

VŠB TECHNICKÁ
UNIVERZITA
OSTRAVA

VSB TECHNICAL
UNIVERSITY
OF OSTRAVA



www.vsb.cz

Kompresa stromových struktur

Semestrální projekt

Marek Beran

VŠB – Technická univerzita Ostrava

marek.beran.st@vsb.cz

23. května 2025

1 Úvod

2 Zpracování dat

3 Architektura

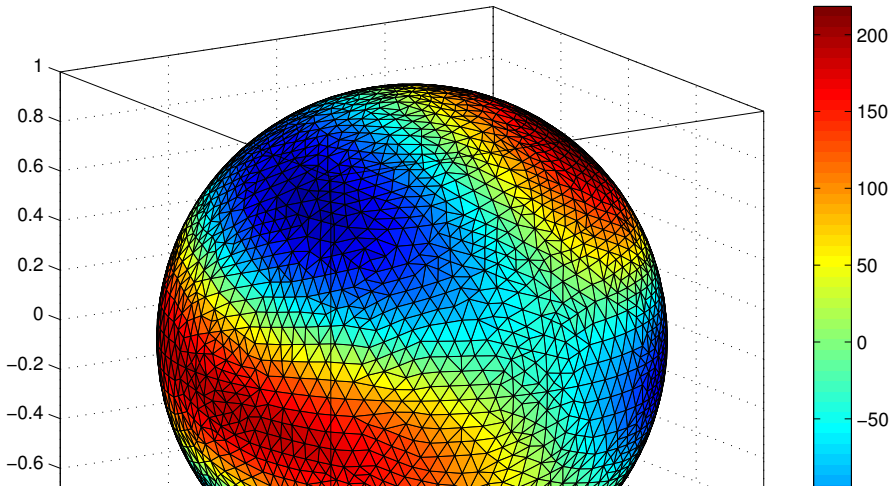
4 Algoritmy

5 Výsledky

Proč komprimovat stromy?



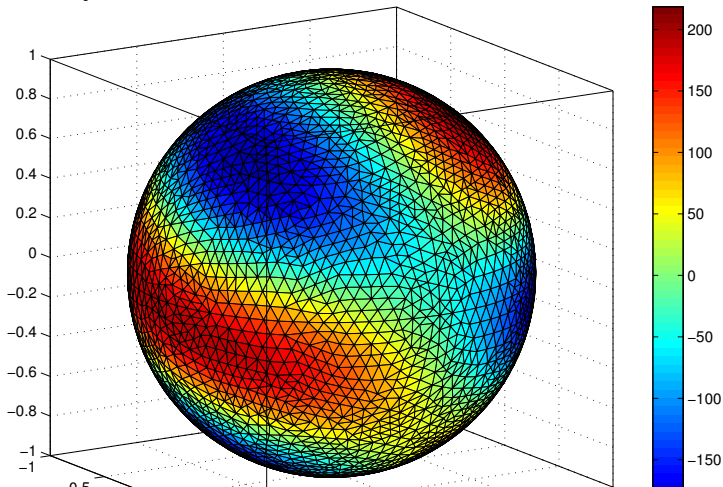
- Stromy = přirozený výstup syntaktické analýzy
- Text obsahuje opakující se vzorce
- Hypotéza: komprese možná díky opakovatelnosti



Cíle projektu



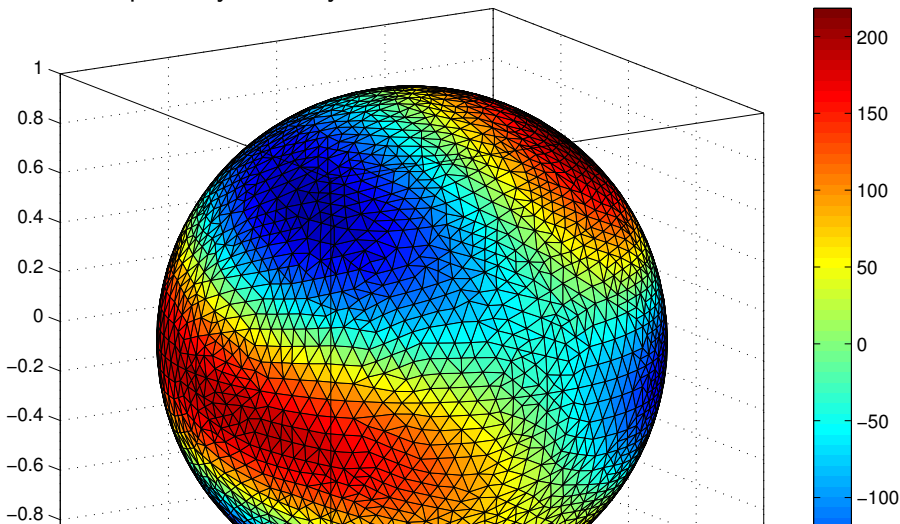
- Komprimovat syntaktické stromy
- Porovnat různé metody
- Vyhodnotit efektivitu



Z textu na strom

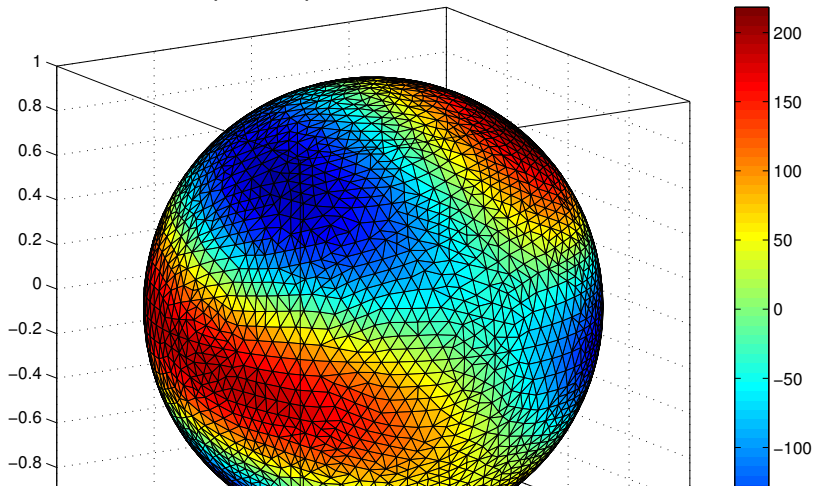


- MorphoDiTa → lemmatizace
- UDPipe → syntaktický strom



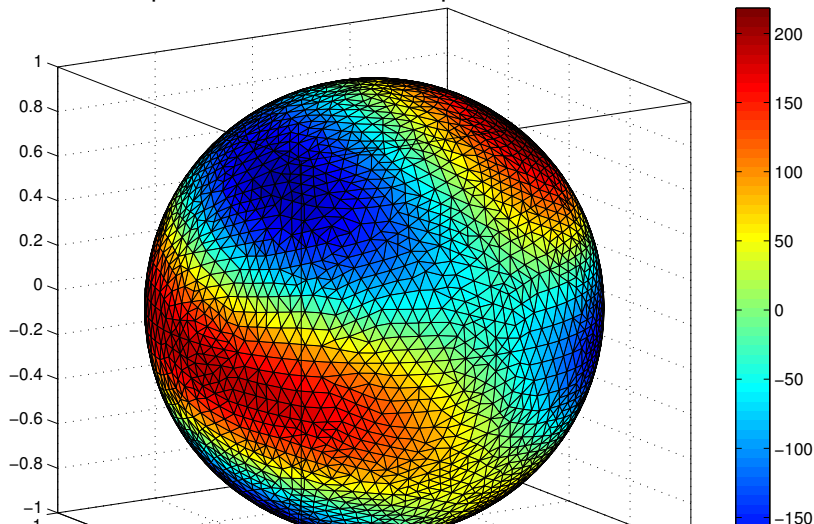


- Vzor Pipes and Filters
- Moduly pro analýzu, kompresi, dekompresi
- Konzolová aplikace pro testování



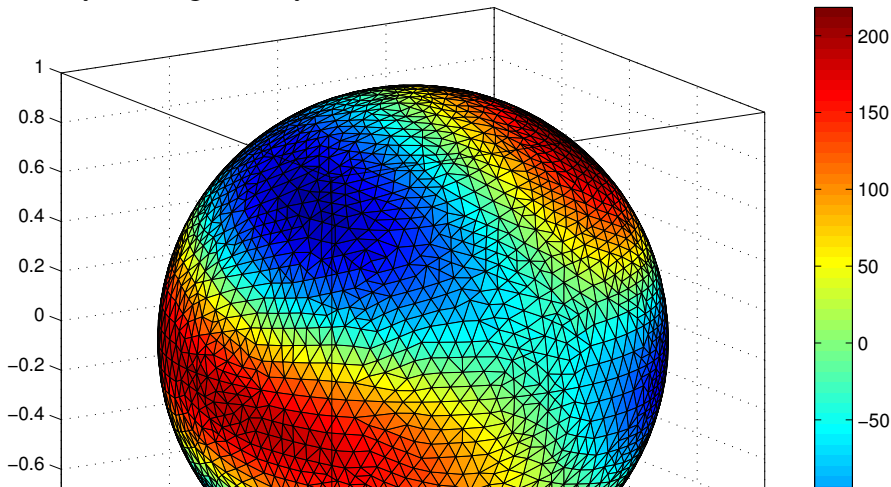


- Komprese: Text \rightarrow Strom \rightarrow Komprese
- Dekompres: Načtení \rightarrow Dekompres \rightarrow Ověření





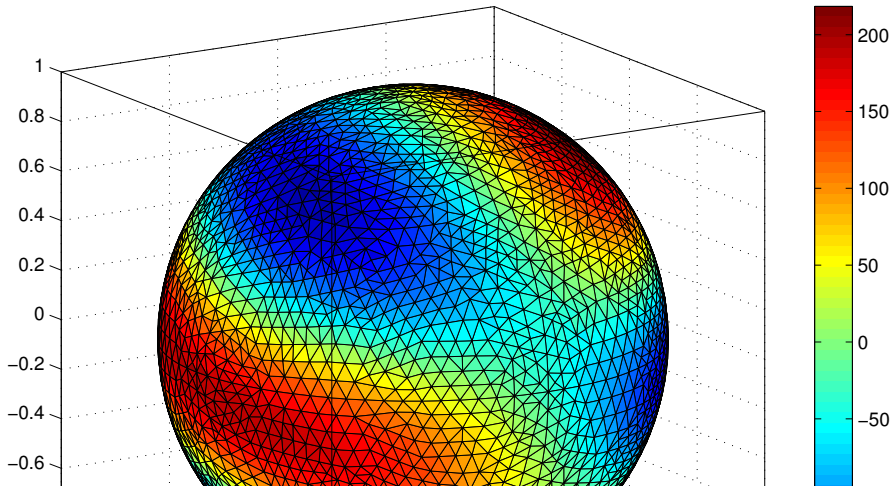
- Hledání opakujících se párů
- Nahrazení neterminály
- Vytvoření gramatiky



TreeRePair – bez linearizace



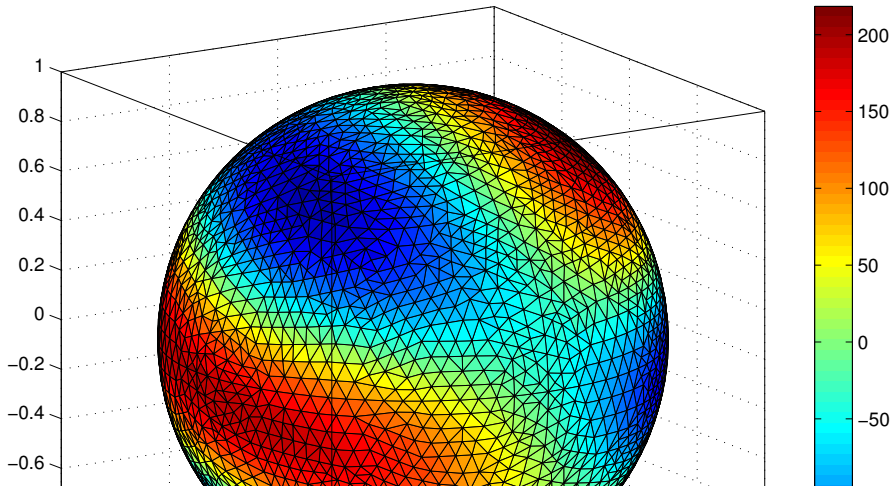
- Hledání opakujících se podstromů
- Hodnocení podle četnosti a velikosti
- Nahrazení pravidly



Linearizace + RePair



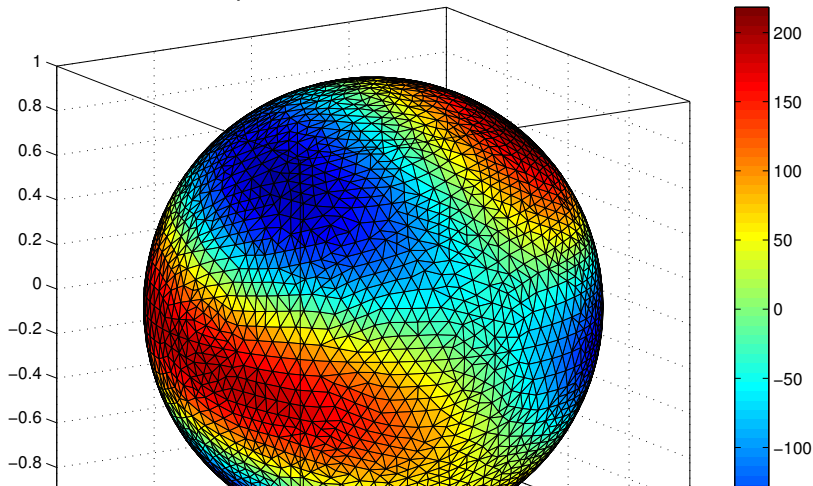
- Preorder průchod \rightarrow sekvence
- RePair \rightarrow pravidla
- Vylepšení: n-gramy, kontext



Shrnutí výsledků

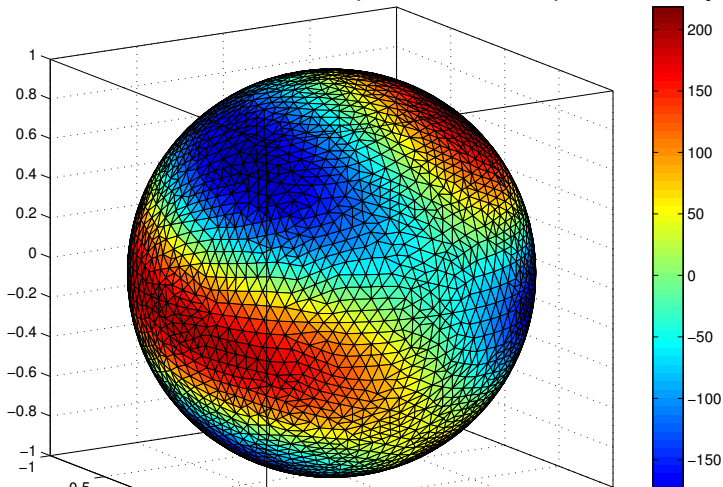


- Linearizace + RePair: nejlepší průměrné výsledky
- TreeRePair: náročnější, méně efektivní
- Skutečná komprese: cca 9





- Kompresce stromů není triviální
- Optimalizace výrazně pomáhají
- Možnosti do budoucna: paralelizace, adaptivní metody



Děkuji za pozornost

Marek Beran

VŠB – Technická univerzita Ostrava

marek.beran.st@vsb.cz

23. května 2025