



Faculty Publications

2005

Analogical Modeling and morphological change: the case of the adjectival negative prefix in English

Don William Chapman

Brigham Young University - Provo, don_chapman@byu.edu

Royal Skousen

Brigham Young University - Provo

Follow this and additional works at: <https://scholarsarchive.byu.edu/facpub>



Part of the [Linguistics Commons](#)

Original Publication Citation

Don Chapman and Royal Skousen. "The Negative Prefix in English and Analogical Modeling of Language." *Journal of English Language and Linguistics* 9 (2005): 333-57.

BYU ScholarsArchive Citation

Chapman, Don William and Skousen, Royal, "Analogical Modeling and morphological change: the case of the adjectival negative prefix in English" (2005). *Faculty Publications*. 6554.

<https://scholarsarchive.byu.edu/facpub/6554>

This Peer-Reviewed Article is brought to you for free and open access by BYU ScholarsArchive. It has been accepted for inclusion in Faculty Publications by an authorized administrator of BYU ScholarsArchive. For more information, please contact ellen_amatangelo@byu.edu.

Analogical Modeling and morphological change: the case of the adjectival negative prefix in English

DON CHAPMAN and ROYAL SKOUSEN

Brigham Young University

(Received 19 November 2003; revised 9 February 2005)

This article examines the usefulness of Skousen's Analogical Modeling (AM) for explaining morphological change. In contrast to previous accounts of analogy, AM constitutes a general unified model of language that accounts for both sporadic and systematic changes. AM also provides explicit constraints on analogy that allow explanation of how morphological changes begin, which forms most likely serve as patterns for analogy, and which forms are most likely to change.

AM is then tested on the case of the adjectival negative prefix in English (*in-*, *un-*, *dis-*, etc.), using the Middle and Early Modern English portions of the Helsinki corpus as a basis for prediction. AM was given the task of using forms containing negative prefixes for one time period to predict the prefixes that adjectives would take in the subsequent time period. For each of the roughly seventy-year periods in the corpus, AM was able to predict valid prefixes about 90 percent of the time.

1 Analogy and morphological change

Language change remains difficult to explain. We can describe details of language that change, but explaining why and how language changes is always a considerable challenge (Lass, 1980). Broadly speaking, we are confronted with two types of changes – those that proceed regularly throughout a language system and those that do not. The impulse in explaining language change, as with other areas of linguistics, has been to concentrate on the regular changes and try to account for the greatest amount of regularity possible. This was the approach of the neogrammarians, who posited completely regular sound changes. Apparent exceptions to the regularity were explained, wherever possible, by regular patterns of exceptionality. Their approach was reasonably successful for sound changes (Hock & Joseph, 1996: 123).

But morphological change has not been as regular. Because of the sporadic propagation of morphological changes, the neogrammarians posited changes in morphology by means of analogy, in contrast to the more regular sound changes (Hock & Joseph, 1996: 154). The lure of analogy has remained strong as an explanation for morphological and lexical change, and claims of analogical change are easy to come by in the literature: OE *fæder* > *father* (cf. *brother*; *mother*), *swollen* > *swelled* (cf. *spelled*, *shelled*, *felled*) (Anttila, 1977; Lehmann, 1992: 219–36; Joseph, 1997: 362–4).

Yet analogy has also remained of limited use as an explanation, because it appears too unsystematic. Until recently, descriptions of analogy have largely been informal and impressionistic, so that any form that resembles any other form can be a candidate for analogy. Even if we acknowledged that some (or many) morphological changes occur

by analogy, what more can we say about the change except to recognize that it occurred? Perhaps we could propose an analogue that prompted the change, but traditional accounts of analogy have not given us any way of telling whether such proposed analogues are likely or not. In other words, traditional accounts lack constraints on the operation of analogy, so that there is no principled means of telling when or where analogy will operate. Without such principled constraints, appeals to analogy are usually little more than post hoc descriptions of a change that do not really explain how the change occurred.

Analogy would be much more useful as an explanation of morphological change if we had some principled way of describing and constraining the operation of analogy. A model of analogy will be the more satisfying the more precisely it can tell us when, where, and how analogy will operate. In explaining *when* analogy will operate, the model should help us identify the occasions that speakers are likely to use analogy to produce new forms. In explaining *where* analogy will operate, the model should give us some way of identifying which forms are most and which are least likely to be influenced, and in explaining *how* analogy will operate, the model should help us identify which forms are most or least likely to influence others.

Recently several approaches to language have come closer to providing a more rigorous approach to analogy. Not all of these approaches characterize themselves as analogy, but like analogy they rely on actual instances of data, instead of rules, for characterizing and predicting language behavior (Aha, Kibler & Albert, 1991; Daelemans, Zavrel, van der Sloot & van den Bosch, 2001; Medin & Schaffer, 1978; Nosofsky, 1988, 1990; Pierrehumbert, 2001; Riesbeck & Schank, 1989). There are numerous differences among the models, some of which have been discussed elsewhere (Daelemans, Gillis & Durieux, 1994; Chandler, 2002; Shanks, 1995). This article will focus on Skousen's Analogical Modeling (AM). While comparisons with other models will occasionally be made, that is not the major aim of this article; instead it aims to examine the usefulness of Skousen's model for explaining morphological change and to illustrate the model with the example of the adjectival negative prefix in English.

2 Analogical Modeling and explaining morphological change

In AM, explanations for when, where, and how analogy will operate follow from the crucial assumption that all language production is instance-based, not rule-based. Instead of storing rules in our grammars, we store instances of language, and when it comes time to produce or interpret language, we review the instances and choose those instances that are similar in some way to the given form that we are trying to produce or interpret. If we need to decide which prefix to use with the adjective *lucky*, for example, chances are very good that we have heard and stored instances of *lucky* and *unlucky*. When it comes time to produce the negative prefix, the stored form *unlucky* will be among the terms selected as possible analogues. Because the stored form *unlucky* is identical to the given (target) form, the probability of choosing it is extremely high. If we have not stored the given form, or if for some other reason we cannot retrieve it (such as momentarily forgetting it), forms close to the given will suggest a suitable choice.

Chances are much less likely that we have heard and stored the negative form for the adjective *consummate*, for example, so instead of retrieving the negative form, we find those forms that are similar to *consummate* and use the negative form that they take. As it turns out, *sensible* and *vincible* are similar in many ways: they have a sequence of a stressed vowel followed by /n/ then /s/, with /ə/ as the vowel in the last syllable. Since *sensible* and *vincible* take *in-* as a prefix, we might predict that *consummate* would also take the *in-* prefix. A more detailed description of how AM chooses analogues will be given in section 2.3 below.

2.1 When AM operates

On one level, the answer to the question of when analogy operates is straightforward in AM: it always applies. The mechanism that produces analogical change is the same mechanism that operates whenever we produce and interpret language. Whenever we produce language, we look for forms that are similar to the forms we need. Often we find exact matches, but often we have to rely on other forms. It is in relying on other forms that change can occur. So Analogical Modeling is not a model to use only when rules fail to describe behavior or change; the model invokes analogy whenever we produce language.

2.2 Where AM operates

The question of where analogy will operate also falls out from the instance-based approach of AM. Analogical Modeling always attempts to find a similar instance to a given context (i.e. a target form). Very often that similar instance will be identical – that is the target form will be among the stored instances, and the speaker will choose it. Those forms with a high likelihood of being recalled by most speakers will not be very susceptible to analogical pressure. Thus words that are frequently encountered in a language should be less likely to change. It is the forms that are not often encountered that should receive heavy analogical pressure and will be most subject to change. This indeed seems to be the case with morphological change and other changes that depend on structure (Phillips, 2001; Bybee, 2001: 12).

2.3 How AM operates

The question of how analogy operates is also addressed directly in Skousen's model, where it is left vague in traditional descriptions. Somehow a model of analogy must constrain the forms that can influence other forms; otherwise we have no principled way to describe the effects of analogy. In Analogical Modeling, three important properties affect the probability of selecting a particular exemplar as an analogue (Skousen, 1989: 4):

1. *proximity*: database items that share more features with the given form will appear in more supracontexts and will therefore have a greater chance of being used as an analogue.

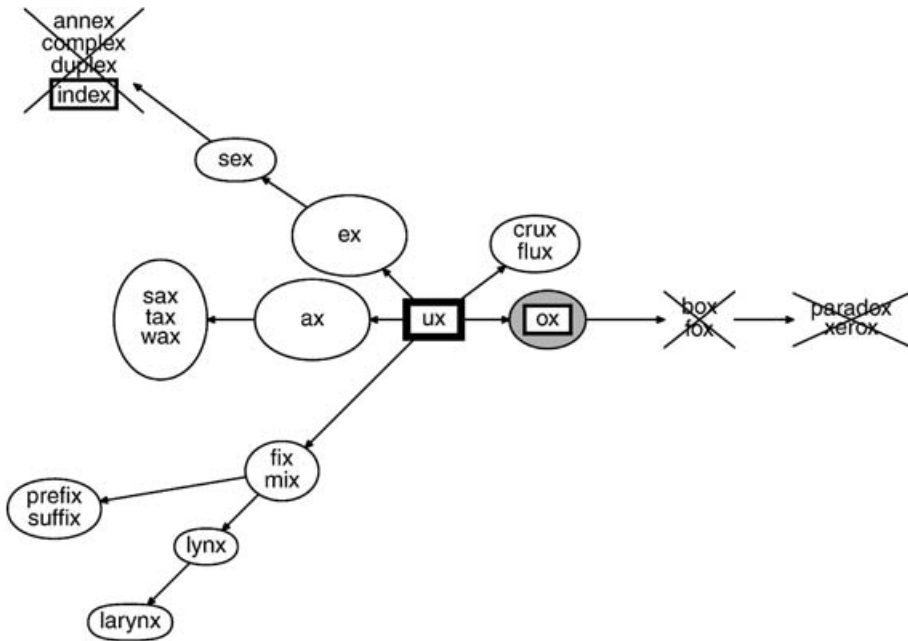


Figure 1

2. *gang effect*: if the example is surrounded by other examples having the same behavior, then the probability of selecting these similarly behaving examples is substantially increased.
3. *heterogeneity*: an example cannot be selected as the analogical model if there are intervening examples with different behavior closer to the given context.

These properties can be illustrated conceptually with an example of determining the plural ending of a newly encountered word *ux*.

As seen in figure 1, *ux* is close to *ox*, so the plural *-en* is possible. This prediction would be the result of proximity. On the other hand, *ex* and *ax* are just as close, and they are joined by very many examples of words, slightly farther away, that use *-es* (*crux*, *flux*, *sax*, *tax*, *wax*, etc.). All these similarly behaving exemplars illustrate the gang effect; in most cases the gang effect will dominate over the choice of *-en*. Finally, two groups cannot exert any influence because their behavior is different from an intervening exemplar. The group of *box*, *fox*, etc. take the *-es* form instead of the *-en* form that *ox* takes, but *ox* intervenes between *ux* and the group *box*, *fox*, etc. The more distant group *index*, *duplex*, etc. mostly take *-es*, but since *index* can also have the form *indices*, that group's behavior is different from the behavior of the intervening examples *sex* and *ex*, which take only the *-es* form. These groups illustrate the principle of *heterogeneity*.

These properties do not simply identify the nearest neighbor to a given word; instead they allow for nonneighbors to act on the analogy as well. The allowance of more distant influence makes the model robust. Analogy is not restricted to a range of obvious

similarities, especially when no obvious similarities exist. We can still make predictions with the best examples we can find. For the adjective *abrogated*, for example, we may not find any term that is obviously close, but several words are close enough to exert some influence. When *abrogated* was tested using AM with data from the Helsinki corpus for the period of 1640–1710, the following set of possible analogues was computed:

<i>Prediction</i>	<i>Adjective</i>	<i>Likelihood of selection as an analogue</i>
dis-	obedient	0.5%
dis-	honored	22.3%
dis-	ordered	22.3%
un-	able	8.8%
un-	doubted	13.0%
un-	approachable	13.0%
un-	habitable	7.0%
un-	spotted	13.0%
<i>Summary for each prefix</i>		
dis-		45.1%
un-		54.9%

The words *obedient*, *honored*, *ordered*, *able*, and *approachable* are all similar because they begin with a vowel. The words *obedient*, *able*, and *habitable* all have a *b* following the initial vowel. The words *honored*, *ordered*, *doubted*, and *spotted* all end with the *-ed* suffix – and of those *doubted* and *spotted* have a /t/ before the *-ed* suffix. As this example illustrates, no single word stands out as an obvious analogue for *abrogated*. The closest terms, in fact, are *honored* and *ordered*, so in a nearest-neighbor approach the predicted prefix would be *dis-*. But in this example the gang effect of the larger group of words that use *un-* narrowly overrides the effect of the closer *dis-* words. And in fact, *unabrogated* is the only form attested in the *Oxford English Dictionary (OED)*, with the first attestation in 1535.

In AM, the creation of the analogical set does not come from a direct pair-wise comparison of a given test item and the other items in the dataset. That means that the similarity between a given test item and the other items in the data set is never directly calculated and is only indirectly relevant to the choice of an analogue. Instead AM chooses analogues according to the number of more general supracontexts that a data item shares with a test item. A supracontext is an increasingly more general representation of an exemplar. Suppose the term *ux* is represented as filling three slots or variables: $\emptyset \underline{u} \underline{x}$, where for the sake of illustration, letters are used, not sounds, and the first variable is filled by ‘nothing’.¹ The more general supracontexts for *ux* would be representations in which the value of one or more of the variables is simply not

¹ In AM, these ‘slots’ are usually referred to as *variables*, as they will be in this article; of course this use of the term *variable* differs from its use in other fields, such as sociolinguistics.

important; the variable slots are wild-cards, as it were. Thus one supracontext would be ‘anything’ (represented as $-$) + u + x . Since a variable can either be counted or unimportant (i.e. filled with a value or a wild-card), there are 2^n supracontexts, where n is the number of variables. For ux , the supracontexts are:

$\emptyset u x$	$- u x$	$-- x$	$---$
	$\emptyset - x$	$\emptyset --$	
	$\emptyset u -$	$- u -$	

The term ax will share four supracontexts with ux : $\emptyset - x$, $-- x$, $\emptyset --$, and $---$. The term fox will only share two: $-- x$, and $---$. The more variables that a data item shares with a test item, the more potential supracontexts can be shared, but since supracontexts are not included when their behavior is heterogeneous (as with the *box* and *fox* group in figure 1), a high correspondence of variables between a data item and test item does not necessarily mean that the data item will be chosen as an analogue, or even have a high probability of being chosen. Probably the best introduction to the computation of analogical sets is Skousen (2002: 11–26). One may also consult Skousen (2003), which is available on the AM website (<http://humanities.byu.edu/am/>).

In answer to the question of which forms can serve as models for analogy, they are those that share similar features with the given form and that are not blocked by intervening forms with a different behavior. The closer the model is to the given context, the stronger the influence, and multiple nearby models behaving alike further strengthen the influence.

2.4 Direction of change

AM also accounts for the direction in which forms can change. Usually the direction of morphological change is from relatively less productive forms to relatively more productive forms. Where it makes sense to speak of regular and irregular forms, as with the English past tense or plurals, changes in morphology have tended to move from irregular to regular forms. In derivational morphology, where terms like regular and irregular may not be as clearly applicable, the changes generally move from those forms with more restriction on productivity to those forms with less restriction. Among negative prefixes, *un-* would have the fewest restrictions, and hence correspond to the regular forms of inflectional morphology. But not all irregular (or less productive) forms change, and sometimes regular (more productive) forms give way to irregular (less productive) forms. If we assume a single ‘regular’ form (the pattern with the fewest restrictions), we can note five possibilities for change:

1. an irregular form becomes regular (healp > helped; infortunate > unfortunate)
2. an irregular form changes to a different irregular form (brought > brang; disobedient > inobedient)
3. a regular form changes to an irregular form (sneaked > snuck; unholy > disholy)
4. an irregular form remains irregular (wrote; intolerable)
5. a regular form remains regular (loved; uneasy)

AM handles all these cases with a single mechanism. In the instance of no change (nos. 4 and 5), the search for an analogue will usually yield a form identical to the given item that is being searched for. There is no change because an exact match is found. Or, if an exact match is not found, similar-behaving forms provide an analogue. The possibilities for change in nos. 1–3 also come when an exact match is not found, but in these cases enough forms close enough to the given item (or target) behave differently from the given form. In most instances the pressure of other forms will tend toward the most productive forms. This is how the model accounts for no. 1 (irregular or less productive forms to regular or more productive forms). Changes from more productive to less productive forms (no. 3) or from one less productive form to another (no. 2) operate from the same mechanism, but the stored instances in the database will contain some less productive forms that are very close to the given. The influence of these close forms will often override whatever pressure more regular, but more distant, forms exert.

2.5 *Role of language input in change*

The role of linguistic data in language change is a crucial question. Somehow speakers incorporate what they hear into their own speech, and any model of language change ought to account for such incorporation. In other words, a theory of language change should account for how new forms are actuated and propagated throughout a speech community (McMahon, 1994: 44–6). Traditionally, linguistic theories have approached linguistic data through a derived, intermediary system. That is, from given data one derives a system: either a set of symbolic rules (with perhaps the help of innate principles, as in some generative approaches, including Albright & Hayes, 2003), an interrelated system of variable factors (as in variationist theories like VARBRUL), or an artificial neural net (as in connectionism). Such derived systems are then used to predict language behavior, rather than the original data.² In derived systems like these, an important question is how and when linguistic data are evaluated to create or alter these systems. In generative models of language change, for instance, the role of language input comes crucially, if not solely, when children form their grammars (Lightfoot, 1999: 77–108); in other words, changes come between generations.

In Analogical Modeling, on the other hand, the data are directly used to make language predictions, not through some intermediary system of relationships. The role of language input is central for all speakers at all times. The role of primary linguistic data is not restricted to the formation of new grammars: since the rules of a grammar are inferred from stored instances of language, there are no grammars to form. Instead, the

² Even TiMBL, as it is usually implemented, is an intermediate system that is derived from the data, since the TiMBL system requires a training stage beyond collecting the data – namely, determining the informational load or gain for each isolated variable, as if each variable has an independent existence.

primary linguistic data are potentially available for every production. The analogical pressure for any given form can change, then, as the data change or expand over time for a speaker.

The appeal of such an approach for explaining language change should be apparent. The same mechanism that produces language produces language change. Speakers invoke analogy every time they produce language, so the potential for creating new forms is present with every utterance. We do not need extra machinery to account for language change, since synchronic productivity and diachronic analogical change proceed from the same mechanism.³ And those analogical changes can proceed toward the more regular forms, as we would expect, but they can also move toward irregular forms, as we see sometimes occurs in language change.

2.6 *Variation and AM*

There is no question that much language change proceeds through stages of variation. The sections above have already demonstrated how variation can be actuated: new forms can be produced when a speaker cannot access a stored exemplar and must therefore choose a close analogue. But another issue of variation is how speakers can know which proportions of variants are appropriate for their given speech community. In other words, how are the proportions of variation maintained? Again, AM accounts for actuation and maintenance with a single mechanism. Somehow speakers behave as if they know the probabilities of the variants; they produce variants in a statistical distribution. But do speakers somehow make some kind of statistical judgment in choosing between variants? Without requiring statistical judgments from speakers, AM accounts for the variable proportions with its assumption of stored instances. If we select our linguistic forms from among stored instances, the chances of selecting one or another variant will vary according to the proportions of the stored instances. In short, the distribution of linguistic variables will fall out automatically from their stored proportions (Skousen, 1989: 77–81).

This description of change by Analogical Modeling does not differ sharply from previous, less formal notions of analogy: high-frequency words will be resistant to analogy, most changes will be toward regular forms, changes can move toward irregular forms if there is a strong model for the change. But AM formalizes these observations and gives us a way to look at the details of analogical change. In short, it allows us to identify and examine the constraints.

The aim of this first section has been to show the characteristics of AM that make it an attractive model for explaining morphological change: AM preserves the intuitive notion that speakers change language forms by relating them to other forms, while at

³ Yet there can be other mechanisms, as there are undoubtedly physical aspects to sound change; cf. Ohala (1974, 1993).

the same time it gives us principled ways of describing and discerning those changes. It remains in the rest of the article to describe the operation of AM in more detail and to test its empirical adequacy. It will be shown how AM works as it is applied to the development of negative prefixes in English.

3 Empirical test of Analogical Modeling

3.1 Test case and hypotheses

We are getting closer to explaining language change when we can predict the changes that will occur in a given period of a language. The test for AM in the next part of the article is for it to predict forms of the negative prefix in the historical periods of English. A significant challenge for a model of language change – in this case AM – is to see whether the model can use present forms in the language to predict future forms. In other words, can AM predict the negative prefixes that future adjectives will take based on the words with negative prefixes that are current in a given stage of the language? Can all the words with negative prefixes that are current in the year 1350, for example, predict the negative prefixes that adjectives like *certain* or *easy* will take in the year 1420?

Predicting which prefix to use with which base is not trivial, and explanations for this aspect of speakers' knowledge have not been sufficient. Generative approaches have been largely silent. Allen (1977: 2–3) notes that *in-* cannot occur with words derived from the suffixes *-ish*, *-ful*, *-ing*, *-ed*, *-some*, *-ous*, *-like*, *-worthy*, *-ly*, *-y*, and accounts for this distribution by claiming a boundary condition of these suffixes that will not allow *in-*.⁴ Allen's analysis, like most other generative studies, otherwise focuses on the boundaries of affixes and their effect on the phonology and morphology of the complex word. The matching of individual affixes with individual bases has been presumably regarded as a lexical issue, with the matches specified in the lexicon. Such an approach does not readily allow for predictions. Even at the synchronic level, we would need a fuller description of these apparently lexical processes to make predictions of the matches. At the diachronic level, the generative models have been even less informative. If the matching between individual prefixes and bases is specified in the lexicon, how do changes in that matching happen? How do the specifications of bases and affixes become established in the lexicon in the first place, and how do they change? These are the questions that would need to be answered to explain examples of language change such as the use of the negative prefix in English.

Less formal descriptions, while probably accurate, do not readily provide for predictions either, because they are insufficiently rigorous. Several have noted that *in-* usually attaches to Latinate bases (Marchand, 1969: 168; Quirk et al., 1985:

⁴ It is unclear how one can independently identify the boundary of these suffixes. Of course, one could also note that each of these suffixes is of native English origin, but it is doubtful that speakers know etymologies.

1540), but how speakers know which bases are of Latin or French origin is not specified.

In a way, these less formal descriptions actually end up appealing loosely to analogy for the historical development of the negative prefixes in English (e.g. *un-*, *in-*, *dis-*). Foreign prefixes entered the language when foreign words were borrowed that contained those negative prefixes. Those words would not originally have been analyzable for speakers who did not know the source language, but as speakers borrowed more prefixed words along with their positive counterparts, like the pair *immortal* and *mortal*, they could begin to analyze the words and regard the originally foreign prefix as an English prefix. At some point the prefix would have felt natural enough to be used as a combinative with other words, including, eventually, native words (Strang, 1970: 188–91; Marchand, 1969: 129; Burnley, 1992: 445–6). The crucial component of this description – the point that speakers recognize a morpheme in borrowed words – really is a description of analogy. Speakers essentially recognize that several words with a foreign prefix are like each other. The speaker can then use that knowledge to extrapolate to other words that are similar.

These accounts suggest three alternative hypotheses for predicting the negative prefixes of adjectives. The first is that the combining of prefixes with stems is not predictable, but instead is a lexical matter. In other words, speakers memorize the prefix and stem combinations. Such a situation would predict no changes to negative prefixes and adjectives. Yet obviously changes have occurred, such as *incomprehensible* > *uncomprehensible*, *unprofitable* > *unprofitable*, *inalterable* > *unalterable*.

The second hypothesis is that the etymological source for the adjective explains the prefix: if English, the original *un-* is maintained; if Latinate, we get *in-* (and its variants), *non-*, etc.; if Greek, we get *a-*, *dis-*, etc. Again, this is not sufficient for the history of the adjectives. There have been numerous cases in which Latinate adjectives have acquired the *un-* prefix: *uncertain*, *unchangeable*, *uncharitable*, etc. There are even a few cases of English words occasionally being used with Latinate prefixes (*disholy*, *insteadfast*).

A third hypothesis would be that the original Old English prefix *un-* is the regular case throughout the history of the language, while all others are treated as individual exceptions and are memorized, thus predicting that if a new item enters the language, it will tend to be regularized; we should not get any examples of change toward any of the exceptional cases. This would be a species of analogy – much like paradigm leveling – but analogy has also been known to move from regular to irregular forms, as the examples above show.

In contrast to these hypotheses is the claim of AM that all forms are available as analogues to all other forms, and thus change can move in all directions outlined in section 2.4: change from irregular to regular, change from one irregular form to another irregular form, change from regular to irregular, no change of irregular form, and no change of regular form. In particular, the test is to see how well AM can predict the use and distribution of historically foreign prefixes in various stages of English. If the predictions are good, we would have evidence confirming that the spread of the prefixes occurred by analogy and that AM accounts for that spread.

Table 1. Occurrences of negative prefixes (Helsinki corpus)

	1	2	3	4	5	6	7	8
Period	1150–1250	1250–1350	1350–1420	1420–1500	1500–70	1570–1640	1640–1710	1710–50
un-	83	63	83	62	35	46	79	149
in-	0	1	12	16	19	23	34	18
im-	0	0	2	6	8	10	15	23
ir-	0	0	0	1	2	1	6	8
il-	0	0	0	0	0	0	0	3
dis-	0	1	5	3	7	10	8	0
non-	0	0	0	0	0	0	2	0
an-	0	0	0	0	0	0	0	0
a-	0	0	0	0	0	0	0	0
Total	83	65	102	88	71	90	144	201

3.2 Method

3.2.1 Dataset

The dataset for this test consists of all the words with negative prefixes from the Helsinki corpus of Middle English and early Modern English.⁵ We grouped the words into their subperiods within the Helsinki corpus as given in table 1.

We then used the prefixed words from one period as the analogical dataset – that is, the stored instances of language from which to make predictions. We used the prefixed words from the subsequent period as the test items – that is, the target forms for which we wished to predict prefixes. The words with negative prefixes current in the period of 1500–70, for example, were used to predict the negative prefixes that adjectives would take which occurred in the period of 1570–1640. We restricted our analysis to those adjectives that would most likely have been analyzable within their period by including only those words for which the corresponding positive form can be found in the Helsinki corpus or *OED*. Thus a form like *innocent* is included for the first few periods, because the term *nocent* can be found in the *OED* for those same time periods, while the term *infant* is not included, because the corresponding term **fant* is not attested at any time in English. Past participles are included if they show evidence of being used as adjectives, but not if they are only used in passive verb phrases (cf. Quirk et al., 1985: sec. 7.15–19). Forms that occur in an early period and later period are included in any intervening periods in which they may have failed to occur, on the assumption that the words would have remained current. Only a handful of such forms have been added. De-adjectival adverbs like *uneasily* were counted as an attestation for the base adjective from which they would have derived (e.g. *uneasy*) only if the corresponding adjective could be found in the positive form (e.g. *easy*).

⁵ The last column (1710–50) consists of data from the Lampeter corpus (1998). Those wishing to inspect the dataset should see the AM website: <http://humanities.byu.edu/am/>

Table 1 lists the prefixes we have examined, namely *a(n)-*, *dis-*, *il-*, *im-*, *in-*, *ir-*, *non-*, and *un-*. As the table shows, *a(n)-* does not occur in the Helsinki corpus, so it will not be examined further. Among the remaining prefixes, a question arises over how to treat the Latinate prefix *in-*, which takes the historically assimilated forms *il-*, *ir-*, and *im-* in accord with the initial sound of the attached stem (thus, *illegal*, *irreverent*, and *impossible*).⁶ These four prefixes can be treated as four distinct and unrelated prefixes, or they can be treated as a single morpheme (represented hereafter as *IN-*), with four allomorphs: *in-*, *il-*, *im-*, and *ir-*. The relationship between the four prefixes could be considered phonological, but only from a historical point of view; that is, one could posit some rule assimilating *nl* to *ll*, *nr* to *rr*, and *nm* to *mm*. Although such an assimilation did originally occur in Latin, its existence in English is dubious, especially since we can pronounce such sequences (for instance, with the prefix *un-*: *unlawful*, *unrepentant*, *unproductive*). Yet even though there is no synchronic phonological connection, this does not mean that the four *IN-* prefixes might not be considered lexically or morphologically related. For instance, speakers consider the indefinite article forms *a* and *an* related, even though there is (no longer) any phonological rule in English deleting final *n* when followed by a consonant (as there was in Middle English). Thus in English today we say *one boy* without any loss of *n*. In the same way, the four *IN-* prefixes might be considered related without specifying an actual rule of phonological assimilation.⁷ Since it is not clear whether the group *il-*, *im-*, *in-*, and *ir-* should be considered as a single prefix or whether each prefix in the group should be considered a separate prefix, we have analyzed the prefixes both ways and give the results of each analysis below.

Another question is whether to count instances of prefixed words from the Helsinki corpus by type or by token. Experimental evidence discussed by Bybee (2001) indicates that lexical prediction seems to operate according to type frequency, not token frequency. Also, Eddington (2004) has found in using Analogical Modeling in general that type frequencies give better results than token frequencies.⁸ In Analogical Modeling, by increasing the number of variables, the token counts go down for specific data items. Given enough variables, most of the token frequencies approach one – that is, we end up basically with a type frequency for most words (Skousen, 1989: 54).

⁶ We also need to be aware of one case of a reanalyzed prefix *ig-* (namely, in the word *ignoble*) as belonging to this prefix as well, although historically this form derives from *in + gnable*, with loss of the *n* in the prefix but loss of the initial *g* in the stem form *noble*. From a semantic perspective in Modern English, *ignoble* is not equivalent to ‘not noble’ and can be ignored.

⁷ We purposely avoid defining different kinds of word and morphological boundaries in order to allow assimilation in one case and block it in another.

⁸ More precise methods of experimentation by Moscoso del Prado Martin, Kostip & Baayen (2004) have provided evidence that family frequencies work even better than type frequencies alone. Family frequency is determined by totaling up the number of different morphological forms that exist in actual databases for a given word or morpheme. Each distinct form provides a count of one, but the token frequencies for each particular form are ignored.

3.2.2 Variables

The basic principle of AM is that analogues are chosen by comparing a test item against every other stored instance. The algorithm does this by basically comparing the values for a given set of features or variables that the test item shares with all the stored instances.⁹ These features can be such things as syllable structure, sounds, morphemes, or practically any other characteristic that might seem important. The model is independent of the variables that are chosen to be compared: the method of comparing forms will be the same, regardless of which features are chosen. That means that AM can be used for a wide variety of phenomena – morphology, semantics, phonology, syntax, though this article considers only morphology and most studies to date have used only phonological and morphological variables.

In principle, any number of variables could be chosen, but because the comparison algorithm involves an exponential increase for every variable, there is usually a computational limit on the number of variables that can be specified.¹⁰ Current computer technology allows for twenty to thirty variables for most applications.¹¹ A crucial step for any problem using AM, then, is to include the variables that are most important for comparison. These variables will obviously differ from application to application. What must be avoided in selecting the variables is to select only variables that we think are crucial to predicting the outcome. If this were done, for instance, in predicting the indefinite article *a/an* in English, only the consonantal/vocalic difference for the immediately following segment would be needed. But then we would consistently predict *an* for every vowel-initial word that follows, even though actual experimental evidence shows a small percentage of leakage toward the *a* form for vowel-initial words. Such leakage can only be derived if we specify some ‘unimportant’ variables. For discussion, see Skousen (2002: 27–34).

We have mainly attempted to use phonological and morphological features to characterize the words in the dataset. Other kinds of features could be used; we know that the negative prefixes entered English attached to words borrowed from other languages, mainly Latin and French, so perhaps a word’s etymology should be an important component in the comparison. But very few speakers explicitly learn a word’s etymology, and whatever sense they might have about a word’s origin likely owes to other factors, like the sounds of the word. Semantic features could be used, but they are harder to identify, whereas phonological and morphological features have the virtue of being readily identifiable. As long as the number of features that can be used remains limited, there will always be the question of whether the most germane variables were chosen. What is remarkable in this case, as will be discussed below, is

⁹ As noted in section 2.3, the comparison is not a simple pair-wise comparison of variables. AM begins with the given item or context and constructs a lattice of all combinations of variables for that context. The items from the dataset are then distributed throughout the lattice, depending on their variables. From this lattice, the analogical set is derived according to the principles given in section 2.3.

¹⁰ The comparison operates by considering every case of a variable’s value being relevant or not. That means there are two possibilities (relevant or not relevant) for n variables. The number of cases examined, then, is 2^n .

¹¹ Since the running of this simulation, the computational limit has been raised to sixty variables.

Table 2. *Variables used in AM*

Variables		Adjective: <i>convenable</i>										Predicted prefix: <i>un-</i>								
Variable	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Value	k	0	k	ə	n	v	n	v	e	n	ə	ə	b	ə	l	0	l	2	V	əb əl
	First sound	First syllable					Second syllable					Last syllable			↑ Last sound	↑ Stressed syllable	Base & suffix			

Explanation of variables

- 1 First sound of the word.
- 2–6 First syllable of the word (if the word consists of more than one syllable). Syllables are weakly defined as being made up of the vowel, the two sounds preceding the vowel, and the two sounds following the vowel.
- 7–11 The second syllable of the word (if the word consists of more than two syllables).
- 12–16 The last syllable of the word.
- 17 The last sound of the word.
- 18 The syllable on which the stress occurs; 0 if the word is single syllable.
- 19–20 Whether the word contains a suffix. If there is a suffix, these variables denote the type of base to which the suffix is attached (noun, adjective, verb, etc.) and the specific form of the suffix.

how well just the phonological and morphological variables do in predicting negative prefixes.

The variables used in the analyses here are listed in table 2. Here the table specifies the values for the word *convenable*, for which the prefix *un-* is expected.

The goal in specifying variables is to provide enough information so that the word is recognizable – that is, the data representation should allow us to distinguish between the different data items. In order to predict the negative prefix for the adjectives in English, we basically specify the first two syllables and the last syllable. Since we are trying to predict a prefix, we make sure that most of our variables are found at the beginning of the word, nearest to the prefix itself. Several psycholinguistic experiments over the years have shown the importance of the initial part of the word in lexical access, with secondary influence from the end of the word (Marslen-Wilson & Zwiterslood, 1989).

We also represent each sound phonemically, not in terms of distinctive features (which would inordinately increase the number of variables). Some aspects of phonetic similarity are therefore missed, but it turns out that phonemic identity is sufficient to get the appropriate predictions. Finally, we do not specify the syllable structure precisely. We do identify the vowels and their surrounding sounds, but do not worry where the syllable boundaries occur, especially since there is in English considerable ambiguity regarding where syllable boundaries should be placed between certain types of consonants (Fallows, 1981). We also specify the basic overall stress pattern by identifying the stressed syllable only, although specifying the stress may be

Table 3. *Phonemic mergers and splits between 1570 and 1710*

Mergers between 1570 and 1640	Mergers between 1640 and 1710
$\varepsilon\bar{u}$ merged with $i\bar{u}$ – <i>dew</i>	some ε : merged with e : – <i>meat</i> (e : > i : later)
$a\bar{i}$ merged with ε : – <i>day</i>	\bar{o} : (from earlier \bar{ou}) merged with > \bar{o} : – <i>low</i>
a : merged with ε : – <i>make</i> (ε : > e : between 1640 and 1710)	$i\bar{u}$ > u : – <i>true</i>
	<u>Splits between 1640 and 1710</u> some \bar{u} split to Λ – <i>cut</i>

somewhat redundant, given that the stressed vowels are specified: in most cases of multiple-syllable words, contrasts between full vowels and the schwa vowel carry enough information to predict the stress.

We also specify two morphological variables that turn out to be largely irrelevant, but we put them in simply because they could be used to distinguish between adjectives and might therefore help in predicting the behavior. On the other hand, there was no obvious way to specify any semantic variables that would reasonably divide up the database, so we decided to ignore such. Most semantic variables would be so restrictive in their applicability that their specific influence on the prediction would be minimal at best. Ultimately, it appears that predicting the negative prefix in English is largely dependent upon the phonemic shape of the word, not its morphology, semantics, or etymology.

For most phonemes, we have assumed a Middle English pronunciation, even for words occurring after the sixteenth century, though all the long vowels and some of the short vowels shifted between the fifteenth and eighteenth centuries. Since the variables in the analogical set are phonemes and not phonetic features, the change in quality of the vowel is unimportant as long as its phonemic status is preserved. The model will predict prefixes for test items based on whether the items in the set had similar or different phonemes, irrespective of what the representation of those phonemes turns out to be. For mergers and splits, however, where vowels changed not only quality but phonemic class, we have represented the change in pronunciation as outlined in Barber (1997), Dobson (1957), Lass (1992), and Jordan (1974). The changes of phonemic class can be summarized as follows in table 3. For unstressed vowels, we have usually favored a schwa pronunciation, even though there is always considerable variation in pronunciation.

We have tried to construct the variables so that each one can be selected independently of any other variable. In other words, we have not built into the variables any structural relationships. If there are n independent variables, then there are 2^n possible sets of variables to consider. If we were dealing with syntactic or phonetic relationships, we would ultimately want to build in restrictions that would reflect the tree structures

necessary for determining whether a structure was a subcontext of a more specific structure. Such interconnections would lead to a much more complex computer program. The current algorithm has each variable act independently of all the others. The results are sufficiently accurate for phonological and morphological issues that for the time being there is no need to build in any more elaborate structure into the representations.

3.2.3 *Role of memory*

The role memory plays in selecting an analogue comes in the process of choosing which items from the dataset will be evaluated for possible inclusion in an analogical set. In AM, remembering an item does not mean selecting an analogue directly from memory – in all cases an analogical set is first constructed and then an analogue is chosen based on its probability. Instead, remembering an item means that it will be included with all the other remembered items in the computation of the analogical set. Once a data item has been included for evaluation, the algorithm for selecting an analogue proceeds the same for both common forms and rare forms. As speakers recall items, those items have the potential to become part of the analogical set, depending on their behavior with respect to the given. As speakers forget items, whether permanently or momentarily, those items will not be available for examination, and consequently cannot become part of the analogical set.

Perfect memory would mean that every item in the dataset would be examined for inclusion in the analogical set. But speakers probably do not have perfect memory: there are very likely forms that speakers sometimes remember and sometimes forget. The algorithm for AM has a way of excluding items in the dataset from consideration for the analogical set, but what makes an item more likely to be recalled and how much more likely than another item has not been worked out. It most likely has to do with frequency, and in an AM simulation, one could exclude infrequent items (either altogether or randomly), but for questions of frequency, our sample was simply too small – many of the words in the dataset occurred only once in the Helsinki corpus; if we had only selected words with multiple occurrences, we would have had too few for an adequate dataset. So for the present simulation, we have assumed perfect memory – that every stored instance could be recalled, evaluated, and potentially included in the analogical set from which an analogue is chosen (cf. Skousen, 1992: 349–53).

We are still left with the question of what to do with test items that have also been stored in the dataset. The adjective *kind*, for example, is a test item in every period, but it also occurs in the dataset from which analogues are selected in every period; an identical form like *kind* has an extremely high likelihood of being selected as the analogue. This situation would correspond to a speaker already having encountered a word and thus knowing which prefix is appropriate. It would make some sense to ‘remember’ forms that are identical to the test item – that is, evaluate them and include them in the analogical set – but because we wished to see how AM would predict a prefix, we have excluded a form when it is identical to a test item.

Table 4. *Analogical set for vincible*

Prediction	Adjective	Likelihood of selection as an analogue
dis-	dishonest	0.01%
un-	uncivil	0.07%
un-	unwilling	0.02%
im-	impatient	0.04%
im-	impossible	5.12%
in-	indifferent	0.14%
in-	indirect	0.03%
in-	indiscreet	0.04%
in-	indissoluble	0.24%
in-	infinite	0.04%
in-	inflexible	10.27%
in-	inconvenient	0.06%
in-	insensible	76.23%
in-	inviolable	7.70%
Statistical summary		
dis-	1	0.01%
un-	18	0.09%
im-	1028	5.16%
in-	18892	94.75%

3.2.4 Choosing an analogue

For each test item, the AM algorithm constructs an analogical set of items that exert analogical pressure, according to the principles of proximity, gang effect, and homogeneity outlined in section 2.3. A sample analogical set is given in table 4 for predicting the negative prefix for *vincible* in the 1640 database. In this example we see *insensible* exerting much pressure. The words *inflexible* and *inviolable* add to the gang effect, so that the probability of choosing *in-* as the prefix is very high.

From the analogical set, an analogue is chosen. In principle, an item with a low probability could be selected, though it would be much less likely than an item with a high probability. In the present simulation, however, we count the selected outcome as that prefix that has the highest probability across all words in the analogical set. In other words, the ‘winner’ is the prefix with the plurality.¹² In the example from table 4, the prefix *in-* had a 94.75 percent chance of being chosen, so it is counted as the prediction for *vincible*. If the prediction agrees with the actual behavior of the given, we call it a successful prediction. In this case, the prediction is successful, since the prefix expected for *vincible* is indeed *in-* in the test set.

¹² In principle, a prediction is not necessarily chosen by plurality. Instead it may be randomly chosen from the analogical set according to its probability within each analogical set. So in principle, *un-* could be chosen for *vincible* (as *unvincible*), even though the probability of its being chosen is extremely low.

Table 5. *Successful predictions (in-, im-, ir, and il- considered as separate prefixes)*

Period	1 1250–1350	2 1350–1420	3 1420–1500	4 1500–70	5 1570–1640	6 1640–1710	7 1710–50	Total
un-	63 (97%)	80 (82%)	61 (74%)	33 (66%)	38 (93%)	60 (77%)	108 (94%)	443 (84%)
in-	-	-	5 (71%)	12 (71%)	16 (57%)	29 (67%)	13 (30%)	75 (54%)
im-	-	-	-	1 (100%)	7 (54%)	10 (77%)	20 (74%)	38 (70%)
ir-	-	-	-	-	1 (100%)	-	5 (83%)	6 (86%)
il-	-	-	-	-	-	-	-	-
dis-	-	0 / 3	-	2 (67%)	4 (67%)	2 (20%)	0 / 8	8 (27%)
non-	-	-	-	-	-	-	0 / 1	0 / 1
Total	63 (97%)	80 (79%)	66 (74%)	48 (68%)	66 (74%)	101 (70%)	146 (73%)	570 (75%)

Table 6. *Successful predictions (in-, im-, ir, and il- considered as a single prefix)*

Period	1 1250–1350	2 1350–1420	3 1420–1500	4 1500–70	5 1570–1640	6 1640–1710	7 1710–50	Total
un-	63 (97%)	80 (82%)	61 (76%)	33 (77%)	37 (97%)	58 (82%)	104 (95%)	436 (87%)
IN-	-	-	8 (89%)	23 (92%)	32 (71%)	48 (74%)	45 (55%)	156 (69%)
dis-	-	-	-	2 (67%)	4 (67%)	1 (13%)	0 / 8	7 (25%)
non-	-	-	-	-	-	-	0 / 1	0 / 1
Total	63 (97%)	80 (80%)	69 (80%)	58 (79%)	73 (77%)	107 (75%)	149 (77%)	599 (80%)

4 Results

As noted above, we consider two cases in our predictions. In the first case, we treat the four *IN-* prefixes as completely unrelated. In the second, we treat them as a single prefix. The summation of all the successful predictions for each prefix in each period is given in tables 5 and 6, with table 5 representing the first case and table 6 the second. The first number in each cell is the number of successful predictions, when using the negative-prefixed words of the succeeding period as the test set. The number in parentheses is the percentage of successful predictions with respect to the total number of predictions for that prefix. If predictions were made for a given prefix in a given period but none were successful, the number following the slash represents the number of predictions made. A dash means that no predictions were made for that prefix.

As the numbers indicate, AM predicts the prefixes quite well. The overall success rates found in tables 5 and 6 are similar, but in several cases the single *IN-* analysis (table 6) performs better. In the latter case, the predictions match the actually occurring forms nearly 80 percent of the time.

The success rate would be higher if we took into account variation among the test items. In the 1640 test set, for example, both *indecent* and *undecent* are found. The prediction for *decent* is *in-* (*indecent*) in both instances, which means it was a successful prediction for *indecent* but not for *undecent*. But the missed predictions should not be considered failures if the prediction is an acceptable variation of the expected forms in the test set. To account for acceptable variation, we looked up all the missed predictions in the *OED*. We count a prediction an acceptable variation if it is attested in the Helsinki corpus or the *OED*. Of the 162 missed predictions made when considering the *IN-* prefixes as a single prefix, 79 (or 49 percent) were indeed attested in the *OED*; for treating *in-* prefixes as separate prefixes, the figures are

Table 7. *Valid predictions (in-, im-, ir, and il- considered as separate prefixes)*

	1	2	3	4	5	6	7	Total
Period	1250–1350	1350–1420	1420–1500	1500–70	1570–1640	1640–1710	1710–50	
un-	63 (97%)	88 (91%)	73 (89%)	49 (98%)	39 (95%)	73 (94%)	114 (99%)	499 (95%)
in-	-	1 (100%)	7 (100%)	13 (76%)	20 (71%)	37 (86%)	22 (50%)	100 (71%)
im-	-	-	-	1 (100%)	7 (54%)	10 (77%)	22 (81%)	40 (74%)
ir-	-	-	-	-	1 (100%)	-	5 (83%)	6 (86%)
il-	-	-	-	-	-	-	-	-
dis-	-	0 / 3	-	2 (67%)	5 (83%)	2 (20%)	2 (25%)	11(37%)
non-	-	-	-	-	-	-	0 / 1	0 / 1
Total	63 (97%)	89 (88%)	80 (90%)	65 (92%)	72 (81%)	122 (85%)	165 (82%)	656 (86%)

Table 8. *Valid predictions (in-, im-, ir, and il- considered as a single prefix, represented by IN)*

	1	2	3	4	5	6	7	Total
Period	1250–1350	1350–1420	1420–1500	1500–70	1570–1640	1640–1710	1710–50	
un-	63 (97%)	88 (91%)	72 (90%)	42 (98%)	38 (100%)	69 (97%)	109 (99%)	481 (95%)
IN-	-	1 (100%)	9 (100%)	25 (100%)	37 (82%)	57 (88%)	55 (67%)	184 (81%)
dis-	-	0 / 3	-	2 (67%)	5 (83%)	1 (13%)	2 (25%)	10 (36%)
non-	-	-	-	-	-	-	0 / 1	0 / 1
Total	63 (97%)	89 (88%)	81 (91%)	69 (97%)	80 (91%)	127 (88%)	166 (84%)	675 (89%)

86 of 190 (or 45 percent). And those are undoubtedly low percentages, since the *OED* is only a sampling of the variation that has occurred in English. When we combine the variations attested in the *OED* with the successful predictions from tables 5 and 6, we see an even higher success rate of predictions using Analogical Modeling, as shown in tables 7 and 8, which again distinguish between analyzing *in-*, *im-*, *ir-*, and *il-* as separate prefixes or as a single prefix. The notation is the same as that in tables 5 and 6.

Of course there would be no point in counting attested variants if there were so much variation that any prediction would have matched a variant. As it turns out, the variation is fairly restrained. Only a few of the adjectives in the missed predictions were attested with more than two prefixes in the *OED* (16 and 18, respectively, for the sets treating *in-* prefixes as one prefix or separately). Even when these items are left out of the calculations for tables 7 and 8, the percentages remain nearly the same (86.71 and 85.98 percent, respectively).

These numbers show that AM predicts a valid prefix for adjectives at a high rate. The percentage for the total number of valid predictions is near 90 percent for nearly all periods when the *in-* prefixes are treated as a single prefix. Even when they are treated as separate prefixes, the percentage is still near 85 percent except for the 1500–70 period. There were no occasions of AM predicting that a form could not occur, when in fact that form was attested.

The results of AM prediction also show that *dis-* is never predicted very well. The sound shape of words beginning with *dis-* does not ever seem to be distinctive enough to predict *dis-* at a high rate, while it is for the other forms. Some predictions of prefixes other than *dis-* are *unadvantageous*, *unagreeable*, *uncontented*, *illoyal*, *unpleased*, *unregarded*, *untempered*. And some invalid predictions with *dis-* are

Table 9. *Directions of change (IN- as a single prefix)*

	1250	1350	1420	1500	1570	1640	1710	Total
1 Stable reg	100.00%	96%	98%	94%	80%	73%	70%	84%
2 Stable irreg	-	-	30%	69%	84%	75%	87%	67%
3 Regularize	-	47%	58%	90%	100%	85%	83%	66%
4 Irregularize	-	33%	100%	50%	56%	38%	29%	36%
5 Irreg to irreg	-	-	-	100%	67%	67%	100%	67%

disconcerned, disexperienced, disaltered, disprepared, disexpected, discorrupted, and disapproachable. The implication is that *dis-* is not a productive prefix; when it is used, speakers are more likely to choose exact analogues than close ones. Again, this analysis agrees with other studies that conclude that *dis-* does not seem to be very productive (Zimmer, 1964: 27–8; Marchand, 1969: 161–2).

5 Discussion

5.1 Hypotheses and directions of change

These rates look good – high enough on their own to warrant confidence in AM. But just how good are the numbers? What do they really mean? It is not the purpose of this article to compare AM with other models for simulation (e.g. TiMBL), though the datasets are available for anyone who would like to compare models; the results of this present study can serve as a baseline.¹³ Instead, this article sets out to examine the feasibility of AM for explaining morphological changes, and for that, the results should be examined against the hypotheses outlined in section 3.1. In theory the model would need to show stability of regular (*un-*) and irregular (*IN-*, *dis-*) forms, movement from irregular to regular forms, movement from regular to irregular forms, and movement from one irregular form to another. The results of AM's predictions for each of these processes are given in table 9.

The measurement for each process is the combination of expected and predicted forms that were actually attested for each process as a percentage of the combinations that should occur. The first row represents the stable maintenance of the regular *un-* form. For this process, every form for which *un-* is expected should also be predicted to take an *un-* prefix. The measurement, then, is the words both expected *and* predicted to have an *un-* prefix as a percentage of all words expected to have an *un-* prefix, whether an *un-* prefix was actually predicted or some other prefix. The second row represents

¹³ For the record, the success rates for other studies using AM, like past-tense forms or indefinite articles, hover in the 80s and 90s. Our test of AM is not far from that. It would not be surprising to see other models achieve higher validity rates, as they have with other problems. One reason is that other models use a training stage. Once the data are analyzed, variables can be selected to make the predictions even stronger. As already noted, this conscious, metalinguistic intervention is not preferred if it can be avoided.

stable maintenance of an irregular form (*IN-* or *dis-*) and is measured similarly, except that the expected and predicted prefixes are irregular. The third row represents the process of regularization; the pattern for this process is the words that were expected to have an *IN-* or a *dis-*, but which were predicted to have an *un-*. The measurement is the *attested* instances of this pattern as a percentage of *all* instances of the pattern, whether attested or not. The fourth row represents the process of irregularization, and is measured similarly, except the pattern would be for words expected to have an *un-*, but predicted to have an *IN-* or *dis-*. The last row represents movement from one irregular form to another, and the pattern for this process is an expected form of *IN-*, but a predicted form of *dis-*, or vice versa. Again, the measurement is the total of attested instances of this pattern as a percentage of all instances, whether attested or not.

AM models four of the five processes reasonably well (regular stable, irregular stable, regularize, and irregular to another irregular). As we would expect, *IN-* forms are not stable initially, but as the number of *IN-* forms increases in the language, the forms become more stable. Conversely, the regular forms are most stable when there is little competition, namely in the early periods. But by 1640 and 1710, they become less stable, meaning that other prefixes are being predicted for adjectives expected to take *un-*. Sometimes this simply represents successful irregularization, but in most cases the predicted irregular forms fail to be attested, so these forms represent both unsuccessful regular stability and unsuccessful irregularization. In practically all cases, these predictions were for Latinate stems, like *intelligent*, *intelligible*, *favorable*, and *governable*. AM essentially recognizes that the stems look like other Latinate words and thus predicts the Latinate prefixes, even though *un-* was expected. Because these were originally Latinate (instead of native) terms, the failure to predict *un-* can be seen as less a failure of regular *un-* words maintaining their *un-* form, than a failure of Latinate stems to fully regularize once they acquire the *un-* prefix. They still look too much like words that take irregular prefixes.

The irregularizing process (row 4) largely reflects the same results. Only six of the eighty-one predictions represented in row 4 were for stems of native or at least Germanic origin: *inspeakable* (twice), *intoward*, *inbecoming*, *insettled*, *nonspotted*. These were all unsuccessful predictions. The rest contained Latinate stems, and among these, sixteen were attested in the *OED* with an *un-* prefix before they were with an *IN-* or *dis-* prefix. These sixteen cases would represent a true irregularization process, if the attestations of the *OED* accurately reflect the history of these words. That leaves at least fifty-nine words that represent a backsliding of regularization, as it were – a regular form had taken root, but not strongly enough to keep the original irregular forms from being predicted.

Much of the problem with the overprediction of irregular *IN-* or *dis-* for expected *un-* most likely comes from the dataset. The prediction of *in-* for *intelligent* and *intelligible*, for example, occurs because there simply are not enough occurrences in the dataset of Latinate stems beginning with *in* to serve as analogues. The only such form in the 1710 dataset is *intermitted*, and while it has the single highest probability of being selected at 27 percent, it is not enough to overcome the gang effect of all the other Latinate

stems that don't begin with *in*. Thus *IN-* is predicted over *un-*, by a margin of 52 to 32 percent, even though the *IN-* form most likely to be selected (*direct*) had only a 14 percent probability of being selected. Just one or two more occurrences of Latinate stems beginning with *in-* would very probably have produced the correct prediction.

Furthermore, for the period that had the most problems maintaining regularity (1710), the dataset came from the Helsinki corpus, while the test set came from the Lampeter corpus. As it happens, the dataset for 1710 has a large number of *IN-* forms while the test set has a large number of Latinate stems that occur with *un-*. For that reason, many of the *un-* forms from the test set are predicted to occur with *IN-*. The effect of differing corpora is noted when the datasets are run against themselves in a 'leave-one-out' test. Whereas only 76.5 percent of the *un-* forms in the 1750 test set were correctly predicted to occur with *un-* when the test set was run against data from 1710 (see table 8), over 95 percent were correctly predicted when the test set was run against itself (i.e. with data from 1750, excluding, of course, the given item when it occurs in the test set). The large number of *IN-* forms in the 1710 dataset appear to have created some anomalies.

So in most respects, AM predicted according to the processes fairly well. Its main weakness was in overpredicting *IN-* forms once the expected form had become *un-*. In other words, it did not capture the regularizing and maintenance of regular forms as well as it did the other processes.

5.2 Effect of excluding the given

One major reason that AM has more trouble representing maintenance of regular forms is that the simulation presented here only partially represents the process of choosing prefixes. The present simulation treats each test item as if it had never occurred in the language before. If the dataset contained the same form as the test item, the stored instance from the dataset was disregarded, and the selection of analogues was made from among the remaining instances. In the 1710 test set, for example, the word *questionable* occurs. *Unquestionable* also occurs in the dataset of 1710 from which analogues are chosen. But when predicting the prefix for *questionable*, the identical stored instance *unquestionable* was not allowed to serve as an analogue. Instead analogues like *estimable*, *placable*, and *movable* were chosen, so that the prefix *IN-* was predicted. The test essentially treats the prediction procedure as if a person had no stored instances of any of the test items and always had to choose from close, not identical, instances. Of course that is not natural. It is more likely that speakers would store instances of many words that would be chosen as exact analogues, as it were, when it comes time to produce those words. For such words, close analogues would only be selected if the speaker could not retrieve the stored instance of the exact analogue for some reason.

In fact, we would expect some items to be predictable *only* by selecting the stored instance as the analogue. How could *went* ever be predicted as the past form of *go*, for example? So the 100 percent validity rate is rightfully unattainable for the present simulation. What this simulation really represents is the degree that analogues can be

successfully chosen for new forms, and that degree seems high enough to show how new prefixes can be adopted over time.¹⁴

Since these results treat every prediction as if the word had not occurred before in the language, they best reflect what happens in times of lexical innovation, when many new words enter the lexicon. Not surprisingly, AM performs best for the periods of high lexical innovation, namely the fifteenth and sixteenth centuries (cf. Nevalainen, 1999: 336–58). Before 1420, not enough *IN-* words occurred to have much analogical effect, so *un-* is predicted very well, but not *IN-*. After 1640, many predicted *IN-* forms failed to occur. In this regard, the prefixes were becoming more fixed and would more often have to be predicted by choosing the exact analogue, rather than a close one. The model suggests that forms produced by their similarity to other forms were most likely to be acceptable in the fifteenth and sixteenth centuries. After that the speakers were more likely to use a form that they had already stored and recalled. This could be another way of saying that *IN-* was most productive in the sixteenth century, but became more fixed by the eighteenth century.

This same tendency was noticed by Heok-Seung Kwon (1997: 28) for words in Chadwyck-Healey's English Poetry Full-Text Database. The variability of *in-* and *un-* was highest during the sixteenth century, but by the eighteenth century, most words had settled on one prefix or the other. Several have remarked the same trend today, noting that the *in-* prefixes do not seem to be used much with new forms (Zimmer, 1964: 28; Marchand, 1969: 168–70), though a search for *in-* forms on the internet revealed that *in-* is still productive (Jewell, 2001: 15–43). The best predictions for negative prefixes in the calculations above, when the possibility of recalling an exact form is disregarded, come precisely when the most variability is tolerated. Indeed we would expect more variation when speakers are not finding the exact matches in their searches. Thus the model suggests that the sixteenth century was when analogy to other forms was used most often for extending negative prefixes.

6 Conclusion

It would appear that AM is not only theoretically rich for explaining language change, but it also performs well on an empirical test case predicting the negative prefix for adjectives in English. This test case suggests that AM can feasibly be used for explaining language change and that analogy can therefore remain a plausible explanation for morphological change. These results should encourage more investigation into AM as an explanation for other morphological changes and even changes in other systems of the language, such as phonology and syntax. The appeal of AM should be apparent: it provides us with a viable, coherent model for morphological change.

¹⁴ A more reliable simulation of how speakers actually choose (and change) prefixes for given words would somehow assign probabilities to each form that would reflect the likelihood of that form being stored among a given speaker's instances of language and being recalled when creating an analogical set. We have not specified such probabilities for this simulation.

Authors' address:

Brigham Young University

4047 JFSB

Provo, UT 84606

USA.

don_chapman@byu.edu, royal_skousen@byu.edu

References

- Aha, D. W., D. Kibler & M. K. Albert (1991). Instance-based learning algorithms. *Machine Learning* **6**: 37–66.
- Albright, A. & B. Hayes (2003). Rules vs. analogy in English past tenses: a computational/experimental study. *Cognition* **90**: 119–61.
- Allen, M. (1977). The morphology of negative prefixes in English. *Proceedings of the Northeast Linguistics Society* **8**: 1–11.
- Anttila, R. (1977). *Analogy*. The Hague: Mouton.
- Barber, C. (1997). *Early Modern English*. 2nd edn. Edinburgh: Edinburgh University Press.
- Burnley, D. (1992). Lexis and semantics. In Blake, N. (ed.) *The Cambridge history of the English language*, vol. 2. Cambridge: Cambridge University Press. 409–99.
- Bybee, J. (2001). *Phonology and language use*. Cambridge: Cambridge University Press.
- Chandler, S. (2002). Skousen's analogical approach as an exemplar-based model of categorization. In Skousen, R., D. Lonsdale & D. B. Parkinson (eds.), *Analogical Modeling: an exemplar-based approach to language*. Amsterdam: John Benjamins. 51–105.
- Daelemans, W., S. Gillis & G. Durieux (1994). Skousen's Analogical Modeling algorithm: a comparison with lazy learning. In Jones, D. (ed.) *Proceedings of the international conference on new methods in language processing*. Manchester: UMIST. **17**: 3–15.
- Daelemans, W., J. Zavrel, K. van der Sloot & A. van den Bosch (2001). *TiMBL: Tilburg memory-based learner, version 4.1 reference guide: induction of linguistic knowledge technical report 0104*. Tilburg: ILK Research Group, Tilburg University.
- Dobson, E. J. (1957). *English pronunciation 1500–1700*. Oxford: Oxford University Press.
- Eddington, D. (2004). Issues in modeling language processing analogically. *Lingua* **114**: 849–71.
- Fallows, D. (1981). Experimental evidence for English syllabification and syllable structure. *Journal of Linguistics* **17**: 309–17.
- The Helsinki corpus of English texts* (1991). Helsinki: Department of English, University of Helsinki.
- Hock, H. H. & B. D. Joseph (1996). *Language history, language change, and language relationship*. Berlin: Mouton de Gruyter.
- Jewell, D. W. (2001). The negative adjectival prefix in English. MA Thesis, Brigham Young University.
- Jordan, R. (1974). *Handbook of Middle English grammar: phonology*. Trans. E. J. Crook. The Hague: Mouton.
- Joseph, B. (1997). Diachronic morphology. In Spencer, A. & A. Zwicky (eds.), *The handbook of morphology*. Oxford: Blackwell. 351–73.
- Kwon, H.-S. (1997). Negative prefixation from 1300 to 1800: a case study in *in-/un-* variation. *ICAME Journal* **21**: 21–42.
- The Lampeter corpus of Early Modern English tracts* (1998). Chemnitz: REAL Centre, Chemnitz University. URL: <http://khnt.hit.uib.no/icame/manuals/LAMPETER/LAMPHOME.HTM>
- Lass, R. (1980). *On explaining language change*. Cambridge: Cambridge University Press.

- Lass, R. (1992). Phonology and morphology. In Blake, N. (ed.), *The Cambridge history of the English language*, vol. 2. Cambridge: Cambridge University Press. 23–155.
- Lehmann, W. P. (1992). *Historical linguistics: an introduction*. London: Routledge.
- Lightfoot, D. (1999). *The development of language*. Oxford: Blackwell.
- Marchand, H. (1969). *The categories and types of present-day English word-formation*. Munich: C. H. Beck'sche.
- Marslen-Wilson, W. & P. Zwiterslood (1989). Accessing spoken words: the importance of word onsets. *Journal of Experimental Psychology: Human Perception and Performance* **15**: 576–85.
- McMahon, A. M. S. (1994). *Understanding language change*. Cambridge: Cambridge University Press.
- Medin, D. L. & M. M. Schaffer (1978). Context theory of classification learning. *Psychological Review* **85**: 207–38.
- Moscoso del Prado Martin, F., A. Kostip, & R. H. Baayen (2004). Putting the bits together: an informational theoretical perspective on morphological processing. *Cognition* **94**: 1–18.
- Nevalainen, T. (1999). Lexis and semantics. In Lass, R. (ed.), *The Cambridge history of the English language*, vol. 3. Cambridge: Cambridge University Press. 332–458.
- Nosofsky, R. M. (1988). Exemplar-based accounts of relations between classification, recognition, and typicality. *Journal of Experimental Psychology: Learning, Memory, and Cognition* **14**: 700–8.
- Nosofsky, R. M. (1990). Relations between exemplar similarity and likelihood models of classification. *Journal of Mathematical Psychology* **34**: 393–418.
- Ohala, J. J. (1974). Phonetic explanation in phonology. *Papers from the Regional Meeting, Chicago Linguistic Society* **10**: 251–74.
- Ohala, J. J. (1993). Sound change as nature's speech perception experiment. *Speech Communication* **13**: 155–61.
- Phillips, B. S. (2001). Lexical diffusion, lexical frequency, and lexical analysis. In Bybee, J. & P. Hooper (eds.), *Frequency and the emergence of linguistic structure*. Amsterdam: John Benjamins. 123–36.
- Pierrehumbert, J. (2001). Exemplar dynamics: word frequency, lenition, and contrast. In Bybee, J. & P. Hooper (eds.), *Frequency and the emergence of linguistic structure*. Amsterdam: John Benjamins. 137–58.
- Pierck, R., S. Greenbaum, G. Leech & J. Svartvik (1985). *A comprehensive grammar of the English language*. London: Longman.
- Riesbeck, C. K. & R. S. Schank (1989). *Inside case-based reasoning*. Hillsdale, NJ: Erlbaum.
- Shanks, D. R. (1995). *The psychology of associative learning*. Cambridge: Cambridge University Press.
- Skousen, R. (1989). *Analogical modeling of language*. Amsterdam: Kluwer.
- Skousen, R. (1992). *Analogy and structure*. Amsterdam: Kluwer.
- Skousen, R. (2002). An overview of analogical modeling. In Skousen, R., D. Lonsdale & D. B. Parkinson (eds.) *Analogical modeling: an exemplar-based approach to language*. Amsterdam: John Benjamins. 11–26.
- Skousen, R. (2003). Analogical modeling: exemplars, rules, and quantum computing. In Nowak, P., C. Yoquelet & D. Mortensen (eds.) *Proceedings of the Twenty-Ninth Annual Meeting of the Berkeley Linguistics Society*. Berkeley: Berkeley Linguistics Society. 425–39.
- Strang, B. M. H. (1970). *A history of English*. London: Methuen.
- Zimmer, K. (1964). *Affixal negation in English and other languages*. London: W. Clowes.