

Dokumentace softwarové architektury pro kompresi stromových struktur

Verze 1.0

30. března 2025

Obsah

1	Úvod	3
1.1	Účel	3
1.2	Rozsah	3
1.3	Definice a zkratky	3
2	Přehled	4
3	Architektonické znázornění	4
3.1	Architektonické cíle a omezení	4
3.2	Use-Case pohled	4
3.3	Logický pohled	5
3.3.1	Popis diagramu tříd	6
3.3.2	Implementační pohled	8
3.4	Procesní pohled	8
3.5	Nasazení	9
4	Velikost a výkon	9
5	Kvalita	9

1 Úvod

1.1 Účel

Tento dokument poskytuje přehled o softwarové architektuře systému navrženého pro kompresi stromových struktur na základě analýzy vět. Cílem projektu je identifikovat opakující se vzorce v textových větách, zpracovat tyto vzorce do stromové struktury a navrhnout metody pro jejich kompresi. Architektura systému je navržena podle architektonického vzoru **Pipes and Filters**, což znamená, že systém bude složen z několika nezávislých filtrů, které zpracovávají data v průběhu několika fází, přičemž každý filtr vykonává specifickou operaci na datech a předává je dalšímu filtru v řetězci.

1.2 Rozsah

Dokument se zaměřuje na návrh a implementaci algoritmů pro kompresi stromových struktur. Systém bude zahrnovat moduly pro analýzu textových vět, detekci opakujících se vzorců, jejich převod na stromové struktury a aplikaci metod pro jejich kompresi. Každý z těchto kroků bude reprezentován samostatným filtrem v rámci architektury Pipes and Filters, čímž se umožní efektivní modulární přístup k vývoji, testování a optimalizaci.

1.3 Definice a zkratky

- **Stromová struktura:** Reprezentace vět s hierarchickými vztahy mezi slovy a frázemi.
- **Kompresní algoritmus:** Algoritmus sloužící k redukci velikosti stromové struktury při zachování všech relevantních informací.
- **Opakující se vzory:** Vzory, které se v textu vyskytují vícekrát a mohou být efektivně komprimovány.
- **Syntaktická analýza:** Proces analýzy textu, který identifikuje gramatické struktury v textových větách.

2 Přehled

3 Architektonické znázornění

3.1 Architektonické cíle a omezení

Cílem architektury je vytvořit systém založený na principu Pipes and Filters, kde každý filtr bude vykonávat specifickou operaci na datech, a výsledky budou postupně procházet celým systémem. Mezi hlavní požadavky na systém patří:

- **Modularita:** Systém bude navržen tak, že každý krok (analýza, identifikace vzorců, komprese) bude reprezentován samostatným filtrem.
- **Výkon:** Optimalizace pro zpracování velkých textových dat a složitých stromových struktur.
- **Flexibilita:** Schopnost přidávat nové filtry pro analýzu textu a kompresní metody.

3.2 Use-Case pohled

Hlavními use-case scénáři pro tento systém jsou:

- **Kompresní algoritmus:** Tento scénář zahrnuje řetězec filtrů, které zahrnují analýzu textu, identifikaci opakujících se vzorců, generování stromové struktury, aplikaci kompresního algoritmu a uložení komprimované struktury.
- **Dekomprese stromové struktury:** Tento proces zahrnuje řetězec filtrů pro dekompresi a ověření správnosti dekomprimované struktury.
- **Analýza textu:** Tento filtr analyzuje text a generuje stromovou strukturu, která bude následně použita pro identifikaci vzorců.
- **Identifikace opakujících se vzorců:** Tento filtr provádí detekci vzorců v textu.
- **Generování stromové struktury:** Tento filtr generuje stromovou strukturu na základě analýzy textu a identifikovaných vzorců.

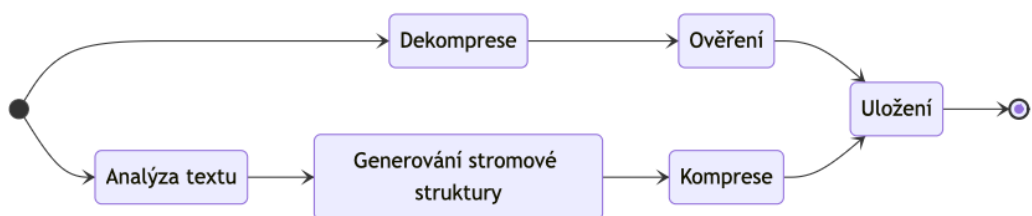
-
- ```

 usecaseDiagram
 actor Uživatel
 usecase UC1 as Kompresní algoritmus
 usecase UC2 as Analýza textu
 usecase UC3 as Generování stromové struktury
 usecase UC4 as Aplikace kompresního algoritmu
 usecase UC5 as Uložení komprimované struktury
 usecase UC6 as Dekomprese stromové struktury
 usecase UC7 as Ověření správnosti dekomprese

 Uživatel --> UC1
 Uživatel --> UC6
 UC1 --> UC2 : +include+
 UC1 --> UC3 : +include+
 UC1 --> UC4 : +include+
 UC1 --> UC5 : +include+
 UC2 --> UC3 : +include+
 UC3 --> UC4 : +include+
 UC4 --> UC5 : +include+
 UC6 --> UC7 : +extend+

```
- The diagram illustrates the functional requirements for a tree structure compression system. It features a central actor, 'Uživatel' (User), and seven use cases. The user interacts with the 'Kompresní algoritmus' (Compression algorithm) and 'Dekomprese stromové struktury' (Decompression of tree structure). The compression process is a sequence of steps: 'Kompresní algoritmus' includes 'Analýza textu' (Text analysis), 'Generování stromové struktury' (Generating tree structure), 'Aplikace kompresního algoritmu' (Applying compression algorithm), and 'Uložení komprimované struktury' (Storing compressed structure). 'Analýza textu' also includes 'Generování stromové struktury', and 'Generování stromové struktury' includes 'Aplikace kompresního algoritmu'. Finally, 'Dekomprese stromové struktury' is extended by 'Ověření správnosti dekomprese' (Verification of decompression correctness).

### 3.3 Logický pohled

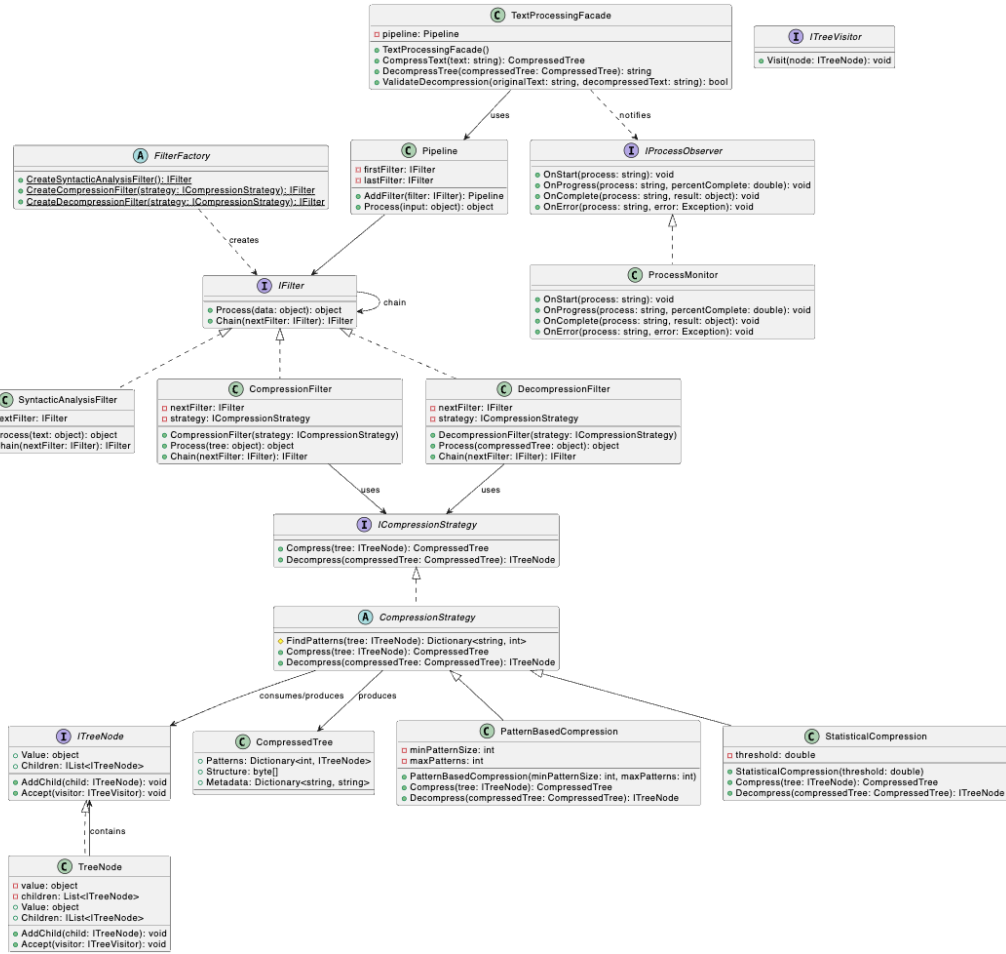


System bude rozdělen do několika filtrů:

- 5

### 3.3.1 Popis diagramu tříd

Diagram tříd znázorňuje architekturu systému pro kompresi stromových struktur pomocí několika návrhových vzorů. Hlavní komponenty a jejich vztahy jsou následující:



Obrázek 3: Třídní diagram nejdůležitějších komponent

### Hlavní rozhraní

- **IFilter**: Definuje metody pro zpracování dat a řetězení filtrů.

- **ITreeNode**: Reprezentuje uzel stromu s hodnotou a seznamem potomků.
- **ITreeVisitor**: Definuje metodu pro návštěvu uzlů stromu.
- **ICompressionStrategy**: Definuje metody pro kompresi a dekompresi stromových struktur.

### Abstraktní továrna pro filtry

- **FilterFactory**: Abstraktní třída pro vytváření různých typů filtrů.

### Konkrétní filtry

- **SyntacticAnalysisFilter**
- **CompressionFilter**
- **DecompressionFilter**

### Stromová struktura (Composite pattern)

- **TreeNode**: Implementuje uzel stromu s hodnotou a seznamem potomků.
- **CompressedTree**: Reprezentuje komprimovanou stromovou strukturu s uloženými vzory a metadaty.

### Strategie pro kompresi (Strategy pattern)

- **CompressionStrategy**: Abstraktní třída pro kompresní strategie.
- **PatternBasedCompression**
- **StatisticalCompression** Třídy sloužící k demonstraci implementace strategií komprese

### Pipeline (Pipes and Filters implementation)

- **Pipeline**: Třída pro řetězení a zpracování filtrů.
- **TextProcessingFacade**: Fasáda pro zjednodušení interakce s klientem.

## Observer Pattern pro monitorování

- **IProcessObserver**: Rozhraní pro sledování průběhu zpracování.
- **ProcessMonitor**: Implementace sledování průběhu zpracování.

## Vztahy mezi třídami

- **Pipeline** používá **IFilter**.
- **IFilter** může být řetězeno s jiným **IFilter**.
- **FilterFactory** vytváří instance **IFilter**.
- **CompressionFilter** a **DecompressionFilter** používají **ICompressionStrategy**.
- **TextProcessingFacade** používá **Pipeline**.
- **TreeNode** obsahuje **ITreeNode**.
- **CompressionStrategy** produkuje **CompressedTree** a zpracovává **ITreeNode**.
- **TextProcessingFacade** notifikují **IProcessObserver**.

### 3.3.2 Implementační pohled

Systém bude využívat jazyk C# a binding knihovny pro jazykové modely od *Digital Research Infrastructure for the Language Technologies, Arts and Humanities*, jako UDPipe a MorphoDiTa. Tyto knihovny jsou doimplementovány a přijímají potřebné modely.

## 3.4 Procesní pohled

Všechny filtry poběží v jednom hlavním vlákně. U kompresního algoritmu se pokusím zamyslet, jakým způsobem by mohla vhodně využívat paralelizaci, případně vhodnou implementaci.



### 3.5 Nasazení

Systém bude nasazen jako konzolová aplikace. Tato aplikace bude implementována v jazyce C# a poběží na platformě .NET. Konzolová aplikace bude umožňovat spouštění jednotlivých filtrů v rámci architektury Pipes and Filters přímo z příkazové řádky. Uživatelé budou moci zadávat vstupní data, specifikovat parametry pro různé filtry a sledovat průběh zpracování v reálném čase. Konzolová aplikace poskytne flexibilní a snadno použitelný způsob interakce se systémem, což usnadní testování, ladění a nasazení v různých prostředích.

## 4 Velikost a výkon

Systém bude optimalizován pro analýzu a kompresi stromových struktur. K dosažení co nejlepšího výkonu bude systém využívat paralelní zpracování pro efektivní analýzu textů a kompresi vzorců.

## 5 Kvalita

Systém bude vysoce kvalitní díky:

- **Rozšiřitelnosti:** Možnost přidávat nové filtry pro kompresi a analýzu.
- **Přesnosti:** Systém bude zajišťovat vysokou přesnost při analýze a kompresi.
- **Modularitě:** Struktura systému umožňuje snadnou údržbu a rozšiřování.