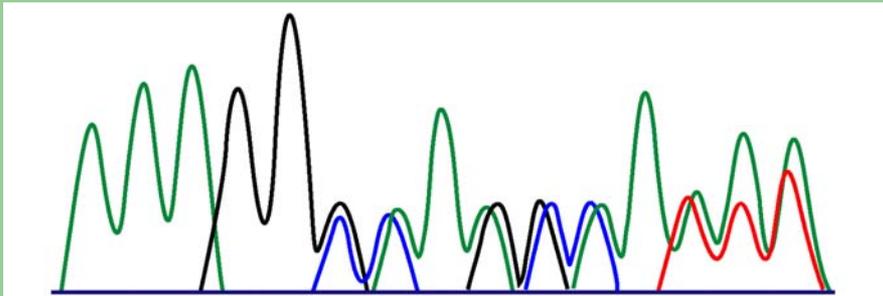


Decoding of Superimposed Traces Produced by Direct Sequencing of Heterozygous Indels



Dmitriev, D.A. & Rakitov, R.A.

Illinois Natural History Survey,
Institute of Natural Resource Sustainability,
University of Illinois at Urbana-Champaign,
1816 S. Oak st.,
Champaign IL, 61820
<http://ctap.inhs.uiuc.edu/dmitriev/indel.asp>



Problem of double peaks on chromatogram

A

```

AAAGGG-CAAGCAATAT
| | | | | | | | | |
AAAGGGGCAAGCAATAT
  
```

B

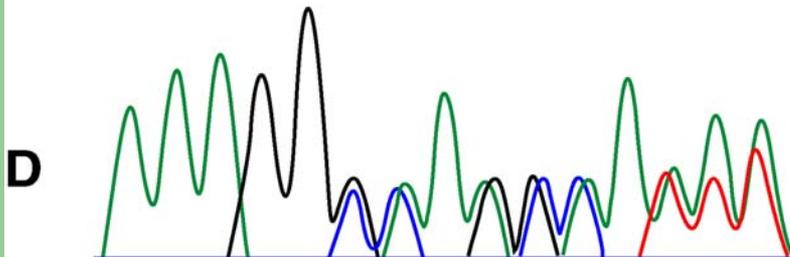
```

AAAGGGCAAGCAATAT.
| | | | | \ \ \ \ \
AAAGGGGCAAGCAATAT
  
```

C

```

AAAGGGSMA RSM AWW
  
```



IUPAC
symbols:

R = A or G
 Y = C or T
 K = G or T
 M = A or C
 S = C or G
 W = A or T
 B = not A
 D = not C
 H = not G
 V = not T
 N = any

A pair of allelic sequences properly aligned (**A**), unaligned (**B**), translated into a consensus (**C**), and resulted chromatogram (**D**).

The following applications could be used to call double peaks PHRED, KB Basecaller, Sequencher

Reasons for double peaks

- Direct sequencing of diploid alleles containing heterozygous insertions/deletions. (Mixed trace downstream of the indel is formed by two allelic traces superimposed onto each other with a phase shift).
- Sequencing of unrelated templates and alternative splicing.
- Single nucleotide polymorphisms (SNPs) or base calling errors due to low quality of chromatogram (individual double picks)

Solutions

- Discard as uninterpretable.
- Use new sequencing technologies, such as pyrosequencing, which works with single DNA molecules.
- Separating the templates prior to sequencing via cloning into a vector or selectively amplifying one allele using allele-specific primers.
- Computational methods to extract information from mixed traces.

Computational methods for decoding of mixed traces

- **Subtracting a reference sequence:** PolyPhred, STADEN, CodonCode Aligner, Mutation Surveyor, InSNP, PolyScan, AutoCSA, and Tenney, A.E et al. (2007) application to automatically resolving double traces aligning the mixed sequence to genomic database
- **Using reverse sequence as a reference:** SeqScape, Varian Reporter, Champuru.
- **Extracting information from an individual mixed trace (without reference):** Shift Detector, CodonCode Aligner Ver. 2, Indelligent.

Dynamic optimization algorithm

A F: M W W M M W Y S K S T K W Y

B

$Z_i=1$	(A/C)	(A/T)	(A/T)	(A/C)	(A/C)	(A/T)	(C/T)	(G/C)	(G/T)	(G/C)	(T/T)	(G/T)	(A/T)	(C/T)
$Z_i=2$	(C/A)	(T/A)	(T/A)	(C/A)	(C/A)	(T/A)	(T/C)	(C/G)	(T/G)	(C/G)	(C/G)	(T/G)	(T/A)	(T/C)
i	1	2	3	4	5	6	7	8	9	10	11	12	13	14

$$p(6, 2, 1) = \max \begin{cases} p'(6, 2, 1, 0) = \max \{p(5, 1, 0), p(5, 2, 0)\} + 1W_m - 1W_m - 1W_{ib} = 2 + 1 - 1 - 1 = 1 \\ p'(6, 2, 1, 1) = \max \{p(5, 1, 1), p(5, 2, 1)\} + 1W_m = 5 + 1 = 6 \\ p'(6, 2, 1, 2) = \max \{p(5, 1, 2), p(5, 2, 2)\} + 1W_m - 1W_m - 1W_{ib} = 4 + 1 - 1 - 1 = 3 \end{cases}$$

C

$k_i=0$	1	0	0	0	1	2	3	4	5	6	7	8	8	8	8
2	0	0	0	0	1	2	3	4	5	6	7	8	8	8	8
$k_i=1$	1	1	3	3	5	5	7	8	8	9	9	10	10	11	
2	1	2	2	4	4	6	6	7	9	9		9	10	10	
$k_i=2$	1	2	2	3	4	4	5	6	7	8	9	9	11	12	
2	1	2	3	4	4	4	5	6	7	8		10	10	11	

order of computation →

D

9	8	7	6	5	4	4	3	2	1	1	0	0	0	0
9	8	7	6	5	4	4	3	2	1		0	0	0	
12	10	10	8	8	6	6	5	4	3	3	2	1	1	
11	11	9	9	7	7	5	4	3	3		2	2	1	
10	10	10	9	7	7	7	6	5	5	4	2	2	1	
10	10	9	8	8	7	7	7	6	4		3	2	1	

← order of computation

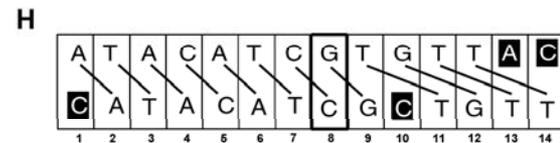
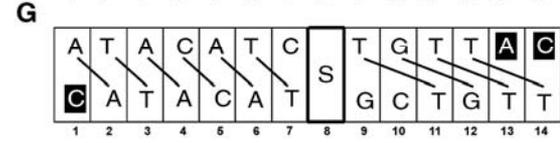
E

9	8	7	7	7	7	8	8	8	8	9	8	8	8	
9	8	7	7	7	7	8	8	8	8		8	8	8	
13	11	13	11	13	11	13	13	12	12	12	12	11	12	
12	13	11	13	11	13	11	11	12	12		11	12	11	
11	12	12	12	11	11	12	12	12	13	13	11	13	13	
11	12	12	12	12	11	12	13	13	12		13	12	12	

$\omega(i, z_i, k_i)$

F

1	(A/C)	(A/T)	(A/C)	(C/T)	(G/C)	(G/C)	(T/T)	(A/T)	(C/T)					
2	(T/A)	(C/A)	(T/A)	(C/G)	(T/G)	(T/G)								
$k_i=1$	1	1	1	1	1	1	?	2	2	2	2	2	2	2
	1	2	3	4	5	6	7	8	9	10	11	12	13	14



I

```

- A T A C A T C G - T G T T A C
  | | | | | | | | | |
  C A T A C A T C G C T G T T - -
    
```

Multiple cooptimal solutions

RWKRWKRWKRWK

ATGATGATGATG
GATGATGATGAT

GATGATGATGAT
ATGATGATGATG

ATGGATATGGAT
GATATGGATATG

A

MKMWKMKKCKKM

ATATGCGTCTGA
CGCATATGCGTC

AGATTCTGGCTTA
CTCAGATTCTGGC

AKATKCGKCTKA
CKCAKATKCGKC

B

MKMWKMKKCKKM

GAGT T TGCCGTC
TTGAG T GTGCCG

GAG T GTGCCGTC
TTGA T TGTGCCG

GAGTKTGCCGTC
TTGAKTGTGCCG

C

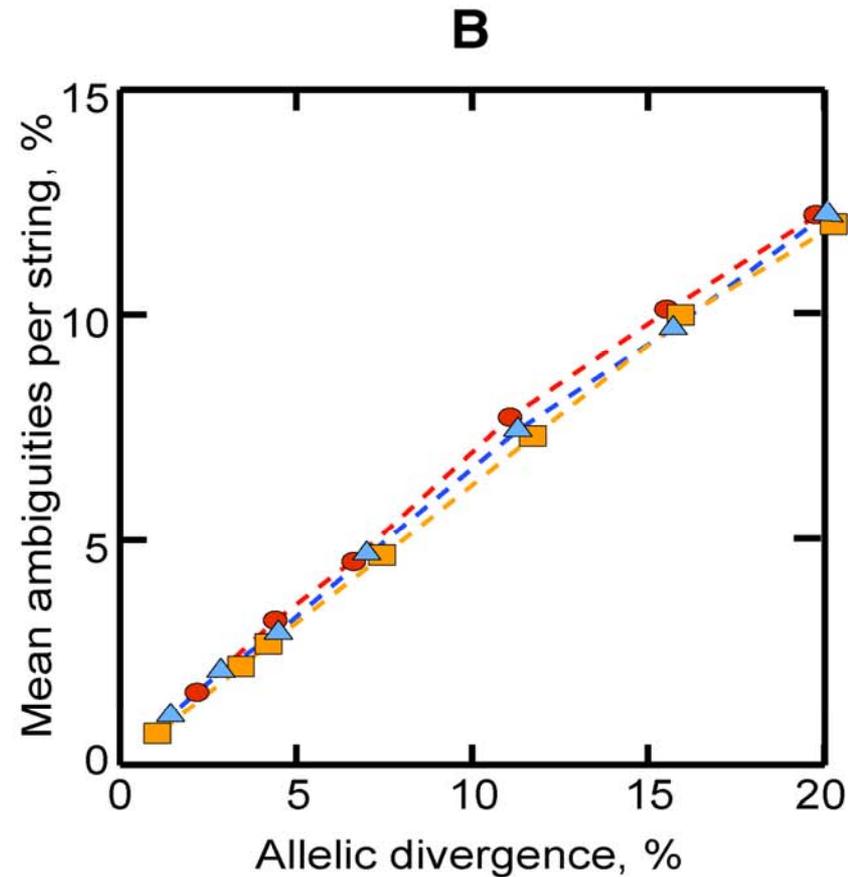
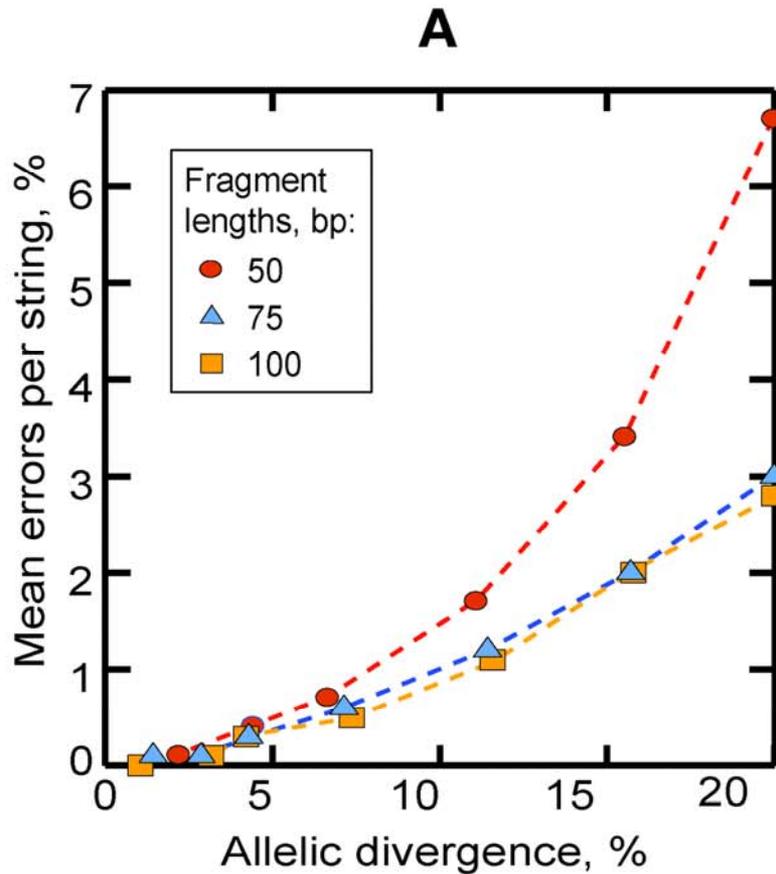
RGKTYYWYCAYS

GGTTCATCACG
AGGTTTCCATC

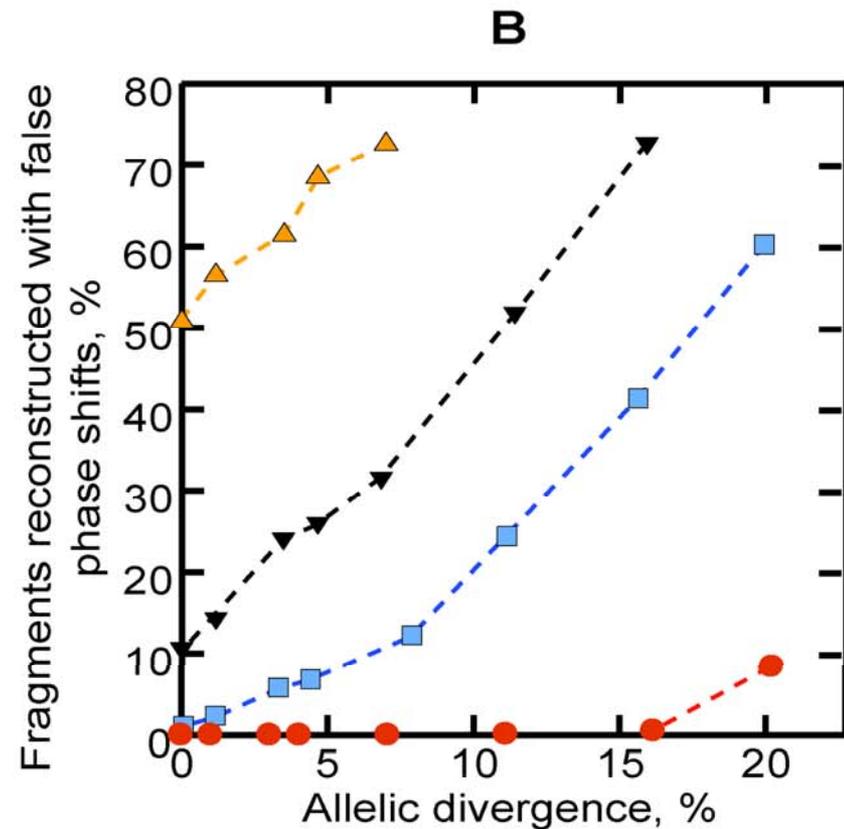
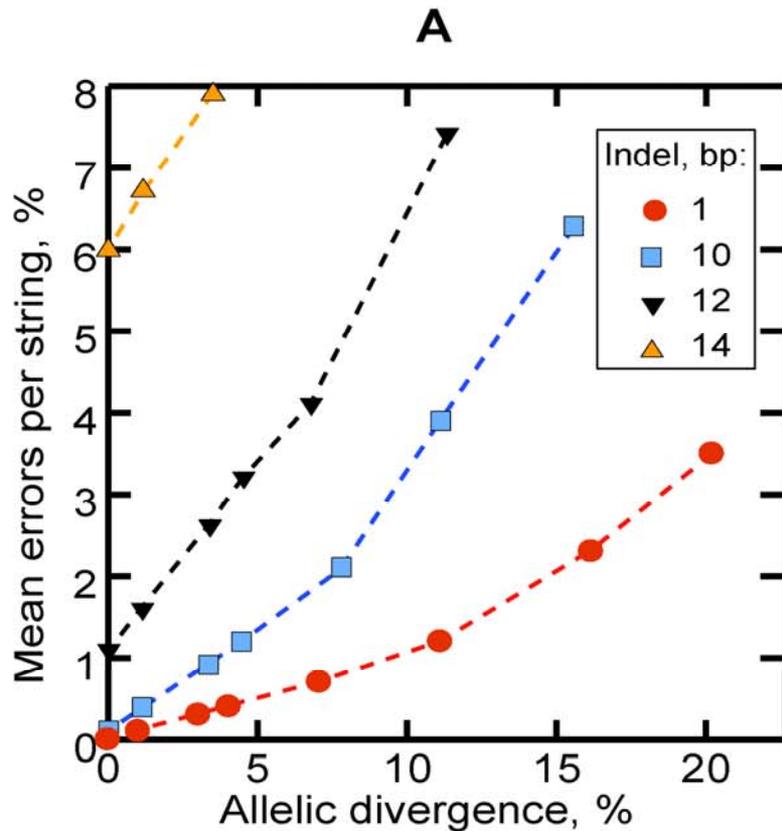
GGTTCATCACG
AGGTTCTCCATC

D

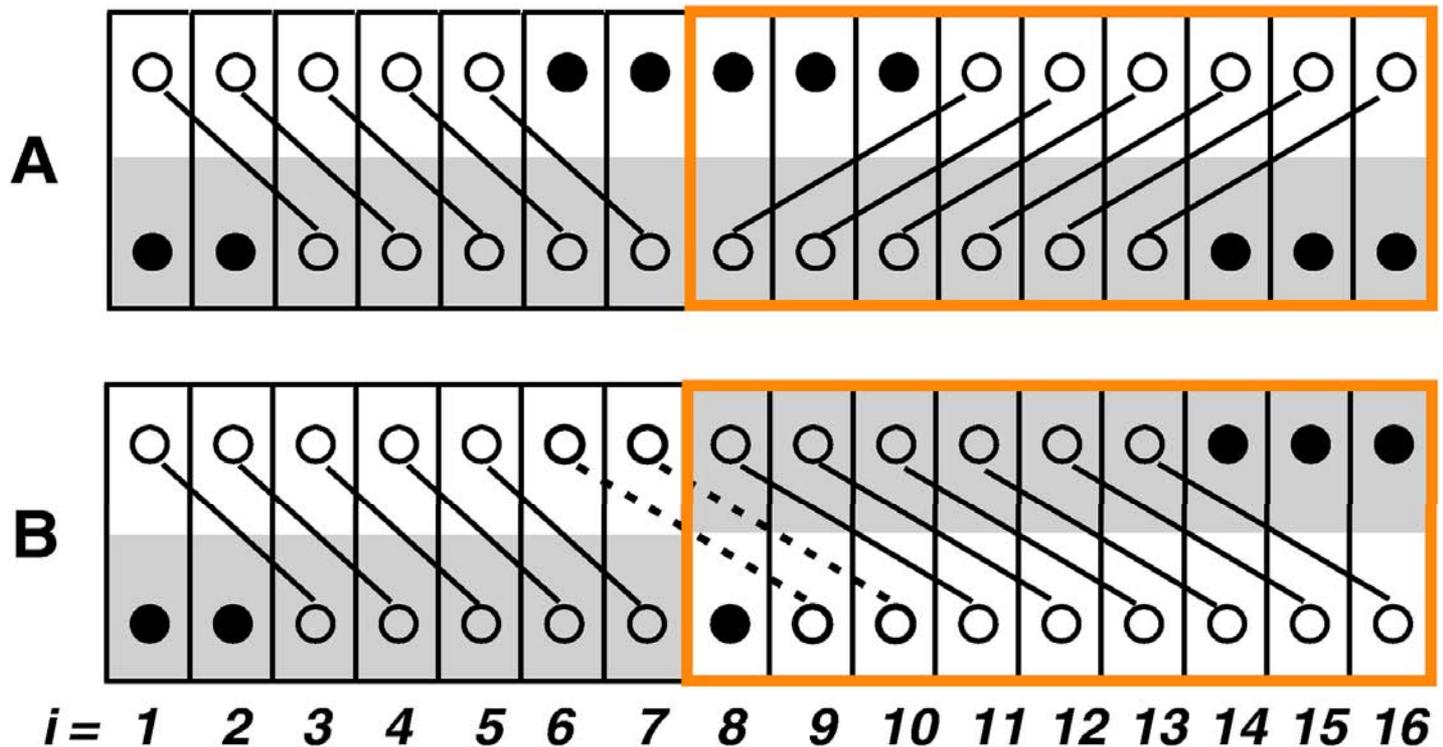
Accuracy of decoding of simulated mixed fragments formed by a single 5bp phase shift (1000 runs for each point)



Accuracy of decoding of simulated mixed 100 bp fragments formed with an insertion of variable size in the middle (1000 runs for each point)



Two aligned solutions of the same mixed fragment representing the transition between two phase sifts (“Long” vs. “Short” indel)



Validation with human traces

- We used 104 (103-677bp, with indels 5-30 bp) traces of 198 recorded by Bhangale et al. (2005) from NCBI Trace Archive as having heterozygous indel.
- Sequencer was used to call second peaks.
- After reconstruction, traces were aligned with best matching human sequences in NCBI Trace Archive.
- 102 traces reconstructed with a single indel, two with two indels.
- 67 traces without errors, 31 – with 1-2 errors, 6 – with 3-7 errors. Half of the fragments reconstructed without ambiguities.
- Mean of $99.1 \pm 1.25\%$ of bases per fragment decoded correctly and unambiguously (in the same conditions ShiftDetector decoded only $72.5 \pm 6.47\%$ of bases).
- About 60% of reconstructed mean of 0.66 errors per fragment were due base calling errors, mostly in low-quality trace regions.

Conditions for reconstructions with 99.5-100% accuracy

- Homologous fragments resulted from indel mutation.
- Analyzed fragment is significantly larger than the indel (at least 10 times; in human 92.3% of indels are 1-10bp).
- Low divergence between mixed traces (<5%; for human noncoding DNA, the average divergence is <0.1%, fruit fly 1-2%, sea squirts 4.5%).
- Multiple indels, if present, are well spaced.
- Methods relies on base calling software.

Indelligent interface

Indelligent v.1.2 [Home page](#) | [Help](#)

© 2008 Dmitry Dmitriev & Roman Rakitov

IUPAC codes

- R = A or G
- Y = C or T
- K = G or T
- M = A or C
- S = C or G
- W = A or T
- B = not A
- D = not C
- H = not G
- V = not T
- N = any

The program reconstructs allelic sequences based on the pattern of calls within an individual convoluted DNA trace produced by direct sequencing of a diploid template containing heterozygous insertions/deletions.

Enter sequence to be analyzed in the window (click on the example: [TKGKKSCMW](#)) and press "Submit". [Input options](#).

Not for use in diagnostic procedures.

```
TTCTCCGAGCTCCAAACAGTTGGAATTGGAAATACTCGAAAGCAAGCCCSWGKYCRKWSCRTAGYAYMGTAMMRTMA  
AAWCMWRRRCMYRRCASMRRCRGYS
```

Parameters

- Max phase shift size, bp.: ?
- Shift change penalty: ?
- Fix shift(s), bp.: ?

Output view options

- Align alleles: ?
- Align floating indels: Left Right ?
- Display "long" indels: ?

[Simulate indels](#) ?

Submit Clear

Number of users: 205.
Fragments analyzed: 7615 (2915960 bp).

Reconstruction results

IUPAC sequence

GTTCCGGATTACCTTTTTAGTGGCTGGTGACGACAACATCGGCAGTGGGGRSRYGMYRYKYKRSYYGWCMRYTMRRSKMGGGKSKKRWGKWRRRRRRM

Reconstructed sequences

GTTCCGGATTACCTTTTTAGTGGCTGGTGACGACAACATCGGCAGTGGGG...ACATGCTGCTTGAGCCGTTCACTAGGGGCGGGTGTAAAGGAGGAGGAAC
GTTCCGGATTACCTTTTTAGTGGCTGGTGACGACAACATCGGCAGTGGGGGGCGACATGCTGCTTGAACCGTTCACTAGGGGCGGGTGTAAAGGAGGA.....

Combined sequence ?

GTTCCGGATTACCTTTTTAGTGGCTGGTGACGACAACATCGGCAGTGGGGGGCGACATGCTGCTTGARCCGTTCACTAGGGGCGGGTGTAAAGGAGGAGGAAC

Statistics

Length of original sequence: 100.

Number of resolved positions: 94 (94 %). ?

Number of mismatches: 1 (1 %). ?

Number of ambiguities: 1 (1 %). ?

Percent of resolved ambiguities: 97.56 %. (Maximum expected by chance: 80.07 %). ?

A floating indel has been reconstructed. ?

Detected phase shifts: ?

0 bp., 50 positions.

5 bp., 50 positions.

Script running time: 0.703125 sec.

Number of users: 205.

Fragments analyzed: 7618 (2916260 bp).

Acknowledgments

We are thankful to Saurabh Sinha for valuable discussion and suggestions

Tushar Bhangale for providing information which enable us to obtain human traces for testing.

Chris Dietrich for helpful comments and support.

The work was partially supported by NSF grants.