



# **Implementing Metadata Standards**

EaaSI Training Module #7

# During This Module

- What kind of metadata does EaaSI (the program of work) gather?
- What metadata has been directly incorporated into the EaaSI platform?
- Which standards are influencing EaaSI system design?
- How does EaaSI share/contribute its metadata to the field?

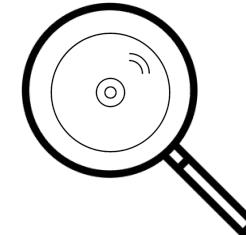




**But first...if someone asked, how would you describe a particular piece of software?**

Image source:  
[https://archive.org/download/Super\\_Cool\\_Shareware\\_PC-World\\_Digital\\_Concepts\\_1999/Super%20Cool%20Shareware%20PC-World%29%28Digital%20Concepts%29%281999%29.page](https://archive.org/download/Super_Cool_Shareware_PC-World_Digital_Concepts_1999/Super%20Cool%20Shareware%20PC-World%29%28Digital%20Concepts%29%281999%29.page)

# Categorizing Software



You could start by asking (and answering)...

- What is the software?
- Who is responsible for the software?
- How is the software distributed?
- What is the software made of?
- What is needed to run the software?

**NB:** The questions above are taken from the Software Preservation Network's [Software Metadata Recommended Format \(SMRF\) Guide!](#)

<b>Property</b>	Quick/Local Definition	<b>Example</b>
<b>Software Product Name</b>	Name of the software product as it appears on distribution mechanism or packaging	PrintMaster Gold Classic
<b>Software Product Version</b>	Version of the software product as it appears within system metadata (or on packaging if system metadata is unavailable)	4.03.43
<b>Distribution Mechanism</b>	Media carrier or file download by which the software is distributed to users	CD-ROM
<b>Installable Languages</b>	Installable languages for the Software Product (note: not for licenses) referenced within the software	English
<b>Executable Location</b>	Where the application is stored in the emulated computer's file system after installation. This must be a full path	C:\pmw\Pmw.exe
<b>Open File Format(s)</b>	File formats that can be opened by the software as represented within the software	Card Files (*.car); Banner Files (*.ban)

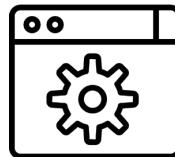
Here is **some** example metadata that the EaaSI program of work gathers about software at the moment

# Every Metadata Field Has a Story



Every field in the prior table was selected because of its value to one or more aspects of our program of work, such as:

- Feature development for EaaS Client



- Search and discovery



- Contributing to Wikidata



- Automated dependency tracking and/or recommendations

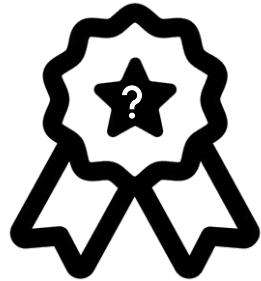


# EaaSI ≠ A Standard



- EaaSI's metadata-related decisions are currently driven by:
  - Functional requirements of emulation, EaaS and/or the EaaSI UI
  - Grant deliverables (reporting benchmarks)
  - Perceived value to global digital preservation efforts in gathering information about legacy software
- EaaSI neither uses a singular centralized authority for software or emulation-relevant metadata - ***nor is trying to be one***
- Given funding and intended community, have focused on the needs of cultural heritage and research data management

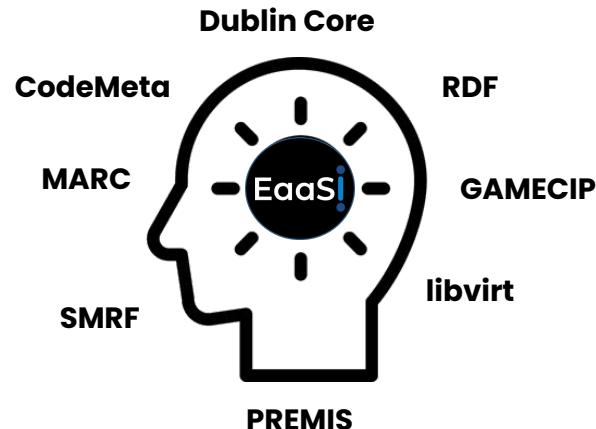
# **What Is a Standard, Then**



- Metadata standards require consensus - or at least common understanding
- Structured requirements that ensure data can be properly interpreted across systems or individual platforms (like the EaaSI stack)
- Evolving demands of digital preservation are currently directing community energy towards creating metadata standards related to software preservation

# Looking for Inspiration

- EaaSI team wants the stack to draw on existing bibliographic and preservation standards...
- ...but also make room for emerging efforts and expertise in standards for software preservation or related domains



# Mapping



- Where possible, EaaSI is making efforts to make crosswalks available between metadata captured by the team to existing standards
- Structural mapping now will make direct integration with repositories and catalogs easier down the road
- Occasionally difficult due to lack of emulation-specific properties in existing preservation and bibliographic standards
- EaaSI needs to describe software both in and of itself - a standalone "work" - and as part of a complex, emulated Environment

# Example Metadata Mapping

EaaSI Metadata Gathering	MARC	MODS	Dublin Core	SMRF	CodeMeta	Wikidata QID
UUID (assigned by EaaS)	<p>0241 - Other Standard Identifier, Universal Product Code  <a href="https://www.loc.gov/marc/bibliographic/bd024.html">https://www.loc.gov/marc/bibliographic/bd024.html</a></p> <p>0248 - Other Standard Identifier, Unspecified type of standard number or code  <a href="https://www.loc.gov/marc/bibliographic/bd024.html">https://www.loc.gov/marc/bibliographic/bd024.html</a></p>	<identifier> <a href="https://www.loc.gov/standards/mods/userguide/identifier.html">https://www.loc.gov/standards/mods/userguide/identifier.html</a>	Identifier <a href="http://purl.org/dc/terms/identifier">http://purl.org/dc/terms/identifier</a>	Identifier	identifier	unique identifier (Q6545185) <a href="https://www.wikidata.org/wiki/Q6545185">https://www.wikidata.org/wiki/Q6545185</a>
Software Product Name	<p>245 - Title Statement  <a href="https://www.loc.gov/marc/bibliographic/bd245.html">https://www.loc.gov/marc/bibliographic/bd245.html</a></p> <p>500 - General Note (to provide source of title)  <a href="https://www.loc.gov/marc/bibliographic/bd500.html">https://www.loc.gov/marc/bibliographic/bd500.html</a></p>	<titleInfo> <a href="https://www.loc.gov/standards/mods/userguide/titleinfo.html">https://www.loc.gov/standards/mods/userguide/titleinfo.html</a>	Title <a href="http://purl.org/dc/elements/1.1/title">http://purl.org/dc/elements/1.1/title</a>	Title		title (Q783521) <a href="https://www.wikidata.org/wiki/Q783521">https://www.wikidata.org/wiki/Q783521</a>
Software Product Version	<p>250 - Edition Statement  <a href="https://www.loc.gov/marc/bibliographic/bd250.html">https://www.loc.gov/marc/bibliographic/bd250.html</a></p> <p>251 - Version Information  <a href="https://www.loc.gov/marc/bibliographic/bd251.html">https://www.loc.gov/marc/bibliographic/bd251.html</a></p>	<originInfo> <edition> <a href="https://www.loc.gov/standards/mods/userguide/origininfo.html#edition">https://www.loc.gov/standards/mods/userguide/origininfo.html#edition</a>	Relation <a href="http://purl.org/dc/terms/relation">http://purl.org/dc/terms/relation</a>	Version	softwareVersion	software version (Q20826013) <a href="https://www.wikidata.org/wiki/Q20826013">https://www.wikidata.org/wiki/Q20826013</a>

# Example Metadata Mapping

<b>EaaSI Metadata Gathering</b>	<b>MARC</b>	<b>MODS</b>	<b>Dublin Core</b>	<b>SMRF</b>	<b>CodeMeta</b>	<b>Wikidata QID</b>
Date Published	<p>260 – Publication, Distribution, etc.  <a href="https://www.loc.gov/marc/bibliographic/bd260.html">https://www.loc.gov/marc/bibliographic/bd260.html</a></p> <p>542 – Information Relating to Copyright Status  <a href="https://www.loc.gov/marc/bibliographic/bd542.html">https://www.loc.gov/marc/bibliographic/bd542.html</a></p> <p>388 – Time period of Creation  <a href="https://www.loc.gov/marc/bibliographic/bd388.html">https://www.loc.gov/marc/bibliographic/bd388.html</a></p>	<originInfo> <dateCreated> <a href="https://www.loc.gov/standards/mods/userguide/origininfo.html">https://www.loc.gov/standards/mods/userguide/origininfo.html</a>	Date <a href="http://purl.org/dc/elements/1.1/date">http://purl.org/dc/elements/1.1/date</a>  Date Created <a href="http://purl.org/dc/terms/created">http://purl.org/dc/terms/created</a>	Date, Date Subtype	datePublished	copyright date (Q59584702) <a href="https://www.wikidata.org/wiki/Q59584702">https://www.wikidata.org/wiki/Q59584702</a>  date of publication (Q1361758) <a href="https://www.wikidata.org/wiki/Q1361758">https://www.wikidata.org/wiki/Q1361758</a>
Required Software (OS)	753 \$c – System Details Access to Computer Files; Operating System <a href="https://www.loc.gov/marc/bibliographic/bd753.html">https://www.loc.gov/marc/bibliographic/bd753.html</a>	<note> <a href="https://www.loc.gov/standards/mods/userguide/note.html">https://www.loc.gov/standards/mods/userguide/note.html</a>	Relation <a href="http://purl.org/dc/elements/1.1/relation">http://purl.org/dc/elements/1.1/relation</a>	Operating System	operatingSystem	operating system (Q9135) <a href="https://www.wikidata.org/wiki/Q9135">https://www.wikidata.org/wiki/Q9135</a>
Required Software (Other)	538 – System Details Note <a href="https://www.loc.gov/marc/bibliographic/bd538.html">https://www.loc.gov/marc/bibliographic/bd538.html</a>	<note> <a href="https://www.loc.gov/standards/mods/userguide/note.html">https://www.loc.gov/standards/mods/userguide/note.html</a>	Relation <a href="http://purl.org/dc/elements/1.1/relation">http://purl.org/dc/elements/1.1/relation</a>	Additional Dependencies		system requirements (Q2275513) <a href="https://www.wikidata.org/wiki/Q2275513">https://www.wikidata.org/wiki/Q2275513</a>

# Example Metadata Mapping

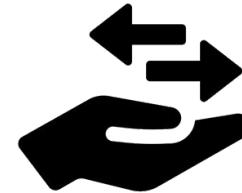
<b>EaaSI Metadata Gathering</b>	<b>MARC</b>	<b>MODS</b>	<b>Dublin Core</b>	<b>SMRF</b>	<b>CodeMeta</b>	<b>Wikidata QID</b>
Selected Base Environment	<p>538 - System Details Note  <a href="https://www.loc.gov/marc/bibliographic/bd538.html">https://www.loc.gov/marc/bibliographic/bd538.html</a></p> <p>753 - System Details Access to Computer Files  <a href="https://www.loc.gov/marc/bibliographic/bd753.html">https://www.loc.gov/marc/bibliographic/bd753.html</a></p>	<note> <a href="https://www.loc.gov/standards/mods/userguide/note.html">https://www.loc.gov/standards/mods/userguide/note.html</a>	Relation <a href="http://purl.org/dc/elements/1.1/relation">http://purl.org/dc/elements/1.1/relation</a>		softwareRequirements	software dependency (Q56859575) <a href="https://www.wikidata.org/wiki/Q56859575">https://www.wikidata.org/wiki/Q56859575</a>
Default Save File Format	<p>256 - Computer File Characteristics  <a href="https://www.loc.gov/marc/bibliographic/bd256.html">https://www.loc.gov/marc/bibliographic/bd256.html</a></p> <p>347 - Digital File Characteristics  <a href="https://www.loc.gov/marc/bibliographic/bd347.html">https://www.loc.gov/marc/bibliographic/bd347.html</a></p> <p>516 - Type of Computer File or Data Note  <a href="https://www.loc.gov/marc/bibliographic/bd516.html">https://www.loc.gov/marc/bibliographic/bd516.html</a></p>	<physicalDescription><form> <a href="https://www.loc.gov/standards/mods/userguide/physicaldescription.html#form">https://www.loc.gov/standards/mods/userguide/physicaldescription.html#form</a>  <physicalDescription><internetMediaType> <a href="https://www.loc.gov/standards/mods/userguide/physicaldescription.html#internetmediatype">https://www.loc.gov/standards/mods/userguide/physicaldescription.html#internetmediatype</a>  <physicalDescription><digitalOrigin> <a href="https://www.loc.gov/standards/mods/userguide/physicaldescription.html#digitalorigin">https://www.loc.gov/standards/mods/userguide/physicaldescription.html#digitalorigin</a>	Format <a href="http://purl.org/dc/elements/1.1/format">http://purl.org/dc/elements/1.1/format</a>	File Format	fileFormat	format (Q2085518) <a href="https://www.wikidata.org/wiki/Q2085518">https://www.wikidata.org/wiki/Q2085518</a>
Open File Format(s)	<p>256 - Computer File Characteristics  <a href="https://www.loc.gov/marc/bibliographic/bd256.html">https://www.loc.gov/marc/bibliographic/bd256.html</a></p> <p>347 - Digital File Characteristics  <a href="https://www.loc.gov/marc/bibliographic/bd347.html">https://www.loc.gov/marc/bibliographic/bd347.html</a></p> <p>516 - Type of Computer File or Data Note  <a href="https://www.loc.gov/marc/bibliographic/bd516.html">https://www.loc.gov/marc/bibliographic/bd516.html</a></p>	<physicalDescription><form> <a href="https://www.loc.gov/standards/mods/userguide/physicaldescription.html#form">https://www.loc.gov/standards/mods/userguide/physicaldescription.html#form</a>  <physicalDescription><internetMediaType> <a href="https://www.loc.gov/standards/mods/userguide/physicaldescription.html#internetmediatype">https://www.loc.gov/standards/mods/userguide/physicaldescription.html#internetmediatype</a>  <physicalDescription><digitalOrigin> <a href="https://www.loc.gov/standards/mods/userguide/physicaldescription.html#digitalorigin">https://www.loc.gov/standards/mods/userguide/physicaldescription.html#digitalorigin</a>	Format <a href="http://purl.org/dc/elements/1.1/format">http://purl.org/dc/elements/1.1/format</a>	File Format	fileFormat	format (Q2085518) <a href="https://www.wikidata.org/wiki/Q2085518">https://www.wikidata.org/wiki/Q2085518</a>

# Metadata is Labor



- Bulk of EaaSI-related metadata is currently gathered by a team of student employees at Yale University Library:  
<https://www.softwarepreservationnetwork.org/emulation-as-a-service-infrastructure/#students>
- Knowing what software metadata is – and how to find it – takes training and experience
- Automation efforts promise future improvements to efficiency, scale, and more time dedicated to creating new Environments rather than describing them

# Metadata Exchange



- Currently Environment and Software resource records can be passed between EaaSI nodes using OAI-PMH (see Module #6)
  - This *does not include* – yet! – much of the descriptive and structural metadata properties about software gathered by the student team
- Experimenting with serialization of expanded Environment and Software metadata via JSON and XML
- Relevant properties gathered by students are sent via spreadsheet to WikiDP

# Wikidata and WikiDP



- Wikidata for Digital Preservation (WikiDP) offers a custom portal for browsing and contributing to the Wikidata knowledge base in categories related to computing and software history
- Partnered with EaaSI to encourage conversion of our metadata into structured, linked open data
- Wikidata contributions are broadly reviewed and consensus-driven, so contributed metadata is purposefully not specific to the EaaSI stack
- Requires some additional evidence to bolster "notability" of new properties (i.e. file formats, software versions, and computing environments need to be formally differentiated)

# Software Preservation Network



- Affiliating with a relevant professional organization taps EaaSI into emerging software preservation process and standards
- Helps EaaSI to advocate for broader adoption and collection of metadata necessary for emulation
- Ensures collective, collaborative feedback on search behavior, future integrations, API design, mapping, etc.



# A Journey in Software Metadata via Emulation



*In this module we have been discussing software metadata for description, automation, discovery, etc.*

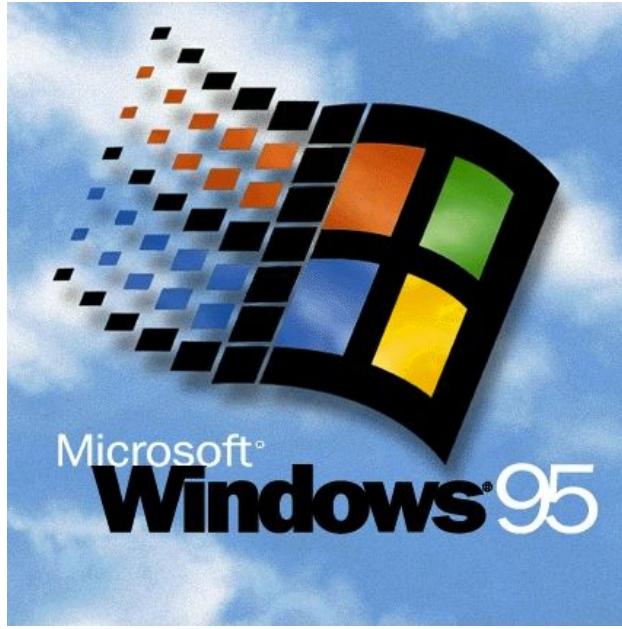
*How does all this actually come into play with emulation and the EaaSI platform? Let's walk through the process of emulating a file from scratch - charting \*some\* of the metadata (and relationships between pieces of metadata) necessary along the way.*

```
ethan@ethan-OptiPlex-5090:~/Downloads$ sf -sig wikidata.sig bobsled.pcx
...
siegfried    : 1.9.4
scandate     : 2022-08-03T16:03:44-04:00
signature    : wikidata.sig
created      : 2022-07-17T20:41:40+02:00
identifiers  :
  - name    : 'archivematica'
    details : 'wikidata-definitions-3.0.0 (2022-07-17); extensions: archivematica-fmt2.xml, archivematica-fmt3.xml, archivematica-fmt4.xml, archivematica-fmt5.xml'
...
filename : 'bobsled.pcx'
filesize : 71486
modified : 2022-08-03T15:44:17-04:00
errors   :
matches   :
  - ns      : 'archivematica'
    id      : 'Q43869672'
    format  : 'PCX, version 3'
    URI     : 'http://www.wikidata.org/entity/Q43869672'
    permalink : 'https://www.wikidata.org/w/index.php?oldid=1423327110&title=Q43869672'
    mime    : 'image/x-pcx; image/vnd.zbrush.pcx'
    basis   : 'extension match pcx; byte match at 0, 2 (Wikidata reference is empty)'
    warning :
```

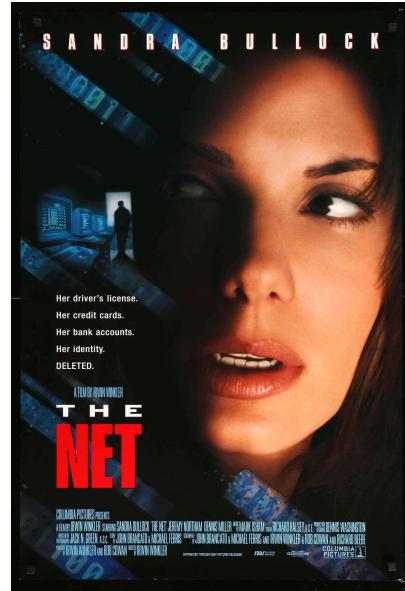
Step 1: I have a file - “**bobsled.pcx**” - that won’t open in contemporary software. Context clues - the extension name, file characterization - leads me to believe this was an image file made for/by a PC, probably sometime in the early to mid-1990s. We can guess it depicts a bobsled, but we can only know for sure by properly rendering it.



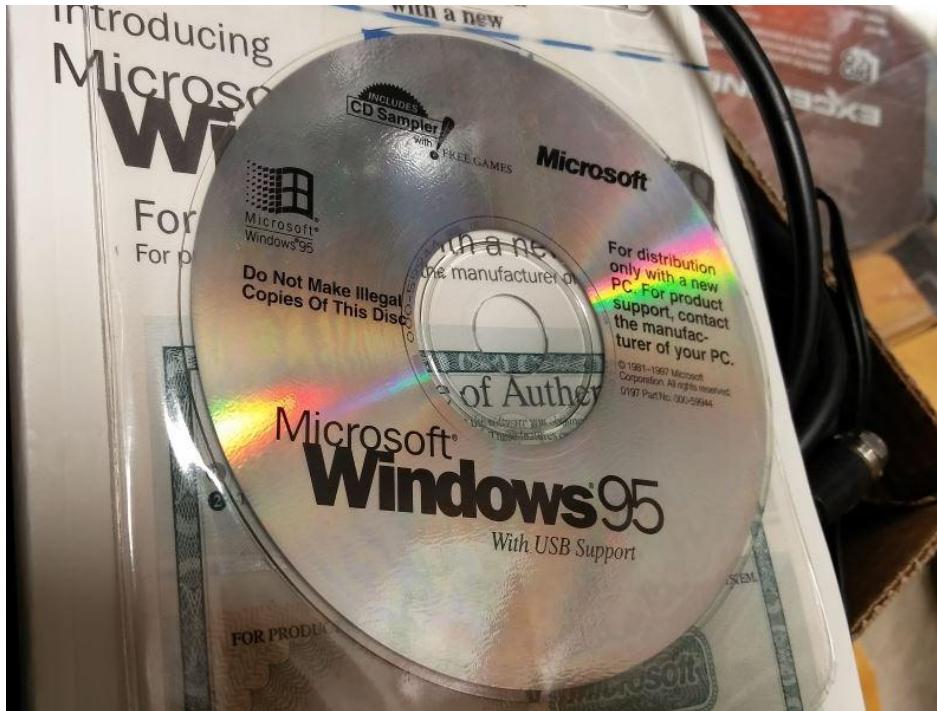
Step 2: I must choose a PC **operating system** to emulate: **Windows**



Step 3: I must also choose a specific **version** to emulate: Windows **95**



Step 4: Windows 95 was released in a particular year: **1995**



Step 5: Windows 95 was distributed using certain ***installation media***:  
**CD-ROM**



Step 6: With a selected operating system version and year we can also assume the type of **processor** we need to emulate: **32-bit x86** (also known as **i386**)



Step 7: This processor requires a particular **emulator**: **QEMU** (specifically **qemu-system-i386**)

ID: 2f57f2f9-81ae-4b97-9b3a-6759372b3b2d

Name: Windows 95 C (OSR 2.5) - Base V2

Handle: --- [create](#)

Description 

Emulation system settings 

Date "2022-08-03T20:32:52.542Z"

Operating System:

Emulator: Qemu

Emulator version: emucon-rootfs/qemu-system|v2.12

Emulator configuration: -m 64 -vga cirrus -soundhw sb16 -net nic,model=pcnet

Linux Runtime:

Step 8: By necessity we must also pick a particular **version** of the emulator as well: **QEMU 2.12**

ID: 2f57f2f9-81ae-4b97-9b3a-6759372b3b2d

Name: Windows 95 C (OSR 2.5) - Base V2

Handle: --- [create](#)

The screenshot shows the QEMU configuration interface. At the top, it displays the ID, Name, and Handle of the configuration. Below this, there are two main sections: "Emulation system settings" and "Configured Drives". The "Emulation system settings" section includes fields for Date, Operating System, Emulator (set to Qemu), Emulator version (set to emucon-roots/qemu-systemv2.12), and Emulator configuration (set to -m 64 -vga cirrus -soundhw sb16 -net nic,model=pcnet). The "Configured Drives" section lists disk, floppy, and cdrom drives, each with checkboxes and an 'X' button. A red oval highlights the Emulator configuration field, and a larger red oval highlights the list of configured drives.

Description >

Emulation system settings ▾

Date "2022-08-03T20:32:52.542Z"

Operating System:

Emulator: Qemu

Emulator version: emucon-roots/qemu-systemv2.12

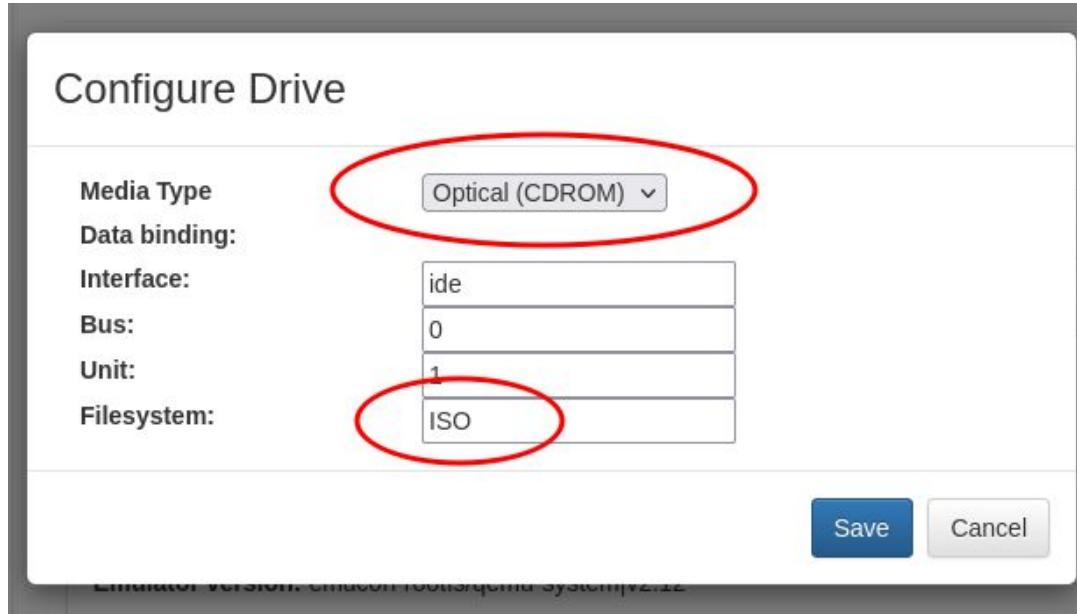
Emulator configuration: -m 64 -vga cirrus -soundhw sb16 -net nic,model=pcnet

Linux Runtime:

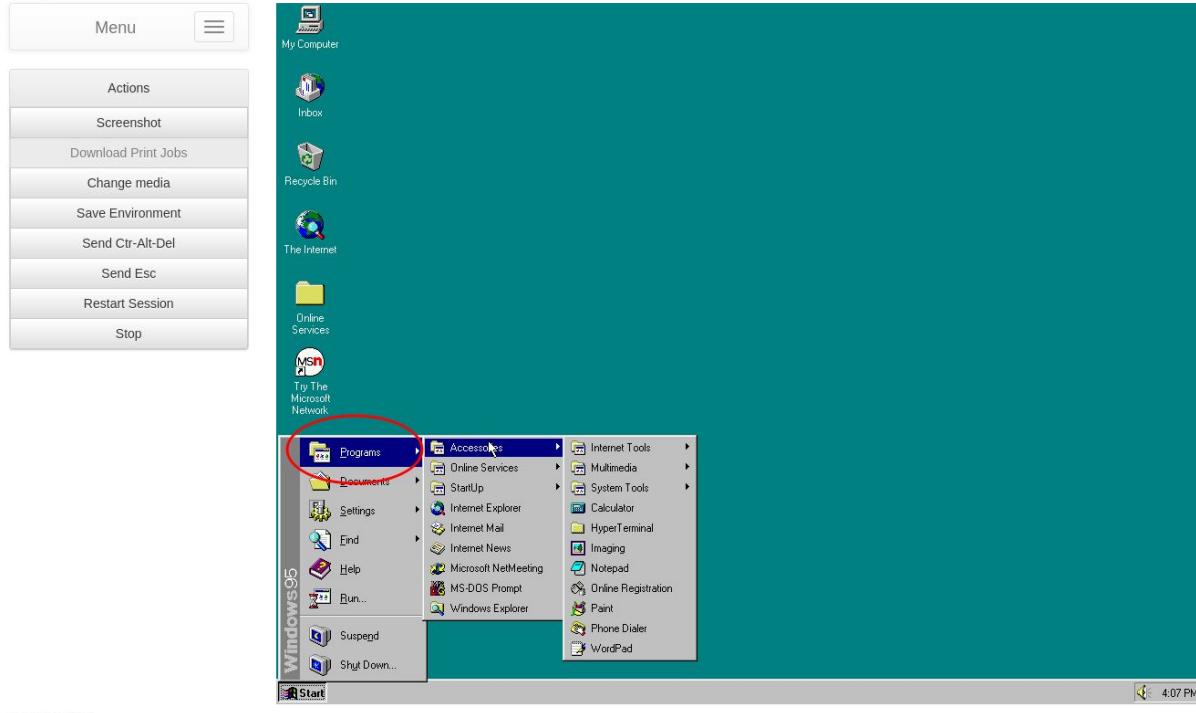
Configured Drives ▾

	Add Drive
disk	<input checked="" type="checkbox"/> <span>X</span>
floppy	<input checked="" type="checkbox"/> <span>X</span>
cdrom	<input checked="" type="checkbox"/> <span>X</span>

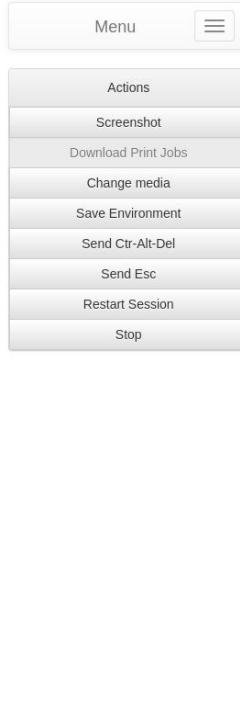
Step 9: This emulator version can recreate a set of **hardware components**: a **generic PS/2 mouse and keyboard**, **64 MB of RAM**, **Cirrus CLGD 5446 VGA card**, **Creative Soundblaster 16 sound card**, **AMD PCNET ethernet card**, **2 PCI IDE interfaces** with **hard disk**, **floppy disk** and **CD-ROM drive support**



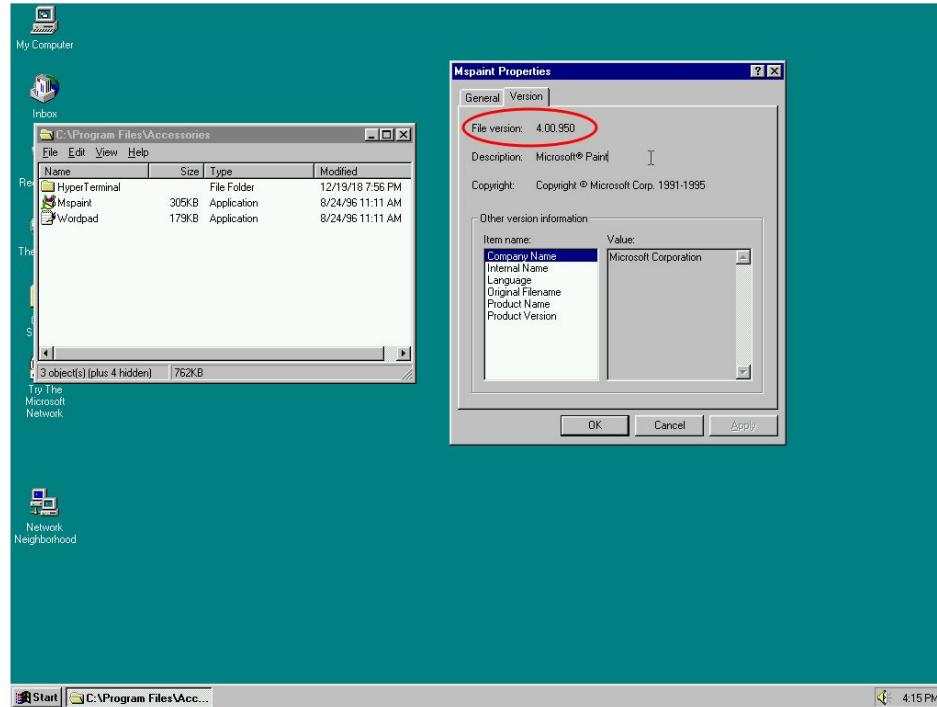
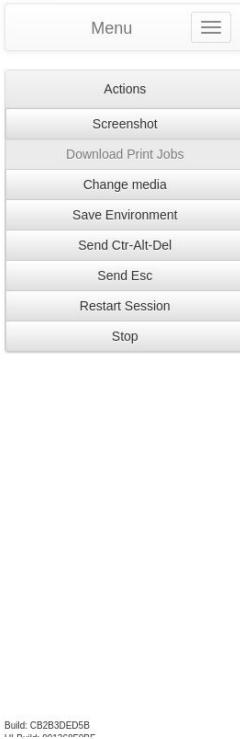
Step 10: QEMU's emulated CD-ROM drive can read certain **types** of **installation media** - like a disk image made from the original CD-ROM (such disk images are commonly referred to as **ISOs**)



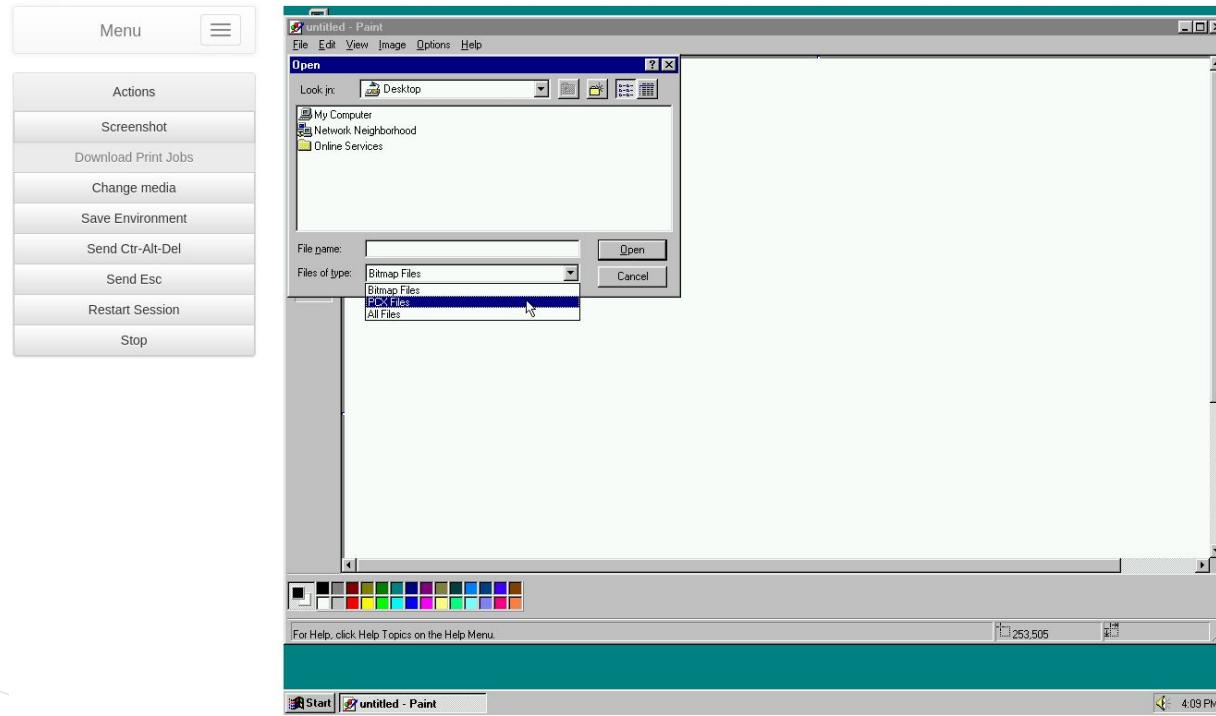
Step 11: The Windows 95 ISO installs a list of **software products** that make up the operating system - **Wordpad, Paint, Calculator, File Manager, Internet Explorer**, etc.



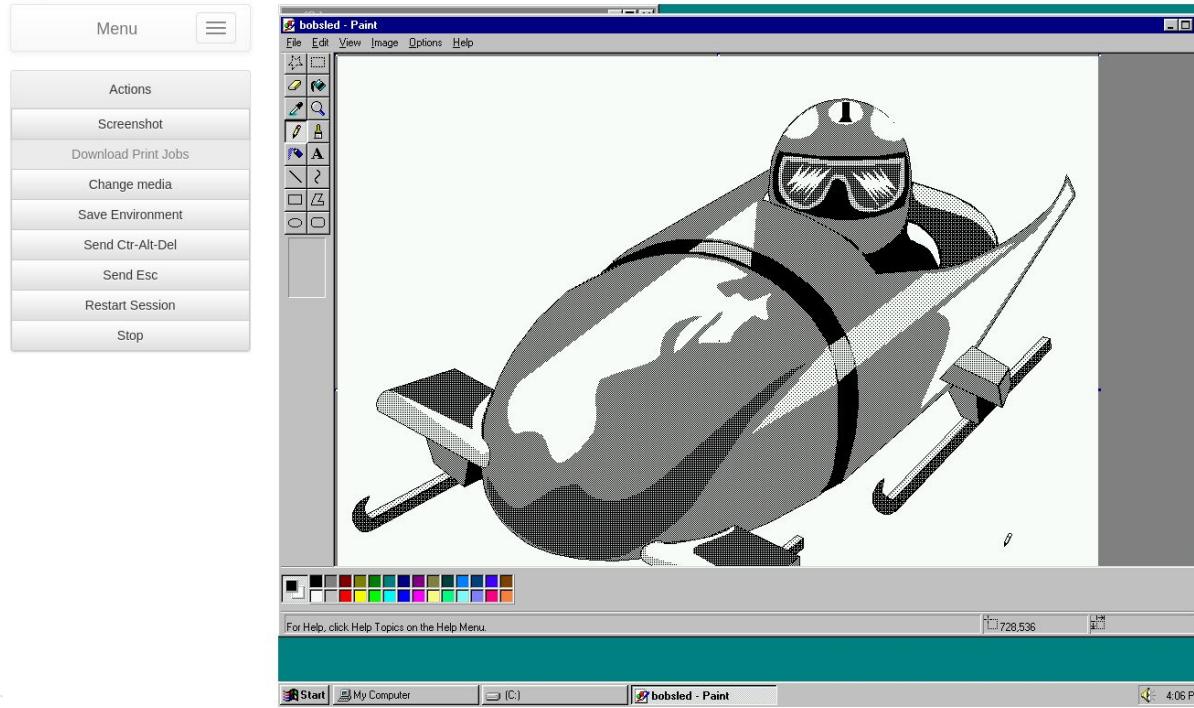
Step 12: Microsoft Paint is opened by running an **executable** at a certain **location** - **C:\Program Files\Accessories\Mspaint.exe**



Step 13: Microsoft Paint has its own **version** - Microsoft Paint **4.00.950**



Step 14: This particular version of this particular software can **open** certain **file formats** - including **.PCX** files (like the one we started with!)



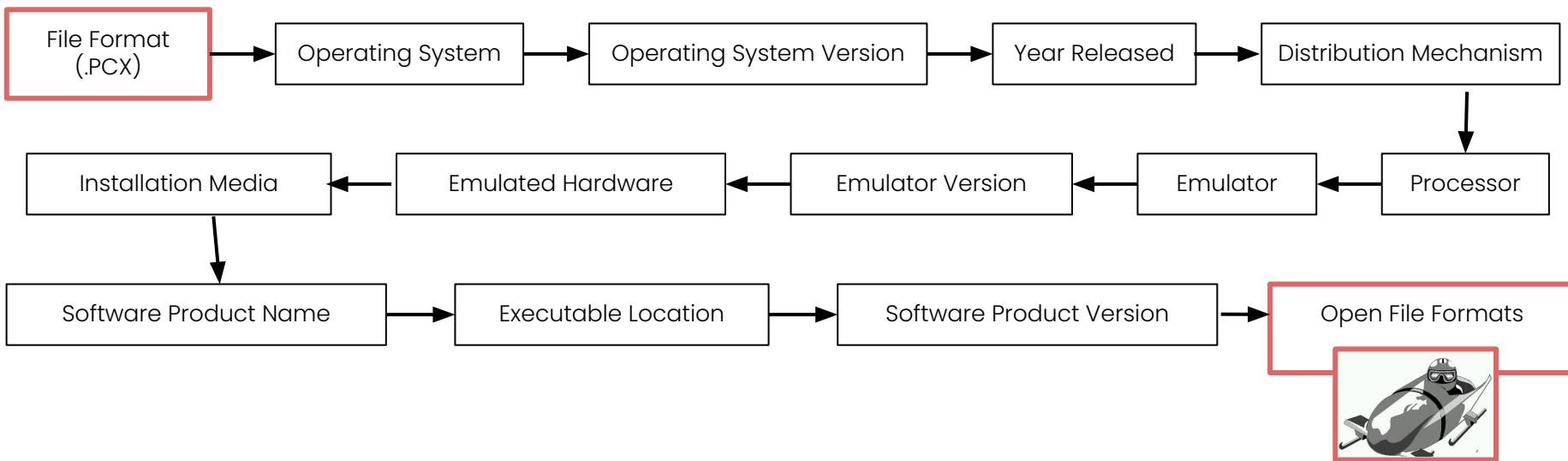
Step 15: Finally, thanks to our emulation of Windows 95 and Microsoft Paint 4.00.950, I can now open my original .PCX file and confirm - it's an **image of a bobsled!**

Build: CB2B3DED5B  
UI-Build: 991368F0BF

# Thanks, Metadata!



- Various pieces of software metadata allowed us to emulate and render this file as accurately as possible
- Pushing standards and investigating automation will help us speed the journey



# Credits

- Training Module written and designed by Ethan Gates, Software Preservation Analyst, Yale University Library & Claire Fox, Metadata Analyst, Yale University Library
- Original photos, screenshots, and videos recorded by Ethan Gates
- Icons sourced from [The Noun Project](#)
- EaaSI program of work sponsored by the Alfred P. Sloan Foundation and the Mellon Foundation, hosted by Yale University Library



**Yale**  
*Principle Partner*



**ALFRED P. SLOAN  
FOUNDATION**

*Sponsor*

 **Mellon  
Foundation**

*Sponsor*