



UNIVERSITY OF TARTU

Elektroonilised andmekogud –  
baas pikaegseteks uuringuteks

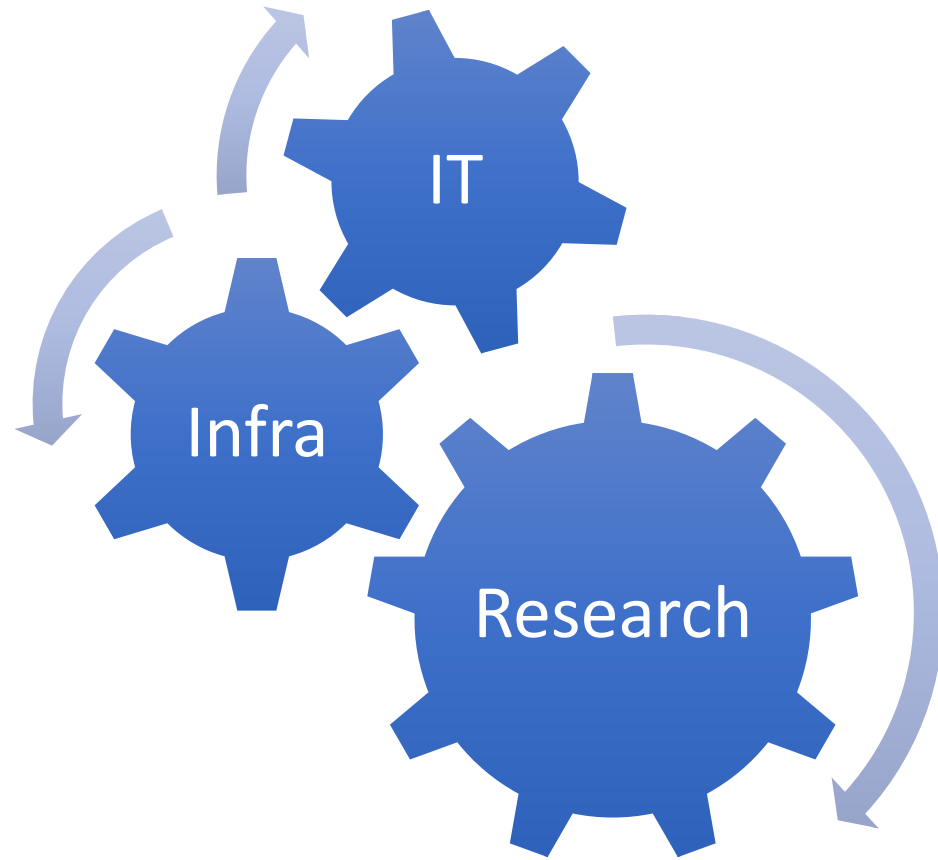
Jaak Vilo

Academic excellence since 1632

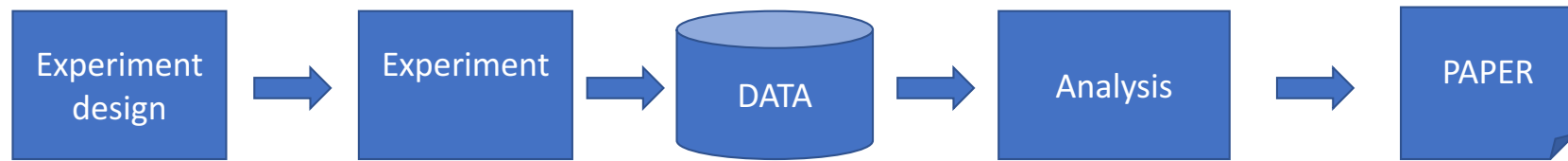


**EMIF**  
European Medical  
Information  
Framework

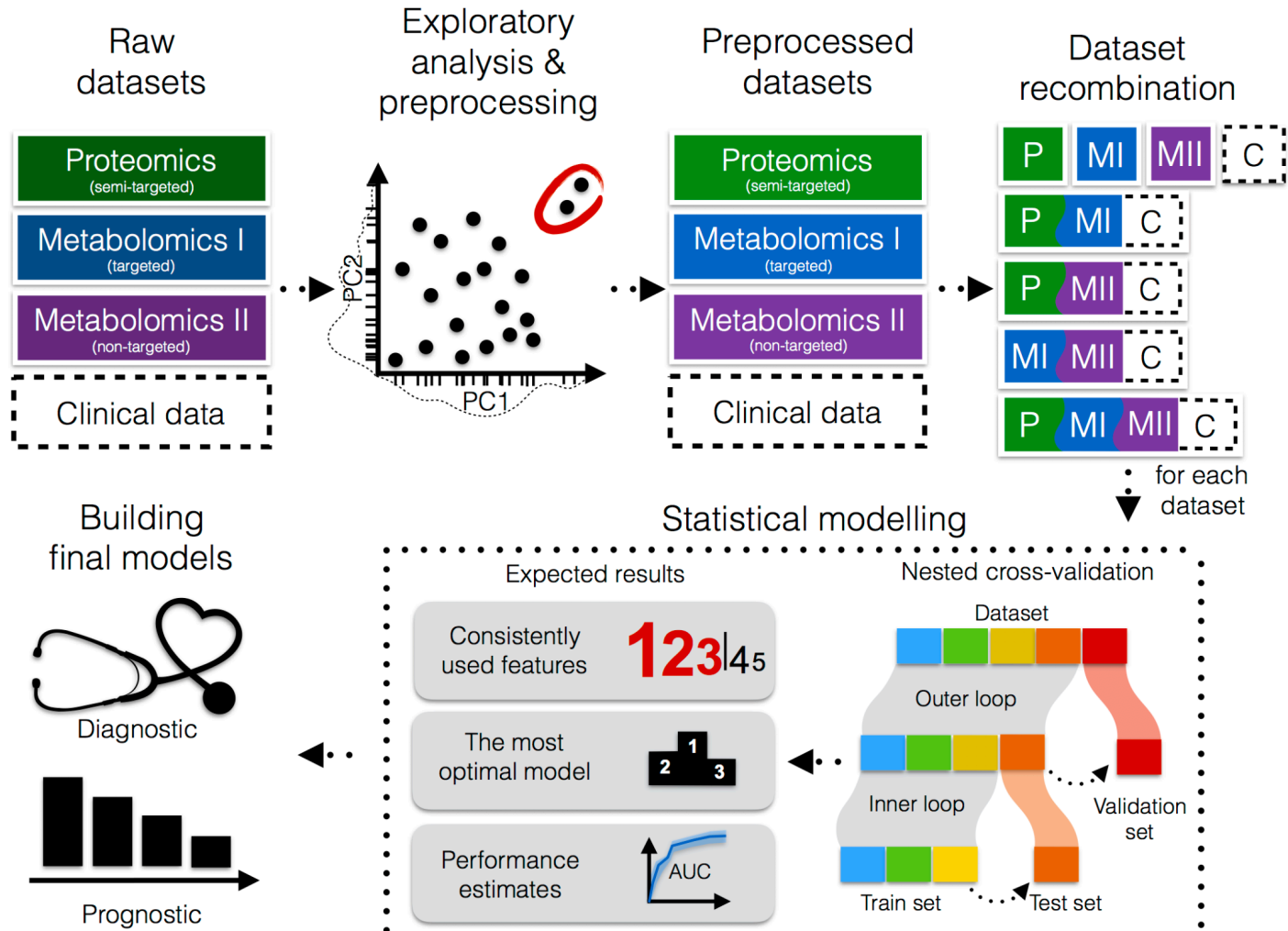




# Research “workflow”



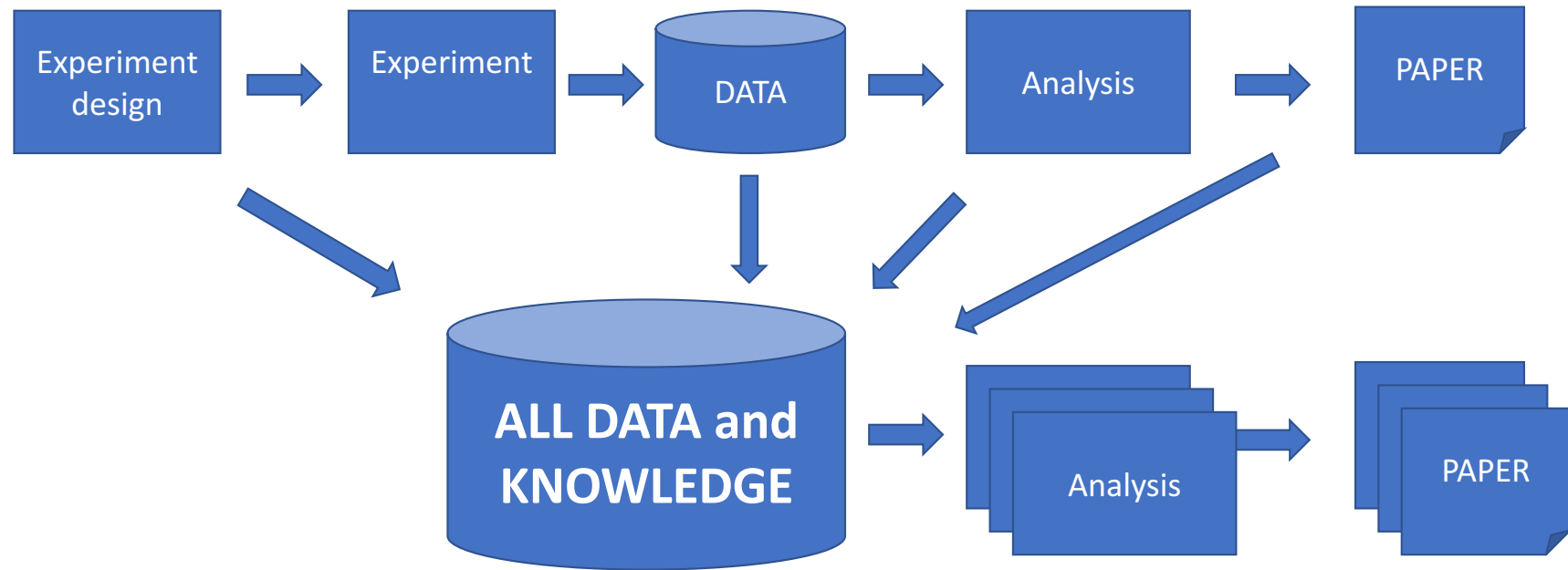
**Annex 7: Strategic plan for data analysis within the BioEndoCar project**



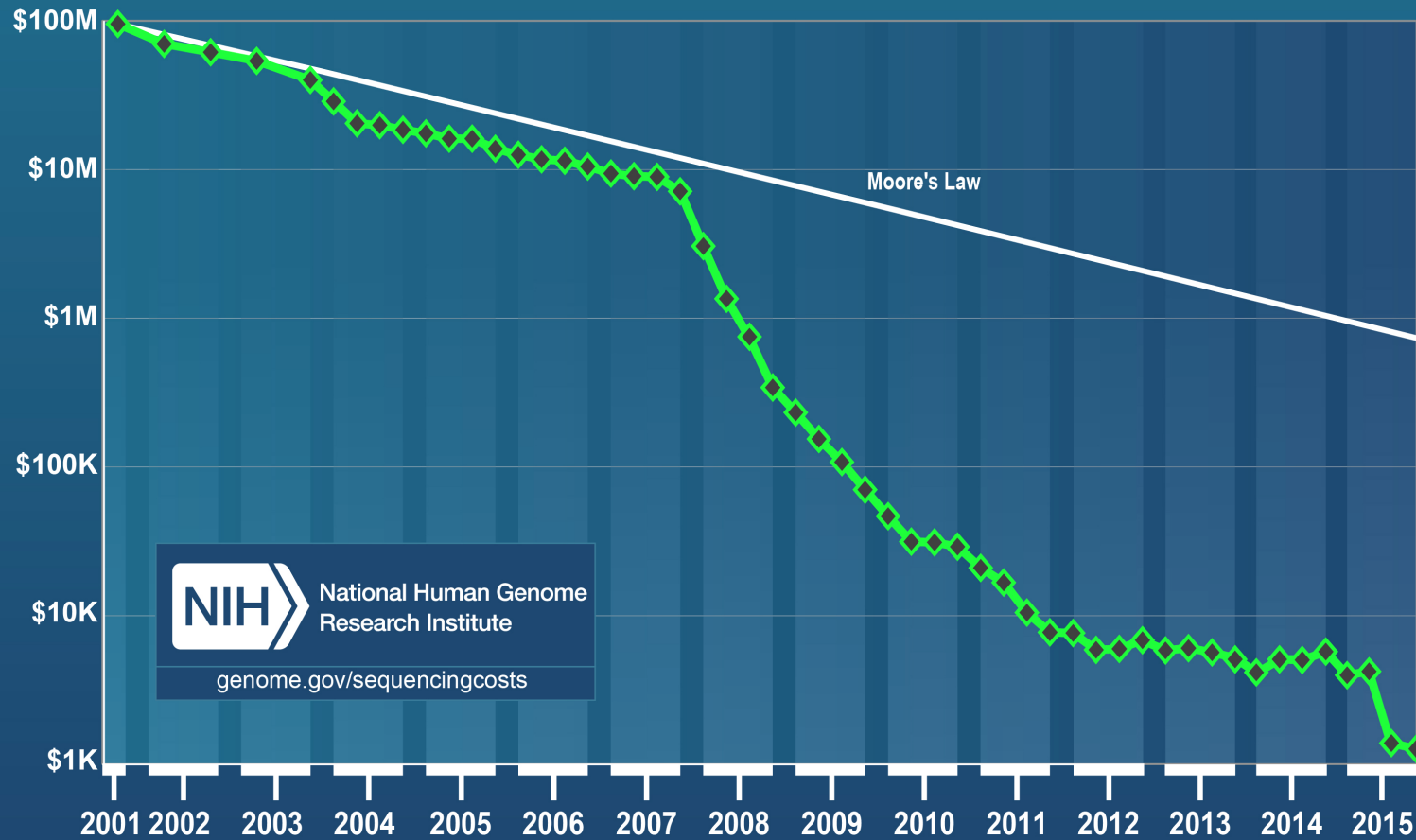
Nested CV visualisation was adopted from Gael Varauquaux

Hypothesis-driven  
vs  
Hypothesis-free(?)

# Research “workflow”



# Cost per Genome







SOLID 5500



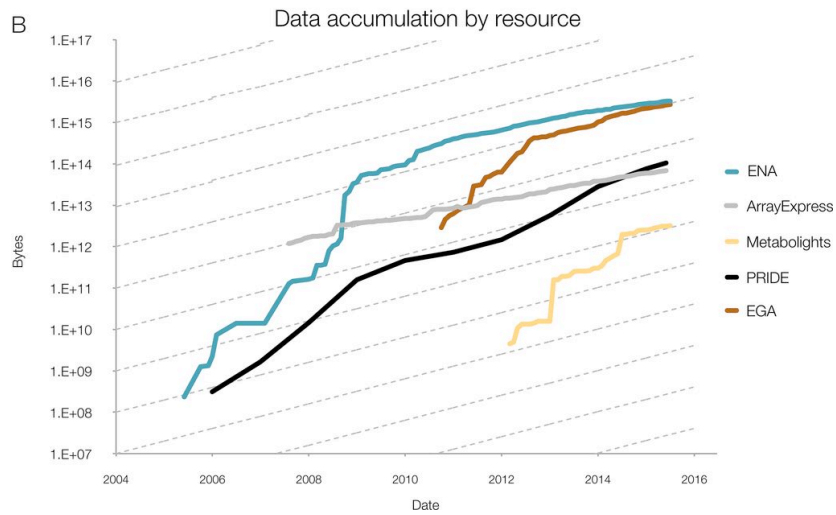
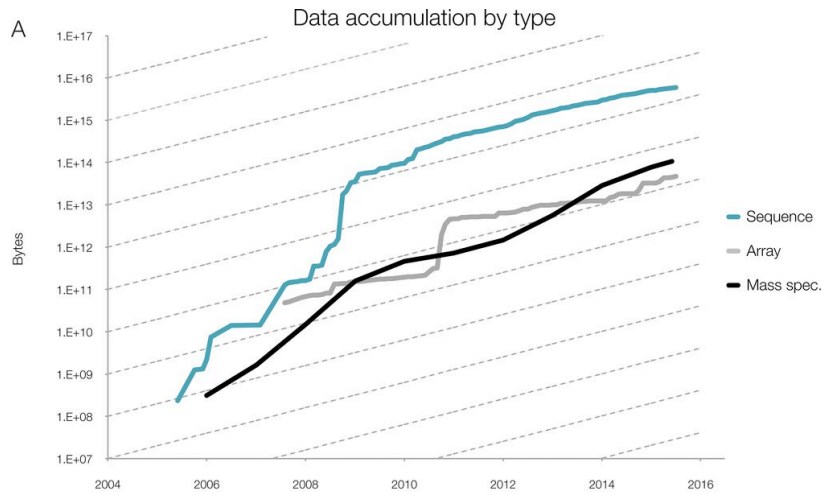
Illumina HiSeq2000



# Data growth in the life sciences

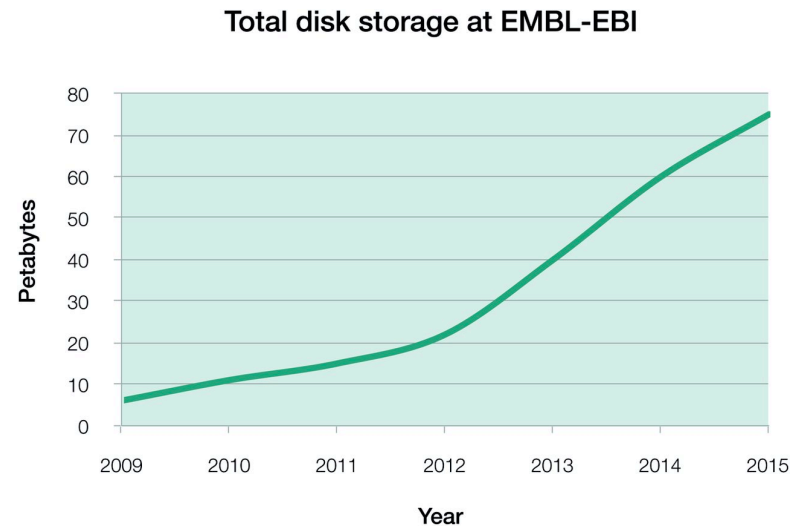
Data growth at EMBL-EBI

Source: Charles E. Cook et al. *Nucl. Acids Res.* 2016;44:D20-D26



# The data challenge: Data growth

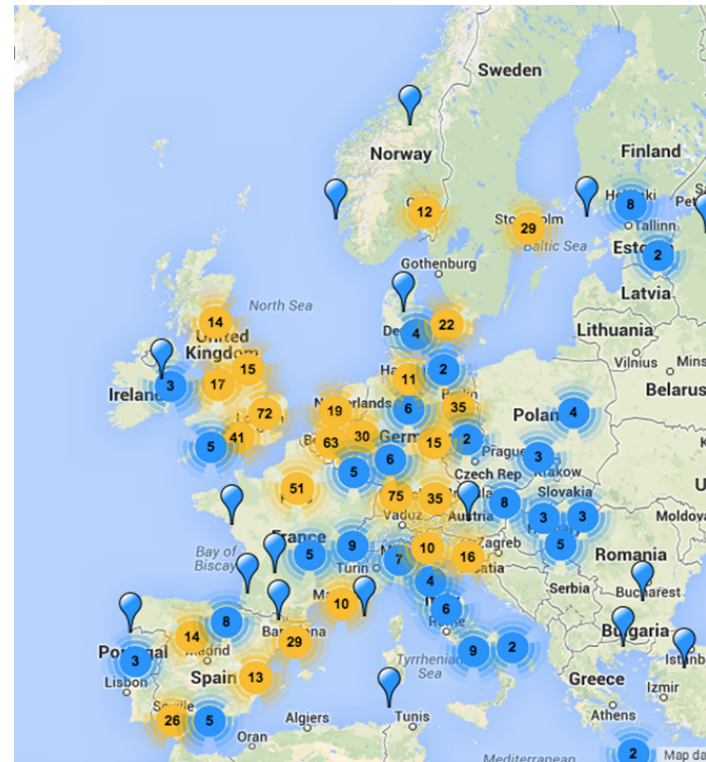
- Computer speed and storage capacity is **doubling every 18 months** and this rate is steady
- DNA sequence data is **doubling every 6-8 months** over the last 3 years and looks to continue for this decade



Source: Charles E. Cook et al. *Nucl. Acids Res.* 2016; 44: D20-D26

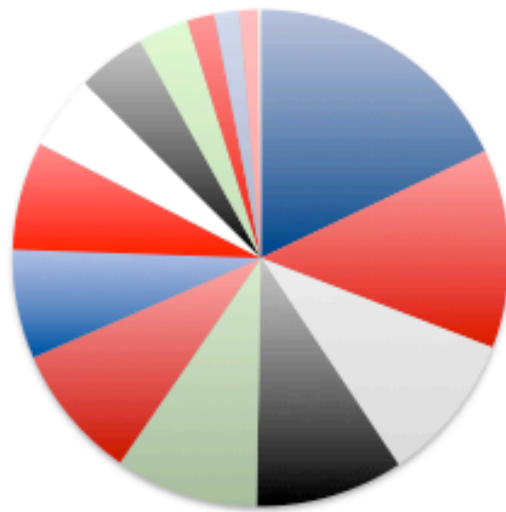
# The data challenge: Geographic spread

- Data production sites increasing across Europe
- European Illumina seq sales up 20% 2014



Source: <http://omicsmaps.com>

# Data resources in life science






- Genomics Databases (non-vertebrate) (17.9%)
- Protein sequence databases (12.9%)
- Human Genes and Diseases (9.8%)
- Structure Databases (9.7%)
- Metabolic and Signaling Pathways (9.3%)
- Nucleotide Sequence Databases (8.8%)
- Human and other Vertebrate Genomes (7.1%)
- Plant databases (7.1%)
- RNA sequence databases (4.9%)
- Microarray and other Gene Expression Databases (4.5%)
- Other Molecular Biology Databases (3.3%)
- Immunological databases (1.8%)
- Organelle databases (1.6%)
- Proteomics Resources (1.2%)
- Cell biology (0.2%)

**molecular biology  
data resources**

**~1800**

Nucleic Acids Research annual Database Issue  
and the NAR online Molecular Biology Database Collection in 2012.  
MY Galperin, GR Cochrane – Nucleic Acids Research, 2011

# We generate data faster than we can deposit it

		Network file transfer rate
<b>24 hours</b>		<b>100 Mb</b>
	DNA sequencing ~100 GB	~5 hours
	Mass spectrometry ~4 TB	~4 days
	Microscopy ~4 TB	~4 days



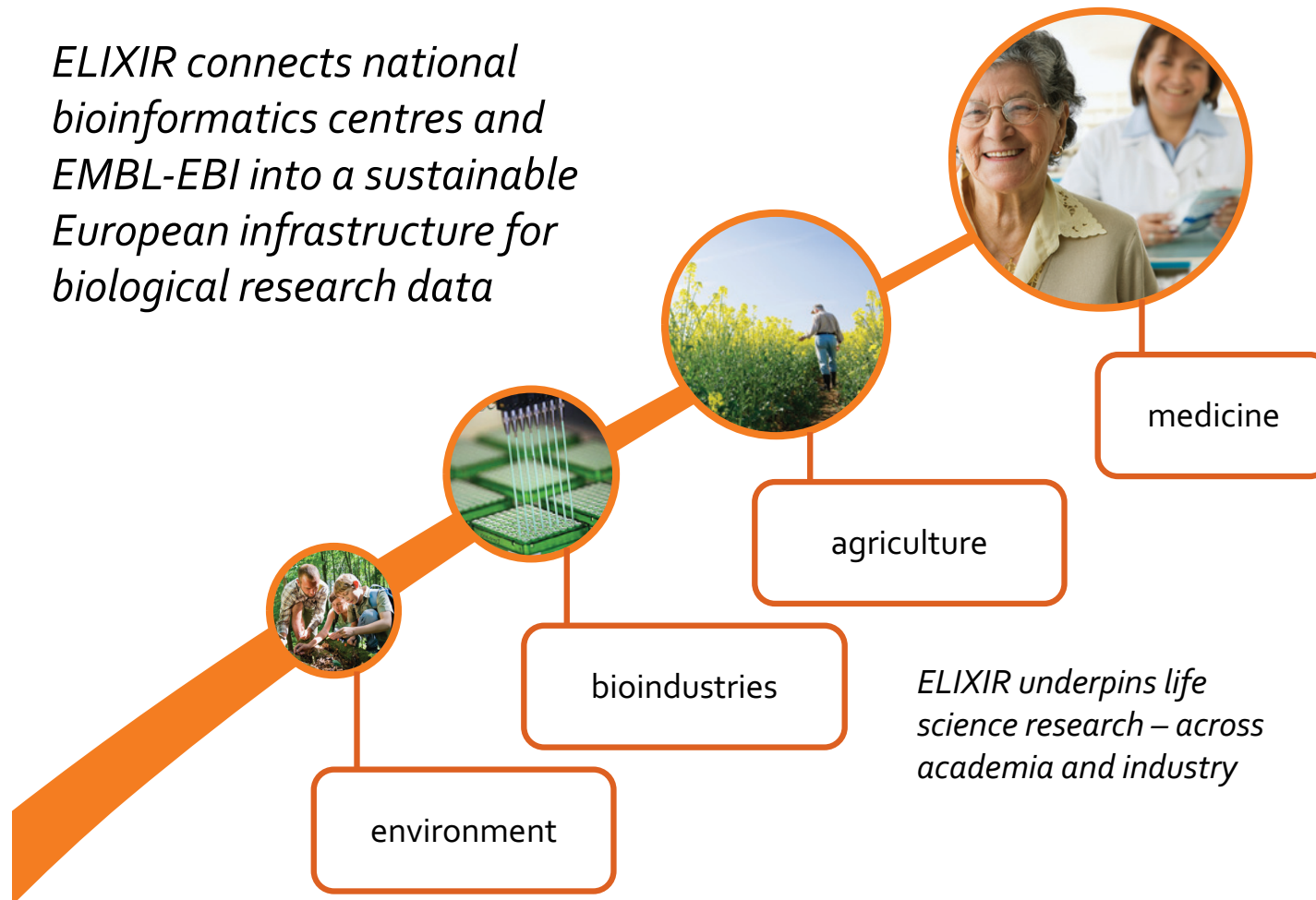
# ELIXIR

*Safeguarding the results of life  
science research in Europe*



[www.elixir-europe.org](http://www.elixir-europe.org)

*ELIXIR connects national bioinformatics centres and EMBL-EBI into a sustainable European infrastructure for biological research data*



*ELIXIR underpins life science research – across academia and industry*



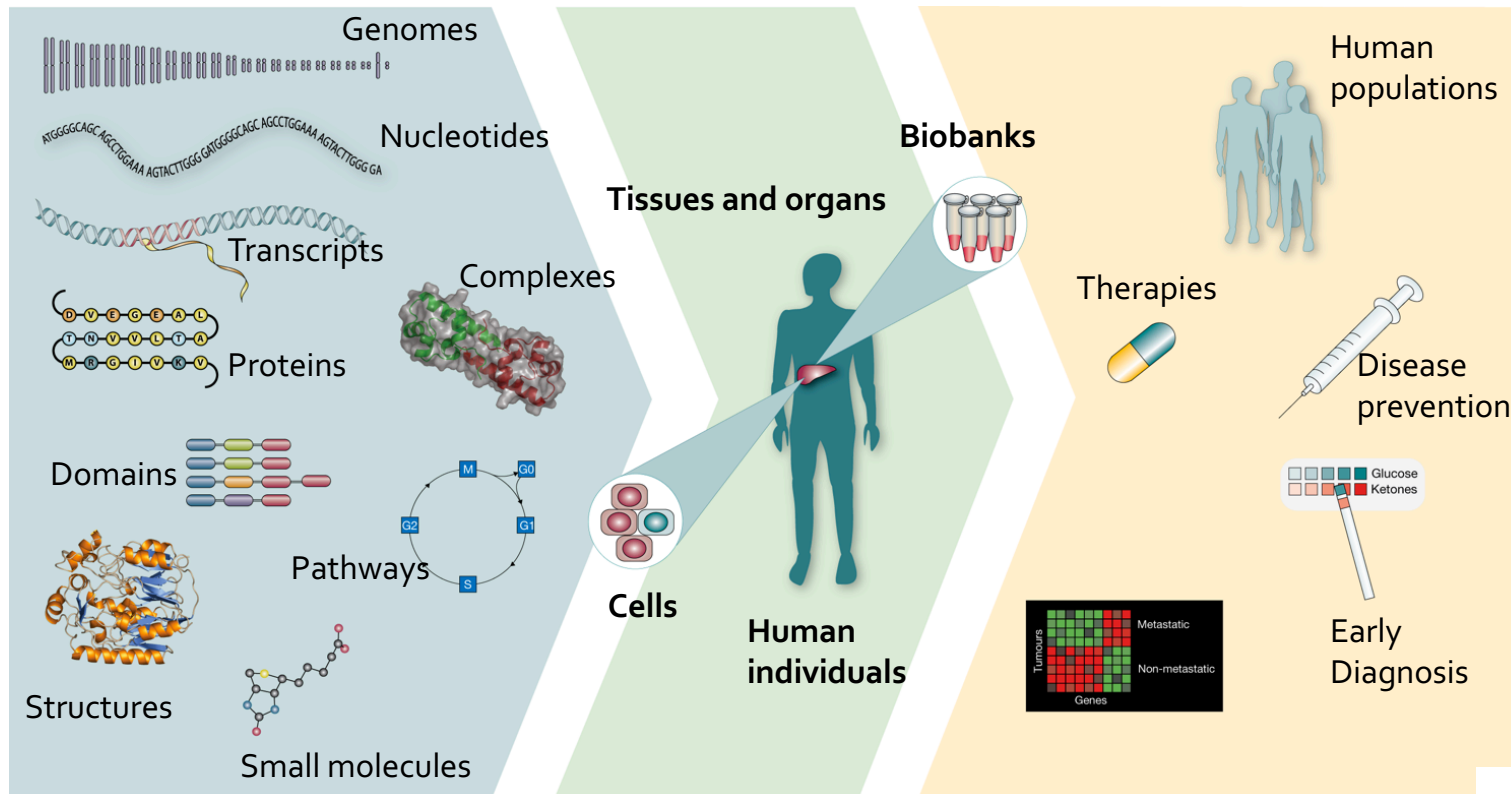


# From molecules to medicine

Molecular components

Integration

Translation



# ELIXIR: European infrastructure for biological information

Data infrastructure for Europe's life-science research:



*Data*



*Interoperability*



*Tools*



*Compute*



*Training*



*Marine metagenomics*



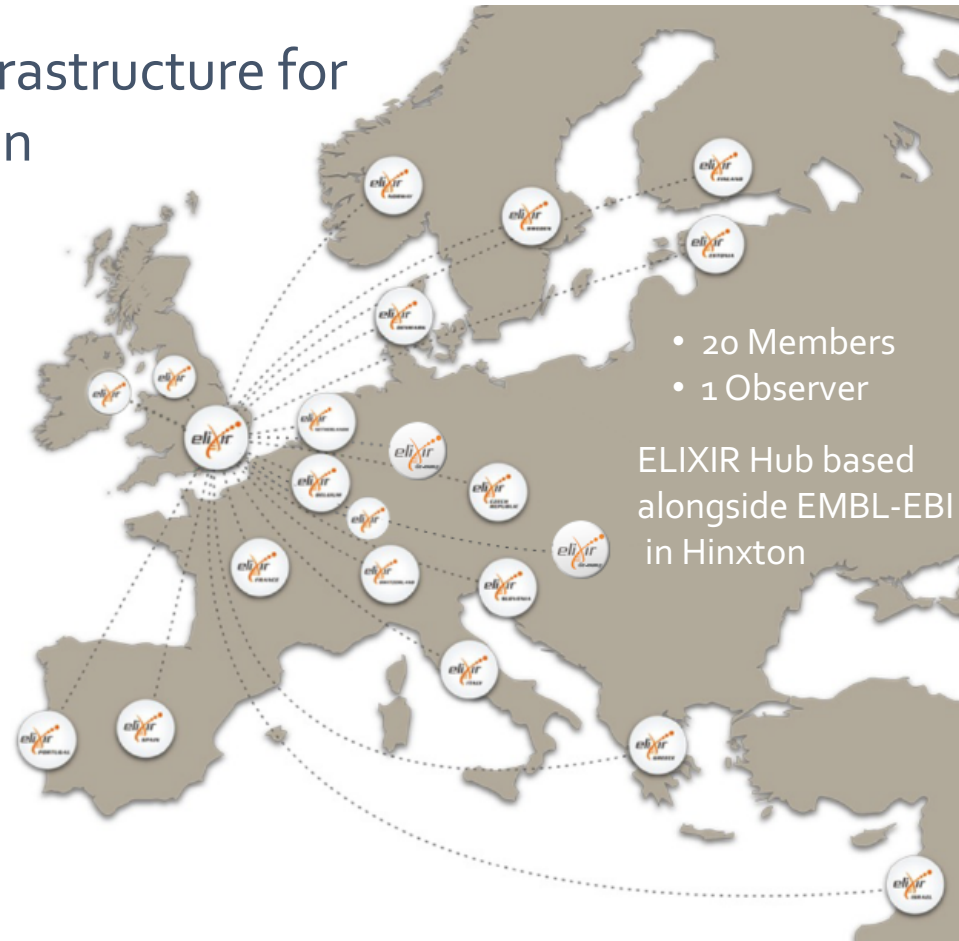
*Crop and forest plants*



*Human data*



*Rare diseases*



[www.elixir-europe.org](http://www.elixir-europe.org)



[@ELIXIREurope](https://twitter.com/ELIXIREurope)



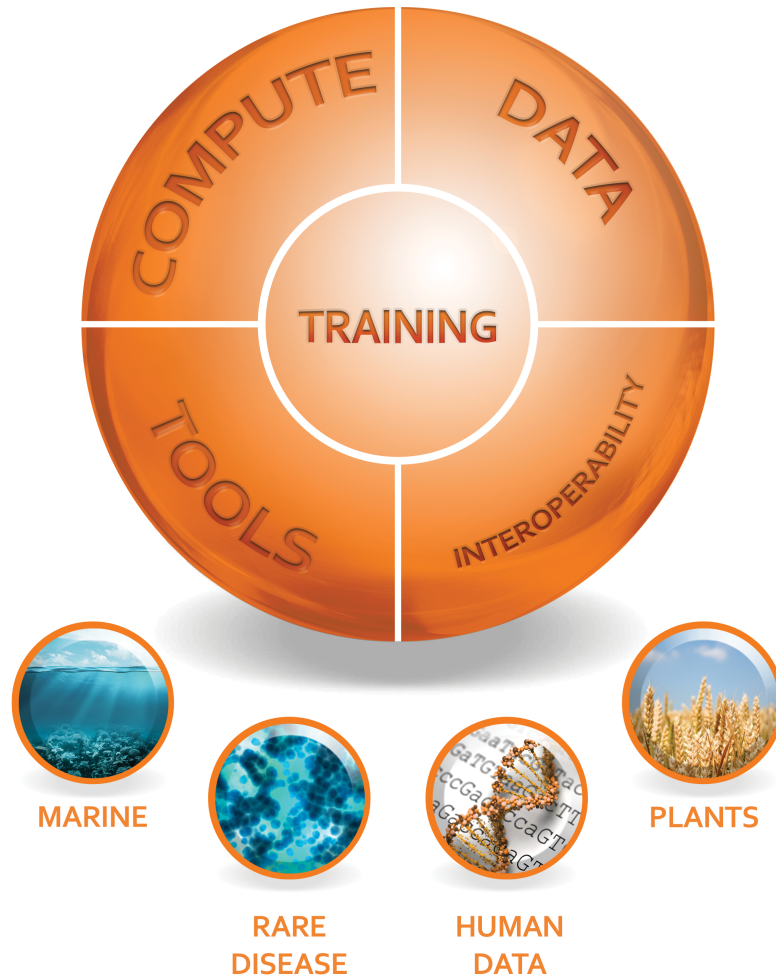
## ELIXIR Members



## ELIXIR Observers



# ELIXIR Structure



Five technical platforms for **Compute**, **Data**, **Tools** and **Interoperability**

Complemented by four use cases for **marine meta-genomics**, **rare diseases**, **human data** and **plants sciences**

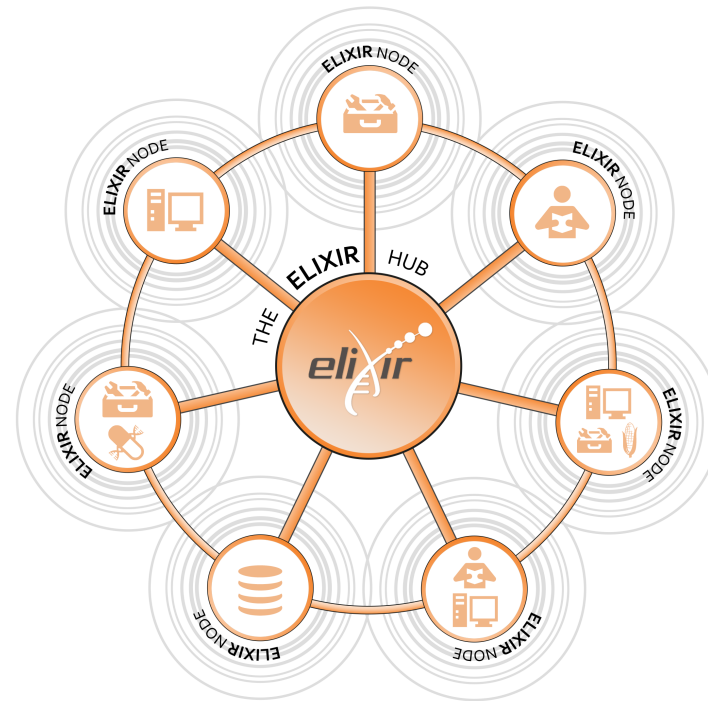


# A distributed infrastructure to scale with the challenge

**ELIXIR** data infrastructure for Europe's life science research sector

**ELIXIR** Nodes build local bioinformatics capacity throughout Europe

**ELIXIR** Nodes build on national strengths and priorities



# FAIR principles for data

**F**indable

**A**ccessible

**I**nteroperable

**R**eusable

Data Management Plan; Data management life-cycle

To be Findable:

F1. (meta)data are assigned a globally unique and eternally persistent identifier.

F2. data are described with rich metadata.

F3. (meta)data are registered or indexed in a searchable resource.

F4. metadata specify the data identifier.

To be Accessible:

A1 (meta)data are retrievable by their identifier using a standardized communications protocol.

A1.1 the protocol is open, free, and universally implementable.

A1.2 the protocol allows for an authentication and authorization procedure, where necessary.

A2 metadata are accessible, even when the data are no longer available.

To be Interoperable:

I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.

I2. (meta)data use vocabularies that follow FAIR principles.

I3. (meta)data include qualified references to other (meta)data.

To be Re-usable:

R1. meta(data) have a plurality of accurate and relevant attributes.

R1.1. (meta)data are released with a clear and accessible data usage license.

R1.2. (meta)data are associated with their provenance.

R1.3. (meta)data meet domain-relevant community standards.

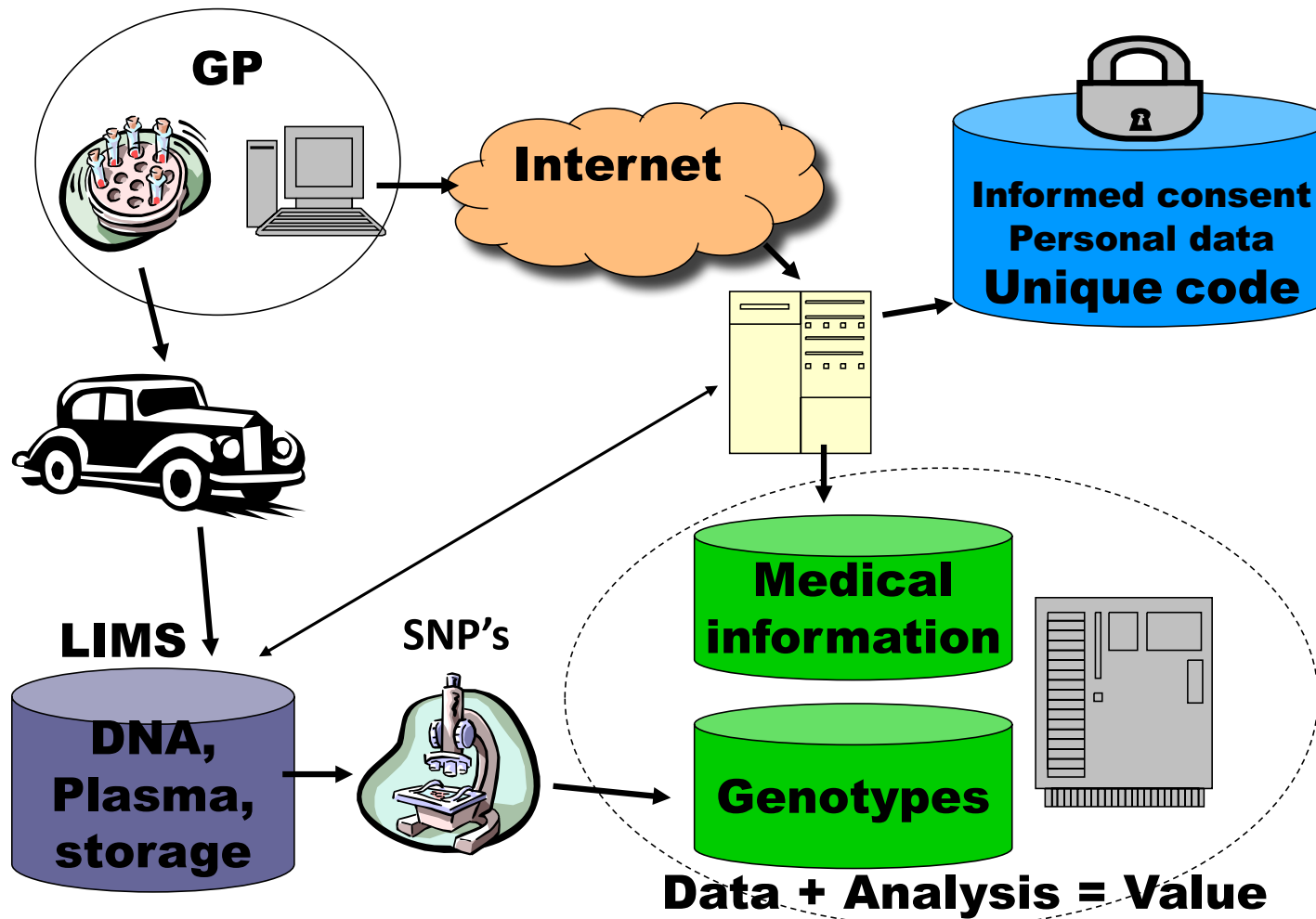
# Rewind 18 years... Estonian Biobank



estonian genome center



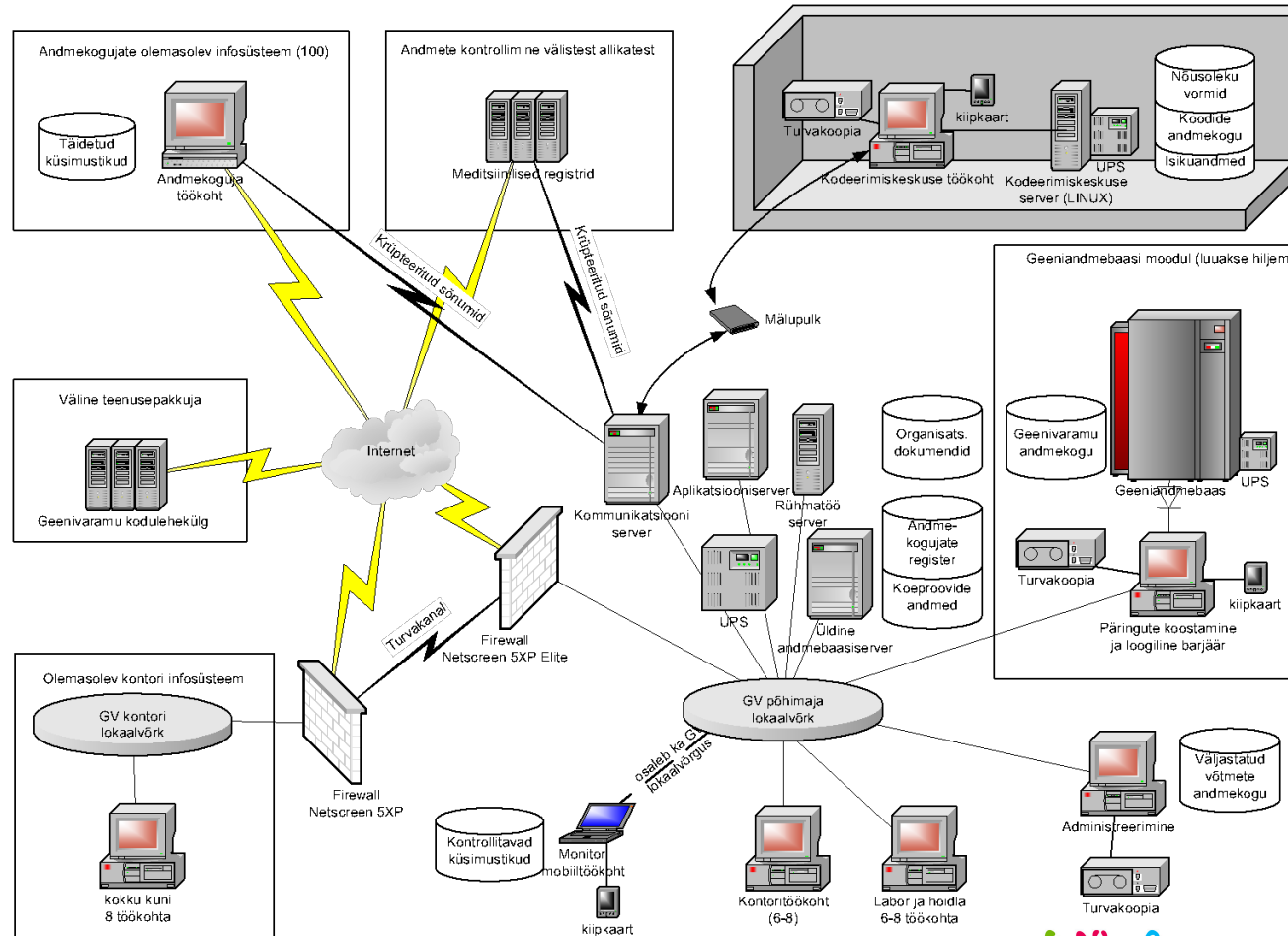
# Process of data collection and handling







# Data collections





## Company Introduction & Estonian Registries

[www.quiretec.com](http://www.quiretec.com)

# Qure Data Management Platform

The core product of Quretec, the Qure Data Management software Platform has been designed for collection, handling, and analysis of complex data like health registries, questionnaires, clinical case report forms and other rich metadata under high quality, security, and robustness requirements.



## Qure Browser – a web based data entry user interface

Firefox | Qure Browser - Tuberculosis Register

[1] Persons > [1] JOHN SCHMIDT (33609010275) 000007 > [A] Case 2012-05-22 primary (00000701) > [A.3] Treatment cards > [B] Treatment card 2012-05-22 primary > [B.1] Data card info (0000070101)

Persons | Case 2012-05-22 primary (00000701) | Death Register data | Search of cases | Treatment cards search

Basic data | Treatment cards | Councils | Comorbidities | Surgery | Medications | Drugs adverse events | Lab notifications

Data card info (0000070101)

[B.1] DATA CARD INFO (0000070101)

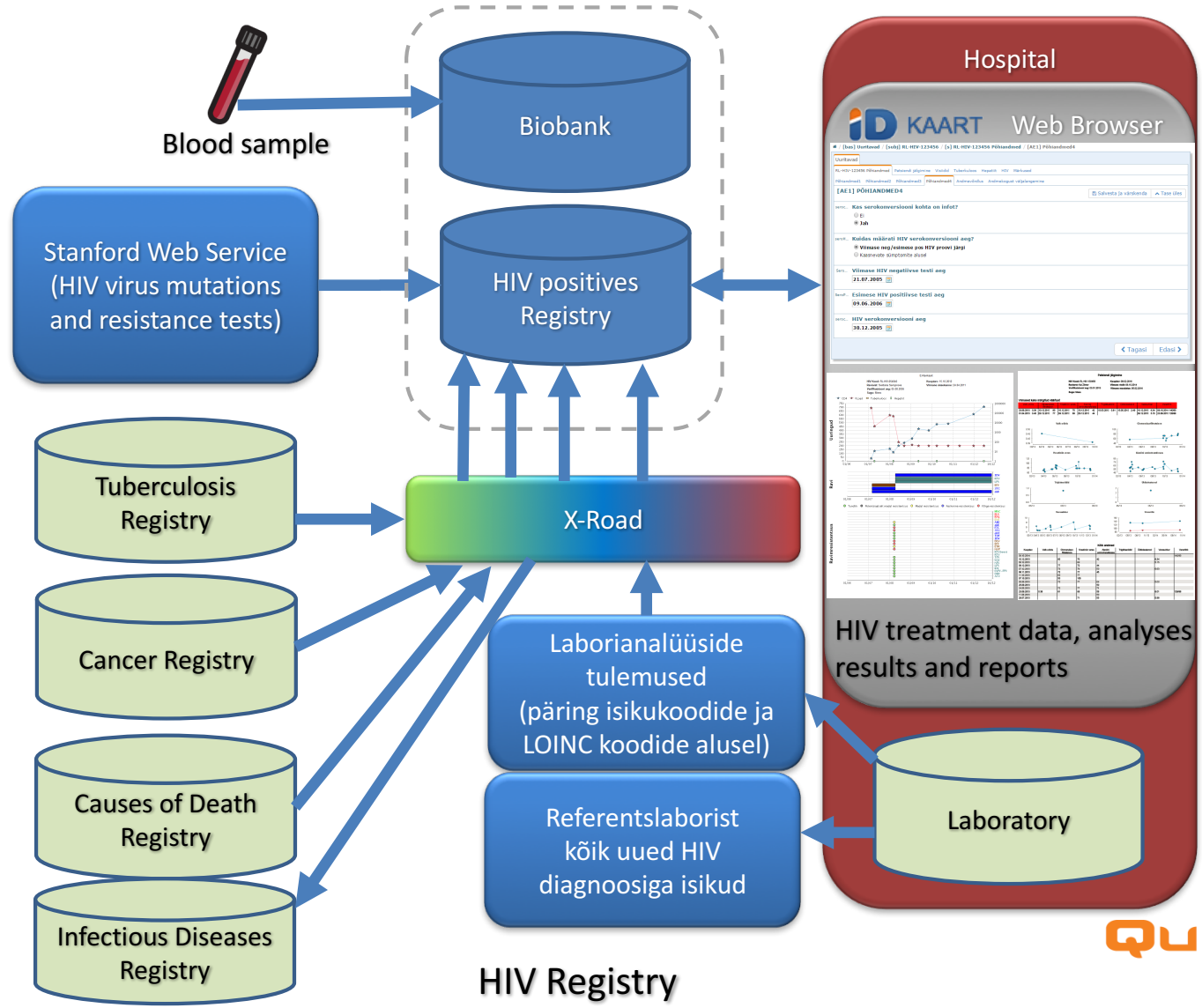
B.1.1 Baseline	22.05.2012	B.1.2 End of treatment	22.05.2012	B.1.3 Date of diagnosis	dd.mm.yyyy	B.1.4 Days of treatment	
B.1.5* Baseline condition	primary	B.1.6 Final Condition	treatment ended	B.1.7 Place of treatment			
B.1.8 Start of outpatient therapy	dd.mm.yyyy	B.1.9 End of outpatient therapy	dd.mm.yyyy	B.1.10 Start of hospital treatment	dd.mm.yyyy	B.1.11 End of hospital treatment	dd.mm.yyyy
B.1.12 Diagnosis of tuberculosis	A15.2 Select... [A15.2] Tuberculosis pulmonum histologic confirmata. Condiciones sub A15.0 datae, histologic	B.1.13 Other diagnosis	Select... (Unselected)				
B.1.14 Diagnosis of tuberculosis 2	Select... (Unselected)	B.1.15 Other diagnosis 2	Select... (Unselected)				
B.1.16 Diagnosis of death	Select... (Unselected)	B.1.17 Date of death	dd.mm.yyyy				
B.1.18* Definition of diagnosis	pulmonary tuberculosis	B.1.19 Form	infiltrative	B.1.20 Location			
B.1.21 Sputum smear	0 not done	B.1.22 Sputum culture	1 +	B.1.23 Destruction	no	B.1.24 Quantiferon test	
B.1.25 HIV		B.1.26 Date of HIV test	dd.mm.yyyy	B.1.27 MDR at baseline	no	B.1.28 MDR ravi lõpul	
B.1.29 BK findings in other material?		B.1.31 Histology					
B.1.33* County date	22.05.2012	B.1.34* County	Tallinn	B.1.35* County medical	Other Doctor (2)		
B.1.36 Medical institution		B.1.37 TOR					
B.1.40 Filler		B.1.41 Filler institution		B.1.42 Filling date	dd.mm.yyyy		
B.1.43 Cause of treatment failure		B.1.44 Basis of		B.1.45 Age	75		
B.1.46* Show statistics	1 yes	B.1.47 Remarks					

# Estonian Health Registries on Qure Data Management Platform



- HIV registry
- Cancer registry
- Tuberculosis registry
- Causes of death registry
- Medical birth registry
- Abortion registry
- Drug treatment database
- Infectious diseases registry
- Estonian Genome Center Biobank
- North Estonia Medical Centre hospital registries
  - Breast Cancer
  - Colon Cancer
  - Pulmonary Arterial Hypertension





### HIV Registry



# QureTEC

- Health registries (Cancer, HIV, Tuberculosis and other registries)
- Biobanks (Estonian Genome Center)
- Clinical Trials
  - EDC
  - Data Management
  - Statistical Analysis
  - Medical Writing



# QureCLINICAL

- Clinical Trials
  - EDC, Data Management, Statistical Analysis, Medical Writing
- 5 active studies at the moment
- Clinical Trial phases I-III
- >100 sites
- >400 active EDC users
- >1000 subjects eCRFs
- Rheumatology, Gynecology, Oncology, Pediatrics

## Näiteid toetatud uuringutest:

CF111 - 2014-2015 uuring Drospirenone farmakokineetika ja -dünaamika hindamiseks rinnapiimas  
H-2011/07-UA - 2013-2014 uuring Herbeprot-P efektiivsuse ja ohutuse hingamiseks troofilise  
jalahaavandiga patsientidel.

DICL001 - 2014-2017 uuring Diclofenac geeli bioekvivalentsuse määramiseks Voltaren geeli suhtes  
põlve osteoartroosiga patsientidel SPIRO - 2014-2015 - uuring Spironolaktooni farmakokineetika  
hindamiseks lastel

HPV-EU-001 - 2015-... uuring HPV vaktsiini turvalisuse ja efektiivsuse määramiseks HPV 16/18  
positiivsetel patsientidel

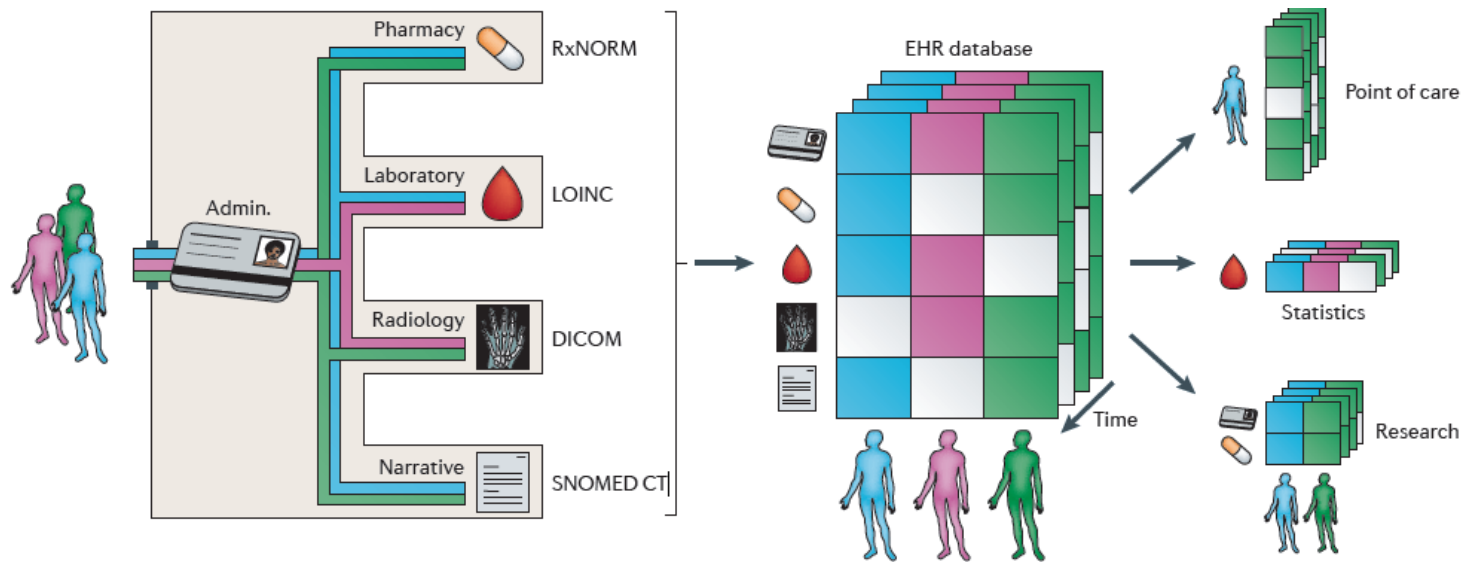
MDP-CLIN-001 - 2015-... uuring Cisplatini ja radioterapia turvalisuse hindamiseks pea- ja kaelavähi  
patsientidel

ARC 2015-... uuring Cefepime farmakokineetika ja -dünaamika hindamiseks lastel

EFIT - 2016-... uuring rTMSi (repetitive transcranial magnetic stimulation) efektiivsuse hindamiseks  
insuldijärgsetel patsientidel

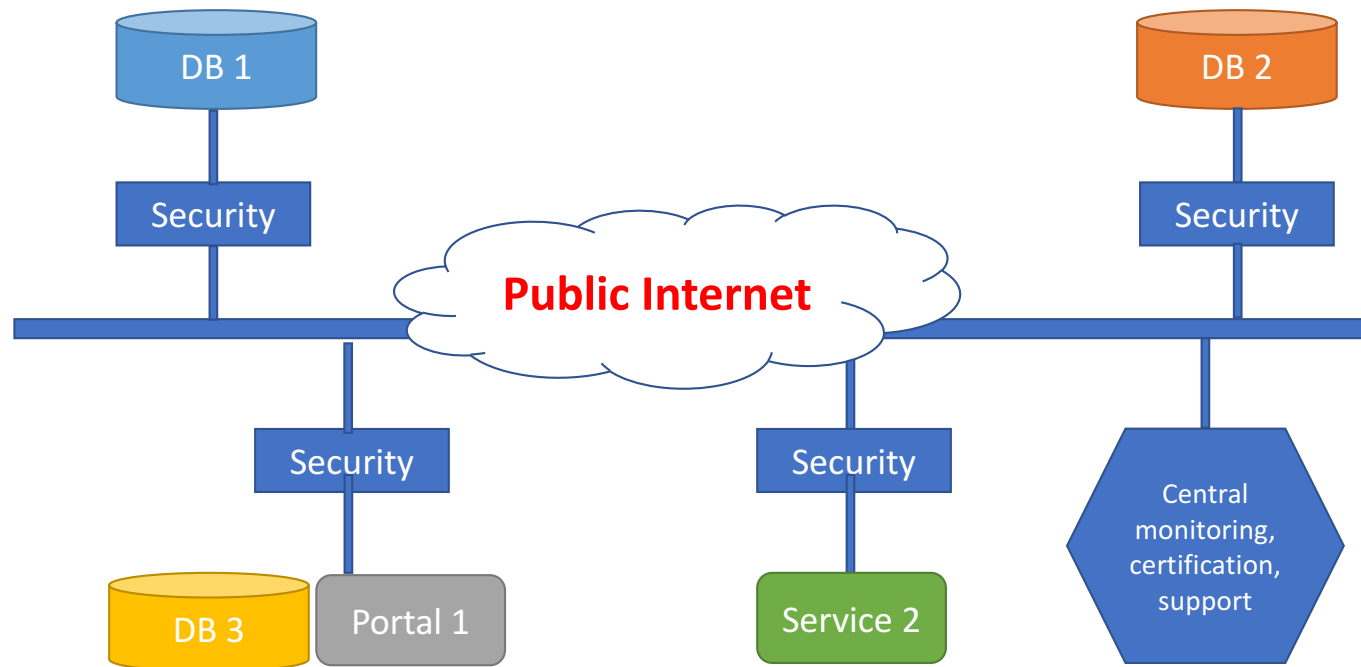
ETVAX - 2017-... uuring ETEC vaktsiini turvalisuse, efektiivsuse ja immunogeensuse hindamiseks  
Lääne-Aafrikat külastavatel patsientidel

# Digitaalne terviselugu ja selle potentsiaal

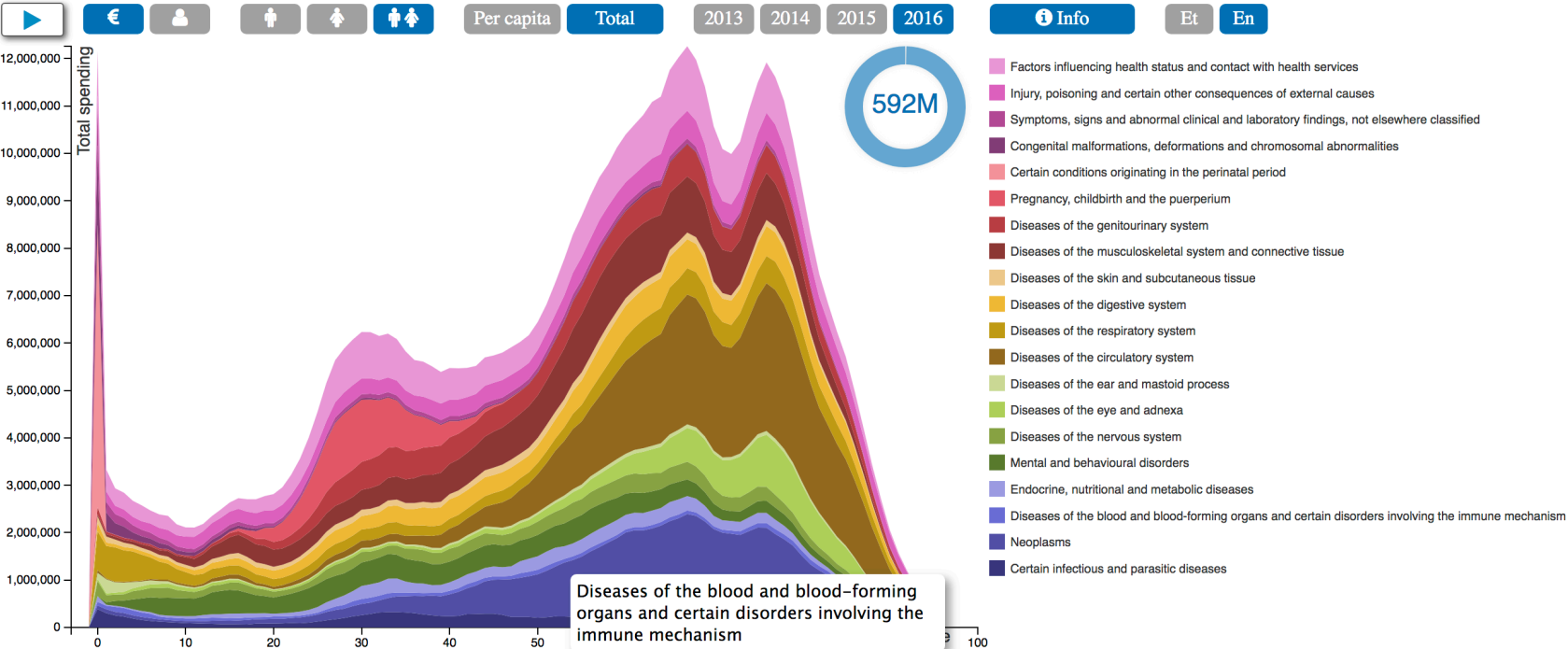


*Peter B. Jensen, Lars J. Jensen and Søren Brunak. 2012. Mining electronic health records: towards better research applications and clinical care. Nature Reviews Genetics 13, 395-405.*

## Linked databases - X-Road common bus



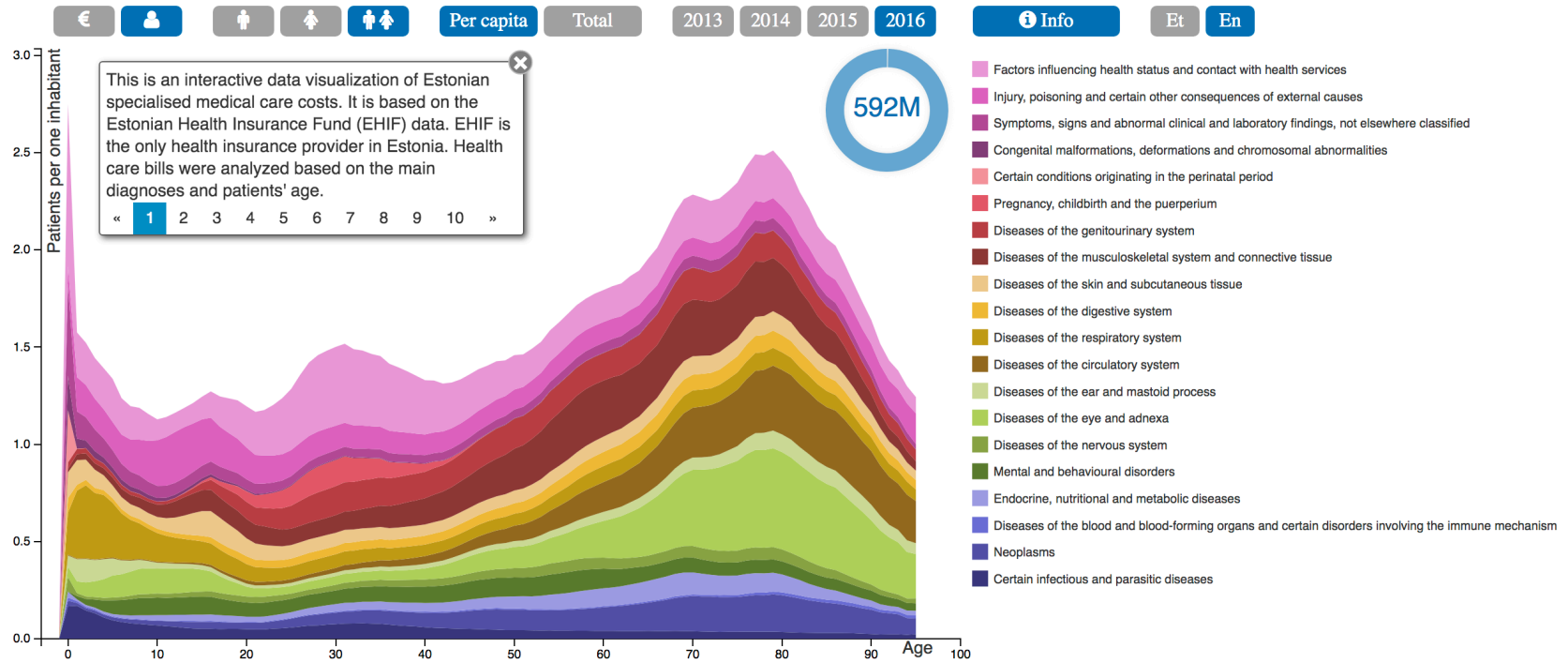
# Overview of the Estonian specialized medical care spending and patient counts by age and diagnosis



Source: Estonian Health Insurance Fund



## Overview of the Estonian specialized medical care spending and patient counts by age and diagnosis

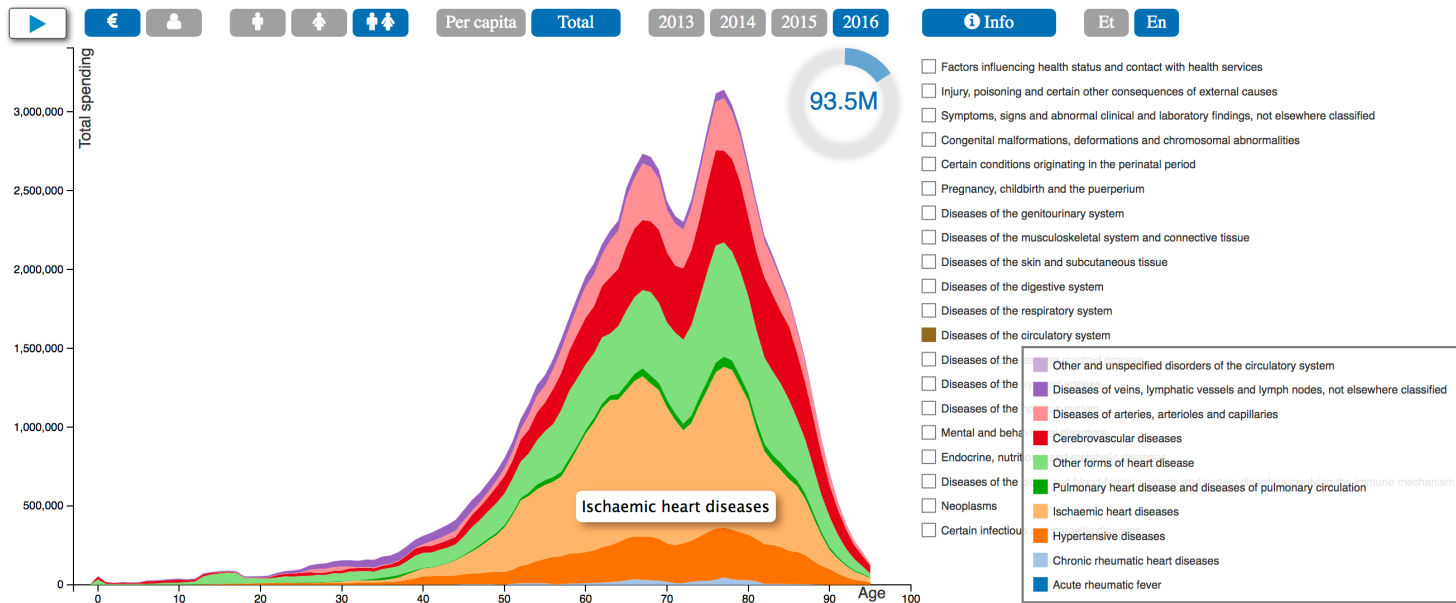


Source: Estonian Health Insurance Fund





## Overview of the Estonian specialized medical care spending and patient counts by age and diagnosis



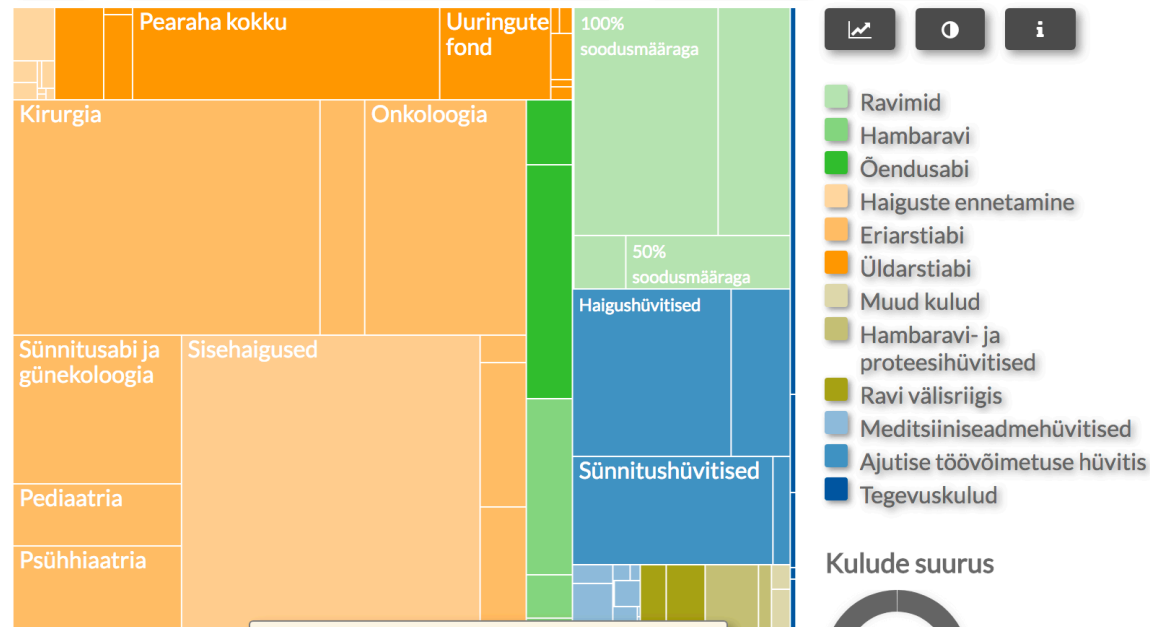
Source: Estonian Health Insurance Fund



# Haigekassa kulud 2016 (Auditeerimata)

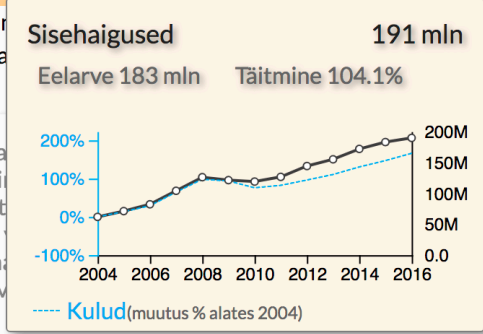
## Overall health expenditure in 2016: 1.06B

Kulud > Ravikindlustuse kulud > Tervishoiuteenuste kulud > Eriarstiabi kulud > Sisehaigused



- Ravimid
- Hambaravi
- Õendusabi
- Haiguste ennetamine
- Eriarstiabi
- Üldarstiabi
- Muud kulud
- Hambaravi- ja proteesihüvitised
- Ravi välisriigis
- Meditsiiniseadme hüvitised
- Ajutise töövõimetuse hüvitis
- Tegevuskulud

Kulude suurus



See graafik annab ülevaate sisehaiguste kulude muutusest aastast 2004. Graafik näitab, et sisehaiguste kulud on suurenenud oluliselt võrreldes 2004. aastaga, kasvades ligikaudu 104% võrra. See graafik annab ülevaate sisehaiguste kulude muutusest aastast 2004. Graafik näitab, et sisehaiguste kulud on suurenenud oluliselt võrreldes 2004. aastaga, kasvades ligikaudu 104% võrra.

# Nr of documents in E-Health DB 1.53 M individuals (06.03.2017)

Document type	Nr. of documents in E-health
Outpatient notes	15,323,163
Inpatient discharge summaries	1,643,296
Development assessment notices	34,308
Immunization side effects	6
Immunization notes	490,624
Growth notes	130,448
Ambulance Cards	340,274
Home nurse reports	5,634
Examination reports	122,945
Advisory notices	105,224
<b>Pointers to "pictures" in PACS</b>	<b>2,866,152</b>
Referrals	991,790
Referral responses	6,902,871

# Tehnoloogia: andmete visualiseerimine

Patsient 90787461031

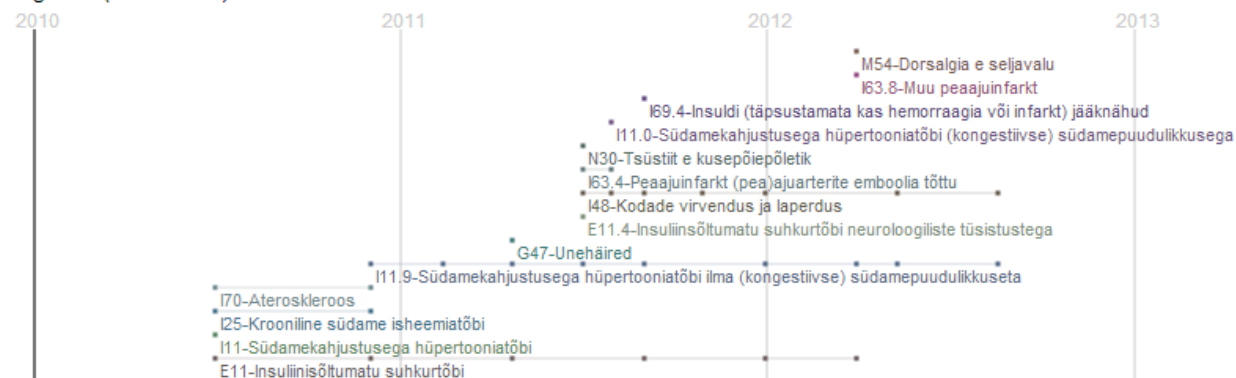
Sünniaasta: 1935

Genereeritud 2014-06-26 15:13:27.

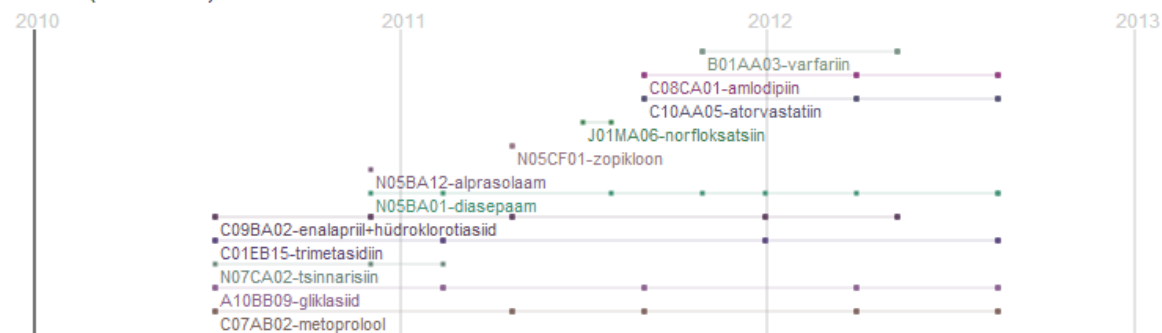
Visiidid (1 erinevat)



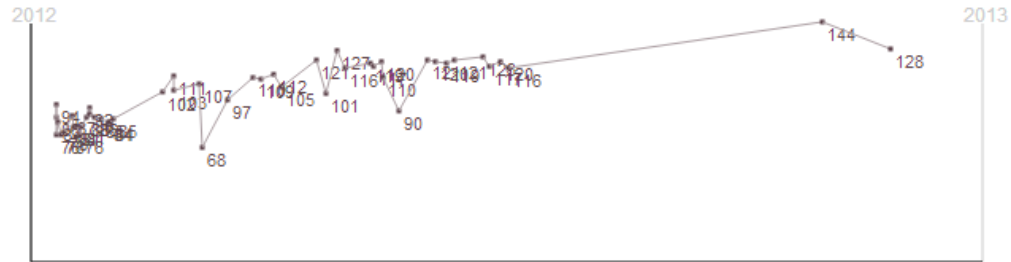
Diagnoos (14 erinevat)



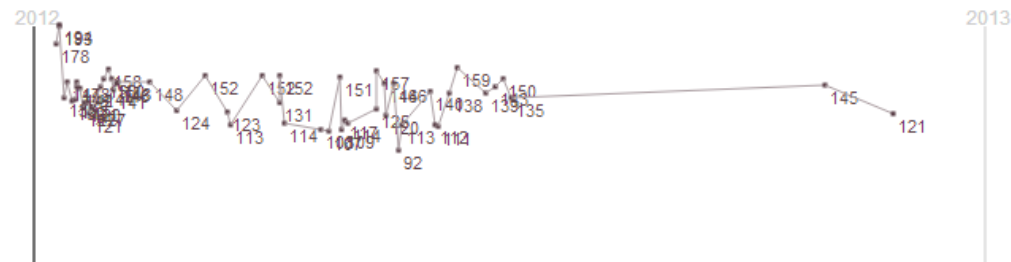
Ravimid (12 erinevat)



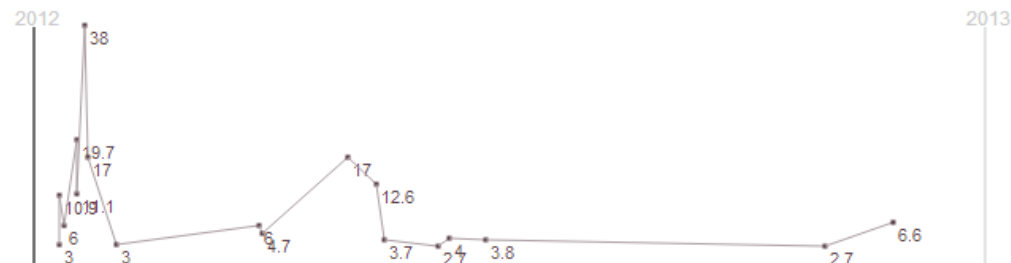
### HEMOGLOBIIN



### KREATINIIN



### VERESUHKUR

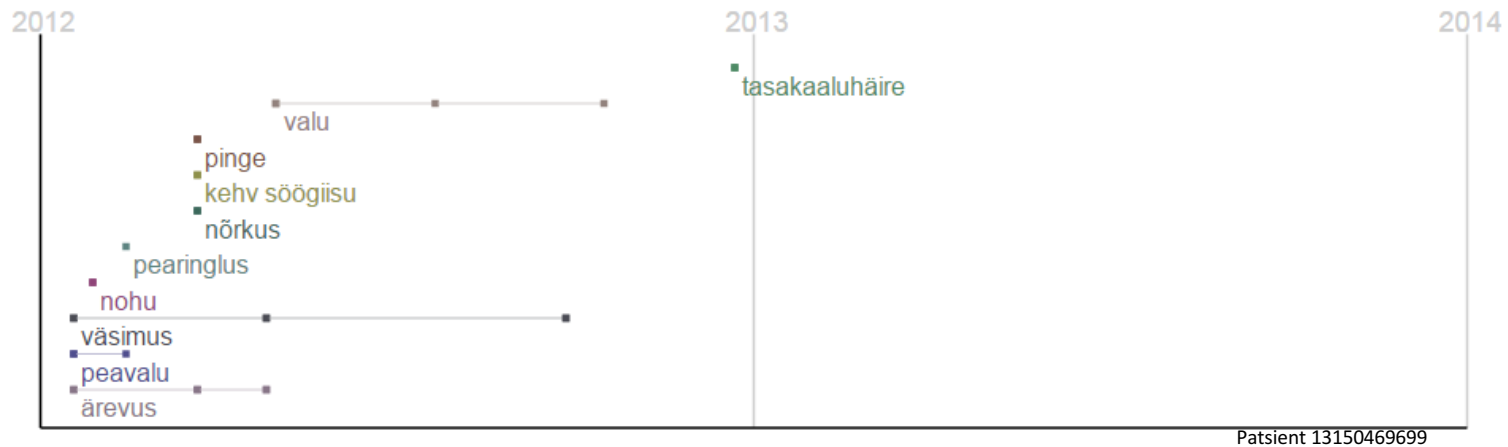


Patsient 16422115157

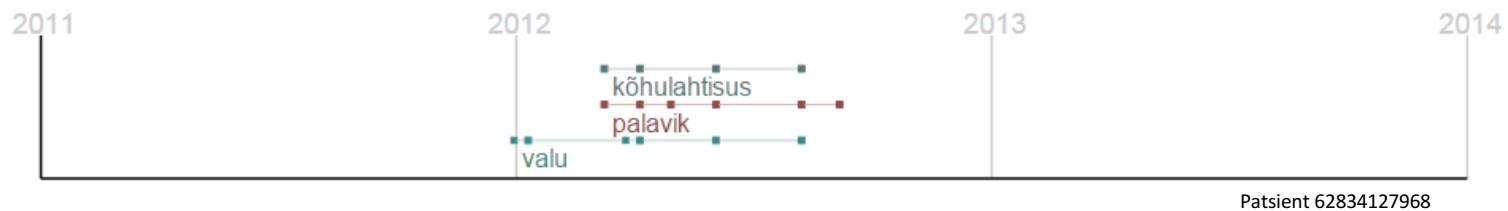
Tekstiosast  
eraldatud  
andmed!

# Näiteid kaebuste automaatsest eraldamisest ja visualiseerimisest

Kaebused (10 erinevat)



Kaebused (3 erinevat)



# Vabatekstides on veel kasulikku informatsiooni

Suunatud PA-lt , probleemiks **kõhuvalu**, mis ca aasta vältel olid 2 x nädalas, aga aprilli algul äge haigus **seedehäiretega**.

Sellest ajast kaebab iga päev, rohkem hommikul ärgates või enne und. **liveldust**, **oksendust** sel ajal **ei ole**.

Iste tavaliselt regulaarne, eile-täna **iste vedel**.

Anamneesis ka **peavalud**.

Saadetud **gastroskoopiasse**.

## **Kaebused:**

- Kõhuvalu
- Seedehäire
- liveldus (neg)
- Oksendamine (neg)
- Iste vedel
- Peavalu


## **Uuringud:**

- gastroskoopia

# Meditsiiniterminoloogia ühtlustamise vajadus

- Arstide keelekasutust iseloomustab variatiivsus
- Vabatekstiväljade kasulikkuse tõstmiseks on vaja see variatiivsus kirjeldada
- Näide epikriisi allergiablokist:

Penicillini  
Penicellin  
Penicilini  
Penicilliin  
PENICILLIN?  
Penicillini?  
PENICILLINID  
PENICLIIN  
Penicllini  
Penitsiliinile  
Penitsilliin  
...



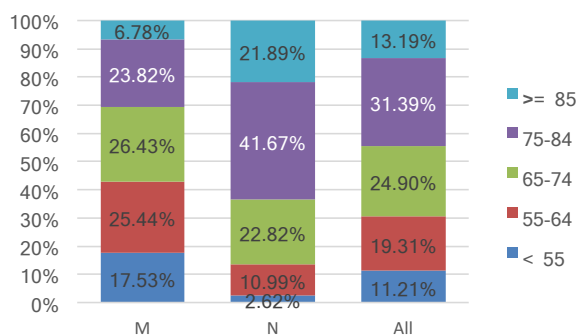
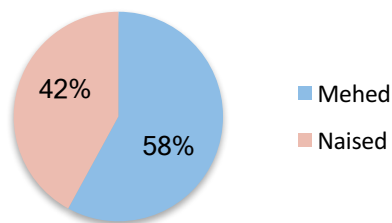
**Penitsilliin**



# Tehnoloogia ja kompetentsi võimekuse test (võrdlus müokardi-infarktiregistriga)

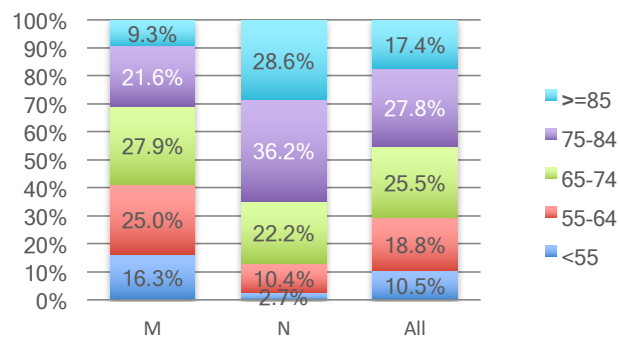
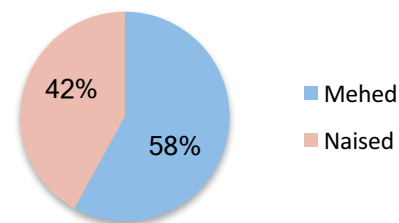
## E-TERVIS

- 2847 haigusjuhtu
- 2754 erinevat patsienti



## Müokardi-infarktiregister

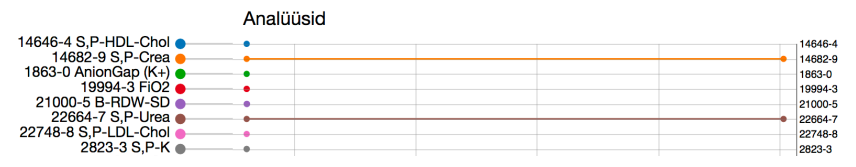
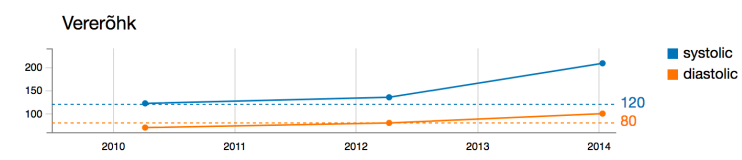
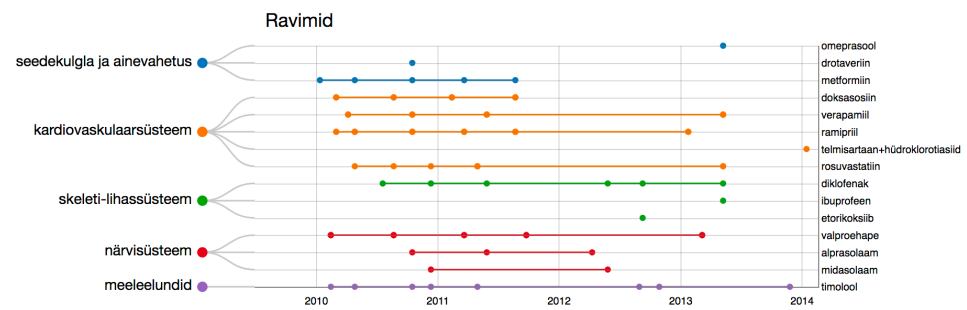
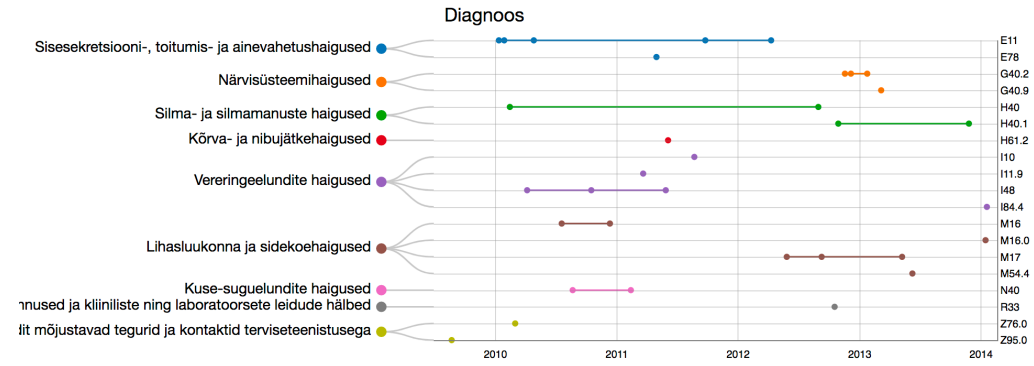
- 2879 infarktijuhtu



Tehnoloogia, mis suudab epikriisidest kätte saada  
praktiliselt kõik infarktijuhtumid!

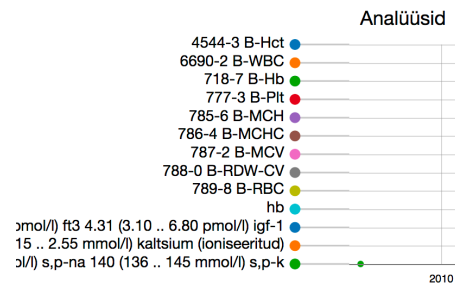
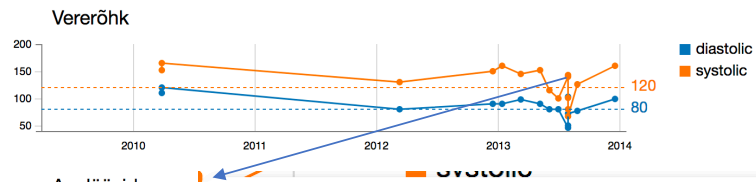
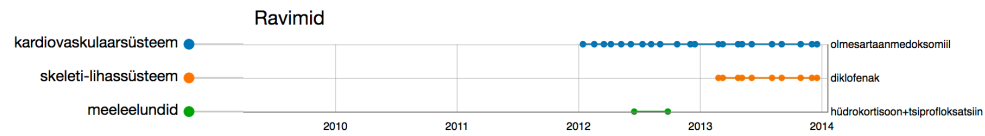
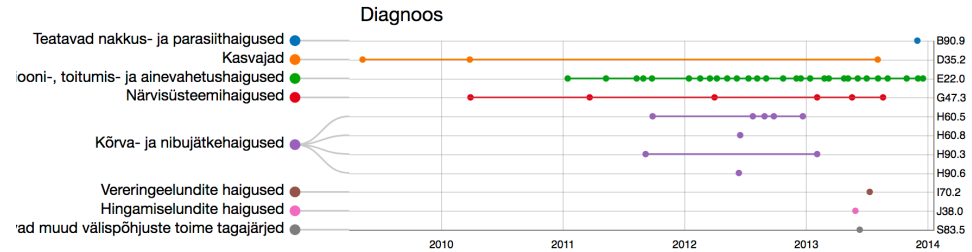
# STACC: Health data analysis for >1M individuals

- **Text Mining** (developing and Estonian NLP)
- **Machine learning** – predicting context, semantic meaning, semantic similarities,
- **Information extraction** from unstructured data
- Disease **co-morbidities** and **trajectories**
- **Quality indicators**
- **Observational data vs registries (e.g. infarctions)**
- IMI – EMIF-Platform: EGCUT data and OMOP
- Linking Genomics data to EHR and the needed information architectures

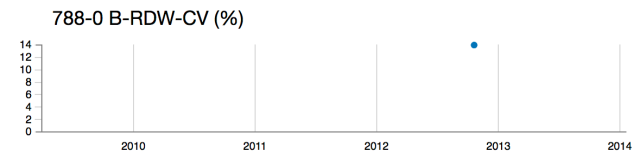


# 70905119764 M92

Soovitused puuduvad



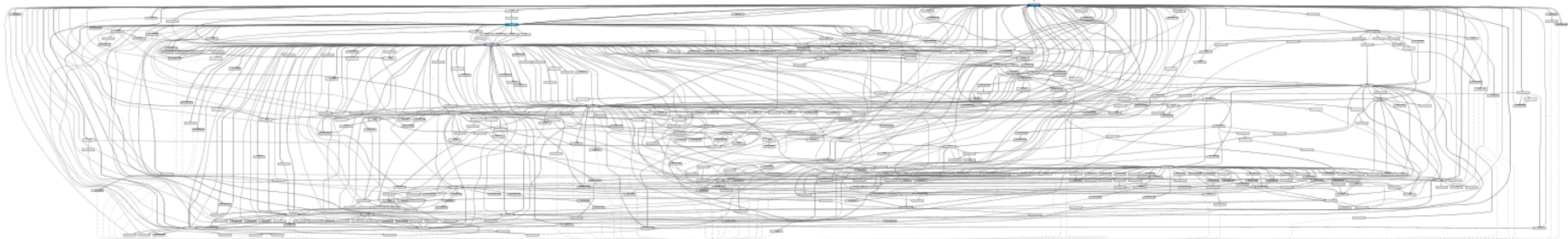
date	31.07.2013
context	väärtusi: 67/46-101/72 mmHg Täna RR=143/103 mmHg ; fr.
type	systolic
value	143



## Mida saab teaduslikus mõttes tuvastada:

- Statistiline juhtumite kokku loendamine (koos **tekstikaevega**)
- Haiguste **trajektoorid** – ajalised sõltuvused; ravi trajektoorid (näide)
- Haiguste **koos-esinemine** (co-morbidity) ja välistamine (riskid)
- Kvaliteedi-indikaatorid
- Kõrvalmõjud, kaebused, tulemused ...

## KOK (COPD) näitel “protsessid” (MSc töö)



*Figure 15: Figure showing the process model used in cluster 3 of diagnosis J44 with all events and edges displayed. This is what is called a “spaghetti model” in the field of process modeling.*

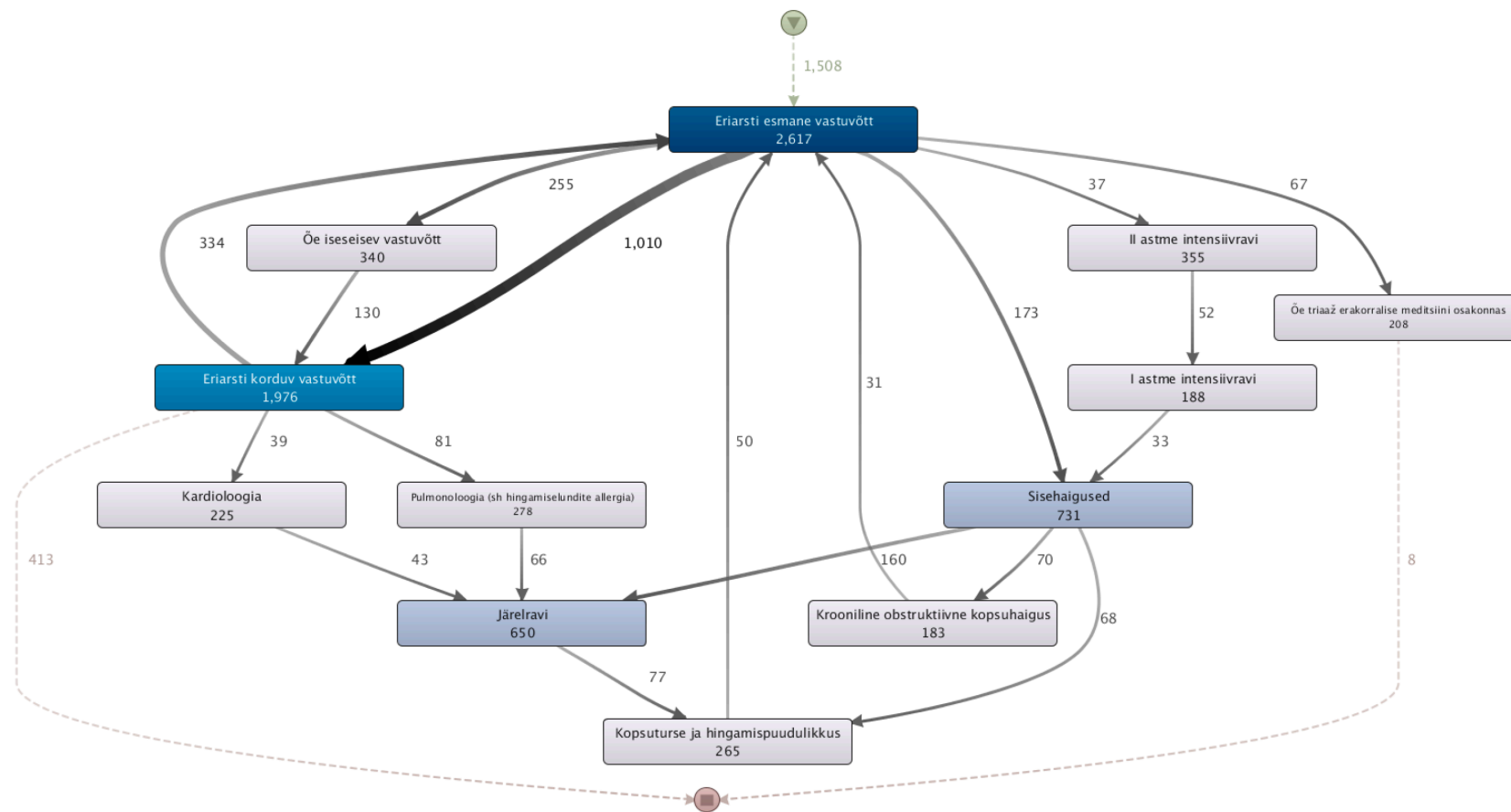


Figure 16: A process model found using Disco. Nodes limited to 2.6% edges to just the most significant ones.

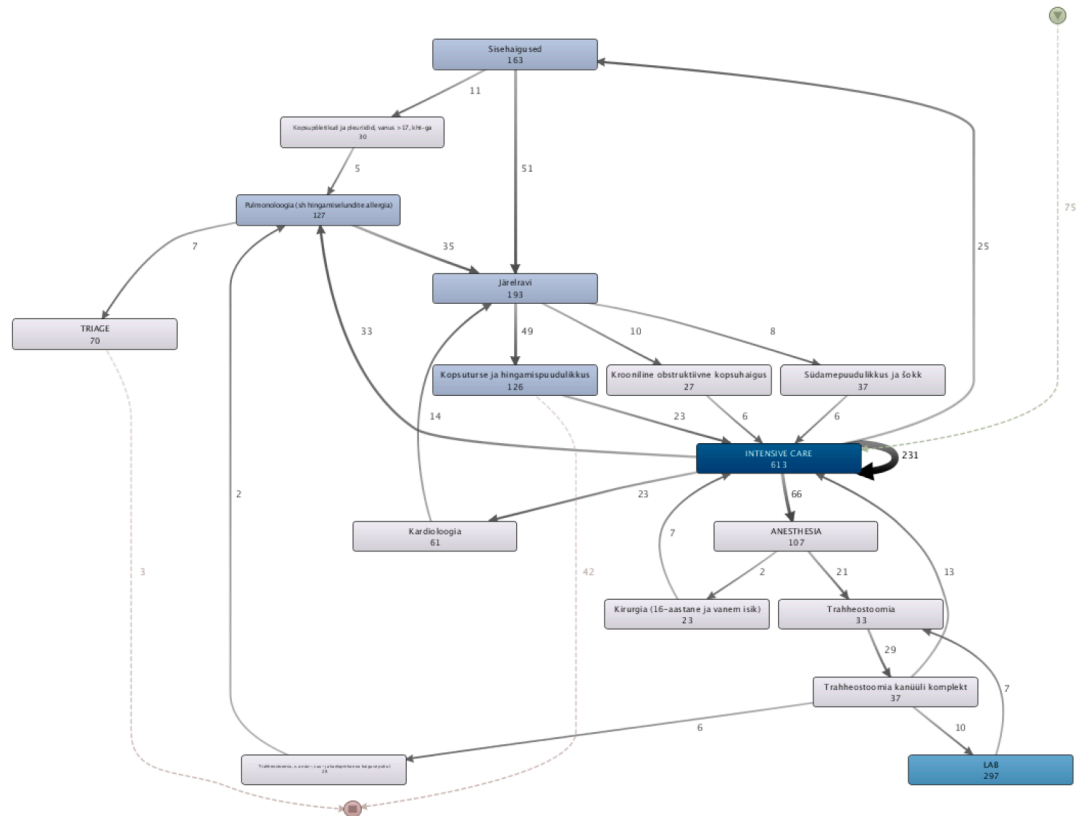
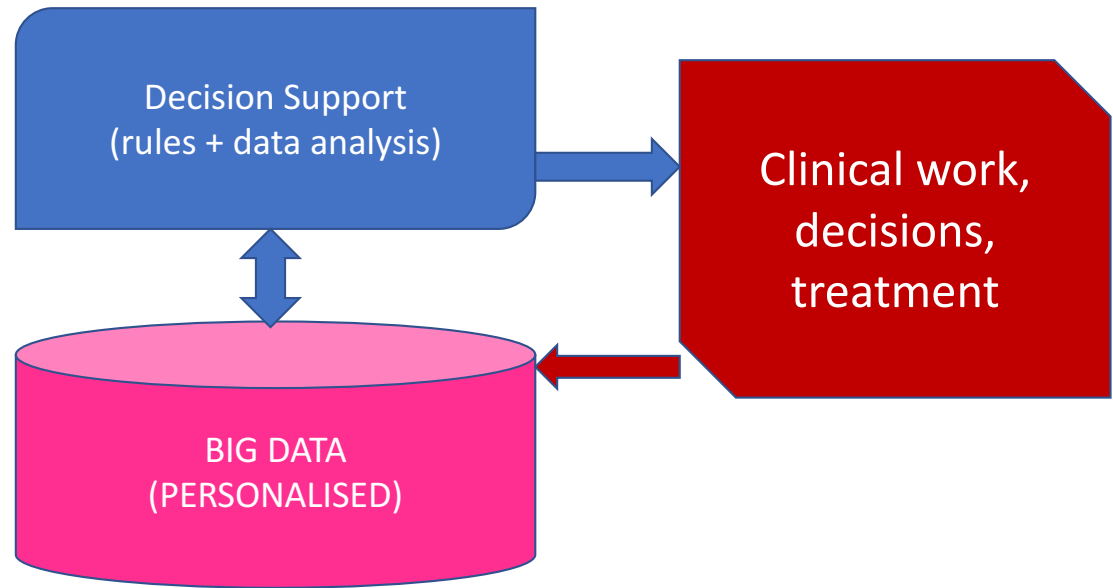
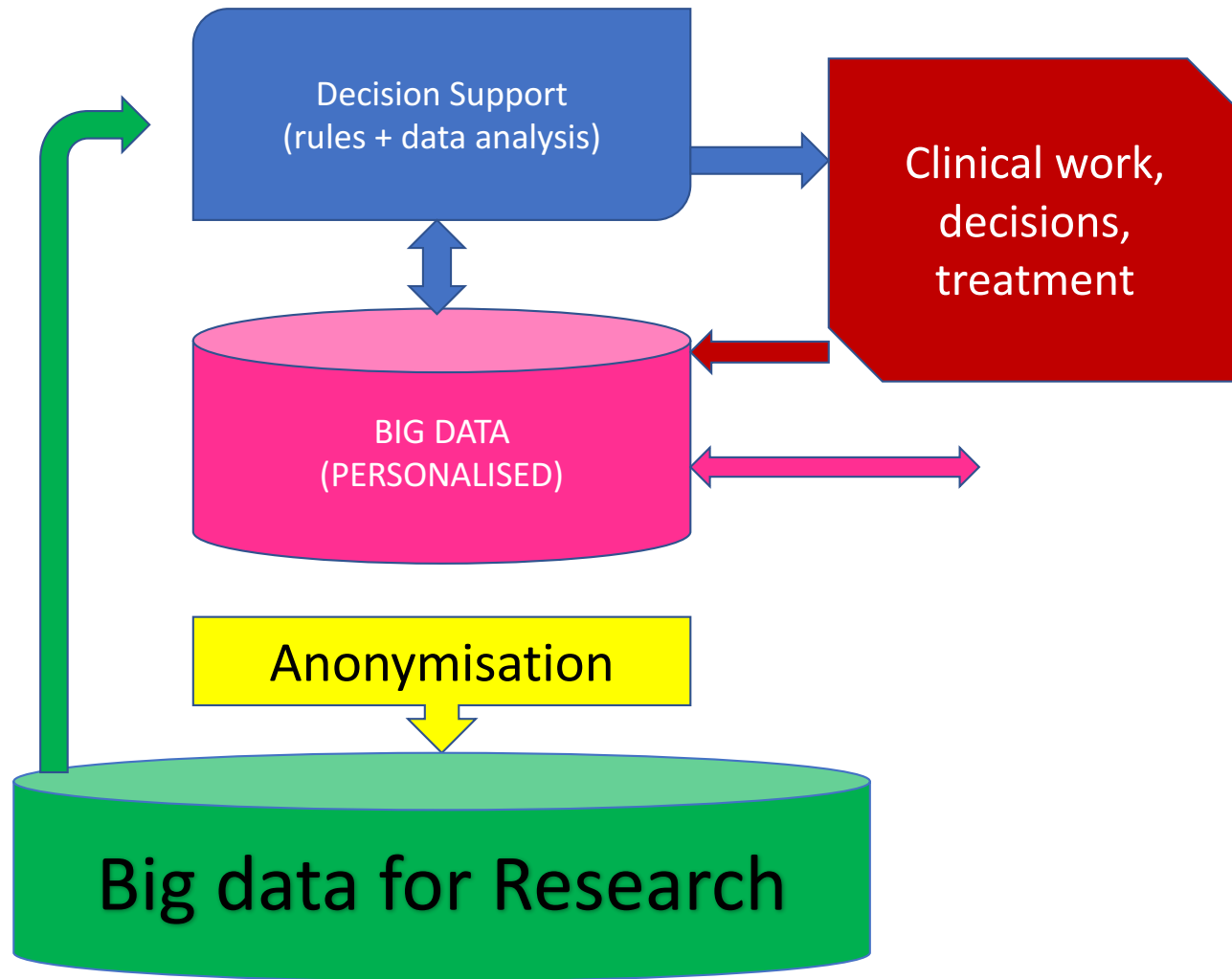


Figure 18: Clinical pathway for cluster 2 of diagnosis J44 on logs filtered with manually defined rules.

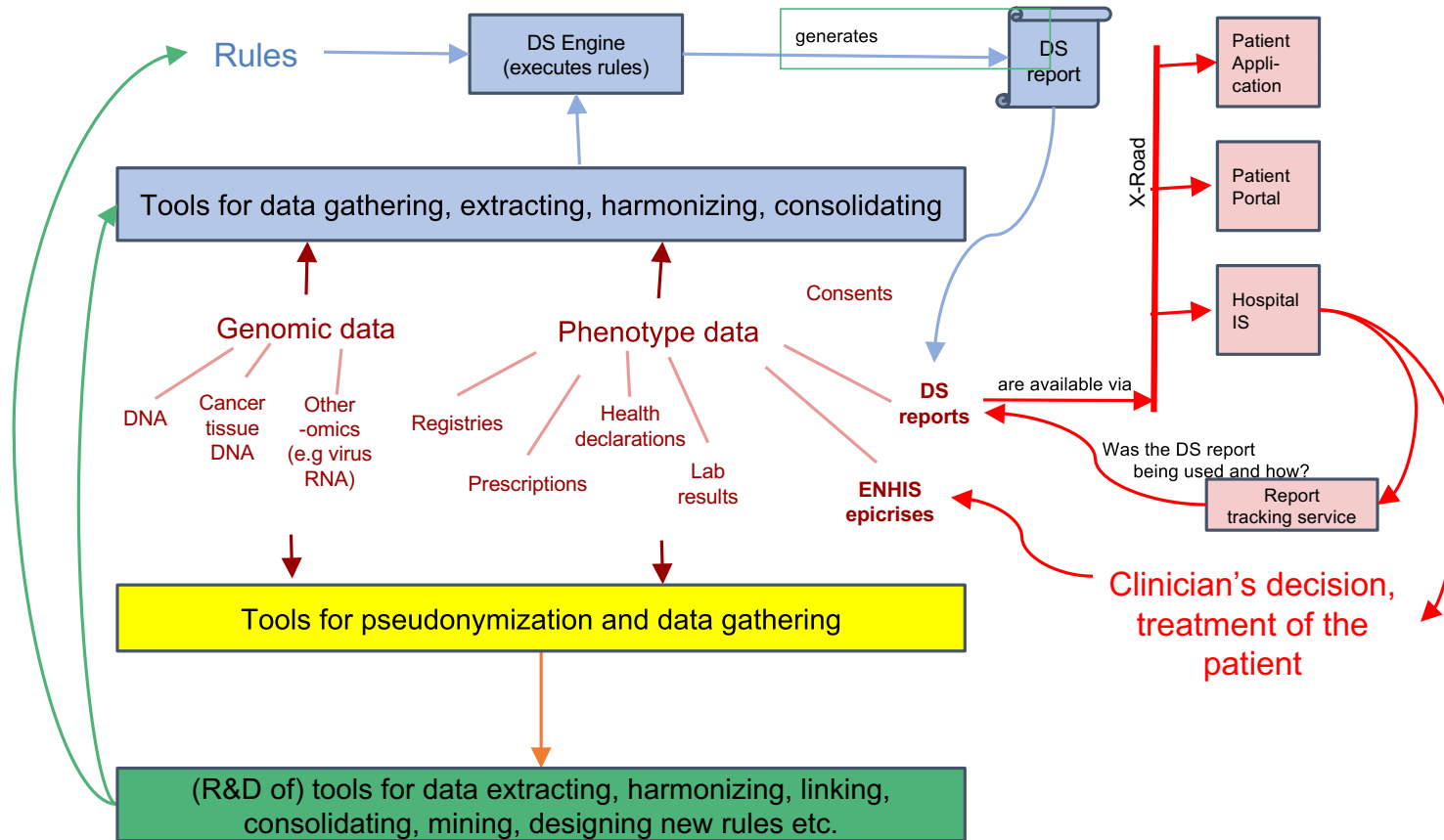
The pathways for this cluster are noticeably more complex, with more complications related to respiratory organs and heart. Interestingly, cancer seems to have much less importance in this cluster.





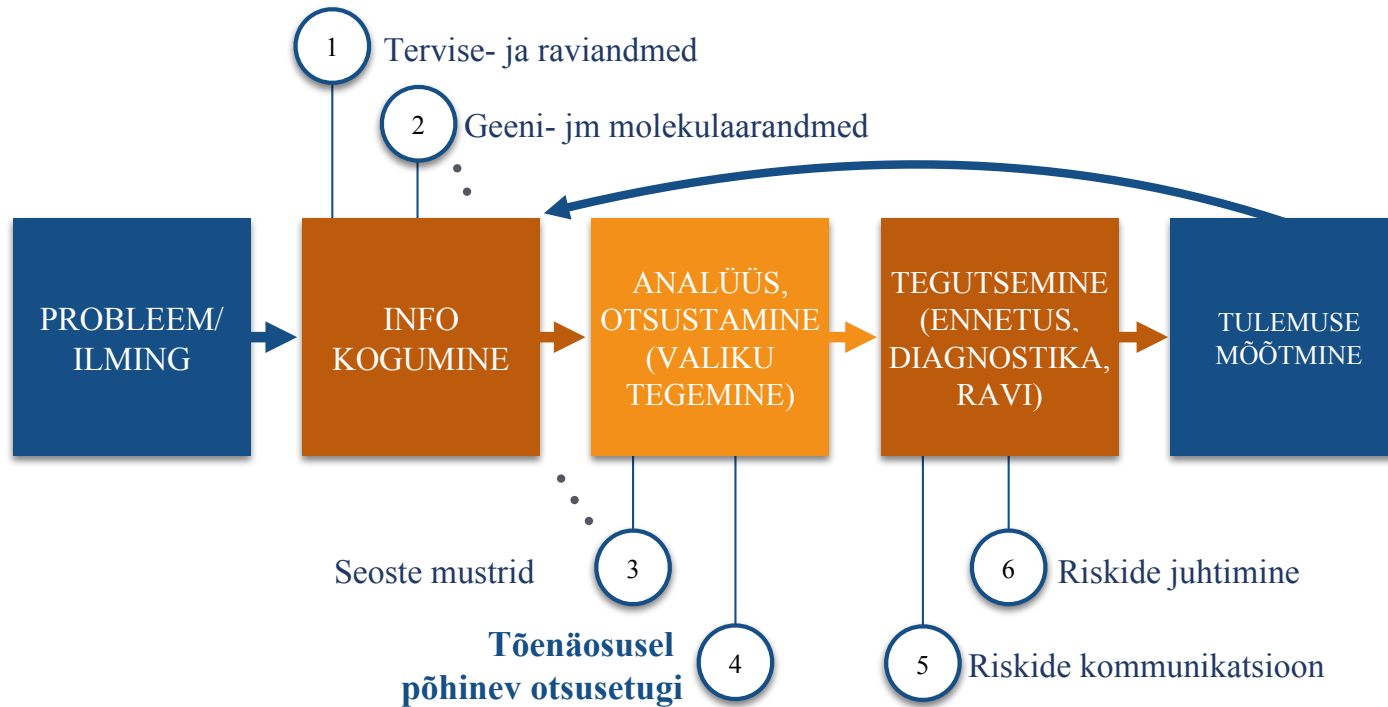


## What tools & data in each process



# Personaalmehitsiin =

Personalised medicine





# UNIVERSITY OF TARTU

STACC

Software Technology and  
Applications Competence Center



estonian genome center

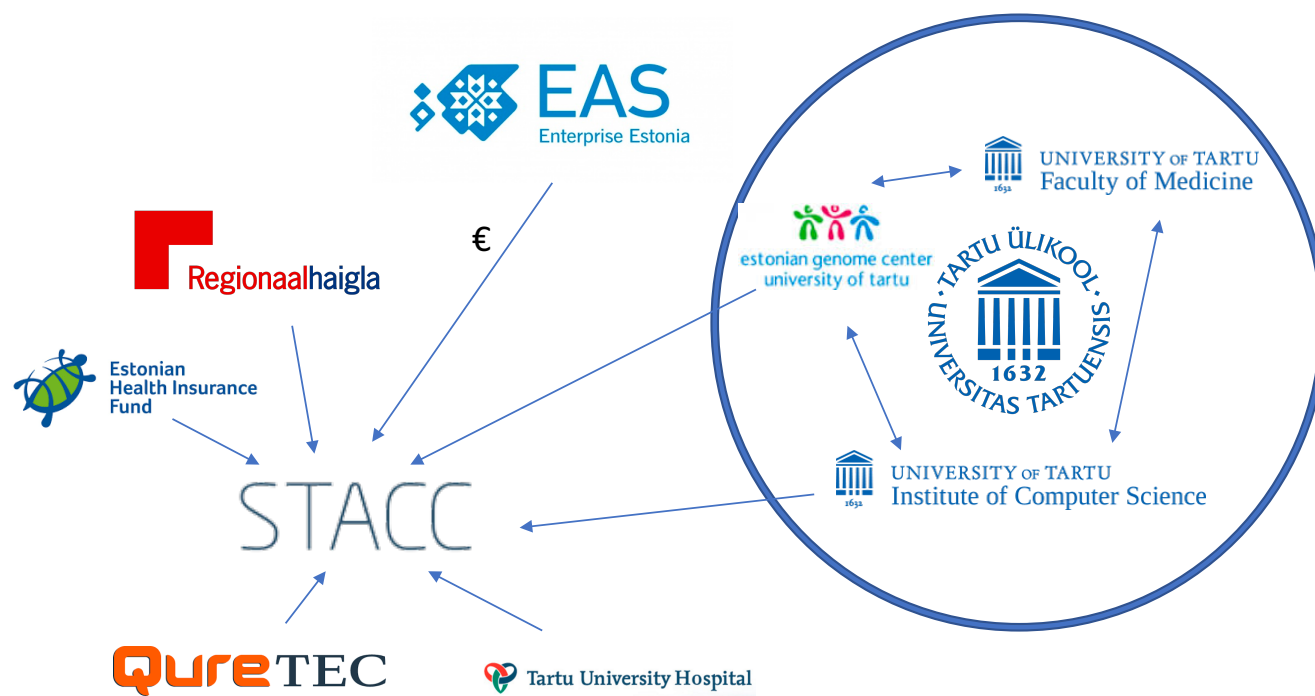


**BBMRI-ERIC**

Biobanking and  
BioMolecular resources  
Research Infrastructure

**eatris**

# Organisations:





**BI&IT**

**QureTEC**

**STACC**

Software Technology and  
Applications Competence Center