

# Krótkie sprawozdanie z projektu

---

*Wojciech Jarosz*

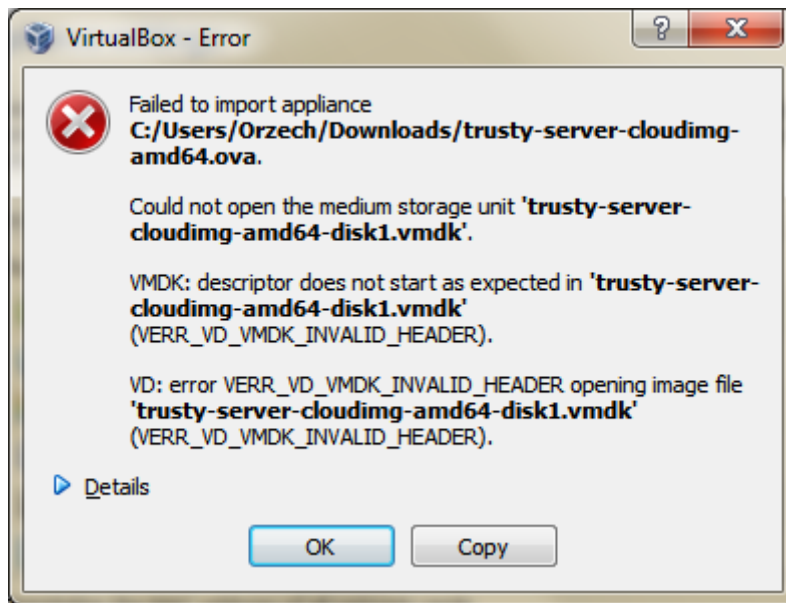
## **Zadanie**

Moim zadaniem było przygotowanie obrazu z nowymi wersjami programów Hadoop next gen (YARN) oraz Apache Spark, który to obraz byłby kompatybilny z platformą OpenStack. Następnie wraz z Jakubem Sawickim mieliśmy przygotować plugin do programu Sahara, który ten obraz miał uruchamiać bez potrzeby dalszej konfiguracji (logowania się do instancji).

## **Część I - Utworzenie obrazu przez skrypt diskimage create oraz program diskimage builder**

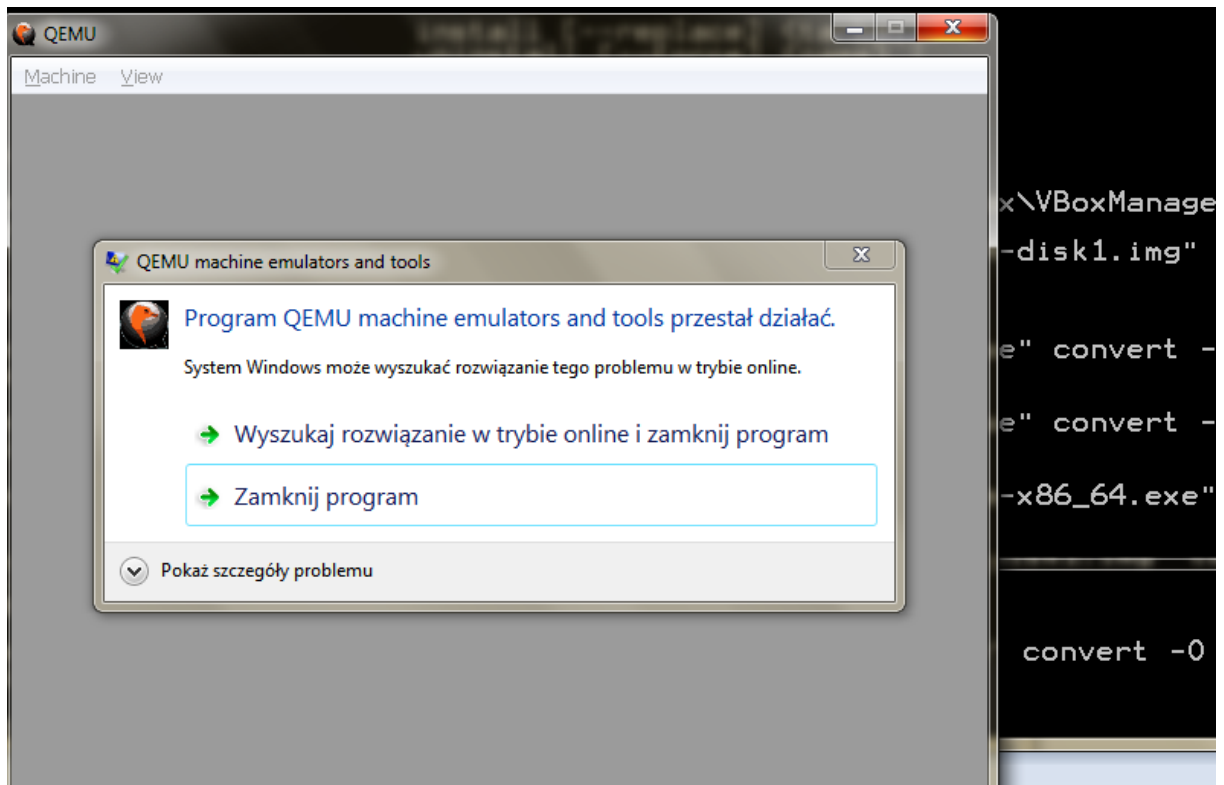
- Próbowałem użyć dostarczonego już skryptu diskimage-create
- Skrypt miał bardzo dużo ograniczeń, np. wspierał jedynie Sparka w wersji standalone, a także Hadoopa w wersji Cloudera (CDH5)
- „Hackując” skrypty udało mi się podnieść wersję Hadoopa na wytworzonym w ten sposób obrazie
- Nie udało mi się podnieść wersji Sparka (choć „hackowanie” skryptów było już na zaawansowanym etapie), ponieważ dostałem informację, że dalsze działanie w tym kierunku jest pozbawione sensu

**Czas poświęcony: ok. 10 godzin – na zapoznanie się ze skryptami, rozbiór ich na czynniki pierwsze, „Zhackowanie” hadoopa oraz próby modyfikacji skryptu tak, aby instalował tylko potrzebne rzeczy**



## Część II – Próba uruchomienia „czystego” obrazu

- Ściągnąłem obraz ubuntu cloud image
- Próbowałem operować na nim za pomocą emulatora qemu, ponieważ zauważyłem że wszystkie obrazy w OpenStack, które mieliśmy, były w formacie QCOW2 naturalnym dla tego emulatora
- Po kilku godzinach wniosek – qemu nie działa na Windowsie



- Odpalenie obrazu na linuxie – sukces, ale nie znam hasła

- Okazuje się, że nie ma hasła – trzeba robić autentykację przy użyciu kluczy – OpenStack je wstrzykuje, co ja mam zrobić? Szukam przez dłuższą chwilę odpowiedzi w Internecie. Nie znajduję jej.
- Pytanie – jak zrobić passwordless ssh do tego obrazu?
- Rozwiązanie – konwersja do VirtualBoxa
- Mija kilka godzin, mimo korzystania z oficjalnej instrukcji [http://docs.openstack.org/image-guide/content/ch\\_converting.html](http://docs.openstack.org/image-guide/content/ch_converting.html) nie jestem w stanie uruchomić obrazu, który byłby w stanie współpracować z VirtualBoxem. Nauczony doświadczeniem, próbuję robić to na Windowsie i Linuxie. Nie daję rady.
- Po wielu próbach i błędzeniu po omacku trafiam na rozwiązanie – należało **najpierw skonwertować obraz do formatu RAW (ok. 10GB) i dopiero następnie do formatu VDI**, pomimo, że oficjalna dokumentacja mówi co innego
- Uruchamiam obraz na VirtualBoxie! Próbuję zrobić passwordless ssh – niestety z marnym rezultatem.
- Podpinam się przy użyciu innych narzędzi do systemu plików obrazu – modyfikuję plik shadow żeby móc się zalogować, używam chroot
- SUKCES, loguję się jako użytkownik ubuntu na maszynie wirtualnej. Konfiguruję dostęp do internetu. Mogę zacząć konfigurować mój obraz!

**Poświęcony czas: ok. 20 godzin, z dobrą dokumentacją byłbym w stanie powyższe zrobić w 10x krótszym czasie**

### Część III - Operacje na obrazie „Trusty cloud image”

Paradoksalnie najkrótsza i najprostsza część

- Instaluję programy, paczki
  - Oracle Java
  - Scala
  - Hadoop 2.6
  - Spark 1.3.1
  - Inne zależne
- Konfiguruję YARNa
  - Zmieniam pliki yarn-site.xml, hdfs-site.xml i podobne. Podaję jako hostname mastera „master” w nadziei, że wpis zostanie dodany do pliku hosts
  - Tworzę zmienne środowiskowe oraz zapisuję je do .bashrc
  - Tworzę skrypty, które na początku polegały na wykorzystaniu pliku /etc/hosts i dodanie do niego wpisu „master”
  - Tworzę skrypt, który zastępuje przestarzałe obecnie start-all.sh

**Poświęcony czas: ok. 4 godziny**

## Część IV – lokalne testy

Manualnie konfiguruje 2 kopie obrazu na 2 różnych maszynach VirtualBox. Kopiuję klucze między maszynami, aby mogły do siebie zrobić passwordless SSH. Test zakończony powodzeniem – wszystkie procesy są aktywne (ResourceManager, NameNode, DataNode itp.)

**Poświęcony czas: ok. 2 godziny**

## Część V – Wgrywanie obrazu na OpenStacka

Próbuję wgrać obraz na OpenStacka

- Pierwszy testowy obraz udało mi się wgrać bez problemu
- Docelowy obraz wgrywam przez przeglądarkę – błąd „błąd”
- Próbuję kilka razy z różnych przeglądarek – obraz jest wgrywany przez jakiś czas (ok. 15 minut na próbę), ale rezultatem jest zawsze nic nie mówiący błąd
- Wgrywam obraz na swój serwer postawiony w Krakowie na RaspberryPi – ponieważ łącze jest nie najlepszej jakości, zajmuje mi to kilka godzin
- Podaję URL do obrazu na RPI. Wygląda na to, że obraz się ściągnie. Dodał się do listy. Mija 10 minut – błąd.
- Pomyślałem, że jest to wina RPI albo łącza w Krakowie. Nielegalnie wgrywam obraz na serwer firmy, której robiłem kiedyś stronę internetową.
- Podaję URL do obrazu na serwerze firmy. Obraz się tworzy – po 10 minutach błąd w rodzaju „nie można stworzyć obrazu”
- Po konsultacji z prowadzącym, mam zapytać się Kuby, jak tworzył swoje obrazy – od niego dostaję informację, że robił to w identyczny sposób, jak ja.
- Po kolejnej konsultacji z prowadzącym, wgrywam obraz na maszynę bezpośrednio, przez SCP. Korzystam z programu glance do wgrania obrazu.
- Dostaję dużo błędów związanych z autentykacją, jednak dzięki współpracy z Łukaszem Sękalskim udało mi się w końcu uruchomić glance.
- 10 prób wgrania obrazu (który jest już na dysku) kończy się fiaskiem.
- Kilka dni później dowiaduję się, że z uwagi na działania innych grup (obrazy po ponad 5GB), na maszynie nie ma już miejsca

**Poświęcony czas: zależny od metodologii liczenia – od kilku do kilkudziesięciu godzin**

## Część VI – obraz wgrany do OpenStacka

- Okazuje się, że obecny plugin do Sparka jest kompletnie niezgodny z moim obrazem, z uwagi na to, że zainstalowałem nowsze wersje oprogramowania oraz np. „czystego” Hadoopa zamiast CDH.
- Plugin wykonuje przestarzałe skrypty, na pewno będzie trzeba go sporo zmienić

- Próby skomunikowania się ze sobą mastera i slave'a kończą się fiaskiem. Manualnie kopiuję klucz, który powinien w przyszłości zostać wstrzyknięty przez plugin
- Duże ostrzeżenie przy próbach modyfikowania /etc/hosts – wymagana będzie inna metoda
- Poprawiam skrypty i piszę nowe, które już są mniej uniwersalne, ale piszą bezpośrednio do konfiguracji (podmieniają placeholders), zamiast do /etc/hosts
- Wygląda na to że po napisaniu pluginu obrazy będą działały. Po kilku manualnych „hackach”, wygląda na to, że odpowiednie procesy się uruchamiają przy użyciu moich skryptów

Poświęcony czas: ok. 5 godzin

## Część VII – plugin

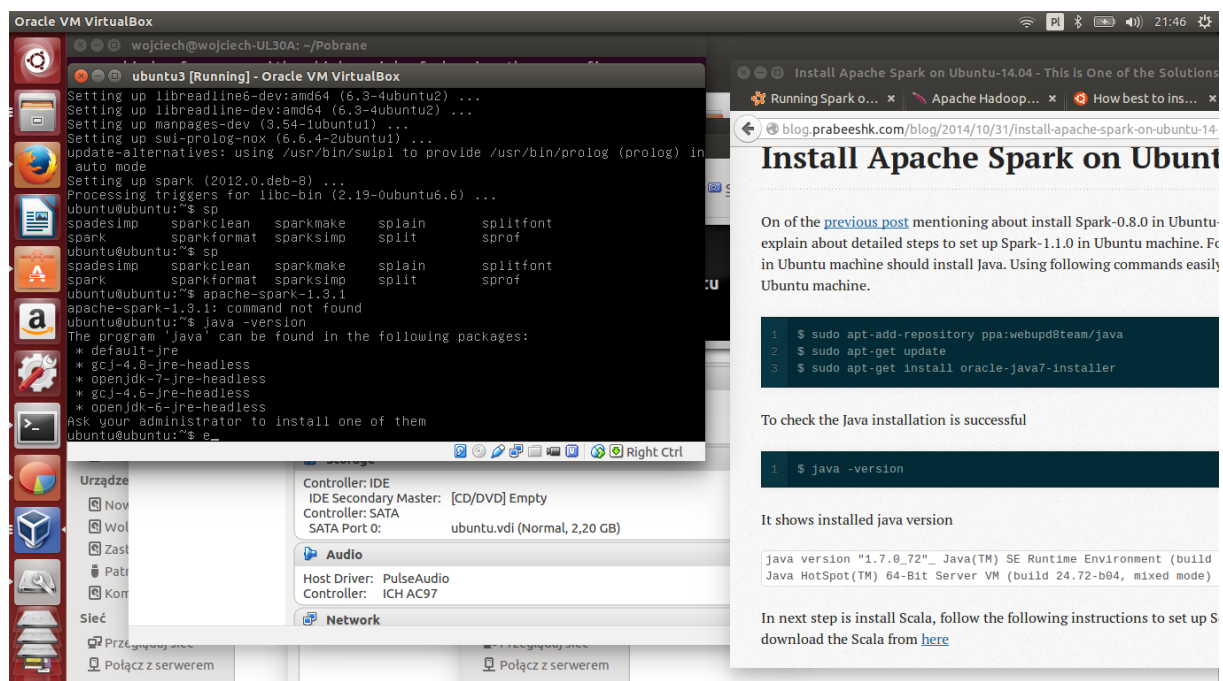
Obecny plugin do Apache Spark jest świetny, ale tylko do edycji Standalone. Przerabiamy go tak, aby nie wołał starych skryptów, tylko ustawiał serwer według moich skryptów oraz konfiguracji, do której doszliśmy metodą prób i błędów.

Główne zasługi na tym polu należą się Kubie, który pomógł mi i Łukaszowi.

Aby usprawnić pracę, przy użyciu narzędzia guestfish, edytuję obraz bezpośrednio na maszynie z OpenStackiem. Powstaje kolejna wersja obrazu z poprawionymi skryptami.

W końcu po wielu udoskonaleniach i przeróbkach instancje wstają z YARNem gotowym do pracy.

Poświęcony czas: ok. 6 godzin.



Ostatecznie jednak plugin jest bardzo ubogi, ale niestety nie mieliśmy czasu, żeby go dokończyć. Przez 2 tygodnie byliśmy wyłączeni przez działania innych grup i braku miejsca na maszynach, a lwia część pracy przypadła na okres egzaminów.