

# Meranie redundancie v dátach

Školiteľ: Ing. Dušan Bernát, PhD.

Matúš Hluch

A large yellow triangle is positioned in the bottom right corner of the slide, pointing towards the top right.

# Motivácia

- Exponenciálne zväčšujúce sa množstvo dát
  - Veľa duplicít
- Súborové systémy
  - Komprimujúce Btrfs, ZFS...
  - Copy-on-write
- Štandardne delené do blokov rovnakej veľkosti
  - 512 bajtov, niektoré nastaviteľné

# Ciele práce

- Preskúmať vlastnosti redundancie v dátach
  - Fixná veľkosť delenia
    - Iné práce - neefektívne
- Podľa typu dát
  - Dokumenty, obrázky, videá, virtuálne stroje...
- Podľa veľkosti okna
  
- Vplyv použitej hešovacej funkcie
  - Nie kryptograficky bezpečné

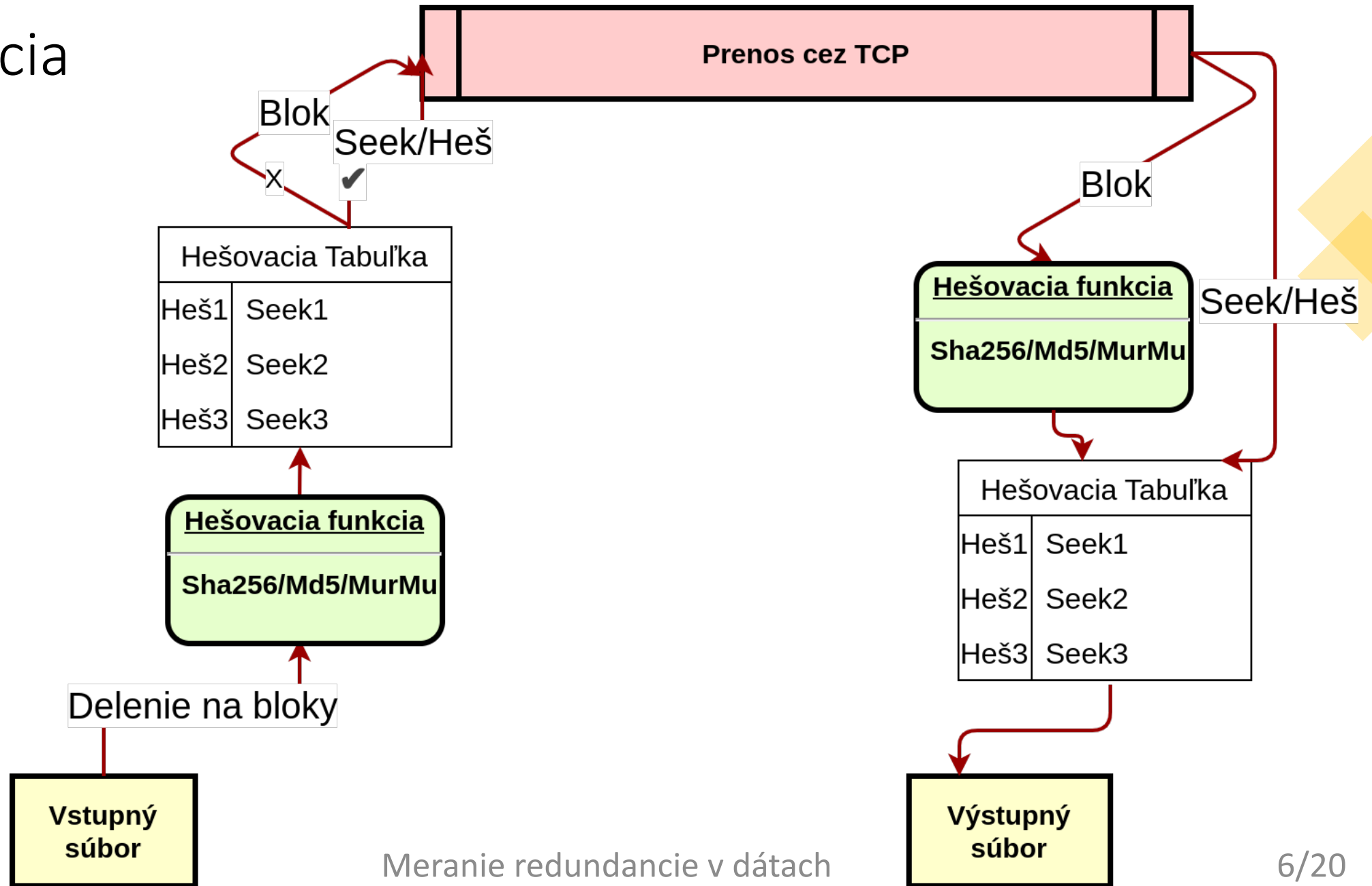
# Rozdelenie práce

- Prvá časť
  - Prakticky využiteľný nástroj s deduplikáciou
  - Výstup dát z nástroja
- Druhá časť práce
  - Analyzovanie výstupu dát z nástroja

# Náš prístup

- Deduplikácia v súborovom systéme
  - Zložitá – väčšie pracovné tímy
- Deduplikácia dát pri prenose cez sieť
  - Rozdelenie súboru na bloky rovnakej veľkosti
  - Hľadanie duplicitných blokov
    - Porovnávanie hešov
  - Blok dát sa pošle iba raz
    - Duplicitné - seek/heš

# Aplikácia



# Aplikácia

Send deduplicated data over TCP.

First parameter is window size (number of bytes into which data will be divided).

Second parameter is host.

Third parameter is port.

Fourth parameter is file name, from which data will be sent (applicable only if not using the stdin flag).

Usage:

```
dedupnetgo send window_size host port *file_name [flags]
```

Flags:

<code>-c, --collisionDetection</code>	Detect hash collisions.
<code>-l, --hash_cache_limit uint</code>	Set limit on number of hashCache entries.
<code>-f, --hash_function string</code>	Choose a hash function: sha256/md5/murmur/xxh/simple. (default "sha256")
<code>-h, --help</code>	help for send
<code>-o, --output_seeks</code>	Output seeks from cache to stderr.
<code>-m, --sendHash</code>	Send duplicates hashes instead of seek (if set, blocks will be stored in memory).
<code>-s, --stdin</code>	Send data from stdin.

# Deduplikácia dát pri prenose cez sieť

- Aplikácia pre príkazový riadok
  - Implementovaná v Go
    - Knižnica Cobra
      - Kubernetes, GitHub CLI...
- Prenos súborov cez TCP
  - Navrhnutý vlastný komunikačný protokol nad TCP





# Aplikácia

- Nastavenia
  - Veľkosť okna
  - Hešovacia funkcia
    - Kolízie, rýchlosť, spotreba pamäte...
  - Detekcia kolízií
    - Vyššia spotreba pamäte
  - Veľkosť pamäte

# Komunikačný protokol nad TCP

- Prvý bajt (hlavička) určuje typ prenášaných dát
- 1. Prenosové nastavenia
  - Nastavuje odosielateľ
- 2. Prenos dát
  - a) Blok dát
  - b) Seek
  - c) Heš
  - d) Posledný blok (môže mať inú veľkosť ako ostatné)
- 3. Koniec prenosu

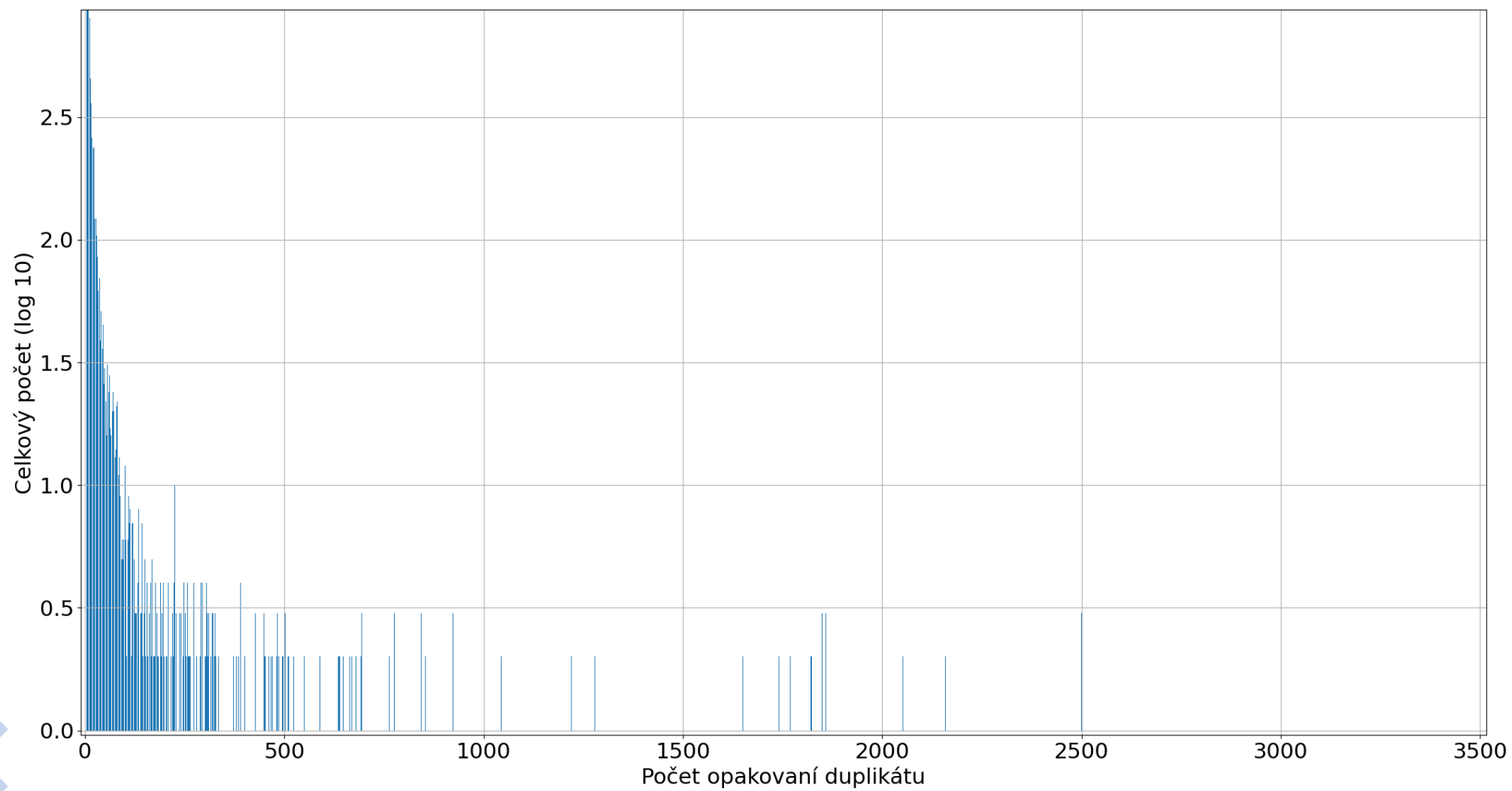
# Analýza prenášaných dát

- Spracovanie výstupu z nástroja
  - Implementované v Pythone
    - Matplotlib, SciPy...
- Analýza viacerých parametrov
  - Kompresia podľa veľkosti okna a typu súboru
  - Vzdialenosti medzi duplikátmi
  - Početnosti duplikátov
  - ...

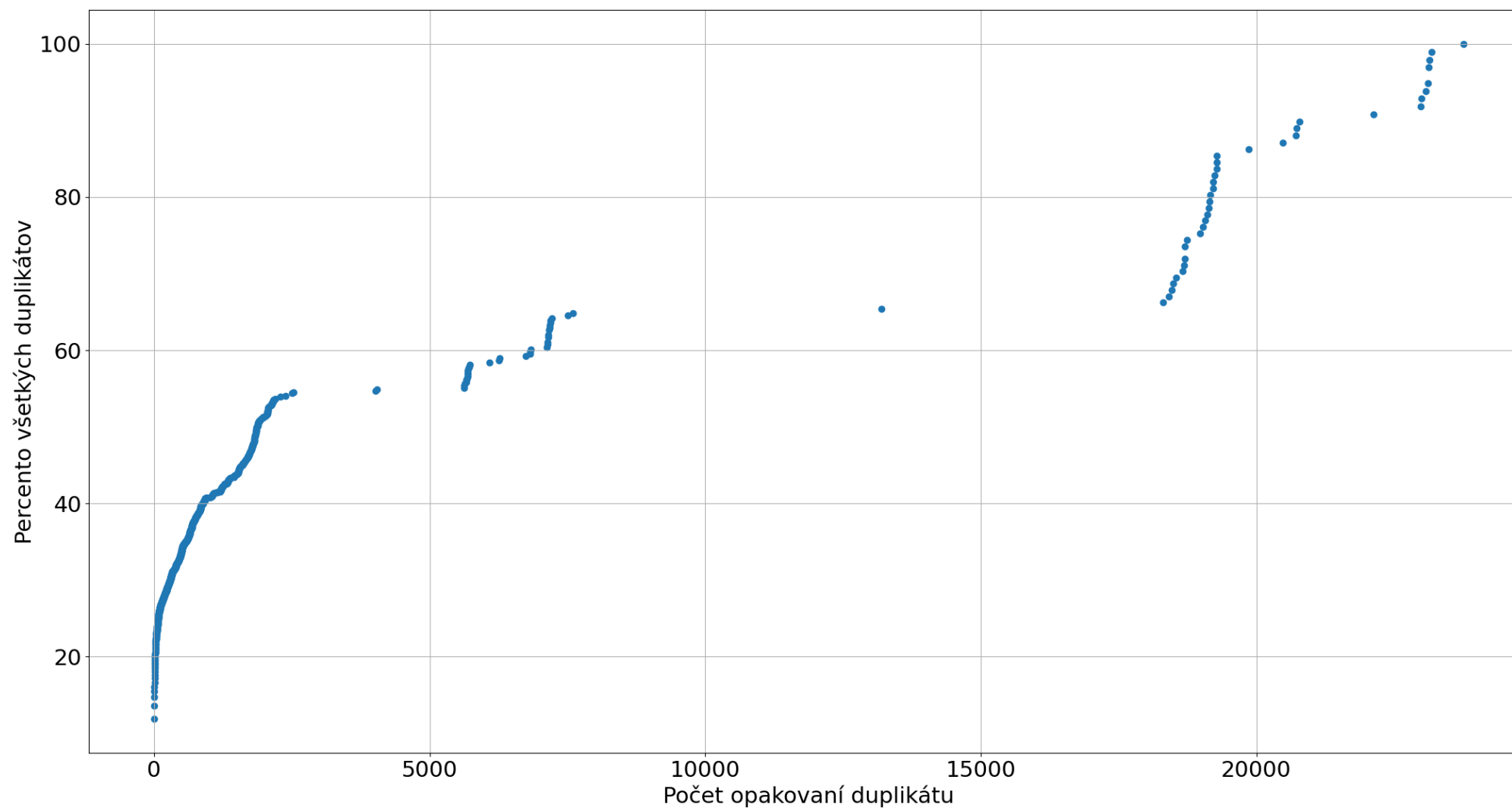
# Vstupné dáta

- Rôzne súborové typy
  - Väčšina datasetov rovnakého typu
    - Určené na strojové učenie, ...
- Govdocs1
  - Milión súborov rôznych typov, voľne šíriteľné
  - Súbory z .gov domén
    - pdf, jpg, txt, html, ppt...
- Využívame
  - Niekoľko GB dát - rozdelených podľa typu
  - Okolo 10 000 súborov

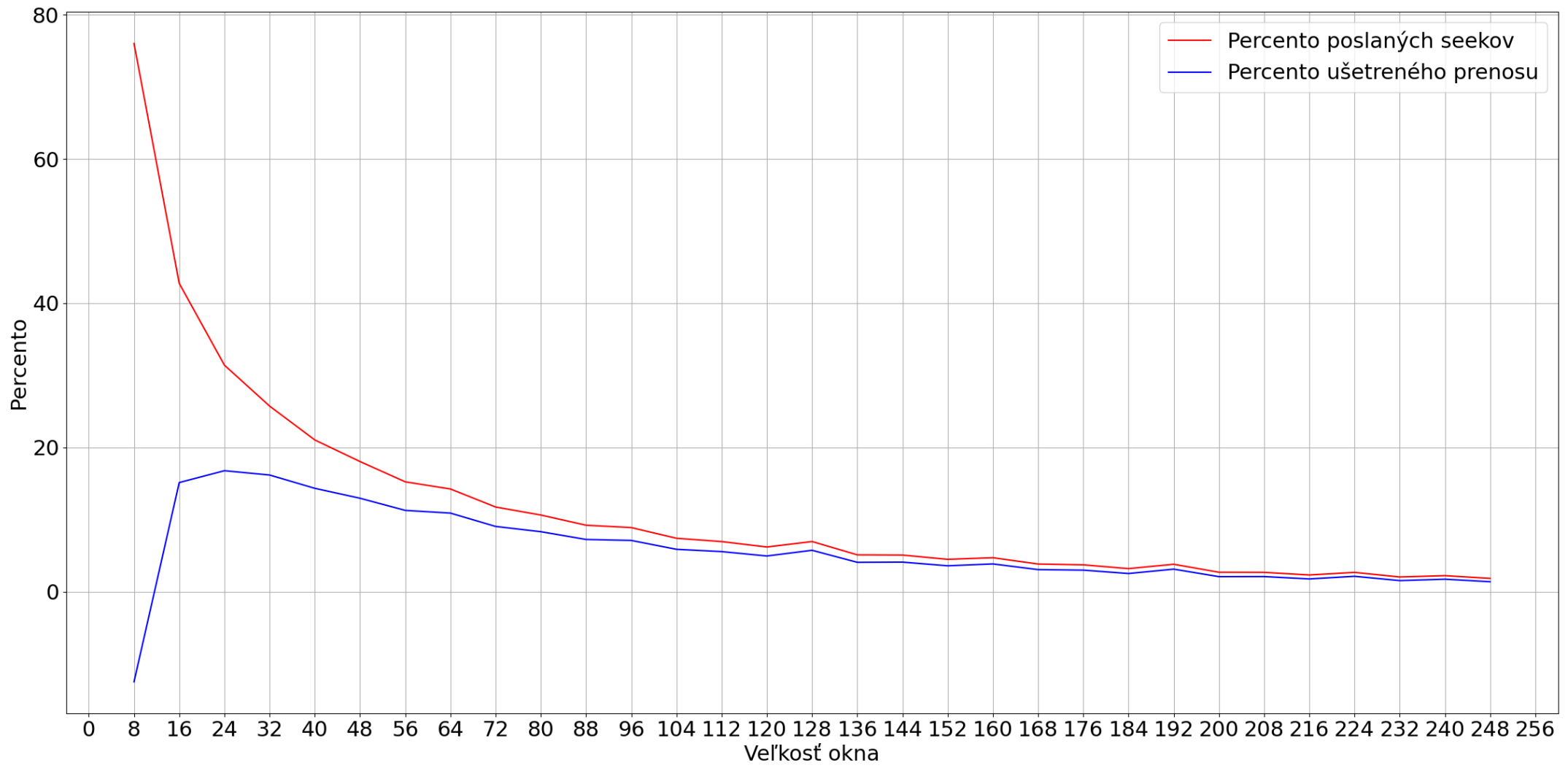
# TXT - veľkosť okna 32 bajtov



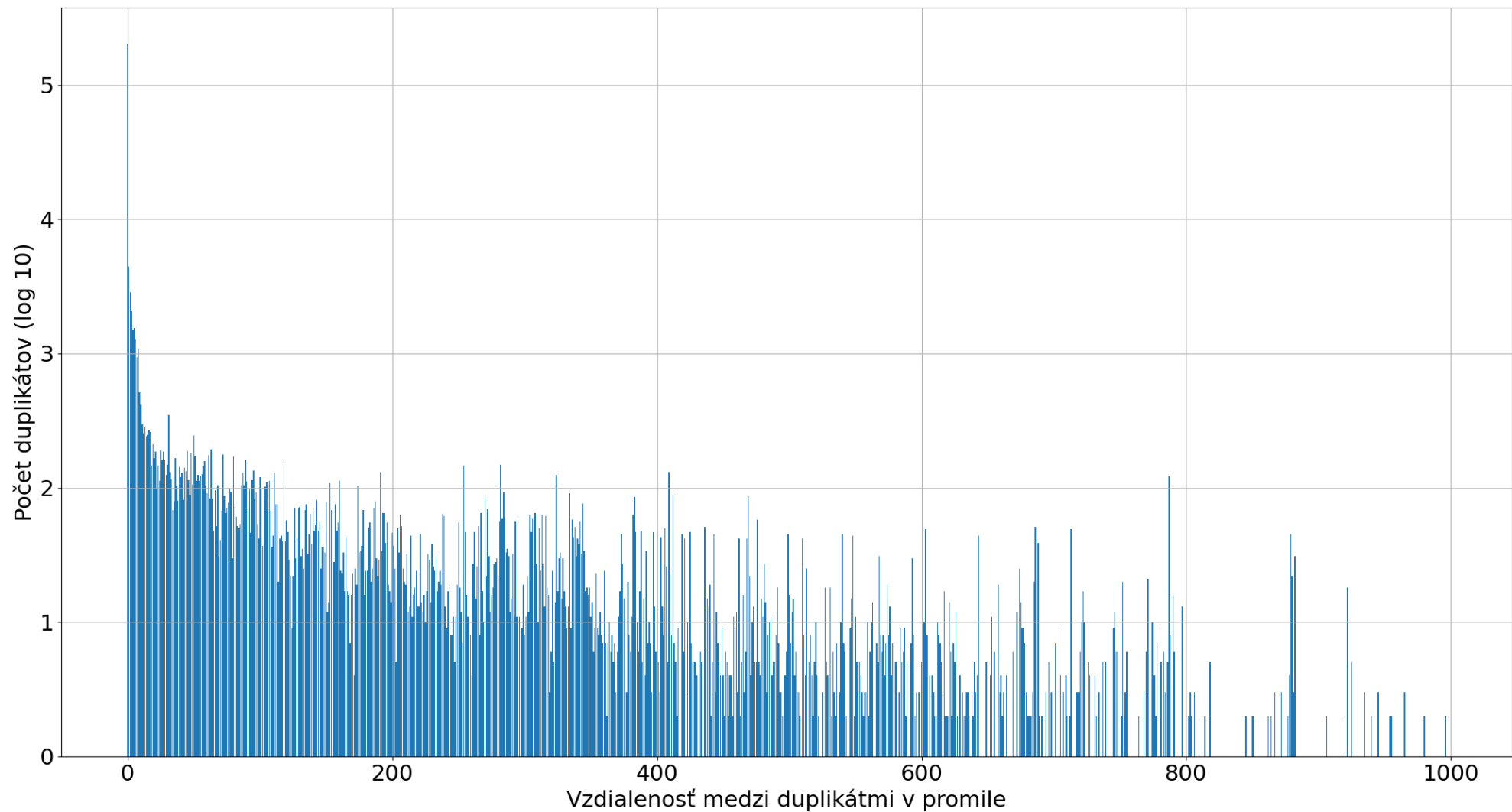
# TXT – veľkosť okna 32 bajtov



# HTML

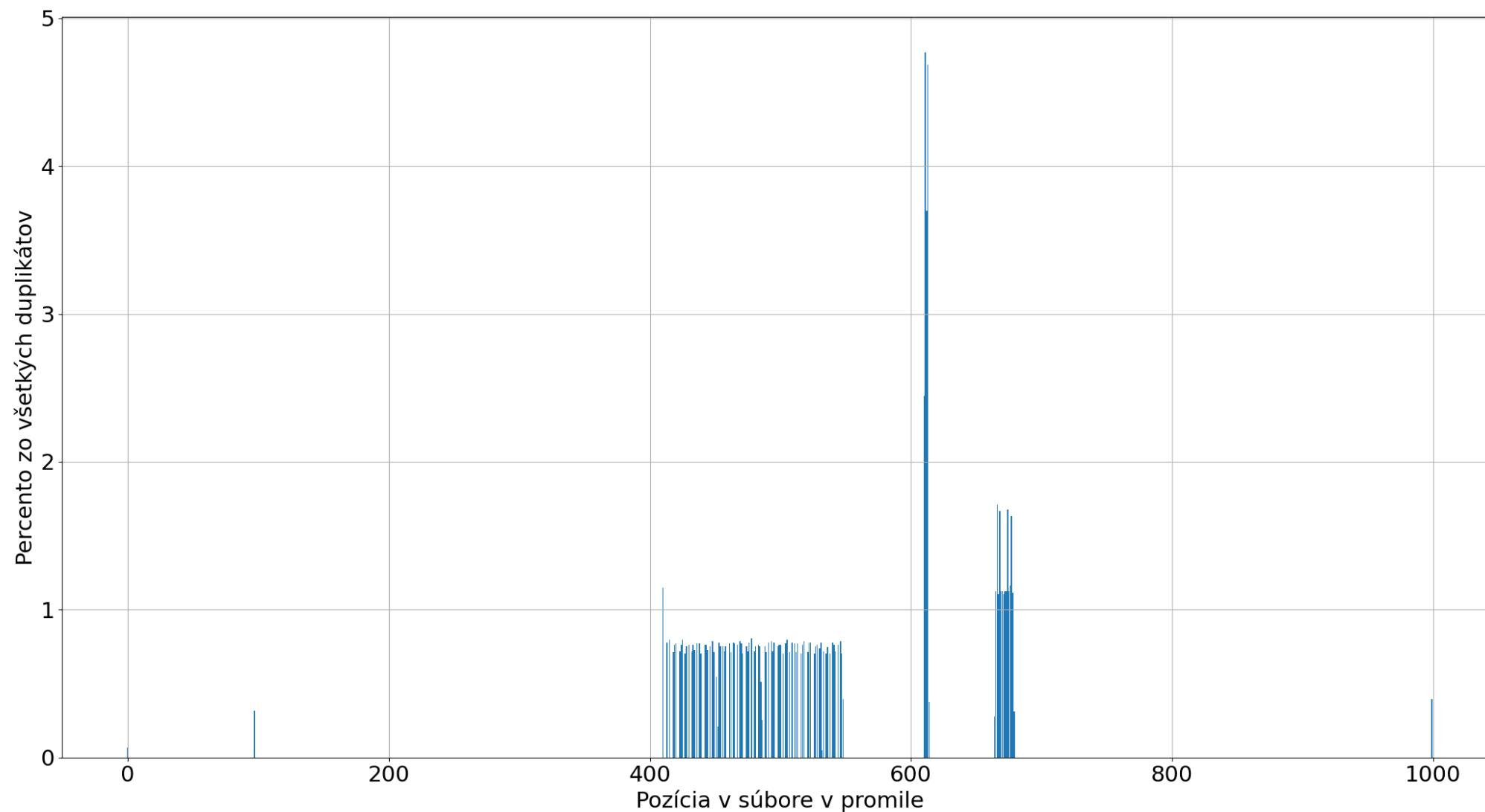


# PDF – veľkosť okna 64 bajtov

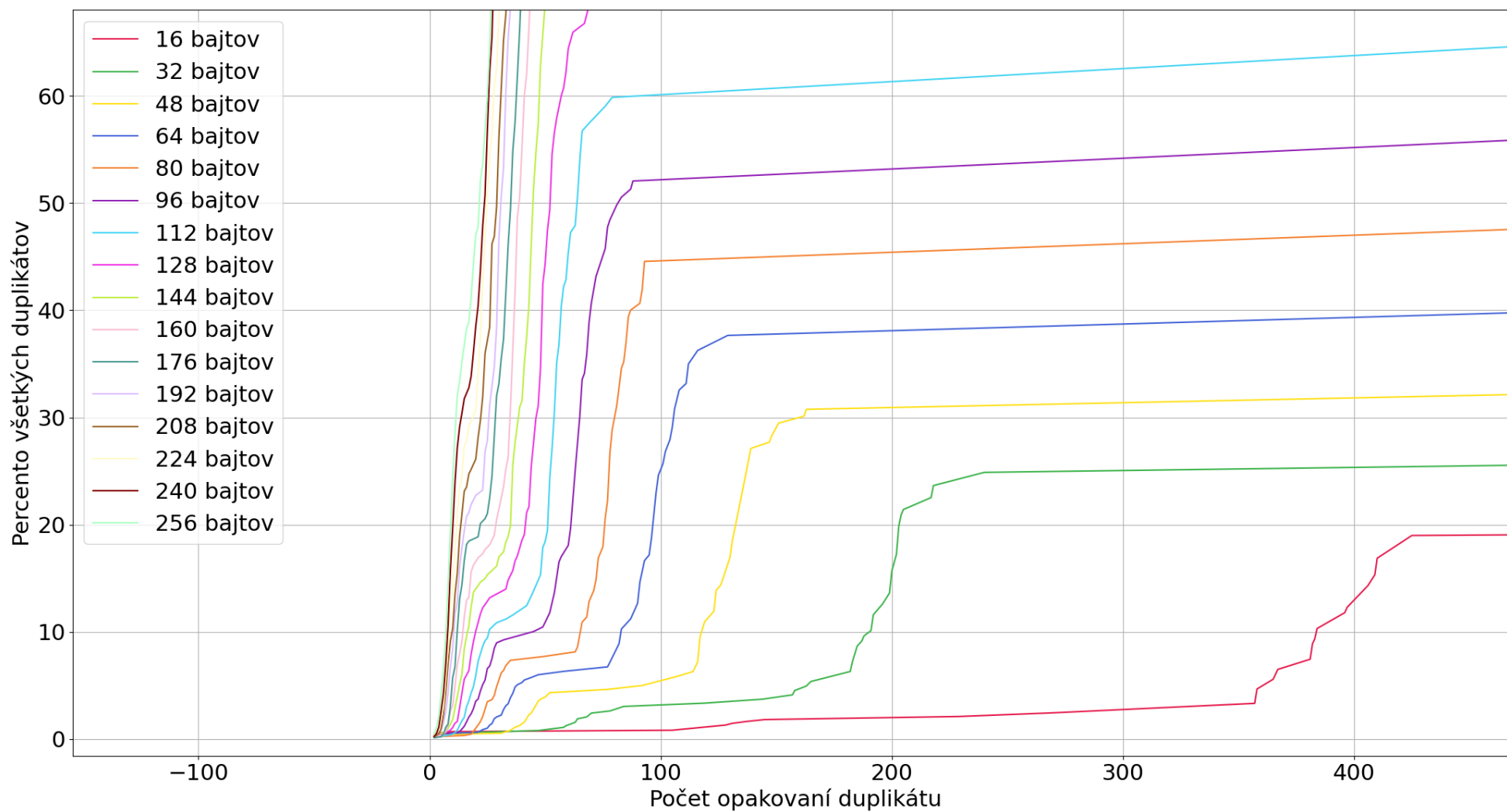




# MP4 – veľkosť okna 128 bajtov



# AVI



# Možné rozšírenia práce

- Využiť dáta
  - Vylepšenia
    - Nástroj
    - Súborové systémy
    - Komprimačné algoritmy
- Preskúmať ďalšie vlastnosti duplikátov
- Premennivá veľkosť okna

# Zhrnutie

- Duplikáty
  - Najviac redundancie
    - nekomprimované dáta
    - html a txt súbory
  - Početnosti duplikátov - Powerlaw
- Veľkosť okna
  - Čo najmenšia
    - Problém - režijné náklady
- Vzdialenosti medzi duplikátmi
  - Väčšina do 1 promile z veľkosti súboru