# Data Management at ARSC

David Newman

(From slides by Tom Logan)

(from Slides from Don Bahls)

# Presentation Overview

1. ARSC storage

2. Data Management within ARSC

3. Additional Notes on Long Term Storage

4. Moving data to/from ARSC to your desktop system

5. Using the queues to manage data

# 1. ARSC storage

- **ARSC provides storage in three primary locations.  Environment variables are defined for each location.**
    - `$HOME`
    - `$CENTER`
    - `$ARCHIVE` or `$ARCHIVE_HOME`

# $HOME

- **Purpose**: location to store configuration files and commonly used executables.

- **Quota**: 8 GB

- **Backed Up**: yes

- **Purged**: no

- **Notes**: Available from computational nodes and login nodes. However, ARSC recommends that you avoid accessing $HOME in parallel jobs.

# $CENTER

- **Purpose**: place to run jobs and store temporary files.

- **Quota**: 750 GB (not enforced at this time).

- **Backed Up**: no

- **Purged**: yes (not enforced at this time).

- **Notes**: Available from computational nodes and login nodes.

# StorageTek Silo & Sun Fire 5440

# $ARCHIVE

- **Purpose**: place to store files long term.

- **Quota**: no quota

- **Backed Up**: yes

- **Purged**: no

- **Notes**: May not be available from all computational nodes.  Available from login nodes.  Files can be offline.  Network Filesystem (NFS) hosted by Sun T5440 system: bigdipper.

# 2) Data Management within ARSC
## Part I

- **Common UNIX commands for local and NFS mounted filesystems.**
    - `mv` move a file or directory
    - `cp` copy a file or directory
    - `rm` remove a file or directory
    - `mkdir` make a directory
    - `rmdir` remove a directory
    - `show_storage` quotas and usage (HPC systems)
    - `quota` quotas and usage (linux workstations)
    - `du` shows disk usage

# A few examples

**Make a directory in $CENTER**

```
f2n1 35% mkdir $CENTER/job1
```

**Copy myfile to $CENTER/job1**

```
f2n1 37% cp myfile $CENTER/job1
```

**Check disk usage (-sk gives summary in kilobytes)**

```
f2n1 38% du -sk $CENTER/job1
16      /center/w/usera/job1
```

**What's in $CENTER/job1**

```
f2n1 39% ls -la $CENTER/job1
total 64
drwx------   2 usera    staff         8192 Jul 17 10:44 .
drwxr-xr-x  19 usera    staff         8192 Jul 17 10:43 ..
-rw-------   1 usera    staff            0 Jul 17 10:44 myfile
```

**Make a directory to store results in $ARCHIVE_HOME (-p makes intermediate directories)**

```
f2n1 40% mkdir -p $ARCHIVE/ICEFLYER/job1/
```

Arctic Region Supercomputing Center

# A few examples continued

**Move myfile from $CENTER to $ARCHIVE**

```
f2n1 41% mv $CENTER/job1/myfile $ARCHIVE/ICEFLYER/job1
```

**Recursive copy**

```
f2n1 43% cp -r $CENTER/job1 $ARCHIVE/ICEFLYER
```

**Using special directories (".") & "..")**

```
f2n1 52% cp -r ../job0 .
```

**Remove myfile**

```
f2n1 53% rm $CENTER/job1/myfile
```

**Recursive remove**

```
f2n1 54% rm -r $CENTER/job1
```

**Checking quotas (use for  linux workstations)**

```
klondike 3% quota -v

Disk quotas for usera (uid 2640):
```

| Filesystem | usage | quota | limit | timeleft | files | quota | limit | timeleft |
|---|---|---|---|---|---|---|---|---|
| /u1 | 79896 | 112400 | 152640 | | 545 | 0 | 0 | |
| /u2 | 0 | 102400 | 112640 | | 0 | 0 | 0 | |
| /tmp | 103896 | 10485760 | 11534336 | | 423 | 0 | 0 | |

Arctic Region Supercomputing Center

# More information

- **All of the aforementioned commands have man pages.**

- **For example:** `man cp, man du, etc.`

- **NOTE: Command options may vary with the operating system.**

- **If you have questions don't forget about the ARSC help desk!**
  - Phone: **(907)450-8602** (**x8602** on campus)
  - Email: **consult@arsc.edu**

# Moving Data between ARSC Systems
## Part II

- **Moving files between systems.**
  - `scp` ssh version of copy
  - `sftp` ssh version of ftp

- **These options are available to users from their local machine (if you are using a UNIX variant).**

- `scp` **supports recursive copies and wildcards (I.e "*","?", etc.)**

- `scp` **requires that you know the path to the files you want.**

# A few examples

**Using scp (be wary of using environment variables!)**

```
f2n1 35% scp -r "iceflyer:/archive/u1/uaf/bahls/ICEFLYER/job1" .
```

**Using sftp (a few commands…)**

```
ftp> open iceflyer.arsc.edu

Connected to iceflyer.arsc.edu.

220 f2n1 FTP server (Version 5.60) ready.

334 Using authentication type GSSAPI; ADAT must follow

GSSAPI accepted as authentication type

GSSAPI authentication succeeded

Name (iceflyer.arsc.edu:fred): usera

...

ftp> get .cshrc

local: .cshrc remote: .cshrc

229 Entering Extended Passive Mode (|||62653|)

150 Opening BINARY mode data connection for .cshrc (2504 bytes).

226 Transfer complete.
```

Arctic Region Supercomputing Center

# A few examples continued

**FTP Commands**

- **`get`** **get a single file from remote system**
- **`put`** **put a single file to remote system**
- **`mget`** **get multiple files from remotes system**
- **`mput`** **put muliple files to remote system**
- **`ls`** **list the contents of a directory on remote system**
- **`cd`** **change remote directory**
- **`lcd`** **change directory on local host**
- **`help`** **shows the ftp help pages**

```
ftp> help

Commands may be abbreviated.   Commands are:


!               cr              mdir            sendport        send

$               delete          mget            put             site

account         debug           mkdir           pwd             size

append          dir             nls             quit            status
```

Arctic Region Supercomputing Center

# 3) Additional Notes on Long

- **Long term storage at ARSC is served by Sun Fire 5440 system.**

- **There are no quotas on the archive filesystem, so there's no need to micro-manage data.**

- **Most of the time there is no need to access the servers directly.**

Tuesday, September 15, 15

# When to Log on to Archive?

- **Large transfers (10's of GB+) to your local machine.   There are several advantages:**
  - Manually issue stage commands (see next slide) to ensure the files to be transferred are online
  - Better overall transfer rates (avoids an extra network transfer).

- **To determine whether or not a file is offline.**

- **When creating big tar files of data on $ARCHIVE.**

# Archive Commands

- `stage` **brings a file or files online.**

- `release` **tells the system to release the on disk copy of the file leaving tape copies only.**

- `sfind` **like** `find` **with flags to determine whether or not a file is online.**

- `sdu` **shows disk usage including offline usage.**

- `sls` **like standard** `ls` **with options to see whether or not a file is online.**

- `batch_stage` **stages a list of files from tape in an orderly manner (ARSC developed).**

Tuesday, September 15, 15

# Archive Examples

**Find all offline files in the current directory.**

```
nanook 10% sfind . -name \* -offline

./my.tar.gz
```

**Check the status of a file using sls.**

```
nanook 11% sls -2 my.tar.gz

-rw-r--r--   1 usera    staff      669944 Jan 27  2005
   my.tar.gz

O--------  guv-- -- --  sg sf
```

**Bring an offline file back online**

```
nanook 13% stage -w my.tar.gz

nanook 14% sls -2 my.tar.gz

-rw-r--r--   1 usera    staff      669944 Jan 27  2005
   my.tar.gz
```

# Archive Examples Cont.

**The batch_stage script was developed at ARSC to improve access to offline files.  If you have large number of files you need to access, consider using it.**
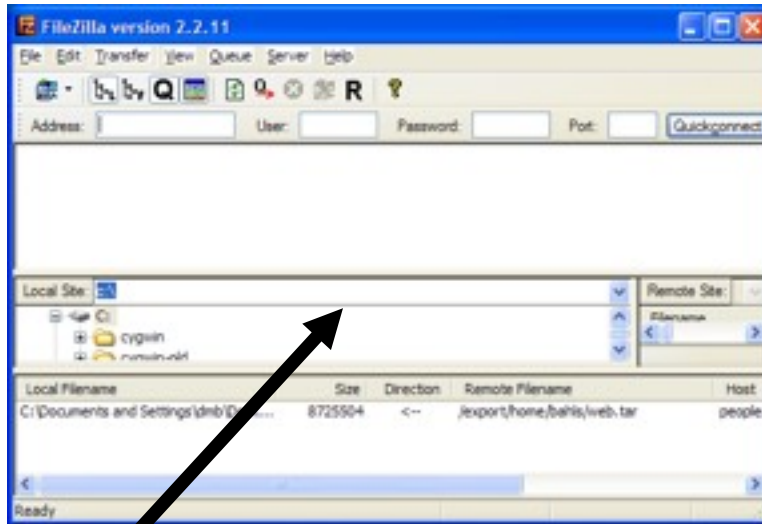
**Staging all files in a directory**
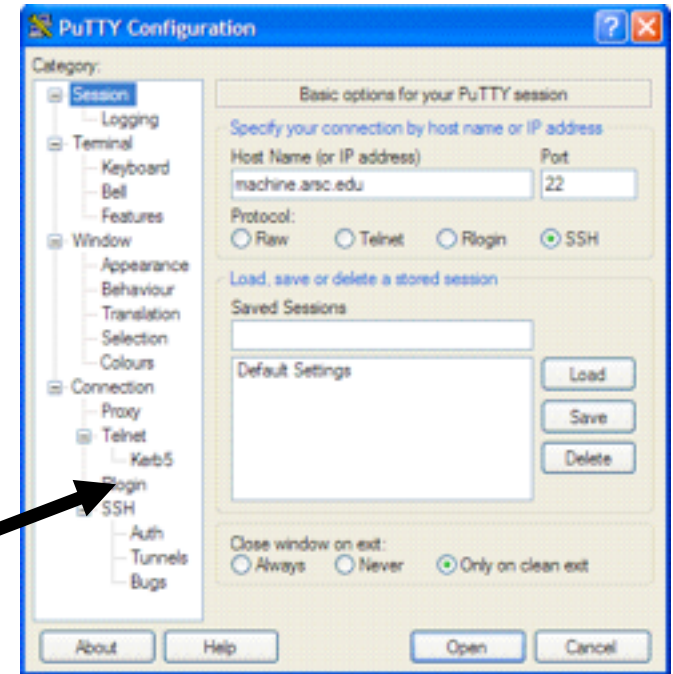
```
nanook 15% batch_stage $ARCHIVE/ICEFLYER/mydata/*
```

**Staging all files in a directory tree**

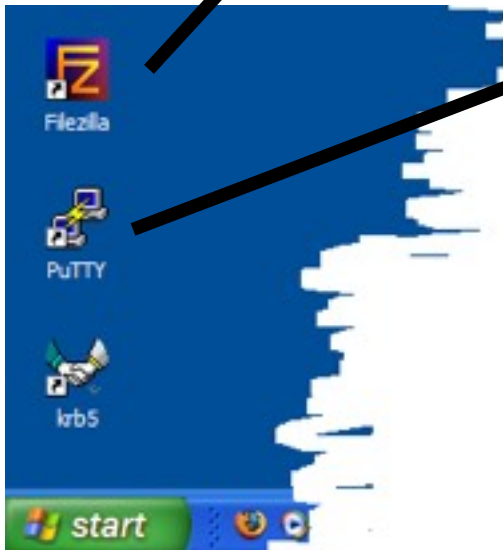# 4) Moving data to/from ARSC to your desktop system

- **Unix-like OS: things are pretty much the same as transferring between machines within ARSC.**

- **Windows Systems**

  - Putty ssh client.

  - Filezilla ftp/sftp client.

  - Others exists as well (e.g. cygwin)

**Filezilla**

**PuTTY**

Arctic Region Supercomputing Center

# Using pscp

1) Open a Windows 'Command Prompt'.

2) Change directory to the directory where your files are

located.

**(e.g.** `cd "C:\Documents and Settings\default\My Documents"` **)**

### 3) Run pscp.exe

```
"C:\Program Files\HPCMP\Putty\pscp.exe" -r mydir
    "username@iceberg.arsc.edu:/u1/uaf/username"
```

# Using the queues to manage

- **As mentioned before $ARCHIVE may not be mounted on computational nodes and is generally not a good place to run your jobs.**

Tuesday, September 15, 15

# Moving data from a job to

- **Job chaining (one job submits the next) PBS**

- **Job dependencies (jobs are dependant on the exit status of previous jobs) PBS**

# Some References

- **Creating Sequences of Batch Jobs in PBS**
    - http://www.arsc.edu/support/news/HPCnews/HPCnews319.shtml
    - http://www.arsc.edu/support/news/HPCnews/HPCnews320.shtml

- **Scripted Chaining of Batch Jobs and File Checks**
    - http://www.arsc.edu/support/news/HPCnews/HPCnews297.shtml

- **Recursive Copies**
    - http://www.arsc.edu/support/news/HPCnews/HPCnews343.shtml#qt

- **Unrelated but maybe useful: X11 on Windows**
    - http://www.arsc.edu/support/howtos/usingcygwin.html

# Need more information?

- **Check out man pages**

- **Call or email the ARSC Help Desk:**
  - **PHONE: 907 450-8602 (x8602 on campus)**
  - **EMAIL: consult@arsc.edu**

- **ARSC website & HPC Users' Newsletter**
  1. http://www.arsc.edu/support
  2. http://www.arsc.edu/support/news/HPCnews.shtml