

Time-saving workflows and easy parallel processing

Thomas Sibley – 21 May 2014 – Mullins Comp Bio Group

I'm going to talk about how to save time when working on computational problems. The first part will be about a method for building reusable workflows that are simple to run, able to grow from simple to complex as your project demands, and self-updating when source data changes. The second part will be about parallel processing using this workflow system and also other tools.

Workflows

- Pipelines
- Scripts
- Data cleanup
- Analysis steps
- Producing graphs and charts
- etc...

2

These are your computational or analytical methods in a paper.

Pretty much whatever you did to go from the raw data to your charts and graphs and final data tables is part of a workflow.

Yours may be usually entirely manual, but hopefully they're automated to some degree. I'm going to talk about automation.

Workflows

- Correct
- Reproducible
- Documented

- Easy is better than tedious
- Fast is better than slow

3

Workflows should be correct.

Workflows should be reproducible. Reproducible doesn't just mean people in other labs. It means people in your lab and even you, a few months or few years later.

Workflows should be documented. Documentation includes software versions, source code, data input/output, and more, but the steps you took to process and analyze the data is the huge one.

If all of those hold about your workflow, then you'd also like it to be easy and fast rather than tedious and slow. Not only do you not spend needless time waiting, but there are tangible benefits in terms of thinking critically about a problem when you can iterate quickly trying out new ideas or fixing problems in the analysis. If it takes you 8 hours to run an analysis, fixing problems is frustrating and you're afraid to make changes and try new things and experiment with different approaches to solving your problems.

Workflows

- Correct*
- Reproducible
- Documented

- Easy is better than tedious
- Fast is better than slow

* <http://software-carpentry.org/blog/2013/02/correctness-isnt-compelling.html>
<http://www.davidhbailey.com/dhbpapers/icerm-report.pdf>

I want to note a caveat: Correctness is obviously important, but it isn't compelling and there was a study to prove it. Incentives are all wrong as few papers are subjected to reproduction attempts.

However, making your analyses reproducible and self-documenting will save you time puzzling over what you did and make it easy to redo analysis on an updated or new data set.

When you do similar analyses in the future or when another member of the lab or your collaborators want to run the same analysis, you'll all save time not reinventing what you already did. And when you reuse what you did previously, you're more certain that it works.

make

- Produces (*makes*) files using recipes
- Recipes are plain text files named *Makefile*
- Language agnostic
- Only does the work necessary
- Stops on error

- Simple to start, allows complexity

5

The tool I'm going to introduce to you is `make`, an old, venerated command-line tool for building workflows. `Make` got its name because it *makes* files using recipes you write. The recipes are plain text in a file named *Makefile* and it's a language agnostic tool. That means it doesn't care if your programs are written in R or Python or Perl or Java or Shell or whatever.

`Make` is also "smart", in that it only does the work necessary and won't regenerate files if it doesn't have to. It also stops on errors in your recipes so you can fix them rather than bluster forward into the unknown.

It's simple to start using but allows for complexity, so it scales to the size of your project and is suited even for the simplest of tasks.

There *are* some sharp corners with `Makefiles` and `make`, but it's a time-tested tool and no software doesn't have sharp corners somewhere. I'll talk about some of the gotchas later.

Makefile recipes

- Filename(s) to make = *targets*
- Necessary input file(s) = *prerequisites*
- A set of commands to run = *actions*

6

A Makefile is a collection of recipes. Recipes just describe what to run, called the actions, the necessary input files, called prerequisites (also called dependencies), and what files are produced, called the targets. You can run core Unix commands, your own Python, R, or Perl scripts, and use features of your shell to pipe data between commands.

I often prototype Makefile recipes straight on the command line before putting them into a Makefile and further test and tweak from there.

The prerequisites are important because make uses those to determine what order to run the recipes. It also uses them to skip over recipes if the files they produce already exist and the input files haven't changed since then. This can save a lot of time when writing complex workflows with steps that take a while. Once you verify that the long running steps are correct, you can run them once and then move on to using them down the road without re-running it every time. Best yet, make will figure this out for you and you don't need to remember what's changed since you last ran it two weeks ago.

Makefile recipes

- Filename(s) to make = *targets*
- Necessary input file(s) = *prerequisites*
- A set of commands to run = *actions*

```
seqs_aa.fasta: seqs_na.fasta
    transeq -sequence seqs_na.fasta \
            -outseq seqs_aa.fasta \
            -frame 1 -clean
```

7

This is what a basic recipe looks like. It says that the file `seqs_na.fasta` is used to produce `seqs_aa.fasta` and how to do it using the command `transeq`. (`transeq` is a program in the EMBOSS toolkit.)

Note the line continuations in the actions. This says that these three lines will be run as a single command. Without those, each line would be a separate command (and wouldn't work because `-outseq` isn't a valid command).

To run this, we'd simply type: `make seqs_aa.fasta` and presto, `make` follows our recipe to produce it from `seqs_na.fasta`. Nifty, though I guess not very impressive yet.

Makefile recipes

- Filename(s) to make = *targets*
- Necessary input file(s) = *prerequisites*
- A set of commands to run = *actions*

```
seqs_aa.fasta: seqs_na.fasta  
    transeq -sequence $<  
            -outseq $@  
            -frame 1 -clean
```

8

This is an equivalent recipe which save some typing. It also provides you flexibility if you rename your targets or inputs later. It's especially useful for recipes which describe how to make multiple files, so I'll use these variables from now on. It's a good practice. Note that recipes can specify multiple targets or prerequisites and \$< and \$@ are only the first of the files being made/input. There are other variables to get all of them, which you'll see later.

I color code targets and prerequisites so it's easier to keep track.


```

seqs_aa.fa: seqs_na.fa
    transeq -sequence $< -outseq $@ \
            -frame 1 -clean

seqs_aa.aligned.fa: seqs_aa.fa
    muscle -quiet < $< > $@

seqs_aa_freq.tsv: seqs_aa.aligned.fa
    # CountAAFreq.pl only takes Nexus
    fasta2nexus < $< > seqs_aa.nxs
    perl CountAAFreq.pl seqs_aa.nxs $@ 0.25 0.5
    rm seqs_aa.nxs

```

9

Here's a longer example showing multiple targets that depend on each other. This takes a fasta of nucleotide sequences, translates it to amino acids, aligns it with muscle, and then runs it through Wenjie's CountAAFreq.pl. It also has to convert the aligned fasta to nexus. Note how each subsequent recipe depends on another. These are ordered sequentially, but they don't have to be and multiple recipes can depend on the same recipe.

To produce a table of amino acid frequencies from a seqs_na.fa file, we can now type: make seqs_aa_freq.tsv and make will follow all the necessary recipes.

It's also worth noting the command line's (shell's) input and output redirection operators, in orange. fasta2nexus, for example, takes an input stream and prints an output stream rather than taking an input filename and output filename itself. CountAAFreq.pl on the other hand takes two filenames along with some cutoff values (gaps and frequency). Do you see why \$< uses <?

Do you notice the temporary nexus file the last recipe creates and then deletes? We can do it that way, but it's a good practice to keep your recipes as short as possible to enable reuse and save time later. With make, this is easy! Let's see how you could modify the recipes.

```
seqs_aa.fa: seqs_na.fa
  transeq -sequence $< -outseq $@ \
    -frame 1 -clean

seqs_aa.aligned.fa: seqs_aa.fa
  muscle -quiet < $< > $@

seqs_aa.nxs: seqs_aa.aligned.fa
  fasta2nexus < $< > $@

seqs_aa_freq.tsv: seqs_aa.nxs
  perl CountAAFreq.pl $< $@ 0.25 0.5
```

10

I just split the fasta2nexus conversion into a separate recipe.

Note that our recipes are simpler by making a separate recipe to describe how to make the nexus file out of the aligned fasta.

We can take it a step further though since there's little use for both a nexus file and the aligned fasta. Why not just produce a nexus file to begin with?

```
seqs_aa.fa: seqs_na.fa
  transeq -sequence $< -outseq $@ \
    -frame 1 -clean

seqs_aa.nxs: seqs_aa.fa
  muscle -quiet < $< | fasta2nexus > $@

seqs_aa_freq.tsv: seqs_aa.nxs
  perl CountAAFreq.pl $< $@ 0.25 0.5
```

11

Now we're getting somewhere! We can pipe the output of muscle directly to fasta2nexus and convert it on the fly. Using pipes not only simplifies the recipe, but it's faster than writing a bunch of files.

Alright. We have a workflow to get amino acid frequencies from a set of nucleotide sequences. That's great!

But it only works for one file, and that's annoying. Ideally we'd like to generalize it so it works for any nucleotide fasta we have without renaming all our files to seqs_na.fa.

```

_aa.fa: _na.fa
    transeq -sequence $< -outseq $@ \
            -frame 1 -clean

seqs_aa.nxs: seqs_aa.fa
    muscle -quiet < $< | fasta2nexus > $@

seqs_aa_freq.tsv: seqs_aa.nxs
    perl CountAAFreq.pl $< $@ 0.25 0.5

```

12

Luckily, make supports this by writing recipes that use wildcards in the target and prerequisites. The % is a wildcard here and represents the same name on each side of the colon.

This means we can take any file named something_na.fa and produce a something_aa.fa just by typing: make something_aa.fa

That's pretty nifty, but notice that our other recipes are still hardcoded. Let's fix that.

```
%_aa.fa: %_na.fa
    transeq -sequence $< -outseq $@ \
            -frame 1 -clean

%.nxs: %.fa
    muscle -quiet < $< | fasta2nexus > $@

seqs_aa_freq.tsv: seqs_aa.nxs
    perl CountAAFreq.pl $< $@ 0.25 0.5
```

13

This recipe will now take any fasta and produce a nexus file. With just a seqs_na.fa file, you can now type: make seqs_aa.nxs and get an aligned amino acid sequences in a nexus file. Make will run the first rule if it needs to, and then run the second rule with the output from the first.

But note that there's no restriction on the filenames other than the extensions, so if want to align your nucleotide sequences instead and get a nexus file of those, you can also do: make seqs_na.nxs. The first rule won't be run since no seqs_aa.fa needs to be made.

Let's keep going to make the whole workflow generalized.

```

_aa.fa: _na.fa
    transeq -sequence $< -outseq $@ \
            -frame 1 -clean

_nxs: %.fa
    muscle -quiet < $< | fasta2nexus > $@

# make pic_aa_freq.tsv

_aa_freq.tsv: _aa.nxs
    perl CountAAFreq.pl $< $@ 0.25 0.5

```

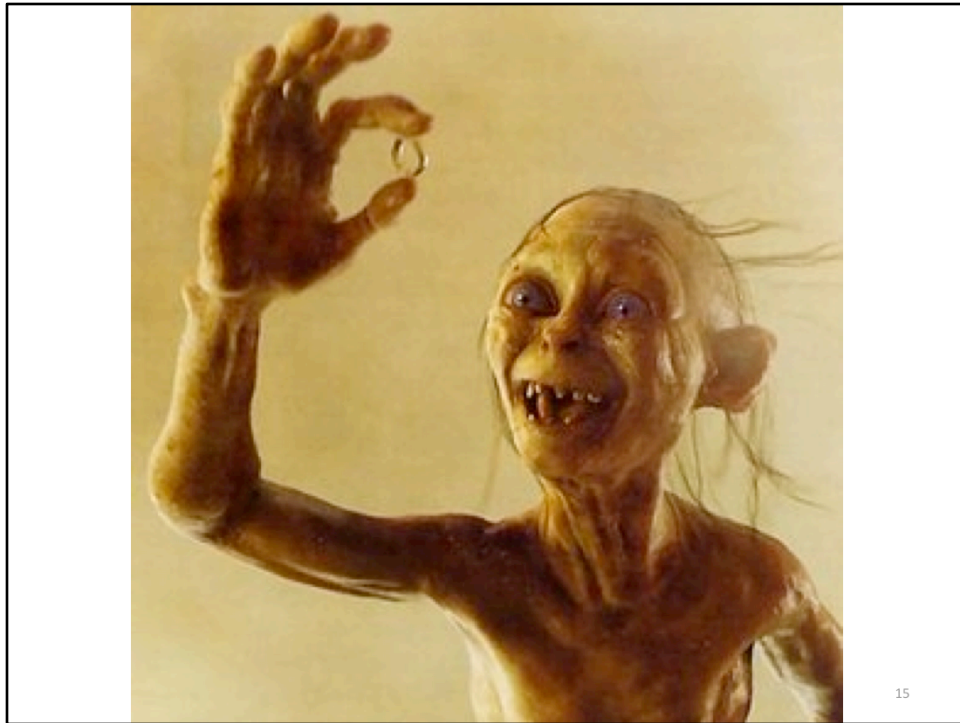
14

This is the last piece of the puzzle. Our rule for counting amino acid frequencies now describes how to take any nexus file of amino acid sequences (using the convention `_aa.nxs`) and runs it through Wenjie's program to produce a frequency table.

If I have a file `pic_na.fa`, I can now run: `make pic_aa_freq.tsv`. It's important to note that if I already have a `pic_aa.fa` file from somewhere else, I can still run ``make pic_aa_freq.tsv`` and make will realize it doesn't need to run the first rule to translate from nucleotides to amino acids.

It's also important to note that the filenames I'm using are just conventions. You can use whatever you want, for example, to distinguish between amino acid and nucleotide fastas, as long as you're consistent within your recipes.

Now, normally make will delete intermediate files after it's done with them. Intermediate files are any files you didn't ask for, but that it had to produce to get from your input to the output you asked for. In the case of going from `pic_na.fa` to `pic_aa_freq.tsv`, there are two intermediate files: `pic_aa.fa` and `pic_aa.nxs` which make will delete when it's done. This is just a cleanliness thing so you have less files to look at in your directory. But sometimes you want to keep those files around, especially if they take a while to produce. `muscle`, for example, might take a long time on a large set of sequences. There's a way to tell make that it shouldn't delete certain intermediate files, that certain files are... precious.



15

```

_aa.fa: _na.fa
    transeq -sequence $< -outseq $@ \
            -frame 1 -clean

.nxs: %.fa
    muscle -quiet < $< | fasta2nexus > $@

# Keep intermediate alignments, for speed
.PRECIOUS: %.nx

_aa_freq.tsv: _aa.nxs
    perl CountAAFreq.pl $< $@ 0.25 0.5

```

16

... using a special target name. The target “.PRECIOUS” does this and any prerequisites you specify won’t be deleted even if they’re intermediate files. Note that the prerequisites of .PRECIOUS should be other targets.

The gray text is a comment, which you can put in your Makefiles by starting a line with a hash or pound sign.

Advanced features

- Variables

```
NAME := Thomas
hello:
    echo 'Hi, my name is $(NAME).'
```

- Using \$ signs in your recipes

```
check_balance:
    echo 'Your balance is $$17.03.'
```

- Targets don't have to be files

17

Make supports a number of more advanced features too, and I'll cover a few now.

Arbitrary variables are one of those features and they help reduce repetition in directory names or other commonly used parameters. Variable names longer than a single character need to be surrounded by parentheses. \$< and \$@ are just automatic variables. := is the assignment operator, and surrounding whitespace doesn't matter.

Since dollar signs introduce a variable in Makefiles, to use an actual dollar sign you need to type it twice.

From these two examples, you can see that targets don't have to be files. Make doesn't create a target file itself, that's up to the recipe. So targets may just be a convenient name for a recipe to run a bunch of commands that don't actually produce a file.

Advanced features

- Recipes don't have to have actions

```
all: gag_aa_freq.tsv env_aa_freq.tsv
```

- Prerequisites don't have to be recipes

```
%_aa_freq.tsv: %_aa.nxs CountAAFreq.pl  
perl CountAAFreq.pl $< @$ 0.25 0.5
```

18

You can also create targets whose sole purpose is to list a bunch of other targets as prerequisites, which is a way of running multiple targets at once which don't depend on each other. In this case, the target "all" will produce a number of specific files from our previous set of generalized recipes.

Prerequisites also don't have to be other recipes. Remember that prerequisites are just the targets or files that a recipe needs to run and that when they change the recipe needs to be re-run. You can list your own programs as prerequisites and then make will know it needs to regenerate the files the next time you ask for them after updating your program. Make will know when you fix bugs or change your analysis change!

Advanced features

- Recipes don't have to have actions

```
all: gag_aa_freq.tsv env_aa_freq.tsv
```

- Prerequisites don't have to be recipes

```
%_aa_freq.tsv: %_aa.nxs CountAAFreq.pl  
perl CountAAFreq.pl $< $@ 0.25 0.5
```

```
%_aa_freq.tsv: CountAAFreq.pl %_aa.nxs  
perl $^ $@ 0.25 0.5
```

19

You could write that last recipe like this too.

`$^` is an automatic variable that lists all the prerequisites, separated by spaces. In this case, it'll expand to the script name and the input filename to make the first two arguments to `perl`.

There are other variables too, such as `$*` which is just the shared stem, or wildcard part, of the filename.

Makefile gotchas

- Hard tabs vs. spaces
- Force updates after changing Makefile
- Change default behavior on error:

```
SHELL := /bin/bash
export SHELLOPTS := errexit:pipefail
.DELETE_ON_ERROR:
```

20

Let's talk about some gotchas when using Makefiles. The most common by far is that the action part of a recipe must use hard tabs for the first indent, not spaces. This is often a source of problems. However, all editors should have a way of highlighting hard tabs vs. spaces.

The next is that make doesn't regenerate files if the input hasn't changed but the recipe itself has. When a Makefile changes, you often need to rerun the recipes. Since the timestamps of the input and output files don't change, just running make won't do that. To get around this, you can run `make -B` to force run a target and all dependent targets. You can also update the timestamps of all your input files using the `touch` command and then rerun your targets with make.

Finally, make's default behaviour on errors is sometimes less than ideal if you're using pipelines, such as when we piped the output of muscle to fasta2nexus. Only the success or failure status of the last command in a pipeline is considered for errors, even if a command in the middle fails partway through the data. You can change this by using the first two lines here which change the command shell make uses to run your recipes and sets an option for the shell. When make *does* catch an error, it leaves any partially made target files around. This can be confusing at first if you don't notice there was an error. You can include the special empty target `".DELETE_ON_ERROR:"` to force make to delete any partially-complete target files if the recipe fails. This avoids manually running other recipes later which may use the partial data left around.

Parallel processing

How to do more than one thing at a time

21

Alright, now it's time to talk about parallel processing, or how to have the computer do more than one thing at a time.

What can be parallelized?

- Easiest is separate input files
- Single input files can be split up
- Independent tasks
- All prerequisites available

22

The easiest thing to parallelize is when you have a bunch of separate input files to do the same thing to and it takes a while to do it. But that said, you can also split up a single input file if you're running tasks over individual things in that input file. For example, and I'll show this in a minute with blast, you might be doing something to every sequence in a single file.

In any multi-step workflow, you must also think about dependencies and what steps of your workflow depend on other steps. For example, you can't run a step expecting amino acids as input when you haven't translated your nucleotide sequences yet. Hey wait... that sounds familiar. We just thought about this for Makefiles!

```

_aa.fa: _na.fa
    transeq -sequence $< -outseq $@ \
            -frame 1 -clean

.nxs: %.fa
    muscle -quiet < $< | fasta2nexus > $@

# Keep intermediate alignments, for speed
.PRECIOUS: %.nx

_aa_freq.tsv: _aa.nxs
    perl CountAAFreq.pl $< $@ 0.25 0.5

```

23

Remember our Makefile from earlier? It calculates amino acid frequencies from a set of nucleotide sequences. We generalized it so that filenames aren't hardcoded and it'll operate on anything named properly.

The work we put in to describe our workflow in a Makefile will pay back again, because make can parallelize your workflow for you.

Parallelizing with make

`gag_na.fa`

`env_na.fa`

`make gag_aa_freq.tsv env_aa_freq.tsv`

24

Given two files, `gag_na.fa` and `env_na.fa`, you'd get amino acid frequencies like this.

Parallelizing with make

```
gag_na.fa
```

```
env_na.fa
```

```
make gag_aa_freq.tsv env_aa_freq.tsv
```

```
make {gag,env}_aa_freq.tsv
```

25

To save yourself some time, you might instead use the shell's combination or alternation syntax. The second make command does the exact same thing as the first.

Parallelizing with make

- What if you had those files for 20 patients?

Pt100_gag_na.fa, Pt100_env_na.fa

Pt101_gag_na.fa, Pt101_env_na.fa

...

Parallelizing with make

- What if you had those files for 20 patients?

Pt100_gag_na.fa, Pt100_env_na.fa

Pt101_gag_na.fa, Pt101_env_na.fa

...

```
make Pt{100,101,...}_{gag,env}_aa_freq.tsv
```

27

Maybe you'd do something like this, listing out all the patient ids. Make is going to process those one after the other though, and that'll take a while. Probably not too long for our simple example, but consider if the alignments were large or you were doing a more computationally intensive task.

Since you described your workflow in a Makefile, make knows what output files require what input files and how to put it all together. All you need to do is tell make how many jobs to run in parallel!

Parallelizing with make

- What if you had those files for 20 patients?

Pt100_gag_na.fa, Pt100_env_na.fa

Pt101_gag_na.fa, Pt101_env_na.fa

...

```
make --jobs=24 \  
    Pt{100,101,...}_{gag,env}_aa_freq.tsv
```

28

make takes a `--jobs` or `-j` option which tells it how many jobs, or tasks, to run at one time. By default it only runs one recipe at a time, but it'll happily run as many as you tell it. Now instead of waiting for `Pt100_gag_aa_freq.tsv` to finish before starting `Pt100_env_aa_freq.tsv`, make will run up to 24 tasks at once! That should cut down on the runtime a lot.


And... it's really that easy if you have a Makefile. It's one of the benefits of describing your workflows with make. Like any software, of course, it's possible for it to break if you missed a prerequisite input file for a recipe. If you're running individual steps by hand, you might not notice until you try to run them in parallel because previously the file from recipe 1 always happened to exist before you ran recipe 2. If you didn't tell make that, however, it might run recipe 2 before recipe 1 when running in parallel.

Note that you shouldn't ask to run more jobs than the number of cores if the tasks are computationally intensive (generally the case for bioinformatics). `themis` has 24 cores, your computer might have 1 or 2 or 4. Other servers will have different numbers. This leads to the question...

How many cores?

- Apple menu → About This Mac → More Info → System Report → “Total number of cores”

Model Name:	iMac
Model Identifier:	iMac10,1
Processor Name:	Intel Core 2 Duo
Processor Speed:	3.06 GHz
Number of Processors:	1
Total Number of Cores:	2



- `lscpu | grep -E 'Core|Socket'`

Core(s) per socket: 12

Socket(s): 2

29

...“How many cores do I have?” You can check yourself on Mac or Linux. The Mac example is from my desktop and the Linux example is from themis.

On Linux, each “socket” is a physical CPU, so multiply Cores per socket by the number of Sockets to get the total number of cores.

I've got 99 problems, but a
Makefile ain't one.

Parallel processing power tools

30

That's about all there is to say about parallel processing with make. The work is frontloaded to writing recipes rather than figuring out yourself how to parallelize it later.

But what if you're not using Makefiles?

No Makefile? No problem.

```
for file in *.fasta; do  
    do_something -in $file -out $file.new  
done
```

31

Maybe you have a command-line loop like this that you run over data files. It's straightforward to turn into a Makefile, but you don't want to. You may know about using an ampersand to run `do_something` in the background which makes the loop complete quickly, and then you wait around for all the `.new` files to pop into existence. That's fine for a handful of files, but if you have more than a couple dozen files, you'll bog down the computer with too many jobs. And if something goes wrong, you may have runaway processes chewing up time.

No Makefile? No problem.

```
for file in *.fasta; do
  do_something -in $file -out $file.new
done

parallel do_something -in {} -out {}.new \
 ::: *.fasta
```

32

Enter a tool called “parallel”. It helps you run other programs in parallel and does a lot to handle input and output to each job. It’s the Leatherman of command-line parallel processing. I’d compare it to a Swiss Army knife, but that’s really not fair to the knife.

In it’s most basic form, parallel just replaces your loops. In more complicated forms, it can split up your input and rejoin the output back together.

Like make, it also has a -j option, but by default it will use as many jobs as the computer has cores. That’s handy! You can run the same command on your desktop as the server and it’ll just magically go faster on the server without thinking about the number of cores involved.

Let me show you more complicated example that will hopefully be immediately useful...

Parallel NCBI BLAST+

```
parallel \  
  --halt 2 \  
  --recstart '>' -N1 \  
  --pipe \  
  blastn \  
    -task blastn \  
    -db ./db/nucleotide/viroverse \  
    -query - \  
    -outfmt 6 \  
    -max_target_seqs 25 \  
< input.fa \  
> results.tsv
```

33

Here's an example parallel invocation I put together to very efficiently run blastn on themis against a copy of the Viroverse blast database (the same thing used by the local Viroblast).

I'll walk through it piece by piece, but first, why should you care? Well, blasting 245 whole and half HIV genomes against Viroverse using our local Viroblast took an hour. Doing the same blast run on themis using this command takes 3 minutes. That's a 20-fold decrease.

Parallel NCBI BLAST+

```
parallel \  
  --halt 2 \  
  --restart '>' -N1 \  
  --pipe \  
  blastn \  
    -task blastn \  
    -db ./db/nucleotide/viroverse \  
    -query - \  
    -outfmt 6 \  
    -max_target_seqs 25 \  
< input.fa \  
> results.tsv
```

34

The first thing I want to note is actually at the end. The file `input.fa` is redirected as the input to `parallel`, and `parallel`'s output is redirected to another file, `results.tsv`.

Parallel NCBI BLAST+

```
parallel \  
  --halt 2 \  
  --recstart '>' -N1 \  
  --pipe \  
  blastn \  
    -task blastn \  
    -db ./db/nucleotide/viroverse \  
    -query - \  
    -outfmt 6 \  
    -max_target_seqs 25 \  
< input.fa \  
> results.tsv
```

35

The next thing to note is what parallel is doing with that input for each job. The `--recstart` option tells parallel to split up the input into multiple “records”. In this case, I’m telling parallel that records start with a “>”, which should be familiar to you as the start of a FASTA sequence. I’m also telling parallel with the `-N1` option to only pass one record at a time to blast. This means we’ll run one blast job for every sequence, but the number of jobs running at the same time is still limited by the number of cores.

Parallel NCBI BLAST+

```
parallel \  
  --halt 2 \  
  --restart '>' -N1 \  
  --pipe \  
  blastn \  
    -task blastn \  
    -db ./db/nucleotide/viroverse \  
    -query - \  
    -outfmt 6 \  
    -max_target_seqs 25 \  
< input.fa \  
> results.tsv
```

36

Since parallel is splitting up the input file, we no longer have a filename for blast to process. Instead, we tell parallel to pipe the query sequence to blast as a standard input stream. And we also tell blast that the query sequences are from the input stream. The dash in place of a filename is a Unix convention to specify “stdin” rather than an actual filename.

Parallel NCBI BLAST+

```
parallel \  
  --halt 2 \  
  --restart '>' -N1 \  
  --pipe \  
  blastn \  
    -task blastn \  
    -db ./db/nucleotide/viroverse \  
    -query - \  
    -outfmt 6 \  
    -max_target_seqs 25 \  
< input.fa \  
> results.tsv
```

37

Finally we have the blast command itself. We're using blastn, specifying a database, choosing a suitable output format and limiting the maximum hits to 25 per query. You can use whatever options make sense for your query and database, such as to adjust the blast scoring or output format.

Output format 6 is tabular data and blastn will print it to stdout, which parallel will capture and print as jobs finish. It's what we redirect into results.tsv.

When parallelizing anything, it's important to choose an output format which can be joined together easily. There are options for joining output formats that can't just be smushed together, but explaining those is for another day. If you do find yourself in this boat, you can always have each blast job produce its own result file instead of joining them all together.

Parallel NCBI BLAST+

```
parallel \  
  --halt 2 \  
  --restart '>' -N1 \  
  --pipe \  
  blastn \  
    -task blastn \  
    -db ./db/nucleotide/viroverse \  
    -query - \  
    -outfmt 0 -out 'results-{}.blastn' \  
    -max_target_seqs 25 \  
< input.fa
```

38

For example, you could do this to get a pairwise output file (format 0) for each blast job. Parallel substitutes {} for the input record number, so you'll get as many files as you have input sequences. Note that we're no longer redirecting parallel's output to a file.

Parallel NCBI BLAST+

```
parallel \  
  --halt 2 \  
  --restart '>' -N1 \  
  --pipe \  
  blastn \  
    -task blastn \  
    -db ./db/nucleotide/viroverse \  
    -query - \  
    -outfmt 0 -out 'results-#{#.blastn' \  
    -max_target_seqs 25 \  
< input.fa
```

39

The last thing to point out is the `--halt 2`, which I skipped over. It tells parallel to abort all jobs and quit with an error if there's a problem running blastn. This makes sure we don't miss problems in the sea of other results. `--halt 2` is what I recommend, but there are other error handling behaviours you can use too.

Getting these tools

- parallel
 - Installed on themis
 - I can install it elsewhere or help you do it
- make
 - Part of Apple's Xcode
 - In the Mac App Store

40

That's about all I have today. To get these tools, you can use parallel on themis or I can help you install it on your desktop. To get make on OS X, you just install Apple's Xcode. You might already have this, so try running make first. Newer versions of OS X may ask if you want to install it if it's not already installed.

Resources

- Manuals
 - <http://www.gnu.org/software/make/manual/make.html>
 - http://www.gnu.org/software/parallel/parallel_tutorial.html
- `man make`
- `man parallel`
- PDF of this talk, with notes, on the wiki
- Ply me with donuts, or just ask nicely

41

For more info, there are online manuals and also the documentation included on any system with the tools.

The make manual is a bit dense, but it's the authoritative reference for how things work. It contains useful things such as descriptions of the automatic variables like `$@`, `$<`, `$*` and more.

The parallel tutorial is a little friendlier and better yet is chock full of examples showing how it works. There are other substitution patterns than just the two I showed, for example, many of which are often useful.