

IThOS

User Guide : version 2.1

Nouri Ben Zakour
Yves Le Loir
Dominique Lavenier

1. Introduction

IThOS is a software package dedicated to the design of primers. The input is one or several genome sequences (FASTA files) and the output is a list of primer candidates that fulfill criteria set by the user. IThOS also determines putative hybridization sites for these primers.

IThOS works in a two-step procedure:

1. Primer design
2. Verification of hybridization sites

Each step relies on a dedicated program:

```
1: ithos_gen <genome> <param> <primers> [-c]
2: ithos_chk <genome> <primers> <param> <pr_out> <pr_hyb> <pr_dic> <pr_pos>
```

Both programs can be used separately. A third program completes the package and enables the visualization of the primer features:

- `ithos_viz <primers> <param>`

Parameters of the different programs are:

`<genome>`
File containing one or several DNA sequences, FASTA format.

`<param>`
File specifying criteria for primer selection.

`<primers>`, `<pr_in>`, `<pr_out>`, `<pr_hyb>`
Files containing a primer list, FASTA format.

`[-c]`
Option enabling the search of primers on the complementary strand.

`<pr_dic>`
Dictionary of primers

`<pr_pos>`
Position of hybridations

2. Design of the primers: ithos_gen

Starting from one or several DNA sequences, the goal is to design primers that fulfill criteria set by the user. The ithos_gen program considers all the words whose size is in an interval that corresponds to the minimum and maximum primer length. For each word, a suite of filters is applied. All oligonucleotides passing successfully through the filters are proposed as primer candidates. For each filter, several parameters can be set by the user to refine the primer selection according to the application. Six filters are implemented. They are described in the following sections.

Filter 1: G+C %

For a primer of size T, filter 1 works as follows:

- Counts the number of G and C nucleotides (# GC)
- Calculates the percentage $\rightarrow P = (\# \text{ GC} * 100) / T$
- Discard the primer if: $P < \text{pcGCMin}$ or $P > \text{pcGCMax}$

The default values are:

- pcGCMin = 40
- pcGCMax = 60

Filter 2: Tm (Melting temperature)

The nearest neighbor method is used to calculate the primer (Santa Lucia et al., 1998). It also takes into account the concentration of nucleotides (dnaConc) and the concentration of salt (saltConc).

This filter removes primers if:

- $T_m < \text{oligoTmMin}$
- $T_m > \text{oligoTmMax}$

The default values are:

- oligoTmMin = 57 °C
- oligoTmMax = 62 °C
- dnaConc = 500 nM
- saltConc = 50 nM

Bibliography

A Unified View of Polymer, Dumbbell, and Oligonucleotide DNA Nearest-Neighbor Thermodynamics John SantaLucia, *Proceedings of the National Academy of Sciences of the United States of America*, Vol. 95, No. 4 (Feb. 17, 1998), pp. 1460-1465

Filter 3: Number of repeats

This filter removes oligonucleotides with N consecutive identical nucleotides or dinucleotides (nbRepeat). For example, if nbRepeat = 4, the following primers will be removed:

```
1: GGGATGGACACGGATTTTTGGACCAGC
2: TTAGCTATATATAGGCAGGGATTAGG
```

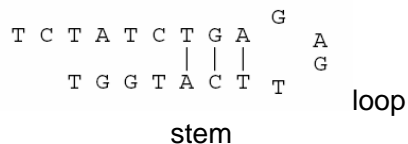
The first one presents a suit of 4 « T ». The second one a suite of 4 « TA »
The default value is:

- nbRepeat = 5

Filter 4: Hairpin

This filter removes oligonucleotides with hairpin loops that present the following features:

- stem size > or = to maxHpDup
- loop size > or = to MaxHpLoop



The default values are:

- maxHpDup = 4
- maxHpLoop = 4

Filter 5 : self-complementarity

This filter checks that a primer will not hybridize with itself during PCR. Thus, primers that form a duplex with their complementary strand are removed. Criteria for selection are as follow:

- For the full length of the pirmer, the authorized deltaG value must not exceed maxDeltaGAuto
- For the 3' end and a distance of sizeDeltaGAuto3, the authorized deltaG value must not exceed maxDeltaGAuto3.
- The minimal size of a tested duplex is sizeDeltaGAuto

The default values are:

- maxDeltaGAuto = -10kcal/mol
- maxDeltaGAuto3p = -7kcal/mol
- sizeDeltaGAuto = 6
- sizeDeltaGAuto3p = 8

Filtre 6 : thermodynamic stability at primer ends

This filter calculates a deltaG on the 5' and 3' ends of the primers. The size of the 5' end to be considered is given by sizeExt5. The size of the 3' end to be considered is given by sizeExt3.

A primer is removed if:

- The value in 5' is above deltaG5
- The value in 3' is out of the interval [deltaG3Min,deltaG3Max]

The default values are:

- sizeExt5 = 5
- sizeExt3 = 5
- deltaG5 = - 4 kcal/mol
- deltaG3Min = -6 kcal/mol
- deltaG3Max = - 4 kcal/mol

Software execution

The software is launched by the command line:

```
ithos_gen <genome> <parameters> <primers> [-c]
```

input files

- <genome> is a text file (FASTA) that contains a genome that can be cut into several sequences.
- <parameters> is a text file that contains the parameters for the different filters

output files

- <primers> is a text file (FASPTA format) that contains all the selected primers. Coordinates of the selected primers are given as comments in between brackets ([]).
- The software generates an additional file <primers.log> that gives a few statistical data on the filtering process.

Option

- [-c] allows researching on the complementary strand of the genome

Example 1 : search for primers on the leading strand

If a file containing the 2 following sequences is considered:

```
>sequence #1
AGACAGATATCACAGATATGAGACCACGAGGATAGAGCCCAGAGTAGACAG
>sequence #2
AATACACGAAGCGAGCAATCAACTTGACCTAGGTGAGGGATAGGACCAGA
```

The software `ithos_gen` will produce the following file:

```
>sequence #1 [25 50]
ACGAGGATAGAGCCCAGAGTAGACAG
>sequence #1 [10 35]
CACAGATATGAGACCACGAGGATAGA
>sequence #1 [12 37]
CAGATATGAGACCACGAGGATAGAGC
>sequence #2 [20 45]
AACTTGACCTAGGTGAGGGATAGGAC
>sequence #2 [3 28]
ACACGAAGCGAGCAATCAACTTGACC
>sequence #2 [21 46]
ACTTGACCTAGGTGAGGGATAGGACC
>sequence #2 [9 34]
AGCGAGCAATCAACTTGACCTAGGTG
>sequence #2 [19 44]
CAACTTGACCTAGGTGAGGGATAGGA
>sequence #2 [6 31]
CGAAGCGAGCAATCAACTTGACCTAG
>sequence #2 [7 32]
GAAGCGAGCAATCAACTTGACCTAGG
>sequence #2 [18 43]
TCAACTTGACCTAGGTGAGGGATAGG
```

The numbers in brackets indicate the coordinates – begin and end – of the primer in the sequence. Note that for each sequence, primers are alphabetically sorted.

Example 2 : search for primers on the complementary strand (option `-c`)

The search for primers is done on the complementary strand of the genome sequence. The software generates exactly the same type of file. The only difference is that the primer coordinates are inverted. Thus, if we consider the previous example, the program will generate:

```
>sequence #1 [36 11]
CTCTATCCTCGTGGTCTCATATCTGT
>sequence #1 [37 12]
GCTCTATCCTCGTGGTCTCATATCTG
>sequence #1 [38 13]
GGCTCTATCCTCGTGGTCTCATATCT
>sequence #1 [39 14]
GGGCTCTATCCTCGTGGTCTCATATC
>sequence #1 [40 15]
TGGGCTCTATCCTCGTGGTCTCATAT
>sequence #2 [43 18]
CCTATCCCTCACCTAGGTCAAGTTGA
```

```

>sequence #2 [42 17]
CTATCCCTCACCTAGGTCAAGTTGAT
>sequence #2 [45 20]
GTCCTATCCCTCACCTAGGTCAAGTT
>sequence #2 [39 14]
TCCCTCACCTAGGTCAAGTTGATTGC
>sequence #2 [44 19]
TCCTATCCCTCACCTAGGTCAAGTTG

```

3. Checking the hybridization sites: ithos_chk

For each primer, the program checks that there are no secondary hybridization sites elsewhere on the whole genome, e.g. for a primer: 5' - T G A - 3', the following hybridization sites must be checked:



Detection of a secondary hybridization site is not based on the percentage of identity but on the calculation of the thermodynamic stability of the duplex (cf. filter auto-complementarity).

A maximal deltaG is calculated on the whole primer length and a maximal deltaG in 3' is calculated on sizeDeltaGHybrid3 nucleotides. A deltaG value is calculated on consecutive matches, including putative mismatches (1 mismatch surrounded by 2 matches). A hybridization site is recognized if one of the two conditions is true:

- deltaG < maxDeltaGHybrid
- deltaG in 3' < maxDeltaGHybrid3

For example:



deltaG = max (G1,G2,G3) and deltaG in 3' = G3

The values for G1, G2 and G3 are the sum of the thermodynamic values between 2 consecutive nucleotide pairs.

The default values are:

- maxDeltaGHybrid = -16 kcal/mol
- maxDeltaGHybrid3 = -9 kcal/mol
- sizeDeltaGHybrid3 = 8

Software execution

The software is launched by the following command line:

```
ithos_chk <genome> <primers> <param> <pr_out> <pr_hyb> <pr_dic> <pr_pos>
```

Input files

- <genome> is a text file (FASTA format) containing a genome sequence
- <primers> is a text file (FASTA format) containing a list of primers
- <param> is a text file containing the parameters for the different filters

Output file

- <pr_out> a text file containing all the primers that have no hybridization sites
- <pr_hyb> a text file containing all the primers that have at least one hybridization site
- The program generates an additional file <pr_hyb.info> that indicates, for each primer of the file <pr_hyb> the positions of hybridization on the genome as well as the deltaG values.
- <pr_dic> a text file containing the primer dictionary (alphabetically sorted)
- <pr_pos> a text file containing all the sorted positions of the hybridization of the primer dictionary.

Example 3: Primer design on a genome and elimination of the primers that hybridize at other positions

If the following genome is considered and memorized in a file named `ex3`

```
>exemple_3
AAGATAGAAATACACGATGCGAGCAATCAAATTTTCAGGTAGAAAGGATAGA
AATACACGAAGCGAGCAATCAACTTGACCTAGGTGAGGGATAGGACCAGA
```

Primer design is launched by the command line:

```
ithos_gen ex3 parameters primer
```

This gives a file named `primer`. This file contains, for example, the following primer list:

```
>exemple_3 [71 96]
AACTTGACCTAGGTGAGGGATAGGAC
>exemple_3 [54 79]
ACACGAAGCGAGCAATCAACTTGACC
>exemple_3 [72 97]
ACTTGACCTAGGTGAGGGATAGGACC
>exemple_3 [60 85]
AGCGAGCAATCAACTTGACCTAGGTG
>exemple_3 [70 95]
CAACTTGACCTAGGTGAGGGATAGGA
>exemple_3 [57 82]
CGAAGCGAGCAATCAACTTGACCTAG
>exemple_3 [58 83]
GAAGCGAGCAATCAACTTGACCTAGG
>exemple_3 [69 94]
TCAACTTGACCTAGGTGAGGGATAGG
```

Checking of the hybridization sites is carried out by the command line:

```
ithos_chk ex3 primer parameters primer_ok primer_hyb primer_dic primer_pos
```

This produces 4 output files:

primer_ok

This file is empty since at least one hybridization site has been found.

primer_hyb

```
>exemple_3 [5 29]
AGAAATACACGATGCGAGCAATCAA
>exemple_3 [11 35]
ACACGATGCGAGCAATCAAATTTCA
>exemple_3 [13 37]
ACGATGCGAGCAATCAAATTTTCAGG
>exemple_3 [44 68]
GGATAGAAATACACGAAGCGAGCAA
>exemple_3 [49 73]
GAAATACACGAAGCGAGCAATCAAC
>exemple_3 [70 94]
ACTTGACCTAGGTGAGGGATAGGAC
```

primer_hyb.info

```
>exemple_3 [71 96]
AACTTGACCTAGGTGAGGGATAGGAC
seq start end dG max dG 3'
5' AACTTGACCTAGGTGAGGGATAGGAC 3' primer
  |||
3' TTGAACTGGATCCACTCCCTATCCTG 5' genome 71 96 -32614 -9574

>exemple_3 [54 79]
ACACGAAGCGAGCAATCAACTTGACC
seq start end dG max dG 3'
5' ACACGAAGCGAGCAATCAACTTGACC 3' primer
  |||
3' TGTGCTTCGCTCGTTAGTTGAACTGG 5' genome 54 79 -35754 -9864

5' ACACGAAGCGAGCAATCAACTTGACC 3' primer
  |||
3' TGTGCTACGCTCGTTAGTTTAAAGTC 5' genome 11 36 -21964 236

>exemple_3 [72 97]
ACTTGACCTAGGTGAGGGATAGGACC
seq start end dG max dG 3'
5' ACTTGACCTAGGTGAGGGATAGGACC 3' primer
  |||
3' TGAACCTGGATCCACTCCCTATCCTGG 5' genome 72 97 -33554 -9574

>exemple_3 [60 85]
AGCGAGCAATCAACTTGACCTAGGTG
seq start end dG max dG 3'
5' AGCGAGCAATCAACTTGACCTAGGTG 3' primer
  |||
3' TCGCTCGTTAGTTGAACTGGATCCAC 5' genome 60 85 -34594 -10434

>exemple_3 [70 95]
CAACTTGACCTAGGTGAGGGATAGGA
seq start end dG max dG 3'
5' CAACTTGACCTAGGTGAGGGATAGGA 3' primer
  |||
3' GTTGAACCTGGATCCACTCCCTATCCT 5' genome 70 95 -33144 -9914

>exemple_3 [57 82]
CGAAGCGAGCAATCAACTTGACCTAG
seq start end dG max dG 3'
5' CGAAGCGAGCAATCAACTTGACCTAG 3' primer
  |||
3' GCTTCGCTCGTTAGTTGAACTGGATC 5' genome 57 82 -34844 -9454

>exemple_3 [58 83]
GAAGCGAGCAATCAACTTGACCTAGG
seq start end dG max dG 3'
5' GAAGCGAGCAATCAACTTGACCTAGG 3' primer
  |||
3' CTTGCTCGTTAGTTGAACTGGATCC 5' genome 58 83 -34624 -10294

>exemple_3 [69 94]
TCAACTTGACCTAGGTGAGGGATAGG
seq start end dG max dG 3'
5' TCAACTTGACCTAGGTGAGGGATAGG 3' primer
  |||
3' AGTTGAACTGGATCCACTCCCTATCC 5' genome 69 94 -33504 -10124
```

primer_dic

0 AACTTGACCTAGGTGAGGGATAGGAC
1 ACACGAAGCGAGCAATCAACTTGACC
2 ACTTGACCTAGGTGAGGGATAGGACC
3 AGCGAGCAATCAACTTGACCTAGGTG
4 CAACTTGACCTAGGTGAGGGATAGGA
5 CGAAGCGAGCAATCAACTTGACCTAG
6 GAAGCGAGCAATCAACTTGACCTAGG
7 TCAACTTGACCTAGGTGAGGGATAGG

primer_pos

11 1
54 1
57 5
58 6
60 3
69 7
70 4
71 0
72 2

4. Visualization of the primers features: ithos_viz

This utility displays the primer features. For each primer, it gives:

- the GC percent
- the melting temperature: Tm
- the maximal suite of identical nucleotides
- the size of the biggest stem-loop structure
- maximal deltaG for the complementary primer
- maximal deltaG in 3' for the complementary primer
- Stability in 5'
- Stability in 3'

Software Execution

The program is launched by the command line:

```
ithos_viz <primers> <parameters>
```

Results are displayed on monitor screen

Input files

- <primers> is a text file (FASTA format) that contains a primers list
- <parameters> is a text file that contains parameters of the different filters

Example 4: visualization of primer_ok file

The execution of the following command line:

```
ithos_viz primer_hyb parameters
```

Displays on the monitor screen:

>exemple_3 [71 96]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
AACTTGACCTAGGTGAGGGATAGGAC	50	58	3	4	1960	1960	-6170	-6440
>exemple_3 [54 79]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
ACACGAAGCGAGCAATCAACTTGACC	50	61	2	3	-210	1960	-7800	-7030
>exemple_3 [72 97]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
ACTTGACCTAGGTGAGGGATAGGACC	53	59	3	4	1960	1960	-6470	-7700
>exemple_3 [60 85]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
AGCGAGCAATCAACTTGACCTAGGTG	50	60	2	2	-280	1960	-8270	-6590
>exemple_3 [70 95]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
CAACTTGACCTAGGTGAGGGATAGGA	50	58	3	4	1960	1960	-6170	-5880
>exemple_3 [57 82]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
CGAAGCGAGCAATCAACTTGACCTAG	50	59	2	3	1960	1960	-7990	-6420
>exemple_3 [58 83]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
GAAGCGAGCAATCAACTTGACCTAGG	50	59	2	3	-280	1960	-7990	-6820
>exemple_3 [69 94]	%GC	Tm	#rep	Hpin	Cp_dG	Cp_dG3'	Sta_5'	Sta_3'
TCAACTTGACCTAGGTGAGGGATAGG	50	58	3	3	1960	1960	-6470	-5880

5. Parameters file

This is a text file that enables modifying the default parameters of the filters. The same file is read by the three programs.

Example

```
# IThOS PARAMETERS

# primers length

lengthMin 25
lengthMax 25

# filter 1: GC percentage

pcGCMin 40
pcGCMax 60

# filter 2: tm

oligoTmMin 57
oligoTmMax 67
dnaConc 500
saltConc 50

# filter 3: hairpin

maxHpDup 4
maxHpLoop 4

# filter 4: repeat

nbRepeat 6

# filter 5: auto complementarity

maxDeltaGAuto -10000
maxDeltaGAuto3 -6000
sizeDeltaGAuto 8
sizeDeltaGAuto3 8

# filter 6: internal stability to 3' & 5' extremities

sizeExt5 5
sizeExt3 5
deltaG5 -4000
deltaG3min -6000
deltaG3max -3000

#site hybridation (only used by ithos_chk)

sizeDeltaGHybrid3 8
maxDeltaGHybrid3 -12000
maxDeltaGHybrid -18000
```