

Jazyky, jako je němčina, holandština, skandinávské jazyky nebo řečtina, na rozdíl od angličtiny netvoří kompozita jako víceslovné výrazy, ale spojením jednotlivých částí do nového slova bez ortografického oddělení. To představuje problém pro řadu úloh zpracování přirozeného jazyka, jako je například statistický strojový překlad či vyhledávání informací. Většina předešlých prací na téma rozkladu složenin na jejich částí se zaměřovala na němčinu. V této práci jsme vyvinuli nový jednoduchý systém založený na nařízeném strojovém učení pro automatický rozklad složenin pro tři reprezentativní kompozitní jazyky: němčina, švédština a maďarština. Součástí práce je vytvoření multilinguální evaluační datové sady z lékařské domény anotováním složenin získaných z korpusu EMEA a vyhodnocení několika variant našeho systému a srovnání s předchozími přístupy.