

# sequencer, v. 1.15: Simulate Shotgun Sequencing

Bernhard Haubold

Max-Planck-Institute for Evolutionary Biology

October 27, 2017

## Introduction

The program `sequencer` simulates shotgun sequencing. It takes as input a set of sequences and produces as output a collection of random reads. These can then be used, for example, to test assembly programs.

## Getting Started

The `sequencer` software was written in standard C on a UNIX computer. It is operated from the command line. To start using the software, carry out the following steps:

- Unpack program

```
tar -xvzf sequencer_xxx.tgz
```

where `xxx` denotes the current version.

- Change into newly created directory

```
cd Sequencer_xxx
```

- Print program options

```
./sequencer -h
```

- Test program

```
./sequencer phage_lambda.fasta
```

## Listing: `sequencer.c`

```
1  **** sequencer.c ****
 * Description: Program for simulating shotgun sequencing.
 * Author: Bernhard Haubold, haubold@evolbio.mpg.de
 * Date: Tue Jun 22 17:30:49 2004.
 * License: GNU General Public
6   ****
#include <stdio.h>
#include <time.h>
#include <string.h>
#include <fcntl.h>
11 #include <stdlib.h>
```

```

#include <math.h>
#include <unistd.h>
#include "sequence_data.h"
#include "eprintf.h"
16 #include "interface.h"
#include "ran.h"

void
run_sequencer (int fd, char *input_file, long *read_id, Args *args);
21 void process_read(FILE *fp, char *seq, long seqLen, long seqPos, int
    readLen, double err, int *dic, char *dicStr, int circular, int **
    profiles);

int
main (int argc, char *argv[])
{
    Args *args;                                /* arguments */
26    char *version;                            /* program version */
    int i, fd;
    long read_id;
    int idum, err;
    FILE *fpra;

31    version = "1.15";
    args = get_args (argc, argv);

    if (args->h == 1)
36    {
        print_usage ();
        return 0;
    }
    else if (args->v)
41    {
        print_version (version);
        return 0;
    }
/* seed for random number generation */
46    if (args->s != 0)
        idum = args->s;
    else if ((fpra = fopen ("random_seed.dat", "r")) != NULL)
    {
51        if ((err = fscanf (fpra, "%d", &idum)) != 0)
            idum = -time (NULL);
        fclose (fpra);
    }
    else
        idum = -time (NULL);
56    init_genrand (idum);

    read_id = 0;
    if (args->num_input_files == 0)
    {
61        fd = 0;
        run_sequencer (fd, "stdin", &read_id, args);
    }

```

```

    else
    {
        for(i = 0; i < args->num_input_files; i++)
        {
            if ((fd = open (args->input_files[i], O_RDONLY, 0)) <= 0)
                eprintf ("ERROR-[sequencer]: cannot open file %s for reading\n"
                         , args->input_files[i]);
            run_sequencer (fd, args->input_files[i], &read_id, args);
            close(fd);
        }
        if (args->s == 0)
        {
            fpra = efopen ("random_seed.dat", "w");
            fprintf (fpra, "%d\n", (int) genrand_int32 ());
            fclose (fpra);
        }
        free (args);
        free (progname ());
        return 0;
    }

/* process_read: process the nucleotide sequence of a sequencing read
   defined by
   * seq: template sequence
   * seqLen: length of the template sequence
   * seqPos: position of the read on the template sequence
   * readLen: length of the sequencing read
   * err: sequencing error per nucleotide
   * dic: map nucleotides to integers
   * dicStr: map integers to nucleotides
   * circular: 1|0 to indicate whether or not the template sequence is
     circular
   * profiles: if this is NULL, print nucleotides, else, count up profile
*/
void process_read(FILE *fp, char *seq, long seqLen, long seqPos, int
    readLen, double err, int *dic, char *dicStr, int circular, int **
    profiles){
long i, n, nucl, nuc2;

n = seqPos + readLen <= seqLen ? readLen : seqLen - seqPos;
for(i=0;i<n;i++){
    nuc1 = dic[(int) seq[seqPos + i]];
    nuc2 = nucl;
    /* sequencing error */
    if (genrand_reali () <= err)
        while ((nuc2 = (int)(genrand_reali () * 4)) == nuc1)
            ;
    if(profiles)
        profiles[seqPos + i][nuc2]++;
    else
        fprintf (fp, "%c", dicStr[nuc2]);
}

```

```

116    if(circular){
117        n = readLen - n;
118        for(i=0;i<n;i++){
119            nuc1 = dic[(int) seq[i]];
120            nuc2 = nuc1;
121            if (genrand_reall () <= err)
122                while ((nuc2 = (int)(genrand_reall () * 4)) == nuc1)
123                    ;
124            if(profiles)
125                profiles[i][nuc2]++;
126            else
127                fprintf(fp,"%c", dicStr[nuc2]);
128        }
129    }
130    if(!profiles)
131        fprintf(fp,"\\n");
132 }

133 /*
134 * run_sequencer: This function is still rather confused. I find it very
135 * hard to integrate paired-end
136 * reads in the current structure while still keeping the option to
137 * generate profiles instead of
138 * sequencing reads. The following steps need to be disentangled:
139 * - Get fragment to sequence; this should perhaps be done explicitly:
140 *   input sequence and output fragment as character array
141 * - Determine length of first read and either
142 *   output read with error
143 *   or
144 *   count nucleotides for profiles with error
145 * - If reads are paired
146 *   determine length of second read and either
147 *   output read with error
148 *   or
149 *   count nucleotides for profiles with error
150 */
151 void
152 run_sequencer (int fd, char *input_file, long *read_id, Args *args)
153 {
154     Sequence *seq; /* sequence data */
155     long residues;
156     long sum, coverage, p, pf, printed_len, i, j;
157     long seqi, seq_ind, max_res;
158     long frag_len, read_len;
159     long *num_res;
160     int **profiles = NULL, *dic = NULL;
161     long start, end;
162     FILE *fout;
163     char *sequence[2], strand;
164     char **forward, **reverse;
165     char *dic_str = "ACGTN";
166
167     seq = NULL;
168     seq = read_fasta (fd);

```

```

fout = stdout;

166 /* generate conveniently accessible forward & reverse strands */
forward = (char **) emalloc (seq->num_seq * sizeof (char *));
reverse = (char **) emalloc (seq->num_seq * sizeof (char *));
num_res = (long *) emalloc (seq->num_seq * sizeof (long));
forward[0] = seq->seq;
num_res[0] = seq->borders[0];
max_res = num_res[0];
171 for(i = 1; i < seq->num_seq; i++){
    forward[i] = seq->seq + seq->borders[i - 1] + 1;
    num_res[i] = seq->borders[i] - seq->borders[i - 1] - 1;
    if (num_res[i] > max_res)
        max_res = num_res[i];
    forward[i][num_res[i]] = '\0';
}
dic = get_restricted_dna_dictionary (dic);
181 if (args->p)
    profiles = initialize_profiles (profiles, max_res);
else
    profiles = NULL;
for (i = 0; i < seq->num_seq; i++)
    reverse[i] = revcomp_string (forward[i]);
coverage = args->c * seq->num_nuc;
if(args->a) /* dealing with an alignment? */
    coverage /= seq->num_seq;
sum = 0;
191 sum = 0;
seq_ind = 0;
while (sum < coverage) {
    /* choose input sequence */
    seq_ind = genrand_reall () * seq->num_seq;
    sequence[0] = forward[seq_ind];
    sequence[1] = reverse[seq_ind];
    residues = num_res[seq_ind];
    (*read_id)++;
    if(args->S) /* use fixed position */
        p = args->S - 1;
    else /* generate random position */
        p = (int) (genrand_reall () * residues);
    pf = p;
    /* get random read length */
    read_len = args->R + (rand_Normal()*args->D) + 0.5;
    /* get random fragment length */
    frag_len = args->l + (rand_Normal()*args->d) + 0.5;
    if(read_len > frag_len)
        read_len = frag_len;
    /* determine strand */
    if (!args->f){
        if (genrand_reall () > 0.5){
            seqi = 1;
            p = residues - p - 1;
        } else

```

```

        seqi = 0;
    } else
        seqi = 0;
221    /* write sequence to output */
    i = 0;
    if (args->r || read_len <= residues - p + 1) /* circular or linear
       without edge*/
        printed_len = read_len;
    else
        printed_len = residues - p;                      /* linear with edge */
226    if (seqi)
        strand = 'r';
    else
        strand = 'f';
231    for (j = 0; j < strlen(seq->headers[i]); j++){
        if (seq->headers[i][j] == '_'){
            seq->headers[i][j] = '\0';
            break;
        }
    }
236    if (!args->p && printed_len){
        if (seqi){
            end = pf + 1;
            start = pf - printed_len + 2;
        } else{
            start = pf + 1;
            end = pf + printed_len;
        }
        fprintf (fout, ">Read#%d_file=%s_sequence=%s_strand=%c_start=%ld_
end=%ld",
241             (int)*read_id, input_file, seq->headers[seq_ind] + 1,
                           strand, start, end);
        if (args->P)
            fprintf(fout, "mate=1\n");
        else
            fprintf(fout, "\n");
    }
251    sum += printed_len;
    process_read(fout, sequence[seqi], residues, p, read_len, args->E,
                 dic, dic_str, args->r, profiles);
    if (args->P){
        /* get new random read length */
256        read_len = args->R + (rand_Normal())*args->D) + 0.5;
        p += frag_len - 1;
        if (args->r)           /* circular genome */
            p %= residues;
        else{                  /* linear genome */
            if (p >= residues){
                if (residues - read_len < p - frag_len) /* does the read fit on
                   the remaining fragment? */
                    read_len = frag_len + residues - p - 1;
                p = residues - 1;
            }
        }
266    }

```

```

if (!args->f){ /* change strand */
    seqi = (seqi == 0) ? 1 : 0;
    if(seqi)
        pf = p;
    p = residues - p - 1;
    if (!seqi)
        pf = p;
}
if (seqi)
    strand = 'r';
else
    strand = 'f';
if (args->r || read_len <= residues - p + 1)
    printed_len = read_len;
else
    printed_len = residues - p;
if (!args->p && printed_len){
    if(seqi){
        end = pf + 1;
        start = pf - printed_len + 2;
    }else{
        start = pf + 1;
        end = pf + printed_len;
    }
    fprintf (fout, ">Read#%d_file=%s_sequence=%s_strand=%c_start=%ld_
end=%ld_mate=2\n",
            (int)*read_id, input_file, seq->headers[seq_ind] + 1,
            strand, start, end);
    process_read(fout, sequence[seqi], residues, p, read_len, args->E
                , dic, dic_str, args->r, profiles);
}
sum += printed_len;
}
if (args->p){
    if(args->C)
        print_profiles_single_contig (fout, profiles, max_res, args->m);
    else
        print_profiles (fout, profiles, max_res, args->m);
}
}

```

## Change History

- v. 1.5 (August 2008)
  - allow multiple input sequences
  - set for read length distribution  $\sigma = \sqrt{\mu}$  (previously:  $\sigma = \mu \times 0.1$ )
- v. 1.6 (April 18, 2009)
  - added more information to the header lines of the reads
  - added the option to print out profiles instead of raw sequences
  - added -f option to sample from forward strand only
  - added -p option to print out profiles

- added `-m` option to specify minimum coverage of valid profile
  - added sequencing error and corresponding `-E` option
  - added random seed handling via file `randomSeed.dat`
- v. 1.7 (May 19, 2009)
  - `-p` option now implies `-f`, which is not a user option any more
- v. 1.15 (October 27, 2017)
  - Reverted from configure/make system to my home-grown style in order to get in line with the other programs published as part of “Bioinformatics for Evolutionary Biologists, A Problems Approach”.