

# Neuerungen im Enterprise Miner 5.2 & Text Miner 2.3

Ulrich Reincke, SAS Deutschland



# Agenda

Der Neue Enterprise Miner 5.2

Der Neue Text Miner 2.3

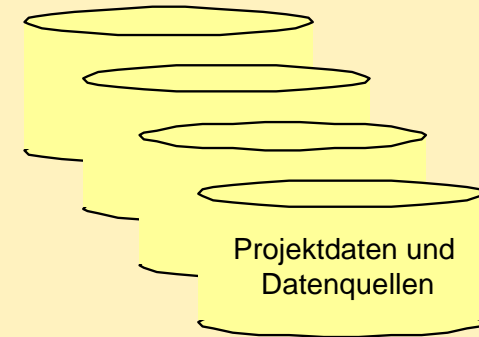
Diskussion

# Data Mining bisher:

Business Analyst  
(Modellentwicklung)



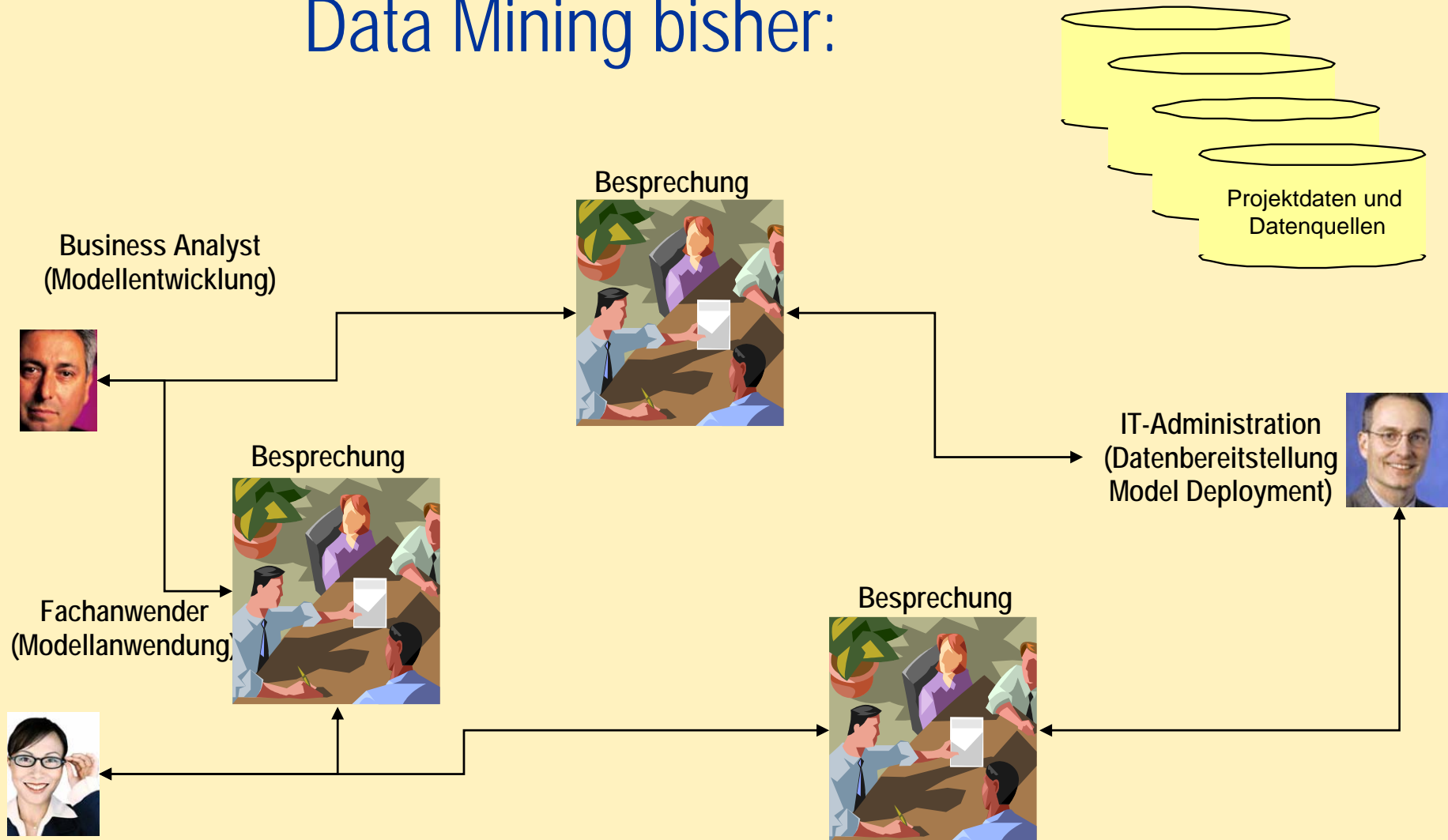
Fachanwender  
(Modellanwendung)



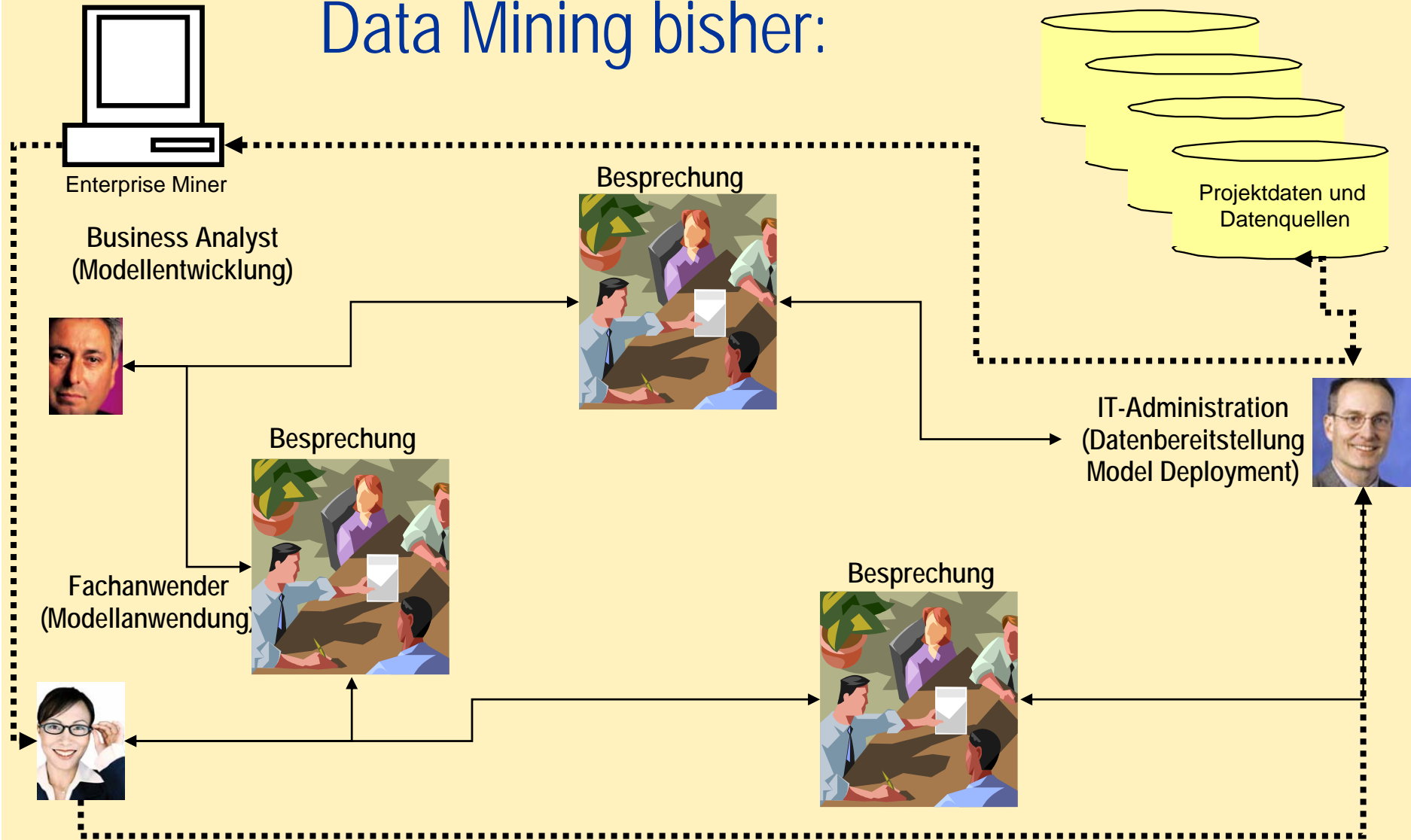
IT-Administration  
(Datenbereitstellung  
Model Deployment)

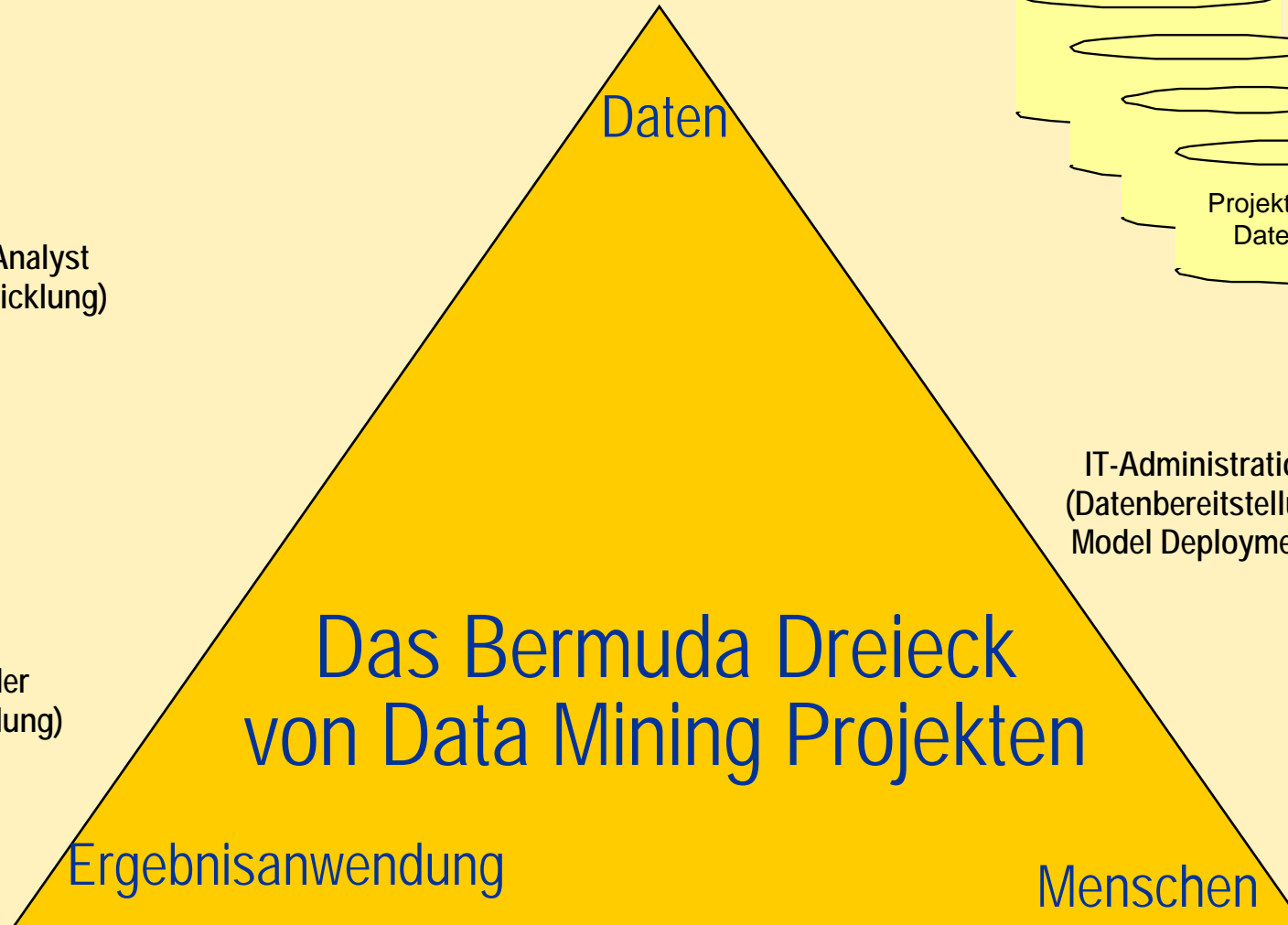


# Data Mining bisher:



# Data Mining bisher:





Business Analyst  
(Modellentwicklung)



IT-Administration  
(Datenbereitstellung  
Model Deployment)

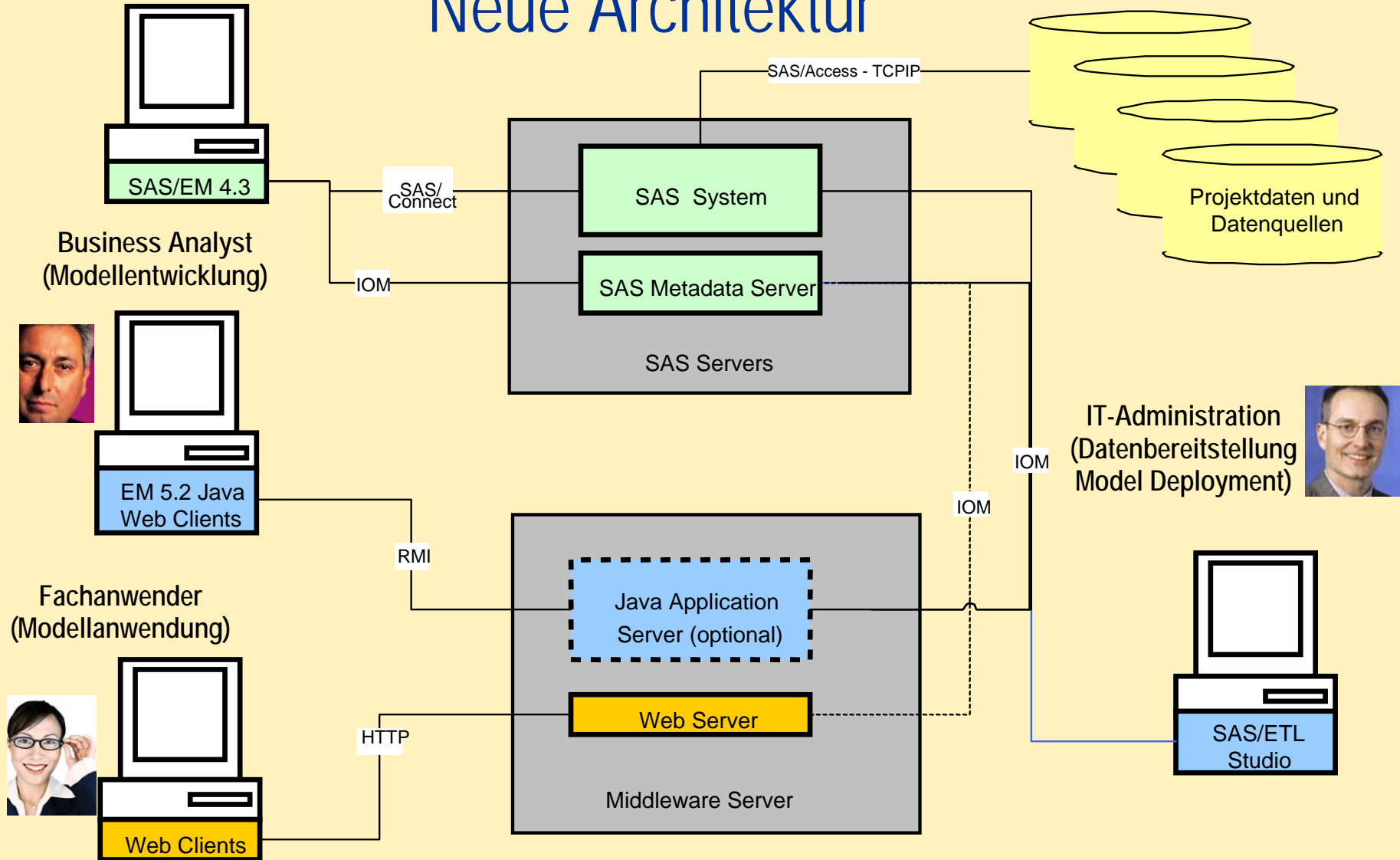


Fachanwender  
(Modellanwendung)

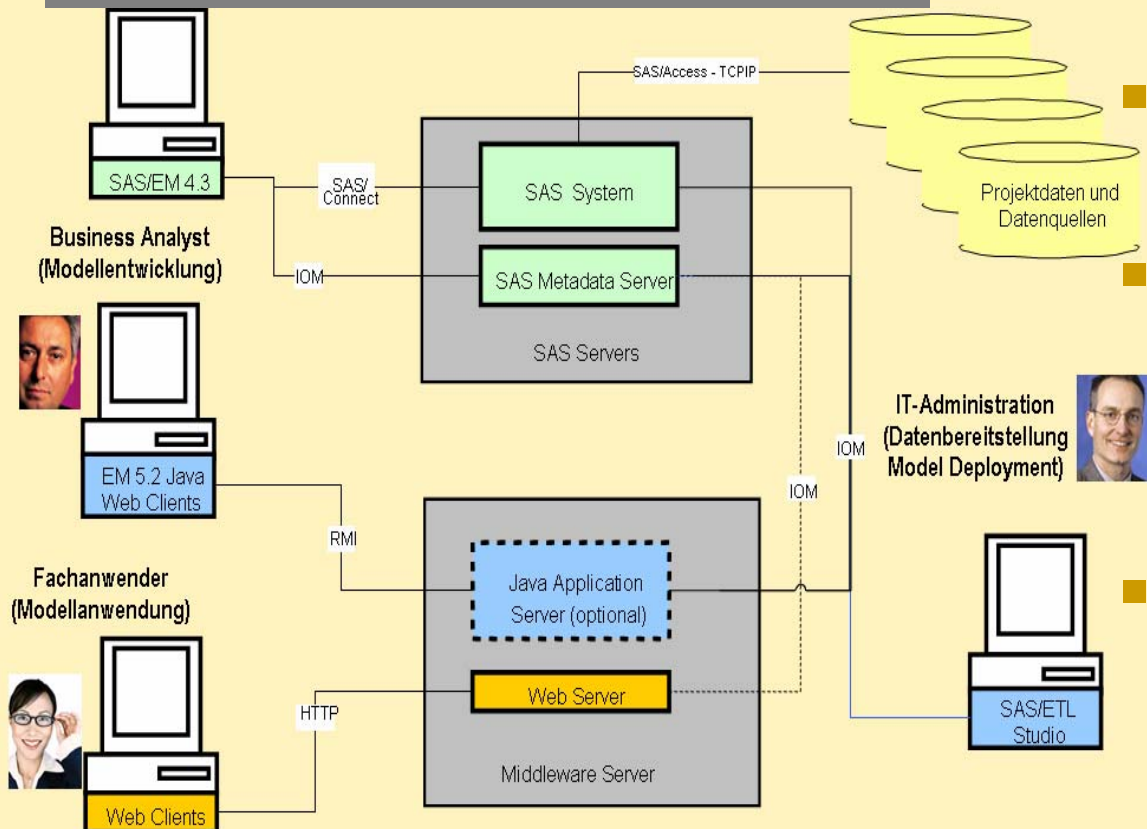




# Neue Architektur



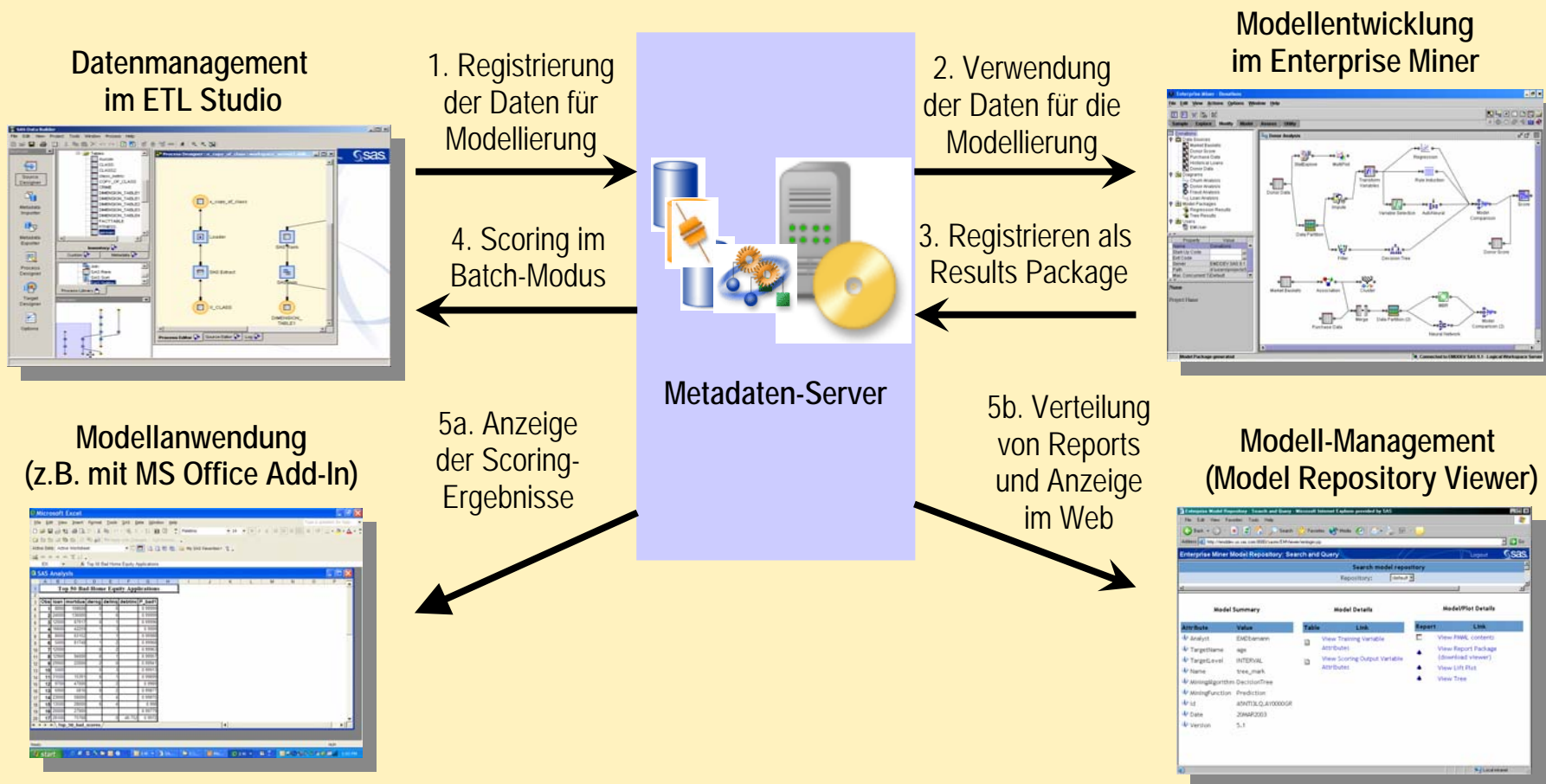
# Skalierbare N-Tier-Architektur



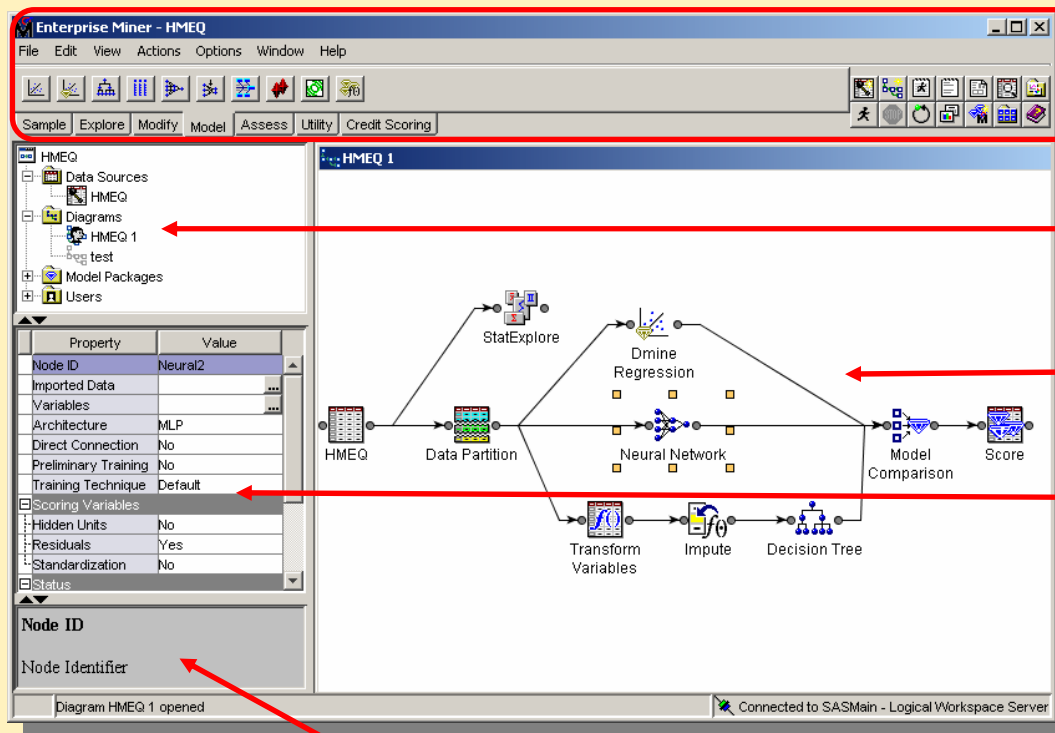
- Neu gestaltete Java-basierte Oberfläche
- Unterstützt n Tier Architektur mit Java Application Server
- Performance & Skalierbarkeit
- Ermöglicht Web-basiertes gemeinsames Arbeiten an Data Mining Projekten
- Durchgehendes Metadatenkonzept



# Einheitliches Metadaten-Konzept von SAS



# Neue Java-basierte Benutzeroberfläche



Neue gestaltete Menü- und Symbolleisten

Projektverwaltung und Datenquellen-Definition

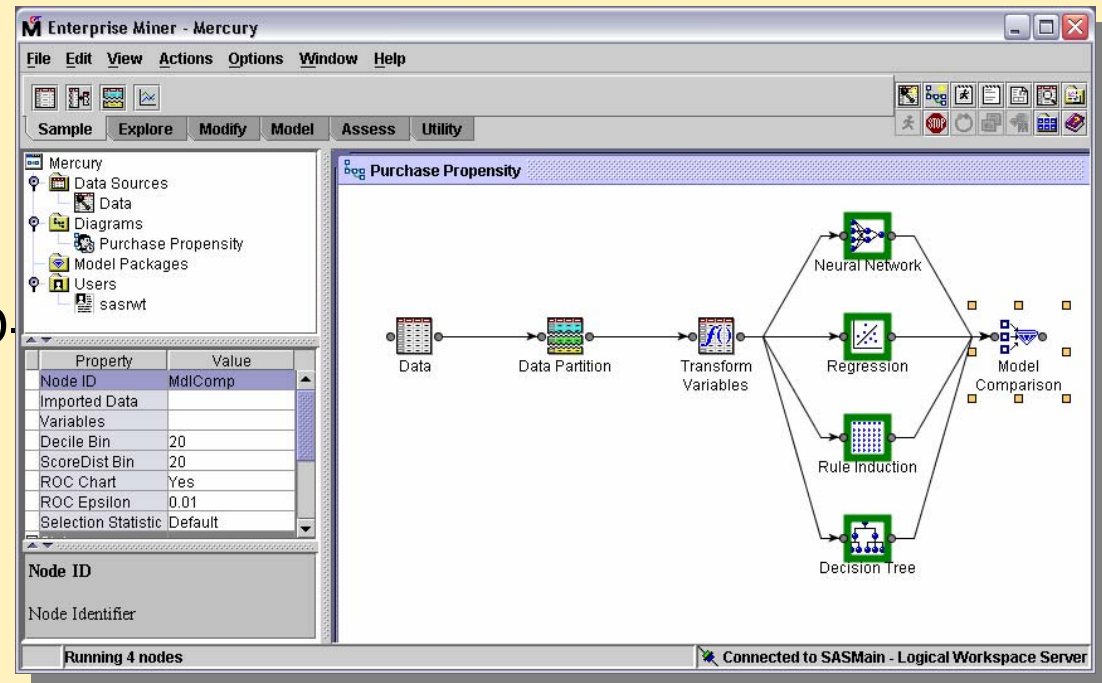
Prozessflussdiagramm als Arbeitsbereich

Eigenschaften für einzelne Objekte (Knoten, Daten...)

Kurzer Hilfetext zur Erläuterung

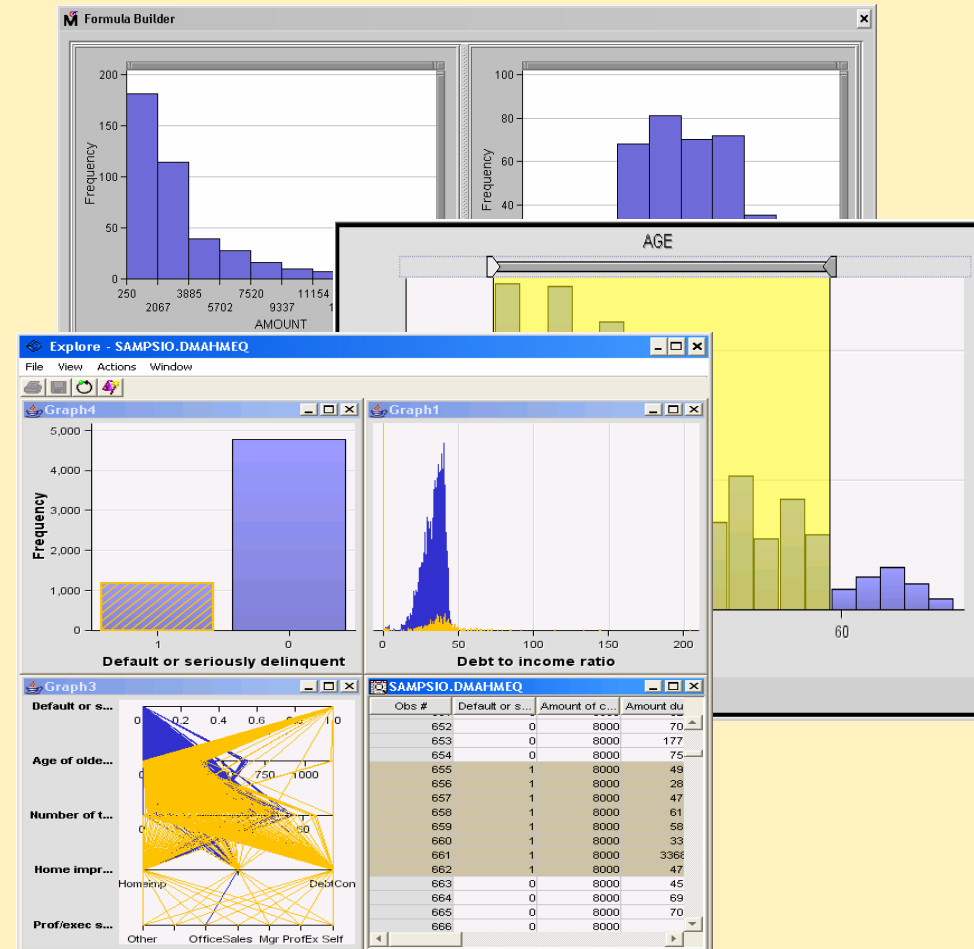
# Verbesserte Performance und Skalierbarkeit

- Server-seitige Projektverwaltung
- Paralleles Verarbeiten mehrerer Modelle
- Unterstützung von GRID-Computing
- Algorithmen mit Multi-Threading
- Asynchrones Modell-Training
- Sauberer Abbruch der Verarbeitung
- Disconnect/Reconnect aus lfd. Verarbeitung

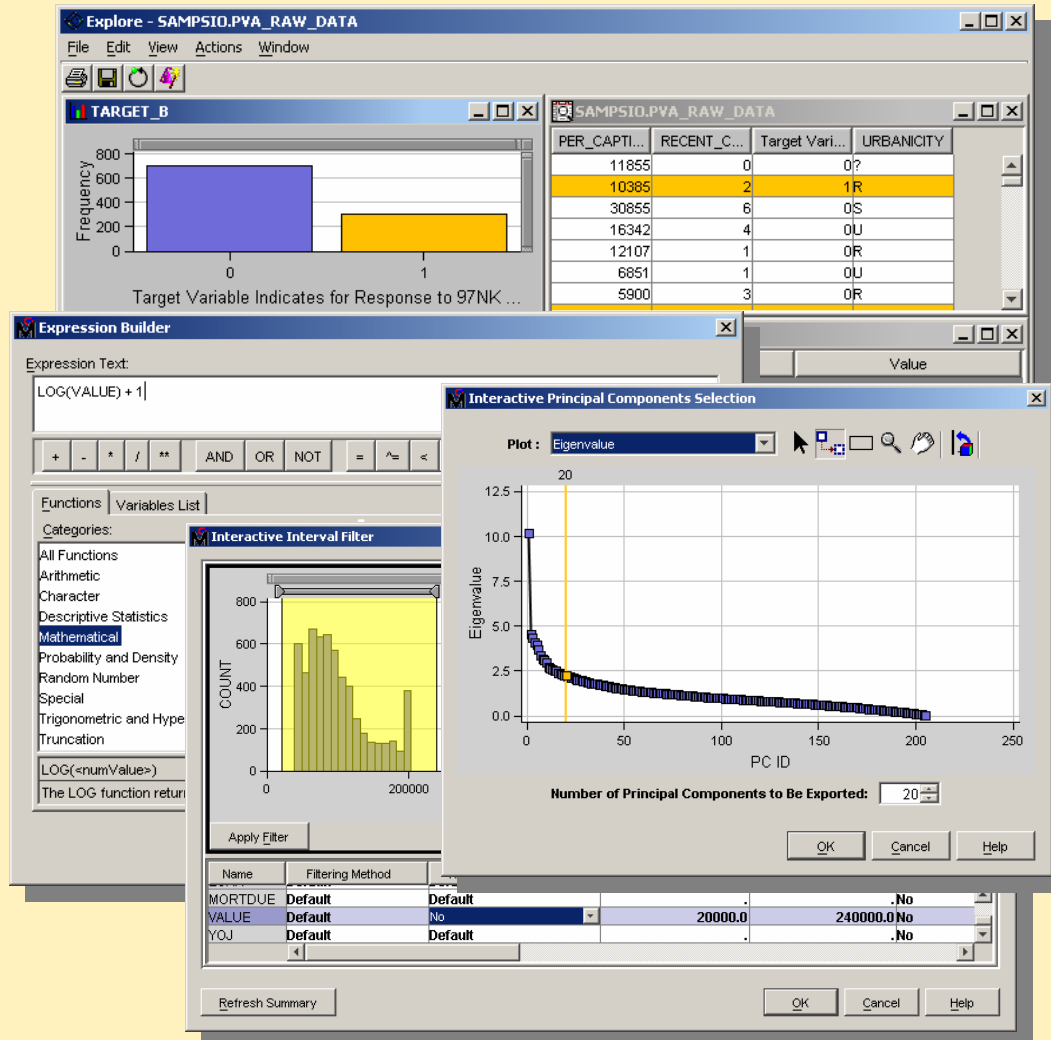


# Datenexploration und Datenmodifikation

- Partitionieren von Daten (Training/Validation/Test)
- Deskriptive Statistik und explorative Grafiken
- Werteersetzung, Behandlung fehlender Werte
- Interaktive Variablentransformationen
- Variablenauswahl und Datenfilter
- Zusammenführen von Dateien und Stichproben
- Definieren von Entscheidungsprofilen für Zielgrößen
- ...



# Interaktive Steuerungsmöglichkeiten

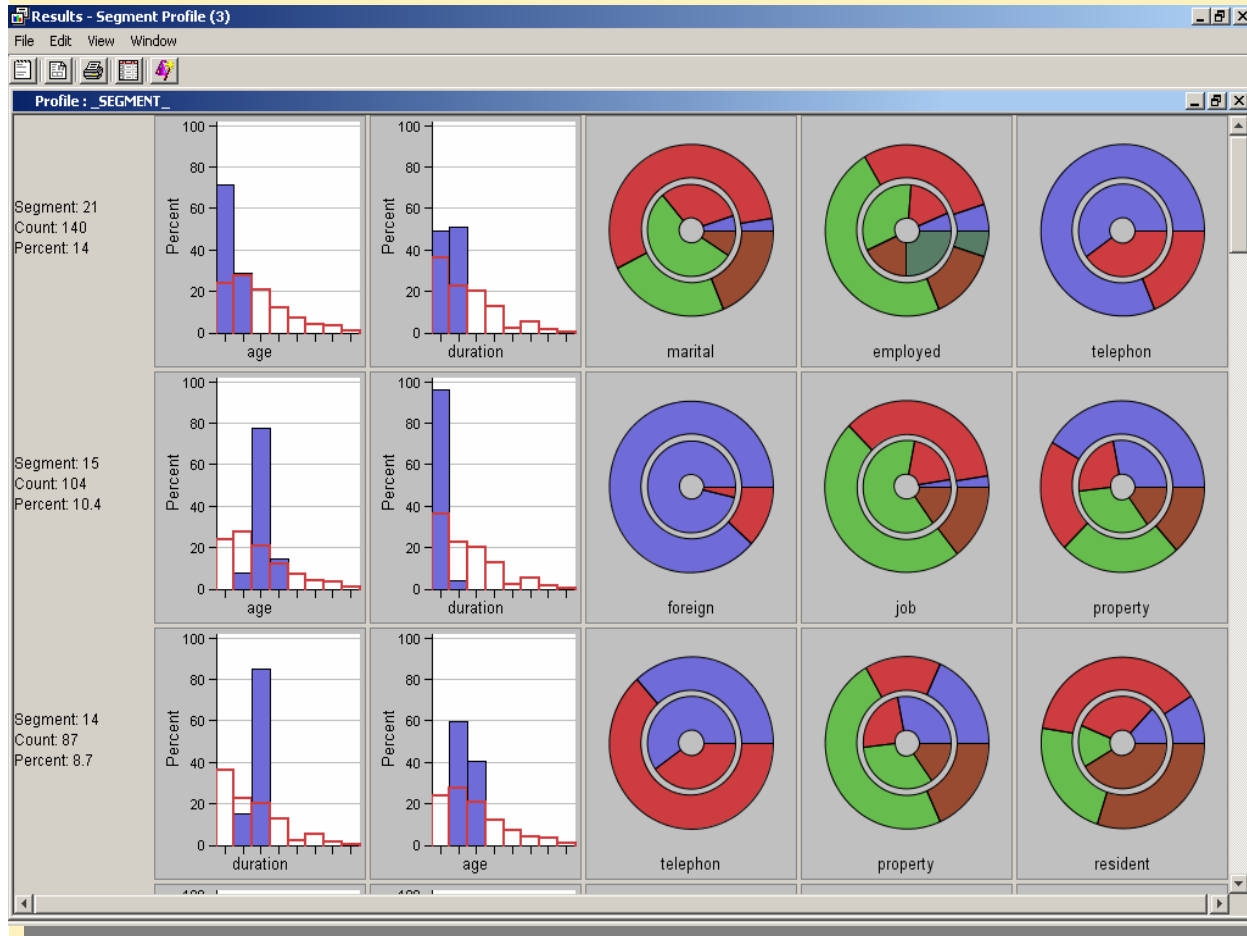


The screenshot displays several SAS windows illustrating interactive capabilities:

- Explore - SAMPPIO.PVA\_RAW\_DATA:** Shows a bar chart for 'TARGET\_B' and a data table with columns: PER\_CAPTI..., RECENT\_C..., Target Vari..., and URBANICITY.
- Expression Builder:** Features an 'Expression Text' field containing `LOG(VALUE) + 1` and a 'Functions' list including Arithmetic, Character, Descriptive Statistics, Mathematical, Probability and Density, Random Number, Special, Trigonometric and Hype, and Truncation.
- Interactive Interval Filter:** Displays a histogram with a yellow selection box over the data.
- Interactive Principal Components Selection:** Shows a plot of 'Eigenvalue' vs 'PC ID' with a vertical line at PC ID 20. Below the plot, it indicates 'Number of Principal Components to Be Exported: 20'.

- Dynamisch verknüpfte Grafik- und Tabelleninhalte
- Flexible Grafiken über neuen Java Graphics Wizard
- Editor für benutzerdefinierte Variablentransformationen
- Benutzerdefinierte Bereichs-  
selektion im Filter-Knoten
- Auswahl beizubehaltender  
Hauptkomponenten in PCA

# Neuer Segment Profiler für Visualisierung von Clusterergebnissen



Ermöglicht leichtere Interpretation von Segmentierungen

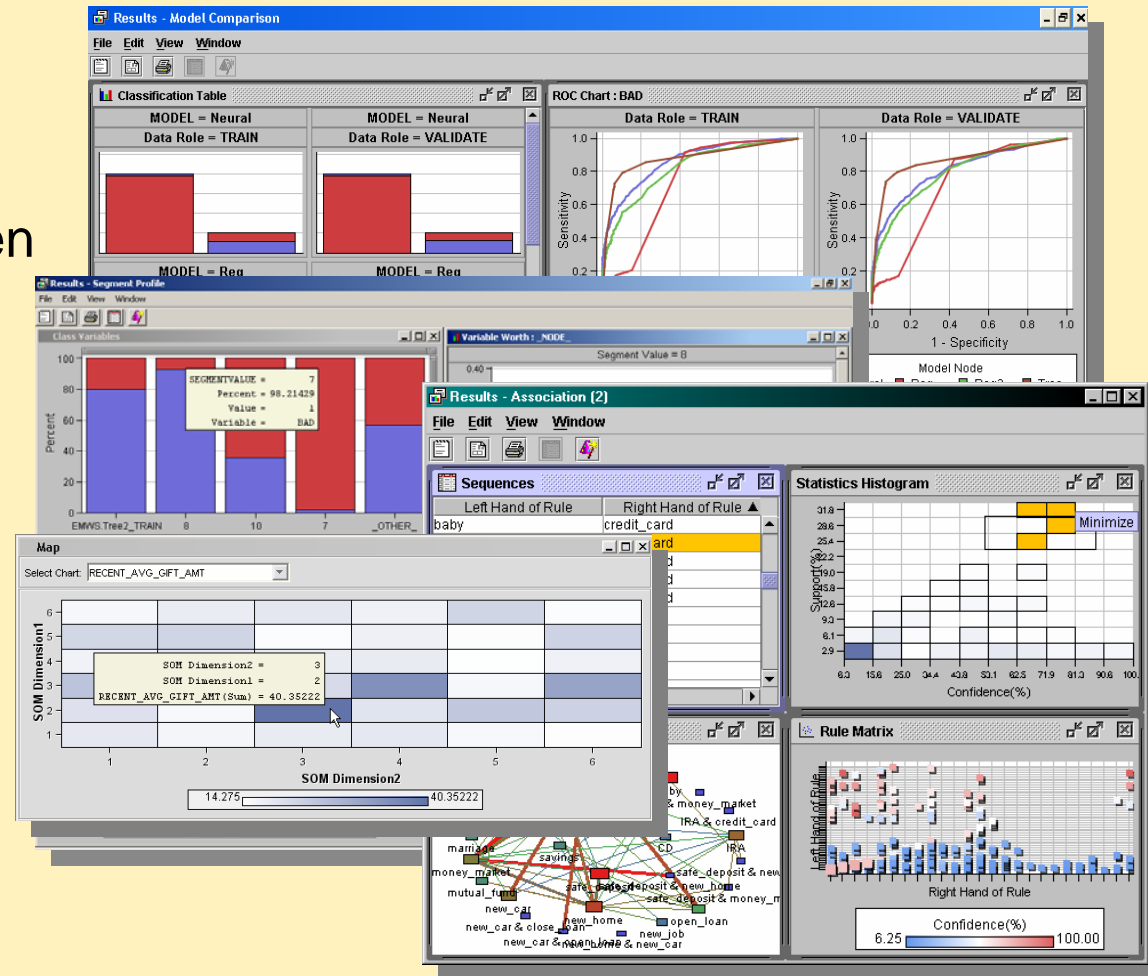
Ermittlung der Wichtigkeit (Worth) von Merkmalen für einzelne Segmente

Merkmalsprofile der jeweiligen Segmente im Vergleich zum Total

- Histogramme für kontinuierliche Variablen
- Kreisdiagramme für kategoriale Variablen

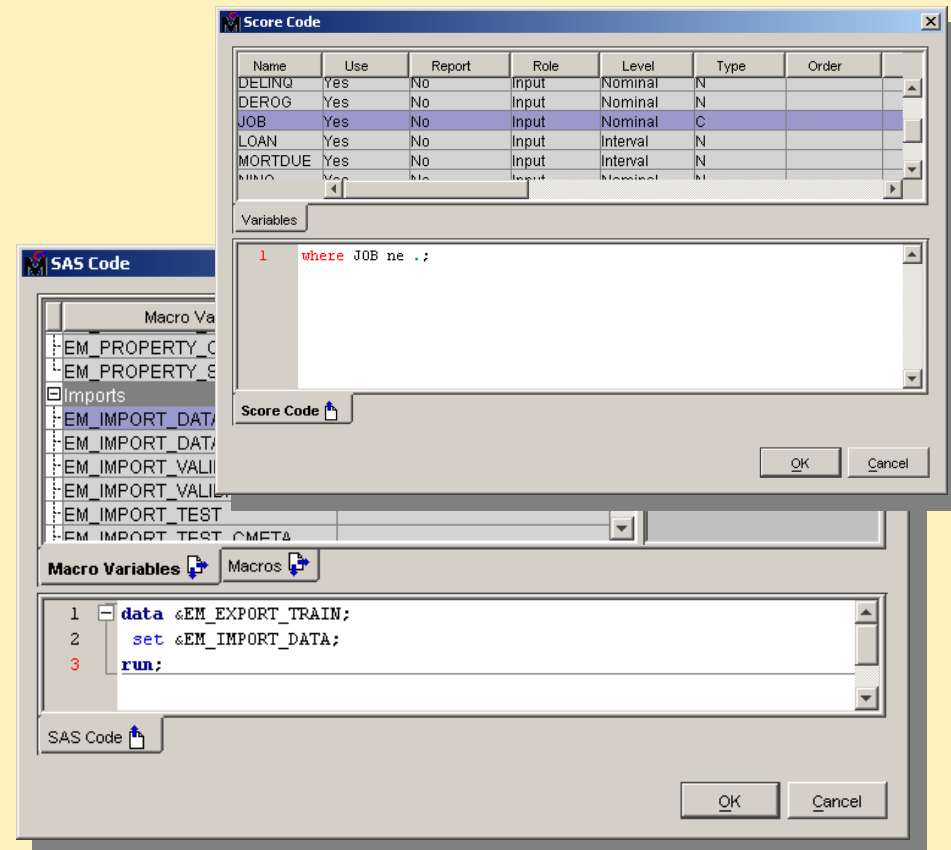
# Modellbewertung und Vergleich

- Lift-Diagramme
- ROC-Kurven
- Modellgüte-Statistiken
- Iterationsverläufe
- Segmentprofile
- Link Graph
- Regelmatrix
- Topologische Karte
- ...



# Erweiterbarkeit durch eigenen SAS Code

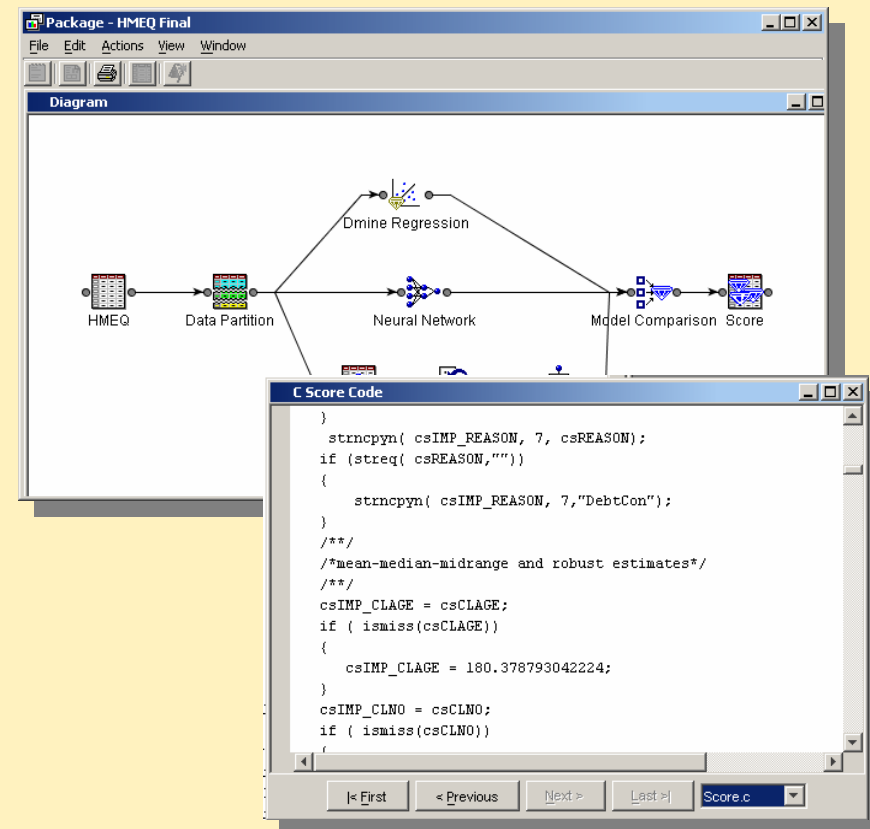
- Code-Knoten für Zugriff auf SAS Prozeduren
- Referenzieren von Makrovariablen per Drag & Drop
- Ermöglicht benutzerdefinierte Funktionalität für
  - Reporting
  - Datentransformationen
  - Analytische Verfahren
- Generieren eigener Knotenerweiterungen





# Score Code für Modellanwendung

- Export von Score Code nach SAS, C, Java und PMML
- Berücksichtigung sämtlicher Verarbeitungsschritte im Code
- Registrieren als Model Package im SAS Metadaten-Server
- Score-Code-Ausführung über Plug-In in SAS ETL Studio





# Agenda

Der Neue Text Miner 2.3



# Sprachunterstützung für 8 Sprachen

- Englisch
- Französisch
- Deutsch
- Niederländisch
- Italienisch
- Portugiesisch
- Spanisch
- Schwedisch



# Macros zur Textbereinigung

```
%textsyn(docds=EMWS.TEXT2_DOCUMENTS,
        termds= EMWS.TEXT2_TERMS,
        outds= EMWS.TEXT2_OUT_T,
        synds=sasuser.autosyns,
        dict=data.engdict,
        textvar=DESCRIPTION);
```

EXAMPLE1	EXAMPLE2	TERM	PARENT	CAT...	CHILDNDO...	NUMDOCS	MINSPEED
REPOSITIONED LEFT SIDE !!A!! C. _H		a.	a		1.0	97.0	-0.0
... DAMPER NOISE REPLACE SHOCK !!A		absorbe	absorber		1.0	5.0	4.0
... REPL. R. RR. SHOCK !!ABSORBR!!		absorbr	absorber		1.0	5.0	7.0
... THE LF FRT WHEN !!ACC!! LYBED THERADIO INOP !!ACC!! SOCK		acc	ac		2.0	20.0	12.0
!!ACCELERATOR!! ...	WHEN PRESSING !!ACCE	accelerator	acceleration		4.0	5.0	13.0
!!ACCELORATOR!! STICKS.		accelorator	accelerator		1.0	4.0	9.0
... SHIFT AND VIBRATES ON !!ACEL!! FRC		acel	accel		1.0	9.0	6.0
... REPL !!ACTUATER!! UNIT. RETEST.		actuater	actuator		1.0	6.0	12.0
REAR HATCH RATTLES _ !!ADJD!! HATC		adjd	adj		1.0	38.0	10.0
... UP MOLYKOATED FRONT PADS !!ADJS... DROVE SEVERAL TIMES		adjsuted	adjusted		2.0	23.0	6.0
REPLACE SEAT !!ADJUSTER!!		adjuster	adjusted		1.0	23.0	12.0
... 50MPH TIRES OUT B !!ALANCE BALAN		al	all		1.0	128.0	12.0
... NEC TO SUBLET FOR !!ALIG!!	... SUBLET FOR 4 WHEEL	alig	align		3.0	154.0	8.0
STEERING WHEEL NOT CENTERED, !!A		aligh	alig		1.0	3.0	8.0
... _ _ _ NOT !!ALIGNED!! AND END CA		alighned	aligned		1.0	7.0	7.0
!!ALRAM!! REMOTE INOP, REPLACED TP		alram	alarm		1.0	7.0	10.0
LEAK FROM THE !!ALXE!! BOOT		alxe	axle		1.0	5.0	12.0
SPLR GRNSH CMNG !!APRT!! _ ROK SI		aprt	apart		1.0	12.0	12.0
WHEEL !!BALAN!!ICE4 WHEN DRIVING 6		balan	balance		1.0	15.0	14.0
... _ DR. , RAD. !!BALANC!!		balanc	balance		1.0	15.0	5.0
... TO TRAILER HITCH WIRING, !!BARE!! V		bare	bar		1.0	3.0	10.0
... CAR WONT START CK !!BAT!! FAILED		bat	batt		1.0	4.0	8.0

# Nutzungsmöglichkeiten von Perl Regular Expressions über den Code Knoten

## Perl Regular Expression Quick Reference 1.03

N.B.: this quick reference is just that - some of the explanations have been simplified. For the authoritative documentation, see the latest edition of *Programming Perl* or `perldoc perlre`.

### Specific characters:

<code>\t</code>	A tab character
<code>\n</code>	A newline character (OS neutral)
<code>\r</code>	A carriage return character
<code>\f</code>	A form feed character
<code>\cX</code>	Control character CTRL-X
<code>\NNN</code>	Octal code for character NNN

### Metacharacters:

The following 12 characters need to be escaped with a backslash - "\ " - because by default, they mean something special.

`\ | ( ) [ { ^ $ * + ? .`

<code>.</code>	Match one character
<code> </code>	Alternation
<code>( )</code>	Group and capture
<code>[ ]</code>	Define character class
<code>\</code>	Modify the meaning of the next char.

### Anchors:

<code>^</code>	Match at the beginning of a string (or line)
<code>\$</code>	Match at the end of a string (or line)
<code>\b</code>	Match at a 'word' boundary
<code>\B</code>	Match at not a 'word' boundary

### Quantifiers:

These quantifiers apply to the preceding *atom*.

<code>*</code>	Match 0 or more times
<code>+</code>	Match 1 or more times
<code>?</code>	Match 0 or 1 times
<code>{N}</code>	Match exactly N times
<code>{N,}</code>	Match at least N times
<code>{N,M}</code>	Match at least N but not more than M times

By default, quantifiers are "greedy". They attempt to match as many characters as possible. In order to make them match as few characters as possible, follow them with a question mark "?".

### Character class metacharacters:

<code>^</code>	If the first character of a class, negates that class
<code>-</code>	Unless first or last character of a class, used for a range

### Character class shortcuts:

<code>\d</code>	<code>[0-9]</code>	A digit
<code>\D</code>	<code>[^0-9]</code>	A non-digit
<code>\s</code>	<code>[\t\n\r\x\v]</code>	A whitespace char.
<code>\S</code>	<code>[^\t\n\r\x\v]</code>	A non-whitespace char.
<code>\w</code>	<code>[a-zA-Z0-9_]</code>	A 'word' char.
<code>\W</code>	<code>[^a-zA-Z0-9_]</code>	A 'non-word' char.

These shortcuts can be used either on their own, or within a character class.

### Metaquote & case translations:

<code>\Q</code>	Quote (de-meta) characters until \E
<code>\U</code>	Uppercase characters until \E
<code>\L</code>	Lowercase characters until \E

### Special variables:

<code>\$`</code>	The characters to the left of the match
<code>\$&amp;</code>	The characters that matched
<code>\$'</code>	The characters to the right of the match
<code>\N</code>	The characters captured by the N <sup>th</sup> set of parentheses (if on the match side)
<code>\$N</code>	The characters captured by the N <sup>th</sup> set of parentheses (if not on the match side)

### Modifiers:

These modifiers apply to the entire pattern

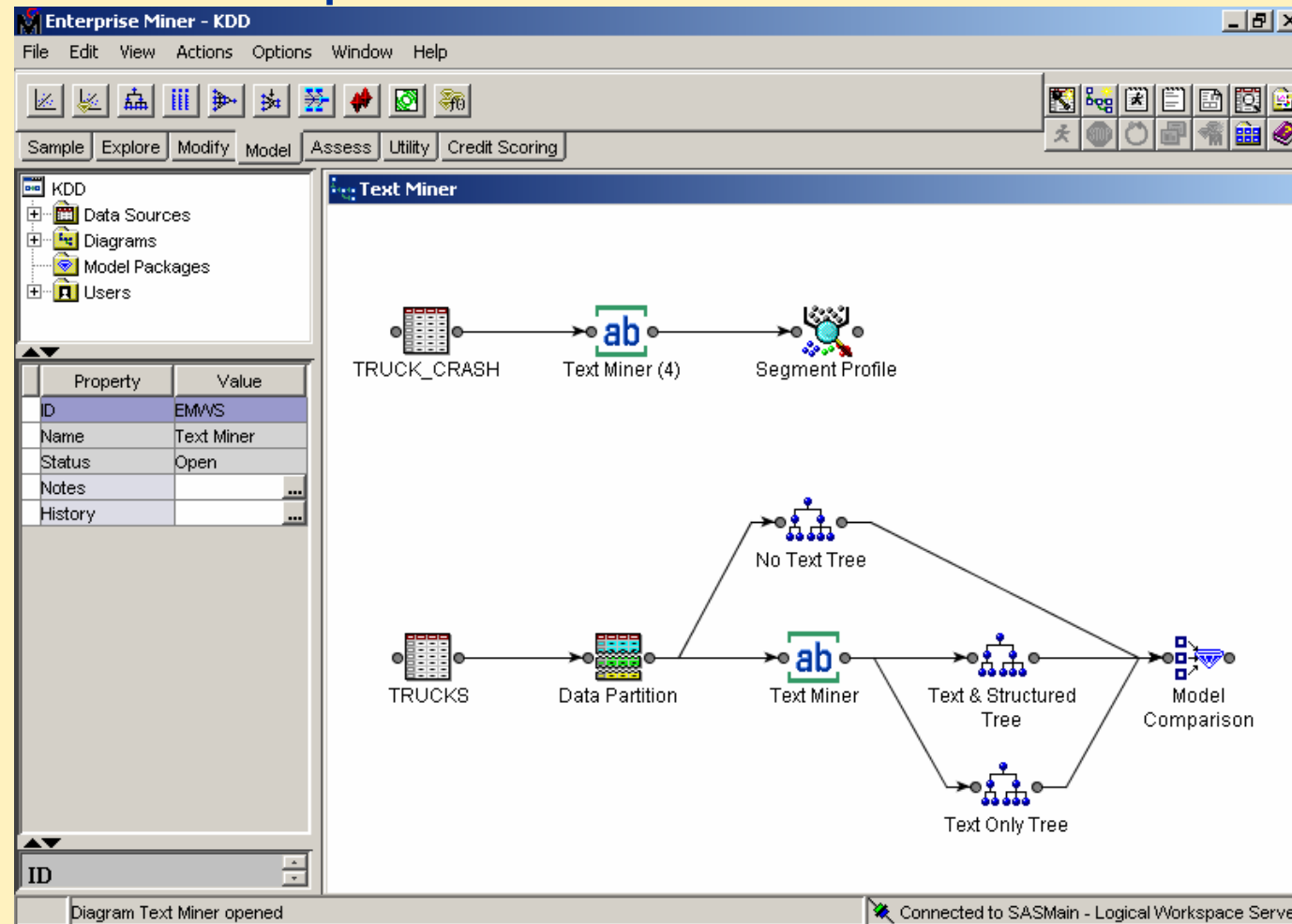
<code>/i</code>	Ignore case
<code>/g</code>	Match globally (all)
<code>/m</code>	Let ^ and \$ match next to embedded \n
<code>/s</code>	Let . match \n
<code>/x</code>	Ignore most whitespace and allow comments
<code>/e</code>	Evaluate right hand side of s/// as an expression

All except /e apply to both `m//` and `s///`.

### Binding operators:

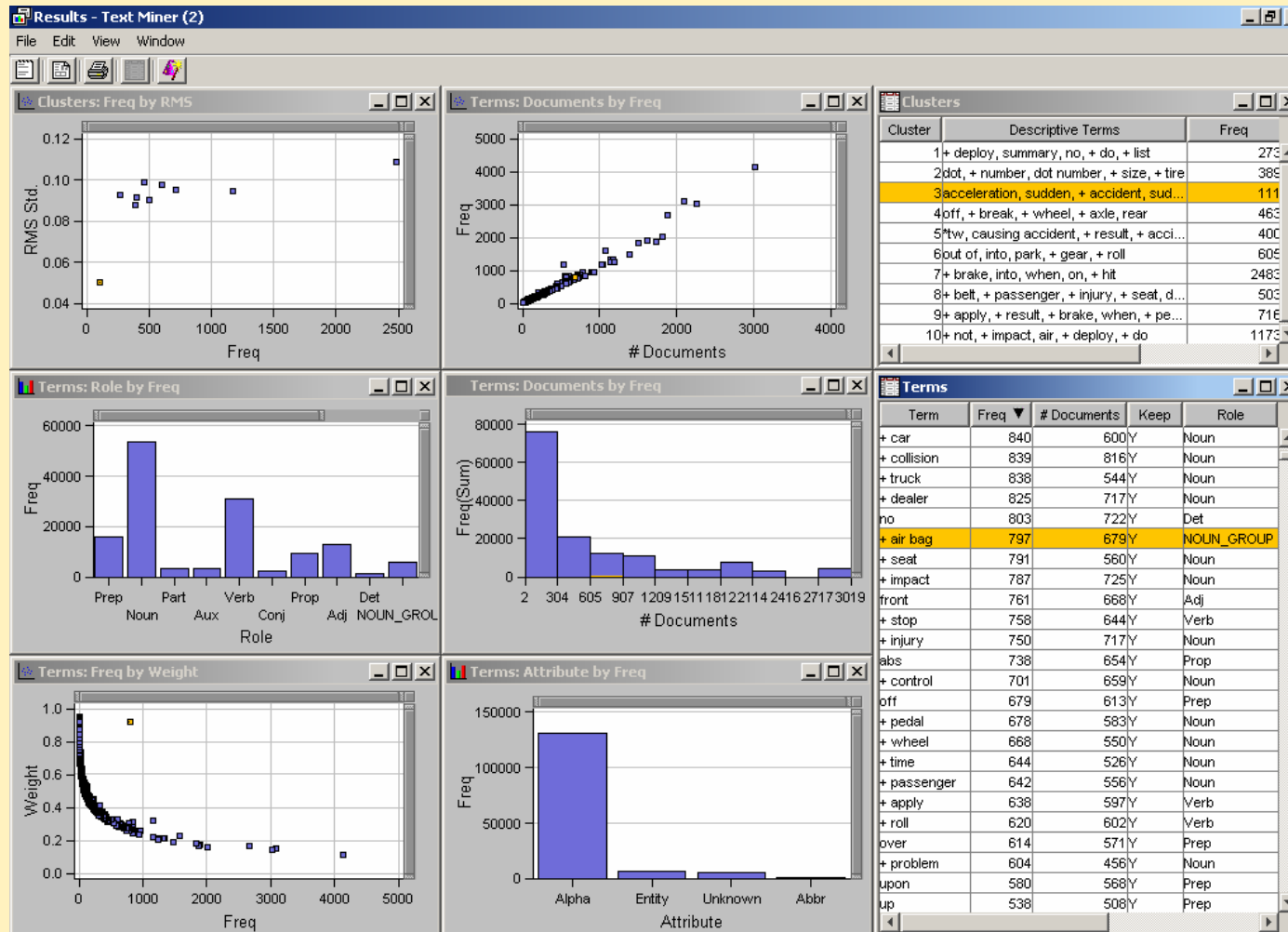
<code>=~</code>	True if the regex matches
<code>!~</code>	True if the regex doesn't match

# Integration in das Prozessflussdiagramm der Enterprise Miner Oberfläche





# Text Miner Static Results Browser





# Interaktives Ergebnisfenster

Clusters  
Fenster



Terms  
Fenster



Results - Interactive

File Edit Tools View Window

Clusters			
#	Descriptive Terms	Freq	Percent
2	dot, + number, dot number, + size, + tire	389	0.054665
3	acceleration, sudden, + accident, sudden acceleration, causing accident	111	0.015598
4	off, + break, + wheel, + axle, rear	463	0.065064
5	*tw, causing accident, + result, + accident, + fail	400	0.056211
6	out of, into, park, + gear, + roll	605	0.085019

Terms				
Term	Freq	# Documents	Keep	Role
the	9983.0	2909.0	<input type="checkbox"/>	Det
+ be	6923.0	3039.0	<input type="checkbox"/>	Verb
and	6137.0	3403.0	<input type="checkbox"/>	Conj
+ vehicle	6009.0	3641.0	<input type="checkbox"/>	Noun
+ a	5308.0	3358.0	<input type="checkbox"/>	Det
in	4128.0	3019.0	<input checked="" type="checkbox"/>	Prep
to	3728.0	2463.0	<input type="checkbox"/>	Part
+ brake	3086.0	2091.0	<input checked="" type="checkbox"/>	Noun
+ not	3027.0	2265.0	<input checked="" type="checkbox"/>	Part
on	2669.0	1885.0	<input checked="" type="checkbox"/>	Prep
of	2644.0	1805.0	<input type="checkbox"/>	Prep
*ak	2583.0	2580.0	<input type="checkbox"/>	Prop
+ have	2383.0	1421.0	<input type="checkbox"/>	Verb
+ accident	2005.0	1817.0	<input checked="" type="checkbox"/>	Noun

Documents

SUMMARY

NO DEPLOYMENT OF DRIVER'S AIRBAG DURING ACCIDENT

NO SUMMARY LISTED FOR ABOVE VEHICLE. \*AK( DOT M

SEPTEMBER 11, 2000. NO SAFETY DEFECT MENTIONE DI

NO SUMMARY LISTED FOR ABOVE VEHICLE. \*AK

AIR BAG DID NOT DEPLOY. TT

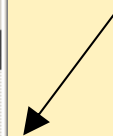
NO DASHLIGHTS, CAUSING POOR VISIBILITY FOR READI

NO SUMMARY LISTED FOR THIS VEHICLE. \*AK

NO DEPLOYMENT OF AIR BAG DURING ACCIDENT RESULTE

AIR BAG DEPLOYMENT RESULTED IN MINOR BURNING OF  
126 CHANGE: ACCID = 0, INJURED = 0. AK

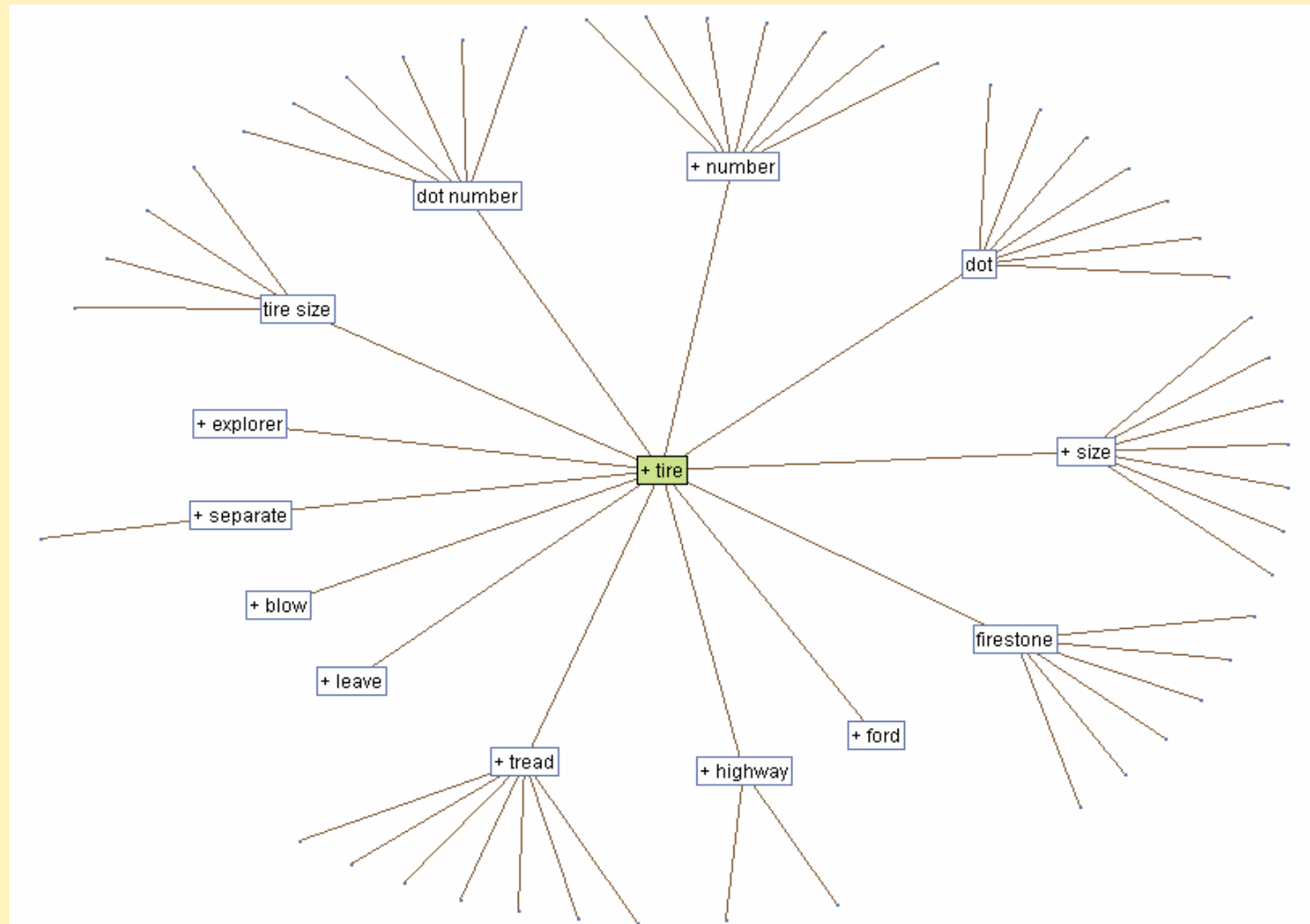
Dokument  
Fenster





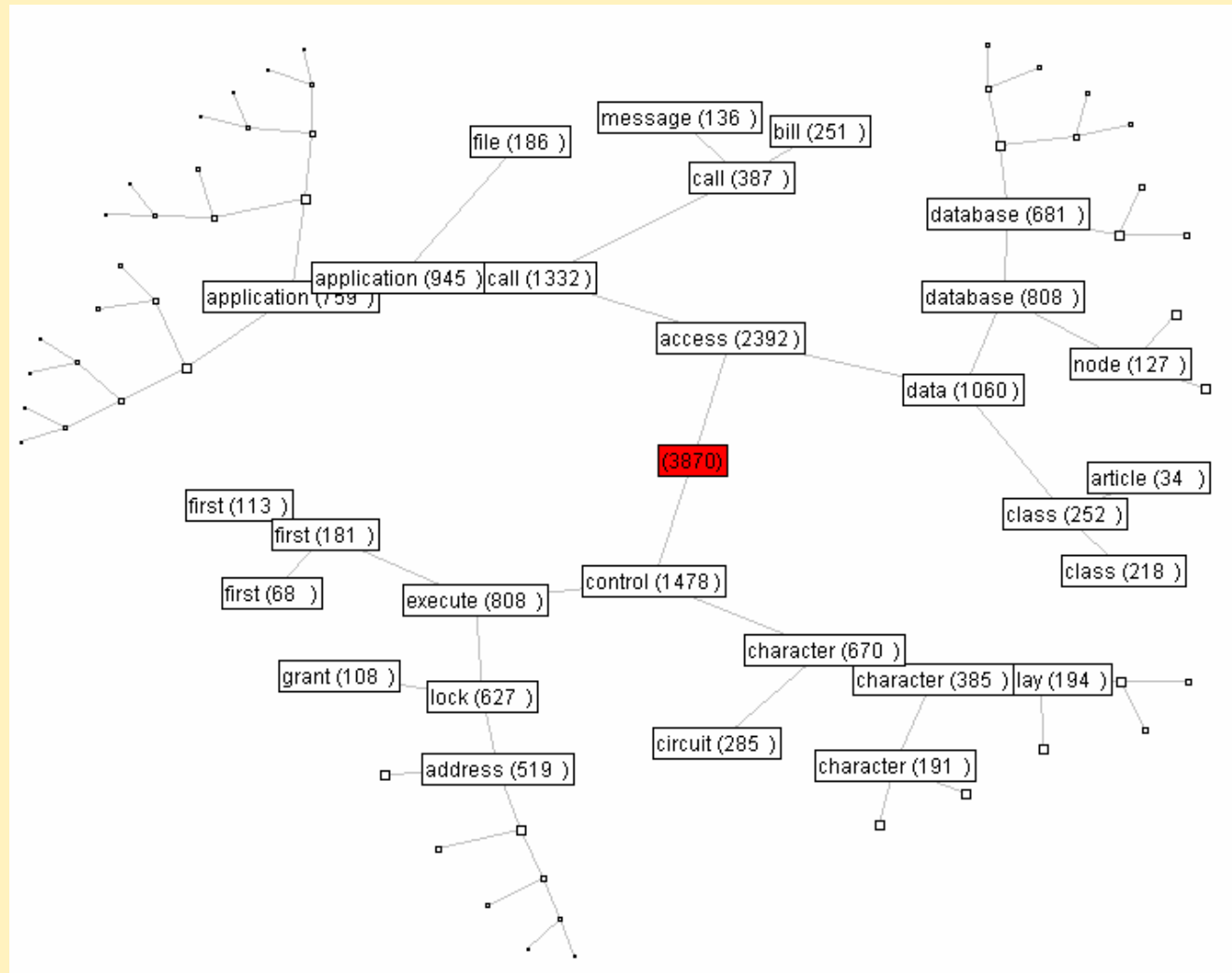


# Konzept Linking: Gefundene Beziehung zwischen dem Begriff "tire" und anderen Begriffen



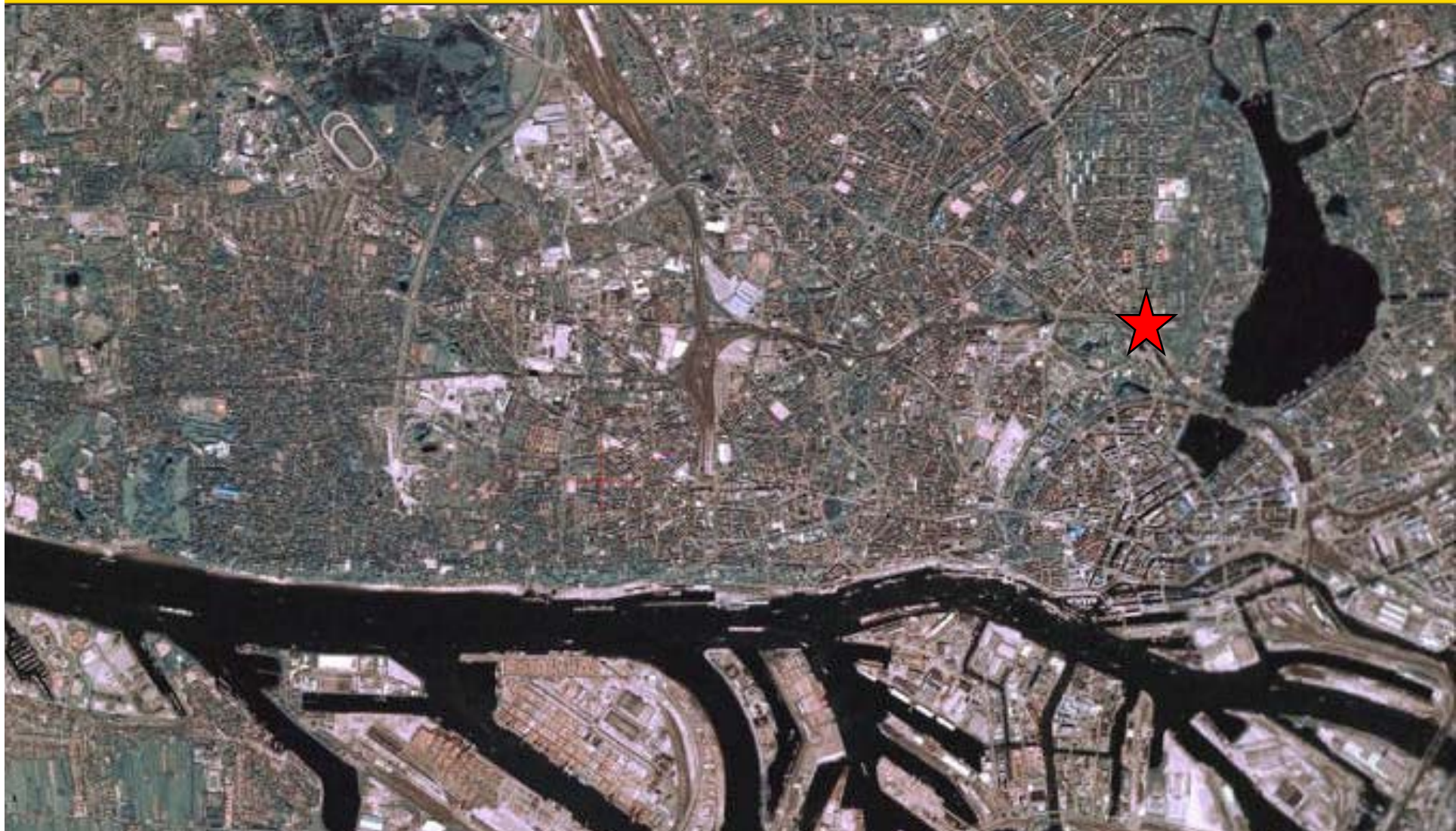


# Hierarchisches Dokument Clustering



# Zusammenfassung

- Enterprise Miner 5.2
  - Einheitliches Metadaten-Konzept
  - Score Code in SAS, C, Java und PMML
  - Bequeme Verwaltung (Model Packages und Model Repository Viewer)
  - Batch Jobs zur Automatisierung von Data Mining-Abläufen
  - Viele neue Analytische Funktionalitäten
  - Performance
- Text Miner 2.3
  - Nahtlose Integration in den Enterprise Miner Prozessfluss
  - Unterstützung von 8 Sprachen
  - Textbereinigung, Perl Regular Expressions
  - Konzept Linking
  - Hierarchisches Dokument Clustering mit Java Tree



# Diskussion