

# A SINGLE-STAGE TOP-DOWN PROBABILISTIC APPROACH TOWARDS UNDERSTANDING SPOKEN AND HANDWRITTEN MATHEMATICAL FORMULAS

Jörg Hunsinger and Manfred Lang

Institute for Human-Machine Communication,  
 Technical University of Munich, D-80290 Munich, Germany  
 {hunsinger, lang}@ei.tum.de

## ABSTRACT

We present a novel approach towards a multimodal analysis of natural speech and handwriting input for entering mathematical expressions into a computer. It utilizes an integrated, multi-level probabilistic architecture with a joint semantic and two distinct syntactic models describing speech and script properties, respectively. Compared to classical multistage solutions our single-stage strategy benefits from an implicit transfer of higher level contextual information into the lower level segmentation and pattern recognition processes involved. For visualization and postprocessing purposes, a transformation into Adobe® FrameMaker® documents is performed.

Fully spoken or handwritten realistic formulas were examined, yielding a structural recognition accuracy of 61.1 % for speech (speaker independent) and 83.3 % for handwriting (writer dependent).

## 1. INTRODUCTION

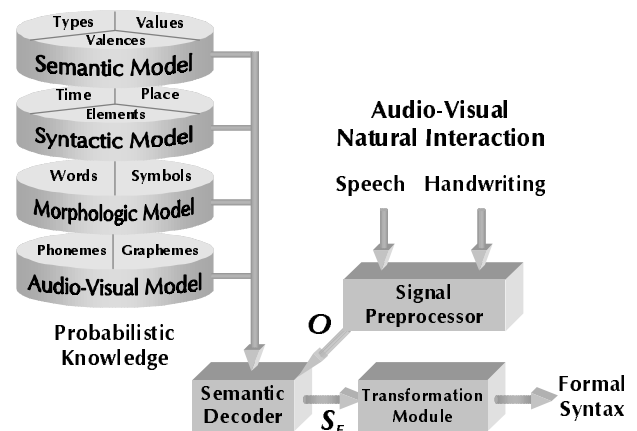
Electronic acquisition of mathematical formulas via conventional tools is a time consuming and complicated task. Therefore it is essential to exploit the capabilities of natural, especially speech and handwriting interaction, both being the fastest and most intuitive channels for registering mathematical expressions [1]. In order to facilitate future data fusion techniques and for uniformity reasons, it makes sense to use common semantic and syntactic representation formalisms for both modalities which may be integrated into a generalized input parsing mechanism. Our context free grammar implementation via probabilistic network structures in conjunction with an extended Earley-type top-down chart parser fulfills these requirements. Fig. 1 shows an overview of the current system components on the different abstraction levels. The following section gives an outline of the applied system architecture.

## 2. SYSTEM OUTLINE

### 2.1 Grammar

The syntactic-semantic attributes of spoken and handwritten mathematical formulas are represented by the parameters of a so-called **Multimodal Probabilistic Grammar**. It combines properties of context free phrase structure grammars with those of graph grammars by allowing for word-type, symbol-type, and position-type terminals. Formally, it is defined by a Chom-

sky set  $G = \langle \Sigma, V, T, P \rangle$  including a start symbol  $\Sigma$ , a set of variables  $V$ , a set of terminals  $T$ , and a set of context free production rules  $P$  [2]. All the production rules are associated with statistical weights obtained from authentic spoken or handwritten training corpora, respectively.



**Figure 1:** System overview. The top-down probabilistic semantic decoder derives a recognized semantic representation  $S_E$  from a preprocessed observation sequence  $O$ . The result is transformed to a formal mathematical description language to be fed into a conventional formula editor.

### 2.2 Semantic Representation

The grammar is implemented into a single stage semantic decoder by means of a compact semantic representation called **Semantic Structure**  $S$  [3]. It is given by an  $N$ -fold hierarchically structured combination out of a predefined inventory of **semuns**  $s$  (semantic units) with corresponding types  $t$ , values  $v$ , and successor attributes, every unit referring to a certain mathematical operator or operand [1]:

$$S = \{s_n\}, 1 \leq n \leq N; \quad s_n = \{t, v, X(t)\}, X \geq 1, \quad (2)$$

where  $X(t)$  denotes the type specific semantic valence, i.e. the number of successor semuns.

The probabilistic nature of this semantic representation is incorporated by three types of statistical weights:

- **root probabilities**  $\eta_0 = P(s_1.t)$  (3)

assigned to every occurring root semun type  $t$  (note: henceforth, the notation  $a.b$  identifies a property  $b$  of the entity  $a$ ),

- **value probabilities**  $\varepsilon_n = P(s_n.v | s_n.t)$  (4)

assigned to every existing semantic value  $v$  of a given semun type  $t$ , and

- **successor probabilities**  $\eta_n = P(r(s_n) | s_n.t)$  (5)

assigned to every allowed combination  $r$  of successor semun types.

Further details of the semantic formalism may be found in [4].

## 2.3 Syntactic Representation

On the syntactic level every semun of a given semantic hypothesis is assigned to a so-called **Syntactic Module**  $SM$ . It consists of an advanced transition network which enables two distinct stochastic processes: 1) transitions from one node to another and 2) emissions of **elements** (i.e. spoken words or handwritten symbols) or local **offsets** between associated symbols or symbol groups. Transitions are responsible for modeling speaking and writing order, whereas emissions account for varying word, symbol, or position choice, respectively. All the  $SM$ s belonging to a complete Semantic Structure form an interconnected so-called **Syntactic Network**  $SN$  which is constituted as follows:

$$SN = \{SM_n\}, 1 \leq n \leq N; \quad (6)$$

$$SM_n = \{St_n, E_n, A_{n1}, \dots, A_{nX}, B_{n1}, \dots, B_{nY}, C_n\}$$

Every  $SM$  is opened via its start node  $St$  and closed via its end node  $E$ . Successor  $SM$ s are connected to their parent  $SM$  via  $X$  individual  $A$  nodes (also called successor nodes), so that the hierarchy of the corresponding Semantic Structure is mapped to the  $SN$ . These  $A$  nodes are also responsible for emitting pairwise positional offset vectors  $\vec{o}$  between all the handwritten symbols or symbol groups belonging to the  $SN$  subbranches connected to them (see below). Further, there is a type specific number  $Y(t) \geq 0$  of  $B$  nodes, each emitting one significant element  $e^+$  of the total input. In case of speech these elements represent spoken words, in case of script they correspond to handwritten symbols. Optionally, a single  $C$  node is entered which emits an insignificant element  $e^-$ . This feature is especially used to model expletive spoken phrases, whereas usually no insignificant elements are found on the syntactic level of handwritten input.

The following three types of probabilistic parameters are needed for a statistical rating of the syntactic contribution to an overall hypothesis score:

- Different paths through a given  $SM_n$  are statistically weighted by means of matrices of **transition probabilities**

$$\Delta_n = \left[ \delta_n^{ij} = P(i \rightarrow j | s_n.t, i) \right], \quad i, j \in SM_n, \quad (7)$$

where  $i \rightarrow j$  denotes any allowed transition from node  $i$  to  $j$ . Additionally, every single  $SM$  node must be passed exactly once – except for the optional  $C$  node – before the end node is reached.

- The type specific **offset emission probabilities**

$$\alpha_{nk} = P(A_{nk} \rightarrow \{\vec{o}\} | s_n.t), \quad 1 \leq k \leq X \quad (8)$$

model the statistical weight for the emission of a set of offset vectors covering the positional relations between all the elements emitted inside the  $SN$  subbranch connected to node  $A_{nk}$ . Since this procedure is performed recursively every time a successor  $SM$  is closed by returning to the next higher  $SM$ 's  $A$  node, a complete pairwise rating of every symbol's position relative to all its syntactic-semantic predecessor symbols is guaranteed. The pairwise offset definition for two consecutive elements is illustrated in Fig.2.

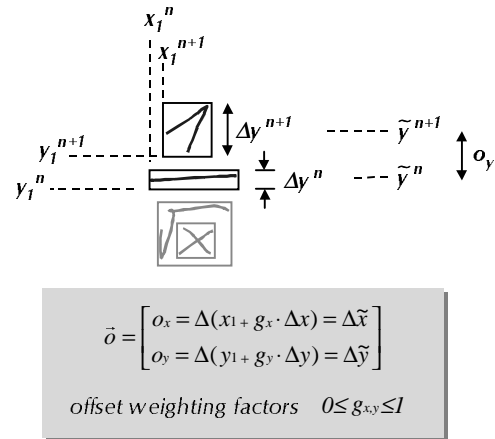
- The so-called **element emission probabilities**

$$\beta_{nl} = P(B_{nl} \rightarrow e^+ | s_n.t, s_n.v), \quad 1 \leq l \leq Y \quad (9)$$

$$\gamma_n = P(C_n \rightarrow e^- | s_n.t) \quad (10)$$

account for a statistical rating of significant (type and value specific) or insignificant (type specific) element emissions, respectively. For consistency reasons,  $\gamma_n$  is set to unity if the corresponding  $C$  node is not passed.

All transition and emission probabilities as well as semantic root, value, and successor probabilities were estimated from training corpora obtained from separate speech and handwriting usability tests (cf. section 2.8). As a summary, a schematic view of a general Syntactic Module with all its attributes is displayed in Fig. 3.



**Figure 2:** Offset vector calculation based on surrounding rectangles. In this example the position of a single handwritten symbol belonging to  $SM_{n+1}$  is charged against that of another single symbol belonging to the predecessor  $SM_n$ . The type specific weighting factors  $g_{x,y}$  account for special constraints due to handwriting conventions.



patterns were derived from this corpus by means of our incidence based iterative training algorithms. We substantially reduced preparative efforts such as manual symbol segmentation and syntactic-semantic annotation by implementing a novel graphical analyzing environment called *StrokeTool* with a universal pen gesture based interface.

The positional parameters (eq. (8)) were estimated by calculating type specific two-dimensional Gaussians over all occurring pairwise symbol (or symbol group) offset vectors. Since only first order dependencies are considered in our approach, any successor semun subbranch was handled as an entity with a unified surrounding rectangle. However, every individual symbol position is included in the resulting parameter set due to recursive processing.

After refining our present position analysis technique, the handwriting knowledge bases will be enlarged in order to improve their statistical significance and to achieve writer independence.

### 3. RESULTS & CONCLUSIONS

For evaluation purposes we performed independent test classifications in either modality. The results for fully spoken or handwritten realistic formulas are summarized in Table 1. Since the positional part of our syntactic model (cf. section 2.8) has not yet been extended to the full range of supported mathematical functions [1], its contribution was neglected in this study. We will present the final recognition results in a subsequent paper.

	Recognition Accuracy	
	Speech	Handwriting
Training Corpus Reclassification	<b>76.2 %</b>	<b>87.5 %</b>
Independent Test Classification	<b>61.1 %</b>	<b>83.3 %</b>

**Table 1:** Recognition results. The numbers refer to full formula structural correctness under toleration of mere character confusions.

For the future we wish to support freely interfering speech and handwriting interactions including mutual coreferencing due to deictic wording and pen gesturing. To this end, the use of speech will presumably be focussed to subterm input and error corrections so that we anticipate a robust and approximately realtime forthcoming system performance.

### 4. REFERENCES

1. Hunsinger J. and Lang M., *A Speech Understanding Module for a Multimodal Mathematical Formula Editor*, Proc. ICASSP 2000, Vol. IV, pp. 2413-2416, Istanbul, Turkey, June 2000.
2. Gazdar G., Klein E., Pullum G.K., and Sag I.A., *Generalized Phrase Structure Grammar*, Basil Blackwell, Oxford, England, 1985.
3. Stahl H., Müller J., and Lang M., *Controlling Limited-Domain Applications by Probabilistic Semantic Decoding of Natural Speech*, Proc. ICASSP 97, pp. 1163-1166, Munich, Germany, April 1997.
4. Müller J. and Stahl H., *Speech Understanding and Speech Translation by Maximum a-posteriori Semantic Decoding*, Artificial Intelligence in Engineering, Vol. 13, No. 4, pp. 373-384, Elsevier Science, 1999.
5. Ruske G., Faltlhauser R., and Pfau T., *Extended linear discriminant analysis (ELDA) for speech recognition*, Proc. ICSLP 98, Vol. 3, pp. 1095-1098, Sydney, Australia, November/December 1998.
6. *Intel Recognition Primitives Library Reference Manual*, Intel Corporation, Order number 637785-007, 1998.
7. Winkler H.-J., Lang M., *Symbol Segmentation and Recognition for Understanding Handwritten Mathematical Expressions*, in A.C. Downton, S. Impedovo, editors, *Progress in Handwriting Recognition*, Proc. 5th International Workshop on Frontiers in Handwriting Recognition IWFHR-5, Essex, England, September 1996, pp. 407-412, World Scientific, 1997.
8. Manke S., Finke M., and Waibel A., *NPen++: A Writer Independent, Large Vocabulary On-Line Cursive Handwriting Recognition System*, Proc. ICDAR 95, Montreal, Canada, August 1995.