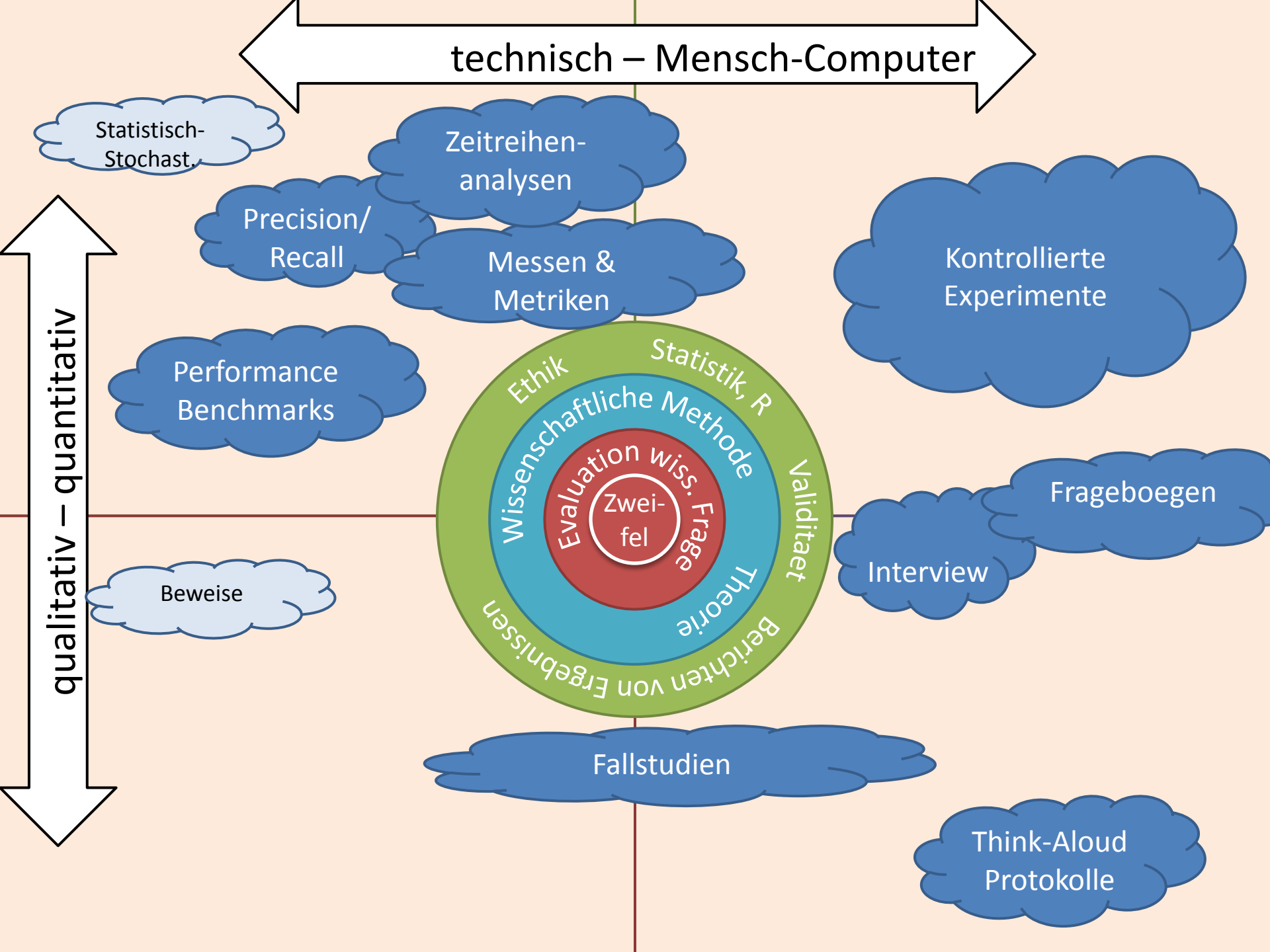


Empirische Methoden für Informatiker

Teil 5: Quantitative Untersuchungen

Christian Kästner



Agenda

- ▶ Weitere Quantitative Untersuchungen
 - ▶ ausser Performance-Messungen (Teil 2)
 - ▶ ausser kontrollierten Experimenten (Teil 6)
- ▶ Metriken
- ▶ Datengewinnung aus Repositories
- ▶ Zeitlichreihen Analyse
- ▶ Frageboegen
- ▶ Fallstudien (erneut)



Messen

Definition: Messen

“Zuweisen von Zahlen zu
Objekten oder Ereignissen
anhand einer eindeutigen Regel”

Skalen

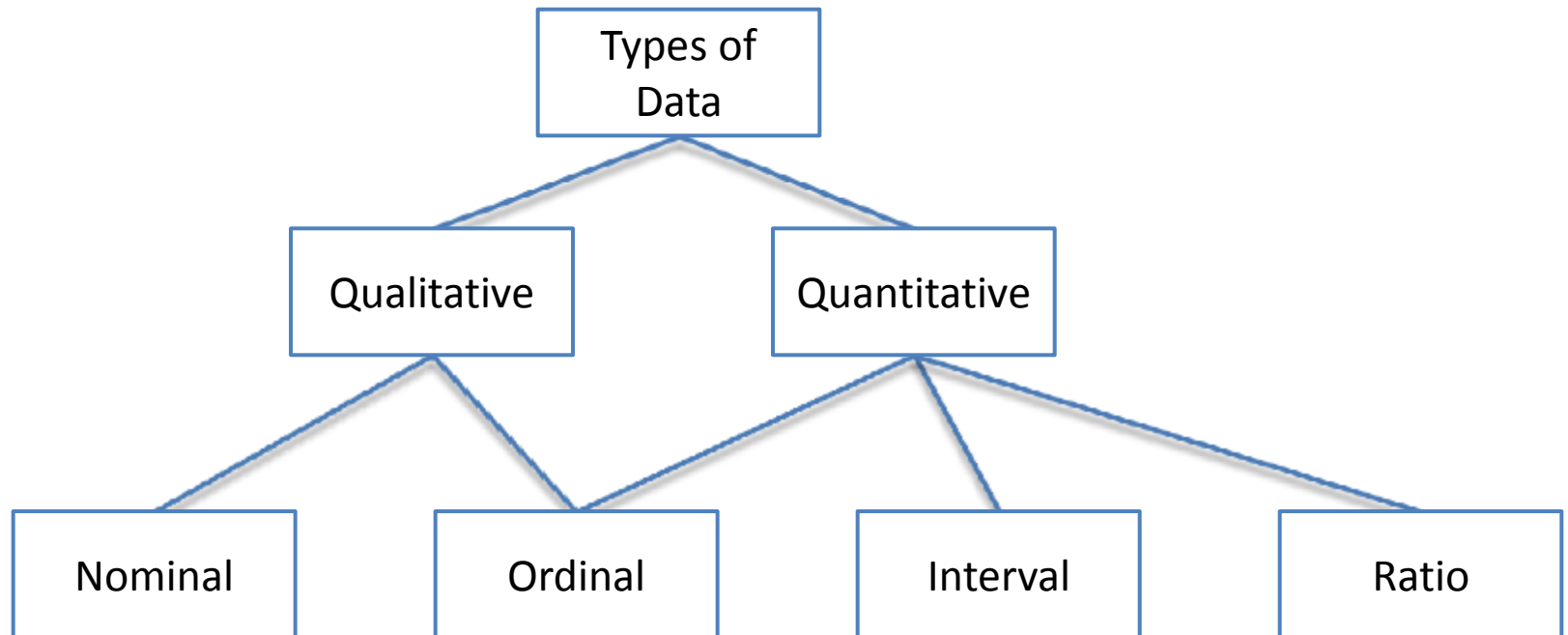
- ▶ Nominalskala
 - ▶ Ordinalskala
 - ▶ Intervallskala
 - ▶ Verhaeltnissskala
- } metrisch

Skalen

- ▶ Nominalskala (z.B. Geschlecht, Studiengang)
- ▶ Ordinalskala (z.B. Gut-Mittel-Schlecht)
- ▶ Intervallskala (z.B. Temperatur)
- ▶ Verhaeltnissskala (z.B. Zeit, Durchsatz, Fehlerrate)

Welche statistischen Verfahren
lassen sich (nicht) anwenden?

Skalen



Messen in der Informatik

- ▶ Performance, Speicherverbrauch, ...
- ▶ Genauigkeit von Vorhersagen, ...
- ▶ Lesbarkeit, Wartbarkeit, ...
- ▶ Stabilität, Fehleranzahl, ...
- ▶ Nützlichkeit, Zufriedenheit, ...

Datenerhebung, Masseinheiten und Messmethoden

- ▶ Stark abhaengig vom Evaluierungsziel
- ▶ Im Folgenden einige typische Beispiele
- ▶ Unvollstaendig, kein fester Katalog
- ▶ Kreativitaet

Kriterien fuer gute Metriken

10TH INTERNATIONAL SOFTWARE METRICS SYMPOSIUM, METRICS 2004

KANER / BOND - 1

Software Engineering Metrics: What Do They Measure and How Do We Know?

Cem Kaner, *Senior Member, IEEE*, and Walter P. Bond

Abstract-

This is dis
rarely in s
attribute w
an attribut
are so sim
evaluating
attributes
quality of
interest in

Index Ter
V&V.

1 INTRODU

W e hear
urement
that many

Folgende Fragen beantworten:

- Welchen Zweck hat die Metrik?
- Welchen Scope hat die Metrik?
- Welches Attribut wird gemessen?
- Was ist die natuerliche Skala des Attributs?
- Was ist die natuerliche Schwankung des Attributs?
- Wie wird gemessen?
- Was ist die natuerliche Skala der Metrik?
- Was ist die natuerliche Schwankung des Messverfahrens? (Messfehler)
- Wie gut passt die Metrik zum Attribut? (Construct validity)
- Welche Seiteneffekte wird es durch die Messung/Evaluierung geben?

Fragebogen

Frageboegen

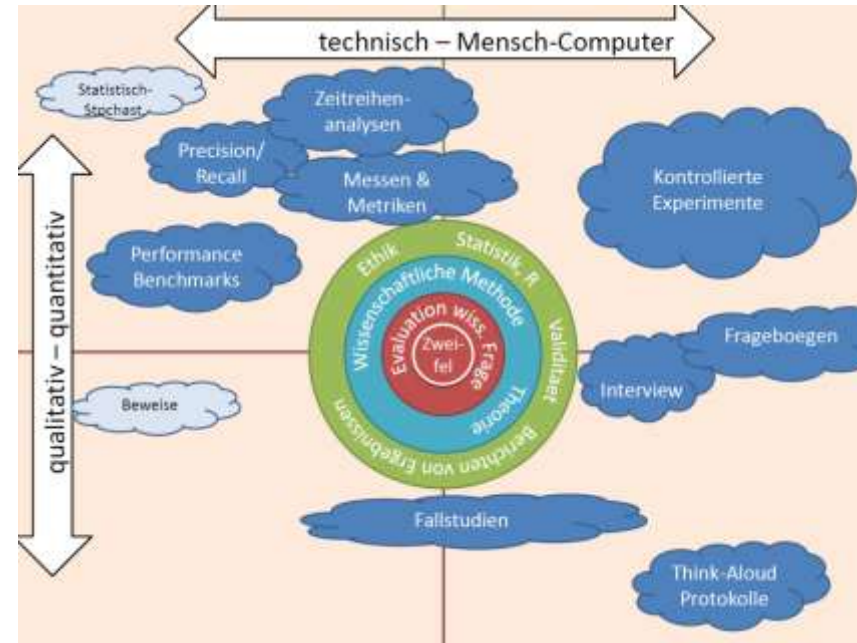
- ▶ Geschlossene Fragen quantitativ Auswerten
- ▶ Likert-Skala 1-5
- ▶ auch Online-Befragung, fertige Software verfuegbar
- ▶ Klar Hypothese erarbeiten
- ▶ Kritisch reflektieren ob Fragebogen geeignetes Mittel
- ▶ In Informatik oft benutzt, aber meist oberflaechlich
- ▶ Vor Beginn Literatur dazu lesen!

Umfangreiche Literatur zu Frageboegen

- ▶ Fragen vs. Behauptungen
 - ▶ “War der Quelltext besser lesbar?” vs. “Der Quelltext war besser lesbar. [Stimme (nicht) zu].”
- ▶ Eindeutige Formulierungen, negative Formulierungen
- ▶ Neutralitaet der Fragen
- ▶ Absolute Formulierungen (“immer besser lesbar”)
- ▶ Antwortmoeglichkeiten disjunkt und erschöpfend
- ▶ Jede Frage notwendig?
- ▶ Angemessen? In Verlegenheit bringen?
- ▶ Aufbau des Fragebogens
- ▶ Anerkennung, Incentives
- ▶ Ruecklaufcharakteristik, Ad-Hoc Stichproben

Fragebogen ergaenzend verwenden

- ▶ Ergaenzend zu Interviews oder kontrollierten Experimenten
- ▶ Beschreiben der Population (Erfahrung/Alter/Geschlecht/Motivation/...)
- ▶ Drittvariablen messen und kontrollieren



- ▶ Wenn moeglich etablierten Fragebogen verwenden

Art und Groesse der Stichprobe

- ▶ siehe Literatur (z.B. Bortz)

Vorhersagen / Suchergebnisse (Precision/Recall)

10.000 Eintraege, davon 1000 Relevant

	Relevant	Nicht Relevant
Gefunden	900	2000
Nicht gefunden	100	7000

Diskussion:
Welcher Algorithmus bietet
das beste Ergebnis?

	Relevant	Nicht Relevant
Gefunden	500	3
Nicht gefunden	500	8997

	Relevant	Nicht Relevant
Gefunden	1000	8500
Nicht gefunden	0	500



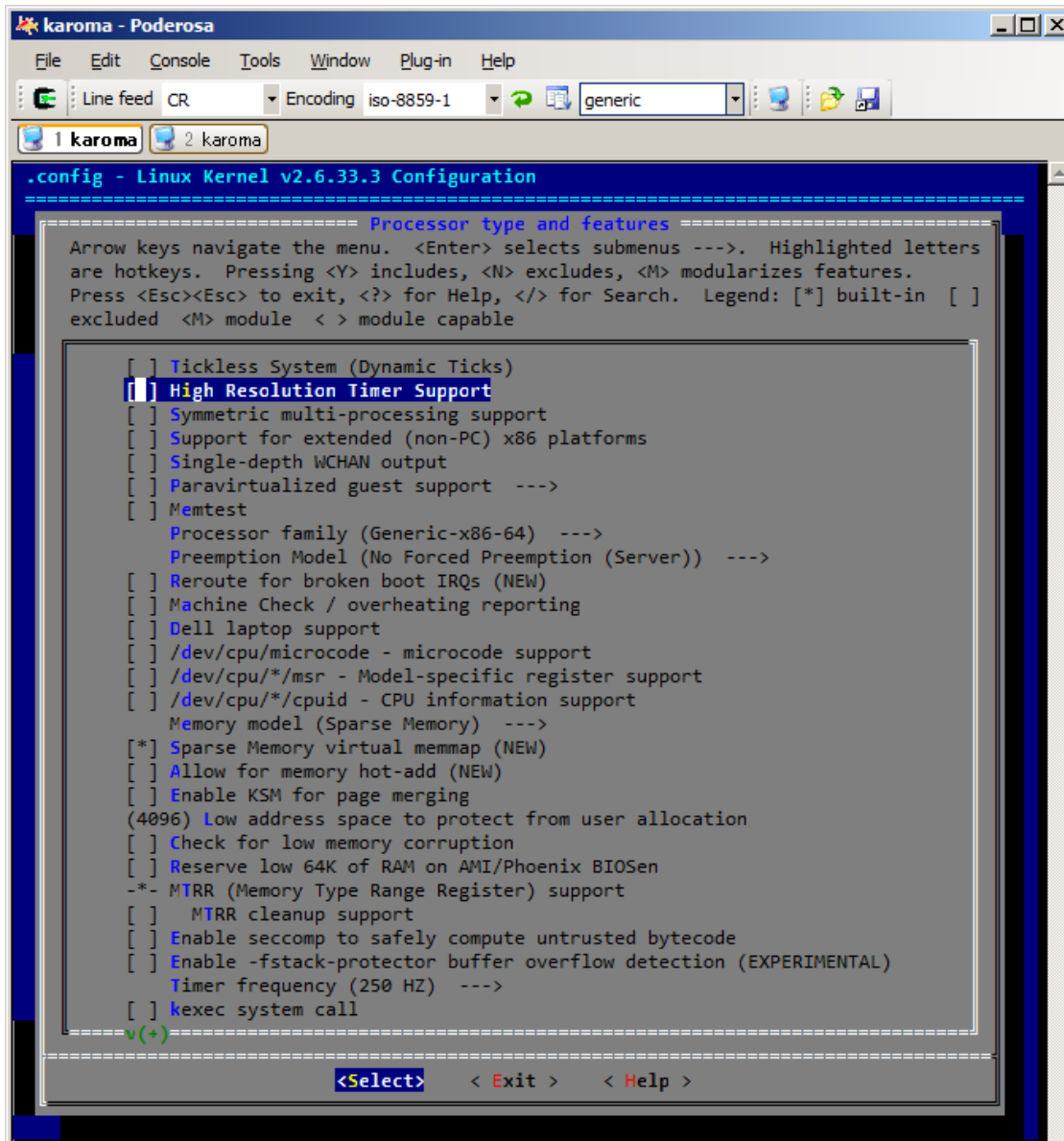
Precision und Recall

- ▶ (Genauigkeit und Trefferquote)
- ▶ Insbesondere bei Information Retrieval und Vorhersagen

- ▶ Hier: Beispiel Variability Mining

Kästner, Dreiling, and Ostermann. **Variability Mining with LEADT**. Technical Report 01/2011, Department of Mathematics and Computer Science, Philipps University Marburg, September 2011.

Exkurs: Produktlinien



VEGETARIAN

WHICH WICH WOULD YOU LIKE?

☐ TRIPLE CHEESE MELT
☐ ELVIS WICH (gn. Honey & Banana)
☐ TOMATO & AVOCADO
☐ BLACK BEAN PATTY
☒ HUMMUS & BELL PEPPERS

CHOOSE YOUR BREAD

☐ WHITE ☒ WHEAT

CHOOSE YOUR CHEESE (Optional)

☐ AMERICAN ☐ SWISS ☐ PROVOLONE
☐ CHEDDAR ☒ PEPPER JACK ☐ MOZZARELLA

How Would You Like Your WICH Worked?

MUSTARDS
☐ Yellow ☐ Dijon ☐ Honey ☒ Dell

MAYOS
☐ Regular ☐ Lite ☐ Horseradish ☒ Spicy

SPREADS & SAUCES
☐ BBQ ☐ Buffalo ☐ Marinara
☐ 1000 Island ☐ Ranch

ONIONS
☒ Red ☐ Grilled ☐ Crispy Strings

VEGGIES
☒ Lettuce ☒ Tomato ☐ Pickles ☒ Jalapenos
☒ Olive Salad ☐ Mushrooms ☐ Sauerkraut
☐ Coleslaw ☐ Bell Peppers

OILS & SPICES
☐ Oil ☐ Vinegar
☒ Salt ☒ Pepper ☐ Oregano ☐ Parmesan

EXTRAS (.75¢ Each)
☐ Bacon ☐ Avocado ☐ Pickle (Whole)
☐ More Meat ☐ More Cheese

Exkurs: Adaption von Produktlinien

- ▶ Oft bereits Quelltext vorhanden
- ▶ Nachtraeglich Variabilitaet hinzufuegen
- ▶ Aufgabe: Finde vorhandenen Quelltext der Feature X implementiert
- ▶ Gegeben Startwert:
Werkzeug schlaegt
Fragmente vor

```
static int __rep_queue_filedone(dbenv,  
rep, rfp)  
    DB_ENV *dbenv;  
    REP *rep;  
    __rep_fileinfo_args *rfp; {  
#ifdef HAVE_QUEUE  
    db_pgno_t first, last;  
    u_int32_t flags;  
    int empty, ret, t_ret;  
#ifdef DIAGNOSTIC  
    DB_MSGBUF mb;  
#endif  
    // over 100 lines of additional code  
#endif  
}
```

Fragestellung:

Hilft das Werkzeug Entwicklern beim
finden von Feature-Quelltext?

Welche Qualitaet haben die Vorschlage?

The screenshot displays the Eclipse IDE interface with several panels and views. The Project Explorer on the left shows a project named 'MobileMedia7-CIDE-Compliant' with a tree structure including 'src', 'lanes.midp.mobilephoto', 'ubc.midp.mobilephoto', 'core', 'comms', 'threads', 'ui', 'util', 'sms', 'JRE System Library [JavaSE-1]', 'Mobile', 'model.colors', and 'model.m'. The Source Editor in the center shows a Java file with code for 'Data().addMediaData(imageData, albumname);' and an 'ImageDataException' handler. The Features view at the bottom left shows a list of features for 'MobileMedia7-CIDE-Compliant', including 'Not_Favourites', 'Not_SMS_Transfer', 'Play_Music', 'SMS_Transfer', 'SMS_or_Copy', 'Single_Media_Mode', 'View_Photo', and 'src/ubc/midp/mobilephot'. The Recommendation Manager at the bottom right shows a table of recommendations with columns for Name, Type, Rec's, T, Value, Reasons, S..., > Val..., Range, and Views.

A Project Explorer

B Source Editor

C Features

D FeatureManager

E FeatureManager

F Recommendation Manager

Code in scenario 06 */

```
Data().addMediaData(imageData, albumname);
```

ImageDataException e) {
= null;
ceof ImagePathNotValidException)
new Alert("Error", "The path is not valid", nul
new Alert("Error", "The image file format is no
isplay(midlet).setCurrent(alert, Display.getDis
Timeout(5000);
newMechanismException e) {

belongs to
+ VaDeFa: albumname
X: Transpose Relations
*requires
local variable access declared
local variable accessed by
local variable transitively acce
parameter accessed by
+ MetInv: addMediaData

base been created for MobileMedia7-CIDE-Compliant

Name	Type	Rec's
Not_Favourites		0
Not_SMS_Transfer		0
Play_Music		0
SMS_Transfer		78
SMS_or_Copy		74
Single_Media_Mode		0
View_Photo		91
src/ubc/midp/mobilephot		91
ExpressionStatement		91
IfStatement		91
+ InsCre: MediaListSc METHO...		2

Name	T	Value	Reasons	S...	> Val...	Range	Views
+ SimNam: albumname	7	1.0	TC:Check Accesses...	1		4392-4...	7
+ VaDeFa: imageData	4	0.76	TPF:Text Match(1), ...	4		2789-2...	0
+ SimNam: imageData	7	0.75	GR:accessed by	1		3001-3...	0
+ SimNam: imageData	7	0.75	GR:accessed by	1		4234-4...	0
+ SimNam: imageData	7	0.37	GR:accessed by	1		4381-4...	1
View_Photo	7	0.05	TPF:Text Match	1		0-00	0
+ MetDec: addMediaData	2	0.55	TC:Check Decl.(1), ...	3	SMS_o...	3203-3...	0
+ TypDec: PhotoViewScri	1	0.55	GR:declared by(5), ...	10	SMS_o...	390-39...	0
+ ImpDec: javax.microed	5	0.5	GR:accessed by(1)	1		198-19...	0
+ ImpDec: ubc.midp.mol	5	0.5	GR:accessed by(2)	2		1077-1...	0
+ VaDeFa: nextcontroller	4	0.5	GR:referenced by(1...	3		15107-...	0

Vorschlagsqualitaet messen

- ▶ **Praezision (Precision)** $\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$
 - ▶ Anteil des gefunden relevanten Quelltexts
 - ▶ Gefundene Element / Korrekte Elemente
 - ▶ Aber: Alles vorschlagen -> 100% Praezision
- ▶ **Genauigkeit (Recall)** $\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}$
 - ▶ Anteil korrekter Vorschlaege
 - ▶ Gefundene Elemente / Vorgeschlagene Elemente
 - ▶ Aber: Nach einem Vorschlag abbrechen: 0/100% Genauigkeit
- ▶ **F-Measure (harmonisches Mittel)** $F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$

Beispiel: Precision/Recall

10.000 Zeilen Code, davon 1000 Relevant

Aufgabe: Berechnen Sie Recall, Precision und F-Measure. Diskutieren Sie welches Verfahren besser ist.

	Relevant	Nicht Relevant
Gefunden	900	2000
Nicht gefunden	100	7000

Precision: $900/1000 = 0.9$

Recall: $900/2900 = 0.31$

F-Measure: 0.46

Precision: $500/1000 = 0.5$

Recall: $500/503 = 0.99$

F-Measure: 0.66

	Relevant	Nicht Relevant
Gefunden	500	3
Nicht gefunden	500	8997



Orakel

- ▶ Woher bekannt welche Daten relevant sind?
- ▶ Möglichst neutrale Referenzdaten nehmen
 - ▶ Im Idealfall schon vorhanden und verifiziert
 - ▶ Bei manchen Problemen fertige Datensätze verfügbar
- ▶ Wenn selber erzeugt, Korrektheit belegen
 - ▶ Interpersonaler Konsens
- ▶ In unserem Beispiel:
 - ▶ MobileMedia: Fertige Produktlinie, code reviewed, von mehreren Forschern verwendet; 3 weitere Orakel
 - ▶ Nachteil: Relativ klein, externe Validität angreifbar
 - ▶ Feature-Markierungen entfernt

Project	Feature	Feature Size			Mining Results		
		LOC	FR	FI	IT	Recall	Prec.
Prevayler	Censor	105 (1 %)	10	5	32	100 %	41 %
	Gzip	165 (2 %)	4	4	27	100 %	18 %
	Monitor	240 (3 %)	19	8	53	100 %	42 %
	Replication	1487 (19 %)	37	28	64	100 %	67 %
	Snapshot	263 (3 %)	29	5	47	81 %	46 %
MobileM.	Copy Media	79 (2 %)	18	6	33	97 %	26 %
	Sorting	85 (2 %)	20	6	36	96 %	46 %
	Favourites	63 (1 %)	18	6	31	100 %	43 %
	SMS Transfer	714 (15 %)	26	14	44	100 %	62 %
	Music	709 (15 %)	38	16	51	99 %	59 %
	Photo	493 (11 %)	35	13	55	99 %	49 %
Lampiro	Media Transfer	153 (3 %)	4	3	25	99 %	13 %
	Compression	5155 (12 %)	33	20	42	100 %	66 %
	TLS Encryption	86 (0 %)	13	6	24	81 %	29 %
Sudoku	Variable Size	44 (2 %)	5	4	24	100 %	29 %
	Generator	172 (9 %)	9	7	29	98 %	42 %
	Solver	445 (23 %)	40	12	46	100 %	58 %
	Undo	39 (2 %)	5	4	29	100 %	21 %
	States	171 (9 %)	26	7	43	99 %	52 %

LOC: lines of code (and percentage of feature code in project's code base);

FR: Number of distinct code fragments; FI: Number of files; IT: Number of iterations

Learning Set vs. Evaluation Set

- ▶ Manche Verfahren “lernen”
- ▶ Manche Verfahren aufgrund Erkenntnisse aus Beispieldaten entwickelt
- ▶ Immer Daten zum Lernen/Entwickeln der Methode von Evaluierungsdaten trennen
 - ▶ Frische Daten fuer Evaluierung
 - ▶ Oder Datensatz teilen: zB. 80% zum Lernen, 20% zum Evalieren

Automatisierung

Automatisierte Datenerhebung (wenn moeglich)

- ▶ i.d.R. selbstimplementierte Infrastruktur
- ▶ Daten sammeln, Werkzeug/Metrik automatisch anwenden
- ▶ Erlaubt Evalierung auf grossen Datenbestaenden
 - ▶ Grosse Quelltextrepositories
 - ▶ Grosse Mailinglisten
 - ▶ Vereinfacht begruendete Auswahl
- ▶ Erlaubt Reproduktion und Variation z.B.:
 - ▶ Zufaelliche Seeds -> Konfidenzintervalle
 - ▶ Verschiedene Algorithmen -> Vergleiche
 - ▶ Verschiedene Reihenfolgen

Automatisierung bei Variability Mining

- ▶ Scripte starten Eclipse und folgen automatisch Vorschlaegen
 - ▶ Ersten Vorschlag ansehen
 - ▶ Entsprechend Orakel akzeptieren oder ablehnen
 - ▶ Weiter zum naechsten Vorschlag
 - ▶ (Weiter zum naechsten Feature)
 - ▶ (Weiter zum naechsten Projekt)
- ▶ Loggen in eine Datenbank
- ▶ Auswertung mit SQL (+Stored Procedures) und R

Vergleiche

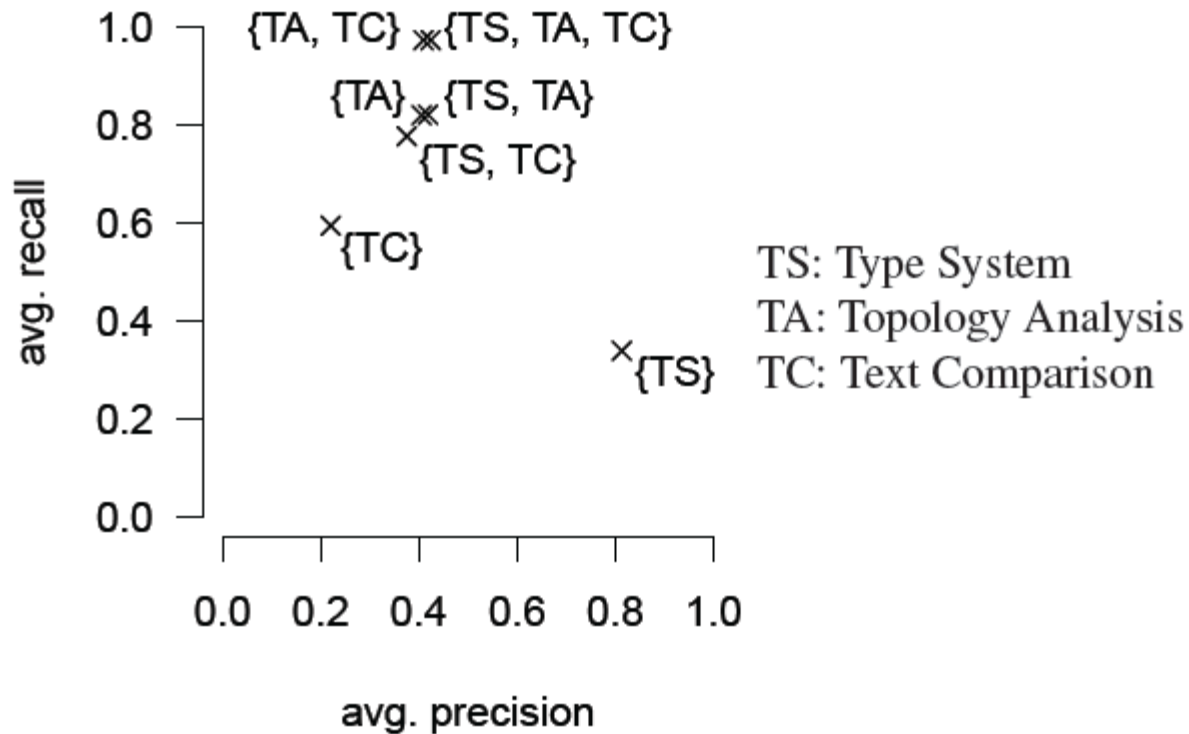


Figure 3: Combining recommendation mechanisms.

Stop-Kriterium

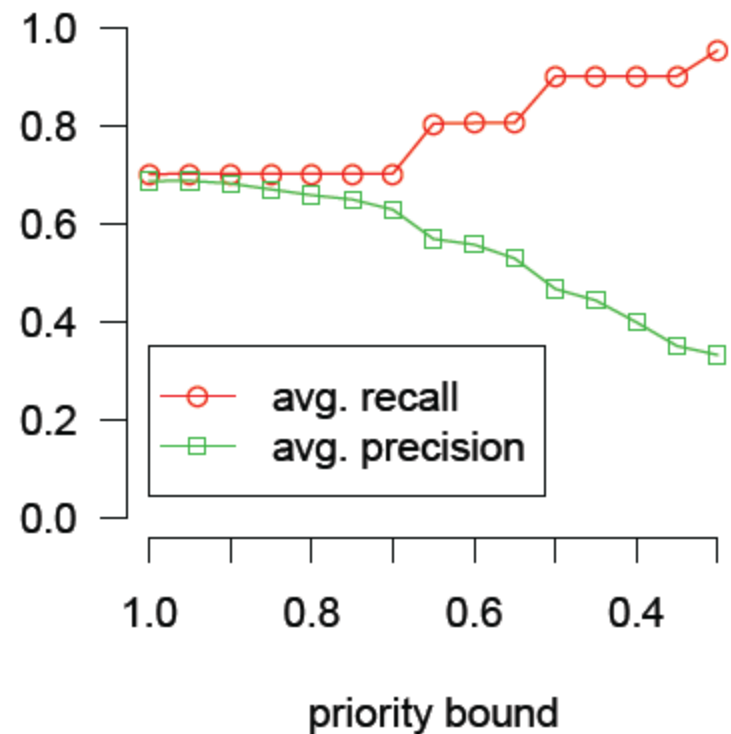
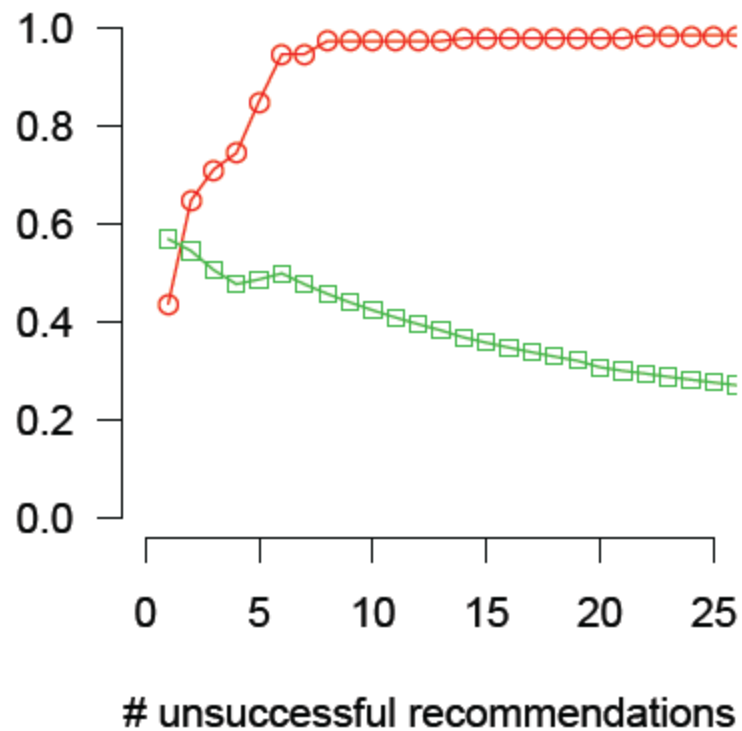


Figure 4: Alternative Stop Criteria.

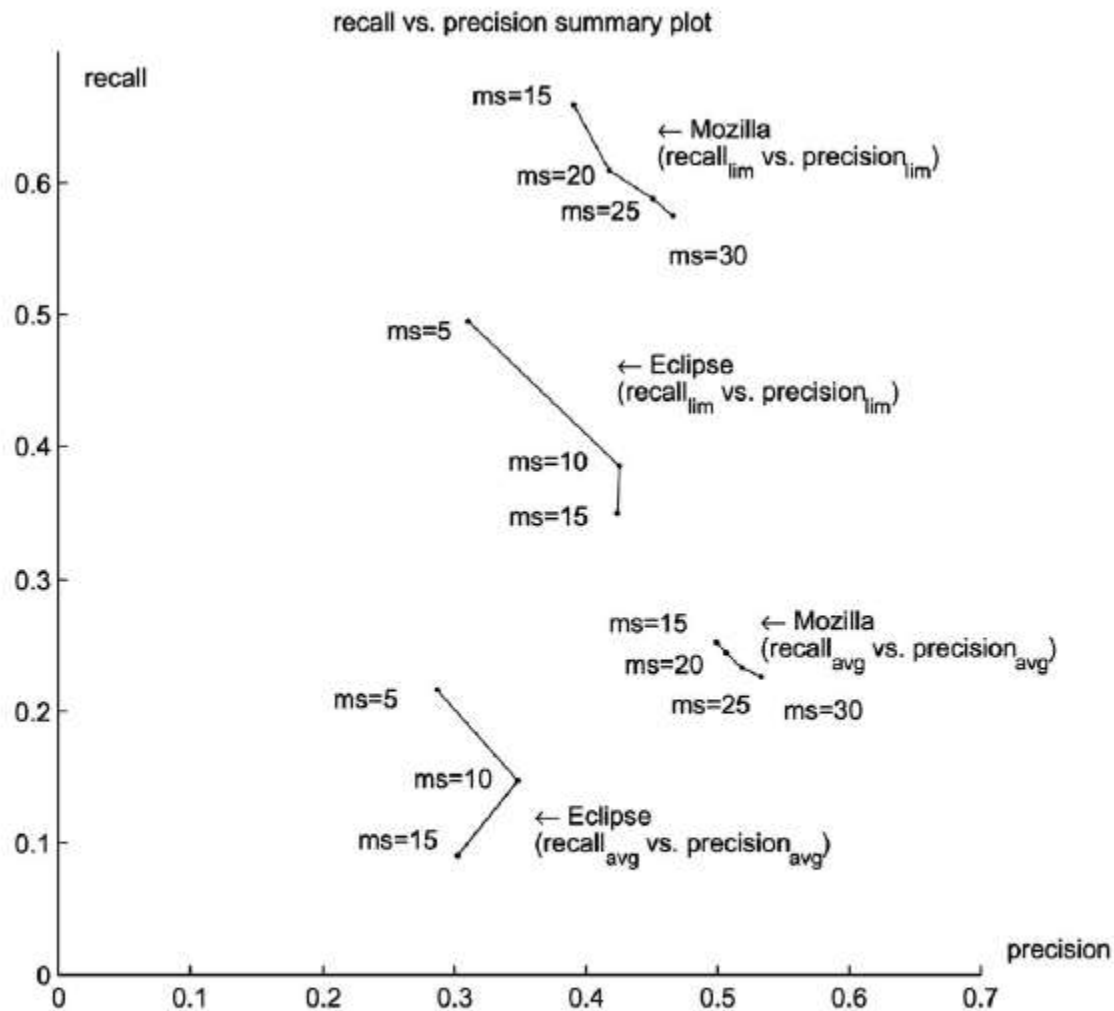


Fig. 2. Recall versus precision plot showing the frequent pattern algorithm applied to the two target systems, Eclipse and Mozilla.

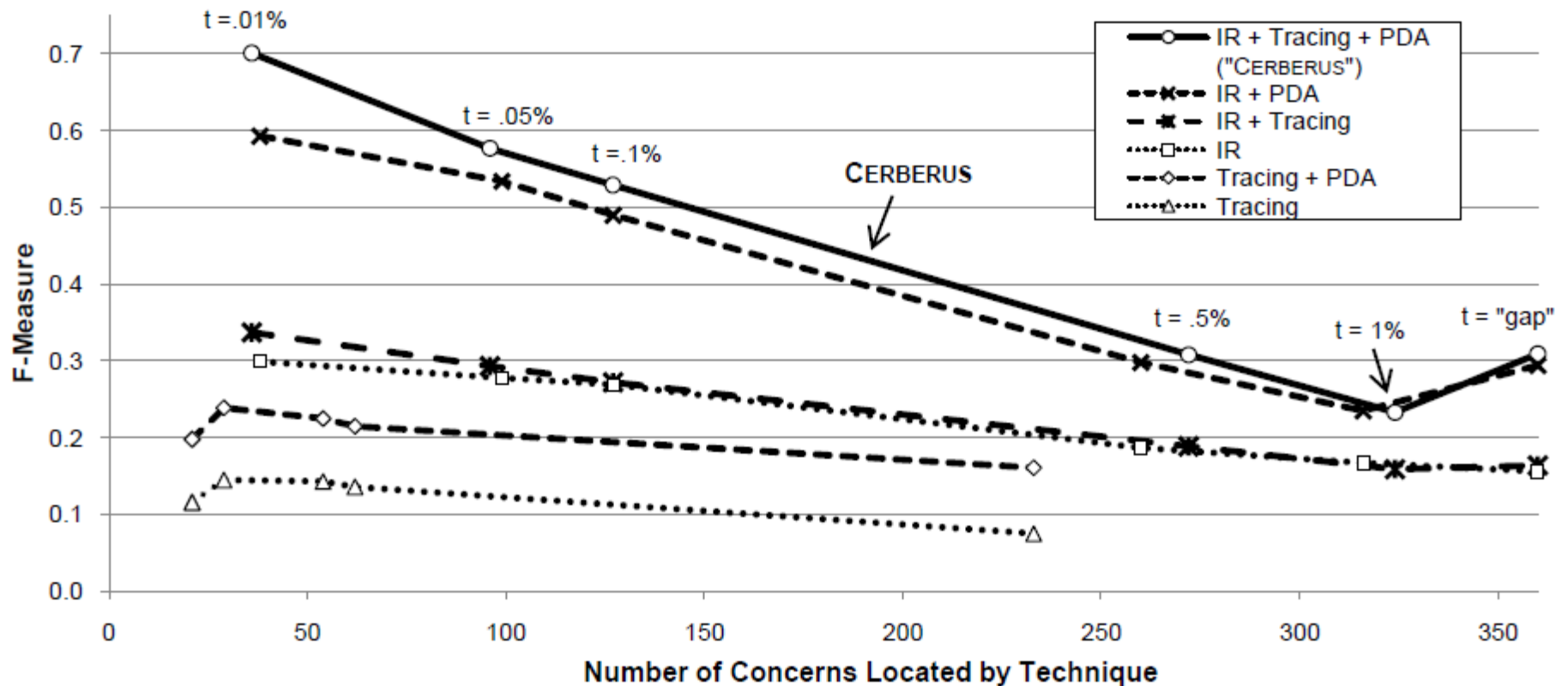


Fig 4. Effectiveness of individual concern location techniques and technique combinations for different threshold values t

Quelle: Eaddy et al. Cerberus: Tracing requirements to source code using information retrieval, dynamic analysis, and program analysis. In ICPC, 2008.

Software-Metriken

Software-Metriken

- ▶ Automatisierte Messung von Quelltexteigenschaften

Eine Softwaremetrik ist eine Funktion, die eine Software-Einheit in einen Zahlenwert abbildet. Dieser berechnete Wert ist interpretierbar als der Erfüllungsgrad einer Qualitätseigenschaft der Software-Einheit (IEEE Standard 1061, 1992)

- ▶ Metriken existieren fuer viele Eigenschaften
- ▶ Fokus typischerweise Qualität (Fehler, Verständlichkeit, ...)
- ▶ Definition eigener Metriken fuer spezielle Evaluierungen nicht unueblich

Metrik fuer Groesse

- ▶ Lines of code (LOC)
- ▶ Anzahl Dateien, Klassen, Methoden
- ▶ Grober Indikator fuer Projektgroesse
- ▶ Einfach zu messen

```
> wc -l file1 file2...
```

LOC	Projekte
450	Expression Evaluator
2.000	Graph Product Line, Sudoku, Functional Graph Library
4.500	MobileMedia, FAME-DBMS
8.000	Prevayler
20.000	MobileRSSReader
40.000	OpenVPN
80-100.000	Berkeley DB, SQLite
150-300.000	Apache, HyperSQL, Busybox, Emacs, Vim, ArgoUML
500-800.000	gimp, glibc, mplayer, php, SVN
1.600.000	gcc
6.000.000	Linux, FreeBSD
8.000.000	Open solaris
45.000.000	Windows XP

Normalisierungsversuche

- ▶ Kommentare und Leerzeilen ignorieren
- ▶ Zeilen mit ≤ 2 Zeichen ignorieren
- ▶ Quelltext auto-formatieren
- ▶ Statements zaehlen (In Java/C: Semikolon zaehlen) -> logical lines of code (LLOC)

```
for (i = 0; i < 100; i += 1) printf("hello"); /* How many lines of code is this? */
```

```
/* How many lines of code is this? */
```

```
for (  
    i = 0;  
    i < 100;  
    i += 1  
){  
    printf("hello");  
}
```

Normalisierungsversuch pro Sprache

Sprache	Statement-Faktor	Zeilen-Faktor
C	1	1
C++	2.5	1
Fortran	2	0.8
Java	2.5	1.5
Perl	6	6
Smalltalk	6	6.25
Python		6.5

Aber

- ▶ Kein Mass fuer Produktivitaet, Komplexitaet
- ▶ Generierter Quelltext
- ▶ Verschiedene Sprachen (assembler vs haskell)
- ▶ Keine allgemein akzeptierte Normalisierung
 - ▶ immer mit angeben
- ▶ Wird fast immer als Kontextinformation geben
- ▶ Nur sehr vorsichtig zusammen mit anderen Informationen interpretieren
- ▶ Eher selten als abhaengige oder unabhaengige Variable

Metrik fuer Komplexitaet

- ▶ Cyclomatic Complexity (zyklomatische Komplexität)
 - ▶ McCabe 1976
 - ▶ Anzahl der Kontrollfluesse durch Methode (oder Anzahl der Verzweigungen + 1)
 - ▶ (Tour de France)

```
if (c1) {  
    f1();  
} else {  
    f2();  
}  
if (c2) {  
    f3();  
} else {  
    f4();  
}
```

```
const String wochentagsName(const int  
nummer) {  
    switch(nummer)  
    {  
        case 1: return "Montag";  
        case 2: return "Dienstag";  
        case 3: return "Mittwoch";  
        case 4: return "Donnerstag";  
        case 5: return "Freitag";  
        case 6: return "Samstag";  
        case 7: return "Sonntag";  
    }  
    return "(unbekannter Wochentag)";  
}
```

Objekt-orientierte Metriken

- ▶ Number of Methods per Class
- ▶ Weighted Methods per Class
- ▶ Depth of Inheritance Tree
- ▶ Number of Child Classes
- ▶ Coupling between Object Classes
- ▶ ...

Weitere Metriken

- ▶ Bugs per line of code
- ▶ Comment density
- ▶ Halstead complexity measures
 - ▶ Derive software complexity from numbers of (distinct) operands and operators
- ▶ Test Coverage
- ▶ Number of classes and interfaces
- ▶ Abstractness = ratio of abstract classes to total number of classes
- ▶ ...

Vorteile Softwaremetriken

- ▶ Mechanisch erfassbar
- ▶ Auch komplexe Metriken in grossen Projekten
- ▶ Deterministisch
 - ▶ im Gegensatz zu Performance
 - ▶ keine Stichproben + Konfidenzintervalle noetig

Aussagekraft

- ▶ I.d.R Plausibilitaetsargumente
- ▶ Oft kontrovers
 - ▶ “What can we conclude about quality if the cyclomatic complexity of our code is 12?”
 - ▶ “Similar to the attempt of measuring the intelligence of a person in terms of the weight or circumference of the brain”
- ▶ Selten empirisch evaluiert
- ▶ Klassengroesse dominiert viele Metriken

Aussagekraft II

- ▶ Vorsicht bei Aussagen ueber menschliche Faktoren wie Verstaendniss und Lesbarkeit
- ▶ Vorsicht bei Wertungen zu Qualitaet
- ▶ Aber geeignet fuer Strukturbeschreibung, fuer interne Eigenschaften

- ▶ Metriken immer im Kontext bewerten
 - ▶ Lines of Code geeignet fuer Messung von Codegroesse
 - ▶ Lines of Code eher ungeeignet fuer Messung von Codequalitaet
 - ▶ Immer Zweck der Messung beruecksichtigen!

A vertical blue bar is located on the left side of the slide, within the same horizontal container as the title.

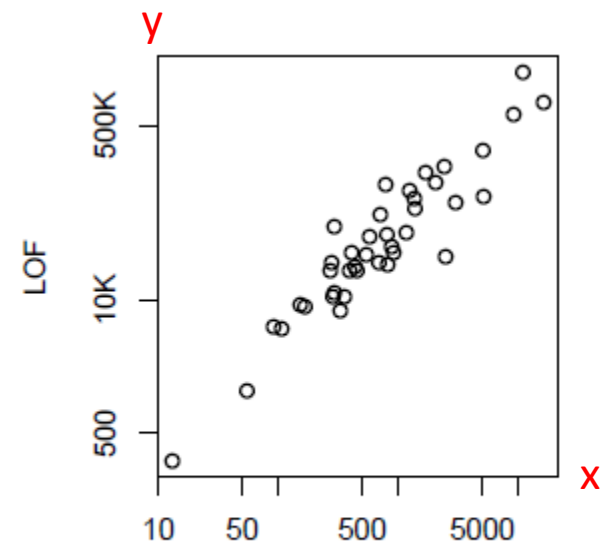
Korrelation

Forschungsfragen mit Metriken

- ▶ Unabhängige Variable X und abhängige Variable Y
- ▶ Einfluss von X auf Y, z.b.
 - ▶ Einfluss von Dateigrosse auf Fehleranfälligkeit
 - ▶ Einfluss von Kommentaren auf Verständlichkeit
 - ▶ Einfluss von GUI auf Benutzergeschwindigkeit
 - ▶ Einfluss von Heap-Groesse auf Ausführungsgeschwindigkeit
 - ▶ Einfluss von Anteil abstrakte Methoden auf Anzahl Testfälle
- ▶ Vergleich zweier Metriken (oder mehr)
 - ▶ Alle Messgrößen separat definieren und begründen
- ▶ Statistischer Zusammenhang?

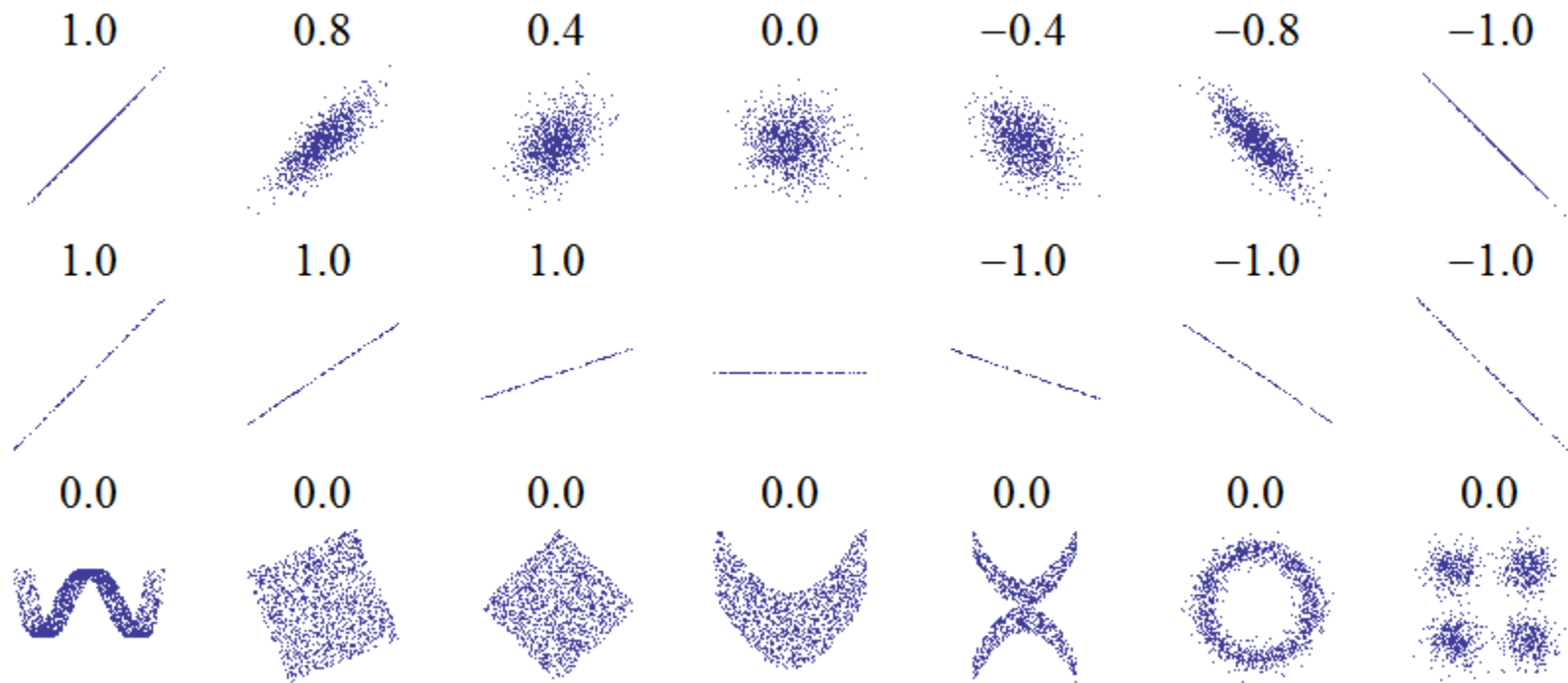
Korrelation

- ▶ Zusammenhang zwischen abhaengeriger Variable (x) und unabhaengeriger Variable (y)
 - ▶ beide min Intervallskalen
- ▶ Typisches Szenario: Es soll gezeigt werden dass x einen Einfluss auf y hat (nicht umgekehrt)
- ▶ Bei perfektem Zusammenhang alle Punkte auf einer Linie



c) NOFC (correlation coefficient: 0.74)

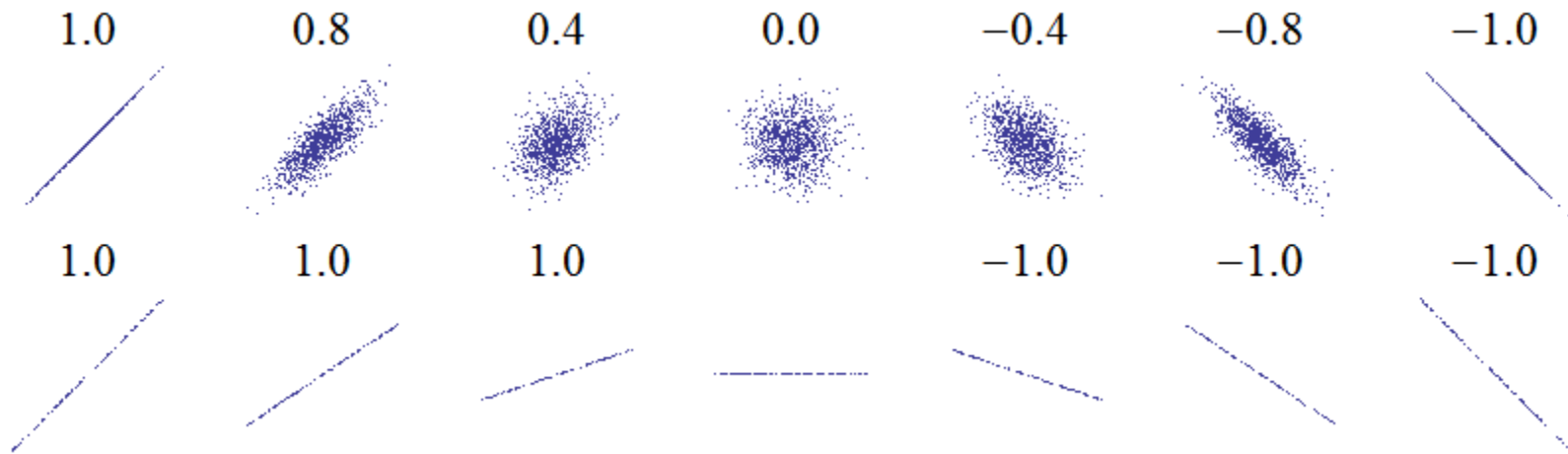
Korrelationskoeffizient



Korrelationskoeffizient

► Zwischen -1 und 1

$$\begin{aligned}\text{Kor}_e(x, y) &:= \varrho_e(x, y) := r_{xy} \\ &:= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}\end{aligned}$$



r	Interpretation
0	Keine Korrelation
<0.4	Schwache Korrelation
>0.6	Starke Korrelation
1	Perfekte Korrelation

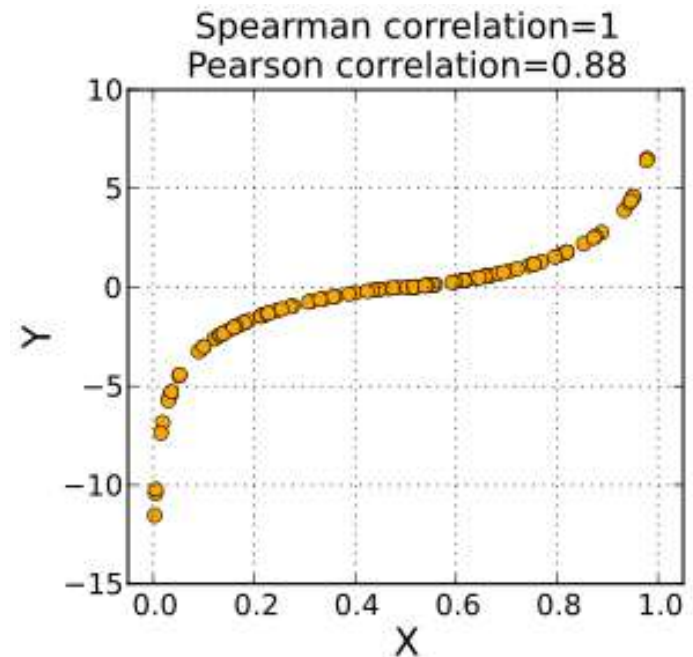
```
> x<-c(1,2,3,4)
> y<-c(4,4,6,8)
> cor(x,y)
[1] 0.9438798
```

Rangkorrelationskoeffizient

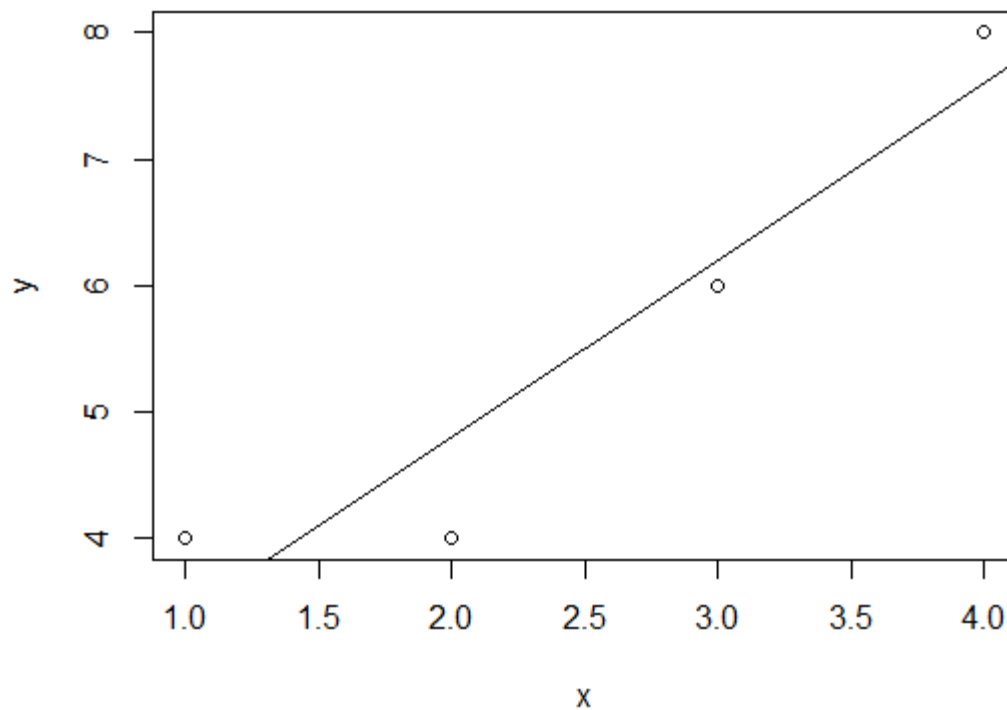
- ▶ Für ordinale Skalen (Abstände werden nicht interpretiert, nur Reihenfolge ist wichtig),
 - ▶ z.B. Likert-Skalen
- ▶ Spearman's Rangkorrelationskoeffizient (Spearman's Rho) und Kendall's Tau

```
> cor(x, y, method="spearman")  
1  
> cor(x, y, method="kendall")  
1
```

- ▶ Details in Statistikbuch



Regression



```
> x<-c(1,2,3,4)
> y<-c(4,4,6,8)
> cor(x,y)
[1] 0.9438798
> plot(y~x)
> abline(lm(y~x))
```

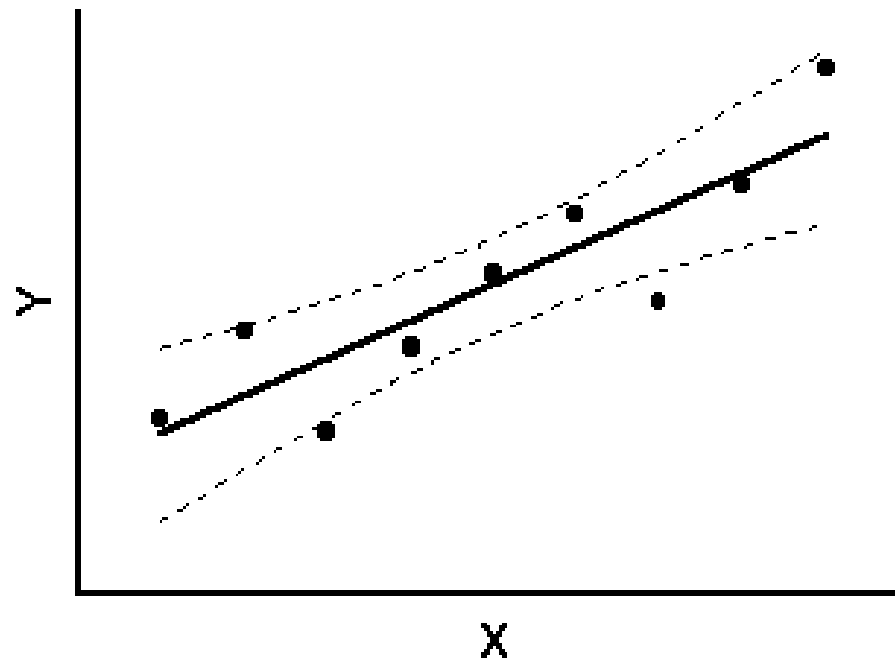
Signifikanter Korrelationskoeffizient

- ▶ Unterschied von r zu 0 statistisch signifikant?
 - ▶ Nullhypothese: Es besteht keine Korrelation
- ▶ Wenn signifikant:
 - ▶ Es besteht ein statistischer Zusammenhang zwischen x und y

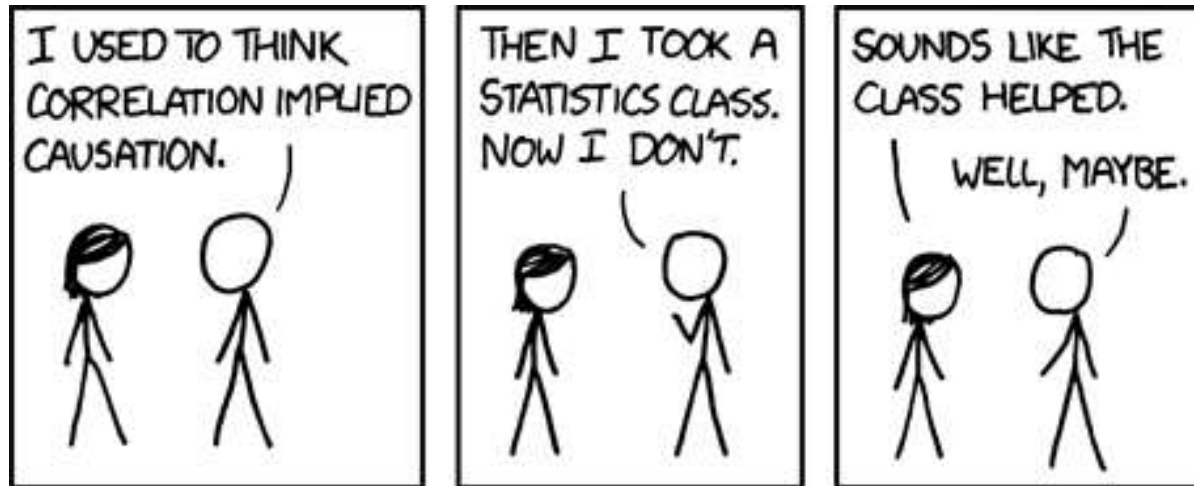
```
> cor.test(x,y)
data: x and y
t = 4.0415, df = 2, p-value = 0.05612
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval: -0.1853248 0.9988549
sample estimates: cor 0.9438798
```

- ▶ Kausalitaet ist damit
nicht belegt!

Regression mit Konfidenzintervallen



Korrelation und Kausalitaet



- ▶ Korrelation impliziert keine Kausalitaet
- ▶ Fuer Kausalitaet
 - ▶ Theorie noetig (aus Domaenenwissen, unabhaengig von den Daten)
 - ▶ Korrelation noetig
 - ▶ Muss neue Faelle **vorhersagen** koennen (Replikation/Validierung)

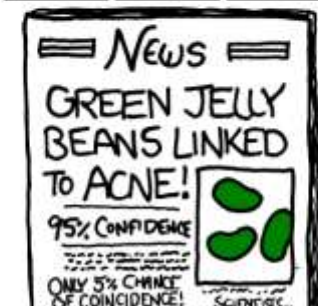
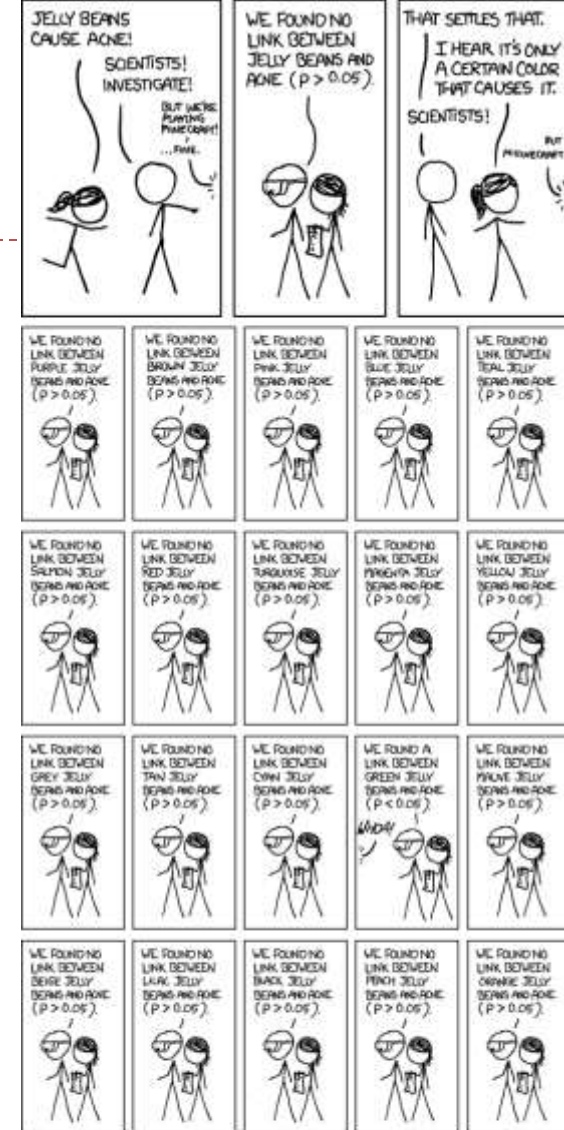
Fishing for Results (Wiederholung)

- ▶ Irgendwelche statistisch signifikanten Zusammenhaenge finden sich immer

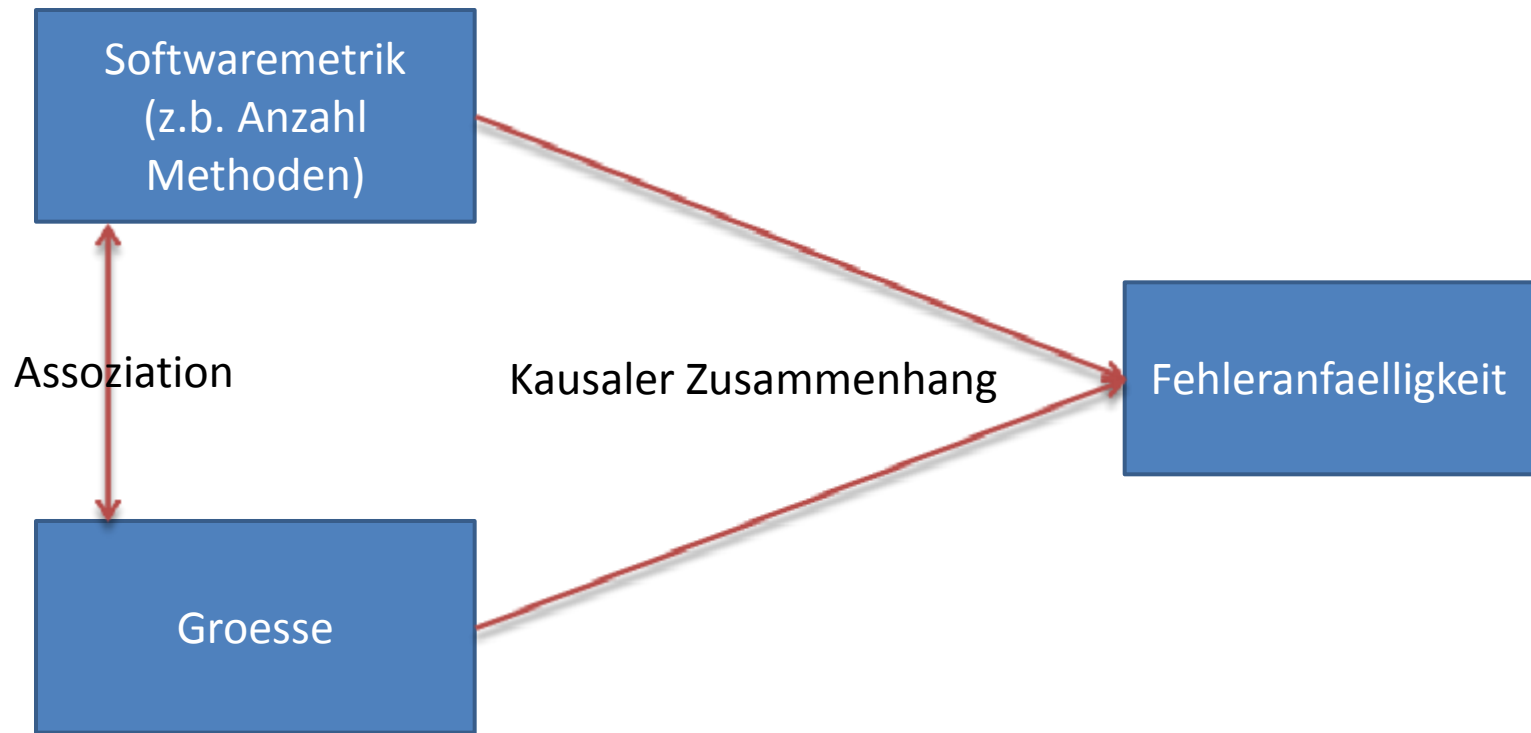
Deshalb:

- ▶ Hypothese vor Beginn des Experiments festlegen und begruenden
- ▶ Hypothesen, die sich beim Auswerten ergeben, als solche klar markieren
 - ▶ Erkundungsexperiment / Explorative Studie

“Fuer den sinnvollen Einsatz der Inferenzstatistik ist es erforderlich, dass vor Untersuchungsbeginn eine theoretisch gut begruendete Hypothese formuliert wurde.”



Stoervariablen



Emam et al. **The Confounding Effect of Class Size on the Validity of Object-Oriented Metrics.** In IEEE TSE 27 (7), 2001

Aufgabe: Metrik entwerfen

- ▶ Entwerfen Sie eine Metrik fuer die (Teil-)Evalierungen folgender Projekte:
 - ▶ Die Qualitaet von Quelltextdokumentation sinkt ueber die Lebenszeit d. Projekts
 - ▶ Quelltext von Projekten waechst linear, die Dokumentation aber weniger als linear
 - ▶ Generics in Java verlangsamten die Ausfuehrungszeit der Programme
 - ▶ Entwurfsmuster (Design Pattern) werden von Entwicklern haeufig benutzt
 - ▶ Die neue Multi-Touch Benutzeroberflaeche ist intuitiv bedienbar
 - ▶ Wie organisieren erfolgreiche Entwicklerteams ihren Quelltext/ihre Zusammenarbeit?
 - ▶ Quelltexte mit vielen Kommentaren sind schwer verstaendlich
 - ▶ Je mehr Entwickler an einem Projekt beteiligt sind, desto eher werden Deadlines ueberzogen
 - ▶ Klassen die viel getestet werden haben weniger Fehler
 - ▶ Fehler werden schneller gefixed, wenn sie einem Entwickler zugewiesen werden
- ▶ Begrunden Sie ihre Metriken anhand der Richtlinien von Kaner und Bond (Paper 12)

Zusammenfassung

- ▶ Metriken erlauben einfache Datensammlung
 - ▶ Einfache Analyse
 - ▶ Schwierige Interpretation
-
- ▶ Fishing for Results vermeiden!

Zeitliche Analyse

Historie

► Historische Entwicklung von Quelltexten/Dokumenten oft

► Versionsverwaltungssysteme

```
* 2009-12-08 575c2d8 bchesneau@gmail.com (bchesneau@gmail.com) Merge remote branch 'thomo/master' into mthomo
* 2009-12-07 a315955 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Merge branch 'master' of git://github.com/couchapp/couchapp
* 2009-12-01 cbd2bc8 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Bug: Trailing carriage returns are not stripped.
* 2009-11-28 c649863 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Merge branch 'master' of git://github.com/couchapp/couchapp
* 2009-11-27 1418fa1 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Fix wrong couchapp path
* 2009-11-27 1944cc0 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Merge branch 'master' of git://github.com/couchapp/couchapp
* 2009-11-27 ae97905 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Fix UTF-8 issue on WinXP
* 2009-11-26 97b5026 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Merge branch 'master' of git://github.com/couchapp/couchapp
* 2009-11-19 00ed4d9 Thomas.Mohaupt@gmail.com (Thomas.Mohaupt@gmail.com) Fix WinXP problem: Attachment name is not UTF-8 encoded
* 2009-11-17 6a52137 mot@sq.s.int (mot@sq.s.int) Fix: USERPROFILE handling
* 2009-12-08 34828d0 bchesneau@gmail.com (bchesneau@gmail.com) fix hooks and compress hook
* 2009-12-08 60dec82 bchesneau@gmail.com (bchesneau@gmail.com) fix compress options search path
* 2009-12-07 3c89e38 jasondavies@apache.org (jasondavies@apache.org) Generate attachment signatures before processing macros.
* 2009-12-07 54c3e83 jasondavies@apache.org (jasondavies@apache.org) Merge branch 'master' of git://github.com/couchapp/couchapp
v20080122-0800a, v20080122-0800 21.01.08 13:04 bbau... 215965 [breadcrumb] Activating editor should set focus into breadcrumb if it was...
v20080115-0800, v20080108-0800, ... 21.06.07 14:44 dmeiert Fixed bug 182108: [api] Add action constants for IDEActionFactory and ActionF...
* 2009-12-04 9b97679 bchesneau@gmail.com (bchesneau@gmail.com) fix __iter__
* 2009-12-04 4b8c9e7 bchesneau@gmail.com (bchesneau@gmail.com) fix typo. thanks to markh.
* 2009-12-04 b37230e bchesneau@gmail.com (bchesneau@gmail.com) patch from @rmg. thanks!
* 2009-11-29 233238a bchesneau@gmail.com (bchesneau@gmail.com) bump version number to 0.5.1.
* 2009-11-29 d8c9709 bchesneau@gmail.com (bchesneau@gmail.com) add update template to generators
* 2009-11-28 0964b63 bchesneau@gmail.com (bchesneau@gmail.com) couchapp standalone for macosx via py2apps
R3_1_2, r312_v20060104-0800, r311... 17.06.05 17:51 dmeiert Updated copyright date to 2005
* 2009-11-28 4def0a6 bchesneau@gmail.com (bchesneau@gmail.com) make sure ocontent type is defined
* 2009-11-28 10b53e3 benoitc@.none (benoitc@.none) new fixes while testing
* 2009-11-28 051e778 bchesneau@gmail.com (bchesneau@gmail.com) rethink hook system. final version. Now like extensions each hooktype is a list of
* 2009-11-28 c47718f bchesneau@gmail.com (bchesneau@gmail.com) rethink extensions. So now extensions are a list of key=value pair. where key is
* 2009-11-27 53a2040 bchesneau@gmail.com (bchesneau@gmail.com) fix compress extension & load extensions like we load hooks
* 2009-11-27 343862d bchesneau@gmail.com (bchesneau@gmail.com) fix template paths with windows
* 2009-11-27 bfa2f8c bchesneau@gmail.com (bchesneau@gmail.com) backport import_module from python 2.7
* 2009-11-27 3380a4c benoitc@.none (benoitc@.none) more windows fix
* 2009-11-27 a049b41 bchesneau@gmail.com (bchesneau@gmail.com) more fixes
trailing whitespace and organized the imports
```

Beispielanalysen ueber Zeit

- ▶ Wachstum Quelltextmenge
- ▶ Adoption von Entwurfsmustern (Design Pattern)?
- ▶ Dauer einen Fehler zu fixen abhaengig von Prioritaet?
- ▶ Performanceverbesserungen?
- ▶ Alter Quelltext sicherer?

Vorgehen

- ▶ Revisionen einzeln analysieren
(automatisieren wenn moeglich!)
- ▶ Deltas analysieren
- ▶ Commits vs. Releases
- ▶ Vergleich von Metriken ueber Revisionen hinweg

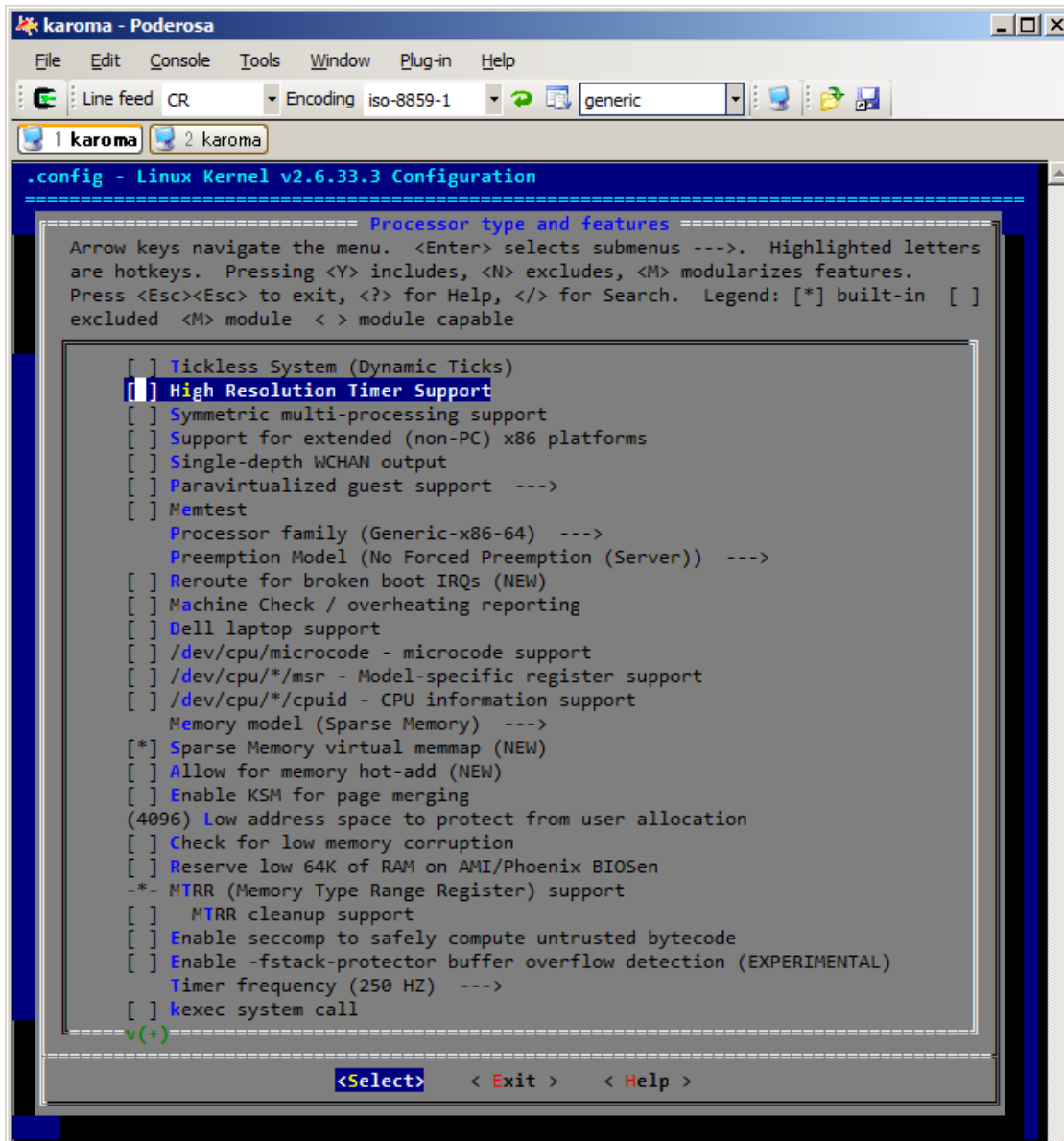
- ▶ Exploration vs. Hypothesen pruefen
- ▶ Stichproben vs. deterministische Messungen

Beispiel: Evolution der Variabilität im Linux Kernel

- ▶ Untersuchen den Linux Kernel
- ▶ Version 2.6.12 bis 2.6.32 (ca. 5 Jahre)
- ▶ Nur Releases
- ▶ Forschungsfragen (Exploration):
 - ▶ Welche Änderungen werden in der Praxis an Variabilitäts-Modellen durchgeführt?
 - ▶ Welche Modell-Refactorings werden in der Praxis eingesetzt?

Lotufo et al. Evolution of the Linux Kernel Variability Model. In SPLC, 2010.

Exkurs: Produktlinien



VEGETARIAN

WHICH WICH WOULD YOU LIKE?

☐ TRIPLE CHEESE MELT
☐ ELVIS WICH (gn. Honey & Banana)
☐ TOMATO & AVOCADO
☐ BLACK BEAN PATTY
☒ HUMMUS & BELL PEPPERS

CHOOSE YOUR BREAD

☐ WHITE ☒ WHEAT

CHOOSE YOUR CHEESE (Optional)

☐ AMERICAN ☐ SWISS ☐ PROVOLONE
☐ CHEDDAR ☒ PEPPER JACK ☐ MOZZARELLA

How Would You Like Your WICH Worked?

MUSTARDS
☐ Yellow ☐ Dijon ☐ Honey ☒ Deli

MAYOS
☐ Regular ☐ Lite ☐ Horseradish ☒ Spicy

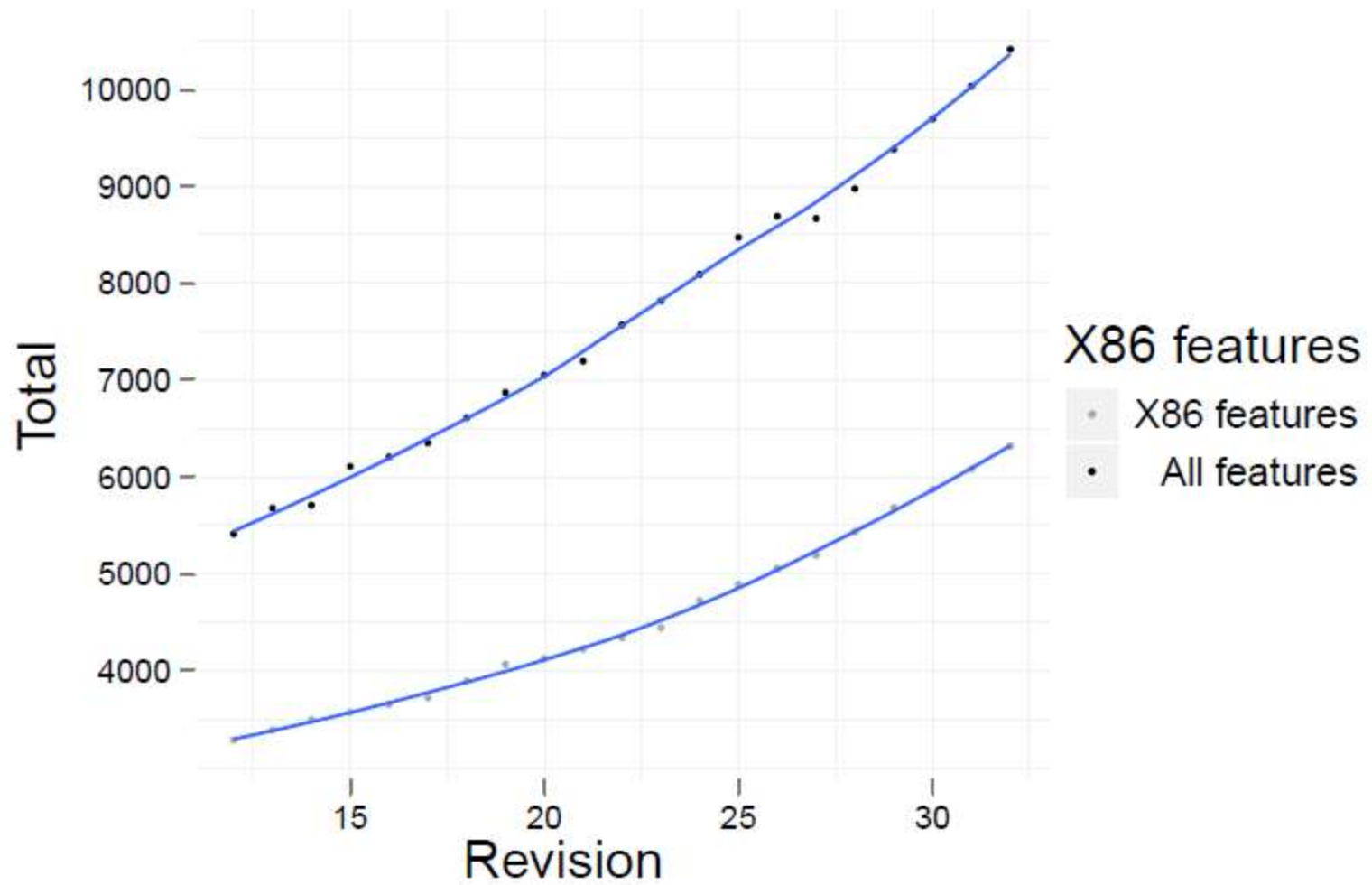
SPREADS & SAUCES
☐ BBQ ☐ Buffalo ☐ Marinara
☐ 1000 Island ☐ Ranch

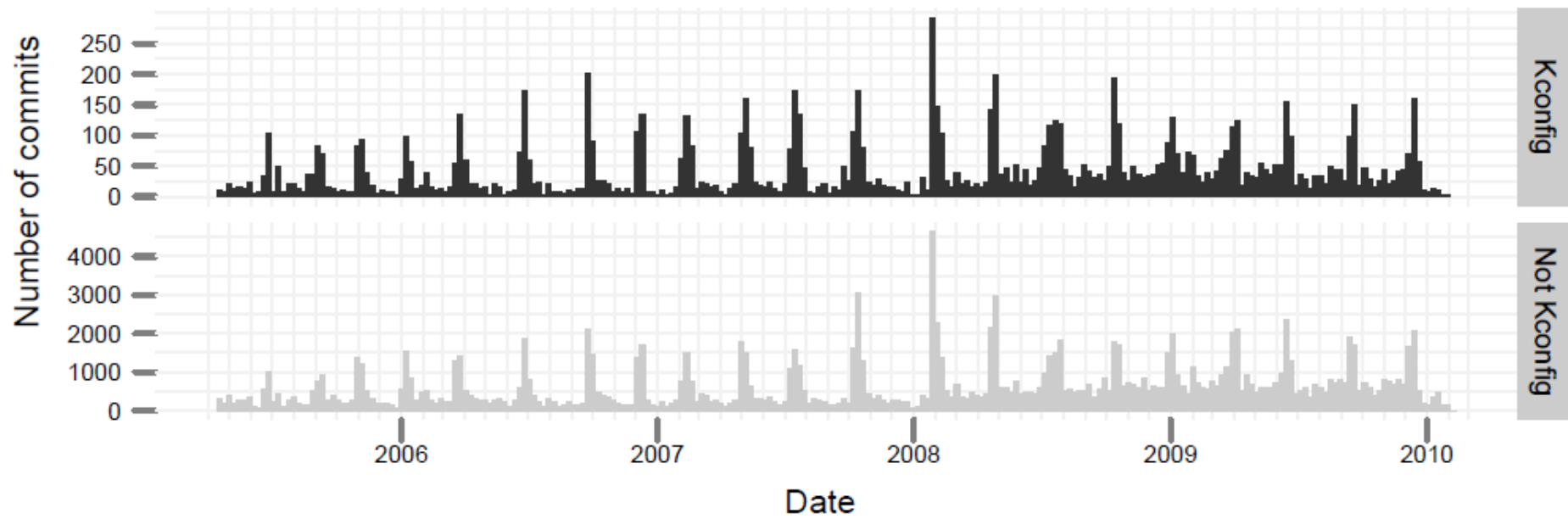
ONIONS
☒ Red ☐ Grilled ☐ Crispy Strings

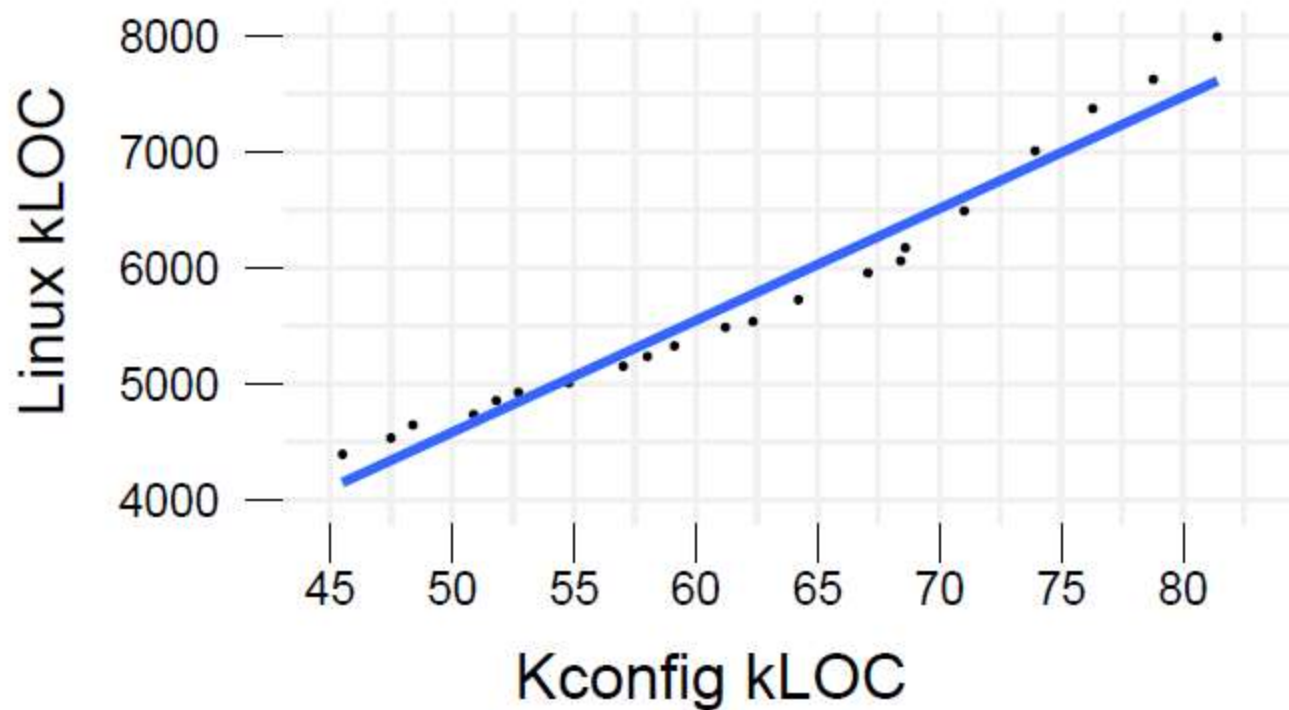
VEGGIES
☒ Lettuce ☒ Tomato ☐ Pickles ☒ Jalapenos
☒ Olive Salad ☐ Mushrooms ☐ Sauerkraut
☐ Coleslaw ☐ Bell Peppers

OILS & SPICES
☐ Oil ☐ Vinegar
☒ Salt ☒ Pepper ☐ Oregano ☐ Parmesan

EXTRAS (.75¢ Each)
☐ Bacon ☐ Avocado ☐ Pickle (Whole)
☐ More Meat ☐ More Cheese



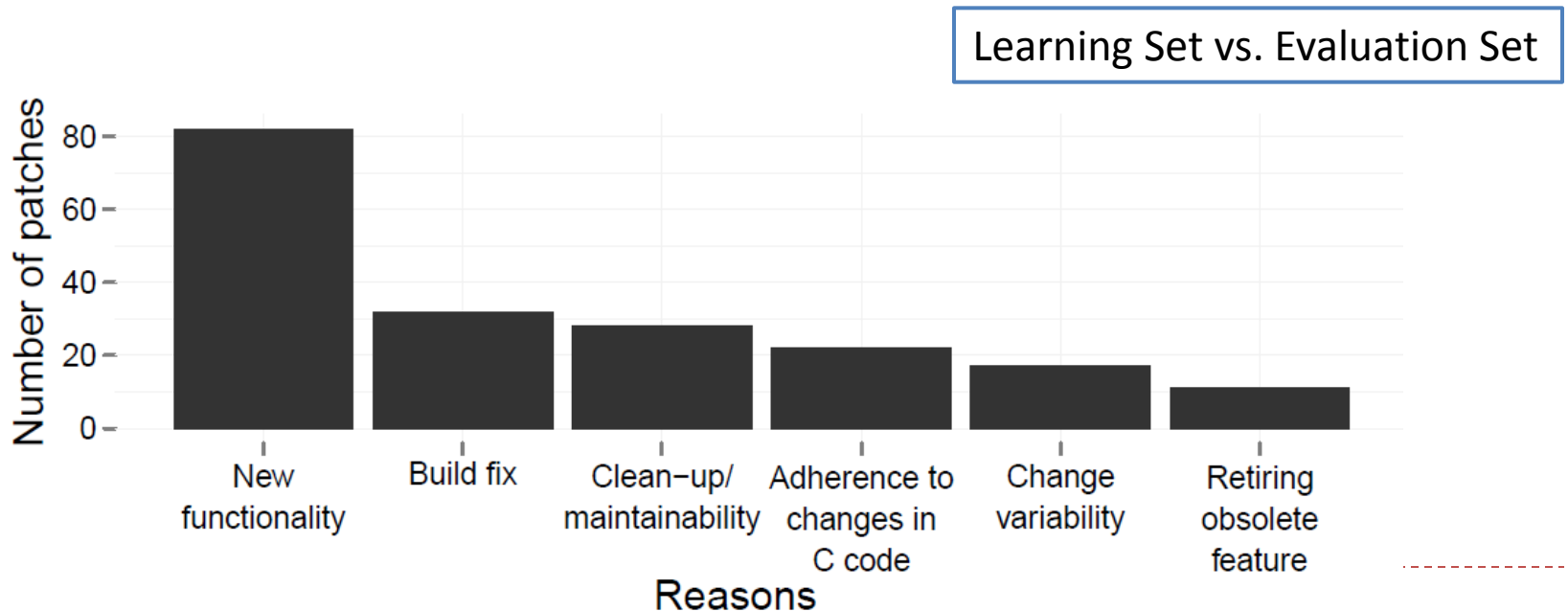




(a) Total kLOC for Linux against kLOC of Kconfig files

Gründe fuer Aenderungen

- ▶ Stichprobe 200 von 8726 Commits die das Modell betreffen
- ▶ Manuelle Analyse von Commit-Nachrichten: Kategorien identifizieren
- ▶ Erneute Stichprobe (in X86) zum Sammeln der Ergebnisse:



Weitere Ergebnisse

- ▶ Hinweise auf die Schwierigkeit durch unuebersichtliche Abhaengigkeiten
 - ▶ Nachtraegliche Korrekturen
- ▶ Hinweise fuer Werkzeugdesigner
 - ▶ Typische Operationen
 - ▶ Synchronisation Quelltext-Modell



Datenherkunft (Kooperationen mit Firmen)

Herkunft der Daten/Quelltexte

- ▶ Selbst entwickelt
- ▶ Uniprojekt
- ▶ Open Source
- ▶ Firmen-Intern

- ▶ Tradeoff
 - ▶ Realismus
 - ▶ Verfügbarkeit, Kosten
 - ▶ Publizierbarkeit
 - ▶ Neutralität

Quelltext Anonymisieren

- ▶ Problem insb. bei Praktika oder Diplomarbeiten in der Industrie
- ▶ Analysieren von vertraulichen betrieblichen Daten
- ▶ Daten aggregieren
- ▶ Namen entfernen (durch A B C ersetzen)
- ▶ Absprechen was veröffentlicht werden darf (**Vor der Untersuchung!**)
- ▶ Abwaegen
 - ▶ Erhoehtes Vertrauen noetig, da nicht extern replizierbar
 - ▶ Realismus

Empfehlungen des Fakultätentags Informatik für Abschlussarbeiten in Informatikstudiengängen an Universitäten

13. Juni 2007



3. Schriftliche Ausarbeitung

Die schriftliche Ausarbeitung ist der wesentliche Repräsentant der Prüfungsleistung und soll im Stil einer wissenschaftlichen Abhandlung angefertigt werden. Der eigene Anteil an den Ergebnissen muss klar erkennbar sein. Alle zur Begutachtung erforderlichen Artefakte (z.B. Quellcode) müssen den Gutachtern zur Verfügung gestellt werden.



5. Veröffentlichung

Alle Abschlussarbeiten sollen der Öffentlichkeit zugänglich gemacht werden. Eine Veröffentlichung wissenschaftlicher Ergebnisse in der einschlägigen Fachliteratur darf nicht durch einen Vertrag ausgeschlossen werden.

[http://www.ft-informatik.de/uploads/tx_sbdownloader/
empfehlungen_abschlussarbeiten.pdf](http://www.ft-informatik.de/uploads/tx_sbdownloader/empfehlungen_abschlussarbeiten.pdf)





Fallstudien (nochmal)

Fallstudien

- ▶ Viele vorgestellte Analysen waren Fallstudien
- ▶ Analysen einzelner Systeme/Faelle
- ▶ Exploration
- ▶ Geziele Quantifizierung der Daten
 - ▶ Aggregation, Kompaktere Darstellung
- ▶ Bedingter Anspruch zur Verallgemeinerung

Quantitative Evaluierung in Papern

- ▶ Studien zu Autovervollstaendigung
 - ▶ Robbes and Lanza. **How Program History Can Improve Code Completion.** In Proc. ASE 2008.
 - ▶ Hou and Pletcher. **An Evaluation of the Strategies of Sorting, Filtering, and Grouping API Methods for Code Completion.** In Proc. ICSM 2011.
 - ▶ Bruch, Monperrus, and Mezini. **Learning from Examples to Improve Code Completion.** In Proc. ESEC-FSE 2009.
- ▶ Was wurde gemessen? Welche Daten, welche Metriken?
- ▶ Ueberzeugend? Positives? Kritikpunkte? Validitaet?

Literatur

- ▶ Singer et al. **Software Engineering Data Collection for Field Studies**. In Guide to Advanced Empirical Software Engineering. Springer 2007.
 - ▶ Beispiele und weiterfuehrende Referenzen zu Instrumentierung
- ▶ Fenton and Pfleeger. **Software Metrics. A Rigorous and Practical Approach**. Thomson Learning 2000.
- ▶ Kaner and Bond. **Software Engineering Metrics: What Do They Measure and How Do We Know?** In Proceedings Int'l Metrics Symposium, 2004.
- ▶ Statistikbuch