**Carnegie Mellon University**

# Introduction to Software Curation and Preservation

Eric Kaltman
*CLIR Fellow for Data Curation in the Sciences*

# Overview of Workshop

1. Why Software Preservation?
2. What is Software Curation?
   a. Scope and Description
   b. Migration
   c. Reproduction and Access
   d. Outreach and Culture
3. Trends
4. Efforts at CMU

**Carnegie
Mellon
University**

# Overview of Workshop

1. Why Software Preservation?
2. What is Software Curation?
   a. Scope and Description
   b. Migration
   c. Reproduction and Access
   d. Outreach and Culture
3. Trends
4. Efforts at CMU

**Carnegie Mellon University**

# What is Software?

Source Code

- Instructions to a computational device to do something

Executable binaries

- Compiled source code that executes on in a target environment

Documentation

- Descriptions of implementation and use

# What is Software?

Data encoded to be executed by a specific computational context ('execution context')

Execution Context - the collection of hardware and software dependencies required to allow for execution

Hardware Dependencies - CPUs, system architecture, peripherals, displays, etc.

Software Dependencies - Compiled libraries, support programs, APIs, etc.

Carnegie
Mellon
University

# Why Software Preservation?

History
- Modern society is built on and from interactions with software

Progress
- Modern research, regardless of the field, is tied to software infrastructures

Pedagogy
- Access to legacy software and systems for teaching and critical engagement.

Practicality
- Software at the root of commercial progress, ties to competition and legal issues

**Carnegie Mellon University**

# History and Reproducibility

Research reproducibility

- Maintaining citable access to previous results, analyzes and data sets

- Shifting research practices toward an archival and curatorial mindset

- Models for sustainable software development

History, legacy and maintenance

- Preserving the historical record, the intellectual history of humanity

- Aligning with a maintenance narrative instead of an innovation one

- Access to the past will be dependent on software for the rest of time

*Access to software is also access to the files produced and interpreted by software*

# Digital Dark Ages

"People think that bits are somehow immortal because somehow they're this ethereal thing in cyberspace...It could be that the format of those bits, the way in which they are interpreted requires a piece of software to figure out what the bits mean. How they should be presented as an image or a video or how you should interact with it in a spreadsheet, but the software doesn't exist. What if the operating system that the old software used to run on doesn't exist anymore? What if the latest software doesn't know how to read the formats of those complex digital objects? Guess what? That information is gone.

There isn't a systematic way to ensure the information that we create today will still be usable 100 years from now. That's why I'm worried about a *digital dark ages*.

*Vint Cerf, co-inventor of TCP/IP*

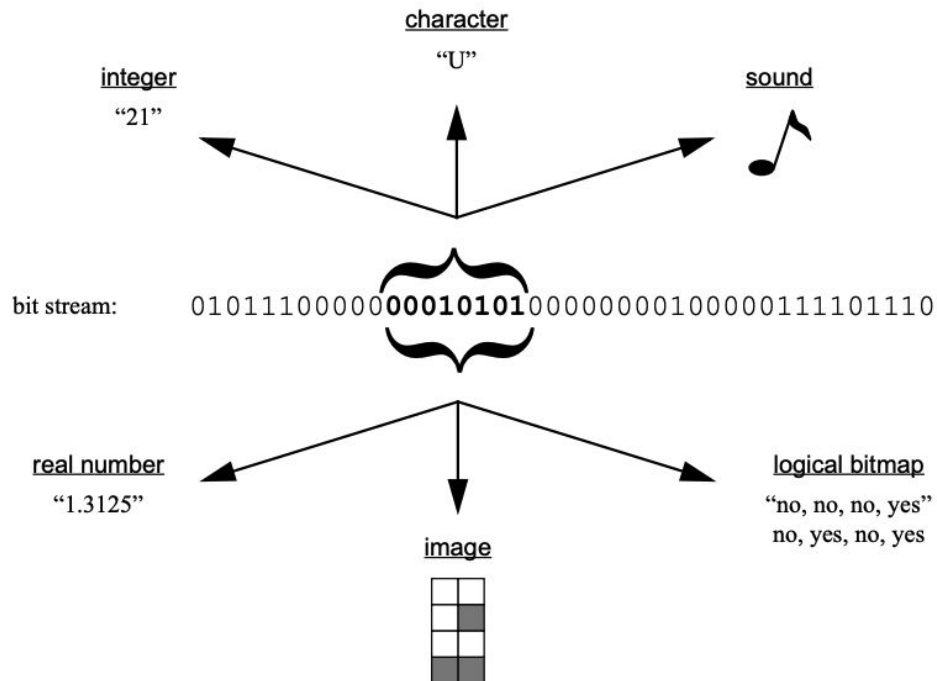# Ensuring the Longevity of Digital Information



Figure 4: A bit stream can represent anything at all

Jeff Rothenberg
Scientific American
1995

# Overview of Workshop

**Carnegie Mellon University**

# Curation

Scope

- What is historically important?

- What will be important for future research efforts?

Description and Standards

- How do we describe software for future access, study and use?

Migration

- How do we transition and preserve data across time?
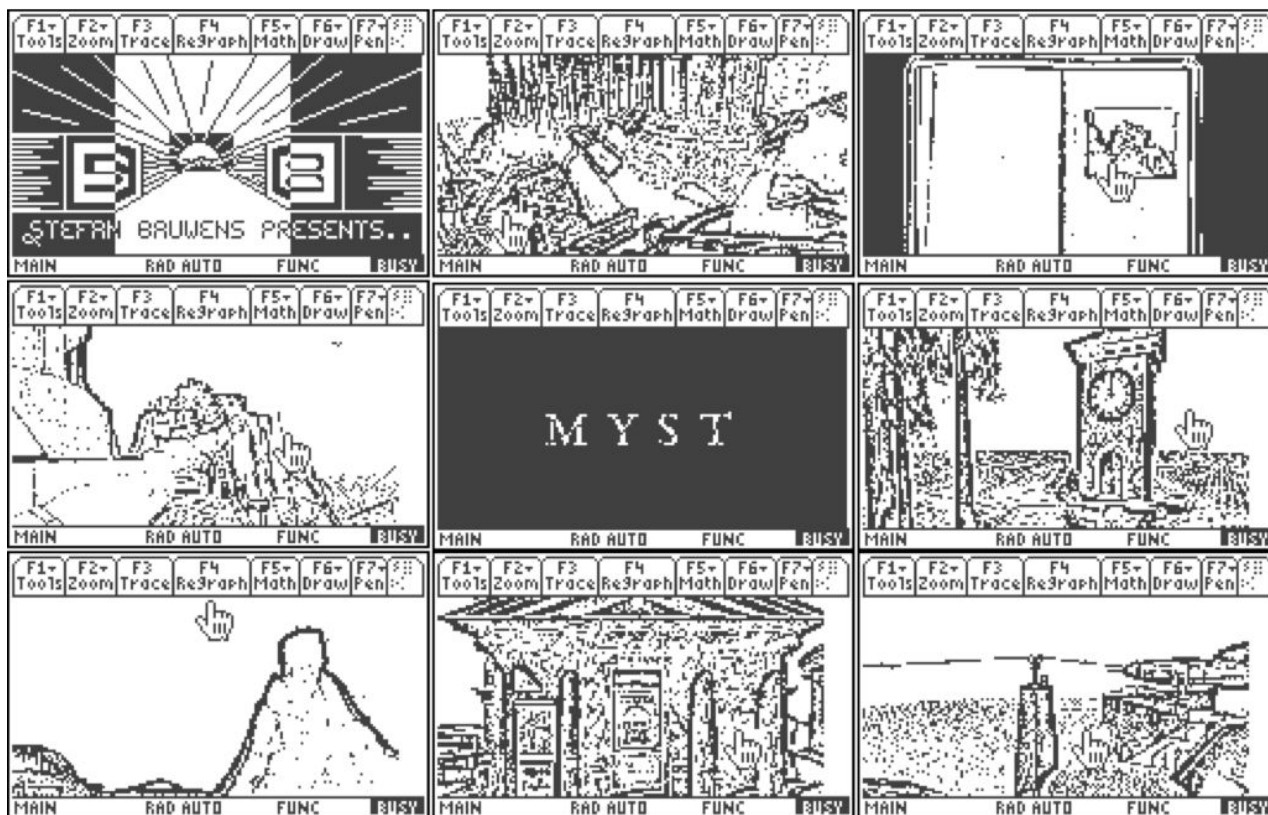
Storage and Access

- Where do we put these things?

- How do we find them again?

# What is Important?

A medical supply company in Miami had received a delivery of botulin, which was to be processed into Botox and distributed. However, it was misprocessed, and a dangerous concentrate was distributed. The FDA had all of the information needed to identify the recipients, but the information was in a file created with a 2003 version of a popular business software application. The 2004 version available to the FDA could not open the data file. The manufacturer of the software was also unable to supply the relevant version.

It so happened that one of the agents involved in the case was familiar with the NSRL, and had in fact provided software to us earlier in the year. He called, explained the situation, and asked if we had the 2003 version of the software. We did! The agent then arranged for an FDA contact to come to NIST, get the software, and put it on a jet to Miami. The people working the case in Miami were able to install the old version, open the data file, and trace the paths of the botulin.
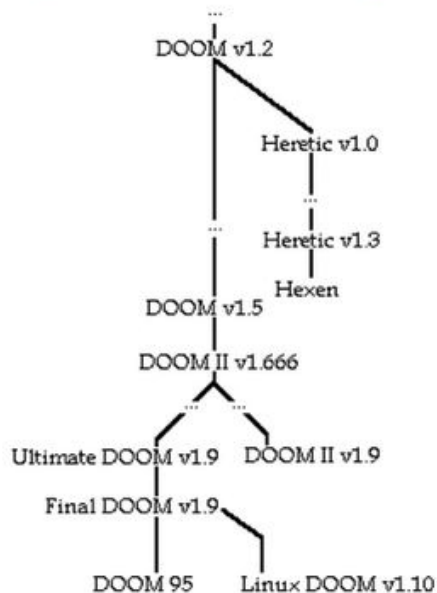
*Preserving.Exe: Toward a National Strategy for Software Preservation - 2013*

MARIO x 6   COINS:32   ▶▶▶▶▷ [P]   0027800

DOOM v1.2
Heretic v1.0
Heretic v1.3
Hexen
DOOM v1.5
DOOM II v1.666
Ultimate DOOM v1.9    DOOM II v1.9
Final DOOM v1.9
DOOM 95    Linux DOOM v1.10

| Filename | File Size | Release Date | notes |
|---|---|---|---|
| doom0_2.zip | 254k | Feb. 4, 1993 | alpha version - doom0_2.zip downloaded by HEL from www.doomworld.com/pageofdoom, 18 July 2008 |
| doom0_4.zip | 950k | Apr. 2, 1993 | alpha version |
| doom0_5.zip | 1264k | May 22, 1993 | alpha version |
| doom_pre_beta.zip | 2555k | Oct. 4, 1993 | press-release beta version |
| doom1_0.zip | 2113k | Dec. 10, 1993 | v0.99 shareware version |
| doom1_1.zip | 2160k | Dec 16, 1993 | v1.1 shareware version |
| doom1_2.zip | 2203k | Feb 17, 1994 | shareware version |
| doom14bt.zip | 2246k | Jun. 28, 1994 | v1.4beta shareware version |
| doom15bt.zip | 2262k | Jul. 8, 1994 | v1.5beta shareware version |
| doom16bt.zip | 2234k | Aug. 3, 1994 | v1.6beta shareware version |
| dm1666sw.zip | 2293k | Sep. 1, 1994 | v1.666 shareware |
| doom_v18.zip | 2423k | Jan 23, 1995 | v1.8 shareware |
| doom19s.zip | 2393k | N/A | v1.9 shareware version |

0.2 Alpha
0.3 Alpha
0.4 Alpha
0.5 Alpha
Press Release
1.0
1.1
1.1 MS-DOS Extenders[1]
NeXTStep
Heretic
1.2
Atari Jaguar
Sega 32X
1.25 Sybex
1.3[2]
1.4 Pre-Release
GameBoy Advance
3DO
1.4 Beta
1.5 Beta
1.6[3]
1.6 Beta
1.666 Early
Strife
Doom II 1.666
1.666
Sony PlayStation
Doom II 1.666 Germany
Doom II 1.7a
Sega Saturn
Doom II 1.7b
1.8
WinDoom
Doom II 1.8 France
1.9
Doom 95
1.9 The Ultimate Doom
Super Nintendo[4]
Final Doom
Final Doom Sony PlayStation
Chex Quest
Final Doom id Anthology
Pocket PC
1.10 Linux Source
Nintendo 64[5]
Microsoft Xbox
Microsoft Xbox 360
Doom Classic Source
Doom 3 BFG Edition
iOS
Doom Classic Complete Sony PlayStation 3

# GNU/Linux Distributions Timeline

**Version 17.10**

© Andreas Lundqvist, Donjan Rodic, Mohammed A. Mustafa
© Konimex, Fabio Loli and contributors
https://github.com/FabioLolix/linuxtimeline
Original source: futurist.se/gldt
Published under the GNU Free Documentation License

- Influence, developer switching
- Rebasing, substantial code flow, project overtaking
- Developer & code sharing, project merging

1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017
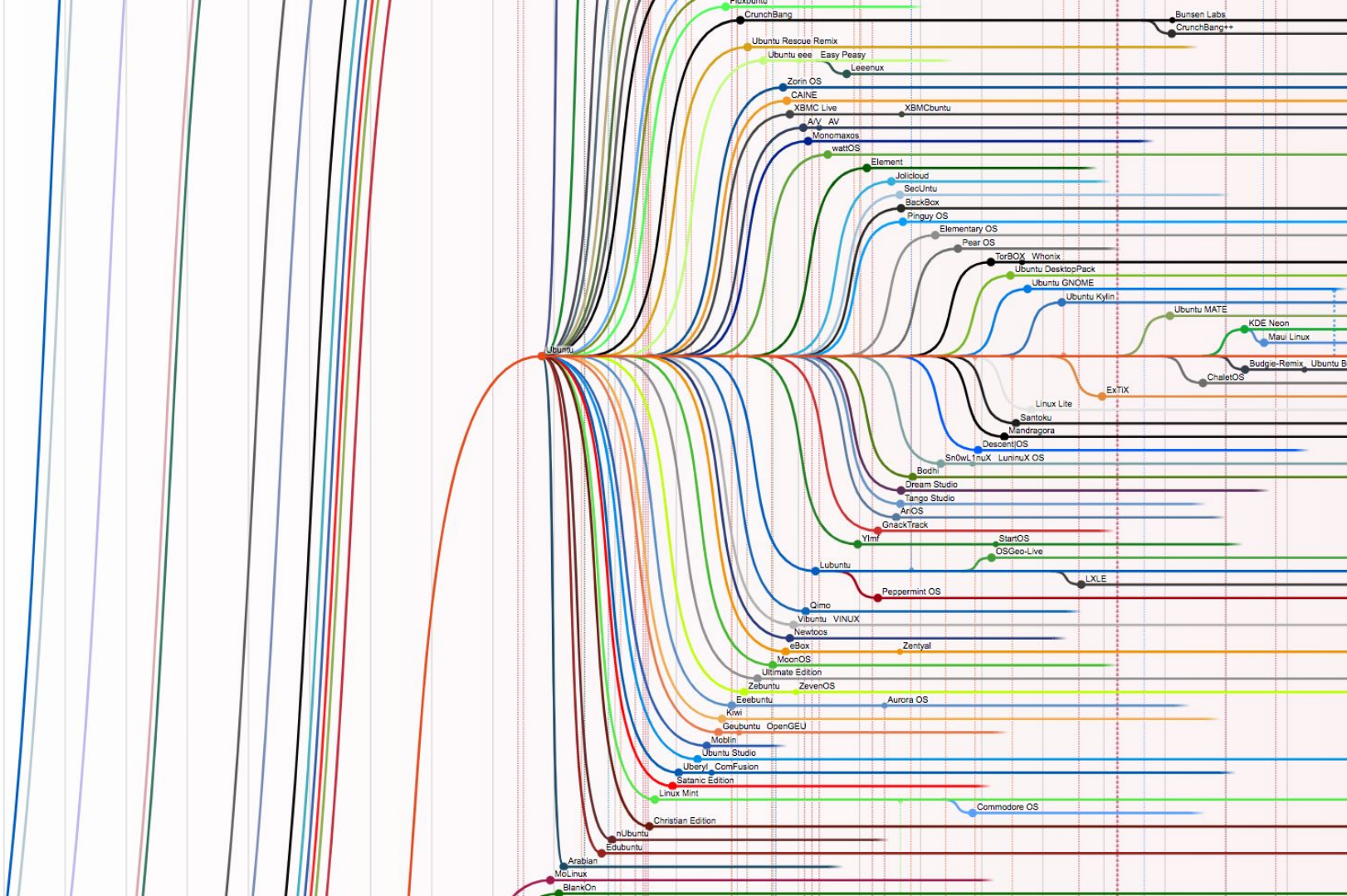
Libranet
Omoikane (Arma)
Quantian
Damn Small Linux
Damn Vulnerable Linux
KnoppMyth
Danix
Parsix
Kanotix
Auditor Security Linux
Backtrack
Kali
B2D
Whoppix
WHAX
WHAX
Symphony OS
Musix
ParallelKnoppix
Knoppix
Kaella
MAX
Feather
Medialinux
Mediainlinux
ArtistX
INSERT
Aquamorph
Dreamlinux
Morphix
ZoneCD
Hiwix  Hiweed
Deepin
Kalango
Kurumin
Poseidon
Dizinha
NeoDizinha  Patinho Faminto
Skolelinux
DebianEdu
Lindows
Linspire
Freespire
Rxart
MEPIS
SimplyMEPIS
antiX
Impi
Swift
Bluewall
K-DEMar
kademar
Euronode
DeadCD
Olive
Underground Desktop
Ulteo
Kubuntu
Polippix
Netrunner
DEFT
Asturix
Bardinux
Gobuntu
Runtu
Voyager
Xubuntu
GalliumOS
Xinutop
Peach OSI
PC/OS
OS4
Black Lab
PUD
xPUD
gNewSense
Muslim Edition
Sabily
Madbox
Mythbuntu
Ubuntulite
U-lite
Greenie

DNALinux
Splack Linux
Tiny
Burapha
Caixa Mágica
Stresslinux
Linkat
S.u.S.E
SuSE
SUSE
Keysoft
EasyNAS
GeckoLinux
OpenSUSE
JOPUX
Astaro
Sun JDS
United Linux
Sophos UTM
Caldera
SCO
UltraPenguin
Redmond Lycoris
Buhawi
Eurielec
ALT
SAM
SAMity
Mandrake
OpenSLS Annvix
Mandriva
Mageia
ROSA
blackPanther
Unity Linux
OpenMandriva Lx
Granular
Phinx
PCLinuxOS
Garuda
TinyMe
Demolinux
KRUD
Eridani
Vine
Armed
Kondara
ELX
Finnix
Miracle
Asianux
Pingo
Rocks
ASP
Independence
HP Secure
EnGarde
BLAG
Berry
ATmission
Momonga
MythDora
Ekaaty
Vixta
Simplis Xange
eZeY
Moblin 2
MeeGo
Mer
Sailfish OS
Tizen
Hanthana
Synergy
Fuduntu
Parsidora
Maui OS
Hawaii OS
Fedora Core
Fedora
Chapeau
Viperr

Synergy
Fuduntu
Parsidora
Maui OS
Hawaii OS
Fedora Core
Fedora
Chapeau
Viperr
Korora
Fusion
Qubes OS
VortexBox
Ojuba
Amahi
FoX
AsianLinux
NST
Elastix
OpenNode
CentOS
NethServer
Baruwa
Rockstor
BlueOnyx
Asterisk@Home
trixbox
CERN
StartCom
Endian
Red Hat Enterprise
ServOS
Oracle Enterprise    Oracle Linux
Tao
Scientific
White Box
PUIAS
Springdale Linux
SulIX
AnNyung
Aurox
Bayanihan
Aurora
K12
SuperRescue
ClarkConnect
ClearOS
Best
SOT
LBA
Happy
BU Linux
Trustix
Linpus
Immunix
Yellow Dog
Red Flag
e-smith
SME Server
Fermi
Turbolinux
PLD
Conectiva
LinuxPPC
Red Hat
WGS Linux Pro

Gentoox
Knopperdisk
epiOS
Kororaa
Papug
Toorox
Funtoo
Chrome OS
NayuOS

Flint OS

CloudReady

Daphile

Enoch    Gentoo    Porteus Kiosk

wtfplay-live

CoreOS    Container Linux

Liberté

Nova

SystemRescueCD

Calculate

Pardus

RR4    Sabayon

Pentoo

VidaLinux    VLOS

Ututo    Ututo XS

Ututo-e

Chakra    KahelOS

Arch Linux ARM

Parabola

Bridge

Cinnarch    Antergos

AudioPhile Linux

Apricity

MorpheusArch Linux

ARCHLabs

Arch    OBRevenge OS    Rev

VeltOS

Obarun

BlackArch

PoliArch

Manjaro    Sonar

Netrunner Rolling

ArchBang

AL-AMLUG    Archie    CTKarchLive    CTKArch

LinHES

MCC Interim

TAMU

Yggdrasil

DLD

LST

Bogus

Xdenu

Linux-FT

Mini

Jurix

Trans-Ameritech

Unifix

Linux Universe

Craftworks

DILINUX    DOSLINUX

mkLinux

Monkey

Linux Router Project

LEAF

Weaver    Nitix

uClinux

ROCK

T2

Project Ballantin

FREESCO

tomsrtbt
Coyote
eiT easyLinux
ELinOS
Peanut
BluePoint
aLinux
NuTyX
Linux From Scratch
ZENIX
AryaLinux
KaarPux
KaeilOS
SmoothWall GPL
SmoothWall Express
IPCop
IPFire
CRUX
Beehive
Midori
Leka Rescue Floppy
Openwall
dyne:bolic
OpenWRT
Lede Project
Ark
NetStation
Thinstation
LPS
TENS
Source Mage
Sorcerer
Lunar
Tinfoil Hat
LinuxConsole
GoboLinux
Yoper
UHU
GeeXboX
Macpup
Simplicity
Puppy
Sage Live CD
TEENpup
Legacy OS
Quirky
Devil
NixOS
GuixSD
QiLinux
Natures Linux
Openfiler
Octoz
Hedinux
Specifix
rPath
Foresight
Paldo
BrazilFW
Jarro Negro
Ophcrack
Alpine
Everest
Qomo
Zeroshell
Parted Magic
SliTaz
Tiny SliTaz
openmamba
Syllable Server
Ångström
PLoP
Exherbo
dCore
Tiny Core

Puppy
TEENpup
Sage Live CD
Legacy OS
Quirky
Devil
NixOS
GuixSD
QiLinux
Natures Linux
Openfiler
Octoz
Hedinux
Specifix    rPath
Foresight
Paldo
BrazilFW
Jarro Negro
Ophcrack
Alpine
Everest    Qomo
Zeroshell
Parted Magic
SliTaz
Tiny SliTaz
openmamba
Syllable Server
Ångström
PLoP
Exherbo
dCore
Tiny Core
Nanolinux
piCorePlayer
Dragora
webOS
LuneOS
LibreELEC
OpenELEC
Lakka
CloudLinux OS
0Linux
4M
noop
Bedrock
AOSC OS
KaOS
Minimal Linux Live
Pisi Linux
Evolve OS   Solus
Clear Linux
RancherOS
Super Grub2 Disk
Void
OviOS
Phoenix OS
Android-x86
Remix OS
Android Wear
Android
Project Brillo    Android Things
Android-IA
LineageOS
CyanogenMod
Cyanogen OS

| 1992 | 1993 | 1994 | 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Carnegie
Mellon
University

# Overview of Workshop

**Carnegie Mellon University**

# Migration

Imaging and Storage

Digital Forensic Workflows

Digital Rights Management
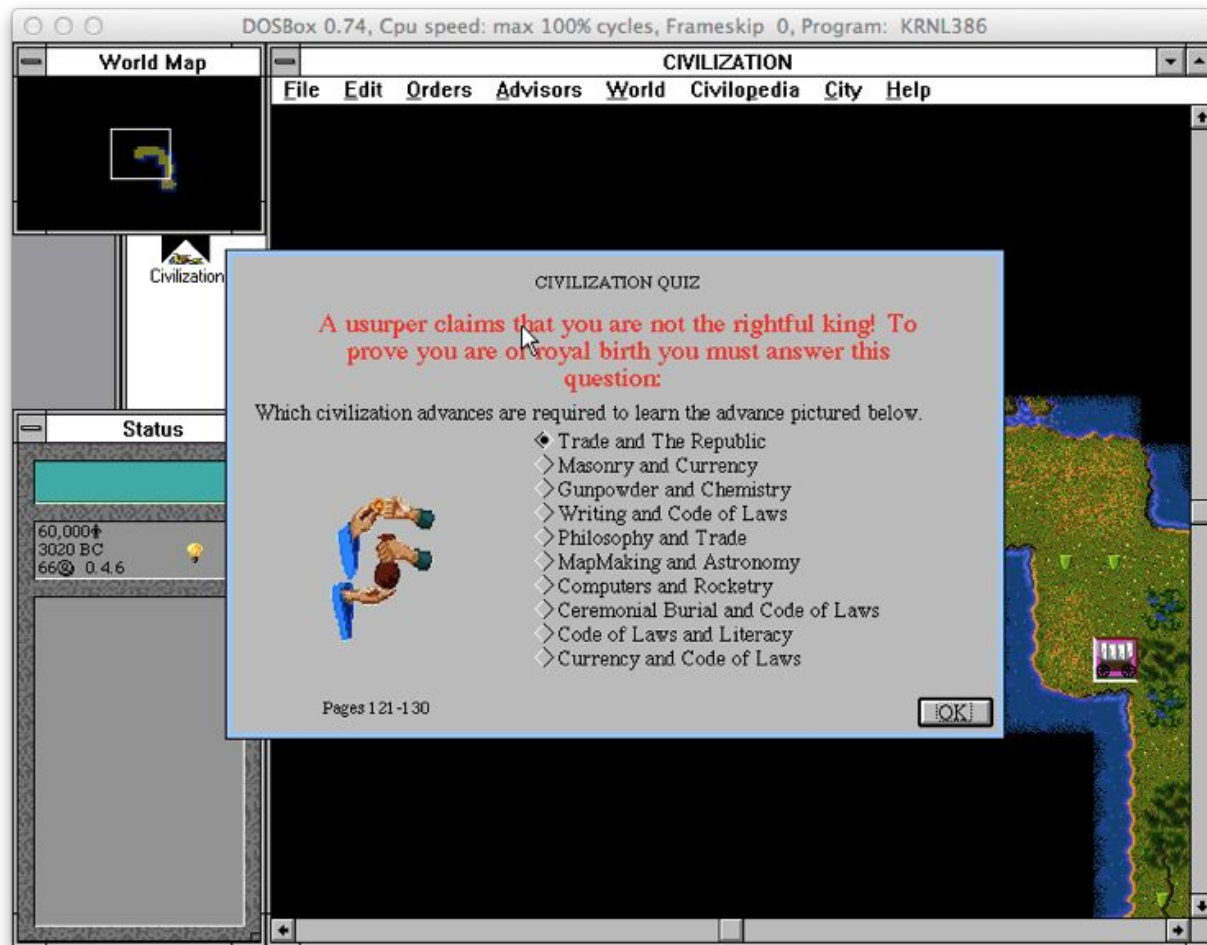
Digital Migration Strategies

# Migration

1. Reset button
2. Status LEDs
3. MOLEX connector
4. Power rail select
5. External power supply connector (rev. B: +5V, rev. C-latest: +7V)
6. Floppy disk drive power connector
7. USB B connector
8. JTAG connector
9. Flash erase (leave open!)
10. Drive select jumpers
11. Write blocker
12. Floppy disk drive connector



Carnegie
Mellon
University

# Digital Migration Strategies

Once imaged or ingested, need a *forever* storage strategy

- Repositories change and evolve, data needs to migrate with those changes

- Provenance information becomes more important

Current repositories are still playing significant catch up

No current long-term solutions available for software data

# Overview of Workshop

**Carnegie
Mellon
University**

# Reproduction Strategies

Virtualization

Containerization

Emulation

Hardware Preservation

Executable

```
010100100010010010
100010101101001011
010101010101100010
010101010101010101
101001010101010010
101010101010101001
101010101010010100
101010110101010010
101010101010101010
```

INTERFACE

Host Platform

Emulator

Container

Virtual Machine

Carnegie
Mellon
University

# Client Side Emulation

Local executable emulation

JavaScript in browser emulation

Pros:
- Lower latency
- Locally inspectable
- If browser based, shareable and single requirement

Cons:
- Dependent on specific system configuration
- Not easily shareable if local executable
- Legally dubious

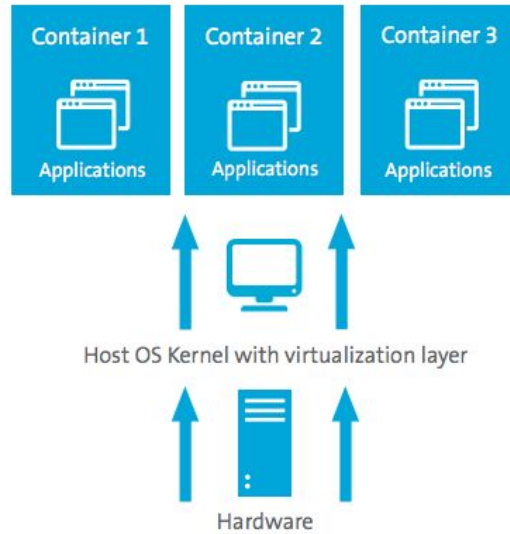# Server Side Emulation

Emulation in cloud

Pros:
- Legally more appealing
- Management and maintenance are centralized
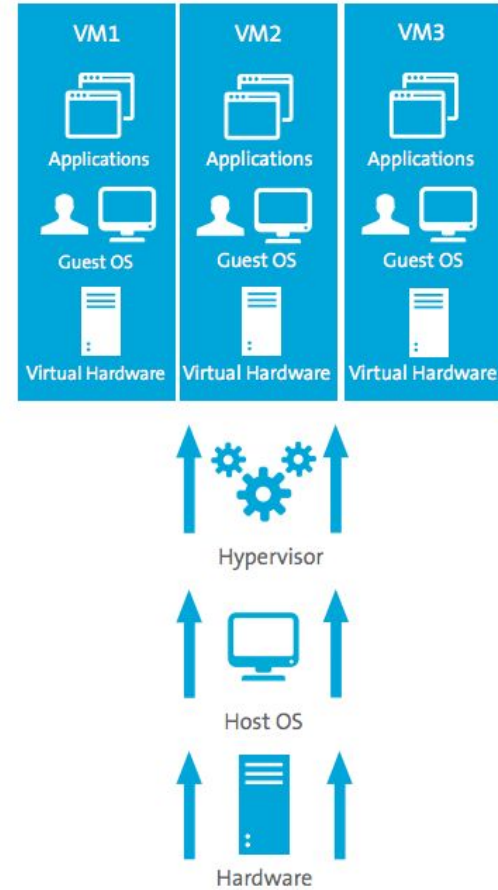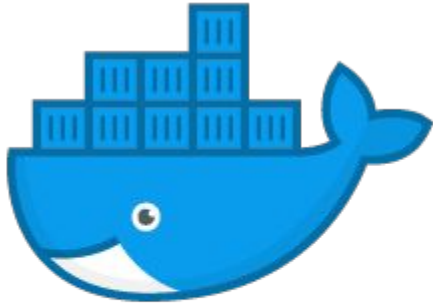- More easy to roll into services

Cons:
- Latency issues
- Much less introspection
- Cloud is a preservation issue in of itself

**Carnegie Mellon University**

# Hardware Preservation

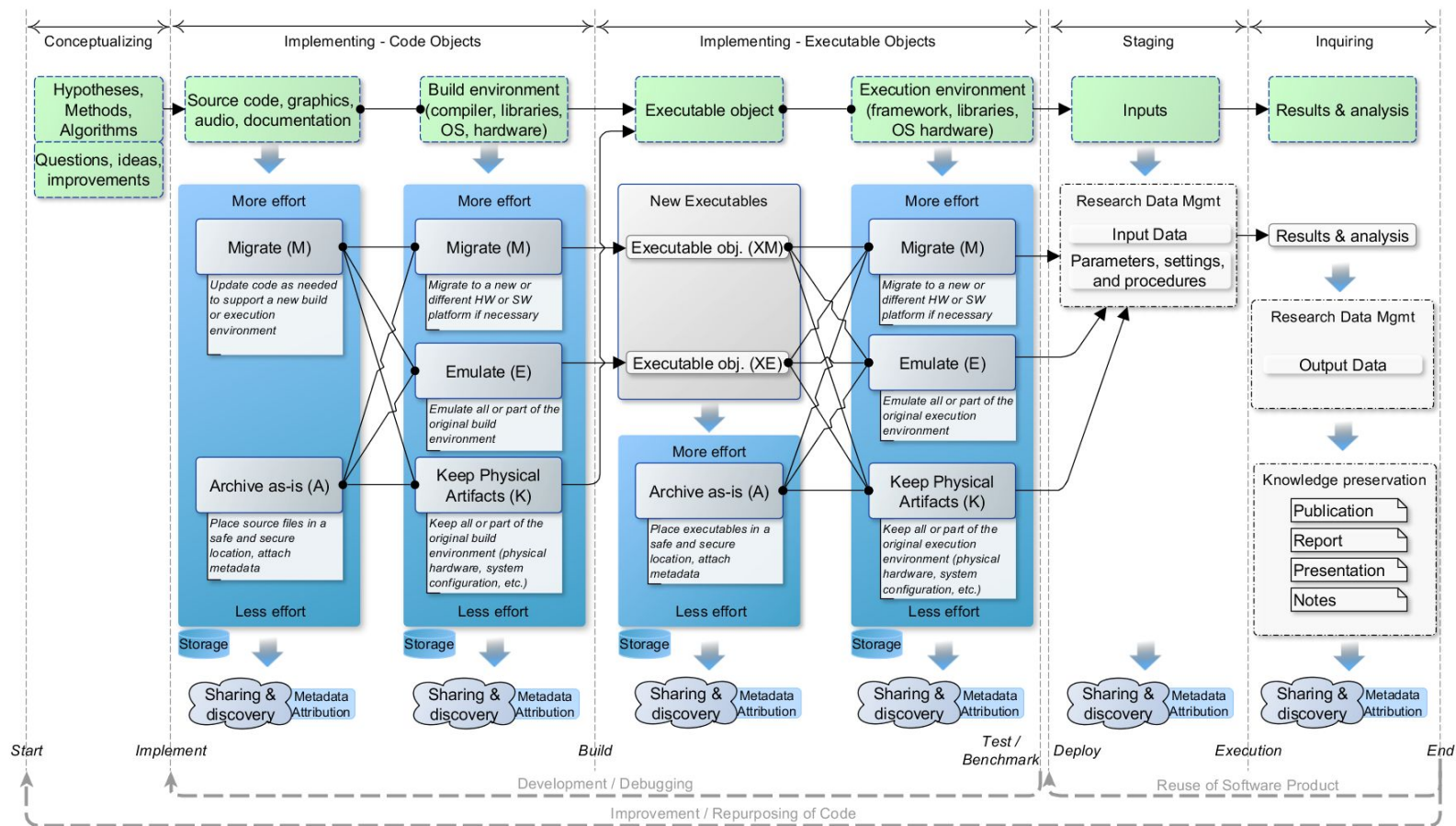Cultural software objects are designed for specific hardware and interactions

- Computer Games
- Interactive Art Installation
- Digital Art
- Other digital media works

Socio-cultural context is not in a VM or container

Hardware peripheral and displays are not replicated

**Hypotheses, Methods, Algorithms**

**Questions, ideas, improvements**

**Source code, graphics, audio, documentation**

**Build environment (compiler, libraries, OS, hardware)**

**Executable object**

**Execution environment (framework, libraries, OS hardware)**

**Inputs**

**Results & analysis**

### More effort

**Migrate (M)**
*Update code as needed to support a new build or execution environment*

**Archive as-is (A)**
*Place source files in a safe and secure location, attach metadata*

### Less effort

Storage

### More effort

**Migrate (M)**
*Migrate to a new or different HW or SW platform if necessary*

**Emulate (E)**
*Emulate all or part of the original build environment*

**Keep Physical Artifacts (K)**
*Keep all or part of original build environment (physical hardware, system configuration, etc.)*

### Less effort

Storage

### New Executables

**Executable obj. (XM)**

**Executable obj. (XE)**

### More effort

**Archive as-is (A)**
*Place executables in a safe and secure location, attach metadata*

### Less effort

Storage

### More effort

**Migrate (M)**
*Migrate to a new or different HW or SW platform if necessary*

**Emulate (E)**
*Emulate all or part of the original execution environment*

**Keep Physical Artifacts (K)**
*Keep all or part of the original execution environment (physical hardware, system configuration, etc.)*

### Less effort

Storage

### Research Data Mgmt
**Input Data**
Parameters, settings, and procedures

**Results & analysis**

### Research Data Mgmt
**Output Data**

### Knowledge preservation
Publication
Report
Presentation
Notes

Sharing & discovery — Metadata Attribution (×6)

Start | Implement | Build | Test / Benchmark | Deploy | Execution | End

Development / Debugging

Reuse of Software Product

Improvement / Repurposing of Code

**Legend**

●—● Non-directional link between two workflow components. Movement can occur in either direction with either node as the start node.

Action boundary

Initiate a preservation or sharing activity (producer's perspective)

→ Directional link between two workflow components. Indicates an action or result.

Original workflow

Preservation activity (P)
(P) is a short identifier for the preservation activity or result

Non-software activities

Software preservation continuum (creator's perspective)

Data Management Services
Johns Hopkins University
http://tiny.cc/s16xay

# Overview of Workshop

**Carnegie Mellon University**

# Outreach

Locating significant materials around CMU community

Implement reproducible practices inside research labs and departments

Sustainable software development

Carnegie
Mellon
University

# Overview of Workshop

1. Why Software Preservation?
2. What is Software Curation?
   a. Scope and Description
   b. Migration
   c. Reproduction and Access
   d. Outreach and Culture
3. Trends
4. Efforts at CMU

# Sustainable Software Development

Testing oriented

Well documented

Instrumentation

- Continuous integration

- Package management

- Configuration management

- Version Control

# Tools and Organizations

Reproducibility Tools Supporting Software

- Code Ocean - http://codeocean.com

- Occam - http://occam.cs.pitt.edu

- Collective Knowledge - http://cknowledge.org/

- Umbrella - http://ccl.cse.nd.edu/software/umbrella/

- ReproZip - https://www.reprozip.org/

Organizations supporting software preservation

- Software Sustainability Institute (SSI) - UK Organization - http://www.software.ac.uk

- Data and Software Preservation for Open Science (DASPOS) - CERN - http://daspos.org

- Software Preservation Network (SPN) - US Memory Institutions

  - http://www.softwarepreservationnetwork.org

Carnegie
Mellon
University

# ACM Reproducibility

- Repeatability (Same team, same experimental setup)

- Replicability (Different team, same experimental setup)

- Reproducibility (Different team, different experimental setup)

**Carnegie Mellon University**

# ACM Badging Levels

Artifacts Available
- Software and data present in publication are available for download and investigation

Artifacts Evaluated - Functional
- Software and data have been audited and validated as working

Artifacts Evaluated - Reusable
- Software and data have been audited by a third party, are functional, and significantly oriented toward reusability through documentation, code / software organization, etc.

Results Replicated
- Software and data have been used to validate results

Results Reproduced
- Different software and data have been used to validate results

**Carnegie Mellon University**

# Overview of Workshop

1. Why Software Preservation?
2. What is Software Curation?
   a. Scope and Description
   b. Migration
   c. Reproduction and Access
   d. Outreach and Culture
3. Trends
4. Efforts at CMU

**Carnegie
Mellon
University**

# Tools and Services at CMU

Kilthub Repository
- Source code
- Research data and software executable binaries

Code Ocean (Beta)
- Targeted for August 2019
- Reproducibility platform

Emulation as a Service (EaaS)
- Currently in research beta
- Distributed containerized execution contexts

Software and Data Carpentries
- Two day courses on basic research support tools like Python, R, and Git

History of Science and Technology at CMU (HOST@CMU)
- Interdisciplinary initiative to locate and celebrate CMU technical history

Carnegie
Mellon
University

Carnegie
Mellon
University

Browse    Search on KiltHub    Log in

# KiltHub

Discover research from Carnegie Mellon University ▾    + Follow

NEW    POPULAR    CATEGORIES    SEARCH

**19,165** posts    **1,750,937** views    **7,196,680** downloads    more stats...

**Spiritual Narratives in Beethoven's Quartet, Op. 132**
John Ito    22/03/2019

**Koch's Metrical Theory and Mozart's Music: A Corpus Study**
John Ito    22/03/2019

**Income Mobility in America**
Manu Navjeevan    21/03/2019

**Monocular Facilitation Implicates Subcortical Involvement in Holistic...**
Rebeka C. Almasi    21/03/2019

University

# Code Ocean
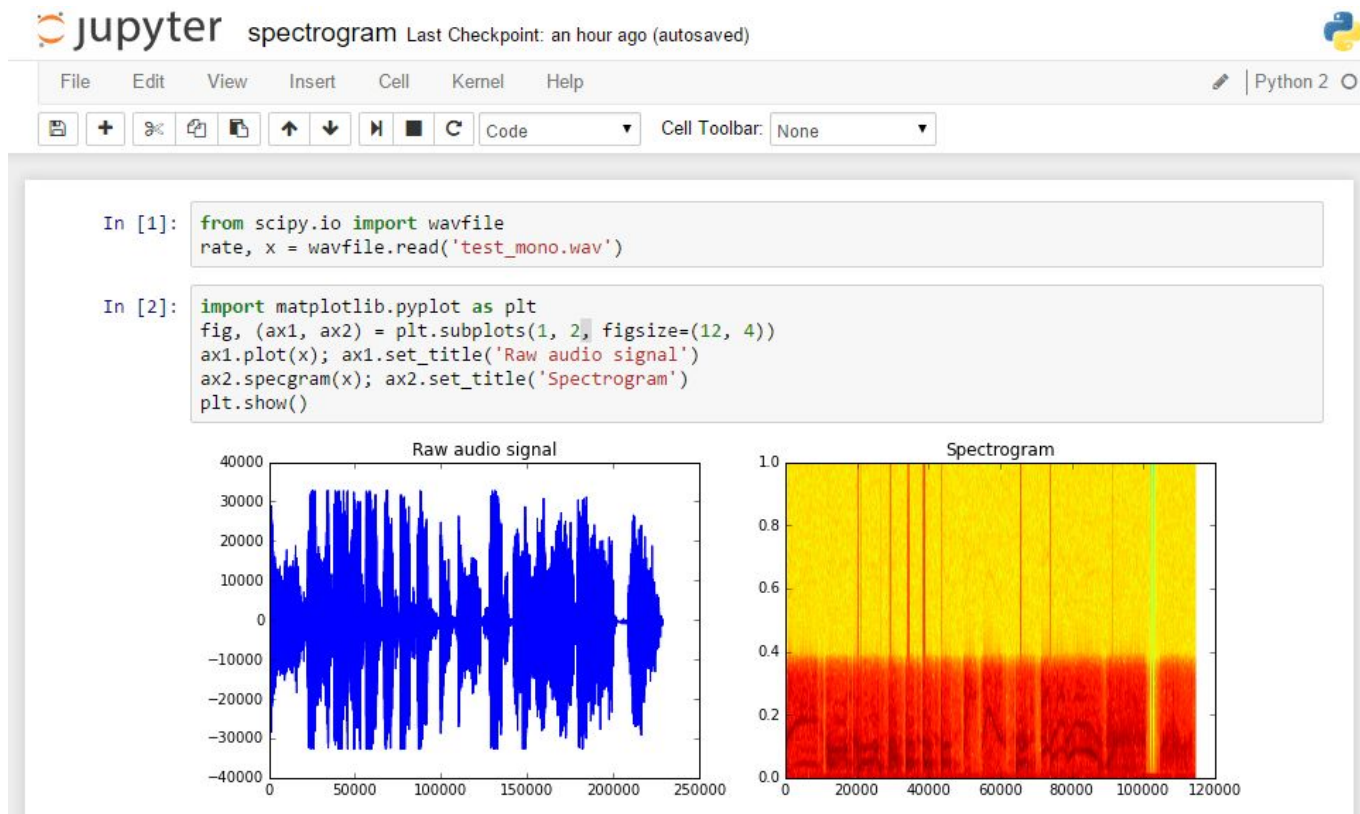
# Jupyter Notebooks

# Emulation as a Service Infrastructure

Interuniversity Distributed Network of
Emulation Nodes

Six partner institutions

Environment Contexts with full software
installation

Access to legacy files and objects



**Carnegie
Mellon
University**

# Questions/Comments

**Eric Kaltman**
*Data Curation Fellow*
ekaltman@andrew.cmu.edu

Carnegie
Mellon
University