# Bio4j: bigger, faster, leaner

Pablo Pareja-Tobes, Alexey Alekhin, Evdokim Kovach,
Marina Manrique, Eduardo Pareja, Raquel Tobes and
Eduardo Pareja-Tobes
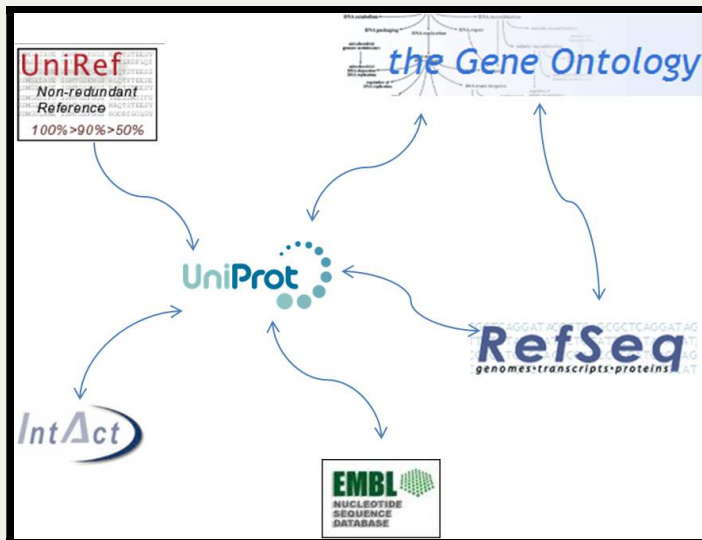
08.04.2014



ohnosequences!
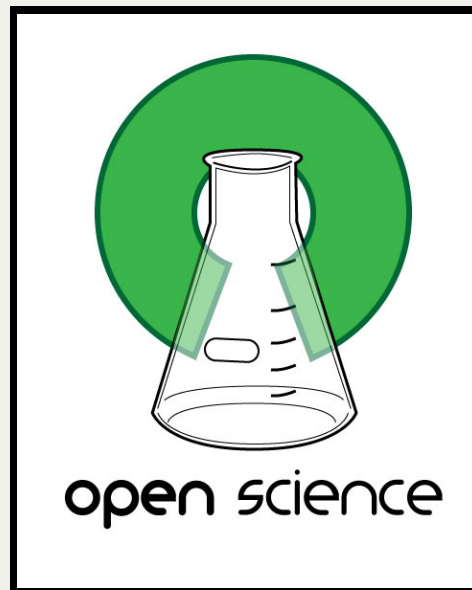
Era7 bioinformatics

# Introduction

# What is Bio4j?

**Bio4j** is a bioinformatics *graph*-based data platform
**integrating** the most representative **open data sources**
around **protein information**

# Data sources



- *UniProt KB* (SwissProt + Trembl)
- *Gene Ontology* (GO)
- *UniRef* (50,90,100)
- *RefSeq*
- *NCBI taxonomy*
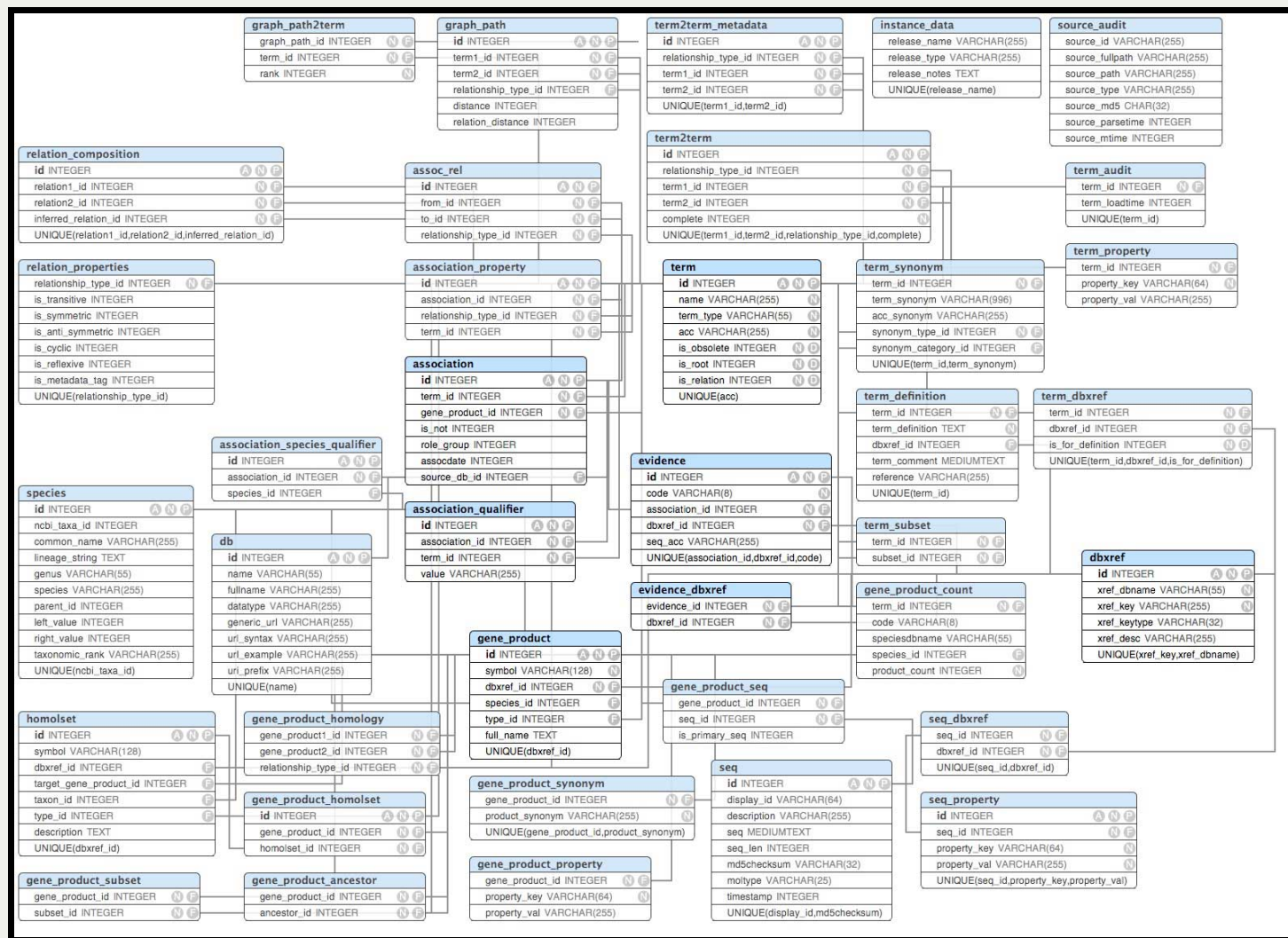- *Expasy Enzyme DB*

# It's open!



- Code is under the **AGPLv3** license
- Only **Open Data** is integrated
- Implementation & release process is **100% public** and totally transparent

# Biology & Databases today

- Highly **interconnected overlapping** knowledge
- spread over **different data sources**
- maintained in the **Relational Databases**
  or sometimes even just as plain **CSV files**

That might be fine for simple scenarios
but as the **amount** and **diversity** of data grows,
**domain models** become *crazily complicated!*

*Doesn't look very compelling right?*

# Relational model

With relational paradigm the double implication

**Entity ⇔ Table**

doesn't go both ways, which implies

- **auxiliary tables**
- **artificial IDs**
- dealing with **raw tables**
  (in spite of entity-relationship diagrams)

**Integrating** new knowledge becomes **difficult**

# Biology ≠ Table

- **Life** in general and **biology** in particular are probably not 100% like a graph…
- but one thing is sure: they *are not a set of tables!*

# Why graph databases?

- Data is stored in a way that **semantically represents its own structure**
- Incorporating new data is easy $\Rightarrow$ it's **scalable**
- **Vertex-centric** *(local)* indices allow to overcome the supernode problem

# Why in the cloud?



## Data as a service

- Services interoperability
- Data distribution
- Backup and storage
- Scalability
- Cost-effectiveness

# Bio4j

=

Bio Data

+

Graph Databases

+

The Cloud

# Details about Bio4j

# How it all started

- Need for **massive access** to *Gene Ontology* annotations
- **BG7** bacterial genome annotation system
- Need for massive direct access to **protein information**

# More and more data!

- As *other* data sources were becoming a *bottleneck* they were integrated into Bio4j
- First it was Uniprot KB, then Uniref, …
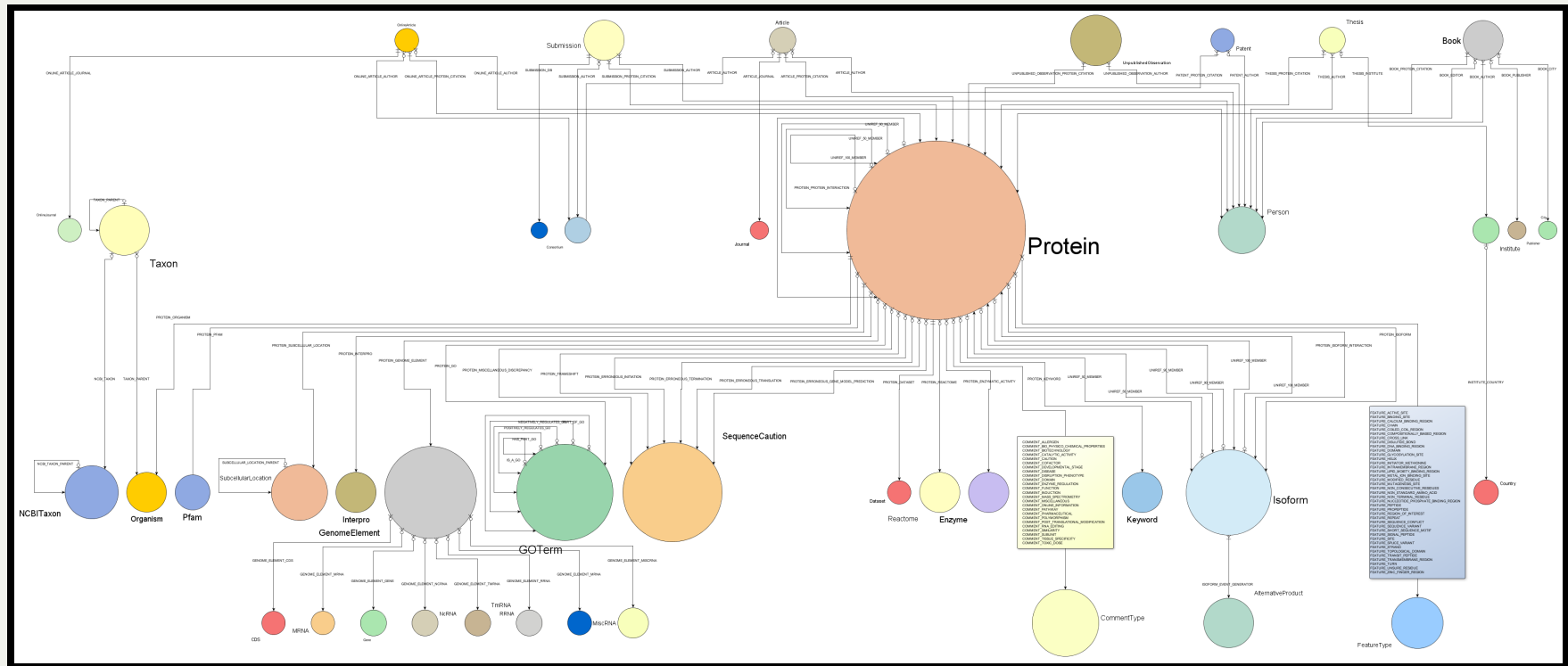- And **we didn't stop yet!**

# Different layers of Bio4j

1. Abstract **domain model** with precise typing
2. Universal Blueprints implementation
3. **Technology-specific** versions:
   - Neo4j
   - Titan (WIP)
   - OrientDB (planned)

*Different* **graph topologies** at the storage level,
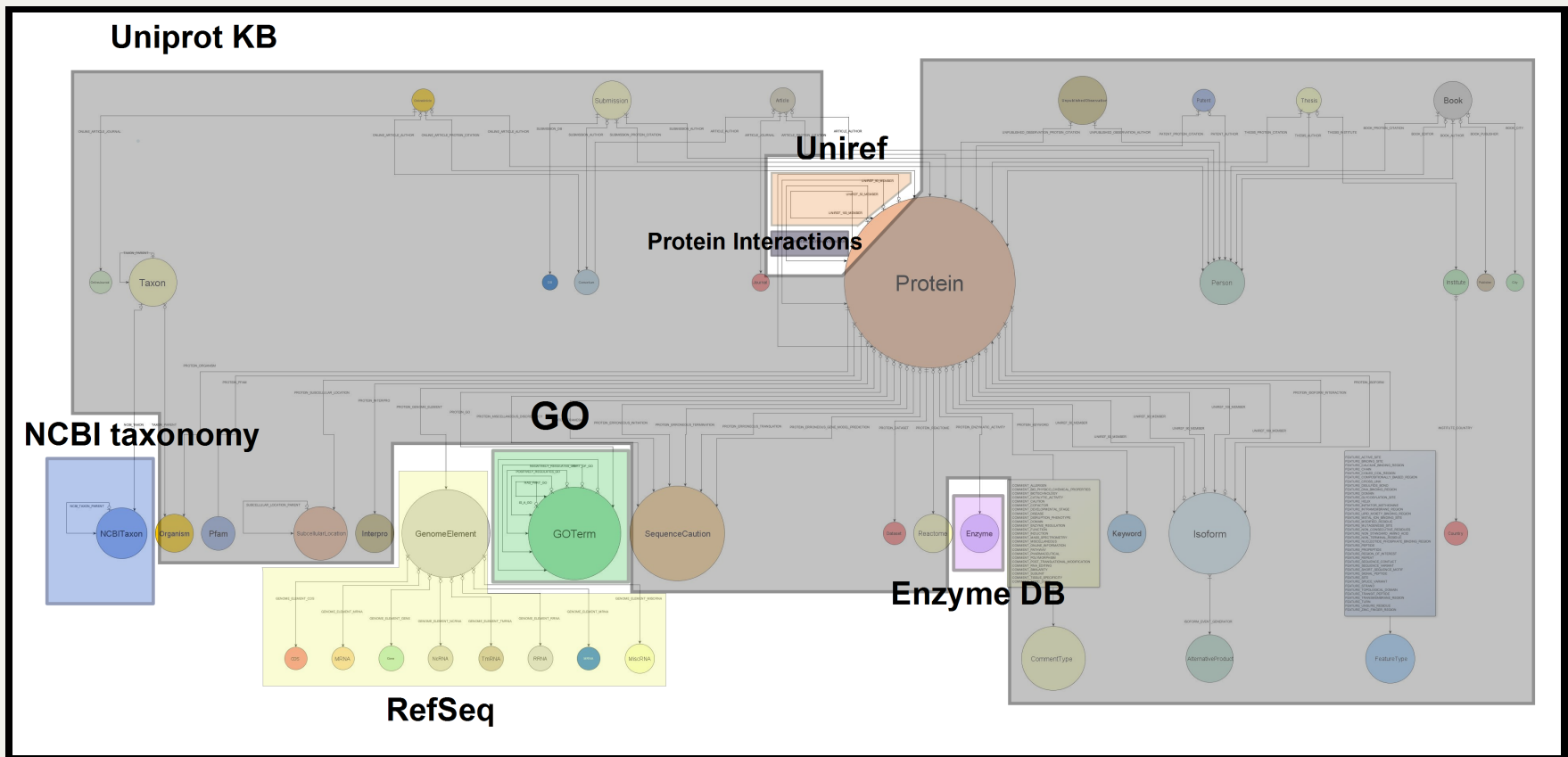*same* **domain model** in the client's code

# Bio4j domain model

- $10^9$ edges of **150 types**
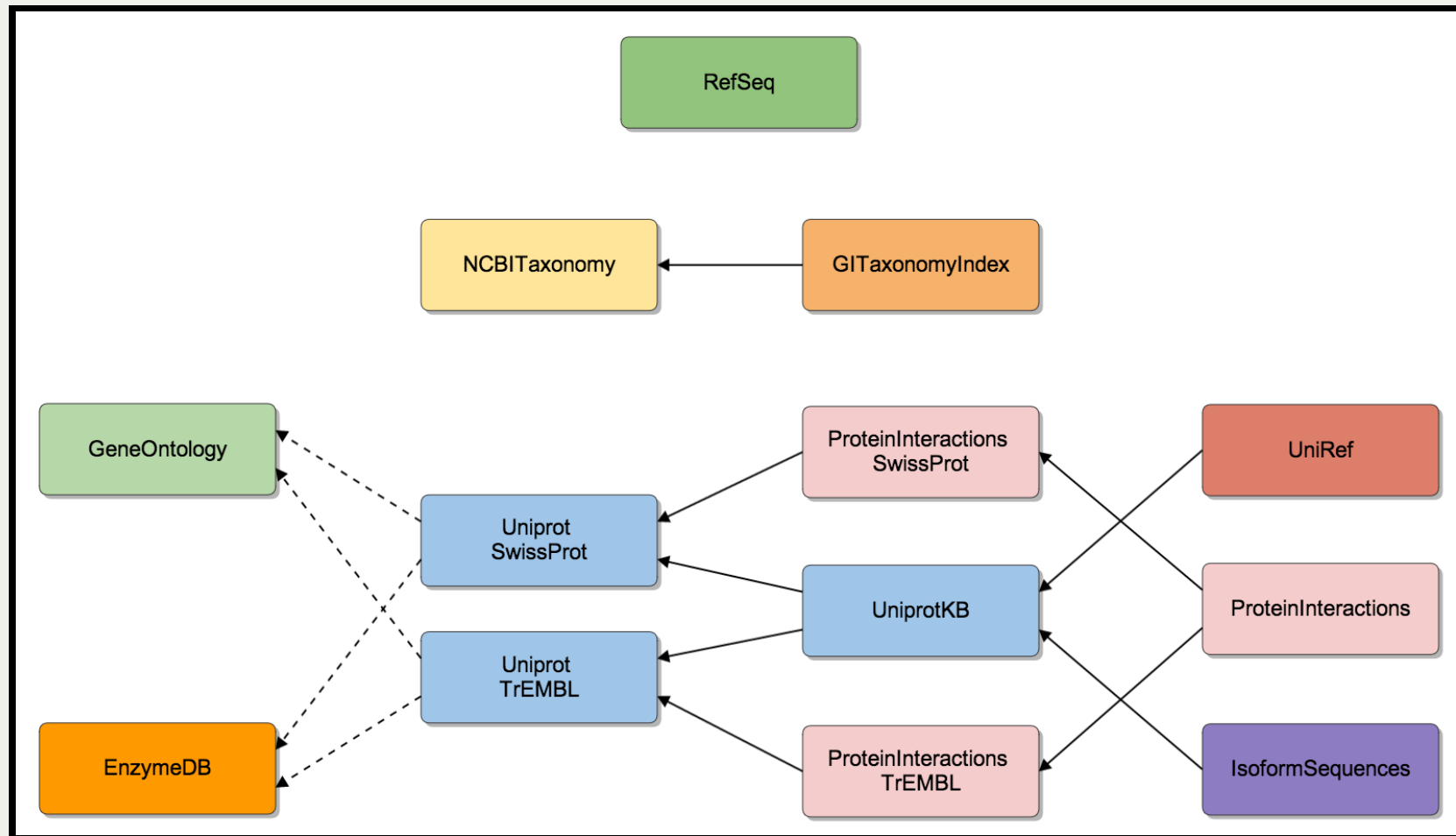- $2 \times 10^8$ nodes of **40 types**
- $6 \times 10^8$ properties

# Bio4j structure

The importing process is **modular** and **customizable** allowing you to import just the data you are interested in

# Bio4j module system

Statika helps to manage dependencies between modules and simplifies import and deployment in the cloud

# Under the hood

# How we use Bio4j in Era7

- **BG7** genome annotation
- **MG7** metagenomics analysis
- Comparative genomics, network analysis, genome assembly, …

# How others use Bio4j

### Ohio State University

- **Integration** and **analysis** of Chip-seq data
- **Modeling** genomic information and **gene regulatory networks**

### Berkeley Phylogenomics Group

- Graph database for *Big Data challenges* in **genomics** developed **on top of Bio4j**

# How we develop Bio4j

- Java + Scala source code
- Statika-based module system
- SBT for building sources
  and automated tests & release
- Git + Github: versioning, docs,
  collaboration, coordination

# Who's doing Bio4j

Ohnosequences!
Era7 bioinformatics R&D group

| | |
|---:|---|
| **Pablo Pareja** | project leader & main developer |
| **Eduardo Pareja-Tobes** | technology & architecture |
| **Raquel Tobes** | bio data integration |
| **Marina Manrique** | bio data integration |
| **Alexey Alekhin** | module system developer |
| **Evdokim Kovach** | developer |

# Contacts

- @bio4j **Twitter** for news
- bio4j **Github** org for the development process
- bio4j-user **Google group** for the user feedback
- bio4j **Linkedin**

# bio4j.com

# Thank you for attention!

*The source and the latest version of these slides can be found at*

`github.com/ohnosequences/IWBBIO-2014`