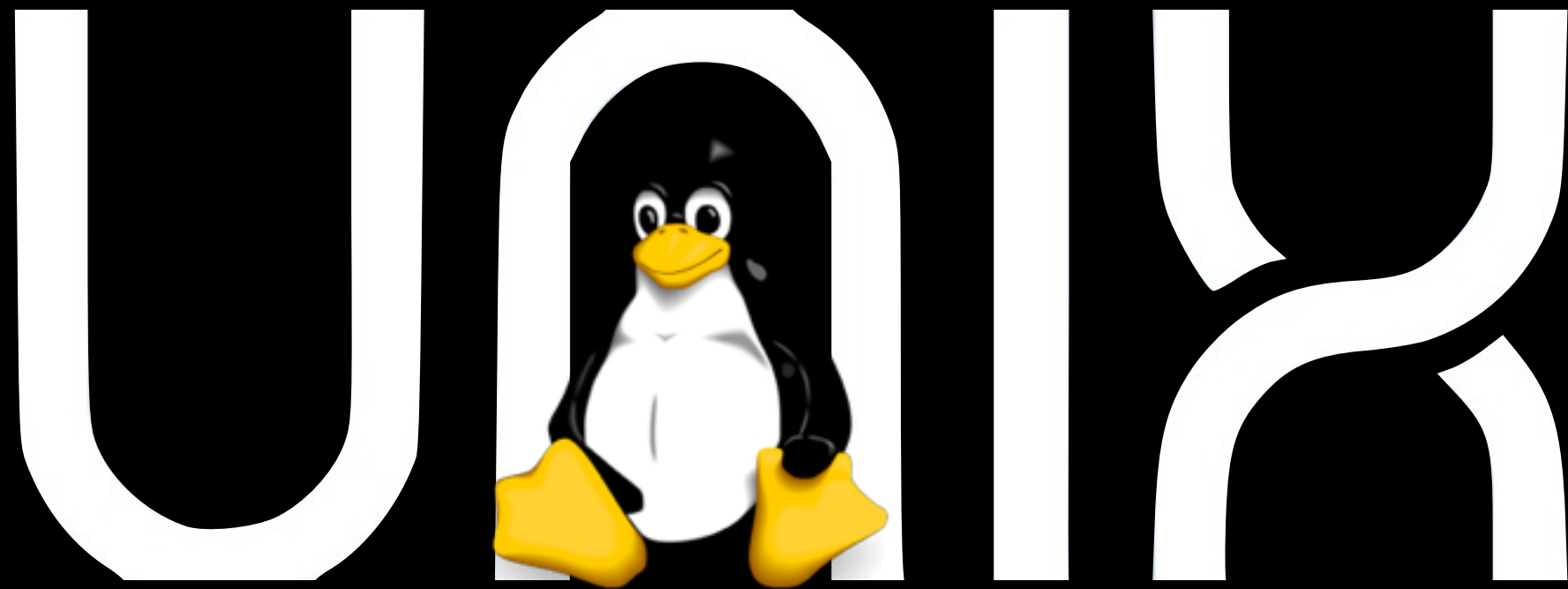


Introduction To



Barry Grant

bjgrant@umich.edu
<http://thegrantlab.org>

Working with Unix

How do we actually use Unix?

Inspecting text files

- **less** - visualize a text file:
 - use arrow keys
 - page down/page up with “space”/“b” keys
 - search by typing “/”
 - quit by typing “q”
- Also see: **head**, **tail**, **cat**, **more**

Creating text files

Creating files can be done in a few ways:

- With a text editor (such as **nano**, **emacs**, or **vi**)
- With the **touch** command (\$ touch a_file)
- From the command line with **cat** or **echo** and redirection (>)
- **nano** is a simple text editor that is recommended for first-time users. Other text editors have more powerful features but also steep learning curves

Creating and editing text files with **nano**

Do it Yourself!

In the terminal type:

> nano yourfilename.txt

^G Get Help	^O WriteOut	^R Read File
^X Exit	^J Justify	^W Where Is
^Y Prev Page	^K Cut Text	^C Cur Pos
^V Next Page	^U UnCut Txt	^T To Spell

^ - Press Control

- There are many other text file editors (e.g. vim, emacs and sublime text, etc.)

Finding the Right Hammer (**man** and **apropos**)

- You can access the manual (i.e. user documentation) on a command with **man**, e.g:

 > man pwd
- The man page is only helpful if you know the name of the command you're looking for.
apropos will search the man pages for keywords.

 > apropos "working directory"

Combining Utilities with **Redirection** (>, <) and **Pipes** (|)

- The power of the shell lies in the ability to combine simple utilities (*i.e.* commands) into more complex algorithms very quickly.
- A key element of this is the ability to send the output from one command into a file or to pass it directly to another program.
- This is the job of >, < and |

Side-Note: Standard Input and Standard Output streams

Two very important concepts that unpin Unix workflows:

- Standard Output (**stdout**) - default destination of a program's output. It is generally the terminal screen.
- Standard Input (**stdin**) - default source of a program's input. It is generally the command line.

Output redirection and piping

Do it Yourself!



> ls /usr/bin # stdin is "/usr/bin"; stdout to screen

Output redirection and piping

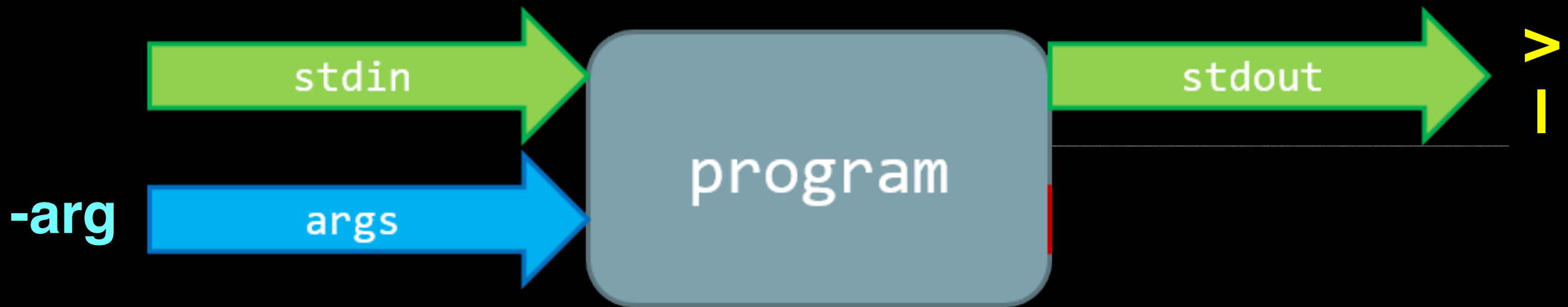


> ls /usr/bin # stdin is "/usr/bin"; stdout to screen

> ls /usr/bin > binlist.txt # stdout **redirected** to file

> ls /usr/bin | less # stdout **piped** to less (no file created)

Output redirection and piping



> ls /usr/bin # stdin is "/usr/bin"; stdout to screen

> ls /usr/bin > binlist.txt # stdout **redirected** to file

> ls /usr/bin | less # stdout **piped** to less (no file created)

> ls -l /usr/bin # extra optional input **argument** "-l"

Output redirection and piping



- > `ls /usr/bin` # `stdin` is `"/usr/bin"`; `stdout` to screen
- > `ls /usr/bin > binlist.txt` # `stdout` **redirected** to file
- > `ls /usr/bin | less` # `stdout` **piped** to `less` (no file created)
- > `ls /noindexexists/` # `stderr` to screen

Output redirection and piping

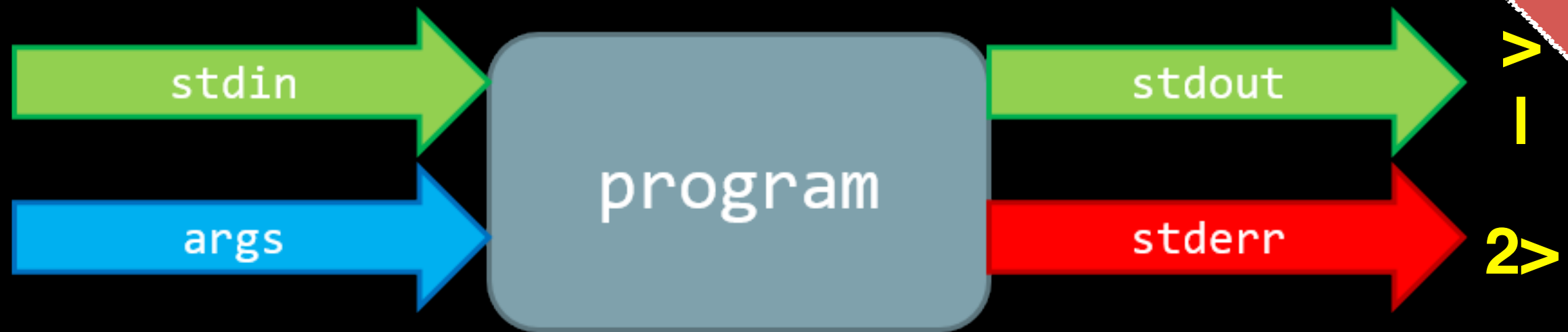
Do it Yourself!



- > ls /usr/bin # stdin is "/usr/bin"; stdout to screen
- > ls /usr/bin > binlist.txt # stdout **redirected** to file
- > ls /usr/bin | less # stdout **piped** to less (no file created)
- > ls /nonexistent/ > binlist.txt # stderr to **screen**

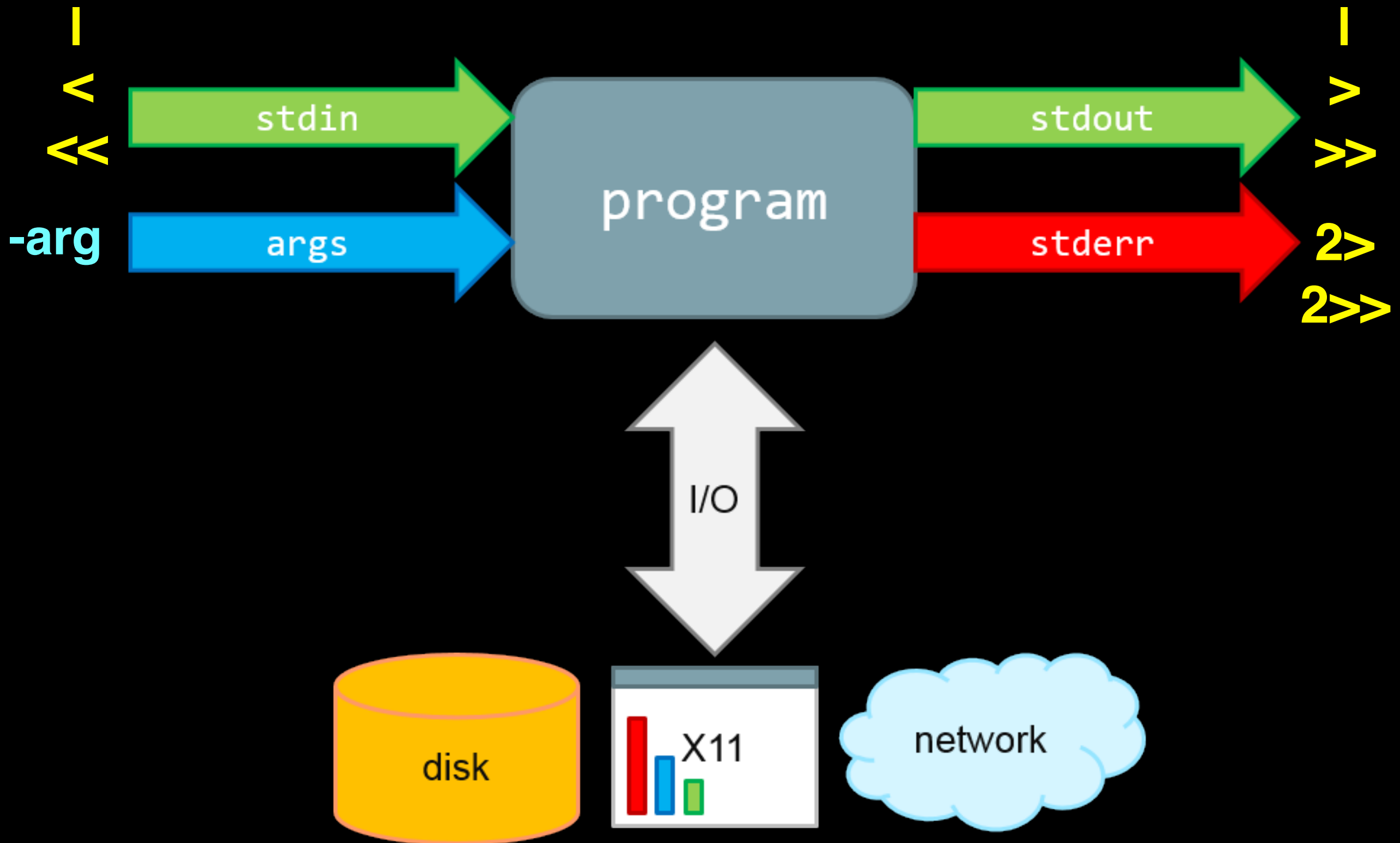
Output redirection and piping

Do it Yourself!

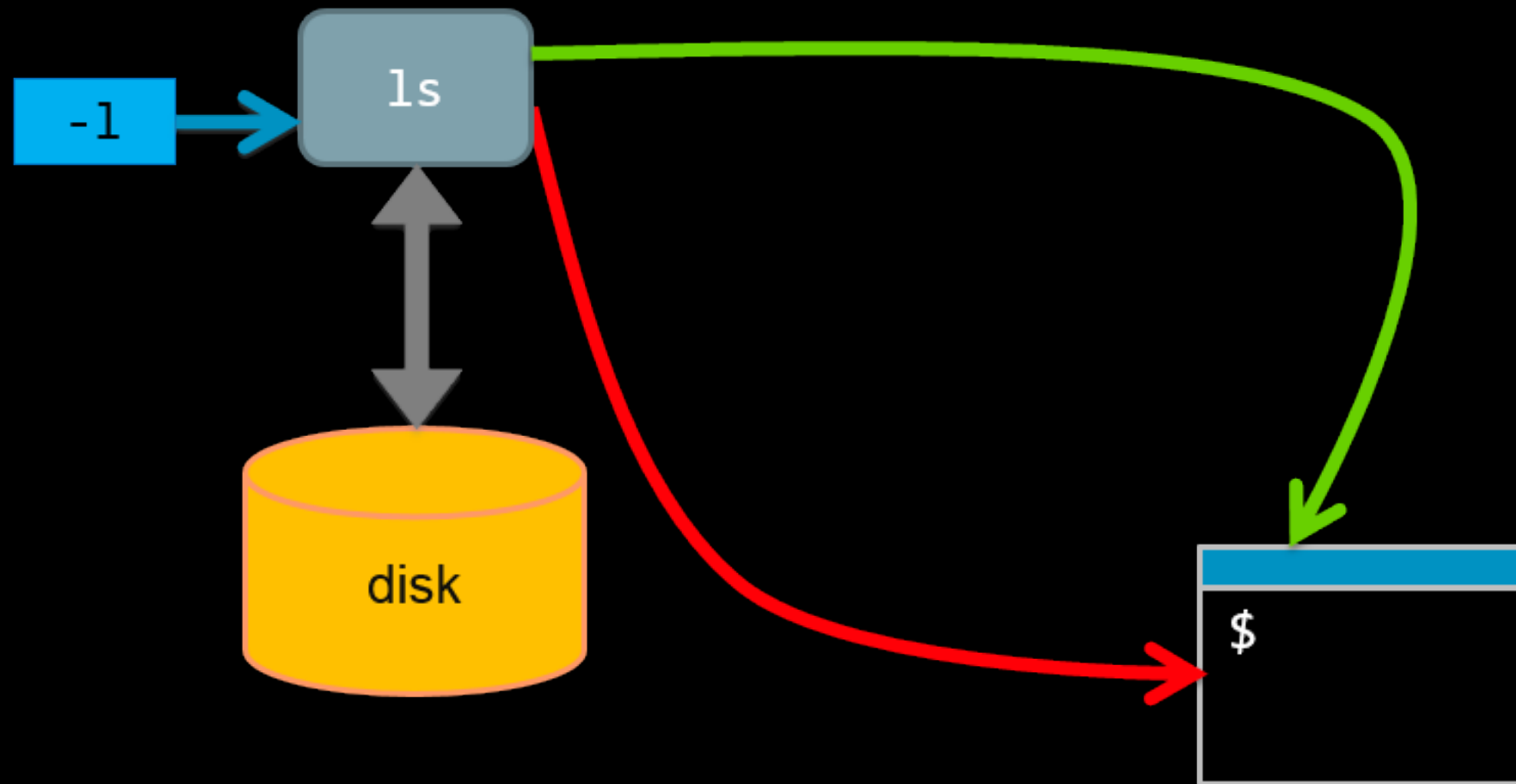


- > ls /usr/bin # stdin is "/usr/bin"; stdout to screen
- > ls /usr/bin > binlist.txt # stdout **redirected** to file
- > ls /usr/bin | less # stdout **pipelined** to less (no file created)
- > ls /nonexistent/ 2> binlist.txt # stderr to **file**

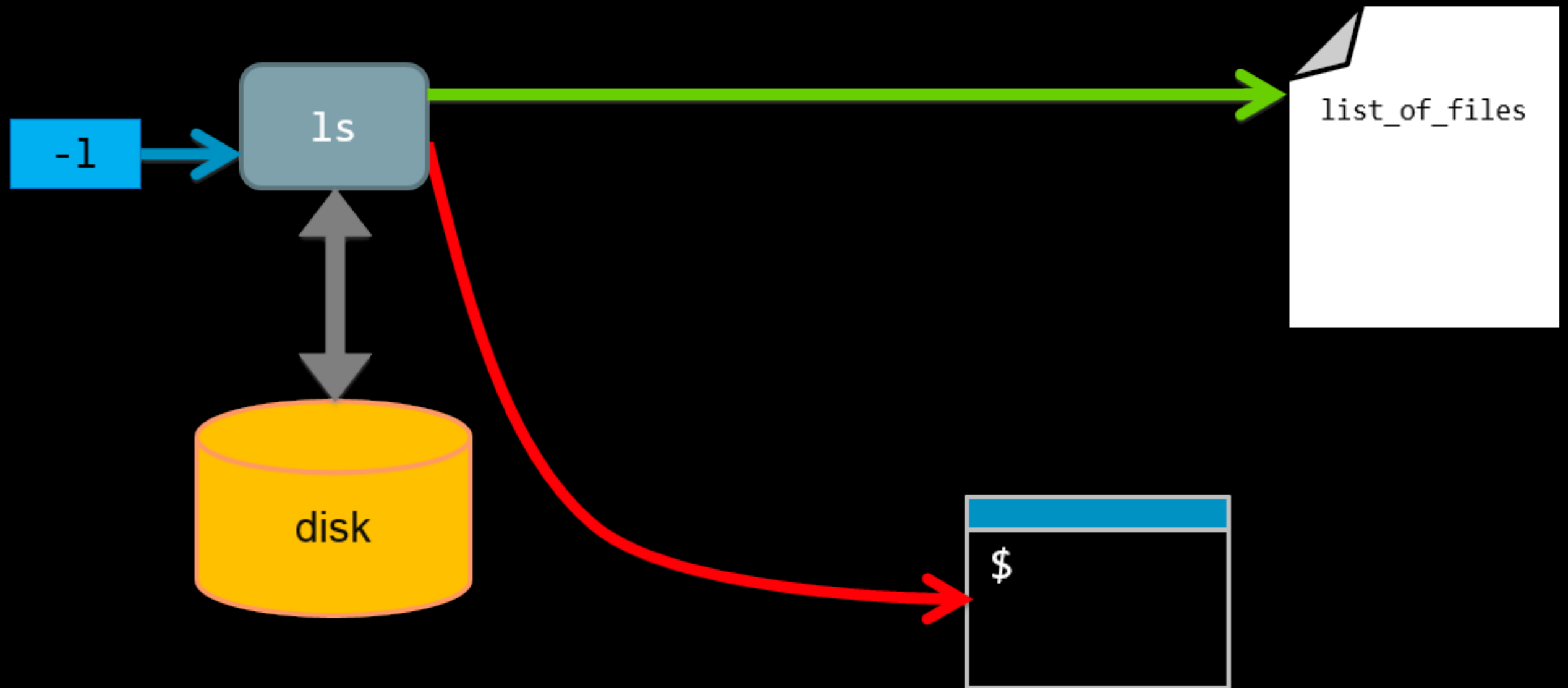
Output redirection summary



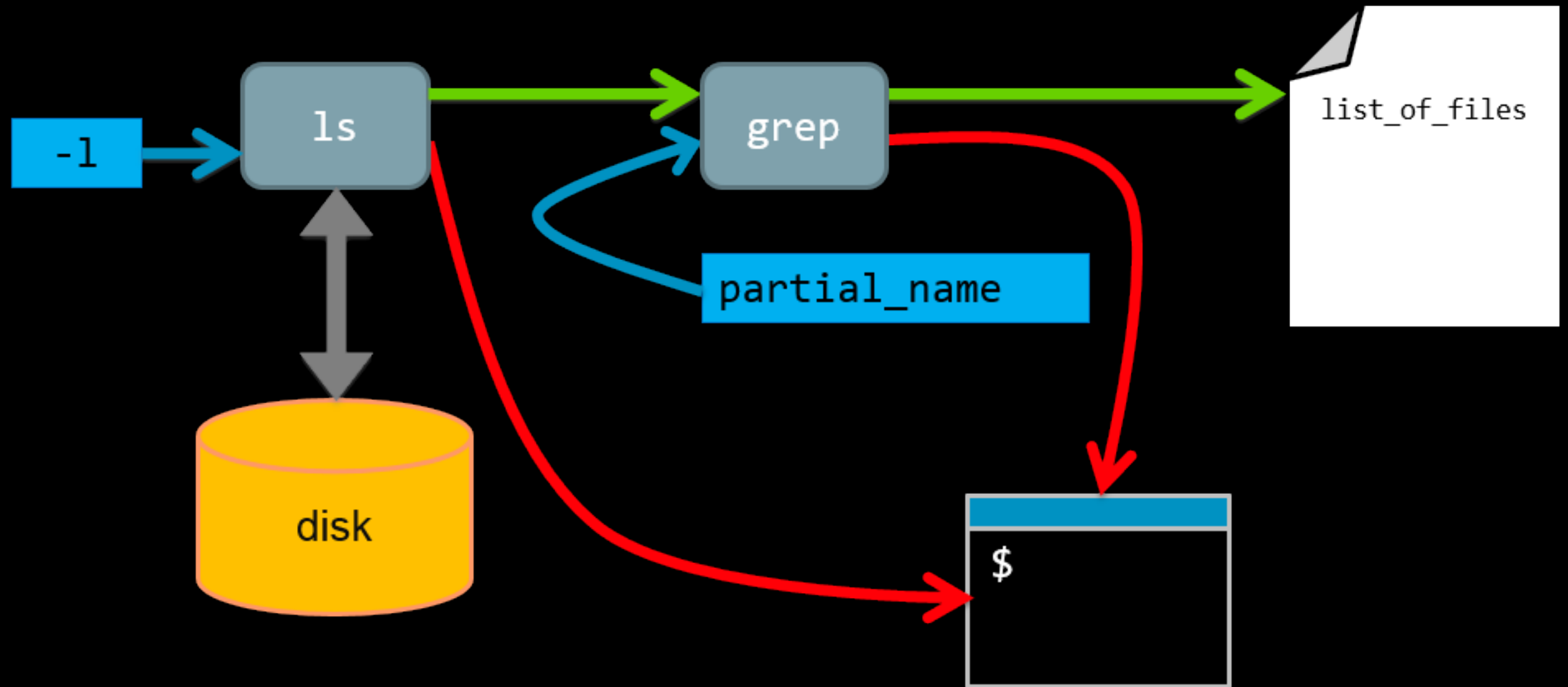
ls -l



ls -l > list_of_files



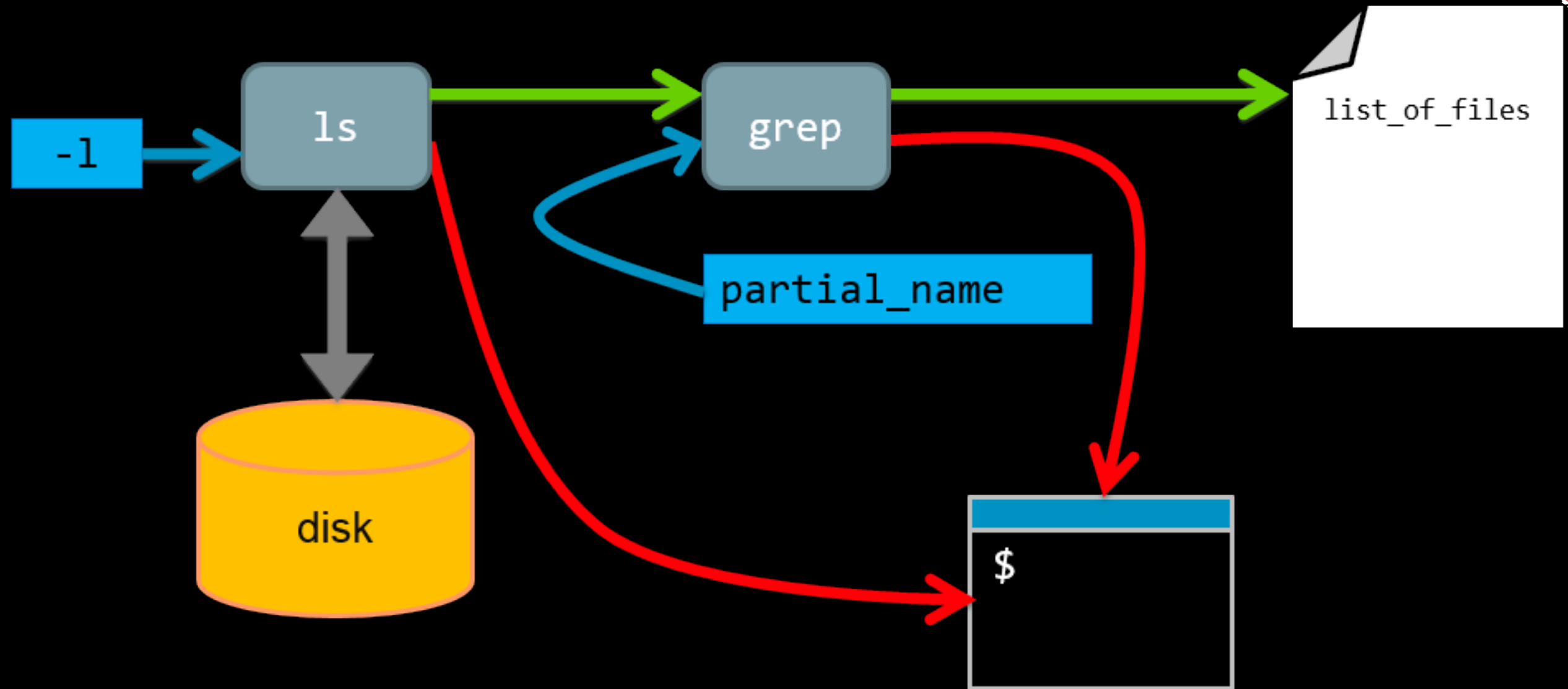
```
ls -l | grep partial_name > list_of_files
```



We have piped (|) the stdout of one command into the stdin of another command!

Do it Yourself!

```
ls -l /usr/bin/ | grep "tree" > list_of_files
```



grep: prints lines containing a string.
Also searches for strings in text files.

Basics	File Control	Viewing & Editing Files	Misc. useful	Power commands	Process related
ls	mv	less	chmod	grep	top
cd	cp	head	echo	find	ps
pwd	mkdir	tail	wc	sed	kill
man	rm	nano	curl	uniq	Ctrl-c
ssh	 (pipe)	touch	source	git	Ctrl-z
	> (write to file)		cat	R	bg
	< (read from file)			python	fg

Side-Note: **grep** ‘power command’

Do it Yourself!

- **grep** - prints lines containing a string pattern. Also searches for strings in text files, e.g.

> **grep --color "GESGKS" sequences/data/seqdump.fasta**

REVKLLLLGAG**GESGKS**TIVKQMKIHEAGYSEEECKQYK

- **grep** is a ‘power tool’ that is often used with pipes as it accepts **regular expressions** as input (e.g. “**G..GK[ST]**”) and has lots of useful options - see the *man page* for details.

grep example using regular expressions

Do it Yourself!

- Suppose a program that you are working with complains that your input sequence file contains non-nucleotide characters. You can eye-ball your file or ...

```
> grep -v “^>” seqdump.fasta | grep --color “[^ATGC]”
```

Exercises:

- (1). Use “man grep” to find out what the **-v** argument option is doing!
- (2). How could we also show line number for each match along with the output?
(tip you can grep the output of “man grep” for ‘line number’)

grep example using regular expressions

Do it Yourself!

- Suppose a program that you are working with complains that your input sequence file contains non-nucleotide characters. You can eye-ball your file or ...

```
> grep -v "^>" seqdump.fasta | grep --color -n "[^ATGC]"
```

- First we remove (with **-v** option) lines that start with a ">" character (these are sequence identifiers).
- Next we find characters that are *not* A, T, C or G. To do this we use ^ symbols second meaning: *match anything but* the pattern in square brackets. We also print line number (with **-n** option) and color output (with **--color** option).

Key Point: Pipes and redirects avoid unnecessary i/o

- Disc i/o is often a bottleneck in data processing!
- Pipes prevent unnecessary disc i/o operations by connecting the stdout of one process to the stdin of another (these are frequently called “**streams**”)

```
> program1 input.txt 2> program1.stderr | \  
program2 2> program2.stderr > results.txt
```

- Pipes and redirects allow us to build solutions from modular parts that work with **stdin** and **stdout streams**.

Unix 'Philosophy' Revisited

“Write programs that do one thing and do it well. Write programs to work together and that encourage open standards. Write programs to handle text streams, because that is a universal interface.”

— Doug McIlory



Pipes provide speed, flexibility and sometimes simplicity...

- In 1986 “*Communications of the ACM magazine*” asked famous computer scientist Donald Knuth to write a simple program to count and print the k most common words in a file alongside their counts, in descending order.
- Knuth wrote a literate programming solution that was 7 pages long, and also highly customized to this problem (e.g. Knuth implemented a custom data structure for counting English words).
- Doug McIlroy replied with one line:

```
> cat input.txt | tr A-Z a-z | sort | uniq -c | sort -rn | sed 10q
```

Key Point:

You can chain any number of programs together to achieve your goal!



This allows you to build up fairly complex workflows within one command-line.

Shell scripting

Do it Yourself!

```
#!/bin/bash  
# This is a very simple hello world script.  
echo "Hello, world!"
```

Exercise:

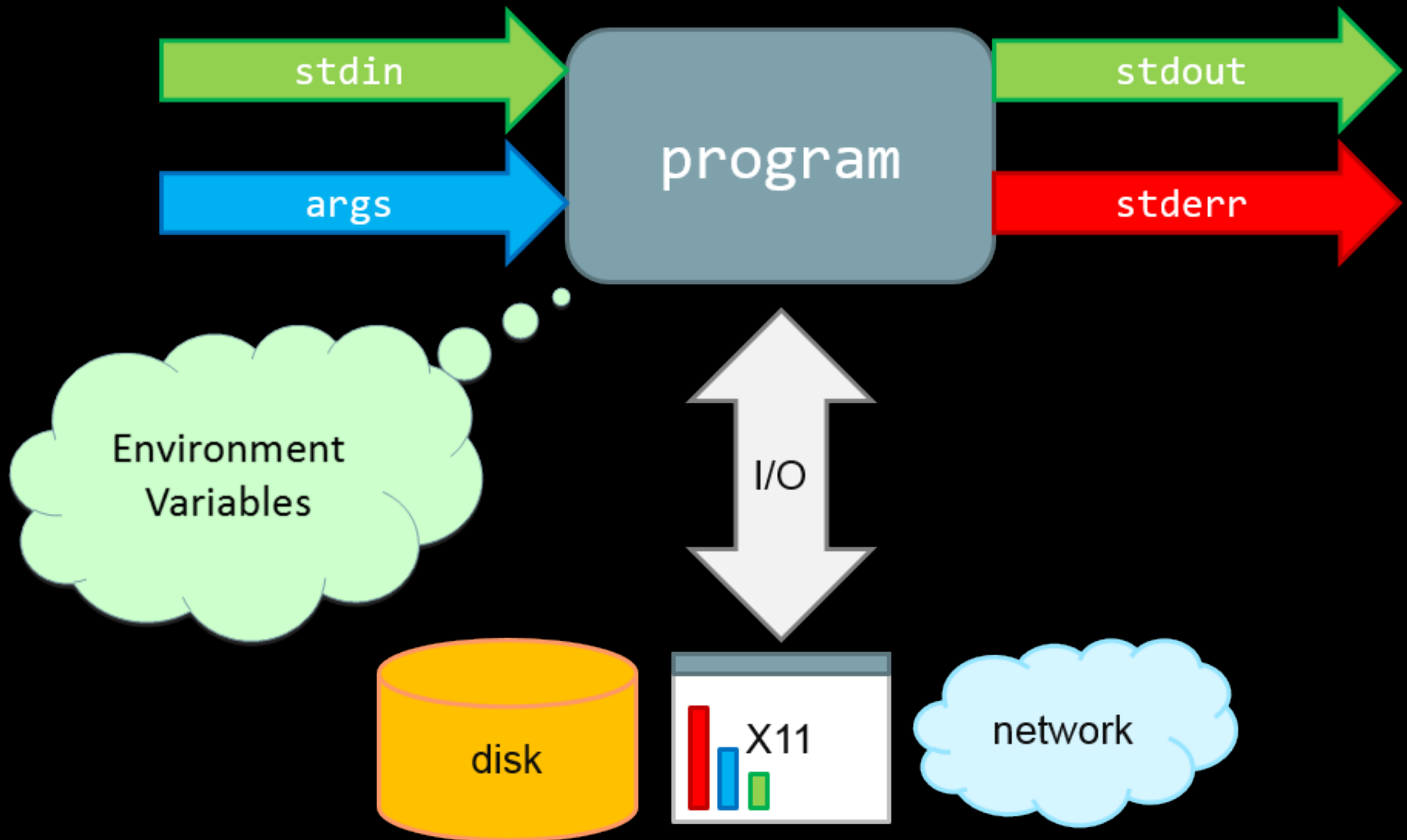
- Create a "Hello world"-like script using command line tools and execute it.
- Copy and alter your script to redirect output to a file using **>** along with a list of files in your home directory.
- Alter your script to use **>>** instead of **>**. What effect does this have on its behavior?

Variables in shell scripts

```
#!/bin/bash  
# Another simple hello world script  
message='Hello World!'  
echo $message
```

- “message” - is a **variable** to which the string 'Hello World!' is assigned
- echo - prints to screen the contents of the variable "\$message"

Side-Note: Environment Variables



\$PATH 'special' environment variable

- What is the output of this command?

 > echo \$PATH
- Note the structure: <path1>:<path2>:<path3>
- **PATH** is an environmental variable which Bash uses to search for commands typed on the command line without a full path.
- Exercise: Use the command **env** to discover more.

Q. Why have we been showing you this?

- On Day-4, we will be talking about how to submit your work to the FLUX high performance computing cluster.
- The scripts you use to submit your work on FLUX are basically bash shell scripts (with some special comments read by the scheduler at the top including instructions where to put stdout and stderr).

Summary

- Built-in unix shell commands allow for easy data manipulation (e.g. sort, grep, etc.)
- Commands can be easily combined to generate flexible solutions to data manipulation tasks.
- The unix shell allows users to automate repetitive tasks through the use of shell scripts that promote reproducibility and easy troubleshooting
- Introduced the 21 key unix commands that you will use during ~95% of your future unix work...

Basics	File Control	Viewing & Editing Files	Misc. useful	Power commands	Process related
ls	mv	less	chmod	grep	top
cd	cp	head	echo	find	ps
pwd	mkdir	tail	wc	sed	kill
man	rm	nano	curl	uniq	Ctrl-c
ssh	 (pipe)	touch	source	git	Ctrl-z
	> (write to file)		cat	R	bg
	< (read from file)			python	fg

Connecting to remote machines (**ssh** & **scp**)

- Most high-performance computing (HPC) resources can only be accessed by **ssh** (Secure SHell)
 - > ssh [user@host.address]
 - > ssh barry@scs.gpcc.itd.umich.edu
 - > ssh -X barry@flux-login.engin.umich.edu
- The **scp** (secure copy) command can be used to copy files and directories from one computer to another.
 - > scp [file] [user@host]
 - > scp localfile.txt bgrant@bigcomputer.net:/remotedir/.