

# *Solanaceae*-centric Linked Data Platform for Plant Breeding & Genomics

Scientific Symposium: FAIR Data Science for Green Life Sciences

netherlands

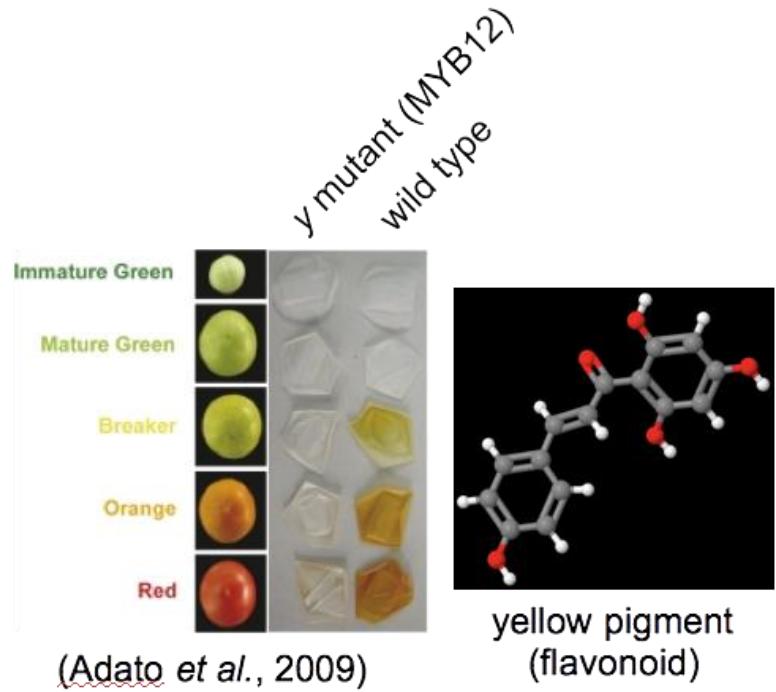
eScience center



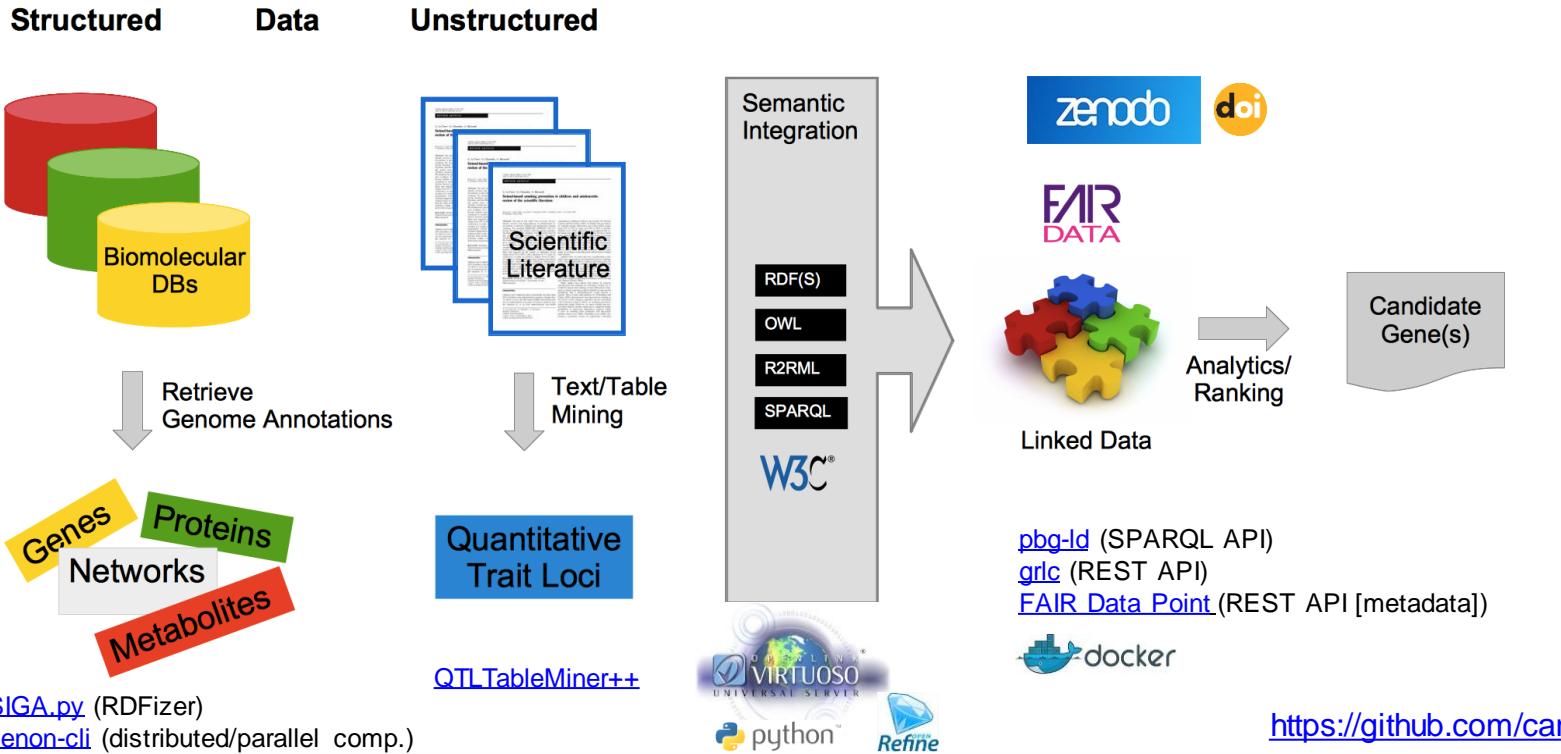
by SURF & NWO

Arnold Kuzniar  
12/12/2018

# Linking genes to traits



# Methodology



# Linked Data approach



Sir Tim Berners-Lee (W3C)

- Use (persistent) URIs/IRIs to identify things
- Use HTTP-resolvable URIs
- Use RDF graph model (triples)
- Link to other resources

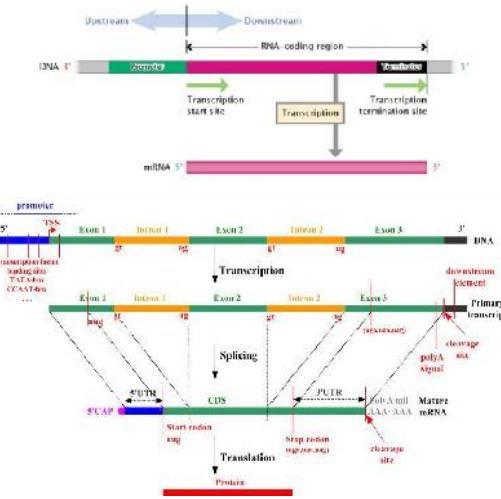
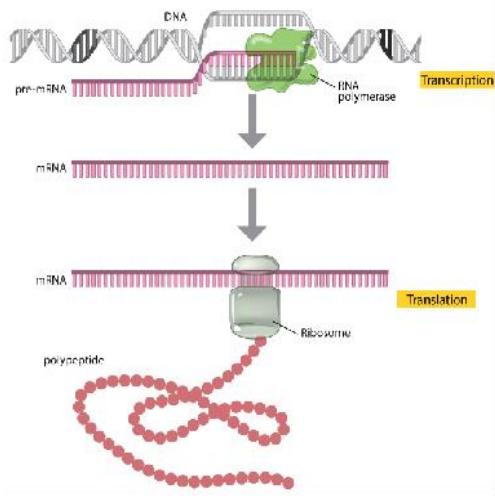


<URI> <URI> <URI | Literal>

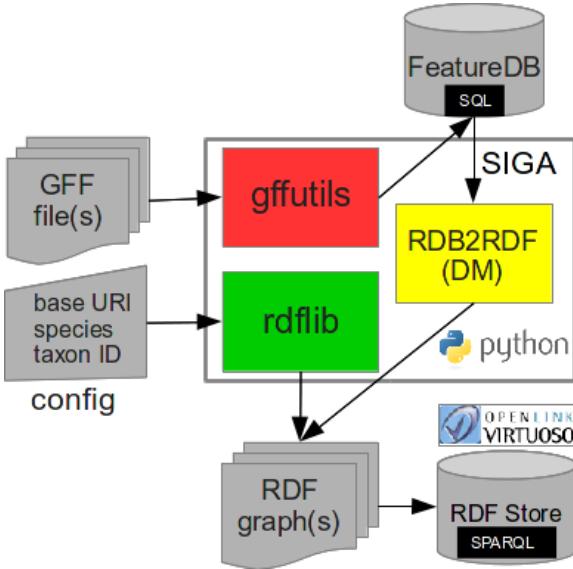
## RDF turtle format

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .  
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .  
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .  
@prefix obo: <http://purl.obolibrary.org/obo/> . Solyc00g005000.2 is a gene  
  
<http://solgenomics.net/genome/Solanum_lycopersicum/gene/Solyc00g005000.2> rdf:type obo:SO_0000704 ;  
rdfs:label "gene Solyc00g005000.2"^^xsd:string ;
```

# Genome annotations

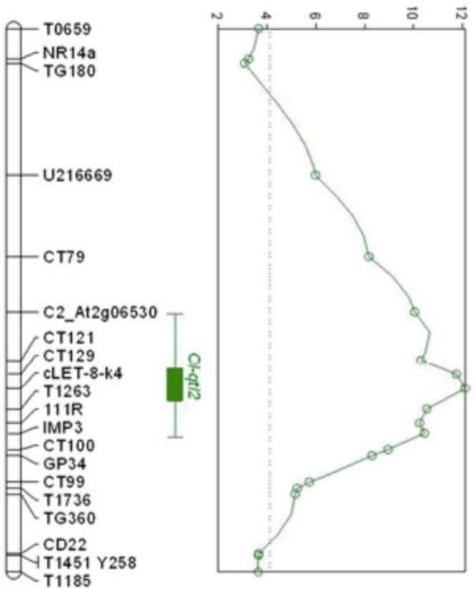


Use controlled vocabularies & ontologies  
(e.g. GenBank Feature Table, SO[FA], FALDO, DCMI)

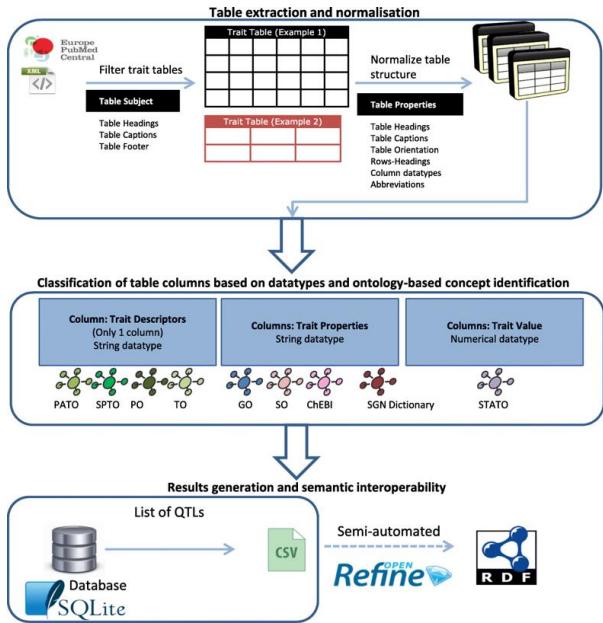
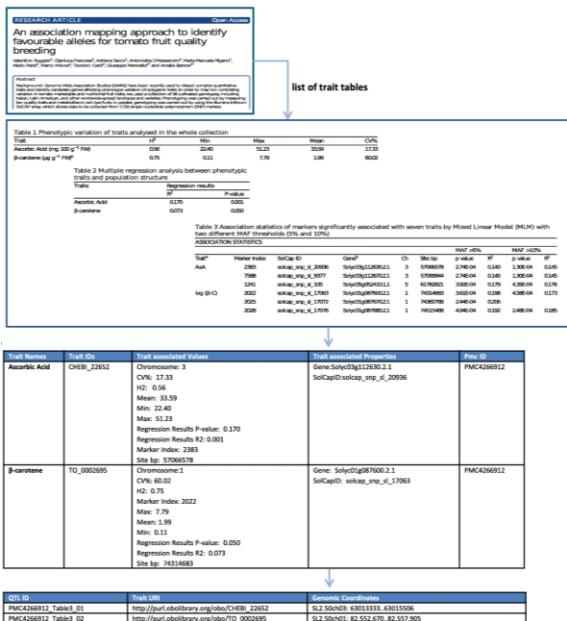


Software architecture

# Mining QTLs from literature



e.g., QTL mapped on chr12  
(Faino *et al.*, 2012)



QTL TableMiner++  
(Singh *et al.* 2018)

# Data sources



## Domestic tomato (*Solanum lycopersicum*)

- gene models (SGN & Ensembl) + proteins (UniProt)
- genetic markers (SGN & SolCAP)
- QTLs (Europe PMC)



## Wild tomato (*S. pennellii*)

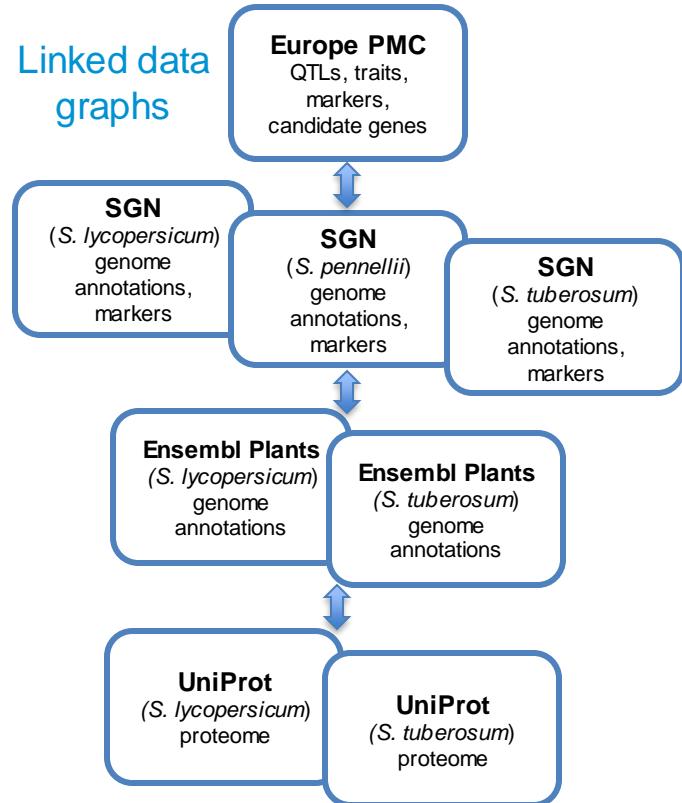
- gene models (SGN)
- genetic markers (SGN)
- QTLs (Europe PMC)



## Potato (*S. tuberosum phureja*)

- gene models (PGSC & Ensembl) + proteins (UniProt)
- genetic markers (SolCAP & DArT)
- QTLs (Europe PMC)

## Linked data graphs



# FAIR research output

Search or jump to... Pull requests Issues Marketplace Explore +

candYgene Software developed as part of the NLeSC & WUR project: Prediction of candidate genes for traits using interoperable genome annotations Netherlands https://www.esciencecenter.nl

Repositories 6 People 6 Teams 0 Projects 0 Settings

Pinned repositories

- = **pbg-Id** Linked Data Platform for Plant Breeding & Genomics
- = **QTM** Forked from PBR/QTM
- = **siga** Semantically Interoperable Genome Annotations

netherlands Science center About the directory

## Research Software Directory

Encouraging the re-use of research software

The Research Software Directory aims to promote the impact, the exchange and re-use of research software. Please use our tools! [Read more](#)

Start typing here to search for software

Sort by: Most mentions

Tags -

- Big data (0)
- GPU (0)
- High performance computing (0)
- Image processing (0)
- Inter-operability & linked data (10)

**QTLTableMiner++** QT pbg-Id pb

A tool to extract Quantitative Trait Locus data from tables in genetics literature and make this data more FAIR.

Provides easy access to integrated plant-specific data on genes associated with traits of interest.

Home About Articles Submission Guidelines Software Open Access

### QTLTableMiner++: semantic mining of QTL tables in scientific articles

Quantitative Trait Loci (QTL) mapping experiments are commonly described in scientific literature but the results presented in these articles are often difficult to reuse. QTLTableMiner++ (QTM) enables to make these results available in machine-readable formats.

DOI: [10.5281/zenodo.1215044](https://doi.org/10.5281/zenodo.1215044)

Quantitative Trait Loci in Solanaceae species

April 8, 2018 Kumar, Arnold Singh, Damodar Contact person(s) Vinkers, Richard

DOI: [10.5281/zenodo.1458169](https://doi.org/10.5281/zenodo.1458169)

Quantitative Trait Loci in Solanaceae species

December 11, 2018 Kumar, Arnold Singh, Damodar Contact person(s) Vinkers, Richard

DOI: [10.5281/zenodo.1193640](https://doi.org/10.5281/zenodo.1193640)

QTLTableMiner++: a tool for mining tables in scientific articles

December 4, 2017 Kumar, Arnold Singh, Damodar Contact person(s) Vinkers, Richard

DOI: [10.5281/zenodo.1076438](https://doi.org/10.5281/zenodo.1076438)

QTLTableMiner++: a tool for mining tables in scientific articles

December 5, 2017 Kumar, Arnold Singh, Damodar Contact person(s) Vinkers, Richard

DOI: [10.5281/zenodo.1083951](https://doi.org/10.5281/zenodo.1083951)

FAIR Data Point: a web service to provide machine-readable descriptions about datasets

December 5, 2017 Kumar, Arnold Singh, Damodar Contact person(s) Vinkers, Richard

DOI: [10.5281/zenodo.1083951](https://doi.org/10.5281/zenodo.1083951)

# Linked Data Platform

Text Search Entity Label Lookup Entity URI Lookup Featured | Demo\_Queries | About

Precision Search & Find

Label fruit shape

Describe

fruit shape [http://purl.obolibrary.org/obo/SP\\_0000038](http://purl.obolibrary.org/obo/SP_0000038)

fruit shape [http://purl.obolibrary.org/obo/QO\\_0002628](http://purl.obolibrary.org/obo/QO_0002628)

fruit shape circular [http://purl.obolibrary.org/obo/SP\\_0000050](http://purl.obolibrary.org/obo/SP_0000050)

fruit shape eccentric [http://purl.obolibrary.org/obo/SP\\_0000058](http://purl.obolibrary.org/obo/SP_0000058)

fruit shape elongated [http://purl.obolibrary.org/obo/SP\\_0000060](http://purl.obolibrary.org/obo/SP_0000060)

fruit shape index [http://purl.obolibrary.org/obo/SP\\_0000070](http://purl.obolibrary.org/obo/SP_0000070)

fruit shape index external 1 [http://purl.obolibrary.org/obo/SP\\_0000040](http://purl.obolibrary.org/obo/SP_0000040)

fruit shape index external 2 [http://purl.obolibrary.org/obo/SP\\_0000046](http://purl.obolibrary.org/obo/SP_0000046)

fruit shape rectangular [http://purl.obolibrary.org/obo/SP\\_0000068](http://purl.obolibrary.org/obo/SP_0000068)

fruit shape triangle [http://purl.obolibrary.org/obo/SP\\_0000049](http://purl.obolibrary.org/obo/SP_0000049)

fruit shape triangle 10% [http://purl.obolibrary.org/obo/SP\\_0000148](http://purl.obolibrary.org/obo/SP_0000148)

fruit shape triangle 20% [http://purl.obolibrary.org/obo/SP\\_0000149](http://purl.obolibrary.org/obo/SP_0000149)

fruit shape triangle 30% [http://purl.obolibrary.org/obo/SP\\_0000150](http://purl.obolibrary.org/obo/SP_0000150)

fruit shape triangle 5% [http://purl.obolibrary.org/obo/SP\\_0000147](http://purl.obolibrary.org/obo/SP_0000147)

OpenLink Virtuoso version 7.0.0.1000 (Build 1000) - 64 bit - 2018-07-10 14:45:00 UTC  
Windows Server 2012 R2 Standard Edition (62 GB total memory)

OPENLINK SOFTWARE

About: **QTL:4321030\_4\_14** Goto Sponge NotDistinct Perma  
An Entity of Type : [http://purl.obolibrary.org/obo/SO\\_0000771](http://purl.obolibrary.org/obo/SO_0000771), within Data Space

Type: QTL Command: Start New Facet Go

Attributes Values

type QTL

label QTL:4321030\_4\_14

dct:identifier QTL:4321030\_4\_14

dct:isReferencedBy <http://identifiers.org/pmc/PMC4321030>

location chromosome 11:33041803-51521123

correlated\_with variation gene165\_0-I23  
variation gene210\_0-I23  
variation gene220\_0-I23  
variation gene221\_0-I23

correlated\_with condition fruit shape  
fruit shape

overlaps gene Solyc11g038340.1  
gene Solyc11g038350.1  
gene Solyc11g039350.1  
gene Solyc11g039360.1  
gene Solyc11g039370.1  
>more>

OPENLINK SOFTWARE

About: **gene Solyc11g038340.1** Goto Sponge NotDistinct Perma  
An Entity of Type : [http://purl.obolibrary.org/obo/SO\\_0001217](http://purl.obolibrary.org/obo/SO_0001217), within Data Space

Type: protein\_coding\_gene Command: Start New Facet Go

Alias: Solyc11g038340; Name: Solyc11g038340.1

Attributes Values

type protein\_coding\_gene

label gene Solyc11g038340.1

comment Alias: Solyc11g038340; Name: Solyc11g038340.1

sameAs gene Solyc11g038340.1

dct:identifier Solyc11g038340.1

location chromosome 11:45259363-45259605

has\_part prim\_transcript Solyc11g038340.1\_1

transcribed\_to prim\_transcript Solyc11g038340.1\_1

is\_overlaps\_of QTL:3464107\_4\_25  
QTL:4321030\_4\_14  
QTL:4321030\_4\_7

Trait



QTL



Gene

<https://github.com/candYgene/pbg-ls>

# Data access & analysis

Virtuoso SPARQL Query Editor

Default Data Set Name (Graph IRI)

Query Text

```
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX owl: <http://www.w3.org/2002/07/owl#>
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX obo: <http://purl.obolibrary.org/obo/>
PREFIX so: <http://purl.obolibrary.org/obo/so#>

SELECT
  str(?qtl_id) AS ?qtl_id
  str(?sgn_gene_id) AS ?sgn_gene_id
  str(?sgn_trans_id) AS ?sgn_trans_id
  str(?annot) AS ?annot
WHERE
  GRAPH <http://europemc.org/articles> {
    ?qtl a obo:SO_0000771 ;
      obo:RO_0003308 ?trait ;
      so:overlaps ?gene ;
      dcterms:identifier ?qtl_id .
    FILTER(?trait = obo:SP_0000366)
  }
  GRAPH <http://solgenomics.net/genome/Solanum_lycopersicum> {
    ?gene so:transcribedTo ?transcript ;
      dcterms:identifier ?sgn_gene_id .
    ?transcript rdfs:comment ?annot ;
      dcterms:identifier ?sgn_trans_id
  }
```

SPARQL API

local

[ View URL | Local host:8890/sparql/local/local/] [editURI]

Count genomic features

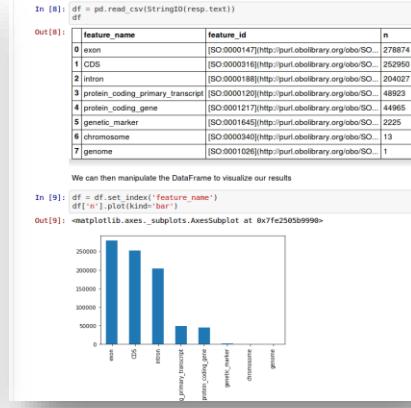
Get gene annotations from SGN given a gene ID

Get genes contained in a QTL

Get QTLs in an article

Get QTLs associated with a trait ID

REST API



Jupyter notebook

# Conclusions & future work

- Use/develop open-source, standards-compliant, well-documented software
- Linked (Open) Data principles ≈ FAIR Data principles (in practice)
- FAIR compliance not only for data but also for software/workflows
- Integrating & "semantifying" (semi-)structured data requires significant time investment & domain expertise
- Improve data access & processing workflow
- Extend *pbg-lid* with algorithms for candidate gene prioritization
- Prepare new software & data releases
- Deploy *pbg-lid* software stack in production environment (HPC Cloud)

# Acknowledgements

- **Gurnoor Singh (WUR)**
- **Richard Finkers (WUR)**
- **Christian Bachem (WUR)**
- **Richard Visser (WUR)**
- **Erik van Mulligen (EMC)**
- **Anand Gavai**
- **Carlos Martinez-Ortiz (NLeSC)**
- **Lars Ridder (NLeSC)**
- **e-infra team (SURFsara)**