

# Aufbereitung und Analyse von Daten in phänotypischen und genotypischen Forschungs- datenbanken mit tranSMART



5. August 2016 | Berlin



## 10.15 Uhr Grundlegende

### Konzepte ■ Begrüßung

- 5' *Dr. Johannes Drepper (TMF), Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen), Matthias Löbe (Universität Leipzig)*
- Einführung in die tranSMART-Plattform
- 15'+3' *Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen)*
- tranSMART Architecture and Roadmap
- 45'+15' *Kees van Bochove (CIO The Hyve, tranSMART-RoadMap-Presentation at tranSMART-Meeting 2015)*
- Open Data Analytics – Gefahren durch p-Hacking
- 15'+3' *Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen)*

## 12.00 Uhr Mittagspause

## 13.00 Uhr Praktische Arbeit

- Klick-a-thon: Durchgehen eines vorbereiteten praktischen Szenarios durch die Workshopteilnehmer
- Lösung von Aufgaben in Eigenarbeit

## 13.00 Uhr Praktische Arbeit

- Klick-a-thon: Durchgehen eines vorbereiteten praktischen Szenarios durch die Workshopteilnehmer
- Lösung von Aufgaben in Eigenarbeit

## 14.15 Uhr Kaffeepause

## 14.45 Uhr Fortgeschrittene Themen

- tranSMART and beyond  
**20'+5'** *Reinhard Schneider / Sascha Herzinger*  
(Universität Luxemburg)
- SmartR – Dynamic Visual Analytics in TranSMART  
**20'+5'** *Reinhard Schneider / Sascha Herzinger*  
(Universität Luxemburg)
- Import und Repräsentation hochdimensionaler Daten  
**10'+3'** *Christoph Knell (Friedrich-Alexander-Universität Erlangen-Nürnberg)*

## 14.45 Uhr Fortgeschrittene Themen

- tranSMART and beyond

**20'+5'** *Reinhard Schneider / Sascha Herzinger*  
(Universität Luxemburg)

- SmartR – Dynamic Visual Analytics in TranSMART

**20'+5'** *Reinhard Schneider / Sascha Herzinger*  
(Universität Luxemburg)

- Import und Repräsentation hochdimensionaler Daten

**10'+3'** *Christoph Knell (Friedrich-Alexander-Universität Erlangen-Nürnberg)*

- Beispiele zu Daten aus der Routineversorgung (§ 21 KHEntG, Intensivmedizin) und aus der Forschung (Mikrobiom)

**10'+3'** *Benjamin Baum (Universitätsmedizin Göttingen)*

- Zusammenfassung und Ausblick

*Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen)*

## 16.15 Uhr Ende des Workshops

# Einführung in die tranSMART-Plattform

## TMF-Workshop

Aufbereitung und Analyse von Daten in phänotypischen und genotypischen  
Forschungsdatenbanken mit tranSMART

Berlin, 05.08.2016

Christian Bauer<sup>1</sup>, Benjamin Baum<sup>1</sup>, Ulrich Sax<sup>1</sup>

<sup>1</sup>Department of Medical Informatics, WG Infrastructure for Translational Research

GEFÖRDERT VOM



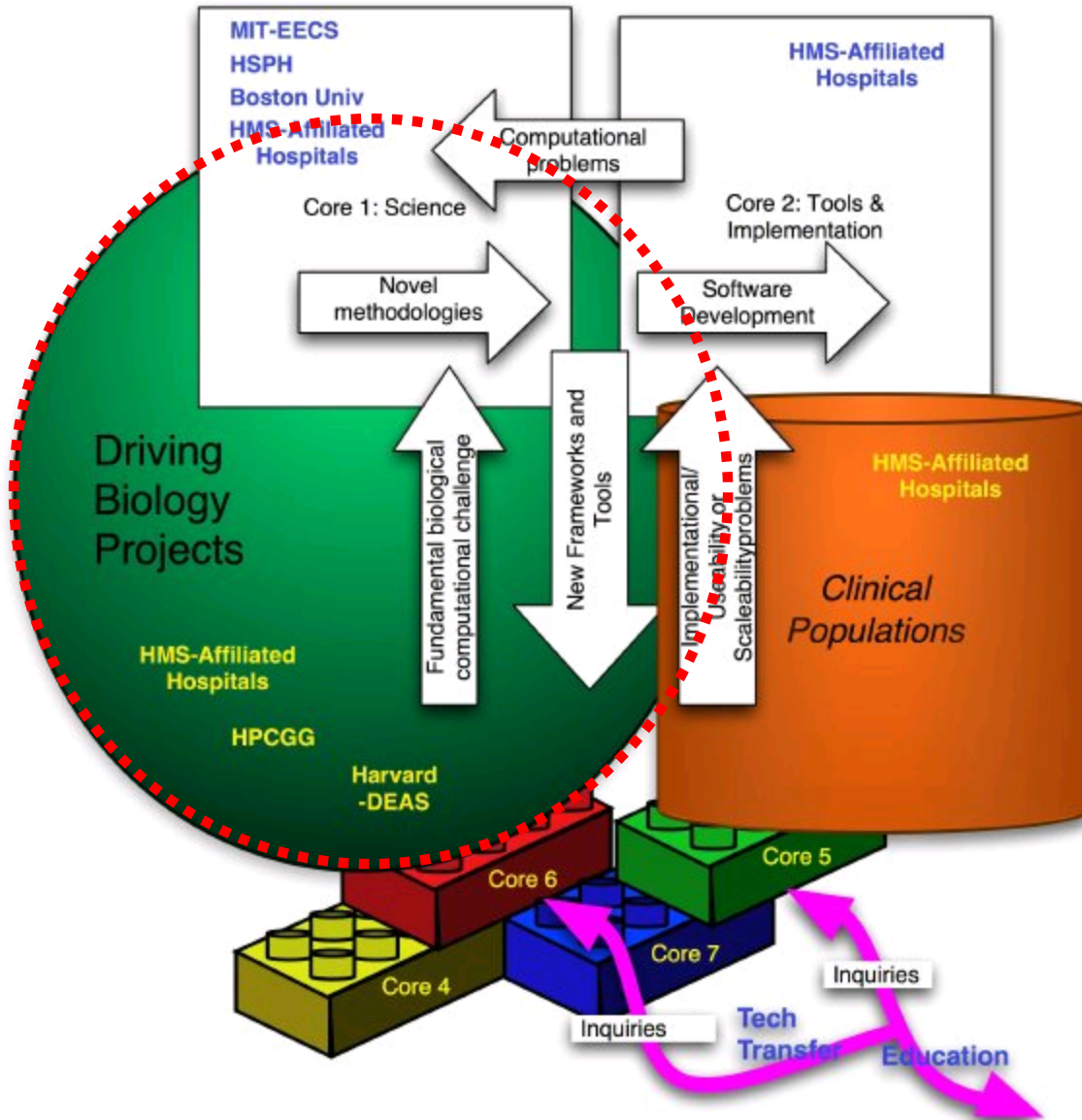
Bundesministerium  
für Bildung  
und Forschung

# Informatics for Integrating Biology and the Bedside (I<sup>2</sup>B<sup>2</sup>)

- ✱ IT-Infrastruktur zur Unterstützung genetischer bzw. genomischer Forschung mit dem Schwerpunkt der Integration klinischer und genomischer Daten
- ✱ Eines von vier Zentren in USA



# Informatics For Integrating Biology And The Bedside (I<sup>2</sup>B<sup>2</sup>)



## Vernetzte Forschung:

- ✳ IT-Infrastruktur zur Unterstützung genetischer bzw. genomischer Forschung
- ✳ Schwerpunkt: Integration klinischer und genomischer Daten
- ✳ Eines von vier Zentren in USA

www.I2B2.org

Proposal to Establish an NIH-Supported National Center for Biomedical Computing  
PI: Isaac Kohane (2004)

Collaboration between the i2b2 (Informatics for Integrating Biology and the Bedside) National Center for Biomedical Computing, the Chair of Medical Informatics, University of Erlangen-Nuremberg and the Department of Information Technology, Göttingen University Hospital

**MEMORANDUM OF UNDERSTANDING**

Date: 19.11.2009

  
\_\_\_\_\_  
Prof. Dr. Hans-Ulrich Prokosch

Head of Department

Chair of Medical Informatics  
University of Erlangen-Nuremberg

  
\_\_\_\_\_  
Prof. Dr. Ulrich Sax

Head of Department

Department of Information  
Technology  
Göttingen University Hospital

\_\_\_\_\_  
Isaac Kohane, M.D., Ph.D.

  
Director

i2b2 National Center for  
Biomedical Computing  
Brigham Women's Hospital Boston



Isaac Kohane, M.D., Ph.D.  
Director  
i2b2 National Center for Biomedical  
Computing  
77 Avenue Louis Pasteur  
Boston, MA 02115  
USA

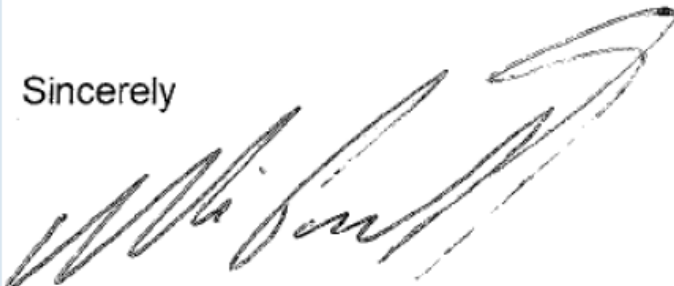
37099 Göttingen Briefpost  
Robert-Koch-Straße 40, 37075 Göttingen, Germany Adresse  
0551 / 39-13148 Telefon  
0551 / 39-8234 Fax  
usax@med.uni-goettingen.de E-Mail

15. Dezember 2009 Datum

**Letter of Support for the „i2b2 National Center of Biomedical Computing” call RFA-RM-09-002**

in response to the Call of the Department of Health and Human Services and the National Institutes of Health (NIH) on the topic: RFA-RM-09-002 National Centers for Biomedical Computing (U54).

Sincerely



Prof. Dr. Hans-Ulrich Prokosch



Prof. Dr. Ulrich Sax



Sebastian C. Semler

# TMF Special Issue: Unlocking Data for Clinical Research - The German i2b2 Experience

T. Ganslandt<sup>1</sup>; S. Mate<sup>2</sup>, Helbing K<sup>3</sup>, U. Sax<sup>3,4</sup>, H.U. Prokosch<sup>1,2</sup>

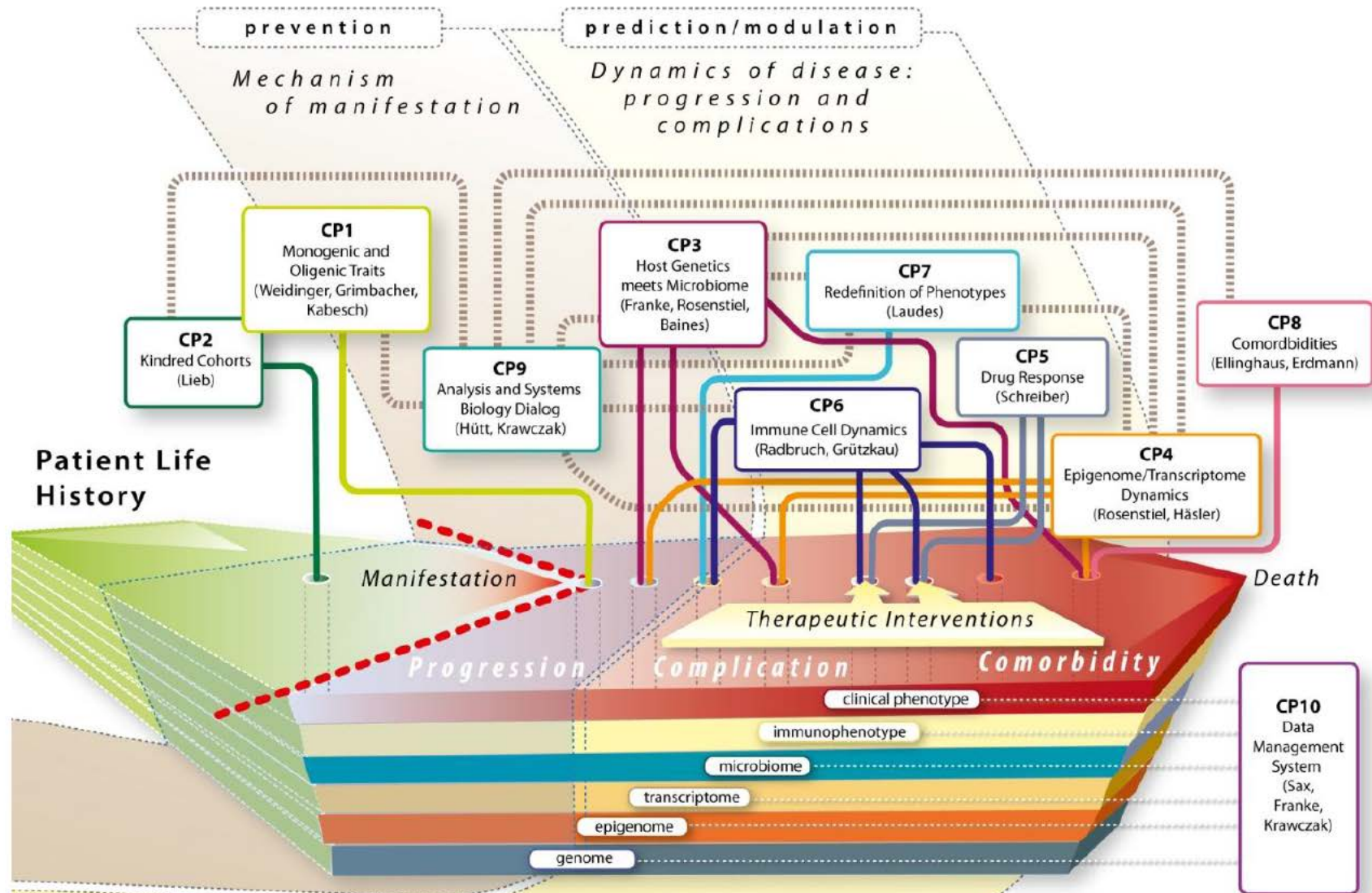
<sup>1</sup> Center for Medical Information and Communication, Erlangen University Hospital, Erlangen, Germany

<sup>2</sup> Chair of Medical Informatics, Friedrich-Alexander-University Erlangen-Nuremberg, Erlangen, Germany

<sup>3</sup> Department of Medical Informatics, University Medical Center Göttingen, Göttingen, Germany

<sup>4</sup> Division of Information Technology, University Medical Center Göttingen, Göttingen, Germany

# e:Med - sysINFLAME



Environment

**Epigenetics**

(Microarrays,  
ChIP-Sequencing)

**Genomics**

(WGS, WES,  
Panels)

**Trans-  
criptomics**

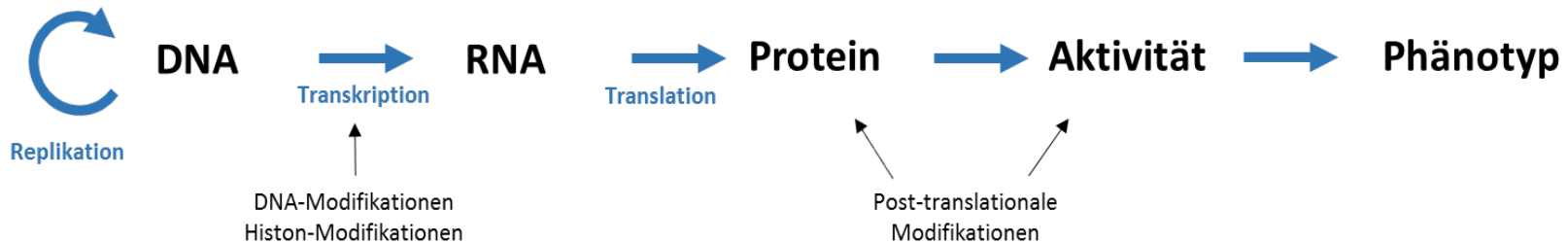
(Microarrays,  
RNA-Sequencing)

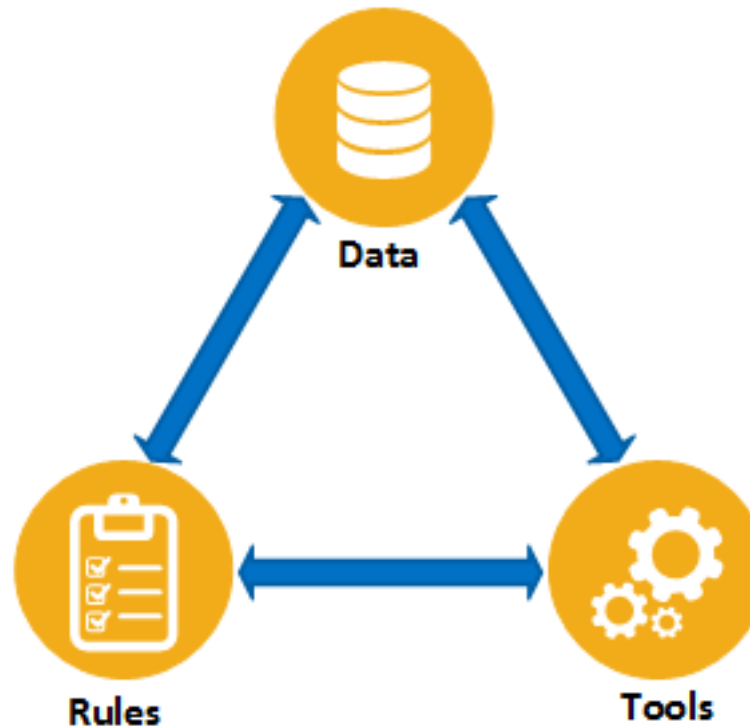
**Proteomics**

(LC-MS)

**Meta-  
bolomics**

(LC-MS)





# WP3: Use and Access Policy (UAP)

sysINFLAME\_Nutzungsordnung v1.5.docx [Schreibgeschützt] [Kompatibilitätsmodus] - Word

🔍 - 🗑️ ✕

## Nutzungsordnung für Daten im Verbund sysINFLAME

v1.5 vom 17.06.2015

1. Ziel dieser Nutzungsordnung ist die Regelung der wissenschaftlichen Verwendung der im Verbund sysINFLAME verfügbaren Daten. Gegenstand der Nutzungsordnung sind allein die in Anhang 1 beschriebenen, bereits existierenden bzw. noch zu erhebenden Datensätze (fortan: „Daten“). Diese Daten werden in einer gemeinsamen Infrastruktur des sysINFLAME-Verbundes zusammengeführt.

2. Jeder Wissenschaftler, der Daten in sysINFLAME einbringt oder nutzt, ist für die Einholung eines eventuell erforderlichen Ethikvotums selbst verantwortlich. In jedem Schritt der Datennutzung sind alle einschlägigen Datenschutzbestimmungen einzuhalten. Daten dürfen auch innerhalb des Verbundes nur in pseudonymisierter Form weitergegeben werden. Alle Einverständniserklärungen und Patientenlisten verbleiben beim Datenhalter.

3. Das Steering Committee (STC) des Verbundes sysINFLAME setzt ein dreiköpfiges Use and Access Committee (UAC) ein, das auf Antrag und nach transparenten Kriterien über die Nutzung der Daten aus Anhang 1 entscheidet. Das UAC wählt aus dem Kreis seiner Mitglieder einen Vorsitzenden. Die Entscheidungskriterien des UAC sind mit dem STC abzustimmen, schriftlich niederzulegen und orientieren sich an den Regeln der guten wissenschaftlichen Praxis. Mit der Unterzeichnung dieser Nutzungsordnung erkennen die Datenhalter in sysINFLAME die Entscheidungskompetenz des UAC über die Nutzung von Daten verbindlich an.

4. Jede Nutzung der Daten aus Anhang 1 bedarf grundsätzlich eines positiven Votums des UAC sowie der Genehmigung des STC, sofern die Nutzung nicht bereits durch eine andere Nutzungsordnung geregelt ist. In diesem Fall ist die Nutzungsordnung dem UAC zur Kenntnis zu geben und ihr entsprechend zu verfahren. Die Nutzung darf weder vom UAC noch vom Datenhalter unbillig verweigert werden. Streitfälle entscheidet das STC mit einfacher Mehrheit seiner Mitglieder.

5. Das UAC hat seine Aufgaben zügig zu erfüllen und spätestens 4 Wochen nach Eingang eines Nutzungsantrags über diesen zu entscheiden. Das UAC beschließt mit einfacher Mehrheit seiner Mitglieder.

6. Jeder Antrag auf Datennutzung enthält die folgenden Angaben:

- verantwortlicher Wissenschaftler
- Vertragspartner
- Titel des Vorhabens
- beabsichtigter Zeitraum
- Ziel des Vorhabens
- wissenschaftlicher Hintergrund
- Vorhabenbeschreibung
- Begründung der Machbarkeit
- Beschreibung möglicher Risiken
- zur Verfügung stehende (materielle und personelle) Ressourcen
- Einzelheiten zu den angeforderten Daten

7. Das UAC prüft den Antrag hinsichtlich folgender Kriterien:

- schlüssige wissenschaftliche Begründung
- Verfügbarkeit eines ausreichenden Datenbestandes (Fallzahl, Einwilligung)
- Erreichbarkeit des Ziels der Nutzung

8. Nach Prüfung des Antrags gibt das UAC schriftlich eine der folgenden drei Empfehlungen:

- Der Antrag soll genehmigt werden
- Der Antrag soll nur unter Auflagen oder nach bestimmten Modifikationen genehmigt werden
- Der Antrag soll abgelehnt werden

Eine Ablehnung oder Genehmigung unter Auflagen ist kurz zu begründen. Das STC entscheidet final über den Antrag und veranlasst ggf. die Bereitstellung der Daten.

# Translational research platforms integrating clinical and omics data: a review of publicly available solutions

*Vincent Canuel\**, *Bastien Rance\**, *Paul Avillach*, *Patrice Degoulet* and *Anita Burgun*

Submitted: 11th November 2013; Received (in revised form): 3rd February 2014

## Briefings in Bioinformatics

- International forum for researchers and educators
- For users of databases and analytical tools
  - Contemporary genetics
  - Molecular and systems biology
  - Provides practical help
- Impact factor: 9.6 (2014)

# Results

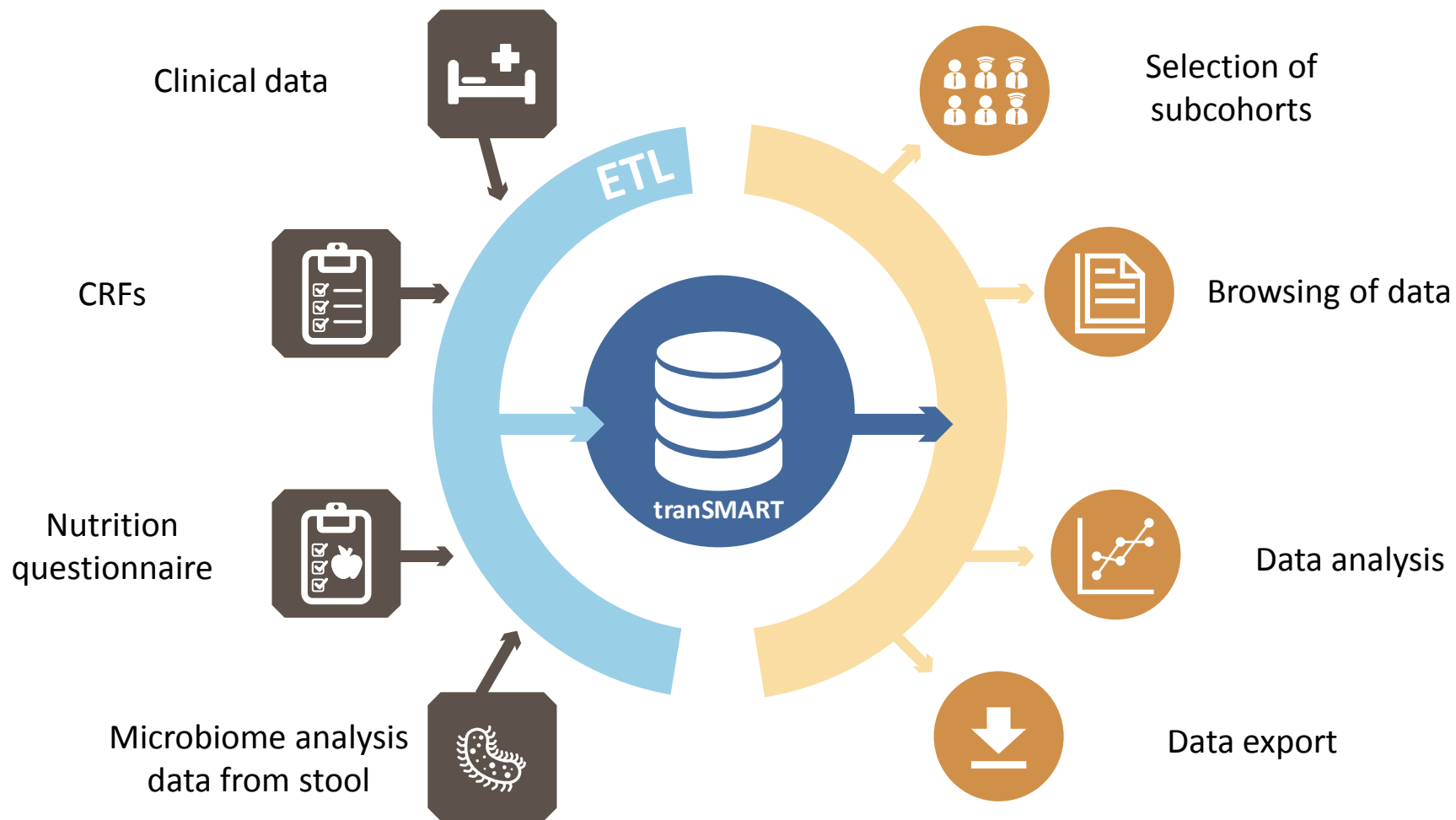
Platform	tranSMART
<b>Community</b>	
Institution/project initiators	Johnson & Johnson, USA
Funding	Initially Johnson & Johnson funded—now public-private consortia
Reference	Szalma <i>et al.</i> 2010
PMID	20642836
Software availability	Open source
Licensing	GPL v3
User mailing list or support URL	Yes <a href="http://transmartfoundation.org">http://transmartfoundation.org</a>
<b>Information content</b>	
<b>Clinical data</b>	
Demographics	Yes
Outcomes	Yes
Biological results	Yes
Images	No
Structured clinical research data	Yes
<b>'omics' data</b>	
mRNA expression	Yes
miRNA expression	No
SNPs	Yes
Copy number variations	Yes
DNA methylation	No
Protein/phosphoprotein expression	Yes
Structural rearrangements	No



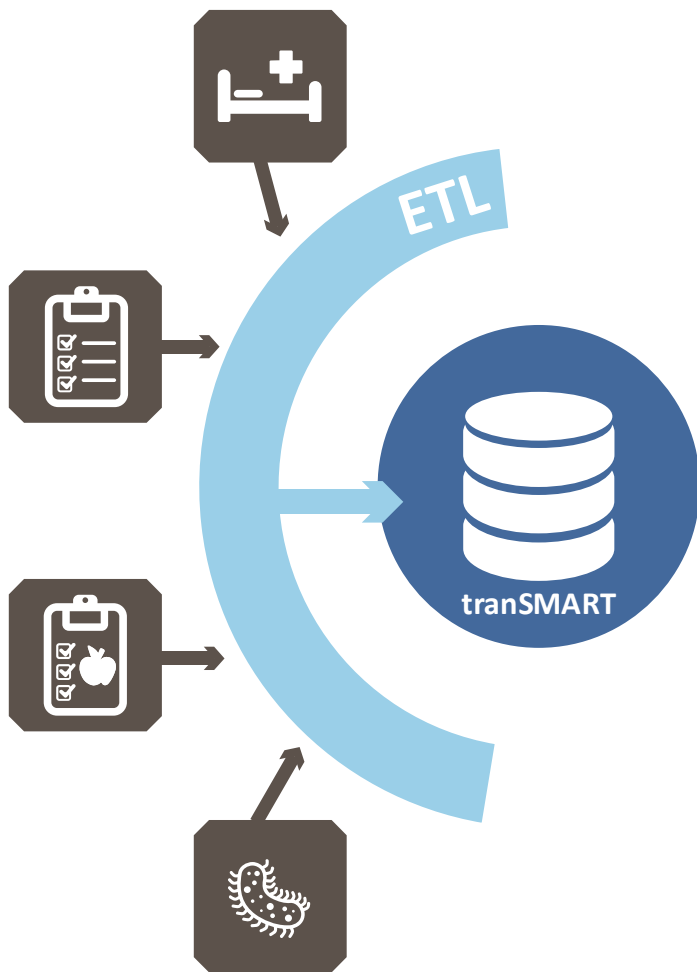
# Results

Platform	tranSMART
Privacy management environment	
Anonymization/de-anonymization	No
Analysis supports	
Statistical framework	R
Analytical features	GenePattern, Bioconductor Plink
Visualization tools	Haploview, IGV and output of analytical tools
Interoperability support	
Ontologies/standard terminologies	i2b2 ontology features
Collaborative environment	Yes
Secure environment	Yes
Support of multisite requests	No
APIs/web services interface	JBoss
System requirements	
Operating system	Linux
Database management system	Oracle II/PostgreSQL
Database type	Centralized
Software dependencies	SOLR, R, i2b2
Main programming language	Java/GRAILS
Server side	Tomcat/JBoss
Client-side interface	Web browser
Platform support	
Installation procedures	Yes
Configuration documentation	Yes
User documentation	Yes

# ShowCase tranSMART



# ShowCase tranSMART



sysINFLAME ETL pipeline: converting source system files to tranSMART format

ToDo: - ClinicalData: ...

preJob sysINFLAME

convert clinical data

updateClinicalData → listClinicalDataFiles → Iterate → init → Iterate → openMappingFile → row2 (Main) → filterCurrentMappingInfos → row1 (Filter)

OnSubjobOk

writeTranSMARTPatientDimensionFile

convert microbiome

updateMicrobiome → listMicrobiomeMappingFiles → Iterate → openMapping → mergeConceptsWithMapping → out1 (Main) → trimConceptCD → writeTranSMARTMicrobiome

OnSubjobOk

listMicrobiomeDataFiles → Iterate → copyDataFiles → openConcepts → row4 (lookup)

convert questionnaires

updateQuestionnaires → status → OnSubjobOk (order:1) → Iterate → listQuestionnaireFiles → openNutriensHierarchy → row6 (Main) → writeNutriensHierarchyToMap

OnSubjobOk (order:2)

status → OnSubjobOk (order:2) → openMapping → row7 (Main) → tJavaRow\_2 → row8 (Main) → tML\_2 → out10 (Main)

OnComponentOk

row9 (lookup)

2

# ShowCase tranSMART

tranSMART v1.2

Comparison Summary Statistics Grid View Advanced Workflow Data Export Export Jobs Workspace Genome Browser

Active Filters and  Clear

Save Subset Clear

Internal Studies (17)

- sysINFLAME (17)
  - Clinical Data (17)
    - B-Cell panel (16)
    - Dauermedikation (16)
    - Diagnosen (17)
    - Familien und Sozialanamnese (17)
      - 12 months (10)
      - affected (13)
        - affected covid flag (3)
          - abc True (3)
        - affected none (8)
        - affected text (2)
        - affected unknown (2)
        - childdeath (13)
          - childdeath none (11)
            - childdeath unknown (1)
            - childdeath yes (1)
        - cosanguinity (11)
          - cosanguinity none (10)
          - unknown (1)
        - employment status (12)
        - family history (9)
        - GDB (13)
        - learned profession (10)
        - social anamnesis (9)
      - Gastroskopie (13)
      - Imaging Ultrasound Abdomen Ascites
      - Immunglobuline (17)
      - Immunglobulinsubstitution (17)
      - Impfantworten (17)
      - Infektanamnese (17)
      - Koloskopie (13)
      - Lymphozyten (17)
      - Special t cell (7)
      - Spezielle Immunologie (15)
      - Stammdaten (17)
      - Stuhlanamnese (17)
      - t cell panel (14)
      - Untersuchungsbefunde (17)
    - Demographics (17)
    - Ernaehrungsfragebogen (17)

Subset 1

Exclude X

...childdeath none

AND Exclude X

AND Exclude X

AND Exclude X

Subset 2

Exclude X

AND Exclude X

AND Exclude X

AND Exclude X

Ready

# ShowCase tranSMART

tranSMART v1.2

Comparison Summary Statistics Grid View Advanced Workflow Data Export Export Jobs Workspace Genome Browser

Browse Analyze Gene Signature/Lists GWAS Upload Data Utilities

Active Filters and  Clear

Navigate Terms Across Trial

Internal Studies

- sysINFLAME (17)
  - (17)
    - Clinical Data (17)
      - B-Cell panel (16)
      - Dauermedikation (16)
      - Diagnosen (17)
      - Familien und Sozialanamnese (17)
        - 12 months (10)
        - affected (13)
          - affected covid flag (3)
            - abc True (3)
          - affected none (8)
          - affected text (2)
          - affected unknown (2)
        - childdeath (13)
          - childdeath none (11)
          - childdeath unknown (1)
          - childdeath yes (1)
        - cosanguinity (11)
          - cosanguinity none (10)
          - unknown (1)
        - employment status (12)
        - family history (9)
        - GDB (13)
        - learned profession (10)
        - social anamnesis (9)
      - Gastroskopie (13)
      - Imaging Ultrasound Abdomen Ascites
      - Immunglobuline (17)
      - Immunglobulinsubstitution (17)
      - Impfantworten (17)
        - vacpeptideresponse (17)
          - flag (17)
            - abc Checked (1)
            - abc Indeterminate (13)
            - abc Unchecked (3)
          - text (1)
            - abc unter SCIG (1)
        - vacpnpsresponse (17)
      - Infektanamnese (17)
      - Koloskopie (13)
      - Lymphozyten (17)
      - ... (7)

Analysis of ...Impfantworten\vacpeptideresponse\flag for subsets:

Subset 1

Concept	Count of Observations
Checked	1
Indeterminate	9
Unchecked	1

Category	Subset 1 (n)	Subset 1 (%n)
Checked	1	9.1%
Indeterminate	9	81.8%
Unchecked	1	9.1%
Total	11	100%

Summary Statistics

Query Summary for Subset 1

(Internal Studies\Internal Studies\sysINFLAME\CP01 (Freiburg) CVID enteropathy\Clinical Data\Familien und Sozialanamnese\childdeath\childdeath none)

Subject Totals		
Subset 1	Both	Subset 2
11	0	0

Histogram of Age

Subset 1	
Mean:	41.64
Median:	44
IQR:	17
SD:	11.6
Data Points:	11

Comparison of Age

Sex

# ShowCase tranSMART

tranSMART v1.2 All

Comparison Summary Statistics Grid View Advanced Workflow Data Export Export Jobs Workspace Genome Browser

Browse Analyze Gene Signature/Lists GWAS Upload Data Utilities

Active Filters and Clear

Analysis

Analysis: Correlation Analysis

Cohorts: Subset 1: (Internal Studies/sysINFLAME Demographics\ )

Variable Selection

Drag two or more numerical concepts from the tree into the box below that you wish to generate correlation statistics on.

...Alkohol [g/Tag]  
...Kochsalz [g/Tag]

Clear

Run Correlation   
Correlation Type

Run

Correlation Table (p-values on top right half, correlation coefficient on bottom left)

	Alkohol g Tag	Kochsalz g Tag
Alkohol g Tag	1.00000	0.44988
Kochsalz g Tag	0.19643	1.00000

Alkohol..g.Tag.

Kochsalz..g.Tag.

# ShowCase tranSMART

tranSMART v1.2 All

Comparison Summary Statistics Grid View **Advanced Workflow** Data Export Export Jobs Workspace Genome Browser Browse Analyze Gene Signature/Lists GWAS Upload Data Utilities

Analysis

Analysis: SYSINFLAME: Diversity

Cohorts:  
Subset 1: (\Internal Studies\sysINFLAME\{

**Variable Selection**

**Microbiome**  
Select a microbiome on which you would like to do the analysis and drag it into the box. This variable is required.

**Variable**  
Select the appropriate categorical variables and drag them into the box. For example, "SEX". This variable is required.

...Achromobacter\  
...Acidovorax\  
...Acinetobacter\  
...Actinomyces\  
...Aeromonas\  
...Alistipes\  
...Anaerococcus

...Laendliche Gegend mit vereinzeltten Haeusern oder Gehoefen\  
...Mittelstadt (20.000 bis kleiner 100.000 Einwohner)

Clear Clear

Correlation Type: Shannon Run

**Alpha-Diversity**

Click on the image to open it in a new window as this may increase readability.

Alpha-diversity Shannon-Index

Group	Min	Q1	Median	Q3	Max
Group 1 (Left)	2.53	2.55	2.56	2.68	2.70
Group 2 (Right)	2.35	2.38	2.39	2.45	2.68

Shannon-Index

transmart01.mi.med.uni-goettingen.de/transmart/datasetExplorer/index#

tranSMART v1.2

Comparison Summary Statistics Grid View Advanced Workflow Data Export Export Jobs Workspace Genome Browser

Analysis

Analysis: sysINFLAME: Outlier

Cohorts: Subset 1: (\Internal Studies\sysINFLAME\ )

**Variable Selection**

Drag numerical concepts from the tree into the box below that you wish to generate outlier statistics on.

...X01..Groesse.in.cm

Clear

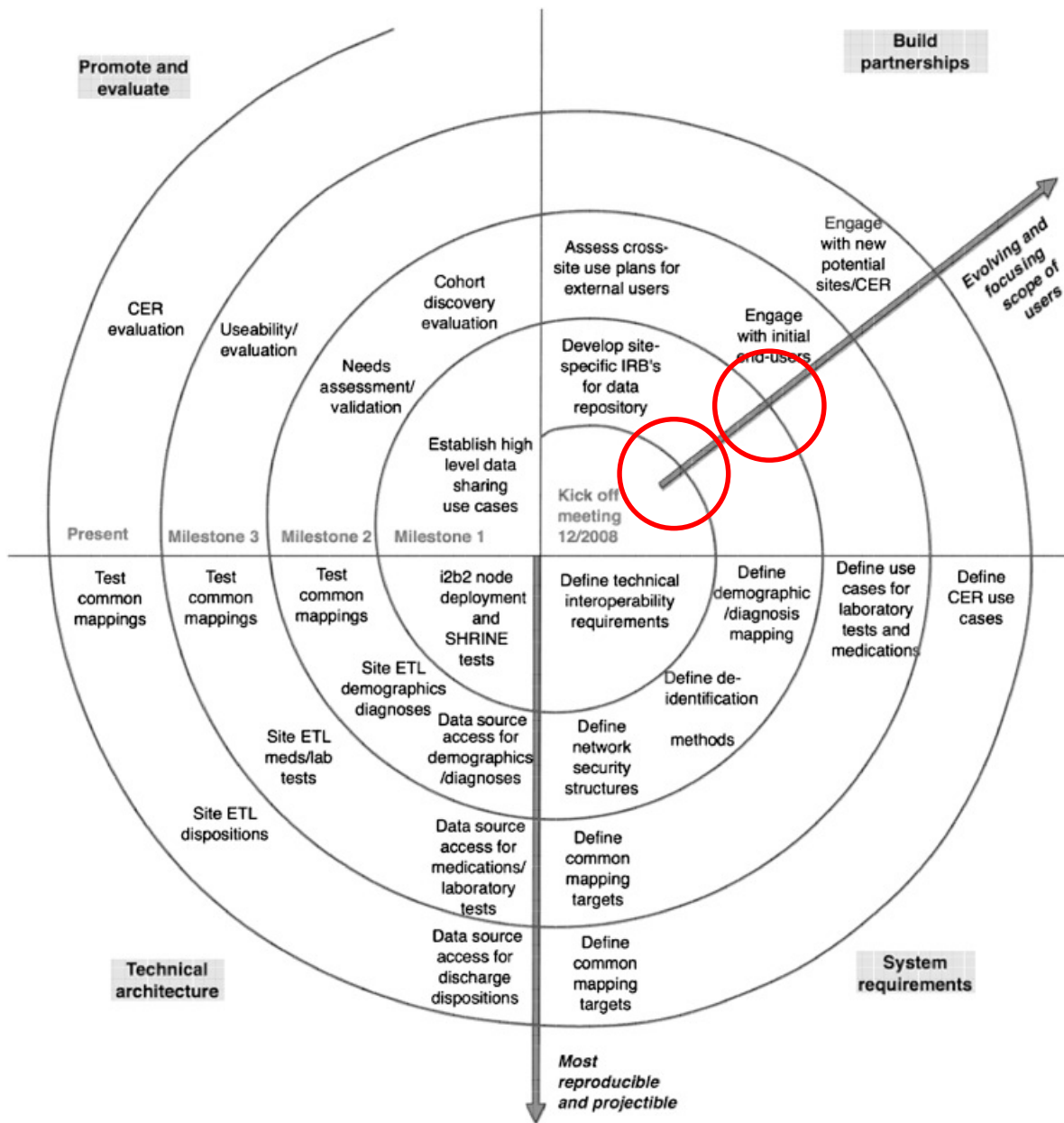
Run

**Outlier**

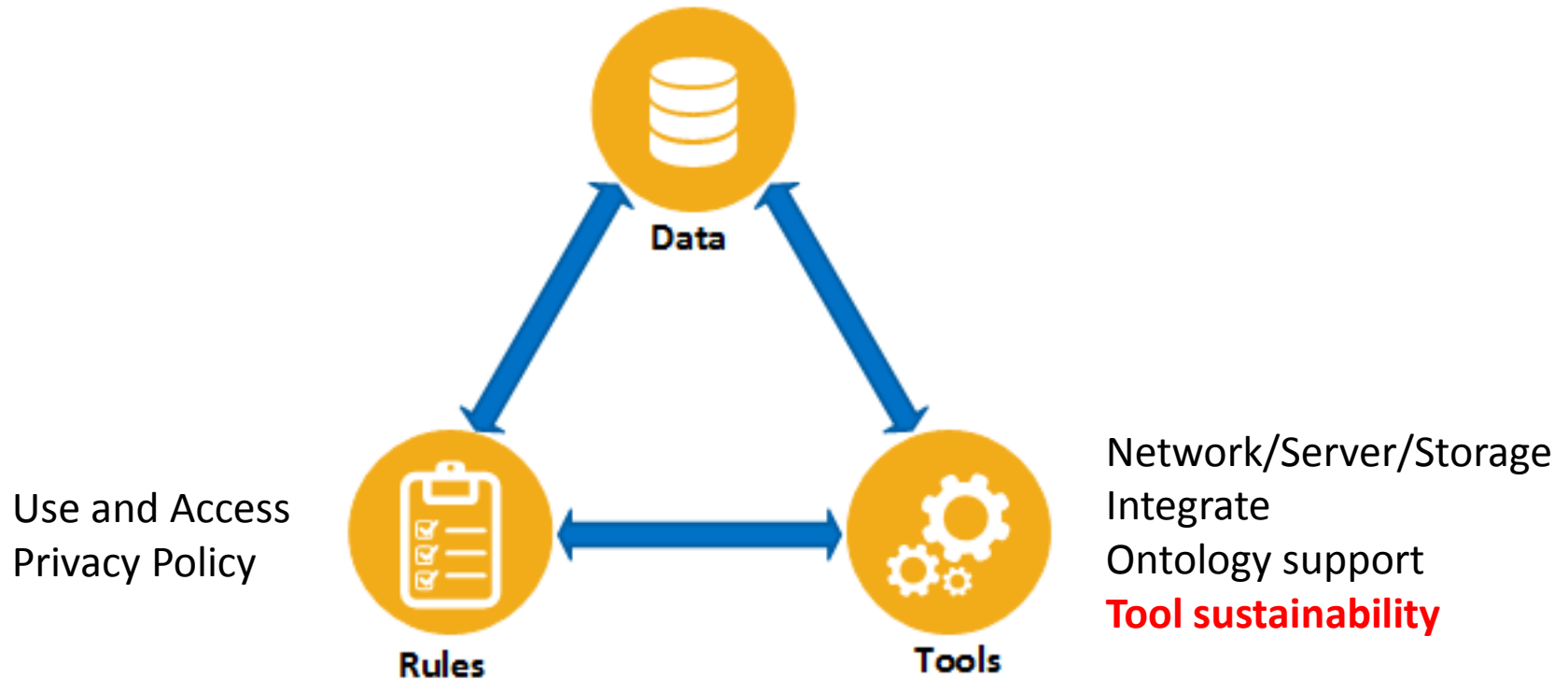
Spaltenname	Samples	Extremwerte	Methode	Bemerkung
X01..Groesse.in.cm	1001	145	Grubbs	

Download raw R data





### Sifting and Cleaning Phenotyping!



## i2b2 und tranSMART

- I2b2-TMF Kooperation seit 2009, Toolkits
- Entwicklung tranSMART auf älterem i2b2-Stand
- „Cousins“
- tranSMART Foundation ermöglicht externe Entwicklungsaufträge, schnelle Releasezyklen

### ABER

- Modifier fehlen (!!!) derzeit in tranSMART
- Re-Integration angestrebt (Avillach, Murphy)
- Zeitschiene?

## 10.15 Uhr Grundlegende

### Konzepte ■ Begrüßung

- 5' *Dr. Johannes Drepper (TMF), Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen), Matthias Löbe (Universität Leipzig)*
- Einführung in die tranSMART-Plattform
- 15'+3' *Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen)*
- tranSMART Architecture and Roadmap
- 45'+15' *Kees van Bochove (CIO The Hyve, tranSMART-RoadMap-Presentation at tranSMART-Meeting 2015)*
- Open Data Analytics – Gefahren durch p-Hacking
- 15'+3' *Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen)*

## 12.00 Uhr Mittagspause

## 13.00 Uhr Praktische Arbeit

- Klick-a-thon: Durchgehen eines vorbereiteten praktischen Szenarios durch die Workshopteilnehmer
- Lösung von Aufgaben in Eigenarbeit

# Open Data Analytics – Gefahren durch p-Hacking

## TMF-Workshop

Aufbereitung und Analyse von Daten in phänotypischen und genotypischen  
Forschungsdatenbanken mit tranSMART

Berlin, 05.08.2016

Prof. Dr. Ulrich Sax

Department of Medical Informatics, WG Infrastructure for Translational Research

GEFÖRDERT VOM



Bundesministerium  
für Bildung  
und Forschung

# Today's Random Medical News

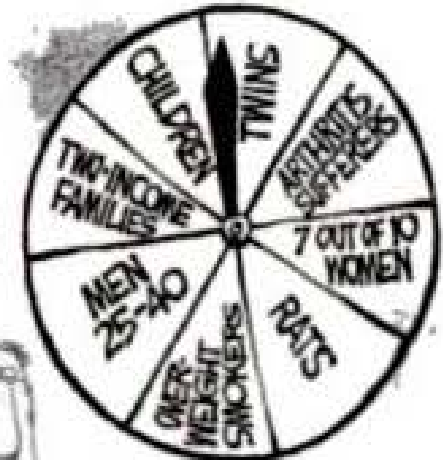
from the New England  
Journal of  
Panic-Inducing  
Gobbledygook



CAN CAUSE



IN



ACCORDING TO A  
REPORT RELEASED  
TODAY...

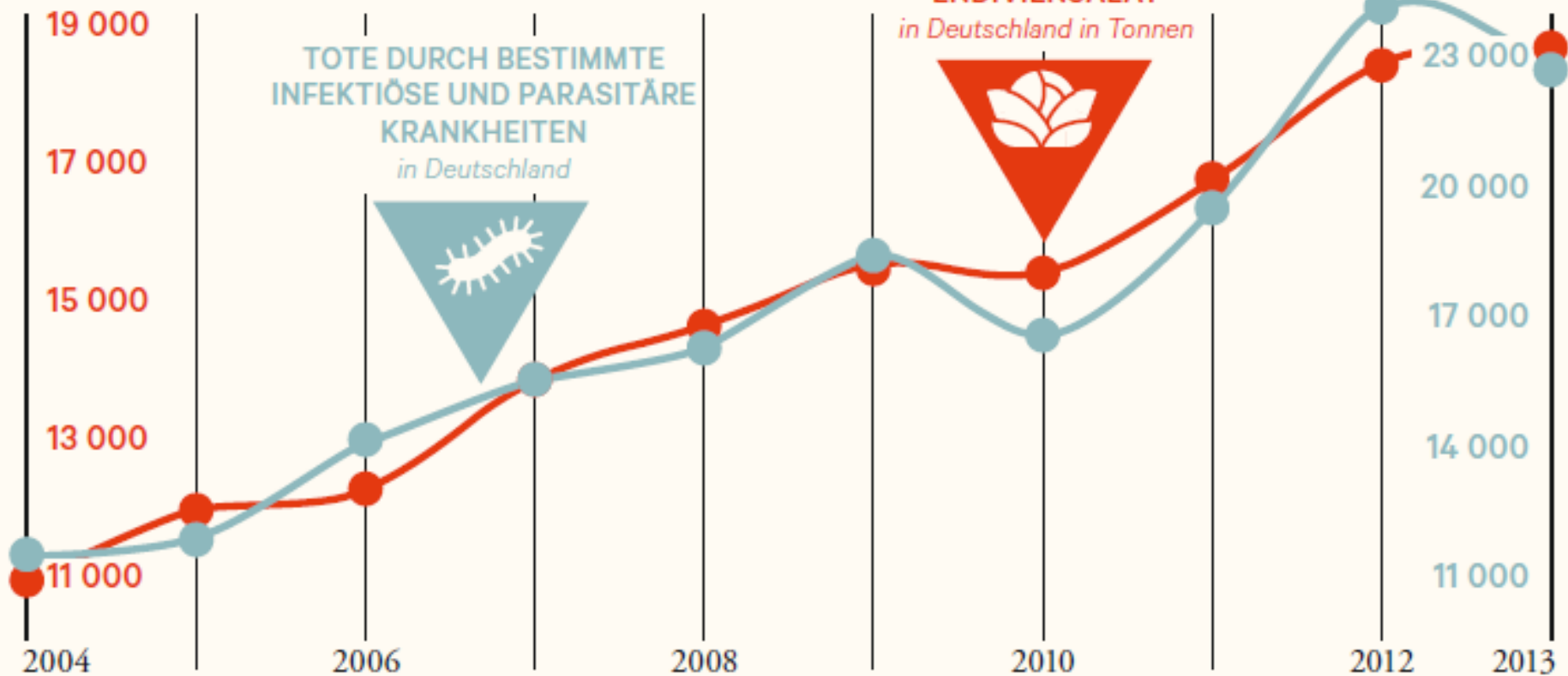
NEWS

**Jim Borgman**  
The Cincinnati Enquirer  
King Features Syndicate

# GIFTIGES GRÜN

Ist Endiviensalat verantwortlich für Infektionskrankheiten?

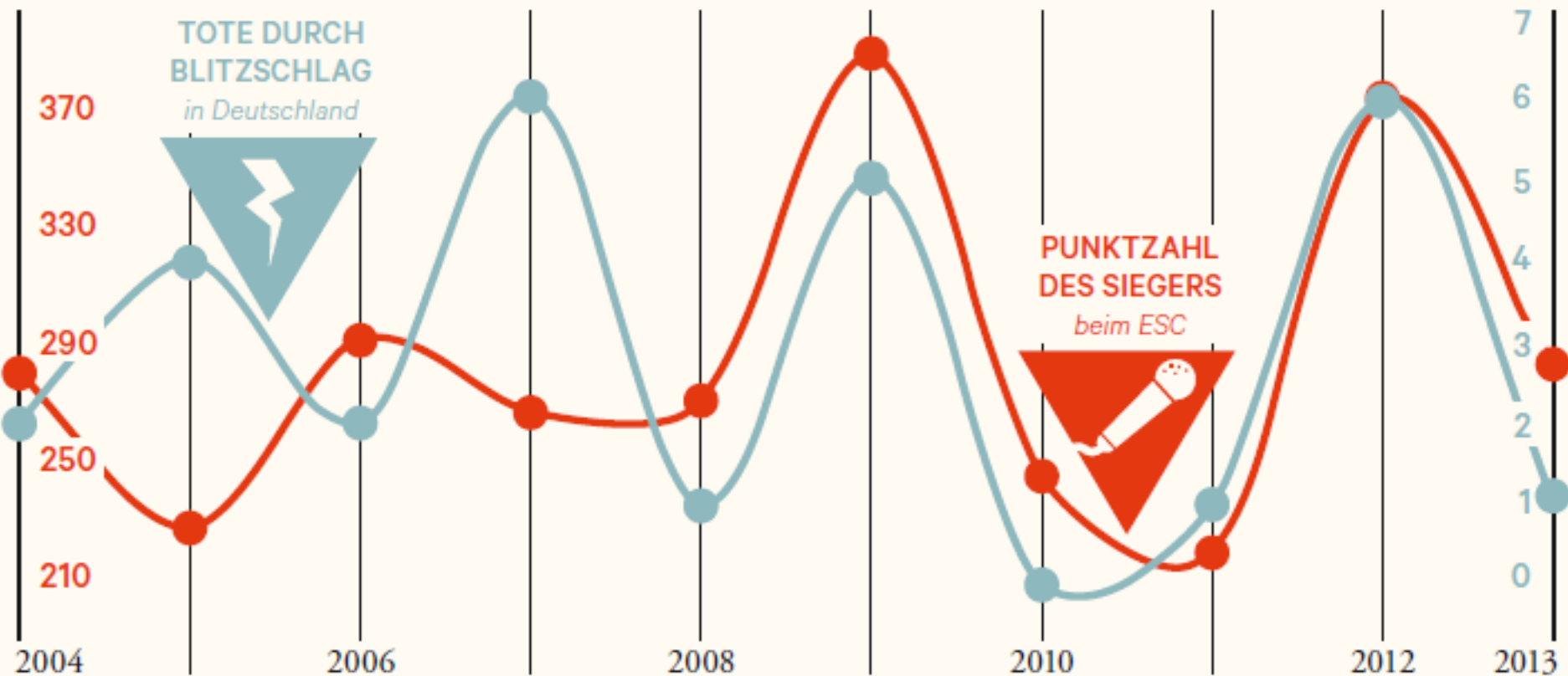
Korrelationskoeffizient: 0,981



# EINSCHLAGENDER ERFOLG

Was hat die Punktzahl des Siegers beim Eurovision Song Contest mit Toten durch Blitzschlag zu tun?

Korrelationskoeffizient: 0,571



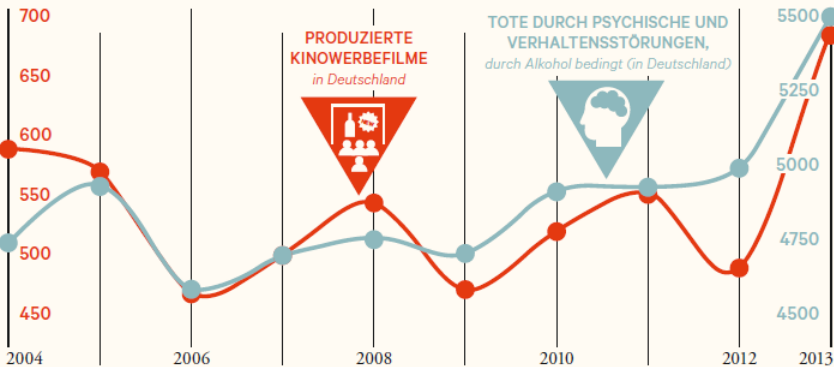


# Es ist nicht so, wie es scheint

Guck mal, die Grafik! Wie ähnlich die beiden Kurven verlaufen, da muss es doch einen Zusammenhang geben ...  
 Nein, muss es nicht. Aber solche Trugschlüsse sind ein häufiges Phänomen. Wir haben in den unterschiedlichsten Statistiken nach solchen Scheinkorrelationen gesucht. Und viel dabei gelacht

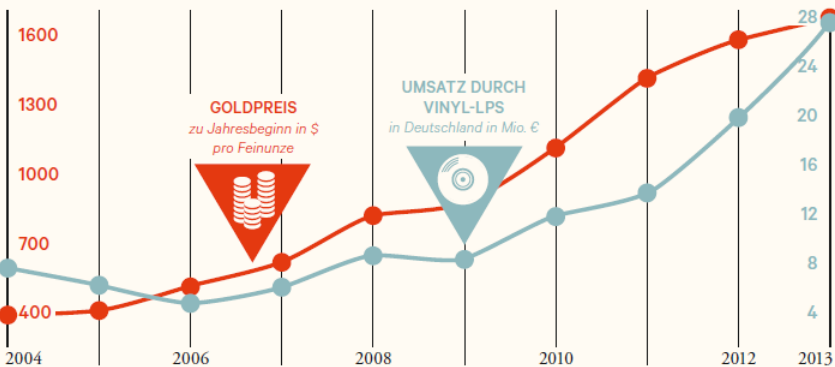
## IRREFÜHRENDE WERBUNG

Können Werbefilme verantwortlich sein für den Tod von Menschen?  
 Korrelationskoeffizient: 0,792 (siehe Erklärung unten rechts)



## GOLDENE SCHALLPLATTEN


Treibt der Verkauf von Vinylplatten den Goldpreis in die Höhe?  
 Korrelationskoeffizient: 0,905





# FOOLING OURSELVES

HUMANS ARE REMARKABLY GOOD AT SELF-DECEPTION



# FOOLING OURSELVES

**HUMANS ARE REMARKABLY GOOD AT SELF-DECEPTION.  
BUT GROWING CONCERN ABOUT REPRODUCIBILITY IS DRIVING MANY  
RESEARCHERS TO SEEK WAYS TO FIGHT THEIR OWN WORST INSTINCTS.**

**BY REGINA NUZZO**

In 2013, five years after he co-authored a paper showing that Democratic candidates in the United States could get more votes by moving slightly to the right on economic policy<sup>1</sup>, Andrew Gelman, a statistician at Columbia University in New York City, was chagrined to learn of an error in the data analysis. In trying to replicate the work, an undergraduate student named Yang Yang Hu had discovered that Gelman had got the sign wrong on one of the variables.

Gelman immediately published a three-sentence correction, declaring that everything in the paper's crucial section should be considered wrong until proved otherwise.

# Gefahren durch p-Hacking

- Was ist p-Hacking?
- Warum ist es gefährlich?

# p-Wert, statistische Signifikanz

- Statistisch signifikant, auch statistisch bedeutsam, wird das Ergebnis einer Untersuchung genannt, wenn die statistische Auswertung der Daten ergibt, dass die Wahrscheinlichkeit für die Annahme, die festgestellten Unterschiede zwischen Messgrößen oder Variablen seien durch Zufall derart zustande gekommen, einen zuvor als Signifikanzniveau festgelegten Schwellenwert nicht überschreitet.

Quelle: [https://de.wikipedia.org/wiki/Statistische\\_Signifikanz](https://de.wikipedia.org/wiki/Statistische_Signifikanz)

1. Stephen Stigler: *Fisher and the 5% level*. In: *CHANCE*. Band 21, Nr. 4, Springer, New York Dezember 2008, S. 12.

# p-Wert, statistische Signifikanz

- Die obere Grenze für die Irrtumswahrscheinlichkeit, also jener Wert, den man für die Wahrscheinlichkeit eines Fehlers 1. Art noch eben zu akzeptieren bereit ist, heißt Signifikanzniveau.
- **Grundsätzlich ist dies frei wählbar; häufig wird ein Signifikanzniveau von 5 % verwendet ( $p \leq 0,05$ ) [1]**
- In der Praxis bedeutet dieses Kriterium, dass im Schnitt eine von 20 Untersuchungen, bei denen die Nullhypothese richtig ist (z. B. ein Medikament tatsächlich wirkungslos ist), zu dem Schluss kommt, sie sei falsch (z. B. behauptet, das Medikament erhöhe die Heilungschancen).

Quelle: [https://de.wikipedia.org/wiki/Statistische\\_Signifikanz](https://de.wikipedia.org/wiki/Statistische_Signifikanz)

1. Stephen Stigler: *Fisher and the 5% level*. In: *CHANCE*. Band 21, Nr. 4, Springer, New York Dezember 2008, S. 12.

[https://de.wikipedia.org/wiki/Statistische\\_Signifikanz](https://de.wikipedia.org/wiki/Statistische_Signifikanz)

- Eine heuristische Motivation des Wertes 5 % ist wie folgt: Eine normalverteilte Zufallsgröße nimmt nur mit einer Wahrscheinlichkeit von weniger als ( $\leq$ ) 5 % einen Wert an, der sich vom Erwartungswert um mehr als die zweifache Standardabweichung unterscheidet:
- Bei einem p-Wert von kleiner oder gleich 5 % [ **$p \leq 0,05$** ] spricht man von einem signifikanten [...] Ergebnis.



Taylor & Francis  
Taylor & Francis Group



---

Spurious Correlation: A Causal Interpretation

Author(s): Herbert A. Simon

Source: *Journal of the American Statistical Association*, Vol. 49, No. 267 (Sep., 1954), pp. 467-479

Published by: Taylor & Francis, Ltd. on behalf of the American Statistical Association

Stable URL: <http://www.jstor.org/stable/2281124>

Accessed: 13-05-2015 09:19 UTC

---





Taylor & Francis  
Taylor & Francis Group



---

# SPURIOUS CORRELATION: A CAUSAL INTERPRETATION\*

HERBERT A. SIMON

*Carnegie Institute of Technology*

To test whether a correlation between two variables is genuine or spurious, additional variables and equations must be introduced, and sufficient assumptions must be made to identify the parameters of this wider system. If the two original variables are causally related in the wider system, the correlation is “genuine.”

**“SCIENCE IS AN ONGOING  
RACE BETWEEN OUR  
INVENTING WAYS TO FOOL  
OURSELVES, AND OUR  
INVENTING WAYS TO AVOID  
FOOLING OURSELVES.”**

Saul Perlmutter, an astrophysicist at the University of California, Berkeley

Failure to understand our own biases has helped to create a crisis of confidence about the reproducibility of published results, says statistician John Ioannidis, co-director of the Meta-Research Innovation Center at Stanford University in Palo Alto, California. The issue goes well beyond cases of fraud. Earlier this year, a large project that attempted to replicate 100 psychology studies managed to reproduce only slightly more than one-third<sup>2</sup>. In 2012, researchers at biotechnology firm Amgen in Thousand Oaks, California, reported that they could replicate only 6 out of 53 landmark studies in oncology and haematology<sup>3</sup>. And in 2009, Ioannidis and his colleagues described how they had been able to fully reproduce only 2 out of 18 microarray-based gene-expression studies<sup>4</sup>.

**Regina Nuzzo** is a freelance writer in Washington DC.

1. Gelman, A. & Cai, C. *J. Ann. Appl. Stat.* **2**, 536–549 (2008).
2. Open Science Collaboration. *Science* <http://dx.doi.org/10.1126/science.aac4716> (2015).
3. Begley, C. G. & Ellis, L. M. *Nature* **483**, 531–533 (2012).
4. Ioannidis, J. P. A. *et al. Nature Genet.* **41**, 149–155 (2009).

From: Finding the Missing Link for Big Biomedical Data

JAMA. 2014;311(24):2479-2480. doi:10.1001/jama.2014.4228

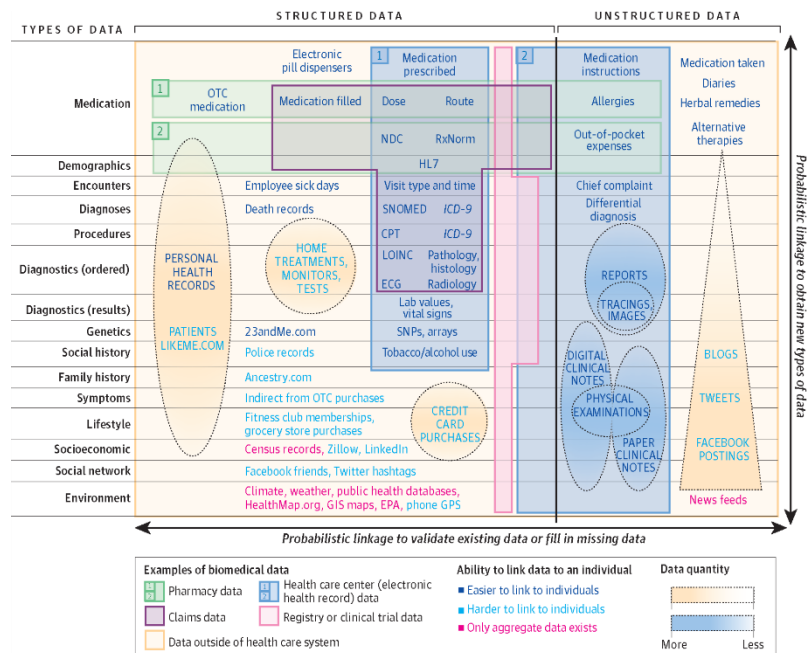


Figure Legend:

The Tapestry of Potentially High-Value Information Sources That May be Linked to an Individual for Use in Health Care CPT indicates current procedural terminology; ECG, electrocardiography; EPA, US Environmental Protection Agency; GIS, geographic information systems; GPS, global positioning system; HL7, Health Level 7 coding standard; ICD-9, Institutional Classification of Diseases, Ninth Revision; LOINC, Logical Observation Identifiers Names and Codes; NDC, National Drug Code; OTC, over-the-counter; SNOMED, Systematized Nomenclature of Medicine; SNP, single-nucleotide polymorphism.

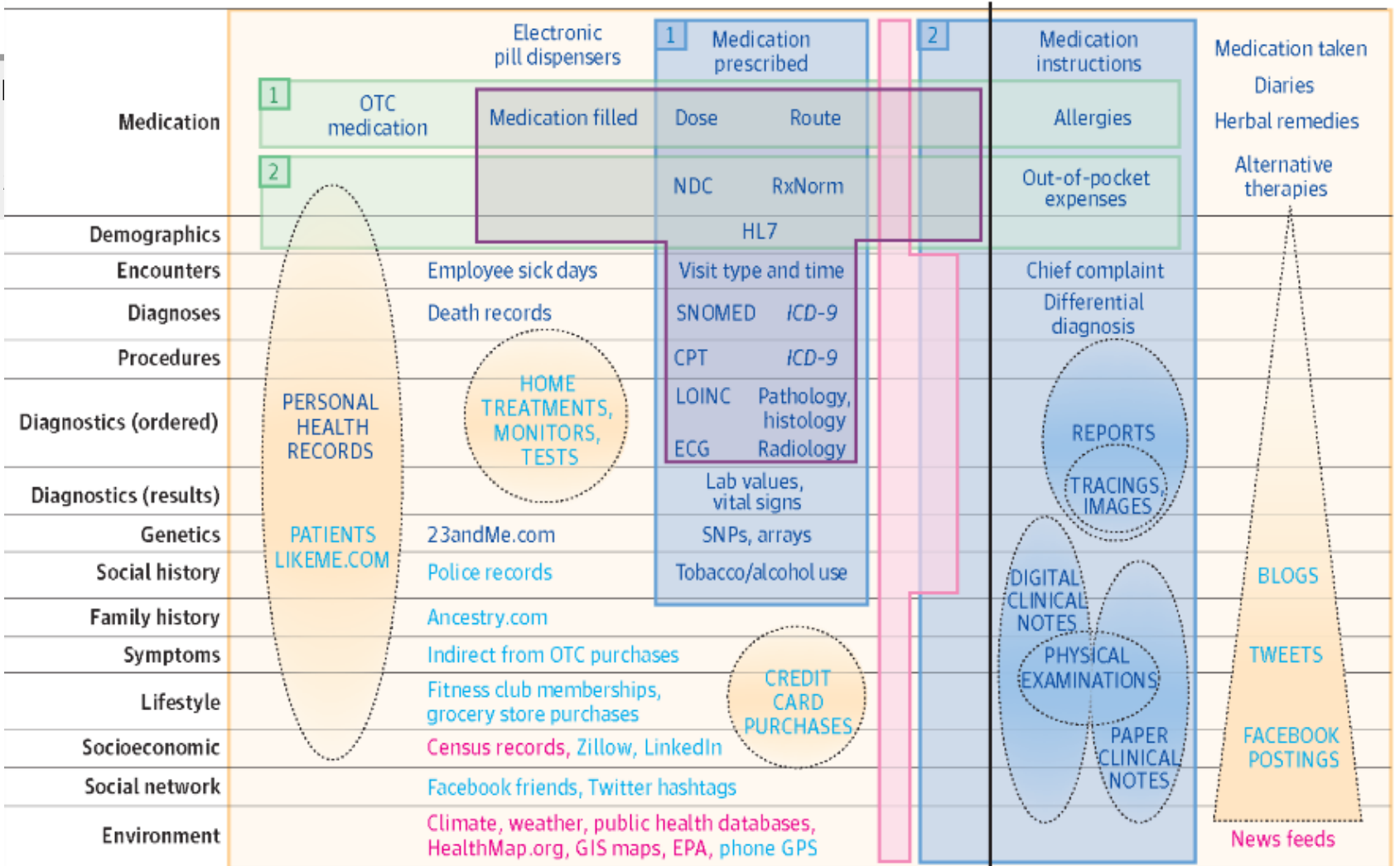


TYPES OF DATA

STRUCTURED DATA

UNSTRUCTURED DATA

From:  
JAMA.



Probabilistic linkage to obtain new types of data

Probabilistic linkage to validate existing data or fill in missing data

Figure 1  
The Tap indicate informa  
Disease  
counter  
Date of

**Examples of biomedical data**

- 1 2 Pharmacy data
- 1 2 Health care center (electronic health record) data
- Claims data
- Registry or clinical trial data
- Data outside of health care system

**Ability to link data to an individual**

- Easier to link to individuals
- Harder to link to individuals
- Only aggregate data exists

**Data quantity**

More Less

# Lösungsansätze

- **Correlation does not replace causation!**
- R. Nuzzo Nature 2015:
  - “One solution that is piquing interest revives an old tradition: explicitly considering competing hypotheses, and if possible working to develop experiments that can distinguish between them. “
  - “Furthermore, when scientists make themselves explicitly list alternative explanations for their observations, they can reduce their tendency to tell just-so stories.”
- Zur Vertiefung:
  - <http://www.laborjournal.de/editorials/981.lasso>
  - [http://www.labtimes.org/editorial/e\\_654.lasso](http://www.labtimes.org/editorial/e_654.lasso)

## Eine neue Wissenschaft-(lichkeit)?

Big Data, Innovation, Personalisierte Medizin und Co. – Sind dies die Markenzeichen einer neuen Wissenschaft-(lichkeit) in der Medizin? Ein Essay von Gerd Antes, Freiburg.



tweet



teilen

g+ +1



Wenn man Editorials, Kommentare oder Meinungsartikel in wissenschaftlichen Zeitschriften liest und ihnen glaubt, so stehen wir am Beginn eines goldenen Zeitalters für Patienten und Gesunde. Patienten werden viel früher und fehlerfrei diagnostiziert und dann mit personalisierter Medizin zielgenau, wirksam und nebenwirkungsfrei behandelt. Gesunde kommen erst gar nicht in die Gefahr, weil sie durch perfekte Vorsorge vor dem Schritt geschützt werden, überhaupt erst krank zu werden.

Erreicht wird dies mit „Systemmedizin“, in der die biologischen Mechanismen der Krankheitsentstehung unter Nutzung der Methoden von „Omics“-Forschung, Systembiologie, Informatik und Netzwerktheorie besser verstanden werden – was dann wiederum durch „Translation“ scheinbar wie von selbst in der personalisierten oder individualisierten Medizin genutzt wird. Überall dabei ist „Big Data“ als

Universalwerkzeug; Hindernisse und Barrieren gibt es praktisch nicht. Erforderlich sind nur unbegrenzte Rechnerleistung, Datenerfassung ohne Behinderung samt deren Speicherung in grenzenlosen Clouds, sowie

ZymoPURE™



Rotilabo®-  
Dispenser



ab 155,00

ILMA



Free ticket with  
PrioCode:

# ShowCase tranSMART



Please login...

Login ID:

Password:

Not a user? Contact [administrator](#) to request an account

## DISCLAIMER

By logging in to this application I acknowledge that, according to the [DFG-Memorandum on Safeguarding Good Scientific Practice](#), I am obliged to ensure good scientific practice by consulting a (bio)medical informaticist (data management) and a statistician or epidemiologist (analysis) prior to publishing any results obtained (partly or in full) through the use of this application.

Noncompliance could lead to a severe lack in reproducibility and may constitute a case of scientific misconduct.



## 10.15 Uhr Grundlegende

### Konzepte ■ Begrüßung

**5'** *Dr. Johannes Drepper (TMF), Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen), Matthias Löbe (Universität Leipzig)*

### ■ Einführung in die tranSMART-Plattform

**15'+3'** *Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen)*

### ■ tranSMART Architecture and Roadmap

**45'+15'** *Kees van Bochove (CIO The Hyve, tranSMART-RoadMap-Presentation at tranSMART-Meeting 2015)*

### ■ Open Data Analytics – Gefahren durch p-Hacking

**15'+3'** *Prof. Dr. Ulrich Sax (Universitätsmedizin Göttingen)*

## 12.00 Uhr Mittagspause

## 13.00 Uhr Praktische Arbeit

■ Klick-a-thon: Durchgehen eines vorbereiteten praktischen Szenarios durch die Workshopteilnehmer

■ Lösung von Aufgaben in Eigenarbeit

## 13.00 Uhr Praktische Arbeit

- Klick-a-thon: Durchgehen eines vorbereiteten praktischen Szenarios durch die Workshopteilnehmer
- Lösung von Aufgaben in Eigenarbeit

## 14.15 Uhr Kaffeepause

## 14.45 Uhr Fortgeschrittene Themen

- tranSMART and beyond

**20'+5'** *Reinhard Schneider / Sascha Herzinger*  
(Universität Luxemburg)

- SmartR – Dynamic Visual Analytics in TranSMART

**20'+5'** *Reinhard Schneider / Sascha Herzinger*  
(Universität Luxemburg)

- Import und Repräsentation hochdimensionaler Daten

**10'+3'** *Christoph Knell (Friedrich-Alexander-Universität Erlangen-Nürnberg)*

# Diskussion

## Zusammenfassung und Ausblick

### TMF-Workshop

Aufbereitung und Analyse von Daten in phänotypischen und genotypischen  
Forschungsdatenbanken mit tranSMART

Berlin, 05.08.2016

Prof. Dr. Ulrich Sax

Department of Medical Informatics, WG Infrastructure for Translational Research

GEFÖRDERT VOM



Bundesministerium  
für Bildung  
und Forschung



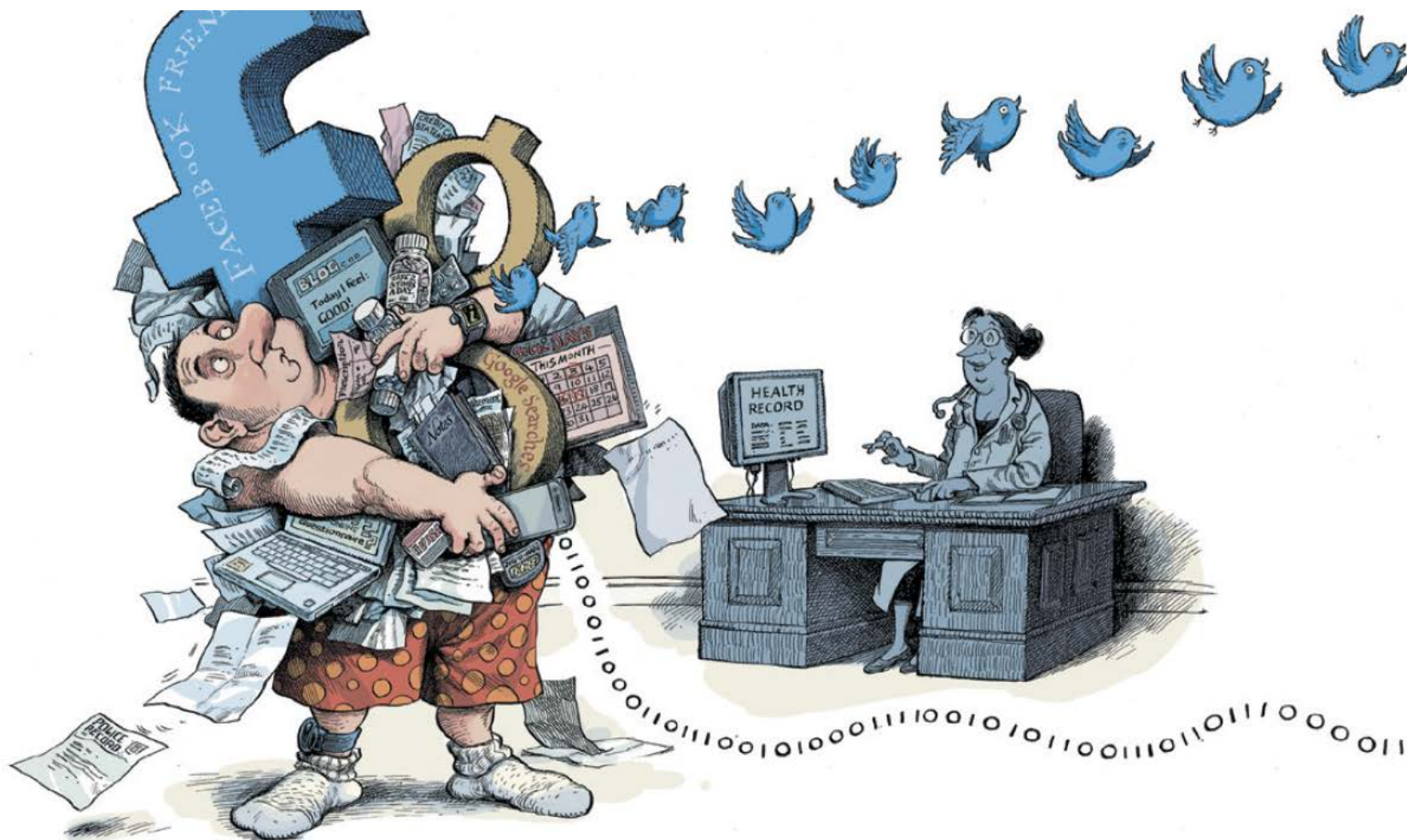
# Wrap up and Future directions

**Was nehmen Sie mit nach Hause?**

**Lösungsansätze für die Praxis**

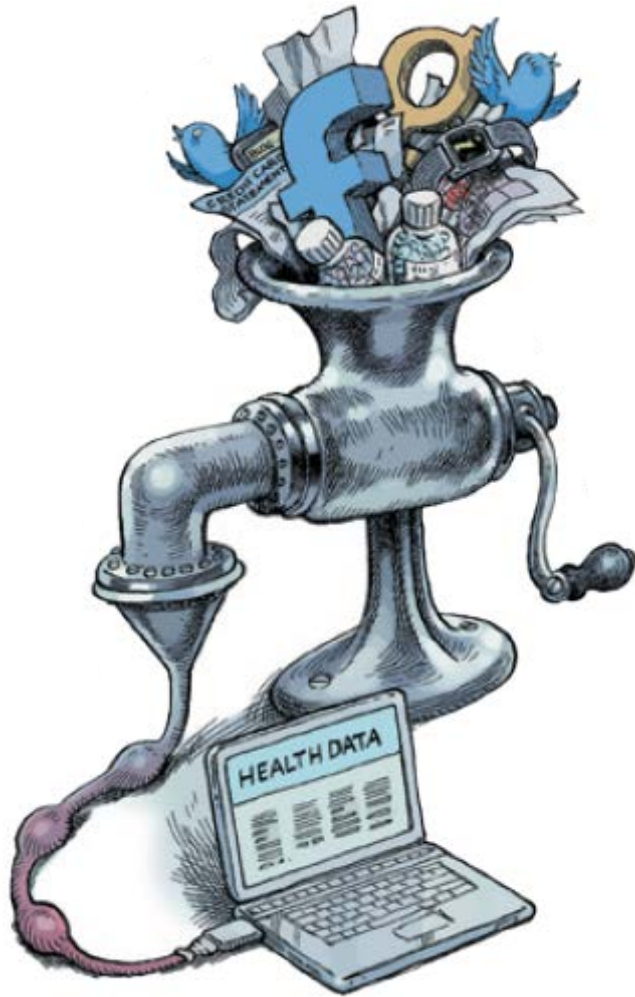
**neue Fragen**

**Ideen und Ansprechpartner**



## Make sense of health data

Develop the science of data synthesis to join up the myriad varieties of health information, insist **Julian H. Elliott, Jeremy Grimshaw** and colleagues.



- Pooling data
- Managing Bias
- Joining the Dots

## Make sense of health data

Develop the science of data synthesis to join up the myriad varieties of health information, insist **Julian H. Elliott, Jeremy Grimshaw** and colleagues.

# Data Provenance: a forgotten concept?

## A Survey of Data Provenance in e-Science

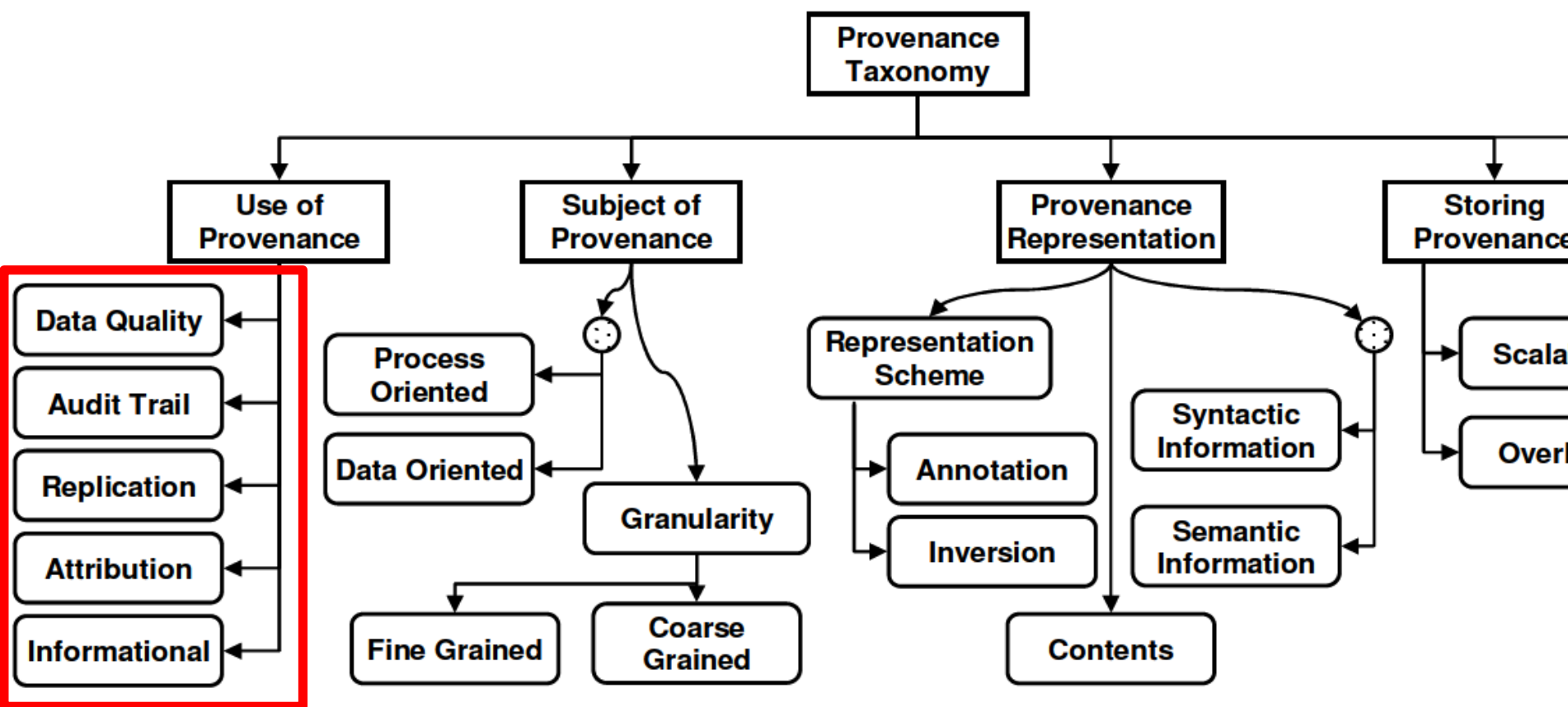
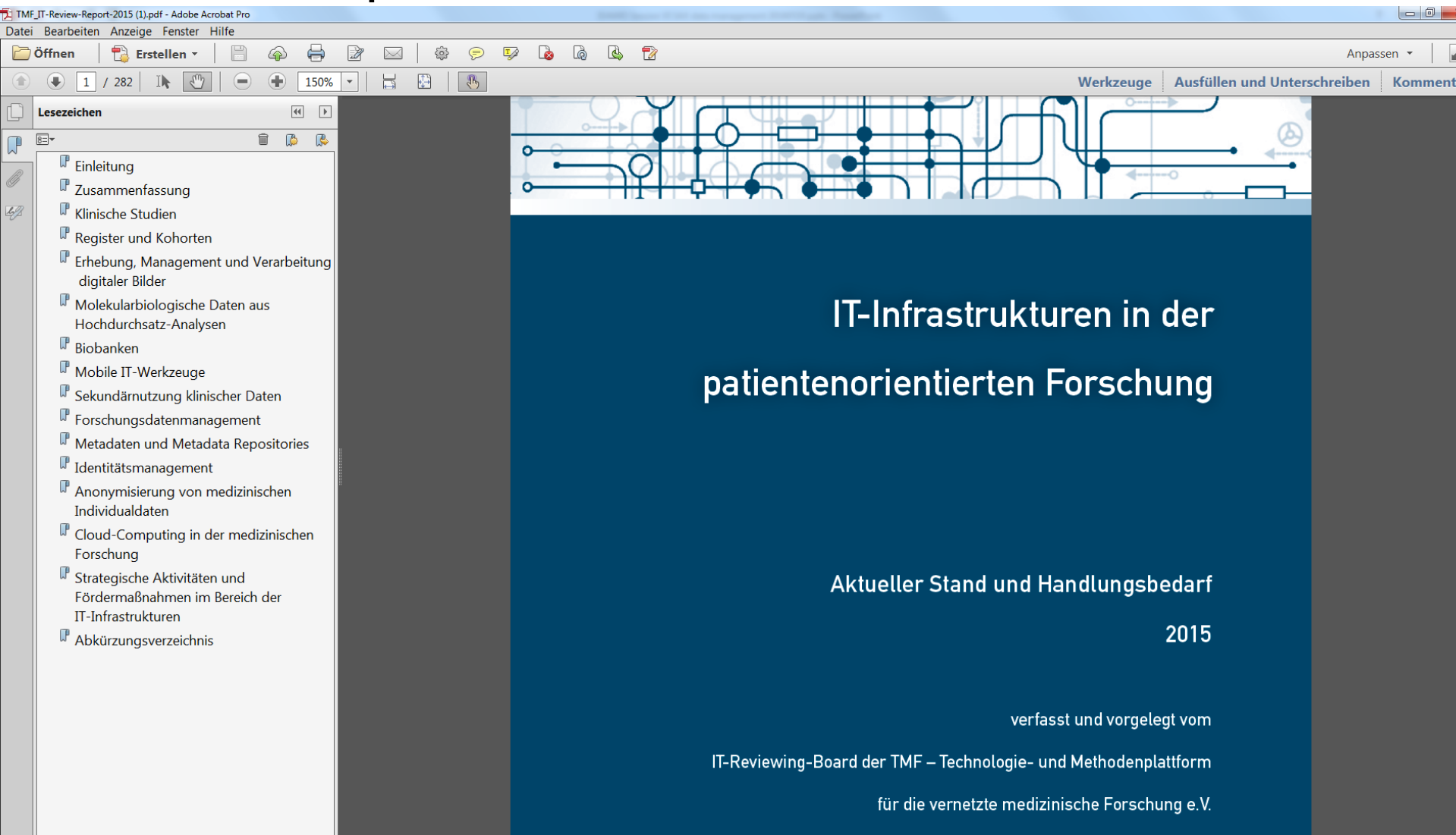


Figure 1 Taxonomy of Provenance

# TMF IT-Report 2015/6



The image shows a screenshot of the Adobe Acrobat Pro interface. The main window displays the cover page of the 'TMF IT-Report 2015/6'. The cover page has a dark blue background with a white circuit-like pattern at the top. The title 'IT-Infrastrukturen in der patientenorientierten Forschung' is centered in white. Below the title, it says 'Aktueller Stand und Handlungsbedarf' and '2015'. At the bottom, it states 'verfasst und vorgelegt vom IT-Reviewing-Board der TMF – Technologie- und Methodenplattform für die vernetzte medizinische Forschung e.V.'.

The left sidebar shows the 'Lesezeichen' (Bookmarks) panel with a list of report sections:

- Einleitung
- Zusammenfassung
- Klinische Studien
- Register und Kohorten
- Erhebung, Management und Verarbeitung digitaler Bilder
- Molekularbiologische Daten aus Hochdurchsatz-Analysen
- Biobanken
- Mobile IT-Werkzeuge
- Sekundärnutzung klinischer Daten
- Forschungsdatenmanagement
- Metadaten und Metadata Repositories
- Identitätsmanagement
- Anonymisierung von medizinischen Individualdaten
- Cloud-Computing in der medizinischen Forschung
- Strategische Aktivitäten und Fördermaßnahmen im Bereich der IT-Infrastrukturen
- Abkürzungsverzeichnis

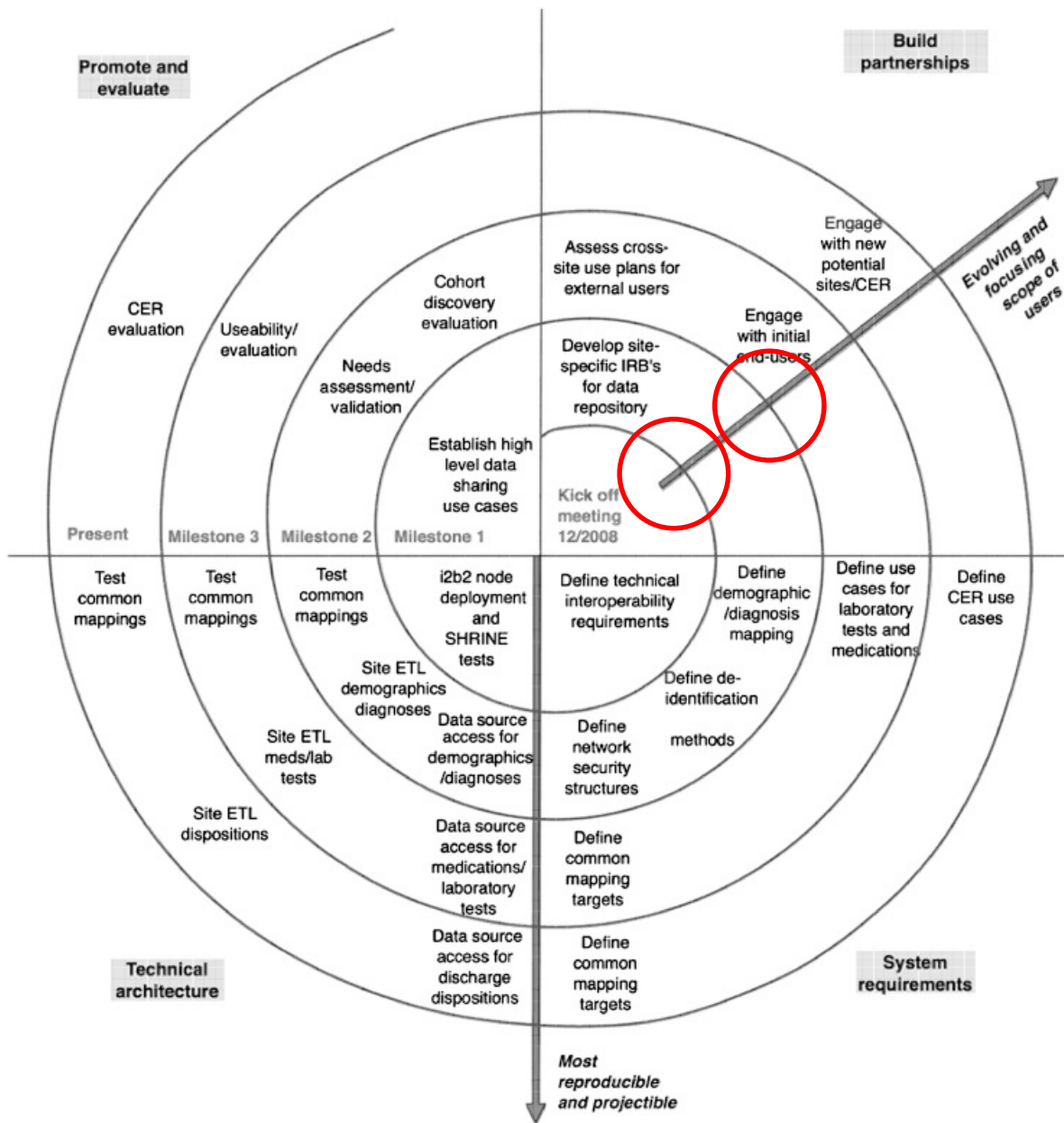


# Wrap up and Future directions

## Was nehmen Sie mit nach Hause?

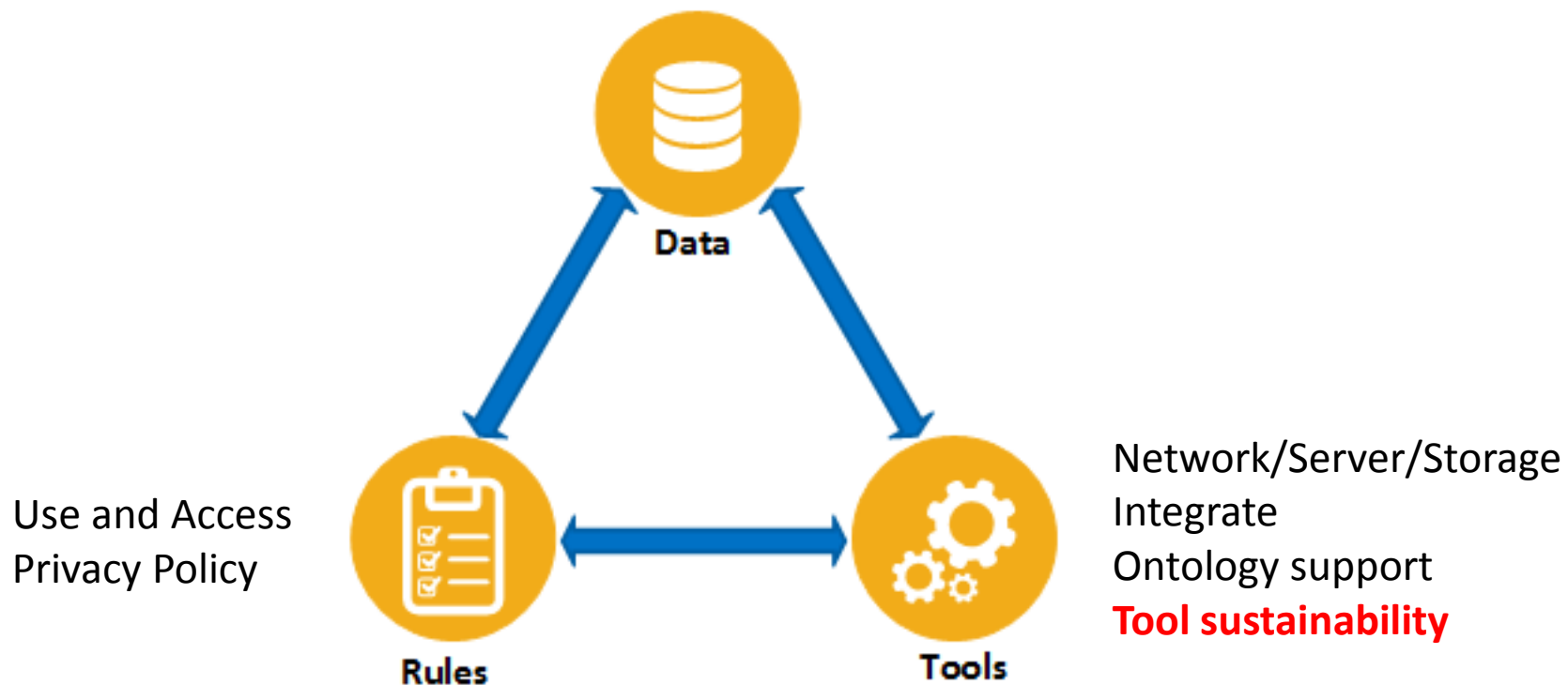
### Lösungsansätze für die Praxis





# Discussion

Sifting and Cleaning  
Phenotyping!



# Aufbereitung und Analyse von Daten in phänotypischen und genotypischen Forschungs- datenbanken mit tranSMART



5. August 2016 | Berlin

