

# Meta-omics

Jasper Koehorst

Laboratory of Systems and Synthetic Biology



16S



Using the shape of the handle to explain the function of the oven

# Who what why... in summary

16S

To understand who is  
there up to genus level

Understand the diversity

\$

Meta-Genomics

To understand what is  
possible

Mostly high abundant  
players

\$\$

Meta-Transcriptomics

To understand why things  
are happening

Mostly high abundant  
players

\$\$



WAGENINGEN  
UNIVERSITY & RESEARCH



100 years  
1918 — 2018

# Generating the data

Sequencing techniques

Current generation

**Illumina (Many short reads, good quality)**

**Pacbio (Few long reads, medium quality)**

Next generation (which is already here)

**Oxford nanopore (Few very long reads, medium quality)**



WAGENINGEN  
UNIVERSITY & RESEARCH



# illumina

- Suitable for meta- samples
- Generates relatively short reads
- Assembly will result in many short contigs
- Cheap

16S

Genomics

Transcriptomics

# PacBio

- Suitable for meta- samples (when not alone)
- Generates relatively long reads ~15 kb
- Assembly will be difficult due to low coverage
- New development on whole rna sequencing
- Expensive

16S

Genomics

Transcriptomics



WAGENINGEN  
UNIVERSITY & RESEARCH



# Oxford

- Suitable for meta- samples (when not alone)
- Yield average per cell 2.3 Gb
- Needs less DNA than PacBio ( $\pm 20x$ ) (rumour)
- Generates long reads
- Assembly will be difficult due to low coverage
- New development on whole rna sequencing
- Expensive
- Still under heavy development and improvements

16S

Genomics

Transcriptomics

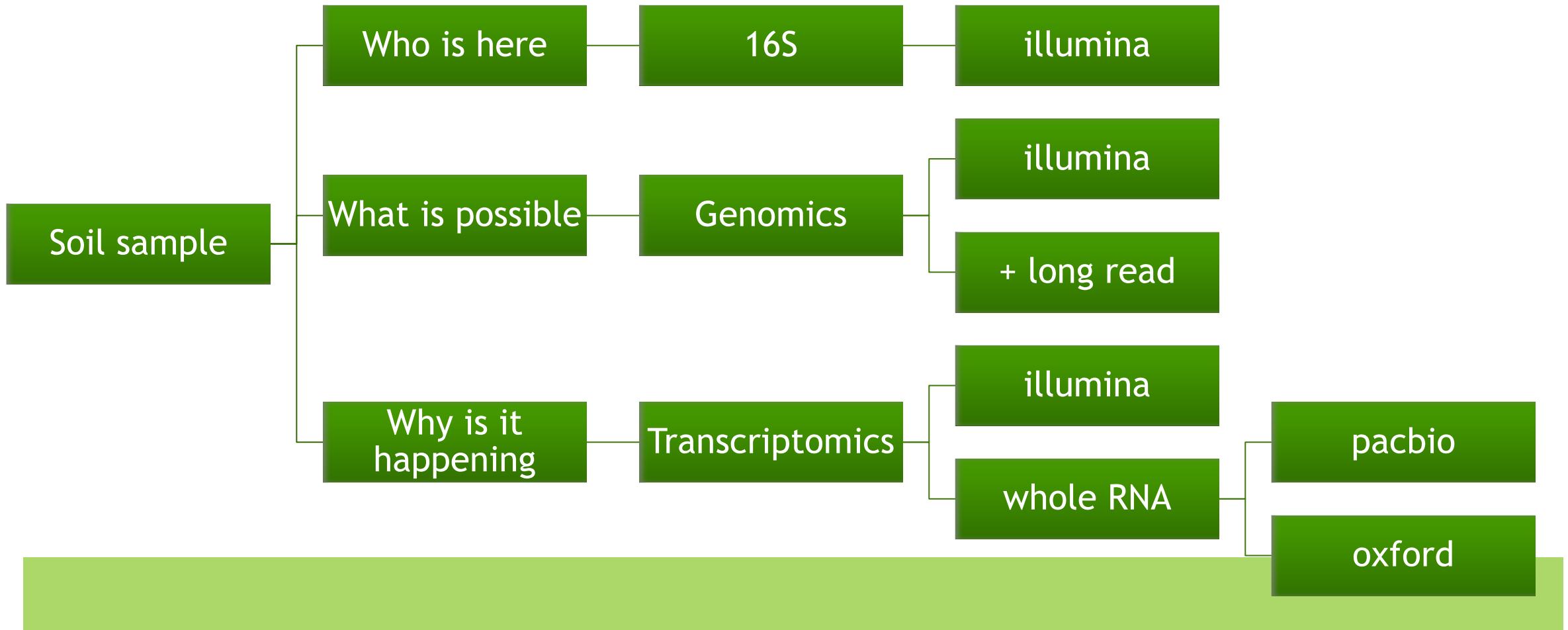


# Congratulations!

The first >1 Mb DNA read, achieved with  
nanopore sequencing

Martin Smith & collaborators, Kinghorn Centre for Clinical  
Genomics, Garvan Institute for Medical Research, Australia  
December 2017

# Sequencing a meta-sample



# What do I want to achieve?

- Have a general understanding who is here
- Extract specific complete genomes from a microbiome sample
- Identify pathways that are expressed under certain conditions
- Is it important to know who is there?



WAGENINGEN  
UNIVERSITY & RESEARCH



# When is illumina good enough?

- When you want to know what is happening
- When you do not want to extract complete genomes
- To perform expression analysis
  - You do **not** need metagenomics for this



WAGENINGEN  
UNIVERSITY & RESEARCH



# When are long reads good enough?

- When you want to extract complete genomes
- For the extraction of complete genomes but its not that easy
  - Use information obtained from short read sequencing
    - Antibiotic resistance, specific growth media, etc to isolate a strain of interest



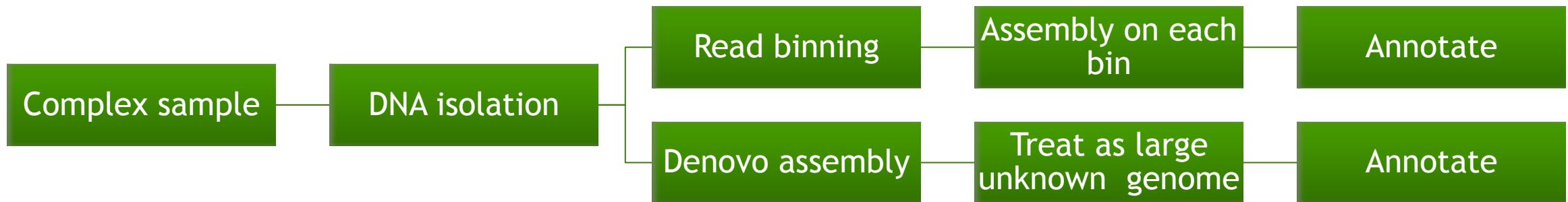
WAGENINGEN  
UNIVERSITY & RESEARCH



# Workflow

- Meta-genomics workflow
- Meta-transcriptomics workflow

# Meta-genomics workflow

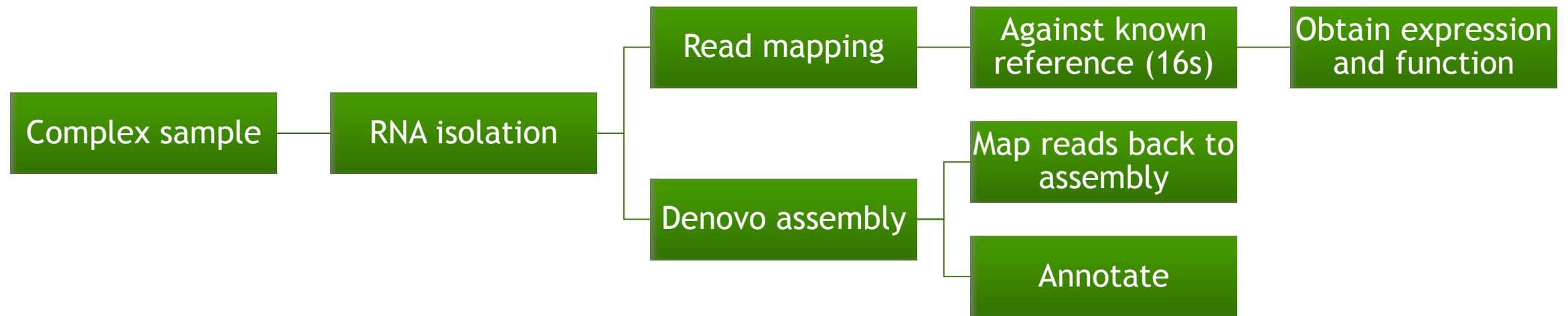


WAGENINGEN  
UNIVERSITY & RESEARCH



100 years  
1918 — 2018

# Meta-transcriptomics workflow



WAGENINGEN  
UNIVERSITY & RESEARCH



100 years  
1918 — 2018

# What to expect?

- 16s
  - Low abundant species will be picked up
- Metagenomics
  - Only high abundant groups are detected (2-4000x coverage required for the same resolution as 16s)
    - 1 Gene vs ±4000 genes
  - Many contigs ( fewer when long reads are used)
  - After annotation
    - You know what the theoretical possibilities are within your sample
    - Identify the landscape variation between different samples
- Transcriptomics
- Only high abundant groups are detected
  - Using a reference (difficult in an unknown or variable community)
    - Map reads, differential expression
    - De-novo assembly, annotate each transcript, ...
    - Obtain the genes of interest, identify the function,
      - Hopefully understand why this is happening

EULER-SR	WhatsHap	GARM	SOAPdenovo		SSPACE	HaploMerger		mip
A5	Telescopier	Contrail	fermi	GABenchmarkToB	QSRA	Opera	RAMPART	
SSAKE	SWAP-Assembler	Newbler		AutoAssemblyD	Dazzler	PCAP	Platanus	Arapan
TIGR	Mapsembler 2	ALLPATHS-LG		HapCompass			Forge	SHEAR
gapfiller	CloudBrush	Cortex	REAPR	VICUNA	Edena	CLC	PERGA	KmerGenie
Arachne	MIRA	dipSPAdes	MetAMOS	Ray	Tedna	TIGRA	Amos	
SCARPA	Celera	VCAKE	GAM	Geneious	SeqMan NGen		Nesoni	ATAC
GRIT	IDBA		PASHA				MetaVelvet-SL	Quast
Phrap	MaSuRCA	H					Scencher	BESST
CGAL	Curtain	SWiPS	KDASsembler	Metassembler		SGA		GGAKE
Pipeline Pilot	SHRAP	Taipan	SILP3		HGAP	PRICE	Pilon	MSR-CA
OMACC	Anchor	Omega	SUTTA	ABySS	IDBA-MTP			
SOPRA	iMetAMOS		DNAexus		HyDA-Vista	SR-ASM	Velvet	Enly
DNA Dragon	CABOG	SAGE		Ragout	SPAdes	Atlas		FRCBam
FALCON	SuccinctAssembly	SHORTY	Cerulean	Monument		SAT-Assembler		image
	GigaAssembler		SHARCGS	GAGM	ngsShoRT	ABBA		
		Lasergene	PBJelly	DecGPU	Khmer	GenoMiner		ELOPER

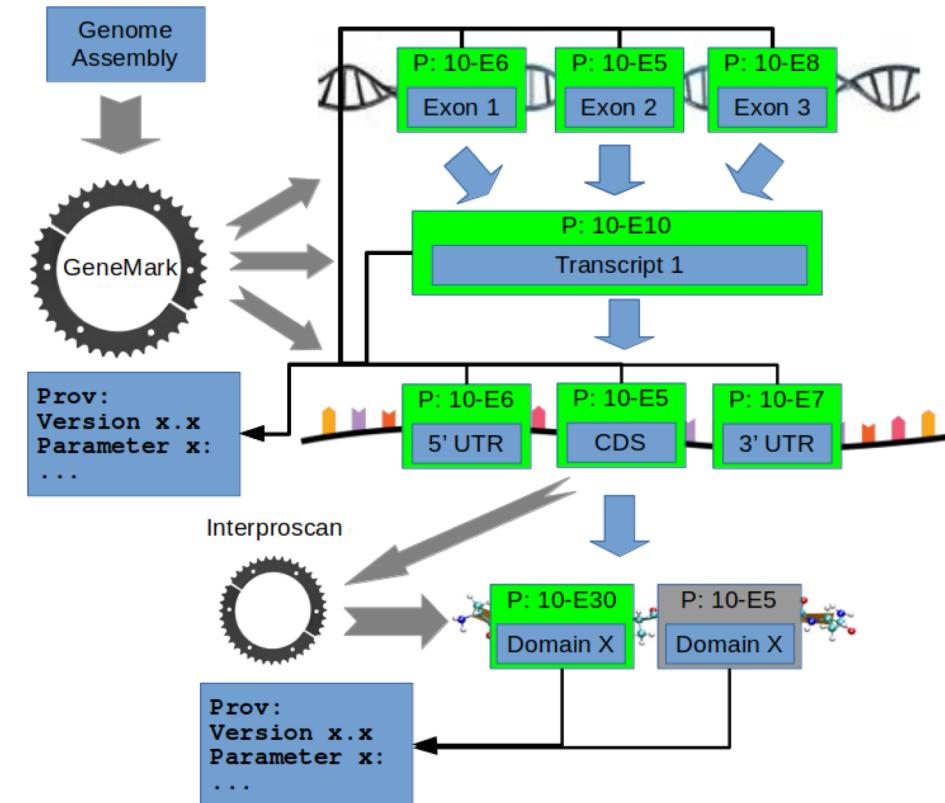
Putting the reads together

# Annotation information storage

Example questions e.g:

- What (functions) distinguish “environmental” samples?
- Which enzymes are there that can catalyze reaction X (maybe with different cofactors?)

Requires a resource of consistently annotated samples that can be easily mined



# Requirements for genome mining

- A semantic annotation platform that incorporates common tools and stores the results in “proper” format. **SAPP**
  - A graph database that can be mined: **SAGERP**
  - A definition of the “proper format”: definitions of biological terms and their relationships: **GBOL ontology**
  - Interface to use the ontology: **GBOL stack**
  - Tools to develop all of these:
    - **Empusa**: code generator
- SAPP is the only thing a user would need to use to annotate a genome
  - Sager-P is the only thing a user would need to mine the data



WAGENINGEN  
UNIVERSITY & RESEARCH



```

0 ##gff-version 3.2.1
1 ##sequence-region ctg123 1 1497228
2 ctg123 . gene      1000  9000 . + . ID=gene00001;Name=EDEN
3 ctg123 . TF_binding_site 1000  1012 . + . ID=tfbs00001;Parent=ge...
4 ctg123 . mRNA      1050  9000 . + . ID=mRNA00001;Parent=ge...
5 ctg123 . mRNA      1050  9000 . + . ID=mRNA00002;Parent=ge...
6 ctg123 . mRNA      1300  9000 . + . ID=mRNA00003;Parent=ge...
7 ctg123 . exon      1300  1500 . + . ID=exon00001;Parent=mRN...
8 ctg123 . exon      1050  1500 . + . ID=exon00002;Parent=mRN...
9 ctg123 . exon      3000  3902 . + . ID=exon0003;Parent=mRN...
10 ctg123 . exon     5000  5500 . + . ID=exon0004;Parent=mRN...
11 ctg123 . exon     7000  9000 . + . ID=exon0005;Parent=mRN...
12 ctg123 . CDS      1201  1500 . + 0 ID=cds00001;Parent=mRN...
13 ctg123 . CDS      3000  3902 . + 0 ID=cds00001;Parent=mRN...
14 ctg123 . CDS      5000  5500 . + 0 ID=cds00001;Parent=mRN...
15 ctg123 LOCUS    SCU49845 5028 bp DNA PLN 21-JUN-1999
16 ctg123 DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Axl2p
17 ctg123 (AXL2) and Rev7p (REV7) genes, complete cds.
18 ctg123 ACCESSION U49845
19 ctg123 VERSION U49845.1 GI:1293613
20 ctg123 KEYWORDS .
21 ctg123 SOURCE Saccharomyces cerevisiae (baker's yeast)
22 ctg123 ORGANISM Saccharomyces cerevisiae
Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
Saccharomycetales; Saccharomycetaceae; Saccharomyces.
23 ctg123 REFERENCE 1 (bases 1 to 5028)
24 ctg123 AUTHORS Torpey,L.E., Gibbs,P.E., Nelson,J. and Lawrence,C.W.
25 ctg123 TITLE Cloning and sequence of REV7, a gene whose function is required for
DNA damage-induced mutagenesis in Saccharomyces cerevisiae
26 ctg123 JOURNAL Yeast 10 (11), 1503-1509 (1994)
27 ctg123 PUBMED 7871890
28 ctg123 REFERENCE 2 (bases 1 to 5028)
29 ctg123 AUTHORS Roemer,T., Madden,K., Chang,J. and Snyder,M.
30 ctg123 TITLE Selection of axial growth sites in yeast requires Axl2p, a novel
plasma membrane glycoprotein
31 ctg123 JOURNAL Genes Dev. 10 (7), 777-793 (1996)
32 ctg123 PUBMED 8846915
33 ctg123 REFERENCE 3 (bases 1 to 5028)
34 ctg123 AUTHORS Roemer,T.
35 ctg123 TITLE Direct Submission
36 ctg123 JOURNAL Submitted (22-FEB-1996) Terry Roemer, Biology, Yale University, New
Haven, CT, USA
37 ctg123 FEATURES Location/Qualifiers
38 ctg123 source 1..5028
39 ctg123 /organism="Saccharomyces cerevisiae"
40 ctg123 /db_xref="taxon:4932"
41 ctg123 /chromosome="IX"
42 ctg123 /map="9"
43 ctg123 <1..206
44 ctg123 /codon_start=3
45 ctg123 /product="TCP1-beta"
46 ctg123 /protein_id="AAA98665.1"
47 ctg123 /db_xref="GI:1293614"
48 ctg123 /translation="SSIYNGISTSGLDLNNTIADMRLQLGIVESYKLKRAVVSSASEA
AEVLLRVDNIIRARPRTANRQHM"

```

# A “proper” format

- Dataset-wise and element-wise provenance
- Mining enabled
- Query enabled

Bioinformatics

Issues Advance Articles Publish ▾ Purchase Alerts About ▾

No cover image available

Volume 3, Issue 4  
November 1987

An access interface for the MS-DOS diskette format of GenBank(R), a gene sequence database

Michael J Weise

Bioinformatics (1987) 3 (4): 313-317. DOI: <https://doi.org/10.1093/bioinformatics/3.4.313>

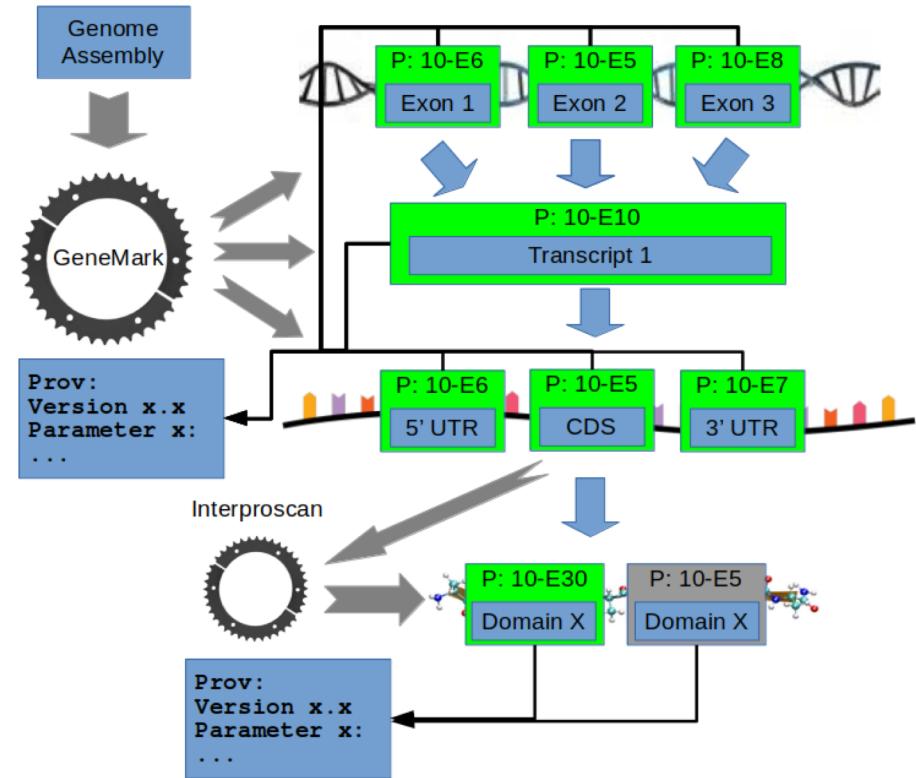
Published: 01 November 1987 Article history ▾

# SAPP: Annotation information storage

- Wrapper to commonly used annotation tools (prokaryotes and eukaryotes) that generates FAIR data
- Current usage:
  - Uniform annotation of over 10 000 bacterial species.
  - Uniform annotation of salmonoids (fish)
  - Uniform annotation of plant genomes
  - Uniform annotation of environmental samples

Koehorst et al Bioinformatics 2017  
<https://gitlab.com/sapp>

Documentation:  
<https://sapp.gitlab.io>



# Modular design

## Conversion types

- EMBL / GenBank
- FASTA
- GFF
- QTL
- VCF
- ...



## Genetic elements

- Gene prediction
- tRNA/rRNA
- Crispr
- ...



## Functional annotation

- BLAST
- Enzyme predictions
- Domain annotation
- Signal peptides
- Transmembrane
- Localization
- ...



WAGENINGEN  
UNIVERSITY & RESEARCH



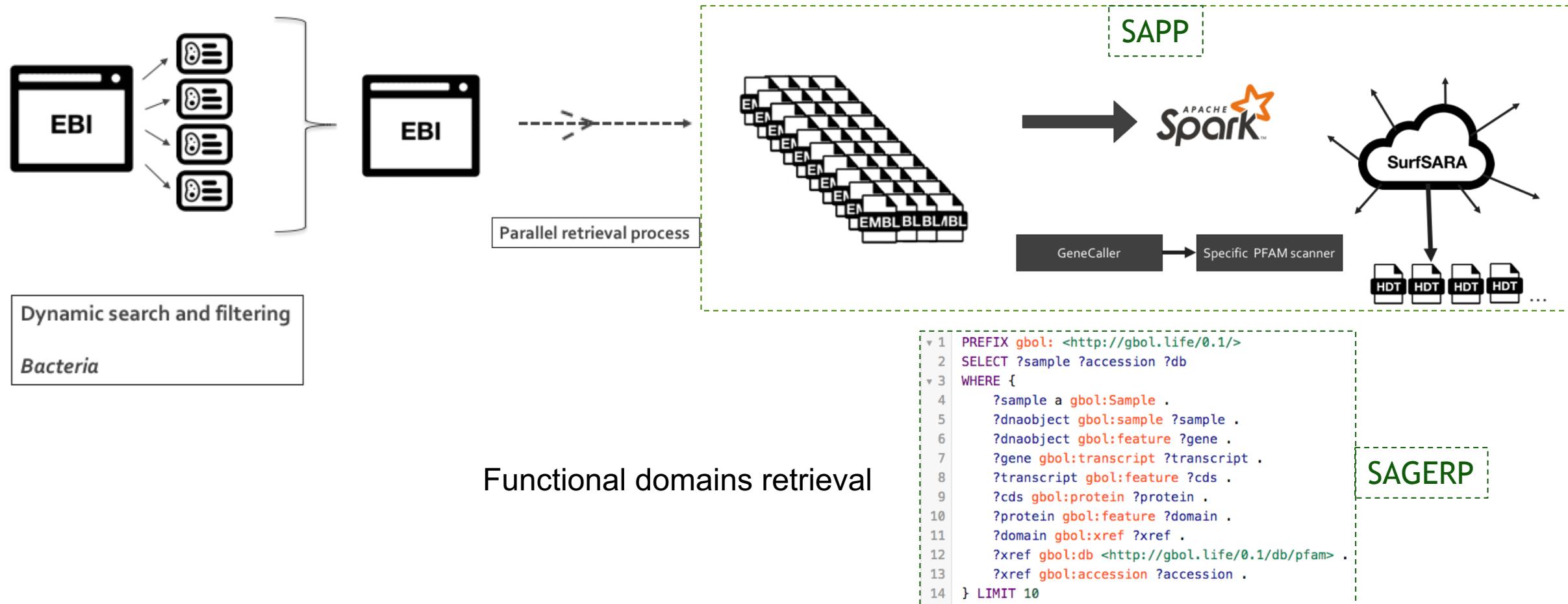
Use cases:  
Computational genomics  
&  
Environmental samples



WAGENINGEN  
UNIVERSITY & RESEARCH



# High Throughput annotation



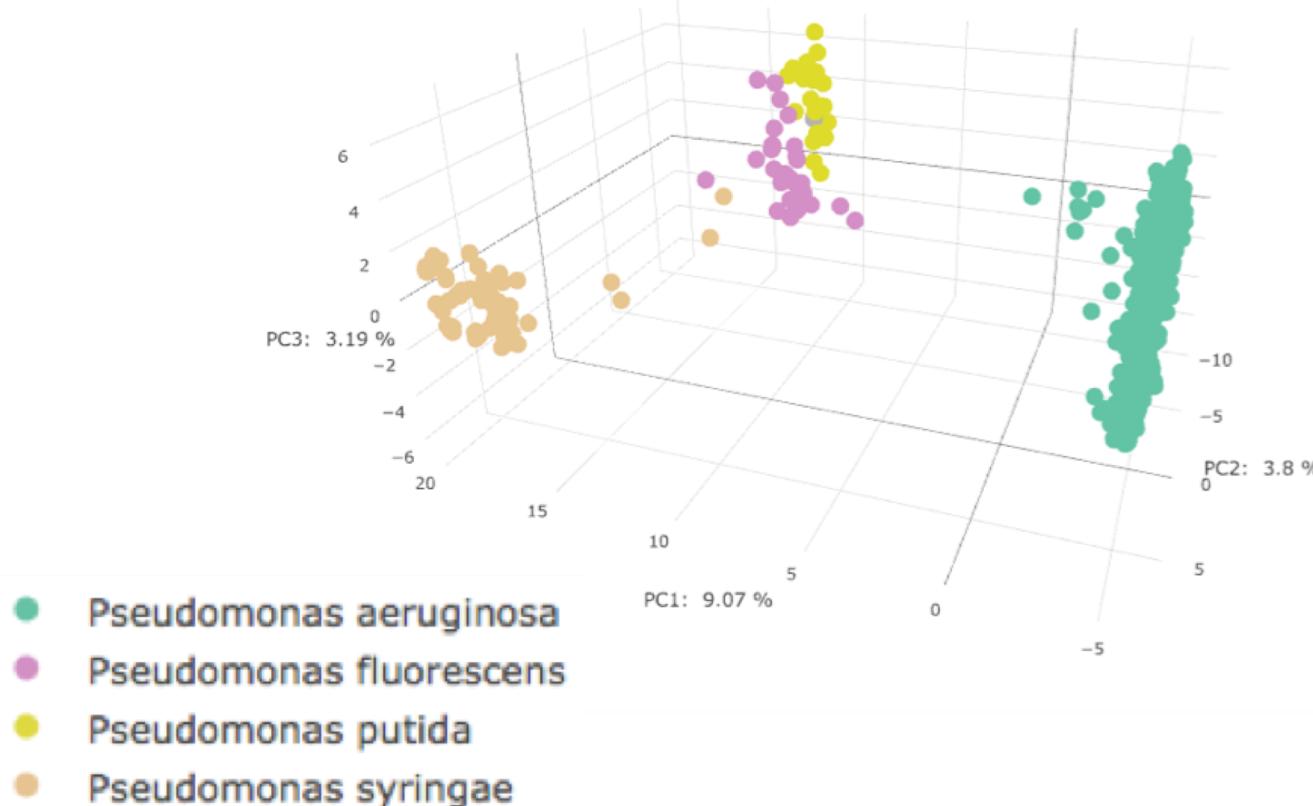
WAGENINGEN  
UNIVERSITY & RESEARCH



100 years  
1918 — 2018

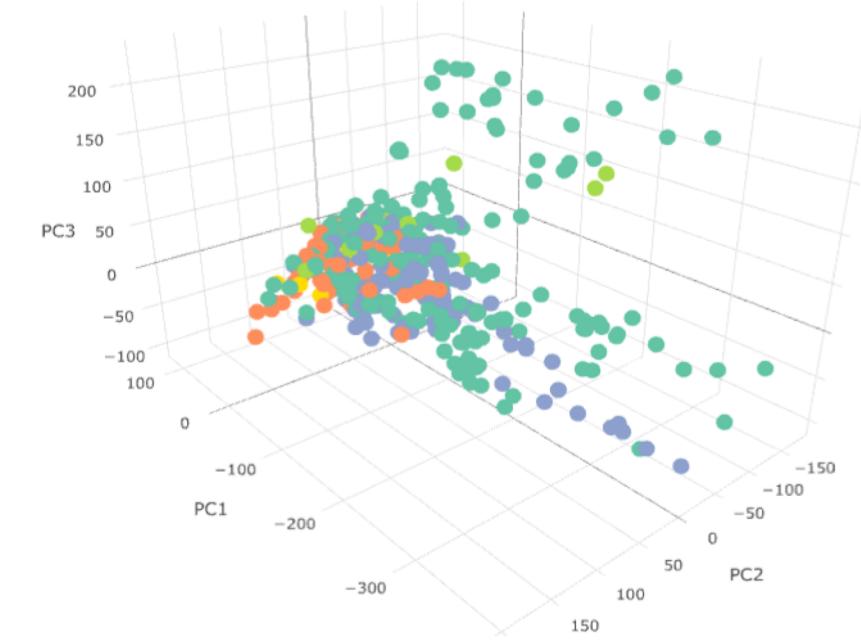
# Functional variation

Phylogeny and phenotype relationships with the functional landscape



Koehorst, Jasper J., et al. *Scientific reports* 6 (2016):

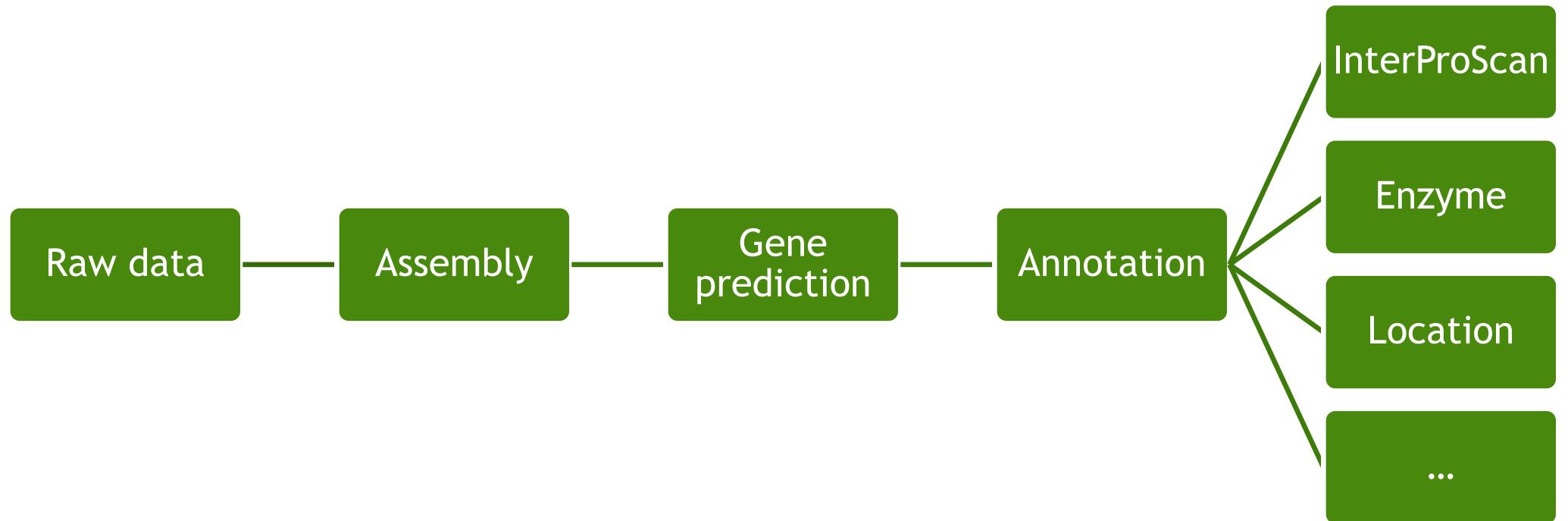
## Scored phenotypes



- Aerobe
- Anaerobe
- Facultative
- Microaerophilic
- Obligate aerobe
- Obligate anaerobe

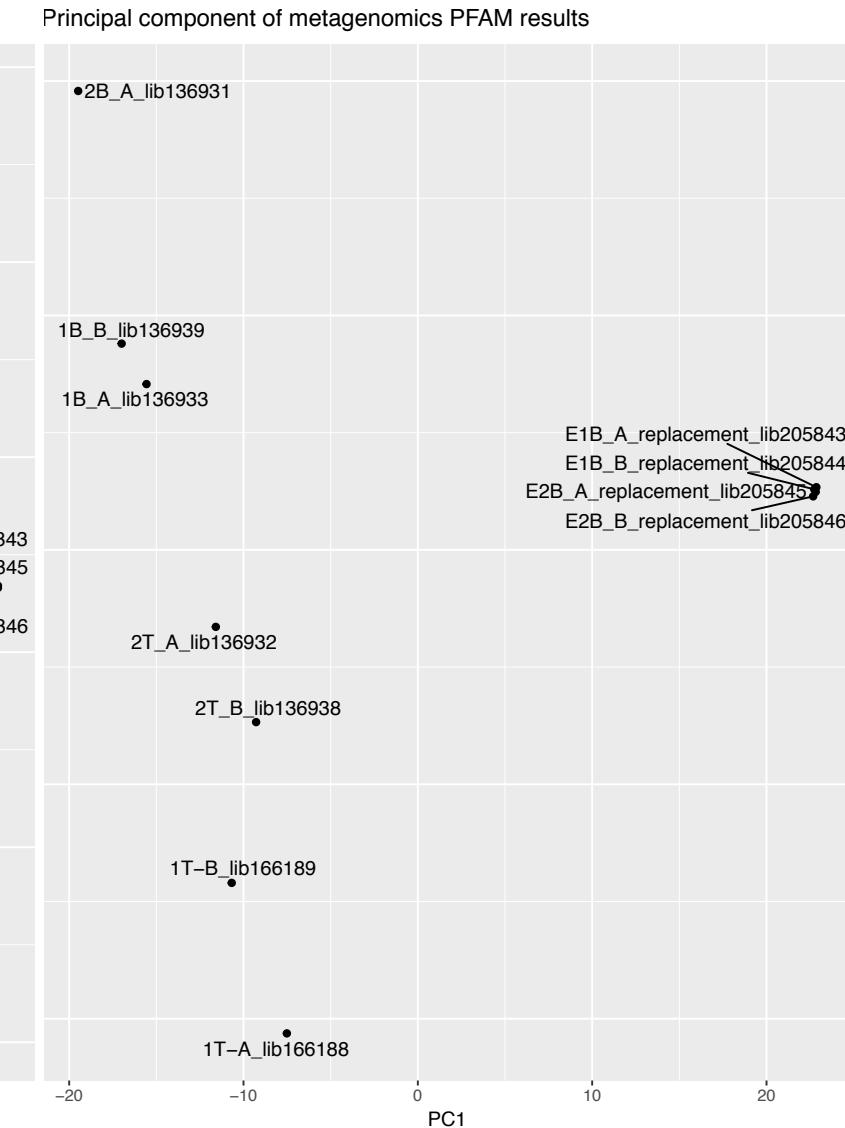
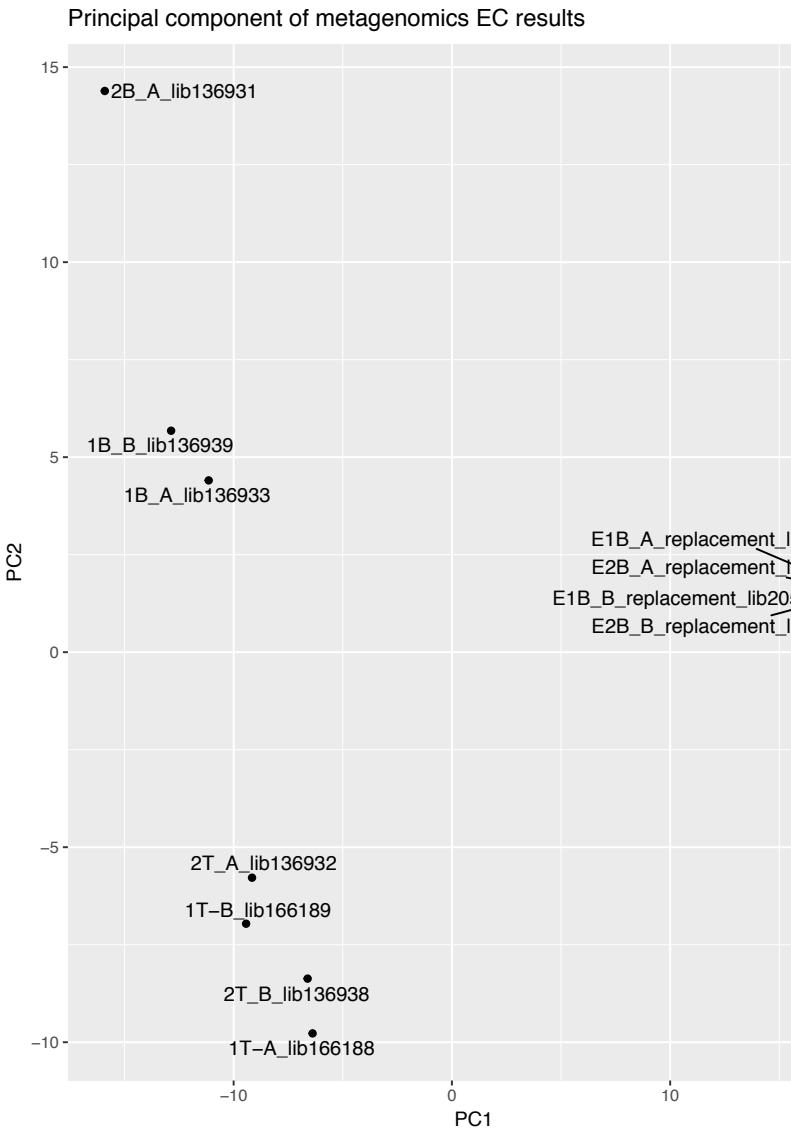


# Environmental samples (current study)

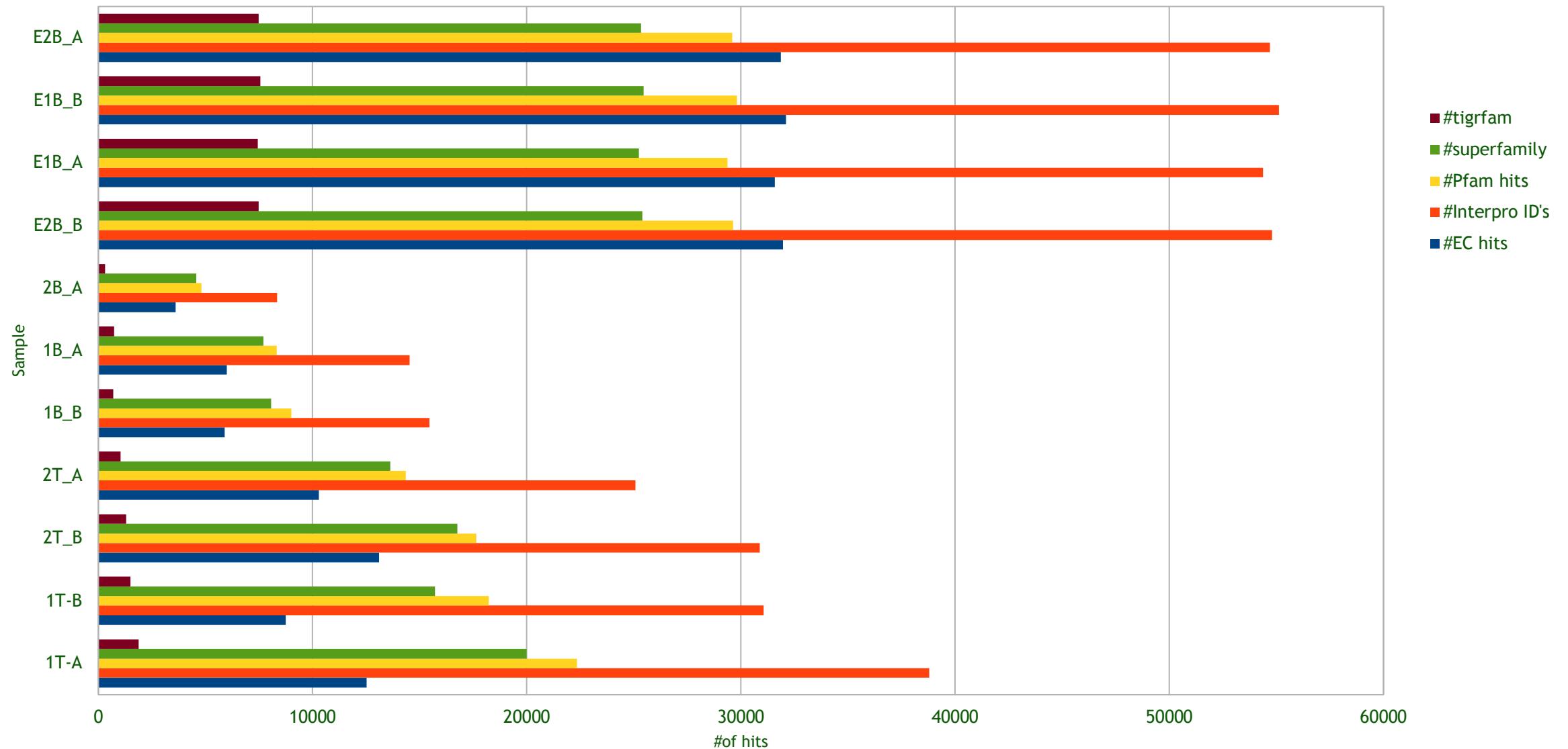


# Functional landscape

518.613 proteins in total



# Domain content per sample



# Current publications using SAPP

1. Comparative genomics highlights symbiotic capacities and high metabolic flexibility of the marine genus *Pseudovibrio*  
D Versluis, B Nijssse, MA Naim, JJ Koehorst, J Wiese, JF Imhoff, ... **Genome biology and evolution** 10 (1), 125-142
2. Concurrent haloalkanoate degradation and chlorate reduction by *Pseudomonas chloritidismutans* AW-1T  
P Peng, Y Zheng, JJ Koehorst, PJ Schaap, AJM Stams, H Smidt, ... **Applied and environmental microbiology** 83 (12), e00325-17
3. Persistence of Functional Protein Domains in Mycoplasma Species and their Role in Host Specificity and Synthetic Minimal Life  
T Kamminga, JJ Koehorst, P Vermeij, SJ Slagman, ... **Frontiers in cellular and infection microbiology** 7, 31
4. Complete Genome Sequence of *Akkermansia glycaniphila* Strain PytT, a Mucin-Degrading Specialist of the Reticulated Python Gut  
JP Ouwerkerk, JJ Koehorst, PJ Schaap, J Ritari, L Paulin, C Belzer, ... **Genome announcements** 5 (1), e01098-16
5. Interoperable genome annotation with GBOL, an extendable infrastructure for functional data mining  
JCJ van Dam, JJ Koehorst, JO Vik, PJ Schaap, M Suarez-Diez **bioRxiv**, 184747
6. Reverse methanogenesis and respiration in methanotrophic archaea  
PHA Timmers, CU Welte, JJ Koehorst, CM Plugge, MSM Jetten, ... **Archaea** 2017
7. Comparison of 432 *Pseudomonas* strains through integration of genomic, functional, metabolic and expression data  
JJ Koehorst, JCJ Van Dam, RGA Van Heck, E Saccenti, ... **Scientific reports** 6, 38699
8. Complete genome sequence of thermophilic *Bacillus smithii* type strain DSM 4216 T  
EF Bosma, JJ Koehorst, SAFT van Hijum, B Renckens, B Vriesendorp, ... **Standards in genomic sciences** 11 (1), 52
9. Protein domain architectures provide a fast, efficient and scalable alternative to sequence-based methods for comparative functional genomics  
JJ Koehorst, E Saccenti, PJ Schaap, VAPM dos Santos, M Suarez-Diez **F1000Research**
10. Assessing the metabolic diversity of streptococcus from a protein domain point of view  
E Saccenti, D Nieuwenhuijse, JJ Koehorst, VAPM dos Santos, PJ Schaap **PloS one** 10 (9), e0137908
11. A genomic view on syntrophic versus non-syntrophic lifestyle in anaerobic fatty acid degrading communities  
P Worm, JJ Koehorst, M Visser, VT Sedano-Núñez, PJ Schaap, ... **Biochimica et Biophysica Acta (BBA)-Bioenergetics** 1837 (12), 2004-2016

# Availability

- **SAPP** Koehorst et al Bioinformatics 2017  
<https://sapp.gitlab.io>
- **Empusa:** <https://gitlab.com/Empusa>
- **GBOL:** Documentation & namespace:  
<http://gbol.life/0.1/>

**SAPP: functional genome annotation and analysis through a semantic framework using FAIR principles** 

Jasper J Koehorst , Jesse C J van Dam, Edoardo Saccenti, Vitor A P Martins dos Santos, Maria Suarez-Diez, Peter J Schaap 

*Bioinformatics*, Volume 34, Issue 8, 15 April 2018, Pages 1401–1403,  
<https://doi.org/10.1093/bioinformatics/btx767>

Published: 23 November 2017 Article history ▾



HOME

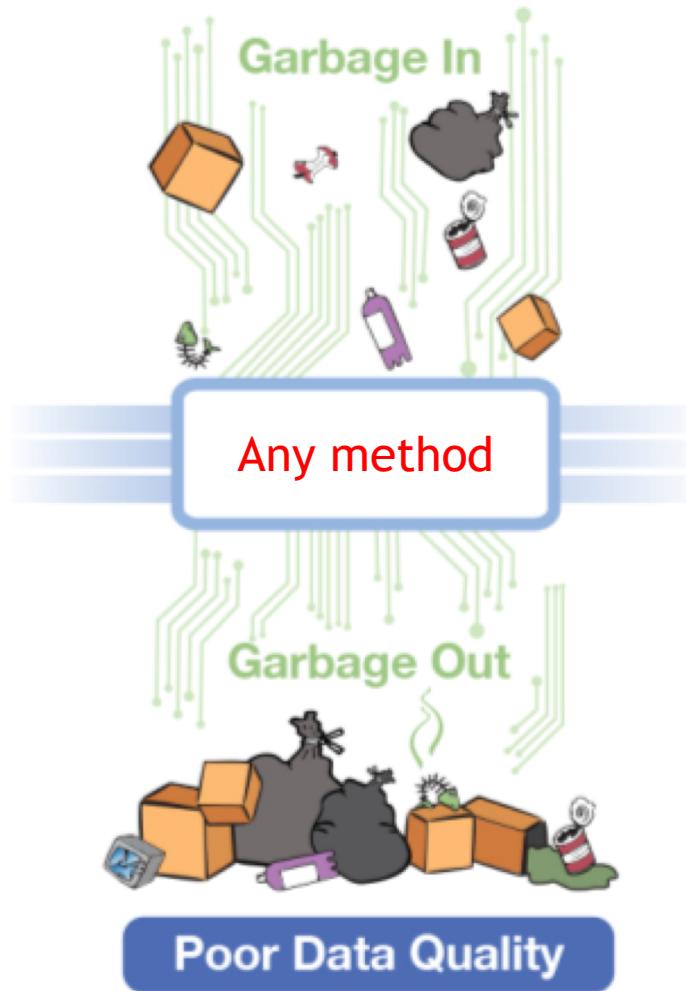
Search

New Results

**Interoperable genome annotation with GBOL, an extendable infrastructure for functional data mining**

 Jesse C.J. van Dam,  Jasper Jan J. Koehorst,  Jon Olav Vik,  Peter J. Schaap,  Maria Suarez-Diez  
**doi:** <https://doi.org/10.1101/184747>

## Poor Data Quality



WAGENINGEN  
UNIVERSITY & RESEARCH

