

2 Tagset

(Author: Adam Przepiórkowski; modified: 2 October 2011)

Each morphosyntactic tag is a sequence of colon-separated values, e.g.: *subst:sg:nom:m1* for the segment *chłopiec* ‘boy’. The first value, e.g., *subst*, determines the *grammatical class* (cf. §2.2), while the values that follow it, e.g., *sg*, *nom* and *m1*, are the values of grammatical categories (cf. §2.1) appropriate for that grammatical class.

2.1 Grammatical categories

The following table presents the repertoire of grammatical categories used in the National Corpus of Polish:

Number: (2 values)		
singular	<i>sg</i>	<i>oko</i>
plural	<i>pl</i>	<i>oczy</i>
Case: (7 values)		
nominative	<i>nom</i>	<i>woda</i>
genitive	<i>gen</i>	<i>wody</i>
dative	<i>dat</i>	<i>wodzie</i>
accusative	<i>acc</i>	<i>wodę</i>
instrumental	<i>inst</i>	<i>wodą</i>
locative	<i>loc</i>	<i>wodzie</i>
vocative	<i>voc</i>	<i>wodo</i>
Gender: (5 values)		
human masculine (virile)	<i>m1</i>	<i>papież, kto, wujostwo</i>
animate masculine	<i>m2</i>	<i>baranek, walc, babsztyl</i>
inanimate masculine	<i>m3</i>	<i>stół</i>
feminine	<i>f</i>	<i>stula</i>
neuter	<i>n</i>	<i>dziecko, okno, co, skrzypce, spodnie</i>
Person: (3 values)		
first	<i>pri</i>	<i>bredzę, my</i>
second	<i>sec</i>	<i>bredzisz, wy</i>
third	<i>ter</i>	<i>bredzi, oni</i>
Degree: (3 values)		
positive	<i>pos</i>	<i>cudny</i>
comparative	<i>com</i>	<i>cudniejszy</i>
superlative	<i>sup</i>	<i>najcudniejszy</i>
Aspect: (2 values)		
imperfective	<i>imperf</i>	<i>iść</i>
perfective	<i>perf</i>	<i>zajść</i>
Negation: (2 values)		
affirmative	<i>aff</i>	<i>pisanie, czytaniego</i>
negative	<i>neg</i>	<i>niepisanie, nieczytaniego</i>
Accentability: (2 values)		
accented (strong)	<i>akc</i>	<i>jego, niego, tobie</i>
non-accented (weak)	<i>nakc</i>	<i>go, -ń, ci</i>
Post-prepositional: (2 values)		
post-prepositional	<i>praep</i>	<i>niego, -ń</i>
non-post-prepositional	<i>npraep</i>	<i>jego, go</i>
Accommodability: (2 values)		
agreeing	<i>congr</i>	<i>dwaj, pięcioma</i>
governing	<i>rec</i>	<i>dwóch, dwu, pięciorgiem</i>
Agglutination: (2 values)		
non-agglutinative	<i>nagl</i>	<i>niósł</i>
agglutinative	<i>agl</i>	<i>niosł-</i>
Vocalicity: (2 values)		
vocalic	<i>wok</i>	<i>-em</i>
non-vocalic	<i>nwok</i>	<i>-m</i>
Fullstoppedness: (2 values)		
with full stop	<i>pun</i>	<i>tzn</i>
without full stop	<i>npun</i>	<i>wg</i>

2.2 Grammatical classes

The scope of traditional parts of speech such as verb, noun, numeral or pronoun is fuzzy and, hence, controversial. For example, are gerundial forms such as *picie* ‘drinking’ and *palenie* ‘smoking’ verbs (they have the category of aspect and they are productively related to verbal forms such as *pić* ‘to drink’ and *palić* ‘to smoke’), or are they nouns (they decline for case, and they have the lexical category of gender)? Are ordinal numerals such as *piąty* ‘fifth’ numerals (semantically, they are numerals), or are they adjectives (they have adjectival inflection)? Are adjectival pronouns such as *taki* ‘such’ pronouns (semantics) or adjectives (inflection)?

Grammatical classes used in the National Corpus of Polish are more precisely delimited and, overall, finer-grained than traditional parts of speech. The classes assumed here are based on the notion of *flexeme*, narrower than the notion of *lexeme*.

The following table contains the rough morphosyntactic characteristics of all flexemic classes assumed in the present tagset. The symbol ⊕ in the table means that, for a given flexemic class, a given grammatical category is a morphological category (flexemes belonging to this class normally inflect for that category), while the symbol ⊙ means that the category is a lexical category (for each flexeme belonging to this class, all forms of that flexeme have the same value of that category, although that value may differ between flexemes, as in the case of the gender of nouns).

	number	case	gender	person	degree	aspect	negation	accentability	post-prep.	accom.	aggl.	vocalicity	fullstop.
noun	⊕	⊕	○										
depreciative form	○	⊕	○										
main numeral	○	⊕	⊕							⊕			
collective numeral	○	⊕	○							⊕			
adjective	⊕	⊕	⊕		⊕								
ad-adj. adjective													
post-prep. adjective													
predicative adjective													
adverb					⊕								
pronoun (non-3rd person)	○	⊕	⊕	○				⊕					
pronoun (3rd person)	⊕	⊕	⊕	○				⊕	⊕				
pronoun SIEBIE		⊕											
non-past form	⊕			⊕		○							
future BYĆ	⊕			⊕		○							
agglut. BYĆ	⊕			⊕		○						⊕	
l-participle	⊕		⊕			○					⊕		
imperative form	⊕			⊕		○							
impersonal form						○							
infinitive						○							
adv. contemp. prtcp.						○							
adv. anter. prtcp.						○							
gerund	⊕	⊕	○			○	⊕						
adj. act. prtcp.	⊕	⊕	⊕			○	⊕						
adj. pass. prtcp.	⊕	⊕	⊕			○	⊕						
winien-like verb	⊕		⊕			○							
predicative													
preposition		○											
coord. conjunction													
subord. conjunction													
particle-adverb													
abbreviation													⊕
bound word													
interjection													
punctuation													
alien													
unknown form													

The following table provides the information about base forms for all grammatical classes, as well as the abbreviations of these classes as used in the National Corpus of Polish.

flexeme	abbreviation	base form	example
noun	<i>subst</i>	singular nominative	<i>profesor</i>
depreciative form	<i>depr</i>	singular nominative form of the corresponding noun	<i>profesor</i>
main numeral	<i>num</i>	inanimate masculine nominative form	<i>pięć, dwa</i>
collective numeral	<i>numcol</i>	inanimate masculine nominative form of the main numeral	<i>pięć, dwa</i>
adjective	<i>adj</i>	singular nominative masculine positive form	<i>polski</i>
ad-adjectival adjective	<i>adja</i>	singular nominative masculine positive form of the adjective	<i>polski</i>
post-prepositional adjective	<i>adjp</i>	singular nominative masculine positive form of the adjective	<i>polski</i>
predicative adjective	<i>adjc</i>	singular nominative masculine positive form of the adjective	<i>zdrowy, ciekawy</i>
adverb	<i>adv</i>	positive form	<i>dobrze, bardzo</i>
non-3rd person pronoun	<i>ppron12</i>	singular nominative	<i>ja</i>
3rd-person pronoun	<i>ppron3</i>	singular nominative	<i>on</i>
pronoun SIEBIE	<i>siebie</i>	accusative	<i>siebie</i>
non-past form	<i>fin</i>	infinitive	<i>czytać</i>
future BYĆ	<i>bedzie</i>	infinitive	<i>być</i>
agglutinate BYĆ	<i>aglt</i>	infinitive	<i>być</i>
l-participle	<i>praet</i>	infinitive	<i>czytać</i>
imperative	<i>impt</i>	infinitive	<i>czytać</i>
impersonal	<i>imps</i>	infinitive	<i>czytać</i>
infinitive	<i>inf</i>	infinitive	<i>czytać</i>
contemporary adv. participle	<i>pcon</i>	infinitive	<i>czytać</i>
anterior adv. participle	<i>pant</i>	infinitive	<i>czytać</i>
gerund	<i>ger</i>	infinitive	<i>czytać</i>
active adj. participle	<i>pact</i>	infinitive	<i>czytać</i>
passive adj. participle	<i>ppas</i>	infinitive	<i>czytać</i>
winien	<i>winien</i>	singular masculine form	<i>powinien, rad</i>
predicative	<i>pred</i>	the only form of that flexeme	<i>warto</i>
preposition	<i>prep</i>	the non-vocalic form of that flexeme	<i>na, przez, w</i>
coordinating conjunction	<i>conj</i>	the only form of that flexeme	<i>oraz</i>
subordinating conjunction	<i>comp</i>	the only form of that flexeme	<i>że</i>
particle-adverb	<i>qub</i>	the only form of that flexeme	<i>nie, -że, się</i>
abbreviation	<i>brev</i>	the full dictionary form	<i>rok, i tak dalej</i>
bound word	<i>burk</i>	the only form of that flexeme	<i>trochu, oścież</i>
interjection	<i>interj</i>	the only form of that flexeme	<i>ech, kurde</i>
punctuation	<i>interp</i>	the only form of that flexeme	<i>;, ,, (,]</i>
alien	<i>xxx</i>	the only form of that flexeme	<i>cool, nihil</i>

unknown form	<i>ign</i>	the only form of that flexeme	
--------------	------------	-------------------------------	--