

Practical Linux Examples

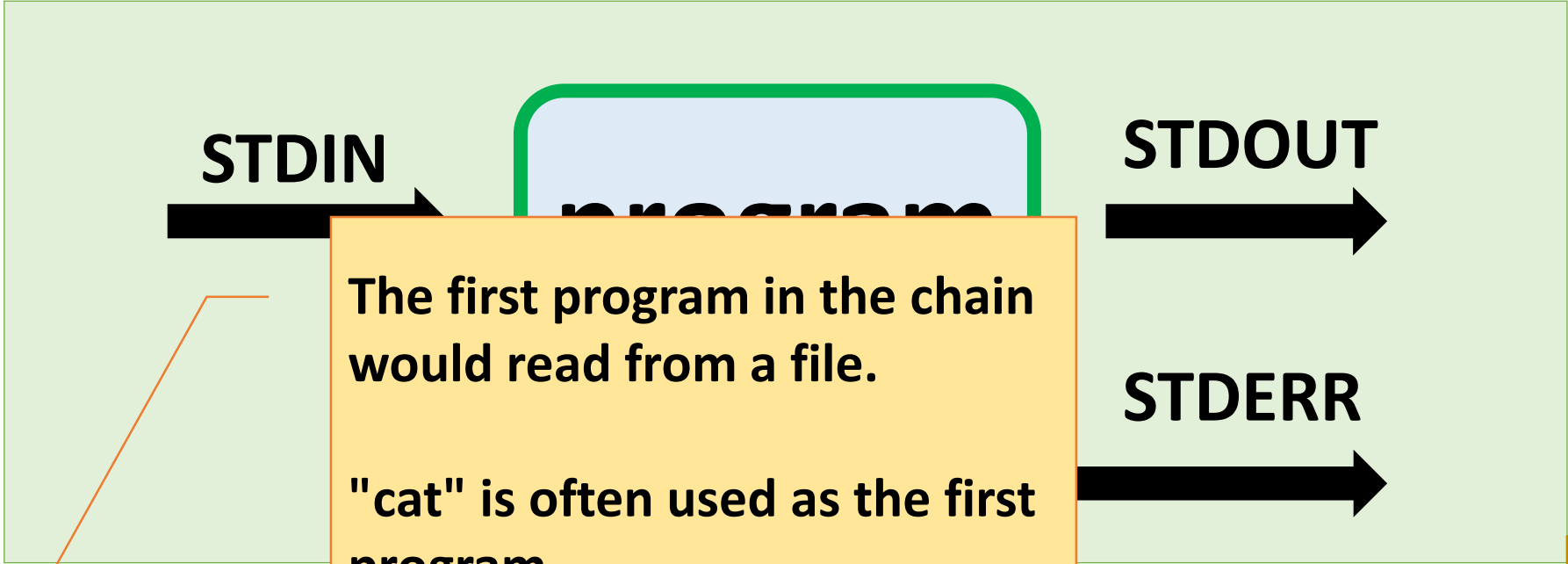
- Processing large text file
- Parallelization of independent tasks

Qi Sun & Robert Bukowski
Bioinformatics Facility
Cornell University

http://cbsu.tc.cornell.edu/lab/doc/linux_examples_slides.pdf

http://cbsu.tc.cornell.edu/lab/doc/linux_examples_exercises.pdf

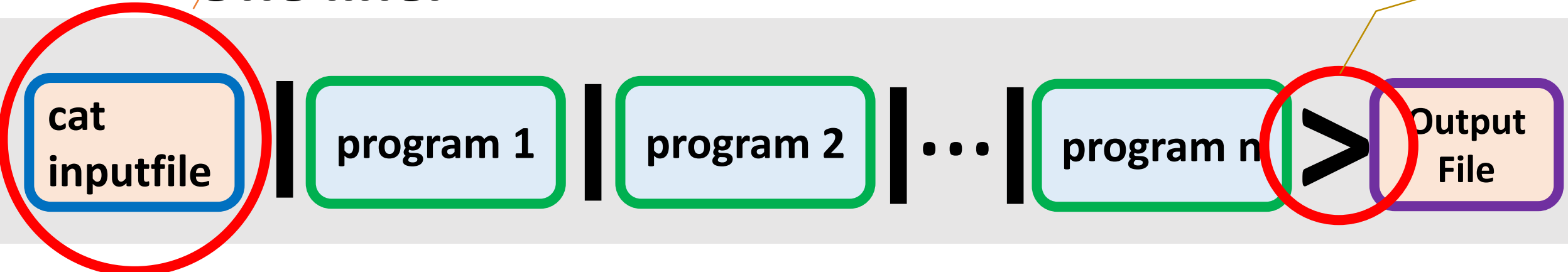
Three streams for a standard Linux program



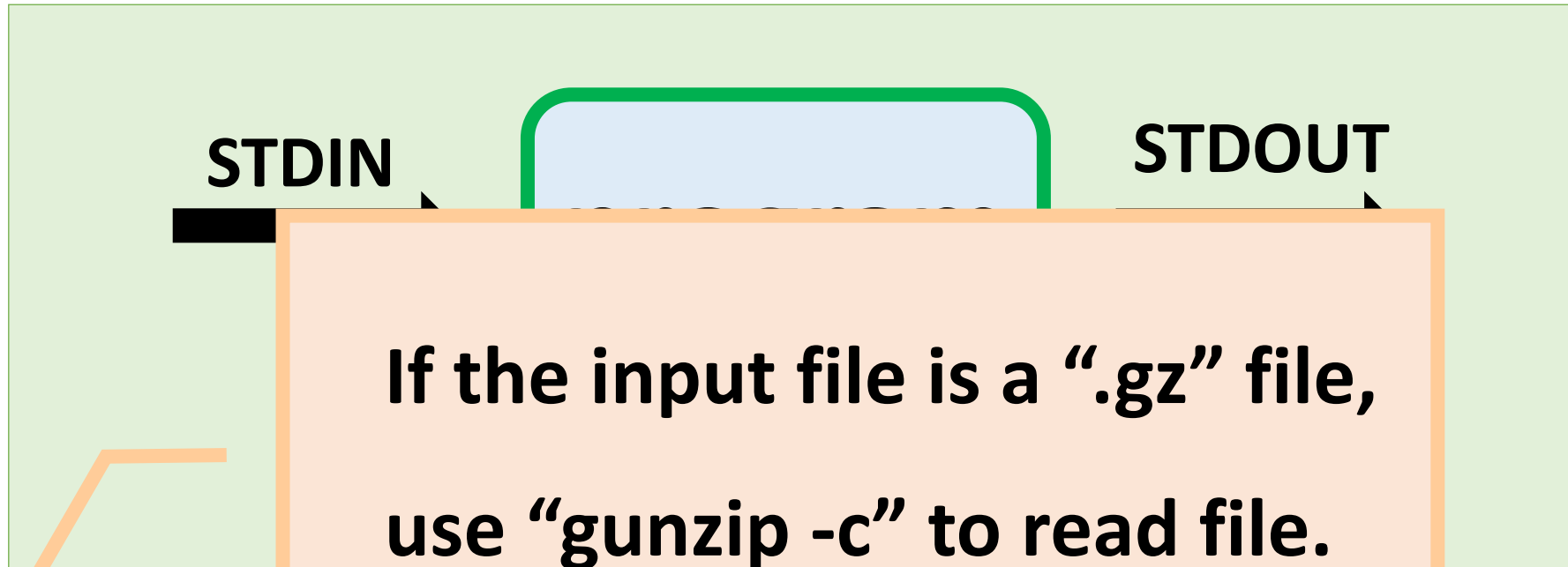
The first program in the chain would read from a file.
"cat" is often used as the first program.

Write to file

One liner



Three streams for a standard Linux program



One line



grep

Search for a pattern and output
matched lines

```
$ cat mydata.txt
```

```
AAGATCAAAAAAGA  
ATTACGAAAAAGA  
ACCTGTTGGATCAAAGTT  
AACTTTCGACGATCT  
ATTTTTTTAGAAAGG
```

```
$ cat mydata.txt | grep '[AC]GATC'
```

```
AAGATCAAAAAAGA  
AACTTTCGACGATCT
```


wc -l Count the number of lines

```
$ cat mydata.txt
```

```
AAGATCAAAAAAGA  
ATTACGAAAAAGA  
ACCTGTTGGATCCAAAGTT  
AAACTTTCGACGATCT  
ATTTTTTTAGAAAGG
```

```
$ cat mydata.txt | grep '[AC]GATC' | wc -l
```

```
2
```

sort

Sort the text in a file

```
$ sort myChr.txt
```

```
Chr1  
Chr10  
Chr2  
Chr3  
Chr4  
Chr5
```

```
$ sort -V myChr.txt
```

```
Chr1  
Chr2  
Chr3  
Chr4  
Chr5  
Chr10
```

```
$ sort -n myPos.txt
```

```
1  
2  
3  
4  
5  
10
```

sort

Sort the text by multiple columns

```
$ sort -k1,1V -k2,2n myChrPos.txt
```

```
Chr1      12  
Chr1     100  
Chr1     200  
Chr2      34  
Chr2     121  
Chr2     300
```

Some extra parameter to set for the “sort” command

LC_ALL=C sort -S 4G -k1,1 myChr.txt

LC_ALL=C vs **LC_ALL=US_en**

AA	aa
BB	a.a
a.a	AA
aa	bb
bb	BB

* BioHPC default locale is C. So you can skip this parameter

-S 4G

Use RAM as buffer 4G

uniq -c

Count the occurrence of unique tags

```
$ cat mydata.txt
```

```
ItemB  
ItemA  
ItemB  
ItemC  
ItemB  
ItemC  
ItemB  
ItemC
```

```
$ cat mydata.txt | sort | uniq -c
```

```
1  ItemA  
4  ItemB  
3  ItemC
```

Mark sure to run
"sort" before "uniq"

Merging files:

cat f1 f2 VS **paste** f1 f2 VS **join** f1 f2

File 1:

Item1

Item2

File2:

Item3

Item4

```
cat File1 File2 > mergedfile1
```

Item1

Item2

Item3

Item4

```
paste File1 File2 > mergedfile2
```

Item1

Item3

Item2

Item4

* Make sure that that two files has same number of rows and sorted the same way. Otherwise, use "join"

join

Merging two files that share a common field

File 1:

Gene1	DNA-binding
Gene2	kinase
Gene3	membrane

File2:

Gene2	764
Gene3	23
Gene4	34

```
join -1 1 -2 1 File1 File2 > mergedfile
```

```
Gene2      Kinase      764  
Gene3      membrane   23
```

```
join -1 1 -2 1 -a1 File1 File2 > mergedfile
```

```
Gene1      DNA-binding  
Gene2      Kinase      764  
Gene3      membrane   23
```

cut

Output selected columns in a table

```
$ cat mydata.txt
```

```
Chr1 1000 2250 Gene1  
Chr1 3010 5340 Gene2  
Chr1 7500 8460 Gene3  
Chr2 8933 9500 Gene4
```

```
$ cat mydata.txt | cut -f 1,4
```

```
Chr1 Gene1  
Chr1 Gene2  
Chr1 Gene3  
Chr2 Gene4
```

sed

Modify text in a file

```
$ cat mydata.txt
```

```
Chr1 1000 2250 Gene1
```

```
Chr1 3010 5340 Gene2
```

```
Chr1 7500 8460 Gene3
```

```
Chr2 8933 9500 Gene4
```

```
$ cat mydata.txt | sed "s/^Chr//"
```

```
1 1000 2250 Gene1
```

```
1 3010 5340 Gene2
```

```
1 7500 8460 Gene3
```

```
2 8933 9500 Gene4
```

awk

Probably the most versatile function
in Linux

```
$ cat mydata.txt
```

```
Chr1 1000 2250 Gene1  
Chr1 3010 5340 Gene2  
Chr1 7500 8460 Gene3  
Chr2 8933 9500 Gene4
```

```
$ cat mydata.txt |\  
awk '{if ($1=="Chr1") print $4}'
```

```
Gene1  
Gene2  
Gene3
```


A Good Practice:

Create a shell script file for the one liner

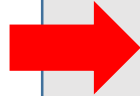
```
cat D7RACXX.fastq | \  
head -n 40000 | \  
grep AGATCGGAAGAGC | \  
wc -l
```

Run the shell script

```
sh checkadapter.sh
```

Debug a one-liner

```
gunzip -c human.gff3.gz | \
```

```
 awk 'BEGIN {OFS = "\t"}; {if ($3=="gene") print $1,$4-1,$5}' | \
```

```
bedtools coverage -a win1mb.bed -b stdin -counts | \
```

```
LC_ALL=C sort -k1,1V -k2,2n > gene.cover.bed
```

```
gunzip -c human.gff3.gz | head -n 1000 > tmpfile
```

```
cat tmpfile | \
```

```
awk 'BEGIN {OFS = "\t"}; {if ($3=="gene") print $1,$4-1,$5}' | head -n 100
```

Many bioinformatics software support STDIN as input

```
bwa mem ref.fa reads.fq | samtools view -bS - > out.bam
```

Use "-" to specify input from STDIN instead of a file

```
gunzip -c D001.fastq.gz | fastx_clipper -a AGATCG
```

```
..... | bedtools coverage -a FirstFile -b stdin
```

BEDtools - An Example

A file: Gene Annotation

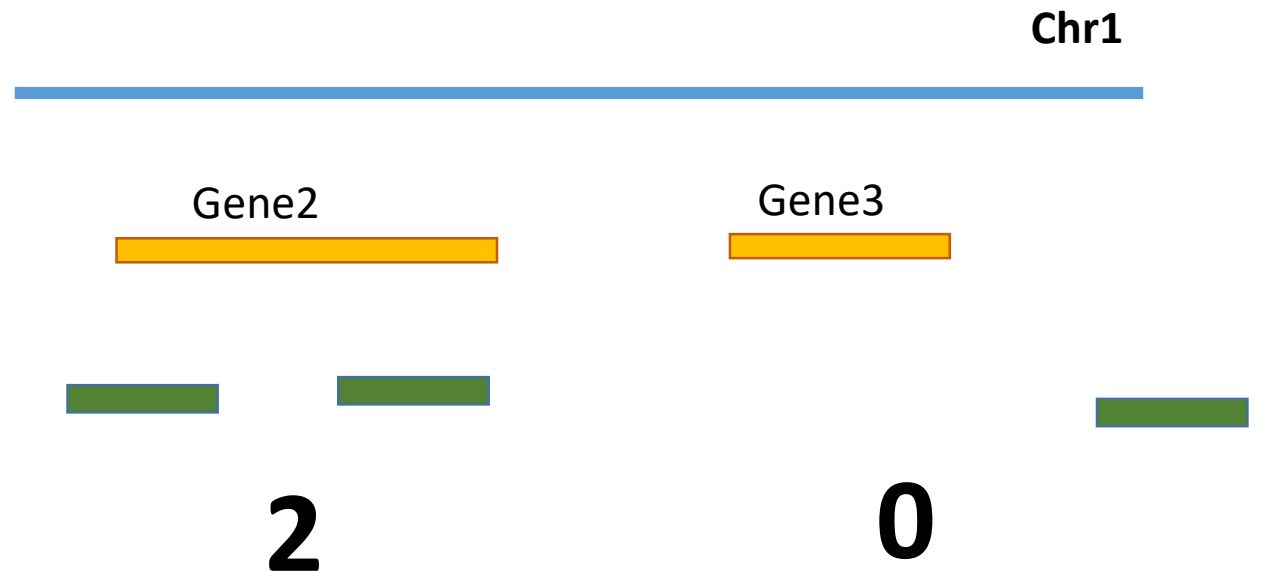
Chr1	1000	2250	Gene1	.	+
Chr1	3010	5340	Gene2	.	-
Chr1	7500	8460	Gene3	.	-
Chr2	8933	9500	Gene4	.	+
Chr2	12000	14000	Gene5	.	+

B file: Recorded Features

Chr1	200	300	Peak1
Chr1	4010	4340	Peak2
Chr1	5020	5300	Peak3
Chr2	8901	9000	Peak4

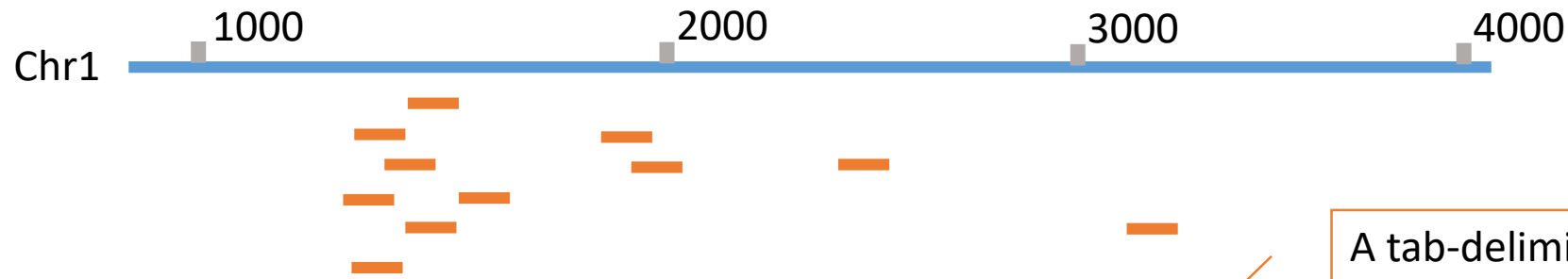
BEDtools coverage

		Number of overlapping peaks
...	Gene1	0
...	Gene2	2
...	Gene3	0
...	Gene4	1



Using BEDtools to process genomics data files

An example: Count the number of reads in each 1kb sliding window of the genome



A tab-delimited text file:

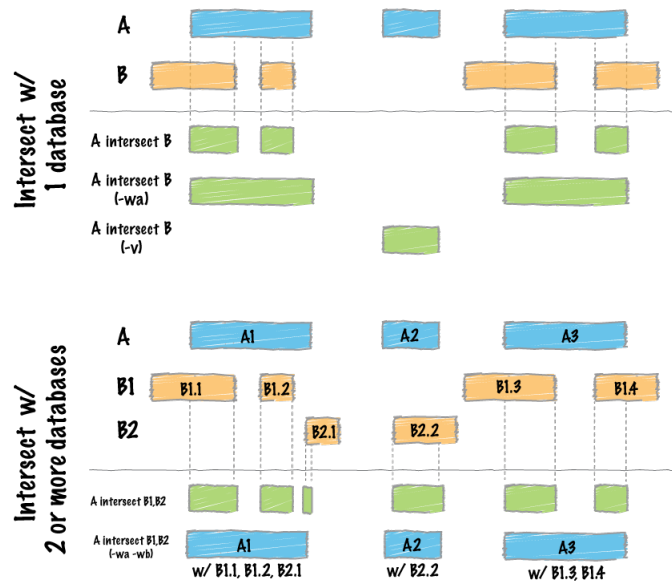
Chr1	21234023
Chr2	19282343
Chr3	15845499

```
bedtools makewindows -g genome.txt -w 1000 -s 1000 > win1000.bed
```

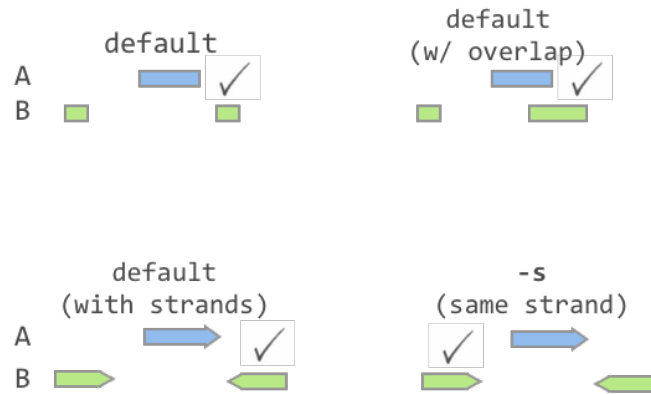
```
bedtools coverage -abam Sample.bam -b win1000.bed -counts> coverage.bed
```


There Are Many Functions in BEDtools that work with: BED, SAM/BAM, VCF, GFF, GTF

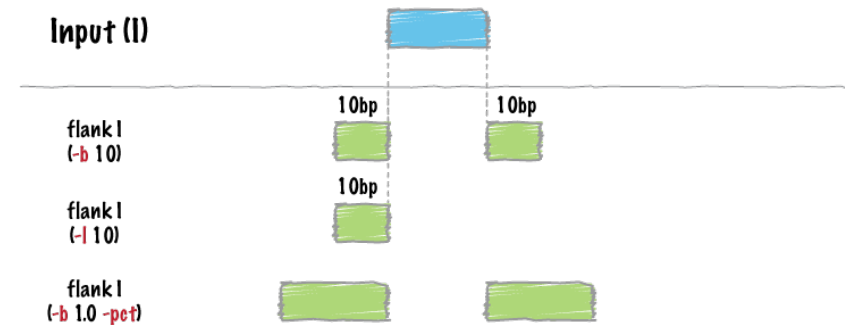
Intersect



Closest



Flanking



Using multi-processor machines

All BioHPC Lab machines feature multiple CPU cores

- ❑ general (cbsum1c*b*): 8 CPU cores
- ❑ medium-memory (cbsumm*): 24 CPU cores
- ❑ marge-memory (cbsulm*): 64+ CPU cores

Using multi-processor machines

Three ways to utilize multiple CPU cores on a machine:

- ❑ Using a given program's built-in parallelization, e.g.:

```
blast+ -num_threads 8 [other options]
```

```
bowtie -p 8 [other options]
```

Typically, all CPUs work together on a single task. Non-trivial, but taken care of by the programmers.

- ❑ Simultaneously executing several programs in the background, e.g.:

```
gzip file1 &
```

```
gzip file2 &
```

```
gzip file3 &
```

- ❑ If the number of independent tasks is larger than the number of CPU cores - use a “driver” program:

```
/programs/bin/perlscripts/perl_fork_univ.pl
```

Multiple independent tasks

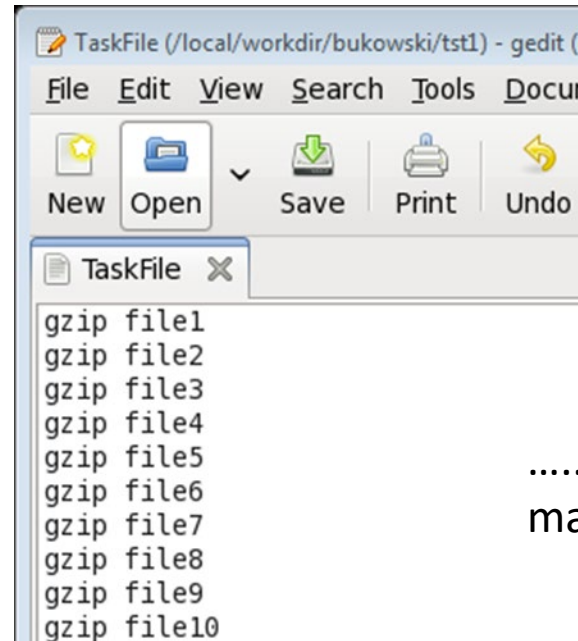
CPU: number of cores

RAM: not exceed the memory

DISK: over load the DISK IO

Using perl_fork_univ.pl

Prepare a file (called, for example, **TaskFile**) listing all commands to be executed. For example,



.....number of lines (i.e., tasks)
may be very large

Then run the following command:

```
/programs/bin/perlscripts/perl_fork_univ.pl TaskFile NP >& log &
```

where **NP** is the number of processors to use (e.g., 10). The file “log” will contain some useful timing information.

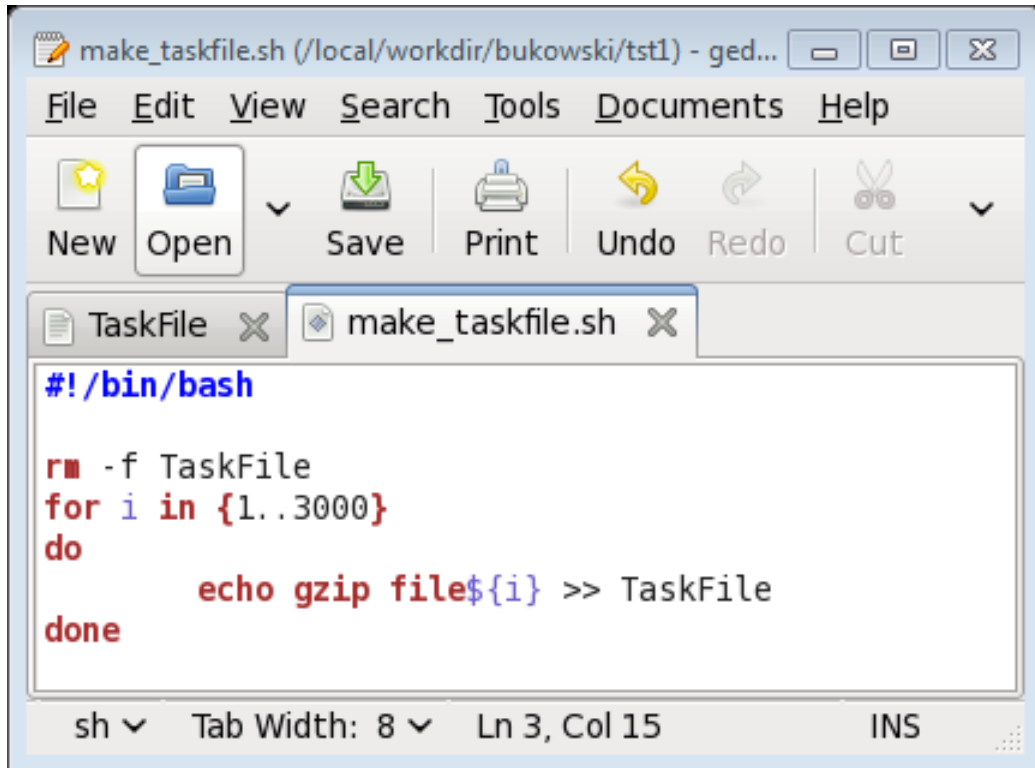
Using perl_fork_univ.pl

What does the script `perl_fork_univ.pl` do?

- ❑ `perl_fork_univ.pl` is an CBSU in-house “driver” script (written in perl)
- ❑ It will execute tasks listed in **TaskFile** using up to **NP** processors
 - The first **NP** tasks will be launched simultaneously
 - The **(NP+1)** th task will be launched right after one of the initial ones completes and a “slot” becomes available
 - The **(NP+2)** nd task will be launched right after another slot becomes available
 - etc., until all tasks are distributed
- ❑ Only up to **NP** tasks are running at a time (less at the end)
- ❑ All **NP** processors always kept busy (except near the end of task list) – **Load Balancing**

Using perl_fork_univ.pl

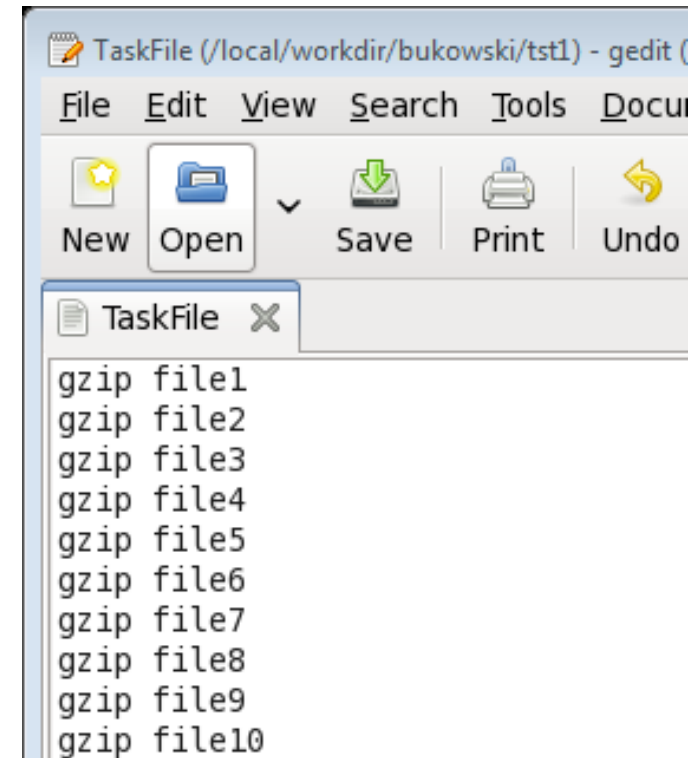
How to efficiently create a long list of tasks? Can use “loop” syntax built into bash:



```
#!/bin/bash
rm -f TaskFile
for i in {1..3000}
do
    echo gzip file${i} >> TaskFile
done
```



TaskFile



```
gzip file1
gzip file2
gzip file3
gzip file4
gzip file5
gzip file6
gzip file7
gzip file8
gzip file9
gzip file10
```

.....

Create and run a script like this, or just type directly from command line, ending each line with RETURN

How to choose number of CPU cores

Typically, determining the right number of CPUs to run on requires some experimentation.

Factors to consider

- total number of CPU cores on a machine: $NP \leq (\text{number of CPU cores on the machine})$
- combined memory required by all tasks running simultaneously should not exceed about 90% of total memory available on a machine; use **top** to monitor memory usage
- disk I/O bandwidth: tasks reading/write to the same disk compete for disk bandwidth. Running too many simultaneously will slow things down
- other jobs running on a machine: they also take CPU cores and memory!