# Coloc-stats

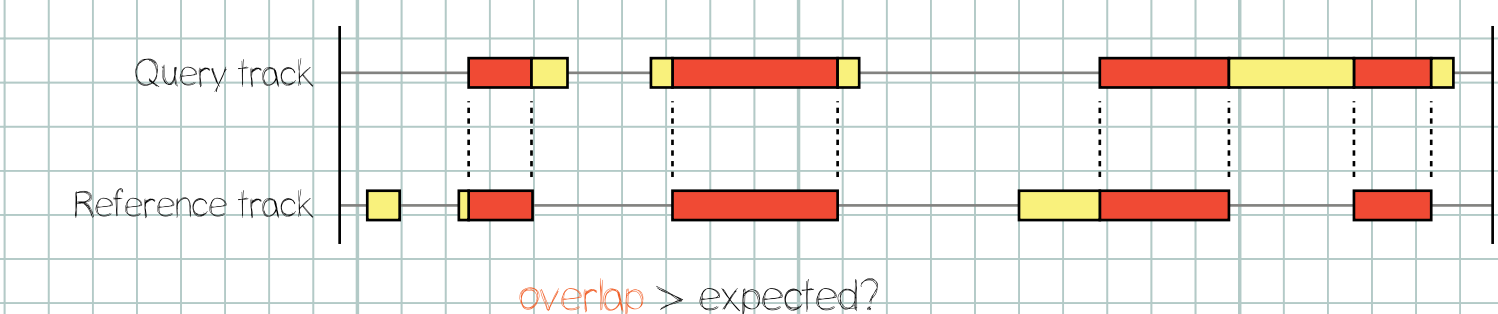## — a unified web interface to perform colocalization analysis of genomic features

Boris Simovski*, Chakravarthi Kanduri*, Sveinung Gundersen*, Dmytro Titov, Diana Domanska, Christoph Bock, Lara Bossini-Castillo,
Maria Chikina, Alexander Favorov, Ryan M Layer, Andrey A Mironov, Aaron R Quinlan, Nathan C Sheffield, Gosia Trynka, and Geir K Sandve
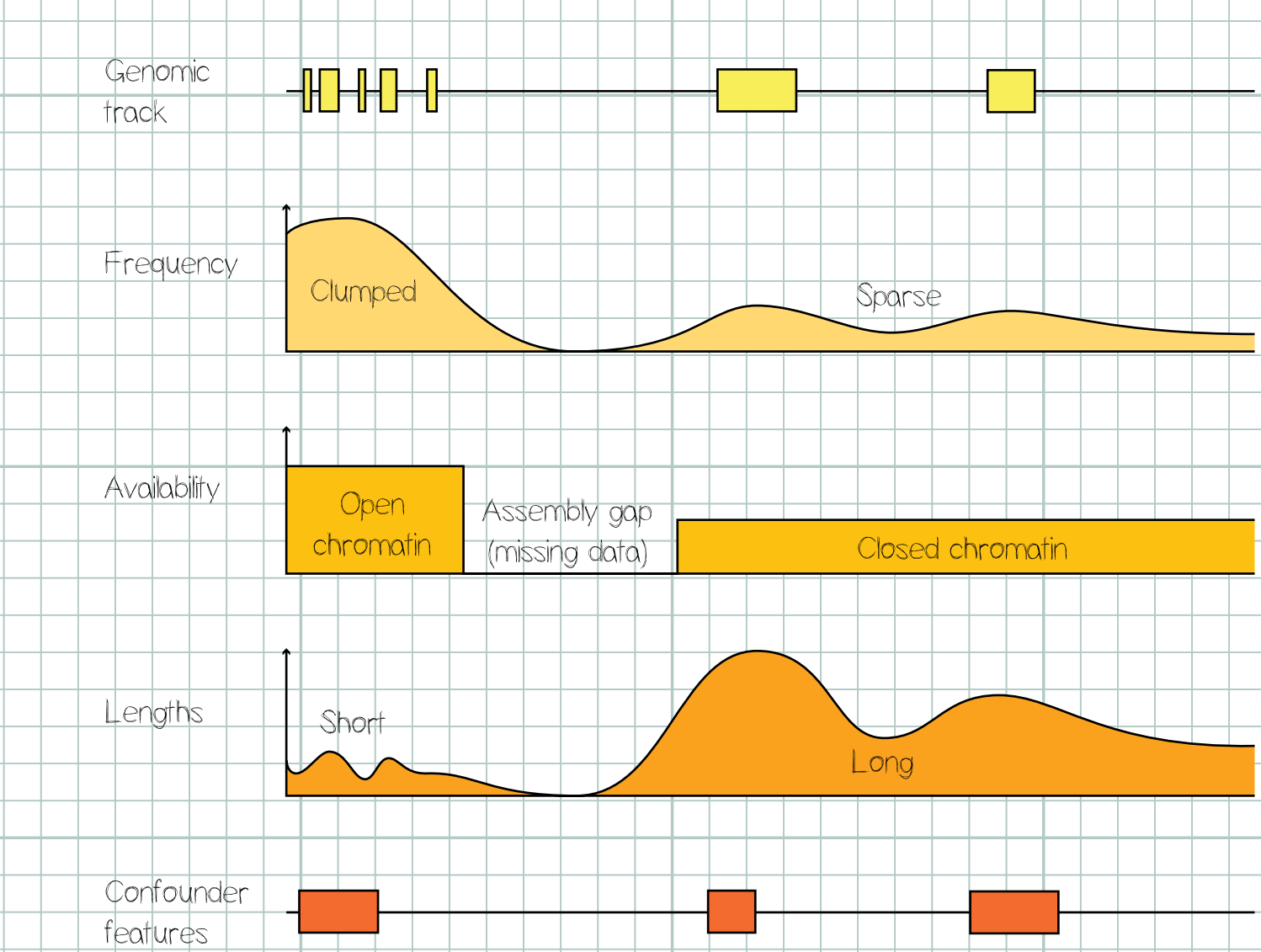
*These authors contributed equally

*2018 Galaxy Community Conference (GCCBOSC2018), June 27–28, Portland, Oregon, USA*

## The problem

*TFs, SNPs, transcripts, histone modifications, enhancers, repeats, DNA methylation, DNAseI HS, ...*

1. Functionally related GENOMIC FEATURES often tend to co-localize:

Query track
Reference track
overlap > expected?

2. The genomic features do not occur independently, but are known to follow various properties:

Genomic track
Frequency — Clumped / Sparse
Availability — Open chromatin / Assembly gap (missing data) / Closed chromatin
Lengths — Short / Long
Confounder features

3. Testing for significant co-localization assumes a model of the biological randomness (the NULL MODEL):

Monte Carlo / Exact
P-value / Test statistic / Mean of null dist.

4. Existing co-localization analysis tools assume quite different null models

and

5. The choice of null model greatly affects the conclusions[1]:

p=0.05
Null model A
Null model B
Null model C
1  2  3  4   -log10(p-value)

Thus, there is a high risk of FALSE POSITIVES!
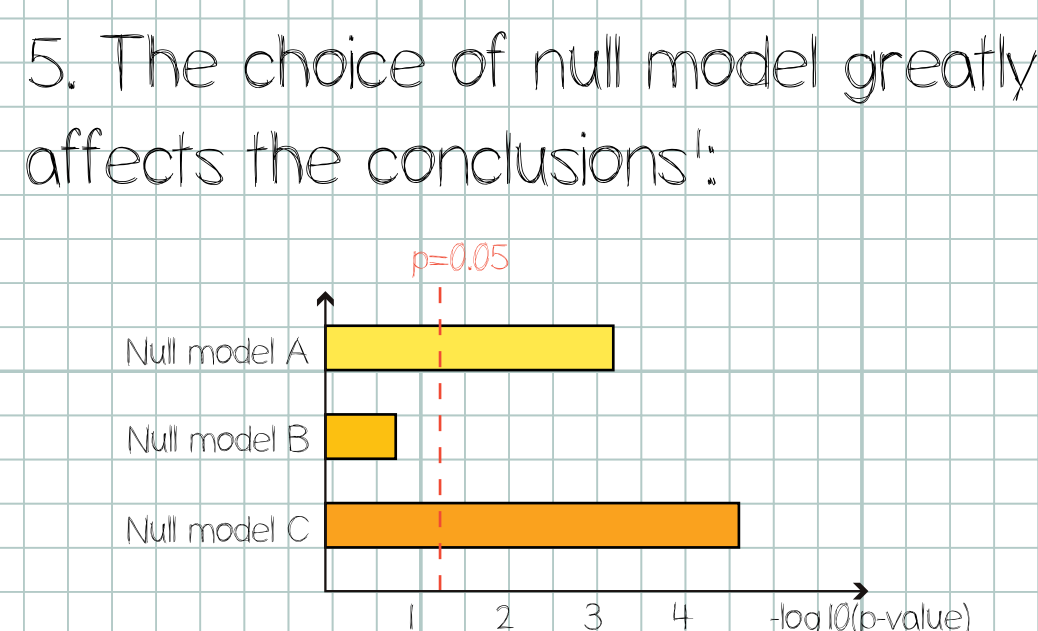
SO, WHAT SHOULD ONE DO?

Run multiple tests with varying methods, parameters, and null models → Compare the results → If they differ, reason about the assumptions, and be conservative
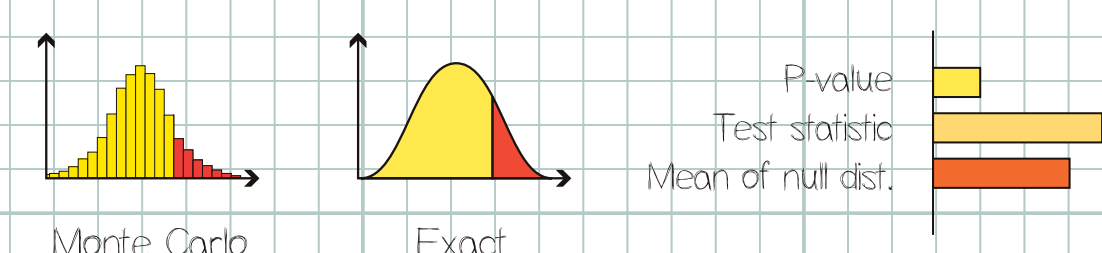
"BUT HOW TO EASILY RUN ALL THE TOOLS?"

## The goal

contain
One tool to rule them all:

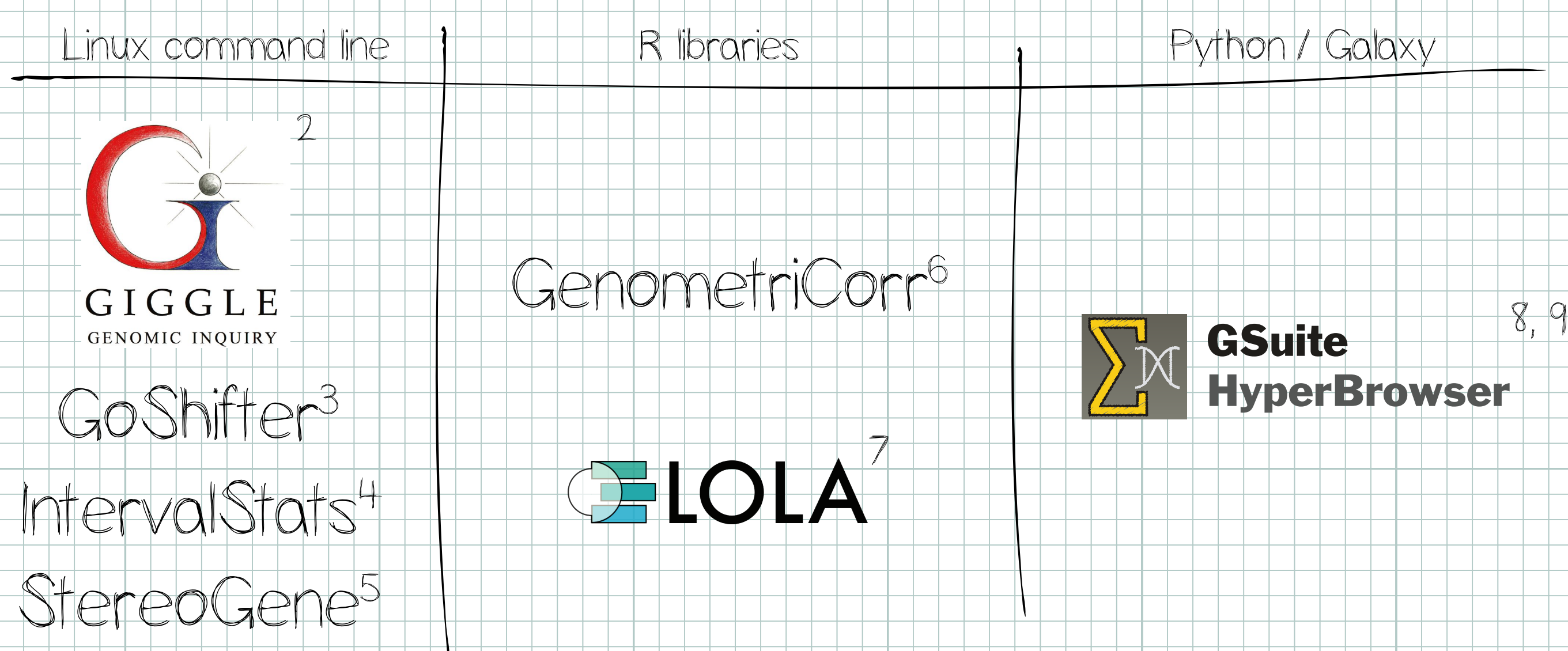### Coloc-stats

THE GRAPHICAL USER INTERFACE (GUI):

Explore and use multiple co-localization analysis tools through a single unified interface

Basic & Advanced mode

Extra feature: use tools created for two tracks also with larger collections of tracks!

Consciously select between several alternative modeling assumptions

THE RESULTS:

### Result page for coloc-stats analysis

Ranking of reference tracks

Query track tested for co-localization 29 – CATA1

| Reference track | Giggle | LOLA | StereoGene | IntervalStats | IntervalStats (v2) | GenometriCorr | HyperBrowser | Consensus rank |
|---|---|---|---|---|---|---|---|---|
| Brg1 | 2 | 1 | 22 | 19 | 14 | 2 | 1 | 4.2 |
| Pol2_UT-A_OpenChrom | 6 | 26 | 21 | 4 | 16 | 22 | 23 | 14.0 |
| c-Myc_Ifna5hUniPk_Yale | 7 | 7 | 18 | 9 | 8 | 4 | 5 | 7.4 |
| RBBP5_(A300-109A) | 16 | 15 | 17 | 15 | 17 | 15 | 13 | 15.4 |
| CCNT2 | 3 | 3 | 7 | 1 | 1 | 1 | 6 | 2.9 |
| c-Jun | 12 | 13 | 10 | 6 | 4 | 8 | 9 | 8.3 |

Arrive at a list of methods that are compatible with chosen modeling assumptions

Examine the robustness of conclusions by comparing results across several different methods/null models

+ Access full results of each method through detailed results pages

## The challenges

1. THE TOOLS ARE AVAILABLE IN DIFFERENT FRAMEWORKS:

| Linux command line | R libraries | Python / Galaxy |
|---|---|---|

GIGGLE GENOMIC INQUIRY [2]
GoShifter [3]
IntervalStats [4]
StereoGene [5]

GenometriCorr [6]
LOLA [7]

GSuite HyperBrowser [8, 9]

2. REQUIRES SINGLE USER INTERFACE WITH REPRODUCIBILITY AND ADVANCED LOGIC:
→ Galaxy framework: good
Galaxy tool XML: hmm...
Galaxy workflow engine: no good

3. TRYING TO COORDINATE SEVEN RESEARCH GROUPS...

4. TOO MANY CHEFS (DEVELOPERS) IN THE SAME CODE
→ Good modularization and interfaces needed to divide responsibilities between low-level (tools, parameters) and high-level (method logic) implementation
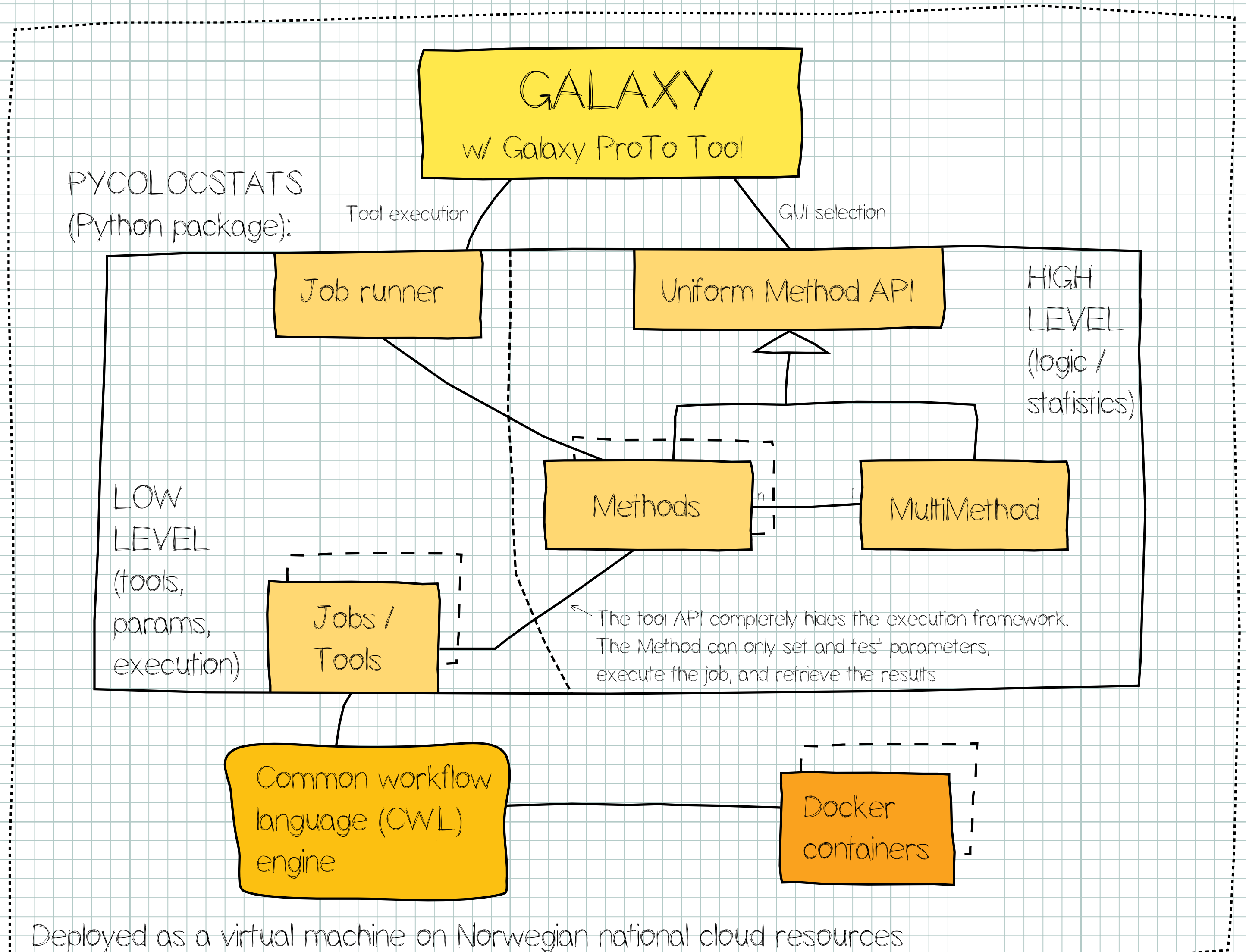
5. NUCLEIC ACIDS WEB SERVER ISSUE[10] DEADLINE:
= Too little time !!

## The solution

GALAXY w/ Galaxy ProTo Tool

PYCOLOCSTATS (Python package):
Tool execution
GUI selection

Job runner
Uniform Method API

HIGH LEVEL (logic / statistics)

LOW LEVEL (tools, params, execution)

Methods
MultiMethod

Jobs / Tools

The tool API completely hides the execution framework. The Method can only set and test parameters, execute the job, and retrieve the results

Common workflow language (CWL) engine

Docker containers

Deployed as a virtual machine on Norwegian national cloud resources

[1] Ferkingstad, E, Holden L. and Sandve, GK. "Monte carlo null models for genomic data." Statistical Science 30.1 (2015): 59-71.
[2] Layer, R.M, et al. "GIGGLE: a search engine for large-scale integrated genome alanysis." Nature methods (2018).
[3] Trynka, G, et al. "Disentangling the effects of colocalizing genomic annotations to functionally prioritize non-coding variants within complex-trait loci." The American Journal of Human Genetics 97.1 (2015): 139-152.
[4] Chikina, MD, and Troyanskaya, OG. "An effective statistical evaluation of ChIPseq dataset similarity." Bioinformatics 28.5 (2012): 607-613.
[5] Stavrovskaya, ED, et al. "StereoGene: rapid estimation of genome-wide correlation of continuous or interval feature data." Bioinformatics 33.20 (2017): 3158-3165.
[6] Favorov, A, et al. "Exploring massive, genome scale datasets with the GenometriCorr package." PLoS computational biology 8.5 (2012): e1002529.
[7] Sheffield, NC, and Bock, C. "LOLA: enrichment analysis for genomic region sets and regulatory elements in R and Bioconductor." Bioinformatics 32.4 (2015): 587-589.
[8] Sandve, GK, et al. "The Genomic HyperBrowser: inferential genomics at the sequence level." Genome biology 11.12 (2010): 1-12.
[9] Sandve, GK, et al. "The Genomic HyperBrowser: an analysis web server for genome-scale data." Nucleic acids research 41.W1 (2013): W133-W141.
[10] Simovski, Boris, et al. "Coloc-stats: a unified web interface to perform colocalization analysis of genomic features." Nucleic acids research 46.W1 (2018).

Background created by Freepik

K.G. JEBSEN Center for Medical Research — ÖAW AUSTRIAN ACADEMY OF SCIENCES — NIH National Institutes of Health — wellcome trust — NIH National Human Genome Research Institute — The Research Council of Norway

ELIXIR NORWAY — HORIZON 2020 — RFBR RUSSIAN FOUNDATION FOR BASIC RESEARCH — MRC Medical Research Council — NATIONAL CANCER INSTITUTE — UiO : University of Oslo