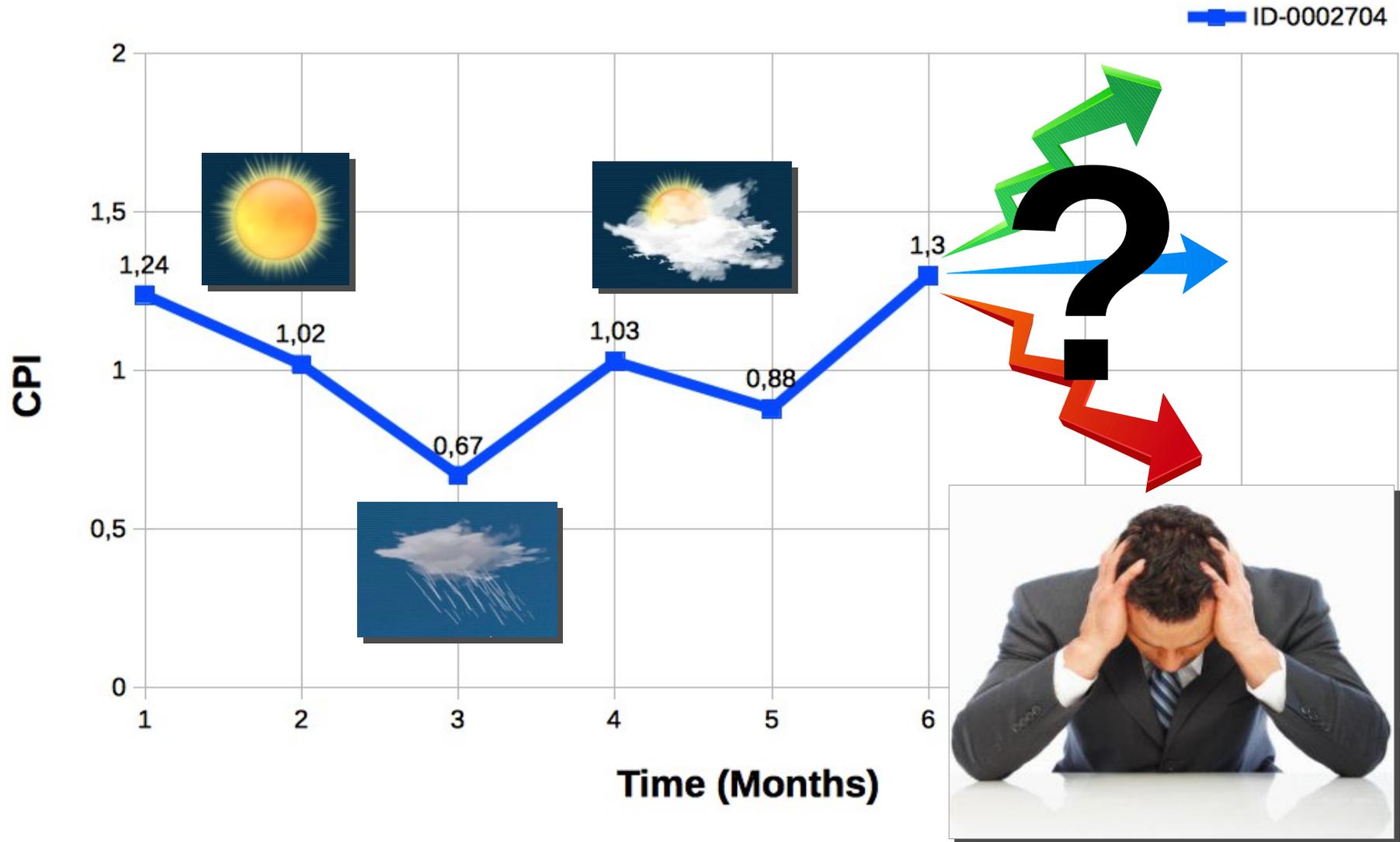


Werkzeugunterstützte Projektprognose

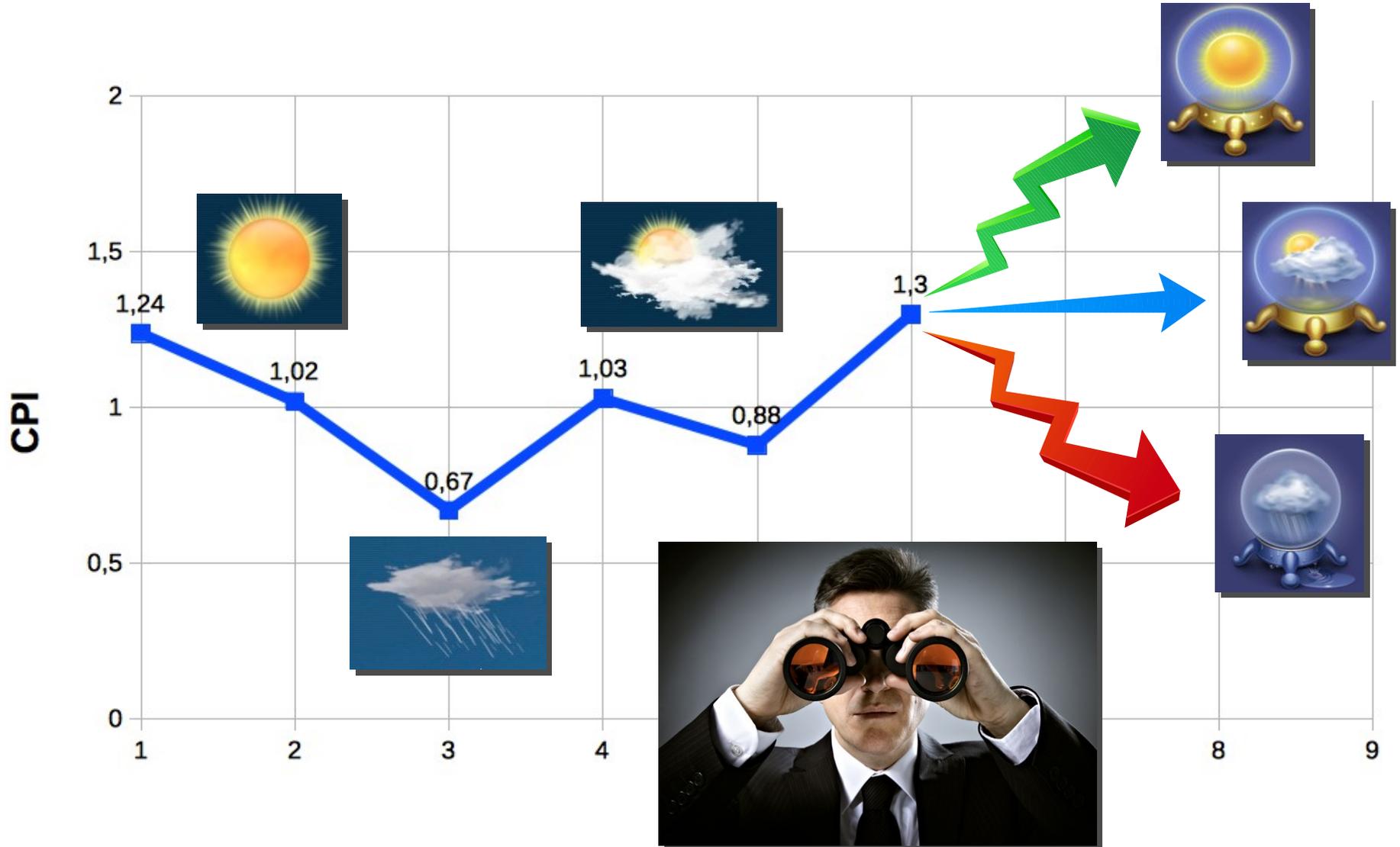
Elena Emelyanova
elena.emelyanova@rwth-
aachen.de

26.06.2015

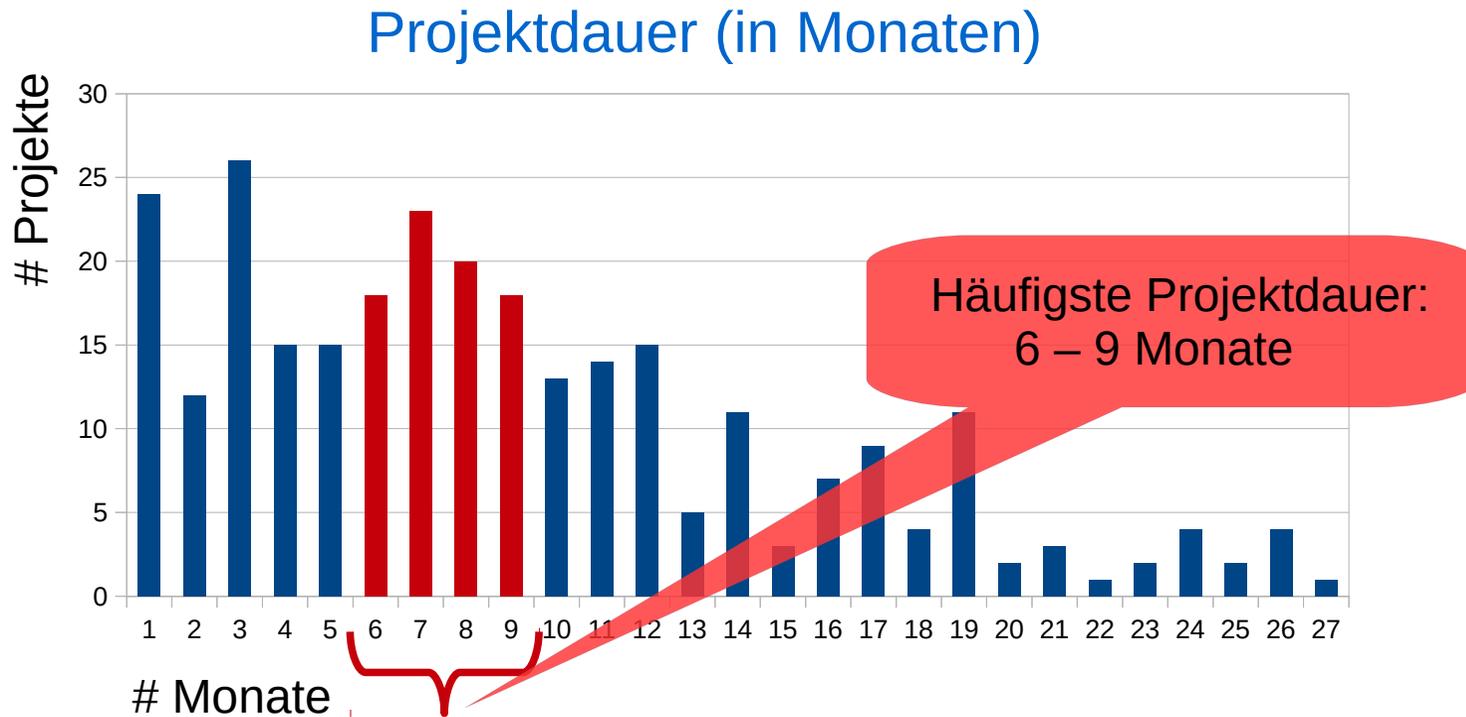
Motivation



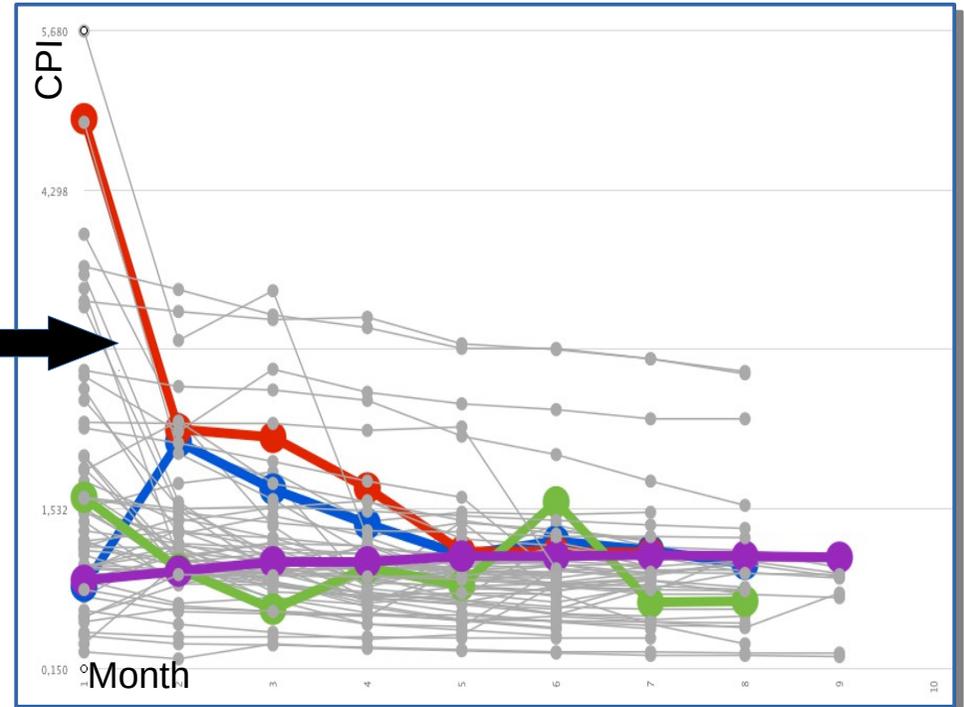
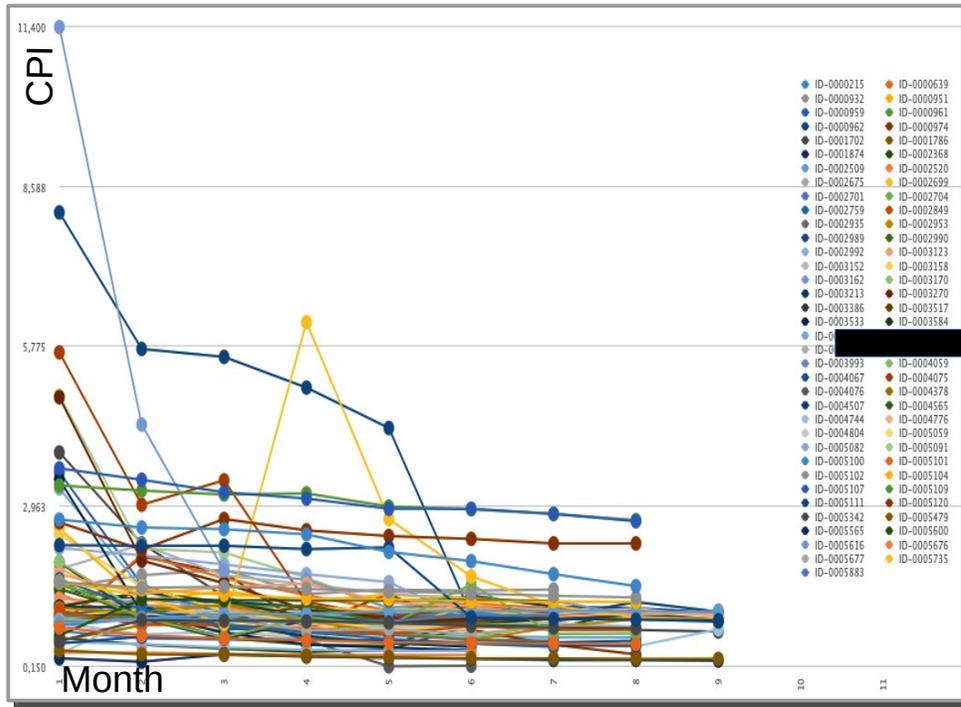
Motivation



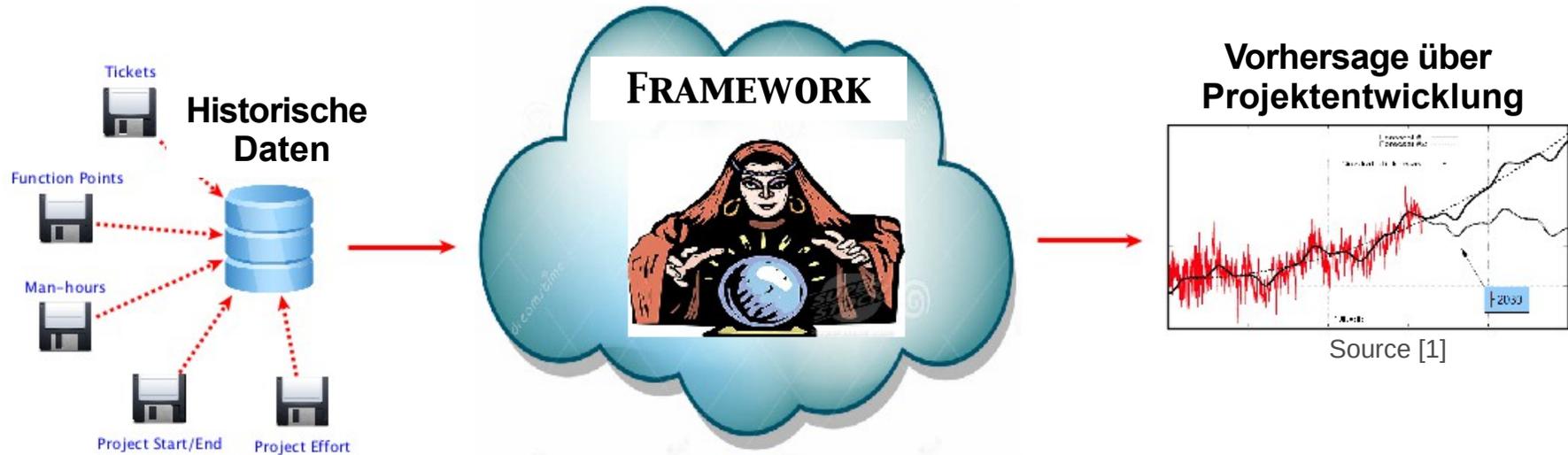
- Daten von externem Kooperationspartner:



- Erste Versuch Muster zu identifizieren

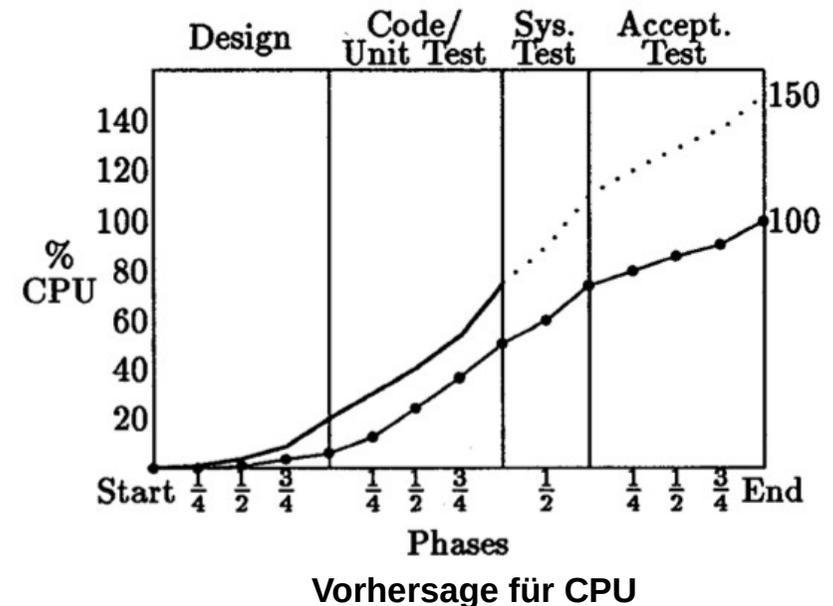


Motivation



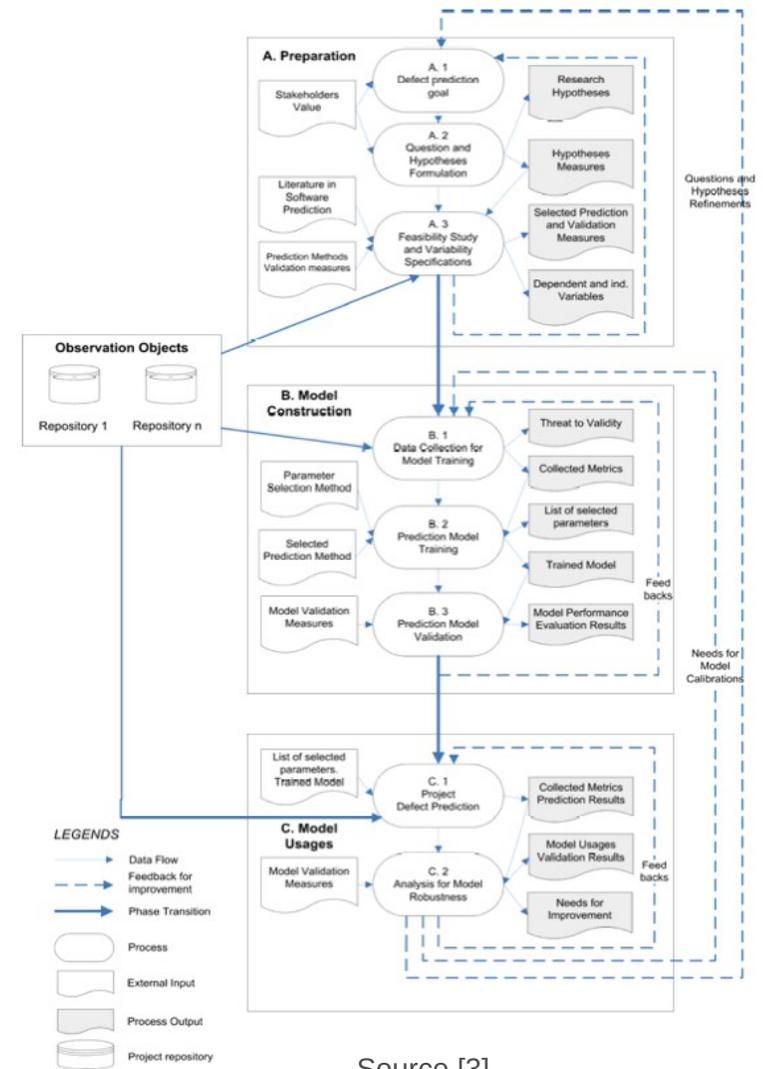
- T. Vi Bach, 2014 [2]
 - **Ansatz:** Konzept des Frameworks zur Projektprognose
 - **Daten:** Historische Daten von großem IT-Unternehmen (CPI, SPI)
 - **Implementierung:** nur Konzept, keine Instrumentalisierung!
 - **Evaluation:** Fehlerrate über 100%

- M.V. Zelkowitz et al., 1993 [4]
 - **Ansatz:** Konstruktion des Prognosemodells für Software Projekte
 - **Daten:** über 100 NASA Projekten (LOC, Reported Errors, CPU, etc.)
 - **Implementierung:** keine Information!



Verwandte Arbeiten

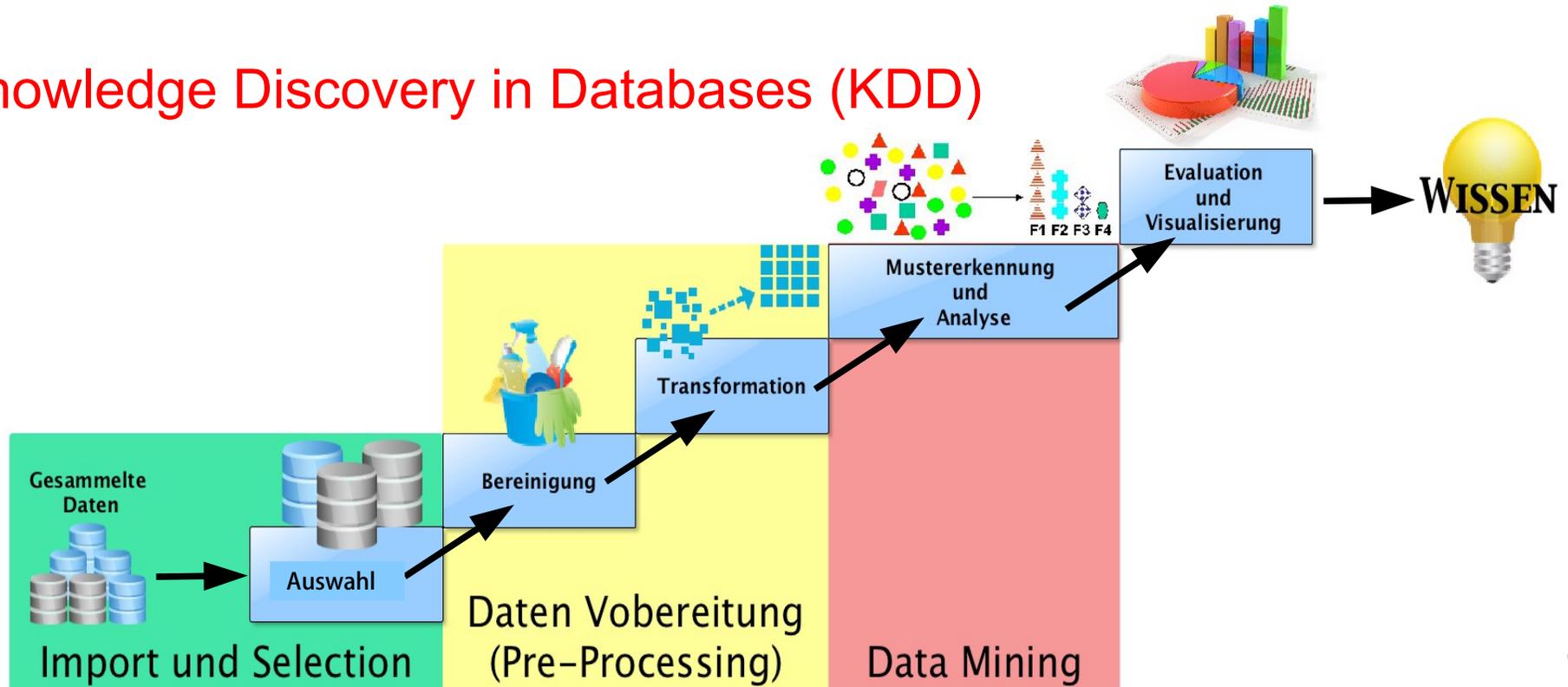
- R. Ramler et al., 2008 [3] and 2009 [5]
 - **Ansatz:** Framework für Vorhersage von Software Defekten basierend auf Data Mining Methoden
 - **Daten:** 6 Versionen von Software System, ca.157 Komponenten / 639 KLOCs
 - **Implementierung:** Definiert 3 Phasen, aber keine Info über Instrumentalisierung!
 - **Evaluation:** 22% < Fehlerrate < 33%



Source [3]

- Was ist **Data Mining**?
 - Entdecken und Extrahieren von Information
 - Finden von Muster / Wissensgewinnung

- **Knowledge Discovery in Databases (KDD)**





**Cluster-Analyse
(Clustering)**

Klassifikation

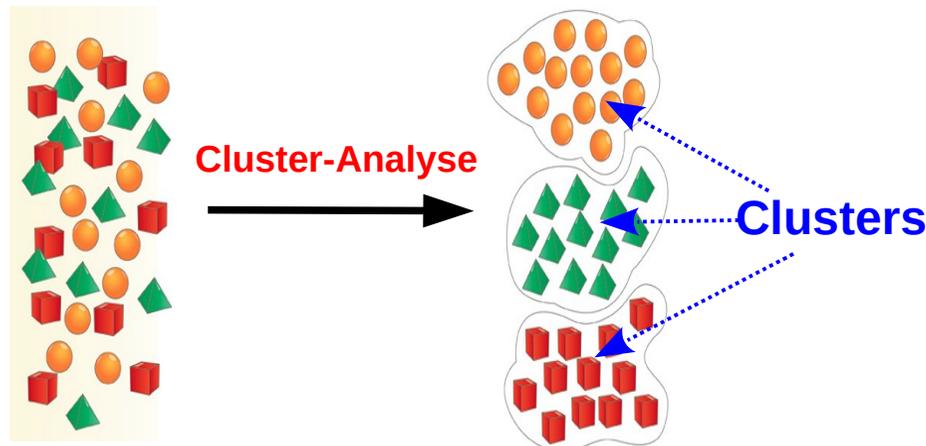
- Zusammenfassen von Objekten zu **Klassen**
- **Klasse** = Gruppe von Objekten mit einer logischen Eigenschaft (**Regel**)
- Benennung der Klassen



Cluster-Analyse (Clustering)

Klassifikation

- Entdeckung von **Ähnlichkeitsstrukturen**
- **Cluster** = Gruppe von “ähnlichen” Objekten gemäß eines **Ähnlichkeitsprinzips** (Distanzmaß)



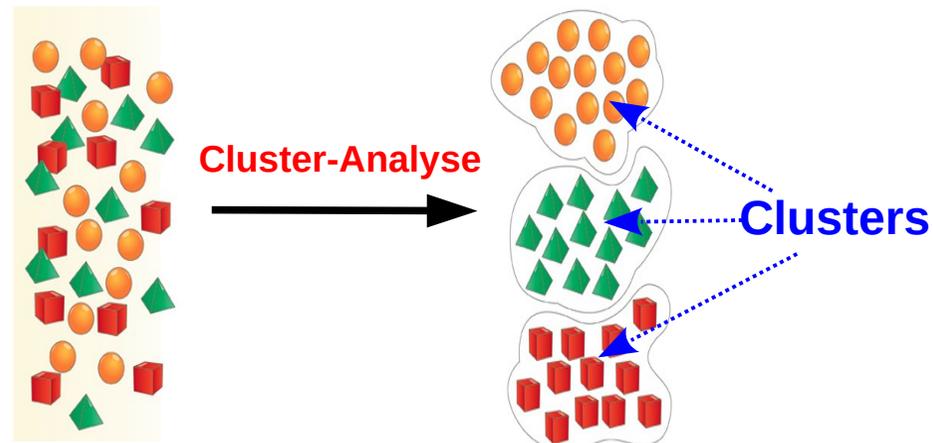


Cluster-Analyse (Clustering)

Klassifikation

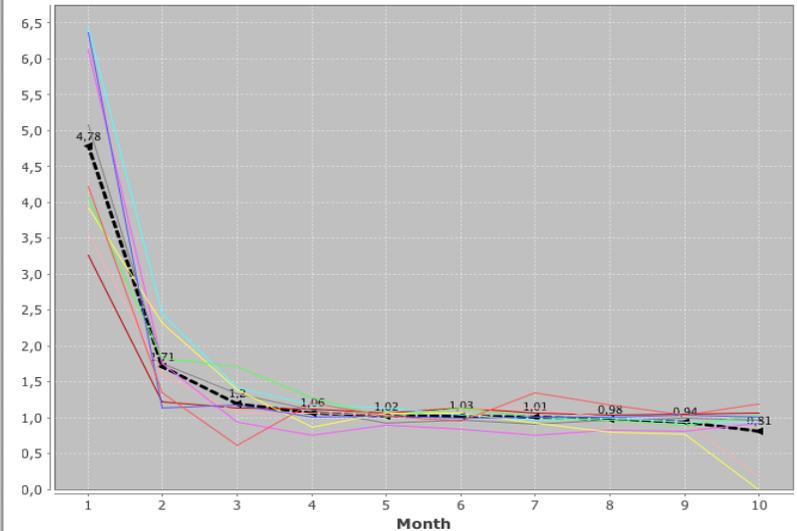
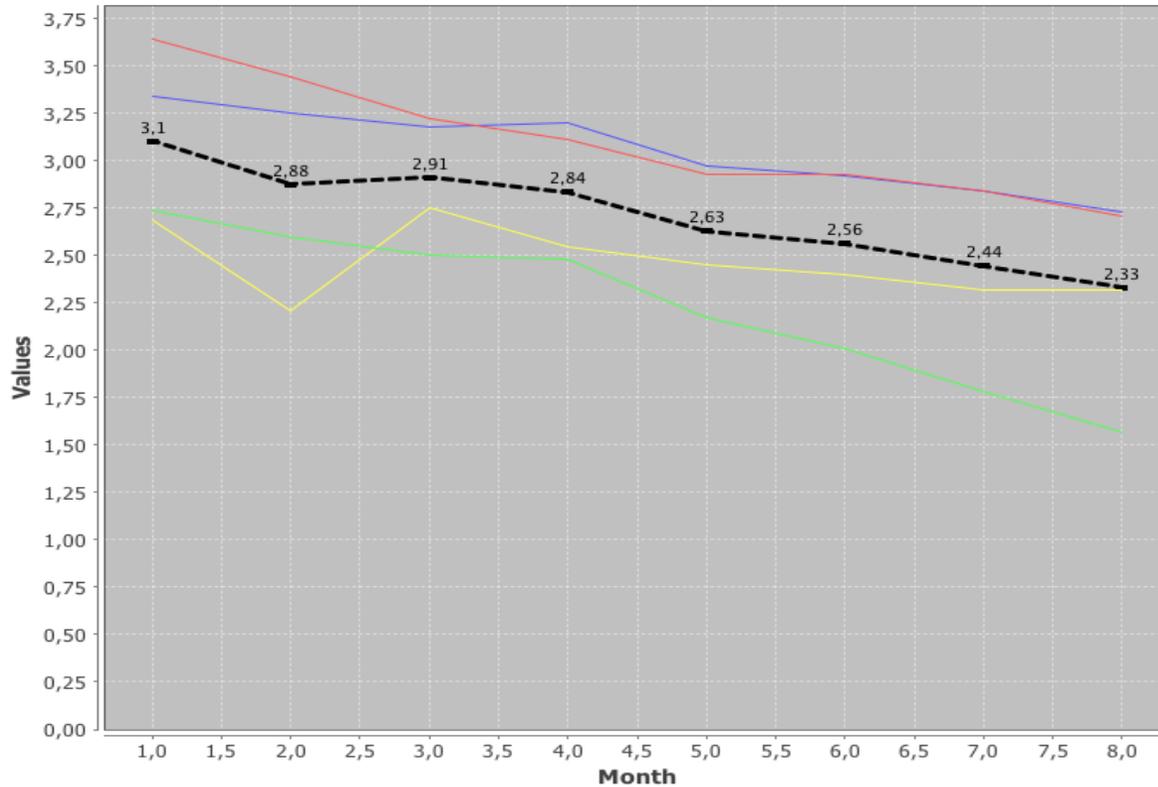
Clustering-Parameter

- Clustering-**Verfahren** (k-Means, Hierarchical Clustering)
- Anzahl von Clusters **k**
- ...



■ Ähnlichkeitsprinzip und Cluster-Repräsentant

Beispiel: Ergebnis der Cluster-Analyse

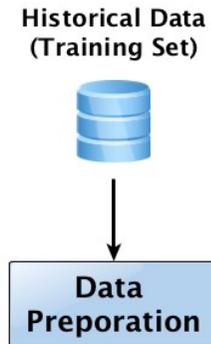


values for Project ID-0005566 — CPI values for Project ID-0004369 — CPI values for Project ID-0004440
 values for Project ID-0005342 — CPI values for Project ID-0005500 — CPI values for Project ID-0004058
 values for Project ID-0005107 — CPI values for Project ID-0005094 — CPI values for Project ID-0005093
 Cluster Representative by Average

— CPI values for Project ID-0000959 — CPI values for Project ID-0000961 — CPI values for Project ID-0000215
 — CPI values for Project ID-0000974 — Cluster Representative by Average

Historical Data
(Training Set)



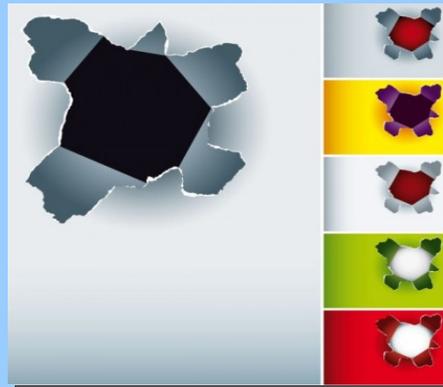


Historical Data
(Training Set)

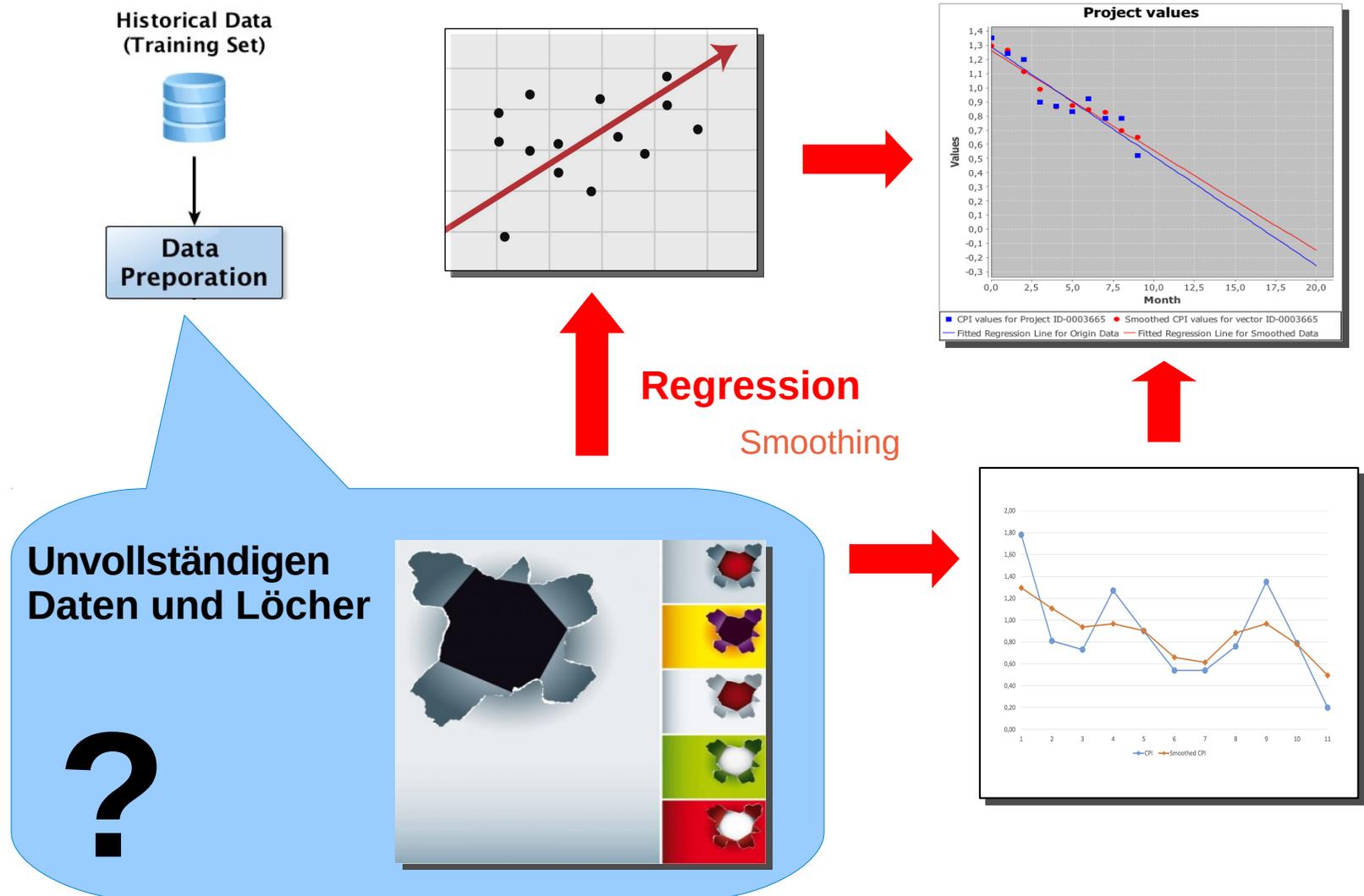


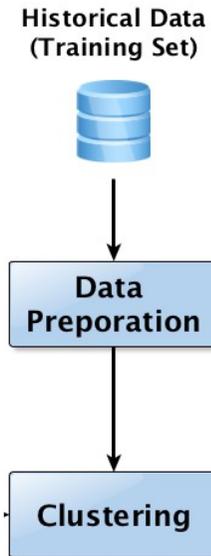
Data
Preparation

Unvollständigen
Daten und Löcher

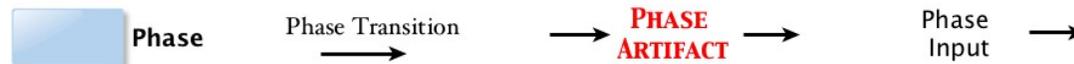


Idee und Konzept

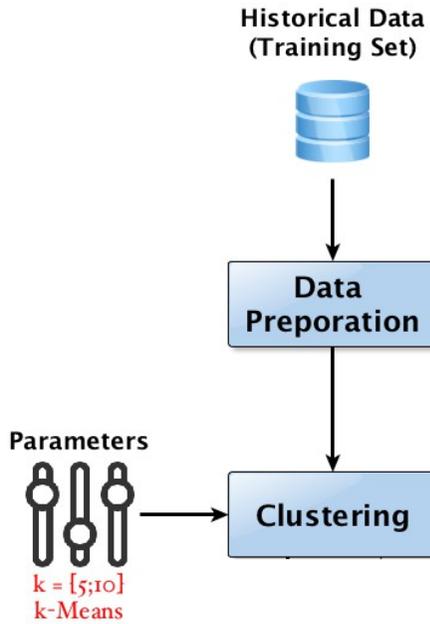




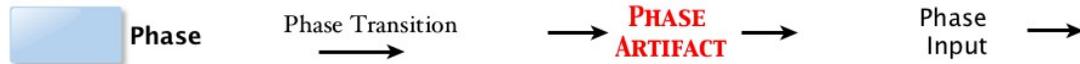
Legend



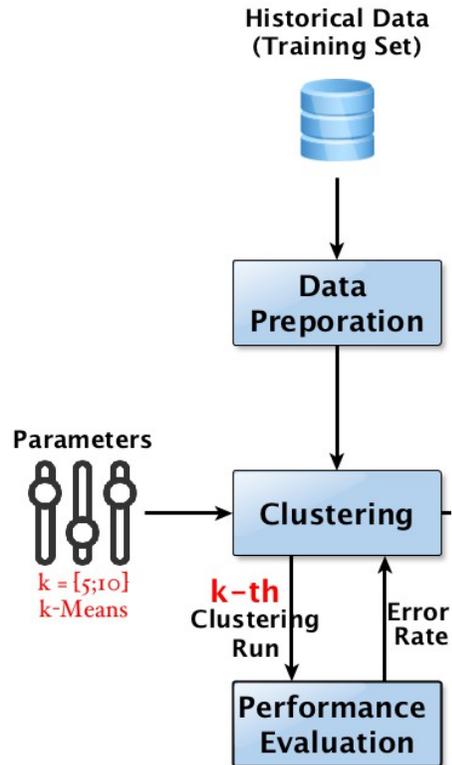
Idee und Konzept



Legend



Idee und Konzept



Legend



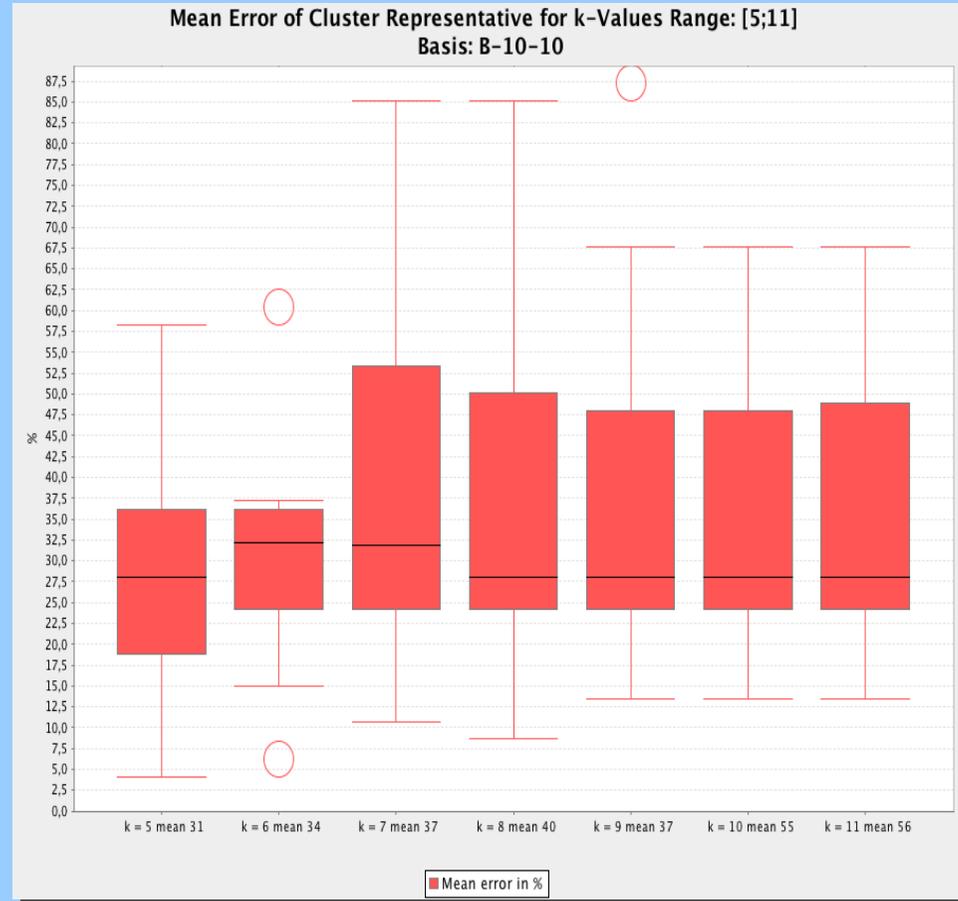
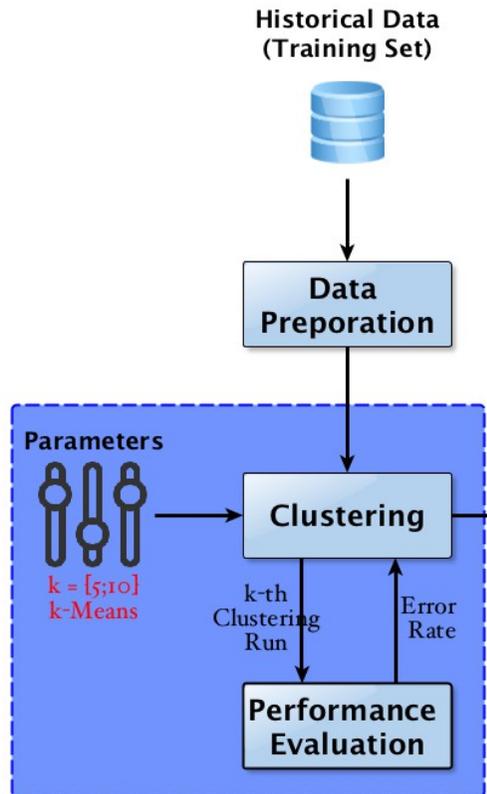
Phase

Phase Transition
→

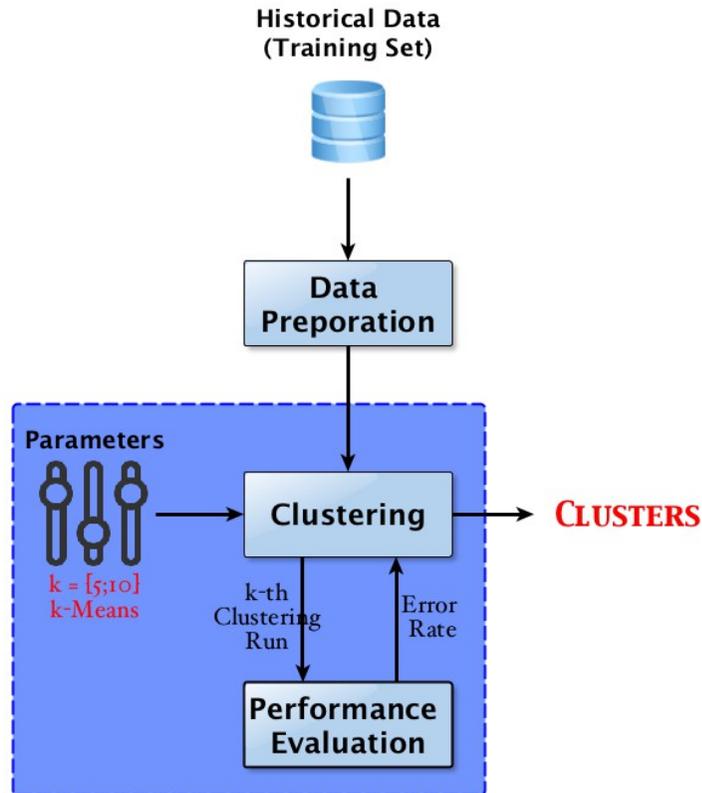
→ **PHASE ARTIFACT** →

Phase Input
→

Idee und Konzept



Idee und Konzept



Legend



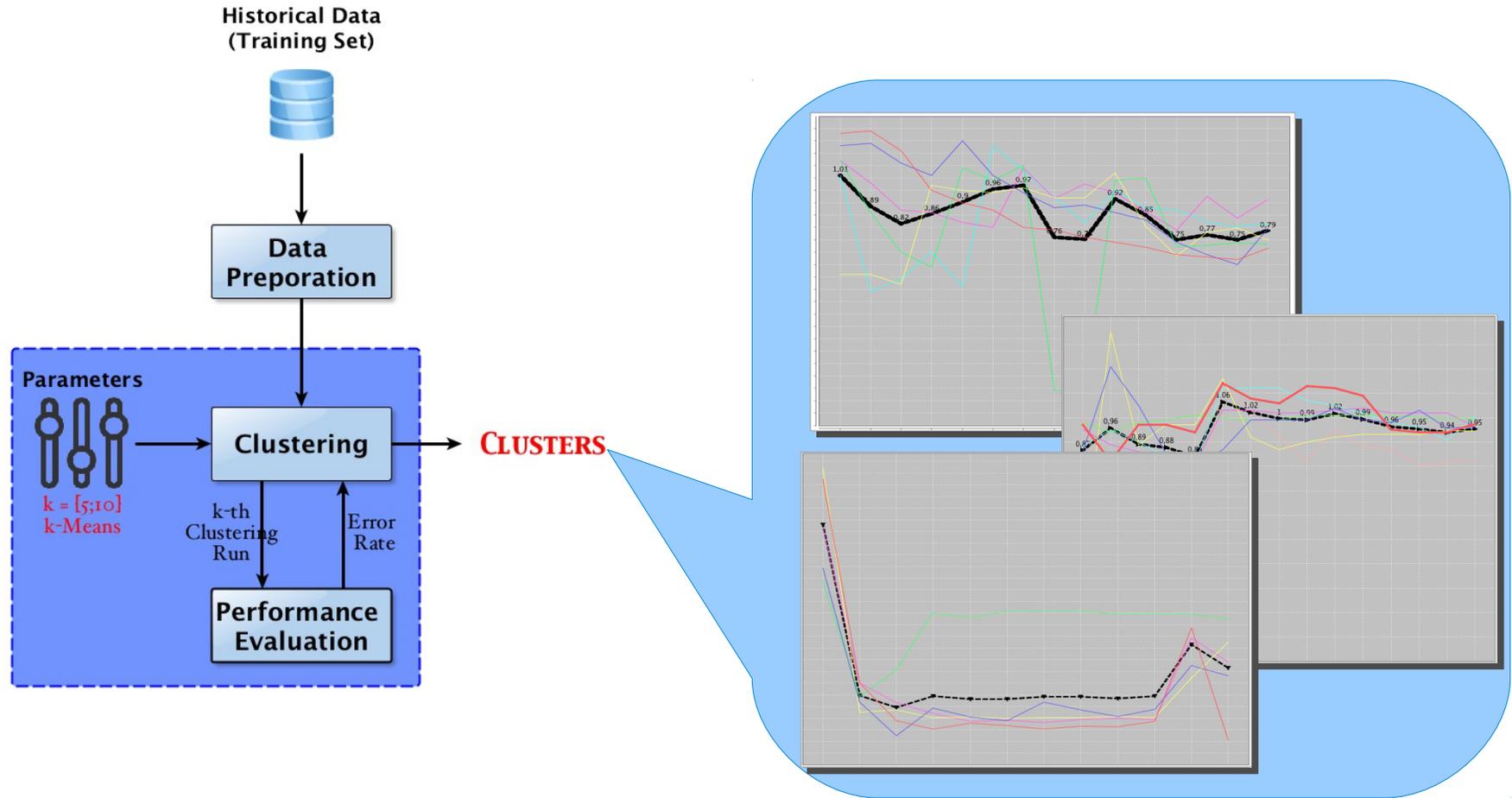
Phase

Phase Transition
→

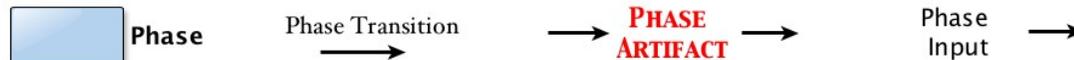
→ **PHASE ARTIFACT** →

Phase Input →

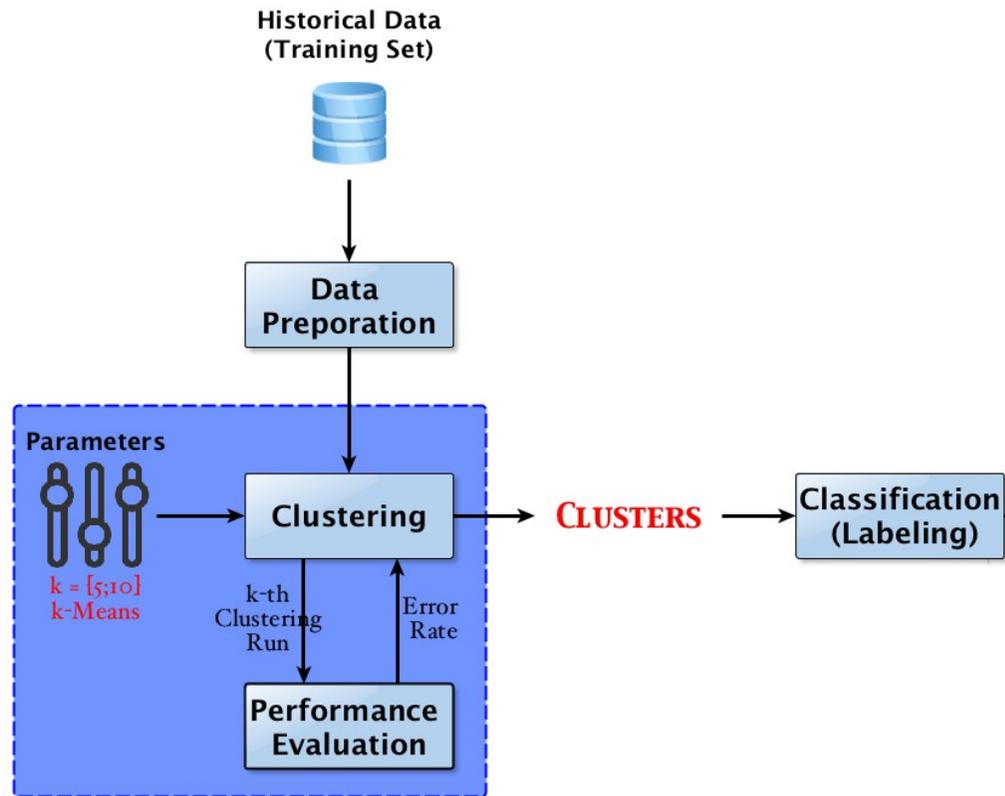
Idee und Konzept



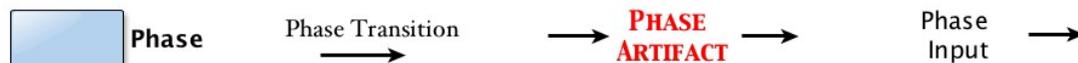
Legend



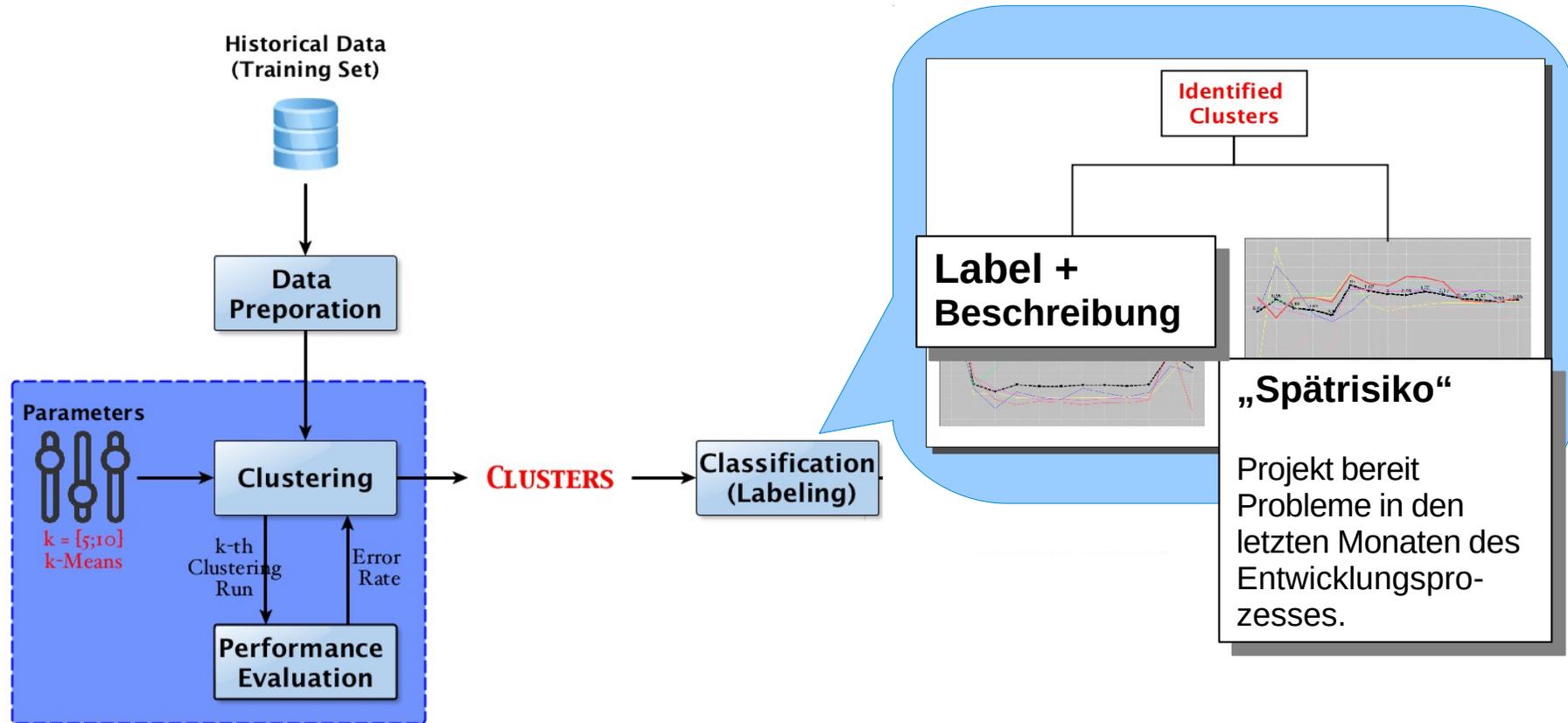
Idee und Konzept



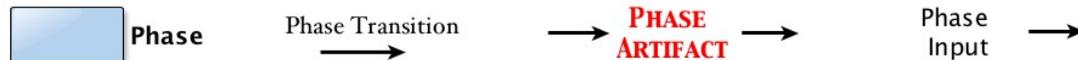
Legend



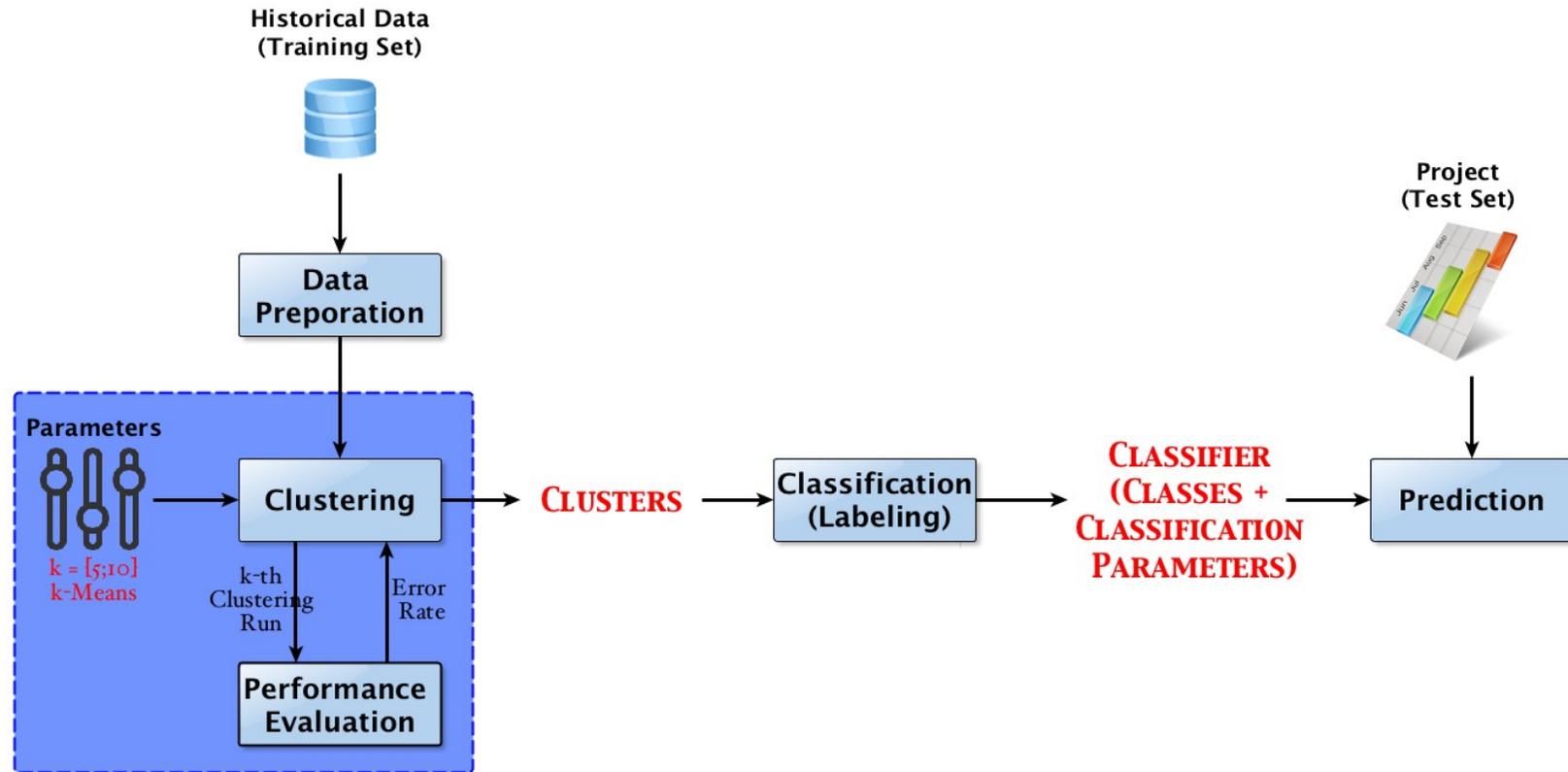
Idee und Konzept



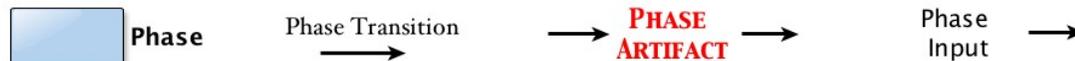
Legend



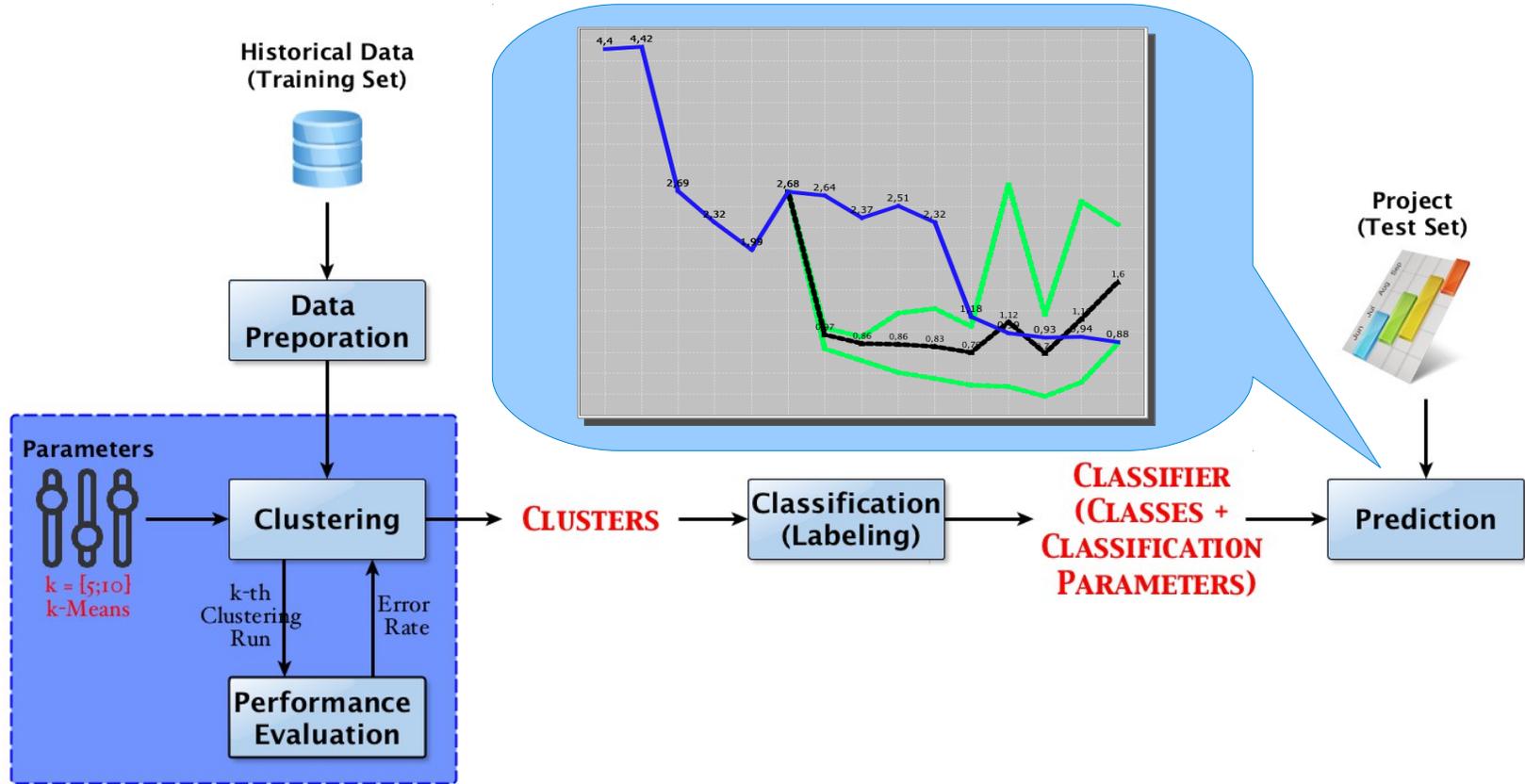
Idee und Konzept



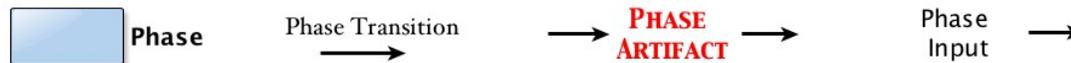
Legend



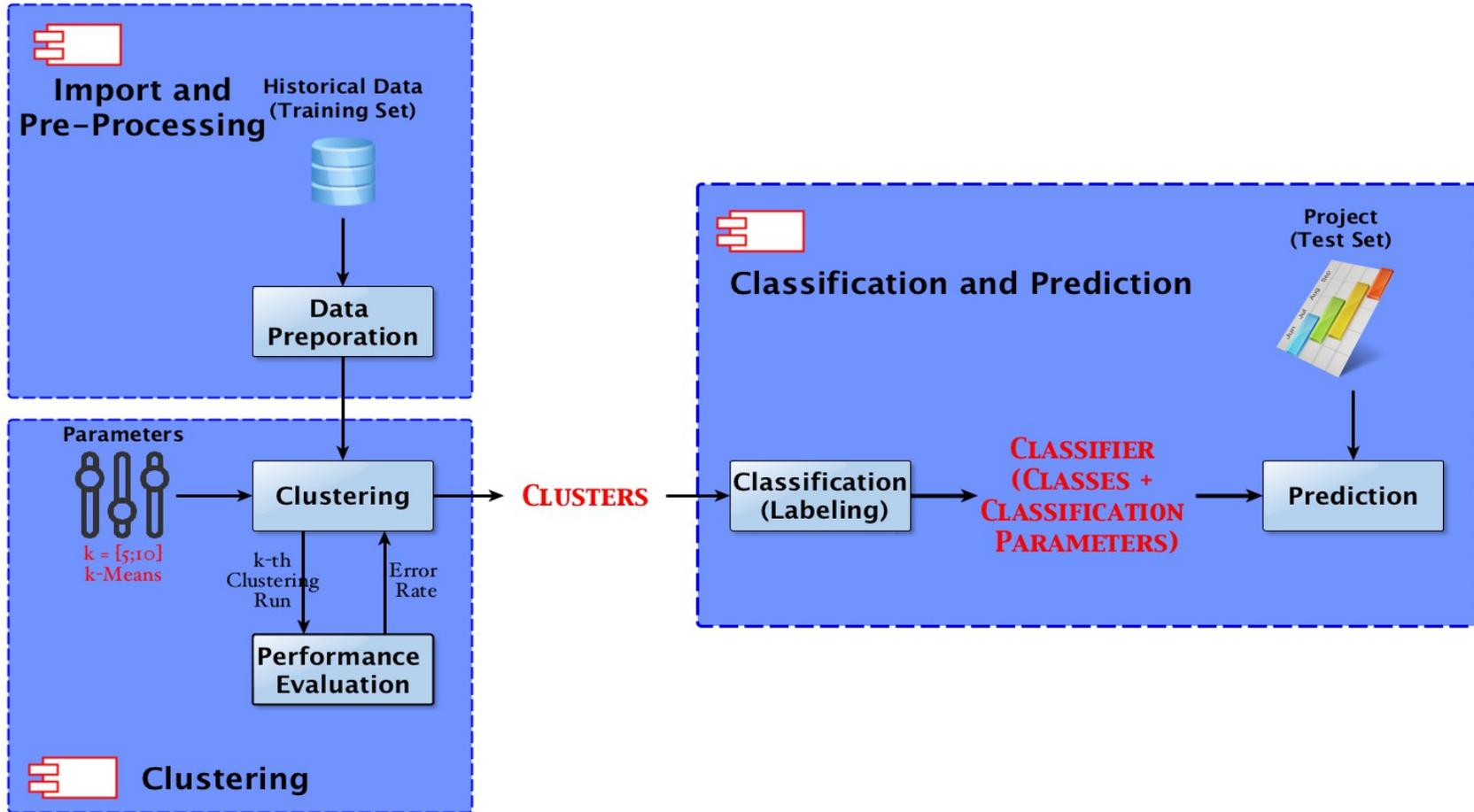
Idee und Konzept



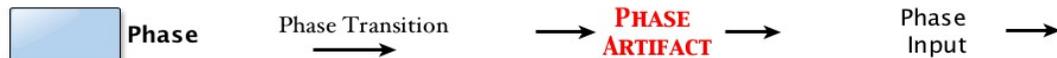
Legend



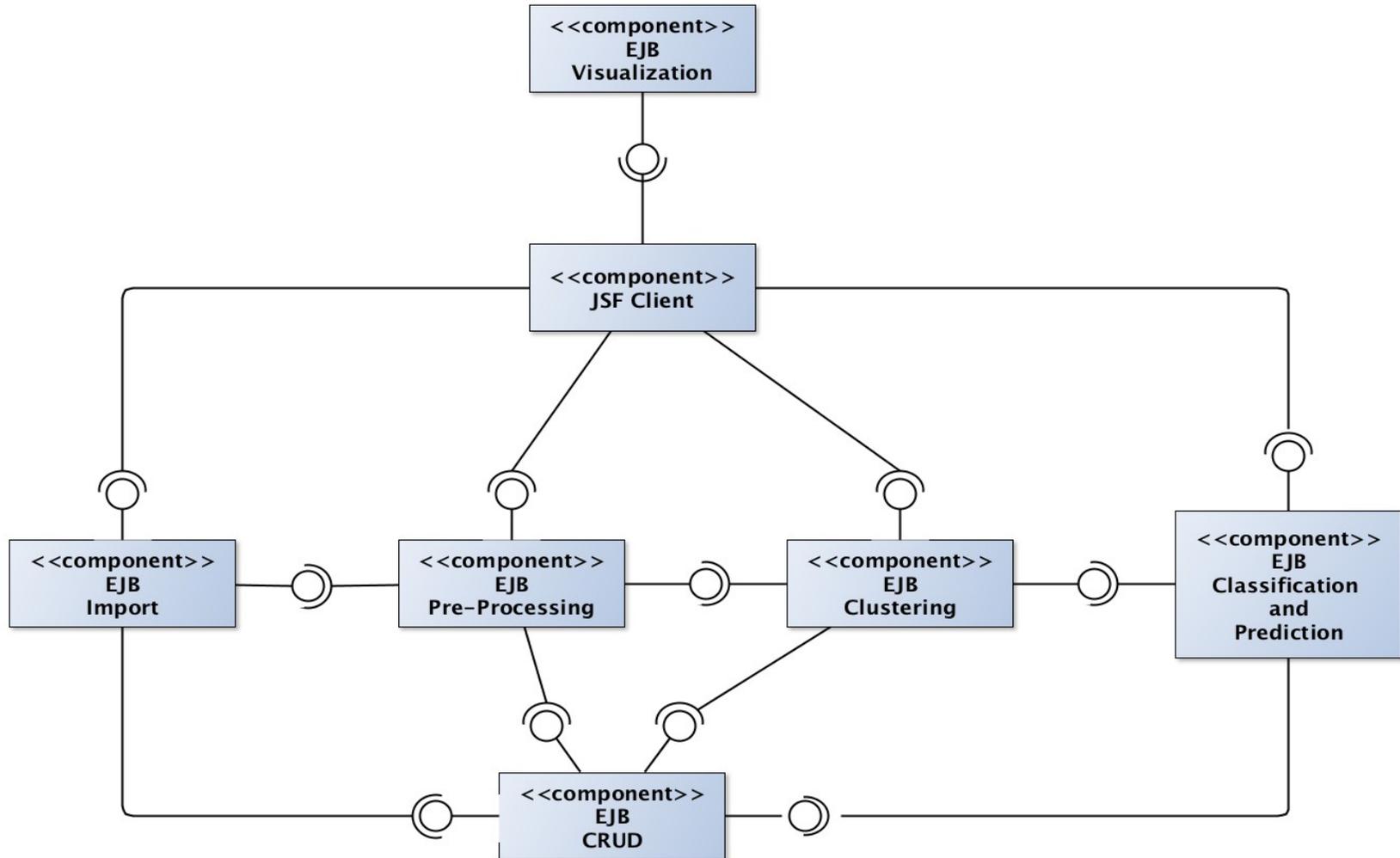
Idee und Konzept



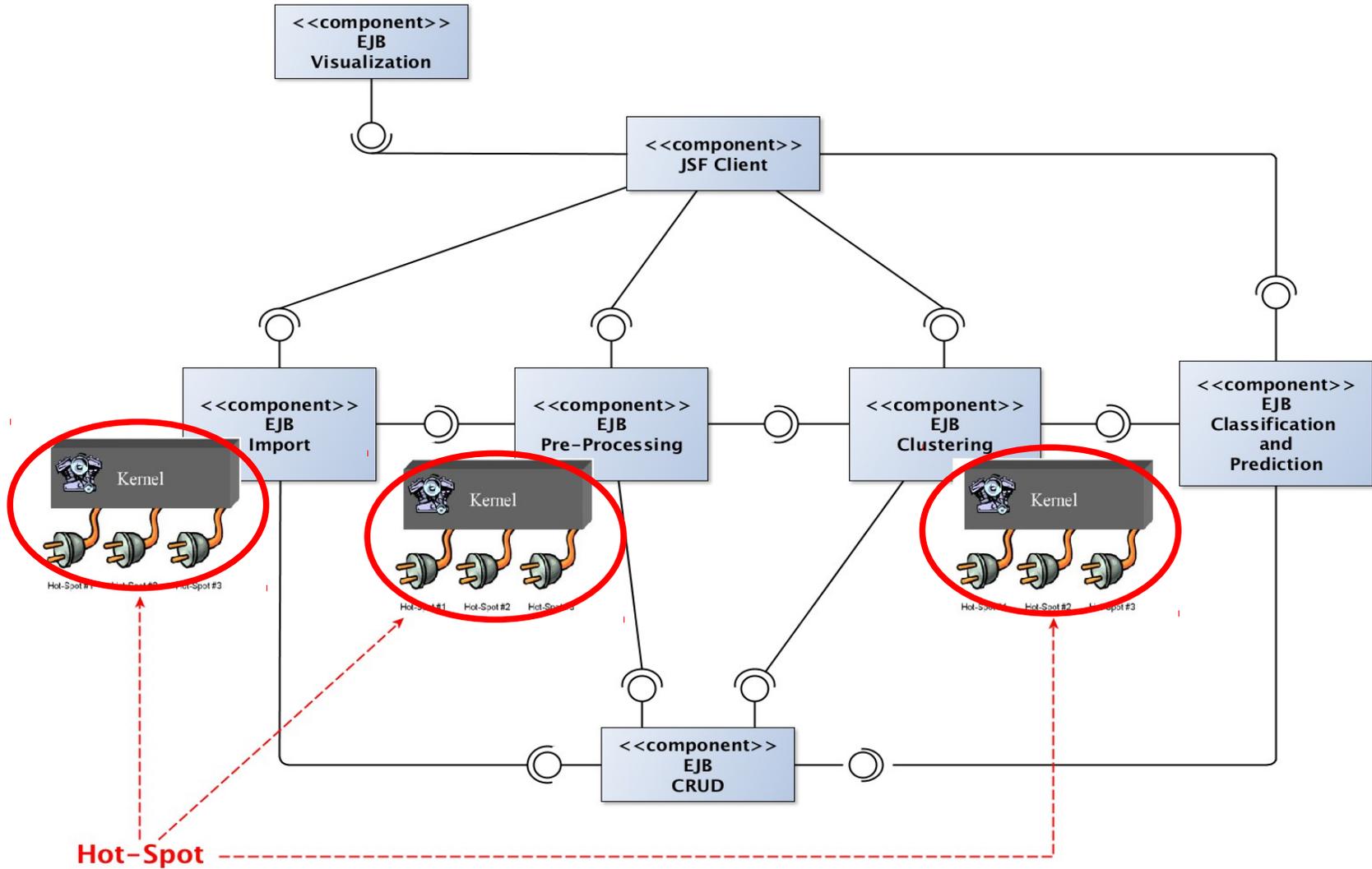
Legend

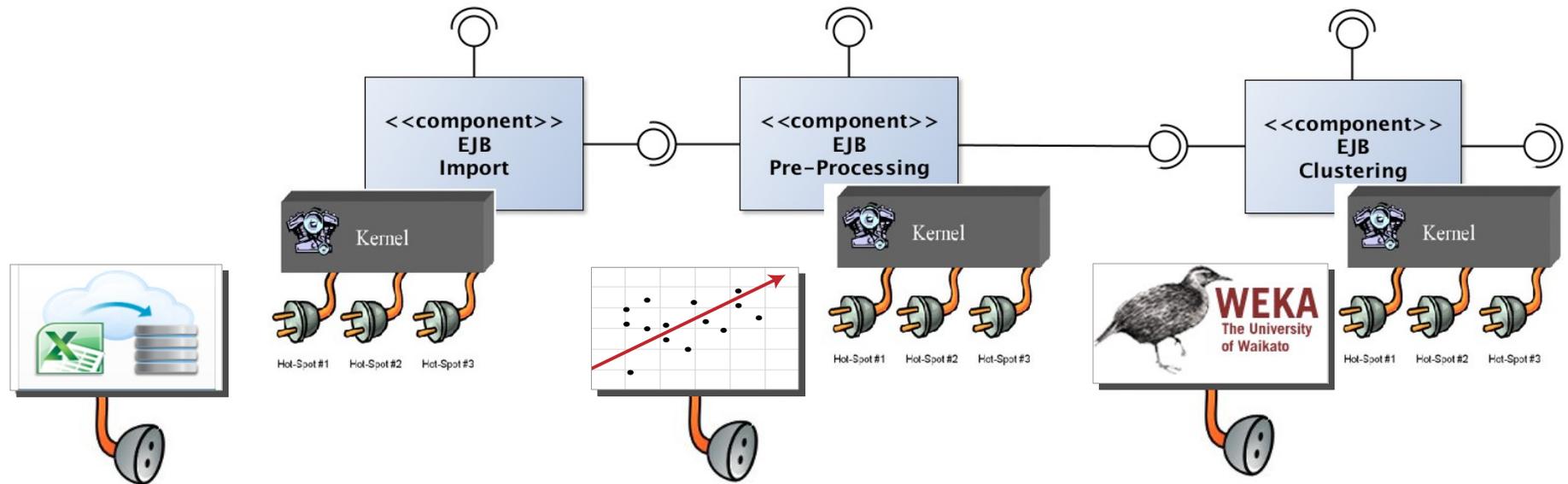


- Entwicklung von **Framework**
 - Lose gekoppelte Architektur
 - Vollständige Implementierung von Data Mining Prozesses (Wissensextraktion)
 - Jeder Schritt der Wissensextraktion ist gekapselt in eigene Komponente
 - Unterstützung von Variabilität und Erweiterbarkeit
 - Web-Basierte Anwendung + JavaEE
- Prototyp Implementierung von **Projekt-Vorhersage-Tool**
 - System “versteht” Daten von externem Kooperationspartner
 - Import, Pre-processing, Clustering
 - Benutzbarkeit und UI
 - Konfiguration ohne Expertenwissen



Hot-Spots



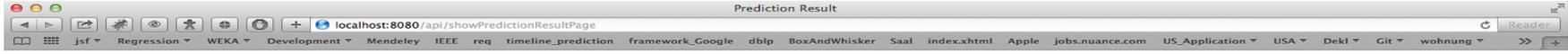


Hot-Spot Implementierung:
Import aus xls/xlsx

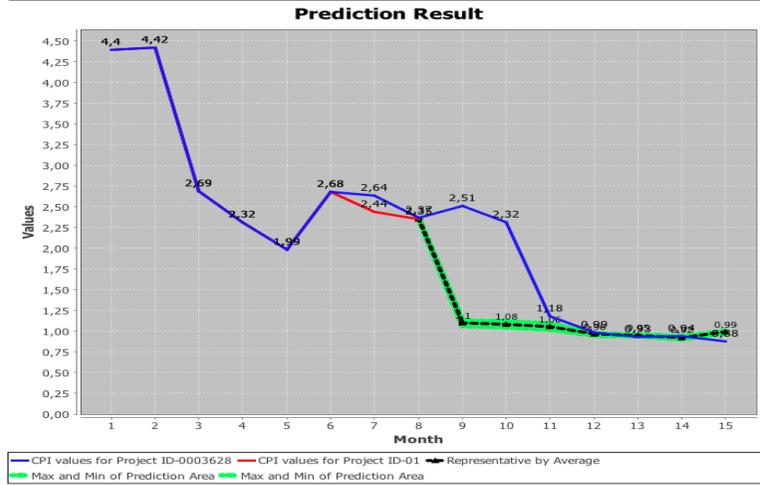
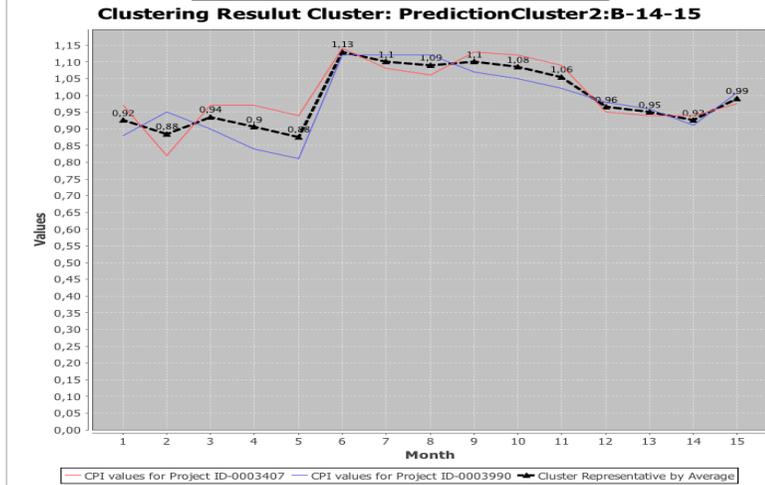
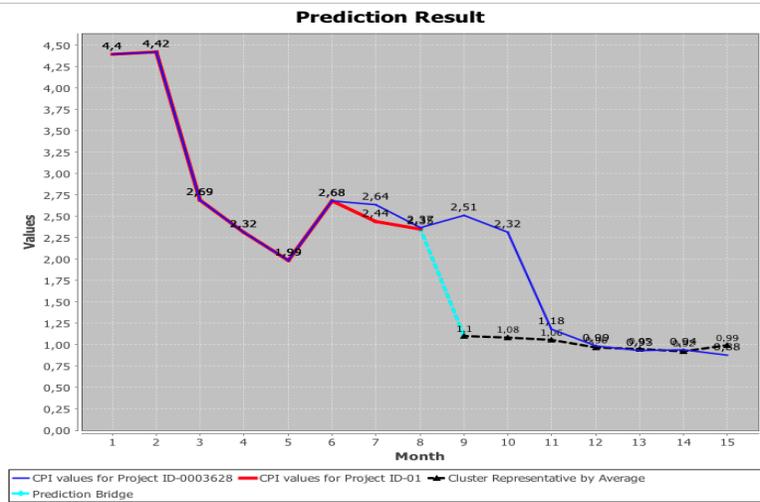
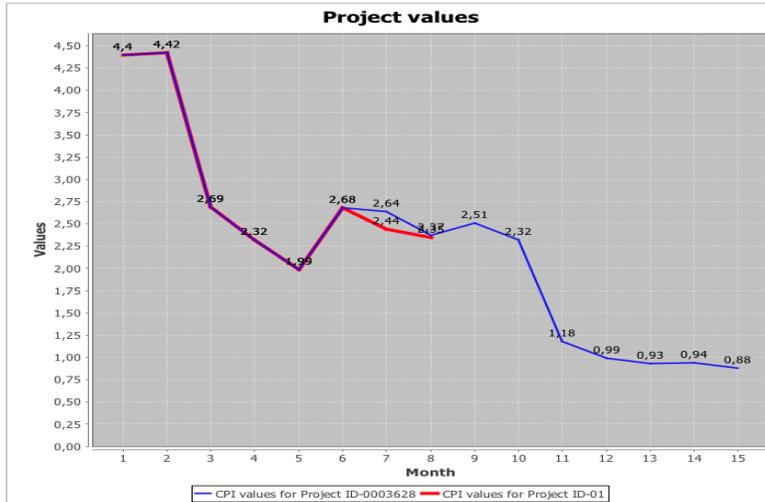
Hot-Spot Implementierung:
Regression Strategies
- Lineare Regression
- Nicht-Lineare Regression
Missing Values Handling

Hot-Spot Implementierung:
Clustering Verfahren
- k-Means
- Hierarchical
Cluster Repräsentant
- Mittelwert
- Max/Min

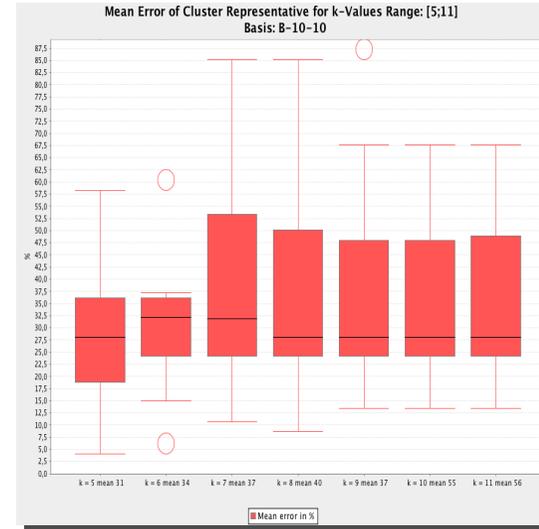
Realisierung: Tool Prototyp



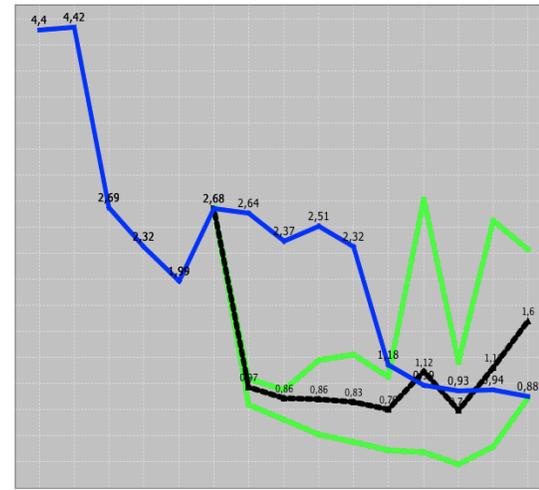
Prediction Result:



- **Szenario 1:**
Bewertung von Cluster-Analyse



- **Szenario 2:**
Bewertung von Vorhersage
(Prediction)

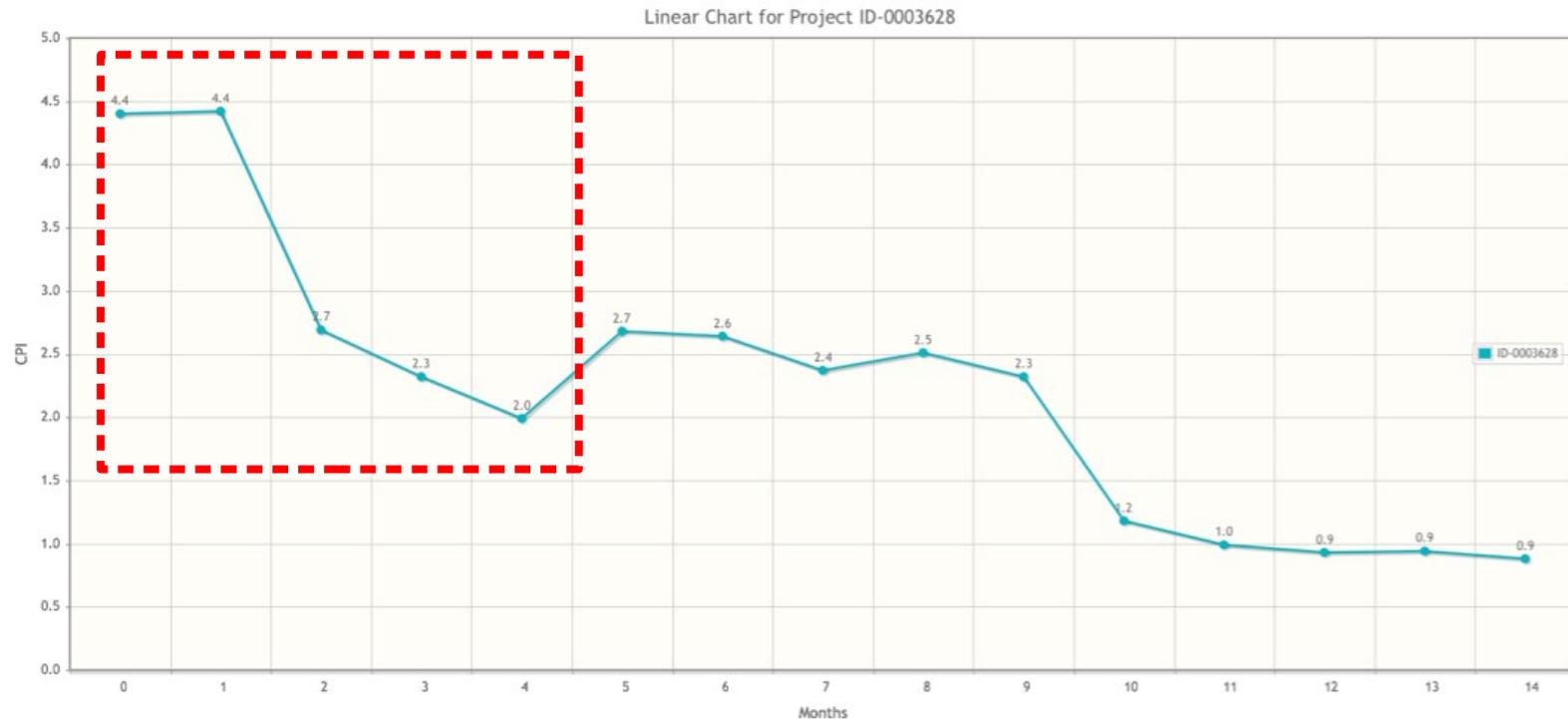


- **Betriebssystem:** Mac OS X, Version 10.9.5, Prozessor 2 GHz
- **Daten:** CPI Werte von 151 Projekten von 5 bis 15 Monaten ($\approx 54\%$)
- **Parameter**
 - Lineare Regression
 - k-Means
 - Cluster-Zahl k aus dem Bereich [5;10]

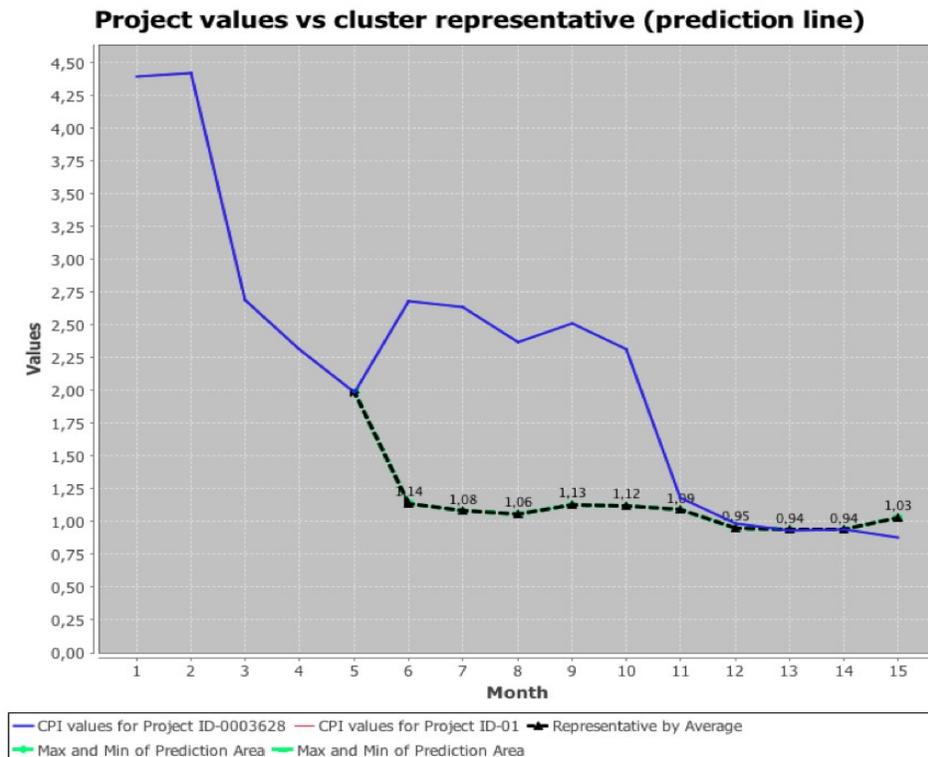
- 10 Datensätze getestet
- **Das beste Ergebnis**
 - Datensatz aus 14 bis 15-monatigen Projekten (12 Proj.)
 - 95% aller Fehler zwischen 3% und 76%
 - Mittelwert der Fehler zwischen 12% und 22%
 - Bestes $k = 8$
- **Das schlechteste Ergebnis**
 - Datensatz aus 8 bis 15-monatigen Projekten (99 Proj.)
 - 95% aller Fehler zwischen 12% und 322%
 - Mittelwert der Fehler zwischen 94% und 135%
 - Bestes $k = 7$

- Regression als Einflussfaktor auf Qualität des Clusterings
- Für jeden Datensatz existiert einen bestes Cluster-Zahl k
- Sehr heterogene Daten haben negative Auswirkung auf Clustering-Qualität

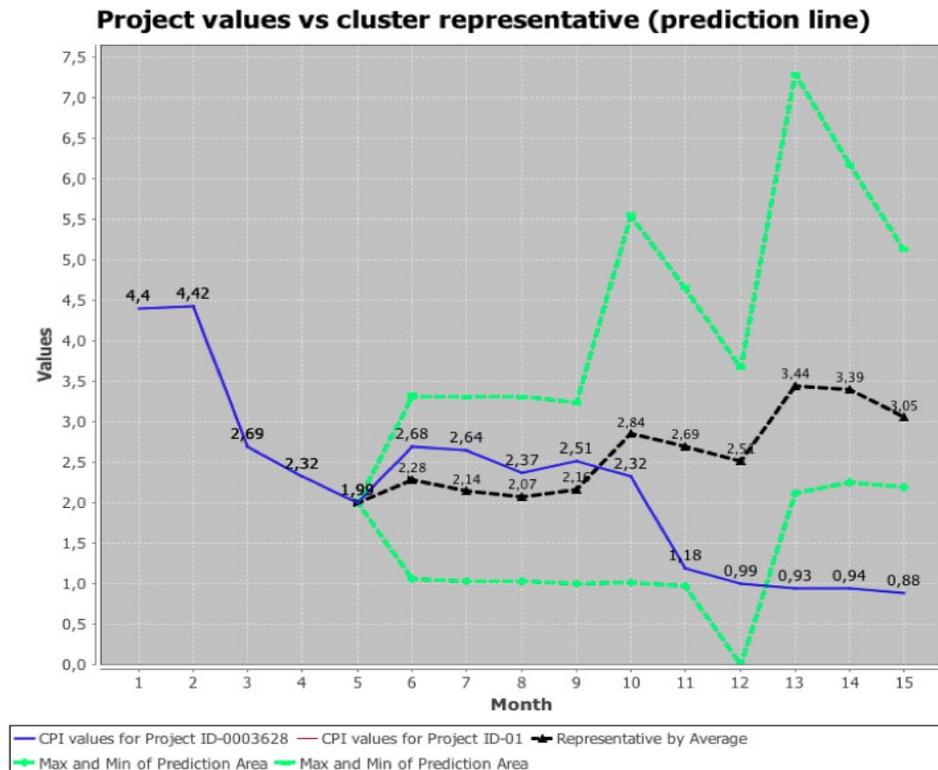
- 2 Experimente
 - Einfache Vorhersage von einem 15-montigen Projekt
 - Schrittweise Vorhersage



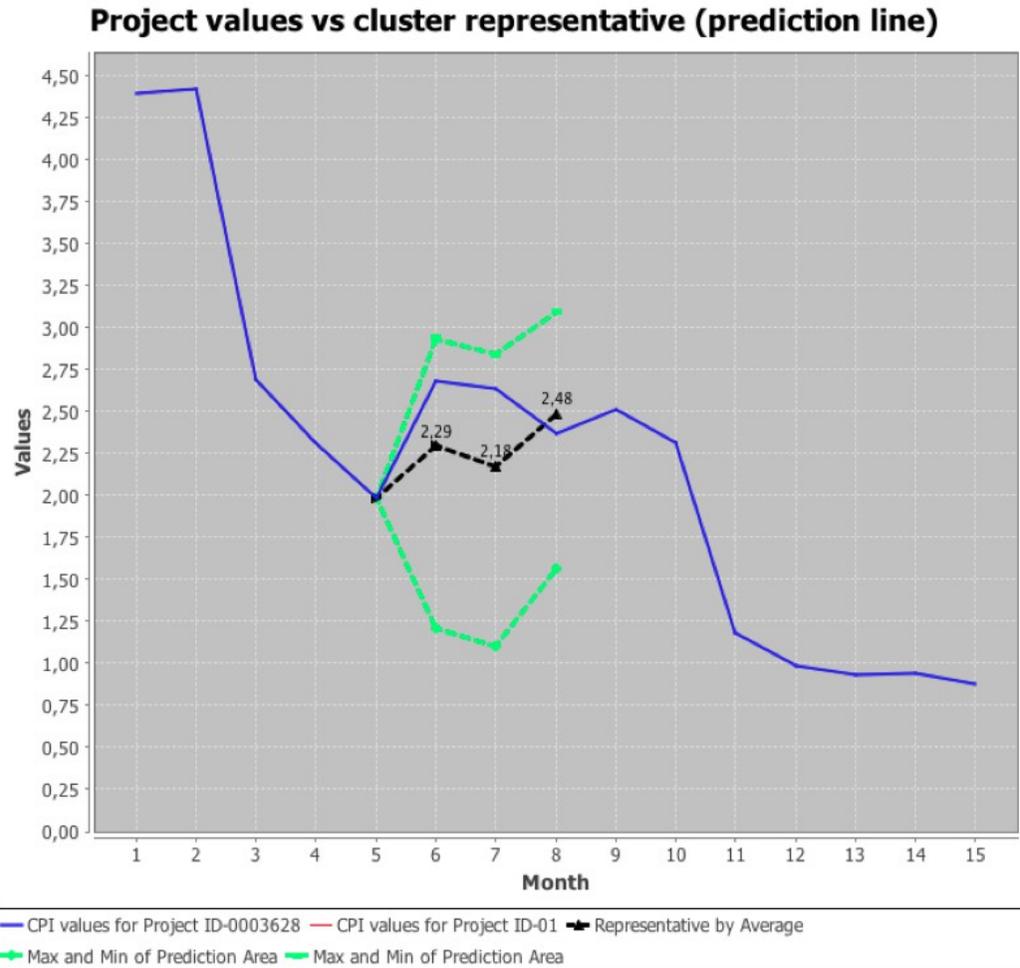
- Einfache Vorhersage
 - Datensatz: 14 bis 15-monatige Projekte
 - Mittelwert der Fehler: **31,83%**



- Einfache Vorhersage
 - Datensatz: 8 bis 15-monatige Projekte
 - Mittelwert der Fehler: 114,12%

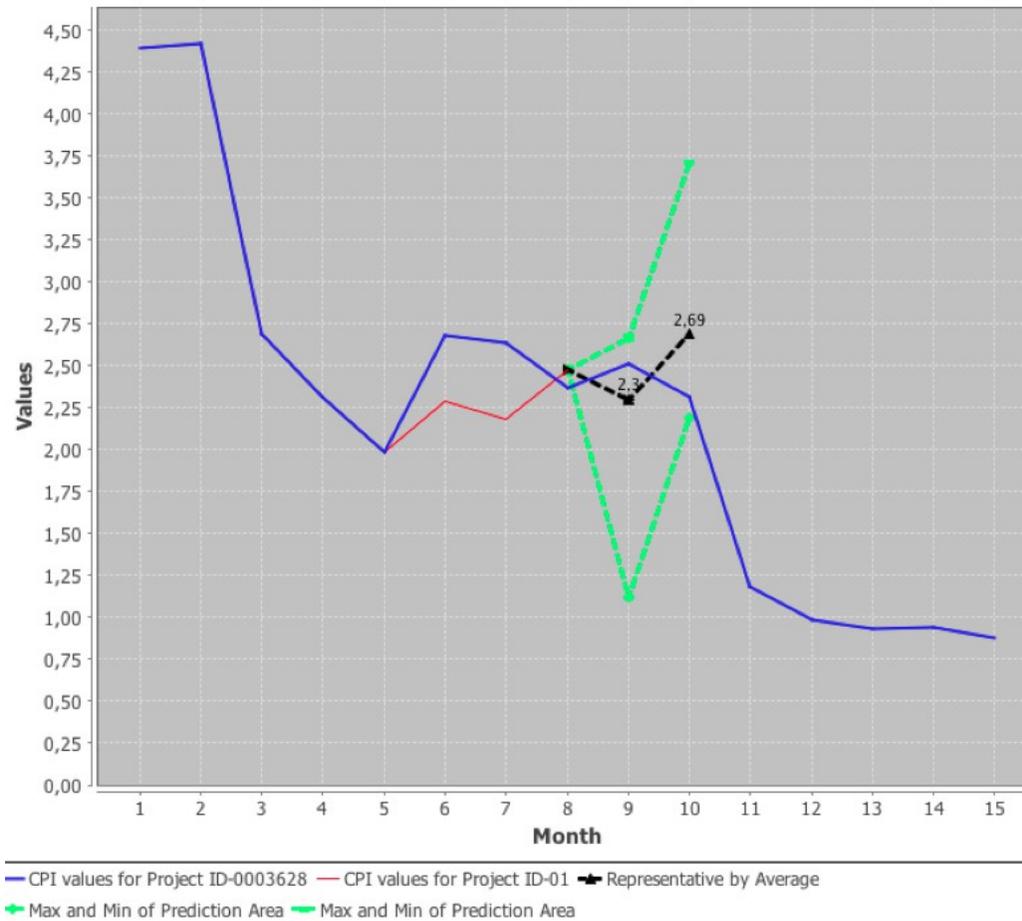


- Schrittweise Vorhersage

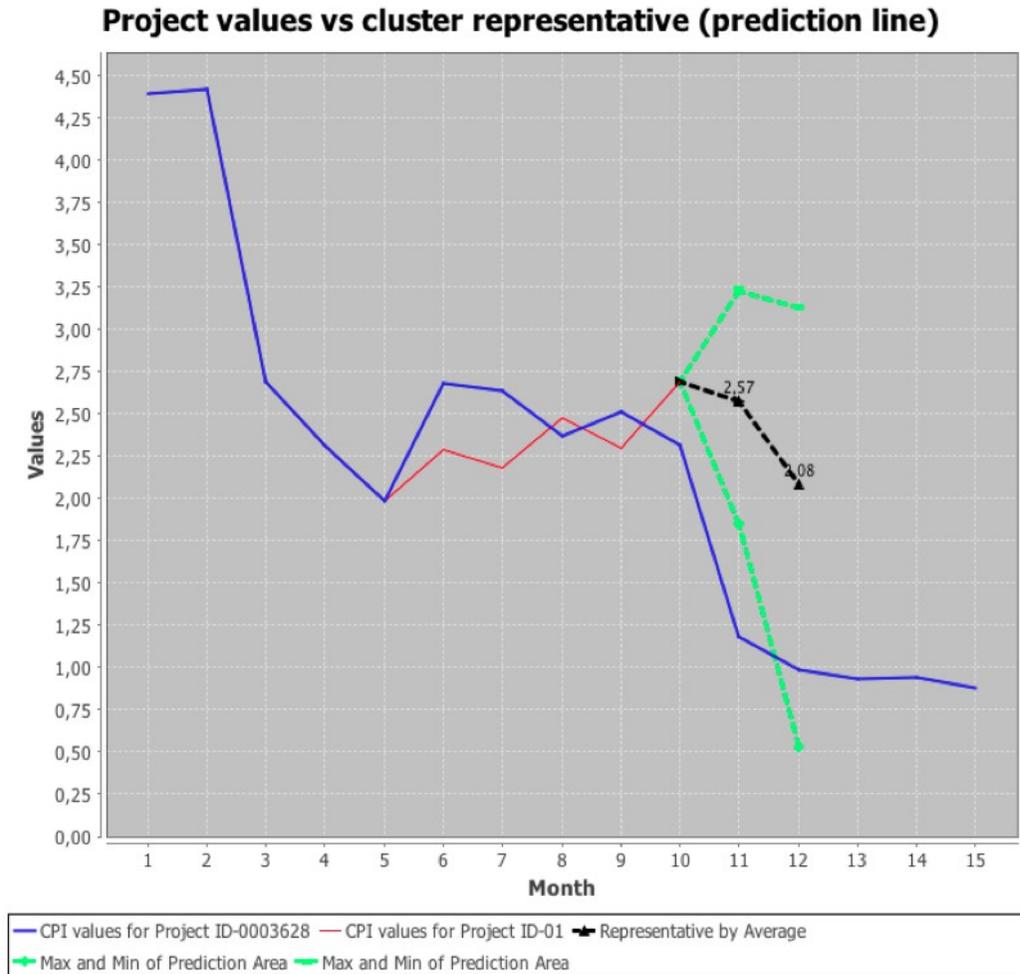


- Schrittweise Vorhersage

Project values vs cluster representative (prediction line)

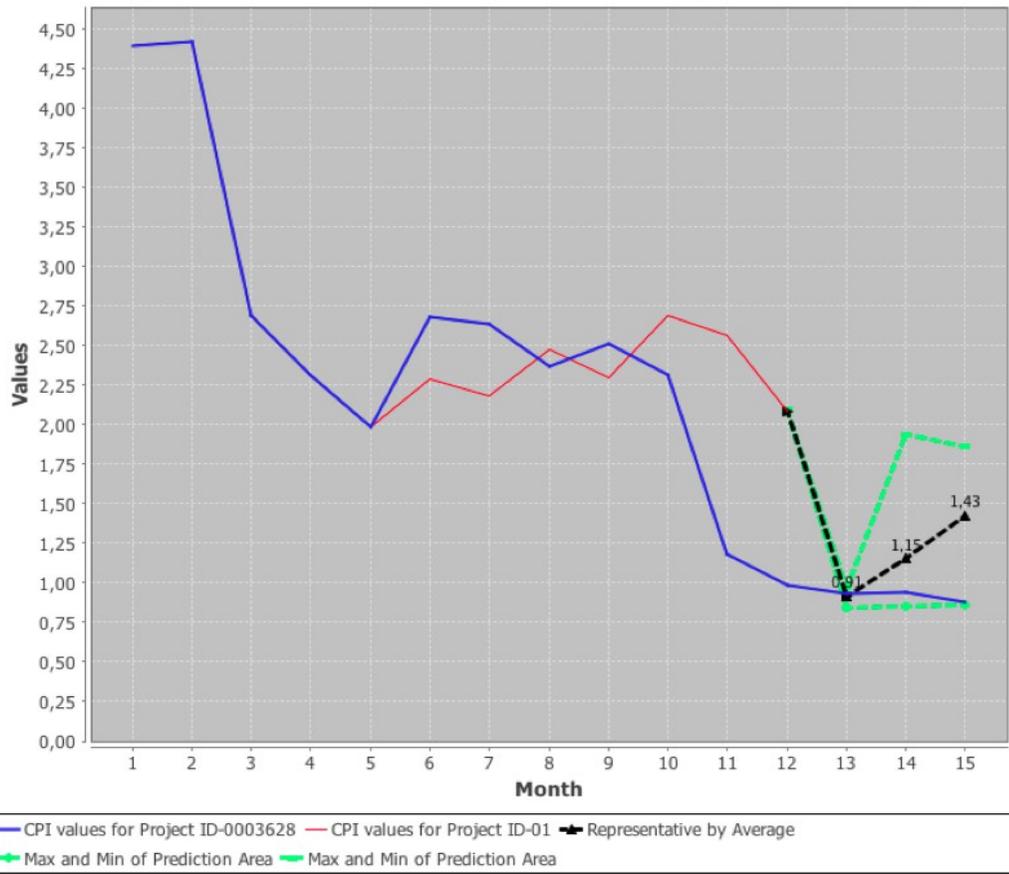


- Schrittweise Vorhersage

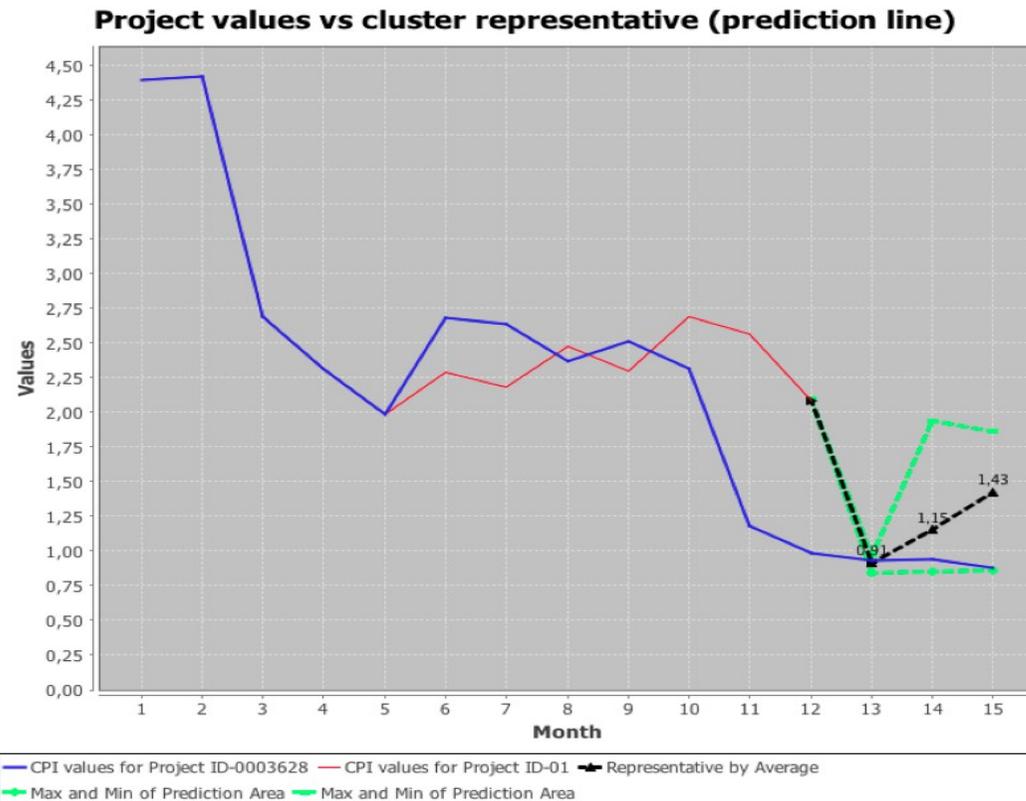


- Schrittweise Vorhersage

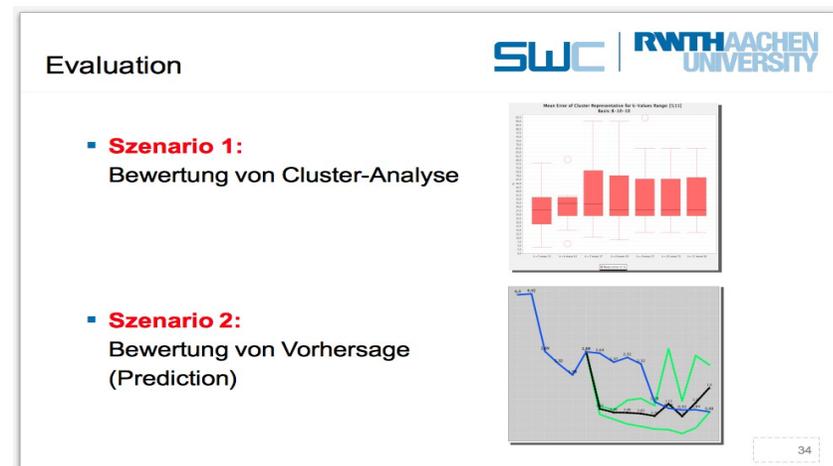
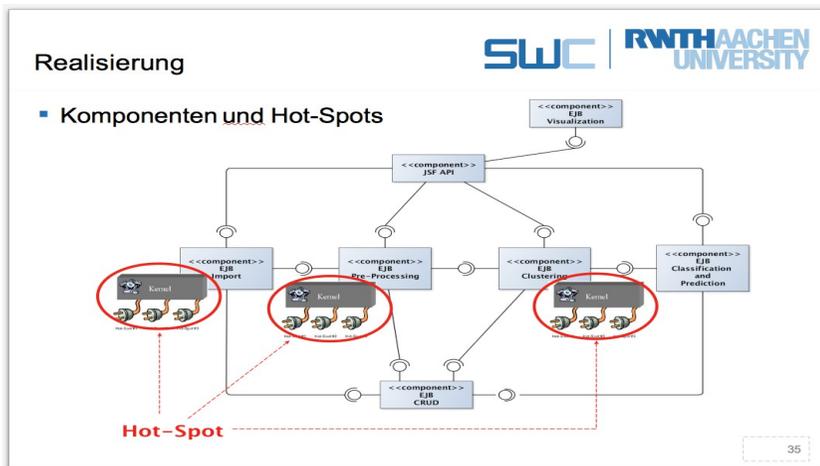
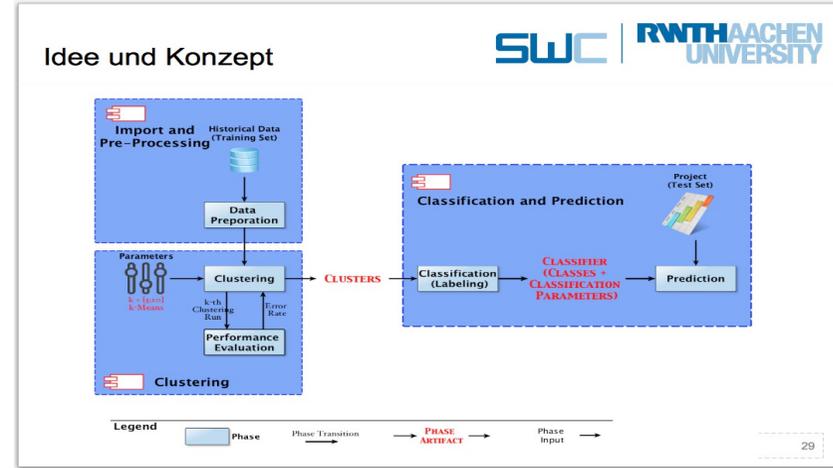
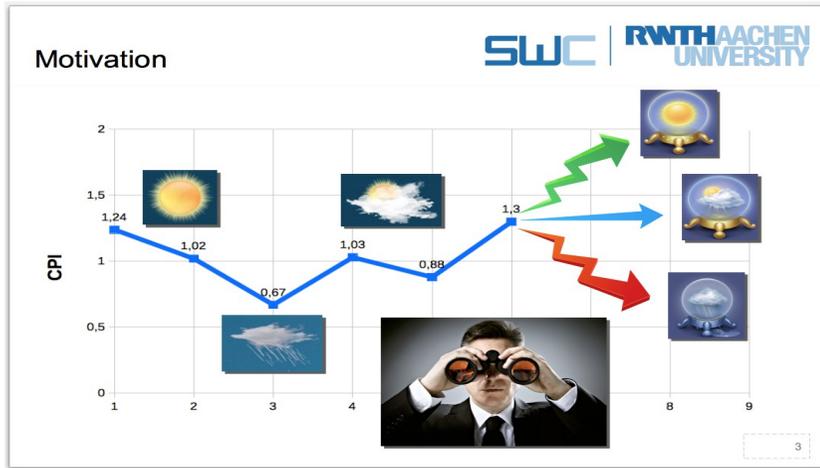
Project values vs cluster representative (prediction line)



- Schrittweise Vorhersage
 - Mittelwert der Fehler: 37,59%



- Weitere Implementierung von Hot-Spots
- Identifizierung von weiteren Einflussfaktoren
- Anwendung und Evaluation in anderen Domänen
- Prototyp-Erweiterung (Client mit AngularJS)
- Integration ins EMI System



- [1] Scafetta Prediction: Scafetta, EPA 2009. – URL <http://vademecum.brandenberger.eu/themen/klima/ursache.php>. – Access date: 10.11.2014
- [2] Thanh Vi Bach: Entwurf prognostischer Softwareprozessmetriken auf Basis iterativen Clusterings, RWTH Aachen University, Master Thesis, 2014.
- [3] Dindin Wahyudin, Rudolf Ramler, and Stefan Biffel. A framework for defect prediction in specific software project contexts. In Zbigniew Huzar, Radek Kocí, Bertrand Meyer, Bartosz Walter, and Jaroslav Zendulka, editors, CEE-SET, volume 4980 of Lecture Notes in Computer Science, pages 261–274. Springer, 2008.
- [4] Ningda R. Li and Marvin V. Zelkowitz. An information model for use in software management estimation and prediction. In Bharat K. Bhargava, Timothy W. Finin, and Yelena Yesha, editors, CIKM 93, Proceedings of the Second International Conference on Information and Knowledge Management, Washington, DC, USA, November 1-5, 1993, pages 481–489. ACM, 1993.
- [5] Rudolf Ramler, Klaus Wolfmaier, Erwin Stauder, Felix Kossak, and Thomas Natschläger. Key questions in building defect prediction models in practice. In Frank Bomarius, Markku Oivo, Päivi Jaring, and Pekka Abrahamsson, editors, PROFES, volume 32 of Lecture Notes in Business Information Processing, pages 14–27. Springer, 2009.