# Appendix A - Tagset and Conversions

Below is a list of all tags used in the final simplified tagset for the models, and lists of what tags from the source materials were converted to that tag in the simplified tagset. For the OHG Referenzkorpus Altdeutsch, the simpler 'posLemma' tags were used as the basis for these conversions where possible, and the more fine-grain 'pos' tags were used elsewhere.

1. **VB (verb)**

OE, OS, and OIce - All tags indicating the verb 'be', the verb 'have', auxiliary verbs, modal verbs, and other verbs in the documentation. (BE, BEI, BEPH, BEPI, BEPS, BEP, BEDI, BEDS, BED, BAG, BEN, HV, HVI, HVPI, HVPS, HVP, HVDI, HVDS, HVD, HAG, HVN, AX, AXI, AXPI, AXPS, AXP, AXDI, AXDS, AXD, AXG, AXN, MD, MDI, MDPI, MDPS, MDP, MDDI, MDDS, MDD, TO, VB, VBI, VBPH, VBPI, VBPS, VBP, VBDI, VBDS, VBD, VAG, VBN)

GO - 'Verb' and 'Participle'

OHG - All tags indicating verbs (''VA', 'VV', 'VM' 'VVPP', 'VVPS', 'VVFIN')

2. **N (noun)**

OE, OS, and OIce - All nominal tags (N, NR, MAN)

GO - 'Noun'

OHG - All tags indicating nominals ('NA', 'NE')

3. **PRO (pronoun)**

OE, OS, and OIce - All pronominal tags (PRO, PRO$)

GO - 'Pronoun'

OHG - All tags indicating pronominals ('PPER', 'PI', 'PRF')

4. **AP (adposition)**

OE, OS, and OIce - All prepositions (P)

GO - 'Preposition'

OHG - Tags for adposition ('AP')

### 5. D (determiner)

OE, OS, and OIce - All determiners (D)

GO - N/A

OHG - All tags indicating determiner forms ('DD', 'DPOS', 'DDS', 'DIS', 'DI')

### 6. ADV (adverb)

OE, OS, and OIce - all forms associated with adverbs (ADV, ADVR, ADVS)

GO - 'Adverb'

OHG - adverbs ('ADV')

### 7. CONJ (conjunction)

OE, OS, and OIce - All conjunctions (CONJ)

GO - 'Conjunction'

OHG - conjunctions ('KO')

### 8. C (complementizer)

OE, OS, and OIce - All complementizers (C)

GO - N/A

OHG - N/A

### 9. ADJ (adjective)

OE, OS, and OIce - all forms associated with adjectives (ADJ, ADJR, ADJS)

GO - 'Adjective'

OHG - Adjectives ('ADJ')

**10. Q (quantifier)**

<u>OE, OS, and OIce</u> - All quantifiers (Q, QR, QS)

<u>GO</u> - N/A

<u>OHG</u> - N/A

**11. NUM (numeral)**

<u>OE, OS, and OIce</u> - All numerals (NUM)

<u>GO</u> - 'Numeral'

<u>OHG</u> - Cardinal Numbers ('CARD')

**12. NEG (negation marker)**

<u>OE, OS, and OIce</u> - All negation (NEG)

<u>GO</u> - N/A

<u>OHG</u> - N/A

**13. PART (other particle)**

<u>OE, OS, and OIce</u> - All particles (FP, RP)

<u>GO</u> - 'Particle'

<u>OHG</u> - Particles ('PTK')

**14. WH (wh word)**

<u>OE, OS, and OIce</u> - All WH-words (WPRO, WADV, WADJ, WQ)

<u>GO</u> - N/A

<u>OHG</u> - Question words ('PW', 'PWG', 'DW')

**15. INTJ (interjection)**

<u>OE, OS, and OIce</u> - All interjections (INTJ)

<u>GO</u> - 'Interjection'

### 16. PUNCT (punctuation)

<u>OE, OS, and OIce</u> - All punctuation

<u>GO</u> - 'Punctuation'

<u>OHG</u> - All tags indicating some form of punctuation ('$,', '$.', '$(', '$\_')

### 17. XX (problematic)

<u>OE, OS, and OIce</u> - All tags marked as FW (foreign word) or not appearing in the source documentation

<u>GO</u> - 'nan', 'unassigned', 'foreign word'

<u>OE, OS, and OIce</u> - All tags marked as 'nan', empty tags, tags not occurring in the documentation.