

# 1. Hadoop 개요 및 Master 서버 설정

## 목 차

1. Hadoop 시스템 개요
2. Master 서버 생성 및 환경 설정

# Hadoop+Hive+Spark 수업 목적

- 최신 버전 Linux server 기반 Cluster 구축
- Hadoop 설치와 환경설정 및 구동
- Hive SQL 활용 Hadoop 데이터 연동
- Spark 기반 분산 데이터 처리

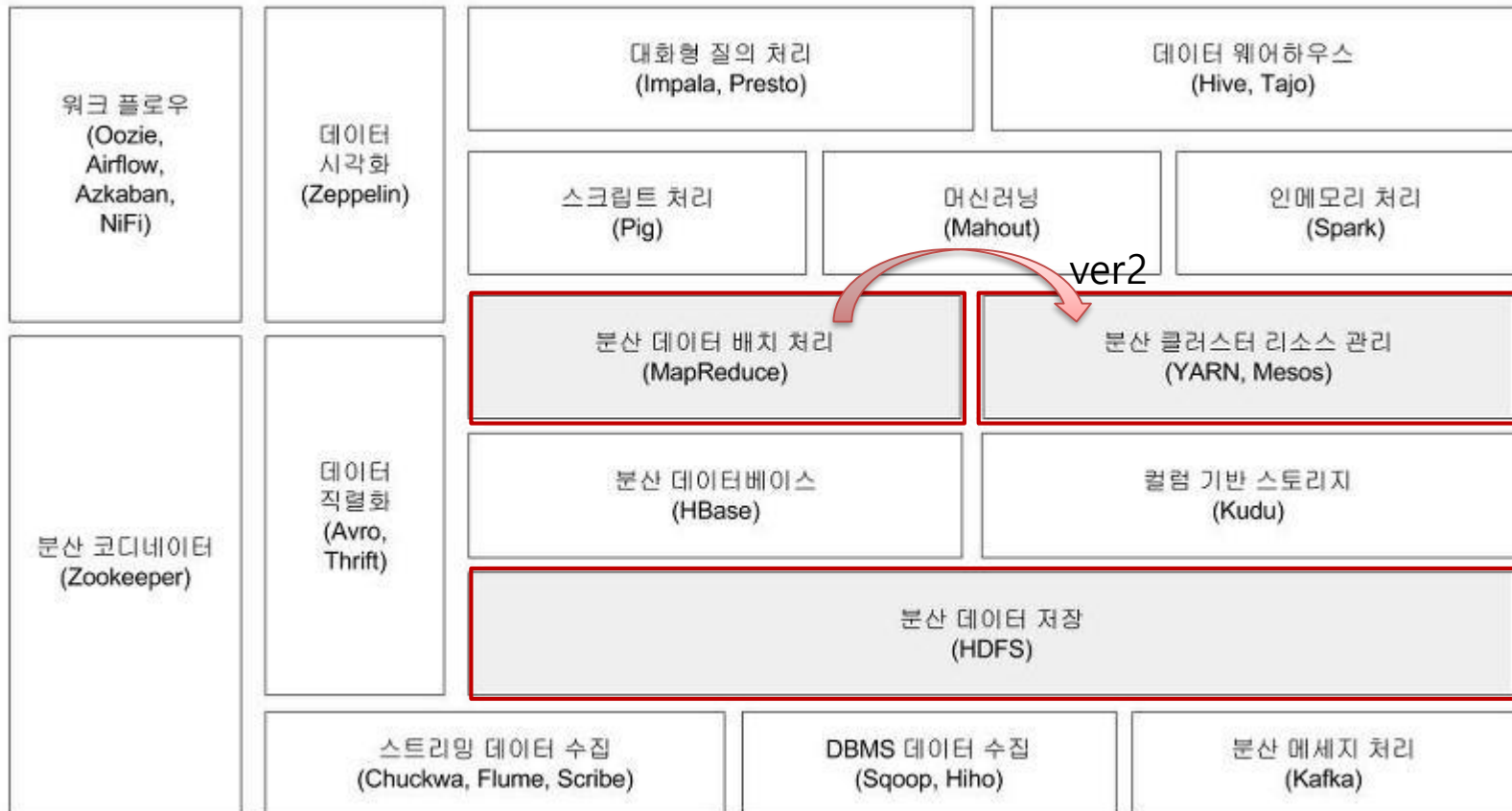
# Hadoop 시스템 개요

- 구글 : GFS(Google File System)와 맵리듀스(MapReduce)
- 2005년 더그커팅 구현
- 2008년 아파치 최상위 프로젝트 승격
- 주요 구성
  - HDFS(Hadoop Distributed File System) : 저장소
  - Map Reduce : 분산처리 시스템

# Hadoop 이점

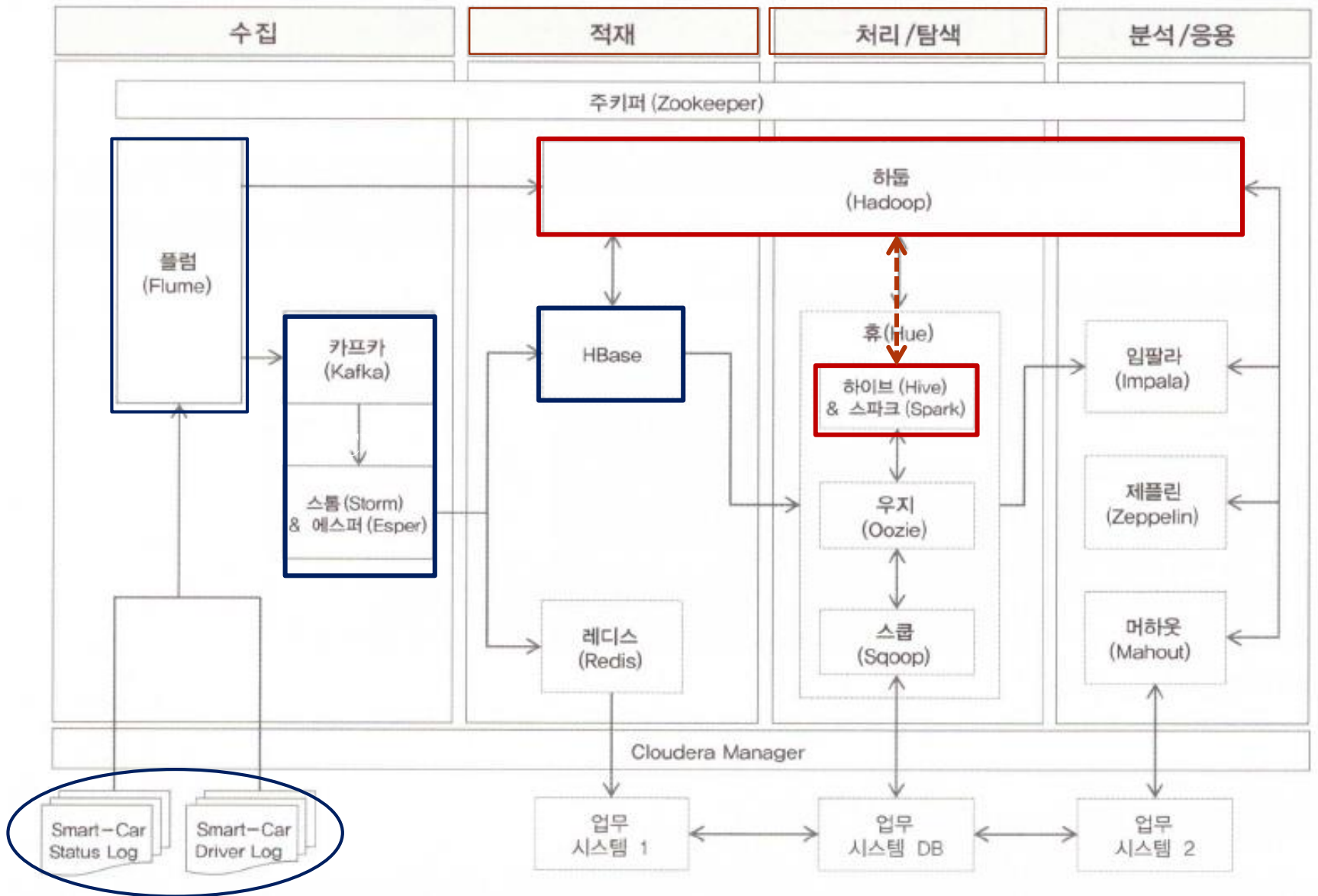
- 여러 대 서버에 데이터 저장
- 각 서버에서 동시 데이터 처리(분산처리)
- 기존의 RDBMS 대체
- 고가 장비 대신 리눅스 서버로 대체
- 오픈 소스
  - SW 라이선스, 고가 HW 비용 절감

# Hadoop 에코 시스템

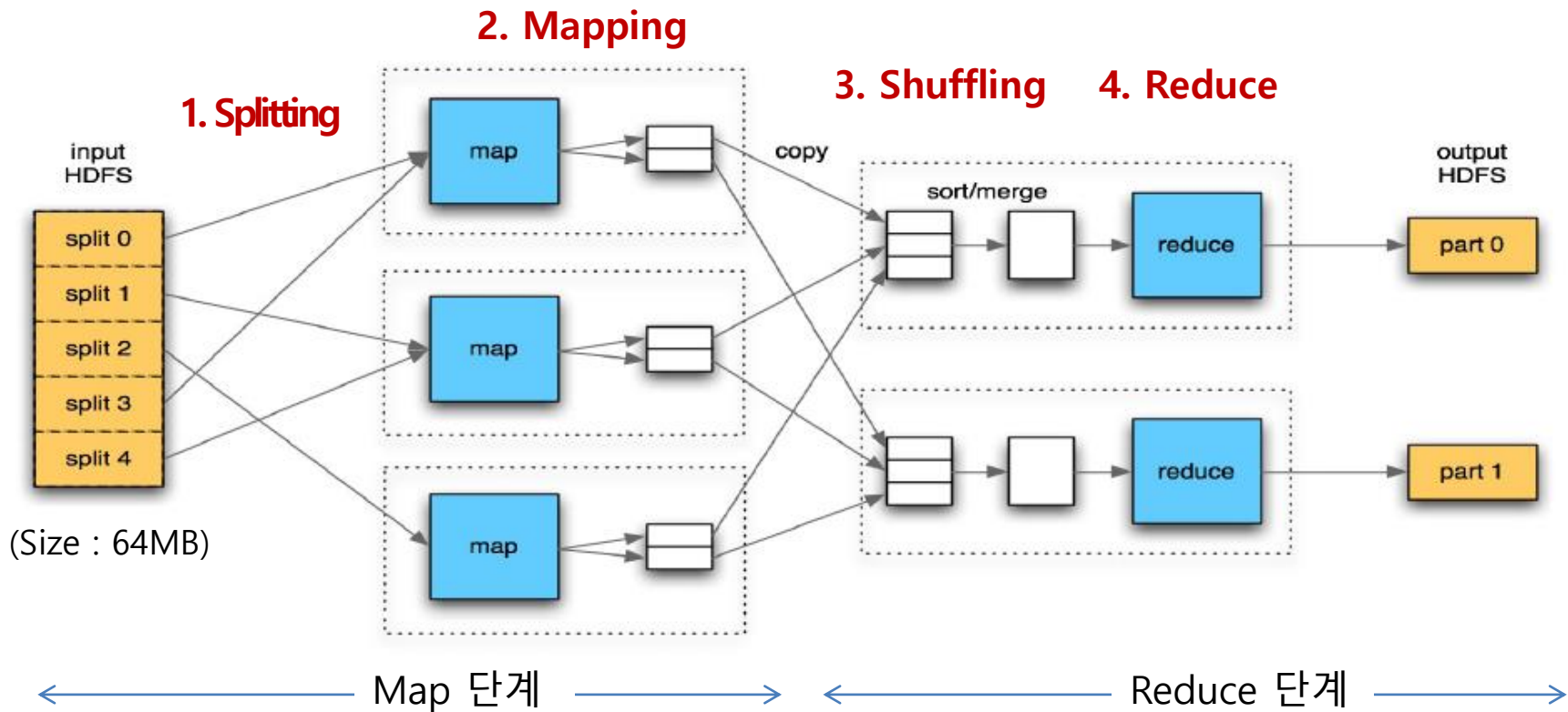


시작하세요. 하둡프로그래밍 – 위키북스 참고

# Hadoop+Hive+Spark 소프트웨어 아키텍처



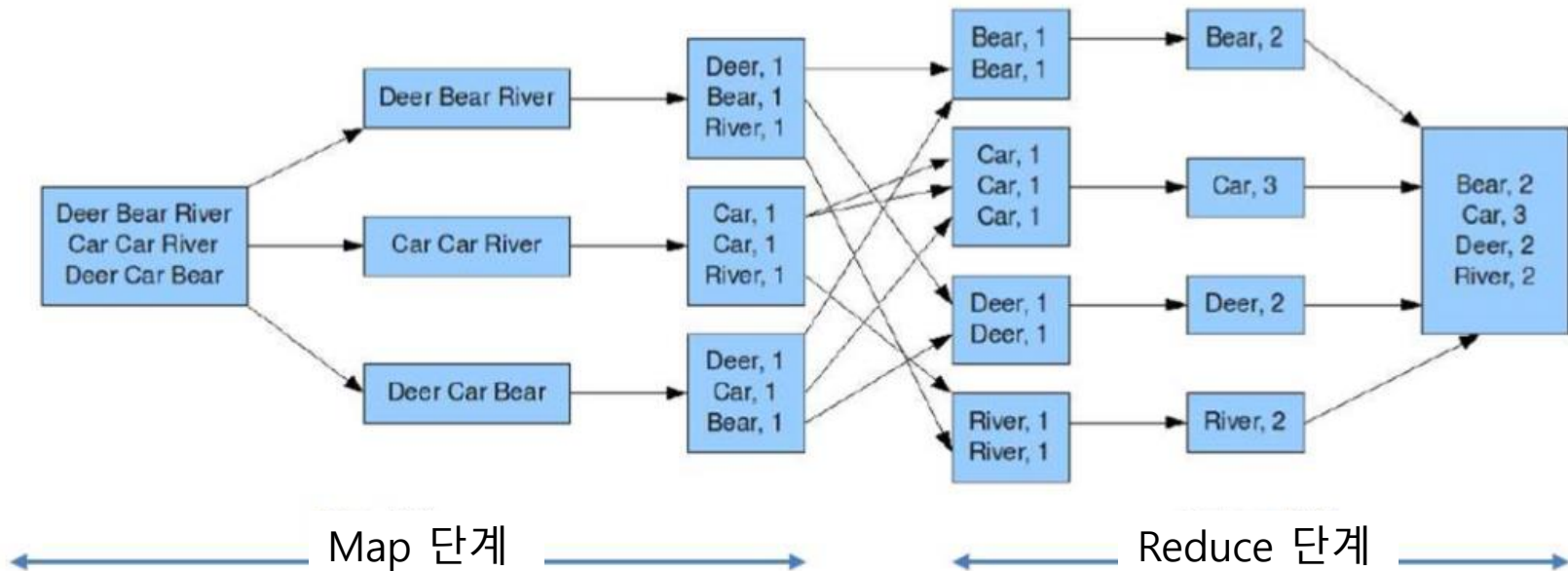
# Hadoop 분산 병렬 처리



# Word Count 예

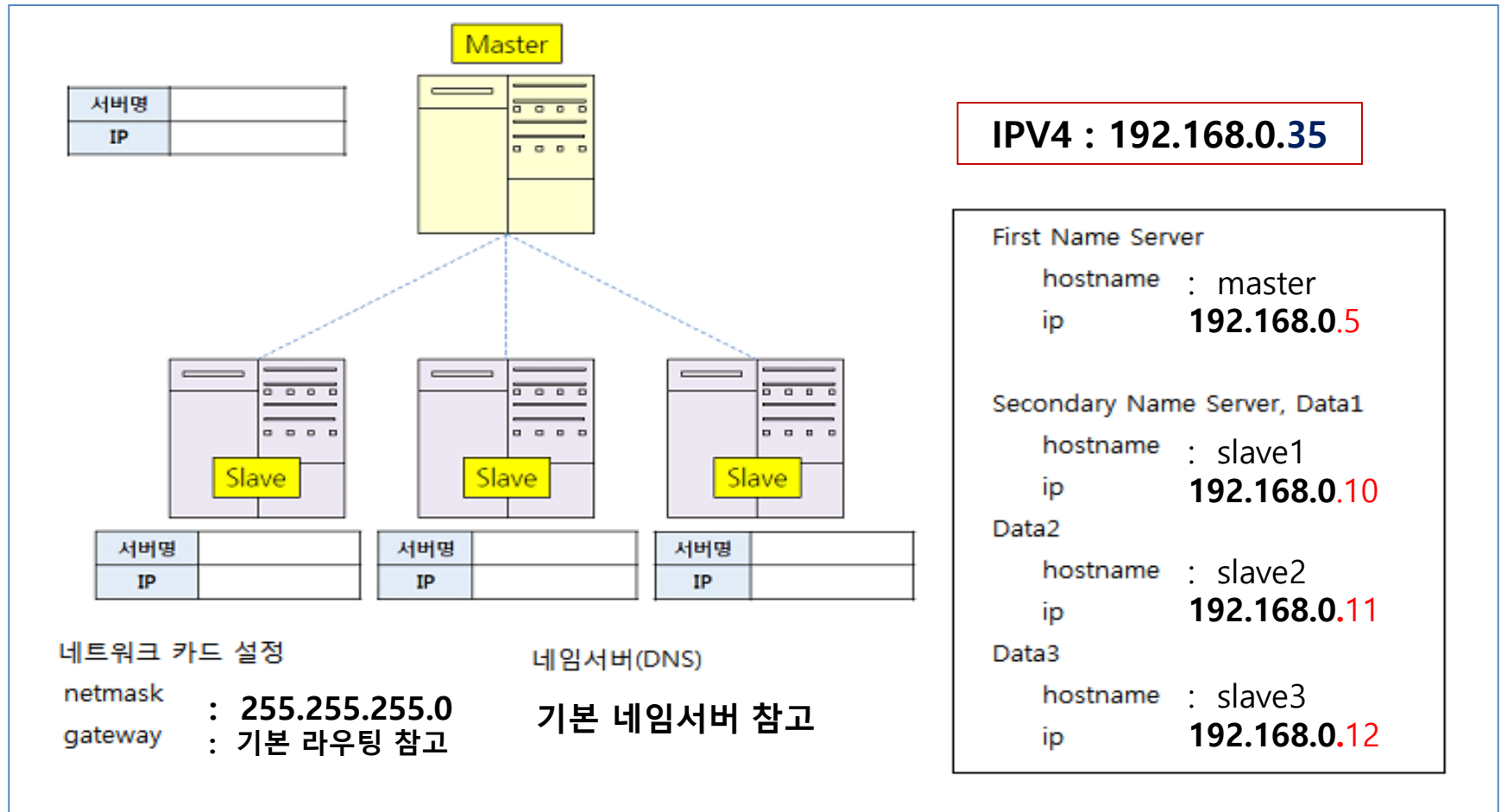
The overall MapReduce word count process

**Input** → **Splitting** → **Mapping** → **Shuffling** → **Reduce** → **Output**



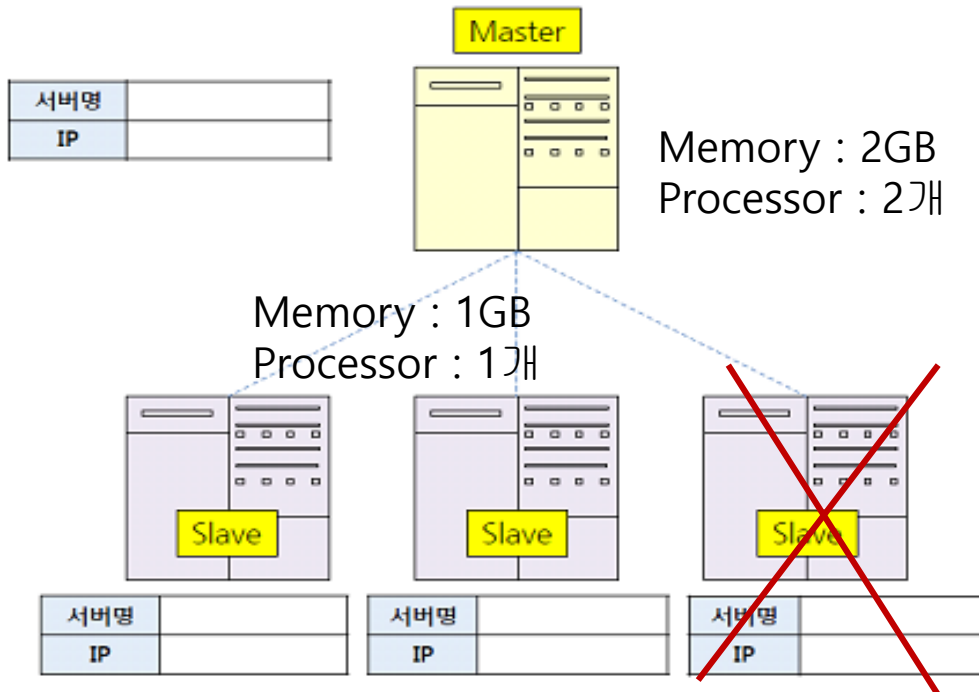


# Hadoop 클러스터 구축 시스템 구성도



# Hadoop 클러스터 구축 시스템 구성도

Linux  
실제 IP



네트워크 카드 설정

netmask : 255.255.255.0  
gateway : 기본 라우팅 참고

네임서버(DNS)

기본 네임서버 참고

IPV4 : 192.168.0.35

First Name Server

hostname : master  
ip : 192.168.0.5

Secondary Name Server, Data1

hostname : slave1  
ip : 192.168.0.10

Data2

hostname : slave2  
ip : 192.168.0.11

Data3

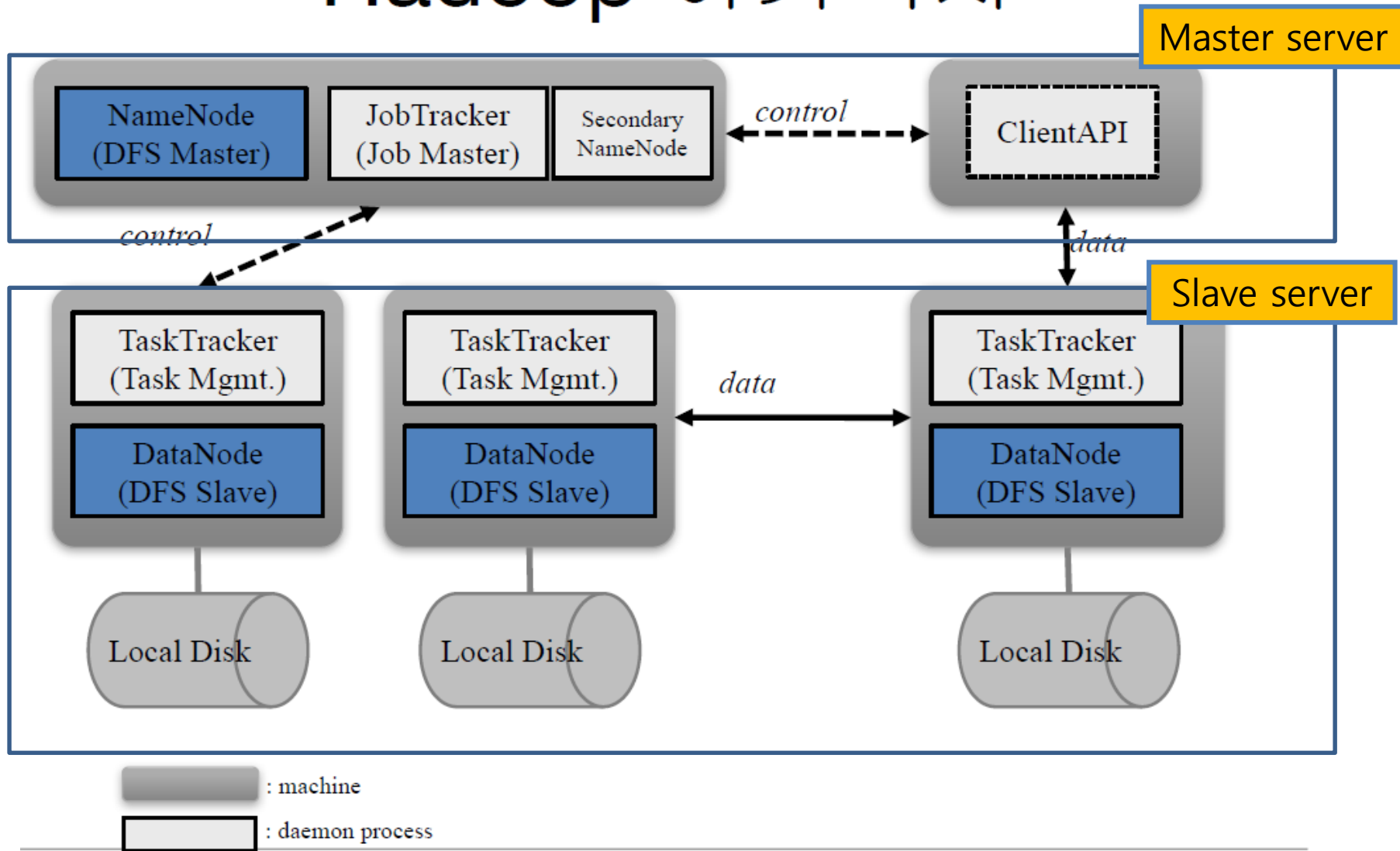
~~hostname : slave3~~  
~~ip : 10.0.2.12~~

# Master vs Slave server

Master Server	Slave Server
<ul style="list-style-type: none"><li>➤ Name node, job Tracker, 보조 네임노드 설치 서버</li><li>➤ <u>디스크 : 2 ~ 4개</u></li><li>➤ <u>메모리 : 32GB ~ 128GB</u><ul style="list-style-type: none"><li>✓ 64GB : 1억 파일 저장</li></ul></li><li>➤ <u>CPU : 코어 수 16~24개</u></li></ul>	<ul style="list-style-type: none"><li>➤ Data node, Task Tracker 설치 서버</li><li>➤ <u>디스크 : 4 ~ 12개(1TB~3TB)</u></li><li>➤ 메모리 : 24GB ~ 48GB</li><li>➤ CPU : 쿼드 코어 2개 이상</li></ul>

- ✓ Job Tracker : Job 실행 요청 받고, Task Tracker에 할당
- ✓ Task Tracker : 실제 Task 실행, Job Tracker 진행 상태 보고

# Hadoop 아키텍처



# Name node(Master Server)

- HDFS에서의 NameNode(master)는 분산환경에서 Job 분배, 지시 및 감독
- 작업의 대상이 되는 파일을 블록(block)단위로 나누어서 slave node 분배
- 메타데이터 저장
  - ✓ 파일 별 블록, 각 블록 별 정보(DataNode 위치 등)
  - ✓ 실제 파일의 블록은 DataNode에 저장됨
- 전체 분산 파일시스템 이상 유무 체크, DataNode(slave) 데이터 입출력 작업 (low-level I/O tasks) 지시
- NameNode 문제 발생 시 Hadoop 클러스터 전체 시스템 차단
  - ✓ Secondary NameNode 이용(NameNode 이중화) 문제 해결

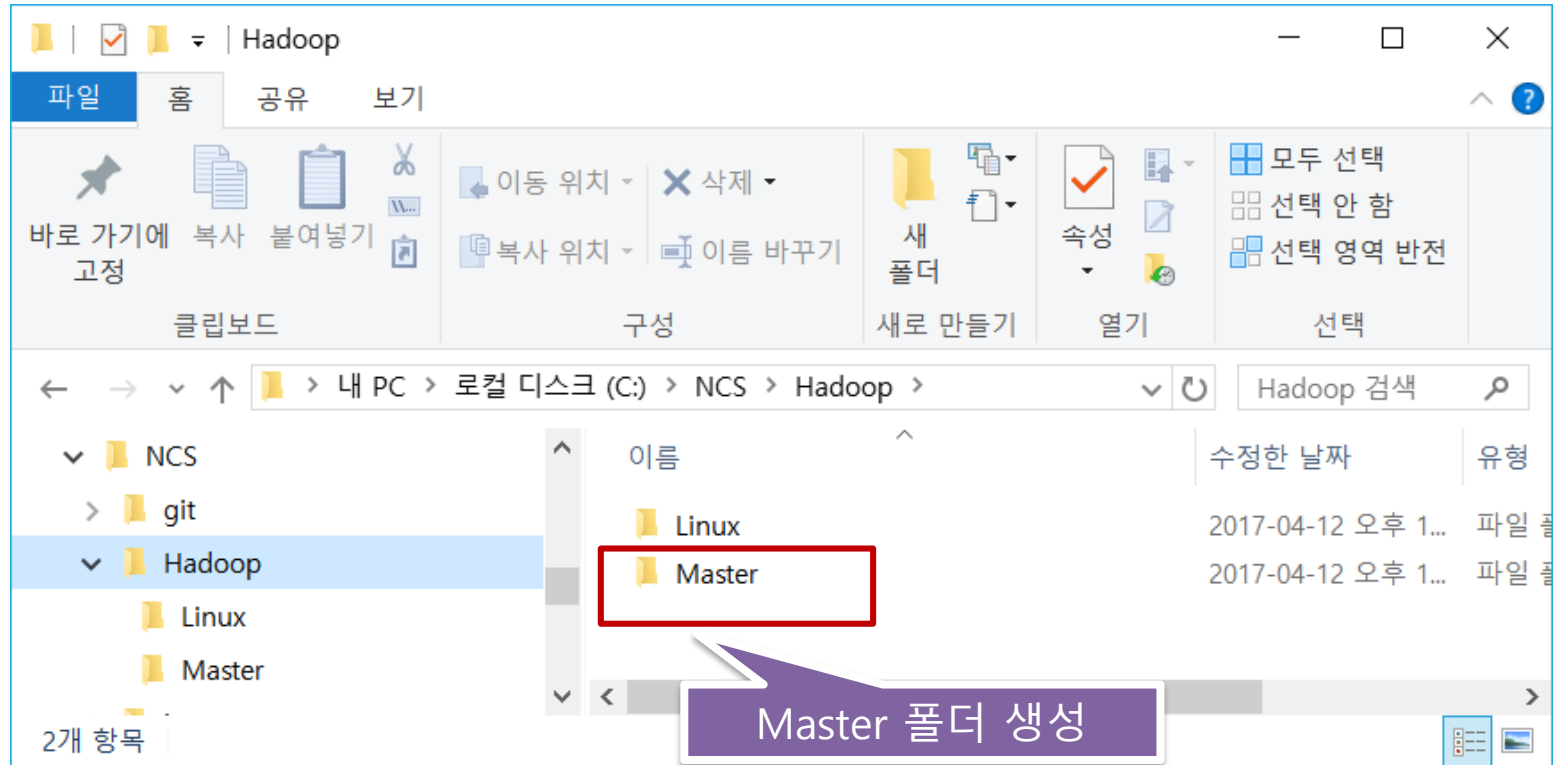
# Data node(Slave Server)

- 실제 HDFS(분산파일시스템) 대상으로 읽기 및 쓰기의 모든 작업 수행
- 모든 작업은 대상 파일을 random하게 블록(block) 단위로 나누어 진행
- DataNode 작업은 NameNode에 수시 보고, 메타데이터 형식 저장
- 데이터를 읽어 들이는 즉시 각 Data node에 분배(Name node)
  - ✓ HDFS는 큰 데이터 파일을 여러 개로 분리시켜서 각 node 처리
- 각각의 조각(chunk) 는 여러 대의 컴퓨터에 중복적으로 복제
  - ✓ 한 컴퓨터에서 장애가 발생해도 다른 컴퓨터를 통해 데이터 이용
- 모든 파일조각들은 하나의 namespace를 공유
  - ✓ 클러스터 내의 모든 node들은 공유된 파일 이용

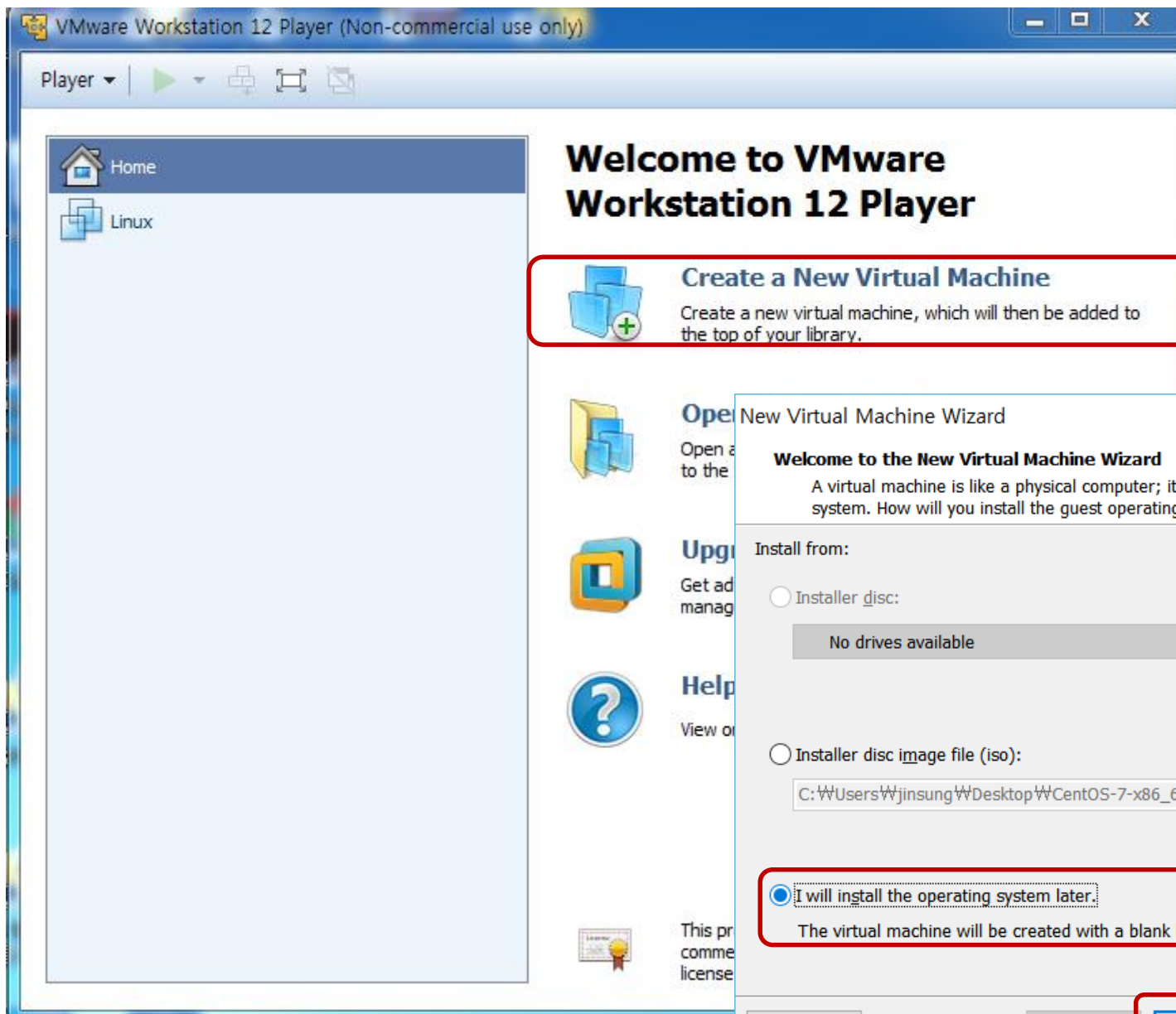
## 2. Master 서버 생성/환경 설정

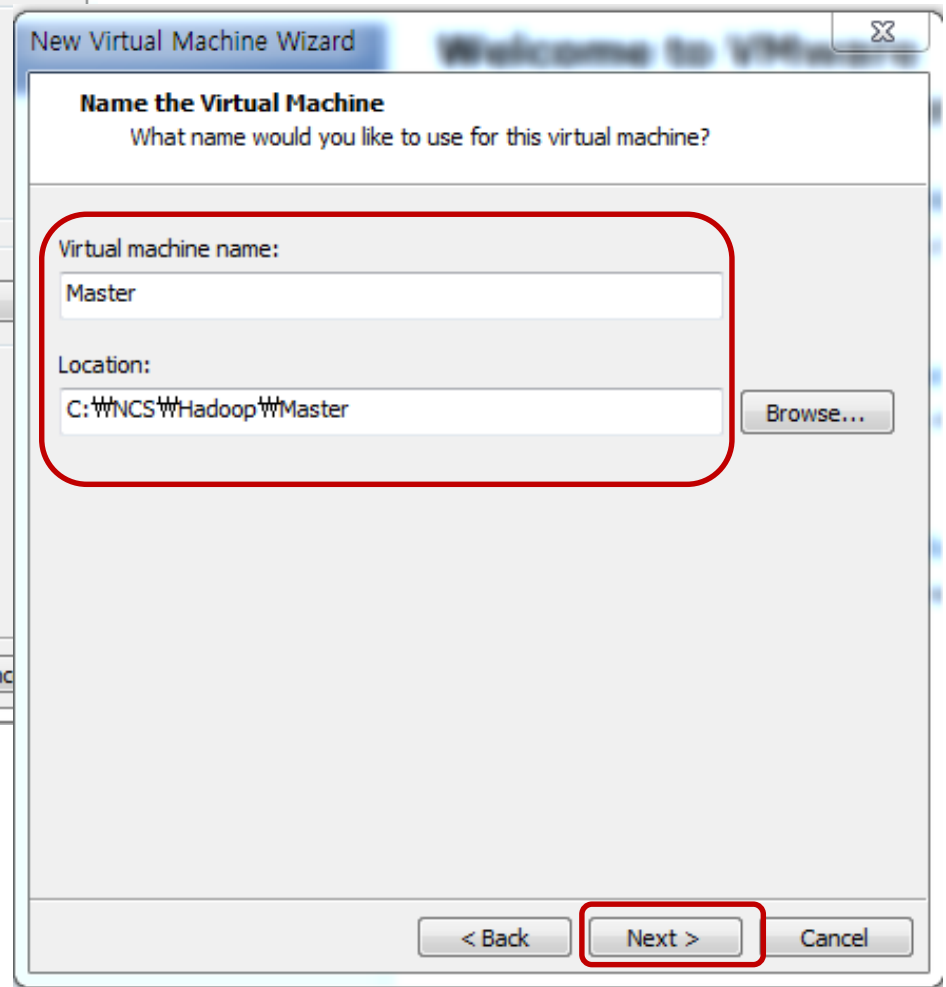
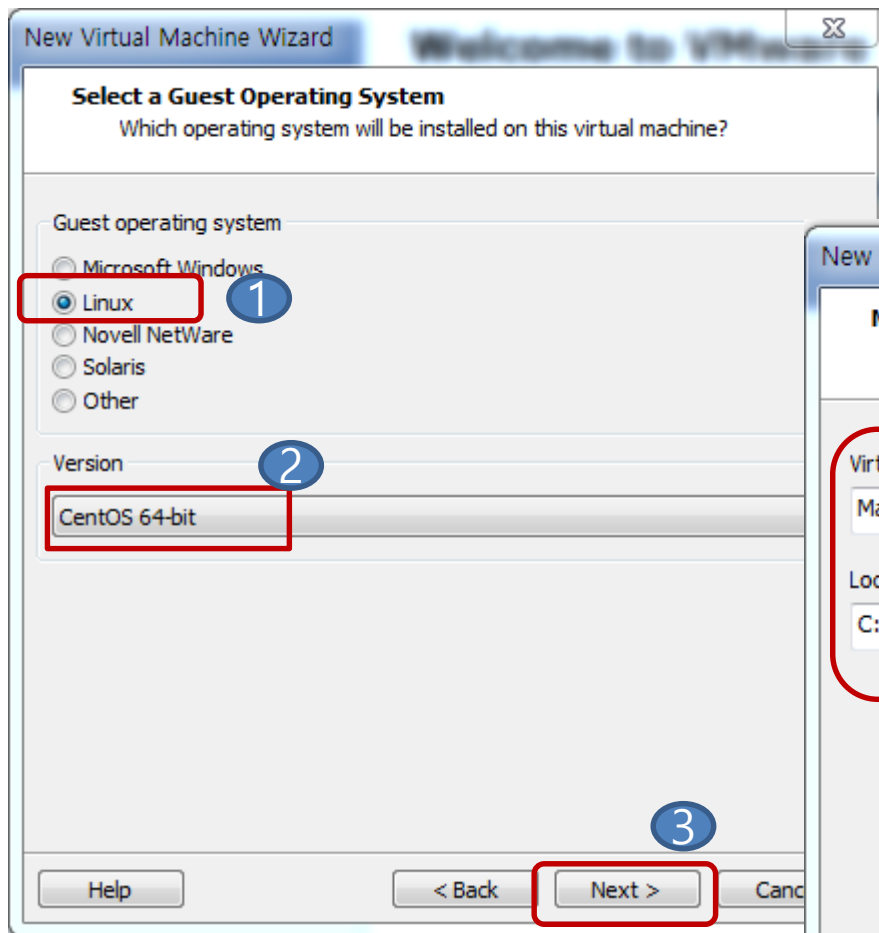
- 1) Master 서버 생성
- 2) 방화벽 제거
- 3) 네트워크 설정
- 4) 호스트 네임설정
- 5) Java 설치
- 6) Slave 폴더 생성(Master 폴더 복사)

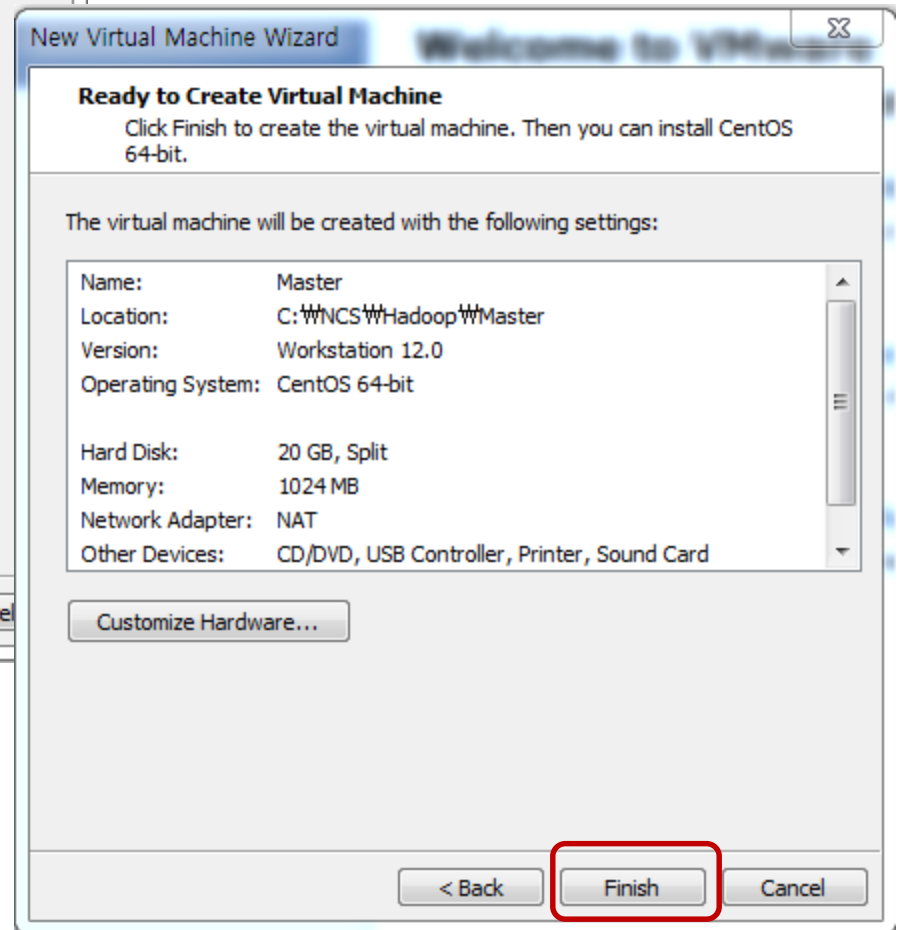
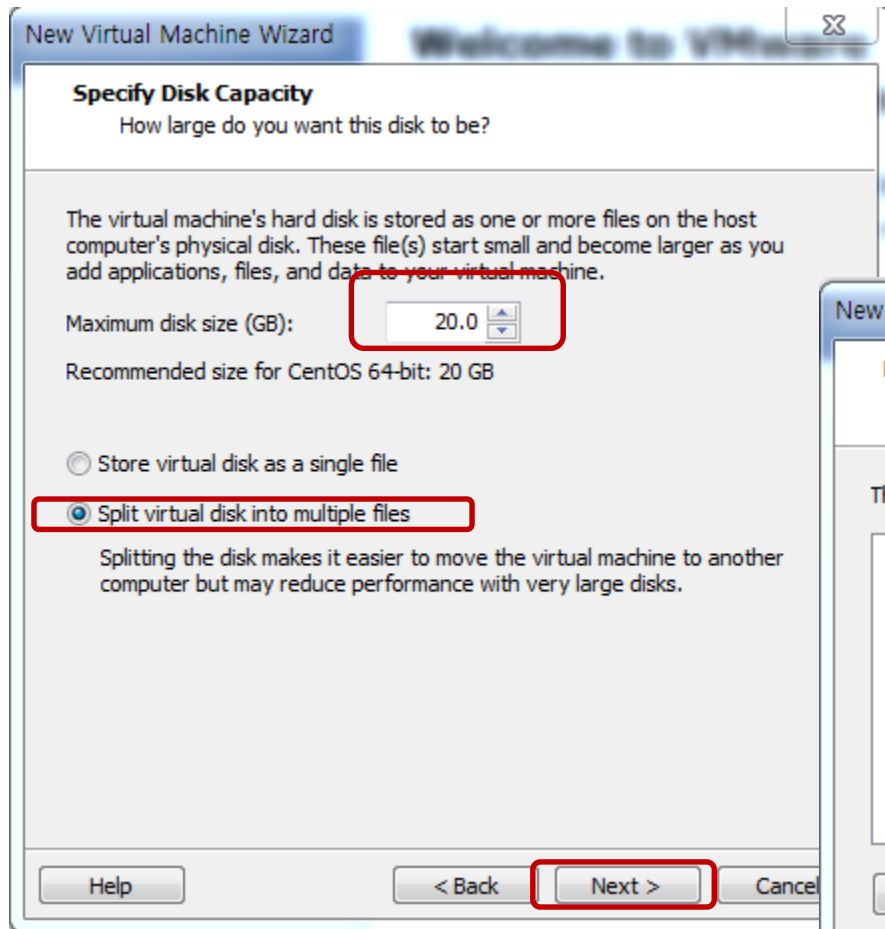
# 1) Master 서버 생성

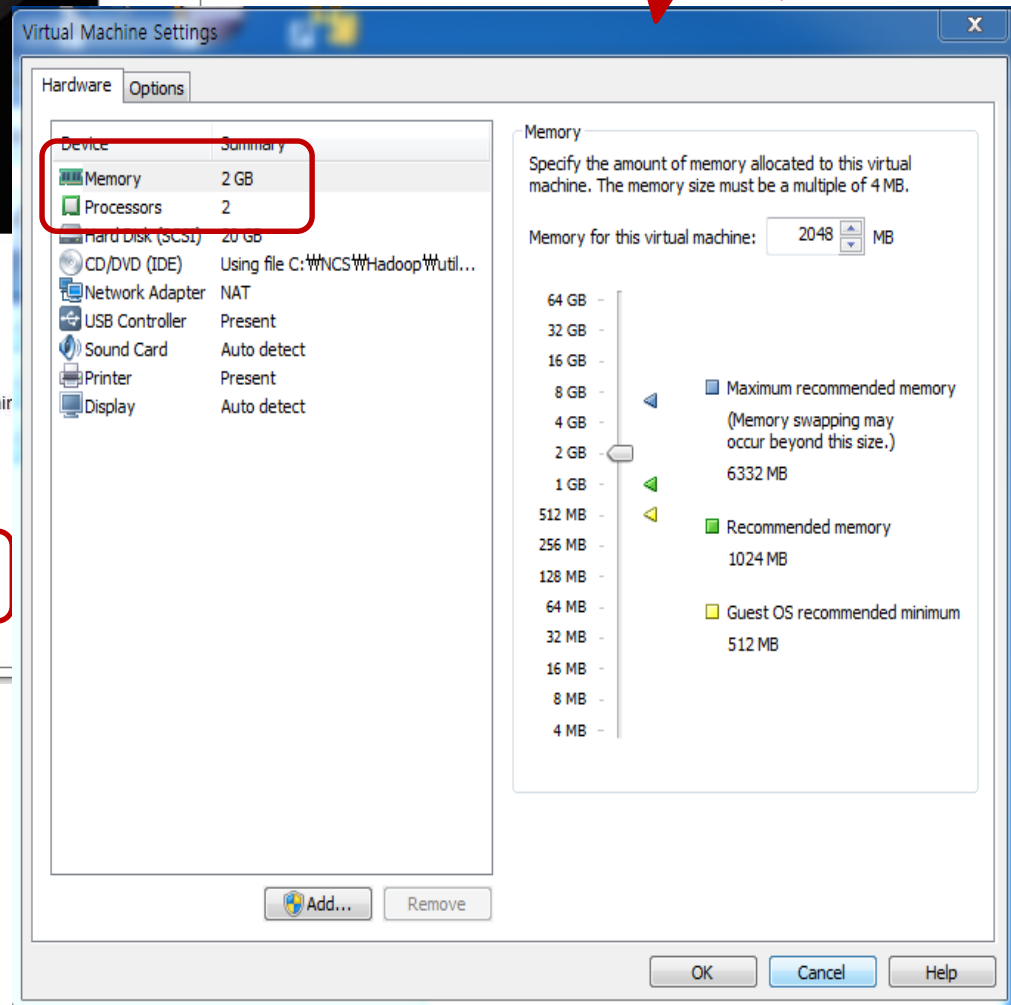
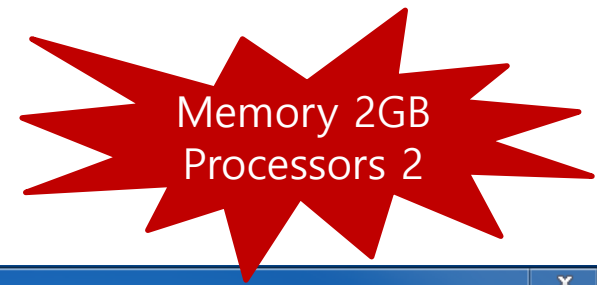
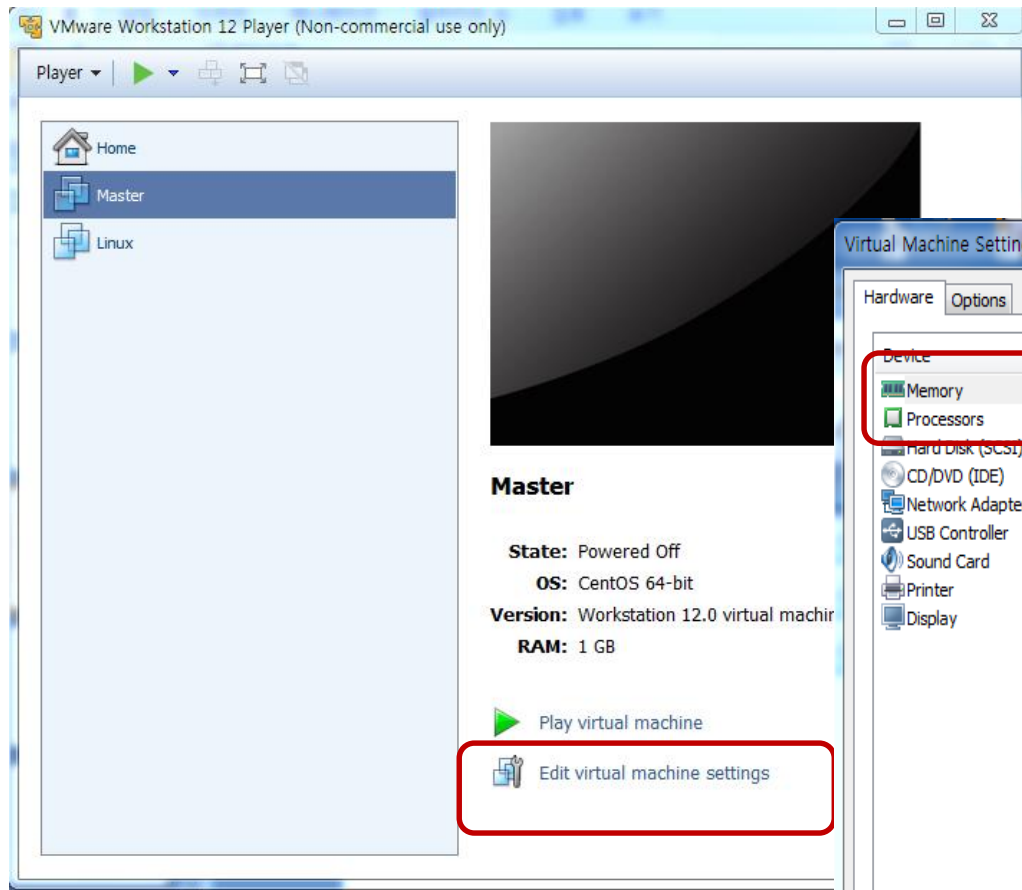


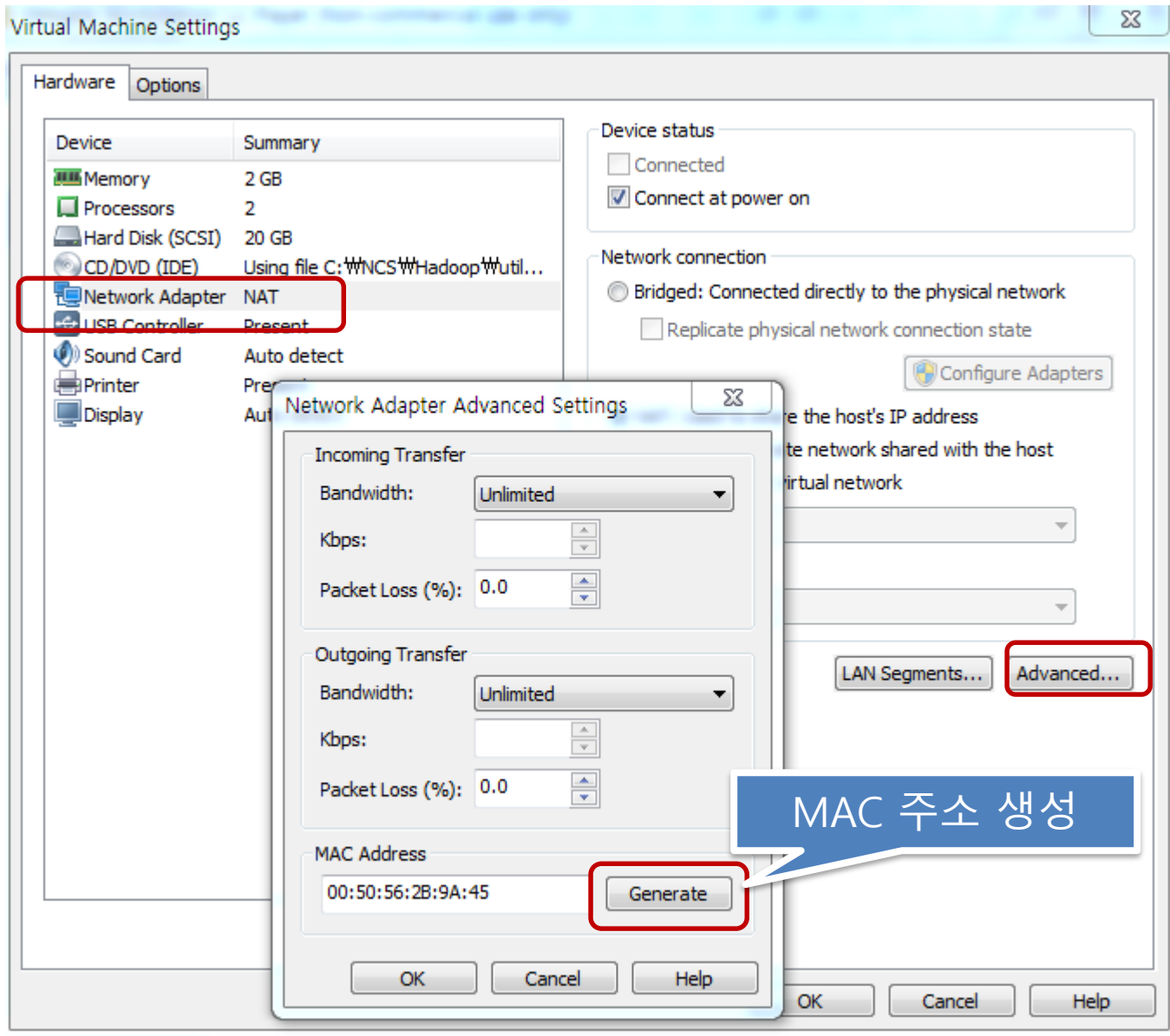


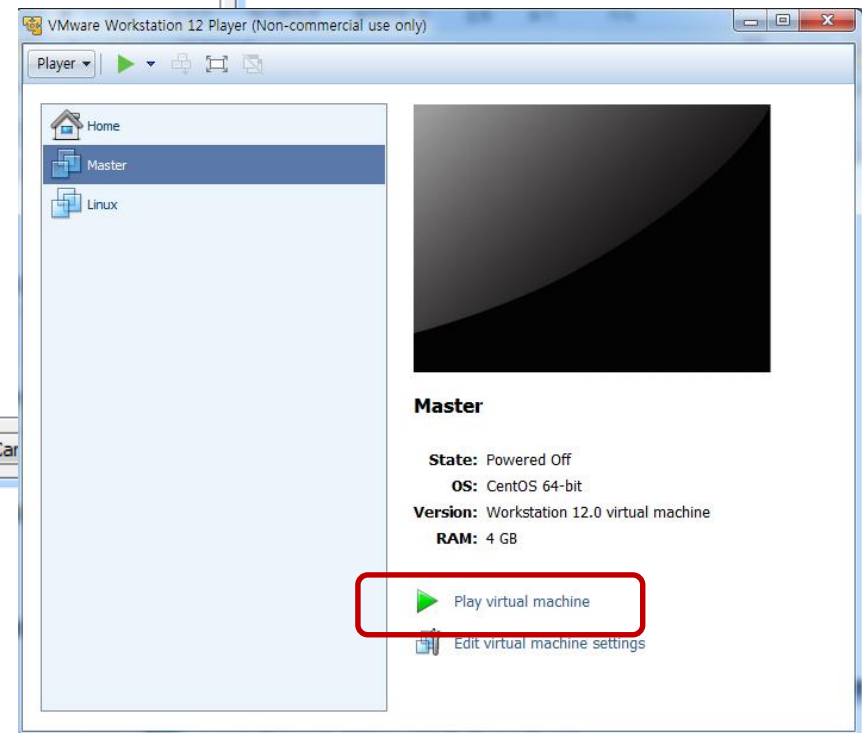
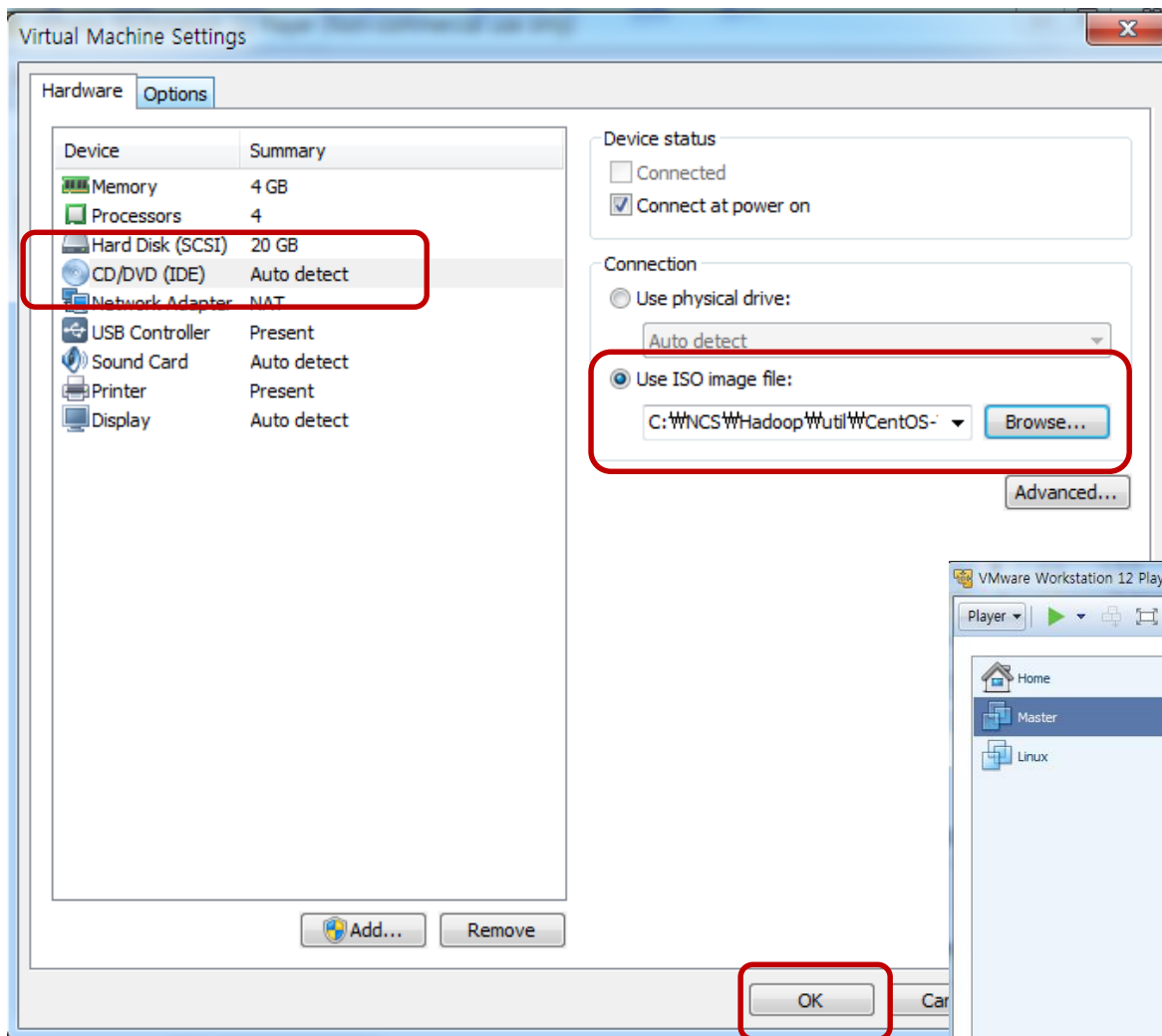


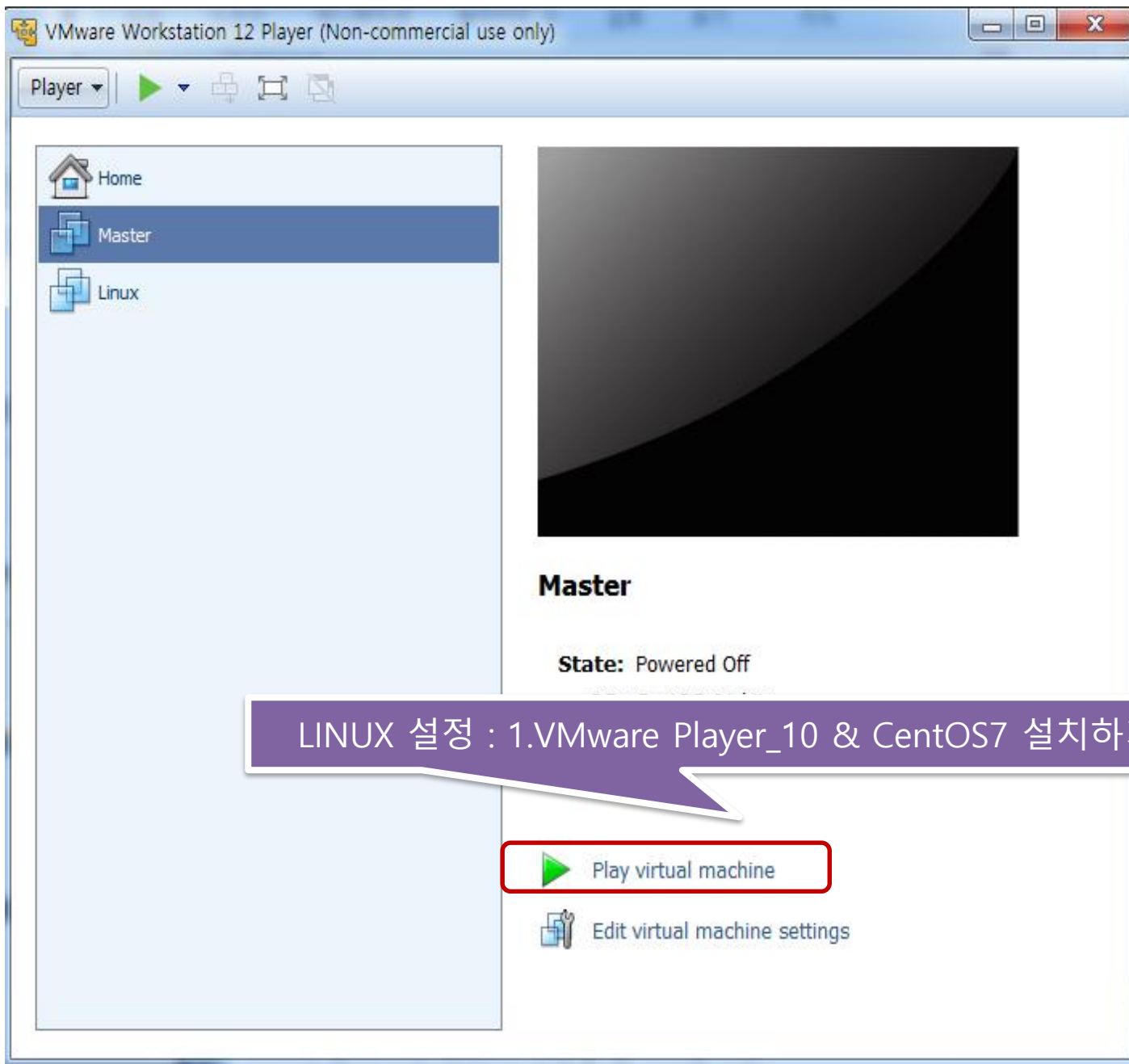










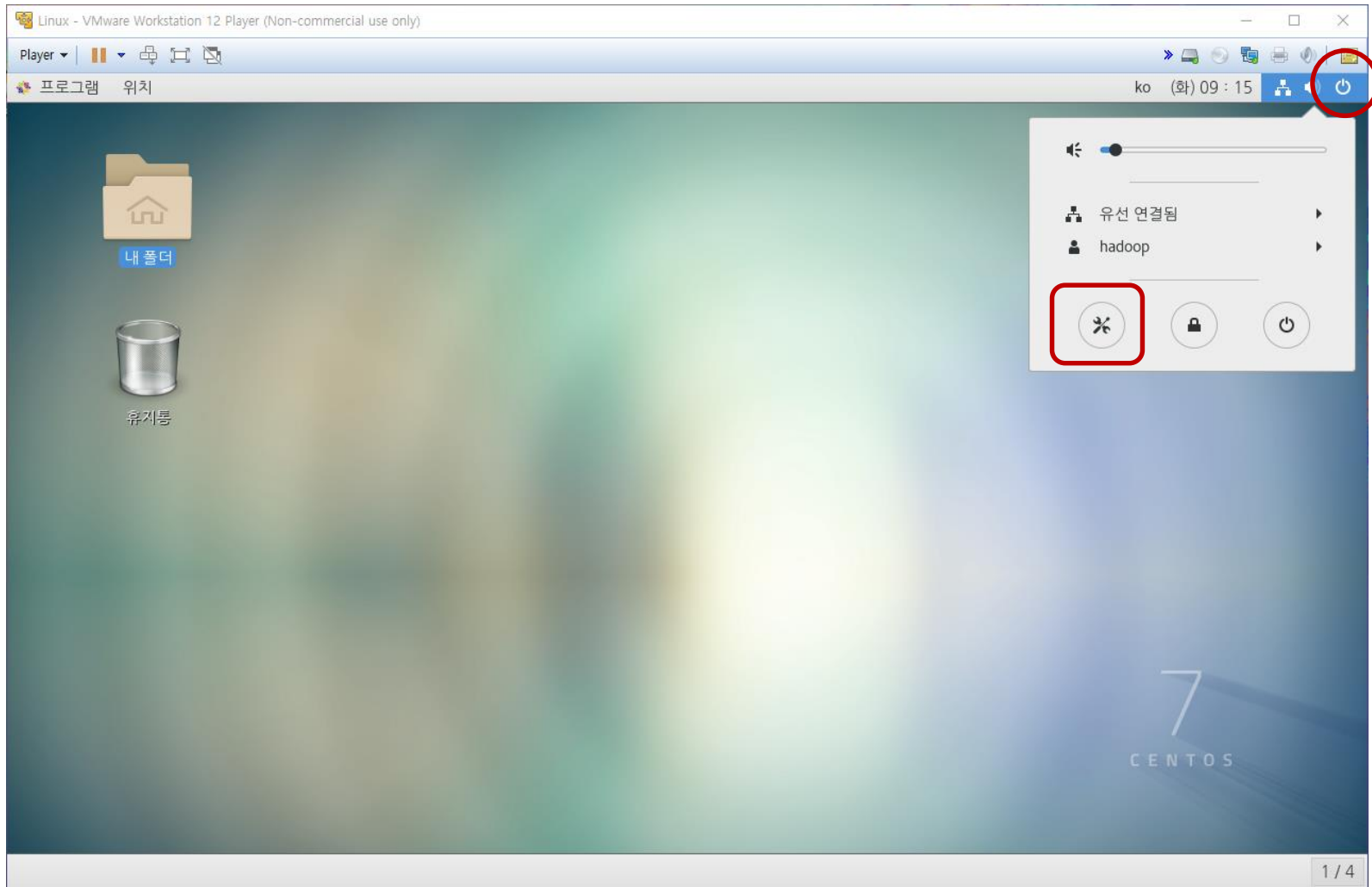


## 2) 방화벽 제거

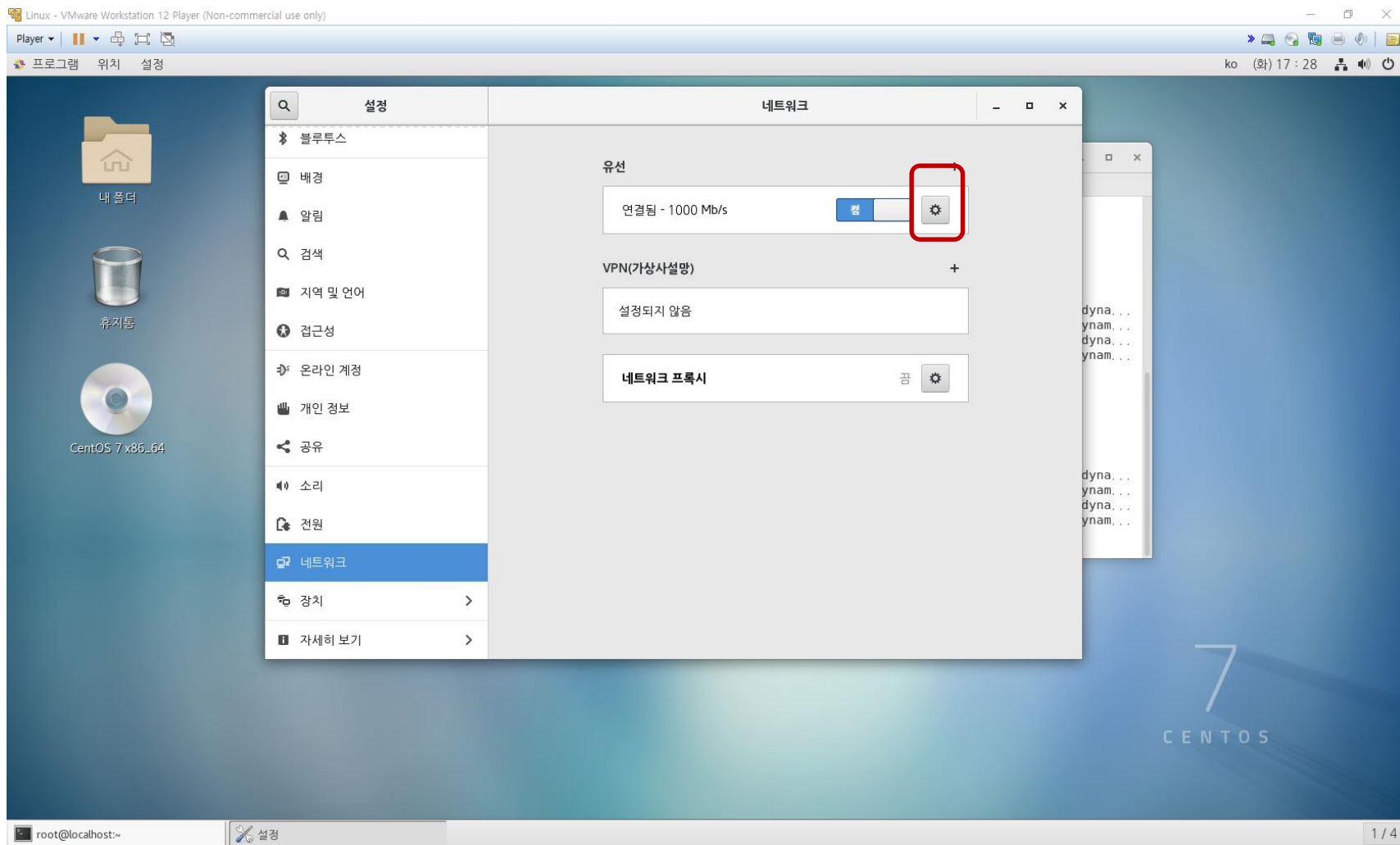
```
root@localhost:~  
파일(F) 편집(E) 보기(V) 검색(S) 터미널(T) 도움말(H)  
[hadoop@localhost ~]$ su -  
암호:  
[root@localhost ~]# systemctl status firewalld.service  
● firewalld.service - firewalld - dynamic firewall daemon  
   Loaded: loaded (/usr/lib/systemd/system/firewalld.service; enabled; vendor preset: enabled)  
   Active: active (running) since 화 2017-03-07 12:54:11 KST; 4h 34min ago  
     Docs: man:firewalld(1)  
  Main PID: 728 (firewalld)  
    CGroup: /system.slice/firewalld.service  
            └─728 /usr/bin/python -Es /usr/sbin/firewalld --nofork --nopid  
  
3월 07 12:54:09 localhost.localdomain systemd[1]: Starting firewalld - dyna...  
3월 07 12:54:11 localhost.localdomain systemd[1]: Started firewalld - dynam...  
Hint: Some lines were ellipsized, use -l to show in full.  
[root@localhost ~]# systemctl stop firewalld  
[root@localhost ~]# systemctl mask firewalld  
Created symlink from /etc/systemd/system/firewalld.service to /dev/null.  
[root@localhost ~]# █
```



### 3) 네트워크 설정



### 3) 네트워크 설정



### 3) 네트워크 설정

Linux - VMware Workstation 12 Player (Non-commercial use only)

Player | 프로그램 위치 설정 ko (화) 17:28

내 폴더  
휴지통  
CentOS 7 x86\_64

설정

네트워크

취소(C) 유선 적용(A)

자세히 보기

전송 속도 1000 MB/s

IPv4 주소 172.16.119.128

IPv6 주소 fe80:a901:9db2:b117:6107

하드웨어 주소 00:0C:29:97:34:4E

기본 라우팅 172.16.119.2

네임서버(DNS) 172.16.119.2

☒ 자동으로 연결(A)

☒ 다른 사용자가 사용할 수 있게 허용(O)

☐ 백그라운드 데이터 사용 제한  
종량제 데이터 요금 또는 데이터가 제한된 연결에서 사용합니다.

연결 프로필 제거

Linux 실제 IP4 주소 확인

기본 라우팅(게이트웨이) 확인

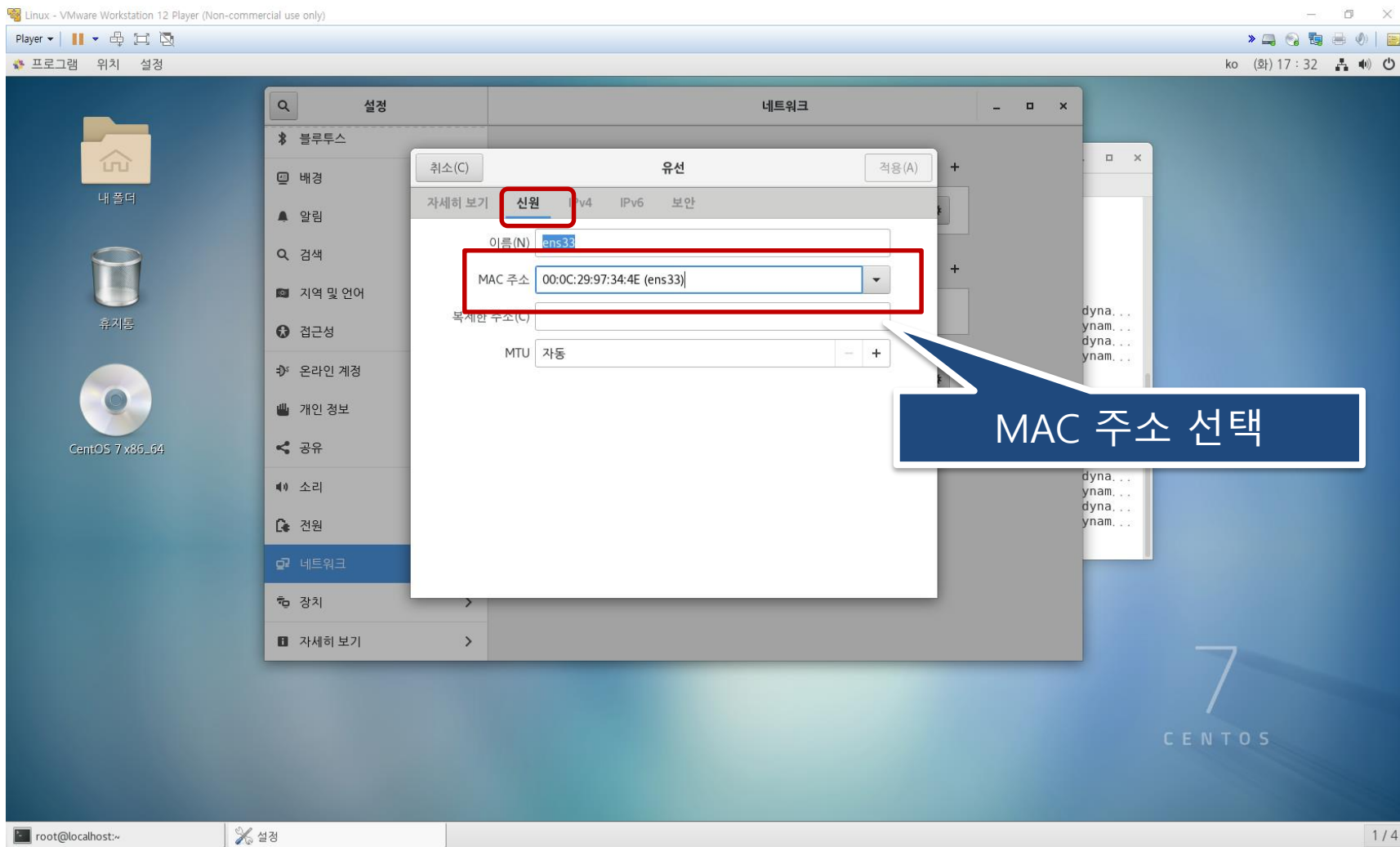
7  
CENTOS

root@localhost:~

설정

1 / 4

### 3) 네트워크 설정



주소 : 실제IP에서 4번째 숫자 변경  
게이트웨이 : 기본 라우팅 동일  
네임서버 : 기본 네임 서버 동일

사용자  
System에  
맞게 지정

The screenshot shows the Windows Network Settings application. The 'Network' (네트워크) tab is selected. The 'IPv4' tab is active, and the 'IPv4 method' (IPv4 방식) is set to 'Manual' (수동). The 'Address' (주소) section is highlighted with a red box and numbered 3, showing the IP address 172.16.119.5, subnet mask 255.255.255.0, and gateway 172.16.119.2. The 'DNS server' (네임서버(DNS)) section is highlighted with a red box and numbered 4, showing the DNS server address 172.16.119.2. The 'Apply' (적용(A)) button is highlighted with a red box and numbered 5. A blue callout box with the text '가상 IP주소 지정' (Virtual IP address specification) points to the IP address field. The background shows the Windows Settings app with various system settings like Bluetooth, Background, Alerts, Search, Region and Language, Proximity, Online status, Personal info, Sharing, Sound, Power, and Network (selected).

설정 네트워크

블루투스 배경 알림 검색 지역 및 언어 접근성 온라인 계정 개인 정보 공유 소리 전원 네트워크 장치

최소(C) 유선 5 적용(A)

자세히 보기 신원 IPv4 IPv6 보안

IPv4 방식 2 자동(DHCP) 수동 링크 로컬만 사용 않기

주소 3

주소	네트마스크	게이트웨이
172.16.119.5	255.255.255.0	172.16.119.2

네임서버(DNS) 4

자동 컴

172.16.119.2

라우팅

주소 네트마스크 게이트웨이 계측

자동 컴

## 4) 호스트 네임 설정

The image shows two overlapping windows. The top window is a terminal running as root on localhost. It shows the execution of `hostnamectl set-hostname master` and `hostname` command, which returns `master`. Then, the `vi /etc/hosts` command is executed. The bottom window is a VMware Workstation 12 Player showing a terminal window for a user named `hadoop` at `master:/home/hadoop`. This terminal shows the contents of the `/etc/hosts` file, which lists IP addresses and their corresponding hostnames: `127.0.0.1 localhost localhost.localdomain localhost4 localhost4.localdomain4` and `::1 localhost localhost.localdomain localhost6 localhost6.localdomain6`. Below these, three entries are added: `192.168.13.5 master`, `192.168.13.10 slave1`, and `192.168.13.11 slave2`. Two red starburst callouts are present: one pointing to the terminal commands with the text 'Root 계정 변경' (Change Root account), and another pointing to the VMware window with the text '사용자 System에 맞게 지정' (Specify according to user System).

```
root@localhost:~  
파일(F) 편집(E) 보기(V) 검색(S) 터미널(T) 도움말(H)  
[ root@localhost ~]# hostnamectl set-hostname master  
[ root@localhost ~]# hostname  
master  
[ root@localhost ~]# vi /etc/hosts  
[ root@localhost ~]#
```

Root 계정  
변경

사용자  
System에  
맞게 지정

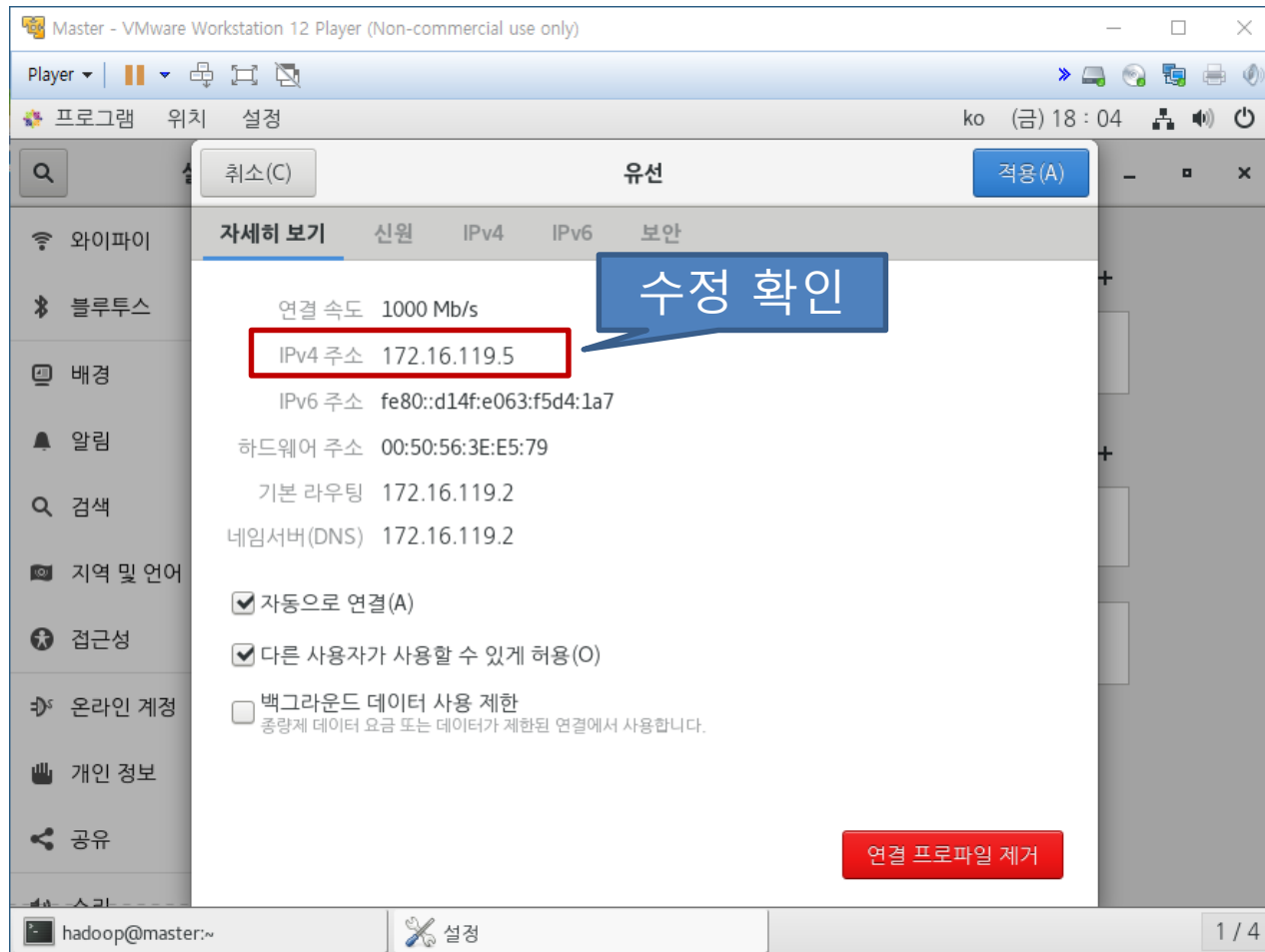
```
hadoop@master:/home/hadoop  
파일(F) 편집(E) 보기(V) 검색(S) 터미널(T) 도움말(H)  
127.0.0.1 localhost localhost.localdomain localhost4 localhost4.localdomain4  
::1 localhost localhost.localdomain localhost6 localhost6.localdomain6  
192.168.13.5 master  
192.168.13.10 slave1  
192.168.13.11 slave2  
~  
~
```

hadoop@master:/home/hadoop 1 / 4

# 인터넷 연결 확인



# 서버 재부팅 후 네트워크 설정 확인





# 5) Java 설치(jdk)

## 1) Java 설치 확인

✓ yum으로 JDK를 설치하기 위해서는 먼저 JRE 설치를 확인 한다.

```
hadoop@nameserver1:/home/hadoop
파일(F) 편집(E) 보기(V) 검색(S) 터미널(T) 도움말(H)
[root@nameserver1 hadoop]# yum list java*jdk-devel
Loaded plugins: fastestmirror, langpacks
Loading mirror speeds from cached hostfile
* base: centos.mirror.cdnetworks.com
* epel: mirror.premi.st
* extras: centos.mirror.cdnetworks.com
* remi-safe: mirror.smartmedia.net.id
* updates: centos.mirror.cdnetworks.com
Installed Packages
java-1.8.0-openjdk-devel.x86_64      1:1.8.0.111-2.b15.el7_3      @updates
Available Packages
java-1.6.0-openjdk-devel.x86_64      1:1.6.0.41-1.13.13.1.el7_3    updates
java-1.7.0-openjdk-devel.x86_64      1:1.7.0.121-2.6.8.0.el7_3     updates
java-1.8.0-openjdk-devel.i686        1:1.8.0.121-0.b13.el7_3       updates
java-1.8.0-openjdk-devel.x86_64      1:1.8.0.121-0.b13.el7_3       updates
[root@nameserver1 hadoop]#
```

현재 1.6, 1.7, 1.8 버전 설치가 가능하다. 여기서는 1.8 버전을 설치한다.

## 2) JDK 설치

### JDK 설치

```
hadoop@nameserver1:/home/hadoop
파일(F) 편집(E) 보기(V) 검색(S) 터미널(T) 도움말(H)
[root@nameserver1 hadoop]# yum install java-1.8.0-openjdk-devel.x86_64
Loaded plugins: fastestmirror, langpacks
Loading mirror speeds from cached hostfile
* base: centos.mirror.cdnetworks.com
* epel: mirror.premi.st
* extras: centos.mirror.cdnetworks.com
* remi-safe: mirror.smartmedia.net.id
* updates: centos.mirror.cdnetworks.com
Resolving Dependencies
--> Running transaction check
---> Package java-1.8.0-openjdk-devel.x86_64 1:1.8.0.111-2.b15.el7_3 will be updated
---> Package java-1.8.0-openjdk-devel.x86_64 1:1.8.0.121-0.b13.el7_3 will be an update
--> Processing Dependency: java-1.8.0-openjdk = 1:1.8.0.121-0.b13.el7_3 for package: 1:java-1.8.0-openjdk-devel-1.8.0.121-0.b13.el7_3.x86_64
--> Running transaction check
---> Package java-1.8.0-openjdk.x86_64 1:1.8.0.111-2.b15.el7_3 will be updated
---> Package java-1.8.0-openjdk.x86_64 1:1.8.0.121-0.b13.el7_3 will be an update
--> Processing Dependency: java-1.8.0-openjdk-headless = 1:1.8.0.121-0.b13.el7_3 for package: 1:java-1.8.0-openjdk-1.8.0.121-0.b13.el7_3.x86_64
--> Running transaction check
---> Package java-1.8.0-openjdk-headless.x86_64 1:1.8.0.111-2.b15.el7_3 will be updated
---> Package java-1.8.0-openjdk-headless.x86_64 1:1.8.0.121-0.b13.el7_3 will be an update
--> Finished Dependency Resolution

Dependencies Resolved
```

yum : application 설치 시 의존관계를 고려하여 설치해준다.

## JDK설치 정보 제공화면

```
hadoop@nameserver1:/home/hadoop
파일(F) 편집(E) 보기(V) 검색(S) 터미널(T) 도움말(H)

=====
Package                                Arch      Version                                Repository  Size
=====
Updating:
java-1.8.0-openjdk-devel             x86_64    1:1.8.0.121-0.b13.el7_3              updates     9.7 M
Updating for dependencies:
java-1.8.0-openjdk                   x86_64    1:1.8.0.121-0.b13.el7_3              updates     232 k
java-1.8.0-openjdk-headless          x86_64    1:1.8.0.121-0.b13.el7_3              updates     31 M

Transaction Summary
=====
Upgrade 1 Package (+2 Dependent packages)

Total size: 41 M
Is this ok [y/d/N]: y
```

### 3) JDK 설치 확인



The image shows a terminal window titled 'root@localhost:~'. The terminal has a menu bar with '파일(F)', '편집(E)', '보기(V)', '검색(S)', '터미널(T)', and '도움말(H)'. Two blue callout boxes are present: 'JDK 설치 확인' pointing to the first command, and 'JRE 버전 확인' pointing to the second command. The terminal output shows the results of these commands.

```
root@localhost:~  
[root@localhost ~]# rpm -qa java*jdk-devel  
java-1.8.0-openjdk-devel-1.8.0.252.b09-2.el7_8.x86_64  
[root@localhost ~]#  
[root@localhost ~]# java -version  
openjdk version "1.8.0_252"  
OpenJDK Runtime Environment (build 1.8.0_252-b09)  
OpenJDK 64-Bit Server VM (build 25.252-b09, mixed mode)  
[root@localhost ~]#  
[root@localhost ~]#
```

현재 JDK1.8.0\_252 최신 버전이 설치되었다.  
JRE는 1.8.0\_252 최신 버전으로 upgrade 되었다.  
➤ **최신 버전은 변경될 수 있음**

## 6) Slave 폴더 생성(Master 폴더 복사)

시스템 종료  
후 작업

