

Představujeme systém na opravu gramatických chyb v českém jazyce. Systém je založen na přístupu neuronového strojového překladu. Požíváme architekturu Transformer, která je závislá na velkém množství anotovaných dat. Vzhledem k tomu, že pro většinu jazyků včetně češtiny není k dispozici dostatek anotovaných dat, volíme syntetické generování dat. Do syntetických chyb zavádíme, jak chyby jednoduché, tak i složitější – typické české chyby. Pro usnadnění experimentování vyvíjíme systém schopný generovat data v reálném čase a rovnou na těchto datech trénovat model. Následně navrhujeme několik vylepšení, jako je převzorkování jazykových domén nebo výběr zdroje dat pro syntetické generování. Náš nejvýkonnější model dosahuje nejlepších výsledků v českém jazyce vůči modelům, které jsou srovnatelně velké. Implementace je zveřejněna na GitHub pod adresou: https://github.com/petrpechman/czech_gec/tree/MasterThesis_PechmanPetr_2024.