

Morfologická segmentace se zabývá dělením slov na morfémy – nejmenší jednotky nesoucí význam. Jedná se o nízkoúrovňový problém z oblasti zpracování přirozeného jazyka. Jelikož se morfologická segmentace někdy používá jako metoda předzpracování dat, její zlepšení může pomoci algoritmům řešícím nejrůznější problémy z oblasti NLP, zejména, pokud v situaci, kdy je nedostatek dat. Zlepšení morfologické segmentace může také pomoci lingvistickému výzkumu, využívajícímu korpusy. V této práci navrhujeme nový ensemble algoritmus pro morfologickou segmentaci Českých lemmat, který používá derivační stromy z datasetu DeriNet. Zároveň vytváříme návrhy na zlepšení tohoto datasetu.