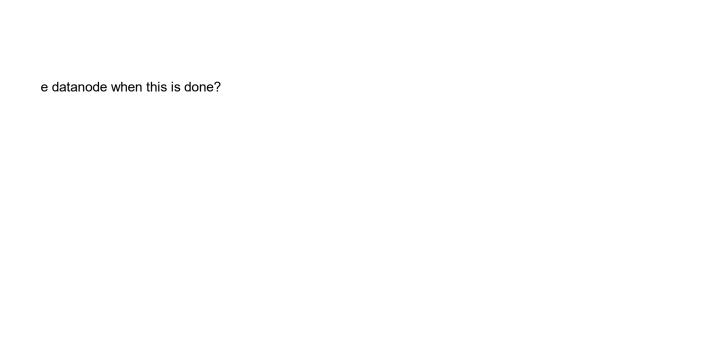
Step	Commands	Done (Y/N)?
set VM		
Update OS on newly created VM	sudo apt-get update	
check if java exists	java -version	
install jdk	sudo apt-get install default-jdk	
•	sudo apt-get install default-juk	
install jre		
create separate users & groups	sudo groupadd hdp1313	
	sudo adduser hdp1313 -ingroup hdp1313	
	sudo adduser hdfs1313 -ingroup hdp1313	
	sudo usermod -aG hdp1313 root	
	sudo adduser mprd1313 -ingroup hdp1313	
add to sudo group	sudo usermod -aG sudo hdfs1313	
	sudo usermod -aG sudo hdp1313	
	sudo usermod -aG sudo mprd1313	
update /etc/security/limits.conf	vi /etc/security/limits.conf	
check kernel parameters	sysctl -a	
check kerner parameters		
	sudo vi /etc/sysctl.conf	
	vm.swappiness=0	
	vm.overcommit_memory=1	
	vm.overcommit_ratio=50	
	sudo sysctl -p	
check disk format	df -Th	
check if ssh is installed	which ssh	
install if not installed	sudo apt-get install openssh-server	
	sudo ss -lnp grep sshd	
	ssh-keygen -t rsa -b 2048	
	cat id_rsa.pub >> authorized_keys	
test ssh to localhost	ssh root@localhost	
	wget https://www-	
download hadoop	us.apache.org/dist/hadoop/common/hadoop-	
unpack	2.7.7/hadoop-2.7.7.tar.gz tar -zxf hadoop-2.7.7.tar.gz	
шраск	sudo chown -R root:root hadoop-2.7.7	
	sudo In -s hadoop-2.7.7 hadoop	
find where is java	update-alternativesconfig java	
set variables in .bashrc	export HADOOP_HOME=/usr/local/hadoop	
	export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop	
	export HADOOP_COMMON_HOME=/usr/local/hadoop	
	export HADOOP_HDFS_HOME=/usr/local/hadoop	
	export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64	(A LIONAT! ""
	export CLASSPATH=\$CLASSPATH:\$JAVA_HOME/lib:\$JAV	
	export PATH=\$JAVA_HOME/bin:\$JAVA_HOME/jre/bin:\$PATH=\$PATH=\$JAVA_HOME/jre/bin:\$PATH=\$	тт.фПАDOOP_Н
source bashrc	. ~/.bashrc	
test - check versions	java -version	
	hadoop version	
create data directories	mkdir -m 775 /home/hdfs1313/data	

	mkdir -m 775 -p /hdata/	'data	
	chown -R hdfs1313:hdp	o1313 /hdata	
	chmod -R 775 /hdata		
add in core-site.xml	<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	<name>hadoop.tmp.dir</name>	<value>/</value>
add in mapred-site.xml	<pre><pre><pre><pre>property><name>map</name></pre></pre></pre></pre>	oreduce.framework.name <val< td=""><td>lue>yarn</td></val<>	lue>yarn
add in hdfs-site.xml	<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	<name>dfs.replication</name>	<value>1<</value>
add in yarn-site.xml	<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	<name>yarn.nodemanager.aux-se</name>	rvices
add in hadoop-env.sh	export JAVA_HOME=/u	usr/lib/jvm/java-8-openjdk-amd64	
format namenode	hdfs namenode -format		
start hadoop	start-dfs.sh		
	jps		
	start-yarn.sh		
	jps		
Configure users	hdfs dfs -chmod -R 777	/tmp	
	hdfs dfs -mkdir -p /user,	/hdp1313	
	hdfs dfs -chown -R hdp	1313:hdp1313 /user/hdp1313	
	hdfs dfsadmin -refreshl	JserToGroupMappings	
Prepare for a test job			
download initial data	mkdir -m 755 ~/scripts		
create hdfs source directory	hdfs dfs -mkdir /userdat		
	hdfs dfs -mkdir /userdat		
copy source file onto hdfs		dp1313/src/iliad.csv /userdata/input	
execute test job	hadoop jar /usr/local/ha	doop/share/hadoop/mapreduce/hadoc	p*examples*.jar v
Discuss observations			
rerun with a variation	file="/home/hdp1313/sr	c/iliad.csv" curl https://www.gutenberg	.org/ebooks/6130
execute test job	hadoop jar /usr/local/ha	doop/share/hadoop/mapreduce/hadoo	p*examples*.jar v
rerun after fixing		doop/share/hadoop/mapreduce/hadoc	

References	Comments
	refer specifications in presentation. Use GUI
a-compliant-cloud-	must be same version on all VMs, down to the micro releasethis is defa
<u>a compliant cloud</u>	index be carrie version an vivie, devin to the finere released ne le dete
	as root; add "@hdp1313 hard nofiles 32768", no quotes
	as root
	as 100t
	shock ash norts if already installed and if required
	check ssh ports if already installed and if required generate keys under /root/.ssh. If running as hadoop user, do it with had
	in .ssh
	ensure works passwordless
	Cristic Works passwordioss
	From under /usr/local
	in /usr/local
	use the value received from command in line 30. Use path until before
)ME/bin:\$HADOOP_HO	ML/SDIN:.
	7 (notice the cub release numbers)
	7 (notice the sub release numbers) 2.7.7.
	Z.1.1.
	, and the state of

nome/hdfs1313/data	/usr/local/hadoop/etc/hadoop	
	/usr/local/hadoop/etc/hadoop	
'value>	/usr/local/hadoop/etc/hadoop	
<pre><value>mapreduce_s</value></pre>	/usr/local/hadoop/etc/hadoop	
export JAVA_HOME=\${JAVA	what happens if you let it stay at default	
check return code, understar	What happens to the namenode when this is done?W	/hat happens to th
check output		
check processes		
check output		
check processes	what is the difference from previous output	
	run as root	
	in your home directory	
	save as big_file.sh under scripts dir	
	as root	
vordcount /userdata/input /us		
	discuss output; save it to a file for later reference	
.txt.utf-8 > \$file	what is the variation here?	
	what happens when this is run? Why does this	
	happen like this? how to fix this?	
vordcount /userdata/input /us	compare output with previous output	
	discuss changes, and relate them to theory	





KlhTxx1209\$n

		Done (Y/N)?
Step	Commands	pri-node
set VM		
edit /etc/hosts file	which are to be part of hadoop cluster	
Update OS on newly created VM	sudo apt-get update	
check if java exists	java -version	
install jdk	sudo apt-get install default-jdk	
install jre	sudo apt-get install default-jre	
create separate users & groups	sudo groupadd hdp1313	
	sudo adduser hdp1313 -ingroup hdp1313	
	sudo adduser hdfs1313 -ingroup hdp1313	
	sudo usermod -aG hdp1313 root	
	sudo adduser mprd1313 -ingroup hdp1313	
add to sudo group	sudo usermod -aG sudo hdfs1313	
add to sade group	sudo usermod -aG sudo hdp1313	
	sudo usermod -aG sudo mprd1313	
undata /ata/aagurity/limita aanf	vi /etc/security/limits.conf	
update /etc/security/limits.conf	•	
check kernel parameters	sysctl -a	
	sudo vi /etc/sysctl.conf	
	vm.swappiness=0	
	vm.overcommit_memory=1	
	vm.overcommit_ratio=50	
	sudo sysctl -p	
check if ssh is installed	which ssh	
install if not installed	sudo apt-get install openssh-server sudo ss -Inp grep sshd	
	ssh-keygen -t rsa -b 2048	
	cat id_rsa.pub >> authorized_keys	
	test ssh to self on localhost	
copy ssh keys from pri-node to others		
test ssh from pri-node to d-node-a		
copy ssh keys from d-node-a to pri-node		NA
test ssh from d-node-a to pri-node check disk format	ldf -Th	NA
CHECK CISK TOTTIAL		
	wget https://www-	
	<pre>us.apache.org/dist/hadoop/commo n/hadoop-2.7.7/hadoop-</pre>	
download hadoop	2.7.7.tar.qz	
unpack	tar -zxf hadoop-2.7.7.tar.gz	
'	sudo mv hadoop-2.7.7 /usr/local	
	sudo chown -R root:root hadoop-2.7.7/	
	sudo In -s hadoop-2.7.7 hadoop	
Construction to the construction of the constr	update-alternativesconfig	
find where is java	java	
set variables in .bashrc	export HADOOP_HOME=/usr/local/hadoop	
	export HADOOP_CONF_DIR=/usr/local/ha export HADOOP_COMMON_HOME=/usr/	
	export HADOOP _COMMON_HOME=/usr/local	
		·

	export JAVA_HOME=/usr/lib/jvm/java-7-o	penjdk-amd64
	export CLASSPATH=\$CLASSPATH:\$JA\	/A_HOME/lib:\$JA\
	export PATH=\$JAVA_HOME/bin:\$JAVA_	HOME/jre/bin:\$PA
	export HADOOP_VERSION=2.7.7	
test - check versions	java -version	
	hadoop version	
	vi	
	/usr/local/hadoop/etc/hadoop/ma	
create masters file	sters	
	vi	
	/usr/local/hadoop/etc/hadoop/sl	
update slaves file	aves	
create data directories	mkdir -m 775 /home/hdfs/data	
create data directories	chown hdfs:hadoop /home/hdfs/data	
	mkdir -m 775 -p /hdata/data	
	chown -R hdfs:hadoop /hdata chmod -R 775 /hdata	
add in core-site.xml		Annua alim almana as
		o.tmp.dir
add in mapred-site.xml	<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	
add in hdfs-site.xml		<name>dfs.replica</name>
update yarn-site.xml on master		demanager.aux-s
update yarn-site.xml on slave	<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	demanager.aux-se
add in hadoop-env.sh	export JAVA_HOME=/usr/lib/jvm/java-7-o	penjdk-amd64
format namenode	hdfs namenode -format	
start hadoop	start-dfs.sh	
	jps	
	<u>start-yarn.sh</u>	
	jps	
Prepare for a test job		
	execute big_file.sh	
create hdfs source directory	hdfs dfs -mkdir /userdata	
	hdfs dfs -mkdir /userdata/input	
copy source file onto hdfs	hdfs dfs -put -f ~/src/iliad.csv /userdata/ing	out
execute test job	hadoop jar /usr/local/hadoop/share/hadoop	o/mapreduce/hado
rerun with a variation	file="/home/akashahuja575/src/iliad.csv" o	curl https://www.gu
evecute test ish		- / /
execute test 100	Inagood jar /usr/jocai/nagood/snare/nagood	o/mapreduce/nadd
execute test job rerun after fixing	hadoop jar /usr/local/hadoop/share/hadoop hadoop jar /usr/local/hadoop/share/hadoop	

(Y/N)?			
d-node-a	References	Comments	
		refer specifications in presentation. Use GUI	
	-		
	a-compliant-cloud-	must be same version on all VMs, down to the micro releasethis is	defa
	-		
	-		
		4 11 11 01 1 1 5 1 00 700 1	
		as root; add "@hadoop hard nofiles 32768", no quotes	
		as root	
		check ssh ports if already installed and if required	
		generate keys under /root/.ssh	
NA			
NA			
		as root	
		in /usr/local	
	+		

	1	1	
/A_HOME/	jre/lib		
	OOP_HOME/bin		
		7 (notice the sub release numbers)	
		2.7.7.	
		final content must be:pri-node (or same dns as entere	ed in /etc/hosts for
		final content must be:pri-noded-node-a (or same dns	as entered in /etc/
			u
		Observe the dfs.replication.factor is now changed to	the number of nod
tion <td>/value></td> <td>Observe the discreption factor is now shanged to</td> <td>the number of ped</td>	/value>	Observe the discreption factor is now shanged to	the number of ped
		Observe the dfs.replication.factor is now changed to	ine number of nod
NA	NA	##### is server name of master node	
ervices <td></td> <td>##### is server name of master node</td> <td></td>		##### is server name of master node	
NA		what happens if you let it stay at default What happens to the namenode when this is done?W	/hat hannana ta th
NA	check output	All discussed deamons shall run	mai nappens io in
INA	check processes	All discussed deallions shall full	
NA	check processes		
INA	check processes	Only datanode and nodemanager must run	
	Check processes	only datanode and nodemanager mast run	
	1		
NA		observe output	
		discuss output; save it to a file for later reference	
NA		what is the variation here?	
		what happens when this is run? Why does this	
NA		happen like this? how to fix this?	
NA		compare output with previous output	
		discuss changes, and relate them to theory	

ault JDK, check others for other versions. Additionally, cloudra RPM needs Oracle Java RPM due to dep	ende

master node)
hosts for master and slave node)
lesHow is this different than a single node install? Is that difference necessary?###### is server name of mas
les
e datanode when this is done?

