

Connectionist Model of Anticipation in Visual Worlds

Marshall R. Mayberry, III, Matthew W. Crocker, and Pia Knoeferle

Department of Computational Linguistics,
Saarland University, Saarbrücken, Germany

{martym, crocker, knoferle}@coli.uni-sb.de

Abstract. Recent “visual worlds” studies, wherein researchers study language in context by monitoring eye-movements in a visual scene during sentence processing, have revealed much about the interaction of diverse information sources and the time course of their influence on comprehension. In this study, five experiments that trade off scene context with a variety of linguistic factors are modelled with a Simple Recurrent Network modified to integrate a scene representation with the standard incremental input of a sentence. The results show that the model captures the qualitative behavior observed during the experiments, while retaining the ability to develop the correct interpretation in the absence of visual input.

Introduction

Two prevalent theories of language acquisition. One view emphasizes syntactic bootstrapping during language acquisition that enable children to learn concepts from mappings between different kinds of information sources [1,2]. The other view emerges from connectionist literature and emphasizes the learning of linguistic structure from purely distributional properties of language usage [3,4]. While the two theories are often taken to be diametrically opposed, both can be seen as crucially dependent on correlations between words and their immediate context, be it the sentence itself or extra-linguistic input, such as a scene.

We combine insights from both distributional and bootstrapping accounts in modeling on-line comprehension of utterances in both the absence and presence of a visual scene. This is an important achievement in at least two regards. First, it emphasizes the complementarity between distributional and bootstrapping approaches—discovering correlations across linguistic and scene contexts [5]. Further, it is an important first step in integrating distributed models of on-line utterance comprehension more tightly to accounts of language acquisition, thus emphasizing the continuity of language processing.

We present results from two simulations on a Simple Recurrent Network (SRN; [3]). The modification of the network to integrate input from a scene together with the characteristic incremental processing of such networks allowed us to model people’s ability to flexibly use the contextual information in order to more rapidly interpret and disambiguate a sentence. The model draws on recent studies that appeal to theories of language acquisition to account for the comprehension of scene-related utterances [6,7].

rapid coordinated interaction of information from the immediate scene, and world knowledge plays a major role in incremental and anticipatory comprehension.

Simulation 1

In Simulation 1, we simultaneously model four experiments that show the rapid influence of linguistic and world knowledge as well as scene information on utterance comprehension. All experiments were conducted in German, a language that allows both subject-verb-object (SVO) and object-verb-subject (OVS) word orders. In the face of word order ambiguity, case marking indicates the subject and object grammatical function, except in the case of feminine and neuter noun phrases where the article does not distinguish the nominative and accusative cases.

Anticipation Depending on Stereotypicality

Two experiments that we modeled examined how linguistic and world knowledge enabled rapid thematic role assignment in unambiguous sentences by determining who-does-what-to-whom in a scene.

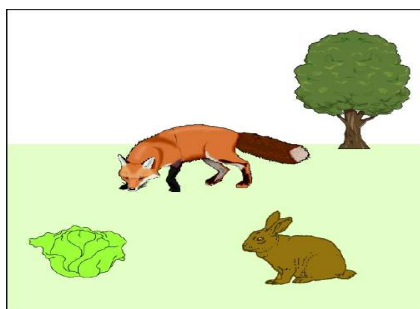


Fig. 1. Selectional Restrictions

Experiment 1: Morphosyntactic and lexical verb information. To examine the influence of morphosyntactic and verb plausibility on thematic role assignment, [8] presented participants with utterances such as (1) or (2) that described a scene showing a hare, a cabbage, a fox, and a tree (see Figure 1) :

(1) *Das Kaninchen frisst gleich den Kohl.*

The hare_{nom} eats shortly the cabbage_{acc}.

(2) *Das Kaninchen frisst gleich der Fuchs.*

The hare_{acc} eats shortly the fox_{nom}.

In Experiment 1, participants heard the sentence “The hare_{nom} eats ...” and “The hare_{acc} eats ...”, people made anticipatory thematic role assignments. This is because the hare is stereotypically associated with eating.

Experiment 2: Verb type information. To further investigate the role of verb information, the authors replaced the agent/patient verbs like *frisst* (“eats”) with experiential verbs like *interessiert* (“interests”). This manipulation interchanged agent (subject) and patient (theme) roles from Experiment 1. For Figure 1 and the subject-object-first sentence (4), participants showed gaze fixations complementary to Experiment 1, confirming that both case and semantic verb information are needed to predict relevant role fillers.

Prinzessin interessiert ganz besonders den Fuchs.
Prinzessin_{nom} interests especially the fox_{acc}.
Der Fuchs interessiert ganz besonders der Kohl.
Fuchs_{acc} interests especially the cabbage_{nom}.

Experiment 3: Anticipation Depending on Depicted Events

This set of experiments investigated whether depicted events showing who-does-what to whom can establish a scene character’s role as agent or patient when syntactic role relations are temporarily ambiguous in the utterance.

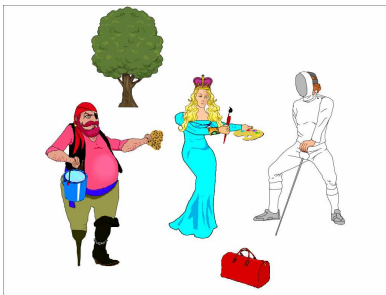


Fig. 2. Depicted Events

Experiment 3: Verb-mediated depicted role relations. [9] presented such initially ambiguous spoken SVO (5) and OVS sentences (6) together with a scene in which a princess paints a fencer and is washed by a pirate (Figure 2):

Prinzessin malt offensichtlich den Fechter.
Prinzessin_{nom} paints obviously the fencer_{acc}.
Prinzessin wäscht offensichtlich der Pirat.
Prinzessin_{acc} washes obviously the pirate_{nom}.

Disambiguation occurred on the second NP; disambiguation prior to the second NP is only possible through use of the depicted events. When the verb identified the agent, the second NP was the patient. When the verb identified the patient, the second NP was the agent.

R. Mayberry, III, M.W. Crocker, and P. Knoeferle

ence of verb-mediated depicted events on the assignment of thematic roles to a sentence-initial noun phrase.

Experiment 4: Weak temporal adverb constraint. [9] also investigated German verb-ambiguous (7) and passive (8) constructions. In this type of sentence, the initial subject is role-ambiguous, and the auxiliary *wird* can have a passive or future reading.

Prinzessin wird sogleich den Pirat waschen.

Prinzessin_{nom} will right away wash the pirate_{acc}.

Prinzessin wird soeben von dem Fechter gemalt.

Prinzessin_{acc} is just now painted by the fencer_{nom}.

Early linguistic disambiguation, temporal adverbs biased the auxiliary *wird* toward the future (“will”) or passive (“is -ed”) reading. Since the verb was sentence-initial, the interplay of scene and linguistic cues (e.g., temporal adverbs) were rather more salient. When the listener heard a future-biased adverb such as *sogleich*, after the auxiliary *wird*, he interpreted the initial NP as agent of a future active construction, as evidenced by anticipatory eye-movements to the patient in the scene. Conversely, listeners interpreted the passive-biased construction *soeben* with these roles exchanged.

Architecture

The Recurrent Network is a type of neural network typically used to process sequences of patterns such as words in a sentence. A common approach is to train the network on prespecified targets, such as verbs and their arguments, that represent what the network is expected to produce upon completing a sentence. Processing is incremental, with each new input word interpreted in the context of the sentence processed so far, represented by a copy of the previous hidden layer as additional input or *context* to the current hidden layer. Because these connectionist models automatically develop correlations among the data they are trained on, they will generally develop expectations about the output even before the sentence is completed because sufficient information occurs early in the sentence to make such predictions. Moreover, during the course of processing a sentence these expectations can be overridden with subsequent input, often abruptly revising an interpretation in a manner reminiscent of how humans seem to process language. Indeed, these characteristics of incremental processing, the automatic development of expectations, seamless integration of multiple sources of information, and nonmonotonic processing have endeared neural network models to cognitive researchers.

In Experiment 1, the four experiments described above have been modelled simultaneously using a single network. The goal of modelling all experimental results by a single architecture required enhancements to the SRN, the development and presentation of training data, as well as the training regime itself. We describe these next.

In the first of the experiments, only three characters are depicted, the representation of

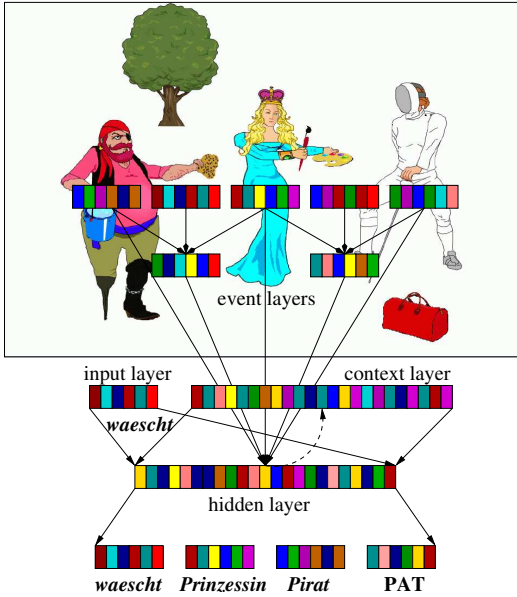


Fig. 3. Scene Integration

in both events, either as an agent or a patient (e.g., **princess**). Only one of the events, however, corresponded to the spoken linguistic input.

The integration of this scene information and its integration into the model's architecture was the primary modification to the SRN. Connections between representations of the depicted characters and the hidden layer were provided. Encoding of the events, when present, required additional links from the characters and actions to **event** layers, and links from these event layers to the SRN's hidden layer. Representations for the events were developed in the event layers by compressing the representations of the involved characters and depicted actions through weights corresponding to the action, its agent and its patient for each event. This event encoding was kept simple and only provided conceptual input to the hidden layer: that to whom was encoded for both events, when depicted; richer grammatical information (e.g., case and gender on articles) only came from the linguistic input.

Standard networks will usually encode any correlations in the data that help to minimize error. In order to prevent the network from encoding regularities in its weights corresponding to the position of the characters and events given in the scene (such as, for example, the central character in the scene corresponds to the first NP in the presented sentence), which are not relevant to the role-assignment task, one set of weights was used for the representations of the characters, and another set of weights used for both events. This weight-sharing meant that the network had to access the information encoded in the event layers, or

put assemblies were the scene representations and the current word from the sentence. The output assemblies were the verb, the first and second nouns, and an assembly that indicated whether the first noun was the agent or patient of the sentence (see Figure 3). Typically, agent and patient assemblies would be fixed in a particular representation without such a discriminator, and the model required to learn to assemble them correctly [10]. However, we found that the model performed much better when the task was recast as having to learn to isolate the nouns in the order in which they are introduced, and separately mark how those nouns relate to the verb. The input assemblies had 100 units each, the event layers contained 200 units each, and the hidden and context layers consisted of 400 units.

Experiment 1: Data, Training, and Experiments

In order for the network to correctly handle sentences involving non-stereotypical events as well as stereotypical ones, both when visual context was present and when it was absent, over half a billion sentence/scene combinations were possible for all of the experiments. To generate training materials, we adopted a grammar-based approach to randomly generate sentences based on the materials from each experiment while holding out the actual materials to be used for testing. Because of the complementary roles that stereotypical and non-stereotypical events play in the two sets of experiments, there was virtually no lexical overlap between the two sets of materials. In order to accurately model the first two experiments involving selectional restrictions on verbs, two additional words were added to the lexicon for each character used by a verb. For example, in the sentence *Der Hase frisst gleich den Kohl*, *Hase1*, *Hase2*, *Kohl1*, and *Kohl2* were used to develop training sentences. *Hase1* and *Hase2* were meant to represent, for example, words such as “rabbit” and “jackrabbit” or “hare” and “cabbage”. With these extra tokens the network could learn that *frisst*, and *Kohl* were correlated without ever encountering all three words in the training sentence. The experiments involving non-stereotypicality did not have this constraint, so training sentences were simply generated to avoid presenting identical items.

In order to facilitate modelling, several standard simplifications to the words have been made. First, multiple, multi-word adverbs such as *fast immer* were treated as one word through the experiments so that sentence length within a given experimental set up is maintained. Second, case markings such as *-n* in *Hasen* were removed to avoid sparse data as these are idiosyncratic, and the case markings on the determiners are more infrequent overall. More importantly, morphemes such as the infinitive marker *-en* and the prefix *ge-* were removed, because, for example, the verb forms *malt*, *malen*, *gemalt*, and *gemalt* would all be treated as unrelated tokens, again contributing unnecessarily to the problem with sparse data. The result is that one verb form is used, and to participate in the network, the network must rely on its position in the sentence (either second or third), as well as whether the word *von* occurs to indicate a participial reading or a non-infinitival. All 326 words in the lexicon for the first four experiments were

is against the held-out test materials for each of the five experiments. Scenes provided half of the time to provide an unbiased approximation to linguistic experience. The network was initialized with weights between -0.01 and 0.01. The learning rate was initially set to 0.05 and gradually reduced to 0.002 over the course of 15000 trials. All four splits took a little less than two weeks to complete on 1.6Ghz PCs.

Results

Table 1 reports the percentage of targets at the network's output layer that the model correctly matches, both as measured at the adverb and at the end of the sentence. The table clearly demonstrates the qualitative behavior observed in all four experiments: the model is able to access the information either from the encoded scene or stereotypicality to disambiguate the sentence and to anticipate forthcoming targets.

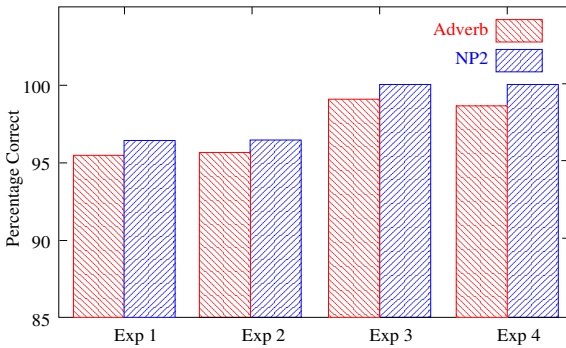


Fig. 4. Results

In the two studies using stereotypical information (experiments 1 and 2), the network achieved just over 96% at sentence end, and anticipation accuracy was just over 96% at the adverb. Because these sentences are unambiguous, the model is able to identify the role of the upcoming argument, but makes errors in token identification for confusing words that are within the selectionally restricted set, such as, for example, *Kohl* and *Kohl2*. Thus, the model has not quite mastered the stereotypical information, particularly as it relates to the presence of the scene.

In the other two experiments using non-stereotypical characters and depicted events (experiments 3 and 4), accuracy was 100% at the end of the sentence. More importantly, the model achieved over 98% early disambiguation on experiment 3, where the sentences were simple, active SVO and OVS. Early disambiguation on experiment 4 was somewhat harder because the adverb is the disambiguating point in the sentence, and it occurred later in the sentence than in the other three experiments. As nonlinear dynamical systems are sensitive to initial conditions, the model's performance is highly sensitive to the initial state of the network.

ifference in performance between the first two experiments and second two experiments can be attributed to the event layer that was only available in experiments. Closer inspection of the model's behavior during processing revealed that finer information was encoded in the links between the event layers and hidden layer than was encoded in the weights between the characters and the hidden layer.

Experiment 2

A series of experiments demonstrated the rapid use of either linguistic knowledge or scene information to anticipate forthcoming arguments in a sentence. A further question is the relative importance of these two informational sources when they conflict. We first review an experimental study by [6] designed to address this issue and report relevant modelling results.



Fig. 5. Scene vs Stored Knowledge

Stored Knowledge. One goal of the study by [6] was to verify that stored knowledge about non-depicted events and information from depicted, but non-stereotypical agents each enable rapid thematic interpretation. Case-marking on the first NP identified the pilot as a patient. After hearing the verb in (9) more inspections to food-serving agent (detective) than to the other agent showed the influence of scene information. In contrast, when people heard the verb in condition two (10), a higher number of anticipatory eye-movements to the only stereotypical agent (wizard) than to the other agent revealed the influence of stereotypical knowledge (see Figure 5).

Piloten verköstigt gleich der Detektiv.

pilot_{acc} serves-food-to shortly the detective_{nom}.

Piloten verzaubert gleich der Zauberer.

pilot_{acc} jinxes shortly the wizard_{nom}.

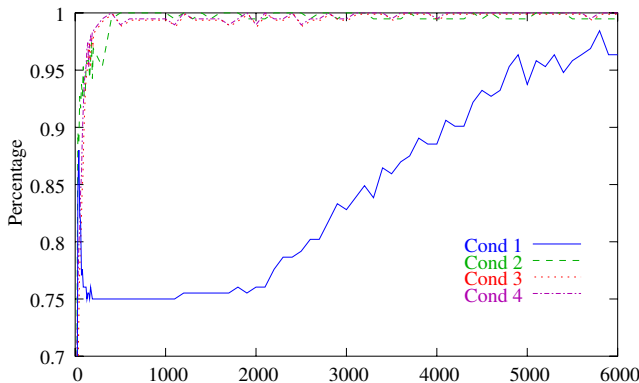
ical (detective) and a depicted agent (wizard). In this case, people preferred to see the immediate event depictions over stereotypical knowledge, and looked more at the wizard, the agent of the depicted event, than at the other, stereotypical agent (the detective).

Piloten bespitzelt gleich der Zauberer.
 pilot_{acc} spies-on shortly the wizard_{nom}.
Piloten bespitzelt gleich der Detektiv.
 pilot_{acc} spies-on shortly the detective_{nom}.

Architecture, Data, Training, and Results

In simulation 1, we modelled experiments that depended on stereotypicality or depicted agent but not both. The experiment modelled in simulation 2, however, was specifically designed to investigate how these two information sources interacted. Accordingly, the network needed to learn to use either information from the scene or stereotypical knowledge when available, and, moreover, favor the scene when the two sources conflicted, as observed in the empirical results. Recall that the network is trained only on the final word of a sentence. Thus, capturing the observed behavior required manipulating the relative frequencies of the four conditions described above during training. In order to encourage the network to develop stereotypical agents for verbs, the frequency that a verb is used with its stereotypical agent, such as *Detektiv* and *bespitzelt* from example (12) needed to be greater than for a non-stereotypical agent. However, the frequency needed to be so great as to override the influence from the scene.

The training resolution we adopted is motivated by theories of language acquisition that take into account the importance of early linguistic experience in a visual environment (see the next section for a detailed Discussion). We found a small range of frequencies that permitted the network to develop an early reliance on the information from the scene while it gradually learned to overcome the stereotypical associations. Figure 6 shows the effect this training regime had over time (epochs) on the ability of the network to accurately anticipate the missing argu-



R. Mayberry, III, M.W. Crocker, and P. Knoeferle

ch of the four conditions described above when the ratio of non-stereotypical typical sentences was 8:1. The network quickly learns to use the scene for condition 4 (examples 10-12), where the action in the linguistic input stream is also allowing the network to determine the relevant event and deduce the missing argument. (Because the graph shows the accuracy of the network at anticipating the argument at the adverb, the lines for conditions 3 and 4 are, in fact, identical.) Condition 1 (sentence 9) requires only stereotypical knowledge. The accuracy on condition 1 remains close to 75% (correctly producing the verb, first NP, and role) until around epoch 1800 or so and then gradually as the network learns the appropriate stereotypical associations.

Results from several separate runs with different training parameters (such as learning rate and stereotypicality ratio) show that the network does indeed model the observed experimental behavior. The best results thus far exceed 99% accuracy in anticipating the proper roles and 100% accuracy at the end of sentence.

In simulation 1, the training corpus was generated by exhaustively combining verbs and actions for all experimental conditions while holding out all test sentences. However, we found that we were able to use a larger learning rate, 0.1, than 0.05 in the first simulation.

The analysis of the network after successful training suggests why this training policy works. Early in training, before stereotypicality has been encoded in the network's weights, patterns are developed in the hidden layer once the verb is read in from the input that enable the network to accurately decode that verb in the output layer. Notably, the network uses these same patterns to encode the stereotypical agent; the constraint for the network is to ensure that the scene can still override this stereotype when the depicted event so dictates.

General Discussion and Future Work

This study demonstrates that reliance on correlations from distributional information in the linguistic input and the scene during training of the model enabled successful modeling of on-line utterance comprehension both in the presence and absence of rich visual contexts. The model that we present acquires stereotypical knowledge from distributional properties of language during training. The mapping from words to the scene is established through cooccurrence of scene-related utterances and depicted events during training. The network that emerges from this training regime successfully models five *visual worlds* eye-tracking experiments in two simulations. A first set of four experiments models the influence of either thematic and syntactic information in the utterance [8], or of depicted events showing who-does-what-to-whom on experimental thematic role assignment [9]. Crucially in modelling the fifth experiment we are able to account for the greater relative priority of depicted events when scene information and event knowledge conflict with each other.

The simple accuracy results belie the complexity of the task in both simulations. For

stic stream when available. This task is rendered more difficult because the event must be extracted from the superimposition of the two events in the which is what is propagated into the model's hidden layer. In addition, it must be able to process all sentences correctly when the scene is not present.

ation 2 is more challenging still. The experiment shows that information from takes precedence when there is a conflict with stereotypical knowledge; oth- which source of knowledge is used when it is available. In the training regime is simulation, the dominance of the scene is established early because it is e frequent than the more particular stereotypical knowledge. As training pro- ereotypical knowledge is gradually learned because it is sufficiently frequent work to capture the relevant associations. As the network weights gradually t becomes more difficult to retune them. But encoding stereotypical knowl- ires far fewer weight adjustments, so the network is able to learn that task g training.

ding to the "Coordinated Interplay" account in [7,6,11], the rapid integration and utterance information and the observed preferred reliance of the compre- system on the visual context over stored knowledge might best be explained ing to bootstrapping accounts of language acquisition. The development of world knowledge occurs in a visual environment, which accordingly plays a role during language acquisition. The fact that the child can draw on two inal sources (utterance and scene) enables it to infer information that it has not ed from what it already knows. Bootstrapping accounts for the fact that a child ate event structure from the world around it with descriptions of events. When ceives an event, the structural information it extracts from it can determine hild interprets a sentence that describes the event in question. The incremental ion of a sentence can in turn direct the child's attention to relevant entities s in the environment. Events are only present for a limited time when utter- r to such events during child language acquisition. This time-limited pres- t determine the tight coordination with which attention in the scene interacts ance comprehension and information extracted from the scene during adult comprehension. This contextual development may have shaped both our cog- nition (i.e., providing for rapid, seamless integration of scene and linguistic on), and comprehension mechanisms (e.g., people rapidly avail themselves of on from the immediate scene when the utterance identifies it).

odel presented in this paper extends current models of on-line utterance com- n when utterances relate to a scene [12] in several ways. Existing models ac- processes of establishing reference in scene-sentence integration when scenes ly objects. Our network accounts for processes of establishing reference, and re models the rapid assignment of thematic roles based on linguistic and world e, as well as scene events. In this way, it achieves rapid scene-utterance inte- r increasingly rich visual contexts, including the construction of propositional tions on the basis of scene events. It models the integration of utterances vely rich scenes (that contain actions and events) in addition to objects. Fur-

R. Mayberry, III, M.W. Crocker, and P. Knoeferle

gh a modification of the training regime that prioritizes scene information.irms suggestions from [7] that a rapid interplay between utterance compreand the immediate scene context during acquisition is one potential cause for e priority of depicted events during on-line comprehension.

ctionist models such as the SRN have been used to model aspects of cogni-ppment, including the time-course of emergent behaviors [13], making them table for simulating developmental stages in child language acquisition (e.g., ng names of objects in the immediate scene, and later proceeding to the acqui-tereotypical knowledge). The finding that modelling this aspect of developides an efficient way to naturally reproduce the observed adult comprehension promises to offer deeper insight into how adult performance is at least partially ence of the acquisition process.

research will focus on combining all of the experiments in one model, and e range of sentence types and fillers to which the network is exposed. The re itself is being redesigned to scale up to much more complex linguistic con-and have greater coverage while retaining the cognitively plausible behavior in this study [14].

Acknowledgements

Two authors were supported by SFB 378 (project “ALPHA”), and the third au-PhD studentship (GRK 715), all awarded by the German Research Foundation

References

Pinker. How could a child use verb syntax to learn verb semantics? In Lila Gleitman Barbara Landau, editors, *The acquisition of the lexicon*, pages 377–410. MIT Press, Cambridge, MA, 1994.

na Fisher, D. G. Hall, S. Rakowitz, and Lila Gleitman. When it is better to receive give: Syntactic and conceptual constraints on vocabulary growth. In Lila Gleitman Barbara Landau, editors, *The acquisition of the lexicon*, pages 333–375. MIT Press, Cambridge, MA, 1994.

L. Elman. Finding structure in time. *Cognitive Science*, 14:179–211, 1990.

Redington, Nick Chater, and Steven Finch. Distributional information: A powerful acquiring syntactic categories. *Cognitive Science*, 22:425–469, 1998.

oy and Alex Pentland. Learning words from sights and sounds: A computational *Cognitive Science*, 26(1):113–146, 2002.

noeferle and Matthew W. Crocker. Stored knowledge versus depicted events: what auditory sentence comprehension. In *Proceedings of the 26th Annual Conference of Cognitive Science Society*. Mahawah, NJ: Erlbaum, 2004. 714–719.

noeferle and Matthew W. Crocker. The coordinated interplay of scene, utterance, and knowledge: evidence from eye-tracking. submitted.

- Hoferle, Matthew W. Crocker, Christoph Scheepers, and Martin J. Pickering. The effect of the immediate visual context on incremental thematic role-assignment: evidence from eye-movements in depicted events. *Cognition*, 95:95–127, 2005.
- Miikkulainen. Natural language processing with subsymbolic neural networks. In Guy Browne, editor, *Neural Network Perspectives on Cognition and Adaptive Robotics*, pages 120–139. Institute of Physics Publishing, Bristol, UK; Philadelphia, PA, 1997.
- Hoferle and Matthew W. Crocker. The coordinated processing of scene and utterance: Evidence from eye-tracking in depicted events. In *Proceedings of International Conference on Cognitive Science*, Allahabad, India, 2004.
- Hoferle and Niloy Mukherjee. Towards situated speech understanding: Visual context priming in language models. *Computer Speech and Language*, 19(2):227–248, 2005.
- David L. Elman, Elizabeth A. Bates, Mark H. Johnson, Annette Karmiloff-Smith, Domenico Borra, and Kim Plunkett. *Rethinking Innateness: A Connectionist Perspective on Development*. MIT Press, Cambridge, MA, 1996.
- Hoferle, R. Mayberry and Matthew W. Crocker. Generating semantic graphs through self-organization. In *Proceedings of the AAAI Symposium on Compositional Connectionism in Cognitive Science*, pages 40–49, Washington, D.C., 2004.