

# NerKor annotálási útmutató

Simon Eszter

2021. szeptember 7.

## 1. Bevezetés

A számítógépes nyelvészet egy interdiszciplináris terület, amelynek célja az emberi nyelv szerkezetének gépi modellálása, valamint a természetes nyelvek számítógépes feldolgozása. Az információkinyerés a számítógépes nyelvészet egyik fontos alterülete; célja, hogy strukturálatlan szövegből automatikusan hozzáfussunk a számunkra értékes információhoz. Mivel ezen információ nagy része tulajdonnevek formájában jelenik meg a szövegben, ezért a tulajdonnévfelismerésnek (named entity recognition, NER) kiemelt jelentősége van.

A NER során egy bemeneti tokensorozatban kell megnevezett entitást (named entity, NE) alkotó intervallumokat kijelölünk, ezeket véges sok kategóriába besorolva. Egy gépi tanuló algoritmus kiértékelése manuálisan annotált korpuszal való összevetés útján történik, és szokásosan maga az algoritmus is egy ilyen korpuszból tanulja meg a paramétereit automatikus módon. Ezért van szükség nagy méretű kézzel annotált korpuszokra. Ez a célja ennek a projektnek is, amelynek a tervezett kimenete a NerKor korpusz, egy egymillió tokenes kézzel NE-annotált korpusz.

A NE-k nem teljesen azonosak a tulajdonnevekkel, de a jelen annotációs sémában nem lépünk ki a nevek világán kívülre, így a 'NE' és a 'név' szavak szinonimaként fognak szerepelni a tárgyalási univerzumunkban.

A nevek egyedi referenciával bírnak, vagyis a világ egy egyedi entitására utalnak (pl. *London*). Ilyet a köznévi frázisok önmagukban állva (pl. *város*) nem tudnak csinálni, csak akkor, ha egyéb nyelvi elemekkel közösen egy olyan főnévi frázist alkotnak, amely már tényleg a világ egy egyedi entitására utal (pl. *az a város, ahol 9 millió brit lakik*). Mi itt csak a neveket annotáljuk, például:

- (1) *Kosztolányi Dezső*
- (2) *Szilas Menti Mezőgazdasági Termelőszövetkezet*
- (3) *United Nations Educational, Scientific and Cultural Organization*
- (4) *Déli-Shetland-szk.*
- (5) *IBM*

- (6) *Kiss János altábornagy utca*
- (7) *Műegyetem*
- (8) *The Coca-Cola Co.*
- (9) *Kovács Pistike*

## 2. Az annotálás alapelvei

Fontos, az annotálás során végig szem előtt tartandó szabályok:

- Csak tulajdonneveket annotálunk. Nem annotálunk olyan frázisokat, amelyek ugyan a világnak valamely egyedi részére utalnak, de nem tulajdonnévvvel. Például a *József Attila Gimnázium* annotálandó, de a szövegben szereplő az *a sul*i frázis nem, hiába derül ki a szövegből, hogy melyik iskolára utal.
- A nevek nem kompozicionálisak. Mivel a nevek jelölete nem a név részeinek a jelöletéből áll össze, ezért a neveket nem bonthatjuk részekre az annotáláskor. Például a *Kossuth Lajos utca* egy földrajzi névként jelölendő, hiába van benne egy személynév. Ebből az következik, hogy mindig a leghosszabb nevet (a legkülsőbbet) jelöljük a jelölhetők közül.
- Nem annotálunk egymást átfedő vagy egymásba ágyazott neveket. Vagyis minden annotációnak be kell fejeződnie, mielőtt egy másik elkezdődik.
- A *tag-for-meaning* elvét követjük. Vagyis egy nevet mindig az aktuális kontextusnak megfelelő referenciája alapján annotálunk.
- Ha az azonosított név ragozott formában szerepel a szövegben, a raggal együtt, a teljes alakot annotáljuk.
- A nevek képzett alakjait nem jelöljük. Nem annotálandók tehát az olyanok, mint *magyarországi*, *fideszes*, *petőfieskedő*.
- Ha a név összetétel előtagja, és az összetétel alaptagja köznévi, például *Horn-kormány*, *Tilos Rádió-hallgatók*, *TA-vezérigazgató*, akkor nem annotálandók névként.
- A névhez nem tartozik hozzá az esetleg előtte álló névelő. Kivétel az az eset, amikor a határozott névelő része a névnek, például *The Hague*, *The Times*.
- A név rövidítése (akronim, mozaikszó, monogram) is névként annotálandó.
- A szöveghez a névannotálás során nem nyúlunk hozzá, vagyis nem javítjuk ki a helyesírási hibákat, nem vonunk egybe különírt szavakat, és nem választunk szét egybeírtakat. Ha valamilyen éktelen hibát látunk az aktuálisan címkézendő névvel kapcsolatban, akkor azt külön fel kell jegyezni

az erre rendszeresített google táblázatban. Ez a szépirodalmi szövegekre nem vonatkozik: ott az az elv, hogy amit a szerző leírt, az sérthetetlen, vagyis nem kell kijavítani, és így a táblázatba sem kell felvenni.

### 3. NE-típusok

A következő típusokat annotáljuk:

**PERSON:** Valós és kitalált személyek neve, becenevek, művésznevek, álnevek. Ide tartoznak a kisebb, kevésbé strukturált embercsoportok, közösségek is.

**ORGANIZATION:** Olyan csoportok nevei, amelyek valamilyen szervezett struktúrával rendelkeznek, mint például intézmények, vállalatok, kormányzati hivatalok, sportcsapatok, múzeumok, egyetemek.

**LOCATION:** Földrajzilag vagy politikailag definiált helyek nevei, úgymint városok, országok, hegyek, völgyek stb. Ide tartoznak az emberalkotta építmények is, mint a repterek, utak, gyárok, épületek stb.

**MISC:** A felsorolt típusok egyikébe sem tartozó nevek.

Az útmutatóban a NE-eket [szögletes zárójelek] közé tesszük. A példánál csak az olyan típusú NE-eket jelöljük, amelyikről éppen szó van. A példákban a személyneveket a PER, a szervezetneveket az ORG, a helyneveket a LOC, a be nem sorolhatókat pedig MISC rövidítésekkel jelöljük.

#### 3.1. Személynevek (PERSON)

Személyekre utalhatnak teljes személynevek, becenevek, művésznevek, álnevek, rövidítések. A kitalált személyek, úgymint mozihősök, mesefigurák, mitológiai alakok, illetve a szentek, bibliai alakok nevei is személynévként annotálандók, például:

- (10) [Ady<sub>PER</sub>] írói álneve [Ida<sub>PER</sub>]
- (11) a legkisebb gyerek, aki gyakran játszik [Mikrobival<sub>PER</sub>]
- (12) zenéjével meglágyította [Hádész<sub>PER</sub>] és [Perszephoné<sub>PER</sub>] szívét
- (13) [Isten<sub>PER</sub>] szellem, legfelsőbb, láthatatlan és halhatatlan lény

Bizonyos bibliai alakok neve jajgatáskor, panaszkodáskor vagy káromkodáskor is felmerül – az ilyen esetekben nem kell névként annotálni:

- (14) *Úristen, milyen emberekkel élünk egy országban?*

A családnevek, az uralkodóházak nevei is személyekre, egészen pontosan személyek csoportjára referálnak, ezért azokat is személynévként kell megjelölni, például:

(15) *a [Széchenyi<sub>PER</sub>] család Nógrád megyéből származik*

(16) *a [Károlyiak<sub>PER</sub>] Apáti nevű faluját felgyújtották*

A zenekarok, együttesek neve is személynévnek minősül, ha az adott kontextusban személyek csoportjaként van rájuk utalva, például:

(17) *A nagyszínpadon fellép többek közt a [Wellhello<sub>PER</sub>]*

Ha viszont kifejezetten egy szervezett csoportként jelenik meg egy együttes a szövegben, aminek tagjai vannak, akkor inkább **ORG**-ként annotálandó, például:

(18) *1993-ban Erik a [Morbid Angel<sub>ORG</sub>] tagja lett.*

Ha emberek bármilyen csoportjáról úgy vélekedünk, hogy az nem vagy kevésbé rendelkezik szervezett struktúrával, inkább informális és kicsi, akkor legyen **PER**-rel jelölve. De ha inkább érezzük dominánsnak a szervezetet, akkor legyen **ORG**. Ez a szabály vonatkozik minden embercsoportra, legyen az együttes, sportcsapat vagy valamilyen közös elv mentén összejött csoport.

A személynévvel együtt csak azok a prefixek vagy szuffixok annotálandók, amelyek hivatalosan is a névhez tartoznak, például magyarban a 'Dr.' prefix vagy a 'PhD' szuffix az igazolványba is bekerülhet, illetve a 'Sir', 'Lord' vagy 'Lady' is hivatalosan a névhez tartozik az angoloknál, valamint a 'Jr.' az amerikaiaknál. De az egyéb rangot, címet, beosztást vagy rokoni viszonyt jelölő köznévi szavak nem tartoznak hozzá, például:

(19) *Lemondott [Heinz-Christian Strache<sub>PER</sub>] osztrák alkancellár*

(20) *Érvcsata [Dr. Boldogkői Zsolt<sub>PER</sub>] molekuláris biológus és [Dr. Gődény György<sub>PER</sub>] gyógyszerész között*

(21) *[Sir Winston Leonard Spencer Churchill<sub>PER</sub>] brit politikus, miniszterelnök.*

(22) *[Gyuri<sub>PER</sub>] bácsi gyógynövényturái*

(23) *Mr. [Joyce<sub>PER</sub>], jöjjön be!*

A köznévi megszólítások nem jelölendők névként, például:

(24) *Nagymama, miért ilyen nagyok a szemeid?*

A @-cal jelölt felhasználónevek is személynévnek minősülnek, így **PER**-rel jelölendők, de a @ nélkül. A # után következő token pedig csak akkor név, ha # nélkül az, és persze jelölni is anélkül kell.

(25) *@ [MazenMahdi<sub>PER</sub>]: A 2 idős férfi, akit letartóztattak tegnap # [karzakanban<sub>LOC</sub>]*

### 3.2. Szervezetnevek (ORGANIZATION)

Azok a tulajdonnevek, amelyek egy szervezett struktúrával rendelkező csoportra referálnak, amelyek aktorként szerepelnek az adott szövegkontextusban, és olyanokat tudnak csinálni, mint döntést hozni, árat emelni, nyilatkozni valamiről stb., szervezetnévként annotálандók. A következők mind ilyenek:

- Cégek, vállalatok

(26) *a [SERCO Kft.<sub>ORG</sub>] az eltelt évek során jelentős fejlődésen ment keresztül*

(27) *1878-ban Grosvenor Lowry-val létrehozzák az [Edison Electric Light Co.-<sub>ORG</sub>]*

- Multinacionális szervezetek

(28) *az [Európai Unió<sub>ORG</sub>] ezen a néven 1992-ben jött létre*

- Politikai pártok

(29) *bántalmazták a [Fidesz<sub>ORG</sub>] egyik ajánlászervevényeket gyűjtő aktivistáját*

- Sportcsapatok

(30) *A [Budapest Black Knights<sub>ORG</sub>] csapata fölényesen legyőzte a [Szolnok Soldiers<sub>ORG</sub>] csapatát.*

- Katonai szervezetek

(31) *Az [Észak-atlanti Szerződés Szervezete<sub>ORG</sub>] székhelye Brüsszelben van.*

- Kórházak, egészségügyi intézmények

(32) *A [Péterfy Kórház-Rendelőintézet Országos Traumatológiai Intézet<sub>ORG</sub>] a Főváros egyik legnagyobb egészségügyi intézménye*

- Hotelek

(33) *A paksiakat is várja az [Erzsébet Nagy Szálloda<sub>ORG</sub>]*

- Színházak, múzeumok

(34) *A [Szépművészeti Múzeum<sub>ORG</sub>] az egyetemes és a magyar művészet emlékeit mutatja be az ókortól a 18. század végéig.*

- Egyetemek

(35) *A [Kossuth Lajos Tudományegyetem<sub>ORG</sub>] Honoris Causa Doktorai.*

- Kormányzati hivatalok

(36) *A [Honvédelmi Minisztérium<sub>ORG</sub>] hivatalos Facebook oldala*

Az online vagy nyomtatott sajtótermékek esetében szervezetként annotáljuk a nevet, amikor aktívan cselekvő szervezetként lépnek fel, például megjelentetnek valamit, vagy kiállnak valami mellett vagy ellen, például:

(37) *Mi, akik a [Telexnél<sub>ORG</sub>] dolgozunk, vállaljuk, hogy*

(38) *Hétfőn a [Sabr Online<sub>ORG</sub>] újság arról számolt be, hogy*

Ha viszont ugyanezt a sajtóterméket úgy említik, mint amit olvasnak, vagy mint amin/amiben megjelent egy cikk, tanulmány, írás, akkor MISC-ként annotálandó, lásd a 3.4. fejezetben.

Az általános intézményneveket, mint *rendőrség* vagy *kormány* nem annotáljuk, mert ezek csak egy főnévi frázis részeként tudnak egy egyedi entitást jelölni, nem egyedi jelölők.

Az olyan hosszú, többtagú intézményneveket, amelyek tartalmazznak köznévi elemet is, teljes egészében névnek kell jelölni, például:

(39) *Az [Állami Privatizációs és Vagyonkezelő Rt.<sub>ORG</sub>] zártkörű részvénytársaságként működő állami vállalat volt.*

Ha egy intézménynév után zárójelbe téve szerepel a rövidítése is, akkor a teljes név és a rövidítés külön nevekként kezelendők, a zárójel pedig nem a név része, például:

(40) *Az autóipari óriás [DaimlerChrysler AG<sub>ORG</sub>] ([DC<sub>ORG</sub>]) amerikai részlege*

### 3.3. Helynevek (LOCATION)

A helynévnek annotálandó entitások közé tartoznak többek között a kontinensek, az országok, a régiók, a városok, a települések, a repterek, az utak, a gyárak, az óceánok, a tengerek, a folyók, a szigetek, a tavak, a nemzeti parkok, a hegyek és a mitikus helyek, például:

(41) *[Franciaországot<sub>LOC</sub>] kilenc ország határolja.*

(42) *[Szihalom<sub>LOC</sub>] község [Heves megye<sub>LOC</sub>] [Füzesabonyi kistérségében<sub>LOC</sub>].*

(43) *A [Bükk Nemzeti Park<sub>LOC</sub>] mintegy 95 százalékát erdő borítja.*

(44) *[Gatwick<sub>LOC</sub>] délre, [Stansted<sub>LOC</sub>] észak-keletre, [Luton<sub>LOC</sub>] észak-nyugatra fekszik [Londontól<sub>LOC</sub>].*

(45) *Platón dialógusaiban részletesen szól [Atlantisz<sub>LOC</sub>] szigetéről.*

Az utak nagy részének csak számozása van, de nincs neve – ezeket nem jelöljük. Névként jelöljük viszont azt, amikor szám helyett vagy mellett neve van az útnak, például:

(46) *Ünnepélyes keretek között tartották szerdán az 54. számú főút [M5<sub>LOC</sub>] autópálya és az 5. számú főút közötti szakasz kapacitásbővítésének átadását*

Ha azonban az út neve után -s melléknévképző áll, akkor nem annotálandó névként, például:

(47) *Kisteherautó égett az M5-ösön*

A szervezeteknek sajátjuk, hogy van székhelyük, és előfordul, hogy a szervezet nevét helymegjelölésként használjuk. Ilyenkor a *tag-for-meaning* elv alapján a kontextusnak megfelelően helynévként annotáljuk őket, például:

(48) *tűz ütött ki a [Kapos Hotelben<sub>LOC</sub>]*

(49) *elbarikádozták magukat az [SZFE-n<sub>LOC</sub>]*

Ha egy geopolitikai entitás (birodalom, ország, megye, város stb.) aktívan cselekszik, valamilyen politikai cselekedeteket hajt végre, mint például megtámad egy másikat, vagy döntést hoz, elítél, nyilatkozik stb., akkor **ORG**-ként annotálandó. Ide tartoznak az ország- vagy városnévvel említett sportcsapatok is, például:

(50) *a [Spanyolország<sub>ORG</sub>]-[Németország<sub>ORG</sub>] összecsapást koordinálta*

(51) *a [Habsburg Birodalom<sub>ORG</sub>] semleges maradt*

(52) *Az esetek 90%-ában [Brüsszel<sub>ORG</sub>] javára dönt ez a bíróság.*

Ha viszont inkább földrajzi a kontextus, vagy egy geopolitikai entitás területén történnek az események, amiket az passzívan elszenved, akkor ezek is **LOC**-ként annotálandók:

(53) *a [Római Birodalom<sub>LOC</sub>] legvégső határára bukkantak*

(54) *[Magyarország<sub>LOC</sub>] közepe [Pusztavacs<sub>LOC</sub>] község területén van.*

(55) *benyomultak a rend területére és elfoglalták [Graudenz<sub>LOC</sub>]*

Ha egy földrajzi név után más nyelven, más írásmóddal is szerepel a név, akkor mindkettőt jelöljük **LOC**-ként, de külön névként, és a köztük levő zárójelet és egyéb részeket kihagyva, például:

(56) *[Bécs<sub>LOC</sub>] (németül: [Wien<sub>LOC</sub>], bajor nyelvjárás szerint: [Wean<sub>LOC</sub>])*

### 3.3.1. Összetett kifejezések

Az olyan összetett kifejezésekben, ahol földrajzi nevek vesszővel elválasztva szerepelnek, és a második név nagyobb helyre referál, tehát egyfajta pontosító funkciót tölt be, a neveket külön annotáljuk, például:

(57)  $[Los\ Angeles_{LOC}], [California_{LOC}]$

(58)  $[Budapest_{LOC}], [Magyarország_{LOC}]$

### 3.3.2. Köznévi tagok

Vannak olyan földrajzi nevek, melyek köznévi utótagot tartalmaznak. A közvetlenül a földrajzi név után álló köznévi frázisok, melyek nélkül a név nem ugyanaz lenne, a névvel együtt annotálандók, mint például az alábbiak:

(59)  $[Váci\ utca_{LOC}]$

(60)  $[Erzsébet\ híd_{LOC}]$

(61)  $[Baranya\ megye_{LOC}]$

(62)  $[Duna-Tisza\ köze_{LOC}]$

Nem tartoznak viszont a földrajzi névhez a magyarázó, pontosító funkciójú elemek, illetve az alkalmi jelzők sem, például:

(63)  $[Kent_{LOC}] \text{ grófság}$

(64)  $[New\ York_{LOC}] \text{ állam}$

(65)  $[Gyöngyös_{LOC}] \text{ város}$

(66)  $[Mátra_{LOC}] \text{ hegység}$

(67)  $[Duna_{LOC}] \text{ folyó}$

(68)  $az\ olasz\ [Alpok_{LOC}]$

(69)  $a\ lengyel\ [Magas-Tátra_{LOC}]$

(70)  $a\ gyönyörű\ [Alpok_{LOC}]$

(71)  $„Mit\ nekem\ te\ zordon\ [Kárpátoknak_{LOC}]...”$



### 3.4. Egyebek (MISC)

Ebbe a kategóriába kerülnek azok, amelyek NE-k, de a felsorolt kategóriák egyikébe sem illenek bele, mint a könyvcímek, újságnevek, konferencianevek, márkanév, tőzsdeindexek nevei, programozási nyelvek, díjak, kampányok, kereskedelmi útvonalak, gázvezetékek, rádiócsatornák nevei, például:

- (72) *A [Le Monde<sub>MISC</sub>] francia napilapban jelent meg, hogy*
- (73) *Az eset azonnal heves vitákat váltott ki a [Sina Weibón<sub>MISC</sub>], a meghatározó kínai közösségi portálon.*
- (74) *fedezze fel a [Fiat<sub>MISC</sub>] modelleket*
- (75) *Érdekel a [Python<sub>MISC</sub>] programozás?*
- (76) *Hitler megtiltotta, hogy német állampolgárok elfogadják a [Nobel-díjat<sub>MISC</sub>].*
- (77) *Újraindult a [Bringázz a munkába!<sub>MISC</sub>] kampány*
- (78) *így kereskedtek az ókorban a [Selyemúton<sub>MISC</sub>]*
- (79) *a [Déli Áramlat<sub>MISC</sub>] gázvezeték éves kapacitása 63 milliárd köbméter*
- (80) *a [Bartók Rádió<sub>MISC</sub>] szeptemberi koncertjei*

A törvények, rendeletek, irányelvek stb. esetében, ha csak számmal és/vagy köznévi elemekkel vannak jelölve, akkor nem jelöljük őket. De ha kaptak valamilyen egyedi nevet, akkor azt névként jelöljük, például:

- (81) *A [TEN-T<sub>MISC</sub>] rendelet (1315/2013 rendelet) 2013 végi hatályba lépését követően*
- (82) *ismételten változott a [Munka törvénykönyve<sub>MISC</sub>] (2012. évi I. tv.)*

Annak eldöntésében, hogy mi micsoda az EU-s jogi szövegekben, hasznos segítség lehet [ez a glosszárúm](#).

Ha egy újság cím úgy van használva, hogy arra referál, amit olvasunk, akkor az MISC. De ha a szerkesztőségre vagy a lapvezetésre utal, mint akik megírják vagy kiadják valamit, akkor az ORG-ként jelölendő.

Az együttesek, zenekarok, zeneszerzők nevét sokszor használjuk az általuk létrehozott mű helyett – ilyenkor a nevüket MISC-ként annotáljuk, például:

- (83) *[Bachot<sub>MISC</sub>] hallgatunk.*

Az események, bajnokságok, csaták MISC-ként jelölendők, ha tulajdonnévvel referálunk rájuk, például:

- (84) *[Waterloo<sub>MISC</sub>] még ma is ott él mindannyiunkban*

(85) *A [Forma-1<sub>MISC</sub>] legfrissebb hírei*

Ha viszont eseményekre köznevekkel van utalva, akkor azokat nem jelöljük névként, például:

(86) *2021-re halasztanák a 2020-as tokiói nyári olimpia megrendezését*

Ha egy nevet metanyelvi szinten használunk, vagyis magáról a névről szól a szöveg, akkor MISC-ként annotálandó, például:

(87) *Az elnevezés továbbra is [Tokió 2020<sub>MISC</sub>] lesz.*

(88) *Nekem rettentően tetszik fiúnévként az [Arion<sub>MISC</sub>]*

(89) *ugyanebben az évben kezdték használni a [Cryptopsy<sub>MISC</sub>] nevet*