taggable entity mention in the document, and label its extent – that is, the string of text that refers to the entity. Mention extents generally begin and end at word (token) boundaries. However, possessive endings ('s) and verbal contractions ('m, 've, 're) should be <u>excluded</u> from the mention extent. As a rule, you should also exclude punctuation characters like commas, periods, and quotation marks unless the same entity mention continues after the punctuation mark. Special rules for entity extents apply depending on whether they are named, nominal, or pronominal mentions, so we will consider each mention level in turn.

Once you have determined the entity mention extent), you must label the entity mention level (NP type in the current tool) and entity type (section 2.3).

### 2.1   Tagging Named Entity Mentions

A named entity mention (NAM) is a mention that uniquely refers to an entity by its proper name, acronym, nickname, alias, abbreviations, or another alternate name. For our purposes, the extent of a name is the entire string representing the name, <u>excluding</u> the preceding definite article ('the') and any other pre-posed or post-posed modifiers. These are excluded because they are not part of the entity's actual name (e.g. Bill Clinton's name is 'Bill Clinton' not 'former president Bill Clinton').

**NOTE:** Mentions of entities with names referred to as "so-called" may also be tagged NAM.

Some examples of named entity mentions follow:

- [Bob Austin]PER
- former president [George W. Bush]PER
- the [Eiffel Tower]LOC
- [IBM]ORG
- the [Yankees]ORG  (sports team)

- [Coca-Cola Bottling Co.]$_{ORG}$
- [Uganda]$_{GPE}$
- [Bowdon]$_{GPE}$, [Georgia]$_{GPE}$
- [Mt. Fuji]$_{GPE}$
- the [Kremlin]$_{ORG}$
- the [Kennedys]$_{PER}$
- [Bill]$_{PER}$ and [Hillary Clinton]$_{PER}$
- [Sean "Puffy" Combs]$_{PER}$, aka [Puff Daddy]$_{PER}$ or [P. Diddy]$_{PER}$
- the [North]$_{GPE}$ (for 'North Korea')
- the so-called [[Northern Cyprus]$_{GPE}$ Chess Federation]$_{ORG}$
- [Land of the Rising Sun]$_{GPE}$ (for 'Japan')
- the famous [Lincoln Memorial]$_{LOC}$
- the incomparable [Steven Spielberg]$_{PER}$
- the [US]$_{GPE}$ [State Department]$_{ORG}$
- the [States] $_{GPE}$ (as a nickname for the US)
- The [Gambia]$_{GPE}$

## 2.2   Tagging Nominal Entity Mentions

A nominal entity mention (NOM) is an entity mention not including the entity's proper name, referring to it by common noun phrase. For our purposes, the extent of a nominal mention is the full mention of the noun or noun phrase, including articles and all pre-posed and post-posed modifiers. This is because modifiers provide information about an entity that could later be used by systems to identify they entity by name.

**NOTE:** Noun phrases beginning with pronominals (see section 2.2.3 below), like "this group", "the other party", "few of the attendees", will be tagged as nominals.

**NOTE:** A good rule for identifying the extent of a nominal mention is that it is the extent of the text that would be replaced by a pronoun (e.g. '[the war-torn country] elected a new president' the GPE mention extent can be replaced by a pronoun '[it] elected a new president'. Replacing part of the mention extent would not make sense 'the war-torn [it]$_{GPE}$' elected a new president').

Some examples of possible nominal mentions are given below:

- [a monument]$_{LOC}$
- [a few well-known monuments]$_{LOC}$
- [some teams]$_{ORG}$
- [the building]$_{LOC}$
- [that city]$_{GPE}$
- [her country]$_{GPE}$
- [the director of the Oscar winning film *Lincoln*]$_{PER}$

- [the family]<sub>PER</sub>
- [another large company whose investors revolted]<sub>ORG</sub>
- [the presidential hopeful from Chicago]<sub>PER</sub>

**NOTE:** Appositives and certain other NAM+NOM combinations expressing identity or categorization should be tagged with care. Some usages of "of" which function to express identity or of which general class something is a particular example. Where possible, do not include nominal mention extents with named mention extents; tag them autonomously:

- [Reuters]<sub>NAM</sub> [international news agency]<sub>NOM</sub>
- [his loudest critic]<sub>NOM</sub>, [Jon Stewart]<sub>NAM</sub>
- [my brother]<sub>NOM</sub> [John]<sub>NAM</sub>
- the [Financial Accounting Standards Board]<sub>NAM</sub>, [the private-sector body based in Norwalk, Conn., that sets the nation's accounting standards]<sub>NOM</sub>
- [the informant]<sub>NOM</sub> called [Deep Throat)<sub>NAM</sub>
- [the London borough]<sub>NOM</sub> of [Greenwich]<sub>NAM</sub>
- [the city]<sub>NOM</sub> of [Denver]<sub>NAM</sub>
- [Google employee]<sub>NOM</sub> [John Doe]<sub>NAM</sub>

When it is not possible to tag NAM+NOM combinations autonomously (such as in cases where a named mention is embedded within a coreferential nominal mention), phrases are still tagged exhaustively:

- [the [Tamil Tigers]<sub>NAM</sub> separatist organization]<sub>NOM</sub>
- [the now-defunct [G17 Plus]<sub>NAM</sub> political party]<sub>NOM</sub>

## 2.3  Tagging Pronominal Entity Mentions

A pronominal entity mention (PRO) is a referring expression that corresponds to a nominal or a named entity. The extent of a pronominal mention is just the single referring unit. Below is a list of pronominal entity mentions (The referring expressions on this list will be tagged as PRO in this task.):

- all
- another
- any
- both
- each
- each other
- either
- everybody
- everyone
- everything

- few
- he
- her
- here
- hers
- herself
- him
- himself
- his
- I
- it
- its
- itself
- little
- many
- me
- mine
- more
- most
- much
- my
- myself
- one
- one another
- other
- others
- our
- ours
- ourselves
- several
- she
- some
- somebody
- someone
- something
- that
- their
- theirs
- them
- themselves
- there
- these
- they
- this

- those
- us
- we
- what
- whatever
- where
- wherever
- which
- whichever
- who
- whoever
- whom
- whomever
- whose
- you
- your
- yours
- yourself
- yourselves

Reflexive pronouns are marked in the same way as other pronouns, e.g.,

- John hurt [himself]<sub>PER</sub>

**NOTE:** Noun phrases beginning with a pronoun listed above, like "this group", "the other party", "few of the attendees", will be tagged as nominals.

Relative pronouns should only be tagged as entities if they are **<u>not</u>** part of a nominal entity mention. For example:

- [John Smith]<sub>PER</sub>, [who]<sub>PER</sub> is a friend of mine, arrived late.
- I work at [the Black Cat]<sub>LOC</sub>, [which]<sub>LOC</sub> is a small restaurant downtown.
- He saw [the students ~~[that]~~<sub>PER</sub> he would be teaching]<sub>PER</sub>.

The possessive ending ('*s*) and verbal endings ('m, 're, 've, etc.) should be <u>excluded</u> from the extent of named, nominal or pronominal mentions. For example:

- [David]<sub>PER</sub>'s house
- the [buildings]<sub>LOC</sub>'s entrance
- [I]<sub>PER</sub>'ve never been there
- [She]<sub>PER</sub>'s a smart girl

## 2.4   General Rules

The following general rules apply at all times:

1. We always tag an expression according to its usage in the context being evaluated. In other words, the annotation of an expression depends on how it is being used. We will call this rule **Tag for Usage**. For example, if we have the sentence '[Kansas] beat [Georgetown] last night' we tag 'Kansas' and 'Georgetown' as ORGs since they are referring to sports teams, even though superficially the strings appear to be referring to a GPE or LOC. See sec. 2.4.3.
2. Entities are tagged regardless of which syntactic function they fulfill. For example, entity mentions may be adjectives ("[Korean] cars"), possessive determiners ("[her] three convictions"), or prepositional complements ("at [the beach]").
3. We allow overlapping or embedded annotations in cases where a modifier of an entity itself refers to a taggable entity. We will call this the **Modifier Referent Rule**. For example, the expression 'the Pakistani people' would have two entity mentions tagged: '[the [Pakistani]$_{GPE}$ people]$_{PER}$' and 'his army' would be tagged as '[[his]$_{PER}$ army]$_{PER}$'.  See sec. 2.4.5 for more.

## 3.    Labeling the Entity Type

Once you have determined and input the entity mention extent, in addition to tagging the entity mention level, you must label the entity type. In this task, we will label X entity types: person (PER), titles (TTL), organizations (ORG), geo-political entities (GPE), locations (LOC), vehicles (VEH), weapons (WEA), medium (MED), and value (VAL). A description of each type follows.

### 3.1   Person Entities (PER)

**NOTE:** For examples in this section, only PER entities are labeled with [square brackets].

Person entities are limited to individual humans or groups of humans identified by a simple referring expression (PER.NOM), a name/nickname/alias (PER.NAM), or pronoun (PER.PRO).

If a group of people meets the definition of an ORG or GPE it should be tagged as such. Otherwise, the group should be tagged as PER. By this standard, family names and ethnic and religious groups that lack formal organizational backing are tagged as PER entities.

**NOTE:** For entities such as movements (e.g. 'Occupy Wall Street', 'the Tea Party', 'rebel movements') which encompass gray areas regarding existence of formal name and structure, use your best judgment as to whether to tag them as ORG or a PER-group. We should usually default to making fewer assumptions and use the less-specific, more conservative entity type (PER).

**NOTE:** Generic PER mentions that reflect GPE names (such as "Americans" in "Americans love fast food") should not be annotated as the PER entity type. They should be tagged as GPE entities. We should default to GPE when a PER mention would be generic and otherwise not tagged.

**NOTE:** Generic PER mentions that reflect ORG names (such as "Democrats" in "The Democrats are all the same") should not be annotated as the PER entity type. They should be tagged as ORG entities. We should default to ORG when a PER mention would be generic and otherwise not tagged.

**NOTE:** Generic PER mentions that reflect LOC names (such as "the west" in "the west still favors him" or "the world" in "the world just never accepted her for who she was") should not be annotated as the PER entity type. They should be tagged as LOC entities. We should default to LOC when a PER mention would be generic and otherwise not tagged.

Fictional characters, religious deities, and non-human characters should <u>not</u> be tagged as PER entities. However, deceased people may be tagged as PER entities (though phrases such as "corpse" or "dead bodies" are not tagged as PER entities). Some examples of PER entities are given below. Recall that counter-examples are given in strikethrough.

- [Bill Clinton]
- [Analysts] told the Guardian that…
- [Judy Garland]
- [Adam West]'s costume from the [~~Batman~~] TV series
- [~~Harry Potter~~]
- [GOP hopeful]
- the [Cartwrights]
- [the squad of [Marines]]
- [the family]
- [the house painters]
- [the Christians]
- [the linguists under the table]
- [the Arabs]
- [[Her] friend] was [[a provincial doctor]'s wife].
- [I]'ve read [~~Jane Eyre~~] 7 times.
- [~~Seinfeld~~]was [my] favorite show.
- [We] have eradicated terrorism from [[our] society].
- He has reported on [65] deaths in the last eight months.

**NOTE:** If a document refers to a religious group or ethnicity in a generic fashion, we do not tag an entity, nor do we default to an ORG or GPE entity type.

- [~~Catholics~~] talk about the Bible like it's the best book ever written.
- They kept talking about how [~~Arabs~~] all eat really well.

You may occasionally encounter an ordinal suffix like 'Jr.', 'Sr.', or 'IV'. These are considered part of a person name and should be included within the mention extent. However, Titles (including honorifics) should NOT be included in a Person entity mention extent for instance:

- Pope [Benedict XVI]
- Mr. [Albert Franklin, Jr.] was on [the research team]

## 3.2  Title Entities (TTL)

**NOTE:** For examples in this section, many entity types are labeled with [square brackets] to illustrate the relationship among the types.

Personal titles and honorifics, official rank or status, and specific employed occupations or professional positions will be labeled as TTLs when they are used to establish a particular relationship between a PER mention and the title, honorific, position, or occupation, both in the same sentence. Annotators will need to use best judgment when determining whether a position or occupation is specific enough to be tagged as a title (e.g., general types such as worker, official, member, employee, will not be tagged TTL).

Because TTLs occur in conjunction with other entities and refer to other entities, they have some special rules. Unlike other entity types, titles do not have a mention level (NAM, NOM or PRO). Title elements are often but not always adjacent to a PER entity, which can be of mention level NAM, NOM or PRO. The most frequent constructions TTLs occur in are title+name, appositives, and copula constructions.

Just as with appositives, all titles, positions, and honorifics will be marked <u>separately</u> from the individual's name. For instance, in the following sentences, there are two separate entities marked:

- [Vice President]$_{TTL}$ [Biden]$_{PER}$
- [spokesperson]$_{TTL}$ [Mary Gillette] $_{PER}$
- [Michelle Obama] $_{PER}$, [First Lady] $_{TTL}$, spoke at the event.
- [Michelle Obama] $_{PER}$, [who]$_{PER}$ is the [First Lady] $_{TTL}$, spoke at the event.
- [She]$_{PER}$ is the [First Lady] $_{TTL}$.
- The famous Greek [philosopher]$_{TTL}$ [Zeno]$_{PER}$ was born in 490 BC.

Sometimes the name of the person is split into two pieces by the title. In these cases, we will annotate a nested title.

- [Alfred [Lord]TTL Tennyson]PER

For our purposes, the extent of a TTL is limited to the string of text that describes the position or rank itself, <u>independent of</u> its organizational circumstances or relationships, and <u>excluding</u> a preceding definite article ('the') and any other pre-posed or post-posed modifiers. TTLs are often mentioned in tandem with the organization (ORG) or geo-political (GPE) entity to which they pertain, e.g.: "Mayor of Boston", "Vice President of Marketing", "the Dover Rotary Club's treasurer", "the commissioner of the National Basketball Association", etc. Thus, if the phrase containing the position title itself also contains another taggable entity type, be sure to tag each entity separately (TTL-ORG relationships will be tagged in the Relations level of annotation). Examples:

- [Treasury]ORG [Secretary]TTL [Jackson]PER
- [Justice]ORG [Minister]TTL [Giovannia Maria Flick]PER
- [Mission Control]ORG [Chief]TTL [Vladimir Solovyov]PER
- [[US]GPE Army]ORG [negotiator]TTL [Harold Norman]PER
- [US]GPE [Secretary]TTL of [State]ORG [Hillary Clinton]PER

Compound and multi-word titles may be tagged so that their mention extents include all words essential to the understanding of the title or position itself, excluding an adjacent organization or GPE mention, e.g.:

- As [Assistant Deputy Sheriff]TTL, it is Bob's duty to inquire into certain deaths.
- Japan's parliament elected him [prime minister]TTL.
- Anne became the [heir apparent]TTL for the remainder of William's reign.
- White House [Chief of Staff]TTL Jack Lew resigned last year.
- former [President]TTL Bill Clinton

The only exception to excluding an ORG or GPE mention from a title mention extent is when one is nested within it:

- [Deputy [Foreign]ORG Minister]TTL.

In some cases, more than one title will be presented for a single person. In case all of these titles are adjacent to the person name, we will tag each of the titles separately as TTL:

- [Karachi]GPE [Mayor]TTL and [GlobalCom]ORG [Chairman]TTL [Mr]TTL [Anwar Ayub]PER

**NOTE:** When a string that may serve as a title or position is being used as a nominal to refer to a person who is not named in relation to it  within the sentence, the entity

mention should be <u>tagged as PER and not as TTL</u>. For instance, consider these two contrasting cases:

- [Joe Biden]<sub>PER</sub> is [Vice President]<sub>TTL</sub>
- [The strongest supporter]<sub>PER</sub> was [the vice president]<sub>PER</sub>.

**NOTE:** Phrases that establish ad hoc, informal, or general roles—as opposed to more specific, formal positions, occupation, titles, and honorifics—should <u>not</u> be tagged as titles:

- Suzie, the ~~[brain]~~<sub>TTL</sub> in [the family]<sub>PER</sub>., explained how the machine worked.
- ~~[Class clown]~~<sub>TTL</sub> John made a joke.
- Amy was always the ~~[philosopher]~~<sub>TTL</sub> in our discussions.

### 3.3   Organization Entities (ORG)

**NOTE:** For examples in this section, only ORG entities are labeled with [square brackets].

Organization entities are groups of people defined by an established organizational structure, identified by a simple referring expression (ORG.NOM), a named expression (ORG.NAM), or a pronoun (ORG.PRO).

**NOTE:** Sets of people who are not formally organized into a unit should be treated as a PER entity rather than an ORG entity. This distinction can sometimes be difficult. <u>If in doubt, label the group as PER instead of ORG</u>. Some examples of entities that should be treated as PER entities instead of ORG entities are:

- [the delegation]<sub>PER</sub>
- [Occupy Wall Street]<sub>PER</sub>
- [Police]<sub>PER</sub> arrested [the group of rebels]<sub>PER</sub>

**NOTE:** Organizations that share their name with a publication (whether printed or digital) should only be tagged as ORGs when it's clear that the organization is being referred to, not the publication. Publications are not, themselves, considered organizations. For instance:

- The [New York Times]<sub>ORG</sub> announced that it has named a new CEO.
- Bob enjoys reading the ~~[New York Times]~~<sub>ORG</sub> on Sunday.
- [Facebook]<sub>ORG</sub> is headquartered in Menlo Park, CA.
- i saw on ~~[facebook]~~<sub>ORG</sub> there was something on the bbc saying the earth had exploded

**NOTE:** Generic PER mentions that reflect ORG names should not be annotated as the

PER entity type. They should be tagged as GPE entities. We should default to ORG when a PER mention would be generic and otherwise not tagged.

- The [Democrats]<sub>ORG</sub> are all the same.
- The [Republicans]<sub>ORG</sub> supported Ryan.

Organizations include the following subtypes: Governmental; Commercial, Educational, Scientific, Medical; Media; Religious, Social, Advocacy; and Sports. Though we will not be labeling these subtypes explicitly, it is useful to consider examples of them:

### *Governmental (includes Political, Quasi-Governmental, Military, and Para-Military Groups)*

- [Republican Party]
- [Labour Party]
- the [Socialist People's Party]
- [Republican National Committee]
- [ACLU]
- The [Cato Institute]
- [NATO]
- The [World Bank]
- three of the [U.N.] workers stationed in East Timor
- [International Monetary Fund] aid
- [Hizbollah]
- [Islamic Resistance]
- [Rally for Congolese Democracy]
- [Institutional Revolutionary Party]
- the [KKK]
- [Al Aqsa Martyr's Brigade]
- [Tamil Tigers]
- the [Caravan of Death], [a military party that killed 73 political prisoners]
- the leading deputy of the [Rally for Congolese Democracy], [one of the biggest rebel movements supported by Uganda]
- [The Salzburg prosecutor's office] is investigating the disaster to determine if criminal charges could be filed.
- Putin, a former [KGB] agent, defended [the court that convicted Pope and [the security services]].
- The [Financial Accounting Standards Board] will take no conclusive action on [its] current project on business combinations until [Congress] has reconvened in 2001.
- [The US navy] now says the USS Cole was being refueled when an explosion ripped through it in Yemen last week, killing 17.

*Commercial*

- the [Roundabout Theater Company] is calling [its] new facility in Times
- [Pixar], [the award-winning animation company]
- the [American Airlines Theater]
- Pope, who owns [TechSource Marine Industries] in State College
- Like [the famous Irish group] the [Chieftains], [Solas] frequently headlines in Celtic festivals.

*Educational, Scientific, Medical*

- [George Washington University]
- [Overseas Chinese Physics Institute]
- [Gulf Coast Research Laboratory]
- [A coalition of medical and health groups from [Massachusetts General Hospital]]
- Pope had worked for the [Applied Research Laboratory] at [Pennsylvania State University].
- [NDSU] and [University of Minnesota] weeds specialist Alan Dexter says 98% of the plants survived.

*Media*

- [Agence France Presse]
- [abc news]
- [Associated Press]
- [Chicago Sun-Times]
- [National Geographic]

*Religious, Social, Advocacy*

- [German Bishops Conference]
- [Rock the Vote]
- [American Medical Association]
- [American Council on Education]
- [National Rifle Association]
- [American Diabetic Association]
- [NAACP]
- [American Bar Association]
- [National Center for Public Policy and High Education]
- The [Red Cross] said about 15 people managed to escape...

*Sports*

- [Taekwondo Association]
- [Philippines Olympic Committee]
- [national hockey league]
- [San Francisco 49ers]
- A group of survivors belonging to [a German ski club in Vilseck, Germany]

## 3.4   Geopolitical Entities (GPE)

**NOTE:** For examples in this section, unless specified, all entity types labeled with [square brackets] are GPE.

Geo-Political Entities are nations or subordinate types of politically-defined territory such as provinces, states, counties, cities, etc.). For something to be taggable as a GPE, it must consist of three elements: political organization, population, and physical territory. Note that sometimes a GPE mention may appear to refer more strictly to the physical location, but in such cases we still tag it as a GPE—for example:

- We went to [France]<sub>GPE</sub> for our vacation.
- They delivered the supplies to [Pakistan]<sub>GPE</sub>

Sometimes the context makes it appear that the mention of the geo-political unit, the capital, or government location is referring specifically to the government itself. In these cases we still tag the mention as a GPE. For instance:

- [Iraq]<sub>GPE</sub> signed a treaty with [Kuwait]<sub>GPE</sub>.
- [Washington]<sub>GPE</sub> discussed economic policies with [Moscow]<sub>GPE</sub> at the summit.
- [The government of [France]<sub>GPE</sub> ]<sub>ORG</sub> welcomed the agreement.
- [India]<sub>GPE</sub> is interested in strengthening economic ties with the [US]<sub>GPE</sub>.
- The Premier said [China]<sub>GPE</sub> would continue on a path of economic liberalization.
- [Turkey]<sub>GPE</sub> regards [Northern Cyprus]<sub>GPE</sub> as a sovereign country.

GPE entities can be single GPEs or groups of GPEs, for example:

- I visited [Britain]<sub>GPE</sub>, [France]<sub>GPE</sub>, and [Germany]<sub>GPE</sub> last summer. I had a great time visiting [these countries]<sub>GPE</sub>.

**NOTE:** When a GPE name or demonym or adjectival GPE name is used to refer to the *people* of a GPE (and the mention is <u>specific</u>), it should be tagged as a PER entity. For example:

- The [Swiss]<sub>PER</sub> have joined us on the bus tour.

- Luckily, the [Australians]<sub>PER</sub> made it to the barbeque on time.

**NOTE:** <u>Generic </u>PER mentions that reflect GPE names should not be annotated as the PER entity type. They should be tagged as named GPE entities. We should default to GPE when a PER mention would be generic and otherwise not tagged. For example:

- [Americans]<sub>GPE</sub> love fast food.
- The [French]<sub>GPE</sub> enjoy wine.
- [Canadians]<sub>GPE</sub> appreciate hockey.

**NOTE:** Use caution with languages. Generally names of languages are not taggable as GPE mentions:

- ~~French~~ is spoken in much of Africa.
- All nations of the former [Yugoslavia]<sub>GPE</sub> have ~~Serbian~~-speaking regions.
- ~~Arabic~~ is a major international language.
- The most widespread [Indian]<sub>GPE</sub> languages are ~~Hindi, Marathi, Tamil, Urdu, Bengali~~, and ~~Telugu~~.

**NOTE:** Sometimes the names of GPE entities may be used to refer to other things associated with a region besides the government, people, or aggregate contents of the region. The most common examples are sports teams:

- [New York]<sub>ORG</sub> defeated [Boston]<sub>ORG</sub> 99-97 in overtime.

As always, we Tag for Usage. So in the example above, both 'New York' and 'Boston' are ORG entities. Note however, that GPE names nested within sports team names should still be tagged as GPEs:

- The [[Philadelphia]<sub>GPE</sub> Eagles]<sub>ORG</sub>

Additional examples of GPEs include the following. In the examples below, only GPE entities are enclosed in [square brackets].

- Hospital officials said all eight survivors were [German]<sub>GPE</sub>.
- the conversion to Christianity of the [Roman]<sub>GPE</sub> emperor Constantine
- [Salzburg] governor Schausberger said…
- Recounts are only just beginning in [Palm Beach]<sub>GPE</sub> and [Volusia counties]<sub>GPE</sub>.
- The economic boom is providing new opportunities for women in [New Delhi]<sub>GPE</sub>.
- …said Norbert Karlsboeck, mayor of [Kaprun]<sub>GPE</sub>, [a town some 50 miles south of [Salzburg]<sub>GPE</sub> in the central [Austrian]<sub>GPE</sub> Alps]<sub>GPE</sub>
- [France]<sub>GPE</sub>'s greatest national treasure

- [France]<sub>GPE</sub> produces better wine than [New Jersey]<sub>GPE</sub>.
- [Israeli]<sub>GPE</sub> troops
- The [Palestinian]<sub>GPE</sub> Authority has banned rallies that are pro-[Iraq]<sub>GPE</sub>

**NOTE:** Countries of countries, such as 'the [European Union]<sub>GPE</sub>' and 'the [United Kingdom]<sub>GPE</sub>' will be annotated as GPEs, since they have all three GPE components (i.e., a population, a government, and a location).

The same formula applies to contested areas like Taiwan; see also sec. 2.4.8 below.

## 3.5    Location Entities (LOC)

**NOTE:** For examples in this section, only LOC entities are labeled with [square brackets].

Location entities are geographically or astronomically defined places that do not have a political component, natural structures like bodies of water and mountains, and human-made structures like airports, factories, and streets. Locations are identified by a simple referring expression (LOC.NOM), a named expression (LOC.NAM), or a pronoun (LOC.PRO).

Examples of place-related strings that are tagged as LOC include heavenly bodies, continents, non-politically-defined regions, airports, highways, street names, factories, cafes, manufacturing plants, street addresses, oceans, seas, straits, bays, channels, sounds, rivers, islands, lakes, national parks, mountains, and monumental structures, such as the Eiffel Tower and Washington Monument. Fictional or mythical locations should <u>not</u> be tagged. Also, rooms or wings within a building are the lowest level of granularity that we will annotate. Objects or places within a room should <u>not</u> be tagged. For instance:

- repairs began on the [Alaskan Pipeline]
- [his driveway] was being refinished
- ~~[The wall]~~ and ~~[coffee table]~~ in [the living room]
- Vice President Cheney visited [the site].
- In Armenia, the three of them will join other, similar delegations from around [the world]...
- The droids landed on ~~[Tatooine]<sub>LOC</sub>~~
- ... eclipse fans are being warned not to look directly at [the sun] ...
- the [Missouri River]
- the collapse of the newly-constructed [Teton Dam]
- [the region where the movement  has found most success recently]
- Police are asking everyone to avoid [the affected area].
- Many people in [North America] saw a partial solar eclipse yesterday.

**NOTE:** Generic PER mentions that reflect LOC names (such as "the west" in "the west still favors him") should not be annotated as the PER entity type. They should be tagged as LOC entities. We should default to LOC when a PER mention would be generic and otherwise not tagged. Other examples include:

- [Europeans]LOC travel all summer.
- [the world]LOC just never accepted her for who she was

## 4.    Difficult Cases and Interactions Among Entity Types

### 4.1   Determiners and Mention Span

The general rule is that determiners are included with nominal mention extents, but not with named mention extents. Determiners are included in the annotation of nominal entities that contain a named entity, as in the following example.

- [a [Gulshan Hotel]LOC driveway]LOC
- [the [Smith]PER's house]LOC

This nesting is particularly common when a NAM entity is adjacent to a NOM entity over which the article has scope.

### 4.2   The Extent of LOC and GPE mentions

There are several issues surrounding the expression of LOC and GPE entities and which parts of a string to tag.

LOC or GPE compound expressions in which place names are typically separated by a comma in English should be tagged as separate entities.

- [Kaohsiung]GPE, [Taiwan]GPE
- [Ford's Theater]LOC, [Washington D.C.]GPE

When a "designator" is customarily used as a regular part of a place name, that word should also be included in the extent of the entity. For example, include in the tagged string the word 'River' in the name of a river, 'Mountain' in the name of a mountain, 'City' in the name of a city, etc., if such words are contained in the string.

- [Mississippi River]LOC
- the [Himalayan Mountains]LOC
- [New York City]GPE

Often times place names are modified by words like 'Southern', 'Lower', 'West', 'the former' and so on. When these modifiers are part of a location's official name they should be tagged as part of the name. For instance:

- [Upper Volta]$_{GPE}$
- [North Dakota]$_{GPE}$

A place name in common use but which does not refer to a region corresponding to a formal GPE should be tagged as a Named location:

- the [Middle East]$_{LOC}$
- the [West Bank]$_{LOC}$
- [Eastern [Europe]$_{LOC}$]$_{LOC}$
- [Siberia]$_{LOC}$

Place names may present difficulties. If you are not sure whether a modifier is part of an official name, you should include the modifier as part of the place name.

Names of regions within GPEs should be tagged as Nominal locations, and the GPE within them should be tagged as well

- the [western [United States]$_{GPE}$]$_{LOC}$
- [southwest [Germany]$_{GPE}$]$_{LOC}$

## 4.3 Entity Types and Tag for Usage

It often happens that the name of one entity is used to refer to another entity. You may also encounter multiple mentions of the same entity that invoke different entity types. Surface forms and meanings may belie actual usage for some entities, so you will need use your judgment in assigning the appropriate entity type—always Tag for Usage, as in the examples below.

- [Wouters]$_{PER}$, 42, died an hour later at **[St. John Macomb Hospital]$_{LOC}$**
- [The suspect]$_{PER}$ died later the same night, **[hospital]$_{ORG}$** [spokeswoman]$_{TTL}$ [Rebecca O'Grady]$_{PER}$ said Thursday.
- **[America]$_{ORG}$** brought home the gold.  (sports team)
- Secretary of Defense William S. Cohen said today that he is satisfied **[Beijing]$_{GPE}$** will not continue sales of anti-ship missiles
- **[The Guggenheim Museum]$_{ORG}$** announced a new acquisition
- **[The Guggenheim Museum]$_{LOC}$** was designed by [Wright]$_{PER}$
- **[Deep Throat]$_{PER}$** was the pseudonym given to [the secret informant to [Woodward]$_{PER}$ and [Bernstein]$_{PER}$]$_{PER}$.
- [He]$_{PER}$ flew into **[JFK]$_{LOC}$** yesterday.

Notice we may need to ignore references to certain entity types within a mention in order to tag the string's basic usage in context. E.g., while in "The Armenians said…", "Armenians"  means "persons who are citizens of the nation of Armenia", it will only

be tagged as PER, and not GPE, because it is being used as a PER entity, and we wish to avoid multiple tags of one string.

## 4.4   Expressions that refer to multiple entities

Care is needed when dealing with coordination in entities. When a phrase refers to multiple, coordinated entities, mark each entity separately where possible. For instance:

- [China]<sub>GPE</sub> and [South Korea]<sub>GPE</sub> signed the agreement.
- [Jimmy]<sub>PER</sub> and [Rosalyn Carter]<sub>PER</sub>
- [North]<sub>LOC</sub> and [South America]<sub>LOC</sub>

But be careful not to split apart proper names that contain a conjunction. For instance:

- [Trinidad and Tobago]<sub>GPE</sub>
- the [Fish and Wildlife Service]<sub>ORG</sub>

The latter example is the name of one organization and should be tagged as a single named entity (it's not 'the Fish Service' and 'the Wildlife Service' as separate names).

When conjunctions are used excessively (mainly when referring to nominal entities), you may split apart the entities at the coordinating conjunctions. For example:
- Instead of tagging: [[my] stepkids and friends and family]
- We can tag: [[my] stepkids] and [friends]and [family]

Some cases of coordination may necessitate a phrase being tagged as a single entity, such as in cases where only a single noun is present but coordinated modifiers might suggest two distinct entities. For instance:

- [American and Canadian soldiers]<sub>PER</sub>
- [the East and South of Iran]<sub>LOC</sub>
- [the CEOs of Google and Youtube]<sub>PER</sub>

Cases where multiple entities are joined together in a single, continuous string are also tagged as a single mention:

- they crossed the [Af-Pak] border
- [Brad&Angelina ]
- [me+you]
- [he/she]
- the [Brogan/Tremolo] wedding

- [USPS] - we're unable to tag [[US]$_{gpe}$PS]$_{org}$ because we cannot split the token
- [Obamacare] – we'd tag this as a named PER because "Obamacare" wouldn't be an otherwise taggable entity. Obama's name is within the token and we should capture this information.

## 4.5 Entities in modifier position

If an entity mention contains another taggable mention nested within it, these nested entities should also be tagged. This applies both to named and nominal entity mentions, for example:

- [the [Clinton]$_{PER.NAM}$ government]$_{ORG.NOM}$
- [[Treasury]$_{ORG.NAM}$ employees]$_{PER.NOM}$
- [[U.S.]$_{GPE.NAM}$ exporters]$_{ORG.NOM}$
- [[Texas]$_{GPE.NAM}$ schools]$_{ORG.NOM}$
- [Kentucky]$_{GPE.NAM}$ Fried Chicken]$_{ORG.NAM}$
- [[government]$_{ORG.NOM}$ offices]$_{ORG.NOM}$
- [[Kurdistan]$_{GPE.NAM}$ Freedom Fighters]$_{ORG.NAM}$
- [the [Midwestern]$_{LOC.NOM}$ bank]$_{ORG.NOM}$
- [the [Russian]$_{GPE.NAM}$ foreign minister]$_{PER.NOM}$
- [the [American]$_{GPE.NAM}$ companies]$_{ORG.NOM}$
- [[Israeli]$_{GPE.NAM}$ troops]$_{PER.NOM}$
- [[Republican]$_{ORG.NAM}$ voters]$_{PER.NOM}$
- [[airline]$_{ORG.NOM}$ regulators]$_{PER.NOM}$
- [[Chrysler]$_{ORG.NAM}$ factories]$_{LOC.NOM}$
- [[union]$_{ORG.NOM}$ leaders]$_{PER.NOM}$
- The [[Chinese]$_{GPE.NAM}$ military]$_{ORG.NOM}$

## 4.6 Possessives

When you encounter a possessive construction, it may contain two taggable entity mentions. Note that when the construction is comprised of two named mentions, such as in the third example below, the two entities are tagged separately (i.e. the possessive entity is not embedded).

- [[Temple University]$_{ORG.NAM}$'s graduate school of business]$_{ORG.NOM}$
- [[Canada]$_{GPE.NAM}$'s parliament]$_{ORG.NOM}$
- [Singapore]$_{GPE.NAM}$'s [Central Narcotics Bureau]$_{ORG.NAM}$

## 4.7 Hyphenated pre-modifiers

Taggable entities that are part of a pre-modifying hyphenated construction should be tagged separately, for example:

- The [GOP]$_{ORG}$-backed candidates toured the area.

## 4.8 Places of contention

Places of contention can be tagged as GPEs long as they have all three components of a GPE (i.e. GPE = population + location + government). If a place of contention does not have all three of these components, it should be tagged as a LOC instead.

Using this rule, 'Palestine' is tagged as a GPE because it has all three GPE components, while 'Gaza strip' is tagged as a LOC, because though it has a population and a location, it doesn't have its own government.

## 4.9 Examples with entities completely annotated

- Videos circulated by [Osama bin Laden]PER.NAM have added to the evidence linking [him]PER.PRO and [the [al-Quaida]ORG.NAM network]PER.NOM to the Sept. 11 terrorist attacks in the [United States]GPE.NAM, [the government]ORG.NOM said Wednesday in an updated dossier on the investigation.
- [Guzman]PER.NAM indicted [Pinochet]PER.NAM, holding [him]PER.PRO responsible for the actions by the "[Caravan of Death]PER.NAM", [a military party that killed [73 political prisoners]PER.NOM shortly after the 1973 coup in which [Pinochet]PER.NAM ousted Marxist [President]TTL [Salvador Allende]PER.NAM]PER.NOM.
- Midway through the hearing, [Chief Justice]TTL [Renquist]PER.NAM seemed to scold [[his]PER.PRO colleagues]PER.NOM for being too talkative when [he]PER.PRO made an unusual offer to [the lawyer representing [[Florida]GPE.NAM's Attorney General]PER.NOM]PER.NOM.
- [Actors and singers also on the flight]PER.NOM held a benefit concert in [Baghdad]GPE.NAM Saturday evening, with most of the $13 cover charge to be donated to support the [Palestinian]GPE.NAM uprising.
- ...said [Archbishop]TTL [Khajag Barasamian]PER.NAM, [head]TTL of the [Diocese of the [[Armenian]GPE.NAM Church]ORG.NAM in [America]GPE.NAM]ORG.NAM, [[whose]ORG.PRO headquarters]LOC.NOM are in [Manhattan]GPE.NAM.

## 5. What NOT to tag

### 5.1 Events

Do not tag event names even if they refer to events that occur on a regular basis and are associated with institutional structures. However, the institutional structures themselves —steering committees, etc. —should be tagged.

- the ~~Pan-American Games~~
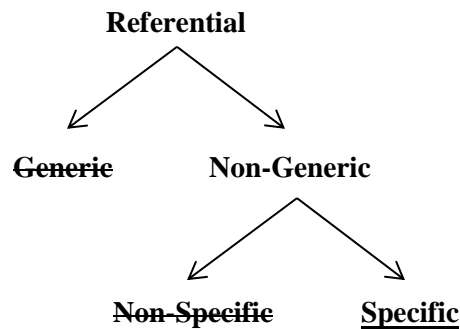- the [Olympic Committee]ORG

## 5.2   Artifacts and Products

Miscellaneous types of proper names that are not to be tagged as named entities include artifacts, other products, and plural names that do not identify a single, unique entity. For instance:

- the ~~Taurus~~ is the latest car model

## 5.3   Generic and Non-Specific Entities

Entity mentions that are generic or non-specific should also not be tagged. In ERE, only specific entities are annotated.

**Referential**

~~Generic~~          **Non-Generic**

~~Non-Specific~~          **Specific**

Do not tag mentions which do not refer to a particular, asserted entity which is understood to be real by the immediate writer, speaker or 'voice' in the text. Consider the following cases:

- I'd like to find ~~[a good dentist]~~$_{PER}$.
- I've got [a good dentist]$_{PER}$.
- Have you ever been to ~~[a trustworthy car dealership]~~$_{ORG}$?
- I once found [a trustworthy car dealership]$_{ORG}$, but it closed.
- ~~[More and more people]~~ are losing jobs to outsourcing.

Names used to refer to a generic category of things and do not refer to a specific entity should not be tagged

- the ~~Campbell Soups~~ of the world

Certain constructions can make more difficult the decision of whether or not to categorize an entity mention as specific. Sometimes a phrase may seem to be more a general description of an entity than an actual mention of an entity. Overly general descriptions are considered underspecified and are not tagged. As is always the case in ERE, tag for meaning/usage. If you think something is only a description and not a mention, don't tag it. Consider the following:

- [John Jacob Marshall]PER.NAM, [the friend who drove me home last night]PER.NOM
- [The guy you saw yesterday]PER.NOM is[Jerry]PER.NAM
- [Dennis R. Beresford]PER.NAM is [that outspoken backgammon enthusiast]PER.NOM
- [He]PER.PRO is ~~[a pain in the neck]~~
- [He]PER.NAM is ~~[my hero]~~
- [John]PER.NAM is a [plumber]TTL
- [My friends]PER.NOM are ~~[great people]~~
- [She]PER.PRO is the [First Lady]TTL
- [France]GPE.NAM is ~~[a nice place to vacation]~~

Similarly, phrases describing an entity's role (other than a Title) in some capacity are generally seen as descriptions, not coreferent entity mentions:

- [Atlanta]GPE.NAM is ~~[a major hub between the Northeast and the Caribbean]~~
- [Bob]PER.NAM is usually ~~[the one who picks up lunch for everyone]~~

The distinction between specific and generic/non-specific entities can be difficult to make, and you will need to use context to make this distinction. E.g.: 'analysts are interesting people' is a generic reference to a category of people in general, while 'analysts told the Guardian that ..' is a reference to specific people. Essentially, if you as a reasonable reader would interpret the entity as specific, go ahead and annotate it. If you are unsure as to whether something is specific or generic, leave it unannotated.

Fully annotated example with explanation:
- The [US]GPE.NAM confined [[Japanese]GPE.NAM Americans]PER.NAM during the war but did not do so with [Italian]GPE.NAM [Americans]GPE.NAM and [German]GPE.NAM [Americans]GPE.NAM.

  - The mentions of "Italian Americans" and "German Americans" are *non-generic* in this example, but they're also *non-specific*. Because of this distinction, we do not tag the mentions as PERs. We default to GPE entities when the PER meaning is non-specific or generic. So, because "Italian America" and "German America" (arguably) aren't valid GPEs, we tag each word as separate GPE entities.

  - Conversely, "Japanese Americans" is *non-generic* as well as *specific*, because the speaker is referring a subset of Japanese Americans that were detained. In this case, we'd tag the string as a PER with Japanese nested within.

  - "Japanese Americans" would also be a Person ARG for a Justice.Arrest-Jail event. As a general rule, we don't change entity types just to grab a

Relation or Event argument. However, thinking of entities in terms of Relations and Events can help us decide how to annotate a mention.

## 6.  Coreference

## 6.1  General Instructions for Coreference Tagging

Once all taggable entities have been labeled, the next task is entity coreference. During coreference labeling we cluster together multiple mentions of the same entity within the document. If two or more mentions refer to the same underlying entity, we must indicate this by coreferencing them, regardless of the entity level (NAM, NOM, PRO).

Coreference can only be done when mentions share the same entity type (PER, ORG, LOC, GPE or TTL). If you want change the type of an entity during coreference, you must go back to the Entity task and change its type before continuing with coreference.

In most cases annotating coreference is very straightforward. In a document about Osama bin Laden, we want all mentions of 'bin Laden' to be lumped together in the same entity of type PER. In the following passage, all the bracketed mentions should be coreferenced as one entity:

- Videos circulated by **[Osama bin Laden]** have added to the evidence linking **[him]** and the al-Qaida network to the Sept. 11 terrorist attacks in the United States, the government said Wednesday in an updated dossier on the investigation. The document, published by Prime Minister Tony Blair's office, said **[the Saudi dissident]** had come "closest to admitting responsibility" for the attacks in an "inflammatory video," allegedly made on Oct. 20 that was not released to the media but circulated to al-Qaida members. "The battle has been moved inside America, and we shall continue until we win this battle, or die in the cause and meet our maker," the document quotes **[bin Laden]** as saying.

The name mentions of 'Osama bin Laden' are easy to spot. However, it is important to annotate all mentions that refer to the entity that is 'bin Laden'. This will include nominal mentions bolded above, such as 'the Saudi dissident' and pronominal mentions such as 'him', which would all be coreferenced together.

**NOTE:** All of these coreferring mentions in the example above have the same entity type: PER.

By contrast, 'Saudi' is a GPE, which means it cannot be coreferential with the PER mentions of 'bin Laden':

- said [the [Saudi]<sub>GPE</sub> dissident]<sub>PER</sub> had come "closest to admitting…"

On the other hand, if we had a sentence like the following:

- The document said [the Saudi]<sub>PER</sub> had come "closest to admitting…"

we would label 'the Saudi' as a PER entity in accordance with the Tag for Usage rule, and this entity would be coreferred with other PER mentions of 'bin Laden'.

## 6.2 Coreference for Titles

Corefer all title mentions of the same role within the same organization, regardless of whether the mentions are associated with different persons. In the example below, the title 'President' refers to the same specific role in both cases (President of the U.S.) so the two mentions are marked as coreferring:

- [President]$_i$ Clinton and [President]$_i$ Bush did not have the same policies.

The following example shows that if the title mentions refer to different roles (President of the US vs. President of France) they should not be marked as coreferring, even if their strings are identical:

- [President]$_k$ Obama greeted [President]$_j$ Hollande.

**NOTE**: Person entity mentions (PER) and title entity mentions (TTL), being different types of entities, will not be coreferenced – but since TTL mentions will only be tagged when establishing a relationship to a PER mention for which the TTL refers, the relationship will not be lost and will be annotated at the Relations level.

**NOTE**: If you are not sure whether two titles refer to the same position (e.g. if you are not sure whether two people are president of the same country), you can look these up online.

Generic titles that are used as part of a form of address, such as *Mr, Mrs, Ms, Dr, Professor* etc. should be marked as coreferring if the same string is used (i.e. we will <u>not</u> coreference these as positions, only as strings)

- [Mrs]$_j$ Jones and [Mrs]$_j$ Smith met for coffee.
- [Dr]$_k$ Mary Arroyo and [Dr]$_k$ Jim Hills published the paper.

## 6.3 Coreference in Questions

When an entity is being questioned coreference can be marked if context makes the identification clear, e.g.:

- Dialog:
    a. A: I went to see [the breeder]$_k$.
    b. B: [Who]$_k$'s [the breeder]$_k$? Is [that]$_k$ [the breeder that you saw yesterday]$_k$?
    c. A: Yes.

## 6.4 Coreferencing Organizations Over Time

When comparing mentions of earlier and later versions of an organization (e.g., "1950s IBM" VS "present-day IBM"; "the Blair government" VS "the Brown government"), we will still coreference them as the same ORG entity.

## 7. Discussion Forums

When annotating discussion forum documents, you should expect to find more colloquial language, including spelling errors, interruptions, unclear expressions and missing punctuation. Annotate each document to the best of your understanding, trying to focus on the author's presumed intent.[1]

## 7.1 Post Metadata

**Post authors:** Each discussion forum post begins with an XML heading similar to the following:

<post author="pollywog" datetime="2009-03-24T11:34:00" id="p3">

In ERE, this data is considered taggable. Therefore, in the above example, there is one taggable entity:

[pollywog]per.nam

**Post author quotes:** XML metadata also signifies the end of discussion forum posts and the boundaries of quotes. It's important to note shifts in post authors, because we will coreference speakers accurately. Take the following example:

```
<post author="Tsukasa" datetime="2011-11-09 id="p188">
<quote orig_author="Schrodinger's Cat">
not that I'm excusing it in any way
</quote>
good! youre starting to make sense to me
</post>
```

---

[1] Note that the policies set down regarding word tokenization in this section are different from Treebank policies on some items.

When a post author quotes another poster, XML displays
<quote orig_author="X"> where X will be the name of whomever is being quoted.
Additionally, the </quote> marker signals the end of quoted text. Similarly,
</post> will mark the end of the post author's post.

So, in the above example, Tsukasa has written "good! youre starting to make
sense to me," while also quoting Schrodinger's Cat, who previously stated,
"not that I'm excusing it in any way."

**Post author nicknames:** Sometimes post authors will use nicknames that are based
on screen names for other post authors. These nicknames should still be tagged as
named Per mentions. Notice the usage in the following example:

> <post author="Scrat" datetime="2003-11-26T23:01:00" id="p173">
> <quote orig_author="Moot">
> C'mon scratsy you gotta tell us the story
> </quote>
> Eh maybe when you're older Mootay.
> </post>

Post author [Scrat] refers to post author [Moot] as "Mootay" after [Moot] refers to
him as "scratsy". It's easy to miss these named PER mentions, especially when they
aren't capitalized.

**Unattributed quotes:** Sometimes discussion forum quotes are unattributed, so
<quote orig_author="X"> will not appear. Instead, <quote> will be the only
indicator. In these cases, coreference can be more difficult, so be sure to create new
equivalence classes for entities that aren't linked with previously mentioned
entities. For instance, the unattributed quotes can be from other discussion forums.

## 7.2 URLs

Potential entity mentions embedded in URLs will generally not be tagged as entities,
as it isn't usually clear in cases such as this that any real-world entity is truly being
referenced (other than perhaps in a generic fashion). For instance, in:

> <a href="http://www.lonelyplanet.com/usa/virginia>

neither "usa" nor "virginia" are tagged as entities. Similarly, in the following,
"whitehouse" is not tagged as an entity:

whitehouse.gov

**Hyperlinks:** While we do not tag entities embedded in URLs, we **do** tag entities within XML hyperlink metadata. Generally appearing before or after a URL, hyperlink metadata surrounds text with <a href="URL"> and </a> markers. For example, a document might contain the following text:

<a href="http://www.flyers.nhl.com/club/schedule.htm">Flyers Captain Claude Giroux Scores with 5 Seconds Remaining to Win Stanley Cup</a>

Within this hyperlink metadata, [Flyers]$_{ORG}$ ,[Captain]$_{TTL}$, and [Claude Giroux]$_{PER}$ are all taggable entities.

### 7.3   Misspellings and Incorrect Punctuation

Annotate misspellings according to the intended meaning, as far as that can be deciphered. In the example below the second "I" is a typo and we can assume that the author intended to write "a". The second "I" should therefore not be marked as PER.PRO.

● I know **I** guy who can help us out

Similarly, incorrect punctuation should be ignored and the text marked according to the author's presumed intent, e.g.,

● I bought two [book's] at the store.

In the case of missing apostrophes, annotate the entire word, even if you would normally exclude the apostrophe from the mention span, e.g.,

● Call me when [your] in town.
● [Im] in town!

In the case of missing spaces, annotate the entire span even if it includes text that you would normally not annotate, e.g.,

● [Iwanna] get out of this town.
● [IDK] who that is.

### 7.4   Repetition and Fragments

In repeated text mark each mention separately, e.g., in the example below both mentions of "there" are marked as separate LOC entities, coreferring with each

other:

- I wanna dive [there]LOC…. *drive [there]LOC I mean

Annotate fragments to the best of your interpretation, e.g., in the example below there are two fragments and one complete sentence mentioning "John". All three mentions should be annotated and marked as coreferring.

- dialogue:
    a. A: [John]
    b. B: [John] was
    c. A: I saw [John]

## 7.5   Coreference in Discussion Forums

Discussion forums contain dialogues between multiple participants. Care must be taken to mark coreference correctly, especially for first and second person pronouns, e.g.,

- dialogue:
    a. [I]$_k$ want to get rid of this one.
    b. Sure [I]$_j$'ll take it off [your]$_k$ hands.