# Predicting Microbial Gene Function on a Massive Scale Reveals Extensive Complementarity between Genome Context Methods

Vedrana Vidulin[1], Tomislav Šmuc[1], Fran Supek[1;2]

[1]Division of Electronics, Ruđer Bošković Institute, Bijenička cesta 54, 10000 Zagreb, Croatia
[2]EMBL/CRG Systems Biology Unit, Centre for Genomic Regulation, Dr. Aiguader 88, 08003 Barcelona, Spain

MAESTRA
LEARNING FROM MASSIVE, INCOMPLETELY ANNOTATED, AND STRUCTURED DATA

InnoMol
Innovation Pipeline

## Abstract

We present a novel pipeline for annotating prokaryotic genes with Gene Ontology functions based on supervised machine learning in a hierarchical multi-label setting. The pipeline includes four genome context methods, which are complementary to a large extent. The four methods are combined and information accretion measured to compare the amount of past vs. newly predicted knowledge on gene function. Results indicate that a comprehensive use of genome context methods allows a sizable increase in knowledge regarding gene function.
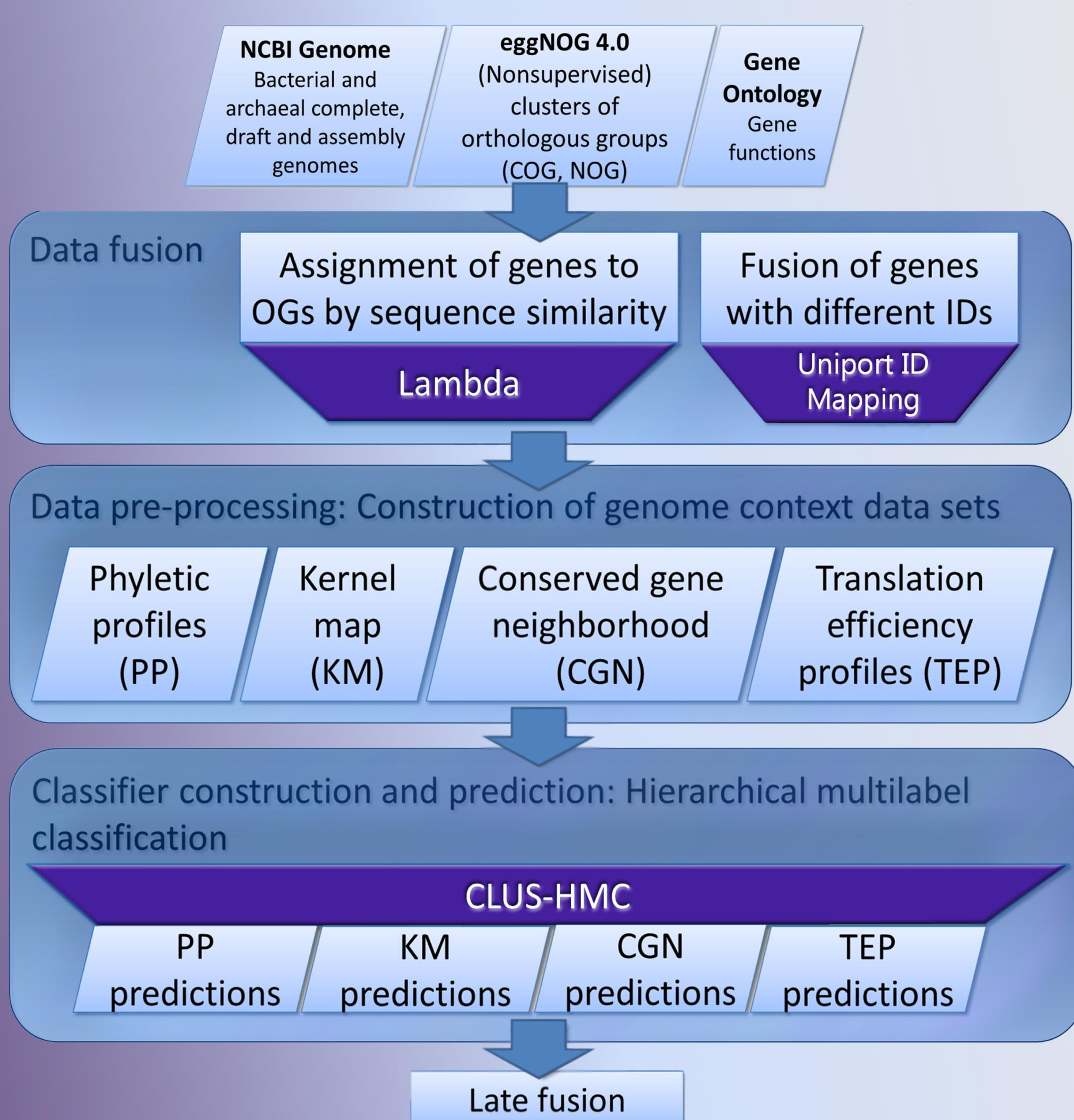
## Methods and Data

**NCBI Genome** — Bacterial and archaeal complete, draft and assembly genomes

**eggNOG 4.0** (Nonsupervised) clusters of orthologous groups (COG, NOG)

**Gene Ontology** Gene functions

Data fusion — Assignment of genes to OGs by sequence similarity (Lambda) | Fusion of genes with different IDs (Uniport ID Mapping)

Data pre-processing: Construction of genome context data sets
- Phyletic profiles (PP)
- Kernel map (KM)
- Conserved gene neighborhood (CGN)
- Translation efficiency profiles (TEP)

Classifier construction and prediction: Hierarchical multilabel classification

**CLUS-HMC**
- PP predictions
- KM predictions
- CGN predictions
- TEP predictions

Late fusion

**Figure 1.** *Gene function prediction pipeline.* Accuracy of classifiers is measured on cross-validation (out-of-bag error in CLUS-HMC Random Forest) and their predictions are combined in a late fusion scheme.
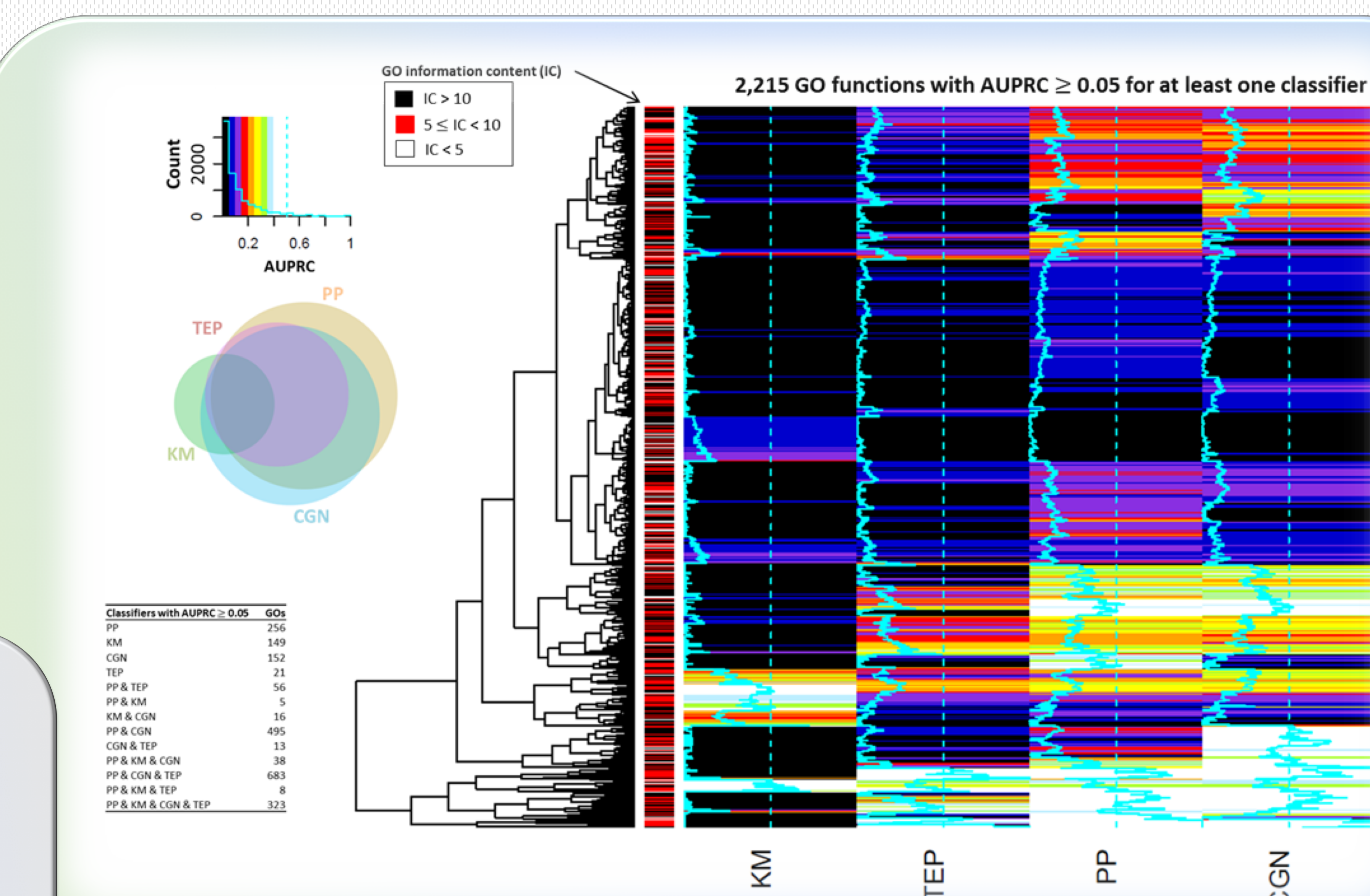
## Results

GO information content (IC)
- IC > 10
- 5 ≤ IC < 10
- IC < 5

2,215 GO functions with AUPRC ≥ 0.05 for at least one classifier

| Classifiers with AUPRC ≥ 0.05 | GOs |
|---|---|
| PP | 256 |
| KM | 149 |
| CGN | 152 |
| TEP | 21 |
| PP & TEP | 56 |
| PP & KM | 5 |
| KM & CGN | 16 |
| PP & CGN | 495 |
| CGN & TEP | 13 |
| PP & KM & CGN | 38 |
| PP & CGN & TEP | 683 |
| PP & KM & TEP | 8 |
| PP & KM & CGN & TEP | 323 |

KM | TEP | PP | CGN

**Figure 3.** *Complementarity between genome context methods.* From 2,215 functions that were successfully learned (cross-validation AUPRC>0.05) 578 were learned exclusively by one method.
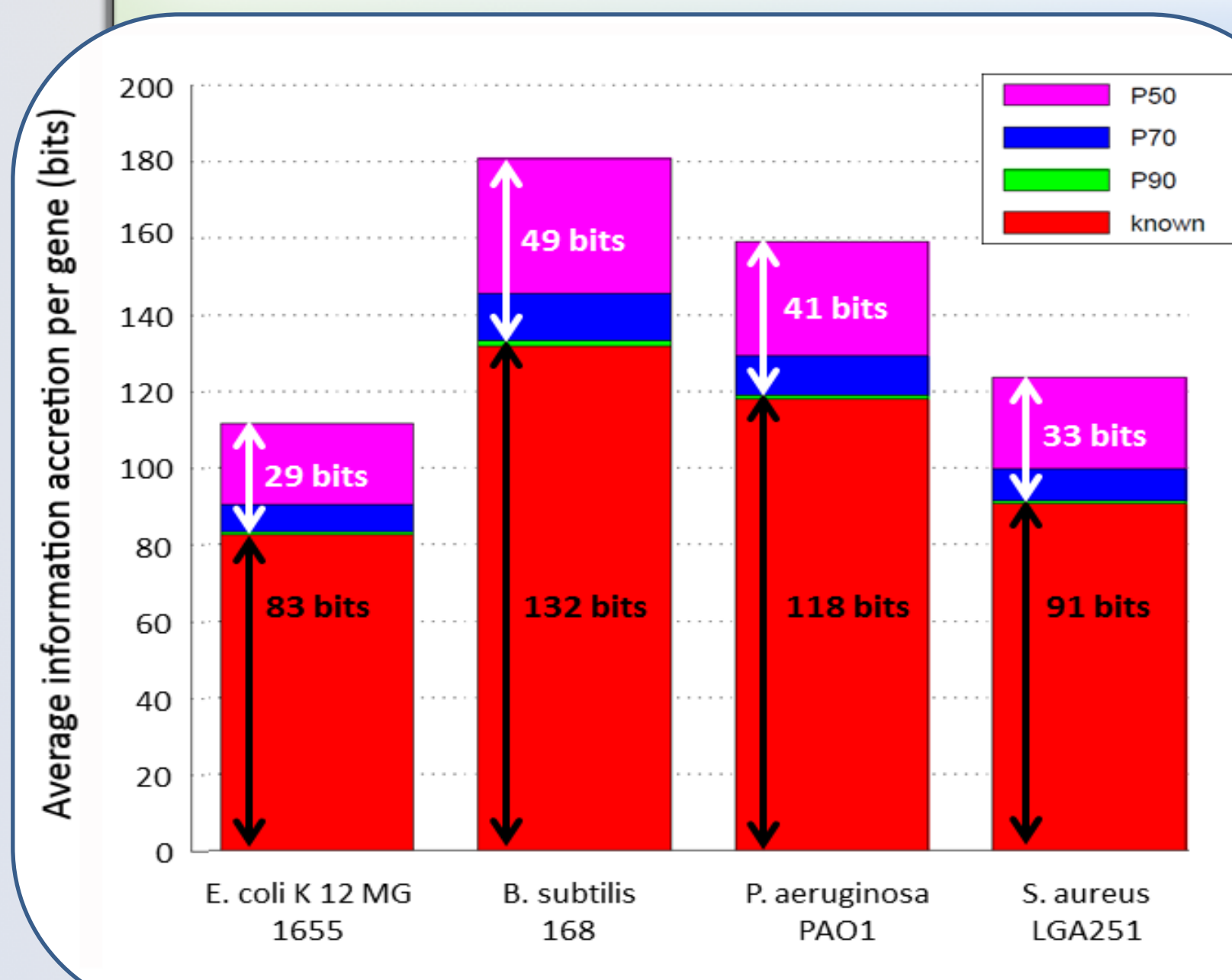
**Figure 4.** *Information accretion indicating the amounts of past vs. newly predicted knowledge on gene functions.*
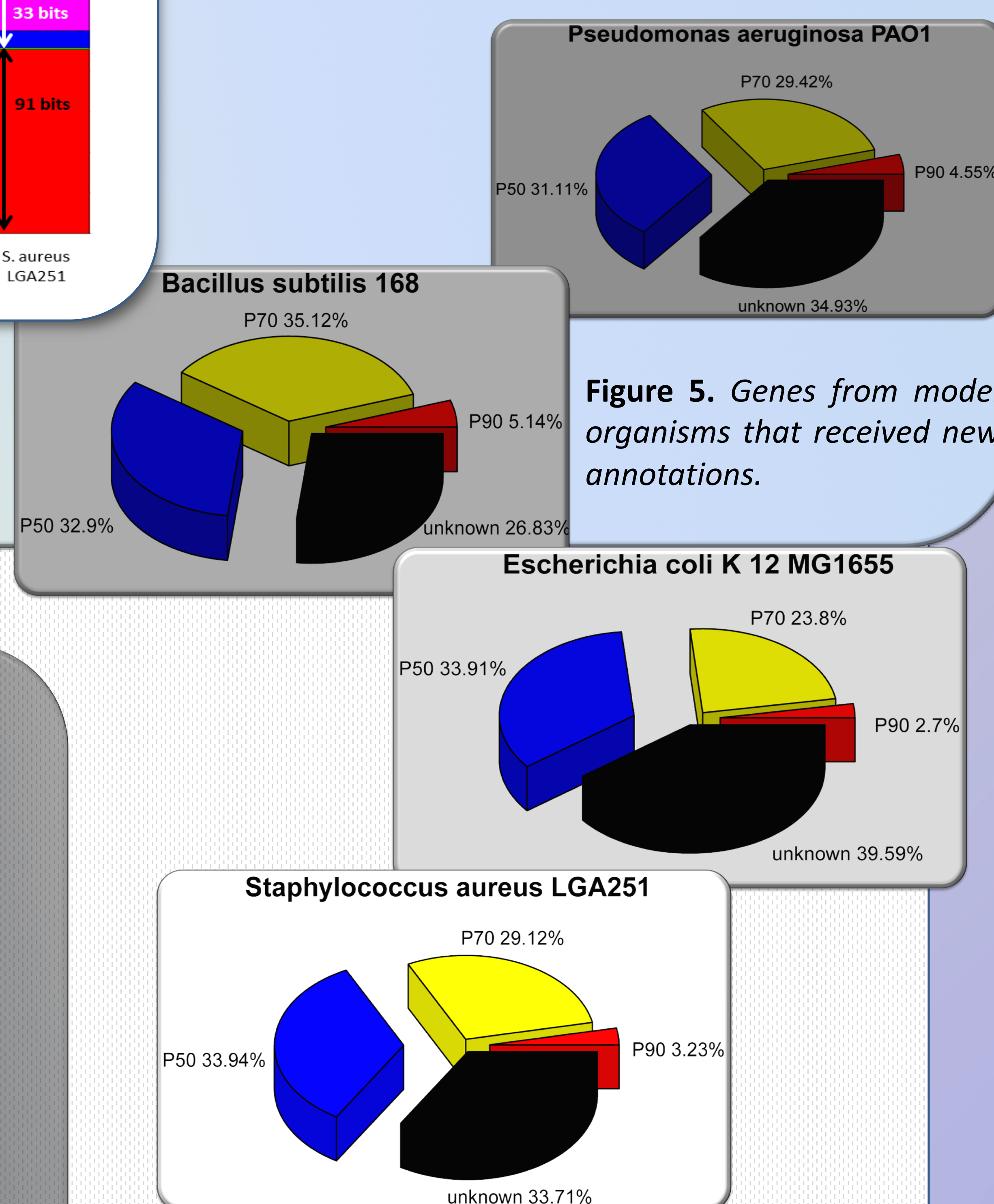
- P50
- P70
- P90
- known

Average information accretion per gene (bits)

| | E. coli K 12 MG 1655 | B. subtilis 168 | P. aeruginosa PAO1 | S. aureus LGA251 |
|---|---|---|---|---|
| known | 83 bits | 132 bits | 118 bits | 91 bits |
| added | 29 bits | 49 bits | 41 bits | 33 bits |

**Figure 5.** *Genes from model organisms that received new annotations.*

**Pseudomonas aeruginosa PAO1**
- P70 29.42%
- P90 4.55%
- unknown 34.93%
- P50 31.11%

**Bacillus subtilis 168**
- P70 35.12%
- P90 5.14%
- unknown 26.83%
- P50 32.9%

**Escherichia coli K 12 MG1655**
- P70 23.8%
- P90 2.7%
- unknown 39.59%
- P50 33.91%

**Staphylococcus aureus LGA251**
- P70 29.12%
- P90 3.23%
- unknown 33.71%
- P50 33.94%

## Conclusions

Evaluation of proposed pipeline's performance shows a sizable increase in knowledge regarding gene function. 60.41% E. coli OGs received at least one novel and likely correct (precision >50%) function. Model bacteria have 83–132 bits/gene of known annotations, while our pipeline annotates on average 38 additional bits per gene.

## Acknowledgement

**Phyletic Profiles (PP)** — Attributes = 2 070 microbial genomes

**Class: 7 956 Gene Ontology (GO) functions**

**Conserved Gene Neighborhoods (CGN)** — Attributes = 5 889 OGs that appear in at least 100 genomes

Examples = 120 199 eggNOG 4.0 OGs

| | g₁ | g₂ | g₃ | g₄ | GO |
|---|---|---|---|---|---|
| OG₁ | 1 | 0 | 0 | 1 | |
| OG₂ | 1 | 1 | 0 | 1 | ? |
| OG₃ | 0 | 1 | 0 | 1 | |
| OG₄ | 1 | 0 | 1 | 1 | ? |

| | OG₁ | OG₂ | OG₃ | OG₄ | GO |
|---|---|---|---|---|---|
| OG₁ | 0 | 19.42 | 13.88 | 7.21 | |
| OG₂ | 19.42 | 0 | 23.81 | 23.81 | ? |
| OG₃ | 13.88 | 23.81 | 0 | 20.38 | |
| OG₄ | 7.21 | 23.81 | 20.38 | 0 | ? |

| | g₁ | g₂ | g₃ | g₄ | GO |
|---|---|---|---|---|---|
| OG₁ | 0.71 | 0.53 | 0.11 | 0.71 | |
| OG₂ | 0.48 | 0.25 | 0.52 | 0.38 | ? |
| OG₃ | 1.22 | 0.56 | 0.27 | 0.44 | |
| OG₄ | 0.66 | 0.56 | 0.34 | 0.59 | ? |

| | OG₁ | OG₂ | OG₃ | OG₄ | GO |
|---|---|---|---|---|---|
| OG₁ | 0 | 0.24 | 6.64 | 6.64 | |
| OG₂ | 0.24 | 0 | -9.87 | 1.32 | ? |
| OG₃ | 6.64 | -9.87 | 0 | 6.64 | |
| OG₄ | 6.64 | 1.32 | 6.64 | 0 | ? |

**Translation Efficiency Profiles (TEP)**

**Kernel Map (KM)** — Attributes = 9 105 from 7 genomes

**Figure 2.** *Genome context representations.* Attribute values represent presence of genes in genomes in PP, predicted gene expression levels in TEP, logarithm of average chromosomal pairwise distances between genes across all genomes in CGN, and logarithm of e-values in KM. 43% of OGs are used for training classifiers.