

---

# Computational Linguistics-1

---

Summer 2021, IIIT Hyderabad

---

## Assignment 1

### Task A

---

#### Original Sentences

COVID-19 പ്രക്ഷേപണം തടയുന്നതിന് സംസ്ഥാന സർക്കാർ ജൂൺ 5 മുതൽ ജൂൺ 9 വരെ കർശനമായ ലോക്ക്ഡൗൺ ഏർപ്പെടുത്തും.

COVID-19 prakṣēpaṇam taṭayunnatin' samstha:na sarkka:r ju:ṇ 5 mutal ju:ṇ 9 vare karfanama:ya lo:kḍauṇ e:rpeṭuttum.

ഈ കാലയളവിൽ അവശ്യ സേവനങ്ങൾ മാത്രമേ ഭരണകൂടം അനുവദിക്കുകയുള്ളൂവെന്ന് മുഖ്യമന്ത്രി പിണറായി വിജയൻ വ്യാഴാഴ്ച ഉന്നതതല പാൻഡെമിക് അവലോകന യോഗത്തിൽ പറഞ്ഞു.

i: ka:laja[avi]l avafya se:vanan[ga]l ma:trame: bharaṇaku:ṭam anuvadikkukayul[u:venn' mukhyamantri piṇaRayi vijayan vya:zhazhca unnatatala pa:ṇḍemik avalo:kana yo:gattil paRaṇṇu.

ജീവിതം, മൊബിലിറ്റി, റീട്ടെയിൽ, വാണിജ്യം എന്നിവയ്ക്കുള്ള ഇളവ് ജൂൺ 4 ന് വൈകുന്നേരം 7 മണിക്ക് അവസാനിക്കും.

ji:vitam, mobiliRRi, Ri:ṭṭeyil, vaṇijyam ennivaykku[la i]lav' ju:ṇ na:lin' vaikunne:ram 7 maṇikk' avasa:nikkum.

അവശ്യസാധനങ്ങൾ, മരുന്നുകൾ, നിർമ്മാണ സാമഗ്രികൾ, പാക്കേജിംഗ്, വ്യവസായിക ഉൽപാദനത്തിനുള്ള അസംസ്കൃത വസ്തുക്കൾ എന്നിവ വിൽക്കുന്ന കടകൾക്ക് മാത്രമേ ബിസിനസിന് തുറക്കാൻ കഴിയൂ.

avafyasa:dhanan[ga]l, marunnuka[ḷ], nirmma:ṇa sa:magrika[ḷ], pa:kke:jim[ga]', vya:vasa:yika ulpa:danattinu[ḷ]a asamskrta vastukka[ḷ] enniva vilkkunna kaṭaka[ḷ]k' ma:trame: bisinasin' turakka:n kazhiyu:.

സ്ക്രാപ്പ് കൈകാര്യം ചെയ്യുന്ന സ്റ്റോറുകൾക്ക് ജൂൺ 4 ന് വൈകുന്നേരം 7 മണി വരെ ബിസിനസ്സ് നടത്താം.

skra:pp' kaika:ryam ceyyunna sRRo:Ruka[ḷ]k' ju:ṇ na:lin' vaikunne:ram 7 maṇi vare bisinass' naṭatta:m.

സർക്കാരും അർദ്ധ സർക്കാരും, പൊതുമേഖലയും സ്ഥാപനങ്ങൾ, കോർപ്പറേഷനുകൾ, വിവിധ കമ്മീഷനുകൾ എന്നിവ അടച്ചിരിക്കും.

sarkka:Rum, arddha sarkka:Rum potume:khalayum stha:panan[ga]l, ko:rppaRe:ṣanuka[ḷ], vividha kammi:ṣanuka[ḷ] enniva aṭaccirikkum.

ജൂൺ 10 മുതൽ 50% ൽ താഴെ തൊഴിലാളികളുമായി അവയ്ക്ക് തുറക്കാൻ കഴിയും.

ju:ṇ 10 mutal 50% il tazhe tozhila:[i]ka[ḷ]uma:yi avaykk' turakka:n kazhiyum.

കേരളത്തിലേക്ക് പ്രവേശിക്കുന്ന ആളുകൾക്ക് COVID നെഗറ്റീവ് സർട്ടിഫിക്കറ്റ് ആവശ്യമാണ്.

ke:ra[ḷ]attile:kk' prave:fikkunna a:[u]ka[ḷ]k' COVID negaRRi:v' sarṭṭiphikkaRR' a:va:fama:ṇ'.

മൂന്നാമത്തെ തരംഗത്തിന് സംസ്ഥാനം സ്വയം ശ്രമിച്ചു.

mu:nnamatte taramgattin' samstha:nam svayam framiccu.

ആദ്യ ഘട്ടമെന്ന നിലയിൽ, ഗസ്റ്റ് വർക്കർമാരെയും മറ്റ് ദുർബല വിഭാഗങ്ങളെയും ഭരണകൂടം പൂർണ്ണമായും കുത്തിവയ്ക്കും.

a:dya ghaṭṭamenna nilayil, gasRR' varkkarma:reyum maRR' durbala vibha:gaṇṇaḷeyum bharaṇaku:ṭam pu:rṇama:yum kuttivaykkum.

---

## Tagged

COVID-19\_NN1 പ്രക്ഷേപണം\_NN1 തടയുന്നതിന്\_VVGN-IF സംസ്ഥാന\_JJ സർക്കാർ\_NN1 ജൂൺ\_NPM 5\_MC മുതൽ\_II ജൂൺ\_NPM 9\_MC വരെ\_II കർശനമായ\_JJ ലോക്ക്ഡൗൺ\_NN1 ഏർപ്പെടുത്തും\_VVF.

COVID-19\_NN1 prakṣe:paṇam\_NN1 taṭayunnatin'\_VVGN-IF samstha:na\_JJ sarkka:r\_NN1 ju:ṇ\_NPM 5\_MC mutal\_II ju:ṇ\_NPM 9\_MC vare\_II karṇanama:ya\_JJ lo:kḍaṇṇ\_NN1 e:rpeṭuttum\_VVF.

ഈ\_DDP കാലയളവിൽ\_NN1-II അവശ്യ\_JJ സേവനങ്ങൾ\_NN2 മാത്രമേ\_RR ഭരണകൂടം\_NN1 അനുവദിക്കുകയുള്ളൂ\_VVF-O എന്ന്\_CST മുഖ്യമന്ത്രി\_NN1 പിണറായി\_NP1 വിജയൻ\_NP1 വ്യാഴാഴ്ച\_NPD1 ഉന്നതതല\_JJ പാൻഡെമിക്\_NN1 അവലോകന\_NN1 യോഗത്തിൽ\_NN-II പറഞ്ഞു\_VVD.

i: DDP ka:layaḷavil\_NN1-II avaḥya\_JJ se:vananṇaḷ\_NN2 ma:trame: RR bharaṇaku:ṭam\_NN1 anuvadikkukayulḷu: VVF-O enn'\_CST mukhyamantri\_NN1 piṇaRayi\_NP1 vijayan\_NP1 vya:zha:zhca\_NPD1 unnatatala\_JJ pa:ṇḍemik\_NN1 avalo:kana\_NN1 yo:gattil\_NN-II paRaṇṇu\_VVD.

ജീവിതം\_NN1, മൊബിലിറ്റി\_NN1, റീട്ടെയിൽ\_NN1, വാണിജ്യം\_NN1 എന്നിവയ്ക്കുള്ള\_PH2NP-IF-EXJ ഇളവ്\_NN1 ജൂൺ\_NPM 4\_MC ന്\_II വൈകുന്നേരം\_NPD1 7\_MC മണിക്ക്\_NN1-II അവസാനിക്കും\_VVF. ji:vitam\_NN1, mobiliRRi\_NN1, Ri:ṭṭeyil\_NN1, vaṇijyam\_NN1 ennivaykkuḷḷa\_PH2NP-IF-EXJ iḷav'\_NN1 ju:ṇ\_NPM na:lin'\_MC-II vaikunne:ram\_NPD1 7\_MC maṇikk'\_NN1-II avasa:nikkum\_VVF.

അവശ്യസാധനങ്ങൾ\_NN2, മരുന്നുകൾ\_NN2, നിർമ്മാണ\_JJ സാമഗ്രികൾ\_NN2, പാക്കേജിംഗ്\_NN1, വ്യാവസായിക\_JJ ഉൽപാദനത്തിനുള്ള\_NN1-IF-EXJ അസംസ്കൃത\_JJ വസ്തുക്കൾ\_NN2 എന്നിവ\_PH2NP വിൽക്കുന്ന\_VVGJ കടകൾക്ക്\_NN2-IF മാത്രമേ\_RR ബിസിനസിന്\_NN1-IF തുറക്കാൻ\_VDI കഴിയൂ\_VVF-O.

avaḥyasa:dhananṇaḷ\_NN2, marunnukaḷ\_NN2, nirmma:ṇa\_JJ sa:magrikaḷ\_NN2, pa:kke:jimg'\_NN1, vya:vasa:yika\_JJ ulpa:danattinulḷa\_NN1-IF-EXJ asamskrta\_JJ vastukkaḷ\_NN2 enniva\_PH2NP vilkkunna\_VVGJ kaṭakaḷkk'\_NN2-IF ma:trame: RR bisinasin'\_NN1-IF turakka:n\_VDI kazhiyu: VVF-O.

സ്ക്രാപ്പ്\_NN1 കൈകാര്യം\_NN1 ചെയ്യുന്ന\_VVGJ സ്റ്റോറുകൾക്ക്\_NN2-IF ജൂൺ\_NPM 4\_MC ന്-IF വൈകുന്നേരം\_NPD1 7\_MC മണി\_NPD1 വരെ\_II ബിസിനസ്സ്\_NN1 നടത്താം\_VVP. skra:pp'\_NN1 kaika:ryam\_NN1 ceyyunna\_VVGJ sRRo:Rukaḷkk'\_NN2-IF ju:ṇ\_NPM na:lin'\_MC-IF vaikunne:ram\_NPD1 7\_MC maṇi\_NPD1 vare\_II bisinass'\_NN1 naṭatta:m\_VVP.

സർക്കാരും\_NN1-CC അർദ്ധ\_JJ സർക്കാരും\_NN1-CC, പൊതുമേഖലയും\_NN1-CC സ്ഥാപനങ്ങൾ\_NN2, കോർപ്പറേഷനുകൾ\_NN2, വിവിധ\_JJ കമ്മീഷനുകൾ\_NN2 എന്നിവ\_PH2NP അടച്ചിരിക്കും\_VVF. sarkka:Rum\_NN1-CC, arddha\_JJ sarkka:Rum\_NN1-CC potume:khalayum\_NN1-CC stha:pananṇaḷ\_NN2, korppaRe:ṣanukaḷ\_NN2, vividha\_JJ kammi:ṣanukaḷ\_NN2 enniva\_PH2NP aṭaccirikkum\_VVF.

ജൂൺ\_NPM 10\_MC മുതൽ\_II 50%\_DB ൽ\_II താഴെ\_RL തൊഴിലാളികളുമായി\_NN2-II അവയ്ക്ക്\_PH2NR തുറക്കാൻ\_VDI കഴിയും\_VVF.

ju:ṇ\_NPM 10\_MC mutal\_II 50%\_DB il-II ta:zhe\_RL tozhila:ḷikaḷuma:yi\_NN2-II avaykk'\_PH2NR turakka:n\_VDI kazhiyum\_VVF.

കേരളത്തിലേക്ക്\_NP1-II പ്രവേശിക്കുന്ന\_VVGJ ആളുകൾക്ക്\_NN2-IF COVID\_NN1 നെഗറ്റീവ്\_JJ സർട്ടിഫിക്കറ്റ്\_NN1 ആവശ്യമാണ്\_JJ-VBZ.  
ke:ra|attile:kk'\_NP1-II prave:řikkunna\_VVGJ a:|uka|kk'\_NN2-IF COVID\_NN1 negaRRi:v'\_JJ sar|t|iphikkaRR'\_NN1 a:va|yama:ŋ'\_JJ-VBZ.

മൂന്നാമത്തെ\_MD തരംഗത്തിന്\_NN1-IF സംസ്ഥാനം\_NN1 സ്വയം\_PPX ശ്രമിച്ചു\_VVD.  
mu:nnamatte\_MD taramgattin'\_NN1-IF samstha:nam\_NN1 svayam\_PPX řramiccu\_VVD.

ആദ്യ\_JJ ഘട്ടം\_NN1 എന്ന\_II നിലയിൽ\_II, ഗസ്റ്റ്\_NN1 വർക്കർമാരെയും\_NN2-II-CC മറ്റ്\_JJ ദുർബല\_JJ വിഭാഗങ്ങളെയും\_NN2-II-CC ഭരണകൂടം\_NN1 പൂർണ്ണമായും\_RR കുത്തിവയ്ക്കും\_VVF.  
a:dya\_JJ gha|řamenna\_NN1 nilayil\_II, gasRR'\_NN1 varkkarma:reyum\_NN2-II-CC maRR'\_JJ durbala\_JJ vibha:gaŋŋa|eyum\_NN2-II-CC bharaŋaku:řam\_NN1 pu:řŋama:yum\_RR kuttivaykkum\_VVF.

---

## Extensions to CLAWS7 Tagset

The following are the extensions that were used in the above corpus, based on which a more complete extension is proposed in task B:

Verbs:

VVF: simple future/present habitual tense (suffix -um “-ഉം”)

VVC: present continuous (suffix -unnu “-ഉന്നു”)

VVP: potential (suffix -a:m “-ആം”)

VVGN: present gerund (noun)

VVGJ: present gerund (adjective)

EXJ: existential gerund (adjective)

-O: “only” (suffix -u: “-ഊ”; accompanies adverb ma:trame: “മാത്രമേ”)

Conjunctions:

CST: quotative particle enn' “എന്ന്” (used as “that” in English)

Pronouns:

PH2NP: neuter 3rd person pronoun plural proximate

PH2NR: neuter 3rd person impersonal pronoun plural remote

PPX: reflexive pronoun (svayam “സ്വയം”)

Note: Tags glossed as “prepositions” in CLAWS7 are treated as postpositional tags (for case markers) here.

---

## Notes on Transliteration

Velar, dental and bilabial non-nasal stops are as in IAST and common transliteration.

The velar nasal is ŋ; the dental nasal is ɳ; and the bilabial nasal is m.

The palatal consonants are c, ch, j, jh, ɟ.

The retroflex consonants are ɭ, ɭh, ɖ, ɖh, ɻ.

The semivowels common to Malayalam and Devanagari are transliterated the common way (y, r, l, v).

The sibilants are respectively ř (retroflex: Devanagari ढ and Telugu ష), ř (palatal: Devanagari श and Telugu శ), s and h.

The retroflex lateral approximant is ɭ (Devanagari ञ and Telugu ఙ) and the voiced retroflex approximant is zh.

The alveolar trill is R (as opposed to the alveolar tap which is r).

RR is a direct transliteration of the letter for R, doubled; it is pronounced as an alveolar voiceless stop (distinct from the retroflex and dental voiceless stops).

m occurring at the end of words or immediately before consonants represents the equivalent of Devanagari ँ (अं) and Telugu ు (అం).

r occurring between two consonants represents the syllabic r (Devanagari ऋ and Telugu ఋ).

' (apostrophe) at the end of words represents the close central unrounded vowel ɪ. In the script it is represented by either the equivalent of the halant or the symbol for 'u'. It does not have an independent vowel sign (without a consonant).

Long vowels are marked by the : symbol.

## Task B

---

### Question 1

No, the tags are not sufficient, since they are designed for English, an analytic language. For example, many tags which ought to be for verb forms in a synthetic language like Malayalam are given over to modals in English.

---

### Question 2

The redundant groups of tags are:

- BCL: before-clause marker (e.g. "in order to"); Malayalam makes use of verb forms (usually non-finite) for these functions.
- CS: subordinating conjunction (e.g. if, because); see BCL.
- DD2: plural determiner (e.g. "these"); determiners have no distinction based on number.
- GE: Germanic genitive marker (' or 's): Malayalam has only one genitive marker, the case marker.
- IF: for (as preposition): The other uses of "for" have different translations. There is no confusion.
- IO: of (as preposition): see IF.
- IW: with, without (as preposition): see IF.
- JJR: general comparative adjective (e.g. "older"): The adjective is not modified; when "than" (equivalent) is present, it is understood.
- PNX1: reflexive indefinite pronoun (oneself): There is no equivalent of "one".
- RA: adjective after nominal head (else, galore): There are no such adjectives.
- RPK: prep. adv., catenative ("about" in "be about to"): This verb form is not a catenative formation in Malayalam.
- RRR: comparative general adverb (e.g. "better"): see JJR.
- TO: infinitive marker (to): There is a specific verb form called the infinitive.
- VM: modal auxiliary (can, will, would): Malayalam uses inflections rather than modals.
- VMK: modal catenative (ought, used): see VM.
- VN: past participle (e.g. given): Malayalam has no equivalent to the past participle.

---

### Question 3

The extra tags needed are:

VVF: simple future/present habitual tense (suffix -um “-ഉം”)

VVAP: present perfect (suffix -ittuṇḍ’ “-ഇട്ടുണ്ട്”)

VVC: present continuous (suffix -unnu “-ഉന്നു” or -ukayaṇ “-ഉക്തയാണു”)

VVP: potential (suffix -a:m “-ആം”)  
VVI: conditional (suffix -a:l “-ആൽ” or -eŋkil “-എങ്കിൽ”)  
VVO: interrogative (suffix -o: “-ഓ”)

VVCVI: primary causative infinitive “to cause to do” (variable suffix)  
VVCV2I: secondary causative infinitive “to cause to cause to do” (variable suffix)  
VVPVI: passive voice infinitive “to be done” (suffix -peṭuka “-പെടുക”)  
VVPCVI: passive causative infinitive “to cause to be done” (suffix -peṭuttuka “-പെടുത്തുക”)  
VVC1PVI: primary causative passive infinitive “to be caused to be done” (variable suffix)  
VVC2PVI: secondary causative passive infinitive “to be caused to be done” (variable suffix)  
[in this manner, each new combination of tense, aspect, mood and voice will need a new tag; this is not an exhaustive list]

VVGJ: present continuous gerund adjective (suffix -unna “-ഉന്ന”)  
VVGJ: present continuous gerund noun (suffix -unnat’ “-ഉന്നത്”)  
[one gerund adjective and one gerund noun for any combination of tense, aspect, mood and voice]

[analogously, all verb tags for the existential uṇḍavuka “ഉണ്ടാവുക” and for the copula a:vuka “ആവുക”, which are irregular]

O: attached to verbs: “only”(-u: “-ഉ”)

PPIS2I: inclusive 1st person plural pronoun (nammaḷ “നമ്മൾ”)  
PPIS2E: exclusive 1st person plural pronoun (naṇṇaḷ “ഞങ്ങൾ”)

PPY1: 2nd personal casual singular pronoun (ta:n “താൻ”)  
PPY2: 2nd person informal singular pronoun (ni: “നീ”)  
PPY3: 2nd person informal plural/formal singular pronoun (niṇṇaḷ “നിങ്ങൾ”)  
PPY4: 2nd person formal singular/plural pronoun (for royalty) (ta:ŋkaḷ “താങ്കൾ”)

PH2NP: neuter 3rd person pronoun plural proximate  
PH2NR: neuter 3rd person pronoun plural remote  
PHF: 3rd person formal pronoun (adde:ham “അദ്ദേഹം”)

PPX: reflexive pronoun (svayam “സ്വയം”)

[for these and all existing pronominal tags, corresponding tags for all 7 cases; the accusative and genitive cases are covered in the existing tagset.]

Because there are so many TAM combinations and cases, the tagset might include an unwieldy number of tags for a synthetic language like Malayalam. The estimated number for a given verb is 144 according to the above categorisation; separate tags for the existential and the copula would more than triple the current size, which is 137.

It might therefore be more efficient to assign tags only to morphemes rather than entire words; thus case and TAM markers would be marked separately from their roots. However, tokenisation would become significantly more complicated. Morpheme boundaries are fuzzy at best and non-existent at worst. There would have to be a tradeoff; possibly some of the categories could be tagged morpheme-wise and the remaining could have lexical tags (e.g. the tagset might include tags for all combinations of tense and aspect, while the morphemes for mood and voice are tagged separately).