

Appendix

Appendix A: spaCy Part of Speech Tag List

- ADJ: adjective
- ADP: adposition
- ADV: adverb
- AUX: auxiliary verb
- CONJ: coordinating conjunction
- DET: determiner
- INTJ: interjection
- NOUN: noun
- NUM: numeral
- PART: particle
- PRON: pronoun
- PROPN: proper noun
- PUNCT: punctuation
- SCONJ: subordinating conjunction
- SYM: symbol
- VERB: verb
- X: other

Appendix B: NLTK Part of Speech Tag List

\$: dollar

\$ -\$ --\$ A\$ C\$ HK\$ M\$ NZ\$ S\$ U.S.\$ US\$

": closing quotation mark

' "

(: opening parenthesis

([{

): closing parenthesis

)] }

,: comma

,

--: dash

--

.: sentence terminator

. ! ?

:: colon or ellipsis

: ; ...

CC: conjunction, coordinating

& 'n and both but either et for less minus neither nor or plus so
therefore times v. versus vs. whether yet

CD: numeral, cardinal

mid-1890 nine-thirty forty-two one-tenth ten million 0.5 one forty-
seven 1987 twenty '79 zero two 78-degrees eighty-four IX '60s .025
fifteen 271,124 dozen quintillion DM2,000 ...

DT: determiner

all an another any both del each either every half la many much nary
neither no some such that the them these this those

EX: existential there

there

FW: foreign word

gemeinschaft hund ich jeux habeas Haementeria Herr K'ang-si vous
lutihaw alai je jour objets salutaris fille quibusdam pas trop Monte
terram fiche oui corporis ...

IN: preposition or conjunction, subordinating

astride among uppon whether out inside pro despite on by throughout
below within for towards near behind atop around if like until below
next into if beside ...

JJ: adjective or numeral, ordinal

third ill-mannered pre-war regrettable oiled calamitous first separable
ectoplasmic battery-powered participatory fourth still-to-be-named
multilingual multi-disciplinary ...

JJR: adjective, comparative

bleaker braver breezier briefer brighter brisker broader bumper busier
calmer cheaper choosier cleaner clearer closer colder commoner costlier
cozier creamier crunchier cuter ...

JJS: adjective, superlative

calmest cheapest choicest classiest cleanest clearest closest commonest
corniest costliest crassest creepiest crudest cutest darkest deadliest
dearest deepest densest dinkiest ...

LS: list item marker

A A. B B. C C. D E F First G H I J K One SP-44001 SP-44002 SP-44005
SP-44007 Second Third Three Two * a b c d first five four one six three
two

MD: modal auxiliary

can cannot could couldn't dare may might must need ought shall should
shouldn't will would

NN: noun, common, singular or mass

common-carrier cabbage knuckle-duster Casino afghan shed thermostat
investment slide humour falloff slick wind hyena override subhumanity
machinist ...

NNP: noun, proper, singular

Motown Venneboerger Czystochwa Ranzer Conchita Trumplane Christos
Oceanside Escobar Kreisler Sawyer Cougar Yvette Ervin ODI Darryl CTCA
Shannon A.K.C. Meltex Liverpool ...

NNPS: noun, proper, plural

Americans Americas Amharas Amityvilles Amusements Anarcho-Syndicalists
Andalusians Andes Andruses Angels Animals Anthony Antilles Antiques
Apache Apaches Apocrypha ...

NNS: noun, common, plural

undergraduates scotches bric-a-brac products bodyguards facets coasts
divestitures storehouses designs clubs fragrances averages
subjectivists apprehensions muses factory-jobs ...

PDT: pre-determiner

all both half many quite such sure this

POS: genitive marker

's

PRP: pronoun, personal

hers herself him himself himself it itself me myself one oneself ours
ourselves oneself self she thee theirs them themselves they thou thy us

PRP\$: pronoun, possessive

her his mine my our ours their thy your

RB: adverb

occasionally unabatingly maddeningly adventurously professedly
stirringly prominently technologically magisterially predominately
swiftly fiscally pitilessly ...

RBR: adverb, comparative

further gloomier grander graver greater grimmer harder harsher
healthier heavier higher however larger later leaner lengthier less-
perfectly lesser lonelier longer louder lower more ...

RBS: adverb, superlative

best biggest bluntest earliest farthest first furthest hardest
heartiest highest largest least less most nearest second tightest worst

RP: particle

aboard about across along apart around aside at away back before behind
by crop down ever fast for forth from go high i.e. in into just later
low more off on open out over per pie raising start teeth that through
under unto up up-pp upon whole with you

SYM: symbol

% & ' " " .)). * + , . < = > @ A[fj] U.S U.S.S.R * * * * *

TO: "to" as preposition or infinitive marker

to

UH: interjection

Goodbye Goody Gosh Wow Jeepers Jee-sus Hubba Hey Kee-reist Oops amen
huh howdy uh dammit whammo shucks heck anyways whodunnit honey golly
man baby diddle hush sonuvabitch ...

VB: verb, base form

ask assemble assess assign assume atone attention avoid bake balkanize
bank begin behold believe bend benefit bevel beware bless boil bomb
boost brace break bring broil brush build ...

VBD: verb, past tense

dipped pleaded swiped regummed soaked tidied convened halted registered
cushioned exacted snubbed strode aimed adopted belied figgered
speculated wore appreciated contemplated ...

VBG: verb, present participle or gerund

telegraphing stirring focusing angering judging stalling lactating

hankerin' alleging veering capping approaching traveling besieging
encrypting interrupting erasing wincing ...

VCN: verb, past participle

multihulled dilapidated aerosolized chaired languished panelized used
experimented flourished imitated reunified factored condensed sheared
unsettled primed dubbed desired ...

VBP: verb, present tense, not 3rd person singular

predominate wrap resort sue twist spill cure lengthen brush terminate
appear tend stray glisten obtain comprise detest tease attract
emphasize mold postpone sever return wag ...

VBZ: verb, present tense, 3rd person singular

bases reconstructs marks mixes displeases seals carps weaves snatches
slumps stretches authorizes smolders pictures emerges stockpiles
seduces fizzes uses bolsters slaps speaks pleads ...

WDT: WH-determiner

that what whatever which whichever

WP: WH-pronoun

that what whatever whatsoever which who whom whosoever

WP\$: WH-pronoun, possessive

whose

WRB: Wh-adverb

how however whence whenever where whereby wherever wherein whereof why

``: opening quotation mark

``

Appendix C: Stopword List

'm	almost	aren	beforehand
're	alone	aren'	begin
's	along	arent	beginning
've	already	arise	beginnings
a	also	around	begins
able	although	as	behind
about	always	aside	being
above	am	ask	believe
absolutely	amazing	asking	below
accordance	among	at	beside
according	amongst	auth	besides
accordingly	an	available	best
across	and	away	better
act	announce	awesome	between
actually	another	awful	beyond
added	any	awfully	big
adj	anybody	back	bigger
affected	anyhow	bad	biggest
affecting	anymore	basically	biol
affects	anyone	be	bit
after	anything	became	both
afterwards	anyway	because	brief
again	anyways	become	briefly
against	anywhere	becomes	but
ago	apparently	becoming	by
ah	approximately	been	ca
all	are	before	came

can	didn'	elsewhere	first
can'	didn't	end	five
can't	different	ending	fix
cannot	do	enough	followed
cause	does	entire	following
causes	doesn	especially	follows
certain	doesn'	et	for
certainly	doesn't	et-al	forever
clearly	doing	etc	former
co	don	even	formerly
com	don't	ever	forth
come	done	every	found
comes	dont	everybody	four
constantly	down	everyone	free
contain	downwards	everything	from
containing	due	everywhere	further
contains	during	ex	furthermore
could	each	excellent	gave
couldn	easi	except	get
couldn'	easier	extremely	gets
couldnt	easy	fantastic	getting
currently	ed	far	give
date	edu	feel	given
day	effect	feeling	gives
days	eg	few	giving
definitely	eight	ff	glad
despite	eighty	fifth	go
did	either	finally	goes
didn	else	fine	going

gone	hereupon	important	kind
good	hers	in	km
goodbye	herself	inc	know
got	hes	incredible	known
gotten	hi	incredibly	knows
great	hid	indeed	largely
greatest	him	index	last
had	himself	information	lately
half	his	initially	later
happens	hither	instead	latter
happy	home	into	latterly
hardly	hope	invention	least
has	hour	inward	less
hasn't	hours	is	lest
have	how	isn	let
haven	howbeit	isn'	lets
haven'	however	isn't	like
haven't	huge	issue	liked
having	hundred	issues	likely
he	i	it	line
hed	i'll	it'll	literally
hello	i've	itd	little
hence	id	its	ll
her	ie	itself	look
here	if	just	looking
hereafter	im	keep	looks
hereby	immediate	keeps	lot
herein	immediately	kept	lots
heres	importance	kg	love

loved	mostly	non	or
ltd	mr	none	ord
luck	mrs	nonetheless	other
made	much	noone	others
mainly	mug	nor	otherwise
make	must	normally	ought
makes	my	nos	our
many	myself	not	ours
matter	n't	noted	ourselves
may	na	nothing	out
maybe	name	now	outside
me	namely	nowhere	over
mean	nay	obtain	overall
means	nd	obtained	owing
meantime	near	obviously	own
meanwhile	nearly	of	page
merely	necessarily	off	pages
mg	necessary	often	part
might	need	oh	particular
million	needs	ok	particularly
minute	neither	okay	past
minutes	never	old	per
miss	nevertheless	omitted	perfect
ml	new	on	perhaps
month	next	once	placed
months	nine	one	please
more	ninety	ones	plus
moreover	no	only	poor
most	nobody	onto	poorly

possible	refs	seeing	slightly
possibly	regarding	seem	so
potentially	regardless	seemed	some
pp	regards	seeming	somebody
predominantly	related	seems	somehow
present	relatively	seen	someone
pretty	research	self	somethan
previously	respectively	selves	something
primarily	resulted	sent	sometime
probably	resulting	seriously	sometimes
problem	results	seven	somewhat
problems	review	several	somewhere
promptly	reviews	shall	soon
proud	right	she	sorry
provides	rock	she'll	specifically
put	rocks	shed	specified
que	run	shes	specify
quickly	s	should	specifying
quite	said	shouldn't	still
qv	same	show	stop
ran	saw	showed	strongly
rather	say	shown	stuff
rd	saying	shows	sub
re	says	significant	substantially
readily	sec	significantly	successfully
really	second	similar	such
recent	seconds	similarly	suck
recently	section	since	sucks
ref	see	six	suddenly

sufficiently	there'll	thus	unto
suggest	there've	til	up
sup	thereafter	till	upon
super	thereby	time	ups
sure	thered	times	us
t	therefore	tip	use
take	therein	to	used
taken	thereof	together	useful
taking	theres	too	usefully
tell	thereto	took	usefulness
tends	thereupon	totally	uses
terrible	these	toward	using
th	they	towards	usually
than	they'll	tried	various
thank	they've	tries	ve
thanks	theyd	true	veri
thanx	theyre	truly	very
that	thing	try	via
that'll	things	trying	viz
that've	think	ts	vol
thats	this	twice	vols
the	those	two	vs
their	thou	un	want
theirs	though	under	wants
them	though	unfortunately	was
themselves	thousand	unless	wasn
then	through	unlike	wasn'
thence	throughout	unlikely	wasnt
there	thru	until	way

we	whereafter	whos	wow
we'll	whereas	whose	wrong
we've	whereby	why	www
wed	wherein	widely	year
week	wheres	will	years
weeks	whereupon	willing	yes
welcome	wherever	wish	yet
well	whether	with	you
went	which	within	you'll
were	while	without	you've
werent	whim	wont	youd
what	whither	words	your
what'll	who	world	youre
whatever	who'll	worse	yours
whats	whod	worst	yourself
when	whoever	would	yourselves
whence	whole	wouldn	zero
whenever	whom	wouldn'	
where	whomever	wouldnt	

Appendix D: spaCy Named Entity Tags

spaCy trains two types of models for Named Entity Recognition. Models trained on the OntoNotes 5 corpus have the following tags:

Type	Description
PERSON	People, including fictional.
NORP	Nationalities or religious or political groups.
FAC	Buildings, airports, highways, bridges, etc.
ORG	Companies, agencies, institutions, etc.
GPE	Countries, cities, states.
LOC	Non-GPE locations, mountain ranges, bodies of water.
PRODUCT	Objects, vehicles, foods, etc. (Not services.)
EVENT	Named hurricanes, battles, wars, sports events, etc.
WORK_OF_ART	Titles of books, songs, etc.
LAW	Named documents made into laws.
LANGUAGE	Any named language.
DATE	Absolute or relative dates or periods.
TIME	Times smaller than a day.
PERCENT	Percentage, including "%".
MONEY	Monetary values, including unit.
QUANTITY	Measurements, as of weight or distance.
ORDINAL	"first", "second", etc.
CARDINAL	Numerals that do not fall under another type.

Models trained on Wikipedia data have the following scheme:

Type	Description
PER	Named person or family.
LOC	Name of politically or geographically defined location (cities, provinces, countries, international regions, bodies of water, mountains).
ORG	Named corporate, governmental, or other organizational entity.
MISC	Miscellaneous entities, e.g. events, nationalities, products or works of art.