

# Analýza a komprese signálů

Reduction Paradox

Vojtěch Prokop

Semestrální práce

Vedoucí práce: Ing. Michal Vašínek, Ph.D.

Ostrava, 2022

## **Abstrakt**

No Czech or Slovak abstract is given

## **Klíčová slova**

No Czech or Slovak keywords are given

## **Abstract**

No English abstract is given

## **Keywords**

No English keywords are given

# Obsah

<b>Seznam obrázků</b>	<b>4</b>
<b>1 Obecný popis</b>	<b>5</b>
1.1 Komprese dat . . . . .	5
1.2 Popis problému . . . . .	5
1.3 Re-Pair . . . . .	6
1.4 Textové soubory . . . . .	7
<b>2 Řešení</b>	<b>8</b>
2.1 Nalezení k-gramů . . . . .	8
2.2 Výpočet entropie . . . . .	9
2.3 Vytvoření nové zprávy . . . . .	11
2.4 Výběr k-gramu . . . . .	11
2.5 Vytvoření pravidla . . . . .	12
2.6 Vygenerování nového znaku . . . . .	12
<b>3 Shrnutí a výsledky</b>	<b>13</b>
3.1 Výsledky souborů . . . . .	13
3.2 Výběr k-gramu . . . . .	17
3.3 Časová náročnost . . . . .	21
3.4 Nalezené extrémy . . . . .	23

# Seznam obrázků

1.1	Bezkontextová gramatika . . . . .	6
1.2	Textové soubory . . . . .	7
2.1	Stav před výběrem k gramu . . . . .	11
3.1	Změny u DNA největší výběr . . . . .	14
3.2	Změny u DNA náhodný výběr . . . . .	15
3.3	Změny u English největší výběr . . . . .	16
3.4	Změny u English náhodný výběr . . . . .	17
3.5	Srovnání metody výběru u DNA . . . . .	18
3.6	Srovnání metody výběru u English . . . . .	19
3.7	Srovnání metody výběru u Sources . . . . .	19
3.8	Srovnání metody výběru u Proteins . . . . .	20
3.9	Časová náročnost algoritmu napříč soubory . . . . .	21
3.10	Časová náročnost algoritmu napříč soubory . . . . .	22
3.11	Počet kroků algoritmu vzhledem k datovému souboru . . . . .	22
3.12	Extrémy po 15 krocích . . . . .	23
3.13	Maximální extrémy . . . . .	23
3.14	Velikost zprávy po určitém běhu . . . . .	24
3.15	Entropie po určitém běhu . . . . .	24

# Kapitola 1

## Obecný popis

### 1.1 Komprese dat

Komprese dat je metodika, která má za účel zpracovat počítačová data, tak aby byla snížena jejich velikosti. Počítačová data lze popsat velikostí, kterou bývá jednotka bajt nebo bit. Důvodem, proč se tedy zabýváme zmenšením objemu dat je:

- Snížení prostorových nároků na archivaci.
- Práce se sítí, kde můžeme dosáhnout například nižší doby pro přenos.
- Propustnost sítě.

Komprese samotná je v dnešní době rozdělena do dvou tříd. Ztrátové komprese, která znemožňuje zpětnou rekonstrukci do originálních dat. Ačkoliv tento fakt tolerujeme, jelikož většinou bývá zmenšení mnohem větší než je tomu u druhé třídy a to bezztrátové komprese. Tato třída má vlastnost, že komprimovaný soubor lze dekomprimací rekonstruovat do původní podoby.

### 1.2 Popis problému

Existují algoritmy, které umožní kompresi dat s pomocí definované gramatiky. Originální zpráva, pak bývá pomocí definovaných pravidel transformována do nové zprávy, která reprezentuje komprimované množství dat. Jeden z takových algoritmů je algoritmus Re-Pair (viz kapitola 1.3), který v každém iteračním kroku nahrazuje nejčtenější dvojici znaků za nový znak, a tímto vytváří pravidlo, dle kterého dochází k transformaci zprávy. V každém kroku je definováno pravidlo, které na levé straně obsahuje neterminální symbol a na pravé kombinaci terminálu a neterminálu. Příklad můžeme vidět na obrázku č. 1.1.

Dále víme, že lze vypočíst entropii nultého řádu, která nám říká kolik bitů je potřeba k zakódování jednoho symbolu ze souboru. Předpokládejme, že máme tedy dvě zprávy  $m_0$  a  $m_1$ , kde

$$\begin{aligned}
S &\rightarrow R_0 R_1 R_2 e \\
R_0 &\rightarrow ab \\
R_1 &\rightarrow R_0 c \\
R_2 &\rightarrow R_1 d
\end{aligned}$$

Obrázek 1.1: Bezkontextová gramatika

$m_0$  reprezentuje originální zprávu a  $m_1$  zprávu po prvním kroku Re-Pair algoritmu. Lze definovat problém, kterému říkáme paradox redukce, který si ukážeme na zprávě  $m_0$  a  $m_1$ . Víme že  $|m_0| > |m_1|$  neboli slovně počet znaků v originální zprávě je větší než počet znaků v transformované zprávě, ačkoliv entropie nultého řádu originální zprávy je menší než entropie nultého řádu transformované zprávy  $Hm_0 < Hm_1$ .

S tímto vědomím můžeme tedy vynaložit snahu provést transformaci zprávy  $n$  krát. Tento pokus reprezentuje nalezení extrému popisujícího kolikrát lze transformovat zprávu s dodržením snížení počtu symbolů, a však zvýšení velikosti entropie oproti předchozí zprávě.

### 1.3 Re-Pair

Jedná se o kompresní algoritmus založený na gramatice, který na základě vstupního textu vytvoří bezkontextovou gramatiku. V každém kroku je nahrazena nejčastější dvojice znaků vyskytující se v textu. Z provedených pokusů víme, že Re-Pair dosahuje vysokých kompresních poměrů. Nevýhodou bohužel je spotřeba paměti, která je přibližně 5krát větší než velikost vstupu.

## 1.4 Textové soubory

Byly využity textové soubory o velikost 50 MB z veřejně dostupné textové kolekce obsahující soubory se zdrojovými kódy, DNA, anglickými texty, XML texty a proteiny. Více informací lze vidět na obrázku č. 1.2.

<i>Collection</i>	<i>Size (bytes)</i>	<i>Alphabet size</i>	<i>Inv match prob</i>
<i>SOURCES</i>	210,866,607	230	24.77
<i>PITCHES</i>	55,832,855	133	39.75
<i>PROTEINS</i>	1,184,051,855	27	17.02
<i>DNA</i>	403,927,746	16	3.91
<i>ENGLISH</i>	2,210,395,553	239	15.25
<i>XML</i>	294,724,056	97	28.73

Obrázek 1.2: Textové soubory

V mém případě jsem pro experimenty využil zdrojové kódy, DNA, proteiny a anglické texty. Níže textové soubory budou více popsány.

### 1.4.1 Zdrojové kódy

Soubor obsahuje sjednocení zdrojových kódů z několika .c, .h, .C a .java souborů. Jedná se tedy o zdrojové soubory s programovacími jazyky Java, C možná C++.

### 1.4.2 DNA

Soubor obsahující sekvence DNA bez metadat. Můžeme zde očekávat vysokou četnost nukleotidů A, C, G, T. Mimo tyto znaky se v souboru vyskytují i některé další speciální znaky, ale v mnohonásobně nižší četnosti než zmíněné znaky popisující první písmeno nukleotidů.

### 1.4.3 Proteiny

Soubor obsahující sekvence proteinů.

### 1.4.4 Anglické texty

Sjednocení textů v anglickém jazyce z vybraných souborů (knih), které jsou ke stažení v Gutenberg Projektu. Je výhodou, že texty neobsahují hlavičku, která se vyskytuje v některých souborech stažených přímo v zmíněném projektu. Hlavička obsahuje metadata k knížce.

## Kapitola 2

# Řešení

Celé řešení problému lze popsat následujícími několika dílčími kroky, které pak budou blíže rozebrány v kapitolách níže.

- Nalezení všech  $k$  gramů.
- Výpočet změny entropie po nahrazení každého  $k$  gramu.
- Výběr  $k$  gramu, která bude využit pro vytvoření pravidla.
- Vygenerování nového znaku.
- Vytvoření pravidla.
- Vytvoření nové zprávy.

### 2.1 Nalezení $k$ -gramů

Textový soubor je vytvořen z posloupnosti znaků. V prvním kroku průběhu algoritmu je potřeba nalézt všechny  $k$ -tice, které budeme moci v sekvenci nahradit. Z těchto nalezených dvojic, pak bude vytvořeno pravidlo, pokud splní dvě podmínky. První podmínkou je zvyšující se entropie po nahrazení každého výskytu  $k$ -gramu. Druhou podmínkou je výběr daného  $k$  gramu námi definovanou metodou. Tento výběr je popsán v kapitole níže (viz 2.4).



Metoda, která nalezne všechny tyto k gramy může vypadat v programovacím jazyce Python následovně:

---

```
def find_k_grams_freq(data, max_size_k=2):
    kgrams_dic = {}
    for k in range(2, max_size_k+1):
        for i in range(len(data) - k):
            n_gram = data[i:i+k]
            kgrams_dic[n_gram] = kgrams_dic.get(n_gram, 0) + 1
    return kgrams_dic
```

---

Vstupními parametry jsou textová data a velikost k-gramu, které chceme v textových datech nalézt. Postupně pak iterujeme přes textový soubor a ukládáme do datové struktury slovníku kolikrát jsme daný k gram viděli.

## 2.2 Výpočet entropie

Z předchozí kroku víme, které k gramy se nám v souboru vyskytují. Této informace využijeme a pro každý k gram vypočteme novou entropii. Entropie bude vycházet z akce vytvoření nového pravidla a nahrazení k-gramu novým symbolem, což zapříčiní rozšíření aktuální abecedy o jeden symbol.

Výpočet entropie lze provést dvěma způsoby a to:

- Vytvořením nové zprávy a výpočtu pomocí shannonského vzorce.
- Využití definovaného vzorce, který vypočte posun vzhledem k entropii.

### 2.2.1 Shannon vzorec

Tento způsob pracuje s přístupem hrubé síly, kdy je vypočtena entropie pro dvě zprávy. A to zprávu, která byla v předchozím kroku a zprávu, která byla vytvořena po aplikaci nově vytvořeného pravidla.

---

```
def calc_H(p):
    H = 0
    for k, v in p.items():
        #Shannon equation!
        H += p[k] * math.log2(p[k])
    return -H

def calc_entropy_for_message(message):
    counter = calc_freq(message)
    n = get_n(counter)
    p = calc_p(counter, n)
    H = calc_H(p)
    return H

def diff_entropy(message1, message2):
    message1_entropy = calc_entropy_for_message(message1)
    message2_entropy = calc_entropy_for_message(message2)
    message1_entropy_size = message1_entropy * len(message1)
    message2_entropy_size = message2_entropy * len(message2)
    diff = message1_entropy_size - message2_entropy_size
    return diff
```

---

## 2.3 Vytvoření nové zprávy

Abychom mohli vypočítat rozdíl entropie je potřeba transformovat zprávu do podoby, která reprezentuje zprávu v následujícím kroku. Takováto akce lze provést v programovacím jazyce Python velice jednoduchým způsobem.

---

```
def transform_message(message, ngram_for_replace):
    current_alphabet_size = np.unique(list(message))
    replace_character = find_not_existing_character(current_alphabet_size)
    return message.replace(ngram_for_replace, replace_character),
        replace_character
```

---

## 2.4 Výběr k-gramu

Po provedení výpočtu změny entropie pro každý nalezený k gram dostaneme několik možností, jak provést samotnou transformaci. Aktuální stav si můžeme představit pomocí obrázku č. 2.1. Diff pak znázorňuje výpočet  $|m0| * Hm0 - |m1| * Hm1$ .

	Counter	Diff
AA	862	-711.287863
AC	612	-700.944705
CG	397	-526.484088
GT	585	-642.013983
TG	609	-624.132715
TT	797	-725.647370
TA	856	-863.100574
GC	457	-493.377886
CT	648	-656.895096
TC	616	-682.016142
CC	476	-450.401292
CA	640	-680.422527
AG	604	-644.015856
GG	424	-397.457487
GA	567	-669.955161
AT	848	-870.225504

Obrázek 2.1: Stav před výběrem k gramu

Z těchto možností je potřeba vybrat jednu transformaci, které lze provést pomocí dvou možností s kterými bylo experimentováno:

- Výběr největší změny.
- Výběr náhodné změny.

## 2.5 Vytvoření pravidla

Jak bylo zmíněno transformace zprávy vychází z nahrazení vybraného k-gramu za nový neexistující znak (viz. kapitola č. 2.6). Tohle pravidlo je následně uloženo do datové struktury slovníku, tak aby bylo možné zpětně rekonstruovat originální zprávu.

## 2.6 Vygenerování nového znaku

Abychom mohli vytvořit nové pravidlo je potřeba použít nový neexistující znak vzhledem k aktuálnímu textovému řetězci. Tento znak lze vygenerovat pomocí dvou způsobů, které byly vyzkoušeny.

- Využití ASCII.
- Využití UNICODE.

Abychom neměli problém s limitovaným prostorem je potřeba využít UNICODE. Ten nabízí větší množství použitelných znaků.

---

```
all_chars_uni = tuple(chr(i) for i in range(32, 0x110000) if chr(i).isprintable())
```

```
all_chars_ascii = list(range(0, 256))
```

```
all_chars_ascii = [chr(ascii_char) for ascii_char in all_chars_ascii]
```

---

## Kapitola 3

# Shrnutí a výsledky

Z důvodů zohlednění běhu algoritmu byly experimenty provedeny s velikostí 10 000 znaků, které byly náhodně vybrány z každého datového zdroje. Pro tyto textové řetězce byl proveden běh algoritmu s limitovaným počtem kroků a to číslem 15. Vizualizace a krátké slovní zhodnocení lze nalézt v kapitole č. 3.1.

Zároveň bylo provedeno porovnání vzhledem k výběrové funkci (největší, náhodný), kde více se můžeme rozvést v kapitole č. 3.2.

Pro každý zdroj byl proveden experiment a s tím změřeno, jak dlouho výpočet běžel. Zároveň bylo provedeno pár optimalizací, které zrychlovali samotný výpočet. Více o tomto se lze dozvědět v kapitole č. 3.3.

Jak dopadlo nalezení samotných extrémů bude rozebráno v kapitole č. 3.4.

### 3.1 Výsledky souborů

V této kapitole budou zobrazeny vizualizace k výsledkům se zmíněnými limity:

- Výběr 10 000 znaků ze souboru.
- 15 provedených kroků.

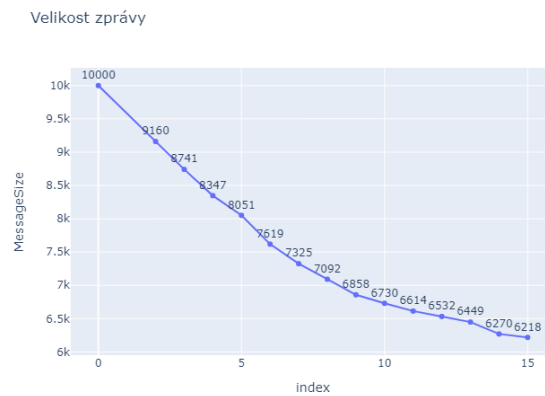
Mimo jiné budou opomenuty některé datové zdroje, jelikož pozorované chování má většinou stejné vlastnosti.

### 3.1.1 DNA

Na obrázku s výběrem největšího lze pozorovat, že počátečních pár výběrů způsobí největší změnu a zmenší zároveň nejvíce velikost. Po provedení těchto transformací, pak zpráva svůj stav mění velice pomalu s logaritmickou tendencí.



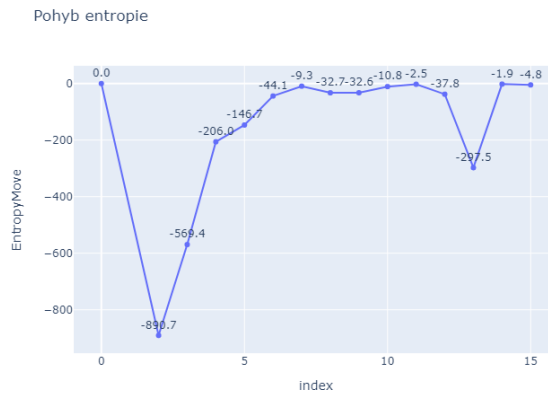
(a) Změna entropie v kroku u DNA při výběru největšího



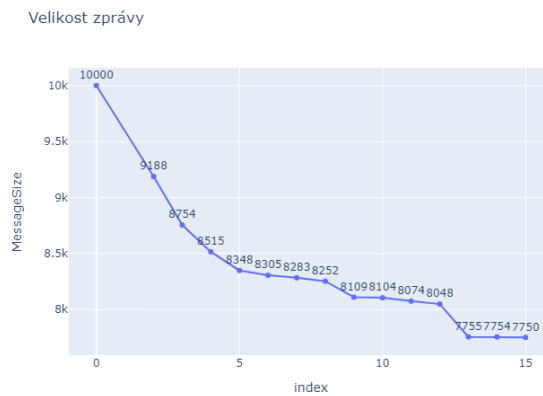
(b) Změna velikost zprávy u DNA

Obrázek 3.1: Změny u DNA největší výběr

V porovnání s výběrem největšího u výběru náhodného lze pozorovat více dynamický graf, kde dochází k náhodným silným skokům změny. Také je nutné podotknout, že tato změna krásně koreluje s změnou velikosti zprávy. Při každé větší změně entropie dochází k markantnější změně ve velikosti zprávy.



(a) Změna entropie v kroku u DNA při náhodném výběru



(b) Změna velikost zprávy u DNA

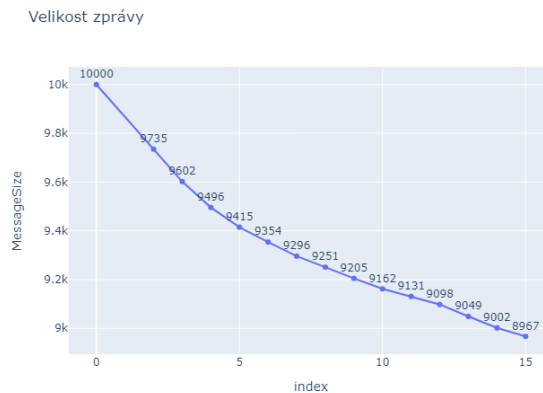
Obrázek 3.2: Změny u DNA náhodný výběr

### 3.1.2 English

Podobně jako u předchozího datové zdroje s DNA můžeme pozorovat, že graf z počátku rychle vstoupne, respektive poklesne velikost zprávy. Přičemž u náhodného výběru jsou grafy více scho-  
dovité. Je potřeba zmínit, že jak uvidíme níže pokles zprávy oproti výběru největšího je za určitý počet kroků mnohem menší, než by mohl být.



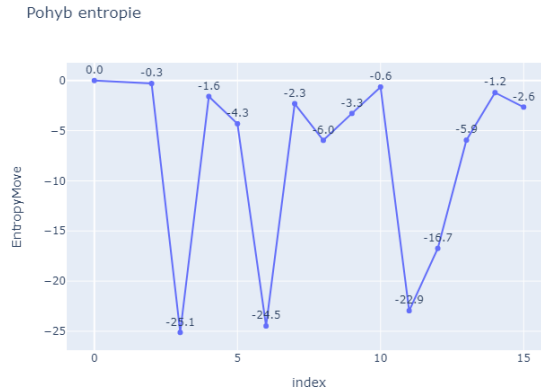
(a) Změna entropie v kroku u English při výběru největšího



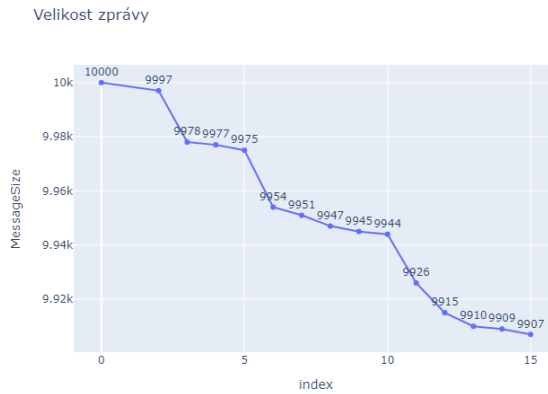
(b) Změna velikost zprávy u English

Obrázek 3.3: Změny u English největší výběr





(a) Změna entropie v kroku u English při náhodném výběru



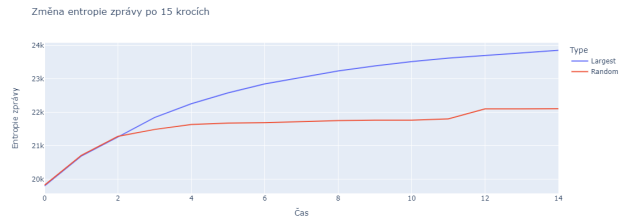
(b) Změna velikost zprávy u English

Obrázek 3.4: Změny u English náhodný výběr

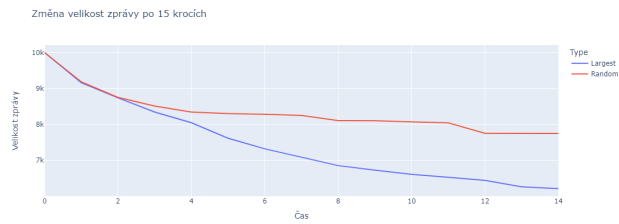
## 3.2 Výběr k-gramu

Implementovány byly dva přístupy, jak v každém kroku vybírat k-gram z kterého bude konstruováno pravidlo. Jedna metody vybírala sofistikovaným způsobem, v kterém došlo prvně ke setřizení dle změny entropie a největší změna byla použita. A druhý způsob fungoval na bázi náhodného výběru k-gramu. Jak lze vidět v podkapitolách níže patřící ke kapitole č. 3.2, metoda náhodného výběru fungovala ne moc efektivním způsobem, což šlo předpokládat. Zároveň je nutné zmínit, že nebylo dosaženo poznatelného zrychlení, jelikož třídící algoritmus má malou časovou složitost. Určitě je třeba se zaměřit na rozdíl změny zde po odběhnutých 15 krocích.

### 3.2.1 DNA



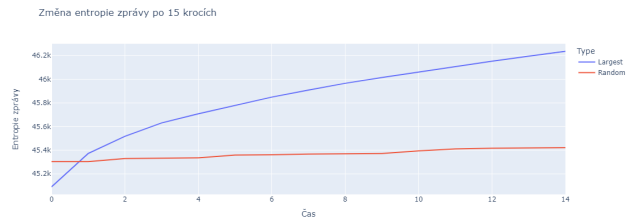
(a) Entropie



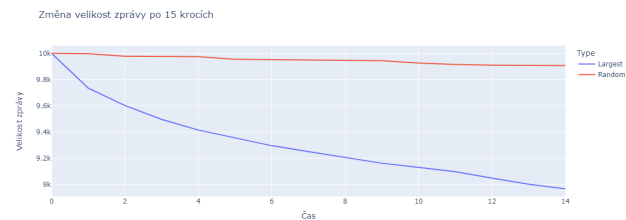
(b) Velikost zprávy

Obrázek 3.5: Srovnání metody výběru u DNA

### 3.2.2 English



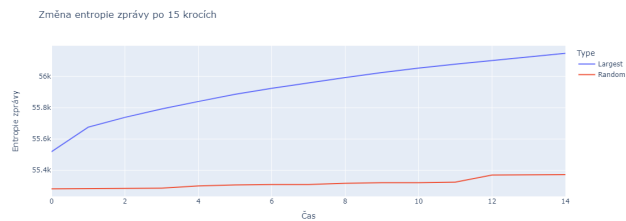
(a) Entropie



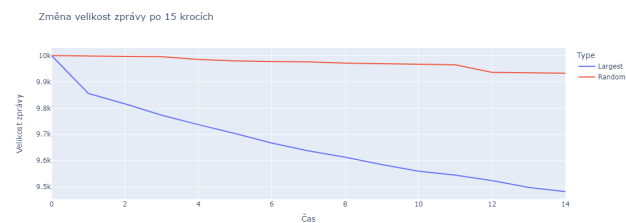
(b) Velikost zprávy

Obrázek 3.6: Srovnání metody výběru u English

### 3.2.3 Sources



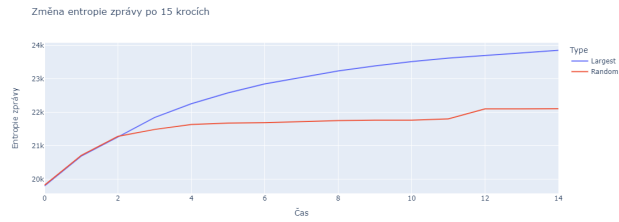
(a) Entropie



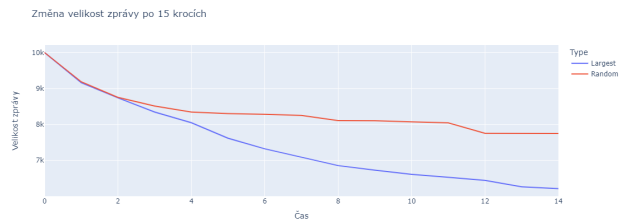
(b) Velikost zprávy

Obrázek 3.7: Srovnání metody výběru u Sources

### 3.2.4 Proteins



(a) Entropie



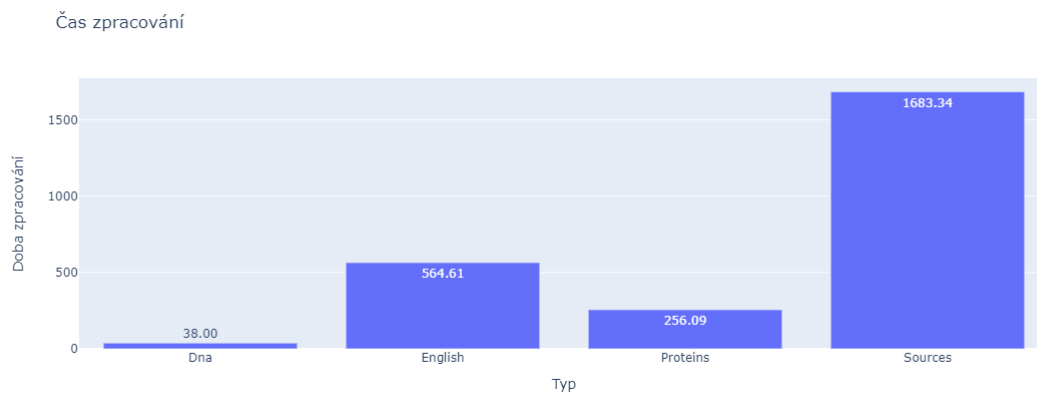
(b) Velikost zprávy

Obrázek 3.8: Srovnání metody výběru u Proteins

## 3.3 Časová náročnost

### 3.3.1 S omezením na 15 kroků

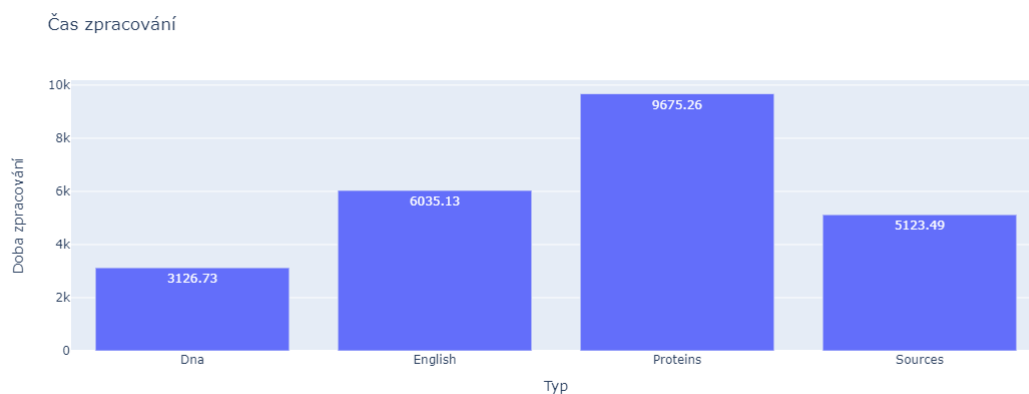
Byly provedeny experimenty nad všemi datovými zdroji, přičemž nejdéle algoritmus běžel na zdrojových kódech, následně na angličtině. Proteiny a DNA soubor proběhly mnohem rychleji. Toto porovnání lze vidět na obrázku č. 3.9. Důvodem, proč ostatní datové zdroje probíhaly mnohem déle je násobně větší abeceda, tudíž více možností, kterými lze sestavit k-gramy.



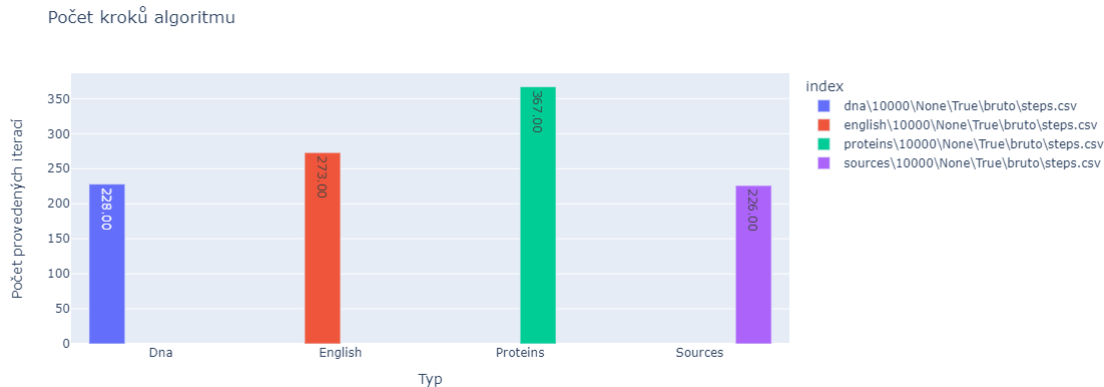
Obrázek 3.9: Časová náročnost algoritmu napříč soubory

### 3.3.2 Bez omezení

Zároveň byly provedeny experimenty bez striktního omezení kroků. Jak můžeme vidět níže na obrázku č. 3.10 doba běhu se lišila v uspořádání oproti předchozímu obrázku s omezeným počtem kroků. Důvodem, proč některé soubory byly zpracovávány delší časový úsek je hlavně počet kroků, které musely být provedeny. Zmíněný počet kroků vizualizuje obrázek č. 3.11.



Obrázek 3.10: Časová náročnost algoritmu napříč soubory



Obrázek 3.11: Počet kroků algoritmu vzhledem k datovému souboru

## 3.4 Nalezené extrémy

### 3.4.1 S omezením na 15 kroků

Na obrázku č. 3.12 lze pozorovat nalezené extrémy po 15 krocích algoritmu na datových zdrojích. Lze pozorovat vlastnosti, které byly zmíněny již výše a to:

- Náhodný výběr nedává smysl, jelikož vrací mnohem slabší výsledky.
- Náhodný výběr je časově stejně náročný jako výběr největšího.
- Velikost abecedy výrazně ovlivňuje běh algoritmu.
- U DNA souboru došlo k výraznému projevení paradoxu redukce.
  - Velikost zprávy se snížila skoro o 40 %.
  - Entropie vstoupila skoro o 50 %.
  - Běh algoritmu je extrémně rychlý, jelikož soubor neobsahuje velké množství znaků.

	LimitSteps	PickMethod	DiffEntropy	DiffSizeEntropy	DiffMessageSize	CalculationTime	MessageSize	Diff %	Entropy	Diff %
Type										
dna	15	Largest	1.856019	4052.715529	3782	37.995292		37.82	48.39	
dna	15	Random	0.870489	2286.866369	2250	32.605225		22.50	30.52	
english	15	Largest	0.647244	1146.008935	1033	564.606717		10.33	12.55	
english	15	Random	0.054385	117.463967	93	489.965724		0.93	1.19	
proteins	15	Largest	0.400220	890.442610	676	256.088325		6.76	8.69	
proteins	15	Random	0.158455	361.133565	281	227.854186		2.81	3.64	
sources	15	Largest	0.369649	629.142635	518	1683.338367		5.18	6.24	
sources	15	Random	0.046409	90.596139	67	1269.716560		0.67	0.83	

Obrázek 3.12: Extrémy po 15 krocích

### 3.4.2 Bez omezení

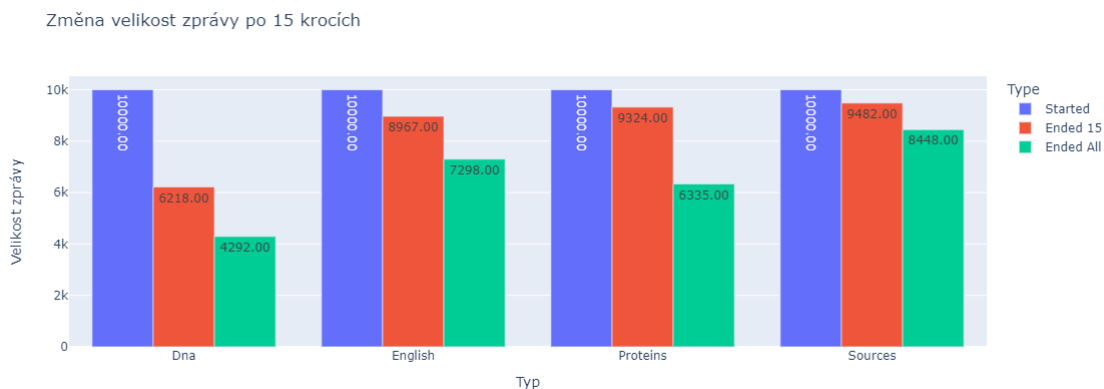
Na obrázku č. 3.13 lze pozorovat nalezené extrémy po maximální počtu kroků. Respektive algoritmus běžel, dokud bylo možné změnu provádět.

	LimitSteps	PickMethod	DiffEntropy	DiffSizeEntropy	DiffMessageSize	CalculationTime	NumberOfSteps	MessageSize	Entropy
Type								Diff %	Diff %
dna	NaN	Largest	4.017859	5888.142596	5708	3126.732789	228	57.08	66.88
english	NaN	Largest	2.148746	3488.149091	2702	6035.128097	273	27.02	32.26
proteins	NaN	Largest	3.123436	4590.176802	3665	9675.256349	367	36.65	42.96
sources	NaN	Largest	1.236022	1972.620438	1552	5123.487468	226	15.52	18.47

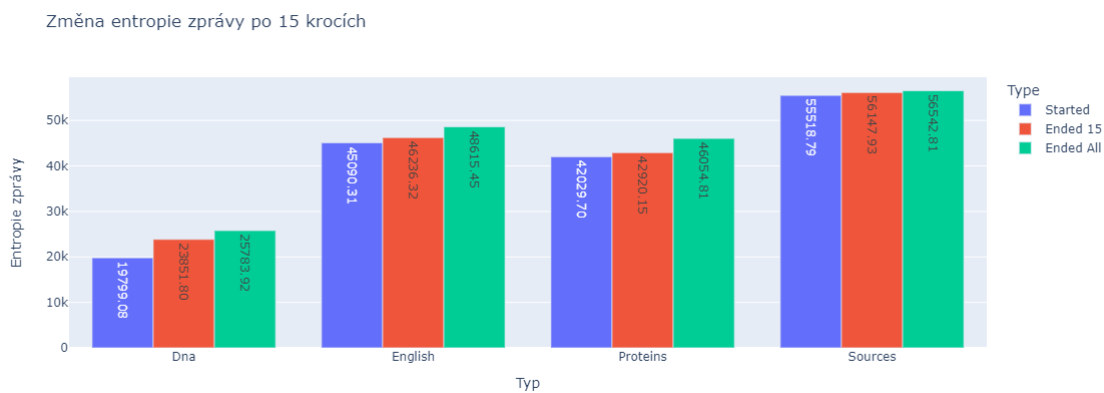
Obrázek 3.13: Maximální extrémy

### 3.4.3 Rozdíl v omezení

Změna měla logaritmické chování, kdy z počátku docházel k největší změně a postupně změna byla menší a menší. Z obrázku č. 3.14 a 3.15 můžeme vidět, že velikost zprávy byla sice poznatelně zmenšena po více provedených krocích, akorát vzhledem k změně entropie už tento rozdíl pozorovat nelze.



Obrázek 3.14: Velikost zprávy po určitém běhu



Obrázek 3.15: Entropie po určitém běhu