

Bioinformatics for metagenomics analysis



John Williams
Parco Tecnologico Padano



Human MicroBiome project

- ~ Sampling 5 sites
- ~ Reference Microbial genomes
- ~ Resource Repositories

- ~ New Computational technologies
- ~ New Computational tools
- ~ Data Analysis center
- ~ Evaluation of data vs disease

- ~ Legal Aspects





MetaGenome analysis within Ruminomics

↪ Samples

- ↪ Extreme 50 samples among 1000 cows from Italy, Sweden and UK.
- ↪ from rumen swap between reindeer and cow

↪ Explore differences in

- ↪ microbial composition
- ↪ gene composition

↪ Impacts of a change in the microflora on

- ↪ biological pathways
- ↪ degradation of specific nutrient sources





Meta-biome challenge

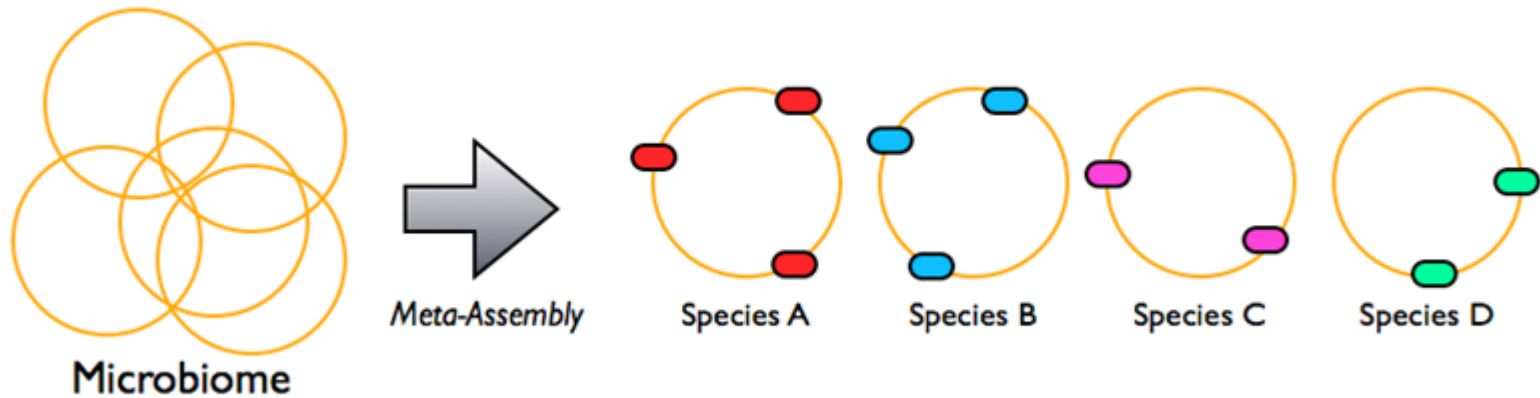
Entrepreneurial research in ag-biotech

~ Resolving micro-organisms present by

~ content

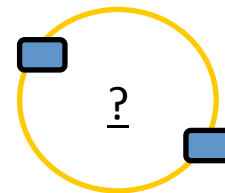
~ proportions

~ Identification of unknowns



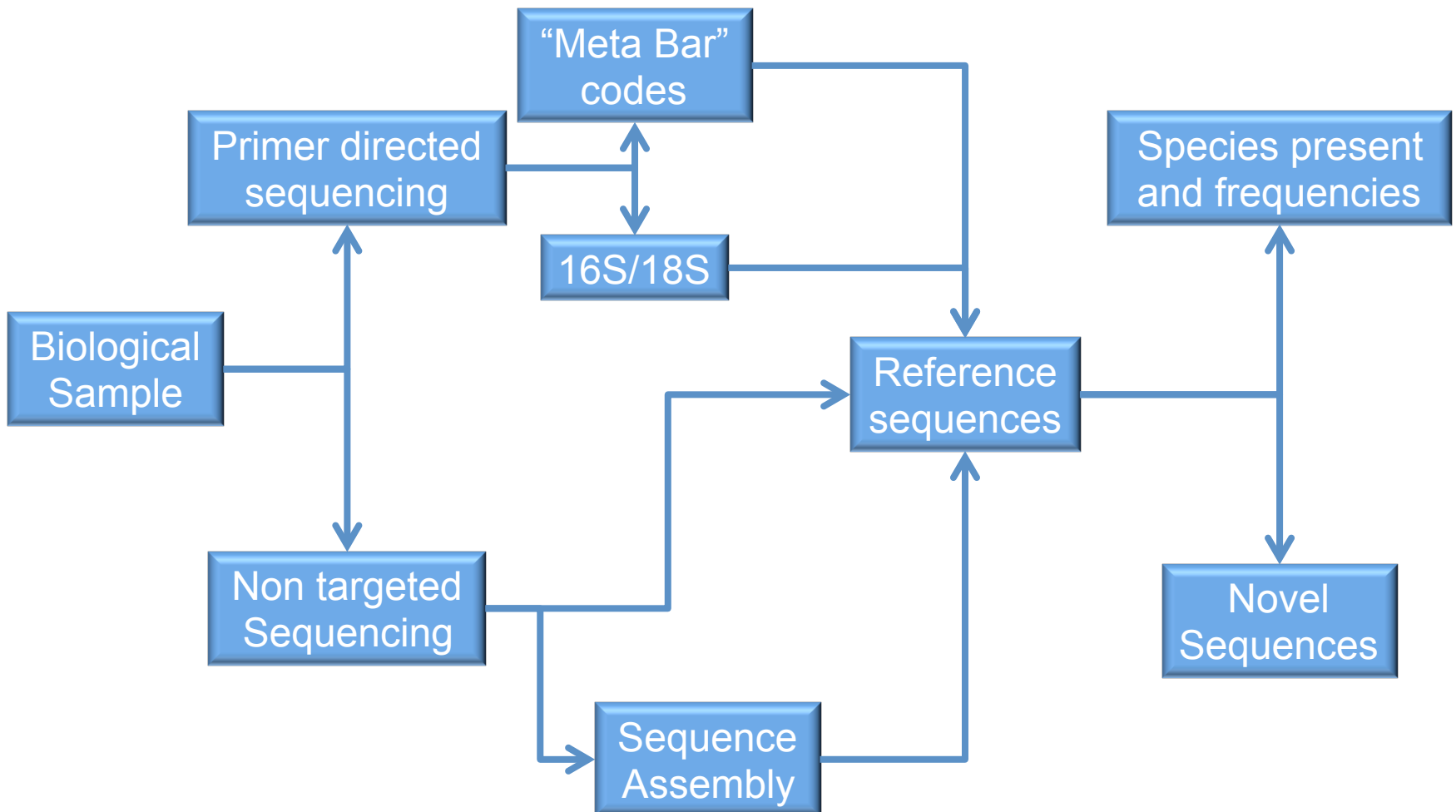
~ Low cost

~ Repeatable





Meta-biome analyses strategies





Considerations

Entrepreneurial research in ag-biotech

↪ Culture based methods

↪ May miss a large proportion of the species...

↪ Primer based methods

↪ Bias as clades missed

↪ Preferential amplification

↪ Sequence based methods

↪ Large data sets

↪ Complex analysis





Advantages of metagenomics

Entrepreneurial research in ag-biotech

- ~ Analysis of the whole consortium of microorganisms
 - ~ Sequencing without targeted amplification
 - ~ Determination of microbial composition and abundances
 - ~ With or without a reference set of genomes

- ~ Expression capability
 - ~ Meta-transcriptome
 - ~ Meta-proteome
 - ~ Impact of environment on the transcription

- ~ sequence to gene and function





Sequence information

Entrepreneurial research in ag-biotech

↪ Classical sequencing (Sanger)

- ↪ *complete genomes*
- ↪ annotation to locate genes, operons etc
- ↪ contained data sets per species

↪ Metagenome sequencing

- ↪ fragmented and incomplete sequences
- ↪ reference required to identify individual or species
- ↪ need to reconstruct genomes
- ↪ large data sets
- ↪ **computational challenges**





Challenges

Entrepreneurial research in ag-biotech

↪ The rate of of genome sequencing has exceeded the bioinformatic capability to analyse the data

↪ This is now the bottleneck

↪ Challenges

↪ More complex than the assembly of single genomes

↪ Variable representation of different species

↪ Genetic diversity within strains

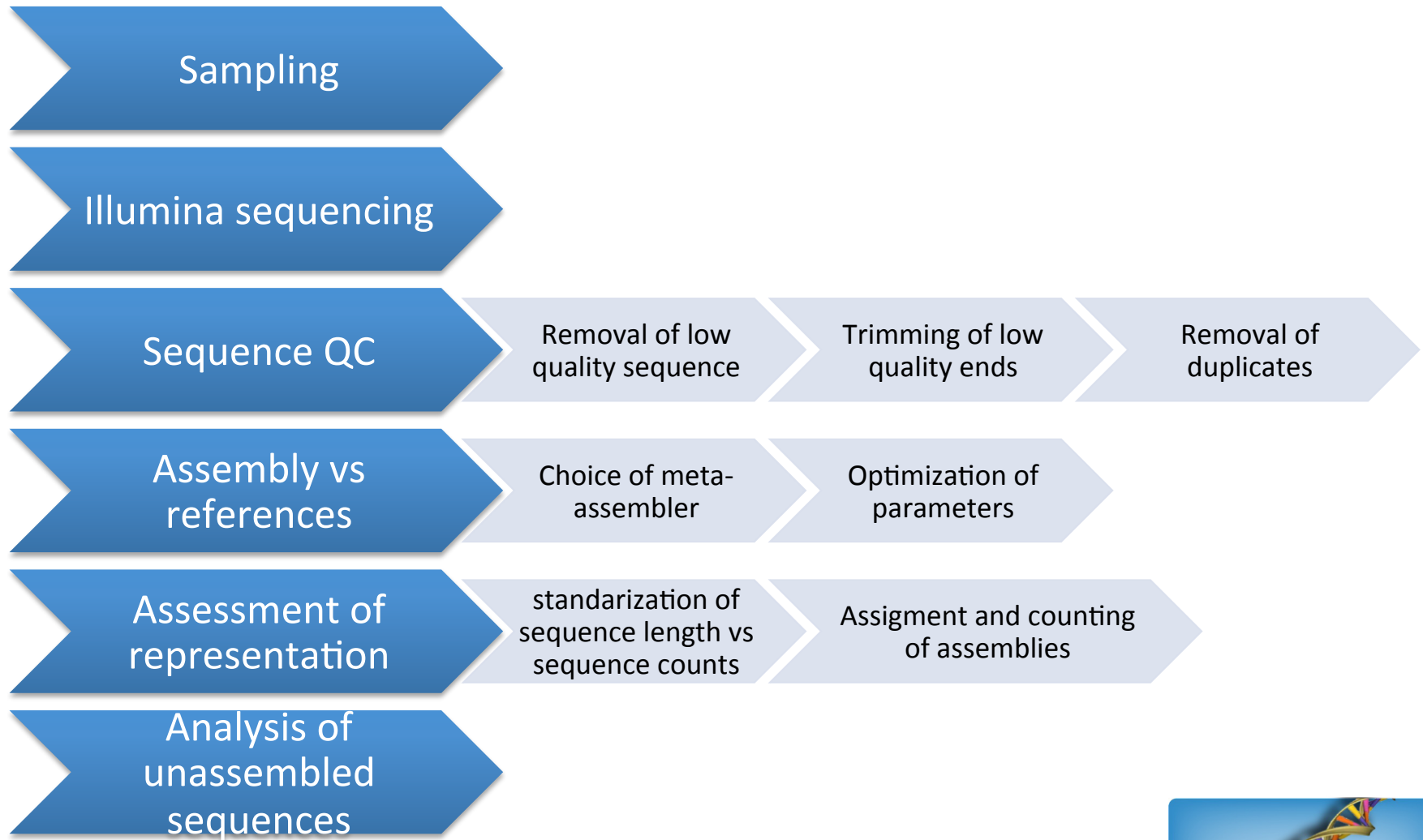
↪ Homology among strains





PTP meta-genome pipeline

Entrepreneurial research in ag-biotech



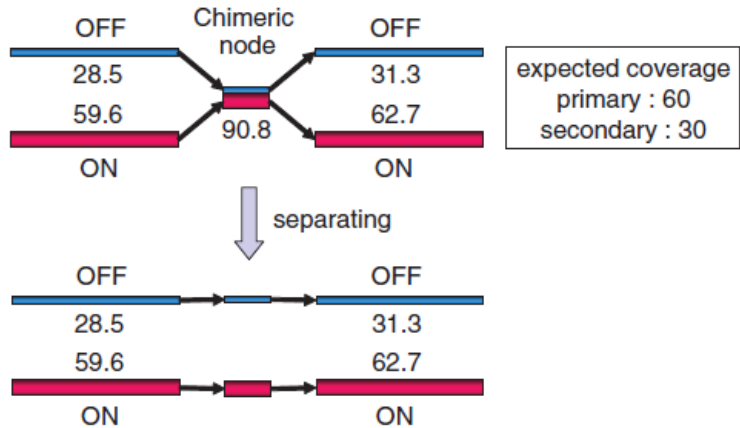


Resolving the graphs, multiple species

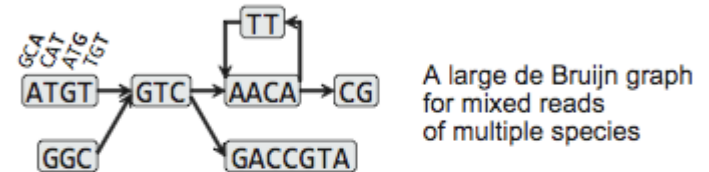
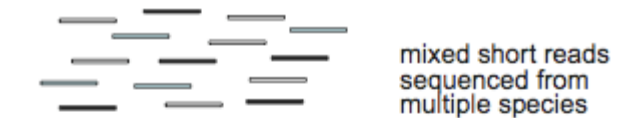
Entrepreneurial research in ag-biotech

MetaVelvet

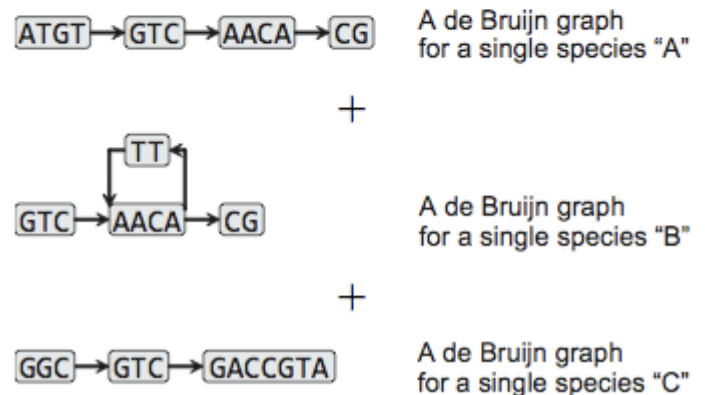
- Velvet (single genomes)
- To meta-genome analysis
- Resolved Debruijn graphs
- But on a single computer



- Use of paired end data



decomposing





MetaRay vs MetaVelvet

Entrepreneurial research in ag-biotech

~ Solutions to computing complexity

~ parallel array computing. eg ABySS ,RAY... for single genomes

~ MetaRay

~ distributed array computing

~ increased data handling time

~ but decreased real time processing

~ uses reference genomes to assign K-mers

~ Better results on the longer contigs

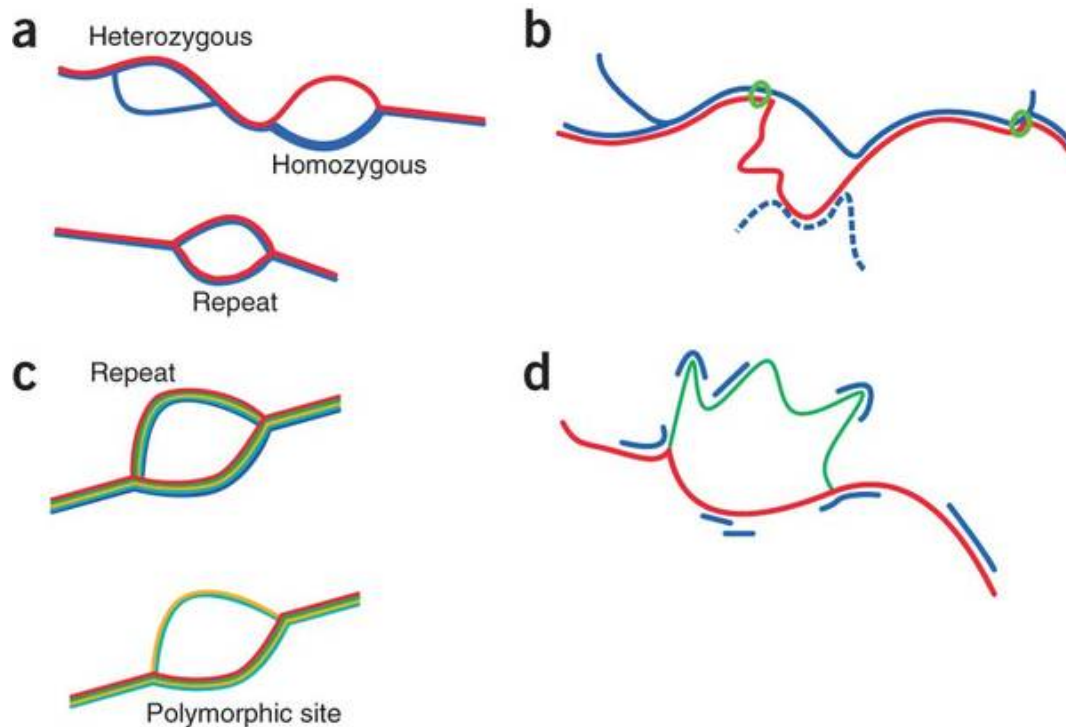
~ With contigs greater than 500bp



MetaRay

Entrepreneurial research in ag-biotech

- ↗ Construction of a de Bruijn graph from the raw sequence
- ↗ de novo assembly
- ↗ colouring with finished bacterial genomes





Simulated performance

Entrepreneurial research in ag-biotech

~ Data set

- ~ 1000 bacterial genomes
- ~ 1% human
- ~ 3×10^9 (Illumina HighSeq Flow Cell)

~ Errors

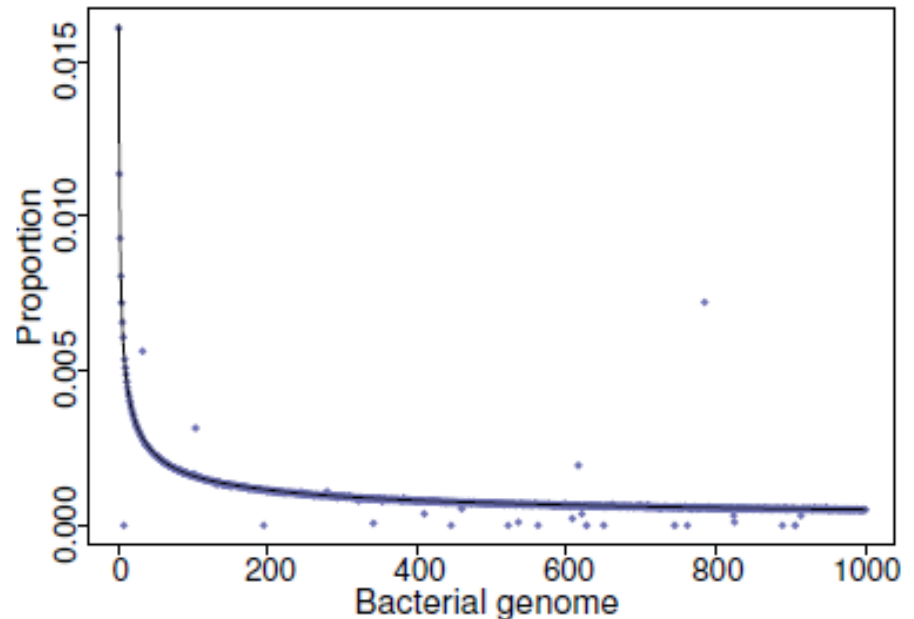
- ~ 2.6K miss-assemblies (1.3%)
- ~ 4 genomes overestimated
- ~ 20 underestimated

~ Analysis

- ~ 15 hours
- ~ 1,024 cores 1.5GB/core

~ Output

- ~ 974K contigs,
- ~ N50 76K
- ~ 974 and 2,894K per genome





Ruminomics data PTP GenHome

↪ Data

↪ 1 sample

↪ HighSeq 32Gb

	MetaRay (k = 31)	MetaRay (k = 63)	MetaVelvet (k = 51)	MetaVelvet (k = 80)
# of contigs	3666701	113241	8871806	594313
N50	228	1069	382	839
Largest contig	190086	57744	51187	43580

↪ Analysis

↪ MetaRay

↪ 512 Gb 48 cores

↪ memory usage 200Gb

↪ 12 hours

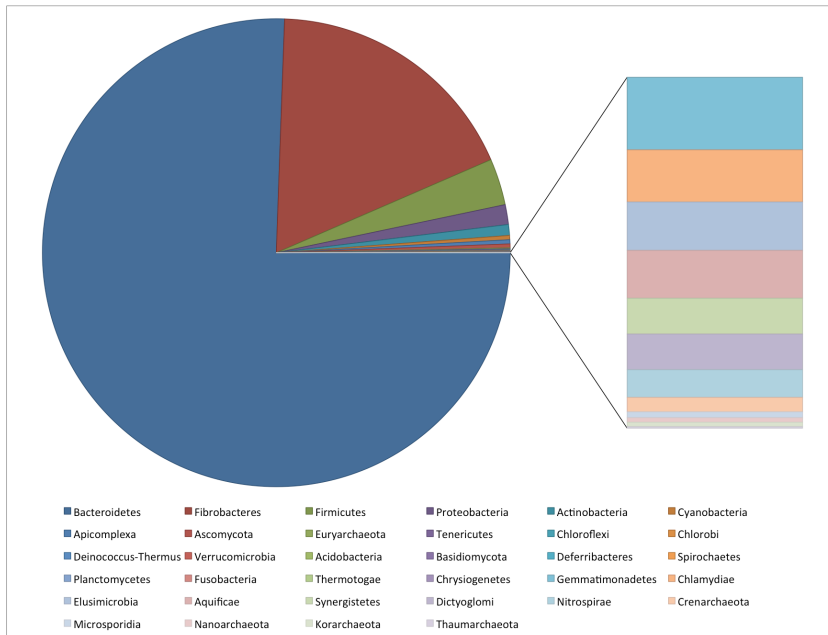




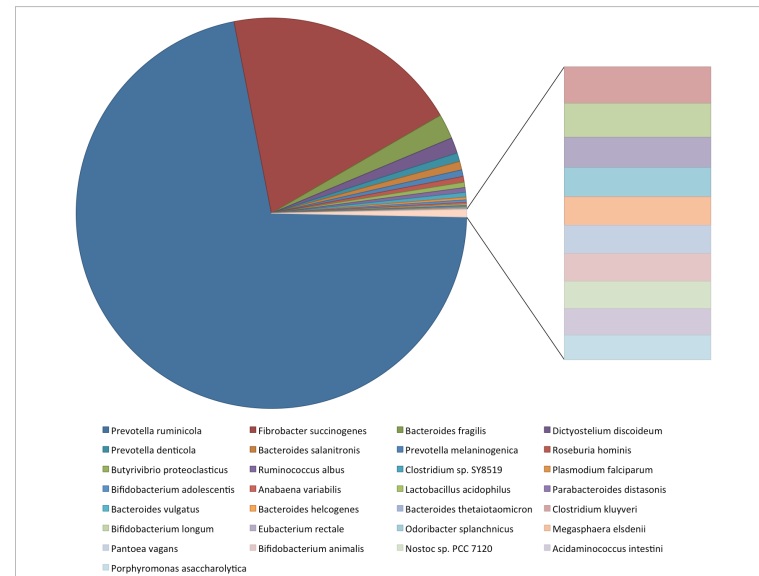
Taxonomic results

Entrepreneurial research in ag-biotech

↗ Phylum



↗ Species



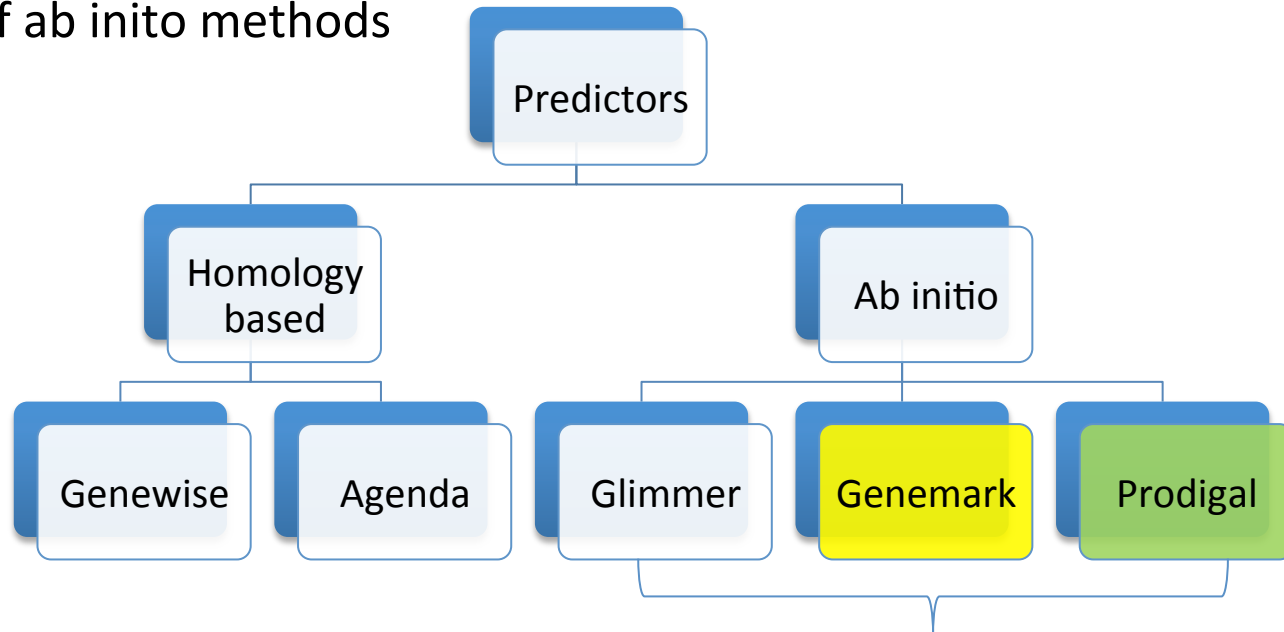


Gene prediction

Entrepreneurial research in ag-biotech

When assemble meta-genomics data will have :

- high fraction of unknown sequences
- Use of ab inito methods



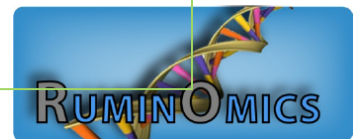
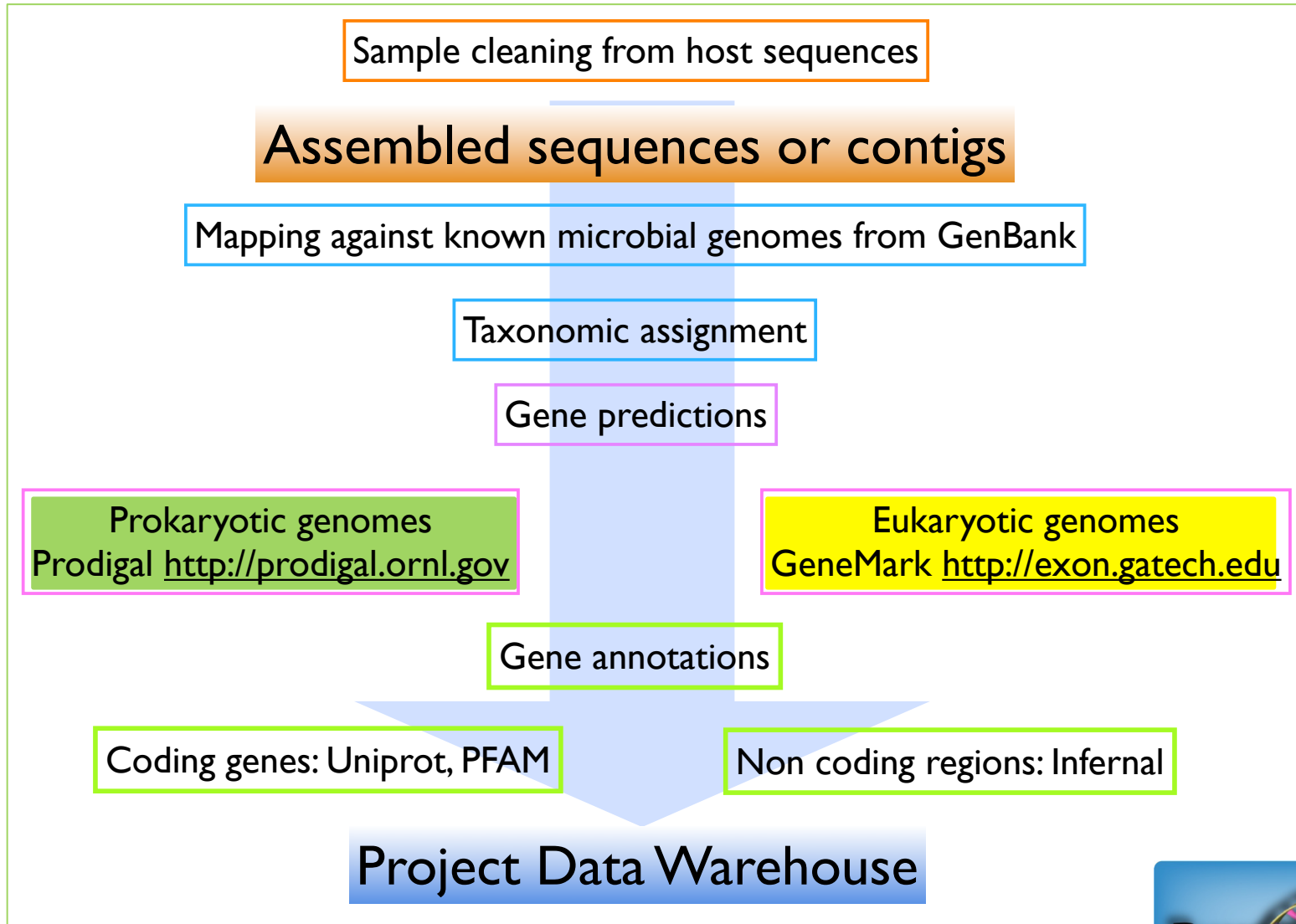
Tools within the NCBI annotation pipeline





From sequence to genes

Entrepreneurial research in ag-biotech



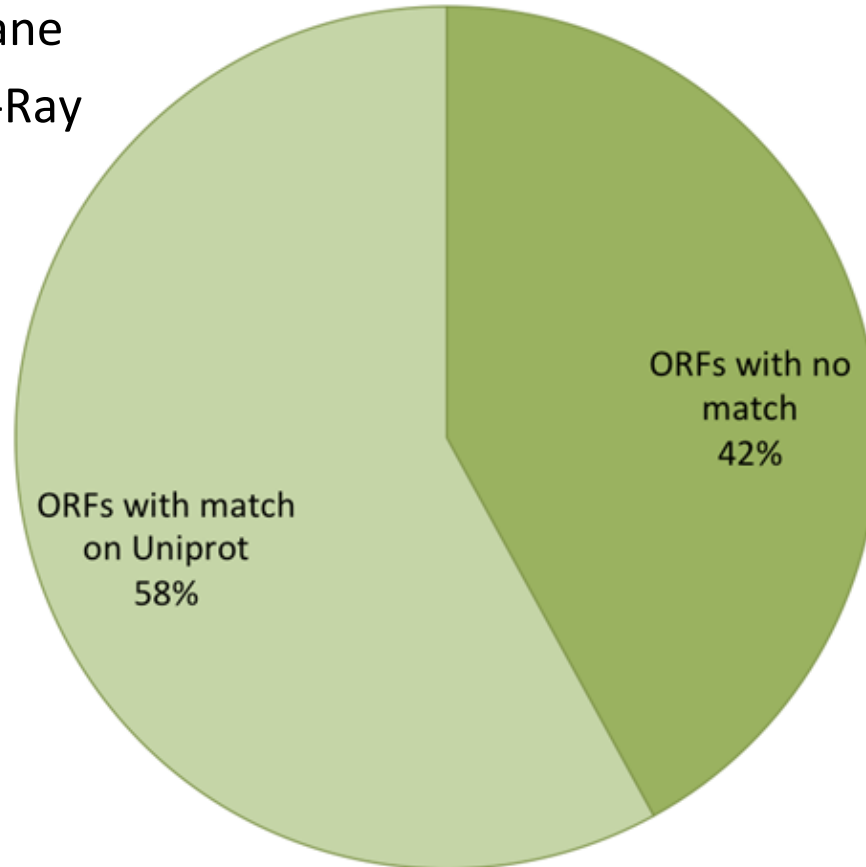


First Ruminomics Meta-data

~ Total ORFs predicted : 206,054

~ 1 sample on 1 HiSeq lane

~ assembled with Meta-Ray



~ Contigs homology

~ BlastN on Uniprot

~ E-value cut-off 1e-10



Functional analyses

Entrepreneurial research in ag-biotech

~ Gene classification by

~ structure

~ function

~ Definition of differences among microbial species

~ Coding vs non-coding genes

~ Assessment of coding capacity

~ metabolic pathways

~ newly identify strains





Statistical analyses

Entrepreneurial research in ag-biotech

- ~ Correlation between functional genomics information and phenotypic information

- ~ Model for analysis includes
 - ~ Factors for environmental effects on phenotype
 - ~ Covariances among complex phenotypes
 - ~ Presence/absence of single species /phila
 - ~ Groups defined by clustering

- ~ Use non parametric methods to solve





In summary

Entrepreneurial research in ag-biotech

- ~ A “Ruminomics” pipeline is available
 - ~ Assembly
 - ~ Composition
 - ~ Function

- ~ Approaches are computationally demanding
 - ~ Tested on a small data set
 - ~ Yet to be tested on a large data set





Acknowledgements

Entrepreneurial research in ag-biotech

↪ Bioinformatics group

↪ Francesco Strozzi

↪ Statistical Genetics

↪ Alessandra Stella

↪ GenHome Unit

↪ Chiara Ferandi





Thank you !

**Parco
Tecnologico
Padano**

Entrepreneurial research in ag-biotech





Impact of Reference Data

Entrepreneurial research in ag-biotech

