



# Software-Systeme datengetrieben analysieren

**Markus Harrer**  
Software Development Analyst

Twitter: @feststelltaste  
Blog: feststelltaste.de

**INNOQ**



THE ORIGINAL HORROR SHOW

LEGACY SYSTEMS

# DAS PROBLEM MIT LEGACY SYSTEME



FACHLICHE  
FEATURES



ERKANNTE  
FEHLER



ARCHITEKTUR

STRANGE  
- THINGS -

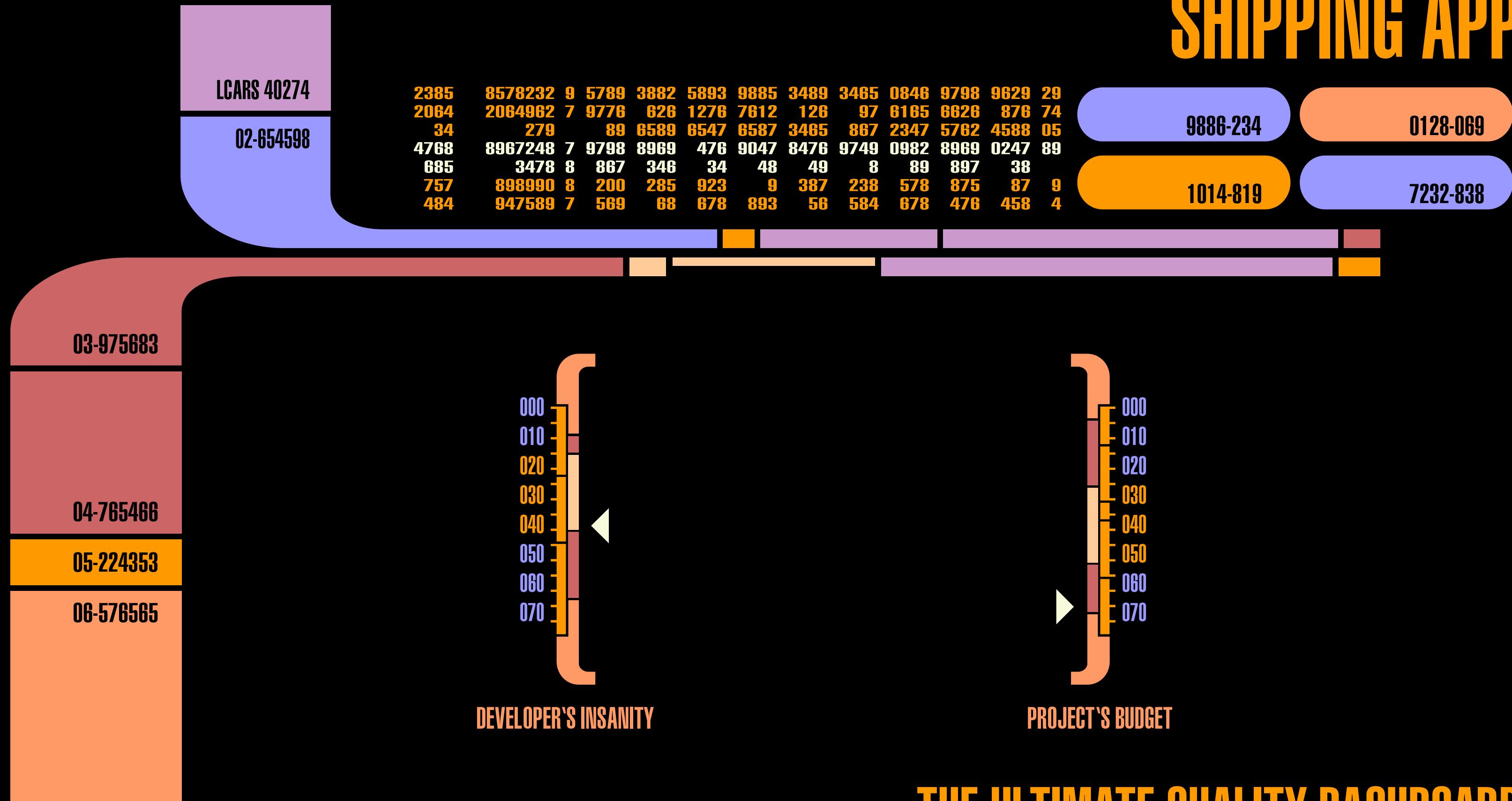
# MANAGEMENT

---

Kommunikationsbarriere

# ENTWICKLUNG

# SHIPPING APP



## THE ULTIMATE QUALITY DASHBOARD

# THE EMPIRIC STRIKES BACK

Not a long time ago, from brains  
not far, far away....

**SOFTWARE  
ANALYTICS**

# SOFTWARE ANALYTICS

A definition of

**MENZIES & ZIMMERMANN**

**Software Analytics**

*is analytics on software data for  
**managers** and **software engineers***

*with the aim of empowering  
**software development individuals**  
and teams*

*to gain and share insight from  
**their data** to make better  
decisions.*



**DIE SOFTWAREDATEN VON  
TEAM  
VON ENTWICKLERN**

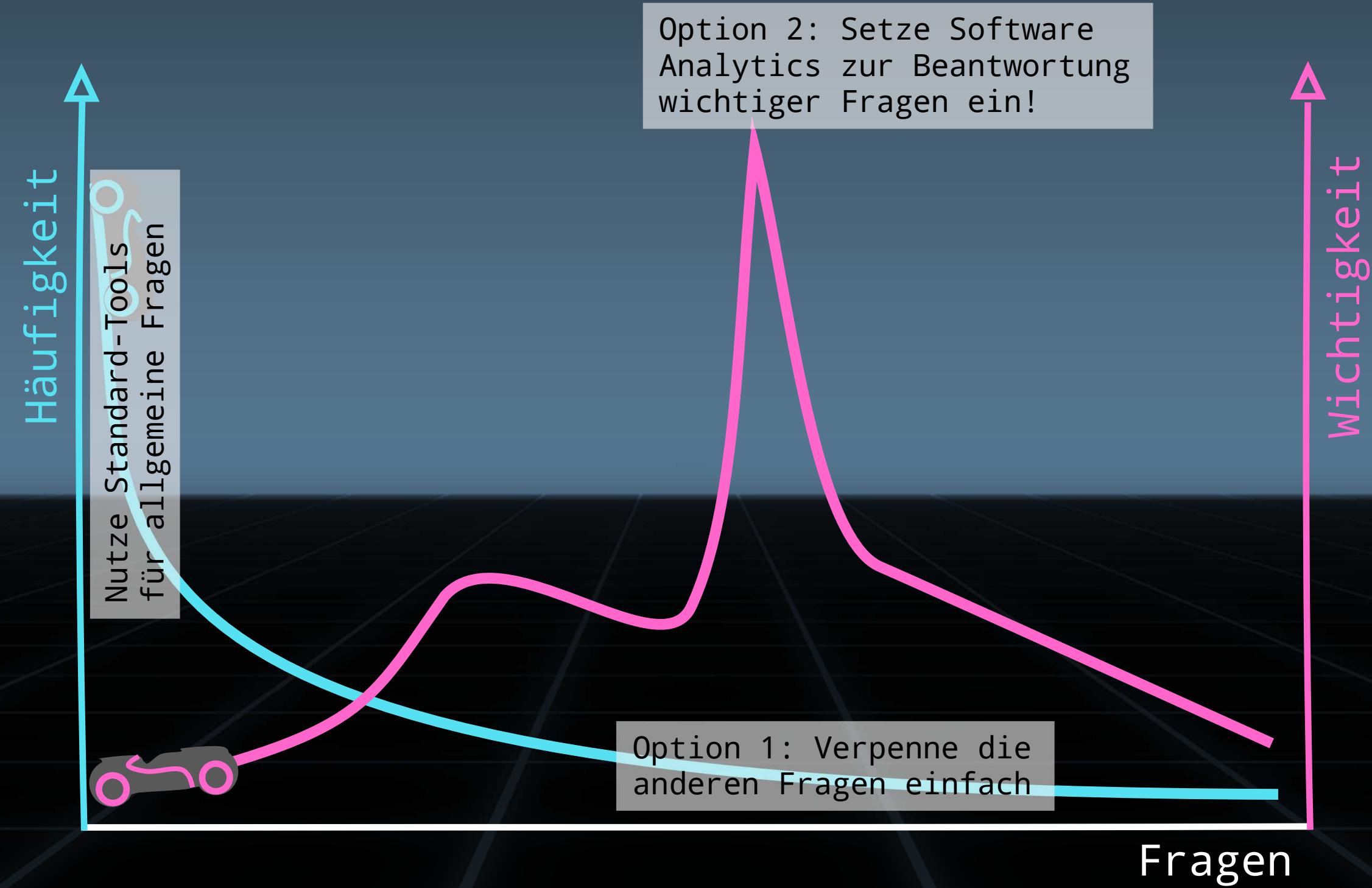
Statische  
Daten

Laufzeit-  
daten

Chronolo-  
gische  
Daten

Daten der  
Community

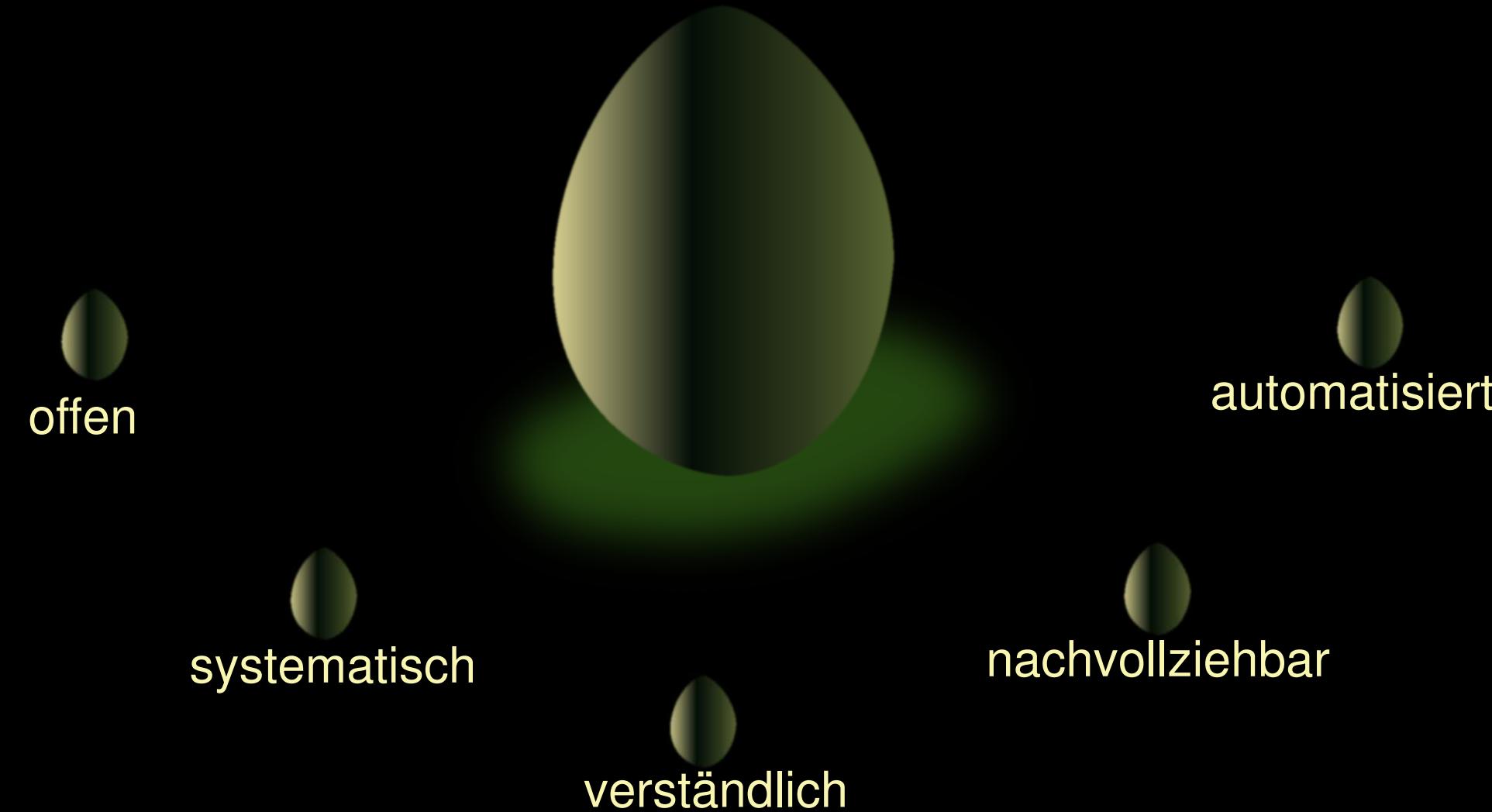
# SOFTWARE ANALYTICS: KEINE OASIS FRAGEN



# BEISPIELE VON DATENANALYSEN ZUM TERMINIEREN VON PROBLEmen

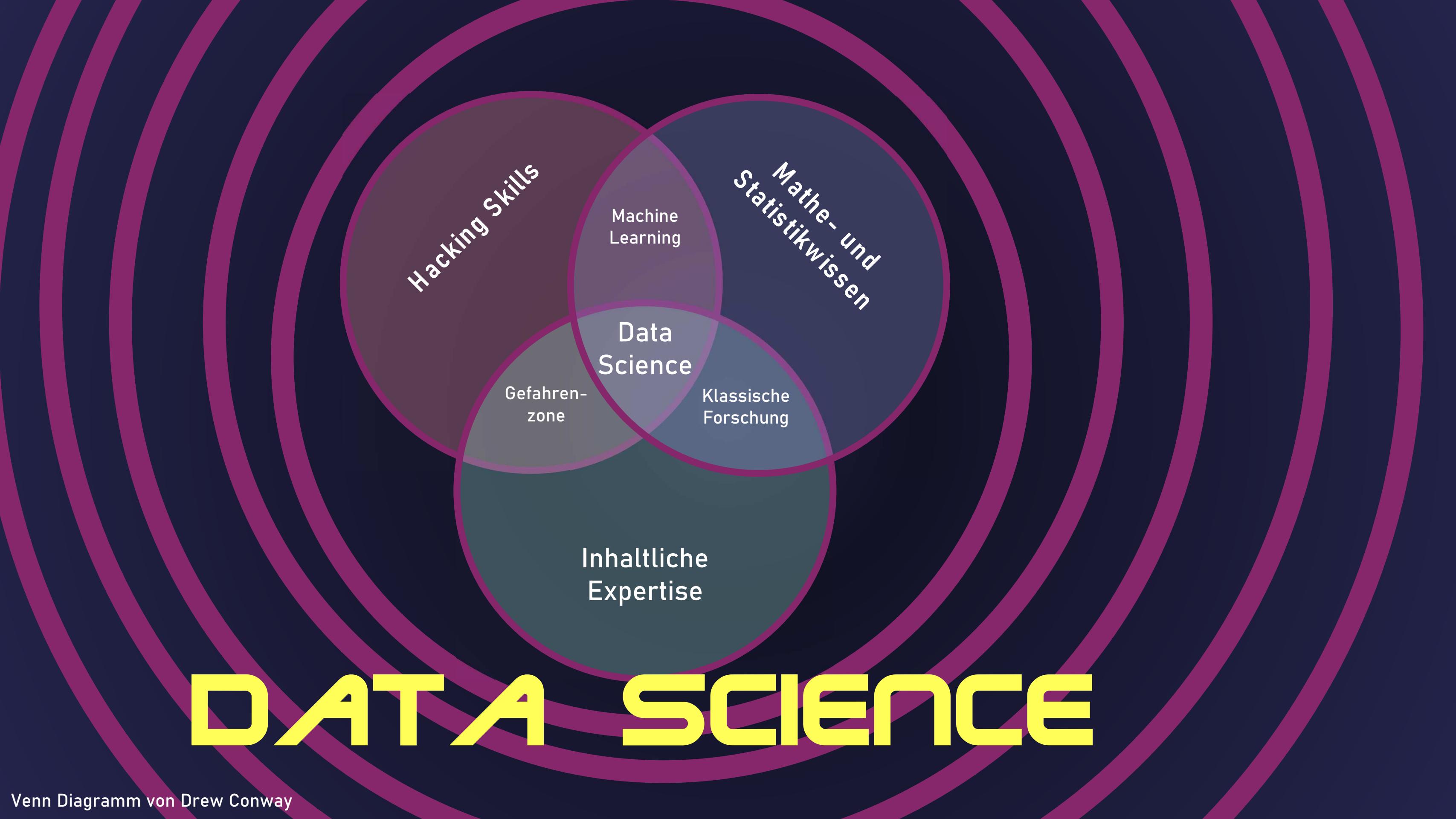
- Quantifizierung des Wissensverlusts bei Entwicklerfluktuation
- Verprobung von Modularisierungsvarianten (“virtuelles Refactoring”)
- Ermittlung von Performance-Hotspots über Call-Tree-Analyse
- Tracking von umfangreichen Code-Umbauarbeiten
- Analyse der Community-Aktivitäten um Open-Source-Software
- <weitere, ganz spezifische Analysen in ganz spezifischen Situationen>

# REPRODUCIBLE DATA ANALYSIS



= EIN WEG ZUR UMSETZUNG VON  
**SOFTWARE ANALYTICS**

# DATA SCIENCE



TOOL TIME

# Python 3



# PANDAS



... and matplotlib, numpy, scikit-learn, NLTK, Pygments, py2neo, requests, BeautifulSoup, Pygal ...

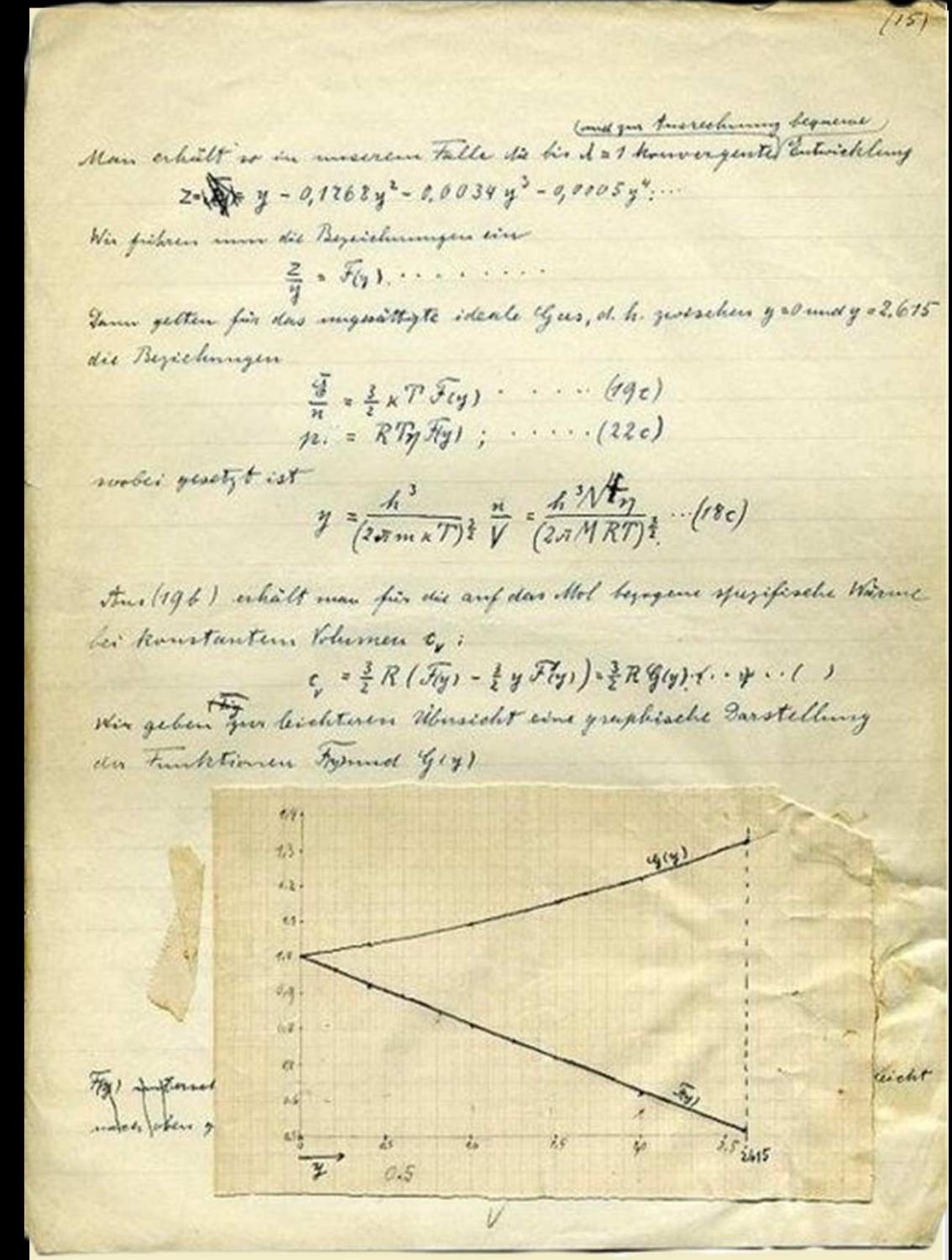
# Computational Notebook

code and data in love

# Computational Notebook

## Jupyter Notebook

- Kontext dokumentiert
- Ideen, Daten, Annahmen und Vereinfachungen aufgeführt
- Berechnungen verständlich dargelegt
- Zusammenfassungen erklärt
- Komplett automatisiert



# Jupyter Notebook

jupyter Production Coverage Demo Notebook Last Checkpoint: 2 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

File Cell Kernel Help

File Cell Kernel Help

## Context

John Doe remarked in [#AP1432](#) that there may be too much code in our application that isn't used at all. Before migrating the application to the new platform, we have to analyze which parts of the system are still in use and which are not.

## Idea

To understand how much code isn't used, we recorded the executed code in production with the coverage tool [JaCoCo](#). The measurement took place between 21st Oct 2017 and 27st Oct 2017. The results were exported into a CSV file using the JaCoCo command line tool with the following command:

```
java -jar jacococli.jar report "C:\Temp\jacoco.exec" --classfiles \
C:\dev\repos\buschmais-spring-petclinic\target\classes --csv jacoco.csv
```

The CSV file contains all lines of code that were passed through during the measurement's time span. We just take the relevant data and add an additional LINES column to be able to calculate the ratio between covered and missed lines later on.



## Context

John Doe remarked in [#AP1432](#) that there may be too much code in our application that isn't used at all. Before migrating the application to the new platform, we have to analyze which parts of the system are still in use and which are not.

## Idea

To understand how much code isn't used, we recorded the executed code in production with the coverage tool [JaCoCo](#). The measurement took place between 21st Oct 2017 and 27st Oct 2017. The results were exported into a CSV file using the JaCoCo command line tool with the following command:

```
java -jar jacococli.jar report "C:\Temp\jacoco.exec" --classfiles \
C:\dev\repos\buschmais-spring-petclinic\target\classes --csv jacoco.csv
```

The CSV file contains all lines of code that were passed through during the measurement's time span. We just take the relevant data and add an additional LINES column to be able to calculate the ratio between covered and missed lines later on.

```
In [1]:  
1 import pandas as pd  
2 coverage = pd.read_csv("../input/spring-petclinic/jacoco.csv")  
3 coverage = coverage[['PACKAGE', 'CLASS', 'LINE_COVERED', 'LINE_MISSED']]  
4 coverage['LINES'] = coverage.LINE_COVERED + coverage.LINE_MISSED  
5 coverage.head(1)
```

Out[1]:

	PACKAGE	CLASS	LINE_COVERED	LINE_MISSED	LINES
0	org.springframework.samples.petclinic	PetclinicInitializer	24	0	24

```
In [1]: 1 import pandas as pd  
2 coverage = pd.read_csv("../input/spring-petclinic/jacoco.csv")  
3 coverage = coverage[['PACKAGE', 'CLASS', 'LINE_COVERED', 'LINE_MISSED']]  
4 coverage['LINES'] = coverage.LINE_COVERED + coverage.LINE_MISSED  
5 coverage.head(1)
```

Out[1]:

	PACKAGE	CLASS	LINE_COVERED	LINE_MISSED	LINES
0	org.springframework.samples.petclinic	PetclinicInitializer	24	0	24

## Analysis

It was stated that whole packages wouldn't be needed anymore and that they could be safely removed. Therefore, we sum up the coverage data per class for each package and calculate the coverage ratio for each package.

```
In [2]: 1 grouped_by_packages = coverage.groupby("PACKAGE").sum()  
2 grouped_by_packages['RATIO'] = grouped_by_packages.LINE_COVERED / grouped_by_packages.LINES  
3 grouped_by_packages = grouped_by_packages.sort_values(by='RATIO')  
4 grouped_by_packages
```

Out[2]:

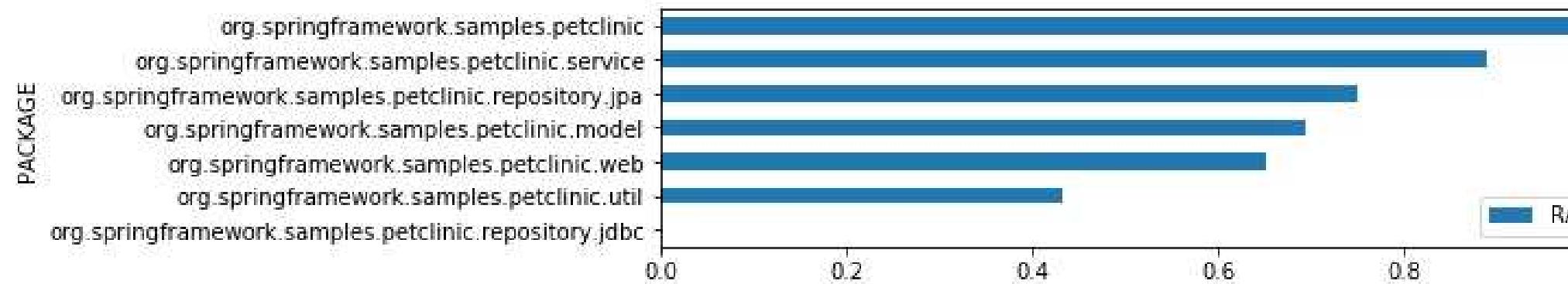
PACKAGE	LINE_COVERED	LINE_MISSED	LINES	RATIO
org.springframework.samples.petclinic.repository.jdbc	0	152	152	0.000000
org.springframework.samples.petclinic.util	13	17	30	0.433333
org.springframework.samples.petclinic.web	75	40	115	0.652174
org.springframework.samples.petclinic.model	75	33	108	0.694444
org.springframework.samples.petclinic.repository.jpa	21	7	28	0.750000
org.springframework.samples.petclinic.service	16	2	18	0.888889
org.springframework.samples.petclinic	24	0	24	1.000000

	LOC	TC	CC	Coverage
org.springframework.samples.petclinic.web	75	40	115	0.652174
org.springframework.samples.petclinic.model	75	33	108	0.694444
org.springframework.samples.petclinic.repository.jpa	21	7	28	0.750000
org.springframework.samples.petclinic.service	16	2	18	0.888889
org.springframework.samples.petclinic	24	0	24	1.000000

We plot the data for the coverage ratio to get a brief overview of the result.

```
In [3]: 1 %matplotlib inline
2 grouped_by_packages[['RATIO']].plot(kind="barh", figsize=(8,2))
```

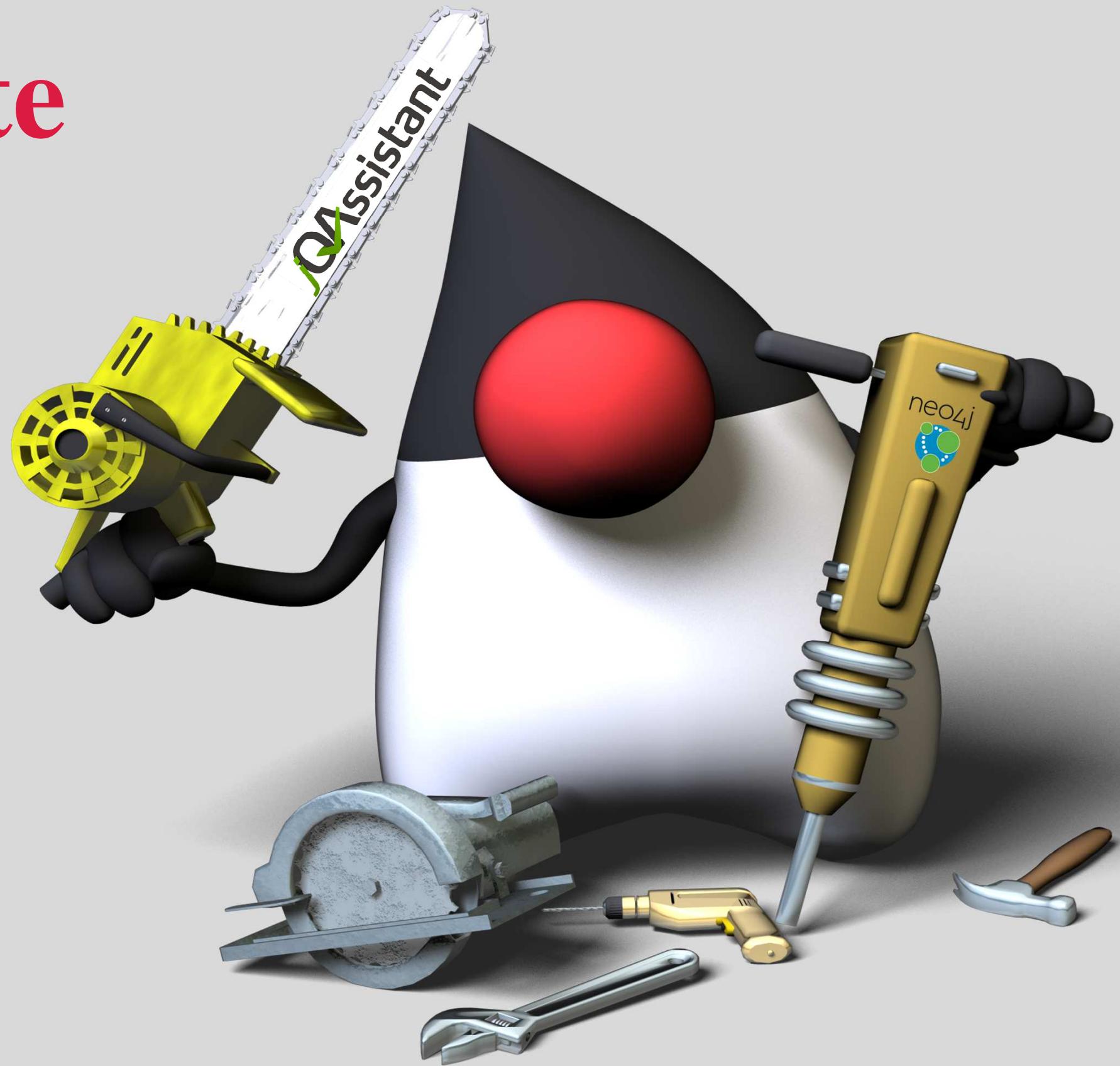
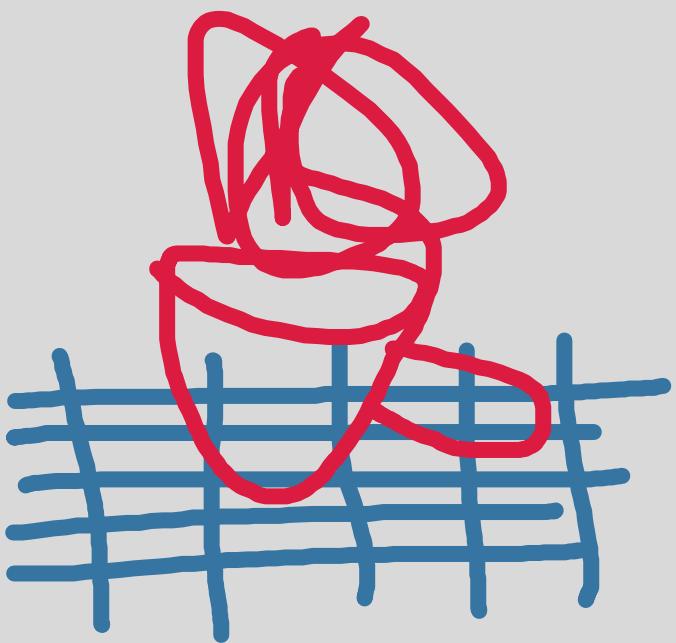
Out[3]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1874cdde9e8>



## Conclusion

The JDBC package `org.springframework.samples.petclinic.repository.jdbc` isn't used at all and can be left out safely when migrating to the new platform.

# Graph-basierte Analyse

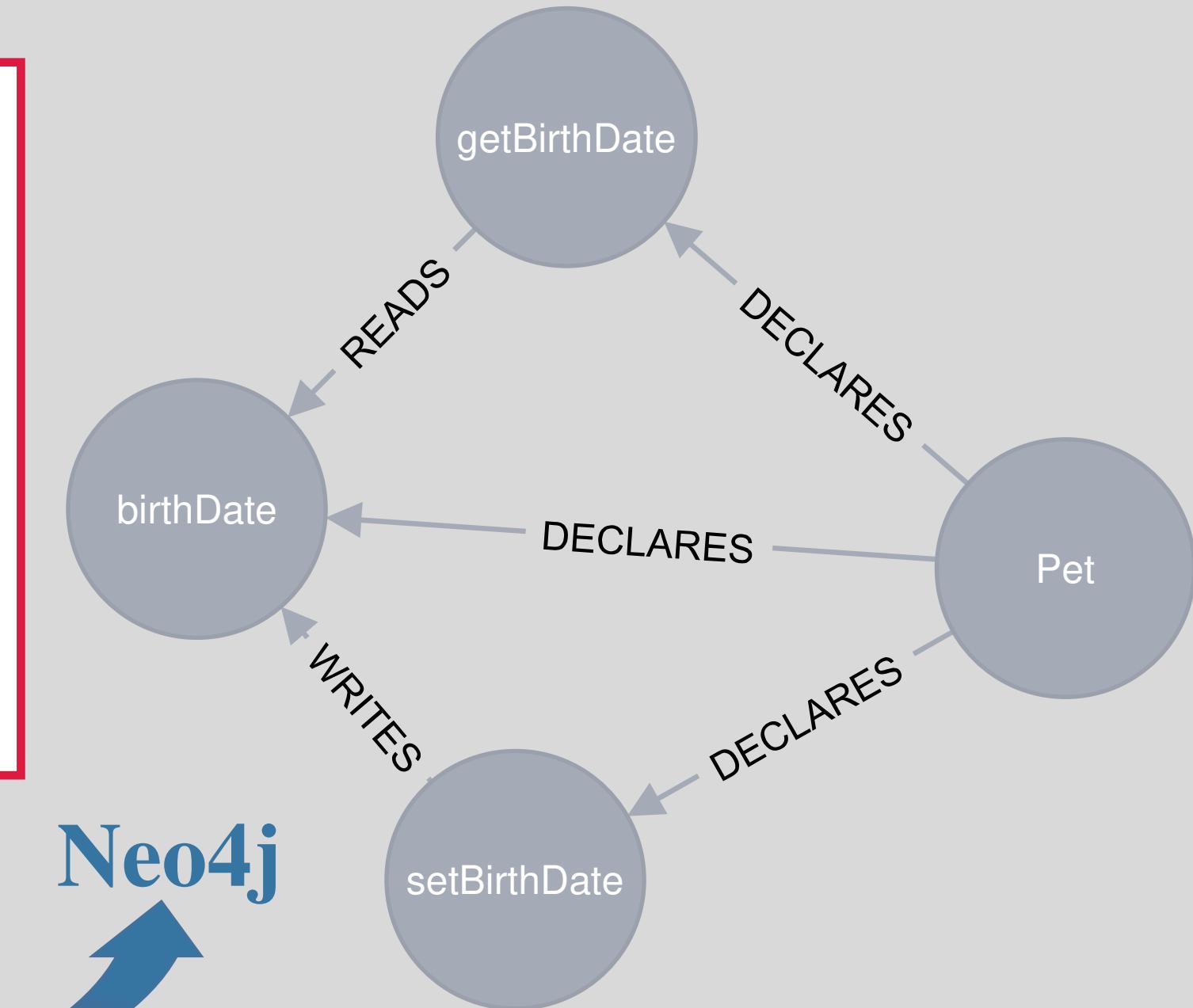


# jQAssistant & Neo4j

```
public class Pet {  
  
    private LocalDate birthDate;  
  
    public LocalDate getBirthDate() {  
        return this.birthDate;  
    }  
  
    public void setBirthDate(LocalDate birthDate) {  
        this.birthDate = birthDate;  
    }  
}
```

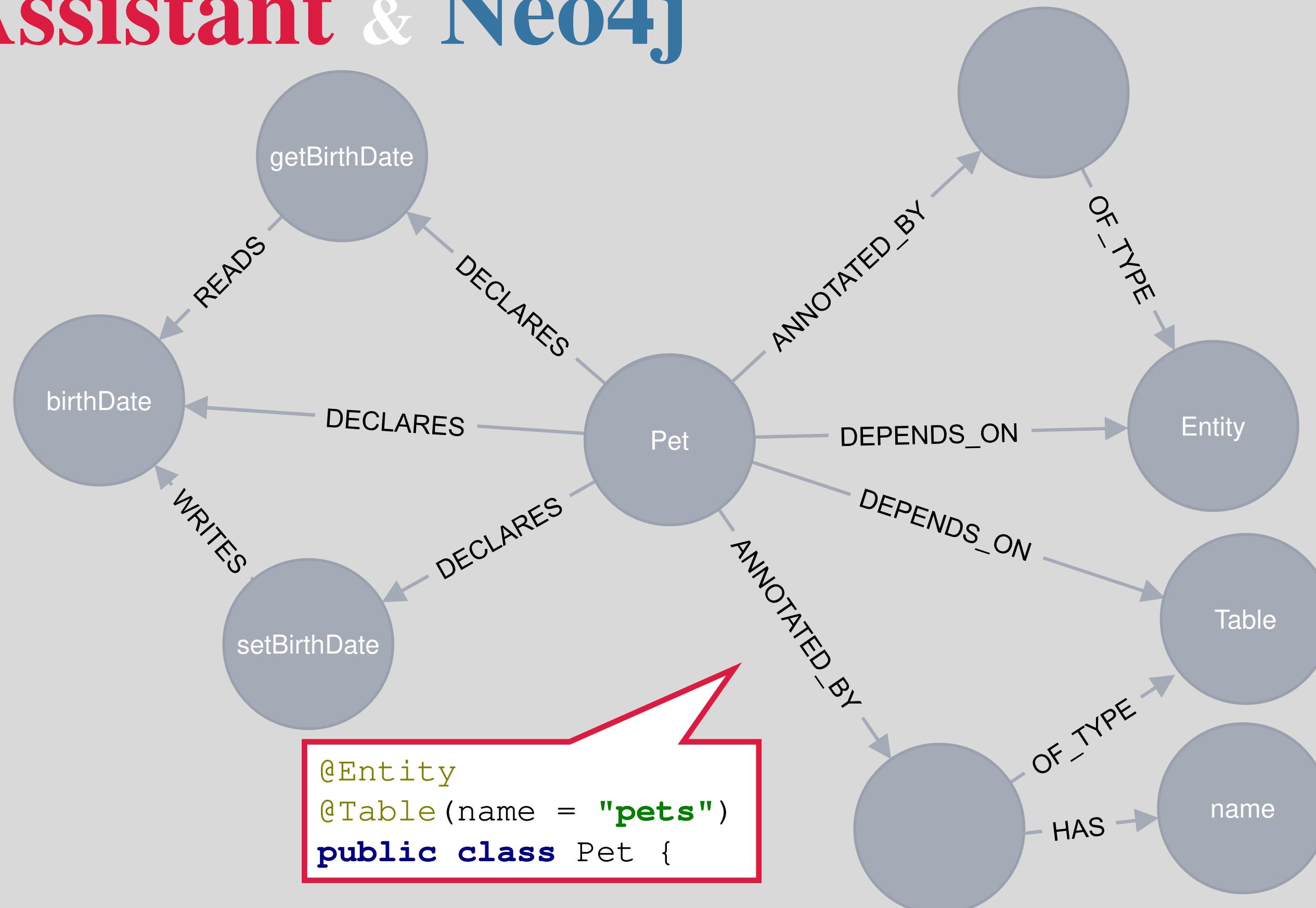
Java Code

jQAssistant



Neo4j

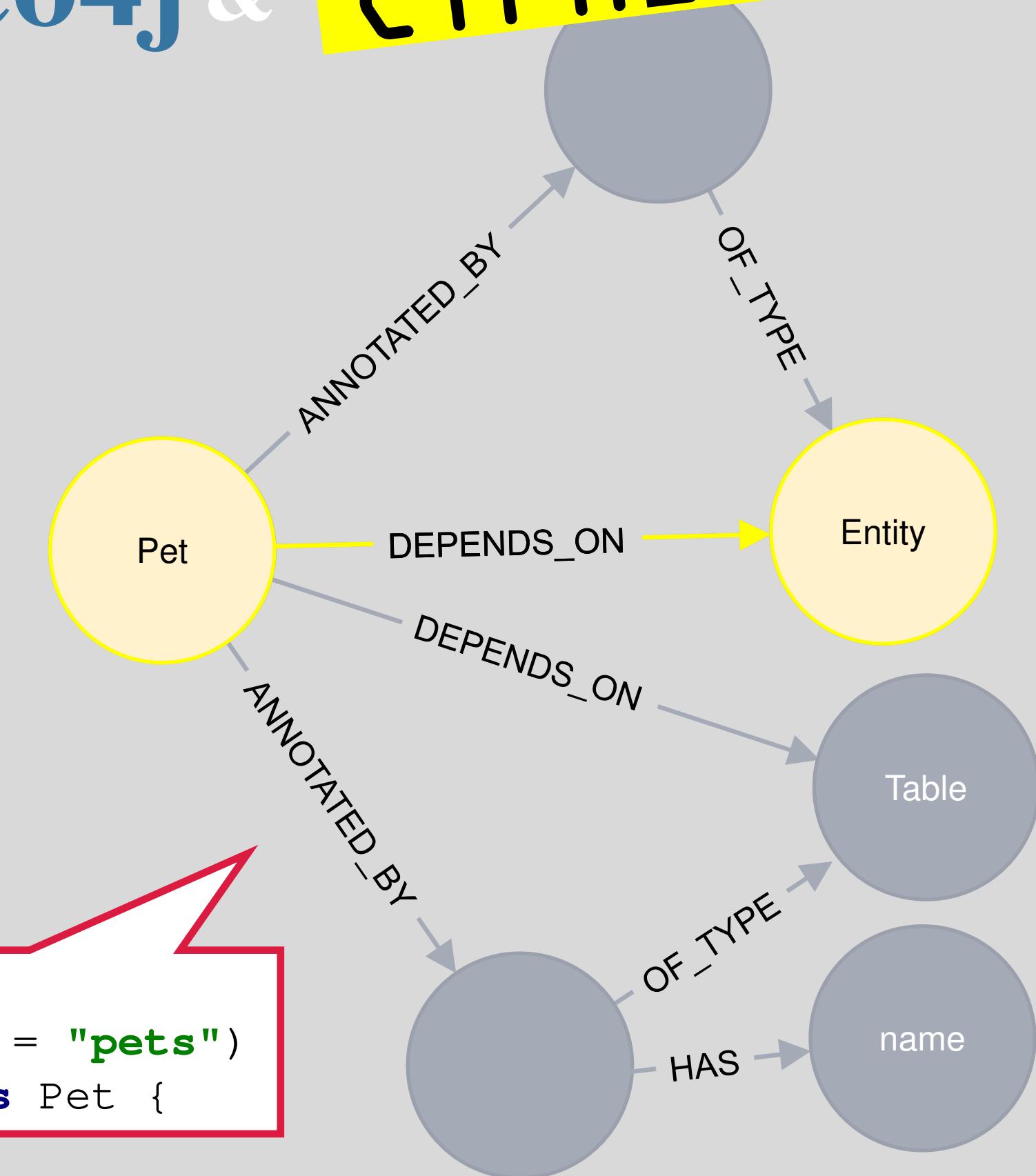
# jQAssistant & Neo4j



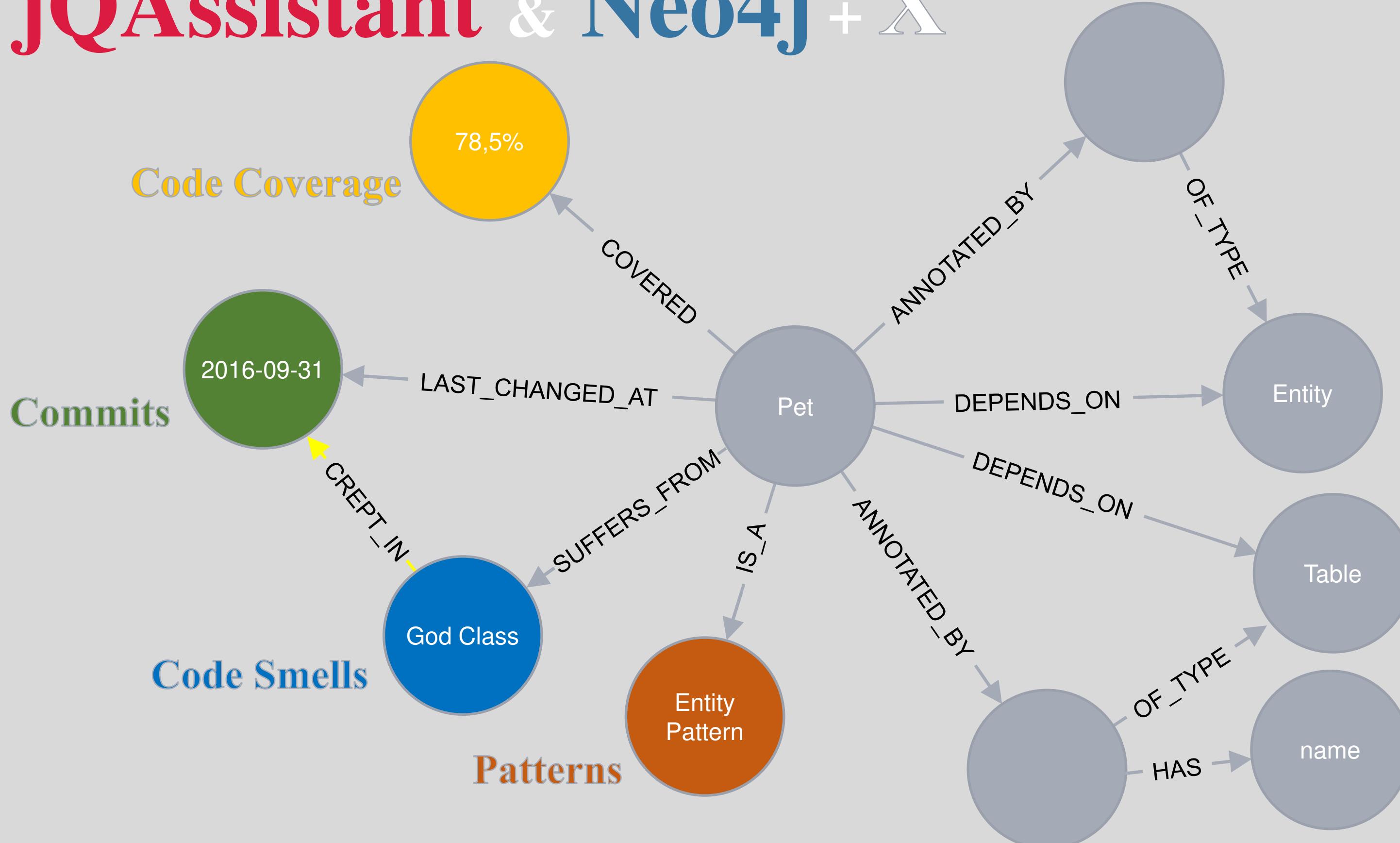
# jQAssistant & Neo4j &

# CYPHER

```
MATCH
  (p:Class) - [:DEPENDS_ON] -> (e:Type)
WHERE
  p.name = "Pet" AND
  e.name = "Entity"
RETURN p, e
```



# jQAssistant & Neo4j + X





WERDE WIEDER  
**HERR DER DINGE**  
MIT DATENGETRIEBENEN SOFTWAREANALYSEN

NO METRIC TO  
RULE THEM ALL

DIE GEFAHREN

DIE ZWEI TIPPS

NACHVOLLZIEHBARKEIT

AUTOMATISIERUNG

ANZAHL DER  
GELÖSTEN PROBLEME

DIE RÜCKKEHR  
DER VERNUNFT



META METRIC

# QUESTION



ASK ' EM ALL

# Infos zu den Demos

## Jupyter Notebook, Python, pandas, matplotlib

### Repo

[https://github.com/feststelltaste/software-analytics/tree/master/demos/20201021\\_INNOQ\\_Technology\\_Lunch](https://github.com/feststelltaste/software-analytics/tree/master/demos/20201021_INNOQ_Technology_Lunch)

### Interaktive Online-Version

[https://mybinder.org/v2/gh/feststelltaste/software-analytics/master?filepath=demos%2F20201021\\_INNOQ\\_Technology\\_Lunch%2FRefactoring.ipynb](https://mybinder.org/v2/gh/feststelltaste/software-analytics/master?filepath=demos%2F20201021_INNOQ_Technology_Lunch%2FRefactoring.ipynb)

## jQAssistant & Neo4j

### Repo Spring PetClinic

<https://github.com/javaonautobahn/spring-petclinic>

### Repo DesignSmells

<https://github.com/feststelltaste/designsmells>



# Einstieg in Software Analytics

## Mein Blog

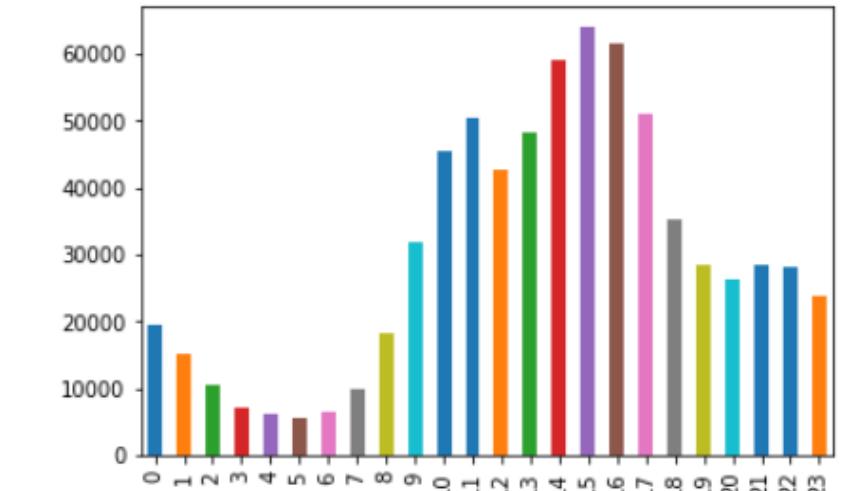
<https://feststelltaste.de>

We can display the result by means of a bar chart and thus get a

```
commits_per_hour.plot.bar();
```

## TOP 5 Software Analytics

<https://www.feststelltaste.de/top-5-software-analytics/>



## Mein Software Analytics Repository

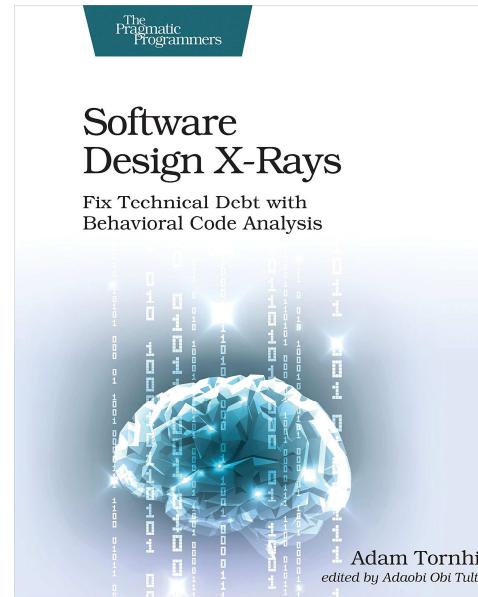
<https://github.com/feststelltaste/software-analytics>

We now additionally label the plot. To do this, we store the return value of the plotting library `matplotlib`, through which we can customize

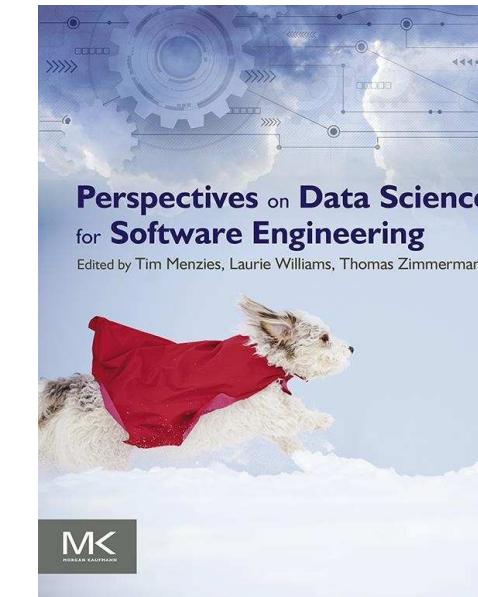
## Mini-Tutorial und mehr zu Software Analytics

<https://github.com/feststelltaste/software-analytics-workshop>

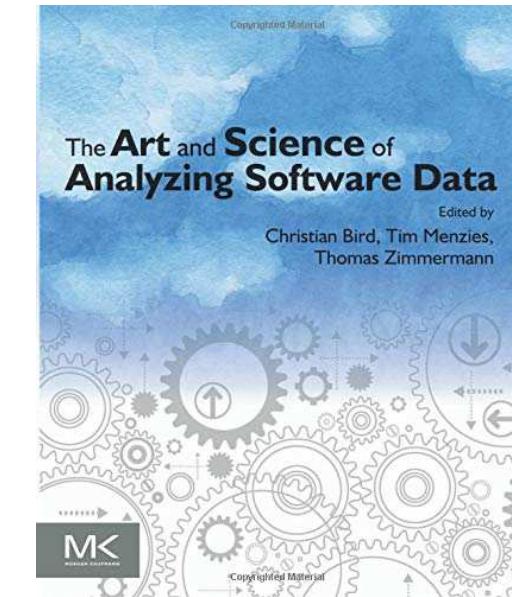
# Mehr Infos zu Software Analytics



Adam Tornhill:  
*Software X-Ray*



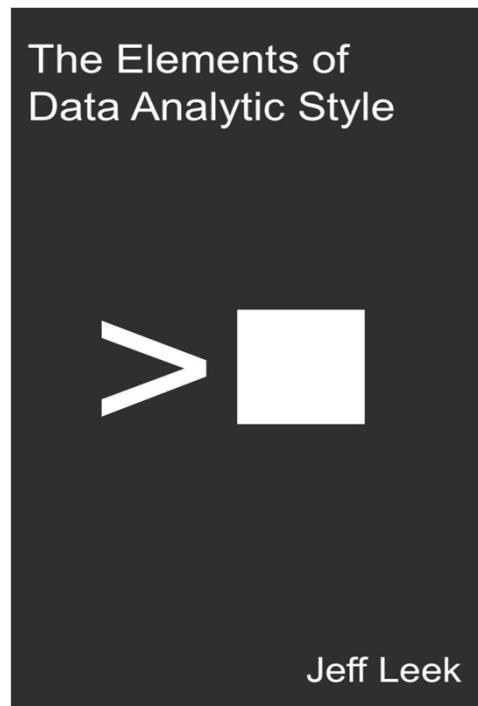
Tim Menzies, Laurie Williams,  
Thomas Zimmermann:  
*Perspectives on Data Science for  
Software Engineering*



Christian Bird, Tim Menzies,  
Thomas Zimmermann:  
*The Art and Science of Analyzing  
Software Data*

Noch mehr: <https://github.com/feststelltaste/awesome-software-analytics>

# Infos zu Data Science

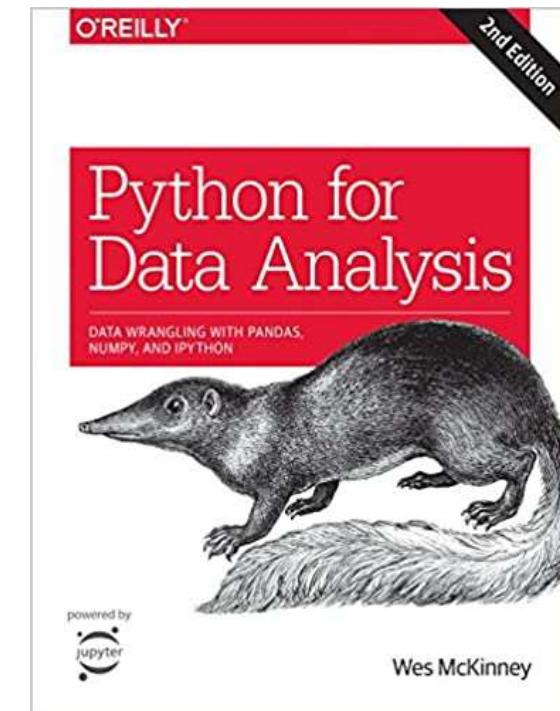


Jeff Leek:  
*The Elements of Data Analytic Style*

Report Writing for Data Science in R



Roger D. Peng



Wes McKinney:  
*Python for Data Analysis*

## Meine TOP5s

- <https://www.feststelltaste.de/top5-jupyter/>
- <https://www.feststelltaste.de/top5-pandas/>
- <https://www.feststelltaste.de/top5-python/>

# Weitere Details zu jQAssistant/Neo4j

<https://easychair.org/publications/preprint/893N>

## Towards an Open Source Stack to Create a Unified Data Source for Software Analysis and Visualization

Richard Müller\*, Dirk Mahler†, Michael Hunger‡, Jens Nerche§ and Markus Harrer¶

\*Leipzig University, Germany

Email: rmueller@wifa.uni-leipzig.de

†buschmais GbR, Dresden, Germany

Email: dirk.mahler@buschmais.com

‡Developer Relations, Neo4j Inc., Malmö, Sweden

Email: michael.hunger@neo4j.com

§Application Development, Kontext E GmbH, Dresden, Germany

Email: j.nerche@kontext-e.de

¶Software Development Analyst, Freelancer, Roth, Germany

Email: contact@markusharrer.de

**Abstract**—The beginning of every software analysis and visualization process is data acquisition. However, there are various sources of data about a software system. The methods used

Creating, storing, and querying the data captured by such graphs is very challenging. Diehl et al. summarize the most important questions in this respect [2].

# Software Analytics Trainings



Öffentliche Termine ab Januar 2021 remote, firmenintern nach Abstimmung

<https://www.innoq.com/de/trainings/software-analytics>

# Vielen Dank!

**INNOQ**  
[www.innoq.com](http://www.innoq.com)

Markus Harrer

[markus.harrer@innoq.com](mailto:markus.harrer@innoq.com)

 @feststelltaste

innoQ Deutschland GmbH

Krischerstr. 100  
40789 Monheim am Rhein  
Germany  
+49 2173 3366-0

Ohlauer Str. 43  
10999 Berlin  
Germany

Ludwigstr. 180E  
63067 Offenbach  
Germany

Kreuzstr. 16  
80331 München  
Germany

innoQ Schweiz GmbH

Gewerbestr. 11  
CH-6330 Cham  
Switzerland  
+41 41 743 01 11

Albulastr. 55  
8048 Zürich  
Switzerland