

1 Introduction: constraints in phonological acquisition

René Kager, Joe Pater, and Wim Zonneveld

This volume presents ten studies in phonological first language acquisition, an area of research that has become one of fast-growing importance in recent years. The reason for this is not just the fruitfulness and linguistic interest of this type of study *per se*: it is also the case that the more we come to know about phonological development, by the analysis of growing numbers of data collections and increasingly sophisticated experiments, the more the field has complied with the notion that acquisition research lies at the heart of the modern study of language. One of the aims of this introduction is to illustrate and discuss these developments. In line with them, the past decade in phonology in particular has witnessed an upswell of productive interaction between empirical acquisition research and theory development. With the arrival and rise of constraint-based models, in particular Prince and Smolensky's (1993) Optimality Theory, phonological theory now provides a framework that meets the desiderata expressed more than two decades ago by Lise Menn (1980: 35–36), who is also a contributor to this volume:

- (1) . . . The child's 'tonguetiedness', that overwhelming reality which Stampe and Jakobson both tried to capture with their respective formal structures, could be handled more felicitously if one represented the heavy articulatory limitations of the child by the formal device of output constraints [. . .] The child's gradual mastery of articulation then is formalized as a relaxation of those constraints.

The rapid emergence of acquisition studies within Optimality Theory reflects the general suitability of constraints for the formalisation of developmental limitations, as well as the usefulness of constraint ranking for expressing the relaxation of these limitations. In this volume, we include several chapters that provide concrete examples of the formalisation of phonological development in terms of constraint ranking. Other chapters address fundamental issues such as learnability, the nature of the initial state, the relation between perception and production, and the contribution the experimental approach can make. They all demonstrate the successes of a constraint-based approach to these issues, but do not ignore the challenges that it faces.

In the first section of this introduction, we show how many issues and ideas that appear in this volume have important precedents in prior research. We provide a brief survey of these issues, as they have been discussed in the research tradition that links phonological theory and phonological acquisition, of which the Optimality theoretic approach is the most recent outgrowth.¹ In the second section, we provide a tutorial on the fundamentals of Optimality Theory, specifically tuned to its application to acquisition. And in the final section we give an overview of some central issues in acquisition and learnability in Optimality Theory, drawing connections between these issues and the contents of the chapters of this volume. This last section includes summaries of all of the chapters except for Lise Menn's contribution, which itself is a summary and discussion of research on phonological acquisition and its relation to Optimality Theory, from a perspective which seems to complement the one taken in this introductory chapter.

1. The young child in generative grammar

It is not *a priori* obvious that the study of phonological theory and that of child speech should be connected in any way whatsoever. After nearly a century of prior diary studies of child utterances, and typological study of the sound systems of the world's languages, it was Roman Jakobson, in his *Child Language, Aphasia and Phonological Universals* (1941),² who first suggested that they were governed by the same principles: Jakobson wanted to

- (2) establish the point that adult language systems are as they are because they necessarily develop in a particular systematic way that can be studied in the form of child language. (Anderson 1985: 130)

These principles Jakobson referred to as 'laws of irreversible solidarity', but they have become known as 'implicational universals'. A typical example is this one (Jakobson 1941/1968: 51): the acquisition of fricatives presupposes the acquisition of stops in child language; and in the linguistic systems of the world the former cannot exist unless the latter exist as well. Alongside such laws,³ Jakobson proposed principles of maximal contrast that governed the gradual emergence of phonological structures in child language, and also determined the content of many cross-linguistic implicational universals. After Jakobson, systematic research on child phonology was largely devoted to testing his claims. Its upshot appears to be that many of his principles are roughly supported (see Menn 1980, and Ingram 1989 for overviews of this stage of research), but that there is considerable unpredicted variability between children. Jakobson was not working with a very explicit ('formal') or fleshed-out phonological theory. One problem that arises is that there is no means of capturing the mapping from adult to child speech. For example, when a child lacks consonant clusters, the

choice it makes between the two consonants appearing in the adult target cluster is not random, though it does vary from child to child. These mappings were brought within the purview of phonological analysis with the emergence of the much more formally oriented generative phonology.

The 'young child' appeared for the first time in the generative literature in Chomsky (1959). This is the passage where (s)he appears:

- (3) The child who learns a language has in some sense constructed [a] grammar for himself on the basis of his observation of sentences and nonsentences (i.e., corrections by the verbal community). Study of the actual observed ability of a speaker to distinguish sentences from nonsentences, detect ambiguities, etc., apparently forces us to the conclusion that this grammar is of an extremely complex and abstract character, and that the young child has succeeded in carrying out what from the formal point of view, at least, seems to be a remarkable type of theory construction. Furthermore, this task is accomplished in an astonishingly short time, to a large extent independently of intelligence, and in a comparable way by all children. [. . .]

The fact that all normal children acquire essentially comparable grammars of great complexity with remarkable rapidity suggests that human beings are somehow specially designed to do this, with data-handling or 'hypothesis-formulating' ability of unknown character and complexity.

These remarks imply a meaningful initial state or *faculté de langage* or Universal Grammar (UG), of Chomsky (1965: 5–6, 1968: 27), the study of which formulates the 'conditions that a system must meet to qualify as a potential human language, conditions that [. . .] constitute the innate organization that determines what counts as linguistic experience and what knowledge of language arises on the basis of that experience'. Universal Grammar is the innate starting point of language acquisition for each human being, and it is a Jakobsonian concept in the sense that its elements are hypothesised also to appear as 'universals' in typological language studies.

The child tries to find its way through the data maze assisted by UG, and constructs an abstract system called a grammar. This grammar is an amalgamation of persisting elements from UG (in the way envisaged by Jakobson) and acquired, language-specific components. It is the linguist's task to make sense of the structure of these grammars, and the role UG plays in their growth. The contributions discussed in this section in some way or other all take up the challenge posed by Jakobson and Chomsky. Although their topics are often intertwined, we chronologically subdivide them as follows: first Smith (1973) and Stampe (1969, 1973a, 1973b); then Braine (1976), Macken (1980), Kiparsky and Menn (1977), and Menn (1978, 1980), and finally the literature

dealing with parameters starting with Chomsky (1981a, 1981b). These works raise issues of fundamental importance, arguing directly from observed properties of the acquisition process. Roughly, these issues can be subdivided into 'formal' ones, and ones of 'substance'. Among the former are formal aspects of the rule-based approach towards the phonological component; and the abstractness of the underlying forms of child speech. Among the latter are: markedness, typology and language acquisition; rules vs. processes; conspiracies and output constraints; parameters; and the perception-production dichotomy.

It is the contention of Chomsky (1965: 16, 28) that the studies of syntax, semantics, and phonology can all proceed along the lines of reasoning just described. Chomsky and Halle (1968) is the first sizeable illustration of this methodology for Generative Phonology. No secret is revealed by the assessment that this work had a very strong 'formal' bias, and that the link between form and substance was underdeveloped. At the heart of this approach were four formal devices: (i) feature representations using a UG-based set of phonological features drawn in from the pre-generative era (see Halle 1983); (ii) derivations mapping lexical representations onto surface ones by rewrite rules, formulated in a UG-based format and mutually ordered along lines also specified in UG ('linear ordering', with a number of further specifications); (iii) morpheme structure rules stating redundancies at the level of the lexicon; and (iv) an evaluation measure also located in UG, evaluating a grammar's formal complexity as a function of the number of symbols in it. This measure selected less complex grammars over more complex ones for any given body of data. Although it was intended both to contribute to an explanation of the behaviour of native speakers, including acquisitional behaviour, in actual practice it was only infrequently called upon given the complexities of the linguist's task to formulate an analysis of the empirical data under investigation at all, given the other three formal devices and their UG aspects.

A seminal acquisitional study in this framework is Smith (1973), which also presents a wealth of original data (assembled in the form of a detailed longitudinal study of the progress of the author's son Amahl between the ages of [2;3] and [3;11]). His monograph has two aims. First, to use child data to argue for the correctness of the view of the grammar as a system of ordered rules deriving an output form from an underlying form. Second, to endorse the idea that child grammars just as much as adult grammars are constrained by the universals of UG. In conformity with the priorities of the times, the principal universal discussed in Smith's study is that of rule ordering; just as rules are (linearly) ordered in an adult grammar, rules can be crucially ordered at certain acquisitional stages, they can be modified, be reordered, or disappear at later stages, until the adult grammar is reached, which contains the final stage of ordered rules. Consider Smith (1973: 158):

- (4) Chomsky [1967] suggested (p. 105) that “the rules of the grammar must be partially ordered,” going on to claim that the principle of rule ordering was an a priori part of the basis which made language acquisition possible (Chomsky 1967, pp. 127–128). To the extent that one can establish the psychological validity of the realisation rules [of Amahl’s grammar] and to the extent that the ordering relations established among these rules are necessary, so is Chomsky’s claim substantiated.

Thus, the way to describe the phonological behaviour of language-acquiring children is by means of rules, and UG says that linear ordering is an essential property of the rule set.

Taking into account current debates, two issues in Smith’s work are worth highlighting. The first is that of ‘opaque’ rule interactions (which become a separate issue, for instance, once translated into Optimality Theory; see McCarthy 1999a, 2001: 163 ff.); and the second that of the nature of the underlying representations in a child grammar.

Smith (1973: 158–161) contains a long list of ordering relations among the rules of Amahl’s grammar, for a variety of successive stages of acquisition. The most interesting ones occur at Stage 1 ([2;2]–[2;4]), when Smith started collecting data. Included here are examples of unusual or ‘marked’ ordering in the sense of Kiparsky (1968): bleeding and counterfeeding orderings, leading to ‘surface opacity’. One of these examples has grown into a celebrated one: that of the ‘chain shift’ exhibited in the *puzzle/puddle* case. Observe that a counterfeeding relation holds between the two rules of (5), coexisting from Amahl’s Stage 1 to approximately Stage 14 [2;8]:

- | | | | | |
|--|---|--------|---------|---------|
| (5) | (a) velarization of coronal stops before [l] (Rule 3) | pedal | → | bɛ[g]u |
| | | bottle | → | bo[k]le |
| | | kennel | → | ke[ŋ]el |
| | | puddle | → | pu[g]le |
| (b) neutralization of coronal fricatives and stops (Rule 24) | zoo | → | [d]oo | |
| | bath | → | b[a:t] | |
| | knife | → | mi[p]e | |
| | puzzle | → | pu[d]le | |

The crucial pair in the data is that of *puddle* and *puzzle*, whose pronunciations require a counterfeeding rule ordering (Kiparsky shows that such orderings are avoided in grammars, and subject to historical change; see also Kiparsky and Menn 1977);⁴ schematically:

- (6) by (5a) puddle → pu[g]le
 by (5b) puzzle → pu[d]le (→ *pu[g]le)

This chain shift also shows that the move away from a coronal in the former case cannot be due to some production inability; the striking characteristic of this case will be returned to below.

In a phonological system in which alternations based on rich morphology are (still) virtually lacking, as in early child systems, the nature of the underlying forms of the system is a potentially controversial issue. Rejecting the perhaps initially plausible view that the child's phonology acts as a self-contained system with a low degree of abstractness in which the underlying forms are very close to the output forms, Smith argues that the child's underlying forms are generally equivalent to the adult surface forms that the child takes in as his day-to-day input. This is not in his case just an *a priori* assumption or just a first approach to the problem: it is based on a range of evidence showing that the child actually operates in this manner (Smith 1973: 11). The conversions in (5a) and (5b) are not just the author's expository manipulations, but represent actual hypothesised characteristics of the child's grammar. Let us briefly review some of this evidence for such analyses (Smith 1973: 133–148).

First, Amahl must have stored the adult forms because he was able to recognise and discriminate items of the adult language which he himself could not or did not produce or discriminate in his production. Thus, he could point correctly to pictures of a *mouse* and a *mouth* 'before he was able to speak at all' and when he was 'still unable to produce the contrast between [s] and [θ]' (p. 134): both are [maut] and then [maus] at exactly the same stages. Second, in a case of so-called 'phonemic overlap', the possibility for some of his *l*-initial output words to alternate with [r-], and for other *l*-initial words not to do so, was entirely based on the adult distinction between words beginning with *r*- and *l*-, respectively (so [rait, lait] for *right*, and [lait] for *light*), implying an optional rule operating on adult-like underlying representations. Third, the development of certain items over time point to the same property of the grammar: 'before consonant+/l/ clusters appeared at all, there was neutralisation of such adult examples as *bed* and *bread* as [bed]. Once clusters appeared these were differentiated as [bed] and [bled], respectively, and likewise for many comparable items' (p. 139). Fourth, Amahl, as soon as he learned a new sound or sound combination, immediately utilised it correctly 'across the board' in all the relevant words, rather than incorporating it separately and slowly into each word. This indicates 'that these sounds and sound sequences must have been stored in the brain "correctly" in order for their appearance to be so consistently right. [Thus,] once [l] appeared for /s/ it appeared in all words containing /s/ nearly at the same time' (p. 139). Finally, evidence for adult-like underlying forms comes from early alternations, such as those involving plural formation (p. 148). Consider the data below, resembling an alternation such as that produced by final devoicing in a language such as Dutch and German.

- (7) sg. cat → [kæt] pl. cats → [kæt]
 horse → [ɔ:t] horses → [ɔ:tid]
 cloth → [klɒt] cloths → [klɒtid]

In the plural, in spite of the wholesale neutralization in the singular, a distinction surfaces between adult stops and fricatives, indicating that the underlying form of these morphemes is much more adult-like than inspection of just the singulars would suggest.⁵

There is a small handful of passages in his work where Smith describes further prospects that go beyond the essentially 'formal' early-Generative Phonology Chomsky and Halle type approach that his work otherwise falls into:

- (8) (a) The use of marking conventions [Chomsky and Halle 1968] in this study has been mainly conspicuous by its absence, though there are two directions in which they might be of relevance to developmental phonology [. . .]. (Smith 1973: 199)
 (b) [The realisation rules of the child phonology] clearly have much in common with the general constraints suggested in Stampe's paper (Stampe, 19[69]). (Smith 1973: 133)
 (c) Kisseberth's [1971] concept of 'functional unity' [...] is clearly relevant here, where there are obvious 'infantile conspiracies' – i.e., sets of rules – implementing the 'general tendencies' underlying the realisation rules discussed above . . . (Smith 1973: 204–205)

The implications of these references can be fleshed out by turning to the references mentioned in (8a–c), and to those listed earlier in this section.

In Chomsky and Halle (1968) a good grammar is one of low cost, in a formal sense. In the final chapter of their book (chapter 9), however, these authors submit that such an approach, at least as presented by them in their monograph, may have been 'overly formal', in particular in the area of the 'intrinsic content' of phonological features. They offer an outline of 'a theory of markedness', based on a proposal for extending the evaluation measure by means of 'marking conventions'. Marking conventions measure lexical representations and phonological rules in terms of their featural content, based on universal markedness. For each feature, an unmarked and a marked value are distinguished, sometimes depending on the segmental context. Two examples are given below, of a context-free and a context-sensitive convention:

- (9) (a) the unmarked value for [nasal] is [–nasal]
 (b) the unmarked value for [round] is [+round], in back vowels

Lexical representations are assumed to be specified in terms of Ms and Us, which were converted into +/– feature values by the conventions. Only M-specifications contribute to the cost of a grammar. With regard to phonological

rules, marking conventions can have a 'linking' function, making some rules more costly than others, hence dispreferred in grammars as well as presumably in the acquisition process. As an example, consider the rounding that often accompanies the backing of front non-round vowels: given the evaluation measure, rounding threatens to be more costly than leaving it out. If markedness convention (9b) links up to the output of a backing rule, it will supply the roundness feature, whereas blocking of the convention must be achieved by specifying the marked value in the structural change of the rule.

Smith (1973: 199–201) makes an attempt at establishing the usefulness of the marking conventions for his results. In spite of the promising beginning of this attempt quoted in (8a), his conclusion is exactly the opposite: 'it seems that with the exception of one or two isolated examples of the type cited, marking conventions are completely irrelevant', adding that '[i]n general it would seem that the present state of ignorance makes it impossible to effect any interesting correlation between acquisitional phenomena and marking conventions'. Anderson (1985: 334–342) makes a more fundamental point. In spite of the expressed aim of being less 'overly formal', Chomsky and Halle's markedness theory 'is in fact an attempt at exhaustively reducing the considerations of phonetic content that might be relevant to phonology to purely formal expression in the notation (now enhanced by its interpretation through the marking conventions). It is thus entirely consistent with the original *SPE* program of reducing all of the theory of phonological structure to a single explicit formal system including a notation and a calculus for manipulating and interpreting expressions within that notation. [...] The revision involved, however, was in a more complete working out of the goal of reducing phonology to a formal system rather than a replacement of that goal with some other' (pp. 333–334). He adds that 'essentially no substantial analyses of phonological phenomena have appeared subsequently in which this aspect of the theory plays a significant role. One MIT dissertation (Kean 1975) was devoted to further elaboration of the theory, but this remained (like chapter 9 of *SPE* [*The Sound Pattern of English*]) at the level of a programmatic statement rather than constituting an extended analysis of the phonology of some language(s) in the terms prescribed by the theory.'

Detailed criticism of Chomsky and Halle's markedness theory was provided by Stampe (1973a, 1973b). At first glance his distinction between (language-specific) *rules* and (universal) *processes* closely resembles the *SPE* one between rules and markedness conventions, but a closer look reveals fundamental differences. Rules 'merely govern alternations' and are more often than not 'phonetically unmotivated', his example being the English rule Velar Softening (Chomsky and Halle 1968) relating *ele[k]tric* and *electri[s]ity*: it has exceptions in the language, and the change from /k/ to [s] is a far from natural one in phonetic terms. Processes, on the other hand, 'reflect genuine limitations on

what we can pronounce'. They are part of the common acquisitional starting point, i.e., of UG (a term not used in his work):

- (10) [I]n its language-innocent state, the innate phonological system expresses the full system of restrictions of speech: a full set of phonological processes, unlimited and unordered. [. . .] A phonological process merges a potential phonological opposition into that member of the opposition which least tries the restrictions of the human speech capacity. [. . .] Each new opposition the child learns to pronounce involves some revision of the innate phonological system. [. . .] The child's task in acquiring adult pronunciation is to revise all aspects of the system which separates his pronunciation from the standard. If he succeeds fully, the resultant system must be equivalent to that of the standard speakers.

In the view I'm proposing, then, the mature system retains all those aspects of the innate system which the mastery of pronunciation has left intact. (Stampe 1969: 443–447)

Language acquisition does not mean acquiring just rules (and ordering them), as in Smith's case in the Chomsky and Halle framework, but in addition to acquiring the phonetically implausible rules such as Velar Softening, the processes supplied by UG are manipulated: they can be suppressed, limited, and ordered, when they are in conflict (conflicts arise, e.g., between absolute and contextually restricted processes: obstruents are voiceless irrespective of context because their oral constriction impedes the airflow required by voicing, but they also prefer being voiced in a voiced environment). Stampe recognises that to a certain extent his processes resemble Jakobsonian implicational laws, but the ground covered by the latter is just a subset of that covered by the 'innate processes', namely that of 'the phonemic inventory [. . .] unaffected by contextual neutralisations' (1969: 446). The laws are static, whereas processes are active, and make predictions about representations as well as inventories; moreover, they may be *contextually conditioned*. His assessment of Chomsky and Halle's markedness theory is that it 'drastically underestimates the number of processes which are innate' and makes 'totally unsupportable claims about the nature of phonology' (1973b: 45–46).

One or two brief examples may clarify how these ideas work in practice, focusing on language acquisition. Consider Stampe's comparison of the behaviour of coronals in Danish and Tamil (1973a: 14–16). In Danish, posttonic /t/ is pronounced as [d], which sounds just as pretonic [d]; in addition, posttonic /d/ is pronounced as [ð]. This is a counterfeeding situation similar to the one from Smith (1973) described in (8): a process of spirantisation precedes one of voicing. In Tamil, both processes have 'unrestricted' applications: /t/ is both

voiced and spirantised to [ð] postvocally. In Stampe's view, the learning going on here is just this: since the processes are part of Universal Grammar, and unrestricted application is the norm, the Danish child, in order to become an adult speaker, will have to learn the ordering of the two processes involved. That is all. The moral he draws from this example is (Stampe 1973a: 15–16):

- (11) What has not been adequately recognized in discussions of ordering is that its recognition depends on our assumption that certain phonological processes are possible and others impossible. In the Danish example the facts could be described by positing a single process which simultaneously voices postvocalic stops and, if they are already voiced, spirantizes them. In fact such an analysis would be simpler from the point of view of Danish in that the postvocalic context of both changes would be expressed just once. But from a language-universal point of view, this analysis is not satisfactory: it merely adds another process to the many we have to explain the existence of. [. . .] Furthermore, this lang[u]age-universal analysis provides us with an interesting prediction: that in no phonological system will [t] become [ð] unless [d] becomes [ð]. This follows from the assumption that the change of [t] to [ð] as in Tamil, involves two distinct processes, voicing and spirantization. If this is so, there is no way that [t] could be voiced and spirantized without [d] also being spirantized. This prediction is not overturned by any language of which I am aware.

His illustrations from language acquisition involve examples of 'one child having ordered two processes which another child has not' and, even, 'examples of a child actually performing the ordering' (Stampe 1969: 447, 1973a: 11–12, 16–17). The structure of two such cases is the following:

- (12) (a) Joan Velten (Velten 1943) pronounces *lamb* as [zab].
This pronunciation results from three distinct processes:
(i) delateralization: /l/ → [j] (this is common in children, cf. *lie* → [jai] by Hildegard in Leopold (1947))
(ii) spirantization: /j/ → [ʒ] (again see Hildegard's *you* → [ʒu])
(iii) depalatalization: /ʒ/ → [z] (as in Joan's own *wash* → [was])
But for Hildegard /l/ remains [j], it does not become spirantized as in Joan's speech. The difference between their pronunciations lies not in the processes but in their order of application. Hildegard does not apply spirantization to the output of delateralization; Joan does. Thus whereas Joan cannot pronounce [j] at all, Hildegard can pronounce it if she attempts to say /l/. This apparent oddity is no stranger than the Dane's inability to say postvocalic [d] except by aiming at [t]. (Stampe 1973a: 16–17)

- (b) Hildegard Leopold at 20 months said [du(ɪ)f] 'juice', [du] 'June', [do:i] 'Joey', beside [dʒuɪʒ] 'church', [dʒudʒu] 'choo-choo', by application of the processes (a) [dʒ] becomes [d], and (b) obstruent voices before a vowel (compare [du] 'to, do'). But at 19 months 'choo-choo' had been [dudu] [. . .] by unordered application of the same processes. (Stampe 1969: 447)

Stampe's ideas enjoyed considerable popularity in the 1970s; they were seen by many as a promising variant of the then strong 'Natural Generative Phonology' (NGP) movement,⁶ and enjoyed considerable popularity among many acquisitionists (see, e.g., Edwards and Shriberg 1983). As indicated (in our (8b)), Smith (1973) alludes to the possibility of reducing (many of) his rules to Stampe's processes.

The wisdom of such a move, however, was seriously doubted by detractors such as Kiparsky and Menn (1977). They argued that a careful look at first language acquisition and its developmental properties and stages does simply not support the main tenets of theories such as those by Jakobson and Stampe, which they lump together as 'rather deterministic': 'In these theories, there is no "discovery", no experimentation, no devising and testing of hypotheses. [T]he child's speech development cannot simply be viewed as a monotonic approximation to the adult model.' Much more typical of the acquisition process is that it proceeds as a 'problem solving' process: 'learning to talk [. . .] is a difficult task' in which 'the child must discover ways to circumvent the difficulties'. They see recognisable speech as a *target* which different children will attempt to reach by a variety of means, including the 'rules' typical of child phonology (Kiparsky and Menn 1977: 56–58):

- (13) [T]here are several ways of dealing with consonant clusters: deletion of all but one of the phonemes (in a stop-X or X-stop cluster, the one preserved is not always the stop), conflation of some of the features of the elements of the cluster (eg., *sm* > *m*, *fl* > *w*), insertion of a vowel to break up the cluster, metathesis (*snow* > *nos*), etc. [. . .]

Different children exclude definable classes of output by different means. When we observe such repeated 'exclusion', we conclude that these classes of outputs (clusters, certain co-occurrences, the 'third position', etc.) represent difficulties to the child, and the various rules of child phonology (substitutions, deletions, etc.) as well as selective avoidance of some adult words, are devices the child finds for dealing with those difficulties. [. . .]

The very diversity and the 'ingenuity' of these devices might indicate that early phonology should be regarded as the result of the child's active 'problem solving'.

They refer to Ingram (1974) for 'extensive consideration and exemplification of rules of child phonology', and add: 'Note that not all children avoid all of these difficulties; we present these as general tendencies rather than as universals.'

Clearly these remarks announce the notion of output constraint that appears in immediately adjacent work by Menn, as in the quote in (1) from Menn (1980), and the even earlier one below from Menn (1978: 162–164):

- (14) many rules are best seen in terms of the satisfaction of output constraints. [. . .] These constraints are interpretable as manifestations of the young child's limited ability to plan and execute a complex motor activity.

Thus, a constraint against consonants of different places of articulation can be met by consonant harmony (as in Amahl's [gɔ:k] for *talk*), but also by consonant deletion: [bu] for *boot*. A prohibition against fricatives in onsets can be met by deletion ([ɪf] for *fish*) or by metathesis ([nos] for *snow*).

These remarks also support Smith's original hunch in (8c), proposing the possible usefulness of Kisseberth's (1970) notion of 'conspiracy', applied to child phonology. In this latter pivotal paper, the author, having shown that in Yawelmani Yokuts rules of consonant deletion and vowel insertion all eliminate triconsonantal clusters, argued two things: first, that Yawelmani phonology appears to contain a 'conspiracy' towards a goal that can be formalised by a negative language-specific 'output constraint', namely *CCC; but second, that now the overall design of the grammar turns out to be very expensive in the sense that a much less well-balanced phonology with just an occasional deletion or insertion rule and no output constraint would be formally much cheaper. Neither Kisseberth nor Smith provide a solution to this problem, but the latter notices a parallel occurring in Amahl's phonology. According to him, an example is constituted by a series of rules in Amahl's speech eliminating consonant clusters in any position in the word:

- (15) a proposed *CC conspiracy in Amahl's speech, Smith (1973: 165–166)
- | | | |
|-----------|-----------|---------|
| (a) tree | → [di:] | Rule 16 |
| (b) taxi | → [gegi:] | Rule 21 |
| (c) meant | → [mɛt] | Rule 1 |
| (d) slip | → [ɪp] | Rule 7 |

Thus, a theoretical proposal by Kisseberth and an observationally based educated guess by Smith converge with Kiparsky and Menn's view of an essential property of the process of first language acquisition: the usefulness of output constraints.

Up until this point the main body of this section has focused on the function of UG in acquisition research, and on the way output constraints were introduced in the literature in that field. In the remainder we shall discuss two

further issues. First, that of the nature of the underlying forms of a child grammar, also in relation to the perception/production dichotomy; and second, the notion of parameter as a (partial) replacement of rules. With regard to the former issue, let us reconsider Smith's contention that the underlying forms of the child's phonological component are constituted, by and large, by the adult surface forms. This view did not go unchallenged for very long. In a review of Smith's monograph, Braine (1976) suspects that the influence of perception on the underlying form must be larger than assumed by Smith. He argues in favour of what he calls a 'partial perception hypothesis', which states that (1976: 492) 'the child's perception of words contains systematic biases, and is therefore only partly accurate'. His support for this contention contains the following elements, among others: there 'is now much evidence, from discrimination testing, that children's ability to discriminate perceptually among phonemic contrasts is far from perfect . . .'. Targeting Smith's *puzzle/puddle* chain shift, he focuses on the fact that the shift away from a coronal stop does not seem to be motivated by a production difficulty because that same coronal stop is in the output of the shift away from the fricative: 'If we drop the assumption that perception inevitably recovers the adult phonemes, [these] phenomena appear in a new light. In *puddle* → pu[g]le, one wonders if A's auditory system could be ranking the flapped or glottalized intervocalics as more similar to normal adult [velars] than normal [coronals]' (1976: 494). Braine's reinterpretation of Amahl's phonological behaviour implies a model with three separate components. First, '[a]uditory encoding laws would state how the child's auditory system transcribes the acoustic input into auditory attributes'. Second, on the output side there would be realisation rules similar to Smith's. Third, the model contains 'correspondence rules that map auditory features into the articulatory features that the child controls (or partially controls). These would specify the articulatory analysis made of words in the lexical store at any point in development, and in effect associate motor commands with the auditory features.' He admits that '[u]nfortunately, the correspondence layer is unlikely to be well-defined, given the "rare mistakes of articulatory coding" exhibited by Smith's son Amahl' (and other children, too, presumably).

A multi-step procedure of a different kind but in the same perception/production area of research was developed in work culminating in Menn (1978). Her proposal has become known as the 'two-lexicon model'. A child enters a perceived form in an input lexicon. These forms undergo reduction processes expressing the child's limited abilities; the output of these processes is stored in an output lexicon, which is completely redundancy-free. For production, a (hopefully small) set of 'production rules' or 'subroutines' convert this representation into the one entering the motor component, containing the articulatory instructions. In the case of *fish*, for instance, the perceived form (usually close to the adult surface form) is entered in the input lexicon. Reduction rules of the

child phonology turn this form into one underlying production: redundancy-free [ɪ, s] is stored in the output lexicon as an unordered pair of vowel and fricative, which must be ordered by the 'production rules'; in this case this ordering process is governed by a prohibition (output constraint) against fricatives in onsets. The principal empirical observation motivating the model is that of the occasional inertia of the rule-learning process. If a new phonological rule enters the child phonology, existing pronunciations sometimes persist, as if output forms serve as independent lexical items (Menn and Matthei 1992: 213):

- (16) An example from Daniel (Menn 1971) makes this clear: fairly early on, Daniel produced 'down' and 'stone', both very frequent words, as [dæwn] and [don], respectively. Sometime later, he began to show a nasal harmony rule, producing 'beans' as [minz] and 'dance' as [næns]. For a while after this point, 'down' and 'stone' were maintained in their nonassimilated forms; then the forms [næwn] and [non] began to appear, in free variation with them. Finally, [næwn] and [non] were triumphant – [dæwn] and [don] disappeared.

It has become clear, however, that this two-lexicon design cannot be maintained in this relatively naive form. Empirical evidence has been put forward that not the purported output lexicon but simply the 'classical' input lexicon is active when the child starts to develop morphophonemic alternations. Smith's data in (6) constitute a case in point, and a similar example, originating from Stemberger (1993), is represented in Bernhardt and Stemberger (1998: 48–49) in the following manner:

- (17) At 2;10, Gwendolyn generally formed the past tense of vowel-final words by adding [d], as in adult English. There were two exceptions to this, however. First, forms that are irregular in adult speech, such as *threw*, could be produced correctly as irregulars or could be regularized (*throwed*), as is common in child language. Second, words that are consonant-final in adult speech but were vowel-final in the child's speech, such as *kiss* /kɪs/ [tʰi:], did not have *-d* added in the past tense [. . .].

The two-lexicon approach seems to predict that all words that are vowel-final in the child's speech should be treated the same for the creation of past tense forms. Since a final /d/ is added to words like *pee*, it should also be added to words like *kiss* [ti:]. This prediction fails, suggesting that the two-lexicon approach is inadequate.

Another argument points out that the model is hard put to account for between-word phonological processes in child language, especially if the same generalisations cover both words and simple syntactic constructions in early child

speech: a lexical approach seems ill-equipped to capture such cases (Menn and Matthei 1992: 223). It seems the explanation of the selective propagation of new child speech rules must lie elsewhere. In the latter paper, Menn and Matthei turn to 'connectionist' models in order to maintain 'what's good' about their approach. Menn (this volume) explains her most recent position.

Finally, in theoretical linguistics of the 1980s attempts were made to replace the view of a grammar component as a rule system by one involving a set of parameters, i.e., a set of choices specified by UG and fixed on a language-particular basis given linguistic experience in that language. Ideally in this approach the grammar becomes a collection of fixed parameter-settings. It was recognised from the very outset that not only does this view provide a framework for the study of language typology but it also has implications for the study of language learning (cf. Chomsky 1981a: 8-9, 1981b: 3-4, 1986: 52, 145-6):

- (18) It is natural to assume, then, that universal grammar consists of a system of principles with a certain degree of intricacy and deductive structure, with parameters that can be fixed in one or another way, given a relatively small amount of experience. Small changes in the values assigned to parameters in a rich system may lead to what appear to be radically different grammars, though at a deeper level, they are all cast in the same mould. [. . .]

The ideal is to reach the point where we can literally deduce a particular human grammar by setting parameters of universal grammar in one or another of the permissible ways. The process of so-called 'language learning' can then be naturally regarded in part as a process of fixing these values; when they are fixed the 'learner' knows the language generated by the grammar that is determined by universal grammar, specified in this way. [. . .]

The parameters must have the property that they can be fixed by quite simple evidence, because this is what is available to the child [. . .]. Once the values of the parameters are set, the whole system is operative. [W]e may think of UG as an intricately structured system, but one that is only partially 'wired up'. The system is associated with a finite set of switches, each of which has a finite number of positions (perhaps two). Experience is required to set the switches. When they are set, the system functions.

As a typological example, consider the parametric theory of syllable structure presented in Kaye (1989: 54-57). In (19), (a) gives the set of parameters and (b) some of the languages that fill the slots of the system (where '0' is no, and '1' is yes):

- (19) [I]f we take all the languages of the world, how many syllable types must we allow for? The answer is surprising. Three parameters suffice to define every extant syllable type. [. . .] They are:
- (10) (1) Does the rime branch? [no/yes]
 (2) Does the nucleus branch? [no/yes]
 (3) Does the onset branch? [no/yes] [. . .]
- (13) (a) 000 No branching rimes, nuclei or onsets: Desano
 (b) 100 Branching rimes, but no branching nuclei
 or onsets: Quechua
 (c) 10[1] Branching rimes and nuclei, no branching
 onsets: Arabic
 (d) 101 Branching rimes and onsets, no branching
 nuclei: Spanish
 (e) 111 Branching rimes, nuclei and onsets: English

The task of the language learning child is to fix the parameter settings. It is usually assumed that parameters are entered in UG with a 'default setting'. Parameter (10.1), for instance, covers the difference between open and closed syllables, but languages actually come in two types: those with open syllables and those with open and closed ones. The 'poverty of the stimulus' argument applied to this case results in the default setting of NO for this parameter: the presence of closed syllables can then be learned on the basis of positive evidence. Among the many findings of Fikkert (1994), an elaborate study of the acquisition of Dutch prosody in this framework, is that in this language, which has branching rhymes, young children first omit word-final consonants ([pʊ:] for *poes* 'pussycat', [ka:] for *klaar* 'ready'), and start producing them some two months later. This can be seen as the setting of this parameter away from the default value.

Similar parameters have been proposed for the typology of word stress systems (Hayes 1980/81, Prince 1983) and the acquisition of word stress has been studied in considerable detail in Dresher and Kaye (1990), Fikkert (1994), and Dresher (1999), both from the acquisitional-empirical and the principled-theoretical angle. In spite of studies such as these, it seems that the parametric approach to acquisition has not (yet) grown to its full potential, and it seems as if parameters are sometimes seen as a poor man's principles: fixed parameter settings are intended to represent alternatives to rules (such as stress rules in the Chomsky and Halle mould) but never fully replace them (see Piggott 1988), and from the UG point of view parameters could be seen as 'failed principles', cf. Archangeli (1997: 26):

- (20) The 'inviolable' principles of syntax have themselves proved to be problematic in that inviolability has been purchased at the cost of

a variety of types of hedges. [S]ome principles are 'parameterized', holding in one way in one language and in another way in another language. Other principles have peculiar restrictions built-in.

2. Linking phonological universals and acquisition through Optimality Theory

It stands to reason that a linguist working on acquisition issues does not know in advance whether her or his hypotheses had best be formulated within the rule-based framework (Chomsky 1965, Chomsky and Halle 1968), within that of Principles and Parameters Theory (Chomsky 1981a, 1981b), a combination of these (e.g., Halle and Vergnaud 1987), or something else, for instance within constraint-based Optimality Theory (Prince and Smolensky 1993). This volume argues that the last-mentioned approach, in itself and in comparison with the other frameworks, allows a very fruitful line of attack.

2.1 *Optimality Theory*

Studies in phonological acquisition can be said to be focused on three priorities. First, to account for universal patterns in phonological acquisition, researchers tried to establish a *substantive* theory of phonological markedness (originating with Jakobson, with incarnations in Chomsky and Halle's theory of markedness, and in Stampe's Natural Phonology). Second, especially in generative approaches, the emphasis was on developing a *formal* theory of phonology that would characterise the child's developing competence as a set of cognitive states leading up to the adult grammar, each state encoding a grammar itself (conceived as a set of linearly ordered rewrite rules in Chomsky and Halle's and Smith's views, a set of unsuppressed natural processes in combination with *ad hoc* rules in Stampe's model, with the need for output constraints being recognised by Menn and others). Finally, recent parameter-based approaches emphasised learnability as a vital issue in acquisition, recognising the need for a theory of *learnability* to explain relations between the learner's input and phonological development.

In this section we will see what Optimality Theory (OT) has to offer in these areas. After presenting an outline of OT, we shall look into the central role of markedness principles in OT. Next, we consider ways in which OT, as a formal theory of grammatical interactions, solves a number of classical problems for derivational rule-based models of phonology, focusing on the duplication problem and the conspiracy problem. This discussion leads us naturally to the notion of typological variation, and how it is captured in OT. Finally, we shall see how OT grammars are formally set up in such a way so as to be learnable. At the end of this section, we are ready to approach the issue of how OT accounts for the

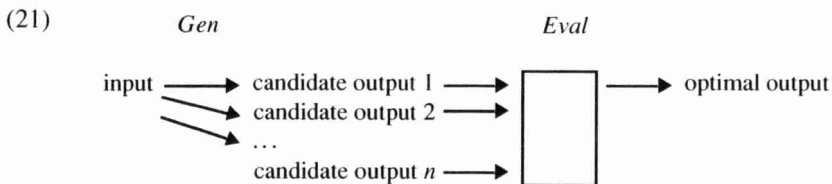
relation between phonological typology and phonological acquisition, which will be taken up in section 3.

2.2 An outline of Optimality Theory

Optimality Theory (Prince and Smolensky 1993, McCarthy and Prince 1993) models linguistic well-formedness by a set of conflicting constraints which state (positive or negative) absolute demands about surface forms.⁷ Constraints are intrinsically in conflict, as they impose hard general structural requirements that cannot be simultaneously satisfied by any logically possible form. Conflicts between constraints are regulated in grammars by imposing a ranking – or constraint hierarchy. The hierarchy is strict, with any constraint taking priority over all lower ranked ones. Consequently, violations of constraints are allowed only to avoid violation of higher ranked ones.

As in earlier derivational theories of phonology, discussed in section 1, inputs (underlying representations) are mapped onto outputs (surface representations). But unlike derivational mappings in earlier theories, which involved a sequence or linearly ordered rules, OT mappings are single-step derivations; for a given input, the grammar selects the ‘optimal’ output form from an infinite set of candidate outputs, which are generated by the constraint component *Gen*. The assumption that the grammar generates and evaluates all logically possible candidate analyses for a given input is called *Freedom of Analysis*.

Each output candidate generated by *Gen* incurs different violation(s) for individual constraints. Accordingly, candidate outputs differ from one another in their ‘harmonic’ well-formedness, that is, the degree to which they meet a set of ranked conflicting constraints – a constraint hierarchy. The evaluation function of the grammar (*Eval*) imposes a harmonic ranking among candidates, with the most harmonic candidate at the top and the least harmonic one at the bottom. The winning (‘optimal’) candidate is the one that best matches the overall constraint hierarchy. Hence, violations of lower ranked constraints will be tolerated in order to satisfy higher ranking ones.



An evaluation of output candidates by a set of ranked constraints can be displayed by a ‘tableau’. The tableau in (22) shows three hypothetical output candidates (a-b-c) in competition, their relative well-formedness measured by three ranked constraints (C_1 - C_2 - C_3). The optimal output is the one that is ‘more

harmonic' in all its pairwise competitions with other candidates; in each pairwise competition, the more harmonic candidate is the one that performs better on the highest-ranking constraint that distinguishes between them (McCarthy 2001: 3). The optimal candidate *b* beats its competitor *a* as it performs better on the highest-ranking constraint distinguishing between them, top-ranked C_1 . The winner also outperforms candidate *c* as it has fewer violations of the highest-ranking constraint distinguishing between C_2 .

(22) Simple constraint interactions

	Constraint 1	Constraint 2	Constraint 3
Candidate a	*!		
☞ Candidate b		*	*
Candidate c		**!	

This tableau shows that the optimal candidate *b* is not the one having no or the smallest number of violation marks across columns. According to such a criterion, candidate *a* would have been the winner. Instead what matters is seriousness of violations, relativised to constraint ranking: Competitor *a* is eliminated due to its single violation of a top-ranked constraint C_1 . Also, it is not the number of constraints violated by a candidate which matters, but rather the distribution of marks over cells: Candidate *c* loses because of its double violation of a single constraint C_2 , even though it has no violations of C_3 .

Many researchers assume that all constraints in grammars of natural languages are part of UG's universal inventory of constraints called *Con* (Prince and Smolensky 1993). According to this view, grammars differ exclusively in the ranking of constraints. The central assumption that typological variation is due to differences in ranking between constraints in a universal inventory has consequences for language acquisition, as it restricts the learner's search space, while establishing a direct relation between phonological typology and acquisition.

The alternative view on the status of constraints is that these emerge from articulatory and perceptual factors which are active during acquisition (Boersma 1998, Hayes 1999). This functional approach also predicts a strong relation between typology and acquisition, because universal functional factors govern the process of selection of constraints by the learner. This implies that constraints, the ingredients of typology, should not differ between languages in arbitrary ways. (See section 3.2 for discussion.)

On the standard view (originating with Prince and Smolensky 1993) constraints in *Con* fall into two broad classes, known as *markedness constraints* and *faithfulness constraints*. Interactions of these constraint types model the extent to which marked structures of certain kinds are allowed in a language.

Markedness (or 'structural') constraints express universal preferences for certain types of structure, such as syllables with (rather than without) onsets, voiceless (rather than voiced) obstruents, or oral (rather than nasal) vowels. The principal motivation for markedness constraints are Jakobsonian implicational universals: a structure S_u is unmarked (with respect to another structure S_m) if, for every language, the presence of S_m implies the presence of S_u . For example, every language allowing onsetless syllables also allows syllables with onsets, while every language allowing voiced obstruents also allows voiceless ones.

Markedness constraints are often subjected to the criterion of being 'grounded' (in the sense of Archangeli and Pulleyblank 1994) in phonetic factors (production and/or perception). As we saw, OT researchers take rather different positions on the degree to which constraints should be functionally motivated, and on the related issue of whether constraints are universal (part of UG) or emergent from functional properties.

Faithfulness constraints make a rather different type of requirement of surface forms: that they match specific properties of other forms, for example their lexical input. Their effect is to prohibit deletions, insertions, featural changes, or other changes in mappings from inputs to outputs. Faithfulness constraints are the natural antagonists of markedness constraints, since the former preserve lexical properties that the latter may ban at the surface.

Correspondence Theory (McCarthy and Prince 1995, 1999) implements faithfulness as constraints on corresponding segments in paired representations, such as input and output. For example, the constraint MAX-IO requires every segment in the input to have a correspondent in the output ('no deletion'). The constraint DEP-IO expresses a mirror-image requirement that every output segment has an input correspondent ('no insertion'), while IDENT[F]-IO requires corresponding segments to share specifications for the feature [F] ('no featural changes').

Besides input-to-output faithfulness, another type of faithfulness relation has been argued for in the literature: constraints requiring identity between an output form and another output form (*output-to-output* faithfulness). While originally motivated for reduplication (McCarthy and Prince 1995, 1999), output-to-output faithfulness has been, controversially, generalised to other domains, comprising, for example, instances of cyclic rule application proposed from early generative phonology onwards, as well as paradigm uniformity (Benua 1997, Burzio 1996, Kenstowicz 1996, Kager 1999; see Kiparsky 1999 for a different view).

In a language, the presence versus absence of marked structures of various types thus depends on the relative ranking of M(arkedness) constraints and F(aithfulness) constraints; when Markedness dominates Faithfulness ($M \gg F$), the unmarked structure surfaces, whereas the opposite ranking $F \gg M$ suppresses

the unmarked structure. This can be illustrated for nasality in vowels, a marked property on typological grounds: all languages have oral vowels, whereas not all languages have nasal ones. The markedness constraint militating against nasal vowels, $*V_{NAS}$ 'no nasal vowels' competes with a faithfulness constraint IDENT-IO(nas), requiring all surface vowels to preserve the specification of [nasal] of their input correspondents. Simple permutation of these constraints produces two grammars, one suppressing nasality where it is specified in the input, and another grammar which allows nasality of input vowels to surface, cf. (23).

(23) Grammar 1: input nasality suppressed

Input: /bã/	$*V_{NAS}$	IDENT-IO(nas)
bã	*!	
\rightarrow ba		*

Grammar 2: input nasality preserved

Input: /bã/	IDENT-IO(nas)	$*V_{NAS}$
\rightarrow bã		*
ba	*!	

Grammar 1, which ranks the markedness constraint $*V_{NAS}$ above the faithfulness constraint IDENT-IO(nas), effectively prohibits a contrast between oral and nasal vowels. This contrast is supported in Grammar 2, which has the reverse ranking. Which features are 'contrastive', and which are 'noncontrastive', then depends on the ranking of specific markedness constraints and faithfulness constraints: $F \gg M$ supports contrasts, while $M \gg F$ neutralises contrasts. (To capture contextual neutralisation, as well as allophonic variation, markedness constraints must be relativised to context.)

As compared to earlier theories (in particular, standard Generative Phonology), OT directly encodes markedness into grammars, with markedness constraints constituting the substance out of which phonologies are built. Consequently, the markedness (or 'naturalness') of phonological processes and segment inventories need no longer be attributed to a grammar-external evaluation measure, as it had been in *SPE*. While deviating from classical Generative Phonology, OT is on a par with Natural Phonology (see section 1) in giving a central function to markedness principles. OT differs from Natural Generative Phonology, however, by taking hierarchically ranked constraints rather than linearly ordered (natural) processes to be the core device.

Consequently, Jakobsonian implicational universals about segment inventories can be brought within the scope of OT's grammatical explanation. The typological generalisation, for example, that all languages have oral vowels (whereas no language has only nasal vowels) is simply due to the logically possible interactions of a markedness constraint ('no nasal vowels', cf. 23) and a faithfulness constraint ('preserve input nasality'): critically, no possible ranking bans oral vowels across the board. Strong typological predictions follow from simple constraint interactions.

Unlike its ancestor theories, OT models phonological generalisations completely at the surface. The assumption that no grammatical restrictions are stated at the level of lexical representation is called *Richness of the Base* (Prince and Smolensky 1993, Smolensky 1996a, Smolensky, Davidson, and Jusczyk, chapter 10, this volume). OT thus abandons morpheme structure rules, a well-known but not uncontroversial device of SPE-type phonological theory whose task it is to state the generalisations holding at the level of the lexicon, where these constraints filled in unspecified (predictable) feature values (thus diminishing the cost of the lexicon; recall the evaluation measure counting symbols, mentioned in section 1). OT thus places the burden of accounting for generalisations about the phoneme inventory and phonotactics on surface constraints in a single component which also accounts for phonological alternations. This surface-oriented architecture offers a radical and principled solution to the *duplication problem*, the phenomenon that morpheme structure rules (capturing 'static phonology') are frequently redundantly duplicated by phonological rules that map lexical representations to surface representations (accounting for alternations – the 'active phonology').⁸

As an example, let us consider Dutch voicing assimilation, a set of processes (of regressive and progressive assimilation) which has the effect of eliminating obstruent clusters of mixed voicing at the surface, most notably within the domain of phonological word. That is, in phonological words no clusters such as [kd] or [bt] occur. Voicing agreement is met dynamically (e.g., in phonological alternations of the past tense suffix /-də/, which assimilates in voicing to the stem-final consonant in *maakte* /ma:k+də/ → [ma:ktə] 'made') as well as 'statically' (obstruent clusters share voicing when belonging to a single lexical item, as in *dokter* [dɔktər] 'doctor'). In classical Generative Phonology, the lack of tauto-morphemic disharmonic clusters such as hypothetical */dɔkdər/ was captured by a morpheme structure rule stating voicing agreement in obstruent clusters at the level of lexical representation. Having both a dynamic and a static version of voicing assimilation amounts to a duplication, however, missing a generalisation.

In OT, such duplication is avoided because a single constraint hierarchy captures the generalisation: the markedness constraint AGREE-VOICE 'Obstruent clusters agree in voicing' dominates the faithfulness constraint

IDENT-IO(voice) 'Output segments preserve their input voicing specification'. This accounts for the voicing alternation:

- (24) Voicing assimilation in 'dynamic' mode (alternating past tense suffix)

Input: /ma:k+də/	AGREE-VOICE	IDENT-IO(voice)
[ma:kdə]	*!	
☞ [ma:ktə]		*

Note how the same ranking accounts for the static generalisation holding for tauto-morphemic clusters (such as *dokter*). Under Richness of the Base conditions on lexical representations are not required, nor can they be stated. Regardless of whether input clusters have (dis-)harmonic voicing specifications, the grammar forces their surface correspondents to be harmonic:

- (25) Voicing assimilation in 'static' mode (tauto-morphemic context)

Input: /dɔktər/	AGREE-VOICE	IDENT-IO(voice)
[dɔkdər]	*!	*
☞ [dɔktər]		
Input: /dɔkdər/	AGREE-VOICE	IDENT-IO(voice)
[dɔkdər]	*!	*
☞ [dɔktər]		

Whereas the harmonic input /dɔktər/ is faithfully mapped onto a licit output, a hypothetical disharmonic input /dɔkdər/ cannot surface unmodified (*[dɔkdər]), and undergoes voicing assimilation.⁹ In sum, voicing assimilation need not be stated twice, but acquires the status of a single grammatical generalisation, which solves the duplication problem.

The surface-oriented architecture of Optimality Theory also offers a solution for another problem for classical Generative Phonology, known as the *conspiracy problem*, noted in section 1. Kisseberth (1970) observed that within grammars different rules conspire towards a common goal: they collectively avoid a 'marked' pattern (for example, CCC clusters are broken up by epenthesis, or reduced by consonant deletion) or establish an 'unmarked' pattern (for example, syllable onsets are created by consonant epenthesis, vowel coalescence, etc.). Conspiracies are a problem for classical derivational phonology: the functional unity between conspiring rules is evident, but left without any

formal expression in the grammar. Optimality Theory accounts for conspiracies since changes to inputs are always triggered by the necessity to avoid violations of a high-ranking markedness constraint; a range of resolution strategies is thus expected for principled reasons. Voicing assimilation in Dutch again serves as an example. Obstruent clusters C_1C_2 which disagree in voicing are avoided by regressive voicing assimilation (in case C_2 is a plosive, e.g., *zakdoek* [zɑgdʊk] ‘handkerchief’), alternatively, by progressive assimilation (in case C_2 is a fricative, e.g., *diepzee* [dipse:] ‘deep sea’, or C_2 belongs to an inflectional affix, e.g., *maakte* [ma:ktə] ‘made’). These repair strategies share a common objective, that is, avoid violation of AGREE-VOICE.

If both lexical representations in (25), /dɔktər/ and /dɔkdər/, lead to the same output, which one is taken to be the correct one? Under *Lexicon Optimisation* (Prince and Smolensky 1993: 192), if a grammar maps multiple distinct lexical representations (say, L_1 and L_2) onto a single surface form S , the lexical representation is selected whose lexical-to-surface mapping is most harmonic in terms of the grammar. For the purpose of selecting lexical representations (a ‘surface-to-lexical’ mapping) the grammar is now used in *backward* mode. Since candidate mappings for a given surface form are all equally marked as to surface well-formedness, all candidates share the same set of violation marks on markedness constraints. Hence, selection of the optimal lexical form is solely carried out by faithfulness: the most harmonic mapping is the one which minimally violates faithfulness constraints. This is the ‘identity’ mapping.

(26) Lexicon Optimisation selecting an underlying representation

Output: [dɔktər]	AGREE-VOICE	IDENT-IO(voice)
/dɔkdər/ → [dɔktər]		*!
☞ /dɔktər/ → [dɔktər]		

Lexicon Optimisation selects lexical representations which equal surface forms only in the case of non-alternating morphemes. For alternating morphemes, such as the past tense suffix [-tə]~[-də] in Dutch, the standard generative assumption of a single underlying representation of surface alternants entails that the lexical-to-surface mapping is unfaithful for at least one alternant. For alternations Lexicon Optimisation needs to be supplemented by other principles (see Tesar and Smolensky 2000: 77–83 and Hayes, chapter 5 in this volume, for some discussion of learning alternations).

The preceding discussion suggests that an OT grammar serves two duties: licensing licit outputs; and filtering out illicit ones. (See Hayes, chapter 5 this volume, for further discussion.)

First, the grammar guarantees that licit outputs are faithfully mapped from segmentally identical input forms. In such cases, the grammar functions solely as a passive filter, licensing lexical items of the language, and allowing the lexicon to be productively extended by items conforming to the language's phonological requirements. When the grammar maps an input onto an unchanged output, it performs what we may refer to as an 'identity' mapping. We may now define the notion 'possible word' in a language as an output that, *when taken as an input*, would undergo the identity mapping. Hence, to subject a form *F* to the 'possible word test', *F* is submitted as an input to the grammar. If *F* is mapped onto an output form *F'* which is non-distinct from its input, *F* passes the test, showing that *F* is a possible word of language *L*.¹⁰

Second, the grammar functions actively as a filtering device by prohibiting illicit forms from surfacing. Note that the grammar does not filter out an illicit form by 'blocking' (prohibiting it from appearing at the surface), but rather by mapping it onto a modified licit form, a *non-identity* mapping. Submitting an illicit form *F* to the 'possible word' test, we feed it into the grammar as an input, which maps it onto an output *F'* which is distinct from *F*. (In a tableau, this shows by violation marks incurred by *F'* on one or more faithfulness constraints.)

There are various sources of evidence for non-identity mappings. The first, classical type of evidence comes from automatic alternations in the shapes of morphemes that depend on phonological context, such as alternations in voicing of obstruents in Dutch, triggered by the markedness constraint *AGREE-VOICE*.

Another type of evidence for non-identity mappings comes from *loanword adaptation*, the familiar phenomenon that words borrowed from another language are modified so as to meet an inviolate phonological requirement of the borrowing language. For example, English speakers tend to repair onset clusters which are phonotactically illicit in their language, such as /kn/, by vowel epenthesis (e.g., *Evel K[ə]nievel*). Loanword adaptations, because of their automatic, forced character and the broad consistency regarding choice of repair strategy (e.g., deletion, insertion, featural change of segments) applied by different speakers, give evidence for the view that it is due to an internalised grammatical system (Hyman 1970).¹¹

Additional evidence for the automatic nature of non-identity mappings comes from *second language phonology*. For example, Dutch learners of English characteristically display final devoicing, neutralising contrasts such as *bit* versus *bid*. This shows the familiar effect of transfer of the first language (L1) into a second language (L2), in this case the *M » F* ranking for final devoicing: **VOICEDCODA » IDENT-IO(voice)*. More strikingly, cases are known in which L2 learners display similar mappings which apparently cannot be explained by transfer from their native language. (See discussion in Stampe 1969, Donegan and Stampe 1979, and section 4.4 below.) For example, many

Mandarin Chinese speakers learning English as their L2 go through a stage in which they display final devoicing (Broselow *et al.* 1998), just like Dutch learners of English. However, final devoicing is not observable in Mandarin as this language lacks words ending in obstruent codas. This *emergence of the unmarked* in L2 phonological acquisition follows from two assumptions about acquisition. On the basis of L1 acquisition, Smolensky (1996a, 1996b) has proposed that in UG's initial state, markedness constraints generally dominate faithfulness constraints. During acquisition, the initial state is transformed into an adult grammar by a step-wise process of constraint re-ranking, which is triggered by positive evidence in the form of input from the target language. (This process will be discussed in more detail in section 2.4.) However, Mandarin speakers never receive any positive evidence to change the initial state's $M \gg F$ ranking of final devoicing, because there are no overt effects of final devoicing in their native language's input, due to a high-ranked constraint banning all coda obstruents. The observations that Mandarin L2 learners of English display final devoicing then follows straightforwardly from a second assumption, namely that of *full transfer* of the learner's L1 grammar into the initial state of the L2 grammar. The $M \gg F$ ranking for final devoicing (covertly present in Mandarin) automatically transfers into the L2 learner's grammar, where it emerges as soon as the markedness constraint banning obstruent codas is demoted.

A general conclusion to be drawn from this example is that constraint rankings of the type $M \gg F$ may be *covertly present* in OT grammars. Such covert rankings may be exposed (i.e., become active in non-identity mappings) when situations change minimally. Exposure may be triggered, for example, by the demotion of a higher ranking constraint that obscured the $M \gg F$ ranking (as we saw in the L2 acquisition case), or by feeding the grammar with data that are not 'normally' fed into it (as in the case of loanword adaptations). This mechanism of exposure of an obscured markedness constraint is known as 'the emergence of the unmarked' (McCarthy and Prince 1994). Schematically, the $M \gg F$ ranking is obscured in (27) by an 'obscuring constraint' C competing with M. Note that C may be a markedness constraint itself, or a faithfulness constraint.

(27) $M \gg F$ ranking obscured by high-ranking constraint C

	C	M	F
Candidate 1	*!		*
☞ Candidate 2		*	

The $M \gg F$ ranking may become activated by a demotion of the obscuring constraint below M, as in (28):

(28) TETU by demotion of obscuring constraint C

	M	F	C
☞ Candidate 1		*	*
Candidate 2	*!		

Alternatively, the M » F ranking may be activated by a vacuous satisfaction of C, leaving the original ranking unaffected, as in (29):

(29) TETU by vacuous satisfaction of obscuring constraint C

	C	M	F
☞ Candidate 1			*
Candidate 2		*!	

Cases of the latter type are well known in reduplication systems (McCarthy and Prince 1995, 1999). In CV reduplication, a typologically common type, the affix copies a segmental portion from the base which equals a single open syllable (CV). For example, in Nootka (Stonham 1990), the reduplicated form *či-čims-'i:h* 'hunting bear' has a CV affix [*či*] which copies a substring of its base [*čims-'i:h*]. This preference for unmarked syllable structure in the affix, showing the activity of NOCODA (constraint M in 29), is not a general property of Nootka, however, a language which otherwise allows for closed syllables. This property is due to domination of NOCODA by MAX-IO (constraint C in 29), which prohibits C-deletion as the general means of attaining open syllables. The CV affix in reduplication is not due to *deletion*, however. Due to the *copying* nature of reduplication, the affix lacks a proper lexical segmental representation. The syllabic shape of the Nootka affix, unchecked by MAX-IO, promptly gravitates to CV to satisfy NOCODA in an emergence of the unmarked. In this example, the role of the dominated faithfulness constraint (F in 29) is taken by MAX-BR, requiring that 'Every segment in the Base have a correspondent in the Reduplicant' (McCarthy and Prince 1995).

(30) The emergence of the Unmarked in Nootka reduplication

Input: /RED-čims-'i:h/	MAX-IO	NO-CODA	MAX-BR
(a) ☞ <i>či-čim.s'i:h</i>		**	****
(b) <i>čim.s'i:h-čim.s'I:h</i>		***!*	
(c) <i>či-či</i>	*!***		

The Emergence of the Unmarked constitutes a powerful argument for OT's central assumption that grammars of natural languages consist of hierarchies of violable universal (markedness and faithfulness) constraints. The Emergence of the Unmarked (TETU) effects of this kind have been observed in a wide range of situations, and most interestingly for our purposes, in (L1 and L2) phonological acquisition; more examples will be discussed in section 3.1.

2.3 *Typological variation in OT*

On the (standard) assumption that constraints are universal – and in fact innate – the view of cross-linguistic variation ('typology') in OT is straightforward. If all languages build their grammars from the same substance – the full content of *Con*, UG's constraint inventory – the locus of typological variation must be in the arrangement of this substance, in the constraint rankings. Accordingly, OT's central claim about typology is that cross-linguistic differences arise by re-rankings of a set of universal constraints.

Whereas the grammars of individual languages are basically free to rank constraints of *Con* in specific hierarchies, it has been argued as early as Prince and Smolensky (1993) that UG may impose restrictions on possible hierarchies. Formally and functionally related constraints may have fixed rankings, which cannot vary cross-linguistically. For example, the assumption is often made that the markedness constraints on place of articulation are universally ranked in a sub-hierarchy with *LABIAL and *DORSAL outranking *CORONAL, which accounts for the cross-linguistically observed unmarkedness of coronals. Other constraints interact with this sub-hierarchy, producing language-specific rankings in which universal markedness relations between places of articulation are enforced in different ways.

The view that typological variation is due to constraint (re-)ranking is notably different from that taken in parametric theory (reviewed in section 1), which explains typological variation in terms of (binary) parameter values. In parametric theory, switching a parameter 'off' implies total inactivity of the requirements involved. For example, a language whose grammar sets the coda parameter to 'off' is predicted to allow codas freely. In contrast, an OT constraint which is dominated is not necessarily inactive: even in its dominated position, it may continue to exert influence on the selection of output candidates. Given the chance, a dominated markedness constraint will jump into activity, a cross-linguistically well-attested effect, and another example of the emergence of the unmarked (McCarthy and Prince 1994), as we saw in the previous section.

The *factorial typology*, a major notion of OT, allows an explicit connection between language typology and acquisition. A factorial typology of a set of constraints is defined as all the logically possible rankings of these constraints. For example, a set of three constraints, C_1 , C_2 , and C_3 , can be ranked in six different ways:

(31) A factorial typology of three constraints

- (a) $C_1 \gg C_2 \gg C_3$ (c) $C_2 \gg C_1 \gg C_3$ (e) $C_3 \gg C_1 \gg C_2$
 (b) $C_1 \gg C_3 \gg C_2$ (d) $C_2 \gg C_3 \gg C_1$ (f) $C_3 \gg C_2 \gg C_1$

Computing the factorial typology for a set of constraints allows testing its adequacy against typological evidence: in principle, every distinct ranking predicted by the factorial typology should match (part of) the grammar of some natural language. Factorial typologies, however, grow quickly with the size of the constraint set, for n constraints can be ranked in $n!$ different ways. Since the factorial typology for a set of the size of *Con* is huge, the task of typologically verifying all predicted grammars poses many problems. Nevertheless, if we consider smaller sets, concentrating on a single typologically variant property (for example, syllable typology or stress typology), factorial typologies usually shrink to sizes small enough to allow for full typological verification.

Consider, for example, a factorial typology of three constraints involved in patterns of obstruent voicing. In addition to the constraints *NO VOICED CODA* and *IDENT-IO(voice)* that were discussed earlier, we assume a third constraint, the *VOICED OBSTRUENT PROHIBITION (VOP)*, a general markedness constraint banning voiced obstruents across the board. There are six logically possible rankings, whose characteristic patterns are indicated. Only three distinct patterns emerge:

(32) A mini-typology of voicing

- | | |
|---|-----------------------|
| (a) <i>NO VOICED CODA</i> » <i>VOP</i> » <i>IDENT-IO(voice)</i> | no voiced obstruents |
| (b) <i>NO VOICED CODA</i> » <i>IDENT-IO(voice)</i> » <i>VOP</i> | final devoicing |
| (c) <i>VOP</i> » <i>NO VOICED CODA</i> » <i>IDENT-IO(voice)</i> | no voiced obstruents |
| (d) <i>VOP</i> » <i>IDENT-IO(voice)</i> » <i>NO VOICED CODA</i> | no voiced obstruents |
| (e) <i>IDENT-IO(voice)</i> » <i>NO VOICED CODA</i> » <i>VOP</i> | full voicing contrast |
| (f) <i>IDENT-IO(voice)</i> » <i>VOP</i> » <i>NO VOICED CODA</i> | full voicing contrast |

All three patterns are typologically attested, in languages such as Finnish (32a), Dutch (32b), and English (32e). The full typology of obstruent voicing patterns is, of course, descriptively richer, implying that more constraints need to be assumed than those considered here (see, e.g., Lombardi 1999). Still, this example serves to illustrate a general point: since many of the rankings in factorial typologies collapse into a single pattern, the class of typologically predicted patterns is much smaller than the number of rankings.

Testing new constraints by calculating their factorial typologies (in interaction with a set of well-established constraints) is methodologically useful, if not imperative. If constraints are universal (i.e., present and ranked in every grammar), then adding a new constraint to the universal inventory may increase the predicted typology. Hence, the merits of a constraint cannot be

exclusively evaluated on the basis of how well it functions in a particular grammar, but need to be projected on a larger, typological, scale. Arguably, every constraint should pass the factorial typology test: it should not overgenerate by predicting systematically unattested grammatical patterns. This establishes a second major criterion for validating a new constraint (the first being grounding). A third criterion, based on language acquisition, will be discussed in section 3.

2.4 *The learnability of Optimality theoretic grammars*

A major question is how grammars (or partial grammars, such as phonologies) can be learned on the basis of positive evidence in the learner's input. In OT, a theory of ranked constraints, any answer to this question must necessarily take into account the learnability of constraint rankings. This important issue was recognised and addressed by Tesar and Smolensky (1993, 1998, 2000). They developed a Constraint Demotion Algorithm (CDA) which ranks a set of constraints on the basis of positive input. The algorithm can deduce information about constraint ranking from surface forms which it is fed with, while an extension of the algorithm can learn to assess the correct representation of a raw surface form given in the input. The second type of algorithm which we shall discuss, the Gradual Learning Algorithm (GLA), can deal with variation in the input, and can also account for gradual well-formedness (Boersma 1997, 1998, Boersma and Hayes 2001).¹²

The idea underlying Tesar and Smolensky's CDA is that the learner can work out the target ranking of the language that (s)he is learning from inspecting patterns of constraint violations in the forms that (s)he encounters. The learner first makes the necessary assumption that all forms which (s)he hears are grammatical, hence optimal under the ranking of the target language. This assumption allows the learner to infer that any constraint violation in observed forms is forced by a high-ranking constraint. Accordingly, the violated constraint itself can be (conservatively) shifted down in the hierarchy. With each new datum encountered, the learner may be able to shift one or more constraints somewhat closer to their eventual positions in the hierarchy. In the *recursive* mode of the algorithm, the learner considers input data one by one, computing their effects on the constraint ranking immediately. Step by step, the intermediate grammar approaches the full target ranking until it finally 'converges' at a steady endpoint. In the *batch* mode, the algorithm takes in all input data in a single sweep, and processes them in a decreasing order of strictness, so that higher ranked constraints are first placed together in a stratum, before the lower ranked constraints are placed.

A simple example of the recursive mode in action clarifies the basic idea of extracting useful information from input data. Assume a learner confronted

with the task of learning the distribution of voiced and unvoiced obstruents in Dutch. The generalisation that holds here is that coda obstruents are devoiced, while elsewhere (that is, in onsets), voicing is contrastive. In Dutch, the VOICED OBSTRUENT PROHIBITION (VOP) is dominated by IDENT-IO(voice) because voicing is contrastive (in onsets). The target ranking is given below:

(33) NOVOICEDCODA » IDENT-IO(voice) » VOP

The following tableau shows the relevant constraint interactions in a single form:

(34) Coda devoicing

Input: /bed/	NOVOICEDCODA	IDENT-IO(voice)	VOP
(a) bɛd	*!		**
(b) b^{h} ɛt		*	*
(c) pɛd	*!	*	*
(d) pɛt		**!	

Note that in the optimal candidate (34b), the coda is devoiced, while the onset consonant is faithful to its input voicing, which shows the activity of input-to-output faithfulness.

Let us abstract away from alternations, and assume that the learner already knows the underlying representation /bed/. (The learning of alternations and underlying representations is discussed by Tesar and Smolensky 2000, and Hayes, chapter 5 this volume.) Under these somewhat simplified conditions, the learning task amounts to inferring the constraint ranking (33) on the basis of forms encountered in the input, such as [bɛt].

The learning process starts from an initial ranking in which all three constraints cluster together in a single stratum. The assumption of a one-stratum initial state will be reconsidered in the next section, in favour of an initial state in which markedness constraints outrank faithfulness constraints, as proposed in Smolensky (1996a, 1996b).

The learner encounters the first datum: [bɛt], which (as we assumed earlier) (s)he knows is based on the underlying representation /bed/. Since the observed output form [bɛt] must be *optimal* for the given input, any other candidates for the same input can be safely assumed to be less harmonic, that is, *sub-optimal*. She starts by arranging the information to be processed by the constraint ranker in the form of *mark-data pairs*, consisting of pairwise comparisons of the optimal candidate (the winner) and a sub-optimal candidate (a loser):

(35) Mark-data pairs

	<i>sub-opt</i> < <i>opt</i>	<i>loser-marks</i>	<i>winner-marks</i>
b < a	[bed] < [bet]	{*NOVOICEDCODA, *VOP* VOP}	{*VOP,*IDENT-IO(voice)}
c < a	[ped] < [bet]	{*NOVOICEDCODA, *VOP,*IDENT-IO(voice)}	{*VOP,*IDENT-IO(voice)}
d < a	[pet] < [bet]	{*IDENT-IO(voice), *IDENT-IO(voice)}	{*VOP,*IDENT-IO(voice)}

Shared constraint violations (between a winner and a loser) are cancelled out. Consequently, the second mark-data pair ceases to be informative, as [ped] contains a superset of violations of [bet], which cannot (by definition) give any information about ranking. New mark-data pairs are drawn up, including (all and only) relevant and non-redundant information:

(36) Mark-data pairs after marks cancellation

	<i>sub-opt</i> < <i>opt</i>	<i>loser-marks</i>	<i>winner-marks</i>
b < a	[bed] < [bet]	{*NOVOICEDCODA,*VOP}	{*IDENT-IO(voice)}
d < a	[pet] < [bet]	{*IDENT-IO(voice)}	{*VOP}

This display can be (informally) interpreted as follows. We begin by interpreting the second harmonic ranking, [pet] < [bet], which is the simplest case, indicating that VOP is violated in the winner [bet]. The only possible ground for violation could be the need to avoid a violation of another, higher ranked constraint, in this case IDENT-IO(voice), which would be violated if the winner had been the loser (and conversely, the loser had been the winner). Accordingly, it must be the case that more priority is given to avoiding violation of IDENT-IO(voice) than to avoidance of violation of VOP. This amounts to a ranking IDENT-IO(voice) » VOP. Similarly, the first harmonic ranking, [bed] < [bet], shows that IDENT-IO(voice) is violated in the winner [bet]. In contrast to the earlier case, there are two (rather than one) constraints in loser-marks, which means that either NOVOICEDCODA or VOP could have been the cause of violation of IDENT-IO(voice) in the winner. To state it differently, it follows that either of the following rankings holds: NOVOICEDCODA » IDENT-IO(voice) or VOP » IDENT-IO(voice).

We find that some information in mark-data pairs can be straightforwardly translated into a sub-ranking, while other information is ambiguous, posing a 'demotion dilemma' that requires a general resolution strategy. The CDA, by applying a procedure which we need not discuss in detail, cautiously extracts reliable information from mark-data pairs and uses this to demote constraints one by one, until the target ranking is reached.

For example, if the learner starts with the second mark-data pair in (36), working from a single-stratum initial state, (s)he can safely decide to demote VOP below IDENT-IO(voice), placing it into a new stratum.

(37) {NO VOICED CODA, IDENT-IO(voice)} » VOP

Continuing with the first mark-data pair, the learner is faced with a dilemma: she can demote IDENT-IO(voice) below NO VOICED CODA, or alternatively, demote it below VOP:

(38) (a) NO VOICED CODA » {VOP, IDENT-IO(voice)}
 (b) NO VOICED CODA » VOP » IDENT-IO(voice)

The correct demotion (38a) is more conservative than the one in (38b), since it does not rank VOP with respect to IDENT-IO(voice), while (38b) establishes an (incorrect) ranking between these constraints. Conservative demotion is promoted to a general strategy in the algorithm, serving the central goal of terminating the recursive demotion process. It is of use whenever the algorithm faces a dilemma about how deep to demote a constraint, that is, whenever two or more constraints occur in loser-marks (as in 36 b < a). It then demotes any constraint which assigns a *winner-mark* to a stratum immediately below the *highest-ranking* constraint which assigns a *loser-mark*. The conservative demotion strategy avoids placing a demoted constraint too deep down the hierarchy, from which it could never escape except by a re-ranking of other constraints, which may produce an eternal re-ranking process, which fails to terminate.

In addition to presenting a constraint-ranking algorithm, Tesar and Smolensky (2000) address two closely related aspects of the learning process: (i) the learning of covert structural descriptions from input to the learner (including their prosodic structure); and (ii) the learning of underlying representations and morphophonological alternations. Tesar and Smolensky's theory of grammatical acquisition, as a whole, closely reflects the architecture of OT itself: maximal emphasis is put on principles such as harmonic ordering and strict domination.

Boersma (1998, 2000) and Tesar and Smolensky (1998) observed that the CDA is vulnerable to variation in the input data, a common feature of natural language. This problem is caused by the major assumption underlying the CDA that all input data are consistent with a single ranking. The CDA, when fed variable input, responds by failure to converge. For example, the recursive mode of the algorithm, when presented with inconsistent input data, never reaches a stable state, eternally going back and forth between rankings. As a solution to this problem, Boersma (1998, 2000) and Boersma and Hayes (2001) propose an alternative algorithm, the Gradual Learning Algorithm (GLA), which is designed to cope with variation in its input, and which constructs grammars reflecting input variation by variable outputs.

Before addressing the topic of learning OT grammars on the basis of noisy input data, let us address the issue of how to model variable outputs in OT.¹³ The insight of Anttila (1997) is that variation can be modelled by leaving some conflicting constraints unranked. In his model, each time the grammar is deployed, it chooses an ordering of the unranked constraints at random. The result is that a single input may variably map to different outputs.¹⁴ Boersma (1998) preserves the view of variation as variable rankings, but enriches the model by making two novel assumptions. First, the grammar ranks each constraint as a point along a *continuous scale*, with a numerical value allowing an exact measure of a constraint's distance from other constraints. When the grammar is deployed in evaluation, however, the constraints are placed in a strict domination order. Second, the order is determined by the constraints' value on the scale, together with a factor of noise, which provides an interval surrounding the mean of a *stochastic distribution*. The closer two constraints are, the more their distributions will overlap, resulting in a larger proportion of rankings in which the lower ranking constraint of the two dominates the higher ranking one.

Boersma's theory of variation provides the basis of the GLA (Boersma 1998, Boersma and Hayes 2001). This learns a grammar on the basis of variable input data, while coping with free variation and noisy inputs. The idea underlying the GLA is that the grammar changes gradually (rather than categorically) under the influence of input data. Each form fed into the learner has a small effect on the grammar, which only grows into a sizeable ranking effect when significantly amplified by many similar forms. In such cases, ranking values of responsible constraints will end up being so far apart that their distributions hardly overlap, a situation corresponding to a categorical ranking. Variable inputs, however, have effects on the grammar that go in opposite directions; in the end, the grammar reflects the distributions of variable input data by placing constraints close enough on the scale as to make their 'noisy' distributions overlap, as discussed earlier. In this way, variable input data are fruitfully processed by the learner, rather than being rejected due to 'inconsistencies'.

After this overview of the architecture of OT and the learnability of OT grammars, we now turn to the relevance of OT for phonological acquisition.

3. Drawing connections between Optimality Theory and child phonology

In section 1, we discussed formal and substantive connections that have been drawn between child phonology and rule-based and parametric theories of phonology. Research on phonological acquisition in Optimality Theory has also brought to light both formal and substantive links between children's sound systems and cross-linguistic phonology. In this section we shall take up each of

these topics in turn, focusing in particular on the contributions made by papers appearing in this volume.

3.1 *Formal connections: constraint interaction*

Optimality Theory shares with other constraint-based theories the virtue of being able to express formally the conspiratorial behaviour of phonological processes. While this sets Optimality Theory apart from a purely rule-based model, generative research on child phonology has incorporated constraints in one way or another for some time now (see, e.g., the discussion of parametric theories in section 1). Thus, the more germane comparison is between Optimality Theory and other theories that make use of constraints. In this context, the main innovation of Optimality Theory lies in the notion that constraints are minimally violable, that a constraint can be violated if and only if the satisfaction of a higher ranked constraint is at issue. This can be contrasted with the usually implicit view that an observed violation of a constraint implies its inactivity: that within a given domain or level, a constraint is either strictly inviolable or completely without force. The advantage of minimally violable constraints is that they allow for a straightforward account of non-uniform constraint application (Prince 1993), in which a constraint is only satisfied, or violated, under particular circumstances.

One type of non-uniform constraint application was discussed in section 2.2: 'the emergence of the unmarked' (McCarthy and Prince 1994). This refers to a situation in which a markedness constraint is generally violated in the language as a whole, but does have effects in a certain context. The examples McCarthy and Prince (1994) discuss are ones in which a reduplicative morpheme is subject to the effects of a markedness constraint that is violated elsewhere in the language. The other type of non-uniformity might be termed 'the emergence of the marked', in which a markedness constraint is generally satisfied, but is violated in a particular context. Both of these are readily captured through constraint ranking. A constraint can be violated only under compulsion of a higher ranked constraint. If the demands of that higher ranked constraint conflict with those of the lower ranked constraint in most, but not all, environments, then we have the emergence of the unmarked. If the demands of the higher ranked constraint only force the violation of the lower ranked constraint in a limited set of contexts, an emergence of the marked situation is produced. Neither of these situations is expected if constraints are inviolable, though they can of course be dealt with, usually by complicating the statement of the constraints in undesirable ways. Those seeking a concrete example might wish to compare the statement of *NONFINALITY* in Prince and Smolensky's (1993) analysis of Kelkar's Hindi with the *Extrametricality Rule* posited for the same data in Hayes (1995).

Starting with Prince and Smolensky (1993), phonological research has turned up a number of cases of non-uniform constraint application, and provided compelling analyses in terms of ranked constraints. Evidence of non-uniform constraint application in child language, and the accompanying analysis in terms of ranked constraints, yields an important formal parallel between phonological theory and child phonology, and a strong argument for an Optimality theoretic approach to the latter domain.

Two such cases are presented in Amalia Gnanadesikan's chapter in this volume, 'Markedness and faithfulness constraints in child phonology'. Here we shall present the simpler of the two so as to provide an explicit example of the role of ranked constraints in child phonology; we take some liberties with Gnanadesikan's analysis for the sake of expository ease. The evidence for non-uniformity comes from the activity of constraints against high sonority onsets in the forms produced by an English-learning child. In cluster reduction, we find that the sonority of the segments determines which consonant is deleted, with the higher sonority segment being lost. For example, a stop-liquid cluster will lose the liquid, rather than the stop (e.g., [piz] *please*). This can be attributed to a constraint against liquid onsets (for our purposes *L-ONS). Outside of cluster reduction, however, approximant onsets freely occur (e.g., [læb] *lab*). This would immediately raise a paradox in a theory of inviolable constraints: how could a constraint that is active in cluster reduction be violated elsewhere? In OT, the answer would be that a higher ranked constraint usually forces *L-ONS to be violated, but that this constraint does not interfere with the satisfaction of *L-ONS in cluster reduction.

One such constraint is MAX-IO, the faithfulness constraint that blocks segmental deletion by requiring every input segment to have an output correspondent (McCarthy and Prince 1999). The tableau in (39) illustrates the effect of ranking this constraint above *L-ONS when the input onset is a singleton:

(39) Marked structure forced by constraint domination

input: /læb/	MAX-IO	*L-ONS
æb	*!	
☞ læb		*

As this tableau shows, the dominance of MAX-IO generally rules out deletion as a means to satisfy *L-ONS. Cluster reduction is different because in this context, MAX-IO must be violated, due to the dominance of a constraint against clusters, *COMPLEX. Regardless of which consonant is deleted, MAX-IO will be violated, so the decision is passed down to the lower ranked *L-ONS constraint:

(40) Emergence of the Unmarked

input: /pliz/	*COMPLEX	MAX-IO	*L-ONS
pliz	*!		
iz		**!	
liz		*	*!
^ɸ piz		*	

Interestingly, Fikkert (1994: 59) documents a stage in the acquisition of Dutch in which only stops are produced, which are the least marked onsets in terms of onset sonority constraints such as those posited by Gnanadesikan:

- (41) Until 1;9.9 adult target onsets other than plosives are either realised with an initial plosive . . . or deleted by Jarmo

At this presumably prior developmental stage, constraints against sonorous onsets are fully satisfied. In terms of Optimality Theory, this would be captured by having the onset sonority constraints dominate all conflicting constraints. In terms of inviolable constraints, one could invoke a constraint, or set of constraints, against non-plosive onsets. However, such inviolable constraint(s) would be of no formal use in describing the stage discussed by Gnanadesikan, in which the constraint applies just in case one member of an onset is deleted. This developmental progression, in which a constraint is at first fully satisfied, then minimally violated, is straightforwardly expressed by Optimality Theory, and forms a second formal bridge between it and the study of phonological acquisition. For further discussion, see Barlow (1997), Pater (1997), and Barlow and Gierut (1999).

Phonological development as constraint re-ranking is discussed from a somewhat different angle by Clara Levelt and Ruben van de Vijver in their chapter 'Syllable types in cross-linguistic and developmental grammars'. Levelt and van de Vijver identify a set of markedness constraints that govern the basic structure of onsets and codas. They show that the factorial typology produced by the interaction of these constraints with a general faithfulness constraint matches the systems attested cross-linguistically. A ranking with all of the markedness constraints above the faithfulness constraint yields the least marked system, with only CV syllables. A hierarchy with faithfulness dominating all of the markedness constraints allows the full range of syllable types, that is, any expansion of the (C)(C)V(C)(C) template. For a child learning Dutch, these grammars correspond to the hypothesised initial state and the final state, respectively. Other languages have faithfulness dominating some subset of the

markedness constraints, and these form possible intermediate grammars for the Dutch-learning child. Levelt and van de Vijver show that only a small number of such possible intermediate stages are attested in a corpus of data on the acquisition of Dutch. They argue that the limited range of developmental paths can be explained by considering the frequency of the various syllable types in caretaker language. The transition between developmental stages involves demotion of a markedness constraint beneath faithfulness. When there is a choice of markedness constraints to be demoted, it is made on the basis of the frequency of the syllable type that will be added to the child's repertory.

3.2 *Substantive connections: the content of constraints*

Substantive connections between cross-linguistic and developmental phonology are captured in Optimality Theory by having the same constraints apply in both domains, in a manner similar to Jakobson's Laws of Irreversible Solidarity, or Stampe's processes. Given the (standard) assumption that constraints are universal – and in fact innate – the relation between acquisition and cross-linguistic patterns ('typology') could not be accidental. Both the adult speaker's and the child's early grammar are built from the same material: UG's constraint inventory; where adult and early grammars differ, the locus of difference must be in the arrangement of the material (the constraint rankings), not in their substance. One may thus expect the typological variation between the phonologies of natural languages to be mirrored, in a general fashion, in the acquisitional variation between early and adult grammars. But not only does one expect the range of variation found in early and adult grammars to be similar; it is also predicted that for each phonological 'process' found in a child's early phonology, there is a counterpart in the phonology of some natural language.

Research in child phonology has continued to find reflections of typologically attested patterns in child sound systems, and the papers both by Gnanadesikan and Levelt and van de Vijver point out such correspondences. The pattern of sonority-based onset selection that Gnanadesikan documents in English child phonology is closely paralleled in Sanskrit reduplication, while Levelt and van de Vijver show that typologically derived syllable structure constraints characterise the stages of development that Dutch children go through. Levelt and van de Vijver do find one stage that lacks a typological correlate, but the subtlety of this restriction (against the co-occurrence of onset and coda clusters) may well explain its absence from extant linguistic descriptions.

The standard interpretation of these connections is that constraints are innate, and that acquisition consists only of constraint re-ranking, and not of constraint construction. This implicit assumption of most Optimality theoretic acquisition and learnability research is made explicit in Gnanadesikan's chapter, as well as in that of Goad and Rose (see also Kager 1999 for discussion). It is not a

necessary assumption, however; one might also claim that constraints emerge in acquisition in response to articulatory and perceptual pressures (see, e.g., Bernhardt and Stemberger 1998, Boersma 1998, Hayes 1999). The universality of such phonetic and cognitive factors would then be held to explain the observed activity of similar constraints across languages, and across developing grammars. It is difficult to tease these innatist and emergentist accounts apart empirically in terms of their predictions about child language. One source of evidence in favour of an (at least partially) emergentist stance may be the occurrence of phenomena in child speech that are unattested typologically. The prototypical case is that of long-distance assimilation of primary place features between non-adjacent consonants, usually referred to as consonant harmony. Given that the pattern is unattested typologically, it would seem unlikely that it is produced by typologically derived constraints. However, Optimality theoretic analyses of consonant harmony have diverged on this issue; while Pater (1997) takes this as evidence for a child-specific articulatorily based constraint, Goad (1997) constructs an analysis using Alignment constraints that are claimed to be active cross-linguistically (see also Levelt 1995, Dinnsen, Barlow, and Morrisette 1997, Bernhardt and Stemberger 1998, and Rose 2000 on consonant harmony in Optimality Theory).

Innatist and emergentist theories appear to make different predictions about the nature of constraints. Under the emergentist view, the constraints should rather directly mirror the articulatory and perceptual factors on which they are based, while an innatist theory would expect that at least some constraints (if not all) should be purely formal in nature, with no direct phonetic motivation (see Boersma 1998 for discussion). The innatist perspective is defended in Heather Goad and Yvan Rose's chapter, 'Input elaboration, head faithfulness, and evidence for representation in the acquisition of left-edge clusters in West Germanic'. They take the sonority-based analysis of onset reduction put forth by Gnanadesikan and others as representative of a relatively phonetically based approach to the phenomenon. They point to another pattern of onset reduction attested in child speech that they term the 'head pattern', which differs from the sonority pattern in that [s]-initial clusters always lose the [s], even when the second member of the cluster is higher in sonority. They argue that an account of this pattern requires a structurally elaborated syllable structure that encodes the difference between all other obstruent-initial clusters and [s]-initial clusters: the former are left-headed branching onsets, while in the latter, [s] is analysed as an adjunct, and head status is trivially assigned to the consonant following it. Faithfulness to the head position favours the preservation of the second member of [s]-initial clusters, but the initial member of all other obstruent initial clusters. Since headedness is in this context a purely formal, rather than functional, principle, Goad and Rose take this to suggest that the substantive content of constraints, and the representations they refer to, is not just functional.

4. **Markedness » Faithfulness: implications and extensions**

An overarching theme of much Optimality theoretic acquisition and learnability research, and one that connects many of the chapters in this volume, is the relative ranking of markedness and faithfulness constraints at the outset of, and through the course of, acquisition.

4.1 *Data from child production and perception*

To express the unmarkedness of child phonology observed by Jakobson and subsequent researchers, it is often maintained that acquisition starts with markedness constraints ranked above faithfulness constraints. This ranking results in output structures that conform to the demands of the markedness constraints, and are thus simple. As faithfulness constraints come to dominate markedness constraints, structures gradually increase in complexity.

Several of the chapters in this volume present data from child production that provide evidence of structures that are unmarked relative to the adult target language, and put forth analyses in which markedness constraints outrank faithfulness constraints, rankings that are reversed in the target language. The chapters by Gnanadesikan, Goad and Rose, and Levelt and van de Vijver all focus on syllable structure, and, in particular, the reductions in complexity of children's syllable structure relative to that of adults (see also Barlow 1997, Ohala 1996, 1999). Other research has examined similar reductions in complexity at higher levels of prosody, evidenced in particular by truncation (Demuth 1995, 1996, 1997, 2000, Pater and Paradis 1996, Kehoe and Stoel-Gammon 1997, Kehoe 1999, Pater 1997, Curtin 2001, 2002, Ota 1999). There has also been some debate on whether early child productions are in fact correctly characterised by a Markedness » Faithfulness ranking; see Bernhardt and Stemberger (1998) and Velleman and Vihman (2000).

Two of the chapters in this volume discuss evidence for Markedness » Faithfulness ranking in experiments on infant speech perception. In the contribution by Lisa Davidson, Peter Jusczyk, and Paul Smolensky, 'The initial and final states: theoretical implications and experimental explorations of Richness of the Base', an experimental paradigm is introduced that is aimed to assess directly the predictions of this initial ranking. Infants from 4.5 to 20 months of age were presented triples of syllables of the form 'A B AB' in which 'AB' was either a faithful concatenation of A and B, or one in which a markedness-reducing sound change had occurred. Under the hypothesis that infants prefer stimuli which conform to their grammar, and that infants interpret the triples as analogous to $/A + B/ \rightarrow [AB]$, the prediction of Markedness » Faithfulness is that sound-change stimuli should be preferred over faithful but marked stimuli.

This was confirmed by the Headturn Preference Procedure for children at 4.5, 10, and 20 months, although no significant difference was found at 15 months.

Joe Pater's chapter 'Bridging the gap between receptive and productive development with minimally violable constraints' draws on previously published studies of infant speech perception to make the case that receptive acquisition follows a course similar to that of productive development: initial stages permit only unmarked structures; more complex structures emerge later. To account for these parallels, Pater develops a model in which markedness constraints apply in perception as well as in production. Since receptive development does typically precede the development of production, it is necessary to allow for differences in the complexity of structures permitted at a single time. Pater's proposal is that faithfulness constraints can be indexed to perception or production, thus allowing for a situation in which perceptual representations are of greater complexity than those created for production.

4.2 *The initial state: an argument from learnability*

Learnability considerations provide an argument for the ranking of markedness constraints above faithfulness constraints at the outset of acquisition. This argument was first made in Smolensky (1996a), who attributes the basic insight to Alan Prince (personal communication); it is further elaborated on in the chapter by Davidson, Jusczyk, and Smolensky in this volume, as well as in the contributions by Hayes, and Prince and Tesar, which will be discussed in section 4.3. Before proceeding with an outline of this argument, we should note that, while widely accepted, it is not universally so. In the somewhat different approach to learnability advocated by Hale and Reiss (1998), an initial state with Faithfulness » Markedness is in fact held to be necessary.

Suppose that a language lacks a particular structure, such as syllable codas. An account of this gap in terms of Optimality Theory requires that a markedness constraint militating against that structure dominate a faithfulness constraint that would prefer its preservation; for the restriction against codas, NOCODA must dominate a faithfulness constraint like DEP ('no epenthesis'). As discussed in section 2.2, this ranking would need to hold even in a language without overt alternations, since in contrast with earlier generative theories of phonology, Optimality Theory countenances no restrictions on the form of inputs, under the principle of Richness of the Base (Prince and Smolensky 1993).

The learnability issue arises in just this case of languages that do not have alternations, since they provide no positive evidence of the need for the markedness constraints to outrank the faithfulness constraint. A learner with DEP » NOCODA would correctly parse the codaless strings of the ambient language, since no constraint would prefer adding a coda. Assuming that learning is

error-driven (Tesar and Smolensky 1998), this ranking would be a trap, and the learner would not be guaranteed to converge on the correct $\text{NoCODA} \gg \text{DEP}$ hierarchy. This can be seen as an instance of the Subset problem discussed in Principles and Parameters theories (e.g., Berwick 1985); the language produced by $M \gg F$ (e.g., $\text{NoCODA} \gg \text{DEP}$; with only V rimes) is a subset of the language produced by $F \gg M$ (e.g., $\text{DEP} \gg \text{NoCODA}$; V and VC rimes). Since all of the data of the subset language are consistent with the superset language, positive evidence alone will not move the learner out of the superset state.

The solution to this problem suggested by researchers such as Demuth (1995), Gnanadesikan (1995, this volume), Levelt (1995), and Smolensky (1996a, 1996b) is to posit an initial state in which all of the markedness constraints outrank the faithfulness constraints. Positive evidence will be available for any Faithfulness \gg Markedness rankings that are inconsistent with this initial state.

4.3 *Persistence of Markedness \gg Faithfulness: learnability issues*

Bruce Hayes's chapter, 'Phonological acquisition in Optimality Theory: the early stages', and Alan Prince and Bruce Tesar's 'Learning phonotactic distributions' both independently argue that an initial ranking does not go far enough, and claim that a bias for low-ranked faithfulness constraints must persist past the initial state, and be incorporated into the learning algorithm itself. While the fundamental insight of the chapters is a shared one, their implementations and extensions of it are quite different, which makes their inclusion in a single volume particularly opportune.

Prince and Tesar introduce an explicit measure of the degree to which a hierarchy possesses $M \gg F$ structure, and investigate the consequences of trying to maximise this measure by low placement of F in suitably biased versions of the Recursive Constraint Demotion Algorithm (Tesar 1995, Tesar and Smolensky 1998). The key issue is deciding which F to demote when there is more than one F constraint to choose from. They suggest that the main desideratum is the 'freeing up' of further M constraints for ranking, though they also show that such decisions have further consequences downstream for the resultant hierarchy that may motivate a certain kind of 'look ahead' in the decision-making process. Prince and Tesar also consider issues arising in the context of constraints that are in a 'special/general' relationship (see Prince and Smolensky 1993 on Panini's Theorem). They suggest that this consideration yields learning-theoretic motivation for resolving the Positional Markedness versus Positional Faithfulness controversy (Beckman 1998, Zoll 1998) and for deeper scrutiny of faithfulness theory as a whole.

Bruce Hayes also develops a version of Tesar and Smolensky's (1998) Constraint Demotion Algorithm that incorporates a bias towards low-ranked faithfulness. He illustrates his algorithm's effectiveness by having it learn the

phonotactic pattern of a simplified language modelled on Korean. Based on literature from infant speech perception, Hayes suggests that infants accomplish much of this phonotactic learning in the first year of life. The later learning of morphological alternations is guided by an additional default ranking: in contrast to input-output faithfulness, output-output faithfulness constraints start out in a dominant position in the hierarchy (see also McCarthy 1999b). Hayes finds empirical evidence from production data showing that children do in fact sometimes assume a higher rank of output-output faithfulness than is necessary in the language.

4.4 *Persistence of Markedness » Faithfulness: production data*

Whether the bias towards low-ranked Faithfulness is inherited from an initial ranking or explicitly maintained throughout learning, the Markedness » Faithfulness ranking is predicted to persist through subsequent developmental stages, and into the adult grammar, when there is no evidence to force a re-ranking. Evidence of this ranking cannot be obtained through simple inspection of the forms of a language. For example, to show that a speaker of a language without codas in fact encodes a restriction against codas in the phonological grammar, one cannot simply point to the fact that the language lacks syllable-final consonants. Instead, to show the productivity of such a restriction, one might invoke data from loanword adaptations or from the production of nonce words with codas. The Markedness » Faithfulness schema also makes predictions beyond the productivity of static restrictions, since it may be that a language provides no evidence at all for the ranking of some constraints. In our language without codas, inspection of overt forms would reveal nothing about the ranking of a markedness constraint against voiced coda obstruents (NOVOICEDCODA) relative to the faithfulness constraint demanding preservation of underlying voice IDENT (voice). Data bearing on just this situation are presented by Broselow *et al.* (1998). They show that when Mandarin learners of English start to acquire codas, they do go through a stage in which they devoice the codas, as the M » F ranking of NOVOICEDCODA » IDENT (voice) would predict (see section 2.2). In Natural Phonology, innate processes are similarly held to persevere into the mature system when they are not contradicted in the language being learned. In support of this position, Stampe (1969) and Nathan (1984) point to several other cases of the emergence of innate processes in second language phonology that are parallel to Broselow *et al.*'s L2 English example. In this volume, instances of persistent Markedness » Faithfulness ranking are provided in data from child productions, loanword adaptation, and second language acquisition.

In chapter 11, 'Child word stress competence: an experimental approach', Wim Zonneveld and Dominique Nouveau report on an experiment conducted

to establish whether Dutch 3- and 4-year-olds have mastery of the Dutch stress system. Based on elicited production data of real and nonsense words, they find that the answer to this question seems to be an affirmative one, and also that a developmental pattern can be detected from one age group to the next. The chapter goes on to show, however, that some subtle patterns in these experimentally collected data are not captured by any existing standard analysis, whether formulated in rules, parameters, or constraints. It is indicated that, by hindsight, these patterns occur in the adult system too, and an analysis is proposed with two properties: irregularity is treated with the aid of deviating hierarchies; and in some of these hierarchies constraints become visible only because they are active in subpatterns first discovered in the child language experiment. In particular, Zonneveld and Nouveau uncover evidence in their experiment for an undominated *CLASH constraint, whose activity is generally masked by other constraints.

Shigeko Shinohara's chapter, 'Emergence of Universal Grammar in foreign word adaptations', presents a study of the adaptation of French loanwords by Japanese speakers. Shinohara discusses patterns of segmental change and insertion as well as accent placement. Some of the phenomena point to the activity of constraints that govern the phonology of the language as a whole, but others, such as avoidance of stressed epenthetic vowels, and stem-syllable alignment, implicate constraints that are uniquely active in the loanword phonology. Since these constraints do have considerable cross-linguistic justification, Shinohara takes them to be part of the constraint set supplied by Universal Grammar. For the most part, their activity in loanword adaptation can be understood as the emergence of latent M » F rankings, but some ranking among the markedness constraints is also required. These rankings of markedness constraints, Shinohara suggests, are potentially universal, and derivable from phonetic scales (see Prince and Smolensky 1993: ch. 5).

Davidson, Jusczyk, and Smolensky present an experimental paradigm that is designed to induce speakers to subject non-native inputs to their English grammars, to quantitatively assess the M » F-based prediction that non-English clusters would be repaired to meet the requirements of English syllable structure. In the condition that best approximated this prediction, non-English clusters divided into several groups which could be ordered according to their probability of repair. They show that appropriate interaction of constraints independently needed to bar non-English clusters can account for the relative markedness of different non-English clusters, suggesting a final English ranking that makes such distinctions, without apparent motivation in the English data in which none of these clusters appears. They go on to suggest possible analyses of how such a ranking might arise, and point out the importance of such 'hidden rankings' in the final state for pursuing the hypothesis that the initial state for second language acquisition is the final state for first language acquisition.

5. Conclusion

In this introduction, we have emphasised the theme of formal and substantive connections between phonological theory and child phonology. Across theories, a basic formal connection is made by positing mappings between underlying and surface representations in child phonology, analogous (but not necessarily equivalent) to such mappings in phonological theory. In rule-based theories this connection is strengthened by arguing that the rules that perform the mappings in both domains are similar in terms of their formal makeup, as well as in how they interact with one another, specifically, through ordering (Stampe 1969, 1973a, 1973b, Smith 1973). Similarly in OT, constraints are argued to interact through ranking in child language as well as in mature grammars. Substantive connections between child and adult phonology were made in rule-based phonology by positing a set of processes that apply in both domains (Stampe 1969, 1973a, 1973b), and in constraint-based phonology, by having constraints that apply to child and adult grammars. Many of the same issues that confronted earlier attempts to connect child phonology and phonological theory continue to apply today; this is particularly obvious from a reading of Lise Menn's contribution to this volume, 'Saving the baby: making sure that old data survive new theories'. At the same time, however, we should not underestimate the progress that has been made on several fronts. Along with the continued discovery of basic formal and substantive parallels between child and adult phonology, recent research has succeeded in providing explicit proposals about difficult issues such as the learnability of phonology, variation, similarities and differences in comprehension and production, and the genesis of constraints. The rapid progress that is being made, combined with the wide range of issues that remain to be explored, makes the intersection between phonological theory and phonological acquisition such an exciting area for ongoing research.

NOTES

1. Those wishing to expand their knowledge of its subject-matter beyond what is discussed in this section, and/or to familiarise themselves with the view of others on similar material and issues, may wish to consult a number of other texts: out of some handfults we recommend Ingram (1989), Ferguson, Menn, and Stoel-Gammon (1992), Fletcher and MacWhinney (1995), Vihman (1996), Jusczyk (1997), Bernhardt and Stemberger (1998), and Tesar and Smolensky (1998, 2000).
2. First published in German in Uppsala, Sweden, when Jakobson was in Norway as a World War Two refugee; this version was reprinted in 1962 in *Selected Writings I*. An English translation, from which we quote here, appeared in 1968.
3. The expected value of the opposition is often called the 'unmarked' one, the unexpected value the 'marked' one (in this example: stop is unmarked, fricative is marked).

4. Kaye (1974) argues that some cases of opaque rule interaction may be interpreted as functionally motivated in that they contribute to the recoverability (in a technical, non-learning sense) of underlying representations, namely if the opaque derivation produces a segment that occurs nowhere else in the language, as when [ŋ] is the unique product of assimilation > velar deletion. Notice, however, that this is not a case of counterfeeding. For further discussion, see McCarthy (1999a: section 3.1).
5. Smith does not literally claim that underlying forms of the child grammar will by definition be *always identical* to the adult form: empirical evidence may suggest otherwise. Consider the rule of Consonant Harmony, whereby a coronal consonant becomes velar under the influence of a velar later in the word; [gɛgi:] for *taxi*, [gɔ:k] for *talk*. The rule also neutralised *take* and *cake*, for instance. When the rule disappeared, in 'hundred or more examples' (p. 144) the completely regular coronal appeared, as expected; as a single exception, the verb *take* remained [gɛik], later [k^heik]. What apparently had happened was that Amahl assumed the underlying form of this verb to be /keik/, only turning to a different assumption, leading to the correct output, in the face of consistent positive evidence.

Recently, and more fundamentally, Macken (1995) distinguishes between three acquisitional rule types (related to perception, articulation, and generalisation) that each have their own functional and developmental characteristics. Her model allows for underlying representations which are the same as surface representations in cases of perception-based neutralisation, which is typically eliminated slowly and word by word. Interestingly, Macken (1980) eliminated Smith's velarisation rule (= 5a) in this latter manner.

6. See, e.g., Anderson (1985: 342 ff.) for a discussion and an assessment; some of NGP's leading ideas have resurfaced in a much different form in OT.
7. For general introductions to OT, we refer to Archangeli and Langendoen (1997), Kager (1999), and McCarthy (2001).
8. The duplication problem was recognised as early as *SPE* (Chomsky and Halle 1968: 382): 'Thus certain regularities are observed within lexical items as well as across boundaries – the rule governing voicing in obstruent sequences in Russian, for example – and to avoid duplication of such rules in the grammar it is necessary to regard them not as redundancy rules but as phonological rules that also happen to apply internally to a lexical item.' This solution to the duplication problem, to order morpheme structure rules among the other phonological rules, was shown to be insufficiently general by Kenstowicz and Kisseberth (1979), for the reason that it cannot handle cases in which a phonological rule is blocked from applying because its output would violate a static condition on the lexicon. After discussing examples from Russian and Tonkawa, Kenstowicz and Kisseberth (1979: 433) conclude: 'In both cases a constraint on UR affects the application of a phonological rule – by adjusting the output of the rule in one case and by preventing application of the rule in the other. An ordering solution works for the former but not the latter. Thus, the ordering solution cannot be accepted as a totally general solution to the duplication problem. It would seem that once a way is found to express the conspirational relation between MSRs and the application of phonological rules in examples such as the Tonkawa one, the duplication involved in examples such as Russian should fall out as a special subcase.' This sketches in essence the approach taken in OT: surface phonological constraints account for generalisations on lexical items, and also

function to trigger and block phonological changes with respect to the input – that is, alternations.

9. There is one case which is not covered by a possible word test based on identity mapping. Chain shifts (see Kirchner 1996 for an OT analysis) are mappings in which an input /A/ is mapped onto [B], while input /B/ maps onto [C]. If 'B' undergoes the test, it will change, and hence fail the test; however, [B] is also the legitimate output of the mapping /A/ ⇒ [B].
10. The fact that Dutch phonology has an alternative way of repairing the input /kd/ (by regressive voicing assimilation) into [gd] is beside the point: both mappings involve a ranking indicated in (25).
11. In the OT literature, loanword adaptations have been argued to bear on various issues, such as the universality of markedness constraints, as opposed to the language-specific nature of rewrite rules. References include Yip (1993), Itô and Mester (1995), Paradis (1996), Paradis and LaCharité (1997), Gussenhoven and Jacobs (2000), LaCharité and Paradis (2000), and the contributions in this volume by Shinohara, and Davidson, Jusczyk, and Smolensky.
12. For another approach to the issue of the learnability of OT grammars, see Pulleyblank and Turkel (1998, 2000).
13. See also Demuth (1997), Curtin (2002), Curtin and Zuraw (forthcoming), and Pater and Werle (2001) on the use of these models to deal with variation in acquisition.
14. See Reynolds (1994) for a slightly different view of variation in OT.

References

- Anderson, S. R. (1985). *Phonology in the Twentieth Century. Theories of Rules and Theories of Representations*. The University of Chicago Press.
- Anttila, Arto (1997). *Variation in Finnish Phonology and Morphology*. Ph.D. dissertation, Stanford University.
- Archangeli, D. (1997). Optimality Theory: an introduction to linguistics in the 1990s. In D. Archangeli and D. T. Langendoen (eds.) *Optimality Theory: an Overview*. Malden, Mass.: Blackwell Publ. 1–32.
- Archangeli, D. and D. T. Langendoen (eds.) (1997). *Optimality Theory: an Overview*. Oxford: Basil Blackwell.
- Archangeli, D. and D. Pulleyblank (1994). *Grounded Phonology*. Cambridge, Mass.: MIT Press.
- Barlow, J. (1997). *A Constraint-Based Account of Syllable Onsets: Evidence from Developing Systems*. Ph.D. dissertation, Indiana University.
- Barlow, J. and J. Gierut (1999). Optimality Theory in phonological acquisition. *Journal of Speech, Language, and Hearing Research* 42. 1482–1498.
- Beckman, J. (1998). *Positional Faithfulness*. Ph.D. dissertation, University of Massachusetts.
- Benua, L. (1997). *Transderivational Identity: Phonological Relations Between Words*. Ph.D. dissertation, University of Massachusetts. [ROA-259] <http://roa.rutgers.edu>
- Bernhardt, B. H. and J. P. Stemberger (1998). *Handbook of Phonological Development*. San Diego: Academic Press.

- Berwick, R. (1985). *The Acquisition of Syntactic Knowledge*. Cambridge, Mass.: MIT Press.
- Boersma, P. (1997). How we learn variation, optionality, and probability. *Proceedings of the Institute of Phonetic Sciences* 21. Amsterdam. 43–58. [Also ROA 221, <http://roa.rutgers.edu>]
- (1998). *Functional Phonology: Formalizing the Interactions between Articulatory and Perceptual Drives*. Ph.D. dissertation, University of Amsterdam, The Hague: Holland Academic Graphics.
- (2000). Learning a grammar in Functional Phonology. In J. Dekkers, F. van der Leeuw and J. van de Weijer (eds.) *Optimality Theory: Syntax, Phonology, and Acquisition*. Oxford University Press. 465–523.
- Boersma, P. and B. Hayes (2001). Empirical tests of the Gradual Learning Algorithm. *Linguistic Inquiry* 32. 45–86.
- Braine, M. D. S. (1976). Review of Smith (1973). *Lg* 52. 489–498.
- Broselow, E., S.-I. Chen, and C. Wang (1998). The Emergence of the Unmarked in second language phonology. *Studies in Second Language Acquisition* 20. 261–280.
- Burzio, L. (1996). Surface constraints versus Underlying Representation. In J. Durand and B. Laks (eds.) *Current Trends in Phonology: Models and Methods*. CNRS, Paris X, and University of Salford. University of Salford Publications.
- Chomsky, N. (1959). A review of B. F. Skinner's *Verbal Behavior*. *Lg* 35. 26–58.
- (1965). *Aspects of the Theory of Syntax*. Cambridge, Mass.: MIT Press.
- (1968). *Language and Mind* (1st edn), page reference to the 1973 enlarged edition. New York: Harcourt, Brace, and Jovanovich.
- (1967). Some general properties of phonological rules. *Lg* 43. 102–128.
- (1981a). *Lectures on Government and Binding*. Dordrecht: Foris.
- (1981b). On the representation of form and function. *The Linguistic Review* 1. 3–40.
- (1986). *Knowledge of Language: Its Nature, Origin and Use*. New York: Praeger.
- Chomsky, N. and M. Halle (1968). *The Sound Pattern of English*. New York: Harper & Row.
- Curtin, S. (2001). Children's early representations: evidence from production and perception. In the *Proceedings of the Holland Institute of Linguistics (HILP 5)*, Potsdam, Germany.
- (2002). *Representational Richness in Phonological Development*. Ph.D. dissertation, University of Southern California.
- Curtin, S. and K. Zuraw (forthcoming). Explaining constraint demotion in a developing system. In the *Proceedings of Boston University Conference on Language Development*, Boston, Mass.
- Demuth, K. (1995). Markedness and the development of prosodic structure. *NELS* 25. 13–25. [Also ROA 50, <http://roa.rutgers.edu>]
- (1996). Alignment, stress, and parsing in early phonological words. In B. Bernhardt, J. Gilbert and D. Ingram (eds.) *Proceedings of the UBC International Conference on Phonological Acquisition*. Somerville, Mass.: Cascadilla Press.
- (1997). Multiple optimal outputs in acquisition. In B. Moren and V. Miglio (eds.) *University of Maryland Working Papers in Linguistics*. College Park: University of Maryland Department of Linguistics.
- (2000). Prosodic constraints on morphological development. In J. Weissenborn and B. Höhle (eds.) *Approaches to Bootstrapping: Phonological, Syntactic and*

Neurophysiological Aspects of Early Language Acquisition, Vol. 2. Amsterdam: John Benjamins.

- Dinnsen, D. A., J. A. Barlow, and M. L. Morrisette (1997). Long-distance place assimilation with an interacting error pattern in phonological acquisition. *Clinical Linguistics and Phonetics* 11. 319–338.
- Donegan, P. J. and D. Stampe (1979). The study of Natural Phonology. In D. A. Dinnsen (ed.) *Current Approaches to Phonological Theory*. Bloomington: Indiana University Press. 126–173.
- Dresher, B. E. (1999). Charting the learning path: cues to parameter setting. *LI* 29. 27–67.
- Dresher, B. E. and J. Kaye (1990). A computational learning model for metrical phonology. *Cognition* 34. 137–195.
- Edwards, M. L. and L. Shriberg (1983). *Phonology: Applications in Communication Disorders*. San Diego, Calif.: College-Hill Press.
- Ferguson, C. A., L. Menn, and C. Stoel-Gammon (eds.) (1992). *Phonological Development: Models, Research, Implications*. Parkton, Md.: York Press.
- Fikkert, P. (1994). *On the Acquisition of Prosodic Structure*. The Hague: HAG.
- Fletcher, P. and B. MacWhinney (eds.) (1995). *The Handbook of Child Language*. Cambridge, Mass.: Blackwell.
- Goad, H. (1997). Consonant harmony in child language: an optimality-theoretic account. In S. J. Hannahs and M. Young-Scholten (eds.) *Focus on Phonological Acquisition*. Amsterdam: John Benjamins. 113–142.
- Hale, M. and C. Reiss (1998). Formal and empirical arguments concerning phonological acquisition. *LI* 31. 157–169.
- Halle, M. (1983). On the origins of the distinctive features. In M. Halle (ed.) *Roman Jakobson: What He Taught Us. International Journal of Slavic Linguistics and Poetics, Vol. XXVII: Supplement*. Columbus, OH: Slavica Publ. 77–86.
- Halle, M. and J.-R. Vergnaud (1987). *An Essay on Stress*. Cambridge, Mass.: MIT Press.
- Hayes, B. (1980/81). *A Metrical Theory of Stress Rules*. Ph.D. dissertation, MIT, Cambridge, Mass. Revised version distributed by Indiana University Linguistics Club, Bloomington, Ind.
- (1995). *Metrical Stress Theory: Principles and Case Studies*. Chicago: University of Chicago Press.
- (1999). Phonetically-driven phonology: the role of Optimality Theory and inductive grounding. In M. Darnell, E. Moravcsik, M. Noonan, F. Newmeyer and K. Wheatly (eds.) *Functionalism and Formalism in Linguistics, Vol. I: General Papers*. Amsterdam: John Benjamins. 243–285.
- Hyman, L. (1970). The role of borrowing in the justification of phonological grammars. *Studies in African Linguistics* 1. 1–48.
- Ingram, D. (1974). Phonological rules in young children. *Journal of Child Language* 1. 49–64.
- (1989). *First Language Acquisition. Method, Description, and Explanation*. Cambridge: Cambridge University Press.
- Itô, J. and R. A. Mester (1995). The core-periphery structure of the lexicon and constraints on reranking. In J. Beckman, L. Walsh Dickey and S. Urbanczyk (eds.) *Papers in Optimality Theory*. [University of Massachusetts Occasional Papers in

- Linguistics 18.] Amherst, Mass.: Graduate Linguistic Student Association. 181–210.
- Jacobs, H. and C. Gussenhoven (2000). Loan phonology: perception, salience, the lexicon and Optimality Theory. In J. Dekkers, F. van der Leeuw and J. van de Weijer (eds.) *Optimality Theory: Phonology, Syntax, and Acquisition*. Oxford University Press. 193–210.
- Jakobson, R. (1941/1968). *Kindersprache, Aphasie und allgemeine Lautgesetze*. Uppsala Universitets Aarskrift. Translated by R. Keiler as *Child Language, Aphasia and Phonological Universals*. The Hague: Mouton.
- Juszyk, P. W. (1997). *The Discovery of Spoken Language*. Cambridge, Mass.: MIT Press.
- Kager, R. (1999). *Optimality Theory*. Cambridge: Cambridge University Press.
- Kaye, J. (1974). Opacity and recoverability in phonology. *Canadian Journal of Linguistics* 19. 134–149.
- (1989). *Phonology: a Cognitive View*. Hillsdale, New Jersey: Lawrence Erlbaum Assoc.
- Kean, M. L. (1975). *The Theory of Markedness in Generative Grammar*. Ph.D. dissertation, MIT, Cambridge, Mass.
- Kehoe, M. (1999). Truncation without shape constraints: the latter stages of phonological acquisition. *Language Acquisition* 8:1. 23–67.
- Kehoe, M. and C. Stoel-Gammon (1997). The acquisition of prosodic structure: an investigation of current accounts of children's prosodic development. *Lg* 73. 113–144.
- Kenstowicz, M. (1996). Base-identity and uniform exponence: alternatives to cyclicity. In J. Durand and B. Laks (eds.) *Current Trends in Phonology: Models and Methods*. CNRS, Paris X, and University of Salford. University of Salford Publications. 363–393. [ROA-103]
- Kenstowicz, M. and C. Kisseberth (1979). *Generative Phonology: Description and Theory*. New York: Academic Press.
- Kiparsky, P. (1968). Linguistic universals and linguistic change. In E. Bach and R. T. Harms (eds.) *Universals in Linguistic Theory*. New York: Holt, Rinehart and Winston. 171–204.
- (1999). Paradigm effects and opacity. MS., Stanford University.
- Kiparsky, P. and L. Menn (1977). On the acquisition of phonology. In J. Macnamara (ed.) *Language Learning and Thought*. New York: Academic Press.
- Kirchner, R. (1996). Synchronic chain shifts in Optimality Theory. *LI* 27. 341–350.
- Kisseberth, C. W. (1970). On the functional unity of phonological rules. *LI* 1. 291–306.
- LaCharité, D. and C. Paradis (2000). Derivational residue: hidden rules in Optimality Theory. In J. Dekkers, F. van der Leeuw and J. van de Weijer (eds.) *Optimality Theory: Syntax, Phonology, and Acquisition*. Oxford: Oxford University Press. 211–233.
- Levelt, C. (1995). Unfaithful kids: Place of Articulation patterns in early vocabularies. Colloquium presented at the University of Maryland.
- Lombardi, L. (1999). Positional faithfulness and voicing assimilation in Optimality Theory. *NLLT* 17. 267–302.
- McCarthy, J. J. (1999a). Sympathy and phonological opacity. *Phonology* 16. 331–399.
- (1999b). Morpheme Structure Constraints and Paradigm Occultation. MS., University of Massachusetts at Amherst. *CLS* 32.

- (2001). *A Thematic Guide to Optimality Theory*. Cambridge: Cambridge University Press.
- McCarthy, J. and A. Prince (1993). Prosodic Morphology I: constraint interaction and satisfaction. MS., University of Massachusetts, Amherst and Rutgers University. [To appear, Cambridge, Mass.: MIT Press. Technical report 3, Rutgers University Center for Cognitive Science.]
- (1994). The emergence of the unmarked: Optimality in prosodic morphology. In M. González (ed.) *NELS* 24. 333–379.
- (1995). Faithfulness and reduplicative identity. In J. Beckman, L. Walsh Dickey and S. Urbanczyk (eds.) *Papers in Optimality Theory*. [University of Massachusetts Occasional Papers in Linguistics 18.] Amherst, Mass.: Graduate Linguistic Student Association. 249–384.
- (1999). Faithfulness and identity in prosodic morphology. In R. Kager, H. van der Hulst and W. Zonneveld (eds.) *The Prosody–Morphology Interface*. Cambridge: Cambridge University Press. 218–309.
- Macken, M. A. (1980). The child's lexical representation. *JL* 16. 1–17.
- (1995). Phonological acquisition. In J. A. Goldsmith (ed.) *The Handbook of Phonological Theory*. Cambridge, Mass.: Blackwell. 671–696.
- Menn, L. (1971). Phonotactic rules in beginning speech. *Lingua* 26. 225–241.
- (1978). Phonological units in beginning speech. In A. Bell and J. B. Hooper (eds.) *Syllables and Segments*. Amsterdam: North-Holland. 157–171.
- (1980). Phonological theory and child phonology. In G. H. Yeni-Komshian, J. F. Kavanagh, and C. A. Ferguson (eds.) *Child Phonology. Vol. 1: Production*. New York: Academic Press. 23–42.
- Menn, L. and E. Matthei (1992). The 'Two-Lexicon' account of child phonology: looking back, looking ahead. In C. A. Ferguson, L. Menn and C. Stoel-Gammon (eds.) *Phonological Development: Models, Research, Implications*. Timonium, Md.: York Press. 211–248.
- Nathan, G. S. (1984). Natural Phonology and interference in second language acquisition. In G. S. Nathan and M. E. Winters (eds.) *The Uses of Phonology: Proceedings of the First Conference on the Uses of Phonology* (Southern Illinois University Occasional Papers in Linguistics 12). Carbondale: Department of Linguistics, Southern Illinois University. 103–113.
- Ohala, D. (1996). *Cluster Reduction and Constraints in Acquisition*. Ph.D. dissertation, University of Arizona.
- (1999). The influence of sonority on children's cluster reductions. *Journal of Communication Disorders* 32. 397–422.
- Ota, M. (1999). *Phonological Theory and the Acquisition of Prosodic Structure: Evidence from Child Japanese*. Ph.D. dissertation, Georgetown University.
- Paradis, C. (1996). The inadequacy of filters and faithfulness in loanword adaptation. In J. Durand and B. Laks (eds.) *Current Trends in Phonology: Models and Methods*. University of Salford Publications. 509–534.
- Paradis, C. and D. LaCharité (1997). Preservation and minimality in loanword adaptation. *JL* 33. 379–430.
- Pater, J. (1997). Minimal violation and phonological development. *Language Acquisition* 6. 201–253.
- Pater, J. and J. Paradis (1996). Truncation without templates in child phonology. In A. Stringfellow, D. Cahana-Amitay, E. Hughes and A. Zukowski (eds.) *Proceedings*

- of the 20th Annual Boston University Conference on Language Development*. Somerville, Mass.: Cascadilla Press.
- Pater, J. and A. Werle (2001). Typology and variation in child consonant harmony. In C. Féry, A. Dubach Green and Ruben van de Vijver (eds.) *Proceedings of HILP5*. University of Potsdam. 119–139.
- Piggott, G. L. (1988). The parameters of nasalization. MS., McGill University.
- Prince, A. S. (1983). Relating to the Grid. *LI* 14. 19–100.
- Prince, A. and P. Smolensky (1993). Optimality Theory: constraint interaction in generative grammar. MS., Rutgers University, New Brunswick and University of Colorado, Boulder. [Technical report 2, Rutgers University Center for Cognitive Science. To appear, Cambridge, Mass.: MIT Press.]
- Pulleyblank, D. and W. J. Turkel (1998). The logical problem of language acquisition in Optimality Theory. In P. Barbosa, D. Fox, P. Hagstrom, M. McGinnis and D. Pesetsky (eds.) *Is the Best Good Enough? Optimality and Competition in Syntax*. Cambridge, Mass.: MIT Press and MITWPL. 399–420.
- (2000). Learning phonology: genetic algorithms and Yoruba tongue root harmony. In J. Dekkers, F. van der Leeuw and J. van de Weijer (eds.) *Optimality Theory: Phonology, Syntax and Acquisition*. Oxford: Oxford University Press. 554–591.
- Reynolds, W. T. (1994). *Variation and Phonological Theory*. Ph.D. dissertation, University of Pennsylvania.
- Rose, Y. (2000). *Headedness and Prosodic Licensing in the L1 Acquisition of Phonology*. Ph.D. dissertation, McGill University.
- Smith, N. V. (1973). *The Acquisition of Phonology: a Case Study*. Cambridge: Cambridge University Press.
- Smolensky, P. (1996a). *The Initial State and 'Richness of the Base' in Optimality Theory*. Technical Report JHU-CogSci-96-4, Cognitive Science Department, Johns Hopkins University. [ROA 154, <http://ruccs.rutgers.edu/roa.html>]
- (1996b). On the comprehension/production dilemma in child language. *LI* 27. 720–731. [ROA 118, <http://ruccs.rutgers.edu/roa.html>]
- Stampe, D. (1969). The acquisition of phonetic representation. In R. I. Binnick *et al.* (eds.) *CLS* 5. 443–454.
- (1973a). *A Dissertation on Natural Phonology*. Ph.D. dissertation, University of Chicago.
- (1973b). On chapter nine. In M. J. Kenstowicz and C. W. Kisseberth (eds.) *Issues in Phonological Theory*. The Hague: Mouton. 44–52.
- Stemberger, J. P. (1993). Rule ordering in child phonology. In M. Eid and G. Iverson (eds.) *Principles and Prediction: the Analysis of Natural Language*. Amsterdam: Benjamins. 305–326.
- Stonham, John (1990). *Current Issues in Morphological Theory*. Ph.D. dissertation, Stanford University.
- Tesar, Bruce (1995). Computational Optimality Theory. Ph.D. dissertation, University of Colorado at Boulder.
- Tesar, Bruce and Paul Smolensky (1993). The Learnability of Optimality Theory: an Algorithm and Some Basic Complexity Results. Technical Report cu-cs-678-93, Department of Computer Science, University of Colorado at Boulder. Rod-2.
- (1998). Learnability in Optimality Theory. *LI* 29. 229–268.
- (2000). *Learnability in Optimality Theory*. Cambridge, Mass.: MIT Press.

- Velleman, S. L. and M. M. Vihman (2000). The optimal initial state. Paper read at the Annual Meeting of the Linguistic Society of America, Chicago, January.
- Vihman, M. M. (1996). *Phonological Development: the Origins of Language in the Child*. Cambridge, Mass.: Blackwell.
- Yip, M. (1993). Cantonese loanword phonology and optimality theory. *Journal of East Asian Linguistics* 2, 262–291.
- Zoll, C. (1998). Positional asymmetries and licensing. MS., Cambridge, Mass.: MIT.