

Eine XML-Kodierung für AVM-Beschreibungen*

*Manfred Sailer und Frank Richter***

Zusammenfassung

Ausgehend von der Beobachtung, dass existierende Plattformen für HPSG-Grammatiken untereinander inkompatible interne Datenformate für AVMs (Attribut-Wert-Matrizen) wählen, schlägt das vorliegende Papier eine XML-basierte Beschreibung von AVMs vor. Diese Normierung von AVMs ist geeignet, als Transferformat an den Schnittstellen existierender HPSG-Systeme zu fungieren und somit die Problematik miteinander unvereinbarer Datenstrukturen zu lindern.

15.1. Einleitung

Im Bereich der HPSG (Head-Driven Phrase Structure Grammar, Pollard und Sag, 1994) ist bereits eine Reihe von Datenbanken vorhanden oder im Aufbau (Flickinger et al., 2000, Hinrichs et al., 2000a,b, Kordoni, 2000, Marciniak et al., 2000), und es gibt Implementierungsplattformen mit Grammatiken mit relativ hoher empirischer und/oder theoretischer Abdeckung wie ALE (Carpenter und Penn, 1996), Babel, LKB (2000), PAGE, ConTroll, Trale (Penn, 2000). Aufgrund der unterschiedlichen internen Datenformate, die diese Systeme verwenden, wäre es wünschenswert, ein einheitliches und der maschinellen Weiterverarbeitung leicht zugängliches Format zu definieren, in das die system-abhängigen Kodierungen übersetzt werden können. Hierfür bietet sich die *Extensible Markup Language* (XML) an.

Für linguistische Datenbanken gibt es bereits Bemühungen, ein XML-Format zu finden, das den Datenaustausch erleichtert. So wird in Mengel und Lezius (2000) gezeigt, wie verschiedene Arten von Baumrepräsentationen in XML enkodiert werden können. Daneben gibt es in Bonhomme und Lopez (1999) einen Vorschlag zur XML-Enkodierung von Strukturen der Tree-Adjoining Grammar (TAG). Dabei handelt es sich um Bäume, die Attribut-Wert-Paare als Knotenlabels haben können; außerdem werden einem Satz zum Teil mehrere Bäume simultan zugeordnet.

Für die HPSG ist keine dieser Kodierungsformen geeignet. Die Ausgabe von HPSG-Parsern sind in der Regel komplexe Attribut-Wert-Matrizen (*attribute-value matrice*, AVM). AVMs sind im HPSG-Verständnis *Beschreibungen* von linguistischen Objekten (vgl. Richter, 2000). Im Gegensatz dazu sind Bäume, auch solche mit komplexen Knotenlabels wie in TAG, *Modelle* von

* Erschienen in: *Proceedings der GLDV-Frühjahrstagung 2001*, Henning Lobin (Hrsg.), Universität Gießen, 28.–30. März 2001, Seite 161–168. <http://www.uni-giessen.de/fb09/asc1/gldv2001/>

** Wir danken Maria Atanasowa, Lothar Lemnitzer, Zdravko Peev, Kiril Simov, Ilona Steiner und Tylman Ule für Kommentare und Diskussion. Eine ausführlichere Version des vorliegenden Papiers ist online erhältlich unter der Adresse <http://www.sfs.nphil.uni-tuebingen.de/~mf/papers/index.html>.

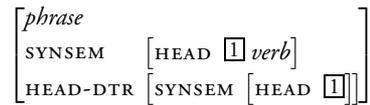


Abbildung 15.1.: Eine einfache AVM

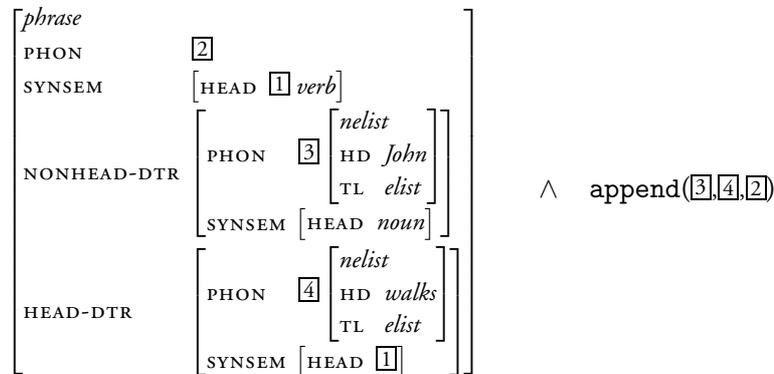


Abbildung 15.2.: Eine AVM mit Junktor und relationalem Ausdruck

linguistischen Objekten. Damit haben Baumbanken prinzipiell einen anderen ontologischen Status als „AVM“-Banken.¹

Die verwendeten AVMs sind rekursive *Klammerstrukturen*, in denen neben atomaren Werten auch AVMs oder Tags auftreten können. Abb. 15.1 zeigt eine einfache AVM, wobei Attributnamen in Kapitälchen erscheinen (SYNSEM), Sortennamen kursiv geschrieben sind (*phrase*) und Tags als Zahlen in Boxen dargestellt werden ($\boxed{1}$). Diese AVM beschreibt all diejenigen Objekte, die das Label *phrase* tragen, und für die gilt, dass die Kante SYNSEM zu einem Objekt führt, von dem aus eine Kante HEAD zu einem weiteren Objekt mit dem Label *verb* führt. Dieses letztgenannte Objekt mit Label *verb* wird ebenfalls erreicht über die Kanten HEAD-DTR, SYNSEM und HEAD. Durch den Tag $\boxed{1}$ wird diese Identität ausgedrückt.

Dass es sich bei einer AVM um eine Beschreibung und nicht um ein linguistisches Objekt selbst handelt, sieht man beispielsweise daran, dass in Abb. 15.1 keine Phonologie angegeben ist, obwohl jedes Objekt der Sorte *phrase* einen PHONOLOGY-Wert hat.

AVMs sind nur ein Teil der HPSG-Beschreibungssprache: Auf Plattformen wie ConTroll und Trale können neben den eigentlichen AVM-Klammerstrukturen auch andere Elemente der Beschreibungssprache im Input und im Output erscheinen. Richter (2000) definiert eine AVM-Beschreibungssprache für HPSG, die auch *relationale Ausdrücke*, *logische Junktoren* und *Quantoren* umfasst. In Abb. 15.2 ist eine solche Beschreibung gegeben, die aus einer Klammerstruktur, dem Junktor „ \wedge “, sowie der Relation *append* besteht.² Das nicht-logische Vokabular von Beschreibungen (Attribut-, Sorten- und Relationsnamen, sowie Stelligkeit von Relationen) wird wie in logischen Sprachen wie z. B. in Sprachen der Prädikatenlogik erster Stufe üblich über eine *Signa-*

¹ Götz und Meurers (1998) argumentiert, dass es für HPSG-basiertes Parsen sinnvoll ist, nicht Modelle zu konstruieren, sondern Beschreibungen zu finden.

² Wir ignorieren Quantoren und interpretieren im Folgenden alle Variablen (Tags) als existentiell abgebunden.

festgelegt. Das logische Vokabular hingegen (Tags und Junktoren) ist signatur-unabhängig.³

In diesem Papier stellen wir die in Abb. 15.3 gezeigte allgemeine Architektur für die XML-Kodierung von AVM-Beschreibungen vor. Da alle Signaturen aufgrund ihrer Definition dieselbe Gestalt haben, geben wir eine DTD für Signaturen vor, die in der Datei `signature.dtd` angegeben ist. Die konkrete Signatur einer Grammatik, das heißt eine Signatur mit bestimmten Attributen, Sorten und Relationen, ist dann als eine Instanz dieser DTD zu verstehen und als XML-Datei `signature.xml` kodiert. Da die Signatur die Beschreibungssprache genau festlegt, können wir aus der Datei `signature.xml` eine DTD für AVM-Beschreibungen, `avm.dtd`, erzeugen. Hierzu verwenden wir XSL-Transformationen, `sig2avm.xsl`. Instanzen der DTD `avm.dtd`, in der Abbildung als `avm.xml` angegeben, sind XML-Kodierungen von AVM-Beschreibungen. Diese Beschreibungen können in ihrer linguistischen Anwendung verschiedenerlei Funktionen erfüllen. Sie können Prinzipien der Theorie einer Grammatik sein; oder relativ allgemeine Beschreibungen linguistischer Objekte, so wie sie als Anfragen an HPSG-Parser gegeben werden; oder relativ exakte Beschreibungen linguistischer Objekte, wie sie als Output von HPSG-Parsern geliefert werden.

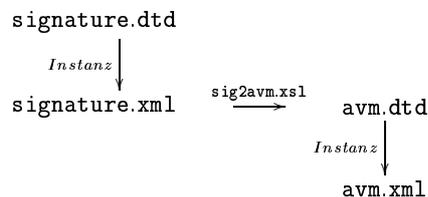


Abbildung 15.3.: Überblick

Im vorliegenden Papier gehen wir nur auf die Dateien `avm.dtd` und `avm.xml` ein, und nehmen eine feste Signatur an. In der ausführlicheren Version des Papiers werden auch die Signaturdateien erläutert.

15.2. Die Kodierung von AVM-Beschreibungen

In Abb. 15.4 (S. 164) ist eine automatisch erzeugte DTD für AVM-Beschreibungen gegeben. Diese DTD besteht aus einem signatur-spezifischen Teil und einem signatur-unabhängigen Teil. Im ersten Teil werden die in der Signatur eingeführten Namen der Sorten, Attribute und Relationen als mögliche Werte der parametrisierten Entitäten `%sorts`; `%attrs`; und `%relation`; aufgeführt. Für jede Relation, hier nur `append`, wird außerdem ein Element eingeführt, das die Stelligkeit der Relation ausdrückt.

Wir definieren ein Element `descr`, das aus einer beliebig langen Folge von Klammerstrukturen, relationalen Ausdrücken oder logischen Verknüpfungen besteht. Klammerstrukturen werden durch `avm`-Elemente kodiert, relationale Ausdrücke durch Elemente der parametrisierten Entität `%relation`; und Junktoren durch Elemente der Entität `%constant`;. Abb. 15.5 (S. 165) zeigt die XML-Kodierung der AVM aus Abb. 15.1.

³ Darüberhinaus enthält eine HPSG-Signatur die Sortenhierarchie und die Attributsdeklarationen. Diese Informationen sind wichtig zur Festlegung der Ontologie linguistischer Objekte, nicht jedoch für die Definition der Beschreibungssprache.

```

<!--      automatisch erzeugte DTD fuer AVM-Beschreibungen      -->

<!-- ##### Signatur-spezifischer Teil ##### -->

<!--      Sorten      -->
<!ENTITY % sorts      "(sign|word|phrase|synsem|
                        list|e_list|ne_list)"      >
<!--      Attribute      -->
<!ENTITY % attrs      "(phon|synsem|h_dtr|n_dtr|hd|tl)"      >
<!--      Relationen      -->
<!ENTITY % relation    "(append)"      >
<!--      einzelne Relationen      -->
<!ELEMENT  append      (arg,arg,arg)      >

<!-- ##### Signatur-unabhaengiger Teil ##### -->

<!--      Beschreibung      -->
<!ELEMENT  descr      (avm|relation;|constants;)+      >
<!--      AVM-Klammerstruktur      -->
<!ELEMENT  avm      (avp*)      >
<!ATTLIST  avm
            top      ID      #REQUIRED
            sort      %sorts;      #IMPLIED      >
<!
            Attribut-Wert-Paar      -->
<!ELEMENT  avp      EMPTY      >
<!ATTLIST  avp
            attr      %attrs;      #REQUIRED
            value      IDREF      #REQUIRED      >
<!--      Argument (fuer relationale Beschr.)      -->
<!ELEMENT  arg      EMPTY      >
<!ATTLIST  arg
            top      IDREF      #REQUIRED      >
<!--      logische Konstanten      -->
<!ENTITY % constant    "(and|or|if|not)"      >
<!--      Konjunktion      -->
<!ELEMENT  and      (descr,descr)      >
<!--      Disjunktion      -->
<!ELEMENT  or      (descr,descr)      >
<!--      Implikation      -->
<!ELEMENT  if      (descr,descr)      >
<!--      Negation      -->
<!ELEMENT  not      (descr)      >

```

Abbildung 15.4.: Eine signatur-spezifische DTD für AVM-Beschreibungen

```

<descr>
<avm sort="phrase">                                <!-- Phrase -->
  <avp attr="synsem" value="id02"/>
  <avp attr="head-dtr" value="id03"/>
</avm>

<avm top="id02">
  <avp attr="head" value="id01"/>
</avm>

<avm top="id01" sort="verb"/>

<avm top="id03">                                    <!-- Kopftochter -->
  <avp attr="synsem" value="id04"/>
</avm>

<avm top="id04">
  <avp attr="head" value="id01"/>
</avm>
</descr>

```

Abbildung 15.5.: XML-Kodierung der AVM in Abb. 15.1

In unserer DTD haben wir die an sich rekursive Gestalt von AVMs aufgebrochen. Dadurch wird die XML-Kodierung flach und spiegelt die Einbettungen von Klammerstrukturen nicht direkt wider. Andererseits können Tags auf einfache Weise als Identifizierer (Werte vom Typ IDREF) kodiert werden. Da Tags in AVMs extrem häufig sind, ist diese Kodierungsform wohl motiviert. Außerdem ermöglicht die flache Kodierung mehr Freiheiten hinsichtlich der graphischen Darstellung von AVMs: die XML-Kodierung in Abb. 15.5 kann graphisch dargestellt werden wie in Abb. 15.1 oder wie in Abb. 15.6.

Abb. 15.7 (S. 166) zeigt die XML-Kodierung der umfangreicheren AVM-Beschreibung von Abb. 15.2. Man beachte den relationalen Ausdruck in der vorletzten Zeile und das Junktorelement `and`.

Alternative 1:
$$\left[\begin{array}{l} \textit{phrase} \\ \text{SYNSEM} \left[\text{HEAD} \boxed{1} \right] \\ \text{HEAD-DTR} \left[\text{SYNSEM} \left[\text{HEAD} \boxed{1} \textit{verb} \right] \right] \end{array} \right]$$

Alternative 2:
$$\left[\begin{array}{l} \textit{phrase} \\ \text{SYNSEM} \left[\text{HEAD} \boxed{1} \textit{verb} \right] \\ \text{HEAD-DTR} \left[\text{SYNSEM} \left[\text{HEAD} \boxed{1} \textit{verb} \right] \right] \end{array} \right]$$

Abbildung 15.6.: Zwei Darstellungsalternativen zu Abb. 15.1

```

<descr>
  <and>
    <descr>
      <avm sort="phrase">
        <!-- Klammerstruktur -->
        <!-- Gesamtphrase -->
        <avp attr="phon" value="id02"/>
        <avp attr="synsem" value="id05"/>
        <avp attr="nonhead-dtr" value="id06"/>
        <avp attr="head-dtr" value="id07"/>
      </avm>
      <avm top="id05"><avp attr="head" value="id01"/></avm>
      <avm top="id01" sort="verb"/>
      <avm top="id06">
        <!-- Nichtkopftochter -->
        <avp attr="phon" value="id03"/>
        <avp attr="synsem" value="id08"/>
      </avm>
      <avm top="id03" sort="nelist">
        <avp attr="hd" value="id09"/>
        <avp attr="tl" value="id10"/>
      </avm>
      <avm top="id09" sort="John"/>
      <avm top="id10" sort="elist"/>
      <avm top="id08">
        <avp attr="head" value="id11"/>
      </avm>
      <avm top="id11" sort="noun"/>
      <avm top="id07">
        <!-- Kopftochter -->
        <avp attr="phon" value="id04"/>
        <avp attr="synsem" value="id12"/>
      </avm>
      <avm top="id04" sort="nelist">
        <avp attr="hd" value="id13"/>
        <avp attr="tl" value="id14"/>
      </avm>
      <avm top="id13" sort="walks"/>
      <avm top="id14" sort="elist"/>
      <avm top="id12">
        <avp attr="head" value="id01"/>
      </avm>
    </descr>
    <descr>
      <!-- relationaler Ausdruck -->
      <append><arg top="id03"/><arg top="id04"/><arg top="id02"/>
    </append>
  </descr>
</and>
</descr>

```

Abbildung 15.7.: XML-Kodierung von Abb. 15.2

15.3. Interaktion mit HPSG-Tools

Die Nützlichkeit der XML-Kodierung ergibt sich aus der Anbindung an bereits existierende HPSG-Werkzeuge. Zur Visualisierung von AVM-Beschreibungen kann aus dem XML-Format automatisch in das Datenformat von Grisu (GRammar vISUalization tool, Wunsch, 2000) konvertiert werden. Grisu erlaubt viele benutzerdefinierte Eingriffsmöglichkeiten in die Darstellungsform einer Klammerstruktur, wie sie bei der Grammatikentwicklung benötigt werden.

Bislang wird Grisu hauptsächlich als Ausgabewerkzeug für die HPSG-Implementierungsplattform Trale (Fouvry und Meurers, 2000, Penn, 2000) verwendet. Um die Ausgabedaten von Trale auch in XML verfügbar zu machen, soll eine Konversion aus dem Grisu-Format in das XML-Format implementiert werden.

Allgemeiner gesehen eignet sich das XML-Format als Datenaustauschformat zwischen verschiedenen Systemen, die AVMs intern in unterschiedlicher Weise repräsentieren. Eine denkbare Anwendung der XML-Kodierung von AVM-Beschreibungen wäre beispielsweise ihre Verwendung beim Erstellen von Grammatik-Testsuiten, die systemunabhängig verwendet werden sollen. Die Testsequenzen können im XML-Format kodiert und dann in das Input-Format des jeweiligen Parsers, zum Beispiel Trale, konvertiert werden. Umgekehrt kann auch aus dem Outputformat von Trale nach XML konvertiert werden. Die hier vorgestellte Kodierung von AVM-Beschreibungen eignet sich damit hervorragend für HPSG-orientiertes NLP.

15.4. Zusammenfassung

In diesem Papier haben wir eine Architektur zur Kodierung von AVM-Beschreibungen vorgestellt und skizziert, wie sie für vorhandene HPSG-Werkzeuge nutzbar gemacht werden kann. Die Besonderheiten sind dabei:

- Die DTD ist signatur-spezifisch. Dadurch ermöglichen Standard-XML-Parser bereits einen Syntaxcheck, der feststellt, ob eine XML-Kodierung einer wohlgeformten AVM-Beschreibung entspricht.
- Die AVM-Beschreibungen beinhalten neben Klammerstrukturen auch Junktoren und relationale Ausdrücke. Damit gibt es eine vollständige XML-Entsprechung zu der in der Literatur als lingua franca der HPSG betrachteten AVM-Beschreibungssprache.
- Die Kodierung von Klammerstrukturen ist nicht-rekursiv. Damit ist sie weitgehend unabhängig von Präferenzen der graphischen Darstellung.

Literaturverzeichnis

BONHOMME, P. UND LOPEZ, P. (1999): "TagML, version 0.4, Wed Nov 24 1999". Online verfügbar: http://www.loria.fr/~lopez/TAG_XML/donnees.html.

CARPENTER, B. UND PENN, G. (1996): "Compiling typed attribute-value logic grammars". In: *Recent Advances in Parsing Technologies*, herausgegeben von Bunt, H. und Tomita, M., Kluwer, S. 145–168.

FLICKINGER, D.; COPESTAKE, A. UND SAG, I. A. (2000): "HPSG Analysis of English". In: *VerbMobil: Foundations of Speech-to-Speech Translation*, herausgegeben von Wahlster, W., Berlin, Heidelberg, New York: Springer, Artificial Intelligence, S. 255–265.

FOUVRY, F. UND MEURERS, D. (2000): "Towards a platform for linearization grammars". In: *On-line Proceedings of the Workshop on Linguistic Theory and Grammar Implementation August 14–18, 2000*. Online verfügbar: <http://www.sfs.nphil.uni-tuebingen.de/~dm/events/ess11i00/proceedings/>.

GÖTZ, T. UND MEURERS, W. D. (1998): "The importance of being lazy – Using lazy evaluation to process queries to HPSG grammars". In: *Lexical and Constructional Aspects of Linguistic Explanation*, herausgegeben von Webelhuth, G.; Koenig, J.-P. und Kathol, A., Stanford: CSLI.

HINRICHS, E. W.; BARTELS, J.; KAWATA, Y.; KORDONI, V. UND TELLJOHANN, H. (2000a): "The verbmobil treebanks". In: *Proceedings of KONVENS 2000*.

HINRICHS, E. W.; BARTELS, J.; KAWATA, Y.; KORDONI, V. UND TELLJOHANN, H. (2000b): "The Tübingen Treebanks for Spoken German, English, and Japanese". In: *Verbobil: Foundations of Speech-to-Speech Translation*, herausgegeben von Wahlster, W., Berlin, Heidelberg, New York: Springer-Verlag, Artificial Intelligence, S. 552–576.

KORDONI, V. (2000): "Stylebook for the English Treebank in VERBMOBIL". Technischer Bericht 241, Verbomobil.

LKB (2000): "The (new) LKB system". System entwickelt von A. Copestake, J. Carroll, R. Malouf, S. Oepen und anderen. Dokumentation verfügbar unter <http://www-csli.stanford.edu/~aac/doc5-2a.pdf>.

MARCINIAK, M.; MYKOWIECKA, A.; KUPŚĆ, A. UND PRZEPIÓRKOWSKI, A. (2000): "An HPSG-annotated test suite for Polish". In: *Proceedings of the Linguistic Resources and Evaluation Conference*.

MENGEL, A. UND LEZIUS, W. (2000): "An XML-based encoding format for syntactically annotated corpora". In: *Proceedings of LREC 2000*. Athen, S. 121–126. Online verfügbar: <http://www.ims.uni-stuttgart.de/projekte/TIGER/paper/lrec200.ps.gz>.

PENN, GERALD (2000): "Applying constraint handling rules to HPSG". In: *Proceedings of the First International Conference on Computational Logic (CL2000), Workshop on Rule-Based Constraint Reasoning and Programming*.

POLLARD, C. UND SAG, I. A. (1994): *Head-Driven Phrase Structure Grammar*. University of Chicago Press.

RICHTER, F. (2000): "A Mathematical Formalism for Linguistic Theories with an Application in Head-Driven Phrase Structure Grammar". Dissertation. Eberhard-Karls-Universität Tübingen.

WUNSCH, H. (2000): "Grisu: Grammar visualization tool". Unveröffentlichter Kode. Eberhard-Karls-Universität Tübingen.

XML (1998): "Extensible Markup Language (XML) Version 1.0". Empfehlung, 2.10.1998, World Wide Web Consortium, Online verfügbar: <http://www.w3.org/TR/1998/REC-xml-19980210>.

XSL (2000): "Extensible Stylesheet Language (XSL) Version 1.0". Working Draft, 27.3.2000, World Wide Web Consortium, Online verfügbar: <http://www.w3.org/TR/2000/WD-xsl-20000327>.