# Allele Intersection Analysis: a novel tool for multi locus sequence assignment in multiply infected hosts

a story about *Wolbachia*, cloning, sequencing and set theory written by

## Wolfgang Arthofer[1], Markus Riegler[2], Karl Moder[3], Daniela Schneider[4], Wolfgang J Miller[4] and Christian Stauffer[5]

[1] Molecular Ecology Group, Institute of Ecology, University of Innsbruck, Austria
[2] Centre for Plants and the Environment, University of Western Sydney, Australia
[3] Institute of Applied Statistics and Computing, BOKU, University of Natural Resources and Applied Life Sciences, Vienna, Austria
[4] Centre of Anatomy and Cell Biology, Medical University of Vienna, Austria
[5] Institute of Forest Entomology, BOKU, University of Natural Resources and Applied Life Sciences, Vienna, Austria
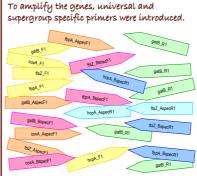
Correspondence: wolfgang.arthofer@uibk.ac.at    http://www.uibk.ac.at/ecology/    http://www.peerart.at/aw

**Back in 2006 ...**

... Baldo et al. published a Multi Locus Sequence Typing (MLST) system for the endosymbiont *Wolbachia*.

The MLST uses sequence information of 5 genes to unambiguously define a *Wolbachia* strain. All information can be found in a web database.

To amplify the genes, universal and supergroup specific primers were introduced.

Today, MLST is the standard for describing a *Wolbachia* strain!

Consider a *Wolbachia* strain as a set containing MLST sequences as elements. And lets name this strain here simply 'red'.

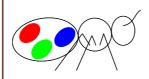Different strains contain different elements – but maybe they even share some: here the blue and the green *Wolbachia* have the same gatB! All sequences of a distinct strain give the MLST sequence type, a unique strain identifier.

A multiply infected insect is again a set, containing different strains as elements. We call the taxative list of strains infecting one individual the 'infection type'. The infection type of this insect is {red, blue, green}.

It's easy to clone and sequence the MLST amplicons of this insect ...

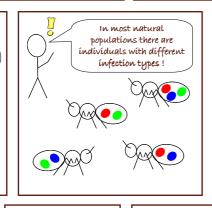| coxA | gatB | fbpA | hcpA | ftsZ |
|---|---|---|---|---|
| Allele 1 | Allele 1 | Allele 1 | Allele 1 | Allele 1 |
| Allele 2 | Allele 2 | Allele 2 | Allele 2 | Allele 2 |
| Allele 3 | | Allele 3 | Allele 3 | Allele 3 |

... and we know now that the insect was triple infected. Unfortunately, by cloning the alleles lost their 'color' – we do not know which alleles belong together.

If there would be only one strain from the A anb B supergroup, specific primers would fix the problem ...

... but the cherry fruit fly, for instance, harbours 5 *Wolbachia* strains, and 3 of them are A-group !

We need some algorithm to define the sequence types !

In most natural populations there are individuals with different infection types !

Let's build a diagnostic system that can easily identify an individual's infection type:

red F    red R
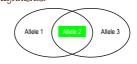blue F    blue R
green F    green R

Strain specific PCR primers can do this job, targeting a highly variable gene – wsp could be a candidate, or a single copy VNTR, or .......

Now we search for two individuals that share only one *Wolbachia* strain.

We amplify, clone and sequence one MLST gene, and compare the alignments.

Allele 1   Allele 2   Allele 3

The allele found in both individuals must belong to the green strain !

Allele 1   Allele 2   Allele 3

Furthermore, it's obvious that allele 1 comes from the red strain, and allele 3 from the blue one !

Any combination of individual infection types that allows the assignment of all strains is called 'informative'.

Look at this informative combination in a 4-fold infected species:
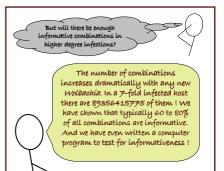
A
B
C

We will first resolve the red strain by intersecting the alignments of individuals A and C. The remaining sequence from A must belong to green.

A                    C
Allele 1   Allele 2   Allele 3
                      Allele 4

From Individual B we know now already all sequences except one. This must belong to the blue strain!

B
Allele 1
Allele 2
Allele 3

Finally, the last unknown sequence in C belongs to yellow.

But will there be enough informative combinations in higher degree infections?

The number of combinations increases dramatically with any new *Wolbachia*. In a 7-fold infected host there are 8935415775 of them ! We have shown that typically 60 to 80% of all combinations are informative. And we have even written a computer program to test for informativeness !

The end.