

Současný exponenciální nárůst genomických dat vyžaduje nové prostorově úsporné algoritmy pro jejich kompresi a vyhledávání. Moderní přístupy často místo původních dat využívají příslušných množin  $k$ -merů, což jsou podřetězce pevné délky  $k$ . Popularita metod založených na  $k$ -merech vedla k vzniku kompaktních textových reprezentací množin  $k$ -merů, jež však stojí na strukturálních předpokladech, které pro data v praxi nemusí platit. V této bakalářské práci ukážeme, že na všechny tyto reprezentace lze nahlížet jako na nadřetězce množin  $k$ -merů a jako takové je zobecníme pomocí uceleného konceptu, kterému říkáme maskované nadřetězce  $k$ -merů. Navrhujeme dva různé hladové algoritmy na jejich výpočet a implementujeme je v nástroji KmerCamel. Dále demonstrujeme, že maskované nadřetězce fungují jako stavební kámen pro nový a jednoduchý index pro množiny  $k$ -merů, který nazýváme FMS-index. Pokud k maskovaným nadřetězcům přiřadíme navíc odmaskovávající funkci  $f$ , výsledný koncept  $f$ -maskovaných nadřetězců umožňuje jednoduché provádění množinových operací s  $k$ -mery. Experimentálně ověříme prostorovou úspornost maskovaných nadřetězců, stejně tak i naší implementace FMS-indexu. Ukážeme, že maskované nadřetězce jsou lépe komprimovatelné v situacích, kde předchozí přístupy byly daleko od optima a že FMS-index je prostorově efektivnější než současné nejlepší přístupy k indexování. Naše výsledky dokládají užitečnost maskovaných nadřetězců jakožto sjednocujícího teoretického rámce a stejně tak i jejich potenciál v návrhu datových struktur pro  $k$ -mery.