# SOFTWARE PROVENANCE: TRACK THE REALITY NOT THE VIRTUAL MACHINE



HTTPS://OCCAM.CS.PITT.EDU

**David Wilkinson**, Luís Oliveira, Daniel Mossé, Bruce Childers
Computer Science Department – University of Pittsburgh
{dwilk, loliveira}@cs.pitt.edu

# Motivation

Preservation crisis

- Substantial amount of new data created everyday
- Both scientific and creative data/work

How to best preserve future **interactivity**?

- What can we do to allow data to not only be *read* but understood or manipulated later?
- Particularly if we discover defects very late!

# Possible Solutions Abound

Preserve older hardware to run older software
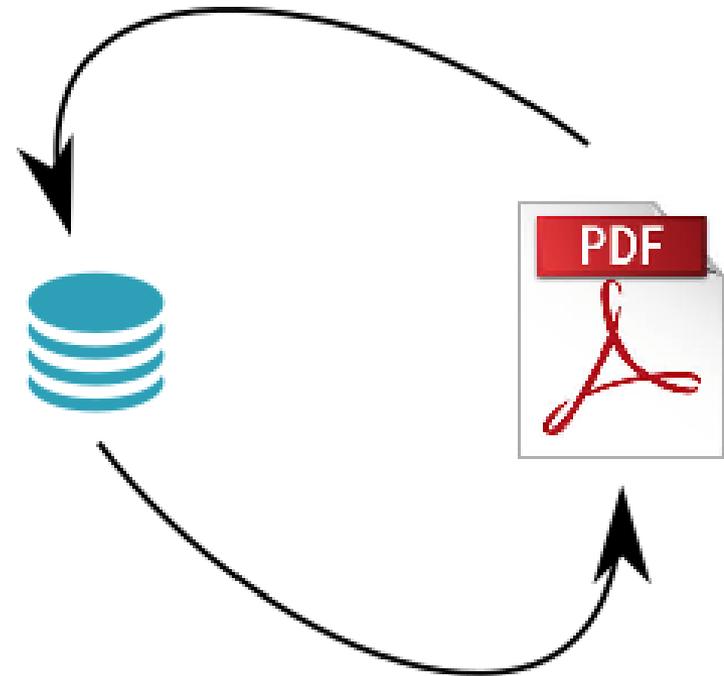
• but how practical is this? How expensive/accessible?

The Byte Cellar | http://www.bytecellar.com/

# Possible Solutions Abound

Transform files into more stable, more modern formats.

- Old Word -> XML, PDF
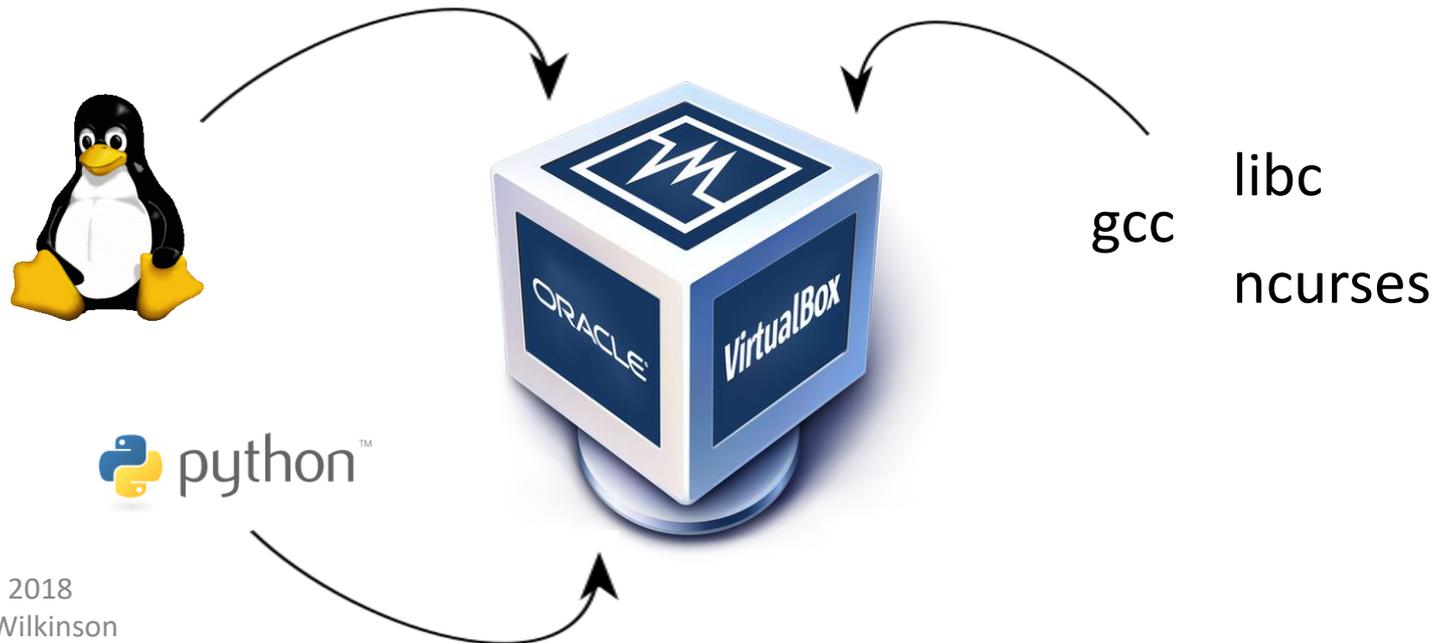
- Old Images -> PNG, JPG

# Transformation Fatigue

The difficulty is within the presumption that the resulting file is actually both better-preserved and complete.

- PDF is not necessarily a stable standard.

- (Do we commit to a subset of PDF? Who decides this?)

- Nothing prevents PDF, PNG, etc from itself becoming obsolete

- Do we convert again in the future?

# The Scientific Method

Datasets, plots, papers are all possible artifacts that may also have these preservation concerns.

Current acceptable approach: Distribute a VM image.

libc

gcc

ncurses

# The Power/Weaknesses of Virtual Machines

Virtual Machines are useful because they are powerful **abstractions**.

Yet, that also means much is hidden or taken as faith.



a human being

hardware

operating system

virtual machine manager

virtual machine

operating system

program

# Generating Virtual Machines

Occam takes the following approach:

Instead of distributing virtual machine images (OVA/OVF), distribute descriptions of virtual machines.

Maintain modular/composable archives of software, much like package managers.

# Example: DOS Games

We want to run some old DOS game: Doom

We would describe Doom as an object that needs a "dos" environment.

The system would prepare a virtual machine for the host machine that can execute a dos environment.

# Example: DOS Games

```
{
  "name": "Doom",  "type": "game",
  "install": [ {
    "type": "resource", "subtype": "application/zip",
    "source": "ftp://ftp.idsoftware.com/doom/doom19s.zip",
    "actions": { "unpack": "." }
  } ],
  "environment": "dos", "architecture": "x86",
  "run": { "command": "DOOM.EXE" }
}
```

# Example: DOS Games

To build the VM, we need to map a "dos" environment to a native "linux" environment.

In Occam, another artifact can describe such a relationship. Such as DOSBox, a DOS emulator for Linux.



linux/x86-64          dos/x86          Commander Keen

EMULATION

# Example: DOS Games

```
{
  "name": "DOSBox", "type": "emulator",
  "provides": [ {
    "environment": "dos", "architecture": "x86"
  } ],
  "dependencies": [ { "name": "x11"}, … ],
  "environment": "linux", "architecture": "x86-64",
  "build": { … },
  "run": { "command": "/usr/bin/dosbox" }
}
```

# Example: DOS Games

Occam then can create a virtual machine or execute a docker container with the emulator (DOSBox), any libraries it has as dependencies, and the DOS game.

This is the **manifest**, which is the main distributable (kilobytes), along with individual software artifacts (megabytes) instead of a complete VM image (gigabytes).

It then executes the virtual machine by running DOSBox, which is then responsible for running the game or application.

# Example: DOS Games

# Software Provenance for Emulation

Through this, we can emulate obsolete software while retaining easy access to how the emulator was built.

# Example: DOS Applications

If we can handle a game, we can handle a word processing application!

P-RECS 2018
David Wilkinson

# Example: Jupyter Notebook

Preserving Jupyter notebooks has been a challenge in the community solved mostly by well-defined common distributions (Anaconda)

Yet, can we provide a more general solution that preserves run-time/build-time provenance?

# Example: Jupyter Notebook

# Example: Jupyter Notebook

```
{
  "name": "jupyter-notebook", "type": "application",
  "views": [ {
    "type": "file", "subtype": "extension/ipynb"
  } ],
  "dependencies": [ { "name": "python"}, ... ],
  "environment": "linux", "architecture": "x86-64",
  "build": { ... },
  "network": { "bind": "8888" },
  "run": { "command": "run.sh" }
}
```
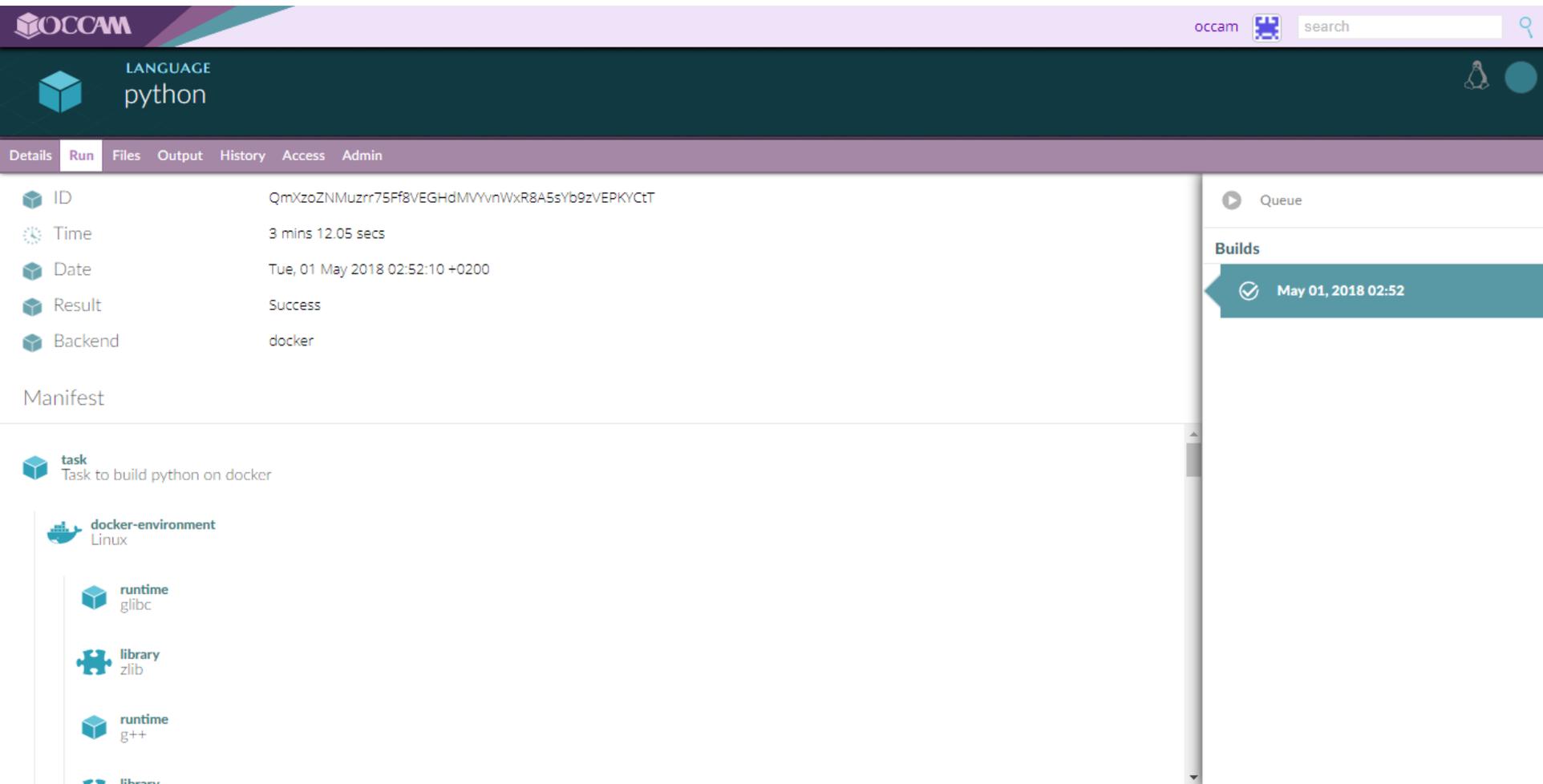
# Example: Jupyter Notebook

# Example: Jupyter Notebook Runtime Provenance via Manifests:

# Example: Jupyter Notebook
# The Build Provenance of Python 3

# Conclusions

In order to provide **interactivity** in the future, we must preserve software (the mechanism of interactivity)

Virtual Machines are great tools, but they are not interested in detailed provenance.

The greatest preservation of both provenance and interactivity is through the generation of virtual machines from manifests.

# Final thoughts

**Need feedback** from the community:

{dwilk,loliveira}@cs.pitt.edu

Current Occam implementation:

- Preserve as much as possible
  - Prevent silent loss of fidelity
  - Improve the longevity of software
- Preserving source code is vital
  - And the ability to build/run it.
- Dependencies are important
  - They may be the source of errors