

在 Windows 上安装 Hadoop 教程

— 见 2010.1.6 www.hadoopor.com/hadoopor@foxmail.com

1. 安装 JDK

不建议只安装 JRE，而是建议直接安装 JDK，因为安装 JDK 时，可以同时安装 JRE。MapReduce 程序的编写和 Hadoop 的编译都依赖于 JDK，光 JRE 是不够的。

JRE 下载地址：http://www.java.com/zh_CN/download/manual.jsp

JDK 下载地址：<http://java.sun.com/javase/downloads/index.jsp>，下载 Java SE 即可。

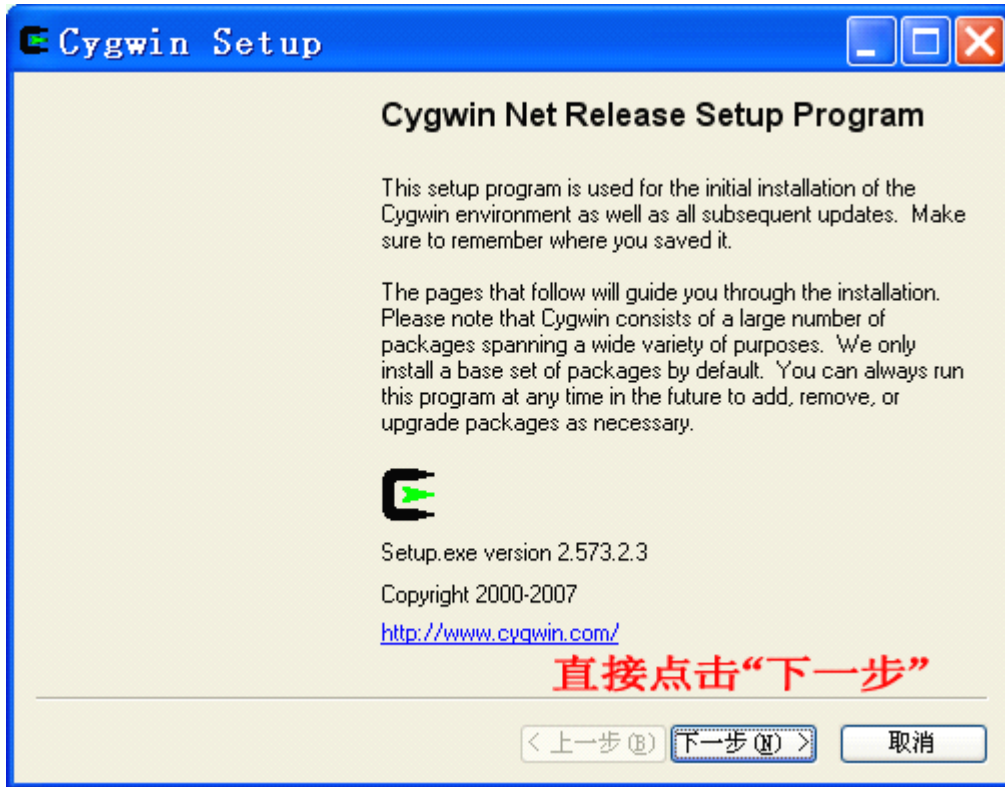
2. 安装 Cygwin

在安装 Cygwin 之前，得先下载 Cygwin 安装程序 setup.exe。

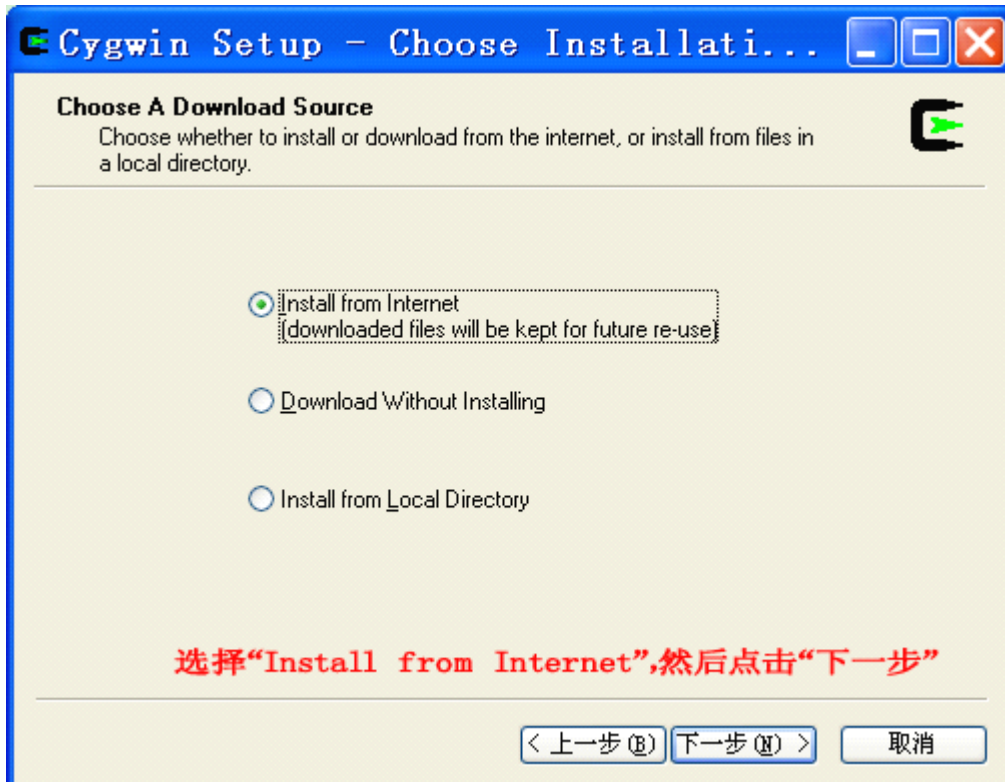
Cygwin 安装程序下载地址：<http://www.cygwin.com/setup.exe>，当然也可以从<http://www.cygwin.cn/setup.exe> 下载 Cygwin 安装程序，不过如果在安装过程中，遇到如下图所示的错误，则只能从<http://www.cygwin.com/setup.exe> 下载，本教程下载的是 Cygwin 1.7.1 版本。



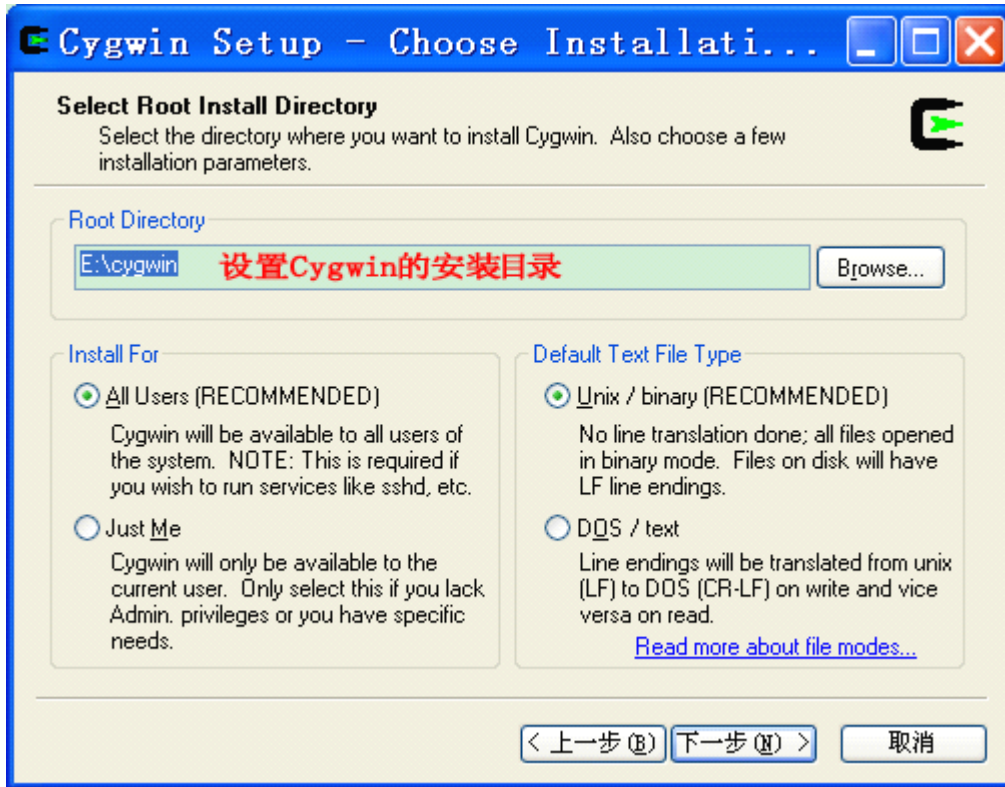
Cygwin 安装程序 setup.exe 的存放目录可随意无要求。当 setup.exe 下载成功后，运行 setup.exe，弹出如下图所示的对话框：



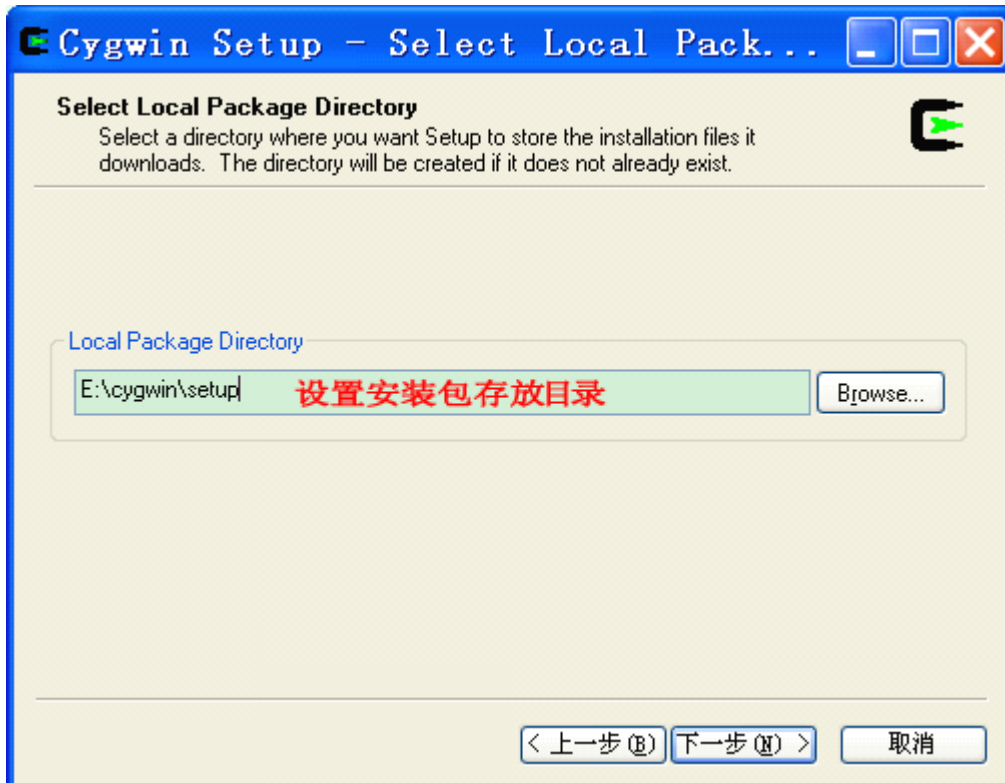
在上图所示的对话框中，直接点击“下一步”，进入如下图所示的对话框：



在上图所示的对话框中，选择“Install from Internet”，然后点击“下一步”，进入如下图所示对话框：



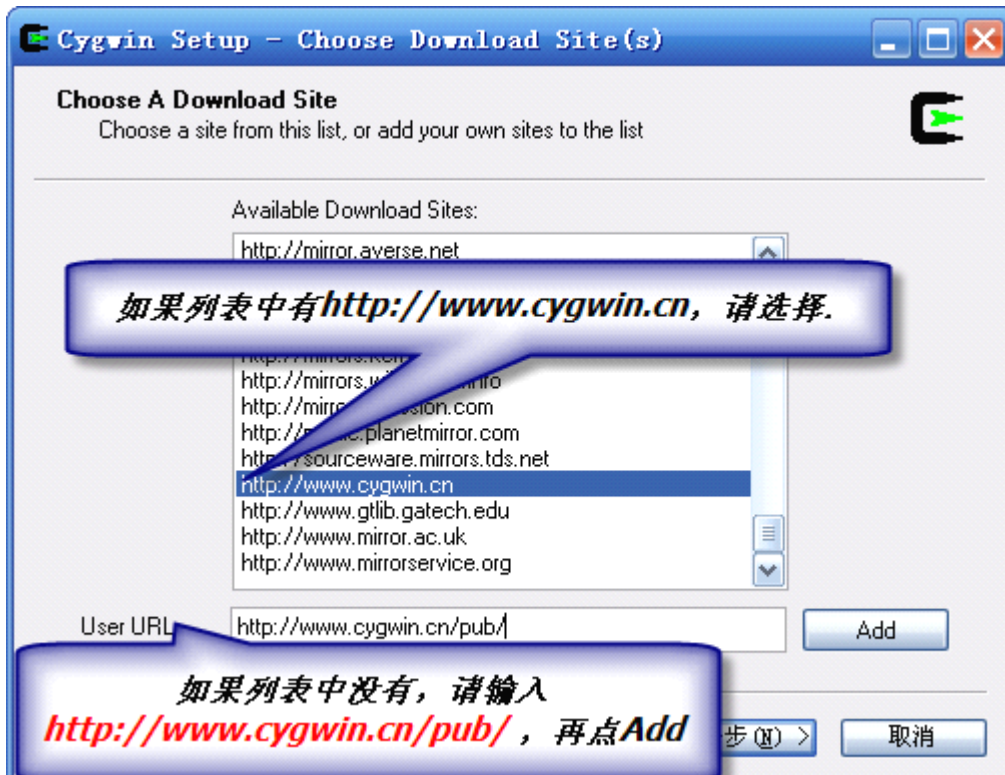
在上图所示的对话框中，设置 Cygwin 的安装目录，Install For 选择 “All Users”，Default Text File Type 选择 “Unix/binary”，然后点击 “下一步”，进入如下图所示对话框：



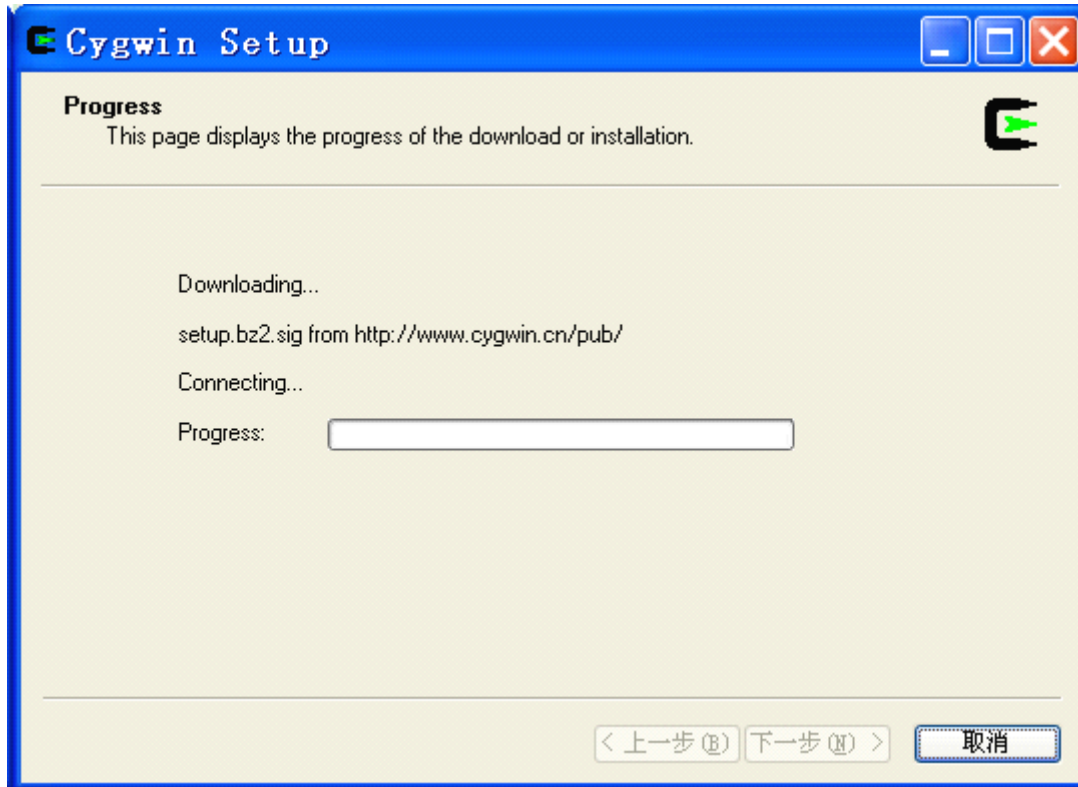
在上图所示的对话框中，设置 Cygwin 安装包存放目录，然后点击 “下一步”，进入如下图所示对话框：



在上图所示的对话框中，选择“Direct Connection”，然后点击“下一步”，进入如下图所示对话框：



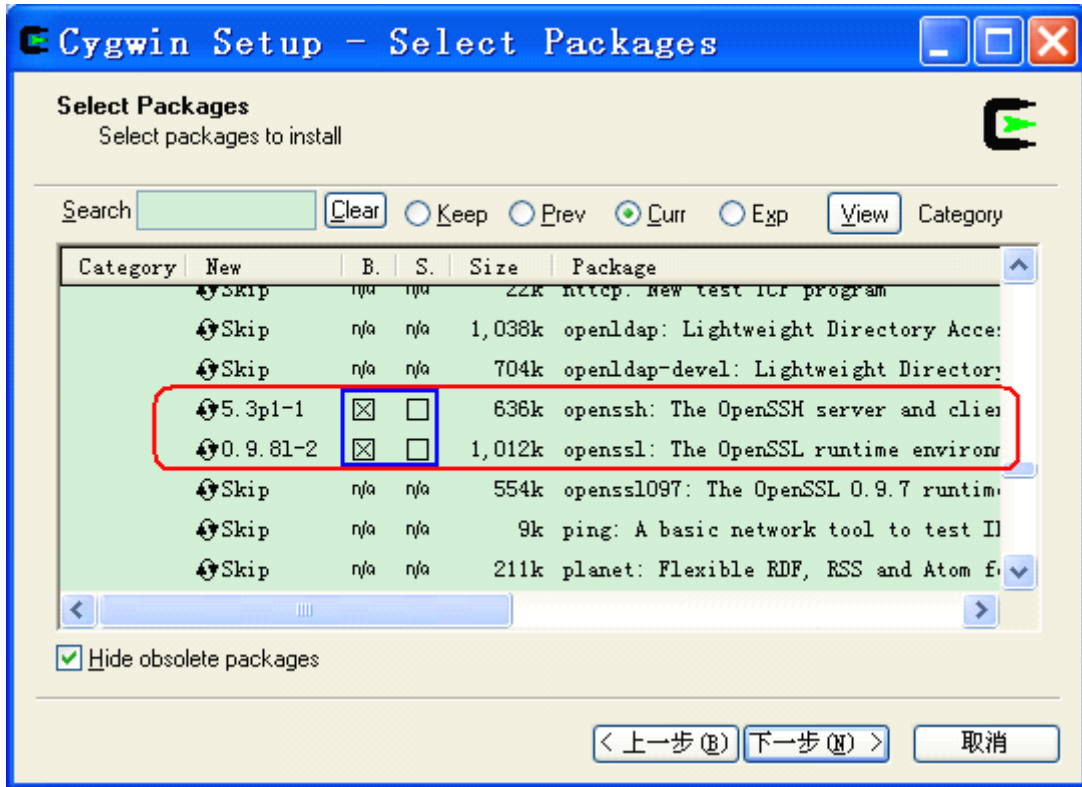
在上图所示的对话框中，点击“下一步”，将进入如下图所示的对话框：



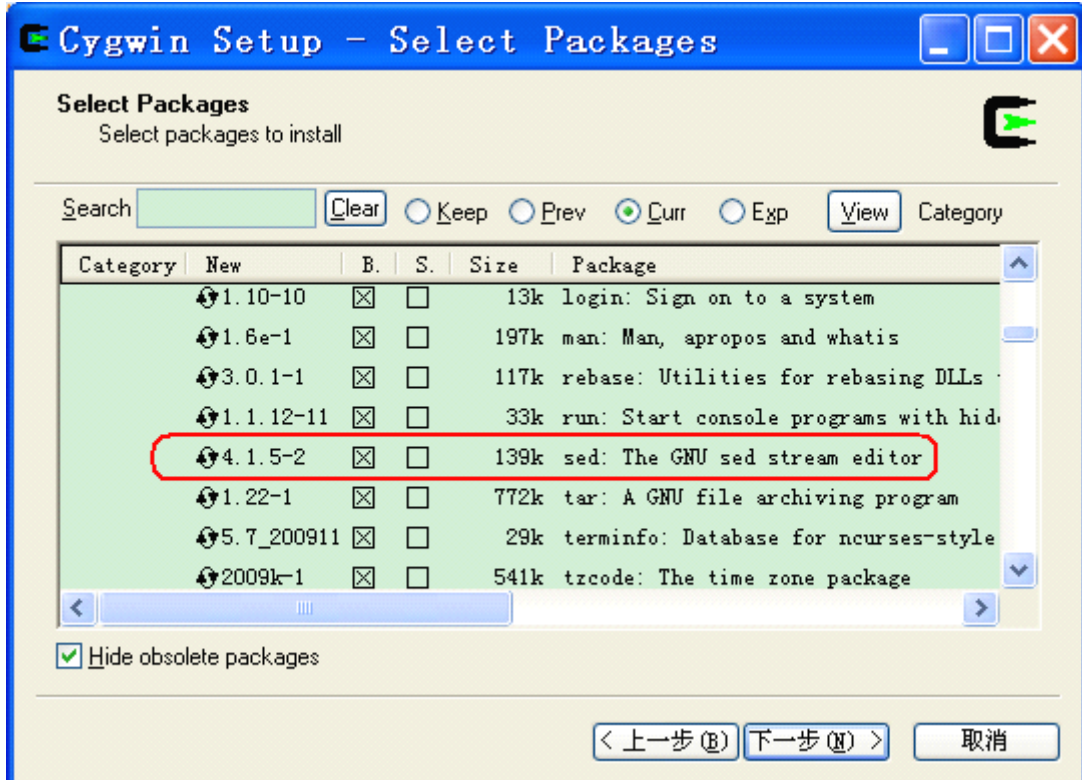
在上图所示的对话框过程中，可能会弹出如下图所示的“Setup Alert”对话框，直接点击“确定”即可。



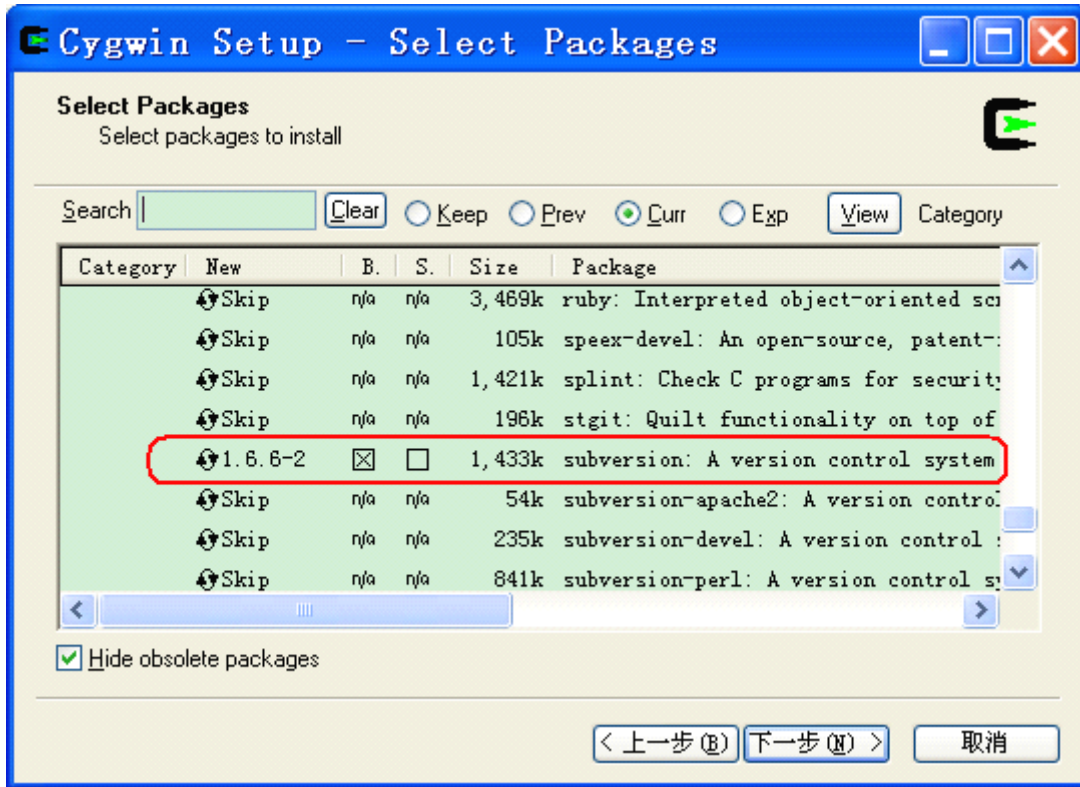
进入“Select Packages”对话框后，必须保证“Net Category”下的“OpenSSL”被安装，如下图所示：



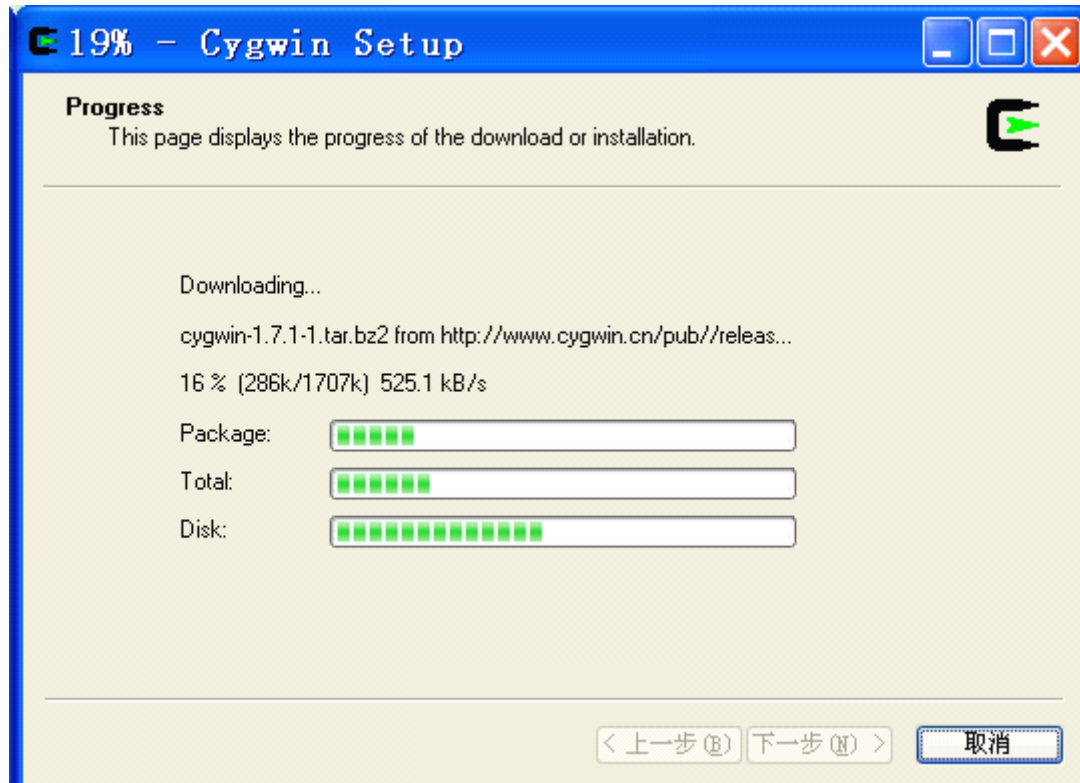
如果还打算在 eclipse 上编译 Hadoop，则还必须安装“Base Category”下的“sed”，如下图所示：



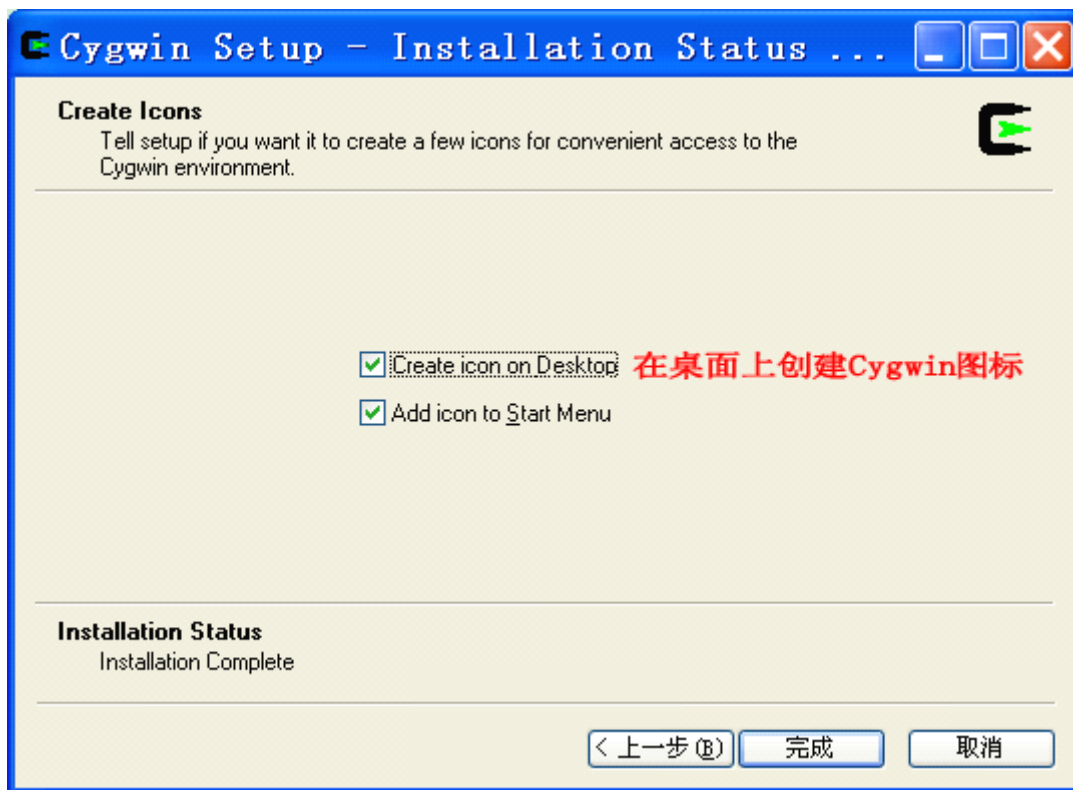
另外，还建议将“Editors Category”下的 vim 安装，以方便在 Cygwin 上直接修改配置文件；“Devel Category”下的 subversion 建议安装，如下图所示：



当完成上述操作后，点击“Select Packages”对话框中“下一步”，进入 Cygwin 安装包下载过程，如下图所示：



等待安装包下载完毕，当下载完后，会自动进入到如下图所示的对话框：

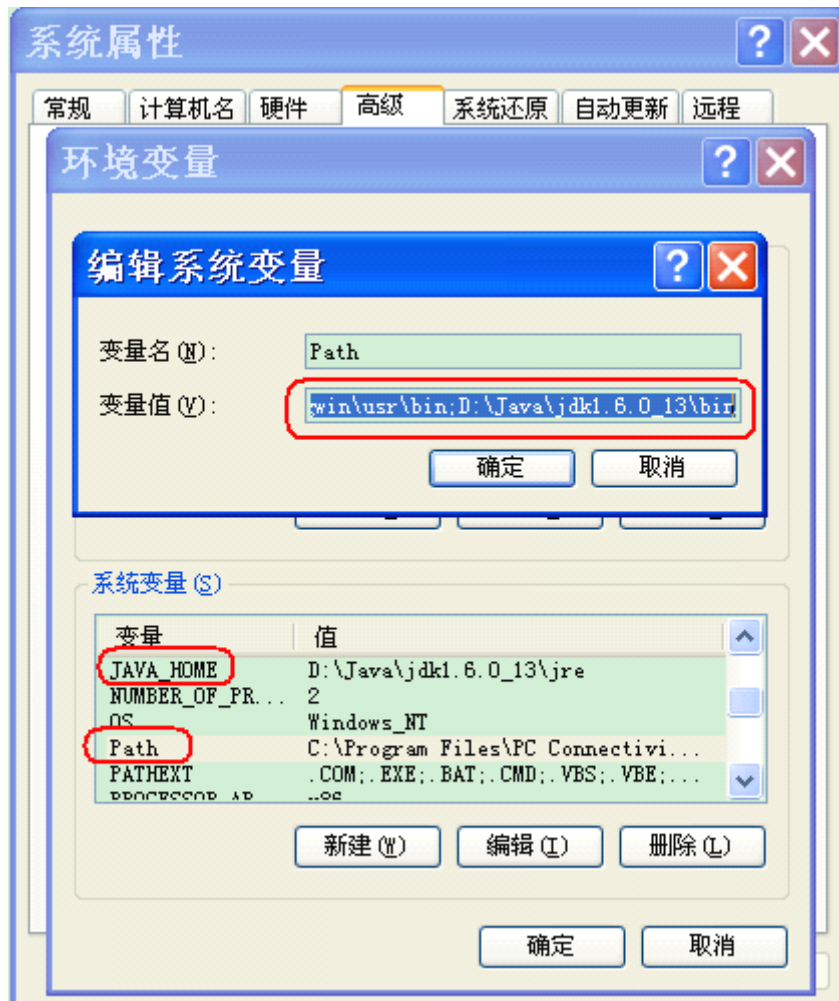


在上图所示的对话框中，选中“Create icon on Desktop”，以方便直接从桌面上启动 Cygwin，然后点击“完成”按钮。至此，Cygwin 已经安装完，安装目录下的内容如下图所示：



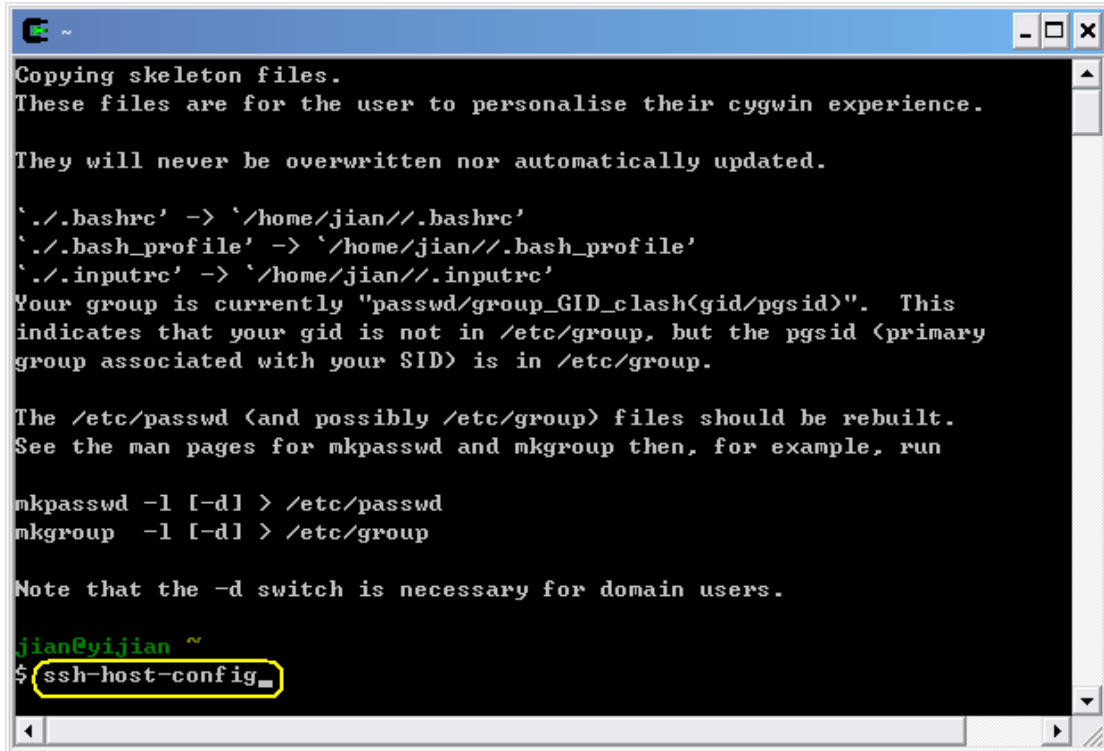
3. 配置环境变量

需要配置的环境变量包括 PATH 和 JAVA_HOME: JAVA_HOME 指向 JRE 安装目录; JDK 的 bin 目录, Cygwin 的 bin 目录, 以及 Cygwin 的 usr\bin 目录都必须添加到 PATH 环境变量中, 如下图所示:



4. 安装 sshd 服务

点击桌面上的 Cygwin 图标，启动 Cygwin，执行 **ssh-host-config** 命令，如下图所示：



```
Copying skeleton files.
These files are for the user to personalise their cygwin experience.

They will never be overwritten nor automatically updated.

`./bashrc' -> `/home/jian//.bashrc'
`./bash_profile' -> `/home/jian//.bash_profile'
`./inputrc' -> `/home/jian//.inputrc'
Your group is currently "passwd/group_GID_clash(gid/pgsid)". This
indicates that your gid is not in /etc/group, but the pgsid <primary
group associated with your SID> is in /etc/group.

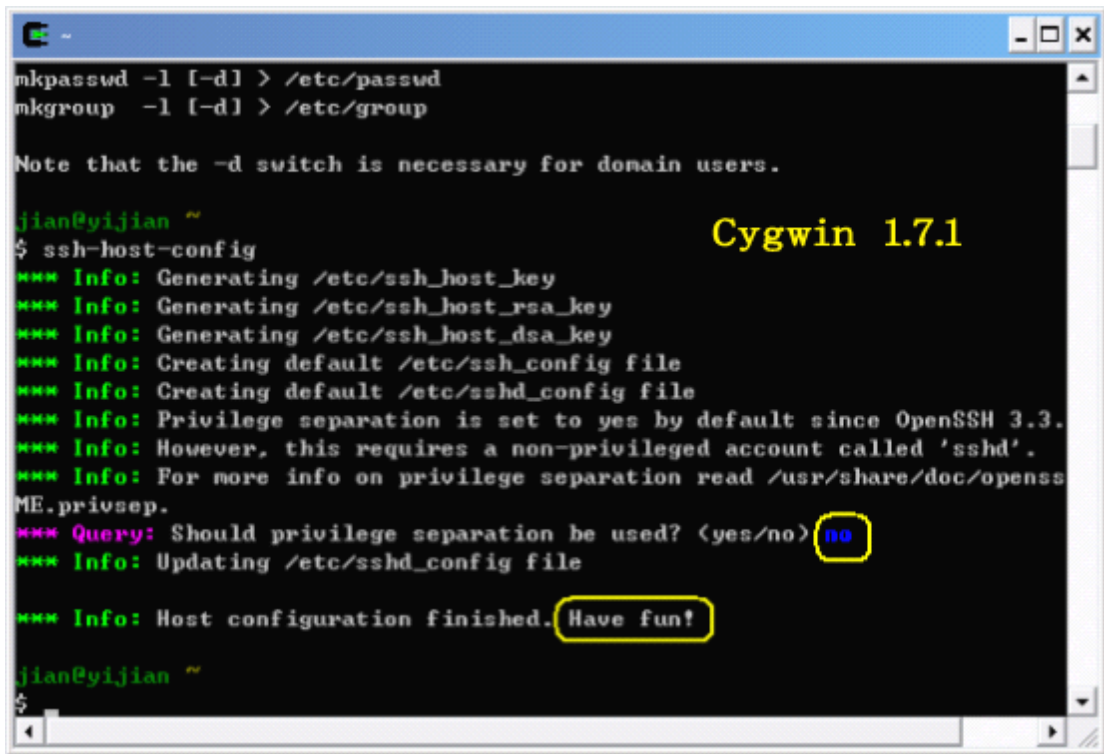
The /etc/passwd <and possibly /etc/group> files should be rebuilt.
See the man pages for mkpasswd and mkgroup then, for example, run

mkpasswd -l [-d] > /etc/passwd
mkgroup -l [-d] > /etc/group

Note that the -d switch is necessary for domain users.

jian@yijian ~
$ ssh-host-config
```

在执行 ssh-host-config 时，当要求输入 yes/no 时，选择输入 no，如下图所示：



```
mkpasswd -l [-d] > /etc/passwd
mkgroup -l [-d] > /etc/group

Note that the -d switch is necessary for domain users.

jian@yijian ~
$ ssh-host-config
                                     Cygwin 1.7.1
*** Info: Generating /etc/ssh_host_key
*** Info: Generating /etc/ssh_host_rsa_key
*** Info: Generating /etc/ssh_host_dsa_key
*** Info: Creating default /etc/ssh_config file
*** Info: Creating default /etc/sshd_config file
*** Info: Privilege separation is set to yes by default since OpenSSH 3.3.
*** Info: However, this requires a non-privileged account called 'sshd'.
*** Info: For more info on privilege separation read /usr/share/doc/openssh
ME.privsep.
*** Query: Should privilege separation be used? <yes/no> no
*** Info: Updating /etc/sshd_config file

*** Info: Host configuration finished. Have fun!

jian@yijian ~
$
```

如果是 Cygwin 1.7 之前的版本，则 ssh-host-config 显示界面如下图所示：

```

$ ssh-host-config
Generating /etc/ssh_host_key
Generating /etc/ssh_host_rsa_key
Generating /etc/ssh_host_dsa_key
Generating /etc/ssh_config file
Privilege separation is set to yes by default since OpenSSH 3.3.
However, this requires a non-privileged account called 'sshd'.
For more info on privilege separation read /usr/share/doc/openssh/R
.

Should privilege separation be used? <yes/no> no
Generating /etc/sshd_config file
Added ssh to C:\WINDOWS\system32\drivers\etc\services

Warning: The following functions require administrator privileges!

Do you want to install sshd as service?
<Say "no" if it's already installed as service> <yes/no> yes

Which value should the environment variable CYGWIN have when
sshd starts? It's recommended to set at least "ntsec" to be
able to change user context without password.
Default is "ntsec". CYGWIN=ntsec

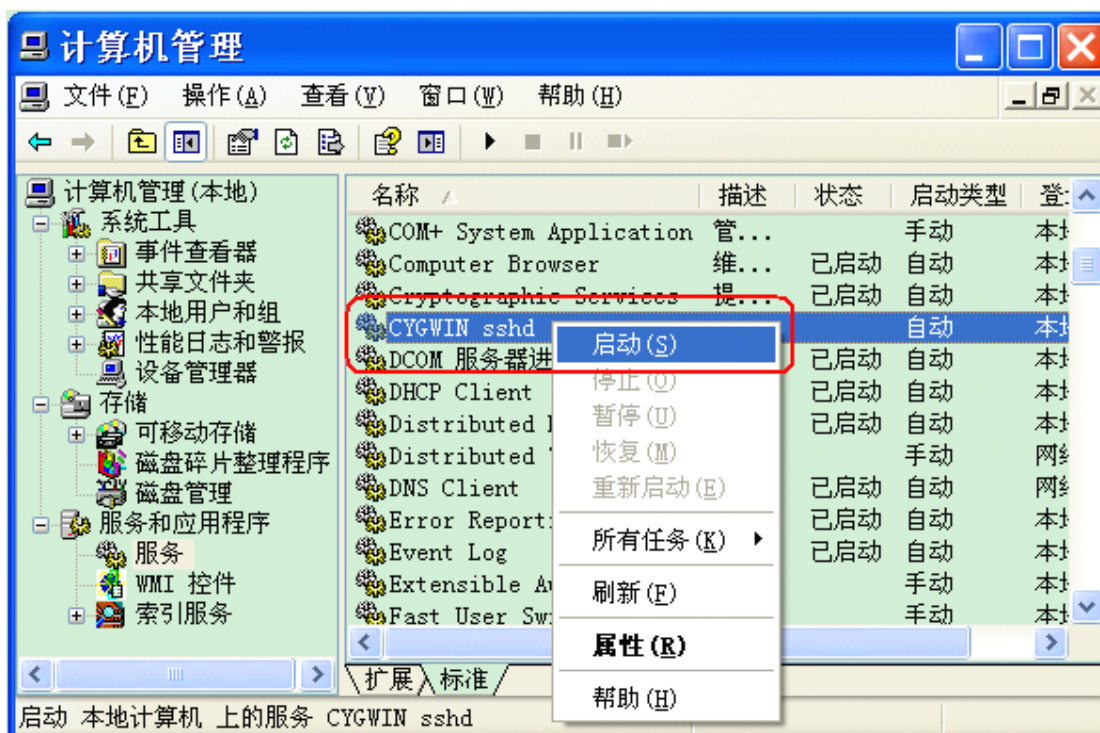
The service has been installed under LocalSystem account.
To start the service, call 'net start sshd' or 'cygrunsrv -S sshd'.

Host configuration finished. Have fun!
    
```

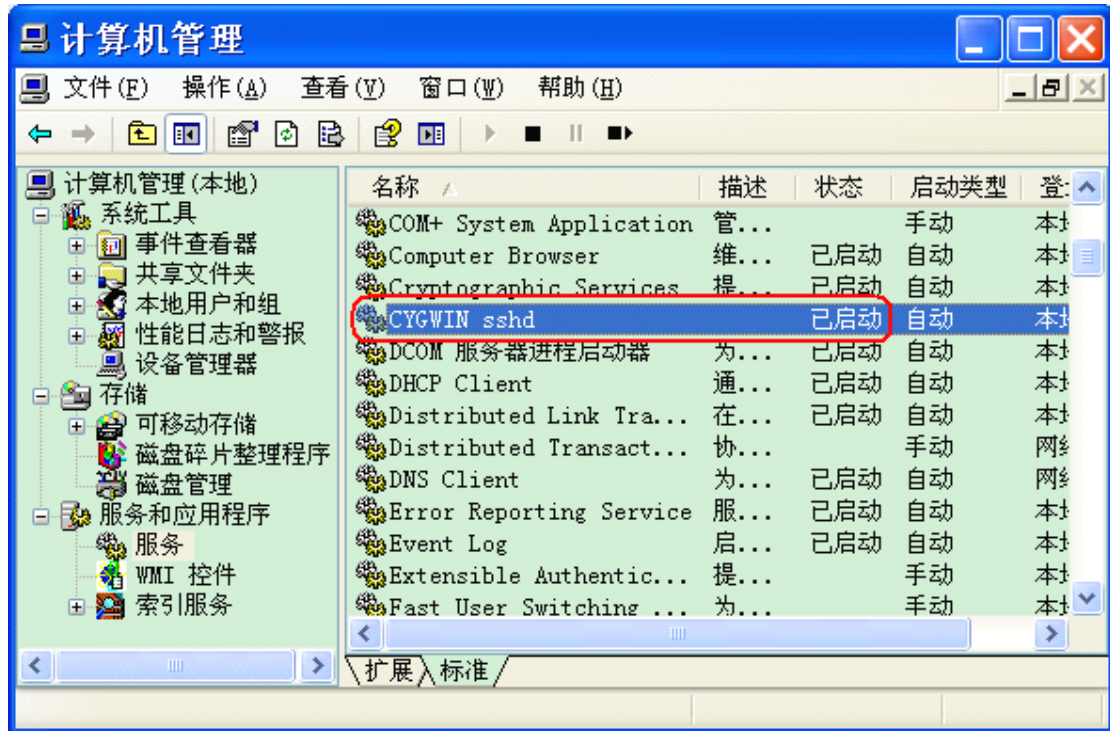
当看到“Have fun”时，一般表示 sshd 服务安装成功了，如上图所示。接下来，需要启动 sshd 服务。

5. 启动 sshd 服务

在桌面上的“我的电脑”图标上单击右键，点击“管理”菜单，进入 Windows 计算机管理，如下图所示：



在上图所示的对话框中，选中“**CYGWIN sshd**”，弹出右键，并启动 CYGWIN sshd 服务，成功后，如下图所示：



当 CYGWIN sshd 的状态为“已启动”后，接下来就是配置 ssh 登录。

6. 配置 ssh 登录

执行 **ssh-keygen** 命令生成密钥文件，如下图所示：

```
~/ssh
jian@yijian ~/ssh
$ ssh-keygen
Generating public/private rsa key pair.
Enter file in which to save the key (/home/jian/.ssh/id_rsa): 直接回车
Enter passphrase (empty for no passphrase): 这里也直接回车
Enter same passphrase again: 这里也直接回车
Your identification has been saved in /home/jian/.ssh/id_rsa.
Your public key has been saved in /home/jian/.ssh/id_rsa.pub.
The key fingerprint is:
e1:a3:9d:4e:b5:3b:95:f0:66:1c:1c:10:58:ca:41:e2 jian@yijian
The key's randomart image is:
+--[ RSA 2048 ]-----+
|      ..oo+o      |
|     - o.o .     |
|      E + . .     |
|      . . . o     |
|      $ .+ o     |
|      o + .B     |
|      . + .+     |
|      o ..       |
|      . ..       |
+-----+
$
```

在上图所示对话框中，需要输入时，直接按回车键即可，如果不出错，应当是需要三次按回车键。接下来生成 **authorized_keys** 文件，按下图所示操作即可：

```
~/ssh
jian@yijian ~/ssh
$ cd ~/.ssh/
jian@yijian ~/ssh
$ ls
id_rsa id_rsa.pub
jian@yijian ~/ssh
$ cp id_rsa.pub authorized_keys
jian@yijian ~/ssh
$ ls
authorized_keys id_rsa id_rsa.pub
jian@yijian ~/ssh
$
```

正如下图所示，只需要两步操作，即可生成 **authorized_keys** 文件：

```
cd ~/.ssh/
```

```
cp id_rsa.pub authorized_keys
```

完成上述操作后，执行 **exit** 命令先退出 Cygwin 窗口，如果不执行这一步操作，下面的操作可能会遇到错误。接下来，重新运行 Cygwin，执行 **ssh localhost** 命令，在第一次执行 ssh localhost 时，会有如下图所示的提示，输入 **yes**，然后回车即可：

```

~/ssh
Your group is currently "passwd/group_GID_clash(gid/pgsid)". This
indicates that your gid is not in /etc/group, but the pgsid (primary
group associated with your SID) is in /etc/group.

The /etc/passwd (and possibly /etc/group) files should be rebuilt.
See the man pages for mkpasswd and mkgroup then, for example, run

mkpasswd -l [-dl] > /etc/passwd
mkgroup -l [-dl] > /etc/group

Note that the -d switch is necessary for domain users.

jian@vijian ~
$ ssh localhost
The authenticity of host 'localhost (127.0.0.1)' can't be establish
RSA key fingerprint is d7:0f:a4:56:be:43:15:9c:f2:02:ac:24:62:5a:ac:
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'localhost' (RSA) to the list of known ho
Last login: Wed Jan  6 22:25:59 2010 from localhost
$
    
```

如果是 Windows 域用户，这步操作可能会遇到问题，错误信息如下：

```

$ ssh localhost
Last login: Thu Jan  7 10:23:09 2010 from localhost
3 [main] -bash 5016 E:\cygwin\bin\bash.exe: *** fatal e
amically determine load address for 'WSAGetLastError' (handle
error 126
Connection to localhost closed.
    
```

这个错误暂无解决办法，问题的解决情况，可关注 Hadoop 技术论坛中的贴：<http://bbs.hadoopor.com/thread-348-1-1.html>(Cygwin 1.7.1 版本 ssh 问题)。否则，如果成功，执行 who 命令时，可以看到如下图所示的信息：

```

$ who
jian      tty0      2010-01-06 22:51 (localhost)
    
```

至此，配置 ssh 登录成功，下面就可以开始安装 hadoop 了。

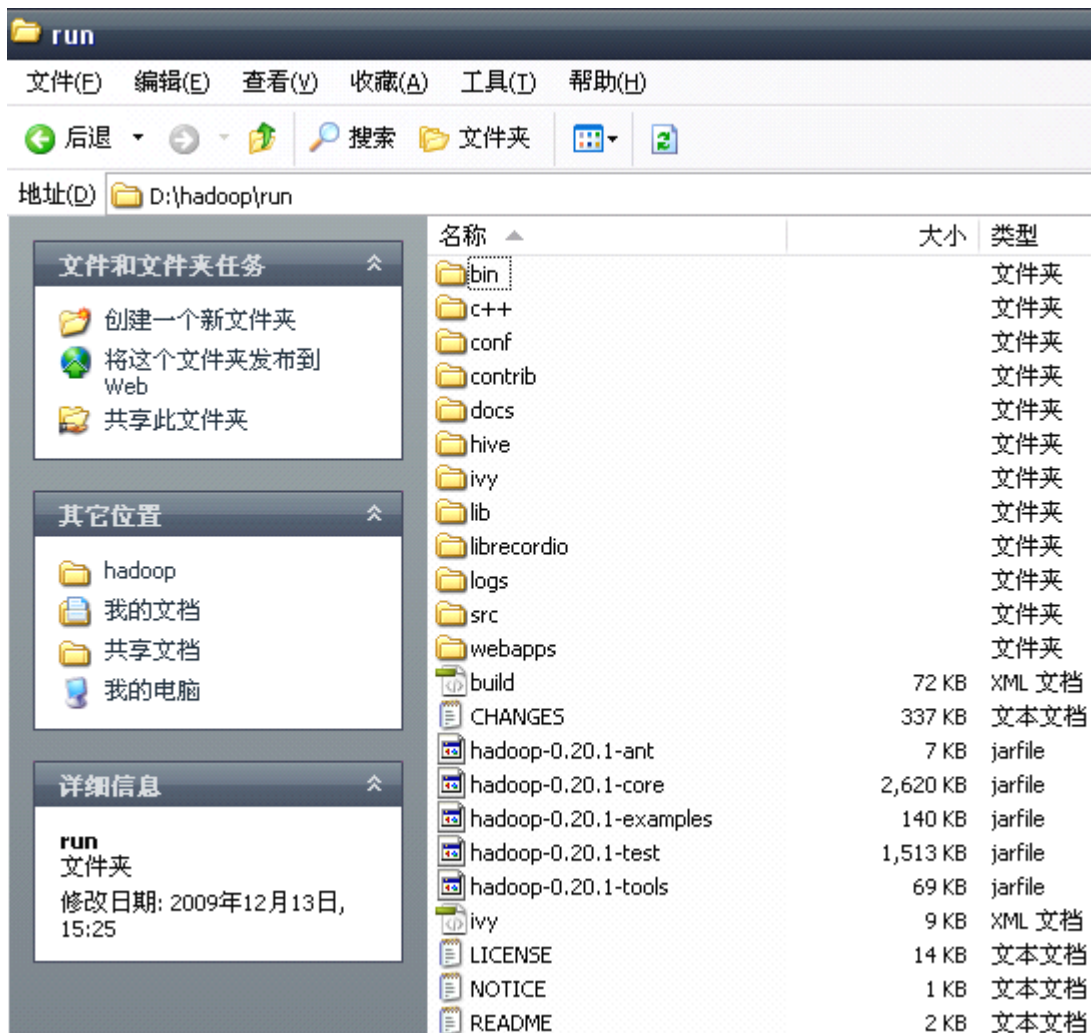
7. 下载 hadoop 安装包

hadoop 安装包下载地址：

<http://labs.xiaonei.com/apache-mirror/hadoop/core/hadoop-0.20.1/hadoop-0.20.1.tar.gz>

8. 安装 hadoop

将 hadoop 安装包 hadoop-0.20.1.tar.gz 解压到 **D:\hadoop\run** 目录(可以修改成其它目录)下, 如下图所示:



接下来, 需要修改 hadoop 的配置文件, 它们位于 conf 子目录下, 分别是 hadoop-env.sh、core-site.xml、hdfs-site.xml 和 mapred-site.xml 共四个文件。在 Cygwin 环境, masters 和 slaves 两个文件不需要修改。

- **修改 hadoop-env.sh**

只需要将 JAVA_HOME 修改成 JDK 的安装目录即可, 请注意 JDK 必须是 1.6 或以上版本。

- **修改 core-site.xml**

为简化 core-site.xml 配置, 将 D:\hadoop\run\src\core 目录下的 core-default.xml 文件复制到 D:\hadoop\run\conf 目录下, 并将 core-default.xml 文件名改成 core-site.xml。修改 fs.default.name 的值, 如下所示:

```
<property>
  <name>fs.default.name</name>
  <value>hdfs://localhost:8888</value>
  <description>The name of the default file system
scheme and authority determine the FileSystem im
uri's scheme determines the config property (fs.
the FileSystem implementation class. The uri's
determine the host, port, etc. for a filesystem.
</property>
```

上图中的端口号 **8888**，可以改成其它未被占用的端口。

- **修改 hdfs-site.xml**

为简化 hdfs-site.xml 配置，将 D:\hadoop\run\src\hdfs 目录下的 hdfs-default.xml 文件复制到 D:\hadoop\run\conf 目录下，并将 hdfs-default.xml 文件名改成 hdfs-site.xml。不需要再做其它修改。

- **修改 mapred-site.xml**

为简化 mapred-site.xml 配置，将 D:\hadoop\run\src\mapred 目录下的 mapred-default.xml 文件复制到 D:\hadoop\run\conf 目录下，并将 mapred-default.xml 文件名改成 mapred-site.xml。

```
<property>
  <name>mapred.job.tracker</name>
  <value>localhost:9999</value>
  <description>The host and port that the MapReduce job
at. If "local", then jobs are run in-process as a si
and reduce task.
</description>
</property>
```

上图中的端口号 **9999**，可以改成其它未被占用的端口。到这里，hadoop 宣告安装完毕，可以开始体验 hadoop 了！

9. 启动 hadoop

在 Cygwin 中，进入 hadoop 的 bin 目录，运行 ./start-all.sh 启动 hadoop，在启动成功之后，可以执行 ./hadoop fs -ls / 命令，查看 hadoop 的根目录，如下图所示：


```
/cygdrive/d/hadoop/run/bin
ondarynamenode-yijian.out' for reading: No such file or directory
starting jobtracker, logging to /cygdrive/d/hadoop/run/bin/../../logs/hadoop-
obtracker-yijian.out
localhost: starting tasktracker, logging to /cygdrive/d/hadoop/run/bin/..
adoop-jian-tasktracker-yijian.out
localhost: /cygdrive/d/hadoop/run/bin/hadoop-daemon.sh: line 117: /cygdr
doop/run/bin/../../logs/hadoop-jian-tasktracker-yijian.out: Permission deni
localhost: head: cannot open `/cygdrive/d/hadoop/run/bin/../../logs/hadoop-
ktracker-yijian.out' for reading: No such file or directory
jian@yijian /cygdrive/d/hadoop/run/bin
$ jps
5332 JobTracker
4192 Jps
4220 NameNode

jian@yijian /cygdrive/d/hadoop/run/bin
$ ./hadoop fs -ls /
Found 4 items
drwxr-xr-x  - jian supergroup          0 2009-12-18 19:28 /tmp
$
```

如果运行 mapreduce，请参考其它文档，本教程的内容到此结束。