



Florian Zipser
Humboldt-Universität zu Berlin

SaltNPepper und das Formatpluriversum

LAUDATIO Workshop

2014-10-07



- Linguistische Daten und Phänomene erfordern viele Annotationsarten

Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Musikcafé	.
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Acc.Sg.Neut	–
ART	NN	APPR	NE	VMFIN	ART	NN	\$.

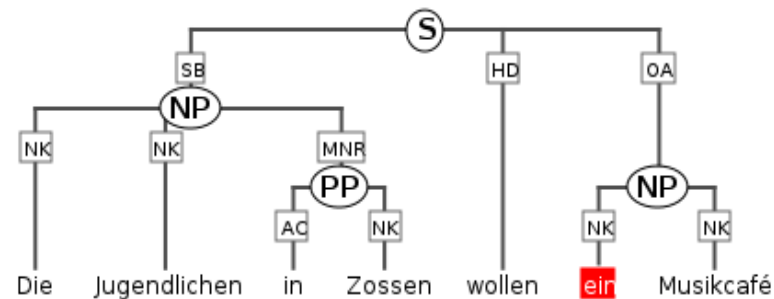
Morphologie



- Linguistische Daten und Phänomene erfordern viele Annotationsarten

Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Mu	.
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Ac	.
ART	NN	APPR	NE	VMFIN	ART	NN	.

Morphologie



Syntax



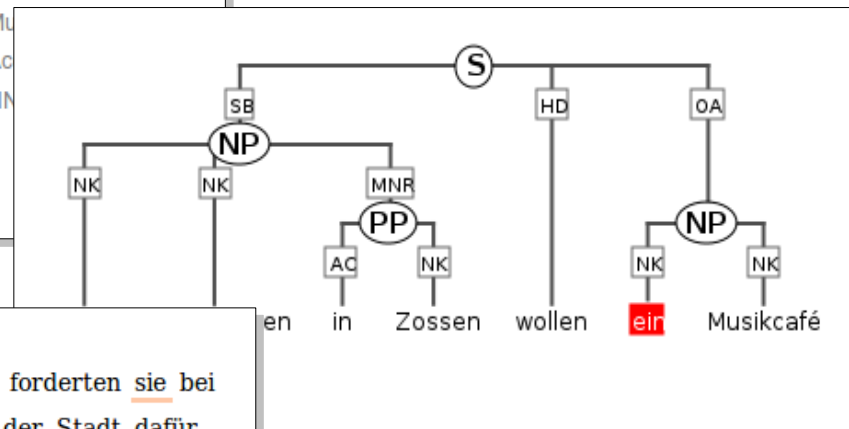
- Linguistische Daten und Phänomene erfordern viele Annotationsarten

Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Mu	
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Ac	
ART	NN	APPR	NE	VMFIN	ART	NN	

Morphologie

Feigenblatt Die Jugendlichen in Zossen wollen ein Musikcafé . Das forderten sie bei der ersten Zossener Runde am Dienstagabend . Dass die Politiker der Stadt dafür

Koreferenz



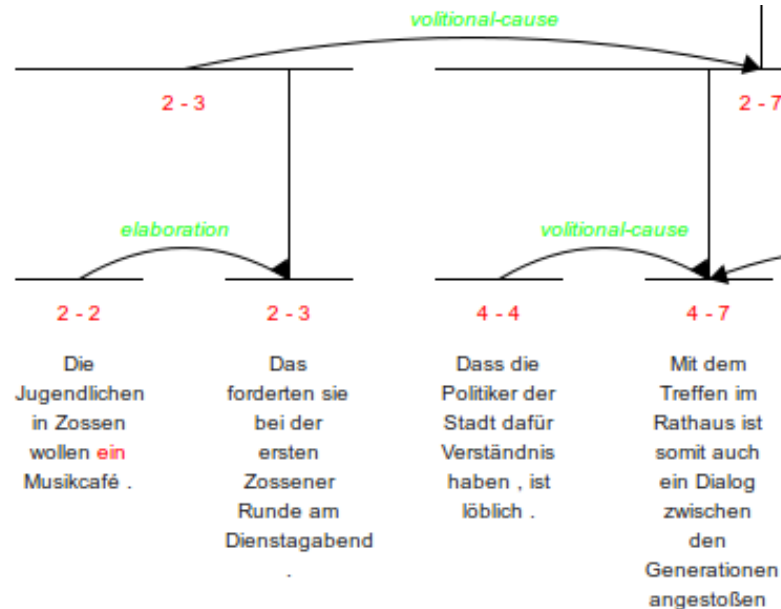


- Linguistische Daten und Phänomene erfordern viele Annotationsarten

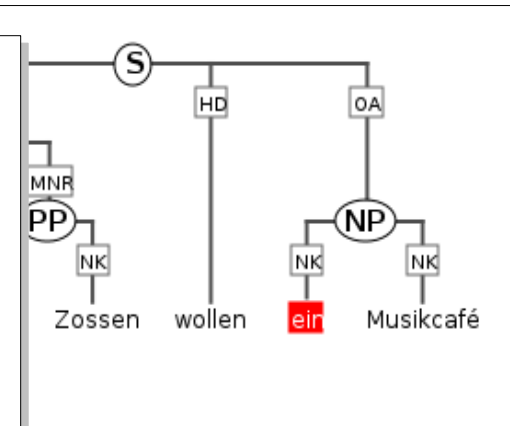
Die Jugendlichen in Zossen wollen ein Musikcafé .
der jugendliche in Zossen wollen ein Mu
Nom.Pl.* Nom.Pl.*
ART NN

Morphologie

Feigenblatt Die Jugend
der ersten Zossener R
Koreferenz



Rhetorische Strukturen



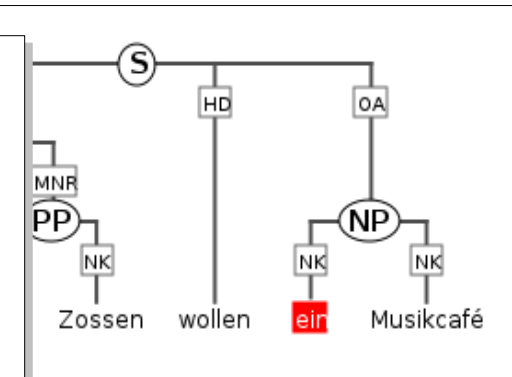
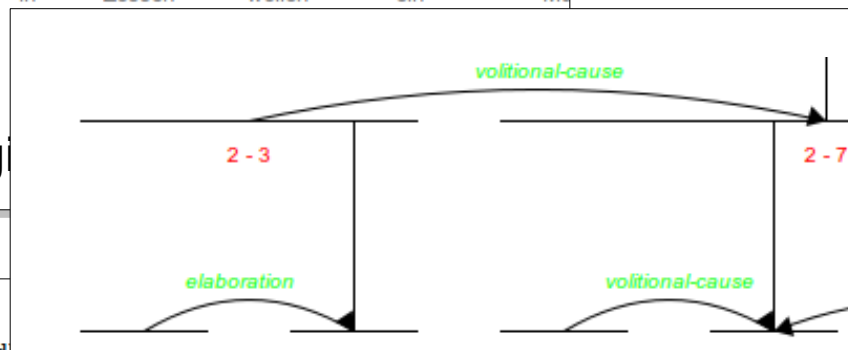


- Linguistische Daten und Phänomene erfordern viele Annotationsarten

Die Jugendlichen in Zossen wollen ein Musikcafé .
der jugendliche in Zossen wollen ein Mu
Nom.Pl.* Nom.Pl.*
ART NN

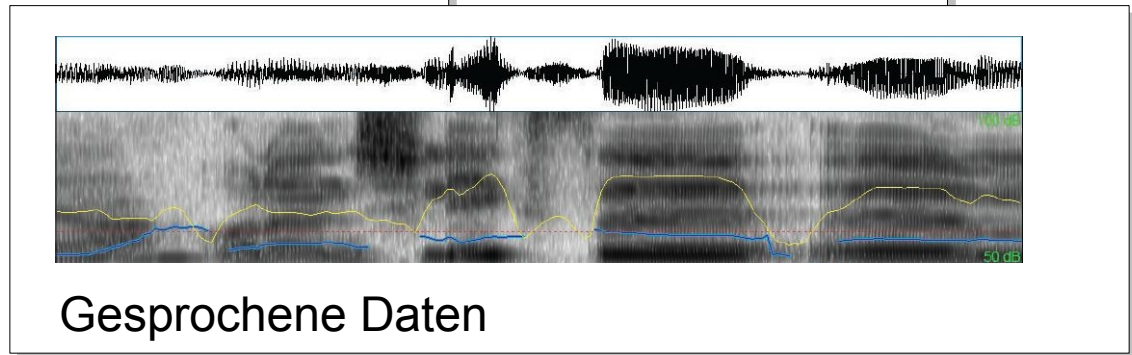
Morphologie

Feigenblatt Die Jugend
der ersten Zossener R
Koreferenz



2 - 2 2 - 3

Die Jugendlichen in Zossen wollen ein Musikcafé .
Das forderten sie bei der ersten Zossener Runde am Dienstagabend



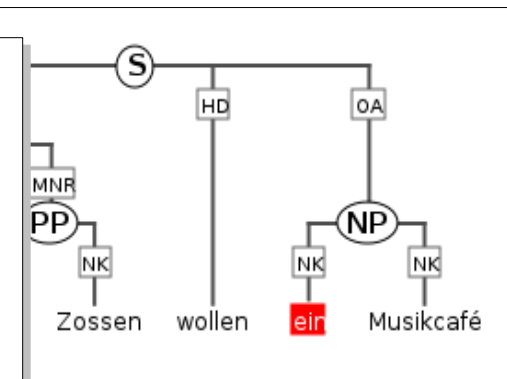
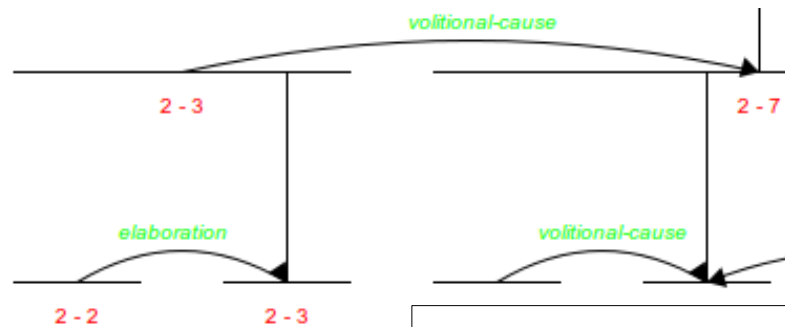
Rhetorische Strukturen



- Linguistische Daten und Phänomene erfordern viele Annotationsarten

Die Jugendlichen in Zossen wollen ein Musikcafé .
der jugendliche in Zossen wollen ein Mu
Nom.Pl.* Nom.Pl.*
ART NN

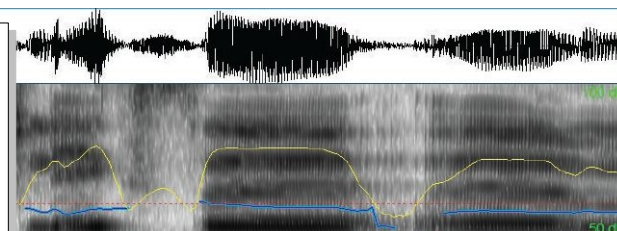
Morphologie



Feigenblatt Die Jugend
der ersten Zossen R

Instructee_dipl									okay	
Instructor_dipl	die	hast	du	dann	rechts	von	dir	genau	und	da

Dialoge

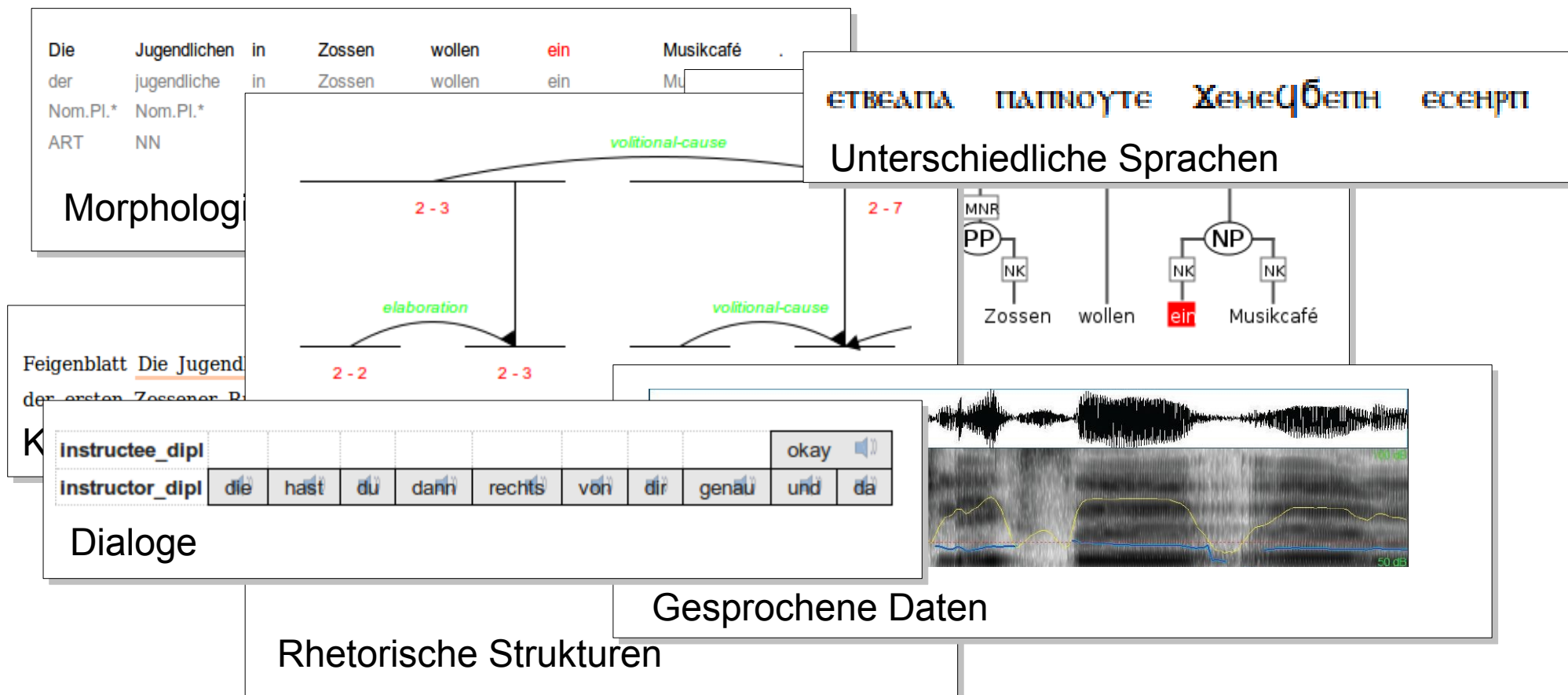


Gesprochene Daten

Rhetorische Strukturen

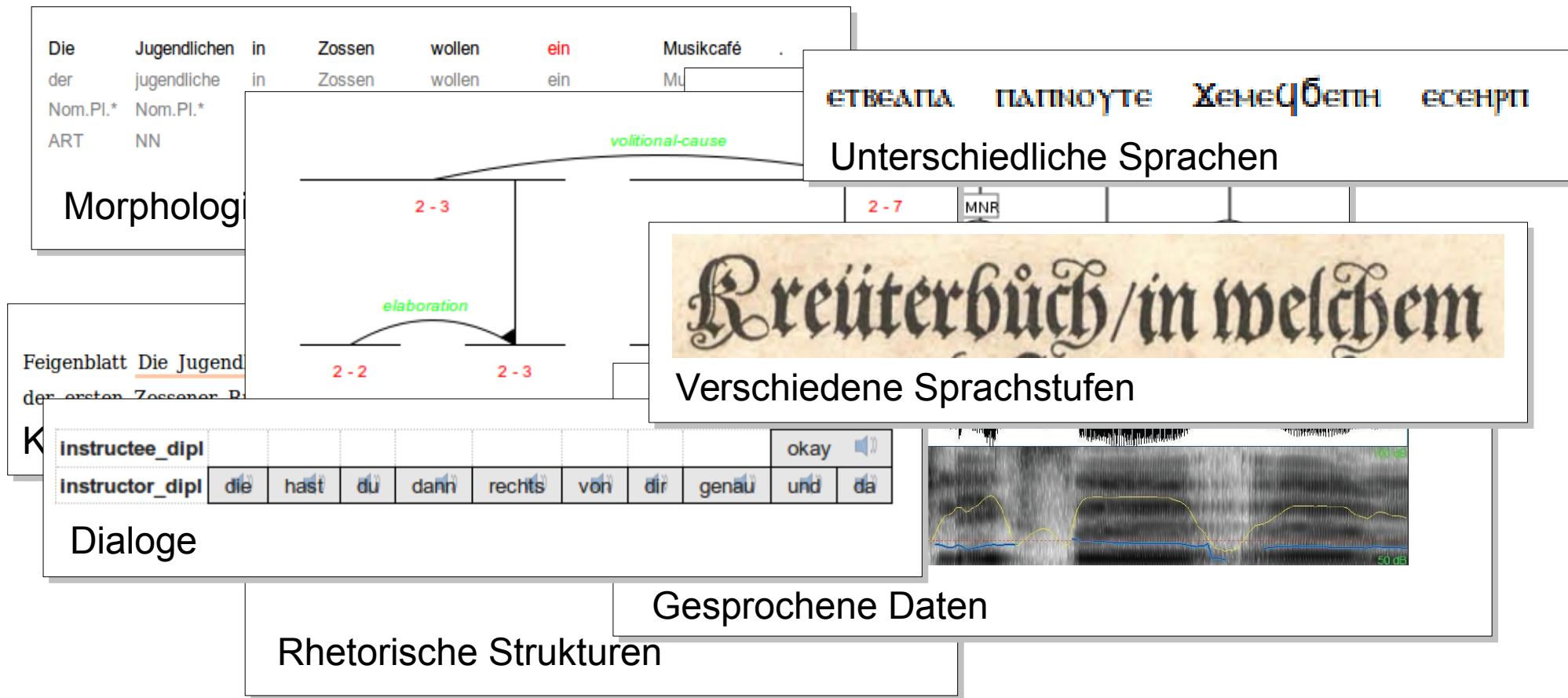


- Linguistische Daten und Phänomene erfordern viele Annotationsarten





- Linguistische Daten und Phänomene erfordern viele Annotationsarten





- Viele Tools, um Daten zu bearbeiten:
 - Manuelle Annotationstools
 - semi-automatische Annotationstools
 - Automatische Annotationstools
 - Suchtools
 - Visualisierungstools

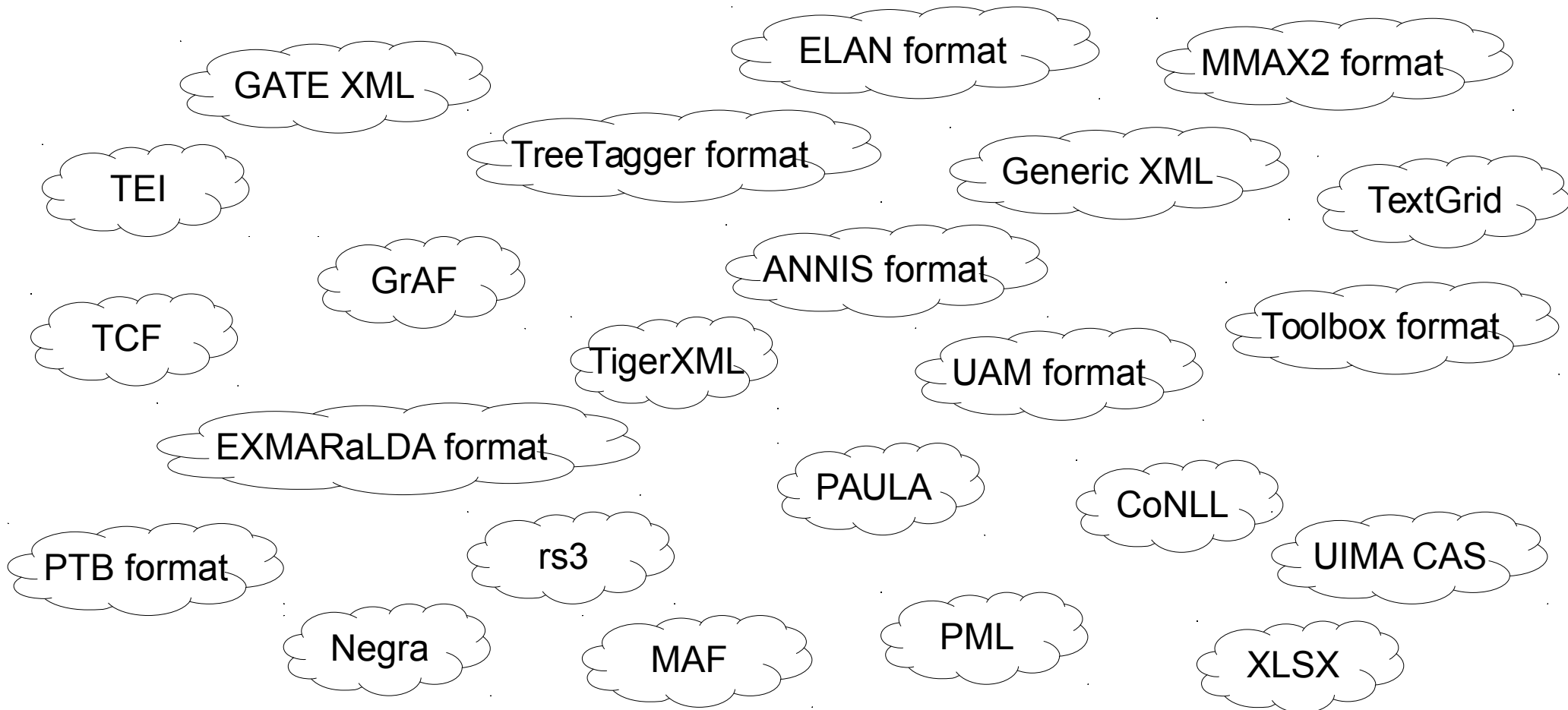


- Viele Tools, um Daten zu bearbeiten:

- EXMARaLDA
- Praat
- ELAN
- Tiger search
- ANNIS
- Gate
- @nnotate
- TrED
- Parser (Berkley, MALT, ...)
- Arborator
- Toolbox
- Synpathie
- TreeTagger
- Weblicht
- MMAX2
- RST
- UIMA
- WebANNO
- ATOMIC
- UAM
- UIMA (dkpro, ...)
- ...



- Viele verschiedene Formate





- Problem 1: Interoperabilität
 - Viele Tools → gut, Nutzer können wählen
 - Aber
 - Tools können nur selten interagieren
 - Primärdaten müssen mehrmals aufbereitet werden (Tokenisierung)



- Problem 2: Mehrebenenkorpora
 - Annotation unterschiedlicher Annotationsarten (Morphologie, Syntax, Koreferenzen) erfordert defacto unterschiedliche Korpora
 - Aber: wir brauchen **ein** Korpus, das alles enthält



- Problem 3: Nachhaltigkeit
 - Einige Tools werden nicht mehr weiterentwickelt
 - Formate werden nicht weiter unterstützt
 - Was ist mit den Daten???



- Nachhaltigkeit der Daten erfordert Nachhaltigkeit der Speicherung
 - Im Web: HTML (W3C)
 - Allgemeine Datenbeschreibung: XML (W3C), JSON
 - Modellierung: UML/ XMI (OASIS)
 - Freitext: PDF bzw. pdf-a



- Es gibt Ideen zur Standardisierung:
 - TEI (TEI consortium)
 - GrAF (ISO)
 - MAF (ISO)
 - SynAF/isoTiger (ISO)

Aber nur wenige Tools arbeiten damit, z.T.

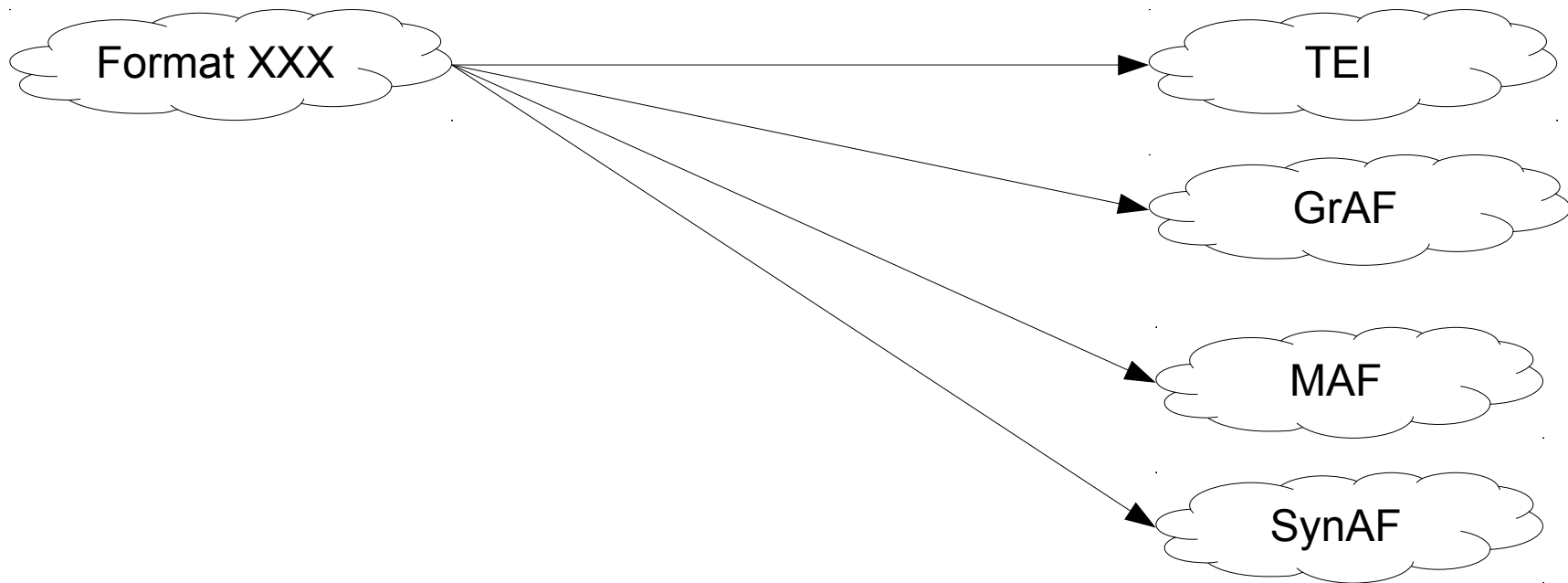
- Sehr komplex
- Unausgereift
- Standards oft jünger als Tool



- Was wir brauchen:
 - Übertragung alter Daten in neue Formate/
Standards (Nachhaltigkeit)
 - Austausch der Daten zwischen unterschiedlichen
Tools (Interoperabilität)
 - Verschmelzen verschiedener Annotationsarten und
-ebenen (Mehrebenenkorpora)



- Nachhaltigkeit:



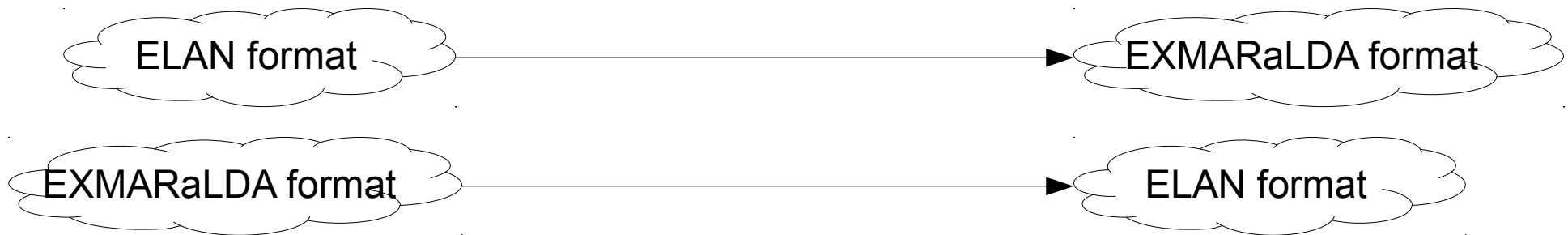


- Austausch:



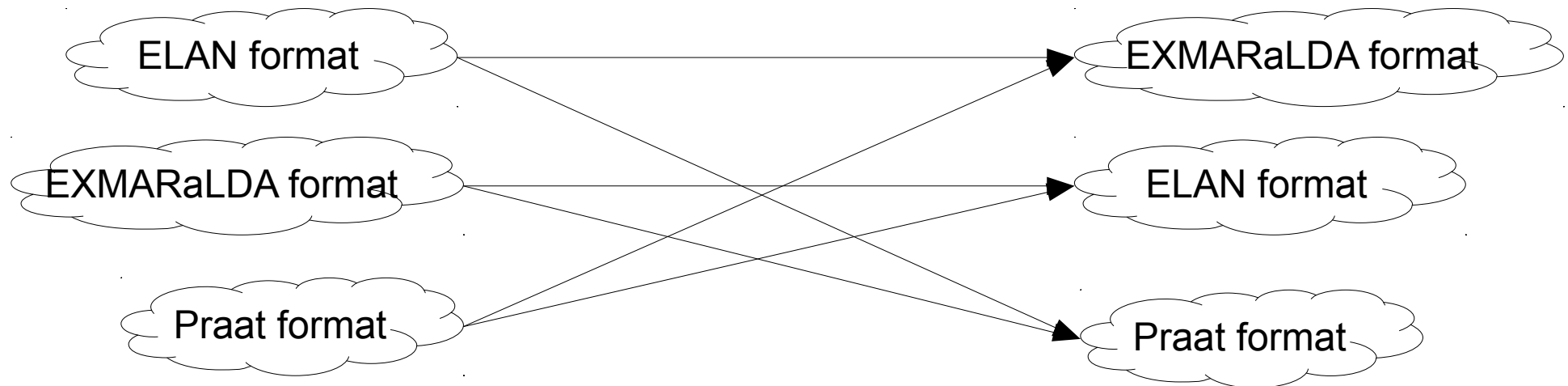


- Austausch:



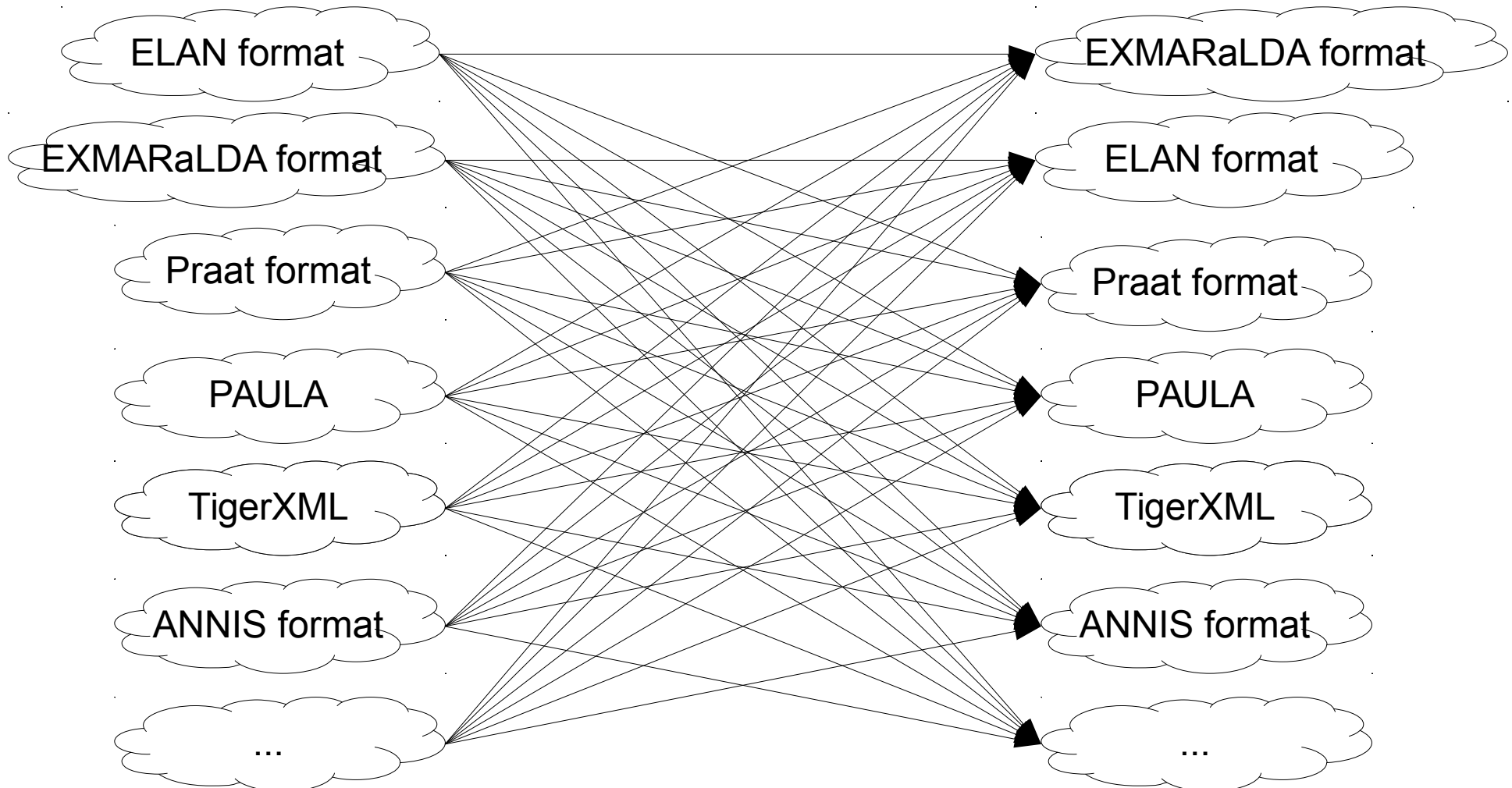


- Austausch:



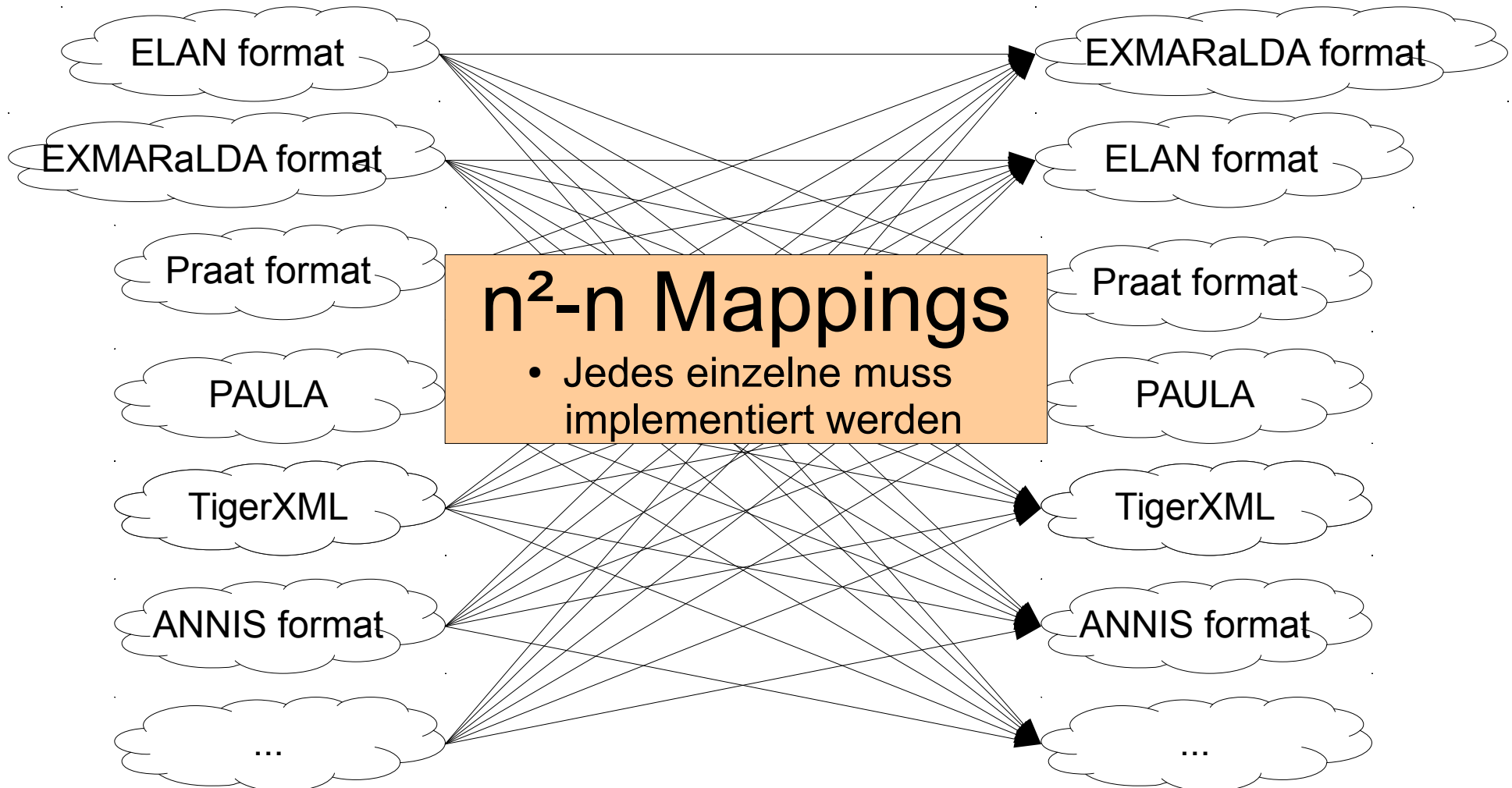


- Austausch:



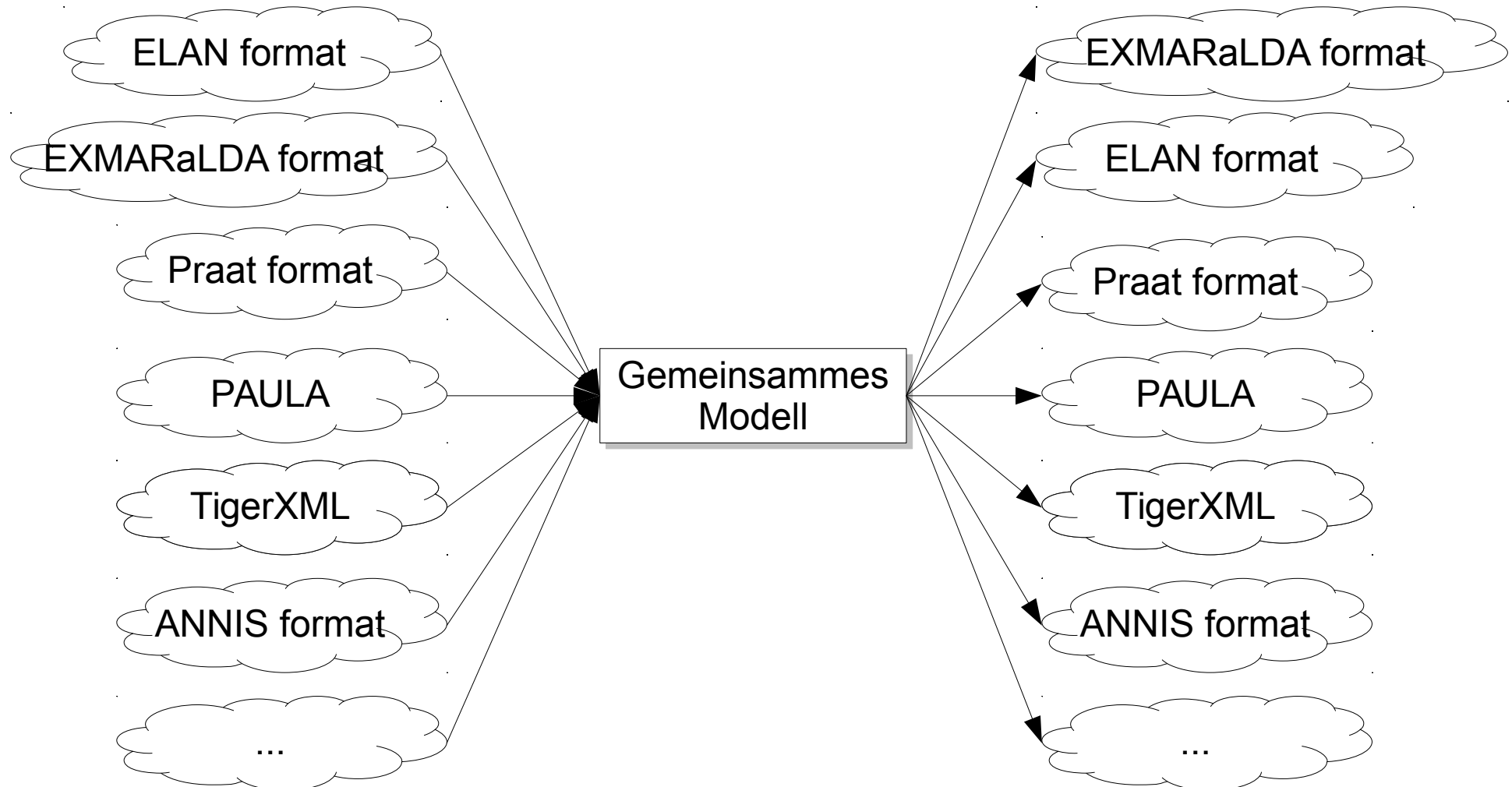


- Austausch:



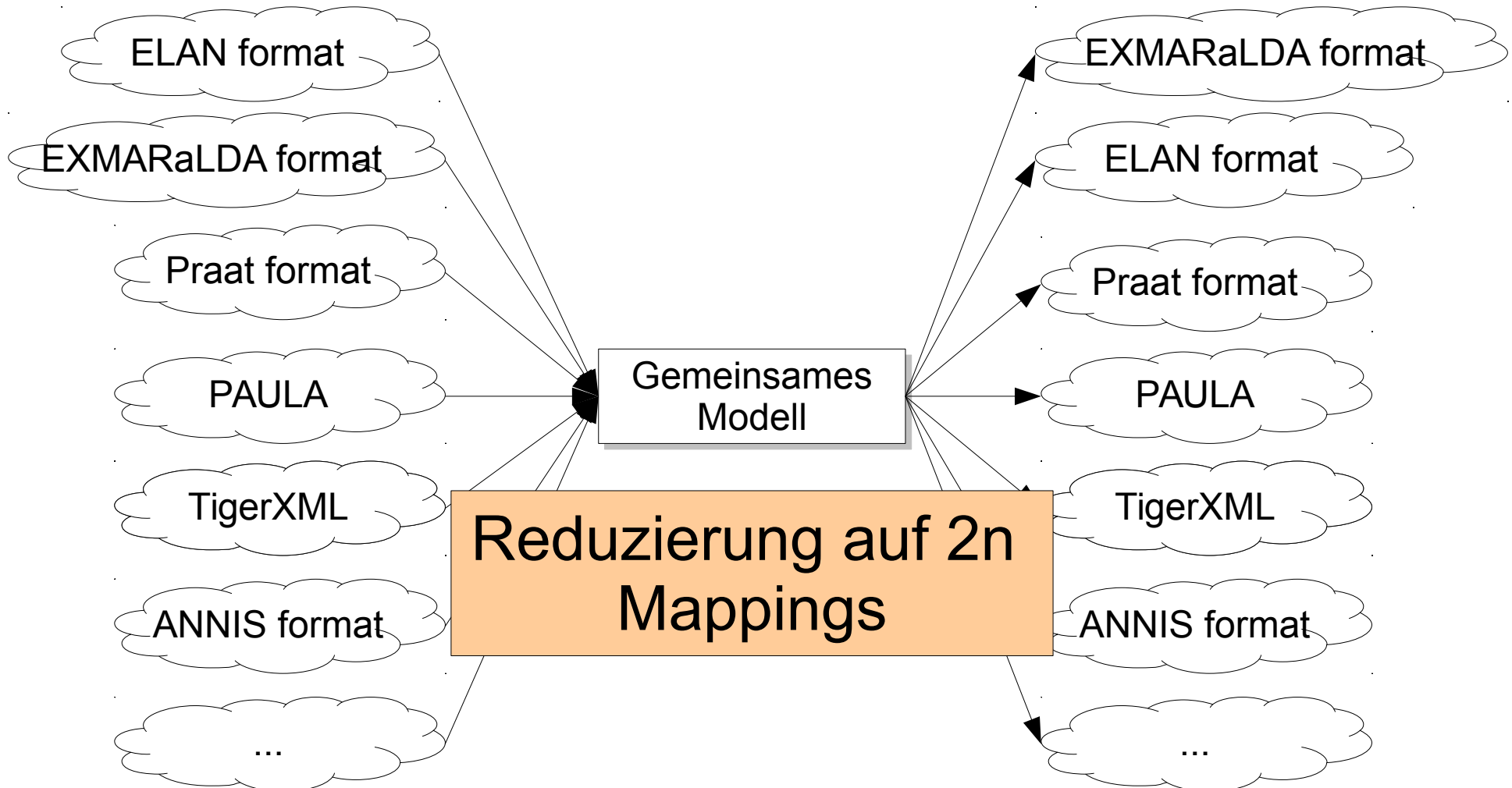


- Austausch:





- Austausch:





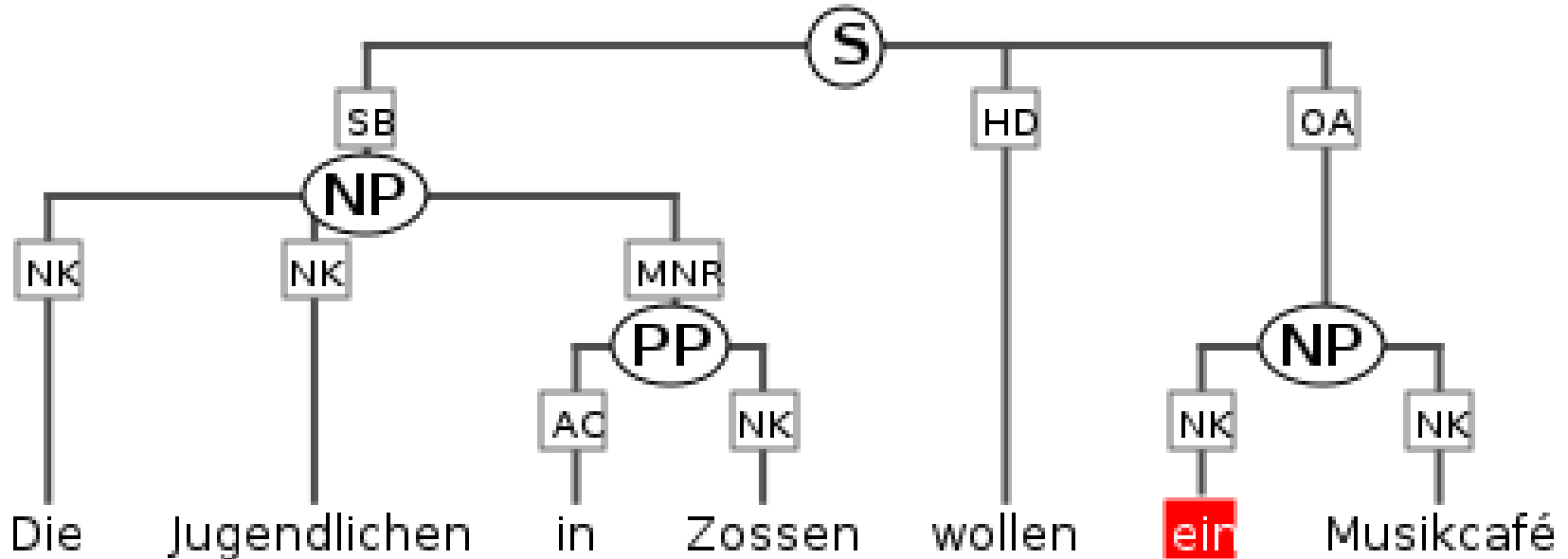
- Anforderungen an Metamodell:
 - Tagsetunabhängig
 - Beliebige Annotationsebenen
 - Unterschiedliche Annotationsarten
 - Theorieneutral



- Salt ist ein Graph?
 - Ein Graph $G = (V, E)$ mit:
 - Einer Menge an Knoten V
 - Einer Menge an Kanten E mit $e = (v_1 \in V, v_2 \in V) \in E$.

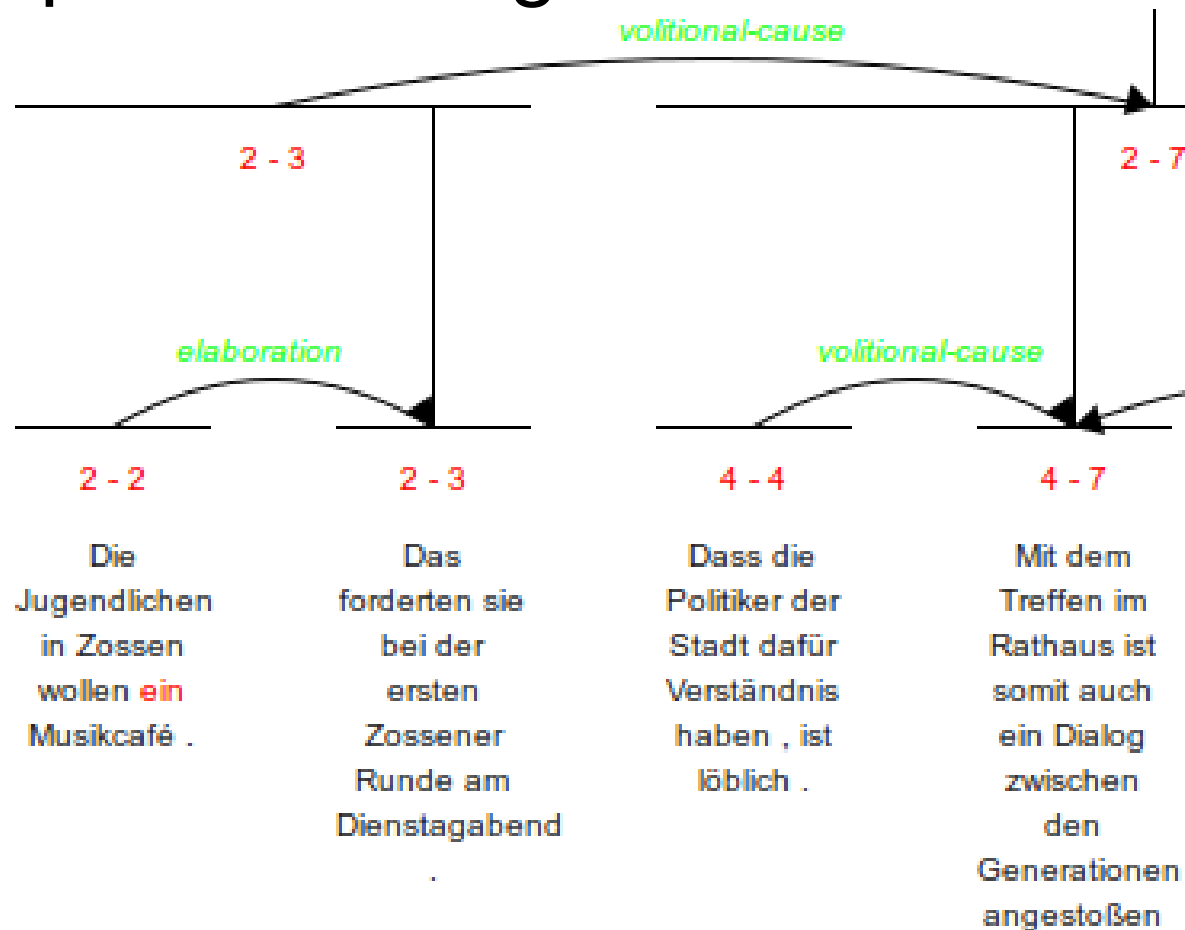


- Ein Graph in der Linguistik?



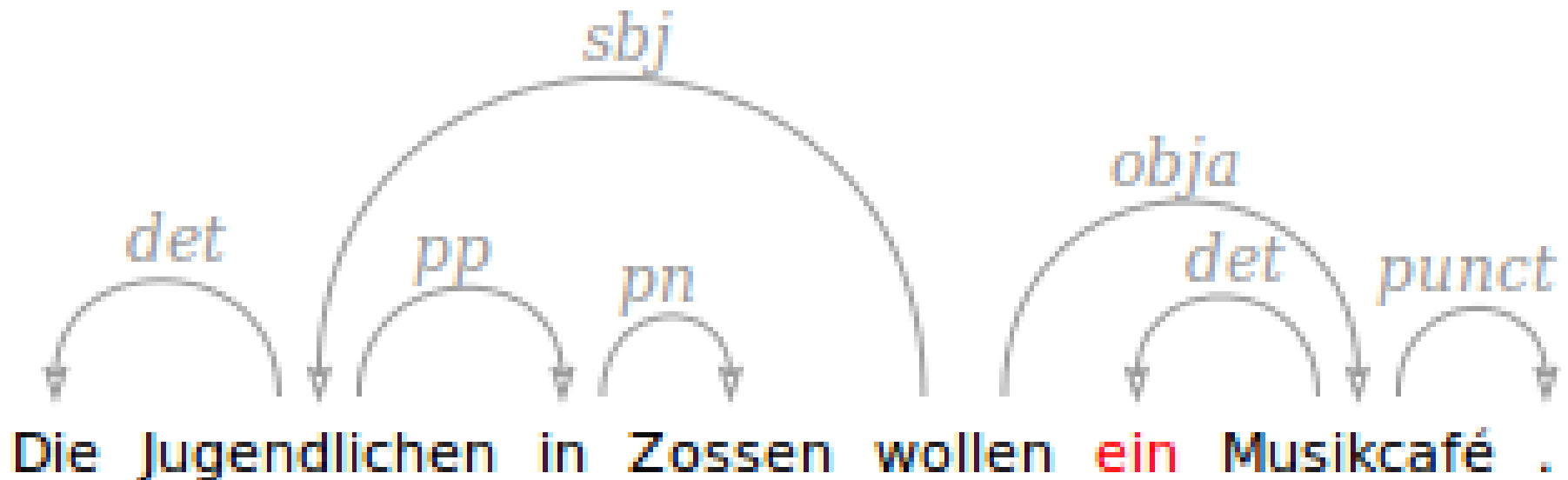


- Ein Graph in der Linguistik?





- Ein Graph in der Linguistik?





- Ist das noch ein Graph?

Feigenblatt Die Jugendlichen in Zossen wollen **ein** Musikcafé . Das forderten sie bei der ersten Zossener Runde am Dienstagabend . Dass die Politiker der Stadt dafür



- Und das?

Focus_newInf				nf-unsol							
Inf-Stat	new			new	new			giv-active		giv-active	
NP	NP			NP	NP			NP		NP	
PP				PP							
Sent	s							s			
Topic	ab										ab
tok	Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.	Das	forderten	sie

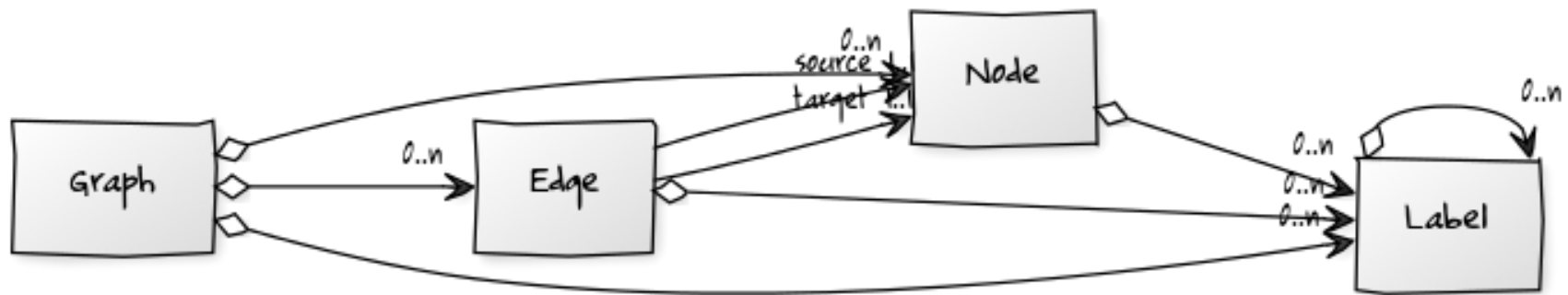


- Oder das?

Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Musikcafé	.
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Acc.Sg.Neut	–
ART	NN	APPR	NE	VMFIN	ART	NN	\$.



- Für Salt ja!





Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Musikcafé	.
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Acc.Sg.Neut	–
ART	NN	APPR	NE	VMFIN	ART	NN	\$.

Primärtext:

Die Jugendlichen in Zossen wollen ein Musikcafé.



Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Musikcafé	.
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Acc.Sg.Neut	–
ART	NN	APPR	NE	VMFIN	ART	NN	\$.

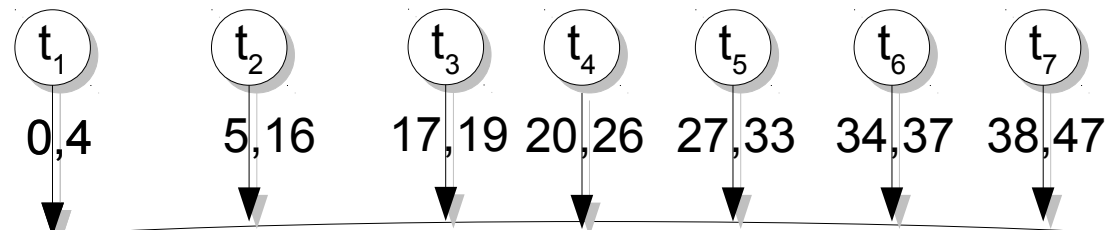
Primärtext:

Die Jugendlichen in Zossen wollen ein Musikcafé.



Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Musikcafé	.
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Acc.Sg.Neut	–
ART	NN	APPR	NE	VMFIN	ART	NN	\$.

Tokenisierung:



Primärtext:

Die Jugendlichen in Zossen wollen ein Musikcafé.

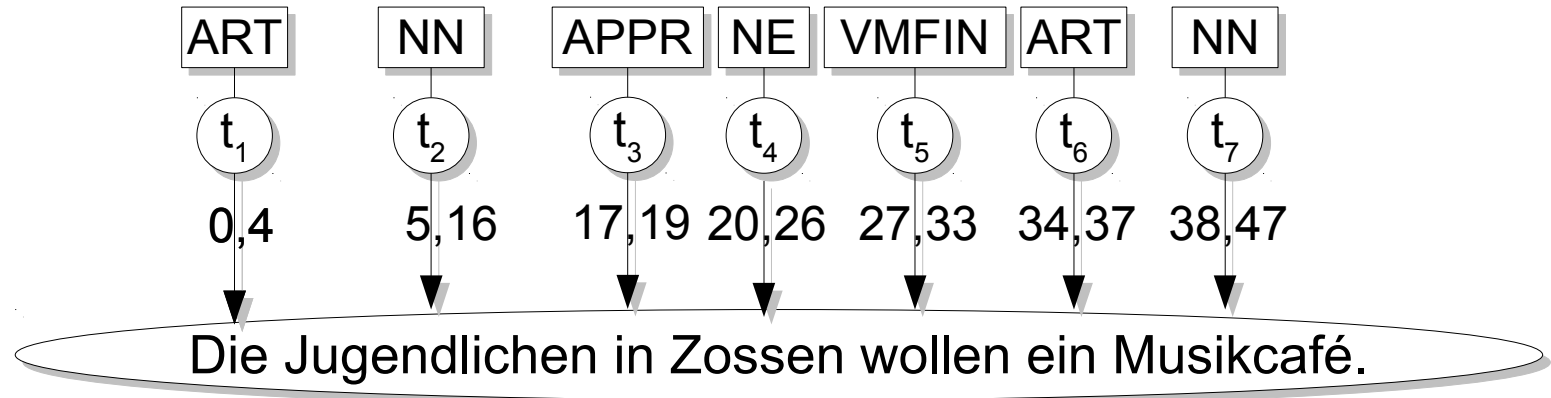


Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.
der	jugendliche	in	Zossen	wollen	ein	Musikcafé	.
Nom.Pl.*	Nom.Pl.*	–	Dat.Sg.Neut	3.Pl.Pres.Ind	Acc.Sg.Neut	Acc.Sg.Neut	–
ART	NN	APPR	NE	VMFIN	ART	NN	\$.

Annotation:

Tokenisierung:

Primärtext:





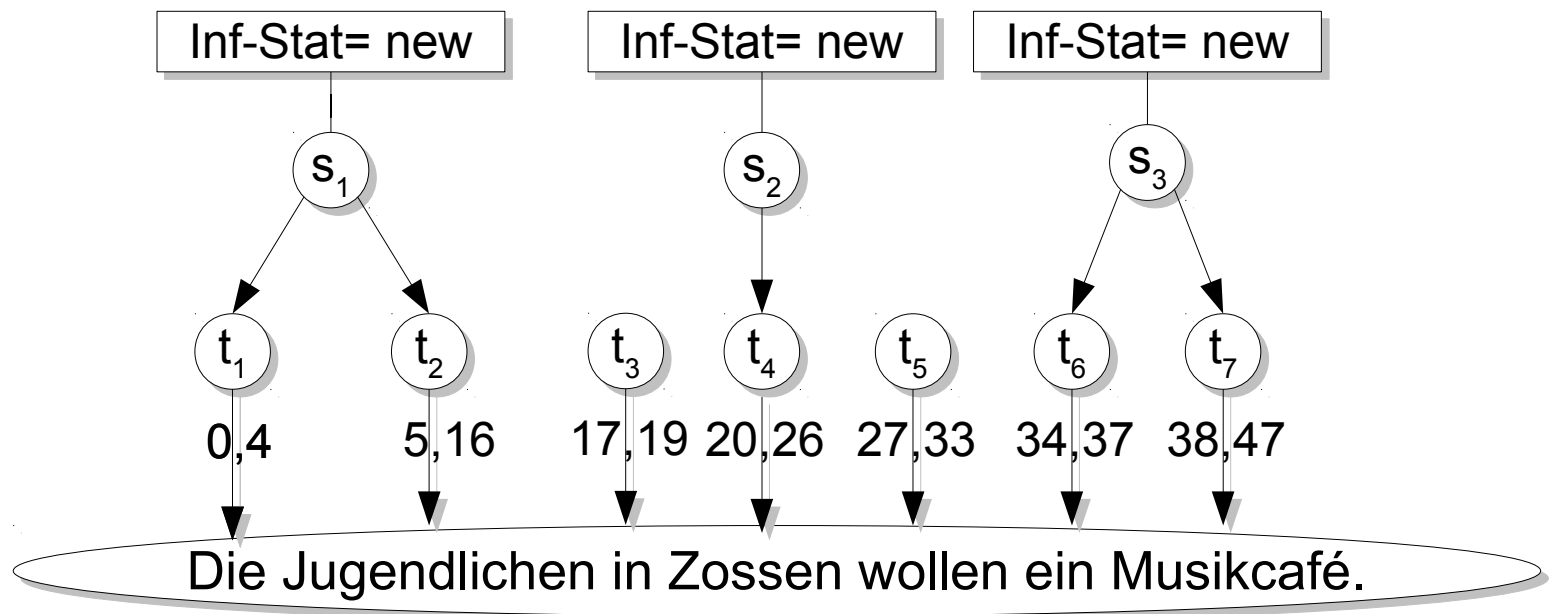
Inf-Stat	new				new			new			giv-active			giv-active	
NP	NP				NP			NP			NP			NP	
PP			PP												
Sent	s										s				
Topic	ab												ab		
tok	Die	Jugendlichen	in	Zossen	wollen	ein	Musikcafé	.	Das		forderten		sie		

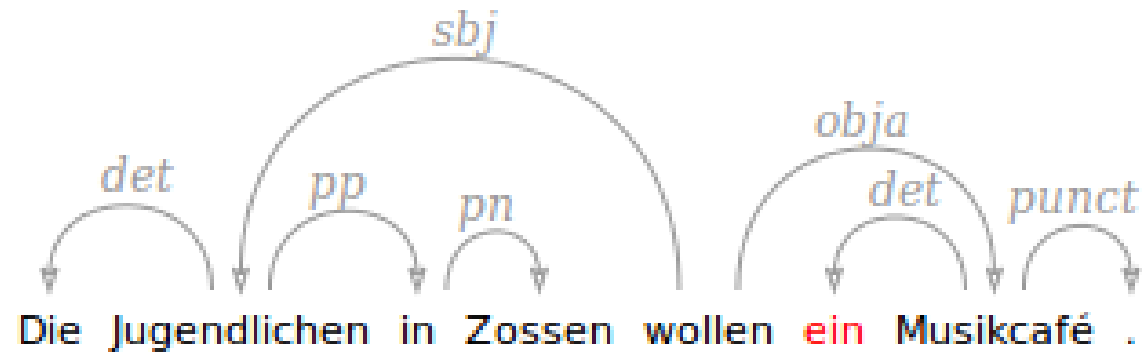
Annotation:

Mengen:

Tokenisierung:

Primärtext:

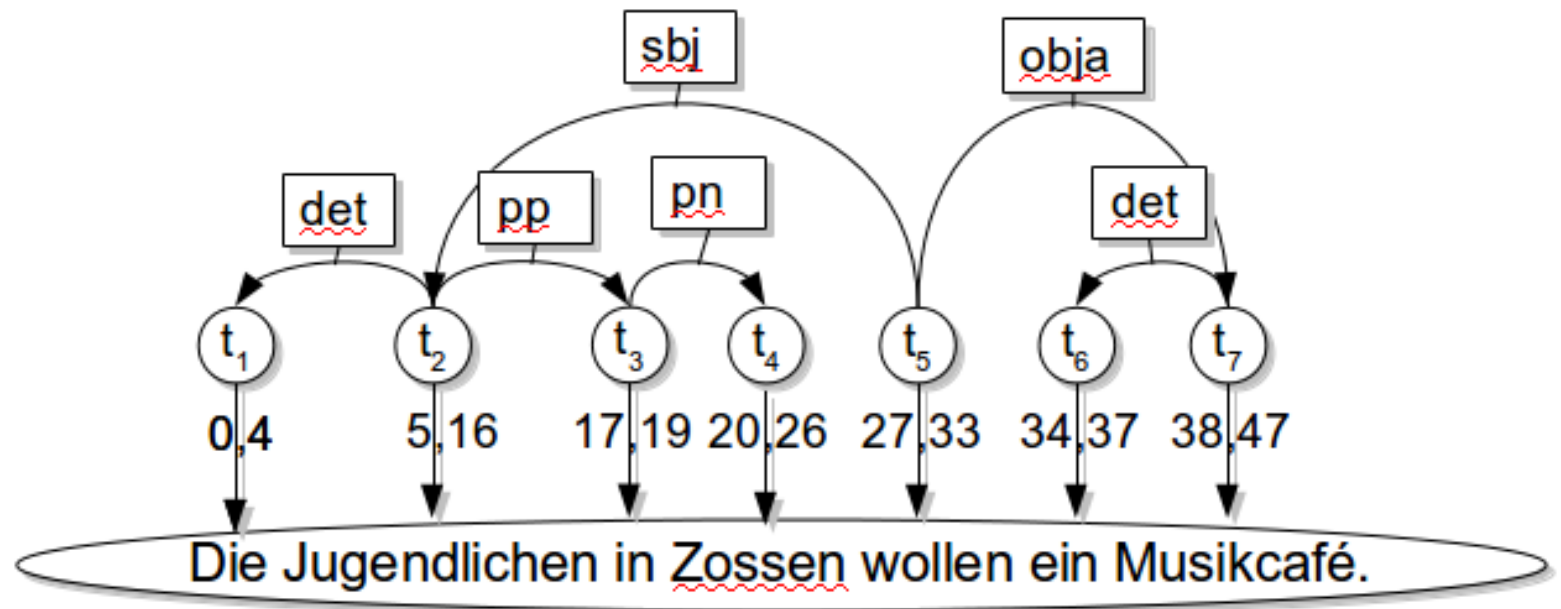




Kanten:

Tokenisierung:

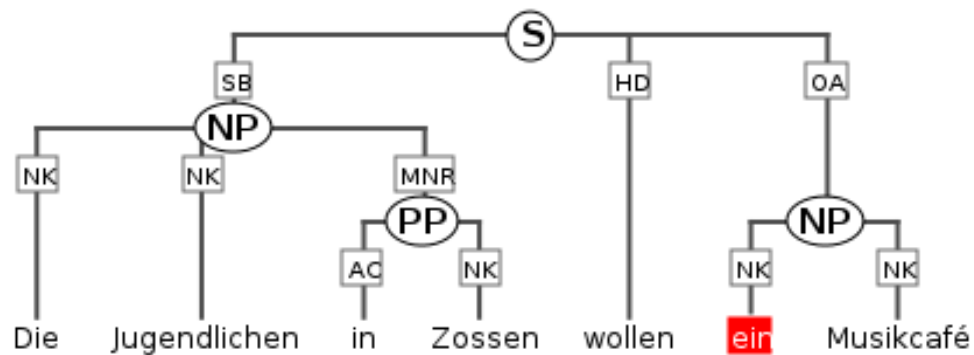
Primärtext:





SaltNPepper und das Formatpluriversum

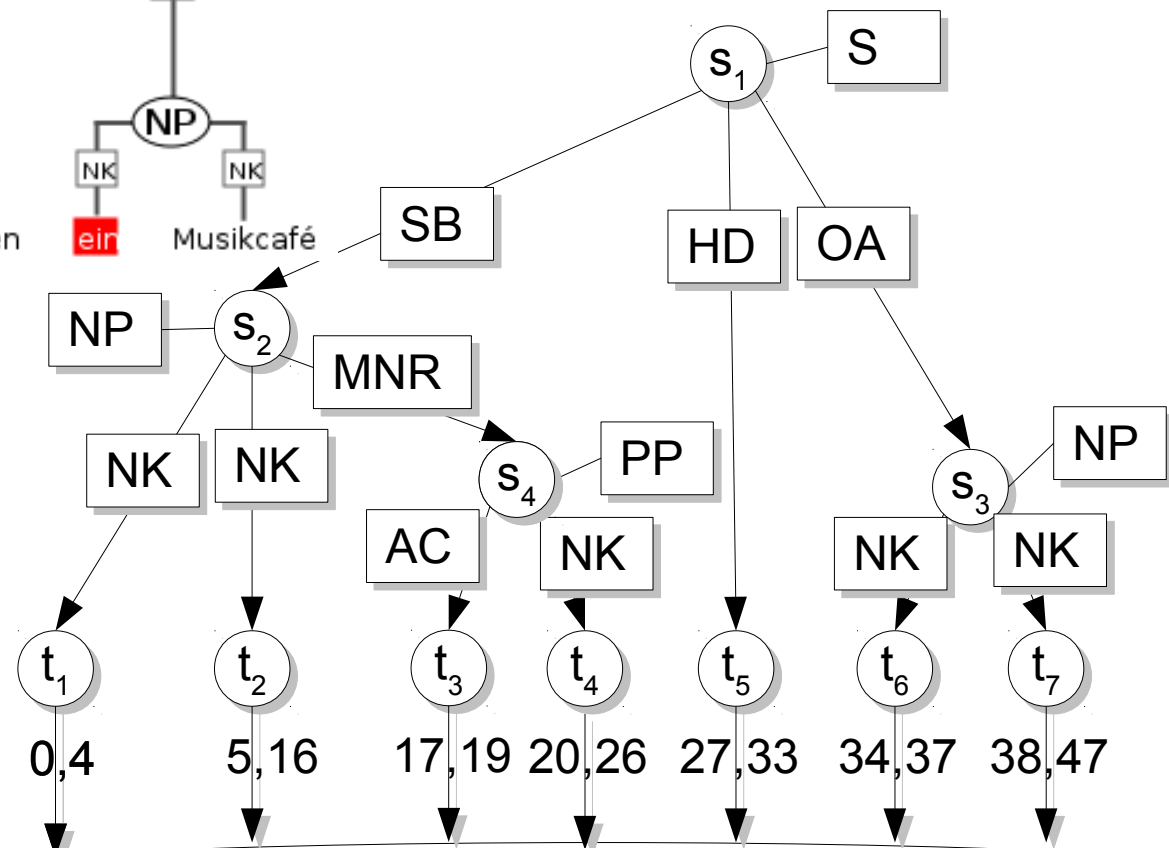
Salt



Hierarchien:

Tokenisierung:

Primärtext:



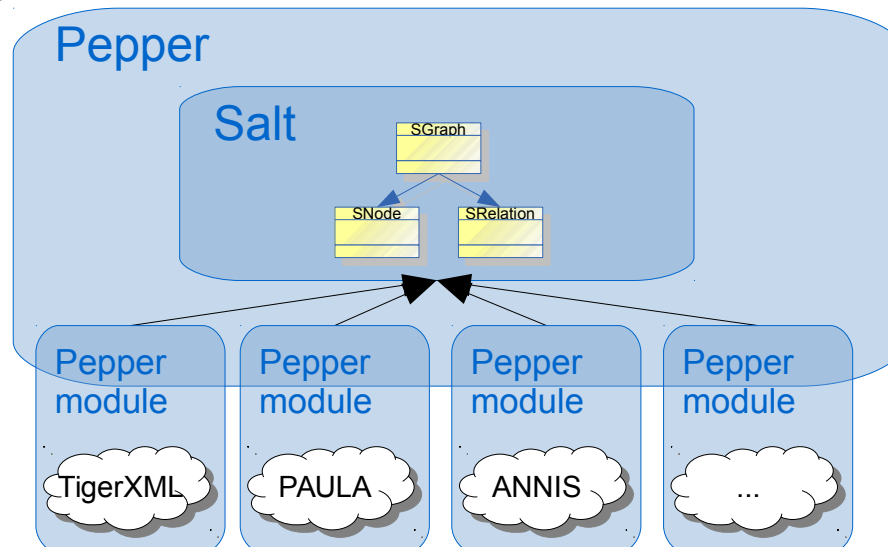
Die Jugendlichen in Zossen wollen ein Musikcafé.



- Anforderungen an Metamodell:
 - ☑ Tagsetunabhängig
frei wählbare Attribut-Wert-Paare für Labels
 - ☑ Beliebige Annotationsebenen
unbegrenzte Anzahl an Labels
 - ☑ Unterschiedliche Annotationsarten
alles, was als Graph darstellbar ist
 - ☑ Theorieneutral
Semantikarmut, Salt kennt nur Zeichenketten

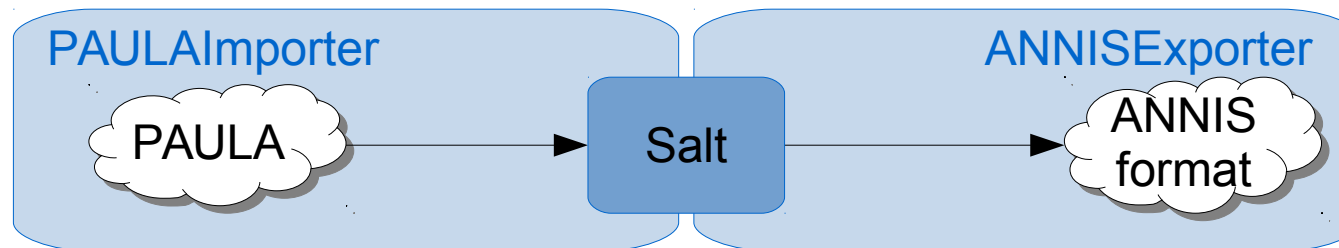


- Pepper
 - Converterframework
 - Basiert auf Salt
 - Nur eine Infrastruktur, die Arbeit machen die Plugins



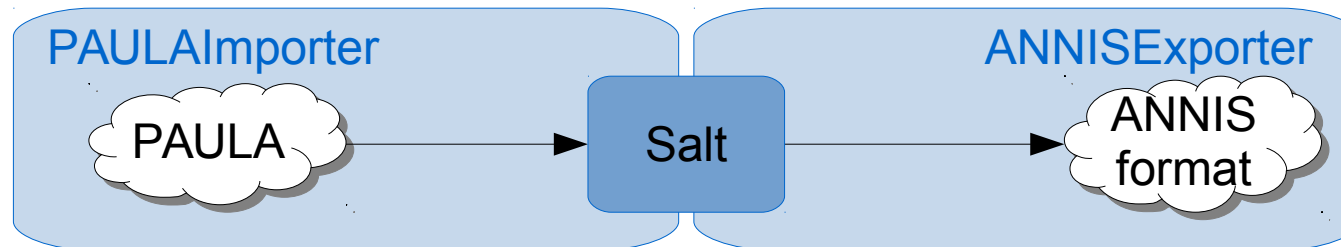


- Drei Arten von Modulen:
 - Importer: Format A \rightarrow Salt
 - Manipulator: Salt \rightarrow Salt
 - Exporter: Salt \rightarrow Format B

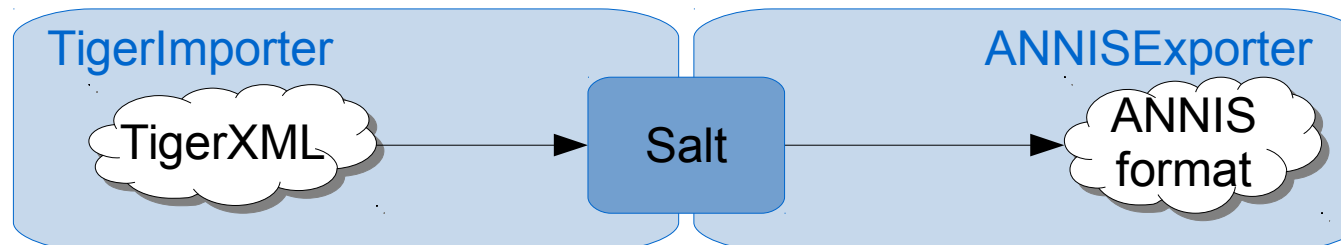




- Drei Arten von Modulen:
 - Importer: Format A \rightarrow Salt
 - Manipulator: Salt \rightarrow Salt
 - Exporter: Salt \rightarrow Format B

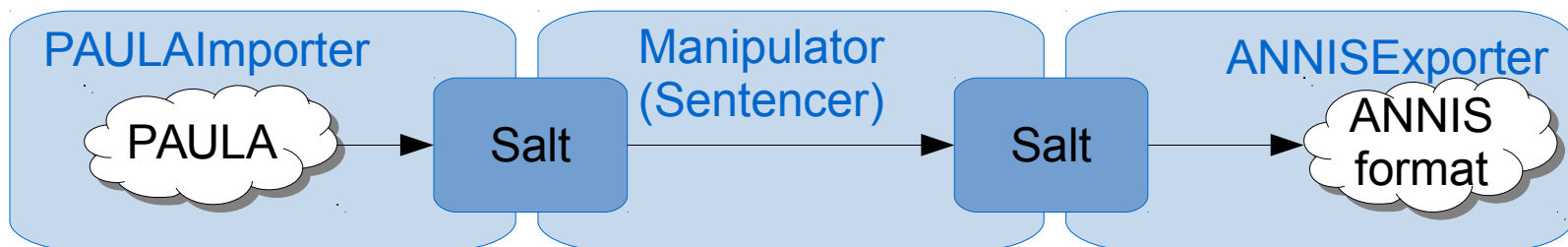


- Kombinierbarkeit



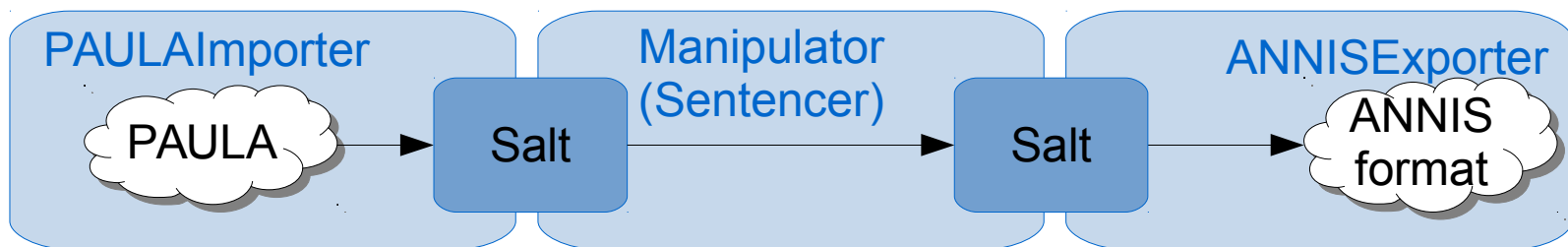


- Manipulation

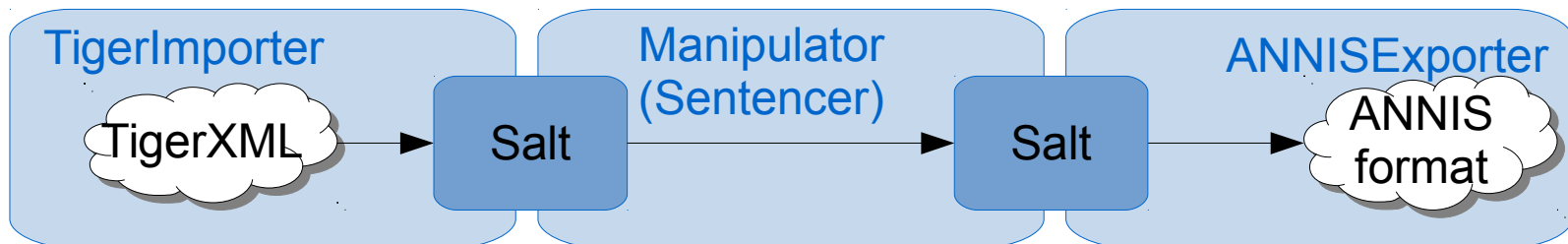




- Manipulation



- Kombinierbarkeit

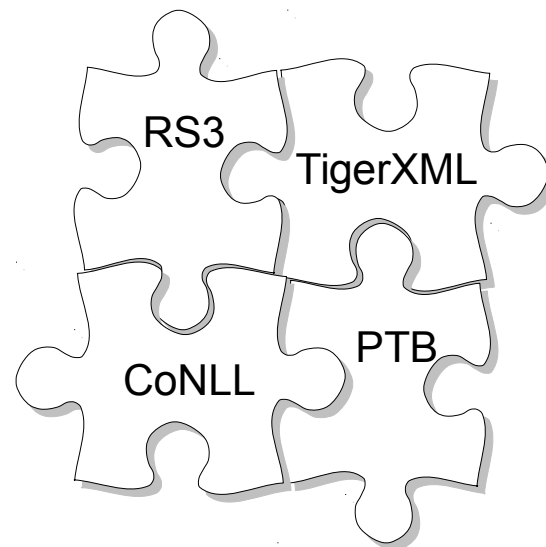




- Was wir brauchen:
 - ☒ Übertragung alter Daten in neue Formate/
Standards (Nachhaltigkeit)
 - ☒ Austausch der Daten zwischen unterschiedlichen
Tools (Interoperabilität)
 - ☐ Verschmelzen verschiedener Annotationsarten und
-ebenen (Mehrebenenkorpora)

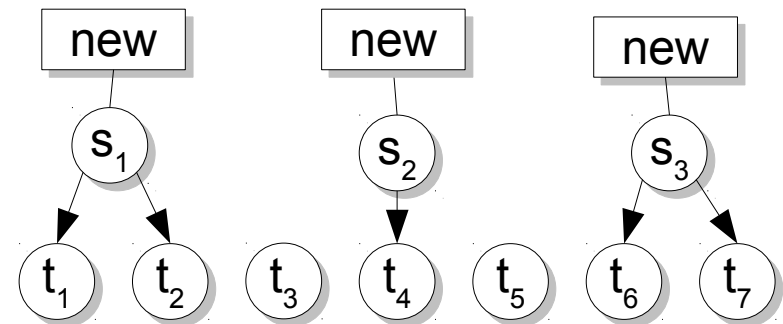
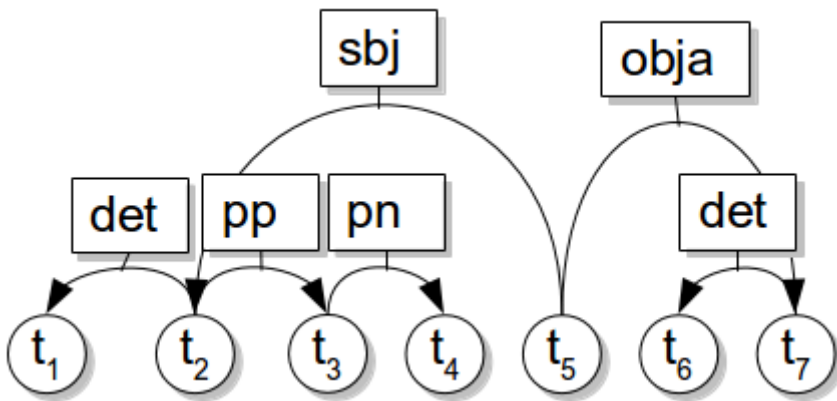


- Problem: es gibt nur wenige Mehrebenenannotationstools (bspw. WebAnno, ATOMIC)
- Idee: Verschmelzen der unterschiedlichen Formate (und somit der Ebenen)



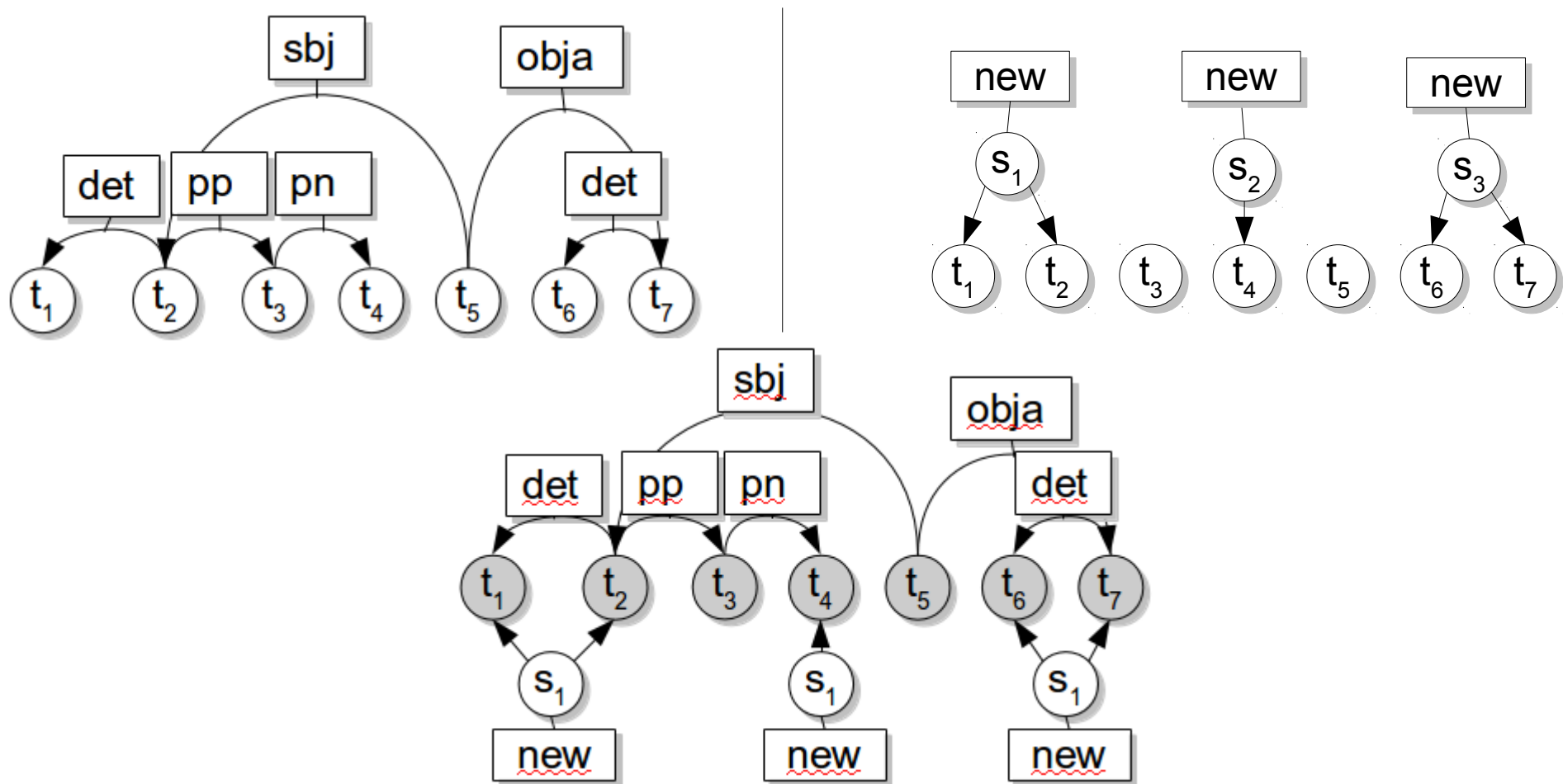


- Salt reduziert Merging zu Graphmerging



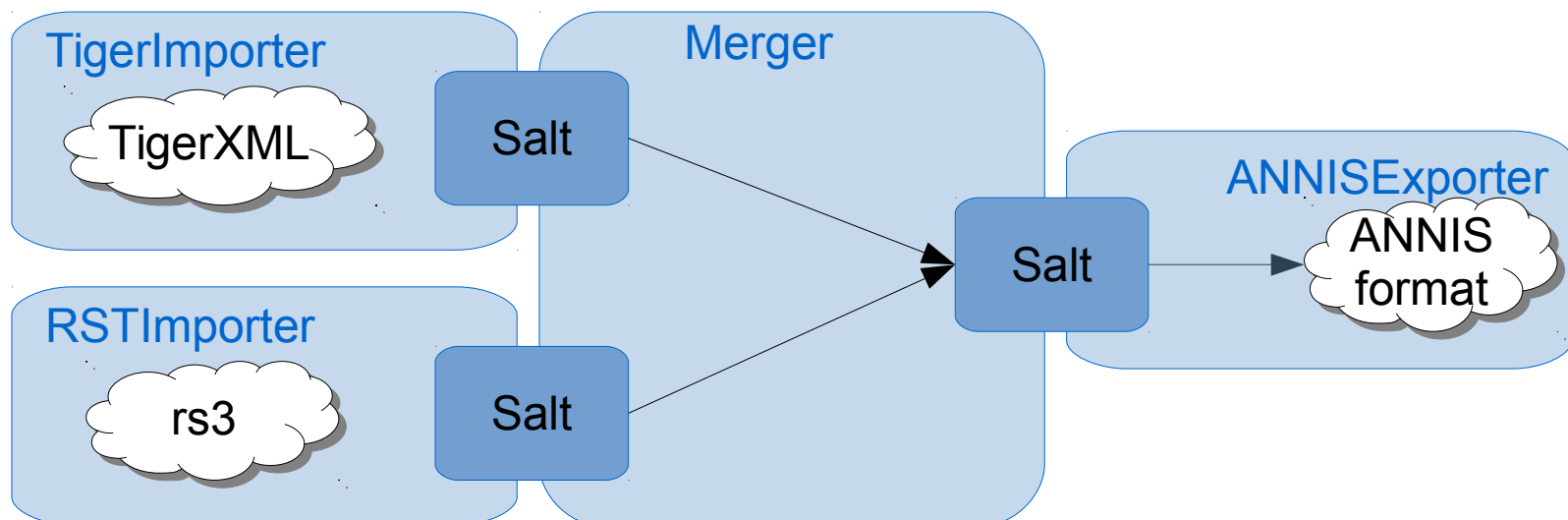


- Salt reduziert Merging zu Graphmerging





- Merger ist Plugin für Pepper (Manipulator)





- Was wir brauchen:
 - ☒ Übertragung alter Daten in neue Formate/
Standards (Nachhaltigkeit)
 - ☒ Austausch der Daten zwischen unterschiedlichen
Tools (Interoperabilität)
 - ☒ Verschmelzen verschiedener Annotationsarten und
-ebenen (Mehrebenenkorpora)



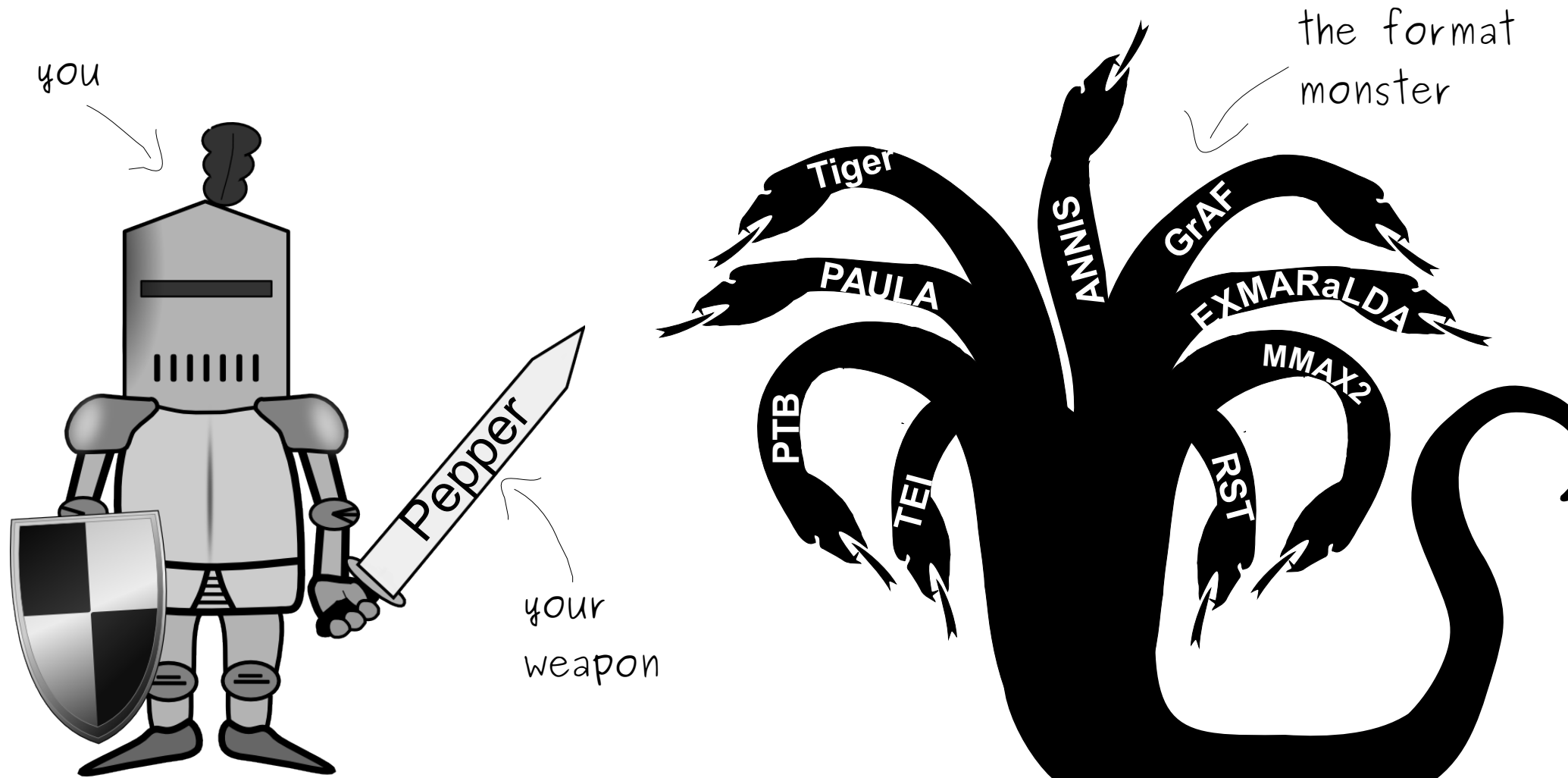
- SaltNPepper
 - Konvertierung von Korpora zwischen Formaten
 - Erweiterbarkeit um neue Formate (Plugins)
 - Open Source (Apache License 2.0)
 - Öffentliche Plattform: GitHub
 - <https://github.com/korpling/pepper>
 - <https://github.com/korpling/salt>



- Nachhaltigkeit von Korpora, Formaten und Software hängt zusammen
- Problem: Projekte sind befristet!
 - Oft stirbt Software nach Ende eines Projektes → Verlust von Geld und Zeit



- Nachhaltigkeit von Korpora, Formaten und Software hängt zusammen
- Problem: Projekte sind befristet!
 - Oft stirbt Software nach Ende eines Projektes → Verlust von Geld und Zeit
- Software braucht zum Überleben:
 - Aktive Entwicklercommunity
 - Open Source
 - Öffentliche Plattform
 - Gute Dokumentation





- Diese Folien wurden erstellt unter Verwendung von:
 - Yuml <http://yuml.me>
 - Openclipart <http://openclipart.org>