

---

# THE GENABEL PROJECT FOR STATISTICAL GENOMICS

**YURII AULCHENKO**

[YuriiA consulting (NL) | ICG SB RAS (RU) | CPHS UoE (UK) | @YuriiAulchenko]

**FOR THE GENABEL PROJECT CONTRIBUTORS**

[ @GenAproj | [www.GemABEL.org](http://www.GemABEL.org) ]

# OUTLINE

---

- Statistical genomics
  - A short history
  - Current state
  - Summary

Why are we different? Why do certain people get a disease?

What are the mechanisms underlying these differences?

**How genetic variation controls the phenotype?**



# STATISTICAL GENOMICS

**Feature 3**

**Sample 1**

	qt1	qt2	rs10	rs18	rs29	rs65	rs73
1	-0.58	4.46	0	0	0	1	0
2	0.80	0.32	0	0	NA	NA	0
3	-0.52	3.26	0	0	NA	1	0
4	-1.55	888.00	0	0	NA	NA	0
5	0.25	5.70	0	1	0	2	0
6	0.15	4.65	0	0	0	0	0
7	-0.56	4.64	1	1	1	2	0
8	NA	5.77	0	0	0	2	NA
9	-2.26	0.71	0	1	0	2	0
10	-1.32	3.26	0	0	0	2	0

# STATISTICAL GENOMICS

			Feature 3						
Sample 1		qt1	qt2		rs10	rs18	rs29	rs65	rs73
	1	-0.58	4.46	1	0	0	0	1	0
	2	0.80	6.32	2	0	0	NA	NA	0
	3	-0.52	3.26	3	0	0	NA	1	0
	4	-1.55	888.00	4	0	0	NA	NA	0
	5	0.25	5.70	5	0	1	0	2	0
	6	0.15	4.65	6	0	0	0	0	0
	7	-0.56	4.64	7	1	1	1	2	0
	8	NA	5.77	8	0	0	0	2	NA
	9	-2.26	0.71	9	0	1	0	2	0
	10	-1.32	3.26	10	0	0	0	2	0

**Traits/  
phenotypes**

?

**Genotypes**

# GENOME-WIDE ASSOCIATION SCANNING (GWAS)

**lm(qt1 ~ rs10)**

	qt1	qt2
1	-0.58	4.46
2	0.80	6.32
3	-0.52	3.26
4	-1.55	888.00
5	0.25	5.70
6	0.15	4.65
7	-0.56	4.64
8	NA	5.77
9	-2.26	0.71
10	-1.32	3.26

	rs10	rs18	rs29	rs65	rs73
1	0	0	0	1	0
2	0	0	NA	NA	0
3	0	0	NA	1	0
4	0	0	NA	NA	0
5	0	1	0	2	0
6	0	0	0	0	0
7	1	1	1	2	0
8	0	0	0	2	NA
9	0	1	0	2	0
10	0	0	0	2	0

**Traits/  
phenotypes**

?

**Genotypes**

# GENOME-WIDE ASSOCIATION SCANNING (GWAS)

	qt1	qt2
1	-0.58	4.46
2	0.80	6.32
3	-0.52	3.26
4	-1.55	888.00
5	0.25	5.70
6	0.15	4.65
7	-0.56	4.64
8	NA	5.77
9	-2.26	0.71
10	-1.32	3.26

**Traits/  
phenotypes**

	rs10	rs18	rs29	rs65	rs73
1	0	0	0	1	0
2	0	0	NA	NA	0
3	0	0	NA	1	0
4	0	0	NA	NA	0
5	0	1	0	2	0
6	0	0	0	0	0
7	1	1	1	2	0
8	0	0	0	2	NA
9	0	1	0	2	0
10	0	0	0	2	0

**Genotypes**



# GENOME-WIDE ASSOCIATION SCANNING (GWAS)

	qt1	qt2
1	-0.58	4.46
2	0.80	6.32
3	-0.52	3.26
4	-1.55	888.00
5	0.25	5.70
6	0.15	4.65
7	-0.56	4.64
8	NA	5.77
9	-2.26	0.71
10	-1.32	3.26

**Traits/  
phenotypes**

	rs10	rs18	rs29	rs65	rs73
1	0	0	0	1	0
2	0	0	NA	NA	0
3	0	0	NA	1	0
4	0	0	NA	NA	0
5	0	1	0	2	0
6	0	0	0	0	0
7	1	1	1	2	0
8	0	0	0	2	NA
9	0	1	0	2	0
10	0	0	0	2	0

**Genotypes**





# GENOME-WIDE ASSOCIATION SCANNING (GWAS)

	qt1	qt2
1	-0.58	4.46
2	0.80	6.32
3	-0.52	3.26
4	-1.55	888.00
5	0.25	5.70
6	0.15	4.65
7	-0.56	4.64
8	NA	5.77
9	-2.26	0.71
10	-1.32	3.26

**Traits/  
phenotypes**

	rs10	rs18	rs29	rs65	rs73
1	0	0	0	1	0
2	0	0	NA	NA	0
3	0	0	NA	1	0
4	0	0	NA	NA	0
5	0	1	0	2	0
6	0	0	0	0	0
7	1	1	1	2	0
8	0	0	0	2	NA
9	0	1	0	2	0
10	0	0	0	2	0

**Genotypes**



# GENOME-WIDE ASSOCIATION SCANNING (GWAS)

	qt1	qt2
1	-0.58	4.46
2	0.80	6.32
3	-0.52	3.26
4	-1.55	888.00
5	0.25	5.70
6	0.15	4.65
7	-0.56	4.64
8	NA	5.77
9	-2.26	0.71
10	-1.32	3.26

**Traits/  
phenotypes**

	rs10	rs18	rs29	rs65	rs73
1	0	0	0	1	0
2	0	0	NA	NA	0
3	0	0	NA	1	0
4	0	0	NA	NA	0
5	0	1	0	2	0
6	0	0	0	0	0
7	1	1	1	2	0
8	0	0	0	2	NA
9	0	1	0	2	0
10	0	0	0	2	0

**Genotypes**



# GENOME-WIDE ASSOCIATION SCANNING (GWAS)

**Few**

	qt1	qt2
1	-0.58	4.46
2	0.80	6.32
3	-0.52	3.26
4	-1.55	888.00
5	0.25	5.70
6	0.15	4.65
7	-0.56	4.64
8	NA	5.77
9	-2.26	0.71
10	-1.32	3.26

**100,000-40,000,000,000**

	rs10	rs18	rs29	rs65	rs73
1	0	0	0	1	0
2	0	0	NA	NA	0
3	0	0	NA	1	0
4	0	0	NA	NA	0
5	0	1	0	2	0
6	0	0	0	0	0
7	1	1	1	2	0
8	0	0	0	2	NA
9	0	1	0	2	0
10	0	0	0	2	0

**1,000-100,000**

**Traits/  
phenotypes**

**?**

**Genotypes**

# SCANNING THROUGH “OMICS” SPACE

**100-100,000**

	qt1	qt2
1	-0.58	4.46
2	0.80	6.32
3	-0.52	3.26
4	-1.55	888.00
5	0.25	5.70
6	0.15	4.65
7	-0.56	4.64
8	NA	5.77
9	-2.26	0.71
10	-1.32	3.26

**100,000-40,000,000,000**

	rs10	rs18	rs29	rs65	rs73
1	0	0	0	1	0
2	0	0	NA	NA	0
3	0	0	NA	1	0
4	0	0	NA	NA	0
5	0	1	0	2	0
6	0	0	0	0	0
7	1	1	1	2	0
8	0	0	0	2	NA
9	0	1	0	2	0
10	0	0	0	2	0

**1,000-100,000**

**Traits/  
phenotypes**

**?**

**Genotypes**

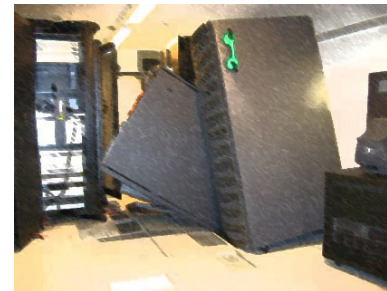
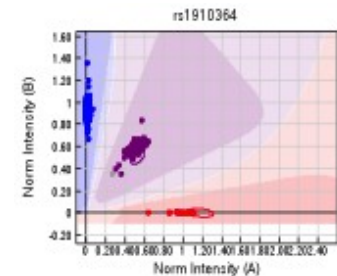
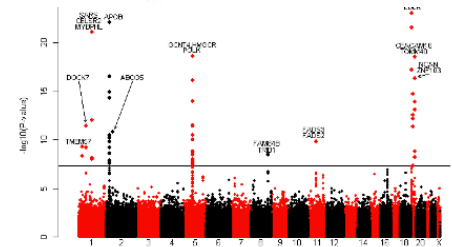
# STATISTICAL GENOMICS:

# WHAT IS SO SPECIAL?

- Rules governing genes & experimental design: analysis methodology and results visualization
- Technological inputs: data formats, quality control, analysis methods
- Analysis is computationally challenging (and IO demanding)



$$y = X\beta + Zu + \epsilon$$



# ANALYSIS SCENARIOS

---

- Classic GWAS scenario
  - One trait – one genetic marker at a time
  - Correlations between phenotypes – mixed models
- Emerging scenarios
  - One trait – multiple genetic markers
  - Multiple traits – single / multiple markers

# OUTLINE

---

- Statistical genomics
  - **A short history**
  - Current state
    - Summary

# A SHORT HISTORY

---

**Package**

**Paper**



**2006**

**GenA**

**2007**

**GenA**

**2008**

**2009**

**2010**

**2011**

**2012**

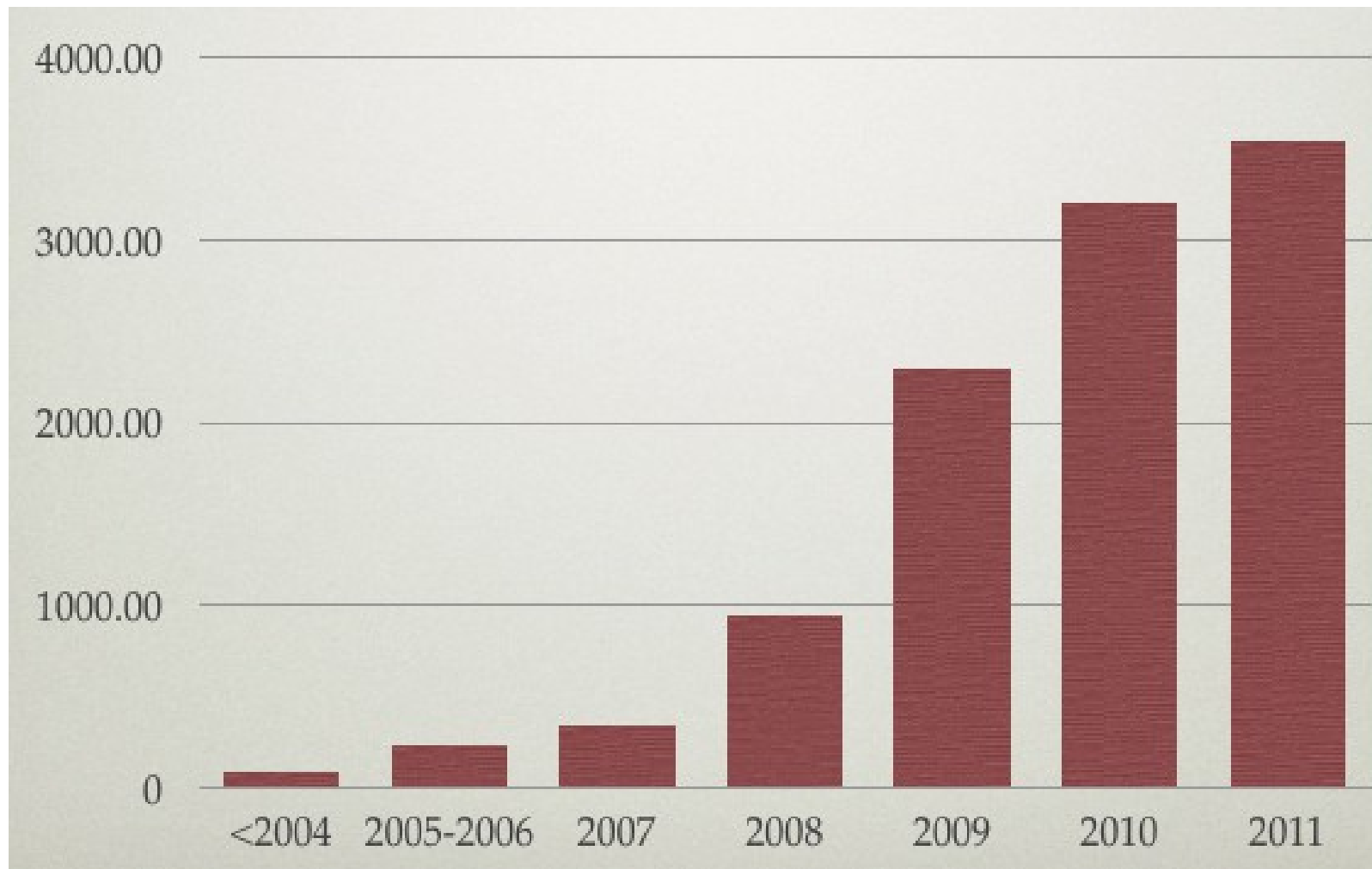
**2013...**

**GenABEL  
package**



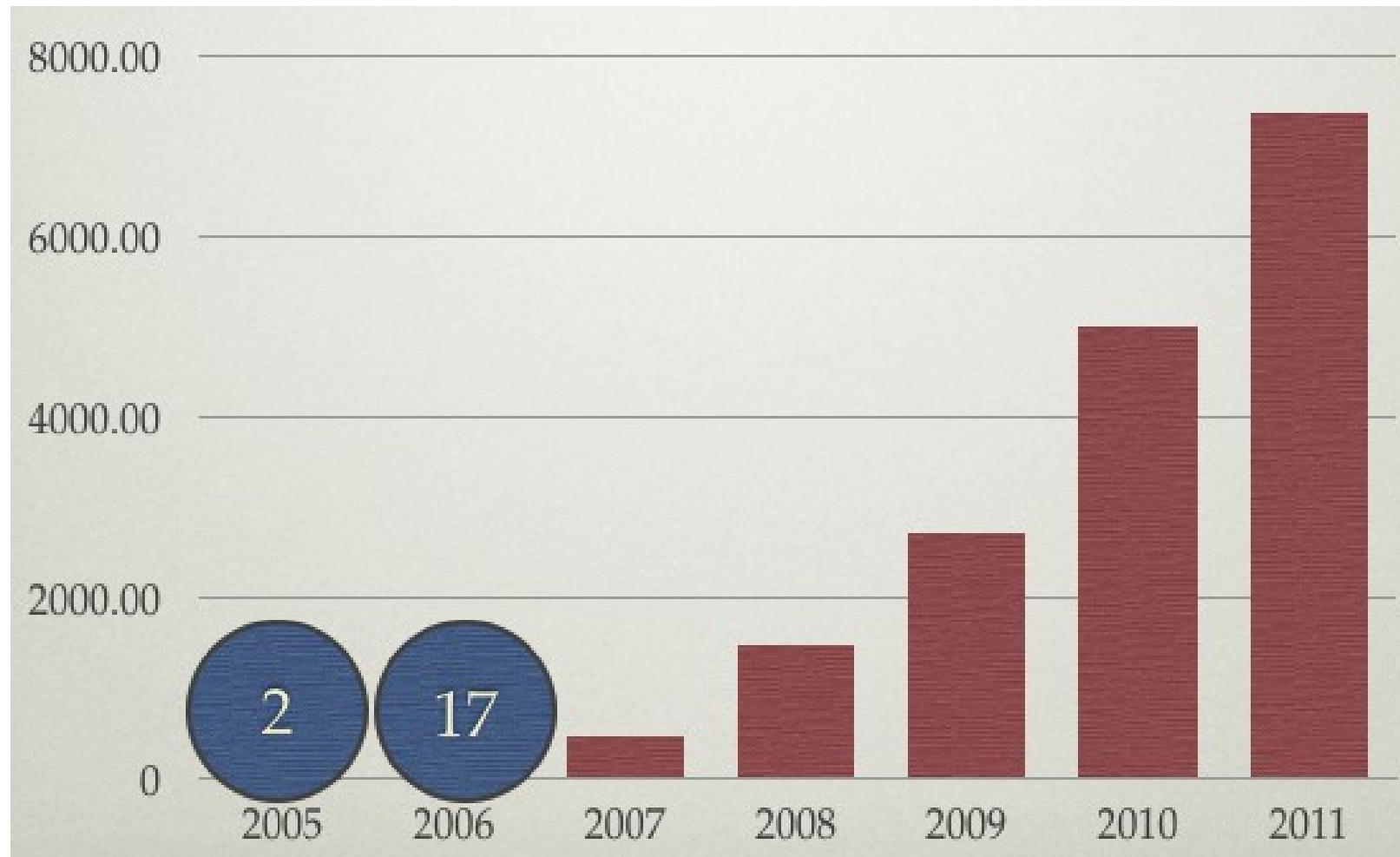
# # GWAS PUBLICATIONS

---



# # LOCI IDENTIFIED IN GWAS

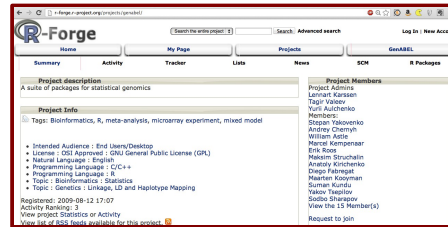
---



# A SHORT HISTORY

**Package**

**Paper**



Data Analysts Captivated by R's Power



**MetA**

**ProbA**

**ParallA**

**MixA**

**2006**

**2007**

**2008**

**2009**

**2010**

**2011**

**2012**

**2013...**

**GenA**

**GenA**

**Data**

**ParallA**

**ProbA**

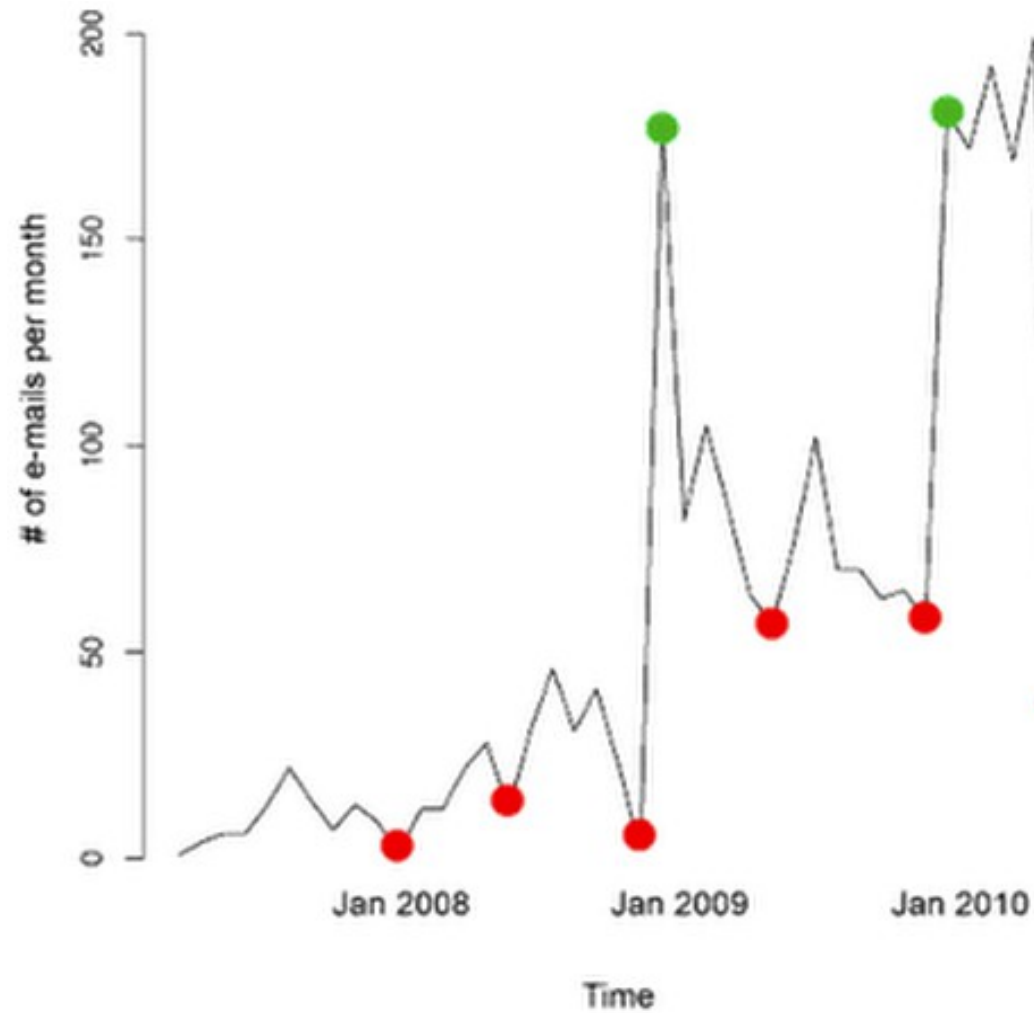
**GenABEL  
package**

**GenABEL  
suite**

**GenABEL**  
STATISTICAL GENOMICS

# TURNING POINT

---



# THE GENABEL PROJECT

---

**Mission:** to provide a framework for development of statistical genomics methodology

**Vision:** collaboration, transparency and free exchange of code, ideas, and data is a key to agile and robust methodology development

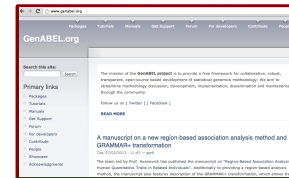
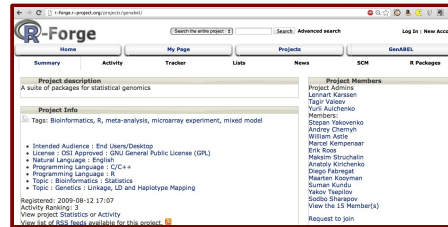
**Strategy:** community-based and driven methodology discussion, development, implementation, dissemination, maintenance, *and application*



# A SHORT HISTORY

**Package**

**Paper**



Data Analysts Captivated by R's Power



**1000 posts  
on forum**

[Genabel-commits] r1000 - pkg/ProbABEL/src

noreply@r-forge.r-project.org  
to genabel-commits

Author: Ickarsen  
Date: 2012-11-05 01:09:45 +0100 (Mon, 05 Nov 2012)  
New Revision: 1000

**Open-source  
tutorial**



nature.com • journal home • archive • issue • technical report • abstract

NATURE GENETICS | TECHNICAL REPORT

日本語要約

Rapid variance components-based method for whole-genome association analysis

**MetaA**

**ProbA**

**ParallA**

**MixA**

**PredictA**

**PredictA**

**GenA**

**2006**

**2007**

**2008**

**2009**

**2010**

**2011**

**2012**

**2013...**

**GenA**

**GenA**

**DataA**

**ParallA**

**ProbA**

**VariA**

**VariA**

**GenA**

**OmicA**

**GenABEL  
package**

**GenABEL  
suite**

**GenABEL  
project**

**GenABEL**  
STATISTICAL GENOMICS

# OUTLINE

---

- Statistical genomics
  - A short history
  - **Current state**
  - Summary

# INFRASTRUCTURE



Project description  
A suite of packages for statistical genomics

Project Members

GenABEL.org

Search this site:

Primary links

- ⊗ Packages
- ⊗ Tutorials
- ⊗ Manuals
- ⊗ Get Support
- ⊗ Forum
- ⊗ For developers
- ⊗ Contribute
- ⊗ People
- ⊗ Showcase
- ⊗ Acknowledgments

The mission of the **GenABEL project** is to provide a free framework for collaborative, robust, transparent, open-source based development of statistical genomics methodology. We aim to streamline methodology discussion, development, dissemination and maintenance; through the GenABEL project.

Follow us on [ Twitter ] [ Facebook ]

**READ MORE**

**A manuscript on a new region-based association analysis method and GRAMMAR+ transformation**

Tue, 07/02/2013 - 11:42 — yurii

The team led by Prof. Axenovich has published the manuscript on "Region-Based Association Analysis of Human Quantitative Traits in Related Individuals". Additionally to providing a region-based analysis method, the manuscript also features description of the GRAMMAR+ transformation, which allows treating

**www.GenABEL.org**

	TOPICS	POSTS	LAST POST
<b>GenABEL</b> Questions about GenABEL (aka *ABEL) suite of packages	242	759	by Nicola Pirastu Thu Jun 27, 2013 12:29 pm
<b>Problem Solving</b> Questions about GenABEL here.	16	6	by yurii Axenovich Sun Jun 23, 2013 2:22 pm
<b>Journal Club for Statistical Genomics</b> Statistical Genomics for dummies and advanced. Discussions, links, useful information.	9	30	by yurii Axenovich Sat Feb 23, 2013 9:44 pm

**forum.GenABEL.org**





# PROJECT IN NUMBERS

## Code of 9 packages

Language	# kLines of code
R	19
C++	19
C	17
Other	2
Rnw/Roxy	20

Estimated  
12 man-years  
\$1,500,000

## Documentation

Manuals	>200 pages
Tutorials	>250 pages
Videos	~10 min

## People

Developers	15 (5)
Forum	430 (71)

## Communications

Devel-list	>700 posts
Forum	>1000 posts

## Publications

Total	7 (4)
# citations	>700 (>500)

# WWW.GENABEL.ORG

---

- ~2,000 visits per month (~1,000 unique visitors)
- Major traffic from Europe (50%) and US (25%)
- ~50% of traffic generated by returning visitors



1 12,936

# GENABEL-PACKAGE

---

Genome-wide analysis of association between directly typed SNPs and quantitative, binary and time-till-event outcomes

## Highlights:

- Converters between different data formats
- Powerful QC organized around the `check.marker()` function
- A line of mixed-models based tools for correction for population stratification

Type of analysis	# functions
Data manipulations	~40
Quality control and descriptives	~10
Analysis	~30
Graphics & data presentation	~5
<b>Total</b>	<b>391</b>

# OTHER R-PACKAGES

---

## GWAS analyses

- *VariABEL* (5): tools for “environmental sensitivity” vGWAS
- *MixABEL* (12): advanced mixed models for GWAS

## Post-GWAS

- *MetABEL* (7): meta-analysis of GWAS results
- *PredictABEL* (111): assessment of (genetic) risk prediction models

## Support

- *DatABEL* (72): out-of-RAM large matrices storage and access
- *ParallABEL* (52): parallelization algorithms for GWAS



# NON-R PACKAGES

---

**ProbABEL:** GWAS of imputed data  
(quantitative, binary, time-till-event traits;  
regression and mixed models)

**Filevector:** C++ base for the DatABEL-package,  
facilitating out-of-core computations on large  
matrices

**OmicABEL:** rapid mixed-model based GWAS  
especially for multiple trait ("omics") analysis.

# OUTLINE

---

- Statistical genomics
  - A short history
  - Current state
  - **Summary**

# SUMMARY

---

- GenABEL is problem-centered project aiming towards agile development of statistical genomics methodology
- The GenABEL suite consist of 9 packages implementing close to 1,000 functions facilitating analyses of polymorphic genomes
- GenABEL suite is widely used for GWAS analyses of human, farm, pet animal, and plant data
- The project runs on enthusiasm and spare time of several people (and \$10 a month from “YuriiA consulting”)



# DIFFICULTIES WE FACE

---

## Core functionality

- The project would gain from re-design and added “core” functionality (e.g. regarding access to different data formats; parallelization)
- This gain is on the project, and not individual developer's level

## Coordination and communication

- Coordination takes time
- It may take a while before problems well-known for developers get through to the end-user (anyone willing to become our PR officer?)



# VACANT ROLES

---

- Project Coordinator
- Lead Developer(s)
- Public Relations Officer

# CURRENT SUPPORT FOR THE PROJECT

---

**Your logo could have  
been here**

# ACKNOWLEDGEMENTS

---

## People behind the GenABEL project

### Coordination

- Yurii Aulchenko (yurii [dot] aulchenko [at] gmail [dot] com; Twitter): project coordinator
- Lennart Karssen: admin for GenABEL@R-forge and GenABEL.org, GPG key ID: 0E1D39E3
- Anatoly Kirichenko (kianvi [at] mail [dot] ru): GenABEL.org web-site admin

### Open methodology discussion

Lars Ronnegard, Gulnara Svischeva, William Astle, Xia Shen, Yurii Aulchenko [in the future, will be automatically generated through 'most active' on the methodological list/forum]

### GenABEL-core coding team

*Current members:* Tagir Valeev, Yurii Aulchenko

*Former members:* Andrey Chernyh, Erik Roos, Maksim Struchalin, Marcel Kempenaar, Stepan Yakovenko

### Maintainers of GenABEL suite packages

A. Cecile J.W. Janssens, Maksim Struchalin, Suman Kundu, Unitsa Sangket, Yurii Aulchenko [in the future, will be automatically generated from 'maintainer' info of the packages]  
*See specific packages for the list of all authors [ Packages ]*

### Code development (commit-team)

Andrey Chernyh, Erik Roos, Lennart Karssen, Maarten Kooyman, Maksim Struchalin, Marcel Kempenaar, Stepan Yakovenko, William Astle, Yurii Aulchenko [in the future, will be automatically generated from GenABEL-commits]

### Code contributions and patches

Nicola Pirastu, Xia Shen, Toby Johnson, John Barnard, Nadezhda Belonogova, Han Chen, Vadim Pinchuk

### Bug reports

Karl Forner, Daniel Tallun, Aron Joon, Richard Ahn, Kati Kristiansson, Ross Fraser, Surakameth Mahasirimongkol, Lorna Lopez, Nadezhda Belonogova [in the future, will be automatically generated from the bug tracker info]

### User forum: moderators

Maria Gonik, Nicola Pirastu

### User forum: most active in support

Maria Gonik, Nicola Pirastu, Lennart Karssen, Yurii Aulchenko

### Note

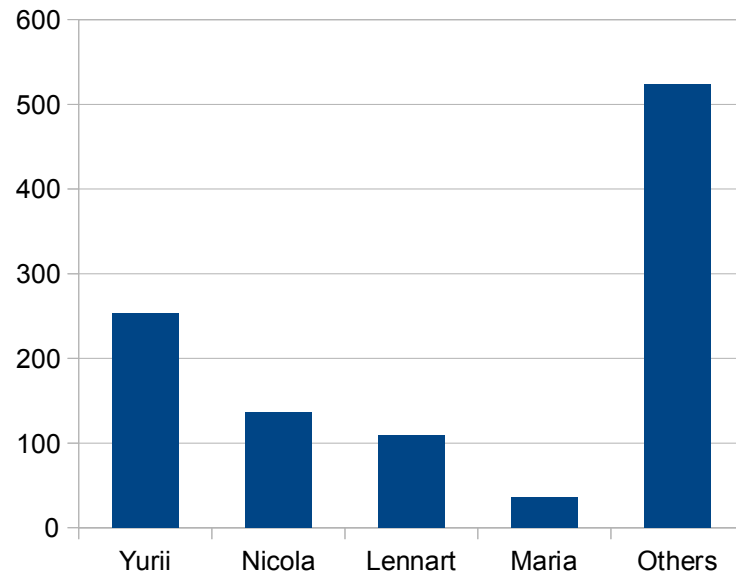
A large part of the GenABEL packages are R-packages, and we greatly appreciate work done by R team, CRAN, and R-forge



# KEY PEOPLE

Lennart Karssen, Nicola Pirastu, Maria Gonik

**Forum (431/70) members**



**Dev-list (45/12) members**

