# In House pipeline: NG-Tax

## Validation based on controls

Open & reproducible microbiome data analysis spring school
Wageningen, The Netherlands, May 28-30, 2018

Gerben DA Hermes, PhD

Laboratory of Microbiology,

Wageningen University & Research

WAGENINGEN
UNIVERSITY & RESEARCH

100years
1918 — 2018

The Graduate School

VLAG

Turun yliopisto
University of Turku

microbiome

# Tested parameters

- <span style="color:red">Primer pair</span>

- <span style="color:red">Barcoding strategy</span>

- Chimera filtering

- Filtering parameters

- OTU picking

- Database

- Etc…

And try to come up with an 'optimal' combination of each

# In house analysis pipeline: NG-Tax



**F1000Research**
Open for Science

BROWSE    SUBJECTS    GATEWAYS    HOW TO PUBLISH ⌄    ABOUT ⌄    BLOG

Search

Check for updates

RESEARCH ARTICLE

NG-Tax, a highly accurate and validated pipeline for analysis of 16S rRNA amplicons from complex biomes [version 1; referees: 2 approved with reservations, 1 not approved]

✉ Javier Ramiro-Garcia[1-3*], Gerben D. A. Hermes[1,2*], Christos Giatsis[4], Detmer Sipkema[2], Erwin G. Zoetendal[1,2], Peter J. Schaap[1,3], Hauke Smidt[2]

* Equal contributors

+ Author details
+ Grant information

METRICS

1897
👁 VIEWS

435
⬇ DOWNLOADS

Get PDF
Get XML

WAGENINGEN
UNIVERSITY & RESEARCH

# NG-Tax analysis & settings
# Galaxy interface

# Validation:

## Positive controls/Synthetic or Mock communities



Qualitative & quantitative

1   2   3   4

- 54 species
- 32 genera

- 7 phyla

**Firmicutes** ■ **Bacteroidetes** ■ **Actinobacteria** □ **Proteobacteria**
■ **Verrucomicrobia** □ **Fusobacteria** □ **Lentisphaerae**

Resolution

How low can we go

WAGENINGEN
UNIVERSITY & RESEARCH
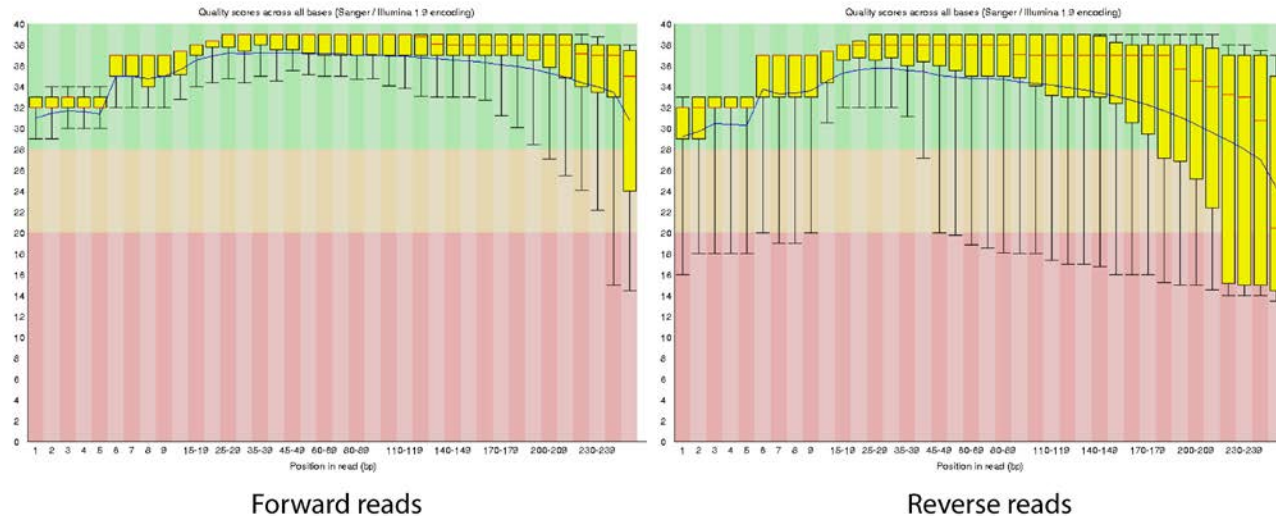
# Effect of quality filtering

- Bokulich 2013, Quality-filtering vastly improves diversity estimates from Illumina amplicon sequencing. Nat methods. supplementals

# Settings optimization

- No quality filtering based on phred score

    Phred 10 = 90% accuracy, 30 99.9%



Forward reads          Reverse reads

- Filter by abundance -> high quality

- "QC" = demultiplexing (matching valid barcodes)

- 'Short' reads: 2x70nt (enough for identification) = less room for error.

WAGENINGEN
UNIVERSITY & RESEARCH

# OTU classification: Classify ratio

# 95% sequence similarity: genus level
# Based on full length 1500nt SILVA alignment

95%    100%    99%

95% sequence similarity: genus level



SILVA: Blautia/Faecalibacterium

Same sequence/different classification

'choose'?

# Real example: very clear



Nr hits that are *Feacalibacterium*

100% match in the SILVA DB

nr hits with 100% match in the SILVA DB

All hits in the SILVA DB with 100%

99.34% of the 100% matches of that sequence (OTU) in the SILVA database are *Feacalibacterium*

# Real example: less clear



Assigned: Enterobacteriacaea

2.7% Pantoea

67% @100% match with 3704 hits -> back to family

25% Citrobacter

# NG-Tax classification accuracy & confidence

## 2x70nt information

# Validation

Sources of error/noise:

- 4 different synthetic community types (each in triplicate)
- 2 different variable regions V4, V5-V6
- 2 primer pairs (515F-806R, 784F-1064R)
- 25-35 PCR cycles
- 7 libraries
- 3 different sequencing runs

# NG-Tax performance: alpha diversity

## Rarefaction plots

# NG-Tax performance: beta diversity



Weighted Unifrac

# NG-Tax performance: beta diversity



Unweighted Unifrac

Retain or retrieve the relevant **high resolution biological information** despite high noise levels

Independent of
- primer type/variable region
- Sample diversity/complexity
- Sequencing run

# Biom (OTU table)

| #OTU ID | Mc.1.t.I56 | Mc.1.1.I01 | Mc.1.2.I01 | Mc.1.3.I01 | Mc.2.t.I56 | Mc.2.2.I01 | Mc.2.3.I01 | Mc.3.t.I56 | Mc.3.2.I01 | Mc.3.3.I01 | Mc.4.t.I56 | Mc.4.2.I01 | Mc.4.3.I01 | taxonomy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 565610 | 1000 | 3242 | 1578 | 1646 | 2000 | 461 | 292 | 2000 | 46 | 87 | 0 | 6 | 12 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 565628 | 0 | 763 | 359 | 408 | 0 | 144 | 56 | 0 | 10 | 20 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 565653 | 0 | 176 | 0 | 79 | 0 | 26 | 17 | 0 | 3 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 565686 | 0 | 0 | 0 | 0 | 2000 | 423 | 139 | 2000 | 24 | 45 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 565688 | 0 | 0 | 0 | 0 | 1000 | 272 | 92 | 1000 | 9 | 35 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 565696 | 0 | 0 | 0 | 0 | 1000 | 153 | 83 | 1000 | 15 | 25 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656106 | 0 | 0 | 0 | 0 | 0 | 64 | 26 | 0 | 0 | 9 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656108 | 0 | 0 | 0 | 0 | 1000 | 83 | 40 | 1000 | 5 | 34 | 2500 | 15 | 42 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656111 | 0 | 0 | 0 | 0 | 1000 | 166 | 104 | 1000 | 3 | 29 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656114 | 0 | 0 | 0 | 0 | 0 | 83 | 35 | 0 | 6 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656124 | 0 | 0 | 0 | 0 | 0 | 42 | 14 | 0 | 6 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656128 | 0 | 0 | 0 | 0 | 0 | 34 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656146 | 0 | 0 | 0 | 0 | 0 | 40 | 29 | 0 | 2 | 11 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656150 | 0 | 0 | 0 | 0 | 0 | 35 | 6 | 0 | 2 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656156 | 0 | 0 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656161 | 0 | 0 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656164 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656174 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656202 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656207 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |
| 5656254 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | k__Bacteria; p__Actinobacteria; c__Actinobacteria; o__Bifidobacteriales; f__Bifidobacteriaceae; g__Bifidobacterium |

# Data interpretation: NG-Tax

- Reads/OTU -> confidence

- ~3 reads not very confident. 30.000 reads very confident

- Classification ratio 0.8, maybe check

- Preferably at genus level, OTUs contain too much risk/error, unless the OTU is very 'clean'

- Preferably use phylogenetic information because OTU based methods are very sensitive to noise

- Check classification of interesting OTUs

- Classify Nas using whatever method (blast is ugly)