Te Puna Mātauranga o Aotearoa
**NATIONAL LIBRARY**
OF NEW ZEALAND

KB ⟩ national library of the netherlands

# IIPC GA 2019 - WCT Workshop

Hands-on Setup Guide

Welcome to the Hands-on Setup Guide for the 2019 IIPC Web Curator Tool Workshop.

*The reality of any hands-on workshop is that things will break. We've tried our best to provide a stable virtual machine that can let you walk through the basics of using the Web Curator Tool.*

If you have any questions, please get in touch via our Slack group (*workshop-iipc-2019 channel*) or have a read of the WCT documentation.

# Setup

## Using our pre-configured VM

The WCT team have prepared a pre-configured virtual machine (VM), that has working WCT, OpenWayback and Heritrix 3 instances.

***To reduce the setup time during the Hands-on session, it is recommended that you download and install Oracle VirtualBox prior to the workshop.***
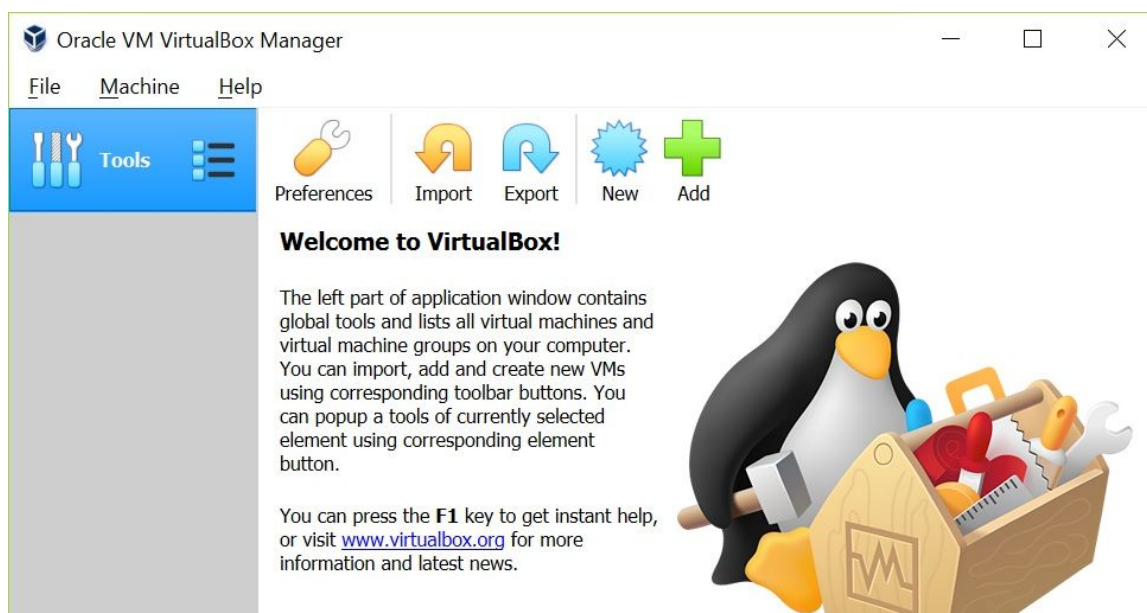
### VirtualBox

This section requires that you install the latest version of Oracle VirtualBox to run the provided VM (.ova file). VirtualBox is available for Windows, OSX and Linux. The version for your operating system can be downloaded here, and installation instructions found here. This VM has been tested on version **5.1** and later of VirtualBox.
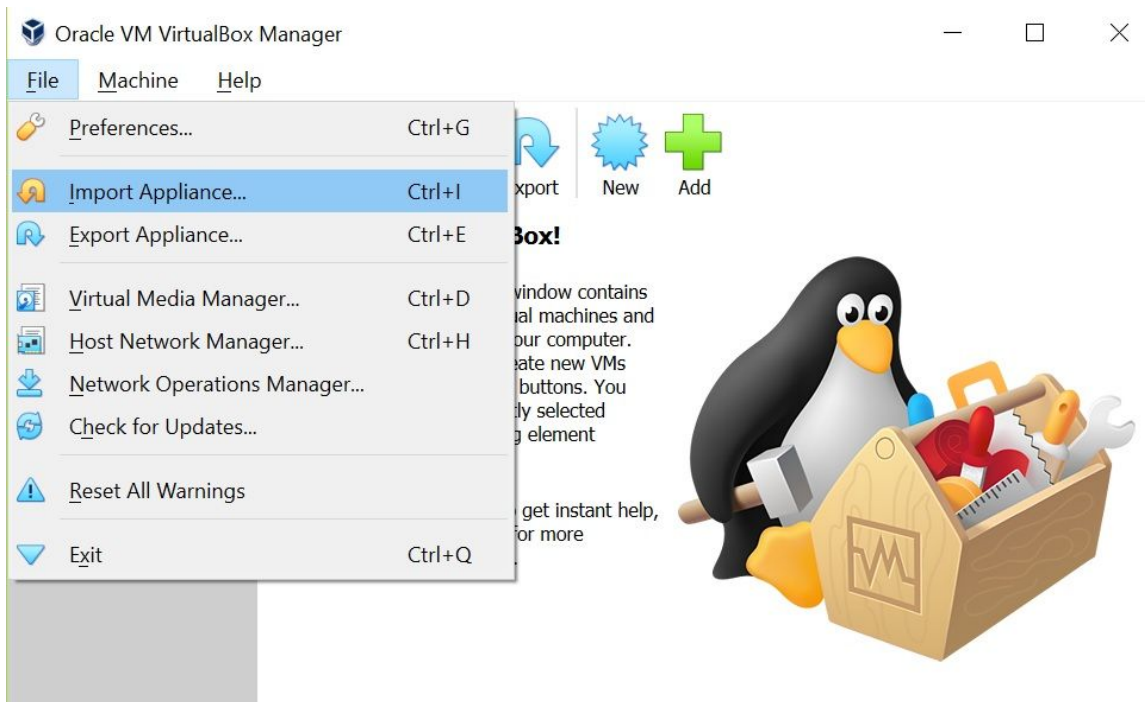
### Downloading the VM

A copy of the VM will be provided on a USB stick during the workshop. *Should you wish to download and familiarise yourself with the VM beforehand, a copy can be found here.*
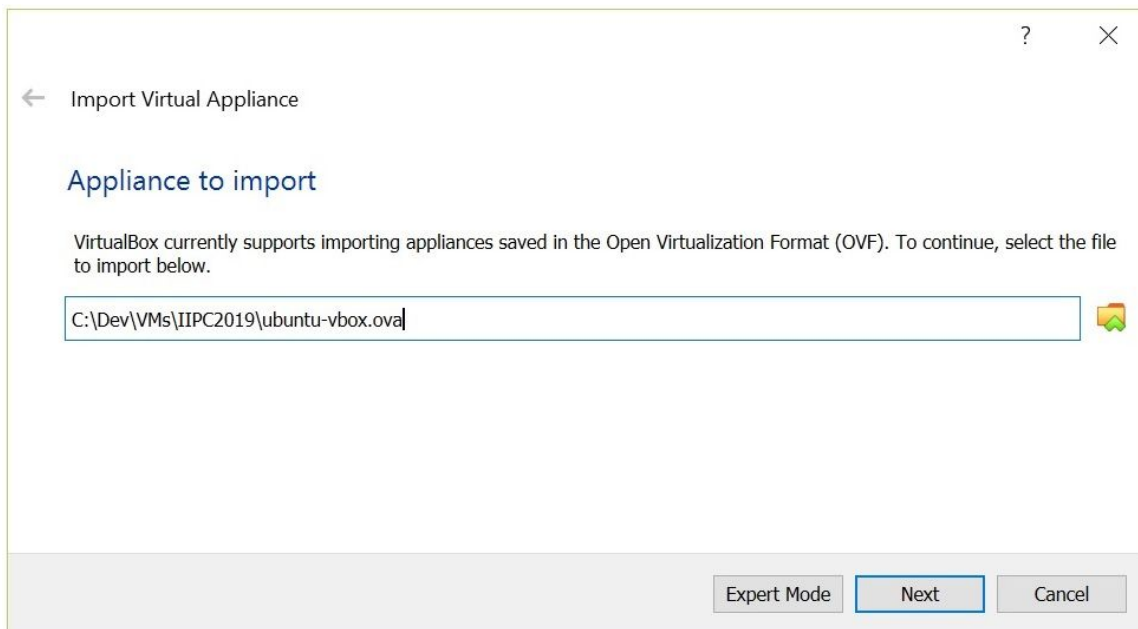
### Opening the VM

Once you have downloaded/copied the VM to an appropriate location on your computer, you can open it using VirtualBox.

Import the VM using File->Import Appliance



Select the VM (.ova file) to import

Specify a name, and the directory to store the imported VM

Import the .ova file, this can take some time



Successfully imported VM

If prompted, allow Firewall access for Virtual Box



Start the VM in headless mode

## Accessing WCT

WCT can be accessed using a browser running on your regular operating system (i.e. outside of the VM). The location is http://localhost:3080/wct/.
Username: wct
Password: Zagreb19

## Logging into the VM

You can log into the VM through your command line, using ssh: **ssh -p 3022 wct@localhost**.
Username: wct
Password: wct
This is also the password for sudo access.


Or using Putty

# Where things are

## WCT

The WCT directories are all under **/opt/wct**.
- The harvest agent stores its files in **/opt/wct/harvest-agent**.
- The store uses **/opt/wct/store**.
- All components log to files in **/opt/wct/logs**.
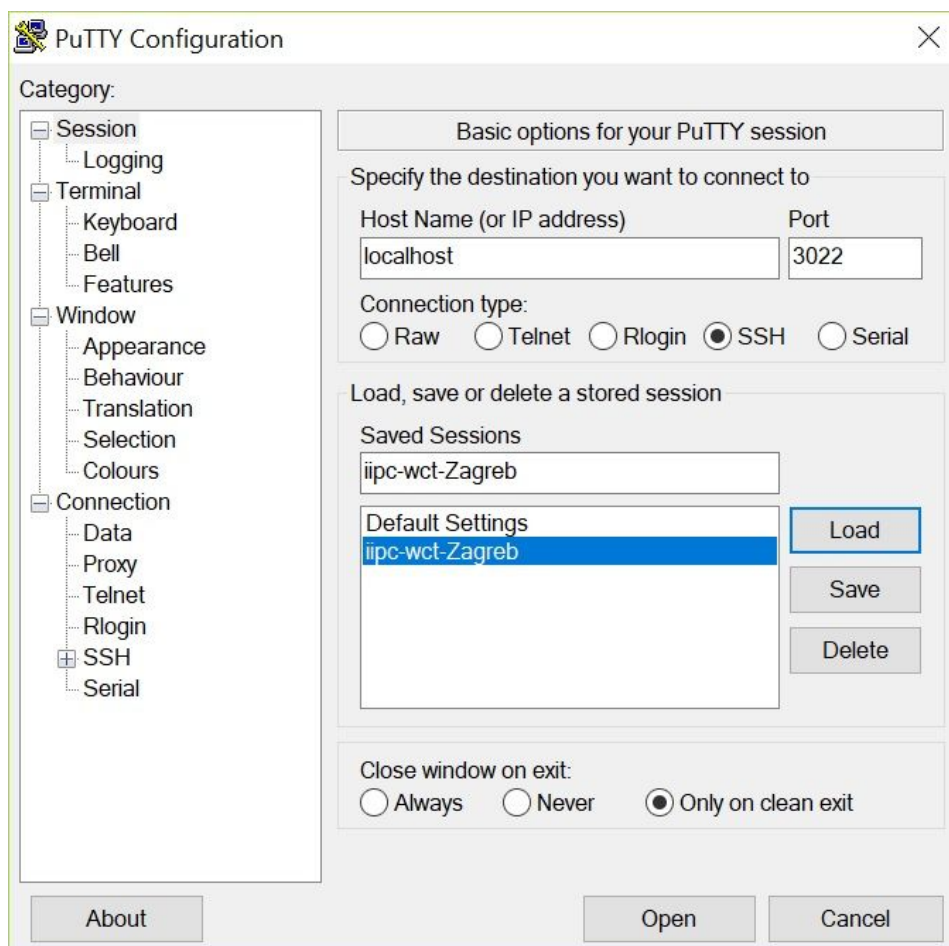- [Heritrix scripts](#) for use within the WCT script console can be placed in **/opt/wct/h3scripts**.

## Apache Tomcat

The Tomcat directory is **/opt/wct/apache-tomcat-9.0.19**. Tomcat is started automatically as a systemd service (called 'tomcat') when the system boots.

## MYSQL

You can log into MariaDB using **mysql -u root -p**. The password is "password". A basic WCT database has been provided, containing everything needed to start configuring and running crawls.

## Heritrix 3

Heritrix can be found in **/opt/wct/heritrix-3.4.0-20190418**. Like Tomcat, Heritrix is also started as a systemd service (called 'heritrix'). You can access the web interface at [https://localhost:3443/engine](https://localhost:3443/engine).
Username: admin
Pasword: admin

## OpenWayback

In this setup WCT uses OpenWayback to display harvest results. This OpenWayback has been deployed in the same Tomcat as WCT itself. It expects its WARC files in **/opt/wct/wayback**. It can be accessed at [http://localhost:3080/owb/wayback/](http://localhost:3080/owb/wayback/).

## Digital repository

This WCT has been configured to use a filesystem-based digital repository, which amounts to just a directory to copy WARC files to. This directory is **/opt/wct/warc-repository**.

## Crawling

The Heritrix 3 instance is configured with the following user agent string:
**Mozilla/5.0 (compatible; heritrix/3.4
+https://docs.google.com/document/d/1SZq2EV5Q0PxCMiU0hcII158gita9LU5CqJ5xKBUtAr0)**

This URL will be valid until the end of this workshop. If you wish to continue using this VM after the workshop, please update the user agent string within the WCT profiles to a URL relevant to your institution.

# Installing your own WCT

Should you wish to install WCT yourself, in your own environment, there are several guides available:
- [a quick start guide to installing WCT](#).
- [a comprehensive System Administrator guide for installing WCT](#).

# Resources

[Documentation](#)
[Slack](#)
[Github](#)