

# Workshop

# Leibniz Bioactives Cloud

2020-01-29

Leibniz Institut für Pflanzenbiochemie  
Halle (Saale)

# Inhalt

- Einführung
  - Zielstellung
  - Architektur
  - Status Quo
- Wo wollen wir hin
  - Multi-Cloud
  - CloudLIMS
  - Semantische Textanalyse?

# Leibniz Forschungsverbund

<https://www.leibniz-wirkstoffe.de>



- Diversität der wissenschaftlichen Themen, Organisationsformen und verteilt über Deutschland
- Exzellenz durch Förderung von Kooperation
- Maßnahmen
  - Konferenzen
  - “Seed money”
  - **Leibniz Bioactives Cloud** als Werkzeug für Wissenschaftler

# Zweck der Leibniz Bioactives Cloud

- Austausch von nichtöffentlichen Dokumenten (pdf)
- Vernetzung der Wissenschaftler
- (Semantische) Analyse der verfügbaren Dokumente
- Austausch von wissenschaftlichen Daten (chemische Strukturen, Assay-Ergebnisse, ...)

# Vorhandene Technologie

Existierende Softwareprojekte:

- Web basiert



- Fortschrittliche Technologie



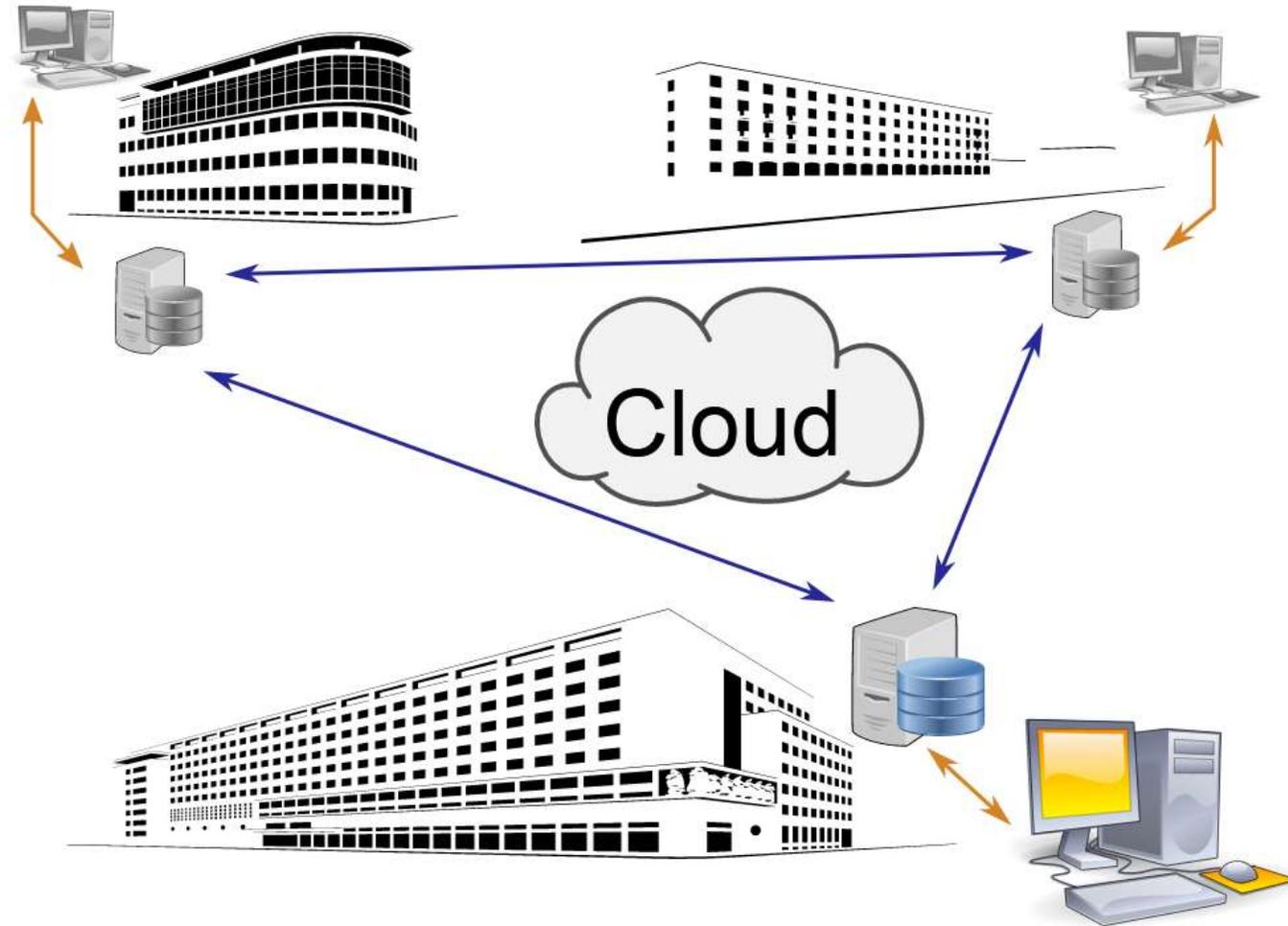
- Verteilte Infrastruktur



- "Chemietauglich"

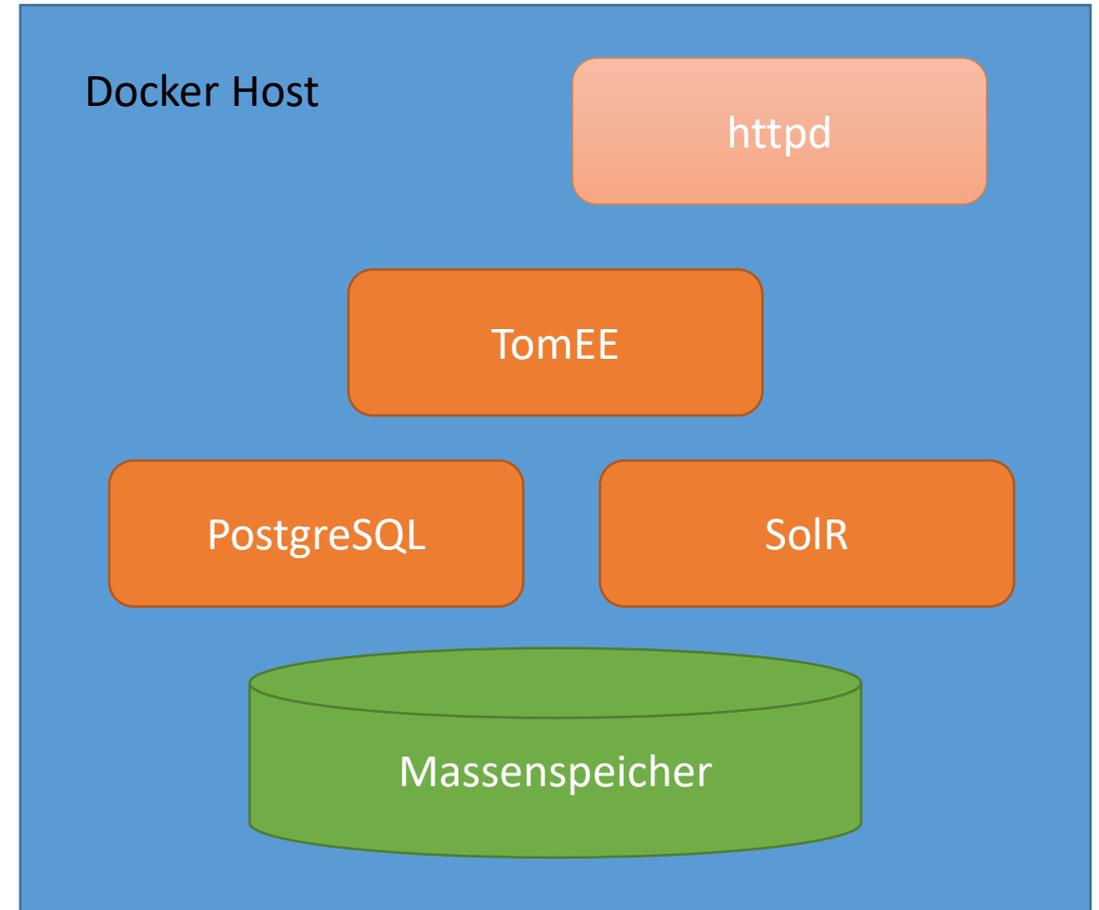


# Verteilte Architektur



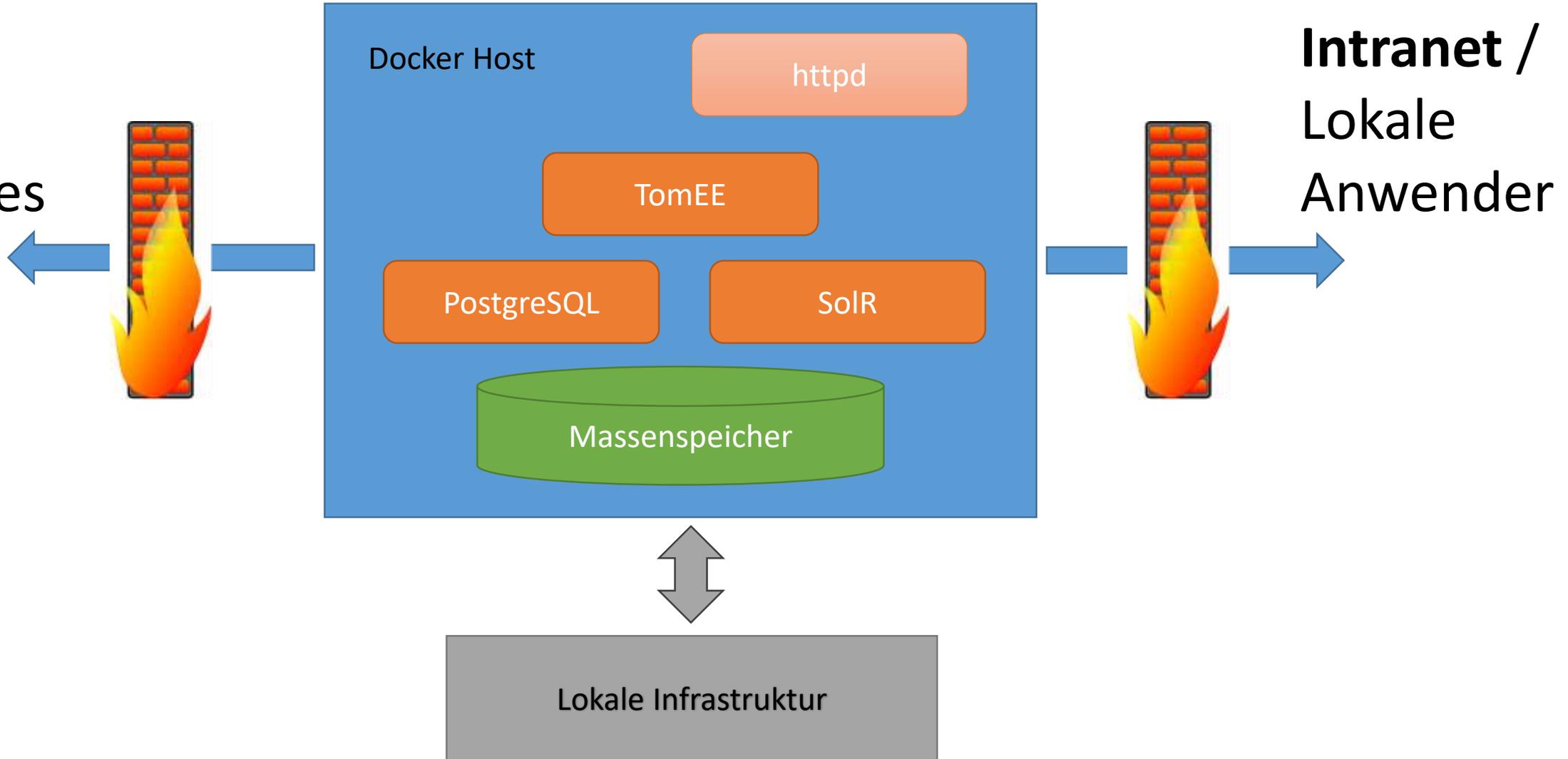
# Architektur / Designprinzipien

- Linux / Docker / Java
- Aufbau aus Standardkomponenten
  - Apache TomEE 7
  - PostgreSQL
  - SolR
  - [Apache httpd]
- Geringe Ressourcenanforderungen (Technik **und** Personal)
- Absicherung durch Verschlüsselung
- Vermeidung unnötiger Netzwerkverbindungen



# Architektur

**Welt /  
Leibniz  
Bioactives  
Cloud**



**Intranet /  
Lokale  
Anwender**

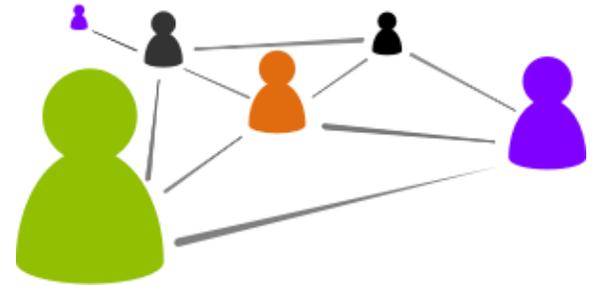
Lokale Infrastruktur

# Entwicklung und Verwaltung

- Build-Umgebung (Compiler, IDE, Build-Werkzeuge, ...)
- Sourcecode Repository, Wiki, Projektmanagement (BitBucket, Confluence Jira)
- Webserver für Newsletter, Softwareverteilung, Dokumentation („Handbücher“): <https://www.leibniz-wirkstoffe.de/>
- Mehrere Testsysteme
- Hilfsprogramme für die Ausstellung von Zertifikaten und Zertifikatssperrrlisten, Softwareverteilung, ...

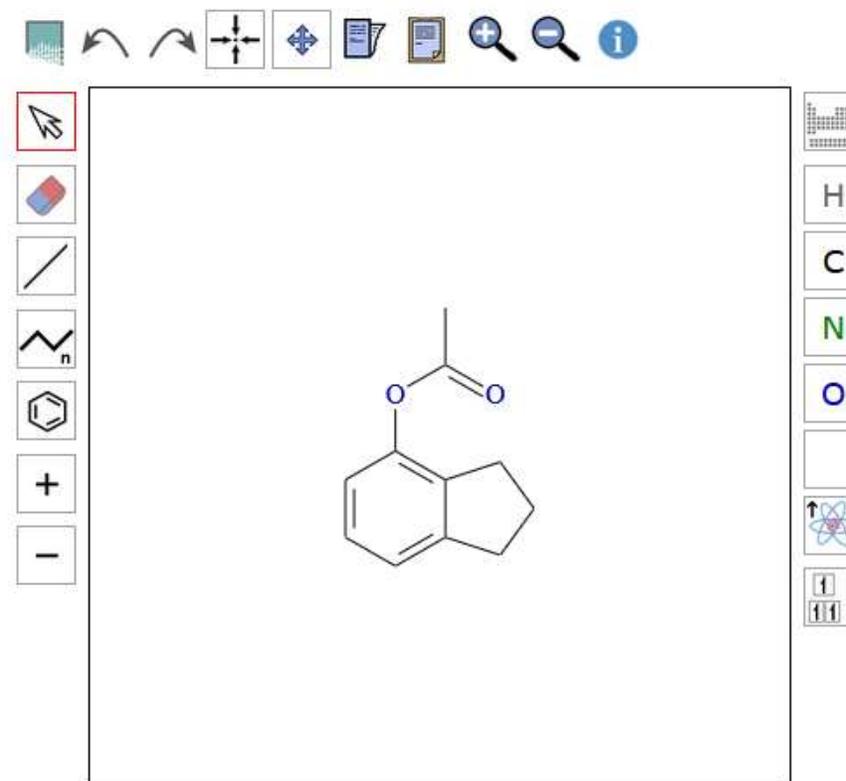
# Leibniz Bioactives Cloud: aktueller Stand

- Volltextsuche nach verteilt gespeicherten Dokumenten
- Mehrsprachige Dokumentenprozessierung mit globaler Relevanzmetrik
- Einfaches (verteiltes) Nutzerforum
- Word-Cloud-Suche (hauptsächlich als Technologiestudie)
- Feingranulare Rechteverwaltung mit optionaler LDAP-Anbindung
- Internationalisiertes Nutzerinterface in „responsivem“ Design



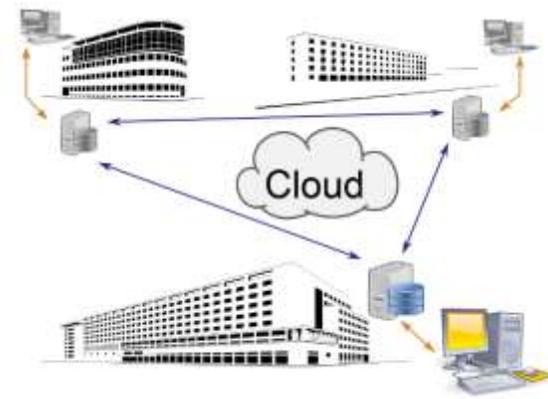
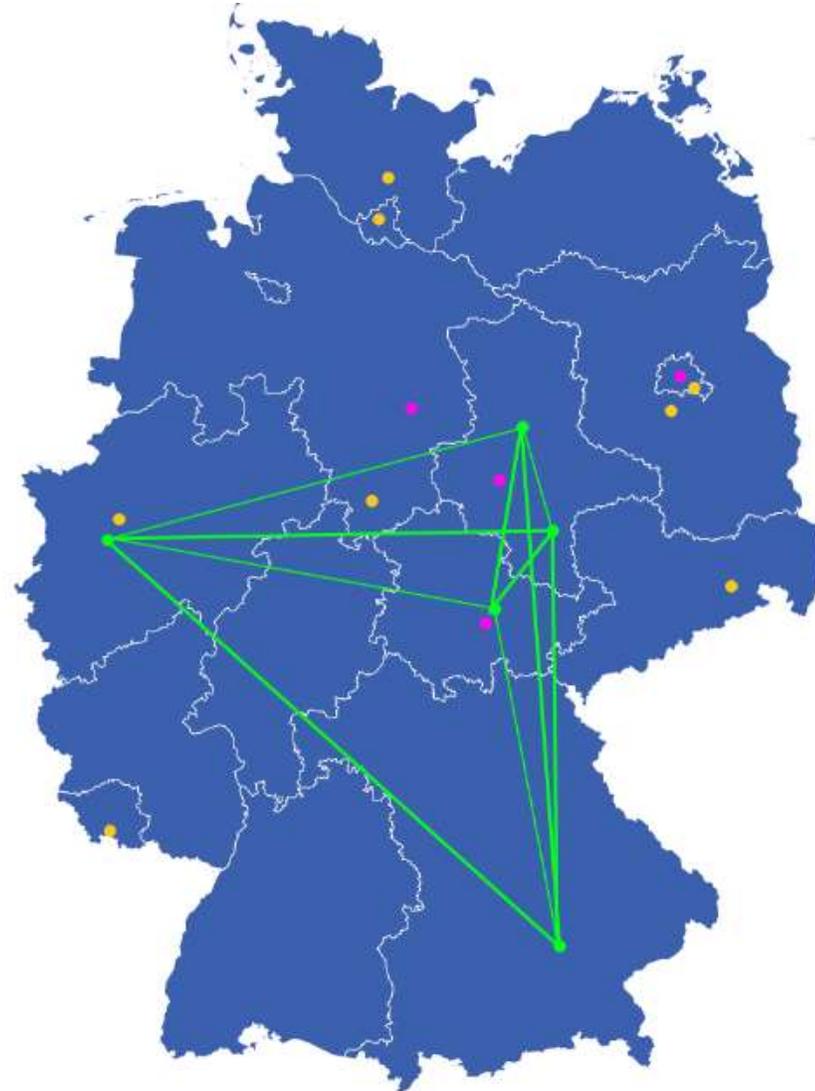
# Leibniz Bioactives Cloud: aktueller Stand

- Rund 70 Unit-Tests zur Qualitätssicherung  
Testabdeckung ca. 40 %
- Weitgehend automatisierte Installation
- Spin-Off: MolPaintJS  
<https://github.com/ipb-halle/MolPaintJS>



# Leibniz Bioactives Cloud: Rollout

- IPB Halle ✓
- HKI Jena ✓
- ISAS Dortmund ✓
- LIN Magdeburg ✓
- Leibniz-LSB@TUM ✓
- IPK Gatersleben 🕒
- FLI Jena 🕒
- FMP Berlin 🕒
- DSMZ 🔴

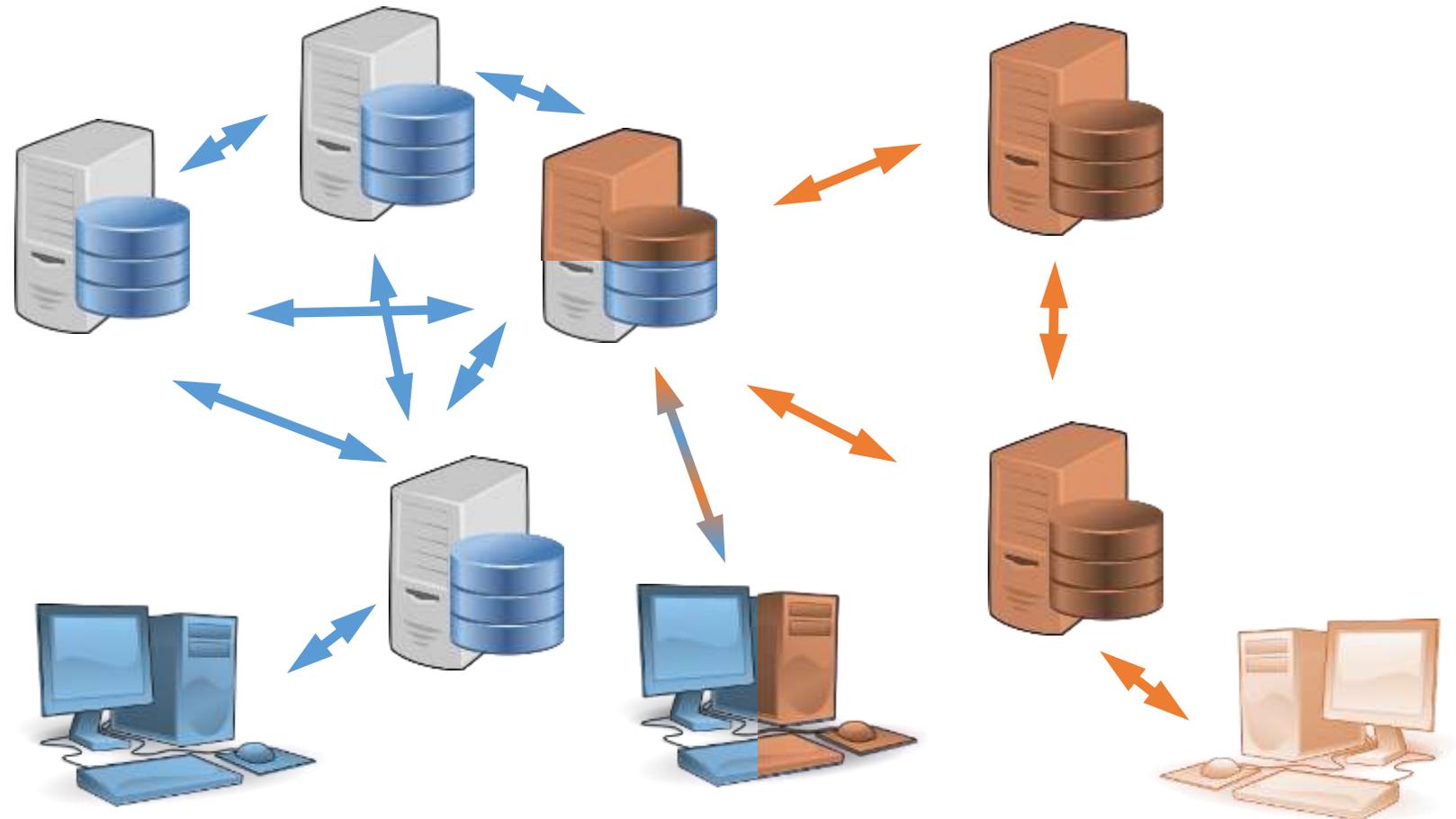


# Zukünftige Entwicklung: Sicherheit

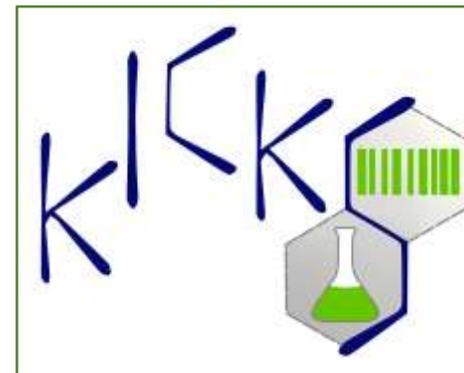
- Erhöhung der Sicherheit
  - HTTP Strict Transport Security (HSTS)
  - Content Security Policy Headers
  - Intruder Lockout
  - Lokale Spiegelung von JavaScript-Quellen (Bootsfaces, JQuery, ...)
- Datenschutz-Compliance
  - Datenschutzerklärung
  - Datensparsamkeit
  - Anonymisierung

# Zukünftige Entwicklung: Multi-Cloud

- Parallelbetrieb mehrerer Clouds
- Verbesserung der Wiederverwendbarkeit, Vermeidung mehrerer Instanzen pro Institut
- Strikte Trennung der einzelnen Clouds
- Zertifikatshierarchie / disjunkte CAs



# Exkurs: KICKS



- Seit 2007: Verwaltung von kommerziellen Chemikalien (Gefahrstoffdatenbank)
- Java & PostgreSQL-basiert
- ca. 120.000 Verbindungen, > 11.000 Gebinde
- Chemikalienbörse, Barcodes, Versionierung / Historie, PDF-Reports
- Chemische Struktursuche und Substruktursuche
- Seit 2016: auch an der Universidade Federal do Rio Grande do Sul
- Codebasis enthält Vorarbeiten für ein LIMS / ELN



Substanzdaten - Microsoft Internet Explorer bereitgestellt von IPB

http://kicks.ipb-halle.de/mf/user/subviews/substlist.jsf

Hilfe  

## Substanzliste

CAS/RN:

Summenformel:

Substanzbezeichnung:

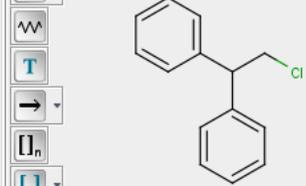
Substanz-ID:

SMILES:

nach verfügbaren Gebinden filtern:

Substruktursuche: [Editor verbergen](#)

Editor



Optionen

Stereochemie berücksichtigen

Nur übereinstimmende Strukturen

Isotopenmuster berücksichtigen

Suchen   Formular zurücksetzen   Zurück   Neue Substanz

Substanz-ID	Summenformel	CAS/RN	Substanz	Gefahreninformation	Aktion
137	C <sub>14</sub> H <sub>9</sub> Cl <sub>5</sub> O	115-32-2	Dicofol Dicofol (ISO) 2,2,2-Trichlor-1,1-bis(4-chlorphenyl)ethan	 R: 21/22-38-43-50/53 S: (2-)36/37-60-61	<a href="#">Details ...</a> <a href="#">Neues Gebinde</a> <a href="#">Gebinde anzeigen</a>
4080	C <sub>14</sub> H <sub>9</sub> Cl <sub>5</sub>	50-29-3	1,1-Bis(4-chlorphenyl)-2,2,2-trichlorethan DDT clofenotane	 R: 25-40-48/25-50/53 S: (1/2-)22-36/37-45-60-61	<a href="#">Details ...</a> <a href="#">Neues Gebinde</a> <a href="#">Gebinde anzeigen</a>

http://kicks.ipb-halle.de/mf/user/subviews/substlist.jsf

# KICKS: ELN / LIMS Screenshots

Taxonomic information - Mozilla Firefox

http://www.ipb-halle.de/mf/.../taxonomicinformation.html

## Taxonomic information

- Taxonomie
  - Animalia (3)
  - Archaeobacteria (1)
  - Chromista (1)
  - Eubacteria (2)
  - Fungi (10)
    - Agaricomycetes (2)
    - Ascomycota
      - Basidiomycetes (20)
        - Boletaceae (13)
          - Boletus chromapes
          - Boletus erythropus
          - Boletus gniseus
          - Boletus longicurvipes
          - Chalciporus piteratus
          - Leccinum chromapes**
          - Porphyrellus pseudoscaber
          - Pulveroboletus auriflammeus
          - Pulveroboletus ravenellii
          - Suillus bovinus
          - Suillus tridentinus
          - Tylopilus felleus
          - Xanthocomium affine
        - Coriolaceae

Organism / Taxon:

### Taxonomic Supercategory (~ies)

Boletaceae

### Category of Taxon

species

### Taxon names

Taxon names
Leccinum chromapes

Container: TS018

Desired action:

### Item Allocation

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y
0																									
1																									
2																									
3																									
4																									
5																									
6																									
7																									
8																									
9																									

ZAC000.637  
 TRAY TR018.A3  
 55-400 mg  
 Aliquots:

PURE: MTP MTP0128.DIOMSO: MTP MTP0127.DI

Experiments - Mozilla Firefox

http://www.ipb-halle.de/mf/.../ExpTable.html

Item Label: ZAC000.637 Purity: below 80%

Amount: 89.5 mg Last solvent:

Place: TRAY TS018.A2 Physical phase: brown crystals

Tara [mg]: 2452 [mg] ...

**ZAC 000.637: 2014-05-12**

**Journal reference:**

**Sample id:**

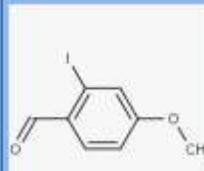
**Organism**

Taxon names	Language
Leccinum chromapes	-

**Substance-id: 115042**

Substance name	Language
2-iodo-4-methoxybenzaldehyde	-

**Components**

Structure	Formula, Concentration, Molar mass (average, exact)
	Empirical formula: C <sub>8</sub> H <sub>7</sub> IO <sub>2</sub> molar mass: 262.044 exact molecular mass: 261.949073 Concentration: 100%

Strucid: 15746

**Substance indexes**

Index description	Index value
Legacy MolId	0018

**Item Data**

Item Label: ZAC000.637 Purity: below 80%

Amount: 55.4 mg Last solvent:

Place: TRAY TS018.A3 Physical phase: white crystals

Tara [mg]: 2429 [mg] ...

[Details](#)

## Search Criteria

Search mode: Expert

Connector	Search Criteria	Action
NO. subara		<input type="button" value="Remove"/>
Taxon	Leccinum%	<input type="button" value="Remove"/>

Empirical formula: C<sub>8</sub>H<sub>7</sub>IO<sub>2</sub>

molar mass: 262.044

exact molecular mass: 261.949073

Concentration: 100%

Query type:  Match stereo  Match whole structure  Match isotope pattern

Add search criteria, search by:

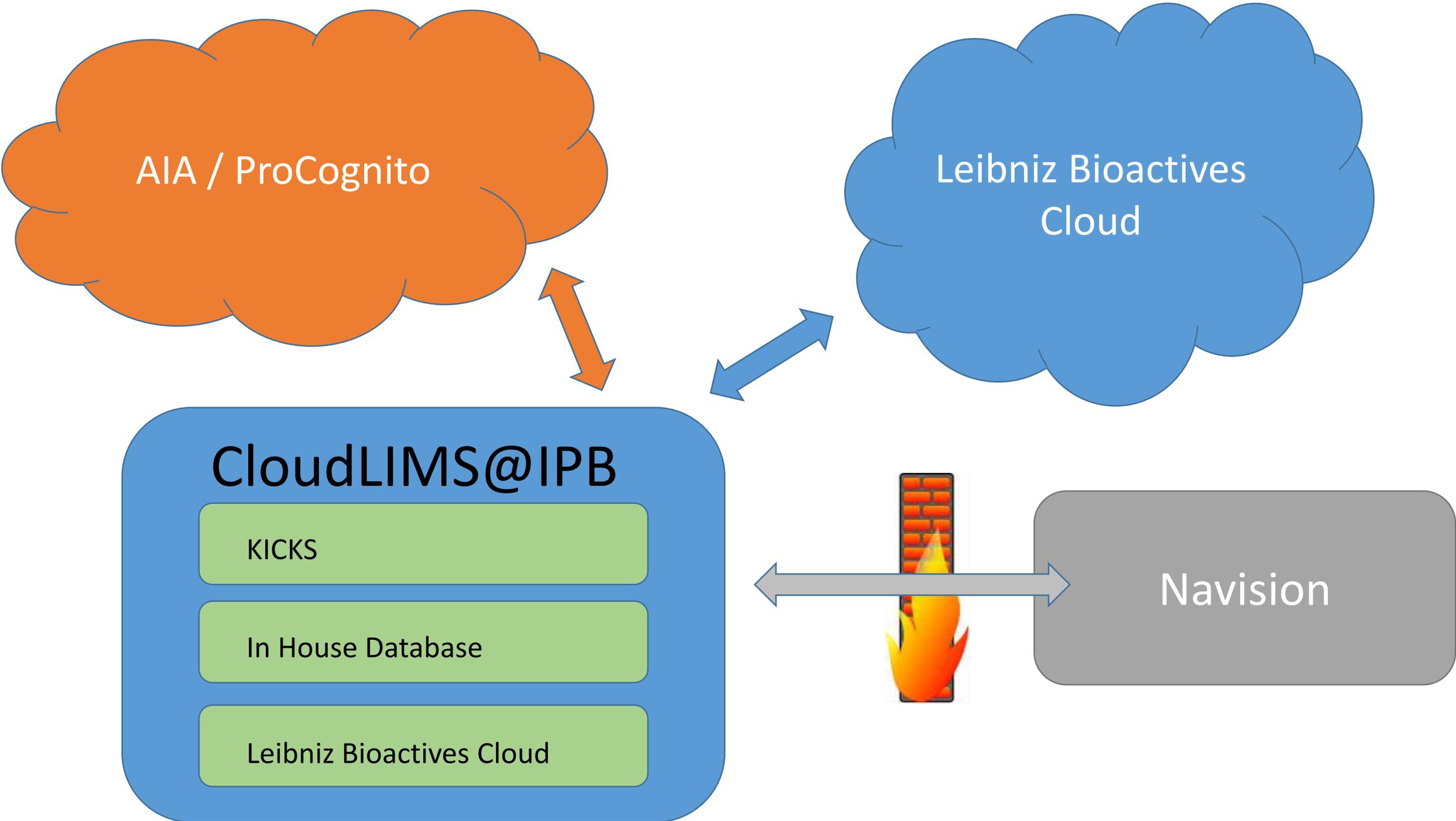
# Exkurs: InhouseDB



- „LIMS“ der Abteilung Natur- und Wirkstoffchemie
- > 20.000 Verbindungen, > 26.800 Substanzproben
- > 3.600 Extrakte aus > 1.200 Organismen
- Verknüpfung zu Laborjournalen
- Registratur für Assay-Ergebnisse
- ChemFinder basiert (MS Access ☹)
- Basis für virtuelles Screening

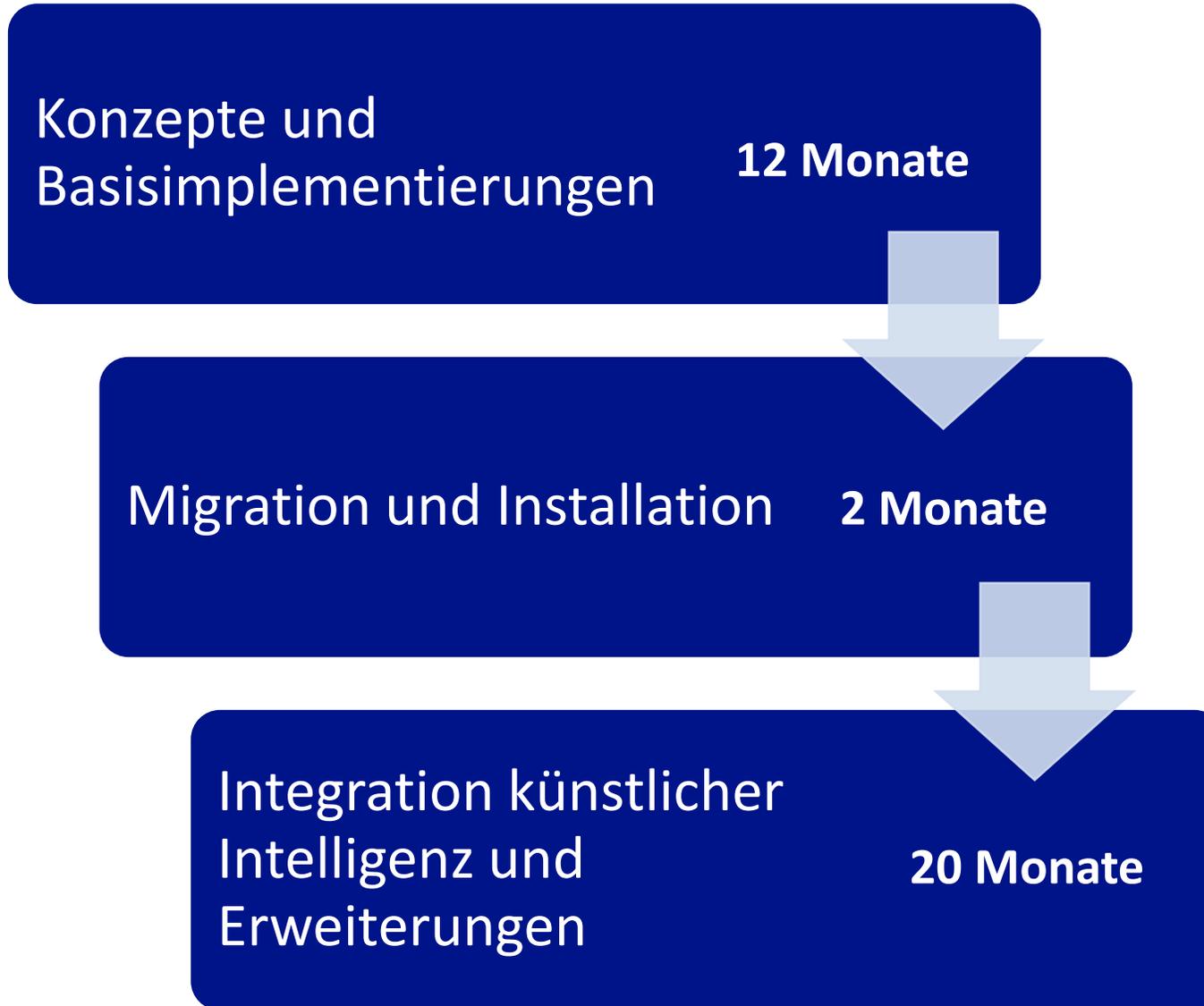


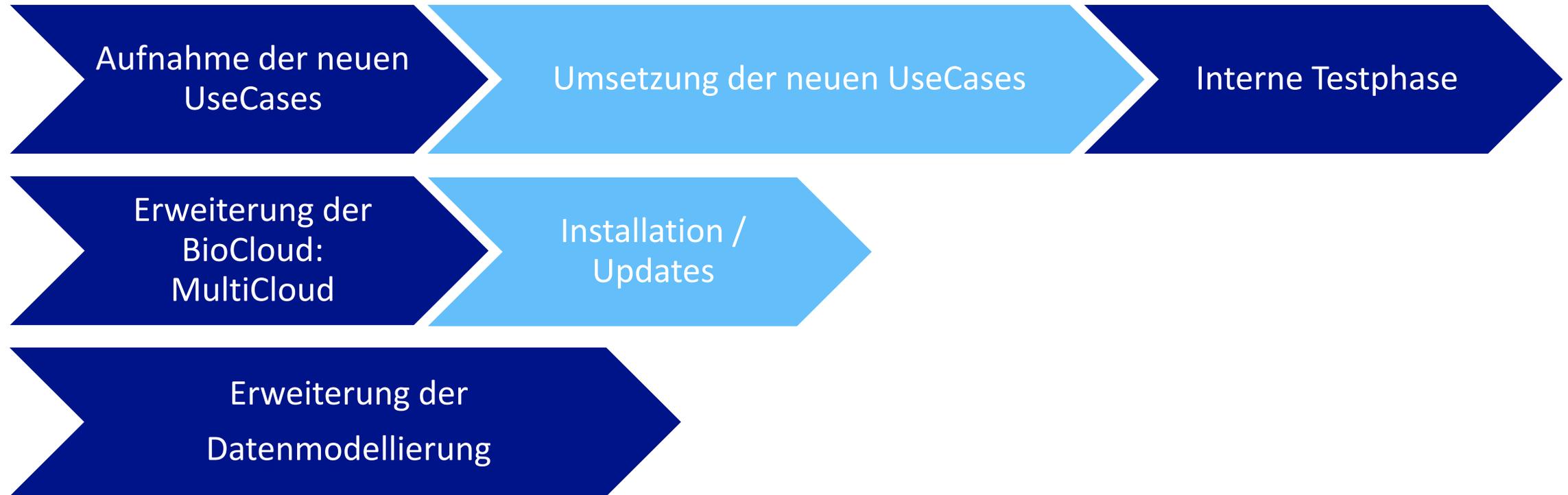




## Unterstützung der wissenschaftlichen Arbeit durch eine innovative Software

- Entscheidungsunterstützung bei der Auswahl geeigneter Biotestsysteme
- Gesicherter Fernzugriff auf ausgewählte Bereiche der Substanzdatenbank durch Mitglieder
- Verbesserung der Textanalyse durch den Einsatz künstlicher Intelligenz





- Die neuen UseCases umfassen hauptsächlich die der bestehenden Softwarebausteine **KICKS** und **InHouseDB**
- Die Unterstützung des Chemikalienlagers ist ebenfalls Bestandteil der neuen UseCases

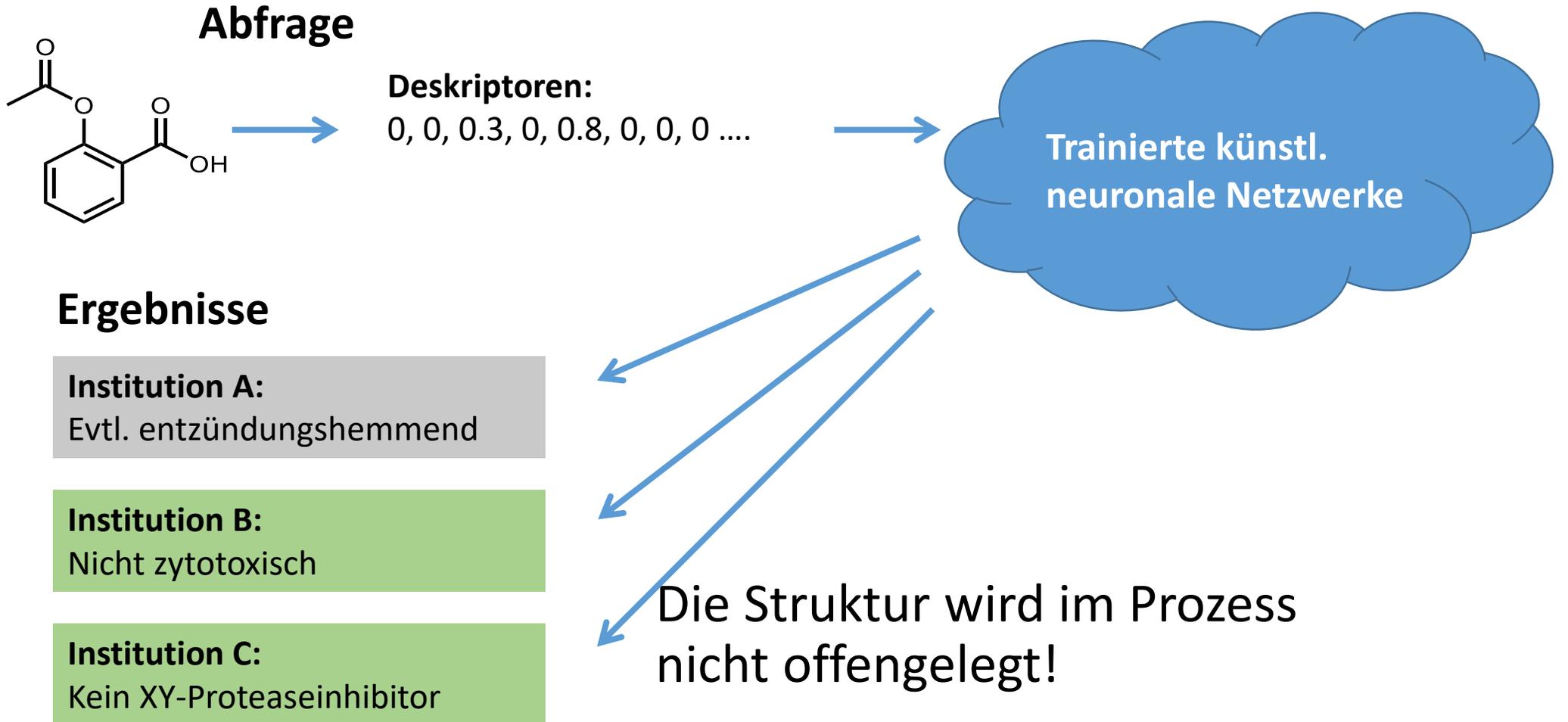


- Die Datenmigration soll verlustfrei und weitestgehend automatisiert geschehen
- Für die Installation ist das Leibniz-Institut für Neurobiologie (LIN) Magdeburg Wunschkandidat
  - Bestehende Installation der Software im Subnetz „BioActivesCloud“
  - Regelmäßige Kontakte mit Projektteilnehmern
  - Geografische Nähe



- Nach der Datenmigration ist zunächst eine Datensichtung mit dem Fokus zur Eignung für künstliche Intelligenz notwendig
- Für die Umsetzung der Entscheidungsunterstützungsfunktionalität müssen verschiedene Ansätze prototypisch getestet werden um eine fundierte Auswahl treffen zu können
- Weiterhin wird an der Erweiterung und Abrundung der Software gearbeitet

# KI: Verteiltes SAR



# KI / Becoming Smart: Textanalyse

- Ontologien
  - Kontrolliertes Vokabular
  - Tripelbasierter Datenspeicher („... ist ein ...“, „... wohnt in ...“, „... ernährt sich von ...“, ...)
  - Abfrage z.B. mit „Apache Jena“
- Datenextraktion: Regelbasierte Ansätze (reguläre Ausdrücke)
  - Liefert mit überschaubarem Aufwand erste Ergebnisse
  - Komplexität nimmt schnell bis zur Unbeherrschbarkeit zu → Sackgasse
- Neuartige Herangehensweisen: Künstliche Neuronale Netzwerke
  - Google BERT
    - Ressourcenbedarf: min. 64 Gbyte RAM
    - Trainingszeit 30 Tage@24 Cores, 1 Tag@GPU Tesla V100

## Context (Wikipedia: Herbal medicines)

Herbal medicine is also called phytomedicine or phytotherapy. Paraherbalism describes alternative and pseudoscientific practices of using unrefined plant or animal extracts as unproven medicines or health-promoting agents. Paraherbalism differs from plant-derived medicines in standard pharmacology because it does not isolate or standardize biologically active compounds, but rather relies on the belief that preserving various substances from a given source with less processing is safer or more effective – for which there is no evidence. Herbal dietary supplements most often fall under the phytotherapy category.

### BERT: Question / Answer

- 1) What are alternative terms for herbal medicine?  
phytomedicine or phytotherapy
- 2) In which category fall herbal dietary supplements?  
phytotherapy
- 3) How would one describe the paraherbalistic practice of using unrefined extracts and unproven medicines?  
alternative and pseudoscientific

# KI: Herausforderungen

- Überführung der Prototypen in anwendungsbereiten Code ist nicht trivial
- Training erfordert hohe Datenmenge **und –qualität**
- Anwendung und Setup werden wahrscheinlich komplex (z.B. viele Meta-Parameter)
- Ressourcenanforderungen zu hoch für normale Cloud-Knoten  
→ Zukünftige Strategie: Offload zu externen Compute-Ressourcen