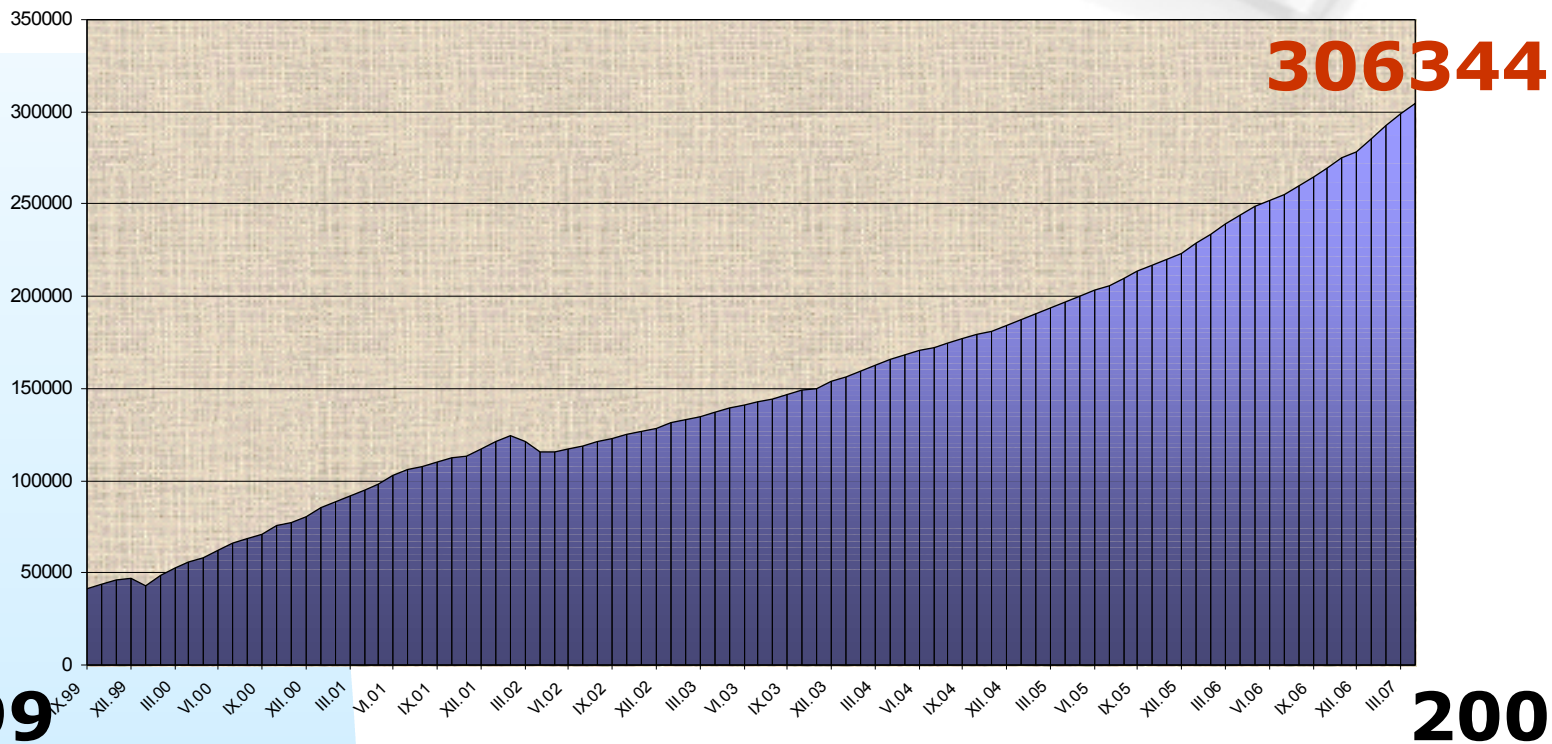




# Software pro archivaci webu

**Ing. Petr Žabička, MZK**  
**Mgr. Lukáš Matějka, MU**  
**Adam Brokeš, MU**

# Registrované domény v zóně .cz



1999

2007



## Software pro archivaci webu

Ing. Petr Žabička, MZK  
Mgr. Lukáš Matějka, MU  
Adam Brokeš, MU

### 1. Specifika webového archivu

- velmi velké množství souborů
- výpočet objemu dat
- velké množství datových formátů
- automatizovatelná archivace dokumentů
- neustálý vývoj standardů, protokolů...

### 2.3 Programy pro archivaci

- Mirador, wget, Teleport...
- nastavení na současně
- vhodné pro jednoválcové a menší projekty
- při velkých objemech dat problematické ze škálovatelnosti

### 2.4.1 Heritrix

- web harvester (nástroj pro stahování webu)
- modulární, rozšiřitelný, prošlá neustálým vývojem (včetně verze 1.12.1)
- robustní podpora pokročilého stahování webu
- datová média (java aplikace)
- snaha o co největší eliminaci stahování duplikátů a nevhodných dat (podpora 40-60% místa ušetřeno díky kompresi)
- obsahuje i parser pracující s modelem dokumentu (kvůli lepší analýze dokumentů, ale poměrně náročný)

### 2.4.3 Nutch

- velké dostupné modulární vyhledávací engine
- umí stáhnout a zpracovat miliony stránek měsíčně, spravovat jejich index, vyhledávat v něm 100x za vteřinu
- stránka škálovatelnost (výběr systémů hardware pro nastavení aplikace na více serverů)

### 2.4.4 NutchWAX

- nadstavba vyhledávacího rozhraní Nutch vytvořená pro potřeby indexování dokumentů archivovanými Heritrixem (ARC formát), přidává do indexu potřebná metadata, především časové razby

### 3 Implementace v rámci projektu WebArchiv

- stahování webu pomocí Heritrixu
- neomezené online zpřístupnění částí archivu pomocí WERA (včetně fulltextu)
- zpřístupnění celého archivu pomocí Wayback (jen URL, přístup k dokumentům jen z důvodu Waybacku)
- vlastní vývoj – databáze obsahující informace o všech souborech v archivu a její integrace s nástroji IPC
- příprava implementace OAI-PMH



### 2. Možnosti řešení

- 2.1 online archiv (Internet Archive)
- 2.2 vlastní online služby
- 2.3 programy pro archivaci „malého“ rozsahu
- 2.4 sada profesionálních nástrojů

### 2.4 Profesionální nástroje pro rozsáhlejší projekty

- konsorcium International Internet Preservation Consortium (IIPC - www.netpreserve.org)
- od 2007 Čestmír Němec a kol. Křehová ČS
- jedním z komerčních členů je Internet Archive, díky převážně velké národní knihovně
- konsorciem vyvíjí vlastní open source nástroje pro stahování, archivaci a zpřístupnění webu
- uživatelská základna Čentř IIPC zahrnuje sloužící i navý vyvíjených nástrojů
- IIPC nabízí mechanismy pro podporu vývoje nových aplikací podle potřeb svých členů
- standardizace – vývoj formátu WARC pro archivaci webu

### 2.4.2 Web Curator Tool

- nástroj pro správu sklizení, první verze uvolněná v září 2006
- nástroj pro Heritrix
- umožňuje správu sklizení méně kvalifikovaným uživatelům prostřednictvím grafické rozhraní a zpracování webových rozhraní
- výborná podpora uživatelských oprávnění
- nepodporuje inkrementální sklizení
- implementovaný workflow nepodporuje možnost daní naší registrací

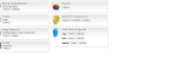
### 2.4.5 WERA

- Web Archive Access – zpřístupnění webového archivu
- velmi snadná navigace a zpracování uživatelské rozhraní (člověk ovládá časová verze dokumentu)
- zpracování archivovaných stránek lze i pomocí zadání přesné URL adresy
- archivované dokumenty v WERA prohledy skrz index NutchWAXe
- problémy s Javascriptem v některých stránkách
- vývoj ukončen, přechod na Wayback



### 2.1 Online archiv

- archive.org, web.archive.cz, ...
- vlastní využití existujících služeb
- omezená dostupnost dokumentů
- závislost na nástroji provozovatele služby



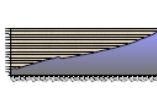
### 2.4.6 Wayback

- aplikace pro zpřístupnění archivu, která v budoucnu nahradí službu Wayback Machine Internet Archive
- dokumenty jsou indexovány a zpřístupňovány pomocí URL a času
- režim zpřístupnění:
  - archivní URL – úprava odkazů na stránce (klik zůstá do archivu)
  - proxy – chová se jako proxy server, ale je pak sice třeba časová verze (WAX Toolbar – slouží pro Firefox)
  - time-line – časová osa, zatím experimentální
  - připravuje se podpora fulltextového vyhledávání
  - dokončena lokalizace verze 0.9 do češtiny

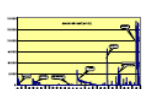


### 4. Statistika

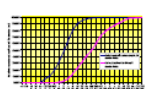
- počet sklizených dokumentů k 6.5.2007 je 134,5 milionů
- objem sklizených dat je 5 465 GB
- počet dokumentů v archivu k 3.5.2007
- probléma analýza zaměřená na sledení archivu podle typu souborů a jejich velikosti, zanalyzována na množství úspory úložného prostoru



Průměr registrací stránek v síti s cílem od října 1993 a 41.000 za více než 300.000



Přehledy souborů WebArchiv - Data



Stránky podle souborů a velikosti archivu při zavěšení stránek na veřejnost

<http://www.webarchiv.cz/>