

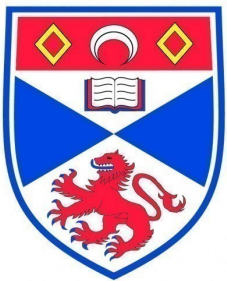
# Active and Guided Learning of Enzyme Function

Luna De Ferrari  
Stuart Aitken  
John Mitchell

[luna.deferrari@st-andrews.ac.uk](mailto:luna.deferrari@st-andrews.ac.uk)

*6th International Workshop  
on Machine Learning in Systems Biology  
@ECCB Basel, Switzerland*

*8 Sep 2012*



University  
of  
St Andrews



# Predict enzyme function

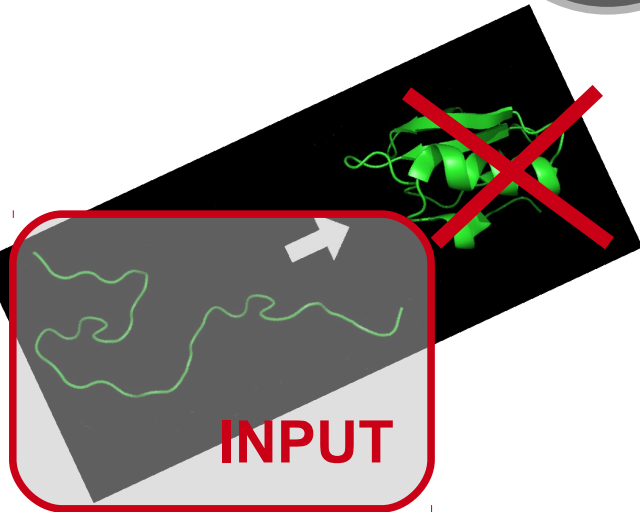
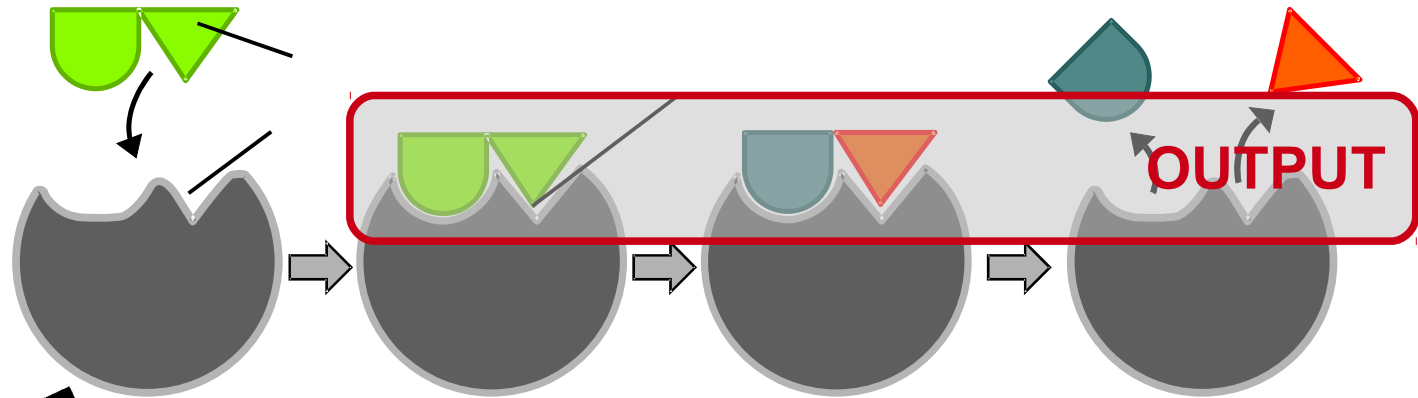


Fig: The induced fit model of enzyme activity. Provided by Tim Vickers. From: <http://en.wikipedia.org/wiki/Enzyme>

Fig: Protein folding. Provided by Dr Kjaergaard. From: [http://en.wikipedia.org/wiki/Protein\\_folding](http://en.wikipedia.org/wiki/Protein_folding)

# Enzyme Commission numbers

## Transferases

transferring nitrogenous groups

other nitrogenous groups

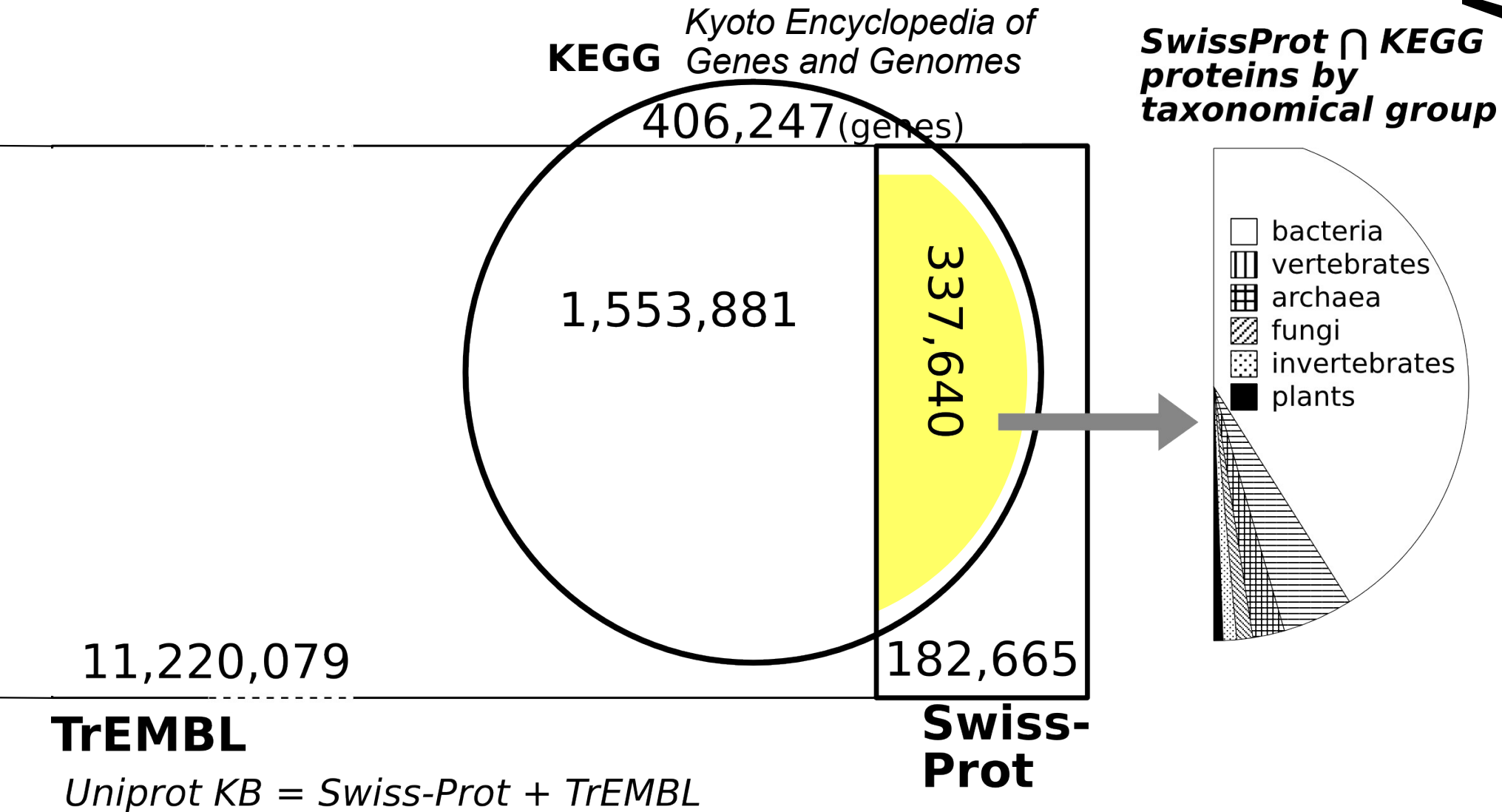


2.6.99.2

**Pyridoxine 5'-  
phosphate synthase**

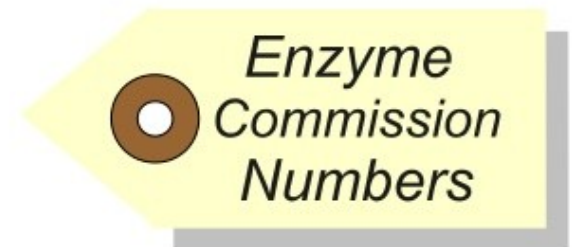
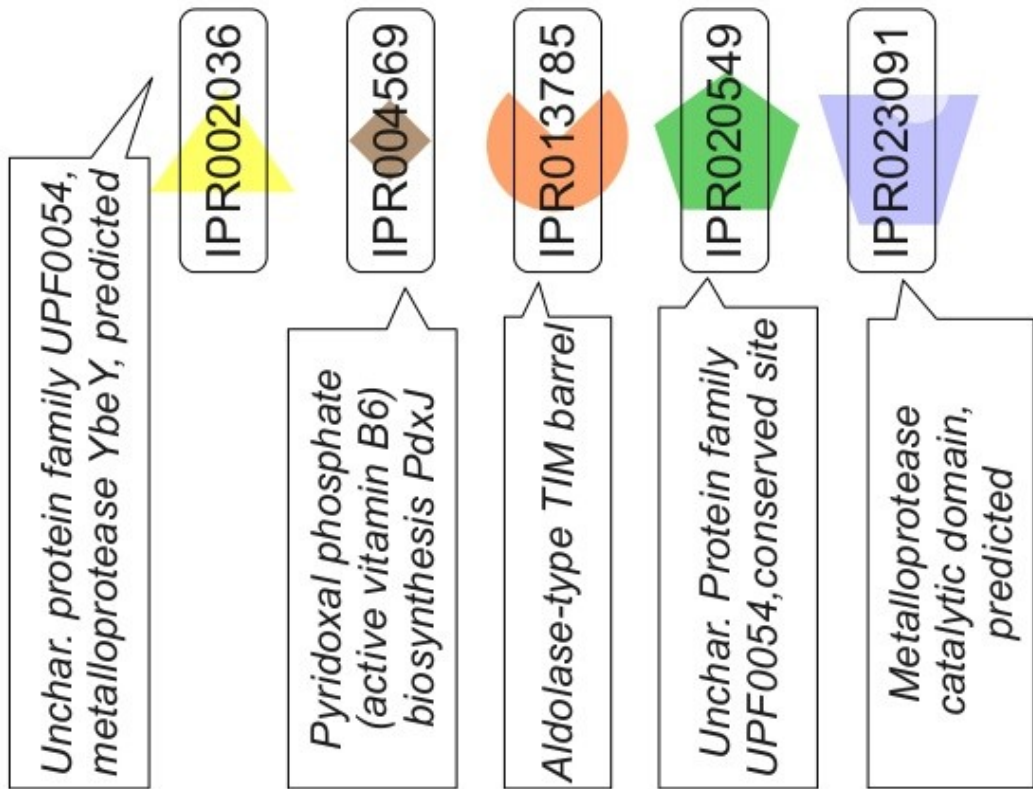
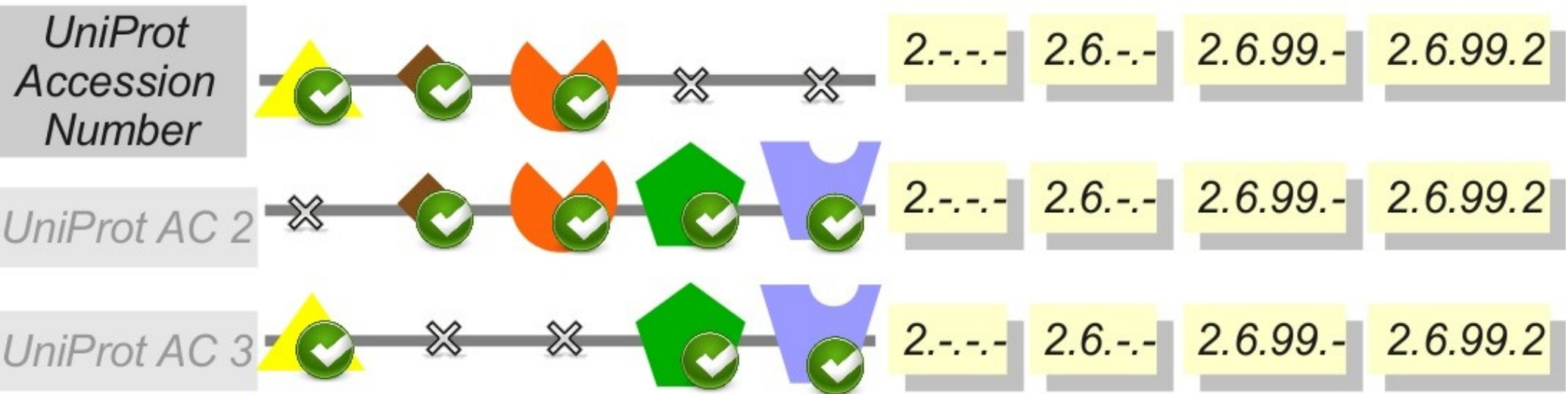
# Data sources

HOW



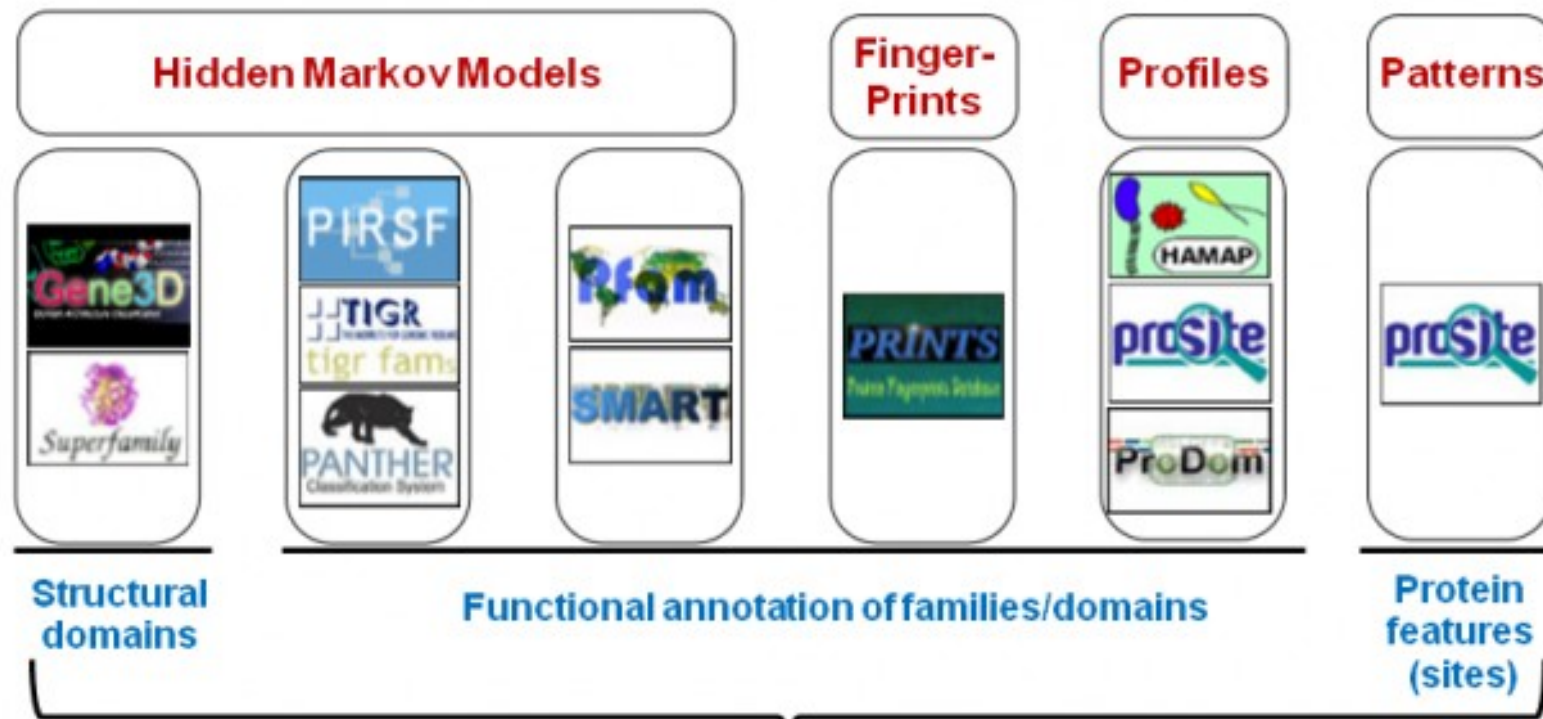
*EnzML: Multi-label prediction of enzyme classes using InterPro signatures,  
L. De Ferrari et al. BMC Bioinformatics 2012, 13:61*





Data  
schema  
(1,108 Uniprot proteins)

# InterPro sequence signatures



From: <http://www.ebi.ac.uk/training/online/course/interpro-quick-tour/what-interpro>



# InterPro signatures

[UniProtKB/Swiss-Prot : P11766](#)

Scale:10aa

**ADHX\_HUMAN**

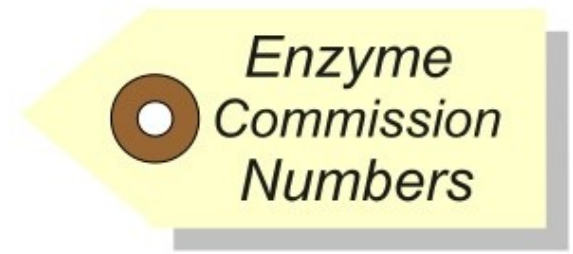
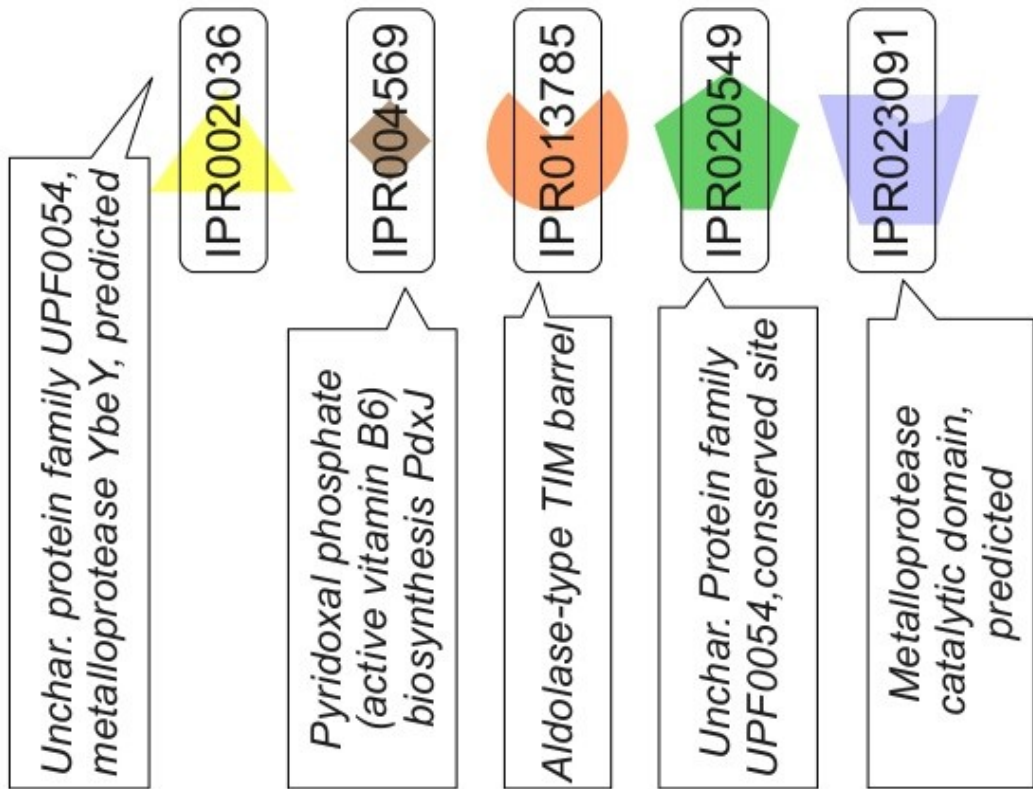
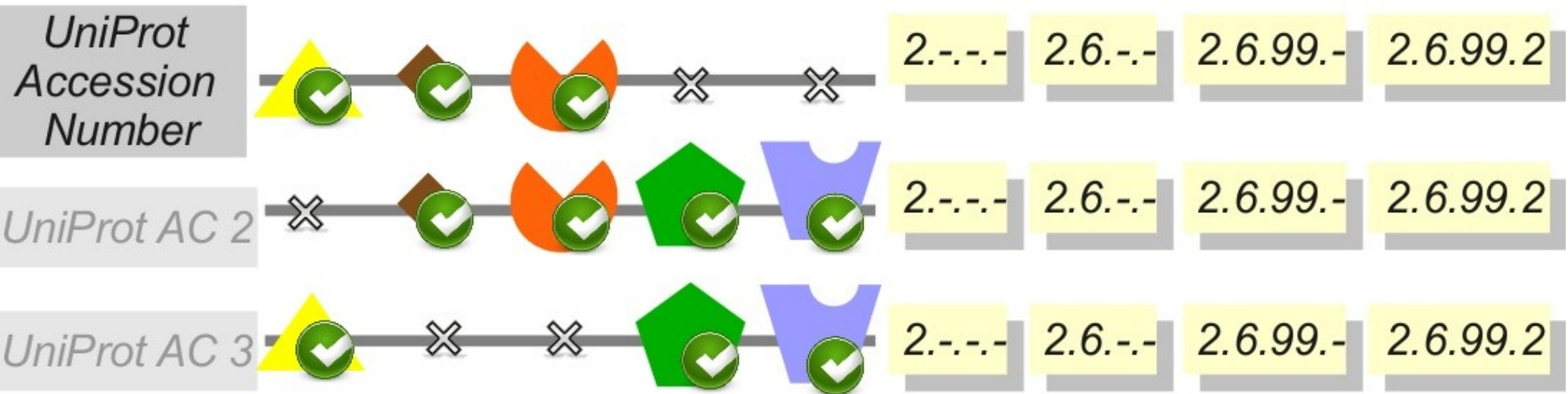
G-H-E-x-{EL}-G-{AP}-x(4)-[GA]-x(2)-[IVSAC]

## InterPro Signatures



From: <http://www.ebi.ac.uk/interpro/ISpy?ac=P11766>





Data  
schema  
(1,108 Uniprot proteins)

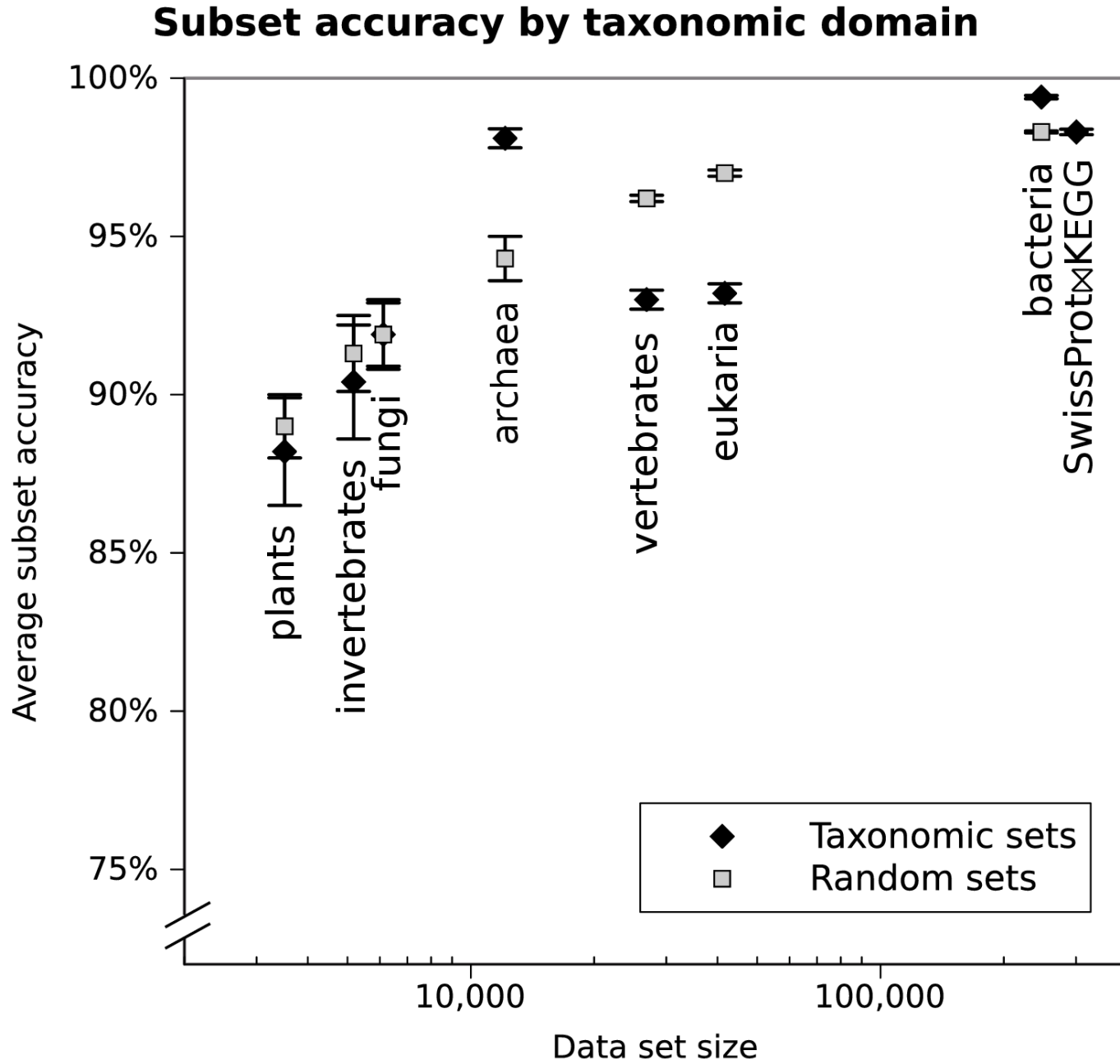
# Algorithm

- Mulan <http://mulan.sourceforge.net/>
  - Multi-Label algorithms library, Weka compatible (Java)
- Fastest/best of the Mulan 1.2.0 algorithms
  - BR-kNN a Multi-Label **Nearest Neighbours** implementation
  - (*E. Spyromitros, et al. An Empirical Study of Lazy Multilabel Classification Algorithms, SETN 2008, Syros, Greece*)
  - k=1 Neighbour

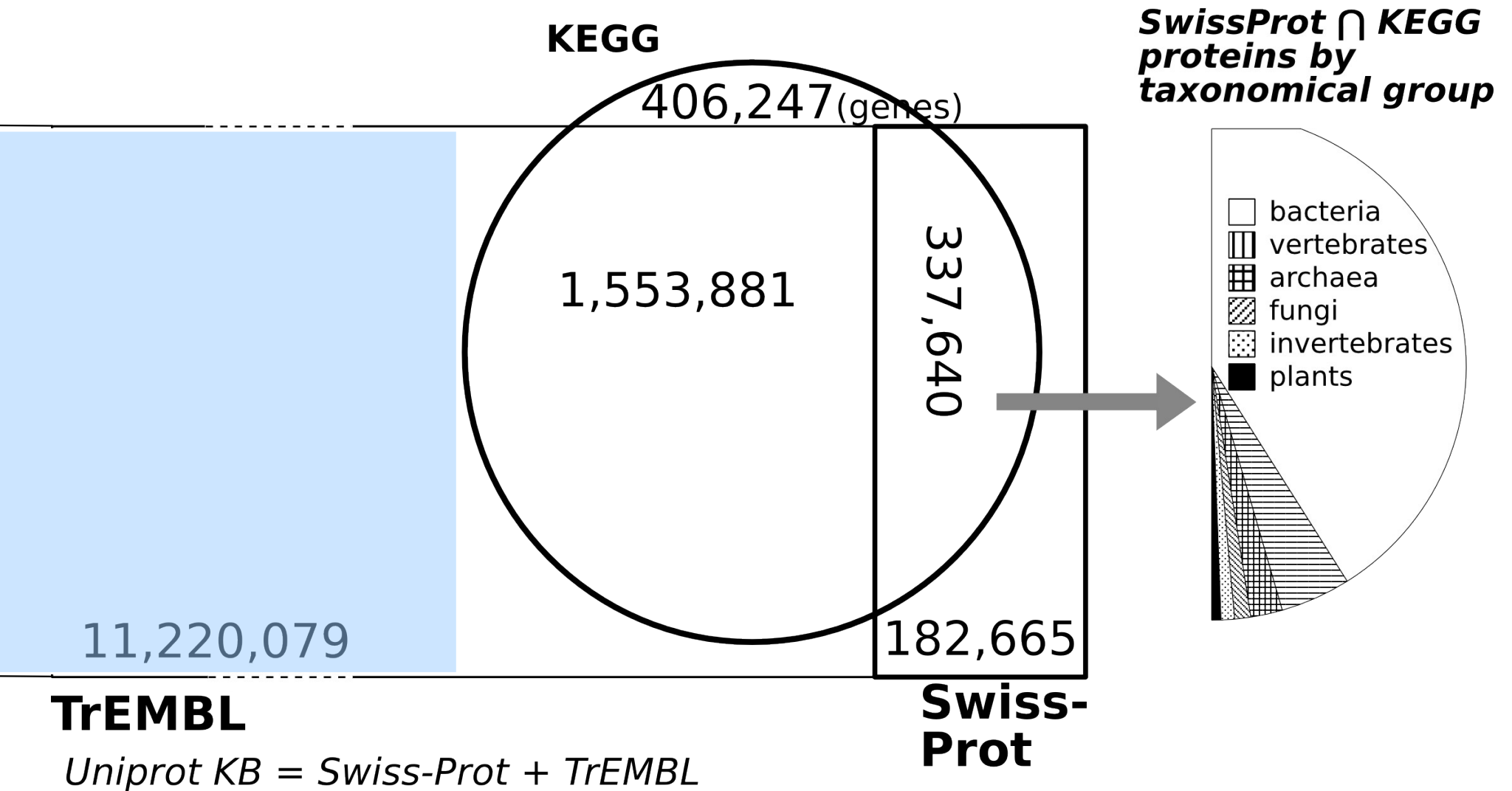
# SwissProt $\bowtie$ KEGG cross-evaluation

# Results

L. De Ferrari et al. BMC Bioinformatics 2012, 13:61



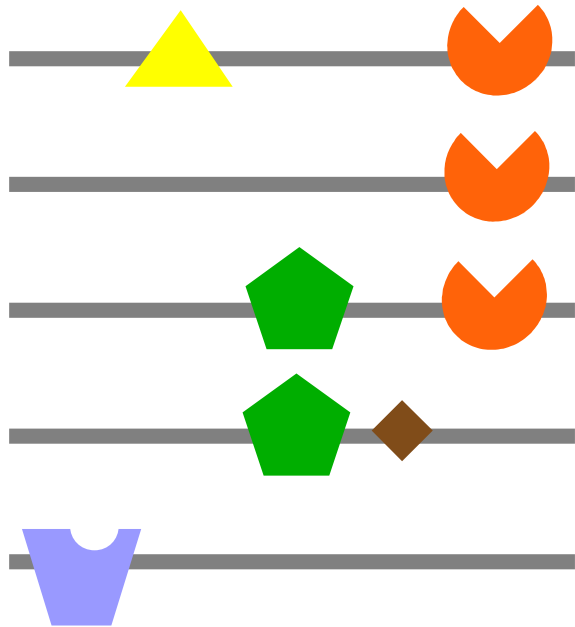
# Biocuration



# Dynamics of collaborative bio-curation + machine learning



EC 1.2.3.4



# Bio-curation + machine learning

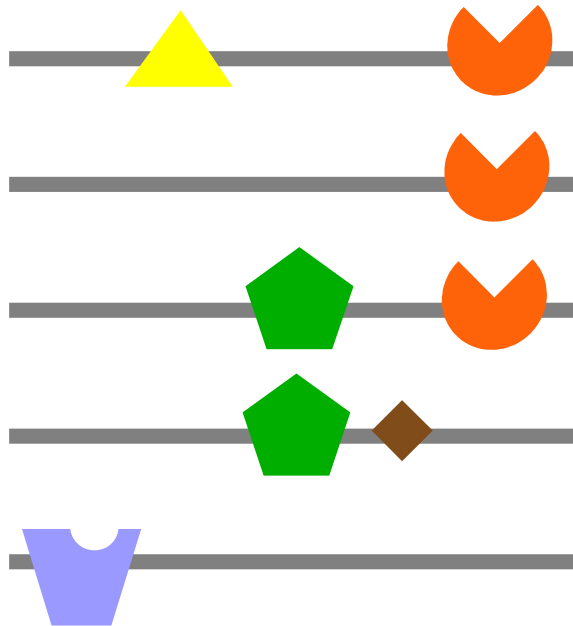


**EC 1.2.3.4**

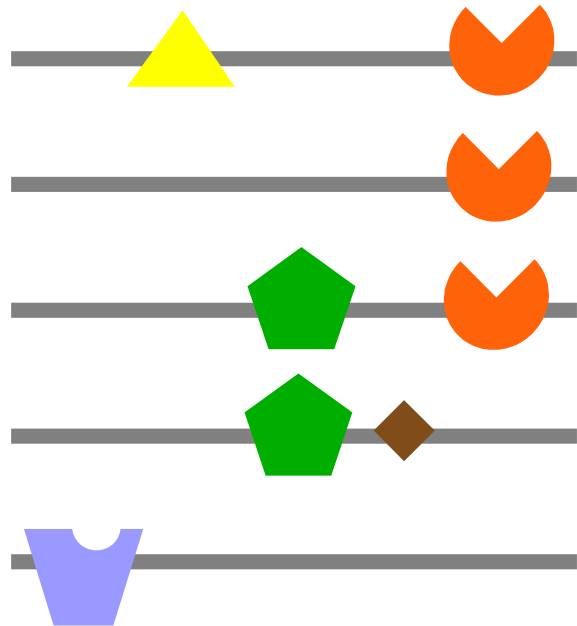
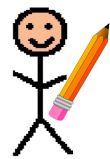
EC 1.2.3.4

EC 1.2.3.4

EC 1.2.3.4



# Bio-curation + machine learning



**EC 1.2.3.4**

EC 1.2.3.4

EC 1.2.3.4

EC 1.2.3.4

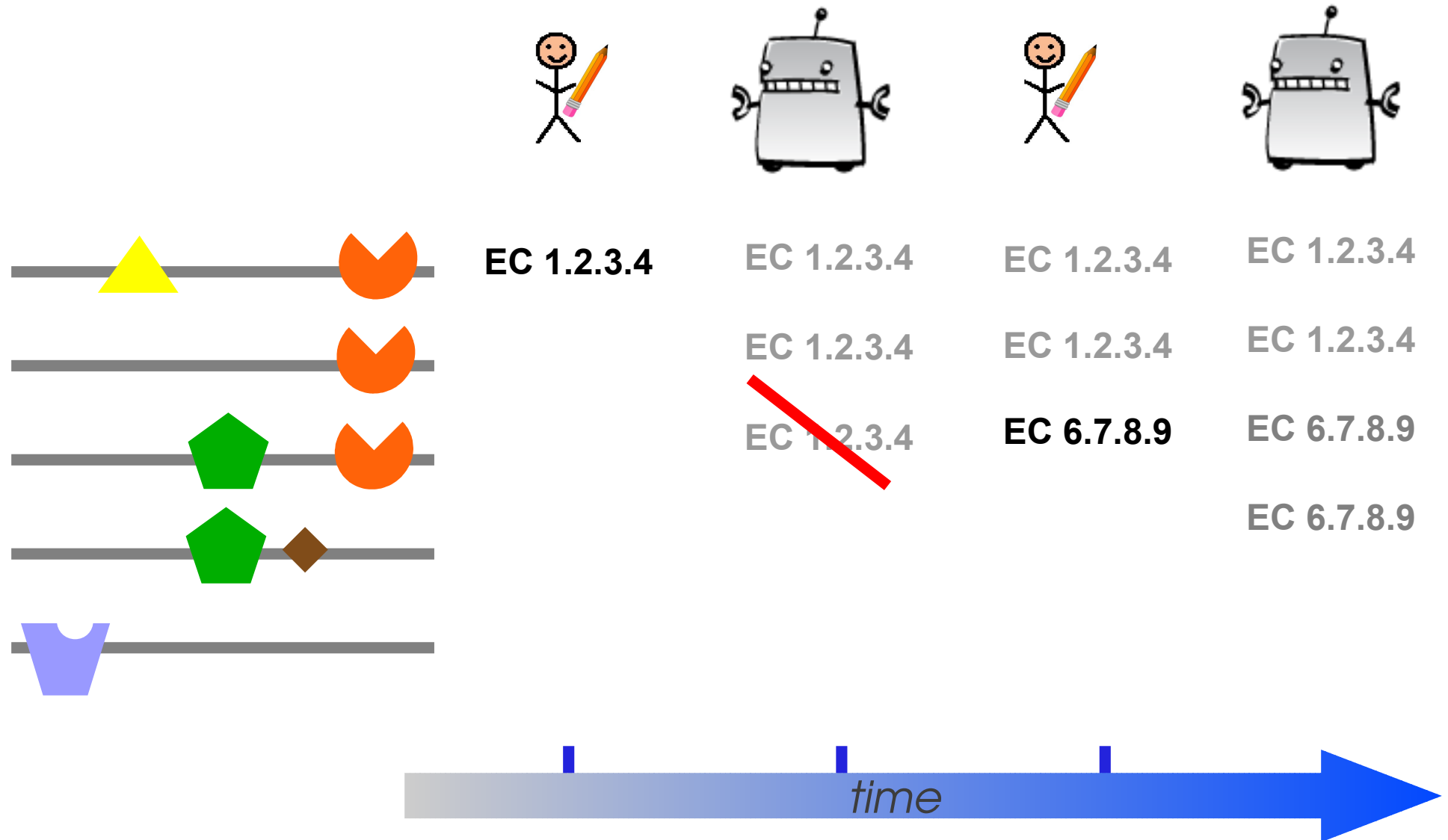
EC 1.2.3.4

~~EC 1.2.3.4~~

**EC 6.7.8.9**

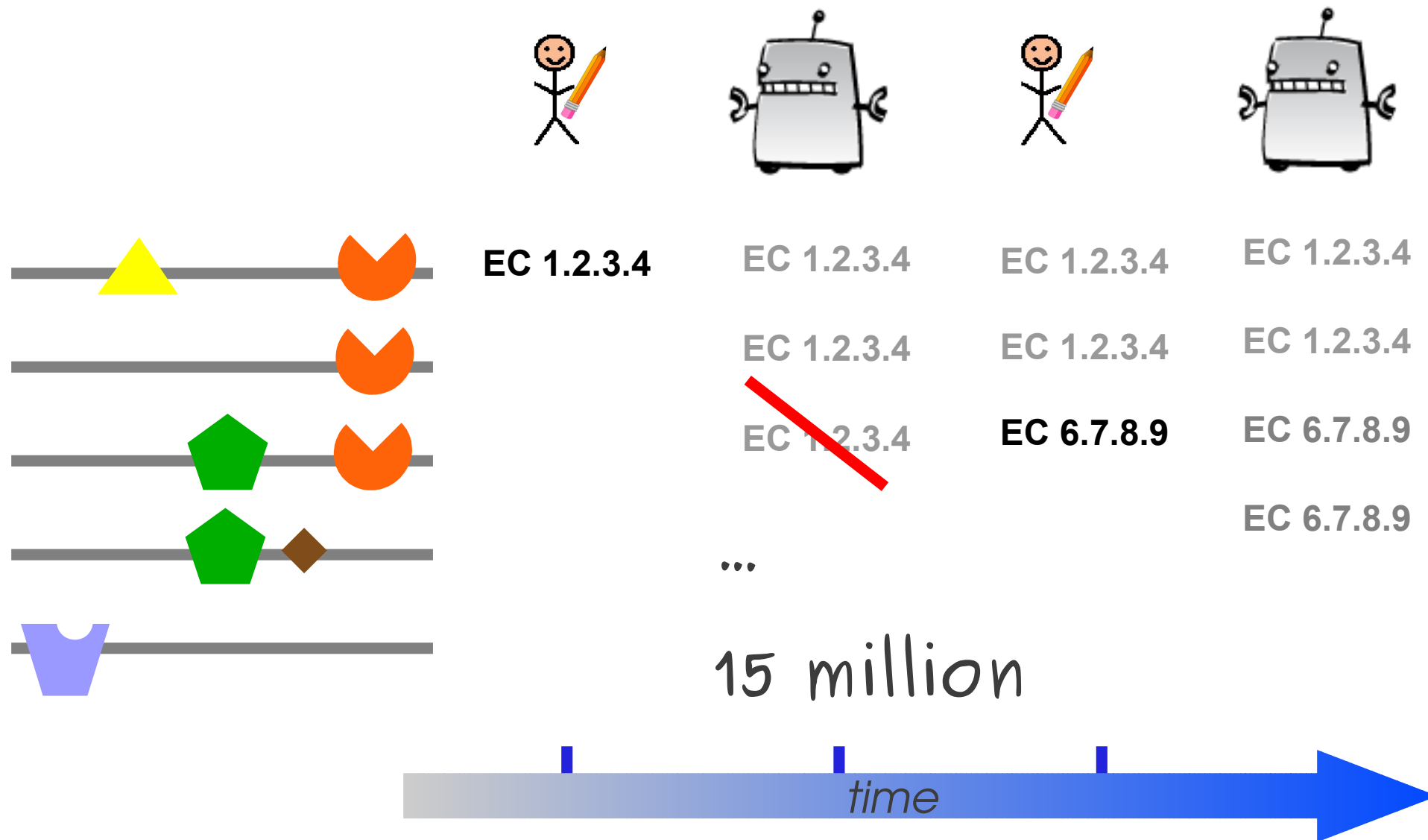


# Bio-curation + machine learning





# Bio-curation + machine learning

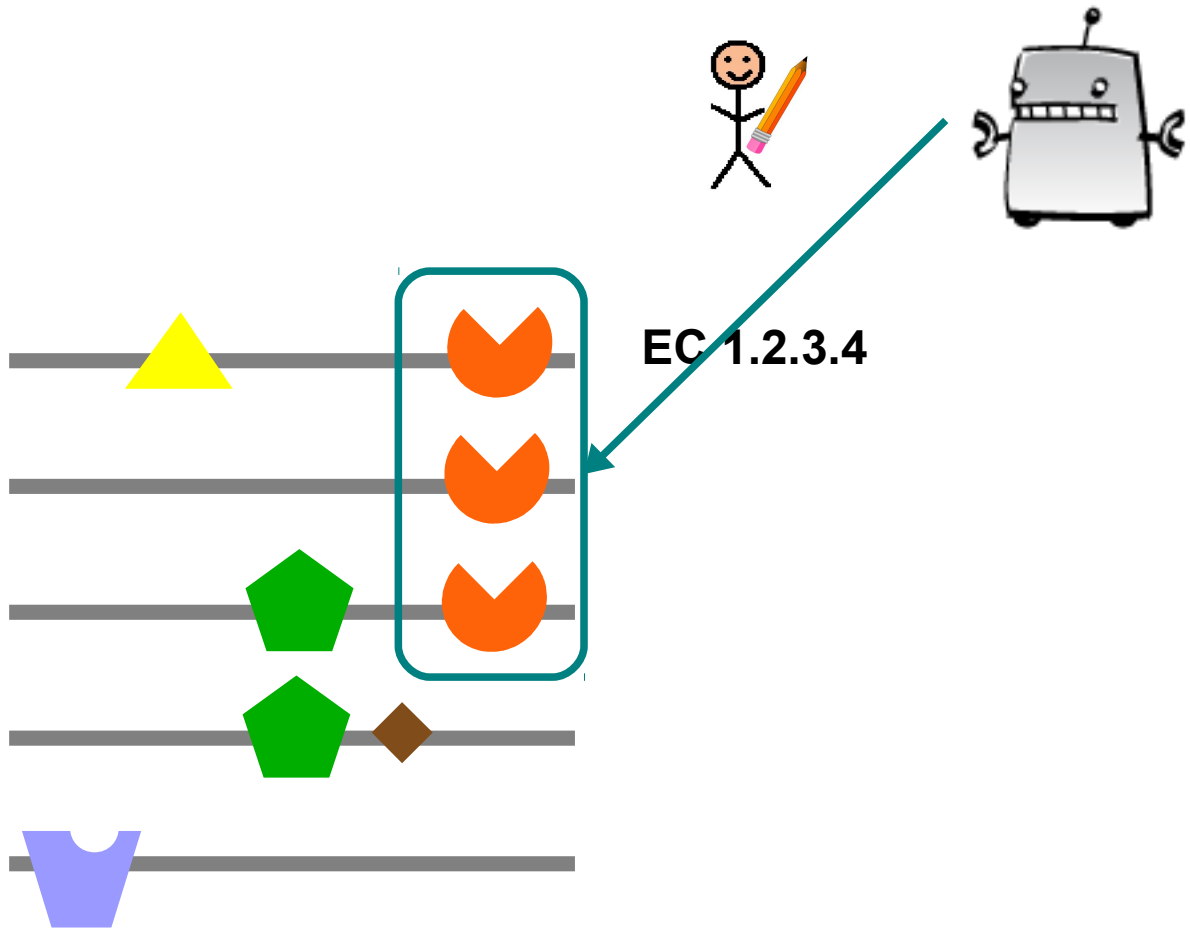


# Active learning

- Abundant unlabelled data
- Labels are expensive to obtain
- Allow the machine learning to **choose** the datum to be labelled next

*Active Learning Literature Survey. Burr Settles. Computer Sciences  
Technical Report 1648, University of Wisconsin–Madison*

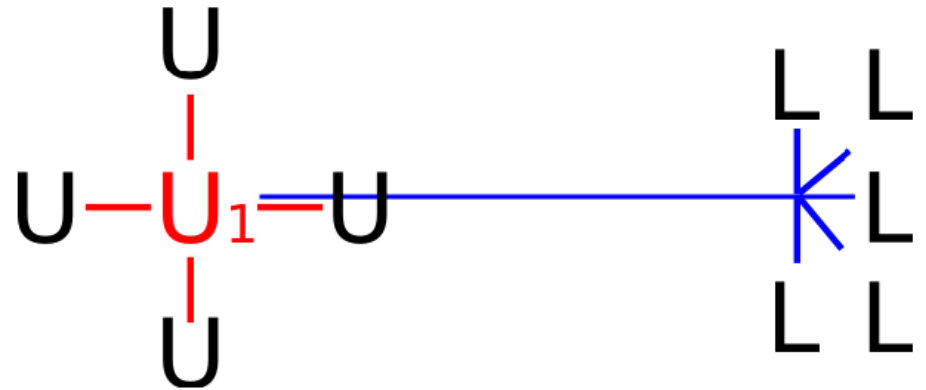
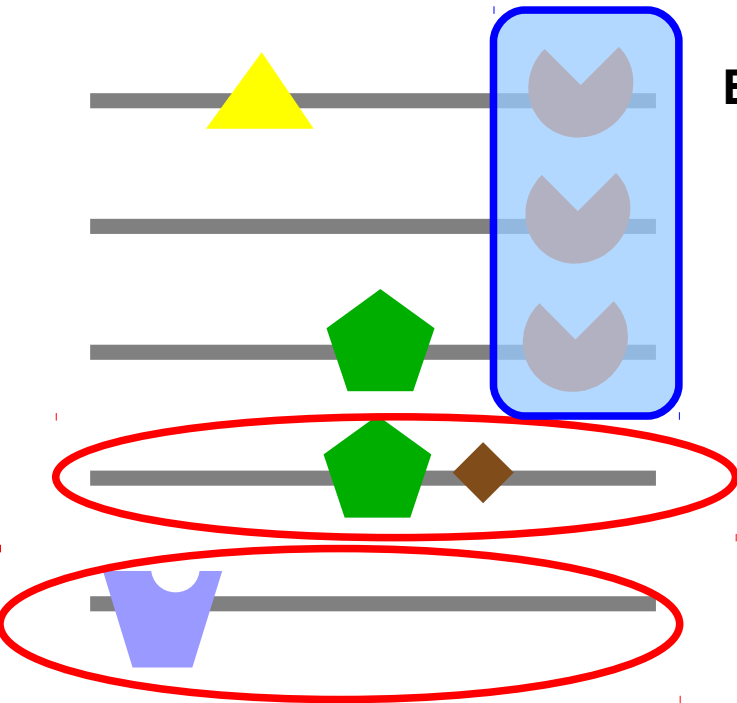
# Active learning: different from labelled & similar to unlabelled



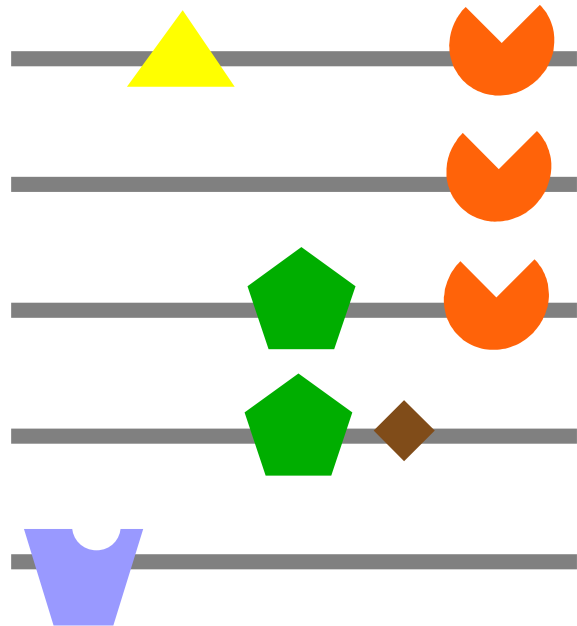
# Active learning: different from labelled & similar to unlabelled



EC 1.2.3.4



# Active learning: confidence



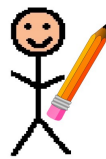
**EC 1.2.3.4**

**EC 1.2.3.4**

EC 1.2.3.4

EC 1.2.3.4

# Active learning: confidence



EC 1.2.3.4



EC 1.2.3.4



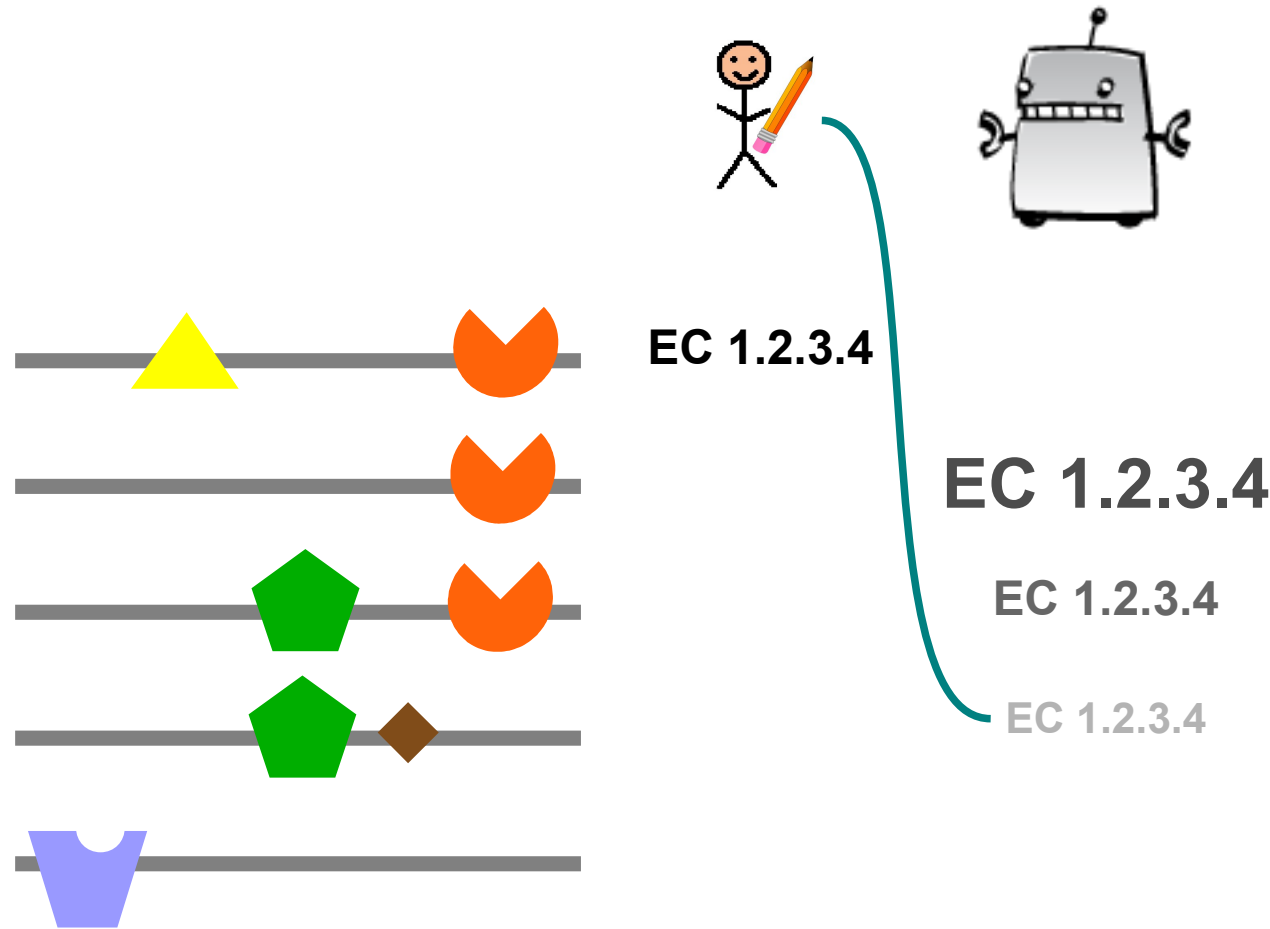
EC 1.2.3.4



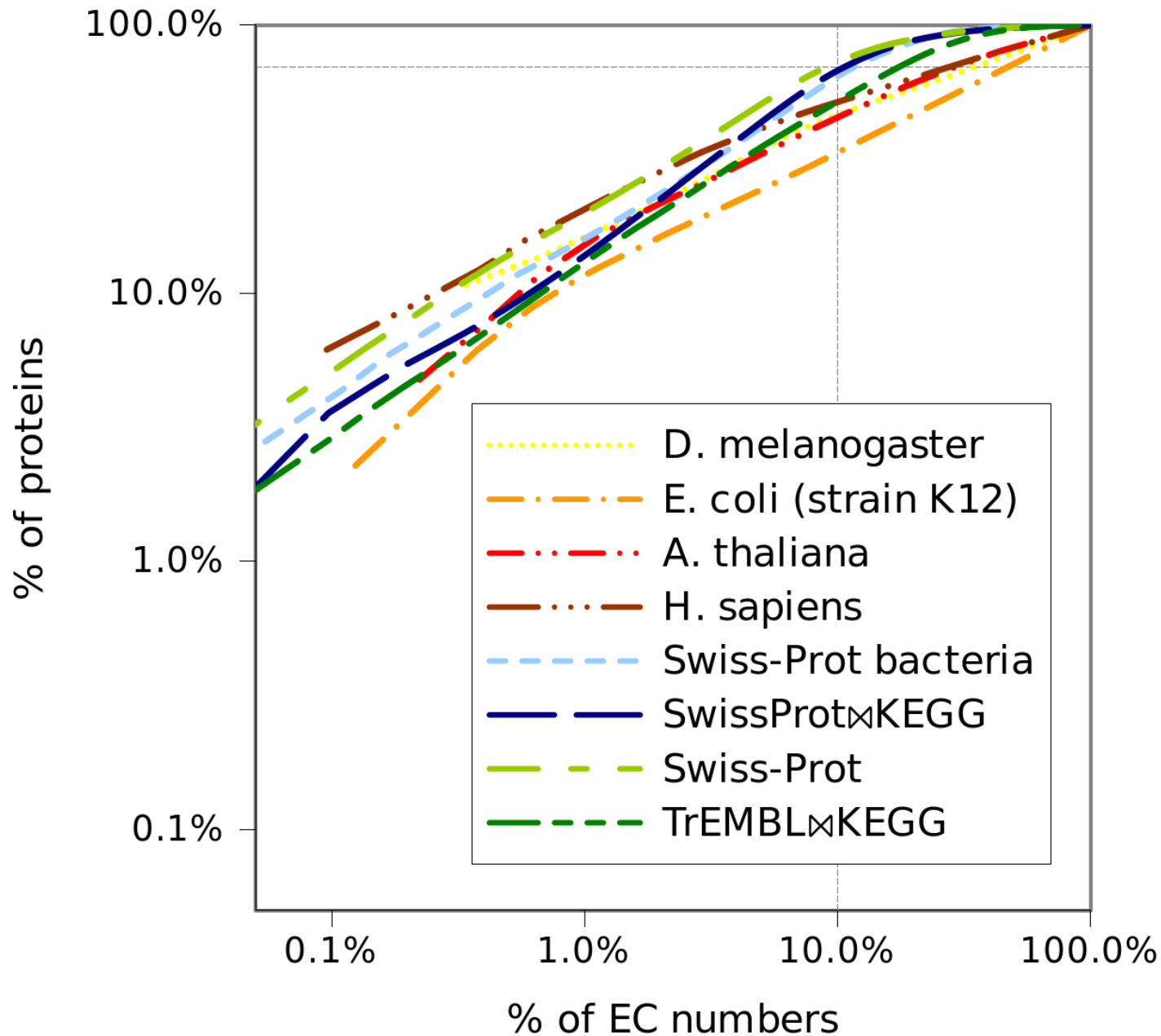
EC 1.2.3.4



# Active learning disadvantages: non-parallel & poor on rare classes



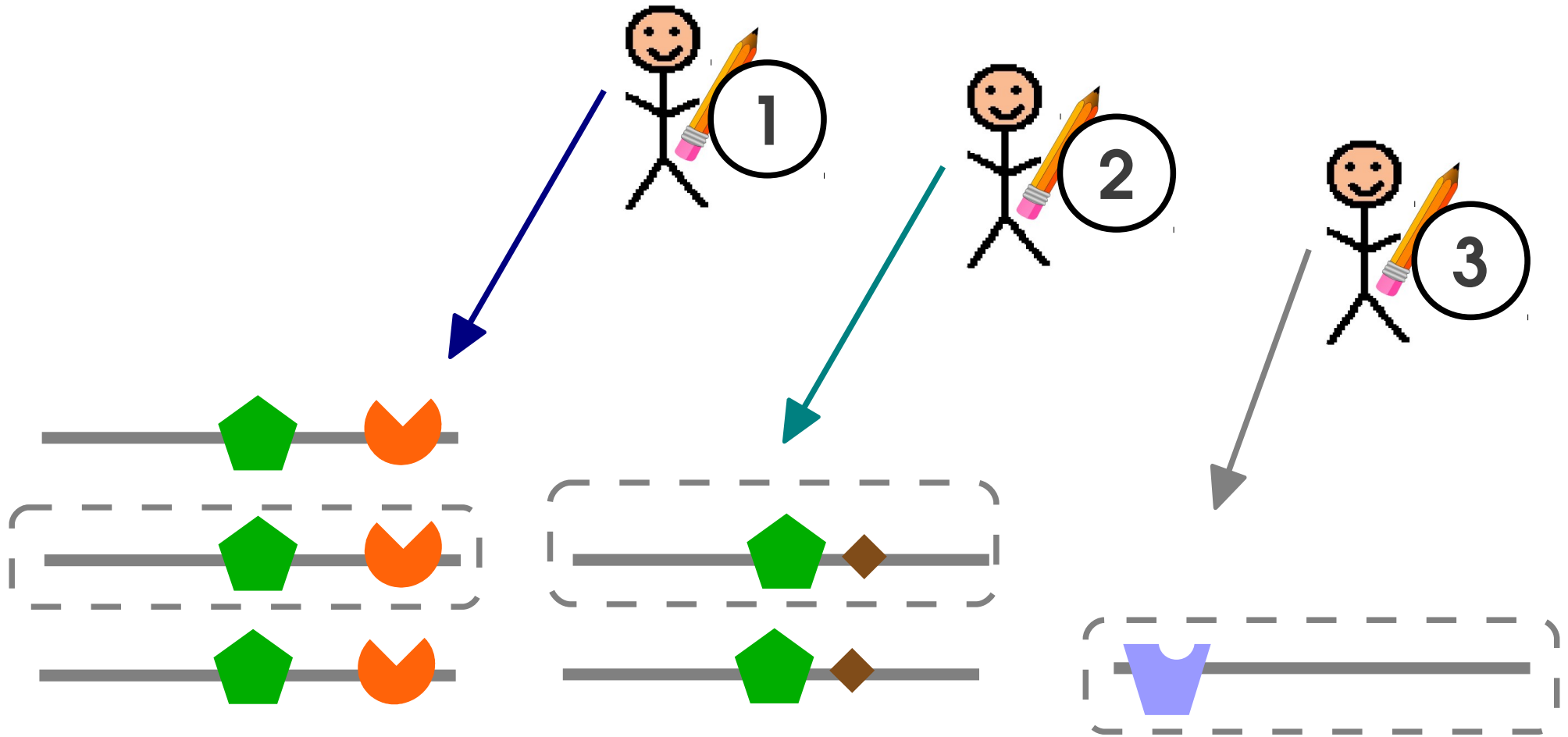
# Distribution of Enzyme Commission numbers





# Guided learning

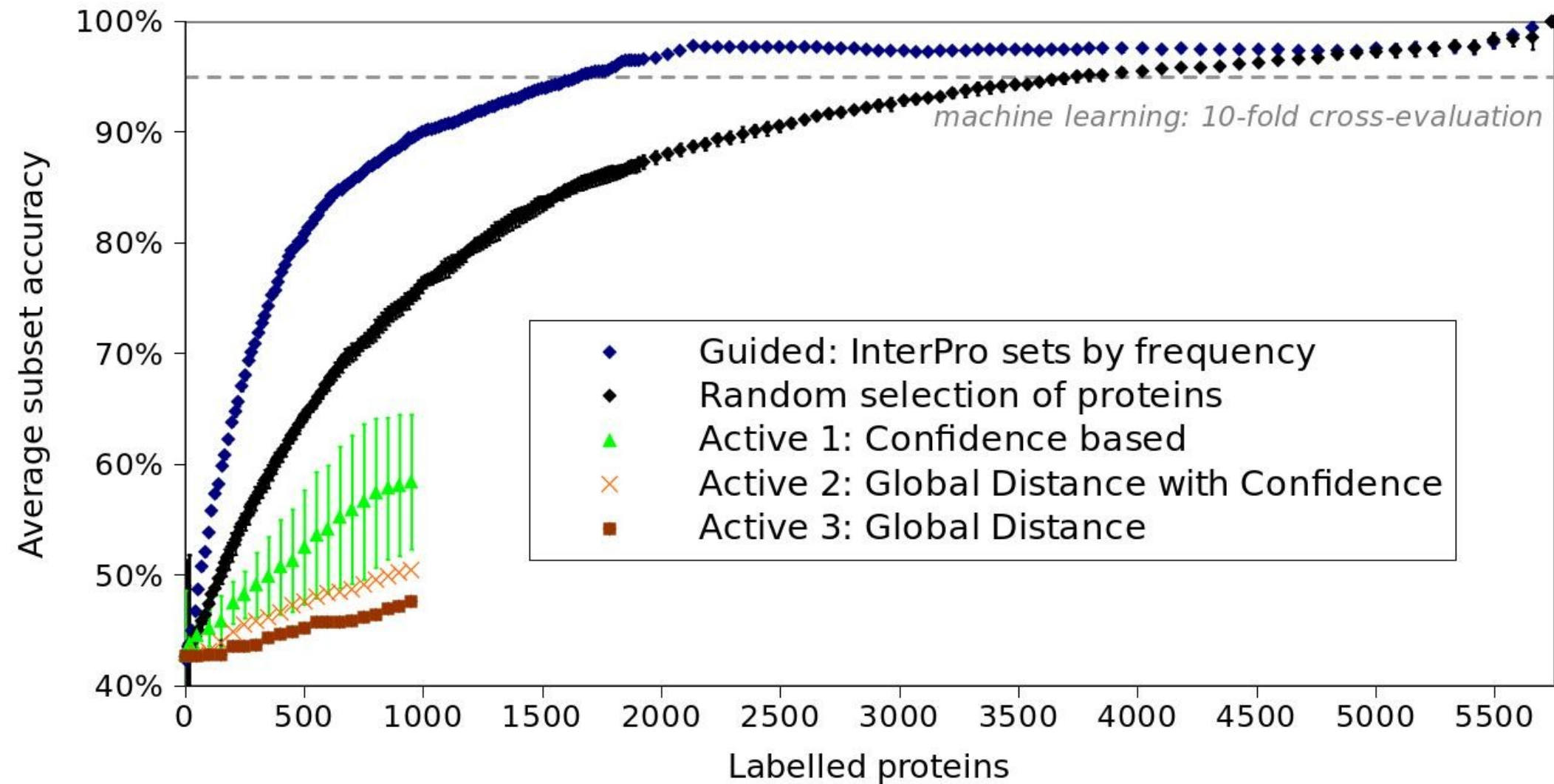
tackle InterPro signature sets by frequency



**Inspired by:** Attenberg, J. & Provost, F.

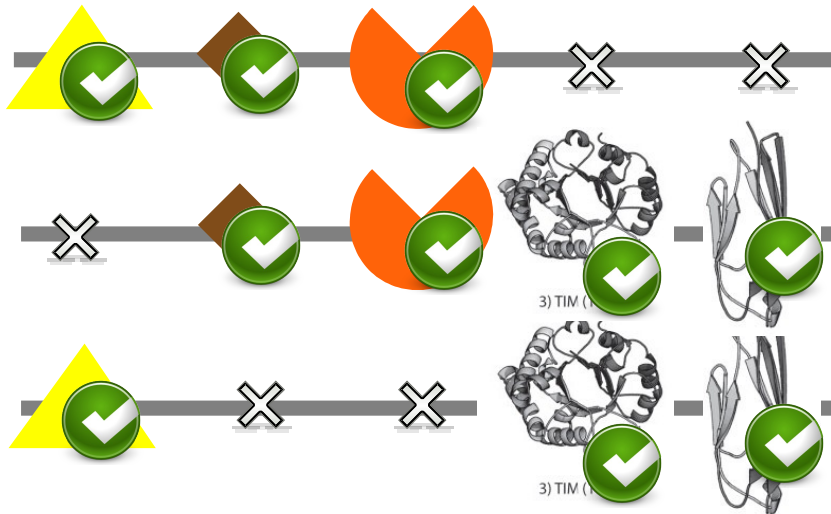
*Why label when you can search?: alternatives to active learning for applying human resources to build classification models under extreme class imbalance.* 16th ACM SIGKDD, ACM, 2010, 423-43

# Guided vs. Active Learning



## Attributes

## Class labels



MACiE MECHANISM  
 pyridoxine  
 5'-phosphate  
 synthase

Unchar. protein family UPF0054,  
 metalloprotease YbeY, predicted

IPR002036

Pyridoxal phosphate  
 (active vitamin B6)  
 biosynthesis PdxJ

IPR004569

Aldolase-type TIM barrel

IPR013785

TIM barrel

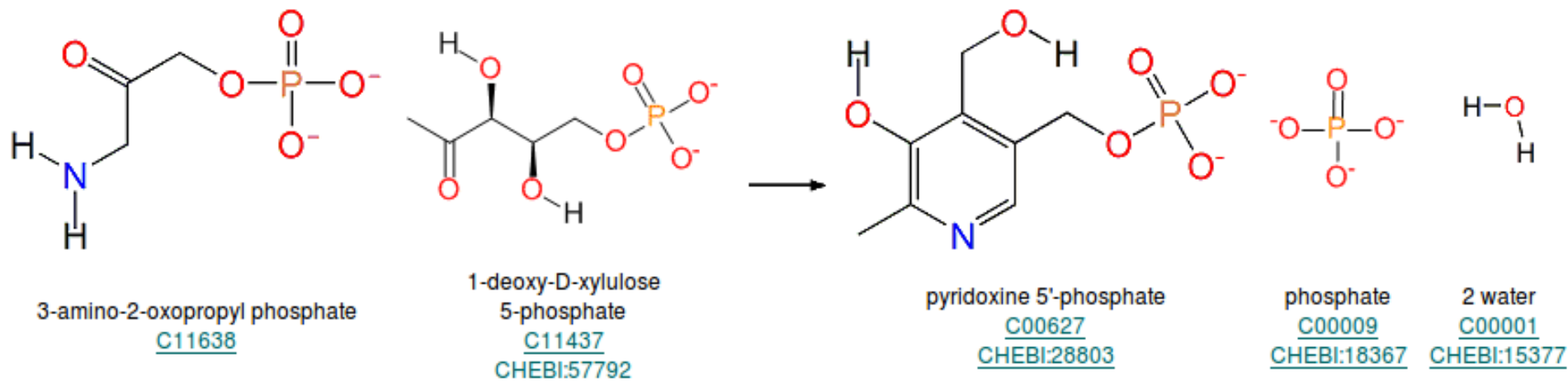


Twitchin



Current work:  
 predict enzyme  
**chemical  
 mechanism**

# MACiE: Mechanism, Annotation and Classification in Enzymes



*Overall Comment:* This enzyme forms part of the vitamin B6 biosynthesis pathway in bacteria.

[View similar reactions](#)

## Stepwise Description of the Reaction

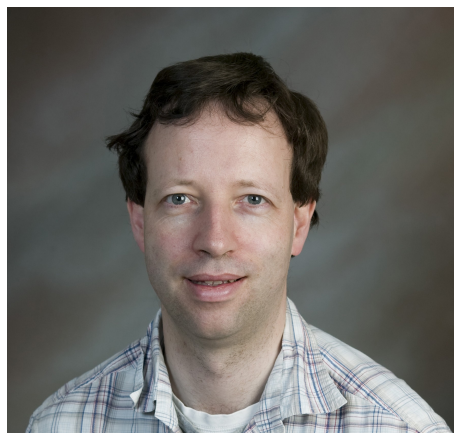
<a href="#">Step 1</a>	The 3-amino-2-oxopropyl phosphate substrate initiates a nucleophilic attack upon the 1-deoxy-D-xylulose 5-phosphate in the first step of a Schiff base formation. The carbonyl oxygen is assumed to deprotonate the amine.
<a href="#">Step 2</a>	The newly formed secondary amine initiates an intramolecular elimination of water, which obtains its proton first from Glu72, and then from water.
<a href="#">Step 3</a>	His45 deprotonates the C3 of the intermediate, which initiates double bond rearrangement, with the Schiff base acting as an electron sink.
<a href="#">Step 4</a>	The Schiff base reforms, initiating a double bond rearrangement and resulting in the elimination of water, which obtains its proton from Glu72.
<a href="#">Step 5</a>	Glu72 deprotonates the remaining hydroxide, which initiates a double bond rearrangement and eliminates phosphate.

From: <http://www.ebi.ac.uk/thornton-srv/databases/cgi-bin/MACiE/entry/getPage.pl?id=M0243>

# Mitchell's group - School of Chemistry, University of St Andrews, United Kingdom



Lazaros Mavridis



John Mitchell



Luna De Ferrari



Neetika Nath



James Mc Donagh



Rosanna Alderson

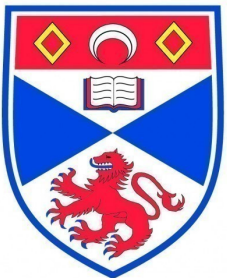
Stuart Aitken – Univ. of Edinburgh, School of Informatics & SynthSys

Thank you



Questions?

*Luna De Ferrari*  
*ldf@st-andrews.ac.uk*

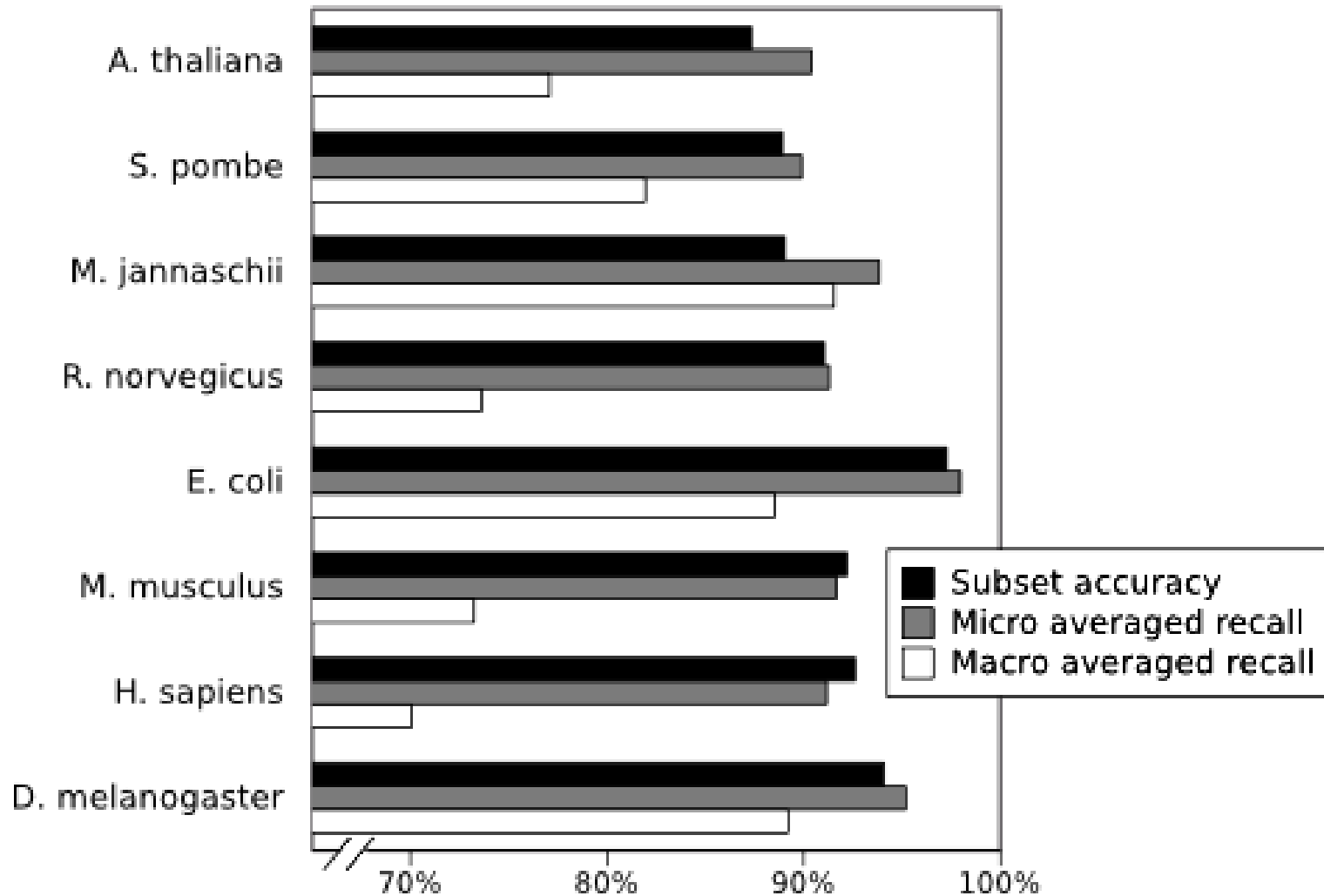


University  
of  
St Andrews

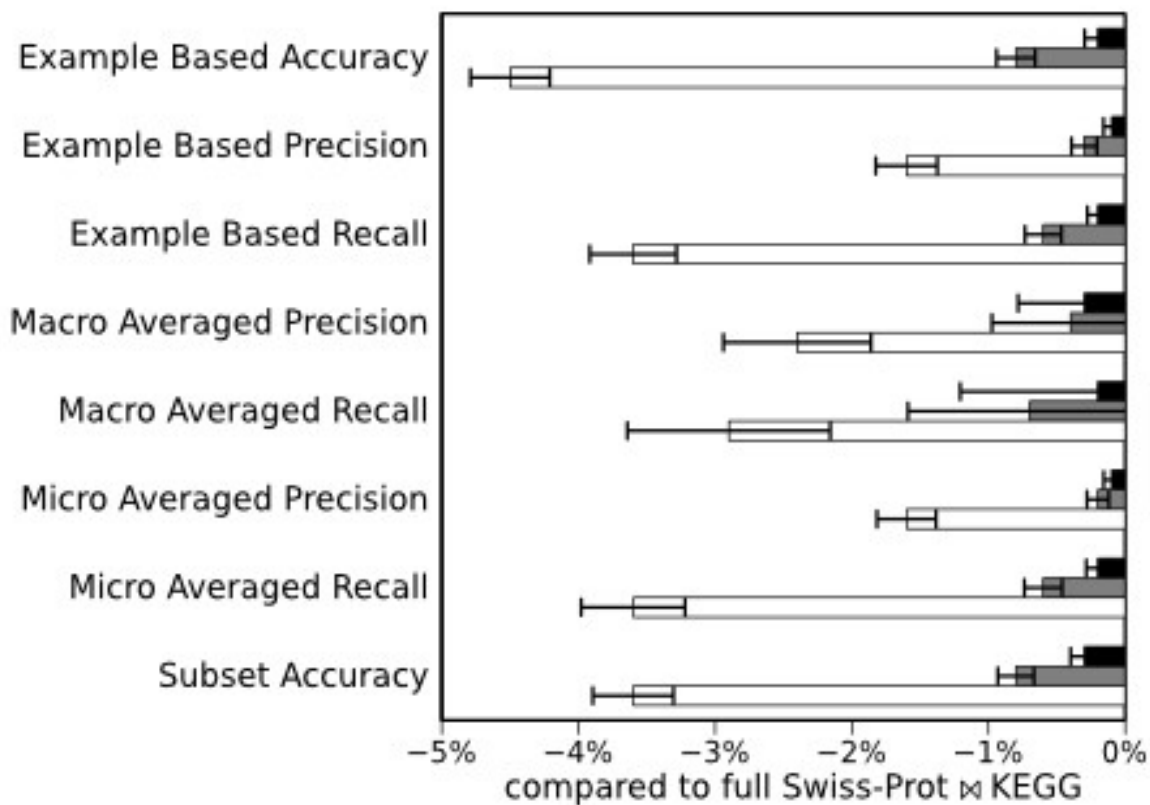


# Proteome re-annotation

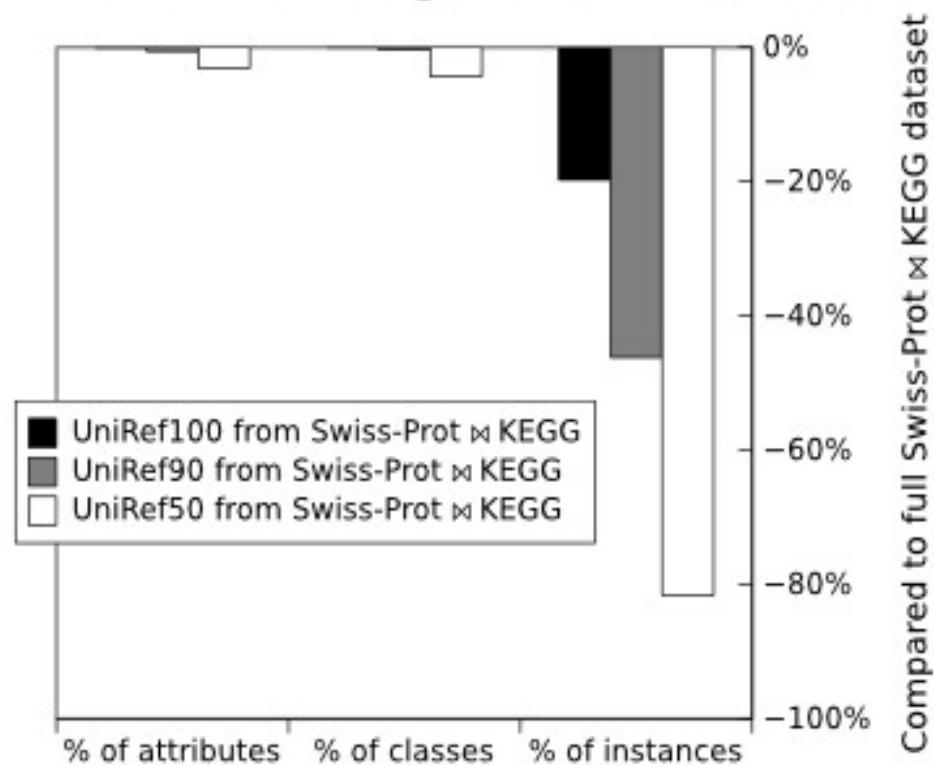
**Reannotation of a single proteome by training on the SwissProt×KEGG dataset**



### Cross evaluation on UniRef reference sequences from Swiss-Prot $\bowtie$ KEGG



### Number of proteins, EC numbers and InterPro signatures in the datasets



*EnzML: Multi-label prediction of enzyme classes using InterPro signatures, L. De Ferrari et al. BMC Bioinformatics 2012, 13:61*



- Why label when you can search?: alternatives to active learning for applying human resources to build classification models under extreme class imbalance
- Authors: Josh Attenberg Polytechnic Institute of NYU, Brooklyn, NY, USA
- Foster Provost NYU Stern School of Business, New York, NY, USA
- KDD '10 Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining
- Pages 423-432
- ACM New York, NY, USA ©2010
- This paper analyses alternative techniques for deploying low-cost human resources for data acquisition for classifier induction in domains exhibiting extreme class imbalance - where traditional labeling strategies, such as active learning, can be ineffective. Consider the problem of building classifiers to help brands control the content adjacent to their on-line advertisements. Although frequent enough to worry advertisers, objectionable categories are rare in the distribution of impressions encountered by most on-line advertisers - so rare that traditional sampling techniques do not find enough positive examples to train effective models. An alternative way to deploy human resources for training-data acquisition is to have them "guide" the learning by searching explicitly for training examples of each class. We show that under extreme skew, even basic techniques for guided learning completely dominate smart (active) strategies for applying human resources to select cases for labeling. Therefore, it is critical to consider the relative cost of search versus labeling, and we demonstrate the tradeoffs for different relative costs. We show that in cost/skew settings where the choice between search and active labeling is equivocal, a hybrid strategy can combine the benefits.