## A. The 'genome model tools' (gmt) command

```
●●●                 ⌂ ssmith — ssmith@blade13-4-4: ~ — ssh — 118×53                    ↗
        ssmith@blade13-4-4: ~            0:1.0 bash        ssh ...
ssmith@blade13-4-10 ~> gmt
Sub-commands for gmt:
  smalt                        Tools to run smalt or work with its output files.
  tophat                       Tools to run Tophat or work with its output files.
  ssaha2                       Tools to run SSAHA2 or work with its output files.
  bfast                        Tools to run Bfast or work with its output files.
  tigra-sv                     Tools to run tigra_sv or work with its output files.
  blat                         Tools to run Blat or work with its output files.
  bwa                          Tools to run BWA or work with its output files.
  crossmatch                   Tools to run Crossmatch or work with its output files.
  picard                       Tools to run the Java toolkit Picard and work with SAM/BAM
                                 format files.
  novocraft                    Tools to run novocraft or work with its output files.
  bsmap                        Tools to run BSMAP or work with its output files.
  mosaik                       Tools to run Mosaik or work with its output files.
```

## B. Each top-level command provides access to a list of tools, or further sub-trees

```
  bowtie                       tools to work with the Bowtie aliger
  breakdancer                  discovers structural variation using breakdancer
  bwa-sw                       tools to work with Ssaha output
  chimera-slayer               Tool to run chimera detector: chimera_slayer
  complete-genomics        ... base class for commands which delegate to sub-commands
  copy-cat                 ... the CopyCat copy number analysis tools
  cufflinks                    Tools to run Cufflinks or work with its output files.
  detect-variants              A selection of variant detectors.
  dgidb                    ... Toolkit for DGIDB related process
  ensembl                  ... Tools to work with the local Ensembl API.
  epitope-prediction       ... Different pipeline steps for Immune Epitope Prediction for
  far                      ... To trim adaptor sequences
  fasta                        Tools for working with FASTA and Qual files
  fastq                        tools for working with FASTQ files
  fastqc                       Tools to run the Java toolkit FastQC and work with the
                                 output reports.
  fastx                        Tools to run Fastx or work with its output files.
  galaxy                   ... the Galaxy web interface
  gatk                         tools to work with Gatk output
  gene-torrent                 no description!!!: define 'doc' in the class definition for
                                 Genome::Model::Tools::GeneTorrent
  gtf                      ... Tools to work with gtf format annotation files.
  htseq                    ... htseq tools (htseq-count and htseq-qa) work with
                                 gene/transcript hit-counts
  lift-over                    wrapper for the UCSC liftOver tool with support for
                                 additional input formats, maintaining additional columns
```

C. The 'gmt fasta' sub-tree contains script-like components for working with FASTA files

```
● ● ●                                    🗀 ss — bash — 116×45
ssmith@blade12-1-1 > gmt fasta
Sub-commands for genome tools fasta:
 apply-diff           --diff=? --input=... applies seq inserts and deletes from a diff file to a fasta
                                            file
 chunk                --chunk-size=? --... Divide fasta into chunk by chunk_size
 concat               --input-files=?[,... Mixin that gives commands color option
 deduplicator         --fasta-file=? [-... remove duplicates from a file of reads
 diff                 [--debug] FILE1 F... use KDiff3 to show differences between fasta data files
 dust                 --dusted-file=? -... Tools for working with FASTA and Qual files
 filter-ids           [--verbose] [--wh... filter sequences from a fasta file based on patterns
                                            applied to the IDs
 orient               --fasta-file=? [-... Orients FASTA (and Quality) files by blastn given sense and
                                            anti-sense sequences
 remove-n             --fasta-file=? [-... remove reads from file containing N
 sanitize             --fasta-file=? [-... Cleans FASTA (and Quality) files
 screen-vector        --fasta-file=? [-... (Fnq = Fasta And Quality) screen for vector
 sliding-windows      --fasta-file=? --... Tools for working with FASTA and Qual files
 sort-by-name         --input-fasta=? [... Sorts a fasta by sequence name
```

D. Each tool has auto-generated help

```
● ● ●          🏠 ssmith — ssmith@linus43: ~ — ssh — 100×40
        ssmith@linus43: ~                          bash

ssmith@linus43 ~> gmt fasta filter-ids -h

USAGE
  gmt fasta filter-ids [--verbose] [--whitelist-regex=?] [--blacklist-regex=?]
      INPUT-FILENAME OUTPUT-FILENAME

SYNOPSIS
      gmt fasta filter-ids in.fa out.fa --whitelist '^(\d+|X,Y)$' --blacklist '6'

REQUIRED INPUTS
   INPUT-FILENAME
      the input file
   OUTPUT-FILENAME
      the path to the file that will be created

OPTIONAL INPUTS
   whitelist-regex
      include only IDs that match this pattern
   blacklist-regex
      exclude any IDs that match this pattern

OPTIONAL PARAMS
   verbose
      more messages

DESCRIPTION
      This tool filters a FASTA sequence file, removing entries based on the ID in the FASTA header.
      If the "whitelist regex" (-w) option is supplied, only IDs that match this regular expression
      will be included. If the "blacklist regex" (-b) option is supplied, only IDs that do NOT match
      this regular expression will be included.


NOTE

      If an ID matches both the black list and the white list, it is skipped.


ssmith@linus43 ~> ▮
```

E. The code for a GMT tool can be as simple as a short script

```perl
1  package Genome::Model::Tools::Fasta::FilterIds;
2  use strict;
3  use warnings;
4  use Genome;
5  use Bio::SeqIO;
6
7  class Genome::Model::Tools::Fasta::FilterIds {
8      is => 'Command::V2',
9      has_input => [
10         input_filename => {
11             is => 'FilesystemPath', shell_args_position => 1,
12             doc => 'the input file',
13         },
14         output_filename => {
15             is => 'FilesystemPath', shell_args_position => 2,
16             doc => 'the path to the file that will be created',
17         },
18         whitelist_regex => {
19             is => 'Text', is_optional => 1,
20             doc => 'include only IDs that match this pattern',
21         },
22         blacklist_regex => {
23             is => 'Text', is_optional => 1,
24             doc => 'exclude any IDs that match this pattern',
25         },
26     ],
27     has_param => [
28         verbose => {
29             is => 'Boolean', is_optional => 1,
30             doc => 'more messages'
31         },
32     ],
33     doc => "filter sequences from a fasta file based on patterns applied to the IDs",
34 };
35
36 sub execute {
37     my $self = shift;
38     my $input_filename = $self->input_filename;
39     my $output_filename = $self->output_filename;
40     my $verbose = $self->verbose;
41
42     my $blacklist_regex = $self->blacklist_regex;
43     my $whitelist_regex = $self->whitelist_regex;
44
45     my $reader = Bio::SeqIO->new( '-file' => '< '.$input_filename, '-format' => 'fasta');
46     my $writer = Bio::SeqIO->new( '-file' => '> '.$output_filename, '-format' => 'fasta');
47
48     while (my $seq = $reader->next_seq) {
49         my $id = $seq->id;
50         if ($blacklist_regex and $id =~ $blacklist_regex) {
51             $self->status_message("skipping $id because it matches the blacklist pattern");
52             next;
53         }
54         elsif ($whitelist_regex and not $id =~ $whitelist_regex) {
55             $self->status_message("skipping $id because it does not match the whitelist pattern");
56             next;
57         }
58         elsif ($verbose) {
59             $self->status_message("keeping $id");
60         }
61         $writer->write_seq($seq);
62     }
63     return 1;
64 }
65
```

F. Additional code can be added to the module to explicitly or dynamically generate other documentation

```perl
65
66 sub help_synopsis {
67 return <<'EOS'
68     gmt fasta filter-ids in.fa out.fa --whitelist '^(\d+|X,Y)$' --blacklist '6'
69 EOS
70 }
71
72 sub help_detail {
73     return <<EOS
74 This tool filters a FASTA sequence file, removing entries based on the ID in the FASTA header.
75 If the "whitelist regex" (-w) option is supplied, only IDs that match this regular expression will be included.
76 If the "blacklist regex" (-b) option is supplied, only IDs that do NOT match this regular expression will be included.
77 EOS
78 }
79
80 sub _additional_help_sections {
81     return (
82         "NOTE" =>
83         "If an ID matches both the black list and the white list, it is skipped."
84   );
85 }
86
87 sub _doc_manual_body {
88     # expect to return POD
89     my $help = shift->help_detail;
90     $help =~ s/\n+$/\n/g;
91     return $help;
92 }
93
94 sub _doc_authors {
95     return <<EOS
96  Scott Smith
97  Edward Belter
98 EOS
99 }
100
101 sub _doc_copyright_years { (2013) }
102
103 sub _doc_license {
104     my $self = shift;
105     my (@y) = $self->_doc_copyright_years;
106     my $range;
107     if (@y == 1) { $range = "$y[0]"; }
108     elsif (@y > 1) { $range = "$y[0]-$y[-1]"; }
109     return <<EOS
110 Copyright (C) $range Washington University in St. Louis.
111
112 It is released under the Lesser GNU Public License (LGPL) version 3.  See the
113 associated LICENSE file in this distribution.
114 EOS
115 }
116
117 sub _doc_credits {
118     return ('','This software was created with funding fromthe National Human Genome Research Institute.');
119 }
120
121 sub _doc_see_also {
122     return <<EOS
123 B<Genome>(3)
124 EOS
125 }
126
127 1;
```