# Defining biologically meaningful molecular operational taxonomic units

M. Göker

# Why molecular taxonomy?
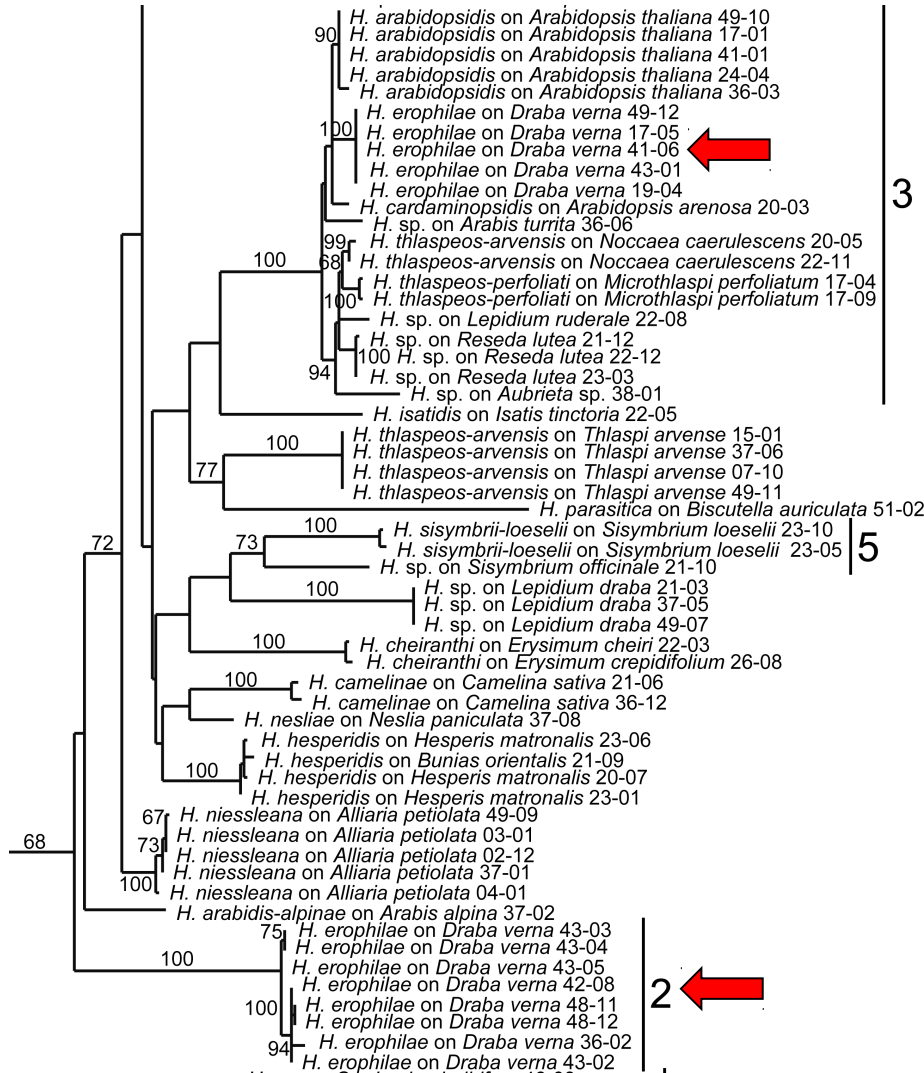
**Definition:**
Establishing a (informal or even formal) taxonomy of organisms based only on molecular sequences
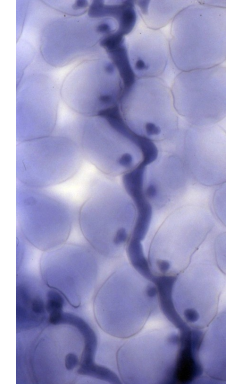
**Uses:**
- Detection of cryptic and pseudocryptic species
- Detection of misidentifications and mislabelled sequences in public databases
- Identification of juvenile specimens
- Analysis of environmental samples (e.g. metagenomics)

# Example: (pseudo-)cryptic species



ITS/LSU rDNA data of the genus *Hyaloperonospora* (Peronosporales, Oomycetes) (Göker et al. 2009)

=> Two genetically distinct but microscopically identical species on *Draba verna* host plants

# Threshold-based clustering

- Calculate distance *d(i,j)* between each pair of sequences *i* and *j*

- Define a threshold *T*

- Principle: if *d(i,j) <= T*, assign *i* and *j* to the same molecular operational taxonomic unit (MOTU)
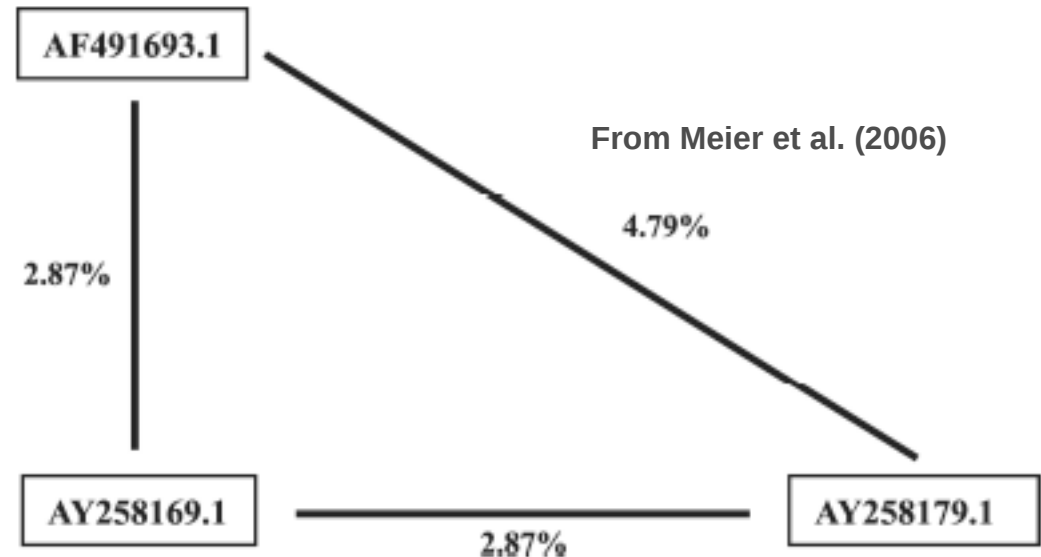


**From Meier et al. (2006)**

FIGURE 1. Pairwise distances for three *Anopheles* sequences (AF491693.1, AY258179.1 = *A. maculipennis*; AY258169.1 = *A. messae*). All belong to the same 3% DNA profile, although one pairwise distance exceeds the threshold.

=> Can lead to inconsistencies if formulated in that way

Optimizing molecular taxonomy

# Impact of the clustering algorithm

- A distance *d(i,j) <= T* is called <u>link</u>

- An additional parameter, the "linkage fraction" *F*, determines how many links between an object and a cluster are necessary to include the object in the cluster
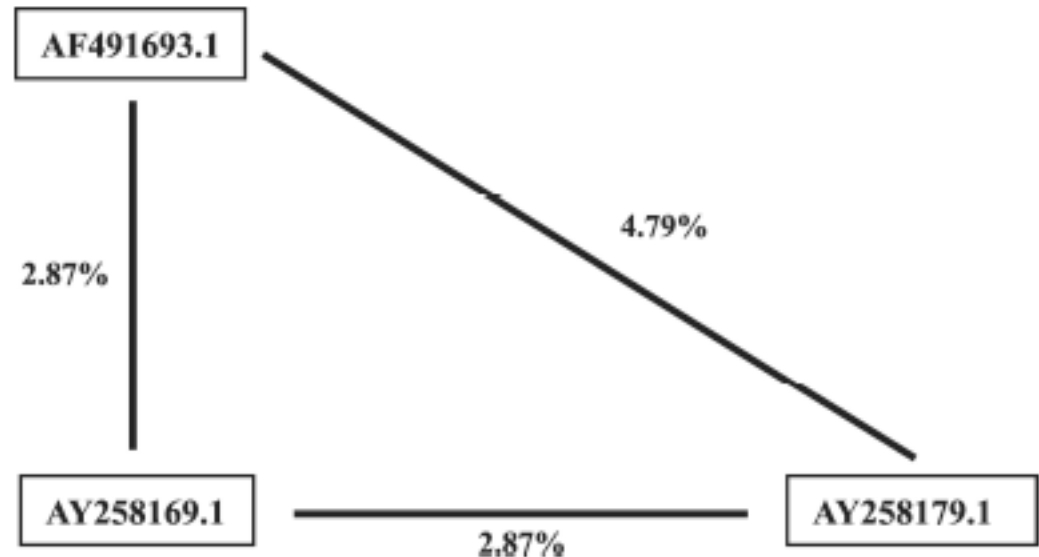


FIGURE 1. Pairwise distances for three *Anopheles* sequences (AF491693.1, AY258179.1 = *A. maculipennis*; AY258169.1 = *A. messae*). All belong to the same 3% DNA profile, although one pairwise distance exceeds the threshold.
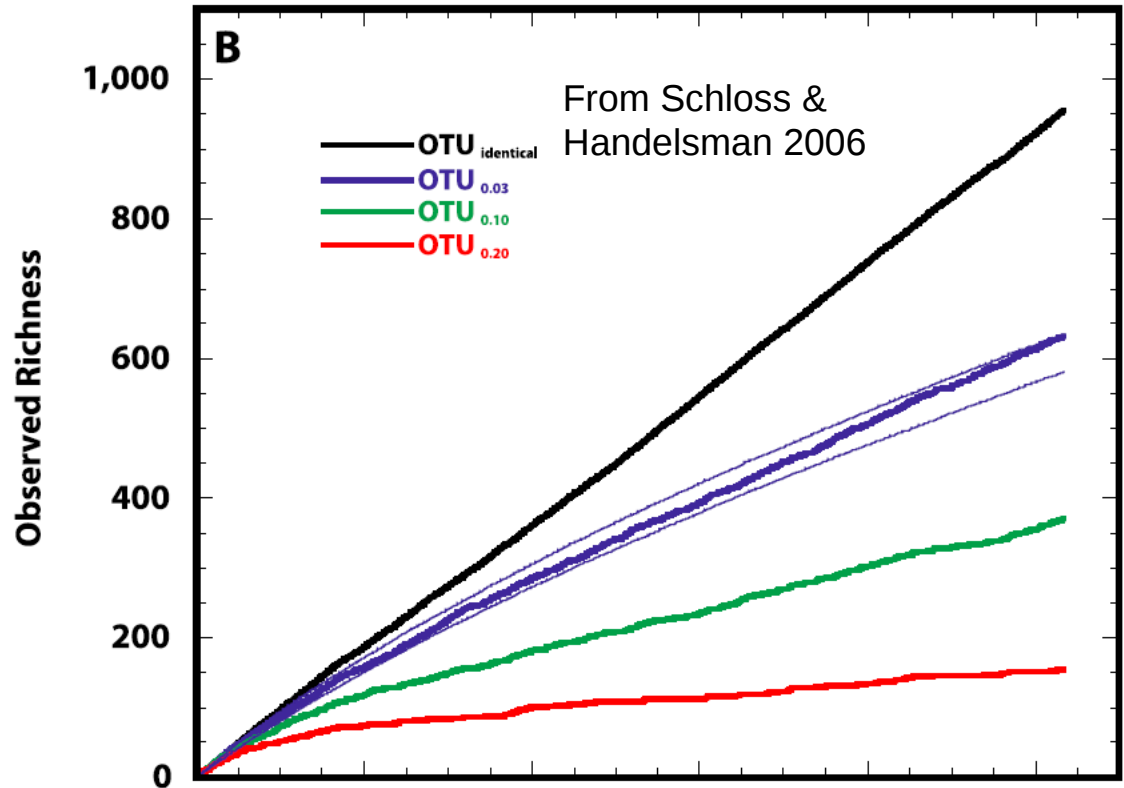
=> Here, 1 cluster for *F* <= 0.5, but 2 clusters otherwise!
=> Lower *F* values allow higher within-cluster divergence

Example:

- Species richness of soil bacteria estimated from 16S rDNA sequences

- Question: Has saturation been obtained?

- Obvious dependency on *T*



From Schloss & Handelsman 2006

=> Choice of parameters has serious consequences for total biodiversity estimates

# The debate between traditional and molecular taxonomists

Ongoing intense (and sometimes hostile) debate between molecular taxonomists and traditional morphologists, particularly in the context of DNA barcoding

**Criticisms of molecular taxonomy:**

• Values of $T$ used for clustering differ in the literature, even if applied to the same groups of organisms and molecular markers

• Values of $T$ are often based on subjective criteria or on a tradition that emerged in recent years for the sake of comparability between studies

• Genetic divergence may differ between morphologically defined lineages

• A smaller distance (or a higher similarity) does not necessarily indicate a closer phylogenetic relationship

=> How can we maximize the agreement between traditional and molecular taxonomy?

# Clustering optimization

- Partition := non-hierarchical, non-overlapping classification

- Many biological data are represented as partitions (e.g. assignment of sequences to species):

- Non-hierarchical clustering also results in a partition, e.g.:

| Accession number | Organism |
|---|---|
| EF050035 | Pseudoperonospora cubensis |
| EF174888 | Peronospora aestivalis |
| EF174890 | Peronospora sepium |
| EF174891 | Peronospora fulva |
| EF174894 | Peronospora lathyri-verni |
| EF174944 | Peronospora orobi |
| ... | ... |

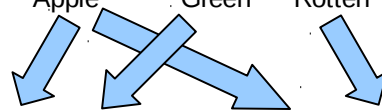| Accession number | Cluster number |
|---|---|
| EF050035 | 29 |
| EF174888 | 26 |
| EF174890 | 25 |
| EF174891 | 24 |
| EF174894 | 27 |
| EF174944 | 27 |
| ... | ... |

Approach:
- Use set of specimens identified using traditional techniques as reference points
- Determine the clustering parameters that maximize the agreement with a reference partition
- Do not require that full agreement can be obtained

Optimizing molecular taxonomy

# Comparing partitions

3 example partitions of 7 objects

| Object | Fruit type | Colour | Condition |
|---|---|---|---|
| A | Apple | Green | Fresh |
| B | Lemon | Yellow | Fresh |
| C | Cherry | Red | Rotten |
| D | Apple | Green | Fresh |
| E | Cherry | Red | Fresh |
| F | Lemon | Yellow | Rotten |
| G | Apple | Green | Rotten |

**Observed values**

| | Same | Different |
|---|---|---|
| Same | 5 | 0 |
| Different | 0 | 16 |

| | Same | Different |
|---|---|---|
| Same | 1 | 4 |
| Different | 8 | 8 |

**Expected values**

| | Same | Different |
|---|---|---|
| Same | 1.19 | 3.81 |
| Different | 3.81 | 12.19 |

| | Same | Different |
|---|---|---|
| Same | 2.14 | 2.86 |
| Different | 6.86 | 9.14 |

**Rand Index**

(5+16)/(5+0+0+16) = 1.0

(1+8)/(1+4+8+8) = 0.43

**Expected Index**

(1.19+12.19)/(5+0+0+16) = 0.64

(2.14+9.14)/(1+4+8+8) = 0.54

**Modified Rand Index (MRI)**

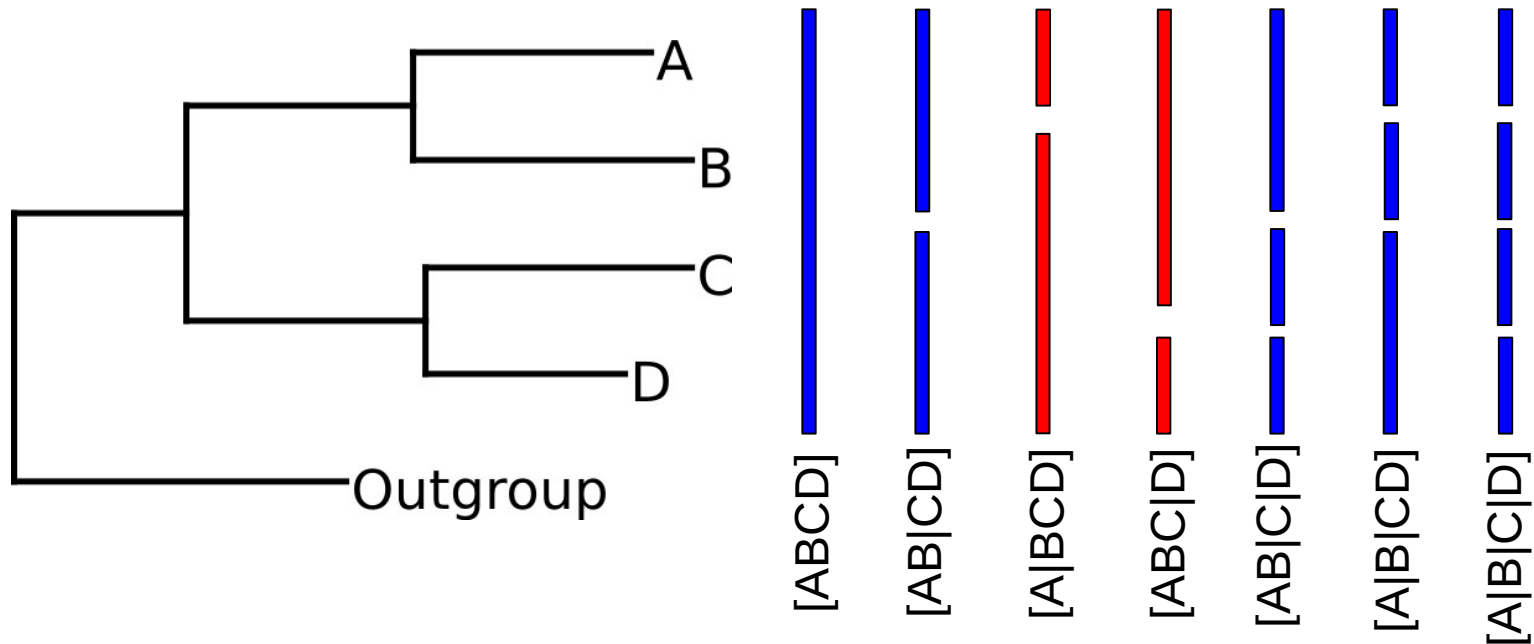(1.0-0.64)/(1.0-0.64) = 1.0

(0.43-0.54)/(1.0-0.54) = -0.24

• Rand index (Rand 1971): traverse all pairs of objects and determine proportion of those being in the same cluster in *both* partitions or in a different cluster in *both* partitions

• Modified Rand index (Hubert & Arabie 1985): corrects for chance (by relating to the expected Rand index for two random partitions with the same cluster number and sizes)

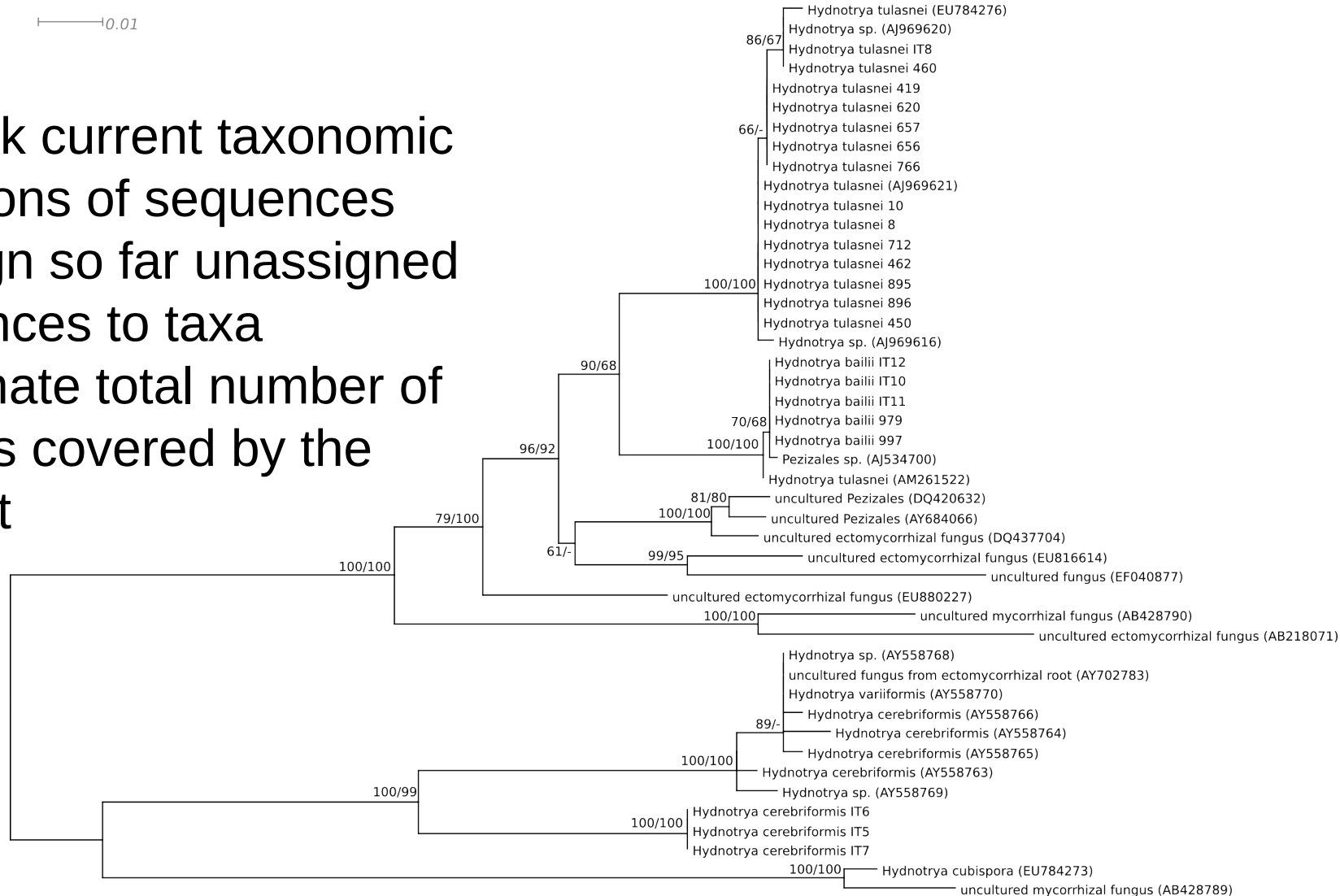Optimizing molecular taxonomy

# Why trees don't help



A phylogenetic tree rules out certain classifications (e.g. red ones), but is compatible with many others (blue ones)

# Example: *Hydnotrya* ITS rDNA

**Tasks**

- Check current taxonomic affiliations of sequences
- Assign so far unassigned sequences to taxa
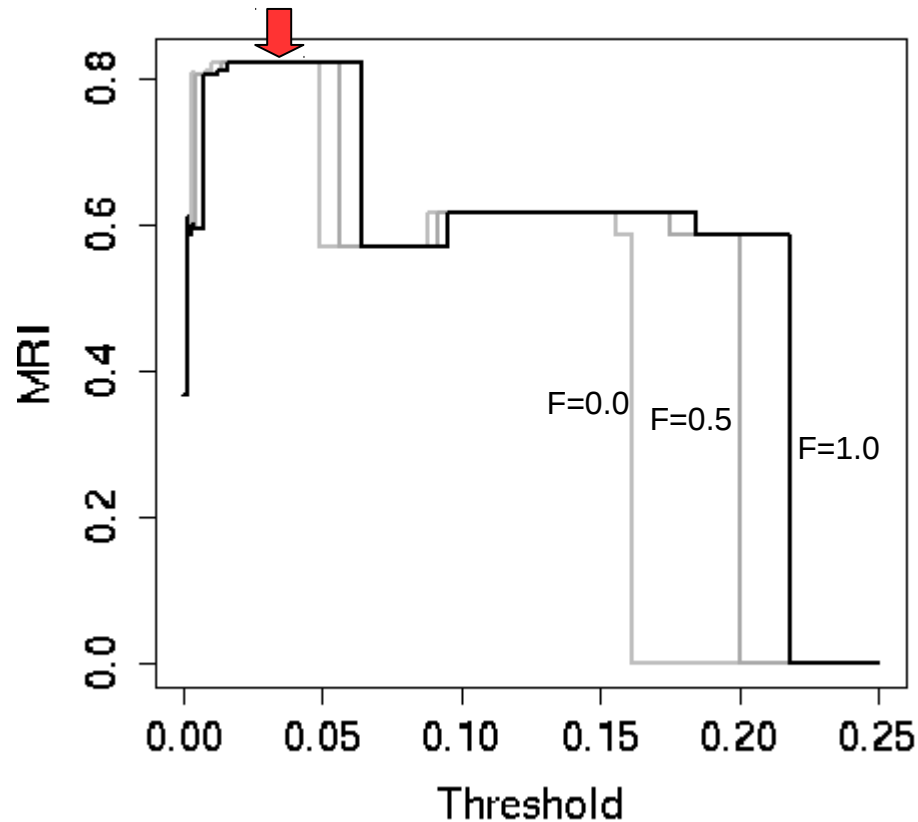- Estimate total number of species covered by the dataset



Optimizing molecular taxonomy

# Example: *Hydnotrya*

## Procedure

1) Restrict dataset to sequences with taxonomic affiliations



| Accession number | Organism | Species name present? |
|---|---|---|
| EU784276 | Hydnotrya tulasnei | Yes |
| AJ969620 | Hydnotrya sp. G-Ht | No |
| AJ969621 | Hydnotrya tulasnei | Yes |
| AJ969616 | Hydnotrya sp. LB-Ht | No |
| AJ534700 | Pezizales sp. B48 | No |
| AM261522 | Hydnotrya tulasnei | Yes |
| DQ420632 | uncultured Pezizales | No |
| ... | ... | ... |

2) Conduct clustering optimization with reduced dataset

3) Place sequences without taxonomic affiliations back in the dataset

4) Conduct clustering with all sequences and optimized parameters

Optimizing molecular taxonomy

# An overlooked hypogeous fungus



*Hydnotrya tulasnei* ascocarp
(picture: G.Hensel)



*Hydnotrya bailii* ascocarp
(picture: G. Hensel)

Stielow et al., under review:
Distinction between *Hydnotrya bailii* Soehner (1959) and *Hydnotrya tulasnei* (Berk.) Berk. & Broome (1846) has been neglected for 50 years!

Optimizing molecular taxonomy

**Revised taxonomy:**

- *H. tulasnei*
- *H. bailii* incl. 1 „*H. tulasnei*"
- *H. cubispora*
- *H. cerebriformis* I incl. 1 „*H. variiformis*"
- *H. cerebriformis* II
- 6 unnamed species
- 7 accessions assigned to a taxon via clustering



Optimizing molecular taxonomy

# Example: *Peronospora* ITS rDNA



*Peronospora* sp. on *Ocimum basilicum*

## Tasks

• Revise nomenclature of all Genbank *Peronospora* ITS rDNA sequences

• Check whether a combination of molecular and host plant characters is sufficient to obtain a consistent species concept

Optimizing molecular taxonomy

# Example: *Peronospora*

## Procedure

1) Restrict dataset to (a) sequences with taxonomic affiliations and (b) sequences with interpretable host names

| Accession number | Organism | Specific host | Species name present? | Host present? |
|---|---|---|---|---|
| EF614964 | Peronospora variabilis | Chenopodium album | Yes | Yes |
| EF614958 | Peronospora sp. SMK20063 | Chenopodium ambrosioides | No | Yes |
| EF614957 | Peronospora sp. DAR45530 | Chenopodium ambrosioides | No | Yes |
| EF614955 | Peronospora farinosa f. sp. chenopodii | Chenopodium hybridum | Yes | Yes |
| EF174939 | Peronospora sp. GG133 | | No | No |
| EF174924 | Peronospora sp. HV956 | | No | No |
| EF174970 | Peronospora trifoliorum | | Yes | No |
| EF174963 | Peronospora trifoliorum | | Yes | No |

2) Conduct clustering optimization with reduced datasets (a) and (b)

3) Check for coincidence of results (i.e. of optimal clustering parameters)

4) Place sequences without taxonomic affiliations or host information back in the dataset

5) Conduct clustering with all sequences and optimized parameters

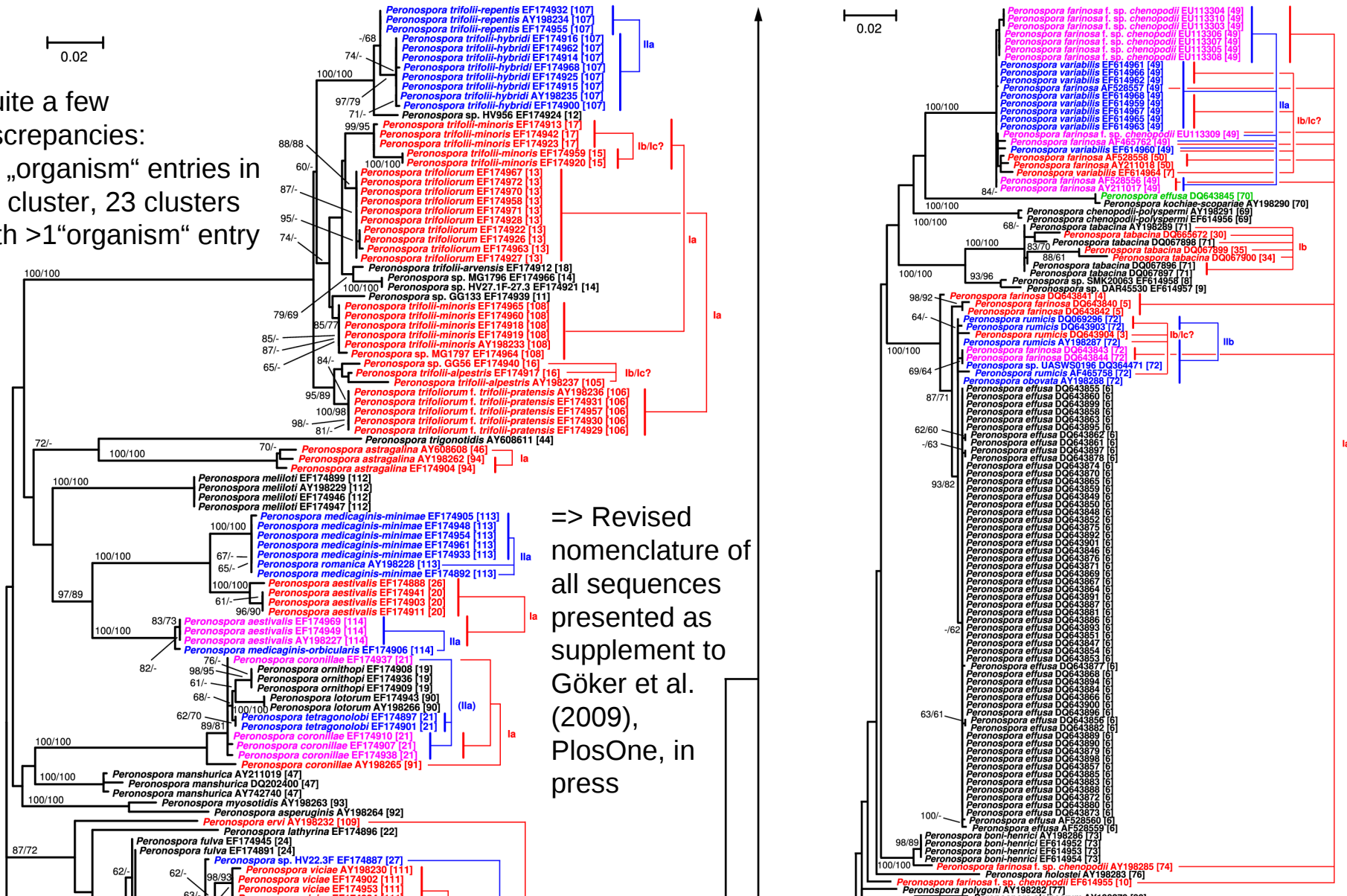Optimizing molecular taxonomy

# Example: *Peronospora* ITS rDNA



- Taxonomy-based optimization: best result (MRI=0.85485) with T=0.0075 and F=1.0 (left picture: thick lines)

- Host-based optimization: best result (MRI=0.85204) with T=0.0075 and F=1.0 (left picture: thin lines) => *exactly the same optimum*

- Resulting in 117 clusters

# Example: *Peronospora* ITS rDNA

Quite a few discrepancies:
20 „organism" entries in >1 cluster, 23 clusters with >1"organism" entry

=> Revised nomenclature of all sequences presented as supplement to Göker et al. (2009), PlosOne, in press

# Robustness against sampling bias



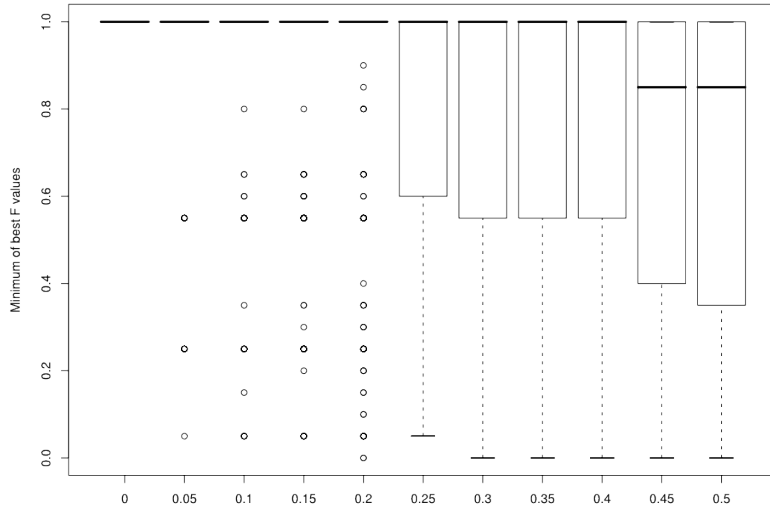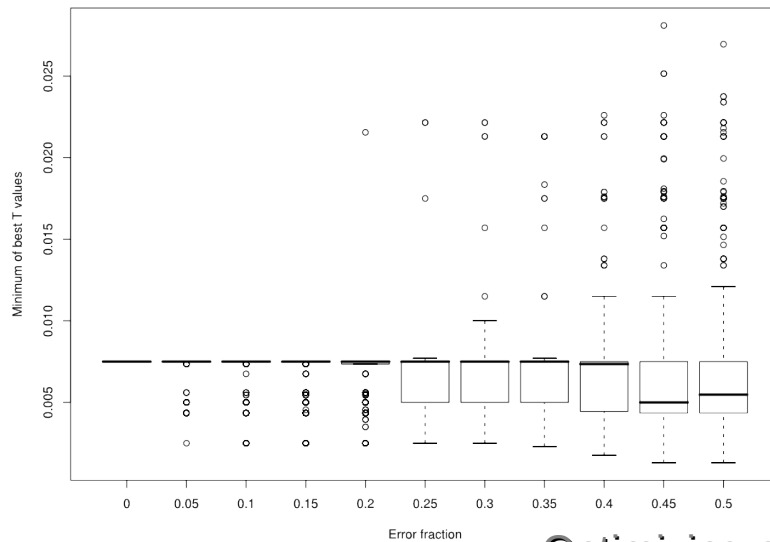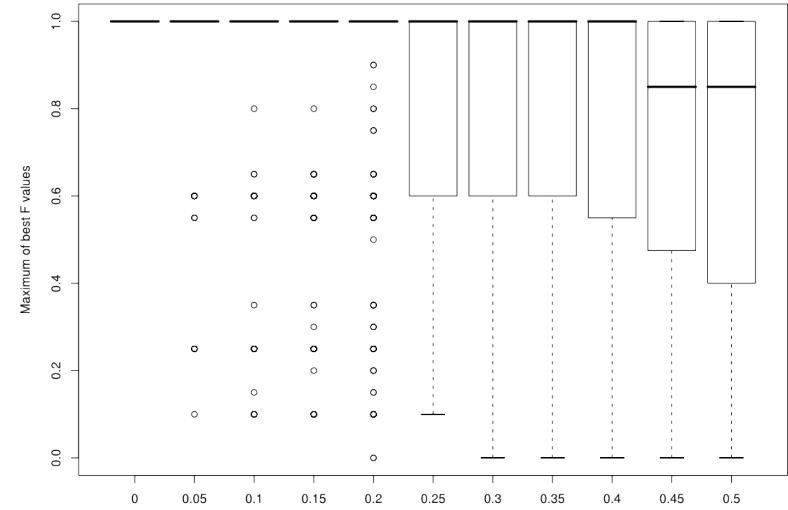Optimizing molecular taxonomy

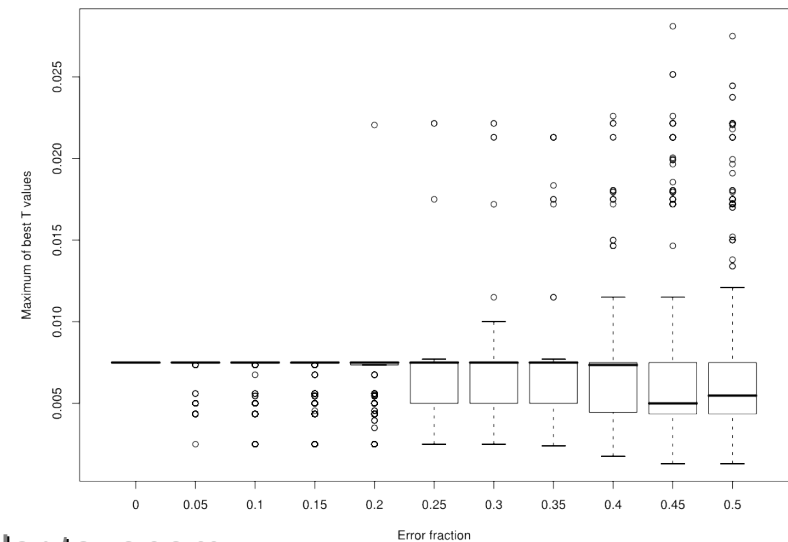Figure 6

Figure 7

Figure 8

Figure 9

Optimizing molecular taxonomy

# Summary

Clustering optimization based on the agreement between partitions...

- leads to MOTUs with highest agreement to traditional taxonomy, but it is robust against errors in such a reference partition

- connects traditional and modern taxonomic disciplines by specifically addressing the issue of how to optimally account for both traditional species concepts and genetic divergence

- can also be used together with different types of reference partitions (e.g. host species of specialized parasites/mutualists)

- leads to biologically reasonable choices for clustering parameters that are also suitable for sequence identification

- is implemented in the OPTSIL software available at http://www.goeker.org/mg/clustering for all major operating systems