

Lecture 4

Extending Bigram Grammars

Bigram grammars are a step in the right direction from the initial list phonology model. Every grammar is still just a list of items, but now we use this list to determine the well-formedness of a word in a compositional manner. The problem of bigram grammars is that they are not powerful enough to capture certain phonological processes like intervocalic voicing — in linguistic parlance, the size of the grammars locality domain is limited to two adjacent symbols, so bigger contexts cannot be represented correctly. The obvious solution is to extend the size of the scanner window so that the grammar uses trigrams, 4-grams, or maybe even something bigger. We will see today that this generalization is a natural one in the sense that it preserves the essential properties of bigram grammars while increasing their empirical coverage.

1 Generalization to Strictly Local Grammars

In order to generalize bigram grammars to grammars with grams of arbitrary length, we just have to change a few parameters in our original definition. We also use this opportunity to slightly change our terminology: instead of *k-gram grammar* and *k-gram languages*, we will speak of *strictly k-local grammars*. This will make it easier later on to distinguish this formalism from a related one that also operates with *k-grams* but interprets them differently.

Definition 4.1 (*k-grams*). A *k-gram*, or *k-factor*, over alphabet Σ is an element of $(\Sigma \cup \{\bowtie, \bowtie\})^k$. Given a string w over Σ , its *k-augmented* counterpart $\hat{w}_k := \bowtie^{k-1} \cdot w \cdot \bowtie^{k-1}$ consists of w with $k - 1$ left edge markers and $k - 1$ right edge markers, and its set of *k-grams* is given by $k\text{-grams}(w) := \{s \in (\Sigma \cup \{\bowtie, \bowtie\})^k \mid \exists u, v \in \Sigma^* \text{ s.t. } u \cdot s \cdot v = \hat{w}\}$.

You can see that this definition of *k-gram* is a natural extension of the concept of bigrams for if we replace k by 2 in the definition, we get exactly the original definition of bigrams. The concept of bigram languages is generalized in the same fashion to *k-gram languages*, which jointly form the class of *strictly local languages*.

Definition 4.2 (Strictly Local Languages). A finite set of *k-grams* is called a *strictly k-local grammar*. A positive strictly *k-local grammar* G generates the language $L(G) := \{w \mid k\text{-grams}(w) \subseteq G\}$. A negative strictly *k-local grammar* G generates the language $L(G) := \{w \mid k\text{-grams}(w) \cap G = \emptyset\}$. A language L is *strictly k-local* iff it is generated by some strictly *k-local grammar*. The class of *strictly local languages* is given by $\bigcup_{k \geq 1} \{L \mid L \text{ is strictly } k\text{-local}\}$.

Notice that the definition of strictly k -local also allows for $k = 0$. Can you list all strictly 0-local languages? *Hint:* There's only two of them.

Intuitively, a language is strictly local iff all its well-formedness conditions are restricted to a locality domain of finitely bounded size.

The phonological processes we have looked at so far — word-final devoicing, nasal assimilation, and intervocalic devoicing — are all strictly local. The first two are strictly 2-local, the third one strictly 3-local.

It is also easy to see that constraints on syllable structure are strictly local. Take some language that only allows for syllables of the form V, VC, CV, and CVC, but not CCV, VCC, CVCC, or CCVCC. In other words, a well-formed word cannot contain more than two consecutive consonants, and these consonants cannot occur at the beginning or the end of the word, where they would necessarily be part of the same syllable. The set of illicit substrings, then, consists of $\times CC$, CCC, and $CC \times$. This can be compiled out into a set of illicit sequences of phones by substituting for each C the phones that are specified for [+cons]. So this example of a restricted syllable template does not exceed the power of a negative strictly 3-local grammar.

If desired, the negative strictly local grammar can be converted into a positive one using our standard procedure — go back to the proof of our theorem that positive and negative bigram grammars are equivalent, you will see that it does not at all rely on the length of the n -grams being 2, so the theorem can easily be generalized to all strictly local grammars.

Feel free to design strictly local grammars for other local processes in phonology, e.g. assimilation across word boundaries, vowel harmony, umlaut, or spirantization. As long as these processes describe surface true generalizations — that is to say, they do not make reference to an underlying form and can be stated purely in terms of which output forms are licit — you should not run into any problems. Notice that these processes must be local but can nonetheless be globally unbounded. Vowel harmony, for instance, may apply throughout an entire word via a sequence of local vowel harmony steps. This kind of behavior is exactly what we expect given how strictly local grammars operate: the n -grams only regulate the shape of local domains, but the well-formedness of the word is evaluated by moving through the word and checking each local domain. So not only are strictly local grammars powerful enough to capture local well-formedness conditions in phonology, the way they enforce them mirrors linguists' intuitions about how local processes can produce global patterns.

The connection between strictly local grammars and local processes in phonology allows us to study the latter through the former. Since strictly local grammars generate strictly local languages, this implies that the properties of strictly local languages can tell us something about local processes in phonology. So now the theoretical task of proving formal properties of strictly local languages has suddenly morphed into an empirically minded investigation of phonology.

2 Exploring Strictly Local Languages

2.1 A Proper Hierarchy of Strictly Local Languages

One may expect that the power of strictly local grammars increases with the size of the locality domain. This can indeed be shown to be the case by recourse to the strictly local languages. If we order the strictly local languages by the size of their locality domain, we get a proper hierarchy: one level properly subsumes the next. Denoting the class of strictly k -local languages by SL_k , we have $SL_k \subsetneq SL_{k+1}$ for all $k \geq 0$.

Lemma 4.3. It holds for every $k \geq 0$ that if language L is strictly k -local, then L is also $k + 1$ -local. \lrcorner

The proof for this lemma is slightly more complicated than anything we have seen so far, but behind the notation lies a very simple idea: the size of the locality domain can be increased from k to $k + 1$ by padding the edge markers and by combining two overlapping k -grams into a single $k + 1$ -gram.

Proof. If L is strictly k -local then it is generated by some positive strictly k -local grammar G over some alphabet Σ . Let G' be the smallest set such that for all k -grams $g_1, g_2 \in G$

- if g_1 starts with \bowtie , then $\bowtie g_1 \in G'$,
- if g_1 ends with \bowtie , then $g_1 \bowtie \in G'$,
- if $g_1 := a_1 a_2 \cdots a_k$ and $g_2 := a_2 \cdots a_k a_{k+1}$, then $a_1 a_2 \cdots a_k a_{k+1} \in G'$.

Clearly G' is finite and can therefore be interpreted as a positive $k + 1$ -local grammar. We show that $L = L(G')$, thus establishing that L is $k + 1$ -local.

If $w \in L(G')$, then $k + 1$ -grams(w) is a subset of G' . All $k + 1$ -grams with multiple edge markers have a corresponding k -gram with one edge marker less, which is contained in G . All other $k + 1$ -grams are split into two k -grams by removing the first or the last symbol. Each k -gram is once again contained in G , so k -grams(w) $\subseteq G$ and hence $w \in L(G) = L$. Since w was arbitrary we have $L(G') \subseteq L$.

If $w \in L$, then k -grams(w) is a subset of G . Assume towards a contradiction that $k + 1$ -grams(w) is not a subset of G' . Then w contains some $k + 1$ -gram g_3 that is not a member of G . If g_3 starts or ends with two edge markers, then the corresponding k -gram with only one of the two markers cannot have been part of G , contradicting our initial assumption. In all other cases, g_3 is built from two overlapping k -grams g_1 and g_2 , at least one of which is not contained in G . But then k -grams(w) is not a subset of G , contradicting once more our initial assumption. It follows, then, that $k + 1$ -grams(w) is a subset of G' after all, wherefore $L \subseteq L(G')$. \square

Lemma 4.4. For every k there is some strictly $k + 1$ -local language that is not strictly k -local. \lrcorner

Proof. Consider the finite language L that contains only the string a^k . It is generated by the strictly $k + 1$ -local grammar $\{\bowtie a^k, a^k \bowtie\}$. However, k -grams(a^k) = $\{\bowtie a^{k-1}, a^k, a^{k-1} \bowtie\} = k$ -grams(a^n) for every $n \geq k$, so a strictly k -local grammar that generates a^k also generates all these a^n and thus a proper superset of L . \square

This proof uses a finite language, but the lemma holds even if we only consider infinite languages. Try to generalize the proof along these lines.

Theorem 4.5. For all $k \geq 0$, $SL_k \subsetneq SL_{k+1}$. \lrcorner

Now we know for sure that the size of the locality domain has a direct effect on generative capacity. This may seem barely surprising, but it is far from trivial. In the coming weeks we will see that for a very similar formalism, all $k \geq 2$ have exactly the same power.

2.2 Relation to Finite Languages

Remember that the list phonology model of Lecture 2 was restricted to lists of finite length, so it could only generate finite languages. This restriction is empirically

inadequate as it conflicts with the assumption that the list of phonological words in a given natural language is infinite and fails to handle nonce words and linguistic creativity in general. But not only is the list phonology model restricted to finite languages, it adds insult to injury with its ability to generate all finite languages. Every finite language is a viable natural language phonology, according to the list model, and we have already seen why this is typologically untenable. So the class of languages that can be generated by the list phonology model is exactly the class of finite languages, which is the class that is least likely to provide an insightful or empirically adequate model of language.

The strictly local languages improve on this significantly, but only if one adopts the right perspective. First, it is obvious that strictly local languages can be infinite, so not every strictly local language is finite. Let us make this claim fully explicit via a proof. We already know that every strictly k -local language is strictly $k + 1$ -local, so all we need is an example of a strictly 1-local language.

Further reflection reveals that almost every 1-local language is infinite. In fact, there are only two 1-local languages that are finite. Can you define them? *Hint:* One of them is strictly 0-local.

Lemma 4.6. There is a strictly 1-local language that is infinite. \lrcorner

Proof. Let ${}^+G := \{\bowtie, a, \bowtie\}$. Then $L({}^+G) = \{\varepsilon, a, aa, aaa, \dots\} = a^*$, which is infinite. \square

Theorem 4.7. For every $k \geq 1$, SL_k contains an infinite language. \lrcorner

This is a welcome result because it shows that no matter what size of locality domain we pick, we are never restricted to just finite languages. An even more appealing property of strictly local languages is that for every level of the infinite hierarchy there are some finite languages that cannot be defined. Consequently, strictly local grammars improve on the list phonology model in that they cannot define any arbitrary phonological systems.

Theorem 4.8. For every $k \geq 1$, there is some finite language that is not contained in SL_k . \lrcorner

Proof. Pick an infinite language that is generated by some strictly k -local grammar G over alphabet Σ . Discard from G all k -grams that start with \bowtie or end in \bowtie . The resulting set is a finite language L . We show that L is not strictly k -local.

Since G generates an infinite set, it must contain (not necessarily distinct) k -grams $a \cdot u$ and $u \cdot b$, where $a, b \in \Sigma$ and $u \in \Sigma^{k-1}$. Note that both $a \cdot u$ and $u \cdot b$ belong to L , whereas $a \cdot u \cdot b$ does not. But $k\text{-grams}(a \cdot u) \cup k\text{-grams}(u \cdot b) = k\text{-grams}(a \cdot u \cdot b)$, and consequently every strictly k -local grammar that generates $a \cdot u$ and $u \cdot b$ also generates $a \cdot u \cdot b$. Hence L is not strictly k -local. \square

The proof above is rather sneaky. It exploits the fact that every strictly k -local grammar is a finite set of strings of length k , which can be viewed as a finite language. This is an interesting perspective: a strictly k -local grammar G is a finite language L_G coupled with a specific algorithm A for creating a new language from L_G . If we try to generate L_G via a strictly k -local grammar, this grammar automatically uses the algorithm A , so for certain choices of L_G the grammar generates additional strings via A that do not belong to L_G .

Careful though: The previous theorem holds only when the size of the locality domain is fixed by some k . The whole class of strictly local languages properly subsumes the class of finite languages.

Theorem 4.9. Every finite language is strictly local. \lrcorner

Proof. Suppose that L is a finite language, the longest string of which has length $k - 1$, $k \geq 1$. We define a k -local grammar G that consists of the k -grams $\bowtie^i \cdot w \cdot \bowtie^j$, where $w \in L$, w has length $l < k$, and $i + j + l = k$. Since every k -gram starts with \bowtie or ends in \bowtie , $L(G) = L$, wherefore L is strictly k -local. \square

The relation between strictly local languages and finite languages thus is more involved than one would expect. Without restrictions on the size of the locality domain, the strictly local languages include all finite languages. However, the class of strictly k -local languages and the class of finite languages are incomparable — they have a non-empty intersection, but neither subsumes the other. So assuming that k is fixed for natural language phonology, e.g. as part of Universal Grammar, strictly local grammars are a good approximation of local processes in phonology. They can handle infinity and do not incorrectly predict languages to vary freely across all dimensions.

2.3 Substring Substitution Closure

In several of the preceding proofs we have used the fact that a strictly local grammar sometimes “overshoots the target”. A finite language L may not be in SL_k because a strictly k -local grammar that tries to generate L will also end up generating other strings outside of L . But a strictly $k + 1$ -local grammar may be able to do a point landing and generate all and only those strings that are members of L . What this shows is that a strictly k -local grammar only has perfect precision within its locality domain of size k , beyond that it has to generalize. This generalization step is what allows strictly local grammars to generate infinite languages and thus the true source of their power.

Crucially, though, strictly local grammars don’t just generalize randomly, quite to the contrary: all strictly local grammars generalize in the same fashion, irrespective of the size of their locality domain. Their generalization strategy is already implicit in the definition of strictly local languages, which categorizes strings as well-formed or ill-formed according to their set of k -grams. If two strings have exactly the same set of k -grams, then either both are well-formed or both are ill-formed.

But it is quite hard to tell what this condition implies for the overall shape of the string languages. Given a string language L , how can we tell whether L is strictly local? Fortunately strictly local languages are uniquely characterized by a property called *substring substitution closure*.

Definition 4.10 (Local Substring Substitution Closure). A language L satisfies *k-local substring substitution closure* iff there is some $k \geq 1$ such that if L contains both $u \cdot x \cdot v$ and $u' \cdot x \cdot v'$, where x has length $k - 1$, then L also contains $u \cdot x \cdot v'$.

Theorem 4.11. A language is in SL_k iff it satisfies k -substring substitution closure. \lrcorner

Let us look at a couple of examples first before wading through the proof of the theorem.

Example 4.1 A Substring Substitution Closed Language

We have already seen that the language $(ab)^+$ is strictly 2-local as it is generated by the grammar $\{\bowtie a, ab, ba, b\bowtie\}$. Now we can also verify this via substring substitution closure. For instance, we can line up $abab$ and $abababab$ to show that the language must also contain $ababab$.

$$\begin{array}{rcccl} & x & & & \\ & ab & a & b & \in L \\ abab & a & bab & & \in L \\ \hline ab & a & bab & & \in L \end{array}$$

Notice how x is a single symbol since its length must be $k - 1 = 2 - 1 = 1$. Also, we could have established the membership of $ababab$ more succinctly using just $abab$ or $abababab$.

$$\begin{array}{rcccl} & x & & & \\ ab & a & b & & \in L \\ & a & bab & & \in L \\ \hline ab & a & bab & & \in L \end{array} \quad \begin{array}{rcccl} & x & & & \\ & ab & a & babab & \in L \\ abab & a & bab & & \in L \\ \hline ab & a & bab & & \in L \end{array}$$

Since substring substitution closure fails if even a single string is missing from the set, it is usually not a good way to of showing that a language is strictly local — if the language in question is infinite, one cannot show via specific substitutions that it is suffix substitution closed. However, suffix substitution closure is an excellent way of showing that a language is not strictly local by giving a single example for how one can find a missing string.

Example 4.2 A Language that Fails Subtree Substitution Closure

Consider the language $(aa)^+$, the variant of $(ab)^+$ where all b s have been replaced by a s. This language is not strictly local as it fails k -local substring substitution closure for any choice of k . Suppose k is an even number:

$$\begin{array}{rcccl} & x & & & \\ a & a \cdots a & a & & \in L \\ & a \cdots a & & & \in L \\ \hline a & a \cdots a & & & \notin L \end{array}$$

A similar pattern can be used if k is odd. So no matter what the value of k , the language is not suffix substitution closed and thus not strictly local. You might find this surprising given that the only difference to $(ab)^+$ is the replacement of b by a . This highlights another property of strictly local languages: the alphabet plays a crucial role in what languages are definable.

Proof. To be done.

□

2.4 Closure Properties

Suffix substitution closure is — as its name implies — a *closure property*. One says that an object o is closed under an operation iff applying this operation to elements of o yields only elements of o . In other words, the operation never takes us outside of o . The natural numbers, for instance, are closed under addition since the sum of two natural numbers is yet again a natural numbers. But they are not closed under subtraction because, say, $2 - 5$ yields -3 , which is an integer but not a natural number. Suffix substitution closure simply means that a language is closed under the operation of substituting suffixes in a specific way.

But of course there are many other operations that can be applied to a language, and it will be interesting to see whether strictly local languages are closed under them. In particular the basic set-theoretic operations of intersection, union, and relative complement are of interest since they tell us how we can build strictly local languages from smaller ones.

Let us look at closure under intersection first. This one is particularly important because of the close correspondence that negative strictly local grammars establish between constraints on the one hand and languages on the other. If every n -gram corresponds to a well-formedness constraints, then one would expect that one can get the intersection of two languages by simply conjoining their respective well-formedness constraints. That is indeed the case.

Lemma 4.12. The class of strictly k -local languages is closed under intersection, $k \geq 0$. ┘

Proof. We prove that for any two strictly local languages generated by (positive) grammars G_1 and G_2 , $L(G_1) \cap L(G_2) = L(G_1 \cap G_2)$. We first show $L(G_1) \cap L(G_2) \subseteq L(G_1 \cap G_2)$. Let w be an arbitrary string belonging to both $L(G_1)$ and $L(G_2)$, i.e. $w \in L(G_1) \cap L(G_2)$. Then every k -gram of w is contained in both G_1 and G_2 , so $w \in L(G_1 \cap G_2)$. Since w is arbitrary, we have $L(G_1) \cap L(G_2) \subseteq L(G_1 \cap G_2)$. The same reasoning can be applied in the other direction, yielding $L(G_1 \cap G_2) \subseteq L(G_1) \cap L(G_2)$. These two facts jointly imply $L(G_1) \cap L(G_2) = L(G_1 \cap G_2)$. □

Give an analogous proof using negative grammars.

Closure under intersection can also be used to prove that a language is not strictly local. Suppose that we know that L is strictly local, but we have a hard time showing that L' is not strictly local. Then we can instead try to show that $L \cap L'$ is not strictly local, as this immediately implies the non-locality of L' , too.

One might expect that closure under union holds too since one can simply take the union of the grammars, but this does not work as expected. The union of two grammars G_1 and G_2 often generates a subset of the union of $L(G_1)$ and $L(G_2)$. Do not even try to look for a smarter strategy to build a grammar for the union of two strictly local languages, there is none that works in all cases.

Lemma 4.13. The class of strictly local languages is not closed under union. ┘

Proof. We give an example of two bigram languages whose union is not a bigram language. The proof can easily be adapted for arbitrary values of k .

Let $L_1 := \{ab\}$ and $L_2 := \{b, bb, bbb, \dots\}$. Assume without loss of generality (w.l.o.g.) that all bigram grammars are positive. Then any bigram grammar that generates L_1 must contain the bigrams $\times a$, ab , and $b \times$. Similarly, a bigram grammar generating L_2 must contain the bigrams $\times b$, bb , and $b \times$. Therefore a bigram grammar

that generates all strings in $L_1 \cup L_2$ must contain at least these bigrams. But such a grammar also generates the string abb , which is not part of $L_1 \cup L_2$. Hence there is no bigram grammar that generates all the strings in $L_1 \cup L_2$ and nothing else, so $L_1 \cup L_2$ is not a bigram language. \square

The previous two lemmata immediately imply that closure under relative complement does not hold either.

Lemma 4.14. The class of strictly local languages is not closed under (relative) complement. \lrcorner

Proof. By De Morgan's law $L_1 \cup L_2 = \overline{\overline{L_1} \cap \overline{L_2}}$. If the class of strictly local languages were closed under both complement and intersection, it would thus be closed under union, too. Since closure under intersection holds while closure under union does not, closure under relative complement cannot hold, either. \square

You may find this result surprising because the complement of a strictly local grammar is a strictly local grammar. But in general $L(G)$ differs from $L(\overline{G})$, just like $L(G_1) \cup L(G_2)$ is not guaranteed to be $L(G_1 \cup G_2)$. These differences illustrate why it is so important to distinguish between grammars and the languages they generate.

We finish with yet another missing closure property.

Lemma 4.15. The class of strictly local languages is not closed under relabelings. \lrcorner

Proof. We have already seen that $(ab)^+$ is strictly 2-local whereas the relabeling $(aa)^+$ is not strictly local at all. But the former is the image of the latter under a relabeling that replaces all bs by as . \square

This is a very important result as it ties directly into the abstractness debate in generative phonology. Suppose that we are allowed to have all kinds of hidden structure in our phonological representation, e.g. a basic syllable template or feet. Then we could assume that the language $(aa)^+$ is underlyingly equipped with a CV-template, so that we should rather think of it as $(a_C a_V)^+$. This language is strictly local (it is a notational variant of $(ab)^+$), so $(aa)^+$ is strictly local under a mapping that removes unpronounced structure. In later lectures we will see that this is a very dangerous route to take: hidden structure pushes strictly local languages to a level of power where all phonological processes are equally simple. This eradicates all distinctions between processes and takes us back to the undesirable egalitarianism of the list phonology model.

3 Implications for Phonology

Our mathematical expedition has taught us a surprising amount about phonology. First of all, all local phonological dependencies (including those spanning across word boundaries) can be handled by strictly local grammars, a very simple formalism with few cognitive requirements. The model comes with a minimal memory burden as it only requires the speaker to keep track of a small number of k -factors. It is also very fast: recognition via a scanner takes linear time and can be carried out in an incremental online fashion. Lookup of k -factors in the grammar takes at most logarithmic time using binary search, and even fast search methods can be implemented.

But strictly local grammars also capture essential properties of phonological competence. They are analytic in nature and thus capable of generating infinite languages, which solves the problems the list phonology model had with linguistic creativity and nonce words. The special relation between strictly k -local languages and finite languages also means that certain typological pitfalls are avoided. Not every random collection of strings is predicted to be a valid phonological system, and grammars generating an infinite language generalize in a linguistically plausible fashion by determining the well-formedness of strings as a composite function of the local domains.

We have also seen that strictly local languages lack closure properties that do not hold of natural languages either. The union or relative complement of a natural language's phonotactics is not guaranteed to be a valid phonological system for a natural language, just like the class of strictly local languages is not closed under these operations. The increase of power brought about by relabelings also cautions us against using hidden structures and highly abstracted alphabets, an issue that has been discussed at length in phonology and that will occupy us at various points for the rest of the course.