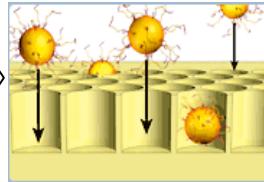
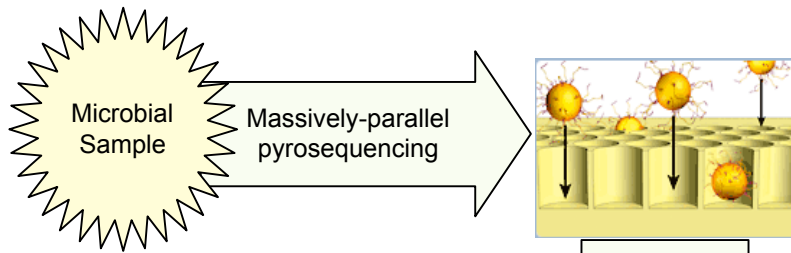


Global Alignment for Sequence Taxonomy (GAST)



Trimming and quality filtering

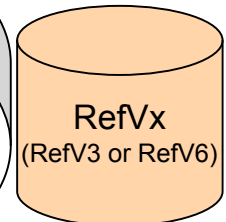
~400,000 hypervariable region tags

Creating RefSSU and RefVx

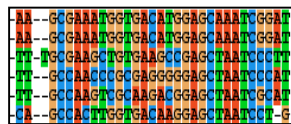
1. Full-length SSU rRNA reference sequences are downloaded from SILVA and low-quality sequences removed.
2. Taxonomy is assigned with RDP, additional taxonomy sources (e.g. reference genomes) are added.
3. The V3 and V6 regions are excised from the ARB alignment using primer locations.
4. Additional filters remove low-quality reference tags to create RefV3 or RefV6.

BLAST

Align each tag to top 100 RefVx BLAST hits



MUSCLE



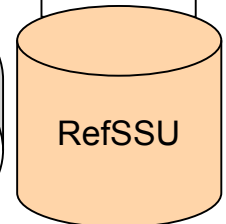
Calculate distance from tag to each RefVx top 100 hits

Excise high-quality variable region reference tags

Best GAST RefVx hit(s)

SQL Query

Retrieve taxonomy of all RefSSU sequences containing RefVx hit(s)



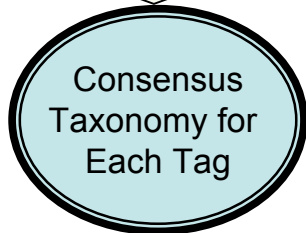
Taxonomic strings

```

2 Bacteria;Firmicutes
107 Bacteria;Firmicutes;Clostridia;Clostridiales
6 Bacteria;Firmicutes;Clostridia;Clostridiales;Clostridiaceae
31 Bacteria;Firmicutes;Clostridia;Clostridiales;Clostridiaceae;Bryantia
1 Bacteria;Firmicutes;Clostridia;Clostridiales;Lachnospiraceae
    
```

Calculate consensus taxonomy at 66% majority

(Bacteria;Firmicutes;Clostridia;Clostridiales)



The GAST Process

1. A microbial sample is sequenced using hypervariable-region specific primers.
2. Each tag has the primers trimmed and quality filters applied.
3. BLAST each high-quality tag against RefVx, a database of reference hypervariable sequences (RefV3, RefV6).
4. The tag is aligned against the top 100 BLAST hits.
5. The RefVx matches having the minimum pairwise distance to the tag are selected.
6. For each best RefVx match, all source RefSSU sequences are selected.
7. A consensus agreement of $\geq 66\%$ of selected RefSSU sources is calculated.
8. The consensus taxonomy is applied to the tag.