

# Estimating the Quality of Ontology-Based Annotations by Considering Evolutionary Changes

Anika Gross, Michael Hartung, Toralf Kirsten,  
Erhard Rahm



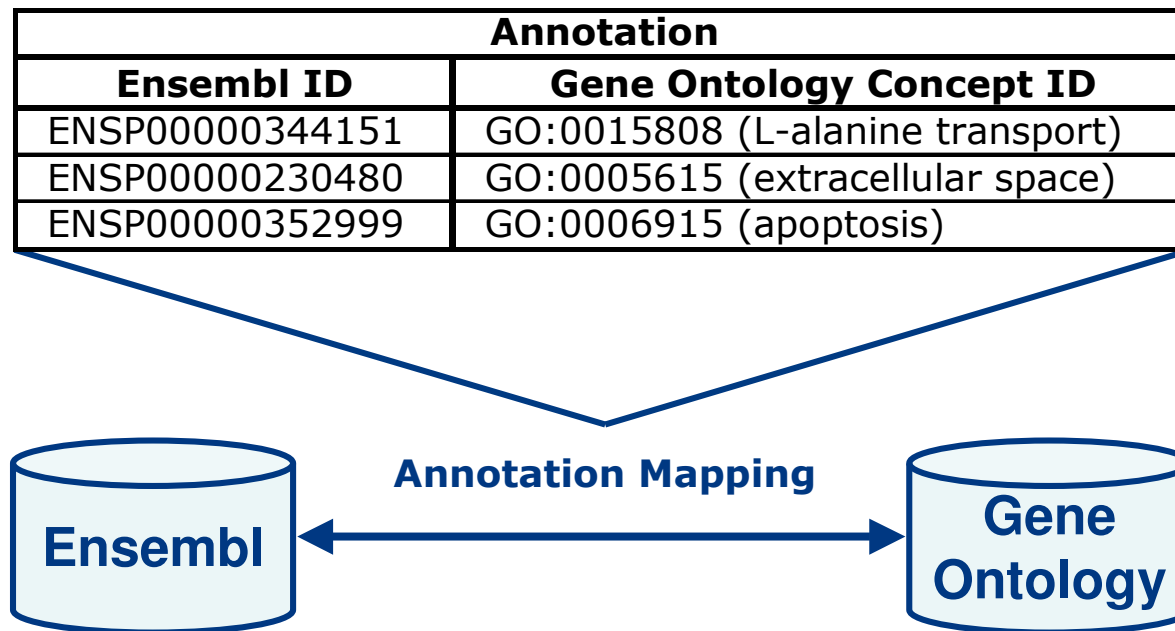
Interdisciplinary Centre for Bioinformatics  
<http://www.izbi.uni-leipzig.de>

Database Group Leipzig  
<http://dbs.uni-leipzig.de>

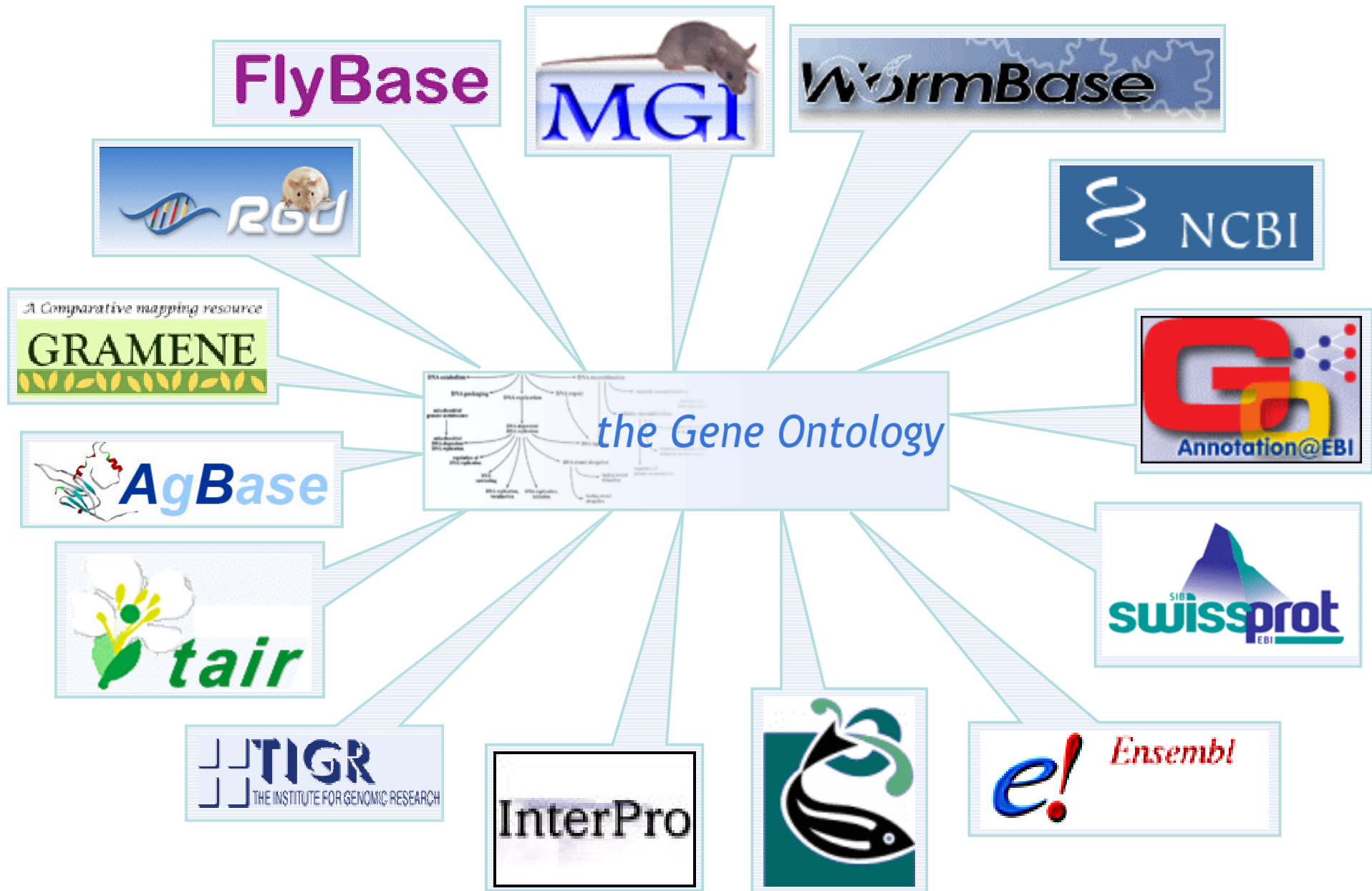


# Annotations in Life Sciences

- Increasing use of ontologies in life sciences, mainly ontology-based annotations
- **Annotations** Semantic descriptions of properties of biological objects, e.g., a protein is associated to a specific biological process

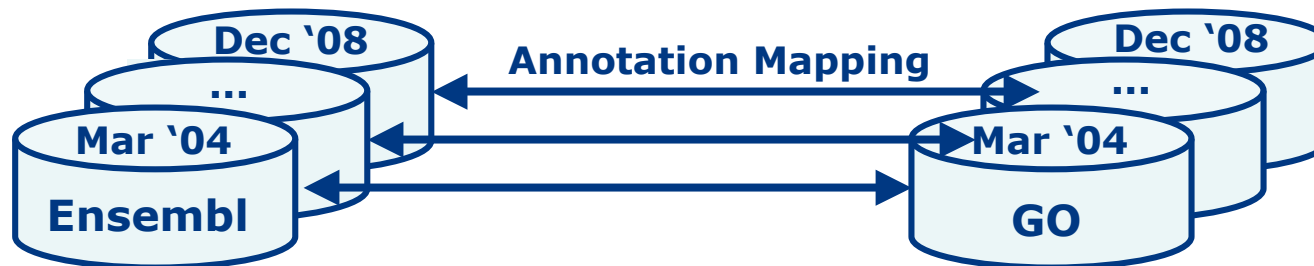


# Usage of Gene Ontology (GO)



# Motivation

- Domain knowledge changes
  - New findings, addition and revision of knowledge
  - Result: evolution of data sources
- First analysis at DILS 2008\*
  - Focused on evolution in ontologies and protein sources
- What about annotations?



Annotation		Provenance				
Ensembl ID	Gene Ontology Concept ID	V <sub>48</sub>	V <sub>49</sub>	V <sub>50</sub>	V <sub>51</sub>	V <sub>52</sub>
ENSP00000344151	GO:0015808 (L-alanine transport)	Green	Green	Green	Green	Green
ENSP00000230480	GO:0005615 (extracellular space)	Yellow	Yellow	Green	Yellow	Red
ENSP00000352999	GO:0006915 (apoptosis)	Green	-	-	-	Green

experimentally verified
author statement
automatically annotated

Dec 2007 – Dec 2008

- Different stability of annotations due to different evolution and provenance changes

# Application of GO Annotations

- Functional profiling of large data sets (e.g. gene expression microarrays) to find significantly shared GO terms

GO Term	Aspect	P-value	Sample frequency	Background frequency	Genes
<a href="#">GO:0002376</a> immune system process	P	1.02e-07	10/14 (71.4%)	1052/19635 (5.4%)	<a href="#">Q9NZ08</a> <a href="#">P42081</a> <a href="#">O15533</a> <a href="#">Q6P179</a> <a href="#">P19838</a> <a href="#">Q9NZQ7</a> <a href="#">P33681</a> <a href="#">Q03519</a>
<a href="#">GO:0048002</a> antigen processing and presentation of peptide antigen	P	3.26e-07	4/14 (28.6%)	18/19635 (0.1%)	<a href="#">Q9NZ08</a> <a href="#">O15533</a> <a href="#">Q6P179</a> <a href="#">Q03519</a>

[http://amigo.geneontology.org/cgi-bin/amigo/term\\_enrichment1](http://amigo.geneontology.org/cgi-bin/amigo/term_enrichment1)

- Unstable input annotations



- Impact on application results (Garbage In/Garbage Out principle)

# Quality of Annotations

## Possible criteria

- Correctness
- Completeness
- **Stability**
- **Provenance**
- ...



# Contributions

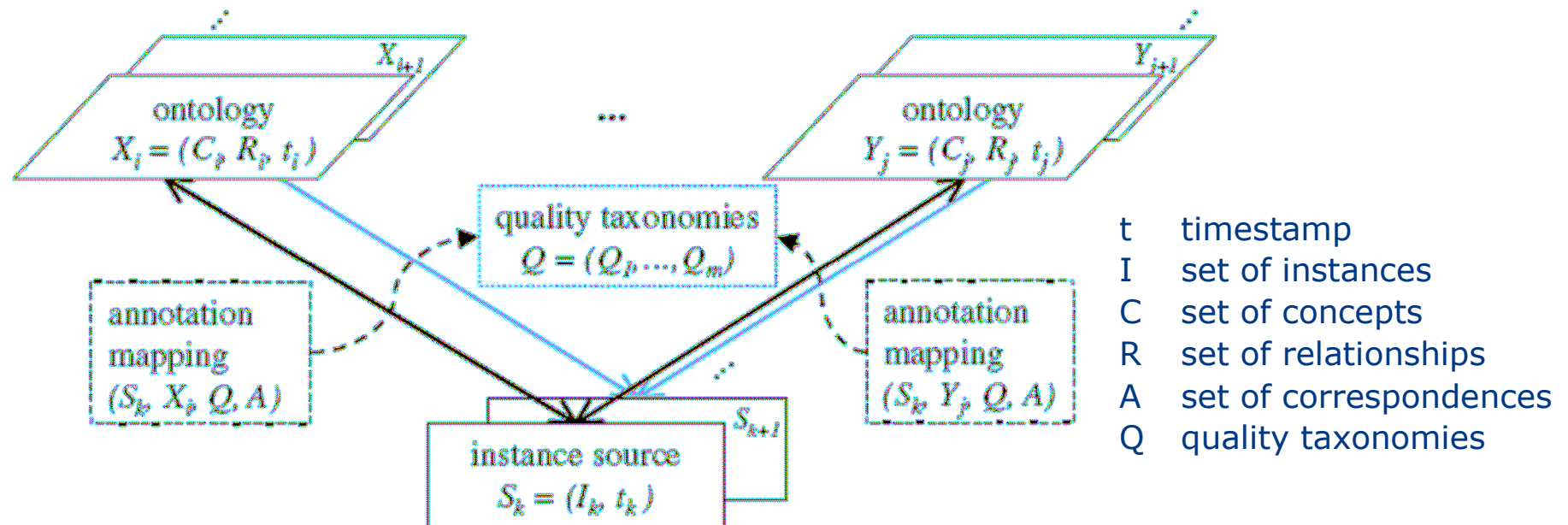
- General approach to analyze annotation mappings considering evolutionary changes
- Evolution-aware annotation model, change operations
- Evolution-based quality measures to identify reliable annotations (stability, provenance)
- Comparative evaluation in two large life science sources (Ensembl, Swiss-Prot)

# Overview

- *Motivation*
- *Annotations and quality*
- Annotation model, change operations
  - Quantitative evaluation of annotation evolution
- Estimating annotation quality
  - Stability measures
  - Evaluation results
- Conclusion



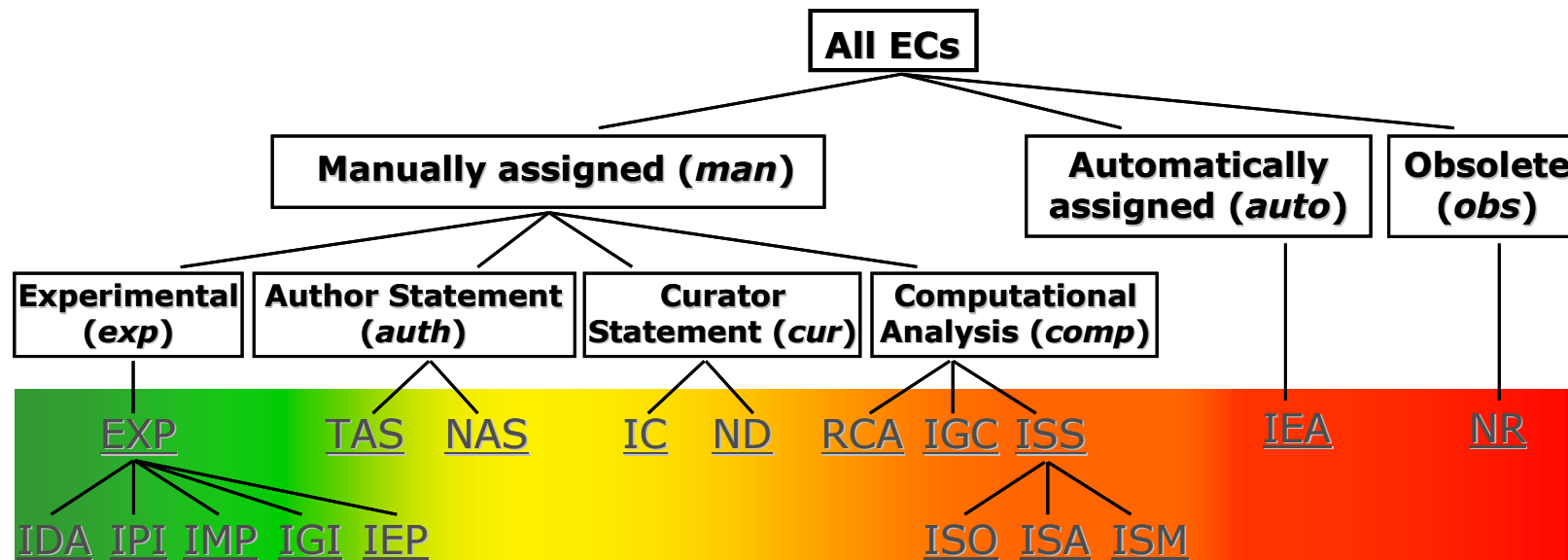
# Annotation Model



- Annotation  $a \in A$ ,  $a = (i, c, \{q\})$
- Linear versioning scheme
- Refers to quality taxonomies (e.g., provenance)

# Provenance Taxonomy - Evidence Codes

- Evidence Code (EC) \* = indicates how the annotation to a particular term has been derived, e.g., by which type of experiment or analysis



\* <http://www.geneontology.org/GO.evidence>

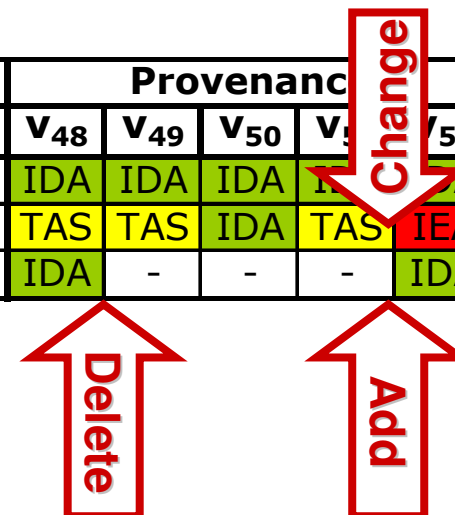
- Other taxonomies are possible
  - age
  - stability
  - ...

# Evolution Operations

- Evolution operations for proteins, ontologies, annotations
  - Add
  - Delete
  - Change
- Propagation of evolutionary changes when ontologies and proteins evolve!

Annotation		Provenance				
Ensembl ID	Gene Ontology Concept ID	V <sub>48</sub>	V <sub>49</sub>	V <sub>50</sub>	V <sub>51</sub>	V <sub>52</sub>
ENSP00000344151	GO:0015808 (L-alanine transport)	IDA	IDA	IDA	IDA	IDA
ENSP00000230480	GO:0005615 (extracellular space)	TAS	TAS	IDA	TAS	IEA
ENSP00000352999	GO:0006915 (apoptosis)	IDA	-	-	-	IDA

Dec 2007 - Dec 2008



# Quantitative Evolution Analysis

- Two large life science sources (Mar 2004 – Dec 2008)
- GO Annotations for human proteins



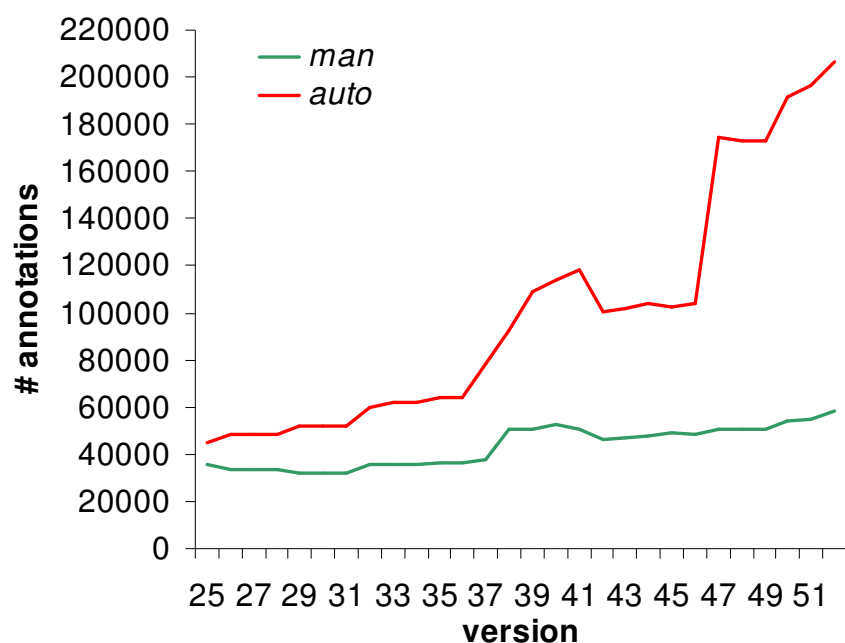
Swiss-Prot v<sub>47</sub>-v<sub>56</sub>



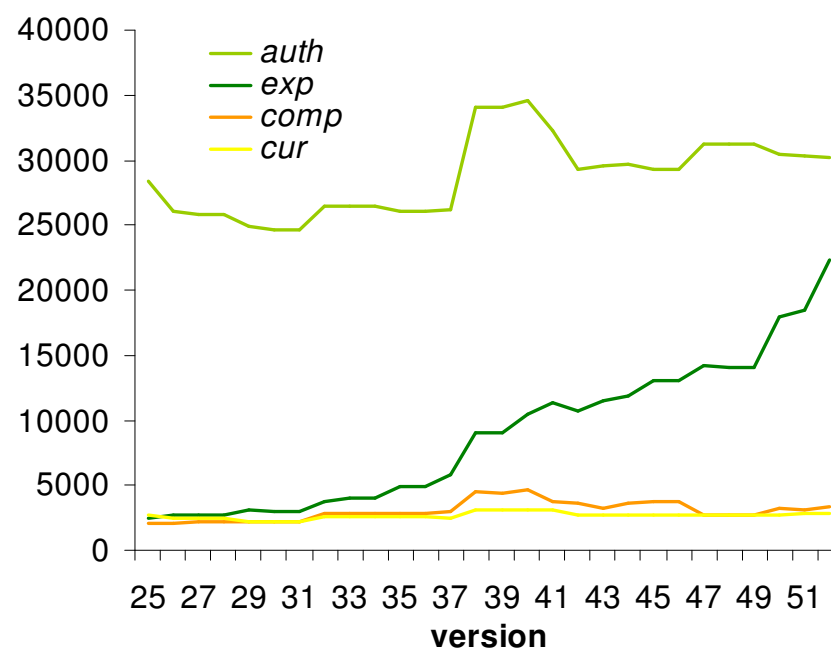
- Analysis of annotation evolution w.r.t.
  - Different provenance
  - Change operations

# Evolution of Annotations

Manually assigned  
vs. automatically assigned



“Subclasses”  
of manually assigned



- 78% of 265,000 automatically assigned
- $\text{growth}_{\text{auto}} 4.6$
- $V_{40} - V_{42}$  considerable number of deletions

- 22% manually assigned
- $\text{growth}_{\text{man}} 1.7$ 
  - $\text{growth}_{\text{exp}} 8.9$
  - $\text{growth}_{\text{auth}} 1.1$

# Evolution Operations

	<i>Add</i>		<i>Change</i>		<i>Delete</i>	
<i>exp</i>	25,979	6.6%	5,826	12.2%	7,575	3.6%
<i>auth</i>	34,046	8.7%	<b>16,381</b>	<b>34.3%</b>	29,148	14.0%
<i>cur</i>	6,362	1.6%	300	0.6%	6,318	3.0%
<i>comp</i>	6,734	1.7%	5,720	12.0%	4,362	2.1%
<i>auto</i>	<b>316,979</b>	<b>80.9%</b>	<b>18,344</b>	<b>38.4%</b>	<b>157,632</b>	<b>75.6%</b>
<i>obs</i>	1,826	0.5%	1,234	2.6%	3,550	1.7%
<b>sum</b>	<b>391,926 (60%)</b>		<b>47,805 (8%)</b>		<b>208,585 (32%)</b>	

>80% of all additions are *auto* annotations

Changes (8%) mainly *auth* and *auto*

Deletions (32%) mainly *auto* annotations

## ➤ Instabilities for *auth* and *auto*

# Provenance Changes

- How many annotations changed **from** one provenance type **to** another?

<i>from / to</i>	<i>exp</i>	<i>auth</i>	<i>cur</i>	<i>comp</i>	<i>auto</i>	<i>obs</i>	<i>Sum</i>	
<i>exp</i>	896	413	11	1,259	2,966	3	<b>5,548</b>	<b>13%</b>
<i>auth</i>	1592	798	73	1,038	11,901	23	<b>15,425</b>	<b>35%</b>
<i>cur</i>	21	27	0	16	182	0	<b>246</b>	<b>1%</b>
<i>comp</i>	1,280	1,206	26	0	3,101	0	<b>5,613</b>	<b>13%</b>
<i>auto</i>	3,311	10,169	228	2,329	0	116	<b>16,153</b>	<b>37%</b>
<i>obs</i>	79	391	9	12	725	0	<b>1,216</b>	<b>3%</b>
<i>Sum</i>	<b>7,179</b>	<b>13,004</b>	<b>347</b>	<b>4,654</b>	<b>18,875</b>	<b>142</b>	<b>44,201</b>	
	16%	<b>29%</b>	1%	11%	<b>43%</b>	0%		

- EC changes predominantly between *auth* and *auto* (in both directions)
- No obvious trend for the rest
- Due to vast amount of *auto* annotations

# Assessing Annotation Quality

## Seen so far

- Most annotation changes are additions of new annotations
- Also many deletions and changes
- Instabilities for *auto* and *auth* annotations

## Idea

- Assessing the quality of annotations based on their history and occurred changes (stability)

## Aim

- Filtering annotations w.r.t. different quality criteria



# Stability Measures

- Existence stability  $a_{\text{age}}$  age of annotation (in #versions)  
 $a_{\text{present}}$  presence within  $a_{\text{age}}$

$$stab_{\text{exis}}(a) = a_{\text{present}} / a_{\text{age}}$$

- Quality stability  $a_{\text{changed}}$  # provenance changes  
 $a_{\text{unchanged}}$  # unchanged provenance

$$stab_{\text{qual}}(a) = a_{\text{unchanged}} / (a_{\text{unchanged}} + a_{\text{changed}})$$

- Combined stability

$$stab_{\text{comb}}(a) = \min ( stab_{\text{qual}}(a), stab_{\text{exis}}(a) )$$

$v_0$	$v_1$	$v_2$	$v_3$	$v_4$	$a_{\text{age}}$	$stab_{\text{exis}}$	$stab_{\text{qual}}$	$stab_{\text{comb}}$
$q_1$	$q_1$	$q_1$	$q_1$	$(i_1, c_1, q_1)$	5	$5/5=1$	$4/(4+0)=1$	1
$q_1$	/	/	$q_1$	$(i_2, c_2, q_1)$	5	$3/5=0.6$	$2/(2+0)=1$	0.6
/	$q_2$	$q_2$	$q_1$	$(i_3, c_3, q_3)$	4	$4/4=1$	$1/(1+2)=0.33$	0.33

# Evaluation Scenario

- All currently available annotations for human proteins within the last three years



Swiss-Prot v<sub>47</sub>-v<sub>56</sub>



- Quality / classification criteria:
  - Provenance *exp, auth, cur, comp, auto*
  - Age
    - old* (> 1.5 years)
    - middle* (0.5 to 1.5 years)
    - novel* (< 0.5 years)
  - Stability
    - stable* (stab ≥ 0.9)
    - unstable* (stab < 0.9)

# Age Analysis

	<i>old</i>	<i>middle</i>	<i>novel</i>
<i>exp</i>	10,114	3,387	8,808
<i>auth</i>	23,445	4,253	2,492
<i>cur</i>	253	67	157
<i>comp</i>	1,913	477	942
<i>auto</i>	92,474	63,954	49,909
<b>sum</b>	<b>128,199</b>	72,138	62,308

49% old annotations

Novel, manual annotations are predominantly *exp*

- Novel and middle aged annotations rarely classified as unstable (results see paper)

# Stability Analysis

	stab <sub>exis</sub>	stab <sub>qual</sub>	stab <sub>comb</sub>
<b>exp</b>	21,659	20,486	20,122
	650	1,880	2,187
<b>auth</b>	29,157	26,862	26,067
	1,033	3,116	4,123
<b>cur</b>	462	399	393
	15	78	84
<b>comp</b>	3,127	2,409	2,317
	205	1,078	1,015
<b>auto</b>	183,127	201,968	179,490
	23,210	4,369	26,847
<b>sum</b>	<b>237,532</b>	<b>252,124</b>	<b>228,389</b>
	<b>25,113</b>	<b>10,521</b>	<b>34,256</b>

stable
unstable

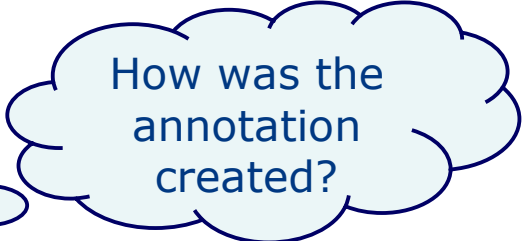
High share of temporal absence

13% unstable, mainly *auto* (80%) and some *auth* (12%)

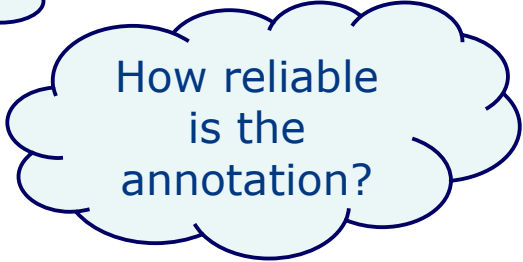
# Use – Quality Analysis

Protein ID	GO Concept ID	Provenance	Age in Years	stab <sub>exis</sub>	stab <sub>qual</sub>	stab <sub>comb</sub>
ENSP00000344151	GO:0015808 (L-alanine transport)	<i>exp</i>	3	1	1	1
ENSP00000230480	GO:0005615 (extracellular space)	<i>auto</i>	2.5	1	<b>0.462</b>	0.462
ENSP00000352999	GO:0006915 (apoptosis)	<i>exp</i>	3	<b>0.824</b>	1	0.824

- Different criteria to assess the quality of annotations w.r.t. provenance, stability, ...
- Filter less/more reliable annotations (e.g. stable, old, manually assigned)



How was the annotation created?



How reliable is the annotation?

# Use – Quality Analysis


- Stable, old, manually assigned:

In Ensembl about  
30,000 (11%)



How many  
high-quality  
annotations are  
available in a  
source?

- Criteria selection is highly dependent on application!
- Annotation instability is not necessarily a negative aspect
- Alternative interpretation



Which  
annotations fit  
best for my  
application?

novel or unstable annotations (in Ensembl 96,000;  
37%) are of special research interest / significant  
biological findings

# Conclusion and Future Work

- Generic approach to estimate the quality of ontology-based annotations by taking their evolution history into account
- Evaluation in two large life sciences sources
  - Instabilities for *auth* or *auto* annotations
- Different quality criteria: provenance, age, and stability to classify annotations
- Investigate other quality aspects
- Explore the impact of unstable annotations on dependent applications