

# OASIS User's Manual

David Robinson  
dgrtwo@princeton.edu

OASIS, or Optimized Annotation System for Insertion Sequences, is a system for automatic annotation of Insertion Sequence (IS) elements in prokaryotic genomes, based on NCBI Genbank files.

OASIS is a free open source package developed on a Mac platform and implemented in Python, available upon request (contact information is above).

## Requirements

- Python 2.4-2.6 (not 3.0)
  - Python can be obtained from <http://www.python.org/>.
- BioPython 1.46 or above
  - BioPython can be obtained from <http://biopython.org/wiki/> and is used by OASIS for sequence input and analysis
- NCBI BLAST, specifically the `blastall` and `formatdb` executables, available from <http://www.ncbi.nlm.nih.gov/BLAST/download.shtml>

## Installation and setup

- If necessary, install the prerequisites of Python and BioPython, and download NCBI BLAST from the above sites
- In the file `Configuration.py`, set the `BLAST_EXE` and `FORMAT_EXE` to the full path to the `blastall` and `formatdb` executables on your computer, respectively.

## Input

OASIS takes as input a microbial genome in Genbank format (<http://www.ncbi.nlm.nih.gov/Genbank/>). The genome must already have annotated transposases, as do most NCBI publicly available prokaryotic genomes, for OASIS to find the insertions sequences. Sequenced microbial genomes in Genbank format are available at <ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria>.

OASIS can alternatively be given a folder of Genbank files, in which case it will output to a folder of annotations.

## Output

OASIS outputs two files for each genome- a fasta file and a text file. The text file is tab-delimited, with one line for each insertion sequence. Multiple copies of a single insertion sequences are divided into blocks by newlines.

Each IS element is described by the following columns:

Chromosome	Start	End	Length	Direction	ORF(s)	Label	Family	Group	LIR	RIR
------------	-------	-----	--------	-----------	--------	-------	--------	-------	-----	-----

- **Chromosome**- the locus name of the chromosome or plasmid on which this IS element appears
- **Start**- the chromosome nucleotide position of the first base pair of the IS element
- **End**- the chromosome nucleotide position of the last base pair of the IS element
- **Length**- IS element length in base-pairs
- **Direction**- Direction of the insertion sequence based on the direction of the transposase Open Reading Frame (ORF). 1 means forward, -1 means reverse.
- **ORF**- The Genbank locus tag (or tags) of the transposase ORFs in this IS element
- **Product**- the identified product of the transposase ORF (in the case of multiple ORFs, the first is taken)
- **Family**- The identified family of the insertion sequence, based on the *ISfinder* database of IS elements (see <http://www-is.biotoul.fr/>).
- **Group**- The group of the insertion sequence (subclass of *ISfinder* family)
- **IRL**- The left (defined as upstream of the coding sequence) inverted repeat sequence, if one is found
- **IRR**- The right (defined as downstream of the coding sequence) inverted repeat sequence, if one is found

## Algorithm

OASIS finds multiple-copy IS elements in each genome by identifying conserved regions surrounding transposase genes. First, groups of already-annotated transposase genes that could compose multiple copies of an IS were identified by length and sequence similarity. Those that fit a high similarity threshold were assumed to be in the same group.

The edges of each element were then found by finding a region of conservation around each group of transposases. The windows upstream and downstream of a group of transposases were compared. We are thus given a set of sequences

and asked to find the edge at which conservation ends, allowing for mismatches within IS elements, in a more computationally efficient matter than a multiple alignment. The edge of conservation was found via the following maximum likelihood approach.

We assumed that upstream and downstream regions each consist of a contiguous conserved region (within the IS) and an unconserved region (outside the IS). For computational efficiency, it was assumed that there were no gap mutations in the conserved region. Consider the set of upstream or downstream  $n$ -windows, in a group that contains  $m$  transposases. Let  $k$  be the length of the conserved region- the length which the IS element extends past the transposase on this side.

The first  $k$  characters of the sequences each come from multinomial distributions that are identical across the  $m$  sequences and are biased heavily towards a single nucleotide (the true sequence of the IS element). Let  $c_j$  be the consensus nucleotide for the  $j$ th character of the true sequence, and let  $x_{i,j}$  represent the  $j$ th character of the  $i$ th sequence.

$$\begin{array}{cccc} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ & \vdots & & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{array}$$

The distribution of  $x_{i,j}$  is assumed to be follows:

$$x_{i,j} \sim \begin{cases} \text{Multinom}(\mathbf{p}_{c_j}) & : i \leq k \\ \text{Multinom}(\mathbf{p}_b) & : i > k \end{cases}$$

Where  $\mathbf{p}_{c_j}$  represents the distribution of a conserved nucleotide with consensus value  $c_j$ , and  $\mathbf{p}_b$  represents a background distribution vector (the distribution of nucleotides in unconserved regions).  $\mathbf{p}_c$  can be any of 4 vectors that are biased towards a consensus nucleotide, with some probability of error. Let  $c, \varepsilon$  be adjustable parameters specifying the probability of a consensus match or an error in a conserved region.

The value of  $k$  is then found that maximizes the likelihood of the set of sequences. The likelihood is found by a dynamic programming algorithm as follows. For each nucleotide  $N \in \{A, G, C, T\}$ , let  $N_i$  be the count of nucleotides  $N$  in  $x_{1\dots m,i}$ . Also let  $M_i$  be the count of the most common nucleotide (the MLE estimate of the consensus nucleotide in a conserved region). The log-likelihood function is then

$$l(k) = \sum_{i=1}^k \log\left(\frac{m!}{A_i!C_i!T_i!G_i!} c^{M_i} \varepsilon^{m-M_i}\right) + \sum_{i=k+1}^n \log\left(\frac{m!}{A_i!C_i!T_i!G_i!} \left(\frac{1}{4}\right)^m\right)$$

Our estimate of  $k$ ,  $\hat{k}$ , is thus

$$\hat{k} = \underset{k}{\operatorname{argmax}} l(k)$$

OASIS computes this through a dynamic programming algorithm.

The putative edge of the IS element is then marked as  $k$  nucleotides from the edge of the transposase. This process is repeated for the windows upstream and downstream of the repeated transposase, for each group of transposases in the genome.

Once the edge of conservation has been found, the edges were then checked for inverted repeats using a Smith-Waterman alignment between the regions surrounding either edge of an IS element. If inverted repeats are found that disagree with the putative edges, the edges are adjusted to a limited extent.

Single-copy IS elements were found separately by finding transposases that were not placed in multiple-copy groups, and checking for inverted repeats in the surrounding regions. A Smith-Waterman alignment of the upstream and reverse-complement of the downstream regions was used to recognize significant inverted repeats and thus possible edges, a method first developed by IScan.

Once groups of IS elements were identified, a sample IS element from the set is selected, and `blastn` (NCBI) was used against the genome sequence to identify missing and partial copies of the IS element.

The family and group of each IS element is identified using `blastp` on the identified IS element ORFs against the ISFinder database, and classifying based on the family and group of the best match.

## Tutorial

What follows is an example use of OASIS for annotation of a prokaryotic genome.

### Input

The example genome is the NC\_004088 chromosome from *Yersinia pestis* KIM, a Gram-negative bacterium belong to the Enterobacteriaceae family. The genome is available to be downloaded from

`ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/Yersinia_pestis_KIM/NC_004088.gbk`

Once the Genbank file has been downloaded, move it into the OASIS folder.

### Usage

To annotate the genome, execute the command

```
./OASIS -g NC_004088.gbk -o Yersinia_Pestis_Annotations
```

### Output

Once the annotation is complete, there are two output files specifying the annotated IS elements in the genome. The text file, `Yersinia_Pestis_Annotations.txt`, will have the format specified in the **Output** section above. The fasta file, `Yersinia_Pestis_Annotations.fasta`, will have the nucleotide sequence and any protein sequences of one example from each group.