

## Tool for Rapid Annotation of Microbial SNPs (TRAMS) Manual

**Licence terms:** This script/program is provided free to academic users, allowing them to use, share or modify them but any commercial use/exploitation will require a written consent from the authors.

**Disclaimer:** The program/scripts have been tested on a variety of datasets and works well. However, the authors *provide no guarantees/warranties about the performance, fitness for a particular purpose, or any other warranties whether expressed or implied. The authors take no responsibility and are not liable for direct, indirect, special, incidental, or consequential damages resulting from the use, misuse, or inability to use this software under any circumstances.*

**Contact:** Vartul Sangal (vartul.sangal@strath.ac.uk); Paul A. Hoskisson (paul.hoskisson@strath.ac.uk)

**Bug reports:** Richard A. Reumerman (richard.reumerman@strath.ac.uk)

**Citation:** Reumerman RA, Tucker NP, Herron PR, Hoskisson PA, Sangal V. Tool for Rapid Annotation of Microbial SNPs (TRAMS): a simple program for rapid annotation of genomic variation in prokaryotes, *Antonie van Leeuwenhoek Journal of Microbiology* (*submitted*)

### 1. Introduction

TRAMS is a program, written using Python 2.7.3, for functional annotation of genomic SNPs as synonymous, nonsynonymous or nonsense. It separates nonsynonymous SNPs in start and stop codons as non-start and non-stop SNPs, respectively and SNPs in multiple overlapping features are annotated separately for each feature. Multiple nucleotide polymorphisms (MNPs) within a codon are combined before annotation.

### 2. Installation

TRAMS was developed using Python version 2.7.3, but works under any version of the 2.7 series. The users need to install python, biopython and Numpy to use this program. These are available to download from:

Python 2.7	<a href="http://www.python.org/download/">http://www.python.org/download/</a> ,
BioPython	<a href="http://biopython.org/wiki/Download">http://biopython.org/wiki/Download</a>
NumPy	<a href="http://sourceforge.net/projects/numpy/files/NumPy/">http://sourceforge.net/projects/numpy/files/NumPy/</a>
TRAMS	<a href="http://sourceforge.net/projects/strathtrams/files/Latest/">http://sourceforge.net/projects/strathtrams/files/Latest/</a>

The distributions of Python 2.7.3, NumPy and BioPython 1.60 are covered by different third party licences.

We provide a quick guide for installation of Python 2.7.3, NumPy and BioPython 1.60 but the users are highly encouraged to carefully read the instructions and detailed procedure from the download pages (mentioned above) for WINDOWS, Linux and Mac OS platforms.

## Windows

TRAMS is available to download as a single file executable for WINDOWS users who do not need to install Python, BioPython or NumPy.

However, if they wish to use TRAMS as Python script, they can download WINDOWS installers for Python 2.7.3, NumPy and BioPython 1.60 from the above websites in the mentioned order.

## Linux

Some Linux distributions such as Ubuntu may have Python 2.7.x pre-installed and the users can find out the version by using the command **python** in the console. The detailed download and installation instructions for different Linux distributions are available at the Python, BioPython and NumPy download pages (mentioned above). Here, we provide quick commands for Ubuntu using apt-get package manager:

```
sudo apt-get install python-numpy python-biopython
```

## Mac OS

Mac OS often have Python pre-installed but if you need to install python, NumPy or BioPython, **fink** may be useful. The detailed instructions to use **fink** are available at <http://finkers.wordpress.com/installing-fink/#install-fast.install>.

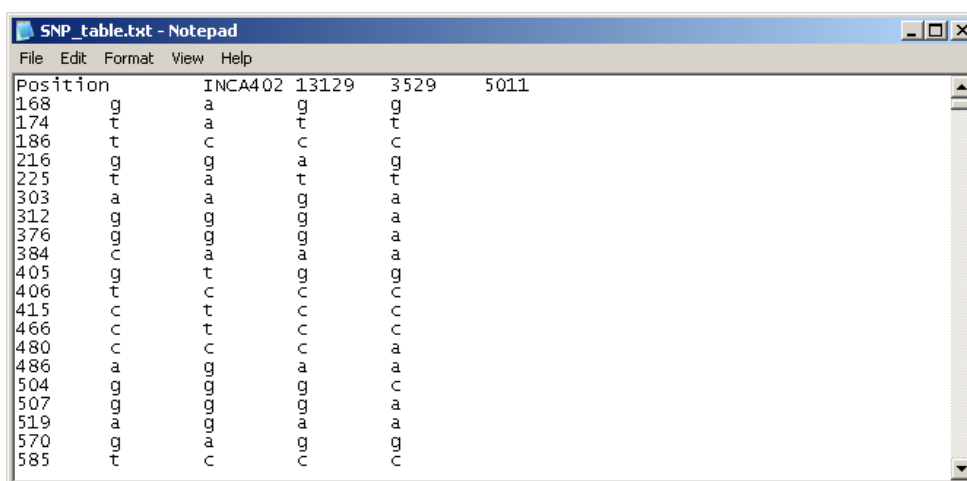
To install Python, NumPy and BioPython on a MacOS, following command can be used:

```
fink install biopython-py27
```

That should be sufficient. See the BioPython installation instructions for more information: <http://biopython.org/DIST/docs/install/Installation.html#htoc28>.

## 3. Using the program

**i. Input file format:** TRAMS uses a tab delimited text file containing SNP locations, reference nucleotide and SNPs in variant strains along with a reference genome sequence in GenBank or EMBL format.



Position	INCA402	13129	3529	5011
168	g	a	g	g
174	t	a	t	t
186	t	c	c	c
216	g	g	a	g
225	t	a	t	t
303	a	a	g	a
312	g	g	g	a
376	g	g	g	a
384	c	a	a	a
405	g	t	g	g
406	t	c	c	c
415	c	t	c	c
466	c	t	c	c
480	c	c	c	a
486	a	g	a	a
504	g	g	g	c
507	g	g	g	a
519	a	g	a	a
570	g	a	g	g
585	t	c	c	c

Figure 1: Input file format for TRAMS

As shown in the Fig, 1, the first row in the SNP input file is the header, followed by the data from the second row containing SNP positions in the reference genome (first column), nucleotide in the

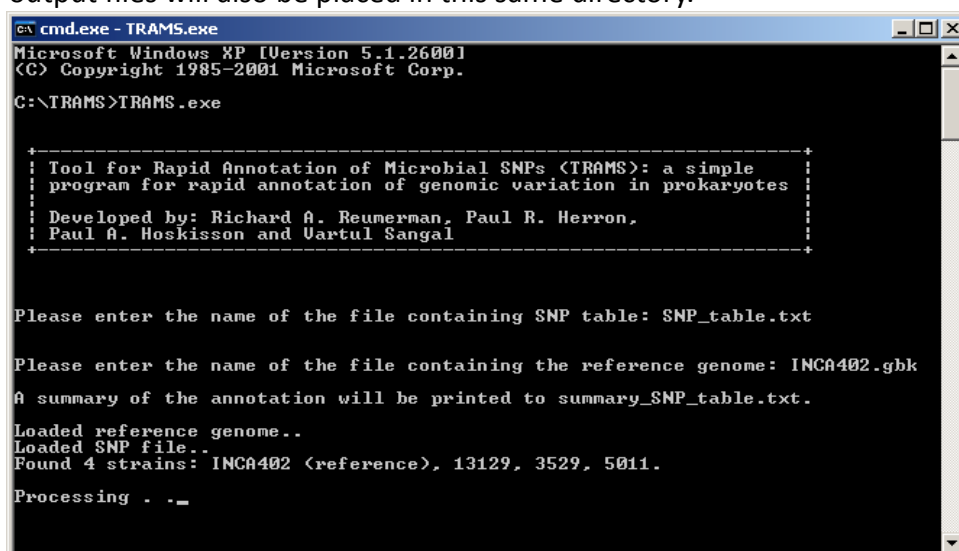
reference strain (second column) and SNP alleles in the variant strains (3 column onwards). SNPs in multiple strains can be added to the same input file as separate columns. If a base is not a SNP in a certain strain, it can either be left blank or kept the same as the reference base. In both cases the base will be ignored for that strain in the annotation file. There is no limit to the number of strains that can be added that may be dependent on the memory of your computer.

The GenBank and EMBL formats are recognized directly by the file extensions (.gbk or .embl, respectively).

We have separately provided two scripts to convert the SNPs tables extracted from Mauve (*mauve2trams.py*) and vcf files (*vcf2trams.py*) into TRAMS input. Please see the user manual for generating TRAMS input files from vcf and Mauve output files for details.

## ii. Running the program:

**Windows:** TRAMS Windows executable can be run without any installation requirements. Double-click on the icon to start the program. A console window will open and prompt for the SNP input table and the reference genome file. Make sure the program is placed in the same directory as the input files. The output files will also be placed in this same directory.



```
cmd.exe - TRAMS.exe
Microsoft Windows XP [Version 5.1.2600]
(C) Copyright 1985-2001 Microsoft Corp.
C:\TRAMS>TRAMS.exe

+-----+
+  Tool for Rapid Annotation of Microbial SNPs (TRAMS): a simple
+  program for rapid annotation of genomic variation in prokaryotes
+-----+
+  Developed by: Richard A. Reumerman, Paul R. Herron,
+  Paul A. Hoskisson and Uartul Sangal
+-----+

Please enter the name of the file containing SNP table: SNP_table.txt

Please enter the name of the file containing the reference genome: INCA402.gbk
A summary of the annotation will be printed to summary_SNP_table.txt.
Loaded reference genome..
Loaded SNP file..
Found 4 strains: INCA402 <reference>, 13129, 3529, 5011.
Processing . -_
```

Figure 2: A snapshot of the TRAMS console from the executable file

You can also run the program, both the executable and the python script, using the following command line:

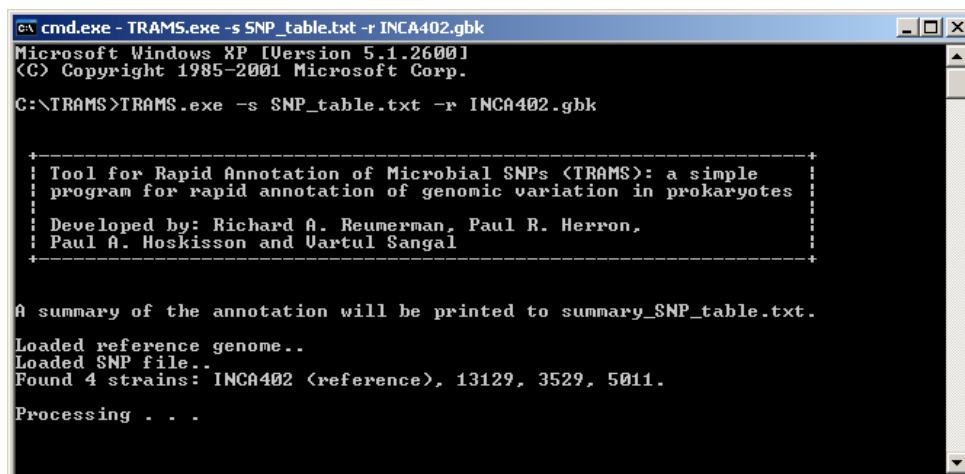
**TRAMS.exe -s SNP\_Table.txt -r Reference.gbk**                      or

**TRAMS.py -s SNP\_Table.txt -r Reference.gbk**

-s        name of the input file

-r        Reference genome in GenBank or EMBL format

The users can use “-h” to display the help message.



```
cmd.exe - TRAMS.exe -s SNP_table.txt -r INCA402.gbk
Microsoft Windows XP [Version 5.1.2600]
(C) Copyright 1985-2001 Microsoft Corp.

C:\TRAMS>TRAMS.exe -s SNP_table.txt -r INCA402.gbk

+-----+
+ Tool for Rapid Annotation of Microbial SNPs (TRAMS): a simple
+ program for rapid annotation of genomic variation in prokaryotes
+-----+
+ Developed by: Richard A. Reumerman, Paul R. Herron,
+ Paul A. Hoskisson and Vartul Sangal
+-----+

A summary of the annotation will be printed to summary_SNP_table.txt.

Loaded reference genome..
Loaded SNP file..
Found 4 strains: INCA402 (reference), 13129, 3529, 5011.

Processing . . .
```

**Figure 3: Running TRAMS using the command line**

TRAMS script will need the correct versions and the necessary modules (Python 2.7.3, BioPython 1.60 and NumPy) for proper functioning. Using the script with different versions (e.g. Python 3.0) cannot be guaranteed to work. The script will not run under Python 2.6, since TRAMS uses the 'argparse' module, which is absent from 2.6. Once the necessary software is installed the script can be run in the same way as the windows executable, using the same command line code. See the installation instructions for more details on setting up the system for TRAMS.

### iii. Output files

TRAMS will produce three output files. These are placed in the working directory of the program is running (where the input files are saved). The file names used are the same as the SNP tables, but with prefixes. In the following, the name "SNP\_table.txt" is used for demonstration purposes.

**Annotation\_SNP\_table.txt:** This file contains the annotation of every SNP in the SNP table in tab delimited format. The first column contains the SNP positions. The following columns describe the feature the SNP is located in (if any), stating feature type, locus tag, gene name, gene product, feature start, feature end, feature strand, base, codon and amino acid. After this the SNP type, mutant base, new codon and new amino acid are given for each strain. If a base in a strain is the same as in the reference strain, the row will be empty for that strain.

If multiple SNPs occur in a single codon, all but the last one in that codon will be labelled "MNP" (Multiple Nucleotide Polymorphism). The annotation of the last SNP in the codon will then take the cumulative effect of all SNPs into account. Mutant types used to annotate SNPs and MNPs in TRAMS are *synonymous* (no change in amino acid), *nonsynonymous* (different amino acid), *nonsense* (normal codon replaced by stop codon), *nonstop* (stop codon replaced by normal codon) and *nonstart* (start codon at start of feature replaced by a normal codon). For this annotation, codon table 11 is used (Bacterial and Plant Plastid) as provided in the BioPython module.

Position	Feature	Locus tag	Gene	Product	Start	End	Strand	Ref. base	Ref. codon	Ref. res.	SNP type	New base	New codon	New res.	SNP type	New base
168 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	TTG	L	synonymous	a	TTa	L				
174 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	GTT	V	synonymous	a	GTa	V				
186 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	CGT	R	synonymous	c	CGc	R			synonymous	c
216 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	GGG	G							synonymous	a
225 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	ATT	I	synonymous	a	ATa	I				
303 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 A	CAA	Q							synonymous	g
312 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	CCG	P								
376 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	GTA	V								
384 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 C	CAC	H	nonsynonymous	a	CAa	Q			nonsynonymous	a
405 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	CCG	A	synonymous	t	GCt	A				
406 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	TCG	S	nonsynonymous	c	cCG	P			nonsynonymous	c
415 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 C	CCG	P	nonsynonymous	t	ICG	S				
466 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 C	CAC	H	nonsynonymous	t	IAC	Y				
480 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 C	ACC	T								
486 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 A	CAA	Q	synonymous	g	CAg	Q				
504 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	CAG	Q								
507 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	TCG	S								
519 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 A	CAA	Q	synonymous	g	CAg	Q				
570 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	CAG	Q	synonymous	a	CAa	Q				
585 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	TCT	S	synonymous	c	TCc	S			synonymous	c
624 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 A	AAA	K	synonymous	g	AAg	K				
667 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	TTG	L							synonymous	a
732 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 C	GTC	V								
756 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	TAT	Y	synonymous	c	TAc	Y				
762 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	CCG	P	synonymous	a	CCa	P			synonymous	a
813 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	GCG	A	synonymous	a	GCa	A				
816 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	ATT	I	synonymous	c	ATc	I				
846 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 A	AAA	K	synonymous	g	AAg	K				
849 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 C	CTC	L							synonymous	t
900 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 G	TCG	S	synonymous	a	TCa	S			synonymous	a
933 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 C	CGC	R								
951 CDS	CDB402_0001	dnaA	chromosomal replication initiation protein	1	1680	1 T	GAT	D								

Figure 4: A snapshot of annotation file

**Summary\_SNP\_table.txt:** The summary file provides a brief overview of the results of the annotation process. Most importantly, it states the number of SNPs in each strain and how they are divided over different feature and SNP types.

```

+-----+
| Tool for Rapid Annotation of Microbial SNPs (TRAMS): a simple |
| program for rapid annotation of genomic variation in prokaryotes |
| Developed by: Richard A. Reuerman, Paul R. Herron, |
| Paul A. Hoskisson and Vartul Sangal |
+-----+

Found 4 strains: INCA402 (reference), 13129, 3529, 5011.

Finished annotation. Total time: 29.5 s

Annotations were written to file 'annotation_snpsbelref.txt'
A list of SNPs in overlapping features is written to file 'overlap_snpsbelref.txt'

+ Strain 13129:
39238 SNPs found
Number of SNPs found CDS features: 34860
(of which in pseudogenes: 148)
- MNPs: 3089
- Synonymous: 22830
- Nonsynonymous: 8880
- Nonsense: 33
- Nonstart: 7
- Nonstop: 21
Intergenic: 4370
tRNA: 8

+ Strain 3529:
27783 SNPs found
Number of SNPs found CDS features: 24517
(of which in pseudogenes: 35)
- MNPs: 1391
- Synonymous: 17143
- Nonsynonymous: 5959
- Nonsense: 15
- Nonstart: 2
- Nonstop: 7
Intergenic: 3260
  
```

Figure 5: A snapshot of the summary output

**Overlap\_SNP\_table.txt:** Since TRAMS annotates SNPs once for each feature they occur in, they will be annotated more than once when features overlap. This is reported in the overlap file, where a

list of SNPs that occur more than once is given. This file is in tab-delimited format and contains the first and last of the range of overlapping SNPs of each feature overlap. It then gives some information on the two overlapping features.

	A	B	C	D	E	F	G	H	I	J
	SNP start	SNP end	Feature 1	Locus tag	Product	Feature 2	Locus tag	Product		
1	20682		CDS	CDB402_0017	hypothetical protein	CDS	CDB402_0018	alpha amylase catalytic region		
2	79562		CDS	CDB402_0077	hypothetical protein	CDS	CDB402_0078	hypothetical protein		
3	152387	152390	CDS	CDB402_0147	putative secreted polysaccharide deacetylase	CDS	CDB402_0148	GntR family regulator		
4	226779		CDS	CDB402_0209	two-component system response regulator	CDS	CDB402_0210	two-component system sensor histidine kinase		
5	293521	293545	CDS	CDB402_0277	putative secreted hydrolase	CDS	CDB402_0278	glucose-1-phosphate thymidyltransferase		
6	340251		CDS	CDB402_0319	putative cation-transporting P-type ATPase	CDS	CDB402_0320	hypothetical protein		
7	350471		CDS	CDB402_0331	hypothetical protein	CDS	CDB402_0332	1,4-dihydroxy-2-naphthoate octaprenyltransferase		
8	665160	665182	CDS	CDB402_0636	putative DNA helicase II	CDS	CDB402_0637	hypothetical protein		
9	763382		CDS	CDB402_0725	phosphonobiosylglycinamide formyltransferase	CDS	CDB402_0726	bifunctional phosphonobiosylaminoimidazolecarboxamide formyltransferase/IMP cyclohydrolase		
10	932556		CDS	CDB402_0890	galactose-1-phosphate uridylyltransferase	CDS	CDB402_0891	galactokinase		
11	1273635		CDS	CDB402_1175	hypothetical protein	CDS	CDB402_1176	hypothetical protein		
12	1409800		CDS	CDB402_1293	putative aminotransferase biotin synthesis related protein	CDS	CDB402_1294	putative secreted protein		
13	1714029		CDS	CDB402_1575	hypothetical protein	CDS	CDB402_1576	putative cobalamin biosynthesis protein		
14	1722968		CDS	CDB402_1581	arabinofuranosyl transferase C	CDS	CDB402_1582	putative alpha/beta hydrolase fold family protein		
15	1764038		CDS	CDB402_1620	putative type B carboxylesterase	CDS	CDB402_1621	isopentenyl-diphosphate delta-isomerase		
16	1786912		CDS	CDB402_1639	copper-sensing transcriptional repressor csoR	CDS	CDB402_1639	putative potassium-efflux system protein		
17	1825692	1825695	CDS	CDB402_1674	valyl-tRNA synthetase	CDS	CDB402_1675	malate dehydrogenase		
18	1842215		CDS	CDB402_1689	hemoglobin-like protein	CDS	CDB402_1689	hypothetical protein		
19	1980866		CDS	CDB402_1813	putative secreted protein	CDS	CDB402_1814	trehalose 6-phosphate phosphatase, biosynthetic		
20	1981604		CDS	CDB402_1814	trehalose 6-phosphate phosphatase, biosynthetic	CDS	CDB402_1815	LacI-family transcriptional regulator		
21	2059676		CDS	CDB402_1896	hypothetical protein	CDS	CDB402_1897	carboxylate-amine ligase		
22	2118663		CDS	CDB402_1962	putative oxidoreductase	CDS	CDB402_1963	putative hydrolase		
23	2210405	2210498	CDS	CDB402_2037	putative integral membrane protein	CDS	CDB402_2038	hypothetical protein		
24	2311601		CDS	CDB402_2119	hypothetical protein	CDS	CDB402_2120	methionine and cysteine biosynthesis regulator		
25										
26										
27										
28										
29										
30										
31										
32										
33										
34										

Figure 6: A snapshot of list containing overlapping SNPs

## Algorithm

The way TRAMS annotates SNPs is outlined in the diagram below. The sequence is divided into regions containing either a feature or a featureless region (intergenic). SNPs are annotated region by region, ensuring that SNPs in overlapping features are annotated once for every feature that spans their position. Multiple SNPs in a single codon are combined, resulting in the annotation of their cumulative effect on the codon rather than their effect in isolation.

