Recursive syntactic pattern learning by songbirds

Timothy Q. Gentner, Kimberly M. Fenn, Daniel Margoliash, and Howard C. Nusbaum

Supplementary Information

Table S1 provides an overview of all the training and testing conditions.

	Condition	Stimulus
1	Baseline Training	A^2B^2 and $(AB)^2$, n=2
2	Transfer 1	Novel A^2B^2 and $(AB)^2$, n=2
3	Probe Test 1	Familiar A ² B ² and (AB) ² , n=2 (80% of trials) Novel A ² B ² and (AB) ² , n=2 (10% of trials) Agrammatical AAAA, BBBB, ABBA, BAAB (10% of trials)
4	Probe Test 2	Familiar A^2B^2 and $(AB)^2$, $n=2$ (80% of trials) Novel A^nB^n and $(AB)^n$, $n=3$, 4 (20% of trials)
5	Probe Test 3	Familiar A ² B ² and (AB) ² , n=2 (80% of trials) A*B* (10% of trials) Novel A ⁿ B ⁿ and (AB) ⁿ , n=2,3,4 (10% of trials)

Table S1. Overview of the training and testing conditions in order to which subjects were exposed. Stimuli marked with (*) comprised speech syllables instead of starling song motifs.

Probe tests with novel A^2B^2 and $(AB)^2$ stimuli

Because the novel transfer stimuli were reinforced with the same response contingencies used in the initial training, the animals' behaviour was changing (i.e. they were learning) as we measured generalization. To examine generalization without contingent reinforcement and thus without additional learning, we tested subjects using a "probe" procedure. The reinforcement regimen during probe sessions allowed us to measure generalization from the classes defined by the familiar CFG and FSG stimuli to novel stimuli without any direct effects of differential reinforcement learning on the latter (see Additional Methods). Following the transfer test (see text), birds were tested by probing with sequences drawn from the same *A/B* vocabularies and CFG/FSG "languages". As in the preceding transfer test, the probe stimuli were novel with respect

to all other grammatical sequences the birds had heard. All subjects accurately classified the novel CFG and FSG probe sequences (Fig. 3b, n=2 examples). The mean d' for all presentations of the probe stimuli was 1.63 ± 0.39 , and the lower bound of the 95% CI around d' was above zero for all subjects (range: 0.07 - 1.96), demonstrating classification well above chance. Accurate classification was apparent in the earliest responses to the novel grammatical stimuli (mean d' over the first five blocks (21 \pm 2 probe stimuli): 4.46 \pm 0.46), and was consistent across the entire probe session (F(3,4) = 0.78, p = 0.5, repeated measures ANOVA for change in d' across sessions). Probe sessions spanned several days during which a total 304 ± 70 (mean ±sem) grammatical probe stimuli were presented to each bird. Thus, subjects had the opportunity to respond (or not) to each exemplar approximately 19 times, which was sufficient to estimate classification behavior reliably. All subjects also maintained classification of the familiar CFG and FSG stimuli during the probe sessions (mean d': 2.32 ± 0.41). The results indicate that birds exhibited accurate classification of these novel probe stimuli throughout all the probe sessions (Fig 3b). This robust generalization to the novel grammatical stimuli is a strong rejection of the rote memorization hypothesis, and provides additional support for the conclusion that subjects learned FSG and CFG pattern information.

Testing finite-state approximations to $A^n B^n$.

The A^*B^* probe condition is crucial to the demonstration that subjects have learned a recursive patterning rule to classify the motif strings because it permits us to rule out the use of many alternative finite-state strategies. One can easily prove that there is no finite-state automaton (FSA) that can correctly identify the non-regular language L= { $A^nB^n : n \ge 0$ }:

Let *M* be any FSA. We will show that *M* does not recognize *L* by finding two strings, $u \in L$ and $v \notin L$, such that *M* ends in the same state reading either *u* or *v*. Imagine that we ask *M* to read strings of a's: a, a^2, a^3, a^4, \ldots , beginning in its start state *s*. Let $q_n = \delta * (s, a^n)$ be its ending state after reading a^n . Since *M* has only a finite number of states, there must be at least two different indices m < n such that $q_m = q_n$. Taking $u = a^m b^m$ and $v = a^n b^m$, we have:

$$\delta * (s, u) = \delta * (s, a^m b^m) = \delta * (q_m, b^m) = \delta * (q_n, b^m) = \delta * (s, a^n b^m) = \delta * (s, v).$$

Thus, *M* either accepts both $u \in L$ and $v \notin L$, or it accepts neither. In either case, *M* is not an FSA for *L*.

Alternate solution strategies

If starlings attend to only the first two motifs when learning the pattern classification, then patterns of the form AAxx and ABxx, where x can be any motif, should be classified similarly to the CFG and FSG, respectively. In other words, the d' value measuring differential classification of the 'AAAA' and 'ABBA' probe stimuli should be similar to that for classification of the novel A^2B^2 and $(AB)^2$ probe stimuli. Alternatively, if subjects attend to only the last two motifs then patterns of the form xxBB and xxAB should be classified similarly to the CFG and FSG, respectively, such that the d' value measuring differential classification of the 'BBBB' and 'BAAB' probe stimuli should be similar to that for classification of the novel A^2B^2 and $(AB)^2$ probe stimuli. Both of these hypotheses are rejected. The mean (± sem) d' value for both the 'AAAA' to 'ABBA' and the 'BBBB' to 'BAAB' comparisons (0.99±0.40 and 0.02±0.26, respectively) was significantly lower than that for the novel A^2B^2 and $(AB)^2$ probe stimuli $(1.63 \pm 0.39; p < 0.05, paired t-test; Fig. 4)$. Thus, neither the primacy nor recency can completely account for the observed classification of novel grammatical patterns. We note that the d' value for the primacy stimuli is significantly greater than that for the recency stimuli (paired t-test, p < 0.05), with the latter at chance, suggesting better classification of the primacy than recency stimuli. Although the d' values associated with the primacy stimuli cannot account for classification of the grammatical probe stimuli, they deserve some consideration. Because subjects were forced to classify an agrammatical stimulus as an FSG or CFG pattern they are likely to adopt a default strategy, if the learned cues for classification (i.e. grammatical patterning) are not present. The primacy probe data suggest that one such default strategy involves a strong bias to attend to the initial motifs in each sequence.

Similar reasoning can be applied to the question of whether subjects attended to the presence or absence of a *B/A* motif transition. The $(AB)^n$ patterns have *n*-1 transitions between *B* and *A* motifs (as presented in time), A^nB^n patterns have none. Among the agrammatical stimuli, *ABBA* and *BAAB* each have one such transition, while *AAAA* and *BBBB* have none. Therefore, if the hypothesized solution strategy is being used, subjects should discriminate *ABBA* and *BAAB* from *AAAA* and *BBBB*, and the *d'* value for the classification of these two stimulus pairs should be similar to that for the novel (n=2) FSG and CFG patterns presented simultaneously. This hypothesis is rejected. The mean (± sem) *d'* for the *B/A* transition pairs (0.51 ± 0.09) was significantly different from that for the novel A^2B^2 and $(AB)^2$ probe stimuli (1.63 ± 0.39; p < 0.05, paired t-test; Fig. 4).

Another possible strategy to solve the FSG/CFG discrimination and generalize to novel sequences is to count the number of transitions between '*A*' and '*B*' motifs in each sequence. The A^nB^n patterns have only one such transition regardless of the value of n, while $(AB)^n$ patterns have n transitions. If subjects count the number of *A*/*B* transitions, then all of the A^*B^* stimuli should be treated similarly to the novel A^nB^n stimuli presented during the same probe sessions (see text) because they all have one *A*/*B* transition. This was not the case. The responses to the A^*B^* stimuli (pooled across all four patterns) were significantly different than the responses to the A^nB^n probe stimuli when n=2, 3, and 4 (X^2 , p < 0.0001, *df* = 3, all cases; see text for individual comparisons).

Finally, if subjects listen for *AA* or *BB* motif pairs (so-called bi-gram strategies) then the following patterns should obtain. If a subject listens for *AA*, then the A^3B , A^2B^3 and A^3B^2 forms of the A^*B^* pattern should be treated the same as the novel A^2B^2 probe stimuli presented in the same sessions. Likewise, if the subject listens for *BB*, then the *AB*³, A^2B^3 and A^3B^2 forms of the A^*B^* pattern should be treated the same as the novel *A2B2* probe stimuli presented in the same sessions. Both hypotheses can be rejected. For each of the four subjects, the pattern of response to the A^2B^2 probe stimuli was significantly different than that for either subset of A^*B^* patterns (X^2 , p < 0.01, df=5, for all 8 comparisons).

In summary, we find no support for putative alternate classification strategies that involve specific and limited attention to either the first or last two motifs in a sequence, counting either *A/B* or *B/A* motif transitions, or attention to *AA* or *BB* motif pairs. We conclude, instead, that subjects learned the patterns defined by the FSG and CFG grammars, and then applied this knowledge to distinguish between novel exemplars derived from the two patterning rules.

Given that our stimulus sets were finite, there must be a finite state grammar that describes them. For example, a language with strings of the form {*aabb, aaabbb, aaaabbbb*} can be generated by the FSG with non-terminal symbols {*C, D, E, F, G, H, J, K, L, M*}, terminal symbols {*a, b*}, the initial symbol {*S*}, and productions {*S*->*aC, C*->*aD, D*->*bE, E*->*b, D*->*aF, F*->*bG, G*->*bH, H*->*b, F*->*aJ, J*->*bK, K*->*bL, L*->*bM, M*->*b*}. We note that this FSG requires learning a total of 13 production rules, only four of which could have been learned during initial operant training. The remaining 11 production rules would therefore need to either generalize from the baseline training or somehow learned along with the four productions that describe A^2B^2 . In contrast, the CFG that produces all strings of the form A^nB^n requires only one non-terminal symbol and one production rule {*S*->*aSb*}, both of which are available during baseline training.

This suggests that using an FSG to solve the classification of $A^n B^n$ strings, while theoretically possible, is less parsimonious than using the CFG.

Individual variation

Of the four birds subjected to extensive testing, one, st218, showed a qualitatively different pattern of responding over the initial agrammatical control tests described in the text. Unlike the other birds, st218 did not treat the novel CFG and agrammatical stimuli in a significantly different way. As noted in the text, this may indicate that the bird learned only one grammar. The failure to find a significant difference in response patterns provides only weak support for such a solution strategy, however, as it does not disprove learning of both patterning rules. In fact, st218's response to one of the four classes of agrammatical sequences was significantly different than the CFG, arguing against a trivial default strategy.

Additional methods

Stimuli. All the starling song motifs were recorded from a single adult male starling. Recording procedures have been described elsewhere¹. From a library of the recorded male's song bouts, we selected 16 different motifs: eight 'rattle' motifs and eight 'warble' motifs (Figs. S1 and S2). These motifs were combined to create the explicit stimuli as described in the text. The $(AB)^n$ and A^nB^n patterns for grammatical stimuli were modelled on those used in an earlier test of grammatical competence in cotton-top tamarins². The complete list of motif sequences for all the stimuli used is given in Table S2.

Condition (pattern	A ⁿ B ⁿ patterns	(AB) ⁿ patterns
order)		
	$a_1 a_3 b_6 b_2$	$a_1 b_6 a_5 b_2$
	$a_2 a_1 b_7 b_5$	$a_2 b_5 a_6 b_7$
	a ₃ a ₄ b ₁ b ₄	$a_3 b_7 a_8 b_3$
Baseline (n=2)	$a_4 a_7 b_3 b_8$	a ₄ b ₃ a ₃ b ₈
	$a_5 a_2 b_5 b_6$	$a_5 b_2 a_2 b_6$
	a ₆ a ₈ b ₈ b ₁	a ₆ b ₄ b ₇ b ₁
	$a_7 a_5 b_2 b_3$	a ₇ b ₁ b ₄ b ₄
	a ₈ a ₆ b ₄ b ₇	a ₈ b ₈ b ₁ b ₅
	$a_1 a_8 b_1 b_3$	$a_1 b_5 a_3 b_3$
	a ₂ a ₄ b ₅ b ₈	a ₂ b ₁ a ₄ b ₆
	a ₃ a ₆ b ₇ b ₆	$a_3 b_7 a_6 b_8$
Transfer (n=2)	a ₄ a ₅ b ₈ b ₅	a₄ b₄ a₅ b₁
	a ₅ a ₁ b ₆ b ₄	a ₅ b ₆ a ₁ b ₄
	$a_6 a_3 b_2 b_7$	a ₆ b ₈ a ₇ b ₇
	$a_7 a_2 b_3 b_2$	$a_7 b_3 a_8 b_2$
	a ₈ a ₇ b ₄ b ₁	$a_8 b_2 a_2 b_5$
	$a_1 a_2 b_5 b_6$	$a_1 b_3 a_3 b_8$
	$a_2 a_7 b_7 b_3$	$a_2 b_4 a_6 b_3$
	$a_3 a_4 b_3 b_4$	$a_3 b_8 a_4 b_4$
Probe (n=2)	$a_4 a_8 b_6 b_5$	$a_4 b_6 a_2 b_7$
	$a_5 a_1 b_2 b_1$	$a_5 b_1 a_8 b_6$
	$a_6 a_5 b_1 b_8$	$a_6 b_5 a_1 b_2$
	$a_7 a_6 b_8 b_7$	$a_7 b_7 a_5 b_5$
	$a_8 a_3 b_4 b_2$	$a_8 b_2 a_7 b_1$
	$a_1 a_7 a_6 b_6 b_2 b_8$	$a_1 b_6 a_2 b_7 a_4 b_5$
	$a_2 a_4 a_3 b_7 b_3 b_5$	$a_2 b_2 a_7 b_3 a_5 b_1$
	$a_3 a_8 a_1 b_2 b_6 b_7$	$a_3 b_8 a_8 b_6 a_1 b_4$
Probe (n=3)	$a_4 a_2 a_8 b_4 b_5 b_3$	$a_4 b_7 a_5 b_1 a_2 b_2$
	$a_5 a_3 a_4 b_5 b_8 b_6$	$a_5 b_5 a_6 b_4 a_3 b_3$
	$a_6 a_5 a_7 b_8 b_4 b_1$	$a_6 b_1 a_1 b_8 a_8 b_6$
	$a_7 a_1 a_5 b_3 b_1 b_2$	a ₇ b ₄ a ₃ b ₅ a ₆ b ₇
	$a_8 a_6 a_2 b_1 b_7 b_4$	$a_8 b_3 a_4 b_2 a_7 b_8$
	$a_1 a_7 a_5 a_2 b_4 b_3 b_8 b_6$	$a_1 b_8 a_7 b_2 a_5 b_3 a_4 b_6$
	$a_2 a_1 a_3 a_8 b_7 b_6 b_4 b_5$	$a_2 b_5 a_1 b_4 a_8 b_7 a_6 b_3$
	$a_3 a_8 a_1 a_5 b_6 b_4 b_2 b_7$	$a_3 b_1 a_4 b_5 a_2 b_6 a_8 b_7$
Probe (n=4)	$a_4 a_2 a_8 a_7 b_1 b_5 b_6 b_3$	$a_4 b_7 a_8 b_6 a_3 b_1 a_5 b_2$
	$a_5 a_6 a_4 a_1 b_8 b_7 b_3 b_2$	$a_5 b_4 a_6 b_8 a_7 b_2 a_3 b_1$
	$a_6 a_4 a_2 a_3 b_5 b_8 b_7 b_1$	$a_6 b_2 a_3 b_7 a_4 b_5 a_1 b_8$
	$a_7 a_3 a_6 a_4 b_2 b_1 b_5 b_8$	$a_7 b_3 a_5 b_1 a_6 b_8 a_2 b_4$
	$a_8 a_5 a_7 a_6 b_3 b_2 b_1 b_4$	$a_8 b_6 a_2 b_3 a_1 b_4 a_7 b_5$

Table S2. Motif patterns for all of the grammatical song stimuli used in this study. Letters and subscripts denote different motifs, as shown in Figure S2.





Gentner et al. Figure S2

Subjects. We used 11 adult European starlings. Subjects were captured on a nearby farm in the fall of 2001, and thus were all adults at the time of testing in spring-summer of 2004. Prior to testing, subjects were housed in large mixed-sex flight cages along with 15 – 20 conspecific birds. The cages were kept in a mixed species aviary containing separately caged zebra finches. The light schedule in the aviary followed local variation in solar day-length. While in the aviary all subjects had free access to food and water. Subjects were naive to all the training and testing stimuli at the start of behavioral training.

Behavioral apparatus. Starlings learned to recognize the training stimuli using an operant apparatus (Fig. S1a), mounted inside a 61 x 96 x 53 cm ID sound attenuation chamber. A cage mounted inside the chamber held the subject, while providing access to a 30 x 30 cm operant panel mounted on one side. The panel contained three circular response 'buttons' spaced 6 cm centre-to-centre, aligned in a row with the centre of each button ~14 cm off the floor of the cage, with the entire row centred on the width of the panel. Each response 'button' was a PVC housed opening in the panel fitted with an IR receiver and transmitter that detected when the bird broke the plane of the opening with its beak. This 'poke-hole' allowed starlings to probe with their beaks, a naturally occurring behavior. Each response opening was illuminated from the rear with an independently controlled LED. Directly below the centre button, in the section of cage floor immediately adjacent to the panel, a fourth PVC lined opening provided access to food. A remotely controlled hopper, positioned behind the panel, moved the food into and out of the subject's reach beneath the opening. Acoustic stimuli were transmitted via USB to each operant station (one station per animal being trained), converted to analogue form via a USB DAC, amplified, and then presented to the subject through a small audio speaker mounted ~30 cm behind the panel, out of the subject's view. We used custom software to monitor the subject's responses, and to control the LEDs, food hoppers, chamber-light and stimulus presentation according to procedural contingencies.

Shaping. Subjects learned to work the apparatus through a series of successive shaping procedures. After learning to feed reliably from the hopper, pecks directed toward a flashing LED behind the centre response port would lead to food reward. Once the subject pecked reliably at the centre port to elicit food reward, the LED ceased flashing, while the requirement to peck at the same location remained in effect. Shortly thereafter, pecks to the centre port initiated the presentation of a song stimulus, and the trial proceeded as described in the text. Although subjects could freely peck at the centre response port throughout stimulus presentation, only the first response following completion of the stimulus triggered reinforcement or punishment.

Probe procedure. Prior to initiation of the first probe session, the rate of food reinforcement for correct responses to S+ stimuli was lowered from 100% (where it had been during baseline training) to 80%, and the rate of "punishment" (dimmed house lights) for incorrect responses to S– stimuli was lowered to 95%. We reinforced all responses to all probe stimuli non-differentially regardless of accuracy as follows: each response to a probe stimulus had a 40% chance of eliciting a food reward, a 40% chance of eliciting punishment (timeout without food), and a 20% chance of eliciting no consequence at all. Because reinforcement of the probe stimuli is random and non-differential with respect to response outcome, subjects have no opportunity to learn to associate a given probe stimulus with a given response. Thus, the correct classification of probe stimuli that are novel exemplars from previously acquired classes is commonly taken as strong evidence that the subject is classifying stimuli based on some set of features common to the class rather than learning rote sets of specific exemplars. In fact, what keeps the subject responding to (and discriminating)

probe stimuli under these non-differential reinforcement contingencies is the generalization with the baseline training stimuli. If there was no generalization, classification accuracy would be at chance and all responses would be the same default.

Analysis. We used d-prime to estimate the sensitivity for discrimination between opposing stimulus classes (e.g. A^2B^2 and $(AB)^2$), such that:

(1)
$$d' = z(H) - z(F)$$

where z(H) is the z-score of the proportion of go responses to the A^2B^2 stimuli and z(F) is the z-score of the proportion of go responses to the $(AB)^2$ stimuli over some common number of trials (typically 100). We used the variance of d'^3 :

(2)
$$\operatorname{var}(d') = \frac{H(1-H)}{N_H[\phi(H)]^2} + \frac{F(1-F)}{N_F[\phi(F)]^2}$$

where N_H and N_F are numbers of trials of each stimulus class and for each class ϕ is given by:

(3)
$$\phi(p) = (2\pi)^{-\frac{1}{2}} \exp\left[-0.5z(p)^2\right]$$

to compute the confidence interval around specific values of d'. Any value of d' for which the lower bound of its 95% confidence interval was greater than 0.0 was considered to indicate significant discrimination of the two stimulus classes.

For each bird, we used Pearson's Chi-square to examine differences in the proportions of responses made to different agrammatical stimuli.

References

1. Gentner, T. Q. & Hulse, S. H. Perceptual mechanisms for individual vocal recognition in European starlings, Sturnus vulgaris. *Animal Behaviour* **56**, 579-594 (1998).

2. Fitch, W. T. & Hauser, M. D. Computational constraints on syntactic processing in a nonhuman primate. *Science* **303**, 377-80 (2004).

3. Macmillan, N. A. & Creelman, C. D. *Detection Theory: A User's Guide* (Cambridge University Press, Cambridge, 1991).