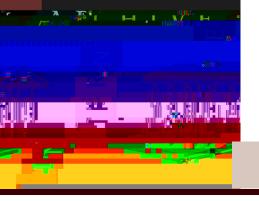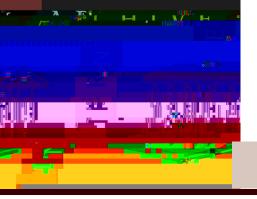Michael Stack

*stack*
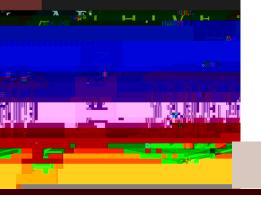
# What is Heritrix?

- Goal: open-source, extensible, web-scale, archival-quality, web crawling software

- Collaboratively developed
  - Internet Archive, IIPC, partner libraries, and others
  - Over seven releases since January 2004

- Technical parameters:
  - Built on many other open source libraries
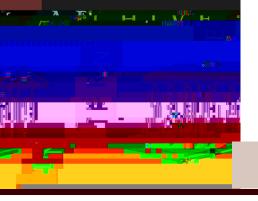  - Tested and supported on Linux

# Last year (IWAW'04)…

- Heritrix introduced at IWAW'04
  - Version 1.0
  - Core architecture, basic features in place
- Primarily useful for focused crawls
  - Hundreds to thousands of specific target web sites
  - Over 20 million collected URIs per crawl
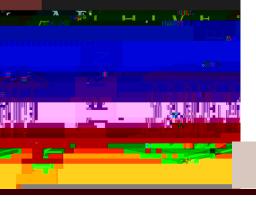  - Crawls run for up to a week

# Heritrix Releases: 1.2

- ## Heritrix 1.2: November 2004

  Motif: "requested features and fixes"

  – Better session-id handling ("URI canonicalization")
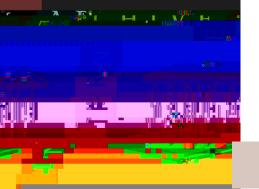
  – IP politen P ~4T"

l   e   d   f

e   r   ta      t–
i  o   n P      ~•Ő azd

# Heritrix Releases: 1.4

- Heritrix 1.4: April 2005

  Motif: "memory efficiency at scale"
  - Much im ̀ovd̀ s̈ingle-machine caq̀ci̊tyÖ
    - Big data structures movd̀ t̄o BDB JE.
  - Customizabt̄e 'decide rule' scoq̀e̊ȯ́P̈Ȯʀ×@
  - Imq̀o̊P̈×bealanced- ̀ogress and jV̀ac̀òn̄-trol ("q ̀V̈P̈V̊ @eting")
  - Exq̀ci̊ P̈Rmental "Adaq̀i̊ Rec̈t̊"ifrontier (b• Ki̊sti̊ꞑ Sgùr̀ðsson, Nì̄onal Libary of Iceland)
  - Exq̀ci̊ P̈Rmental ̀ogrammatic remote control (v̄a ̀"I̊MX" Jav ̀management extensio ̀Ö̊à
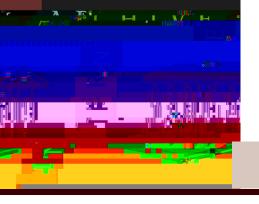
# Heritrix: Current Wo'''

# Heritrix: Future Plans (2.0+)

- Bigger & faster crawling
  - Automatic multi-machine coordination
  - Improved prioritizatza

# Das Ende

Thank you

*stack @ archive.org*