

Optimality and diachronic adaptation*

Martin Haspelmath

Max-Planck-Institut für evolutionäre Anthropologie, Inselstr. 22,
04103 Leipzig, haspelmath@eva.mpg.de

Abstract

In this programmatic paper, I argue that the universal constraints of Optimality Theory (OT) and the functional explanations of functionalists need to be complemented by a theory of diachronic adaptation. OT constraints are traditionally stipulated as part of Universal Grammar, but this misses the generalization that the grammatical constraints normally correspond to constraints on language use. As in biology, observed adaptive patterns in language can be explained through diachronic evolutionary processes, as the unintended cumulative outcome of numerous individual intentional actions. The theory of diachronic adaptation also provides a solution to the teleology problem, which has often been used as an argument against functional explanations. Finally, I argue against the view that the grammatical constraints could be due to accident, and I conclude that an explanatory theory of grammatical structure needs a theory of adaptation.

1. Preferences in competition: an old and new concept

There is a long tradition in theoretical linguistics which holds that structural patterns of grammar are determined by highly general preferences or constraints that may come into conflict with each other. Gabelentz (1901:256) was very clear about the tension between the “striving for ease” (*Bequemlichkeitsstreben*) and the “striving for clarity” (*Deutlichkeitsstreben*). The neogrammarians were primarily concerned with the conflict between phonological tendencies (leading to exceptionless sound changes) and the tendency toward morphological

* I received useful comments on an earlier version of this paper from Susanne Michaelis, Bill Croft, Joan Bybee, Rudi Keller, Thomas Müller-Bardey, Esa Itkonen, John Hawkins, an anonymous *Zeitschrift für Sprachwissenschaft* referee, and from audiences at the Konstanz meeting of the Deutsche Gesellschaft für Sprachwissenschaft (February 1999) and at the Max Planck Institute for Evolutionary Anthropology. I am grateful to everybody who in this way helped me to improve this paper.

analogy. Havers (1931 : 191 ff) discusses in great detail the interaction of various general “conditions and forces” in syntax.

With the advent of structuralism and its rigid synchrony/diachrony separation, this kind of thinking went out of fashion, as the focus was now on explicit and elegant descriptions of individual languages, rather than on highly general (if often vague) explanatory principles. But after several decades of abstention, linguists again began to become interested in highly general principles; and since principles can be formulated in a more general way if they are violable, this meant that the idea of conflicting preferences resurfaced. Within one tradition, such competing preferences were called *naturalness principles in conflict* (e. g. Dressler 1977 : 13, Dressler et al. 1987 : 7, 93); in another, *competing motivations* (Haiman 1983 : 812, Du Bois 1985, Croft 1990 : § 7.4). Langacker (1977 : 102) used the term *optimality*:

“I believe we can isolate a number of broad categories of linguistic optimality. Languages will tend to change so as to maximize optimality in each of these categories... The tendencies toward these various types of optimality will often conflict with one another.”

In this line of thinking, it has always been assumed that the competing preferences are not just highly general, but in fact constitute universal properties, or design features of human language. Prince & Smolensky's (1993) formal framework of Optimality Theory (OT) uses this idea also in synchronic descriptions of grammatical structures (originally, in phonology, but later also in syntax) by introducing the notion of language-specific preference ranking (or “constraint ranking”, in their terminology).¹ Much work in Optimality Theory has shown that the availability of violable constraints often yields more elegant and appealing descriptions than accounts in terms of inviolable rules.

This general framework has been enormously successful and popular, perhaps not only because it allows linguists to formulate much more general principles than were hitherto possible, but also because many of the constraints are intuitively plausible, and because the description in terms of the “best” and “worse” candidates often corresponds to our pretheoretical feelings. Consider, as a simple example, the distribution of the plural allomorphs /-z/, /-əz/, and /-s/ in English. This may be accounted for in the OT framework by postulating the four constraints SAMEVOICE (“Sequences of obstruents within a syllable must agree for voicing”), OCP(SIBILANT) (“Sequences of sibilants are prohibited within the word”), DEPIO (“Insertion of segments is prohibited”), and

1 Ranking of naturalness principles has been widely discussed in Natural Morphology, but mainly in terms of universal ranking (e. g. Wheeler 1993) and type-specific ranking (e. g. Dressler 1985 a). Language-specific differences were attributed to other factors by these authors.

Table 1

	SAME VOICE	OCP(SIB)	DEPIO	IDENT (VOICE)
input: kæt-z				
kæt-z	*!			
kæt-əz			*!	*
☞ kæt-s				
input: buʃ-z				
buʃ-z	*!			
☞ buʃ-əz			*	
buʃ-s		*!		
input: stoun-z				
☞ stoun-z				
stoun-əz			*!	*!
stoun-s				

IDENTITY(VOICE) (“Input and output are identical for voicing”) (cf. Gussenhoven & Jacobs 1998: 48–49). Assuming an underlying /-z/ for the plural -s, we get the constraint tableau in 1, where the three relevant cases *cat-s* [kæt-s], *bush-es* [buʃ-əz], and *stone-s* [stoun-z] are shown.

Only *stone-s* [stoun-z] shows no constraint violation at all. In *bush-es* [buʃ-əz], DEPIO (the constraint against epenthesis) is violated, but SAMEVOICE and OCP(SIBILANT) are ranked higher, so [buʃ-əz] is the optimal candidate. In *cat-s* [kæt-s], IDENTITY(VOICE) is violated, but again, the two competing candidates [kæt-z] and [kæt-əz] violate higher-ranked constraints. In informal OT parlance, [kæt-s] is “better” than [kæt-z] and [kæt-əz], and this quasi-technical terminology coincides nicely with our feeling that indeed [kæt-s] “sounds better” than its competitors in that it is easier to pronounce.

However, this intuitive coincidence between “good” in the sense of “optimal with respect to OT constraints” and “good” in the sense of “good for the language user” has not been captured in mainstream versions of OT. I will argue in this paper that by capturing this correspondence between **grammatical optimality** and **user optimality**, we are able to reach a significantly higher level of explanatory adequacy.

2. Why are the constraints the way they are?

The OT framework does two things very well. On the one hand, it allows descriptions of language-specific facts that are more principled than those in the previous frameworks of generative linguistics. For instance, the constraint set in Tableau 1 is more general than phonological rules like “z → [-voice]/[-voice]__” and “z → əz/[+strident, +coronal]__”, from whose formulation it is not immediately clear that they are by no means arbitrary.

On the other hand, the OT framework allows an elegant statement of typological options, called **factorial typology**. The typology of possible languages is given by the set of possible rankings of the constraints. Consider a simple example, again from phonology (Prince & Smolensky 1993: §6.1): The two widely applicable syllable structure constraints ONSET (“A syllable must have an onset”) and NoCODA (“A syllable must not have a coda”), together with the constraint FAITHFULNESS (“The output must not contain fewer or more segments than the input”) allow three types of languages, depending on their mutual ranking ($X \succcurlyeq Y$ means ‘X is ranked higher than Y’):

(1)	ONSET \succcurlyeq FAITH	FAITH \succcurlyeq ONSET
NoCODA \succcurlyeq FAITH	CV (e. g. Hua)	(C)V (e. g. Cayuvava)
FAITH \succcurlyeq NoCODA	CV(C) (e. g. Tharrkari)	(C)V(C) (e. g. Mokilese)

However, this cannot be the whole story yet. We must ask further: Why are there no constraints such as CODA or NoONSET, which are opposite to NoCODA and ONSET? Nothing in standard OT prohibits these constraints, so if it is true (as it seems to be) that they do not exist, this can only be achieved by stipulation. Such an account may be satisfactory for linguists who limit their goal to an elegant description of particular languages. But the theoretically minded linguist will be more ambitious and ask a further why question: **Why are the constraints the way they are?**

It could of course turn out that this question is unanswerable, and that the constraints are not more than accidents of history. The usual assumption is that the OT constraints are innate, and it might be that they arose as an accidental side-effect of some adaptive modification of the brain (cf. §7 below for further discussion of this possibility). But there seems to be a widespread feeling among OT practitioners that this is not the whole story. Otherwise there would be no need to justify new constraints with reference to non-distributional evidence. But this is what one commonly finds. For instance, Bresnan (1997) postulates a constraint PROAGR (“Pronominals have the referentially classificatory properties of agreement”) and states: “The functional motivation for the present constraint could be that pronouns ... bear classificatory features to aid in

reference tracking, which would reduce the search space of possibilities introduced by completely unrestricted variable reference.” Similarly, Casali (1997: 500) justifies his constraint **MAXLEX**, which he uses to express the fact that vowel elision in hiatus contexts typically does not affect roots and content words, by noting that it “arises from a more general functional motivation, a preference for maintaining phonological material belonging to elements that typically encode greater semantic content”. And Morelli (1998: 7) introduces the constraint ***STOP-OBSTRUENT** (“A tautosyllabic sequence containing a stop followed by any obstruent is disallowed”) and states: “It is justified both phonetically and phonologically. Phonetically, it reflects the preference for stops to be released into more sonorous segments ...”

Strictly speaking, such justifications are irrelevant in a theory that assumes innate constraints. But the fact that they are mentioned by OT practitioners indicates that they have the intuition that the constraints are not arbitrary and are in principle susceptible of (or even in need of) further explanation. However, to my knowledge nobody has so far made an attempt to explain OT constraints in a systematic fashion. In the next two sections, we will see what such an explanation might look like.

3. User optimality and adaptation

What the justifications of constraints by Bresnan, Casali, and Morelli have in common is that they portray the constraints as being good for speakers and hearers in one way or another, i. e. as exhibiting **user optimality** (to use the term introduced in §1). To see this more clearly, in (2) I reformulate some of the constraints that we have seen so far in terms of the language users’ needs:

(2) User optimality of grammatical constraints

name	grammatical constraint	corresponding user constraint
MAXLEX (Casali 1997: 501)	"Every input segment in a lexical word or morpheme must have a corresponding segment in the output."	Preserving phonological material of elements with greater semantic content helps the hearer to identify the most important parts of a discourse.
IDENTITY (McCarthy & Prince 1995)	"Input and output are identical."	Input-output identity, i.e. uniformity of morphemes across environments, helps the hearer to identify morphemes.
SAMEVOICE (Gussenhoven & Jacobs 1998: 48)	"Sequences of obstruents within a syllable must agree for voicing."	Obstruent sequences with different phonation types are difficult to pronounce because the phonation is impeded by the obstruent occlusion.

Many further constraints that have been used in the literature, including the literature on OT in syntax, can be reformulated in terms of user optimality as well. Some further examples are shown in (3).

(3) User optimality of further grammatical constraints

name	grammatical constraint	corresponding user constraint
STAY (Grimshaw 1997, Speas 1997)	"Do not move."	Leaving material in canonical positions helps the hearer to identify grammatical relationships and reduces processing costs for the speaker.
TELEGRAPH (Pesetsky 1998)	"Do not pronounce function words."	Leaving out function words reduces pronunciation costs for the speaker in a way that is minimally disruptive for understanding by the hearer.
RECOVERABILITY (Pesetsky 1998)	"A syntactic unit with semantic content must be pronounced unless it has a sufficiently local antecedent."	Omitting a meaning-bearing element in pronunciation makes the hearer's task of extracting the intended meaning from the speech signal very difficult unless it can be inferred from the context.

There is probably no need to go into the details of what exactly makes language structures “good” for speakers and hearers, i. e. what constitutes user optimality. I take it as evident that the best option among a range of alternatives is the one which promises the highest net benefit to speaker and hearer. The most important cost factors are motor costs and cognitive processing costs, and the most important benefits are informativeness and persuasiveness (cf. Keller 1998: 189 ff. for some general discussion).

Not all of the constraints that OT practitioners have been working with can be rephrased in terms of user optimality so easily. Sometimes more discussion is required. For instance, Hawkins (1999) gives compelling arguments for the view that filler-gap dependencies of the island type are difficult to process. There is thus good motivation for rephrasing Pesetsky’s (1998) ISLAND CONDITION constraint in terms of user optimality, although this is not as straightforward as in (2)–(3). In other cases, a proposed OT constraint is so highly specific that it seems unlikely that a direct reformulation in user-optimality terms will ever be possible. Examples are Grimshaw’s (1997: 374) constraint NO LEXICAL HEAD MOVEMENT (“A lexical head cannot move”) and Pesetsky’s (1998) constraint LEFT EDGE(CP) (“The first pronounced word in CP is a function word related to the main verb of that CP”). However, these constraints clearly have the flavor of theoretical constructs that help make the particular analysis work, but that would be the first candidates for elimination if this becomes possible. Sometimes OT analyses also posit language-specific constraints, and these clearly cannot have a counterpart in terms of user optimality: User optimality is necessarily universal. Finally, constraints that are the direct opposites of each other cannot be rephrased as user constraints, because opposite user constraints would cancel each other out and have no effect. But again, it seems that most OT practitioners consider analyses superior that avoid language-specific constraints and operate entirely with highly general, plausibly universal constraints. It is my impression that most of the widely used, non-ephemeral constraints can be reformulated in user-optimality terms in one way or another. I cannot of course demonstrate that this is indeed the case, but in addition to the OT constraints in (2) and (3), I will mention further constraints together with their user-optimality counterpart later in this paper. Readers who are not well-versed in the functionalist literature will in this way get an idea of why I am optimistic in this respect, even if I should not succeed in convincing them.

Thus, there is a generalization here that has not been captured so far: Loosely speaking, what is “good” from the point of view of the theory is good from the point of view of language users. Grammatical optimality and user optimality are largely parallel. The obvious way of accounting for this striking match between grammatical structures and speaker needs is the notion of **adaptation**. Grammatical structures are adapted to the needs of language users (cf. Croft 1990: 252). By making use of the notion of adaptation, we achieve two things. First, we can account for the parallels between grammatical constraints and constraints on

speakers observed in this section. Second, we can answer the question of §2, why the grammatical constraints are the way they are: The grammatical constraints are ultimately based on the constraints on language users.

The concept of adaptation is familiar from evolutionary biology. For instance, consider the fact that various fish species living in the Arctic and Antarctic regions have antifreeze proteins in their blood. These proteins constitute a structural fact about several unrelated species living far apart, which is obviously to the benefit of these fish. It would be completely mysterious without assuming either a benevolent Creator (the almost universally accepted view until the 19th century) or a historical process of adaptation. It was Charles Darwin's insight that a long-term evolutionary process of successive modified replications combined with environmental selection can account not only for the origin of species, but can also explain the highly complex adaptations found in biological organisms. In short: Arctic and Antarctic fish have antifreeze proteins in their blood because at some point antifreeze proteins arose accidentally (by random genetic mutation). This genetic feature spread because it allowed its bearers to enter a previously unoccupied ecological niche.

I argue in this paper that linguistic adaptation is in many ways analogous to biological adaptation.

4. A mechanism for adaptation: diachronic change

Although historical processes are typically associated with the social sciences and the humanities, they are in fact central to evolutionary biology. Evolutionary biology, in turn, is central to theoretical biology, as is expressed in Theodosius Dobzhansky's well-known remark that "nothing in biology makes sense except in the light of evolution". If biologists restricted their attention to purely synchronic phenomena (as they did well into the 19th century), they would understand very little of what they observe.

I will now argue that historical (or, as linguists say, diachronic) processes are of equally central importance for linguistic theory. Just like biological adaptation, linguistic adaptation requires time. We need to consider diachronic change if we want to understand why the OT constraints are the way they are, i.e. in what sense they are based on the user constraints of §3.

Of course, I am not the first to argue that grammatical structures are "based" on "user constraints" (or "performance constraints", or "functional pressures"). There is a long tradition of functionalist thinking in linguistics that attempts to explain properties of language structure with reference to properties of language use (e.g. Jespersen 1894, Horn 1921, Hawkins 1994, Givón 1995). However, the functionalists have generally paid little attention to possible mechanisms for adaptation – they have usually taken adaptation for granted.

Consider as a concrete example Dik's (1997: 30–34) discussion of Berlin & Kay's famous hierarchy of color terms (black/white > red > green/yellow > blue > brown > others), which embodies the claim that if a language has a basic term for a color somewhere on the hierarchy, then it also has terms for all the colors to the left of this color. Dik observes that this hierarchy is also relevant for the frequency with which color terms are used: 'black' and 'white' are the most frequently used color terms, followed by 'red', and so on. And he continues: "This suggests a functional explanation for the existence of hierarchies of this type: the more frequent the need for referring to some colour, the higher the chance that there will be a separate lexical item for indicating that colour." (Dik 1997: 33).

This is an interesting suggestion, but it is not an explanation.² Useful or needed things are not sufficiently explained by their usefulness or the need for them. Again, biology provides the appropriate analogy: Antifreeze proteins are surely useful for polar fish, indeed necessary for their survival, but this does not suffice as an explanation for their presence. Taking functional statements as sufficient explanation can be called the **Teleological Fallacy**, which is just a special case of humans' general tendency to think in anthropomorphic terms. When speaking about human artifacts, functional or teleological statements are unproblematic: "A bicycle saddle is softer than other parts of the bicycle **in order for** cyclers to sit comfortably." This statement suffices as an explanation for the softness of the saddle because it can easily be converted into a purely causal statement: "A bicycle saddle is soft **because** the bicycle makers have made it soft **in order for** cyclers to sit comfortably." This can be considered a full explanation because the purpose clause depends on an action verb, and the purpose can be attributed to goal-oriented human design. Similarly, antifreeze proteins in polar fish can be fully explained with reference to goal-oriented, purposeful divine design, if one has no concept of evolution or rejects this concept ("Polar fish have antifreeze proteins in their blood **because** God created polar fish with antifreeze proteins **in order to** help them survive in freezing water").

For obvious reasons, neither human design nor divine design are available in linguistics to convert functional statements into full explanations. But linguists have often fallen victim to the Teleological Fallacy (if only in their rhetoric), and we often find statements such as those in (4) (emphasis added).

- (4) a. "Case is formed for reasons of ambiguity, because at some point in history speakers must have talked without *cases* (*Cato interficit Caesar*). Then inflection was added, **in order for** the meaning of the sentence to become clear." (Scaliger 1584: Book 4, ch. 77: 169–80, cited after Bрева-Claramonte 1983: 66)

2 I propose an explanation below in §6.6 (vi).

- b. “[C]oding devices tend to be employed strategically [by grammars] so as to guarantee, at minimal formal expense, the distinguishability only of those grammatical relations which would otherwise be too difficult to distinguish by the addressee.” (Plank 1987: 177)
- c. “[S]yntactically relevant morphemes tend to occur at the periphery, in order to be visible for the syntax.” (Booij 1998 b: 21)
- d. “Of was introduced in order to Case-mark a NP/DP which would not otherwise be Case-marked.” (Lightfoot 1999: 121)

Critics of functionalism in linguistics have rightly pointed out that such explanations are not viable. Haider (1998: 98) observes that “the fact that the design is good for a function is not the driving force that led to the design”, and Tooby & Cosmides (1990: 762) remark: “It is magical thinking to believe that the “need” to solve a problem automatically endows one with the equipment to solve it”.

However, the fact that functionalists rarely provide an explicit mechanism for functional adaptation in language structure does not mean that none exists and that functional explanation in adaptationist terms is possible only in biology. I will now argue that linguistic change is sufficiently similar to biological change that we can transfer some key notions of evolutionary biology to linguistics (see also Croft (1996) and (2000), Kirby (1999), Nettle (1999) for evolutionary accounts that are close in spirit to mine). That linguists have largely ignored this possibility may be due to the fact that in the 20th century the prestige of diachronic studies has not been very high. But as in biology, we cannot understand synchronic language structure without taking into account its diachronic evolution.

5. Variation and selection in language

Let us briefly recapitulate how adaptive explanations work in biology. In ordinary colloquial speech, quasi-teleological statements such as (5a) are very common. They are accepted because everybody knows how teleological statements are translated into purely causal statements in Darwinian evolutionary theory (cf. 5 b).

- (5) a. Giraffes have long necks in order to be able to feed on the leaves of high trees.
- b. At some earlier time, there was genetic variation: There were giraffes with somewhat longer necks and giraffes with somewhat shorter necks. Because giraffes with somewhat longer necks had the addition-

al food source of high trees, they had greater reproductive success. **Therefore** the long-neck gene spread throughout the whole population.

I propose that the translation from teleological to causal statements works very similarly in linguistics. The quasi-teleological, functionalist statement in (6 a) is insufficient on its own, but it becomes quite acceptable when we realize that it can be thought of as just an abbreviation of the purely causal statement in (6 b).

- (6) a. In *cat-s* [kæts], the suffix consonant is voiceless **in order to** satisfy the SAMEVOICE constraint. (Or: ... **in order to** facilitate the pronunciation of this obstruent cluster.)
- b. At some earlier time, there was structural variation: The suffix *-s* could be pronounced [z] or [s]. **Because** [kæts] required less production effort than [kætz], speakers chose it increasingly often (in order to save production energy). After some time, the form [kæts] had become very frequent and **therefore** was reanalyzed as obligatory, while [kætz] was no longer acquired and dropped out of the language.³

On the analogy of the biological term “natural selection”, this process can be called “functional selection” (cf. Nettle (1999: 30–35) for this term and some discussion; Kirby’s (1999: 36) equivalent term is “linguistic selection”). The application of the evolutionary scenario in linguistics presupposes three hypotheses: (i) Languages show structural **variation** in all areas of grammar, and language change is unthinkable without structural variation; (ii) frequency of use is determined primarily by the **usefulness** (or “**user optimality**”) of linguistic structures; and (iii) **high-frequency** structures may become **obligatory**, and **low-frequency** items may be **lost** as a result of their (high or low) frequencies. In the remainder of this section, I will briefly motivate these hypotheses (a full justification is of course beyond the scope of this paper).

The insight that there is constant variation in species was one of the key ingredients in Darwin’s evolutionary theory – before Darwin, species had been thought of only in terms of their properties, as immutable eternal essences (Mayr 1982). Only Darwin’s shift to a population-based view of species, which allowed

3 It must be admitted that this example is not ideal because [kæts] and [kætz] cannot have occurred as variants side by side for a very long time. [kætz] is very difficult to pronounce, so it was presumably eliminated very soon. I chose this example because it was mentioned in a different context earlier. (A better example would have been the choice between [stounz] and [stounəz], which presumably occurred side by side for a long time. However, since [stounz] arose by vowel loss from the earlier [stounəz], rather than the latter by epenthesis from an earlier [stounz], the parallel with the OTE analysis of Tableau 1 would not be so clear.)

for variation and historical change, made evolutionary theory possible. That languages show constant variation has been commonplace in linguistics for a long time, and students of diachronic change routinely assume that every change begins with variation, both at the individual and at the social level. Like descriptive anatomists, descriptive grammarians have usually worked with idealized systems, for good reasons. However, one of the consequences of the present approach is that variation is highly relevant for the theoretical grammarian.

The second hypothesis probably does not need any further justification: That speakers use more user-friendly structures more often than less user-friendly ones can easily be derived from unchallenged common-sense knowledge about human nature.

The third hypothesis is perhaps not so obvious, but there is of course ample evidence for the crucial role of frequency of exposure in establishing cognitive patterns, and more specifically grammatical patterns. The establishment of grammatical structures in the mind is called **entrenchment** by Langacker (1987):

“Every use of a structure has a positive impact on its degree of entrenchment, whereas extended periods of disuse have a negative impact. With repeated use, a novel structure becomes progressively entrenched, to the point of becoming a unit; moreover, units are variably entrenched depending on the frequency of their occurrence (*driven*, for example, is more entrenched than *thriven*).” (Langacker 1987: 59)

The psycholinguistic evidence for frequency as a relevant factor for mental representation is of course enormous. What is less clear is how high frequency of use can turn a linguistic variant into the only possible option. Here further research is needed, but in any event some such mechanism must exist (see also Kirby 1999: ch.2 for discussion). Entrenchment due to frequency thus corresponds to selection in biology. Just like the useful genes spread in a species because of the greater reproductive capacities of their bearers, linguistic features may spread in a speech community because of their usefulness (combined with their social value), and they may become obligatory in grammars because of their high degree of entrenchment. Croft (1996) puts it as follows:

“The proper equivalent [of the perpetuation of genes] is that the perpetuation of a particular utterance structure is directly dependent on the survival of the cognitive structures in a grammar that are used by the speaker in producing utterances of that structure. I suggest that the interactive-activation model used by cognitive grammar and by Bybee (1985) provides a mechanism by which cognitive structures can “survive” – become entrenched in the mind – or “become extinct” – decay in their entrenchment.” (Croft 1996: 115–16)

Of course, the correlation between frequency of use and certain linguistic structures has often been noted, e. g. by Jespersen (1894), Horn (1921), Zipf

(1935), Greenberg (1966), Du Bois (1987). However, linguists have generally been vague about the mechanism by which frequency of use influences language structure. Zipf (1935: 29) claimed that “high frequency is the cause of small magnitude”, but he did not explain how frequency shrinks linguistic units. Du Bois (1987) observed that “grammars code best what speakers do most”, but he did not explain how this marvelous fit of form to function comes about. If entrenchment, i. e. the establishment of patterns in speakers’ mental grammars, is frequency-sensitive, we can actually explain such frequency-based generalizations.

Although language change is of course not intended by speakers, linguistic evolution as outlined above has intentional aspects. Speakers speak and listen intentionally, and their choices of specific expressions from a range of options can also be said to be intentional, although these are usually fairly automatic and unconscious (cf. Itkonen’s 1983: 185 ff. concept of “unconscious rationality”). But unlike (6 a), which cannot be literally true (and therefore has to be translated into (6 b)), a statement such as (7) can be taken as literally true.

- (7) Speakers often chose [kæts] rather than [kætz] in order to save production energy.

Processes of language change like the one outlined in (6 b) are thus neither completely intentional processes (clearly languages don’t change because speakers want to change them) nor completely mechanical processes in which human intention plays no role. Keller (1994) has exposed the frequent fallacy of dichotomizing all processes into the disjoint classes of human actions and natural processes. Processes of linguistic change and selection do not fit into either of these two categories: They are the cumulative outcome of a large number of intentional actions of type (7), an outcome that is not among the goals of these actions. A recent example of this fallacy is Haider’s (1998: 97) characterization of the functionalist view as “the hypothesis that grammar might indeed be a human artifact, that is, a tool shaped by the human mind for a highly desirable function, namely effective and efficient communication”. But on the present view, grammar is neither a human artifact nor a biological entity which can be studied in depth without any regard for human actions or choices. It is the unintended product of a complex but reasonably constrained and regular historical process, linguistic evolution.

There is one important difference between biological evolution and linguistic evolution that should be mentioned at this point: While the source of genetic variation in biology is restricted to random mutations, the source of linguistic variation, innovations in the speech of individual speakers, is often non-random. For instance, the introduction of the variant pronunciation [kæts] (*cats*) in addition to the older [kætz] was clearly motivated by the same user constraint that led to the increasing use of this variant and its eventual

obligatoriness. In this sense, the evolution of linguistic structures is in part “Lamarckian”, like the evolution of other conventional mental structures (generically called “memes” by Dawkins 1976). This difference does not mean that linguistic evolution cannot be regarded as an evolutionary process (cf. Keller 1994: §6.1, Croft 1996). In biology, “Lamarckian” evolution does not work because acquired characters are not inherited, but in linguistic evolution, acquired features can evidently be passed on. The main argument I have made here, that synchronically adaptive structures can be understood in terms of a diachronic process of variation and selection, is not affected by this difference in the mechanism of replication.

One might even go so far as to attribute functional adaptation in language exclusively to the functional factors influencing speaker innovations. This is done implicitly by Croft (2000), who draws the analogy between biological and linguistic evolution in the following way: Mutation is analogous to innovation, and selection is analogous to propagation of a change. Croft maintains that linguistic variants are selected/propagated by speakers because of their social value, i. e. the social status and relationships of the people using the selected variant⁴. To use Nettle’s (1999) terms, Croft attributes the propagation of linguistic features exclusively to “social selection” and sees no role for “functional selection”. Even if this turned out to be correct, the main point of this paper would not be affected (Bill Croft, p.c.). Linguistic evolution would then be entirely “Lamarckian”, but such an evolutionary scenario would be equally capable of transforming teleological statements into causal statements. I do not doubt that social selection is extremely important in linguistic diachrony. To a large extent, the fact that languages differ in their structure, i. e. conventionally assign different weights to different constraints, must ultimately be attributed to social selection. In order to avoid the difficult consonant cluster in [kætz], speakers could also have selected other options, e. g. they could have modified the stem consonant (yielding, e. g., [kædz]); that they did not do this must have been due to social selection. Whatever the precise roles of social and functional selection, structural adaptation in language must be due the effect of constraints on performance combined with a mechanism that turns preferred options of language use into structural patterns of grammar.

In the next section, I will make this general approach more concrete by examining a number of proposed (theory) optimality constraints and by showing how they can be understood as resulting ultimately from user optimality.

4 “[I]n general, differences in functional utility do not play a role in the propagation of a variant; only differences in social utility do.” (Croft (in press), manuscript 132-33)

6. Grammatical optimality reduced to user optimality

In the preceding sections, I proposed that the correspondence between grammatical optimality and user optimality can be explained in terms of a theory of diachronic adaptation. This is a very strong claim which can be falsified easily by showing that a particular synchronically adaptive structure could not have arisen through a diachronic process of adaptation as sketched in §5. In this section, I will examine a number of proposed OT constraints and show in each case how they have arisen from the corresponding user constraints. I should perhaps emphasize here that these case studies are not intended as substantive contributions to the respective areas of linguistics, and that I necessarily gloss over many controversies in the brief accounts given here. My only purpose in this section is to illustrate in a concrete fashion how the general program of diachronic adaptation in linguistics might work.

6.1 No Voice Coda

Let us begin with optimality constraints that have been proposed in phonology. German syllable-final devoicing is generally accounted for by invoking a constraint NO VOICE CODA (e. g. Golston 1996, Raffelsiefen 1998).

- (8) NO VOICE CODA
Voiced coda obstruents are forbidden. (Golston 1996: 717)

The diachronic origin of this constraint is fairly clear. Old High German records show no evidence of this constraint: The spelling consistently has voiced obstruents in syllable-final position, e. g. *tag* 'day', genitive *tages* 'day's'. But by the Middle High German period, the spelling typically indicates that the pronunciation was voiceless (*tac*, genitive *tages*). So at some point in the Middle Ages, the devoiced pronunciation must have become an obligatory part of the grammar.

Obligatory devoicing was in all likelihood preceded by a period of variation in which both the voiced and the unvoiced pronunciation of obstruents in coda position was possible (as well as indefinitely many degrees of voicing in between). How this variation came about in the first place is clear: Voiced obstruents are difficult to pronounce in coda position for well-understood phonetic reasons (cf. Keating et al. 1983), so speakers of all languages with voiced coda obstruents have a tendency to devoice these in pronunciation, thus introducing phonetic variation. In German these devoiced pronunciations became prevalent at some point, and speakers came to treat them as part of the conventionalized grammatical pattern.

Thus, the user constraint corresponding to NO VOICE CODA can be formulated as in (9).

- (9) “User-optimal NO VOICE CODA”:
Coda obstruents should be pronounced voiceless in order to avoid articulatory difficulties.

6.2 MAXLEX

The constraint MAXLEX is proposed by Casali (1997) to account for the fact that vowel elision is less likely to affect roots and content words than affixes and function words:

- (10) MAXLEX
“Every input segment in a lexical word or morpheme must have a corresponding segment in the output.” (Casali 1997: 501)

For example, in Etsako (Niger-Congo) vowel elision in hiatus contexts generally affects the first of two adjacent vowels, i. e. the initial vowel of the second word is preserved (e. g. /owa ɔda/ ‘a different house’ → [ow’ ɔda]). But when the second word is a function word, its vowel is elided (e.g. /ɔna aru ɔli/ ‘that louse (lit. the louse that)’ → [ɔn’ aru ‘li]).

While no direct diachronic evidence is available in this case, it is easy to reconstruct how the current distribution must have come about. Originally, the underlying sequence /ɔna aru ɔli/ could be pronounced with all its vowels intact, and at some point speakers began to drop vowels to avoid the hiatus. Initially any vowel could be elided, but speakers more often elided the final vowel to aid word recognition (words are more easily recognized by their initial segments). However, in function words such as /ɔli/ ‘that’, speakers tended to elide the first vowel, because due to their high frequency and predictability, function words can be recognized more easily than content words. Thus, [ɔn’ aru ‘li] was used significantly more often than [ɔn’ ar’ ɔli], and as a result it became fixed (i. e. entrenched) in speakers’ grammars. Thus, speakers make use of the user-optimal counterpart to (10):

- (11) “User-optimal MAXLEX”
Lexical morphemes should be pronounced fully because they are relatively rare and unpredictable, while functional morphemes can be reduced phonetically without a major threat to comprehensibility.

MAXLEX is of course an old insight. Jespersen (1922: 271) observed that “[i]t has often been pointed out ... that stem or root syllables are generally better

preserved than the rest of the word: the reason can only be that they have greater importance for the understanding of the idea as a whole than other syllables". Jespersen was also aware that this match between function and form must somehow lie in language use,⁵ but like most other functionalists of the 19th and 20th centuries, he did not make the causal connection between constraints on language use and constraints on language structure explicit.

6.3 DROPTOPIC

Let us go on to syntactic constraints now. The constraint DROPTOPIC is proposed by Grimshaw & Samek-Lodovici (1998) to account for the fact that subject pronouns are omitted in many languages (e. g. Italian *ha cantato* 'he has sung', not ??*lui ha cantato*) when they convey topical information.

(12) DROPTOPIC

"Leave arguments coreferent with the topic structurally unrealized."
(Grimshaw & Samek-Lodovici 1998)

In non-null-subject languages like English, DROPTOPIC is dominated by the constraint PARSE (or MAXIO), which requires the underlying topical pronoun to be present overtly. Like the constraint MAXLEX of the preceding subsection, DROPTOPIC corresponds to speakers' tendency to use overt material economically. While MAXLEX specifies that lexical (i. e. relatively unpredictable) information should be preserved, DROPTOPIC specifies that topical arguments, i. e. relatively predictable information, should be omitted. A more general statement of this requirement is Pesetsky's (1998) TELEGRAPH: "Do not pronounce function words." If one considers topical personal pronouns to be function words,⁶ then TELEGRAPH subsumes DROPTOPIC.

In this case the diachronic scenario is so well known that I need not say much here: As a general (though not exceptionless) rule, languages with rich subject agreement do not allow a personal pronoun when it conveys topical information (cf. Gilligan 1987). However, the pronoun may be used occasionally for reasons of extravagance or "expressiveness" (cf. Haspelmath (to appear)), thus introducing variation. Now in languages that are losing their rich subject agreement

5 Cf. Jespersen (1922: 271): "In ordinary conversation one may frequently notice how a proper name or technical term, when first introduced, is pronounced with particular care, while no such pains is taken when it recurs afterwards: the stress becomes weaker, the unstressed vowels more indistinct, and this or that consonant may be dropped." Here he refers to first mention vs. later mention of a rare word, but similar considerations apply to rare vs. frequent words.

6 At the very least, personal pronouns are normally omitted in "telegraphic speech", just like other function words.

morphology on the verb (as has happened in English and French, for instance), speakers will increasingly tend to choose the option of using the personal pronoun, because the verbal agreement does not provide the information required for referent identification in a sufficiently robust way. At some point the use of personal pronouns becomes so frequent that it is reanalyzed as obligatory and the frequent performance pattern comes to be reflected in a competence pattern. This scenario gives rise to the English and French situation, in which PARSE dominates DROP TOPIC. Conversely, speakers of older Italian did not use the overt pronoun much because the full subject agreement on the verb made this unnecessary, and as a result the pronounless pattern is (still) obligatory in Italian today. Thus, the user constraint corresponding to DROP TOPIC in (12) is as shown in (13).

(13) “User-optimal DROP TOPIC”:

A topical subject pronoun should be omitted to save production energy when it is relatively predictable, e.g. in a language with rich subject agreement. (It should not be omitted when no robust information from agreement is available.)

6.4 Stay

The constraint STAY was proposed by Grimshaw (1997: 374) and is technically formulated as ECONOMY OF MOVEMENT (“Trace is not allowed.”). For most purposes, this amounts to the same as (14), which is Speas’s formulation.

(14) STAY

“Do not move.” (Speas 1997: 176)

Grimshaw uses this constraint to account for the ungrammaticality of multiple *wh*-questions with multiple *wh*-movement in English (**What will where they put?*). Since I know too little about the diachronic evolution of this particular construction, I will choose as my example another construction where it would seem natural to invoke STAY as well. SVO languages with rigid word order (such as English) typically show NP-PP word order in postverbal position, i.e. a sentence like (15 a) is the only possibility, and (15 b) is ungrammatical.

- (15) a. *I introduced Kostya to Toshio.*
 b. **I introduced to Toshio Kostya.*

If (15 a) shows the underlying order (V-NP-PP), then (15 b) is ruled out because it violates STAY: The direct object NP has moved to the right of the PP (or conversely).

Again, this constraint in English has its roots in earlier diachronic variation. And again, the facts are too well known to need much discussion: Word order in Old English was much less constrained than in modern English, and the equivalent of (15b), with V-PP-NP word order, was unproblematic. But as morphological case was being lost, there was an increasing need to identify syntactic relations of phrases by their surface positions.⁷ What speakers did was to vary word order much less in performance (relying on other means to convey information-structural information), generalizing the most common order V-NP-PP until it became obligatory.⁸ So again, frequent occurrence in speech gives rise to a grammatical pattern. The performance constraint analogous to STAY is formulated in (16).

(16) “User-optimal STAY”:

Syntactic elements should not be linearized in a non-canonical way if that creates potential ambiguity for the hearer.

6.5 ANIMATE INANIMATE

The constraint ANIM(ATE) > INANIM(ATE) is used by Aissen (1997) to account for various animacy effects, for instance the restriction in Tzotzil (a Mayan language of Mexico) that prohibits the active voice when the patient is inanimate (cf. 17a).⁹ In such cases, the passive voice must be used (cf. 17b), because the non-subject must not outrank the subject on the hierarchy “animate > inanimate”.

(17) Tzotzil (Aissen 1997: 725, 727)

a. **I-x-poxta Xun li pox-e.*
 ASPECT-3P.AGENT-cure Juan the medicine-ENCLITIC
 ‘The medicine cured Juan.’

b. *Ipoxta-at ta pox li Xun.*
 cure-PASSIVE by medicine the Juan
 ‘Juan was cured by the medicine.’

7 This is not the only possibility, as Bill Croft has reminded me. It is equally possible that word order became fixed “spontaneously” and that this in turn facilitated the loss of case distinctions. In this case, the functional, user-optimal motivation for the change would be much less obvious (cf. Lehmann 1992 for a proposed answer).

8 Why V-NP-PP rather than, say, V-PP-NP was generalized is a separate issue that is irrelevant here. See Hawkins (1994) for a theory of syntactic processing that explains the preference for NP-PP order over PP-NP order in VO languages.

9 Cf. also Müller (1997a: 315) for the use of a similar animacy-based constraint in a different context.

Constraints having to do with animacy are of course very familiar from the functionalist literature (cf. Comrie 1989: ch.9), so this is a particularly bad candidate for an innate constraint. A much more plausible scenario again begins with frequency in performance. Universally, there is a strong statistical correlation between topicality and animacy: We tend to talk about humans and other animates, and our sentences usually predicate additional information about them. In those languages that have a strong association between topicality and subjecthood, most subjects will therefore be animate, and most inanimates will be non-subjects. In some languages, these skewed frequencies may become categorical distinctions, i. e. the most frequent patterns may become the only possible ones. This is what must have happened in Tzotzil at some point in the past.

Thus, Aissen's competence constraint ANIM > INANIM corresponds to a very general preference of speakers to talk about animates more than about inanimates. The corresponding "user constraint" cannot really be called a constraint in the sense of a restriction put on speakers – it is what people naturally tend to do.

(18) "User-optimal ANIM > Inanim":

An animate referent should be chosen as topic because the hearer is more likely to be interested in getting more information about animates than about inanimates.

6.6 Further cases

It would not be difficult to continue this list of grammatical optimality constraints that can be shown to have arisen as a result of selection from the variation introduced through language change. As I observed earlier, not all constraints that have been used in OT analyses can be reduced to user constraints in a straightforward fashion, but it seems to me that most widely used constraints can be so reduced. This is of course particularly true of the most general constraints whose names evoke a long earlier literature, such as RECOVERABILITY (e. g. Pesetsky 1998), SALIENCE (e. g. Müller 1997 a), SONORITY (e. g. Raffelsiefen 1998), OCP (e. g. Booij 1998 a), ANIM > INANIM (e. g. Aissen 1997), (LEXICAL) INTEGRITY (e. g. Anderson 1997). They are most obviously adaptive, but these are also the constraints for which an innateness assumption is the least plausible. Their use in constraint tableaux is often very convenient, but it is clear that this cannot be the whole story of explanation. In each case we need a diachronic scenario of conventionalization that links the constraints on language use to the observed patterns of grammar.

The same is of course true for classical cases of functional explanations evoking a highly general theoretical construct which is intended to explain an

observed grammatical pattern, but is not really sufficient as an explanation. Examples include the following:

(i) **Iconicity:** Haiman (1983) notes that there is an iconic relationship between the form and the meaning of, for instance, causative constructions: Causatives expressed by more closely bound items tend to express more direct causation. But this correlation becomes an explanation only if it can be shown that speakers are constrained by iconicity in language use and that patterns of use become grammatical patterns.

(ii) **Economy:** Many linguists have stressed the role of economy in explaining grammatical patterns, especially the shortness of frequent expressions, or the omission of redundant expressions (cf. Zipf 1935, Greenberg 1966, Haiman 1983, Werner 1989). But again, pointing out a correlation is not sufficient: We also have to show how frequent use leads to shortness (e.g. by increased diachronic reduction in frequent items).

(iii) **Phonetic efficiency:** Gussenhoven & Jacobs (1998: 34) note the tendency for phonetic inventories to lack a [p] in the series [p, t, k], and a [g] in the series [b, d, g], and they relate this to the relative inefficiency of [p] and [g]. For instance, [g] “is relatively inefficient from the point of view of the speaker, because the relatively small air cavity behind the velar closure causes the air to accumulate below it, thus increasing the supraglottal air pressure and diminishing the glottal airflow, and thereby causing voicing to stop. That is, a [g] is relatively hard to say.” But the authors do not say how it might be explained that languages tend to lack inefficient stop consonants. They merely suggest that “languages somehow monitor the development of their phonologies”, as if it were obvious what the literal translation of this metaphorical way of speaking should be.

(iv) **Compensation:** Nettle (1995) argues that languages with large phonemic inventories have the compensatory advantage of allowing shorter linguistic units. We thus have a tradeoff relation between paradigmatic costs and syntagmatic economy (and vice versa). Nettle notes that this is explained if “language is functionally adapted to the needs of efficient communication” (1995: 359), as if functional adaptation were not an explanandum itself. He also hints that languages should be seen as “dynamical, self-organizing systems” (1995: 365), but of course we need to know how the self-organization works (but see Nettle 1998, 1999, where Nettle does provide the needed background theory).

(v) **Early Immediate Constituents:** Hawkins (1990, 1994) shows that the principle of Early Immediate Constituents makes correct predictions both about the distribution of word order patterns in performance (where word order is mandated by grammatical rules) and about universals of grammatical word order rules. Hawkins vaguely talks about the “grammaticalization” of word order patterns, but he does not elaborate on this. Clearly, what is needed is a theory of how frequent word order choices in performance tend to become fixed in diachronic change (cf. Kirby 1994, 1999).

(vi) **Frequency:** We saw above that Dik (1997) attempts to explain the color term hierarchy with reference to the frequency of color terms. This correlation between frequency and cross-linguistic occurrence can be turned into an adaptive explanation in the following way: First, basic color terms that a language already possesses are the less likely to be lost from the lexicon the more often they are used by speakers, because high frequency of use leads to a high degree of entrenchment. Second, of the colors for which a language does not have basic terms, those that are the most frequently referred to by non-basic terms will be the most likely to acquire basic terms, for instance by change of a polymorphemic non-basic color term to a basic color term. The fusion of a polymorphemic word to a monomorphemic word is facilitated by high frequency of use. Thus, because of speakers' tendencies in language use, we obtain the universal hierarchy of basic color terms.

7. Are grammatical constraints due to accident?

I have argued so far that the grammatical constraints employed in Optimality Theory are the way they are because they arise from universal constraints on language use through a diachronic adaptive process. But of course, it is theoretically possible that the recurring correspondence between grammatical constraints and user constraints is "a mere coincidence, a serendipitous outcome that speakers may exploit" (to use Durie's 1995: 278 phrase). This would be an astonishing coincidence indeed (and I doubt that anybody would seriously defend such a view), but it is nevertheless possible. In fact, recently a number of linguists have tended to emphasize the dysfunctional aspects of language structure (e. g. Chomsky 1991: 448, Haider 1998, Uriagereka 1998, Lightfoot 1999: Ch. 9), and the view that OT constraints or all of UG are accidental properties of the human mind is more than just a straw man. "UG may have evolved as an accidental side-effect of some other adaptive mutation" (Lightfoot 1999: 249; cf. also Haider 1998: 106). Persuasive evidence for this view would be a widely attested OT constraint that is dysfunctional, but proponents of this view have so far only presented far less convincing cases.

Lightfoot (1999) mentions the example of the constraint that traces must be governed lexically, which prohibits complementizer deletion in (19 b), but not in (19 a).

- (19) a. *Fay believes that*∅ *Kay left.*
 b. *Fay believes, but Ray doesn't, that*/*∅ *Kay left.*

Now according to Lightfoot the same condition also prohibits straightforward subject extraction in a variety of languages: (20) is a problem not just for English.

(20) **Who_i do you think e_i that e_i saw Fay?*

Lightfoot claims that this constraint is dysfunctional because clearly structures like (20) are needed by speakers, as is shown by auxiliary structures employed in diverse languages to “rescue” the structure.

But such a case shows nothing about the dysfunctionality of UG constraints. Lightfoot’s fundamental error is that he does not distinguish the **functional** effects of the constraints from their **incidental** effects. This distinction has been widely discussed by philosophers (e.g. Wright 1973, Millikan 1984): For example, pumping blood is a functional effect of the heart, but throbbing noises are incidental effects. The heart both pumps blood and makes throbbing noises, but it is only the former effect that the heart has been designed by selection to produce. The throbbing noises may sometimes be inconvenient, but these incidental effects cannot be used as an argument that the heart is dysfunctional or is an accidental side-effect of some other adaptation. Lightfoot (1999: 249) admits that the condition on movement traces “may well be functionally motivated, possibly by parsing considerations”, so the ungrammaticality of (20) in English only shows that grammatical constraints may have incidental effects, not that they may be non-adaptive or dysfunctional.

Even less impressive is Haider’s (1998) case for the dysfunctionality of superiority effects in English *wh*-movement. Haider notes that in some languages (e.g. German) the counterpart of (21 b) is grammatical.

- (21) a. *Who bought what?*
 b. **What did who order?*
 c. *What was ordered by whom?*

However, as Haider notes, Kuno & Takami (1993) have proposed a usage-based explanation for the contrast between (21 a) and (21 b), which starts from the observation that sentences like (21 b–c), in which agents are sorted on the basis of themes, are “unnatural in normal circumstances”. This is exactly the kind of situation in which we would expect a grammatical constraint (“WH-SUBJECT > WH-OBJECT”) to arise in the process of diachronic adaptation (analogous to the constraint ANIM > INANIM of §6.5). Haider’s objection against the functional explanation is that not all languages show its effects, but this reveals a fundamental misunderstanding of the way in which user optimality and grammatical optimality work: In languages like German, the universal constraint is simply violated, and the counterpart of (21 b) is grammatical because other constraints are ranked higher. Thus, far from being dysfunctional, the constraint against (21 b) is functionally motivated, and the fact that it prohibits some potentially useful structures is in no way special (for instance, nobody would suggest that a constraint against morphological repetition is dysfunctional just because it rules out potentially useful words like **friendlyly* or **monthlily*).

I conclude that the case for dysfunctionality of grammatical constraints is very weak. As we have seen, many grammatical constraints correspond directly to user constraints, and the likelihood that there is no causal connection between the two sets of constraints is infinitesimally small.

One possibility is of course that the grammatical constraints arose in some way as an adaptive response to the user constraints in biological evolution, not in diachronic linguistic evolution. This has been proposed by various authors (e. g. Pinker & Bloom 1990, Newmeyer 1991), and it is a possibility that must be taken very seriously. However, a full discussion of this possibility is beyond the scope of this paper. The main practical problem with the biological-adaptation scenario is that it is necessarily more speculative than my scenario of diachronic linguistic evolution. I think it is a sound methodological principle to try the more empirically constrained explanations first, before speculating about prehistoric events that have left no direct trace. Moreover, the violability of the optimality constraints makes them poor candidates for innate devices, whereas violability follows automatically if the constraints arise in diachronic adaptation. But even so, I expect the argument made in this paper to be challenged primarily from the direction of biological evolution, so theoretical linguists are well advised to watch developments in biological evolutionary linguistics closely.

8. Conclusion

My main argument in this paper has been that optimality constraints of the type postulated in Optimality Theory, which are usually conceived of as stipulated elements in a pure competence theory, need to be further analyzed in terms of constraints on language use. Otherwise it remains mysterious why the constraints that we find applicable in languages are the way they are, and why many logically possible constraints play no role in any language (e. g. NOONSET, OBLIGATORYCODA, DON'TSTAY, MAXFUNC, INANIM > ANIM, and so on, i. e. the exact opposites of the constraints we have seen).

The mechanism proposed here for linking grammatical constraints to user constraints is diachronic adaptation: In language change, variants are created from which speakers may choose. Being subject to various constraints on language use, speakers tend to choose those variants that suit them best. These variants then become increasingly frequent and entrenched in speakers' minds, and at some point they may become obligatory parts of grammar. In this way, grammars come to be adapted to speakers' needs, although speakers cannot shape language actively and voluntarily. Grammatical constraints are thus the way they are because they have arisen from user constraints in a diachronic

process of adaptation.¹⁰ Diachronic adaptation in language is in many ways analogous to adaptation in biological change.

That grammatical structures are typically adapted to language users' needs in a highly sophisticated way is an old insight, but how exactly this adaptation should be explained is rarely discussed even by functionalists. Croft (1993: 21–22) notes that “the philosophical analogy between linguistic functional explanations and biological adaptation is not always fully worked out in linguistics”. The Teleological Fallacy appears to be so powerful that linguists have rarely seen the necessity of providing a theory of diachronic adaptation. But that such a theory is needed has been recognized by other authors as well (e.g. Bybee 1988, Kirby 1994, 1997, 1999, Durie 1995, Nettle 1998). Hall (1988) observes that in addition to finding “underlying principles, probably of a psychological or functional nature”, we must

“attempt to establish the mechanism by which the underlying pressure or pressures actually instantiate in language the pattern under investigation. This latter requirement will involve the investigation of diachronic change for some properties and of phylogenetic evolution for others.” (Hall 1988: 323)

Diachronic adaptation provides an account of the paradoxical situation that intentional actions of individuals, which have nothing to do with grammatical optimality, can have the cumulative effect of creating an adapted grammar, consisting of constraints that are good not only in a theory-internal sense, but also from the language users' point of view. Situations of this kind, in which a large number of micro-events give rise to a macro-structure in a surprising way, go by different names in the literature: “emergence” (Kirby 1997, 1999), “invisible-hand process” (Keller 1994), “spontaneous order” (Keller 1997), “self-organization” (Lindblom et al. 1983), “synergetic process” (Köhler 1986). So far there is no unified conceptual framework and terminology in linguistics for such phenomena, but it seems clear to me that this is a very promising paradigm.

If my proposal is correct, then the grammatical constraints are not innate, and are not part of Universal Grammar. They arise from general constraints on language use, which for the most part are in no way specific to language. This does not, of course, mean that there is no UG, no innate mental organ that is

10 Note that I am not claiming that all of language change is adaptive and motivated by user optimality of one kind or another (contra Vennemann 1993). For instance, grammaticalization changes, which probably account for the great majority of morpho-syntactic changes, can hardly be described as adaptive (cf. Haspelmath to appear). I tend to agree with Dahl (1999), who describes grammaticalization as a kind of counter-adaptive inflationary process in which forms lose their functions and thus need to be replaced. My claim here is only that whatever adaptive structures we find synchronically must have their origin in an adaptive diachronic change.

specialized for linguistic skills. Clearly, there are universal properties of language that probably cannot be derived from constraints on language use, e. g. the fact that grammars generally do not contain numerical specifications (e. g. “a word may be at most 15 segments long”); or indeed the fact that humans use fairly rigid grammatical rules to begin with, rather than arranging morphemes in a random way and leaving interpretation to pragmatics (cf. Durie 1995: 279). But these features of language are so general that they have little to do with the grammarian’s everyday work.¹¹

The language-particular aspects of grammar that occupy most linguists most of the time can largely be accounted for in terms of conventionalized constraints on language use. The highly general constraints of OT have thus opened up new possibilities of (functional) explanation that were not available earlier in generative grammar.

Thus, by incorporating a theory of diachronic adaptation, linguistics can answer why questions, and is not limited to how questions (cf. Nettle 1998: 460). In this respect, linguistics is more like biology than like physics, more Darwinian than Galilean. Ridley (1994) puts it as follows:

“In physics, there is no great difference between a why question and a how question. How does the earth go round the sun? By gravitational attraction. Why does the earth go round the sun? Because of gravity. Evolution, however, causes biology to be a very different game, because it includes contingent history... Every living creature is a product of its past. When a neo-Darwinian asks ‘Why?’, he is really asking ‘How did this come about?’ He is a historian.” (Ridley 1994: 16–17)

In much the same way, I argue, a linguist who asks ‘Why?’ must be a historian.¹²

Eingereicht: 25. 2. 99

Überarb. Fassung eingereicht: 31. 5. 1999

11 In OT, they correspond to the components GEN and EVAL; the former has been largely ignored, apparently because of the implicit presupposition that it is not very interesting. However, from an innatist perspective it is the most interesting part of the theory, because it is the part that is the most likely to be innate.

12 Of course, whether one finds why questions interesting or not is a subjective matter. Hoekstra & Kooij (1988) argue that explaining language universals is not an important goal of generative linguistics. But I have doubts whether one can reach the goal of generative linguistics (discovering the nature of UG, i.e. answering a how question) while ignoring the question why linguistic structures are the way they are.