



UPPSALA  
UNIVERSITET



Karolinska  
Institutet



e-Science for Cancer  
Prevention and Control

## Enabling Translational Medicine with e-Science

Ola Spjuth

[ola.spjuth@ki.se](mailto:ola.spjuth@ki.se)

[ola.spjuth@farmbio.uu.se](mailto:ola.spjuth@farmbio.uu.se)

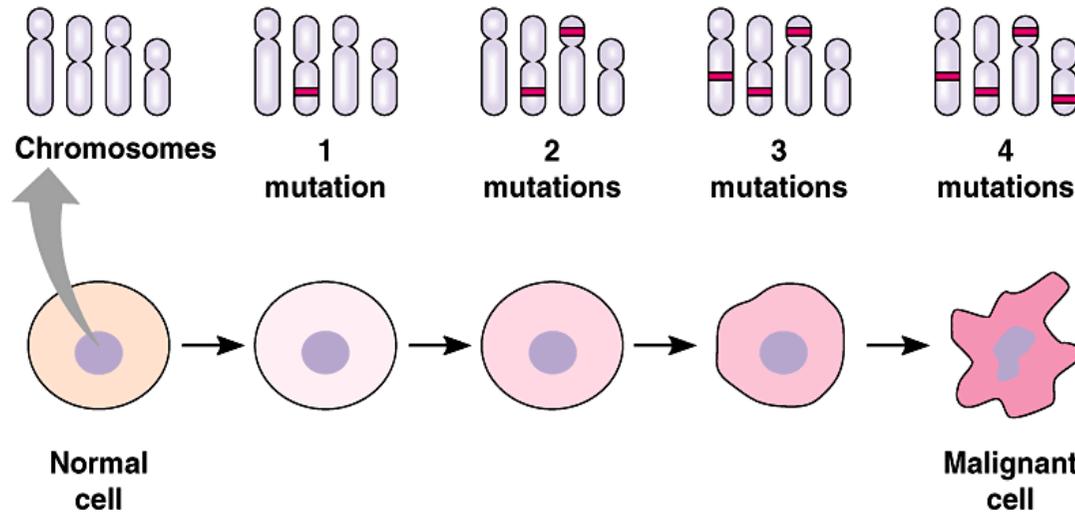
SciLifeLab

eScience  
THE E-SCIENCE COLLABORATION

SERC

# Cancer

Normal cells accumulate mutations

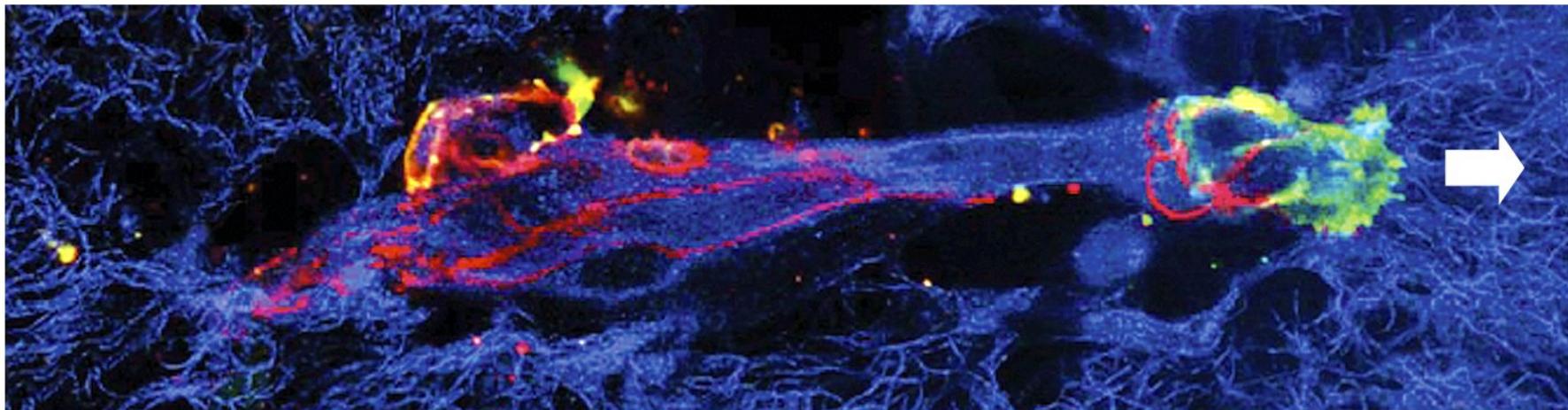


Copyright © 2003 Pearson Education, Inc., publishing as Benjamin Cummings.

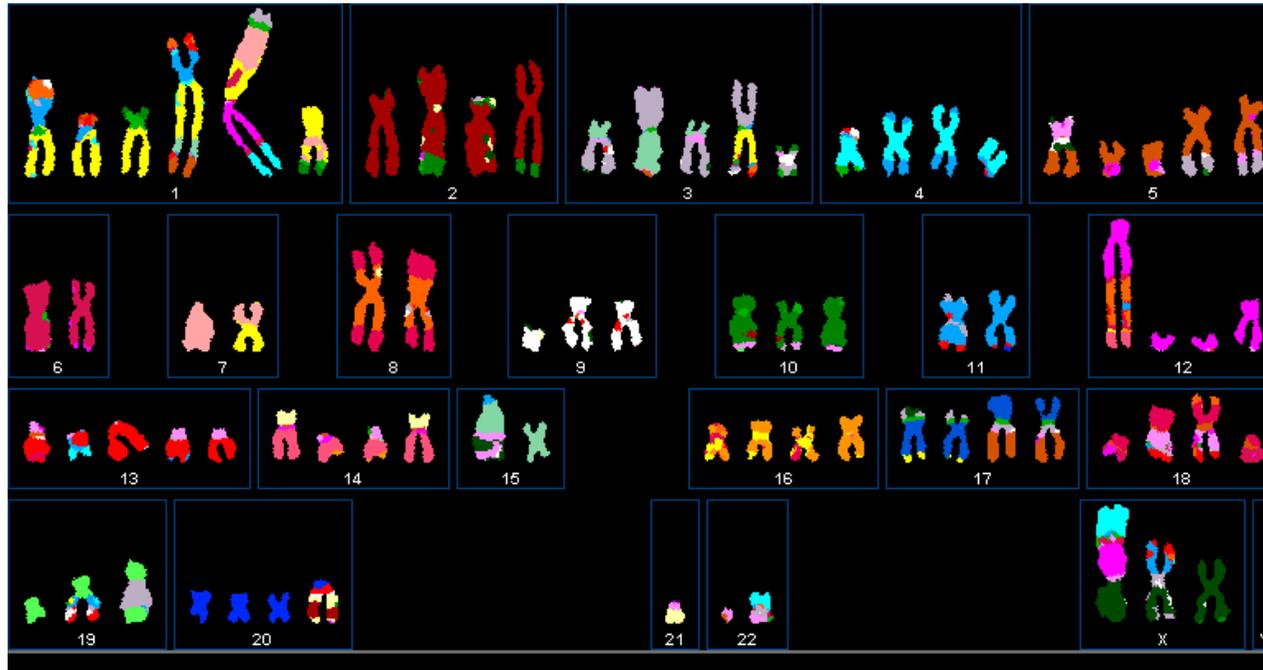
Cancer cells divide excessively and form tumors



Malignant cancer cells can invade other tissues



# Cancer cells have altered genomes



- Each cancer has its own genome  
→ Calls for individualized/stratified diagnostics and treatment



# New challenges: Data management and analysis

- Storage
- Analysis methods, pipelines
- Scaling
- Automation
- Data integration, security
- etc.



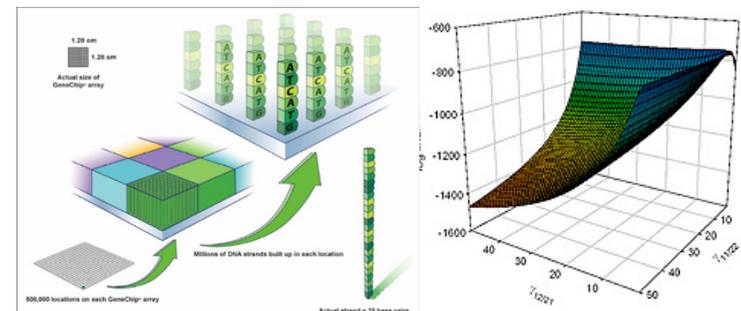
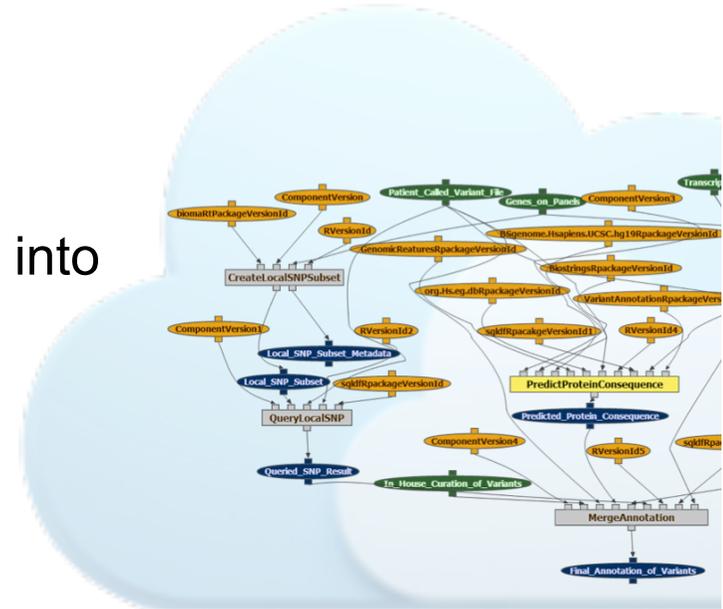
# Translational medicine

- Aim: **Turn Basic Research into Medicines and Treatments**
  - “from bench to bedside”
  - “from laboratory to clinic”
- Traditionally a slow process
  - May take 10 -20 years for original research to translate to routine medical practice



# e-Science in translational medicine

- e-Science solutions are desperately needed in order to translate high-throughput technologies into clinical settings (diagnostics, treatment,...)
  - e-infrastructure (computers, storage, networks, frameworks)
  - Methods, workflows/pipelines
  - Operations, experience and expertise
  
- e-Science can aid translational medical research
  - Simulation and prediction models
  - Large-scale collaborative, integrative research



## 2012-2014: e-Science for Cancer Prevention and Control

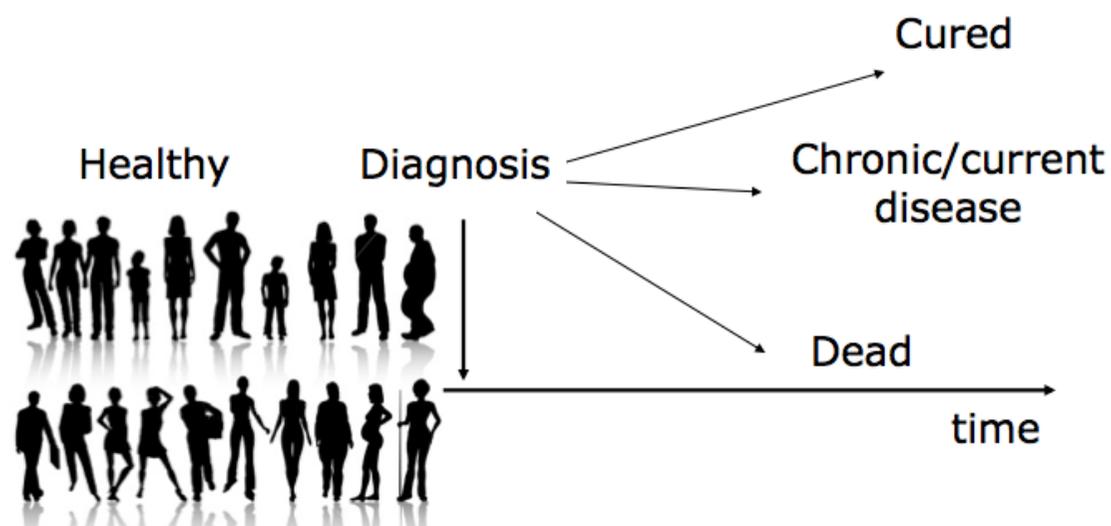
- a SeRC flagship project

Use statistical modeling and data integration in cancer research:

- Individualized prevention strategies
- Individualized treatments



PIs: Jun Palmgren & Jan-Eric Litton



# Cancer Risk Prediction Centre, CRiSP VR Linnaeus 2008-2018



**Breast cancer** is the most common cancer among women in Sweden with almost 8,000 new cases annually. In Sweden 1,500 women die from breast cancer yearly but there is a remarkable difference between

outcomes of localized vs advanced disease.



**Prostate cancer** is the most common cancer among men in Sweden today and yearly almost 10,000 new cases are diagnosed. Despite the old age of onset, the morbidity and mortality of this cancer is substantial with more

than 2,500 deaths annually.

We know that cancer mortality can be reduced if cases are detected and treated early, but there is a problem with over-diagnosis and over-treatment. **What if we instead could predict the risk for aggressive cancers?** Our research focuses on understanding cancer risk and how to design individualized prevention strategies.



PI: Per Hall

PI: Henrik Grönberg



## Personalized Cancer Prevention!

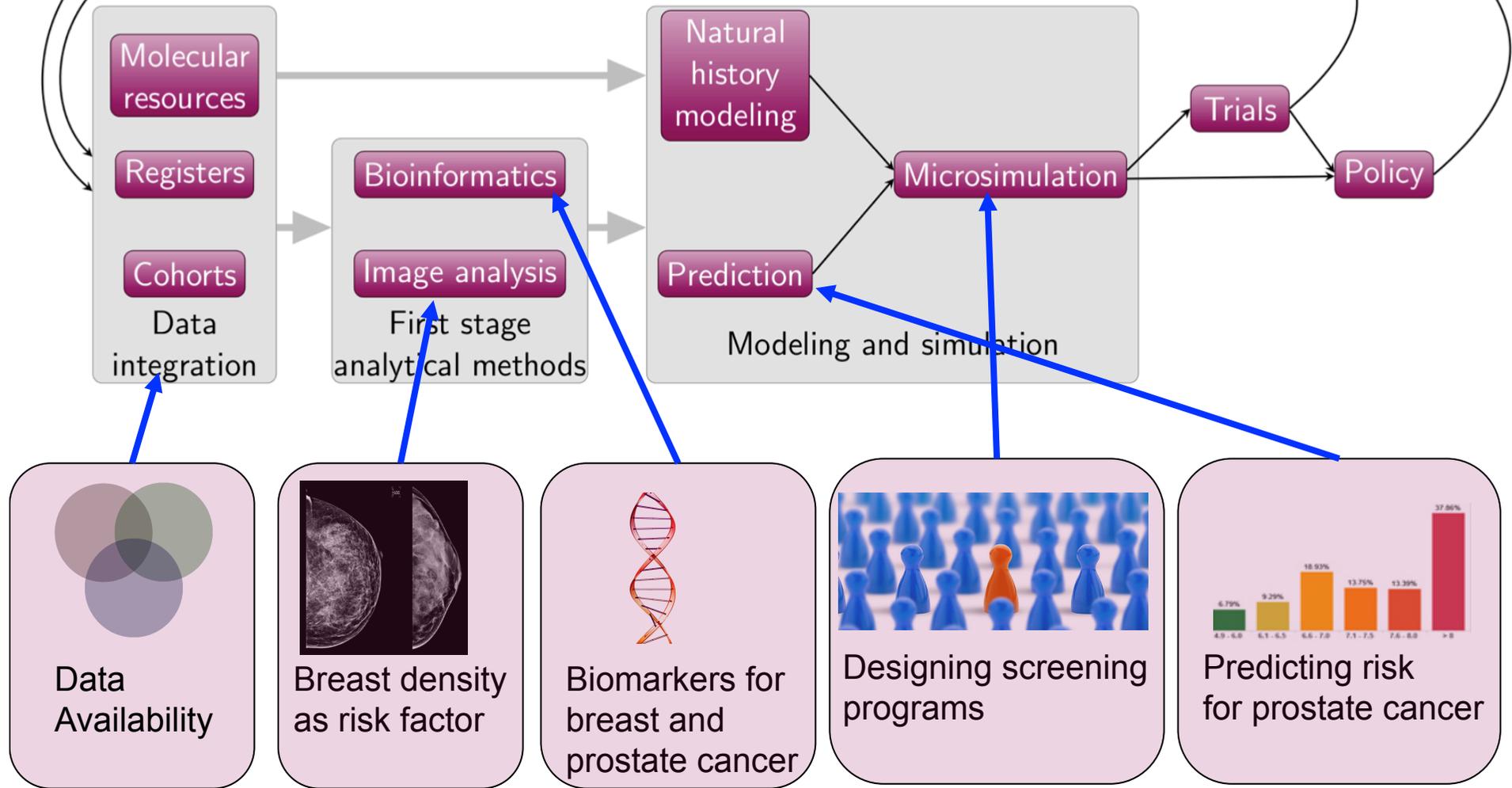
# Personalized screening; Why?

- Early detection is important.
- Efficiency of current practice for early detection of breast and prostate cancer is questioned!

We have:

- Organized mammographic screening
- Widespread opportunistic PSA testing
- Rates of detection of slow growing cancers increase
  - Beware of over-diagnosis, overtreatment, increasing cost and increase in side effects!
- Aggressive cancers and mortality do not decrease enough

# e-Science components in research on personalized screening

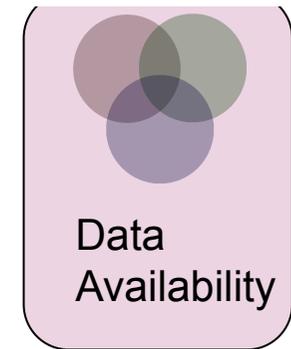


highlight 5 eCPC case studies

# Data Integration

## SAIL – Sample Availability System

- Overview of data – what is available
- Plan studies, investigate data available for subset of patients/samples across data archives such as biobanks



The screenshot shows the SAIL web application interface. The browser address bar displays <http://www.ebi.ac.uk/Tools/sail/>. The interface includes a navigation bar with tabs for "Welcome", "Summary", and "Report constructor". Below this, there are dropdown menus for "Study: [ANY]" and "Collection: [ANY]".

The main content area is divided into two sections. On the left, the "Parameter list" section shows a table of parameters. On the right, the "Report request" section shows predefined queries and a request builder.

Code	Name	Description	Filter	Records	V	E
GW_AFFY_100k	Affymetrix Genome-wide Hun	Affymetrix Genome-wide Hum		0	1	0
GW_AFFY_5	Affymetrix Genome-wide Hun	Affymetrix Genome-wide Hum		0	1	0
GW_AFFY_500k	Affymetrix Genome-wide Hun	Affymetrix Genome-wide Hum		3142	1	0
GW_AFFY_6	Affymetrix Genome-wide Hun	Affymetrix Genome-wide Hum		0	1	0
GW_AFFY	Affymetrix Genome-wide gen	Affymetrix Genome-wide gen		3142	1	0
AGE	Age	Age		231543	1	0
ALC	Alcohol	Alcohol		18963	1	1
ALCQ	Alcohol quantity	grams absolute ethanol / week		60415	1	0
ANTIHYPR	Antihypertensives	Antihypertensive treatment		206255	1	1
APOB	Apo B mg/l	Biochemistry Apolipoprotein B		2077	1	1
BMI	BMI	Body Mass Index, kg/m2		226985	1	0
BASO	Basophils (0.02-0.1)	Blood Basophils		69	1	0
BICEPS	Biceps mm	Thickness of a skinfold on the l		69	1	0
BYR	Birth Year	Birth Year		213457	1	0
BP	Blood pressure	Blood pressure (svstolic, diast		220949	2	0

The "Report request" panel on the right shows "Predefined queries" including "MetS\_IDF" and "MetS\_WHO". Below this, a "Request" section shows a logical query: "ANTIHYPR (Antihypertensives) OR CHD (Coronary Heart Disease) OR BMI (BMI)". At the bottom, there are checkboxes for "Use split by collection" (checked) and "Use relations" (unchecked).

**SERC**  
Swedish e-Science Research Centre



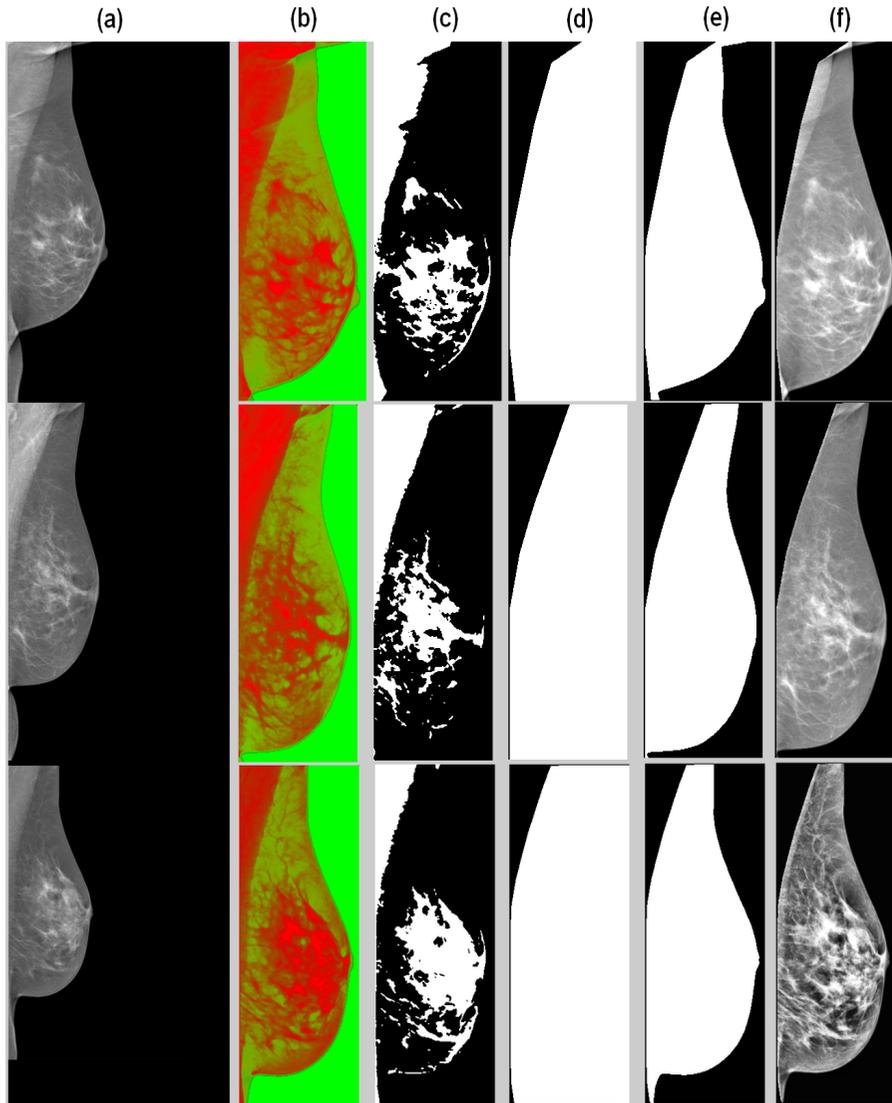
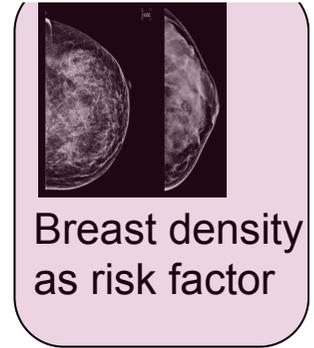
EMBL-EBI 

**FIMM**  
Institute for Molecular Medicine Finland  
Nordic EMBL Partnership for Molecular Medicine

O. Spjuth et al.  
“Enabling integrative cross-biobank research: the SAIL method for harmonizing and linking biomedical and clinical data across disparate data archives”  
*Eur J Hum Genet.* 2015 Aug 26. [Epub ahead of print]

# Breast density as risk factor

## Imaging technology and computational techniques



- Applying e-Science methods to image analysis data can feed into prediction modeling and help guide screening and preventative strategies.
- eCPC developed new measures of mammographic density, with the aim of providing stronger risk factors for breast cancer.



Keith Humphreys

# Biomarker discovery and validation

- Validated genetic and protein markers conveying prostate cancer risk
- Studied genomic profile of breast cancers

→ Data-intensive bioinformatics making use of national HPC e-Infrastructures (SNIC)



Biomarkers for  
breast and  
prostate cancer

Stora prostata-  
cancerstudien:

---

**STHLM 2**

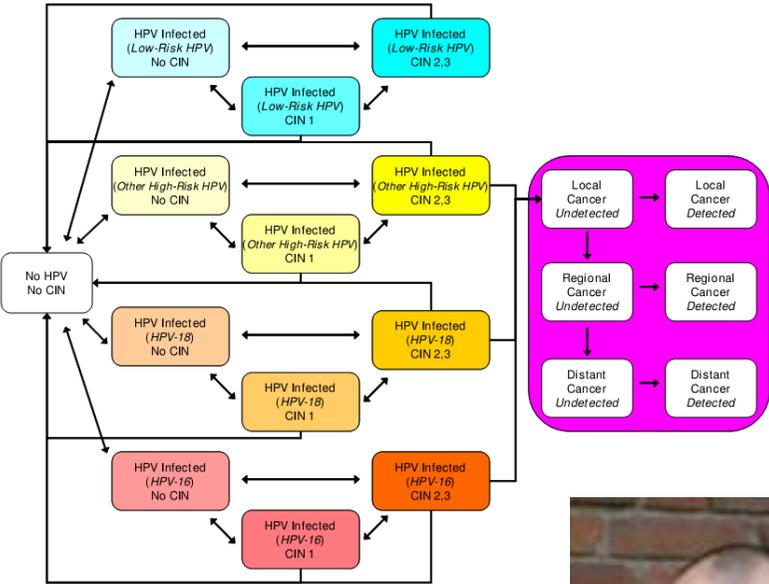
Stora nationella  
bröstcancerstudien:

---

**Karma**

# eCPC Microsimulation Model

- Simulate individual event histories
- Aggregate to population level
- Assess the model fit (calibration)
- Evaluate effects of intervention

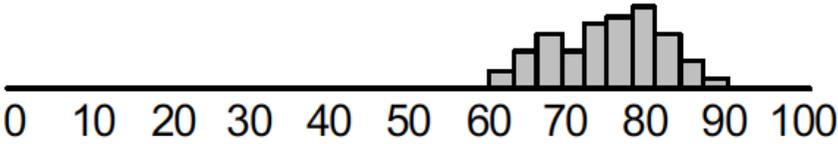


Mark Clements

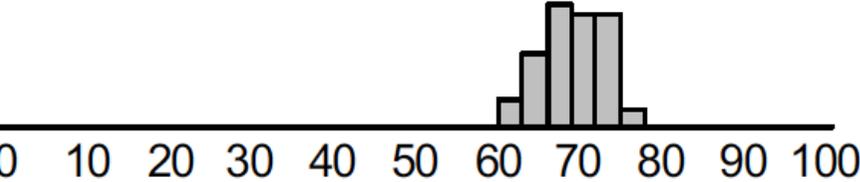
## Example: Reduction in Cervical Cancer Incidence

(Bodhager Diebert et al. 2007)

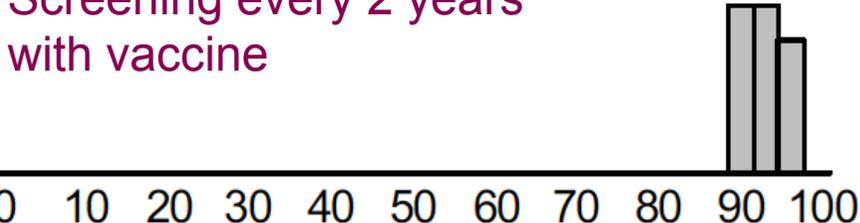
Vaccination only



Screening every 2 years only



Screening every 2 years with vaccine



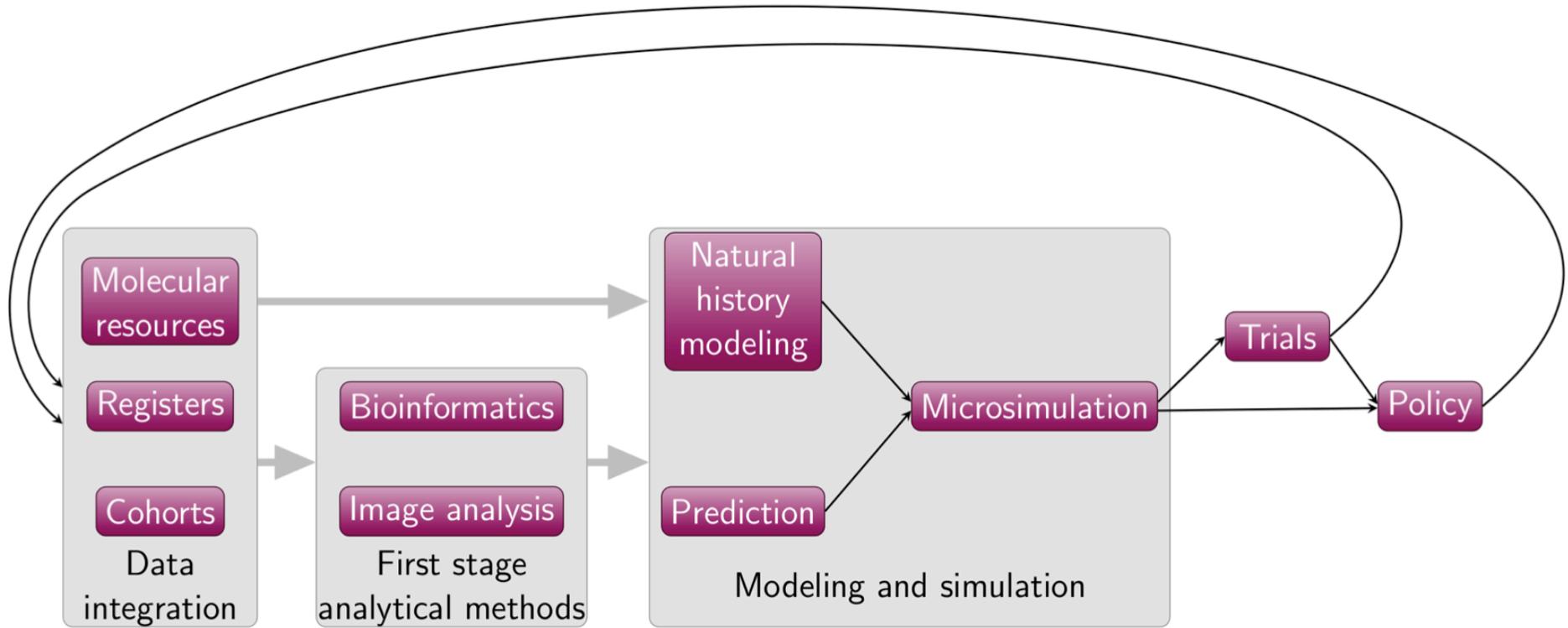
## eScience and eCPC microsimulation

- Used OpenMP/MPI or PGAS for Bayesian calibration of the microsimulation on HPC
- Compared HPC with MapReduce for the microsimulation
- Trans-compile from C++ to WebAssembly/JavaScript for microsimulation on the browser
- Compared client-server and client-only web interfaces for the microsimulation

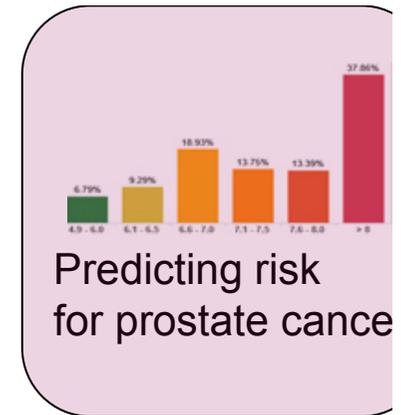
## Microsimulation and STHLM-2



- **Stockholm 2 study:** Find out how heredity, environment and genetic changes can be used to predict the risk of developing prostate cancer.
  - 50,000 participants
- Microsimulation was developed and used in Sthlm 2 for evaluating and planning screening strategies
  - Resulted in the design of Stockholm 3 study: Develop strategies to improve risk predictions



## Sthlm 3 study: Predicting risk for prostate cancer based on PSA + biomarkers



- Constructed and validated Stockholm 3 model (S3M)
  - Prediction algorithm for predicting a man's risk for having prostate cancer, which shows improved specificity (given constant sensitivity) over prostate specific antigen (PSA) testing alone in the population-based STHLM3 diagnostic trial
  - 50 000 men in Study Q1 2015
  - Thermo Fisher customized chip
- The developed microsimulation model will be used for health economics calculations



## 2015: eCPC sharpens the focus on translational medicine

- Adding new prioritized field: **Clinical sequencing**
  - How can new high-throughput sequencing technologies aid cancer diagnostics and treatments?
- Joint collaboration, SeRC and eSENCE



PI: Junii Palmgren



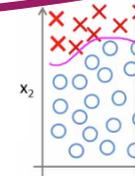
# eCPC research on clinical sequencing

Applied e-Science research



Individualized diagnostics

e-Science methods development



Posters:  
18 + 28

Prediction models,  
machine learning

e-Infrastructure development



Posters:  
41 + 46

Automation, Big Data

# Clinical sequencing supported by eCPC



## Stockholm: ClinSeq

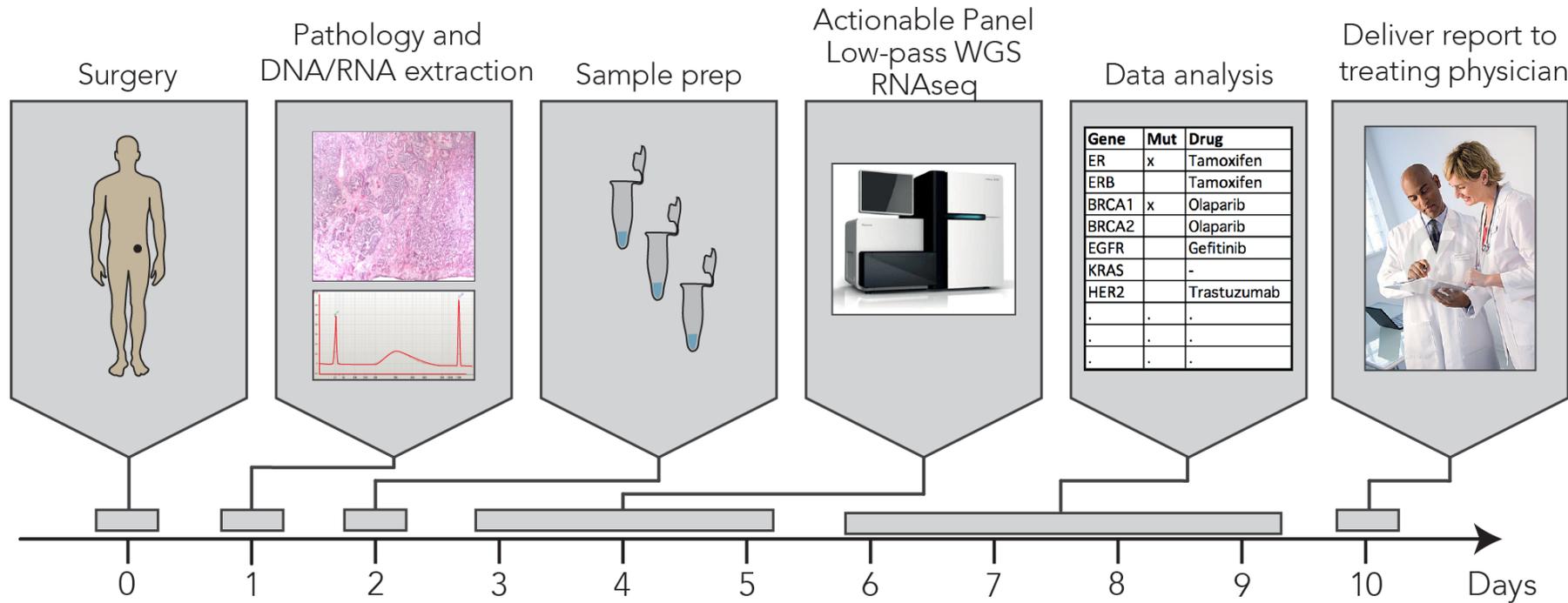
- Short read technology
- Pan-cancer approach
- [www.clinseq.se](http://www.clinseq.se)

## Uppsala

- Long read technology
- Targeted approaches



# ClinSeq pipeline



**!** Patient value  
Research opportunities

**Pan-cancer approach - the same pipeline for all cancers**

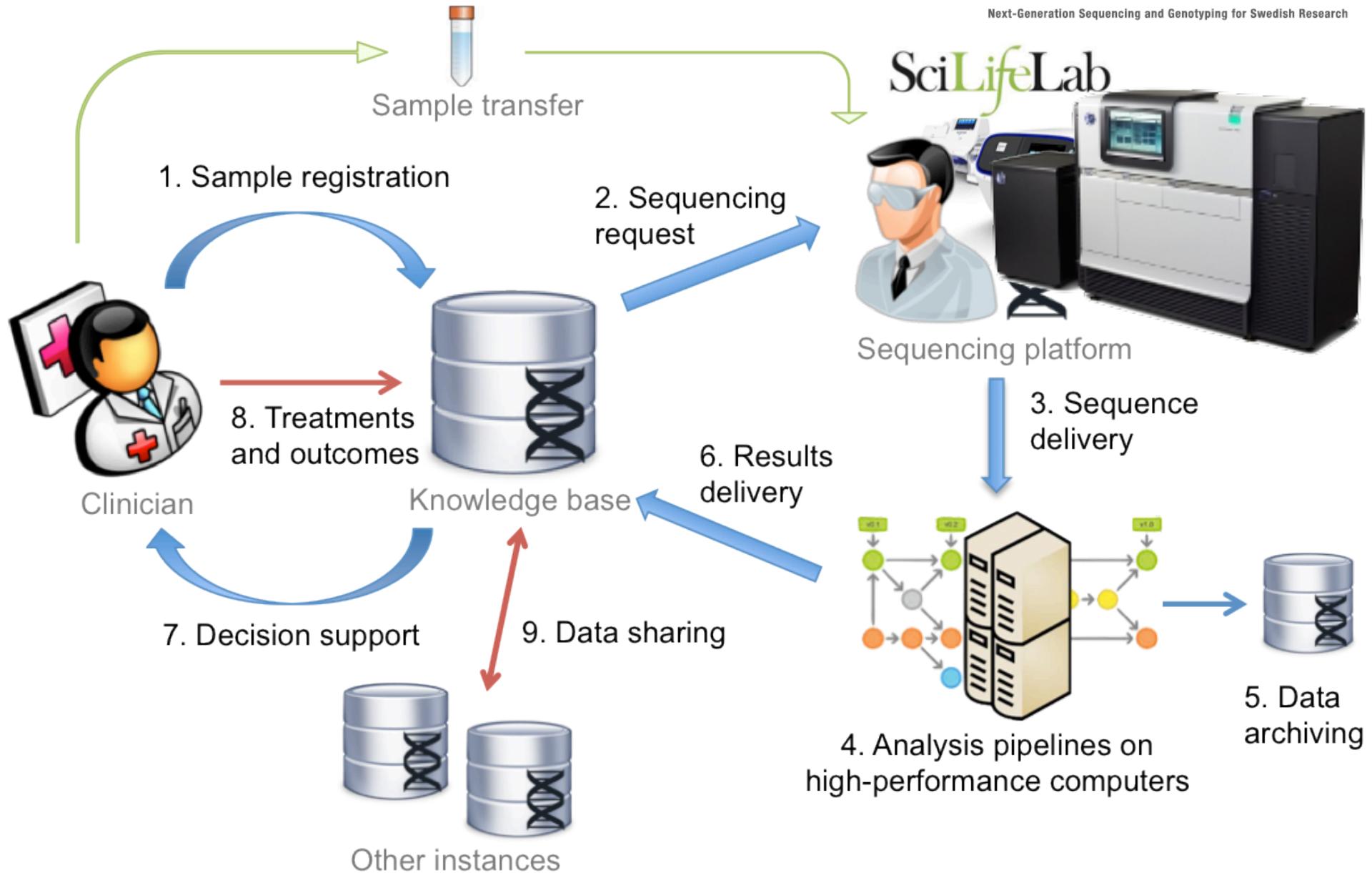
## E-science in ClinSeq

- **HPC / distributed computing**
  - Bioinformatic pre-processing in CLINSEQ: 15 Gb raw data per patient (~100 CPU hours / patient)
  - Automated bioinformatic pipeline that run in cluster environment (e.g. UPPMAX)
- **Machine learning and computational statistics:** Methods development and application
  - Supervised and unsupervised learning across multiple types of high-dimensional data
  - Biomarker discovery (robust methods for variable selection): definition of robust sets of clinically relevant biomarkers
- Integration of in-house data with **public DBs and data sets**
  - The Cancer Genome Atlas (TCGA), clinically relevant mutations (COSMIC, ClinVar)

## ClinSeq in breast cancer

- Results suggests that routine clinical biomarkers could be replaced with DNA and RNA sequencing-based diagnostics
- The ClinSeq profile add value by providing more detailed diagnostic information (subtype, mutations, transcriptomic grade)
- Prospective validation study is planned (N=500)

# Supporting clinical sequencing at Uppsala Academic Hospital





## e-Science in eCPC clinical sequencing in Uppsala

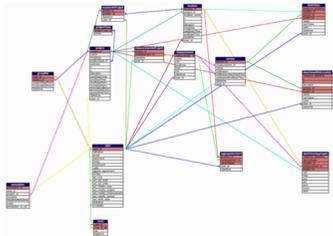
- Construct and automate pipelines on HPC and Cloud
- Transfer results into database system
  - Structured data (variants, tables, visualizations, data files etc.)
  - Inside hospital domain with clinical phenotypes (future)
- Security – sensitive personal information
- Encrypt and archive raw data
- Build up knowledge base for future predictive modeling
  
- *Now:* In pilot production at UAH
- *Future:* Extend to other cancers and genomic regions

## eCPC involvements

- NIASC – The Nordic Information for Action e-Science Center of Excellence (<http://nordicehealth.se/>, 2014-1018)
  - Integrated Nordic research in population-based cancer screening
  - generic eScience infrastructure
  - eScience-based predictive algorithm within a national screening program
- PhenoMeNal H2020 project (start Sept. 1<sup>st</sup> 2015)
  - clinical metabolomics

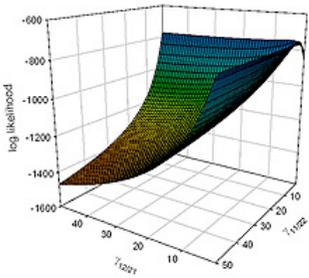
# Sweden and Nordic countries

Enormous potential for eScience in medical research



Reliable demographics and healthcare registers

Clinical and population cohorts



National biobanks

Biotechnology and Information technology

High quality epidemiology and clinical research

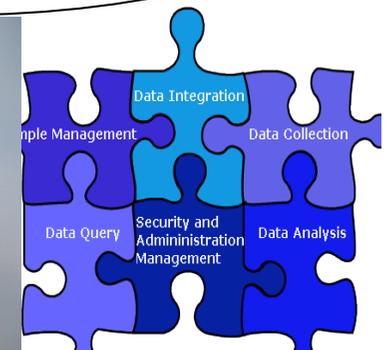
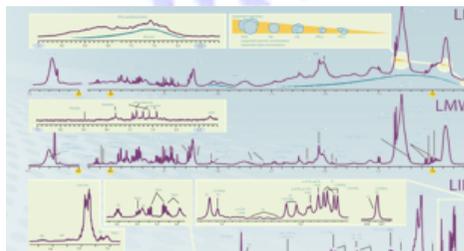
Bioinformatics, computational biology, biostatistics



illumina

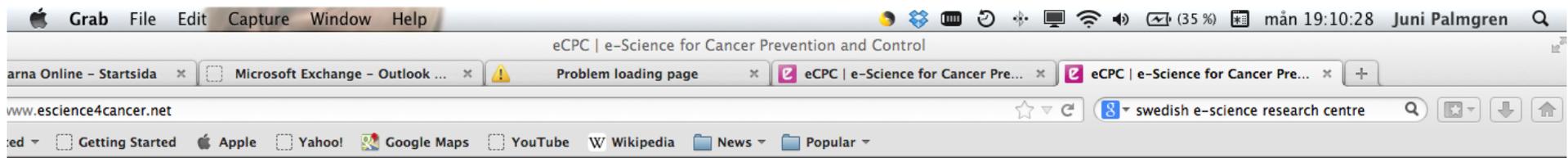


Applied Biosystems



- Thank you -

[www.ecpc.e-science.se](http://www.ecpc.e-science.se)



[Start](#) | [Work Packages](#) | [Targeted Diseases](#) | [Applications](#) | [Downloads](#) | [About](#) | [Interviews](#)

